

Federated Robustness Propagation: Sharing Adversarial Robustness in Federated Learning

Junyuan Hong*

Haotao Wang†

Zhangyang Wang†

Jiayu Zhou*

*Department of Computer Science and Engineering
Michigan State University, East Lansing, MI 48823, USA
{hongju12, jiayuz}@msu.edu

†Department of Electrical and Computer Engineering
University of Texas at Austin, Austin TX 78712, USA
{htwang, atlaswang}@utexas.edu

Abstract

Federated learning (FL) emerges as a popular distributed learning schema that learns a model from a set of participating users without requiring raw data to be shared. One major challenge of FL comes from heterogeneity in users, which may have distributionally different (or *non-iid*) data and varying computation resources. Just like in centralized learning, FL users also desire model robustness against malicious attackers at test time. Whereas **adversarial training** (AT) provides a sound solution for centralized learning, extending its usage for FL users has imposed significant challenges, as many users may have very limited training data as well as tight computational budgets, to afford the data-hungry and costly AT. In this paper, we study a novel learning setting that propagates adversarial robustness from high-resource users that can afford AT, to those low-resource users that cannot afford it, during the FL process. **We show that existing FL techniques cannot effectively propagate adversarial robustness among *non-iid* users**, and propose a simple yet effective propagation approach that transfers robustness through carefully designed batch-normalization statistics. We demonstrate the rationality and effectiveness of our method through extensive experiments. Especially, the proposed method is shown to grant FL remarkable robustness even when only a small portion of users afford AT during learning. Codes will be published upon acceptance.

1 Introduction

Federated learning (FL) [1] is a learning paradigm that trains models from distributed users or participants (e.g., mobile devices) without requiring raw training data to be shared, alleviating the rising concern of privacy issues when learning with sensitive data and facilitating learning deep models by enlarging the amount of data to be used for training. In a typical FL algorithm, each user trains a model locally using their own data and a server iteratively aggregates users' incremental updates or intermediate models, converging to a model that fuses training information from all users. A major challenge in FL comes from the heterogeneity of users. One source of heterogeneity is distributional differences in training data collected by users from diverse user groups [2, 3]. Yet another source is the difference of computing resources, as different types of hardware used by users usually result in varying computation budgets. For example, consider an application scenario of FL from mobile phones [4], where different types of mobile phones (e.g., generations of the same brand) may have drastically different computational power.

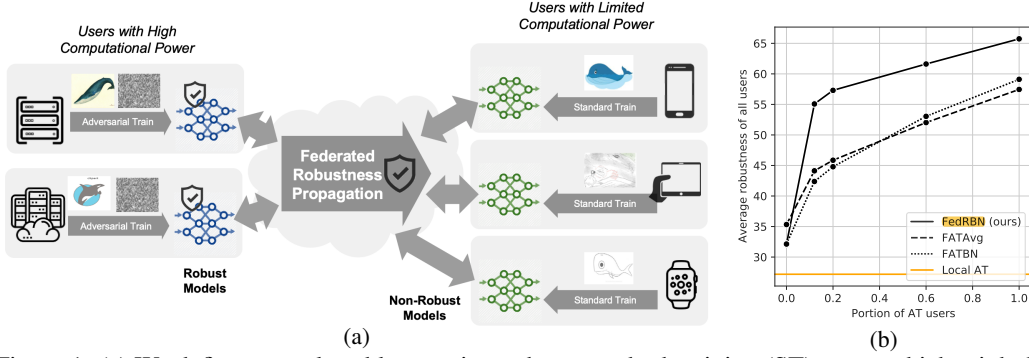


Figure 1: (a) We define a novel problem setting, where standard-training (ST) users, which might be limited by data or computational resources, can “share” robustness from adversarial-training (AT) users who can afford it. (b) Comparison of robustness on a varying portion of AT users, where a 5-domain digit recognition dataset is distributed to 50 users in total and details are in Section C.7.

Data heterogeneity should be carefully handled during the learning as a single model trained by FL may fail to accommodate the differences [5]. A variety of approaches have been proposed to address the issue, such as customizing network structures [6, 7] or tailoring training strategies [2, 8] for each user. Even though hardware heterogeneity is ubiquitous, their impacts to FL processes and possible solutions have received very limited attention so far.

The impacts of the two types of heterogeneity become aggravated when participating users’ desire adversarial robustness during the inference stage, against imperceptible noise that can significantly mislead model predictions. To address this issue, a straightforward extension of FL, federated adversarial training (FAT), can be adopted, which idea was explored in [9, 10]. Locally, each user trains models with adversarially augmented samples, namely adversarial training (AT) [11]. As studied in central settings, the AT is data-thirsty and computationally expensive [12]. Therefore, involving a fair amount of users in FAT is essential, given the fact that each individual user may not have enough data to perform AT. However, this implies an increasing difficulty for fitting diverse data distributions and more intensive computation for each user, which could be 3 – 10 times more costly than the standard equivalent [12, 13]. The computation overhead can be prohibitive for FL users with limited computational budget such as mobile devices. As such, it is often unrealistic to enforce *all* users in a FL process to conduct AT locally, despite the fact that the robustness is indeed a strongly desired or even required property for all users. This conflict raises a challenging yet interesting question: Is it possible to *propagate adversarial robustness in FL* so that budget-limited users can benefit from robustness training of users with abundant computational resources?

Motivated by the question above, we formulate a novel problem setting called Federated Robustness Propagation (FRP), as depicted in Fig. 1a. We consider a rather common non-iid FL setting that involves budget-sufficient users (AT users) that conduct adversarial training, and budget-limited ones (ST users) that can only afford standard training. The goal of FRP is to propagate the adversarial robustness from AT users to ST users. Note that sharing adversarial data is prohibited for mitigating the overhead of adversarial-data generation, due to the privacy consideration of the FL framework. In Fig. 1b, we show that independent AT by users without FL (local AT) will not yield a robust model since each user has scarce training data. Directly extending an existing FL algorithm *FedAvg* [1] or a heterogeneity-mitigated one *FedBN* [6] with AT treatments, named as FATAvg and FATBN, give very limited capability of propagating robustness. The limitation may arise from the following two reasons: **1)** Heterogeneity of users results in distinct noise behavior in different users, degrading the transferability of model robustness. **2)** Structural differences between clean and adversarial samples may further impede robustness propagation [14].

To address the aforementioned challenges, we propose a novel method Federated Robust Batch-Normalization (FedRBN) to facilitate propagation of adversarial robustness among FL users. **1)** We design a surprisingly simple linear method that transmits the robustness by copying batch-normalization (BN) statistics, inspired by the strong connection between model robustness and statistic parameters in the BN layer [15]; **2)** To efficiently propagate robustness among non-iid users, we weight and average multiple AT users’ statistics as BN for every ST user; **3)** Facing the structural difference between the clean and adversarial data, we train two separate BNs for each data type, which are adaptively chosen at the inference stage. Our method is communication-efficient as it

only incurs an one-time additional communication after training. We conduct extensive experiments demonstrating the feasibility and effectiveness of the proposed method. In Fig. 1b, we highlight some experimental results from Section 5. When only 20% of non-iid users used AT during learning, the proposed FedRBN yields robustness, competitive with the best all-AT-user baseline (FATBN) by only a 2% drop (out of 59%) on robust accuracy. Note that even if our method with 100% AT users increase the upper bound of robustness, such a bound is usually not attainable because of the presence of resource-limited users that cannot afford AT during learning.

2 Related Work

Federated learning for robust models. The importance of adversarial robustness in the context of federated learning, i.e., federated adversarial training (FAT), has been discussed in a series of recent literature [9, 10, 16]. Zizzo *et al.* [9] empirically evaluated the feasibility of practical FAT configurations (e.g., ratio of adversarial samples) augmenting FedAvg with AT but only in *iid* and label-wise non-*iid* scenarios. The adversarial attack in FAT was extended to a more general affine form, together with theoretical guarantees of distributional robustness [10]. It was found that in a communication-constrained setting, a significant drop exists both in standard and robust accuracies, especially with non-*iid* data [17]. In addition to the challenges investigated above, this work studies challenges imposed by hardware heterogeneity in FL, which was rarely discussed. Especially, when only limited users have devices that afford AT, we strive to efficiently share robustness among users, so that users without AT capabilities can also benefit from such robustness.

Robust federated optimization. Another line of related work focuses on the robust aggregation of federated user updates [16, 18]. Especially, Byzantine-robust federated learning [19] aims to defend malicious users whose goal is to compromise training, e.g., by model poisoning [20, 21] or inserting model backdoor [22]. Various strategies aim to eliminate the malicious user updates during federated aggregation [23, 19, 24, 25]. However, most of them assume the normal users are from similar distributions with enough samples such that the malicious updates can be detected as outliers. Therefore, these strategies could be less effective on attacker detection when a finite dataset is given [26]. Even though both the proposed FRP and Byzantine-robust studies work with robustness, but they have fundamental differences: the proposed work focus on *the robustness during inference*, i.e., after the model is learned and deployed, whereas Byzantine-robust work focus on the robust learning process. As such, all Byzantine-robust techniques can be combined with the proposed approach to provide training robustness.

3 Problem Setting: Federated Robustness Propagation (FRP)

In this section, we will review AT, present the unique challenges from hardware heterogeneity in FL and formulate the problem of federated robustness propagation (FRP). In this paper, we assume that a dataset D includes sampled pairs of images $x \in \mathbb{R}^d$ and labels $y \in \mathbb{R}^c$ from a distribution \mathcal{D} . Though we limit the data as images in this paper, our discussion could be generalized to other data forms. We model a classifier, mapping from the \mathbb{R}^d data/input space to classification logits $f : \mathbb{R}^d \rightarrow \mathbb{R}^c$, by a deep neural network (DNN). Whenever not causing confusing, we use the symbol of a model and its parameters interchangeably. For brevity, we slightly abuse $\mathbb{E}[\cdot]$ for both empirical average and expectation and use $[N]$ to denote $\{1, \dots, N\}$.

3.1 Standard training and adversarial training

An **adversarial attack** applies a bounded noise $\delta_\epsilon : \|\delta_\epsilon\| \leq \epsilon$ to an image x such that the perturbed image $A_\epsilon(x) \triangleq x + \delta_\epsilon$ can mislead a well-trained model to give a wrong prediction. The norm $\|\cdot\|$ can take a variety of forms, e.g., L_∞ -norm for constraining the maximal pixel scale. A model f is said to be *adversarially robust* if it can predict labels correctly on a perturbed dataset $\tilde{D} = \{(A_\epsilon(x), y) | (x, y) \in D\}$, and the standard accuracy on D should not be greatly impacted.

Consider the following general learning objective:

$$\min_f L(f, D) = \min_f \frac{1}{|D|} \sum_{(x, y) \in D} [(1 - q) \ell_c(f; x, y) + q \ell_a(f; x, y)], \quad (1)$$

where ℓ_c is a standard classification loss on clean images and ℓ_a is an adversarial loss promoting robustness. Eq. (1) performs *standard training* if $q = 0$, and *adversarial training* if $q \in (0, 1]$.

Without loss of generality, we limit our discussion for q_a as 0 or 0.5. A popular instantiation of Eq. (1) is based on PGD attack [27, 28]: $\ell_c(f; x, y) = \ell(f(x), y)$, $\ell_a(f; x, y) = \max_{\|\delta\| \leq \epsilon} \ell(f(x + \delta), y)$, where $\|\cdot\|$ is the L_∞ -norm, ℓ can be the cross-entropy loss, i.e., $\ell(f(x), y) = -\sum_{t=1}^c y_t \log(f(x)_t)$ where t is the class index and $f(x)_t$ represents the t -th output logit.

3.2 Problem setup and challenges

We start with a typical FL setting: a finite set of distributions \mathcal{D}_i for $i \in [C]$, from which a set of datasets $\{D_k\}_{k=1}^K$ are sampled and distributed to K users' devices. The users from distinct domains related with \mathcal{D}_i expect to optimize objectives like Eq. (1). Some users can afford AT training (*AT users*, $q_k = 0.5$) when the remaining users cannot afford and use standard training (*ST users*, $q_k = 0$). If the two groups of users train models separately, the models of ST users will be much less robust than those of AT ones. Note that data exchange among users is forbidden according to the FL setting for privacy concerns. The goal of *federated robustness propagation (FRP)* is to transfer the robustness from AT users to ST users at minimal computation and communication costs while preserve data locally. Formally, the FRP objective minimizes:

$$\text{FRP}(\{f_k\}; \{D_k | D_k \sim \mathcal{D}_i\}, \mathbf{q}) \triangleq \sum_{k \in [K]} \frac{1}{|D_k|} \sum_{(x, y) \in D_k} [(1 - q_k)\ell_c(f_k) + q_k\ell_a(f_k)], \quad (2)$$

where $\mathbf{q} \triangleq [q_1, \dots, q_K]$. Note that different from FAT [9], FRP assumes that D_k sampled from different distributions and that there are at least one zero entry in \mathbf{q} . In the federated setting, each user's model is trained separately when initialized by a global model, and is aggregated to a global model at the end of each epoch. A popular aggregation technique is FedAvg [1], which averages parameters by $f = \frac{1}{K} \sum_{k=1}^K a_k f_k$ with normalization coefficients a_k proportional to $|D_k|$.

Remarkably, Eq. (2) formalizes two types of common user heterogeneity in FL. The first one is the *hardware heterogeneity* where users' varying computation budgets result in sparsity of \mathbf{q} . A node of tight computation budget, e.g., smartphone, may join FL with $q = 0$, while a powerful one, e.g., desktop computer, uses $q = 0.5$ [4]. Besides, *data heterogeneity* is represented as \mathcal{D}_i differing by i . We limit our discussion as the common feature distribution shift (on x) in contrast to the label distribution shift (on y), as previously considered in [6]. Such distribution shift often happens when users are distributed across different environments, e.g., sensor data collected indoor and outdoor.

New Challenges. We emphasize that *jointly* addressing the two types of heterogeneity in Eq. (2) forms a new challenge, distinct from either of them considered exclusively. First, the sparsity in \mathbf{q} worsens the data heterogeneity as additional distribution shift in the hidden representations from adversarial augmentation [14]. That means even if two users sample from the same distribution, their classification layers may operate on different distributions.

Second, the data heterogeneity makes the transfer of robustness non-trivial [29]. Hendrycks *et al.* discussed the transfer of models adversarially trained on multiple domains and massive samples [30]. In [29], Shafahi *et al.* firstly studied the transferability of adversarial robustness from one data domain to another without the data-hungry problem. They proposed fine-tuning the robustness-sensitive layers in a neural network on a target domain. Distinguished from Shafahi *et al.*'s work, the FRP problem focuses on propagating robustness from multiple AT users to multiple ST users who have diverse distributions. Thus, fine-tuning all source models in ST users is often not possible due to prohibitive computation costs.

4 Method: Federated Robust Batch-Normalization (FedRBN)

4.1 Robustness propagation by copying debiased BN layers

In centralized learning, an important observation is that robustness is highly correlated with the BN statistics [14]. We extend this investigation to the FL setting, where we assume all other parameters are shared besides BN layers. There are significant differences in BN parameters (mean and variance) between ST and AT users from the same domain, as shown in Fig. 2a. This observation indicates that directly using local BN statistics can hardly grant robustness to an ST user, and suggests a possible way to transfer robustness through leveraging the BN layers from AT users in ST users upon predicting possible adversarial input images. However, the distributions of users from distinct

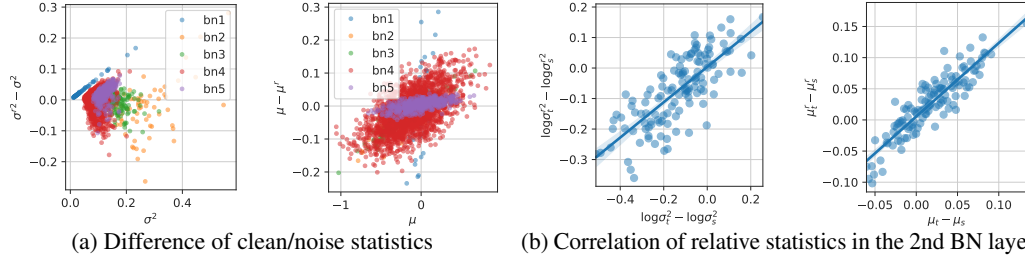


Figure 2: Models are trained with decoupled BN layers on Digits dataset. bn1 is the first BN layer in the network. (a) Results on SVHN. (b) The relative statistics are compared on MNIST versus SVHN.

domains can be quite different [31], and therefore directly copying BN among users can suffer from the distribution shift by domains. This motivates us to develop a shift-aware debiasing method.

To capture the domain bias, we can leverage the BN layers as they are modeling local distributions. Ideally, differential modeling two users will yield BN statistics for the corresponding two distribution, separately. However, one challenge here is that the BN statistics in non-iid AT and ST users are *biased simultaneously* by the domain difference and the adversarial noise. As such, directly differential modeling will capture the mixture of both types of bias and therefore would not be effective to infer the domain bias. Instead, we incrementally propose to simultaneously model clean data and noise data by BN for all AT users, since clean data are also available during training. To do so, we replace standard BN layers by ones that use the *dual batch-normalization (DBN)* structure [11], which keep two sets of BN statistics: one for clean data and one for noise data. The channel-wise mapping of DBN from the input layer x to output z is

$$z = w \left[(1 - h) \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon_0}} + h \frac{x - \mu^r}{\sqrt{\sigma^{r2} + \epsilon_0}} \right] + b, \quad (3)$$

where μ and σ^2 are the mean and variance over all non-channel dimensions, μ^r and σ^{r2} are their corresponding noised statistics, h serves as a model-wise switch, which is 0 for clean inputs or 1 for noised inputs, ϵ_0 is a small constant for the numerical stability. Different from prior work, e.g., [11, 14], we use a shared affine weights (w) and bias (b) for efficiency considerations. In the user-side training, we explicitly choose the clean or noise BN based on the input. Though we introduce DBN for a bias-inference purpose, the DBN still merits the model performance as it normalize representations more precisely as prior work investigated [11].

Transfer robustness from a single user via copying BN layers. With the clean BN embedded in DBN, we can estimate the distributional difference by an AT clean BN statistic tuple (μ_s, σ_s^2) and an ST (clean) BN statistic tuple (μ_t, σ_t^2) . Formally, we propose a debiased statistic estimation by

$$\hat{\mu}_t^r = \mu_s^r + \lambda(\mu_t - \mu_s), \quad \hat{\sigma}_t^{r2} = \sigma_s^{r2} \left(\sigma_t^2 / (\sigma_s^2 + \epsilon_0) \right)^\lambda, \quad (4)$$

where λ is a hyper-parameter in $[0, 1]$. Note that when the distributions are matched, i.e., $\mu_t = \mu_s$, then debiasing is not necessary and is thus vanished. Note that when λ is 0, the debias term is ignored. On the other hand, since the debias term is roughly estimated, we may not want to choose an over-confident value either, e.g., set $\lambda = 1$. Hence, we use $\lambda = 0.1$ generally in our experiments. To justify the rationality of Eq. (4), we contrast the $\mu_s - \mu_t$ with $\mu_s^r - \mu_t^r$ in Fig. 2b. The clean and noise BN statistics are estimated during training DBN, and we observe a strong correlation of the relative difference among domains both for the mean and variance.

To understand the debiasing method, we provide a principled analysis on a simplified one-dimensional example. We assume the noised inputs to a BN layer in user s can be approximated by an affine-noise model $\tilde{x}_s = \lambda x_s + \nu$, $\tilde{x}_t = \lambda x_t + \nu$, where $x_s \sim \mathcal{N}(\mu_s, \sigma_s^2)$, $x_t \sim \mathcal{N}(\mu_t, \sigma_t^2)$ and $\nu \sim \mathcal{N}(\mu', \sigma'^2)$ is domain-independent noise. λ is a constant scalar. We further assume ν is independent from x_s and x_t . Taking expectation gives $\mu_s^r = \lambda\mu_s + \mu'$, $\mu_t^r = \lambda\mu_t + \mu'$; $\sigma_s^{r2} = \lambda^2\sigma_s^2 + \sigma'^2$, $\sigma_t^{r2} = \lambda^2\sigma_t^2 + \sigma'^2$. Due to the invariance assumption of (μ', σ'^2) , we have: $\hat{\mu}_t^r = \mu_s^r + \lambda(\mu_t - \mu_s)$, $\hat{\sigma}_t^{r2} = \sqrt{\sigma_s^{r2} + \lambda^2(\sigma_t^2 - \sigma_s^2)}$. However, $\hat{\sigma}_t^{r2}$ is meaningless when $\sigma_s^{r2} + \lambda^2(\sigma_t^2 - \sigma_s^2) < 0$. To fix this, we use a division instead of subtraction to represent the relative relation in Eq. (4). This simplified example by no means serves as rigorous analysis, which is an open problem outside the scope of this paper.

Multiple-source propagation by weighted BN averaging. In FL with multiple AT users available, we propose a strategy to transfer robustness from multiple sources. Given N source BN statistics, we

use a weighted average to estimate the noise BN of a target ST user $\hat{\mu}_t^r = \sum_i^N \alpha_i \hat{\mu}_{t,s_i}^r$, where $\hat{\mu}_{t,s_i}^r$ is the estimated value from user s_i by Eq. (4). Likewise, $\hat{\sigma}_t^r$ can be estimated. However, the difference between the s_i -th adversarial distribution and the t -th counterpart is unknown. To tackle the issue, we first present the following result:

Lemma 4.1 (Informal). *Suppose the divergence between any data distribution \mathcal{D} and its adversarial distribution $\tilde{\mathcal{D}}$ is bounded, i.e., $d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}, \mathcal{D}) \leq d_\epsilon$ where $d_{\mathcal{H}\Delta\mathcal{H}}$ is $\mathcal{H}\Delta\mathcal{H}$ -divergence in hypothesis space \mathcal{H} . If a target model is formed by α_i -weighted average of models from D_{s_i} , the summation $\sum_i \alpha_i d_{\mathcal{H}\Delta\mathcal{H}}(D_{s_i}, D_t)$ of divergence between a set of source standard datasets $\{D_{s_i}\}$ and the target adversarial dataset D_t weighted by α_i upper-bounds the generalization error of the target adversarial data distribution \tilde{D}_t .*

The lemma extends an existing bound for federated domain adaptation [32], and shows that the generalization error on the unseen target noised distribution \tilde{D}_t is bounded by the α_i -weighted distribution gaps. Motivated by the analysis, we set α_i to be reversely proportional to the divergence between D_{s_i} and D_t . Specifically, we use a layer-averaged RBF-kernel to approximate the weight, i.e., $\alpha_i = \frac{1}{L} \sum_{l=1}^L \exp[-\gamma_{\text{rbf}} d_W^l(D_{s_i}, D_t)/p^l]$, where p^l is the number of channels in layer l . The distribution discrepancy can be approximately measured by Wasserstein distance as $d_W^l(D_s, D_t) = \|\mu_{s_i}^l - \mu_t^l\|_2 + \|\sigma_{s_i}^l - \sigma_t^l\|_2^2$. We use a large γ_{rbf} , i.e., $100 \times \max_l p^l$, to contrast the in-distribution and out-distribution difference. Lastly, we normalize α_i such that $\sum_i \alpha_i = 1$. The formal analysis can be found in Section B.1.

4.2 FedRBN algorithm and its efficiency

We are now ready to present the proposed the two-stage Federated Robust Batch-Normalization (FedRBN) algorithm, as summarized in Fig. 3. During training (Algorithm 1), we train models locally with decoupled clean and noise BNs for each user. Locally trained models excluding BN are then aggregated by federated parameter averaging. After training (Algorithm 2), we attentively copy BN parameters from multiple AT models into ST ones.

Algorithm 1 FedRBN: user training

Input: An initial model f from the server, adversary $A(\cdot)$, dataset D , adversarial hyper-parameter q

- 1: **for** mini-batch $\{(x, y)\}$ in D **do**
- 2: Set f to use clean BN by $h \leftarrow 0$
- 3: $L \leftarrow \mathbb{E}_{(x,y)}[\ell(f, (x, y))]$
- 4: **if** $q > 0$ **then**
- 5: Perturb data $\tilde{x} \leftarrow A(x)$
- 6: Set f to use noise BN by $h \leftarrow 1$
- 7: $L \leftarrow (1 - q)L + q\mathbb{E}_{(\tilde{x}, y)}[\ell(f, (\tilde{x}, y))]$
- 8: Update f by one-step gradient descent
- 9: **Upload** parameters of layers except BN layers

Algorithm 2 FedRBN: post-training

Input: Source AT users $\{s_i\}$, adversary $A(\cdot)$, local validation dataset D_{val}

- 1: Copy debiased noise-BN statistics attentively from users $\{s_i\}$ if current user is a ST user
- 2: $D_a \leftarrow \emptyset$
- 3: Set f to use clean BN by $h \leftarrow 0$
- 4: **for** mini-batch $\{(x, y)\}$ in D_{val} **do**
- 5: $\tilde{x} \leftarrow A(x)$
- 6: $D_a \leftarrow D_a \cup \{(f(x), 0), (f(\tilde{x}), 1)\}$
- 7: Fit a noise detector $g(\cdot)$ on D_a
- 8: **Return** noise detector $g(\cdot)$, modified f

Inference-stage BN selection. One issue of FedRBN is the choice of BN at inference time. To balance accuracy and robustness, we need a strategy to automatically select the right BN for each sample, i.e., applying noise BN when the input is found to be an adversarial one. Although the differences between clean data and adversarial are subtle in raw images, recent advances have shown promising results with the help of network structures. For example, in [33], the authors find that the non-maximal entropy of their logits are quite distinct. In [34], their representations are separable under the similarity measurement of RBF-kernel (or K-density). In Figs. 8 and 9, we visually show that adversarial examples can be identified in the space of logit vectors.

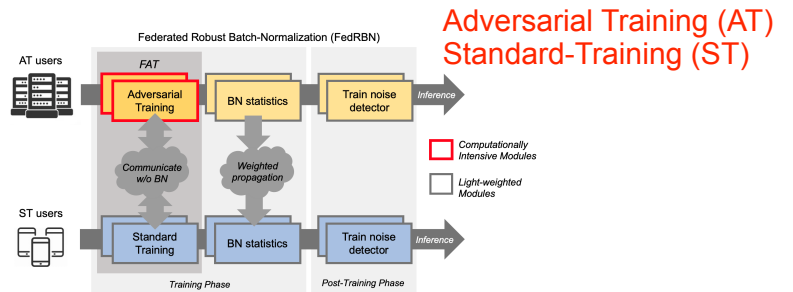


Figure 3: Overview of our FedRBN algorithm.

During training (Algorithm 1), we train models locally with decoupled clean and noise BNs for each user. Locally trained models excluding BN are then aggregated by federated parameter averaging. After training (Algorithm 2), we attentively copy BN parameters from multiple AT models into ST ones.

As described in [Algorithm 2](#), we propose to fit a support vector machine (SVM) [35], denoted as $g(\cdot)$, with RBF kernels on the validation set of each user. At inference, we predict h in [Eq. \(3\)](#) by $\hat{h} = g(f(x))$. The use of RBF kernels is partially inspired by [34], which used kernels on representations instead of logits.

We are aware of efforts that eliminate the detected adversarial samples from test/inference, e.g., [33]. However, simply refusing to predict suspicious samples may break up the downstream services, especially in the challenging scenarios when the detection becomes sensitive and many samples are suspicious. Instead, our method detects and predicts the suspicious samples using the robust model (with noise BN’s), and does not refuse any predictions.

BN operations. Since the BN statistics are only a small portion of any networks and do not require back-propagation, an additional BN statistic will not significantly affect the efficiency [36]. During training, because users do not send out BN layers, the communication cost is the same as a non-iid FL method (FedBN [6]) and less than other fully-shared methods like FedAvg [1]. After the training is done, BN layers will be copied to transfer robustness, and such one-time cost of transferring marginal components of networks will be neglectable, compared to the total cost incurred in FL.

The training and inference of noise detector. Introducing a noise detector results in an additional training cost, but such overhead is marginal in practice. First, the training is only done once after the network training is done. We only need to forward the whole network once for each sample to obtain the logit required for training the noise detector. Therefore, the overall overhead for receiving robustness is extremely efficient as compared to the AT overhead. Suppose adversarial samples are collected at real time which could be private. Training a noise detector only requires approximately $a/T \times 100\%$ of the training adversarial samples where T is the number of epochs and a is the ratio of validation set versus the training set.

5 Experiments

Datasets and models. To implement a non-iid scenario, we adopt a close-to-reality setting where users’ datasets are sampled from different distributions. We used two multi-domain datasets for the setting. The first is a subset (30%) of DIGITS, a benchmark for domain adaption [32]. DIGITS has 28×28 images and serves as a commonly used benchmark for FL [37, 1, 38]. DIGITS includes 5 different domains: MNIST (MM) [39], SVHN (SV) [40], USPS (US) [41], SynthDigits (SY) [42], and MNIST-M (MM) [42]. The second dataset is DOMAINNET [43] processed by [6], which contains 6 distinct domains of large-size 256×256 real-world images: Clipart (C), Infograph (I), Painting (P), Quickdraw (Q), Real (R), Sketch (S). For DIGITS, we use a convolutional network with BN (or DBN) layers following each conv or linear layers. For the large-sized DOMAINNET, we use AlexNet [44] extended with BN layers after each convolutional or linear layer [6].

Training and evaluation. For AT users, we use n -step PGD (projected gradient descent) attack [27] with a constant noise magnitude ϵ . Following [27], we use $\epsilon = 8/255$, $n = 7$, and attack inner-loop step size $2/255$, for training, validation, and test. We uniformly split the dataset for each domain into 10 subsets for DIGITS and 5 for DOMAINNET, following [6], which are distributed to different users, respectively. Accordingly, we have 50 users for DIGITS and 30 for DOMAINNET. Each user trains local model for one epoch per communication round. We evaluate the federated performance by standard accuracy (SA), classification accuracy on the clean test set, and robust accuracy (RA), classification accuracy on adversarial images perturbed from the original test set. All metric values are averaged over users.

We defer other details of experimental setup such as hyper-parameters to [Section C](#), and focus on discussing the results.

5.1 Comprehensive study

To further understand the role of each component in the proposed FedRBN, we conduct a comprehensive study on its properties. In experiments, we use three representative federated baselines combined with AT: FedAvg [1], FedProx [38], and FedBN [6]. We use FATAvg to denote the AT-augmented FedAvg, and similarly FATProx and FATBN. To implement hardware heterogeneity, we let 20%-per-domain users from 3/5 domains conduct AT.

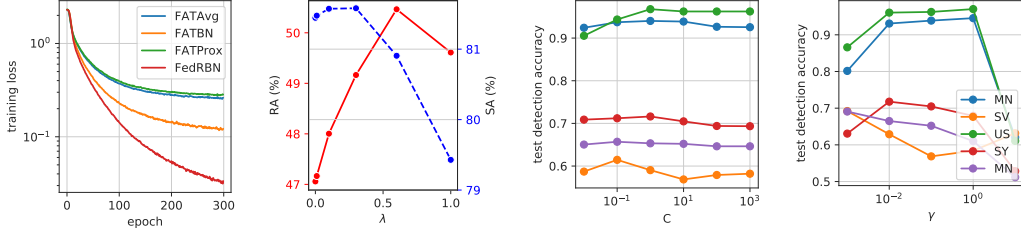


Figure 4: The convergence curves and parameters sensitivity of λ , C and γ . C is for regularization and γ is for RBF-kernel used in SVM whose performance is evaluated on Digits domains.

Table 1: Ablation of different FedRBN components. Standard deviations are enclosed in brackets.

AT users	metric	base	+DBN	+detector	+copy	+debias	+reweight
20% in 5/5 domains	RA	41.2 (1.3)	38.8 (1.4)	42.4 (1.4)	55.7 (1.0)	55.7 (1.3)	56.2 (1.0)
	increment		-2.4	+3.6	+13.3	+0.0	+0.5
	SA	86.4 (0.4)	86.5 (0.4)	86.2 (0.4)	85.2 (0.6)	85.2 (0.4)	85.3 (0.4)
100% in 1/5 domain	RA	36.5 (1.8)	34.3 (2.0)	37.3 (1.8)	48.1 (1.6)	49.8 (1.6)	49.8 (1.5)
	increment		-2.2	+3.0	+10.8	+1.7	+0.0
	SA	86.4 (0.4)	86.4 (0.4)	86.3 (0.4)	84.3 (0.5)	84.4 (0.4)	84.3 (0.5)

Convergence. The first plot in Fig. 4 shows convergence curves of different competing algorithms. Since FedRBN only differs from FATBN by a DBN structure, FATBN and FedRBN have similar convergence rates that are faster than others. We see that FedRBN converges even faster than FATBN. A possible reason is that DBN decouples the normal and adversarial samples, the representations after BN layers will be more consistently distributed among non-iid users.

Parameter Sensitivity of the λ , C and γ . The second plot in Fig. 4 shows a preferred λ is neither too small or too close to 1. Since λ is critical when heterogeneity is severer, we evaluate the sensitivity as only one domain (MNIST in DIGITS) is adversarially trained. We find that a larger λ is more helpful for the RA, as the estimation is closer to the true robust one. The rest plots in Fig. 4 demonstrate the stability of the noise detector when a choice of $C = 10$ and $\gamma = 1/10$ for SVM generally works.

Ablation Study. We use FATBN as the base method and incrementally add FedRBN components: +DBN, +detector (for adaptively BN selection), +copy BN statistics, +debias copying, and +reweight average before copying. Table 1 shows that even though DBN is a critical structure, simply adding DBN does not help unless the noise detector is applied. Also, the most influential component is copying, supporting our key idea. The +debias is more important in single AT domain case where domain gaps varies by different ST domains, whereas +reweight matters more when more AT domains are available. Other than +copy, all other components barely affect the SA.

Impacts from Data Heterogeneity. To study the influence of different AT domains, we set up an experiment where AT users only reside on one single domain. For simplicity, we let each domain contains a single user as in [6] and utilize only 10% of DIGITS dataset. The single AT domain plays the central role in gaining robustness from adversarial augmentation and propagating to other domains. The task is hardened by the non-singleton of gaps between the AT domain and multiple ST domains and a lack of the knowledge of domain relations. Results in Fig. 5a shows the superiority of the proposed FedRBN, which improves the RA more than 10% in all cases with small drops in SA. We see that the RA is the worst when MNIST serves the AT domain, whereas RA propagates better when the AT domain is SVHN or SynthDigits. The possible reason is that SVHN and SynthDigits are more visually different from the rest domains (see Fig. 7), forming larger domain gaps at test.

Impacts from Hardware Heterogeneity. We vary the number of AT users in training from $1/N$ (most heterogeneous) to N/N (homogeneous) to compare the robustness gain. Fig. 5b shows that our method consistently improves the robustness. Even when all domains are noised, FedRBN is the best due to the use of DBN. When not all domains are AT, our method only needs half of the users to be noised such that the RA is close to the upper bound (fully noised case).

5.2 Comparison to baselines

To demonstrate the effectiveness of the proposed FedRBN, we compare it with baselines on two benchmarks. We repeat each experiment for three times with different seeds. We introduce two more baselines: personalized meta-FL extended with FAT (FATMeta) [2] and federated robust training (FedRob) [10]. Because FedRob requires a project matrix of the squared size of image and the matrix is up to $256^2 \times 256^2$ on DOMAINNET which does not fit into a common GPU, we exclude it from comparison. Given the same setting, we constrain the computation cost in the similar scale

Table 2: Benchmarks of robustness propagation, where we measure the computation time (T) by counting $\times 10^{12}$ times of float add-or-multiplication operations (FLOPs).

AT users	Digits									DomainNet								
	All			20%			MNIST			All			20%			Real		
	RA	SA	T	RA	SA	T	RA	SA	T	RA	SA	T	RA	SA	T	RA	SA	T
FedRBN (ours)	66.7	87.3	2218	56.2	85.3	665	49.8	84.3	665	30.6	53.7	38490	24.2	61.5	11547	18.0	59.5	10425
FATBN	60.0	87.3	2211	41.2	86.4	663	36.5	86.4	663	29.9	52.6	38363	20.3	63.2	11509	11.3	60.5	10390
FATAvg	58.3	86.1	2211	42.6	84.6	663	38.4	84.1	663	24.6	47.4	38363	15.4	57.8	11509	9.4	57.4	10390
FATProx	58.5	86.3	2211	42.8	84.5	663	38.1	84.1	663	24.8	47.1	38363	14.5	57.3	11509	9.4	56.6	10390
FATMeta	43.6	71.6	2211	35.0	72.6	663	35.3	72.2	663	6.0	23.5	38363	0.0	37.2	11509	0.1	38.1	10390
FedRob	13.1	13.1	2211	20.6	59.3	1032	17.7	48.9	645	-	-	-	-	-	-	-	-	-

Table 3: Compare FedRBN versus efficient federated AT on Digits.

	20% 3/5 AT domains	100% Free AT [12]	
	FedRBN	FATAvg	FATBN
RA	56.1	44.9	47.1
SA	86.2	85.6	63.6
T	273	271	276

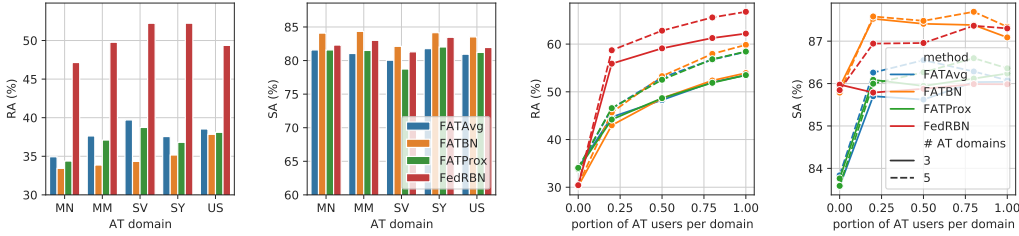
Table 4: Evaluation of RA with various attacks on Digits. n and ϵ are the step number and the magnitude of attack.

Attack (n, ϵ)	PGD (20,16)	PGD (100,8)	MIA [45] (20,16)	MIA (100,8)	LSA [46] (7, -)	SA
FedRBN	42.8	54.5	39.9	52.2	73.5	84.2
FATBN	28.6	41.6	27.0	39.7	64.0	84.6
FATAvg	31.5	43.4	30.0	41.5	63.3	84.2

for cost-fair comparison. We evaluate methods on two FRP settings. **1) Propagate from a single domain.** In reality, a powerful computation center may join the FL with many other users, e.g., mobile devices. Therefore, the computation center is an ideal node for the computation-intensive AT. Due to limitations of data collection, the center may only have access to a single domain, resulting gaps to most other users. We evaluate how well the robustness can be propagated from the center to others. **2) Propagate from a few multi-domain AT users.** In this case, we assume that to reduce the total training time, ST users are exempted from the AT tasks in each domain. Thus, an ST user wants to gain robustness from other same-domain users. The different-domain users may hinder the robustness performance due to the domain gaps in adversarial samples.

Table 2 shows that our method outperforms all baselines for all tasks, while it associates to only marginal overhead (for fitting noise detector). Importantly, we show that only 20% users are enough to achieve robustness comparable to the best fully-trained baseline. Even evaluated by different attacks (see Table 4), our method still outperforms others.

Compare to full efficient AT. In Table 3, we show that when computation time is comparable, our method can achieve both better RA and SA than full-AT baselines. For results to be comparable, we train FedRBN for limited 150 epochs while Free AT for 300 epochs. Although Free AT improves the robustness compared to FedAvg, it also greatly sacrifices SA performance. Thanks to stable convergence and decoupled BN, FedRBN maintains both accurate and robust performance though the AT is not ‘free’ for a few users.



(a) FRP from a single AT domain

(b) FRP from partial AT users per domain

Figure 5: Evaluating FRP performance with different FRP settings.

6 Conclusion

In this paper, we investigate a novel problem setting, federate propagating robustness, and propose a FedRBN algorithm that transfers robustness in FL through robust BN statistics. Extensive experiments demonstrate the rationality and effectiveness of the proposed method, delivering both generalization and robustness in FL. We believe such a client-wise efficient robust learning can broaden the application scenarios of FL to users with diverse computation capabilities.

Broader impact. The vulnerability of well-trained deep neural networks has drawn great attention in the research community. In the scope of federated learning, different users have different needs and capabilities for such a robust learning schema. For this sake, we propose a novel research task for propagating robustness among users in federated learning. Our methods can be applied in many high-stakes real-world scenarios, including self-driving [47], face recognition [48], natural language prediction [4], medical image analysis [49, 50] and computer-aided diagnosis [51].

References

- [1] Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Agüera y Arcas. Communication-Efficient Learning of Deep Networks from Decentralized Data. In *Artificial Intelligence and Statistics*, pages 1273–1282, April 2017.
- [2] Alireza Fallah, Aryan Mokhtari, and Asuman Ozdaglar. Personalized Federated Learning: A Meta-Learning Approach. In *Advances in Neural Information Processing Systems*, June 2020.
- [3] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-Free Knowledge Distillation for Heterogeneous Federated Learning. *arXiv:2105.10056 [cs]*, May 2021.
- [4] Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated Learning for Mobile Keyboard Prediction. *arXiv:1811.03604 [cs]*, February 2019.
- [5] Tao Yu, Eugene Bagdasaryan, and Vitaly Shmatikov. Salvaging Federated Learning by Local Adaptation. *arXiv:2002.04758 [cs, stat]*, February 2020.
- [6] Xiaoxiao Li, Meirui Jiang, Xiaofei Zhang, Michael Kamp, and Qi Dou. FedBN: Federated Learning on Non-IID Features via Local Batch Normalization. In *International Conference on Learning Representations*, September 2020.
- [7] Manoj Ghuhan Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated Learning with Personalization Layers. *arXiv:1912.00818 [cs, stat]*, December 2019.
- [8] Canh T. Dinh, Nguyen H. Tran, and Tuan Dung Nguyen. Personalized Federated Learning with Moreau Envelopes. In *Advances in Neural Information Processing Systems*, June 2020.
- [9] Giulio Zizzo, Amrith Rawat, Mathieu Sinn, and Beat Buesser. FAT: Federated Adversarial Training. *arXiv:2012.01791 [cs]*, December 2020.
- [10] Amirhossein Reisizadeh, Farzan Farnia, Ramtin Pedarsani, and Ali Jadbabaie. Robust Federated Learning: The Case of Affine Distribution Shifts. In *Advances in Neural Information Processing Systems*, June 2020.
- [11] Cihang Xie, Mingxing Tan, Boqing Gong, Jiang Wang, Alan Yuille, and Quoc V. Le. Adversarial Examples Improve Image Recognition. *arXiv:1911.09665 [cs]*, April 2020.
- [12] Ali Shafahi, Mahyar Najibi, Amin Ghiasi, Zheng Xu, John Dickerson, Christoph Studer, Larry S. Davis, Gavin Taylor, and Tom Goldstein. Adversarial Training for Free! *Advances in Neural Information Processing Systems*, 2019.
- [13] Dinghui Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You Only Propagate Once: Accelerating Adversarial Training via Maximal Principle. *Advances in Neural Information Processing Systems*, page 12, 2019.
- [14] Cihang Xie and Alan Yuille. Intriguing properties of adversarial training at scale. *International Conference on Learning Representations*, December 2019.
- [15] Steffen Schneider, Evgenia Rusak, Luisa Eck, Oliver Bringmann, Wieland Brendel, and Matthias Bethge. Improving robustness against common corruptions by covariate shift adaptation. *Advances in Neural Information Processing Systems*, October 2020.
- [16] Raouf Kerkouche, Gergely Ács, and Claude Castelluccia. Federated Learning in Adversarial Settings. *arXiv:2010.07808 [cs]*, October 2020.
- [17] Devansh Shah, Parijat Dube, Supriyo Chakraborty, and Ashish Verma. Adversarial training in communication constrained federated learning. *arXiv:2103.01319 [cs]*, March 2021.
- [18] Shuhao Fu, Chulin Xie, Bo Li, and Qifeng Chen. Attack-Resistant Federated Learning with Residual-based Reweighting. September 2019.

- [19] Peva Blanchard, El Mahdi El Mhamdi, Rachid Guerraoui, and Julien Stainer. Machine Learning with Adversaries: Byzantine Tolerant Gradient Descent. *Advances in Neural Information Processing Systems*, 30:119–129, 2017.
- [20] Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing Federated Learning through an Adversarial Lens. In *International Conference on Machine Learning*, November 2018.
- [21] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local Model Poisoning Attacks to Byzantine-Robust Federated Learning. In *29th {USENIX} Security Symposium ({USENIX} Security 20)*, pages 1605–1622, 2020.
- [22] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How To Backdoor Federated Learning. *International Conference on Machine Learning*, July 2018.
- [23] Yudong Chen, Lili Su, and Jiaming Xu. Distributed Statistical Machine Learning in Adversarial Settings: Byzantine Gradient Descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(2):44:1–44:25, December 2017.
- [24] Dong Yin, Yudong Chen, Ramchandran Kannan, and Peter Bartlett. Byzantine-Robust Distributed Learning: Towards Optimal Statistical Rates. In *International Conference on Machine Learning*, pages 5650–5659. PMLR, July 2018.
- [25] Krishna Pillutla, Sham M. Kakade, and Zaid Harchaoui. Robust Aggregation for Federated Learning. In *Advances in Neural Information Processing Systems*, 2020.
- [26] Zhaoxian Wu, Qing Ling, Tianyi Chen, and Georgios B. Giannakis. Federated Variance-Reduced Stochastic Gradient Descent With Robustness to Byzantine Attacks. *IEEE Transactions on Signal Processing*, 68: 4583–4596, 2020.
- [27] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards Deep Learning Models Resistant to Adversarial Attacks. *arXiv:1706.06083 [cs, stat]*, September 2019.
- [28] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness May Be at Odds with Accuracy. *ICLR*, September 2019.
- [29] Ali Shafahi, Parsa Saadatpanah, Chen Zhu, Amin Ghiasi, Christoph Studer, David Jacobs, and Tom Goldstein. Adversarially robust transfer learning. *International Conference on Learning Representations*, May 2019.
- [30] Dan Hendrycks, Kimin Lee, and Mantas Mazeika. Using Pre-Training Can Improve Model Robustness and Uncertainty. *arXiv:1901.09960 [cs, stat]*, October 2019.
- [31] Joaquin Quiñero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. *Dataset Shift in Machine Learning* | The MIT Press. MIT Press, 2008.
- [32] Xingchao Peng, Zijun Huang, Yizhe Zhu, and Kate Saenko. Federated Adversarial Domain Adaptation. In *International Conference on Learning Representations*, September 2019.
- [33] Tianyu Pang, Chao Du, Yinpeng Dong, and Jun Zhu. Towards Robust Detection of Adversarial Examples. *Advances in Neural Information Processing Systems*, November 2018.
- [34] Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner. Detecting Adversarial Samples from Artifacts. *arXiv:1703.00410 [cs, stat]*, November 2017.
- [35] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 20(3):273–297, September 1995.
- [36] Haotao Wang, Tianlong Chen, Shupeng Gui, Ting-Kuei Hu, Ji Liu, and Zhangyang Wang. Once-for-All Adversarial Training: In-Situ Tradeoff between Robustness and Accuracy for Free. *Advances in Neural Information Processing Systems*, November 2020.
- [37] Sebastian Caldas, Sai Meher Karthik Duddu, Peter Wu, Tian Li, Jakub Konečný, H. Brendan McMahan, Virginia Smith, and Ameet Talwalkar. LEAF: A Benchmark for Federated Settings. *arXiv:1812.01097 [cs, stat]*, December 2019.
- [38] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated Optimization in Heterogeneous Networks. *MLSys*, April 2020.

- [39] Y. Lecun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, November 1998.
- [40] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, 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.
- [41] J.J. Hull. A database for handwritten text recognition research. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 16(5):550–554, May 1994.
- [42] Yaroslav Ganin and Victor Lempitsky. Unsupervised Domain Adaptation by Backpropagation. In *International Conference on Machine Learning*, pages 1180–1189. PMLR, June 2015.
- [43] Xingchao Peng, Qinxun Bai, Xide Xia, Zijun Huang, Kate Saenko, and Bo Wang. Moment Matching for Multi-Source Domain Adaptation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1406–1415, 2019.
- [44] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In *Proceedings of the 25th International Conference on Neural Information Processing Systems - Volume 1, NIPS’12*, pages 1097–1105, Red Hook, NY, USA, December 2012. Curran Associates Inc.
- [45] Yinpeng Dong, Fangzhou Liao, Tianyu Pang, Hang Su, Jun Zhu, Xiaolin Hu, and Jianguo Li. Boosting Adversarial Attacks with Momentum. *arXiv:1710.06081 [cs, stat]*, March 2018.
- [46] Nina Narodytska and Shiva Prasad Kasiviswanathan. Simple Black-Box Adversarial Perturbations for Deep Networks. *arXiv:1612.06299 [cs, stat]*, December 2016.
- [47] Xinle Liang, Yang Liu, Tianjian Chen, Ming Liu, and Qiang Yang. Federated Transfer Reinforcement Learning for Autonomous Driving. *arXiv:1910.06001 [cs]*, October 2019.
- [48] Divyansh Aggarwal, Jiayu Zhou, and Anil K. Jain. FedFace: Collaborative Learning of Face Recognition Model. *arXiv:2104.03008 [cs]*, April 2021.
- [49] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *Medical Image Computing and Computer-Assisted Intervention – MICCAI 2015*, Lecture Notes in Computer Science, pages 234–241, Cham, 2015. Springer International Publishing.
- [50] Dianwen Ng, Xiang Lan, Melissa Min-Szu Yao, Wing P. Chan, and Mengling Feng. Federated learning: A collaborative effort to achieve better medical imaging models for individual sites that have small labelled datasets. *Quantitative Imaging in Medicine and Surgery*, 11(2):852–857, February 2021.
- [51] Awais Mansoor, Ulas Bagci, Brent Foster, Ziyue Xu, Georgios Z. Papadakis, Les R. Folio, Jayaram K. Udupa, and Daniel J. Mollura. Segmentation and Image Analysis of Abnormal Lungs at CT: Current Approaches, Challenges, and Future Trends. *Radiographics: A Review Publication of the Radiological Society of North America, Inc*, 35(4):1056–1076, 2015 Jul-Aug.
- [52] Eric Wong, Leslie Rice, and J. Zico Kolter. Fast is better than free: Revisiting adversarial training. In *International Conference on Learning Representations*, September 2019.
- [53] Alvin Chan, Yi Tay, and Yew-Soon Ong. What it Thinks is Important is Important: Robustness Transfers through Input Gradients. *CVPR*, March 2020.
- [54] Chuanbiao Song, Kun He, Liwei Wang, and John E Hopcroft. Improving the generalization of adversarial training with domain adaptation. In *International Conference on Learning Representations*, page 14, 2019.

Appendices

A Related Work

Efficient centralized adversarial training. A line of work has been motivated by similar concerns on the high time complexity of adversarial training. For example, Zhang *et al.* proposed to adversarially train only the first layer of a network which is shown to be more influential for robustness [13]. Free AT [12] trades in some standard iterations (on different mini-batches) for estimating a cached adversarial attack while keeping the total number of iterations unchanged. Wong *et al.* proposed to randomly initialize attacks multiple times, which can improve simpler attacks more efficiently [52]. Most of existing efforts above focus on speed up the local training by approximated attacks that trade in either RA or SA for efficiency. Instead, our method relocated the computation cost from budget-limited users to budget-sufficient users who can afford the expansive standard AT. As result, the computation expense is indeed exempted for the budget-limited users and their standard performance is not significantly influenced.

Robustness transferring. Our work is related to transferring robustness from multiple AT users to ST users. For example, a new user can enjoy the transferrable robustness of a large model trained on ImageNet [30]. In order to improve the transferability, some researchers aim to align the gradients between domains by criticizing their distributional difference [53]. A similar idea was adopted for aligning the logits of adversarial samples between different domains [54]. By fine-tuning a few layers of a network, Shafahi *et al.* shows that robustness can be transferred better than standard fine-tuning [29]. Rather than a central algorithm gathering all data or pre-trained models, our work considers a distributed setting when samples or their corresponding gradients can not be shared for distribution alignment. Meanwhile, a large transferrable model is not preferred in the setting, because of the huge communication cost associating to transferring models between users. Because of the non-iid nature of users, it is also hard to pick a proper user model, that works well on all source users, for fine-tuning on a target user.

B Method

The structure of FedRBN is depicted in Fig. 6. The DBN layer has the same input/output interface as an ordinary BN layer. Therefore, it can be plugged into most network structures whenever BN can be adopted. Except for DBN layers, FedRBN can be extended to other dual normalization layers, for instance, dual instance normalization [36].

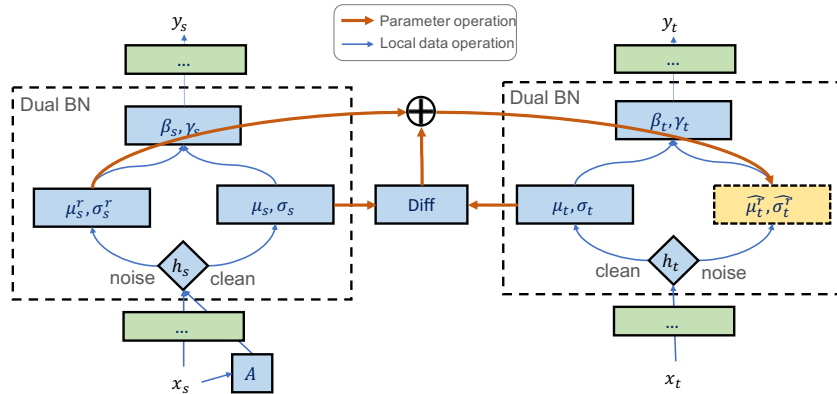


Figure 6: Illustration of the Dual-BN (DBN) layer and the copying operation for robustness propagation.

B.1 Proof of Lemma 4.1

In this section, we use the notation D for a dataset containing images and excluding labels. To provide supervisions, we define a ground-truth labeling function g that returns the true labels given images. So as for distribution \mathcal{D} .

First, in Definition B.1, we define the $\mathcal{H}\Delta\mathcal{H}$ -divergence that measures the discrepancy between two distributions. Because the $\mathcal{H}\Delta\mathcal{H}$ -divergence measures differences based on possible hypotheses (e.g., models), it can help relating model parameter differences and distribution shift.

Definition B.1. Given a hypothesis space \mathcal{H} for input space \mathcal{X} , the \mathcal{H} -divergence between two distributions \mathcal{D} and \mathcal{D}' is $d_{\mathcal{H}}(\mathcal{D}, \mathcal{D}') \triangleq 2 \sup_{S \in \mathcal{S}_{\mathcal{H}}} |\Pr_{\mathcal{D}}(S) - \Pr_{\mathcal{D}'}(S)|$ where $\mathcal{S}_{\mathcal{H}}$ denotes the collection of subsets of \mathcal{X} that are the support of some hypothesis in \mathcal{H} . The $\mathcal{H}\Delta\mathcal{H}$ -divergence is defined on the symmetric difference space $\mathcal{H}\Delta\mathcal{H} \triangleq \{f(x) \oplus h'(x) | h, h' \in \mathcal{H}\}$ where \oplus denotes the XOR operation.

Then, we introduce Assumption B.1 to bound the distribution differences caused by adversarial noise. The reason for introducing such an assumption is that the adversarial noise magnitude is bounded and the resulting adversarial distribution should not differ from the original one too much. Since all users are (or expected to be) noised by the same adversarial attacker $A_{\epsilon}(\cdot)$ during training, we can use d_{ϵ} to universally bound the adversarial distributional differences for all users.

Assumption B.1. Let d_{ϵ} be a non-negative constant governed by the adversarial magnitude ϵ . For a distribution \mathcal{D} , the divergence between \mathcal{D} and its corresponding adversarial distribution $\tilde{\mathcal{D}} \triangleq \{A_{\epsilon}(x) | x \sim \mathcal{D}\}$ is bounded as $d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}, \mathcal{D}) \leq d_{\epsilon}$.

Now, our goal is to analyze the *generalization error* of model \tilde{f}_t on the target adversarial distribution $\tilde{\mathcal{D}}$, i.e., $L(\tilde{f}_t, \tilde{\mathcal{D}}_t) = \mathbb{E}_{\tilde{x} \sim \tilde{\mathcal{D}}_t} [|\tilde{f}_t(\tilde{x}) - g(\tilde{x})|]$. Since we estimate \tilde{f}_t by a weighted average, i.e., $\sum_i \alpha_i \tilde{f}_{s_i}$ where \tilde{f}_{s_i} is the robust model on D_{s_i} , we can adapt the generalization error bound from [32] for adversarial distributions. For consistency, we assume the AT users reside on the *source* clean/noised domains while ST users reside on the *target* clean/noised domains. Without loss of generality, we only consider one target domain and assume one user per domain.

Theorem B.1 (Restated from Theorem 2 in [32]). Let \mathcal{H} be a hypothesis space of VC-dimension d and $\{\tilde{D}_{s_i}\}_{i=1}^N$, \tilde{D}_t be datasets induced by samples of size m drawn from $\{\tilde{\mathcal{D}}_{s_i}\}_{i=1}^N$ and $\tilde{\mathcal{D}}_t$, respectively. Define the estimated hypothesis as $\tilde{f}_t \triangleq \sum_{i=1}^N \alpha_i \tilde{f}_{s_i}$. Then, $\forall \alpha \in \mathbb{R}_+^N$, $\sum_{i=1}^N \alpha_i = 1$, with probability at least $1 - p$ over the choice of samples, for each $f \in \mathcal{H}$,

$$L(f, \tilde{\mathcal{D}}_t) \leq L(\tilde{f}_t, \tilde{D}_s) + \sum_{i=1}^N \alpha_i \left(\frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t) + \tilde{\xi}_i \right) + C, \quad (5)$$

where $C = 4\sqrt{\frac{2d \log(2Nm) + \log(4/p)}{Nm}}$, $\tilde{\xi}_i$ is the loss of the optimal hypothesis on the mixture of \tilde{D}_{s_i} and \tilde{D}_t , and \tilde{D}_s is the mixture of all source samples with size Nm . $d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t)$ denotes the divergence between domain s_i and t .

Based on Theorem B.1, we may choose a weighting strategy by $\alpha_i \propto 1/d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t)$. However, the divergence cannot be estimated due to the lack of the target adversarial distribution $\tilde{\mathcal{D}}_t$. Instead, we provide a bound by clean-distribution divergence in Lemma B.1.

Lemma B.1 (Formal statement of Lemma 4.1). Suppose Assumption B.1 holds. Let \mathcal{H} be a hypothesis space of VC-dimension d and $\{D_{s_i}\}_{i=1}^N$, D_t be datasets induced by samples of size m drawn from $\{\mathcal{D}_{s_i}\}_{i=1}^N$ and \mathcal{D}_t . Let an estimated target (robust) model be $\tilde{f}_t = \sum_i \alpha_i \tilde{f}_{s_i}$ where \tilde{f}_{s_i} is the robust model trained on D_{s_i} . Let \tilde{D}_s be the mixture of source samples from $\{\tilde{\mathcal{D}}_{s_i}\}_{i=1}^N$. Then, $\forall \alpha \in \mathbb{R}_+^N$, $\sum_{i=1}^N \alpha_i = 1$, with probability at least $1 - p$ over the choice of samples, for each $f \in \mathcal{H}$, the following inequality holds:

$$L(f, \tilde{\mathcal{D}}_t) \leq L(\tilde{f}_t, \tilde{D}_s) + d_{\epsilon} + \sum_{i=1}^N \alpha_i \left(\frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t) + \xi_i \right) + C,$$

where C and ξ_i are defined in Theorem B.1. D_s is the mixture of all source samples with size Nm . $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$ is the divergence over clean distributions.

Proof. Notice that Eq. (5) is a loose bound as $d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t)$ is neither bounded nor predictable. Differently, $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$ can be estimated by clean samples which is available for all users. Thus, we can bound $d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t)$ with $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$. By Assumption B.1, it is easy to attain

$$\begin{aligned} d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \tilde{\mathcal{D}}_t) &\leq d_{\mathcal{H}\Delta\mathcal{H}}(\tilde{\mathcal{D}}_{s_i}, \mathcal{D}_{s_i}) + d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t) + d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_t, \tilde{\mathcal{D}}_t) \\ &\leq 2d_\epsilon + d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t), \end{aligned} \quad (6)$$

where we used the triangle inequality in the space measured by $d_{\mathcal{H}\Delta\mathcal{H}}(\cdot, \cdot)$. Substitute Eq. (6) into Eq. (5), and we finish the proof. \square

In Lemma B.1, we discussed the bound for a $f \in \mathcal{H}$ (which also generalize to \tilde{f}_t) estimated by the linear combination of $\{\tilde{f}_{s_i}\}_i$. In our algorithm, \tilde{f}_t and \tilde{f}_{s_i} both represent the models with noise BN layers, and they only differ by the BN layers. Therefore, Lemma B.1 guides us to re-weight BN parameters according to the domain differences. Specifically, we should upweight BN statistics from user s_i if $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$ is large, vice versa. Since $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$ is hard to estimate, we may use the divergence over empirical distributions, i.e., $d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{s_i}, \mathcal{D}_t)$ instead.

C Experiments

In this section, we provide more details about our experiments and additional evaluation results.

C.1 Experiment details

Data. By default, we use 30% data of Digits for training. Datasets for all domains are truncated to the same size following the minimal one. In addition, we leave out 50% (60% for DomainNet) of the training set for validation for Digits. Test sets are preset according to the benchmarks in [6]. Models are selected according to the validation accuracy. To be efficient, we validate robust users with RA while non-robust users with SA. We use a large ratio of the training set for validation, because the very limited sample size for each user will result in biased validation accuracy. When selecting a subset of domains for AT users, we select the first n domains by the order: (MN, SV, US, SY, MM) for DIGITS, and (R, C, I, P, Q, S) for DOMAINNET. Some samples are plotted in Fig. 7 to show the visual difference between domains.

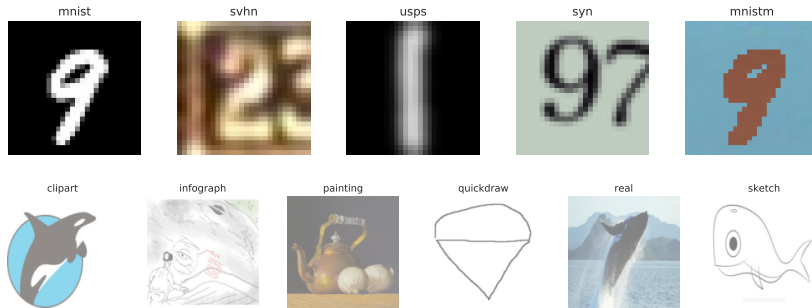


Figure 7: Visualization of samples.

Hyper-parameters. We use a fixed parameters tuned by the DIGITS dataset and adopt it also for the DOMAINNET dataset: $\lambda = 0.1$, $C = 10$ and $\gamma = 1/10$. The same tuning strategy is applied for other baseline parameters.

Network architectures for DIGITS and DOMAINNET are listed in Tables 5 and 6. For the convolutional layer (Conv2D or Conv1D), the first argument is the number channel. For a fully connected layer (FC), we list the number of hidden units as the first argument.

Training. Following [6], we conduct federated learning with 1 local epoch and batch size 32, which means users will train multiple iterations and communicate less frequently. Input images are resized to 256×256 for DOMAINNET and 28×28 for DIGITS. SGD (Stochastic Gradient Descent) is utilized to optimize models locally with a constant learning rate 10^{-2} . Models are trained for 300

Table 5: Network architecture for Digits dataset.

Layer	Details
feature extractor	
conv1	Conv2D(64, kernel size=5, stride=1, padding=2)
bn1	DBN2D, RELU, MaxPool2D(kernel size=2, stride=2)
conv2	Conv2D(64, kernel size=5, stride=1, padding=2)
bn2	DBN2D, ReLU, MaxPool2D(kernel size=2, stride=2)
conv3	Conv2D(128, kernel size=5, stride=1, padding=2)
bn3	DBN2D, ReLU
classifier	
fc1	FC(2048)
bn4	DBN2D, ReLU
fc2	FC(512)
bn5	DBN1D, ReLU
fc3	FC(10)

Table 6: Network architecture for DomainNet dataset.

Layer	Details
feature extractor	
conv1	Conv2D(64, kernel size=11, stride=4, padding=2)
bn1	DBN2D, ReLU, MaxPool2d(kernel size=3, stride=2)
conv2	Conv2D(192, kernel size=5, stride=1, padding=2)
bn2	DBN2D, ReLU, MaxPool2d(kernel size=3, stride=2)
conv3	Conv2D(384, kernel size=3, stride=1, padding=1)
bn3	DBN2D, ReLU
conv4	Conv2D(256, kernel size=3, stride=1, padding=1)
bn4	DBN2D, ReLU
conv5	Conv2D(256, kernel size=3, stride=1, padding=1)
bn5	DBN2D, ReLU, MaxPool2d(kernel size=3, stride=2)
avgpool	AdaptiveAvgPool2d(6, 6)
classifier	
fc1	FC(4096)
bn6	DBN1D, ReLU
fc2	FC(4096)
bn7	DBN1D, ReLU
fc3	FC(10)

epochs by default. For FedMeta, we use the 0.001 learning rate for the meta-gradient descent and 0.02 for normal gradient descent following the published codes from [8]. We fine-tune the parameters for DOMAINNET such that the model can converge fast. FedMeta converges slower than other methods, as it uses half of the batches to do the one-step meta-adaptation. We do not let FedMeta fully converge since we have to limit the total FLOPs for a fair comparison. FedRob fails to converge because locally estimated affine mapping is less stable with the large distribution discrepancy.

We implement our algorithm and baselines by PyTorch. The FLOPs are computed by thop package in which the FLOPs of common network layers are predefined¹. Then we compute the times of forwarding (inference) and backward (gradient computing) in training. Accordingly, we compute the total FLOPs of the algorithm. Because most other computation costs are relatively minor compared to the network forward/backward, these costs are ignored in our reported results.

C.2 Logits of adversarial and clean examples

Visualization of logits by t-SNE is presented in Fig. 9 (DIGITS) and Fig. 8 (DOMAINNET). We generally observe that the clean and adversarial logits are separable with generalizable decision

¹Retrieve the thop python package from <https://github.com/Lyken17/pytorch-OpCounter>.

boundaries. Moreover, MNIST and USPS domains turn out to be the most separable cases as they are easier classification tasks compared to the rest ones. Though some domains have a few mixed samples in the visualization, their noise detection accuracies are mostly higher than 90%. Thus, it is rational to fit a noise detector on the validation set for helping BN selection at test time.

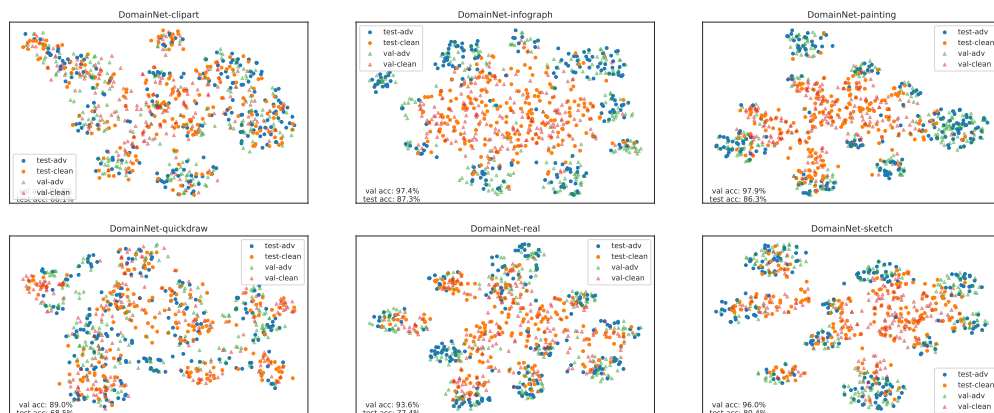


Figure 8: Logits of standard-trained models visualized by t-SNE on DOMAINNET.

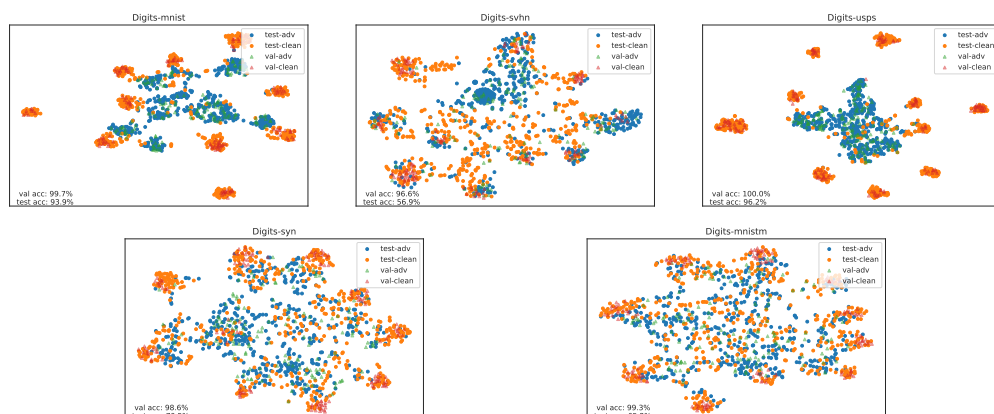


Figure 9: Logits of standard-trained models visualized by t-SNE on DIGITS.

C.3 Results on the Office-Caltech10 Dataset

Following the same setting as DOMAINNET experiments, we extend our experiments to a smaller dataset, Office-Caltech10 dataset preprocessed by [6] with images acquired by different cameras. The dataset includes 4 domains: Amazon, Caltech, DSLR, Webcam. Because the dataset has very few samples, we only generate 2 users per domain such that each user has at least 100 samples. In Table 7, we see that our method outperforms baselines as only one domain is adversarially trained. As the training set is rather small, the RAs are generally worse than the ones on DIGITS or DOMAINNET.

Table 7: Comparison to baselines on the Office-Caltech10 dataset. Standard deviations are reported in brackets.

AT users metric	Amazon		All	
	RA	SA	RA	SA
FedRBN (ours)	9.2 (3.4)	62.9 (3.4)	29.1 (2.4)	68.7 (1.7)
FedBN	5.1 (1.1)	65.9 (2.4)	30.8 (2.5)	67.2 (2.1)
FedAvg	0.6 (0.5)	54.7 (3.8)	13.3 (2.3)	56.0 (2.6)
FedProx	0.6 (0.6)	55.3 (4.7)	13.6 (1.8)	56.2 (2.1)

C.4 Evaluation with different FL configurations

In this section, we evaluate our method against FedBN with different federated configurations of local epochs E and batch size B . We constrain the parameters by $E \in \{1, 4, 8\}$ and $B \in \{10, 50, 100\}$. The 20% 3/5 domain FRP setting is adopted with DIGITS dataset. In Table 8, we report the average performance over 3 repetitions. The competition results are consistent that our method significantly promotes robustness over FedBN. We also observe that both our method and FedBN prefer a smaller batch size and fewer local epochs for better RA and SA. In addition, our method drops less RA when E is large or batch size increases.

Table 8: Evaluation with different FL configurations

B	E	method	RA	SA
10	1	FedBN	50.9	83.9
		FedRBN	60.0	82.8
10	4	FedBN	42.0	75.8
		FedRBN	56.3	76.1
10	8	FedBN	30.9	63.1
		FedRBN	53.4	68.4
50	1	FedBN	37.0	85.8
		FedRBN	53.2	84.5
100	1	FedBN	35.7	85.3
		FedRBN	53.0	83.8

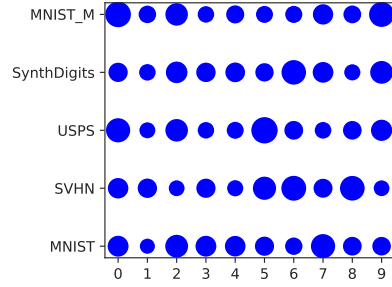


Figure 10: Dataset sizes for users when the global seed is set as 1. Larger circles indicate more training samples. The x-axis represents the user index.

C.5 Evaluation with unequal dataset sizes

In Section 5, we sub-sample the same number of data for each user, which may not be realistic. Therefore, we extend the experiment setting such that the sample sizes for users are not equal. We assume a user samples a variable ratio of data, which follows a Dirichlet distribution. We plot the different sample sizes for users in Fig. 10. Due to the varying dataset sizes, we let each user run a fixed number of iterations which is calculated by the average number of the per-epoch iterations of all users. In Table 9, we summarize the 3-repetition-averaged comparison results on the 20% 3/5 domain FRP setting on the DIGITS dataset. We see that our method is still most competitive with non-uniform dataset sizes.

Table 9: Comparison with unequal user-dataset sizes.

	RA	SA
FedRBN (ours)	53.1	84.4
FedBN	37.3	85.7
FedAvg	39.6	83.4
FedProx	39.5	83.4

Table 10: Comparison to robustness transferring by fine-tuning (FT).

	# FT iterations	RA	SA
FedRBN	-	53.1	84.4
FedAvg	-	44.7	85.7
FedAvg+FT	20	40.6	79.6
FedAvg+FT	100	40.6	83.4
FedAvg+FT	200	39.2	83.6

C.6 Comparison to robustness transferring by fine-tuning

As an alternative to FRP, fine-tuning (FT) the federated-trained models on target users can enjoy even better efficiency than FedRBN. Here, we first train AT users by FedAvg for 300 epochs. Note that we do not adopt FedBN because FedBN will not output a single model for adapting to new users. Then, the model is used for initializing the models for ST users. These ST users will be trained by FedAvg for a given number of FT iterations. Still, we adopt the 20% 3/5 domain FRP setting on the DIGITS dataset. In Table 10, we see that such a fine-tuning does not improve the robustness (RA).

C.7 Experiments in Fig. 1b

Though the results in Fig. 1b have been reported in previous experiments, we re-summarize the results in Table 11 for ease of reading. The basic setting follows the previous experiments on the Digits dataset. We construct different portions of AT users by *in-domain* or *out-domain* propagation settings. When robustness is propagated in domains, we sample AT users in each domain by the same portion and leave the rest as ST users. When robustness is propagated out of domains, all users from the last two domains will not be adversarially trained and gain robustness from other domains. Concretely, we add the FedRBN without copy propagation (FedRBN w/o prop) in the table, to show the propagation effect. FedRBN w/o prop outperforms the baselines only when the AT-user portion is more than 60%. Meanwhile, due to the lack of copy propagation, the RA is much worse than the propagated FedRBN. Unless no AT user presents in the federated learning, FedRBN always outperforms baselines.

Table 11: Results and detailed configurations of Fig. 1b on the 5-domain Digits dataset. FedAvg and FedBN corresponds to FATAvg and FATBN in the figure.

AT-user ratio	propagation	method	RA	SA	# AT domain	per-domain AT ratio
0%	none	FedRBN (ours)	32.1	84.3	0	0.0
		FedRBN w/o prop	32.1	84.3	0	0.0
		FedAvg	35.3	82.0	0	0.0
		FedBN	32.1	84.3	0	0.0
12%	mix	FedRBN (ours)	55.1	84.6	3	20%
		FedRBN w/o prop	38.7	84.5	3	20%
		FedAvg	44.1	84.1	3	20%
		FedBN	42.4	86.0	3	20%
20%	in-domain	FedRBN (ours)	57.3	85.3	5	20%
		FedRBN w/o prop	46.1	86.1	5	20%
		FedAvg	45.9	84.7	5	20%
		FedBN	44.8	86.0	5	20%
60%	out-domain	FedRBN (ours)	61.6	85.0	3	100%
		FedRBN w/o prop	56.2	85.5	3	100%
		FedAvg	52.0	84.2	3	100%
		FedBN	53.0	85.5	3	100%
100%	none	FedRBN (ours)	65.7	85.9	5	100%
		FedRBN w/o prop	65.8	85.9	5	100%
		FedAvg	57.5	84.7	5	100%
		FedBN	59.1	85.9	5	100%

C.8 Impact of data size and validation ratio

To investigate the impact of data size, we conduct experiments with varying training dataset sizes and validation ratios. Experiments follow previous protocols on the Digits dataset. Following the training/testing split in [6], we first sample a percentage of data for training. From the training set, we randomly select a subset for validation. We denote the two subsampling ratios as training percentage and validation ratio, respectively. When varying the training percentage, we fix the validation ratio at 10%. When varying the validation ratio, we use 30% training data. As shown in Fig. 11a, more training samples can improve the robustness and our method outperforms baselines consistently. In Fig. 11b, the ratio of validation set is less influential for the robustness performance of our FedRBN, but a larger validation ratio can reduce the time complexity of training as less samples are used for gradient computation. Though baseline methods obtain higher robust accuracies with smaller validation ratios, our method still introduces large gains in all cases.

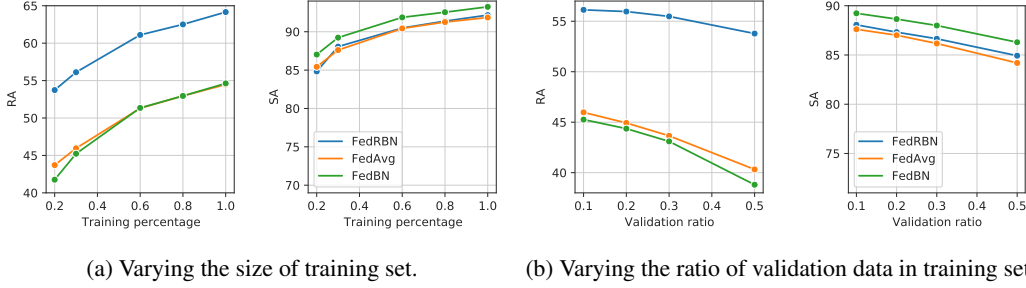


Figure 11: Experiments with varying data size.