

# Adversarial Machine Learning - II

Heng Hsu

NIC Lab Group Meeting

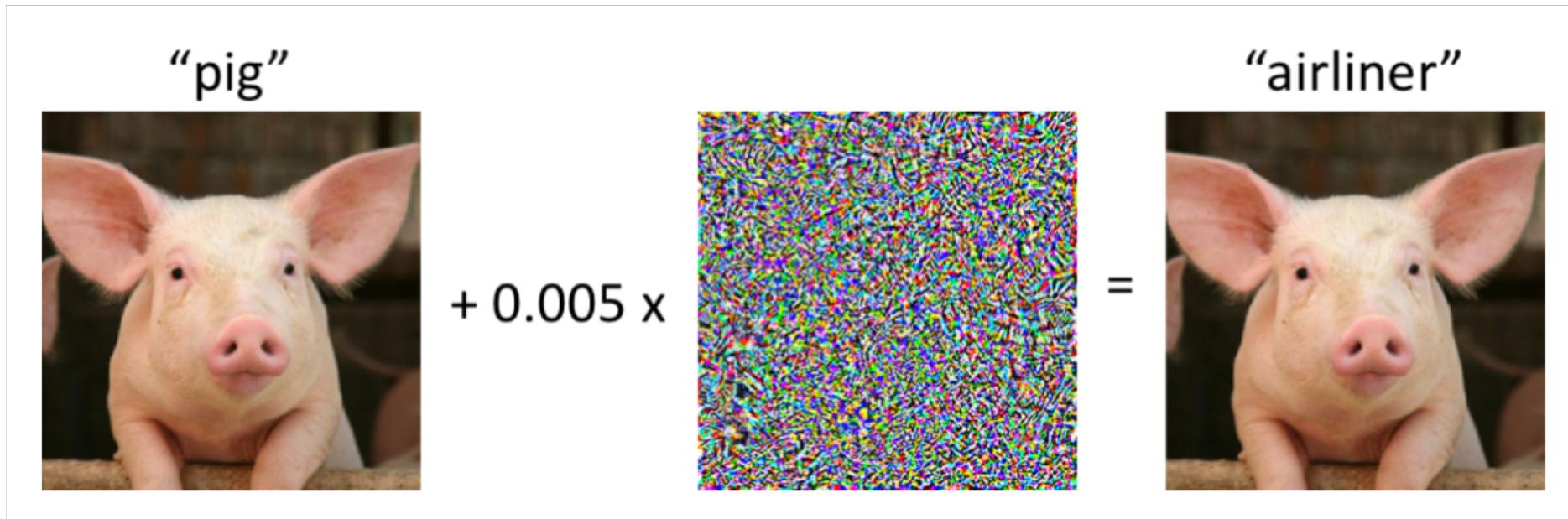
2018/10/30

# Outline

- Quick review
- Two state-of-the-art methods
  - Ensemble Adversarial Training
  - Adversarial Logit Pairing
- Summary

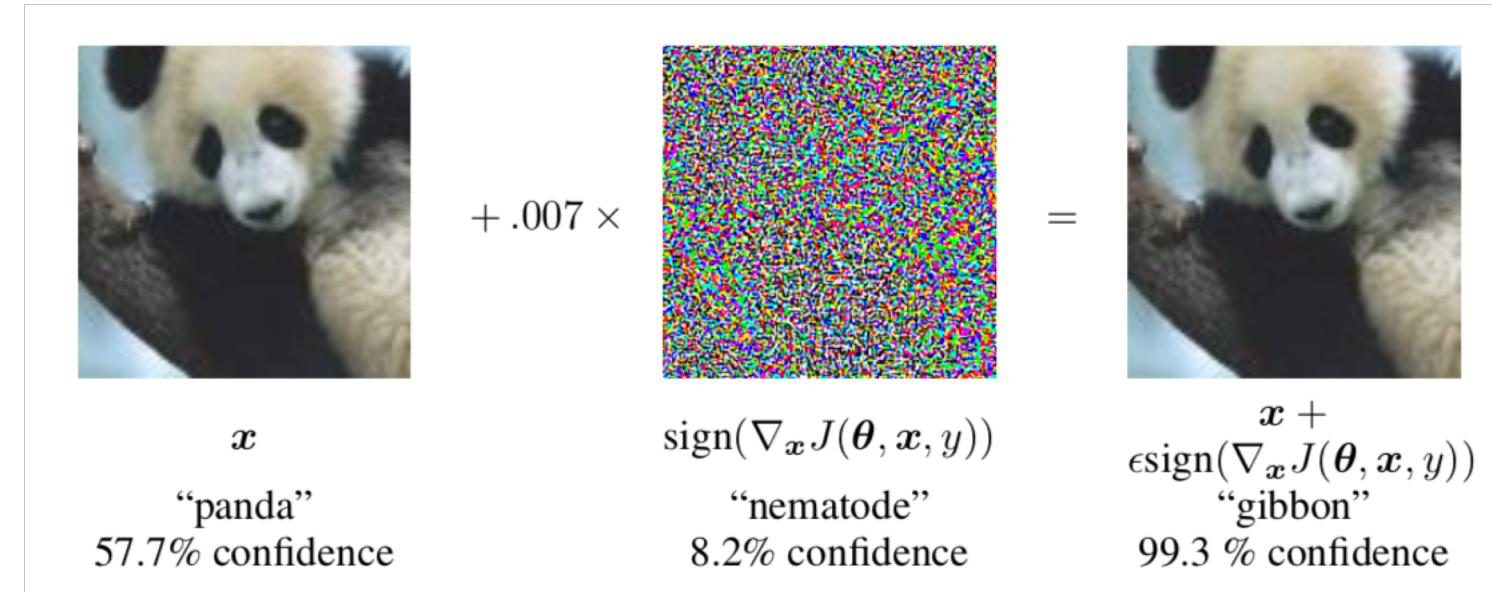
# Adversarial examples

- Adversarial examples are inputs to machine learning models that an attacker has **intentionally designed** to cause the model to make a mistake.<sup>[1]</sup>



# FGSM<sup>[2]</sup>

$$\boldsymbol{\eta} = \epsilon \text{sign}(\nabla_{\mathbf{x}} J(\boldsymbol{\theta}, \mathbf{x}, y))$$



# Gradient Masking (Obfuscated Gradients)

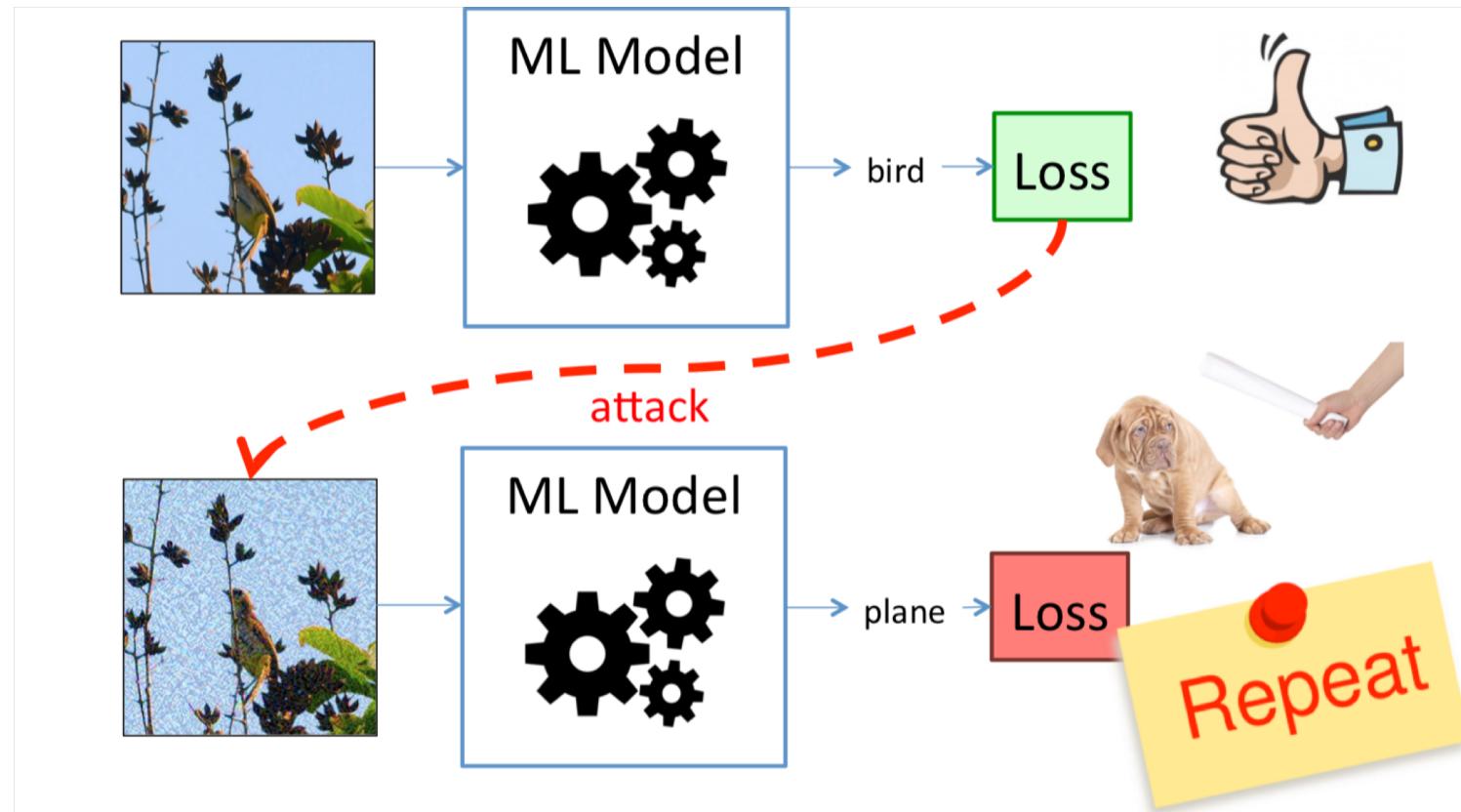
- Some defenses look like they work because they break gradient-based white box attacks.
- The defense denies the attacker access to a useful gradient but these defenses are still vulnerable.
- Obfuscated Gradients Give a False Sense of Security: Circumventing Defenses to Adversarial Examples<sup>[3]</sup>

# Two state-of-the-art methods

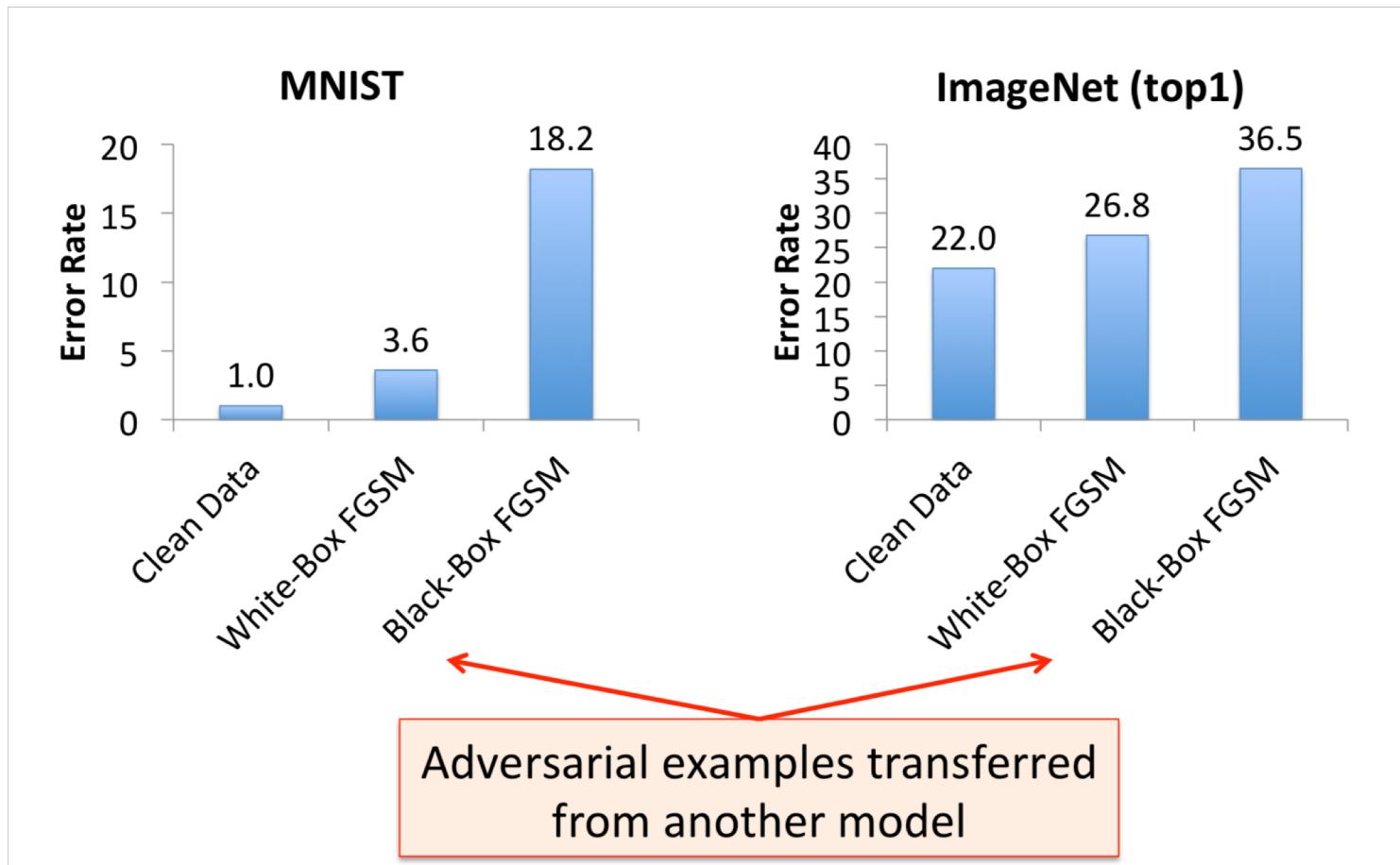
- Ensemble Adversarial Training
- Adversarial Logit Pairing

# Ensemble Adversarial Training<sup>[4]</sup>

- Vanilla adversarial training:

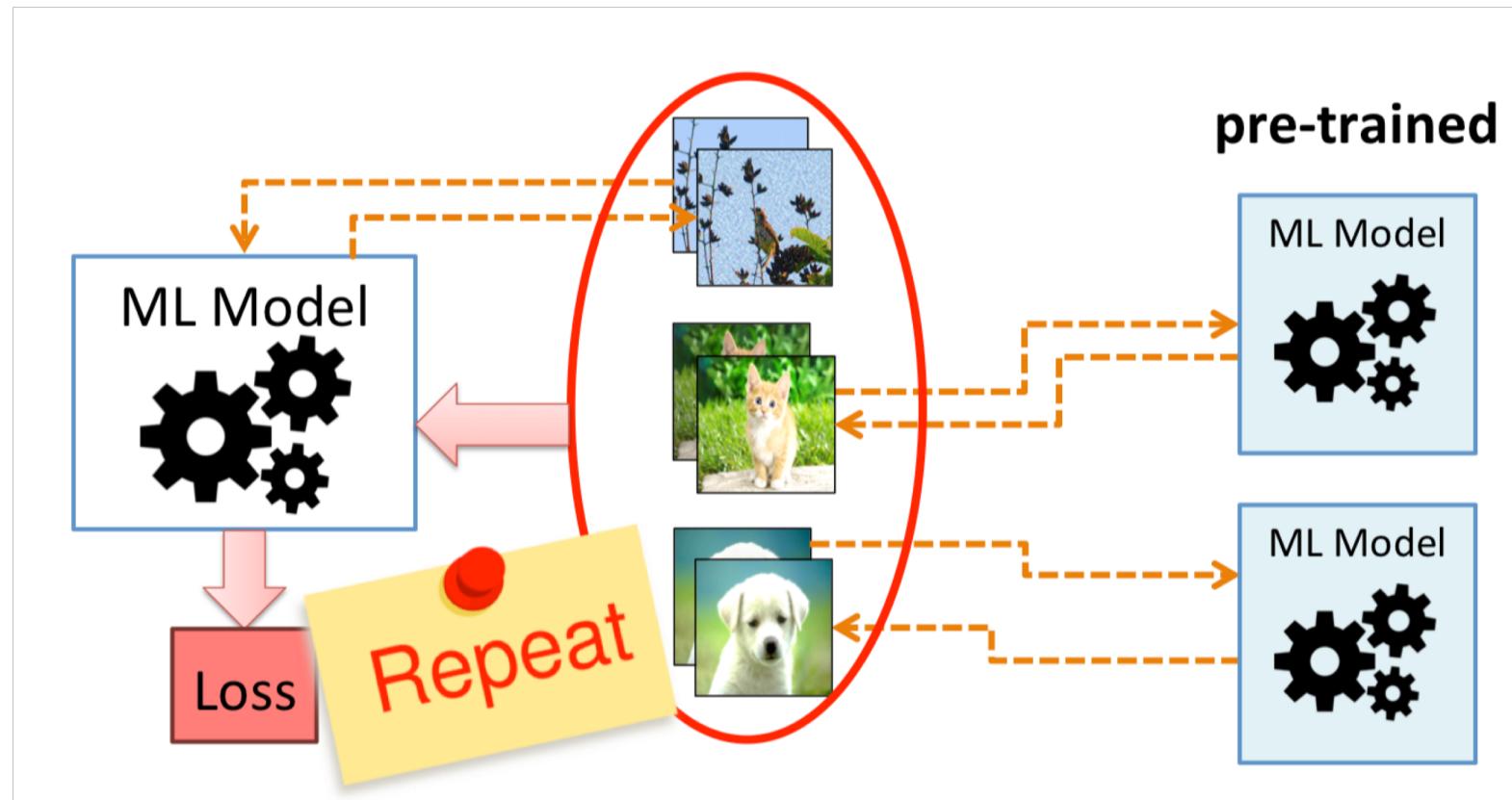


# Vanilla Adversarial Training

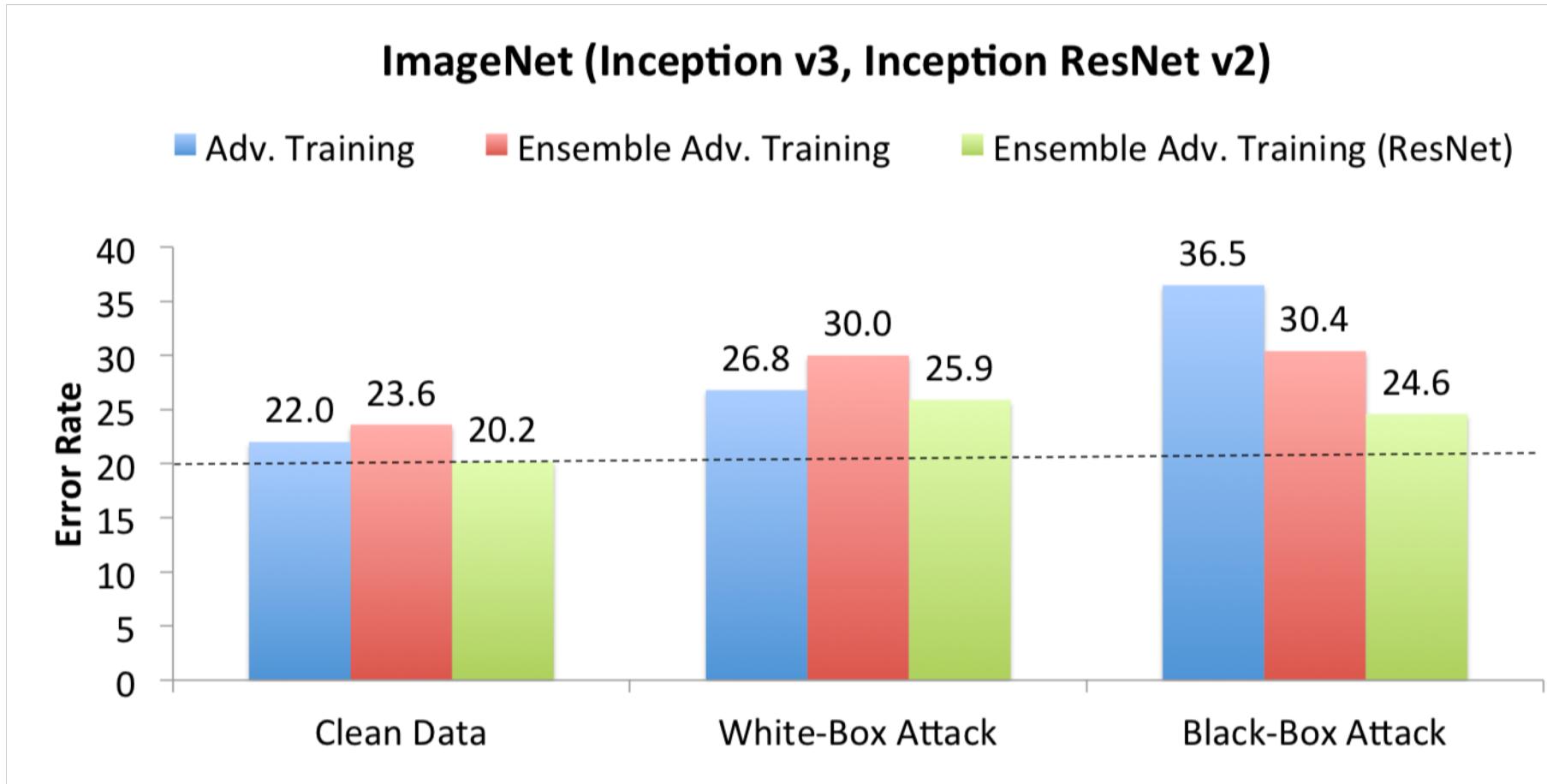


# Defense – Ensemble Adversarial Training

- To *decouple* attack and defense

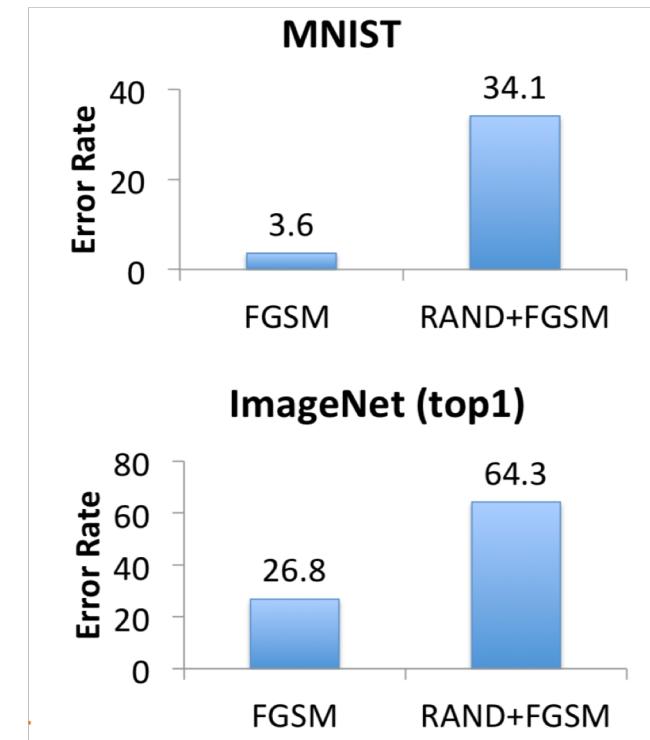
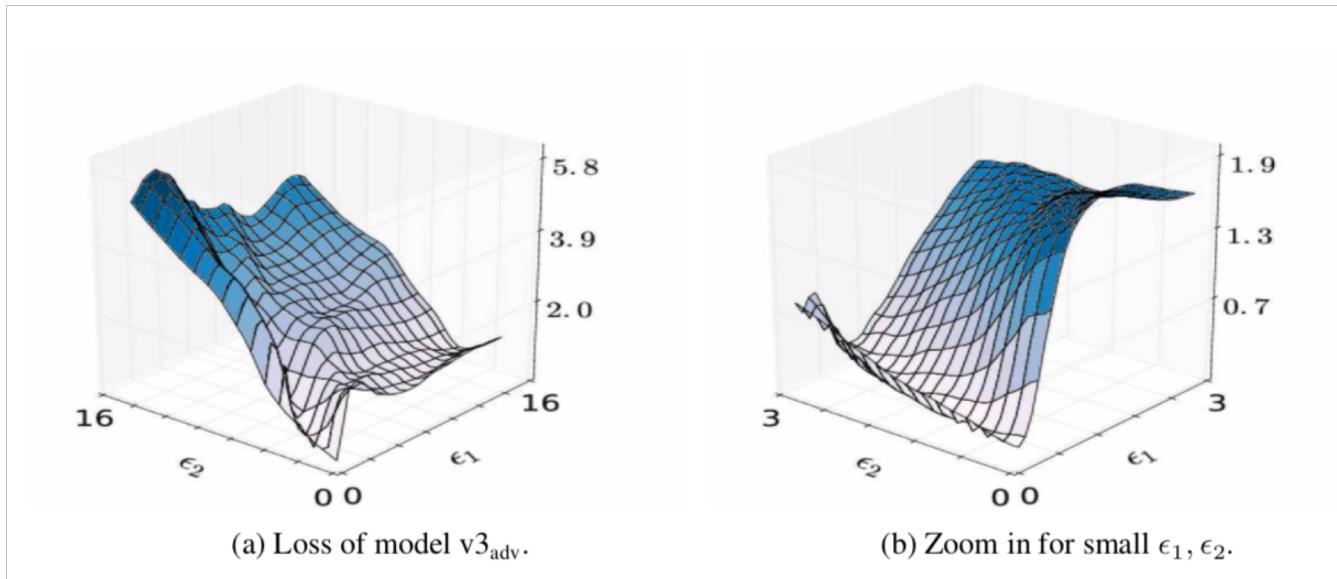


# Results



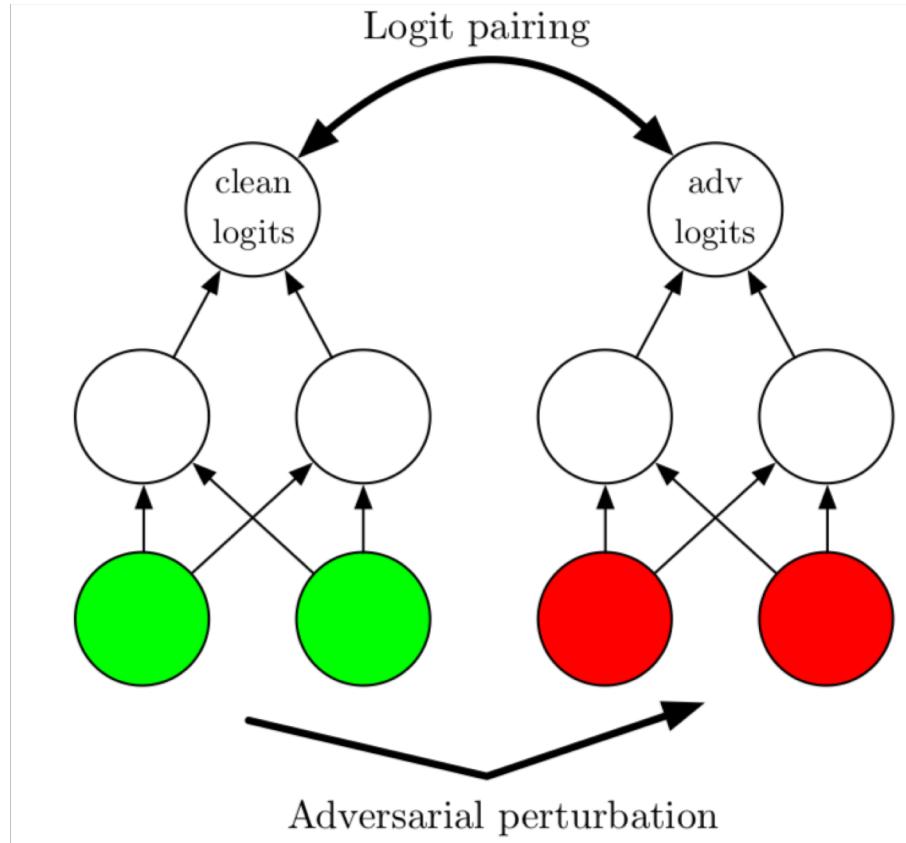
# Attack – RAND + FGSM

- From previous slide, we find adversarial training can defend white-box attack but black-box attack.
- They proposed a white-box attack to beat adversarial training:
  1. Small random step
  2. Step in direction of gradient



# Adversarial Logit Pairing<sup>[5]</sup>

- Logit pairing:



# Adversarial Logit Pairing

- Vanilla adversarial training:

$$\theta^* = \arg \min_{\theta} [\mathcal{L}_{\text{cross entropy}}(x, t) + \mathcal{L}_{\text{cross entropy}}(x', t)]$$

- Adversarial logit pairing:

$$\theta^* = \arg \min_{\theta} [\mathcal{L}_{\text{cross entropy}}(x, t) + \mathcal{L}_{\text{cross entropy}}(x', t) + \lambda \mathcal{L}_{\text{logit pairing}}(f(x), f(x'))]$$

# Results

- MNIST
- SVHN

Method	White Box	Black Box	Clean
M-PGD	93.2%	96.0%	98.5%
ALP	<b>96.4%</b>	<b>97.5%</b>	<b>98.8%</b>

Method	White Box	Black Box	Clean
M-PGD	44.4%	55.4%	<b>96.9%</b>
ALP	<b>46.9%</b>	<b>56.2%</b>	96.2%

# Summary

- Today: enhanced adversarial training
- State-of-the-art but quite **heuristic**
- Proposed research direction:
  - Mathematical perspective explanation ?

**THANK YOU**