

# Training Robust Neural Networks

## attaqueoudefense

Victor Lemaître, Azammat Charaf Zadah, Louis Carron

IASD - Université PSL - Paris Dauphine

Academic Year 2025-2026

# Outline

1. Setup
2. Improving Model Capacities
3. Randomness
4. Towards Black Box attack
5. Other Failed Improvement

# 1. Setup

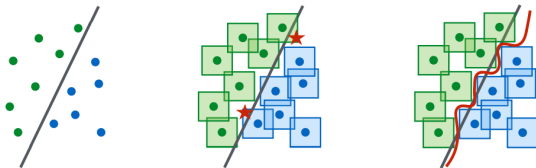
- ▶ Implemented attacks FGSM & PGD
- ▶ Standard adversarial training:
  - Perturb a fraction of each batch
  - Trade-off: clean accuracy drops too much

**Key idea: Curriculum Adversarial Training**

$$\left\{ \begin{array}{l} t = \frac{\text{epoch}}{\text{total epochs}} \\ \epsilon(t) = 0.01(1 - t) + 0.1 t \\ \text{PGD iterations}(t) = \text{int}(10(1 - t) + 40 t) \end{array} \right.$$

## 2. Improving Model Capacities

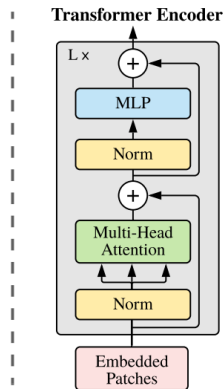
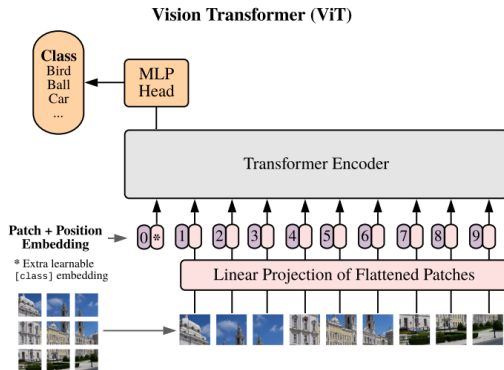
**Key idea:** stronger models improve the clean/attacked trade-off.



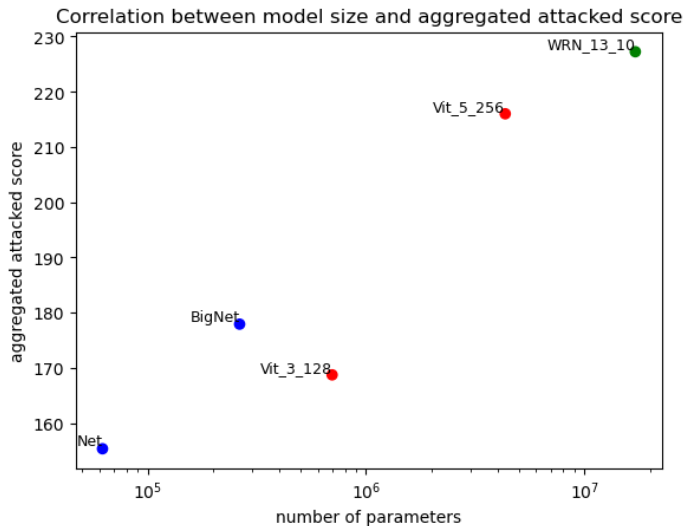
**Practical improvements:**

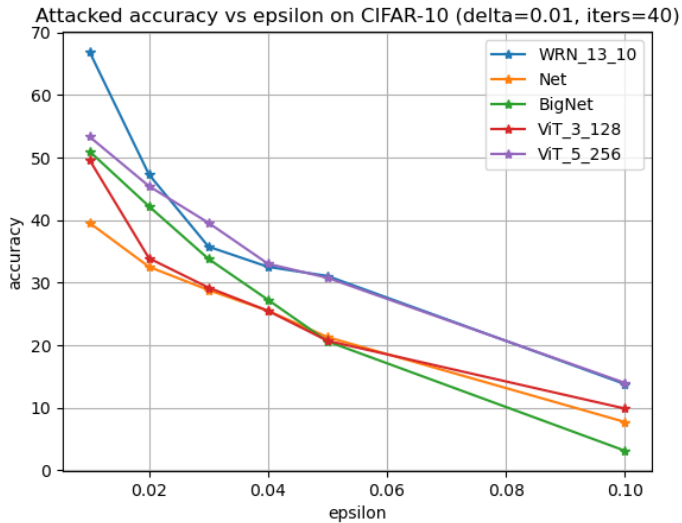
- ▶ Switch from SGD+momentum to AdamW
- ▶ Use GELU activations (+5–7% natural accuracy on larger models)

## 2. Improving Model Capacities

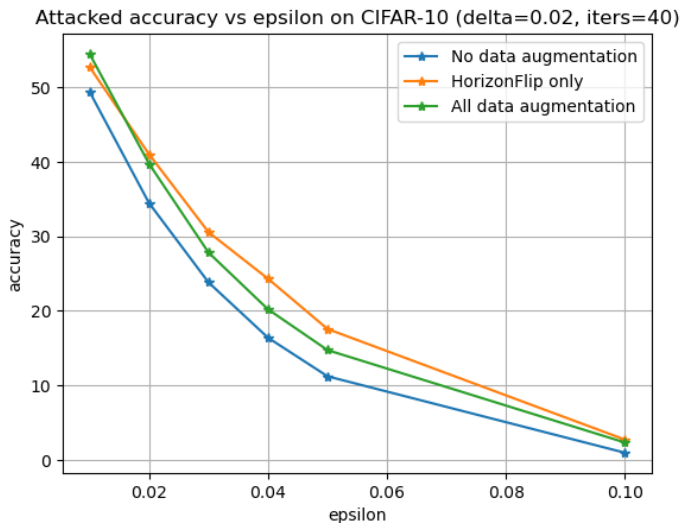


## 2.1 Model Size





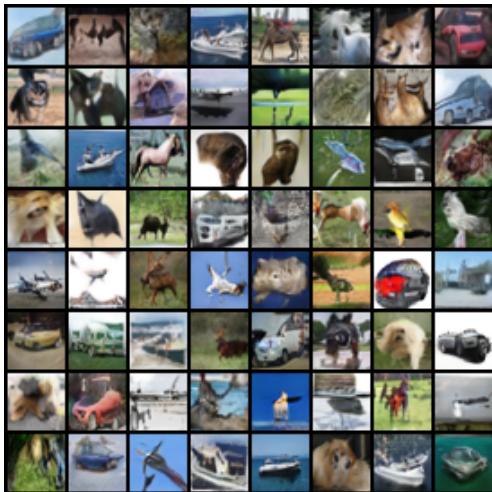
## 2.2 Data Augmentation



Impact of data augmentation on BiggerNet (horizontal flip + random cropping + color jittering)



## 2.2 Synthetic data



The best we could get with a large conditional GAN, spectral normalization on the discriminator, Hinge loss, and exponential moving average of the generator's weights.

### 3. Randomness

**Key idea:** explore stochastic defenses to improve robustness.

**Methods implemented:**

► **Dropout**

► **PNI — Parametric Noise Injection**

- Gaussian noise added to weight tensor
- Learnable scaling factor  $\alpha_l$

$$\tilde{v}_{l,i} = v_{l,i} + \alpha_l \eta_{l,i}, \quad \eta_{l,i} \sim \mathcal{N}(0, \sigma^2)$$

► **RSE — Random Self-Ensemble**

- Gaussian noise layer before each convolution
- Two fixed noise levels:  $\sigma_{\text{init}}$  and  $\sigma_{\text{inner}}$
- Noise parameters are *not* learned

### 3. Randomness — Experimental Results

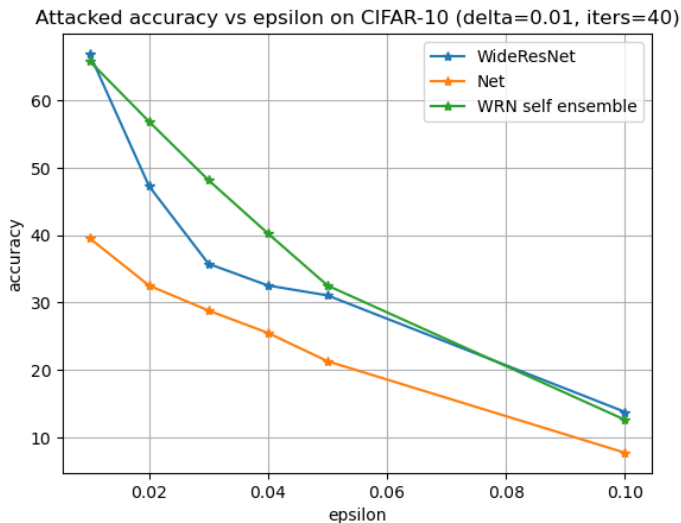
#### PGD-based adversarial training:

Model	BigNet	BigNet-PNI	BigNet-RSE
Natural acc.	53.32%	51.66%	49.41%
Attacked acc.	31.15%	31.07%	<b>32.50%</b>

*Results of adversarial training for BigNet and randomized defense variants. (10 epochs,  $\epsilon = 0.03$ ,  $\delta = 0.008$ )*

- ▶ Small scale experiments first
- ▶ PNI: no significant improvement
- ▶ RSE: best attacked accuracy, slight clean accuracy drop

### 3. Randomness — Experimental Results



Base Net, WRN, WRN-RSE against PGD attacks (no data augmentation)

## 4. Towards Black Box attack

**Setting:** True gradients are not accessible.

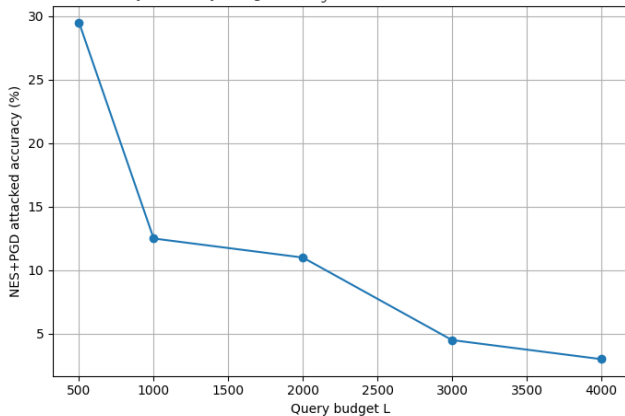
**Method:** NES Gradient Estimator

$$\hat{\nabla} P(y|x) = \frac{1}{2n\sigma} \sum_{i=1}^n [P(y|x + \sigma u_i) - P(y|x - \sigma u_i)] u_i.$$

- ▶  $P(y|x)$  classifier probability for class  $y$  given  $x$
- ▶  $u_i \sim \mathcal{N}(0, I)$  are isotropic Gaussian directions
- ▶  $\sigma$  the search variance.
- ▶ Given a total query budget  $L$  and a gradient budget  $N$ , perform  $\frac{L}{N}$  steps of PGD.

## 4. Towards Black Box attack

Attacked Accuracy vs Query Budget ( $N_{images} = 200$ ,  $\delta = 0.02$ ,  $\varepsilon = 0.1$ ,  $N = 100$ ,  $\sigma = 0.001$ )



## 5. Other Failed Improvement

- ▶ Spectral normalization
- ▶ Low-rank PCA projection
- ▶ Gaussian noise on inputs