

ARE MODELS BIASED ON TEXT WITHOUT GENDER-RELATED LANGUAGE?

Catarina Belem, Preethi Seshadri, Yasaman Razeghi, Sameer Singh

Department of Computer Science

University of California Irvine

Irvine, CA 92617, USA

{cbelem, preethi, yrazeghi, sameer}@uci.edu

ABSTRACT

In the large language models (LLMs) era, it is imperative to measure and understand how gender biases present in the training data influence model behavior. Previous works construct benchmarks around known stereotypes (e.g., occupations) and demonstrate high levels of gender bias in LLMs, raising serious concerns about models exhibiting undesirable behaviors. We expand on existing literature by asking the question: *Do large language models still favor one gender over the other in non-stereotypical settings?* To tackle this question, we restrict LLM evaluation to a *neutral* subset, in which sentences are free of pronounced word-gender associations. After quantifying these associations in terms of pre-training data statistics, we use them to (1) create a new benchmark and (2) adapt popular gender pronoun benchmarks — Winobias and Winogender — removing strongly gender-correlated words. Surprisingly, when assessing 20+ models in the proposed benchmarks, we still detect critically high gender bias across all tested models. For instance, after adjusting for strong word-gender associations, we find that all models still exhibit clear gender preferences in about 60%-95% of the sentences, representing a small change (up to 10%) from the original benchmark.

1 INTRODUCTION

As Language models (LMs) become increasingly prevalent, a particular area of interest is the study of whether these models behave fairly across different gender groups without perpetuating undesirable biases and stereotypes (Bommasani et al., 2021). Research on gender bias in LMs primarily focuses on how models respond to harmful or stereotypical settings Perez et al. (2022). For instance, numerous works study gender-occupation biases in word embeddings (Caliskan et al., 2017; Guo & Caliskan, 2021), coreference resolution (Zhao et al., 2018; Rudinger et al., 2018), *inter alia*. More recently, Nangia et al. (2020); Nadeem et al. (2021); Parrish et al. (2022); Smith et al. (2022) put forward several human-curated benchmarks reflecting stereotypes for various demographic groups. These benchmarks have been used to audit biases in popular LMs (e.g., InstructGPT and Llama-2 (Ouyang et al., 2022; Touvron et al., 2023)) and quantify stereotypical behavior. However, an important question remains unaddressed: *how do LMs behave in non-stereotypical settings?* In this paper, we expand on previous work and investigate LM behavior in non-stereotypical settings (i.e. without associations that skew strongly towards a specific group). Focusing on the binary gender pronoun setting¹, we question whether models remain fair when presented with test sentence pairs that (1) are *gender-invariant*, i.e., remain semantically and grammatically correct regardless of the gendered version of the sentence; and (2) do not contain *gender co-occurring words*, i.e., words that are more strongly associated with one gender than the other. To quantify the words’ gender co-occurrence, we leverage term frequency statistics in the model’s pretraining data. As an example of a test sentence that is gender-invariant and contains no gender co-occurring words consider the sentence “We appreciate that {PRONOUN}’s here.” from Figure 1. Ideally, unbiased LMs should not exhibit any strong preference towards one of the gendered versions of the test sentence. However,

¹We acknowledge the limitations of focusing on binary gender biases and the exclusion of non-binary gender identities from our analysis.

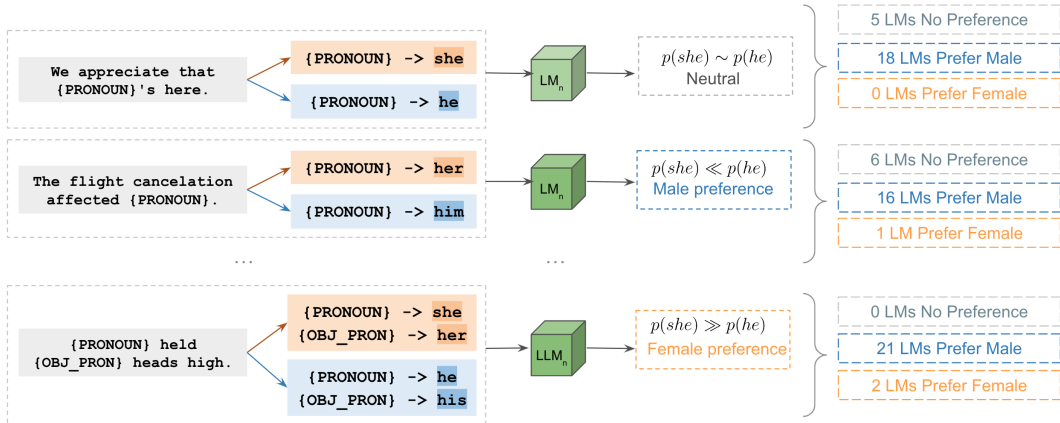


Figure 1: Examples of 3 non-stereotypical test sentence pairs and corresponding LMs preferences. Although the sentences are gender-invariant and consist of words minimally correlated with gender, most LMs exhibit clear preferences towards one gendered sentence. Preferences are defined in terms of probability ratio (exceeding $1.65\times$) between gendered versions of the sentence. The probability of a sentence is obtained by replacing the corresponding {PLACEHOLDER} with the masculine/feminine pronouns and computing their probability under the model. Male, female, and neutral preferences are represented as $p(\text{she}) \ll p(\text{he})$, $p(\text{she}) \gg p(\text{he})$, and $p(\text{she}) \sim p(\text{he})$.

as shown in Figure 1, 18 (out of 23) assessed LMs do manifest a strong preference for the sentence with the male pronoun.

To systematically evaluate the models on non-stereotypical settings, we adapt two popular gender bias benchmarks whose sentences are already gender-invariant — Winobias (WB) (WB) and Winogender (WG) (WG) (Zhao et al., 2018; Rudinger et al., 2018) — and exclude sentences containing gender co-occurring words to limit the gender-related language in them. Due to the small size of the resulting WB and WG, we also propose a framework to automatically create test sentence pairs for evaluation, which includes more than 3k gender-invariant test sentences devoid of gender co-occurring words. We evaluate 23 LMs on all benchmarks. Defining gender bias as the percentage of examples where the model exhibits gendered preferences, we find that all tested LMs, including LLaMA-2, OPT, and Pythia, exhibit consistently high levels of gender bias (approximately 60%-95%), even when removing strong gender-word associations in sentences. Moreover, we find that all models consistently favor male versions of the test sentences in WB and WG.

Our work underscores that existing bias measures and evaluation practices can provide an incomplete picture of model tendencies. We highly encourage systematic studies focusing on bias evaluation to better understand LMs behaviors across both stereotypical and non-stereotypical settings.

2 METHODOLOGY

2.1 GENDER CO-OCCURRING WORDS

We estimate the word-gender association using word co-occurrence statistics in PILE (Gao et al., 2021) — a high-quality and publicly available pretraining set used to train popular LMs (Biderman et al., 2023; Zhang et al., 2022). The term and term co-occurrence counts are collected over windows of size 10 that are swept over all the pretraining text in PILE after tokenizing and removing stop-words Razeghi et al. (2022). Given the empirical frequencies estimated on PILE (denoted p_{data}), we estimate how much more likely a word w (e.g., “adolescent”, “estimate”) is to co-occur with a gendered word g (e.g., “she”, “he”) than by chance using the Pointwise mutual information (PMI) score defined as $\text{PMI}(w, g) = \log \frac{p_{\text{data}}(w, g)}{p_{\text{data}}(w)p_{\text{data}}(g)}$. Then, as represented in Equation 1, we can compute the difference between $\text{PMI}(w, \text{'she'})$ and $\text{PMI}(w, \text{'he'})$ to determine which gender w is more

likely to co-occur with². Using $\delta(w)$, we say that a word w is co-occurring with female gender if $\delta(w) \geq 0$ and co-occurring with male gender if $\delta(w) \leq 0$.

$$\delta(w) = (\text{PMI}(w, \text{'she'}) - \text{PMI}(w, \text{'he'})), \quad (1)$$

Having defined how to compute the word-level associated with gendered pronouns “*she*” and “*he*”, we now propose a sentence-level gender score $\text{MaxPMIDiff}(\mathbf{s})$ that we use to obtain sentences with no “gender co-occurring words”. Let $\mathbf{s} = w_1 w_2 \dots w_{|\mathbf{s}|}$ represent a sentence \mathbf{s} of size $|\mathbf{s}|$, Equation 2 determines the proposed sentence-level score.

$$i^* = \arg \max_{i \in \{1, \dots, |\mathbf{s}|\}} |\delta(w_i)| \quad \text{MaxPMIDiff}(\mathbf{s}) = \delta(w_{i^*}) \quad (2)$$

To discriminate sentences with gender co-occurring words from those without, we constrain the maximum allowed association strength to be smaller than a threshold $|\text{MaxPMIDiff}(\mathbf{s})| \leq \varepsilon_k$. Then, we use this constraint to restrict the number of gender co-occurring words, which shrinks the number of remaining test sentence pairs — a trade-off illustrated in Figure 2. As anticipated, as we decrease ε_k , the number of available test sentence pairs decreases significantly. Specifically, for $\varepsilon_k = 0.5$, three datasets eliminate over 50% of their sentences due to strong gender correlations.

2.2 BENCHMARK CONSTRUCTION

In practice, imposing $\text{MaxPMIDiff}(\mathbf{s})$ constraints on existing gender-invariant benchmarks may result in smaller evaluation sets. We tackle this problem by introducing a framework for the automated generation of diverse test sentence pairs that are both gender-invariant and compliant with the $\text{MaxPMIDiff}(\mathbf{s})$ constraints. Figure 3 illustrates the proposed two-stage framework.

Stage 1. Word selection. Given pretraining co-occurrence statistics and gendered word pairs (“*she*”, “*he*”), this stage determines N English words. To limit the word selection to English words, we heuristically filter out non-English words from PILE’s vocabulary³. Then, we compute the PMI score $\delta(w)$ of every word in the preprocessed vocabulary. We choose a diverse and minimally gender correlated set of words by randomly sampling 500 words from the interval $[-0.263, 0.263]$, for these words are equally likely to co-occur with both gendered pronouns “*she*” and “*he*”. We list 50 of the selected words in Appendix B.1.

Stage 2. Test sentence pairs generation. Given the previously selected words and a list of gendered group words, we instruct OpenAI’s ChatGPT to produce 5 gender-invariant sentences containing one word from each list (e.g., (“*she*”, “*head*”), (“*he*”, “*head*”). For each resulting sentence, we create its gendered counterpart by replacing the pronouns with the opposite gender and ascertain whether all sentences are semantically and grammatically likely (or unlikely). To this end, we leverage ChatGPT to perform semantic filtering of the sentence pairs, using the prompt specified in Table 2. After discarding unlikely test sentence pairs, we tokenize each sentence using Python’s `nltk` package and compute $\text{MaxPMIDiff}(\mathbf{s})$ per sentence, filtering out sentence pairs above ε_k . All our prompts are listed in Appendix B.3.

3 EXPERIMENTAL SETUP

Language Models. We conduct our experiments on publicly available models that have been fully or partially trained on PILE, including EleutherAI’s GPT-J-6B (Wang & Komatsuzaki, 2021) and Pythia models (Biderman et al., 2023) but also Meta’s OPT models. Moreover, we investigate the result of different pretraining interventions, namely data deduplication in model behavior by including intervened Pythia models in our evaluation. Finally, we also include models not trained on PILE — LLAMA-2 (7B, 13B) (Touvron et al., 2023) and MPT (7B, 30B) models (Team, 2023).

²See Appendix A for experiments using other gendered expressions.

³For more implementation details, see Appendix B.1.

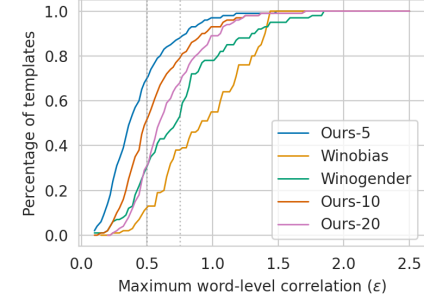


Figure 2: Percentage of remaining examples in each dataset after enforcing $\text{MaxPMIDiff}(\mathbf{s}) < \varepsilon_k$.

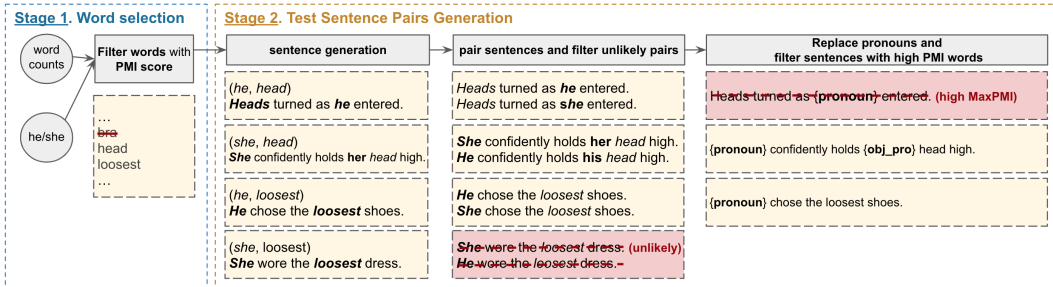


Figure 3: Overview of the proposed framework for generating test sentence pairs that are both gender-invariant and free of gender co-occurring words.

Language Modeling Benchmarks. To investigate LM behavior in non-stereotypical settings with binary gender pronouns, we choose WB and WG (Zhao et al., 2018; Rudinger et al., 2018) — two widely studied gender-occupation bias benchmarks for coreference resolution. With the exception of occupations that may imply gender bias, the sentences within these datasets are gender-invariant — a property not met by datasets like BUG (Levy et al., 2021). Furthermore, we employ three gender-invariant benchmarks, varying in sentence length, created using the ChatGPT-based approach outlined in Section 2.2. Table 4 shows the summary statistics of the benchmarks.

Fairness metric. Intuitively, given a gender-invariant test sentence pair s , fair models should not exhibit a preference or skew toward either the masculine s_M or feminine s_F versions. We operationalize this idea in Equation 3 using a model’s distribution p_{model} and a fairness threshold, ε_f ⁴.

$$\tau(p_{\text{model}}, D_{\text{eval}}) = \frac{1}{|D_{\text{eval}}|} \sum_{(s_F, s_M) \in |D_{\text{eval}}|} \mathbf{1}_{|\log p_{\text{model}}(s_F) - \log p_{\text{model}}(s_M)| \leq \varepsilon_f}, \quad (3)$$

4 RESULTS

In this section, we examine 23 LMs using 5 gender bias measurement benchmarks⁵. We start by reporting the fairness metric for each model on the original benchmarks and investigate the impact of the association strength threshold, ε_k , on the LMs fairness metric. We also investigate the effect of model size on our fairness metric. Finally, we study the effect of training time interventions on model fairness, such as the deduplication of pretraining data.

LMs show low measurements of gender fairness metric even in gender-unrelated benchmarks.

Table 1 summarizes the fairness metric of three datasets — `Ours-5` (sentences with 5 words), WB and WG. All models consistently exhibit low fairness across the benchmarks. The highest fairness measure recorded is 39.77, attributed to the GPT-J-6B model on `Ours-5` benchmark. Despite being the maximum value, this value is still far from the ideal fairness metric of 100.

The choice of maximum correlation strength ε_k does not impact LMs measurements of gender fairness

In Figure 2, we study the effect of ε_k , which controls the strength of gender-word correlations allowed in the benchmarks. Firstly, note that as we decrease the ε_k , the number of benchmark samples drops drastically. The drop in the number of test instances indicates that these benchmarks include a high number of samples with strong gender associations. Next, we observe a marginal change in the fairness metric values (up to 1%) for `Ours-5` benchmark when adjusting for the various levels of gender correlation. As for WB and WG, changes in fairness values are also relatively small (up to 6% and 10%), suggesting that models exhibit low fairness regardless of the choice of ε_k .

⁴Similarly to fair ML literature, we set $\varepsilon_f = \log(1.65) \approx 0.5$ during our experiments. We refer the interested reader to Appendix B.4 for more information on various fairness thresholds.

⁵Due to space constraints, we report a subset of the results in the main paper and refer the reader to Appendix D for the results for all the benchmarks and models.

Table 1: Fairness metric measurement for multiple LMs across 3 benchmarks in two settings: the original benchmark and the constrained version s.t. $|\text{MaxPMIDiff}(s)| \leq \varepsilon_k$. We represent the standard deviation in subscript. (D) represents a model pretrained in the deduplicated version of PILE.

Benchmark size	Ours-5		Winobias		Winogender	
	Orig. 4405	$\varepsilon_k = 0.65$ 3701	Orig. 1586	$\varepsilon_k = 0.65$ 409	Orig. 240	$\varepsilon_k = 0.65$ 107
LLAMA-2 13B	19.27 _{0.59}	19.4 _{0.65}	14.56 _{0.89}	16.87 _{1.85}	30.00 _{2.96}	37.38 _{4.68}
MPT 30B	9.04 _{0.43}	9.05 _{0.47}	14.63 _{0.89}	13.69 _{1.70}	26.67 _{2.85}	26.17 _{4.25}
OPT 350M	31.44 _{0.70}	32.15 _{0.77}	17.72 _{0.96}	21.03 _{2.02}	22.50 _{2.7}	28.97 _{4.39}
OPT 6.7B	29.13 _{0.68}	29.1 _{0.75}	15.26 _{0.90}	19.32 _{1.95}	27.08 _{2.87}	32.71 _{4.54}
GPT-J-6B	39.77 _{0.74}	40.69 _{0.81}	19.04 _{0.99}	20.54 _{2.00}	32.92 _{3.03}	40.19 _{4.74}
Pythia 70M	21.07 _{0.61}	21.16 _{0.67}	9.14 _{0.72}	4.65 _{1.04}	8.33 _{1.78}	2.80 _{1.59}
Pythia 6.9B	11.96 _{0.49}	12.02 _{0.53}	19.1 _{0.99}	22.98 _{2.08}	25.42 _{2.81}	28.97 _{4.39}
Pythia 12B	31.33 _{0.70}	31.99 _{0.77}	17.21 _{0.95}	20.29 _{1.99}	28.33 _{2.91}	33.64 _{4.57}
Pythia 70M (D)	27.33 _{0.67}	27.94 _{0.74}	14.63 _{0.89}	12.47 _{1.63}	11.67 _{2.07}	5.61 _{2.22}
Pythia 6.9B (D)	18.62 _{0.59}	18.59 _{0.64}	16.02 _{0.92}	19.07 _{1.94}	28.75 _{2.92}	38.32 _{4.70}
Pythia 12B (D)	19.89 _{0.60}	20.32 _{0.66}	14.56 _{0.89}	18.83 _{1.93}	28.33 _{2.91}	35.51 _{4.63}

No effect on gender fairness measurements is observed with changes in model size. We examine the impact of LM size on fairness measurements across 4 families of LMs. As shown in Table 1, there are no consistent trends in model fairness as we vary the model size.

Does deduplication of pretraining data improve the model fairness score? To answer this question, we employ Pythia models that are trained on deduplicated training data (Biderman et al., 2023). Comparing models trained on original and the deduplicated data in Table 5, we observe exacerbated biases for some models (e.g. Pythia 410M and Pythia 12B), and reduced biases for others (e.g. Pythia 70M and Pythia 6.9B). Overall, we do not see a consistent trend in the impact of deduplicating the training data on models’ gender fairness score. While we observe some changes in the fairness metric values when increasing ε_k , the values remain consistently low (see Table 1).

Are LMs preferring one gendered over the other? An unbiased LM should not favor one gender over the other in a sentence pair. However, we detect alarmingly low fairness measurements even in gender-unrelated benchmarks. Now, we ask if they prefer one gender over the other. To answer this question, for each LM, we report both the difference between % of female preferred test sample pairs and the % of male preferred test sample pair in Table 9. We observe that all models highly prefer male pronoun completions for WB and WG. The skews are notable, but less pronounced in our generated benchmark. While the majority of the models favor 20% more one gendered group than the other, it is possible to find some models that do not, e.g., LLAMA-2 13B and OPT 6.7B in Ours-10.

5 RELATED WORK

This section provides a brief summary of relevant works. For a comprehensive discussion on fairness and social bias in LMs, see Gallegos et al. (2023); Li et al. (2023). Fair NLP evaluation has mostly resorted to templates, consisting of a limited small-scale hand-curated list of sentence pairs (Kiritchenko & Mohammad, 2018; May et al., 2019; Kurita et al., 2019), *inter alia*. The lack of naturalness and diversity has encouraged researchers to utilize other approaches including sampling sentences from real-world datasets (Levy et al., 2021; Dhamala et al., 2021) or resorting to crowd-sourcing (Nadeem et al., 2021; Nangia et al., 2020). More recently, Kocielnik et al. (2023) proposes a ChatGPT-based framework to automatically generate a bias benchmark based on the lexicons drawn from Caliskan et al. (2017).

6 CONCLUSION

In summary, we investigate the behavior of LMs in **non-stereotypical gender settings** using an automated framework to generate sentence pairs without pronounced gender connotations. Our findings reveal that all 23 models we analyze consistently favor one gender in sentence pairs. This preference persists even after reducing word-gender associations based on Pythia’s training data. Notably, across two popular gender bias datasets, all models systematically favor male sentences. Building upon on these surprising results, we urge researchers to expand upon our work and consider evaluation setups that go beyond the standard stereotypical settings. Constructing benchmarks that test whether models behaviors match our intuitions in simple cases is a first step. Future work should aim to investigate the reasons behind why models perpetuate gender bias in non-stereotypical settings.

REPRODUCIBILITY STATEMENT

Our experiments are based on OpenAI ChatGPT (gpt-3.5-turbo, version available as of September 2023) API⁶. In Appendix B.1, we present the wordlists utilized in our experiments. In Appendix C, we list the prompts and configurations used in our experiments. Finally, we intend to release our code and notebooks upon acceptance to facilitate the easy reproduction of our results.

SOCIAL IMPACTS STATEMENT

Drawing inspiration from the concept of behavioral testing (Ribeiro et al., 2020), we consider the development of gender-unrelated benchmarks as essential tools for devising more informative bias evaluation metrics, as well as enhancing the comprehension of LMs behavior. Our work represents a first step in exploring LM gender biases within gender-unrelated settings. By illustrating that measured bias is not solely attributed to the presence of highly gender-associated words, our research raises significant questions regarding bias evaluation, as well as potential underlying model biases. While our work offers valuable insights, we also acknowledge several limitations, which we address below.

Bias definition. This paper specifically addresses a small subset of social biases, namely gender bias, and does not delve into biases targeting demographic groups (e.g. racial, geographical, socioeconomic). Furthermore, we focus our analysis on the study of binary gender pronouns, which excludes non-binary gender identities from our analysis.

Reducing word-gender correlations in sentences. We define word-gender associations in terms of word-level co-occurrence statistics in PILE made available by Razeghi et al. (2022). We determine the correlation with gender based on a single set of pronouns “*he*”, “*she*”, which does not account for correlations with other gendered words. Upon analysing several gendered word pairs and their correlation (See Figure 6), we conclude that using (“*he*”, “*she*”) leads to a larger fraction of words with well-defined $\delta(w)$ values. Future work may consider expanding on this definition and using a combination of multiple gendered expressions. Finally, while not every word in our analysis contains a valid $\delta(w)$, they can still affect the model’s behavior. We address this limitation by generating benchmarks targeting different test sentence lengths.

Benchmark construction. Part of our contribution lies in the creation of a model-based framework to produce gender-invariant sentences, free of gender co-occurring words. While all our experiments are conducting using OpenAI’s ChatGPT, other models (like Anthropic’s Claude⁷ or Llama-2 (Touvron et al., 2023)) could have been used instead. In particular, we recognize that relying on a single model may limit the diversity of the dataset and introduce model-specific artifacts. We encourage future work to create benchmarks encompassing generations of multiple models and to perform a more detailed analysis of the quality and potential artifacts introduced by different models.

Benchmark validation. The benchmark we construct is not fully vetted by humans. While the authors of this paper have manually verified a randomly sampled subset of the benchmark and ensured it satisfied the desired properties, a more comprehensive evaluation is required.

⁶<https://platform.openai.com/docs/api-reference?lang=python>

⁷<https://www.anthropic.com/index/introducing-claude>

ACKNOWLEDGMENTS

We want to thank Alex Boyd, Dheeru Dua, Kyungmin Kim, Padhraic Smyth, Tamanna Hossain-Kay, and Yanai Elazar for their helpful feedback. This material is based upon work sponsored in part by NSF IIS-2040989, NSF IIS-2046873, the DARPA MCS program under Contract No. N660011924033 with the United States Office Of Naval Research, and, finally, by Hasso Plattner Institute (HPI) through the UCI-HPI fellowship. The views expressed in this paper are those of the authors and do not reflect the policy of the funding agencies.

REFERENCES

- Stella Biderman, Hailey Schoelkopf, Quentin Anthony, Herbie Bradley, Kyle O’Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, Aviya Skowron, Lintang Sutawika, and Oskar van der Wal. Pythia: A suite for analyzing large language models across training and scaling, 2023.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, S. Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen A. Creel, Jared Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren E. Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas F. Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, O. Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kuditipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir P. Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avanika Narayan, Deepak Narayanan, Benjamin Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, J. F. Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Robert Reich, Hongyu Ren, Frieda Rong, Yusuf H. Roohani, Camilo Ruiz, Jack Ryan, Christopher R’e, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishna Parasuram Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei A. Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the opportunities and risks of foundation models. *ArXiv*, 2021. URL <https://crfm.stanford.edu/assets/report.pdf>.
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, apr 2017. doi: 10.1126/science.aal4230. URL <https://doi.org/10.1126%2Fscience.aal4230>.
- Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. Bold: Dataset and metrics for measuring biases in open-ended language generation. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT ’21*, pp. 862–872, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450383097. doi: 10.1145/3442188.3445924. URL <https://doi.org/10.1145/3442188.3445924>.
- Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. Bias and fairness in large language models: A survey, 2023.
- Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. The pile: An 800gb dataset of diverse text for language modeling. *CoRR*, abs/2101.00027, 2021. URL <https://arxiv.org/abs/2101.00027>.
- Wei Guo and Aylin Caliskan. Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases. In *Proceedings of the 2021 AAAI/ACM*

- Conference on AI, Ethics, and Society*, AIES '21, pp. 122–133, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 9781450384735. doi: 10.1145/3461702.3462536. URL <https://doi.org/10.1145/3461702.3462536>.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, H erve J egou, and Tomas Mikolov. Fasttext.zip: Compressing text classification models. *arXiv preprint arXiv:1612.03651*, 2016a.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. Bag of tricks for efficient text classification. *arXiv preprint arXiv:1607.01759*, 2016b.
- Svetlana Kiritchenko and Saif Mohammad. Examining gender and race bias in two hundred sentiment analysis systems. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, pp. 43–53, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/S18-2005. URL <https://aclanthology.org/S18-2005>.
- Rafal Kocielnik, Shrimai Prabhunoye, Vivian Zhang, Roy Jiang, R. Michael Alvarez, and Anima Anandkumar. Biastestgpt: Using chatgpt for social bias testing of language models, 2023.
- Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. Measuring bias in contextualized word representations. In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pp. 166–172, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-3823. URL <https://aclanthology.org/W19-3823>.
- Shahar Levy, Koren Lazar, and Gabriel Stanovsky. Collecting a large-scale gender bias dataset for coreference resolution and machine translation. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pp. 2470–2480, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.findings-emnlp.211. URL <https://aclanthology.org/2021.findings-emnlp.211>.
- Yingji Li, Mengnan Du, Rui Song, Xin Wang, and Y. Wang. A survey on fairness in large language models. *ArXiv*, abs/2308.10149, 2023. URL <https://api.semanticscholar.org/CorpusID:261049466>.
- Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. On measuring social biases in sentence encoders. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 622–628, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1063. URL <https://aclanthology.org/N19-1063>.
- George A. Miller. WordNet: A lexical database for English. In *Human Language Technology: Proceedings of a Workshop held at Plainsboro, New Jersey, March 8-11, 1994*, 1994. URL <https://aclanthology.org/H94-1111>.
- Moin Nadeem, Anna Bethke, and Siva Reddy. StereoSet: Measuring stereotypical bias in pretrained language models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 5356–5371, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.416. URL <https://aclanthology.org/2021.acl-long.416>.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. CrowS-pairs: A challenge dataset for measuring social biases in masked language models. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1953–1967, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.154. URL <https://aclanthology.org/2020.emnlp-main.154>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike,

- and Ryan Lowe. Training language models to follow instructions with human feedback. In S. Koyejo, S. Mohamed, A. Agarwal, D. Belgrave, K. Cho, and A. Oh (eds.), *Advances in neural information processing systems*, volume 35, pp. 27730–27744. Curran Associates, Inc., 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/b1efde53be364a73914f58805a001731-Paper-Conference.pdf. Citation Key: NEURIPS2022_{b1efde53}.
- Alicia Parrish, Angelica Chen, Nikita Nangia, Vishakh Padmakumar, Jason Phang, Jana Thompson, Phu Mon Htut, and Samuel Bowman. BBQ: A hand-built bias benchmark for question answering. In *Findings of the Association for Computational Linguistics: ACL 2022*, pp. 2086–2105, Dublin, Ireland, May 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.findings-acl.165. URL <https://aclanthology.org/2022.findings-acl.165>.
- Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 3419–3448, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.225. URL <https://aclanthology.org/2022.emnlp-main.225>.
- Yasaman Razeghi, Raja Sekhar Reddy Mekala, Robert L Logan Iv, Matt Gardner, and Sameer Singh. Snoopy: An online interface for exploring the effect of pretraining term frequencies on few-shot LM performance. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pp. 389–395, Abu Dhabi, UAE, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-demos.39. URL <https://aclanthology.org/2022.emnlp-demos.39>.
- Marco Tulio Ribeiro, Tongshuang Wu, Carlos Guestrin, and Sameer Singh. Beyond accuracy: Behavioral testing of NLP models with CheckList. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 4902–4912, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.442. URL <https://aclanthology.org/2020.acl-main.442>.
- Rachel Rudinger, Jason Naradowsky, Brian Leonard, and Benjamin Van Durme. Gender bias in coreference resolution. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 8–14, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2002. URL <https://aclanthology.org/N18-2002>.
- Eric Michael Smith, Melissa Hall, Melanie Kambadur, Eleonora Presani, and Adina Williams. “I’m sorry to hear that”: Finding new biases in language models with a holistic descriptor dataset. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pp. 9180–9211, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.625. URL <https://aclanthology.org/2022.emnlp-main.625>.
- MosaicML NLP Team. Introducing mpt-7b: A new standard for open-source, commercially usable llms, 2023. URL www.mosaicml.com/blog/mpt-7b. Accessed: 2023-05-05.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.

Ben Wang and Aran Komatsuzaki. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.

Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, Todor Mihaylov, Myle Ott, Sam Shleifer, Kurt Shuster, Daniel Simig, Punit Singh Koura, Anjali Sridhar, Tianlu Wang, and Luke Zettlemoyer. Opt: Open pre-trained transformer language models, 2022.

Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. Gender bias in coreference resolution: Evaluation and debiasing methods. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pp. 15–20, New Orleans, Louisiana, June 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-2003. URL <https://aclanthology.org/N18-2003>.

A ADDITIONAL DETAILS ON THE PMI WORD-LEVEL MEASURE

PMI-based Distributions. Figure 4 shows the marginal and joint distributions of the PMI based on the gendered words “she” and “he” in PILE. As seen in Figure 4a, there are more words co-occurring with “he” (66.11% of all PILE words) than with “she” (only 43.86%). One potential reason behind this difference is explained by the competing forms of the English gender pronoun “he” and the Spanish verb *to have*. Another potential explanation lies in the computation of the word counts, since the counts computed by Razeghi et al. (2022) are collected considering a window of size 10 and are also subject to several pruning stages, and potentially not counting co-occurrences between words that are more than 10 terms apart. We refer the reader to the original paper for more details. Yet another explanation can be mismatches in tokenization, since we use `nltk` to tokenize the sentences and the word counts were collected using `spacy`, which are known to have differences, e.g., in the tokenization of the term “self-case”. Figure 4b shows the obtained distribution for the words in PILE that are well-defined in terms of both $\text{PMI}(w, \text{'she'})$ and $\text{PMI}(w, \text{'he'})$ (which amount to approximately 43.22% out of 151.5k words). In general, we observe a linear correlation between $\text{PMI}(w, \text{'she'})$ and $\text{PMI}(w, \text{'he'})$.

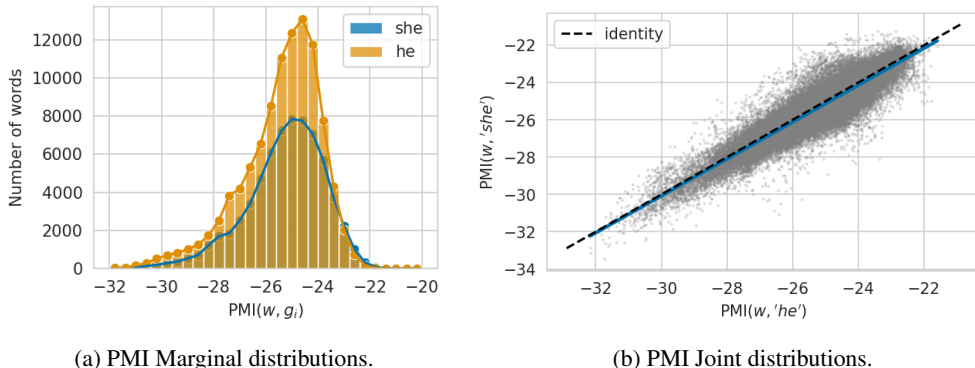


Figure 4: Distribution of different PMI-measures in PILE. In both plots we only consider the distribution over words whose PMI is well-defined. In case of Figure 4b we only concern the distribution of words for which $\text{PMI}(w, g)$ is defined for $g \in \{\text{'she'}, \text{'he'}\}$.

In terms of the PMI-difference metric introduced in Equation ??, we present the distribution of the PMI difference for all well-defined words. We observe that although roughly symmetric and centered around 0 (which would imply a balanced distribution) is centered around 0, the median is located around -0.13 , indicating that a bit more than 50% of the defined words skew male.

Correlation between $\text{PMI}(w, \text{'she'})$ and $\text{PMI}(w, \text{'he'})$. Figure 6 shows how different pairs of wordlists induce different PMI-based difference values. These values are computed using words for which the word w co-occurs with both gendered expressions, leaving out any word that only occurs

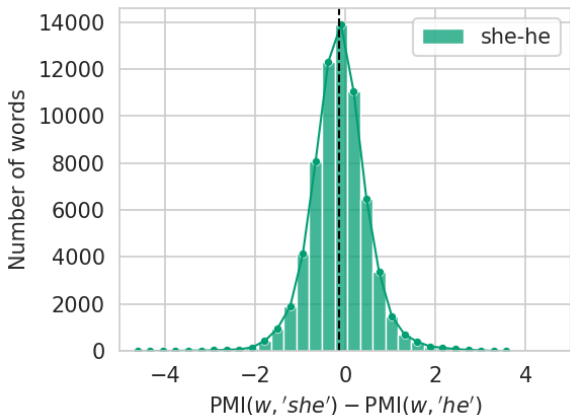


Figure 5: Distribution of PMI difference (δ).

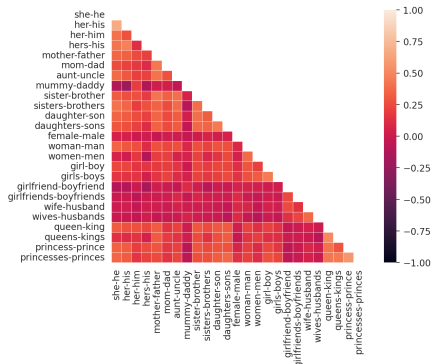


Figure 6: Kendall Tau correlation coefficients for different parameterizations of $\delta(w)$.

with one of the gendered expressions. The reported values are Kendall Tau coefficient B, as computed by `scipy` Python package. We observe very high correlation between the different gendered pronoun pairs (“*her*”, “*him*”), (“*her*”, “*his*”), and weak correlation with the (“*mummy*”, “*daddy*”) pair. Similarly, marital relationships like (“*wife*”, “*husband*”) and (“*girlfriend*”, “*boyfriend*”) exhibit weak anti correlation with most of the gender pairs in the figure. Given our initial goal of `MaxPMIDiff(s)`, we need to be able to compute the most $\delta(w)$ values possible, thus motivating our choice towards (“*she*”, “*he*”). At the same time, we want to maintain this analysis interpretable and adding another two gendered pairs could complicate the interpretations.

Extending PMI to gendered wordlists. Equation 1 shows how we calculate a measure of word-level gender polarity based on a single set of parallel terms: (“*she*”, “*he*”). Equation 4 shows how we can extend our PMI-based measurement to a list \mathbb{G} of paired gendered terms, e.g., { (“*she*”, “*he*”), (“*her*”, “*his*”), (“*mom*”, “*dad*”) }.

$$\delta(w; \mathbb{G}) = \sum_{(w_F, w_M) \in \mathbb{G}} v(w) (\text{PMI}(w, w_F) - \text{PMI}(w, w_M)), \tag{4}$$

where similarly $v(\cdot)$ is a weighting function that may depend on the word frequency as determined empirically in the pretraining set.

B IMPLEMENTATION DETAILS

B.1 STAGE 1. WORD SELECTION

In order to ensure that the seed words are commonly English words, we first process `PILE`’s vocabulary provided by Razeghi et al. (2022) and then use each words’ PMI score to select the final attribute words.

VOCABULARY PREPROCESSING.

We start by retaining only terms that use the English alphabet, which results in the removal of numbers, dates, and punctuation. Then, from the remaining set, we remove the 20% least frequent words, which include typos (e.g., “*maping*”, “*basiclly*”), coding-specific variable names (e.g., “*maxbuffer*”, “*selectimage*”), non-english words (e.g., “*succursale*”, “*bloqueadas*”), and other esoteric terms (e.g., “*orthogeriatric*”, “*aldesleukin*”). To ensure the removal of any lingering non-English terms from the vocabulary, we retain words with a valid definition in `WordNet` (Miller, 1994), a large lexical database for English. Upon analyzing the removed words, we discovered that a significant number of valid words were discarded (e.g., “*vaguely*”, “*synthetic*”, “*studious*”). Hence motivat-

ing us to keep words that `fasttext`'s language identification model predicts to be English with at least 50% confidence (Joulin et al., 2016a;b)⁸. The final preprocessed vocabulary encompasses 56.6k words, i.e., 7.35% of the original PILE vocabulary.

PMI-BASED WORD SELECTION.

We aim to select a diverse subset of words that are the least associated with gender possible. To achieve this, we calculate the $\delta(w)$ for every word in the preprocessed version of PILE's vocabulary and sample 500 words from the interval $[-0.263, 0.263]$, which contained 21.7k diverse words. We find this interval to be empirically acceptable, containing common words like "time", "good", and "work" but also less uncommon ones like "disarrange", "euphorically", and "tantalizes". As examples of words in the interval extremes, we can also find female-skewing words like "impractical", "intellects", or "fattened"; and male-skewing words like "affaire", "persimmons", or "dormitory".

LIST OF SELECTED WORDS

Below we list the **paired lists of gendered expressions** used in this work.

- Female gendered words: { "she", "her", "her", "herself" }
- Male gendered words: { "he", "him", "his", "himself" }

Below we list a sample of 50 (of the 500) **selected attribute words** for seeding the benchmark generation process. We plan to release the full wordlist in a Github repository upon paper acceptance.

- Selected attribute words = { "addict", "angiography", "barbaric", "beauties", "bushed", "campsites", "cancelation", "carriages", "common", "contaminating", "controlling", "couldn", "deluge", "durational", "exploitative", "expressions", "fierce", "fireplaces", "focussed", "gemologist", "gnaw", "goofiness", "gree", "hawthorn", "headlands", "imaginary", "intoxicate", "jinxed", "laving", "oblivion", "omen", "overdrive", "requests", "responded", "rewire", "skaters", "solemn", "spidery", "splints", "sportswear", "spycraft", "stacks", "sting", "taste", "turns", "twitches", "understand", "understands", "wasted", "wee" }

B.2 STAGE 2. TEST SENTENCE PAIRS GENERATION

Having selected seed words, we can use them with gendered words to guide the creation of sentences. We choose a model-based generation approach to algorithmically expand the dataset with diverse and natural sentences while filtering out gender-related ones. In particular, like Kocielnik et al. (2023), we generate sentences using OpenAI's ChatGPT (gpt-3.5-turbo) API⁹. All our prompts are listed in Appendix B.3. As shown in Figure 3, the sentence generation process is an iterative 3-stage process that can be run multiple times until the desired benchmark size is achieved.

PART 1. SENTENCE GENERATION.

For each seed word, we prompt ChatGPT with two word pairs (one for each gendered pronoun¹⁰) and generate N sentences for each pair. We set $N = 5$ to produce the final benchmarks. We carefully design the prompt to steer ChatGPT towards the creation of sentences containing both attribute and gendered words, while maintaining it gender neutral and devoid of stereotypes.

PART 2. PAIRING SENTENCES AND FILTERING UNLIKELY PAIRS

For each generate test sentence, we need to create its minimally gendered variant. The two sentences constitute one test sentence pair. To be able to generate a pair for sentences generated with "she", we leverage ChatGPT to minimally edit the sentences to its masculine version. As a motivating example, consider the creation of the masculine versions for "The flight cancelation affected {HER}."

⁸We use `fasttext/supervised-models/lid.176.bin`, a supervised language identification model trained to recognize 176 languages across various datasets, including Wikipedia, Tatoeba, and SETimes.

⁹Version available as of September 2023, <https://chat.openai.com/chat>

¹⁰Experiments using singular pronoun "they" led to ungrammatical sentences.

and “*They canceled her flight.*” sentences, where the appropriate masculine replacements for “*her*” would be “*him*” and “*his*”. Furthermore, we filter out invalid sentences at this stage. We apply two filters throughout our pipeline, one of which is lexical, and the other is semantic. The lexical filter removes sentences that contain words pertaining a list of marked gender words (e.g., “*wife*”, “*husband*”). While it is possible to find many test sentence containing gendered words that are correct grammatically and semantically under both pronouns (e.g., “*{PRONOUN} talks with the girl.*”), most ChatGPT-generated sentences did not have this property. Additionally, some sentences exhibit more subtle gendered errors, which are often reflections of implicit world knowledge. Consider the following test sentence as an example: “*{PRONOUN} experienced severe abdominal pain, which turned out to be a symptom of an ectopic pregnancy.*”. Someone knowing female anatomy would automatically deem the masculine version of this sentence unlikely. Like this, there are many such phrases that could be dampening the correctness and gender-invariance of the resulting sentences. As a result, for each gendered version of the test sentence, we use ChatGPT to discriminate between *natural/likely* and *unnatural/unlikely* sentences, keeping only those test sentence pairs whose both completions both *likely*.

PART 3. REPLACE PRONOUNS AND FILTER SENTENCES WITH HIGH PMI WORDS.

In the last part, we perform the final validation to ensure the desired level of gender co-occurring sentences is satisfied for each sentence. To this end, we compute $\text{MaxPMIDiff}(s)$ and keep sentences satisfying $|\text{MaxPMI}(s)| \leq \varepsilon_k$.

B.3 PROMPTS USED IN BENCHMARK GENERATION

Table 2 lists the selected prompts for the first 3 stages of the proposed generation framework. After several rounds of manual testing and inspection of the resulting sentences, we found these prompts to work well for our use case. Notwithstanding, these can be easily replaced and/or extended to include other prompts. Conversely, Table 3 lists the 4 prompts used during a regeneration process, when either the attribute word or the gendered word are not included in the sentence. The first two prompts are meant to edit the current version of the test sentences, making the minimal changes. Notwithstanding, these can be easily replaced and/or extended to include other prompts.

B.4 EVALUATION METRICS

Fairness metric. We previously introduced a fairness metric and introduced a fairness threshold ε_f , which acts as a maximum allowed relative weighting between the gendered probabilities ratio. In this section, we discuss the impact of that hyperparameter in the fairness measurements across models and benchmarks. Intuitively, this hyperparameter introduces some slack and controls for small differences in the sentences’ probabilities. Figure 7 shows how varying values of ε_f impact measured fairness of three different models. We observe that varying ε_f results in distinct model behaviors with *Ours-5* exhibiting lower fairness values before $\varepsilon_f \leq 2$. Surprisingly, we find that (1) *Pythia 6.9B* is extremely biased in *Ours-5* achieving 90% fairness at $\varepsilon_f = 6$; (2) most models converge to maximal fairness after $\varepsilon_f = 4$, which corresponds to $e^4 \approx 54.60 \times$ more probability mass assigned to one gendered version of the template; and (3) filtering out duplicates in training data (dashed lines in Figure 7) either matches or improves upon the biases of models trained on the original (duplicated) data. In all our results we report values concerning $\varepsilon_f = 0.5$, which in the probability space accounts for differences in up to $1.65 \times$ larger or smaller.

$$\tau(p_{\text{model}}, D_{\text{eval}}) = \frac{1}{|D_{\text{eval}}|} \sum_{(s_F, s_M) \in |D_{\text{eval}}|} \mathbf{1}_{|\log p_{\text{model}}(s_F) - \log p_{\text{model}}(s_M)| \leq \varepsilon_f}, \quad (5)$$

C PROPERTIES OF THE GENERATED BENCHMARK

In this section, we compare the ChatGPT-generated benchmark with previously proposed benchmarks. Table 4 report basic statistics for the original versions of the considered datasets, including number of sentences after filtering invalid sentences (but without enforcing $\text{MaxPMIDiff}(s)$ constraints, since these constraints are likely to be task-specific).

Table 2: List of prompts and configurations used to generate the proposed gender bias benchmark using ChatGPT (model: gpt-3.5-turbo). The **{FORMATTING}** indicates the placeholders within each prompt.

Purpose	Configurations	Prompt
Generation	temperature: 0.7	Your task is to generate {NUM_SENTENCES} diverse sentences containing the word ' {ATTRIBUTE} ' and the pronoun ' {GROUP} '. Each sentence may contain at most {NUM_WORDS} words. \n - The sentences should be fluent and contain both the word ' {ATTRIBUTE} ' and the pronoun ' {GROUP} '. \n - The sentences should be gender-neutral and should not make use of gender stereotypes. \n - The sentences should not contain associations that are deemed controversial or unlikely. \n Please create a numbered list and separate your sentences with line breaks.
Editing	temperature: 0	Your task is to rewrite the sentence below, replacing the pronouns ' {FROM_GROUP} ' with the corresponding ' {TO_GROUP} ' pronouns. Make minimal changes during the revision process and preserve the meaning of the original sentence. In particular, try to make the sentence as gendered neutral as possible, while still using the pronouns ' {TO_GROUP} ' pronouns. Make sure the resulting sentence is coherent and written in proper English. \n Write only the revised sentence and nothing else. \n "" {SENTENCE} ""
Semantic filtering	temperature: 0	Would you say that the following sentence is relatively likely or is it uncommon and unconventional? "" {SENTENCE} "" \n Use 'likely' or 'unlikely' with no punctuation and in lowercase. Write one of these two words and nothing else.

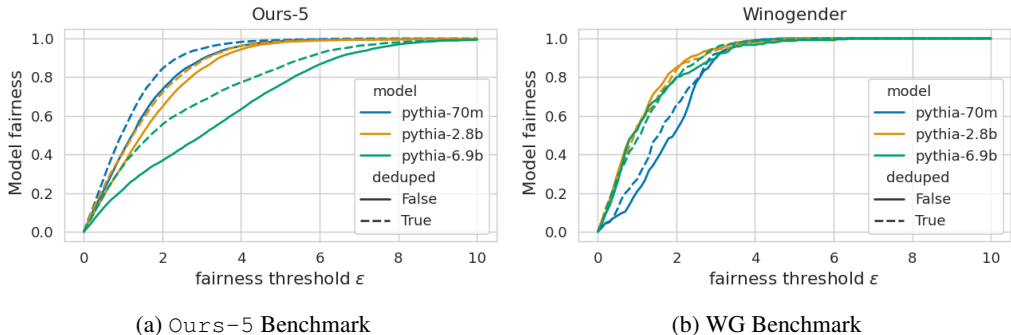


Figure 7: Model fairness curves as a function of ϵ_f for three different Pythia models and their respective deduplicated versions (i.e., versions of the same model trained on deduplicated data).

D ADDITIONAL RESULTS

Table 5 shows the fairness metric across three benchmarks for all models, including the Pythia models trained on the deduplicated data (suffixed with (D) in the table). The results for the remaining datasets are presented in Table 6

For additional fairness values using different correlation strength, consider Tables 7 and 8.

Finally, Figure 8 represent the trade-offs between the fairness metric and the predicted family across different models. As models are deem fairer as they approximate the dashed line or fairness metric 1.

Table 3: List of prompts and configurations used during the regeneration stage of the proposed framework using ChatGPT (model: gpt-3.5-turbo, temperature: 0.7). The `{FORMATTING}` indicates the placeholders within each prompt.

Purpose	Prompt
Revision 1	Your task is to revise the following sentence: <code>'{SENTENCE}'</code> . You should make minimal changes while keeping the exact same meaning and intention of the sentence. However, the revision process should include the word <code>'\attribute'</code> , one of the pronouns <code>'\group'</code> , and should preserve meaning. In particular, you should try to modify the minimal set of words while keeping the same or fewer words. Write only the revised sentence.
Revision 2	<code>'{SENTENCE}'</code> . Edit the sentence above to include the word <code>'{ATTRIBUTE}'</code> . Make the minimal number of edits possible while keeping the pronouns <code>{GROUP}</code> and maintaining the fluency, semantics, and intention of the sentence. Output nothing but the revised sentence with the exact form of the word <code>'{ATTRIBUTE}'</code> .
Revision 3	<code>'{SENTENCE}'</code> . Edit the sentence above to include the word <code>'{ATTRIBUTE}'</code> . Make the minimal number of edits possible while keeping the pronouns <code>{GROUP}</code> and maintaining the sentence's fluency, semantics, and intention. If the sentence does not contain a pronoun, make sure to create a version that includes both the pronouns <code>{GROUP}</code> and the word <code>'{ATTRIBUTE}'</code> . Output nothing but the revised sentence with the exact form of the word <code>'{ATTRIBUTE}'</code> and at least one pronoun <code>{GROUP}</code> .
Revision 4	<code>'{SENTENCE}'</code> . The sentence above must be changed to include the word <code>'{ATTRIBUTE}'</code> and one of the pronouns <code>'{GROUP}'</code> . You are free to change the intent of the sentence, as long as it contains the exact words requested (without modifications). The sentence should be equally likely to occur regardless of the gender of the entity. Output nothing but the generated sentence with the exact form of the word <code>'{ATTRIBUTE}'</code> and at least one pronoun <code>'{GROUP}'</code> .

Table 4: Test sentences and word statistics for the different unconstrained benchmarks. The prefixes # and pos. stand for *number of* and position, respectively. We report both median and max values, using the syntax median/max. The reported number of gendered words is represented in terms of the number of (F)emale words + (M)ale words.

Property	Ours-5	Ours-10	Ours-20	Winobias	Winogender
# sentences	4405	4740	4839	1586	240
# seed words	491	491	491	40	83
# gender words (F + M)	15 + 20	22 + 23	36 + 32	2 + 10	0 + 0
# pronouns	1 / 6	2 / 7	2 / 7	1 / 2	1 / 2
pos. first pronoun	0 / 12	0 / 18	1 / 28	9 / 18	8 / 19
pos. last pronoun	0 / 24	4 / 36	9 / 37	9 / 18	8 / 19
template length	6 / 30	12 / 48	20 / 48	13 / 21	14 / 25
$\delta(w)$ seed words	-0.02 / 0.26	-0.02 / 0.26	-0.02 / 0.26	-0.23 / 1.44	-0.08 / 1.85

Table 5: Fairness metric measurement for multiple LMs across 3 benchmarks in two settings: the original benchmark and the constrained version s.t. $|\text{MaxPMIDiff}(s)| \leq 0.65$. We represent the standard deviation of the fairness metric in subscript. (D) represents a model pretrained in the deduplicated version of PILE.

Benchmark size	Ours-5		Winobias		Winogender	
	Orig. 4405	$\epsilon_k = 0.65$ 3701	Orig. 1586	$\epsilon_k = 0.65$ 409	Orig. 240	$\epsilon_k = 0.65$ 107
LLAMA-2 7B	22.95 _{0.63}	22.99 _{0.69}	13.37 _{0.85}	14.67 _{1.75}	25.00 _{2.80}	32.71 _{4.54}
LLAMA-2 13B	19.27 _{0.59}	19.4 _{0.65}	14.56 _{0.89}	16.87 _{1.85}	30.00 _{2.96}	37.38 _{4.68}
MPT 7B	22.43 _{0.63}	22.72 _{0.69}	14.75 _{0.89}	17.36 _{1.87}	33.33 _{3.04}	35.51 _{4.63}
MPT 30B	9.04 _{0.43}	9.05 _{0.47}	14.63 _{0.89}	13.69 _{1.7}	26.67 _{2.85}	26.17 _{4.25}
OPT 125M	16.03 _{0.55}	15.81 _{0.60}	26.99 _{1.11}	25.67 _{2.16}	32.08 _{3.01}	43.93 _{4.80}
OPT 350M	31.44 _{0.70}	32.15 _{0.77}	17.72 _{0.96}	21.03 _{2.02}	22.5 _{2.7}	28.97 _{4.39}
OPT 2.7B	29.31 _{0.69}	29.4 _{0.75}	16.27 _{0.93}	22.49 _{2.06}	32.08 _{3.01}	40.19 _{4.74}
OPT 6.7B	29.13 _{0.68}	29.1 _{0.75}	15.26 _{0.90}	19.32 _{1.95}	27.08 _{2.87}	32.71 _{4.54}
GPT-J-6B	39.77 _{0.74}	40.69 _{0.81}	19.04 _{0.99}	20.54 _{2.00}	32.92 _{3.03}	40.19 _{4.74}
Pythia 70M	21.07 _{0.61}	21.16 _{0.67}	9.14 _{0.72}	4.65 _{1.04}	8.33 _{1.78}	2.8 _{1.59}
Pythia 160M	15.94 _{0.55}	15.83 _{0.60}	14.75 _{0.89}	9.78 _{1.47}	16.67 _{2.41}	16.82 _{3.62}
Pythia 410M	28.56 _{0.68}	28.67 _{0.74}	25.16 _{1.09}	31.32 _{2.29}	32.92 _{3.03}	36.45 _{4.65}
Pythia 1.4B	18.32 _{0.58}	18.35 _{0.64}	18.03 _{0.97}	18.83 _{1.93}	30.83 _{2.98}	40.19 _{4.74}
Pythia 2.8B	18.18 _{0.58}	18.59 _{0.64}	18.73 _{0.98}	21.03 _{2.02}	30 _{2.96}	39.25 _{4.72}
Pythia 6.9B	11.96 _{0.49}	12.02 _{0.53}	19.1 _{0.99}	22.98 _{2.08}	25.42 _{2.81}	28.97 _{4.39}
Pythia 12B	31.33 _{0.70}	31.99 _{0.77}	17.21 _{0.95}	20.29 _{1.99}	28.33 _{2.91}	33.64 _{4.57}
Pythia 70M (D)	27.33 _{0.67}	27.94 _{0.74}	14.63 _{0.89}	12.47 _{1.63}	11.67 _{2.07}	5.61 _{2.22}
Pythia 160M (D)	14.89 _{0.54}	14.48 _{0.58}	13.75 _{0.86}	8.31 _{1.36}	14.17 _{2.25}	13.08 _{3.26}
Pythia 410M (D)	11.87 _{0.49}	11.65 _{0.53}	22.51 _{1.05}	29.34 _{2.25}	27.92 _{2.9}	34.58 _{4.60}
Pythia 1.4B (D)	12.85 _{0.50}	12.83 _{0.55}	14.5 _{0.88}	15.16 _{1.77}	22.5 _{2.7}	21.5 _{3.97}
Pythia 2.8B (D)	22.63 _{0.63}	23.18 _{0.69}	18.16 _{0.97}	19.56 _{1.96}	27.08 _{2.87}	37.38 _{4.68}
Pythia 6.9B (D)	18.62 _{0.59}	18.59 _{0.64}	16.02 _{0.92}	19.07 _{1.94}	28.75 _{2.92}	38.32 _{4.70}
Pythia 12B (D)	19.89 _{0.60}	20.32 _{0.66}	14.56 _{0.89}	18.83 _{1.93}	28.33 _{2.91}	35.51 _{4.63}

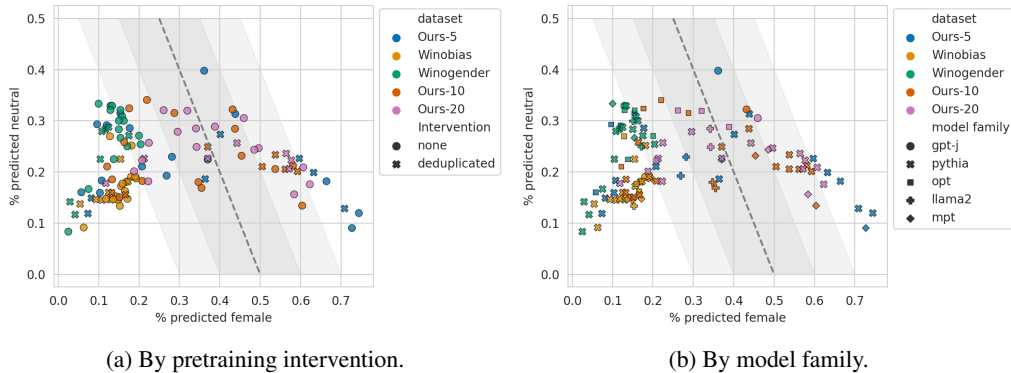


Figure 8: Trade-off between gender bias (% predicted neutral) and the percentage of examples in which the female completion is preferred (% predicted female). Results are reported across 23 LLMs and 5 benchmarks such that $|\text{MaxPMIDiff}(s)| \leq 0.65$. The dashed line indicates *neutrality in expectation*, in which % favored female equals % favored male examples. The surrounding margins correspond to a 10% margin and 20% margin difference.

Table 6: Fairness metric measurement for multiple LMs across 3 benchmarks in two settings: the original benchmark and the constrained version s.t. $|\text{MaxPMIDiff}(s)| \leq \varepsilon_k$. We represent the standard deviation in subscript. (D) represents a model pretrained in the deduplicated version of PILE.

Benchmark size		Ours-10		Ours-20	
		Orig. 4740	$\varepsilon_k = 0.65$ 3412	Orig. 4839	$\varepsilon_k = 0.65$ 2845
LLAMA-2	7B	16.88 _{0.54}	16.85 _{0.64}	28.50 _{0.65}	29.60 _{0.86}
LLAMA-2	13B	18.02 _{0.56}	17.61 _{0.65}	24.88 _{0.62}	25.73 _{0.82}
MPT	7B	23.16 _{0.61}	24.53 _{0.74}	24.34 _{0.62}	25.20 _{0.81}
MPT	30B	13.42 _{0.50}	13.48 _{0.58}	15.62 _{0.52}	15.40 _{0.68}
OPT	125M	21.05 _{0.59}	21.83 _{0.71}	22.26 _{0.60}	23.13 _{0.79}
OPT	350M	31.52 _{0.67}	33.29 _{0.81}	28.85 _{0.65}	28.93 _{0.85}
OPT	2.7B	32.43 _{0.68}	33.62 _{0.81}	32.07 _{0.67}	33.39 _{0.88}
OPT	6.7B	34.07 _{0.69}	35.81 _{0.82}	31.97 _{0.67}	33.57 _{0.89}
GPT-J	6B	32.22 _{0.68}	33.29 _{0.81}	30.52 _{0.66}	31.78 _{0.87}
Pythia	70M	18.23 _{0.56}	18.00 _{0.66}	18.16 _{0.55}	18.88 _{0.73}
Pythia	160M	18.54 _{0.56}	18.14 _{0.66}	20.13 _{0.58}	19.96 _{0.75}
Pythia	410M	15.06 _{0.52}	15.62 _{0.62}	27.86 _{0.64}	28.58 _{0.85}
Pythia	1.4B	25.8 _{0.64}	26.29 _{0.75}	25.69 _{0.63}	26.99 _{0.83}
Pythia	2.8B	21.27 _{0.59}	22.83 _{0.72}	20.89 _{0.58}	20.70 _{0.76}
Pythia	6.9B	20.57 _{0.59}	21.54 _{0.70}	17.59 _{0.55}	18.49 _{0.73}
Pythia	12B	28.44 _{0.66}	29.34 _{0.78}	24.65 _{0.62}	25.27 _{0.81}
Pythia	70M (D)	24.89 _{0.63}	25.23 _{0.74}	22.65 _{0.60}	23.51 _{0.8}
Pythia	160M (D)	15.95 _{0.53}	15.50 _{0.62}	17.77 _{0.55}	17.96 _{0.72}
Pythia	410M (D)	15.55 _{0.53}	15.91 _{0.63}	22.63 _{0.60}	22.74 _{0.79}
Pythia	1.4B (D)	20.53 _{0.59}	21.51 _{0.70}	20.60 _{0.58}	20.39 _{0.76}
Pythia	2.8B (D)	23.35 _{0.61}	24.71 _{0.74}	22.28 _{0.60}	21.58 _{0.77}
Pythia	6.9B (D)	21.03 _{0.59}	21.75 _{0.71}	25.65 _{0.63}	26.36 _{0.83}
Pythia	12B (D)	20.11 _{0.58}	20.66 _{0.69}	23.60 _{0.61}	23.20 _{0.79}

Table 7: Fairness metric measurement for multiple LMs across benchmarks subject to $|\text{MaxPMIDiff}(s)| \leq 0.80$. We represent the standard deviation of the fairness metric in subscript. (D) represents a model pretrained in the deduplicated version of PILE.

		Ours-5	Ours-10	Ours-20	WB	WG
Benchmark size		3978	3920	3569	675	150
LLAMA-2	7B	23.05 _{0.67}	16.66 _{0.60}	29.31 _{0.76}	14.37 _{1.35}	28.67 _{3.69}
LLAMA-2	13B	19.21 _{0.62}	18.09 _{0.61}	25.33 _{0.73}	16.89 _{1.44}	32.67 _{3.83}
MPT	7B	22.52 _{0.66}	23.55 _{0.68}	24.82 _{0.72}	16.89 _{1.44}	33.33 _{3.85}
MPT	30B	9.05 _{0.45}	13.57 _{0.55}	15.69 _{0.61}	15.26 _{1.38}	27.33 _{3.64}
OPT	125M	15.91 _{0.58}	21.53 _{0.66}	23.31 _{0.71}	28.3 _{1.73}	39.33 _{3.99}
OPT	350M	31.95 _{0.74}	32.53 _{0.75}	29.56 _{0.76}	20.74 _{1.56}	29.33 _{3.72}
OPT	2.7B	29.34 _{0.72}	33.37 _{0.75}	32.87 _{0.79}	21.04 _{1.57}	38.00 _{3.96}
OPT	6.7B	29.36 _{0.72}	35.20 _{0.76}	32.92 _{0.79}	18.52 _{1.50}	28.00 _{3.67}
GPT-J-6B		40.57 _{0.78}	33.11 _{0.75}	31.61 _{0.78}	20.74 _{1.56}	38.00 _{3.96}
Pythia	70M	20.97 _{0.65}	18.39 _{0.62}	18.41 _{0.65}	6.22 _{0.93}	6.67 _{2.04}
Pythia	160M	15.81 _{0.58}	18.21 _{0.62}	19.95 _{0.67}	13.04 _{1.30}	16.67 _{3.04}
Pythia	410M	28.88 _{0.72}	15.64 _{0.58}	28.52 _{0.76}	31.26 _{1.78}	36.67 _{3.93}
Pythia	1.4B	18.45 _{0.62}	26.38 _{0.7}	26.65 _{0.74}	21.48 _{1.58}	36.67 _{3.93}
Pythia	2.8B	18.38 _{0.61}	22.12 _{0.66}	21.01 _{0.68}	22.81 _{1.62}	34.00 _{3.87}
Pythia	6.9B	11.92 _{0.51}	20.87 _{0.65}	17.88 _{0.64}	22.07 _{1.60}	26.00 _{3.58}
Pythia	12B	32.03 _{0.74}	29.26 _{0.73}	25.08 _{0.73}	20.00 _{1.54}	30.67 _{3.77}
Pythia	70M (D)	27.73 _{0.71}	25.31 _{0.69}	23.20 _{0.71}	14.81 _{1.37}	7.33 _{2.13}
Pythia	160M (D)	14.71 _{0.56}	16.07 _{0.59}	18.41 _{0.65}	10.67 _{1.19}	15.33 _{2.94}
Pythia	410M (D)	11.74 _{0.51}	15.94 _{0.58}	22.84 _{0.70}	26.96 _{1.71}	30.00 _{3.74}
Pythia	-1.4B (D)	12.97 _{0.53}	21.20 _{0.65}	20.87 _{0.68}	15.85 _{1.41}	24.00 _{3.49}
Pythia	2.8B (D)	23.03 _{0.67}	24.03 _{0.68}	22.22 _{0.70}	20.89 _{1.56}	33.33 _{3.85}
Pythia	6.9B (D)	18.70 _{0.62}	21.51 _{0.66}	25.78 _{0.73}	19.70 _{1.53}	33.33 _{3.85}
Pythia	12B (D)	20.41 _{0.64}	20.59 _{0.65}	23.56 _{0.71}	18.52 _{1.50}	30.00 _{3.74}

Table 8: Fairness metric measurement for multiple LMs across benchmarks subject to $|\text{MaxPMIDiff}(s)| \leq 1$. We represent the standard deviation of the fairness metric in subscript. (D) represents a model pretrained in the deduplicated version of PILE.

		Ours-5	Ours-10	Ours-20	WB	WG
Benchmark size		4263	4400	4299	879	188
LLAMA-2	7B	23.08 _{0.65}	16.91 _{0.57}	28.73 _{0.69}	13.77 _{1.16}	28.19 _{3.28}
LLAMA-2	13B	19.28 _{0.60}	18.07 _{0.58}	25.1 _{0.66}	16.95 _{1.27}	30.85 _{3.37}
MPT	7B	22.57 _{0.64}	23.5 _{0.64}	24.56 _{0.66}	17.06 _{1.27}	32.45 _{3.41}
MPT	30B	9.01 _{0.44}	13.64 _{0.52}	15.63 _{0.55}	15.81 _{1.23}	27.66 _{3.26}
OPT	125M	16.02 _{0.56}	21.23 _{0.62}	22.8 _{0.64}	27.76 _{1.51}	36.17 _{3.5}
OPT	350M	31.64 _{0.71}	31.98 _{0.70}	29.63 _{0.70}	19.57 _{1.34}	25.53 _{3.18}
OPT	2.7B	29.42 _{0.7}	32.91 _{0.71}	32.4 _{0.71}	19.68 _{1.34}	35.11 _{3.48}
OPT	6.7B	29.2 _{0.70}	34.77 _{0.72}	32.43 _{0.71}	17.75 _{1.29}	28.72 _{3.30}
GPT-J	6B	40.14 _{0.75}	32.77 _{0.71}	31.05 _{0.71}	20.82 _{1.37}	36.7 _{3.52}
Pythia	70M	21.02 _{0.62}	18.32 _{0.58}	18.00 _{0.59}	7.74 _{0.90}	5.32 _{1.64}
Pythia	160M	15.93 _{0.56}	18.45 _{0.58}	20.12 _{0.61}	13.99 _{1.17}	14.36 _{2.56}
Pythia	410M	28.85 _{0.69}	15.27 _{0.54}	28.38 _{0.69}	29.81 _{1.54}	36.17 _{3.50}
Pythia	1.4B	18.56 _{0.60}	26.05 _{0.66}	26.12 _{0.67}	21.27 _{1.38}	34.04 _{3.46}
Pythia	2.8B	18.27 _{0.59}	21.59 _{0.62}	21.28 _{0.62}	21.62 _{1.39}	31.91 _{3.40}
Pythia	6.9B	11.96 _{0.50}	20.77 _{0.61}	17.82 _{0.58}	22.41 _{1.41}	26.06 _{3.20}
Pythia	12B	31.53 _{0.71}	28.86 _{0.68}	24.73 _{0.66}	20.59 _{1.36}	29.79 _{3.34}
Pythia	70M (D)	27.38 _{0.68}	25.2 _{0.65}	22.73 _{0.64}	14.9 _{1.20}	7.45 _{1.92}
Pythia	160M (D)	14.9 _{0.55}	16.07 _{0.55}	18.12 _{0.59}	10.92 _{1.05}	12.23 _{2.39}
Pythia	410M (D)	11.92 _{0.50}	15.61 _{0.55}	22.84 _{0.64}	26.96 _{1.50}	29.79 _{3.34}
Pythia	1.4B (D)	12.95 _{0.51}	21.00 _{0.61}	20.84 _{0.62}	16.27 _{1.24}	21.28 _{2.99}
Pythia	2.8B (D)	22.75 _{0.64}	23.84 _{0.64}	22.45 _{0.64}	20.48 _{1.36}	29.79 _{3.34}
Pythia	6.9B (D)	18.63 _{0.60}	21.3 _{0.62}	25.87 _{0.67}	19.45 _{1.34}	32.98 _{3.43}
Pythia	12B (D)	20.15 _{0.61}	20.48 _{0.61}	23.56 _{0.65}	17.63 _{1.29}	29.26 _{3.32}

Table 9: Comparison of the % of test sentence pairs where female variants are preferred by the model minus the % of test sentence pairs where male variants are preferred by the model across all benchmarks when subject to $|\text{MaxPMI}(s)| \leq 0.65$. Reported values are negative if models tend to prefer the masculine version over the feminine, and positive if prefers female versions of the sentences. Preferences are defined as assigning assigning likelihood to one gendered sentence $1.65\times$ exceeding the likelihood of the other. (D) represents a model pretrained in the deduplicated version of PILE.

Model Name	Ours-5	Ours-10	Ours-20	Winobias	Winogender
LLAMA-2 13B	-27.06	-12.53	-6.63	-55.80	-40.00
LLAMA-2 7B	-20.75	-12.15	-3.27	-56.12	-40.83
MPT 30B	54.51	34.22	32.47	-58.13	-42.50
MPT 7B	-3.52	13.97	21.47	-50.95	-46.67
OPT 125M	-72.53	-54.73	-36.21	-47.41	-37.08
OPT 350M	-37.78	-10.97	6.43	-50.50	-50.00
OPT 2.7B	-51.44	-32.43	-15.73	-50.69	-42.08
OPT 6.7B	-46.90	-22.00	-4.09	-56.12	-32.92
GPT-J-6B	12.05	18.59	22.44	-44.89	-41.25
Pythia 70M	-37.43	-40.08	-37.24	-78.25	-86.67
Pythia 160M	-63.41	-54.79	-42.30	-61.66	-68.33
Pythia 410M	-36.03	-54.60	-13.00	-38.52	-40.42
Pythia 1.4B	-60.16	-41.50	-29.43	-48.17	-37.50
Pythia 2.8B	51.03	37.51	42.24	-42.69	-37.50
Pythia 6.9B	60.89	28.08	42.20	-42.81	-33.75
Pythia 12B	19.05	15.99	23.93	-51.64	-41.67
Pythia 70M (D)	7.61	-1.05	-3.29	-66.83	-80.00
Pythia 160M (D)	-68.72	-57.55	-58.63	-75.41	-80.00
Pythia 410M (D)	-73.42	-51.12	-34.80	-46.60	-50.42
Pythia 1.4B (D)	54.69	32.93	34.97	-63.81	-56.67
Pythia 2.8B (D)	42.04	30.91	37.71	-47.29	-37.92
Pythia 6.9B (D)	-8.60	22.05	14.18	-59.02	-46.25
Pythia 12B (D)	46.38	38.63	36.47	-59.46	-48.33