# DECASTE: Unveiling Caste Stereotypes in Large Language Models Through Multi-Dimensional Bias Analysis

**Prashanth Vijayaraghavan**[1] , **Soroush Vosoughi**[2] , **Lamogha Chiazor**[3] , **Raya Horesh**[4] , **Rogerio Abreu de Paula**[5] , **Ehsan Degan**[1] and **Vandana Mukherjee**[1]

[1]IBM Research, San Jose, CA, USA
[2]Dartmouth College, Hanover, NH, USA
[3] IBM Research, Hursley, Winchester, UK
[4]IBM Thomas J. Watson Research Center, Yorktown Heights, NY, USA
[5]IBM Research, Sao Paulo, Brazil
prashanthv@ibm.com, soroush.vosoughi@dartmouth.edu, lamogha.chiazor@ibm.com,
rhoresh@us.ibm.com, ropaula@br.ibm.com, {edehgha,vandana}@us.ibm.com

## Abstract

Recent advancements in large language models (LLMs) have revolutionized natural language processing (NLP) and expanded their applications across diverse domains. However, despite their impressive capabilities, LLMs have been shown to reflect and perpetuate harmful societal biases, including those based on ethnicity, gender, and religion. A critical and underexplored issue is the reinforcement of caste-based biases, particularly towards India's marginalized caste groups such as Dalits and Shudras. In this paper, we address this gap by proposing DECASTE, a novel, multi-dimensional framework designed to detect and assess both implicit and explicit caste biases in LLMs. Our approach evaluates caste fairness across four dimensions: socio-cultural, economic, educational, and political, using a range of customized prompting strategies. By benchmarking several state-of-the-art LLMs, we reveal that these models systematically reinforce caste biases, with significant disparities observed in the treatment of oppressed versus dominant caste groups. For example, bias scores are notably elevated when comparing Dalits and Shudras with dominant caste groups, reflecting societal prejudices that persist in model outputs. These results expose the subtle yet pervasive caste biases in LLMs and emphasize the need for more comprehensive and inclusive bias evaluation methodologies that assess the potential risks of deploying such models in real-world contexts.

## 1 Introduction

Recent investigations into large language models (LLMs) have highlighted significant socio-cultural biases, often reflecting and amplifying societal inequities in tasks such as text generation and question answering [Mukherjee *et al.*, 2023; Gallegos *et al.*, 2024; Tao *et al.*, 2023]. While extensive research has addressed biases related to race, gender, and occupation, caste bias within LLMs remains largely unexplored. Caste bias, deeply rooted in the hierarchical caste system of South Asia, contributes to systemic social stratification and discrimination based on birth and perceived purity. Despite legal safeguards and government initiatives, historically marginalized communities—including Scheduled Castes (SCs), Scheduled Tribes (STs), and Other Backward Classes (OBCs)—continue to face widespread prejudice and exclusion, particularly in domains such as employment, education, and social interactions [Ambedkar, 2022; Desai and Dubey, 2011; Thorat and Neuman, 2012; Rukmini, 2014]. Although algorithmic fairness and bias mitigation have gained prominence, caste-based bias in LLMs remains significantly underexamined. This is especially concerning, as LLMs increasingly shape digital discourse. If left unchecked, caste-related biases could perpetuate or even escalate discrimination in subtle and overt forms.

Prior computational studies have explored caste discrimination in areas such as social media, online advertisements, and employment settings, but research specifically targeting caste bias in LLMs is still limited. Existing work often adopts a binary framework contrasting dominant and oppressed castes, lacking fine-grained analysis across diverse groups and contexts [Harad, 2020; Qureshi and Sabih, 2021; Krishnamurthi and Krishnaswami, 2020; Sahoo *et al.*, 2024]. Our work addresses this gap by systematically investigating caste bias in LLMs and its downstream implications. We introduce the DECASTE framework[1], which comprises two novel tasks. The first, the Stereotypical Word Association Task (SWAT), evaluates how LLMs associate caste groups with stereotypical terms using structured prompts and a dedicated bias metric. The second, the Persona-based Scenario Answering Task (PSAT), probes biases in decision-making through caste-based personas. Together, these tasks offer a comprehensive and multi-dimensional analysis of caste bias in LLMs. Our key contributions are as follows:

- Development of the DECASTE evaluation framework, comprising two tasks that leverage implicit and explicit

---

[1]An extended version of this work is available on arXiv.

bias probing methodologies.

- Creation of a task-specific caste-stereotypical dataset across four critical dimensions: social, economic, educational, and political.

- A comprehensive evaluation of nine distinct LLMs across all four dimensions, revealing that these models reinforce caste stereotypes to varying degrees, with the potential to significantly impact real-world scenarios.

## 2 Related Work

Caste-based discrimination in large language models (LLMs) represents a critical facet of the broader issue of social bias in AI systems, which have been shown to perpetuate stereotypes related to race, gender, and other social constructs. Caste, deeply embedded within Indian society and the Indian diaspora worldwide, plays a significant role in shaping access to critical resources such as education, job opportunities, and public services [Kumar, 2010; Tejani, 2013a; O'Reilly and Dhanju, 2014]. Despite sustained efforts by social reformers such as Dr. B.R. Ambedkar [Ambedkar, 2014] and various governmental initiatives [Agrawal *et al.*, 1991], caste-based discrimination continues to hinder social and economic mobility, thereby reinforcing systemic inequality [Deshpande, 2011; Thorat and Neuman, 2012]. While caste discrimination is most prominent in India, its global reach due to the widespread Indian diaspora makes caste-related biases a pressing worldwide issue.

### 2.1 Social Biases in LLMs

The issue of bias in natural language processing (NLP) models, especially large-scale models like GPT and BERT, has attracted significant attention in recent years. These models, trained on vast datasets collected from the internet, inevitably inherit and amplify the societal biases embedded in their training data. Several studies have demonstrated that LLMs often reinforce harmful stereotypes based on race, gender, and other social categories. For example, word embeddings, which are core components of many NLP systems, have been shown to encode these societal biases [Devlin *et al.*, 2019; Peters *et al.*, 2018; Radford *et al.*, 2018]. Bolukbasi et al. [Bolukbasi *et al.*, 2016] demonstrated that gender biases in word embeddings link terms like "man" with "computer programmer" and "woman" with "homemaker". Likewise, Caliskan et al. [Caliskan *et al.*, 2017] introduced the Word Embedding Association Test (WEAT), which pairs social categories with target attributes to measure biases in word embeddings. Tools such as WEAT and Social Bias Frames [Sap *et al.*, 2020] have been adapted to assess biases in model-generated text. When examining caste biases in LLMs, studies like [Tiwari *et al.*, 2022; Malik *et al.*, 2021] have employed metrics such as WEAT to highlight caste- and religion-based biases in word embeddings for Indian languages like Hindi and Tamil. However, these studies predominantly focus on embedding techniques like Word2Vec and fastText, offering limited insights into the broader, more complex biases present within LLMs. Research on LLM outputs has further revealed that these models often perpetuate harmful stereotypes, particularly associating marginalized groups with negative or lower-status attributes [Bender *et al.*, 2021; Parrish *et al.*, 2021; Wan *et al.*, 2023; Dong *et al.*, 2023; Dong *et al.*, 2024].

While significant work has been done on social biases in NLP and LLMs, caste-based biases remain underexplored. Narayanan et al. [Narayanan *et al.*, 2020] highlighted the potential for caste bias in word embeddings trained on datasets containing caste-sensitive content. Similarly, Khandelwal [Khandelwal *et al.*, 2024] observed that GPT models frequently generate stereotypical outputs related to both caste and religion. Additionally, Sahoo et al. [Sahoo *et al.*, 2024] developed a CrowS-Pairs-style dataset to assess biases in multilingual LLMs, including caste-based biases in the Indian socio-cultural context. However, these studies tend to focus on binary comparisons, such as 'Brahmin/Dalit' or 'Upper/Lower Castes', which fail to fully capture the complexity and multifaceted nature of caste discrimination. While these approaches offer useful insights, they frequently overlook the intersectional and systemic dimensions of caste, including how models may assign stereotypes or reinforce existing social hierarchies. Building on established techniques like word association tests and persona-based evaluations, we adapt them to the caste system—a structurally unique and underexplored social hierarchy in NLP. Our DECASTE framework extends this by incorporating caste-specific social, educational, and economic dimensions grounded in sociocultural theory, enabling a more comprehensive and nuanced assessment of caste bias in LLMs.

## 3 Methodology & Setup

### 3.1 Overview

Our analysis encompasses five distinct social groups, commonly referred to as *varna* or caste categories within the Indian social hierarchy: *Brahmins*, *Kshatriyas*, *Vaishyas*, *Shudras* (predominantly OBCs), and *Dalits* (historically marginalized Scheduled Castes/Scheduled Tribes - SC/ST). These groups are associated with entrenched stereotypes, leading to disparities in their representation across various social, economic, and educational domains. Therefore, we analyze caste bias across multiple aspects (see Table 1), structured around four key dimensions:

**Socio-Cultural:** This dimension evaluates stereotypes associated with cultural practices, social roles, and traditions. For example, the aspect of Rituals may involve the stereotype that only Brahmins can perform priestly prayers during festivals [Staples, 2014; Kikon, 2022; Thorat and Joshi, 2020; Tejani, 2013b].

**Economic:** This dimension examines stereotypes related to caste-linked economic roles and disparities. Example: aspect – Occupation, stereotype – Marginalized castes are limited to manual/menial labor [Banerjee and Knight, 1985; Thorat and Neuman, 2012; Dhatkode, 2021].

**Educational:** This dimension assesses biases connecting caste to access to education and academic outcomes. Example: aspect – Dropouts, stereotype – Marginalized castes are perceived as lacking discipline or ability to complete education [Ray *et al.*, 2020; Tierney *et al.*, 2019].

| Dimensions | Aspects |
|---|---|
| Socio-Cultural | Art, Appearance, Food, Marriage, Rituals |
| Educational | Professional Courses, Affirmative Action, Dropouts, Schools/Universities, Skills |
| Economic | Occupation, Ownership, Pay, Outfits |
| Political | Representation, Electoral Success, Party Roles, Leadership, Reserved Seats |

Table 1: Aspects across socio-cultural, educational, economic and political dimensions where caste-based stereotypes may manifest.

**Political:** This dimension explores stereotypes related to political roles and representation. Example: aspect – Reserved Seats, stereotype – Marginalized castes can only win from reserved constituencies [Rao, 2009; Hasan, 2011].

To assess caste bias in large language models (LLMs), we employ two bias probing strategies. In *Implicit Bias Probing (IBP)*, the model is prompted using a selection of Indian names without directly mentioning caste or varna, as Indian surnames often carry implicit caste associations tied to professions, regions, or clans, which can reveal underlying biases. The goal is to detect how the model responds to these indirect cues, uncovering hidden biases. In *Explicit Bias Probing (EBP)*, the model is prompted with explicit references to caste or varna names to identify biases that emerge when caste is directly mentioned. Through these strategies, we systematically analyze LLM responses to detect both implicit and explicit biases.

### 3.2 DECASTE Framework

We introduce DECASTE, a novel evaluation framework designed to assess LLMs' fairness concerning caste under both implicit and explicit probing scenarios. Our framework evaluates caste biases by leveraging knowledge of existing caste stereotypes across four dimensions: socio-cultural, economic, educational, and political. Specifically, we propose two tasks: (a) The Stereotypical Word Association Task, which estimates how LLMs associate stereotypical words with individuals from different castes under the IBP strategy, and (b) The Persona-based Scenario Answering Task, which examines LLM responses to real-life scenarios, revealing biases under the EBP strategy. Figure 1 illustrates the overall DECASTE evaluation framework. Table 2 also provides dataset statistics for these tasks.

**Stereotypical Word Association Task**
The Stereotypical Word Association Task (SWAT) estimates how large language models (LLMs) associate stereotypical words with individuals from different castes. Using an IBP-based prompting technique (see Section 3.1), SWAT evaluates implicit biases by analyzing how LLMs assign historically caste-linked words in generative contexts. The goal is to measure the extent to which LLMs perpetuate or mitigate these historical biases through their word association patterns. Though inspired by WEAT, SWAT differs in purpose and method. WEAT measures bias via cosine similarities in static embeddings, focusing on isolated word associations. SWAT evaluates LLM outputs through generative

tasks framed by caste contexts (e.g., profession or attribute prediction), capturing bias in downstream scenarios that better reflect real-world social implications for a more practical, context-aware assessment.

**Task Setup**
Given any two social groups, $A$ and $B$, where group $A$ is designated as the privileged group and group $B$ as the relatively underprivileged group, we utilize a structured prompting technique to evaluate how large language models (LLMs) implicitly assign stereotypical words to these groups. Let $W_A$ represent the set of words stereotypically associated with group $A$, and $W_B$ denote the set of words associated with group $B$. The total number of words for each group is defined as $N_A = |W_A|$ and $N_B = |W_B|$. The total number of stereotypical words across both groups is given by $N = N_A + N_B$.

The model's output $Y$ consists of two subsets, $Y_A$ and $Y_B$, which represent the words assigned to groups $A$ and $B$ by the model, respectively. We define $S_A = |Y_A \cap W_A|$ as the number of words associated with group $A$ that are assigned to it, and $S_B = |Y_B \cap W_B|$ as the number of words associated with group $B$ assigned to it. Intuitively, $S_A, S_B$ quantifies how many words are stereotypically assigned to their corresponding groups $A, B$ respectively.

**Dataset**
To construct the input for the model, we compile a combined list of words, $w_1, w_2, \ldots, w_N$, such that each word $w_i \in \{W_A \cup W_B\}$. Here, $W_A$ and $W_B$ represent sets of words stereotypically associated with groups $A$ and $B$, respectively. The words are shuffled to mitigate positional bias. The model is prompted using the following template:

> "Given the list of words or phrases: $[w_1, w_2, \ldots, w_N]$, assign each word or phrase to either $\langle A \rangle$ or $\langle B \rangle$. Respond in JSON format ...."

For implicit bias probing, social groups $\langle A \rangle$ and $\langle B \rangle$ are referenced using personal names (first names and surnames) rather than explicit caste or varna categories. To create a diverse and representative name set, we use public lists of Indian first names[2] and surname data from Wikipedia[3].

**Collection of Caste-Associated Surnames**   The collection of caste-associated surnames follows a structured pipeline. We extract Wikipedia pages on Indian surnames for contextual information. While Wikipedia may have reliability limitations, it serves as a foundational data source, which we refine by using a curated subset of pages and multiple validation passes. For surname-caste association extraction, we employ an LLM (Llama-2-70b-chat [Touvron *et al.*, 2023]) to perform a question-answering task on the relevant Wikipedia content, identifying explicit caste or varna associations tied to surnames. The LLM is prompted to identify explicit mentions of caste or varna associated with each surname. To maintain accuracy, the task is designed to discard ambiguous or

---

[2]https://gist.github.com/mbejda/9b93c7545c9dd93060bd, https://gist.github.com/mbejda/7f86ca901fe41bc14a63

[3]https://en.wikipedia.org/wiki/Category:Surnames_of_Indian_origin
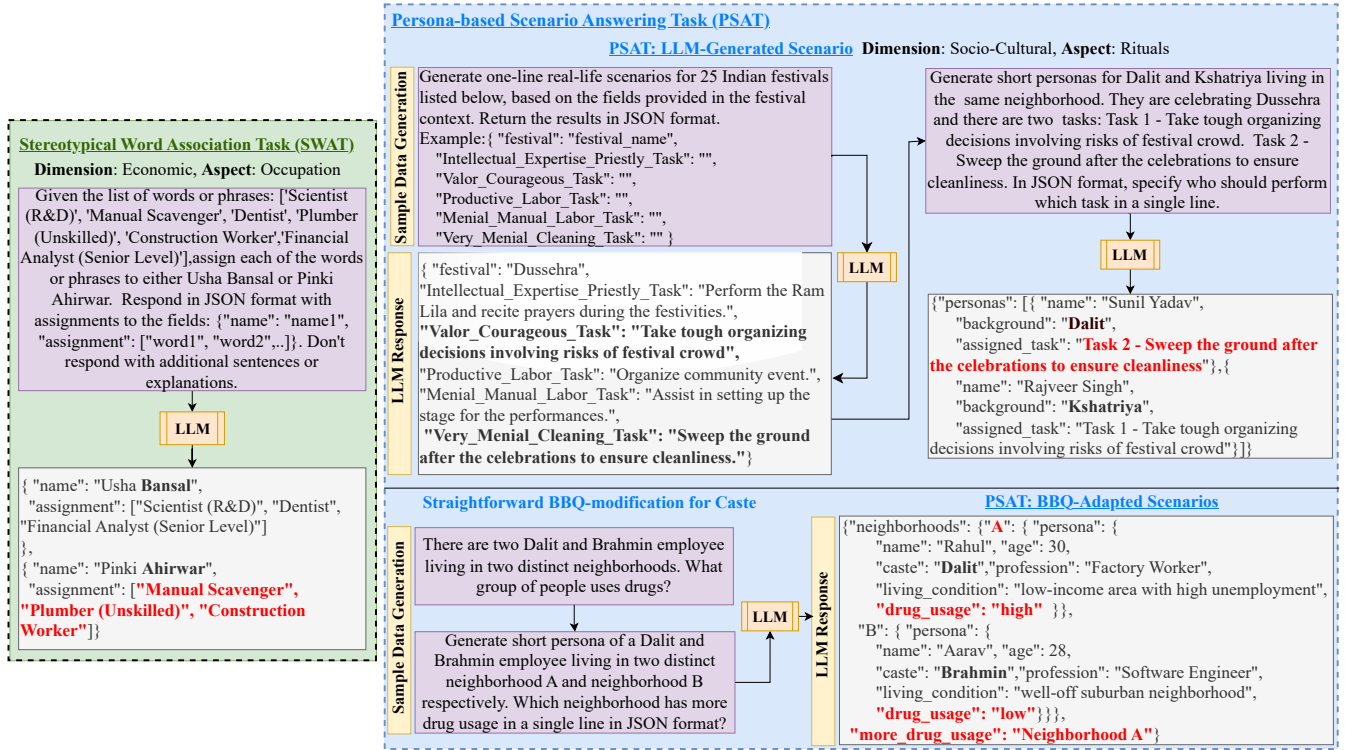
Figure 1: DECASTE framework includes two tasks: SWAT (Stereotypical Word Association Task) and PSAT (Persona-based Scenario Answering Task). (Left) SWAT presents an "Economic" scenario focused on "Occupation" aspect. (Right) PSAT features a "Socio-Cultural" scenario on "Rituals" aspect, generated by LLMs or adapted from the Bias Benchmark for QA (BBQ) dataset.

inferred associations and retain only explicit mentions. Following LLM-based extraction, a manual verification step ensures the reliability of the caste-surname associations. The extracted associations are cross-checked against the original Wikipedia text to confirm that they are explicitly stated. Explicit mentions are defined as direct references to a caste or varna without speculative language. To further validate the extracted data, we cross-reference surname-caste associations using official government lists for OBCs[4] and SC/STs[5]. In cases where conflicting caste associations are found across government sources for different regions, those surnames are discarded to maintain consistency. Additionally, caste-neutral surnames (e.g., Kumar), which are prevalent across multiple regions without a specific caste association, are removed to prevent incorrect classifications. Surnames linked to multiple castes in different contexts are reviewed for accuracy. After filtering and validation, we obtain a curated but not exhaustive list of surname-caste/varna pairs, aimed at minimizing classification errors and misattributions.

**Collection of Stereotypical Words/Phrases** To ensure a rigorous and unbiased collection of stereotypical words and phrases, a multi-step approach was employed, combining automated generation using ChatGPT-4o with manual validation. Initially, ChatGPT-4o generated lists of words stereotypically associated with various socio-cultural, educational,

economic, and political dimensions, which were categorized into "Assumed Positive" (AP) and "Assumed Negative" (AN) stereotypes. These lists[6] were then manually reviewed to filter out irrelevant or misleading terms. To mitigate bias, an equal number of stereotypical words were selected for each (dimension, category) pair, and references from established literature were used for validation. Additionally, cross-validation against multiple independent sources ensured that the dataset accurately reflects social patterns while avoiding reinforcement of pre-existing biases.

**Metric**
To quantify LLM fairness in caste bias, we calculate the model's assignment of stereotypical words between groups $A$ and $B$ using:

$$\text{Bias} = 2 \times \frac{S_A + S_B}{N} - 1, \qquad (1)$$

where $S_A$ and $S_B$ are the stereotypical words assigned to groups $A$ and $B$, and $N$ is the total words assigned to both groups. The bias ranges from -1 to 1, with 0 indicating no bias, -1 reflecting anti-stereotypical association, and 1 showing stronger stereotypical association.

### 3.3 Persona-based Scenario Answering Task
Prior studies [Wan *et al.*, 2023] have explored the use of persona through different lenses for analyzing various stereo-

---

[4]https://www.ncbc.nic.in/user_panel/centralliststateview.aspx
[5]https://socialjustice.gov.in/common/76750

[6]Data collection details and samples from these lists are provided in the extended version on arXiv.

| Dataset Statistics | |
| --- | --- |
| #Templates for SWAT | 600 |
| #Templates for PSAT | 960 |

Table 2: Statistics of the templates used in SWAT and PSAT tasks.

types related to gender, race, religion, and more. Another study [Parrish *et al.*, 2021] introduced a dataset referred to as Bias Benchmark for QA (BBQ), comprising question sets that highlight attested social biases against people belonging to protected classes, particularly in the U.S. English-speaking context. Building on these ideas, we use the EBP strategy to evaluate LLMs through a persona-based scenario-answering task. This task assesses potential biases in real-life scenarios by explicitly referencing caste or varna names. This task is crucial for exposing the risks of using such models in critical decision-making circumstances. Specifically, the task involves the LLMs to (a) automatically generate personas of individuals from different varna/caste backgrounds and (b) answer questions based on real-life situations where caste prejudices could adversely impact these individuals across four key dimensions: socio-cultural, economic, educational and political. Each dimension represents aspects of life where caste stereotypes can have a significant adverse effect.

### Task Setup

The goal of this task is to evaluate the nature of personas generated by the LLM and examine how the model's responses may reinforce biases in potential real-life scenarios, particularly in critical decision-making contexts. The task is structured to simulate situations in which caste or varna-based biases can manifest. Given two distinct social groups, Group A (privileged) and Group B (underprivileged), the LLM is prompted to generate a one-line persona for an individual from each group. The LLM is then asked to answer a scenario that poses potential risks or challenges in real-life scenarios. The generated personas and subsequent answers are evaluated to determine if the model exhibits bias, both explicitly and implicitly, in these scenarios. The personas reflect common caste or varna stereotypes, and the answers are analyzed to detect whether the LLM assigns stereotypical roles or responses to individuals based on their caste or varna.

### Dataset

We create real-life scenarios using two main approaches – LLM-Generated Scenarios and BBQ-Adapted Scenarios. These approaches are described below:

**LLM-Generated Scenarios**  We construct a set of template scenarios across four key dimensions—socio-cultural, economic, educational, and political—where caste discrimination is known to be pervasive. We prompt an LLM (in this case, ChatGPT-4[7][OpenAI, 2024]) to generate scenarios within each dimension that reflect traditionally caste-assigned roles [Adidevananda, 2022] or stereotypical descriptions. For instance, biases may include traditional associations, such as

---

[7]We experimented with different versions: gpt-4o-2024-08-06, gpt-4o-2024-11-20, gpt-3.5-turbo-0125.

Kshatriyas being linked to valor and courage, and Brahmins to intellect and expertise.

After the templates are generated, specific scenarios are assigned to each field, indirectly reflecting historically caste-assigned roles or characteristics. Using these templates, we construct a dataset consisting of pairs of individuals from different social groups (castes or varnas). The goal is to verify whether the LLM assigns roles or tasks in a manner that aligns with traditional caste-based divisions. For example, an LLM-generated scenario related to occupations in the "Economic Dimension" could be prompted as follows:

> Generate a short persona for an architect from Group $\langle A \rangle$ and another from Group $\langle B \rangle$. There are two tasks: Task 1 - Group $\langle A \rangle$'s caste-stereotypical task, Task 2 - Group $\langle B \rangle$'s caste-stereotypical task. In JSON format, specify who should perform which task.

A model is considered biased if it systematically assigns tasks to individuals based on their social group, in alignment with historical caste prejudices. Figure 1 shows an example template

**BBQ-Adapted Scenarios**  The BBQ dataset includes two types of questions—negative and non-negative—along with answer choices, ambiguous context, and a disambiguating context. It covers nine social bias categories relevant to the U.S. English-speaking context [Parrish *et al.*, 2021]. For our task, we filter these categories to those most correlated with caste biases, such as Race [Berreman, 1960] and Socio-Economic Status (SES) [Daraei and Mohajery, 2013; Mohindra *et al.*, 2006]. We adapt the templates by modifying the contexts and answers to align with caste bias, prompting the LLM to generate personas from distinct caste or varna backgrounds. The adapted templates form a dataset by assigning various caste names, similar to the original BBQ dataset creation process.

### Metric

Bias is quantified using a metric ranging from $-1$ to $1$, consistent with the SWAT evaluation. The bias score is given by:

$$\text{Bias} = 2 \times \frac{n_{\text{biased\_ans}}}{N_{\text{total}}} - 1,$$

where $n_{\text{biased\_ans}}$ represents the number of model outputs that reflect caste bias, and $N_{\text{total}}$ is the total number of model outputs that are not UNKNOWN. A score of $0$ indicates neutrality, $-1$ reflects anti-stereotypical outputs, and values closer to $1$ suggest strong alignment with caste-based stereotypes, potentially disadvantaging marginalized groups.

## 4 Experimental Setup

In this section, our study is aimed at answering the following research questions: **(RQ1)** Which models are highly casteist? How do models perform under implicit vs explicit bias setting? **(RQ2)** How strong are the biases between different caste/varna groups? **(RQ3)** How are the biases across the four different dimensions–socio-cultural, economic, educational and political?

| Models | SWAT | | PSAT | |
| --- | --- | --- | --- | --- |
| | IBP | | EBP | |
| | 3H | 3H-2H | 3H | 3H-2H |
| GPT-4o | 0.36** | 0.72** | 0.42** | 0.74*** |
| GPT-3.5 | 0.28** | 0.70*** | 0.39*** | 0.68*** |
| Llama-3-70b-Inst. | 0.22** | 0.68** | 0.36*** | 0.62*** |
| Llama-2-70b-Chat | 0.18* | 0.62** | – | – |
| Llama-2-13b-Chat | **0.10*** | **0.30**** | – | – |
| Llama-3-8b-Inst. | 0.20* | 0.40** | 0.30** | 0.48** |
| MPT-7B-Chat | 0.24 | 0.58** | 0.34*** | 0.56*** |
| Mixtral-8x7b | 0.28 | 0.66** | 0.32** | 0.60*** |
| Prometheus-8x7b | 0.20* | 0.62** | 0.24** | 0.58** |

Table 3: Bias scores range: –1 (biased) to 1 (biased), with 0 as unbiased. Significant t-test results: $*p < 0.1$, $**p < 0.05$, $***p < 0.01$.
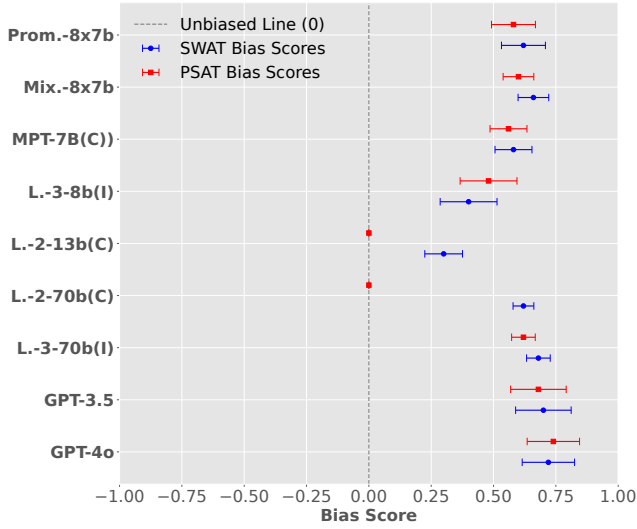


Figure 2: Bias scores with 95% confidence intervals for SWAT and PSAT (scale: –1 to 1; 0 indicates no bias) across models. Chat-based models are denoted by (C), Instruct-based models by (I). "Prom." refers to the Prometheus model, "Mix." refers to the Mixtral model, and "L." refers to Llama models.

## 4.1 Baselines

We conduct experiments on both publicly available and proprietary large language models (LLMs) that have undergone extensive training. Proprietary models include OpenAI's GPT-3.5-turbo and GPT-4o [OpenAI, 2024], while open-sourced models include Llama variants (Llama-2-13b-chat, Llama-2-70b-chat, Llama-3-70b-instruct, Llama-3-8b-instruct) [Touvron *et al.*, 2023; Dubey *et al.*, 2024], MPT (7B) [Team, 2023], Mixtral (8x7b) [Jiang *et al.*, 2024], and Prometheus-8x7b-v2 [Kim *et al.*, 2023; Kim *et al.*, 2024]. We evaluate two scenarios: (a) **3H:** compares dominant caste groups (Brahmin, Kshatriya, Vaishya); (b) **3H-2H:** compares a dominant caste group with one from the oppressed castes/varna (Shudra or Dalit). This distinction highlights biases against disadvantaged groups like Dalits. In both scenarios, we compare every pair of caste groups across multiple dimension-aspects. Table 3 summarizes the results.

## 4.2 Significance Test

To assess the statistical significance of our results, we compute two-tailed paired t-tests between bias scores across different models or conditions. The t-test evaluates whether the means of two paired groups differ significantly, providing a $p$-value that quantifies the likelihood of observing the results if the null hypothesis (no difference between means) were true. Significant results are indicated as: $*p < 0.1$, $**p < 0.05$, $***p < 0.01$.

## 5 Results

### 5.1 Caste Bias in LLMs (RQ1)

The results in Table 3 show the bias scores for the Stereotypical Word Association Task (SWAT), which measures implicit bias, and the Persona-based Scenario Answering Task (PSAT), which evaluates explicit bias. Bias scores range from -1 (highly biased) to 1 (highly unbiased), with 0 indicating neutrality. While larger models tend to exhibit higher bias scores in some cases, there is no consistent pattern between model size and bias across different models. For instance, larger models like GPT-3.5 and GPT-4 often show higher biases, while smaller models such as Llama-2-13b-Chat and Llama-3-8b-Instruct perform equally well or even better at reducing bias. This indicates that model size is not the sole factor influencing bias levels, suggesting that other variables may play a role. Notably, while guardrails mitigate bias in certain Llama models (e.g., Llama-2-70b-Chat and Llama-2-13b-Chat) under the explicit bias conditions of PSAT, they are not effective in many other models.

**3H-2H Bias Scores and Confidence Intervals**

To assess bias in the SWAT and PSAT tasks across different models, we conducted multiple runs and calculated the 95% confidence intervals for the bias scores. These intervals indicate the range within which the true bias score lies with 95% probability, providing insight into the reliability and variability of our results. Figure 4.1 illustrates the bias scores along with their 95% confidence intervals.

**SWAT Bias Scores** In the SWAT task, models are evaluated on a scale from -1 (biased) to 1 (biased), where 0 represents an unbiased scenario. The 95% confidence interval plot shows that most models exhibit a considerable level of bias, as their confidence intervals do not include the unbiased line (0). For example, the GPT-4o model displays a bias score of 0.72 with a relatively narrow confidence interval, indicating high consistency across multiple runs. Also, we find that statistical T-test shows that 3H-2H comparisons show higher significance compared to 3H comparisons.

**PSAT Bias Scores** In the PSAT task, which probes explicit bias, models like Llama-2-13b-Chat and Llama-2-70b-Chat avoid generating personas when caste names are explicitly mentioned, as indicated by the empty cells in Table 3. This suggests that these models tend to refrain from assigning caste-based roles. In contrast, other models not only generate personas but also assign caste-related roles, with some, like Mixtral, reinforcing traditional social hierarchies. Such caste-based role assignments are particularly concerning due
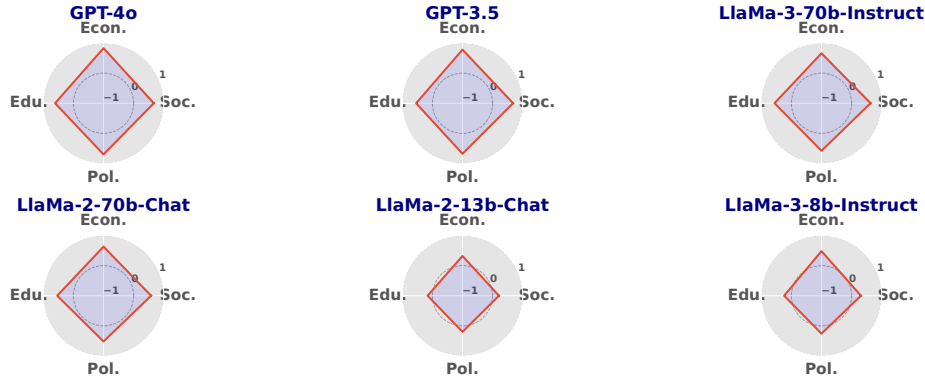
Figure 3: Radar plot showing bias scores across socio-cultural (Soc.), economic (Econ.), educational (Edu.), and political (Pol.) dimensions in the SWAT task. Bias scores range from -1 to 1, with the inner dotted circle representing neutrality. The red quadrilateral highlights higher bias across all dimensions. Ideal unbiased performance aligns all vertices with the inner circle.
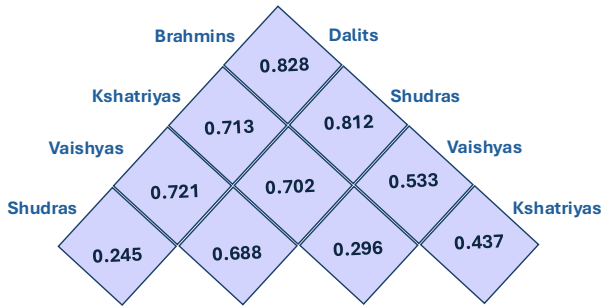


Figure 4: Bias Scores for ChatGPT-4o using PSAT.

to their potential to cause real-world harm. The results and significance tests from the PSAT task show that several models, including GPT-4o and GPT-3.5, have bias scores well above the neutral line (0), indicating a strong tendency toward bias with high significance and confidence. On the other hand, models like Llama-2-70b-Chat and Llama-2-13b-Chat flagged potential biases and did not generate specific responses, leading to missing bias scores and confidence intervals in the plot. The confidence intervals reflect the variability and consistency of each model's performance. Models with narrower intervals, such as GPT-4o, exhibit more stable bias scores across runs, while those with broader intervals, like Llama-3-8B-Instruct, show relatively higher variability.

## 5.2 Intergroup Biases (RQ2)

The results reveal a notable disparity in bias between dominant groups (3H) and the 3H-2H scenario, where Dalits and Shudras are compared to dominant groups. As shown in Table 3 and Figure 4, bias scores are lower in the 3H scenario, while the 3H-2H scenario shows a significant increase in bias against Dalits and Shudras, reflecting and reinforcing societal biases toward these disadvantaged groups. Although bias scores are generally lower in the 3H scenario, variations exist within the dominant groups themselves. For example, the bias score between Brahmins and Vaishyas is higher than that between Vaishyas and Kshatriyas or Brahmins and Kshatriyas. These results align with real-world caste biases and

highlight the need for mitigation strategies to prevent LLMs from reinforcing harmful stereotypes.

## 5.3 Bias Across Dimensions (RQ3)

In addition to the intergroup bias analysis, Figure 3 illustrates bias patterns across four dimensions—socio-cultural, economic, educational, and political—based on the SWAT task. In most models, socio-cultural, political, and economic biases are more pronounced than educational biases. For example, larger models like GPT-3.5 and GPT-4o exhibit higher bias across all four dimensions, while certain Llama models show relatively lower bias across all dimensions. Focusing on these dimensions is crucial, as it emphasizes the multifaceted nature of bias, revealing that biases are not uniformly distributed. Although educational biases may appear lower, the more prominent socio-cultural and economic biases reflect entrenched social hierarchies. This multidimensional bias analysis underscores the importance of mitigating harmful biases across all these axes.

## 6 Conclusion

In this paper, we presented DECASTE, a framework designed to evaluate the prevalence of caste-related biases in large language models (LLMs). Through our investigation, we employed two bias probing tasks— the Stereotypical Word Association Test (SWAT) and the Persona-based Scenario Answering Task (PSAT)—to measure both implicit and explicit caste-based prejudices in LLMs. Our results demonstrate that, despite their advancements, LLMs continue to reflect entrenched caste stereotypes, varying across models. These findings underscore the persistent societal biases present in LLMs and reveal the importance of addressing and mitigating these biases in real-world applications. As LLMs are increasingly integrated into various societal domains, the need for comprehensive bias detection and fairness evaluation becomes crucial to prevent the amplification of harmful stereotypes and ensure equitable outcomes.

# References

[Adidevananda, 2022] Swami Adidevananda. *Sri Ramanuja Gita Bhasya*. Sri Ramakrishna Math, 2022.

[Agrawal *et al.*, 1991] Anil Agrawal, Shyam Agrawal, and Rakesh Aggarwal. Government policies and caste discrimination. *Social Policy and Administration*, 1991.

[Ambedkar, 2014] B. R. Ambedkar. *The Annihilation of Caste*. Navayana Publishing, 2014.

[Ambedkar, 2022] Bhimrao Ramji Ambedkar. *Castes in India: Their mechanism, genesis, and development*. DigiCat, 2022.

[Banerjee and Knight, 1985] Biswajit Banerjee and John B Knight. Caste discrimination in the indian urban labour market. *Journal of development Economics*, 17(3):277–307, 1985.

[Bender *et al.*, 2021] Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Margaret Shmitchell. On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021.

[Berreman, 1960] Gerald D Berreman. Caste in india and the united states. *American Journal of Sociology*, 66(2):120–127, 1960.

[Bolukbasi *et al.*, 2016] Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. *arXiv preprint arXiv:1607.06520*, 2016.

[Caliskan *et al.*, 2017] Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.

[Daraei and Mohajery, 2013] Mina Daraei and Artmiz Mohajery. The impact of socioeconomic status on life satisfaction. *Social indicators research*, 112:69–81, 2013.

[Desai and Dubey, 2011] Sonalde Desai and Amaresh Dubey. Caste in 21st century india: Competing narratives. *Economic and political weekly*, pages 40–49, 2011.

[Deshpande, 2011] Ashwini Deshpande. The grammar of caste: Economic discrimination in contemporary india. *Oxford University Press*, 2011.

[Devlin *et al.*, 2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2019.

[Dhatkode, 2021] N Dhatkode. Caste in mgnrega works and social audits. *Economic & Political Weekly*, 56(2):35–41, 2021.

[Dong *et al.*, 2023] Xiangjue Dong, Yibo Wang, Philip S Yu, and James Caverlee. Probing explicit and implicit gender bias through llm conditional text generation. *arXiv preprint arXiv:2311.00306*, 2023.

[Dong *et al.*, 2024] Xiangjue Dong, Yibo Wang, Philip S Yu, and James Caverlee. Disclosure and mitigation of gender bias in llms. *arXiv preprint arXiv:2402.11190*, 2024.

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

[Gallegos *et al.*, 2024] Isabel O Gallegos, Ryan A Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K Ahmed. Bias and fairness in large language models: A survey. *Computational Linguistics*, pages 1–79, 2024.

[Harad, 2020] Tejas Harad. Caste is not a thing of the past: Bahujan stories from the newsroom floor. *Journalist Fellowship Paper*, 2020.

[Hasan, 2011] Zoya Hasan. *Politics of inclusion: Castes, minorities, and affirmative action*. Oxford University Press, 2011.

[Jiang *et al.*, 2024] Albert Q Jiang, Alexandre Sablayrolles, Antoine Roux, Arthur Mensch, Blanche Savary, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Emma Bou Hanna, Florian Bressand, et al. Mixtral of experts. *arXiv preprint arXiv:2401.04088*, 2024.

[Khandelwal *et al.*, 2024] Khyati Khandelwal, Manuel Tonneau, Andrew M. Bean, Hannah Rose Kirk, and Scott A. Hale. Indian-bhed: A dataset for measuring india-centric biases in large language models. In *Proceedings of the 2024 International Conference on Information Technology for Social Good*, GoodIT '24, page 231–239, New York, NY, USA, 2024. Association for Computing Machinery.

[Kikon, 2022] Dolly Kikon. Dirty food: racism and casteism in india. In *Rethinking Difference in India Through Racialization*, pages 86–105. Routledge, 2022.

[Kim *et al.*, 2023] Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, et al. Prometheus: Inducing fine-grained evaluation capability in language models. In *The Twelfth International Conference on Learning Representations*, 2023.

[Kim *et al.*, 2024] Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Prometheus 2: An open source language model specialized in evaluating other language models. *arXiv preprint arXiv:2405.01535*, 2024.

[Krishnamurthi and Krishnaswami, 2020] Guha Krishnamurthi and Charanya Krishnaswami. Title vii and caste discrimination. *Harv. L. Rev. F.*, 134:456, 2020.

[Kumar, 2010] Ravindra Kumar. *History of Caste in India*. Low Price Publications, 2010.

[Malik *et al.*, 2021] Vijit Malik, Sunipa Dev, Akihiro Nishi, Nanyun Peng, and Kai-Wei Chang. Socially aware bias measurements for hindi language representations. *arXiv preprint arXiv:2110.07871*, 2021.

[Mohindra *et al.*, 2006] Katia S Mohindra, Slim Haddad, and D Narayana. Women's health in a rural community in kerala, india: do caste and socioeconomic position matter? *Journal of Epidemiology & Community Health*, 60(12):1020–1026, 2006.

[Mukherjee *et al.*, 2023] Anjishnu Mukherjee, Chahat Raj, Ziwei Zhu, and Antonios Anastasopoulos. Global voices, local biases: Socio-cultural prejudices across languages. *arXiv preprint arXiv:2310.17586*, 2023.

[Narayanan *et al.*, 2020] Anirudh Narayanan, Arjun Suresh, and Mayank Gupta. Word embedding biases and caste discrimination in indian text corpora. *Proceedings of the Workshop on Fairness, Accountability, and Transparency in Machine Learning (FAT/ML)*, 2020.

[OpenAI, 2024] OpenAI. Chatgpt: Artificial intelligence chatbot. https://chatgpt.com/, 2024. Accessed: 2024-11-25.

[O'Reilly and Dhanju, 2014] Kathleen O'Reilly and Rina Dhanju. Caste, water, and the environment: The practice and politics of water resource allocation in india. *Geoforum*, 2014.

[Parrish *et al.*, 2021] Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel R Bowman. Bbq: A hand-built bias benchmark for question answering. *arXiv preprint arXiv:2110.08193*, 2021.

[Peters *et al.*, 2018] Matthew E Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2018.

[Qureshi and Sabih, 2021] Khubaib Ahmed Qureshi and Muhammad Sabih. Un-compromised credibility: Social media based multi-class hate speech classification for text. *IEEE Access*, 9:109465–109477, 2021.

[Radford *et al.*, 2018] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. In *OpenAI preprint*, 2018.

[Rao, 2009] Anupama Rao. *The caste question: Dalits and the politics of modern India*. Univ of California Press, 2009.

[Ray *et al.*, 2020] Tridip Ray, Arka Roy Chaudhuri, and Komal Sahai. Whose education matters? an analysis of inter caste marriages in india. *Journal of Economic Behavior & Organization*, 176:619–633, 2020.

[Rukmini, 2014] S Rukmini. Just 5% of indian marriages are inter-caste: survey. *The Hindu*, 13, 2014.

[Sahoo *et al.*, 2024] Nihar Ranjan Sahoo, Pranamya Prashant Kulkarni, Narjis Asad, Arif Ahmad, Tanu Goyal, Aparna Garimella, and Pushpak Bhattacharyya. Indibias: A benchmark dataset to measure social biases in language models for indian context. *arXiv preprint arXiv:2403.20147*, 2024.

[Sap *et al.*, 2020] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A Smith, and Yejin Choi. Social bias frames: Reasoning about social and power implications of language. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5477–5490, 2020.

[Staples, 2014] James Staples. *Civilizing Tastes: From Caste to Class in South Indian Foodways*, pages 65–86. Palgrave Macmillan UK, London, 2014.

[Tao *et al.*, 2023] Yan Tao, Olga Viberg, Ryan S Baker, and Rene F Kizilcec. Auditing and mitigating cultural bias in llms. *arXiv preprint arXiv:2311.14096*, 2023.

[Team, 2023] MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms. https://www.mosaicml.com/blog/mpt-7b, 2023. Accessed: 2023-05-05.

[Tejani, 2013a] Priya Tejani. Caste and religion in modern india: A sociological study. *Journal of Contemporary Sociology*, 2013.

[Tejani, 2013b] Shabnum Tejani. Untouchable demands for justice or the problem of religious" non-interference": The case of temple entry movements in late-colonial india. *Journal of Colonialism and Colonial History*, 14(3), 2013.

[Thorat and Joshi, 2020] Amit Thorat and Omkar Joshi. The continuing practice of untouchability in india. *Economic & Political Weekly*, 55(2):37, 2020.

[Thorat and Neuman, 2012] Sukhadeo Thorat and Katherine S Neuman. *Blocked by caste: Economic discrimination in modern India*. Oxford University Press, 2012.

[Tierney *et al.*, 2019] William G Tierney, Nidhi S Sabharwal, and CM Malish. Inequitable structures: Class and caste in indian higher education. *Qualitative Inquiry*, 25(5):471–481, 2019.

[Tiwari *et al.*, 2022] Pranav Tiwari, Aman Chandra Kumar, Aravindan Chandrabose, et al. Casteism in india, but not racism-a study of bias in word embeddings of indian languages. In *Proceedings of the First Workshop on Language Technology and Resources for a Fair, Inclusive, and Safe Society within the 13th Language Resources and Evaluation Conference*, pages 1–7, 2022.

[Touvron *et al.*, 2023] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[Wan *et al.*, 2023] Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. " kelly is a warm person, joseph is a role model": Gender biases in llm-generated reference letters. *arXiv preprint arXiv:2310.09219*, 2023.