

# Towards Understanding the Regularization of Adversarial Robustness on Neural Network

Yuxin Wen\*, Shuai Li\*, Kui Jia

South China University of Technology (SCUT)



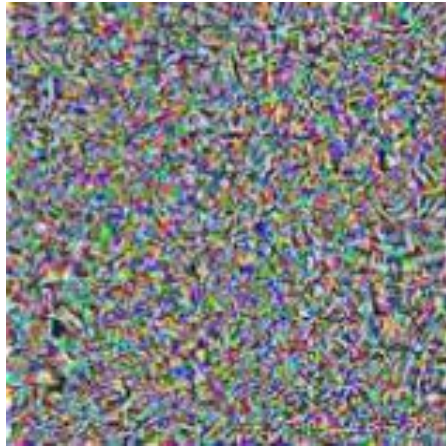
\* indicates equal contribution

# Adversarial examples



pig

+



texture changes

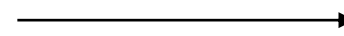
=



airliner

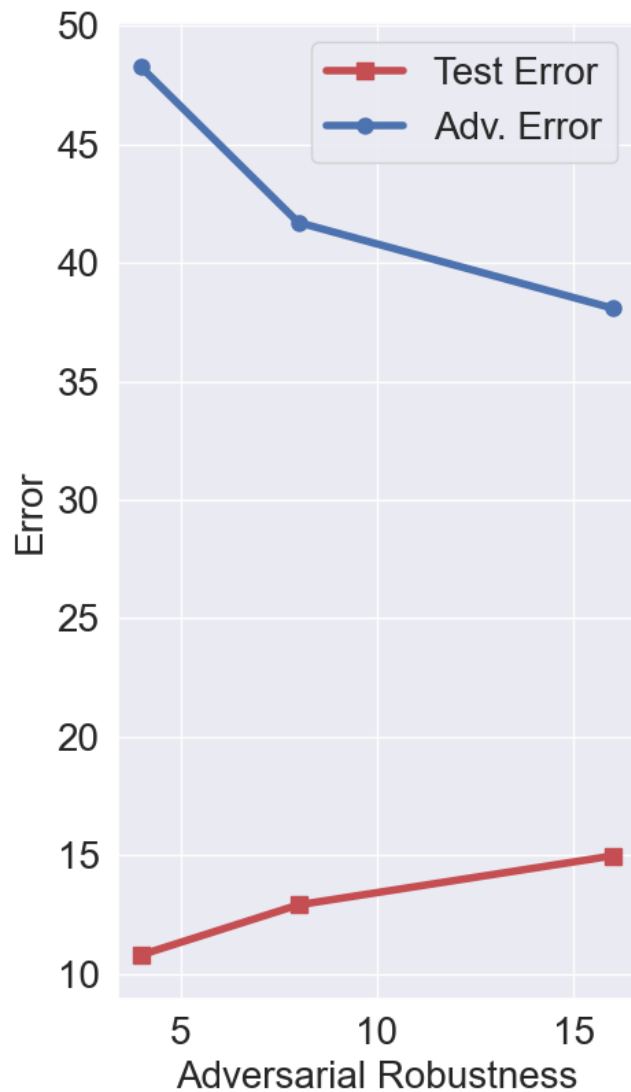
# Adversarial training

adversarial examples generated via  
multi-step method



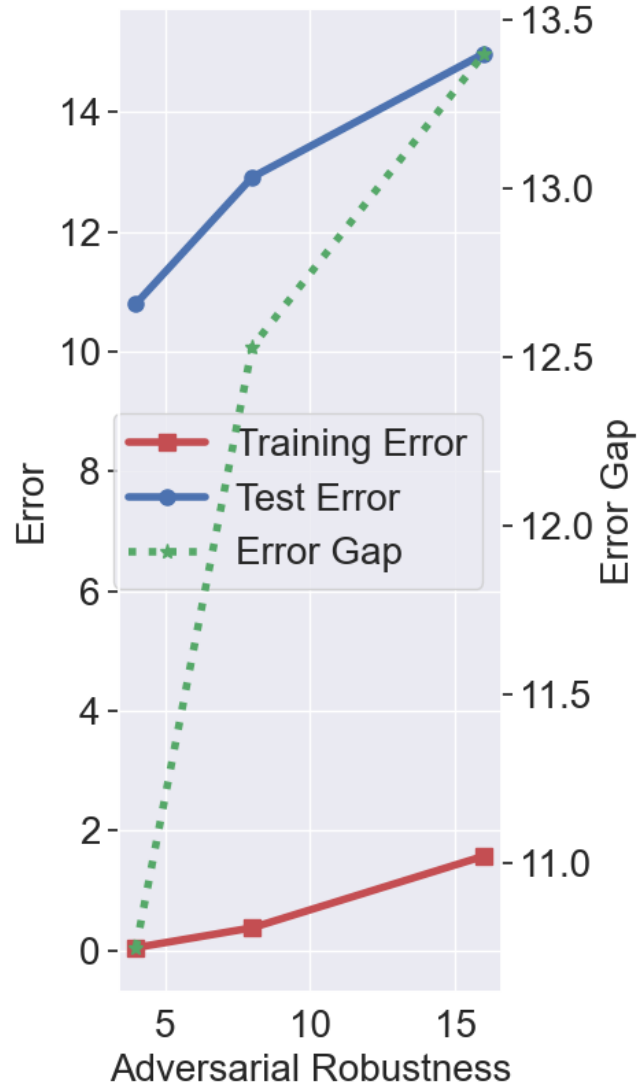
used for training

# Puzzling observations



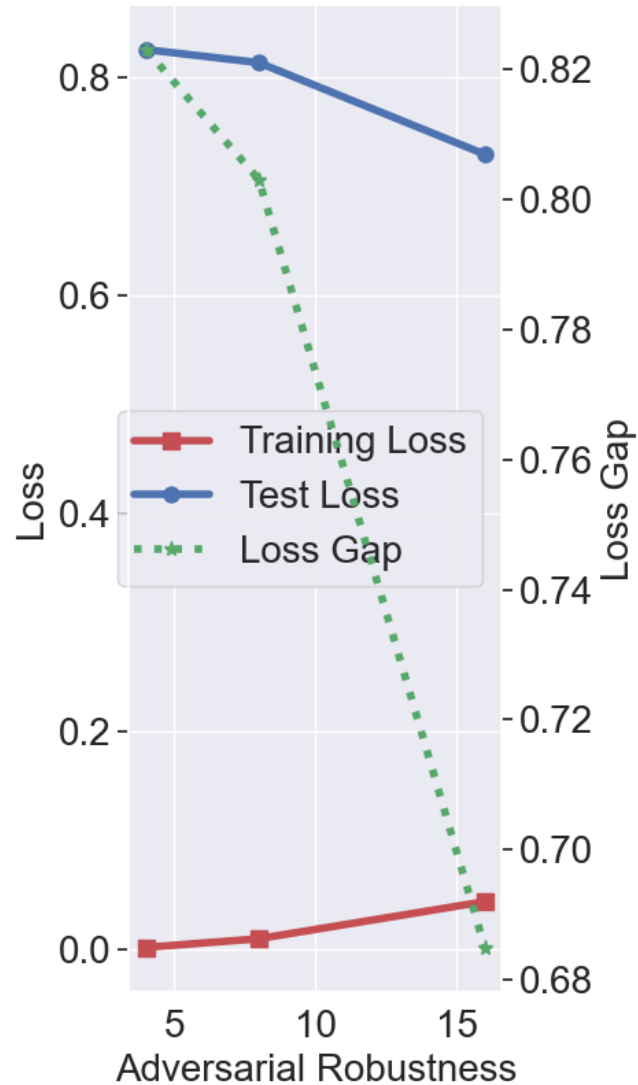
There exists a trade-off between adversarial robustness and accuracy

# Puzzling observations



Adversarial training exacerbate overfitting (the error gap is enlarged) ?

# Puzzling observations



Or an effective regularization (loss gap is reduced) ?

We study why effective regularization of adversarial robustness leads to poorer performance.

# Key results

- Establish a generalization bound that characterizes the generalization errors through the **margin**, **adversarial robustness radius** and **singular values of weight matrices** of neural networks.
- Empirical results:
  - For NNs with **high adversarial robustness**, the **singular values of weight matrices** has **low variance**.
  - The reduced variance of singular values of weight matrices results in **concentration of examples** around decision boundaries.
- The concentration of examples around decision boundaries smoothens **sudden changes** induced by perturbations, but also **increases indecisive misclassifications**.



# The generalization bound

**Theorem 3.1.** *Let  $T$  denote a NN with ReLU and MaxPooling nonlinear activation functions (the definition is put at eq. (6) for readers' convenience),  $l_\gamma$  the ramp loss defined at definition 4, and  $\mathcal{Z}$  the instance space assumed in section 3. Assume that  $\mathcal{Z}$  is a  $k$ -dimensional regular manifold that accepts an  $\epsilon$ -covering with covering number  $(\frac{C_X}{\epsilon})^k$ , and assumption assumption 3.1 holds. If  $T$  is  $\epsilon_0$ -adversarially robust (defined at definition 2),  $\epsilon \leq \epsilon_0$ , and denote  $v_{\min}$  the smallest IM margin in the covering balls that contain training examples (defined at definition 6),  $\sigma_{\min}^i$  the smallest singular values of weight matrices  $\mathbf{W}_i, i = 1, \dots, L - 1$  of a NN,  $\{\mathbf{w}_i\}_{i=1, \dots, |y|}$  the set of vectors made up with  $i$ th rows of  $\mathbf{W}_L$  (the last layer's weight matrix), then given an i.i.d. training sample  $S_m = \{z_i = (\mathbf{x}_i, y_i)\}_{i=1}^m$  drawn from  $\mathcal{Z}$ , its generalization error  $GE(l \circ T)$  (defined at eq. (1)) satisfies that, for any  $\eta > 0$ , with probability at least  $1 - \eta$*

$$GE(l_\gamma \circ T) \leq \max\left\{0, 1 - \frac{u_{\min}}{\gamma}\right\} + \sqrt{\frac{2 \log(2) C_X^k}{\epsilon^k m} + \frac{2 \log(1/\eta)}{m}}$$

$$\text{where } u_{\min} = \min_{y \neq \hat{y}} \|w_y - w_{\hat{y}}\|_2 \prod_{i=1}^{L-1} \sigma_{\min}^i v_{\min}$$

# The generalization bound

smallest margin on feature space:  $u_{\min} = \min_{y \neq \hat{y}} \|w_y - w_{\hat{y}}\|_2 \prod_{i=1}^{L-1} \sigma_{\min}^i v_{\min}$

(an example of margin :  $f_y(x) - \max_{i \neq y} f_i(x)$ )

$GE(l_\gamma \circ T) \leq \underbrace{\max\left\{0, 1 - \frac{u_{\min}}{\gamma}\right\}}_{\text{constant value}} + \underbrace{\sqrt{\frac{2 \log(2) C_X^k}{\epsilon^k m} + \frac{2 \log(1/\eta)}{m}}}_{\text{standard term}}$

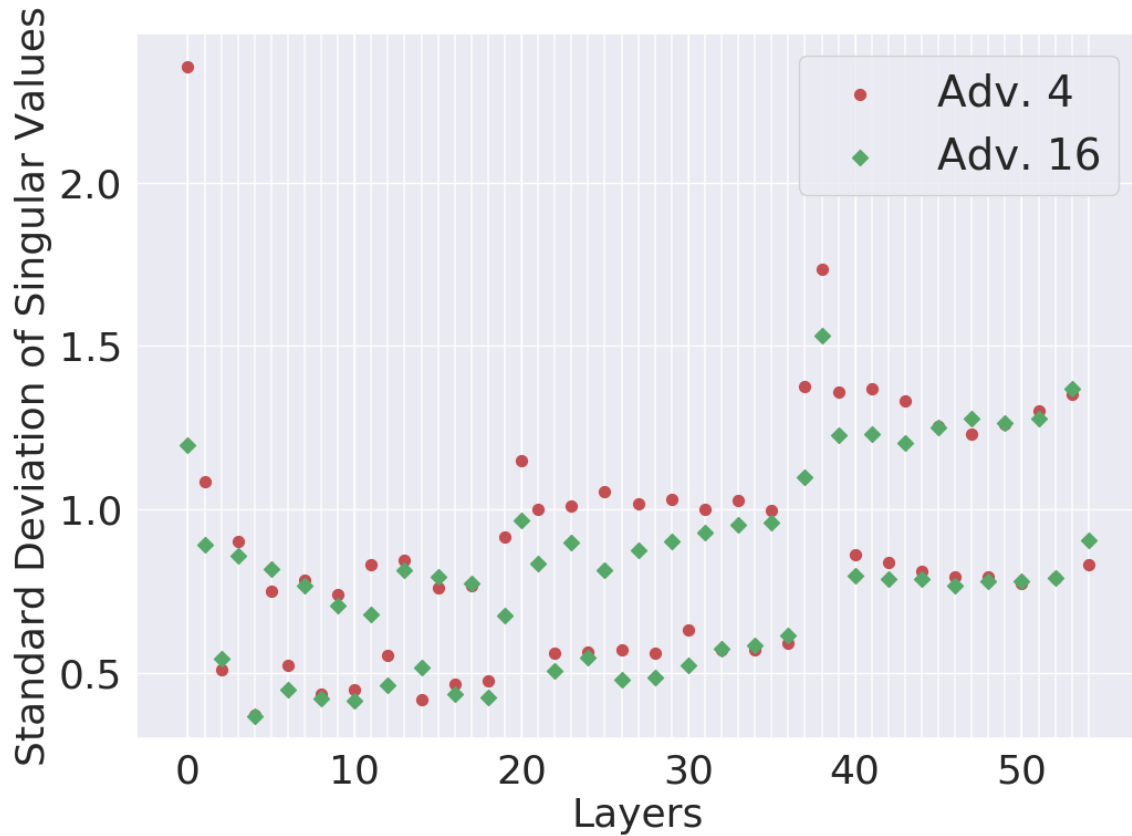
last layers' weight matrix

smallest singular values of neural networks' weight matrices

smallest instance-space margin, depends on the adversarial robustness radius

# Reduced variance of singular values

(a) Results from CIFAR10

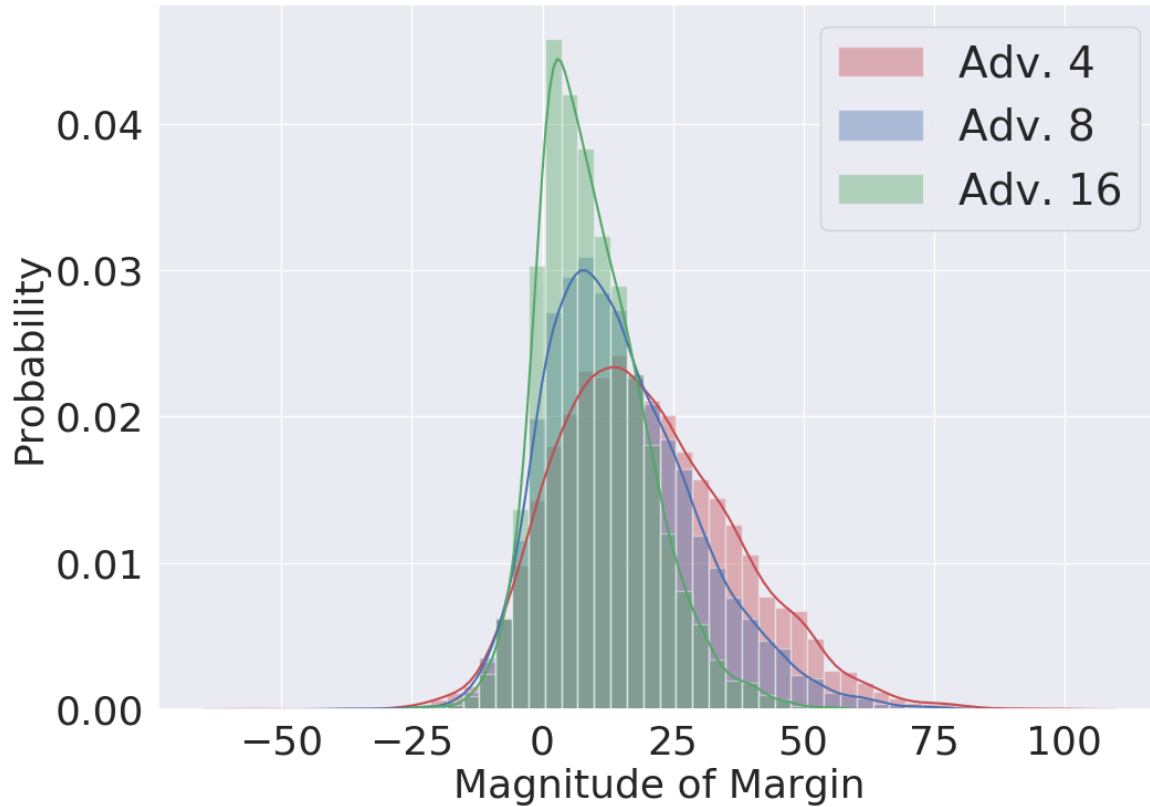


$$u_{\min} = \min_{y \neq \hat{y}} \|w_y - w_{\hat{y}}\|_2 \prod_{i=1}^{L-1} \sigma_{\min}^i v_{\min}$$

Take away message: stronger adversarial robustness reduces variance of singular values; and the reduced variance results in reduced variance of the norms of the activation outputs.

# Concentration of margin

(a) Results from CIFAR10



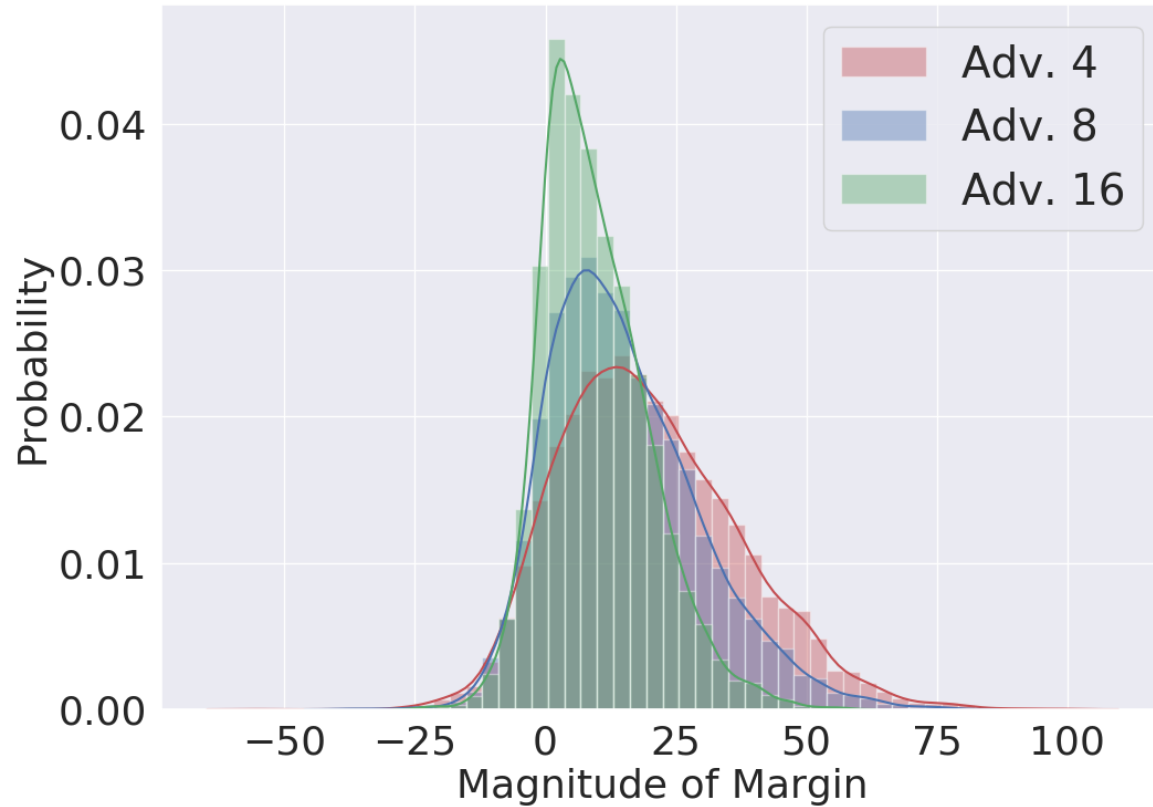
$$GE(l_\gamma \circ T) \leq \max \left\{ 0, 1 - \frac{u_{\min}}{\gamma} \right\}$$

$$+ \sqrt{\frac{2 \log(2) C_X^k}{\varepsilon^k m} + \frac{2 \log(1/\eta)}{m}}$$

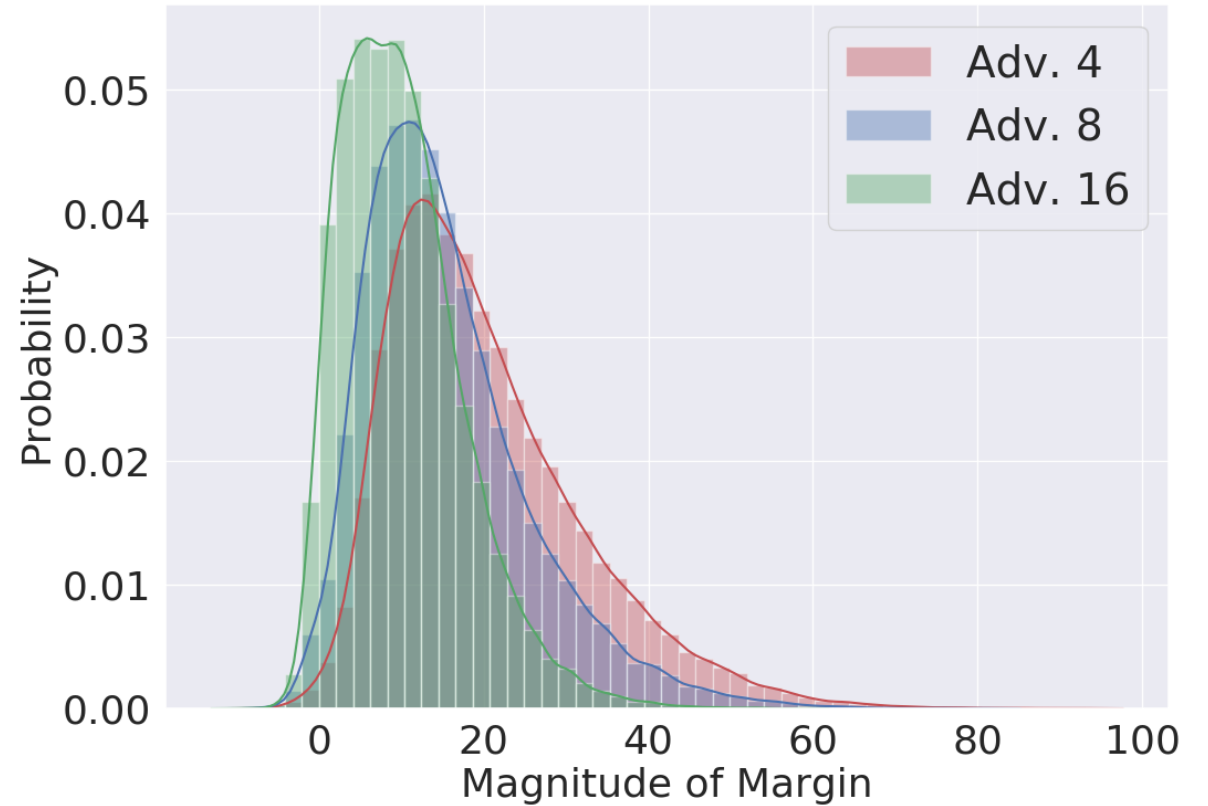
Take away message: reduced variance of the norms of the activation outputs results in concentration of examples; and the concentration results in different output of model.

# Concentration of margin

(a) Results from CIFAR10 test set



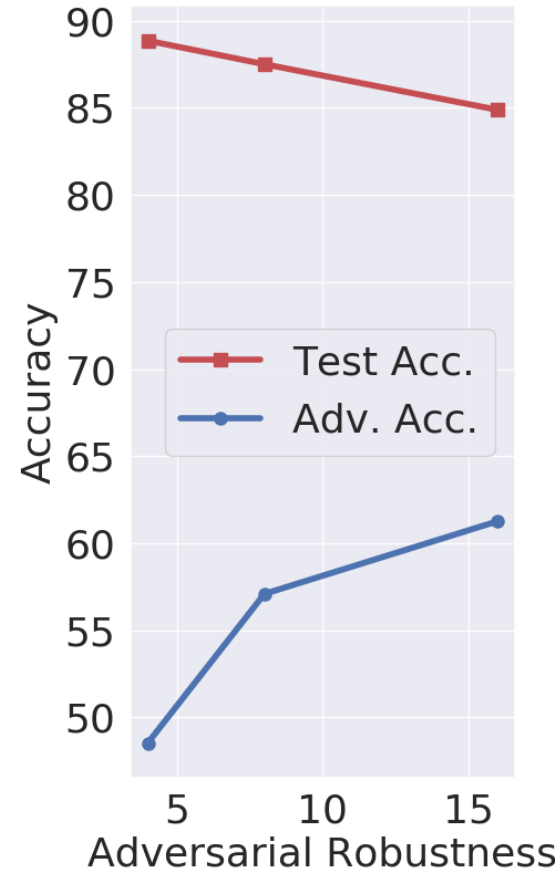
(b) Results from CIFAR10 training set



Take away message: the concentration also results in reduced loss/GE gap.

# The overall effect

(a) Results from CIFAR10



Take away message: the sample concentration around decision boundaries smoothens sudden changes induced perturbations, but also increases indecisive misclassification

# Take away message

Adversarial training indeed regularizes NNs, however, it does so by hurting the capacity of the NN hypothesis space.

# Future works

- Study the possible hypothesis: The concentration phenomena in NNs induced by AR suggests that to reduce the effects of adversarial noise, a NN might sacrifice its ability to distinguish inter-class difference.



Thanks for your attention