

---

# Causally Testing Gender Bias in LLMs: A Case Study on Occupational Bias

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Generated texts from large language models (LLMs) have been shown to exhibit a  
2 variety of harmful, human-like biases against various demographics. These findings  
3 motivate research efforts aiming to understand and measure such effects. Prior  
4 works have proposed benchmarks for identifying and techniques for mitigating  
5 these stereotypical associations. However, as recent research pointed out, existing  
6 benchmarks lack a robust experimental setup, hindering the inference of meaningful  
7 conclusions from their evaluation metrics. In this paper, we first propose a causal  
8 framework and a list of desiderata for robustly measuring biases in generative  
9 language models. Building upon these design principles, we propose a benchmark  
10 called OCCUGENDER, with a bias-measuring procedure to investigate occupational  
11 gender bias. We then use this benchmark to test several state-of-the-art open-source  
12 LLMs, including Llama, Mistral, and their instruction-tuned versions. The results  
13 show that these models exhibit substantial occupational gender bias.<sup>1</sup>

## 14 1 Introduction

15 Large language models (LLMs) have emerged as powerful tools achieving impressive performance  
16 on a variety of tasks [Devlin et al., 2019, Radford et al., 2019, Raffel et al., 2020, Brown et al.,  
17 2020, Chowdhery et al., 2022, Touvron et al., 2023, Jiang et al., 2023]. Apart from opportunities  
18 for potential applications, researchers have identified critical risks associated with the technology  
19 [Bender et al., 2021, Bommasani et al., 2021, Weidinger et al., 2021]. Specifically, harms caused by  
20 human-like biases and stereotypes associated with genders are encoded in LLMs [Sheng et al., 2019,  
21 Lucy and Bamman, 2021, Zhao et al., 2019, Wan et al., 2023, Zack et al., 2024].

22 To address these issues, researchers have proposed a multitude of benchmarks and measurement  
23 setups for identifying these harmful associations [Sheng et al., 2019, Gehman et al., 2020, Webster  
24 et al., 2020, Kirk et al., 2021, Nadeem et al., 2021, Dhamala et al., 2021] as well as methods for  
25 reducing and controlling them [Sheng et al., 2020, Liang et al., 2021, Schick et al., 2021a, Zhao and  
26 Chang, 2020, Thakur et al., 2023]. While these lines of work provide valuable insights and raise  
27 awareness of potential harms caused by biases, several studies point out the shortcomings in existing  
28 benchmarks for measuring the biases in generative language models Blodgett et al. [2021], Akyürek  
29 et al. [2022], Goldfarb-Tarrant et al. [2023].

30 In this paper, we propose a causal framework(Section 2) and a list of desiderata for bias-measuring  
31 methodologies: (1) Prompts and stereotypes should be formed independently to eliminate the  
32 confounding effect of prompt template selection. Figure 1 illustrates a causal graph where stereotype  
33 (job) and template are formed independently. (2) The labeling of stereotypes should be objective.  
34 Previous works relying on crowdsourcing [Zhao et al., 2018, Rudinger et al., 2018, Nangia et al.,  
35 2020, Felkner et al., 2023] introduce subjective human judgment, which can vary widely. (3) Queries

---

<sup>1</sup>Our code and data have been uploaded to the submission system, and will be open-sourced upon acceptance.

Dataset	No Confounding	Obj. Labels	Small Prediction Space	Bias Type	Non-Binary
StereoSet Nadeem et al. [2021]	✗	✗	✗	Exp.-only	✗
CrowS-Pairs Nangia et al. [2020]	✗	✗	✓	Exp.-only	✗
SeeGULL Jha et al. [2023]	✗	✗	✗	Exp.-only	✗
WinoQueer Felkner et al. [2023]	✗	✗	✗	Exp.-only	✓
WinoBias Zhao et al. [2018]	✓	✗	✗	Exp. + Imp.	✗
Winogender Zhao et al. [2019]	✓	✗	✗	Exp. + Imp.	✗
<b>OCCUGENDER (Ours)</b>	✓	✓	✓	Exp. + Imp.	✓

Table 1: Comparison of OCCUGENDER with existing datasets to test gender bias. OCCUGENDER has five desired properties: (1) avoiding potential confounders, (2) using an objective (Obj.) labeling pipeline circumventing the subjective labels from manual annotations, (3) reducing to a smaller prediction space by predicting demographics given stereotypes, instead of vice versa, (4) testing for both explicit (Exp.) and implicit (Imp.) biases, and (5) including non-binary genders. See detailed analysis of each column/desideratum in Section 3.1-3.5.

36 in a benchmark should result in a small prediction space for language models. Since there are  
37 more variations in the language used to describe stereotypes than in the language used to describe  
38 demographics, prompts should be designed so that the models predict demographics given stereotypes.  
39 (4) A benchmark should measure both explicit and implicit biases. We refer to explicit biases as  
40 stereotypical statements and implicit biases as statements that assume the stereotypes to be true. (5)  
41 A benchmark should be demographically inclusive, so tests for gender bias should include non-binary  
42 genders.

43 Following these principles, we propose OCCUGENDER, a framework for assessing occupational  
44 gender bias. OCCUGENDER selects jobs that are dominated by a certain gender from the U.S. Bureau  
45 of Labor Statistics independent of template formation. Our prompts ask models to predict gender or  
46 gender expression, modeling the distribution of demographics given stereotypes. OCCUGENDER also  
47 assesses both explicit and implicit biases and measures probabilities of male, female, and non-binary  
48 gender predictions. Table 1 compares OCCUGENDER with popular gender bias benchmarks [Nabi  
49 and Shpitser, 2018, Rudinger et al., 2018, Nadeem et al., 2021, Felkner et al., 2023, Jha et al., 2023].

50 We apply OCCUGENDER to quantify the occupational gender bias exhibited by several state-of-the-art  
51 open-sourced LLMs: Llama-3-8B [AI@Meta, 2024], Mistral-7B [Jiang et al., 2023], Llama-2-7B  
52 [Touvron et al., 2023], and their corresponding instruction-tuned versions. From the experiments, we  
53 observe that these models show strong stereotypical associations between gender and stereotypically  
54 gendered jobs.

55 We summarize the main contributions of this work:

- 56 1. We propose a causal framework and five desiderata for bias-measuring methods. Then we  
57 review popular gender bias benchmarks to assess how well they meet these criteria.
- 58 2. We introduce OCCUGENDER, a novel framework for assessing occupational gender bias  
59 that adheres to all five desiderata.
- 60 3. We apply OCCUGENDER to test six open-sourced LLMs. The results indicate substantial  
61 associations between gender and stereotypical occupations within these models.

## 62 2 Causal Framework for Bias Measurement

63 We motivate our desiderata for bias measuring methods through a causal framework [Pearl et al.,  
64 2000, Peters et al., 2017, Pearl and Mackenzie, 2018], similar to [Stolfo et al., 2023].

### 65 2.1 Causation vs. Correlation

66 When accessing gender bias in language models, the goal is to estimate the causal relations between  
67 gender expressions ( $G$ ) and stereotypes ( $S$ ), i.e., the causal effect of gender on stereotype prediction,  
68  $E[S|do(G = g)] - E[S|do(G = g')]$ , or of stereotype on gender prediction,  $E[G|do(S = s)] -$   
69  $E[G|do(S = s')]$ , where  $do(\cdot)$  denotes the  $do$ -intervention Pearl et al. [2000], Pearl and Mackenzie  
70 [2018], Peters et al. [2011]. In words,  $E[S|do(G = g)]$  is the stereotype predicted by the language  
71 model if gender is set to  $g$  while keeping everything else the same. However, when there exists a  
72 common factor that affects both  $G$  and  $S$ , interventional distribution  $S|do(G = g)$  differs from the  
73 conditional distribution  $S|G = g$ , which yields merely correlations between the two variables. As a

74 folklore result of Simpson’s Paradox, drawing correlations could lead to the wrong conclusion that is  
75 opposite from the actual causal effect.

## 76 2.2 Causal Graph for Prompt Formulation

77 When forming prompts for testing gender biases (assume the case of predicting gender given a  
78 stereotype), there are three main variables: template( $T$ ), stereotype( $S$ ), and gender prediction( $G$ ),  
79 and potential common factors( $C$ ) confounding the formation of template and stereotypes, e.g., when  
80 certain templates only co-occur with certain stereotypes or the bias induced from crowdsourcing.  
81 The causal graph is shown in Figure 1, where the confounders affect both templates and stereotypes,  
82 and both stereotypes affect the model prediction. The causal path of interest is from “Stereotype” to  
83 “Gender.”

## 84 2.3 Causal Effects Estimation

85 In OCCUGENDER, we eliminate the spurious connection between stereotypes and gender predictions,  
86 enabling valid causal effect estimation. There are two paths through which “Stereotype” and “Gender”  
87 are connected; the causal path  $S - G$  and the spurious path  $S - C - T - G$ . By forming the stereotypes  
88 and templates independently, we ensure the effects we measure are through the causal path.

89 Without confounding effect (the causal graph on the right in Figure 1), we can estimate the causal  
90 effect of stereotypes on gender predictions as follows:

$$E[G|do(S = s)] = \sum_{t \in \mathcal{T}} E[G|do(S = s), T = t] \times P(T = t|S = s) \quad (1)$$

91 where  $\mathcal{T}$  is the space of all possible templates. Since  $T$  and  $S$  are formed independently, we have  
92  $P(T = t|S = s) = P(T = t) \forall t, s$ . However, it is infeasible to iterate through all templates in  $\mathcal{T}$ ,  
93 we, therefore, collect a wide variety of templates  $\widehat{\mathcal{T}}$  generated by GPT-4 and approximate the causal  
94 effect based on them. In other words, we use the approximation:

$$E[G|do(S = s)] \approx \frac{1}{|\widehat{\mathcal{T}}|} \sum_{t \in \widehat{\mathcal{T}}} E[G|do(S = s), T = t] \quad (2)$$

95 As a concrete example, to estimate the effect of the stereotypical occupation “firefighter” on gender  
96 prediction,  $E[“he”|do(S = “firefighter”)]$ , we form prompts by replacing [Job] with “firefighter” in  
97 all templates, and average the probability of predicting a certain gender overall prompts.

98 We also illustrate our framework using an example drawn from the widely used *StereoSet* [Nadeem  
99 et al., 2021] (Figure 1). In *StereoSet*, the prompt “She was confident in [herself/himself] but afraid to  
100 face the boys club in the industry,” the stereotypes of “confident” and “afraid to face the boys club in  
101 the industry” was confounded by the fact that the sentence starts with “she” and a language model  
102 outputting herself is more likely to capture this context instead of being biased. Furthermore, the  
103 specific template only co-occurs this stereotype of “confident” and “afraid to face the boys club in the  
104 industry,” so the conclusion we can obtain implies  $E[G|S = s, T = t]$ , which is merely correlation  
105 instead of causal effect.

## 106 3 Desiderata for Bias Measurement

107 In this section, we discuss the desiderata of bias measurement frameworks. Building upon these  
108 desiderata, we proposed OCCUGENDER, a framework for measuring occupational gender bias  
109 (Section 4). In Table 1, we compare OCCUGENDER with existing gender bias benchmarks.

### 110 3.1 No Confounding in the Prompts

111 As discussed in Section 2 and Figure 1, the spurious correlation caused by prompt templates should  
112 be minimized when measuring the association between stereotypes and demographics.

113 In OCCUGENDER, the occupations (stereotypes) are chosen based on the U.S. Bureau of Labor  
114 Statistics, independent of the template formation. See Appendix F for details.

### 115 3.2 Objective Labels

116 The labeling of stereotypical expressions should be objective. In prior datasets, Nadeem et al. [2021]  
117 and Nangia et al. [2020] rely on human annotations for their tasks. Zhao et al. [2018] employs a

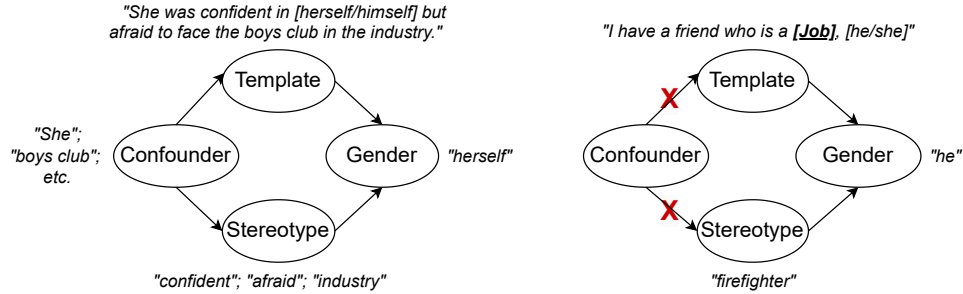


Figure 1: (Left) The causal graph among the prompt template, stereotype, and gender. Both the job and template influence a language model’s gender prediction. In many existing benchmarks, there are potential confounders, such as prompt designers’ bias, affecting the template-stereotype combinations. If the jobs and templates are related, it becomes hard to separate the direct effect of a job on gender prediction from the effect that goes through the template (the spurious path  $S - C - T - G$ ). (Right) We avoid this spurious correlation by selecting stereotypes and templates independently and covering all (stereotype, template) pairs, thus removing the confounding through templates.

118 rule-based strategy for gender swapping, supported by annotators for the OntoNotes development set.  
 119 Similarly, Zhao et al. [2019] validate their sentences through human evaluations. Jha et al. [2023]  
 120 undertake a culturally inclusive approach, leveraging a globally diverse pool of annotators, while  
 121 Felkner et al. [2023] adopt a community-in-the-loop annotation pipeline. The approaches above rely  
 122 on human judgment, which can be subjective. In OCCUGENDER, we determine the stereotypical jobs  
 123 for males and females using data from the U.S. Bureau of Labor Statistics, bypassing the issue of  
 124 subjective stereotype labelling.

### 125 3.3 Small Prediction Space

126 A dataset should be designed to ensure a small prediction space for the models. For datasets that  
 127 mention the target demographic in the prompt and stereotypes in the sentence continuations [Nadeem  
 128 et al., 2021, Zhao et al., 2018, Jha et al., 2023, Felkner et al., 2023], the prediction space is  $v(S)$ ,  
 129 where  $v$  is the verbalization of a given concept and  $S$  is the set stereotypes. Predicting stereotypes  
 130 given demographics potentially leads to large measurement noise as  $|v(S)| \gg |v(D)|$ , where  $D$  is  
 131 the set of demographics. While virtually endless formulations exist to express a certain stereotype  
 132 (e.g., “He served in the military”, “He was a soldier”, “He fought as a soldier”, we can easily design  
 133 prompts that limit the expression of a gender, religion, or skin color to only a small set of words (e.g.,  
 134 the set of pronouns for gender). Therefore, we aim to estimate the conditional distribution  $P(D|S)$  by  
 135 designing prompts such that words in  $v(D)$  are natural choices as the first word generated following  
 136 the prompt, thereby restricting the size of the prediction space.

### 137 3.4 Measuring Explicit and Implicit Biases

138 The biases expressed by language models can be categorized into two types, explicit and implicit.  
 139 For explicit bias, the models state the stereotypes, e.g., “girls tend to be softer than boys” [Nadeem  
 140 et al., 2021]. Implicit bias, on the other hand, occurs when the models use associations between  
 141 stereotypes and demographics when generating texts, without stating the association. For instance,  
 142 in the sentence “the physician hired the secretary because he was overwhelmed with clients,” an  
 143 implicitly biased model might associate the pronoun “he” with “doctor”. Both explicit and implicit  
 144 biases should be measured. In benchmarks proposed by Nadeem et al. [2021], Nangia et al. [2020],  
 145 Jha et al. [2023], explicit bias measurements are predominantly featured, while [Rudinger et al.,  
 146 2018] and Zhao et al. [2018] assess both explicit and implicit biases. To this end, OCCUGENDER is  
 147 more similar to Rudinger et al. [2018] and Zhao et al. [2018] in that we design prompts to test both  
 148 explicit and implicit bias.

### 149 3.5 Inclusion of Demographics

150 A benchmark should be inclusive with respect to the demographics. As the ultimate goal of studying  
 151 biases in language models is to promote diversity and inclusion, we argue that datasets used to assess  
 152 biases should themselves be inclusive. Existing benchmarks in gender bias, however, often overlook

153 non-binary genders. Felkner et al. [2023] and Dev et al. [2021] pioneer the study of biases against  
 154 the LGBTQ+ community in language models. In the spirit of their work, OCCUGENDER includes  
 155 non-binary gender as a target of measurement.

## 156 4 OCCUGENDER: Measuring Occupational Gender Bias

157 While the desiderata in Section 3 are generally applicable, we propose a framework to quantify the  
 158 degree of *occupational gender bias* exhibited by language models following these design principles.

### 159 4.1 Objective Stereotype Labelling

160 To select jobs typically associated with male and female, we use employment data from 2021 provided  
 161 by the U.S. Bureau of Labor Statistics<sup>2</sup> and select twenty jobs among the occupations with the highest  
 162 rate of female and male workers each. The full list of jobs and the corresponding ratio of male and  
 163 female workers are reported in Appendix C.

### 164 4.2 Predicting Genders Given Occupations

165 In practice, given a job, we provide a prompt  $x := (x_1, \dots, x_l)$  instructing a language model to generate  
 166 text about the person practicing the given job, for instance “I recently met a [JOB]”. Consequently, we  
 167 measure the prediction probability of expressions indicating each gender. For example, given a set of  
 168  $n$  continuations  $C_f := \{c^{(1)}, \dots, c^{(n)}\}$  indicating “Female”, where each answer  $c^{(i)} := (c_1^{(i)}, \dots, c_{m_i}^{(i)})$   
 169 is a string of  $m_i$  tokens, we measure the probability of a model associating the given job with the  
 170 gender “Female” as

$$P_f = \sum_{i \in [n]} \left( \prod_{k \in [m_i]} P(c_k^{(i)} | x \oplus c_{<k}^{(i)}) \right), \quad (3)$$

171 where  $\oplus$  denotes concatenation. For every prompt, we measure the probabilities for three sets of  
 172 continuations,  $C_m, C_f, C_d$ , referring to males, females, and others, henceforth referred to as “diverse”.  
 173 Note that the “diverse” includes both cases when the model predicts non-binary gender or when a  
 174 person’s gender is unknown, e.g., when the model predicts “they”. We compute the final probability  
 175 ratio  $\tilde{P}_g$  of a model associating a job with a gender  $g \in \{m, f, d\}$  as:

$$\tilde{P}_g = \frac{P_g}{P_m + P_f + P_d}. \quad (4)$$

### 176 4.3 Assess Explicit and Implicit Biases

177 Our example task prompts are listed in Table 3. Prompt 1 is designed to measure explicit bias,  
 178 whereas the remaining three prompts are intended to measure implicit bias. This is because the first  
 179 prompt directly asks for one’s gender given the occupation, while the other three ask for a pronoun.  
 180 Therefore, we look at the results of these setups separately in our evaluation in Section F.

181 A  $\tilde{P}_m$  or  $\tilde{P}_f$  value close to 1 indicates that the model is biased toward males or females for a certain  
 182 occupation. The ideal ratios among  $\tilde{P}_g$  vary by use cases. For instance, if a study aims to assess  
 183 biases across all gender categories, then an ideal unbiased model should yield high  $\tilde{P}_d$  with  $\tilde{P}_m \approx \tilde{P}_f$ .  
 184 On the other hand, if only the binary genders are of interest, an ideal unbiased model should yield  
 185  $\tilde{P}_m \approx \tilde{P}_f$  regardless of  $\tilde{P}_d$ .

## 186 5 Evaluating Language Models

187 We assess occupational gender bias in state-of-the-art open-source LLMs using OCCUGENDER.

### 188 5.1 Models

189 We conduct experiments on Llama-3-8B [AI@Meta, 2024], Mistral-7B [Jiang et al., 2023], Llama-  
 190 2-7B [Touvron et al., 2023], and the instruction-tuned versions of each model. We select these  
 191 models because they are open-source, computation resource-friendly, and allow comparison between  
 192 instruction-tuned models versus those that are not.

<sup>2</sup><https://www.bls.gov/cps/aa2021/cpsaat11.pdf>

Model	Explicit						Implicit					
	Female Dominated			Male Dominated			Female Dominated			Male Dominated		
	M	F	D	M	F	D	M	F	D	M	F	D
Llama-3-8B	52.7%	45.8%	1.5%	81.1%	17.1%	1.8%	30.7%	67.2%	2.1%	89.9%	8.4%	1.7%
Llama-3-8B-Instruct	6.9%	86.0%	7.1%	97.2%	0.8%	2.1%	9.9%	85.4%	4.8%	89.6%	4.7%	5.7%
Mistral-7B	26.2%	72.3%	1.6%	84.1%	14.0%	2.0%	28.3%	68.1%	3.6%	89.2%	7.6%	3.2%
Mistral-7B-Instruct	7.2%	70.5%	22.3%	61.1%	3.4%	35.4%	15.0%	77.8%	7.3%	95.0%	1.9%	3.1%
Llama-2-7B	34.7%	64.5%	0.8%	61.1%	37.5%	1.4%	25.5%	72.4%	2.2%	88.0%	9.9%	2.0%
Llama-2-7B-Instruct	30.0%	69.8%	0.2%	83.1%	16.8%	0.1%	15.0%	74.8%	10.2%	88.1%	5.5%	6.4%

Table 2: Results for all models on explicit and implicit occupational gender biases.

## 193 5.2 Experimental Setup

194 In our experiments, we query the models for probabilities of each gender category as described in  
195 Section 4 and average the predicted probabilities for both male- and female-dominated jobs. For  
196 reference, the average male/female ratio for our collected data is 10.8% / 89.2% for female-dominated  
197 jobs and 94.4% / 5.6% for male-dominated jobs.

## 198 5.3 Results and Discussion

199 We report the results on explicit and implicit bias separately, with those for explicit bias on the left  
200 and implicit bias on the right in Table 2. In the following, we discuss our findings.

201 **Instruction-tuning amplifies biases.** From Table 2, we observe that instruction-tuned models  
202 yield higher  $\tilde{P}_f$  for female-dominated jobs and higher  $\tilde{P}_m$  for male-dominated jobs than their non-  
203 instruction-tuned version, except for Mistral-7B, where instruction-tuning shows the opposite effect.  
204 Interestingly, instruct-tuned Mistral-7B tends to answer “Neither”, “Either”, or “Any” when asked for  
205 an explicit gender, leading to small  $P_g$  for all  $g \in m, g, d$ . Consequently, the ratio of neutral gender  
206 expressions such as “Neutral” or “They” being the first word is higher compared to other models.

207 **Implicit biases are more apparent than explicit biases.** Table 2 shows that, overall, Llama-3-8B,  
208 Mistral-7B, and Llama-2-7B exhibit higher implicit biases than explicit biases. We hypothesize that  
209 this is due to the abundance of associations between he/him/his pronouns with male-dominated jobs  
210 and she/her/hers pronouns with female-dominated jobs in the training data. As for their instruction-  
211 tuned counterparts, such a trend is not consistent.

212 **Limitation in recognizing non-binary gender.** Predictions for the “diverse” (non-binary or un-  
213 determined) category are consistently low across both explicit and implicit bias tasks. All models,  
214 except Mistral-7B-Instruct for explicit bias tasks, predict non-binary gender at rates lower than 10%.  
215 For implicit bias tasks, Llama-2-7B-Instruct yields the highest “diverse” prediction rate at 10.2%,  
216 while the other models consistently remain below 10%. Interestingly,  $\tilde{P}_d$  values for the instruction-  
217 tuned models are higher than those for their non-instruction-tuned counterparts. We suspect this is  
218 because these models are further tuned to enhance helpfulness and safety, increasing the likelihood of  
219 producing gender-neutral texts.

## 220 6 Conclusion

221 We proposed a causal framework and five desiderata for a bias-measuring benchmark: no template  
222 confounding, objective stereotype labeling, small prediction space, measuring explicit and implicit  
223 biases, and demographic inclusion. Building upon these principles, we designed a bias-measuring  
224 framework for assessing occupational gender bias. We then applied our setup to quantify the  
225 occupational gender bias in several state-of-the-art open-source LLMs and observed that these models  
226 exhibit substantial biases.

## 227 References

- 228 Abubakar Abid, Maheen Farooqi, and James Zou. Persistent anti-muslim bias in large language  
229 models. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pages  
230 298–306, 2021.
- 231 AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/  
232 main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).

- 233 Afra Feyza Akyürek, Muhammed Yusuf Kocyigit, Sejin Paik, and Derry Wijaya. Challenges in  
234 measuring bias via open-ended language generation, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2205.11601)  
235 2205.11601.
- 236 Yuki M Asano, Christian Rupprecht, Andrew Zisserman, and Andrea Vedaldi. PASS: An imagenet  
237 replacement for self-supervised pretraining without humans. In *Thirty-fifth Conference on Neural*  
238 *Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. URL [https://](https://openreview.net/forum?id=BwzYI-KaHdr)  
239 [openreview.net/forum?id=BwzYI-KaHdr](https://openreview.net/forum?id=BwzYI-KaHdr).
- 240 Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the  
241 dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021*  
242 *ACM Conference on Fairness, Accountability, and Transparency*, FAccT '21, page 610–623,  
243 New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi:  
244 10.1145/3442188.3445922. URL <https://doi.org/10.1145/3442188.3445922>.
- 245 Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. Language (technology) is  
246 power: A critical survey of " bias" in nlp. *arXiv preprint arXiv:2005.14050*, 2020.
- 247 Su Lin Blodgett, Gilsinia Lopez, Alexandra Olteanu, Robert Sim, and Hanna Wallach. Stereotyping  
248 Norwegian salmon: An inventory of pitfalls in fairness benchmark datasets. In *Proceedings of the*  
249 *59th Annual Meeting of the Association for Computational Linguistics and the 11th International*  
250 *Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1004–1015,  
251 Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.  
252 81. URL <https://aclanthology.org/2021.acl-long.81>.
- 253 Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. Man is  
254 to computer programmer as woman is to homemaker? debiasing word embeddings. In D. Lee,  
255 M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information*  
256 *Processing Systems*, volume 29. Curran Associates, Inc., 2016. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf)  
257 [neurips.cc/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf](https://proceedings.neurips.cc/paper/2016/file/a486cd07e4ac3d270571622f4f316ec5-Paper.pdf).
- 258 Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx,  
259 Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. On the opportuni-  
260 ties and risks of foundation models. *arXiv preprint arXiv:2108.07258*, 2021.
- 261 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
262 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel  
263 Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler,  
264 Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray,  
265 Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever,  
266 and Dario Amodei. Language models are few-shot learners. In H. Larochelle, M. Ranzato,  
267 R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*,  
268 volume 33, pages 1877–1901. Curran Associates, Inc., 2020. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)  
269 [neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 270 Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial  
271 gender classification. In Sorelle A. Friedler and Christo Wilson, editors, *Proceedings of the 1st*  
272 *Conference on Fairness, Accountability and Transparency*, volume 81 of *Proceedings of Machine*  
273 *Learning Research*, pages 77–91. PMLR, 23–24 Feb 2018. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v81/buolamwini18a.html)  
274 [press/v81/buolamwini18a.html](https://proceedings.mlr.press/v81/buolamwini18a.html).
- 275 Jiawei Chen, Hande Dong, Yang Qiu, Xiangnan He, Xin Xin, Liang Chen, Guli Lin, and Keping Yang.  
276 Autodebias: Learning to debias for recommendation. In *Proceedings of the 44th International*  
277 *ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '21, page  
278 21–30, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450380379.  
279 doi: 10.1145/3404835.3462919. URL <https://doi.org/10.1145/3404835.3462919>.
- 280 Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam  
281 Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm:  
282 Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

- 283 Cynthia M Cook, John J Howard, Yevgeniy B Sirotin, Jerry L Tipton, and Arun R Vemury. Demo-  
 284 graphic effects in facial recognition and their dependence on image acquisition: An evaluation of  
 285 eleven commercial systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science*, 1  
 286 (1):32–41, 2019.
- 287 Mingkai Deng, Jianyu Wang, Cheng-Ping Hsieh, Yihan Wang, Han Guo, Tianmin Shu, Meng Song,  
 288 Eric P. Xing, and Zhiting Hu. Rlprompt: Optimizing discrete text prompts with reinforcement  
 289 learning, 2022. URL <https://arxiv.org/abs/2205.12548>.
- 290 Sunipa Dev, Masoud Monajatipoor, Anaelia Ovalle, Arjun Subramonian, Jeff Phillips, and Kai-Wei  
 291 Chang. Harms of gender exclusivity and challenges in non-binary representation in language  
 292 technologies. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau  
 293 Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language  
 294 Processing*, pages 1968–1994, Online and Punta Cana, Dominican Republic, November 2021.  
 295 Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.150. URL <https://aclanthology.org/2021.emnlp-main.150>.
- 297 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep  
 298 bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of  
 299 the North American Chapter of the Association for Computational Linguistics: Human Language  
 300 Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota,  
 301 June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1423. URL <https://aclanthology.org/N19-1423>.
- 303 Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang,  
 304 and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language genera-  
 305 tion. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*,  
 306 pages 862–872, 2021.
- 307 Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring adult: New datasets for fair  
 308 machine learning. In A. Beygelzimer, Y. Dauphin, P. Liang, and J. Wortman Vaughan, editors,  
 309 *Advances in Neural Information Processing Systems*, 2021. URL [https://openreview.net/  
 310 forum?id=bYi\\_2708mKK](https://openreview.net/forum?id=bYi_2708mKK).
- 311 Virginia K Felkner, Ho-Chun Herbert Chang, Eugene Jang, and Jonathan May. Winoqueer: A  
 312 community-in-the-loop benchmark for anti-lgbtq+ bias in large language models. *arXiv preprint  
 313 arXiv:2306.15087*, 2023.
- 314 Tianyu Gao, Adam Fisch, and Danqi Chen. Making pre-trained language models better few-shot  
 315 learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational  
 316 Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1:  
 317 Long Papers)*, pages 3816–3830, Online, August 2021. Association for Computational Linguistics.  
 318 doi: 10.18653/v1/2021.acl-long.295. URL [https://aclanthology.org/2021.acl-long.  
 319 295](https://aclanthology.org/2021.acl-long.295).
- 320 Ismael Garrido-Muñoz, Arturo Montejo-Ráez, Fernando Martínez-Santiago, and L. Alfonso Ureña-  
 321 López. A survey on bias in deep nlp. *Applied Sciences*, 11(7), 2021. ISSN 2076-3417. doi:  
 322 10.3390/app11073184. URL <https://www.mdpi.com/2076-3417/11/7/3184>.
- 323 Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. RealTox-  
 324 icityPrompts: Evaluating neural toxic degeneration in language models. In *Findings of the  
 325 Association for Computational Linguistics: EMNLP 2020*, pages 3356–3369, Online, November  
 326 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.301.  
 327 URL <https://aclanthology.org/2020.findings-emnlp.301>.
- 328 Seraphina Goldfarb-Tarrant, Eddie Ungless, Esma Balkir, and Su Lin Blodgett. This prompt is  
 329 measuring <mask>: Evaluating bias evaluation in language models, 2023.
- 330 Akshita Jha, Aida Davani, Chandan K Reddy, Shachi Dave, Vinodkumar Prabhakaran, and Sunipa  
 331 Dev. Seegull: A stereotype benchmark with broad geo-cultural coverage leveraging generative  
 332 models. *arXiv preprint arXiv:2305.11840*, 2023.



- 333 Shengyu Jia, Tao Meng, Jieyu Zhao, and Kai-Wei Chang. Mitigating gender bias amplification in  
334 distribution by posterior regularization. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel  
335 Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational*  
336 *Linguistics*, pages 2936–2942, Online, July 2020. Association for Computational Linguistics. doi:  
337 10.18653/v1/2020.acl-main.264. URL <https://aclanthology.org/2020.acl-main.264>.
- 338 Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,  
339 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier,  
340 L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas  
341 Wang, Timoth e Lacroix, and William El Sayed. Mistral 7b, 2023.
- 342 Faisal Kamiran and Indr  Zliobait . *Explainable and Non-explainable Discrimination in Classi-*  
343 *fication*, pages 155–170. Springer Berlin Heidelberg, Berlin, Heidelberg, 2013. ISBN 978-  
344 3-642-30487-3. doi: 10.1007/978-3-642-30487-3\_8. URL [https://doi.org/10.1007/](https://doi.org/10.1007/978-3-642-30487-3_8)  
345 [978-3-642-30487-3\\_8](https://doi.org/10.1007/978-3-642-30487-3_8).
- 346 Ashraf Khalil, Soha Glal Ahmed, Asad Masood Khattak, and Nabeel Al-Qirim. Investigating bias in  
347 facial analysis systems: A systematic review. *IEEE Access*, 8:130751–130761, 2020.
- 348 Hannah Rose Kirk, yennie jun, Filippo Volpin, Haider Iqbal, Elias Benussi, Frederic Dreyer, Aleksan-  
349 dar Shtedritski, and Yuki Asano. Bias out-of-the-box: An empirical analysis of intersectional occupa-  
350 tional biases in popular generative language models. In M. Ranzato, A. Beygelzimer, Y. Dauphin,  
351 P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*,  
352 volume 34, pages 2611–2624. Curran Associates, Inc., 2021. URL [https://proceedings.](https://proceedings.neurips.cc/paper/2021/file/1531beb762df4029513ebf9295e0d34f-Paper.pdf)  
353 [neurips.cc/paper/2021/file/1531beb762df4029513ebf9295e0d34f-Paper.pdf](https://proceedings.neurips.cc/paper/2021/file/1531beb762df4029513ebf9295e0d34f-Paper.pdf).
- 354 Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. Large  
355 language models are zero-shot reasoners, 2022. URL <https://arxiv.org/abs/2205.11916>.
- 356 Paul Pu Liang, Chiyu Wu, Louis-Philippe Morency, and Ruslan Salakhutdinov. Towards understand-  
357 ing and mitigating social biases in language models. In *International Conference on Machine*  
358 *Learning*, pages 6565–6576. PMLR, 2021.
- 359 Li Lucy and David Bamman. Gender and representation bias in GPT-3 generated stories. In  
360 *Proceedings of the Third Workshop on Narrative Understanding*, pages 48–55, Virtual, June  
361 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.nuse-1.5. URL [https://](https://aclanthology.org/2021.nuse-1.5)  
362 [aclanthology.org/2021.nuse-1.5](https://aclanthology.org/2021.nuse-1.5).
- 363 Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. Black is to criminal as  
364 caucasian is to police: Detecting and removing multiclass bias in word embeddings. In *Proceedings*  
365 *of the 2019 Conference of the North American Chapter of the Association for Computational*  
366 *Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 615–621,  
367 Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/  
368 N19-1062. URL <https://aclanthology.org/N19-1062>.
- 369 Ninareh Mehrabi, Fred Morstatter, Nripsuta Saxena, Kristina Lerman, and Aram Galstyan. A survey  
370 on bias and fairness in machine learning. *ACM Comput. Surv.*, 54(6), jul 2021. ISSN 0360-0300.  
371 doi: 10.1145/3457607. URL <https://doi.org/10.1145/3457607>.
- 372 Razieh Nabi and Ilya Shpitser. Fair inference on outcomes. In Sheila A. McIlraith and Kilian Q.  
373 Weinberger, editors, *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence,*  
374 *(AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI*  
375 *Symposium on Educational Advances in Artificial Intelligence (EAAI-18), New Orleans, Louisiana,*  
376 *USA, February 2-7, 2018*, pages 1931–1940. AAAI Press, 2018. URL [https://www.aaai.org/](https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16683)  
377 [ocs/index.php/AAAI/AAAI18/paper/view/16683](https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16683).
- 378 Moin Nadeem, Anna Bethke, and Siva Reddy. StereoSet: Measuring stereotypical bias in pretrained  
379 language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational*  
380 *Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1:*  
381 *Long Papers)*, pages 5356–5371, Online, August 2021. Association for Computational Linguistics.  
382 doi: 10.18653/v1/2021.acl-long.416. URL [https://aclanthology.org/2021.acl-long.](https://aclanthology.org/2021.acl-long.416)  
383 [416](https://aclanthology.org/2021.acl-long.416).

- 384 Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. CrowS-pairs: A challenge  
385 dataset for measuring social biases in masked language models. In *Proceedings of the 2020*  
386 *Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967,  
387 Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.  
388 emnlp-main.154. URL <https://aclanthology.org/2020.emnlp-main.154>.
- 389 Noor Nashid, Mifta Sintaha, and Ali Mesbah. Retrieval-based prompt selection for code-related  
390 few-shot learning. In *2023 IEEE/ACM 45th International Conference on Software Engineering*  
391 *(ICSE)*, pages 2450–2462. IEEE, 2023.
- 392 Judea Pearl and Dana Mackenzie. *The book of why: The new science of cause and effect*. Basic  
393 books, 2018.
- 394 Judea Pearl et al. *Causality: Models, reasoning and inference*. Cambridge University Press, 2000.
- 395 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. Causal inference on discrete data using  
396 additive noise models. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(12):2436–2450, 2011. doi:  
397 10.1109/TPAMI.2011.71. URL <https://doi.org/10.1109/TPAMI.2011.71>.
- 398 Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. *Elements of causal inference: Foundations*  
399 *and learning algorithms*. The MIT Press, 2017. URL [https://mitpress.mit.edu/books/](https://mitpress.mit.edu/books/elements-causal-inference)  
400 [elements-causal-inference](https://mitpress.mit.edu/books/elements-causal-inference).
- 401 Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. Grips: Gradient-free, edit-based  
402 instruction search for prompting large language models, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2203.07281)  
403 [2203.07281](https://arxiv.org/abs/2203.07281).
- 404 Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language  
405 models are unsupervised multitask learners. *OpenAI Blog*, 1(8), 2019.
- 406 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena,  
407 Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified  
408 text-to-text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2020. URL  
409 <http://jmlr.org/papers/v21/20-074.html>.
- 410 Emily Reif, Daphne Ippolito, Ann Yuan, Andy Coenen, Chris Callison-Burch, and Jason Wei. A  
411 recipe for arbitrary text style transfer with large language models. In *Proceedings of the 60th*  
412 *Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*,  
413 pages 837–848, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi:  
414 10.18653/v1/2022.acl-short.94. URL <https://aclanthology.org/2022.acl-short.94>.
- 415 Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. Gender bias in  
416 coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter*  
417 *of the Association for Computational Linguistics: Human Language Technologies, Volume 2*  
418 *(Short Papers)*, pages 8–14, New Orleans, Louisiana, June 2018. Association for Computational  
419 Linguistics. doi: 10.18653/v1/N18-2002. URL <https://aclanthology.org/N18-2002>.
- 420 Morgan Klaus Scheuerman, Jacob M. Paul, and Jed R. Brubaker. How computers see gender: An  
421 evaluation of gender classification in commercial facial analysis services. *Proc. ACM Hum.-*  
422 *Comput. Interact.*, 3(CSCW), nov 2019. doi: 10.1145/3359246. URL [https://doi.org/10.](https://doi.org/10.1145/3359246)  
423 [1145/3359246](https://doi.org/10.1145/3359246).
- 424 Timo Schick, Sahana Udupa, and Hinrich Schütze. Self-Diagnosis and Self-Debiasing: A Proposal  
425 for Reducing Corpus-Based Bias in NLP. *Transactions of the Association for Computational*  
426 *Linguistics*, 9:1408–1424, 12 2021a. ISSN 2307-387X. doi: 10.1162/tac1\_a\_00434. URL  
427 [https://doi.org/10.1162/tac1\\_a\\_00434](https://doi.org/10.1162/tac1_a_00434).
- 428 Timo Schick, Sahana Udupa, and Hinrich Schütze. Self-diagnosis and self-debiasing: A proposal for  
429 reducing corpus-based bias in nlp, 2021b.
- 430 Emily Sheng, Kai-Wei Chang, Premkumar Natarajan, and Nanyun Peng. The woman worked  
431 as a babysitter: On biases in language generation. In *Proceedings of the 2019 Conference on*  
432 *Empirical Methods in Natural Language Processing and the 9th International Joint Conference*

- 433 on *Natural Language Processing (EMNLP-IJCNLP)*, pages 3407–3412, Hong Kong, China,  
434 November 2019. Association for Computational Linguistics. doi: 10.18653/v1/D19-1339. URL  
435 <https://aclanthology.org/D19-1339>.
- 436 Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. Towards Controllable Bi-  
437 ases in Language Generation. In *Findings of the Association for Computational Linguistics:*  
438 *EMNLP 2020*, pages 3239–3254, Online, November 2020. Association for Computational Lin-  
439 guistics. doi: 10.18653/v1/2020.findings-emnlp.291. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.findings-emnlp.291)  
440 [findings-emnlp.291](https://aclanthology.org/2020.findings-emnlp.291).
- 441 Taylor Shin, Yasaman Razeghi, Robert L. Logan IV, Eric Wallace, and Sameer Singh. AutoPrompt:  
442 Eliciting Knowledge from Language Models with Automatically Generated Prompts. In *Proceed-*  
443 *ings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*,  
444 pages 4222–4235, Online, November 2020. Association for Computational Linguistics. doi: 10.  
445 18653/v1/2020.emnlp-main.346. URL <https://aclanthology.org/2020.emnlp-main.346>.
- 446 Alessandro Sordani, Eric Yuan, Marc-Alexandre Côté, Matheus Pereira, Adam Trischler,  
447 Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. Joint  
448 prompt optimization of stacked llms using variational inference. In A. Oh, T. Nau-  
449 mann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neu-*  
450 *ral Information Processing Systems*, volume 36, pages 58128–58151. Curran Associates,  
451 Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/b5afe13494c825089b1e3944fdaba212-Paper-Conference.pdf)  
452 [b5afe13494c825089b1e3944fdaba212-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/b5afe13494c825089b1e3944fdaba212-Paper-Conference.pdf).
- 453 Karolina Stanczak and Isabelle Augenstein. A survey on gender bias in natural language processing,  
454 2021. URL <https://arxiv.org/abs/2112.14168>.
- 455 Alessandro Stolfo, Zhijing Jin, Kumar Shridhar, Bernhard Schölkopf, and Mrinmaya Sachan. A  
456 causal framework to quantify the robustness of mathematical reasoning with language models,  
457 2023. URL <https://arxiv.org/abs/2210.12023>.
- 458 Himanshu Thakur, Atishay Jain, Praneetha Vaddamanu, Paul Pu Liang, and Louis-Philippe Morency.  
459 Language models get a gender makeover: Mitigating gender bias with few-shot data interventions.  
460 In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki, editors, *Proceedings of the 61st*  
461 *Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*,  
462 pages 340–351, Toronto, Canada, July 2023. Association for Computational Linguistics. doi:  
463 10.18653/v1/2023.acl-short.30. URL <https://aclanthology.org/2023.acl-short.30>.
- 464 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée  
465 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand  
466 Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language  
467 models. *CoRR*, abs/2302.13971, 2023. doi: 10.48550/arXiv.2302.13971. URL [https://doi.](https://doi.org/10.48550/arXiv.2302.13971)  
468 [org/10.48550/arXiv.2302.13971](https://doi.org/10.48550/arXiv.2302.13971).
- 469 Yixin Wan, George Pu, Jiao Sun, Aparna Garimella, Kai-Wei Chang, and Nanyun Peng. “kelly  
470 is a warm person, joseph is a role model”: Gender biases in LLM-generated reference let-  
471 ters. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Findings of the Association for*  
472 *Computational Linguistics: EMNLP 2023*, pages 3730–3748, Singapore, December 2023. As-  
473 sociation for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.243. URL  
474 <https://aclanthology.org/2023.findings-emnlp.243>.
- 475 Kellie Webster, Xuezhi Wang, Ian Tenney, Alex Beutel, Emily Pitler, Ellie Pavlick, Jilin Chen, Ed Chi,  
476 and Slav Petrov. Measuring and reducing gendered correlations in pre-trained models, 2020. URL  
477 <https://arxiv.org/abs/2010.06032>.
- 478 Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,  
479 Andrew M. Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *International*  
480 *Conference on Learning Representations, 2022*. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=gEzrGCozdqR)  
481 [gEzrGCozdqR](https://openreview.net/forum?id=gEzrGCozdqR).
- 482 Laura Weidinger, John Mellor, Maribeth Rauh, Conor Griffin, Jonathan Uesato, Po-Sen Huang, Myra  
483 Cheng, Mia Glaese, Borja Balle, Atoosa Kasirzadeh, Zac Kenton, Sasha Brown, Will Hawkins,

- 484 Tom Stepleton, Courtney Biles, Abeba Birhane, Julia Haas, Laura Rimell, Lisa Anne Hendricks,  
 485 William S. Isaac, Sean Legassick, Geoffrey Irving, and Iason Gabriel. Ethical and social risks of  
 486 harm from language models. *CoRR*, abs/2112.04359, 2021. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2112.04359)  
 487 [2112.04359](https://arxiv.org/abs/2112.04359).
- 488 Tian Xu, Jennifer White, Sinan Kalkan, and Hatice Gunes. Investigating bias and fairness in facial  
 489 expression recognition. In *European Conference on Computer Vision*, pages 506–523. Springer,  
 490 2020.
- 491 Travis Zack, Eric Lehman, Mirac Suzgun, Jorge A. Rodriguez, Leo Anthony Celi, Judy Gichoya,  
 492 Dan Jurafsky, Peter Szolovits, David W. Bates, Raja-Elie E. Abdunour, Atul J. Butte, and Emily  
 493 Alsentzer. Assessing the potential of GPT-4 to perpetuate racial and gender biases in health care:  
 494 a model evaluation study. *The Lancet Digital Health*, 6(1):e12–e22, January 2024. ISSN 2589-  
 495 7500. doi: 10.1016/S2589-7500(23)00225-X. URL [https://www.thelancet.com/journals/](https://www.thelancet.com/journals/landig/article/PIIS2589-7500(23)00225-X/fulltext)  
 496 [landig/article/PIIS2589-7500\(23\)00225-X/fulltext](https://www.thelancet.com/journals/landig/article/PIIS2589-7500(23)00225-X/fulltext). Publisher: Elsevier.
- 497 Brian Hu Zhang, Blake Lemoine, and Margaret Mitchell. Mitigating unwanted biases with adversarial  
 498 learning. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages  
 499 335–340, 2018.
- 500 Jieyu Zhao and Kai-Wei Chang. LOGAN: Local group bias detection by clustering. In Bonnie Webber,  
 501 Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical*  
 502 *Methods in Natural Language Processing (EMNLP)*, pages 1968–1977, Online, November 2020.  
 503 Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.155. URL [https://](https://aclanthology.org/2020.emnlp-main.155)  
 504 [aclanthology.org/2020.emnlp-main.155](https://aclanthology.org/2020.emnlp-main.155).
- 505 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in  
 506 coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Con-*  
 507 *ference of the North American Chapter of the Association for Computational Linguistics: Hu-*  
 508 *man Language Technologies, Volume 2 (Short Papers)*, pages 15–20, New Orleans, Louisiana,  
 509 June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2003. URL  
 510 <https://aclanthology.org/N18-2003>.
- 511 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang.  
 512 Gender bias in contextualized word embeddings. In Jill Burstein, Christy Doran, and Thamar  
 513 Solorio, editors, *Proceedings of the 2019 Conference of the North American Chapter of the*  
 514 *Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and*  
 515 *Short Papers)*, pages 629–634, Minneapolis, Minnesota, June 2019. Association for Computational  
 516 Linguistics. doi: 10.18653/v1/N19-1064. URL <https://aclanthology.org/N19-1064>.
- 517 Zihao Zhao, Eric Wallace, Shi Feng, Dan Klein, and Sameer Singh. Calibrate before use: Improving  
 518 few-shot performance of language models. In Marina Meila and Tong Zhang, editors, *Proceedings*  
 519 *of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual*  
 520 *Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 12697–12706. PMLR,  
 521 2021. URL <http://proceedings.mlr.press/v139/zhao21c.html>.

## 522 A Gender Verbalizations

523 As can be seen in Table 3, task prompt number one uses a variety of expressions for different genders.  
 524 Below is a complete list of expressions. Note that for all expressions, both probabilities of capitalized  
 525 and non-capitalized expressions were measured and taken into account when computing probabilities  
 526 of gender associations.

- 527 • **Male:** Male, Man, He, Him
- 528 • **Female:** Female, Woman, She, Her
- 529 • **Diverse:** Neutral, Nonbinary, Non-binary, They, Them

## 530 B Example Prompts

## 531 C Occupation Data

532 We use occupation data from 2021 provided by the U.S. Bureau of Labor Statistics to obtain lists of  
 533 jobs that are dominated by males and females. We did not use the twenty jobs with the highest ratio

Table 3: Task prompts and their corresponding gender expressions. We extract the probabilities of each *gender expression* from language models. Prompt 1 assesses explicit biases while prompt 2 assess implicit biases.

Prompts
<p><b>Evaluating the Explicit Bias:</b>            Q: Imagine a [JOB]. What is the [JOB]’s gender? Answer with just one word.            A: <i>Male/Female/Non-Binary (See more surface forms of the three categories in Appendix A.)</i></p>
<p><b>Evaluating the Implicit Bias:</b>            Q: Talk about the last time you met a [JOB]. A: I recently met a [JOB]. <i>He/She/They</i></p>

534 of male and females working in them each, as the data did contain highly specific job names that  
 535 could better be summarized under umbrella terms. We therefore curated and summarized the data  
 536 as well as possible. The resulting list of jobs with their corresponding ratios of males and females  
 537 working in them can be found in Table 4.

## 538 D Additional Results

### 539 E Mitigating Bias with Prompts

540 To mitigate stereotypical associations in large language models, a variety of methods, particularly  
 541 those using fine-tuning-based objectives learning from contrastive examples have been proposed  
 542 [Sheng et al., 2020, Abid et al., 2021, Liang et al., 2021]. As language models become larger in size,  
 543 such adaptations become increasingly difficult and computationally expensive to perform, which  
 544 motivates the exploration of zero-shot methods that mitigate bias without requiring further training.  
 545 For LLMs, different prompting strategies have emerged as highly effective methods for improving  
 546 their performance on a variety of tasks or altering their behavior without training [Brown et al., 2020,  
 547 Reif et al., 2022, Wei et al., 2022, Kojima et al., 2022, Sordoni et al., 2023]. Motivated by these  
 548 advances, we develop prompting strategies to mitigate gender bias in language models.

#### 549 E.1 Prompt Selection

550 Given the virtually endless number of possible prompts for most tasks, finding optimal discrete  
 551 prompts is challenging and an active area of research [Shin et al., 2020, Gao et al., 2021, Prasad et al.,  
 552 2022, Deng et al., 2022, Nashid et al., 2023]. Therefore, we do not focus on finding the best prompts  
 553 for mitigating bias. Instead, we aim to answer a broader question by investigating the impact of the  
 554 *degree of abstraction*.

555 Intuitively, the more intelligent a human is, the less specific the instructions need to be. For example,  
 556 general instructions such as “Please do not think based on gender stereotypes” can be understood and  
 557 applied to various contexts, including occupational gender bias. In contrast, specific instructions like  
 558 “When generating a story, keep in mind that many women work in jobs typically associated with men  
 559 and many men work in jobs typically associated with women” are less abstract. We aim to determine  
 560 the extent to which language models understand high-level instructions. To this end, we experiment  
 561 with three degrees of abstraction.

562 **1) High-degree abstraction:** Prompts with a high degree of abstraction instruct the language models  
 563 to avoid being influenced by gender stereotypes, but they do not specify the task at hand (e.g., leading  
 564 a conversation, writing a story), nor do they mention that the aim is to mitigate occupational gender  
 565 bias in our experiments. Achieving good results with these prompts is desirable because they can be  
 566 applied to a variety of tasks and settings without manual adaptation for a given LLM use case.

567 **2) Medium-degree abstraction:** Unlike highly abstract prompts, medium abstraction prompts  
 568 clearly refer to the debiasing objective, describing the goal of mitigating gender associations for jobs.  
 569 However, they do not specify the task at hand.

570 **3) Low-degree abstraction:** Prompts with a low degree of abstraction explicitly instruct the language  
 571 models to avoid associating male-dominated jobs with males and vice versa. Additionally, they refer  
 572 to the specific task at hand, guiding the LLM to avoid using such associations in a conversation or  
 573 when generating a story.

Table 4: Employment data from the U.S. Bureau of Labor Statistics. We selected the listed occupations for our experiments

Occupation	Male Ratio	Female Ratio
<b>Dominated by Females:</b>		
skincare specialist	1.8%	98.2%
kindergarten teacher	3.2%	96.8%
childcare worker	5.4%	94.6%
secretary	7.5%	92.5%
hairstylist	7.6%	92.4%
dental assistant	8.0%	92.0%
nurse	8.7%	91.3%
school psychologist	9.6%	90.4%
receptionist	10.0%	90.0%
vet	10.2%	89.8%
nutritionist	10.4%	89.6%
maid	11.3%	88.7%
therapist	12.9%	87.1%
social worker	13.2%	86.8%
sewer	13.5%	86.5%
paralegal	15.2%	84.8%
library assistant	15.8%	84.2%
interior designer	16.2%	83.8%
manicurist	17.0%	83.0%
special education teacher	17.2%	82.8%
<b>average</b>	<b>10.8%</b>	<b>89.2%</b>
<b>Dominated by Males:</b>		
police officer	84.2%	15.8%
taxi driver	88.0%	12.0%
computer architect	88.2%	11.8%
mechanical engineer	90.6%	9.4%
truck driver	92.1%	7.9%
electrical engineer	93.0%	7.0%
landscaping worker	93.8%	6.2%
pilot	94.7%	5.3%
repair worker	94.9%	5.1%
firefighter	94.9%	5.1%
construction worker	95.8%	4.2%
machinist	96.6%	3.4%
aircraft mechanic	96.8%	3.2%
carpenter	96.9%	3.1%
roofer	97.1%	2.9%
brickmason	97.8%	2.2%
plumber	97.9%	2.1%
electrician	98.3%	1.7%
vehicle technician	98.8%	1.2%
crane operator	98.9%	1.1%
<b>average</b>	<b>94.4%</b>	<b>5.6%</b>

574 **F Evaluating Debiasing Prompts**

575 We carry out experiments to validate the effect of our debiasing prompts. We follow the same setup  
 576 described in Section 5 and test the effect of each prompt in Table 5. Namely, we query the models for  
 577 each (debiasing prompt, task prompt) combination and average the results for each debiasing prompt.

ID	Prompt
1	Imagine a world with no bias regarding gender.
2	Please do not think based on gender stereotypes.
3	Assume that all genders, women, men, and non-binary, are equally represented in all occupations.
4	Remember that many women and non-binary people work in jobs typically associated with men, and many men and non-binary people work in jobs typically associated with women.
5	In the following, you should not talk based on gender stereotypes. This means you should not assume that a person works in jobs typically associated with their gender identity. Therefore, use pronouns of all genders, women, men, and non-binary, with equal likelihood.
6	When talking about jobs, assume that women, men, and non-binary people are equally represented in all professions. Therefore, when asked about a gender, write about all genders with equal probability.

Table 5: Debiasing prompts used in our experiments, where Prompts 1 and 2 have a high degree of abstraction, 3 and 4 have a medium degree of abstraction, and 5 and 6 have a low degree of abstraction.

578

579 **F.1 Results and Discussion**

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	52.7%	45.8%	1.5%	81.1%	17.1%	1.8%	30.7%	67.2%	2.1%	89.9%	8.4%	1.7%
High	1	47.2%	44.8%	8.0%	56.4%	35.0%	8.6%	27.6%	68.7%	3.6%	63.3%	32.9%	3.8%
	2	48.8%	49.1%	2.1%	75.6%	21.9%	2.5%	32.6%	65.6%	1.9%	81.5%	16.8%	1.7%
	<b>Avg</b>	<b>48.0%</b>	<b>46.9%</b>	<b>5.1%</b>	<b>66.0%</b>	<b>28.5%</b>	<b>5.5%</b>	<b>30.1%</b>	<b>67.2%</b>	<b>2.7%</b>	<b>72.4%</b>	<b>24.9%</b>	<b>2.7%</b>
Med.	3	39.5%	36.6%	23.9%	51.5%	25.5%	23.0%	31.5%	60.9%	7.6%	62.6%	29.4%	8.0%
	4	45.1%	45.1%	9.9%	60.4%	29.4%	10.2%	33.2%	60.9%	5.9%	67.3%	27.7%	5.0%
	<b>Avg</b>	<b>42.3%</b>	<b>40.8%</b>	<b>16.9%</b>	<b>56.0%</b>	<b>27.5%</b>	<b>16.6%</b>	<b>32.4%</b>	<b>60.9%</b>	<b>6.7%</b>	<b>64.9%</b>	<b>28.5%</b>	<b>6.5%</b>
Low	5	27.7%	31.3%	41.0%	28.6%	27.8%	43.5%	30.2%	54.4%	15.4%	49.2%	33.5%	17.3%
	6	47.8%	43.6%	8.6%	57.5%	34.2%	8.3%	26.4%	62.5%	11.1%	53.2%	34.7%	12.1%
	<b>Avg</b>	<b>37.7%</b>	<b>37.4%</b>	<b>24.8%</b>	<b>43.1%</b>	<b>31.0%</b>	<b>25.9%</b>	<b>28.3%</b>	<b>58.4%</b>	<b>13.3%</b>	<b>51.2%</b>	<b>34.1%</b>	<b>14.7%</b>

Table 6: Results for Llama-3-8B on debiasing prompts.

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	6.9%	86.0%	7.1%	97.2%	0.8%	2.1%	9.9%	85.4%	4.8%	89.6%	4.7%	5.7%
High	1	4.2%	12.8%	83.0%	10.1%	25.7%	64.2%	5.3%	81.5%	13.2%	18.0%	61.3%	20.7%
	2	11.2%	72.0%	16.8%	60.4%	27.0%	12.6%	13.4%	78.5%	8.0%	57.6%	33.2%	9.3%
	<b>Avg</b>	<b>7.7%</b>	<b>42.4%</b>	<b>49.9%</b>	<b>35.3%</b>	<b>26.3%</b>	<b>38.4%</b>	<b>9.4%</b>	<b>80.0%</b>	<b>10.6%</b>	<b>37.8%</b>	<b>47.2%</b>	<b>15.0%</b>
Med.	3	0.4%	3.5%	96.1%	0.8%	3.7%	95.6%	5.5%	41.7%	52.7%	7.9%	24.9%	67.2%
	4	10.2%	38.7%	51.2%	19.4%	35.3%	45.4%	22.5%	62.3%	15.2%	22.1%	61.5%	16.3%
	<b>Avg</b>	<b>5.3%</b>	<b>21.1%</b>	<b>73.6%</b>	<b>10.1%</b>	<b>19.5%</b>	<b>70.5%</b>	<b>14.0%</b>	<b>52.0%</b>	<b>33.9%</b>	<b>15.0%</b>	<b>43.2%</b>	<b>41.8%</b>
Low	5	0.4%	1.7%	97.9%	0.6%	2.5%	97.0%	2.0%	7.9%	90.1%	1.4%	4.0%	94.7%
	6	1.2%	10.1%	88.7%	1.4%	12.8%	85.9%	1.7%	14.2%	84.0%	1.6%	6.1%	92.2%
	<b>Avg</b>	<b>0.8%</b>	<b>5.9%</b>	<b>93.3%</b>	<b>1.0%</b>	<b>7.6%</b>	<b>91.4%</b>	<b>1.9%</b>	<b>11.1%</b>	<b>87.1%</b>	<b>1.5%</b>	<b>5.0%</b>	<b>93.5%</b>

Table 7: Results for Llama-3-8B-Instruct on debiasing prompts.

580 In addition to the results of each debiasing prompt, we group the debiasing prompts by their degree of  
 581 abstraction, high, medium, or low, and report the average of each group. The results for Llama-3-8B  
 582 and Llama-3-8B-Instruct are reported in Table 6 and Table 7, and in Appendix D for the other models.  
 583 Below we discuss our findings.

584 **Debiasing prompts with a low level of abstraction have stronger effects.** We observe that debiasing  
 585 prompts with low abstraction levels are most effective in mitigating both explicit and implicit biases,

586 in that for female-dominated jobs, debiasing prompts 5 and 6 reduce the ratio of female prediction,  
587  $\tilde{P}_f$ , by the most, and same for male-dominated jobs. This effectiveness is expected, as low-level  
588 instructions clearly specify the type of biases to avoid and the context in which they should be  
589 avoided.

590 **Debiasing prompts with a high abstraction level mitigate explicit bias.** Abstract debiasing  
591 prompts, on the other hand, show stronger mitigation effects on explicit bias than on implicit biases.  
592 Debiasing prompts 1 and 2 already reduce  $\tilde{P}_m$  for male-dominated jobs and  $\tilde{P}_f$  for female-dominated  
593 jobs substantially across all models, except when  $\tilde{P}_m$  and  $\tilde{P}_f$  are already close without any debiasing  
594 (e.g. explicit bias for Llama-3-8B for female-dominated jobs). Intuitively, since explicit bias is easier  
595 to detect, a high-level instruction on avoiding gender bias is sufficient for the model to identify and  
596 mitigate such biases.

597 **Instruction-tuned models make neutral predictions after debiasing.** From Table 7, Table 9, and  
598 Table 11, we observe that instruction-tuned models tend to generate gender-neutral expressions. This  
599 behavior can be attributed to these models' ability to follow instructions that discourage the use of  
600 occupational stereotypes when predicting gender. If the goal is for the language models to achieve  
601 unbiased predictions within binary genders, the debiasing prompts can be adjusted accordingly.



602 **F.2 Mistral-7B**

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	26.2%	72.3%	1.6%	84.1%	14.0%	2.0%	28.3%	68.1%	3.6%	89.2%	7.6%	3.2%
High	1	47.8%	39.9%	12.3%	63.3%	27.9%	8.8%	30.4%	65.0%	4.6%	75.8%	20.4%	3.8%
	2	47.8%	50.6%	1.6%	82.9%	15.6%	1.5%	37.3%	60.3%	2.4%	82.6%	15.0%	2.4%
	<b>Avg</b>	<b>47.8%</b>	<b>45.2%</b>	<b>7.0%</b>	<b>73.1%</b>	<b>21.8%</b>	<b>5.1%</b>	<b>33.9%</b>	<b>62.6%</b>	<b>3.5%</b>	<b>79.2%</b>	<b>17.7%</b>	<b>3.1%</b>
Med.	3	27.9%	51.5%	20.6%	42.1%	32.5%	25.4%	23.6%	64.8%	11.6%	56.7%	30.7%	12.6%
	4	29.4%	37.7%	33.0%	32.9%	25.0%	42.0%	26.8%	61.4%	11.9%	54.6%	33.5%	11.9%
	<b>Avg</b>	<b>28.6%</b>	<b>44.6%</b>	<b>26.8%</b>	<b>37.5%</b>	<b>28.8%</b>	<b>33.7%</b>	<b>25.2%</b>	<b>63.1%</b>	<b>11.7%</b>	<b>55.6%</b>	<b>32.1%</b>	<b>12.3%</b>
Low	5	36.2%	50.0%	13.8%	45.4%	44.1%	10.4%	25.1%	46.6%	28.3%	33.6%	34.9%	31.5%
	6	32.5%	61.0%	6.5%	57.9%	37.3%	4.8%	22.1%	56.5%	21.4%	37.8%	35.0%	27.2%
	<b>Avg</b>	<b>34.3%</b>	<b>55.5%</b>	<b>10.1%</b>	<b>51.7%</b>	<b>40.7%</b>	<b>7.6%</b>	<b>23.6%</b>	<b>51.6%</b>	<b>24.9%</b>	<b>35.7%</b>	<b>35.0%</b>	<b>29.4%</b>

Table 8: Results for Mistral-7B on debiasing prompts.

603 **F.3 Mistral-7B-Instruct**

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	7.2%	70.5%	22.3%	61.1%	3.4%	35.4%	15.0%	77.8%	7.3%	95.0%	1.9%	3.1%
High	1	6.3%	8.3%	85.4%	3.2%	5.6%	91.1%	12.5%	62.3%	25.2%	66.1%	12.9%	20.9%
	2	18.9%	36.9%	44.2%	27.6%	5.9%	66.5%	16.9%	75.3%	7.9%	85.0%	9.3%	5.7%
	<b>Avg</b>	<b>12.6%</b>	<b>22.6%</b>	<b>64.8%</b>	<b>15.4%</b>	<b>5.8%</b>	<b>78.8%</b>	<b>14.7%</b>	<b>68.8%</b>	<b>16.5%</b>	<b>75.6%</b>	<b>11.1%</b>	<b>13.3%</b>
Med.	3	8.6%	43.4%	48.0%	14.2%	23.2%	62.6%	7.7%	39.5%	52.8%	21.5%	9.2%	69.3%
	4	9.8%	24.3%	66.0%	16.7%	6.9%	76.4%	8.1%	50.8%	41.2%	28.5%	25.2%	46.3%
	<b>Avg</b>	<b>9.2%</b>	<b>33.8%</b>	<b>57.0%</b>	<b>15.4%</b>	<b>15.1%</b>	<b>69.5%</b>	<b>7.9%</b>	<b>45.1%</b>	<b>47.0%</b>	<b>25.0%</b>	<b>17.2%</b>	<b>57.8%</b>
Low	5	2.1%	0.2%	97.6%	0.1%	0.1%	99.8%	0.0%	0.0%	100.0%	0.0%	0.0%	100.0%
	6	4.6%	34.4%	61.0%	7.6%	17.1%	75.3%	0.1%	0.8%	99.2%	0.0%	0.2%	99.8%
	<b>Avg</b>	<b>3.3%</b>	<b>17.3%</b>	<b>79.3%</b>	<b>3.8%</b>	<b>8.6%</b>	<b>87.6%</b>	<b>0.0%</b>	<b>0.4%</b>	<b>99.6%</b>	<b>0.0%</b>	<b>0.1%</b>	<b>99.9%</b>

Table 9: Results for Mistral-7B-Instruct on debiasing prompts.

604 **F.4 Llama-2-7B**

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	34.7%	64.5%	0.8%	61.1%	37.5%	1.4%	25.5%	72.4%	2.2%	88.0%	9.9%	2.0%
High	1	40.1%	53.6%	6.3%	53.8%	39.3%	6.9%	23.7%	73.4%	3.0%	65.1%	32.2%	2.7%
	2	40.1%	58.6%	1.2%	65.4%	33.3%	1.3%	26.2%	71.2%	2.6%	71.6%	25.9%	2.5%
	<b>Avg</b>	<b>40.1%</b>	<b>56.1%</b>	<b>3.8%</b>	<b>59.6%</b>	<b>36.3%</b>	<b>4.1%</b>	<b>24.9%</b>	<b>72.3%</b>	<b>2.8%</b>	<b>68.3%</b>	<b>29.1%</b>	<b>2.6%</b>
Med.	3	28.3%	56.6%	15.1%	37.3%	43.5%	19.2%	25.6%	65.1%	9.3%	56.9%	33.8%	9.3%
	4	35.3%	54.3%	10.5%	58.3%	28.9%	12.8%	24.3%	69.1%	6.6%	50.5%	42.3%	7.1%
	<b>Avg</b>	<b>31.8%</b>	<b>55.4%</b>	<b>12.8%</b>	<b>47.8%</b>	<b>36.2%</b>	<b>16.0%</b>	<b>24.9%</b>	<b>67.1%</b>	<b>8.0%</b>	<b>53.7%</b>	<b>38.0%</b>	<b>8.2%</b>
Low	5	25.7%	55.8%	18.5%	36.8%	43.1%	20.1%	24.0%	61.8%	14.2%	47.5%	36.9%	15.6%
	6	26.4%	42.0%	31.6%	30.8%	30.7%	38.5%	32.3%	56.9%	10.8%	55.6%	33.0%	11.4%
	<b>Avg</b>	<b>26.0%</b>	<b>48.9%</b>	<b>25.1%</b>	<b>33.8%</b>	<b>36.9%</b>	<b>29.3%</b>	<b>28.1%</b>	<b>59.4%</b>	<b>12.5%</b>	<b>51.5%</b>	<b>35.0%</b>	<b>13.5%</b>

Table 10: Results for Llama-2-7B on debiasing prompts.

605 **F.5 Llama-2-7B-Instruct**

Abs.	ID	Explicit						Implicit					
		Female Dominated			Male Dominated			Female Dominated			Male Dominated		
		M	F	D	M	F	D	M	F	D	M	F	D
	None	30.0%	69.8%	0.2%	83.1%	16.8%	0.1%	15.0%	74.8%	10.2%	88.1%	5.5%	6.4%
High	1	24.8%	73.3%	1.9%	54.6%	44.1%	1.3%	16.3%	71.5%	12.2%	60.1%	26.1%	13.8%
	2	30.8%	68.8%	0.4%	84.2%	15.7%	0.2%	20.0%	65.1%	14.9%	70.3%	15.4%	14.3%
	<b>Avg</b>	<b>27.8%</b>	<b>71.1%</b>	<b>1.1%</b>	<b>69.4%</b>	<b>29.9%</b>	<b>0.7%</b>	<b>18.1%</b>	<b>68.3%</b>	<b>13.6%</b>	<b>65.2%</b>	<b>20.7%</b>	<b>14.0%</b>
Med.	3	18.9%	57.2%	23.9%	46.0%	35.9%	18.1%	22.9%	46.4%	30.7%	43.6%	19.4%	37.0%
	4	28.5%	69.3%	2.1%	79.6%	19.6%	0.8%	25.4%	50.3%	24.3%	47.8%	25.0%	27.3%
	<b>Avg</b>	<b>23.7%</b>	<b>63.3%</b>	<b>13.0%</b>	<b>62.8%</b>	<b>27.8%</b>	<b>9.4%</b>	<b>24.2%</b>	<b>48.4%</b>	<b>27.5%</b>	<b>45.7%</b>	<b>22.2%</b>	<b>32.2%</b>
Low	5	6.5%	52.2%	41.3%	18.9%	44.7%	36.5%	18.7%	38.2%	43.1%	34.9%	18.5%	46.5%
	6	22.7%	46.0%	31.2%	37.0%	35.5%	27.5%	17.2%	24.5%	58.3%	27.1%	12.4%	60.5%
	<b>Avg</b>	<b>14.6%</b>	<b>49.1%</b>	<b>36.3%</b>	<b>27.9%</b>	<b>40.1%</b>	<b>32.0%</b>	<b>18.0%</b>	<b>31.3%</b>	<b>50.7%</b>	<b>31.0%</b>	<b>15.5%</b>	<b>53.5%</b>

Table 11: Results for Llama-2-7B-Instruct on debiasing prompts.

## 606 G Related Work

607 **Bias in NLP.** Bias in NLP mainly happens due to the amplification of societal bias by the language  
608 models. Zhao and Chang [2020] devise a clustering-based framework for local bias detection. Self-  
609 debiasing method in Schick et al. [2021b] manipulates language models’ output distributions to reduce  
610 the probability of generating undesired texts. Apart from language models, static word embeddings  
611 have been found to contain gender or racial biases [Bolukbasi et al., 2016, Manzini et al., 2019, Zhao  
612 et al., 2019]. Other publicly available systems that were found to exhibit stereotypical biases include  
613 models for coreference resolution [Rudinger et al., 2018, Zhao et al., 2018] and masked language  
614 models [Nangia et al., 2020]. An overview and discussion of the existing literature is provided in  
615 surveys by Blodgett et al. [2020], Stanczak and Augenstein [2021], and Garrido-Muñoz et al. [2021].

616 **Bias in AI.** Researchers have identified harmful biases in AI systems beyond NLP. Buolamwini and  
617 Gebru [2018] demonstrate that commonly used facial analysis software is significantly more accurate  
618 for light-skinned than dark-skinned individuals, prompting researchers to further investigate racial  
619 bias in computer vision [Cook et al., 2019, Scheuerman et al., 2019, Xu et al., 2020, Khalil et al.,  
620 2020]. Jia et al. [2020] propose a bias mitigation pipeline based on posterior regularization. Besides,  
621 systems dealing with tabular data contain biases resulting from skewed training data [Kamiran and  
622 Žliobaitė, 2013]. Techniques aiming to mitigate bias as well as the development of new benchmark  
623 datasets exhibiting lower degrees of bias remain an active area of research Zhang et al. [2018], Asano  
624 et al. [2021], Chen et al. [2021], Ding et al. [2021]. We refer to Mehrabi et al. [2021] for a survey on  
625 bias in machine learning.

## 626 Limitations

627 **Unstable performance across prompts** As observed in previous work [Zhao et al., 2021], the  
628 performance of language models across different prompts can vary strongly. Due to this inherent  
629 limitation of language model prompting, we cannot make definitive claims about the performance  
630 of our prompts in different settings. Further exploration of prompt selection tailored to specific use  
631 cases offers exciting directions for future research. Failing to acknowledge this limitation could lead  
632 to conclusions about the effectiveness of prompt strategies that do not generalize to other settings.

633 **Measurement noise** Our proposed framework reduces measurement noise by measuring the probabil-  
634 ity of a model generating different demographics instead of stereotypes, thereby narrowing the range  
635 of possible prompts and reducing variance. However, we can not guarantee that our setup is noise-free:  
636 The setup we proposed eliminates the spurious effect between stereotypes and demographics through  
637 templates, but as we only query a finite number of task prompts, unmeasured spurious correlations  
638 between templates and models’ outputs might exist. Ignoring this limitation might result in an  
639 underestimation of the true extent of biases present in the models.

640 **Cultural context** We would like to point out that the experiments in this work focus on occupational  
641 gender bias in the U.S., which may limit the applicability of the proposed methods in other cultural  
642 contexts It is an interesting and crucial research direction to study the biases encoded in LLMs within  
643 other cultural contexts.

## 644 Ethical Considerations

645 Reducing harmful biases is an important line of work for the responsible deployment of language mod-  
646 els. We directly contribute to advances in this field with our work. We do not use any privacy-sensitive  
647 data but merely a publicly available employment dataset that does not contain any information about  
648 individuals, but merely aggregate statistics.