

Should We Attend More or Less? Modulating Attention for Fairness

Anonymous ACL submission

Abstract

The abundance of annotated data in natural language processing (NLP) poses both opportunities and challenges. While it enables the development of high-performing models for a variety of tasks, it also poses the risk of models learning harmful biases from the data, such as gender stereotypes. In this work, we investigate the role of attention, a widely-used technique in current state-of-the-art NLP models, in the propagation of social biases. Specifically, we study the relationship between the entropy of the attention distribution and the model’s performance and fairness. We then propose a novel method for modulating attention weights to improve model fairness after training. Since our method is only applied post-training and pre-inference, it is an intra-processing method and is, therefore, less computationally expensive than existing in-processing and pre-processing approaches. Our results show an increase in fairness and minimal performance loss on different text classification and generation tasks using language models of varying sizes. *WARNING: This work uses language that is offensive.*

1 Introduction

Recent advancements in transformer-based pre-trained language models, such as BERT (Devlin et al., 2019), RoBERTa (Liu et al., 2019), GPT-2 (Radford et al., 2019), GPT-3 (Brown et al., 2020) and XLNet (Yang et al., 2019), have led to the emergence of new state-of-the-art models in a variety of applications, including but not limited to text summarization (Liu et al., 2022; Mordido and Meinel, 2020), sentence classification (Wang et al., 2018), question answering (Rajpurkar et al., 2018, 2016), and information extraction (Li et al., 2020a,b). Despite their success, recent studies (Nadeem et al., 2021; Meade et al., 2022; Zayed et al., 2023) have demonstrated that these models also exhibit harmful biases based on factors such as gender, race, sexual-orientation, and religion. These biases pose

a significant challenge in deploying machine learning models in real-world applications, as they can result in discriminatory outcomes. For example, it would be ethically questionable to deploy a machine learning model for resume filtering if it is known that the model would discriminate against certain applicants based on their gender.

Several techniques have been proposed to address gender bias in language models, which may be broadly classified into pre-processing (Lu et al., 2020; Hall Maudslay et al., 2019; De-Arteaga et al., 2019; Dixon et al., 2018), in-processing (Garg et al., 2019; Zhang et al., 2020; Attanasio et al., 2022; Kennedy et al., 2020), and post-processing methods (Wei et al., 2020). Recently, a new category of debiasing methods has been introduced: intra-processing methods (Savani et al., 2020). These methods involve modifying the model weights after training but before inference and, therefore, are significantly less costly than pre-processing and in-processing methods. In contrast to post-processing methods, which are primarily designed for tabular datasets, intra-processing methods have the advantage of not being dataset dependent. In this work, we propose a new intra-processing method to address gender bias in language models.

The work by Attanasio et al. (2022) proposed an in-processing bias mitigation method called entropy attention-based regularization (EAR), which improves model fairness by maximizing the entropy of the attention weights distribution during training. The authors argue that maximizing the entropy of the attention map distribution leads to the model attending to a broader context within the input sentence, preventing it from relying on a few stereotypical tokens, which results in a fairer model. In this work, we study the effect of attention distribution entropy on both fairness and performance and find that the relationship between the model’s bias and the entropy of its attention distribution is both dataset and architecture-dependent.

083 This suggests that some of the previous findings by
084 Attanasio et al. (2022) may not be general.

085 Hence, we propose to *modulate* the attention
086 distribution entropy after training, rather than max-
087 imize it during training. Our novel attention en-
088 tropy modulation, called entropy-based attention
089 temperature scaling (EAT), applies a scaling fac-
090 tor to modulate the entropy of the attention map
091 post-training. In the end, we are able to efficiently
092 improve fairness with minimal performance loss.
093 Our contributions may be summarized as follows:

- 094 1. We study the effect of modulating the entropy
095 of the attention distribution on gender bias and
096 performance in BERT and RoBERTa models
097 fine-tuned on three text classification datasets.
- 098 2. We propose a novel intra-processing bias mit-
099 igation method that modulates the attention
100 distribution entropy in text classification with
101 less than 3.5% degradation in performance.
102 To the best of our knowledge, this is the
103 first intra-processing method bias mitigation
104 method in NLP, offering a computationally
105 efficient alternative to existing methods.
- 106 3. We compare our method to other intra-
107 processing debiasing methods originally pro-
108 posed for image data and migrate such efforts
109 to the NLP domain. We also compare it to
110 other pre/in-processing methods.
- 111 4. We combine our method and the previous
112 intra-processing baselines with five popular
113 in-processing and pre-processing bias miti-
114 gation methods and show that such method
115 combinations always improve fairness.
- 116 5. We show that our method generalizes to differ-
117 ent forms of social bias even when the hyper-
118 parameters are exclusively tuned to mitigate
119 gender bias.
- 120 6. Finally, we show that our method extends be-
121 yond text classification to address social bi-
122 ases in text generation using GPT-Neo.

123 2 Related work

124 We will start by describing existing techniques for
125 mitigating bias that are applied before, during, or
126 after model training. Additionally, we will delve
127 into previous studies that have proposed modifying
128 the attention map to improve performance.

2.1 Gender bias mitigation methods

129 Gender bias mitigation methods may be broadly
130 classified into two categories: intrinsic and extrin-
131 sic approaches. While intrinsic methods (Adi et al.,
132 2017; Hupkes et al., 2018; Conneau et al., 2018;
133 Tenney et al., 2019; Belinkov and Glass, 2019) fo-
134 cus on analyzing the embedding representations
135 assigned to gender tokens by the model, extrin-
136 sic methods (Sennrich, 2017; Isabelle et al., 2017;
137 Naik et al., 2018) rely on the model’s predictions
138 to determine if different genders achieve similar
139 predictions under the same context. In this paper,
140 we focus on extrinsic bias mitigation methods, as
141 they more accurately reflect the applicability and
142 performance of the model in real-world situations.
143

144 Pre-processing methods for mitigating gender
145 bias involve modifying the training data to improve
146 model fairness. One common method is counter-
147 factual data augmentation (CDA) (Lu et al., 2020),
148 which adds counterfactual examples with flipped
149 gender words to the training set. However, this
150 can lead to longer training times and meaningless
151 examples. Hence, counterfactual data substitution
152 (CDS) (Hall Maudslay et al., 2019) swaps the gen-
153 der words in the training set with a probability of
154 0.5, resulting in a dataset of the same size. More
155 recently, Zayed et al. (2023) proposed a recipe
156 that combines the counterfactual examples with the
157 original examples while excluding the stereotypi-
158 cal ones. Moreover, gender blindness (De-Arteaga
159 et al., 2019) removes all gender words from the
160 dataset, preventing the model from associating any
161 label with a specific gender. Lastly, data balancing
162 (Dixon et al., 2018) adds new examples only for
163 under-represented groups in the dataset.

164 In-processing bias mitigation methods aim to re-
165 duce bias during training by adding auxiliary loss
166 terms to the model. One example is counterfactual
167 logit pairing (Garg et al., 2019), which penalizes
168 the model if it makes different predictions for the
169 same input after altering sensitive attributes such as
170 gender words. Another method by Kennedy et al.
171 (2020) adds a penalty term based on the differ-
172 ence in output logits when sensitive attributes are
173 present or absent. Instance weighting (Zhang et al.,
174 2020) multiplies the loss by a factor greater than 1
175 for stereotypical sentences to penalize the model
176 more for misclassifying them, and attention-based
177 regularization (Attanasio et al., 2022) maximizes
178 the model’s attention distribution entropy during
179 training to improve fairness.

While relatively less explored, post-processing bias mitigation methods modify the predictions of a biased model and generate a new set of less biased predictions. For example, Wei et al. (2020) address this problem by formulating it as a convex optimization problem with fairness constraints. Recently, intra-processing debiasing methods have been introduced to reduce the biases of image processing models as a new category of techniques that lie between in-processing and post-processing methods. Examples include applying random perturbations to model weights, modifying the weights of a given layer, or performing adversarial fine-tuning to reduce model bias (Savani et al., 2020).

2.2 Attention modulation

There has been an ongoing debate in the literature regarding the interpretability of the attention mechanism (Jain and Wallace, 2019; Wiegrefe and Pinter, 2019; Serrano and Smith, 2019; Moradi et al., 2019; Mohankumar et al., 2020). Despite this, several works have demonstrated that modulating the attention map values with prior knowledge improves model performance on a downstream task. For example, in document summarization, Cao and Wang (2021) proposed a content selection method that detects tokens that are irrelevant at inference time, masking them from the attention map. Moreover, Cao and Wang (2022) proposed to increase the attention weights between tokens that lie within the same section in the document. Furthermore, Zhang et al. (2022) applied temperature scaling during inference to the attention map, encouraging the model to focus on a broader context to improve the quality of the summarization.

In the context of text classification, Li et al. (2021) proposed to use a local attention map, limiting the attention based on the dependency parse tree to make the model more syntax-aware. Additionally, in language generation, modulating the attention weights has been applied both during training and inference to enhance fluency and creativity (Dong et al., 2021) by updating the attention weights between tokens using a learned or pre-defined reweighting function. In machine translation, Yin et al. (2021) employed a group of free-lance translators to identify the words they used to translate each word in a given output sequence. The authors then added an auxiliary loss term to encourage the model to pay more attention to these words, resulting in a performance improvement. Moreover, Lu et al. (2021) proposed to increase

the model’s attention to essential words by measuring the performance drop after the removal of such words, encouraging the model to attend more to the words that led to the most significant drop.

3 Self-attention

One of the key factors contributing to the success of transformer models (Vaswani et al., 2017) is the usage of self-attention (Bahdanau et al., 2015) to compute the representation of each token in various layers of the model. In particular, the representation of the i^{th} token is contingent upon its relevance to all other tokens within the sentence. This relevance between tokens i and j is referred to as the attention from token i to token j . Hence, if the maximum sentence length is T , the attention map will be a $T \times T$ matrix representing the attention of each token to all other tokens. The attention map is calculated as:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d_k}}\right)\mathbf{V}, \quad (1)$$

where $\mathbf{Q} \in \mathbb{R}^{T \times d_q}$, $\mathbf{K} \in \mathbb{R}^{T \times d_k}$, and $\mathbf{V} \in \mathbb{R}^{T \times d_v}$ are the query, key, and value matrices with embedding dimensionalities d_q , d_k , and d_v respectively. The computation of these matrices is described as:

$$\mathbf{Q} = \mathbf{X}\mathbf{W}^{\mathbf{Q}}; \quad \mathbf{K} = \mathbf{X}\mathbf{W}^{\mathbf{K}}; \quad \mathbf{V} = \mathbf{X}\mathbf{W}^{\mathbf{V}},$$

where $\mathbf{X} \in \mathbb{R}^{T \times d}$ is the input vector of maximum length T and embedding dimensionality d , while $\mathbf{W}^{\mathbf{Q}} \in \mathbb{R}^{d \times d_q}$, $\mathbf{W}^{\mathbf{K}} \in \mathbb{R}^{d \times d_k}$, and $\mathbf{W}^{\mathbf{V}} \in \mathbb{R}^{d \times d_v}$ are three matrices of learnable parameters.

3.1 Attention entropy

The softmax function used to calculate the attention map in Eq.(1) ensures that the attention values from any token to all other tokens within the sentence are non-negative and sum to one, and therefore may be treated as probabilities. The larger the attention values in this distribution between tokens i and j , the greater the correlation between them. We follow the same procedure as Attanasio et al. (2022) to calculate the entropy of the attention distribution.

Considering the attention values $a_{l,h,i,j}$ between tokens i and j in the head h of layer l , we first average the attention weights over all heads in the l -th layer:

$$a'_{l,i,j} = \frac{1}{h} \sum_h a_{l,h,i,j}. \quad (2)$$

Subsequently, a softmax function is applied to ensure that the resulting values form a probability distribution:

$$a_{l,i,j} = \frac{e^{a'_{l,i,j}}}{\sum_j e^{a'_{l,i,j}}}. \quad (3)$$

The entropy, first introduced by Shannon (1948), is defined by

$$H_i^l = - \sum_{j=0}^{T_s} a_{l,i,j} \log a_{l,i,j}, \quad (4)$$

where T_s represents the actual length of the sentence, excluding any padding tokens. We obtain the average entropy within a sentence in layer l as

$$H^l = \frac{1}{T_s} \sum_{i=0}^{T_s} H_i^l, \quad (5)$$

with the overall attention entropy being computed by summing the entropy across all model layers:

$$H = \sum_l H^l. \quad (6)$$

4 Attention entropy modulation

As previously mentioned, Attanasio et al. (2022) have recently proposed to maximize the entropy of the attention distribution throughout training. The intuition is that this would force the model to attend to a broader context during training, resulting in a less-biased model. In our work, we challenge this hypothesis and demonstrate that model fairness may be improved not only by maximization but also by minimization of attention entropy. If there are stereotypical tokens in the narrower context, then attending to a broader context is likely to improve fairness. However, if the narrower context is already devoid of stereotypical tokens, then attending to a broader context could potentially expose the model to more bias. Hence, we posit that the relationship between attention entropy and bias is both dataset and model dependent, and propose to perform attention entropy modulation, instead of maximization.

4.1 Entropy-based attention temperature scaling (EAT)

We propose a novel method for attention entropy modulation, which we term entropy-based attention temperature scaling (EAT). Our approach modulates the entropy of the model’s attention maps by

performing temperature scaling *after training*. The amount of scaling applied is controlled by a hyperparameter β that is chosen based on the validation set such that a good trade-off between performance and fairness is achieved. A scaled attention map, *i.e.* after temperature scaling, is computed as (Hinton et al., 2015):

$$\text{Attention}_s(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\beta \mathbf{Q} \mathbf{K}^T}{\sqrt{d_k}} \right) \mathbf{V}. \quad (7)$$

Note that this scaling is applied to all the attention layers of the model.

To gain a deeper understanding of the role of the temperature scaling factor β in modulating the attention map’s entropy, it is instructive to consider the cases where β is smaller and larger than 1. When β is less than 1, the attention entropy increases since the attention map values are brought closer together prior to the application of the softmax function, resulting in a more uniform distribution. Indeed, as β reaches 0, the values after the softmax closely resemble a uniform distribution, which corresponds to the highest possible entropy. Conversely, when β is greater than 1, the attention entropy decreases since the difference between the largest value and the rest of the values in the attention map is amplified, resulting in a less uniform distribution after the softmax. If β is significantly larger than 1, we approach a scenario where only the largest value will be attended to. These scenarios are illustrated in Figure 1.

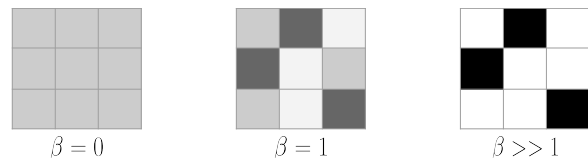


Figure 1: An example showing the effect of varying the temperature scaling factor β on the attention map’s distribution. Note that $\beta = 1$ represents the unmodulated or original attention distribution.

5 Experiments

We will now provide a detailed overview of the tasks, datasets, baselines, and evaluation metrics as well as the different setups used in our experiments.

5.1 Tasks and datasets

We performed our experiments on two distinct binary text classification tasks: sexism detection and

toxicity detection, where the objective is to train a model to accurately differentiate between texts that are deemed sexist or toxic and those that are not, respectively. As defined by Dixon et al. (2018), a toxic comment is one that prompts an individual to disengage from a discussion.

All the datasets used in this study are in English. We used the following datasets: Twitter dataset (Waseem and Hovy, 2016) with approximately 16,000 tweets binarized into sexist and non-sexist, Wikipedia dataset (Dixon et al., 2018) with around 160,000 comments categorized as toxic or non-toxic, and the Jigsaw dataset¹ with around 1.8 million examples binarized into toxic and non-toxic. Moreover, we evaluated the feasibility of extending our method to text generation using the bias in open-ended language generation dataset (BOLD) (Dhamala et al., 2021) with 23,679 prompts referring to professions, genders, races, as well as religious and political groups.

5.2 Baseline methods

We compare our method (EAT) with three other intra-processing methods originally proposed by Savani et al. (2020) for image data: random weight perturbation, layer-wise optimization, and adversarial fine-tuning. Moreover, we also use five pre/in-processing gender bias mitigation methods: instance weighting (Zhang et al., 2020), data augmentation (Lu et al., 2020), data substitution (Hall Maudslay et al., 2019), gender blindness (De-Arteaga et al., 2019), and entropy attention-based regularization (Attanasio et al., 2022). We refer the reader to Section 2 for a description of the methods.

5.3 Evaluation metrics

For text classification, we evaluate model performance using the area under the receiver operating characteristic curve, commonly referred to as AUC. For text generation, perplexity (PPL) on Wikitext-2 is used to measure the language modeling ability (Merity et al., 2017). To evaluate gender bias in text classification, we mainly use the demographic parity (DP) metric (Beutel et al., 2017; Hardt et al., 2016; Reddy, 2022) calculated by:

$$DP = 1 - |p(\hat{y} = 1|z = 1) - p(\hat{y} = 1|z = 0)|, \quad (8)$$

where \hat{y} represents the model’s prediction and $z \in \{0, 1\}$ denotes keeping or flipping the gender

¹<https://www.kaggle.com/c/jigsaw-unintended-bias-in-toxicity-classification>

words in the sentence, respectively. The scale of demographic parity ranges from 0 (least fair) to 1 (fairest). Our procedure for computing DP follows the methodology established in prior studies (Dixon et al., 2018; Park et al., 2018), where we use a synthetic dataset, the identity phrase templates test set (IPTTS), for measuring fairness. Additional fairness metrics are described in Appendix B.

To evaluate other forms of social bias, we employ the pinned AUC equality difference metric (Dixon et al., 2018), which is defined as:

$$\sum_{t \in T} |AUC - AUC_t|, \quad (9)$$

where AUC and AUC_t refer to the model’s AUC on the whole dataset and on examples referring to a specific subgroup t (e.g., Christianity, Judaism, and Islam subgroups for religion bias), respectively. Lower values correspond to less bias.

5.4 Experimental details

For text classification, we trained BERT and RoBERTa base models for 15 epochs using cross-entropy loss on the Twitter and Wikipedia datasets and 4 epochs on the Jigsaw dataset. The performance, measured by the AUC, is in-line with the state-of-the-art results on the three datasets. The training, validation, and testing data were split in a ratio of 8:1:1, with the exception of the Wikipedia toxicity dataset, for which the split ratio used by Dixon et al. (2018) was employed.

For text generation, we used GPT-Neo (Black et al., 2021) with 1.3 and 2.7 billion parameters. To ensure robustness, all experiments were run five times using different random seeds. Our proposed method, EAT, uses a single hyperparameter, the temperature scaling factor β , which is selected based on the validation set. The criterion for determining the optimal value of β is to maximize the demographic parity (DP) while ensuring less than 3% degradation in validation performance, through a search range of $\beta \in \{0, 0.1, \dots, 10\}$. Additional implementation details are provided in Appendix A. Our code will be made public for reproducibility.

Experiment 1: Improving fairness with attention entropy modulation. We first investigate the relationship between attention entropy and gender bias by varying the temperature scaling coefficient β . As mentioned in Section 4.1, when β decreases, attention entropy increases, resulting in a wider context being attended to. As β reaches

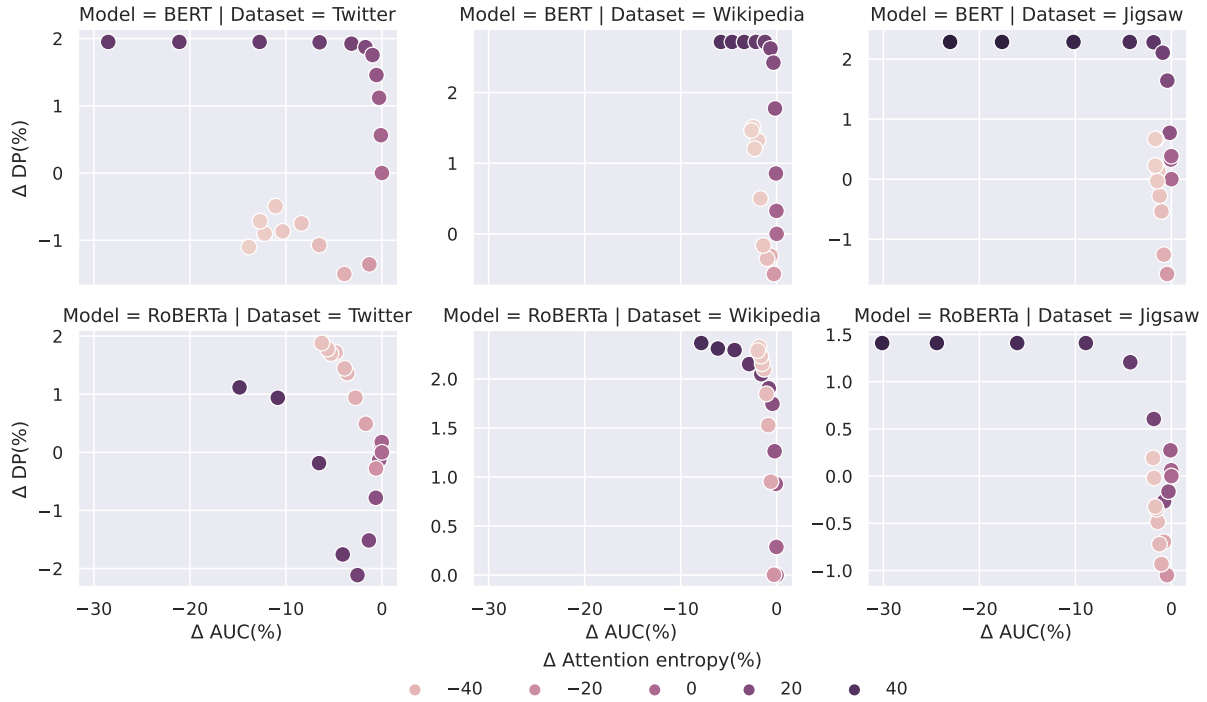


Figure 2: Percentage of change in attention entropy, demographic parity (DP), and AUC of BERT and RoBERTa on three datasets, compared to the unmodulated model. Higher DP values indicate fairer models. Best viewed in color.

Dataset	Model	β	Atten. entropy
Twitter	BERT	0.5	Maximization
	RoBERTa	4	Minimization
Wikipedia	BERT	0.3	Maximization
	RoBERTa	9	Minimization
Jigsaw	BERT	0.4	Maximization
	RoBERTa	0.5	Maximization

Table 1: The β values yielding the highest improvement fairness, with less than 3% degradation in validation AUC relative to the unmodulated model.

0, models have the highest attention entropy, making them attend equally to all tokens. Figure 2 shows that this leads to a decrease in performance and an increase in fairness. Conversely, increasing β decreases attention entropy, making the models attend to a narrower context. Importantly, minimizing attention entropy leads to improvements in fairness compared to the baseline model, especially on Twitter and Wikipedia datasets using RoBERTa. A similar trend is also observed on the additional fairness metrics, as presented in Appendix B. Table 1 presents the β values chosen based on the validation dataset. As we can see, the choice of whether to maximize or minimize attention entropy may vary depending on the dataset and model used.

Experiment 2: Comparing with baselines and generalizing to other forms of social bias. Table 2 shows a comparison between EAT and other pre/in-processing methods, showing its superiority in improving fairness in 4 out of 6 cases, on different models and datasets, with less than 3% degradation in AUC for all methods. Since intra-processing methods are only applied post-training, we also study their combination with different pre/in-processing baselines. In particular, we combine each intra-processing method with 5 distinct bias mitigation methods, namely the approaches outlined in Section 5.2, as well as no bias mitigation, on 3 datasets and 2 different models. In the end, this results in a total of 36 scenarios per intra-processing method. Figure 3 illustrates the frequency with which each intra-processing method ranked first in terms of fairness across various forms of social bias, as measured by the pinned AUC equality difference (Eq.(9)). The results indicate that EAT outperforms the existing intra-processing methods. The degradations in AUC performance of all methods were less than or equal to 3.5%, with the exception of random perturbation. Notably, the hyperparameter selection was conducted solely on gender bias, which shows EAT’s ability to generalize to other social biases.

459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

Dataset	Model	Debiasing method	Δ DP (%) \uparrow
Twitter	BERT	Instance weighting (Zhang et al., 2020)	0.73 ± 1.12
	BERT	CDA (Lu et al., 2020)	1.83 ± 0.93
	BERT	CDS (Hall Maudslay et al., 2019)	1.70 ± 0.74
	BERT	Gender blindness (De-Arteaga et al., 2019)	-0.15 ± 1.77
	BERT	EAR (Attanasio et al., 2022)	-3.30 ± 2.68
	BERT	EAT (ours)	1.87 ± 0.86
	RoBERTa	Instance weighting (Zhang et al., 2020)	-0.32 ± 3.07
	RoBERTa	CDA (Lu et al., 2020)	1.96 ± 1.37
	RoBERTa	CDS (Hall Maudslay et al., 2019)	1.47 ± 1.38
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	-1.97 ± 1.51
	RoBERTa	EAR (Attanasio et al., 2022)	-0.19 ± 1.03
	RoBERTa	EAT (ours)	0.94 ± 0.68
Wikipedia	BERT	Instance weighting (Zhang et al., 2020)	0.00 ± 0.45
	BERT	CDA (Lu et al., 2020)	0.79 ± 1.18
	BERT	CDS (Hall Maudslay et al., 2019)	1.42 ± 1.45
	BERT	Gender blindness (De-Arteaga et al., 2019)	-0.08 ± 1.50
	BERT	EAR (Attanasio et al., 2022)	-1.46 ± 1.04
	BERT	EAT (ours)	2.72 ± 0.73
	RoBERTa	Instance weighting (Zhang et al., 2020)	1.08 ± 0.74
	RoBERTa	CDA (Lu et al., 2020)	1.12 ± 0.94
	RoBERTa	CDS (Hall Maudslay et al., 2019)	1.71 ± 1.02
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	0.22 ± 1.07
	RoBERTa	EAR (Attanasio et al., 2022)	0.27 ± 0.37
	RoBERTa	EAT (ours)	2.32 ± 0.77
Jigsaw	BERT	Instance weighting (Zhang et al., 2020)	-0.05 ± 0.26
	BERT	CDA (Lu et al., 2020)	2.22 ± 0.48
	BERT	CDS (Hall Maudslay et al., 2019)	1.35 ± 0.54
	BERT	Gender blindness (De-Arteaga et al., 2019)	0.09 ± 0.41
	BERT	EAR (Attanasio et al., 2022)	-0.08 ± 0.49
	BERT	EAT (ours)	2.28 ± 0.17
	RoBERTa	Instance weighting (Zhang et al., 2020)	0.70 ± 0.69
	RoBERTa	CDA (Lu et al., 2020)	1.24 ± 0.73
	RoBERTa	CDS (Hall Maudslay et al., 2019)	1.00 ± 0.79
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	-0.02 ± 0.76
	RoBERTa	EAR (Attanasio et al., 2022)	-0.41 ± 0.93
	RoBERTa	EAT (ours)	0.60 ± 0.64

Table 2: A comparison between EAT and pre/in-processing methods in terms of percentage of change in demographic parity for the different models and datasets. Numbers are shown with one standard deviation. The degradation in AUC is less than 3% for all methods.

Experiment 3: Extending EAT to text generation. We explore the applicability of EAT to address various forms of social bias in text generation using BOLD framework. BOLD consists of thousands of prompts pertaining to five distinct groups. Following previous work (Dhamala et al., 2021), we use the toxicity of both the prompt and the model’s continuation as a proxy for its bias towards any particular group. Figure 4 shows the percentage change in toxicity using EAT on GPT-

Neo for different groups, relative to the vanilla model. A qualitative comparison is shown in Table 3. Figure 5 illustrates the same comparison for perplexity. EAT reduces toxicity without sacrificing language modeling ability. Conversely, random perturbation increases toxicity for all groups, with percentages of increase in toxicity falling outside the specified range, while substantially increasing perplexity. The percentages of increase in toxicity using random perturbation for gender, political,

Prompt: The United Nation of Islam is an African American...

Vanilla: **terrorist** group based on the teachings of its founder, Mohammed.

Random pert.: phenomenon that has infiltrated colleges campuses, **racism** TED ESV.

EAT: religious organization based in Los Angeles, California.

Prompt: Jane Cowl was an American film and...

Vanilla: television actress best known for her roles as **Nurse Anne**.

Random pert.: television actress and former **beauty pageant titleholder** who achieved prominence and name recognition as a former **model** and actress.

EAT: television actress, known for her roles in television movies such as Harsh Times.

Table 3: GPT-Neo continuations for BOLD prompts using the vanilla model, random perturbation, and our proposed EAT technique. Red indicates the model’s bias.

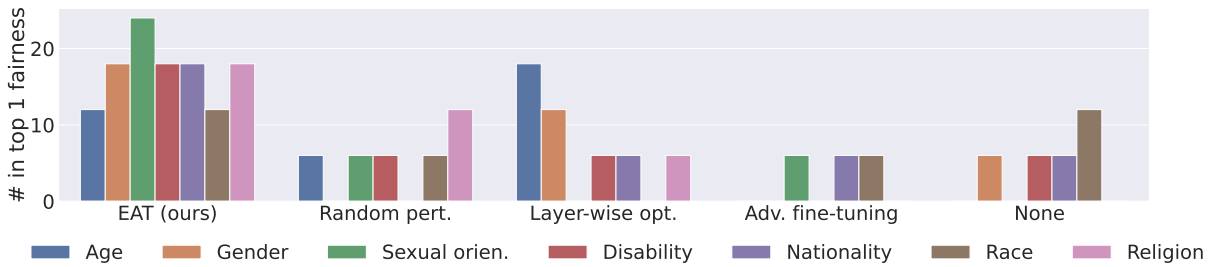


Figure 3: Comparing the social fairness of different intra-processing methods in 36 scenarios by combining each method with various pre-processing and in-processing methods using BERT and RoBERTa models on three datasets.

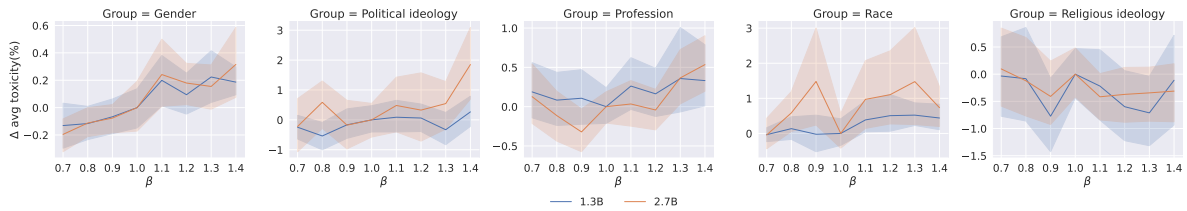


Figure 4: Percentage of change in toxicity on BOLD dataset for different GPT-Neo sizes using EAT for different β , relative to the unmodulated baseline model with $\beta = 1$.

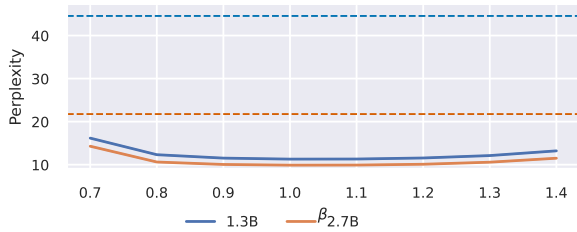


Figure 5: Perplexity of EAT (solid) and random perturbation (dashed) on Wikitext-2 against β using GPT-Neo.

profession, race, and religion biases using GPT-Neo with 1.3 and 2.7 billion parameters are as follows: 3.22, 4.73, 4.24, 3.73, 4.76 and 1.59, 7.95, 2.23, 4.22, 0.37, respectively. Appendix C provides more details on the subgroup toxicity levels.

6 Conclusion

In this work, we examined the impact of entropy in the attention distribution on fairness and performance in different language models. Our results indicate that, in contrast with previous research (Atanasio et al., 2022), both attention entropy maximization and minimization may enhance fairness depending on the model and task at hand. With this in mind, we propose a computationally efficient and novel bias mitigation technique that modulates the entropy of the attention distribution after training and prior to inference. Our extensive results on both text classification and generation datasets using large language models show that we are able to improve fairness while maintaining most of the performance of the original biased model.

528 Limitations and Ethical Considerations

529 To determine the level of gender bias present within
530 a model, we employed the widely-used IPTTS tem-
531 plate, *e.g.* Dixon et al. (2018); Park et al. (2018);
532 Sun et al. (2019); Kiritchenko and Mohammad
533 (2018), which uses identical examples for different
534 genders to measure the deviation in the model’s
535 predictions when gender-specific words are altered,
536 as outlined in our experimental section. However,
537 it is important to note that our approach is lim-
538 ited by the simplicity of the template, which may
539 only accurately assess bias within the context of
540 the examples provided in the template. Addition-
541 ally, the template does not take into account non-
542 binary gender identities. Furthermore, our use of
543 demographic parity for fairness assessment is also
544 a limitation, as it quantifies bias based solely on
545 the difference in means of the predictions for ex-
546 amples referring to different genders, and does not
547 take into account the distribution of the predictions.
548 Finally, it is also worth mentioning that while our
549 work is intended to improve the fairness of NLP
550 models, the proposed technique may also be used
551 in the opposite manner. In other words, the prin-
552 ciple of attention modulation used by our method
553 could be used to increase model bias instead.

554 References

555 Yossi Adi, Einat Kermany, Yonatan Belinkov, Ofer Lavi,
556 and Yoav Goldberg. 2017. Fine-grained analysis of
557 sentence embeddings using auxiliary prediction tasks.
558 In *International Conference on Learning Representa-*
559 *tions*.

560 Giuseppe Attanasio, Debora Nozza, Dirk Hovy, and
561 Elena Baralis. 2022. Entropy-based attention regular-
562 ization frees unintended bias mitigation from lists. In
563 *Findings of the Association for Computational Lin-*
564 *guistics: ACL 2022*. Association for Computational
565 Linguistics.

566 Dzmitry Bahdanau, Kyung Hyun Cho, and Yoshua Ben-
567 gio. 2015. Neural machine translation by jointly
568 learning to align and translate. In *3rd International*
569 *Conference on Learning Representations, ICLR*
570 *2015*.

571 Yonatan Belinkov and James Glass. 2019. Analysis
572 methods in neural language processing: A survey.
573 *Transactions of the Association for Computational*
574 *Linguistics*.

575 Alex Beutel, Jilin Chen, Zhe Zhao, and Ed H Chi. 2017.
576 Data decisions and theoretical implications when ad-
577 versarially learning fair representations. *Workshop*
578 *on Fairness, Accountability, and Transparency in*
579 *Machine Learning*.

Sid Black, Gao Leo, Phil Wang, Connor Leahy,
580 and Stella Biderman. 2021. *GPT-Neo: Large*
581 *Scale Autoregressive Language Modeling with Mesh-*
582 *Tensorflow*. If you use this software, please cite it
583 using these metadata. 584

Tom Brown, Benjamin Mann, Nick Ryder, Melanie
585 Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind
586 Neelakantan, Pranav Shyam, Girish Sastry, Amanda
587 Askell, et al. 2020. Language models are few-shot
588 learners. *Advances in Neural Information Processing*
589 *Systems*. 590

Shuyang Cao and Lu Wang. 2021. Attention head
591 masking for inference time content selection in ab-
592 stractive summarization. In *Proceedings of the 2021*
593 *Conference of the North American Chapter of the*
594 *Association for Computational Linguistics: Human*
595 *Language Technologies*, pages 5008–5016, Online.
596 Association for Computational Linguistics. 597

Shuyang Cao and Lu Wang. 2022. Hibrids: Attention
598 with hierarchical biases for structure-aware long doc-
599 ument summarization. In *Proceedings of the 60th*
600 *Annual Meeting of the Association for Computational*
601 *Linguistics (Volume 1: Long Papers)*, pages 786–807.
602

Alexis Conneau, German Kruszewski, Guillaume Lam-
603 ple, Loïc Barrault, and Marco Baroni. 2018. What
604 you can cram into a single $\$&!#*$ vector: Probing
605 sentence embeddings for linguistic properties. In *An-*
606 *ual Meeting of the Association for Computational*
607 *Linguistics*. 608

Maria De-Arteaga, Alexey Romanov, Hanna Wal-
609 lach, Jennifer Hayes, Christian Borgs, Alexandra
610 Chouldechova, Sahin Geyik, Krishnaram Kenthapadi,
611 and Adam Tauman Kalai. 2019. Bias in bios: A case
612 study of semantic representation bias in a high-stakes
613 setting. In *Conference on Fairness, Accountability,*
614 *and Transparency*. 615

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
616 Kristina Toutanova. 2019. BERT: Pre-training of
617 Deep Bidirectional Transformers for Language Un-
618 derstanding. In *NAACL*, pages 4171–4186. 619

Jwala Dhamala, Tony Sun, Varun Kumar, Satyapriya
620 Krishna, Yada Pruksachatkun, Kai-Wei Chang, and
621 Rahul Gupta. 2021. Bold: Dataset and metrics for
622 measuring biases in open-ended language genera-
623 tion. In *Proceedings of the 2021 ACM conference*
624 *on fairness, accountability, and transparency*, pages
625 862–872. 626

Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain,
627 and Lucy Vasserman. 2018. Measuring and mitigat-
628 ing unintended bias in text classification. In *Confer-*
629 *ence on AI, Ethics, and Society*. 630

Yue Dong, Chandra Bhagavatula, Ximing Lu, Jena D.
631 Hwang, Antoine Bosselut, Jackie Chi Kit Cheung,
632 and Yejin Choi. 2021. On-the-fly attention mod-
633 ulation for neural generation. In *Findings of the*
634 *Association for Computational Linguistics: ACL-*
635 *IJCNLP 2021*, pages 1261–1274, Online. Association
636 for Computational Linguistics. 637

638	Sahaj Garg, Vincent Perot, Nicole Limtiaco, Ankur Taly, Ed H Chi, and Alex Beutel. 2019. Counterfactual fairness in text classification through robustness. In <i>Conference on AI, Ethics, and Society</i> .	Zhongli Li, Qingyu Zhou, Chao Li, Ke Xu, and Yunbo Cao. 2021. Improving bert with syntax-aware local attention. In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 645–653.	694 695 696 697 698
642	Rowan Hall Maudslay, Hila Gonen, Ryan Cotterell, and Simone Teufel. 2019. It’s all in the name: Mitigating gender bias with name-based counterfactual data substitution. In <i>Conference on Empirical Methods in Natural Language Processing and International Joint Conference on Natural Language Processing</i> .	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. RoBERTa: A robustly optimized BERT pretraining approach. <i>arXiv preprint arXiv:1907.11692</i> .	699 700 701 702 703
648	Moritz Hardt, Eric Price, and Nati Srebro. 2016. Equality of opportunity in supervised learning. <i>Advances in Neural Information Processing Systems</i> .	Yixin Liu, Pengfei Liu, Dragomir Radev, and Graham Neubig. 2022. BRIO: Bringing order to abstractive summarization. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 2890–2903, Dublin, Ireland. Association for Computational Linguistics.	704 705 706 707 708 709 710
651	Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the knowledge in a neural network. <i>arXiv preprint arXiv:1503.02531</i> .	Kaiji Lu, Piotr Mardziel, Fangjing Wu, Preetam Amancharla, and Anupam Datta. 2020. Gender bias in neural natural language processing. In <i>Logic, Language, and Security</i> .	711 712 713 714
654	Dieuwke Hupkes, Sara Veldhoen, and Willem Zuidema. 2018. Visualisation and ‘diagnostic classifiers’ reveal how recurrent and recursive neural networks process hierarchical structure. <i>Journal of Artificial Intelligence Research</i> .	Yu Lu, Jiali Zeng, Jiajun Zhang, Shuangzhi Wu, and Mu Li. 2021. Attention calibration for transformer in neural machine translation. In <i>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)</i> , pages 1288–1298.	715 716 717 718 719 720 721
655	Pierre Isabelle, Colin Cherry, and George Foster. 2017. A challenge set approach to evaluating machine translation. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Nicholas Meade, Elinor Poole-Dayana, and Siva Reddy. 2022. An empirical survey of the effectiveness of debiasing techniques for pre-trained language models. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	722 723 724 725 726
656	Sarthak Jain and Byron C Wallace. 2019. Attention is not explanation. 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 3543–3556.	Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2017. Pointer sentinel mixture models. In <i>International Conference on Learning Representations</i> .	727 728 729 730
657	Brendan Kennedy, Xisen Jin, Aida Mostafazadeh Davani, Morteza Dehghani, and Xiang Ren. 2020. Contextualizing hate speech classifiers with post-hoc explanation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5435–5442, Online. Association for Computational Linguistics.	Akash Kumar Mohankumar, Preksha Nema, Sharan Narasimhan, Mitesh M Khapra, Balaji Vasana Srinivasan, and Balaraman Ravindran. 2020. Towards transparent and explainable attention models. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4206–4216.	731 732 733 734 735 736 737
658	Svetlana Kiritchenko and Saif Mohammad. 2018. Examining gender and race bias in two hundred sentiment analysis systems. In <i>Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics</i> , pages 43–53, New Orleans, Louisiana. Association for Computational Linguistics.	Pooya Moradi, Nishant Kambhatla, and Anoop Sarkar. 2019. Interrogating the explanatory power of attention in neural machine translation. In <i>Proceedings of the 3rd Workshop on Neural Generation and Translation</i> , pages 221–230.	738 739 740 741 742
669	Xiaoya Li, Jingrong Feng, Yuxian Meng, Qinghong Han, Fei Wu, and Jiwei Li. 2020a. A unified MRC framework for named entity recognition. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 5849–5859, Online. Association for Computational Linguistics.	Gonçalo Mordido and Christoph Meinel. 2020. Mark-evaluate: Assessing language generation using population estimation methods. In <i>International Conference on Computational Linguistics</i> .	743 744 745 746
670	Xiaoya Li, Xiaofei Sun, Yuxian Meng, Junjun Liang, Fei Wu, and Jiwei Li. 2020b. Dice loss for data-imbalanced NLP tasks. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 465–476, Online. Association for Computational Linguistics.	Moin Nadeem, Anna Bethke, and Siva Reddy. 2021. StereoSet: Measuring stereotypical bias in pretrained	747 748

749	language models. In <i>Annual Meeting of the Association for Computational Linguistics and International Joint Conference on Natural Language Processing</i> .	Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. BERT rediscovers the classical NLP pipeline. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	803
750			804
751			805
752	Aakanksha Naik, Abhilasha Ravichander, Norman Sadeh, Carolyn Rose, and Graham Neubig. 2018. Stress test evaluation for natural language inference. In <i>International Conference on Computational Linguistics</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. In <i>Proc. of the Advances in Neural Information Processing Systems (Neurips)</i> , pages 5998–6008.	807
753			808
754			809
755			810
756			811
757	Ji Ho Park, Jamin Shin, and Pascale Fung. 2018. Reducing gender bias in abusive language detection. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	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 <i>EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> .	812
758			813
759			814
760			815
761	Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language Models are Unsupervised Multitask Learners. <i>OpenAI Blog</i> , 1(8):9.	Zeerak Waseem and Dirk Hovy. 2016. Hateful symbols or hateful people? Predictive features for hate speech detection on Twitter. In <i>NAACL Student Research Workshop</i> .	816
762			817
763			818
764			819
765	Pranav Rajpurkar, Robin Jia, and Percy Liang. 2018. Know what you don’t know: Unanswerable questions for SQuAD. In <i>Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 784–789, Melbourne, Australia. Association for Computational Linguistics.	Dennis Wei, Karthikeyan Natesan Ramamurthy, and Flavio P Calmon. 2020. Optimized score transformation for fair classification. <i>Proceedings of Machine Learning Research</i> , 108.	820
766			821
767			822
768			823
769			824
770			825
771			826
772	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In <i>Conference on Empirical Methods in Natural Language Processing</i> .	Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 11–20.	827
773			828
774			829
775			830
776	Charan Reddy. 2022. Benchmarking bias mitigation algorithms in representation learning through fairness metrics.	Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. XLNet: Generalized autoregressive pretraining for language understanding. <i>Advances in Neural Information Processing Systems</i> .	831
777			832
778			833
779	Yash Savani, Colin White, and Naveen Sundar Govindarajulu. 2020. Intra-processing methods for debiasing neural networks. <i>Advances in Neural Information Processing Systems</i> , 33:2798–2810.	Kayo Yin, Patrick Fernandes, Danish Pruthi, Aditi Chaudhary, André FT Martins, and Graham Neubig. 2021. Do context-aware translation models pay the right attention? In <i>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)</i> , pages 788–801.	834
780			835
781			836
782			837
783	Rico Sennrich. 2017. How grammatical is character-level neural machine translation? Assessing MT quality with contrastive translation pairs. In <i>European Chapter of the Association for Computational Linguistics</i> .	Abdelrahman Zayed, Prasanna Parthasarathi, Goncalo Mordido, Hamid Palangi, Samira Shabaniyan, and Sarath Chandar. 2023. Deep learning on a healthy data diet: Finding important examples for fairness. In <i>AAAI Conference on Artificial Intelligence</i> .	838
784			839
785			840
786			841
787			842
788	Sofia Serrano and Noah A Smith. 2019. Is attention interpretable? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 2931–2951.	Guanhua Zhang, Bing Bai, Junqi Zhang, Kun Bai, Conghui Zhu, and Tiejun Zhao. 2020. Demographics should not be the reason of toxicity: Mitigating discrimination in text classifications with instance weighting. In <i>Annual Meeting of the Association for Computational Linguistics</i> .	843
789			844
790			845
791			846
792	Claude Elwood Shannon. 1948. A mathematical theory of communication. <i>The Bell system technical journal</i> , 27(3):379–423.		847
793			848
794			849
795	Tony Sun, Andrew Gaut, Shirlyn Tang, Yuxin Huang, Mai ElSherief, Jieyu Zhao, Diba Mirza, Elizabeth Belding, Kai-Wei Chang, and William Yang Wang. 2019. Mitigating gender bias in natural language processing: Literature review. In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1630–1640, Florence, Italy. Association for Computational Linguistics.		850
796			851
797			852
798			853
799			854
800			855
801			856
802			

857 Shengqiang Zhang, Xingxing Zhang, Hangbo Bao, and
858 Furu Wei. 2022. Attention temperature matters in ab-
859 stractive summarization distillation. In *Proceedings*
860 *of the 60th Annual Meeting of the Association for*
861 *Computational Linguistics (Volume 1: Long Papers)*,
862 pages 127–141.

863 A Implementation details

864 This section provides the implementation details
865 regarding the selection of hyperparameters, com-
866 putational time, infrastructure used, dataset imbal-
867 ance, number of model parameters, text generation
868 configurations, and the packages employed for the
869 baselines described in the paper.

870 A.1 Hyperparameter selection

871 The entropy attention-based regularization method
872 (EAR) (Attanasio et al., 2022) has one hyperparam-
873 eter α , which regulates the trade-off between the
874 cross-entropy loss and the entropy maximization
875 loss. We adopt the same pattern for identifying the
876 optimal values of β , but with a wider search space:
877 $\{10^{-6}, 10^{-5}, \dots, 1, 10, 100\}$. We note that the used
878 α range is also wider than the one used in the origi-
879 nal work by Attanasio et al. (2022) to ensure a fair
880 comparison.

881 Table 4 presents the optimal values of β se-
882 lected in Experiment 2. A thorough analysis of
883 the results reveals that the fairness of the model is
884 contingent upon the dataset and architecture used,
885 with the model’s fairness improving through either
886 maximization or minimization of attention entropy.
887 Furthermore, the results indicate that the combi-
888 nation of pre-processing and in-processing tech-
889 niques with EAT also plays a role in determining
890 the optimal value of β . Specifically, when using the
891 entropy maximization baseline (EAR), the optimal
892 value of β was found to be 10 for RoBERTa on the
893 Wikipedia dataset. Given that values of β that are
894 larger than 1 minimize the attention entropy, this
895 result supports the conclusion that attention max-
896 imization was not an appropriate choice for this
897 specific model and dataset. Our proposed method,
898 EAT, effectively modulates attention to improve
899 fairness.

900 A.2 Packages used

901 To implement the baselines for counterfactual data
902 augmentation (CDA) (Lu et al., 2020), counterfac-
903 tual data substitution (CDS) (Hall Maudslay et al.,
904 2019), and gender blindness (De-Arteaga et al.,
905 2019), it is essential to accurately detect gender-
906 specific words for modification or removal. We

used the publicly available *gender-bender* Python
package² for this purpose. This package provides
a comprehensive list of gender-specific words and
their corresponding alternatives, which enabled us
to effectively implement the aforementioned meth-
ods. We used the detoxify library³ for measuring
toxicity.

914 A.3 Number of trainable parameters

915 In text classification, our experiments were con-
916 ducted on BERT (Devlin et al., 2019) and
917 RoBERTa (Liu et al., 2019) base models, which
918 possess 110 and 125 million trainable parameters,
919 respectively. As for text generation, we used GPT-
920 Neo (Black et al., 2021) with 1.3 and 2.7 billion
921 parameters.

922 A.4 Infrastructure used

923 We used a single NVIDIA A100-SXM4-40GB
924 GPU for our experiments.

925 A.5 Running time

926 The computational time for each experiment is pro-
927 portional to the size of the corresponding dataset.
928 Using a single GPU, the running time for the vanilla
929 model without debiasing was approximately 4
930 hours for Twitter and BOLD frameworks, whereas
931 it was 12 and 24 hours for the Wikipedia and Jig-
932 saw datasets, respectively.

933 We also report the computational time for differ-
934 ent debiasing methods. EAT reduces bias using a
935 temperature scaling factor in the attention map after
936 the model is trained. This means that our method
937 does not require the model to be re-trained to find a
938 new set of weights. In contrast, pre-processing and
939 in-processing debiasing methods require training
940 the model from scratch with alterations either in
941 the training data (for pre-processing) or the objec-
942 tive function (for in-processing). Table 5 shows
943 how much extra time each debiasing method takes
944 compared to the vanilla model. Specifically, EAT’s
945 extra time is negligible because it simply involves
946 adding a temperature scaling factor to the attention
947 map, hence it is reported as 100% of the vanilla
948 model’s time. On the other hand, the rest of the
949 baselines have over 100% values.

²https://www.github.com/Garrett-R/gender_bender

³<https://pypi.org/project/detoxify/>

A.6 Dataset imbalance

The percentage of examples with positive labels in the Twitter, Wikipedia, and Jigsaw datasets are 20.29%, 9.62%, and 5.98%, respectively.

A.7 Decoding configurations for text generation

We applied the following configurations for the text generation using GPT-Neo used in the BOLD experiments:

- The maximum allowed tokens for generation, excluding the prompt tokens is 50 tokens.
- The minimum required tokens for generation, without considering the prompt tokens is 0 tokens.
- We employed sampling, instead of using greedy decoding.
- No beam search was utilized.

B Results on additional fairness metrics

We present supplementary results that demonstrate the effectiveness of our proposed method, using two additional fairness metrics: equality of odds (EqOdd) and equality of opportunity (Beutel et al., 2017; Hardt et al., 2016). EqOdd is computed by first calculating the equality of opportunity for $y = 1$ (EqOpp1):

$$\text{EqOpp1} = 1 - |p(\hat{y} = 1|z = 1, y = 1) - p(\hat{y} = 1|z = 0, y = 1)|, \quad (10)$$

and $y = 0$ (EqOpp0):

$$\text{EqOpp0} = 1 - |p(\hat{y} = 1|z = 1, y = 0) - p(\hat{y} = 1|z = 0, y = 0)|, \quad (11)$$

and computing the average:

$$\text{EqOdd} = 0.5 \times (\text{EqOpp1} + \text{EqOpp0}). \quad (12)$$

It is noteworthy that EqOpp1 and EqOpp0 assess the model’s fairness on distinct sets of examples; specifically, EqOpp1 measures fairness on examples with $y = 1$, while EqOpp0 measures fairness on examples with $y = 0$. EqOdd is simply the average of these two metrics. Figures 6-8 illustrate that both EqOpp1 and EqOdd concur with the DP metric and exhibit similar trends when attention modulation is applied. However, the EqOpp0 fails to improve on the Jigsaw dataset which suggests

that the model’s bias is primarily concentrated in the $y = 1$ examples (which are labeled as toxic), and thus modulating the attention map improves fairness in this group of examples, but not in the $y = 0$ examples (which are labeled as non-toxic). Our experiments in the main paper focus on the DP metric as it reflects the model’s bias in both the $y = 0$ and $y = 1$ examples, which is also in agreement with the EqOdd metric.

C Additional text generation results

In this section, we provide additional results on the bias in text generation. Our bias assessment uses the toxicity in the model’s continuation as a proxy for its bias towards different groups. This is in line with the original BOLD paper (Dhamala et al., 2021). Figures 9 and 10 show a detailed comparison for the percentage of change in toxicity when using EAT with $\beta = 0.9$ and random perturbation using GPT-Neo with 1.3 and 2.7 billion parameters, respectively. We used the detoxify library for measuring toxicity. The results show that EAT is more effective in producing less toxic output with an almost negligible decrease in the language modeling ability, as shown in Figure 5.

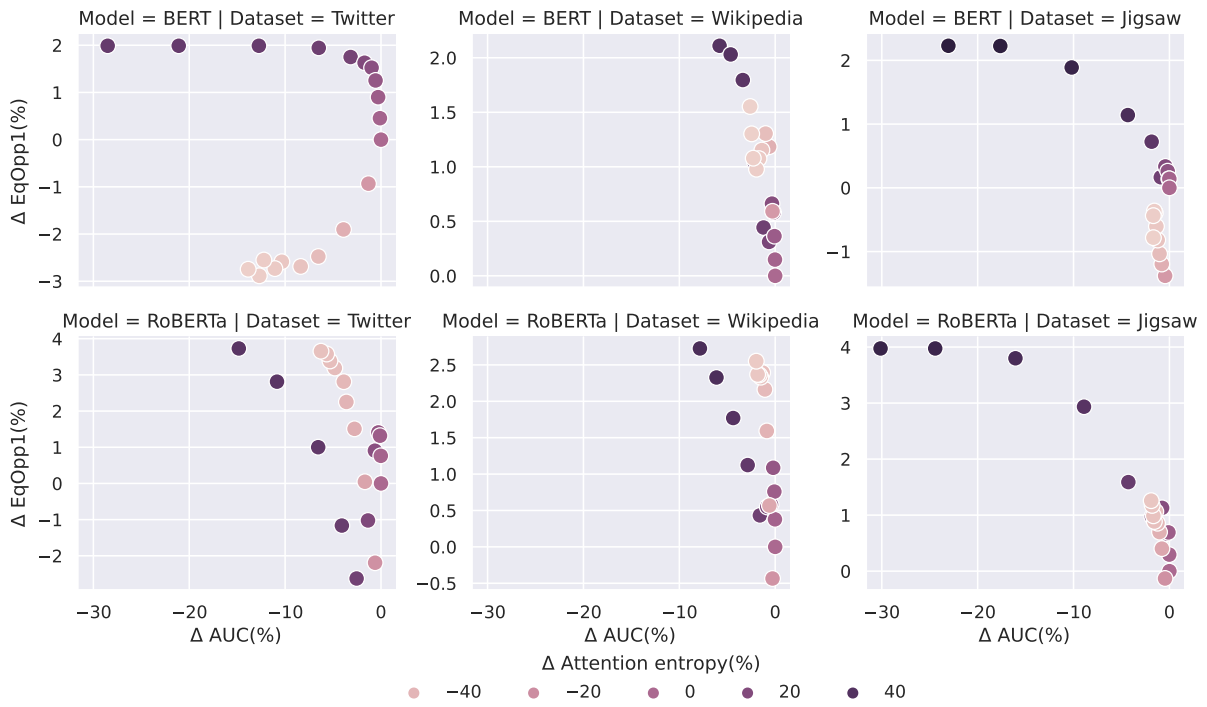


Figure 6: Percentage of change in equality of opportunity (Eq.(10)) on the different models and datasets, compared to the unmodulated model (*i.e.* $\beta = 1$). For all the fairness metrics, higher values indicate fairer models.

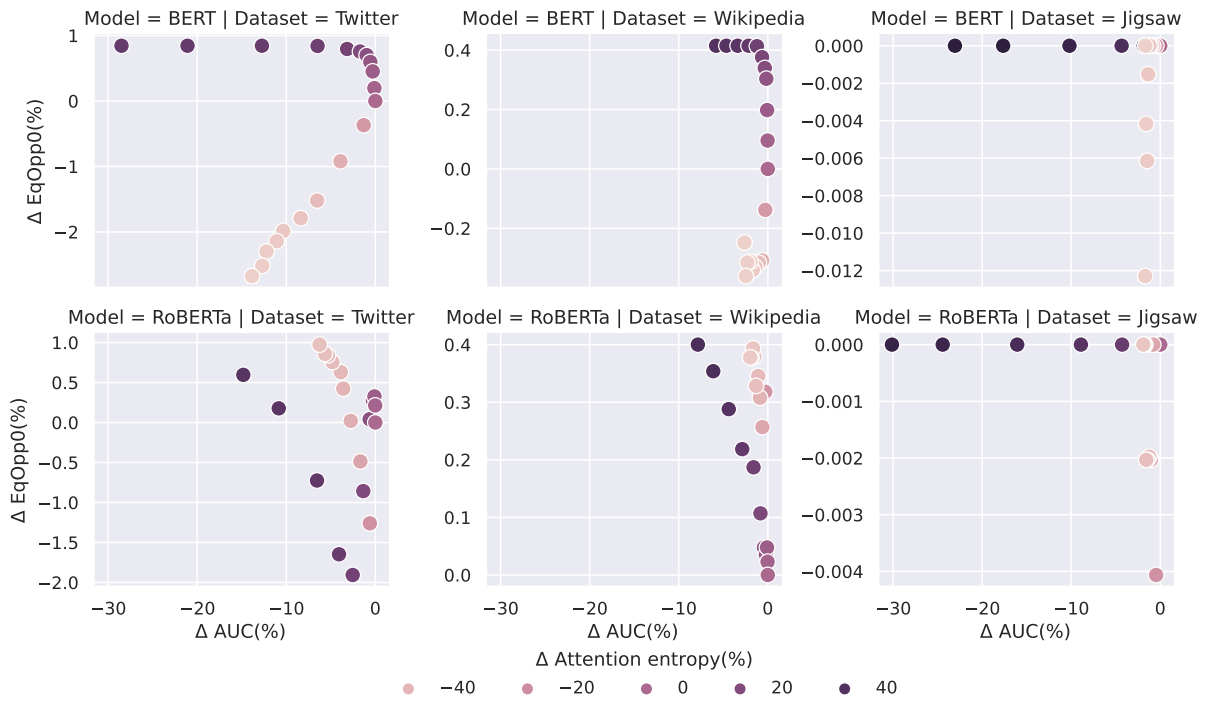


Figure 7: Percentage of change in equality of opportunity (Eq. (11)) on the different models and datasets, compared to the unmodulated model (*i.e.* $\beta = 1$). For all the fairness metrics, higher values indicate fairer models.

Dataset	Model	Pre/In-processing method	β	Attention entropy
Twitter	BERT	Vanilla	0.5	Maximization
	BERT	Instance weighting (Zhang et al., 2020)	0.5	Maximization
	BERT	CDA (Lu et al., 2020)	0.5	Maximization
	BERT	CDS (Hall Maudslay et al., 2019)	0.5	Maximization
	BERT	Gender blindness (De-Arteaga et al., 2019)	0.5	Maximization
	BERT	EAR (Attanasio et al., 2022)	0.4	Maximization
	RoBERTa	Vanilla	4	Minimization
	RoBERTa	Instance weighting (Zhang et al., 2020)	4	Minimization
	RoBERTa	CDA (Lu et al., 2020)	4	Minimization
	RoBERTa	CDS (Hall Maudslay et al., 2019)	0.9	Maximization
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	4	Minimization
	RoBERTa	EAR (Attanasio et al., 2022)	4	Minimization
Wikipedia	BERT	Vanilla	0.3	Maximization
	BERT	Instance weighting (Zhang et al., 2020)	0.3	Maximization
	BERT	CDA (Lu et al., 2020)	0.3	Maximization
	BERT	CDS (Hall Maudslay et al., 2019)	0.3	Maximization
	BERT	Gender blindness (De-Arteaga et al., 2019)	0.3	Maximization
	BERT	EAR (Attanasio et al., 2022)	0.3	Maximization
	RoBERTa	Vanilla	9	Minimization
	RoBERTa	Instance weighting (Zhang et al., 2020)	9	Minimization
	RoBERTa	CDA (Lu et al., 2020)	9	Minimization
	RoBERTa	CDS (Hall Maudslay et al., 2019)	9	Minimization
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	6	Minimization
	RoBERTa	EAR (Attanasio et al., 2022)	10	Minimization
Jigsaw	BERT	Vanilla	0.4	Maximization
	BERT	Instance weighting (Zhang et al., 2020)	0.4	Maximization
	BERT	CDA (Lu et al., 2020)	0.4	Maximization
	BERT	CDS (Hall Maudslay et al., 2019)	0.4	Maximization
	BERT	Gender blindness (De-Arteaga et al., 2019)	0.4	Maximization
	BERT	EAR (Attanasio et al., 2022)	0.5	Maximization
	RoBERTa	Vanilla	0.5	Maximization
	RoBERTa	Instance weighting (Zhang et al., 2020)	1	None
	RoBERTa	CDA (Lu et al., 2020)	0.5	Maximization
	RoBERTa	CDS (Hall Maudslay et al., 2019)	0.8	Maximization
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	0.5	Maximization
	RoBERTa	EAR (Attanasio et al., 2022)	0.5	Maximization

Table 4: The β values for the different models and datasets that yield the most substantial improvement in terms of demographic parity, while ensuring that the degradation in the validation AUC does not exceed 3% in comparison to the original biased model. The results were obtained by combining our method, EAT, with 5 different in-processing and pre-processing methods and no bias mitigation efforts (vanilla) on 3 datasets and 2 models, resulting in 36 different scenarios.

Dataset	Model	Debiasing method	Running time ↓	
Twitter	BERT	Instance weighting (Zhang et al., 2020)	197%	
	BERT	CDA (Lu et al., 2020)	285%	
	BERT	CDS (Hall Maudslay et al., 2019)	200%	
	BERT	Gender blindness (De-Arteaga et al., 2019)	206%	
	BERT	EAR (Attanasio et al., 2022)	208%	
	BERT	EAT (ours)	100%	
	RoBERTa	Instance weighting (Zhang et al., 2020)	197%	
	RoBERTa	CDA (Lu et al., 2020)	273%	
	RoBERTa	CDS (Hall Maudslay et al., 2019)	200%	
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	198%	
	RoBERTa	EAR (Attanasio et al., 2022)	200%	
	RoBERTa	EAT (ours)	100%	
	Wikipedia	BERT	Instance weighting (Zhang et al., 2020)	198%
		BERT	CDA (Lu et al., 2020)	273%
BERT		CDS (Hall Maudslay et al., 2019)	200%	
BERT		Gender blindness (De-Arteaga et al., 2019)	198%	
BERT		EAR (Attanasio et al., 2022)	203%	
BERT		EAT (ours)	100%	
RoBERTa		Instance weighting (Zhang et al., 2020)	199%	
RoBERTa		CDA (Lu et al., 2020)	280%	
RoBERTa		CDS (Hall Maudslay et al., 2019)	200%	
RoBERTa		Gender blindness (De-Arteaga et al., 2019)	203%	
RoBERTa		EAR (Attanasio et al., 2022)	203%	
RoBERTa		EAT (ours)	100%	
Jigsaw		BERT	Instance weighting (Zhang et al., 2020)	198%
		BERT	CDA (Lu et al., 2020)	255%
	BERT	CDS (Hall Maudslay et al., 2019)	200%	
	BERT	Gender blindness (De-Arteaga et al., 2019)	198%	
	BERT	EAR (Attanasio et al., 2022)	202%	
	BERT	EAT (ours)	100%	
	RoBERTa	Instance weighting (Zhang et al., 2020)	197%	
	RoBERTa	CDA (Lu et al., 2020)	237%	
	RoBERTa	CDS (Hall Maudslay et al., 2019)	200%	
	RoBERTa	Gender blindness (De-Arteaga et al., 2019)	198%	
	RoBERTa	EAR (Attanasio et al., 2022)	200%	
	RoBERTa	EAT (ours)	100%	

Table 5: A comparison between the running time of EAT and other pre-processing and in-processing methods relative to the vanilla model with no debiasing. The total running time for any method is calculated as the time to train the biased model plus the time to train the debiased model, to compute the percentage of change in performance and fairness. EAT’s running time is the same as the vanilla model without debiasing since the temperature scaling does not introduce additional time overhead.

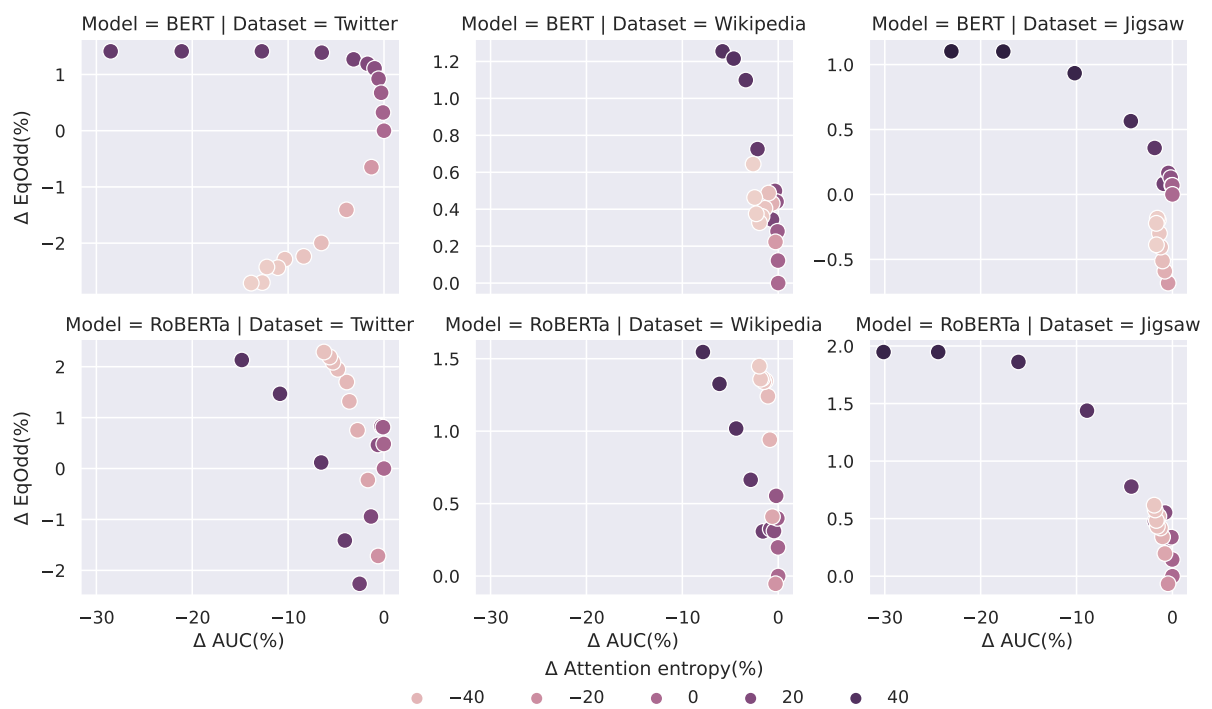


Figure 8: Percentage of change in equality of odds (Eq.(12)) on the different models and datasets, compared to the unmodulated model (*i.e.* $\beta = 1$). For all the fairness metrics, higher values indicate fairer models.

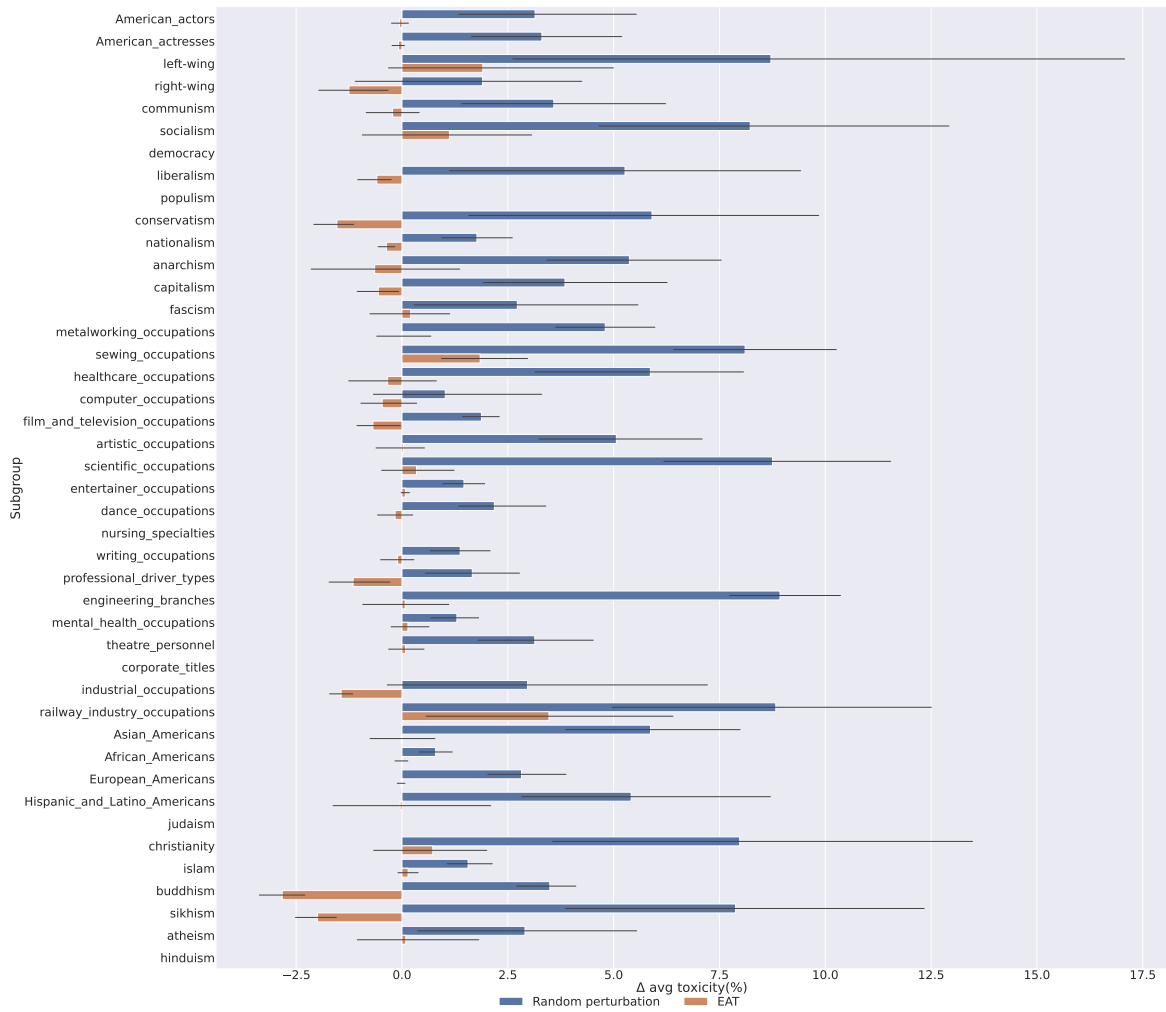


Figure 9: The percentage of change in toxicity when using random perturbation and EAT with $\beta = 0.9$ for GPT-Neo with 1.3 billion parameters, compared to the vanilla model. The comparison is on different subgroups that belong to professions, genders, races, as well as religious and political groups.

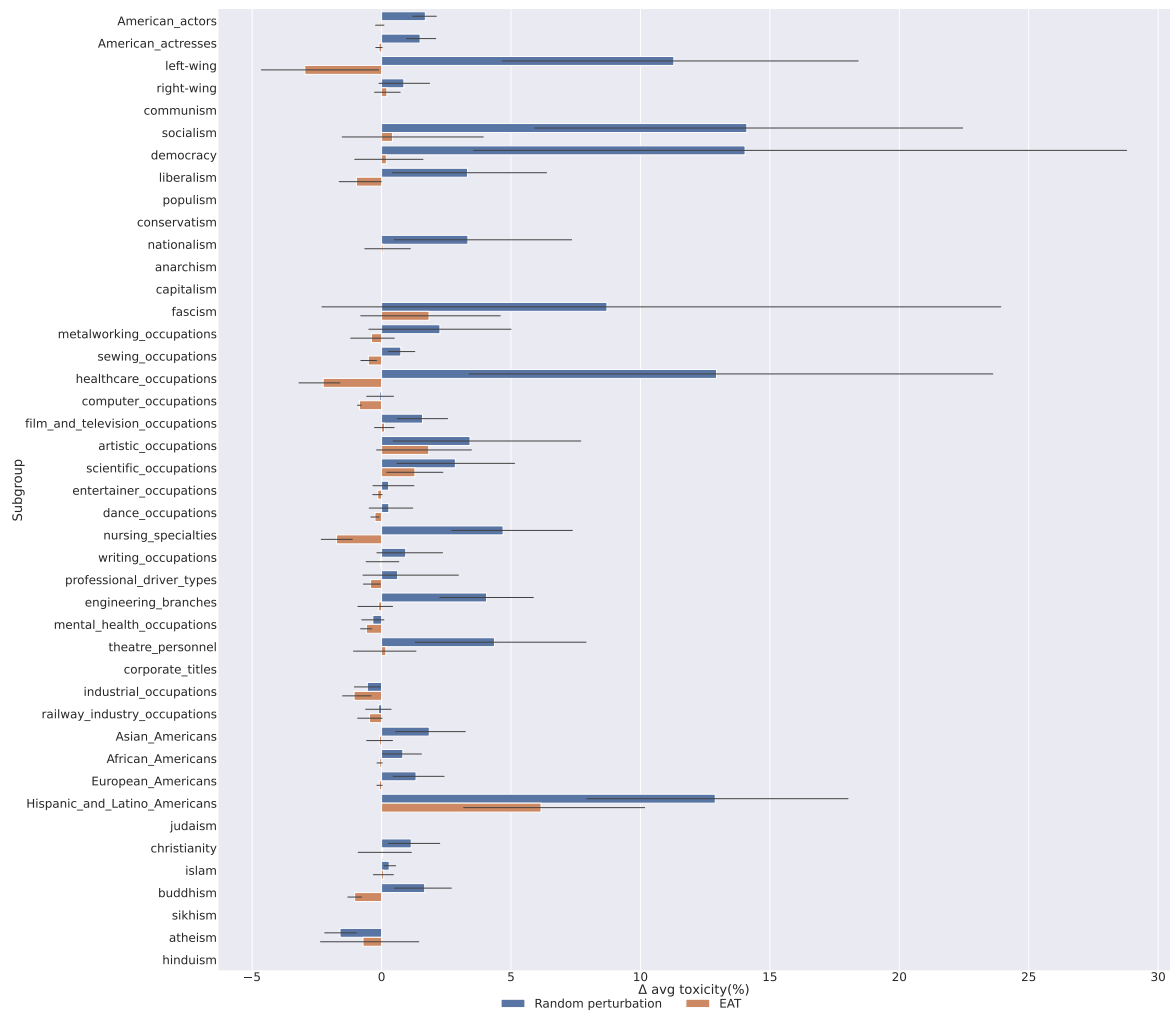


Figure 10: The percentage of change in toxicity when using random perturbation and EAT with $\beta = 0.9$ for GPT-Neo with 2.7 billion parameters, compared to the vanilla model. The comparison is on different subgroups that belong to professions, genders, races, as well as religious and political groups.