

CAT: Customized Adversarial Training for Improved Robustness

Minhao Cheng¹, Qi Lei², Pin-Yu Chen⁵, Inderjit Dhillon^{3,6} and Cho-Jui Hsieh⁴

¹Department of Computer Science and Engineering, HKUST

²Department of Electrical and Computer Engineering, Princeton University

³Department of Computer Science, UT Austin

⁴Department of Computer Science, UCLA

⁵IBM Research AI

⁶Amazon

minhaocheng@ust.hk, qilei@princeton.edu, pin-yu.chen@ibm.com, chohsieh@cs.ucla.edu, inderjit@cs.utexas.edu

Abstract

Adversarial training has become one of the most effective methods for improving robustness of neural networks. However, it often suffers from poor generalization on both clean and perturbed data. Current robust training method always use a uniform perturbation strength for every samples to generate adversarial examples during model training for improving adversarial robustness. However, we show it would lead worse training and generalization error and forcing the prediction to match one-hot label. In this paper, therefore, we propose a new algorithm, named Customized Adversarial Training (CAT), which adaptively customizes the perturbation level and the corresponding label for each training sample in adversarial training. We first show theoretically the CAT scheme improves the generalization. Also, through extensive experiments, we show that the proposed algorithm achieves better clean and robust accuracy than previous adversarial training methods. The full version of this paper is available at <https://arxiv.org/abs/2002.06789>.

1 Introduction

Deep neural networks (DNNs) have proved their effectiveness on a variety of domains and tasks. However, it has been found that DNNs are highly vulnerable to adversarial examples [Szegedy *et al.*, 2014]. To enhance the robustness of DNNs against adversarial examples, adversarial training [Goodfellow *et al.*, 2015; Madry *et al.*, 2018] has become one of the most effective and widely used methods. Given a pre-defined perturbation tolerance, denoted as ϵ , adversarial training aims to minimize the robust loss, defined as the worst-case loss within ϵ -ball around each example, leading to a min-max optimization problem. [Madry *et al.*, 2018] shows that applying a multi-step projected gradient descent (PGD) attack to approximately solve the inner maximization leads to a robust model, and several recent research has proposed various ways to improve adversarial training [Zhang *et al.*, 2019;

Wang, 2019; Wang *et al.*, 2019a; Balaji *et al.*, 2019; Ding *et al.*, 2018].

However, standard adversarial training methods still have a hypothetical and possibly problematic assumption: the perturbation tolerance ϵ is a large and fixed constant throughout the training process, which ignores the fact that every data point may have different intrinsic robustness. Intuitively, some examples are naturally closer to the decision boundary, and enforcing large margin on those examples will force the classifier to give up on those examples, leading to a distorted decision surface. This intuition may explain the known issue of the undesirable robustness-accuracy tradeoff in adversarial robustness [Su *et al.*, 2018; Tsipras *et al.*, 2019]. Furthermore, with a different perturbation tolerance, it is questionable whether we should still force the model to learn to fit the one-hot label as in the original adversarial training formulation. In the extreme case, if an example is perturbed to the decision boundary, a good classifier yielding the binary class prediction probabilities should output $[0.5, 0.5]$ instead of $[1, 0]$. This aspect becomes crucial when each example is associated with a different level of perturbation. Although some recent papers have started to address the uniform ϵ issue by treating correctly and incorrectly classified examples differently [Ding *et al.*, 2018] or assigning non-uniform perturbation level [Balaji *et al.*, 2019], none of them have tried to incorporate customized training labels in this process.

Motivated by these ideas, we propose a novel Customized Adversarial Training (CAT) framework that can substantially improve the performance of adversarial training. Throughout the adversarial training process, our algorithm dynamically finds a non-uniform and effective perturbation level and the corresponding customized target label for each example. This leads to better generalization performance and furthermore, with a careful design on adaptive ϵ tuning, our algorithm has only negligible computational overhead and runs as fast as the original adversarial training algorithm. Furthermore, we theoretically explain why the proposed method could lead to improved generalization performance.

Our method significantly outperforms existing adversarial training methods on the standard CIFAR-10 defense task.

With Wide-ResNet structure on CIFAR-10, under $8/255 \ell_\infty$ perturbation, our method achieves 73% robust accuracy under PGD attack and 71% robust accuracy under Carlini and Wagner (C&W) attack [Carlini and Wagner, 2017], while the current best model only achieves 58.6% under PGD attack and 56.8% under C&W attack. Furthermore, our method only degrades the clean accuracy from 95.93% (standard test accuracy) to 93.48%, while other adversarial training methods have clean accuracy below 91.34%.

2 Background and Motivation

2.1 Preliminaries

Adversarial training can be formulated as a min-max optimization problem. For a K -class classification problem, let $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1, \dots, n}$ denote the set of training samples in the dataset with $\mathbf{x}_i \in \mathbb{R}^d, y_i \in \{1, \dots, K\} := [K]$, we consider a classification model $f_\theta(\mathbf{x}) : \mathbb{R}^d \rightarrow [K]$ parameterized by θ . We denote by $h_\theta(\mathbf{x}) : \mathbb{R}^d \rightarrow [0, 1]^K$ as the prediction output for each class, i.e., $f_\theta(\mathbf{x}) = \operatorname{argmax}_i [h_\theta(\mathbf{x})]_i$.

Adversarial training can be formulated as:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon)} \ell(f_\theta(\mathbf{x}'_i), y_i), \quad (1)$$

where $\mathcal{B}_p(\mathbf{x}_i, \epsilon)$ denotes the ℓ_p -norm ball centered at \mathbf{x}_i with radius ϵ . The inner maximization problem aims to find an adversarial version of a given data point \mathbf{x}_i that achieves the highest loss. In general, one can define $\mathcal{B}_p(\mathbf{x}_i, \epsilon)$ based on the threat model, but the ℓ_∞ ball is the most popular choice adopted by recent works [Madry *et al.*, 2018; Zhang *et al.*, 2019; Wang, 2019; Ding *et al.*, 2018; Wang *et al.*, 2019b], which we also use in this paper. For a deep neural network model, the inner maximization does not have a closed form solution, so adversarial training methods typically use a gradient-based iterative solver to approximately solve the inner problem. The most commonly used choice is the multi-step PGD [Madry *et al.*, 2018] and C&W attack [Carlini and Wagner, 2017].

2.2 Motivation

Intuitively, if adversarial training can always find a model with close-to-zero robust error, one should always use a large ϵ for training because it will automatically imply robustness to any smaller ϵ . Unfortunately, in practice a uniformly large ϵ is often harmful. In the following we empirically explain this problem and use it to motivate our proposed algorithm.

We use a simple linear classification case to demonstrate why a uniformly large ϵ is harmful. In Figure 1a, we generate a synthetic linearly separable dataset with the margin set to be 1.75 for both classes, where the correct linear boundary can be easily obtained by standard training. In Figure 1b, we run adversarial training with $\epsilon = 1$, and since this ϵ is smaller than the margin, the algorithm can still obtain near-optimal results. However, when we use a large $\epsilon = 4$ for adversarial training in Figure 1c, the resulting decision boundary becomes significantly worse. It is because adversarial training cannot correctly fit all the samples with a margin up to 4, so it will sacrifice some data samples, leading to distorted and

Testing ϵ	Error Type	Training ϵ		
		0.01	0.02	0.03
0.01	Train	99.96%	99.99%	99.16%
	Test	69.79%	69.06%	66.04%

Table 1: The influence of different fixed ϵ values used in adversarial training on the robust accuracy with $\epsilon = 0.01$.

undesirable decision boundary. This motivates the following two problems:

- We shouldn't set the same large ϵ uniformly for all samples. Some samples are intrinsically closer to the decision boundary and they should use a smaller ϵ . Without doing this, adversarial training will give up on those samples, which leads to worse training and generalization error (see more discussions in Section 3.3 on the generalization bounds).
- The adversarial training loss is trying to force the prediction to match the one-hot label (e.g., $[1, 0]$ in the binary classification case) even after large perturbations. However, if a sample is perturbed, the prediction shouldn't remain one-hot. For instance, if a sample is perturbed to the decision boundary of a binary classification problem, the prediction of a perfect model should be $[0.5, 0.5]$ instead of $[1, 0]$, which also makes adversarial training fail to recover a good decision hyperplane.

Furthermore, we observe that even if adversarial training can obtain close-to-zero training error with large ϵ (e.g., [Gao *et al.*, 2019] proves that this will happen for overparameterized network with large-enough margin), a uniformly large ϵ will lead to larger generalization gap. This could be partially explained by the theoretical results provided by [Yin *et al.*, 2018], which shows that the adversarial Rademacher complexity has a lower bound with an explicit dependence on the perturbation tolerance. The empirical results in Table 1 also illustrate this problem. When conducting adversarial training with $\epsilon = 0.3$ on CIFAR10 VGG-16, we found that the model achieves close-to-zero robust training error on all $\epsilon \leq 0.3$, but it suffers larger generalization gap compared to training with smaller ϵ . This also demonstrates that a uniformly large ϵ is harmful even when it achieves perfect training error.

3 CAT (Customized Adversarial Training)

In this section, we propose the Customized Adversarial Training (CAT) framework that improves adversarial training by addressing the above-mentioned problems. First, our algorithm has an auto-tuning method to customize the ϵ used for each training example. Second, instead of forcing the model to fit the original label, we customize the target label for each example based on its own ϵ . In the following we will describe these two components in more detail.

3.1 Auto-tuning Perturbation Strength for Adversarial Training

The first component of our algorithm is an ϵ auto-tuning method which adaptively assigns a suitable ϵ for each example during the adversarial training procedure. Let ϵ_i be the

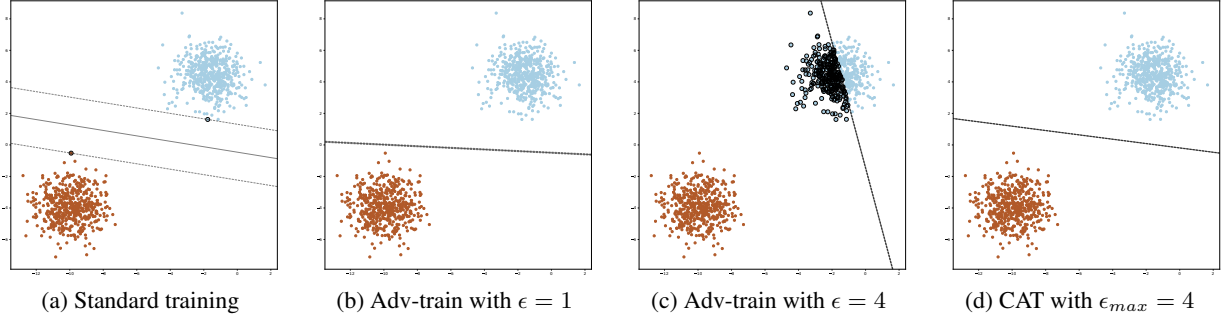


Figure 1: Different training methods on a linearly separable binary classification dataset with 1.75 margin for both classes. Adversarial training with small ϵ works fine, but for a large ϵ beyond the true margin, adversarial training would ruin the classifier’s classification performance, while our proposed adaptive customized adversarial training method still keeps a good generalization performance.

perturbation level assigned to example i . Based on the intuition mentioned in Section 2.2, we do not want to further increase ϵ if we find the classifier does not have capacity to robustly classify the example, which means we should set

$$\epsilon_i = \operatorname{argmin}_{\epsilon} \left\{ \max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon)} f_{\theta}(\mathbf{x}'_i) \neq y_i \right\} \quad (2)$$

and the adversarial training objective becomes

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon_i)} \ell(f_{\theta}(\mathbf{x}'_i), y_i). \quad (3)$$

Note that ϵ_i in (2) depends on θ , while θ in (3) also depends on ϵ_i . We thus propose an alternative update scheme — conducting one SGD update on θ , and then updating the ϵ_i in the current batch. However, finding ϵ_i exactly requires brute-force search for every possible value, which adds significant computational overhead to adversarial training.

Therefore, we only conduct a simplified update rule on ϵ_i as follows. Starting from an initial perturbation level of zero, at each iteration we conduct adversarial attack (e.g., PGD attack) with perturbation tolerance $\epsilon_i + \eta$ where η is a constant. If the attack is successful, then we reset current ϵ_i to 0 to encourage model learning a more robust classifier towards those examples. While if the attack is unsuccessful, which means an attacker still cannot find an adversarial example that satisfies $\max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon_i + \eta)} f_{\theta}(\mathbf{x}'_i) \neq y_i$, then we increase $\epsilon_i = \epsilon_i + \eta$. The attack results will also be used to update the model parameter θ , so this adaptive scheme does not require any additional cost. In practice, we also have an upper bound on the final perturbation to ensure that ϵ_i remains bounded for each i .

3.2 Adaptive Label Uncertainty for Adversarial Training

As mentioned in Section 2.2, the standard adversarial training loss is trying to enforce a sample being classified as the original one-hot label after ϵ perturbation. However, this may not be ideal. In the extreme case, if a sample is perturbed to the decision boundary, the prediction must be far away from one-hot. This problem is more severe when using non-uniform ϵ_i , since each different ϵ_i will introduce a different bias to the

loss, and that may be one of the reasons that purely adaptive ϵ -scheduling does not work well (see our ablation study in Section 5 and also the results reported in [Balaji *et al.*, 2019]).

In the following, we propose an adaptive label smoothing approach to reflect different perturbation tolerance on each example. Szegedy [2016] introduced label smoothing that converts one-hot label vectors into one-warm vectors representing low-confidence classification, in order to prevent the model from making over-confident predictions. Specifically, with a one-hot encoded label y , the smoothed version is

$$\tilde{y} = (1 - \alpha)y + \alpha u,$$

where $\alpha \in [0, 1]$ is the hyperparameter to control the smoothing level. In the adaptive setting, we set $\alpha = c\epsilon_i$ so that a larger perturbation tolerance would receive a higher label uncertainty and c is a hyperparameter. A common choice of u is $u = \frac{1}{K}$. However, this strict requirement tries to enforce every other labels having the same probability, which may not make sense in practice. On the other hand, as shown Section 2.2, adversarial training is easy to overfit and generate a large generalization gap. To better address these issues, we sample from a distribution instead. Specifically, we use $u = \text{Dirichlet}(\beta)$ where $\text{Dirichlet}(\cdot)$ refers to the Dirichlet distribution and $\beta \in \mathbb{R}^K$ is concentration hyperparameter. With different perturbation tolerance, the adaptive version of label smoothing is

$$\tilde{y}_i = (1 - c\epsilon_i)y_i + c\epsilon_i \text{Dirichlet}(\beta). \quad (4)$$

The Final Objective Function. Combining the two aforementioned techniques, our Customized Adversarial Training (CAT) method attempts to minimize the following objective:

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon_i)} \ell(f_{\theta}(\mathbf{x}'_i), \tilde{y}_i) \\ \text{s.t.} \quad & \epsilon_i = \operatorname{argmin}_{\epsilon} \left\{ \max_{\mathbf{x}'_i \in \mathcal{B}_p(\mathbf{x}_i, \epsilon)} f_{\theta}(\mathbf{x}'_i) \neq y_i \right\}, \end{aligned} \quad (5)$$

where \tilde{y}_i is defined in (4). As described in Section 3.1, we approximately minimize this objective with an alternative update scheme, which incurs almost no additional cost compared to the original adversarial training algorithm. The detailed algorithm is presented in Algorithm 1.

Algorithm 1 CAT algorithm

Input: Training dataset (X, Y) , cross entropy loss or mix loss ℓ , scheduling parameter η , weighting factor c , perturbation upperbound ϵ_{max}
 Initial every sample's ϵ_i with 0
for epoch=1, ..., N **do**
 for $i=1, \dots, B$ **do**
 $\tilde{y}_i \leftarrow (1 - c\epsilon_i)y_i + c\epsilon_i \text{Dirichlet}(\beta)$
 $\epsilon_i \leftarrow \epsilon_i + \eta$
 $\delta_i \leftarrow 0$
 for $j = 1 \dots m$ **do**
 $\delta_i \leftarrow \delta_i + \alpha \cdot \text{sign}(\nabla_{\delta} \ell(f_{\theta}(\mathbf{x}_i + \delta_i), \tilde{y}_i))$
 $\delta_i \leftarrow \max(\min(\delta_i, \epsilon_i), -\epsilon_i)$
 end for
 if $f_{\theta}(\mathbf{x}_i + \delta_i) \neq y_i$ **then**
 $\epsilon_i \leftarrow 0$
 end if
 $\tilde{\epsilon}_i \leftarrow \min(\epsilon_{max}, \epsilon_i)$
 $\tilde{y}_i \leftarrow (1 - c\tilde{\epsilon}_i)y_i + c\tilde{\epsilon}_i \text{Dirichlet}(\beta)$
 $\theta \leftarrow \theta - \gamma_{\theta} \nabla_{\theta} \ell(f_{\theta}(\mathbf{x}_i + \delta_i), \tilde{y}_i)$
 end for
end for
return θ

Choice of Loss Function. In general, our framework can be used with any loss function $\ell(\cdot)$. In the previous works, cross entropy loss is commonly used for ℓ . However, the model trained by smoothing techniques tends to have a smaller logit gap between true label and other labels. Therefore, in order to encourage model to generate a larger logit gap, we propose a mixed loss to enhance the defense performance towards C&W $_{\infty}$ attack. That is,

$$\text{CE}(f_{\theta}(\mathbf{x}'_i), \tilde{y}_i) + \max\{\max_{j \neq y_0} \{Z(\mathbf{x}'_i)_j - [Z(\mathbf{x}'_i)]_{y_0}\}, -\kappa\},$$

where $Z(\mathbf{x}) \in \mathbb{R}^K$ is the final (logit) layer output, and $[Z(\mathbf{x})]_i$ is the prediction score for the i -th class and y_0 is the original label. The parameter κ encourages the adversary to find higher confident adversarial examples in training.

3.3 Theoretical Analysis

To better understand how our scheme improves generalization, we now provide some theoretical analysis. Recall we denote by $h_{\theta}(\mathbf{x}) : \mathbb{R}^d \rightarrow [0, 1]^K$ as the prediction probability for the K classes. We define the bilateral margin that our paper is essentially maximizing over as follows.

Definition 1 (Bilateral margin). *We define the bilateral perturbed network output by $H_{\theta}(\mathbf{x}, \delta^i, \delta^o)$:*

$$H_{\theta}(\mathbf{x}, \delta^i, \delta^o) := h_{\theta}(\mathbf{x} + \delta^i \|\mathbf{x}\|) + \|\mathbf{x} + \delta^i \|\mathbf{x}\|\| \cdot \delta^o.$$

The bilateral margin is now defined as the minimum norm of (δ^i, δ^o) required to cause the classifier to make false predictions:

$$\begin{aligned} m_F(\mathbf{x}, y) &:= \min_{\delta^i, \delta^o} \sqrt{\|\delta^i\|^2 + \|\delta^o\|^2} \\ \text{s.t. } &\max_{y'} H_{\theta}(\mathbf{x}, \delta^i, \delta^o)_{y'} \neq y. \end{aligned} \quad (6)$$

This margin captures both the relative perturbation on the input layer δ^i and the soft-max output δ^o .

Theorem 2. *Suppose the parameter space Θ we optimize over has covering number that scales as $\log \mathcal{N}_{\|\cdot\|_{op}}(\eta, \Theta) \leq \lfloor \mathcal{C}^2 / \eta^2 \rfloor$ for some complexity \mathcal{C} . Then with probability $1 - \delta$ over the draw of the training data, any classifier $f_{\theta}, \theta \in \Theta$ which achieves training error zero satisfies:*

$$\mathbb{E}[f_{\theta}(\mathbf{x}) = y] \lesssim \frac{\mathcal{C} \log^2 n}{\sqrt{n}} \sqrt{\frac{1}{n} \sum_{i=1}^n \frac{1}{m_F(\mathbf{x}_i, y_i)}} + \zeta,$$

where ζ is of small order $O(\frac{1}{n} \log(1/\delta))$.

We defer the proof to the Appendix, which is adapted from Theorem 2.1 of [Wei and Ma, 2019]. We observe the population risk is bounded by two key factors, the average of $\frac{1}{m_F(\mathbf{x}_i, y_i)}$ and \mathcal{C} , the covering number of the parameter space. On one side, the average of $\frac{1}{m_F(\mathbf{x}_i, y_i)}$ is dominated by the samples with the smallest margin. Therefore when we do adversarial training, it is important that we not only achieve higher overall accuracy, but also make sure the samples closer to the decision boundary have large enough margin. This can not be achieved by simply using constant and large ϵ that will maintain a large margin for most samples but sacrifice the accuracy of a small portion of data. On the other hand, the covering number of the network's parameter space can be roughly captured by a bound of product of all layers' weight norms. We hypothesize that with more flexibility in choosing ϵ , our algorithm will converge faster than using larger constant ϵ and will have more implicit regularization effect. To testify this hypothesis, we roughly measure the model complexity \mathcal{C} by the product of the weight norms of different models. In comparison to our model, when training with constant $\epsilon = 0.01, 0.02$ and 0.03 , it respectively yields \mathcal{C} as large as 2.54, 3.53 and 1.39 times of that of our model, which means our model indeed has more implicit regularization effect among others.

4 Related Work

Adversarial Training. To enhance the adversarial robustness of a neural network model, a natural idea is to iteratively generate adversarial examples, add them back to the training data, and retrain the model. For example, [Goodfellow et al., 2015] use adversarial examples generated by FGSM to augment the data, and [Kurakin et al., 2017] propose to use a multiple-step FGSM to further improve the performance. [Madry et al., 2018] show that adversarial training can be formulated as a min-max optimization problem, and propose to use PGD attack (similar to multi-step FGSM) to find adversarial examples for each batch. After that, many defense algorithms are based on a similar min-max framework. However, each of them uses slightly different loss functions.

We see that except for natural training which directly minimizes the cross entropy loss (denoted as CE), all training techniques involve the use of the min-max framework. TRADES and MMA use the unperturbed data's cross entropy loss as an additional regularization term to achieve a better trade-off between clean and robust error.

Similar to our method, both MMA and IAAT have sample-wise adaptive ϵ during training. They also utilize the adaptive ϵ to find the largest possible ϵ_i for every sample x_i . However, they do not consider the adaptive label technique mentioned in Section 3.2. As a result, they can only achieve better clean accuracy while the improvements in robust accuracy are limited. Our CAT algorithm (CAT with CE loss) is more general than IAAT and MMA. CAT reduce to IAAT when we set $c = 0$ in adaptive label smoothing. Moreover, MMA could be treated as a special case of CAT when we use a line search scheme to find the ϵ_i and $c = 0$. Also, in Section 5, we will show the importance of the adaptive label uncertainty step in CAT. Recently, a concurrent work [Stutz *et al.*, 2019] combines label smoothing with adversarial training. However, in the adversarial training process, they still use the same ϵ for all the examples, which is quite different from our instance-wise auto-tuning ϵ . Since their model jointly performs detection (dropping low confident examples) and prediction, the results and formulation are not directly comparable to other adversarial training methods.

5 Performance Evaluation

In this section, we conduct extensive experiments to show that CAT achieves a strong result on both clean and robust accuracy. We include the following methods into our comparison:

- Customized Adversarial Training (CAT): Our proposed method.
- Adversarial training: The adversarial training method proposed in [Madry *et al.*, 2018] where they use a K-step PGD attack as adversary.
- TRADES: TRADES [Zhang *et al.*, 2019] improves adversarial training by an additional loss on the clean examples and achieves the state-of-art performance on robust accuracy.
- Natural: the natural training which only minimizes the cross entropy loss.

Furthermore, since many recently proposed adversarial training methods have considered CIFAR-10 with Wide-ResNet structure as the standard setting for reporting their numbers, we also compare our performance with 7 previous methods on this specific setting.

5.1 Experimental Setup

Dataset and Model Structure. We use CIFAR-10 dataset for performance evaluation. We use both standard VGG-16 [Simonyan and Zisserman, 2015] and Wide ResNet that is used in both vanilla adversarial training [Madry *et al.*, 2018] and TRADES [Zhang *et al.*, 2019]. For VGG-16, we implement adversarial training with the standard hyper-parameters and train TRADES with the official implementation. For Wide ResNet, since the model has become standard for testing adversarial training methods, we use exactly the same model structure provided by [Madry *et al.*, 2018; Zhang *et al.*, 2019]. We use the models’ checkpoint released by TRADES official repository and implement the Madry’s adversarial training using the standard hyper-parameters. All

our experiments were implemented in Pytorch-1.4 and conducted using a GTX 2080 TI GPU.

Implementation Details. We set the number of iterations in adversarial attack to be 10 for all methods during training. Adversarial training and TRADES are trained on PGD attacks setting $\epsilon = 8/255$ with cross entropy loss (CE). All the models are trained using SGD with momentum 0.9, weight decay 5×10^{-4} . For VGG-16/Wide ResNet models, we use the initial learning rate of 0.01/0.1, and we decay the learning rate by 90% at the 80th, 140th, and 180th epoch. For CAT, we set epsilon scheduling parameter $\eta = 0.005$, $\epsilon_{max} = 8/255$ and weighting parameter $c = 10$. We set $\beta = 1$ for the distribution Dirichlet(β), which is equal to a uniform distribution. Also, we set $\kappa = 10$. Our code is publicly available at <https://github.com/cmhcbb/CAT-Customized-Adversarial-Training-for-Improved-Robustness>.

5.2 Robustness Evaluation and Analysis

White-box Attacks Results. For CIFAR10, we evaluate all the models under white-box $\epsilon = 8/255$ ℓ_∞ -norm bounded non-targeted PGD and C&W attack. Specifically, we use both PGD^X (X-step PGD with step size $\epsilon/5$) and C&W _{∞} . As suggested, we test our model under different steps PGD and multiple random restarts.

The experimental results are shown in Table 2, where we can easily see that CAT clearly outperforms other methods. CAT achieves a significant better robust accuracy at the standard $8/255$ perturbation threshold considered in the literature, and also have better clean accuracy. We also test the performance of CAT under attacks with 20 restarts and 1,000 iterations to confirm the robustness of the model. Furthermore, we visualize the loss landscape in the ablation study.

Wide-ResNet has become a standard structure for comparing adversarial training methods, and it’s standard to train and evaluate with $8/255$ ℓ_∞ norm perturbation. For this setting, we collect the reported accuracy from 7 other adversarial training methods, with several of them published very recently, to have a detailed full comparison. As shown in Table 3, our method achieves state-of-art robust accuracy while maintaining a high clean accuracy.

Black-box Transfer Attacks Results. We follow the criterion of evaluating transfer attacks as suggested by [Athalye *et al.*, 2018] to inspect whether the models trained by CAT will cause the issue of obfuscated gradients and give a false sense of model robustness. We generate 10,000 adversarial examples of CIFAR-10 from natural models with $\epsilon = 8/255$ and evaluate their attack performance on the target model. Table 4 shows that CAT achieves the best accuracy compared with adversarial training and TRADES, suggesting the effectiveness of CAT in defending both white-box and transfer attacks.

5.3 Ablation Study

The Importance of Adaptive Label Uncertainty. Here we discuss and perform an ablation study using VGG-16 and CIFAR-10 on the importance of adaptive label uncertainty and adaptive instance-wise ϵ . In Table 5, Adv train denotes the original adversarial training, Adv+LS denotes adversarial training with label smoothing (setting y by Eq (4)), Adp-Adv denotes adversarial training with adaptive instance-wise

Methods	No attack	Deepfool	PGD ¹⁰⁰	C&W ¹⁰⁰	20 PGD ¹⁰⁰⁰	20 C&W ¹⁰⁰⁰
Natural train	93.34%	16.39%	0.6%	0.0%	0.0%	0.0%
Adv train [Madry <i>et al.</i> , 2018]	80.32%	44.65%	36.36%	37.89%	36.12%	36.8%
TRADES [Zhang <i>et al.</i> , 2019]	84.85%	48.37%	38.81%	39.49%	37.95%	38.94%
CAT (ours)	85.44%	70.19%	75.54%	51.81%	75.17%	50.08%

Table 2: The clean and robust accuracy of VGG-16 models trained by various defense methods. All robust accuracy results use $\epsilon = 8/255$ ℓ_∞ ball. ^(X) denotes using a X -step PGD attack. X random denotes X times random restart.

Methods	Clean accuracy	PGD accuracy	C&W accuracy
Natural training	95.93%	0%	0%
Adversarial training [Madry <i>et al.</i> , 2018]	87.30%	52.68%	50.73%
Dynamic adversarial training [Wang <i>et al.</i> , 2019a]	84.51%	55.03%	51.98%
TRADES [Zhang <i>et al.</i> , 2019]	84.22%	56.40% ⁽²⁰⁾	51.98%
Bilateral Adv Training [Wang, 2019]	91.00%	57.5% ^(*20)	56.2% ^(*20)
MMA [Ding <i>et al.</i> , 2018]	84.36%	47.18%	X
MART [Wang <i>et al.</i> , 2019b]	84.17%	58.56% ⁽²⁰⁾	54.58%
IAAT [Balaji <i>et al.</i> , 2019]	91.34%	48.53% ^(*10)	56.80%
CAT (ours)	89.61%	73.16% ^(*20)	71.67% ^(*20)

Table 3: The clean and robust accuracy of Wide Resnet models trained by various defense methods. All robust accuracy results use $\epsilon = 8/255$ ℓ_∞ ball. We reported the best performance listed in the papers. ^(*) denotes random-restart is applied in the testing attack. ^(X) denotes using a X -step PGD attack. **X** denotes not reported.

Method	VGG 16	Wide ResNet
Adv train	79.13%	85.84%
TRADES	83.53%	83.90%
CAT	86.58%	88.66%

Table 4: Robust accuracy under transfer attack on CIFAR-10

Methods	Clean Acc	PGD Acc
Adv train	80.32%	36.63%
Adv+LS	80.25%	43.0%
Adp-Adv	87.91%	38.59%
CAT	84.22%	75.54%

Table 5: Ablation study on CAT by changing the loss function and removing Label Adaption (LA). All robust accuracy results use $\epsilon = 8/255$ ℓ_∞ ball.

ϵ , and CAT is the proposed method which is a combination of these two tricks. We found that only applying adaptive instance-wise ϵ or label smoothing cannot significantly boost the robust accuracy over standard adversarial training, but the proposed method, by nicely combining these two ideas, can significantly improve the performance.

Loss Landscape Exploration. To further verify the superior robustness using CAT, we visualize the loss landscape of different training methods in Figure 2. Following the implementation in [Engstrom *et al.*, 2018], we divide the data input along a linear space grid defined by the sign of the input gradient and a random Rademacher vector, where the x - and y - axes represent the magnitude of the perturbation added in each direction and the z -axis represents the loss. As shown in Figure 2, CAT generates a model with a lower and smoother loss landscape. Also, it could be taken as another strong evidence that we have found a robust model through CAT.

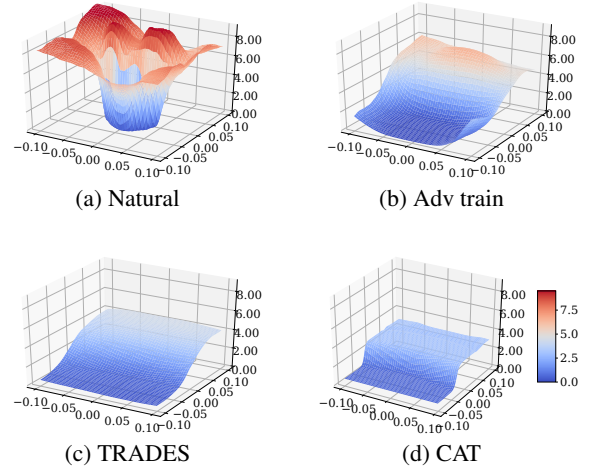


Figure 2: Loss landscape comparison of different adv training methods

6 Conclusions

In this paper, we propose CAT, a customized adversarial training method that is designed to have better generalization for both clean and robust performance. We also provide a theoretical analysis that explains the performance of our algorithm. Experimental results show that CAT has achieved state-of-art robust accuracy and a high clean accuracy while keeping similar running time as standard adversarial training. The success of CAT indicates that it is crucial to customize the perturbation level on both data sample side and its label in adversarial training.

References

- [Athalye *et al.*, 2018] Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *International Conference on Machine Learning*, 2018.
- [Balaji *et al.*, 2019] Yogesh Balaji, Tom Goldstein, and Judy Hoffman. Instance adaptive adversarial training: Improved accuracy tradeoffs in neural nets. *arXiv preprint arXiv:1910.08051*, 2019.
- [Carlini and Wagner, 2017] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *IEEE Symposium on Security and Privacy*, pages 39–57, 2017.
- [Ding *et al.*, 2018] Gavin Weiguang Ding, Yash Sharma, Kry Yik Chau Lui, and Ruitong Huang. Max-margin adversarial (mma) training: Direct input space margin maximization through adversarial training. *arXiv preprint arXiv:1812.02637*, 2018.
- [Engstrom *et al.*, 2018] Logan Engstrom, Andrew Ilyas, and Anish Athalye. Evaluating and understanding the robustness of adversarial logit pairing. *arXiv preprint arXiv:1807.10272*, 2018.
- [Gao *et al.*, 2019] Ruiqi Gao, Tianle Cai, Haochuan Li, Chojui Hsieh, Liwei Wang, and Jason D Lee. Convergence of adversarial training in overparametrized neural networks. In *Advances in Neural Information Processing Systems*, pages 13009–13020, 2019.
- [Goodfellow *et al.*, 2015] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations*, 2015.
- [Kurakin *et al.*, 2017] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale. *International Conference on Learning Representations*, 2017.
- [Madry *et al.*, 2018] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *International Conference on Learning Representations*, 2018.
- [Simonyan and Zisserman, 2015] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*, 2015.
- [Stutz *et al.*, 2019] David Stutz, Matthias Hein, and Bernt Schiele. Confidence-calibrated adversarial training: Towards robust models generalizing beyond the attack used during training. *arXiv preprint arXiv:1910.06259*, 2019.
- [Su *et al.*, 2018] Dong Su, Huan Zhang, Hongge Chen, Jinfeng Yi, Pin-Yu Chen, and Yupeng Gao. Is robustness the cost of accuracy?—a comprehensive study on the robustness of 18 deep image classification models. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 631–648, 2018.
- [Szegedy *et al.*, 2014] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *International Conference on Learning Representations*, 2014.
- [Szegedy *et al.*, 2016] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016.
- [Tsipras *et al.*, 2019] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *International Conference on Learning Representations*, 2019.
- [Wang *et al.*, 2019a] Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanquan Gu. On the convergence and robustness of adversarial training. In *International Conference on Machine Learning*, pages 6586–6595, 2019.
- [Wang *et al.*, 2019b] Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, and Quanquan Gu. Improving adversarial robustness requires revisiting misclassified examples. In *International Conference on Learning Representations*, 2019.
- [Wang, 2019] Jianyu Wang. Bilateral adversarial training: Towards fast training of more robust models against adversarial attacks. *International Conference on Computer Vision*, 2019.
- [Wei and Ma, 2019] Colin Wei and Tengyu Ma. Improved sample complexities for deep networks and robust classification via an all-layer margin. *arXiv preprint arXiv:1910.04284*, 2019.
- [Yin *et al.*, 2018] Dong Yin, Kannan Ramchandran, and Peter Bartlett. Rademacher complexity for adversarially robust generalization. *arXiv preprint arXiv:1810.11914*, 2018.
- [Zhang *et al.*, 2019] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I Jordan. Theoretically principled trade-off between robustness and accuracy. *International Conference on Machine Learning*, 2019.