

Discover and Mitigate Unknown Biases with Debiasing Alternate Networks

Zhiheng Li¹, Anthony Hoogs², and Chenliang Xu¹

¹University of Rochester ²Kitware, Inc.

{zhiheng.li,chenliang.xu}@rochester.edu anthony.hoogs@kitware.com

Abstract. Deep image classifiers have been found to learn biases from datasets. To mitigate the biases, most previous methods require labels of protected attributes (e.g., age, skin tone) as full-supervision, which has two limitations: 1) it is infeasible when the labels are unavailable; 2) they are incapable of mitigating unknown biases—biases that humans do not preconceive. To resolve those problems, we propose Debiasing Alternate Networks (DebiAN), which comprises two networks—a Discoverer and a Classifier. By training in an alternate manner, the discoverer tries to find multiple unknown biases of the classifier without any annotations of biases, and the classifier aims at unlearning the biases identified by the discoverer. While previous works evaluate debiasing results in terms of a single bias, we create Multi-Color MNIST dataset to better benchmark mitigation of multiple biases in a multi-bias setting, which not only reveals the problems in previous methods but also demonstrates the advantage of DebiAN in identifying and mitigating multiple biases simultaneously. We further conduct extensive experiments on real-world datasets, showing that the discoverer in DebiAN can identify unknown biases that may be hard to be found by humans. Regarding debiasing, DebiAN achieves strong bias mitigation performance.

Keywords: Bias Identification, Bias Mitigation, Fairness, Unsupervised Debiasing

1 Introduction

Many studies have verified that AI algorithms learn undesirable biases from the dataset. Some biases provide shortcuts [18] for the network to learn superficial features instead of the intended decision rule causing robustness issues, *e.g.*, static cues for action recognition [7,11,40]. Other biases make AI algorithms discriminate against different protected demographic groups such as genders* [3,25–27,58,61,66] and skin tones [9,23], leading to serious fairness problems. Therefore, it is imperative to mitigate the biases in AI algorithms. However, most previous bias mitigation methods [4,54,60,64,66] are supervised methods—requiring annotations of the biases, which has several limitations: First, bias mitigation cannot be performed when labels are not available due to privacy concerns. Second,

*In this work, “gender” denotes visually perceived gender, not real gender identity.

they cannot mitigate *unknown* biases—biases that humans did not preconceive, making the biases impossible to be labeled and mitigated.

Since supervised debiasing methods present many disadvantages, in this work, we focus on a more challenging task—unsupervised debiasing, which mitigates the *unknown* biases in a learned classifier without any annotations. Without loss of generality, we focus on mitigating biases in image classifiers. Solving this problem contains two steps [2,14,36,45,52]: bias identification and bias mitigation.

Due to the absence of bias annotations, the first step is to assign the training samples into different bias groups as the pseudo bias labels, which is challenging since the biases are even unknown. The crux of the problem is to define the unknown bias. Some previous works make strong assumptions about the unknown biases based on empirical observations, such as biases are easier to be learned [45], samples from the same bias group are clustered in feature space [52], which can be tenuous for different datasets or networks. Other works quantify the unknown biases by inversely using the debiasing objective functions [2,14], which can face numerical or convergence problems (more details in Sec. 2). Unlike previous works, we follow an axiomatic principle to define the unknown biases—classifier’s predictions that violate a fairness criterion [13,17,21,22,35,46,55]. Based on this definition, we propose a novel *Equal Opportunity Violation* (EOV) loss to train a *discoverer* network to identify the classifier’s biases. In specific, it shepherds the *discoverer* network to predict bias group assignments such that the *classifier* violates the Equal Opportunity [22,46] fairness criterion (Figs. 1, 2).

Regarding debiasing as the second step, most previous approaches [2,14,52] preprocess the identified biases into pseudo bias labels and resort to other supervised bias mitigation methods [6,47] for debiasing. In contrast, we propose a novel *Reweighted Cross-Entropy* (RCE) loss that leverages soft bias group assignments predicted by the *discoverer* network to mitigate the biases in the *classifier* (Fig. 1). In this way, the *classifier* is guided to meet the Equal Opportunity.

In addition, many previous works [2,14,52] treat bias identification and bias mitigation as two isolated steps. In [2,14], the biases are identified from an undertrained classifier, which is suboptimal since the classifier may learn different biases at different training stages. Consequently, these two-stage methods fail to mitigate other biases learned by the classifier at later training stages. In contrast, we employ an alternate training scheme to carry out bias identification and bias mitigation simultaneously. We jointly update the *discoverer* and *classifier* in an interleaving fashion (Figs. 1 and 2). In this way, the *discoverer* can repetitively inspect multiple biases that the *classifier* learns at the entire training stage.

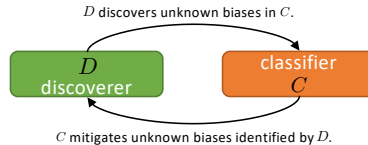


Fig. 1: DEBIASING ALTERNATE NETWORKS (DebiAN). We alternately train two networks—a *discoverer* and a *classifier*. *Discoverer* actively identifies *classifier*’s unknown biases. At the same time, the *classifier* mitigates the biases identified by the *discoverer*.

We integrate our novel losses and training scheme into a unified framework—DEBIASING ALTERNATE NETWORKS (DebiAN), which contains two networks—a *discoverer* D and a *classifier* C (see Fig. 1). We jointly train the two networks in an alternate manner. Supervised by our novel EOv loss, D tries to discover C ’s multiple unknown biases that violate the Equal Opportunity fairness criterion. Trained with our RCE loss, C aims at mitigating multiple biases identified by the *discoverer* D to satisfy Equal Opportunity. After the alternate training, the unknown biases in *classifier* C are mitigated, leading to a fairer and more robust classification model. Besides, when employed with other network explanation methods [49,50,67], the *discoverer* is helpful to interpret the discovered unknown biases, facilitating dataset curators to locate dataset biases [53].

While previous works [7,31,41,45,47] only evaluate debiasing results in terms of a single bias, we create Multi-Color MNIST dataset with two biases in the dataset, which benchmarks debiasing algorithms in the multi-bias setting. Our new dataset surfaces the problems in previous methods (*e.g.*, LfF [45]) and demonstrates the advantage of DebiAN in discovering and mitigating multiple biases. We further conduct extensive experiments to verify the efficacy of DebiAN in real-world image datasets. In the face image domain, DebiAN achieves better gender bias mitigation results on CelebA [43] and bFFHQ [32] datasets. On the gender classification task, DebiAN achieves better debiasing results on CelebA w.r.t. multiple bias attributes. We further show an interesting unknown bias discovered by DebiAN in gender classification—**visible hair area**. Lastly, we show that DebiAN applies to other image domains for broader tasks, such as action recognition and scene classification. Our method not only achieves better debiasing results, but also identifies interesting unknown biases in scene classifiers.

Our contributions are summarized as follows: (1) We propose a novel objective function, *Equal Opportunity Violation* (EOv) loss, for identifying unknown biases of a classifier based on Equal Opportunity. (2) We propose a *Reweighted Cross-Entropy* (RCE) loss to mitigate the discovered unknown biases by leveraging the soft bias group assignments. (3) We create Multi-Color MNIST dataset to benchmark debiasing algorithms in a multi-bias setting. (4) Our DEBIASING ALTERNATE NETWORKS (DebiAN) outperforms previous unsupervised debiasing methods on both synthetic and real-world datasets.

2 Related Work

Bias Identification Most previous works identify *known* biases based on bias labels. In [9], face images are labeled with gender and skin tone to identify the performance gaps across intersectional groups. Balakrishnan *et al.* [8] further synthesize intersectional groups of images and analyze the biases with additional labels. Beyond face images, recent works [44,56] compute the statistics of labels based on the rule mining algorithm [1] or external tools. [34] uses clustering on image embeddings to discover unknown biases. [37,42] discovers *unknown* biases without labels. However, these works rely on GAN [20,29] to synthesize images,

which suffers from image quality issues. In contrast, DebiAN directly classifies real images into different bias attribute groups to discover the unknown biases.

Supervised Debiasing Supervised debiasing methods use bias labels for debiasing. [28] proposes a supervised reweighing method. Wang *et al.* [62] benchmark recent supervised debiasing methods [4,54,64,66]. [15] lets the model be flexibly fair to different attributes during testing. [48] uses disentanglement for debiasing. Singh *et al.* [51] propose a feature splitting approach to mitigate contextual bias. [16,19] use adversarial training to mitigate biases in face recognition.

Known Bias Mitigation with Prior knowledge Without using labels, some works use prior knowledge to mitigate certain known biases. ReBias [7] uses model capacity as the inductive bias to mitigate texture bias and static bias in image and video classification. HEX [57] introduces a texture extractor to mitigate the texture bias. Beyond image classification, RUBi [10] and LearnedMixIn [12] mitigate unimodal bias for visual question answering [5] with prior knowledge.

Unsupervised Debiasing In the field of mitigating unknown biases, Sohoni *et al.* [52] apply clustering on samples in each class and use the clustering assignment as the predicted bias labels, which could be inaccurate due to its unsupervised nature. Li *et al.* [39,40] fix the parameters of feature extractors and focus on mitigating the representation bias. LfF [45] identifies biases by finding easier samples in the training data through training a bias-amplified network supervised by GCE loss [65], which up-weights the samples with smaller loss values and down-weights the samples with larger loss values. In other words, GCE loss does not consider the information of the classifier, *e.g.*, the classifier’s output. Therefore, LfF’s bias-amplified network blindly finds the biases in the data samples instead of the classifier. Unlike LfF, the EOv loss in DebiAN actively identifies biases in the classifier based on the classifier’s predictions, leading to better debiasing performance. Following LfF, BiaSwap [32] uses LfF to discover biases and generate more underrepresented images via style-transfer for training. Other works [2,14,36,59] inversely use the debiasing objective function to maximize an unbounded loss (*e.g.*, gradient norm penalty in IRMv1 [6]) for bias identification, which may encounter numerical or convergence problems. As a comparison, our EOv loss (Eq. (2)) minimizes negative log-likelihood, which is numerically stable and easier to converge.

3 Method

Overview The overview of our proposed DEBIASING ALTERNATE NETWORKS (DebiAN) is shown in Fig. 2. It contains two networks—a *discoverer* D and a *classifier* C . As shown in Fig. 2 (a), the *discoverer* D tries to discover the unknown biases in the *classifier* C by optimizing our proposed EOv loss (\mathcal{L}_{EOv}) and UA penalty (\mathcal{L}_{UA}) (Sec. 3.1). As shown in Fig. 2 (b), the *classifier* C ’s goal is to mitigate the biases identified by D via a novel *Reweighted Cross-Entropy* loss (\mathcal{L}_{RCE}) (Sec. 3.2). Lastly, we train the two networks in an alternate manner as the full model for discovering and mitigating the unknown biases (Sec. 3.3).

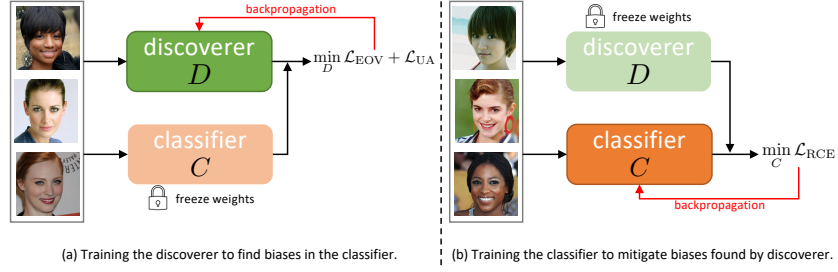


Fig. 2: Overview of DEBIASING ALTERNATE NETWORKS (DebiAN). DebiAN consists of two networks—a *discoverer* D and a *classifier* C . D is trained with \mathcal{L}_{EOV} and \mathcal{L}_{UA} (Sec. 3.1) to find the unknown biases in C . C is optimized with \mathcal{L}_{RCE} (Sec. 3.2) to mitigate the biases identified by D .

Background To better explain our motivation for discovering the *unknown* biases (without manual annotations of biases), let us first revisit the traditional approach for identifying *known* biases when labels of biases (*e.g.*, protected attributes) are available, which is illustrated in Fig. 3 (a). The following are given for identifying *known* biases—a well-trained *classifier* C for predicting a target attribute, n testing images $\{\mathbf{I}_i\}_{i=1}^n$, target attribute labels of each image $\{y_i\}_{i=1}^n$, and bias attribute labels $\{b_i\}_{i=1}^n$. We denote the i -th image target attribute as $y_i \in \{1, 2, \dots, K\}$ and K is the number of classes. We consider the bias attribute that is binary or continuously valued (*i.e.*, $b_i \in \{0, 1\}$ or $b_i \in [0, 1]$), such as biological gender (*e.g.*, female and male) and skin tones (*e.g.*, from dark skin tones to light skin tones in Fitzpatrick skin type scale). We leave bias attributes with multi-class values for future works. Then, the given *classifier* C is tested for predicting the target attribute \hat{y}_i for each testing image \mathbf{I}_i . Finally, we check whether the predictions meet a fairness criterion, such as Equal Opportunity [22]:

$$\Pr\{\hat{y} = k \mid b = 0, y = k\} = \Pr\{\hat{y} = k \mid b = 1, y = k\}, \quad (1)$$

where the LHS and RHS are true positive rates (TPR) in negative ($b = 0$) and positive ($b = 1$) bias attribute groups, respectively. $k \in \{1 \dots K\}$ is a target attribute class. Equal Opportunity requires the same TPR across two different bias attribute groups. That is, if the TPR is significantly different in two groups of the bias attribute, we conclude that *classifier* C contains the bias of attribute b because C violates the Equal Opportunity fairness criterion. For example, as shown in Fig. 3 (a), although all images are female, a gender classifier may have a larger TPR for the group of long-hair female images than the group of short-hair female images. Thus the gender classifier is biased against different hair lengths.

3.1 Unknown Bias Discovery

As for identifying unknown biases, we do not have the labels to assign images into two groups for comparing TPR since 1) we do not assume images come with bias

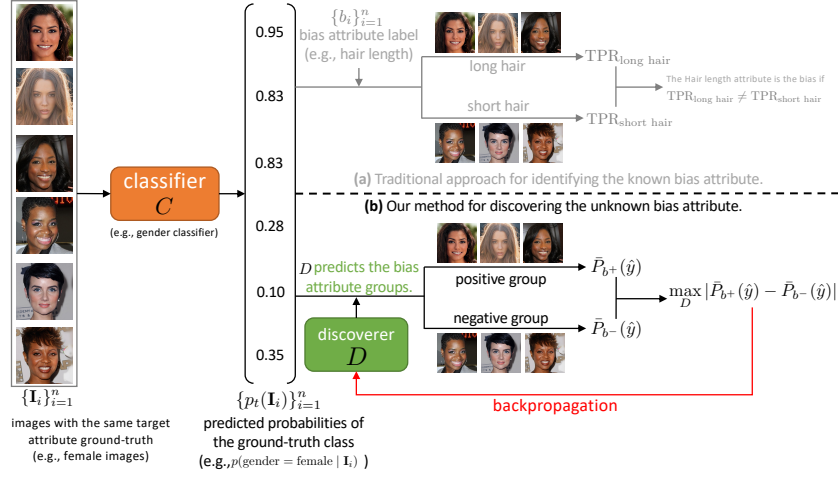


Fig. 3: (a): The traditional approach for identifying the *known* bias attribute (e.g., hair length) by comparing true positive rates (TPR) of the target attribute (e.g., gender) in two groups of bias attributes (e.g., long hair and short hair), where the group assignment of bias attribute is based on the labels of the bias attribute. (b): Our method trains a *discoverer* (D) to predict the groups of the unknown bias attribute such that the difference of averaged predicted probabilities on the target attribute (e.g., gender) in two groups are maximized (see Eq. 2)

attribute labels, and 2) the type of bias is even unknown. However, we can compare the difference in TPR for any group assignments based on speculated biases—a significant difference in TPR hints that the Equal Opportunity fairness criterion is most violated (our method mainly focuses on the Equal Opportunity fairness criterion, and we leave other fairness criteria for future work). Motivated by this finding, instead of using labels of bias attribute $\{b_i\}_{i=1}^n$ for group assignment, we train a *discoverer* D to predict the group assignment for each image, i.e., $p(\hat{b} | \mathbf{I}_i) := D(\mathbf{I}_i)$. By optimizing loss functions below, we find the most salient bias of the *classifier* C that violates the Equal Opportunity fairness criterion, which is illustrated in Fig. 3 (b).

Equal Opportunity Violation (EOV) Loss To shepherd the *discoverer* D to find the group assignment where *classifier* C violates the Equal Opportunity fairness criterion, we propose the *Equal Opportunity Violation* (EOV) loss, denoted by \mathcal{L}_{EOV} , as the objective function to train D . For computing \mathcal{L}_{EOV} , we sample a set of n images $\{\mathbf{I}_i\}_{i=1}^n$ with the *same* target attribute labels (i.e., $\forall_i y_i = k$), e.g., all images in Fig. 3 (b) are female. The *classifier* C has been trained for predicting the target attribute y of the images (i.e., $p(\hat{y} | \mathbf{I}_i) := C(\mathbf{I}_i)$). For simplicity, we denote p_t as C 's prediction on images of the ground-truth class (i.e., $p_t(\mathbf{I}_i) = p(\hat{y} = y_i | \mathbf{I}_i)$). Meanwhile, the same set of images $\{\mathbf{I}_i\}$ are fed to the *discoverer* D for predicting the binary bias attribute group assignment:

$p(\hat{b} \mid \mathbf{I}_i) := D(\mathbf{I}_i)$. Finally, we define the EOV loss as:

$$\mathcal{L}_{\text{EOV}} = -\log \left(\left| \bar{P}_{b+}(\hat{y}) - \bar{P}_{b-}(\hat{y}) \right| \right), \quad (2)$$

where $\bar{P}_{b+}(\hat{y})$ and $\bar{P}_{b-}(\hat{y})$ are defined by:

$$\begin{aligned} \bar{P}_{b+}(\hat{y}) &= \frac{\sum_{i=1}^n p(\hat{b} = 1 \mid \mathbf{I}_i) p_t(\mathbf{I}_i)}{\sum_{i=1}^n p(\hat{b} = 1 \mid \mathbf{I}_i)}, \\ \bar{P}_{b-}(\hat{y}) &= \frac{\sum_{i=1}^n p(\hat{b} = 0 \mid \mathbf{I}_i) p_t(\mathbf{I}_i)}{\sum_{i=1}^n p(\hat{b} = 0 \mid \mathbf{I}_i)}. \end{aligned} \quad (3)$$

Intuitively, $\bar{P}_{b+}(\hat{y})$ and $\bar{P}_{b-}(\hat{y})$ are the weighted average predicted probabilities of the target attribute in two bias attribute groups, which can be regarded as a relaxation to Equal Opportunity’s true positive rate (Eq. 1) where the predicted probabilities are binarized into predictions with a threshold (*e.g.*, 0.5). Minimizing \mathcal{L}_{EOV} leads D to maximize the discrepancy of averaged predicted probabilities of target attributes in two bias attribute groups (*i.e.*, see $\max_D |\bar{P}_{b+}(\hat{y}) - \bar{P}_{b-}(\hat{y})|$ in Fig. 3), thus finding the bias attribute group assignments where C violates the Equal Opportunity fairness criterion. For example, in Fig. 3 (b), if the gender classifier C is biased against different hair lengths, then by optimizing \mathcal{L}_{EOV} , D can assign the female images into two bias attribute groups (*i.e.*, short hair and long hair) with the predicted bias attribute group assignment probability $p(\hat{b} \mid \mathbf{I}_i)$, such that the difference of averaged predicted probabilities on gender in these two groups is maximized.

Unbalanced Assignment (UA) penalty However, we find that optimizing \mathcal{L}_{EOV} alone may let the *discoverer* D find a trivial solution—assigning all images into one bias attribute group. For example, suppose D assigns all images to the positive bias attribute group (*i.e.*, $\forall_i, p(\hat{b} = 1 \mid \mathbf{I}_i) = 1$). In that case, $\bar{P}_{b-}(\hat{y})$ becomes zero since the negative group contains no images. $\bar{P}_{b+}(\hat{y})$ becomes a large positive number by simply averaging $p_t(\mathbf{I}_i)$ for all of the n images, which can trivially increase $|\bar{P}_{b+}(\hat{y}) - \bar{P}_{b-}(\hat{y})|$, leading to a small \mathcal{L}_{EOV} . To prevent this trivial solution, we propose the *Unbalanced Assignment* (UA) loss denoted by:

$$\mathcal{L}_{\text{UA}} = -\log \left(1 - \frac{1}{n} \left| \sum_{i=1}^n p(\hat{b} = 1 \mid \mathbf{I}_i) - p(\hat{b} = 0 \mid \mathbf{I}_i) \right| \right). \quad (4)$$

Intuitively, minimizing \mathcal{L}_{UA} penalizes the unbalanced assignment that leads to large difference between $\sum_{i=1}^n p(\hat{b} = 1 \mid \mathbf{I}_i)$ and $\sum_{i=1}^n p(\hat{b} = 0 \mid \mathbf{I}_i)$, which can be regarded as the numbers of images assigned into positive and negative bias attribute groups, respectively. Therefore, \mathcal{L}_{EOV} is jointly optimized with \mathcal{L}_{UA} to prevent the trivial solution. We acknowledge a limitation of the UA penalty. Although it resolves the trivial solution, it introduces a trade-off since the bias attribute groups are usually spuriously correlated with the target attribute (*e.g.*, more long-hair females than the short-hair females in the dataset). Hence encouraging balanced assignments may make the *discoverer* harder to find the correct assignment. However, our ablation study shows that the benefits of using \mathcal{L}_{UA} outweigh its limitations. The results are shown in Sec. 4.1 and Tab. 1.

3.2 Unknown Bias Mitigation by Reweighting

We further mitigate C 's unknown biases identified by D . To this end, we propose a novel *Reweighted Cross-Entropy* loss that adjusts the weight of each image's classification loss. Based on the bias attribute group assignment $p(\hat{b} \mid \mathbf{I}_i)$ predicted by D , we define the weight $\mathcal{W}(\mathbf{I}_i)$ of classification loss for each image \mathbf{I}_i as:

$$\begin{aligned} \mathcal{W}(\mathbf{I}_i) = & \mathbb{1} [\bar{P}_{b+}(\hat{y}) \geq \bar{P}_{b-}(\hat{y})] p(\hat{b} = 0 \mid \mathbf{I}_i) \\ & + \mathbb{1} [\bar{P}_{b+}(\hat{y}) < \bar{P}_{b-}(\hat{y})] p(\hat{b} = 1 \mid \mathbf{I}_i), \end{aligned} \quad (5)$$

where $\mathbb{1}$ is an indicator function. Then, the *Reweighted Cross-Entropy* loss (\mathcal{L}_{RCE}) is defined by:

$$\mathcal{L}_{\text{RCE}} = -\frac{1}{n} \sum_{i=1}^n (1 + \mathcal{W}(\mathbf{I}_i)) \log p_t(\mathbf{I}_i). \quad (6)$$

For example, when C performs better on images from the positive bias attribute group (*i.e.*, $\bar{P}_{b+}(\hat{y}) \geq \bar{P}_{b-}(\hat{y})$), we use $p(\hat{b} = 0 \mid \mathbf{I}_i)$ as the weight, which up-weights the images from the negative bias attribute group, where classifier C is worse-performed. At the same time, it down-weights the images from the positive bias attribute group where C is already better-performed. Adding one to the weight in Eq. (6) lets the loss function degenerate to standard cross-entropy loss when $\mathcal{W}(\mathbf{I}_i) = 0$. By minimizing the *Reweighted Cross-Entropy* loss, C is guided to meet Equal Opportunity.

3.3 Full Model

We summarize the proposed losses in Sec. 3.1 and Sec. 3.2 for the full model of DEBIASING ALTERNATE NETWORKS (DebiAN), which is shown in Fig. 2. When the task is to only discover (*i.e.*, not mitigate) the unknown biases of a given classifier, the classifier's parameters are fixed and we only train the *discoverer* D by minimizing \mathcal{L}_{EOV} (Eq. 2) and \mathcal{L}_{UA} (Eq. 4) on the classifier's training data. When the task is to mitigate the unknown biases, we jointly train two networks in an alternate fashion:

$$\min_D \mathcal{L}_{\text{EOV}} + \mathcal{L}_{\text{UA}}, \quad (7)$$

$$\min_C \mathcal{L}_{\text{RCE}}. \quad (8)$$

In Eq. 7, C 's parameters are fixed, and D is optimized to identify C 's unknown biases where C violates the Equal Opportunity. Through Eq. 8, C is optimized for mitigating the unknown biases discovered by D to satisfy the Equal Opportunity while D 's parameters are frozen. After the alternate training, C 's unknown biases identified by D are mitigated, leading to a fairer and more robust *classifier*. The pseudocode of the complete algorithm is in Appendix A.

4 Experiment

We conduct extensive experiments to verify the efficacy of DebiAN. First, we evaluate the results on our newly created Multi-Color MNIST dataset (Sec. 4.1) in a multi-bias setting. We further conduct experiments on real-world datasets in multiple image domains—face (Sec. 4.2) and other image domains (*e.g.*, scene, action recognition) (Sec. 4.3). More details (*e.g.*, evaluation metrics) are introduced in each subsection. The code and our newly created Multi-Color MNIST dataset are released at <https://github.com/zhihengli-UR/DebiAN>.

Comparison Methods We mainly compare with three unsupervised debiasing methods: 1) LfF [45] uses Generalized Cross-entropy (GCE) loss [65] to train a “biased model” for reweighing the classifier; 2) EIIL [14] identifies the bias groups by optimizing bias group assignment to maximize the IRMv1 [6] objective function. The identified bias groups will serve as pseudo bias labels for other supervised debiasing methods to mitigate the biases. Following [14], IRM [6] is used as the debiasing algorithm for EIIL. 3) PGI [2] follows EIIL to identify the biases by training a small multi-layer perceptron for bias label predictions. Concerning debiasing, PGI minimizes the KL-divergence of the classifier’s predictions across different bias groups. We use the officially released code of LfF, EIIL, and PGI in our experiment. Besides, we also compare with vanilla models, which do not have any debiasing techniques (*i.e.*, only using standard cross-entropy loss for training). On bFFHQ [32] and BAR [45] datasets, we also compare with BiaSwap [32], which follows LfF to identify unknown biases, and then uses style-transfer to generate more underrepresented images for training. Since its code has not been released, we cannot compare DebiAN with BiaSwap on other datasets. All results shown below are the mean results over three random seeds of runs, and we also report the standard deviation as the error bar.

4.1 Experiment on Multi-Color MNIST

Many previous works use synthetic datasets to benchmark bias mitigation performance. For example, Colored MNIST [7,31,41] adds color bias to the original MNIST [38] dataset, where each digit class is spuriously correlated with color (see Fig. 4 (a)). We compare DebiAN with other methods on the Colored MNIST dataset in Appendix E. However, we believe that the single-bias setting is an oversimplification of the real-world scenario where multiple biases may exist. For instance, Lang *et al.* [37] find that gender classifiers are biased with multiple independent bias attributes, including wearing lipsticks, eyebrow thickness, nose width, *etc.* The benchmarking results on such a single-bias synthetic dataset may not help us to design better debiasing algorithms for real-world usage.

To this end, we propose **Multi-Color MNIST** dataset to benchmark debiasing methods under the multi-bias setting. In the training set, each digit class is spuriously correlated with *two* bias attributes—**left color** and **right color** (Fig. 4 (b)). Following the terms used in LfF [45], we call samples that can be correctly predicted with the bias attribute as bias-aligned samples. Samples that cannot be correctly predicted with the bias attribute are called bias-conflicting



Fig. 4: Comparison between (a) previous Colored MNIST [7, 31, 41] with a single color bias and (b) our new Multi-Color MNIST dataset that contains *two* bias attributes—**left color** and **right color**

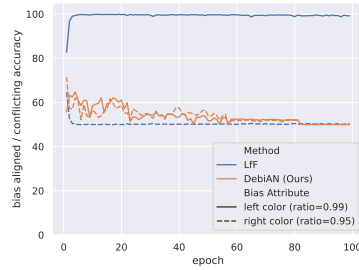


Fig. 5: Evaluating bias discovery w.r.t. **left color**, **right color** biases throughout the training epochs on Multi-Color MNIST. LfF only finds the more salient **left color** bias (ratio=0.99), whereas DebiAN’s *discoverer* finds both biases at the early training stage. Then accuracies gradually converge to 50% as debiasing is performed in the *classifier*, making the *discoverer* harder to find biases

samples. For example, if most digit “0” images are in red **left color** in the training set, we call them bias-aligned samples w.r.t. **left color** attribute, and we regard digit “0” images in a different **left color** (e.g., yellow) as bias-conflicting samples. Since the dataset contains two bias attributes, there exist images that are bias-aligned w.r.t. to **left color** and bias-conflicting w.r.t. **right color** simultaneously, or vice versa. Following [45], we use the *ratio* of the bias-aligned samples for each bias attribute to indicate how strong the spurious correlation is in the training set. The two ratios for two bias attributes can be different, which is more common in the real-world scenario. The images in the testing set also contain two background colors, but the testing set has a balanced distribution of bias-aligned and bias-conflicting samples w.r.t. each bias attribute.

Evaluation Metrics and Settings Following [45], we report the accuracy results in bias-aligned and bias-conflicting samples on the testing set. Since Multi-Color MNIST contains two bias attributes, we report the four accuracy results in the combination of (bias-aligned, bias-conflicting) \times (**left color**, **right color**), e.g., middle four rows in Tab. 1 for each method. We also report the unbiased accuracy, which averages the four results above. Here, we choose 0.99 as the ratio of bias-aligned samples w.r.t. **left color** and 0.95 as the ratio of bias-aligned samples w.r.t. **right color**. In this way, the **left color** is a *more salient* bias than the **right color**. We report the results of other ratio combinations in Appendix C.3. We strictly use the same set of hyperparameters used in (single) Colored MNIST in LfF. More details are in Appendix B.

Debiasing Results on Multi-Color MNIST The debiasing results are shown in Tab. 1. Except for LfF, all other methods achieve higher accuracy results on **left color** bias-aligned samples (1st and 2nd rows) than **right color** bias-aligned samples (1st and 3rd rows), indicating that most methods are more biased w.r.t. the more salient bias, i.e., **left color** (ratio=0.99) in the multi-bias setting.

Table 1: Debiasing results on Multi-Color MNIST dataset. The accuracy results in the four combinations of two bias attributes, (*i.e.*, **left color** and **right color**) and (bias-aligned and bias-conflicting) are reported. Unbiased accuracy averages the results over all four combinations. We bold top-2 results and underline lowest results

left color ratio = 0.99	right color ratio = 0.95	vanilla	LfF	EiIL	PGI	w/o \mathcal{L}_{UA} (Ours)	DebiAN (Ours)
bias-aligned	bias-aligned	100.0 \pm 0.0	99.6 \pm 0.5	100.0 \pm 0.0	98.6 \pm 2.3	100.0 \pm 0.0	100.0 \pm 0.0
bias-aligned	bias-conflicting	97.1 \pm 0.5	<u>4.7</u> \pm 0.5	97.2 \pm 1.5	82.6 \pm 19.6	97.2 \pm 0.5	95.6 \pm 0.8
bias-conflicting	bias-aligned	27.5 \pm 3.6	98.6 \pm 0.4	70.8 \pm 4.9	<u>26.6</u> \pm 5.5	71.6 \pm 0.7	76.5 \pm 0.7
bias-conflicting	bias-conflicting	5.2 \pm 0.4	<u>5.1</u> \pm 0.4	10.9 \pm 0.8	9.5 \pm 3.2	13.8 \pm 1.1	16.0 \pm 1.8
unbiased accuracy		57.4 \pm 0.7	<u>52.0</u> \pm 0.1	69.7 \pm 1.0	54.3 \pm 4.0	70.6 \pm 0.3	72.0 \pm 0.8

Unlike all other methods, LfF gives abnormal results—high accuracy results (*e.g.*, 99.6, 98.6) for the **right color** bias-aligned samples and low accuracy results (*e.g.*, 4.7, 5.1) for the **right color** bias-conflicting samples. Consequently, LfF achieves the worst unbiased accuracy (52.0). The results indicate that LfF only mitigates the more salient **left color** bias, rendering the classifier to learn the less salient **right color** bias (ratio=0.95). Compared with all other methods, DebiAN achieves better unbiased accuracy results (72.0). More importantly, DebiAN achieves much better debiasing result (16.0) in bias-conflicting samples w.r.t. both **left color** and **right color** attributes, where neither color can provide the shortcut for the correct digit class prediction, demonstrating better debiasing results of DebiAN for mitigating multiple biases simultaneously in the multi-bias setting, which is closer to the real-world scenarios.

Bias Discovery: LfF vs. DebiAN We further evaluate the bias discovery results throughout the entire training epochs, which helps us better understand LfF’s abnormal results and DebiAN’s advantages. We use LfF’s “biased model” and DebiAN’s *discoverer* to predict if a given image is bias-aligned or bias-conflicting w.r.t. a bias attribute (*i.e.*, binary classification, more details in Appendix D.1). We show the accuracy results of bias discovery w.r.t. each bias attribute at the end of each epoch in Fig. 5, which shows that LfF only discovers the more salient **left color** bias attribute (100% accuracy), but completely ignores the less salient **right color** bias (50% accuracy) throughout the entire training stage. It reveals the problem of LfF’s definition of the unknown bias—an attribute in the dataset that is easier, which only holds in the single-bias setting but does not generalize to the multi-bias setting. In contrast, DebiAN uses the principled definition to define the bias—classifier’s predictions that violate equal opportunity, enabling *discoverer* to find both biases accurately at the beginning (it achieves about 60% to 70% accuracy because debiasing is simultaneously performed before the end of the first epoch). At the same time, DebiAN’s alternate training scheme lets the classifier mitigate both biases, making the *discoverer* harder to predict the biases, *e.g.*, accuracies of both bias attributes gradually converge to 50%. More discussions are in Appendix D.4.

Ablation Study on UA penalty We conduct an ablation study to show the effectiveness of Unbalanced Assignment (UA) penalty (Sec. 3.1). Tab. 1 shows

Table 2: Results of mitigating the gender bias of **Blond Hair** classifier on CelebA [43]

	vanilla	LfF	EHL	PGI	DebiAN (Ours)
Avg Group Acc.	79.8 \pm 0.3	80.9 \pm 1.4	82.0 \pm 1.1	81.6 \pm 0.3	84.0\pm1.4
Worst Group Acc.	37.9 \pm 1.1	43.3 \pm 3.0	46.1 \pm 4.9	40.9 \pm 6.4	52.9\pm4.7

Table 3: Accuracy results on bias-conflicting samples on bFFHQ [32]

vanilla	LfF	PGI	EHL	BiaSwap	DebiAN
51.03	55.61	55.2 \pm 5.3	59.2 \pm 1.9	58.87	62.8\pm0.6

that \mathcal{L}_{UA} improves the debiasing results (see w/o \mathcal{L}_{UA}). Besides, we also conduct ablation studies on different batch sizes, which are included in Appendix C.2.

4.2 Experiments on Face Image Dataset

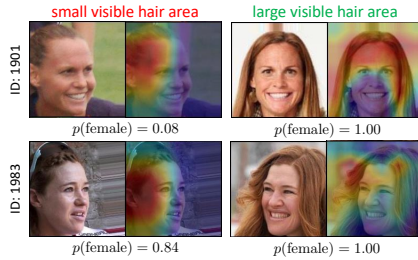
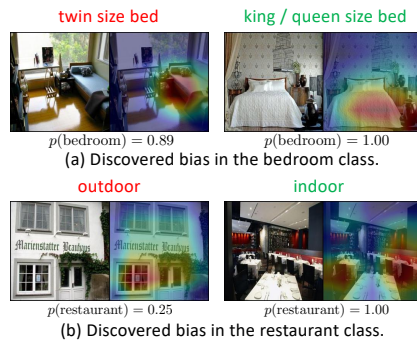
Gender Bias Mitigation In the face image domain, we conduct experiments to evaluate gender bias mitigation results on CelebA [43] dataset, which contains 200K celebrity faces annotated with 40 binary attributes. The dataset has spurious correlations between gender and **Blond Hair**, leading to gender biases when performing hair color classification. We follow most of the settings used in LfF, such as using ResNet-18 [24] as the backbone, using Adam [33] optimizer, *etc.* The only difference is that LfF reports the results on the validation set of CelebA, whereas we use the validation set to select the epoch with the best validation set accuracy (bias labels in the validation set are not used) to report the results on the testing set. All methods (including LfF) are benchmarked under the same setting. We report results in two evaluation metrics: 1) Average Group Accuracy (Avg. Group Acc.), which calculates the unweighted average of accuracies in four groups between target attribute and bias attribute, *i.e.*, (male, female) \times (blond, not blond); 2) Worst Group Accuracy (Worst Group Acc.) [47], which takes the lowest accuracy in the four groups. As shown in Tab. 2, DebiAN achieves better Average and Worst Group accuracy results, which shows that DebiAN can better mitigate gender bias without labels. We also conduct experiments on bFFHQ [32] where the training data contains the spurious correlation between age and gender. We compare DebiAN with other methods of gender bias mitigation. We strictly follow the setting in [32]. We report the age accuracy results on the bias-conflicting samples in the testing set in Tab. 3. The results of vanilla, LfF, and BiaSwap are from [32] and [32] does not provide the error bars. DebiAN achieves the best unsupervised results for mitigating gender bias.

Mitigating Multiple Biases in Gender Classifier The results on Multi-Color MNIST dataset suggest that DebiAN better mitigates multiple biases in the classifier. In the face image domain, a recent study [37] shows that gender classifier is biased by multiple attributes, such as **Heavy Makeup** and **Wearing Lipstick**. Hence, we train gender classifiers on CelebA dataset and evaluate Average Group Accuracy and Worst Group Accuracy w.r.t. these two bias attributes. As shown in Tab. 4, DebiAN achieves better debiasing results w.r.t. both bias attributes, proving that the *discoverer* can find multiple biases in the classifier C during the alternate training, enabling *classifier* to mitigate multiple biases simultaneously.

Identified Unknown Bias in Gender Classifier Gender classifier can have more biases beyond **Wearing Lipstick** and **Heavy Makeup**. For example, Bal-

Table 4: Results of mitigating multiple biases (*i.e.*, **Wearing Lipstick** and **Heavy Makeup**) in gender classifier on CelebA dataset

bias attribute	metric	vanilla	LfF	PGI	EIIL	DebiAN (Ours)
Wearing Lipstick	Avg. Group Acc.	86.6 \pm 0.4	87.0 \pm 0.9	86.9 \pm 3.1	86.3 \pm 1.0	88.5\pm1.1
	Worst Group Acc.	53.9 \pm 1.2	55.3 \pm 3.6	56.0 \pm 11.7	52.4 \pm 3.2	61.7\pm4.2
Heavy Makeup	Avg. Group Acc.	85.1 \pm 0.0	85.5 \pm 0.6	85.4 \pm 3.4	84.0 \pm 1.2	87.8\pm1.3
	Worst Group Acc.	45.4 \pm 0.0	46.9 \pm 2.6	46.9 \pm 13.1	40.9 \pm 4.5	56.0\pm5.2

**Fig. 6:** Discovered bias of gender classifier: **visible hair area** based on *discoverer*'s saliency map. $p(\text{female})$ is vanilla classifier's predicted probability of the face is female. In the two groups predicted by D , the visible hair areas are different, where the classifier has different confidences on gender for the same identity**Fig. 7:** Discovered biases in Places [68] dataset. We apply CAM on *discoverer* to generate saliency map. The value $p(\text{bedroom})$ ($p(\text{restaurant})$) is vanilla classifier's predicted probability of the scene image is bedroom (restaurant)

akrishnan *et al.* [8] leverages StyleGAN2 [30] to generate high-quality synthesized images and identify the **hair length** bias of the gender classifier, *e.g.*, longer hair length makes the classifier predict the face as female. Related to their finding, the *discoverer* D in DebiAN identifies an interesting unknown bias: **visible hair area**. We use D to predict the bias attribute group assignment on images in CelebA. To better interpret the bias attribute, we further use the identity labels in CelebA to cluster images with the same identity. Fig. 6 shows that D assigns images of the same identity into two distinct groups based on the visible hair area, which is verified by D 's CAM [67] saliency maps. Strictly speaking, all females in Fig. 6 have long hair. However, due to the hairstyle, pose, or occlusion, visible hair areas differ between the two groups. As a result, the gender classifier has lower predicted probabilities on the female images with smaller visible hair areas. More visualizations are shown in Appendix G.1.

4.3 Experiments on Other Image Domains

Our method is not limited to synthetic and face image domains. Here we conduct experiments on action recognition and scene classification tasks.

Table 5: Results on Biased Action Recognition (BAR) [45] dataset

vanilla	LfF	PGI	EIIL	BiaSwap	DebiAN
51.85 \pm 5.92	62.98 \pm 2.76	65.19 \pm 1.32	65.44 \pm 1.17	52.44	69.88\pm2.92

Table 6: Scene classification accuracy results on the *unseen* LSUN [63] dataset

vanilla	LfF	PGI	EIIL	DebiAN (Ours)
79.3 \pm 0.3	71.1 \pm 1.0	74.1 \pm 1.9	79.4 \pm 0.2	80.0\pm0.4

Mitigating Place Bias in Action Recognition We conduct experiments on Biased Action Recognition (BAR) dataset [45], an image dataset with the spurious correlation between action and place in the training set. The testing set only contains bias-conflicting samples. Hence, higher accuracy results on the testing set indicate better debiasing results. The accuracy results in Tab. 5 show that DebiAN achieves better debiasing results than other methods.

Improving Cross-dataset Generalization on Scene Classification We conduct experiments on the more challenging scene classification task, where datasets are more complex and may contain multiple unknown biases. The biases in this task are underexplored by previous works partly due to the lack of attribute labels. Due to the absence of attribute labels, we use cross-dataset generalization [53] to evaluate the debiasing results. Concretely, models are trained on Places [68] with ten classes overlapped with LSUN [63] (*e.g.*, bedroom, classroom, *etc.*), and evaluated on the *unseen* LSUN dataset. The results are shown in Tab. 6. DebiAN achieves the best result on the unseen LSUN dataset, showing that DebiAN unlearns the dataset biases [53] in Places to improve the robustness against distributional shifts between different datasets.

Identified Unknown Biases in Scene Classifier DebiAN discovers Places dataset’s unknown biases that humans may not preconceive. In Fig. 7, the *discoverer* separates bedroom and restaurant images based on **size** of beds and **indoor/outdoor**. The vanilla classifier performs worse on bedroom images with twin-size beds and outdoor restaurant images (see more in Appendix G.2).

5 Conclusion

We propose DEBIASING ALTERNATE NETWORKS to discover and mitigate the unknown biases. DebiAN identifies unknown biases that humans may not preconceive and achieves better unsupervised debiasing results. Our Multi-Color MNIST dataset surfaces previous methods’ problems and demonstrates DebiAN’s advantages in the multi-bias setting. Admittedly, our work has some limitations, *e.g.*, DebiAN focuses on binary or continuously valued bias attributes, not multi-class ones. We hope our work can facilitate research on bias discovery and mitigation.

Acknowledgment This work has been partially supported by the National Science Foundation (NSF) under Grant 1764415, 1909912, and 1934962 and by the Center of Excellence in Data Science, an Empire State Development-designated Center of Excellence. The article solely reflects the opinions and conclusions of its authors but not the funding agents.

References

1. Agrawal, R., Srikant, R.: Fast Algorithms for Mining Association Rules in Large Databases. In: International Conference on Very Large Data Bases (1994)
2. Ahmed, F., Bengio, Y., van Seijen, H., Courville, A.: Systematic generalisation with group invariant predictions. In: International Conference on Learning Representations (2021)
3. Albiero, V., K. S., K., Vangara, K., Zhang, K., King, M.C., Bowyer, K.W.: Analysis of Gender Inequality In Face Recognition Accuracy. In: The IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW) (2020)
4. Alvi, M., Zisserman, A., Nellaaker, C.: Turning a Blind Eye: Explicit Removal of Biases and Variation from Deep Neural Network Embeddings. In: The European Conference on Computer Vision Workshop (ECCVW) (2018)
5. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: Visual Question Answering. In: The IEEE International Conference on Computer Vision (ICCV) (2015)
6. Arjovsky, M., Bottou, L., Gulrajani, I., Lopez-Paz, D.: Invariant Risk Minimization. arXiv:1907.02893 [cs, stat] (2020)
7. Bahng, H., Chun, S., Yun, S., Choo, J., Oh, S.J.: Learning De-biased Representations with Biased Representations. In: International Conference on Machine Learning (2020)
8. Balakrishnan, G., Xiong, Y., Xia, W., Perona, P.: Towards causal benchmarking of bias in face analysis algorithms. In: The European Conference on Computer Vision (ECCV) (2020)
9. Buolamwini, J., Gebru, T.: Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In: ACM Conference on Fairness, Accountability, and Transparency (2018)
10. Cadene, R., Dancette, C., Ben younes, H., Cord, M., Parikh, D.: RUBi: Reducing Unimodal Biases for Visual Question Answering. In: Advances in Neural Information Processing Systems (2019)
11. Choi, J., Gao, C., Messou, J.C.E., Huang, J.B.: Why Can't I Dance in the Mall? Learning to Mitigate Scene Bias in Action Recognition. In: Advances in Neural Information Processing Systems (2019)
12. Clark, C., Yatskar, M., Zettlemoyer, L.: Don't Take the Easy Way Out: Ensemble Based Methods for Avoiding Known Dataset Biases. In: Empirical Methods in Natural Language Processing (2019)
13. Corbett-Davies, S., Pierson, E., Feller, A., Goel, S., Huq, A.: Algorithmic Decision Making and the Cost of Fairness. In: Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (2017)
14. Creager, E., Jacobsen, J.H., Zemel, R.: Environment Inference for Invariant Learning. In: International Conference on Machine Learning (2021)
15. Creager, E., Madras, D., Jacobsen, J.H., Weis, M., Swersky, K., Pitassi, T., Zemel, R.: Flexibly Fair Representation Learning by Disentanglement. In: International Conference on Machine Learning (2019)
16. Dhar, P., Gleason, J., Roy, A., Castillo, C.D., Chellappa, R.: PASS: Protected Attribute Suppression System for Mitigating Bias in Face Recognition. In: The IEEE International Conference on Computer Vision (ICCV) (2021)
17. Dwork, C., Hardt, M., Pitassi, T., Reingold, O., Zemel, R.: Fairness through awareness. In: Proceedings of the 3rd Innovations in Theoretical Computer Science Conference (2012)

18. Geirhos, R., Jacobsen, J.H., Michaelis, C., Zemel, R., Brendel, W., Bethge, M., Wichmann, F.A.: Shortcut learning in deep neural networks. *Nature Machine Intelligence* (2020)
19. Gong, S., Liu, X., Jain, A.K.: Jointly De-biasing Face Recognition and Demographic Attribute Estimation. In: *The European Conference on Computer Vision (ECCV)* (2020)
20. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: *Advances in Neural Information Processing Systems* (2014)
21. Grgic-Hlaca, N., Zafar, M.B., Gummadi, K.P., Weller, A.: The case for process fairness in learning: Feature selection for fair decision making. In: *NIPS Symposium on Machine Learning and the Law* (2016)
22. Hardt, M., Price, E., Srebro, N.: Equality of Opportunity in Supervised Learning. In: *Advances in Neural Information Processing Systems* (2016)
23. Hazirbas, C., Bitton, J., Dolhansky, B., Pan, J., Gordo, A., Ferrer, C.C.: Towards Measuring Fairness in AI: The Casual Conversations Dataset. *arXiv:2104.02821 [cs]* (2021)
24. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
25. Hendricks, L.A., Burns, K., Saenko, K., Darrell, T., Rohrbach, A.: Women also Snowboard: Overcoming Bias in Captioning Models. In: *The European Conference on Computer Vision (ECCV)* (2018)
26. Jia, S., Meng, T., Zhao, J., Chang, K.W.: Mitigating Gender Bias Amplification in Distribution by Posterior Regularization. In: *Annual Meeting of the Association for Computational Linguistics* (2020)
27. Joo, J., Kärkkäinen, K.: Gender Slopes: Counterfactual Fairness for Computer Vision Models by Attribute Manipulation. In: *International Workshop on Fairness, Accountability, Transparency and Ethics in Multimedia* (2020)
28. Kamiran, F., Calders, T.: Data preprocessing techniques for classification without discrimination. *Knowledge and Information Systems* (2012)
29. Karras, T., Laine, S., Aila, T.: A Style-Based Generator Architecture for Generative Adversarial Networks. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
30. Karras, T., Laine, S., Aittala, M., Hellsten, J., Lehtinen, J., Aila, T.: Analyzing and Improving the Image Quality of StyleGAN. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2020)
31. Kim, B., Kim, H., Kim, K., Kim, S., Kim, J.: Learning Not to Learn: Training Deep Neural Networks With Biased Data. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
32. Kim, E., Lee, J., Choo, J.: BiaSwap: Removing dataset bias with bias-tailored swapping augmentation. In: *The IEEE International Conference on Computer Vision (ICCV)* (2021)
33. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization. In: *International Conference on Learning Representations* (2015)
34. Krishnakumar, A., Prabhu, V., Sudhakar, S., Hoffman, J.: UDIS: Unsupervised Discovery of Bias in Deep Visual Recognition Models. In: *British Machine Vision Conference, BMVC* (2021)
35. Kusner, M.J., Loftus, J., Russell, C., Silva, R.: Counterfactual Fairness. In: *Advances in Neural Information Processing Systems* (2017)

36. Lahoti, P., Beutel, A., Chen, J., Lee, K., Prost, F., Thain, N., Wang, X., Chi, E.: Fairness without Demographics through Adversarially Reweighted Learning. In: *Advances in Neural Information Processing Systems* (2020)
37. Lang, O., Gandselman, Y., Yarom, M., Wald, Y., Elidan, G., Hassidim, A., Freeman, W.T., Isola, P., Globerson, A., Irani, M., Mosseri, I.: Explaining in Style: Training a GAN to explain a classifier in StyleSpace. In: *The IEEE International Conference on Computer Vision (ICCV)* (2021)
38. Lecun, Y., Bottou, L., Bengio, Y., Haffner, P.: Gradient-based learning applied to document recognition. *Proceedings of the IEEE* (1998)
39. Li, W., Zhang, P., Zhang, L., Huang, Q., He, X., Lyu, S., Gao, J.: Object-Driven Text-To-Image Synthesis via Adversarial Training. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
40. Li, Y., Li, Y., Vasconcelos, N.: RESOUND: Towards Action Recognition without Representation Bias. In: *The European Conference on Computer Vision (ECCV)* (2018)
41. Li, Y., Vasconcelos, N.: REPAIR: Removing Representation Bias by Dataset Resampling. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
42. Li, Z., Xu, C.: Discover the Unknown Biased Attribute of an Image Classifier. In: *The IEEE International Conference on Computer Vision (ICCV)* (2021)
43. Liu, Z., Luo, P., Wang, X., Tang, X.: Deep Learning Face Attributes in the Wild. In: *The IEEE International Conference on Computer Vision (ICCV)* (2015)
44. Manjunatha, V., Saini, N., Davis, L.S.: Explicit Bias Discovery in Visual Question Answering Models. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2019)
45. Nam, J., Cha, H., Ahn, S., Lee, J., Shin, J.: Learning from Failure: Training Debaised Classifier from Biased Classifier. In: *Advances in Neural Information Processing Systems* (2020)
46. Pleiss, G., Raghavan, M., Wu, F., Kleinberg, J., Weinberger, K.Q.: On Fairness and Calibration. In: *Advances in Neural Information Processing Systems* (2017)
47. Sagawa*, S., Koh*, P.W., Hashimoto, T.B., Liang, P.: Distributionally Robust Neural Networks for Group Shifts: On the Importance of Regularization for Worst-Case Generalization. In: *International Conference on Learning Representations* (2020)
48. Sarhan, M.H., Navab, N., Albarqouni, S.: Fairness by Learning Orthogonal Disentangled Representations. In: *The European Conference on Computer Vision (ECCV)* (2020)
49. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: *The IEEE International Conference on Computer Vision (ICCV)* (2017)
50. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization. *International Journal of Computer Vision* (2019)
51. Singh, K.K., Mahajan, D., Grauman, K., Lee, Y.J., Feiszli, M., Ghadiyaram, D.: Don't Judge an Object by Its Context: Learning to Overcome Contextual Bias. In: *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2020)
52. Sohoni, N.S., Dunnmon, J.A., Angus, G., Gu, A., Ré, C.: No Subclass Left Behind: Fine-Grained Robustness in Coarse-Grained Classification Problems. In: *Advances in Neural Information Processing Systems* (2020)

53. Torralba, A., Efros, A.A.: Unbiased look at dataset bias. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2011)
54. Tzeng, E., Hoffman, J., Darrell, T., Saenko, K.: Simultaneous Deep Transfer Across Domains and Tasks. In: The IEEE International Conference on Computer Vision (ICCV) (2015)
55. Verma, S., Rubin, J.: Fairness Definitions Explained. In: 2018 IEEE/ACM International Workshop on Software Fairness (FairWare) (2018)
56. Wang, A., Narayanan, A., Russakovsky, O.: REVISE: A Tool for Measuring and Mitigating Bias in Image Datasets. In: The European Conference on Computer Vision (ECCV) (2020)
57. Wang, H., He, Z., Lipton, Z.C., Xing, E.P.: Learning Robust Representations by Projecting Superficial Statistics Out. In: International Conference on Learning Representations (2019)
58. Wang, J., Liu, Y., Wang, X.E.: Are Gender-Neutral Queries Really Gender-Neutral? Mitigating Gender Bias in Image Search. In: Empirical Methods in Natural Language Processing (2021)
59. Wang, T., Yue, Z., Huang, J., Sun, Q., Zhang, H.: Self-Supervised Learning Disentangled Group Representation as Feature. In: Advances in Neural Information Processing Systems (2021)
60. Wang, X., Ang, M.H., Lee, G.H.: Cascaded Refinement Network for Point Cloud Completion. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
61. Wang, Z., Liu, X., Li, H., Sheng, L., Yan, J., Wang, X., Shao, J.: CAMP: Cross-Modal Adaptive Message Passing for Text-Image Retrieval. In: The IEEE International Conference on Computer Vision (ICCV) (2019)
62. Wang, Z., Qinami, K., Karakozis, I.C., Genova, K., Nair, P., Hata, K., Russakovsky, O.: Towards Fairness in Visual Recognition: Effective Strategies for Bias Mitigation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2020)
63. Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., Xiao, J.: LSUN: Construction of a Large-scale Image Dataset using Deep Learning with Humans in the Loop. arXiv:1506.03365 [cs] (2016)
64. Zhang, B.H., Lemoine, B., Mitchell, M.: Mitigating Unwanted Biases with Adversarial Learning. In: AAAI/ACM Conference on AI, Ethics, and Society (2018)
65. Zhang, Z., Sabuncu, M.: Generalized Cross Entropy Loss for Training Deep Neural Networks with Noisy Labels. In: Advances in Neural Information Processing Systems (2018)
66. Zhao, J., Wang, T., Yatskar, M., Ordonez, V., Chang, K.W.: Men Also Like Shopping: Reducing Gender Bias Amplification using Corpus-level Constraints. In: Empirical Methods in Natural Language Processing (2017)
67. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning Deep Features for Discriminative Localization. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
68. Zhou, B., Lapedriza, A., Khosla, A., Oliva, A., Torralba, A.: Places: A 10 million image database for scene recognition. IEEE Transactions on Pattern Analysis and Machine Intelligence (2018)