

A Appendix

A.1 Restricted Imagenet Result

In addition to CIFAR, we also test the performance on the Restricted Imagenet dataset, which has been used in previous papers such as (Sinha et al., 2019). This dataset consists of a subset of imagenet classes which have been grouped into 9 different classes. The experimental results are shown in Table 5. All the results except our methods are reported in (Sinha et al., 2019).

Table 7: The clean and robust accuracy of Resnet50 models trained by various defense methods. All robust accuracy are measured under $\epsilon = 8/255 \ell_\infty$ ball. We reported the best performance listed in the papers. ^(*) denotes random-restart is applied in the testing attack. ^(X) denotes it use a X -step PGD attack

Methods	Clean accuracy	PGD accuracy	C&W accuracy
Adversarial training [20]	91.83%	17.52%	X
FAT (Sinha et al., 2019)	91.59 %	18.81%	X
LAT (Sinha et al., 2019)	89.86 %	22.00%	X
CAT (ours)	88.63%	58.4% ^(*20)	58.4% ^(*20)

A.2 Omitted Proofs

In this section we provide the omitted proof for Theorem 3.2, which is adapted from Theorem 2.1 from [33]. They defined the all layer margin for a k -layer network $h_\theta(\mathbf{x}) = f_k \circ f_{k-1} \circ \dots \circ f_1(\mathbf{x})$ and perturbation $\delta = (\delta_1, \delta_2, \dots, \delta_k)$ as follows:

$$\begin{aligned} h_1(\mathbf{x}, \delta) &= f_1(\mathbf{x}) + \delta_1 \|\mathbf{x}\|_2 \\ h_i(\mathbf{x}, \delta) &= f_i(h_{i-1}(\mathbf{x}, \delta)) + \delta_i \|h_{i-1}(\mathbf{x}, \delta)\|_2 \\ H_\theta(\mathbf{x}, \delta) &= h_k(\mathbf{x}, \delta). \end{aligned}$$

They define the all-layer margin as the minimum norm of $\delta = (\delta_i)_{i=1}^k$ required that causes the classifier to make a false prediction.

$$\begin{aligned} m_F(\mathbf{x}, y) &:= \min_{\delta^i, \delta^o} \sqrt{\|\delta^i\|^2 + \|\delta^o\|^2} \\ &\text{subject to } \max_{y'} H_\theta(\mathbf{x}, \delta^i, \delta^o)_{y'} \neq y. \end{aligned} \tag{8}$$

They consider the function class $\mathcal{F} = \{f_k \circ f_{k-1} \circ \dots \circ f_1 : f_i \in \mathcal{F}_i\}$ be the class of compositions of functions from function classes $\mathcal{F}_1, \dots, \mathcal{F}_k$. They achieve the generalization bound as follows:

Theorem A.1 (Theorem 2.1 from [33]) *In the above setting, with probability $1 - \delta$ over the draw of the data, all classifiers $F \in \mathcal{F}$ which achieve training error 0 satisfy*

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

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

For our problem, we define $h_\theta(\mathbf{x}) := f_2 \circ f_1(\mathbf{x})$, where f_1 is identity mapping, and f_2 is the original function h_θ . Therefore the all layer margin is reduced to our bilateral margin:

$$\begin{aligned} h_1(\mathbf{x}, \delta^i) &= f_1(\mathbf{x}) + \delta^i \|\mathbf{x}\|_2 = \mathbf{x} + \delta^i \|\mathbf{x}\| \\ H_\theta(\mathbf{x}, \delta) &= h_2(\mathbf{x}, \delta^i, \delta^o) = f_2(h_1(\mathbf{x}, \delta^i)) + \delta^o \|h_1(\mathbf{x}, \delta^i)\| \\ &= h_\theta(\mathbf{x} + \delta^i \|\mathbf{x}\|) + \delta^o \|\mathbf{x} + \delta^i \|\mathbf{x}\|. \end{aligned}$$

Next, notice since f_1 is identity mapping, and composition with h_θ doesn't affect the overall complexity. We apply Theorem A.1 and get our result.

420 A.3 Performance under different PGD strength

421 For CIFAR10, we evaluate all the models under different tolerance of white-box ℓ_∞ -norm bounded
 422 non-targeted PGD and C&W attack.

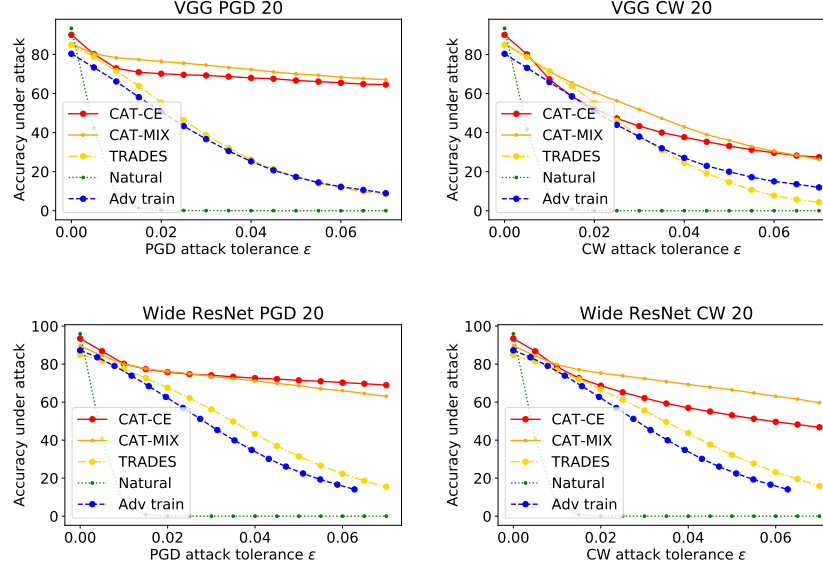


Figure 2: Robust accuracy under different levels of attacks on CIFAR-10 datasets with VGG and Wide-ResNet architectures. CAT-CE and CAT-MIX clearly outperform TRADES and adversarial training.

423 A.4 Loss Landscape Exploration

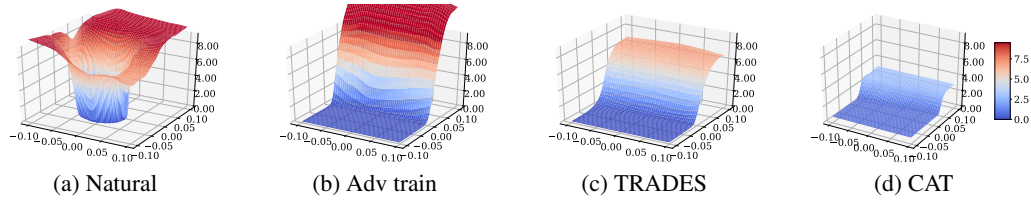


Figure 3: Loss landscape comparison of different adversarial training methods

424 To further verify the superior robustness using CAT, we visualize the loss landscape of different
 425 training methods in Figure 3. Following the implementation in [11], we divide the data input along a
 426 linear space grid defined by the sign of the input gradient and a random Rademacher vector, where
 427 the x- and y- axes represent the magnitude of the perturbation added in each direction and the z-axis
 428 represents the loss. As shown in Figure 3, CAT generates a model with a lower and smoother loss
 429 landscape. Also, it could be taken as another strong evidence that we have found a robust model
 430 through CAT training.