



北京大学  
PEKING UNIVERSITY

# Rethinking the Min-max Problem for Adversarial Robustness

Yisen Wang  
Peking University

[yisen.wang@pku.edu.cn](mailto:yisen.wang@pku.edu.cn)  
<https://yisenwang.github.io/>

Guest Lecture for CS498@UIUC  
Mar 17, 2021



# ML is Everywhere



Speech recognition

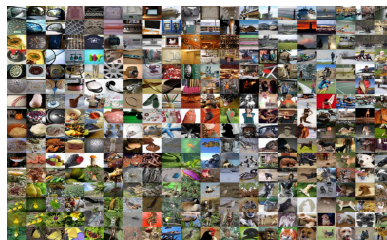
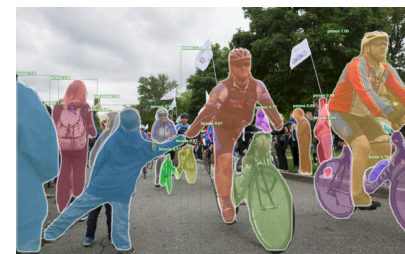
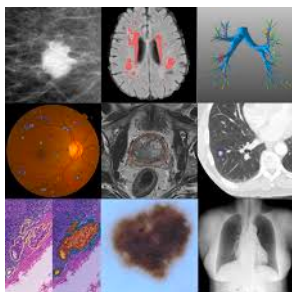


Image classification

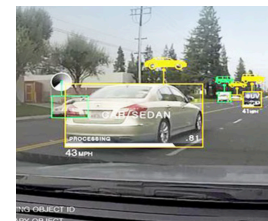


Object detection



Medical diagnosis

**Machine Learning**



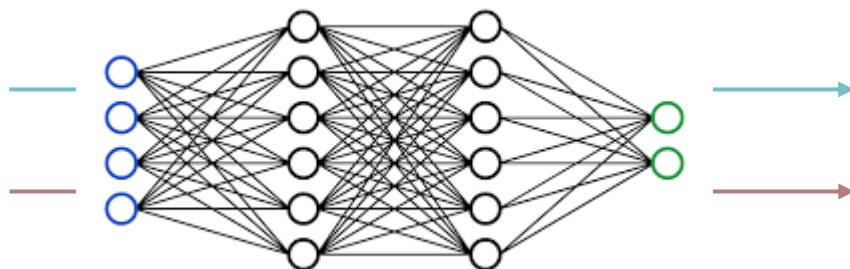
Autonomous driving



Playing games

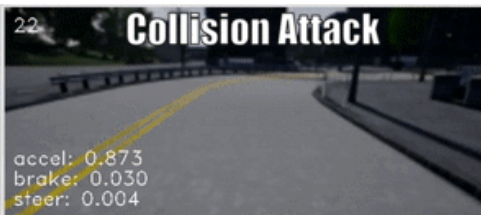


# However



Dog,  
82% confidence

Ostrich,  
98% confidence





Are we doomed?  
(Is ML inherently not reliable?)

**NO!** But we need to re-think how we do ML  
(adversarial aspects = stress-testing our solutions)





# Adversarial Example

Model training:

$$\min_{\theta} \sum_{(x_i, y_i) \in D_{train}} L(f_{\theta}(x_i), y_i)$$

$D_{train}$ : training data  
 $x_i$ : training sample  
 $y_i$ : class label  
 $L$ : loss function  
 $f_{\theta}$ : model

Adversarial attack:

$$\max_{x'} L(f_{\theta}(x'), y) \quad \text{st. } \|x' - x\|_p \leq \epsilon \quad \text{for } x \in D_{test}$$

increase error

small change

test time attack

$$\|x' - x\|_{\infty} \leq \epsilon = \frac{8}{255} \approx 0.031$$

- Fast Gradient Sign Method (FGSM) (*Goodfellow et al., 2014*):

$$x' = x + \epsilon \cdot \text{sign } \nabla_x L(f_{\theta}(x), y)$$

$x'$ : adv examples

- Projected Gradient Descent (PGD) is an iterative version of FGSM (*Madry et al., 2018*)

$$x'^{(k+1)} = \Pi_{\epsilon} (x'^{(k)} + \alpha \cdot \text{sign } \nabla_x L(f_{\theta}(x'^{(k)}), y))$$



How to obtain adversarially **robust models**?



# Adversarial Training

Adversarial training is a **min-max optimization** process:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \overbrace{\max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta}(x'_i), y_i)}^{\text{attacking}}$$

$L$ : loss,  $f_{\theta}$ : model,  $x_i$ : clean example,  $y_i$ : class,  
 $x'_i$ : adversarial example.

## 1. Inner Maximization:

- This is to **generate adversarial examples**, by maximizing the loss  $L$ .
- It is a constrained optimization problem:  $\|x'_i - x_i\|_p \leq \epsilon$ .

## 2. Outer Minimization:

- A typical process to **train a model**, but on adversarial examples  $x'_i$  generated by the inner maximization.



# Convergence Score of the Maximization

**Question: How well the inner maximization is solved?**

## Definition ( First-Order Stationary Condition (FOSC))

Given a data sample  $x^0 \in X$ , let  $x^k$  be an intermediate example found at the  $k^{\text{th}}$  step of the inner maximization. The First-Order Stationary Condition of  $x^k$  is

$$c(x^k) = \max_{x \in \chi} \langle x - x^k, \nabla_x f(\theta, x^k) \rangle,$$

where  $\chi = \{x \mid \|x - x^0\|_\infty \leq \epsilon\}$  is the input domain of the  $\epsilon$ -ball around normal example  $x^0$ ,  $f(\theta, x^k) = \ell(h_\theta(x^k), y)$ , and  $\langle \cdot \rangle$  is the inner product.

## FOSC:

- A smaller value of  $c(x^k)$  indicates a better solution of the inner maximization, or equivalently, better convergence quality of the adversarial example  $x^k$ .
- To help Danskin's Theorem hold.



# Convergence Theorem

## Theorem 1

Under certain assumptions, let  $\Delta = L_S(\theta^0) - \min_{\theta} L_S(\theta)$ . If the step size of the outer minimization is set to  $\eta_t = \min\left(\frac{1}{L}, \sqrt{\frac{\Delta}{L\sigma^2 T}}\right)$ . Then the output of **Adversarial**

**Training** satisfies

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[\|\nabla L_S(\theta^t)\|_2^2] \leq 4\sigma \sqrt{\frac{L\Delta}{T}} + \frac{5L_{\theta x}^2 \delta}{\mu},$$

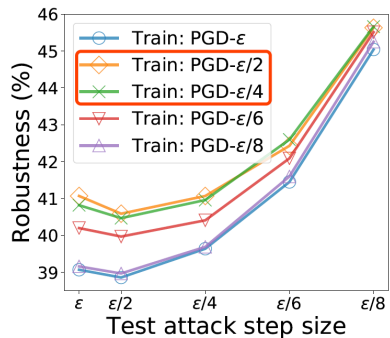
where  $L = \left(\frac{L_{\theta x} L_{\theta x}}{\mu} + L_{\theta\theta}\right)$ .

- Inner maximization:  $\text{FOSC} \leq \delta$ , adversarial training can converge to a first-order stationary point up to a precision of  $\frac{5L_{\theta x}^2 \delta}{\mu}$
- If  $\delta$  is sufficiently small such that  $\frac{5L_{\theta x}^2 \delta}{\mu}$  small enough, adversarial training can find a robust model  $\theta^T$ .

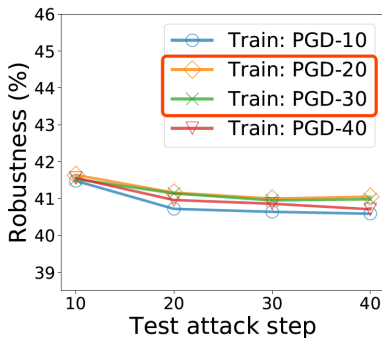




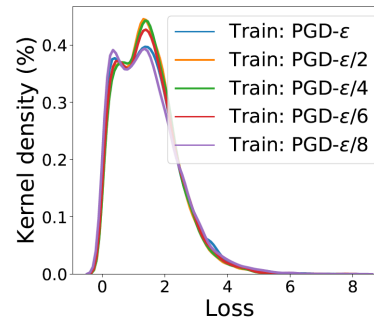
# Why do we need FOSC?



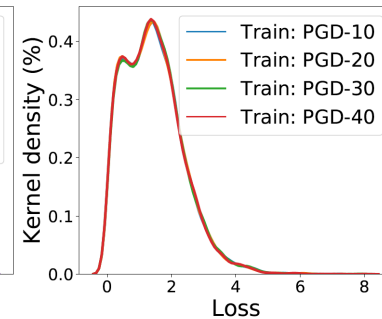
(a) Robustness vs. Step size



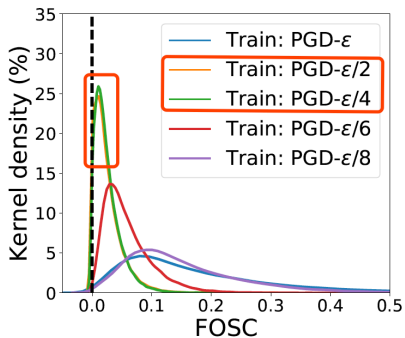
(b) Robustness vs. Step number



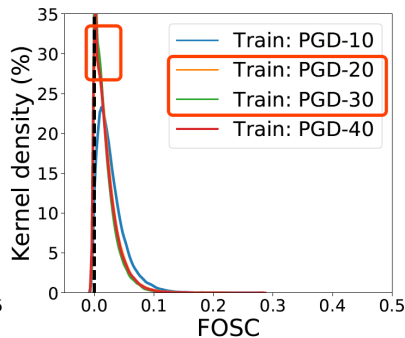
(e) Loss vs. Step size



(f) Loss vs. Step number



(c) FOSC vs. Step size



(d) FOSC vs. Step number

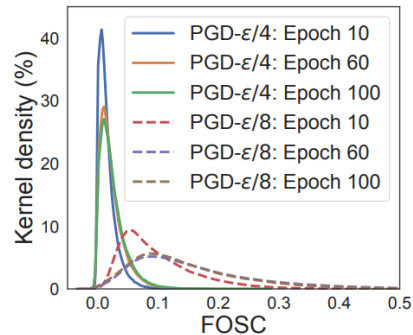
## Adversarial Training with different settings for PGD-based inner maximization.

- **PGD step size:**  $\text{PGD}-\frac{\epsilon}{2}$  /  $\text{PGD}-\frac{\epsilon}{4}$  produces the best robustness, their FOSC values are also concentrated around 0.
- **PGD step number:** similar robustness, with PGD-20/30 are slightly better, reflected by the distribution of FOSC.
- **Loss distributions** are similar for different robustness.

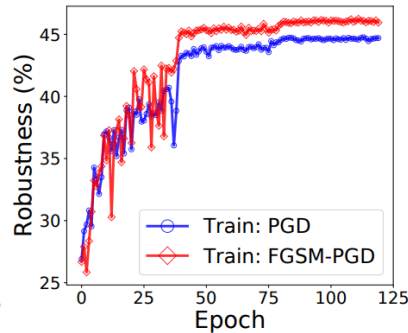
**FOSC is a good and reliable indicator of the final robustness**



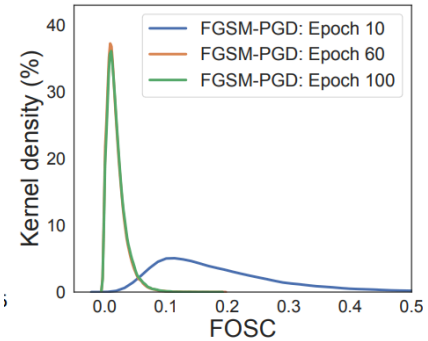
# FOSC View of Adversarial Training



(a) FOSC



(b) Robustness



(c) FOSC

- Standard adversarial training **overfits** to strong PGD adversarial examples at the **early stage**.
- Simply use **weak attack FGSM** at the **early stage** can improve robustness.
- Improvement in robustness is also reflected in FOSC distribution.

**The principle behind warm-up techniques**



Warm-up is a method to solve max better,  
is there other options?



# Rethinking the Robust Generalization Gap

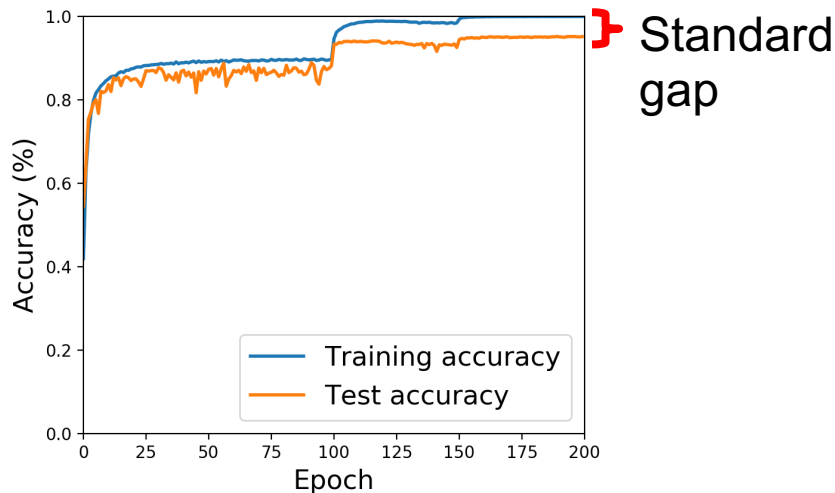
Adversarial training is a **min-max optimization** process:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta}(x'_i), y_i)$$

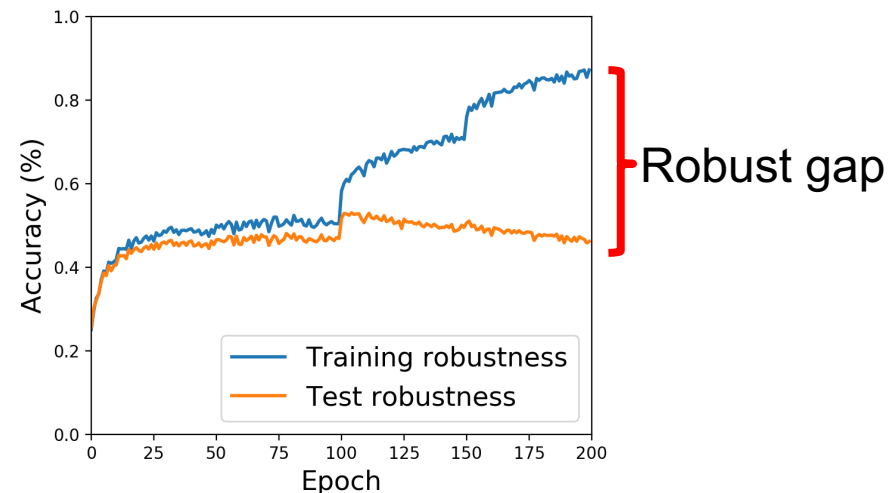
can be rewritten as

$$\min_{\mathbf{w}} \rho(\mathbf{w}), \quad \text{where} \quad \rho(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} \ell(f_{\mathbf{w}}(x'_i), y_i),$$

Standard training



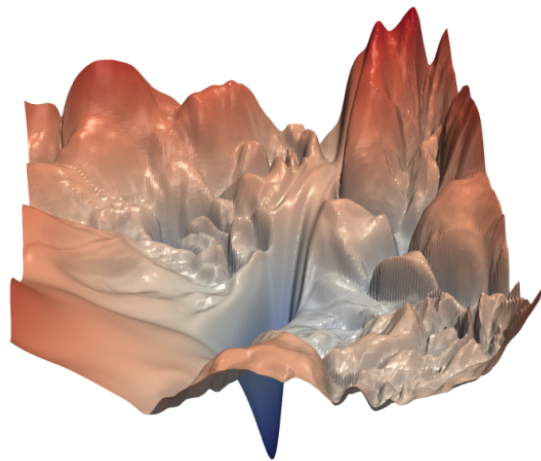
Adversarial training



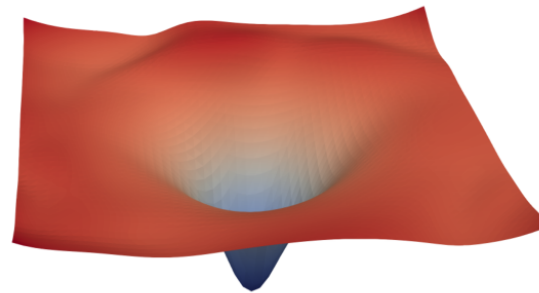


# View from weight loss landscape

- Inspiring from standard Training:
  - flatter weight loss landscape, smaller standard generalization gap



(a) without skip connections



(b) with skip connections

**Is this conclusion still existing in adversarial training?**





# Adapted Visualization Method

- Inspiring from standard Training:
  - flatter weight loss landscape, smaller standard generalization gap
- Is this conclusion still existing in adversarial training?

Visualization method in *Hao Li et al.* NeurIPS2018

Standard training

$$g(\alpha) = L(f_{w+\alpha d}(\mathbf{x}_i), y_i)$$



Adversarial training

$$g(\alpha) = L(f_{w+\alpha d}(\mathbf{x}'_i), y_i) ?$$

$\mathbf{x}'_i$  is from **pre-generated** adversarial examples<sup>[1,2]</sup>

Failed to draw the conclusion

The correct way:

$$g(\alpha) = \rho(\mathbf{w} + \alpha \mathbf{d}) = \frac{1}{n} \sum_{i=1}^n \max_{\|\mathbf{x}'_i - \mathbf{x}_i\|_p \leq \epsilon} \ell(f_{\mathbf{w} + \alpha \mathbf{d}}(\mathbf{x}'_i), y_i),$$

Generating adversarial examples **on-the-fly**

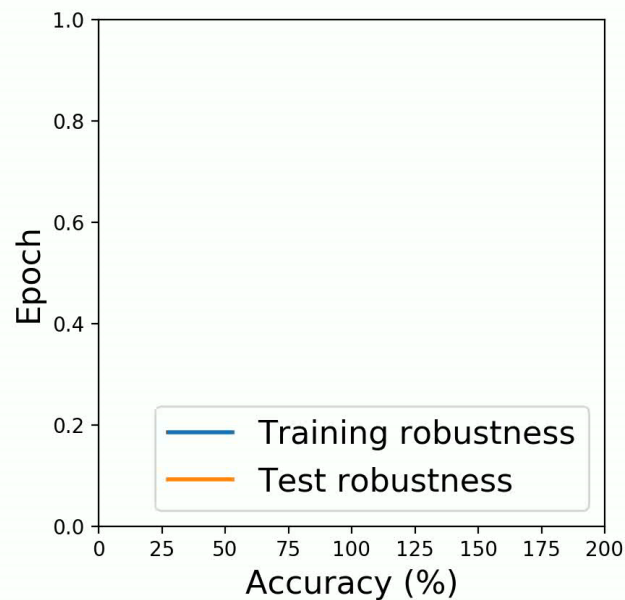
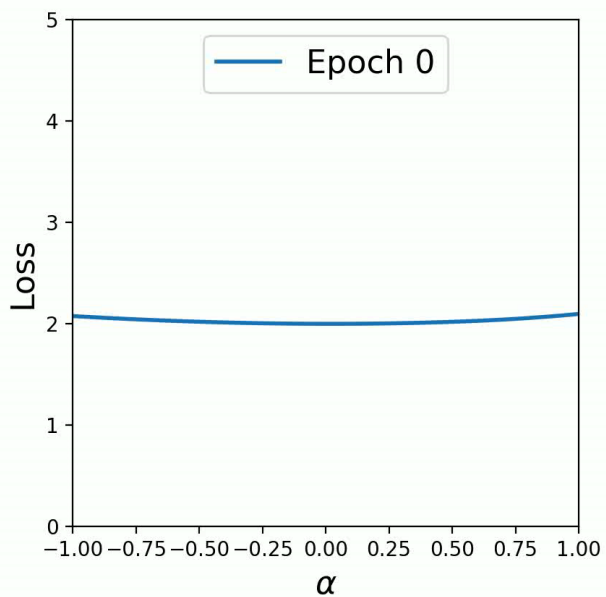
[1] Understanding adversarial robustness through loss landscape geometries, *arxiv* 2019.

[2] Interpreting adversarial robustness: A view from decision surface in input space, *arxiv* 2018



# Weight loss landscape

In the learning process of adversarial training

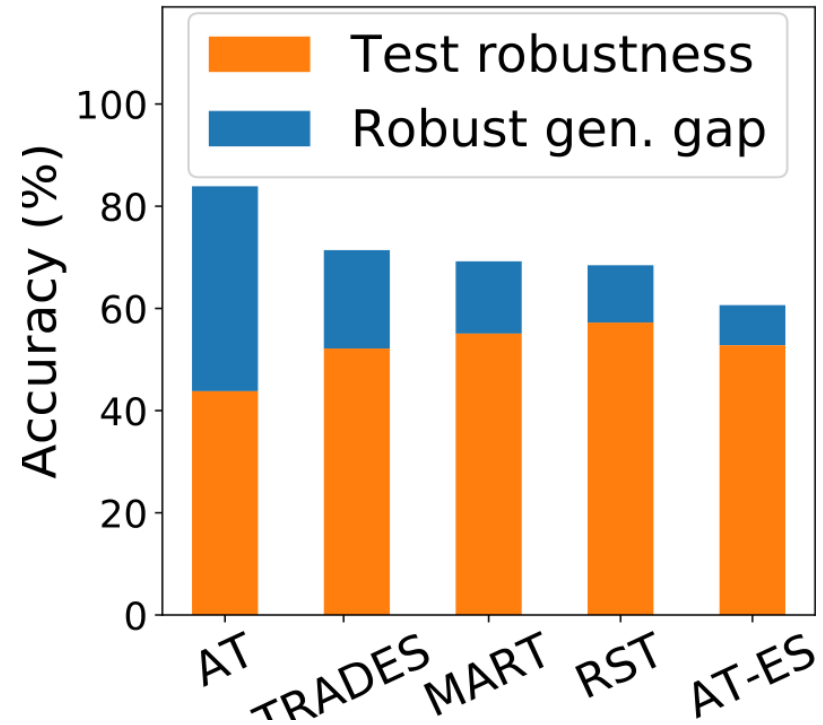
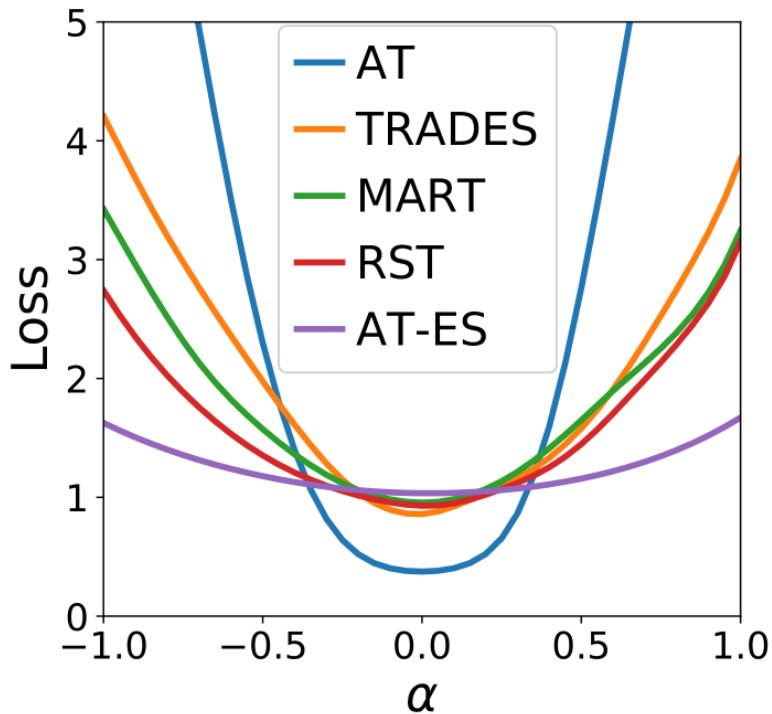


**Weight loss landscape has a strong correlation with robust generalization gap**



# Weight loss landscape

Across different adversarial training methods



**Weight loss landscape has a strong correlation with robust generalization gap**



# Theoretical view

- Informally from PAC-Bayesian bound

$$\mathbb{E}_{\{\mathbf{x}_i, y_i\}_{i=1}^n, \mathbf{u}}[\rho(\mathbf{w} + \mathbf{u})] \leq \rho(\mathbf{w}) + \underbrace{\{\mathbb{E}_{\mathbf{u}}[\rho(\mathbf{w} + \mathbf{u})] - \rho(\mathbf{w})\}}_{\text{flatness of weight loss landscape}} + 4\sqrt{\frac{1}{n}KL(\mathbf{w} + \mathbf{u} \| P) + \ln \frac{2n}{\delta}}.$$

flatness of weight loss landscape

- Explicitly flattening the weight loss landscape via replacing expectation by maximization

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta}(x'_i), y_i) \quad \longrightarrow \quad \min_{\theta} \max_{\|v\|_p \leq \gamma \|\theta\|_p} \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta+v}(x'_i), y_i)$$

- Two max makes the maximization (min-max) solve better
- How to intuitively understand these two perturbations?
  - Input perturbation is **local** worst for **each example**
  - Weight perturbation is **global** worst for **multiple examples**



# Implementation

AWP-based Adversarial training (AT-AWP)

$$\min_{\theta} \max_{\|v\|_p \leq \gamma \|\theta\|_p} \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta+v}(x'_i), y_i)$$

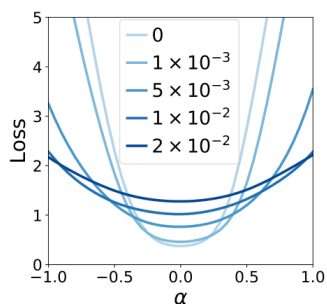
- An empirical implementation:
  1. craft adversarial examples  $x'_i$ ;
  2. calculate AWP based on  $x'_i$  using one extra forward and backward propagation;
  3. update the parameter using the gradient based on  $x'_i$  and AWP.
- Only ~8% time overhead in our implementation of AT-AWP.
- AWP is easily extended to other methods, such as TRADES, MART and RST.



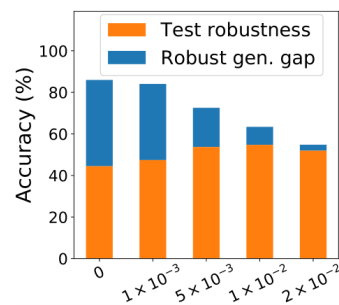


# Real robustness improvement

- AWP indeed flattens weight loss landscape, and reduces the robust generalization gap.

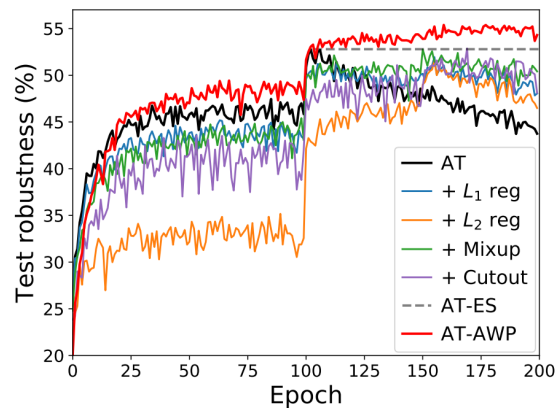


(c) Weight loss landscape



(d) Generalization gap

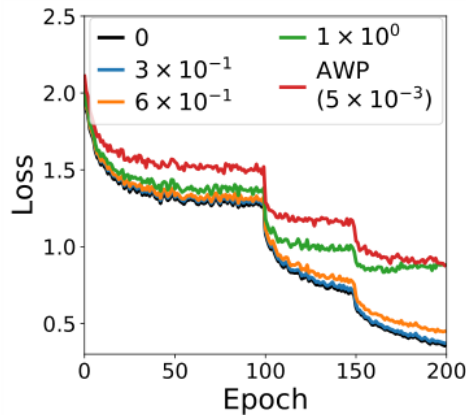
- AWP really improves both the **last** and **best** robustness during training.



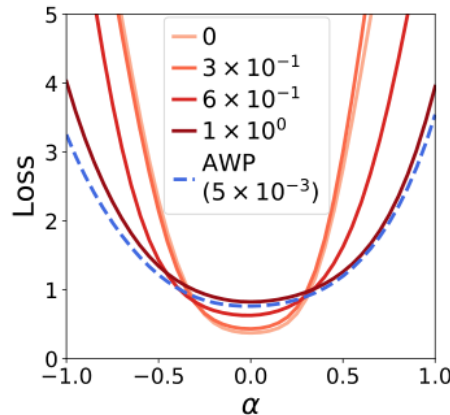


# AWP vs. Random WP

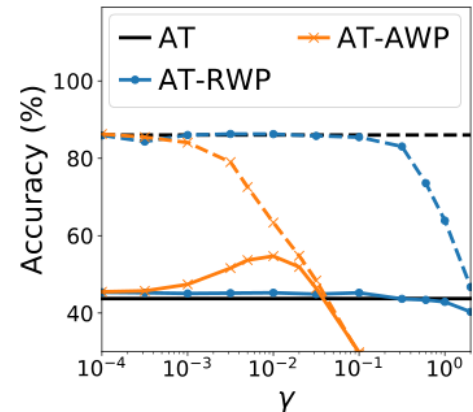
- AWP easily finds the worst-case perturbation, while RWP needs a relatively large perturbation;
- AWP obtain a flatter weight loss landscape using smaller perturbations;
- AWP balances the training robustness and robust gap well.



(a) Loss curve



(b) Weight loss landscapes



(c) Robustness



# Universal robustness improvement

Table 2: Test robustness (%) on CIFAR-10 using WideResNet under  $L_\infty$  threat model. We omit the standard deviations of 5 runs as they are very small ( $< 0.40\%$ ), which hardly effect the results.

Defense	Natural	FGSM	PGD-20	PGD-100	$CW_\infty$	SPSA	AA
AT	<b>86.07</b>	61.76	56.10	55.79	54.19	61.40	52.60 <sup>4</sup>
AT-AWP	85.57	<b>62.90</b>	<b>58.14</b>	<b>57.94</b>	<b>55.96</b>	<b>62.65</b>	<b>54.04</b>
TRADES	84.65	61.32	56.33	56.07	54.20	61.10	53.18
TRADES-AWP	<b>85.36</b>	<b>63.49</b>	<b>59.27</b>	<b>59.12</b>	<b>57.07</b>	<b>63.85</b>	<b>56.17</b>
MART	84.17	61.61	58.56	57.88	54.58	58.90	51.10
MART-AWP	<b>84.43</b>	<b>63.98</b>	<b>60.68</b>	<b>59.32</b>	<b>56.37</b>	<b>62.75</b>	<b>54.23</b>
Pre-training	87.89	63.27	57.37	56.80	55.95	62.55	54.99
Pre-training-AWP	<b>88.33</b>	<b>66.34</b>	<b>61.40</b>	<b>61.21</b>	<b>59.28</b>	<b>65.55</b>	<b>57.39</b>
RST	<b>89.69</b>	67.94	62.60	62.22	60.47	67.60	59.65
RST-AWP	88.25	<b>69.60</b>	<b>63.73</b>	<b>63.58</b>	<b>61.62</b>	<b>68.72</b>	<b>61.10</b>

Table 3: Test robustness (%) on CIFAR-10 using WideResNet under  $L_\infty$  threat model. In brackets, + indicates improvements over Pre-training.

Defense	PGD-20	$CW_\infty$	AA
Pre-training	57.37	55.95	54.92
TRADES-AWP	59.27 (+1.90)	57.07 (+1.12)	56.17 (+1.25)
Pre-training-AWP	61.40 (+4.03)	59.28 (+3.33)	57.39 (+2.47)



Except max process, how about min process?



# Revisiting the Input Examples

Adversarial training is a **min-max optimization** process:

$$\min_{\theta} \frac{1}{n} \sum_{i=1}^n \max_{\|x'_i - x_i\|_p \leq \epsilon} L(f_{\theta}(x'_i), y_i)$$

$L$ : loss,  $f_{\theta}$ : model,  $x_i$ : clean example,  $y_i$ : class,  $x'_i$ : adversarial example.

Adversarial examples are only defined on **correctly classified examples**

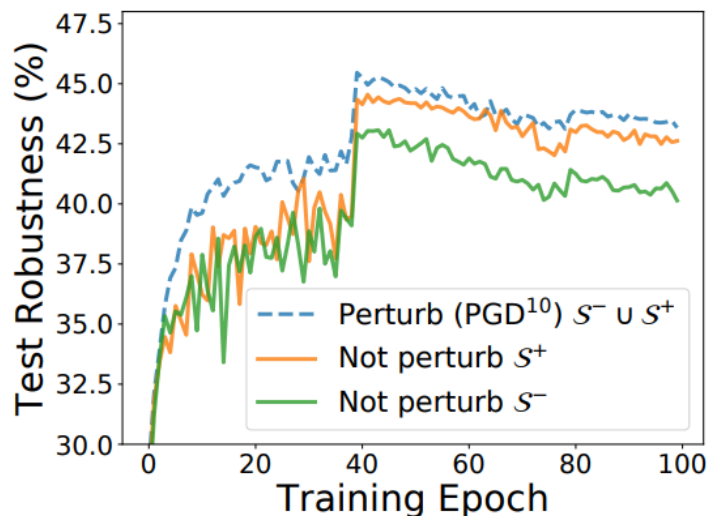
**How about misclassified examples?**





# Misclassified vs. correctly classified examples

- A pre-trained network to select the **same size (13%)**
  - Subset of misclassified examples  $S^-$
  - Subset of correctly classified examples  $S^+$

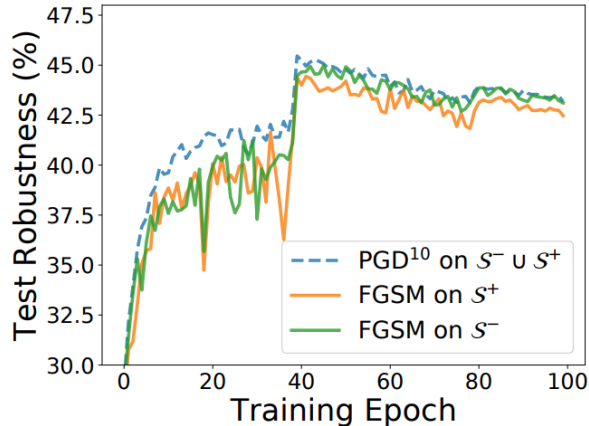


**Misclassified examples have a significant impact on the final robustness**



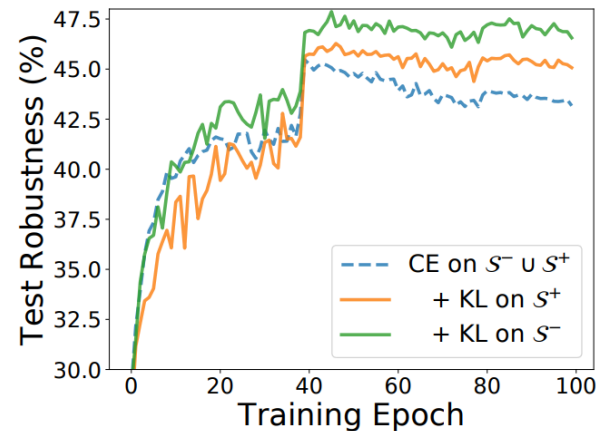
# Delving into the max and min processes

- For inner maximization process:
  - Weak attack on misclassified examples  $S^-$
  - Weak attack on correctly classified examples  $S^+$
- For outer minimization process:
  - Regularization on misclassified examples  $S^-$
  - Regularization on correctly classified examples  $S^+$



(b) Inner maximization

different maximization techniques have **negligible** effect



(c) Outer minimization

different minimization techniques have **significant** effect



# Misclassification aware adversarial risk

- Adversarial risk:

$$\mathcal{R}(h_{\theta}) = \frac{1}{n} \sum_{i=1}^n \max_{\mathbf{x}'_i \in \mathcal{B}_{\epsilon}(\mathbf{x}_i)} \mathbb{1}(h_{\theta}(\mathbf{x}'_i) \neq y_i),$$

- Correctly classified and misclassified examples:

$$\mathcal{S}_{h_{\theta}}^+ = \{i : i \in [n], h_{\theta}(\mathbf{x}_i) = y_i\} \quad \text{and} \quad \mathcal{S}_{h_{\theta}}^- = \{i : i \in [n], h_{\theta}(\mathbf{x}_i) \neq y_i\}$$

- Misclassification aware adversarial risk:

$$\begin{aligned} \min_{\theta} \mathcal{R}_{\text{misc}}(h_{\theta}) &:= \frac{1}{n} \left( \sum_{i \in \mathcal{S}_{h_{\theta}}^+} \mathcal{R}^+(h_{\theta}, \mathbf{x}_i) + \sum_{i \in \mathcal{S}_{h_{\theta}}^-} \mathcal{R}^-(h_{\theta}, \mathbf{x}_i) \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left\{ \mathbb{1}(h_{\theta}(\hat{\mathbf{x}}'_i) \neq y_i) + \mathbb{1}(h_{\theta}(\mathbf{x}_i) \neq h_{\theta}(\hat{\mathbf{x}}'_i)) \cdot \mathbb{1}(h_{\theta}(\mathbf{x}_i) \neq y_i) \right\} \end{aligned}$$



# Misclassification Aware adveRsarial Training (MART)

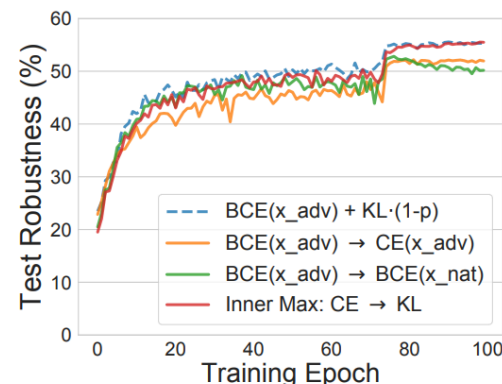
- Surrogate loss functions (existing methods and MART):

Defense Method	Loss Function
<i>Standard</i>	$\text{CE}(\mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}), y)$
ALP	$\text{CE}(\mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}), y) + \lambda \cdot \ \mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}) - \mathbf{p}(\mathbf{x}, \boldsymbol{\theta})\ _2^2$
CLP	$\text{CE}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}), y) + \lambda \cdot \ \mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}) - \mathbf{p}(\mathbf{x}, \boldsymbol{\theta})\ _2^2$
TRADES	$\text{CE}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}), y) + \lambda \cdot \text{KL}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}) \parallel \mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}))$
MMA	$\text{CE}(\mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}), y) \cdot \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x}) = y) + \text{CE}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}), y) \cdot \mathbb{1}(h_{\boldsymbol{\theta}}(\mathbf{x}) \neq y)$
<b>MART</b>	$\text{BCE}(\mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta}), y) + \lambda \cdot \text{KL}(\mathbf{p}(\mathbf{x}, \boldsymbol{\theta}) \parallel \mathbf{p}(\hat{\mathbf{x}}', \boldsymbol{\theta})) \cdot (1 - \mathbf{p}_y(\mathbf{x}, \boldsymbol{\theta}))$

$$\text{BCE}(\mathbf{p}(\hat{\mathbf{x}}'_i, \boldsymbol{\theta}), y_i) = -\log(\mathbf{p}_{y_i}(\hat{\mathbf{x}}'_i, \boldsymbol{\theta})) - \log\left(1 - \max_{k \neq y_i} \mathbf{p}_k(\hat{\mathbf{x}}'_i, \boldsymbol{\theta})\right)$$



(a) Removing



(b) Replacing

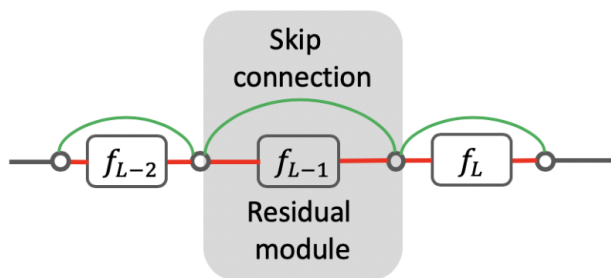


Beyond training objective,  
is **model architecture** related to robustness?



# Skip connection matters

- Neural network architectures:
  - Skip connection, activation, batch normalization, ...
- Skip connection

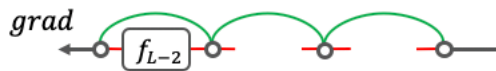


white-box / black-box

100% / 52.52%



100% / 55.24%



100% / **62.10%**

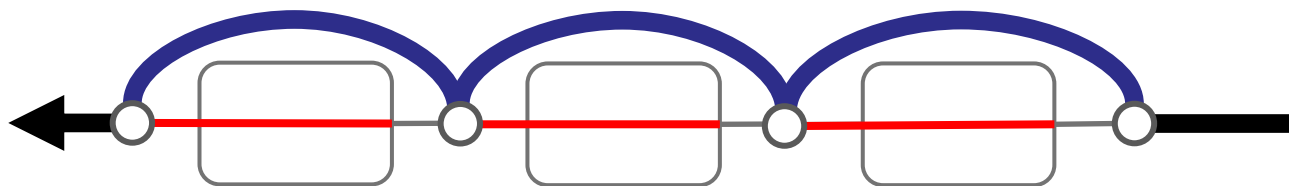


99.86% / 47.70%

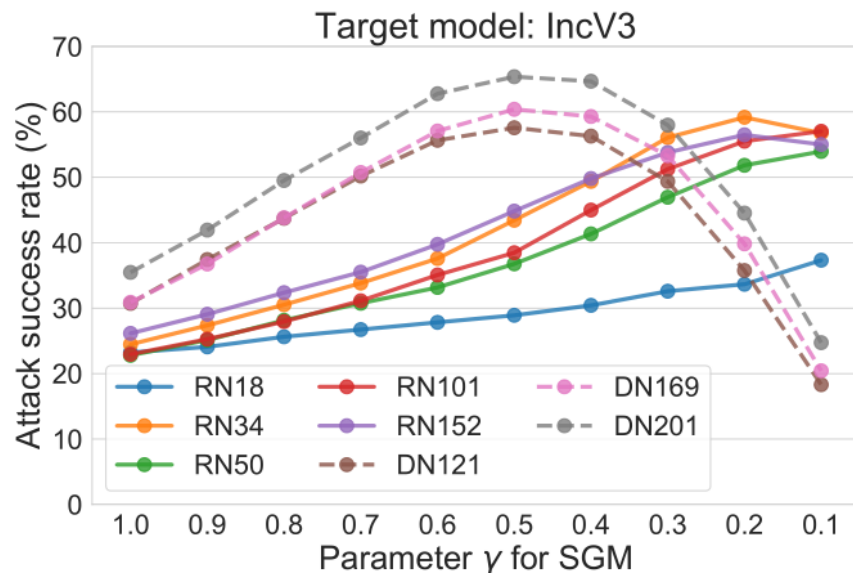
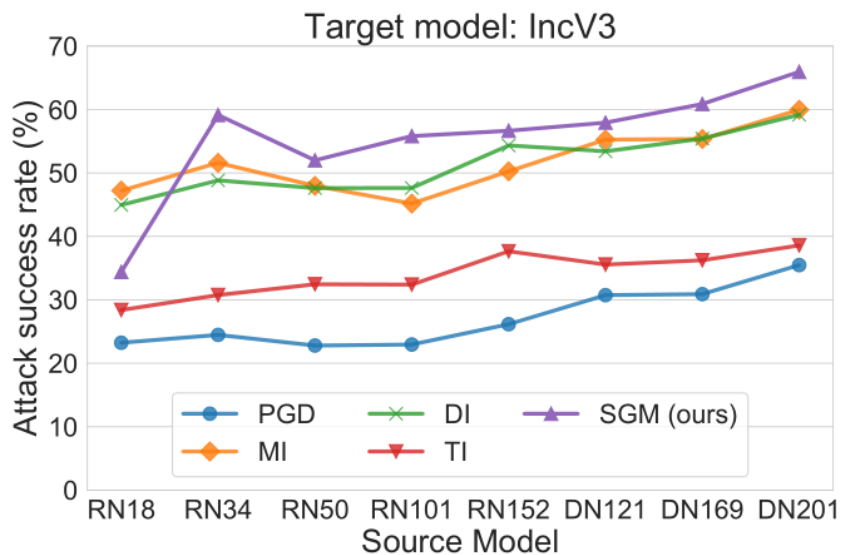
**Skip connections expose more transferable information !**



# Skip Gradient Method (SGM)



$$\nabla_x \ell = \frac{\partial \ell}{\partial \mathbf{z}_L} \prod_{i=0}^{L-1} \left( \gamma \frac{\partial f_{i+1}}{\partial \mathbf{z}_i} + 1 \right) \frac{\partial \mathbf{z}_0}{\partial \mathbf{x}}$$





# Takehome Message

- For the min-max problem, the following aspects are essential:
  - how to make max solves better
  - How to make min process easily
- Model architecture is also important for adversarial research





# Related Papers

- Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, Quanquan Gu, “*On the Convergence and Robustness of Adversarial Training*”, **ICML 2019 Long Talk**
- Dongxian Wu, Shu-Tao Xia, Yisen Wang, “*Adversarial Weight Perturbation Helps Robust Generalization*”, **NeurIPS 2020**
- Yisen Wang, Difan Zou, Jinfeng Yi, James Bailey, Xingjun Ma, Quanquan Gu, “*Improving Adversarial Robustness Requires Revisiting Misclassified Examples*”, **ICLR 2020**
- Dongxian Wu, Yisen Wang, Shu-Tao Xia, James Bailey, Xingjun Ma, “*Skip Connections Matter: On the Transferability of Adversarial Examples Generated with ResNets*”, **ICLR 2020 Spotlight**
- Hanxun Huang, Xingjun Ma, Sarah Monazam Erfani, James Bailey, Yisen Wang, “*Unlearnable Examples: Making Personal Data Unexploitable*”, **ICLR 2021 Spotlight**
- Yang Bai, Yuyuan Zeng, Yong Jiang, Shu-Tao Xia, Xingjun Ma, Yisen Wang, “*Improving Adversarial Robustness via Channel-wise Activation Suppressing*”, **ICLR 2021 Spotlight**

Building ML one can truly rely on



北京大学  
PEKING UNIVERSITY

# Thanks!

[yisen.wang@pku.edu.cn](mailto:yisen.wang@pku.edu.cn)  
<https://yisenwang.github.io/>