

# CERTIFIABLY-ROBUST FEDERATED ADVERSARIAL LEARNING VIA RANDOMIZED SMOOTHING

**Cheng Chen**

Department of ECE  
University of Utah  
u0952128@utah.edu

**Bhavya Kailkhura**

Lawrence Livermore National Laboratory  
Livermore, US  
kailkhura1@llnl.gov

**Ryan Goldhahn**

Lawrence Livermore National Laboratory  
Livermore, US  
goldhahn1@llnl.gov

**Yi Zhou**

Department of ECE  
University of Utah  
yi.zhou@utah.edu

## ABSTRACT

Federated learning is an emerging data-private distributed learning framework, which, however, is vulnerable to adversarial attacks. Although several heuristic defenses are proposed to enhance the robustness of federated learning, they do not provide certifiable robustness guarantees. In this paper, we incorporate randomized smoothing techniques into federated adversarial training to enable data-private distributed learning with certifiable robustness to test-time adversarial perturbations. Our experiments show that such an advanced federated adversarial learning framework can deliver models as robust as those trained by the centralized training. Further, this enables provably-robust classifiers to  $\ell_2$ -bounded adversarial perturbations in a distributed setup. We also show that one-point gradient estimation based training approach is  $2 - 3\times$  faster than popular stochastic estimator based approach without any noticeable certified robustness differences.

## 1 INTRODUCTION

Federated learning is an emerging distributed learning framework that enables edge computing at a large scale (Konečný et al., 2016; Li et al., 2020b; McMahan et al., 2017; Chen et al., 2020), and has been successfully applied to various areas such as Internet of Things (IoT), autonomous driving, health care (Li et al., 2020b), etc. In particular, federated learning aims to exploit the distributed computation and heterogeneous data of a large number of edge devices to perform distributed learning while preserving full data privacy. The original federated learning framework proposed the federated averaging (FedAvg) algorithm (McMahan et al., 2017). In each learning round, a subset of edge devices are selected to download a global model from the cloud server, based on which the selected devices train their local models using local data for multiple stochastic gradient descent (SGD) iterations. Then, these devices upload the trained local models to the server, where the local models are aggregated and averaged to obtain an updated global model that will be used in the next learning round. Throughout the federated learning process, all data are kept privately on the local devices.

However, as modern federated learning often adopts over-parameterized models (e.g., deep neural networks) that have been proven to be vulnerable to adversarial perturbations to the test data (Szegedy et al., 2014; Goodfellow et al., 2015; Bulusu et al., 2020), there is a rising concern about the adversarial robustness of the federated learning models used by massive number of edge devices. As an example, if a federated-trained model is vulnerable to adversarial examples, then its performance on edge devices solving safety-critical tasks can be significantly degraded in turn having serious consequences. To defend such adversarial attacks in federated learning, many studies propose to include standard adversarial training in the local training steps of federated learning (Zhou et al., 2021; Zizzo et al., 2020; Kerkouche et al., 2020; Bhagoji et al., 2019). However, these approaches may not be able to defend strong adversaries and do not have certifiable adversarial robustness guar-

antee. To address these issues, some studies proposed the randomized smoothing technique that can train certifiably robust models at scale (Lecuyer et al., 2019; Cohen et al., 2019; Li et al., 2020a).

Specifically, randomized smoothing procedure uses a smoothed version of the original classifier  $f$  and certifies the adversarial robustness of the new classifier. The smoothed classifier is defined as  $g(x) = \arg \max_c \mathbb{P}_{\delta \sim \mathcal{N}(0, \sigma^2 I)}(f(x + \delta) = c)$ , meaning the label of a data sample  $x$  corresponds to the class whose decision region  $\{x' \in \mathbb{R}^d : f(x') = c\}$  has the largest measure under the distribution  $\mathcal{N}(x, \sigma^2 I)$ , where  $\sigma$  is used for smoothing. Suppose that while classifying a point  $\mathcal{N}(x, \sigma^2 I)$ , the original classifier  $f$  returns the class  $c_A$  with probability  $p_A = \mathbb{P}(f(x + \delta) = c_A)$ , and the “runner-up” class  $c_B$  is returned with probability  $p_B = \max_{c \neq c_A} \mathbb{P}(f(x + \delta) = c)$ , then the prediction of the point  $x$  under the smoothed classifier  $g$  is robust within the radius  $r(g; \sigma) = \frac{\sigma}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$ , where  $\Phi^{-1}$  is the inverse CDF of the standard Normal distribution. In practice, Monte Carlo sampling is used to estimate a lower bound on  $p_A$  and an upper bound on  $p_B$  as its difficult to estimate the actual values for  $p_A$  and  $p_B$ . Since standard training of the base classifier does not achieve high robustness guarantees, (Cohen et al., 2019) proposed to use Gaussian data augmentation based training in which the base classifier is trained on Gaussian noise corruptions of the clean data. Recently, the authors in (Salman et al., 2019) combined adversarial training approach with randomized smoothing to obtain significantly improved certification guarantees.

Such a smoothed model has been shown to outperform other existing certifiably robust models (Cohen et al., 2019) and the randomized smoothing scheme is applicable to deep networks and large datasets. To further enhance certifiable robustness of deep models, (Salman et al., 2019) proposed an adversarial training approach that uses strong attacks generated against the smoothed model to train the smoothed model. In particular, (Salman et al., 2019) demonstrated that such an adversarial training approach can substantially improve the robustness of smoothed models. However, these certifiably-robust training approaches are only applied to centralized learning setup, and similar provably-robust approaches in a federated learning setup is virtually non-existent. To bridge this gap, in this paper, we incorporate the randomized smoothing (with adversarial training) approach into the paradigm of federated learning to develop certifiably robust federated learning models.

**Our contributions.** We apply the randomized smoothing (with adversarial training) approach to enable the certifiable robustness of federated learning to adversarial perturbations. Specifically, in the local training phase, each device applies adversarial training to train a robust smoothed local model to defend  $\ell_2$  adversarial attacks. These local models are further aggregated by the central server to obtain a robust global model. To the best of our knowledge, this is the first work in the direction of enabling certifiable robust federated learning.

## 2 FEDERATED ADVERSARIAL LEARNING WITH RANDOMIZED SMOOTHING

### 2.1 ADVERSARIAL LEARNING WITH RANDOMIZED SMOOTHING

Consider a standard soft classifier  $F_\theta$  that is parameterized by  $\theta$  and maps an input data  $x \in \mathbb{R}^d$  to a probability mass of class labels  $\mathcal{Y}$ . Then, its corresponding smoothed soft classifier  $G_\theta$  is defined as

$$G_\theta(x) := \mathbb{E}_{\delta \sim \mathcal{N}(0, \sigma^2 I)}[F_\theta(x + \delta)]. \quad (1)$$

Intuitively, the smoothed classifier  $G_\theta$  perturbs the input sample with Gaussian noises and averages the predicted class distributions of all corrupted samples. In particular, the standard deviation  $\sigma$  of the Gaussian noise controls the level of certifiable robustness of the smoothed classifier.

To improve the performance, in (Salman et al., 2019), the authors proposed to leverage adversarial examples of the input data against the smoothed classifier  $G_\theta$  (instead of  $F_\theta$ ). Specifically, (Salman et al., 2019) proposed the following adversarial training problem, where the training uses the adversarial data  $\hat{x}$  that is found within an  $\ell_2$  ball of the original data  $x$  by attacking  $G_\theta$ .

$$\mathbf{SmoothAdv} : \min_{\theta} \max_{\|\hat{x}-x\|_2 \leq \epsilon} J_\theta(\hat{x}) := -\log[G_\theta(\hat{x})]_y, \quad (2)$$

where  $[G_\theta(\hat{x})]_y$  denotes the  $y$ -th entry of the predicted classification probability mass. This approach is referred to as **SmoothAdv** and the objective function is highly stochastic and non-convex. To solve the above adversarial optimization problem, two approaches were proposed in (Salman et al., 2019). For the first approach, the authors approximate the gradient of the above objective function using

stochastic samples as follows

$$\text{(Stochastic estimator)} \quad \nabla_x J(\hat{x}) \approx -\nabla_x \log \left( \frac{1}{m} \sum_{i=1}^m [F_\theta(\hat{x} + \delta_i)]_y \right), \quad (3)$$

where  $\delta_i, i = 1, \dots, m$  are drawn i.i.d from  $\mathcal{N}(0, \sigma^2 I)$ . Then, standard projected gradient ascent is applied to find adversarial samples. While the above stochastic gradient estimator provides an accurate gradient estimation, it is computational expensive as for every sample  $x$  we need to perform back-propagation on a mini-batch of  $m$  corrupted samples.

To avoid performing back-propagation, (Salman et al., 2019) discussed another gradient-free (Liu et al., 2020) approach. Specifically, note that the adversarial optimization problem is equivalent to  $\hat{x} = \arg \min_{\|\hat{x}-x\|_2 \leq \epsilon} [G_\theta(\hat{x})]_y$ . In particular, the gradient of  $[G_\theta(\hat{x})]_y$  can be conveniently characterized using the following one-point gradient-free estimator.

$$\text{(One-point estimator)} \quad \nabla_x [G_\theta(\hat{x})]_y \approx \frac{1}{m} \sum_{i=1}^m \left[ \frac{\delta_i}{\sigma^2} \cdot [F_\theta(\hat{x} + \delta_i)]_y \right]. \quad (4)$$

The above estimator only involves function values that can be efficiently computed via forward-propagation. In particular, each gradient estimate  $\frac{\delta_i}{\sigma^2} \cdot [F_\theta(\hat{x} + \delta_i)]_y$  only needs to evaluate the function value at a single point  $\hat{x} + \delta_i$ . Compared to the gradient-based stochastic estimator, this one-point estimator is computation lighter but induces a higher estimation variance. In (Salman et al., 2019), the performance of the one-point estimator was not evaluated for **SmoothAdv**, and its comparison with the stochastic estimator was not comprehensive.

## 2.2 FEDERATED ADVERSARIAL LEARNING

In this section, we incorporate the **SmoothAdv** method into the federated learning framework. Our proposed algorithm is referred to as **Fed-SmoothAdv** is presented in Algorithm 1.

---

### Algorithm 1: Federated Adversarial Learning (**Fed-SmoothAdv**)

---

**Central-server executes:** # *Run on the central server*

**for** learning round  $t = 1, 2, \dots$  **do**  
  Sample a subset  $S_t$  of clients  
  **for** each client  $k \in S_t$  in parallel **do**  
     $\theta_{t+1}^k \leftarrow \text{LocalTrain}(k, \theta_t)$   
    Send  $\theta_{t+1}^k$  to the server  
  Server aggregates  $\theta_{t+1} \leftarrow \sum_{k \in S_t} \frac{n_k}{n} \theta_{t+1}^k$

---

**LocalTrain** ( $k, \theta$ ): # *Local training of client k*

**for** local iteration  $i = 1, 2, \dots, E$  **do**  
  Sample a minibatch of data  $b$   
   $\theta \leftarrow \text{SmoothAdv}(\theta, b)$  # *Use one of the two gradient estimators*

---

**SmoothAdv** ( $\theta, b$ ): # *Adversarial training with randomized smoothing*

Data samples  $(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(b)}, y^{(b)})$

Generate noises  $\{\delta_i^{(j)}\}_{i=1}^m \sim \mathcal{N}(0, \sigma^2 I)$  for any  $x^{(j)}, j = 1, \dots, b$

$L \leftarrow []$  # *List of adversarial examples*

**for**  $1 \leq j \leq b$  **do**  
  Generate adversarial sample  $\hat{x}^{(j)}$  for  $x^{(j)}$  by attacking the smoothed classifier using one of the gradient estimators in eqs.(3,4) and noises  $\{\delta_i^{(j)}\}_{i=1}^m$ .  
  Append  $\{(\hat{x}^{(j)} + \delta_1^{(j)}, y^{(j)}), \dots, (\hat{x}^{(j)} + \delta_m^{(j)}, y^{(j)})\}$  to list  $L$ .

---

Train model  $\theta$  using adversarial samples in  $L$  for multiple SGD steps.

---

To elaborate, the hierarchical structure of **Fed-SmoothAdv** is the same as that of standard federated learning, i.e., a subset of edge devices is sampled in every round to perform local training, and then their local models are aggregated by the cloud server. However, in our federated adversarial learning, each client uses **SmoothAdv** to perform local adversarial training using strong adversarial samples generated by attacking the smoothed local model.

### 3 EXPERIMENTS

We compare the certified robustness of **Fed-SmoothAdv** with the baseline method **SmoothAdv** in training an AlexNet (Krizhevsky et al., 2012) on CIFAR-10 (Krizhevsky, 2009). Here, certified robustness is defined as the fraction of the test samples that are correctly classified (without abstaining) by  $G_\theta$  and are certified within an  $\ell_2$  radius of  $r$ . We set the smoothing parameter  $\sigma = \{0.12, 0.25, 0.5\}$  and the perturbation bound  $\epsilon = \{64, 128, 256\}$ , and use the same  $\sigma$  for certification as that used in the training. For both methods, we apply both the stochastic estimator and the one-point estimator. Moreover, we test **Fed-SmoothAdv** under different levels of device data heterogeneity  $\gamma_{\text{device}}$  (the higher the more heterogeneous). Please refer to Appendix A for all the other hyperparameters used in the experiments.

In Figure 1, we plot the certified accuracy of both **SmoothAdv** and **Fed-SmoothAdv** (with heterogeneity  $\gamma_{\text{device}} = 0.1, 0.5$ ) with  $\sigma = 0.25, \epsilon = 128$ . It can be seen that the certified accuracy of **Fed-SmoothAdv** is slightly lower than that of **SmoothAdv**, but is reasonably close. Also, the data heterogeneity  $\gamma_{\text{device}}$  does not affect the certified accuracy of **Fed-SmoothAdv**, which implies that **SmoothAdv** can be effectively applied to enhance the adversarial robustness of heterogeneous federated learning. Moreover, we note that while the performance of the one-point estimator is almost the same as that of the stochastic estimator, the training time is significantly reduced by 2-3 times due to avoidance of backpropagation. All these results show that applying **SmoothAdv** with the one-point estimator to federated learning can efficiently enhance the certified model accuracy. In Figure 2, we plot the certified accuracy results under  $\sigma = 0.5$  and  $\epsilon = 128$ . One can observe a similar comparison between the two methods as that in Figure 1. In particular, with a larger  $\sigma$ , the certified accuracy is lower but spans over a wider range of  $\ell_2$  radius. The results corresponding to other choices of  $\sigma, \epsilon$  can be found in Appendix C.

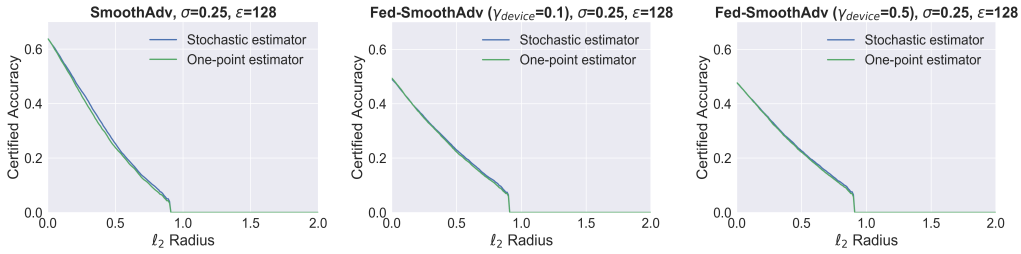


Figure 1: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.25$  and  $\epsilon = 128$ .

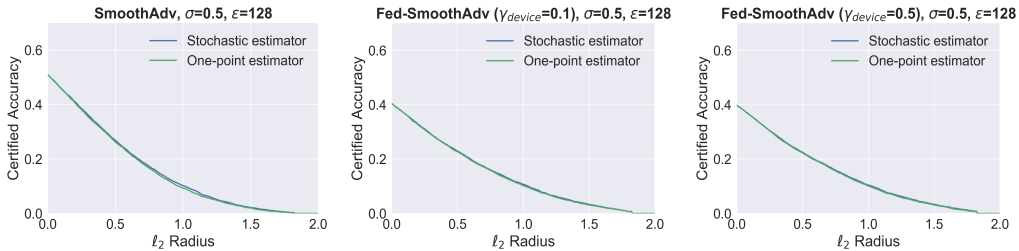


Figure 2: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.5$  and  $\epsilon = 128$ .

We further explore the certified accuracy of **Fed-SmoothAdv** under ablation settings. Figure 3 plots the result under  $\sigma = 0.25, \epsilon = 128$  with heterogeneous data. Please see Appendix B for more details and results. It can be seen that adversarial training of smoothed classifiers is critical for achieving a high certified accuracy. Standard training and adversarial training of original classifier (without smoothing) perform poorly in terms of certified robustness. This demonstrates the necessity of smoothed classifier to enable certifiably-robust federated learning.

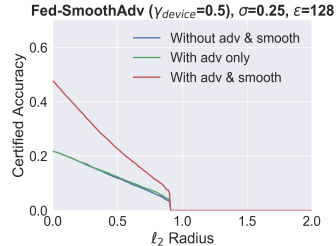


Figure 3: Ablation study of **Fed-SmoothAdv**.

### 4 CONCLUSION

In this paper, we incorporated the randomized smoothing techniques into the federated adversarial learning framework to enable certifiable robustness to test-time adversarial perturbations. We

demonstrated through extensive experiments that our adversarially smooth federated learning models could successfully achieve similar certified robustness as the centralized models. Meanwhile, we empirically proved that the device data heterogeneity and type of gradient estimator did not affect the performance much. The attempt in this paper is crucial for the applications of federated learning because of the adversarial attacks on its user’s devices and the resulting strong demand for user’s data privacy and security in the real world. In the future, we will apply randomized smoothing to more complex federated learning frameworks (Chen et al., 2020) and theoretically study its performance.

## ACKNOWLEDGEMENTS

This work was performed under the auspices of the U.S. Department of Energy by the Lawrence Livermore National Laboratory under Contract No. DE-AC52-07NA27344, Lawrence Livermore National Security, LLC. This document was prepared as an account of the work sponsored by an agency of the United States Government. Neither the United States Government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or Lawrence Livermore National Security, LLC. The views and opinions of the authors expressed herein do not necessarily state or reflect those of the United States Government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes. This work was supported by LLNL Laboratory Directed Research and Development project 20-SI-005 and released with LLNL tracking number LLNL-CONF-820514.

## REFERENCES

- Arjun Nitin Bhagoji, Supriyo Chakraborty, Prateek Mittal, and Seraphin Calo. Analyzing federated learning through an adversarial lens. In *Proc. International Conference on Machine Learning*, volume 97, pp. 634–643, 2019.
- Saikiran Bulusu, Bhavya Kailkhura, Bo Li, Pramod K Varshney, and Dawn Song. Anomalous example detection in deep learning: A survey. *IEEE Access*, 8:132330–132347, 2020.
- Cheng Chen, Ziyi Chen, Yi Zhou, and Bhavya Kailkhura. Fedcluster: Boosting the convergence of federated learning via cluster-cycling. *arXiv preprint arXiv:2009.10748*, 2020.
- Jeremy Cohen, Elan Rosenfeld, and Zico Kolter. Certified adversarial robustness via randomized smoothing. In *Proc. International Conference on Machine Learning*, volume 97, pp. 1310–1320, 09–15 Jun 2019.
- Ian Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. In *International Conference on Learning Representations*, 2015.
- Raouf Kerkouche, Gergely Ács, and Claude Castelluccia. Federated learning in adversarial settings. *arXiv:2010.07808*, 2020.
- Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. *ArXiv:1610.02527*, 2016.
- A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, volume 25, 2012.
- M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana. Certified robustness to adversarial examples with differential privacy. In *IEEE Symposium on Security and Privacy (SP)*, pp. 656–672, 2019.

- Linyi Li, Maurice Weber, Xiaojun Xu, Luka Rimanic, Tao Xie, Ce Zhang, and Bo Li. Provable robust learning based on transformation-specific smoothing. *arXiv preprint arXiv:2002.12398*, 2020a.
- Tian Li, Anit Kumar Sahu, Ameet Talwalkar, and Virginia Smith. Federated learning: Challenges, methods, and future directions. *IEEE Signal Processing Magazine*, 37(3):50–60, 2020b.
- Sijia Liu, Pin-Yu Chen, Bhavya Kailkhura, Gaoyuan Zhang, Alfred O Hero III, and Pramod K Varshney. A primer on zeroth-order optimization in signal processing and machine learning: Principals, recent advances, and applications. *IEEE Signal Processing Magazine*, 37(5):43–54, 2020.
- Brendan McMahan, Eider Moore, Daniel Ramage, Seth Hampson, and Blaise Aguera y Arcas. Communication-efficient learning of deep networks from decentralized data. In *Proc. International Conference on Artificial Intelligence and Statistics (AISTATS)*, volume 54, pp. 1273–1282, 20–22 Apr 2017.
- Hadi Salman, Jerry Li, Ilya Razenshteyn, Pengchuan Zhang, Huan Zhang, Sebastien Bubeck, and Greg Yang. Provably robust deep learning via adversarially trained smoothed classifiers. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *International Conference on Learning Representations*, 2014.
- Yao Zhou, Jun Wu, and Jingrui He. Adversarially robust federated learning for neural networks. 2021. URL <https://openreview.net/forum?id=5xaInvrGWp>.
- Giulio Zizzo, Amrith Rawat, Mathieu Sinn, and Beat Buesser. Fat: Federated adversarial training. *arXiv:2012.01791*, 2020.

## SUPPLEMENTARY MATERIAL

### A EXPERIMENT SETUP AND HYPERPARAMETERS

For **Fed-SmoothAdv**, we simulate 1000 edge devices and only 10% of them are sampled in each learning round. Each device holds 500 data samples. To control the data heterogeneity of each device, we define a data heterogeneity ratio  $\gamma_{\text{device}}$  in  $(0, 1)$ . Specifically, we randomly assign one class label as the major class of each device. Then, for each device,  $\gamma_{\text{device}}$  portion of samples are sampled from the major class, and the rest  $(1 - \gamma_{\text{device}})$  portion of samples are drawn from the remaining classes uniformly at random. In the experiments, we set  $\gamma_{\text{device}} = 0.1, 0.5$  that correspond to homogeneous data and heterogeneous data, respectively.

In the experiments, we set the number of Gaussian noise samples to be  $m = 2$ , and use 2 projected gradient descent steps for generating the adversarial samples. We set the inner-learning-rate for generating adversarial samples to 0.01, and the outer-learning-rate for updating the model parameters to 0.01. We set batch-size to 30 for each activated device of **Fed-SmoothAdv** and 60 for **SmoothAdv**. Moreover, each activated device of **Fed-SmoothAdv** uses 20 batches of data in the local training of a learning round, and this is equivalent to 1000 batches of data used by the centralized **SmoothAdv**. The total number of learning rounds is 150. In the certification phase, we set  $\alpha = 0.001$ , which means that there is at most 0.1% chance that the certification falsely certifies a non-robust input.

### B ABLATION STUDIES

In this section, we explore the certified accuracy of **Fed-SmoothAdv** under the following ablation settings: (1) Without adv & smooth: standard training of original classifier; (2) With adv only: adversarial training of original classifier (stochastic estimator); and (3) With adv & smooth: adversarial training of smoothed classifier, i.e., **Smoothadv** (stochastic estimator).

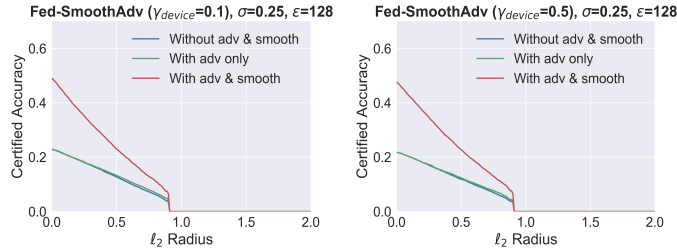


Figure 4: Ablation study of **Fed-SmoothAdv** with  $\sigma = 0.25$  and  $\epsilon = 128$ .

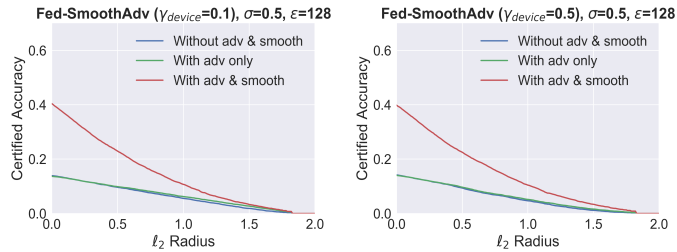


Figure 5: Ablation study of **Fed-SmoothAdv** with  $\sigma = 0.5$  and  $\epsilon = 128$ .

All of the plots yield following conclusions. First, the certified accuracy of **Fed-SmoothAdv** is much higher than that of standard and adversarial training of original (non-smoothed) classifier, which indicates that randomized smoothing is very helpful to improve the performance of **Fed-SmoothAdv**. Second, adversarial training does not achieve significantly higher certified accuracy than standard training, which again indicates that importance of having a smoothed classifier.

## C ADDITIONAL EXPERIMENTAL RESULTS

In this section, we present the certified accuracy results of both **SmoothAdv** and **Fed-SmoothAdv** (with heterogeneity  $\gamma_{\text{device}} = 0.1, 0.5$ ) under other choices of  $\sigma$  and  $\epsilon$ . From Figure 6-Figure 12, we observe the same comparison and make the same conclusions as those in Section 3.

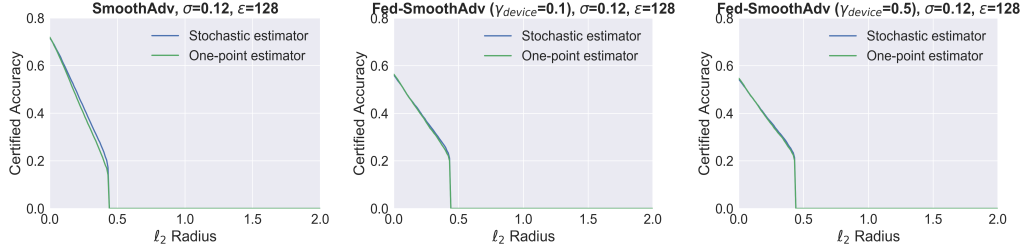


Figure 6: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.12$  and  $\epsilon = 128$ .

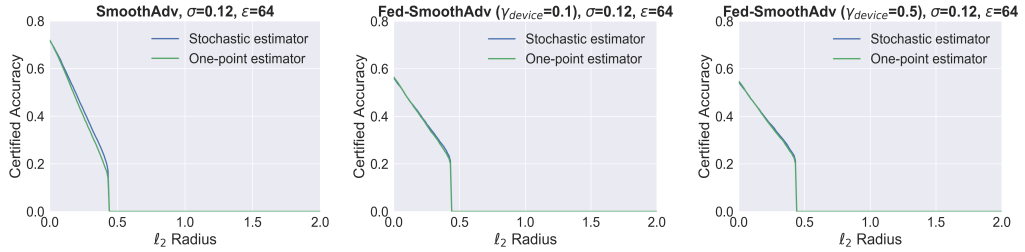


Figure 7: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.12$  and  $\epsilon = 64$ .

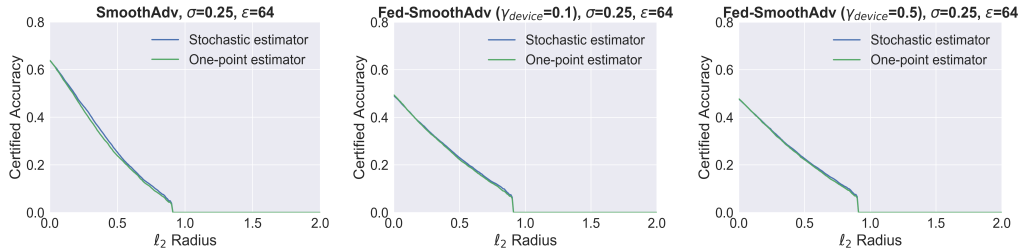


Figure 8: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.25$  and  $\epsilon = 64$ .



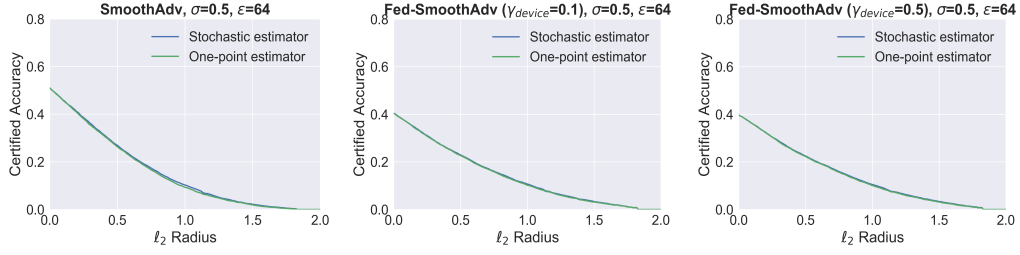


Figure 9: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.5$  and  $\epsilon = 64$ .

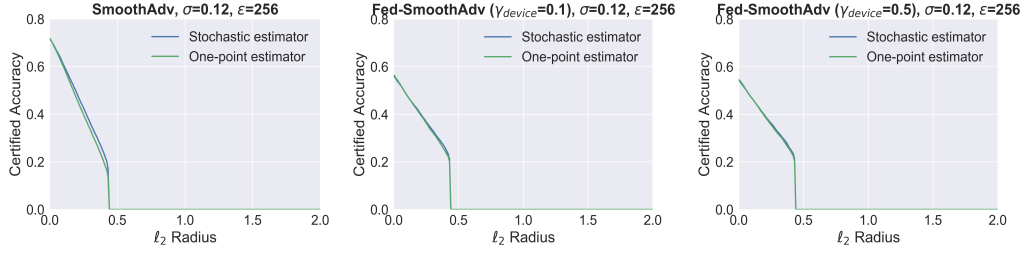


Figure 10: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.12$  and  $\epsilon = 256$ .

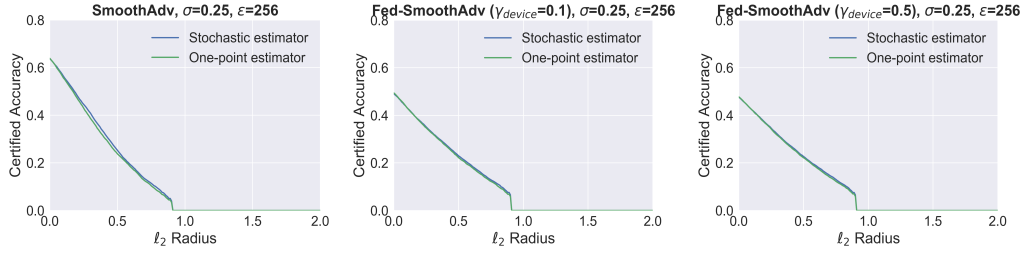


Figure 11: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.25$  and  $\epsilon = 256$ .

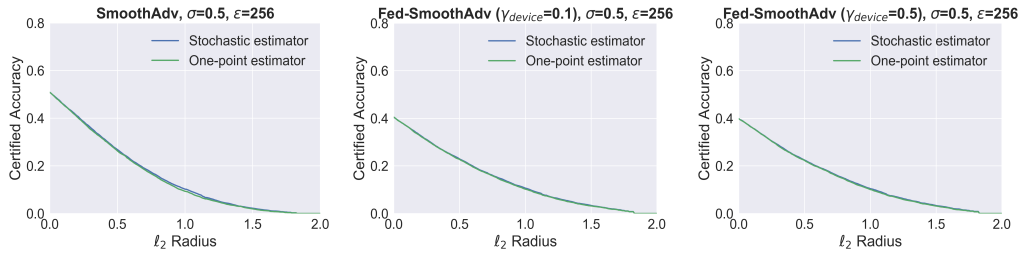


Figure 12: Certified accuracy of **SmoothAdv** and **Fed-SmoothAdv** with  $\sigma = 0.5$  and  $\epsilon = 256$ .