

# Adversarial Robustness

M. Soleymani

Sharif University of Technology

Spring 2024

Most slides have been adapted from:

Zico Kolter, Aleksander Madry, Adversarial ML Tutorial

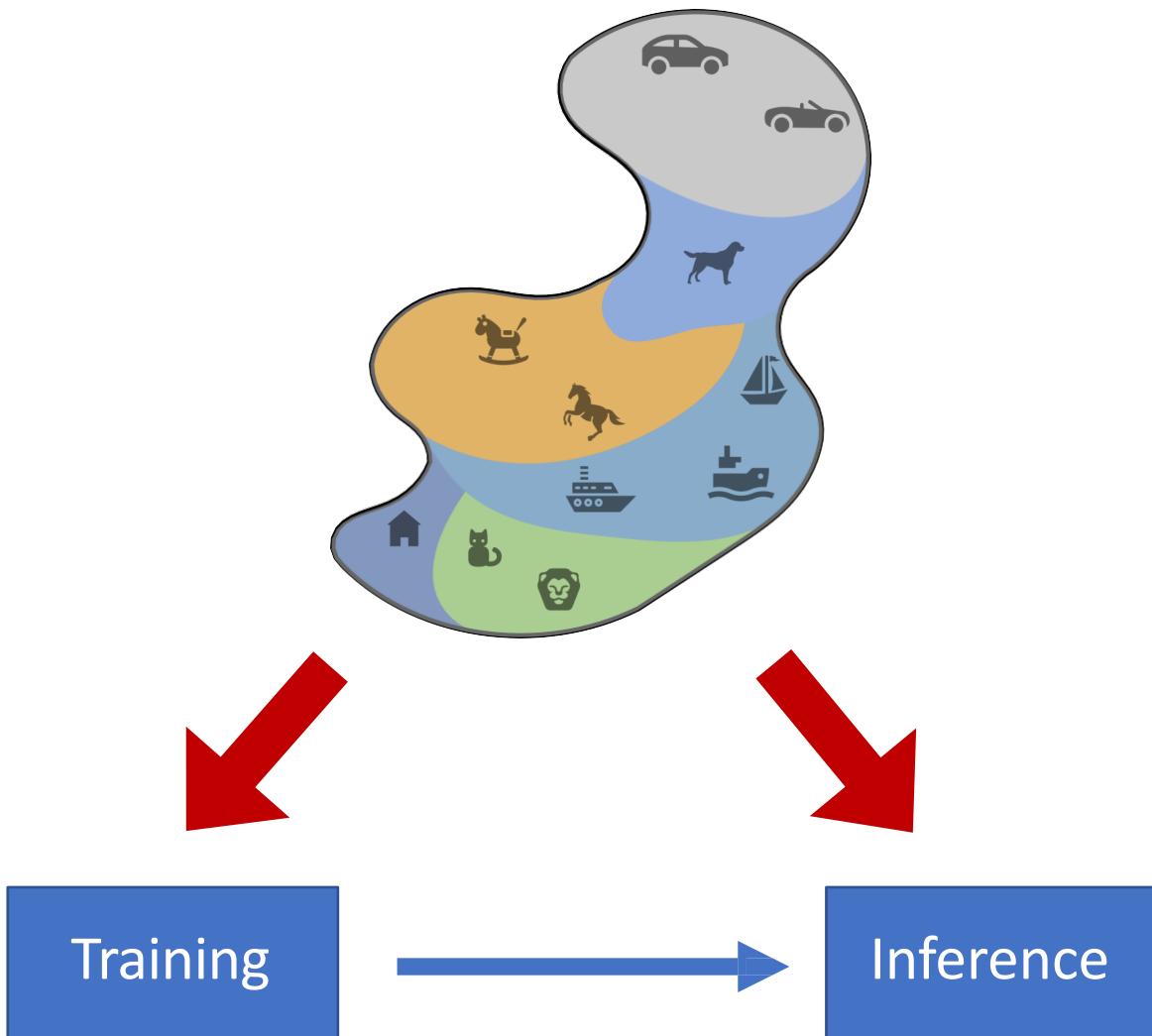
Anupam Datta, John Mitchell; Stanford CS329T

# Outline

- Machine learning models can be very brittle
- What are adversarial attacks? How to generate adversarial examples?
- There are systematic methods to improve robustness
- There are limitations in the current state of the art/science
- Other types of threats for ML methods

Is ML **truly** ready for  
real-world  
deployment?

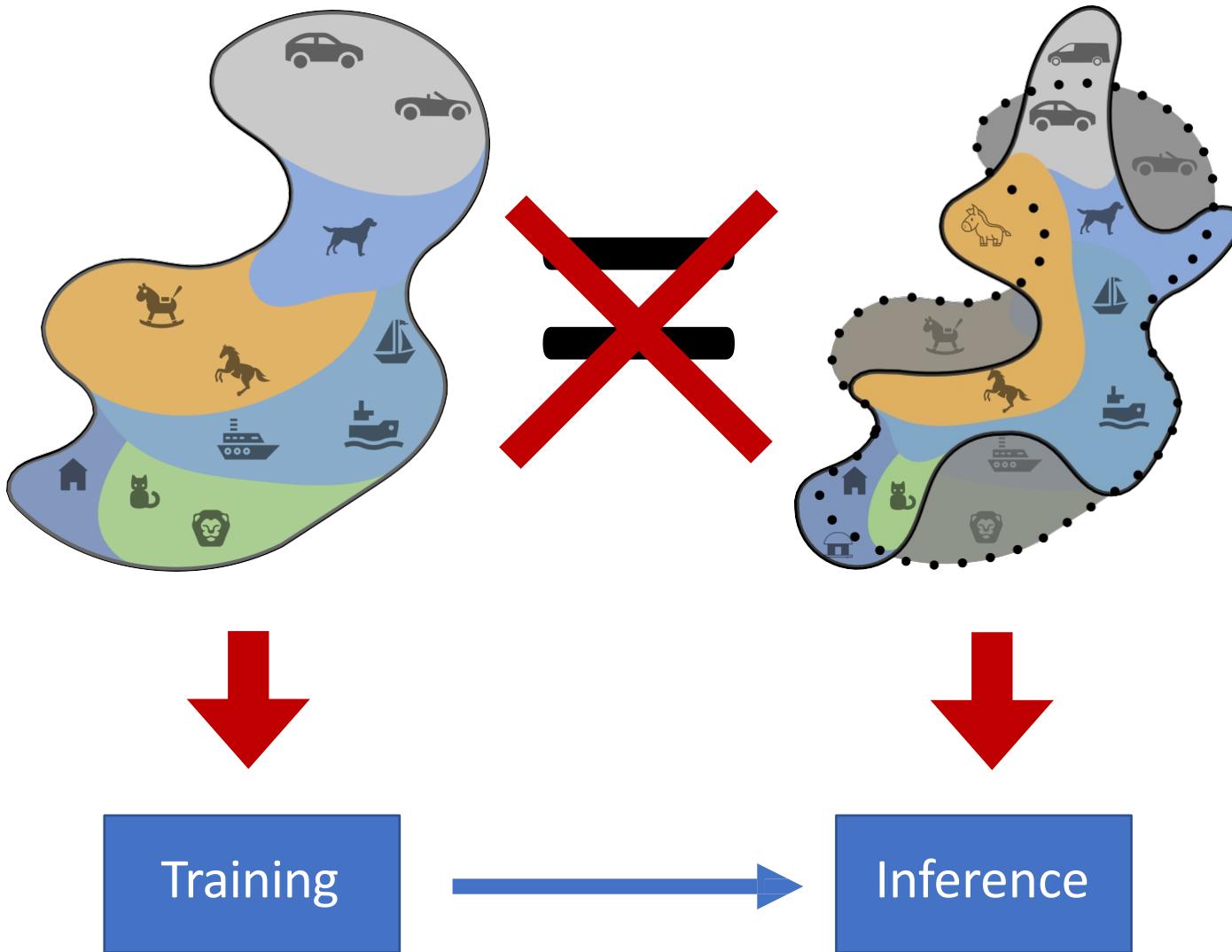
# A Limitation of the (Supervised) ML Framework



**Measure of performance:**  
Fraction of mistakes during testing

**But:** In reality, the distributions we **use** ML on are NOT the ones we **train** it on

# A Limitation of the (Supervised) ML Framework



**Measure of performance:**  
Fraction of mistakes during testing

**But:** In reality, the distributions we **use** ML on are NOT the ones we **train** it on

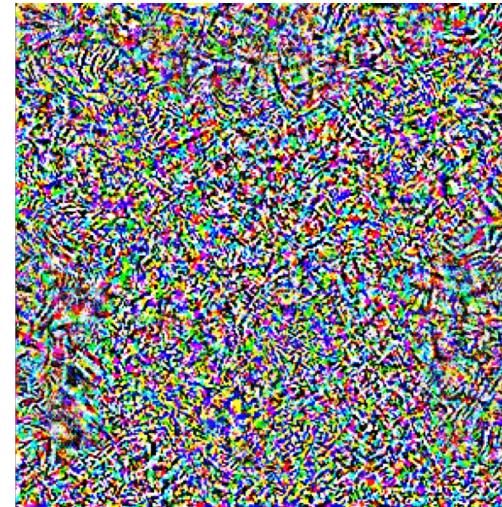
What can go wrong?

# ML Predictions Are (Mostly) Accurate but Brittle

“pig” (91%)



noise (NOT random)



“airliner” (99%)



+ 0.005 x

=

Adversarial examples are imperceptible to humans

[Szegedy Zaremba Sutskever Bruna Erhan Goodfellow Fergus 2013]

[Biggio Corona Maiorca Nelson Srndic Laskov Giacinto Roli 2013]

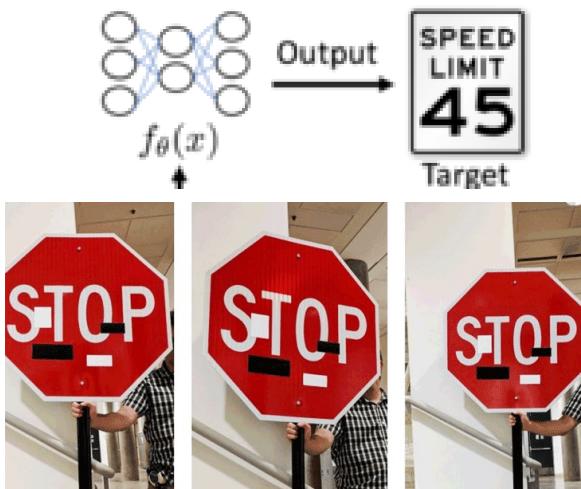
# ML Predictions Are (Mostly) Accurate but Brittle



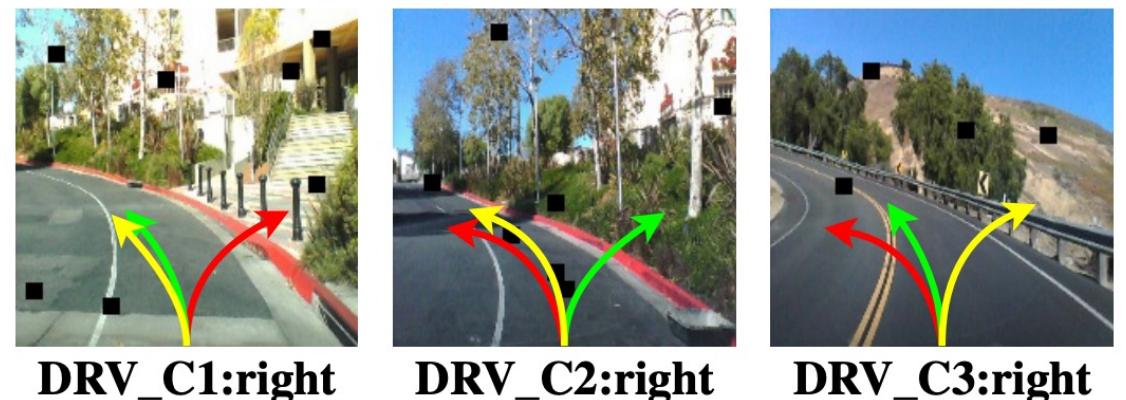
[Kurakin, Goodfellow, Bengio, 2017]



[Athalye et al., 2017]

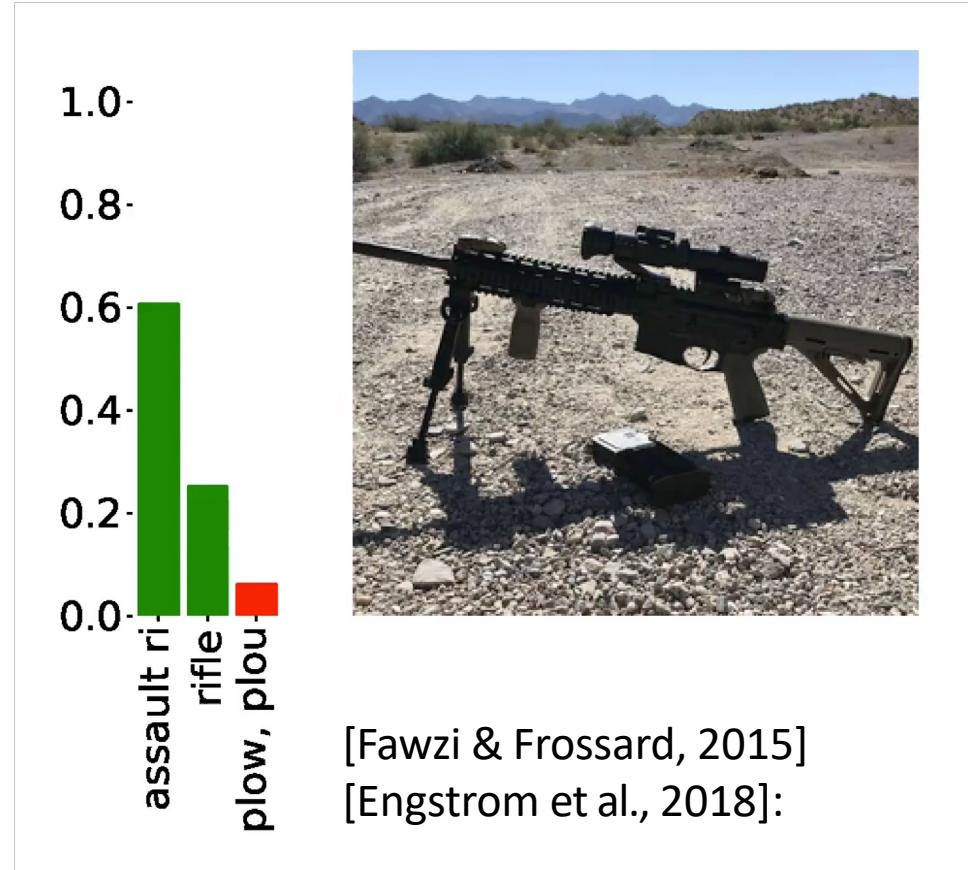


[Eykholt et al., 2017]



[Pei et al., DeepXplore, 2017]

# ML Predictions Are (Mostly) Accurate but Brittle



Rotation + Translation suffices to fool state-of-the-art vision models

→ Data augmentation does not seem to help here either

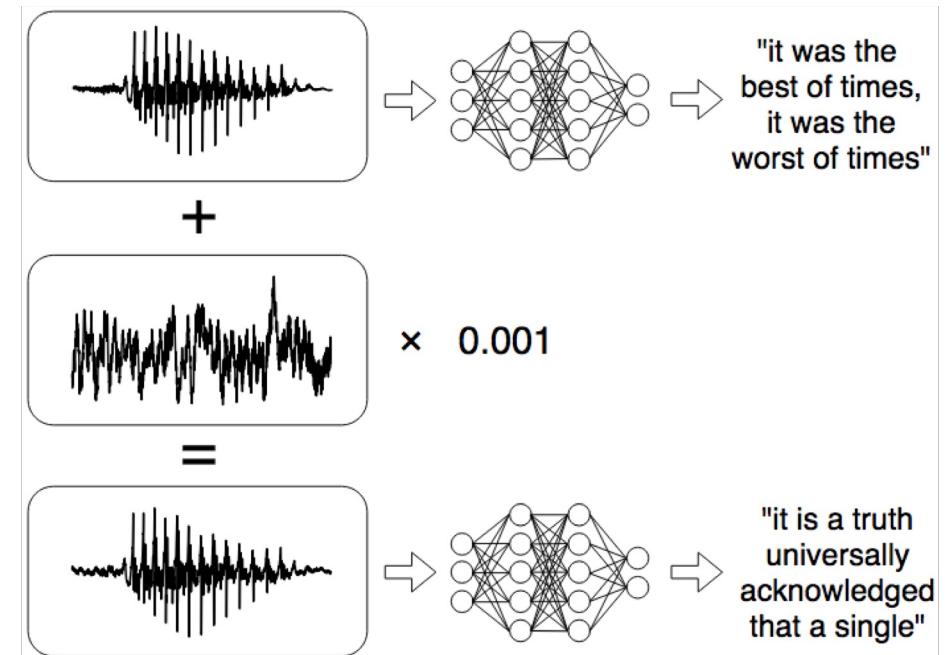
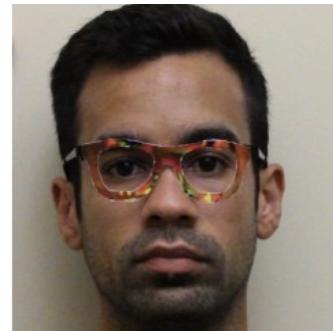
So: Brittleness of ML is a thing

Should we be worried?

# Why Is This Brittleness of ML a Problem?

→ Security

[Carlini Wagner 2018]:  
Voice commands that are  
unintelligible to humans

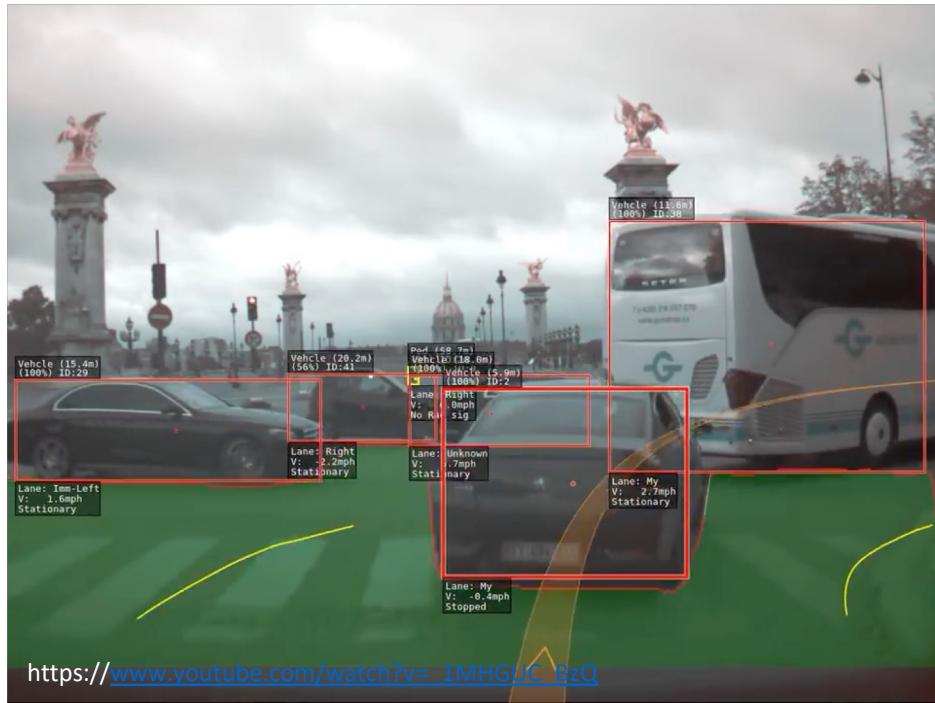


[Sharif et al., 2016]:  
Glasses that fool face recognition

# Why Is This Brittleness of ML a Problem?

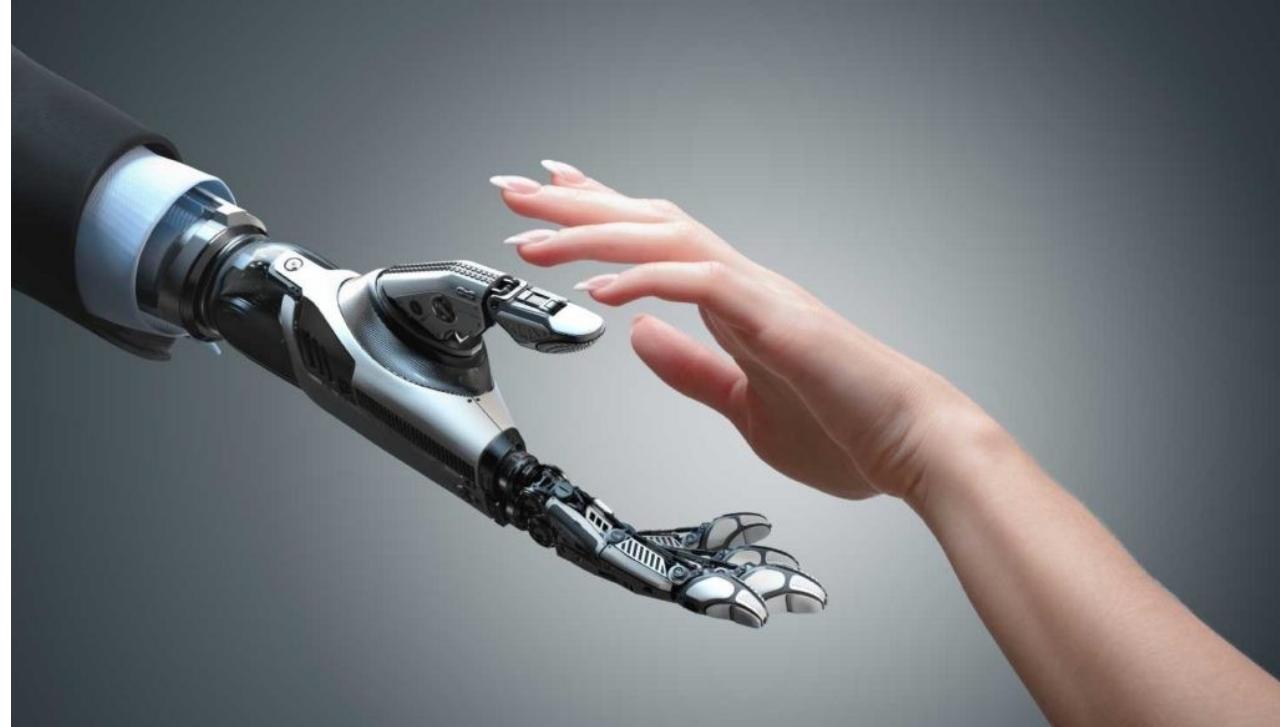
→ Security

→ Safety



# Why Is This Brittleness of ML a Problem?

- Security
- Safety
- ML Alignment



Need to understand the  
“failure modes” of ML

# Adversarial attack

- Given an input image  $X$  and a label  $t$ , find  $X' \approx X$  with  $f(X') = t$



Dog



Hummingbird

This is a strong form of attack: for *any*  $X$  and *any*  $t$ , find  $X'$  with predicted class  $t$

# Targeted Attack vs. untargeted Attack

- Targeted Attack aims to misclassify the input to a specific target label
  - Input: a deep model  $C_\theta(\cdot)$ , an input  $x$ , and a target class  $t$
  - Output: A perturbation  $\delta$  such that  $C_\theta(x + \delta) = t$
- Untargeted Attack aims to misclassify the input to any wrong label
  - Input: a deep model  $C_\theta(\cdot)$  and an input  $x$
  - Output: A perturbation  $\delta$  such that  $C_\theta(x + \delta) \neq C_\theta(x)$

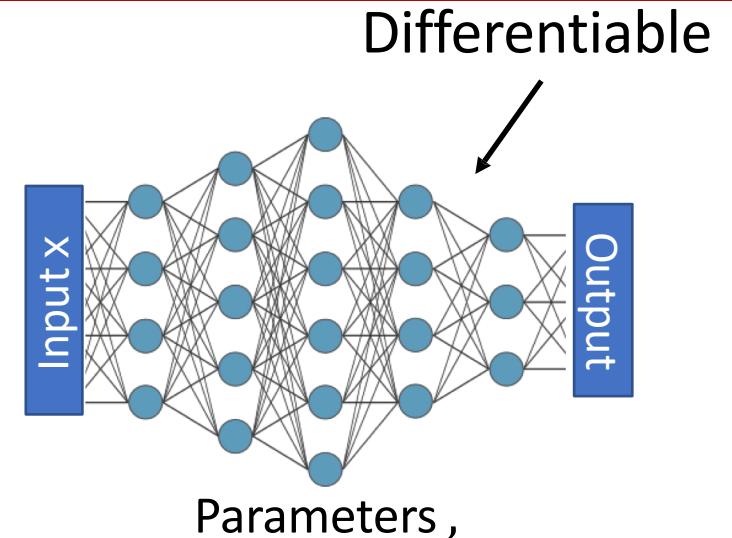
# Where Do Adversarial Examples Come From?

To get an adv. example

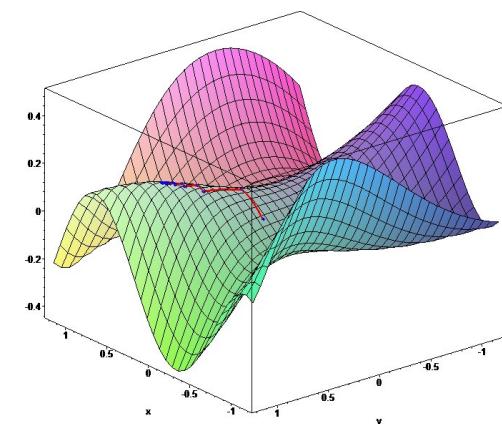
**Goal of training:**

Input      Parameters      Correct Label

$\min_{\theta} \text{loss}(f(x; \theta), y)$



Can use gradient descent method to find good  $\theta$

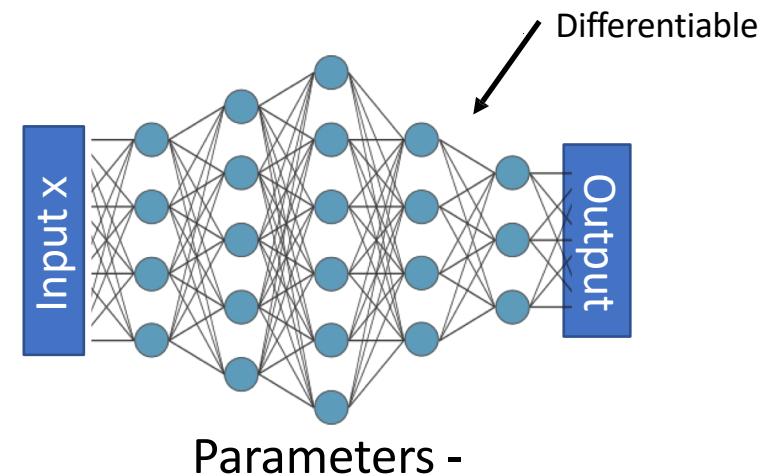


# Where Do Adversarial Examples Come From?

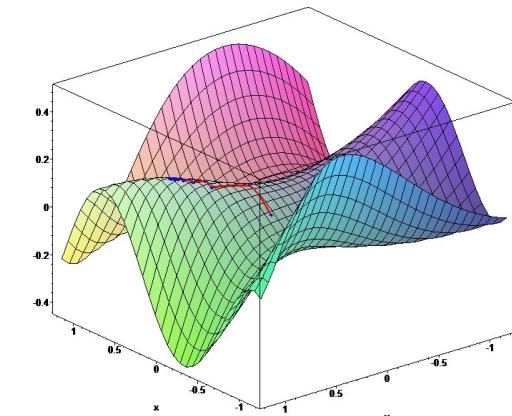
To get an adv. example

Goal of training:

$$\max_{\delta \in \Delta} \text{loss}(f(x + \delta; \theta), y)$$



Can use gradient descent method to find good  $\delta$



# Where Do Adversarial Examples Come From?

To get an adv. example

Goal of training:

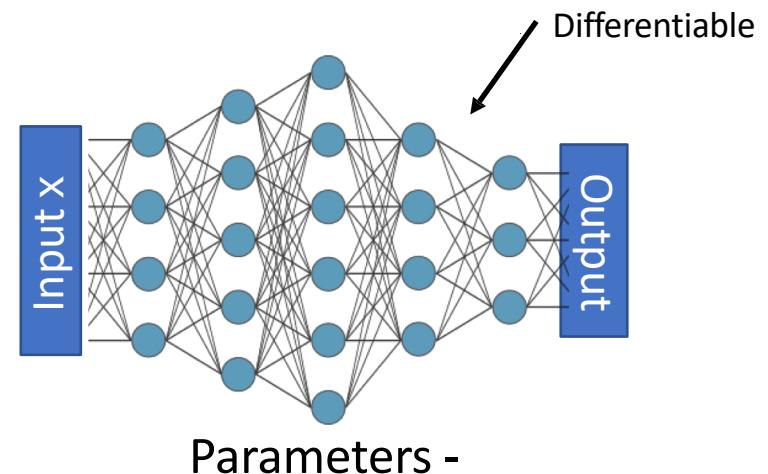
$$\max_{\delta \in \Delta} \text{loss}(f(x + \delta; \theta), y)$$

Which  $\delta$  are allowed?

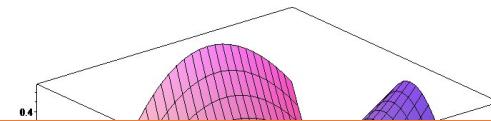
Examples: that is small wrt:

- $L_P$ -norm
- Rotation and/or translation
- VGG feature perturbation

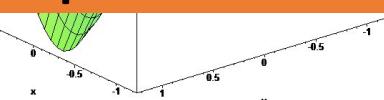
(add the perturbation you need here)



Can use gradient descent method to find good  $\delta$



**Still:** We have to confront  
(small)  $L_P$  norm perturbations



# Fast Gradient Sign Method (FGSM)

- After training the network, for each  $(x, y)$ , do one gradient step to increase  $loss(\theta; x + \delta, y)$  (untargeted attack):

$$x_{adv} = x + \epsilon \text{sign}(\nabla_x loss(\theta; x + \delta, y))$$

or specifically consider the target class  $t$  (targeted attack):

$$x_{adv} = x - \epsilon \text{sign}(\nabla_x loss(\theta; x + \delta, t))$$

Attacks can be targeted or untargeted  
(the class of  $x_{adv}$  is given or not)

- $L_\infty$  bound on the perturbation magnitude

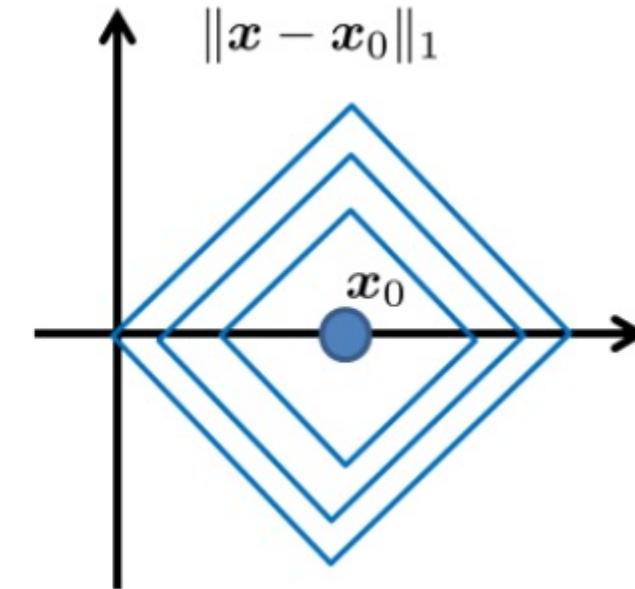
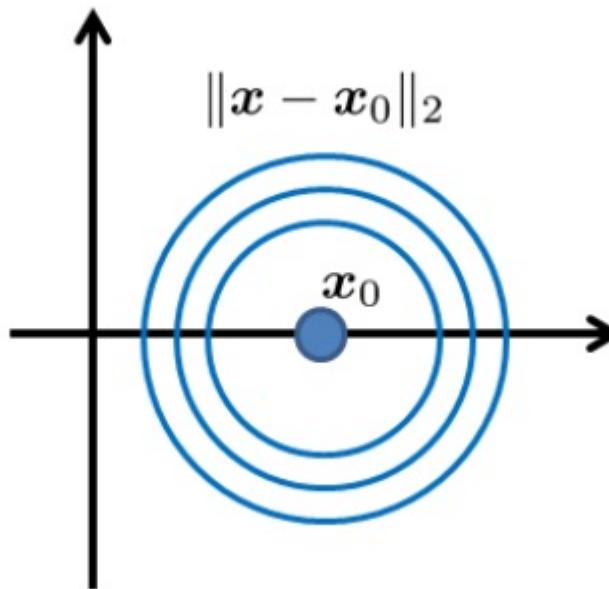
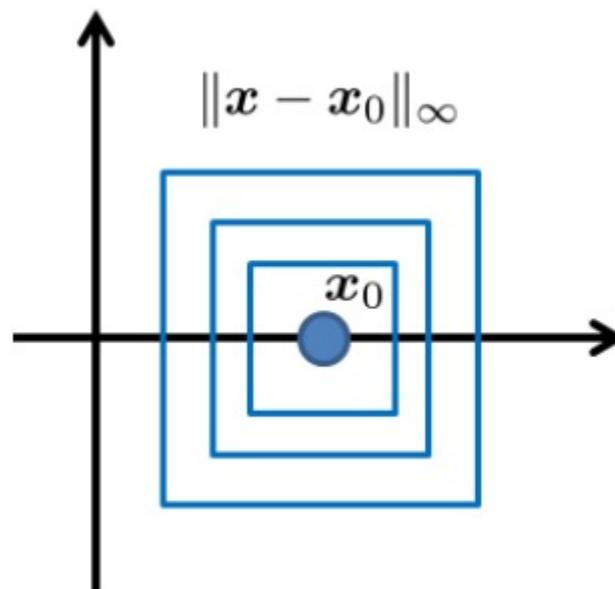
- It is 100% successful on IMAGENET even for small  $\epsilon$ .

# Human Alignment: Metrics used for robustness

- Do metrics match human perceptibility?

- $L_0$ : how many pixels changed at all
- $L_1$ : total absolute pixel change
- $L_2$ : Euclidean distance
- $L_\infty$ : max change of any single pixel

$$\|x - x'\|_p = \sqrt[p]{\sum_i (x_i - x'_i)^p}$$

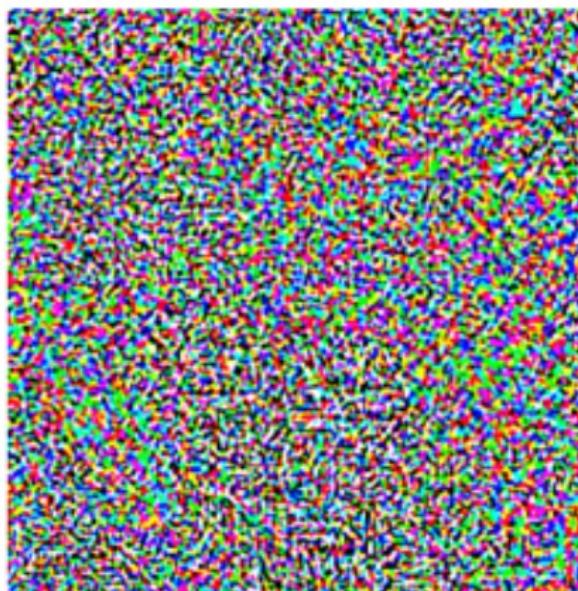


# FGSM attack: Example



$\mathbf{x}$   
“panda”  
57.7% confidence

$$+ .007 \times$$



$\text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$   
“nematode”  
8.2% confidence

=



$\mathbf{x} +$   
 $\epsilon \text{sign}(\nabla_{\mathbf{x}} J(\theta, \mathbf{x}, y))$   
“gibbon”  
99.3 % confidence

# FGSM

- FGSM is a 1-step attack
- $x + \delta$  is guaranteed to stay inside the box  $[x - \epsilon, x + \epsilon]$
- FGSM was designed to be fast, not optimal
  - may not compute minimal perturbation
  - is not a strong attack as shown later

# Carlini & Wagner attack

- Basic idea: Given input  $x$  and desired label  $t$ , find adversarial input  $x + \delta$

$$\min_{\delta} D(x + \delta, x)$$

$$s.t. C(x + \delta) = t$$

$$x + \delta \in [0,1]^n$$

- Problem: This constraints in the optimization problem does not work well
- Solution: Define  $g$  such that  $C(x + \delta) = t$  if and only if  $g_t(x + \delta) \leq 0$

$$\min_{\delta} D(x + \delta, x) + \gamma g_t(x + \delta)$$

$$s.t. x + \delta \in [0,1]^n$$

$$g_t(x) = \left( \max_{i \neq t} f_i(x) - f_t(x) \right)^+$$

# Carlini & Wagner attack: Example

	0	1	2
0	0	0	0
1	1	1	1
2	2	2	2

$L_\infty$  attack

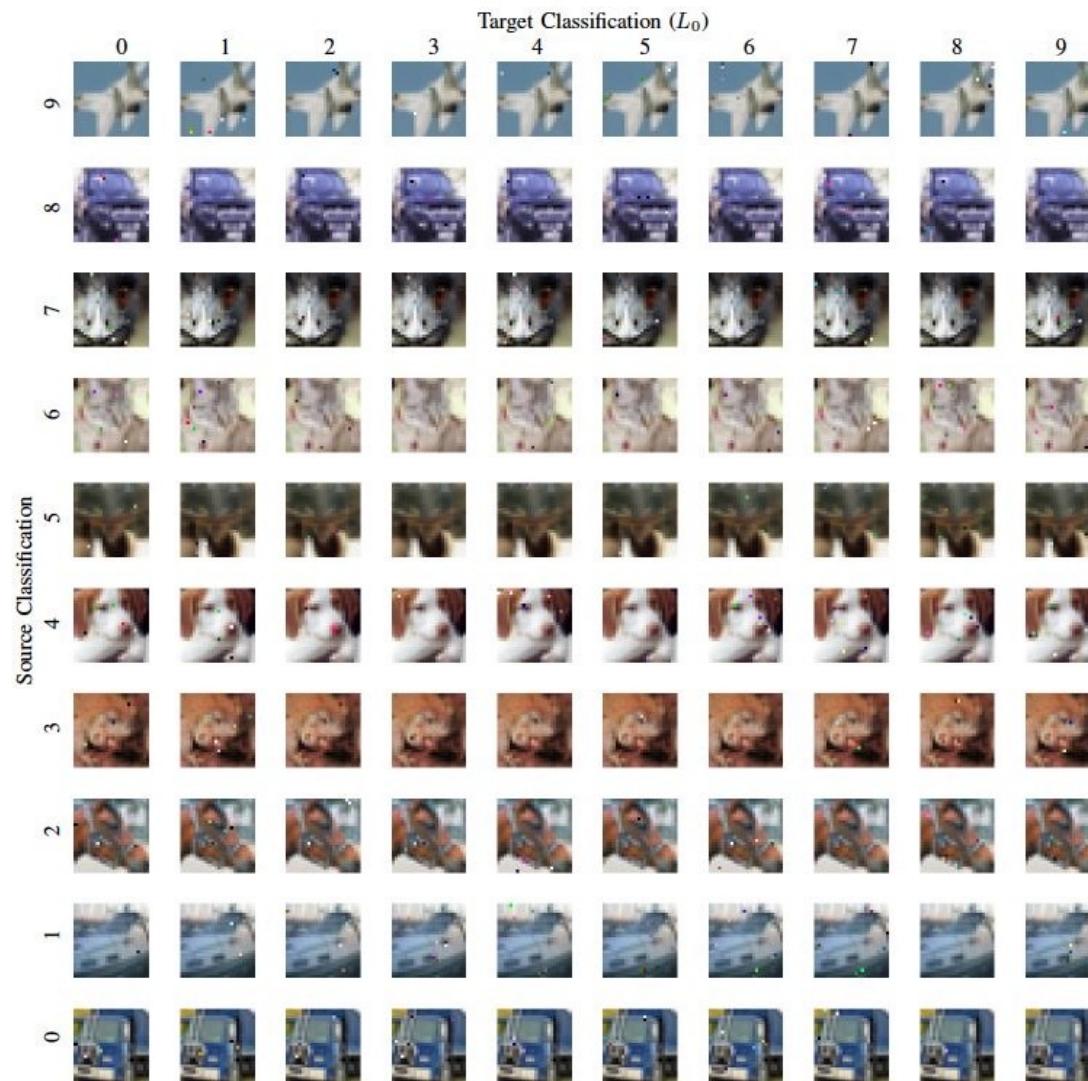
$L_2$  attack

	0	1	2
0	0	0	0
1	1	1	1
2	2	2	2

$L_0$  attack

	0	1	2
0	0	0	0
1	1	1	1
2	2	2	2

# Carlini & Wagner attack: Example



# Projected gradient descent (PGD) attack

- It is an extension of FGSM
- After each step of perturbation, the adversarial example is projected back onto the  $\epsilon$ -ball of  $x$  using a projection function  $\Pi$ :

$$x_{adv}^t = \Pi_\epsilon \left( x_{adv}^{t-1} + \alpha \text{ sign } \nabla_x \text{loss}(\theta; x_{adv}^{t-1} + \delta, y) \right)$$

is smaller than  $\epsilon$

- PGD is regarded as the strongest  $L_\infty$  attack

# Comparing attacks

- **Strongness:** Attack 1 is stronger than Attack 2 in the  $L_p$  distance if it has more ability to generate successful  $x_{adv} \in B_p(x, \epsilon)$
- CW, PGD are the most powerful attacks.
  - There has been very slight improvements since then.
- $L_0, L_2, L_\infty$  are generally imperceptible.

# Attacks: Summary

Attack	Region	Optimization	Outcome
<b>FGSM</b>	Fixed to $\delta \in \{-\epsilon, \epsilon\}^n$	Take one $\epsilon$ -sized step	On the region boundary
<b>PGD</b>	Any region one can project to	Take many steps	Inside region
<b>C&amp;W</b>	No restriction (only $x + \delta \in [0,1]^n$ )	Want to produce $\delta$ with small $L_\infty$ . Takes many steps using LBFGS-B to ensure $x + \delta$ stays in $[0,1]^n$	Hopefully small $L_\infty(\delta)$

# White-box vs. Black-box Attacks

- White box attacks: the attacker knows the model, the parameters, and the architecture
- Black box attacker: the attacker knows the architecture (e.g., the layers) but not its parameters (e.g., weights)
- Adversarial examples are transferrable:
  - Given the same training data as the original network
    - an attacker can train their own mirror network of the black box original network
    - can attack the mirror network with white-box techniques.
  - If attack on mirror network succeeds, it will likely succeed on the original

# Robustness

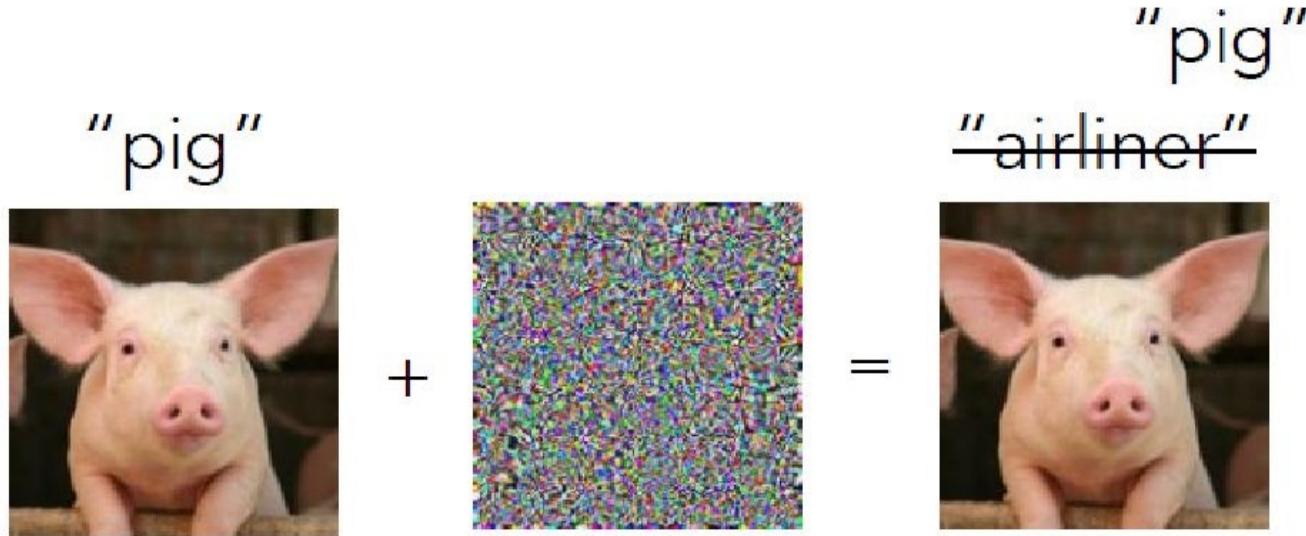
- Robustness: A network is robust if it returns correct output on all inputs
- Impractical: the input space is too large to be covered
- Local Robustness (informal): A learning model is locally-robust if it returns the correct output on inputs similar to inputs in training data

# Accuracy vs. robustness tradeoff

- Accuracy: fraction of inputs which are correctly classified
- $\epsilon$ -robustness score: fraction of inputs  $x$  such that for all  $x' \in B_p(x, \epsilon)$ ,  $\text{class}(x) = \text{class}(x')$
- Highly accurate models may have very low robustness scores

# Systematic methods can increase robustness

- Robust training
  - Train so that small changes in input do not change the classifier output



# Towards ML Models that Are Adv. Robust

**Key observation:** Lack of adv. robustness is **NOT** at odds with what we currently want our ML models to achieve

Standard generalization:

$$E_{(x,y) \sim D} [\text{loss}(\theta; x, y)]$$

**But:** Adversarial noise is a “needle in a haystack”

# Towards ML Models that Are Adv. Robust

**Key observation:** Lack of adv. robustness is **NOT** at odds with what we currently want our ML models to achieve

~~Standard~~ generalization:  $E_{(x,y) \sim D} [\max_{\delta \in \Delta} \text{loss}(\theta; x + \delta, y)]$

Adversarially robust

**But:** Adversarial noise is a “needle in a haystack”

# Use similar framework to train robustly

$$\min_{\theta} \mathbb{E}_{x,y \sim D} [\max_{\delta \in \Delta} \text{loss}(\theta, x + \delta, y)]$$

↑  
finding a robust model

↑  
finding a worst-case perturbation

Improve robustness: Train on perturbed inputs

(aka “adversarial training” [Goodfellow et al. 2015])

Actually leads to **robust models** (with some care)

# Questions for discussion

- Do you think this goal captures “robustness” correctly?
  - A small change in input should not produce arbitrary large changes in output
- Are there other goals you associate with “robustness”?
- Does you think it is always better to build modules by perturbing input slightly, as shown in last few slides?
- What kind of input changes qualify as “slight perturbation”?

$$\max_{\delta \in \Delta} \text{loss}(\theta, x + \delta, y)$$



Allowed perturbations: pixel-wise, rotations, ...

# Adversarial Robustness Beyond Security

## ML via Adversarial Robustness Lens

**Overarching question:**

How does adv. robust ML differ from “standard” ML?

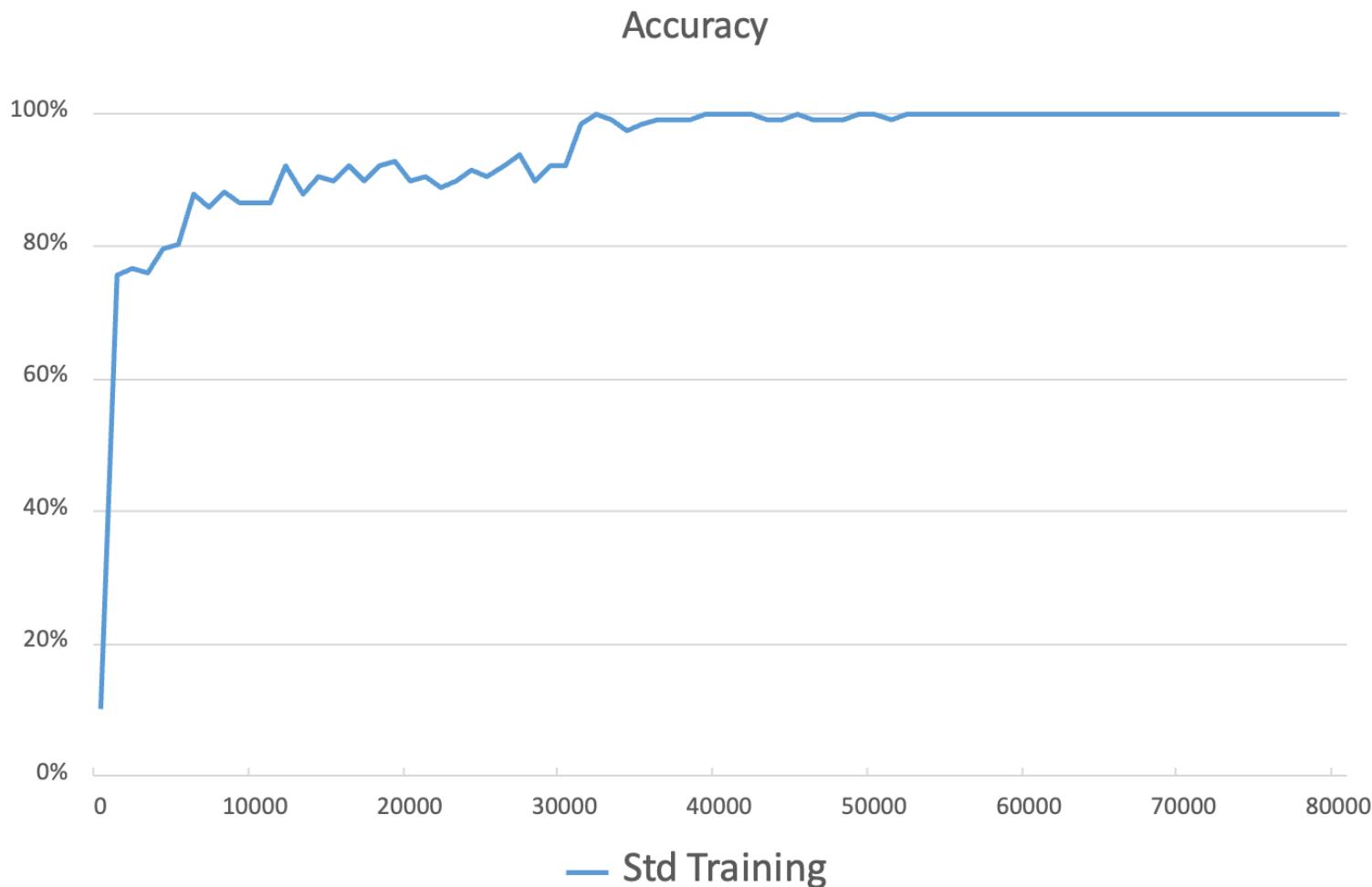
$$E_{(x,y) \sim D} [\text{loss}(\theta; x, y)]$$

vs

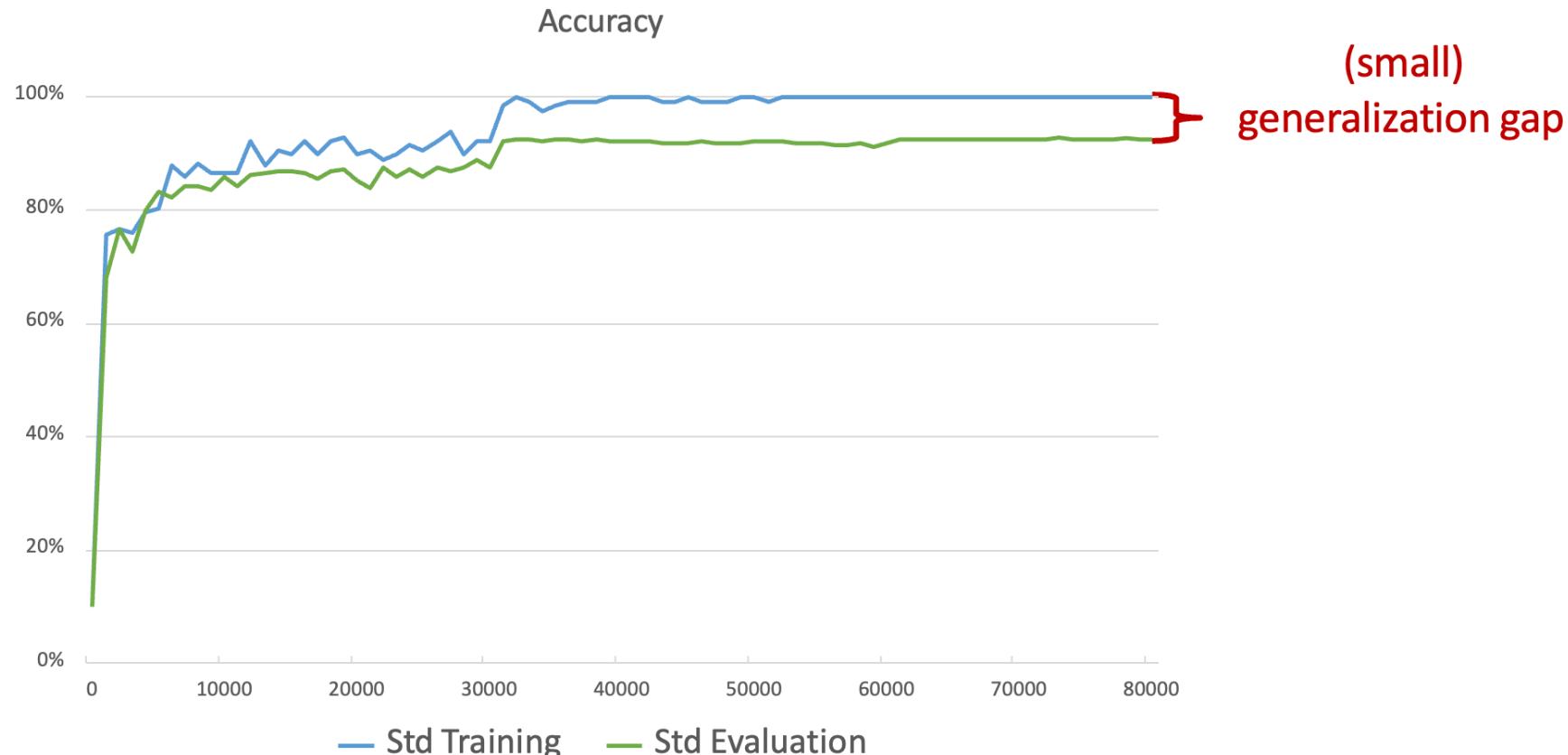
$$E_{(x,y) \sim D} [\max_{\delta \in \Delta} \text{loss}(\theta; x + \delta, y)]$$

(This goes **beyond** deep learning)

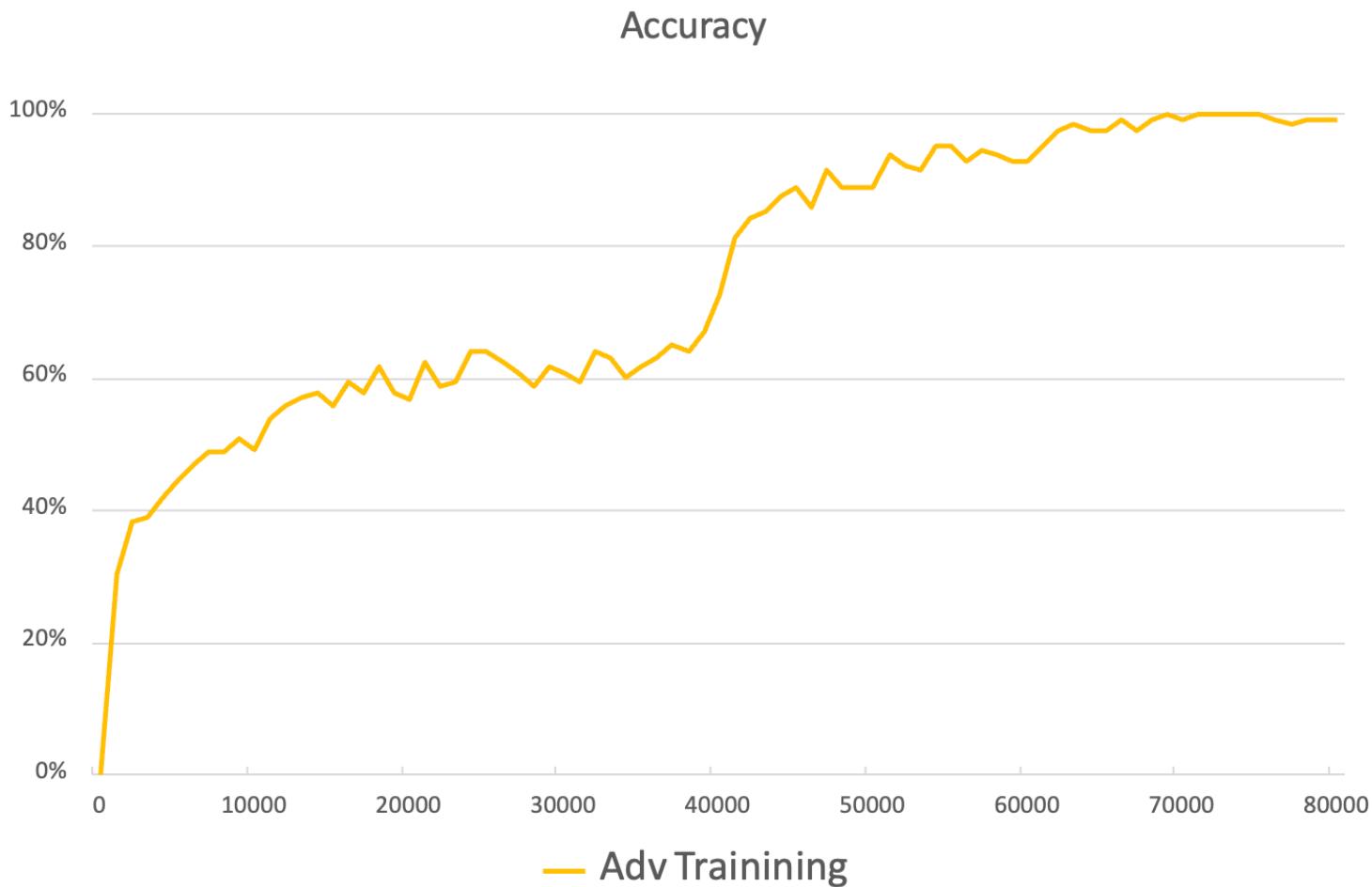
# Do Robust Deep Networks Overfit?



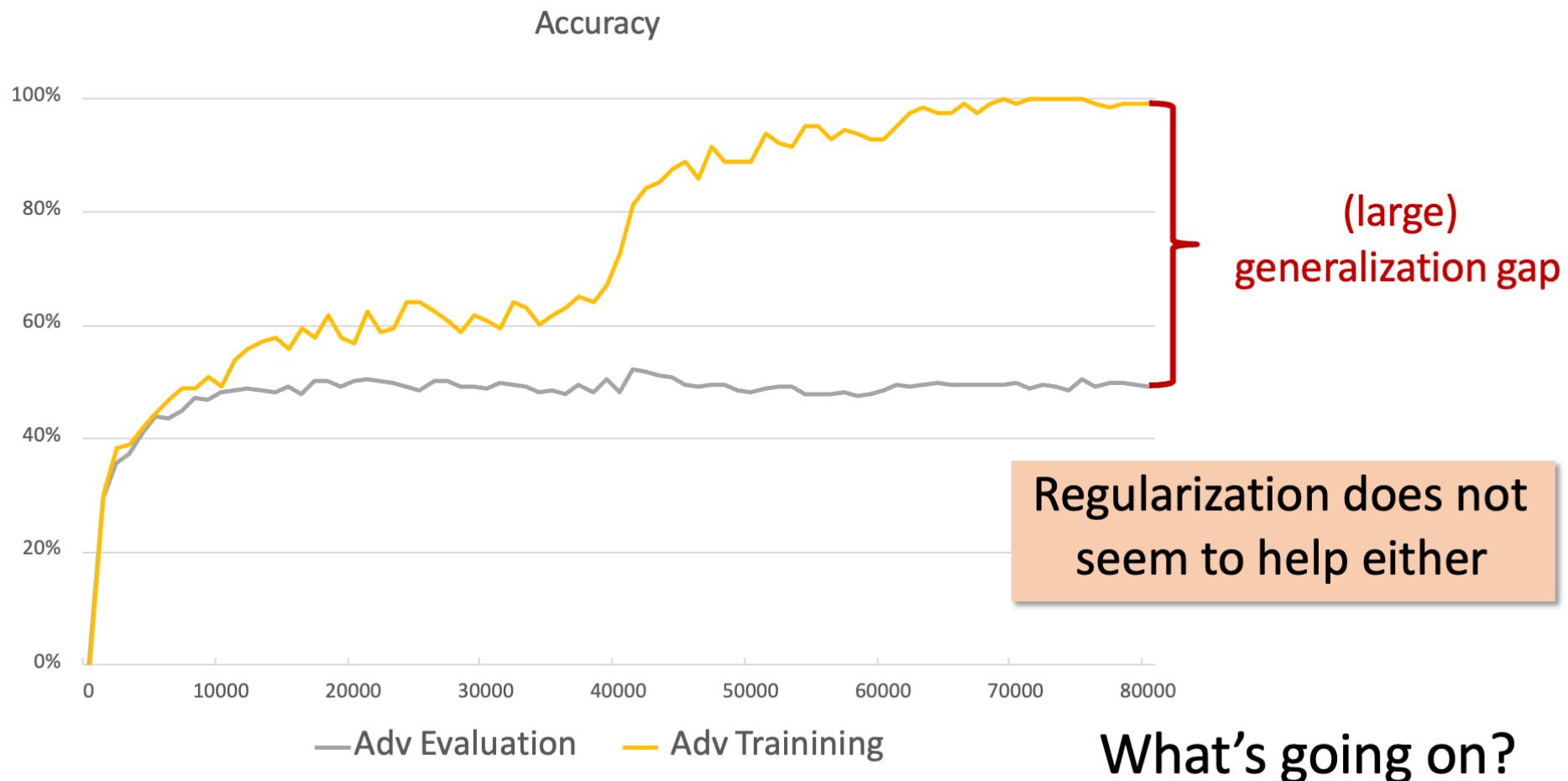
# Do Robust Deep Networks Overfit?



# Do Robust Deep Networks Overfit?

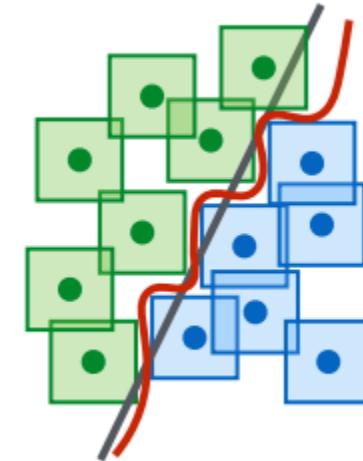
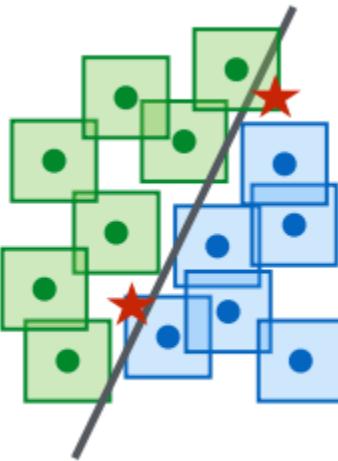
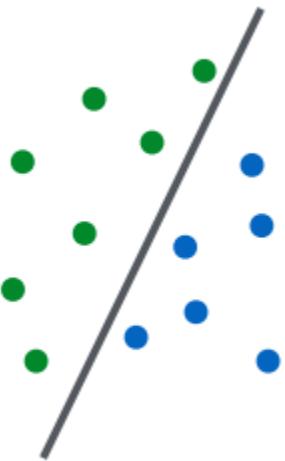


# Do Robust Deep Networks Overfit?



# Adversarial robustness

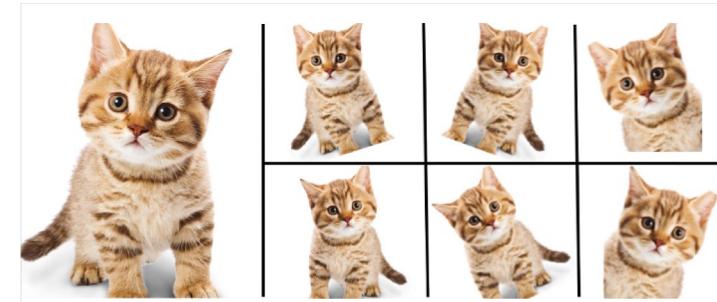
- More capacity is required



- More training data might be required
- Optimization during training may be more difficult

# Does Being Robust Help “Standard” Generalization?

**Data augmentation:** An effective technique to improve “standard” generalization



Adversarial training

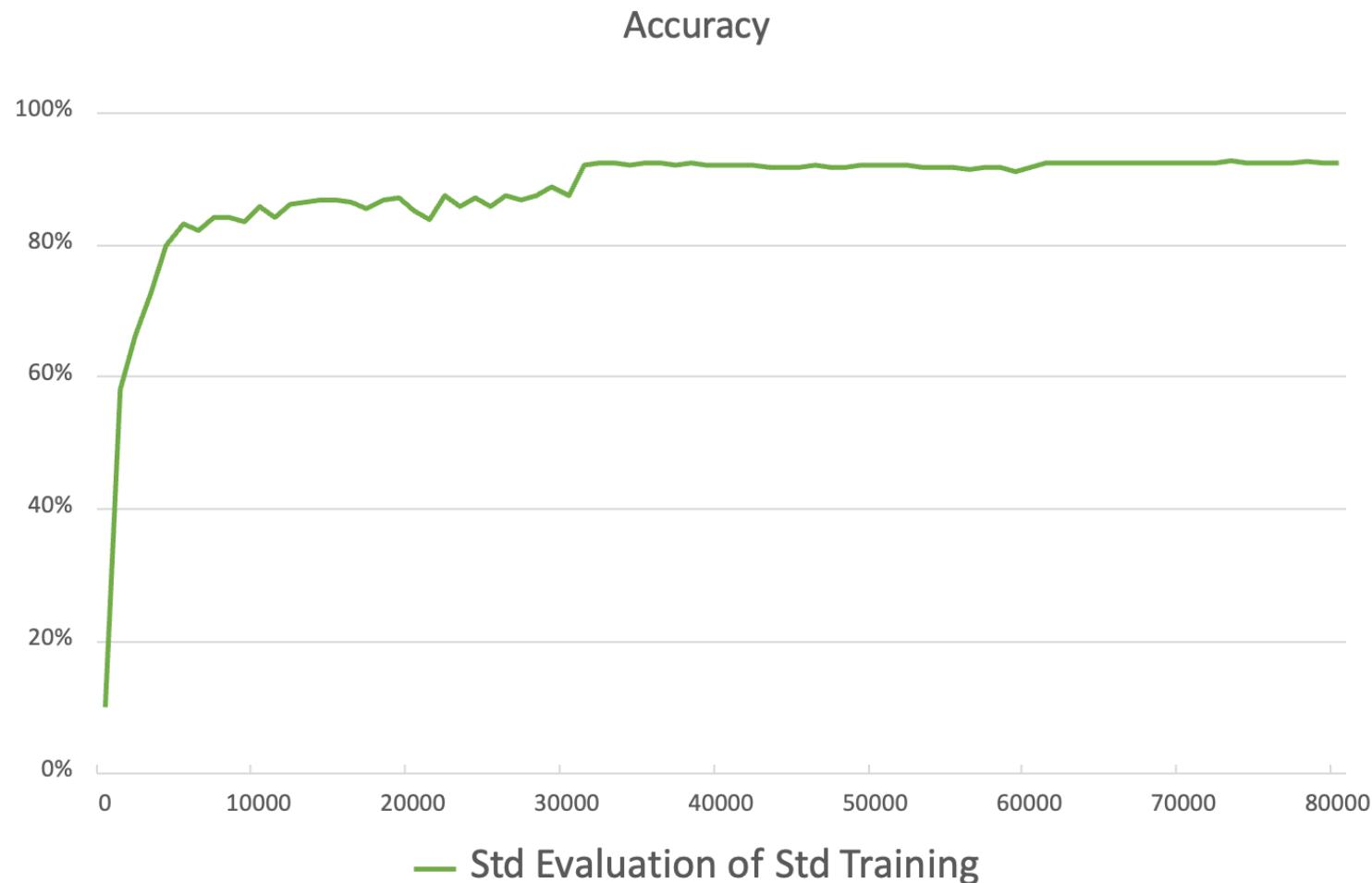
=

An “ultimate” version of data augmentation?

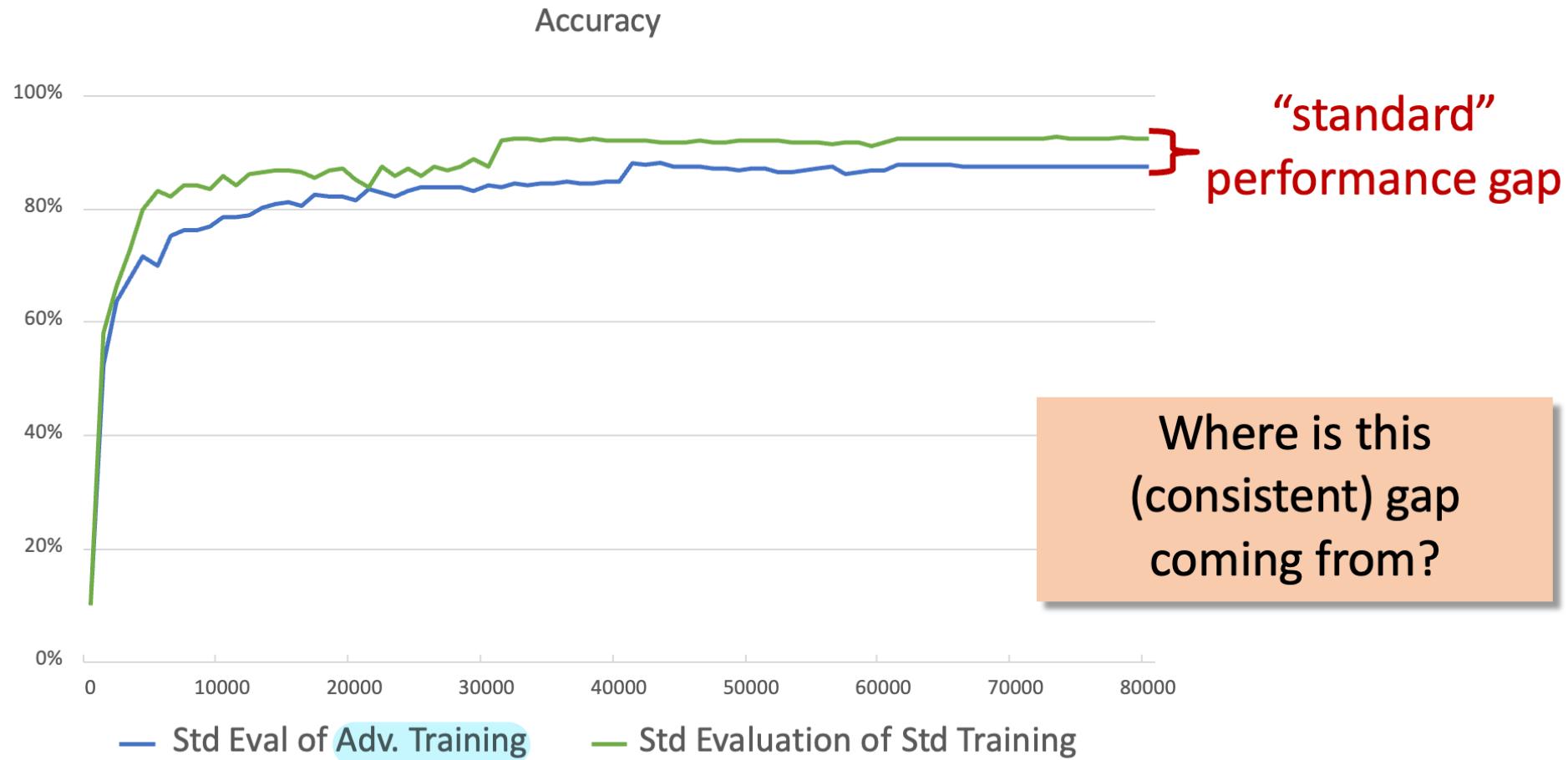
(since we train on the “most confusing” version of the training set)

Does adversarial training always improve  
“standard” generalization?

# Does Being Robust Help “Standard” Generalization?



# Does Being Robust Help “Standard” Generalization?

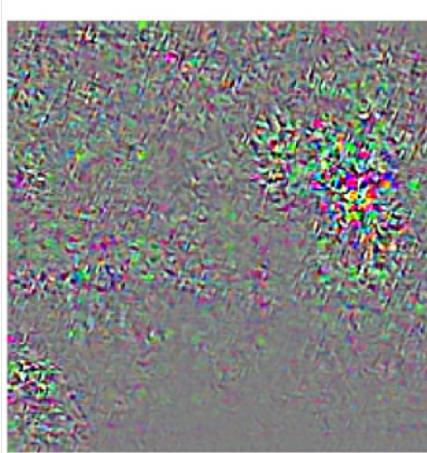


# But There Are (Unexpected?) Benefits Too

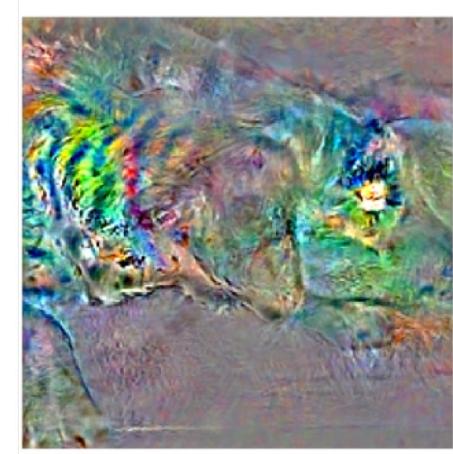
Models become more **semantically meaningful**



Input



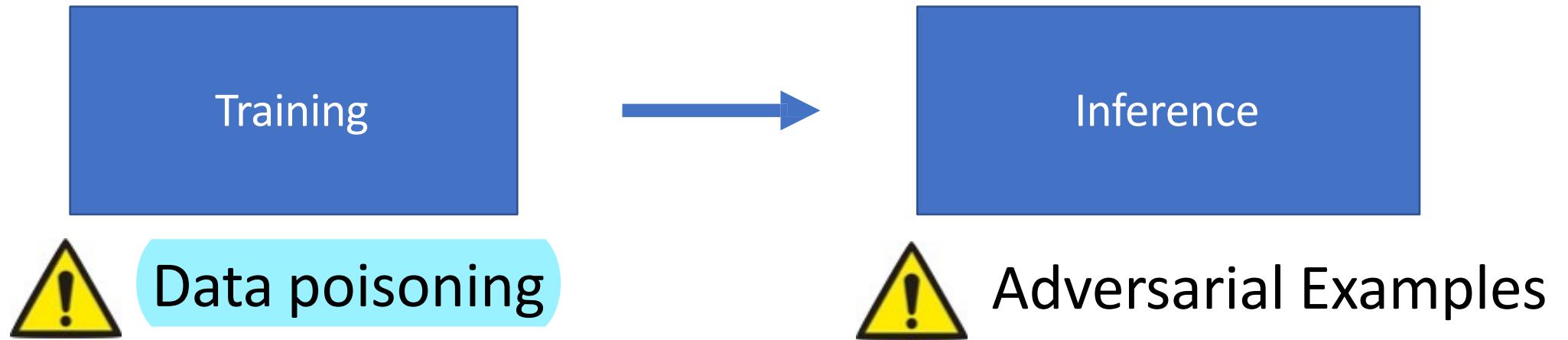
Gradient of  
standard model



Gradient of  
adv. robust model

What are other types of attacks?

# Training time attack

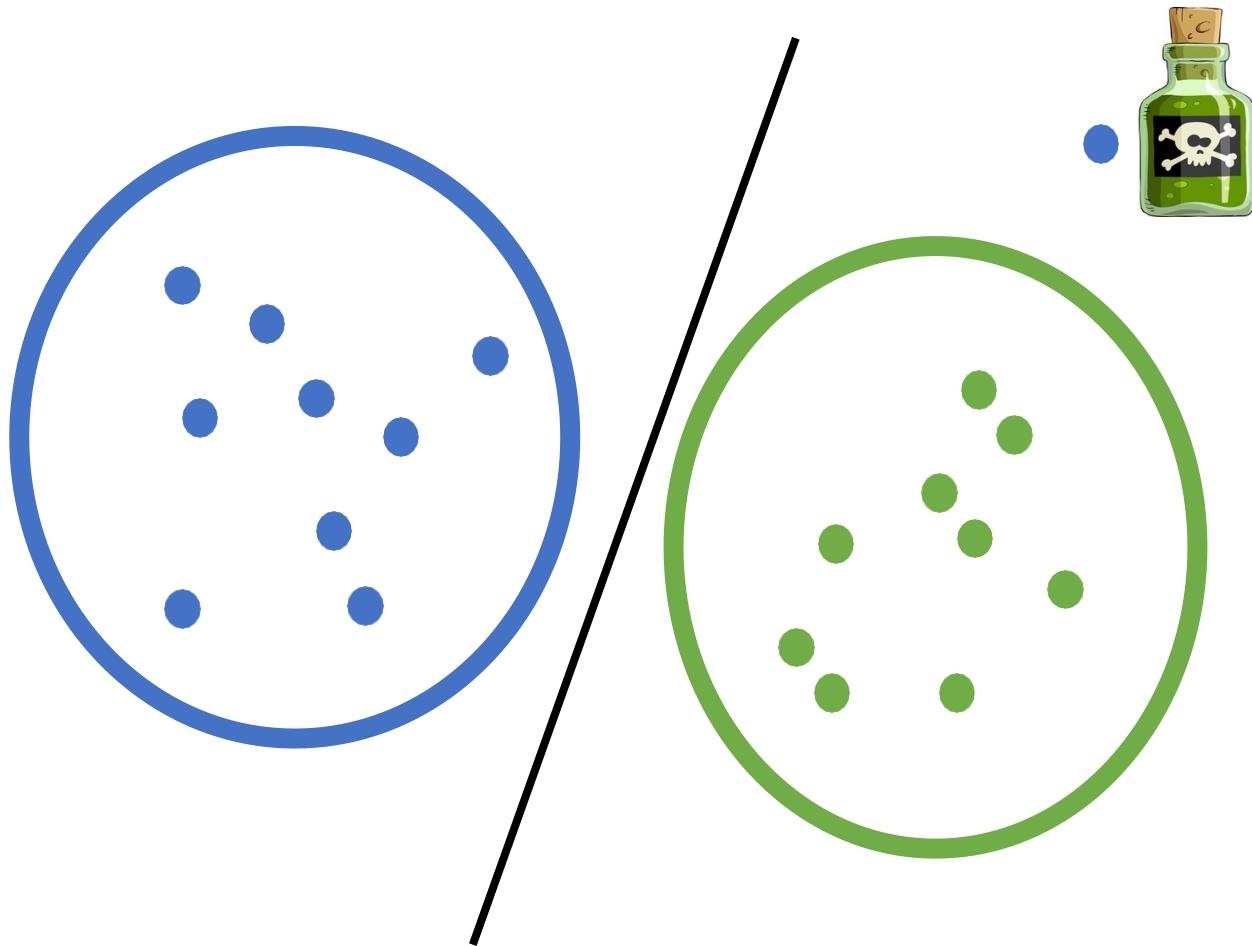


**(Deep) ML is “data hungry”**  
→ Can’t afford to be too picky about where we get the training data from

What can go wrong?

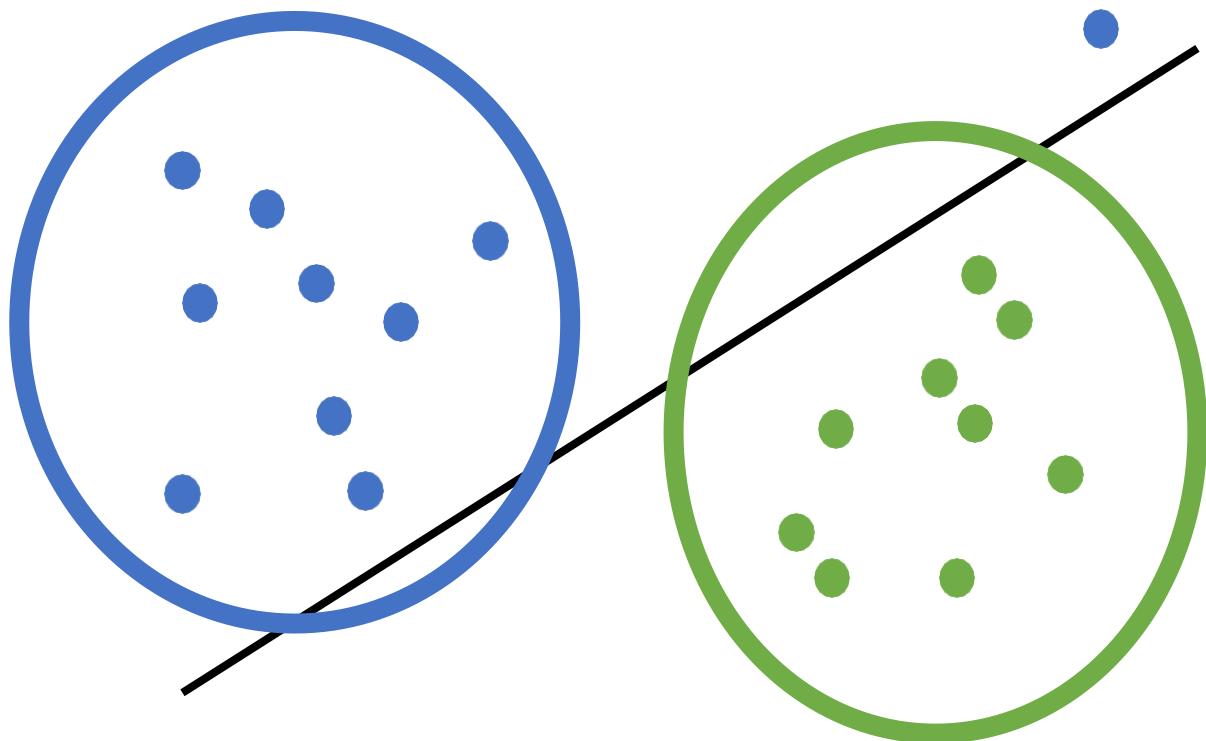
# Data Poisoning

**Goal:** Maintain training accuracy but hamper generalization



# Data Poisoning

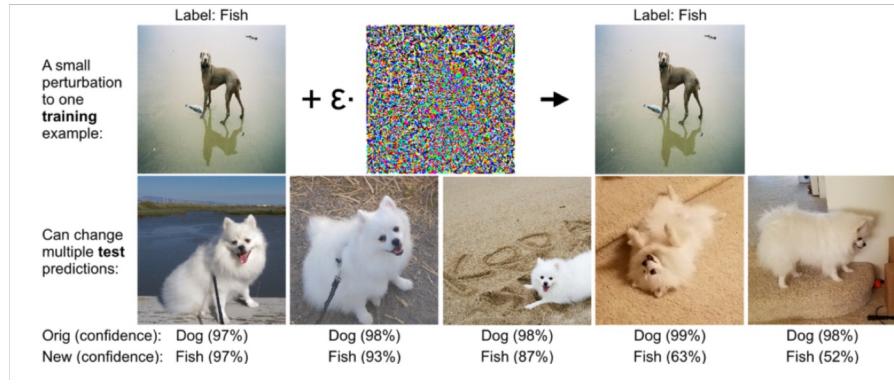
**Goal:** Maintain training accuracy but hamper generalization



- Fundamental problem in “classic” ML (robust statistics)
- **But:** seems less so in deep learning
- **Reason:** Memorization?

# Data Poisoning

**Goal:** Maintain training accuracy but hamper generalization



“van”



“dog”

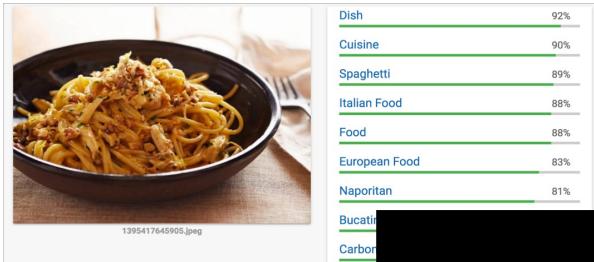
[Koh Liang 2017]: Can manipulate many predictions with a single “poisoned” input

[Gu Dolan-Gavitt Garg 2017][Turner Tsipras M 2018]:  
Can plant an undetectable backdoor that gives an almost total control over the model

**But:** This gets (much) worse

# Model steeling

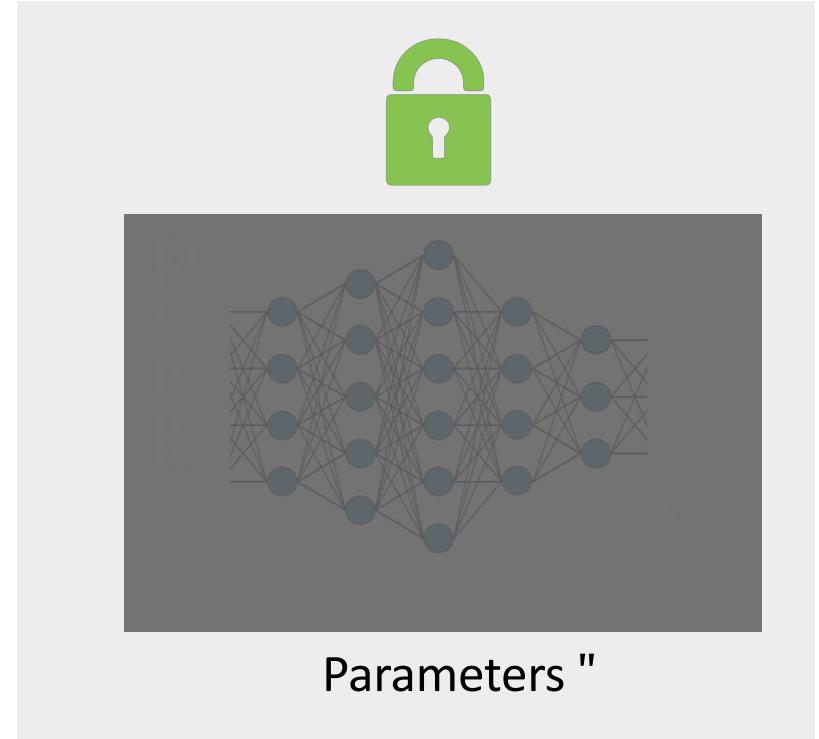
## Google Cloud Vision API



## Microsoft Azure (Language Services)

A screenshot of the Microsoft Azure Language Services API documentation page. It features four main sections:

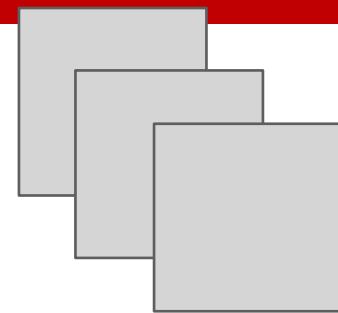
- Language Understanding (LUIS)**: A card with a description of teaching apps to understand commands, a try button, and a link to use with an Azure subscription.
- Text Analytics API**: A card with a description of evaluating sentiment and topics, a try button, and a link to use with an Azure subscription.
- Bing Spell Check API**: A card with a description of detecting and correcting spelling mistakes, a try button, and a link to use with an Azure subscription.
- Translator Text API**: A card with a description of conducting machine translation, a try button, and a link to use with an Azure subscription.



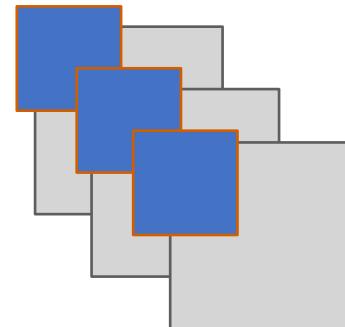
# Model steeling

Does limited access  
give security?

In short: No



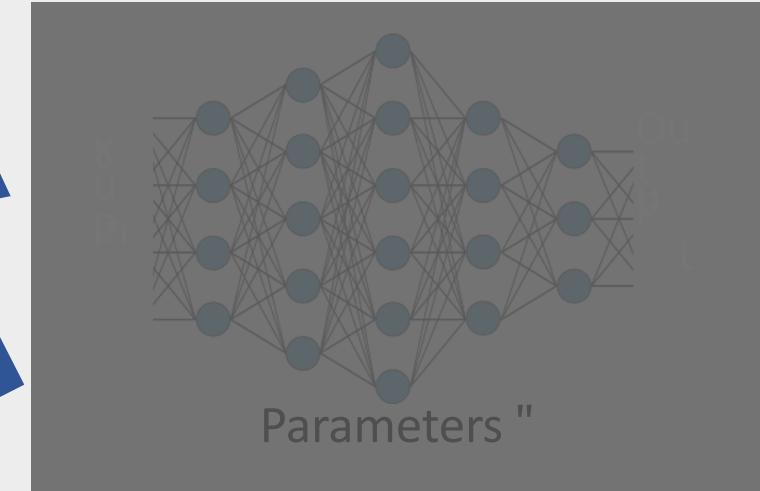
Data



Predictions



Black box attacks



# Black-box attacks

Does limited access give security?

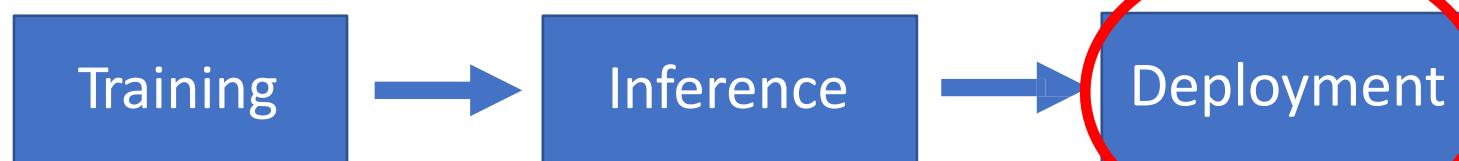
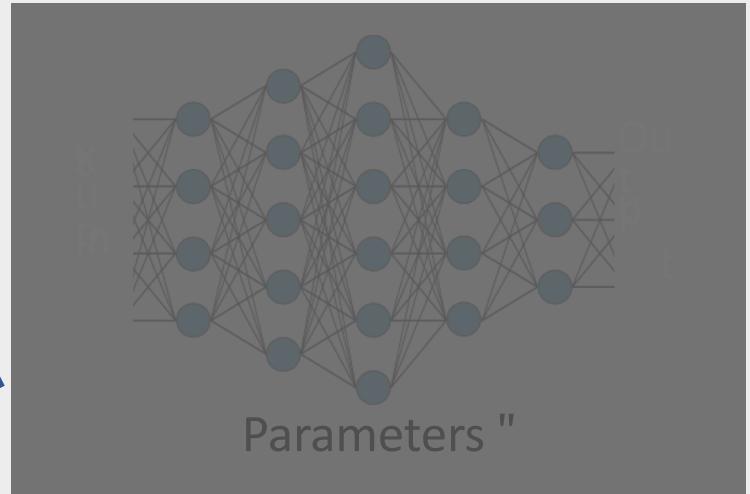
**Model stealing:** “Reverse engineer” the model

[Tramer et al., 2016]

**Black box attacks:** Construct adv. examples from queries

Data

Predictions



Black box attacks

# Three commandments of Secure/Safe ML

- shall not train on data you don't fully trust  
(because of data poisoning)
- shall not let anyone use your model (or observe its outputs) unless you completely trust them  
(because of model stealing and black box attacks)
- shall not fully trust the predictions of your model  
(because of adversarial examples)

# Summary

- Machine learning models can be very brittle
  - A small change in the input can produce a very different model output
- Adversarial examples can be produced systematically by well-known attacks like FGSM, PGD, and C&W
- There are systematic methods to improve robustness
- There are limitations in the current state of the art/science
  - More robust models may be less accurate, depending on circumstances
  - Intuitive measures of “small change in input” are tricky - no clear metric for human perception
- Other types of threats for ML methods