# Lecture 21: Adversarial Networks

## CS109B Data Science 2
### Pavlos Protopapas and Mark Glickman

# How vulnerable are Neural Networks?
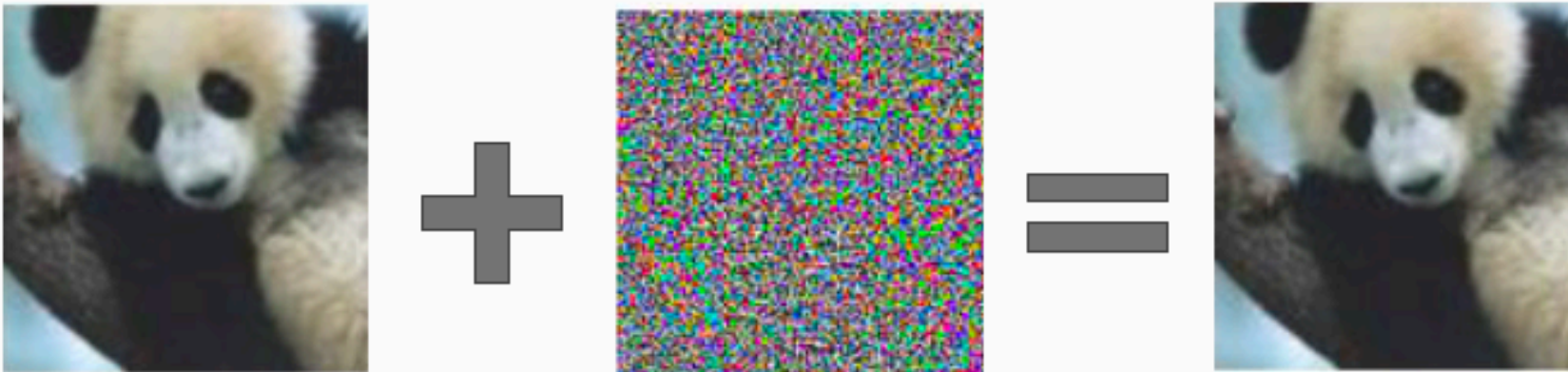
## Uses of Neural Networks

# How vulnerable are Neural Networks?

# Explaining Adversarial Examples

[Goodfellow et. al '15]

1. Robust attacks with FGSM
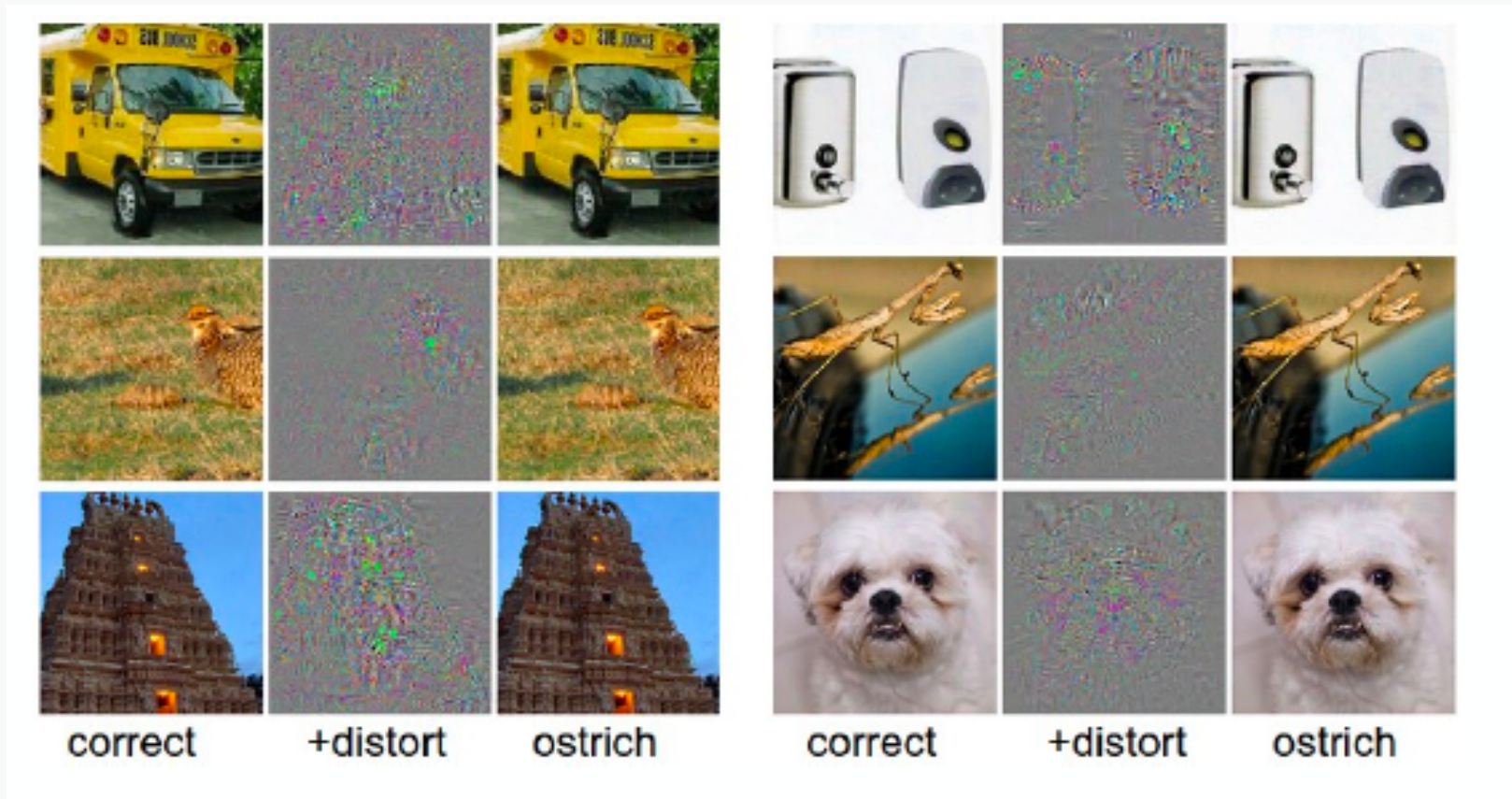2. Robust defense with Adversarial Training



"Panda"
57.7%

Strategic
Noise

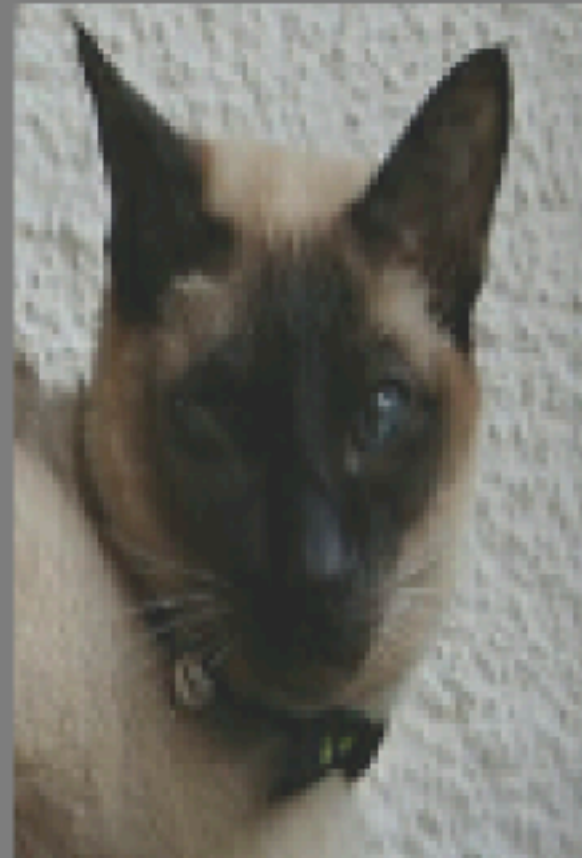# Explaining Adversarial Examples



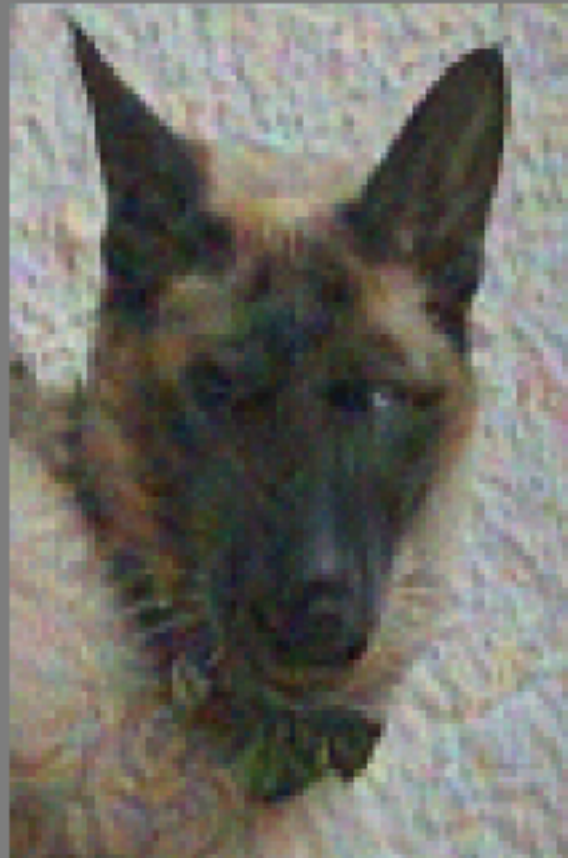correct     +distort     ostrich        correct     +distort     ostrich
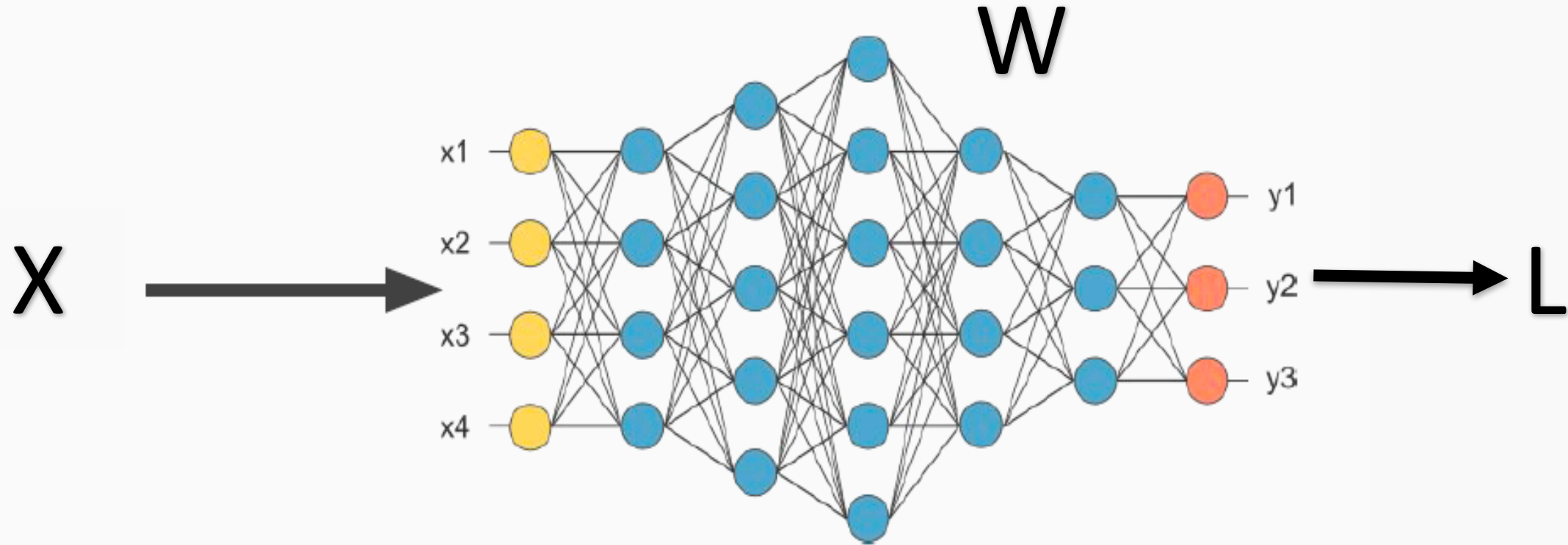
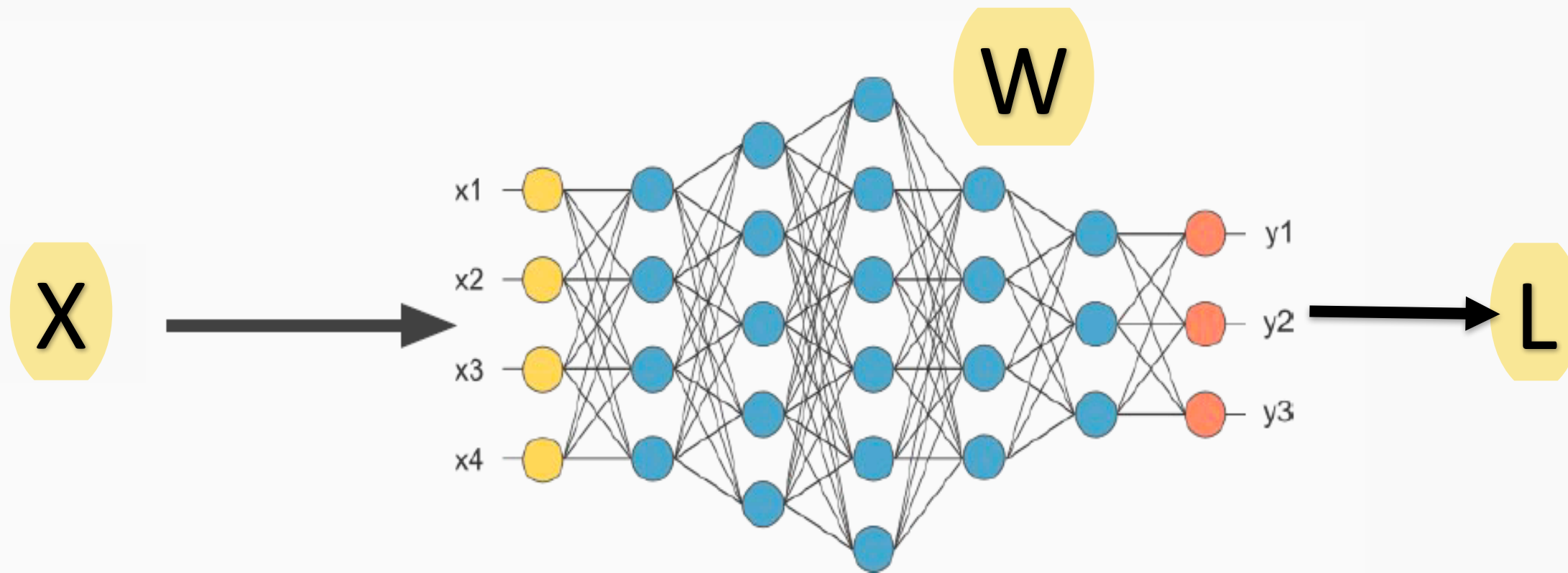# Some of these adversarial examples can even fool humans:

# Attacking with **F**ast **G**radient **S**ign **M**ethod (FGSM)



$$x + \lambda \cdot \text{sign}(\nabla_x L) \Rightarrow x^*$$

# Attacking with **F**ast **G**radient **S**ign **M**ethod (FGSM)



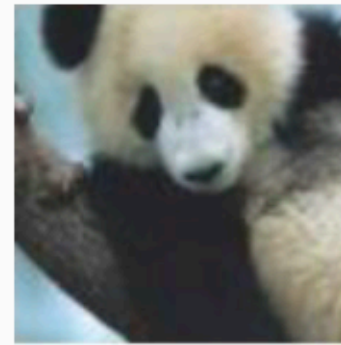$$x + \lambda \cdot \text{sign}(\nabla_{x}L) \Rightarrow x^{*}$$

$$x + \lambda \cdot \text{sign}(\nabla_{\text{x}}\text{L}) \Rightarrow \text{x}^*$$
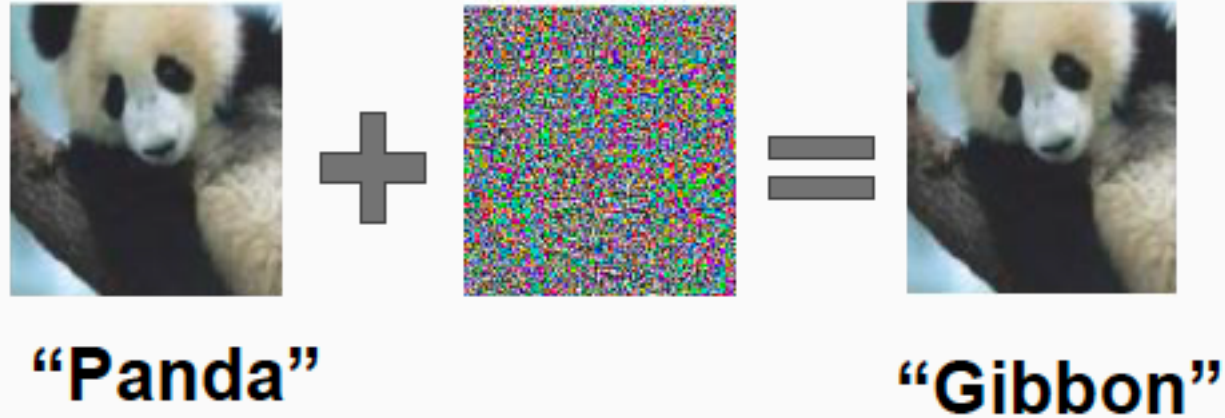
# Defending with Adversarial Training



**"Panda"** **+** **=** **"Gibbon"**

1. Generate adversarial examples
2. Adjust labels

# Defending with Adversarial Training



"Panda" + [noise] = "Gibbon" ➡ "Panda"

1. Generate adversarial examples
2. Adjust labels

# Defending with Adversarial Training



"Panda" + = "~~Gibbon~~" ➡ "**Panda**"

1. Generate adversarial examples
2. Adjust labels
3. Add them to the training set
4. Train new network
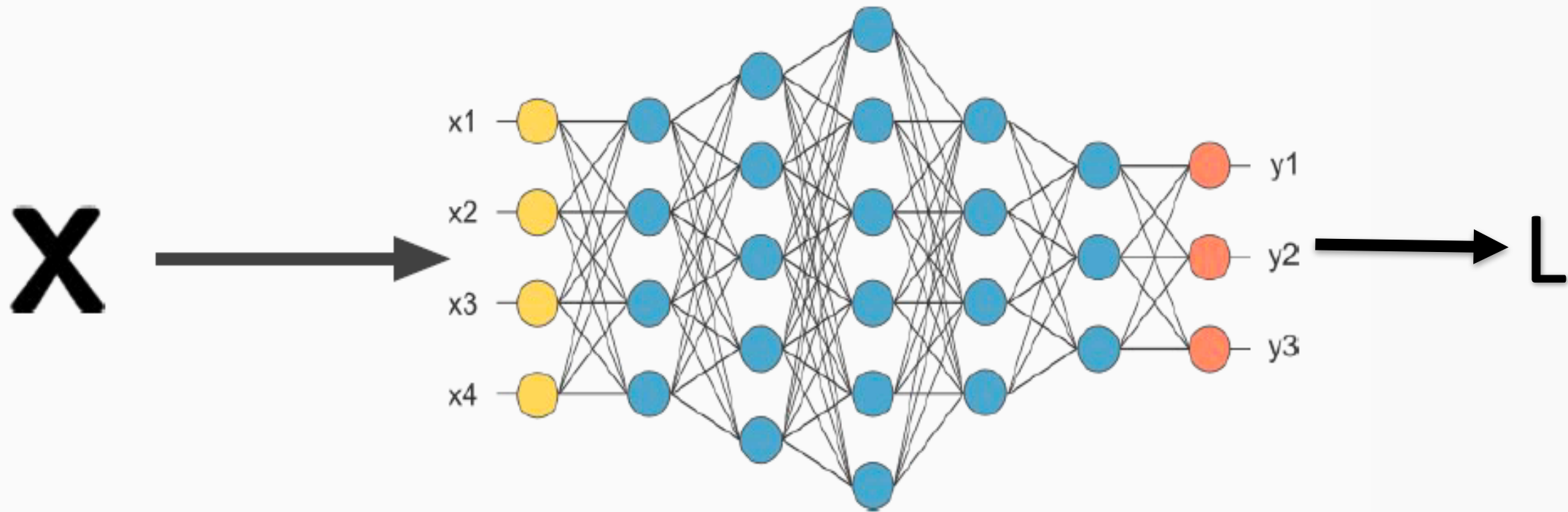
# Attack methods post GoodFellow 2015

- FGSM [Goodfellow et. al '15]

- JSMA [Papernot et. al '16]

- C&W [Carlini + Wagner '16]

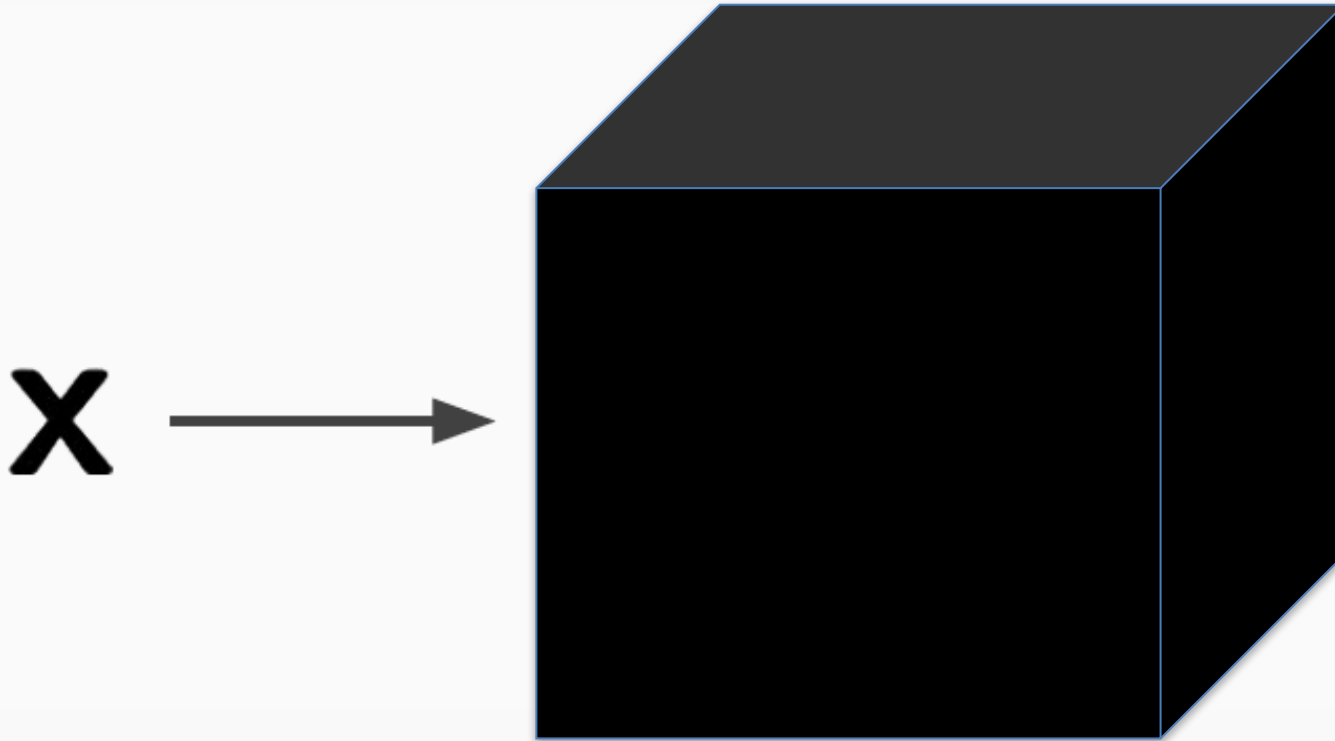- Step-LL [Kurakin et. al '17]

- I-FGSM [Tramer et. al '18]

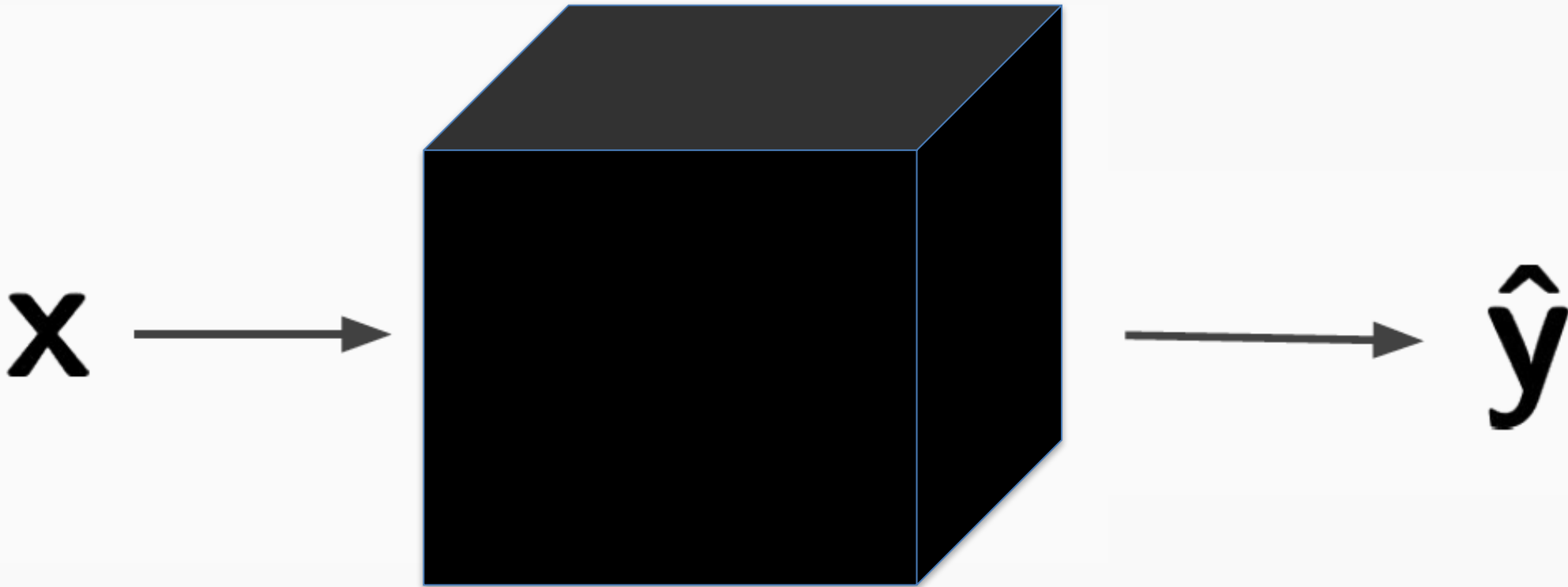$$x + \lambda \cdot \mathrm{sign}(\nabla_x \mathrm{L}) \Rightarrow \mathrm{x}^*$$

"Black Box" Attacks [Papernot et. al '17]

# "Black Box" Attacks

Examine inputs and outputs of the model

$$\mathbf{x} \longrightarrow \qquad \longrightarrow \mathbf{\hat{y}}$$

# "Black Box" Attacks



**Panda**

x → ⬛ → ŷ

**Panda**
**Gibbon**

$x$

$\hat{y}$

# "Black Box" Attacks



X

Panda
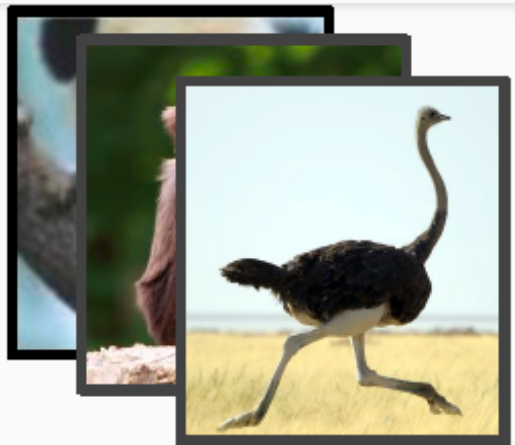Gibbon
Ostrich

$\hat{y}$

# "Black Box" Attacks
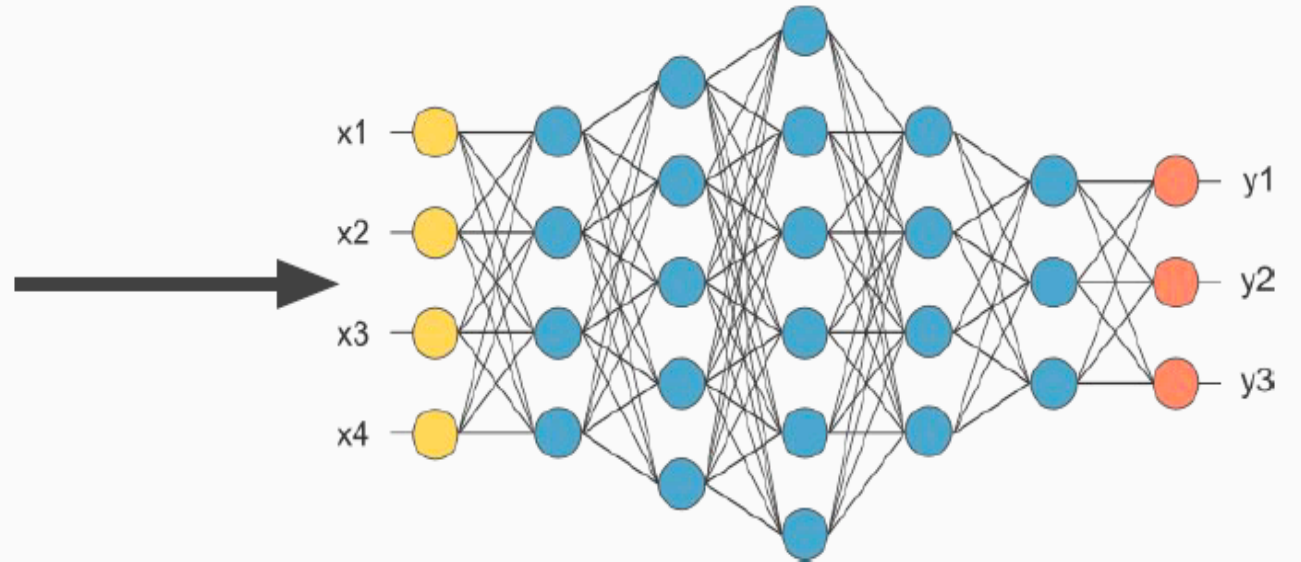
Train a model that performs the same as the black box
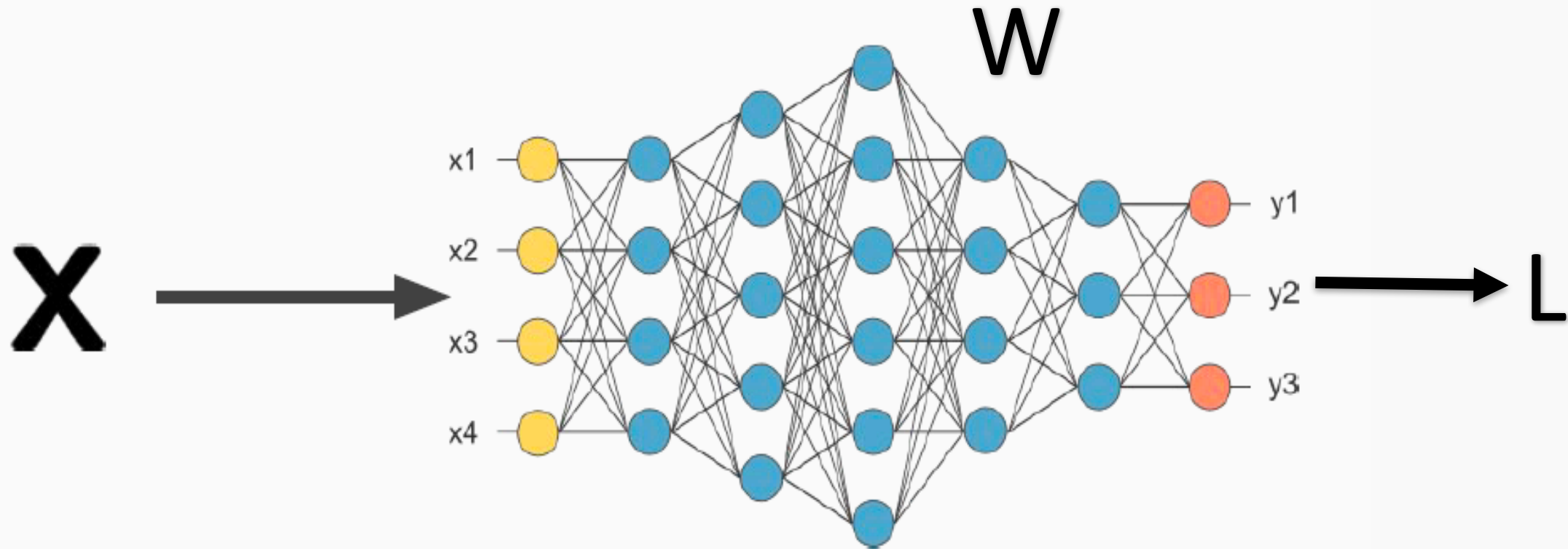
Train a model that performs the same as the black box

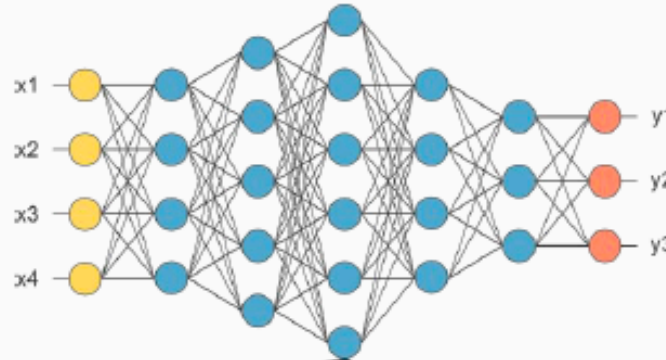

**Panda**

**Gibbon**

**Ostrich**

# "Black Box" Attacks

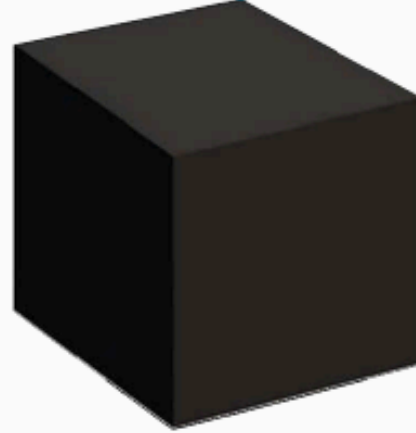Now attack the model you just trained with "white" box attack

**W**

**X** $\longrightarrow$

x1
x2
x3
x4

y1
y2
y3

$\longrightarrow$ **L**

$$x + \lambda \cdot \text{sign}(\nabla_x L) \Rightarrow \dot{x}^*$$

# "Black Box" Attacks
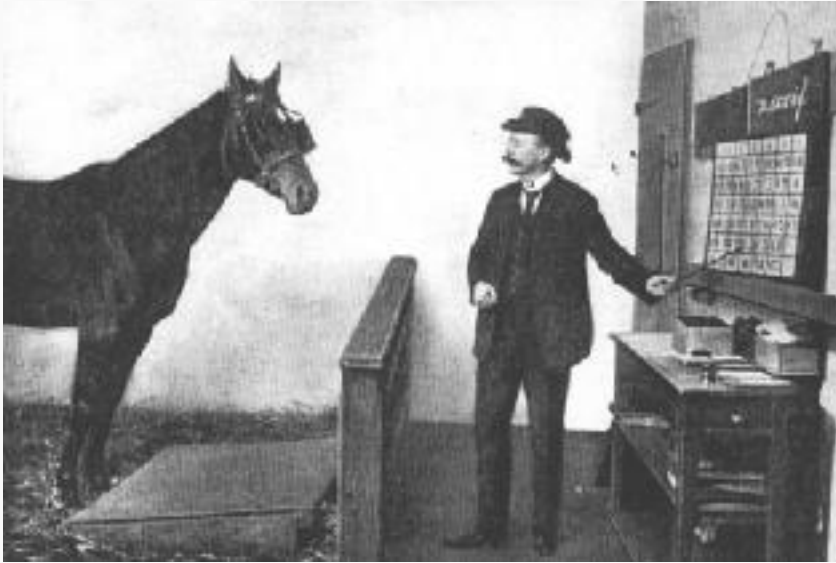
Use those adversarial examples to the "black" box



*"Gibbon"*

*"Gibbon"*

# CleverHans



A Python library to benchmark machine learning systems' vulnerability to adversarial examples.

https://github.com/tensorflow/cleverhans
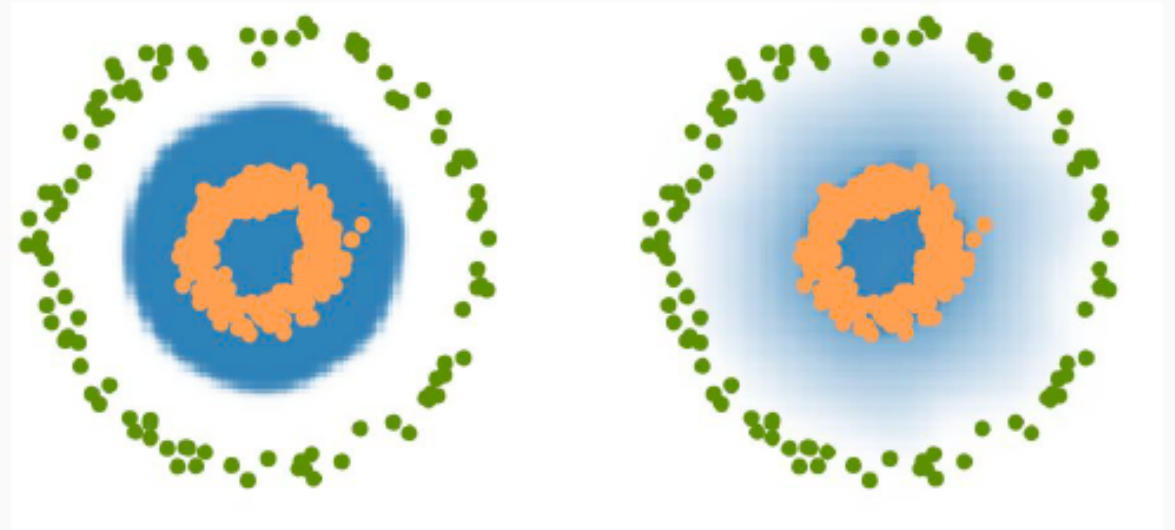http://www.cleverhans.io/

**Mixup:**

- Mix two training examples

- Augment training set

$$\tilde{x} = \lambda x_i + (1 - \lambda)x_j$$
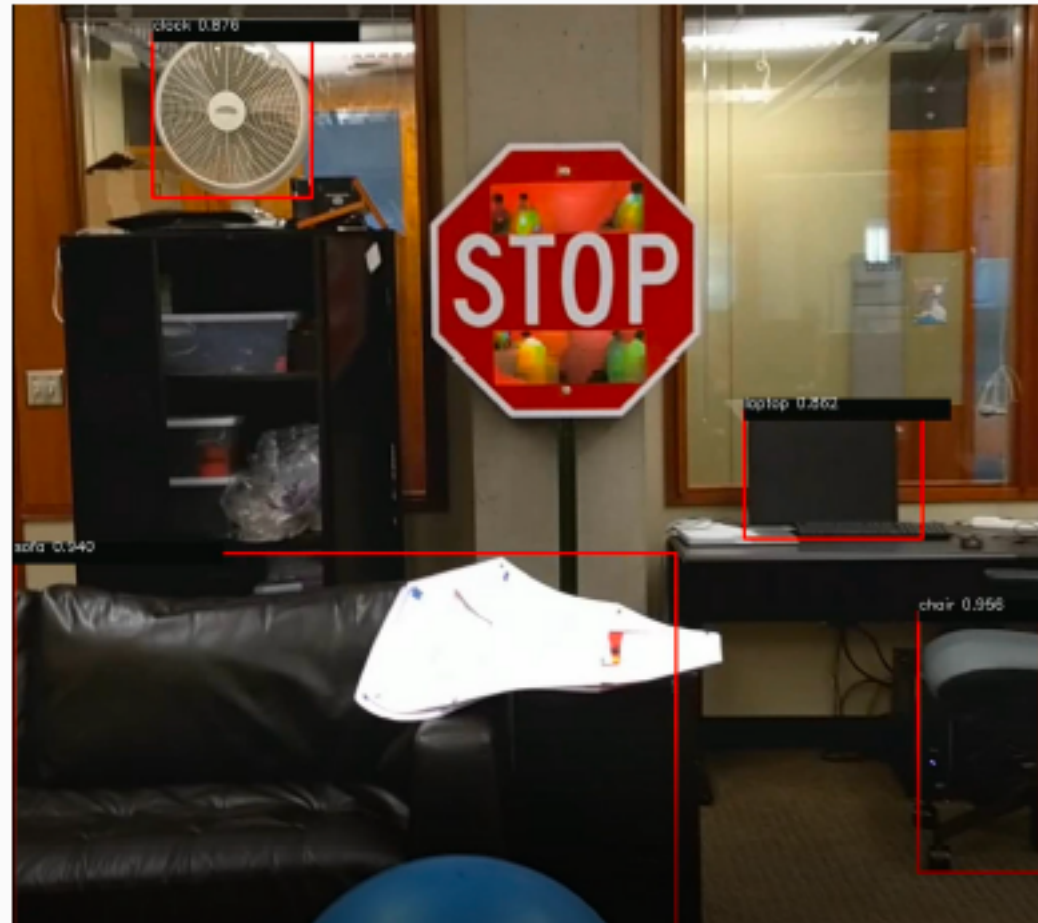$$\tilde{y} = \lambda y_i + (1 - \lambda)y_j$$

**Smooth decision boundaries:**

- Regularize the derivatives wrt to x

# Physical attacks

- Object Detection
- Adversarial Stickers

# Thank you.