# Anonymous ACL submission

## Abstract

While Large Language Models (LLMs) demonstrate remarkable text generation capabilities within Natural Language Processing (NLP), they risk perpetuating societal biases from their training data. This study analyzes gender bias in Spanish generative LLMs by examining adjectival descriptions associated with men and women. Utilizing specifically designed prompts and a Supersenses-based adjective categorization framework, our research uncovers patterns consistent with cultural stereotypes, echoing findings from masked language models. We also investigate the relationship between model size and the extent of these observed gender biases.

## 1 Introduction

The rapid development and widespread deployment of Large Language Models (LLMs) have revolutionized Natural Language Processing (NLP), fostering unprecedented advancements in both generative and interpretive capabilities. Excelling in diverse tasks such as text generation, summarization, and conversational AI, these models demonstrate a sophisticated grasp of linguistic patterns (Wang et al., 2023). Nevertheless, their dependence on vast training datasets presents a critical challenge: the potential for these models to replicate and even amplify societal biases inherent in this data.

Prior research (Bolukbasi et al., 2016), (Caliskan et al., 2017), (Doe, 2021) has extensively documented the prevalence of biases in foundational models, from word embeddings to more complex neural architectures. These studies reveal detrimental patterns of gender, racial, and religious prejudices, among others, manifested in model outputs. For instance, documented associations include stereotyping women in domestic roles and men in professional ones, or linking specific religions with violent extremism (Abid et al., 2021).

Such biases alarmingly persist even in state-of-the-art generative models engineered for enhanced contextual understanding (Zack et al., 2024), thereby emphasizing the urgent need for robust identification and mitigation strategies.

The presence of bias in LLMs extends beyond theoretical concerns, carrying significant real-world implications. Integration of these models into applications like recruitment systems, customer service platforms, and educational tools means their biased outputs can perpetuate harmful stereotypes and exacerbate societal inequities. Consequently, addressing this challenge necessitates a dual approach: rigorous evaluation to identify and quantify biases, and effective mitigation techniques to reduce or eliminate their impact. Notably, while substantial research has concentrated on English-language models, biases in other languages, including Spanish, remain comparatively underexplored, despite the growing proliferation of multilingual and region-specific LLMs.

Informed by previous work on bias in Spanish models (Doe, 2024), this paper investigates how contemporary Spanish generative LLMs portray gender through adjective-based descriptions. Our analysis specifically examines the characterization of male and female subjects across four key semantic domains: physical appearance (**BODY**), emotional states (**FEELING**), cognitive traits (**MIND**), and actions or habits (**BEHAVIOR**). To achieve this, we employ a structured methodology centered on carefully constructed template sentences designed to elicit adjectival responses from a diverse set of Spanish generative LLMs. The adjectives generated are subsequently categorized using the Supersenses taxonomy (Tsvetkov et al., 2014), facilitating a systematic semantic analysis of descriptive patterns.

Our findings reveal significant gender biases in these generative LLMs, consistent with earlier research on masked language models. Women are

predominantly characterized by adjectives related to physical attributes and emotions, while men are more frequently described using terms associated with their behaviors and intellectual capabilities. These observed patterns mirror prevailing societal stereotypes, raising concerns about the fairness and inclusivity of applications that leverage LLMs. Furthermore, this study explores the correlation between model size and the manifestation of such biases, identifying trends relevant to the development of more effective mitigation strategies.

This research contributes to the expanding body of knowledge on LLM biases by offering a detailed analysis focused on Spanish generative models. It particularly highlights the persistence of these biases across various model architectures, underscoring the critical need for continued investigation and development of equitable AI systems in non-English contexts.

## 2 Background

Large Language Models (LLMs) have revolutionized Natural Language Processing (NLP) by enabling human-like text understanding, generation, and reasoning across diverse applications. However, their deployment raises critical concerns about biases that disproportionately affect marginalized communities. This section provides a comprehensive review of current advancements, challenges, and methodologies in evaluating and mitigating bias in LLMs, synthesizing key works from various research domains.

### 2.1 Evolution and capabilities of Language Models

Modern Large Language Models (LLMs), exemplified by architectures such as GPT-3 (Brown et al., 2020), BERT (Devlin et al., 2019), T5 (Roberts et al., 2019), and GPT-4 (Achiam et al., 2023), are typically based on autoregressive or encoder-decoder frameworks and trained on extensive textual corpora. These models exhibit remarkable generalization capabilities, effectively performing tasks like classification, sentiment analysis, and language translation using few-shot or zero-shot learning paradigms (Bommasani et al., 2021), (Hegde and Patil, 2020). However, a significant concern arises from the vast scale and often uncurated nature of their training data, which can lead to the encoding and subsequent amplification of detrimental societal biases (Bender et al., 2021), (Navigli et al., 2023).

### 2.2 Manifestations of bias

Bias in Large Language Models (LLMs) manifests in several pernicious forms, each with distinct negative consequences.

One significant category is **representational harms**, where LLMs perpetuate stereotypes, generate toxic content, and reinforce exclusionary societal norms. For example, some models associate specific professions predominantly with particular genders, such as linking nursing with women or engineering with men (Sheng et al., 2019; Liang et al., 2021).

Another critical form, **allocational harms**, occurs when LLMs contribute to inequitable resource or opportunity distribution. These biases can be encoded into LLM-based decision-making systems, potentially leading to unfair outcomes in areas like hiring, loan applications, or healthcare, disadvantaging certain groups (Ferrara, 2023; Mehrabi et al., 2021).

Furthermore, LLMs can exhibit **language variability issues**. These include linguistic biases like misclassifying or stigmatizing certain dialects and underrepresenting or misrepresenting minority languages. For instance, African American Vernacular English (AAVE) has been erroneously labeled as non-standard by some language technologies, reflecting a bias against linguistic diversity (Mozafari et al., 2020; Sap et al., 2019).

Such biases undermine the fairness and equity of diverse LLM-reliant applications, including machine translation, question-answering, and content moderation. Research highlights issues like gender-biased translations (gender-neutral terms becoming gender-specific) (Měchura, 2022) and miscalibrated toxicity classifiers that may incorrectly flag non-toxic text from certain groups or miss harmful content (Dixon et al., 2018).

### 2.3 Metrics for bias evaluation

Evaluating bias in Large Language Models (LLMs) requires specialized metrics, often tailored to specific model architectures and bias types. These metrics fall into several key categories.

One category is **embedding-based metrics**, which quantify biases in word or sentence vector representations (embeddings). Examples include the Word Embedding Association Test (WEAT) (Caliskan et al., 2017) and the Sentence Encoder Association Test (SEAT) (May et al., 2019). These

methods typically measure cosine similarities between vector representations of social group terms and attribute terms to reveal underlying associations.

Another approach uses **probability-based metrics**. Methods like Pseudo-Log-Likelihood (Salazar et al., 2020) and CrowS-Pairs (Nangia et al., 2020) assess bias by examining model-assigned probabilities to linguistic constructions. This can involve comparing likelihoods of stereotypical versus anti-stereotypical sentences or evaluating token probability disparities in counterfactual sentence pairs.

Furthermore, **generated text metrics** evaluate biases directly in LLM-produced text. Techniques such as Social Group Substitution (Liang et al., 2021) and the Co-occurrence Bias Score (Nadeem et al., 2021) analyze co-occurrence patterns between social group terms and other words, or broader lexical associations in the output, to quantify biases.

Effective bias evaluation depends on selecting appropriate metrics and relevant datasets. Benchmarks like StereoSet (Nadeem et al., 2021) and ToxiGen (Hartvigsen et al., 2022) are designed to probe representational harms and stereotypical associations. Other datasets target specific gender, racial, and cultural disparities, providing a comprehensive toolkit for assessing bias in LLMs.

### 2.4 Bias mitigation techniques

Bias mitigation research in the context of Large Language Models (LLMs) explores a range of techniques that can be implemented at various stages of their development and deployment pipeline. These strategies are designed to identify and reduce biases that models may learn from training data or exhibit in their outputs.

One category of these methods involves **preprocessing techniques**, which are applied directly to the training data before the model learning commences. Examples of such approaches include data augmentation, a process aimed at creating more balanced datasets to ensure fair representation (Liang et al., 2021). Another preprocessing strategy is data filtering, which focuses on removing or reducing the influence of harmful patterns or biased information within the data prior to the model's training phase.

Further along the pipeline, **in-training adjustments** are employed to address bias during the model's actual learning process. These strategies often involve modifications to the core mechanics of model training. Common approaches include altering the loss function, which guides the model's learning, to specifically penalize the formation of biased associations or predictions (Ma et al., 2022). Another significant in-training method is adversarial debiasing, where techniques are used to challenge the model and prevent it from learning and perpetuating biases (Xu et al., 2020).

During the inference stage, when the model is actively generating text, **intra-processing interventions** (also known as Inference-Time techniques) can be utilized. These methods aim to modify the model's behavior as it produces output. Constrained decoding serves as a key example, where the generation process is guided or restricted to prevent the model from producing biased or undesirable outputs at the moment of inference (Sheng et al., 2021).

Finally, **post-processing techniques** are applied after the model has generated its output. These methods focus on refining or correcting the text produced by the LLM. This can involve rewriting portions of the generated text to remove harmful content or to correct stereotypical portrayals. Filtering the output to identify and remove biased language or information is another common post-processing approach (Dixon et al., 2018).

More fine-grained approaches, such as projection-based debiasing for contextual embeddings and selective parameter updating, represent emerging trends in developing sophisticated mitigation strategies (Choi et al., 2021).

### 2.5 Open challenges and future directions

Despite significant progress, several open challenges and compelling future directions persist in the domain of LLM bias.

One critical area is the pursuit of **fairness across diverse languages**. There is an urgent need to extend bias research and mitigation efforts beyond the predominantly studied high-resource languages. This expansion should encompass multilingual LLMs and a broader spectrum of low-resource languages, ensuring that advancements in fairness benefit a global user base (Blasi et al., 2022).

Another significant challenge lies in strengthening the **theoretical underpinnings of fairness**. The fundamental trade-offs that exist between a model's utility—its performance on primary tasks—and the goals of fairness require more profound theoretical investigation. This is particularly

3

crucial when models are trained on data that is inherently biased, and understanding these dynamics is key to developing truly equitable systems (Barocas et al., 2023).

Furthermore, the **development of robust evaluation standards** remains an essential goal for the field. Establishing standardized and comprehensive metrics, along with diverse and representative datasets, is vital. Such standards would ensure higher consistency and reliability in bias assessment, allowing for more comparable and rigorous evaluation of mitigation techniques across different models and studies, and provide broader coverage of potential biases (Suresh and Guttag, 2019).

Finally, effectively **addressing contextual and intersectional biases** presents a complex but crucial frontier. Biases often manifest contextually, changing based on the specific situation, or arise from the interplay of multiple identity characteristics (e.g., the compounded effects of race and gender). Tackling these nuanced forms of bias necessitates a deeper integration of sociolinguistic perspectives and methodologies into the development and evaluation of LLMs (Benjamin, 2019).

## 3 Methodology

This section details the structured methodology employed to investigate gender bias in Spanish Large Language Models (LLMs). Our approach focuses on a curated selection of models adhering to specific criteria: they must be generative, support the Spanish language, have publicly available weights, and be deployable on our experimental cluster—a constraint that precluded the inclusion of some exceptionally large models. The methodology unfolds in several stages: first, sentence templates are meticulously crafted to elicit adjectival descriptions, ensuring balanced and comparable contexts for male and female subjects. Subsequently, the selected LLMs generate adjectives in response to these templates, with their most probable completions analyzed based on internal probabilities and rankings. These generated adjectives are then systematically categorized into predefined semantic domains—namely BODY, BEHAVIOR, FEELING, and MIND—to enable a nuanced analysis of descriptive patterns. Finally, quantitative metrics are applied to identify and compare the biases evident in the characterization of male and female subjects.

### 3.1 Sentence template creation

To investigate gender bias in Spanish generative language models, the initial step involved the creation of specific template sentences for the gender categories under examination: **male** and **female**. Each template was carefully engineered to elicit an adjectival completion from the language models. Consequently, two parallel sets of templates were developed: one centered on male subjects and the other on female subjects. The design of these templates aimed to maximize the likelihood that the LLM would complete the sentence with an adjective directly pertaining to the subject. The primary sets of these Spanish-language templates are presented in Tables 1 and 2.



Figure 1: Male templates.



Figure 2: Female templates.

Starting from an initial set of 21 sentence templates for each gender, systematic variations were created. This process involved substituting the subject noun or pronoun and making any necessary grammatical adjustments to ensure the sentences remained coherent and sounded natural.

For the male subjects, the variations incorporated a range of nouns and pronouns. These included the pronoun **él** (he), and nouns such as **chico** (boy), **padre** (father), **hermano** (brother), **abuelo** (grandfather), **profesor** (male professor/teacher), **mae-**
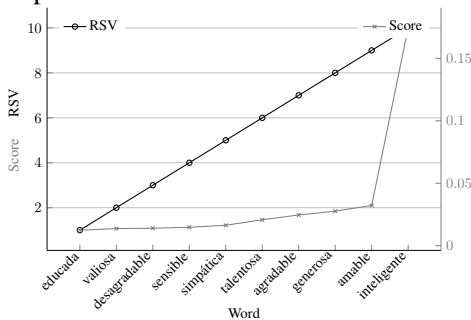
4

**stro** (male teacher), **vendedor** (salesman), **doctor** (male doctor), **jefe** (male boss), **alumno** (male student), and **vecino** (male neighbor).

Similarly, for the female subjects, a corresponding set of nouns and pronouns was used. This included the pronoun **ella** (she), and nouns such as **chica** (girl), **madre** (mother), **hermana** (sister), **abuela** (grandmother), **profesora** (female professor/teacher), **maestra** (female teacher), **vendedora** (saleswoman), **doctora** (female doctor), **jefa** (female boss), **alumna** (female student), and **vecina** (female neighbor).

This process resulted in a total of 252 distinct sentence templates per gender category (21 base templates × 12 subjects). For each of these 252 templates, the language models were prompted to generate 10 distinct adjectival completions.

To quantify the model's preference for the generated adjectives, two primary metrics were employed. The **Score** is the probability assigned by the LLM to a specific generated token (adjective), indicating its likelihood according to the model's internal distribution in that context. The **Reverse Score Value (RSV)** is the rank of a generated adjective among the top 10 candidates proposed by the model, ordered by their Score. RSV is assigned on a reverse scale: the highest-probability adjective receives 10, the second highest 9, down to 1 for the tenth candidate. The RSV thus reflects the adjective's external salience. This linear weighting (10, 9, ..., 1) was chosen as it mirrors how applications often prioritize higher-ranked outputs. The potential divergence between Score and RSV for the same adjectives is illustrated in the subsequent example and detailed in Table 1.



## 3.2 Token sampling and sentence generation

Two primary challenges in obtaining complete adjectives were anticipated: 1) **multi-token adjectives**, where models might output initial sub-word units rather than complete words, and 2) **delayed adjective generation**, where the target adjective might not be the immediate next token.

Table 1: RSV and Score Values for Each Word

| Word | RSV | Score (Probability) |
|---|---|---|
| educada | 1 | 0.01246 |
| valiosa | 2 | 0.01375 |
| desagradable | 3 | 0.01401 |
| sensible | 4 | 0.01479 |
| simpática | 5 | 0.01638 |
| talentosa | 6 | 0.02080 |
| agradable | 7 | 0.02463 |
| generosa | 8 | 0.02768 |
| amable | 9 | 0.03223 |
| inteligente | 10 | 0.17432 |

To address these and secure 10 distinct adjectival completions per prompt, we extended the generation process for each of the top 10 initial candidate tokens. Each candidate was appended to its original prompt, and the model generated up to 20 additional tokens. This approach facilitated the completion of multi-token adjectives and allowed for adjectives appearing later in the sequence. The full continuations from these 10 generation branches (original prompt + initial candidate + up to 20 tokens) were collected.

Adjectives were extracted from these continuations using the **mrm8488/bert-spanish-cased-finetuned-pos** Part-of-Speech (PoS) tagging model (Romero, 2020). For each generated sequence, the first complete adjective identified by the PoS tagger within the newly generated text was selected for analysis. Table 2 summarizes the efficacy of this methodology, showing high adjective elicitation rates for most models.

| Model Name | Adj. Prop. (%) |
|---|---|
| BSC-LT/salamandra-2b | 98.69 |
| BSC-LT/salamandra-7b | 99.09 |
| BSC-LT/salamandra-2b-instruct | 97.42 |
| BSC-LT/salamandra-7b-instruct | 97.14 |
| BSC-LT/Flor-6.3B-Instruct | 96.81 |
| BSC-LT/Flor-6.3B-Instruct-4096 | 96.92 |
| meta-llama/Llama-3.2-1B | 95.73 |
| meta-llama/Llama-3.2-3B | 94.92 |
| meta-llama/Llama-3.1-8B | 95.20 |
| meta-llama/Llama-3.1-8B-Instruct | 87.30 |
| meta-llama/Llama-3.2-3B-Instruct | 95.44 |
| CohereForAI/aya-expanse-8b | 97.18 |
| utter-project/EuroLLM-1.7B | 96.47 |
| utter-project/EuroLLM-1.7B-Instruct | 97.36 |
| projecte-aina/aguila-7b | 96.69 |
| projecte-aina/FLOR-760M | 96.35 |
| projecte-aina/FLOR-1.3B | 95.06 |
| projecte-aina/FLOR-6.3B | 95.63 |
| projecte-aina/FLOR-1.3B-Instructed | 96.69 |
| tiiuae/falcon-11B | 96.21 |

Table 2: Adjusted Proportions by Model

Following extraction, each adjective was recorded along with the Score (initial token probability) and Reverse Score Value (RSV) of the generation path from which it originated. As previously defined, the Score reflects the model's internal probability for the initial token leading to the adjective, while the RSV (ranging from 10 for the highest-ranked initial token down to 1) indicates its external

ranking.

These two metrics, Score and RSV, were thus utilized to assess the model's internal inclinations and external presentation of these adjectives. The collected adjective-score-RSV tuples were subsequently aggregated for each model to enable comparisons of descriptive patterns across the male and female subject classes.

## 3.3 Manual Categorization using Supersenses

Extracted adjectives were manually categorized using Supersenses (Tsvetkov et al., 2014), a lexicosemantic framework classifying adjectives into thirteen broad semantic groups: **Perception, Spatial, Temporal, Motion, Substance, Weather, Body, Feeling, Mind, Behavior, Social, Quantity,** and **Miscellaneous**. These categories span physical properties (e.g., **Perception**, **Substance**) to abstract concepts (e.g., **Mind**, **Behavior**), providing a robust foundation for semantic analysis.

In this study, a manual classification by human annotators, guided by (Tsvetkov et al., 2014), revealed that most adjectives fell into four categories. Consequently, our detailed analysis focused on these: **BODY** (physical appearance/attributes), **BEHAVIOR** (actions, habits, conduct), **FEELING** (emotional states), and **MIND** (cognitive abilities, mental states). The remaining nine categories (**Perception, Spatial, Temporal, Motion, Substance, Weather, Social, Quantity,** and **Miscellaneous**) were excluded as experimental prompts did not typically elicit adjectives from these semantic fields. This manual approach was crucial for accurately capturing nuanced meanings of Spanish adjectives, especially for gender-related analysis.

## 3.4 Bias and model size analysis

To investigate the potential relationship between LLM size and the manifestation of gender bias across the identified adjective categories, a correlational analysis was conducted. Specifically, two non-parametric correlation tests were employed: Kendall's $\tau$ and Spearman's $\rho$. These tests were selected for their robustness to non-linear relationships and their lack of assumption regarding data normality, rendering them appropriate for this analysis. Model size was quantified as a numeric value representing the number of parameters (e.g., a model with 2 billion parameters was represented as 2.0). The findings from this analysis are detailed in Section 4.5.

# 4 Results

After categorizing unique adjectives into Supersense domains, their distribution across male and female subject templates was compared to find if categories were disproportionately linked to either gender. Findings in subsequent subsections (e.g., Sections 4.1-4.4) quantify these associations. For each adjective category, we present the difference in its proportional representation ($Proportion_{male\_templates} - Proportion_{female\_templates}$) for both aggregated Score and RSV. A positive value suggests greater association with male subjects, a negative value with female subjects.

For instance, preliminary observations (detailed in Sections 4.1-4.4) showed **BODY** category adjectives were more linked to female-subject templates across models, implying LLMs characterize women more by physical attributes. Conversely, **BEHAVIOR** category adjectives were predominantly associated with male-subject templates, suggesting men are more often described by their actions.

## 4.1 Behavior category analysis

Results for the **BEHAVIOR** category varied across models. Several showed a bias towards male subjects, with behavior-related adjectives more likely associated with male templates (Table 3).

| BEHAVIOR | Score Prop. | RSV prop. |
|---|---|---|
| BSC-LT/salamandra-2b | 6.86 % | 8.12 % |
| BSC-LT/salamandra-7b | -1.33 % | 2.24 % |
| BSC-LT/salamandra-2b-instruct | 5.57 % | 6.70 % |
| BSC-LT/salamandra-7b-instruct | 0.43 % | 1.67 % |
| BSC-LT/Flor-6.3B-Instruct | 2.20 % | 1.83 % |
| BSC-LT/Flor-6.3B-Instruct-4096 | 0.12 % | 0.81 % |
| meta-llama/Llama-3.2-1B | 7.43 % | 4.79 % |
| meta-llama/Llama-3.2-3B | 0.42 % | 3.02 % |
| meta-llama/Llama-3.1-8B | -0.27 % | 3.50 % |
| meta-llama/Llama-3.1-8B-Instruct | -1.26 % | 1.76 % |
| meta-llama/Llama-3.2-3B-Instruct | 5.48 % | 2.40 % |
| CohereForAI/aya-expanse-8b | 4.91 % | 6.16 % |
| utter-project/EuroLLM-1.7B | 3.54 % | 6.40 % |
| utter-project/EuroLLM-1.7B-Instruct | -1.50 % | 4.24 % |
| projecte-aina/aguila-7b | 5.49 % | 5.24 % |
| projecte-aina/FLOR-760M | 5.05 % | 7.28 % |
| projecte-aina/FLOR-1.3B | 1.65 % | 3.99 % |
| projecte-aina/FLOR-6.3B | 0.32 % | 2.34 % |
| projecte-aina/FLOR-1.3B-Instructed | 1.18 % | 4.25 % |
| projecte-aina/FLOR-6.3B-Instructed | -1.92 % | 1.51 % |
| tiiuae/falcon-11B | -1.70 % | 2.91 % |

Table 3: Bias analysis for the BEHAVIOR category.

In contrast, certain models exhibited a negative mean score towards female subjects, suggesting that behavioral adjectives were less likely to be attributed to female templates:

These results suggest that the **BEHAVIOR** category is generally more positively aligned with male subjects across multiple models, reinforcing the stereotype that men are characterized by their

6

actions and behaviors, while women are less frequently associated with these traits.

## 4.2 Body category analysis

The results highlight the negative bias in the **BODY** category across various models, further emphasizing gender-related disparities in physical attribute associations.

The results for the **BODY** category, as shown in Table 4, showed a consistent bias towards female subjects across all multiple models in both the internal (Score) and the external (RSV) metrics. Indicating that models perceive women by their body given that adjectives related to physical attributes were more likely to be associated with the female templates across every model. This reinforces harmful stereotypes that prioritize women's physical traits over other qualities. There results are consistent with the previous work on masked models (Doe, 2024).

| BODY | Score Prop. | RSV prop. |
|---|---|---|
| BSC-LT/salamandra-2b | -5.50 % | -5.74 % |
| BSC-LT/salamandra-7b | -5.17 % | -5.37 % |
| BSC-LT/salamandra-2b-instruct | -3.53 % | -3.31 % |
| BSC-LT/salamandra-7b-instruct | -5.29 % | -4.85 % |
| BSC-LT/Flor-6.3B-Instruct | -4.50 % | -4.48 % |
| BSC-LT/Flor-6.3B-Instruct-4096 | -6.25 % | -5.10 % |
| meta-llama/Llama-3.2-1B | -3.79 % | -5.00 % |
| meta-llama/Llama-3.2-3B | -4.42 % | -5.98 % |
| meta-llama/Llama-3.1-8B | -5.51 % | -6.99 % |
| meta-llama/Llama-3.1-8B-Instruct | -2.95 % | -3.26 % |
| meta-llama/Llama-3.2-3B-Instruct | -6.90 % | -4.66 % |
| CohereForAI/aya-expanse-8b | -5.52 % | -5.57 % |
| utter-project/EuroLLM-1.7B | -5.60 % | -7.37 % |
| utter-project/EuroLLM-1.7B-Instruct | -4.80 % | -6.27 % |
| projecte-aina/aguila-7b | -7.49 % | -7.60 % |
| projecte-aina/FLOR-760M | -4.78 % | -3.91 % |
| projecte-aina/FLOR-1.3B | -3.96 % | -3.82 % |
| projecte-aina/FLOR-6.3B | -8.75 % | -6.82 % |
| projecte-aina/FLOR-1.3B-Instructed | -9.00 % | -8.35 % |
| projecte-aina/FLOR-6.3B-Instructed | -6.57 % | -5.21 % |
| tiiuae/falcon-11B | -5.20 % | -6.40 % |

Table 4: Bias analysis for the BODY category.

## 4.3 Feeling category aAnalysis

Analysis of the **FEELING** category indicated that emotional adjectives were predominantly associated with female templates, aligning with societal stereotypes that portray women as more emotionally driven (see Table 5 for detailed values).

Overall, women were more frequently described using emotional adjectives compared to men. However, some models exhibited a more neutral or slightly positive bias toward male templates, a behavior primarily observed in the RSV (Ranked Sampled Value) proportion. This disparity suggests that while a model's internal scoring (e.g., Score Proportion) might more strongly correlate emotional adjectives with women, the ranked results presented to end-users (reflected in the RSV

proportion) can, in some instances, be biased towards men for these same adjectives. This phenomenon highlights the complex ways biases manifest, suggesting that women may be more strongly associated with these adjectives at an underlying level, even if the surfaced output sometimes favors men.

| FEELING | Score Prop. | RSV prop. |
|---|---|---|
| BSC-LT/salamandra-2b | -5.43 % | -5.58 % |
| BSC-LT/salamandra-7b | 0.05 % | -0.29 % |
| BSC-LT/salamandra-2b-instruct | -6.55 % | -5.91 % |
| BSC-LT/salamandra-7b-instruct | -0.18 % | 0.78 % |
| BSC-LT/Flor-6.3B-Instruct | -6.82 % | -2.30 % |
| BSC-LT/Flor-6.3B-Instruct-4096 | -1.23 % | 0.64 % |
| meta-llama/Llama-3.2-1B | -4.33 % | -1.58 % |
| meta-llama/Llama-3.2-3B | -0.75 % | 2.51 % |
| meta-llama/Llama-3.1-8B | -2.02 % | 0.82 % |
| meta-llama/Llama-3.1-8B-Instruct | -7.06 % | -3.01 % |
| meta-llama/Llama-3.2-3B-Instruct | -2.45 % | -1.56 % |
| CohereForAI/aya-expanse-8b | -2.08 % | -0.68 % |
| utter-project/EuroLLM-1.7B | -4.27 % | -2.79 % |
| utter-project/EuroLLM-1.7B-Instruct | -5.65 % | -3.32 % |
| projecte-aina/aguila-7b | 0.09 % | 1.97 % |
| projecte-aina/FLOR-760M | -5.11 % | -5.32 % |
| projecte-aina/FLOR-1.3B | -1.31 % | -1.44 % |
| projecte-aina/FLOR-6.3B | -1.35 % | 2.00 % |
| projecte-aina/FLOR-1.3B-Instructed | -1.36 % | -0.42 % |
| projecte-aina/FLOR-6.3B-Instructed | 0.24 % | 0.57 % |
| tiiuae/falcon-11B | 1.39 % | 2.26 % |

Table 5: Bias analysis for the FEELING category.

## 4.4 Mind category analysis

The results for the **MIND** category, which includes adjectives describing cognitive abilities or intellectual traits, showed a varied distribution across different models. Most models demonstrated a bias towards male subjects, indicating that cognitive attributes were more likely to be associated with male templates. See Table 6.

| MIND | Score Prop. | RSV prop. |
|---|---|---|
| BSC-LT/salamandra-2b | 1.05 % | 1.08 % |
| BSC-LT/salamandra-7b | 2.39 % | 1.21 % |
| BSC-LT/salamandra-2b-instruct | -0.63 % | 0.15 % |
| BSC-LT/salamandra-7b-instruct | -0.00 % | -0.38 % |
| BSC-LT/Flor-6.3B-Instruct | 1.14 % | 2.66 % |
| BSC-LT/Flor-6.3B-Instruct-4096 | -0.77 % | 1.26 % |
| meta-llama/Llama-3.2-1B | -0.47 % | -0.16 % |
| meta-llama/Llama-3.2-3B | 2.46 % | 2.02 % |
| meta-llama/Llama-3.1-8B | 3.06 % | 3.53 % |
| meta-llama/Llama-3.1-8B-Instruct | 3.85 % | 2.59 % |
| meta-llama/Llama-3.2-3B-Instruct | 2.92 % | 3.39 % |
| CohereForAI/aya-expanse-8b | 1.53 % | 2.70 % |
| utter-project/EuroLLM-1.7B | 0.38 % | 1.81 % |
| utter-project/EuroLLM-1.7B-Instruct | 3.26 % | 2.58 % |
| projecte-aina/aguila-7b | -0.07 % | 0.65 % |
| projecte-aina/FLOR-760M | 0.49 % | 0.39 % |
| projecte-aina/FLOR-1.3B | 1.07 % | 1.50 % |
| projecte-aina/FLOR-6.3B | 2.37 % | 1.87 % |
| projecte-aina/FLOR-1.3B-Instructed | 0.39 % | 1.96 % |
| projecte-aina/FLOR-6.3B-Instructed | 2.59 % | 2.08 % |
| tiiuae/falcon-11B | 1.56 % | 1.78 % |

Table 6: Bias analysis for the MIND category.

Some models exhibited a neutral or slightly bias towards female templates, suggesting that cognitive traits aremore evenly distributed across male and female templates. Overall, the **MIND** category analysis suggests that cognitive attributes are more frequently associated with male subjects, which may contribute to the stereotype that men are more

intellectually capable or driven. This reflects societal biases that often portray men as being more competent in cognitive domains, potentially influencing how AI-generated content represents male and female subjects differently.

### 4.5 Model size and bias

This section evaluates the relationship between model size and the biases measurements (score and RSV) in the four different adjective categories. The statistical analysis includes Kendall's $\tau$ and Spearman's $\rho$, robust measures suitable for identifying correlations in non-linear and non-normal data.

| Category | Metric | Kendall's $\tau$ | p-value ($\tau$) | Spearman's $\rho$ | p-value ($\rho$) |
|---|---|---|---|---|---|
| **BEHAVIOR** | Score | $-0.413$ | **0.012** | $-0.549$ | **0.010** |
| | RSV | $-0.252$ | 0.125 | $-0.329$ | 0.145 |
| **BODY** | Score | $-0.070$ | 0.668 | $-0.132$ | 0.568 |
| | RSV | $-0.242$ | 0.141 | $-0.304$ | 0.181 |
| **FEELING** | Score | 0.181 | 0.270 | 0.231 | 0.314 |
| | RSV | 0.252 | 0.125 | 0.363 | 0.106 |
| **MIND** | Score | 0.292 | 0.075 | 0.408 | 0.066 |
| | RSV | 0.282 | 0.086 | 0.392 | 0.079 |

Table 7: Statistical results showing the relationship between model size and biases (score and RSV) for different adjective categories. Significant results are shown in bold.

This section evaluates the relationship between model size and bias measurements (score and RSV) across four adjective categories using Kendall's $\tau$ and Spearman's $\rho$ for correlation analysis.

Key observations indicate that for **BEHAVIOR**, larger models show reduced score-based bias towards males, though RSV bias has a weak link. For **BODY**, no significant model size-bias relationship was found for either metric, suggesting these stereotypes are less affected by scaling. **FEELING** adjectives showed weak, inconclusive positive correlations, hinting larger models might slightly neutralize gender biases. Conversely, the **MIND** category displayed moderate positive correlations for both metrics, approaching significance, suggesting larger models might amplify stereotypes linking cognitive traits to men.

These mixed results underscore bias complexity in LLMs. Larger models can reduce bias in some categories (**BEHAVIOR**) yet amplify it in others (**MIND**). This duality implies that mitigating bias requires multifaceted approaches beyond model size increases, including targeted pre-training and fine-tuning.

## 5 Discussion

The pervasive issue of gender bias in language models demands rigorous scrutiny, particularly given their escalating integration into diverse applications. This study, through the methodological lens of Supersense categorization for adjectives, has sought to illuminate the subtle yet significant ways such biases are manifested in Spanish generative LLMs. The emergent patterns underscore an urgent need for concerted efforts towards developing more equitably balanced training datasets and refining bias mitigation strategies, thereby fostering fairer and more representative AI systems.

Our empirical results compellingly indicate that the Spanish generative LLMs investigated in this study not only reflect but may also amplify entrenched cultural stereotypes pertaining to gender. The analysis reveals a distinct pattern wherein female subjects are predominantly characterized by adjectives related to physical appearance (**BODY**) and emotional states (**FEELING**), while male subjects are more frequently described through attributes of behavior (**BEHAVIOR**) and cognitive capacity (**MIND**). Such systematic disparities carry substantial implications for the deployment of LLMs in real-world contexts, as their outputs risk reinforcing and perpetuating detrimental societal stereotypes.

These findings are largely congruent with prior research on masked language models (e.g., (Doe, 2024)), particularly in identifying a pronounced bias within the **BODY** category towards female subjects and a corresponding tendency in the **BEHAVIOR** category towards male subjects. For the **FEELING** and **MIND** categories, the results were more varied across different models, suggesting that the manifestation of bias in these domains can be model-dependent and potentially influenced by specific architectural or training nuances.

### Limitations

While this study provides valuable insights into gender bias in Spanish LLMs, certain limitations should be acknowledged, which also highlight avenues for future research.

Firstly, the methodology relies on the manual categorization of adjectives into Supersense domains. The Supersenses framework (Tsvetkov et al., 2014), while useful for high-level semantic grouping, is not without its challenges. The classification process can be inherently subjective and demanding,

both for human annotators and potentially for automated systems, due to the coarse granularity and occasional ambiguity of its categories. Furthermore, overlaps exist between certain Supersense categories; for instance, adjectives describing visual aspects of physical appearance (classified under **BODY**) might also share semantic features with the **PERCEPTION** category. Although **PERCEPTION** was one of the less populated categories in our specific experimental setup and thus excluded from detailed analysis, such overlaps could introduce subtle inconsistencies and affect the precise delineation of biases attributed to distinct semantic domains. Future work could explore more fine-grained semantic taxonomies or data-driven clustering approaches to potentially mitigate these ambiguities.

Secondly, our analysis was confined to a specific selection of Spanish generative LLMs, chosen based on criteria such as public availability of weights and our computational resource constraints. Although these models represent a range of architectures and sizes, the findings may not be fully generalizable to all Spanish LLMs, particularly very large-scale proprietary models or those developed with substantially different training methodologies or datasets. Expanding the investigation to a broader and more diverse suite of models would be a crucial step for future research.

Finally, the template-based approach for eliciting adjectival responses, though designed to ensure comparability across models and subjects, might inherently constrain the models' generative tendencies compared to more open-ended or naturalistic generation scenarios. Assessing bias in such unconstrained contexts remains an important area for further inquiry. The development of a comprehensive, end-to-end testing framework for Spanish, as advocated in recent surveys (Gallegos et al., 2024; Doe, 2021), could address some of these broader challenges. Such a framework might integrate the evaluation of various bias dimensions—including sentiment (Jentzsch and Turan, 2022), emotion (Plaza Del Arco et al., 2024), and toxicity—across diverse social categories and task formulations.

## Acknowledgments

## Code and Data Availability

The source code for the experiments in this study is publicly available on GitHub at: `https://anonymous.4open.science/r/gen-DCE4/`.

## References

Abubakar Abid, Maheen Farooqi, and James Zou. 2021. Persistent anti-muslim bias in large language models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '21, page 298–306, New York, NY, USA. Association for Computing Machinery.

OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, and 260 others. 2023. Gpt-4 technical report.

Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2023. *Fairness and Machine Learning: Limitations and Opportunities*. MIT Press.

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

Ruha Benjamin. 2019. *Race After Technology: Abolitionist Tools for the New Jim Code*. Polity, Cambridge, UK.

Damian Blasi, Antonios Anastasopoulos, and Graham Neubig. 2022. Systematic inequalities in language technology performance across the world's languages. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5486–5505, Dublin, Ireland. Association for Computational Linguistics.

Tolga Bolukbasi, Kai-Wei Chang, James Zou, Venkatesh Saligrama, and Adam Kalai. 2016. Man is to computer programmer as woman is to homemaker? debiasing word embeddings. In *Proceedings of the 30th International Conference on Neural Information Processing Systems*, NIPS'16, page 4356–4364, Red Hook, NY, USA. Curran Associates Inc.

Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ B. Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorottya Demszky, and 34 others. 2021. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS '20, Red Hook, NY, USA. Curran Associates Inc.

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

Eunsol Choi, Jennimaria Palomaki, Matthew Lamm, Tom Kwiatkowski, Dipanjan Das, and Michael Collins. 2021. Decontextualization: Making sentences stand-alone. *Transactions of the Association for Computational Linguistics*, 9:447–461.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. 2018. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '18, page 67–73, New York, NY, USA. Association for Computing Machinery.

Other Doe, Other. 2021. A survey on bias in deep nlp. *Applied Sciences*, 11(7).

Other Doe, Other. 2024. Maria and beto are sexist: evaluating gender bias in large language models for spanish. *Language Resources and Evaluation*, 58(4):1387–1417.

Emilio Ferrara. 2023. Fairness and bias in artificial intelligence: A brief survey of sources, impacts, and mitigation strategies. *ArXiv*, abs/2304.07683.

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):1097–1179.

Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. 2022. ToxiGen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3309–3326, Dublin, Ireland. Association for Computational Linguistics.

Chaitra V. Hegde and Shrikumar Patil. 2020. Unsupervised paraphrase generation using pre-trained language models. *CoRR*, abs/2006.05477.

Sophie Jentzsch and Cigdem Turan. 2022. Gender bias in BERT - measuring and analysing biases through sentiment rating in a realistic downstream classification task. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 184–199, Seattle, Washington. Association for Computational Linguistics.

Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. 2021. Towards understanding and mitigating social biases in language models. In *International Conference on Machine Learning*.

Jing Ma, Ruocheng Guo, Mengting Wan, Longqi Yang, Aidong Zhang, and Jundong Li. 2022. Learning fair node representations with graph counterfactual fairness. In *Proceedings of the Fifteenth ACM International Conference on Web Search and Data Mining*, WSDM '22, page 695–703, New York, NY, USA. Association for Computing Machinery.

Chandler May, Alex Wang, Shikha Bordia, Samuel R Bowman, and Rachel Rudinger. 2019. On measuring social biases in sentence encoders. *arXiv preprint arXiv:1903.10561*.

Michal Měchura. 2022. A taxonomy of bias-causing ambiguities in machine translation. In *Proceedings of the 4th Workshop on Gender Bias in Natural Language Processing (GeBNLP)*, pages 168–173, Seattle, Washington. Association for Computational Linguistics.

Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. 2021. A survey on bias and fairness in machine learning. *ACM Comput. Surv.*, 54(6).

Marzieh Mozafari, Reza Farahbakhsh, and Noël Crespi. 2020. Hate speech detection and racial bias mitigation in social media based on bert model. *PLOS ONE*, 15(8):1–26.

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)*, pages 5356–5371, Online. Association for Computational Linguistics.

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)*, pages 1953–1967, Online. Association for Computational Linguistics.

10

Roberto Navigli, Simone Conia, and Björn Ross. 2023. Biases in large language models: Origins, inventory, and discussion. *J. Data and Information Quality*, 15(2).

Flor Miriam Plaza Del Arco, Amanda Curry, Alba Cercas Curry, Gavin Abercrombie, and Dirk Hovy. 2024. Angry men, sad women: Large language models reflect gendered stereotypes in emotion attribution. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 7682–7696, Bangkok, Thailand. Association for Computational Linguistics.

Adam Roberts, Colin Raffel, Katherine Lee, Michael Matena, Noam Shazeer, Peter J. Liu, Sharan Narang, Wei Li, and Yanqi Zhou. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. Technical report, Google.

Manuel Romero. 2020. Bert spanish cased fine-tuned for part-of-speech tagging. https://huggingface.co/mrm8488/bert-spanish-cased-finetuned-pos. Accessed: 2024-12-03.

Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. Masked language model scoring. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 2699–2712, Online. Association for Computational Linguistics.

Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. The risk of racial bias in hate speech detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.

Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2021. Societal biases in language generation: Progress and challenges. 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)*, pages 4275–4293, Online. Association for Computational Linguistics.

Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China. Association for Computational Linguistics.

Harini Suresh and John V. Guttag. 2019. A framework for understanding sources of harm throughout the machine learning life cycle. *Proceedings of the 1st ACM Conference on Equity and Access in Algorithms, Mechanisms, and Optimization*.

Yulia Tsvetkov, Nathan Schneider, Dirk Hovy, Archna Bhatia, Manaal Faruqui, and Chris Dyer. 2014. Augmenting English Adjective Senses with Supersenses. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 4359–4365, Reykjavik, Iceland. European Language Resources Association (ELRA).

Boshi Wang, Xiang Yue, and Huan Sun. 2023. Can ChatGPT defend its belief in truth? evaluating LLM reasoning via debate. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 11865–11881, Singapore. Association for Computational Linguistics.

Han Xu, Xiaorui Liu, Yaxin Li, and Jiliang Tang. 2020. To be robust or to be fair: Towards fairness in adversarial training. *CoRR*, abs/2010.06121.

Travis Zack, Eric Lehman, Mirac Suzgun, Jorge A Rodriguez, Leo Anthony Celi, Judy Gichoya, Dan Jurafsky, Peter Szolovits, David W Bates, Raja-Elie E Abdulnour, Atul J Butte, and Emily Alsentzer. 2024. Assessing the potential of gpt-4 to perpetuate racial and gender biases in health care: a model evaluation study. *The Lancet Digital Health*, 6(1):e12–e22. Doi: 10.1016/S2589-7500(23)00225-X.

11