

# Measuring Social Biases in Masked Language Models by Proxy of Prediction Quality

**Rahul Zalkikar\***  
Independent Researcher  
rayzck9@gmail.com

**Kanchan Chandra**  
New York University  
kc76@nyu.edu

## Abstract

Innovative transformer-based language models produce contextually-aware token embeddings and have achieved state-of-the-art performance for a variety of natural language tasks, but have been shown to encode unwanted biases for downstream applications. In this paper, we evaluate the social biases encoded by transformers trained with the masked language modeling objective using proposed proxy functions within an iterative masking experiment to measure the quality of transformer models' predictions and assess the preference of MLMs towards disadvantaged and advantaged groups. We find that all models encode concerning social biases. We compare bias estimations with those produced by other evaluation methods using benchmark datasets and assess their alignment with human annotated biases. We extend previous work by evaluating social biases introduced after retraining an MLM under the masked language modeling objective and find proposed measures produce more accurate and sensitive estimations of biases based on relative preference for biased sentences between models, while other methods tend to underestimate biases after retraining on sentences biased towards disadvantaged groups.

*Warning: This paper contains explicit statements of biased stereotypes and may be upsetting.*

## 1 Introduction

Word embeddings have proven useful in a variety of Natural Language Processing (NLP) tasks due to their ability to efficiently model complex semantic and syntactic word relationships. Token-level embeddings, such as those produced by Word2Vec (Mikolov et al., 2013) and GloVe (Pennington et al., 2014) algorithms, learn a non-contextualized text representation and produce static word embeddings that can uncover linear semantic or syntactic relationships between tokens.

Masked language models (MLM(s); Brown et al., 2020; Devlin et al., 2019; Liu et al., 2019a; Radford et al., 2019) such as transformers BERT and RoBERTa incorporate bidirectionality and self-attention, producing contextually-aware embeddings. Unlike BERT<sub>unc</sub>, RoBERTa was pretrained solely on the masked language modeling objective (MLMO) on a larger corpus of text. RoBERTa uses a dynamic masking strategy for a diverse set of representations during training and has achieved state-of-the-art results on GLUE, RACE, and SQuAD (Liu et al., 2019b; Rajpurkar et al., 2016; Lai et al., 2017; Wang et al., 2018). Distilled model variants have been shown to train significantly faster for minor decreases in performance (Sanh et al., 2019). MLMs have produced state-of-the-art results for masked language modeling, named entity recognition, and intent or topic classification, but also encode concerning social biases against disadvantaged groups that are undesirable in production settings. As MLMs become increasingly prevalent, researchers have been working on methods to measure biases embedded in these models (Nangia et al., 2020; Kaneko and Bollegala, 2022; Salutari et al., 2023).

To address the issue of social biases in MLMs at its source, we measure biases of MLMs while focusing on an MLM's key pretraining objective: masked language modeling. In this work, we focus on proposing and assessing bias evaluation measures for MLMs and *not* proposing methods for debiasing MLMs. However, when assessing evaluation measures, we consider prior research that involves retraining or fine-tuning under the MLMO with debiased or counterfactual data to reduce biases in MLMs, such as Zhao et al., 2018 and Zhao et al., 2019 which use data augmentation to swap gendered words with their opposites prior to retraining. Motivated by this, we assess whether proposed measures satisfy an important criterion for improvement over previously proposed ones: align-

\*Corresponding author.

ment with biases introduced by MLM retraining under the MLMO (Section 4.5.1).

We represent MLM bias through a model’s relative preference for ground truth tokens between two paired sentences along a social bias axis. In each sentence pair, one contains bias against disadvantaged groups (stereotypical) and the other contains bias against advantaged groups (anti-stereotypical). We assess the relative preferences of MLMs towards sentences biased against disadvantaged and advantaged groups using proxy measures for prediction quality. In particular, **we quantify MLM preference by proxy of masked token prediction quality given unmasked token context between encoded sentences within pairs.**

To measure the preference of an MLM using the (attention-weighted) quality of its predictions under the MLMO, we propose and validate proxy functions  $\Delta_{PA}$  (Equation 7; Probability Difference with Attention) and  $C_{RRA}$  (Equation 6; Complementary Reciprocal Rank with Attention) and introduce a modified  $\Delta_P$  (Equation 3) to measure the likelihood an MLM will select a ground truth token to replace a masked one, and extend these definitions for a sentence. We apply per-model indicator function  $BSPT$  (Equation 9) to estimate the encoded social biases in pretrained MLMs. Our approach differs from prior research in measuring biases in MLMs by using attention-weights under an Iterative Masking Experiment (IME; Section 4) to probe MLM preferences.

We compare pre- and retrained MLMs within the same model class to recover the nature of biases introduced by retraining. In particular, we define and apply a proxy for the relative preference between two MLMs with **model-comparative indicator function BSRT (Bias Score for MLM Re-training; Equation 10) to estimate social biases introduced by retraining MLMs under the MLMO.** These "introduced" biases must be represented by the relative change in biases after retraining an MLM under the MLMO.

In summary, the primary contributions of this work are as follows:

- We explore MLM bias through a model’s relative preference for ground truth tokens between two minimally distant sentences with contrasting social bias under the IME. We measure this using the attention-weighted quality of predictions.
- We propose the model-comparative indicator

function BSRT to estimate social biases against disadvantaged groups for a retrained MLM relative to its pretrained base, and assess bias evaluation measures for alignment with biases introduced by MLM retraining under the MLMO.

- We evaluate social biases for four transformer models available through the Transformers library (Wolf et al., 2020). We use proxy measures for MLM prediction quality with model-comparative function BSRT to estimate social biases introduced by retraining MLMs under the MLMO. We find that proposed measures  $\Delta_{PA}$  and  $C_{RRA}$ , along with modified  $\Delta_P$  and  $C_{RR}$ , produce more accurate estimations of biases introduced by MLM retraining than previously proposed ones, which can produce concerning underestimations of biases after retraining MLMs on sentences biased against disadvantaged groups. We observe that proposed measures  $\Delta_{PA}$  and  $C_{RRA}$ , along with modified  $\Delta_P$ , show greater sensitivity than  $C_{RR}$ ,  $C_{SPS}$ ,  $AUL$ , and  $AULA$  to relative changes in MLM bias due to retraining, indicated by larger and smaller bias scores and more frequently significant relative difference in proportions of bias between pre- and retrained transformers.

Our methodology could help others evaluate social biases encoded in an MLM after it is retrained on the MLMO, such as for any downstream fill-mask task, or after it is fine-tuned for other objectives that alter weights and change MLMO performance. We release a package for measuring biases in MLMs to enable bias score computation on user-supplied or benchmark datasets and easy integration into existing evaluation pipelines.

## 2 Related Work

### 2.1 Biases in Static Word Embeddings

Static word embeddings (Pennington et al., 2014; Mikolov et al., 2013) can be shifted in a direction to decompose bias embedded in learned text data representations. These could be analogies or biases along an axis, such as gender<sup>1</sup> (Bolukbasi et al., 2016) or race (Manzini et al., 2019). The Word Embedding Association Test (WEAT; Caliskan et al., 2017) measures association between targets and

<sup>1</sup>For example, doctor - man + woman = nurse.

attributes using cosine similarity between static word embeddings, but has been shown to overestimate biases by [Ethayarajh et al., 2019](#), which proposed the robust relational inner product association (RIPA) method, derived from the subspace projection method to debias vectors in [Bolukbasi et al., 2016](#).

While WEAT has shown that token-level embeddings produced by GloVe and Word2Vec encode biases based on gender and race ([Caliskan et al., 2017](#)), the Sentence Encoder Association Test (SEAT; [May et al., 2019](#)) extends on WEAT to measure social biases in sentence-level encoders such as ELMo and BERT using template sentences with masked target tokens<sup>2</sup>, averaging token embeddings to form sentence-level embeddings on which cosine similarity is applied as a proxy for semantic association. As an alternative evaluation method under a different objective, [Liang et al., 2020](#) assessed differences in log-likelihood between gender pronouns in a template sentence<sup>3</sup> where occupations can uncover the directionality of the bias encoded by an MLM.

## 2.2 Evaluating Biases in MLMs

When considering a sentence  $s$  containing tokens  $\{t_1, t_2, \dots, t_{l_s}\}$ , where  $l_s$  is the number of tokens in  $s$ , (modified) token(s) of  $s$  can characterize its bias towards either disadvantaged or advantaged groups. For a given sentence  $s$  with tokens  $t \in s$  we denote all tokens besides  $t_x$  as  $s \setminus t_x$  (where  $1 \leq x \leq l_s$ ), and we denote modified tokens as  $M$  and unmodified tokens as  $U$  ( $s = U \cup M$ ). For a given MLM with parameters  $\theta$ , we denote a masked token as  $t_m$  and a predicted token as  $t_p$ .

[Salazar et al., 2020](#) uses pseudo-log-likelihood scoring to approximate  $P(U|M, \theta)$ , or the probability of unmodified tokens conditioned on modified ones. Similarly, [Nangia et al., 2020](#) reports CrowS-Pairs Scores (CSPS; Appendix B), a pseudo-log-likelihood score for an MLM selecting unmodified tokens given modified ones. [Nadeem et al., 2021](#) reports a StereoSet Score (SSS; Appendix C), a pseudo-log-likelihood score for an MLM selecting modified tokens given unmodified ones.

[Salutari et al., 2023](#) tests MLMs in an iterative fill mask setting where the model outputs a set of tokens (or the (log)softmax of model logits mapped to tokens) to fill the masked one, starting with the

token of highest probability  $P(t_p|c)$  and, as such, first rank  $\rho(t_p|c) = 1$  in the set of model token predictions (which is limited by the MLM’s embedding space).

$$\begin{aligned} \Delta P(t|s \setminus t_m; \theta) &= P(t_p|s \setminus t_m; \theta) - P(t_m|s \setminus t_m; \theta) \\ &= \Delta P(w; \theta) \end{aligned} \quad (1)$$

$\Delta P(t|s \setminus t_m; \theta)$  (Equation 1) represents the difference between the probability of a predicted token  $t_p$  and a masked token  $t_m$  in a sentence  $s$ . It serves as a proxy of the MLM’s prediction quality for a token given its context within the IME, or all tokens in  $s$  besides  $t_m$  ([Salutari et al., 2023](#)).

$$\begin{aligned} \text{CRR}(t|s \setminus t_m; \theta) &= (\rho(t_p|s \setminus t_m; \theta))^{-1} - \rho(t_m|s \setminus t_m; \theta)^{-1} \\ &= 1 - \rho(t_m|s \setminus t_m; \theta)^{-1} \\ &= \text{CRR}(w; \theta) \end{aligned} \quad (2)$$

[Salutari et al., 2023](#) proposes the Complementary Reciprocal Rank ( $\text{CRR}(t|s \setminus t_m; \theta)$ ; Equation 2) for a masked token given its context, where  $\rho(t_p|s \setminus t_m)^{-1}$  is the reciprocal rank of the predicted token (and always equal to 1) and  $\rho(t_m|s \setminus t_m)^{-1}$  is the reciprocal rank of the masked token. Thus,  $\rho(t_m|s \setminus t_m)^{-1}$  provides a likelihood measure for  $t_m$  being chosen by the model as a candidate token to replace the ground truth (masked) one.

[Salutari et al., 2023](#) defines  $\Delta P(s)$  (Appendix D) as the probability difference for a sentence  $s$  and  $\text{CRR}(s)$  (Appendix E; average of all single masked token CRRs with token ordering preserved) as the complementary reciprocal rank for a sentence  $s$ , and uses them as proxies for an MLM’s prediction quality or preference, claiming metrics based on CRR for a sentence  $s$  are necessary to fully capture the biases embedded in MLMs.

[Kaneko and Bollegala, 2022](#) propose evaluation metrics All Unmasked Likelihood (AUL; Appendix F) and AUL with Attention weights (AULA; Appendix G), where both are generated by requiring the MLM to predict all tokens (unmasked input). By requiring the MLM to simultaneously predict all of the tokens in a given unmasked input sentence  $s$ , [Kaneko and Bollegala, 2022](#) aim to eliminate biases associated with masked tokens under previously proposed pseudo-likelihood-based scoring methods ([Nadeem et al., 2021](#), [Nangia et al., 2020](#)), which assumed that masked tokens are statistically independent, and selectional biases from masking a subset of input tokens, such as high frequency

<sup>2</sup>For example, "This is <mask>".

<sup>3</sup>For example, "<mask> is a/an [occupation]".

words (which are masked more often during training).

Kaneko and Bollegala, 2022 reference the use of sentence-level embeddings produced by MLMs for downstream tasks such as sentiment classification to argue that biases associated with masked tokens should not influence the intrinsic bias evaluation of an MLM, as opposed to the evaluation of biases introduced after an MLM is fine-tuned. In contrast, we focus on an MLM’s key pretraining objective, masked language modeling, to measure social biases of the MLM.<sup>4</sup> In addition, we measure relative changes in biases w.r.t. the intrinsic biases of the same base MLM after retraining under the MLMO (as opposed to fine-tuning). Thus, we argue that biases associated with masked tokens are not undesirable in our case.

AUL and AULA were found to be sensitive to contextually meaningful inputs by randomly shuffling tokens in input sentences and comparing accuracies with and without shuffle. CRR is also conditional to the unmasked token context by definition. We argue measure sensitivity to unmasked token contexts is desirable when evaluating a given MLM’s preference under a fill-mask task, or when estimating biases using contextualized token-level embeddings produced by an MLM.

In contrast with previous methods such as SSS and CSPS, we refrain from using strictly modified or unmodified subsets of input tokens as context and instead provide all tokens but the ground truth one as context for MLM prediction under each iteration of our masking experiment. In this sense, CRR and  $\Delta P$  consider all tokens equally, and like AUL, they might also benefit from considering the weight of MLM attention as a proxy for token importance when probing for MLM preferences for ground truth tokens between two paired sentences along a social bias axis.<sup>5</sup>

Existing benchmark datasets such as CPS are limited to one ground truth per masked token, so an important consideration is an MLM’s ability to predict multiple plausible tokens for a context that could qualify for concerning social bias but

<sup>4</sup>Masked language modeling was a pretraining objective for all transformers considered in this work. Next sentence prediction was not used for RoBERTa and variants due to relatively lower performance with its inclusion (Liu et al., 2019b).

<sup>5</sup>By incorporating attention weights,  $\Delta PA$  and CRR consider token-level contributions to MLM predictions, leveraging the attention mechanism’s role in weighting contextually important tokens.

goes unrecognized during evaluation using previously proposed measures. CRR could perform better than pseudo-(log)likelihood-based measures for this sensitivity that yields larger relative differences in  $\Delta P(t|c)$  as opposed to  $CRR(t|c)$ <sup>6</sup>, and is deemed critical for evaluation by Salutari et al., 2023 due to the possible uniformity of probabilities generated by a particular MLM w.r.t. others.

### 3 Methodology

#### 3.1 Biases in Pretrained MLMs

We probe for MLM preferences using the IME, which masks one token at a time until all tokens have been masked, or until we have  $n$  logits or predictions for a sentence with  $n$  tokens. Appendix I shows an IME example for one model and text.

Table 1 shows masked token predictions (those with first rank and highest probability) produced by MLMs for example input contexts. The ground truth masked token rank, or  $\rho(t_m|c)$ , represents the quality of the models’ predictions relative to a sentence biased against an advantaged or disadvantaged group. Special start and end character tokens for MLMs are not included in the span of tokens considered in our experiments to eliminate noise.<sup>7</sup>

$$\begin{aligned} \Delta P(t|s_{\setminus t_m}; \theta) &= (\log P(t_p|s_{\setminus t_m}; \theta) - \log P(t_m|s_{\setminus t_m}; \theta)) \\ &= \Delta P(w; \theta) \end{aligned} \quad (3)$$

Motivated by the typically long-tail distribution of probabilities output by MLMs, we redefine Equation 1 and propose a modified version  $\Delta P(w)$  as shown in Equation 3. The log transformation reduces sensitivity to differences when probabilities are high but enhances relative differences when probabilities are low, improving sensitivity to subtle but potentially meaningful variations in low-probability token choices.

$$\begin{aligned} CRR_A(t|s_{\setminus t_m}; \theta) &= a_m (1 - \log \rho(t_m|s_{\setminus t_m}; \theta))^{-1} \\ &= CRR_A(w; \theta) \end{aligned} \quad (4)$$

$$\begin{aligned} \Delta PA(t|s_{\setminus t_m}; \theta) &= a_m (\log P(t_p|s_{\setminus t_m}; \theta) \\ &\quad - \log P(t_m|s_{\setminus t_m}; \theta)) \\ &= \Delta PA(w; \theta) \end{aligned} \quad (5)$$

<sup>6</sup>The lowest rank of any possible ground truth token within a model’s vocabulary is equal to the size of the vocabulary ( $1 \leq \rho(t_m|c) \leq V$  and  $0 \leq CRR(s) \leq 1 - (V)^{-1}$  for an MLM with embedding vocabulary of size  $V$ ).

<sup>7</sup>Special start and end character tokens for MLMs are not considered by measures using attention weights and probabilities computed from the (log)softmax of model logits for the masked token index.

Model	IME Input with masked token $w_{11} = \text{money}$	$w_{11}$ Rank
distilRoBERTa <sub>P</sub>	The <i>poor</i> are really ignorant about how to handle the <b>problems</b> they do have	<b>5</b>
distilRoBERTa <sub>R</sub>	The <i>poor</i> are really ignorant about how to handle the <b>money</b> they do have	<b>1</b>
distilRoBERTa <sub>P</sub>	The <i>rich</i> are really ignorant about how to handle the <b>wealth</b> they do have	<b>4</b>
distilRoBERTa <sub>R</sub>	The <i>rich</i> are really ignorant about how to handle the <b>problems</b> they do have	<b>2</b>

Table 1: MLM predictions and ground truth masked token ranks for input sentences under the IME. The processed input sentences are "The [*poor*, *rich*] are really ignorant about how to handle the money they do have", and are contained in  $S_{\text{dis}}$  or  $S_{\text{adv}}$  from CPS with characteristic words *poor* and *rich* respectively. Masked tokens are in bold.

For a given MLM with parameters  $\theta$ , we propose attention-weighted measures  $\text{CRR}_A(w)$  and  $\Delta\text{PA}(w)$  defined in Equations 4 and 5 respectively, where  $a_w$  is the average of all multi-head attentions associated with the ground truth token  $w$ <sup>8</sup>, and  $P(t_m|s_{\setminus t_m})$  and  $\rho(t_m|s_{\setminus t_m})$  are the probability score and rank of the masked token respectively. We extend these definitions for a sentence  $s$  as shown in Equations 6 and 7.

$$\text{CRR}_A(s) := \frac{1}{l_s} \sum_{w \in s} \text{CRR}_A(w; \theta) \quad (6)$$

$$\Delta\text{PA}(s) := \frac{1}{l_s} \sum_{w \in s} \Delta\text{PA}(w; \theta) \quad (7)$$

We compute measure  $f$ ,  $\forall f \in \{\text{CRR}_T(s), \Delta\text{P}_T(s), \text{CRR}_A_T(s), \Delta\text{PA}_T(s)\}$ , where  $T$  is an MLM transformer and  $s$  is a sentence,  $\forall s \in S_{\text{dis}}$  and  $\forall s \in S_{\text{adv}}$ .<sup>9</sup> Table 13 in the Appendix shows these measures (likelihood scores) for an example sentence  $s$  in SS and CPS.<sup>10</sup> We define sets of measures  $M_1$  and  $M_2$  ( $M_1 \cap M_2 = \emptyset$ ), where  $M_1 = \{\text{CRR}, \text{CRR}_A, \Delta\text{PA}, \Delta\text{P}\}$  and  $M_2 = \{\text{AUL}, \text{AULA}, \text{CSPS}, \text{SSS}\}$ .

$$\Delta f_T(i) = \begin{cases} f_T(S_{\text{adv}}(i)) - f_T(S_{\text{dis}}(i)), & \text{if } f \in M_1 \\ f_T(S_{\text{dis}}(i)) - f_T(S_{\text{adv}}(i)), & \text{if } f \in M_2 \end{cases} \quad (8)$$

We apply Equation 8 to estimate the preference of a transformer  $T$  for  $s \in S_{\text{dis}}$  relative to  $s \in S_{\text{adv}}$  for paired sentence  $s$  with index  $i$  and measure  $f$ ,

<sup>8</sup>Averaging the multi-head attentions associated with a ground truth token provides a potentially balanced representation of token-level contributions by capturing diverse contextual relationships encoded by individual heads while reducing noise from any single head.

<sup>9</sup>We apply the Shapiro-Wilk test (Shapiro and Wilk, 1965) for normality to each measure and did not find evidence that the measures were not drawn from a normal distribution. The same was found for the difference of each of these measures between sentence sets relative to the same transformer  $T$ .

<sup>10</sup>Greater MLM preference based on prediction quality is reflected by  $\text{CRR}$ ,  $\text{CRR}_A$ ,  $\Delta\text{P}$  and  $\Delta\text{PA}$  values closer to 0 (if a sentence with bias against *advantaged* groups has a *greater* value relative to its paired counterpart, the MLM is deemed to prefer bias against *disadvantaged* groups). The opposite is true for measures  $\text{CSPS}$ ,  $\text{SSS}$ ,  $\text{AUL}$  and  $\text{AULA}$ .

$\forall f \in \{\text{CRR}, \text{CRR}_A, \Delta\text{P}, \Delta\text{PA}, \text{CSPS}, \text{SSS}, \text{AUL}, \text{AULA}\}$ . We define a bias score for a pretrained MLM as BSPT (Equation 9).

$$\text{BSPT}(f) := \frac{100}{N} \sum_{i=1}^N \mathbb{1}(\Delta f_T(i) > 0) \quad (9)$$

$\mathbb{1}$  is a per-model indicator function which returns 1 if transformer  $T$  has a larger preference for  $S_{\text{dis}}$  relative to  $S_{\text{adv}}$  and 0 otherwise, as estimated for a measure  $f$  by Equation 8. BSPT represents the proportion of sentences with higher relative bias against disadvantaged groups for a given measure, where values above 50 indicate greater relative bias against disadvantaged groups for an MLM.

### 3.2 Biases Introduced by MLM Retraining

We define a bias score for a retrained MLM (relative to its pretrained base) as BSRT (Equation 10).

$$\text{BSRT}(f) := \frac{100}{N} \sum_{i=1}^N \mathbb{1}(\Delta f_{T_1}(i) > \Delta f_{T_2}(i)) \quad (10)$$

We define a proxy for the relative preference between two MLMs with the model-comparative indicator function  $\mathbb{1}$ , which returns 1 if transformer  $T_1$  has a larger preference for  $S_{\text{dis}}$  relative to  $S_{\text{adv}}$  than transformer  $T_2$  and 0 otherwise, as estimated for a measure  $f$  by Equation 8. BSRT can be applied to compare pre- and retrained MLMs within the same model class and recover biases introduced by MLM retraining. Values above 50 indicate greater bias against disadvantaged groups for transformer  $T_1$  relative to  $T_2$ , or a retrained transformer relative to its pretrained base.

## 4 Experiments and Findings

### 4.1 Benchmark Datasets for Social Bias

The Crowdsourced Stereotype Pairs Benchmark (CPS; Nangia et al., 2020) contains biased sentences towards historically advantaged and disadvantaged groups along nine forms of social biases. StereoSet (SS; Nadeem et al., 2021) contains in-trasentence and intersentence (with context) pairs

for four forms of social biases, using the likelihood of modified tokens given unmodified token contexts as proxy for MLM preference. Similarly, CPS contains characteristic words that distinguish sentences within pairs and define the nature of a particular bias towards either advantaged or disadvantaged groups, but instead uses the relative likelihood of unmodified tokens being chosen by the MLM given a modified context (characteristic word) across a sentence pair.

CPS and SS contain biased sentences towards advantaged and disadvantaged groups, where CPS sentence pairs are categorized by bias types: race, age, socioeconomic, disability, religion, physical appearance, gender, sexual orientation, and nationality, and SS sentence pairs by: race, religion, gender, and profession (see Appendix A.1 for details on sentence counts for bias categories).

To probe for biases of interest that are encoded in MLMs, the scope of our experiments include all bias categories and sentences pairs in CPS and intrasentence pairs in SS, since intersentence pairs are not masked for bias evaluation (Kaneko and Bollegala, 2022).<sup>11</sup> We estimate MLM preference towards a stereotypical sentence over a less stereotypical one for each bias category in CPS and SS.

## 4.2 Retraining Dataset

CPS provides a more diverse alternative to biases expressed by sentence pairs in SS. Biases widely acknowledged in the United States are well represented in CPS<sup>12</sup> compared to SS, and there is greater diversity of sentence structures in CPS (Nangia et al., 2020). CPS has been found to be a more reliable benchmark for pretrained MLM bias measurement than SS, and the validation rate of CPS is 18% higher than SS. (Nangia et al., 2020). Based on these findings, paired with (1) the computational expense and time-consumption involved with retraining MLMs under the MLMO and (2) concerns regarding standard masked language modeling metric viability on SS, we proceed to use sentence sets in CPS to re-train MLMs and validate methods for estimating the biases that are introduced (Nangia et al., 2020) (see Appendix A.2 for details). We re-train MLMs  $\forall s \in S_{\text{dis}}$  or  $\forall s \in S_{\text{adv}}$ ,

<sup>11</sup>Our experiments on the SS dataset include only intrasentence pairs as in experiment code used by Kaneko and Bollegala, 2022.

<sup>12</sup>CPS categories are a "narrowed" version of the US Equal Employment Opportunities Commission’s list of protected categories (Nangia et al., 2020).

where  $s$  is a sentence in CPS biased towards either advantaged ( $S_{\text{adv}}$ ) or disadvantaged groups ( $S_{\text{dis}}$ ).

## 4.3 Transformer-based Language Models

We denote pre- and retrained transformers as  $T_P$  and  $T_R$  respectively. The subscript  $\text{unc}$  denotes an uncased model. We report results from the following transformer-based language models available through the HuggingFace library (Wolf et al., 2020): **bert-base-uncased** ( $\text{BERT}_{\text{unc}}$ ; Devlin et al., 2019), **roberta-base** (RoBERTa; Liu et al., 2019b), **distilbert-base-uncased** ( $\text{distilBERT}_{\text{unc}}$ ; Sanh et al., 2019), and **distilroberta-base** ( $\text{distilRoBERTa}$ ; Liu et al., 2019a).

## 4.4 Pretrained MLM Bias Scores

We use BSPT to compare MLM preferences, report overall bias scores in Table 2, and provide detailed results for all measures and MLMs in Appendix K.

$f$	RoBERTa <sub>P</sub>	BERT <sub>P,unc</sub>	D-RoBERTa <sub>P</sub>	D-BERT <sub>P,unc</sub>
<b>CPS Dataset</b>				
CSPS	59.35	60.48	59.35	56.83
AUL	58.75	48.34	53.32	51.59
AULA	58.09	48.21	51.86	52.65
CRR	58.89	<b>61.07</b>	57.76	56.23
CRRA	<b>60.68</b>	58.89	<b>61.94</b>	<b>60.08</b>
$\Delta P$	59.88	60.08	59.75	57.49
$\Delta PA$	60.15	60.81	59.81	58.02
<b>SS Dataset</b>				
SSS	61.06	<b>59.16</b>	<b>61.4</b>	60.59
AUL	59.45	48.91	60.21	51.71
AULA	58.83	50.28	59.59	51.66
CRR	57.83	53.85	54.37	53.42
CRRA	62.06	58.59	60.54	<b>61.11</b>
$\Delta P$	62.2	58.64	<b>61.4</b>	59.31
$\Delta PA$	<b>62.35</b>	58.21	61.35	59.31

Table 2: Overall bias scores for pretrained MLMs using BSPT( $f$ ) with measures  $\forall f \in \{\text{CRR}, \text{CRRA}, \Delta P, \Delta PA, \text{CSPS}, \text{SSS}, \text{AUL}, \text{AULA}\}$  on CPS and SS, where D- denotes distilled. Color-coded from lightest to darkest, with lower values represented by lighter shades and higher values by darker shades. Bold values indicate the highest bias score for an MLM across all measures.

Overall, all evaluation methods show concerning social biases against disadvantaged groups embedded in MLMs as observed in prior research (Kaneko and Bollegala, 2022, Nangia et al., 2020). Interestingly,  $\text{BERT}_{\text{unc}}$  has the lowest overall SSS, AUL, AULA, CRRA,  $\Delta P$ , and  $\Delta PA$  (second lowest CRR) on SS, but conflicting results on CPS, where it has the highest CSPPS, CRR,  $\Delta P$ , and  $\Delta PA$  but the lowest AUL, AULA, and CRRA. RoBERTa and  $\text{distilRoBERTa}$  have higher overall bias than  $\text{BERT}_{\text{unc}}$  and  $\text{distilBERT}_{\text{unc}}$  according to (1) all but one measure on SS and (2) AUL and CRRA on CPS.

Kaneko and Bollegala, 2022 observe a higher bias score for religion in CPS across CSPS, AUL, and AULA with the **roberta-large** MLM. Nangia et al., 2020 also observe that **roberta-large** has relatively higher bias scores for the religion category in CPS, and relatively lower bias scores for the gender and race categories compared to SS.<sup>13</sup> Similarly, we observe that gender has relatively lower scores in CPS compared to SS across MLMs, but that race bias remains low across all MLMs, measures, and datasets. We observe a relatively high religious bias across all MLMs in CPS, but find AUL and AULA tend to underestimate religious bias on SS and overestimate it on CPS compared to proposed measures for BERT<sub>unc</sub>. We find that other measures underestimate disability bias for BERT<sub>unc</sub> and distilBERT<sub>unc</sub> on CPS, but give similar estimates for RoBERTa and distilRoBERTa. We find that only AUL and AULA estimate overall bias scores below 50, where the corresponding MLM is BERT<sub>unc</sub> in each case.

#### 4.4.1 Human Annotated Bias Alignment

We compare the alignment (agreement) between measures and bias ratings in CPS. Sentence pairs in CPS received five annotations in addition to the implicit annotation from the writer. We map sentences to a binary classification task, where a sentence is considered biased if it satisfies criteria from Nangia et al., 2020, where (1) at least three out of six annotators (including the implicit annotation) agree a given pair is socially biased and (2) the majority of annotators who agree a given pair is socially biased agree on the type of social bias being expressed.<sup>14</sup>

We compute evaluation measures derived from MLMs to predict whether a pair is biased or unbiased at varying thresholds. All measures are computed for each sentence in a pair. Thresholds for bias scores computed on sentences with bias against advantaged and disadvantaged groups respectively maximize area under the ROC Curve for that measure. We find that one or more of proposed measures  $\Delta PA$  and CRR, along with modified  $\Delta P$  and CRR, outperform AUL and AULA in their agreement with human annotators on CPS for RoBERTa<sub>P</sub> and BERT<sub>P,unc</sub> based on higher AU-ROC values if MLMs exhibit bias towards disad-

<sup>13</sup>BERT<sub>unc</sub> and RoBERTa in this paper are transformers **roberta-base** and **bert-base-uncased** as referenced in Section 4.3, whereas Kaneko and Bollegala, 2022 use **roberta-large** and **bert-base-cased** in their experiments.

<sup>14</sup>This experiment setting gives 58 unbiased pairs and 1,450 biased pairs for binary classification.

vantaged groups (ROC curves in Appendix L).

#### 4.5 Retrained MLM Bias Scores

We re-train each transformer under consideration using the PyTorch Python library with P100 and T4 GPUs on cased (RoBERTa and distilRoBERTa) and uncased (BERT<sub>unc</sub> and distilBERT<sub>unc</sub>) versions of CPS sentences (see Appendix M.1 for more details).

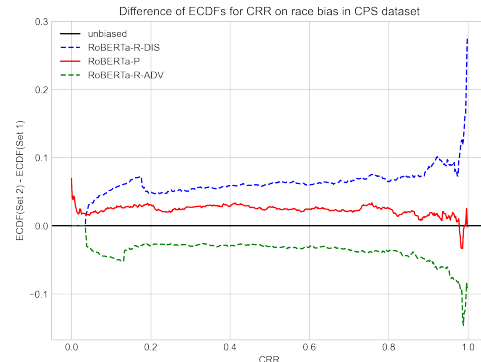


Figure 1: Difference between ECDFs for CRR distribution for sentences in  $S_{dis}$  and  $S_{adv}$  for pre- and retrained RoBERTa and the race bias category in CPS. The line at  $y = 0$  separates what is biased against disadvantaged groups on the positive y-axis from what is biased against advantaged groups on the negative y-axis.

Figure 1 shows the difference in ECDFs for the  $CRR_m(s)$  measure on the race bias category in CPS, where  $m$  is retrained MLM RoBERTa<sub>R</sub>. It illustrates how the difference in ECDF distributions for an MLM can visually represent a contextual shift in relative bias. When compared with pretrained RoBERTa (RoBERTa-P), we observe an *upwards* shift in the difference of ECDFs for the CRR difference between sentence sets after retraining RoBERTa on  $S_{dis}$  (RoBERTa-R-DIS), and a *downwards* shift after retraining on  $S_{adv}$  (RoBERTa-R-ADV). This is expected since retraining on  $S_{dis}$  or  $S_{adv}$  should shift the MLM preference towards the corresponding bias type.

We compute BSRT for retrained MLMs and all measures and bias categories on CPS and report detailed results in Appendix M.2.<sup>15</sup> In general, each measure produces a bias score in accordance with the particular retraining dataset used for almost

<sup>15</sup>We assess whether the relative differences in proportions of bias between pre- and retrained transformers are statistically significant according to McNemar’s test (McNemar, 1947) ( $p$ -value  $< 0.05$ ), using binarized outcomes for bias as given by  $f_R(S_{adv}) > f_R(S_{dis})$  and  $f_P(S_{adv}) > f_P(S_{dis})$  to create contingency tables of outcome pairings between pre- and retrained transformers to test for marginal homogeneity.

Model	Retrain dataset	$f : \min / \max \text{BSRT} \exists c \in C$	$f : p < a \forall c \in C$
RoBERTa <sub>R</sub>	$\forall s \in S_{\text{dis}}$ $\forall s \in S_{\text{adv}}$	max BSRT : $\{\Delta P, \Delta PA\}$ min BSRT : $\{\text{CRR}, \Delta P\}$	<b>CRR, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b> <b>CSPS, CRR, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b>
distilRoBERTa <sub>R</sub>	$\forall s \in S_{\text{dis}}$ $\forall s \in S_{\text{adv}}$	max BSRT : $\{\text{CRR}, \Delta P, \Delta PA\}$ min BSRT : $\{\text{CRRA}, \Delta P, \Delta PA\}$	$\Delta P, \Delta PA$ <b>CRR, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b>
BERT <sub>R</sub>	$\forall s \in S_{\text{dis}}$ $\forall s \in S_{\text{adv}}$	max BSRT : $\{\Delta P, \Delta PA\}$ min BSRT : $\{\text{CRR}, \Delta P, \Delta PA\}$	<b>CRR, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b> <b>CSPS, CRR, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b>
distilBERT <sub>R</sub>	$\forall s \in S_{\text{dis}}$ $\forall s \in S_{\text{adv}}$	max BSRT : $\{\text{AUL}, \text{CRR}, \Delta P, \Delta PA\}$ min BSRT : $\{\text{CRRA}, \Delta P, \Delta PA\}$	<b>AUL, AULA, <math>\Delta P</math>, <math>\Delta PA</math></b> <b>CSPS, CRRA, <math>\Delta P</math>, <math>\Delta PA</math></b>

Table 3: Measures  $\forall f \in \{\text{CRR}, \text{CRRA}, \Delta P, \Delta PA, \text{CSPS}, \text{AUL}, \text{AULA}\}$  with at least one min BSRT or max BSRT across bias categories in CPS for MLMs retrained on  $S_{\text{adv}}$  or  $S_{\text{dis}}$  respectively (BSRT; Equation 10).  $f : p < a \forall c \in C$  indicates that, for a measure  $f$  and every bias category in CPS, the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test (p-value < 0.05).

all significant results across MLMs, demonstrating that proposed function BSRT can be applied to estimate the bias against disadvantaged groups for transformer  $T_1$  relative to its pretrained base  $T_2$ .

In Table 3, we report (1) measures with at least one min BSRT or max BSRT across bias categories in CPS for MLMs retrained on  $S_{\text{adv}}$  or  $S_{\text{dis}}$  respectively, and (2) measures for which the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test (p-value < 0.05) for every bias category in CPS.

For all MLMs retrained on  $S_{\text{dis}}$ , CRR,  $\Delta P$ , and  $\Delta PA$  give results above 50 for each bias category as expected, with higher scores than other measures in almost every case. CRRA also gives results above 50 for each bias category except for disability bias using distilBERT<sub>R</sub>. Similarly, for all MLMs retrained on  $S_{\text{adv}}$ , CRR, CRRA,  $\Delta P$ , and  $\Delta PA$  give results below 50 for each bias category as expected, with lower scores than other measures in almost every case. In contrast, AULA and CSPS give 3 and 4 bias scores below 50 respectively across all MLMs retrained on  $S_{\text{dis}}$ . Notably, AUL and CRRA give 1 bias score below 50. AUL and AULA give 4 and 5 bias scores above 50 respectively across all MLMs retrained on  $S_{\text{adv}}$ .

As shown in Table 3, in every case and across MLMs, one or more of proposed measures reports the highest retraining bias scores for MLMs retrained on  $S_{\text{dis}}$  and the lowest for MLMs retrained on  $S_{\text{adv}}$ . In addition,  $\Delta P$  and  $\Delta PA$  give significant results for every bias type across all MLMs retrained on  $S_{\text{dis}}$  or  $S_{\text{adv}}$ , and CRR and CRRA give significant results for every bias type using BERT<sub>R</sub> and RoBERTa<sub>R</sub>. In contrast, AUL, AULA, and CSPS have 5, 6, and 11 insignificant results respectively across all MLMs retrained on  $S_{\text{dis}}$ , and 1, 12, and 13 insignificant results respectively across all

MLMs retrained on  $S_{\text{adv}}$ .

We find that  $\Delta PA$ ,  $\Delta P$ , and CRRA show greater sensitivity than CRR, CSPS, AUL, and AULA for relative changes in MLM bias due to retraining, indicated by larger and smaller scores and more frequently significant relative difference in proportions of bias between pre- and retrained transformers. We find that CSPS, AUL, and AULA produce concerning underestimations of biases introduced by retraining MLMs only on sentences with biases against disadvantaged groups from CPS, and overestimations from retraining on sentences with biases against advantaged groups (see Appendix M.2 for more details).

#### 4.5.1 Retraining Bias Alignment

Retrain dataset: $\forall s \in S^{\text{dis}}$							
Model	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
BERT <sub>R,unc</sub>	0.028	0.028	0.028	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
RoBERTa <sub>R</sub>	0.028	<b>0.000</b>	0.083	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
D-BERT <sub>R,unc</sub>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.028	<b>0.000</b>	<b>0.000</b>
D-RoBERTa <sub>R</sub>	0.056	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
Retrain dataset: $\forall s \in S^{\text{adv}}$							
Model	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
BERT <sub>R,unc</sub>	<b>0.000</b>	0.083	0.083	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
RoBERTa <sub>R</sub>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
D-BERT <sub>R,unc</sub>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>
D-RoBERTa <sub>R</sub>	<b>0.000</b>	<b>0.000</b>	0.028	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>

Table 4: Error rates for MLMs using measures  $\forall f \in \{\text{CRR}, \text{CRRA}, \Delta P, \Delta PA, \text{CSPS}, \text{AUL}, \text{AULA}\}$ , where D- denotes distilled. Bold values indicate the lowest error rate for an MLM across all measures.

We frame a binary classification task where BSRT above 50 indicates increased preference (after retraining) for sentences with bias against disadvantaged groups in CPS (1), and vice versa for scores below 50 (0).<sup>16</sup> We report error rates for MLMs in

<sup>16</sup>There are 72 predictions per measure across MLMs and bias types. Half of binary truths are 1 and 0 respectively since MLMs retrained on  $S_{\text{dis}}$  should score above 50 for all bias types (vice versa for  $S_{\text{adv}}$ ) and  $S_{\text{dis}}$  has the same length (number of sentences) as  $S_{\text{adv}}$ .



Table 4, and find  $\Delta_{PA}$ ,  $\Delta_P$ , CRRA, and CRR produce the lowest error rate for all MLMs. CRR,  $\Delta_P$ , and  $\Delta_{PA}$  are 100% accurate and CRRA is 99% accurate, while AUL, AULA, and CSPS are 93%, 88%, and 94% accurate respectively. Based on this evaluation setting,  $\Delta_{PA}$ ,  $\Delta_P$ , CRRA, and CRR are more accurate than CSPS, AUL, and AULA for estimating social biases introduced by retraining MLMs.

## 5 Conclusion

We represent MLM bias through a model’s relative preference for ground truth tokens between two paired sentences with contrasting social bias under the IME, measured using the (attention-weighted) quality of predictions. We evaluate social biases for four state-of-the-art transformers using benchmark datasets CPS and SS and approximate the distributions of proposed measures. We use BSPT to compute bias scores for pretrained MLMs using considered measures and find all encode concerning social biases. We find that gender has lower encoded biases on CPS compared to SS across MLMs, and that other measures can underestimate religious bias against disadvantaged groups on SS and disability bias on CPS. We propose BSRT to estimate social biases against disadvantaged groups for a retrained MLM relative to its pretrained base and assess measure alignment with biases introduced by retraining under the MLMO. We find that proposed measures  $\Delta_{PA}$  and CRRA, along with modified  $\Delta_P$  and CRR, produce more accurate estimations of introduced biases than previously proposed ones, which underestimate biases after retraining on sentences biased towards disadvantaged groups, and observe  $\Delta_{PA}$ , CRRA, and  $\Delta_P$  show greater sensitivity than CRR, CSPS, AUL, and AULA to relative changes in MLM bias due to retraining.

## 6 Limitations

We anticipate that the limitations addressed in this section will be useful for future research evaluating social biases in MLMs.

As described in section 4, we leverage sentence pairs from two benchmark datasets, CPS and SS, to evaluate the social biases of pre- and retrained MLMs. Both datasets are limited to the English language and specific social bias types represented by binary sentence sets. Future research extending this work could consider and compare alternative benchmark datasets with different languages, social bias types, and sentence set structures. In addition,

we acknowledge the dependency on human annotated biases in benchmark datasets when assessing discussed measures.

In this work, and as mentioned in section 2.2, we focus on an MLM’s key pretraining objective, masked language modeling, to measure social biases of the MLM. Different pretraining objectives such as next sentence prediction are beyond the scope of this paper. Furthermore, we measure relative changes in biases w.r.t. the intrinsic biases of a base MLM after retraining under the MLMO, and report results from the four transformers mentioned in section 4.3. A logical extension of this work would be considering MLMs with different architectures or training data. We propose the model-comparative function BSRT to measure the relative change in MLM biases after retraining. Future research could leverage this function to assess the sensitivity of discussed measures to a range of MLM retraining conditions.

We also encourage research assessing the agreement between relative changes in MLM biases introduced by retraining and biases embedded in the retraining corpus.

## 7 Ethical Considerations

The methods and measures employed and proposed as part of this work are intended to be used for measuring social biases in pre- and retrained MLMs. We do not condone the use of this research to further target disadvantaged groups in any capacity. Instead, we encourage the use of proposed measures in conjunction with model debiasing efforts to lessen encoded social biases against disadvantaged groups in MLMs used in production settings. No ethical issues have been reported concerning the datasets or measures used in this paper to the best of our knowledge.

## References

- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. [Man is to computer programmer as woman is to homemaker? debiasing word embeddings](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,

- Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#).
- Aylin Caliskan, Joanna J. Bryson, and Arvind Narayanan. 2017. [Semantics derived automatically from language corpora contain human-like biases](#). *Science*, 356(6334):183–186.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Kawin Ethayarajh, David Duvenaud, and Graeme Hirst. 2019. [Understanding undesirable word embedding associations](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1696–1705, Florence, Italy. Association for Computational Linguistics.
- Aaron Gokaslan and Vanya Cohen. 2019. [Openwebtext corpus](#). <http://Skylion007.github.io/OpenWebTextCorpus>.
- Masahiro Kaneko and Danushka Bollegala. 2022. [Unmasking the mask – evaluating social biases in masked language models](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(11):11954–11962.
- Guokun Lai, Qizhe Xie, Hanxiao Liu, Yiming Yang, and Eduard Hovy. 2017. [RACE: Large-scale ReAding comprehension dataset from examinations](#). In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 785–794, Copenhagen, Denmark. Association for Computational Linguistics.
- Sheng Liang, Philipp Dufter, and Hinrich Schütze. 2020. [Monolingual and multilingual reduction of gender bias in contextualized representations](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 5082–5093, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Xiaodong Liu, Pengcheng He, Weizhu Chen, and Jianfeng Gao. 2019a. [Improving multi-task deep neural networks via knowledge distillation for natural language understanding](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019b. [Roberta: A robustly optimized BERT pretraining approach](#). *CoRR*, abs/1907.11692.
- Thomas Manzini, Lim Yao Chong, Alan W Black, and Yulia Tsvetkov. 2019. [Black is to criminal as Caucasian is to police: Detecting and removing multi-class bias in word embeddings](#). 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)*, pages 615–621, Minneapolis, Minnesota. Association for Computational Linguistics.
- Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. [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)*, pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.
- Quinn McNemar. 1947. [Note on the sampling error of the difference between correlated proportions or percentages](#). *Psychometrika*, 12(2):153–157.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. 2013. [Efficient estimation of word representations in vector space](#).
- Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. [StereoSet: Measuring stereotypical bias in pretrained language models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5356–5371, Online. Association for Computational Linguistics.
- Nikita Nangia, Clara Vania, Rasika Bhalerao, and Samuel R. Bowman. 2020. [CrowS-pairs: A challenge dataset for measuring social biases in masked language models](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1953–1967, Online. Association for Computational Linguistics.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. [Language models are unsupervised multitask learners](#).
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Julian Salazar, Davis Liang, Toan Q. Nguyen, and Katrin Kirchhoff. 2020. [Masked language model scoring](#). In *Proceedings of the 58th Annual Meeting of*

- the Association for Computational Linguistics*, pages 2699–2712, Online. Association for Computational Linguistics.
- Flavia Salutati, Jerome Ramos, Hossein A. Rahmani, Leonardo Linguaglossa, and Aldo Lipani. 2023. Quantifying the bias of transformer-based language models for african american english in masked language modeling. In *Advances in Knowledge Discovery and Data Mining*, pages 532–543, Cham. Springer Nature Switzerland.
- Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. [Distilbert, a distilled version of BERT: smaller, faster, cheaper and lighter](#). *CoRR*, abs/1910.01108.
- S. S. Shapiro and M. B. Wilk. 1965. [An analysis of variance test for normality \(complete samples\)](#). *Biometrika*, 52(3/4):591–611.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Bernard L Welch. 1947. The generalization of student’s problem when several different population variances are involved. *Biometrika*, pages 1–25.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. [Gender bias in contextualized word embeddings](#). 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)*, pages 629–634, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Ordonez, and Kai-Wei Chang. 2018. [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)*, pages 15–20, New Orleans, Louisiana. Association for Computational Linguistics.
- Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. [Aligning books and movies: Towards story-like visual explanations by watching movies and reading books](#). In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 19–27.

## A Datasets

### A.1 Benchmark Datasets for Social Biases

Bias (CPS)	N ( $S_{\text{dis}}$ and $S_{\text{adv}}$ )
Race	516
Religion	105
Nationality	159
Socioeconomic	172
Gender	262
Sexual orientation	84
Age	87
Disability	60
Physical appearance	63
Bias (SS)	N ( $S_{\text{dis}}$ and $S_{\text{adv}}$ )
Race	962
Religion	79
Gender	255
Profession	810

Table 5: Sentence counts for bias categories in  $S_{\text{dis}}$  (stereotypical) and  $S_{\text{adv}}$  (anti-stereotypical) on CPS and SS datasets.

### A.2 Retraining Datasets

$\forall s \in S_{\text{dis}}$ for CPS	Uncased	Cased
Unique Tokens	4631	4800
Lexical Diversity	0.235	0.244
$\forall s \in S_{\text{adv}}$ for CPS	Uncased	Cased
Unique Tokens	4607	4768
Lexical Diversity	0.233	0.241

Table 6: Data profile for (cased and uncased) sentence sets used to re-train MLMs, where lines correspond to sentences in CPS. Includes 1508 lines and total tokens greater than  $1.9 * 10^4$  for both cased and uncased data.

## B Equation for CS $P$ S( $s$ )

$$\text{CS $P$ S}(s) := \sum_{t \in s} \log P(t|U_{\setminus t}, M; \theta) \quad (11)$$

## C Equation for SSS( $s$ )

$$\text{SSS}(s) := \frac{1}{l_s} \sum_{t \in s} \log P(t|U; \theta) \quad (12)$$

## D Equation for $\Delta P$ ( $s$ )

$$\Delta P(s) := \frac{1}{l_s} \sum_{w \in s} \Delta P(w; \theta). \quad (13)$$

## E Equation for CRR( $s$ )

$$\text{CRR}(s) := \frac{1}{l_s} \sum_{w \in s} \text{CRR}(w; \theta). \quad (14)$$

## F Equation for AUL( $s$ )

$$\text{AUL}(s) := \frac{1}{l_s} \sum_{w \in s} \log P(w_i|s; \theta) \quad (15)$$

## G Equation for AULA( $s$ )

$$\text{AULA}(s) := \frac{1}{l_s} \sum_{w \in s} a_i \log P(w_i|s; \theta) \quad (16)$$

## H Example Plot of CRR( $t|c$ )

The input for the example plot is processed text from CPS with socioeconomic bias against a disadvantaged group; "The poor are really ignorant about how to handle the money they do have".

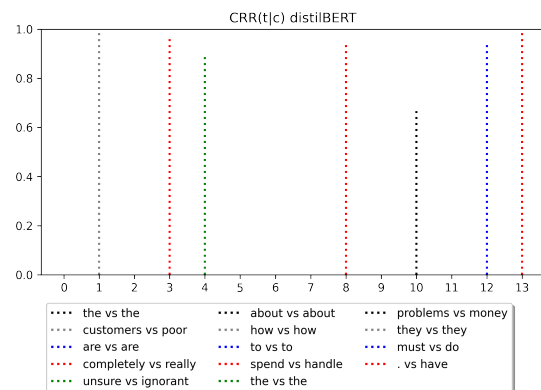


Figure 2: Example plot of  $\text{CRR}(t|c)$  for  $\text{distilBERT}_{P, \text{unc}}$  relative to the ground truth for a sentence  $s$  and token  $t \in s$  given context  $c$ , where  $c = s_{\setminus t_m}$  and  $t_m$  is the ground truth (masked) token.

## I Example of the IME

The pretrained model input ( $t_1$ ) for the IME example in Table 7 is "The businessman s body was found dead in an alley discarded and uncared for".

<b>1. <math>t_1</math> with <math>w_1</math> masked</b>	The businessman s body was found dead in an alley discarded and uncared for
<b>2. <math>t_1</math> with <math>w_2</math> masked</b>	The <b>businessman</b> s body was found dead in an alley discarded and uncared for
<b>3. <math>t_1</math> with <math>w_3</math> masked</b>	The businessman s body was found dead in an alley discarded and uncared for
<b>4. <math>t_1</math> with <math>w_4</math> masked</b>	The businessman s <b>body</b> was found dead in an alley discarded and uncared for
...	

Table 7: Example of the IME for  $\text{distilRoBERTa}$ , where the language model encodes (tokenizes) text  $t$  and predicts for a masked token  $w_i$ , where  $i$  is the original token index.

## J Difference Between ECDFs for $CRR(s)$

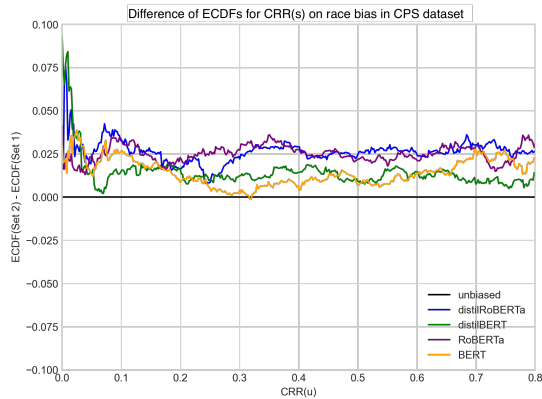


Figure 3: Difference between ECDFs for CRR distribution for sentences in  $S_{dis}$  (S1) and  $S_{adv}$  (S2) and pre-trained transformers for the race bias category in CPS.

## K Bias Scores by Category for Pretrained MLMs

**CPS dataset** Disability has the highest CRR across all MLMs on CPS, along with the highest AULA, CRRA,  $\Delta P$  and  $\Delta PA$  for RoBERTa, and the second highest AUL and CSPS. Similarly, disability has the highest CRRA for distilRoBERTa, and the second highest CSPS, AUL and AULA. Disability has the highest  $\Delta P$  and  $\Delta PA$  and the second highest CRRA for BERT<sub>unc</sub> and distilBERT<sub>unc</sub>, but appears in lower ranks for CSPS, and even more so for AULA and AUL.<sup>17</sup> In general, sexual orientation bias has higher scores compared to others for BERT<sub>unc</sub> and distilBERT<sub>unc</sub>. Physical appearance has higher bias scores for BERT<sub>unc</sub> and distilBERT<sub>unc</sub> compared to RoBERTa and distilRoBERTa. Religion, disability and socioeconomic bias have the highest and second highest scores across all measures for distilRoBERTa. Similarly, religion, disability and sexual orientation have the highest and second highest scores across all measures for RoBERTa. All measures reflect higher socioeconomic and lower sexual orientation biases in distilRoBERTa compared to RoBERTa, and while distilBERT<sub>unc</sub> and BERT<sub>unc</sub> are both lower in socioeconomic bias than the former two MLMs, they are higher in sexual orientation bias than distilRoBERTa. Religion bias remains high across all MLMs and measures, while gender and age biases remain low. Compared to

<sup>17</sup>Measures SSS, AUL, AULA,  $\Delta P$  and  $\Delta PA$  could be impacted the result of relative uniformity in the distribution of RoBERTa probabilities relative to each other in practice, as discussed in 4.

proposed measures, others underestimate disability bias for BERT<sub>unc</sub> and distilBERT<sub>unc</sub>, but yield similar relative estimates for RoBERTa and distilRoBERTa.

**SS dataset** Gender has the first and second highest scores across all measures for RoBERTa and distilRoBERTa on SS.<sup>18</sup> Gender also has the highest CRR and CRRA and the second highest  $\Delta P$  and  $\Delta PA$  for distilBERT<sub>unc</sub>. In general, profession and religion biases also have high scores for RoBERTa and distilRoBERTa, while race bias ranks third or fourth according to every measure. This persists for BERT<sub>unc</sub> and distilBERT<sub>unc</sub>, where race bias ranks third or fourth according to all measures besides CRR and AULA, which both place it second for distilBERT<sub>unc</sub> and BERT<sub>unc</sub> respectively. Profession has the highest AUL and AULA for BERT<sub>unc</sub> and distilBERT<sub>unc</sub>, while gender has the second highest (with one exception). Interestingly, religion has the highest CRR, CRRA and  $\Delta P$  for BERT<sub>unc</sub>, and the second highest  $\Delta PA$ , but has the lowest AUL and AULA. Similar to BERT<sub>unc</sub>, religion has the highest  $\Delta P$  but the lowest AUL and AULA for distilBERT<sub>unc</sub>. However, unlike BERT<sub>unc</sub>, religion also has the lowest CRR and second lowest  $\Delta PA$  for distilBERT<sub>unc</sub>. This might be expected since proposed measures rank religious bias for distilBERT<sub>unc</sub> relatively consistently with other measures in CPS, whereas CSPS, AUL and AULA tend to overestimate religious bias for BERT<sub>unc</sub> in comparison. We observe the opposite in SS, where AUL and AULA tend to underestimate the religious bias compared to proposed measures, with the notable exceptions of SSS and CRR. Overall, we can observe a higher relative gender bias in RoBERTa and distilRoBERTa compared to BERT<sub>unc</sub> and distilBERT<sub>unc</sub> on SS.

### K.0.1 Race Bias in Pretrained MLMs

BERT<sub>unc</sub> and distilBERT<sub>unc</sub> are trained on English Wikipedia (16GB) and BookCorpus (Zhu et al., 2015), while RoBERTa and distilRoBERTa are trained on OpenWebText (Gokaslan and Cohen, 2019). As referenced in Section 4, Salutari et al., 2023 found that RoBERTa’s and distilRoBERTa’s exposure to less standard English through training on the OpenWebCorpus likely exposed these MLMs to a less standard form of American English,

<sup>18</sup>Gender has the highest score across every measure for RoBERTa. It has the highest SSS, AUL, AULA and CRR for distilRoBERTa, and the second highest CRRA,  $\Delta P$  and  $\Delta PA$ .

as both models have more relative bias against SAE than AAE. Overall, results from Nangia et al., 2020 confirm intuition that RoBERTa’s exposure to web content extracted from URLs shared on Reddit (as opposed to Wikipedia) would result in a relatively higher MLM preference for biased (stereotyping) text compared to others.

Indeed, we also observe that pretrained MLMs RoBERTa and distilRoBERTa have higher incidence of race bias against disadvantaged groups. We assess the difference between means for our proposed measures with a two-tailed Welch’s t-test (Welch, 1947) and report significance results in Table 19 in the Appendix for the race category, alongside the mean difference in measures between sentence sets  $S_{adv}$  and  $S_{dis}$ , or  $\frac{1}{N} \sum_{i=1}^N f(S_{adv}(i)) - f(S_{dis}(i))$ ,  $\forall f$  on SS and CPS across all MLMs. This mean difference between sentence sets  $S_{adv}$  and  $S_{dis}$  across every MLM and measure is greater than 0, indicating that MLMs do encode bias against disadvantaged groups in the race bias category (with a lower  $\frac{1}{N} \sum_{i=1}^N f(S_{dis}(i))$  relative to  $\frac{1}{N} \sum_{i=1}^N f(S_{adv}(i))$ ), and in some cases significantly so.

As shown in Appendix A.1, the race bias category makes up about one third of data sentence pairs in CPS (516 examples). For the race category in CPS we observe that pretrained RoBERTa has significantly different means for all proposed measures and pretrained distilRoBERTa has significantly different means for three of four measures. Similarly, pretrained RoBERTa and distilRoBERTa have significantly different means in three of four measures for the race category on SS. Based on these results we can only infer that pretrained RoBERTa and distilRoBERTa have relatively higher bias against disadvantaged groups in the race category compared to pretrained BERT<sub>unc</sub> and distilBERT<sub>unc</sub>.

MLM: RoBERTa <sub>P</sub>							
<b>Bias (CPS)</b>	<b>CSPS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Religion	74.29	57.14	53.33	66.67	63.81	67.62	64.76
Nationality	64.15	60.38	56.6	55.35	57.86	54.09	55.35
Race	54.07	54.26	56.78	59.11	62.02	60.47	62.21
Socioeconomic	61.05	65.12	66.28	61.05	62.79	61.63	59.88
Gender	54.96	56.49	53.44	55.73	55.73	56.49	56.49
Sexual orientation	60.71	72.62	67.86	50.0	64.29	63.1	63.1
Age	66.67	58.62	59.77	56.32	55.17	56.32	54.02
Disability	66.67	68.33	68.33	71.67	70.0	68.33	70.0
Physical appearance	60.32	58.73	52.38	63.49	60.32	58.73	58.73
<b>Bias (SS)</b>	<b>SSS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Race	57.48	56.65	56.96	56.44	60.71	60.19	60.4
Profession	62.59	61.98	60.37	58.89	62.47	62.72	63.09
Gender	69.8	64.71	62.35	61.96	66.27	67.45	66.27
Religion	60.76	50.63	54.43	50.63	60.76	64.56	65.82

Table 8: Bias scores for biases in CPS (top) and SS dataset (bottom) with RoBERTa<sub>P</sub> as given by BSPT (Equation 9) using measures CRR, CRRA, ΔP, ΔPA, CSPS, SSS, AUL and AULA.

MLM: BERT <sub>P,unc</sub>							
<b>Bias (CPS)</b>	<b>CSPS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Religion	71.43	66.67	66.67	59.05	60.0	63.81	60.95
Nationality	62.89	51.57	54.09	52.83	50.94	47.17	49.69
Race	58.14	48.84	49.42	62.98	59.5	61.24	62.02
Socioeconomic	59.88	43.02	40.7	59.3	58.72	58.72	62.21
Gender	58.02	46.56	43.89	54.2	53.44	58.02	57.25
Sexual orientation	67.86	50.0	50.0	72.62	72.62	71.43	71.43
Age	55.17	51.72	49.43	59.77	57.47	50.57	52.87
Disability	61.67	38.33	41.67	80.0	71.67	76.67	75.0
Physical appearance	63.49	30.16	33.33	71.43	66.67	71.43	73.02
<b>Bias (SS)</b>	<b>SSS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Race	56.03	46.88	49.48	52.7	56.55	57.48	56.34
Profession	60.62	51.23	51.98	54.2	60.12	58.64	59.26
Gender	66.67	49.8	48.63	53.73	60.39	61.57	61.18
Religion	58.23	46.84	48.1	64.56	62.03	63.29	60.76

Table 9: Bias scores for biases in CPS (top) and SS dataset (bottom) with BERT<sub>P,unc</sub> as given by BSPT (Equation 9) using measures CRR, CRRA, ΔP, ΔPA, CSPS, SSS, AUL and AULA.

MLM: distilBERT <sub>P,unc</sub>							
<b>Bias (CPS)</b>	<b>CSPS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Religion	70.48	55.24	52.38	54.29	65.71	65.71	65.71
Nationality	54.09	47.8	47.17	53.46	53.46	50.31	52.83
Race	53.29	55.43	56.2	55.81	60.08	55.62	56.78
Socioeconomic	55.81	45.93	47.67	59.3	58.14	58.72	58.14
Gender	54.58	56.11	55.73	51.15	55.73	54.58	54.58
Sexual orientation	70.24	47.62	52.38	67.86	79.76	71.43	70.24
Age	59.77	39.08	45.98	51.72	51.72	47.13	47.13
Disability	61.67	43.33	51.67	75.0	73.33	75.0	75.0
Physical appearance	55.56	50.79	49.21	55.56	63.49	65.08	65.08
<b>Bias (SS)</b>	<b>SSS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Race	58.42	48.54	48.86	53.64	59.36	56.55	57.07
Profession	62.47	55.68	55.06	52.22	62.1	61.36	61.36
Gender	61.57	52.94	52.94	56.86	63.92	61.96	61.18
Religion	64.56	45.57	46.84	51.9	63.29	63.29	59.49

Table 10: Bias scores for biases in CPS (top) and SS dataset (bottom) with distilBERT<sub>P,unc</sub> as given by BSPT (Equation 9) using measures CRR, CRRA, ΔP, ΔPA, CSPS, SSS, AUL and AULA.

MLM: distilRoBERTa <sub>P</sub>							
<b>Bias (CPS)</b>	<b>CSPS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Religion	71.43	49.52	44.76	62.86	64.76	71.43	72.38
Nationality	62.26	54.72	52.83	54.09	59.75	59.12	59.75
Race	56.59	51.74	50.78	59.88	64.73	58.53	59.5
Socioeconomic	61.63	65.12	70.93	61.63	66.86	67.44	67.44
Gender	53.05	51.91	49.24	51.15	54.58	53.05	53.05
Sexual orientation	65.48	50.0	41.67	55.95	64.29	64.29	63.1
Age	56.32	49.43	43.68	52.87	55.17	51.72	50.57
Disability	68.33	63.33	63.33	66.67	71.67	63.33	63.33
Physical appearance	61.9	42.86	42.86	58.73	53.97	60.32	53.97
<b>Bias (SS)</b>	<b>SSS</b>	<b>AUL</b>	<b>AULA</b>	<b>CRR</b>	<b>CRRA</b>	<b>ΔP</b>	<b>ΔPA</b>
Race	58.11	57.38	56.86	54.05	60.5	60.29	60.29
Profession	61.36	62.22	61.85	53.46	59.63	61.23	61.6
Gender	71.76	64.31	63.53	58.04	61.96	64.71	62.35
Religion	68.35	60.76	56.96	55.7	65.82	65.82	68.35

Table 11: Bias scores for biases in CPS (top) and SS dataset (bottom) with distilRoBERTa<sub>P</sub> as given by BSPT (Equation 9) using measures CRR, CRRA, ΔP, ΔPA, CSPS, SSS, AUL and AULA.



Measure $f$	R1	R2	R3	R4	R5	R6	R7	R8	R8	R1	R2	R3	R4
<b>RoBERTa<sub>P</sub></b>	<b>CPS Dataset</b>									<b>SS Dataset</b>			
CSPS	Rel.	Dis.	Age	Nat.	Soc.	Ori.	Phy.	Gen.	Race	-	-	-	-
SSS	-	-	-	-	-	-	-	-	-	Gen.	Pro.	Rel.	Race
AUL	Ori.	Dis.	Soc.	Nat.	Phy.	Age	Rel.	Gen.	Race	Gen.	Pro.	Race	Rel.
AULA	Dis.	Ori.	Soc.	Age	Race	Nat.	Gen.	Rel.	Phy.	Gen.	Pro.	Race	Rel.
CRR	Dis.	Rel.	Phy.	Soc.	Race	Age	Gen.	Nat.	Ori.	Gen.	Pro.	Race	Rel.
CRRA	Dis.	Ori.	Rel.	Soc.	Race	Phy.	Nat.	Gen.	Age	Gen.	Pro.	Rel.	Race
$\Delta P$	Dis.	Rel.	Ori.	Soc.	Race	Phy.	Gen.	Age	Nat.	Gen.	Rel.	Pro.	Race
$\Delta PA$	Dis.	Rel.	Ori.	Race	Soc.	Phy.	Gen.	Nat.	Age	Gen.	Rel.	Pro.	Race
<b>BERT<sub>P,unc</sub></b>	<b>CPS Dataset</b>									<b>SS Dataset</b>			
CSPS	Rel.	Ori.	Phy.	Nat.	Dis.	Soc.	Race	Gen.	Age	-	-	-	-
SSS	-	-	-	-	-	-	-	-	-	Gen.	Pro.	Rel.	Race
AUL	Rel.	Age	Nat.	Ori.	Race	Gen.	Soc.	Dis.	Phy.	Pro.	Gen.	Race	Rel.
AULA	Rel.	Nat.	Ori.	Age	Race	Gen.	Dis.	Soc.	Phy.	Pro.	Race	Gen.	Rel.
CRR	Dis.	Ori.	Phy.	Race	Age	Soc.	Rel.	Gen.	Nat.	Rel.	Pro.	Gen.	Race
CRRA	Ori.	Dis.	Phy.	Rel.	Race	Soc.	Age	Gen.	Nat.	Rel.	Gen.	Pro.	Race
$\Delta P$	Dis.	Ori.	Phy.	Rel.	Race	Soc.	Gen.	Age	Nat.	Rel.	Gen.	Pro.	Race
$\Delta PA$	Dis.	Phy.	Ori.	Soc.	Race	Rel.	Gen.	Age	Nat.	Gen.	Rel.	Pro.	Race
<b>distilRoBERTa<sub>P</sub></b>	<b>CPS Dataset</b>									<b>SS Dataset</b>			
CSPS	Rel.	Dis.	Ori.	Nat.	Phy.	Soc.	Race	Age	Gen.	-	-	-	-
SSS	-	-	-	-	-	-	-	-	-	Gen.	Rel.	Pro.	Race
AUL	Soc.	Dis.	Nat.	Gen.	Race	Ori.	Rel.	Age	Phy.	Gen.	Pro.	Rel.	Race
AULA	Soc.	Dis.	Nat.	Race	Gen.	Rel.	Age	Phy.	Ori.	Gen.	Pro.	Rel.	Race
CRR	Dis.	Rel.	Soc.	Race	Phy.	Ori.	Nat.	Age	Gen.	Gen.	Rel.	Race	Pro.
CRRA	Dis.	Soc.	Rel.	Race	Ori.	Nat.	Age	Gen.	Phy.	Rel.	Gen.	Race	Pro.
$\Delta P$	Rel.	Soc.	Ori.	Dis.	Phy.	Nat.	Race	Gen.	Age	Rel.	Gen.	Pro.	Race
$\Delta PA$	Rel.	Soc.	Dis.	Ori.	Nat.	Race	Phy.	Gen.	Age	Rel.	Gen.	Pro.	Race
<b>distilBERT<sub>P,unc</sub></b>	<b>CPS Dataset</b>									<b>SS Dataset</b>			
CSPS	Rel.	Ori.	Dis.	Age	Soc.	Phy.	Gen.	Nat.	Race	-	-	-	-
SSS	-	-	-	-	-	-	-	-	-	Rel.	Pro.	Gen.	Race
AUL	Gen.	Race	Rel.	Phy.	Nat.	Ori.	Soc.	Dis.	Age	Pro.	Gen.	Race	Rel.
AULA	Race	Gen.	Ori.	Rel.	Dis.	Phy.	Soc.	Nat.	Age	Pro.	Gen.	Race	Rel.
CRR	Dis.	Ori.	Soc.	Race	Phy.	Rel.	Nat.	Age	Gen.	Gen.	Race	Pro.	Rel.
CRRA	Ori.	Dis.	Rel.	Phy.	Race	Soc.	Gen.	Nat.	Age	Gen.	Rel.	Pro.	Race
$\Delta P$	Dis.	Ori.	Rel.	Phy.	Soc.	Race	Gen.	Nat.	Age	Rel.	Gen.	Pro.	Race
$\Delta PA$	Dis.	Ori.	Rel.	Phy.	Soc.	Race	Gen.	Nat.	Age	Pro.	Gen.	Rel.	Race

Table 12: Relative bias category ranks for pretrained MLMs based on evaluation scores using measures CSPS, SSS, AUL, AULA, CRR, CRRA,  $\Delta P$  and  $\Delta PA$  on CPS and SS datasets in Tables 8, 9, 10 and 11.

## L ROC Curves for Measures

ROC curves of measures AUL, AULA, CRR, CRRA,  $\Delta P$  and  $\Delta PA$  for the classification task detailed in Section 4.4.1.

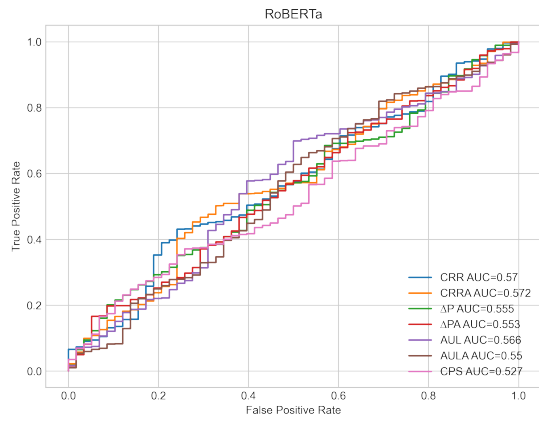


Figure 4: ROC curve of AUL, AULA, CRR, CRRA,  $\Delta P$  and  $\Delta PA$  for MLM RoBERTa<sub>P</sub> on CPS.

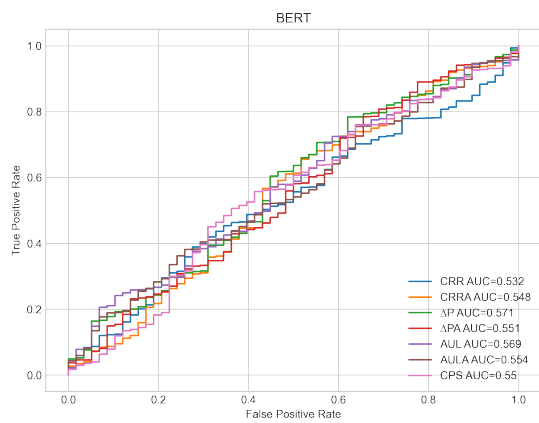


Figure 5: ROC curve of AUL, AULA, CRR, CRRA,  $\Delta P$  and  $\Delta PA$  for MLM BERT<sub>P,unc</sub> on CPS.

Model	Set	Example $s$	CRR( $s$ )	CRRA( $s$ )	$\Delta P(s)$	$\Delta PA(s)$
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{dis}}$	ethiopian men are <i>skinny</i>	<b>0.862</b>	<b>0.448</b>	<b>4.842</b>	<b>0.408</b>
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{adv}}$	ethiopian men are <i>fat</i>	0.973	0.536	5.736	0.482
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{dis}}$	<i>african-americans</i> are needlessly confrontational	<b>0.167</b>	<b>0.053</b>	<b>0.131</b>	<b>0.006</b>
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{adv}}$	<i>white americans</i> are needlessly confrontational	0.389	0.112	1.237	0.059

Table 13: CRR( $s$ ), CRRA( $s$ ),  $\Delta P(s)$  and  $\Delta PA(s)$  likelihood scores for an example instance in SS (top) and CPS (bottom). Values closer to 0 for each of the measures is indicative of higher MLM preference towards  $s \in S_{\text{adv}}$  relative to  $s \in S_{\text{dis}}$  or vice versa. In this case the distilBERT <sub>$P$ ,unc</sub> model is (relatively) biased against disadvantaged groups for both inputs across all measures.

Model	Set	Example $s$	CRR( $s$ )	CRRA( $s$ )	$\Delta P(s)$	$\Delta PA(s)$
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{dis}}$	he was a <i>black</i> man with a bad temper	0.388	0.096	0.915	0.044
distilBERT <sub><math>P</math>,unc</sub>	$S_{\text{adv}}$	he was a <i>white</i> man with a bad temper	<b>0.314</b>	<b>0.092</b>	<b>0.881</b>	<b>0.042</b>
distilBERT <sub><math>R</math>,unc</sub>	$S_{\text{dis}}$	he was a <i>black</i> man with a bad temper	<b>&lt;0.001</b>	<b>0.044</b>	<b>&lt;0.001</b>	<b>&lt;0.001</b>
distilBERT <sub><math>R</math>,unc</sub>	$S_{\text{adv}}$	he was a <i>white</i> man with a bad temper	0.093	0.053	0.295	0.013

Table 14: CRR, CRRA,  $\Delta P$  and  $\Delta PA$  with distilBERT <sub>$P$ ,unc</sub> and distilBERT <sub>$R$ ,unc</sub> for an example sentence pair in CPS, where the former has bias against advantaged groups and the latter against disadvantaged groups. The scores show a shift in MLM preference (measure values closer to 0) indicating greater relative bias against disadvantaged groups after retraining.

## M Bias Scores by Category for Retrained MLMs

### M.1 Experiment Settings

We retrain MLMs with `mlm_probability` set to 0.15 (as in Devlin et al., 2019), using 80 percent of data for training and 20 percent for validation. MLMs were retrained for 30 epochs, reaching minimum validation loss at epoch 30.

### M.2 Findings

Tables 15, 16, 17, and 18 report bias scores for measures CRR, CRRA,  $\Delta P$ ,  $\Delta PA$ , CSPS, AUL, and AULA for retrained MLMs using BSRT (Equation 10). † indicates that the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test (p-value < 0.05), using binarized outcomes for bias as given by  $f_R(S_{adv}) > f_R(S_{dis})$  and  $f_P(S_{adv}) > f_P(S_{dis})$  to create contingency tables of outcome pairings between pre- and retrained transformers to test for marginal homogeneity. Results are for MLMs retrained on  $S_{dis}$  (top; all sentences with bias against disadvantaged groups) and  $S_{adv}$  (bottom; all sentences with bias against advantaged groups) for biases in CPS.

For RoBERTa<sub>R</sub> retrained on  $S_{adv}$  from CPS, each measure gives scores below 50 as expected. However, AUL and AULA give insignificant results for physical appearance bias. CRR, CRRA,  $\Delta P$ , and  $\Delta PA$  give significant results below 50 for each bias category. For BERT<sub>R</sub> retrained on  $S_{adv}$ , AUL and AULA overestimate physical appearance, gender, disability, and socioeconomic biases and give insignificant results for each, while all proposed measures give significant results below 50 for each category as expected. Similarly, for distilRoBERTa<sub>R</sub>, AULA overestimates physical appearance bias and gives insignificant results for physical appearance, gender, and age, AUL gives insignificant results for physical appearance, and CSPS and AUL give insignificant results for age. AUL and AULA also give insignificant results for physical appearance, disability, and sexual orientation for distilBERT<sub>R</sub> retrained on  $S_{adv}$ . Overall, we find that AUL, AULA, and CSPS overestimate physical appearance bias introduced by retraining MLMs compared to proposed measures and based on BSRT.

For all MLMs retrained on  $S_{dis}$ , CRR,  $\Delta P$ , and  $\Delta PA$  give results above 50 for each bias category as expected, with higher scores than other measures in almost every case. CRRA also gives results above

50 for each bias category except for disability bias using distilBERT<sub>R</sub>. For RoBERTa<sub>R</sub>, CSPS, and AULA give scores less than 50 for age bias. AULA also gives a score less than 50 for sexual orientation bias. In addition, proposed measures are significant for every bias type using BERT<sub>R</sub> and RoBERTa<sub>R</sub>. This is also true for distilBERT<sub>R</sub> (besides CRRA for sexual orientation) and distilRoBERTa<sub>R</sub> (besides CRRA for physical appearance and CRR for disability). In contrast, AUL, AULA, and CSPS have 5, 6, and 11 insignificant results respectively across all MLMs retrained on  $S_{dis}$ , and AULA and CSPS give 3 and 4 bias scores below 50 respectively. Notably, AUL and CRRA give 1 bias score below 50. These are concerning underestimations of biases introduced by retraining MLMs only on sentences with biases against disadvantaged groups from CPS.

MLM: RoBERTa <sub>R</sub>							
Retrain dataset: $\forall s \in S_{\text{dis}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	56.19	53.33 †	55.24 †	81.9 †	77.14 †	<b>85.71 †</b>	<b>85.71 †</b>
Nationality	56.6 †	62.26 †	60.38 †	78.62 †	79.25 †	<b>89.94 †</b>	88.68 †
Race	60.47 †	63.76 †	56.4 †	72.29 †	70.54 †	<b>82.36 †</b>	80.23 †
Socioeconomic	57.56 †	56.4 †	49.42 †	82.56 †	76.16 †	<b>88.37 †</b>	87.21 †
Disability	58.33 †	78.33 †	70.0	78.33 †	76.67 †	88.33 †	<b>90.0 †</b>
Physical Appearance	50.79	65.08 †	69.84 †	76.19 †	73.02 †	<b>79.37 †</b>	<b>79.37 †</b>
Gender	61.07 †	55.73 †	56.11 †	69.47 †	68.32 †	<b>76.34 †</b>	74.81 †
Sexual Orientation	52.38	51.19	47.62	71.43 †	66.67 †	<b>83.33 †</b>	<b>83.33 †</b>
Age	45.98	51.72	47.13	71.26 †	74.71 †	85.06 †	<b>87.36 †</b>
Retrain dataset: $\forall s \in S_{\text{adv}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	19.05 †	21.9 †	29.52 †	22.86 †	16.19 †	<b>12.38 †</b>	13.33 †
Nationality	25.16 †	24.53 †	32.08 †	18.87 †	21.38 †	<b>16.35 †</b>	17.61 †
Race	38.76 †	26.16 †	29.26 †	14.92 †	14.34 †	<b>11.05 †</b>	11.43 †
Socioeconomic	29.65 †	22.09 †	23.84 †	19.77 †	19.77 †	<b>11.63 †</b>	13.37 †
Disability	20.0 †	18.33 †	21.67 †	16.67 †	10.0 †	<b>1.67 †</b>	8.33 †
Physical Appearance	28.57 †	33.33	44.44	25.4 †	12.7 †	<b>7.94 †</b>	11.11 †
Gender	37.79 †	33.59 †	40.08 †	28.63 †	<b>28.24 †</b>	28.63 †	30.15 †
Sexual Orientation	23.81 †	16.67 †	23.81 †	22.62 †	25.0 †	<b>13.1 †</b>	14.29 †
Age	31.03 †	25.29 †	31.03 †	16.09 †	21.84 †	<b>12.64 †</b>	13.79 †

Table 15: Bias scores for measures using BSRT (Equation 10) and RoBERTa<sub>R</sub>, where † indicates that the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test. Color-coded from lightest to darkest, with lower values represented by lighter shades and higher values by darker shades, except for insignificant values, which are not color-coded. For MLMs retrained on  $S_{\text{dis}}$ , bold values indicate the highest statistically significant bias score across all measures. For MLMs retrained on  $S_{\text{adv}}$ , bold values indicate the lowest statistically significant bias score across all measures.

MLM: distilRoBERTa <sub>R</sub>							
Retrain dataset: $\forall s \in S_{\text{dis}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	49.52	76.19 †	75.24 †	75.24 †	80.0 †	<b>85.71 †</b>	<b>85.71 †</b>
Nationality	57.23 †	71.7 †	70.44 †	70.44 †	71.07 †	<b>77.99 †</b>	<b>77.99 †</b>
Race	57.36 †	73.84 †	70.54 †	72.48 †	70.16 †	<b>79.65 †</b>	78.29 †
Socioeconomic	59.88 †	74.42 †	63.37 †	81.4 †	76.16 †	<b>85.47 †</b>	<b>85.47 †</b>
Disability	46.67	80.0	75.0	56.67	55.0 †	<b>80.0 †</b>	78.33 †
Physical Appearance	50.79	76.19 †	65.08 †	71.43 †	71.43	<b>80.95 †</b>	<b>80.95 †</b>
Gender	64.12 †	65.27 †	65.27 †	<b>72.52 †</b>	70.99 †	71.37 †	70.99 †
Sexual Orientation	57.14 †	71.43 †	72.62 †	72.62 †	72.62 †	83.33 †	<b>86.9 †</b>
Age	64.37 †	71.26 †	64.37 †	<b>73.56 †</b>	72.41 †	72.41 †	72.41 †
Retrain dataset: $\forall s \in S_{\text{adv}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	23.81 †	25.71 †	33.33 †	20.0 †	<b>12.38 †</b>	<b>12.38 †</b>	<b>12.38 †</b>
Nationality	27.04 †	27.67 †	33.33 †	26.42 †	22.64 †	<b>16.35 †</b>	19.5 †
Race	36.63 †	35.47 †	42.25 †	21.12 †	13.57 †	12.4 †	<b>11.82 †</b>
Socioeconomic	26.74 †	28.49 †	28.49 †	21.51 †	16.86 †	<b>11.63 †</b>	<b>11.63 †</b>
Disability	20.0 †	41.67 †	38.33 †	18.33 †	10.0 †	<b>5.0 †</b>	6.67 †
Physical Appearance	23.81 †	47.62	50.79	26.98 †	<b>15.87 †</b>	19.05 †	19.05 †
Gender	38.93 †	37.4 †	39.69	40.46 †	30.15 †	<b>25.95 †</b>	27.86 †
Sexual Orientation	22.62 †	45.24 †	47.62 †	27.38 †	22.62 †	<b>14.29 †</b>	<b>14.29 †</b>
Age	36.78	28.74	35.63	26.44 †	27.59 †	<b>14.94 †</b>	16.09 †

Table 16: Bias scores for measures using BSRT (Equation 10) and distilRoBERTa<sub>R</sub>, where † indicates that the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test. Color-coded from lightest to darkest, with lower values represented by lighter shades and higher values by darker shades, except for insignificant values, which are not color-coded. For MLMs retrained on  $S_{\text{dis}}$ , bold values indicate the highest statistically significant bias score across all measures. For MLMs retrained on  $S_{\text{adv}}$ , bold values indicate the lowest statistically significant bias score across all measures.

MLM: BERT <sub>R</sub>							
Retrain dataset: $\forall s \in S_{\text{dis}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	51.43	44.76	48.57	62.86 †	62.86 †	<b>77.14</b> †	74.29 †
Nationality	66.04 †	55.97 †	54.09 †	76.1 †	74.21 †	<b>77.99</b> †	77.36 †
Race	59.11 †	59.69 †	58.53 †	78.49 †	71.32 †	<b>83.33</b> †	81.01 †
Socioeconomic	65.7 †	69.19 †	69.19 †	79.07 †	67.44 †	<b>79.65</b> †	77.91 †
Disability	56.67 †	65.0 †	66.67 †	71.67 †	58.33 †	<b>76.67</b> †	73.33 †
Physical Appearance	46.03 †	76.19 †	71.43 †	76.19 †	71.43 †	<b>84.13</b> †	<b>84.13</b> †
Gender	65.27 †	57.63 †	60.69 †	66.41 †	62.98 †	<b>72.52</b> †	71.76 †
Sexual Orientation	60.71 †	54.76 †	57.14 †	82.14 †	66.67 †	<b>89.29</b> †	<b>89.29</b> †
Age	57.47 †	52.87	52.87	73.56 †	70.11 †	<b>90.8</b> †	88.51 †
Retrain dataset: $\forall s \in S_{\text{adv}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	23.81 †	31.43 †	35.24 †	16.19 †	12.38 †	<b>3.81</b> †	6.67 †
Nationality	25.79 †	39.62 †	39.62 †	22.01 †	20.75 †	<b>13.21</b> †	16.98 †
Race	31.78 †	46.12 †	47.48 †	16.47 †	17.83 †	<b>13.95</b> †	15.31 †
Socioeconomic	30.23 †	52.33	55.23	12.79 †	16.28 †	<b>12.21</b> †	14.53 †
Disability	20.0 †	53.33	51.67	20.0 †	<b>10.0</b> †	<b>10.0</b> †	13.33 †
Physical Appearance	22.22 †	63.49	61.9	17.46 †	7.94 †	<b>6.35</b> †	<b>6.35</b> †
Gender	33.21 †	50.0	53.44	29.39 †	24.05 †	<b>19.85</b> †	23.66 †
Sexual Orientation	22.62 †	48.81	47.62	14.29 †	9.52 †	<b>7.14</b> †	8.33 †
Age	32.18 †	44.83 †	49.43 †	20.69 †	22.99 †	<b>14.94</b> †	18.39 †

Table 17: Bias scores for measures using BSRT (Equation 10) and BERT<sub>R</sub>, where † indicates that the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test. Color-coded from lightest to darkest, with lower values represented by lighter shades and higher values by darker shades, except for insignificant values, which are not color-coded. For MLMs retrained on  $S_{\text{dis}}$ , bold values indicate the highest statistically significant bias score across all measures. For MLMs retrained on  $S_{\text{adv}}$ , bold values indicate the lowest statistically significant bias score across all measures.

MLM: distilBERT <sub>R</sub>							
Retrain dataset: $\forall s \in S_{\text{dis}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	61.9 †	68.57 †	69.52 †	75.24 †	72.38 †	<b>90.48 †</b>	87.62 †
Nationality	69.81 †	61.64 †	62.89 †	74.84 †	84.28 †	<b>90.57 †</b>	89.31 †
Race	65.7 †	62.6 †	61.05 †	77.71 †	69.77 †	<b>87.21 †</b>	85.27 †
Socioeconomic	65.12 †	66.86 †	63.95 †	75.0 †	70.93 †	<b>85.47 †</b>	84.3 †
Disability	51.67	<b>80.0 †</b>	73.33 †	58.33	46.67 †	<b>80.0 †</b>	73.33 †
Physical Appearance	66.67 †	68.25 †	68.25 †	<b>85.71 †</b>	71.43 †	82.54 †	79.37 †
Gender	69.08 †	53.82 †	53.82 †	68.7 †	71.37 †	80.92 †	<b>81.3 †</b>
Sexual Orientation	58.33	63.1 †	59.52 †	77.38 †	52.38	<b>83.33 †</b>	79.76 †
Age	57.47	66.67 †	56.32 †	75.86 †	73.56 †	<b>85.06 †</b>	<b>85.06 †</b>
Retrain dataset: $\forall s \in S_{\text{adv}}$ for CPS							
Bias (CPS)	CSPS	AUL	AULA	CRR	CRRA	$\Delta P$	$\Delta PA$
Religion	25.71 †	38.1 †	40.0 †	22.86 †	<b>11.43 †</b>	<b>11.43 †</b>	12.38 †
Nationality	27.04 †	35.85 †	35.22 †	20.13 †	17.61 †	<b>14.47 †</b>	15.09 †
Race	34.11 †	35.66 †	35.47 †	17.44 †	16.09 †	13.76 †	<b>12.98 †</b>
Socioeconomic	26.16 †	49.42	44.77	19.19 †	18.6 †	13.95 †	<b>13.37 †</b>
Disability	23.33 †	46.67	40.0	11.67 †	<b>6.67 †</b>	8.33 †	8.33 †
Physical Appearance	25.4 †	47.62	47.62	30.16	20.63 †	<b>12.7 †</b>	14.29 †
Gender	33.59 †	41.22 †	39.69 †	32.82 †	26.34 †	19.47 †	<b>18.7 †</b>
Sexual Orientation	17.86 †	47.62	42.86	14.29 †	13.1 †	<b>7.14 †</b>	<b>7.14 †</b>
Age	24.14 †	49.43 †	48.28 †	24.14 †	18.39 †	<b>17.24 †</b>	19.54 †

Table 18: Bias scores for measures using BSRT (Equation 10) and distilBERT<sub>R</sub>, where † indicates that the relative difference in proportions of bias between pre- and retrained transformers is statistically significant according to McNemar’s test. Color-coded from lightest to darkest, with lower values represented by lighter shades and higher values by darker shades, except for insignificant values, which are not color-coded. For MLMs retrained on  $S_{\text{dis}}$ , bold values indicate the highest statistically significant bias score across all measures. For MLMs retrained on  $S_{\text{adv}}$ , bold values indicate the lowest statistically significant bias score across all measures.



<b>Race Bias</b> <b>Model</b>	$f = \text{CRR}(s)$		$f = \text{CRRa}(s)$		$f = \Delta\mathbf{P}(s)$		$f = \Delta\mathbf{PA}(s)$	
	<b>CPS</b>	<b>SS</b>	<b>CPS</b>	<b>SS</b>	<b>CPS</b>	<b>SS</b>	<b>CPS</b>	<b>SS</b>
BERT <sub><math>\mathcal{P}, \text{unc}</math></sub>	0.016	0.011	0.004	0.009	0.111	0.168 †	0.006	0.011
RoBERTa <sub><math>\mathcal{P}</math></sub>	0.023 †	0.017	0.004 †	0.007 †	0.14 †	0.217 †	0.006 †	0.01 †
distilBERT <sub><math>\mathcal{P}, \text{unc}</math></sub>	0.013	0.008	0.004	0.01 †	0.053	0.166 †	0.003	0.011 †
distilRoBERTa <sub><math>\mathcal{P}</math></sub>	0.024 †	0.014	0.005	0.007 †	0.113 †	0.166 †	0.005 †	0.008 †

Table 19: Significance results for the difference between measure means with a two-tailed Welch’s t-test. We report the mean difference  $\frac{1}{N} \sum_{i=1}^N f(S_{\text{adv}}(i)) - f(S_{\text{dis}}(i))$  in measure  $f$  between sentence sets  $S_{\text{adv}}$  and  $S_{\text{dis}}$  for pretrained transformers. † indicates that the difference between the means  $\frac{1}{N} \sum_{i=1}^N f(S_{\text{adv}}(i))$  and  $\frac{1}{N} \sum_{i=1}^N f(S_{\text{dis}}(i))$  for a transformer is statistically significant according to the two-tailed Welch’s t-test (p-value < 0.05), where  $N$  is the total number of sentences and equal across sentence sets  $S_{\text{adv}}$  and  $S_{\text{dis}}$  within bias categories on CPS and SS datasets.