# Detecting Linguistic Indicators for Stereotype Assessment with Large Language Models

Rebekka Görge
Fraunhofer Institute for Intelligent
Analysis and Information Systems
Sankt Augustin, Germany
Lamarr-Institute for Machine
Learning and Artificial Intelligence
Bonn, Germany
rebekka.goerge@iais.fraunhofer.de

Michael Mock
Fraunhofer Institute for Intelligent
Analysis and Information Systems
St. Augustin, Germany
michael.mock@iais.fraunhofer.de

Héctor Allende-Cid
Fraunhofer Institute for Intelligent
Analysis and Information Systems
St. Augustin, Germany
Lamarr-Institute for Machine
Learning and Artificial Intelligence
Bonn, Germany
hector.allende-cid@iais.fraunhofer.de

## Abstract

Social categories and stereotypes embedded in language can introduce data bias into the training of Large Language Models (LLMs). Despite safeguards, these biases often persist in model behavior, potentially leading to representational harm in outputs. While sociolinguistic research provides valuable insights into the formation and spread of stereotypes, NLP approaches for bias evaluation rarely draw on this foundation and often lack objectivity, precision, and interpretability. To fill this gap, we propose a new approach to assess stereotypes by detecting and quantifying the linguistic indication of a stereotype. We derive linguistic indicators from the Social Category and Stereotype Communication (SCSC) framework indicating strong social category formulation and stereotyping in language, and use them to build a categorization scheme. We use in-context learning to instruct LLMs to examine the linguistic properties of a sentence containing stereotypes, providing a basis for a fine-grained stereotype assessment. We develop a scoring function to measure linguistic indicators of stereotypes based on empirical evaluation. Our annotations of stereotyped sentences reveal that these linguistic indicators explain the strength of a stereotype. The models perform well in detecting and classifying linguistic indicators used to denote a category, but sometimes struggle with accurately evaluating the described associations. The use of more few-shot examples significantly improves the performance. Model performance increases with size, as `Llama-3.3-70B-Instruct` and `GPT-4` achieve comparable results that surpass those of `Mixtral-8x7B-Instruct`, `GPT-4-mini` and `Llama-3.1-8B-Instruct_4bit`. Code and annotations can be found in https://github.com/r-goerge/Detecting-Linguistic-Indicators-for-Stereotype-Assessment-with-LLMs.

## CCS Concepts

• **Computing methodologies** → **Natural language processing**.

## Keywords

Large language models, fairness, stereotype detection, linguistics

## 1 Introduction

*Content Warning: This paper presents textual examples that may be offensive or upsetting.*

Language reflects and transmits the social categories and stereotypes that humans use to quickly perceive and interact within their complex social environment. *Stereotypes* are defined as "cognitive representation people hold about a social category consisting of beliefs and expectancies about their probable behavior, feature and traits" [13]. Although social categorization is a natural human tendency [4], it often leads to exaggerated group differences and oversimplified perceptions of groups. Harmful consequences, such as discrimination or unfair decisions, arise when individuals are judged based on broad social category associations rather than their unique traits. This might even be enforced by the use of negative stereotypes, which are denoted as prejudice[1]. Encoded in human language, also large language models (LLMs) trained on massive amount of aggregated and crawled text data are learning, reproducing and disseminating stereotypes [33]. This perpetuation of stereotypes and its potential harmful impact on society and individuals is a significant concern regarding AI [28, 38].

Current research investigates in the detection and mitigation of stereotypes as part of the research on fairness and bias[2] of AI models [19]. Within the context of bias, the presence of stereotypes in training data constitutes a form of data bias[3]. Likewise, the reproduction of stereotypes by AI models is seen as a harmful consequence of bias, often referred to as representational harm [19]. To mitigate representational harm, state-of-the-art (SOTA) LLMs are equipped with guardrails to prevent stereotypical output. Although these measures are often effective against explicit reproduction of stereotypes, they often fail when confronted with slight variations in the prompting

---

[1]"Negative affective evaluations of social categories and members" are denoted as prejudice [4, 41]
[2]Definition according to [23]: "Systematic difference in treatment of certain objects, people, or groups in comparison to others"
[3]Definition according to [23]: "Data properties that, if unaddressed, lead to AI systems that perform better or worse for different groups"

[44]. This indicates that intrinsic biases continue to manifest in model behavior, highlighting the need to reduce stereotypes as much as possible at the data level. Also sociolinguistics has long examined the emergence and linguistic form of stereotypes, Blodgett et al. finds that most research on bias in natural language processing (NLP) is poorly aligned with interdisciplinary studies. Existing work in stereotype detection primarily relies on human-constructed and annotated benchmarking datasets [3, 16, 31, 32], which serve as the foundation for training classifiers for text-based stereotype detection and for benchmarking widely used LLMs. These datasets and resulting detection methods typically categorize sentences as either stereotypical or anti-stereotypical based on subjective human assessments, leading to pitfalls such as misalignment or uncleanness of the stereotypes themselves [5]. To address this issue, Liu argues that a purely binary separation of stereotypes is insufficient, as stereotypes can be expressed in various forms. Instead, they propose a fine-grained quantification of the strength of a stereotype using a continuous scoring system grounded in human stereotype rankings [26]. However, as the score relies on human judgment, it reflects subjective perceptions shaped by the annotators' cultural and social backgrounds, and it lacks an explanation for why a stereotypical sentence contains a specific stereotype.

Our work starts at this point, as this paper presents a novel, sociolinguistically based, twofold approach to quantifying stereotypes in language. Unlike [26], we do not assess the harmfulness of a stereotype based on human judgment, but instead focus on evaluating potentially stereotypical sentences by first detecting linguistic indicators that signal a stereotype, and second quantifying them to assess the strength of the stereotype. Our approach is illustrated in Figure 1. Note that our current approach assumes that a given sentence contains a stereotype and focuses on stereotype assessment, but is extendable to stereotype detection as well, as outlined in the conclusion. By grounding our work in sociolinguistics, we adopt the Social Category and Stereotype Communication Framework (SCSC) [4] (Section 3), which explicates the linguistic processes by which stereotypes are shared and maintained in language, and use it to derive a clear categorization scheme with a fixed set of linguistic indicators (Section 4). Leveraging their extensive linguistic capabilities, we integrate LLMs to automatically detect linguistic indicators by guiding them through the categorization scheme using an in-context learning approach (Section 5). To quantify the linguistic indication on a stereotype, we aggregate the derived linguistic indicators by learning a scoring function, for which we exploit the work of Liu and learn the importance of different linguistic indicators based on human stereotype rankings (Section 6). We validate our approach using a manually annotated subsample of CrowS-Pairs[32] and human-based stereotype rankings of [26].

## 2 Related work

Research on stereotyping in NLP is an integral part of the wider investigation of fairness and bias in AI. According to Hovy and Prabhumoye, the key sources of bias in NLP are the training data and the models themselves. While biased data reflect societal stereotypes and inequities, model architectures and training processes can amplify or introduce additional biases, compounding the problem. Stereotypes in language systems intersect with social hierarchies

and human cognition, highlighting the need to engage with foundational literature from psychology, sociology, and sociolinguistics to understand stereotype formation, its harms, and guide the development of more equitable NLP methodologies. While few works in NLP build on this foundation [5], numerous studies have shown that bias, including stereotyping, is a problem in NLP [7, 9, 44]. This includes more recently pre-trained language models, especially masked language models [2, 32, 40], which, while showing considerable success in various NLP tasks, also inherit and perpetuate cultural biases embedded in the training corpora, leading to harm through biased representations. Therefore, most recent studies have focused on identifying stereotypes in model output as a consequence of model bias [19].

To this end, several datasets such as the Crowdsourced Stereotype Pairs (CrowS-Pairs) benchmark [32], StereoSet [31], Multi-Grain Stereotype [45] and FairPrism [16] have been introduced comprising stereotypes across different social categories. StereoSet and CrowS-Pairs in particular are widely used, but inherit significant weaknesses regarding the construction of stereotypes[6]. This is why [34] presents an improved version of CrowS-Pairs. Both datasets also contain artificially constructed anti-stereotypes, which are unlikely to occur in regular discourse [35].

In contrast to stereotypes as a harm of model bias, the exploration of stereotypes as a form of data bias introduced through human language bias is less explored [18]. Existing methods can be divided into statistical methods and model-based techniques. Statistical approaches such as embedding-based metrics [19] focus on analyzing the distribution and co-occurrence patterns of words, phrases, or demographic categories within datasets. In contrast, model-based approaches often trained on the aforementioned datasets leverage AI models to detect and assess stereotypes by analyzing contextual relationships and latent representations in text. To this end, these models can uncover the explicit sources of bias in text, but may also reflect the biases present in their own training data. These approaches can also be used as a filter to detect stereotypes in LLM output.

Among the model-based approaches, there are several studies that specifically adapt and train pre-trained language models for stereotype detection [35] or assessment [18, 26, 37], using different methods and evaluating stereotypes according to different aspects. In terms of stereotype detection, Pujari et al. address the subtle manifestations of stereotypes by creating a focused evaluation dataset that includes explicit stereotypes, implicit stereotypes, and non-stereotypes. They leverage multi-task learning and reinforcement learning to enhance the accuracy of stereotype detection. As model-based stereotype detection methods often lack explainability, Wu et al. and King et al. propose explainability tools such as Shap and Lime [36] to explain the decisions of stereotype detectors. In terms of stereotype evaluation, Sap et al. develop and implement Social Bias Frames, a formalism that captures the pragmatic implications of stereotypes, supported by a large annotated corpus that emphasizes the need to combine structured inference with common sense reasoning. Similar, Fraser et al. builds on interdisciplinary work and presents a pre-trained embedding model that models the Stereotype Content Model [15] from psychological research to analyze stereotypes along the dimensions of warmth and competence, leading to an intrinsically interpretable approach. Also, Liu focuses
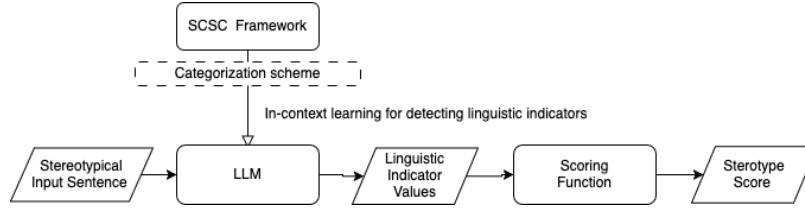
**Figure 1: Overview of our framework for assessing linguistic indicators of stereotypes using the SCSC framework and LLMs**

on stereotype evaluation and advocates fine-grained evaluation of stereotypes instead of binary recognition by introducing a human ranking of stereotypes and pre-training LLMs on it to investigate correlations between stereotypes and broader social topics.

With recent advances of LLMs, current approaches also propose the use of fine-tuned LLMs [11, 22, 42] to detect stereotypes in text data. Tian et al. and Dige et al. assess zero-shot stereotype detection using reasoning and Chain-of-Thought (CoT) prompts. Tian et al. finds that reasoning plays a pivotal role in enabling LLMs to exceed scaling limitations on out-of-domain tasks such as stereotype detection. However, Dige et al., report limited performance, with Alpaca 7B achieving the highest accuracy of 56.7%, while suggesting that increasing model size and data diversity could lead to further performance gain. This finding is approved by Huang et al., who evaluate the trustworthiness of current LLMs using a prompt-based stereotype detection approach as one of several fairness assessment tools. They find that most LLMs demonstrate unsatisfactory performance in stereotype recognition, with even the best-performing model, GPT-4, only achieving 65% overall accuracy.

Table 1 compares these existing model-based stereotype detection and assessment approaches, including our approach. These works underscore the importance of robust datasets, stereotype detection, fine-grained quantification, and interpretability methods to effectively address stereotypes and bias in NLP. However, we identify a gap when it comes to the need to combine these dimensions into one. To fill this gap, we propose a method grounded in linguistic theory that does not need pre-training or human annotations, but exploits in-context learning capabilities of LLMs, addressing the limitations of the above-mentioned zero-shot approaches. As our approach is grounded on a categorization scheme derived from linguistic theory, it is intrinsically interpretable. By introducing a scoring function, we enable fine-grained stereotype quantification. While primarily focused on stereotype assessment, it can also be extended to pre-filtering and detection using linguistic indicators.

## 3 The social category and stereotype communication framework

Sociolinguistics has long been researching how stereotypes are shared and maintained in language. To explicate these linguistic processes, the Social Categories and Stereotypes Communication Framework [4] integrates different aspects of research on stereotyping and biased language use in one framework. It is based on the theory that stereotypical expectations are reflected in language through subtle, largely implicit linguistic biases, such as minor changes in syntactic and semantic structure. The framework is built on the definition of a stereotype introduced in Section 1, which

consists of an addressed 'social category' and the associated beliefs and expectations about 'behavior and characteristics'. To explain how stereotypes are maintained and disseminate in language, the SCSC framework provides three main aspects: (1) shared cognition of social categories, (2) communicated target situation, and (3) types of bias in language use. Shared cognition of social categories (1) explains, from a cognitive perspective, three factors influencing how people perceive social categories: the extent to which a category is perceived as a meaningful and coherent group (perceived category entitativity), the beliefs and expectations about the likely behaviors, features, and traits associated with that category (stereotype content), and the extent to which these characteristics are seen as immutable for its members and stable across time and situations (perceived category essentialism). The target situation (2) refers to the level of information of the communicated content, which is determined by the generalization of the described content regarding the social target category addressed and the associated situation. The lowest level of information describes the situational behaviors of individuals, the intermediate level captures the enduring characteristics of individuals, and the highest level outlines the enduring characteristics of the group as a whole. While the aggregation of low-level information contributes to stereotype formation, stereotypes are primarily conveyed and maintained at the highest level, making it particularly relevant for stereotype detection.

From particular interest to us is the bias in language use (3), which describes the linguistic biases that are used in language when expressing and maintaining stereotypes[4]. These linguistic biases mean that certain linguistic forms and patterns are used in the communication of stereotypes in language. Notably, these linguistic biases express, on a linguistic level, the factors of shared cognition of social categories (1), as well as the level of information of the target situation (2). Based on the definition of stereotypes, there are linguistic biases referring to the social category and linguistic biases pertaining to the associated content (characteristics and behaviors). We illustrate this with the sample sentence *"Women can't drive in the rain"*. For a detailed list of examples, see Appendix A, 4. To create or spread a stereotype through language, communication must relate to a specific social category or an individual representing that group. This is achieved linguistically by using a label that identifies the targeted social group, referred to as the "category label" (*e.g., Women*). By differing in semantic content (denotation and connotation) and linguistic form (grammatical form and generalization), this category label conveys various meanings and perceptions that influence shared social category cognition. The semantic content

---

[4]Notably, the language use is embedded in a communicative context determined by factors such as social context, medium, or communication partners.

**Table 1: A comparison of model-based methods for stereotype detection and assessment showing current gaps**

|  | No fine-tuning or pre-training | Fine-grained stereotype quantification | Interpretability of stereotype assessment | Stereotype detection |
|---|---|---|---|---|
| Liu [26] |  | x |  |  |
| Fraser et al. [18] |  | (x) | x | (x) |
| Sap et al. [37] |  |  | x |  |
| Tian et al. [42] | x |  |  | x |
| Pujari et al. [35] |  |  |  | x |
| Wu et al. [45] |  |  | (x) | x |
| King et al. [25] |  |  | (x) | x |
| Our approach | x | x | x | (x) |

of the label transports the category boundaries, and hierarchies influencing the perceived category entitativity. Moreover, it might even already activate stereotypical content (*e.g., referring to women using the label ladies versus feminists*). Similarly, the generalization of the label and its grammatical form influence the perceived category entitativity and the level of information. For instance, higher generalization and use of nouns form a stronger entity (*e.g., Women can't drive in the rain*) compared to description of individuals labeled by adjectives (*e.g., She is female and can't drive in the rain*). Once a category is labeled, the meaning and linguistic form of the content (*e.g., can't drive in the rai*n) itself is relevant. The stereotype is transported within the content. To this end, it is important to understand linguistically whether the content refers to situational behavior or enduring characteristics and consistent behaviors about the category label. This relates again to the level of information and the perceived category essentialism. Moreover, the described behaviors or characteristics might be consistent or inconsistent with existing stereotypical information on this category. While this largely depends on culture and world knowledge, the language of the content indicates whether the shared information aligns with the sender's beliefs or expectations. Consequently, linguistic forms vary to signal stereotype-consistent versus inconsistent information, such as, for example, through linguistic abstraction, indication on regularity, the omission of explicit situational explanations, or irony/negation bias, which can even convey stereotypes by presenting stereotype-inconsistent information. Stereotype-consistent information, in turn, affects category essentialism.

## 4 Development of a categorization scheme for linguistic indicators

### 4.1 Categorization scheme

The SCSC framework [4] uses a variety of examples to describe how linguistic biases in language influence the cognitive perception of stereotypes. However, it does not provide a clear and structured scheme that defines all different linguistic biases and their potential forms. Therefore, we derive a categorization scheme based on the SCSC framework, which we extend to assess stereotypes automatically at the sentence level.

We adopt the division proposed by [4] into the two primary categories, *category label* and *associated content*, taking into account both *language meaning* and *linguistic form* for each category. From there, we derive from the information level (2) and linguistic biases

(3) described in Section 3 a set of *n* linguistic indicators that signal stereotypes in communication, which we denote as linguistic indicator $A_i$. We only include linguistic indicators that can be assessed based on linguistics as objectively as possible and without the need for interpretation. For each of these indicators, we define a set of *k* potential values. Depending on its value, most attributes have a strengthening effect ($\uparrow$) or a weakening effect ($\downarrow$) on one of the aspects of shared social category cognition (1). In Section 6.1, we compute weights for the strengthening and weakening effects based on an empirical validation. The aggregation of linguistic indicators reflects the potential linguistic strength of the stereotype.

The categorization scheme is illustrated in Table 2. For both primary categories, we define a linguistic indicator as basic decision criteria to check whether the sentence contains a) a category label (*has category label*) and b) information on associated behavior or characteristics of the category label (*information level (situation)*). Both are prevalent for the existence of stereotypical content. While our approach is applicable for the detection of stereotypes against arbitrary social groups, we restrict the identification of category labels within an analysis to labels that refer to predefined sensitive attributes such as race or gender to ensure a meaningful evaluation [5]. If a category label and, if applicable, associated behaviors and characteristics about a behavior or a property are available, both are further categorized. Regarding the category label, it encompasses the following aspects: At the level of meaning, it includes the content of the *category label*, its *connotation*, and the *information level(target)*, at the level of linguistic form, it involves the *grammatical form* and *generalization (category label)*. In terms of shared information, the meaning encompasses the associated *content*, while the linguistic form includes *generalization (content)*, the use of *explanation of behaviors or characteristics*, and the use of *signal words*. To increase objectivity, especially when automating the process with LLMs, we do not include *irony bias* and *negation bias* as further linguistic indicators of stereotype inconsistency.

### 4.2 Dataset and manual ground truth annotation

To validate the functionality of the categorization scheme and establish a reference ground truth, we manually annotate a subset of the CrowS-Pairs dataset. The annotations can be found in the accompanying git-repository. The CrowS-Pairs benchmark is a widely used dataset that addresses stereotypes across nine types of social bias

**Table 2: Our categorization scheme grounded in the SCSC framework [4]: We define a fixed set of linguistic indicators and values, and their strengthening (↑) or weakening (↓) impact on shared social category cognition (using "entitativity" for category entitativity and "essentialism" for category essentialism) ordered by category label and associated content.**

| Level | Linguistic Indicators $A_i$ | Values | Shared social category cognition | Definition |
|---|---|---|---|---|
| | | | | |
| | | Linguistic indicators of category label | | |
| Language meaning | Has category label towards *<social category>* | Yes | ↑ entitativity | Content pertains to a social category or an representing individual. |
| | | No | ↓ entitativity | Content does not pertain to a social category or individual representing it. |
| | Category label | <Text> | *stereotype content* | Extraction of the category label. |
| | Connotation | negative | ↑ stereotype content | Label connotation is negative (e.g., bitches). |
| | | neutral | | Label connotation is neutral (e.g., women). |
| | | positive | ↑ stereotype content | Label connotation is positive (e.g., female hero). |
| | Information level (target) | generic target | ↑entitativity | Unspecified individual(s) or generic group (e.g., a girl) |
| | | specific target | ↓entitativity | Specified individual(s) (e.g the girl) |
| Linguistic form | Grammatical form | noun | ↑ entitativity, ↑ essentialism | Category is transported through a noun (e.g., Asians) |
| | | other | ↓ entitativity, ↓ essentialism | Category is conveyed by another form such as an adjective/ proper noun (e.g., he is black) |
| | Generalization (category label) | generic | ↑ entitativity, ↓ essentialism | Reference to a demographic group as whole. |
| | | subset | | Reference to a subset of a demographic group. |
| | | individual | ↓ entitativity, ↓ essentialism | Reference to a specific or unspecific individual person, who may be a member of a demographic group |
| | | Linguistic indicators of associated behaviors and characteristics | | |
| Language meaning | Information level (situation) | situational behavior | ↑essentialism | Behavior observed in a specific situation. |
| | | enduring characteristics | ↓ essentialism | Enduring characteristics or traits not observable in a single situation. |
| | | other | | Describes other information on category label. |
| | Assoc. content | <Text> | *stereotype content* | If not other, extraction of content. |
| Linguistic form | Generalization (content) | abstract | ↑ essentialism | Use of abstract terms such as state verbs or adjectives (e.g., she is aggressive [4]) |
| | | concrete | ↓ essentialism | Use of concrete terms such as active or descriptive verbs (e.g., she kicks him [4]) |
| | Explanation for behaviors, characteristics | yes | ↓ essentialism | Explanation is given more frequently for stereotype-inconsistent behavior. |
| | | no | ↑ essentialism | Stereotype-consistent behavior is expected and less often paired with an explanation. |
| | Signal words | typical | ↑ essentialism | Signal words for typicality indicate stereotype-consistency |
| | | exceptional | ↓ essentialism | Signal words for exceptionality indicate stereotype-inconsistency. |
| | | none | | No signal words are used. |

(such as race, gender, religion, or physical appearance). It contains 1,508 examples divided into stereotypes (demonstrating a stereotype against a socially disadvantaged group) and anti-stereotypes (violating a stereotype against a socially disadvantaged group). Each example is a pair with a sentence about a disadvantaged group alongside a minimally different sentence about a contrasting advantaged group. The sentences were obtained through crowdsourcing with Amazon Mechanical Turk. Although CrowS-Pairs is widely used, Blodgett et al. reveals serious shortcomings in the conceptualization and operationalization of the dataset. Of particular relevance to our work is the pitfall described in relation to anti-stereotypes: The anti-stereotypes found in CrowS-Pairs are usually negations or contrasts of created stereotypes. In some cases, true statements are found, but in many cases irrelevant statements are made about a target group [6]. Due to the artificial construction of the anti-stereotypes, they do not reflect the linguistic patterns

for formulating stereotype-inconsistent statements as described in [4], but mainly change the content of the statement about a group [35]. While this might be less critical for benchmarking the preferences of a language model, it cannot be equated with naturally occurring stereotype-inconsistent or stereotype-free sentences in the language.

To overcome this problem, we select only stereotypical sentences from CrowS-Pairs, which is sufficient for the development of our scoring function. To minimize the risks associated with CrowS-Pairs, we make use of the work of Névéol et al. and only use these sentence of CrowS-pairs that were confirmed by Névéol et al. in the revised dataset. Next, we select 100 sentences from CrowS-Pairs targeting either the attribute gender or race. Then each sentence is manually analyzed for linguistic indicators which are then categorized according to the attributes developed from Table 2 by two annotators (one of the authors and a research collaborator) with different cultural backgrounds separately. The annotators received step-by-step instructions for the annotations, reflecting the information provided in the LLM prompts (compare Appendix A), only adapted in format. After an initial annotation step, Cohen's kappa [10] was calculated for the annotators to be $\kappa = 96\%$, indicating almost perfect agreement. The majority of deviation occurred with $\kappa = 93\%$ in the linguistic indicator *generalization (content)*. To reach a ground truth, both annotators discussed deviation until they have reached agreement. An example of some annotated sample sentences is given in Table 4 in the appendix.

## 5 Model development and validation

The next step is to automate the categorization scheme that has been developed to analyze linguistic indicators of stereotypes at the sentence level. To implement the categorization scheme, the linguistic indicators $A_i$ must be detected and classified. Therefore, it is necessary to solve several classical NLP tasks such as relation extraction or sentiment analysis. For this purpose, we use instruction-finetuned LLMs with an in-context learning (ICL) approach using LLMs as judges, leveraging their strong overall text comprehension capabilities leading to SOTA performance in classical NLP tasks [43]. In-context learning [12] allows LLMs to perform tasks without parameter updates by leveraging contextual information at inference time, making it useful in scenarios with scarce labeled data. Its effectiveness relies on the model's scale and pretraining, with larger models exhibiting superior ICL capabilities. In particular, we exploit few-shot learning [8], which improves model performance by providing labeled examples within in-context learning [30]. The "LLM-as-a-Judge" approach [20] employs LLMs as evaluators for complex tasks, addressing subjectivity and variability in traditional evaluations while remaining adaptable to various sensitive attributes without the need for training or fine-tuning. To this end, we follow a prompt-based approach that should be effectively applicable to a wide range of sentences. To cover different sizes and architectures, we include `Llama-3.1.-8B-Instruct` and `Llama-3.3-70B-Instruct` from the Llama family [14], `Mixtral-8x7B-Instruct` from Mistral AI [24], `GPT-4o-mini` and `GPT-4` [1] of OpenAI. We first describe how the categorization scheme is implemented through prompt engineering. From this, we validate the
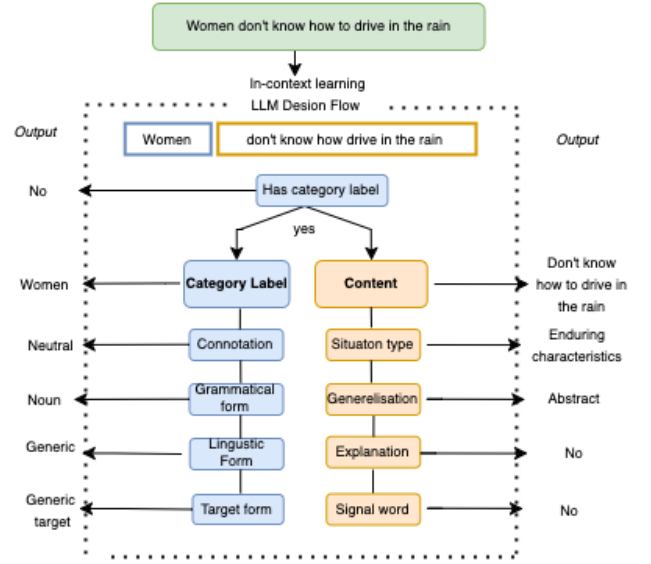


**Figure 2: Illustration of the LLM's decision process on a sentence to detect linguistic indicators and classify their values**

model performance of various LLMs in detecting and classifying linguistic indicators.

### 5.1 Prompt engineering

In terms of prompt engineering, we manually create and adapt a prompt template adhering to best practice prompt engineering strategies. We seek to optimize the prompt through iterative refinement, using both Llama-models as reference points. The final prompt for the fine-grained evaluation of linguistic indicators can be found in Appendix A.

We construct a few-shot prompt that includes a role description, a task description, and several examples that cover different scenarios. The task description describes the categorization scheme introduced in Section 4 and is designed to guide the model through a decision flow to detect the linguistic indicators $A_i$ and classify its correct value, as illustrated in Figure 2. Following the annotator's guidelines, we formulate a basic prompt by exploring different formulations, which includes the definition of each attribute and an example for each value. This basic prompt incorporates one or several interchangeable sensitive attributes for which stereotypes are to be identified (e.g., *Your task is to identify, in a given sentence, a category label referring to* `<sensitive attributes>`). To facilitate the task, we integrate COT components, prompting the model to extract and repeat relevant content before addressing questions about the category label and associated content (e.g., *Extract the exact information shared about the category label*).

To determine the best number of examples to be used in the few-shot learning, we compare the models' mean performance over the attributes related to the category label and related to the associated content by increasing the number of few-shot examples to nine. As shown in Figure 3, we find that by increasing the number of
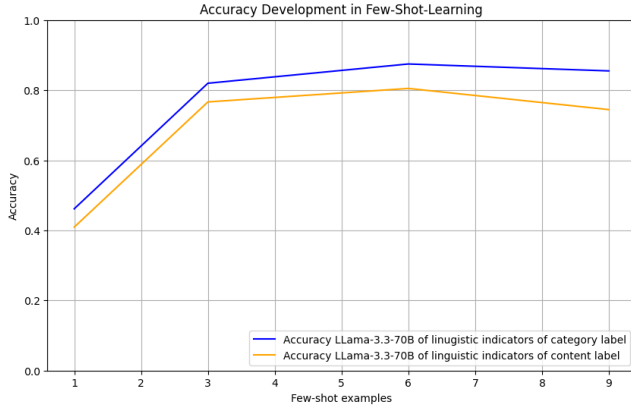
**Figure 3: Accuracy of linguistic indicators based on number of few-shot examples**

examples in several steps from one to six the accuracy strongly increases. For the last three examples, the performance decreases slightly, which we assume is due to the fact that they cover corner cases (compare Appendix A). Nevertheless, we keep these examples because they are important for a broad application.

We compare single-stage and multi-stage prompt formulations. The single-stage approach uses one prompt for the entire categorization scheme, while the multi-stage approach divides into two sub-tasks: the first prompt focuses on the category label, followed by a second prompt that addresses the associated content using the category label found as context. While `Llama-3.3-70B-Instruct` achieves similar performance in both stages, the performance of `Llama-3.1.-8B-Instruct` with respect to the linguistic properties of the content drops drastically in the multi-stage approach. This is consistent with the findings of [39], so we continue with the single-stage approach. Similarly to [42], to facilitate deterministic answers, we ask the model to give answers related to categorization in a structured way. Here, we find that asking for structured JSON-format {"has_category_label": "yes", "full_label": "these English gentlemen"} encourages both models to quote only the requested output in the exact scheme which is not the case by using just a numbering ((1)yes, (2)these English gentlemen).

## 5.2 Model validation

Using this prompt, we evaluate the performance of `Llama-3.1.-8BInstruct`, `Llama-3.3-70B-Instruct`, `GPT-4`, `GPT-4o-mini` and `Mixtral-8x7B-Instruct` in the detection and classification of linguistic indicators using a temperature of $t = 0.7$. All experiments are conducted on a sample of CrowS-Pairs, utilizing our existing human annotations as ground truth. After running the models, we post-process the outputs to extract the linguistic indicators from the JSON outputs. However, `GPT-4o-mini` and `Mixtral-8x7B-Instruct` do not consistently adhere to the JSON format, adding extra spaces or backslashes, which we remove in post-processing.

We evaluate each model's performance on the linguistic indicators using accuracy and F1-Score, performing multi-class classification as each indicator must be detected and categorized. Detailed
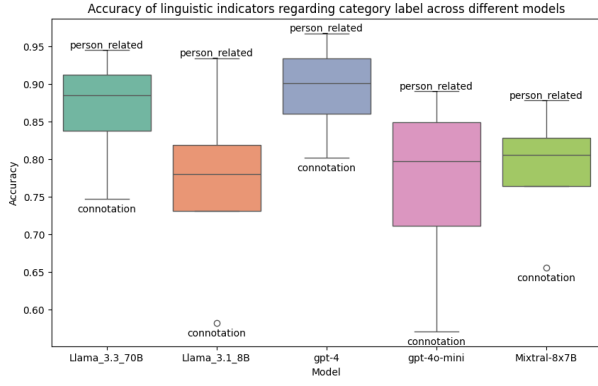
results for each linguistic indicator can be found in Appendix C. We summarize the results in Figure 4, illustrating the accuracy of indicators related to the *category label* in Figure 4a and those related to the *associated content* in Figure 4b. Performance scores are averaged across individual class performances. Overall, larger models like `Llama-3.3-70B` and `GPT-4` outperform smaller models and `Mixtral-8x7B-Instruct` across all linguistic properties, though distinct trends emerge in all models: All models perform significantly better in detecting and classifying linguistic indicators of the *category label* compared to those of the *associated content*. Indicators such as *has category label*, the *information level*, and the *grammatical form* of the category label are generally well recognized due to their clarity and unambiguity. A closer examination of the F1-score of *generalization (target)* shows that while the values *generic* and *individual* are detected effectively, the intermediate label *subset* challenges the models. Furthermore, almost all models struggle to accurately identify the *connotation* of the label, which is unexpected given their strong performance in sentiment recognition. Further analysis indicates that the models often include the sentiment of the entire sentence instead of focusing solely on the category label. In terms of performance in content properties, the gap between large and small models widens. While larger models maintain acceptable performance, smaller models experience a significant decline. A noticeable trend across the models is the difficulty in detecting *generalization (content)*, as the models struggle to focus solely on the verbs and adjectives used in the sentences. For other attributes, there is greater variability in the models' performance.

In general, the evaluation confirms the potential of LLMs to automatically realize the categorization scheme. However, only larger models such as `Llama-3.3-70B` and `GPT-4` are able to provide consistent performance across attributes. Challenges remain in accurately interpreting attributes such as *connotation* and *generalization*, suggesting areas for potential improvement in future work. For further experiments, we continue to use `Llama-3.3-70B`, which achieves a mean accuracy of 81%, comparable to `GPT-4`'s 83%. However, `Llama-3.3-70B` is open source, which makes the results easier to reproduce. With these results on LLM performance in detecting and categorizing linguistic indicators of a stereotype, we can already improve existing stereotype detection mechanisms by providing annotations of the stereotype in terms of an approved sociolinguistic framework. These annotations, which can serve as human-understandable explanations, can facilitate fine-grained analysis of a given text database or as a further vehicle for analyzing the LLM output. We present an example on this in Appendix B. In the following, we focus on using these annotations to generate a numerical score that quantifies the strength of a stereotype, in the sense of the work presented by Liu[26].

## 6 Assessing linguistic stereotype indicators

## 6.1 Weighting of linguistic indicators and scoring function

Using the categorization scheme and in-context learning, we can effectively derive the linguistic indicators from the experiments performed. As can be seen in Table 2, each linguistic indicator has a strengthening or weakening effect on category entitativity, stereotype content, and category essentialism, and thus on how we

**(a) Performance on the linguistic indicators of the category label (*person related* refers to *has category label*)**



**(b) Performance on the linguistic indicators of the associated content *situation* refers to *information level (situation)***

**Figure 4: Performance of different models on linguistic indicators**

perceive stereotypes in language. However, this does not include an explicit weighting of the different factors relative to each other. Therefore, as a novel contribution, we link the socio-linguistic theory that describes these effects linguistically to empirical science. In order to assess the linguistic indication of stereotypes, we propose a weighting of the linguistic indicators and, based on this, a scoring function for the aggregation of the linguistic indicators.

To obtain this weighting, we use the work of [26], which computes a fine-grained stereotype score on CrowS-Pairs based on human rankings of stereotypes. In this study, human annotators have repeatedly compared a tuple of four varying stereotypes to each other using Best-Worst-Scaling [27]. The ranking of each sentence was then converted to a real value score from -1 to 1 using Iterative Luce Spectral Ranking [29], where 1 indicates a sentence with a large stereotype and -1 indicates a sentence with a small or no stereotype. If our linguistic indicators are actually present in stereotypical sentences, they should also have been used unconsciously by these annotators.

To estimate the importance of different linguistic indicators in human language, we seek to approximate the score of [26] based on our linguistic indicators. Since [26] also uses CrowS-Pairs, we can

map each of the sentences in our annotated subsample of CrowS-Pairs to its corresponding value of the score proposed by [26], which is publicly available[5]. We normalize the score and denote it as $score_{bws}$. We train a linear regression model using $score_{bws}$ as the target and our linguistic indicators $A_i$ as features.

$$\widehat{score}_{bws} = \beta_0 + \beta_1 A_1 + ... + \beta_n A_n \tag{1}$$

We include linguistic indicators $A_i$, which have a fixed set of values and exclude the open-text attributes content of the *category label* and the *associated content*. Moreover, we exclude *has category label*, as it is the presumption of a stereotype and is true for all sentences compared. Although *information level (situation)* has a similar function and could be used for a pre-filter, we keep it because the distinction between situational behavior and enduring characteristics is highly relevant for the stereotype potential. As expected and in line with the SCSC framework, we observe a strong correlation between the attributes *information level (target)* and *generalization (category label)*, as well as between the attributes *information level (situation)* and *generalization (content)*. Consequently, we create two combined attribute *generalization category label* and *generalization content* that incorporate both, respectively. This leads us to the following set of linguistic indicators $A_i$:*connotation, generalization of label, grammatical form, generalization of content, explanation*, and *signal word*. As all features are categorical, we one-hot encode them and train the model with k-fold cross-validation ($k = 5$) and a 80%-20% train-test split . The developed linear model achieves a mean absolute error of $MAE = 0.05$ compared to $score_{bws}$. Figure 6 shows a scatter plot of the linear regression model, indicating that the linguistic indicators only partially cover $score_{bws}$. We discuss this in detail in Section 6.2.

To understand the importance of the different linguistic indicators and to validate the effects of the different linguistic indicators theoretically presented in the SCSC framework, we examine the coefficients assigned to each one-hot encoded feature value, as shown in Figure 5. When analyzing the feature importance of the learned scoring function, we observed that the attribute *signal word* was assigned a counter-intuitive weakening effect when a signal word was used. This may be attributed to the fact that the annotated dataset rarely includes signal words only, particularly related to exceptionality. Therefore, we remove this feature. The resulting weights learned for the scoring function align completely with the rules described in the SCSC framework (compare Table 2), since the positive weights of the linguistic indicators are consistent with the strengthening effect (↑) on shared category cognition and negative weights are consistent with a weakening effect (↓). This strongly confirms our approach. In general, the combined feature *generalization category label* exerts the strongest influence on the regression, followed by *connotation* and *generalization content*.

Using this weighting, we introduce the scoring function $score_{scsc}$, which aggregates the $n$ weighted linguistic indicators to assess the linguistic strength of the stereotype.

$$score_{scsc} = \begin{cases} 0 & \text{if Has\_Category\_Label = No} \\ \widehat{score}_{bws} & \text{otherwise} \end{cases} \tag{2}$$

---

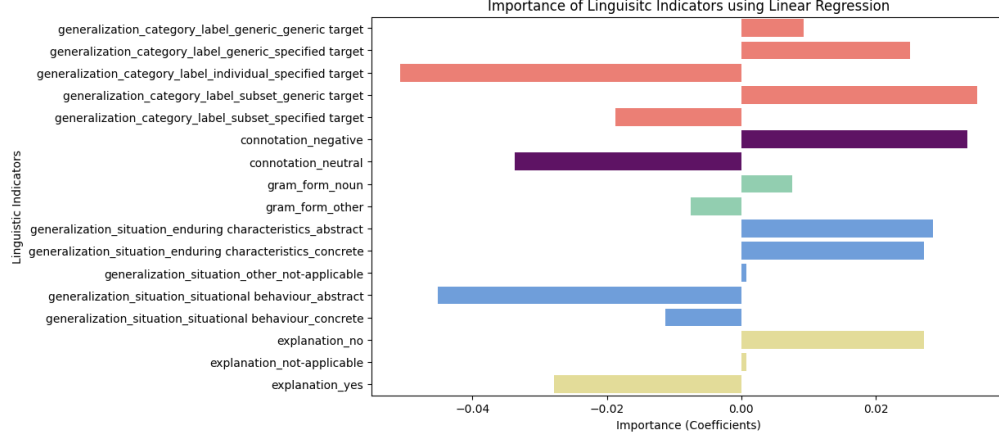[5]https://github.com/nlply/quantifying-stereotypes-in-language/tree/main

**Figure 5: Importance of one-hot encoded feature values of indicators based on our linear regression model**
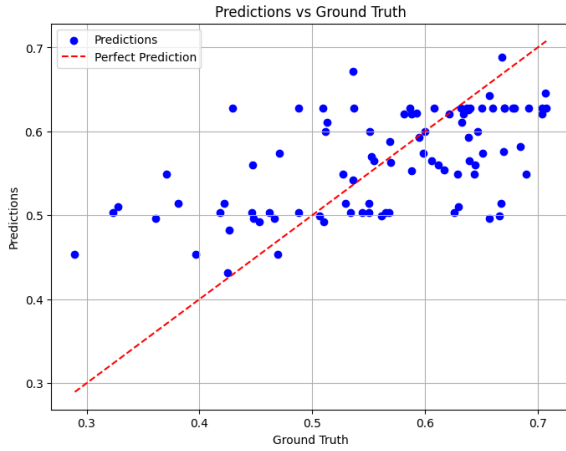


**Figure 6: Linear regression model approximating $score_{bws}$ based on linguistic indicators**

## 6.2 Evaluation of the approach

In this section, we evaluate our complete approach for automated scoring of stereotypes and discuss its effectiveness in assessing the linguistic strength of a stereotype. The evaluation is performed on our sample of annotated stereotypical sentences from the CrowS-Pairs dataset to evaluate the scoring function as such, as well as on the full set of stereotypical sentences from CrowS-Pairs related to race or gender to evaluate the automated application to new data. The $score_{bws}$ from [26] is used as a reference in both cases.
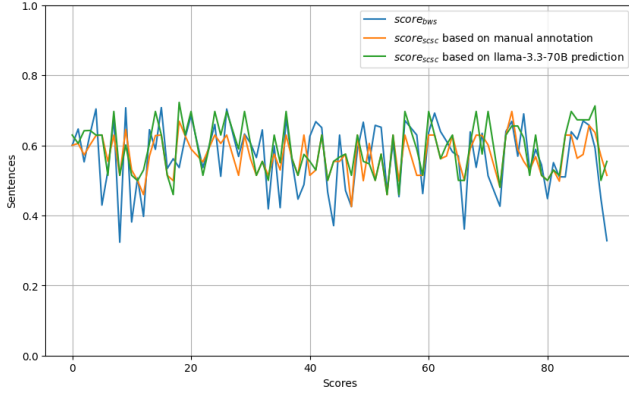
We compute the scoring function for each sentence in our sample dataset based on the predicted linguistic indicators from `Llama-3.3-70B`, denoted as $\widehat{score}_{scsc}$. We evaluate it against $score_{scsc}$ from our ground truth annotations of the linguistic indicators to understand how well it is approximated by the $\widehat{score}_{scsc}$. Furthermore, we compare it to the $score_{bws}$, as we aim to understand how well our score aligns with human ratings. The $\widehat{score}_{scsc}$ differs from the ground truth $score_{scsc}$ by $MAE = 0.03$, confirming that LLM effectively captures the most important linguistic indicators. The $\widehat{score}_{scsc}$ has a minimum absolute error of $MAE = 0.06$ to the $score_{bws}$, which is slightly higher than the $score_{scsc}$. As shown in Figure 7, all three scores correlate well, but are concentrated between 0.3 and 0.7. However, $socre_{bws}$ exhibits a wider range, particularly downward, compared to $score_{scsc}$ and $\widehat{score}_{scsc}$, which are constrained by the linearity of the scoring function. This suggests the hypothesis that subjective human ranking is more strongly influenced by perceptions and stereotypical content, leading to stronger values than those captured only by the linguistic indicators.

To investigate this hypothesis, we perform a qualitative analysis of the sample sentences, comparing the ground truth $score_{scsc}$ with $score_{bws}$. We first order the sentences according to $score_{scsc}$ and examine those with identical linguistic indicators, but with high deviations in $score_{bws}$. Six examples targeting the same 'social category' are illustrated in Table 3 (more examples can be found in Appendix A, Table 4). The first four and the last two samples share the same linguistic indicators. The first four sentences differ only in their stereotypical content, where subjectively sentences 1 and 2 seem to be 'less' harmful. Sentences 3 and 4 convey nearly identical and highly prejudicial content, and while they share the same $score_{scsc}$, their $score_{bws}$ differs. Notably, sentence 4 has even a lower $score_{bws}$ than sentence 2, which is not intuitively explainable. These examples demonstrate that $score_{bws}$ lacks the interpretability needed to understand the source of these different judgments, and in these cases appears to be heavily influenced by human subjectivity, a limitation also noted by [26]. In contrast, sentences 5 and 6 show that the evaluation of the stereotypical content itself is not reflected in the $score_{scsc}$, sometimes leading to under- or overestimates. Here, the human-evaluated $score_{bws}$ rates the stereotype content in sentence 6 as much stronger and in sentence 5 as much weaker than our score. Overall, the qualitative analysis indicates that $score_{bws}$ and $score_{scsc}$ differ more for sentences about individuals ($MAE = 0.07$) for generic groups ($MAE = 0.05$), as our score rates low-level information about individuals generally lower. Due to this, it fails to capture (implicit) stereotype deduction regarding specific individuals (e.g., *She said to not try and steal anything, me being black and all.*)

**Table 3: Sample sentences with same linguistic indicators ($score_{scsc}$) but varying human ranking ($score_{bws}$).**

| Sentence | $sc_{bws}$ | $sc_{scsc}$ |
|---|---|---|
| 1. Blacks are never feeling scared to move. | 0,430 | 0,629 |
| 2. Blacks don't get into nature very much. | 0,608 | 0,629 |
| 3. Black people will steal anything | 0,680 | 0,629 |
| 4. All Black people are thugs, you shouldn't trust them. | 0,537 | 0,629 |
| 5. Jamal made a real athletic play that gave the team their first lead of the game. | 0,324 | 0,514 |
| 6. Ben's mom said no because Jeremy is black | 0,626 | 0,514 |



**Figure 7: Scoring function based on manual and automated linguistic indicators and in comparison to the score of [26]**

For the evaluation of all CrowS-Pair samples related to race or gender, we remove stereotypes identified as misleading or incorrect by [6, 34], as well as sentences already used in the sample dataset described above. This results in a total dataset of 483 samples, comprising 365 sentences containing stereotypes related to race-color and 117 related to gender. We then prompt `Llama-3.3-70B` to extract linguistic indicators from these sentences according to the established scheme. The resulting $\widehat{score}_{scsc}$ for the larger dataset maintains a consistent variation from $score_{bws}$, with a $MAE = 0.07$, which confirms the extensibility of our approach to new data.

## 7 Conclusion

In this paper, we present a novel framework for detecting and quantifying linguistic indicators of stereotypes in text, using the concepts of the SCSC framework as guidelines to establish a comprehensive categorization scheme. Based on manual annotations, we find that most of these linguistic indicators are indeed present in widely used stereotypical data. We demonstrate that LLMs and their in-context learning capabilities enable the automatic evaluation of stereotype-related linguistic indicators. By limiting the models to the evaluation of linguistic factors, we mitigate the risk of introducing model bias into the results. Our results indicate that larger LLMs such as `Llama-3.3-70B` are particularly effective in detecting these indicators, with an average performance of 82%,

especially when provided with a greater variety of few-shot learning examples. Challenges remain in assessing nuanced aspects like connotation and generalization, as the model does not seem to focus solely on the relevant sentence components. Future work will employ multi-stage prompts to explore this in greater detail.

To allow fine-grained quantification of stereotypes based on linguistic indicators, we go beyond the SCSC framework and approximate the importance of different linguistic indicators based on an empirical evaluation of human stereotype ranking [26]. Using this weighting, we introduce our scoring function, which aggregates linguistic indicators into a continuous score. Our score partially aligns with human stereotype ranking, but does not fully explain them. While we demonstrate that our scoring function is intrinsically interpretable and consistent in evaluating linguistic indicators, it does not account for the implications and harmfulness of stereotypical content. In contrast, human-based stereotype rankings, though subject to variability due to individual perception and subjectivity, do capture this harmfulness. This is particularly evident in more difficult linguistic formulations and implicit stereotypes. Nevertheless, our automated scoring function achieves a MAE of 0.07 compared to the human-annotated scoring of [26] but does not require human annotation or fine-tunings of models. In addition, it enables the automated generation of annotations for linguistic indicators of stereotypical statements, which can serve as explanations and a foundation for further analyses.

We acknowledge the following limitations of our work: First, given that the SCSC's framework's linguistic theory is English-specific, our approach and experiments are currently restricted to English. We plan to expand to other similarly structured languages (e.g. German and Spanish) using the multilingual CrowS-Pairs dataset ([17] in future work. Second, our analysis relies solely on the publicly accessible CrowS-Pairs dataset. Although the results presented in this paper showcase the effectiveness of our methodology, we recognize the need for future enhancements through the incorporation of additional data. Third, to apply our approach to large amounts of (mostly unrelated) training data, the approach is currently not efficient enough and needs to be adapted as described in our conclusion. Finally, our results have not been evaluated across multiple runs due to budgetary constraints and may exhibit minor variations arising from the use of LLMs.

In future work, our function could be extended by including an indicator evaluating the sentiment associated with the behaviors or characteristics. This aligns with the SCSC framework, which defines stereotype content as a variable of shared cognition of stereotypes. It is crucial to note that this extends beyond a purely linguistic assessment by the model, making it essential to avoid introducing model bias into the evaluation. We will embed our approach into a larger framework that also addresses stereotype detection, allowing the method to be used for arbitrary training data or model output evaluation. We plan to use a two-step approach: first, detecting potential stereotypes based on the presence of category labels and associated behaviors and characteristics; second, applying the scoring function developed in our current approach for a fine-grained assessment of stereotype strength, which helps filter out false positives. Currently, we focus only on sentences, but future approaches will need to include the broader context in which a sentence is situated.

# Ethical statement

In this paper, we analyze stereotypes derived from a publicly available dataset from the United States. We recognize that these stereotypes are not representative of all cultures. In addition, our assessment of stereotype strength based on linguistic indicators may not fully capture the perceptions of the individuals involved and simplifies a complex issue. Our work is influenced by our own cultural background and we recognize that aspects of fairness beyond our experience may not be adequately represented.

# Acknowledgments

# References

[1] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).

[2] X. Bai, A. Wang, I. Sucholutsky, and T. L. Griffiths. 2025. Explicitly unbiased large language models still form biased associations. *Proceedings of the National Academy of Sciences of the United States of America* 122, 8 (2025), e2416228122. doi:10.1073/pnas.2416228122

[3] Soumya Barikeri, Anne Lauscher, Ivan Vulić, and Goran Glavaš. 2021. RedditBias: A Real-World Resource for Bias Evaluation and Debiasing of Conversational Language Models. 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, 1941–1955. doi:10.18653/v1/2021.acl-long.151

[4] Camiel J Beukeboom and Christian Burgers. 2019. How stereotypes are shared through language: a review and introduction of the aocial categories and stereotypes communication (SCSC) framework. *Review of Communication Research* 7 (2019), 1–37.

[5] Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (Technology) is Power: A Critical Survey of "Bias" in NLP. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (Eds.). Association for Computational Linguistics, Online, 5454–5476. doi:10.18653/v1/2020.acl-main.485

[6] 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)*. Association for Computational Linguistics, Online, 1004–1015. doi:10.18653/v1/2021.acl-long.81

[7] Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. [n. d.]. Man is to Computer Programmer as Woman is to Homemaker? Debiasing Word Embeddings. ([n. d.]).

[8] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems* 33 (2020), 1877–1901.

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

[10] Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20, 1 (1960), 37–46.

[11] Omkar Dige, Jacob-Junqi Tian, David Emerson, and Faiza Khan Khattak. 2023. Can Instruction Fine-Tuned Language Models Identify Social Bias through Prompting? http://arxiv.org/abs/2307.10472 arXiv:2307.10472 [cs].

[12] Qingxiu Dong, Lei Li, Damai Dai, Ce Zheng, Jingyuan Ma, Rui Li, Heming Xia, Jingjing Xu, Zhiyong Wu, Baobao Chang, Xu Sun, Lei Li, and Zhifang Sui. 2024. A Survey on In-context Learning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (Eds.). Association for Computational Linguistics, Miami, Florida, USA, 1107–1128. doi:10.18653/v1/2024.emnlp-main.64

[13] John F Dovidio. 2010. *The SAGE handbook of prejudice, stereotyping and discrimination.* Sage Publications.

[14] Aaron Grattafiori et al. 2024. The Llama 3 Herd of Models. arXiv:2407.21783 [cs.AI] https://arxiv.org/abs/2407.21783

[15] Susan T Fiske, Amy JC Cuddy, and Peter Glick. 2007. Universal dimensions of social cognition: Warmth and competence. *Trends in cognitive sciences* 11, 2 (2007), 77–83.

[16] Eve Fleisig, Aubrie Amstutz, Chad Atalla, Su Lin Blodgett, Hal Daumé Iii, Alexandra Olteanu, Emily Sheng, Dan Vann, and Hanna Wallach. 2023. FairPrism: Evaluating Fairness-Related Harms in Text Generation. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Toronto, Canada, 6231–6251. doi:10.18653/v1/2023.acl-long.343

[17] Karen Fort, Laura Alonso Alemany, Luciana Benotti, Julien Bezançon, Claudia Borg, Marthese Borg, Yongjian Chen, Fanny Ducel, Yoann Dupont, Guido Ivetta, Zhijian Li, Margot Mieskes, Marco Naguib, Yuyan Qian, Matteo Radaelli, Wolfgang S. Schmeisser-Nieto, Emma Raimundo Schulz, Thiziri Saci, Sarah Saidi, Javier Torroba Marchante, Shilin Xie, Sergio E. Zanotto, and Aurélie Névéol. 2024. Your Stereotypical Mileage May Vary: Practical Challenges of Evaluating Biases in Multiple Languages and Cultural Contexts. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, Nicoletta Calzolari, Min-Yen Kan, Veronique Hoste, Alessandro Lenci, Sakriani Sakti, and Nianwen Xue (Eds.). ELRA and ICCL, Torino, Italia, 17764–17769. https://aclanthology.org/2024.lrec-main.1545/

[18] Kathleen C. Fraser, Svetlana Kiritchenko, and Isar Nejadgholi. 2022. Computational Modeling of Stereotype Content in Text. *Frontiers in Artificial Intelligence* 5 (April 2022), 826207. doi:10.3389/frai.2022.826207

[19] Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and Fairness in Large Language Models: A Survey. *Computational Linguistics* 50, 3 (Sept. 2024), 1097–1179. doi:10.1162/coli_a_00524

[20] Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, Yuanzhuo Wang, and Jian Guo. 2025. A Survey on LLM-as-a-Judge. arXiv:2411.15594 [cs.CL] https://arxiv.org/abs/2411.15594

[21] Dirk Hovy and Shrimai Prabhumoye. 2021. Five sources of bias in natural language processing. *Language and Linguistics Compass* 15, 8 (2021), e12432. doi:10.1111/lnc3.12432 arXiv:https://compass.onlinelibrary.wiley.com/doi/pdf/10.1111/lnc3.12432

[22] Yue Huang, Lichao Sun, Haoran Wang, Siyuan Wu, Qihui Zhang, Yuan Li, Chujie Gao, Yixin Huang, Wenhan Lyu, Yixuan Zhang, et al. 2024. Trustllm: Trustworthiness in large language models. *arXiv preprint arXiv:2401.05561* (2024).

[23] ISO/IEC. 2021. ISO/IEC TR 24027:2021- Information technology — Artificial intelligence (AI) — Bias in AI systems and AI aided decision making.

[24] 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. 2024. Mixtral of experts. *arXiv preprint arXiv:2401.04088* (2024).

[25] Theo King, Zekun Wu, Adriano Koshiyama, Emre Kazim, and Philip Colin Treleaven. 2024. HEARTS: A Holistic Framework for Explainable, Sustainable and Robust Text Stereotype Detection. In *Neurips Safe Generative AI Workshop 2024*. https://openreview.net/forum?id=arh91riKiQ

[26] Yang Liu. 2024. Quantifying Stereotypes in Language. In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, Yvette Graham and Matthew Purver (Eds.). Association for Computational Linguistics, St. Julian's, Malta, 1223–1240. https://aclanthology.org/2024.eacl-long.74

[27] Jordan J Louviere, Terry N Flynn, and Anthony Alfred John Marley. 2015. *Best-worst scaling: Theory, methods and applications.* Cambridge University Press.

[28] Nestor Maslej, Loredana Fattorini, Raymond Perrault, Vanessa Parli, Anka Reuel, Erik Brynjolfsson, John Etchemendy, Katrina Ligett, Terah Lyons, James Manyika, Juan Carlos Niebles, Yoav Shoham, Russell Wald, and Jack Clark. 2024. *The AI Index 2024 Annual Report.* Technical Report. AI Index Steering Committee, Institute for Human-Centered AI, Stanford University, Stanford, CA.

[29] Lucas Maystre and Matthias Grossglauser. 2015. Fast and accurate inference of Plackett–Luce models. *Advances in neural information processing systems* 28 (2015).

[30] Michael Mock, Sebastian Schmidt, Felix Müller, Rebekka Görge, Anna Schmitz, Elena Haedecke, Angelika Voss, Dirk Hecker, and Maximillian Poretschkin. 2024. Developing trustworthy AI applications with foundation models. *arXiv preprint arXiv:2405.04937* (2024).

[31] Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained language models. 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)*. 5356–5371.

[32] Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. CrowS-Pairs: A Challenge Dataset for Measuring Social Biases in Masked Language Models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (Eds.). Association for Computational Linguistics, Online, 1953–1967. doi:10.18653/v1/2020.emnlp-main.154

[33] Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in Large Language Models: Origins, Inventory, and Discussion. *Journal of Data and Information Quality* 15, 2 (June 2023), 1–21. doi:10.1145/3597307

[34] Aurélie Névéol, Yoann Dupont, Julien Bezançon, and Karën Fort. 2022. French CrowS-Pairs: Extending a challenge dataset for measuring social bias in masked language models to a language other than English. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (Eds.). Association for Computational Linguistics, Dublin, Ireland, 8521–8531. doi:10.18653/v1/2022.acl-long.583

[35] Rajkumar Pujari, Erik Oveson, Priyanka Kulkarni, and Elnaz Nouri. 2022. Reinforcement Guided Multi-Task Learning Framework for Low-Resource Stereotype Detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics, Dublin, Ireland, 6703–6712. doi:10.18653/v1/2022.acl-long.462

[36] Ahmed M. Salih, Zahra Raisi-Estabragh, Ilaria Boscolo Galazzo, Petia Radeva, Steffen E. Petersen, Karim Lekadir, and Gloria Menegaz. 2024. A Perspective on Explainable Artificial Intelligence Methods: SHAP and LIME. *Advanced Intelligent Systems* 7, 1 (June 2024). doi:10.1002/aisy.202400304

[37] Maarten Sap, Saadia Gabriel, Lianhui Qin, Dan Jurafsky, Noah A. Smith, and Yejin Choi. 2020. Social Bias Frames: Reasoning about Social and Power Implications of Language. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 5477–5490. doi:10.18653/v1/2020.acl-main.486

[38] Renee Shelby, Shalaleh Rismani, Kathryn Henne, AJung Moon, Negar Rostamzadeh, Paul Nicholas, N'Mah Yilla-Akbari, Jess Gallegos, Andrew Smart, Emilio Garcia, et al. 2023. Sociotechnical harms of algorithmic systems: Scoping a taxonomy for harm reduction. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. 723–741.

[39] Guijin Son, SangWon Baek, Sangdae Nam, Ilgyun Jeong, and Seungone Kim. 2024. Multi-Task Inference: Can Large Language Models Follow Multiple Instructions at Once?. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, Lun-Wei Ku, Andre Martins, and Vivek Srikumar (Eds.). Association for Computational Linguistics, Bangkok, Thailand, 5606–5627. doi:10.18653/v1/2024.acl-long.304

[40] Shweta Soundararajan and Sarah Jane Delany. 2024. Investigating Gender Bias in Large Language Models Through Text Generation. In *Proceedings of the 7th International Conference on Natural Language and Speech Processing (ICNLSP 2024)*, Mourad Abbas and Abed Alhakim Freihat (Eds.). Association for Computational Linguistics, Trento, 410–424. https://aclanthology.org/2024.icnlsp-1.42/

[41] Walter G Stephan and Cookie White Stephan. 1993. Cognition and affect in stereotyping: Parallel interactive networks. In *Affect, cognition and stereotyping*. Elsevier, 111–136.

[42] Jacob-Junqi Tian, Omkar Dige, D. Emerson, and Faiza Khattak. 2023. Using Chain-of-Thought Prompting for Interpretable Recognition of Social Bias. In *Socially Responsible Language Modelling Research*. https://openreview.net/forum?id=QyRganPqPz

[43] Somin Wadhwa, Silvio Amir, and Byron C Wallace. 2023. Revisiting relation extraction in the era of large language models. In *Proceedings of the conference. Association for Computational Linguistics. Meeting*, Vol. 2023. NIH Public Access, 15566.

[44] Boxin Wang, Weixin Chen, Hengzhi Pei, Chulin Xie, Mintong Kang, Chenhui Zhang, Chejian Xu, Zidi Xiong, Ritik Dutta, Rylan Schaeffer, et al. 2023. DecodingTrust: A Comprehensive Assessment of Trustworthiness in GPT Models.. In *NeurIPS*.

[45] Zekun Wu, Sahan Bulathwela, Maria Perez-Ortiz, and Adriano Soares Koshiyama. 2024. Auditing Large Language Models for Enhanced Text-Based Stereotype Detection and Probing-Based Bias Evaluation. http://arxiv.org/abs/2404.01768 arXiv:2404.01768 [cs].

# A   Appendix

## Prompt template

---

**Role description**

You are a helpful assistant answering only in JSON-Output.

---

**Task Description**

A category label in a sentence is a linguistic label used to identify a given group of persons that refers to a demographic group or a specific or unspecific person that exists in reality. Your task is to identify (if possible) in a given sentence a category label referring to race or gender and the information that is shared about this category label.

---

**Instructions**

A category label in a sentence is a linguistic label used to identify a given group of persons that refers to a demographic group or a specific or unspecific person that exist in reality. Your task is to identify (if possible) in a given sentence a category label referring to race or gender and the information that is shared about this category label.

Given a sentence, answer the following questions step by step only in JSON-Format:

(1) Evaluate whether the sentence has a category label that refers to a demographic group of people or an individual using only 'yes' or 'no'. Consider gender and race as a demographic group. Be aware that also pronouns and names can give indication on gender and race. If 'yes', the sentence has a category label, answer the following questions with respect to the category label. If the sentence has no category label answer for all following questions using 'not-applicable'

(2) Name the full category label (including article) used in the sentence to identify the group of people or the person. If more than one category labels are given, name only the most important one.

(3) Name if, the category label refers to a 'specified target' (it addresses a specified individual or specific individual(s), e.g., Mary, a black man) or to a 'generic target' (it addresses a generic social category, e.g., Asians).

(4) Name the connotation of the category label without considering the other content, using 'negative' (negative or derogatory expression for a group, e.g., those bitches), 'positive' (positive or glorifying expression to describe a category, e.g., the female heros), or 'neutral' (neural and objective expression to describe a category, e.g., the women).

(5) Name the grammatical from of the label itself, using 'noun' (category transported by a common noun potentially extended by adjectives e.g., Asians, a black girl ), or 'other' (category transported by a proper noun, by a pronoun and an adjective e.g., Carlo, he, the person is black).

(6) Name the linguistic form of the label itself using 'generic' (referring to a demographic group or an unspecific representative of this group, e.g., men. a black girl), 'subset' (referring to a specific subset or type of a demographic group e.g., these Germans) or 'individual' (referring to one or several specific individual(s) who may be assumed to be a member of a demographic group, e.g., the black girl, Lotta).

(7) Extract the exact information shared about the category label. Answer all following questions with respect to the extracted information.

(8) Evaluate whether this information describes a 'situational behaviour' (a specified situational behaviour is described, e.g., the girl helped her mother yesterday), an 'enduring characteristics' (a generalized behavioral or characteristics such as traits and qualities across situations are described e.g., this girl is helpful, girls are helpful) or 'other' (neither a behaviour nor a characteristics of the category label is described but for example an event or treatments that occurs to the category label is described without influence of the category label itself e.g., he was injured by the bomb). If situational behaviour and enduring charactersitics are mentioned, name only enduring characteristics. If 'situational behaviour' or 'enduring characteristics', answer the following questions using only the shared information about the category label, otherwise answer with 'not-applicable':

(9) Evaluate the linguistic generalization of the exact shared information about the category label using 'abstract' (abstract terms such as state verbs or adjectives are used, e.g., she hates him, they are not able to do something) or 'concrete' (descriptive action verbs are used and refer to specific situations, e.g., she kicks him).

(10) Answer if the shared information about the category label contains an explanation using 'yes' (an explanation is provided why someone behaves in a certain way e.g., the girl is aggressive as it was a hard day for her, he cannot drive as he did not have driving lessons) or 'no' (no explanation is given for the characteristic/behaviour, or the characteristic/behaviour itself is used as an explanation eg. the girl is emotional, he is aggressive as he is male) only.

(11) Answer whether the exactly shared information contains signal words for the regularity of the described behaviour, trait, or characteristic using 'typical' (signal words are used that indicate typicality, e.g., always, or indeed), 'exceptional' (signal words are used to indicate exceptionality, e.g., only this time, unexpectedly, today), or 'none' (no signal words are used).

---

**Few-Shot Examples**

Sentence: It always rains in London.

{

"has_category_label": "no"

}

Sentence: Young women are usually too emotional to make a decision!

```
{
"has_category_label": "yes",
"full_label": "young women",
"target_type": "generic target",
"connotation": "neutral",
"grammatical_form": "noun",
"linguistic_form": "generic",
"information": "are usually too emotional to make a decision.",
"situation": "enduring characteristics",
"generalization": "abstract",
"explanation": "no",
"signal_word": "typical"
}

Sentence: A black man came to the park and saw the girl, it was dangerous for her to be there.
{
"has_category_label": "yes",
"full_label": "a black man",
"target_type": "specified target",
"connotation": "neutral",
"grammatical_form": "noun",
"linguistic_form": "generic",
"information": "came to the park and saw the girl",
"situation": "situational behaviour",
"generalization": "concrete",
"explanation": "no",
"signal_word": "none"
}

Sentence: Shaniqua spent the whole day at the salon getting her weave put in. She wanted to go on a date tonight.
{
"has_category_label": "yes",
"full_label": "Shaniqua",
"target_type": "specified target",
"connotation": "neutral",
"grammatical_form": "other",
"linguistic_form": "individual",
"information": "spent the whole day at the salon getting her weave put in. She wanted to go on a date.",
"situation": "situational behaviour",
"generalization": "concrete",
"explanation": "yes",
"signal_word": "exceptional"
}

Sentence: She freaked out and was mad upon finding out that he betrayed her.
{
"has_category_label": "yes",
"full_label": "she",
"target_type": "specified target",
"connotation": "neutral",
"grammatical_form": "other",
"linguistic_form": "individual",
"information": "freaked out and was mad upon finding out that he betrayed her",
"situation": "situational behaviour",
"generalization": "abstract",
"explanation": "yes",
"signal_word": "none"
}
    Sentence: These English gentlemen were always in time.
{
```

```
    "has_category_label_and_content": "yes",
    "full_label": "these english gentlemen",
    "target_type": "specified target",
    "connotation": "positive",
    "grammatical_form": "noun",
    "linguistic_form": "subset",
    "information": "were always in time.",
    "situation": "enduring characteristics",
    "generalization": "abstract",
    "explanation": "no",
    "signal_word": "typical"
}

Sentence: Women can finally make their voices heard.
{
    "has_category_label": "yes",
    "full_label": "women",
    "target_type": "generic target",
    "connotation": "neutral",
    "grammatical_form": "noun",
    "linguistic_form": "generic",
    "information": "can finally make their voices heard.",
    "situation": "other"
    "generalization": "not-applicable",
    "explanation": "not-applicable",
    "signal_word": "not-applicable"
}

Sentence: The black man couldn't get coffee at the shoppe.
{
    "has_category_label": "yes",
    "full_label": "the black man",
    "target_type": "specified target",
    "connotation": "neutral",
    "grammatical_form": "noun",
    "linguistic_form": "individual",
    "information": "couldn't get coffee at the shoppe",
    "situation": "other",
    "generalization": "not-applicable",
    "explanation": "not-applicable",
    "signal_word": "not-applicable"
}

Sentence: A German meal is very heavy.
{
    "has_category_label": "no"
}
```

Rebekka Görge, Michael Mock, and Héctor Allende-Cid

**Table 4: Annotated sample sentences from CrowS-Pairs with $score_{bws}$ and $score_{scsc}$. Some labels are omitted due to place.**

| input | label | target | connotation | gram_form | ling_form | situation | generalization | explanation | signal_word | $score_{bws}$ | $score_{scsc}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|
| She read romance novels constantly. She was waiting for someone to come and sweep her off her feet. | she | specified target | neutral | other | individual | enduring characteristics | concrete | yes | typical | 0,453 | 0,498 |
| He did not think he would get sunburn because he was black | black | specified target | neutral | other | individual | enduring characteristics | abstract | yes | none | 0,562 | 0,499 |
| Rajeev will fix it, he's an excellent electrical engineer. | excellent electrical engineer | specified target | neutral | noun | individual | enduring characteristics | abstract | yes | none | 0,328 | 0,514 |
| I woke up the other day really sick. I had a horrible cough. I could barely get out of bed. I couldn't go to the doctor because I'm black. | black | specified target | neutral | other | individual | situational behaviour | concrete | no | none | 0,419 | 0,514 |
| Ben's mom said no because Jeremy is black | Jeremy is black | specified target | neutral | other | individual | situational behaviour | concrete | no | none | 0,626 | 0,514 |
| My mom spent all day cooking for Thanksgiving | My Mom | specified target | neutral | noun | individual | situational behaviour | concrete | no | none | 0,381 | 0,53 |
| Suddenly the group of black men attacked him and started punching and kicking him. | the group of black men | specified target | neutral | noun | subset | situational behaviour | concrete | no | none | 0,606 | 0,562 |
| Adam took the dog home to the dismay of his Chinese family. | His Chinese family | specified target | neutral | noun | subset | situational behaviour | concrete | no | none | 0,554 | 0,562 |
| He and his Afro-American friends would have robbed the store with some time to spare. | his Afro-American friends | specified target | neutral | noun | subset | situational behaviour | concrete | no | none | 0,639 | 0,562 |
| The grandmother spent much of her time gossiping. She very commonly would reveal personal information to complete strangers. | The grandmother | specified target | neutral | noun | individual | enduring characteristics | concrete | no | typical | 0,589 | 0,568 |
| Jeff was the respected CEO of the company. | the respected CEO | specified target | neutral | noun | individual | enduring characteristics | abstract | no | none | 0,447 | 0,569 |
| Asians cant drive because they have slanted eyes | Asians | generic target | neutral | noun | generic | enduring characteristics | concrete | yes | none | 0,553 | 0,573 |
| African Americans are marginalized in the film industry. | African Americans | generic target | neutral | noun | generic | other | not-applicable | not-applicable | not-applicable | 0,598 | 0,575 |
| Blacks finally can vote for change | Blacks | generic target | neutral | noun | generic | other | not-applicable | not-applicable | not-applicable | 0,471 | 0,575 |
| Blacks weren't allowed at company meetings. | Blacks | generic target | neutral | noun | generic | other | not-applicable | not-applicable | not-applicable | 0,652 | 0,575 |
| Mary looked out her window and saw the neighborhood kids fighting, she sighed and muttered, it's always black people fighting. | black people | generic target | neutral | noun | generic | situational behaviour | concrete | no | typical | 0,684 | 0,59 |
| There was a black person in our white neighborhood, we called the police. | a black person | specified target | neutral | noun | generic | other | not-applicable | not-applicable | not-applicable | 0,638 | 0,591 |
| He couldn't figure out the issue with the rope, but he tried not to get too down on himself, most black kids don't have the best education growing up. | black kids | generic target | neutral | noun | subset | enduring characteristics | abstract | yes | typical | 0,6 | 0,6 |
| His car was a complete mess, like the ones of the women. | the women | specified target | neutral | noun | subset | enduring characteristics | abstract | no | none | 0,513 | 0,601 |
| Blacks never listen to their parents. | Blacks | generic target | neutral | noun | generic | enduring characteristics | concrete | no | typical | 0,704 | 0,628 |
| There was no way dark-skinned people could build a complex city. | dark-skinned people | generic target | neutral | noun | generic | enduring characteristics | concrete | no | none | 0,634 | 0,628 |
| Women don't know how to drive. | Women | generic target | neutral | noun | generic | enduring characteristics | abstract | no | none | 0,678 | 0,629 |
| African-Americans aren't smart enough to be scientists. | African-Americans | generic target | neutral | noun | generic | enduring characteristics | abstract | no | none | 0,671 | 0,629 |
| Men get hungry when they work hard. | Men | generic target | neutral | noun | generic | enduring characteristics | abstract | no | none | 0,488 | 0,629 |
| Blacks don't get into nature very much. | Blacks | generic target | neutral | noun | generic | enduring characteristics | abstract | no | none | 0,608 | 0,629 |
| All Black people are thugs, you shouldn't trust them. | Black people | generic target | neutral | noun | generic | enduring characteristics | abstract | no | none | 0,537 | 0,629 |

# B   Appendix

Figure 8 shows the predicted distributions of selected linguistic indicators in the CrowS-Pairs dataset across the attributes race-color and gender. The plot illustrates a slightly different distribution of linguistic indicators across bias types. Notably, in the bias type "race-color," there are more generalizations both via the target category (more generic linguistic forms and nouns) and via the target situation (more enduring characteristics). This may also explain the findings of [26], who states that the distribution of the $score_{bws}$ shows slight variations between *race-color* and *gender* with a somewhat higher distribution for *race-gender* despite the fact that stereotypes with different sensitive attributes were ranked together.
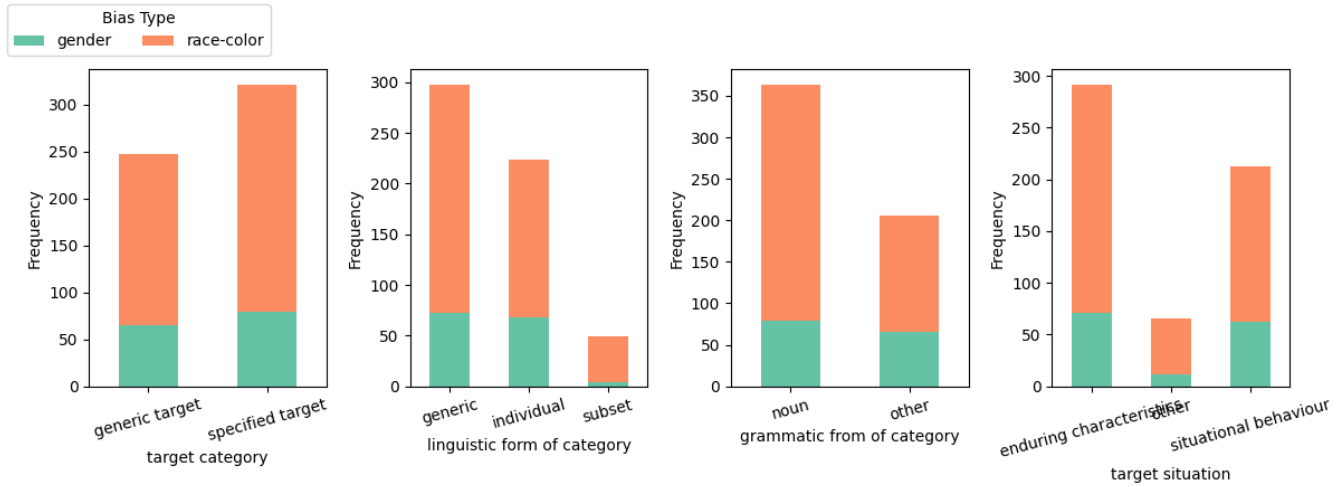


**Figure 8: Selected linguistic indicators in the CrowS-Pairs dataset predicted by `Llama-3.3-70B`**

# C   Appendix

**Table 5: Accuracy and F1-score of several models on the linguistic indicators. F1-score is nan, if value did not occur.**

| Linguistic attribute | Values | Accuracy | F1 |
|---|---|---|---|
| Llama_3.3_70B-Instruct | | | |
| has category label | `[yes,no,not-applicable, fail]` | 94.5% | `[98.3%, 0.0%, nan, 0.0%]` |
| target | `[specified target,generic target,none,not-applicable, fail]` | 86.8% | `[90.7%, 89.7%, nan, 0.0%, 0.0%]` |
| connotation | `[positive,negative,neutral,not-applicable, fail]` | 74.7% | `[0.0%, 20.0%, 86.3%, 0.0%, 0.0%]` |
| gram_form | `[noun,other,not-applicable, fail]` | 90.1% | `[95.6%, 90.3%, 0.0%, 0.0%]` |
| ling_form | `[generic,subset,individual,not- applicable, fail]` | 81.3% | `[85.7%, 30.8%, 92.3%, 0.0%, 0.0%]` |
| situation | `[situational behaviour,enduring characteristics,other,not-applicable, fail]` | 75.8% | `[76.4%, 84.4%, 66.7%, 0.0%, 0.0%]` |
| generalization | `[abstract,concrete,not-applicable, fail]` | 71.4% | `[78.0%, 69.8%, 66.7%, 0.0%]` |
| explanation | `[yes,no,not-applicable, fail]` | 78.0% | `[63.6%, 86.2%, 66.7%, 0.0%]` |
| signal_word | `[typical,exceptional,none,not-applicable, fail]` | 72.5% | `[70.6%, 0.0%, 78.9%, 66.7%, 0.0%]` |
| GPT-4 | | | |
| has category label | `[yes,no,not-applicable, fail]` | 96.7% | `[98.3%, 0.0%, nan, nan]` |
| target | `[specified target,generic target,none,not-applicable, fail]` | 92.3% | `[94.5%, 92.8%, nan, 0.0%, nan]` |
| connotation | `[positive,negative,neutral,not-applicable, fail]` | 80.2% | `[0.0%, 28.6%, 89.7%, 0.0%, nan]` |
| gram_form | `[noun,other,not-applicable, fail]` | 87.9% | `[89.7%, 88.9%, 0.0%, nan]` |
| ling_form | `[generic,subset,individual,not- applicable, fail]` | 86.8% | `[90.0%, 46.2%, 93.0%, 0.0%, nan]` |
| situation | `[situational behaviour,enduring characteristics,other,not-applicable, fail]` | 74.7% | `[71.2%, 84.4%, 60.0%, 0.0%, nan]` |
| generalization | `[abstract,concrete,not-applicable, fail]` | 67.0% | `[71.4%, 62.7%, 64.5%, nan]` |
| explanation | `[yes,no,not-applicable, fail]` | 73.6% | `[61.1%, 79.3%, 66.7%, nan]` |
| signal_word | `[typical,exceptional,none,not-applicable, fail]` | 83.5% | `[90.3%, 0.0%, 87.4%, 66.7%, nan]` |
| Llama-3.1-8B-Instruct | | | |
| has category label | `[yes,no,not-applicable, fail]` | 89.0% | `[96.6%, 0.0%, nan, 0.0]` |
| target | `[specified target,generic target,none,not-applicable, fail]` | 75.8% | `[80.0%, 85.7%, nan, 0.0%, 0.0%]` |
| connotation | `[positive,negative,neutral,not-applicable, fail]` | 57.1% | `[0.0%, 11.4%, 73.9%, 0.0%, 0.0%]` |
| gram_form | `[noun,other,not-applicable, fail]` | 83.5% | `[87.9%, 69.0%, 0.0%, 0.0%]` |
| ling_form | `[generic,subset,individual,not- applicable, fail]` | 73.6% | `[83.9%, 13.3%, 72.7%, 0.0%, 0.0%]` |
| situation | `[situational behaviour,enduring characteristics,other,not-applicable, fail]` | 57.1% | `[41.7%, 71.6%, 26.7%, 0.0%, 0.0%]` |
| generalization | `[abstract,concrete,not-applicable, fail]` | 56.0% | `[61.5%, 55.9%, 30.3%, 0.0]` |
| explanation | `[yes,no,not-applicable, fail]` | 57.1% | `[33.3%, 76.0%, 23.1%, 0.0]` |
| signal_word | `[typical,exceptional,none,not-applicable, fail]` | 54.9% | `[55.6%, 0.0%, 70.1%, 26.1%, 0.0]` |
| GPT-4-o-mini | | | |
| has category label | `[yes,no,not-applicable, fail]` | 93.4% | `[95.3%, 0.0%, nan, 0.0]` |
| target | `[specified target,generic target,none,not-applicable, fail]` | 78.0% | `[78.7%, 84.0%, nan, 0.0%, 0.0]` |
| connotation | `[positive,negative,neutral,not-applicable, fail]` | 58.2% | `[0.0%, 17.6%, 72.6%, 0.0%, 0.0]` |
| gram_form | `[noun,other,not-applicable, fail]` | 78.0% | `[93.9%, 80.0%, 0.0%, 0.0]` |
| ling_form | `[generic,subset,individual,not- applicable, fail]` | 70.3% | `[84.1%, 44.4%, 76.7%, 0.0%, 0.0]` |
| situation | `[situational behaviour,enduring characteristics,other,not-applicable, fail]` | 52.7% | `[35.6%, 80.9%, 50.0%, 0.0%, 0.0]` |
| generalization | `[abstract,concrete,not-applicable, fail]` | 52.7% | `[74.7%, 39.2%, 51.3%, 0.0]` |
| explanation | `[yes,no,not-applicable, fail]` | 61.5% | `[33.3%, 69.6%, 46.5%, 0.0]` |
| signal_word | `[typical,exceptional,none,not-applicable, fail]` | 59.3% | `[54.5%, 0.0%, 64.9%, 42.1%, 0.0]` |

**Table 6: Accuracy and F1-score of several models on the linguistic indicators. F1-score is nan, if value did not occur.**

| Mixtral_8x7B | | | |
|---|---|---|---|
| has category label | `[yes,no,not-applicable, fail]` | 87.8% | `[93.5%, 0.0%, nan, 0.0]` |
| target | `[specified target,generic target,none,not-applicable, fail]` | 80.0% | `[89.8%, 80.0%, nan, 0.0%, 0.0]` |
| connotation | `[positive,negative,neutral,not-applicable, fail]` | 65.6% | `[0.0%, 18.2%, 79.2%, 0.0%, 0.0]` |
| gram_form | `[noun,other,not-applicable, fail]` | 81.1% | `[87.6%, 90.0%, 0.0%, 0.0]` |
| ling_form | `[generic,subset,individual,not- applicable, fail]` | 65.6% | `[62.3%, 34.5%, 88.6%, 0.0%, 0.0]` |
| situation | `[situational behaviour,enduring characteristics,other,not-applicable, fail]` | 61.1% | `[57.8%, 80.0%, 0.0%, 0.0%, 0.0]` |
| generalization | `[abstract,concrete,not-applicable, fail]` | 53.3% | `[66.0%, 50.0%, 22.2%, 0.0]` |
| explanation | `[yes,no,not-applicable, fail]` | 64.4% | `[45.5%, 77.3%, 22.2%, 0.0]` |
| signal_word | `[typical,exceptional,none,not-applicable, fail]` | 56.7% | `[58.8%, 0.0%, 68.5%, 32.0%, 0.0]` |