

# Modular and On-demand Bias Mitigation with Attribute-Removal Subnetworks

Lukas Hauzenberger, Shahed Masoudian, Deepak Kumar,  
Markus Schedl, Navid Rekabsaz

Johannes Kepler University Linz, Austria  
Linz Institute of Technology, AI Lab  
{first\_name.family\_name}@jku.at

## Abstract

Societal biases are reflected in large pre-trained language models and their fine-tuned versions on downstream tasks. Common in-processing bias mitigation approaches, such as adversarial training and mutual information removal, introduce additional optimization criteria, and update the model to reach a new debiased state. However, in practice, end-users and practitioners might prefer to switch back to the original model, or apply debiasing only on a specific subset of protected attributes. To enable this, we propose a novel modular bias mitigation approach, consisting of stand-alone highly sparse debiasing subnetworks, where each debiasing module can be integrated into the core model on-demand at inference time. Our approach draws from the concept of *diff* pruning, and proposes a novel training regime adaptable to various representation disentanglement optimizations. We conduct experiments on three classification tasks with gender, race, and age as protected attributes. The results show that our modular approach, while maintaining task performance, improves (or at least remains on-par with) the effectiveness of bias mitigation in comparison with baseline finetuning. Particularly on a two-attribute dataset, our approach with separately learned debiasing subnetworks shows effective utilization of either or both the subnetworks for selective bias mitigation.

## 1 Introduction

A large body of research evidences the existence of societal biases and stereotypes in pre-trained language models (PLMs) (Zhao et al., 2019; Sheng et al., 2019; Rekabsaz et al., 2021), and their potential harms when used in down-stream tasks (Blodgett et al., 2020; De-Arteaga et al., 2019; Rekabsaz and Schedl, 2020; Stanovsky et al., 2019). Common in-processing approaches to bias mitigation update a model’s (typically all) parameters to satisfy specific attribute erasure criteria through optimization methods such as adversarial train-

ing (Elazar and Goldberg, 2018; Rekabsaz et al., 2021), and mutual information reduction (Colombo et al., 2021). These methods are shown to be effective in reducing the footprint of protected attributes (e.g., gender, race, etc.) in the resulting model.

However when using such debiasing models in practice and in specific use-cases, system designers or end-users might still prefer to instead use the original model, or a debiased variation in respect to a particular subset of protected attributes. This can be due to various reasons such as the nature of a given input, preference of an individual end-user, or fairness-utility trade-off considerations. For instance, while a bias-aware model should indeed be agnostic to genders when the input is about gender-neutral occupations (such as *nurse* or *CEO*), certain topics like *pregnancy* may specifically require gender information for a correct model decision.<sup>1</sup> Also as shown in previous studies (Zerveas et al., 2022; Biega et al., 2018; Rekabsaz et al., 2021), since improving fairness on specific tasks may come with the cost of performance degradation, it is necessary to provide on-demand control over whether to impose fairness/debiasing criteria. Using existing approaches, this would require maintaining and deploying multiple large parallel models for every protected attribute, resulting in overly complex and resource-heavy pipelines and increased latency.

To address this, we introduce a novel modular bias mitigation approach using sparse weight-difference networks. In our approach, the required changes in a model’s parameters for erasing a bias attribute are stored in a decoupled subnetwork, trained simultaneously by a debiasing and a sparsification objective. At inference time, adding each debiasing module to the core model results in delivering debiasing qualities to a model’s prediction in respect to the corresponding protected attribute. Our approach extends the principle idea

<sup>1</sup>See also the discussion in Krieg et al. (2023) about the need to separate bias-sensitive queries from “normal” ones.

of *diff pruning* (Guo et al., 2021) introduced for parameter-efficient task training to bias mitigation, by viewing the objective of erasing a protected attribute as a stand-alone *diff* module. This module replaces fine-tuning by training only a small set of parameters added to the corresponding PLM’s parameters, to deliver bias mitigation of a specific protected attribute. We further propose a novel procedure to train such debiasing subnetworks separately, and to selectively add an arbitrary set of them to the core model at inference time (§3).

Our approach can be applied to any debiasing and representation disentanglement method, provided that its objective has separate learning signals for the task and each protected attribute. In comparison with adapter networks (Rebuffi et al., 2017; Houlsby et al., 2019), as shown by Guo et al. (2021) and also evidenced in our experiments, even more parameter-efficiency can be provided. Additionally, since our approach extends the base model with *diff* subnetworks, the resulting model is expected to perform (at least) as good as the fine-tuning variation, avoiding possible performance degradations. The modularity of our approach supports separating the process of developing debiasing solutions for a task from using them, such that stand-alone debiasing modules can be created and shared, and later be utilized in a final system on-demand.

We evaluate our approach on three bias mitigation tasks: occupation prediction from biographies involving gender (De-Arteaga et al., 2019), hate speech detection with dialect-based race as sensitive attribute (Founta et al., 2018); and mention prediction in tweets with two attributes of gender and age of the authors (Pardo et al., 2016). The last dataset particularly enables the study of combining independently trained debiasing modules (details in §4). The evaluation results show that our approach, due to learning the subnetworks specialized on the narrow functionality of debiasing an attribute, provides on par or better debiasing performance in comparison with strong baselines. Additionally, we observe that on the mention detection task, learning the debiasing subnetworks post-hoc to model training provides effective results when combining the two (independently trained) subnetworks at inference time. Remarkably, these results are achieved with debiasing subnetworks of maximum 1% size of the core model (BERT-Base), in some cases (e.g., for gender attribute) even only 0.01% (details in §5).

## 2 Related Work

### 2.1 Parameter-efficient and Modular Training

The concept of subnetworks is grounded in the lottery ticket hypothesis (Frankle and Carbin, 2019), stating that in deep neural networks, one can find many sparse subnetworks with a capacity comparable to that of the base network. Zhou et al. (2019) show that spotting such subnetworks through binary masks can in fact be seen as a form of model training. Zhao et al. (2020) further consider the magnitude of the parameters for pruning and filter out the ones lower than a specific threshold. Guo et al. (2021) use  $L_0$ -regularization (Louizos et al., 2018) to reduce the number of active neurons. Hu et al. (2021) propose low-rank adaptation via rank decomposition matrices. These methods are extended by structural pruning approaches with additional architectural constraints, commonly led to a higher parameter-efficiency (Lagunas et al., 2021; Jaszczur et al., 2021).

Recent studies exploit subnetworks in the context of multi-task learning. Wortsman et al. (2020) isolate the learning signal of each task in a separate masking subnetwork, while Ben-Zaken et al. (2021) only finetune the bias weights. Xu et al. (2021) learn a subset of parameters via masking out the gradients of other parameters during backward-pass. Guo et al. (2021) suggest learning a sparse *diff* subnetwork for each task, whose parameter values are added to the corresponding parameters of the base network. An alternative architecture are adapter networks, first introduced in the context of multi-task learning (Rebuffi et al., 2017; Houlsby et al., 2019; Stickland and Murray, 2019), and are then extended in respect to their parameter efficiency (Rücklé et al., 2021; Han et al., 2021a), architectural variations (Mahabadi et al., 2021), and transfer learning capacity (Pfeiffer et al., 2021). As stated by (Sung et al., 2021), adapters (in their original form) slightly increase the inference cost of a model in comparison with the original form or pruning-based variations. Our proposed approach contributes to this line of research by extending this concept to on-demand bias mitigation.

### 2.2 Fairness & Bias Mitigation in NLP

Several studies explore methods for debiasing PLMs, such as linearly projecting embeddings into the space with minimum correlations to protected features (Ravfogel et al., 2020; Kaneko and Bolukbasi, 2021; Bolukbasi et al., 2016), utilizing a

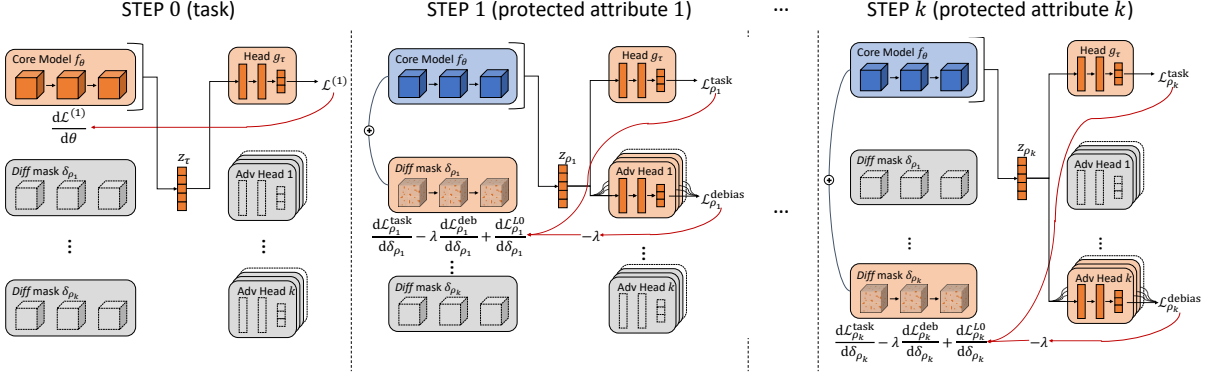


Figure 1: Training procedure of MODDIFFY-PAR in one batch with adversarial bias removal for  $k$  protected attributes. The color orange indicates the trainable parameters during each step, while blue shows the frozen ones. The core model learns the task in Step 0, and the subnetworks — each responsible for debiasing a protected attribute — are trained in the next steps. In MODDIFFY-POST training regime, after finetuning the core model in Step 0, the task head parameters (together with the ones of the core model) remain frozen in the next steps.

distribution alignment loss (Guo et al., 2022), or penalizing bias by utilizing the encoded information in models (Schick et al., 2021). Adversarial training, originally introduced in domain adaptation (Ganin et al., 2016; Ganin and Lempitsky, 2015) is utilized in the context of fair representation learning (Xie et al., 2017; Madras et al., 2018), and later to erase demographic data from text classifiers (Elazar and Goldberg, 2018; Barrett et al., 2019; Han et al., 2021b; Wang et al., 2021), information retrieval models (Rekabsaz et al., 2021), and recommendation systems (Ganhör et al., 2022). Mutual information removal is an alternative approach, which minimizes the approximate upper bound of the common information between the task and protected attributes (Cheng et al., 2020; Colombo et al., 2021). Our work utilizes these optimizations to learn a novel modular on-demand bias mitigation approach.

Few recent studies explore parameter-efficient training for bias mitigation. Lauscher et al. (2021) approach debiasing PLMs using a stack of adapters. While shown effective in practice, the adapters in the higher levels inherently depend on the ones in the lower levels and cannot be learned nor utilized stand-alone. Zhang et al. (2021) approach bias mitigation with binary masks applied to the base network. More recently and in the context of removing spurious shortcuts in natural language understanding datasets, Meissner et al. (2022) train sparse binary masks on a finetuned model. Our work extends this line of research by encapsulating concept erasure modules in separate *diff* subnetworks for each protected attribute, and selectively

applying them to the base model at inference time.

### 3 Modular Debiasing with *Diff* Subnets

We start with defining the general approach to model bias mitigation. We consider an arbitrary PLM denoted by  $f_\theta$  with the set of parameters  $\theta$ . The model learns the task  $\tau$  using the loss function  $\mathcal{L}_\tau$ . The predictions of  $f_\theta$  might be sensitive to the variations in any of the  $k$  protected attributes of the set  $P = \{\rho_1, \dots, \rho_k\}$ . The bias mitigation objective is to make the model invariant to these variations while maintaining the effectiveness on the task, approached by defining the debiasing loss  $\mathcal{L}_{\rho_i}$  for the protected attribute  $\rho_i$ . We discuss two realizations of this loss function in Section 3.2. A debiased model in respect to attributes  $P$  is achieved by training on the following loss function:

$$\mathcal{L}_{total} = \mathcal{L}_\tau + \mathcal{L}_{\rho_1} + \dots + \mathcal{L}_{\rho_k} \quad (1)$$

Using  $\mathcal{L}_{total}$ , one can finetune all parameters of the model (Elazar and Goldberg, 2018; Rekabsaz et al., 2021), or utilize any parameter-efficient training such as adapters (Lauscher et al., 2021), *diff* pruning (Guo et al., 2021), or binary masks (Zhang et al., 2021). We use some of these methods as baselines, explained in the following sections. In the remainder of this section, we introduce our *Modular Debiasing with Diff Subnetworks* (MODDIFFY), explain its training and inference procedure, and describe two debiasing optimization methods.

#### 3.1 MODDIFFY

We aim to encapsulate the debiasing functionality of the protected attribute  $\rho_i$ , provided by the signal

of the corresponding loss  $\mathcal{L}_{\rho_i}$ , into the sparse *diff* subnetwork characterized by the set of parameters  $\delta_{\rho_i}$ . Each parameter in  $\delta_{\rho_i}$  corresponds to a parameter in  $\theta$ , such that adding  $\delta_{\rho_i}$  to the corresponding parameters in  $\theta$  results in debiasing the model  $f$ .

To learn the model, let us consider a training data item in form of  $\langle x, y_\tau, y_{\rho_1}, \dots, y_{\rho_k} \rangle$ , where  $x$  is the input,  $y_\tau$  the task label, and  $y_{\rho_i}$  denotes the label of the corresponding protected attribute  $\rho_i$ . Figure 1 depicts the MODDIFFY approach with adversarial bias mitigation optimization. We first optimize for the task by encoding  $x$  into the vector  $z_\tau$  using  $f(\cdot; \theta)$ . An arbitrary decoder network denoted as  $g_\tau$  uses  $z_\tau$  to predict the task output, and task loss is calculated using cross entropy (CE), as formulated below:

$$z_\tau = f(x; \theta), \quad \mathcal{L}_{total}^{(0)} = \text{CE}(g_\tau(z_\tau), y_\tau) \quad (2)$$

This loss updates  $\theta$  as well as the parameters of  $g_\tau$ . While in our experiments we opt for fully finetuning  $f_\theta$ , in practice the model can be trained using any parameter-efficient method.

Next, we iterate over the protected attributes, and in each step learn the corresponding debiasing *diff* subnetwork. Specifically in step  $i$  dedicated to the protected attribute  $\rho_i$ , we learn the set of sparse additional parameters  $\delta_{\rho_i}$  such that by being added to  $\theta$ , the resulting  $f(\cdot; \theta + \delta_{\rho_i})$  model is debiased. Learning  $\delta_{\rho_i}$  is characterized by three loss functions. The first loss,  $\mathcal{L}_{\rho_i}^{\text{task}}$  maintains the task performance when the task output is predicted from the altered encoder formulated below:

$$z_{\rho_i} = f(x; \theta + \delta_{\rho_i}), \quad \mathcal{L}_{\rho_i}^{\text{task}} = \text{CE}(g_\tau(z_{\rho_i}), y_\tau) \quad (3)$$

The second is the debiasing loss  $\mathcal{L}_{\rho_i}^{\text{debias}}$ , defined based on  $z_{\rho_i}$  and the label of the corresponding protected attribute  $y_{\rho_i}$ . We discuss adversarial bias removal and mutual information reduction as two possible realization of this representation disentanglement loss in Section 3.2. The third loss imposes the sparsity constraint, defined as the  $L_0$  regularization of  $\delta_{\rho_i}$ . The  $L_0$  loss aims to reduce the number of non-zero parameters, namely the term  $\sum_{j=1}^{|\delta_{\rho_i}|} \mathbb{1}\{\delta_{\rho_i,j} \neq 0\}$ , and is realized with the differentiable approximation proposed by Louizos et al. (2018). Following Guo et al. (2021),  $\delta_{\rho_i}$  is decomposed into the element-wise multiplication of two sets of parameters:  $\delta_{\rho_i} = \mathbf{m}_{\rho_i} \odot \mathbf{w}_{\rho_i}$ . The parameter set  $\mathbf{w}_{\rho_i}$  stores the magnitude changes (*diff* values), and  $\mathbf{m}_{\rho_i}$  learns to mask out the parameters.

$\mathbf{m}_{\rho_i}$  is characterized by the hard concrete distribution (Guo et al., 2021) with  $(\log \alpha_{\rho_i}, 1)$  parameters, and  $\gamma < 0$  and  $\zeta > 1$  hyperparameters. This results in the following formulation:

$$\mathcal{L}_{\rho_i}^{L_0} = \sum_{j=1}^{|\delta_{\rho_i}|} \sigma \left( \log \alpha_{\rho_i,j} - \log \left( -\frac{\gamma}{\zeta} \right) \right) \quad (4)$$

where  $\sigma$  denotes the sigmoid function, and as in Guo et al. (2021)  $\beta$  is set to 1. The masking network can be simply reduced by assigning a mask to a group of parameters (such as a weight matrix, or a layer) instead of each individual one.

Putting all together, the objective of MODDIFFY at step  $i > 0$  is defined as:

$$\mathcal{L}_{total}^{(i)} = \mathcal{L}_{\rho_i}^{\text{task}} + \mathcal{L}_{\rho_i}^{\text{debias}} + \mathcal{L}_{\rho_i}^{L_0} \quad (5)$$

We should note that while each subnetwork is trained independently from the others, the output embeddings  $z_{\rho_i}$  are passed to the same decoder network  $g_\tau$ . This design choice forces the learned embeddings to remain in the same embedding space, making it possible to add multiple (independently trained) subnetworks together to the core model. Another aspect is that the sparsity rate of each resulting subnetwork is not fixed and may vary due to the factors such as the hyperparameter setting, and the extent of encoded information content. To achieve a fixed sparsity rate, we further apply magnitude pruning by only keeping a fixed portion of the parameters with the largest absolute values, and finetuning the resulting parameters with the task and debiasing loss terms.

We conduct optimization with  $\mathcal{L}_{total}$  under two training regimes. In the first one referred to as MODDIFFY-PAR, we repeat the mentioned training steps for each training batch. The second approach referred to as MODDIFFY-POST trains debiasing subnetworks post-hoc to training the core model. While MODDIFFY-PAR accommodates for a higher flexibility in optimization by training the networks in parallel, MODDIFFY-POST provides the practical benefits of learning various debiasing solutions for an already trained core model.

Finally at *inference time*, one can use the core model in its original form  $f(x; \theta)$ , or in combination with any of the debiasing subnetworks in the form of  $f(x; \theta + \delta_{\rho_i})$ . Our approach also enables simultaneously debiasing all or any subset of the protected attributes by adding their corresponding subnetworks to the core network, for instance



as in  $f(x; \theta + \delta_{\rho_1} + \dots + \delta_{\rho_k})$ . We should note that, although these subnetworks reside in the same distributional space, they might affect each others functionality, depending on the inherent nature of the biases as well as the possible correlations between the protected attributes. We will examine this case in the next sections on a dataset with two protected attributes of gender and age.

### 3.2 Bias Mitigation Objectives

We explain two bias mitigation optimization methods used in MODDIFFY in the following.

**Adversarial Bias Removal** This method first defines a new classification head  $h_{\rho_i}$  for each protected attribute  $\rho_i$ . This head receives  $z_{\rho_i}$  as input, predicts the corresponding protected attribute, and calculates the cross entropy loss function  $\mathcal{L}_{\rho_i}^{\text{debias}}$ . This loss needs to remove the information of  $\rho_i$  from  $f$  but train  $h_{\rho_i}$  such that it can effectively predict the protected attribute. This optimization forms the min-max game:  $\mathcal{L}_{\rho_i}^{\text{debias}} = \min_f \max_{h_{\rho_i}} \text{CE}(h_{\rho_i}(z_{\rho_i}), y_{\rho_i})$ . A common approach to turn this loss into a minimization problem is by using gradient reversal layer (GRL) (Ganin and Lempitsky, 2015) added before the debiasing heads. GRL multiplies the gradient of  $\mathcal{L}_{\rho_i}^{\text{debias}}$  with a factor of  $-\lambda_i$ , and thereby simplifies the learning process to a standard gradient-based optimization, formulated below:

$$\mathcal{L}_{\rho_i}^{\text{debias}} = \min_{f, h_{\rho_i}} \text{CE}(h_{\rho_i}(z_{\rho_i}), y_{\rho_i})$$

**Mutual Information (MI) Reduction** This approach represents a family of algorithms that aims to remove the mutual information of the encoded embeddings of the task and protected attributes. Maximum Mean Discrepancy (MMD), first introduced in the context of domain adaptation (Gretton et al., 2012; Tzeng et al., 2014), offers a realization of MI reduction by minimizing the ability to separate the subsets belonging to two protected attributes. In particular, given a set of data points  $X$  split into two subsets  $X_{\rho_i}^A$  and  $X_{\rho_i}^B$  according to the values of the (binary) protected attribute  $\rho_i$ , MMD minimizes the distance between the encoded embeddings of the subgroups with the following loss formulation:

$$\mathcal{L}_{\rho_i}^{\text{debias}} = \left( \frac{\sum_{x^A \in X_{\rho_i}^A} \phi(f(x^A))}{|X_{\rho_i}^A|} - \frac{\sum_{x^B \in X_{\rho_i}^B} \phi(f(x^B))}{|X_{\rho_i}^B|} \right)^2$$

where  $\phi$  is the feature map kernel defined as a linear combination of multiple Gaussian kernels.

## 4 Experiment Design

**Datasets** We evaluate our approach on three datasets on the tasks of occupation prediction, hate speech detection, and mention prediction; involving protected attributes of gender, age, and race dialect. The first dataset is **BIOS** (De-Arteaga et al., 2019) which contains short biographies used to predict a person’s job, where the name and any indication of the person’s gender (such as pronouns) in the biography are omitted. The BIOS dataset contains around 430K data points with 28 occupations, and two protected attribute classes (female/male). The second dataset is **FDCL18** (Founta et al., 2018) for hate speech detection, containing a set of tweets each classified as hateful, abusive, spam, or none. As discussed in Xia et al. (2020), hate speech might have a strong correlation with dialect-based racial bias. Following previous studies (Sap et al., 2019; Ravfogel et al., 2020), we assign race dialect labels of *African American* and *White American* to FDCL18 using the probabilistic model developed by Blodgett et al. (2016), resulting in the dataset of approximately 62K data points. The third dataset is **PAN16** (Rangel et al., 2016) containing a set of tweets accompanied with the labels of gender and age of the authors. The task’s objective is to predict mentions (whether another user is mentioned in a tweet). PAN16 provides approximately 200K data points with binary task classes (*mention, no mention*), as well as two gender labels and five age groups. Further details on the three datasets are provided in Appendix A.

**Models and Baselines** We conduct the experiments on the following models and baselines. **FINETUNE**: finetuning all parameters of the PLM on the task without any bias mitigation objective. **FINETUNE-DEBIAS**: the same model as FINETUNE but with the bias mitigation objective. **ADAPTER**: learning the task with an adapter network while the rest of the PLM’s parameters are kept frozen. **ADAPTER-DEBIAS**: the same model as ADAPTER but the adapter is trained on both task and bias mitigation objectives. **DIFFPRUN**: using a *diff* network to learn the task while PLM parameters remain unchanged. **DIFFPRUN-DEBIAS**: same as DIFFPRUN but the *diff* network is trained on both task and bias mitigation objectives. **MODDIFFY-POST**: our introduced post-hoc approach

Model	BIOS (gender)				FCDL18 (race-dialect)			
	Adversarial		MI Reduction		Adversarial		MI Reduction	
	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$
FINETUNE	84.1 <sub>0.3</sub>	67.2 <sub>0.7</sub>	84.1 <sub>0.3</sub>	67.2 <sub>0.7</sub>	82.0 <sub>0.4</sub>	93.0 <sub>0.4</sub>	82.0 <sub>0.4</sub>	93.0 <sub>0.4</sub>
ADAPTER	84.4 <sub>0.1</sub>	65.9 <sub>0.1</sub>	84.4 <sub>0.1</sub>	65.9 <sub>0.1</sub>	80.9 <sub>0.1</sub>	82.1 <sub>4.1</sub>	80.9 <sub>0.1</sub>	82.1 <sub>4.1</sub>
FINETUNE-DEBIAS	84.2 <sub>0.1</sub>	56.9 <sub>0.9</sub>	84.1 <sub>0.7</sub>	61.7 <sub>0.8</sub>	81.9 <sub>0.7</sub>	84.5 <sub>3.9</sub>	81.6 <sub>0.4</sub>	87.3 <sub>0.5</sub>
ADAPTER-DEBIAS	84.6 <sub>0.2</sub>	60.8 <sub>0.4</sub>	84.4 <sub>0.0</sub>	65.6 <sub>0.2</sub>	80.5 <sub>0.7</sub>	65.6 <sub>0.1</sub> <sup>♣</sup>	80.1 <sub>0.4</sub>	82.1 <sub>1.3</sub> <sup>♣</sup>
DIFFPRUN	84.6 <sub>0.1</sub>	68.9 <sub>0.2</sub>	84.6 <sub>0.1</sub>	68.9 <sub>0.2</sub>	81.3 <sub>0.3</sub>	93.2 <sub>0.3</sub>	81.3 <sub>0.3</sub>	93.2 <sub>0.3</sub>
DIFFPRUN-DEBIAS	84.5 <sub>0.1</sub>	62.4 <sub>0.3</sub>	84.2 <sub>0.4</sub>	63.2 <sub>1.4</sub>	81.6 <sub>0.4</sub>	66.8 <sub>3.7</sub>	81.3 <sub>0.2</sub>	91.8 <sub>0.2</sub>
MODDIFFY-POST	84.5 <sub>0.1</sub>	61.6 <sub>0.7</sub>	84.3 <sub>0.1</sub>	64.5 <sub>0.4</sub>	81.3 <sub>0.2</sub>	<b>66.0</b> <sub>1.2</sub>	81.2 <sub>0.2</sub>	91.1 <sub>0.9</sub>
MODDIFFY-PAR	84.2 <sub>0.2</sub>	<b>53.7</b> <sub>1.5</sub> <sup>♣</sup>	84.5 <sub>0.2</sub>	<b>58.8</b> <sub>1.2</sub> <sup>♣</sup>	81.2 <sub>0.6</sub>	75.4 <sub>3.7</sub>	81.3 <sub>0.7</sub>	<b>85.5</b> <sub>0.4</sub>

Table 1: Results of the BIOS and FCDL18 datasets on **BERT-Base** with adversarial bias removal and mutual information (MI) reduction. Task performance is measured with accuracy, and bias mitigation with balanced accuracy of the probes. The protected attribute is gender for BIOS, and race-dialect for FCDL18. The results with the best bias mitigation performance (lowest values) among the models that use *diff* subnetworks (lower part of the table) are shown in **bold**, and among all models with the ♣ symbol. Subscript values indicate standard deviation.

where the debiasing modules are learned after training the model. We use FINETUNE as the base model for MODDIFFY-POST, whose parameters are kept frozen during post-hoc training. **MODDIFFY-PAR**: our introduced parallel approach where the debiasing modules are learned together with the base model finetuned on the task. Additionally, to provide a comprehensive view on bias mitigation methods, we evaluate the datasets on the **INLP** (Ravfogel et al., 2020) approach using the implementation and suggested hyperparameter setting. The evaluation results of the INLP method are separately reported in Table 8 in Appendix B. As the PLM encoder for all models, we using two versions of BERT (Devlin et al., 2019) with different sizes, namely BERT-Mini (Turc et al., 2019) and BERT-Base. This provides us a more comprehensive picture regarding the effect of encoder size and number of involved parameters on the bias mitigation methods.

All debiasing models are separately trained according to the adversarial bias mitigation and mutual information reduction methods. We particularly opt for a non-linear adversarial head with two fully connected layers and the tanh activation. Furthermore, to improve the capacity of adversarial learning, we initialize 5 instances of  $h_{\rho_i}$  and calculate the average of the loss for the backward pass. In MI reduction, in the case of protected attributes with more than two classes, we turn the multi-classes setting to multiple one-versus-rest splits. For the models with *diff* subnetworks, we conduct preliminary experiments to find proper thresholds for magnitude pruning that improves

sparsity as much as possible without sacrificing performance. For BERT-Base and BERT-Mini, we set the minimum sparsity threshold to 99% and 95% (maximum size of 1% and 5%), respectively. We note that these ratios are the lower bounds, and the  $L_0$  regularization may by itself reach a higher sparsity. The complete details of our hyperparameters setting and training procedure are explained in Appendix A. Our code and trained resources are available in <https://github.com/CPJKU/ModularizedDebiasing>.

**Evaluation Metrics** We evaluate the performance of the classifiers on the core task using the accuracy metric. We evaluate bias mitigation based on the concept of *fairness through blindness*, namely by examining whether models are agnostic about the protected attributes. Concretely following previous works (Elazar and Goldberg, 2018; Barrett et al., 2019), we report the leakage of a protected attribute in terms of the performance of a strong probe (or attacker) network. To train the probe, we freeze the model’s parameters, and train a new classification head (two-layer feed-forward layer with a tanh activation) to predict the protected attribute from the  $z$  encoding vector. For each evaluation, we train an ensemble of 5 probes for 40 epochs with early stopping if validation loss does not increase over 5 epochs, and report the results of the best performing probe. We report the performance of the probe in terms of balanced/macro accuracy (average of per-class accuracy scores). Balanced accuracy has the benefit of better reflecting the performance of the methods when considering minority groups, particularly given the unbalanced

Model	Adversarial Bias Mitigation			MI Reduction		
	Task $\uparrow$	Probe $_G\downarrow$	Probe $_A\downarrow$	Task $\uparrow$	Probe $_G\downarrow$	Probe $_A\downarrow$
FINETUNE	93.5 <sub>0.2</sub>	70.4 <sub>2.2</sub>	53.6 <sub>3.2</sub>	93.5 <sub>0.2</sub>	70.4 <sub>2.2</sub>	53.6 <sub>3.2</sub>
ADAPTER	87.3 <sub>0.1</sub>	67.8 <sub>0.6</sub>	37.1 <sub>1.7</sub>	87.3 <sub>0.1</sub>	67.8 <sub>0.6</sub>	37.1 <sub>1.7</sub>
FINETUNE-DEBIAS $_G$	93.7 <sub>0.0</sub>	52.3 <sub>0.5</sub>	34.8 <sub>1.1</sub>	93.7 <sub>0.1</sub>	62.2 <sub>0.1</sub>	41.9 <sub>1.4</sub>
FINETUNE-DEBIAS $_A$	92.9 <sub>0.2</sub>	52.3 <sub>0.3</sub>	31.4 <sub>2.6</sub>	93.7 <sub>0.1</sub>	61.5 <sub>0.7</sub>	41.2 <sub>2.1</sub>
FINETUNE-DEBIAS $_{G\&A}$	93.0 <sub>0.2</sub>	51.7 <sub>1.4</sub>	33.9 <sub>1.8</sub>	93.5 <sub>0.1</sub>	61.4 <sub>0.7</sub>	40.8 <sub>1.0</sub>
ADAPTER-DEBIAS $_G$	86.9 <sub>0.1</sub>	54.9 <sub>0.4</sub>	29.2 <sub>1.7</sub>	87.0 <sub>0.1</sub>	67.4 <sub>0.3</sub>	37.0 <sub>2.2</sub>
ADAPTER-DEBIAS $_A$	86.1 <sub>0.2</sub>	58.5 <sub>0.2</sub>	25.6 <sub>2.0</sub> <sup>♣</sup>	87.1 <sub>0.3</sub>	67.1 <sub>0.5</sub>	36.6 <sub>1.2</sub>
ADAPTER-DEBIAS $_{G\&A}$	92.2 <sub>0.2</sub>	51.6 <sub>2.4</sub> <sup>♣</sup>	27.5 <sub>1.9</sub>	92.8 <sub>0.1</sub>	64.2 <sub>0.4</sub>	34.9 <sub>0.2</sub>
DIFFPRUN	93.0 <sub>0.1</sub>	76.1 <sub>0.4</sub>	62.4 <sub>0.5</sub>	93.0 <sub>0.1</sub>	76.1 <sub>0.4</sub>	62.4 <sub>0.5</sub>
DIFFPRUN-DEBIAS $_G$	93.4 <sub>0.1</sub>	56.3 <sub>0.8</sub>	48.5 <sub>1.2</sub>	93.3 <sub>0.0</sub>	73.0 <sub>0.7</sub>	58.2 <sub>0.5</sub>
DIFFPRUN-DEBIAS $_A$	92.9 <sub>0.1</sub>	64.5 <sub>0.8</sub>	34.1 <sub>1.1</sub>	93.4 <sub>0.0</sub>	73.1 <sub>0.8</sub>	57.1 <sub>0.6</sub>
DIFFPRUN-DEBIAS $_{G\&A}$	92.9 <sub>0.2</sub>	54.7 <sub>2.6</sub>	29.7 <sub>2.0</sub>	93.6 <sub>0.2</sub>	73.3 <sub>0.3</sub>	57.2 <sub>0.9</sub>
MODDIFFY-POST $_G$	93.7 <sub>0.1</sub>	57.0 <sub>1.0</sub>	45.4 <sub>1.3</sub>	93.6 <sub>0.1</sub>	66.2 <sub>0.2</sub>	46.2 <sub>1.4</sub>
MODDIFFY-POST $_A$	93.7 <sub>0.1</sub>	63.7 <sub>0.9</sub>	31.7 <sub>2.8</sub>	93.6 <sub>0.0</sub>	66.6 <sub>0.2</sub>	46.5 <sub>2.8</sub>
MODDIFFY-POST $_G + ditto_A$	92.3 <sub>0.6</sub>	57.7 <sub>1.2</sub>	32.0 <sub>2.8</sub>	93.4 <sub>0.1</sub>	66.4 <sub>0.8</sub>	46.6 <sub>1.8</sub>
MODDIFFY-POST $_{G\&A}$	93.6 <sub>0.1</sub>	52.6 <sub>2.2</sub>	30.9 <sub>0.5</sub>	93.6 <sub>0.1</sub>	66.9 <sub>0.4</sub>	47.4 <sub>0.1</sub>
MODDIFFY-PAR $_G$	93.5 <sub>0.2</sub>	53.0 <sub>1.6</sub>	32.2 <sub>1.1</sub>	93.6 <sub>0.0</sub>	59.3 <sub>0.4</sub>	35.7 <sub>0.4</sub>
MODDIFFY-PAR $_A$	93.5 <sub>0.2</sub>	53.8 <sub>1.9</sub>	30.1 <sub>1.0</sub>	93.7 <sub>0.2</sub>	55.0 <sub>0.2</sub>	29.9 <sub>0.9</sub>
MODDIFFY-PAR $_G + ditto_A$	93.5 <sub>0.2</sub>	52.8 <sub>1.7</sub>	30.2 <sub>0.8</sub>	93.6 <sub>0.2</sub>	56.1 <sub>1.1</sub>	30.1 <sub>1.0</sub>
MODDIFFY-PAR $_{G\&A}$	93.8 <sub>0.2</sub>	<b>52.3</b> <sub>1.4</sub>	<b>28.3</b> <sub>1.6</sub>	93.6 <sub>0.2</sub>	<b>52.7</b> <sub>1.3</sub> <sup>♣</sup>	<b>29.4</b> <sub>0.6</sub> <sup>♣</sup>

Table 2: Results of the PAN16 dataset on **BERT-Base**. The subscripts G and A refer to the protected attributes gender and age, respectively. The sign  $G\&A$  denotes that the bias mitigation loss is the sum of the debiasing loss terms of both gender and age. The MODDIFFY models in the form of “+ditto” refer to the case where first two debiasing subnetworks with the same core model are trained separately, and then they are added to the base model at inference time. The results with the best bias mitigation performance (lowest values) among the models that use *diff* subnetworks (lower part of the table) are shown in **bold**, and among all models with the ♣ symbol. Subscript values indicate standard deviation.

distributions over protected labels in the datasets. To account for possible variabilities, we repeat every experiment five times and report the mean and standard deviation.

## 5 Results and Analysis

**Single-attribute Evaluation** Table 1 reports the evaluation results of the BIOS and FCDL18 datasets on BERT-Base using adversarial bias removal and mutual information (MI) reduction. The results of the same experiments on BERT-Mini are shown in Table 6 in Appendix B.

Starting from task accuracy, the models using subnetworks (variations of DIFFPRUN and MODDIFFY shown at the lower part of the table) consistently perform the same as the fully finetuned models on both datasets and debiasing methods. We observe a slight decrease in performance for the adapter-based models on FCDL18.

Looking at the leakage probing performance of bias attributes, MODDIFFY models show better

performance (lower values) among the subnetwork-based models on both datasets, and overall on BIOS.<sup>2</sup> In particular, MODDIFFY models outperform the directly comparable baselines FINETUNE-DEBIAS and DIFFPRUN-DEBIAS on all configurations, indicating the benefits of learning separate debiasing modules on bias mitigation performance. This indeed comes with the core advantages of MODDIFFY models in proving modularized and on-demand bias mitigation.

The results also show that MODDIFFY-POST, while (as expected) slightly weaker than MODDIFFY-PAR, provides competitive bias mitigation performance (i.e., on par with DIFFPRUN-DEBIAS). Finally, comparing between debiasing optimizations, MI reduction shows consistently worse bias

<sup>2</sup>Particularly on FCDL18, we observe high variations, which requires us to more cautiously interpret results. We assume that this is due to the small size of this dataset especially in the learning regimes with many parameters on BERT-Base, as this effect is less pronounced on BERT-Mini (Table 6 in Appendix B).

mitigation performance in comparison with adversarial training, particularly on FCDDL18 where the protected attribute has more than two labels.

**Two-attribute Evaluation** The results on PAN16 using BERT-Base are reported in Table 2, and the same experiments on BERT-Mini in Table 7 in Appendix B. In this experiment, every debiasing model is trained based on either gender, age, or simultaneously on both gender and age, shown with the subscripts  $G$ ,  $A$ ,  $G\&A$ , respectively. An additional evaluation is indicated with  $\text{MODDIFFY-}^*_G + \text{ditto}_A$ , which refers to adding the two separately trained gender and age debiasing subnetworks to the core model at inference time. In fact, for the experiments of each MODDIFFY model indicated with “+ditto”,  $G$ , and  $A$ , we train only one model with gender and age subnetworks, and then add the respective subnetwork(s) to the core model.

Looking at task performance results, similar to the previous datasets the models perform on par with fully finetuning, except the adapter-based models which in this case significantly underperform in both optimization methods. Regarding bias mitigation performance, we observe similar patterns to the ones discussed on the other datasets: MODDIFFY models particularly MODDIFFY-PAR show the least attribute leakage among the subnetwork-based models consistently, and also over all models with MI reduction. This reaffirms the benefits of modularizing debiasing of the attributes separately from the task. The results also indicate the existence of a correlation between the gender and age attributes in this dataset, such that debiasing each attribute also results in a decrease in leakage of the other attribute. The overall best results are achieved on the models that simultaneously optimize on both attributes ( $G\&A$ ).

Finally, let us have a closer look at the results of applying the two independently trained subnetworks. On both training regimes, we observe on par debiasing performance between  $\text{MODDIFFY-}^*_G + \text{ditto}_A$  and the corresponding results with one subnetwork, namely the ones of gender and age debiasing in  $\text{MODDIFFY-}^*_G$  and  $\text{MODDIFFY-}^*_A$ , respectively. These results indicate the viability of our approach to effectively merge subnetworks at inference time.

**Subnetworks analysis** We further investigate the achieved sparsity rate of the subnetworks, keeping

	PAN16		BIOS	FCDDL18
	Gender	Age	Gender	Dialect
Overall	0.01%	1.00%	0.27%	0.18%
Layer 12	0.04%	7.00%	0.26%	0.92%
Layer 11	0.02%	3.59%	0.30%	0.52%
Layer 10	0.01%	1.76%	0.33%	0.30%
Layer 9	0.00%	0.49%	0.20%	0.18%
Layer 8	0.00%	0.38%	0.22%	0.11%
Layer 7	0.01%	0.21%	0.30%	0.14%
Layer 6	0.00%	0.15%	0.28%	0.13%
Layer 5	0.01%	0.03%	0.28%	0.07%
Layer 4	0.01%	0.25%	0.43%	0.08%
Layer 3	0.01%	0.25%	0.56%	0.08%
Layer 2	0.01%	0.38%	0.41%	0.07%
Layer 1	0.01%	0.21%	0.44%	0.06%
Embeddings	0.00%	0.04%	0.05%	0.03%

Table 3: The percentage of non-masked parameters in the debiasing subnetworks of MODDIFFY-PAR.

in mind that the maximum capacity of the debiasing subnetworks on BERT-Base is set to 1%. The size of a subnetwork indicates the amount of information or in fact modifications needed to be applied, in order to debias a protected attribute. Table 3 reports the percentage of the number of non-masked parameters in every layer, and also overall, in the subnetworks of the MODDIFFY-PAR models regarding the protected attributes. The results show interesting patterns in respect to various protected attributes: the gender attribute on both PAN16 and BIOS dataset require much smaller subnetworks, such that the subnetwork on PAN16 is only 0.01% of the size of the core model. The age attribute appears to be a more complex topic in the underlying PLM, as it fully uses the 1% maximum capacity. Looking across the layers, the results show that debiasing the gender attribute is mostly handled in the lower transformer layers (particularly on BIOS), while debiasing age and dialect attributes mostly happens at the higher layers.

In Appendix C, we further discuss this topic for all models on the level of individual weight matrices. Moreover, in Appendix D we investigate to what extent the subnetworks of a model across several runs affect on the same set of parameters.

## 6 Conclusion

We propose MODDIFFY, a novel bias mitigation approach which enables integration of an arbitrary subset of the debiasing modules at inference time. Our method encapsulates the functionality of bias mitigation in respect to a protected attribute into a separate magnitude-difference subnetwork, which can then be applied to the core model on-demand. Our experiments on three classification tasks show



that MODDIFFY improves bias mitigation achieved by separating the debiasing from the task network, and effectively mitigates the bias of two (and possibly more) attributes when their respective subnetworks are simultaneously utilized.

## 7 Limitations

An important limitation of our work concerns the definition of the protected attributes in the datasets used for evaluation. In particular, gender in BIOS and PAN16 is limited to the binary female/male, lacking an inclusive and nuanced definition of gender. Similarly in FDCL18, we consider only two dialects of *African American* and *White American*, while clearly this definition is limited and non-inclusive. Furthermore as in previous work (Sap et al., 2019; Ravfogel et al., 2020; Zhang et al., 2021), the labels of this protected attribute are assigned through a probabilistic model, and hence the dataset might not represent the nuances and traits of the real-world.

The second limitation regards reaching strong conclusions on the generalizability of the multi-attribute setting for MODDIFFY over any possible number of protected attributes or subset of them. Our multi-attribute experiments are conducted on one dataset with two attributes of gender and age, particularly due to the lack of available suitable datasets. Hence, Further studies (as well as more suitable datasets) are required for achieving a more comprehensive picture on the topic.

Finally, we should also highlight two general limitations, shared with the other related studies in the area of model bias mitigation. First, we should consider that the aim of representation disentanglement optimizations is to reduce the existing *correlations* in the model with the protected attributes based on the *observed data*. These data-oriented approaches might lack effective generalization, particularly when the model is evaluated in other domains or out-of-distribution data. Second, our bias mitigation evaluation is grounded in the notion of *fairness through blindness*, and the debiasing optimization methods are designed to support this form of fairness. The effects of our method on other possible definitions of fairness are therefore left for future work.

## 8 Acknowledgment

This work received financial support by the Austrian Science Fund (FWF): P33526 and DFH-23;

and by the State of Upper Austria and the Federal Ministry of Education, Science, and Research, through grants LIT-2020-9-SEE-113 and LIT-2021-YOU-215.

## References

- Maria Barrett, Yova Kementchedjheva, Yanai Elazar, Desmond Elliott, and Anders Søgaard. 2019. Adversarial removal of demographic attributes revisited. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6331–6336.
- Elad Ben-Zaken, Shauli Ravfogel, and Yoav Goldberg. 2021. Bitfit: Simple parameter-efficient fine-tuning for transformer-based masked language-models. *ArXiv*, abs/2106.10199.
- Asia J Biega, Krishna P Gummadi, and Gerhard Weikum. 2018. Equity of attention: Amortizing individual fairness in rankings. In *The 41st international ACM SIGIR Conference on Research & Development in Information Retrieval*, pages 405–414.
- Su Lin Blodgett, Solon Barocas, Hal Daumé III, and Hanna Wallach. 2020. Language (technology) is power: A critical survey of “bias” in nlp. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5454–5476.
- Su Lin Blodgett, Lisa Green, and Brendan O’Connor. 2016. [Demographic dialectal variation in social media: A case study of African-American English](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1119–1130, Austin, Texas. Association for Computational Linguistics.
- Tolga Bolukbasi, Kai-Wei Chang, James Y Zou, Venkatesh Saligrama, and Adam T Kalai. 2016. [Man is to computer programmer as woman is to homemaker? debiasing word embeddings](#). In *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.
- Tianlong Chen, Jonathan Frankle, Shiyu Chang, Sijia Liu, Yang Zhang, Zhangyang Wang, and Michael Carbin. 2020. The lottery ticket hypothesis for pre-trained BERT networks. In *Proceedings of NeurIPS*.
- Pengyu Cheng, Martin Renqiang Min, Dinghan Shen, Christopher Malon, Yizhe Zhang, Yitong Li, and Lawrence Carin. 2020. [Improving disentangled text representation learning with information-theoretic guidance](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7530–7541, Online. Association for Computational Linguistics.
- Pierre Colombo, Pablo Piantanida, and Chloé Clavel. 2021. A novel estimator of mutual information for learning to disentangle textual representations. 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 6539–6550.
- Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. 2019. [Bias in bios: A case study of semantic representation bias in a high-stakes setting](#). page 120–128, New York, NY, USA. Association for Computing Machinery.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [Bert: Pre-training of deep bidirectional transformers for language understanding](#). In *Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, volume 1, pages 4171–4186. Association for Computational Linguistics.
- Yanai Elazar and Yoav Goldberg. 2018. [Adversarial removal of demographic attributes from text data](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 11–21. Association for Computational Linguistics.
- Antigoni Maria Founta, Constantinos Djouvas, Despoina Chatzakou, Ilias Leontiadis, Jeremy Blackburn, Gianluca Stringhini, Athena Vakali, Michael Sirivianos, and Nicolas Kourtellis. 2018. Large scale crowdsourcing and characterization of twitter abusive behavior. In *Twelfth International AAAI Conference on Web and Social Media*.
- Jonathan Frankle and Michael Carbin. 2019. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *Proceedings of International Conference on Learning Representations, ICLR*.
- Christian Ganhör, David Penz, Navid Rekabsaz, Oleg Lesota, and Markus Schedl. 2022. [Unlearning protected user attributes in recommendations with adversarial training](#). In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22*, page 2142–2147, New York, NY, USA. Association for Computing Machinery.
- Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *International conference on machine learning*, pages 1180–1189. Proceedings of Machine Learning Research.
- Yaroslav Ganin, Evgeniya Ustinova, Hana Ajakan, Pascal Germain, Hugo Larochelle, François Laviolette, Mario Marchand, and Victor Lempitsky. 2016. Domain-adversarial training of neural networks. In *Journal of Machine Learning Research*, volume 17, pages 1–35.
- Arthur Gretton, Karsten M Borgwardt, Malte J Rasch, Bernhard Schölkopf, and Alexander Smola. 2012. A kernel two-sample test. *The Journal of Machine Learning Research*, 13(1):723–773.
- Demi Guo, Alexander Rush, and Yoon Kim. 2021. Parameter-efficient transfer learning with diff pruning. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics*, pages 4884–4896.
- Yue Guo, Yi Yang, and Ahmed Abbasi. 2022. Auto-debias: Debiasing masked language models with automated biased prompts. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1012–1023.
- Wenjuan Han, Bo Pang, and Ying Nian Wu. 2021a. Robust transfer learning with pretrained language models through adapters. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 854–861.
- Xudong Han, Timothy Baldwin, and Trevor Cohn. 2021b. [Diverse adversaries for mitigating bias in training](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2760–2765, Online. Association for Computational Linguistics.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for NLP](#). In *International Conference on Machine Learning*, volume 97, pages 2790–2799. Proceedings of Machine Learning Research.
- Edward Hu, Yelong Shen, Phil Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).
- Sebastian Jaszczur, Aakanksha Chowdhery, Afroz Mohiuddin, Łukasz Kaiser, Wojciech Gajewski, Henryk Michalewski, and Jonni Kanerva. 2021. Sparse is enough in scaling transformers. *Advances in Neural Information Processing Systems*, 34.
- Masahiro Kaneko and Danushka Bollegala. 2021. Debiasing pre-trained contextualised embeddings. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1256–1266.
- Klara Krieg, Emilia Parada-Cabaleiro, Gertraud Medicus, Oleg Lesota, Markus Schedl, and Navid Rekabsaz. 2023. Grep-biasir: A dataset for investigating gender representation-bias in information retrieval results. In *Proceeding of the ACM SIGIR Conference On Human Information Interaction And Retrieval (CHIIR)*.
- François Lagunas, Ella Charlaix, Victor Sanh, and Alexander M Rush. 2021. Block pruning for faster transformers. In *Proceedings of the 2021 Conference*

- on *Empirical Methods in Natural Language Processing*, pages 10619–10629.
- Anne Lauscher, Tobias Lueken, and Goran Glavaš. 2021. [Sustainable modular debiasing of language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 4782–4797, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Christos Louizos, Max Welling, and Diederik P. Kingma. 2018. Learning sparse neural networks through l0 regularization. In *International Conference on Learning Representations*.
- David Madras, Elliot Creager, Toniann Pitassi, and Richard Zemel. 2018. Learning adversarially fair and transferable representations. In *Proceedings of the International Conference on Machine Learning*, pages 3384–3393. PMLR.
- Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. 2021. [Compacter: Efficient low-rank hypercomplex adapter layers](#). In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 1022–1035.
- Johannes Mario Meissner, Saku Sugawara, and Akiko Aizawa. 2022. Debiasing masks: A new framework for shortcut mitigation in NLU. In *Proceeding of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Francisco Manuel Rangel Pardo, Paolo Rosso, Ben Verhoeven, Walter Daelemans, Martin Potthast, and Benno Stein. 2016. Overview of the 4th author profiling task at pan 2016: Cross-genre evaluations. In *CLEF*.
- Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. 2021. Adapterfusion: Non-destructive task composition for transfer learning. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 487–503.
- Francisco Rangel, Paolo Rosso, Ben Verhoeven, Walter Daelemans, Martin Potthast, and Benno Stein. 2016. [Pan16 author profiling](#). Zenodo.
- Shauli Ravfogel, Yanai Elazar, Hila Gonen, Michael Twiton, and Yoav Goldberg. 2020. Null it out: Guarding protected attributes by iterative nullspace projection. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7237–7256.
- Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. 2017. Learning multiple visual domains with residual adapters. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30.
- Navid Rekabsaz, Simone Kopeinik, and Markus Schedl. 2021. Societal biases in retrieved contents: Measurement framework and adversarial mitigation of BERT rankers. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 306–316.
- Navid Rekabsaz and Markus Schedl. 2020. Do neural ranking models intensify gender bias? In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2065–2068.
- Andreas Rücklé, Gregor Geigle, Max Glockner, Tilman Beck, Jonas Pfeiffer, Nils Reimers, and Iryna Gurevych. 2021. AdapterDrop: On the efficiency of adapters in transformers. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7930–7946.
- Maarten Sap, Dallas Card, Saadia Gabriel, Yejin Choi, and Noah A. Smith. 2019. [The risk of racial bias in hate speech detection](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1668–1678, Florence, Italy. Association for Computational Linguistics.
- Timo Schick, Sahana Udupa, and Hinrich Schütze. 2021. Self-diagnosis and self-debiasing: A proposal for reducing corpus-based bias in NLP. *Transactions of the Association for Computational Linguistics*, 9:1408–1424.
- Emily Sheng, Kai-Wei Chang, Prem Natarajan, and Nanyun Peng. 2019. The woman worked as a babysitter: On biases in language generation. 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 3398–3403.
- Gabriel Stanovsky, Noah A Smith, and Luke Zettlemoyer. 2019. Evaluating gender bias in machine translation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1679–1684.
- Asa Cooper Stickland and Iain Murray. 2019. [BERT and PALs: Projected attention layers for efficient adaptation in multi-task learning](#). In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5986–5995. PMLR.
- Yi-Lin Sung, Varun Nair, and Colin Raffel. 2021. Training neural networks with fixed sparse masks. *ArXiv*, abs/2111.09839.
- Iulia Turc, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Well-read students learn better: On the importance of pre-training compact models. *arXiv preprint arXiv:1908.08962v2*.
- Eric Tzeng, Judy Hoffman, Ning Zhang, Kate Saenko, and Trevor Darrell. 2014. Deep domain confusion: Maximizing for domain invariance. *arXiv preprint arXiv:1412.3474*.

- Liwen Wang, Yuanmeng Yan, Keqing He, Yanan Wu, and Weiran Xu. 2021. Dynamically disentangling social bias from task-oriented representations with adversarial attack. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3740–3750.
- Mitchell Wortsman, Vivek Ramanujan, Rosanne Liu, Aniruddha Kembhavi, Mohammad Rastegari, Jason Yosinski, and Ali Farhadi. 2020. Supermasks in superposition. *Advances in Neural Information Processing Systems*, 33.
- Mengzhou Xia, Anjalie Field, and Yulia Tsvetkov. 2020. [Demoting racial bias in hate speech detection](#). In *Proceedings of the Eighth International Workshop on Natural Language Processing for Social Media*, pages 7–14, Online. Association for Computational Linguistics.
- Qizhe Xie, Zihang Dai, Yulun Du, Eduard Hovy, and Graham Neubig. 2017. Controllable invariance through adversarial feature learning. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 585–596.
- Runxin Xu, Fuli Luo, Zhiyuan Zhang, Chuanqi Tan, Baobao Chang, Songfang Huang, and Fei Huang. 2021. [Raise a child in large language model: Towards effective and generalizable fine-tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9514–9528, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- George Zerveas, Navid Rekabsaz, Daniel Cohen, and Carsten Eickhoff. 2022. [Mitigating bias in search results through contextual document reranking and neutrality regularization](#). In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '22*, page 2532–2538, New York, NY, USA. Association for Computing Machinery.
- Xiongyi Zhang, Jan-Willem van de Meent, and Byron Wallace. 2021. [Disentangling representations of text by masking transformers](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 778–791, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Jieyu Zhao, Tianlu Wang, Mark Yatskar, Ryan Cotterell, Vicente Ordonez, and Kai-Wei Chang. 2019. Gender bias in contextualized word embeddings. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 629–634.
- Mengjie Zhao, Tao Lin, Fei Mi, Martin Jaggi, and Hinrich Schütze. 2020. [Masking as an efficient alternative to finetuning for pretrained language models](#). In *Empirical Methods in Natural Language Processing*, pages 2226–2241. Association for Computational Linguistics.
- Hattie Zhou, Janice Lan, Rosanne Liu, and Jason Yosinski. 2019. Deconstructing lottery tickets: Zeros, signs, and the supermask. *Advances in Neural Information Processing Systems*, 32:3597–3607.



## A Experiment Settings – Additional Details

In FDCL18 dataset, we use the TwitterAAE model (Blodgett et al., 2016) to assign racial dialect classes. The TwitterAAE model predicts four racial classes, *African American*, *White American*, *Hispanic*, and *Others*. We labeled a tweet as *African American* or *White American* if the prediction score was greater than 0.5. For PAN16 dataset, following (Sap et al., 2019) we balanced the task labels and sampled 200K data. The age groups of this dataset are 18-24, 25-34, 35-49, 50-64, and 65+.

We randomly split the dataset into train, validation, and test set with the proportions 63:12:15 for BIOS, 63:12:15 for FDCL18, and 80:5:15 on PAN16. We use the validation set for hyperparameter tuning, and the best result on the validation set is evaluated on test set for the final results. The validation and test sets in all datasets follow the same distribution as the whole dataset. To address the unbalancedness of the dataset and the potential problems in adversarial learning, we apply up-sampling only on the *training sets* of BIOS and FDCL18 datasets, to balance the protected attribute labels within each task label. For instance, genders are balanced in the dentist class by repeating the data items of the minority subgroup.

Adversarial heads consist of five classifiers with different initialization. The loss for the five classifiers is averaged and accuracy is measured via majority vote. All baseline models are trained for 20 epochs. All DIFFPRUN and MODDIFFY model variants are trained for 30 epochs as they need to account for the two phases of *diff* pruning, where a model requires more training to recover its performance after the magnitude pruning step. We fix the learning rate of BERT weights to  $2e-5$  and the learning rate for the classifier heads to  $1e-4$ . We set the batch size to 64 for all experiments. We keep other *diff*-specific hyperparameters the same as suggested by Guo et al. (2021). Adapter baselines follow Pfeiffer et al. (2021) with reduction factor of two. Rest of the hyperparameters are same as the subnetwork-based models. Table 4 reports the hyperparameters of our experiments.

## B Additional Results

The results of all models using BERT-Mini are shown in Table 6 for BIOS and FCDL18 datasets, and in Table 7 for PAN16. Table 8 reports the evaluation results of the INLP method.

training	
batch_size	64
structured_diff_pruning	True
alpha_init	5
concrete_samples	1
concrete_lower	-1.5
concrete_upper	1.5
num_epochs	20
num_epochs_finetune	15
num_epochs_fixmask	15
learning_rate	2e-05
learning_rate_task_head	0.0001
learning_rate_adv_head	0.0001
learning_rate_alpha	0.1
task_dropout	0.3
task_n_hidden	0
adv_dropout	0.3
adv_n_hidden	1
adv_count	5
adv_lambda	1.0
sparsity_pen	1.25e-07
max_grad_norm	1.0
adv attack	
batch_size	64
num_epochs	40
learning_rate	0.0001
adv_n_hidden	1
adv_count	5
adv_dropout	0.3

Table 4: Hyperparameters used for training

Parameter Name	Size
word_embeddings	117,204,480
intermediate.dense	11,811,840
output.dense	11,800,320
attention.self.query	2,952,960
attention.self.key	2,952,960
attention.self.value	2,952,960
attention.output.dense	2,952,960
position_embeddings	1,966,080
output.adapter	590,976
others	7680

Table 5: BERT-Base number of parameters

## C Sparsity rates of subnetworks

We visualize the percentage of non-masked parameters of the subnetworks in BERT-Base for each parameter matrix in Figures 2, 3 and 4 for MODDIFFY-PAR, MODDIFFY-POST, and DIFFPRUN-DEBIAS, respectively. In addition to the discussion in Section 5, we observe in these detailed figures that the LayerNorm module of the last Transformer block generally has a high density. We assume

Model	BIOS				FCDL18			
	Adversarial		MI Reduction		Adversarial		MI Reduction	
	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$
FINETUNE	82.9 <sub>0.1</sub>	65.5 <sub>0.4</sub>	82.9 <sub>0.1</sub>	65.5 <sub>0.4</sub>	82.1 <sub>0.2</sub>	90.3 <sub>0.9</sub>	82.1 <sub>0.2</sub>	90.3 <sub>0.9</sub>
ADAPTER	81.6 <sub>0.2</sub>	65.7 <sub>0.2</sub>	81.6 <sub>0.2</sub>	65.7 <sub>0.2</sub>	81.9 <sub>0.0</sub>	79.2 <sub>0.5</sub>	81.9 <sub>0.0</sub>	79.2 <sub>0.5</sub>
FINETUNE-DEBIAS	81.6 <sub>0.2</sub>	56.4 <sub>1.7</sub>	81.6 <sub>0.1</sub>	59.9 <sub>0.9</sub>	80.0 <sub>0.5</sub>	73.7 <sub>2.7</sub>	79.7 <sub>0.5</sub>	87.3 <sub>2.6</sub>
ADAPTER-DEBIAS	81.5 <sub>0.1</sub>	63.4 <sub>0.1</sub>	81.7 <sub>0.1</sub>	65.4 <sub>0.1</sub>	81.1 <sub>0.2</sub>	64.6 <sub>1.3</sub> <sup>♣</sup>	81.8 <sub>0.1</sub>	78.5 <sub>0.5</sub> <sup>♣</sup>
DIFFPRUN	83.5 <sub>0.1</sub>	65.8 <sub>0.3</sub>	83.5 <sub>0.1</sub>	65.8 <sub>0.3</sub>	82.8 <sub>0.1</sub>	92.6 <sub>0.9</sub>	82.8 <sub>0.1</sub>	92.6 <sub>0.9</sub>
DIFFPRUN-DEBIAS	83.3 <sub>0.1</sub>	59.1 <sub>0.7</sub>	82.3 <sub>0.1</sub>	65.5 <sub>0.9</sub>	82.3 <sub>0.4</sub>	65.8 <sub>2.8</sub>	82.5 <sub>0.3</sub>	91.9 <sub>0.9</sub>
MODDIFFY-POST	83.1 <sub>0.0</sub>	57.3 <sub>0.7</sub>	83.1 <sub>0.0</sub>	63.4 <sub>1.4</sub>	81.6 <sub>0.4</sub>	69.7 <sub>1.9</sub> <sup>♣</sup>	82.3 <sub>0.1</sub>	89.3 <sub>0.4</sub>
MODDIFFY-PAR	81.5 <sub>0.2</sub>	<b>55.7</b> <sub>0.6</sub> <sup>♣</sup>	81.4 <sub>0.2</sub>	<b>58.8</b> <sub>0.8</sub> <sup>♣</sup>	80.2 <sub>0.3</sub>	73.8 <sub>6.7</sub>	79.9 <sub>80.6</sub>	<b>85.4</b> <sub>1.6</sub>

Table 6: Results of the BIOS and FCDL18 datasets on **BERT-Mini** with adversarial bias removal and mutual information (MI) reduction. Task performance is measured with accuracy, and bias mitigation with balanced accuracy of the probes. The protected attribute is gender for BIOS, and race-dialect for FCDL18. The results with the best bias mitigation performance (lowest values) among the models that use *diff* subnetworks (lower part of the table) are shown in **bold**, and among all models with the ♣ symbol.

Model	Adversarial Bias Mitigation			MI Reduction		
	Task $\uparrow$	Probe <sub>G</sub> $\downarrow$	Probe <sub>A</sub> $\downarrow$	Task $\uparrow$	Probe <sub>G</sub> $\downarrow$	Probe <sub>A</sub> $\downarrow$
FINETUNE	91.5 <sub>0.2</sub>	64.8 <sub>0.6</sub>	46.8 <sub>0.3</sub>	91.5 <sub>0.2</sub>	64.8 <sub>0.3</sub>	88.4 <sub>1.1</sub>
ADAPTER	78.4 <sub>0.2</sub>	65.8 <sub>0.2</sub>	35.3 <sub>0.8</sub>	78.4 <sub>0.2</sub>	65.8 <sub>0.2</sub>	35.3 <sub>0.8</sub>
FINETUNE-DEBIAS <sub>G</sub>	91.7 <sub>0.2</sub>	54.6 <sub>0.6</sub>	42.2 <sub>0.6</sub>	91.6 <sub>0.2</sub>	61.3 <sub>0.6</sub>	42.4 <sub>0.5</sub>
FINETUNE-DEBIAS <sub>A</sub>	91.1 <sub>0.2</sub>	61.9 <sub>0.6</sub>	39.1 <sub>0.8</sub>	91.4 <sub>0.6</sub>	61.9 <sub>0.1</sub>	43.2 <sub>0.0</sub>
FINETUNE-DEBIAS <sub>G&amp;A</sub>	91.2 <sub>0.1</sub>	57.0 <sub>1.0</sub>	38.5 <sub>0.6</sub>	91.9 <sub>0.0</sub>	62.5 <sub>0.1</sub>	42.1 <sub>0.9</sub>
ADAPTER-DEBIAS <sub>G</sub>	78.1 <sub>0.1</sub>	59.6 <sub>0.4</sub>	32.0 <sub>1.3</sub>	78.5 <sub>0.2</sub>	65.9 <sub>0.2</sub>	34.8 <sub>0.0</sub>
ADAPTER-DEBIAS <sub>A</sub>	77.3 <sub>0.1</sub>	60.5 <sub>0.9</sub>	27.3 <sub>0.9</sub>	78.3 <sub>0.1</sub>	65.5 <sub>0.2</sub>	34.1 <sub>0.3</sub> <sup>♣</sup>
ADAPTER-DEBIAS <sub>G&amp;A</sub>	80.9 <sub>0.6</sub>	55.6 <sub>0.4</sub>	25.5 <sub>1.2</sub> <sup>♣</sup>	82.0 <sub>0.1</sub>	64.2 <sub>0.4</sub>	34.9 <sub>0.2</sub>
DIFFPRUN	90.0 <sub>0.1</sub>	67.2 <sub>0.3</sub>	49.4 <sub>0.8</sub>	90.0 <sub>0.1</sub>	67.2 <sub>0.3</sub>	49.4 <sub>0.8</sub>
DIFFPRUN-DEBIAS <sub>G</sub>	90.1 <sub>0.1</sub>	<b>54.1</b> <sub>1.1</sub> <sup>♣</sup>	44.1 <sub>0.7</sub>	78.5 <sub>0.2</sub>	65.9 <sub>0.2</sub>	34.8 <sub>0.0</sub>
DIFFPRUN-DEBIAS <sub>A</sub>	89.2 <sub>0.3</sub>	64.8 <sub>0.6</sub>	39.4 <sub>1.9</sub>	78.3 <sub>0.1</sub>	65.5 <sub>0.2</sub>	<b>34.1</b> <sub>0.3</sub> <sup>♣</sup>
DIFFPRUN-DEBIAS <sub>G&amp;A</sub>	89.1 <sub>0.1</sub>	57.9 <sub>1.7</sub>	35.8 <sub>2.2</sub>	86.3 <sub>0.0</sub>	68.8 <sub>0.1</sub>	52.3 <sub>0.9</sub>
MODDIFFY-POST <sub>G</sub>	91.7 <sub>0.2</sub>	55.7 <sub>1.2</sub>	42.7 <sub>0.6</sub>	91.5 <sub>0.0</sub>	62.6 <sub>0.0</sub>	43.7 <sub>1.3</sub>
MODDIFFY-POST <sub>A</sub>	91.4 <sub>0.2</sub>	62.3 <sub>0.5</sub>	34.8 <sub>0.8</sub>	91.5 <sub>0.0</sub>	62.8 <sub>0.3</sub>	44.0 <sub>0.0</sub>
MODDIFFY-POST <sub>G</sub> + <i>ditto</i> <sub>A</sub>	91.0 <sub>0.4</sub>	59.3 <sub>0.6</sub>	37.3 <sub>1.2</sub>	!!! 91.5 <sub>0.5</sub>	62.8 <sub>0.4</sub>	43.9 <sub>0.8</sub>
MODDIFFY-POST <sub>G&amp;A</sub>	91.5 <sub>0.2</sub>	56.7 <sub>1.3</sub>	35.1 <sub>2.1</sub>	91.5 <sub>0.0</sub>	63.0 <sub>0.3</sub>	43.9 <sub>0.4</sub>
MODDIFFY-PAR <sub>G</sub>	91.6 <sub>0.2</sub>	55.9 <sub>0.8</sub>	41.7 <sub>0.7</sub>	91.6 <sub>0.1</sub>	60.9 <sub>0.4</sub>	40.3 <sub>0.8</sub>
MODDIFFY-PAR <sub>A</sub>	91.3 <sub>0.2</sub>	61.3 <sub>0.5</sub>	37.6 <sub>1.2</sub>	91.5 <sub>0.2</sub>	60.8 <sub>0.6</sub>	41.6 <sub>0.7</sub>
MODDIFFY-PAR <sub>G</sub> + <i>ditto</i> <sub>A</sub>	91.4 <sub>0.4</sub>	60.7 <sub>1.0</sub>	39.6 <sub>1.6</sub>	91.7 <sub>0.1</sub>	60.2 <sub>0.5</sub> <sup>♣</sup>	39.7 <sub>0.1</sub>
MODDIFFY-PAR <sub>G&amp;A</sub>	91.3 <sub>0.3</sub>	55.5 <sub>1.1</sub>	<b>33.7</b> <sub>1.6</sub>	91.5 <sub>0.3</sub>	<b>60.7</b> <sub>1.3</sub>	41.3 <sub>1.1</sub>

Table 7: Results of the PAN16 dataset on **BERT-Mini**. The subscripts G and A refer to the protected attributes gender and age, respectively. The sign G&A denotes that the bias mitigation loss of the model consists of the debiasing loss terms of both gender and age. The MODDIFFY models in the form of “+*ditto*” refer to the case, where first two debiasing subnetworks with the same core model are trained separately, and then they are added to the base model at inference time. The results with the best bias mitigation performance (lowest values) among the models that use *diff* subnetworks (lower part of the table) are shown in **bold**, and among all models with the ♣ symbol.

Model	BIOS (gender)		FDCL18 (race)	
	Task $\uparrow$	Probe $\downarrow$	Task $\uparrow$	Probe $\downarrow$
BERT-Mini	71.1 <sub>0.2</sub>	59.8 <sub>0.9</sub>	73.6 <sub>1.6</sub>	57.9 <sub>0.4</sub>
BERT-Base	66.2 <sub>0.4</sub>	50.5 <sub>0.3</sub>	71.6 <sub>2.2</sub>	50.2 <sub>0.1</sub>

(a) BIOS and FDCL18

Model		Task $\uparrow$	Probe $_G$ $\downarrow$	Probe $_A$ $\downarrow$
BERT-Mini	$G$	69.3 <sub>0.1</sub>	60.1 <sub>0.3</sub>	29.2 <sub>0.8</sub>
	$A$	66.7 <sub>1.8</sub>	60.6 <sub>0.1</sub>	25.6 <sub>0.2</sub>
BERT-Base	$G$	69.4 <sub>0.1</sub>	54.4 <sub>0.1</sub>	25.5 <sub>0.2</sub>
	$A$	49.8 <sub>0.8</sub>	54.6 <sub>0.6</sub>	27.5 <sub>0.2</sub>

(b) PAN16

Table 8: Evaluation results of INLP.

that this is due to the additive nature of these *diff*-based methods, as changing the weight magnitude through adding a subnetwork requires rescaling the final output.

## D Consistency in finding subnetworks

Figures 5, 6 and 7 show the percentage of common non-masked parameters in subnetworks on a particular weight matrix/layer across 5 runs for MODDIFFY-PAR, MODDIFFY-POST, and DIFF-PRUN-DEBIAS, respectively. We report the percentage of overlap between the subnetworks of two, three, four, and five runs, separated with the “/”. The results show that the equivalent subnetworks across various runs (with different initialization seeds) seem to be largely separated. This results are consistent with the observations on the lottery ticket hypothesis on large neural networks (Chen et al., 2020).

	0.00%	0.00%	0.02%	0.02%	0.04%	0.01%	0.01%	0.09%	0.01%
Layer 11	0.00%	0.00%	0.08%	0.10%	0.04%	0.03%	0.07%	0.14%	0.04%
Layer 10	0.00%	0.00%	0.03%	0.04%	0.13%	0.02%	0.03%	0.29%	0.02%
Layer 9	0.00%	0.00%	0.01%	0.01%	0.08%	0.01%	0.01%	0.17%	0.01%
Layer 8	0.00%	0.00%	0.01%	0.00%	0.03%	0.01%	0.00%	0.23%	0.00%
Layer 7	0.00%	0.00%	0.01%	0.00%	0.00%	0.01%	0.00%	0.00%	0.00%
Layer 6	0.00%	0.00%	0.01%	0.00%	0.00%	0.01%	0.00%	0.10%	0.01%
Layer 5	0.00%	0.00%	0.01%	0.01%	0.00%	0.01%	0.00%	0.05%	0.00%
Layer 4	0.00%	0.00%	0.00%	0.01%	0.00%	0.01%	0.00%	0.36%	0.01%
Layer 3	0.00%	0.00%	0.01%	0.01%	0.08%	0.01%	0.00%	0.00%	0.01%
Layer 2	0.00%	0.00%	0.01%	0.05%	0.00%	0.01%	0.00%	0.09%	0.01%
Layer 1	0.00%	0.00%	0.02%	0.02%	0.13%	0.01%	0.00%	0.00%	0.01%
Layer 0	0.00%	0.00%	0.02%	0.03%	0.00%	0.00%	0.00%	0.35%	0.01%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.00%	0.06%	1.28%	0.12%	0.00%				0.00%
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(a) PAN16 - Gender

	0.72%	0.69%	1.14%	1.44%	0.63%	1.55%	1.31%	3.27%	1.00%
Layer 11	2.30%	2.95%	7.26%	8.87%	4.60%	8.73%	8.65%	66.98%	7.00%
Layer 10	1.73%	1.65%	3.44%	4.20%	1.46%	4.09%	3.94%	0.48%	3.59%
Layer 9	1.04%	0.79%	1.35%	1.67%	0.68%	2.18%	1.90%	0.53%	1.76%
Layer 8	0.55%	0.26%	0.22%	0.63%	0.18%	0.90%	0.17%	0.18%	0.49%
Layer 7	0.79%	0.40%	0.54%	0.49%	0.07%	0.48%	0.11%	0.10%	0.38%
Layer 6	0.49%	0.39%	0.27%	0.27%	0.10%	0.18%	0.10%	0.03%	0.21%
Layer 5	0.70%	0.17%	0.00%	0.19%	0.10%	0.15%	0.03%	0.12%	0.15%
Layer 4	0.01%	0.13%	0.13%	0.14%	0.07%	0.00%	0.00%	0.20%	0.03%
Layer 3	0.09%	0.73%	0.00%	0.35%	0.10%	0.35%	0.12%	0.18%	0.25%
Layer 2	0.30%	0.00%	0.28%	0.15%	0.05%	0.56%	0.00%	0.10%	0.25%
Layer 1	0.38%	0.38%	0.00%	0.26%	0.04%	0.87%	0.00%	0.00%	0.38%
Layer 0	0.22%	0.48%	0.14%	0.01%	0.09%	0.16%	0.27%	0.07%	0.21%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.04%	0.00%	0.00%	0.04%	0.00%				0.04%
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(b) PAN16 - Age

	0.11%	0.11%	0.68%	1.06%	1.73%	0.35%	0.34%	1.96%	0.27%
Layer 11	0.02%	0.03%	0.84%	0.95%	1.81%	0.17%	0.22%	3.70%	0.26%
Layer 10	0.03%	0.06%	1.13%	1.12%	1.90%	0.18%	0.14%	3.24%	0.30%
Layer 9	0.05%	0.09%	1.30%	1.31%	0.83%	0.19%	0.11%	1.81%	0.33%
Layer 8	0.12%	0.16%	0.19%	0.91%	1.50%	0.20%	0.07%	0.16%	0.20%
Layer 7	0.11%	0.17%	0.00%	0.91%	1.22%	0.27%	0.09%	2.04%	0.22%
Layer 6	0.06%	0.15%	0.28%	0.85%	0.64%	0.38%	0.17%	2.23%	0.30%
Layer 5	0.13%	0.09%	0.01%	0.66%	1.32%	0.49%	0.14%	1.82%	0.28%
Layer 4	0.13%	0.10%	0.00%	0.59%	2.62%	0.48%	0.15%	2.23%	0.28%
Layer 3	0.21%	0.16%	1.07%	1.14%	1.51%	0.48%	0.16%	4.09%	0.43%
Layer 2	0.16%	0.17%	1.38%	1.74%	2.21%	0.60%	0.20%	2.55%	0.56%
Layer 1	0.18%	0.11%	0.57%	1.11%	2.29%	0.47%	0.26%	1.56%	0.41%
Layer 0	0.11%	0.05%	1.43%	1.38%	2.93%	0.31%	0.27%	0.78%	0.44%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.04%	0.60%	0.00%	2.41%	0.05%				0.05%
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(c) BIOS - Gender

	0.03%	0.03%	0.42%	0.71%	2.46%	0.20%	0.29%	2.74%	0.18%
Layer 11	0.07%	0.05%	1.62%	2.64%	9.36%	0.83%	1.07%	6.65%	0.92%
Layer 10	0.04%	0.02%	1.14%	1.03%	3.70%	0.48%	0.51%	8.92%	0.52%
Layer 9	0.05%	0.05%	0.70%	0.99%	4.36%	0.21%	0.23%	4.04%	0.30%
Layer 8	0.07%	0.11%	0.26%	0.84%	1.26%	0.14%	0.06%	1.28%	0.18%
Layer 7	0.01%	0.01%	0.15%	0.30%	2.51%	0.14%	0.05%	2.16%	0.11%
Layer 6	0.05%	0.03%	0.14%	0.67%	2.33%	0.14%	0.06%	3.39%	0.14%
Layer 5	0.04%	0.03%	0.24%	0.35%	1.73%	0.18%	0.05%	2.46%	0.13%
Layer 4	0.02%	0.02%	0.17%	0.29%	0.90%	0.04%	0.04%	1.43%	0.07%
Layer 3	0.02%	0.01%	0.23%	0.48%	0.55%	0.04%	0.02%	1.15%	0.08%
Layer 2	0.01%	0.01%	0.21%	0.29%	1.73%	0.09%	0.03%	0.78%	0.08%
Layer 1	0.01%	0.01%	0.11%	0.25%	1.12%	0.07%	0.03%	1.45%	0.07%
Layer 0	0.00%	0.01%	0.12%	0.34%	0.00%	0.05%	0.02%	2.43%	0.06%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.01%	0.94%	0.00%	0.64%	0.03%				0.03%
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(d) FCDL18 - Dialect

Figure 2: Sparsity rate of the subnetworks in MODDIFFY-PAR on BERT-Base. Each value shows the sparsity rate on the specific block/matrix.



	0.21%	0.23%	1.19%	1.44%	2.78%	0.69%	0.52%	4.02%		0.46%
Layer 11	0.01%	0.01%	0.40%	0.46%	2.94%	0.15%	0.28%	<b>28.44%</b>		0.21%
Layer 10	0.02%	0.02%	0.35%	0.40%	1.03%	0.11%	0.05%	2.04%		0.12%
Layer 9	0.04%	0.08%	0.44%	0.50%	2.08%	0.14%	0.04%	2.36%		0.15%
Layer 8	0.15%	0.28%	0.07%	0.71%	3.19%	0.31%	0.04%	2.72%		0.22%
Layer 7	0.18%	0.27%	1.49%	1.35%	0.98%	0.71%	0.12%	0.87%		0.55%
Layer 6	0.38%	0.44%	1.64%	1.76%	2.88%	1.11%	0.28%	1.07%		0.82%
Layer 5	0.34%	0.42%	1.32%	1.59%	2.53%	1.14%	0.43%	4.60%		0.83%
Layer 4	0.38%	0.42%	0.68%	1.80%	3.66%	1.12%	0.44%	4.28%		0.80%
Layer 3	0.30%	0.25%	2.22%	2.28%	5.20%	1.17%	0.51%	4.73%		0.98%
Layer 2	0.28%	0.22%	2.32%	2.29%	2.86%	1.01%	0.47%	4.26%		0.92%
Layer 1	0.24%	0.16%	1.51%	1.80%	3.48%	0.78%	0.42%	3.91%		0.71%
Layer 0	0.16%	0.15%	1.86%	2.34%	2.55%	0.51%	0.32%	3.72%		0.65%
attention.self.query										
attention.self.key										
attention.self.value										
attention.output.dense										
attention.output.LayerNorm										
intermediate.dense										
output.dense										
output.LayerNorm										
Layer embeddings	0.03%	1.42%	0.00%	5.13%						0.05%
word_embeddings										
position_embeddings										
token_type_embeddings										
LayerNorm										

(a) PAN16 - Gender

	0.45%	0.47%	1.42%	2.62%	3.57%	1.25%	1.14%	5.76%		0.86%
Layer 11	0.02%	0.02%	1.82%	3.32%	6.72%	1.43%	2.67%	<b>63.19%</b>		1.67%
Layer 10	0.05%	0.04%	1.28%	1.70%	3.70%	0.34%	0.32%	6.07%		0.48%
Layer 9	0.14%	0.07%	0.33%	0.89%	4.77%	0.26%	0.08%	3.27%		0.23%
Layer 8	0.28%	0.51%	0.00%	1.73%	3.07%	0.30%	0.16%	0.25%		0.36%
Layer 7	0.36%	0.59%	0.31%	2.06%	1.90%	1.25%	0.13%	1.93%		0.74%
Layer 6	0.70%	0.84%	2.03%	2.90%	0.00%	1.76%	0.63%	1.30%		1.33%
Layer 5	0.74%	0.88%	0.21%	3.03%	3.75%	1.88%	0.81%	2.96%		1.30%
Layer 4	0.74%	0.74%	1.06%	2.88%	4.32%	1.87%	0.94%	5.59%		1.39%
Layer 3	0.72%	0.63%	1.92%	3.40%	4.04%	1.89%	0.98%	3.79%		1.51%
Layer 2	0.63%	0.49%	2.66%	3.03%	2.01%	1.72%	0.95%	1.74%		1.46%
Layer 1	0.54%	0.41%	2.48%	2.88%	3.40%	1.40%	0.82%	1.26%		1.27%
Layer 0	0.42%	0.36%	2.89%	3.55%	5.13%	0.92%	0.70%	4.11%		1.14%
attention.self.query										
attention.self.key										
attention.self.value										
attention.output.dense										
attention.output.LayerNorm										
intermediate.dense										
output.dense										
output.LayerNorm										
Layer embeddings	0.04%	1.23%	0.00%	6.88%						0.06%
word_embeddings										
position_embeddings										
token_type_embeddings										
LayerNorm										

(b) PAN16 - Age

	1.19%	1.21%	0.94%	1.35%	0.46%	1.48%	1.08%	2.00%		1.00%
Layer 11	0.45%	0.58%	2.03%	3.03%	2.15%	2.22%	2.85%	<b>39.61%</b>		2.04%
Layer 10	1.06%	1.28%	0.46%	1.66%	0.40%	1.50%	1.09%	1.17%		1.24%
Layer 9	1.78%	1.68%	0.42%	1.57%	0.30%	1.31%	0.67%	0.40%		1.11%
Layer 8	1.38%	1.37%	0.39%	1.18%	0.27%	1.45%	0.56%	0.05%		1.03%
Layer 7	1.21%	1.18%	0.82%	1.26%	0.56%	1.52%	0.70%	0.31%		1.11%
Layer 6	1.18%	1.22%	1.06%	1.18%	0.21%	1.50%	0.80%	0.00%		1.15%
Layer 5	1.07%	1.13%	1.03%	1.17%	0.10%	1.41%	0.79%	0.05%		1.10%
Layer 4	1.14%	1.20%	0.77%	1.06%	0.23%	1.33%	0.78%	0.27%		1.05%
Layer 3	1.28%	1.31%	0.93%	1.06%	0.35%	1.40%	0.88%	0.13%		1.14%
Layer 2	1.11%	0.91%	0.89%	0.64%	0.14%	1.36%	0.93%	0.00%		1.06%
Layer 1	1.25%	1.25%	0.98%	1.05%	0.25%	1.41%	0.94%	0.22%		1.16%
Layer 0	1.42%	1.37%	1.53%	1.29%	0.49%	1.37%	1.11%	0.31%		1.29%
attention.self.query										
attention.self.key										
attention.self.value										
attention.output.dense										
attention.output.LayerNorm										
intermediate.dense										
output.dense										
output.LayerNorm										
Layer embeddings	0.23%	0.61%	0.00%	1.48%						0.23%
word_embeddings										
position_embeddings										
token_type_embeddings										
LayerNorm										

(c) BIOS - Gender

	0.62%	0.66%	2.33%	2.65%	2.22%	1.36%	1.23%	2.64%		1.00%
Layer 11	0.06%	0.09%	2.96%	3.64%	4.24%	1.24%	1.70%	<b>21.71%</b>		1.43%
Layer 10	0.11%	0.09%	1.84%	2.19%	3.55%	0.44%	0.46%	4.32%		0.65%
Layer 9	0.33%	0.45%	1.93%	2.15%	2.16%	0.49%	0.35%	0.27%		0.68%
Layer 8	0.73%	0.89%	1.90%	2.55%	1.13%	0.95%	0.42%	0.52%		0.96%
Layer 7	0.46%	0.58%	2.22%	2.50%	1.74%	1.39%	0.63%	0.66%		1.15%
Layer 6	0.92%	1.08%	2.80%	3.06%	1.95%	1.74%	0.87%	0.49%		1.52%
Layer 5	0.91%	1.01%	2.59%	2.89%	1.08%	1.86%	0.97%	0.53%		1.56%
Layer 4	0.98%	1.05%	2.51%	2.67%	1.69%	1.80%	0.98%	0.99%		1.53%
Layer 3	0.91%	0.95%	2.56%	2.84%	2.24%	1.93%	1.11%	1.67%		1.62%
Layer 2	0.71%	0.64%	2.39%	2.53%	2.01%	1.73%	1.09%	1.63%		1.46%
Layer 1	0.73%	0.60%	1.84%	2.14%	2.37%	1.61%	1.03%	2.12%		1.32%
Layer 0	0.54%	0.51%	2.40%	2.68%	2.51%	1.10%	0.84%	1.76%		1.16%
attention.self.query										
attention.self.key										
attention.self.value										
attention.output.dense										
attention.output.LayerNorm										
intermediate.dense										
output.dense										
output.LayerNorm										
Layer embeddings	0.07%	0.52%	0.00%	2.42%						0.08%
word_embeddings										
position_embeddings										
token_type_embeddings										
LayerNorm										

(d) FCDL18 - Dialect

Figure 3: Sparsity rate of the subnetworks in MODDIFFY-POST on BERT-Base. Each value shows the sparsity rate on the specific block/matrix.

	1.14%	1.22%	1.30%	1.54%	0.84%	1.54%	1.07%	0.66%	1.00%
Layer 11	0.06%	0.09%	0.87%	1.13%	1.74%	0.17%	0.17%	0.04%	0.27%
Layer 10	0.08%	0.18%	0.95%	0.83%	1.02%	0.18%	0.13%	1.33%	0.27%
Layer 9	0.31%	0.44%	0.64%	0.80%	0.52%	0.34%	0.21%	0.34%	0.37%
Layer 8	1.20%	1.53%	0.27%	1.32%	0.21%	0.77%	0.36%	0.59%	0.74%
Layer 7	1.39%	1.58%	0.86%	1.60%	0.62%	1.73%	0.71%	0.43%	1.26%
Layer 6	1.86%	2.07%	1.89%	1.88%	0.18%	2.35%	1.03%	0.52%	1.77%
Layer 5	1.86%	1.97%	2.06%	2.04%	0.61%	2.43%	1.32%	0.52%	1.91%
Layer 4	1.68%	1.80%	1.86%	1.93%	1.47%	2.34%	1.35%	0.65%	1.84%
Layer 3	1.55%	1.59%	1.67%	1.77%	1.05%	2.40%	1.60%	0.01%	1.88%
Layer 2	1.21%	1.13%	1.52%	1.55%	0.79%	2.22%	1.57%	0.62%	1.71%
Layer 1	1.24%	1.10%	1.08%	1.55%	0.76%	1.94%	1.41%	0.25%	1.53%
Layer 0	1.19%	1.13%	1.99%	2.09%	1.13%	1.64%	1.51%	0.44%	1.58%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.08%	0.29%	0.00%	1.71%	0.09%				
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(a) PAN16 - Gender

	1.16%	1.23%	1.21%	1.52%	0.77%	1.53%	1.07%	2.01%	1.00%
Layer 11	0.00%	0.07%	1.71%	2.44%	1.84%	0.98%	0.71%	33.75%	0.85%
Layer 10	0.03%	0.08%	0.82%	1.01%	0.65%	0.32%	0.39%	0.68%	0.40%
Layer 9	0.13%	0.23%	0.14%	0.86%	0.23%	0.25%	0.26%	0.70%	0.28%
Layer 8	0.90%	1.16%	0.29%	1.33%	0.16%	0.43%	0.22%	0.18%	0.52%
Layer 7	1.19%	1.34%	0.40%	1.18%	0.61%	1.31%	0.50%	0.38%	0.94%
Layer 6	1.95%	2.07%	1.45%	1.50%	0.12%	2.02%	0.75%	0.36%	1.50%
Layer 5	1.95%	2.10%	1.76%	1.76%	0.55%	2.29%	1.12%	0.33%	1.77%
Layer 4	1.82%	1.95%	1.65%	1.70%	0.91%	2.33%	1.29%	0.70%	1.80%
Layer 3	1.75%	1.86%	1.57%	1.63%	0.64%	2.38%	1.54%	0.36%	1.87%
Layer 2	1.44%	1.35%	1.47%	1.47%	0.73%	2.26%	1.62%	0.17%	1.77%
Layer 1	1.40%	1.26%	1.28%	1.42%	0.86%	2.03%	1.49%	0.66%	1.62%
Layer 0	1.36%	1.28%	1.94%	1.98%	1.91%	1.75%	1.61%	0.66%	1.67%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.11%	0.29%	0.00%	1.12%	0.11%				
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(b) PAN16 - Age

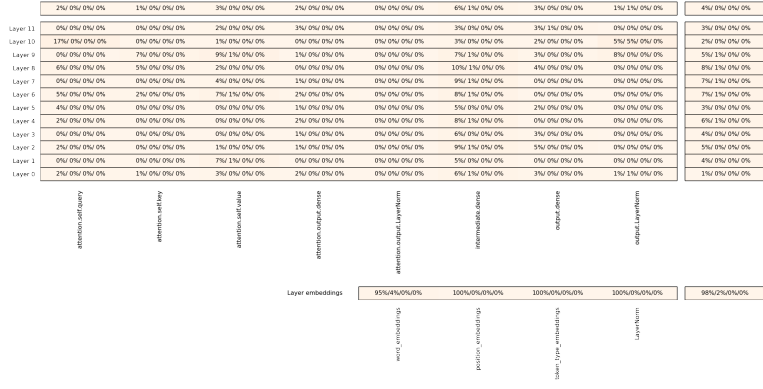
	1.36%	1.37%	0.90%	1.28%	0.29%	1.53%	0.96%	0.27%	1.00%
Layer 11	0.63%	0.87%	1.34%	1.74%	0.60%	1.42%	1.35%	0.64%	1.21%
Layer 10	1.66%	1.54%	0.00%	1.41%	0.16%	1.45%	0.81%	0.85%	1.14%
Layer 9	1.92%	1.86%	0.98%	1.45%	0.14%	1.35%	0.56%	0.16%	1.15%
Layer 8	1.54%	1.54%	0.00%	1.26%	0.26%	1.59%	0.53%	0.25%	1.07%
Layer 7	1.39%	1.37%	0.80%	1.36%	0.39%	1.70%	0.73%	0.20%	1.22%
Layer 6	1.33%	1.39%	1.17%	1.28%	0.21%	1.67%	0.83%	0.14%	1.26%
Layer 5	1.22%	1.27%	1.13%	1.25%	0.36%	1.55%	0.82%	0.13%	1.20%
Layer 4	1.27%	1.16%	1.02%	1.15%	0.42%	1.47%	0.81%	0.26%	1.15%
Layer 3	1.39%	1.46%	0.77%	1.12%	0.00%	1.54%	0.91%	0.29%	1.21%
Layer 2	1.15%	1.18%	0.93%	0.78%	0.17%	1.50%	0.98%	0.00%	1.17%
Layer 1	1.35%	1.35%	0.99%	1.10%	0.21%	1.55%	0.99%	0.01%	1.25%
Layer 0	1.54%	1.44%	1.67%	1.44%	0.59%	1.52%	1.23%	0.13%	1.42%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.26%	0.70%	0.00%	1.15%	0.27%				
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(c) BIOS - Gender

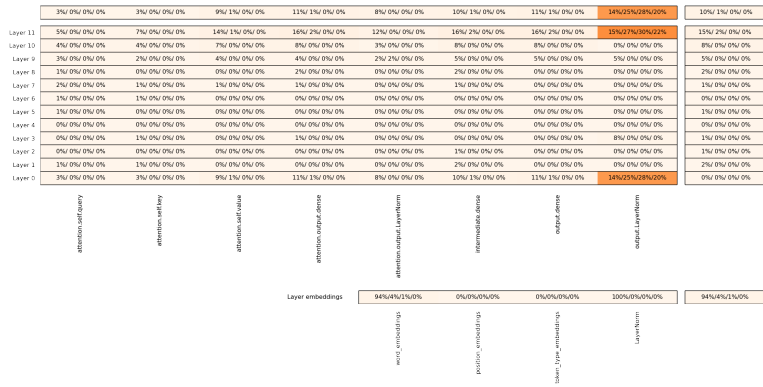
	0.90%	1.01%	1.98%	2.17%	1.10%	1.35%	1.12%	1.20%	1.00%
Layer 11	0.01%	0.02%	2.91%	3.52%	1.47%	1.04%	1.30%	1.38%	1.22%
Layer 10	0.05%	0.05%	2.12%	2.38%	1.52%	0.59%	0.46%	1.77%	0.74%
Layer 9	0.06%	0.15%	1.42%	1.57%	0.59%	0.46%	0.32%	2.07%	0.53%
Layer 8	0.10%	0.23%	1.01%	1.67%	0.25%	0.40%	0.23%	0.49%	0.46%
Layer 7	0.09%	0.22%	0.62%	1.42%	0.87%	0.52%	0.25%	0.95%	0.45%
Layer 6	0.44%	0.75%	1.50%	1.60%	0.81%	0.68%	0.26%	1.09%	0.67%
Layer 5	0.64%	0.98%	1.58%	1.51%	0.52%	0.98%	0.49%	1.08%	0.88%
Layer 4	1.57%	1.89%	1.76%	1.83%	1.24%	1.38%	0.59%	1.56%	1.24%
Layer 3	2.18%	2.43%	2.40%	2.17%	1.15%	2.48%	1.08%	0.98%	1.95%
Layer 2	1.90%	1.82%	2.63%	2.44%	1.80%	2.77%	1.70%	0.83%	2.22%
Layer 1	2.05%	1.99%	2.41%	2.46%	1.37%	2.73%	1.80%	1.46%	2.25%
Layer 0	1.69%	1.60%	3.43%	3.41%	1.68%	2.20%	1.77%	1.80%	2.17%
attention.self.query									
attention.self.key									
attention.self.value									
attention.output.dense									
attention.output.LayerNorm									
intermediate.dense									
output.dense									
output.LayerNorm									
Layer embeddings	0.16%	0.60%	0.00%	3.24%	0.17%				
word_embeddings									
position_embeddings									
token_type_embeddings									
LayerNorm									

(d) FCDL18 - Dialect

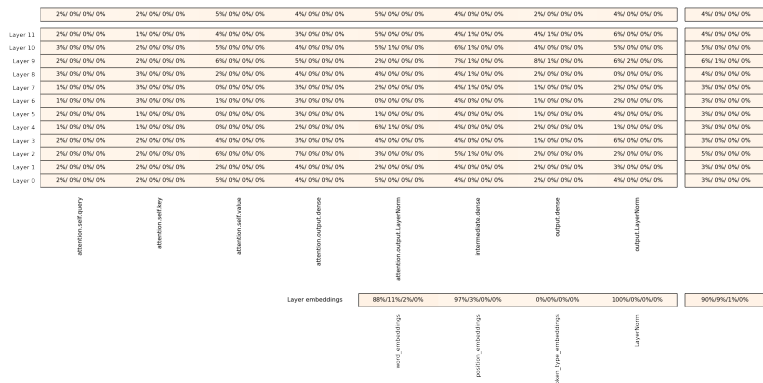
Figure 4: Sparsity rate of the subnetworks in DIFFPRUN-DEBIAS using BERT-Base. Each value shows the sparsity rate on the specific block/matrix.



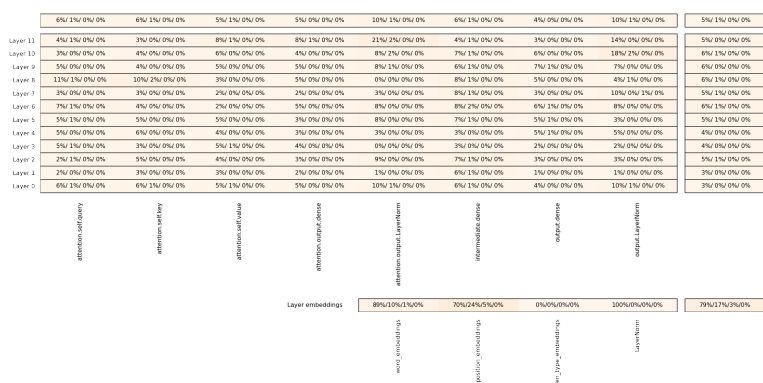
(a) PAN16 - Gender



(b) PAN16 - Age

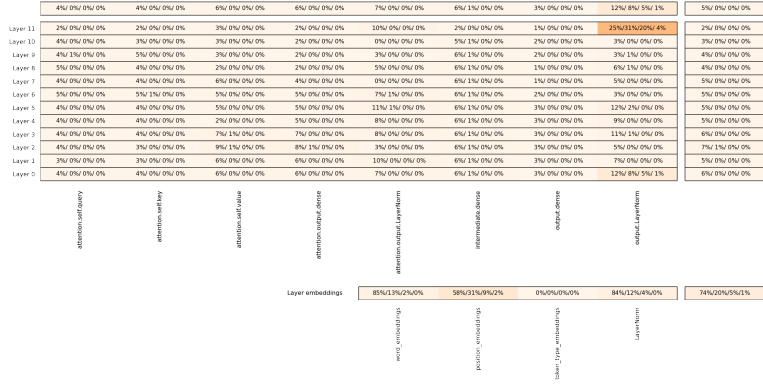


(c) BIOS - Gender

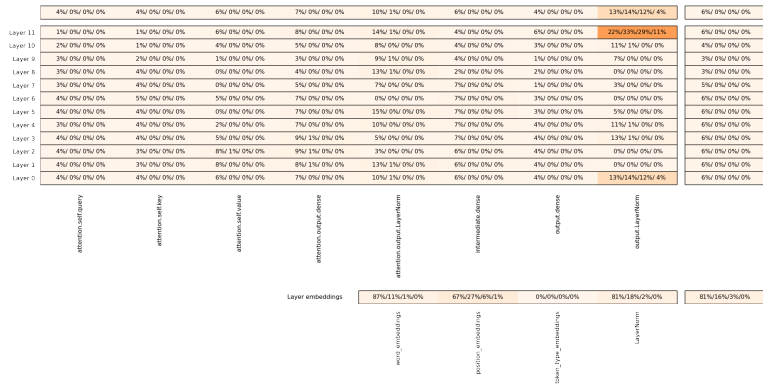


(d) FCDL18 - Dialect

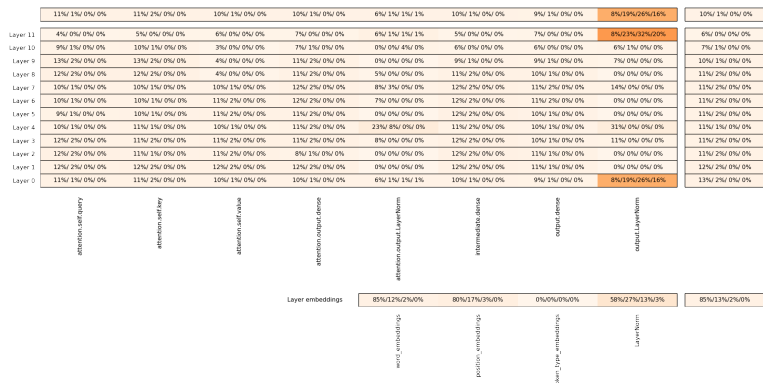
Figure 5: The percentage of common non-masked parameters for MODDIFFY-PAR across 5 runs with different initialization seeds. For each block, the numbers indicate the percentage of the number of common parameters across two, three, four, and five runs of a subnetwork, respectively.



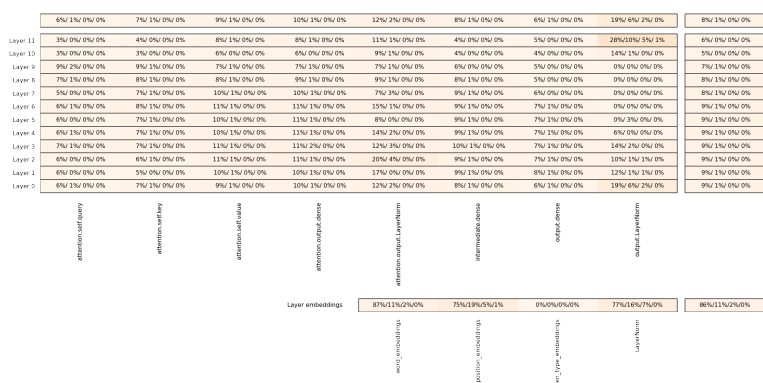
(a) PAN16 - Gender



(b) PAN16 - Age



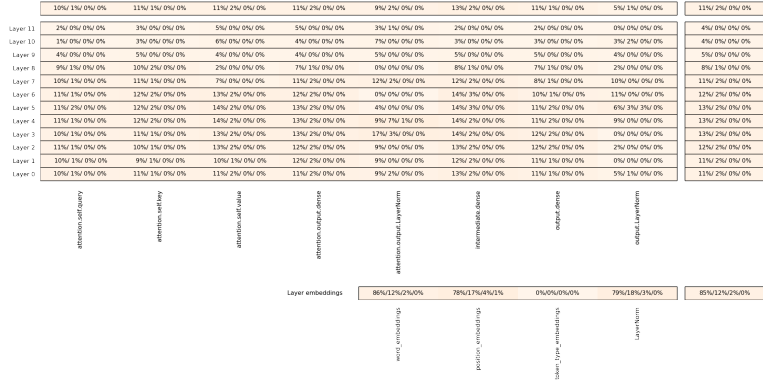
(c) BIOS - Gender



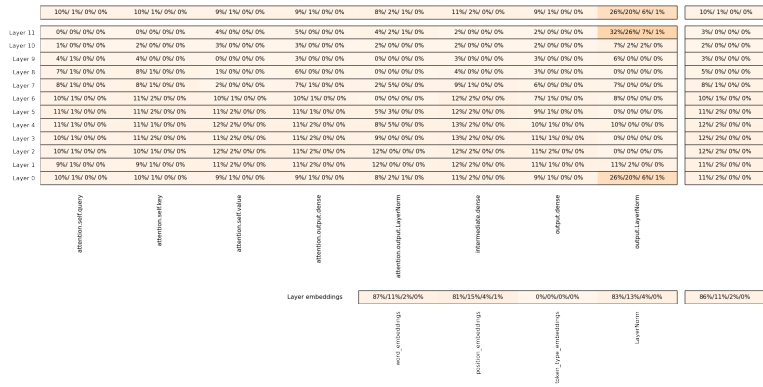
(d) FCDL18 - Dialect

Figure 6: The percentage of common non-masked parameters for MODDIFFY-POST across 5 runs with different initialization seeds. For each block, the numbers indicate the percentage of the number of common parameters across two, three, four, and five runs of a subnetwork, respectively.

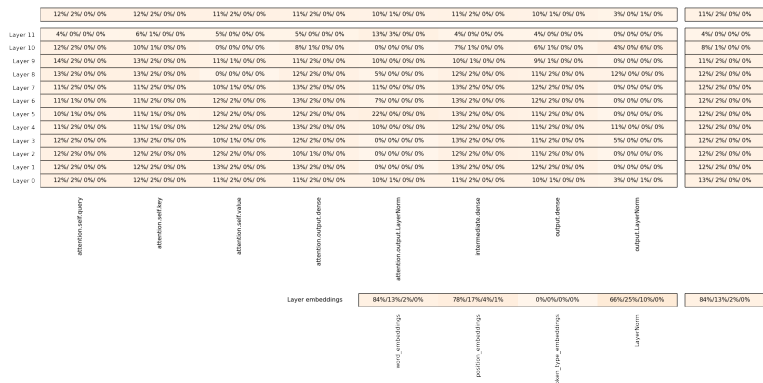




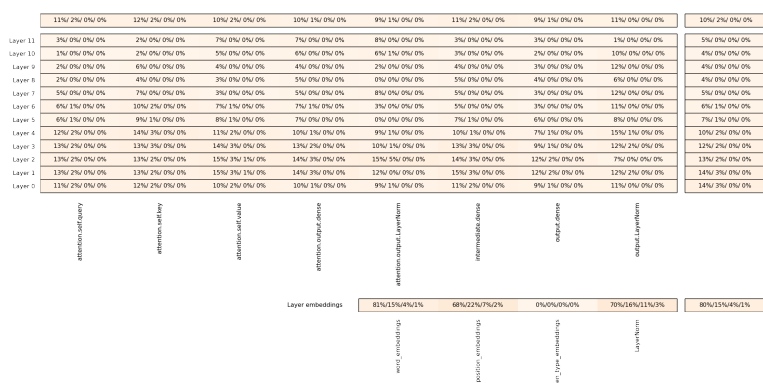
(a) PAN16 - Gender



(b) PAN16 - Age



(c) BIOS - Gender



(d) FCDL18 - Dialect

Figure 7: The percentage of common non-masked parameters for DIFFPRUN-DEBIAS across 5 runs with different initialization seeds. For each block, the numbers indicate the percentage of the number of common parameters across two, three, four, and five runs of a subnetwork, respectively.

## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?  
*Left blank.*
- A2. Did you discuss any potential risks of your work?  
*Left blank.*
- A3. Do the abstract and introduction summarize the paper’s main claims?  
*Left blank.*
- A4. Have you used AI writing assistants when working on this paper?  
*Left blank.*

### B Did you use or create scientific artifacts?

*Section 4, 7 and Appendix A*

- B1. Did you cite the creators of artifacts you used?  
*Section 4*
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?  
*Left blank.*
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?  
*Section 4, 7 and Appendix A*
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?  
*Section 4 and Appendix A*
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?  
*Section 4 and Appendix A*
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.  
*Section 4 and Appendix A*

### C Did you run computational experiments?

*Section 4, 5, and Appendix A*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?  
*Appendix A*

---

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

*Section 4 and Appendix A*

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

*Section 5, and Appendix B*

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

*No response.*

**D  Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

*No response.*

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

*No response.*

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

*No response.*

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

*No response.*

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*No response.*