

# Bias-Eliminating Augmentation Learning for Debaised Federated Learning

Yuan-Yi Xu<sup>1</sup>    Ci-Siang Lin<sup>1</sup>    Yu-Chiang Frank Wang<sup>1,2</sup>  
<sup>1</sup>National Taiwan University, Taiwan    <sup>2</sup>NVIDIA, Taiwan  
 {r10942094, d08942011, ycwang}@ntu.edu.tw

## Abstract

Learning models trained on biased datasets tend to observe correlations between categorical and undesirable features, which result in degraded performances. Most existing debaised learning models are designed for centralized machine learning, which cannot be directly applied to distributed settings like federated learning (FL), which collects data at distinct clients with privacy preserved. To tackle the challenging task of debaised federated learning, we present a novel FL framework of **Bias-Eliminating Augmentation Learning (FedBEAL)**, which learns to deploy **Bias-Eliminating Augmenters (BEA)** for producing client-specific bias-conflicting samples at each client. Since the bias types or attributes are not known in advance, a unique learning strategy is presented to jointly train BEA with the proposed FL framework. Extensive image classification experiments on datasets with various bias types confirm the effectiveness and applicability of our FedBEAL, which performs favorably against state-of-the-art debiasing and FL methods for debaised FL.

## 1. Introduction

Deep neural networks have shown promising progress across different domains such as computer vision [14] and natural language processing [8]. Their successes are typically based on the collection of and training on data that properly describe the inherent distribution of the data of interest. However, in real-world scenarios, biased data [24] are often observed during data collection. Biased datasets [10, 22, 42] contain features that are highly correlated to class labels in the training dataset but not sufficiently describing the inherent semantic meaning. Training on such biased data thus result in degraded model generalization capability. Take Fig. 1 for example; when addressing the cat-dog classification task, training images collected by users might contain only orange cats and black dogs. Their color attributes are strongly correlated with the image labels during training, but such attributes are not necessarily relevant to the classification task during inference. As

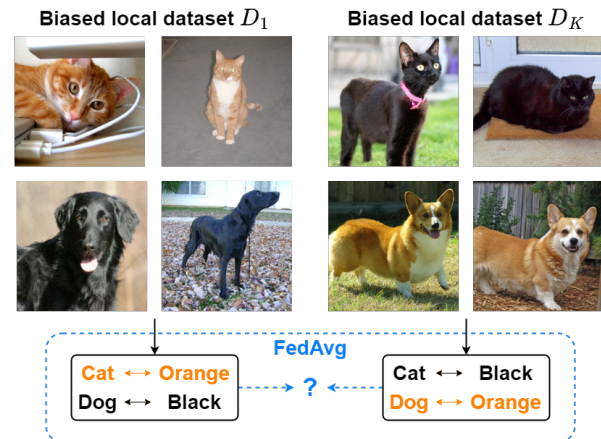


Figure 1. **Example of local data bias in FL.** When deploying FL to train a cat-dog classifier with image datasets collected by multiple pet owners, most of the local images are obtained with their pets with specific colors. Therefore, the models trained with each local dataset are likely to establish decision rules on biased attributes (e.g., fur color), which prevents the aggregated model from learning proper representation for classification.

pointed out in [10, 42], deep neural networks trained with such biased data are more likely to make decisions based on *bias* attributes instead of semantic attributes. As a result, during inference, performances of the learned models would dramatically drop when observing *bias-conflicting* samples (i.e., data containing semantic and bias attributes that are rarely correlated in the training set).

To tackle the data bias problem, several works have been proposed to remove or alleviate data bias when training deep learning models [6, 11, 18, 24, 27, 32, 36, 40]. For example, Nam *et al.* [36] train an intentionally biased auxiliary model while enforcing the main model to go against the prejudice of the biased network. Lee *et al.* [27] utilize the aforementioned biased model to synthesize diverse bias-conflicting hidden features for learning debaised representations. Nevertheless, the above techniques are designed for centralized datasets. When performing distributed training of learning models, such methods might fail to generalize.

For distributed learning, federated learning (FL) [35]

particularly considers data collection and training conducted at each client, with data privacy needing to be preserved. When considering privately distributed datasets, real-world FL applications are more likely to suffer data heterogeneity issues [20, 28, 51], *i.e.*, data collected by different clients are not independent and identically distributed (IID). Recently, several works [19, 21, 29–31, 34, 47] propose to alleviate performance degradation caused by data heterogeneity. However, existing methods typically consider data heterogeneity in terms of label distribution skew [21, 29, 30, 34, 47] or domain discrepancy [19, 31] among clients. These FL methods are not designed to tackle potential data bias across different clients, leaving the debiased FL a challenging task to tackle.

To mitigate the local bias in **Federated learning**, we propose a novel FL scheme of **Bias-Eliminating Augmentation Learning (FedBEAL)**. In FedBEAL, we learn a Bias-Eliminating Augmenter (BEA) for each client, with the goal of producing bias-conflicting samples. To identify and introduce the desirable semantic and bias attributes to the augmented samples, our FedBEAL uniquely adopts the global server model and each client model trained across iterations without prior knowledge of bias type or annotation. With the introduced augmenter and the produced bias-conflicting samples, debiased local updates can be performed at each client, followed by simple aggregation of such models for deriving the server model.

We now summarize the contributions of this work below:

- To the best of our knowledge, We are among the first to tackle the problem of debiased federated learning, in which local yet distinct biases exist at the client level.
- We present FedBEAL for debiased FL, which introduces Bias-Eliminating Augmenters (BEA) at each client with the goal of generating bias-conflicting samples to eliminate local data biases.
- Learning of BEA can be realized by utilizing the global server and local client models trained across iterations, which allows us to identify and embed desirable semantic and bias features for augmentation purposes.

## 2. Related Work

### 2.1. Debiasing in Centralized Machine Learning

With the presence of biased datasets, neural networks are prone to relying on simpler features (*e.g.*, color information) and remaining invariant to other predictive complex features [10, 42] (*e.g.*, semantic information), which limit the performances of the learned models. Several works [6, 11, 43, 46] propose debiasing techniques to improve the robustness of the model trained on such biased datasets. However, they either assume the bias type to be

known (*e.g.*, texture bias) in advanced [11] or require auxiliary annotations of the bias attributes (*e.g.*, color information) for each sample [6, 43, 46]), which might not be practically available. To alleviate this concern, some research works [18, 27, 36, 44] focus on mitigating dataset biases without presuming bias categories or involving additional annotations. For instance, Nam *et al.* [36] train a biased model by repeatedly amplifying its prejudice. Based on the assumption that biased models fail on bias-conflicting samples, they further upweight the bias-conflicting samples so that a debiased model can be trained accordingly. Lee *et al.* [27] follow the above approach to debias the main model by disentangling the semantic and bias features. On the other hand, Hong *et al.* [18] apply contrastive learning [13, 23] to encourage intra-bias feature dissimilarities. Although the above methods have shown promising performances, they are mainly applicable to centralized learning schemes. For distributed learning like federated learning, these methods cannot be directly applied.

### 2.2. Federated Learning with Data Heterogeneity

**Label distribution skew.** Under the heterogeneous label distribution, existing methods [21, 29, 30, 33, 34, 47, 53] focus on correcting client drift using global information. For example, FedProx [30] adds a regularization term to preserve consistency between local updates and the global model. SCAFFOLD [21] mitigates gradient dissimilarity using control variates. MOON [29] addresses non-IID problems by applying contrastive learning at the model level.

**Distribution shift across clients.** As for feature distribution drift (also known as domain shift), previous FL works [19, 31] are designed to bridge the domain gap between different clients. For instance, FedBN [31] choose to fix the parameters for local Batch Normalization and do not synchronize them with the global model. As for FCCL [19], it views domain shift as a catastrophic forgetting problem and approaches it by using knowledge distillation techniques.

**Debiased federated learning.** Recently, a number of FL works [1–4, 9] are proposed to eliminate *local biases* from the training data. In [3, 4], such biases are referred to as label distribution skew. For example, [4] uses the term *local learning bias* to describe decision boundaries discrepancy among networks trained on heterogeneous data. As for [1, 9], additional efforts are made to take care of underprivileged or sensitive data subsets (*e.g.*, racial, gender groups). For example, Ezzeldin *et al.* [9] propose a fairness-aware FL framework for preventing the trained model from being biased toward an underlying demographic group, aiming to produce a fair model across clients while maintaining high utility. It can be seen that we are among the first

to address the learning task of *debiased federated learning*, in which undesirable correlations of bias attributes and class labels are observed at each client.

### 3. Method

#### 3.1. Problem Definition and Method Overview

**Problem definition** For the sake of completeness, we first define the problem setting and notations used in this paper. We assume that training image data are privately distributed in  $K$  clients  $D = \{D_1, D_2, \dots, D_K\}$ , each containing a set of image-label pairs  $D_k = \{(x, y) \mid P_k(X = x, Y = y)\}$ . To formulate local data biases, we follow Hong *et al.* [18] and assume that images  $X$  can be decomposed into semantic attributes  $A_{sem}$  and hidden bias attributes  $A_{bias}$ . Note that  $A_{sem}$  is expected to describe categorical information, while  $A_{bias}$  contains undesirable features highly correlated with  $Y$ . As depicted in Figure 1, we assume each client with disparate bias-label correlations (*i.e.*,  $\forall_{k \neq k'} P_k(Y|A_{bias}) \neq P_{k'}(Y|A_{bias})$ ). On the other hand, since this work focuses on mitigating local client bias instead of the bias of the global dataset  $D$ , we assume the union of all local training datasets shares the same bias distribution with the test dataset  $D_{test}$  (*i.e.*,  $P(Y|A_{bias}) = P_{test}(Y|A_{bias})$ ). With a total of  $T$  communication rounds, the goal of debiased FL is to derive a model  $f$  that satisfies

$$\arg \min_f \sum_{k=1}^K \frac{|D_k|}{|D|} \mathcal{L}_k(f), \quad (1)$$

where  $\mathcal{L}_k(f) = \mathbb{E}_{(x,y) \sim D_k} [\ell(f(x), y)]$  represents the empirical loss of client  $k$ .

**Method overview** Based on FedAvg [35], our proposed **Bias-Eliminating Augmentation Learning (FedBEAL)** trains a network robust to data bias observed at each client. Similar to standard FL, training of FedBEAL requires alternative optimization between the two stages. More specifically, *debiased local update* is performed at the client side, and *global aggregation* is conducted at the server side. To address local bias problems, we uniquely propose to learn a **Bias-Eliminating Augmenter (BEA)**  $g_k$  for each client  $k$ . As depicted in Figure 2, BEA is deployed to generate bias-conflicting samples and allows updates of each  $f_k$ . As for the global aggregation stage, each  $f_k$  will be uploaded to the server for producing a debiased global model  $f$ . We now detail our proposed learning scheme below.

#### 3.2. Bias-Eliminating Augmenter

To eliminate the local bias in FL, we propose to deploy Bias-Eliminating Augmenters at each client. Since the bias information is unknown, how to design BEA for creating bias-conflicting samples within each local client would be

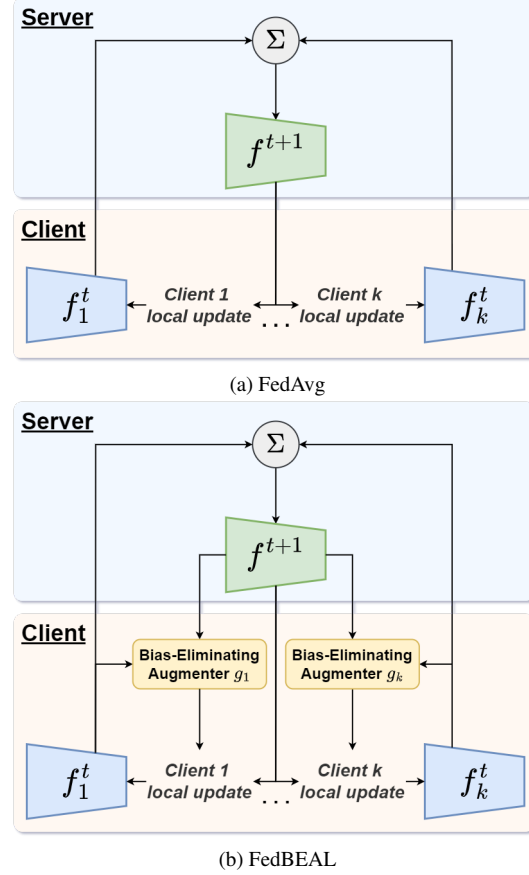


Figure 2. **Comparisons between (a) FedAvg and (b) FedBEAL.** Our FedBEAL learns Bias-Eliminating Augmenters (BEA) to produce bias-conflicting samples at each client, allowing the learned model to produce improved debiased representations.

challenging. With local image data and their class labels observed, we now explain how our BEA can be learned in an FL scheme.

##### 3.2.1 Design and architecture

As depicted in Figure 3, for each client  $k$ , we randomly sample two samples  $x^i$  and  $x^j$  with *distinct* labels from the local dataset  $D_k$ . Inspired by recent mixed sample data augmentation (MSDA) techniques [7, 12, 16, 17, 38, 49, 52], we produce the mixed bias-conflicting sample  $\tilde{x}$  by utilizing U-Net as the backbone, with a modulator  $M \in [0, 1]^{H \times W \times 3}$  deployed. With the concatenation of  $x^i$  and  $x^j$  as the input to BEA, the output  $\tilde{x}$  can be expressed as:

$$\tilde{x} = M \odot x^i + (1 - M) \odot x^j, \quad (2)$$

where  $\odot$  indicates the element-wise multiplication, and we have  $\tilde{y} = y^i$  for the manipulated output (*i.e.*, the class label of  $\tilde{x}$  is identical to that of  $x^i$ ).

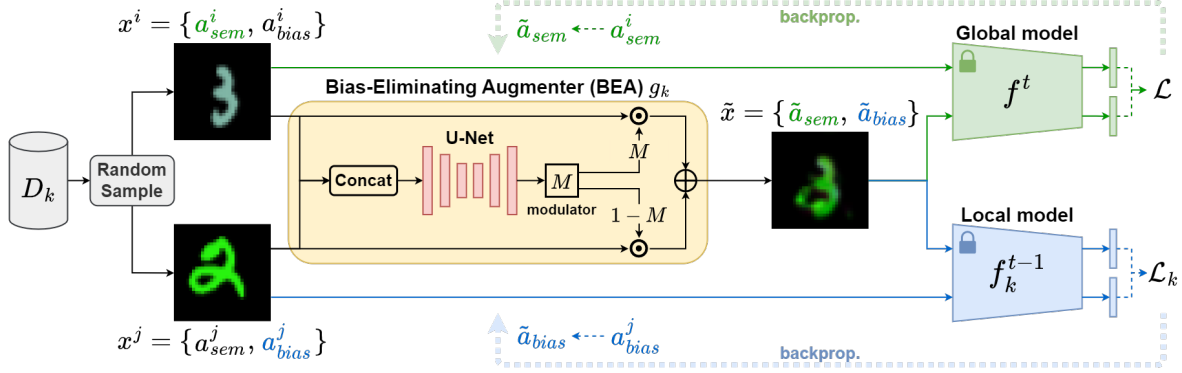


Figure 3. **Design and learning of Bias-Eliminating Augmenter.** Given two randomly selected images  $x^i$  and  $x^j$  at client  $k$ , the Bias-Eliminating Augmenter (BEA) learns to produce a bias-conflicting sample  $\tilde{x}$ . That is, the semantic attribute  $a_{sem}$  of  $\tilde{x}$  is expected to be close to that of  $x^i$ , while the bias attribute  $a_{bias}$  of  $\tilde{x}$  would be extracted from  $x^j$ . Note that  $f^t$  and  $f_k^{t-1}$  denote the server and client models learned at  $t$ -th and  $(t-1)$ -th iterations, respectively.

For  $\tilde{x}$  being a bias-conflicting example, it would be desirable for  $\tilde{x}$  to share the same semantic attribute with  $x^i$  (i.e.,  $\tilde{a}_{sem}$  of  $\tilde{x}$  to be closed to  $a_{sem}^i$  of  $x^i$ ), while sharing the same bias attribute with  $x^j$  (i.e.,  $\tilde{a}_{bias}$  of  $\tilde{x}$  to be closed to  $a_{bias}^j$  of  $x^j$ ). Once such bias-conflicting samples are obtained, one can train the associated client model and update the global model accordingly, which is expected to produce debiased representations.

### 3.2.2 Learning of BEA

Without prior knowledge of bias types, providing guidance to train the BEA would not be straightforward. In order to have BEA identify desirable intrinsic semantic and inherent bias attributes for manipulating bias-conflicting samples, we propose a unique learning scheme utilizing the global server model  $f^t$  and local client model  $f_k^{t-1}$ .

**Extracting semantic attributes via unbiased global prediction.** For a bias-conflicting sample  $\tilde{x}$ , its semantic attribute  $\tilde{a}_{sem}$  is expected to be similar to  $a_{sem}^i$  of  $x^i$ . In FL, since the global server model  $f^t$  is derived by global aggregation,  $f^t$  can be considered relatively unbiased when compared to the local model  $f_k^{t-1}$  produced at the previous iteration. Thus, it would be desirable for  $\tilde{a}_{sem}$  and  $a_{sem}^i$  to exhibit large similarity, which can be derived from the difference between the predictions of  $\tilde{x}$  and  $x^i$  derived from the global model  $f^t$ . To be precise, the loss function for encouraging such semantic attribute consistency is defined as:

$$\mathcal{L} = d_{KL}(f^t(\tilde{x}), f^t(x^i)), \quad (3)$$

where  $d_{KL}$  calculates the KL divergence between the predictions using  $f^t$ .

### Producing bias attributes via biased local prediction.

On the other hand, for a bias-conflicting sample  $\tilde{x}$ , its bias attribute  $\tilde{a}_{bias}$  is expected to be similar to  $a_{bias}^j$  of  $x^j$ , which is sampled from an instance from a different category (as described in Sec. 3.2.1).

To identify and relate such bias attributes, we take the local client model  $f_k^{t-1}$  as the guidance. Note that, compared to the aggregated server model, client models produced at prior iterations are considered to be affected more by local biased data, which is more likely to predict the output  $f_k^{t-1}(x)$  based on its hidden bias attributes. Therefore, we define the similarity between the bias attributes  $\tilde{a}_{bias}$  and  $a_{bias}^j$ , which is now calculated and guided by the difference between the predictions of  $\tilde{x}$  and  $x^j$  using  $f_k^{t-1}$ . Specifically, we minimize:

$$\mathcal{L}_k = d_{KL}(f_k^{t-1}(\tilde{x}), f_k^{t-1}(x^j)), \quad (4)$$

where  $(t-1)$  denotes the training round.

From the above design and derivation, we have the objective for training BEA as:

$$\mathcal{L}_{total} = \mathcal{L} + \mathcal{L}_k. \quad (5)$$

As depicted in Figure 3, via minimization of  $\mathcal{L}$ , BEA will be optimized so that the semantic attribute  $\tilde{a}_{sem}$  of  $\tilde{x}$  will be updated and be close to  $a_{sem}^i$  of  $x^i$ . On the other hand, minimizing  $\mathcal{L}_k$  encourages the bias attribute  $\tilde{a}_{bias}$  of  $\tilde{x}$  to be updated as  $a_{bias}^j$  of  $x^j$ . In other words, optimization of BEA would encourage the generated samples whose semantic and bias attributes are extracted from training data of distinct classes.

While our BEA can be viewed as performing mixed-sample data augmentation, existing MSDA methods [7, 12, 16, 17, 38, 49, 52] are only designed to produce handcrafted augmentation outputs, which may not necessarily to be



---

**Algorithm 1: Training of FedBEAL**

---

**Input:**  $T, T_w, K, D = \{D_1, D_2, \dots, D_K\}, p, g_1, g_2, \dots, g_K, f^0$ , local epochs  $E_g$  and  $E_f$ , learning rate  $\eta_g$  and  $\eta_f$

**for**  $t = 0, 1, \dots, T - 1$  **do**

**for**  $k = 1, 2, \dots, K$  **in parallel do**

**if**  $t \geq T_w$  **then**

            TrainBEA( $f^t, f_k^{t-1}$ )

$f_k^t \leftarrow \text{LocalUpdate}(f^t)$

$f^{t+1} \leftarrow \sum_{k=1}^K \frac{|D_k|}{|D|} f_k^t$

**Output:** return  $f^T$

**TrainBEA**( $f^t, f_k^{t-1}$ )

**for**  $e = 1, 2, \dots, E_g$  **do**

**for**  $(x^i, x^j)$  of  $D_k$  **do**

$\tilde{x} \leftarrow g_k(x^i, x^j)$

$L \leftarrow d_{KL}(f^t(\tilde{x}), f^t(x^i))$

$L_k \leftarrow d_{KL}(f_k^{t-1}(\tilde{x}), f_k^{t-1}(x^j))$

$L_{total} \leftarrow L + L_k$

$g_k \leftarrow g_k - \eta_g \nabla L_{total}$

**LocalUpdate**( $f^t$ )

$f_k^t \leftarrow f^t$

**for**  $e = 1, 2, \dots, E_f$  **do**

**for**  $(x^i, x^j, y^i, y^j)$  of  $D_k$  **do**

**if**  $t \geq T_w$  and  $Uniform(0, 1) < p$  **then**

$\tilde{x}, \tilde{y} \leftarrow g_k(x^i, x^j), y^i$

$L_{cls} \leftarrow CrossEntropy(f_k^t(\tilde{x}), \tilde{y})$

**else**

$L_{cls} \leftarrow CrossEntropy(f_k^t(x^i), y^i)$

$f_k^t \leftarrow f_k^t - \eta_f \nabla L_{cls}$

    return  $f_k^t$

---

bias-conflicting. For example, spatial location-based augmentations (e.g., CutMix [49], FMix [12]) only fuse two images by replacing a region of one image with that from another, alleviating only high-level bias (e.g., background bias [36]). Style-based augmentations [16, 17, 52] are only capable of alleviating low-level biases by mixing style and content from distinct images. As verified in Section 4, learning of BEA would be desirable for debiased FL.

### 3.3. Training of FedBEAL

**Debiased local update.** After BEA is learned and deployed at each client  $k$ , we perform debias local updates by training each local model  $f_k^t$  using additionally produced bias-conflicting data pairs (i.e.,  $\tilde{x}$  and  $\tilde{y}$ ). To further improve the robustness of our framework, we follow [49] to consider several techniques at this local update stage. That is, we define  $p \in [0, 1]$  as the probability of augmenting

each data batch to control the degree of debiasing. Moreover, we define the warm-up round  $T_w$  (i.e., BEA is activated after round  $T_w$ ) to avoid undesirable augmentation outputs harmful to local training happening in the beginning stage. With bias-conflicting data and the introduced learning techniques, we are able to enforce the local model to be better guided by the semantic information while suppressing the bias.

**Global aggregation.** For each training iteration, once the debiased local updates are performed, we then collect and aggregate the learned weights of each local model (weighted by the size of the corresponding local dataset [35]). To be more specific, the global model for the next round  $f^{t+1}$  can be calculated as follows:

$$f^{t+1} = \sum_{k=1}^K \frac{|D_k|}{|D|} f_k^t. \quad (6)$$

With the convergence of the overall training process, the final global model can be applied to perform classification on unbiased test data. The pseudo-code of our complete FedBEAL framework is summarized in Algorithm 1.

## 4. Experiments

### 4.1. Datasets and Implementation Details

**Datasets.** To evaluate the effectiveness and applicability of our learning scheme in different types of bias, we consider three commonly used biased datasets, including Colored MNIST [5] (with color bias), Corrupted CIFAR-10 [15] (with corruption bias), and Collage CIFAR-10 [44] (with collaged images as bias). Colored MNIST contains images of hand-written digits colored with different colors. Corrupted CIFAR-10 includes images applied with random corruptions (e.g., noises, blurring, brightness/contrast adjustment). In Collage CIFAR-10, a sample is combined with four images originating from four different datasets, including MNIST [26], Fashion MNIST [48] and SVHN [37] that jointly serve as bias attributes, and CIFAR-10 [25] as the semantic information. As noted in Section 3.1, we distribute the training set to  $K$  clients, where  $K$  is set to 10 across all our experiments. To quantify the severity of local bias in training data, we further define the ratio for the amount of biased local data  $\beta$ .

**Implementation details.** For Colored MNIST, Corrupted CIFAR-10, and Collage CIFAR-10, input images are resized to  $28 \times 28$ ,  $32 \times 32$ , and  $64 \times 64$  pixels. For simplicity, we use LeNet [26] as the classifier  $f$  for Colored MNIST and ResNet-18 [14] for Corrupted CIFAR-10 and Collage CIFAR-10. A U-Net [39] with the encoder of ResNet-18 is adopted as the augments  $g$ . The communication round  $T$

Dataset	Colored MNIST		Corrupted CIFAR-10		Collage CIFAR-10	
Bias ratio $\beta$	0.99	0.999	0.99	0.999	0.99	0.999
<b>Baselines</b>						
SOLO	46.90	14.46	16.80	13.19	12.28	10.58
FedAvg [35]	93.90	72.67	49.03	40.28	52.93	36.91
<b>Centralized Debiasing Methods</b>						
LfF [36]	87.64	55.27	53.47	42.25	46.53	26.96
SoftCon [18]	96.75	86.39	55.38	47.61	54.19	42.98
Lee <i>et al.</i> [27]	90.28	61.35	54.86	45.90	41.02	22.58
<b>Data Heterogeneous Federated Learning</b>						
FedProx [30]	94.51	73.07	44.06	34.01	41.87	25.94
SCAFFOLD [21]	95.01	68.41	41.73	34.35	38.37	33.85
MOON [29]	93.33	69.37	36.79	26.06	34.71	19.97
FedBN [31]	N/A	N/A	48.46	36.52	46.51	32.53
Ours	<b>98.58</b>	<b>91.99</b>	<b>59.18</b>	<b>49.09</b>	<b>69.53</b>	<b>64.53</b>

Table 1. **Comparisons to SOTA federated learning and debiasing algorithms.** **Bold** denotes the best result, while underline denotes the second best. Note that in Colored MNIST, FedBN is not applicable due to disregard of Batch Normalization layers.

is set to 100. For each round, each client train their  $g$  and  $f$  sequentially for 5 epochs using the SGD optimizer, with the batch size of 64, the learning rate of 0.01, the momentum of 0.9, and the weight decay of 0.00001. We implement our model using PyTorch, and conduct training on a single NVIDIA 3090 GPU with 24GB memory.

## 4.2. Quantitative Evaluation

### 4.2.1 Comparisons to debiasing and FL methods

We first compare proposed learning scheme with existing centralized debiasing [18, 27, 36] and heterogeneous federated learning [21, 29–31, 35] methods. In our experiments, SOLO and FedAvg [35] are viewed as baselines. The former only performs local training without global averaging of client models, while the latter is the fundamental framework for all the other methods reported in this section. Note that we report the mean accuracy of each local model in SOLO. As shown in Table 1, we evaluate state-of-the-art methods with the three datasets with  $\beta$  set from 0.99 to 0.999. From the upper half of Table 1, we applied existing debiasing methods designed for centralized machine learning [18, 27, 36] to debias local update at each client. For example, SoftCon [18] enabled each client to preserve intra-bias features dissimilarities to debias the model, which improved the results of Colored MNIST with  $\beta$  of 0.999 by 13.72%. On the other hand, from the bottom half of Table 1, existing FL approaches tackled data heterogeneity by preserving the consistency between the local and global models. It can be observed that such constraints were not sufficient to mitigate severe local bias and only slightly improved the performance (*e.g.*, FedProx [30] improved the accuracy by 0.4% on Colored MNIST with  $\beta$  of 0.999). Instead, our FedBEAL performed favorably against the above methods on all datasets (*e.g.*, accuracy improvements of

Dataset	Colored MNIST		Corrupted CIFAR-10		Collage CIFAR-10		Avg.
Bias ratio $\beta$	0.99	0.999	0.99	0.999	0.99	0.999	
FedAvg [35]	93.90	72.67	49.03	40.28	52.93	36.91	57.62
Mixup [50]	91.38	74.76	53.98	40.85	50.13	37.65	58.13
CutMix [49]	82.73	59.55	41.39	31.69	<b>71.26</b>	<u>63.98</u>	58.43
MixStyle [52]	<b>99.13</b>	<b>99.20</b>	<u>58.99</u>	<u>46.27</u>	49.75	34.09	<u>64.57</u>
Ours	<u>98.58</u>	<u>91.99</u>	<b>59.18</b>	<b>49.09</b>	<u>69.53</u>	<b>64.53</b>	<b>72.15</b>

Table 2. **Comparisons to MSDA methods for debiased FL.** **Bold** denotes the best result, while underline denotes the second best.

19.32% on Colored MNIST with  $\beta$  of 0.999). These quantitative comparisons verify that our proposed FL approach removes local biases across different clients for improved classification performances.

### 4.2.2 Comparisons to MSDA methods

To further verify the effectiveness of our augmentation scheme, we further compare our method with state-of-the-art mixed sample data augmentation algorithms [49, 50, 52]. Existing handcrafted MSDA methods are generally designed to handle particular types of bias and cannot easily generalize to bias types not defined in advance. As shown in Table 2, MixStyle [52] benefited low-level biases (*e.g.*, color or corruption bias) by transferring style information of the images and improved the accuracies from 5.99% to 26.53% on Colored MNIST and Corrupted CIFAR-10. However, such augmentations was not able to mitigate high-level biases (*e.g.*, background bias [36]), as the performance of MixStyle dropped from 2.82% to 3.18% on Collage CIFAR-10. On the other hand, CutMix significantly improved the accuracy by 27.07% on Collage CIFAR-10 with  $\beta$  of 0.999 since the cut-and-paste operation efficiently removed high-level regional bias. However, it failed to handle low-level color and corruption biases in Colored MNIST and Corrupted CIFAR-10 and degraded the performance from 7.64% to 13.12%. Compared to such MSDA methods, our approach learns to find the optimal BEA and thus exhibits more robust debiasing effects on various bias types. As shown in the last column of Table 2, our method performed favorably compared to MSDA approaches and achieved improved accuracy by 14.53%, indicating the robustness and generalization capability of our augmentation scheme to different bias types.

### 4.2.3 Debiasing server and client models

As indicated in Section 3.2.2, the design and learning objectives for our BEA are based on the assumption that local models are relatively biased compared to the global aggregated one at each iteration. To verify this assumption, we quantitatively compare the *bias level* of the global and local models in FedAvg and FedBEAL on the Colored MNIST

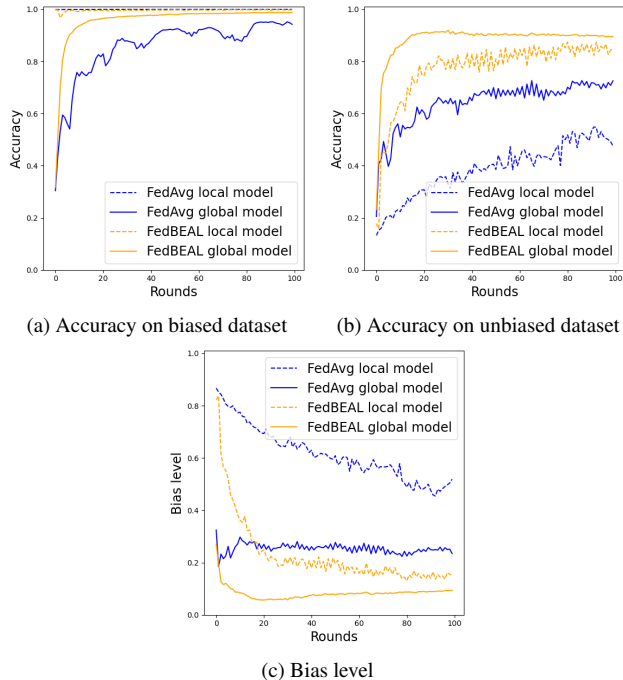


Figure 4. **Learning curve comparisons of global/local models from FedAvg and FedBEAL on Colored MNIST.** (a) and (b) show accuracies for biased and unbiased datasets, and (c) compares the bias level  $\mathcal{S}$  (see Equation (7)). Note that the local model with FedBEAL shows improved debiased performance, and its global model also exhibits improved unbiased ability over FedAvg.

dataset. Given a biased dataset  $D_k$  from client  $k$  and an unbiased testing dataset  $D_{test}$ , we first define the bias level  $\mathcal{S}$  of the local model  $f_k^t$  and the global model  $f^t$  as follows:

$$\mathcal{S} = 1 - \frac{Acc_{unbias}}{Acc_{bias}}, \quad (7)$$

where  $Acc_{bias}$  and  $Acc_{unbias}$  are the accuracies evaluated on  $D_k$  and  $D_{test}$ , respectively. In other words, the model is biased (*i.e.*,  $\mathcal{S}$  is higher) if the model achieves high accuracy on the biased dataset while performing relatively unfavorably on the unbiased dataset.

Based on the above setting, we train our model on Colored MNIST with the bias ratio  $\beta$  of 0.999. As illustrated in Figure 4c, while the local model of FedBEAL was relatively biased compared to the global model (see orange curves), we were able to gradually debias such models for improved performances when comparing to FedAvg. The above results support the design and learning scheme for the proposed BEA.

### 4.3. Qualitative Evaluation

**Representation visualization.** We now qualitatively assess the ability of FedBEAL to derive semantic-aware and debiased feature representations. As shown in Figure 5, we

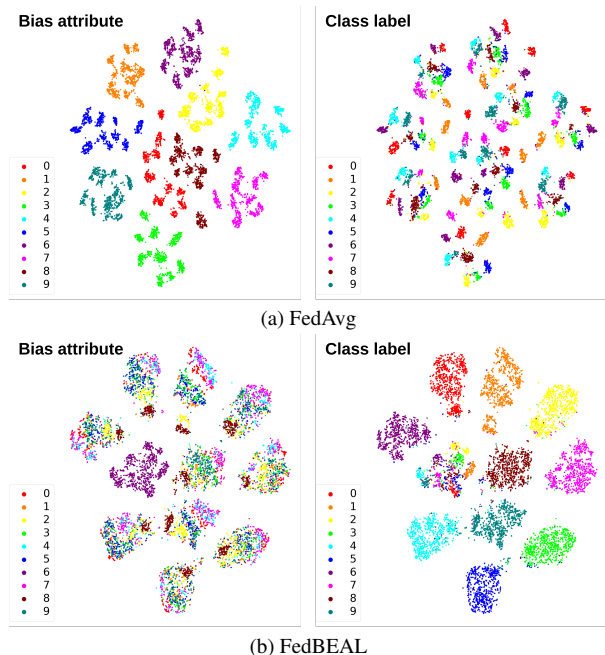


Figure 5. **t-SNE comparisons between FedAvg and FedBEAL on Colored MNIST.** Data points in the left column are colored based on the bias attributes (*i.e.*, color), while those in the right column are colored based on the class labels.

apply t-SNE [45] to compare the hidden representation derived by the global model of FedAvg and our approach on Colored MNIST. In Figure 5a, we see that features extracted by FedAvg were grouped according to *bias* attributes and were not properly separated with respect to the class labels. In contrast, features derived by our model remained relatively uncorrelated in terms of the bias attributes, and the separation between different class clusters was more significant. The above observation indicates that our proposed bias-eliminating augmentation learning allows the derivation of discriminative and debiased features.

### Visualization of augmented bias-conflicting samples.

We now show example augmented images  $\tilde{x}$  produced by our method, which is expected to preserve the categorical information of  $x^i$  and impose the bias from  $x^j$ . In Figure 6, we first see example images for Colored MNIST, and we observe  $\tilde{x}$  obtained the digit color from  $x^j$  while preserving the original digit shape as of  $x^i$ . From the second image set of Corrupted CIFAR-10,  $\tilde{x}$  inherited the high chromatic *impluse noise* from  $x^j$  while still maintaining semantically recognizable foreground objects. As for Collage CIFAR-10, our modulators  $M$  successfully captured the unbiased bottom-left image region for augmentation. From the examples, we confirm that our proposed BEA is capable of capturing inherent dataset bias while preserving desirable semantic attributes for augmenting bias-conflicting samples.

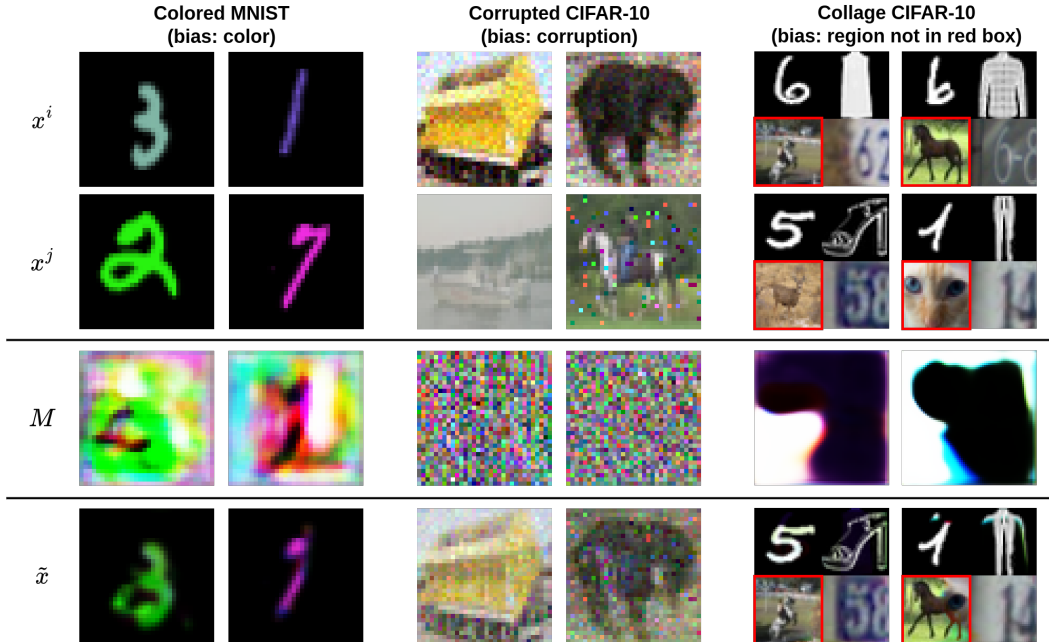


Figure 6. **Visualization of images produced by BEA.** Based on the mask learned by the BEA modulator  $M$ , the augmented bias-conflicting  $\tilde{x}$  can be seen as the mixture of the content from  $x^i$  and the bias from  $x^j$ .

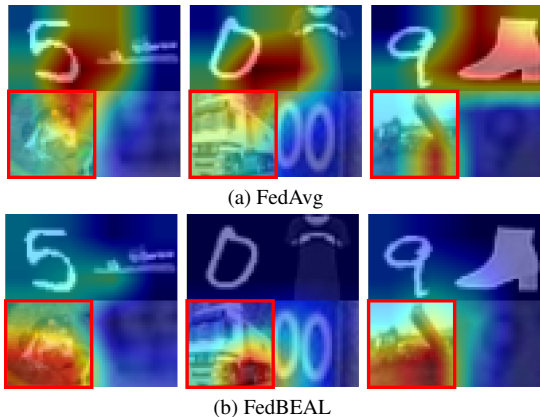


Figure 7. **Grad-CAM comparisons between FedAvg and FedBEAL on Collage CIFAR-10.** Compared to FedAvg results in (a), the attention map for FedBEAL in (b) better identify the object region of interest for classification (in red rectangles).

**Grad-CAM visualization.** Grad-CAM [41] is commonly used to visually explain how deep learning models make classification decisions. To verify the effectiveness of the proposed learning scheme, we consider FedAvg and our proposed method on the Collage CIFAR-10 dataset with  $\beta$  of 0.99, and we apply Grad-CAM to interpret the trained global models during classification (see Figure 7). From Figure 7a, we see that the global model trained with FedAvg attended to ambiguous or irrelevant image regions, implying the lack of ability to indicate regions with proper semantic features for classification. In Figure 7a), we see

that the global model trained by our proposed FedBEAL attended image regions on the augmented samples, which are correlated to the categorical information of interest. This also explains the reason why our FedBEAL is able to achieve satisfactory performances on debiased FL tasks.

## 5. Conclusion

In this paper, we addressed the challenging problem of debiased FL and proposed FedBEAL for mitigating local biases. By introducing and learning Bias-Eliminating Augmenters at each client, bias-conflicting samples can be automatically learned. The learning of BEA can be simply utilized by the global server and local client models obtained during the training progress, and thus no prior knowledge of bias type or annotation would be required. We conducted extensive experiments, including comparisons to state-of-the-art debiasing, FL, and MSDA methods, and visualization of augmented images, which quantitatively and qualitatively confirmed the effectiveness and robustness of our proposed approach in discovering and solving unknown dataset bias in federated learning schemes.

## Acknowledgment

This work is supported in part by the National Science and Technology Council under grant 111-2634-F-002-020. We also thank to National Center for High-performance Computing (NCHC) for providing additional computational and storage resources.



## References

- [1] Annie Abay, Yi Zhou, Nathalie Baracaldo, Shashank Rajamoni, Ebube Chuba, and Heiko Ludwig. Mitigating bias in federated learning. *arXiv preprint arXiv:2012.02447*, 2020. [2](#)
- [2] Durmus Alp Emre Acar, Yue Zhao, Ruizhao Zhu, Ramon Matas, Matthew Mattina, Paul Whatmough, and Venkatesh Saligrama. Debiasing model updates for improving personalized federated training. In *International Conference on Machine Learning*, pages 21–31. PMLR, 2021. [2](#)
- [3] Anonymous. Feddebias: Reducing the local learning bias improves federated learning on heterogeneous data. In *Submitted to The Eleventh International Conference on Learning Representations*, 2023. under review. [2](#)
- [4] Anonymous. Learning to aggregate: A parameterized aggregator to debias aggregation for cross-device federated learning. In *Submitted to The Eleventh International Conference on Learning Representations*, 2023. under review. [2](#)
- [5] Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant risk minimization. *arXiv preprint arXiv:1907.02893*, 2019. [5](#)
- [6] Hyojin Bahng, Sanghyuk Chun, Sangdoon Yun, Jaegul Choo, and Seong Joon Oh. Learning de-biased representations with biased representations. In *International Conference on Machine Learning*, pages 528–539. PMLR, 2020. [1, 2](#)
- [7] Ali Dabouei, Sobhan Soleymani, Fariborz Taherkhani, and Nasser M Nasrabadi. Supermix: Supervising the mixing data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13794–13803, 2021. [3, 4](#)
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018. [1](#)
- [9] Yahya H Ezzeldin, Shen Yan, Chaoyang He, Emilio Ferrara, and Salman Avestimehr. Fairfed: Enabling group fairness in federated learning. *arXiv preprint arXiv:2110.00857*, 2021. [2](#)
- [10] Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias Bethge, and Felix A Wichmann. Shortcut learning in deep neural networks. *Nature Machine Intelligence*, 2(11):665–673, 2020. [1, 2](#)
- [11] Robert Geirhos, Patricia Rubisch, Claudio Michaelis, Matthias Bethge, Felix A. Wichmann, and Wieland Brendel. Imagenet-trained CNNs are biased towards texture; increasing shape bias improves accuracy and robustness. In *International Conference on Learning Representations*, 2019. [1, 2](#)
- [12] Ethan Harris, Antonia Marcu, Matthew Painter, Mahesan Niranjan, Adam Prügél-Bennett, and Jonathon Hare. Fmix: Enhancing mixed sample data augmentation. *arXiv preprint arXiv:2002.12047*, 2020. [3, 4, 5](#)
- [13] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. [2](#)
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1, 5](#)
- [15] Dan Hendrycks and Thomas G Dietterich. Benchmarking neural network robustness to common corruptions and surface variations. *arXiv preprint arXiv:1807.01697*, 2018. [5](#)
- [16] Dan Hendrycks, Norman Mu, Ekin D Cubuk, Barret Zoph, Justin Gilmer, and Balaji Lakshminarayanan. Augmix: A simple data processing method to improve robustness and uncertainty. *arXiv preprint arXiv:1912.02781*, 2019. [3, 4, 5](#)
- [17] Minui Hong, Jinwoo Choi, and Gunhee Kim. Stylemix: Separating content and style for enhanced data augmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14862–14870, 2021. [3, 4, 5](#)
- [18] Youngkyu Hong and Eunho Yang. Unbiased classification through bias-contrastive and bias-balanced learning. *Advances in Neural Information Processing Systems*, 34:26449–26461, 2021. [1, 2, 3, 6](#)
- [19] Wenke Huang, Mang Ye, and Bo Du. Learn from others and be yourself in heterogeneous federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10143–10153, 2022. [2](#)
- [20] Peter Kairouz, H Brendan McMahan, Brendan Avent, Aurélien Bellet, Mehdi Bennis, Arjun Nitin Bhagoji, Kallista Bonawitz, Zachary Charles, Graham Cormode, Rachel Cummings, et al. Advances and open problems in federated learning. *Foundations and Trends® in Machine Learning*, 14(1–2):1–210, 2021. [2](#)
- [21] Sai Praneeth Karimireddy, Satyen Kale, Mehryar Mohri, Sashank Reddi, Sebastian Stich, and Ananda Theertha Suresh. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, pages 5132–5143. PMLR, 2020. [2, 6](#)
- [22] Kimmo Karkkainen and Jungseock Joo. Fairface: Face attribute dataset for balanced race, gender, and age for bias measurement and mitigation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1548–1558, January 2021. [1](#)
- [23] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. Supervised contrastive learning. *Advances in Neural Information Processing Systems*, 33:18661–18673, 2020. [2](#)
- [24] Byungju Kim, Hyunwoo Kim, Kyungsu Kim, Sungjin Kim, and Junmo Kim. Learning not to learn: Training deep neural networks with biased data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9012–9020, 2019. [1](#)
- [25] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical Report 0, University of Toronto, Toronto, Ontario, 2009. [5](#)
- [26] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. [5](#)

- [27] Jungsoo Lee, Eungyeup Kim, Juyoung Lee, Jihyeon Lee, and Jaegul Choo. Learning debiased representation via disentangled feature augmentation. *Advances in Neural Information Processing Systems*, 34:25123–25133, 2021. 1, 2, 6
- [28] Qinbin Li, Yiqun Diao, Quan Chen, and Bingsheng He. Federated learning on non-iid data silos: An experimental study. In *2022 IEEE 38th International Conference on Data Engineering (ICDE)*, pages 965–978. IEEE, 2022. 2
- [29] Qinbin Li, Bingsheng He, and Dawn Song. Model-contrastive federated learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10713–10722, 2021. 2, 6
- [30] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2:429–450, 2020. 2, 6
- [31] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. Fedbn: Federated learning on non-iid features via local batch normalization. *arXiv preprint arXiv:2102.07623*, 2021. 2, 6
- [32] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9572–9581, 2019. 1
- [33] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020. 2
- [34] Mi Luo, Fei Chen, Dapeng Hu, Yifan Zhang, Jian Liang, and Jiashi Feng. No fear of heterogeneity: Classifier calibration for federated learning with non-iid data. *Advances in Neural Information Processing Systems*, 34:5972–5984, 2021. 2
- [35] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Artificial intelligence and statistics*, pages 1273–1282. PMLR, 2017. 1, 3, 5, 6
- [36] Junhyun Nam, Hyuntak Cha, Sungsoo Ahn, Jaeho Lee, and Jinwoo Shin. Learning from failure: De-biasing classifier from biased classifier. *Advances in Neural Information Processing Systems*, 33:20673–20684, 2020. 1, 2, 5, 6
- [37] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on Deep Learning and Unsupervised Feature Learning 2011*, 2011. 5
- [38] Viktor Olsson, Wilhelm Tranheden, Juliano Pinto, and Lennart Svensson. Classmix: Segmentation-based data augmentation for semi-supervised learning. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1369–1378, 2021. 3, 4
- [39] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5
- [40] Shiori Sagawa\*, Pang Wei Koh\*, Tatsunori B. Hashimoto, and Percy Liang. Distributionally robust neural networks. In *International Conference on Learning Representations*, 2020. 1
- [41] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*, pages 618–626, 2017. 8
- [42] Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The pitfalls of simplicity bias in neural networks. *Advances in Neural Information Processing Systems*, 33:9573–9585, 2020. 1, 2
- [43] Enzo Tartaglione, Carlo Alberto Barbano, and Marco Granetto. End: Entangling and disentangling deep representations for bias correction. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 13508–13517, 2021. 2
- [44] Damien Teney, Ehsan Abbasnejad, Simon Lucey, and Anton van den Hengel. Evading the simplicity bias: Training a diverse set of models discovers solutions with superior ood generalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16761–16772, 2022. 2, 5
- [45] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008. 7
- [46] Haohan Wang, Zexue He, Zachary C Lipton, and Eric P Xing. Learning robust representations by projecting superficial statistics out. *arXiv preprint arXiv:1903.06256*, 2019. 2
- [47] Jianyu Wang, Qinghua Liu, Hao Liang, Gauri Joshi, and H Vincent Poor. Tackling the objective inconsistency problem in heterogeneous federated optimization. *Advances in neural information processing systems*, 33:7611–7623, 2020. 2
- [48] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 5
- [49] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. Cutmix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6023–6032, 2019. 3, 4, 5, 6
- [50] Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017. 6
- [51] Yue Zhao, Meng Li, Liangzhen Lai, Naveen Suda, Damon Civin, and Vikas Chandra. Federated learning with non-iid data. *arXiv preprint arXiv:1806.00582*, 2018. 2
- [52] Kaiyang Zhou, Yongxin Yang, Yu Qiao, and Tao Xiang. Domain generalization with mixstyle. *arXiv preprint arXiv:2104.02008*, 2021. 3, 4, 5, 6
- [53] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In Marina Meila and Tong Zhang, editors, *Proceedings of the*

38th International Conference on Machine Learning, volume 139 of *Proceedings of Machine Learning Research*, pages 12878–12889. PMLR, 18–24 Jul 2021. [2](#)