# A THEORETICAL LOOKS AT ADVERSARIAL EXAMPLES

## Tom Goldstein

…and also…

**Ali Shafahi, Ronny Huang,
Mahyar Najibi, Octavian Suciu,
Christoph Studer, Soheil Feizi, Tudor Dumitras**

UNIVERSITY OF MARYLAND

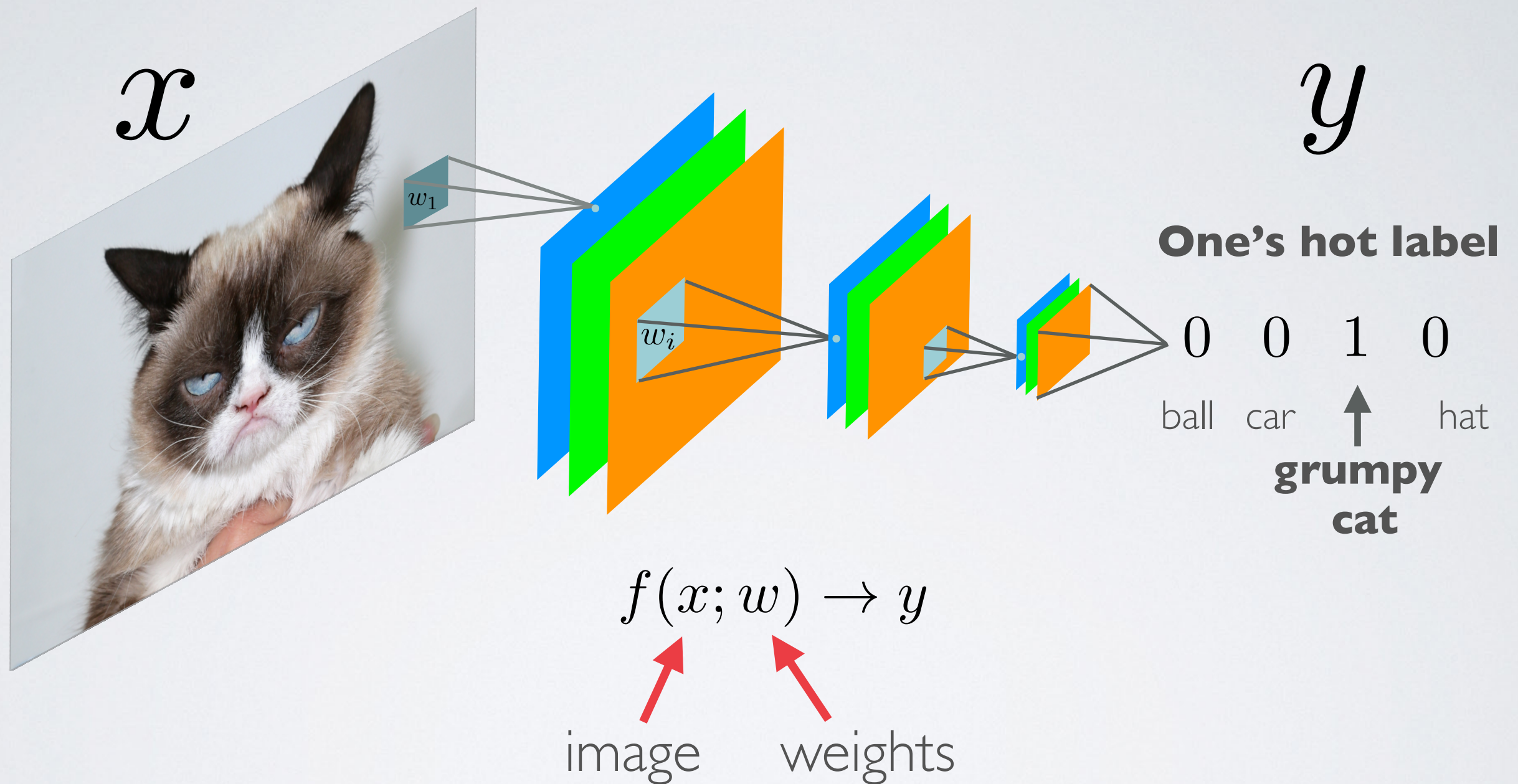# OVERVIEW

**Why is optimization so easy on neural nets?**

**What are adversarial examples,
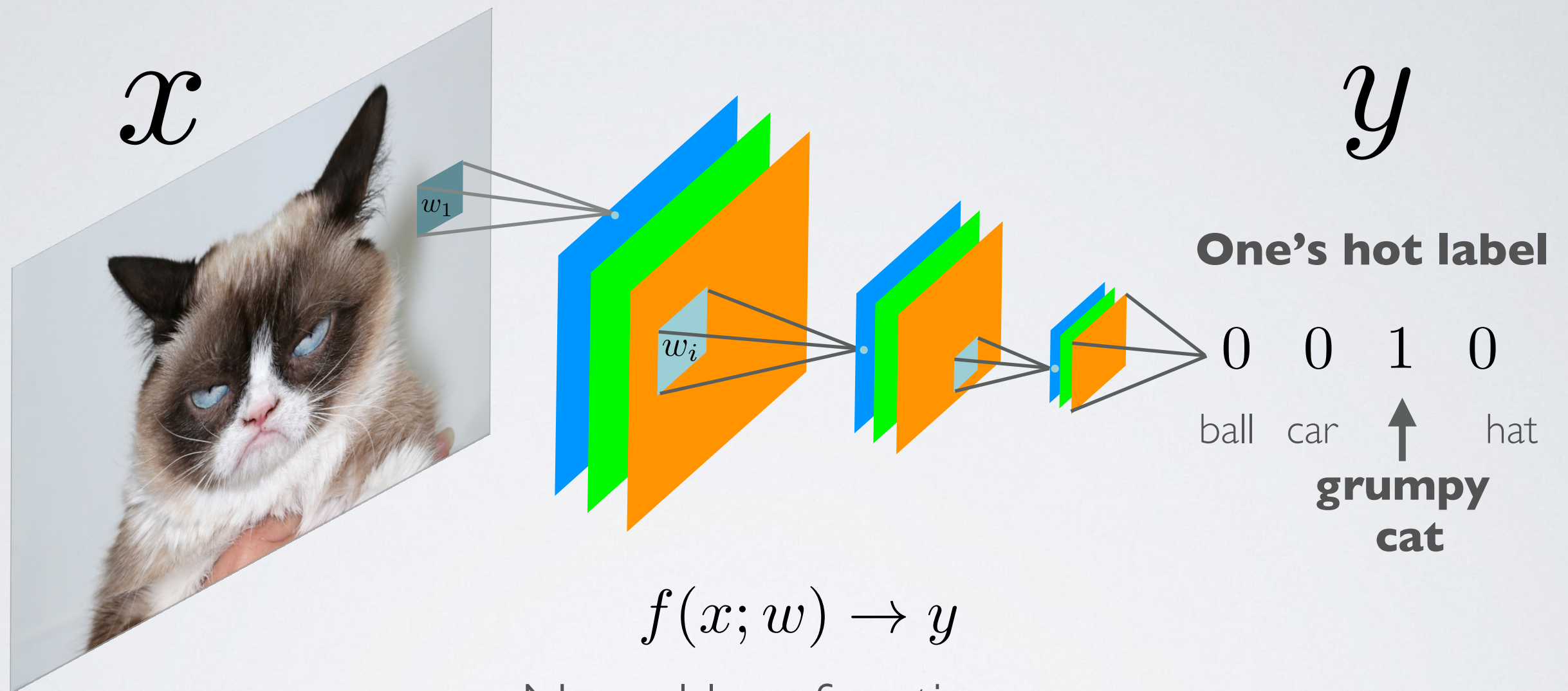and what are their risks?**

**Poison attacks**

**Are they an escapable problem?**

# CONVOLUTIONAL NET

$x$

$y$

**One's hot label**

$$0 \quad 0 \quad 1 \quad 0$$

ball    car    $\uparrow$    hat

**grumpy cat**

$w_1$

$w_i$

$$f(x; w) \to y$$

image     weights

# CONVOLUTIONAL NET

$x$

$y$



**One's hot label**

$0 \quad 0 \quad 1 \quad 0$

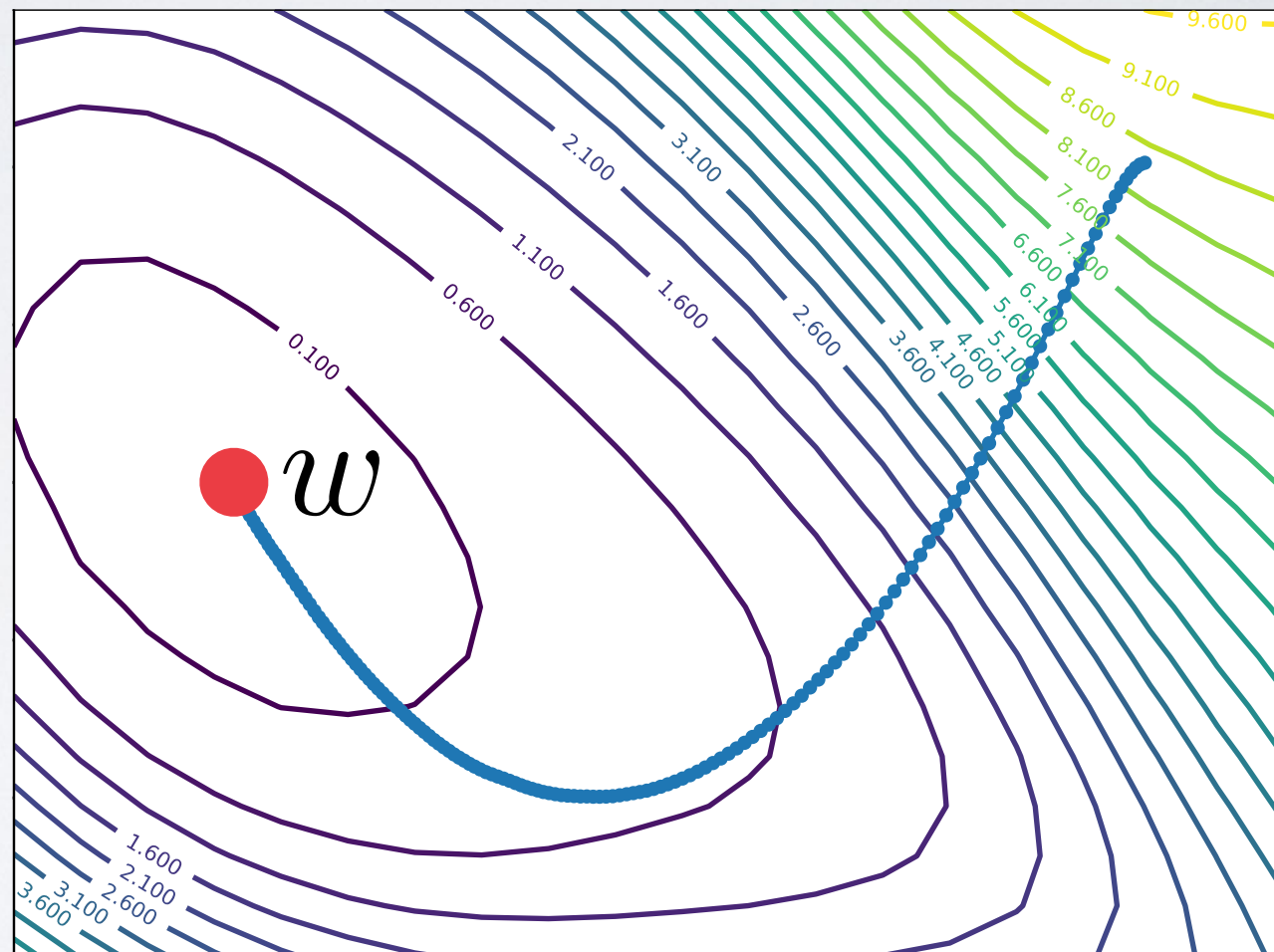ball   car   ↑   hat

**grumpy cat**

$$f(x; w) \rightarrow y$$

Neural loss function

$$L(w) = \min_w \sum_i \| f(x_i; w) - y_i \|^2$$

**Non-convex?**

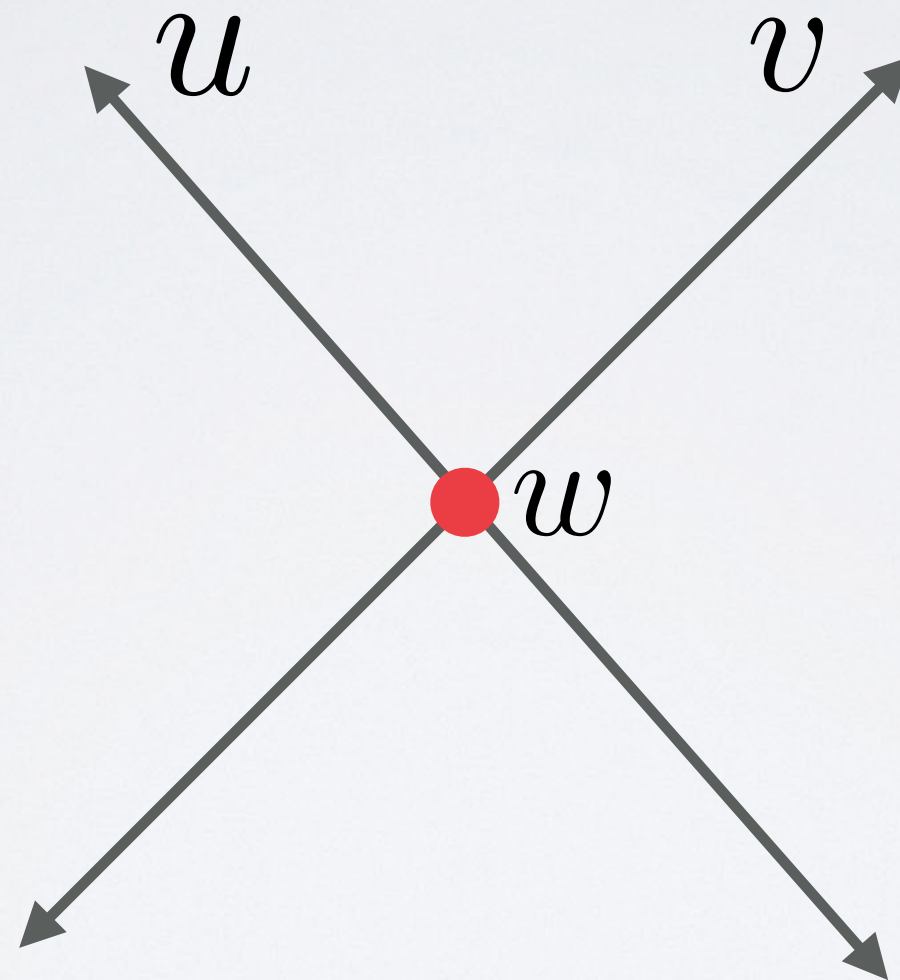# VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION

**Step 1:**

**Find minimizer**



**30 million dimensions**

# VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION
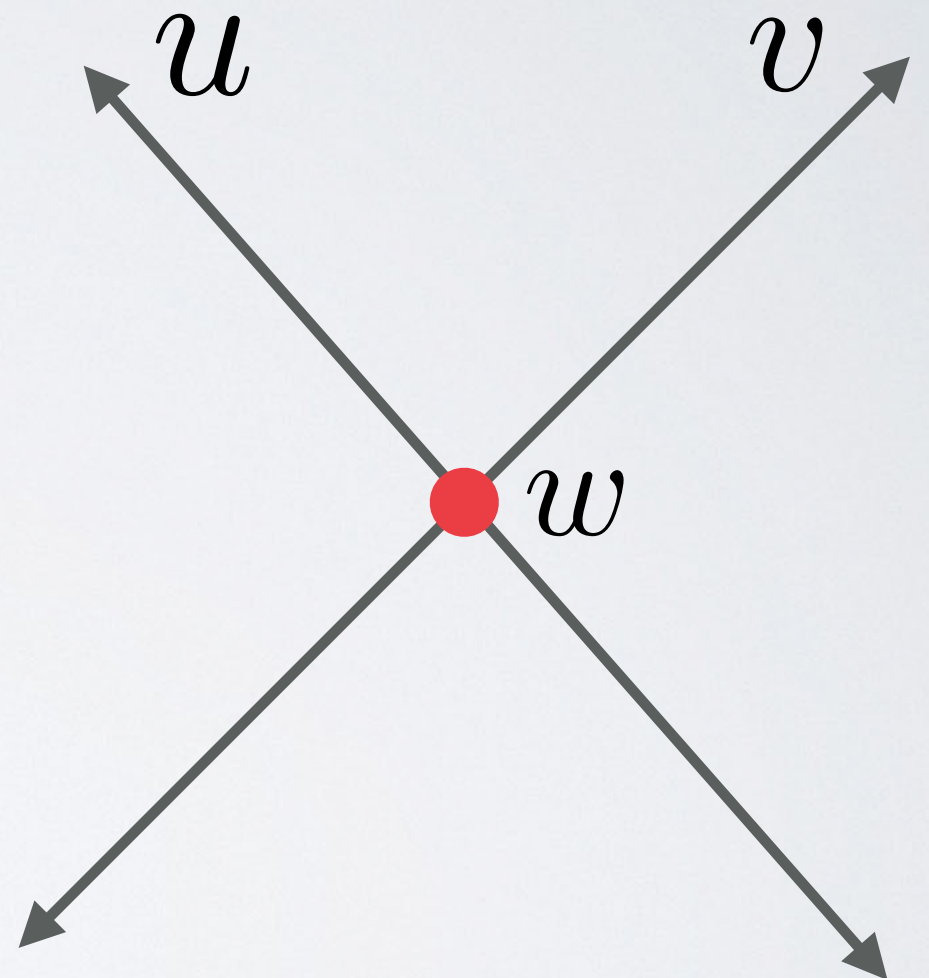
**Step 2:**

**Random directions**

$u, v$

$u$

$v$

$w$

# VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION
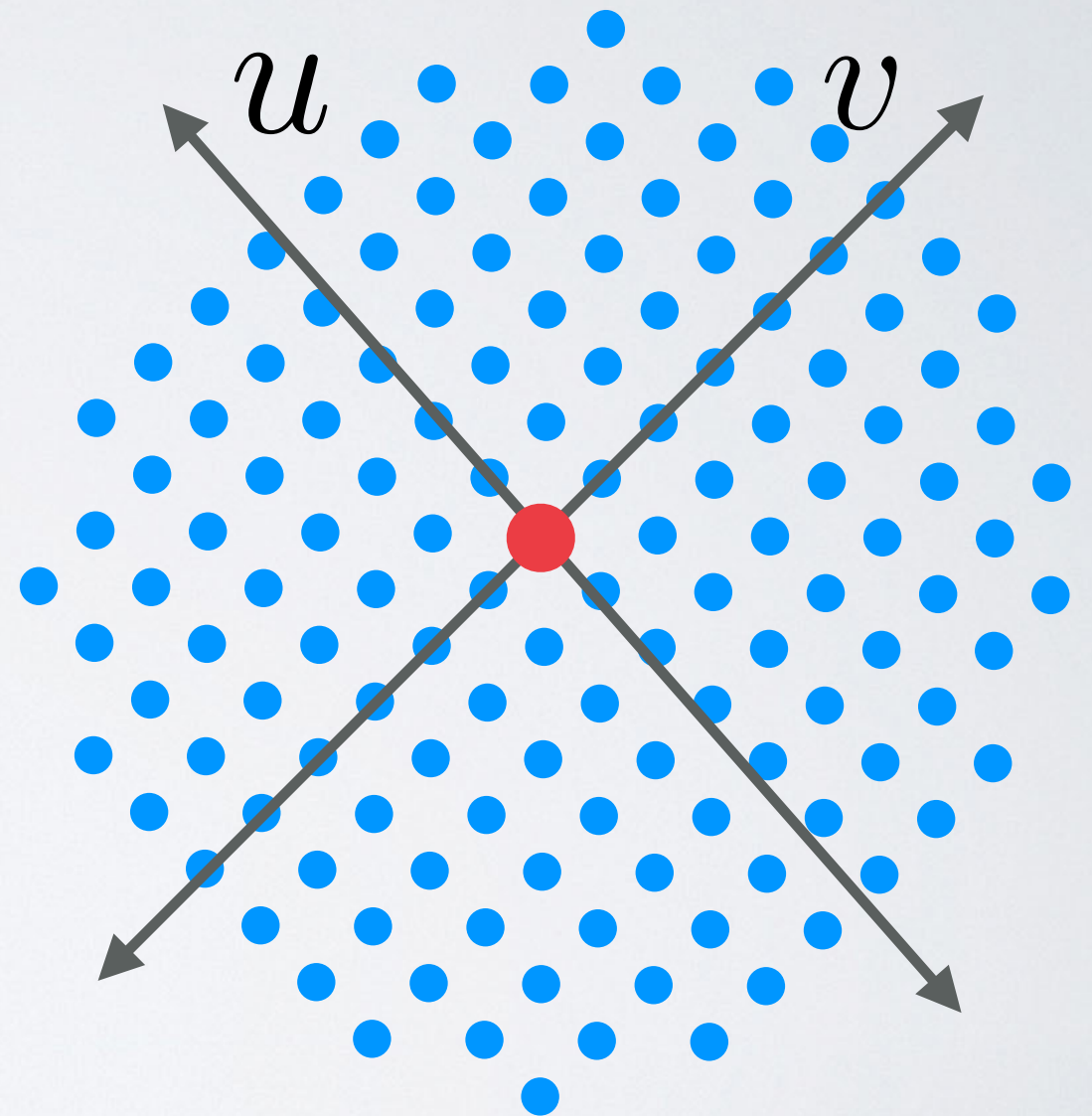
**Step 3:**

**Filter normalization**

$$u_i \leftarrow u_i \cdot \frac{\|w_i\|}{\|u_i\|}$$

$u$ $v$

$w$

**Li, Xu, Taylor, Studer, G. "Visualizing the loss landscape of neural nets."**

# VISUALIZING LOSS FUNCTIONS: FILTER NORMALIZATION

**Step 4**

Plot

$u$ $v$

**Li, Xu, Taylor, Studer, G. "Visualizing the loss landscape of neural nets."**

# 56 LAYER "VGG" NET

## CIFAR-10

# 56 LAYER NEURAL NET

## CIFAR-10

# 56 LAYER NEURAL NET

## CIFAR-10

VGG-like

ResNet

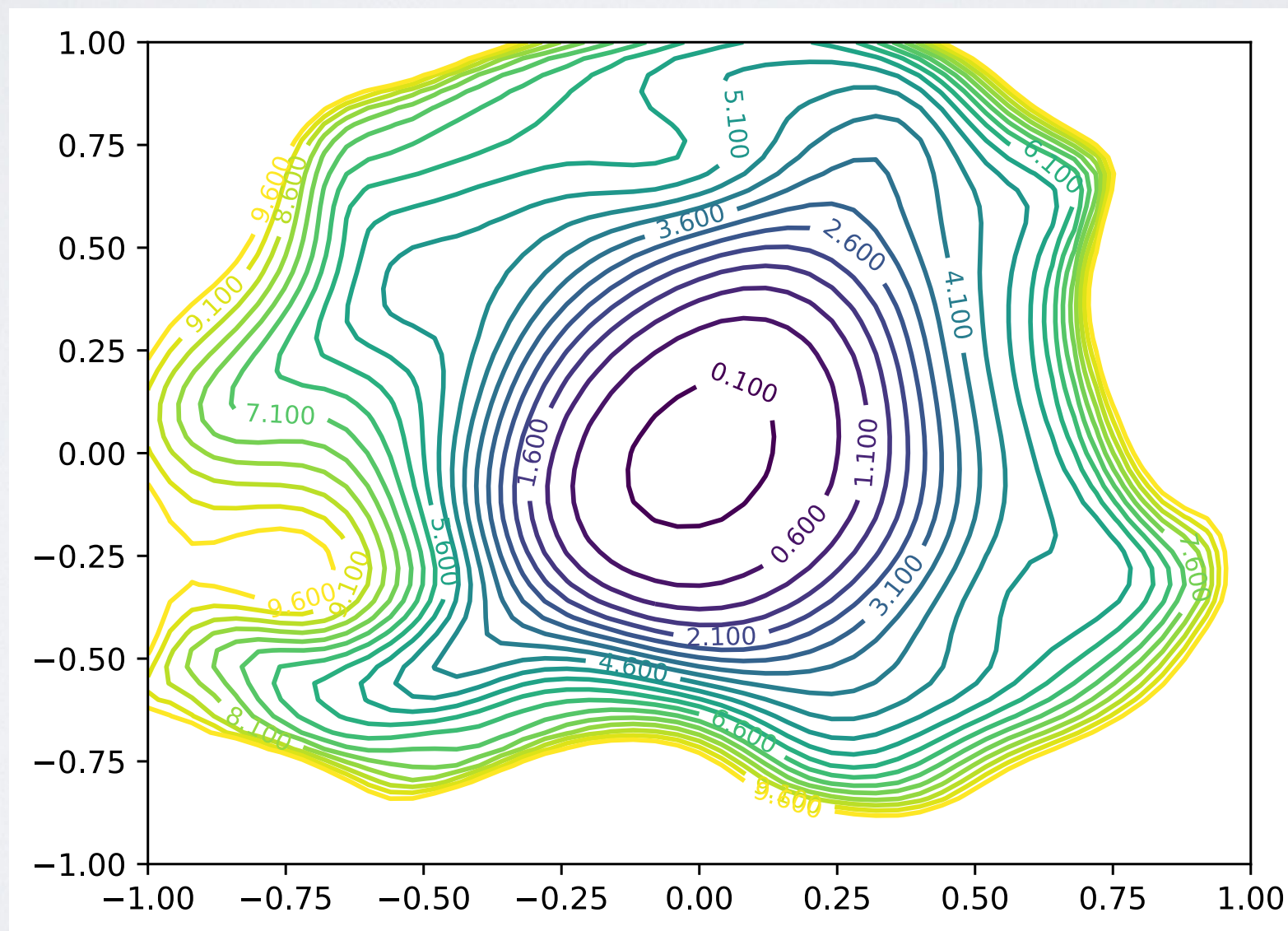skip connections

# CHAOTIC TRANSITIONS

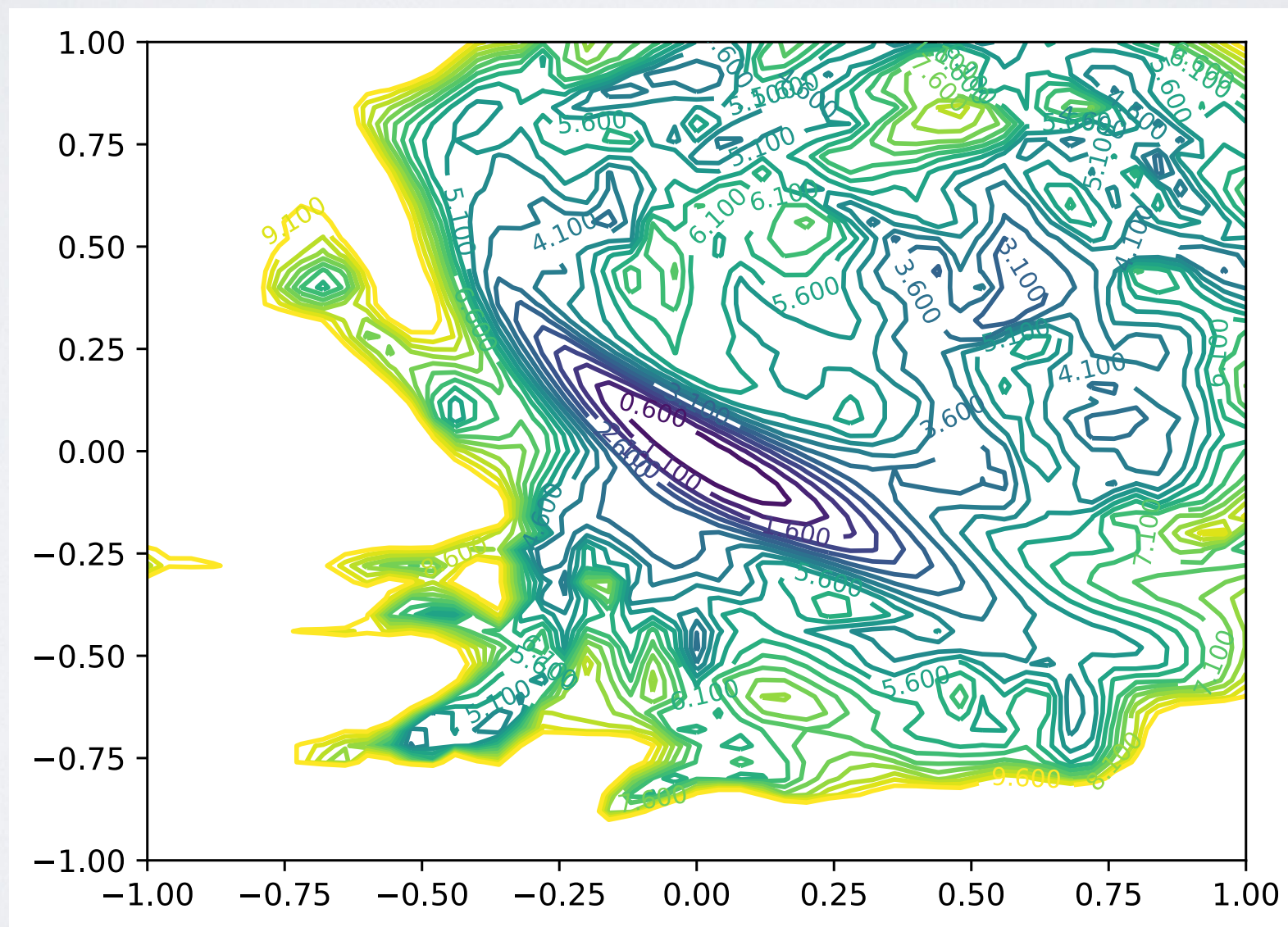## VGG-9

# CHAOTIC TRANSITIONS

## VGG-20

# CHAOTIC TRANSITIONS

## VGG-56

# CHAOTIC TRANSITIONS

## VGG-110

# CHAOTIC TRANSITIONS

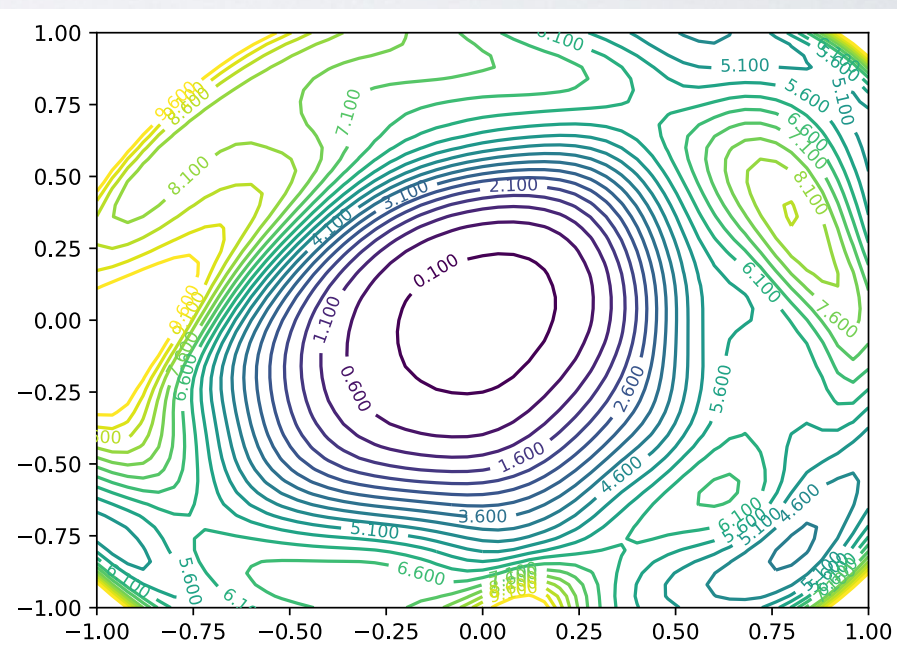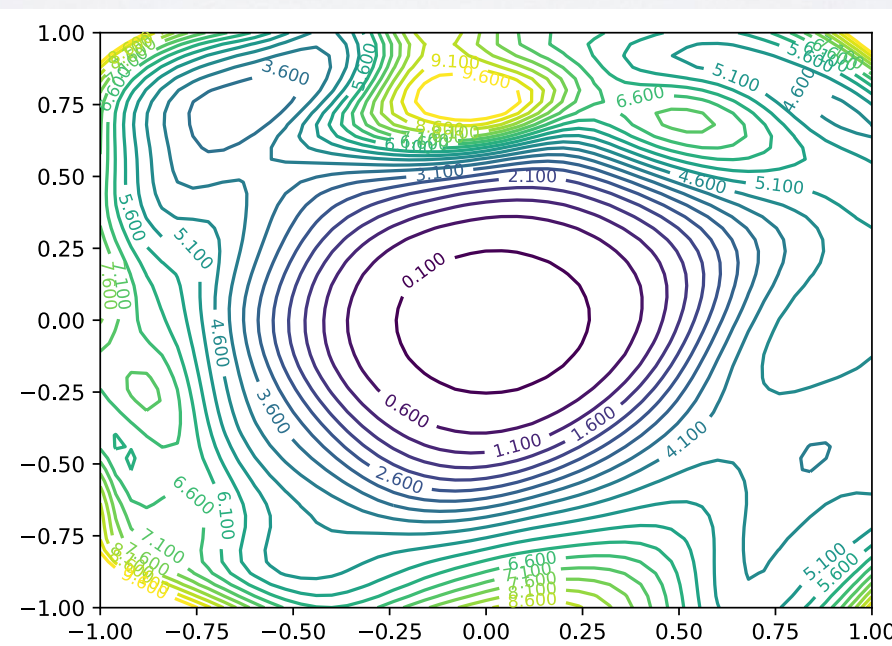**VGG-20**    **VGG-56**    **VGG-110**
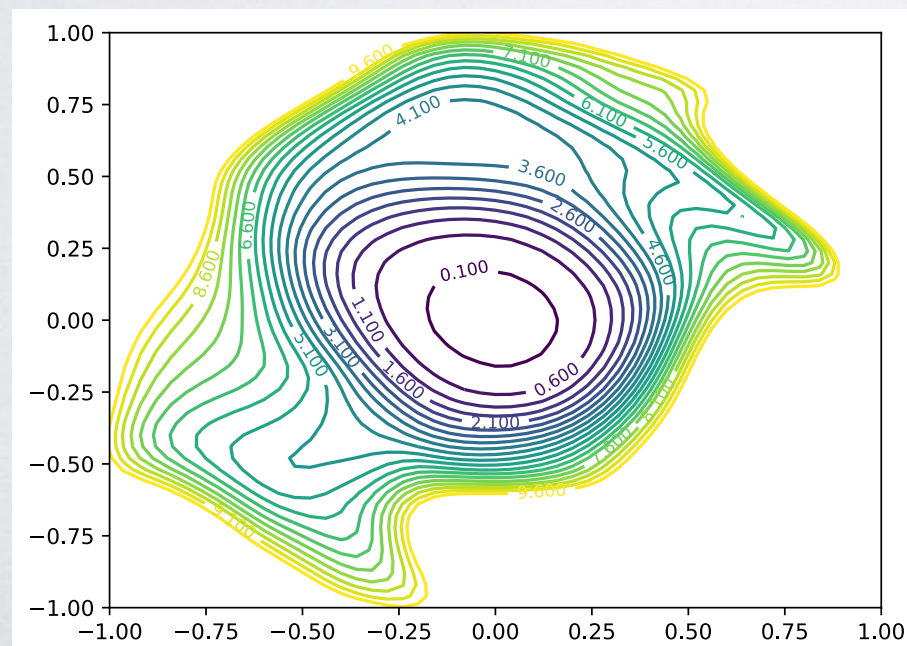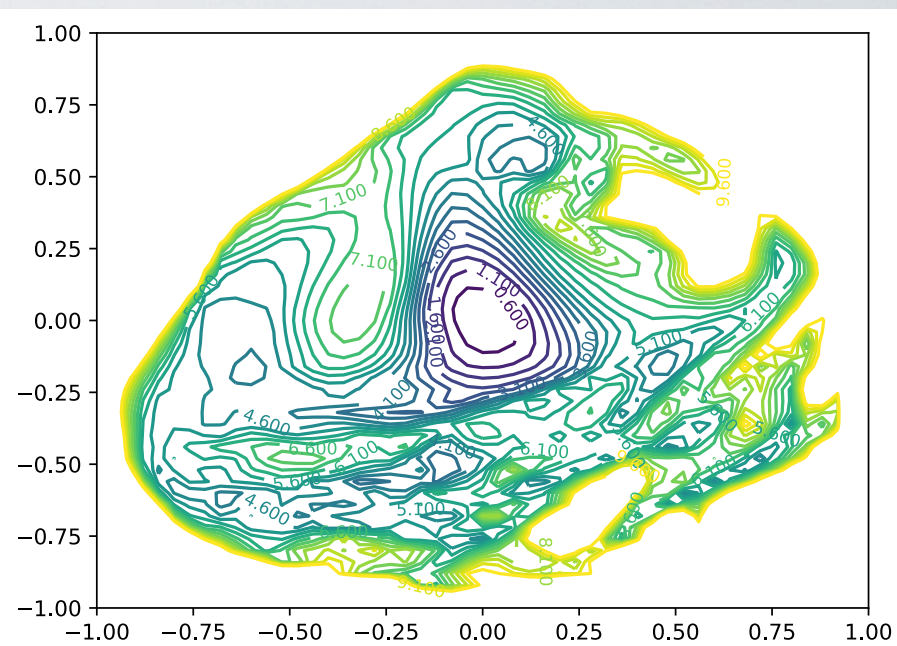


**Convexity**    **Chaos**

# CHAOTIC TRANSITIONS

**VGG-20**  **VGG-56**  **VGG-110**

**ResNet-20**  **ResNet-56**  **ResNet-110**

Optimization on neural nets is easy!

That's great for ML.

…but bad for security.

# ADVERSARIAL EXAMPLES

$x$

$y$

**One's hot label**

0   0   1   0

ball   car   ↑   hat

**grumpy cat**

add "tweaks" to image

compute "tweak" to each pixels

# ADVERSARIAL ATTACKS

"Egyptian Cat" 28%

"Traffic Light" 97%

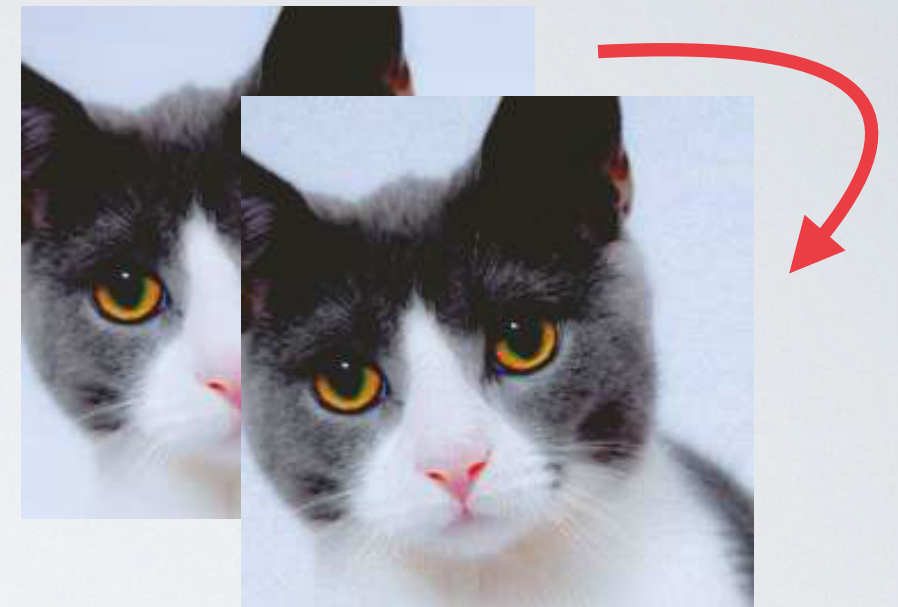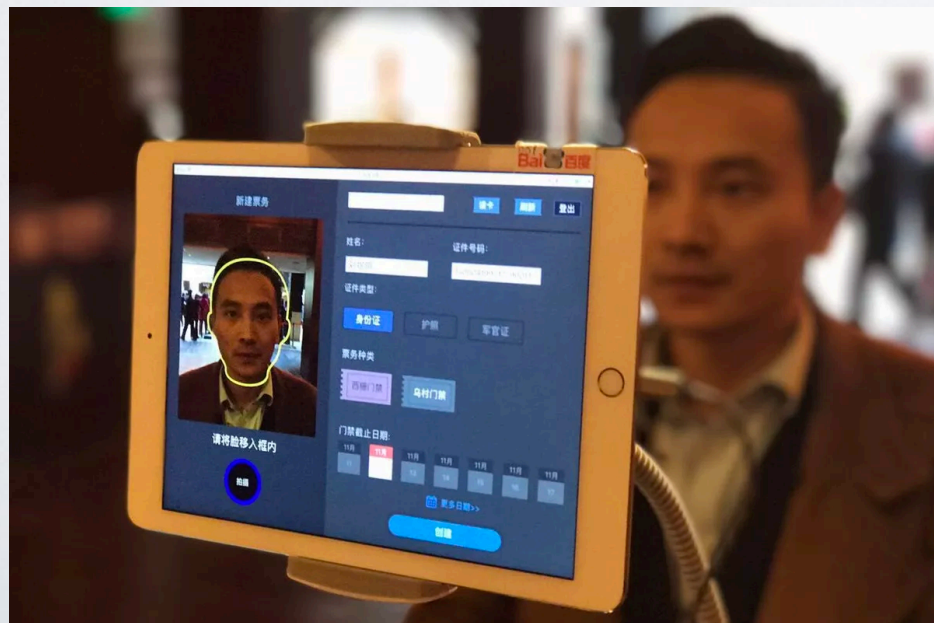# ADVERSARIAL ATTACKS

"Ox" 85%

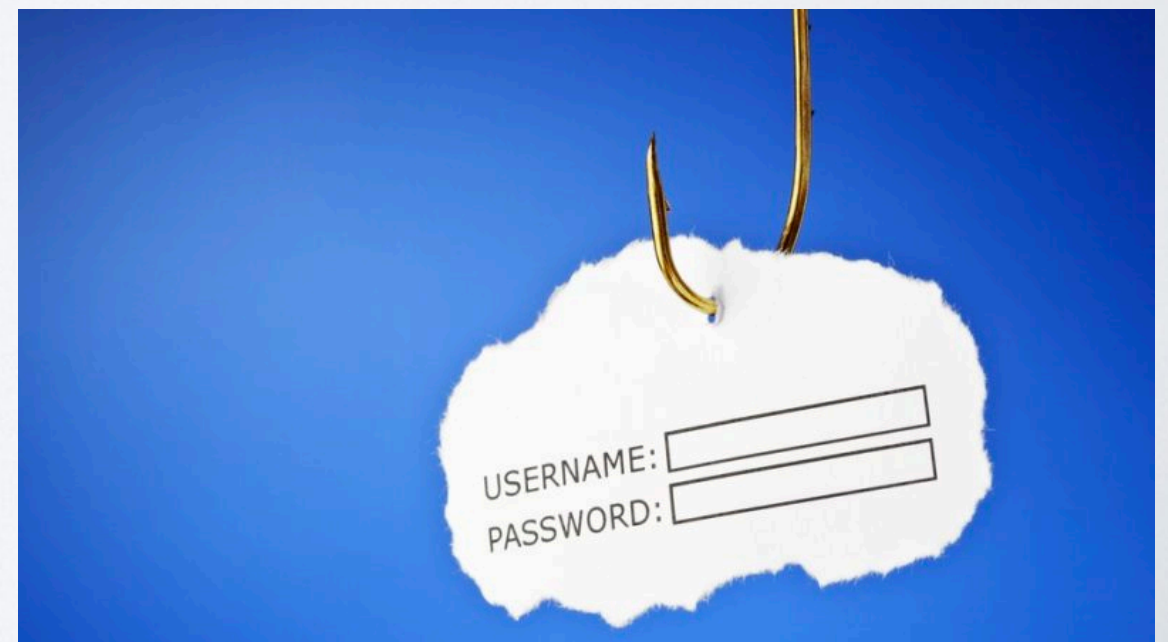"Traffic Light" 96%

# THREAT MODEL: EVASION

**Test-time attacks: adversary controls inputs**

**Fails when...**

Supervised security desk

Phishing email/
Competitor email

# THREAT MODEL: POISON

**Train-time attacks:
adversary controls training data**
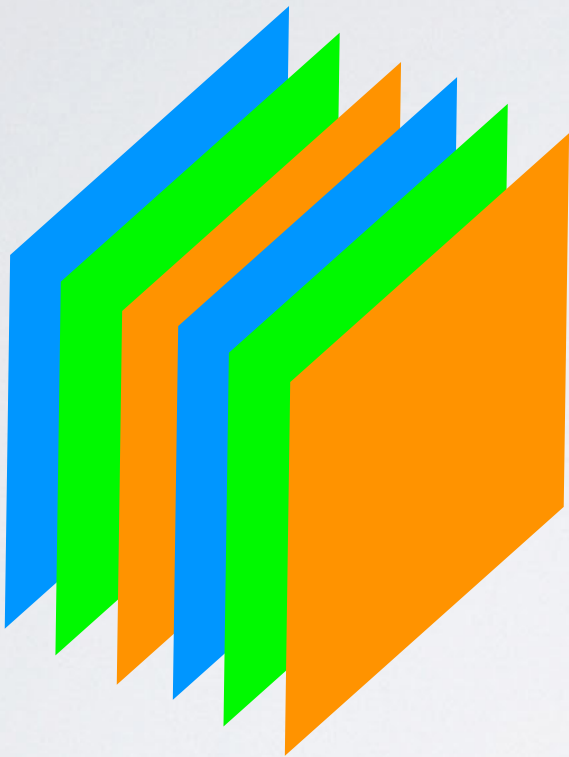
**Does this *actually* happen?**

Scraping images from the web

Harvesting system inputs (spam detector)

Bad actors/inside agents
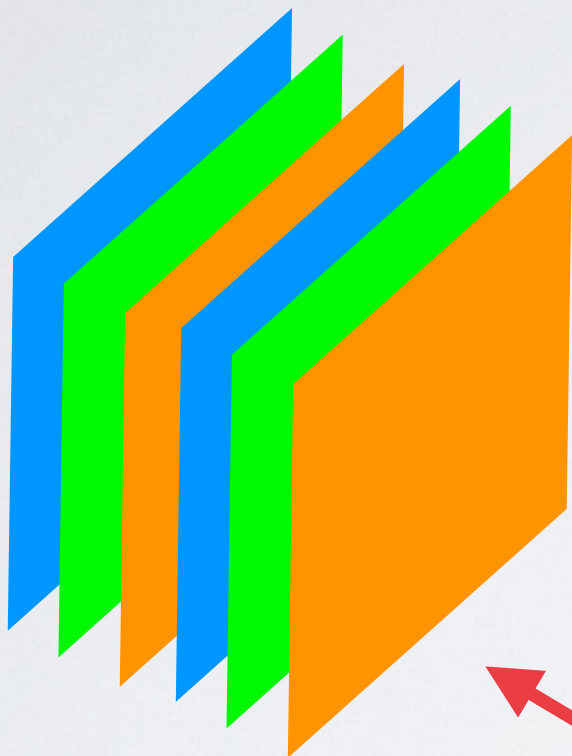
# HOW POISONING WORKS
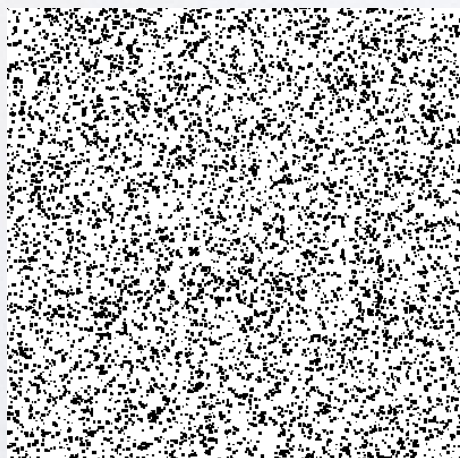
Training data

Testing example

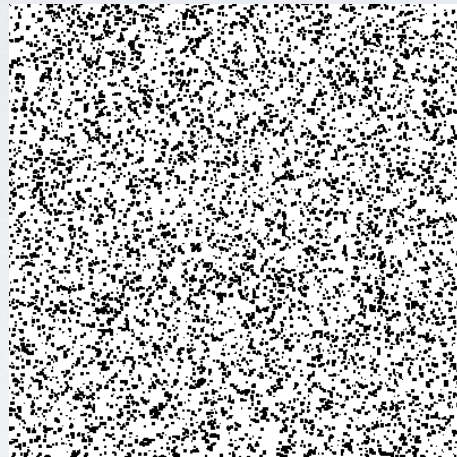Plane          Frog

Base     +     =     Poison!

# CLEAN-LABEL + TARGETED

**Base**        **Poison!**



**Attacks are hard to detect**

Clean label: poisons are labeled "correctly"

Performance only changes on selected target
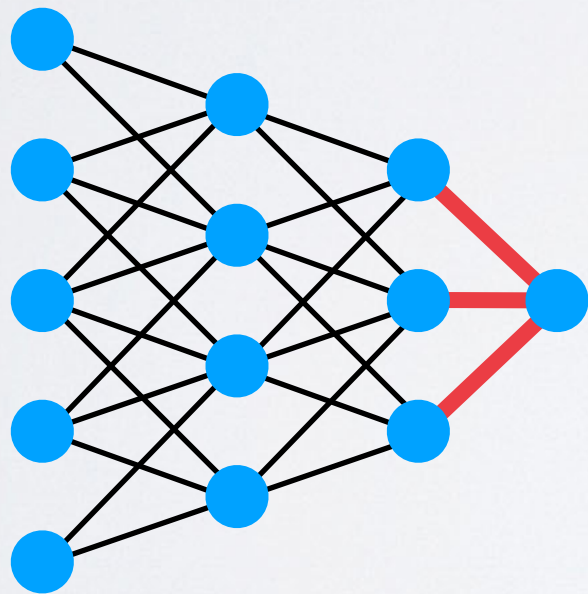
**Attacks can be executed by outsider**

Poison data can be placed on the web

Poison data can be sent/emailed to data collectors
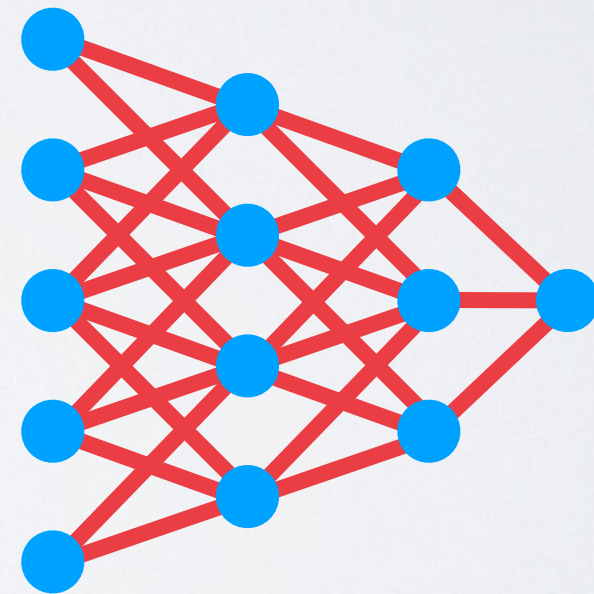
# TWO CONTEXTS

## Transfer learning

- Standard, pre-trained net is used
- "Feature extraction" layers frozen
- Classification layers re-trained
- Common practice in industry

## End-to end re-training

- Pre-trained net is used
- All-layers are re-trained



"One-shot kill" possible

Multiple poisons required

# COLLISION ATTACK

$$\mathbf{p} = \underset{\forall \mathbf{x}}{\arg\min} \; \|f(\mathbf{x}) - f(\mathbf{t})\|^2 + \beta \|\mathbf{x} - \mathbf{b}\|^2 \qquad (1)$$
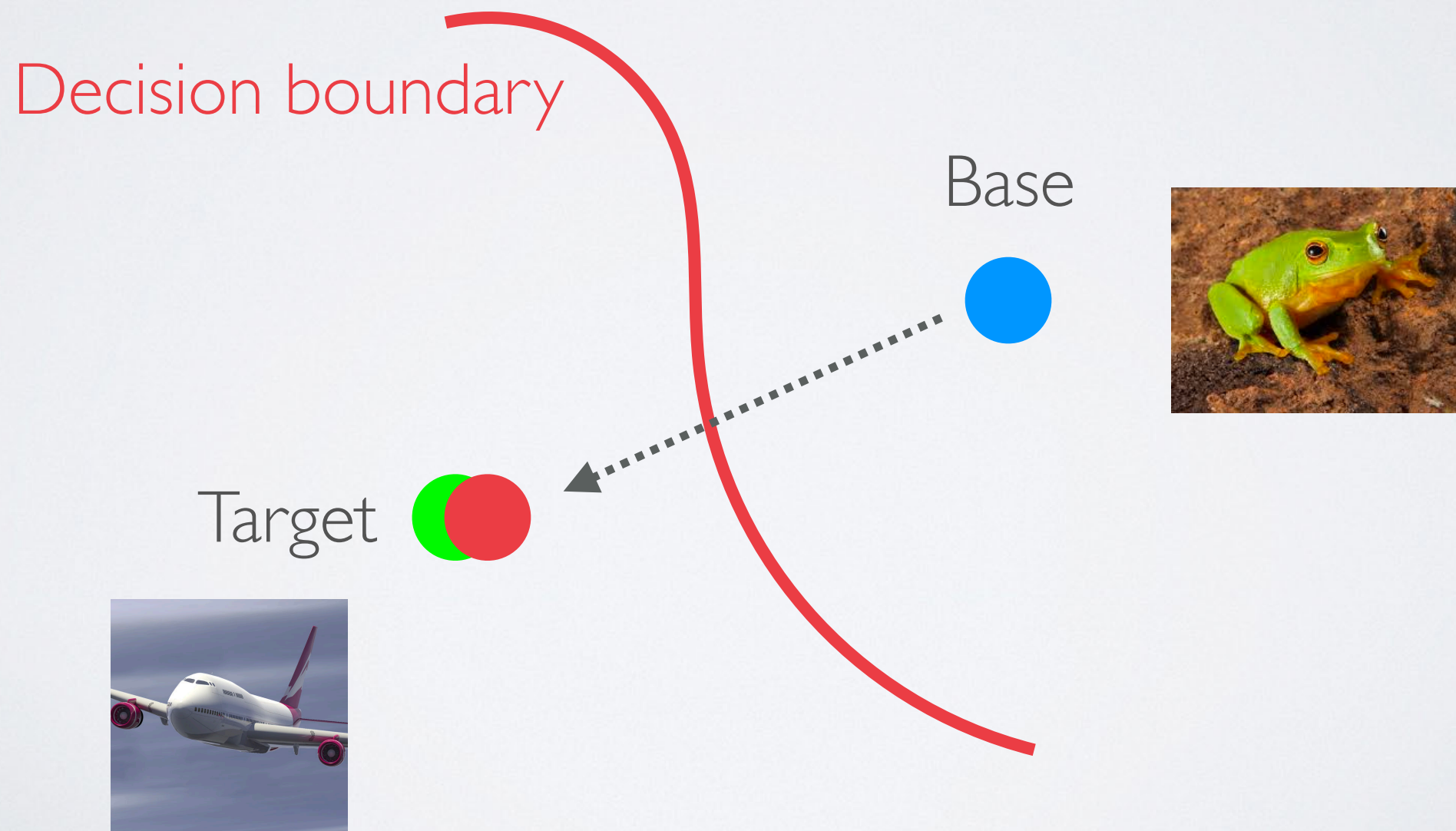
Decision boundary

Base

Target

# COLLISION ATTACK
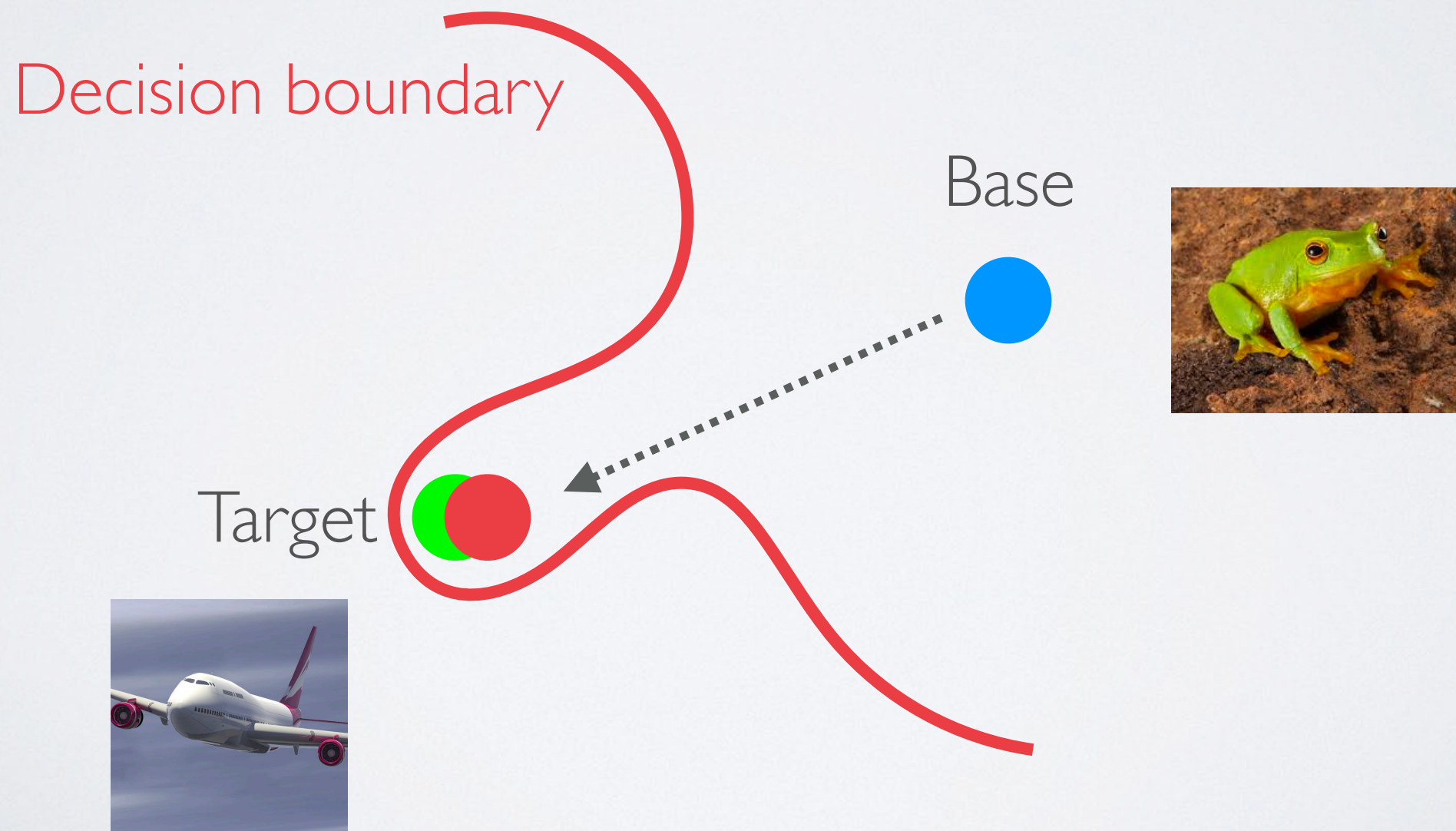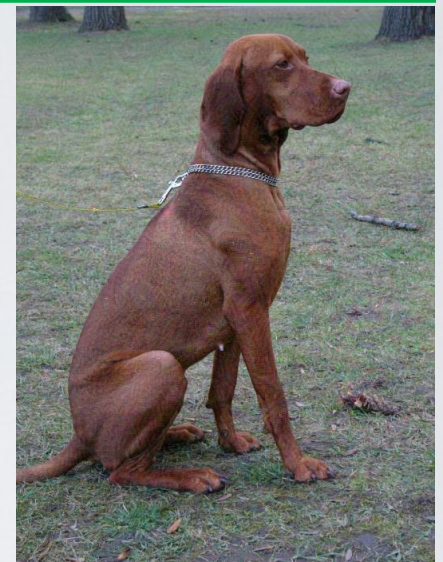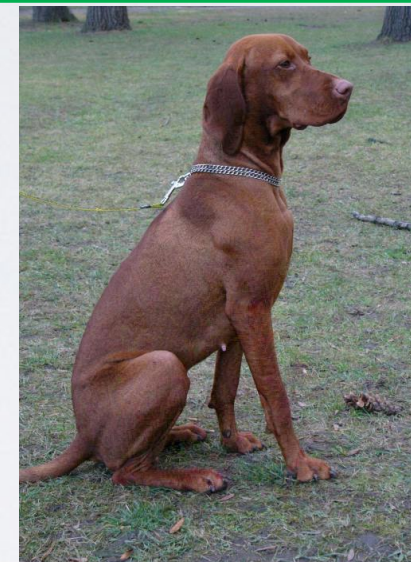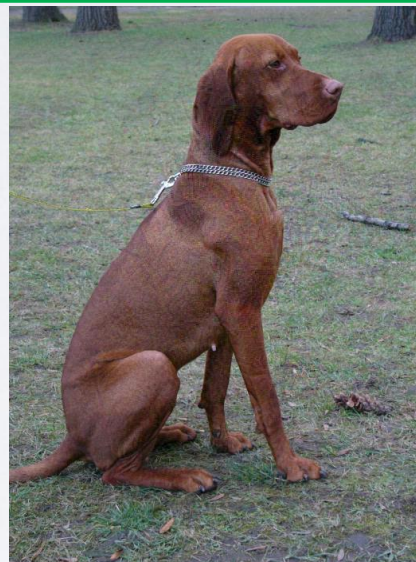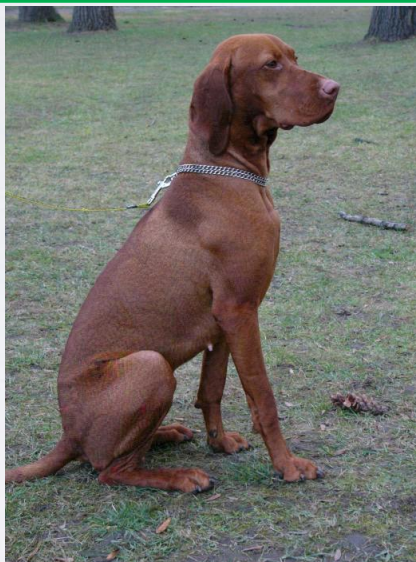
$$\mathbf{p} = \operatorname*{argmin}_{\forall \mathbf{x}} \ \|f(\mathbf{x}) - f(\mathbf{t})\|^2 + \beta \|\mathbf{x} - \mathbf{b}\|^2 \qquad (1)$$

Decision boundary

Base

Target

# COLLISION ATTACK

$$\mathbf{p} = \operatorname*{argmin}_{\forall \mathbf{x}} \ \|f(\mathbf{x}) - f(\mathbf{t})\|^2 + \beta \|\mathbf{x} - \mathbf{b}\|^2 \qquad (1)$$

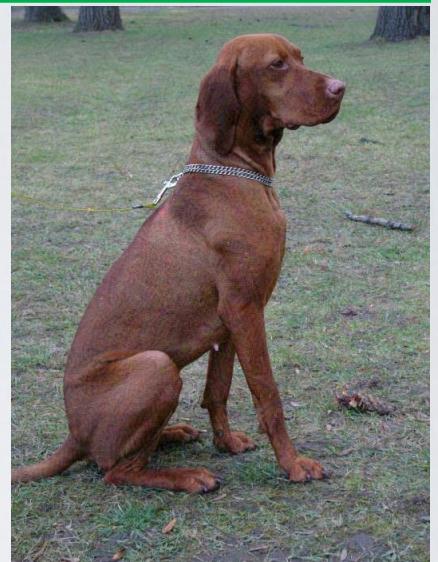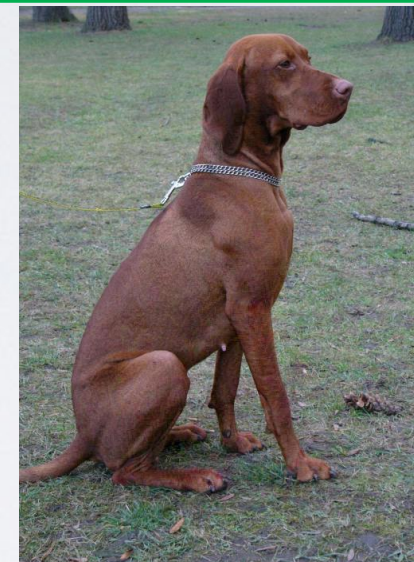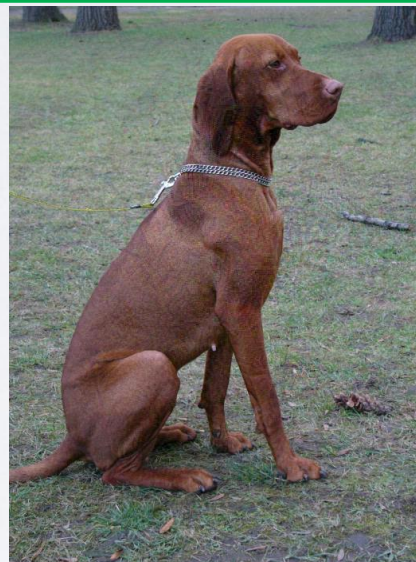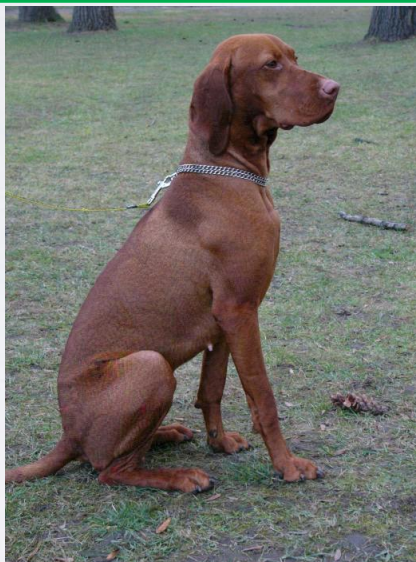Decision boundary

Base

Target

Target instances from Fish class

Clean Base

Original image

Clean Base

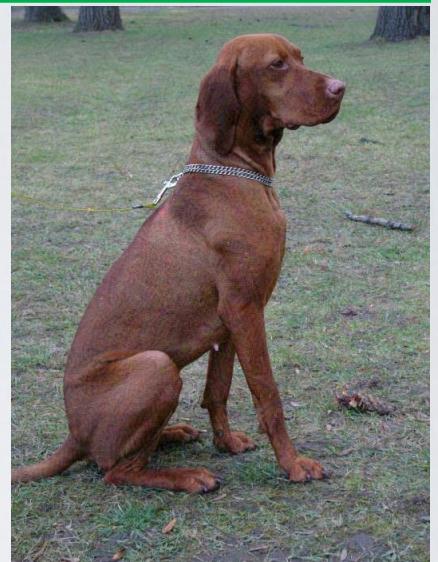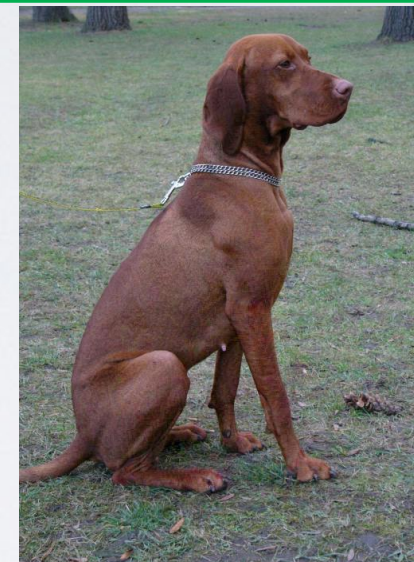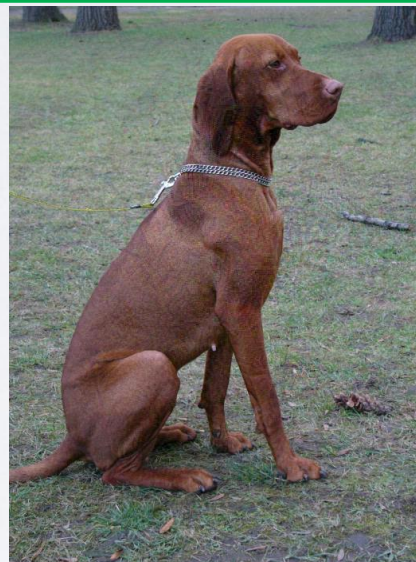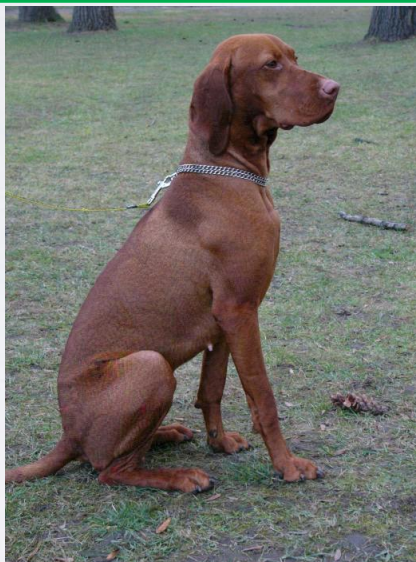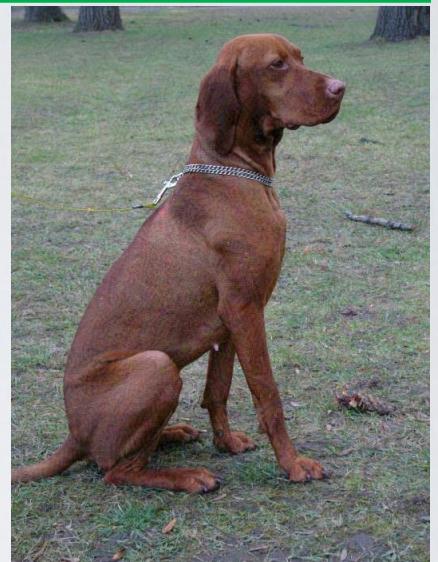Target instances from Fish class
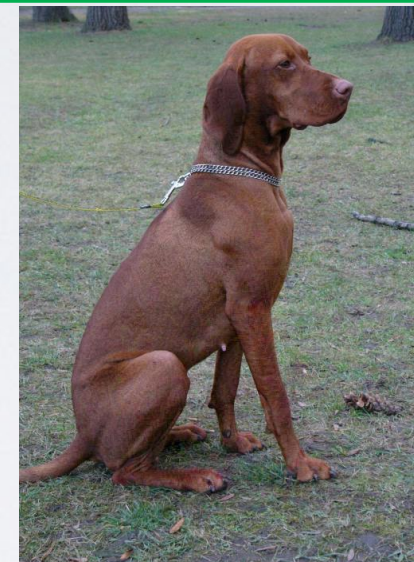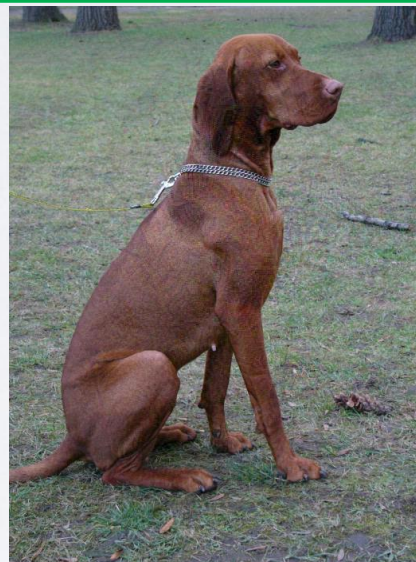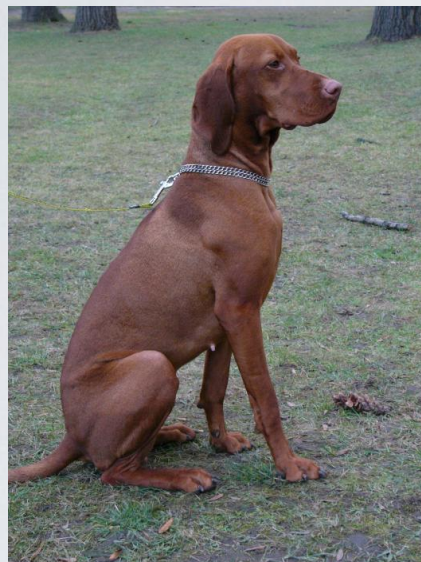
poison

Target instances from Fish class

Clean Base
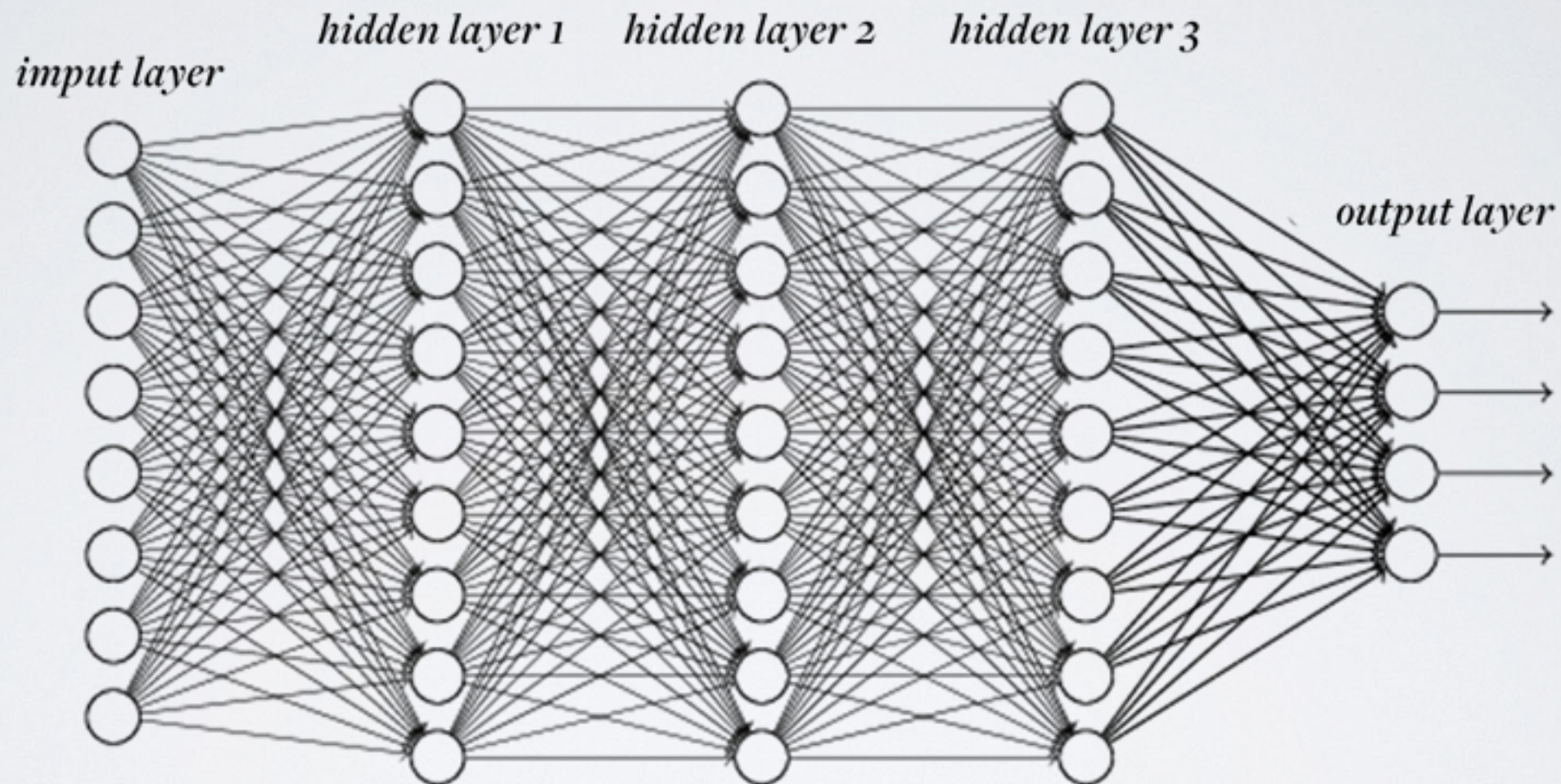
poison

Clean Base

Target instances from Fish class

poison

# END-TO-END TRAINING?

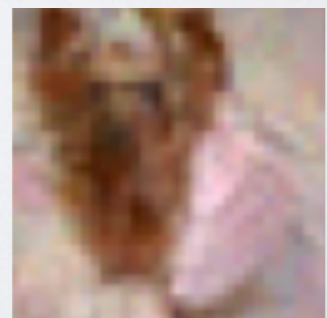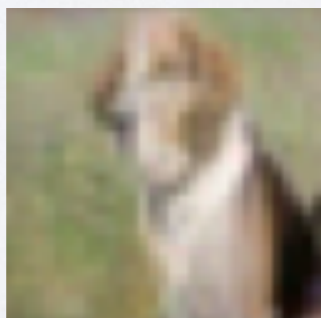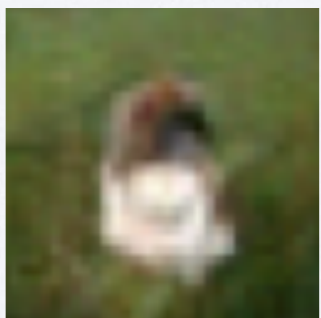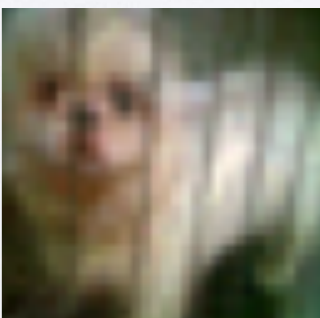**Feature extractors learn to ignore adversarial perturbation**



**Feature extraction layers**

# OH NO! POISON DOGS!

60 poison dogs cause a bird to be mis-classified

# THEORY OF ADVERSARIAL EXAMPLES

# ATTACK & DEFENSES

**Adversarial attacks**
Szegedy et al, 2013
Biggio et al, 2013

**Adversarial training**
Goodfellow et al 2015

**Multi-stage attacks**
Kurakin et al, 2016
Tramer et al, 2017

**Distillation**  Papernot '16
**Bounded relu**  Zantedeschia '16
**MagNet**  Meng & Chen '17

**Optimization attacks**
Carlini & Wagner '17

**Thermometer**  Buckman '18
**Detection**  Ma et al, '18
**Compression**  Guo, '18
**GANs**  Samangouei, '18

**Approximation attacks**
Athalye et al, 2018

…and **LOTS** more

# ARE ADVERSARIAL EXAMPLES **INEVITABLE**?

# RELATED WORK

**K-nearest neighbors classifier**

"Analyzing the Robustness of Nearest Neighbors to Adversarial Examples"
Wang, Jha, Chaudhuri, 2017

**Datasets produced by GAN-type generator**

"Adversarial vulnerability for any classifier"
Fawzi, Fawzi, Fawzi, 2018

**Classes lie on concentric spheres**

"Adversarial spheres"
Gilmer, Metz, Faghri, Schoenholz, Raghu, Wattenberg, Goodfellow, 2018

**Most similar to ours…**

"The Curse of Concentration in Robust Learning"
Mahloujifar, Diochnos, Mahmoody, 2018

# ARE ADVERSARIAL EXAMPLES
# INEVITABLE?

**spoiler alert**

**...and the answer is...**

# YES!

...if the adversary is strong enough.

# ARE ADVERSARIAL EXAMPLES INEVITABLE?

...but computer scientists think...

# NO!

Common assumptions…

**Human perception is not exploitable**

**High dimensional spaces aren't too weird**

**Adversarial example**

$$\|x - \hat{x}\|_p < \epsilon.$$

# TOY PROBLEM

**Dimension**

3

# TOY PROBLEM

**Dimension**

3

**Surface area**

50%



$\mathcal{A}$  $\mathcal{B}$

Adversarial
examples?

# TOY PROBLEM

**Dimension**

3

**Surface area**

55%

$\epsilon = 0.1$

$\mathcal{A}$ $\mathcal{B}$

# TOY PROBLEM

**Dimension**

100

**Surface area**

84%

# TOY PROBLEM

$\epsilon = 0.1$

**Dimension**

1000

**Surface area**

99.8%

$\mathcal{A}$    $\mathcal{B}$

**random sampling** $\longrightarrow$ **adversarial susceptibility**

# Theorem (Levy & Pellegrino, 1951)

The $\epsilon$-expansion of *any* set that occupies half the sphere is at least as big as the $\epsilon$-expansion of a semi-sphere.



**This classifier** is worse than **this classifier**

# WHAT ABOUT *REALISTIC* MODELS?
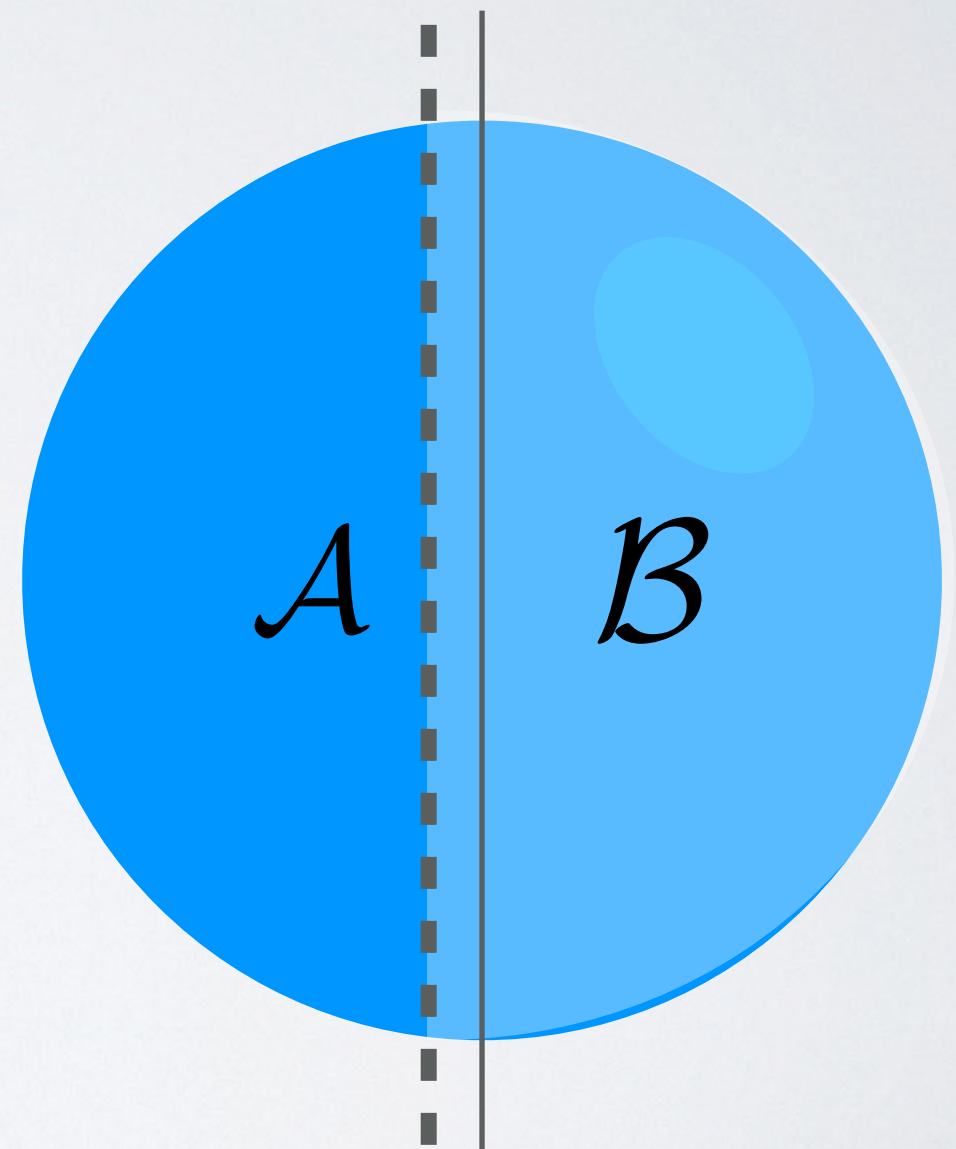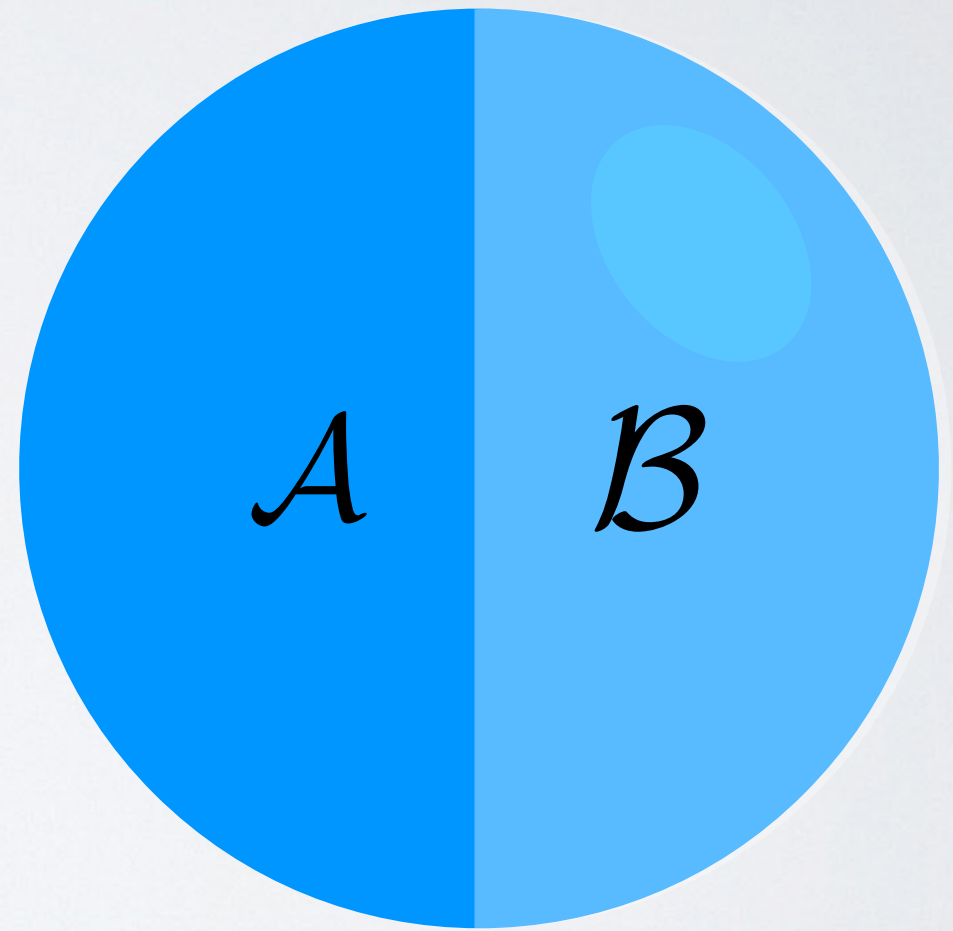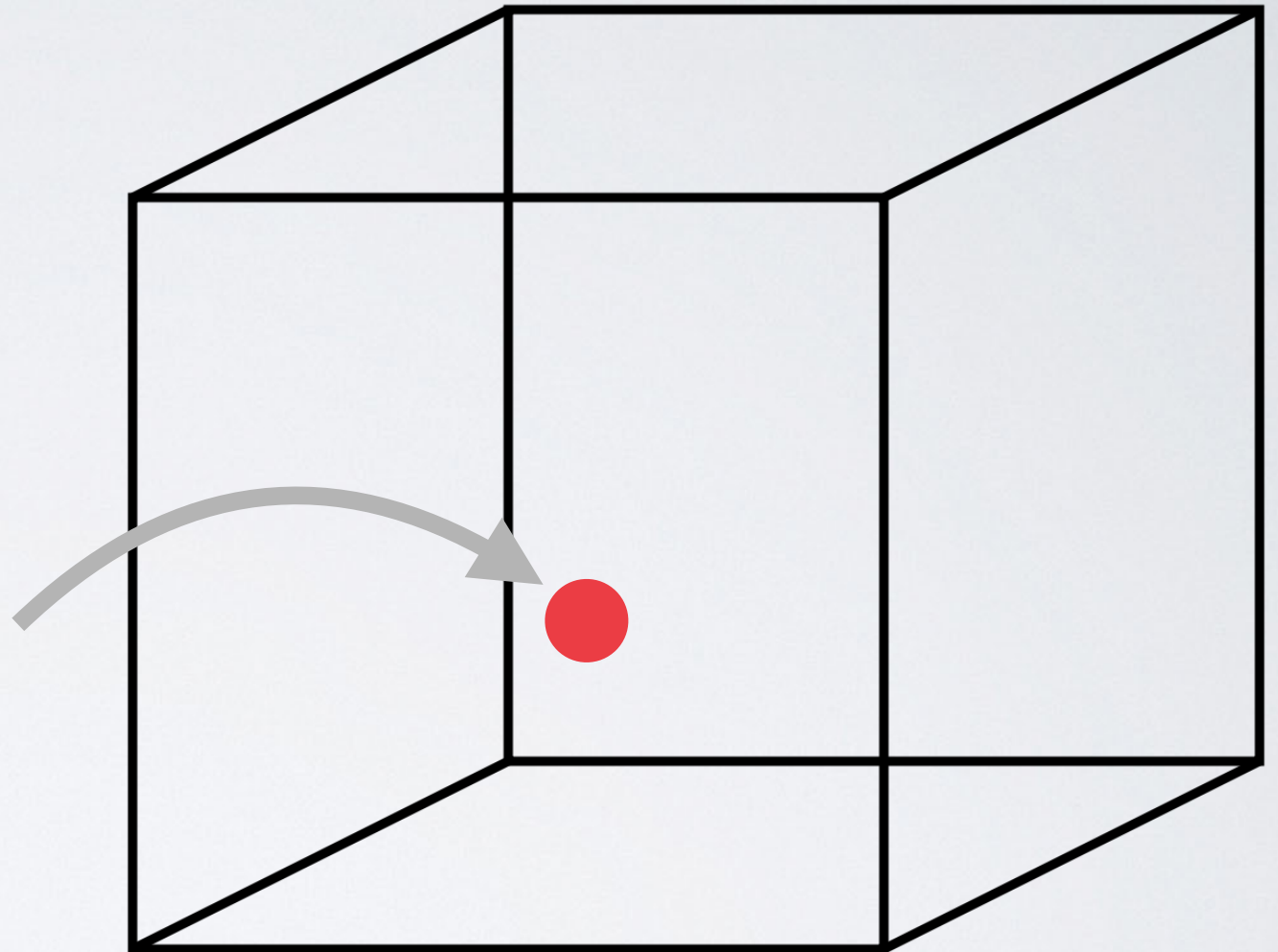
# THE SETUP
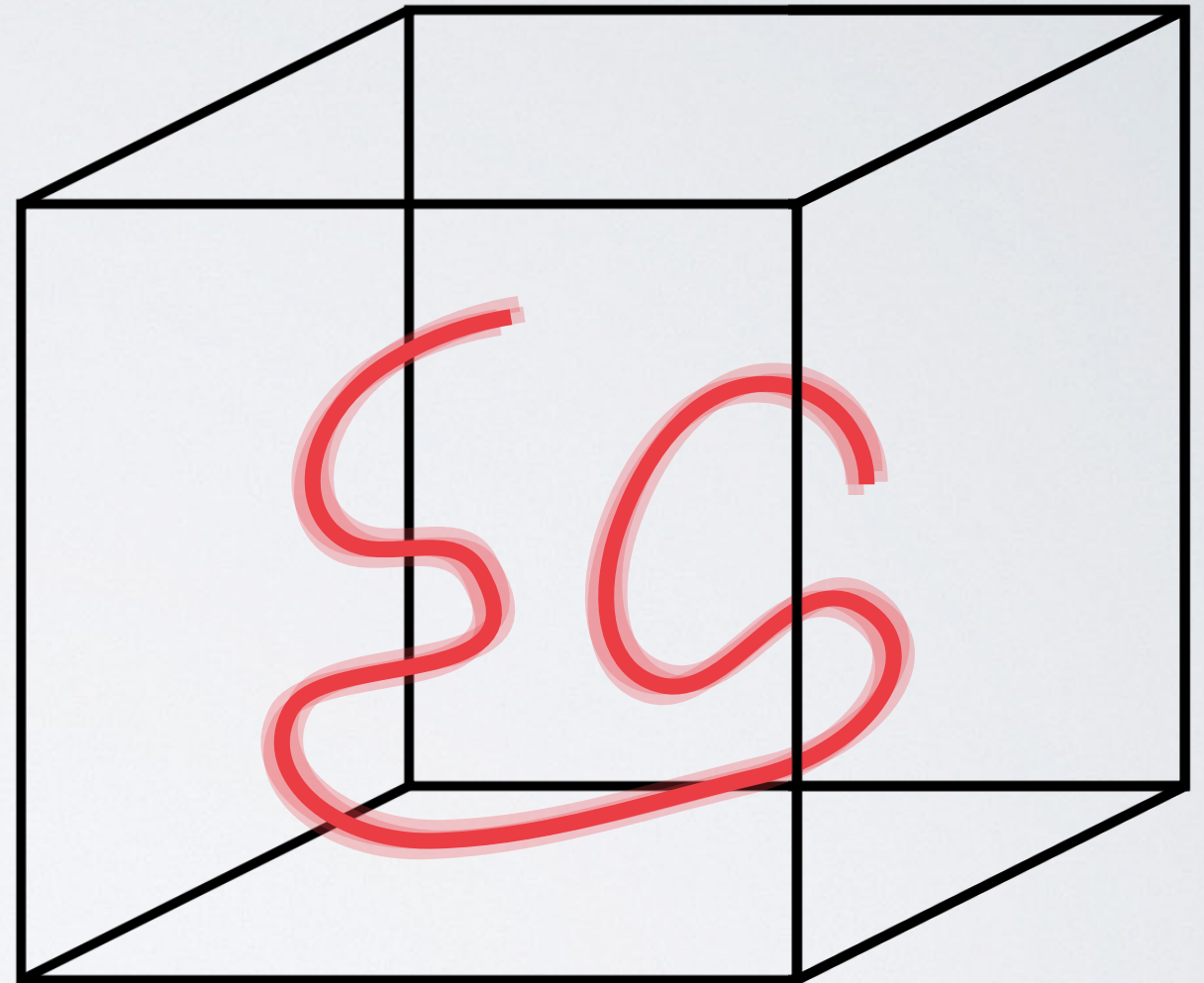
**Images**

Points in a unit cube

# THE SETUP

**Images**

Points in a unit cube

**Class**

Probability density
function on cube
(bounded by $U_c$)
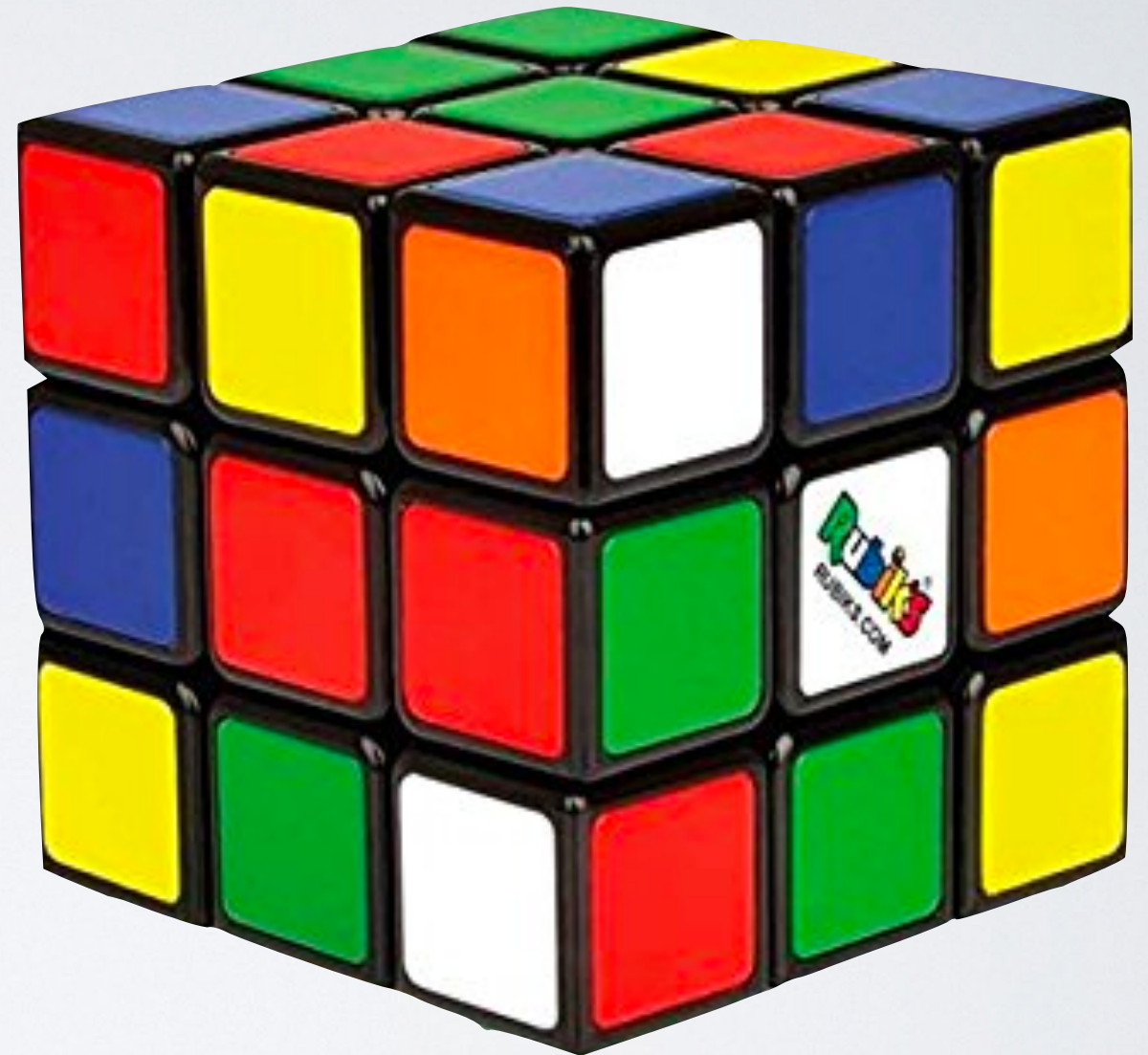
# THE SETUP

**Images**

Points in a unit cube

**Class**

Probability density function on cube (bounded by $U_c$)

**Classifier**
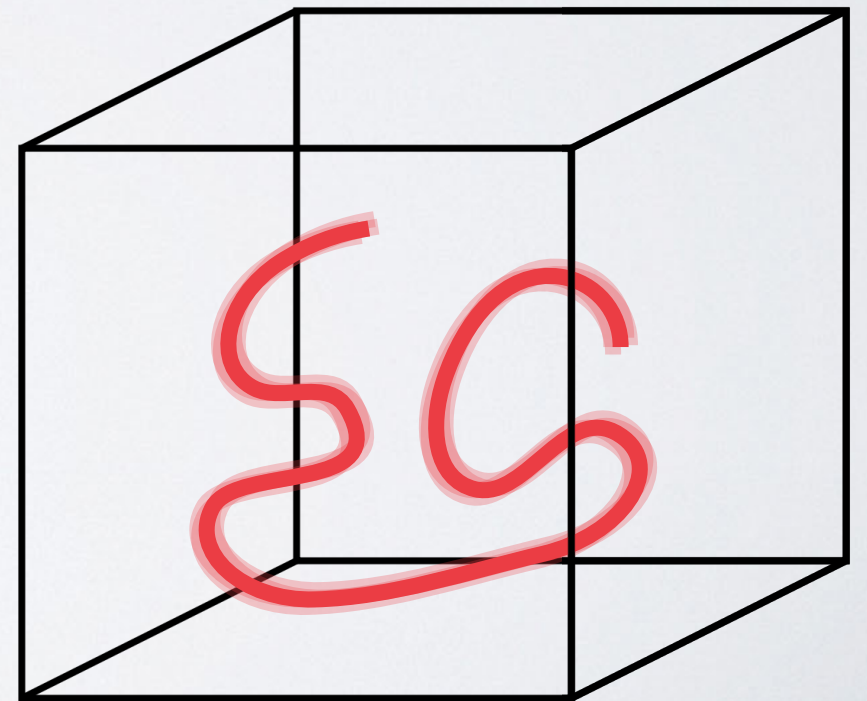
Partitions cube into disjoint sets (measurable)

# "MOST" THINGS ARE ADVERSARIAL

**Theorem**

Choose a class $c$ that occupies less than half the cube according to the classifier. Define...

$U_c$ : supremum of the density function for class $c$

## Theorem

Choose a class $c$ that occupies less than half the cube according to the classifier. Define...

$U_c$ : supremum of the density function for class $c$
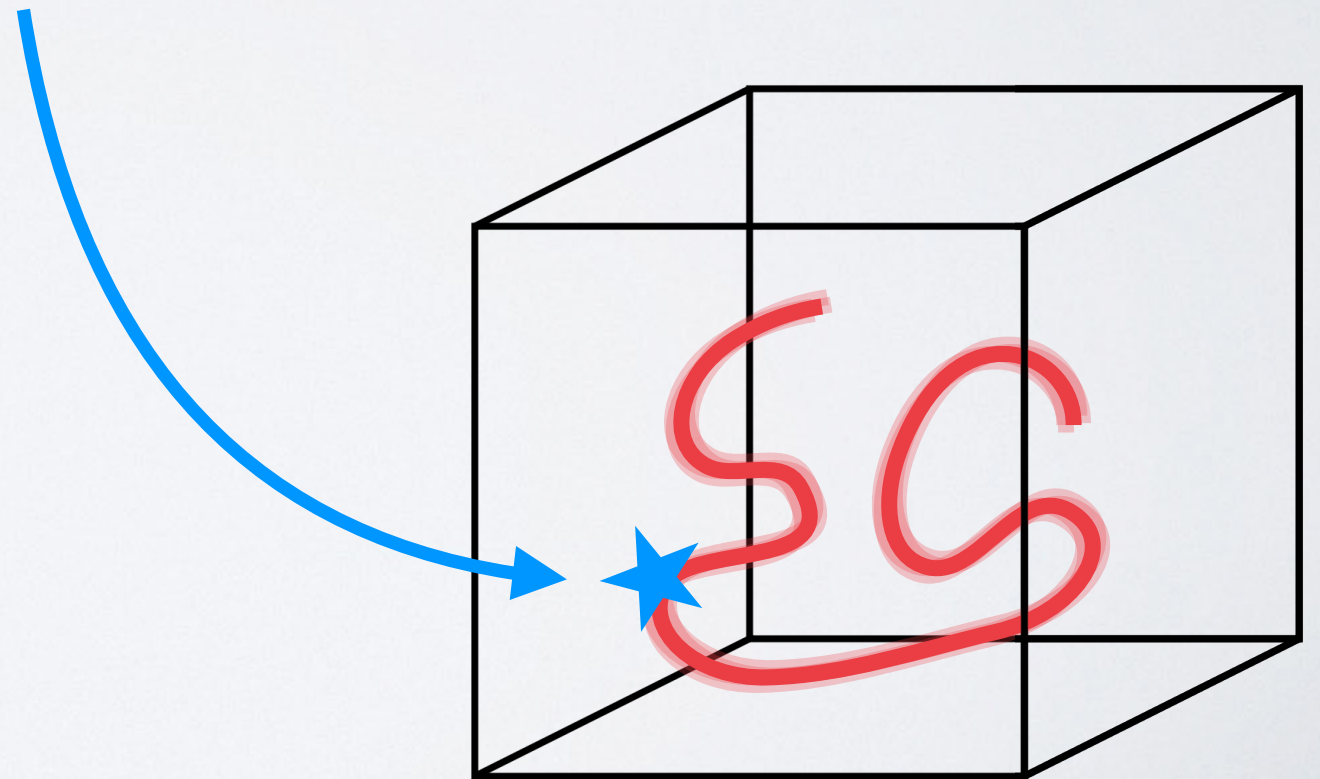
Sample a random point $x$ from the class distribution. With probability at least
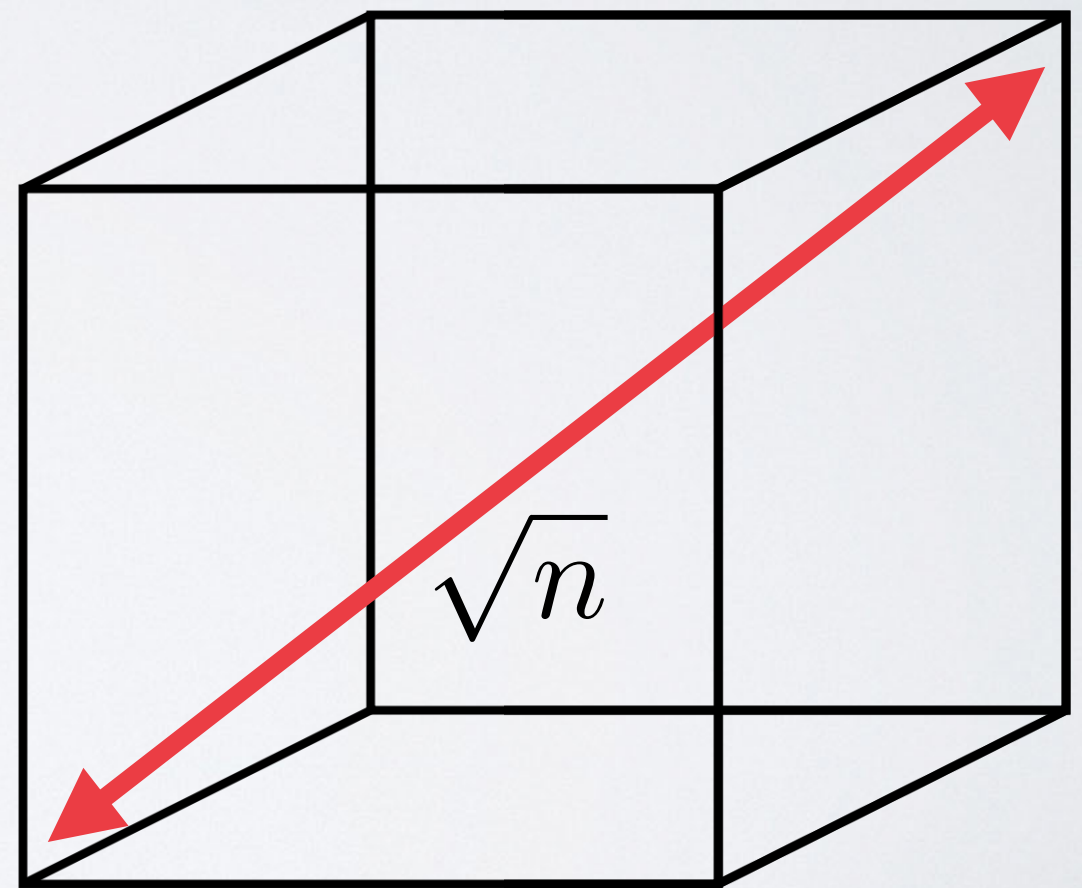
$$1 - U_c \exp(-\pi \epsilon^2)$$

One of the following conditions holds:

- $x$ is misclassified by the classifier
- $x$ has an adversarial example $\hat{x}$ with $\|x - \hat{x}\|_2 < \epsilon$.

# "MOST" THINGS ARE ADVERSARIAL

$$1 - U_c \exp(-\pi\epsilon^2)$$

$$\epsilon = 10$$

**Adversarial example**

An image $x$ has an $\epsilon$-adversarial example in the $p$ norm if there is a point $\hat{x}$ in a different class with

$$\|x - \hat{x}\|_p < \epsilon.$$

$p = 0$

$$\|x - \hat{x}\|_0 = \mathrm{card}\{i | x_i \neq \hat{x}_i\}$$
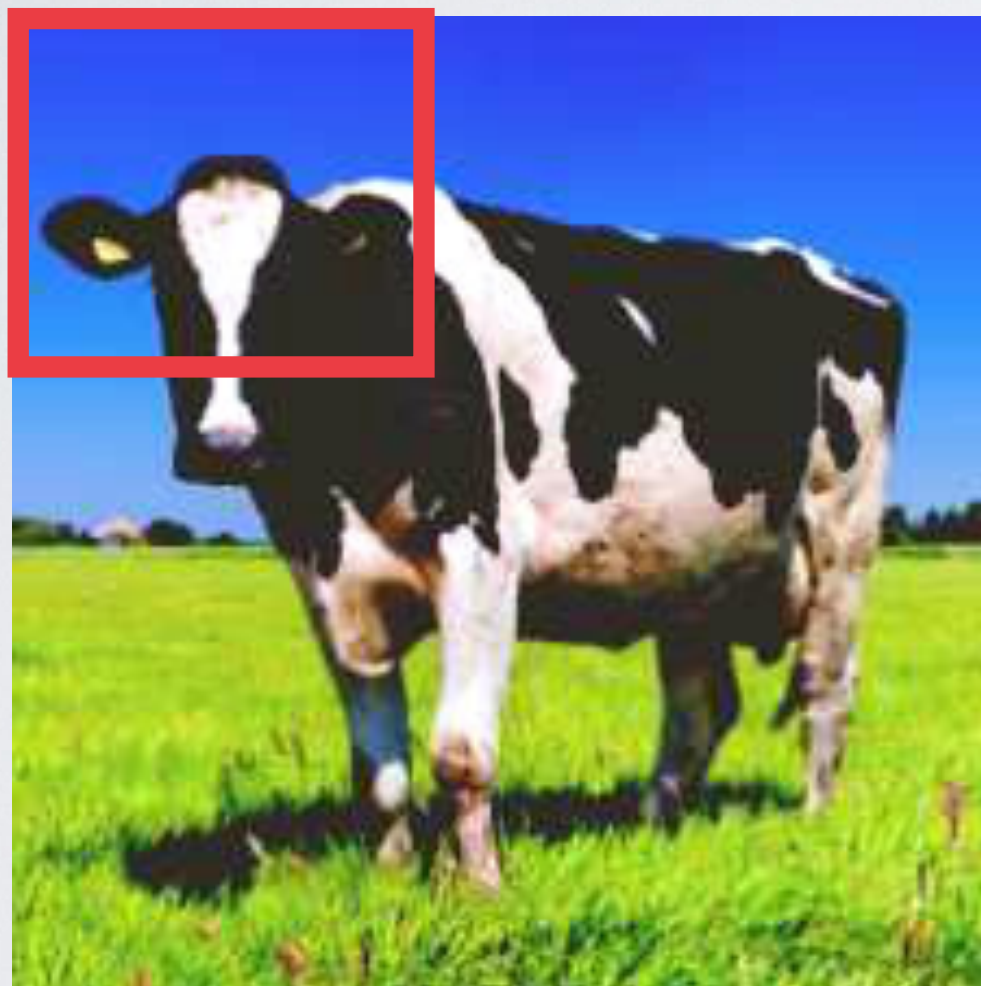
**Sparse** adversarial example

# SPARSE ATTACKS

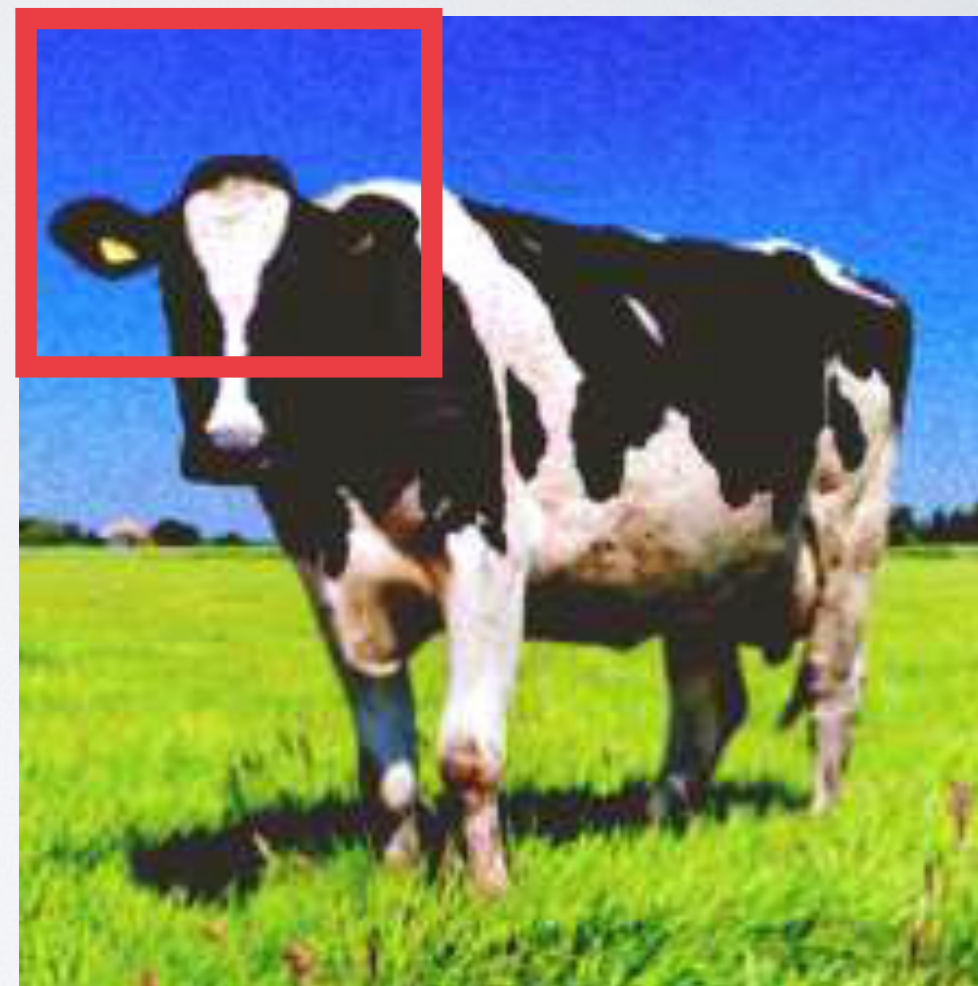2-norm attack



"Ox"

"Traffic Light"
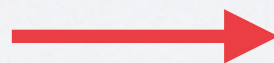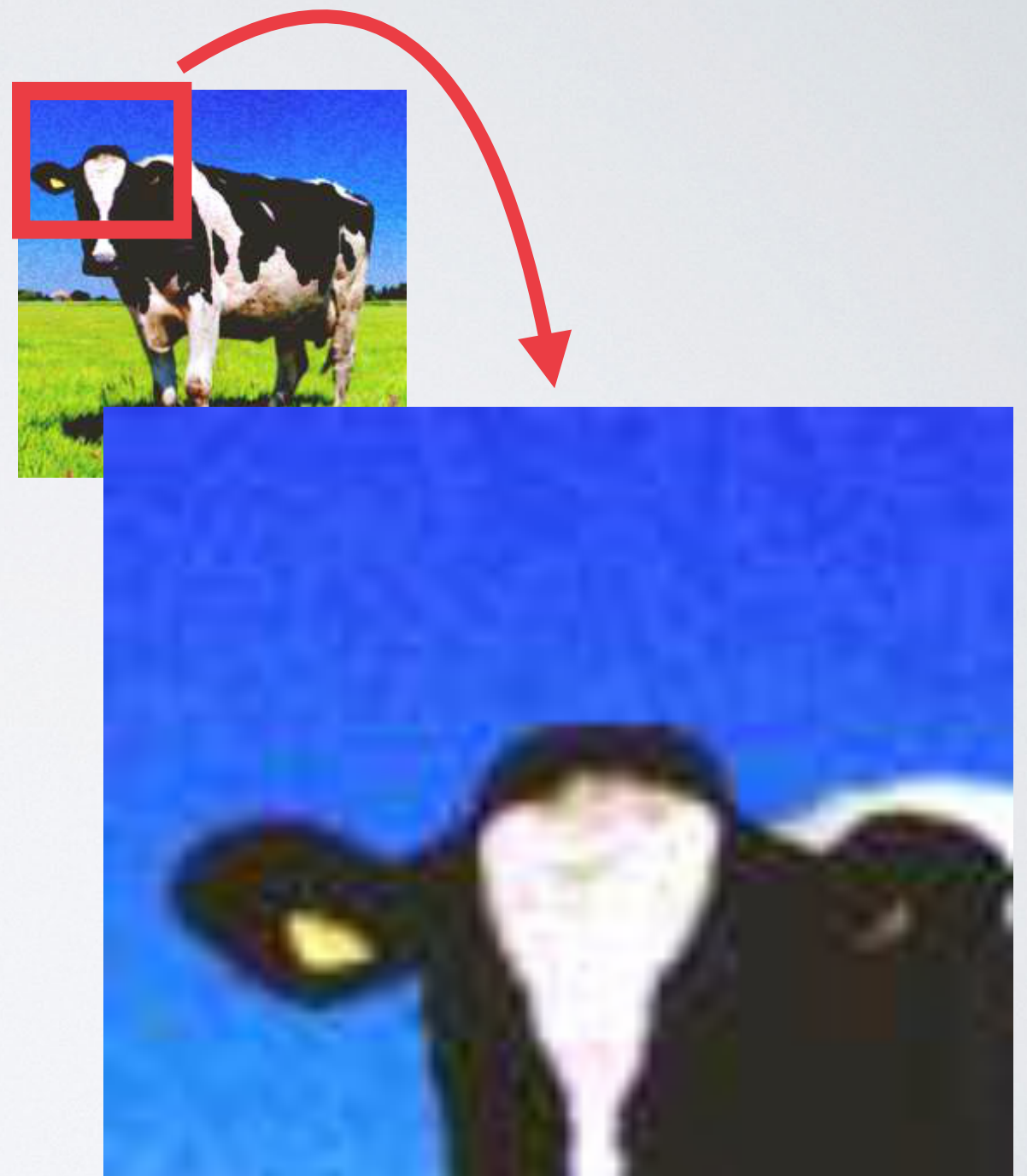
# SPARSE ATTACKS

2-norm attack



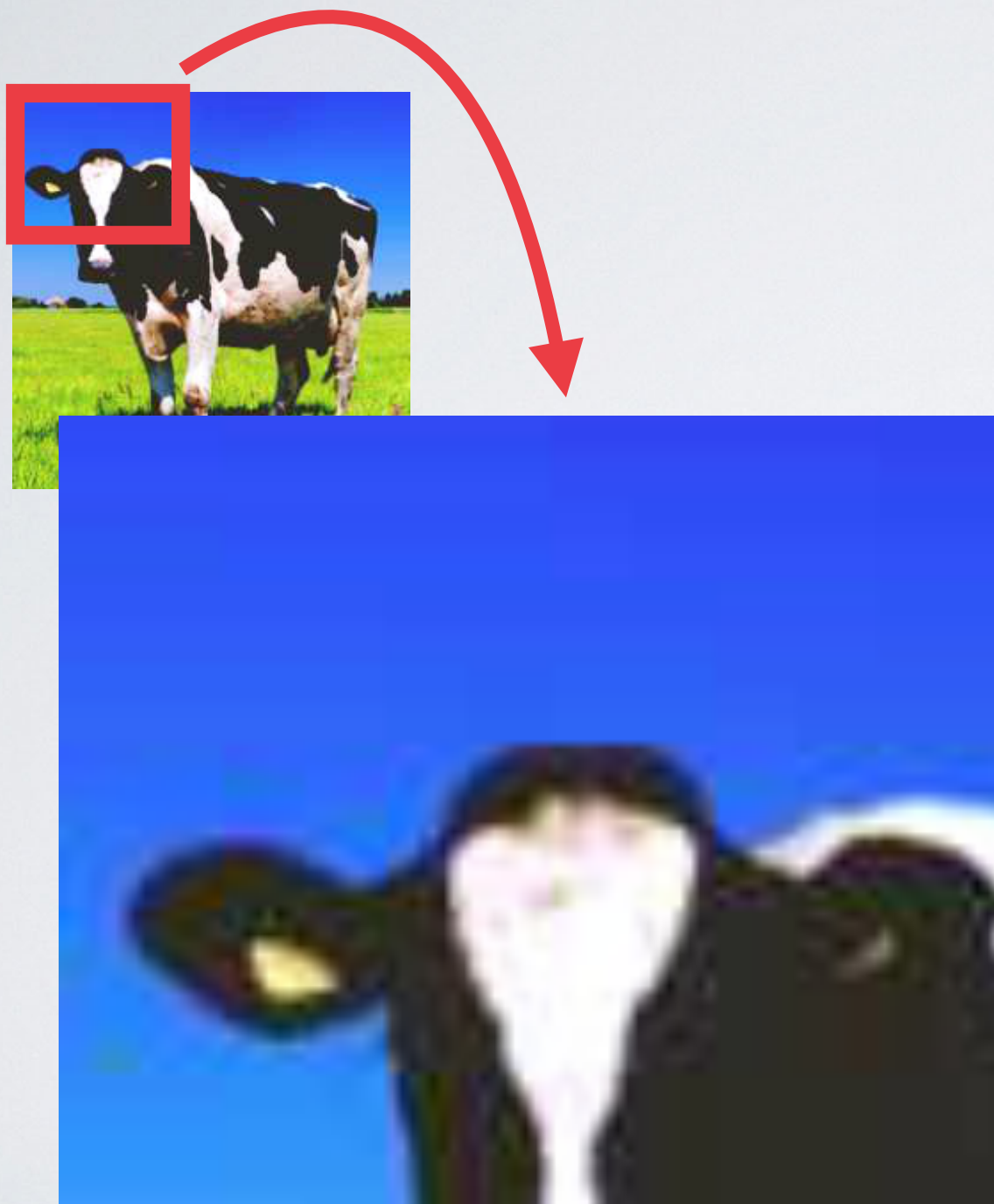"Ox"

"Traffic Light"

# SPARSE ATTACKS

# SPARSE ATTACKS

3% pixels changed



"Ox"                    "Traffic Light"

# SPARSE ADVERSARIAL EXAMPLES

## Theorem

Choose a class $c$ that occupies less than half
the cube according to the classifier. Define...

$U_c$ : supremum of the density function for class $c$

Sample a random point $x$ from the class distribution.
With probability at least

$$1 - 2U_c \exp(-k^2/n)$$

# of pixels
changed
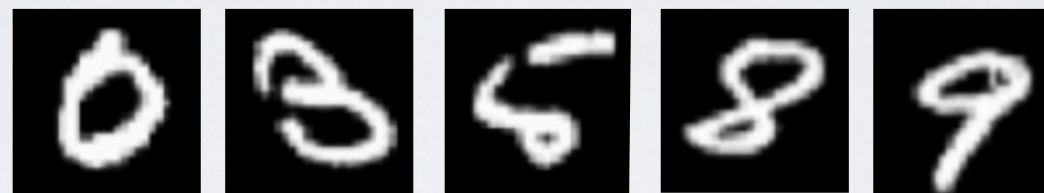
One of the following conditions holds:

- $x$ is misclassified by the classifier

- The label of $x$ can be changed by
  modifying at most $k$ pixels.

# WHAT ABOUT HIGH DIMENSIONS?

# WHAT ABOUT HIGH DIMENSIONS?
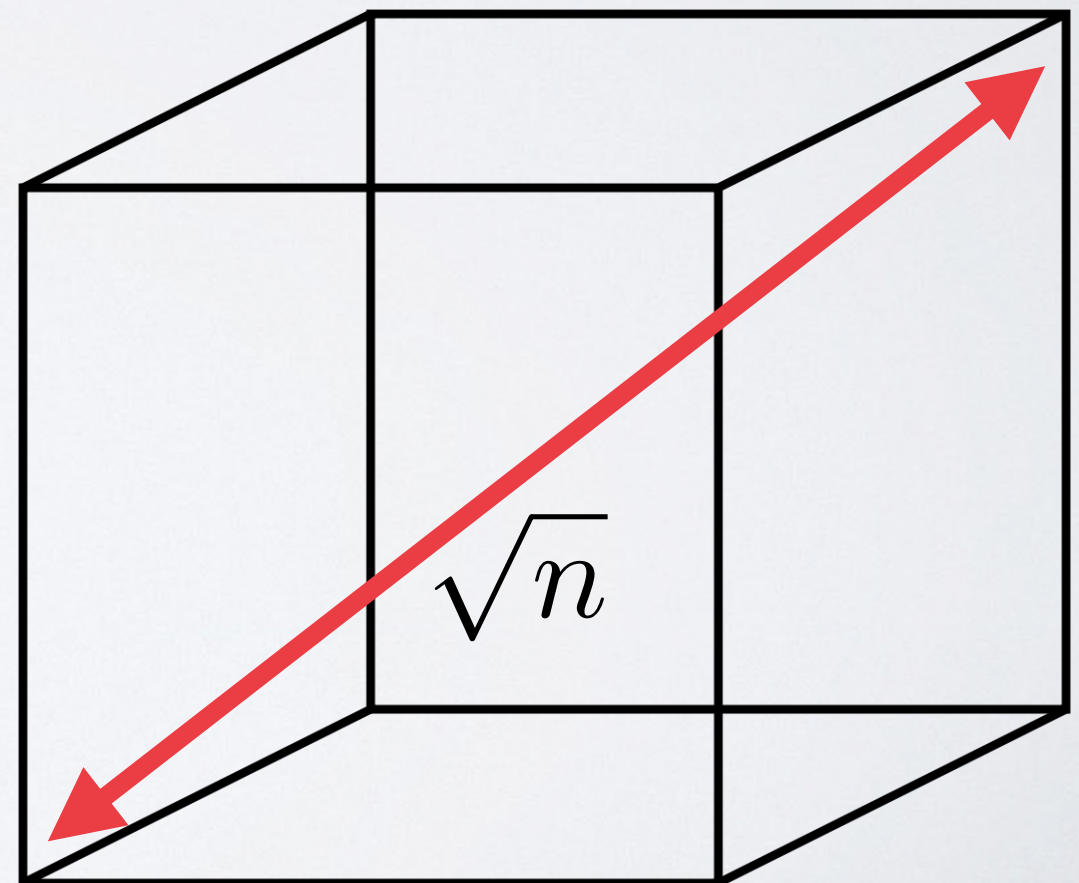


Clean

Adversarial

"dog" 9%

"traffic light" 97%

# BOUNDS IN HIGH DIMENSIONS

$$\epsilon = O(\sqrt{n})$$

$$1 - U_c \exp(-\pi\epsilon^2)$$

Does this stay the same for large n?
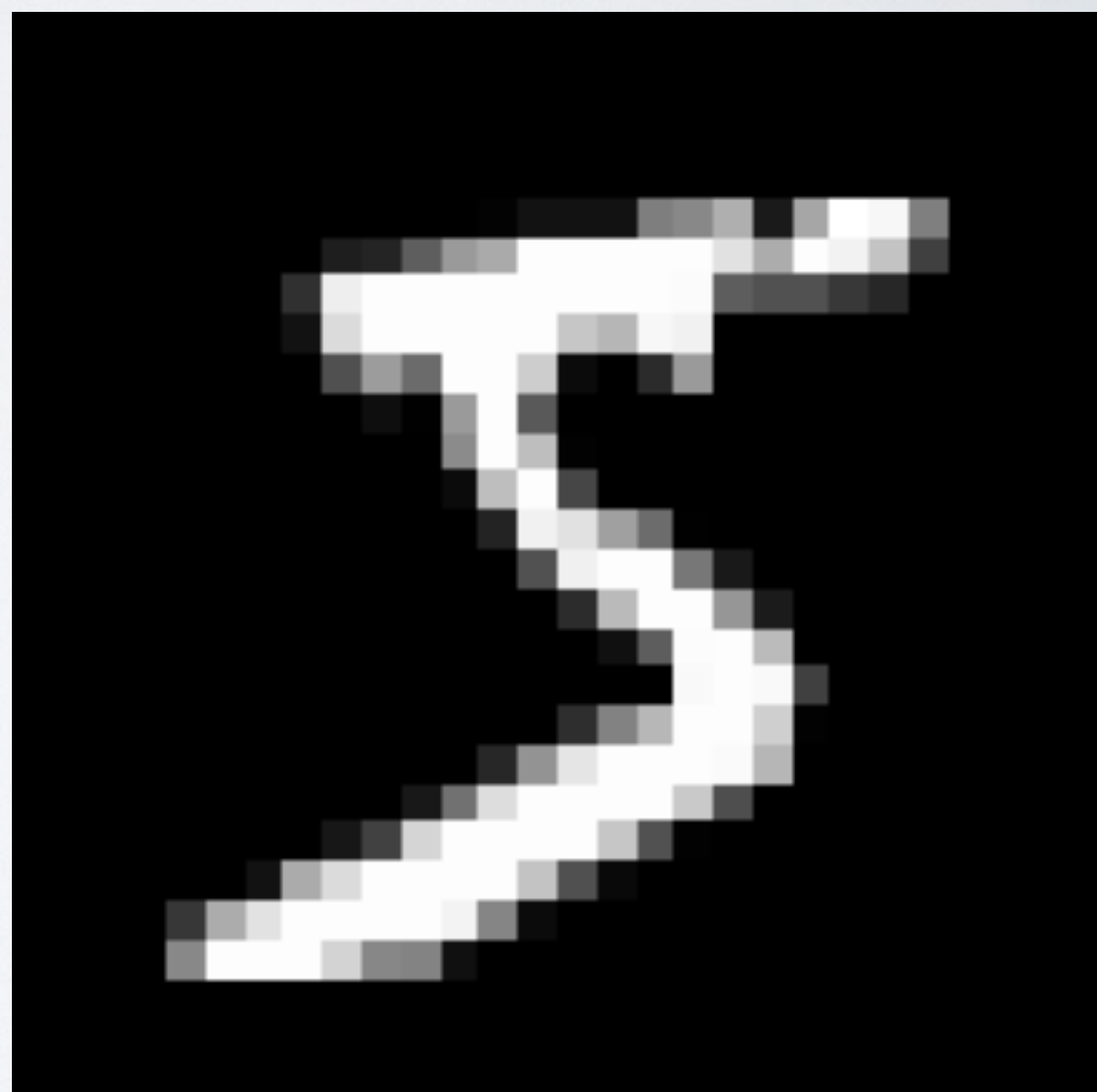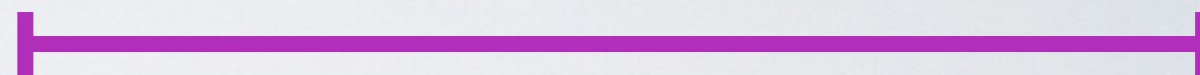
**NOPE!**
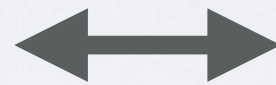
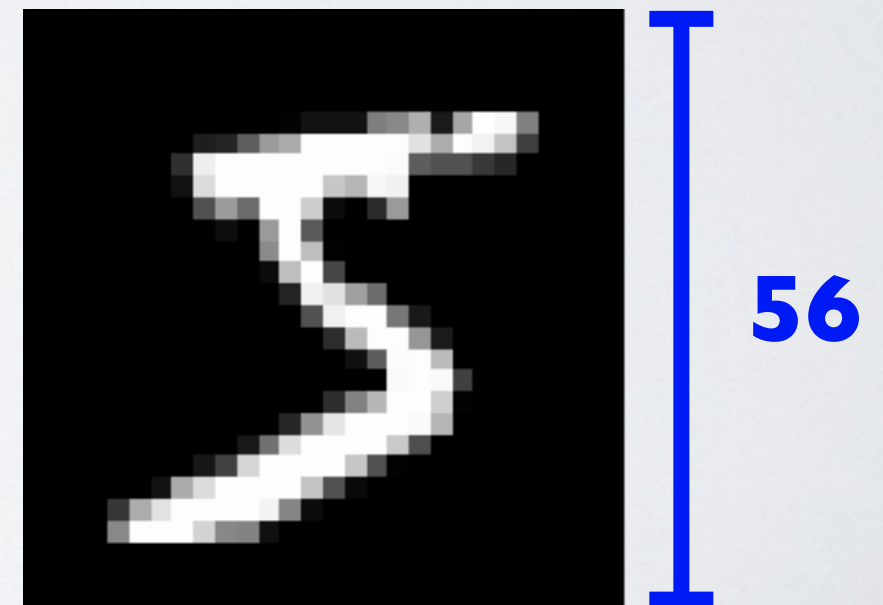$$\sqrt{n}$$

# BIG MNIST

**28**

**56**

**112**

# Theorem

## 28x28 MNIST

For all classifiers, a random image has an $\epsilon$-adversarial example with probability $p$.
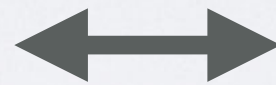
## 56x56 MNIST

For all classifiers, a random image has an $2\epsilon$-adversarial example with probability $p$.

$\longleftrightarrow$

28

56

# Theorem

## 28x28 MNIST

For all classifiers, a random image has an $\epsilon$-adversarial example with probability $p$.
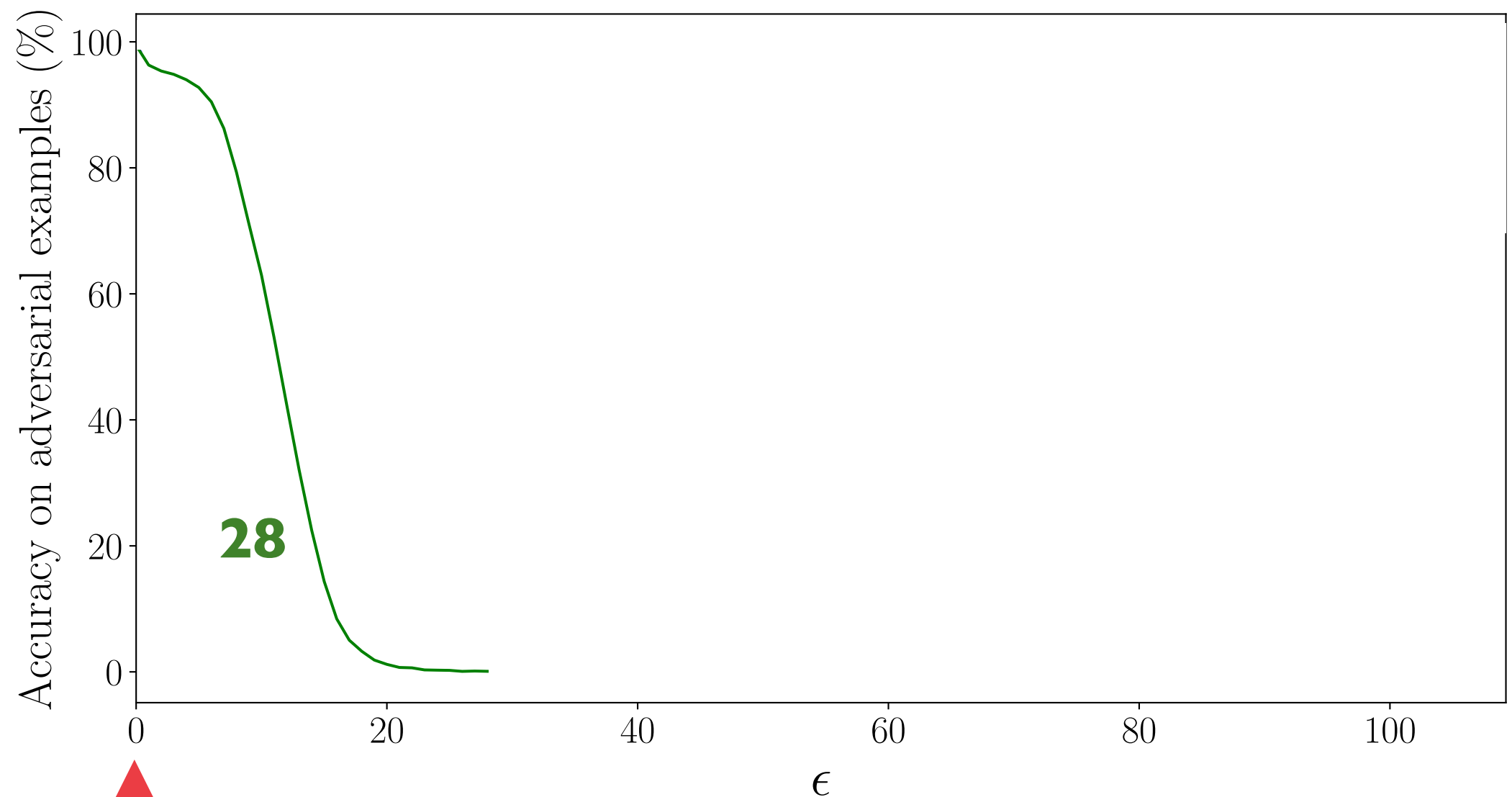
$\longleftrightarrow$

## 56x56 MNIST

For all classifiers, a random image has an $2\epsilon$-adversarial example with probability $p$.

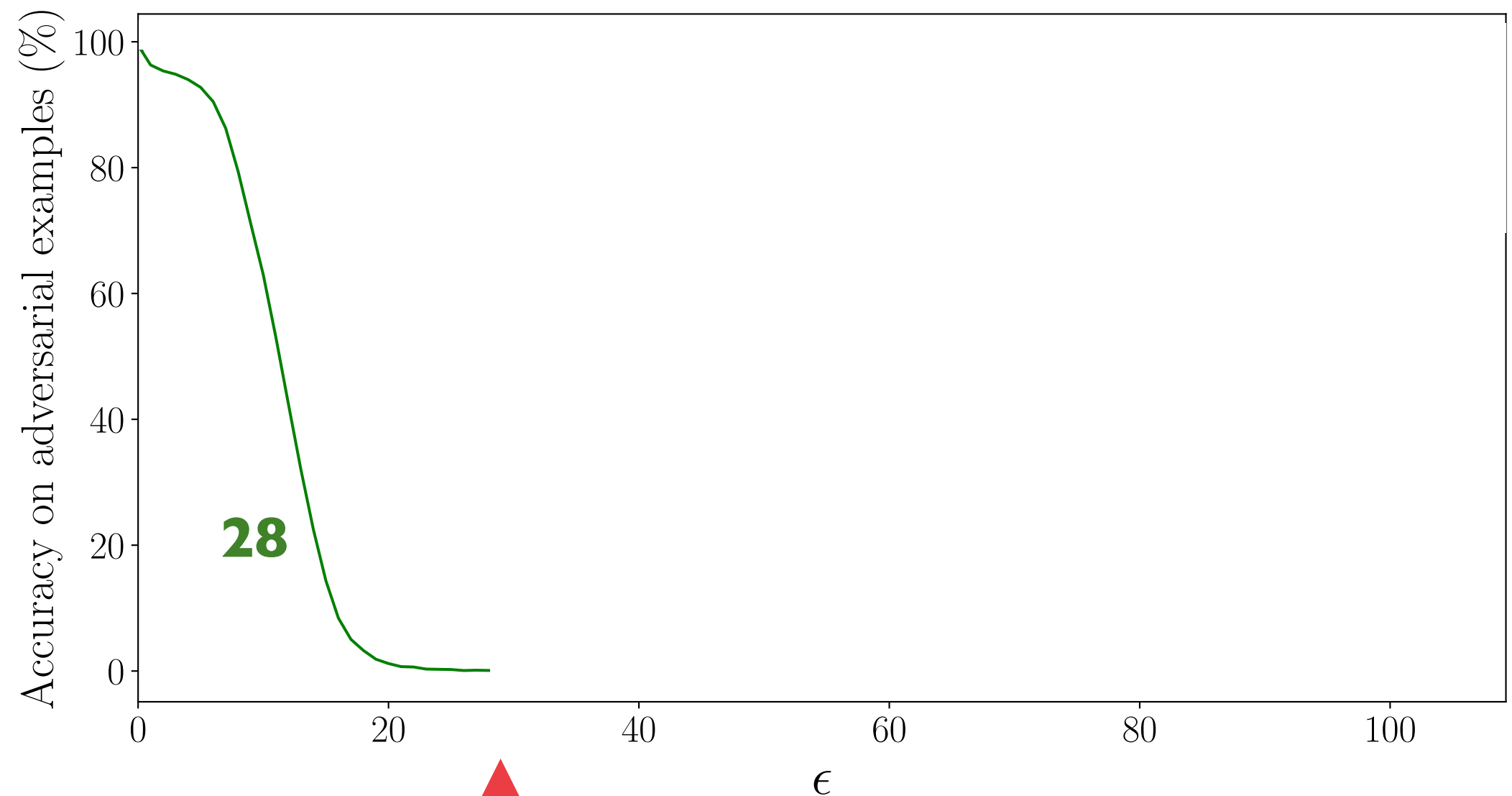**There is no relation between dimensionality and robustness!**

# ADVERSARIAL TRAINING

MNIST hardened using PGD (30 steps)



**High accuracy**

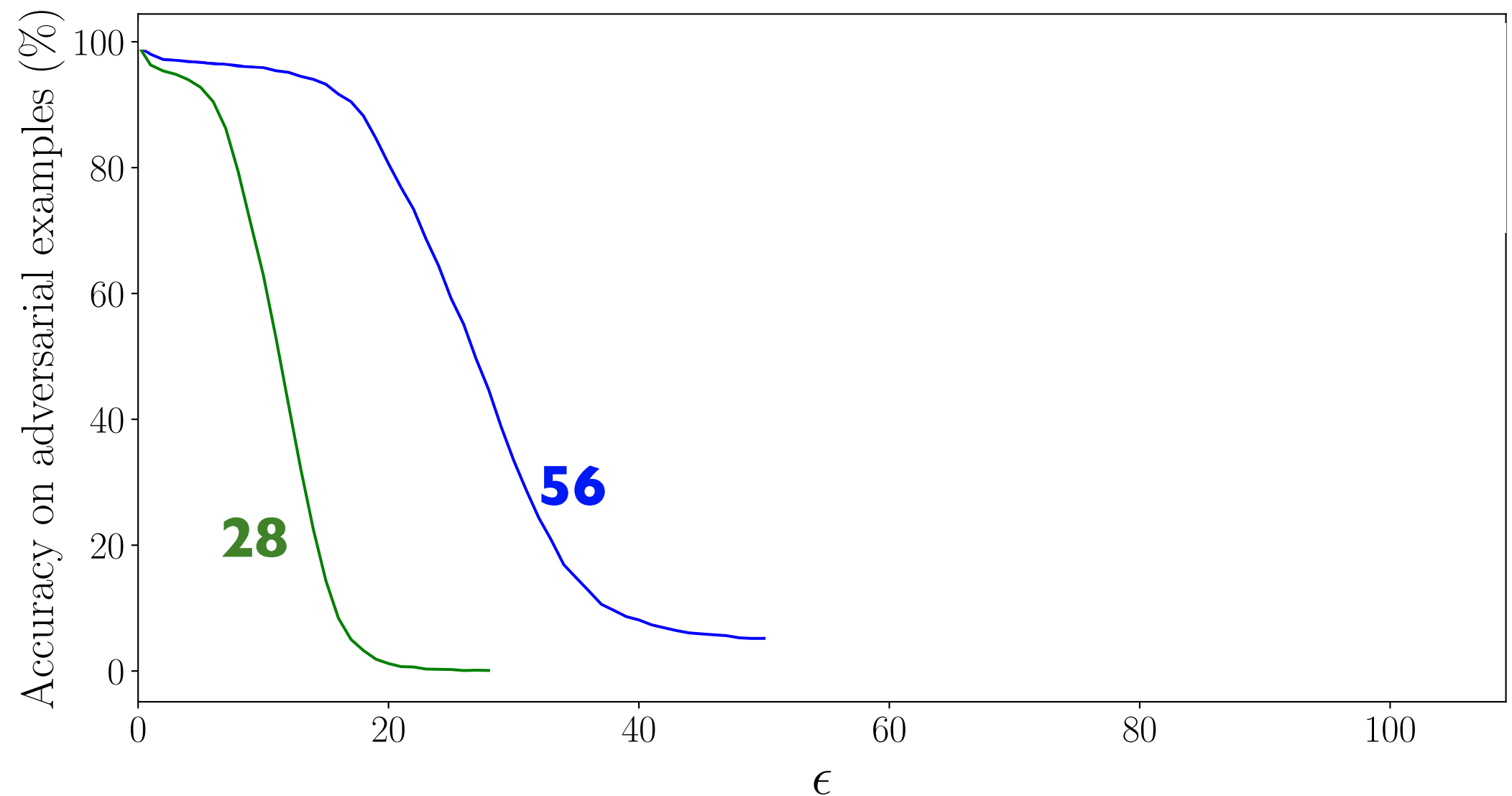# ADVERSARIAL TRAINING

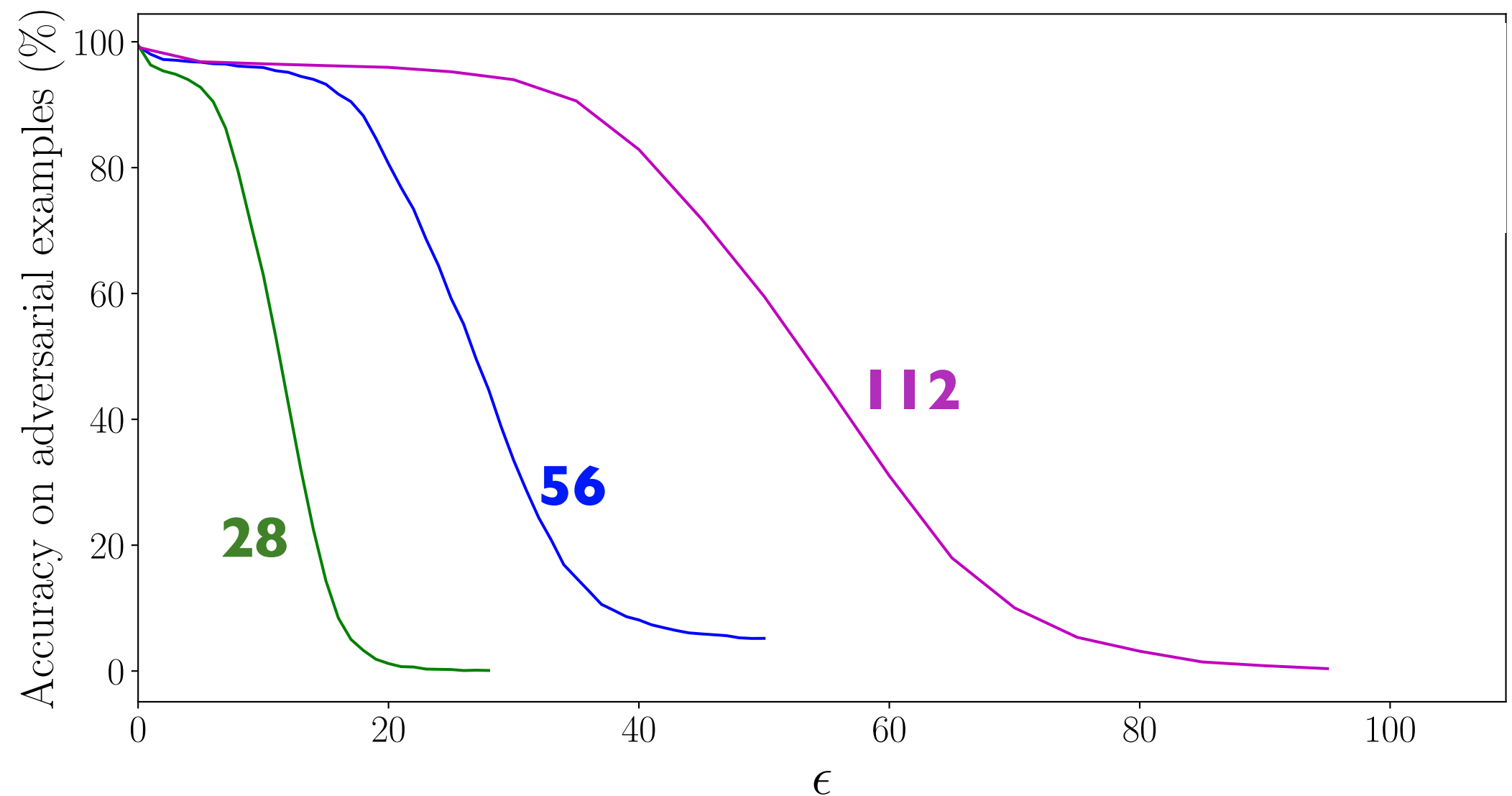MNIST hardened using PGD (30 steps)

# ADVERSARIAL TRAINING

MNIST hardened using PGD (30 steps)
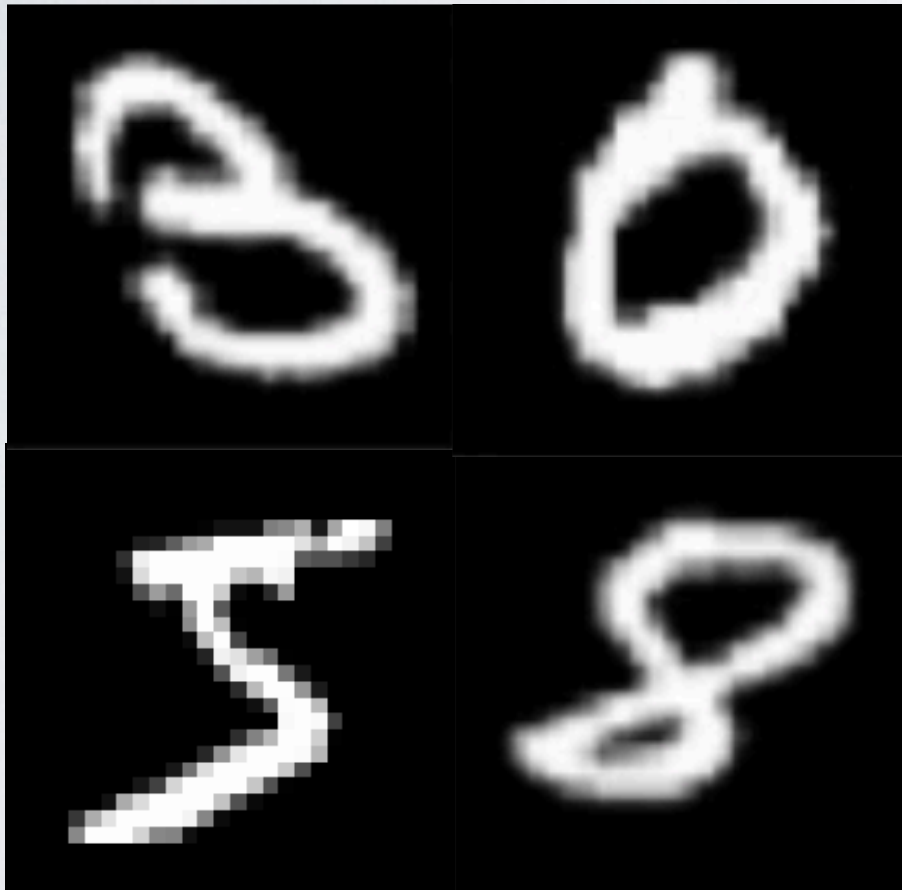
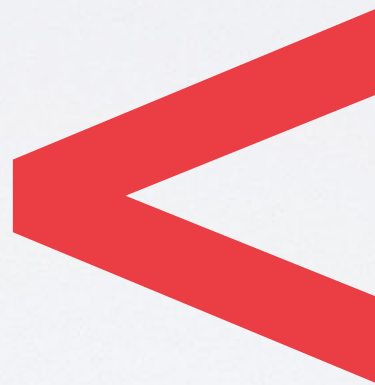# ADVERSARIAL TRAINING

MNIST hardened using PGD (30 steps)
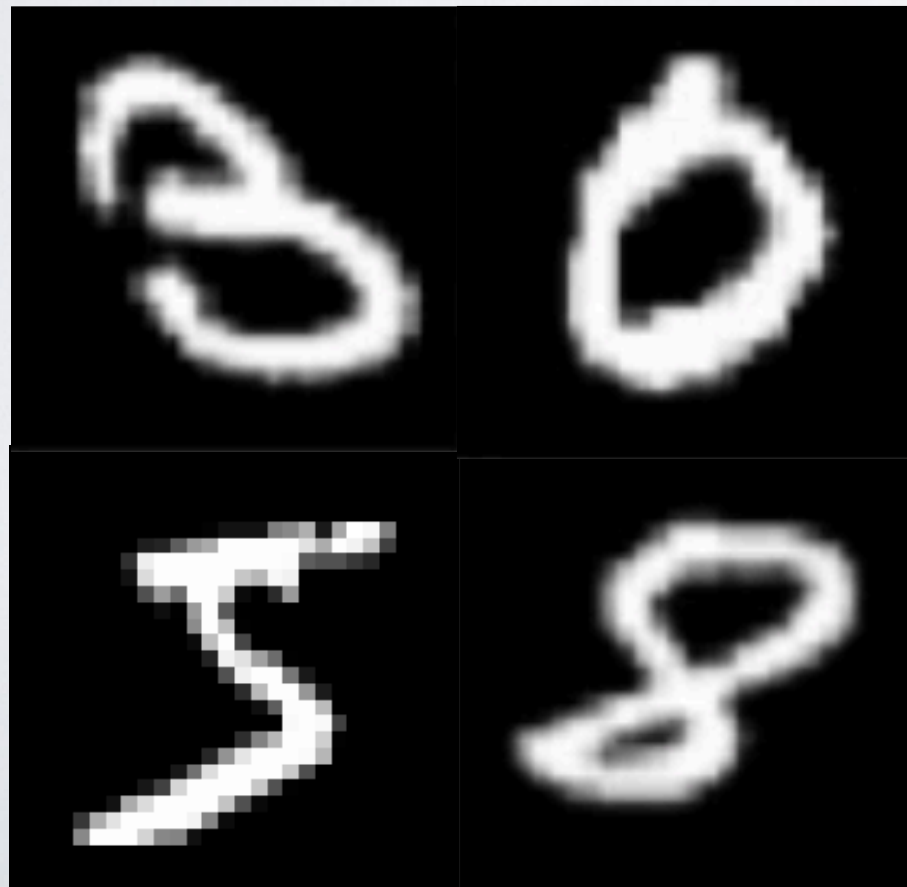
# WHAT AFFECTS ROBUSTNESS?

MNIST

CIFAR



**susceptibility**

# WHAT AFFECTS ROBUSTNESS?

$$1 - U_c \exp(-\pi\epsilon^2)$$

**concentration**

pixels correlated
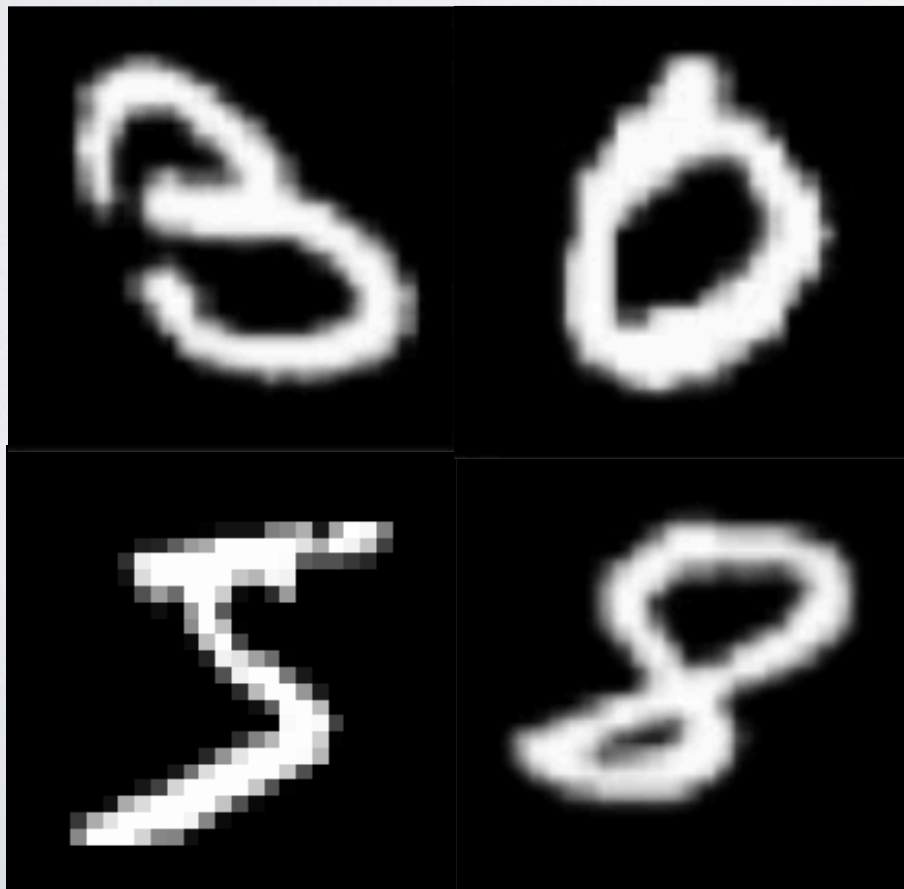low-dimensional

low pixel correlations
high-dimensional

# WHAT AFFECTS THE BOUND?
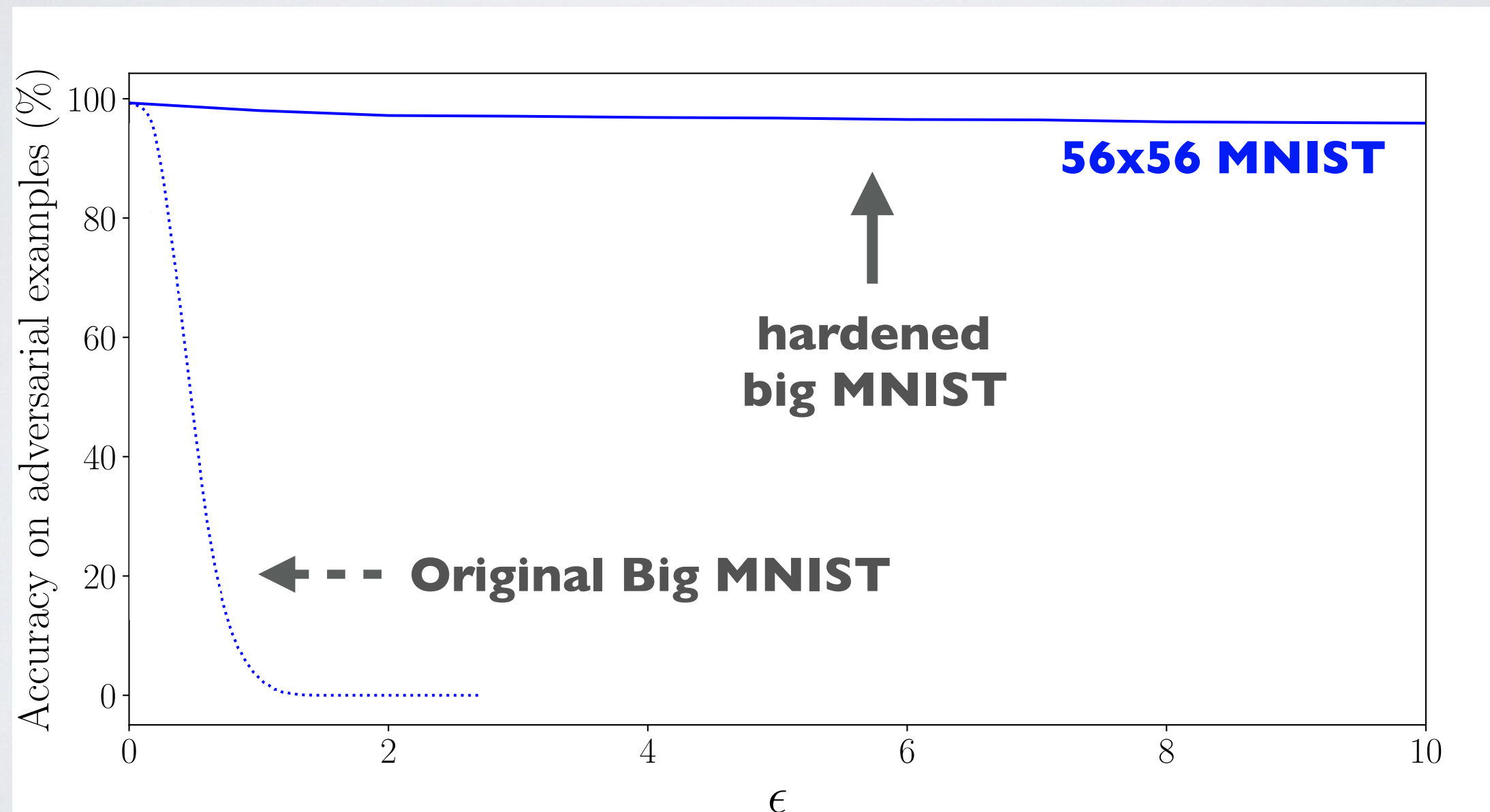
**56x56 MNIST**

3136 features
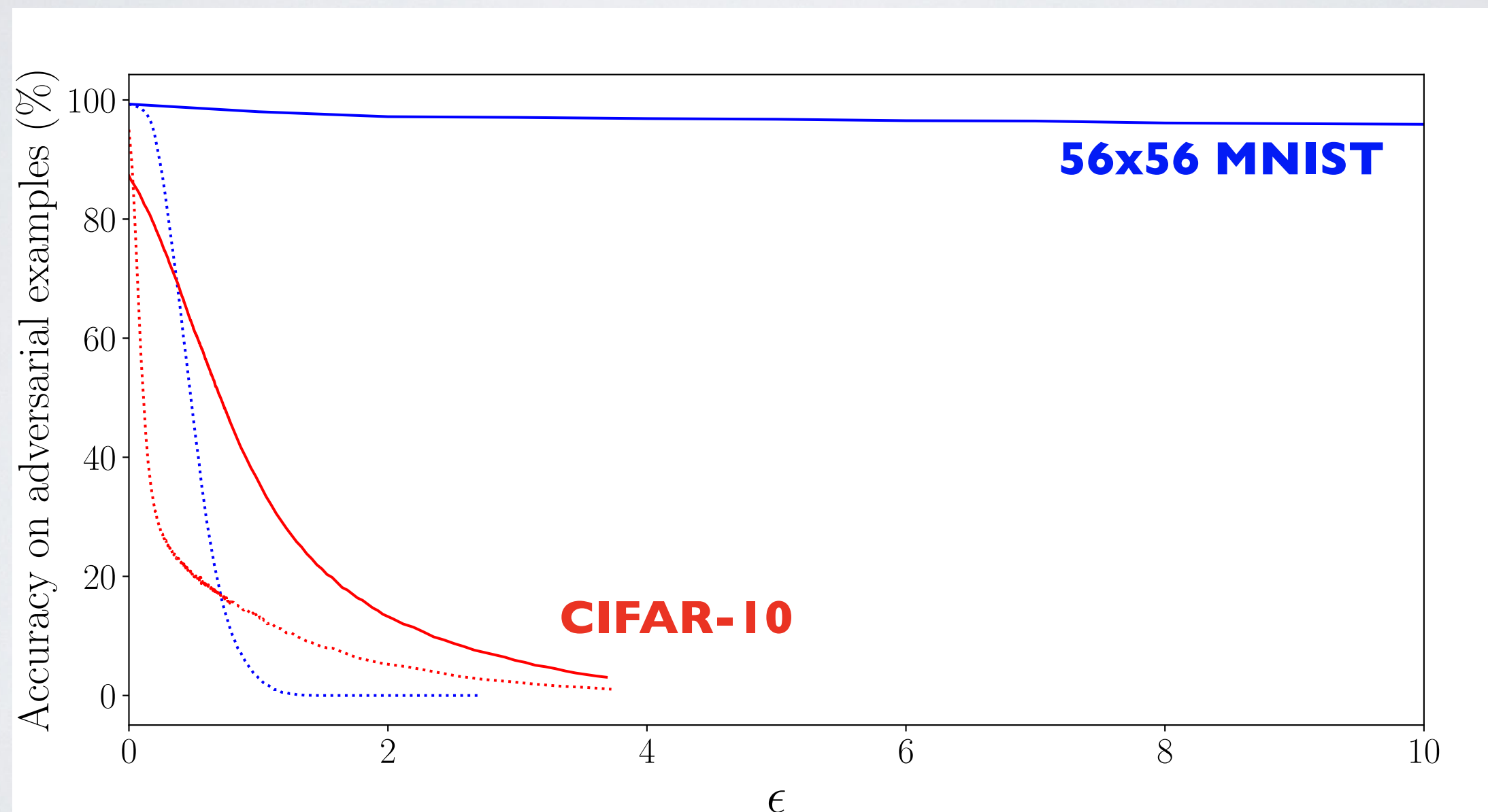
10 classes



**CIFAR-10**

3072 features

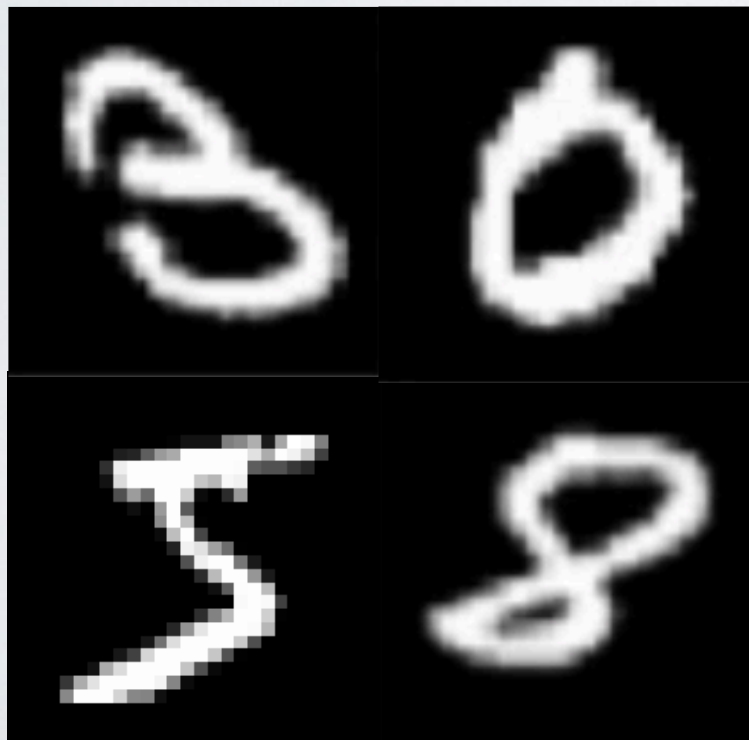10 classes

# ADVERSARIAL TRAINING

# ADVERSARIAL TRAINING

# IMAGE COMPLEXITY LOWERS ROBUSTNESS

$$1 - U_c \exp(-\pi\epsilon^2)$$

**"Complex" image classes have low density**

lower pixel correlations
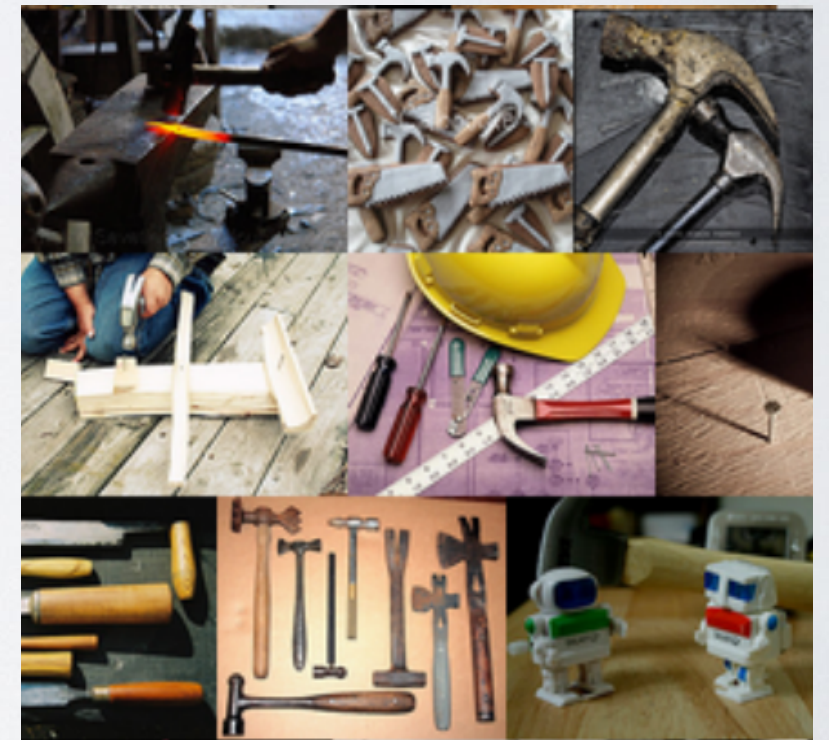higher-dimensional manifolds

MNIST

CIFAR

ImageNet



complexity

# TAKEAWAYS

Robustness has *fundamental* limits

Not specific to neural nets

Can't escape by being clever

**Robustness limit for neural nets might be far worse than intuition tells us!**