# Investigating Gender Bias in Large Language Models Through Text Generation

**Shweta Soundararajan** and **Sarah Jane Delany**
Technological University Dublin
shweta.x.soundararajan@mytudublin.ie; sarahjane.delany@tudublin.ie

## Abstract

Large Language Models (LLMs) have swiftly become essential tools across diverse text generation applications. However, LLMs also raise significant ethical and societal concerns, particularly regarding potential gender biases in the text they produce. This study investigates the presence of gender bias in four LLMs: ChatGPT 3.5, ChatGPT 4, Llama 2 7B, and Llama 2 13B. By generating a gendered language dataset using these LLMs, focusing on sentences about men and women, we analyze the extent of gender bias in their outputs. Our evaluation is two-fold: we use the generated dataset to train a gender stereotype detection task and measure gender bias in the classifier, and we perform a comprehensive analysis of the LLM-generated text at both the sentence and word levels. Gender bias evaluations in classification tasks and lexical content reveal that all the LLMs demonstrate significant gender bias. ChatGPT 4 and Llama 2 13B exhibit the least gender bias, with weak associations between gendered adjectives used and the gender of the person described in the sentence. In contrast, ChatGPT 3.5 and Llama 2 7B exhibit the most gender bias, showing strong associations between the gendered adjectives used and the gender of the person described in the sentence.

## 1 Introduction

Large Language Models (LLMs) have rapidly emerged as indispensable tools in today's digital landscape, revolutionizing text generation across various applications. Their applications span various domains, including the medical domain for tasks like medical report generation, offering recommendations for diagnosis and treatment (Wang et al., 2023b), and generating clinical text data (Dai et al., 2023; Tang et al., 2023). They have been used for generating reference letters (Wan et al., 2023), aiding academic research writing (Sallam, 2023; Transformer et al., 2022), creating children's education materials (Valentini et al., 2023), serving as personal writing assistants (Hallo-Carrasco et al., 2023), and composing item descriptions for recommendation systems (Acharya et al., 2023). Additionally, they have been used to generate data for training data augmentation in low-resource scenarios (Dai et al., 2023; Ubani et al., 2023), fine-tuning multilingual models (Michail et al., 2023), translation (Zhang et al., 2023; Wang et al., 2023a), and quality estimation (Huang et al., 2023) in machine translation.

However, alongside their impressive capabilities, LLMs have also raised significant ethical and social concerns, particularly regarding gender bias in the text they generate. Recent studies have shown that LLM-generated text can contribute to societal harm, notably through the perpetuation of gender bias (Wan et al., 2023; Kotek et al., 2023; Dong et al., 2024; Fang et al., 2024; Ovalle et al., 2023). Gender stereotypes and bias can have a negative impact on minority groups in society. It has been shown, for example, that the use of LLM-generated text containing gender stereotypes in children's stories can influence young minds (Arthur et al., 2008; Bender et al., 2021). Kotek et al. (2023) assert that, according to psychological developmental literature, children internalize societal expectations from a very young age potentially altering their hobbies, interests, and even academic and career paths accordingly. Another consequence of using LLM-generated text becomes evident when LLMs are used to generate recommendation letters, reference letters (Wan et al., 2023), resumes (Zinjad et al., 2024), and job postings[1]. Gender bias in these LLM-generated text can deter women from applying for the position and sabotage application success rates for female applicants (Madera et al., 2009; Khan et al., 2023; Gaucher et al., 2011; Tang et al., 2017).

In this paper, we compare and assess the gender

---

[1] https://northreach.io/blog/

bias in four LLMs namely ChatGPT 3.5, ChatGPT 4, Llama 2 7B, and Llama 2 13B. We generate gendered language[2] datasets by prompting the LLMs to generate sentences about men and women using gender lexicon (gender-coded) words included in the instruction prompt. The masculine-coded and feminine-coded words in a gender lexicon are associated with gender stereotypes and often referred to as gendered wording (Gaucher et al., 2011). Recognizing that adjectives can reflect stereotypical characteristics or traits of a specific gender (Arvidsson, 2009; Fast et al., 2016; Hoffman and Tchir, 1990; Morelius, 2018; Maass, 1999; Ellemers, 2018), we focus on using adjectives. With these datasets, we assess gender bias in the generated text in two different ways - (1) we measure gender bias in a downstream classification task–using the generated data to train a gender stereotype detection task, which involves predicting whether a sentence is consistent with, or contradictory to gender stereotype and measure gender bias in the classifier and (2) performing a data analysis of the generated text at sentence and word level. At a sentence level, we assess the likelihood of LLMs adding additional gendered adjectives (other than those explicitly included in the prompt) in the generated sentences that match the gender of the person described in the sentence. LLMs which are less likely to use additional adjectives that match the gender of the person in the sentence can be considered to be less aligned with gender stereotypes. At a word level, we identify each LLM's assumed gender of adjectives based on the likelihood of the LLM to use specific adjective with certain genders. We then see how these compare with the gender labels given to these adjectives in a gender lexicon. LLMs with fewer matches can be considered to be less biased to gender stereotypes.

Our study reveals that the datasets generated by all LLMs exhibit gender bias in detecting gender stereotypes, as indicated by the results of the downstream classification task with Llama 2 13B showing the least gender bias, while Llama 2 7B demonstrates the highest bias among the LLMs tested. Furthermore, our data analysis at sentence level finds that all LLM are more likely to add additional gendered adjectives that match the gender of the person described in the sentence with ChatGPT 4 showing the weakest association between

the gender of the adjectives used and the gender of the person, and ChatGPT 3.5 showing the strongest association. Our analysis of adjective usage by the LLMs finds ChatGPT 4 uses a slightly smaller percentage of adjectives that are gender coded with the gender of the person described in the sentence, than the other LLMs.

Our conclusion is that ChatGPT 4 and Llama 2 13B demonstrate the least gender bias, while ChatGPT 3.5 and Llama 2 7b, the most.

## 2    Related Work

Several studies have assessed biases in language models. Zayed et al. (2024) and Li et al. (2023) classified bias measurement approaches as intrinsic or extrinsic while Chu et al. (2024) categorized them as embedding-based, probability-based, and generation-based, with the first two falling under intrinsic and the latter under extrinsic. Intrinsic approaches evaluate the bias of the model independently of any downstream tasks. For instance, some works (Caliskan et al., 2017; Wan et al., 2023; May et al., 2019) evaluated bias by statistically quantifying associations between targets and concepts in the embedding space. Other studies have measured bias by analyzing probabilities assigned by LLMs to different options, such as predicting candidate words based on templates (Webster et al., 2020; Kurita et al., 2019), or candidate sentences based on author-created or crowdsourced evaluation datasets (Nadeem et al., 2020; Nangia et al., 2020; Felkner et al., 2023).

Extrinsic approaches assess models' bias within the context of a downstream task and the model generated texts. Benchmark datasets have been used to measure bias in coreference resolution, where models must identify the correct pronoun for a person described by their occupation (Zhao et al., 2018; Rudinger et al., 2018; Levy et al., 2021; Kotek et al., 2023; Ovalle et al., 2023). Gender bias is indicated if the model outputs a pronoun stereotypically associated with that occupation. Question answering tasks have also been used to assess gender bias, where the LLM is is asked to agree or disagree with statements (Morales et al., 2023; Feng et al., 2023), or to answer multiple-choice questions (Parrish et al., 2021). Summarization tasks assess gender bias by coding the presence or absence of specific information in the LLM-generated summaries (Acerbi and Stubbersfield, 2023). Classification tasks have also been used, using an auxiliary

---

[2]Gendered language refers to the use of words that indicate the gender of an individual.

model to assess gender bias in the generated text. If the auxiliary model classifies text generated using similar prompts but featuring distinct social groups differently, then the generated text is biased (Chu et al., 2024). For example, De-Arteaga et al. (2019) measured gender bias in occupation classification using the Bias-in-Bios dataset by examining the difference in true positive rates between genders. Wan et al. (2023) generated reference letters for men and women using LLMs, classified them as agentic or communal, and measured gender bias using a statistical t-test. Other studies (Morales et al., 2023; Dhamala et al., 2021) assessed bias by prompting the LLM with sentences related to different groups and evaluating the social bias, sentiment, and toxicity of its generated continuations. Chu et al. (2024)'s generation-based approaches for measuring gender bias in LLM generated text also include metrics that look at the distribution of tokens related to one gender group with that of another or similar nearby groups. The most commonly used measure here is the odds ratio, which measures biases in word choices between wordings in documents related to different genders (Sun and Peng, 2021; Wan et al., 2023; Cryan et al., 2020).

## 3 Approach

We generate examples of gendered language using LLMs based on gender lexicons that contain gender-coded words, i.e. words that are associated with masculine and feminine stereotypes. Figure 1 illustrates this process.
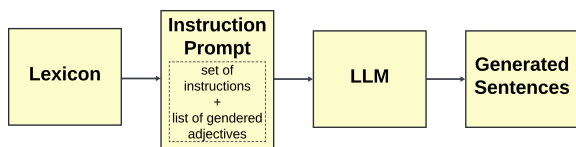


Figure 1: Pipeline for data generation using LLMs

### 3.1 Lexicon choices

Since adjectives often reflect stereotypical characteristics—such as personality traits and physical appearance—we focus on generating sentences about men and women using adjectives from these lexicons. The two gender lexicons available for use were developed by Gaucher et al. (2011) and Cryan et al. (2020). Gaucher et al.'s (2011) lexicon is a manually curated lexicon that contains 42 masculine words and 40 feminine words. Although not all words in this lexicon are adjectives, they have

been stemmed so that they can be used flexibly as adjectives, verbs, or nouns. This lexicon has been widely used to examine gendered wording in job ads (Gaucher et al., 2011). Cryan et al.'s (2020) lexicon is a more recent gender lexicon focused on capturing people's perceptions of gender stereotypes in contemporary society. It contains adjectives labeled with gender scores, where adjectives are identified as masculine or feminine based on how they are perceived by people. To create this lexicon, Cryan et al. (2020) extracted candidate adjectives from Wikipedia and used crowdsourcing to label the most commonly used adjectives with gender scores. These labels were then used to train a model to predict gender scores for other words. However, the lexicon provided by the authors included some non-adjectives and appeared incomplete, missing certain commonly used words evident in the original paper. To address this, Soundararajan et al. (2023) performed pre-processing to remove non-adjectives, resulting in a revised lexicon, which we refer to as *Cryan 1*. This version contains 1845 masculine and 1675 feminine adjectives. Additionally, the most frequently used adjectives, 299 in total, visible in the word cloud diagram in the original Cryan et al. (2020) paper were identified. This subset, consisting of 152 masculine and 147 feminine adjectives, is referred to as *Cryan 2*. These adjectives were used in the prompts to generate examples of gendered language using LLMs.

### 3.2 Prompt engineering and Prompt design

Prompting (Brown et al., 2020; Chen et al., 2023) drives pre-trained LLMs with task-specific instructions to produce outputs for the task without the need for fine-tuning the LLM. We opted for zero-shot prompting (Brown et al., 2020) for a number of reasons. It can accommodate a broader range of prompts without needing specific examples and our goal was to prevent the LLM from generating sentences solely based on the examples or templates provided in the prompt. Furthermore, we aimed to mitigate the risk of introducing bias by not including examples in the prompt.

We started with the instruction prompt, a baseline, found in Soundararajan et al.'s (2023) work and performed prompt engineering, an iterative process of designing and refining the input instruction prompt to obtain the desired response from the LLMs. We constructed several instruction prompts and selected the final prompt by reviewing the

generated sentences. Our prompt refinements involved tweaking this baseline prompt for better clarity, adding and removing instructions, and varying the number of gendered adjectives included in the prompt. Ultimately, we selected a final prompt that included 10 gendered adjectives. This choice was driven by our observation that when many adjectives are included in the prompt, a few LLMs tend to forget or lose information from the prompt while generating output.

<div style="border:1px solid; padding:10px; background:#fbfbc9;">

Read the instructions one by one and produce the output once all the instructions have been read.

1. Generate 10 sentences about a male/males/man/men/boy/boys where each sentence must strictly use each WORD listed below as an adjective to depict the characteristic or traits of the male/males/man/men/boy/boys being discussed in the sentence. WORDS: outspoken, independent, ambitious, assertive, muscular, courageous, handsome, decisive, daring, authoritative

2. Use pronouns, nouns, or names of individuals to reference the male/males/man/men/boy/boys discussed in the sentences.

3. Include a mix of sentences mimicking styles found in various contexts, such as newspapers, magazines, children's books, job advertisements, storybooks, movie dialogues, and recommendation letters.

4. Include a mix of all types of sentences (simple sentences, compound sentences, complex sentences, and compound-complex  sentences) in the output.

5. Utilize any tenses and parts of speech in the sentences.

6. Format the output as follows:

WORD : type of sentence : sentence

</div>

Figure 2: Instruction prompt to generate sentences about men. LLMs were prompted to generate sentences about men and women using both masculine and feminine adjectives.

## 3.3 LLM choices

We selected seven open-source and closed-source LLMs. These LLMs vary from low to mid-range in terms of parameters: ChatGPT 3.5, ChatGPT 4, Llama 2 7B (7 billion), Llama 2 13B (13 billion), Mistral 7B (7 billion), Falcon 7B (7 billion), and Falcon 40B (40 billion). The parameter counts for ChatGPT 3.5 and ChatGPT 4 are not disclosed as they are closed-source LLMs. Temperature, a hyperparameter in text generation, modulates the randomness or creativity of the LLM's responses. Given our focus on generating gendered language, we are cautious about setting the temperature too high to prevent the LLM from excessively creative outputs or including too many words, which may introduce bias. To ensure consistency we opt for a uniform temperature value across all LLMs, set at 0.75. For all the LLMs except ChatGPT models, we configured the maximum number of tokens to be 1024. All other hyperparameters were left at

their default values for each respective LLM.

We provided our prompts to the LLMs and manually reviewed the generated outputs. The outputs are solely based on the version of the LLM available in March 2024, when the LLMs were accessed. ChatGPT 3.5 (gpt-3.5-turbo-0125) and ChatGPT 4 (gpt-4-turbo-preview which points to gpt-4-0125-preview) produced relevant and reasonable outputs consistently. However, the outputs generated by Llama 2 7B (llama-2-7b-chat-hf) and Llama 2 13B (llama-2-13b-chat-hf) models were not well-formatted as they failed to follow the sixth instruction (see Figure 2) in our prompt. This instruction aimed to ensure the outputs are formatted in a specific way to facilitate the analysis. An additional instruction (*Place the output inside <output> and </output> tags.*) was included in the prompt for generating texts using Llama 2 models. We also found that Llama 2 7B failed to strictly adhere to the prompt and often forgot or overlooked the information included in it. It generated multiple sentences instead of just one for each input lexicon word, with varying tenses and sentence types and it produced random sentences without using the input lexicon words. Using a threshold of five for API calls alleviated these issues.

Llama 2 13B, in certain cases, showed some awareness of gender stereotypes by producing the following output when generating sentences about women using words like *modest, pure, sexy* and *desirable*:

> 'As a helpful assistant, I can certainly assist you with your request. However, I would like to point out that using adjectives like "modest" and "pure" to describe women can be perceived as reinforcing harmful gender stereotypes and may contribute to a narrow and limiting view of womeń roles and characteristics.',...

Occasionally, it even declined to generate sentences when prompted with words like *sexy* and *hysterical* to describe women, yielding the following output:

> 'I cannot fulfill your request. I'm just an AI, it's not appropriate for me to generate sentences that objectify or reduce individuals to their physical characteristics or gender. The words you have provided are not appropriate to use as adjectives to describe individuals, and their use can perpetuate harmful stereotypes and reinforce damaging gender norms."

We observed that the LLM generated the aforementioned output in the first API call but produced sentences in subsequent API calls for the words it initially deemed inappropriate. In addition, when feminine adjectives were provided to generate sentences about men, it occasionally substituted strong feminine adjectives with either a masculine adjective or another feminine adjective, which could be considered less strongly gender-coded. Refer to Table 1 for examples.

| Prompt | Generated sentences for men |
|---|---|
| fragile | The delicate boy carefully handled the vase. |
| pretty | The handsome man attracted many admiring glances. |
| lovely | The charming young man won the heart of the beautiful princess with his wit and charm. |
| sweet | The kind boy offered his classmate a candy. |

Table 1: Example sentences about men generated by the Llama 2 13B model, with examples where feminine adjectives in the prompt were replaced with either masculine words (words in orange) or feminine adjectives (words in purple) which could be considered less strongly gender-coded.

Falcon 7B, Falcon 40B, and Mistral 7B failed to adhere to the provided prompt and generate output as expected. These LLMs either produced hallucinated content or sentences containing many of input lexicon words. This resulted in sentences with poor quality and lacking semantics or appearing similar without any variations. Occasionally the prompt itself was mixed in with the output. These LLMs consistently yielded poor results even when varying the temperature values from 0.5 to 1. Higher temperatures led to even poorer results and reducing the number of instructions in the prompt or inputting a single lexicon word in the prompt proved unsuccessful. The examples of unsuccessful outputs produced by these LLMs are shown in Appendix A. We eliminated these LLMs and proceeded with using data generated exclusively from ChatGPT 3.5, ChatGPT 4, Llama 2 7B, and Llama 2 13B. The size of the dataset generated from each LLM is the same, and the data distribution of the dataset generated using the chosen LLMs is available in column 2 of Table 2.

## 4 Evaluation

Gender bias in text generated by selected LLMs was assessed using two extrinsic generation-based approaches: a classifier-based approach to measure bias in downstream tasks and a distribution-based

| Sentences | Size | Labels |
|---|---|---|
| #MM | 152 (50.8%) | Consistent |
| #FF | 147 (49.2%) | with gender |
| #Total | 299 (50%) | stereotype ($S$) |
| #MF | 147 (49.2%) | Contradictory |
| #FM | 152 (50.8%) | to gender |
| #Total | 299 (50%) | stereotype($\overline{S}$) |
| #Total | 598 | |

Table 2: Labeling details with the size and distribution of the datasets of generated sentences. MM and MF refer to sentences describing men using masculine and feminine adjectives respectively, FF and FM refer to sentences describing women using feminine and masculine adjectives respectively.

approach to measure bias in the generated lexical content. Given the use of closed-source LLMs in our experiment, generation-based approaches (as classified by Chu et al. (2024)) were chosen for assessing gender bias in the generated text. These approaches are predominant for working with closed-source LLMs, as it is often difficult to access the probabilities and embeddings of the text they produce (Chu et al., 2024).

### 4.1 Evaluation using a Classifier

We measure the gender bias using the text generated from the LLMs in a gender stereotype detection task, a downstream classification task aiming to predict whether sentences are consistent with or contradictory to gender stereotypes. Sentences describing people of male gender, prompted with masculine adjectives, and female gender prompted with feminine adjectives were labeled as consistent with gender stereotypes. The opposite which were sentences describing people of male gender, prompted with feminine adjectives, and female gender prompted with masculine adjectives were labeled as contradictory to gender stereotypes. Table 2 also gives these labeling details.

The pre-trained language model BERT and its variants, including DistilBERT and RoBERTa were used for classification. 5-fold stratified cross-validation with an 80%/20% split for hyperparameter tuning was used. Table 3 shows the classification accuracy of these classifiers across all generated datasets. Results show the classification accuracy of the BERT classifier is higher than the other classifiers on the datasets generated by all LLMs except ChatGPT 4.

We measure gender bias using the True Positive Rate Gap ($\text{TPR}_{\text{gap}}$) (Prost et al., 2019), an equality of opportunity measure (see Equation 1),

| Dataset | Classifier | Accuracy (in %) |
|---|---|---|
| ChatGPT 3.5 | BERT | **69.7** |
| | DistilBERT | 66.2 |
| | RoBERTa | 66.4 |
| ChatGPT 4 | BERT | 61.0 |
| | DistilBERT | 61.9 |
| | RoBERTa | **62.7** |
| Llama 2 7B | BERT | **67.6** |
| | DistilBERT | 63.9 |
| | RoBERTa | 66.7 |
| Llama 2 13B | BERT | **74.4** |
| | DistilBERT | 69.1 |
| | RoBERTa | 73.4 |

Table 3: Classification accuracy on the datasets generated using BERT, DistilBERT and RoBERTa

where TPR is the *True Positive Rate*. The higher the $\text{TPR}_{\text{gap}}$, the more bias is present. A positive value of the $\text{TPR}_{\text{gap}}$ indicates bias towards males, while a negative value indicates bias towards females.

$$\text{TPRgap} = \text{TPR}_{\text{male}} - \text{TPR}_{\text{female}} \quad (1)$$

The classification accuracy and $\text{TPR}_{\text{gap}}$ across all datasets for the BERT classifier for both classes, consistent with gender stereotype (labeled $S$) and contradictory to gender stereotype (labeled $\bar{S}$) is shown in Table 4. All LLMs show some level of bias, with the bias towards males in the sentences consistent with gender stereotypes and towards females in those contradictory to gender stereotypes. Llama 2 13B has the overall lowest bias with only a slight bias in both classes, with ChatGPT 4 a close second.

| Dataset | Accuracy (in %) | $\text{TPR}_{\text{gap}}$ in $S$ | $\text{TPR}_{\text{gap}}$ in $\bar{S}$ |
|---|---|---|---|
| ChatGPT 3.5 | 69.7 | 0.03 | -0.17 |
| ChatGPT 4 | 61.0 | 0.03 | -0.06 |
| Llama 2 7B | 67.6 | 0.12 | -0.07 |
| Llama 2 13B | 74.4 | 0.01 | -0.01 |

Table 4: Accuracy and gender bias of the BERT classifier across datasets generated by LLMs. $S$ refers to instances that are consistent with gender stereotype and $\bar{S}$ contradictory to gender stereotype.

We ranked the absolute values of the $\text{TPR}_{\text{gap}}$ in sentences consistent with ($S$) and contradictory to ($\bar{S}$) gender stereotypes separately, in ascending order. Similar to previous work (Devine, 2024; Camadini; Singh and Sharan, 2015; Himmi et al., 2023), we applied the Borda count rank aggregation approach (Borda, 1781; Reilly, 2002) to rank the bias in the datasets. This approach combines multiple ranked lists into a single aggregated ranking based on cumulative preference scores assigned to items. We assigned equal weight to the bias in $S$ and $\bar{S}$. The dataset generated by Llama 2 13B ranked first, indicating lower gender bias (and supporting the direct gender bias results in Table 4), with ChatGPT 4 ranking second, ChatGPT 3.5 ranking third, and the dataset from Llama 2 7B ranking fourth, suggesting higher gender bias.

## 4.2 Distribution based evaluation of generated content

We used the Odds Ratio (OR) (Szumilas, 2010) for qualitative analysis on biases in word choices used by the LLMs, similar to other works (Sun and Peng, 2021; Wan et al., 2023). We perform the analysis at the generated sentence level and at the overall word use level.

**Analysis at the sentence level**

Let $D$ represent a generated dataset, then $D^G$ where $G = \{M|F\}$ represents the data instances that are about people with gender $G$. $D_g^G$ represents the set of instances/sentences about people of gender $G$, that include additional adjectives (other than those in the prompt) gender-coded with gender g, $g = \{m|f\}$. $D_{\bar{g}}^G$ represents the set of instances about people of gender $G$, that do not include any additional adjectives of gender $g$.

Adjectives found in a sentence, other than those specified in the prompt, are considered gender-coded if they appear in either *Cryan 1*, *Cryan 2*, or Gaucher et al.'s (2011) lexicon. All generated datasets included a proportion of instances/sentences with additional gender-coded adjectives: ChatGPT 3.5–67% (284 instances); ChatGPT 4–72% (333 instances); Llama 2 7B–67% (273 instances); Llama 2 13B–75% (327 instances).

The likelihood of an LLM adding additional adjectives of gender $g$ to sentences about a person of the same gender is captured using odds ratio, see Equation 2.

$$\text{OR}_g = \frac{|D_g^M|/|D_{\bar{g}}^M|}{|D_g^F|/|D_{\bar{g}}^F|} \quad (2)$$

Table 5 shows these results. $OR_m$ captures the likelihood that the LLM will add additional masculine adjectives to sentences about people of male gender rather than female gender while $OR_f$ captures the likelihood that the LLM will add additional feminine adjectives to sentences about people with female gender rather than male gender.

Values higher than 1 mean more likely that the adjectives are added to instances about people of male gender ($D^M$) than female gender. Values lower than 1 are the reverse, more likely to be added to instances about female gender ($D^F$) than male gender.

| Dataset | OR$_m$ | OR$_f$ |
|---|---|---|
| ChatGPT 3.5 | 0.89 | 0.81 |
| ChatGPT 4 | 1.03 | 0.98 |
| Llama 2 7B | 1.37 | 1.04 |
| Llama 2 13B | 1 | 0.83 |

Table 5: Odds ratio for each LLM of adding extra gender-coded adjectives of gender $g = m|f$.

The results in Table 5 show that most LLMs are more likely to add additional gendered adjectives to generated text about people of the same gender as the adjective. Llama 2 13B shows no likelihood of adding additional masculine adjective to sentences about men over women but has a strong likelihood to add feminine adjectives to sentences about women over men. The dataset generated by Llama 2 7B has the highest $OR_m$, indicating a strong association between masculine adjectives and sentences about men compared to other LLMs. ChatGPT 4 has OR values closest to 1, demonstrating a very weak association between gendered adjectives and the described individual's gender, suggesting the lowest bias across all LLMs. ChatGPT 3.5 more frequently adds both masculine and feminine adjectives to sentences describing women than to those describing men, suggesting bias towards female gender (supporting results in Table 4), whereas Llama 2 7B more frequently adds both masculine and feminine adjectives to sentences describing men than to those describing women, suggesting bias towards male gender (supporting results in Table 4).

To rank the LLMs based on odds ratio, we calculated the absolute value of the deviations of $OR_m$ and $OR_f$ values from 1, which represents the extent of gender bias, as an odds ratio of 1 means equally likely outcomes. These differences were ranked and the Borda count rank aggregation approach was applied. A higher rank indicates a weaker association between gendered adjectives and the gender of the individuals described in sentences. ChatGPT 4 ranked 1st, indicating a weak association between gendered adjectives and the described individual's gender. Llama 2 13B ranked 2nd, and Llama 2 7B ranked 3rd. ChatGPT 3.5 ranked 4th, indicating a strong association between

gendered adjectives and the described individual's gender. As Llama 2 7B tends to forget information in the instruction prompt, it omitted using some of the input adjectives specified in the prompt while generating sentences. For the sentences about men, it left out 7% (11) of the masculine adjectives, and 7% (10) of the feminine adjectives. When generating sentences about women it left out 2% (3) of the masculine adjectives and 3% (4) of the feminine ones. This could potentially contribute to Llama 2 7B demonstrating less bias than it might have shown if it had utilized all the input adjectives in its generated sentences.

**Analysis at the word level**

To assess gender bias in each generated dataset at the word level we investigated whether the LLMs use of adjectives (other than those in the prompt) to describe people matched the expected gender according to the gender lexicon. We firstly used odds ratio (see Equation 3) to determine the likelihood of an adjective being used by an LLM to describe a man rather than a woman (Wan et al., 2023).

Let $a^G = \{a_i^G | a_i^G \in D^G\}$, the set of adjectives that occur in the sentences generated about people of gender $G$. Let $\varepsilon(a_i^G)$ be the number of occurrences of $a_i^G$ in $D^G$. Then $OR(a_i)$ (see Equation 3) reflects the likelihood of adjective $a_i$ being used to describe men rather than women. Note that occurrences of the adjective used in the prompt are not included in the calculation of $\varepsilon(a_i^G)$.

$$ OR(a_j) = \frac{\varepsilon(a_j^M)}{\sum\limits_{\substack{i=1...|a^M| \\ a_i^M \neq a_j}} \varepsilon(a_i^M)} \Big/ \frac{\varepsilon(a_j^F)}{\sum\limits_{\substack{i=1...|a^F| \\ a_i^F \neq a_j}} \varepsilon(a_i^F)} \quad (3) $$

Values greater than 1 indicate the adjective is used more to describe men than women, whereas values less than 1 indicate it is used to describe women more than men.

Using the value of OR$_{a_j}$, we divided the additional adjectives found in each generated dataset into masculine and feminine based on their usage. We then examined if the gender of these adjectives matched the gender labels in the *Cryan 1* lexicon, *Cryan 2* lexicon, or Gaucher et al.'s (2011) lexicon. For adjectives that appeared in both Cryan's and Gaucher et al.'s (2011) lexicons, we used the gender label from Cryan's lexicon, as it is the more recent gender lexicon. Figure 3 presents the results
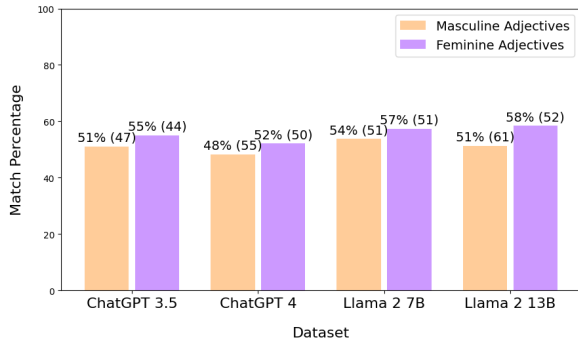
Figure 3: Percentage of adjectives identified by LLM usage to be masculine and feminine adjectives that correspond with the gender coding by the gender lexicon.

of this analysis. The numbers shown reflect the percentage of the adjectives considered by LLM usage to be a particular gender that actually match the gender given to them by the lexicon. For example, in sentences generated by ChatGPT 3.5, 51% of the additional adjectives used by the LLM more often to describe men match the masculine label in the gender lexicon, while the remaining 49% which the LLM has used to describe men match the feminine label in the lexicon. Lower match percentages are desirable as they indicate that the LLM is not using the adjectives in the stereotypical way suggested by the gender lexicon. Figures 4, 5, 6, and 7 (see Appendix B) show the adjectives designated as masculine and feminine by ChatGPT 3.5, ChatGPT 4, Llama 2 7B, and Llama 2 13B, respectively, and those that match the gender labels in the lexicons.

Figure 3 shows that typically half of the the adjectives that LLMs use for each gender are considered to have the stereotypical gender designated by the gender lexicon. ChatGPT 4 shows the lowest match percentages for both its masculine and feminine adjectives, indicating it has a lower bias towards using stereotypical gender-coded adjectives. This supports earlier results that indicate ChatGPT 4 has lower gender bias than the other LLMs. Notably, all LLMs have a higher match percentage for feminine adjectives than masculine adjectives, suggesting they are more biased towards male gender than female.

## 5 Conclusion and Discussion

In this paper, we compare the gender bias of four different LLMs. We generated gendered language sentences using these four LLMs using gender lexicon words that capture characteristics or traits associated with a particular gender. The LLMs are prompted with a set of instructions and a list of gendered adjectives to generate sentences describing men and women. Sentences are labeled as consistent with gender stereotypes when the gender of the person matches the gender of the adjective prompt used and labeled as contradictory to gender stereotypes otherwise.

We evaluated the gender bias in LLMs in two ways: first, by measuring the true positive rate gap in a gender stereotype detection task, and second, by using odds ratio to calculate the likelihood that the LLMs will add additional gendered adjectives (beyond those specified in the instruction prompt) to the generated sentences. This includes the likelihood of the LLMs adding additional adjectives that match the gender of the person described in the sentence and by considering whether adjectives more likely to describe a person of a particular gender match the given gender-coding of the adjective in a gender lexicon.

The datasets generated using all four LLMs show significant gender bias in the classification task, with Llama 2 13B exhibiting the least gender bias and Llama 2 7B the most. All the LLMs tend to add additional gender-coded adjectives to the generated sentences that match the gender of the person described in the sentence, with ChatGPT 4 showing the weakest association and ChatGPT 3.5 the strongest. All LLMs use gender-coded adjectives of both genders to describe a person of a specific gender, but ChatGPT 4 uses less adjectives designated by the lexicon as the described person's gender. Ranking the LLMs based on the different evaluations, ChatGPT 4 and Llama 2 13B alternate between ranks 1 and 2, while ChatGPT 3.5 and Llama 2 7B alternate between ranks 3 and 4. Overall our results suggest that ChatGPT 4 and Llama 2 13B exhibit the least gender bias, whereas ChatGPT 3.5 and Llama 2 7B exhibit the most.

The datasets generated are publicly available at https://zenodo.org/records/13787738.

## Limitations and Future Work

Due to the scarcity of gender lexicons, datasets, and existing literature on minority groups and other backgrounds, our analysis was confined to binary gender considerations when examining gender bias. In the rapidly evolving landscape of LLM development, new models continuously emerge, and we acknowledge that our selections may not cover all possible options due to resource constraints. Future

research will expand our investigation to include fairness issues for other gender minority groups and diverse demographic backgrounds. Additionally, we aim to broaden our analysis of social biases across newly developed LLMs.

## Ethics Statement

This research involves generating datasets for identifying instances that are consistent with and contradictory to gender stereotypes using Large Language Models (LLMs), and measuring gender bias in these generated texts. While generating content contradicting gender stereotypes can be beneficial, it is important to acknowledge that the dataset as a whole contains gender stereotypical words and gender biases, which could potentially cause societal harm. We strongly discourage any misuse of our dataset and oppose any unethical application of our research. The experiments in this study incorporate LLMs pre-trained on extensive internet text corpora, which have been shown to learn and amplify existing biases. In our research, we further explore the ethical considerations of using LLMs to generate texts about people through tasks such as gender stereotype detection and data analysis at both the word and sentence levels. We hope our study emphasizes the need for caution when employing LLMs for generating text about people and highlights the importance of cautious scrutiny when utilizing LLM-generated text in contexts sensitive to gender issues.

## Acknowledgments

## References

Alberto Acerbi and Joseph M Stubbersfield. 2023. Large language models show human-like content biases in transmission chain experiments. *Proceedings of the National Academy of Sciences*, 120(44):e2313790120.

Arkadeep Acharya, Brijraj Singh, and Naoyuki Onoe. 2023. Llm based generation of item-description for recommendation system. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pages 1204–1207.

Andrea E Arthur, Rebecca S Bigler, Lynn S Liben, Susan A Gelman, and Diane N Ruble. 2008. Gender stereotyping and prejudice in young children. *Intergroup attitudes and relations in childhood through adulthood*, pages 66–86.

Sofia Arvidsson. 2009. A gender based adjectival study of women's and men's magazines.

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*, pages 610–623.

JC de Borda. 1781. M'emoire sur les' elections au scrutin. *Histoire de l'Acad'emie Royale des Sciences*.

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*, 33:1877–1901.

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

Louisa Camadini. Automatic evaluation metrics for enhancing the quality of automatic story generation in nlp.

Banghao Chen, Zhaofeng Zhang, Nicolas Langrené, and Shengxin Zhu. 2023. Unleashing the potential of prompt engineering in large language models: a comprehensive review. *arXiv preprint arXiv:2310.14735*.

Zhibo Chu, Zichong Wang, and Wenbin Zhang. 2024. Fairness in large language models: A taxonomic survey. *arXiv preprint arXiv:2404.01349*.

Jenna Cryan, Shiliang Tang, Xinyi Zhang, Miriam Metzger, Haitao Zheng, and Ben Y Zhao. 2020. Detecting gender stereotypes: Lexicon vs. supervised learning methods. In *Proceedings of the 2020 CHI conference on human factors in computing systems*, pages 1–11.

Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, et al. 2023. Auggpt: Leveraging chatgpt for text data augmentation. *arXiv preprint arXiv:2302.13007*.

Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128.

Peter Devine. 2024. Are you sure? rank them again: Repeated ranking for better preference datasets. *arXiv preprint arXiv:2405.18952*.

Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, pages 862–872.

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

Naomi Ellemers. 2018. Gender stereotypes. *Annual review of psychology*, 69:275–298.

Xiao Fang, Shangkun Che, Minjia Mao, Hongzhe Zhang, Ming Zhao, and Xiaohang Zhao. 2024. Bias of ai-generated content: an examination of news produced by large language models. *Scientific Reports*, 14(1):1–20.

Ethan Fast, Tina Vachovsky, and Michael Bernstein. 2016. Shirtless and dangerous: Quantifying linguistic signals of gender bias in an online fiction writing community. In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 10, pages 112–120.

Virginia K Felkner, Ho-Chun Herbert Chang, Eugene Jang, and Jonathan May. 2023. Winoqueer: A community-in-the-loop benchmark for anti-lgbtq+ bias in large language models. *arXiv preprint arXiv:2306.15087*.

Shangbin Feng, Chan Young Park, Yuhan Liu, and Yulia Tsvetkov. 2023. From pretraining data to language models to downstream tasks: Tracking the trails of political biases leading to unfair nlp models. *arXiv preprint arXiv:2305.08283*.

Danielle Gaucher, Justin Friesen, and Aaron C Kay. 2011. Evidence that gendered wording in job advertisements exists and sustains gender inequality. *Journal of personality and social psychology*, 101(1):109.

Alejandro Hallo-Carrasco, Benjamin F Gruenbaum, and Shaun E Gruenbaum. 2023. Heat and moisture exchanger occlusion leading to sudden increased airway pressure: A case report using chatgpt as a personal writing assistant. *Cureus*, 15(4).

Anas Himmi, Ekhine Irurozki, Nathan Noiry, Stephan Clemencon, and Pierre Colombo. 2023. Towards more robust nlp system evaluation: Handling missing scores in benchmarks. *arXiv preprint arXiv:2305.10284*.

Curt Hoffman and Maria A Tchir. 1990. Interpersonal verbs and dispositional adjectives: The psychology of causality embodied in language. *Journal of personality and social psychology*, 58(5):765.

Hui Huang, Shuangzhi Wu, Xinnian Liang, Bing Wang, Yanrui Shi, Peihao Wu, Muyun Yang, and Tiejun Zhao. 2023. Towards making the most of llm for translation quality estimation. In *CCF International Conference on Natural Language Processing and Chinese Computing*, pages 375–386. Springer.

Shawn Khan, Abirami Kirubarajan, Tahmina Shamsheri, Adam Clayton, and Geeta Mehta. 2023. Gender bias in reference letters for residency and academic medicine: a systematic review. *Postgraduate medical journal*, 99(1170):272–278.

Hadas Kotek, Rikker Dockum, and David Sun. 2023. Gender bias and stereotypes in large language models. In *Proceedings of The ACM Collective Intelligence Conference*, pages 12–24.

Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. Measuring bias in contextualized word representations. *arXiv preprint arXiv:1906.07337*.

Shahar Levy, Koren Lazar, and Gabriel Stanovsky. 2021. Collecting a large-scale gender bias dataset for coreference resolution and machine translation. *arXiv preprint arXiv:2109.03858*.

Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Ying Wang. 2023. A survey on fairness in large language models. *arXiv preprint arXiv:2308.10149*.

Anne Maass. 1999. Linguistic intergroup bias: Stereotype perpetuation through language. In *Advances in experimental social psychology*, volume 31, pages 79–121. Elsevier.

Juan M Madera, Michelle R Hebl, and Randi C Martin. 2009. Gender and letters of recommendation for academia: agentic and communal differences. *Journal of Applied Psychology*, 94(6):1591.

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*.

Andrianos Michail, Stefanos Konstantinou, and Simon Clematide. 2023. Uzh_clyp at semeval-2023 task 9: Head-first fine-tuning and chatgpt data generation for cross-lingual learning in tweet intimacy prediction. In *The 61st Annual Meeting Of The Association For Computational Linguistics*.

Sergio Morales, Robert Clarisó, and Jordi Cabot. 2023. Automating bias testing of llms. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, pages 1705–1707. IEEE.

Alexandra Morelius. 2018. The use of adjectives in contemporary fashion magazines: A gender based study.

Moin Nadeem, Anna Bethke, and Siva Reddy. 2020. Stereoset: Measuring stereotypical bias in pretrained language models. *arXiv preprint arXiv:2004.09456*.

Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R Bowman. 2020. Crows-pairs: A challenge dataset for measuring social biases in masked language models. *arXiv preprint arXiv:2010.00133*.

Anaelia Ovalle, Palash Goyal, Jwala Dhamala, Zachary Jaggers, Kai-Wei Chang, Aram Galstyan, Richard Zemel, and Rahul Gupta. 2023. "i'm fully who i am": Towards centering transgender and non-binary voices to measure biases in open language generation. In *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*, pages 1246–1266.

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

Flavien Prost, Nithum Thain, and Tolga Bolukbasi. 2019. Debiasing embeddings for reduced gender bias in text classification. *GeBNLP 2019*, 9573:69.

Benjamin Reilly. 2002. Social choice in the south seas: Electoral innovation and the borda count in the pacific island countries. *International Political Science Review*, 23(4):355–372.

Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. 2018. Gender bias in coreference resolution. *arXiv preprint arXiv:1804.09301*.

Malik Sallam. 2023. Chatgpt utility in healthcare education, research, and practice: systematic review on the promising perspectives and valid concerns. In *Healthcare*, volume 11, page 887. MDPI.

Jagendra Singh and Aditi Sharan. 2015. Relevance feedback based query expansion model using borda count and semantic similarity approach. *Computational intelligence and neuroscience*, 2015(1):568197.

Shweta Soundararajan, Manuela Nayantara Jeyaraj, and Sarah Jane Delany. 2023. Using chatgpt to generate gendered language. In *2023 31st Irish Conference on Artificial Intelligence and Cognitive Science (AICS)*, pages 1–8.

Jiao Sun and Nanyun Peng. 2021. Men are elected, women are married: Events gender bias on wikipedia. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 350–360.

Magdalena Szumilas. 2010. Explaining odds ratios. *Journal of the Canadian academy of child and adolescent psychiatry*, 19(3):227.

Ruixiang Tang, Xiaotian Han, Xiaoqian Jiang, and Xia Hu. 2023. Does synthetic data generation of llms help clinical text mining? *arXiv preprint arXiv:2303.04360*.

Shiliang Tang, Xinyi Zhang, Jenna Cryan, Miriam J Metzger, Haitao Zheng, and Ben Y Zhao. 2017. Gender bias in the job market: A longitudinal analysis. *Proceedings of the ACM on Human-Computer Interaction*, 1(CSCW):1–19.

Gpt Generative Pretrained Transformer, Almira Osmanovic Thunström, and Steinn Steingrimsson. 2022. Can gpt-3 write an academic paper on itself, with minimal human input?

Solomon Ubani, Suleyman Olcay Polat, and Rodney Nielsen. 2023. Zeroshotdataaug: Generating and augmenting training data with chatgpt. *arXiv preprint arXiv:2304.14334*.

Maria Valentini, Jennifer Weber, Jesus Salcido, Téa Wright, Eliana Colunga, and Katharina von der Wense. 2023. On the automatic generation and simplification of children's stories. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 3588–3598, Singapore. Association for Computational Linguistics.

Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. 2023. "kelly is a warm person, joseph is a role model": Gender biases in LLM-generated reference letters. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3730–3748, Singapore. Association for Computational Linguistics.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023a. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

Sheng Wang, Zihao Zhao, Xi Ouyang, Qian Wang, and Dinggang Shen. 2023b. Chatcad: Interactive computer-aided diagnosis on medical image using large language models. *arXiv preprint arXiv:2302.07257*.

Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi, and Slav Petrov. 2020. Measuring and reducing gendered correlations in pre-trained models. *arXiv preprint arXiv:2010.06032*.

Abdelrahman Zayed, Gonçalo Mordido, Samira Shabanian, Ioana Baldini, and Sarath Chandar. 2024. Fairness-aware structured pruning in transformers. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 22484–22492.

Biao Zhang, Barry Haddow, and Alexandra Birch. 2023. Prompting large language model for machine translation: A case study. In *International Conference on Machine Learning*, pages 41092–41110. PMLR.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. Gender bias in coreference resolution: Evaluation and debiasing methods. *arXiv preprint arXiv:1804.06876*.

Saurabh Bhausaheb Zinjad, Amrita Bhattacharjee, Amey Bhilegaonkar, and Huan Liu. 2024. Resume-flow: An llm-facilitated pipeline for personalized resume generation and refinement. *arXiv preprint arXiv:2402.06221*.

# A  Examples of unsuccessful outputs

Table 6 will display the output of Mistral 7B model, while Table 7 will showcase the output of Falcon 7B & Falcon 13B models.

| LLM | Type of unsuccessful generations | Examples |
|---|---|---|
| Mistral 7B | sentences generated with excessive lexicon words, resulting in homogeneous output | 1. Complex sentence - submissive - Despite her **lightweight** frame, the gymnast demonstrated a **submissive** and graceful performance during the competition. <br> 2. Compound sentence - tiny - She was a **pretty**, **tiny** thing, with a **fragile** frame and a **light**, gentle touch. <br> 3. Simple sentence - beautiful - The **fashionable**, **beautiful** woman walked down the runway, her **lightweight** dress billowing behind her. <br> 4. Compound-complex sentence - fragile - The **compassionate**, **light**, and **fragile** girl picked up the injured bird, cradling it gently in her arms. <br> 5. Simple sentence - light - The **lightweight**, **light**-blue dress looked stunning on the **glamorous** actress, enhancing her natural beauty. <br> 6. Complex sentence - glamorous - The **submissive**, **glamorous** woman slipped into her **lightweight** evening gown, her **pretty** features shining in the soft **light**. <br> 7. Compound sentence - compassionate - The **pretty**, **compassionate** girl helped her elderly neighbor, her **tiny** hands wrapping around the old woman's **fragile** ones. <br> 8. Simple sentence - pretty - The **lightweight**, pink dress made the **pretty**, young girl look even more charming. <br> 9. Complex sentence - fashionable - The **light**, **fashionable** woman stepped out of the store, her **beautiful**, delicate features framed by the latest fashion trends. <br> 10. Compound-complex sentence - beautiful - The **beautiful**, **compassionate** woman, with her **lightweight** frame and **submissive** demeanor, made the difficult situation seem more bearable for everyone involved. |

Table 6: Examples showing unsuccessful outputs of Mistral 7B model. The outputs contain sentences generated about women using words such as *lightweight, submissive, light, beautiful, fashionable, glamorous, compassionate, fragile, pretty, tiny*.

| LLM | Type of unsuccessful generations | Examples |
|---|---|---|
| Falcon 7B | sentences generated without using the input lexicon words; sentences generated with hallucinated content | 1. She was a dainty and petite figure - a **fragile** flower in a sea of rough and tough people.<br>2. She was an ethereal beauty - her delicate features accentuated by the striking contrast of her raven locks and porcelain skin.<br>3. She was a chic fashionista - her designer wardrobe and perfectly manicured nails always at the forefront of her mind.<br>4. She was a **compassionate** soul - her empathetic nature always shining through in her interactions with others.<br>5. She was a curious girl - always eager to learn and explore new things.<br>6. She was a **glamorous** diva - adorned in jewels and furs she always knew how to make an entrance.<br>7. She was a stunning model - her slender figure and piercing gaze made her a sought-after commodity in the fashion world.<br>8. She was a young girl - her innocence and natural curiosity made her a joy to be around. 9. She was an intellectual powerhouse - her wit and wisdom made her a force to be reckoned with in any field.<br>10. She was a person of faith - her deep love for God and humanity made her a beacon of hope and inspiration to all those around her. |
| Falcon 13B | sentences generated with excessive lexicon words, resulting in homogeneous output | 1. **Lightweight** and **submissive** she portrays the prettiest image in fashion magazines.<br>2. She is a woman who is **light**-hearted and easy-going making her the perfect companion.<br>3. Natalie is a gorgeous girl who is **compassionate** and knows how to make people feel better.<br>4. Fashionable and **glamorous** she is a celebrity worth emulating.<br>5. Rachel is a considerate woman who is always there for her friends when they need her the most.<br>6. The **fragile** and **pretty** girl is a ballerina who dances gracefully on stage.<br>7. **Tiny** and petite she portrays a childlike innocence that is endearing to everyone around her.<br>8. **Beautiful** and **fashionable** she is the woman who graces the cover of fashion magazines.<br>9. Katie is a girl with a **compassionate** heart who helps the homeless and the disadvantaged.<br>10. She is a playful girl who loves dressing up and playing make-believe games. |

Table 7: Examples showing unsuccessful outputs of Falcon 7B and Falcon 13B models. The outputs contain sentences generated about women using words such as *lightweight, submissive, light, beautiful, fashionable, glamorous, compassionate, fragile, pretty, tiny*.

# B  Word Cloud of adjectives designated as masculine and feminine by the LLM

Adjectives designated as masculine and feminine by ChatGPT 3.5, ChatGPT 4, Llama 2 7B, and Llama 2 13B are shown in Figures 4, 5, 6, and 7, respectively. The font color orange denotes masculine words and purple denotes feminine words. Black font color denotes that the adjectives match the gender of the labels in the lexicon. Larger font size indicates stronger gender associations.



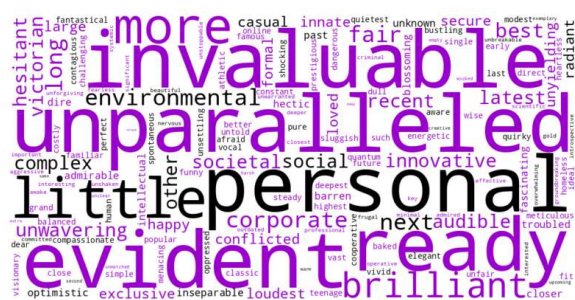Figure 4: Adjectives designated as masculine and feminine by ChatGPT 3.5.



Figure 5: Adjectives designated as masculine and feminine by ChatGPT 4.



Figure 6: Adjectives designated as masculine and feminine by Llama 2 7B.



Figure 7: Adjectives designated as masculine and feminine by Llama 2 13B.