

# How Gender Debiasing Affects Internal Model Representations, and Why It Matters

Anonymous ACL submission

## Abstract

Common studies of gender bias in NLP focus either on extrinsic bias measured by model performance on a downstream task or on intrinsic bias found in models' internal representations. However, the relationship between extrinsic and intrinsic bias is relatively unknown. In this work, we illuminate this relationship by measuring both quantities together: we debias a model during downstream fine-tuning, which reduces extrinsic bias, and measure the effect on intrinsic bias, which is operationalized as bias extractability with information-theoretic probing. Through experiments on two tasks and multiple bias metrics, we show that our intrinsic bias metric is a better indicator of debiasing than (a contextual adaptation of) the standard WEAT metric, and can also expose cases of superficial debiasing. Our framework provides a comprehensive perspective on bias in NLP models, which can be applied to deploy NLP systems in a more informed manner. Our code will be made publicly available.

## 1 Introduction

Efforts to identify and mitigate gender bias in Natural Language Processing (NLP) systems typically target one of two notions of bias. *Extrinsic* evaluation methods and debiasing techniques focus on the bias reflected in a downstream task (De-Arteaga et al., 2019; Zhao et al., 2018), while *intrinsic* methods focus on a model's internal representations, such as word or sentence embedding geometry (Caliskan et al., 2017; Bolukbasi et al., 2016; Guo and Caliskan, 2021). Despite an abundance of evidence pointing towards gender bias in pre-trained language models (LMs), the extent of harm caused by these biases is not clear when it is not reflected in a specific downstream task (Barocas et al., 2017; Kate Crawford, 2017; Blodgett et al., 2020; Bommasani et al., 2021). For instance, while the word embedding proximity of "doctor" to "man" and "nurse" to "woman" is intuitively normatively

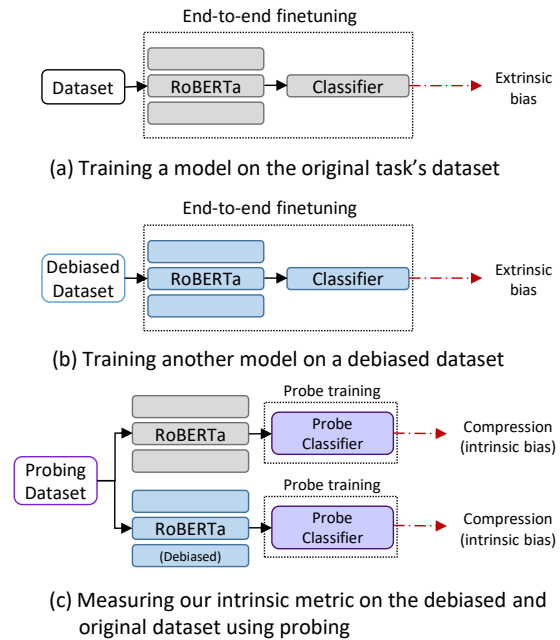


Figure 1: Our proposed framework. Black arrows mark forward passes, red arrows mark things we measure. We first (a) train a model on a downstream task, then (b) train another model on the same task using a debiased dataset, and finally (c) measure intrinsic bias in both models and compare.

wrong, it is not clear when such phenomena would lead to downstream predictions manifesting in social biases. Recently, Goldfarb-Tarrant et al. (2021) have shown that debiasing static embeddings intrinsically is not correlated with extrinsic gender bias measures, but the nature of the reverse relationship is unknown: how are extrinsic interventions reflected in intrinsic representations? Furthermore, Gonen and Goldberg (2019a) demonstrated that a number of intrinsic debiasing methods applied to static embeddings only partially remove the bias and that most of it is still hidden within the embedding. Complementing their view, we examine *extrinsic* debiasing methods, as well as demonstrate the possible harm this could cause. Contrary to their conclusion, we do not claim that these debias-

ing methods should not be trusted, *as long as they are utilized with care.*

Our goal is to gain a better understanding of the relationship between a model’s internal representations and its extrinsic gender bias by examining the effects of various debiasing methods on the model’s representations. Specifically, we fine-tune models with and without gender debiasing strategies, evaluate their external bias using various bias metrics, and measure intrinsic bias in the representations. We operationalize intrinsic bias via two metrics: First, we use CEAT (Guo and Caliskan, 2021), a contextual adaptation of the widely used intrinsic bias metric WEAT (Caliskan et al., 2017). Second, we propose to use an information-theoretic probe to quantify the degree to which gender can be extracted from the internal model representations. Then, we examine how these intrinsic metrics correlate with a variety of extrinsic bias metrics that we measure on the model’s downstream performance. Our approach is visualised in Figure 1.

We perform extensive experiments on two downstream tasks (occupation prediction and coreference resolution); several debiasing strategies that involve alterations to the training dataset (such as removing names and gender indicators, or balancing the data by oversampling or downsampling); and a multitude of extrinsic bias metrics. Our analysis reveals new insights into the way language models encode and use information on gender:

- The effect of debiasing on internal representations is reflected in gender extractability, while not always in CEAT.
- In cases of high gender extractability but low extrinsic bias metrics, the debiasing is superficial, and the internal representations are a good indicator for this: The bias is still present in internal representations and can be restored by retraining the classification layer.
- The two tasks show different patterns of correlation between intrinsic and extrinsic bias. The coreference task exhibits a high correlation. The occupation prediction task exhibits a lower correlation, but it increases after retraining (a case of superficial debiasing). Gender extractability shows higher correlations to extrinsic metrics than CEAT.

## 2 Methodology

In this study, we investigate the relationship between extrinsic bias metrics of a task and a model’s internal representations, under various debiasing conditions, for two datasets in English. We perform extrinsic debiasing, evaluate various extrinsic and intrinsic bias metrics before and after debiasing, and examine correlations.

**Dataset.** Let  $D = \{\mathcal{X}, \mathcal{Y}, \mathcal{Z}\}$  be a dataset consisting of input data  $\mathcal{X}$ , labels  $\mathcal{Y}$  and protected attributes  $\mathcal{Z}$ .<sup>1</sup> This work focuses on gender as the protected attribute  $Z$ . In all definitions,  $F$  and  $M$  indicate female and male gender, respectively, as the value of the protected attribute  $Z$ .

**Trained Model.** The model is optimized to solve the downstream task posed by the dataset. It can be formalized as  $f(g(x)) \mapsto \mathbb{R}^{|\mathcal{Y}|}$ , where  $g(\cdot)$  is the feature extractor, implemented by a language model, e.g., RoBERTa (Liu et al., 2019), and  $f(\cdot)$  is the classification function.

### 2.1 Bias Metrics

Each bias evaluation method described in the literature can be categorized as extrinsic or intrinsic. In all definitions,  $\mathcal{R}$  indicates the model’s predictions.

#### 2.1.1 Extrinsic Metrics

Extrinsic methods involve measuring the bias of a model solving a downstream problem. The extrinsic metric is a mapping:

$$E(\mathcal{X}, \mathcal{Y}, \mathcal{R}, \mathcal{Z}) \mapsto \mathbb{R}$$

The output represents the quantity of bias measured; the further from 0 the number is, the larger the bias is. Our analysis comprises a wide range of extrinsic metrics, including some that have been measured in the past on the analyzed tasks (Zhao et al., 2018; De-Arteaga et al., 2019; Ravfogel et al., 2020; Goldfarb-Tarrant et al., 2021) and some that have never been measured before, and shows our results apply to many of them. For illustration, we will consider occupation prediction, a common task in research on gender bias (De-Arteaga et al., 2019; Ravfogel et al., 2020; Romanov et al., 2019). The input  $X$  is a biography and the prediction  $Y$  is the profession of the person described in it. The protected attribute  $Z$  is the gender of that person.

<sup>1</sup> $\mathcal{Z}$  is by convention used for attributes for which we want to ensure fairness, such as gender, race, etc. It is purposefully broad, and depending on the task and data could refer to the gender of an entity in coreference, the subject of a text, the demographics of the author of a text, etc.

**Performance gap.** This is the difference in performance metric for two different groups, for instance two groups of binary genders, or a group of pro-stereotypical and a group of anti-stereotypical examples. We measure the following metrics: True Positive Rate (TPR), False Positive Rate (FPR), and Precision. In occupation prediction, for instance, the TPR gap for each profession  $y$  expresses the difference in the percentage of women and men whose profession is  $y$  and are correctly classified as such. We also measure F1 of three standard clustering metrics for coreference resolution.

We compute two types of performance gap metrics: (1) the sum of absolute gap values over all classes; (2) the Pearson correlation between the performance gap for a class and the percentage of women in that class. For instance, if  $y$  is a profession, we measure the correlation between performance gaps and percentages of women in each profession.<sup>2</sup> The two metrics are closely related but answer slightly different questions: the sum quantifies how a model behaves differently on different genders, and the correlation shows the relation of model behaviour to social biases (in the world or the data) without regard to actual gap size.

**Statistical metrics.** For breadth of analysis, we examine three additional statistical metrics (Barocas et al., 2019), which correspond to different notions of bias. All three are measured as differences ( $d$ ) between two probability distributions, and we then obtain a single bias quantity per metric by summing all computed distances.

- *Independence:*  $d(P(R|Z = z), P(R)) \forall z \in \{F, M\}$ . For instance, we measure the difference between the distribution of model’s predictions on women and the distribution of all predictions.
- *Separation:*  $d(P(R|Y = y, Z = z), P(R|Y = y)) \forall y \in \mathcal{Y}, z \in \{F, M\}$ . For instance, we measure the difference between the distribution of a model’s predictions on women who are teachers and the distribution of predictions on all teachers.
- *Sufficiency:*  $d(P(Y|R = r, Z = z), P(Y|R = r))$ . For instance, we measure the difference between the distribution of gold labels on women classified as teachers by the model and the distribution of gold labels on all individuals classified as teachers by the model.

<sup>2</sup>Percentages for coreference resolution are taken from labour statistics, following Zhao et al. (2018). For occupation prediction we use training set statistics following De-Arteaga et al. (2019), before balancing.

## 2.1.2 Intrinsic Metrics

Intrinsic methods are applied to the representation obtained from the feature extractor. These methods are independent of any downstream task. The intrinsic metric is a mapping:

$$I(g(X), Z) \mapsto \mathbb{R}$$

**Compression.** Our main intrinsic metric is the *compression* of gender information evaluated by a minimum description length (MDL) probing classifier (Voita and Titov, 2020), trained to predict gender from the model’s representations. Probing classifiers are widely used for predicting various properties of interest from frozen model representations (Belinkov and Glass, 2019). MDL probes were proposed because a probe’s accuracy may be misleading due to memorization and other issues (Hewitt and Liang, 2019; Belinkov, 2021). We use the MDL online code, where the probe is trained in timesteps, on increasing subsets of the training set, then evaluated against the rest of it. Higher compression indicates greater gender extractability.

**CEAT.** We also measure CEAT (Guo and Caliskan, 2021), which is a contextualized version of WEAT (Caliskan et al., 2017), a widely used bias metric for static word embeddings. WEAT defines sets  $X$  and  $Y$  of target words, and sets  $A$  and  $B$  of attribute words. For instance,  $A, B$  contain males and females names, while  $X, Y$  contain career and family related words, respectively. The bias is operationalized as the geometric proximity between the target and attribute word embeddings, and is quantified in CEAT by the Combined Effect Size (CES) and a p-value for the null hypothesis of having no biased associations. For more information on CEAT refer to Appendix A.4.3.

## 2.2 Debiasing Techniques

We debias models by modifying the downstream task’s training data before fine-tuning. *Scrubbing* (De-Arteaga et al., 2019) removes first names and gender-specific terms (“he”, “she”, “husband”, “wife”, “Mr”, “Mrs”, etc.). *Balancing* subsamples or oversamples examples such that each gender is equally represented in the resulting dataset w.r.t each label. *Anonymization* (Zhao et al., 2018) removes named entities. *Counterfactual Augmentation* (Zhao et al., 2018) involves replacing male entities in an example with female entities, and adding the modified example to the training set.

As some of these are dataset/task-specific, we give more details in the following section.

### 3 Experiments

In each experiment, we fine-tune a model for a downstream task. For training, we use either the original dataset or a dataset debiased with one of the methods from Section 2.2. Figure 2 presents examples of debiasing methods for the two downstream tasks. We measure two intrinsic metrics by probing that model’s inner representations for gender extractability (as measured by MDL) and by CEAT, and test various extrinsic metrics. The relation between one intrinsic and one extrinsic metric becomes one data point, and we repeat over many random seeds (for both the model and the probe). Further implementation details are in appendix A.

#### 3.1 Occupation Prediction

The task of occupation prediction is to predict a person’s occupations (from a closed set), based on their biography. We use the Bias in Bios dataset (De-Arteaga et al., 2019). Regardless of the training method, the test set is subsampled such that each profession has equal gender representation.

**Model.** Our model is a RoBERTa model (Liu et al., 2019) topped with a linear classifier, which receives the [CLS] token embedding as input and generates a probability distribution over the professions. In addition, we train a baseline classifier layer on top of a frozen, non-finetuned RoBERTa.

**Debiasing Techniques.** Following De-Arteaga et al. (2019) we experiment with scrubbing the training dataset. Figure 2 shows an example biography snippet and its scrubbed version. We also conduct balancing (per profession, subsampling and oversampling to ensure an equal number of males and females per profession), which has not previously been used on this dataset and task.

**Metrics.** We measure all bias metrics from Section 2.1 except for F1.

**Probing.** The probing dataset for this task is the test set, and the gender label of a single biography is the gender of the person described in it. We probe the [CLS] token representation of the biography. In addition to the models described above, we measure baseline extractability of gender information from a randomly initialized RoBERTa model.

#### 3.2 Coreference Resolution

The task of coreference resolution is to find all textual expressions referring to the same real-world

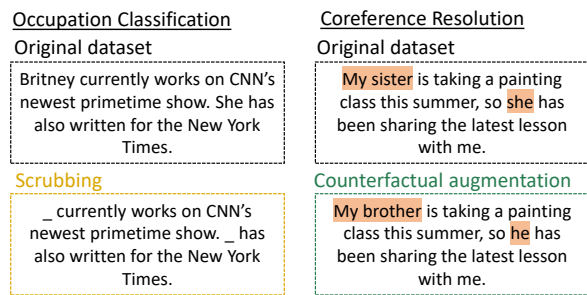


Figure 2: Examples of two debiasing methods performed on the data.

entities. We train on Ontonotes 5.0 (Weischedel et al., 2013) and test on the Winobias challenge dataset (Zhao et al., 2018). Winobias consists of sentence pairs, pro- and anti-stereotypical variants, with individuals referred to by their profession. For example, “The physician hired the secretary because *he/she* was busy.” is pro/anti-stereotypical, based on US labor statistics.<sup>3</sup> A coreference system is measured by the performance gap between the pro- and anti-stereotypical subsets.

**Model.** We use the model presented in Lee et al. (2018a) with RoBERTa as a feature extractor.

**Debiasing Techniques.** Following Zhao et al. (2018), we apply anonymization (denoted as Anon) and counterfactual augmentation (CA) on the training set. These techniques were used jointly in previous work; we examine each individually as well.

**Metrics.** Following Zhao et al. (2018), we measure the F1 difference between anti- and pro-stereotypical examples.<sup>4</sup> We also interpret the task as a classification problem, and measure all metrics from Section 2.1. For more details refer to Appendix A.4.2.

**Probing.** We probe the representation of a profession word as extracted from Winobias sentences, after masking out the pronouns. We define a profession’s gender as the stereotypical gender for this profession. To prevent memorization by the probe—given the small number of professions—the dataset is sorted so that professions are gradually added to the training set, so a success on the validation set is on previously unseen professions.

<sup>3</sup>Labor Force Statistics from the Current Population Survey, <https://www.bls.gov/cps/cpsaat11.htm>

<sup>4</sup>We combined the T1 and T2 datasets, as well as the dev and test datasets, to create a single held-out challenge dataset.

Debiasing Strategy	Extrinsic									
	Intrinsic		Before				After			
	Compression	CEAT	TPR (P)	FPR (S)	Sep	Suff	TPR (P)	FPR (S)	Sep	Suff
Random	5.61*	0.12†	-	-	-	-	-	-	-	-
Pre-trained	10.12	0.49*	-	-	-	-	-	-	-	-
None	4.12	0.22	0.76	0.08	0.33	9.45	0.78	0.073	0.33	9.70
Oversampling	8.52*	0.29	0.73	0.09*	0.31	8.32*	0.81*	0.068*	0.33	10.91*
Subsampling	3.57	0.22	<b>0.32*</b>	<b>0.03*</b>	<b>0.20*</b>	<b>1.22*</b>	<b>8.37*</b>	0.08*	0.30*	1.32*
Scrubbing	<b>1.70*</b>	0.23	0.70*	0.06*	0.30	4.93*	0.71*	<b>0.06*</b>	<b>2.56*</b>	<b>0.81*</b>

(a) Occupation classification: Comparison of intrinsic and extrinsic metrics before and after retraining of classification layer, over 10 seeds per fine-tuned model and per retrained classification model.

Debiasing Strategy	Extrinsic								
	Intrinsic		Before				After		
	Compression	CEAT	F1 diff	FPR (S)	Sep	Suff	F1 diff	FPR (S)	Sep
Random	0.83*	0.12†	-	-	-	-	-	-	-
Pre-trained	0.96	0.49*	-	-	-	-	-	-	-
None	1.98	0.35	6.63	0.12	1.25	8.69	6.07	0.11	1.19
Anon	2.07*	0.31*	7.26	0.13	1.34	8.82	7.42*	0.13*	1.34*
CA	<b>1.50*</b>	0.27*	<b>2.30*</b>	0.05*	<b>0.54*</b>	1.67*	3.67*	0.06*	0.67*
Anon + CA	1.54*	<b>0.25*</b>	2.42*	<b>0.049*</b>	0.56*	<b>1.56*</b>	<b>2.86*</b>	<b>0.05*</b>	<b>0.59*</b>

(b) Coreference resolution: Comparison of intrinsic and extrinsic metrics before and after retraining of classification layer, over 10 seeds per fine-tuned model and 5 seeds per retrained classification model.

Table 1: Results on both tasks. \* marks significant reduction or increase in bias ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy None). The lowest bias score in each column is marked with **bold**. P = Pearson; S = Sum. † was computed only on 3 out of 10 tests for which CEAT’s  $p < 0.05$ .

## 4 Results

Tables 1a and 1b present intrinsic and extrinsic metrics on the occupation prediction and coreference resolution tasks, respectively. We present a representative subset of the measured metrics that demonstrate the observed phenomena; full results are found in Appendix B.

### 4.1 Compression Reflects Debiasing Effects

As shown in the tables, compression captures differences in models that were debiased differently. CEAT, however, cannot differentiate between occupation prediction models. For example, in occupation prediction (Table 1a) the compression rate varies significantly between a non-debiased and a debiased model via scrubbing and oversampling, while CEAT detects no difference between the models. In coreference resolution (Table 1b), both compression and CEAT are able to identify differences between the non-debiased model and the others, such as CA, which has both a lower compression and CEAT effect. But the CEAT effect sizes are small (below 0.5), which implies no bias in contrast to the extrinsic metrics.

### 4.2 High Gender Extractability Implies Superficial Debiasing

**Extrinsic and intrinsic effects of debiasing.** In occupation classification (Table 1a), somewhat surprisingly, subsampling the training data has the strongest effect on extrinsic metrics, but not on compression rate. Scrubbing reduces both intrinsic and extrinsic metrics, although its effect on extrinsic metrics is limited compared to subsampling. Training with oversampling caused less reduction in extrinsic bias metrics. A consequence of oversampling is that some metrics are less biased, but compression rates are increased, so gender information is more accessible. The effectiveness of subsampling over other metrics is further discussed in appendix C. In coreference resolution (Table 1b), while both CA and CA with anonymization reduced gender extractability as well as external bias metrics, anonymization alone *increased* intrinsic bias without affecting external bias metrics significantly.

**Debiasing without fine-tuning.** As the effect on extrinsic bias did not match the effect on intrinsic bias in several cases, we examined the role of the classification layer. We trained a model for occupa-

Metric	Occupation Classification				Coreference Resolution			
	$R^2$ Compression		$R^2$ CEAT		$R^2$ Compression		$R^2$ CEAT	
	Before	After	Before	After	Before	After	Before	After
F1 diff ( <i>pro</i> – <i>anti</i> )	-	-	-	-	0.821	0.709	0.246	0.005
TPR gap (P)	0.046	0.304	0.042	0.049	0.222	0.006	0.008	0.012
TPR gap (S)	0.049	0.449	0.022	0.036	0.817	0.752	0.297	0.003
FPR gap (P)	0.001	0.120	0.008	0.002	0.021	0.054	0.002	0.000
FPR gap (S)	0.353	0.046	0.079	0.001	0.844	0.773	0.263	0.004
Precision gap (P)	0.004	0.063	0.006	0.002	0.223	0.008	0.009	0.013
Precision gap (S)	0.150	0.291	0.031	0.054	0.817	0.752	0.296	0.003
Independence gap (S)	0.251	0.382	0.050	0.005	0.778	0.732	0.355	0.001
Separation gap (S)	0.066	0.165	0.046	0.009	0.835	0.776	0.261	0.005
Sufficiency gap (S)	0.202	0.567	0.040	0.034	0.825	0.753	0.287	0.002

Table 2: Coefficient determination of the regression line taken on the compression rate and each metric, before and after retraining of the classification layer. P = Pearson; S = Sum.

tion prediction without fine-tuning the underlying RoBERTa model. Training on a subsampled dataset also reduced the extrinsic metrics (0.15, 0.03, 0.20, and 0.31, respectively, on TPR gaps Pearson, FPR gaps sum, separation sum, and sufficiency sum). Detailed results of this experiment can be found in Appendix B. Since no updates were made to the LM, the internal representations could not be debiased, thus the debiasing observed in this model can only be superficial.

**Retraining the classification layer.** Fine-tuning of both tasks revealed that lower extrinsic metrics did not always lead to lower compression. Does this indicate cases where the debiasing process is only superficial, and the internal representations remain biased? To test this hypothesis, we froze the previously fine-tuned LM’s weights, and retrained the classification layer. We used the original (non-debiased) training set for retraining.

Tables 1a and 1b also compare extrinsic metrics before and after retraining. All models show bias restoration, due to the classification layer being trained on the biased dataset.<sup>5</sup> The amount of bias restored varies between models in a way that is predictable by the compression metric.

In the occupation prediction task, comparing Before and After numbers in Table 1a, the model fine-tuned using a scrubbed dataset—which has the

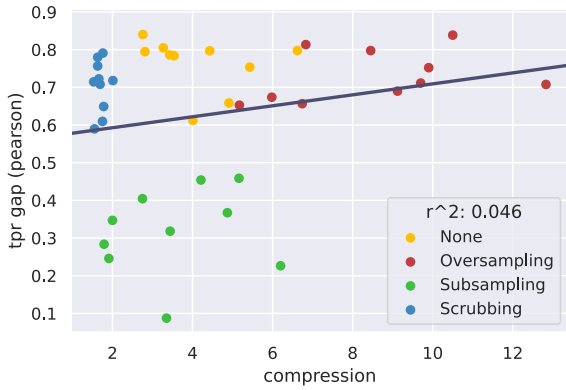
<sup>5</sup>The training datasets contain bias. The occupation prediction set has an unbalanced amount of males and females per profession (for example 15% of software engineers are females). The coreference resolution training set has more male than female pronouns, and males are more likely to be referred to by their profession (Zhao et al., 2018).

lowest compression rate—displays the least bias restoration, confirming that the LM absorbed the process of debiasing. The model fine-tuned on subsampled data has higher extrinsic bias after retraining. Hence, the debiasing was primarily cosmetic, and the representations within the LM were not debiased. The model fine-tuned on oversampled data—which has the highest compression—has the highest extrinsic bias (except for FPR), even though this was not true before retraining.

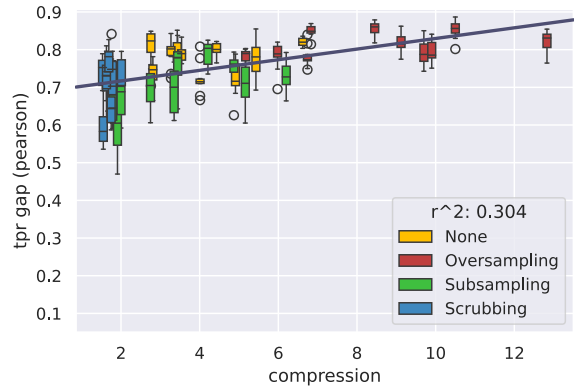
In coreference resolution, comparing Before and After numbers in Table 1b, models with the least extrinsic bias (CA and CA+Anon) are also least biased after retraining. Compression rate predicted this; these models also had lower compression rates than non-debiased models. Interestingly, the model fine-tuned with an anonymized dataset is the most biased after retraining, consistent with its high compression rate relative to the other models. As with subsampling and oversampling in occupation prediction, anonymization’s (lack of) effect on extrinsic metrics was cosmetic (compare None and Anon in Before block, Table 1b). Anonymization actually had a biasing effect on the LM, which was realized after retraining. We conclude that compression rate is a useful indicator of superficial debiasing.

### 4.3 Correlation between Extrinsic and Intrinsic Metrics

Table 2 shows correlations between compression rate and various extrinsic metrics before and after retraining. In occupation prediction, certain extrinsic metrics have a weak correlation with compres-

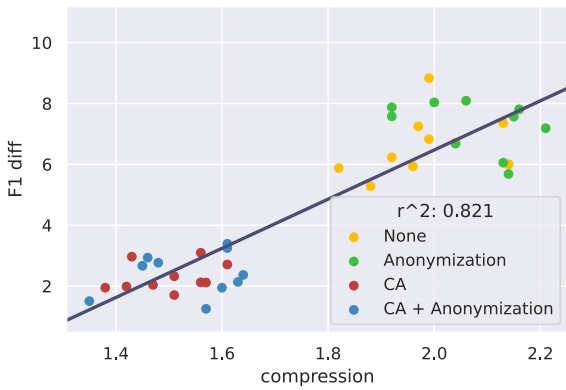


(a) Fine-tuned models. Each point is a single seed for training and testing the model.

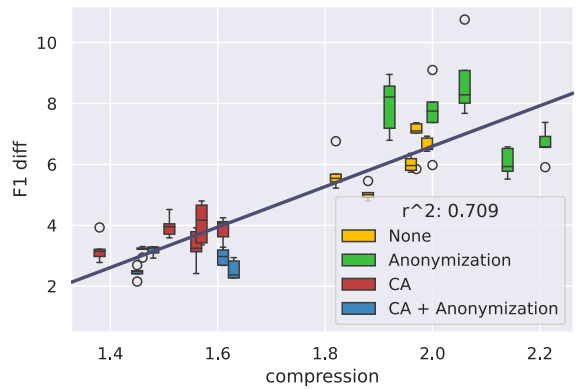


(b) After retraining. Each box represents 10 runs of retraining on the same fine-tuned feature extractor.

Figure 3: Occupation prediction: Compression vs. TPR-gap (Pearson) after various debiasing strategies.



(a) Fine-tuned models. Each point is a single seed for training and testing the model.



(b) After retraining. Each box represents 5 runs of retraining on the same fine-tuned feature extractor.

Figure 4: Coreference resolution: Compression vs. F1 difference after various debiasing strategies.

423 sion rate, while others do not. Except one metric  
 424 (FPR gap sum), the compression rate and the extrinsic  
 425 metric correlate more after retraining. Figure 3  
 426 illustrates this for TPR-gap (Pearson). The increase  
 427 is due to superficial debiasing, especially by sub-  
 428 sampling data, which prior to retraining had low  
 429 extrinsic metrics and relatively high intrinsic met-  
 430 rics. This shows that correlation between extrinsic  
 431 metrics and compression rate for certain metrics  
 432 is stronger than it appeared before retraining. It is  
 433 unsurprising that CEAT does not correlate with any  
 434 extrinsic metrics, since CEAT could not distinguish  
 435 between different models.

436 Coreference resolution shows stronger correla-  
 437 tions between compression rate and extrinsic met-  
 438 rics, but low correlations between Pearson metrics.  
 439 We further discuss cases of no correlation in ap-  
 440 pendix D. Correlations decrease after retraining,  
 441 but metrics that were highly correlated remain so  
 442 ( $> 0.7$  after retraining). The correlations are visu-

443 alized for F1 difference metrics in Figure 4. CEAT  
 444 and extrinsic metrics correlate much less than com-  
 445 pression rate (Table 2). Our results are in line with  
 446 those of Goldfarb-Tarrant et al. (2021), who found  
 447 a lack of correlation between extrinsic metrics and  
 448 WEAT, the static-embedded version of CEAT.

## 5 Related Work 449

450 There are few studies that examine both intrinsic  
 451 and extrinsic metrics. Previous work by Goldfarb-  
 452 Tarrant et al. (2021) showed that debiasing static  
 453 embeddings intrinsically is not correlated with ex-  
 454 trinsic bias, challenging the assumption that intrin-  
 455 sic metrics are predictive of bias. We examine the  
 456 other direction, exploring how extrinsic debiasing  
 457 affects intrinsic metrics. We also extend their work  
 458 to contextualized embeddings, a wider range of  
 459 extrinsic metrics, and a new, more effective intrin-  
 460 sic metric based on information-theoretic probing.  
 461 Studies that inspect extrinsic metrics include either

a challenge dataset curated to expose differences in model behavior by gender, or a test dataset labelled by gender. Among these datasets are Winobias (Zhao et al., 2018), Winogender (Rudinger et al., 2018) and GAP (Webster et al., 2018) for coreference resolution, WinoMT (Stanovsky et al., 2019) for machine translation, EEC (Kiritchenko and Mohammad, 2018) for sentiment analysis, BOLD (Dhamala et al., 2021) for language generation, gendered NLI (Sharma et al., 2020) for natural language inference and Bias in Bios (De-Arteaga et al., 2019) for occupation prediction.

Studies that measure gender bias intrinsically in static word or sentence embeddings measure characteristics of the geometry, such as the proximity between female- and male-related words to stereotypical words, or how embeddings cluster or relate to a gender subspace (Bolukbasi et al., 2016; Caliskan et al., 2017; Gonen and Goldberg, 2019b; Ethayarajh et al., 2019). However, metrics and debiasing methods for static embeddings do not apply directly to contextualized ones. Several studies use sentence templates to adapt to contextual embeddings (May et al., 2019; Kurita et al., 2019; Tan and Celis, 2019). This templated approach is difficult to scale, and lacks the range of representations that a contextual embedding offers. Other work extracts embedding representations of words from natural corpora (Zhao et al., 2019; Guo and Caliskan, 2021; Basta et al., 2019). These studies often adapt the WEAT method (Caliskan et al., 2017), which measures embedding geometry. None measure the effect of the presumably found “bias” on a downstream task.

There is a growing conversation in the field (Barocas et al., 2017; Kate Crawford, 2017; Blodgett et al., 2020; Bommasani et al., 2021) about the importance of articulating the harms of measured bias. In general, extrinsic metrics have clear, interpretable impacts for which harm can be defined. Intrinsic metrics have an unclear effect. Without evidence from a concrete downstream task, a found intrinsic bias is only theoretically harmful. Our work is a step towards understanding whether intrinsic metrics provide valuable insights about bias in a model.

## 6 Discussion and Conclusions

This study examined whether bias in internal representations is related to extrinsic bias. We designed a new framework in which we debias a model on

a downstream task, and measure its intrinsic bias. We found that gender extractability from internal representations, measured by compression rate via MDL probing, reflects bias in a model. Compression was much more reliable than an alternative intrinsic metric for contextualised representations, CEAT. Compression correlated well—to varying degrees—with many extrinsic metrics.

Our results show that when a debiasing method reduces extrinsic metrics but not compression, it indicates that the language model remains biased. When such superficial debiasing occurs, the debiased language model may be reapplied to another task, as in Jin et al. (2021), resulting in unexpected biases and nullifying the supposed debiasing. Our findings suggest that practitioners of NLP should take special care when adopting previously debiased models and inspect them carefully, perhaps using our framework.

Our work also highlighted the importance of the classification layer. Using a debiased objective, such as a balanced dataset, the classification layer can provide significant debiasing. This holds even if the internal representations are biased and the classifier is a single linear layer, as shown in the occupation prediction task. Bias stems in part from internal LM bias and in part from classification bias. Practitioners should focus their efforts on both parts when attempting to debias a model.

We used a broader set of extrinsic metrics than is typically used, and found that the bias metrics behaved differently: some decreased more than others after debiasing, and they correlated differently with compression rate. Debiasing efforts may not be fully understood by testing only a few extrinsic metrics. MDL probing can indicate meaningful debiasing of internal model representations even when not all metrics are easily measurable, since it correlates well with many extrinsic metrics.

A major limitation of this study is the use of gender as a binary variable, which is trans-exclusive. Cao and Daumé III (2020) made the first steps towards inclusive gender bias evaluation in NLP, revealing that coreference systems fail on gender-inclusive text. Further work is required to adjust our framework to non-binary genders, potentially revealing insights about the poor performance of NLP systems in that area.



560  
561  
562  
563  
564  
565  
566  
  
567  
568  
569  
  
570  
571  
572  
573  
574  
575  
  
576  
577  
578  
  
579  
580  
581  
582  
  
583  
584  
585  
  
586  
587  
588  
589  
590  
591  
592  
  
593  
594  
595  
596  
597  
598  
  
599  
600  
601  
  
602  
603  
604  
605  
606  
  
607  
608  
609  
610  
611  
612

## References

Solon Barocas, Kate Crawford, Aaron Shapiro, and Hanna Wallach. 2017. The problem with bias: Allocative versus representational harms in machine learning. In *9th Annual Conference of the Special Interest Group for Computing, Information and Society*.

Solon Barocas, Moritz Hardt, and Arvind Narayanan. 2019. *Fairness and Machine Learning*. fairmlbook.org. <http://www.fairmlbook.org>.

Christine Basta, Marta R. Costa-jussà, and Noe Casas. 2019. [Evaluating the underlying gender bias in contextualized word embeddings](#). In *Proceedings of the First Workshop on Gender Bias in Natural Language Processing*, pages 33–39, Florence, Italy. Association for Computational Linguistics.

Yonatan Belinkov. 2021. [Probing classifiers: Promises, shortcomings, and alternatives](#). *Computational Linguistics 2021*.

Yonatan Belinkov and James Glass. 2019. [Analysis Methods in Neural Language Processing: A Survey](#). *Transactions of the Association for Computational Linguistics*, 7:49–72.

Su Lin Blodgett, Solon Barocas, Hal Daum’è, and H. Wallach. 2020. Language (technology) is power: A critical survey of “bias” in nlp. In *ACL*.

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

Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258*.

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

Yang Trista Cao and Hal Daumé III. 2020. [Toward gender-inclusive coreference resolution](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4568–4595, Online. Association for Computational Linguistics.

Maria De-Arteaga, Alexey Romanov, H. Wallach, J. Chayes, C. Borgs, A. Chouldechova, S. C. Geyik, K. Kenthapadi, and A. Kalai. 2019. Bias in bios: A case study of semantic representation bias in a high-stakes setting. *Proceedings of the Conference on Fairness, Accountability, and Transparency*.

J. Dhamala, Tony Sun, Varun Kumar, Satyapriya Krishna, Yada Pruksachatkun, Kai-Wei Chang, and Rahul Gupta. 2021. [Bold: Dataset and metrics for measuring biases in open-ended language generation](#). *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*. 613  
614  
615  
616  
617  
618

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. 619  
620  
621  
622  
623  
624

Seraphina Goldfarb-Tarrant, Rebecca Marchant, Ricardo Muñoz Sánchez, Mugdha Pandya, and Adam Lopez. 2021. [Intrinsic bias metrics do not correlate with application bias](#). 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 1926–1940, Online. Association for Computational Linguistics. 625  
626  
627  
628  
629  
630  
631  
632  
633

Hila Gonen and Y. Goldberg. 2019a. [Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them](#). In *NAACL-HLT*. 634  
635  
636  
637

Hila Gonen and Yoav Goldberg. 2019b. [Lipstick on a pig: Debiasing methods cover up systematic gender biases in word embeddings but do not remove them](#). 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 609–614, Minneapolis, Minnesota. Association for Computational Linguistics. 638  
639  
640  
641  
642  
643  
644  
645  
646

Wei Guo and Aylin Caliskan. 2021. 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*, pages 122–133. 647  
648  
649  
650  
651

John Hewitt and Percy Liang. 2019. [Designing and interpreting probes with control tasks](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, Hong Kong, China. Association for Computational Linguistics. 652  
653  
654  
655  
656  
657  
658

Xisen Jin, Francesco Barbieri, Brendan Kennedy, Aida Mostafazadeh Davani, Leonardo Neves, and Xiang Ren. 2021. [On transferability of bias mitigation effects in language model fine-tuning](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3770–3783, Online. Association for Computational Linguistics. 659  
660  
661  
662  
663  
664  
665  
666

Kate Crawford. 2017. The trouble with bias. keynote at neurips. 667  
668

669	Svetlana Kiritchenko and Saif Mohammad. 2018. <a href="#">Examining gender and race bias in two hundred sentiment analysis systems</a> . In <i>Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics</i> , pages 43–53, New Orleans, Louisiana. Association for Computational Linguistics.	726
670		727
671		728
672		
673		729
674		730
		731
675	Keita Kurita, Nidhi Vyas, Ayush Pareek, Alan W Black, and Yulia Tsvetkov. 2019. <a href="#">Measuring bias in contextualized word representations</a> . In <i>Proceedings of the First Workshop on Gender Bias in Natural Language Processing</i> , pages 166–172, Florence, Italy. Association for Computational Linguistics.	732
676		733
677		734
678		735
679		736
680		737
681		738
682	Kweku Kwegyir-Aggrey, Rebecca Santorella, and Sarah M. Brown. 2021. Everything is relative: Understanding fairness with optimal transport. <i>ArXiv</i> , abs/2102.10349.	
683		
684		
685	Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018a. <a href="#">Higher-order coreference resolution with coarse-to-fine inference</a> . In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)</i> , pages 687–692, New Orleans, Louisiana. Association for Computational Linguistics.	739
686		740
687		741
688		742
689		743
690		744
691		745
692		746
693	Kenton Lee, Luheng He, and Luke Zettlemoyer. 2018b. Higher-order coreference resolution with coarse-to-fine inference. <i>arXiv preprint arXiv:1804.05392</i> .	
694		
695		
696	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <a href="#">Roberta: A robustly optimized BERT pretraining approach</a> . <i>CoRR</i> , abs/1907.11692.	
697		
698		
699		
700		
701	Chandler May, Alex Wang, Shikha Bordia, Samuel R. Bowman, and Rachel Rudinger. 2019. <a href="#">On measuring social biases in sentence encoders</a> . In <i>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)</i> , pages 622–628, Minneapolis, Minnesota. Association for Computational Linguistics.	747
702		748
703		749
704		750
705		
706		
707		
708		
709	Michael Mendelson and Yonatan Belinkov. 2021. <a href="#">De-biasing methods in natural language understanding make bias more accessible</a> . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> . Association for Computational Linguistics.	751
710		752
711		753
712		754
713		755
714		756
715	F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. <i>Journal of Machine Learning Research</i> , 12:2825–2830.	757
716		758
717		
718		
719		
720		
721		
722	Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. <a href="#">Null it out: Guarding protected attributes by iterative nullspace projection</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7237–7256, Online. Association for Computational Linguistics.	759
723		760
724		761
725		762
		763
		764
		765
		766
		767
		768
		769
		770
		771
		772
		773
		774
		775
		776
		777
		778
		779

- 780 Kellie Webster, Marta Recasens, Vera Axelrod, and Ja-  
781 son Baldrige. 2018. [Mind the GAP: A balanced](#)  
782 [corpus of gendered ambiguous pronouns](#). *Transac-*  
783 *tions of the Association for Computational Linguis-*  
784 *tics*, 6:605–617.
- 785 R. Weischedel, E. Hovy, M. Marcus, and Martha Palmer.  
786 2013. Ontonotes : A large training corpus for  
787 enhanced processing. *LDC2013T19, Philadelphia,*  
788 *Penn.: Linguistic Data Consortium.*
- 789 Liyan Xu and Jinho D. Choi. 2020. [Revealing the myth](#)  
790 [of higher-order inference in coreference resolution](#).  
791 In *Proceedings of the 2020 Conference on Empirical*  
792 *Methods in Natural Language Processing (EMNLP)*,  
793 pages 8527–8533. Association for Computational  
794 Linguistics.
- 795 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell,  
796 Vicente Ordonez, and Kai-Wei Chang. 2019. [Gender](#)  
797 [bias in contextualized word embeddings](#). In *Proceed-*  
798 *ings of the 2019 Conference of the North American*  
799 *Chapter of the Association for Computational Lin-*  
800 *guistics: Human Language Technologies, Volume*  
801 *1 (Long and Short Papers)*, pages 629–634, Min-  
802 neapolis, Minnesota. Association for Computational  
803 Linguistics.
- 804 Jieyu Zhao, Tianlu Wang, Mark Yatskar, Vicente Or-  
805 donez, and Kai-Wei Chang. 2018. [Gender bias in](#)  
806 [coreference resolution: Evaluation and debiasing](#)  
807 [methods](#). In *Proceedings of the 2018 Conference*  
808 *of the North American Chapter of the Association for*  
809 *Computational Linguistics: Human Language Tech-*  
810 *nologies, Volume 2 (Short Papers)*, pages 15–20, New  
811 Orleans, Louisiana. Association for Computational  
812 Linguistics.

813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859

## A Implementation Details

We used RoBERTa in all models which has 120M parameters. We use following random seeds in all repeated experiments: 0, 5, 11, 26, 42, 46, 50, 63, 83, 90.

### A.1 Occupation Classification

We fine-tuned a RoBERTa-base model with a linear classification layer on top. Training was done for 10 epochs at a learning rate of 5e-5, batch size of 64. The input to RoBERTa was the biography tokens, which is limited to the first 128 tokens. The resulting [CLS] token embedding is fed to the classifier to predict the occupation. The probing task involves using the same [CLS] token and training the probing classifier to predict the gender of the person in the biography. The experiments without fine-tuning included either a pre-trained or a previously fine-tuned RoBERTa. We first extracted the pre-trained RoBERTa’s embeddings of tokens from the [CLS] and then trained a linear classifier on them. The learning rate was 0.001 and the batch size was 64. We trained the classification layer with pre-trained RoBERTa on 300 epochs, but with fine-tuned RoBERTa, 10 epochs were sufficient. For all training processes, the epoch with the greatest validation accuracy was saved. Fine-tuning took 7 hours on GeForce RTX 2080 Ti GPU. Bias in Bios contains almost 400k biographies, and we obtain validation (10%) and test set (25%) by splitting with Scikit-learn’s (Pedregosa et al., 2011) test\_train\_split with our random seeds.

### A.2 Coreference Resolution

We use the implementation of Xu and Choi (2020), a model that was introduced by Lee et al. (2018b) and has been adopted by many coreference resolution models. Coreference resolution is the process of clustering different mentions in a text that refer to the same real-world entities. The task is solved by detecting mentions through text spans and then predicting for each pair of spans if they represent the same entity. The span representations were extracted with a RoBERTa model, which is fine-tuned throughout the training process, except in the retraining experiment. Fine-tuning took 3 hours on NVIDIA RTX A6000 GPU. Ontonotes 5.0 has 625k sentences and we use the standard validation and test splits.

### A.3 Probing Classifier

We use the MDL probe (Voita and Titov, 2020) implementation by Mendelson and Belinkov (2021). In all experiments, we use a linear probe and train it with a batch size of 16 and a learning rate of 1e-3. The timestamps used, meaning the accumulating fractions of data that the probe is trained on, are 2.0%, 3.0%, 4.4%, 6.5%, 9.5%, 14.0%, 21.0%, 31.0%, 45.7%, 67.6%, 100%.

### A.4 Metrics

#### A.4.1 Fairness-Based Metrics Implementation

All three statistical fairness metrics measure the difference between two probability distributions, where this difference describes a notion of bias. We calculate Independence and Separation via Kullback–Leibler (KL) divergence, using the AllenNLP implementation (<https://github.com/allenai/allennlp>). We calculate Sufficiency via Wasserstein distance instead, which is motivated by Kwegyir-Aggrey et al. (2021). In this case, we cannot use KL divergence, since there are some classes that do not occur in model predictions for both male and female genders. This causes the probability distributions to not have the same support, and KL divergence is unbounded. Wasserstein distance lacks the requirement for equal support.

#### A.4.2 Classification Metrics Interpretation in Winobias

Winobias datasets contain pairs of stereotypical and anti-stereotypical sentences. The stereotypes are derived from the US labor statistics (for instance, a profession with a majority of males is stereotypically male). Since coreference resolution is viewed as a clustering problem, it is usually measured via clustering evaluation metrics. Coreference resolution is commonly measured as the average F1 score of these, and the same is true for Winobias. Nevertheless, coreference resolution is accomplished by making a prediction for each pair of mentions, so it can be seen as a classification task. Winobias can be viewed as a simpler task than general coreference resolution, as it contains exactly two mentions of professions and one pronoun, which refers to exactly one profession. Therefore, we reframe it as a classification problem. In a Winobias sentence with two professions  $x$  and  $y$ , as well as a pronoun  $p$ , where  $p$  is referring to  $x$ , a true positive would be to cluster  $x$  and  $p$  together, while a false positive would be to cluster  $y$  and  $p$  together. Our

860  
861  
862  
863  
864  
865  
866  
867  
868  
869  
870  
871  
872  
873  
874  
875  
876  
877  
878  
879  
880  
881  
882  
883  
884  
885  
886  
887  
888  
889  
890  
891  
892  
893  
894  
895  
896  
897  
898  
899  
900  
901  
902  
903  
904  
905  
906  
907  
908

classification metrics are derived based on these definitions. For instance, the TPR gap for profession “teacher”, which is a stereotypical female occupation, is the TPR rate on pro-stereotypical sentences (with a female pronoun) minus the TPR rate on anti-stereotypical sentences (with a male pronoun).

### A.4.3 CEAT

The Word Embedding Association Test (WEAT) developed by (Caliskan et al., 2017) is a method for evaluating bias in static word embeddings. The test is defined as follows: given two sets of target words  $X, Y$  (e.g., ‘executive’, ‘management’, ‘professional’ and ‘home’, ‘parents’, ‘children’) and two sets of attribute words (e.g., male names and female names), and using  $\vec{w}$  to represent the word embedding for word  $w$ , the effect size is:

$$ES = \text{mean}_{x \in X} s(x, A, B) - \text{mean}_{y \in Y} s(y, A, B)$$

where

$$s(x, A, B) = \frac{\text{mean}_{a \in A} \cos(\vec{x}, \vec{a}) - \text{mean}_{a \in B} \cos(\vec{x}, \vec{a})}{\text{std-dev}_{w \in X \cup Y} s(w, A, B)}$$

In essence, the effect size measures how different are the distances between the embedding vectors of each target group and the attribute groups. Specifically, if  $s(x, A, B) > 0$ ,  $\vec{x}$  is more similar to attribute words  $B$  and vice versa. For instance, a larger effect size is observed if target words  $X$  are more similar to attribute words  $A$  and target words  $Y$  are more similar to attribute words  $B$ .  $|ES| > 0.5$  and  $|ES| > 0.8$  are considered medium and large effect sizes, respectively (Rice and Harris, 2005). The null hypothesis holds that there is no difference between the two sets of target words in terms of their relative similarity to the two sets of attribute words, indicating that there are no biased associations. Statistical significance is defined by the p-value of WEAT, which reflects the probability of observing the effect size under the null hypothesis.

Since a word can take on a great variety of vector representations in a contextual setting,  $ES$  varies according to the sentences used to extract word representation. Thus, to adopt WEAT to contextualized representations, the Combined Effect Size (CES) (Guo and Caliskan, 2021) is derived as the

distribution of WEAT effect sizes over many possible contextual word representations:

$$CES(X, Y, A, B) = \frac{\sum_{i=1}^N v_i ES_i}{\sum_{i=1}^N v_i}$$

where  $ES_i$  denotes the WEAT effect size of the  $i$ ’th choice of word representations from a large corpus, and  $v_i$  is the inverse of the sum of in-sample variance  $V_i$  and between-sample variance in the distribution of random-effects. As in Guo and Caliskan (2021), the representation for each word is derived from 10,000 random sentences extracted from a corpus of Reddit comments.

The combined effect size of each of the models is examined on WEAT stimulus 6, which contains target words of career/family and attribute words of male/female names. This was the only one that detected bias on a pre-trained RoBERTa (CES close to 0.5 and  $p < 0.05$ ). The points that we kept in our analysis are those where  $p < 0.05$ , which make up 90% of the points in occupation prediction and 95% of the points in coreference resolution.

## B Full Results

In this section we provide the full results of a RoBERTa model trained on the downstream task. The results for the occupation prediction task after fine-tuning are presented in Table 3 and Table 4 presents the retrained model results. Figure 5 illustrates the correlations between extrinsic metrics and compression rate before and after retraining. Table 5 presents the complete results of the model trained without fine-tuning, meaning that the RoBERTa model is the pretrained version from Liu et al. (2019) and only the classification layer was updated. Subsampling the dataset has significant debiasing effects, which suggests that this debiasing method can achieve low extrinsic bias even when internal bias exists.

Regarding the coreference resolution task, Table 6 displays the results on a finetuned model and Table 7 displays the retraining results. Figure 6 shows the correlations between compression rate and extrinsic metrics before and after the retraining.

Metric	Debiasing Strategy			
	None	Oversampling	Subsampling	Scrubbing
Compression	4.121 ± 1.238	8.522* ± 2.354	3.568 ± 1.516	<b>1.699*</b> ± 0.138
Accuracy	<b>0.861</b> ± 0.005	0.852* ± 0.004	<b>0.861</b> ± 0.003	0.851* ± 0.003
TPR gap (P)	0.763 ± 0.071	0.729 ± 0.067	<b>0.319*</b> ± 0.114	0.704* ± 0.068
TPR gap (S)	2.391 ± 0.257	2.145* ± 0.220	<b>1.598*</b> ± 0.273	2.019* ± 0.262
FPR gap (P)	0.591 ± 0.052	0.491* ± 0.059	<b>0.087*</b> ± 0.094	0.552 ± 0.063
FPR gap (S)	0.075 ± 0.010	0.085* ± 0.011	<b>0.030*</b> ± 0.006	0.057* ± 0.007
Precision gap (P)	0.398 ± 0.053	0.327* ± 0.044	<b>0.166*</b> ± 0.055	0.347* ± 0.050
Precision gap (S)	0.015 ± 0.001	0.015 ± 0.001	<b>0.011*</b> ± 0.001	0.013* ± 0.001
Independence gap (S)	0.009 ± 0.002	0.008 ± 0.002	<b>0.001*</b> ± 0.000	0.005* ± 0.001
Separation gap (S)	0.327 ± 0.051	0.305 ± 0.030	<b>0.204*</b> ± 0.032	0.296 ± 0.053
Sufficiency gap (S)	9.451 ± 1.945	8.324* ± 1.537	<b>1.218*</b> ± 0.330	4.930* ± 0.927

Table 3: Occupation Prediction: Results on a RoBERTa-based model trained over 10 seeds. Significant reduction or increase in a metric ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy is None), is marked with \*. The lowest bias score or highest performance metric in each column is marked with **bold**. P = Pearson; S = Sum.

Metric	Debiasing Strategy			
	None	Oversampling	Subsampling	Scrubbing
Compression	4.121 ± 1.238	8.522 ± 2.354	3.568 ± 1.516	1.699 ± 0.138
Accuracy	0.859 ± 0.004	0.856 ± 0.003	0.853 ± 0.003	0.854 ± 0.003
TPR gap (P)	0.777 ± 0.047	0.813* ± 0.040	<b>0.704*</b> ± 0.075	0.714* ± 0.068
TPR gap (S)	2.482 ± 0.238	2.593* ± 0.240	2.164* ± 0.284	<b>1.989*</b> ± 0.227
FPR gap (P)	0.596 ± 0.041	0.603 ± 0.047	0.602 ± 0.041	<b>0.536*</b> ± 0.038
FPR gap (S)	0.073 ± 0.008	0.068* ± 0.007	0.081* ± 0.012	<b>0.059*</b> ± 0.005
Precision gap (P)	0.373 ± 0.067	0.356* ± 0.070	0.338* ± 0.054	<b>0.309*</b> ± 0.053
Precision gap (S)	0.016 ± 0.002	0.017* ± 0.002	0.015* ± 0.002	<b>0.014*</b> ± 0.002
Independence gap (S)	0.009 ± 0.002	0.010* ± 0.002	0.009 ± 0.003	<b>0.005*</b> ± 0.001
Separation gap (S)	0.334 ± 0.050	0.328 ± 0.048	0.300* ± 0.049	<b>0.274*</b> ± 0.041
Sufficiency gap (S)	9.701 ± 1.305	10.908* ± 1.354	8.370* ± 2.558	<b>5.239*</b> ± 0.798

Table 4: Occupation Prediction after retraining: Results on a RoBERTa-based model after retraining of the classification layer with 10 seeds for each pre-trained model. Significant reduction or increase in a metric ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy is None), is marked with \*. The lowest bias score or highest performance metric in each column is marked with **bold**. P = Pearson; S = Sum.

Metric	Debiasing Strategy			
	None	Oversampling	Subsampling	Scrubbing
Accuracy	0.824 ± 0.003	0.815* ± 0.005	<b>0.831*</b> ± 0.001	0.807* ± 0.003
TPR gap (P)	0.839 ± 0.011	0.443* ± 0.053	<b>0.158*</b> ± 0.156	0.814 ± 0.029
TPR gap (S)	3.088 ± 0.192	<b>1.545*</b> ± 0.177	1.621* ± 0.088	3.154 ± 0.332
FPR gap (P)	0.598 ± 0.016	0.369* ± 0.029	<b>0.067*</b> ± 0.050	0.550* ± 0.012
FPR gap (S)	0.087 ± 0.004	0.041* ± 0.004	<b>0.027*</b> ± 0.003	0.112* ± 0.005
Precision gap (P)	0.476 ± 0.027	0.163* ± 0.025	<b>0.134*</b> ± 0.065	0.479 ± 0.038
Precision gap (S)	0.017 ± 0.001	0.012* ± 0.001	<b>0.010*</b> ± 0.001	0.016* ± 0.002
Independence gap (S)	0.014* ± 0.002	0.001* ± 0.000	<b>0.000*</b> ± 0.000	0.022* ± 0.001
Separation gap (S)	0.336* ± 0.044	0.214* ± 0.038	<b>0.203*</b> ± 0.024	0.432* ± 0.048
Sufficiency gap (S)	12.019* ± 1.721	2.105* ± 0.576	<b>1.478*</b> ± 0.394	13.798* ± 0.966

Table 5: Occupation Prediction: Results on a RoBERTa-based model trained without fine-tuning, over 5 seeds. The compression rate computed on a pre-trained RoBERTa model is 10.122. Significant reduction or increase in a metric ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy is None), is marked with \*. The lowest bias score or highest performance metric in each column is marked with **bold**. P = Pearson; S = Sum.

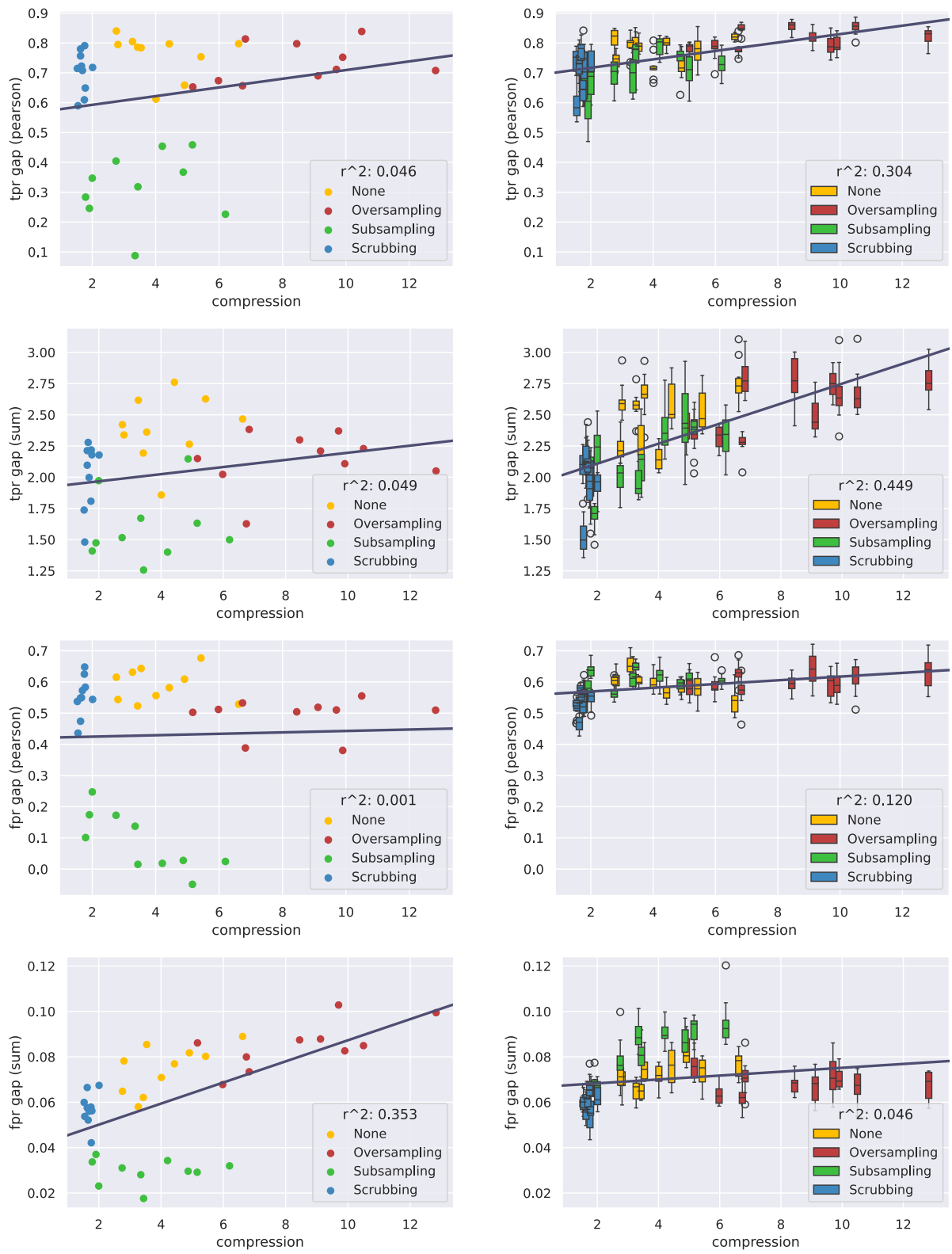


Figure 5: Occupation prediction: Before (left) and after (right) plots of compression rate versus and extrinsic metric.



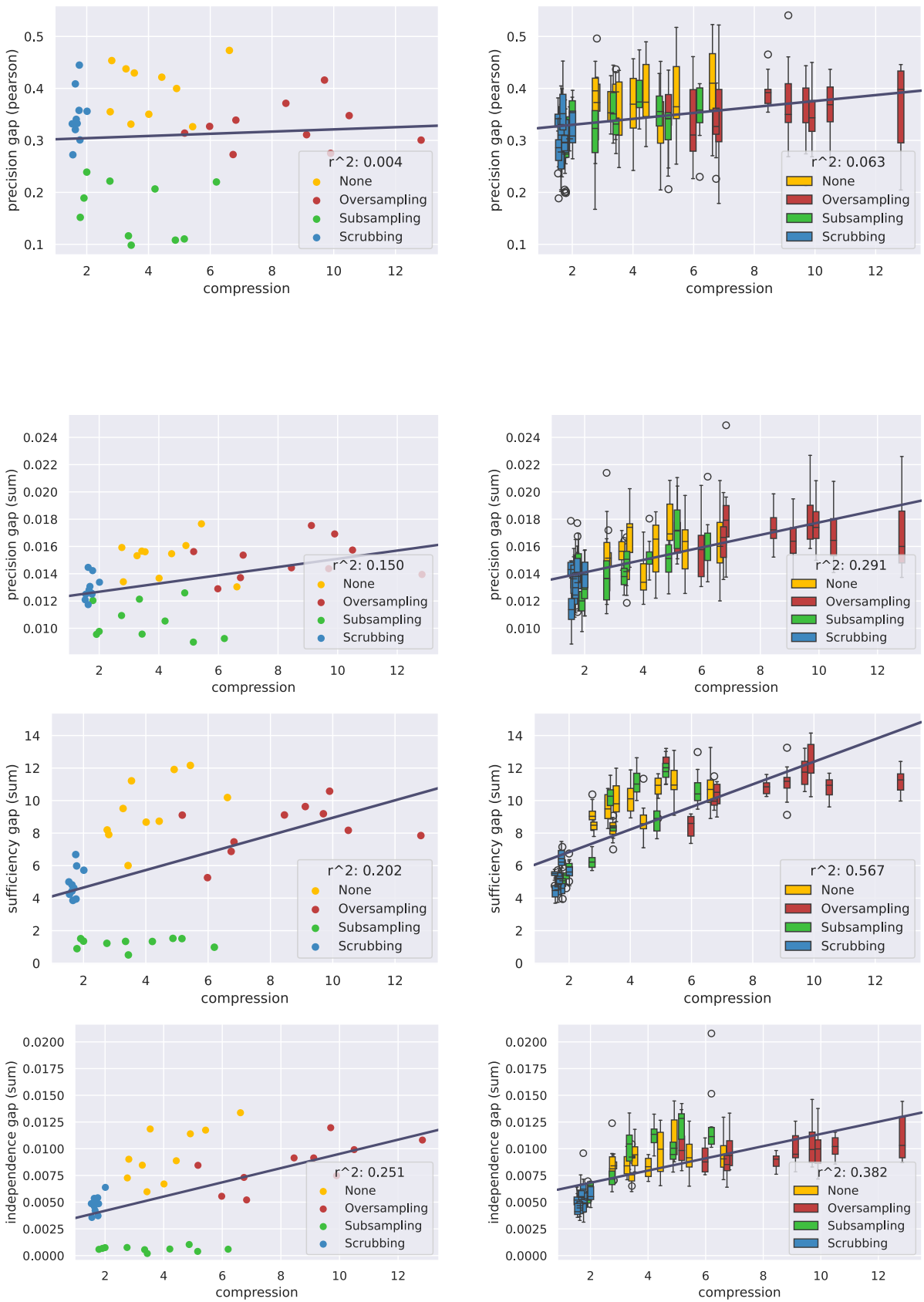


Figure 5: Occupation prediction: Before (left) and after (right) plots of compression rate versus and extrinsic metric.

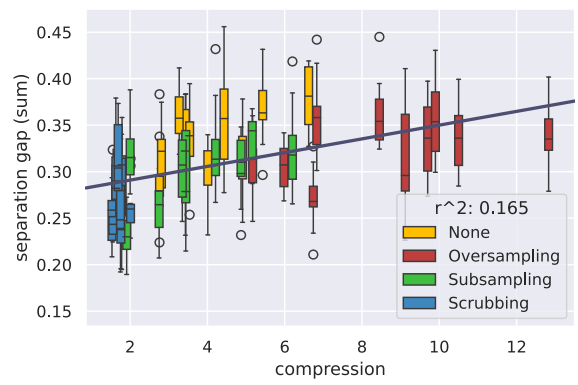
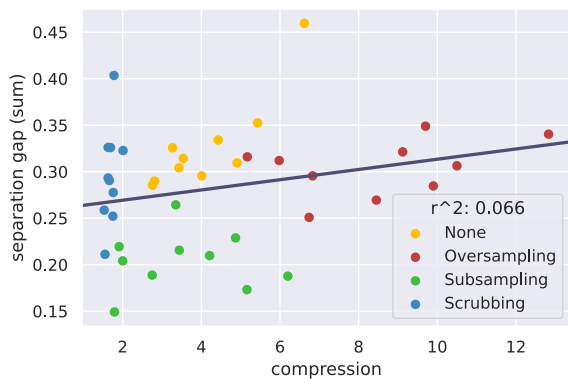


Figure 5: Occupation prediction: Before (left) and after (right) plots of compression rate versus and extrinsic metric.

Metric	Debiasing Strategy			
	None	Anon	CA	Anon + CA
Compression	1.984 ± 0.101	2.073* ± 0.102	<b>1.502*</b> ± 0.075	1.540* ± 0.098
F1 (Ontonotes test)	76.406 ± 0.165	76.538 ± 0.176	77.187* ± 0.071	<b>77.246*</b> ± 0.230
F1 diff ( <i>pro</i> – <i>anti</i> )	6.631 ± 1.013	7.256 ± 0.846	<b>2.302*</b> ± 0.466	2.422* ± 0.714
TPR gap (P)	0.654 ± 0.069	0.710* ± 0.047	<b>0.607</b> ± 0.082	0.627 ± 0.100
TPR gap (S)	4.884 ± 0.698	4.870 ± 0.509	2.041* ± 0.228	<b>2.014*</b> ± 0.286
FPR gap (P)	0.602 ± 0.036	0.620 ± 0.056	<b>0.572</b> ± 0.078	0.629 ± 0.107
FPR gap (S)	0.120 ± 0.015	0.128 ± 0.011	0.050* ± 0.006	<b>0.049*</b> ± 0.007
Precision gap (P)	0.654 ± 0.068	0.710* ± 0.048	<b>0.607</b> ± 0.083	0.627 ± 0.099
Precision gap (S)	0.061 ± 0.009	0.061 ± 0.006	0.026* ± 0.003	<b>0.025*</b> ± 0.004
Independence gap (S)	0.027 ± 0.008	0.025 ± 0.004	<b>0.004*</b> ± 0.001	<b>0.004*</b> ± 0.001
Separation gap (S)	1.247 ± 0.150	1.344 ± 0.137	<b>0.537*</b> ± 0.061	0.557* ± 0.070
Sufficiency gap (S)	8.684 ± 1.883	8.816 ± 1.544	1.673* ± 0.354	<b>1.557*</b> ± 0.384

Table 6: Coreference resolution: results on Ontonotes test set and Winobias challenge set. Each model was trained over 10 seeds. \* Marks significant reduction or increase in bias ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy None). The lowest bias score or highest performance metric in each column is in **bold**. P = Pearson; S = Sum.

Metric	Debiasing Strategy			
	None	Anon	CA	Anon + CA
Compression	1.984 ± 0.065	2.073* ± 0.104	<b>1.502*</b> ± 0.081	1.540* ± 0.079
F1 (Ontonotes test)	76.40* ± 0.16	76.48* ± 0.22	76.72* ± 0.15	<b>76.91*</b> ± 0.19
F1 diff ( <i>pro</i> – <i>anti</i> )	6.072 ± 0.789	7.417* ± 1.280	3.674* ± 0.599	<b>2.858*</b> ± 0.382
TPR gap (P)	<b>0.635</b> ± 0.053	0.688* ± 0.052	0.679* ± 0.062	0.654 ± 0.049
TPR gap (S)	4.561 ± 0.414	5.143* ± 0.713	2.590* ± 0.420	<b>2.178*</b> ± 0.201
FPR gap (P)	<b>0.579</b> ± 0.046	0.637* ± 0.055	0.620* ± 0.070	0.692* ± 0.075
FPR gap (S)	0.113 ± 0.011	0.126* ± 0.016	0.063* ± 0.010	<b>0.052*</b> ± 0.004
Precision gap (P)	<b>0.636</b> ± 0.052	0.690* ± 0.052	0.679* ± 0.062	0.652 ± 0.050
Precision gap (S)	0.057 ± 0.005	0.064* ± 0.009	0.032* ± 0.005	<b>0.027*</b> ± 0.003
Independence gap (S)	0.022 ± 0.003	0.026* ± 0.006	0.006* ± 0.002	0.004* ± 0.001
Separation gap (S)	1.188 ± 0.114	1.336* ± 0.175	0.670* ± 0.111	<b>0.594*</b> ± 0.057
Sufficiency gap (S)	7.350 ± 0.914	8.655* ± 1.726	0.2.401* ± 0.610	<b>1.653*</b> ± 0.294

Table 7: Coreference resolution after retraining: results on Ontonotes test set and extrinsic bias metrics on Winobias challenge set. Each model finetuned over 10 seeds and re-trained over 5 seeds. \* Marks significant reduction or increase in bias ( $p < 0.05$  on Pitman’s permutation test), compared to the non-debiased model (debiasing strategy None). The lowest bias score or highest performance metric in each column is in **bold**. P = Pearson; S = Sum.

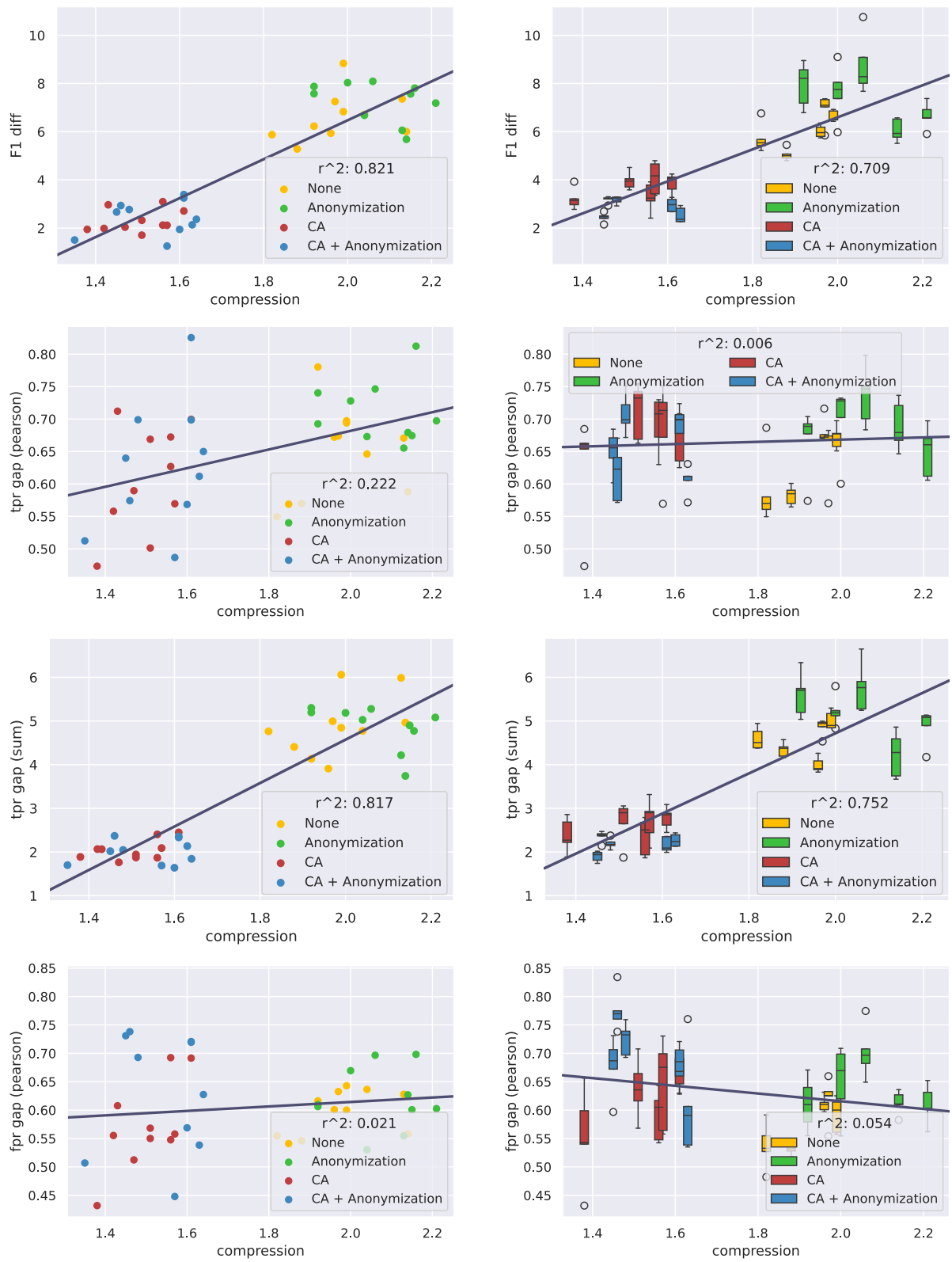


Figure 6: Coreference resolution: Before (left) and after (right) plots of compression rate versus and extrinsic metric.

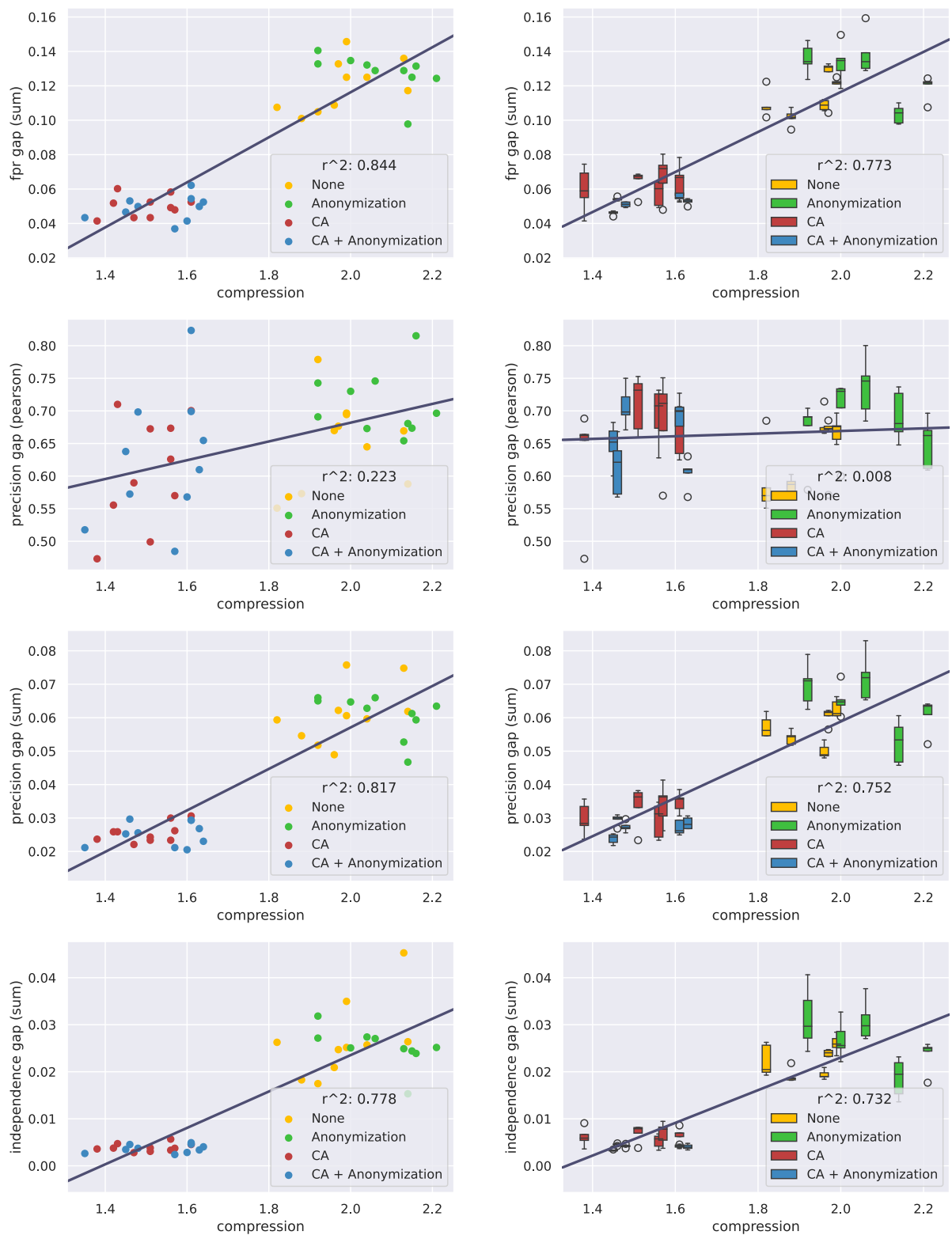


Figure 6: Coreference resolution: Before (left) and after (right) plots of compression rate versus and extrinsic metric.

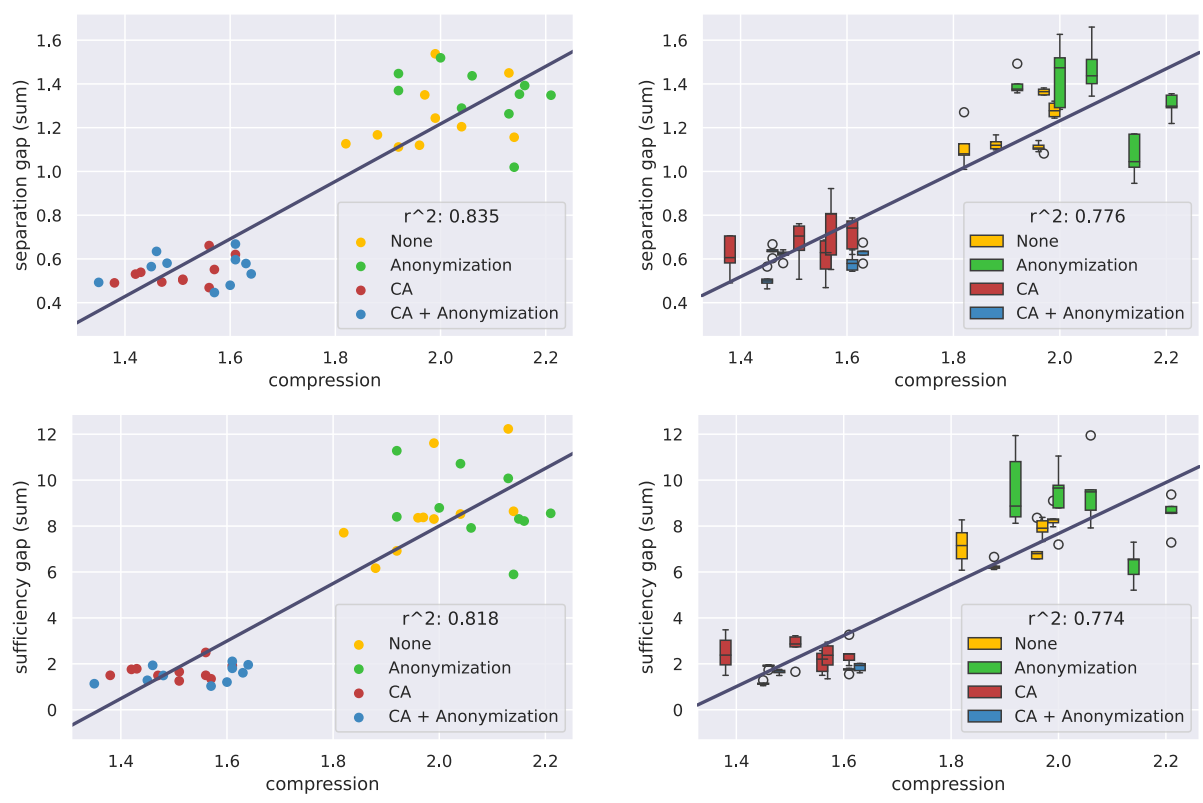


Figure 6: Coreference resolution: Before (left) and after (right) plots of compression rate versus and extrinsic metric.

Female Words	Male Words
husband, women, gender, listed, practices, nurse, specializes, children, ba, child, reading, families, location, place, affiliated, family, experiences, spanish, love, justice	chief, companies, computer, applications, md, accepts, known, doctors, npi, sports, philosoph', problems, rating, no, systems, theory, practicing, software, security, major

Table 8: Top 20 significant words used to predict gender on all biographies, as obtained from a logistic regression model trained on predicting the gender of a person described in a biography. The words are sorted by importance.

Female Words	Male Words
husband , women, midwife , providing book , includes, joining, faculty	holds , emergency, vanderbilt, forces, registered, mental, assistant, president

Table 9: Top 8 words used to predict gender of female and male nurses, as obtained from a logistic regression model trained on predicting the gender of a person described in a biography. The words are sorted by importance.

## C Why is scrubbing not as effective as subsampling?

The debiasing method of subsampling significantly reduced external biases in the occupation prediction task. Although compression rates show that scrubbing reduced more gender information, subsampling outperforms it as a debiasing method. We find that in spite of the scrubbing, a probe is able to correctly identify the gender from an internal representation with 68.8% accuracy compared to 90.7% on the original, non-scrubbed data. This means that although the scrubbing process reduces extrinsic bias significantly, gender information is still embedded in the [CLS] token embeddings.

To investigate the source of gender information after scrubbing, we use logistic regression (LR)

model to predict the gender from the Bag-of-Words of the scrubbed biographies. We perform an iterative process for automatic extra scrubbing: in each iteration we (1) train a LR model for gender prediction (2) scrub the  $n$  most significant words for each gender according to the LR weights. The most relevant words among 5 seeds of training with  $n=10$  words scrubbed per iteration are displayed in Table 8. The model learns indirect correlations to gender in the absence of explicit gendered words. Because the significant words are related to male- or female-dominated professions, we conducted the process on a specific profession. Table 9 presents the most significant words for biographies of nurses. There are differences in wording even between females and males in the same profession. The results of this study are in line with the results of other studies that have been conducted on the way biographies are written for men and women (Wagner et al., 2016; Sun and Peng, 2021).

Subsampling is therefore more effective even when gender information is present since it prevents the model from learning correlations between gender information and a profession whereas scrubbing only attempts to remove gender indicators without removing correlations. On the other hand, it is possible that oversampling is less effective for debiasing since seeing more non-unique examples an unrepresented group encourages learning correlations.

## D A closer look into no-correlation cases

### D.1 Occupation Prediction

Although compression has the ability to identify bias in most cases, some metrics still show little or no correlation with compression rate. These results suggest that gender information comprises only one facet of embedded bias in the representations. Other factors that may influence these metrics are not considered or measured, such as the connection between a name and a profession.

For example, as can be see in Tables 3 and 4, LMs finetuned on subsampled data have the largest FPR gaps after retraining, despite being the least biased before retraining, while those finetuned on oversampled data have the next-to-lowest FPR gaps after retraining. The information encoded in the internal representations may have been encoded in a manner that allowed the classification layer to exhibit a smaller FPR gap when trained on a balanced dataset. However, when the classification

1061 layer was retrained on biased training data, it used  
1062 the same features to make biased predictions.

## 1063 **D.2 Coreference Resolution**

1064 The cases where there is no correlation between  
1065 our intrinsic metric and an extrinsic metric are the  
1066 cases where the metric is based on Pearson corre-  
1067 lation. Unlike occupation prediction, coreference  
1068 resolution seems to exhibit no correlation between  
1069 those metrics and compression rate. These metrics  
1070 are computed as the Pearson correlation between a  
1071 performance gap for a specific profession and the  
1072 percentage of women in that profession, however  
1073 the percentages are computed differently in each  
1074 task: in occupation prediction, the percentages are  
1075 computed from the train set, focusing on the rep-  
1076 resentation each gender has in the data. In Wino-  
1077 bias, the percentages are taken from the US labor  
1078 statistics, and are unrelated to the training dataset  
1079 statistics. We note that the two statistics can be dif-  
1080 ferent - the real-world representation of women in a  
1081 profession does not have to be equal to their repre-  
1082 sentation in written text (Suresh and Gutttag, 2021).  
1083 We thus decided to test what happens if we change  
1084 the statistics used in Winobias to dataset statistics,  
1085 but Ontonotes 5.0 has very little representation to  
1086 each profession and the statistics extracted from  
1087 it would not be reliable. We thus took a different  
1088 approach and computed the Pearson correlations  
1089 for occupation prediction with real world statistics  
1090 instead of dataset statistics. To do this, we mapped  
1091 the professions appearing in this dataset to pro-  
1092 fessions from the US labor statistics, and dropped  
1093 those who could no be mapped (6 out of 29 of the  
1094 professions which is 21.4%). We then repeated  
1095 all experiments on the Pearson metrics using these  
1096 statistics. Figure 7 shows the results. Correlations  
1097 are very different when computed with respect to  
1098 real-world statistics. TPR-gap has no correlation at  
1099 all although it had with training data statistics, the  
1100 correlation for FPR-gap after retraining exists but  
1101 is negative, and the correlation with precision-gap  
1102 does not exist after retraining. We thus conclude  
1103 that the Pearson metrics are less reliable as they are  
1104 heavily dependent on the statistics with respect to  
1105 which they are calculated.



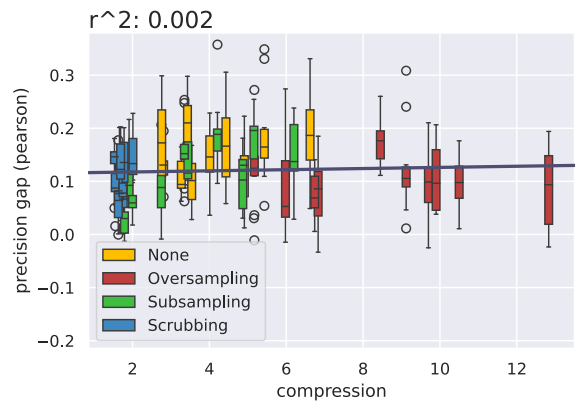
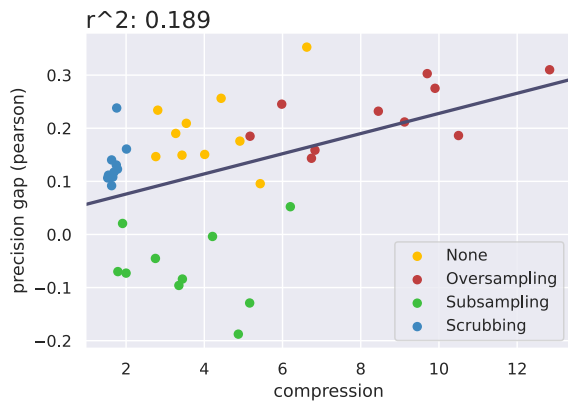
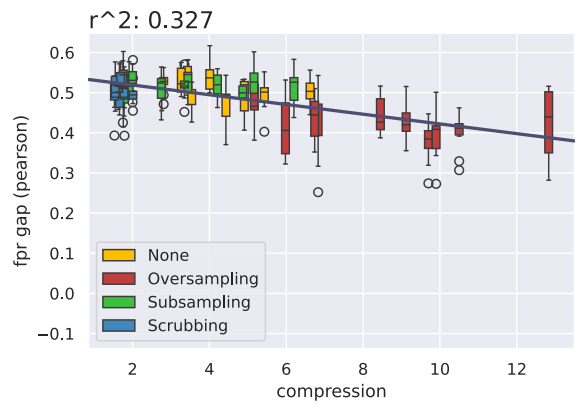
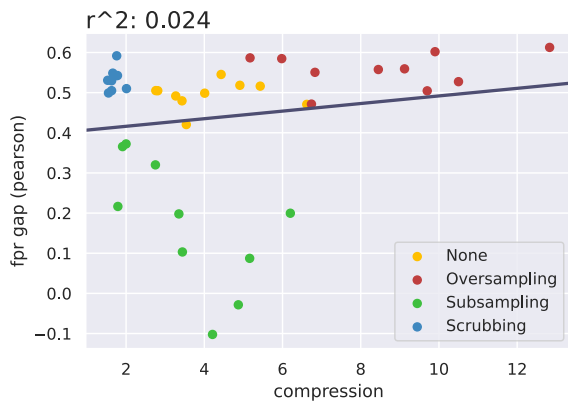
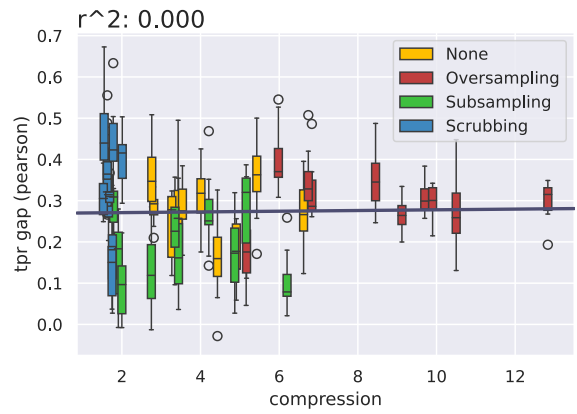
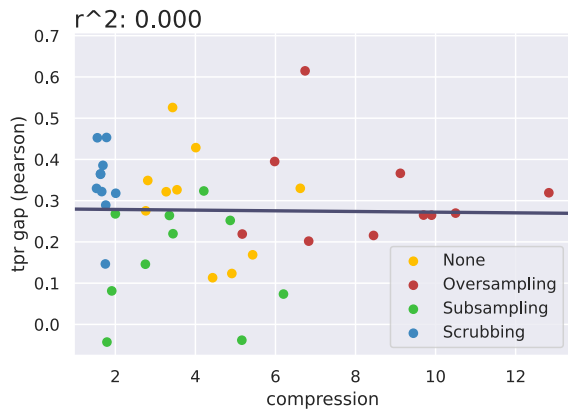


Figure 7: Occupation prediction: Before (left) and after (right) plots of compression rate versus and Pearson metrics as computed from real-world statistics.