

Fault Tolerance of Neural Networks in Adversarial Settings

Vasisht Duddu¹, N. Rajesh Pillai², D. Vijay Rao³, Valentina E. Balas⁴

¹ Indraprastha Institute of Information Technology, Delhi, India

² Scientific Analysis Group, Delhi, India

³ Institute for Systems Studies and Analyses, Delhi, India

⁴ Aurel Vlaicu University of Arad, Arad, Romania

vduddu@tutamail.com, rpillai@sag.drdo.in, vijayrao@issa.drdo.in, valentina.balas@uav.ro

ABSTRACT

Applications using Artificial Intelligence techniques demand a thorough assessment of different aspects of trust, namely, data and model privacy, reliability, robustness against adversarial attacks, fairness, and interpretability. While each of these aspects has been extensively studied in isolation, an understanding of the trade-offs between different aspects of trust is lacking. In this work, the trade-off between fault tolerance, privacy, and adversarial robustness is evaluated for Deep Neural Networks, by considering two adversarial settings under security and a privacy threat model. Specifically, this work studies the impact of training the model with input noise (Adversarial Robustness) and gradient noise (Differential Privacy) on Neural Network's fault tolerance. While adding noise to inputs, gradients or weights enhances fault tolerance, it is observed that adversarial robustness lowers fault tolerance due to increased overfitting. On the other hand, $(\epsilon_{dp}, \delta_{dp})$ -Differentially Private models enhance the fault tolerance, measured using generalisation error, which theoretically has an upper bound of $e^{\epsilon_{dp}} - 1 + \delta_{dp}$. This novel study of the trade-offs between different aspects of trust is pivotal for training trustworthy Machine Learning models.

KEYWORDS

Trustworthy Machine Learning, Differential Privacy, Fault Tolerance, Adversarial Robustness, Deep Learning.

1 INTRODUCTION

There is a growing reliance on Artificial Intelligence (AI) techniques in safety-critical real-time applications with high-stake decision-making such as autonomous vehicles, criminal justice, and healthcare. These applications demand satisfying different aspects of trust: fairness among disparate groups, the privacy of individuals in the training data, robustness against adversarially perturbed inputs, fault tolerance for safety and model interpretability. Training the models to incorporate and optimize for all the aspects of trust is difficult and hence, it is crucial to understand the trade-offs between different aspects of trust for designing efficient Pareto-optimal solutions for trustworthy AI systems. Prior research has indicated that privacy and explainability [32], adversarial robustness and membership privacy [35] are at odds. On the other hand, other aspects of trust go hand in hand: fairness and differential privacy [18], adversarial robustness and explainability [13]. However, the impact on fault tolerance due to robust and private models has not been explored yet. To address this requirement, this research analyses the impact of training Machine Learning models, specifically Deep Neural Networks, for adversarial robustness (security) and differential privacy on the model's fault tolerance (reliability).

Fault Tolerance is a crucial property for Neural Networks to ensure reliable computation for a long duration with graceful degradation over time. Typically, well generalized models have the parameters with low variance ensuring equal computational weight to all nodes in the network [10]. Hence, the loss of some of the nodes can be compensated by other nodes without a significant loss in performance [8][7][27]. Practically, this is achieved by adding noise during training to the inputs, gradients or the weights. The noise added to the inputs can be modeled as Tikhonov regularization which enhances the generalization [4].

However, in an adversarial setting under a security threat model, carefully crafted *imperceptible* noise can be added to input images by an adversary to force the model to misclassify the image, violating the integrity of the model prediction. The design of Neural Networks within such an adversarial setting requires training on inputs with adversarial noise to ensure robustness against adversary's worst-case perturbation. The goal of this work, in this setting, is to address the following research question,

What is the impact of adversarially robust (input noise) training on fault tolerance?

Alternatively, in an adversarial setting under a privacy threat model, an adversary performs inference attacks to identify training data attributes or membership details for sampled data point [15][33][31]. This poses a serious privacy risk for sensitive training data such as financial and medical records, personal photos, location history, and user preferences. Differential Privacy provides a provable guarantee on the maximum privacy leakage by making the data points indistinguishable using gradient noise during training [12][1]. In such a setting, this work addresses the following research question,

What is the impact of training models with Differential Privacy (gradient noise) on fault tolerance?

Main Contributions. This work makes the following novel contributions:

- Evaluate the fault tolerance of provably robust Neural Networks and compare them with the theoretical equivalent Tikhonov regularized model.
- Evaluate the fault tolerance of Differentially Private models under a privacy threat model and compare it with regularised and naturally (without noise) trained models.
- Prove theoretically that generalization error, used to measure fault tolerance, has an upper bound of $e^{\epsilon_{dp}} - 1 + \delta$ on training the model using $(\epsilon_{dp}, \delta_{dp})$ -Differential Privacy, thus giving provable guarantees on fault tolerance metric.

The performance of the training algorithms is evaluated using two common benchmarking datasets: CIFAR10 and FashionMNIST,

with different Neural Network architectures. It is observed that adversarial robustness and fault tolerance are at odds with each other, i.e., training a model with adversarial input noise results in overfitting which lowers fault tolerance. On the other hand, noise added for $(\epsilon_{dp}, \delta_{dp})$ -Differential Privacy is an alternate approach for enhancing fault tolerance, while guaranteeing privacy, with a theoretical bound on the generalization error in terms of the privacy parameters: $\epsilon_{dp}, \delta_{dp}$. To the best of our knowledge, this is the first work that evaluates the fault tolerance of adversarially robust and differentially private Deep Neural Networks. Such an analysis is crucial for a unified framework for trustworthy Machine Learning combining security, privacy, and reliability for real-world deployment.

Table 1: Variable notations for Reliability, Adversarial Robustness and Differentially Privacy.

Symbol	Description
\mathcal{F}	Target Machine Learning Model
\mathcal{D}_{train}	Training Dataset
\mathcal{D}_{test}	Testing Dataset
\mathcal{X}	Space of Input Data Points
\mathcal{Y}	Space of Output Labels
(x, y)	Data Sample with input and output label
$\mathbf{P}(x, y)$	Data Distribution over all samples
\mathcal{L}	Loss Function for Training
Fault Tolerance	
ϵ_{ft}	Fault Tolerance Bound
\mathcal{N}	Ideal Neural Network without Faults
\mathcal{N}_{fault}	Neural Network model with Faults
$\mathcal{H}_{\mathcal{N}}$	Trained Ideal (faultless) Neural Network
$\mathcal{H}_{\mathcal{N}_{fault}}$	Trained Faulty Neural Network
\mathcal{G}_{error}	Generalization Error; Metric for measuring fault tolerance
Adversarial Robustness	
ϵ_{adv}	Adversarial Noise Bounds
Δ	Set of Adversarial Noise Values bounded by ϵ_{adv}
δ_{adv}	Value of Adversarial Noise in Δ
Differentially Privacy	
ϵ_{dp}	Differential Privacy Leakage Bound
δ_{dp}	Differential Privacy failure probability
\mathbb{E}_{train}	Error on training set
\mathbb{E}_{test}	Error on test set
\mathcal{D}'	Dataset after adding/removing single data point

2 BACKGROUND

2.1 Fault Tolerance in Neural Networks

DEFINITION 1. A Neural Network \mathcal{N} performing computations $\mathcal{H}_{\mathcal{N}}$ is said to be fault tolerant if the computation performed by a faulty network $\mathcal{H}_{\mathcal{N}_{fault}}$ is close to $\mathcal{H}_{\mathcal{N}}$. Formally, a Neural Network is ϵ fault tolerant if,

$$\left\| \mathcal{H}_{\mathcal{N}}(\mathcal{X}) - \mathcal{H}_{\mathcal{N}_{fault}}(\mathcal{X}) \right\| \leq \epsilon_{ft} \quad (1)$$

for $\epsilon_{ft} > 0$ and $\mathcal{X} \in \mathcal{D}$.

Fault Tolerance Metric. Improving the generalization results in enhancing the fault tolerance and vice versa [3]. Hence, capturing the overfitting of the model provides a way to compare the relative fault tolerance of multiple models. This has been extensively used in literature to measure fault tolerance of Neural Networks [36] [10].

Formally, given a dataset \mathcal{D}_{test} such that $\mathcal{D}_{train} \cap \mathcal{D}_{test} = \phi$, the accuracy of the model is estimated on the training set (R_{train}) and on the testing set (R_{test}). The generalization error is given by the difference between training accuracy and the testing accuracy,

$$\mathcal{G}_{err} = \mathcal{R}_{train} - \mathcal{R}_{test} \quad (2)$$

This gives the estimate of overfitting in the Neural Network ($R_{train} > R_{test}$), i.e., higher the generalization error more the overfitting. This estimate of fault tolerance is used for the evaluation of different Neural Networks throughout the paper.

Fault Model. In this work, the faults occurring in the hardware are simulated in the form of stuck at “0” errors during the Neural Network computation. Further, multiple faults can occur simultaneously for which the performance degradation is measured using test accuracy. The faults are simulated in two ways: firstly, these faults can manifest in the form of random node crashes in the Neural Network due to which the output of the node is forced to zero. Secondly, the parameters of the Neural Networks can be stuck at “0” which includes weights in the case of Multilayer Perceptron, and filter and kernel values in case of Convolutional Neural Networks. This is a common fault model frequently used in evaluating the reliability of systems [10].

Related Work. A widely used approach for enhancing fault tolerance is to penalize large values of the parameters using a regularization function [8][36]. Alternatively, constraint optimisation approaches using minimax constraint [27] as well as quadratic programming [7] can be used for small networks. Unlike simple regularization functions, unsupervised pre-training of the initial network layers followed by supervised fine-tuning can significantly enhance the fault tolerance [10]. Traditional techniques to enhance reliability such as additional redundancy by adding nodes and synapses provides partial fault tolerance [29]. Further, reliability of axonal transport has been explored using Hammock Networks [2][6]. Detection of faults and enhancing tolerance in software has been explored for fuzzy control systems [21].

2.2 Adversarial Robustness

Within the adversarial setting with a security adversary, the problem of adversarial robustness is modeled as a game between the attacker and defender with conflicting interests. Here, the adversary wants to force the target model to misclassify by adding carefully calculated noise to input, while the defender wants to train the model to defend against such inputs with adversarial noise [9].

Attacker Knowledge. In this work, the adversary has no knowledge about the target model. In other words, the adversary has remote access to the target black-box model and can query the model and receive corresponding predictions through an API. This is typically the black box setting seen in Machine Learning as a Service.

Attacker Goal. The goal of the adversary is to find the noise to maximize the loss of the target model ($F_{\theta}()$) and force the model to misclassify the perturbed input. Formally, given an input sampled from the underlying data distribution $x \sim P(X, Y)$, the adversary computes the worst-case adversarial noise δ_{adv} to maximize the loss (l) between predicted output and true output y ,

$$\delta^*(x)_{adv} = \underset{\|\delta_{adv}\| \leq \epsilon_{adv}}{\operatorname{argmax}} \ell(F_{\theta}(x + \delta_{adv}), y) \quad (3)$$

subject to a bound on the perturbation computed using parameter ϵ_{adv} ,

$$\Delta = \{\delta_{adv} : \|\delta_{adv}\|_p \leq \epsilon_{adv}\} \quad (4)$$

The optimization is subjected to the condition that the noise is imperceptible by restricting δ within a perturbation region Δ defined by a l_p norm, more commonly l_∞ [22].

Defender Strategy. To defend against the worst attack possible, the model is trained using adversarial inputs as part of the training data. Formally, this empirical defense can be modeled as a minimax optimization problem given below,

$$\min_{\theta} \left(\frac{1}{|\mathcal{D}|} \sum_{x,y \in \mathcal{D}} \max_{\|\delta\| \leq \epsilon} \ell(F_{\theta}(x + \delta), y) \right) \quad (5)$$

Here, instead of minimising the expected loss over the data points sampled from the distribution, the optimisation minimises the worst case loss over the data with adversarial noise. In other words, the defender minimises the loss corresponding to the worst case adversarial attack.

In this work, TRADES algorithm is considered as a defense since it provides provable bounds against adversarial examples [38]. TRADES algorithm decomposes the prediction error for adversarial example as the sum of natural classification error and the boundary error to provide a tight differentiable upper bound. This defense minimizes the maximum Kullback-Leibler (KL) Divergence between the output prediction corresponding to a benign sample (X) and adversarial sample ($X_{adv} \leftarrow X + \delta_{adv}$). This is used to generate adversarial examples within the inner maximization.

$$\min_{\theta} (\mathbb{E}[l(\theta, F(X), Y)] + \max_{X_{adv} \in (X, \epsilon)} d_{kl}(F(X), F(X_{adv}))) \quad (6)$$

Algorithm 1 Adversarial Training by adding noise to inputs

Input: $D_{train} = \{(x_i, y_i), \dots, (x_N, y_N)\}$ and Loss: $L(\theta_t, x_i)$

Input: \mathbb{A} : Algorithm to generate adversarial noise

- 1: **for** each epoch **do**
 - 2: Sample a random batch $B \in D_{train}$
 - 3: **for** each $(x_i, y_i) \in B$ **do**
 - 4: **Compute Adversarial Noise:**
 $x_{adv}^* \leftarrow \mathbb{A}(x_i, y_i)$
 - 5: **Compute gradient on adversarial inputs:**
 $\frac{\partial L_i}{\partial \theta_t} \leftarrow \nabla_{\theta_t} L(\theta_t, x_{adv}^*)$
 - 6: **Update θ using Gradient Descent:**
 $\theta_{t+1} \leftarrow \theta_t - \alpha \frac{\partial L_i}{\partial \theta_t}$
 - 7: **end for**
 - 8: **Output:** Parameters θ of trained model with adversarial robustness
 - 9: **end for**
-

Related Work. While only the defense algorithm with tight upper bound and provable robustness guarantees is considered in this work (TRADES), other empirical approaches use Projected Gradient Descent [22] and Wasserstein norm [34] to ensure robustness. Alternatively, verification based defenses use function transformations to compute the worst-case loss to express the adversarial perturbations [24].

2.3 Differential Privacy

Differential Privacy is the de facto privacy standard that provides a strong privacy definition with provable bounds on information leaked by a mechanism in terms of the privacy budget ϵ_{dp} [12]. The output of a randomized mechanism should not allow the adversary to learn any more about an individual in the dataset than that could be learned via the same analysis without the individual in the dataset. In this sense, this definition of privacy captures the individual’s membership in the dataset.

DEFINITION 2. For a randomized mechanism $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{Y}$ is $(\epsilon_{dp}, \delta_{dp})$ differentially private on two neighbouring datasets \mathcal{D} and \mathcal{D}' differing by an individual data point, then for all outputs $\mathcal{O} \subseteq \mathcal{Y}$,

$$\mathcal{P}[\mathcal{M}(\mathcal{D}) \in \mathcal{O}] \leq e^{\epsilon_{dp}} \mathcal{P}[\mathcal{M}(\mathcal{D}') \in \mathcal{O}] + \delta_{dp} \quad (7)$$

Here, the parameters ϵ_{dp} is considered as the privacy budget and δ_{dp} is considered as the failure probability [11].

A tighter and accurate estimation of the privacy loss can be computed using the Renyi Differential Privacy [25] which uses the Renyi divergence metric \mathcal{D}_α which applies to any moment of the privacy loss random variable.

DEFINITION 3. For a randomized mechanism $\mathcal{M} : \mathcal{X} \rightarrow \mathcal{Y}$ is ϵ_{dp} -Renyi differentially private of the order α , on two neighbouring datasets \mathcal{D} and \mathcal{D}' differing by a individual data point, then for all outputs $\mathcal{O} \subseteq \mathcal{Y}$,

$$\mathcal{D}_\alpha(\mathcal{M}(\mathcal{D}) \parallel \mathcal{M}(\mathcal{D}')) \leq \epsilon \quad (8)$$

Attacker Strategy. The goal of the attacker within the privacy setting is to use inference attacks to leak training data details resulting in privacy violations where the training data is sensitive. Membership inference attacks infer whether a given data point was used in the training data or not based on the difference in model performance on training data and testing data [33]. On the other hand, attribute inference attacks extract particular features of the training data [14] or reconstruct the entire training data [30]. Another class of attacks exploits the memorization capacity of the model to infer the sensitive attributes in the data by querying the model [5].

Defender Strategy. The defender, in order, to reduce the success of the inference attacks utilizes the notion of differential privacy to train the model with provable privacy leakage guarantees. To this extent, the defender samples a noise from either a Laplace or Gaussian distribution proportional to the sensitivity of the model (S).

$$\mathcal{M}(D) \simeq \mathcal{H}_{\mathcal{N}}(\mathcal{X}) + \mathcal{N}(0, S^2, \sigma^2) \quad (9)$$

In case of Deep Neural Networks, this noise is added to the gradients of the model and mitigates several of the attacker’s inference strategies, however, at the cost of utility (performance) [19].

Related Work. To preserve privacy against membership inference attacks, several empirical defenses exist. For instance, the adversary’s inference attack can be modeled as a minimax optimization problem, where the target model is trained to minimize the adversary’s best attack [26]. Another line of defense is to add noise to the output of the model, force the inference attack machine learning model to misclassify while ensuring the utility does not degrade

[20]. While these approaches are empirical, all of them face utility privacy trade-off and do not provide a theoretical guarantee on the maximum leakage of the model. Hence, in this work, Differential Privacy based private training is used since it provides provable guarantees on the leakage of the model about the training data [1]. An alternative Differential Privacy based framework uses a teacher-student ensemble approach [28].

Algorithm 2 Differentially Private Training by adding noise to gradients during Backpropagation

Input: $D_{train} = \{(x_i, y_i), \dots, (x_N, y_N)\}$ and Loss function: $L(\theta_t, x_i)$

- 1: **for** each epoch **do**
- 2: Sample a random batch $B \in D_{train}$
- 3: **for** each $(x_i, y_i) \in B$ **do**
- 4: **Compute gradient:**
 $\frac{\partial L_i}{\partial \theta_t} \leftarrow \nabla_{\theta_t} L(\theta_t, x_i)$
- 5: **Gradient Clipping:**
 $g_t(x_i) \leftarrow \frac{\partial L_i}{\partial \theta_t} / \max(1, \frac{\|\frac{\partial L_i}{\partial \theta_t}\|_2}{C})$
- 6: **Add Noise:**
 $g'_t \leftarrow \frac{1}{L} (\sum_i g_t(x_i) + \mathcal{N}(0, \sigma^2 C^2 I))$
- 7: **Update θ using Gradient Descent:**
 $\theta_{t+1} \leftarrow \theta_t - \alpha g'_t$
- 8: **end for**
- 9: **Output:** Parameters θ of model trained using (ϵ, δ) -Differential Privacy.
- 10: **end for**

3 EXPERIMENT SETUP

The code for training using TRADES adversarial training algorithm is based on the official source code from the authors¹. The official code from Tensorflow Privacy Library² for Differentially Private training is adapted to different architectures and datasets.

3.1 Datasets

The evaluation and training of adversarially robust and differentially private models are done on two major benchmarking datasets, namely, FashionMNIST and CIFAR10.

FashionMNIST. The dataset is similar to the MNIST dataset and consists of a training set of 60,000 examples and a test set of 10,000 examples. Each data sample is a 28×28 grayscale image associated with a label from 10 classes such as boots, shirt, bag and so on.

CIFAR10. The CIFAR-10 dataset consists of 60000 32x32 colour images in 10 classes, with 6000 images per class. There are 50000 training images and 10000 test images.

3.2 Architectures

For the TRADES robust training algorithm on CIFAR10 dataset, an Neural Network architecture with nine convolutional layers followed by fully connected layers is used. In the case of FashionMNIST dataset with adversarial robust training, the FMNIST CNN 2

Table 2: Architectures for CIFAR10 and FashionMNIST datasets.

CIFAR10 Architecture	FMNIST CNN 1	FMNIST CNN 2
Convolution 32 (3x3)(1)	Convolution 20 (5x5)(1)	Convolution 16 (8)(1)
Convolution 32 (4x4)(2)	MaxPool (2) (2)	MaxPool (2) (1)
Convolution 32 (3x3)(1)	Convolution 50 (5x5)(1)	Convolution 32 (4)(2)
Convolution 64 (3x3)(2)	MaxPool (2) (2)	MaxPool (2) (1)
Convolution 64 (3x3)(1)	Dense 500	Dense 32
Convolution 128 (3x3)(2)	Dense 10	Dense 10
Convolution 128 (3x3)(1)		
Convolution 256 (3x3)(2)		
Convolution 256 (3x3)(1)		FMNIST DNN
Dense 512		Dense 512
Dense 512		Dense 512
Dense 10		Dense 10

architecture based on LeNet architecture is used. For Differential Private training, FMNIST CNN 1 architecture is used with minor differences in hyperparameters to the FMNIST CNN 2 architecture. Further, for evaluating on a simple Multilayer Perceptron Network, a Neural Network with two hidden layers of sizes 512 nodes each is used. The details of the exact architectures used in the experiments are given in Table 2.

4 FAULT TOLERANCE AND ADVERSARIAL ROBUSTNESS

4.1 Input Noise as a Regularizer

Adding noise to inputs has been shown to provide a regularization effect that theoretically is equivalent to Tikhonov regularization [4]. To understand the effect on input noise to generalization, a simple architecture of 500 hidden layers is considered for a binary classification problem and differentiating two types of circles³. Here, two types of noise are considered, namely, additive Gaussian noise and the multiplicative Gaussian noise which are commonly used for enhancing the fault tolerance of Neural Networks by improving the generalization [16] [17] [23]. This indicates that adding noise enhances for fault tolerance for small noise values beyond which the model overfits and the performance degrades.

Table 3: Generalization Error for additive gaussian noise and multiplicative gaussian noise added to inputs.

σ	Training Accuracy	Testing Accuracy	Generalization Error
No Noise			
0	100.00%	75.7%	24.3%
Additive Gaussian Noise			
0.01	100.00%	77.1%	22.9%
0.1	93.3%	67.1%	26.2%
0.5	73.3	44.3	29.0%
Multiplicative Gaussian Noise			
0.03	100.00%	77.1%	22.9%
0.1	96.7%	71.4%	25.3%
0.5	80.0	48.6	31.4%

¹<https://github.com/yaodongyu/TRADES>

²<https://github.com/tensorflow/privacy>

³sklearn.datasets.make_circles

In the case of additive noise, on adding a small noise of a standard deviation of 0.01, the generalization performance on the binary classification problem improves (Table 3). However, on increasing the values of the standard deviation, the model starts to overfit the noisy training data and the performance starts to decline. A similar phenomenon is observed on training models in the presence of worst-case adversarial noise as shown in the subsequent sections.

4.2 Adversarial Noise

On adding adversarial noise, the goal is to estimate the extent of overfitting for different training algorithms which will indicate the impact on fault tolerance. As shown in Figure 1, the model on training using adversarially robust algorithm overfits significantly compared to model trained using natural Stochastic Gradient Descent. In this particular architecture, the generalization error of robust models is about 17% compared to 9% error of naturally trained model.

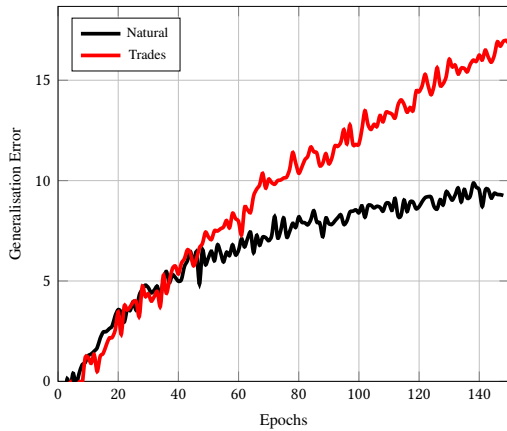


Figure 1: The training curve for naturally training a Neural Network compared to robust training using TRADES algorithm.

For CIFAR10 dataset, the performance for robust training using TRADES is evaluated and compared with the baseline of Tikhonov regularization. Tikhonov regularized models are used as a baseline since, training with input noise is shown as theoretically equivalent to Tikhonov function in the objective function. As shown in Table 4, adding adversarial noise results in increase in the overall generalization error (extent of model overfitting).

Table 4: Generalization Error of Adversarially Robust Models on CIFAR10.

Algorithm	Training Accuracy	Testing Accuracy	Generalization Error
Natural	95.50%	86.24%	9.26%
Tikhonov	89.19%	82.24%	6.95%
TRADES	93.26%	76.04%	17.22%

For the FashionMNIST dataset, the generalization error for the TRADES adversarial training algorithm is evaluated. As shown in Table 5, the generalization error for the models trained using adversarial noise is significantly higher compared to the generalization

of the Tikhonov regularization and naturally trained model without any additional optimizations. This indicates that adversarially computed noise, in fact, has a negative impact on the fault tolerance. On evaluating the generalization error on Multilayer Perceptron, the error increases from 3.29% for regularized model to 8.83% for robust models. A similar pattern is observed for a more complex Convolutional Neural Network, with the generalization error increases from 4.90% for a regularized model to 8.01% for robust models.

Table 5: Generalization Error of Adversarially Robust Models on FashionMNIST on both CNN and MLP architectures.

Multilayer Perceptron Architecture			
Algorithm	Training Accuracy	Testing Accuracy	Generalization Error
Natural	98.83%	90.23%	8.60%
Tikhonov	90.99%	87.70%	3.29%
TRADES	98.76%	89.93%	8.83%
Convolutional Neural Network Architecture			
Algorithm	Training Accuracy	Testing Accuracy	Generalization Error
Natural	98.60%	90.59%	8.01%
Tikhonov	95.94%	91.04%	4.90%
TRADES	99.59%	90.73%	8.86%

4.3 Comparing Fault Tolerance through Parameter Distribution

Alternatively, another approach to estimate the fault tolerance is by evaluating the standard deviation of the trained model’s parameter distribution [10]. The standard deviation of the parameters θ of the model is written as,

$$\sigma = \sqrt{\frac{\sum_i |\theta_i - \bar{\theta}|}{N}} \quad (10)$$

where $\bar{\theta}$ is the average of all the parameter values and N is the total number of parameters in the model. Higher the standard deviation of the parameter distribution, more varied are the parameter values due to which some nodes are given more importance over the others. Here, the loss of those important nodes in case of random faults results in a significant drop in accuracy. For low σ , the parameters give equal weightage to all the nodes and hence, the loss of a few nodes does not impact the overall model performance.

In Figure 2 (left) for the CIFAR10 dataset, the standard deviation of the adversarially robust model is 0.125378 compared to 0.032451 of the regularized model parameter distribution. For FashionMNIST dataset (Figure 2 (right)), the standard deviation of the parameter distribution for robust model is 0.17901 while the deviation for regularized model is 0.023772. This indicates that the fault tolerance of models trained using adversarial noise is significantly less than Tikhonov regularized model or naturally trained model.

4.4 Impact of Varying ϵ_{adv}

An important study is to evaluate the impact of increasing the overall range of perturbation added to the inputs, i.e., ϵ_{adv} . Increasing ϵ_{adv} , results in increasing the overall noise region from which the noise can be sampled. Thus, this results in increasing the strength of the noise added to the input.

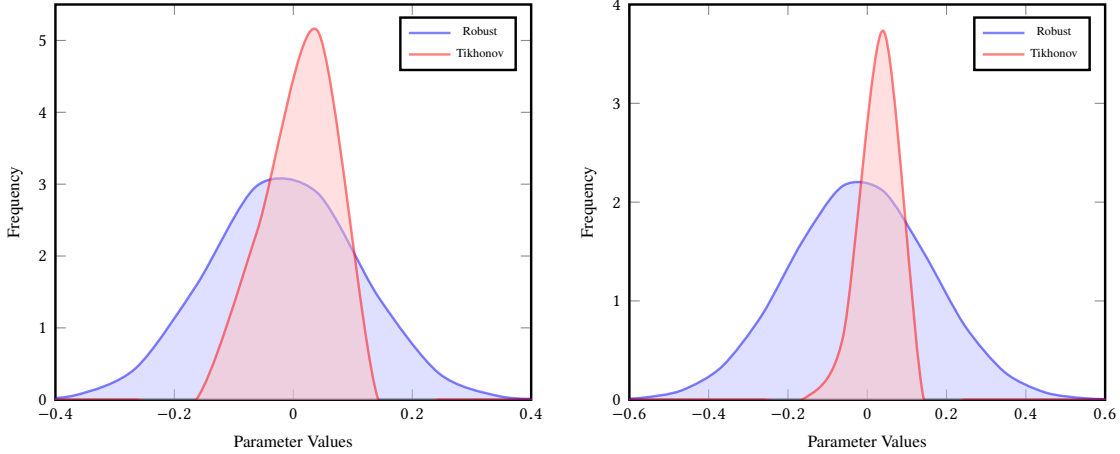


Figure 2: (Left) CIFAR10. (Right) FashionMNIST. For both the datasets, the distribution for the model trained using Tikhonov regularisation has a lower standard deviation of the parameter distribution compared to the models trained using adversarial input noise. This difference in the distributions indicate that training with noise enhances fault tolerance only for small values of noise and it is not equivalent to Tikhonov regularisation for adversarial noise.

Table 6: The impact of varying ϵ_{adv} values on the fault tolerance for CIFAR10 dataset.

ϵ_{adv}	Training Accuracy	Testing Accuracy	Generalization Error
2/255	90.12%	81.93%	8.19%
4/255	91.02%	80.05%	10.97%
8/255	93.26%	76.04%	17.22%

As seen in Table 6, the overall fault tolerance measured as the difference in training and testing accuracy increases with an increase in the noise budget ϵ_{adv} . Specifically, for the case of Convolutional Neural Network trained on CIFAR10 dataset, increasing ϵ_{adv} from 2/255 to 8/255, the generalization error increases (equivalently fault tolerance decreases) from 8.19% to 17.22%.

4.5 Simulation of Faults for Adversarially Robust Models

Since the fault model considered in the paper is random faults resulting in stuck at "0" values, these faults are simulated for adversarially robust models and compared with Tikhonov regularized models. As seen in Figure 3, on increasing the faults into the parameters from 50% to 90%, the accuracy drop in the case of the regularized model is 1.14% compared to 26.93% of the robust model. This confirms the above analysis, indicating that adversarially robust models are less fault-tolerant to regularized models.

5 DIFFERENTIAL PRIVACY AND FAULT TOLERANCE

In this section, the impact of adding noise to the gradients, to mitigate inference attacks via Differential Privacy, on fault tolerance is considered.

For the Convolutional Neural Network (Table 7), a clear trade-off between the generalization error (fault tolerance) and accuracy

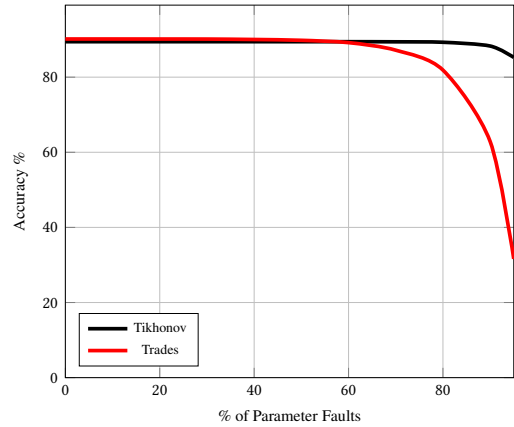


Figure 3: Comparison of performance drop in model trained using the theoretically equivalent Tikhonov regularisation (fault tolerant) and model trained using adversarial input noise.

can be seen. Higher fault tolerance can be achieved at the cost of low test accuracy. This trade-off is also observed by other standard functions such as L1 and L2 regularizers, however, they do not provide privacy guarantees. As the privacy leakage bound ϵ_{dp} is increased from 0.49 to 10^6 , the generalization error increases from 0.75% to 4.40%. An increase of ϵ_{dp} indicates more information leakage. This indicates that fault tolerance and privacy are highly correlated with each other, i.e., increasing the privacy (lowering ϵ_{dp}) will also increase the overall fault tolerance at the cost of test accuracy.

As seen in Table 8 for MLP based model, a similar pattern is observed where the generalization error increases from 1.01% to 8.29% as the values of ϵ_{dp} increases.

Table 7: CNN: Generalization Error of Differentially Private Models on FashionMNIST.

	Training Accuracy	Testing Accuracy	Generalization Error
Natural	97.06%	89.92%	7.14%
Tikhonov	90.12%	89.43%	0.69%
ϵ_{dp}	Training Accuracy	Testing Accuracy	Generalization Error
0.49	76.84%	76.09%	0.75%
2.97	84.35%	83.22%	1.13%
24.66	87.72%	86.26%	1.46%
2×10^6	94.21%	89.81%	4.40%

Table 8: MLP: Generalization Error of Differentially Private Models on FashionMNIST.

	Training Accuracy	Testing Accuracy	Generalization Error
Natural	99.58%	89.60%	9.98%
Tikhonov	88.96%	87.45%	1.51%
ϵ_{dp}	Training Accuracy	Testing Accuracy	Generalization Error
0.49	80.44%	79.43%	1.01%
2.97	85.70%	83.64%	2.06%
24.66	88.21%	85.43%	2.78%
2×10^6	96.64%	88.35%	8.29%

THEOREM 1. *Given a Machine Learning Model trained using $(\epsilon_{dp}, \delta_{dp})$ -Differential Privacy, the model's fault tolerance metric, given by the generalization error, is bounded by $e^{\epsilon_{dp}} - 1 + \delta_{dp}$.*

Proof Sketch. Differential Privacy is a strong notion of stability where the change in the data point in the training data should not change the final output. Further, fault tolerance is also a notion of stability where a change in the model architecture should not change the final output. A Differentially private mechanism is also uniform RO stable and the generalization error of the mechanism can be bounded by $e^{\epsilon_{dp}} - 1 + \delta_{dp}$ [37]. Since generalization error is used to measure the relative fault tolerance between different models, the corresponding fault tolerance is bounded by $e^{\epsilon_{dp}} - 1 + \delta_{dp}$.

Proof. Given the data population P of all possible input and output pairs, the model is trained on a subset of data D_{train} sampled from P by minimising the training error,

$$E_{train} = \frac{1}{n_{train}} \sum_i l(f(\theta, x_i), y_i) \quad (11)$$

In order to evaluate the performance on any possible sample that the model might encounter, we evaluate the error on the testing dataset D_{test} sampled from P , where $D_{test} \cap D_{train} = \emptyset$.

$$E_{test} = \frac{1}{n_{test}} \sum_{i'} l(f(\theta, x_{i'}), y_{i'}) \quad (12)$$

The generalization error is given by the difference between the testing (E_{test}) and training error (E_{train}).

A mechanism which satisfies $(\epsilon_{dp}, \delta_{dp})$ -Differential Privacy also satisfies uniform RO stability [37]. Hence, for datasets D and D' differing by a single data point,

$$|E_D - E_{D'}| \leq e^{\epsilon_{dp}} - 1 + \delta_{dp} \quad (13)$$

Further, generalizing this result for the training dataset and testing dataset,

$$|E_{train} - E_{test}| \leq e^{\epsilon_{dp}} - 1 + \delta_{dp} \quad (14)$$

Since the fault tolerance is measured as the difference in the training and testing error, we can see that this is bounded by $e^{\epsilon_{dp}} - 1 + \delta_{dp}$ on training the model with $(\epsilon_{dp}, \delta_{dp})$ -Differential Privacy. This result on provable bound on generalization error is based on the folklore theorem by Frank McSherry. For small values of ϵ_{dp} , $e^{\epsilon_{dp}} \approx 1 + \epsilon_{dp}$ and hence, $e^{\epsilon_{dp}} - 1 + \delta_{dp}$ can be written as $\epsilon_{dp} + \delta_{dp}$ which is agreement with folklore theorem. Hence, training for privacy objective using Differential privacy provides an alternate approach for enhancing fault tolerance with a provable bound on the generalization error.

6 CONCLUSIONS

Designing a trustworthy Machine Learning system requires to understand the trade-offs between different aspects of trust. This work highlights the trade-offs between three such aspects of trust in Machine Learning, namely, reliability, privacy, and adversarial robustness. This work considers two adversarial settings, with a security threat model where the adversary aims to force the model to misclassify by adding adversarial noise to the input, and a privacy threat model where the adversary aims to infer whether a data point was part of the sensitive training data or not. Under the security threat model, the impact of fault tolerance on adversarially robust Neural Networks is evaluated and robust Neural Networks are observed to have lower the fault tolerance due to overfitting. Under the privacy threat model, it is shown that Differentially Private models exhibit fault tolerance for a careful choice of privacy parameters $(\epsilon_{dp}, \delta_{dp})$. Hence, fault tolerance can be achieved by training models with privacy objective. Theoretically, the bound on the model's generalization error is shown in terms of the parameters for Differential Privacy. This study is a crucial step towards understanding the design of trustworthy Machine Learning systems.

ACKNOWLEDGEMENT

Valentina E. Balas would like to thank the European Research Development Fund under the Competitiveness Operational Program (BioCell-NanoART = Novel Bio-inspired Cellular Nano-architectures, POC-A1-A1.1.4-E nr. 30/2016) for supporting the research.

REFERENCES

- [1] Martin Abadi, Andy Chu, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Kunal Talwar, and Li Zhang. 2016. Deep Learning with Differential Privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security (CCS '16)*. ACM, New York, NY, USA, 308–318. <https://doi.org/10.1145/2976749.2978318>
- [2] V. Beiu, N. C. Rohatnovici, L. Dăuş, and V. E. Balas. 2017. Transport reliability on axonal cytoskeleton. In *2017 14th International Conference on Engineering of Modern Electric Systems (EMES)*. 160–163. <https://doi.org/10.1109/EMES.2017.7980404>
- [3] Jose Bernier, Julio Ortega, Eduardo Vidal, Ignacio Rojas, and Alberto Prieto. 2001. A Quantitative Study of Fault Tolerance, Noise Immunity, and Generalization Ability of MLPs. *Neural Computation* 12 (01 2001), 2941–2964. <https://doi.org/10.1162/089976600300014782>
- [4] Chris M. Bishop. 1995. Training with Noise is Equivalent to Tikhonov Regularization. *Neural Comput.* 7, 1 (Jan. 1995), 108–116. <https://doi.org/10.1162/neco.1995.7.1.108>
- [5] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. 2019. The Secret Sharer: Evaluating and Testing Unintended

- Memorization in Neural Networks. In *28th USENIX Security Symposium (USENIX Security 19)*. USENIX Association, Santa Clara, CA, 267–284. <https://www.usenix.org/conference/usenixsecurity19/presentation/carlini>
- [6] S. R. Cowell, V. Beiu, L. Dăuş, and P. Poulin. 2018. On the Exact Reliability Enhancements of Small Hammock Networks. *IEEE Access* 6 (2018), 25411–25426. <https://doi.org/10.1109/ACCESS.2018.2828036>
 - [7] D. Deodhare, M. Vidyasagar, and S. Sathya Keethi. 1998. Synthesis of fault-tolerant feedforward neural networks using minimax optimization. *IEEE Transactions on Neural Networks* 9, 5 (Sep. 1998), 891–900. <https://doi.org/10.1109/72.712162>
 - [8] P. Dey, K. Nag, T. Pal, and N. R. Pal. 2018. Regularizing Multilayer Perceptron for Robustness. *IEEE Transactions on Systems, Man, and Cybernetics: Systems* 48, 8 (Aug 2018), 1255–1266. <https://doi.org/10.1109/TSMC.2017.2664143>
 - [9] Vasisht Duddu. 2018. A Survey of Adversarial Machine Learning in Cyber Warfare. *Defence Science Journal* 68, 4 (Jun. 2018), 356–366. <https://doi.org/10.14429/dsj.68.12371>
 - [10] Vasisht Duddu, D Vijay Rao, and Valentina E Balas. 2019. Adversarial Fault Tolerant Training for Deep Neural Networks. *arXiv preprint arXiv:1907.03103* (2019).
 - [11] Cynthia Dwork. 2008. Differential Privacy: A Survey of Results. In *Theory and Applications of Models of Computation*, Manindra Agrawal, Dingzhu Du, Zhenhua Duan, and Angsheng Li (Eds.). Springer Berlin Heidelberg, Berlin, Heidelberg, 1–19.
 - [12] Cynthia Dwork and Aaron Roth. 2014. The Algorithmic Foundations of Differential Privacy. (2014). <http://www.cis.upenn.edu/~aaroht/Papers/privacybook.pdf>
 - [13] Christian Etmann, Sebastian Lunz, Peter Maass, and Carola Schoenlieb. 2019. On the Connection Between Adversarial Robustness and Saliency Map Interpretability. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), Vol. 97. PMLR, Long Beach, California, USA, 1823–1832. <http://proceedings.mlr.press/v97/etmann19a.html>
 - [14] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart. 2015. Model Inversion Attacks That Exploit Confidence Information and Basic Countermeasures. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security (CCS '15)*. ACM, New York, NY, USA, 1322–1333. <https://doi.org/10.1145/2810103.2813677>
 - [15] Karan Ganju, Qi Wang, Wei Yang, Carl A. Gunter, and Nikita Borisov. 2018. Property Inference Attacks on Fully Connected Neural Networks Using Permutation Invariant Representations. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (CCS '18)*. ACM, New York, NY, USA, 619–633. <https://doi.org/10.1145/3243734.3243834>
 - [16] K. I. Ho, C. Leung, and J. Sum. 2010. Convergence and Objective Functions of Some Fault/Noise-Injection-Based Online Learning Algorithms for RBF Networks. *IEEE Transactions on Neural Networks* 21, 6 (June 2010), 938–947. <https://doi.org/10.1109/TNN.2010.2046179>
 - [17] L. Holmstrom and P. Koistinen. 1992. Using additive noise in back-propagation training. *IEEE Transactions on Neural Networks* 3, 1 (Jan 1992), 24–38. <https://doi.org/10.1109/72.105415>
 - [18] Matthew Jagielski, Michael Kearns, Jieming Mao, Alina Oprea, Aaron Roth, Saeed Sharif Malvajerdi, and Jonathan Ullman. 2019. Differentially Private Fair Learning. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), Vol. 97. PMLR, Long Beach, California, USA, 3000–3008. <http://proceedings.mlr.press/v97/jagielski19a.html>
 - [19] Bargav Jayaraman and David Evans. 2019. Evaluating Differentially Private Machine Learning in Practice. In *28th USENIX Security Symposium (USENIX Security 19)*. USENIX Association, Santa Clara, CA, 1895–1912. <https://www.usenix.org/conference/usenixsecurity19/presentation/jayaraman>
 - [20] Jinyuan Jia, Ahmed Salem, Michael Backes, Yang Zhang, and Neil Zhenqiang Gong. 2019. MemGuard: Defending against Black-Box Membership Inference Attacks via Adversarial Examples. *arXiv preprint arXiv:1909.10594* (2019).
 - [21] Cong Jin and Shu-Wei Jin. 2014. Applications of Fuzzy Integrals for Predicting Software Fault-prone. *J. Intell. Fuzzy Syst.* 26, 2 (March 2014), 721–729. <http://dl.acm.org/citation.cfm?id=2596370.2596386>
 - [22] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. 2018. Towards Deep Learning Models Resistant to Adversarial Attacks. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rJzIBZAb>
 - [23] K. Matsuoka. 1992. Noise injection into inputs in back-propagation learning. *IEEE Transactions on Systems, Man, and Cybernetics* 22, 3 (May 1992), 436–440. <https://doi.org/10.1109/21.155944>
 - [24] Matthew Mirman, Timon Gehr, and Martin Vechev. 2018. Differentiable Abstract Interpretation for Provably Robust Neural Networks. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Jennifer Dy and Andreas Krause (Eds.), Vol. 80. PMLR, Stockholm, Sweden, 3578–3586. <http://proceedings.mlr.press/v80/mirman18b.html>
 - [25] I. Mironov. 2017. Rényi Differential Privacy. In *2017 IEEE 30th Computer Security Foundations Symposium (CSF)*. 263–275. <https://doi.org/10.1109/CSF.2017.11>
 - [26] Milad Nasr, Reza Shokri, and Amir Houmansadr. 2018. Machine Learning with Membership Privacy Using Adversarial Regularization. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security (CCS '18)*. ACM, New York, NY, USA, 634–646. <https://doi.org/10.1145/3243734.3243855>
 - [27] C. Neti, M. H. Schneider, and E. D. Young. 1992. Maximally fault tolerant neural networks. *IEEE Transactions on Neural Networks* 3, 1 (Jan 1992), 14–23. <https://doi.org/10.1109/72.105414>
 - [28] Nicolas Papernot, Shuang Song, Ilya Mironov, Ananth Raghunathan, Kunal Talwar, and Ulfar Erlingsson. 2018. Scalable Private Learning with PATE. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=rkZB1XbRZ>
 - [29] D. S. Phatak and I. Koren. 1995. Complete and partial fault tolerance of feedforward neural nets. *IEEE Transactions on Neural Networks* 6, 2 (March 1995), 446–456. <https://doi.org/10.1109/72.363479>
 - [30] Ahmed Salem, Apratim Bhattacharyya, Michael Backes, Mario Fritz, and Yang Zhang. 2019. Updates-Leak: Data Set Inference and Reconstruction Attacks in Online Learning. *CoRR* abs/1904.01067 (2019). [arXiv:1904.01067](https://arxiv.org/abs/1904.01067)
 - [31] Ahmed Salem, Yang Zhang, Mathias Humbert, Mario Fritz, and Michael Backes. 2019. ML-Leaks: Model and Data Independent Membership Inference Attacks and Defenses on Machine Learning Models. In *Annual Network and Distributed System Security Symposium (NDSS)* (2019-02-24), preliminary. <https://arxiv.org/abs/1806.01246>
 - [32] Reza Shokri, Martin Strobel, and Yair Zick. 2019. Privacy Risks of Explaining Machine Learning Models. *arXiv preprint arXiv:1907.00164* (2019).
 - [33] R. Shokri, M. Stronati, C. Song, and V. Shmatikov. 2017. Membership Inference Attacks Against Machine Learning Models. In *2017 IEEE Symposium on Security and Privacy (SP)*. 3–18. <https://doi.org/10.1109/SP.2017.41>
 - [34] Aman Sinha, Hongseok Namkoong, and John Duchi. 2018. Certifiable Distributional Robustness with Principled Adversarial Training. In *International Conference on Learning Representations*. <https://openreview.net/forum?id=Hk6kPgZA->
 - [35] Liwei Song, Reza Shokri, and Prateek Mittal. 2019. Privacy Risks of Securing Machine Learning Models against Adversarial Examples. *arXiv preprint arXiv:1905.10291* (2019).
 - [36] J. Wang, Q. Chang, Q. Chang, Y. Liu, and N. R. Pal. 2018. Weight Noise Injection-Based MLPs With Group Lasso Penalty: Asymptotic Convergence and Application to Node Pruning. *IEEE Transactions on Cybernetics* (2018), 1–19. <https://doi.org/10.1109/TCYB.2018.2864142>
 - [37] Yu-Xiang Wang, Jing Lei, and Stephen E. Fienberg. 2016. Learning with Differential Privacy: Stability, Learnability and the Sufficiency and Necessity of ERM Principle. *Journal of Machine Learning Research* 17, 183 (2016), 1–40. <http://jmlr.org/papers/v17/15-313.html>
 - [38] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric Xing, Laurent El Ghaoui, and Michael Jordan. 2019. Theoretically Principled Trade-off between Robustness and Accuracy. In *Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research)*, Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.), Vol. 97. PMLR, Long Beach, California, USA, 7472–7482. <http://proceedings.mlr.press/v97/zhang19p.html>