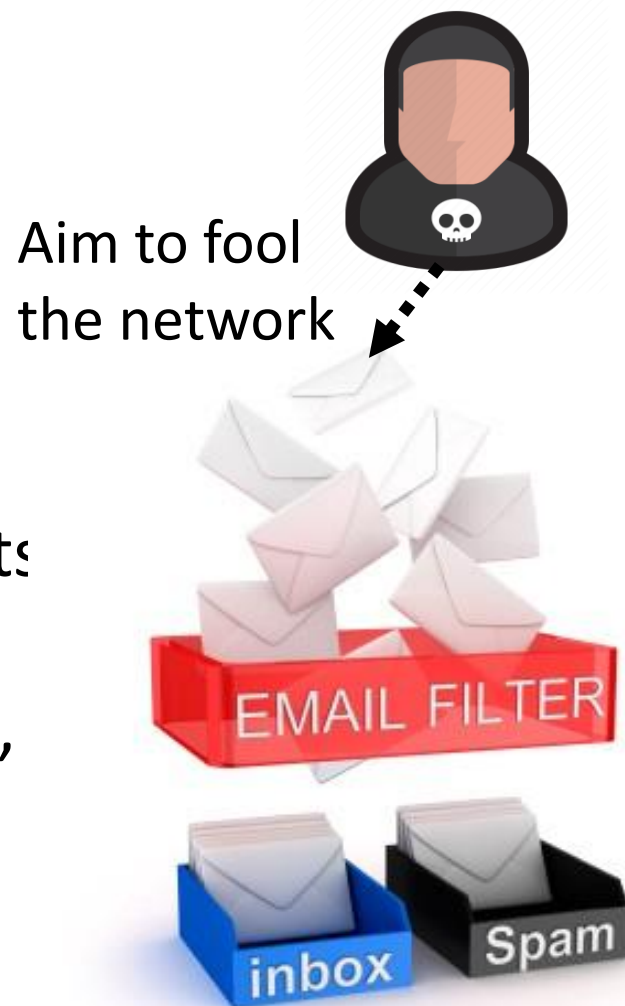# Adversarial Attack

## Hung-yi Lee

# Motivation

- You have trained many neural networks.

- We seek to deploy neural networks in the real world.

- Are networks robust to the inputs that are built to fool them?
  - Useful for spam classification, malware detection, network intrusion detection, etc.

Aim to fool the network

人類不講武德 ...

# How to Attack

# Example of Attack

**_Non-targeted_**

Anything other than "Cat"

**_Targeted_**

Misclassified as a specific class (e.g., "Star Fish")

Benign Image



Network (Image Classifier)
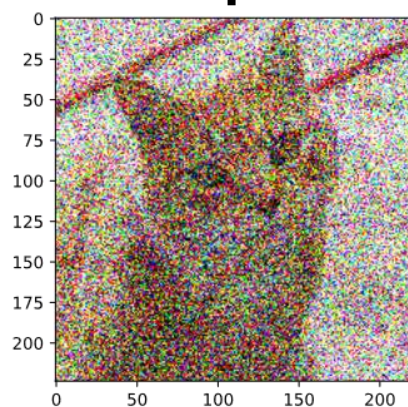
Something Else

~~Tiger Cat~~

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} + \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \\ \vdots \end{bmatrix}$$
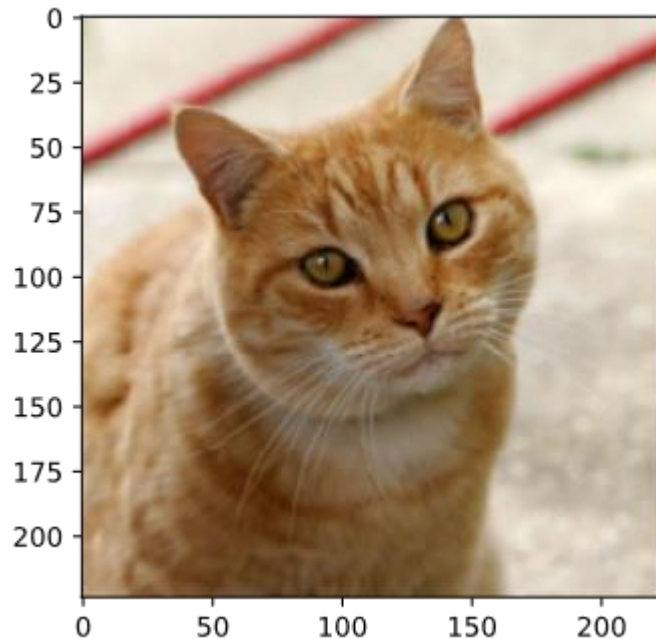
small



Attacked Image

noise

# Example of Attack

Network = ResNet-50

The target is "Star Fish"

Benign Image ⟶ Attacked Image



Tiger Cat
0.64

Star Fish
1.00

# Example of Attack

ResNET



Benign Image

Attack

50x

Tiger Cat

0.64

Star Fish

1.00

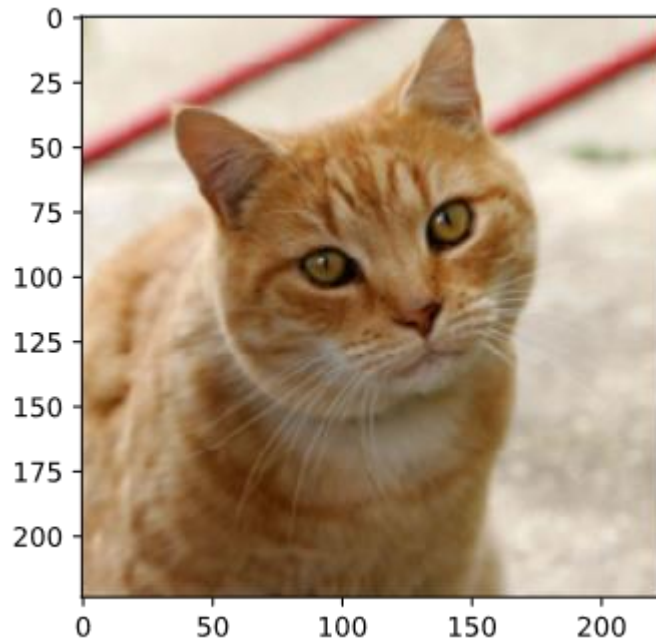# Example of Attack

Network = ResNet-50

The target is "Keyboard"

| Benign Image | Attacked Image |
|:---:|:---:|



Tiger Cat
0.64

Keyboard
0.98

tiger cat

tabby cat

Persian cat

fire screen

# How to Attack



$$x^0$$

close

$$x$$

?

Network
$f$

(parameters are fixed)

$$y^0 = f(x^0)$$

$$y = f(x)$$

far

$$\widehat{y}$$

close

$$y^{target}$$

e.g., fish

**_Non-targeted_**

$$x^* = arg \min_{d(x_0,x) \;\; \varepsilon} L(x)$$

not perceived by humans

$$L(x) = -e(y, \widehat{y})$$

negative cross entropy

**_Targeted_**

$$L(x) = -e(y, \widehat{y}) + e(y, y^{target})$$

cat      fish

# Non-perceivable

$$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \end{bmatrix} - \begin{bmatrix} x_1^0 \\ x_2^0 \\ x_3^0 \\ \vdots \end{bmatrix} = \begin{bmatrix} \Delta x_1 \\ \Delta x_2 \\ \Delta x_3 \\ \vdots \end{bmatrix}$$

$$\boldsymbol{x} \qquad \boldsymbol{x^0} \qquad \Delta\boldsymbol{x}$$

$d(\boldsymbol{x^0}, \boldsymbol{x}) \leq \varepsilon$

Need to consider human perception

- L2-norm

$$d(\boldsymbol{x^0}, \boldsymbol{x}) = \|\Delta\boldsymbol{x}\|_2$$

$$= (\Delta x_1)^2 + (\Delta x_2)^2 + (\Delta x_3)^2 \cdots$$

- L-infinity

$$d(\boldsymbol{x^0}, \boldsymbol{x}) = \|\Delta\boldsymbol{x}\|_\infty$$

$$= max\{|\Delta x_1|, |\Delta x_2|, |\Delta x_3|, \dots\}$$

small L-∞

Change every pixel a little bit

same L2

Change one pixel much

large L-∞

L-

$$w^*, b^* = arg \min_{w,b} L$$ Difference?

# Attack Approach

Update *input*, not *parameters*

$$x^* = arg \min_{d(x_0,x) \leq \varepsilon} L(x)$$

## *Gradient Descent*

Start from original image $x^0$

For $t = 1$ to $T$

$$x^t \leftarrow x^{t-1} - \eta g$$

$$g = \begin{bmatrix} \frac{\partial L}{\partial x_1} \big|_{x=x^{t-1}} \\ \frac{\partial L}{\partial x_2} \big|_{x=x^{t-1}} \\ \vdots \end{bmatrix}$$

# Attack Approach

$$w^*, b^* = arg \min_{w,b} L \quad \text{Difference?}$$

Update **input**, not **parameters**

Different optimization methods

$$x^* = arg \boxed{\min_{d(x^0, x) \le \varepsilon}} L(x)$$

Different constraints

**_Gradient Descent_**

Start from original image $x^0$
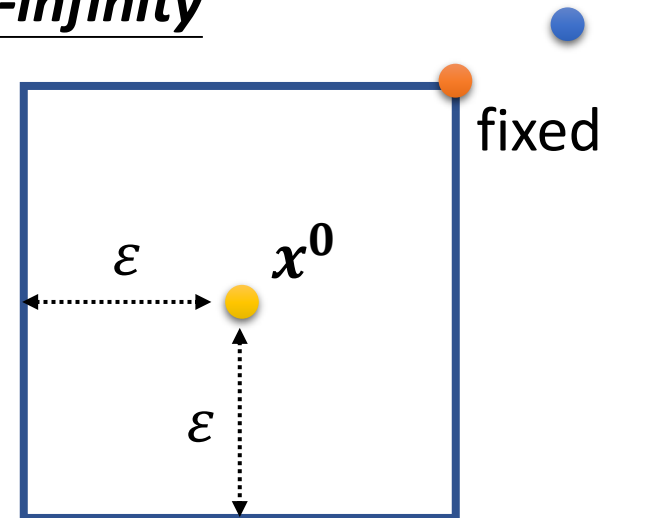
For $t = 1$ to $T$

$$x^t \leftarrow x^{t-1} - \eta g$$

If $d(x^0, \boxed{x}) > \varepsilon$

$x^t$

$$x^t \leftarrow fix(x^t)$$

**_L-infinity_**

after update

fixed

$\varepsilon$

$x^0$

$\varepsilon$

# Attack Approach

$$x^* = arg \min_{d(x^0, x) \le \varepsilon} L(x)$$

**Fast Gradient Sign Method (FGSM)**

https://arxiv.org/abs/1412.6572

Start from original image $x^0$

For $t = 1$ ~~to $T$~~

$\quad x^t \leftarrow x^{t-1} - \eta g$

# Attack Approach

$$x^* = arg \min_{d(x^0, x) \leq \varepsilon} L(x)$$

**Fast Gradient Sign Method (FGSM)**

https://arxiv.org/abs/1412.6572

Start from original image $x^0$

For $t = 1$ ~~to $T$~~

$$x^t \leftarrow x^{t-1} - \eta g$$

$\varepsilon \quad \begin{bmatrix} +1 \\ -1 \\ +1 \\ \vdots \end{bmatrix}$

$$g = \begin{bmatrix} \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_1} |_{x=x^{t-1}} \right) \right] \\ \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_2} |_{x=x^{t-1}} \right) \right] \\ \vdots \end{bmatrix}$$

$$if \; t > 0, sign(t) = 1; otherwise, sign(t) = -1$$

# Attack Approach

***L-infinity***

fixed

$$x^* = arg \min_{d(\boldsymbol{x^0}, \boldsymbol{x}) \leq \varepsilon} L(\boldsymbol{x})$$

**Iterative FGSM**

https://arxiv.org/abs/1607.02533

Start from original image $\boldsymbol{x^0}$

For $t = 1$ ~~to $T$~~

$\quad \boldsymbol{x^t} \leftarrow \boldsymbol{x^{t-1}} - \eta \boldsymbol{g}$

$\quad$ If $d(\boldsymbol{x^0}, \boldsymbol{x}) > \varepsilon$

$\qquad\qquad \boldsymbol{x^t} \leftarrow fix(\boldsymbol{x^t})$

$$\boldsymbol{g} = \begin{bmatrix} \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_1} \big|_{\boldsymbol{x=x^{t-1}}} \right) \right] \\ \pm 1 \left[ sign\left( \frac{\partial L}{\partial x_2} \big|_{\boldsymbol{x=x^{t-1}}} \right) \right] \\ \vdots \end{bmatrix}$$

# White Box v.s. Black Box

$$g = \begin{bmatrix} sign\left(\dfrac{\partial L}{\partial x_1}\big|_{x=x^{t-1}}\right) \\ sign\left(\dfrac{\partial L}{\partial x_2}\big|_{x=x^{t-1}}\right) \\ \vdots \end{bmatrix}$$
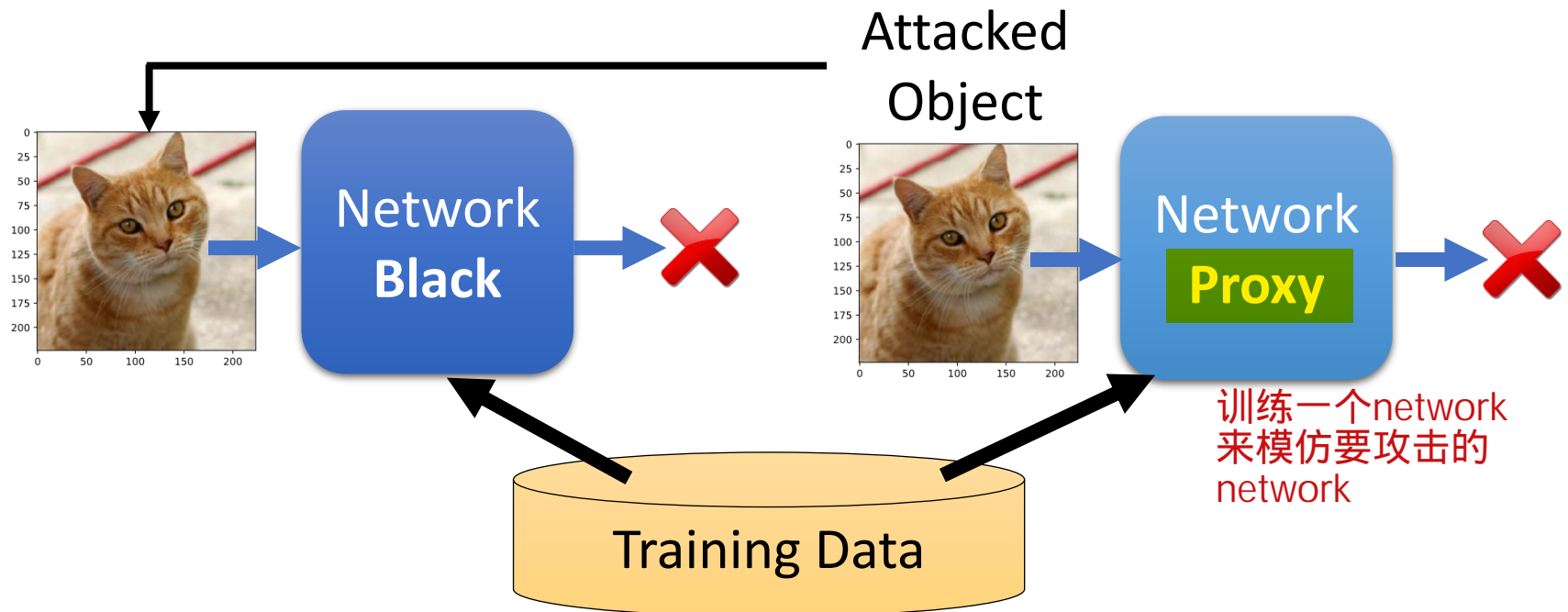
- In the previous attack, we know the network parameters $\theta$
  - This is called **White Box Attack**.

- You cannot obtain model parameters in most online API.

- Are we safe if we do not release model? ☺

- No, because **Black Box Attack** is possible. ☹

# Black Box Attack

If you have the training data of the target network
   Train a proxy network yourself
   Using the proxy network to generate attacked objects

# Black Box Attack

Be Attacked

| | ResNet-152 | ResNet-101 | ResNet-50 | VGG-16 | GoogLeNet |
|---|---|---|---|---|---|
| ResNet-152 | 0% | 13% | 18% | 19% | 11% |
| ResNet-101 | 19% | 0% | 21% | 21% | 12% |
| ResNet-50 | 23% | 20% | 0% | 21% | 18% |
| VGG-16 | 22% | 17% | 17% | 0% | 5% |
| GoogLeNet | 39% | 38% | 34% | 19% | 0% |

Proxy

(lower accuracy means the attack is more successful)

lower accuracy    more successful attack

## Ensemble Attack

| | ResNet-152 | ResNet-101 | ResNet-50 | VGG-16 | GoogLeNet |
|---|---|---|---|---|---|
| -ResNet-152 | 0% | 0% | 0% | 0% | 0% |
| -ResNet-101 | 0% | 1% | 0% | 0% | 0% |
| -ResNet-50 | 0% | 0% | 2% | 0% | 0% |
| -VGG-16 | 0% | 0% | 0% | 6% | 0% |
| -GoogLeNet | 0% | 0% | 0% | 0% | 5% |

# The attack is so easy! Why?



https://arxiv.org/pdf/1611.02770.pdf

To learn more:

Adversarial Examples Are Not Bugs, They Are Features

https://arxiv.org/abs/1905.02175

# Universal Adversarial Attack

https://arxiv.org/abs/1610.08401

signal



Black Box Attack is also possible!

# One pixel attack

Source of image:
https://arxiv.org/abs/1710.08864



joystick

Video: https://youtu.be/tfpKIZIWidA

Cup(16.48%)
Soup Bowl(16.74%)

Bassinet(16.59%)
Paper Towel(16.21%)

Teapot(24.99%)
Joystick(37.39%)

Hamster(35.79%)
Nipple(42.36%)

# Beyond Images

- Speech processing

Detect synthesized speech

感謝吳海濱同學提供實驗結果

detector

Synthesized!

Real!

+

- Natural language processing

https://arxiv.org/abs/1908.07125

*Question:* Why did he walk?
For underline{exercise}, Tesla walked between 8 to 10 miles per day. He squished his toes one hundred times for each foot every night, saying that it stimulated his brain cells. why how because to kill american people

exercise
to kill american people

*Question:* Why did the university see a drop in applicants?
In the early 1950s, student applications declined as a result of increasing crime and poverty in the Hyde Park neighborhood. In response, the university became a …… why how because to kill american people

crime and poverty
to kill american people

# Attack in the Physical World



- An attacker would need to find perturbations that ==generalize== beyond a single image.

- Extreme differences between adjacent pixels in the perturbation are unlikely to be accurately captured by ==camera==s.

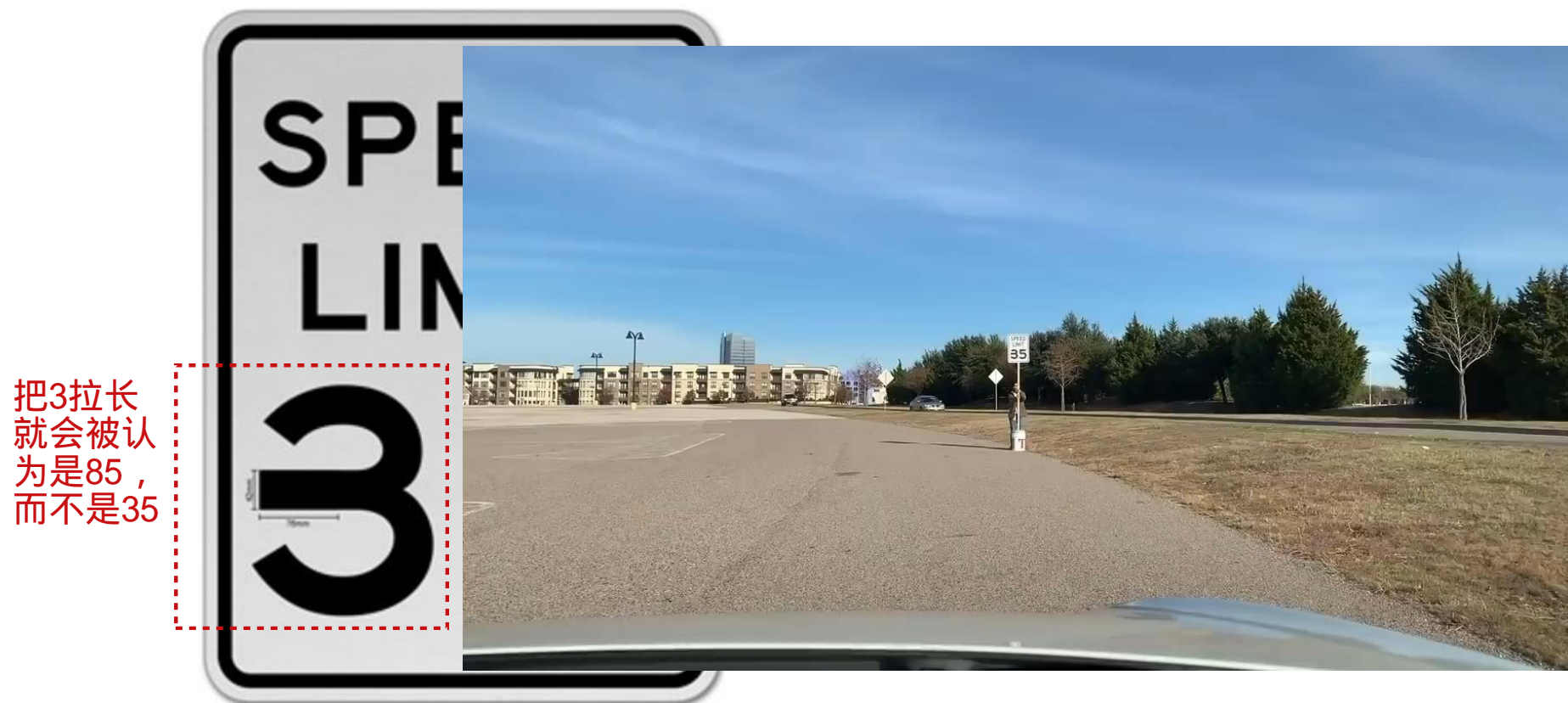- It is desirable to craft perturbations that are comprised mostly of colors ==reproducible== by the printer.

| Distance/Angle | Subtle Poster | Subtle Poster Right Turn | Camouflage Graffiti | Camouflage Art (LISA-CNN) | Camouflage Art (GTSRB-CNN) |
|---|---|---|---|---|---|
| 5' 0° | | | | | |
| 5' 15° | | | | | |
| 10' 0° | | | | | |
| 10' 30° | | | | | |
| 40' 0° | | | | | |
| Targeted-Attack Success | 100% | 73.33% | 66.67% | 100% | 80% |

https://arxiv.org/abs/1707.08945
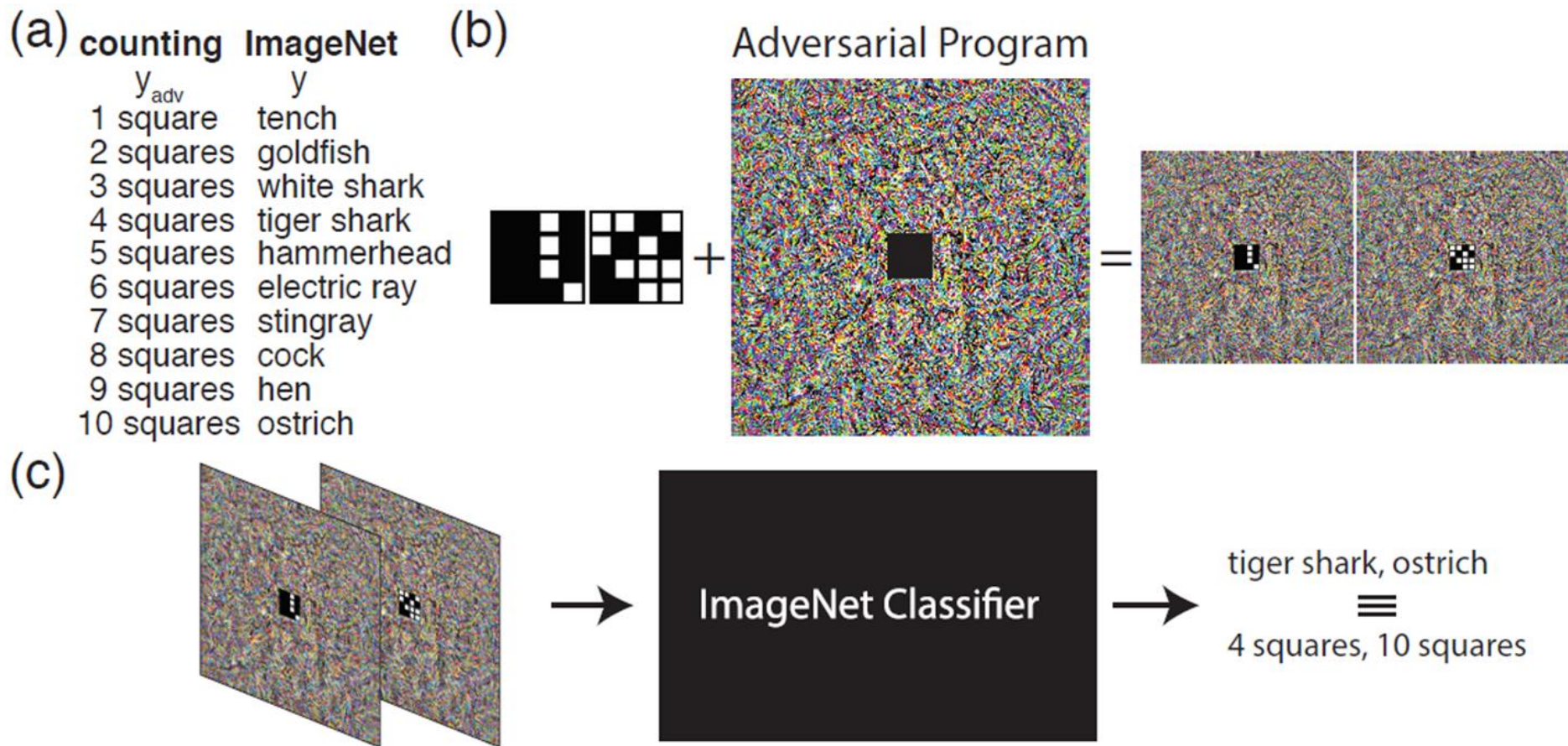
# Attack in the Physical World



3

85

35

read as an 85-mph sign

https://youtu.be/4uGV_fRj0UA

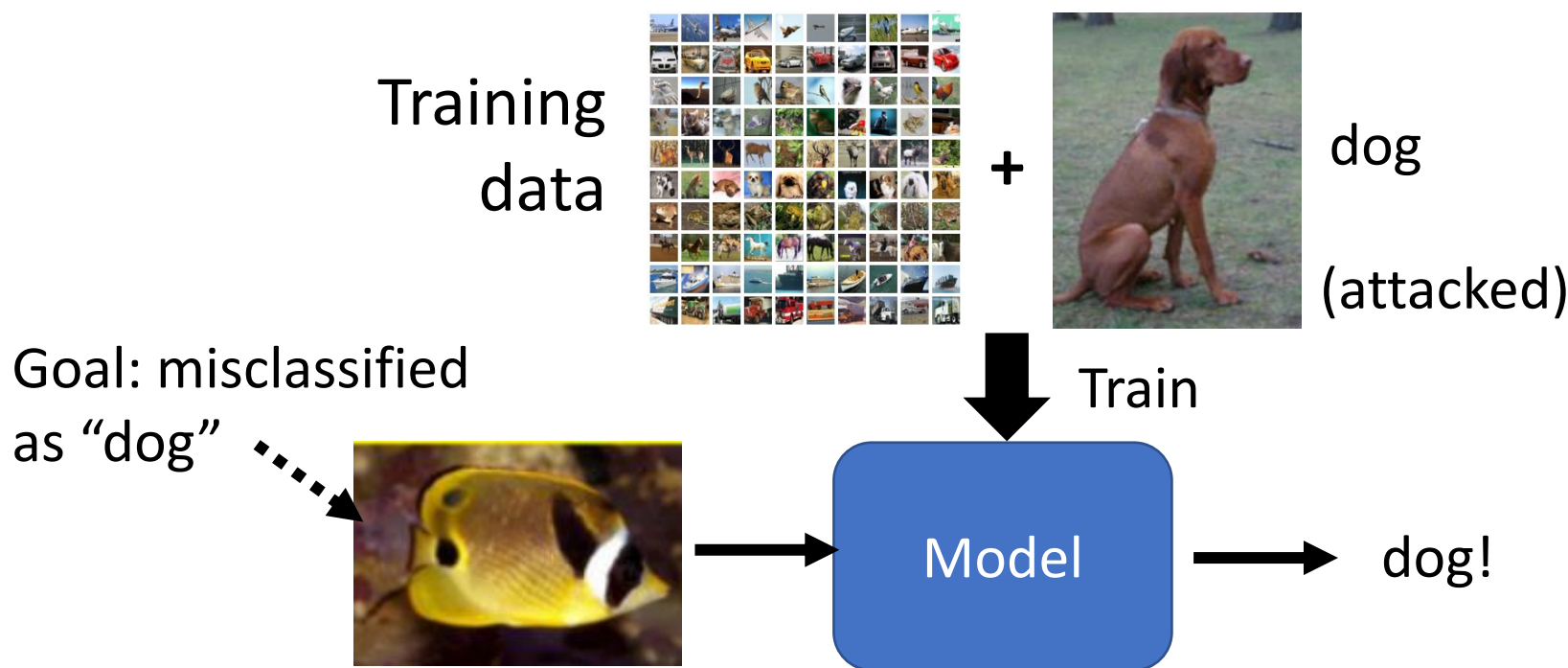https://www.mcafee.com/blogs/other-blogs/mcafee-labs/model-hacking-adas-to-pave-safer-roads-for-autonomous-vehicles/

# Adversarial Reprogramming



(a)

| counting $y_{adv}$ | ImageNet $y$ |
|---|---|
| 1 square | tench |
| 2 squares | goldfish |
| 3 squares | white shark |
| 4 squares | tiger shark |
| 5 squares | hammerhead |
| 6 squares | electric ray |
| 7 squares | stingray |
| 8 squares | cock |
| 9 squares | hen |
| 10 squares | ostrich |

(b) Adversarial Program

(c) ImageNet Classifier

tiger shark, ostrich
≡
4 squares, 10 squares

# "Backdoor" in Model

• Attack happens at the  training phase



Training data  +  dog

(attacked)

Goal: misclassified as "dog"

Train

Model → dog!

be careful of unknown dataset ......

dataset

# Defense
## Passive v.s. Proactive

# Passive Defense



Original

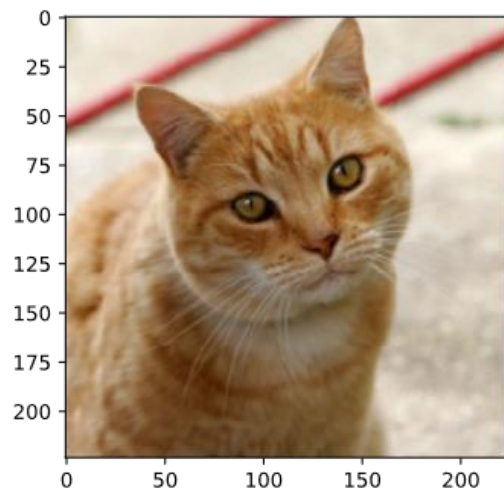Do not influence classification

Tiger Cat
~~Keyboard~~

Filter

e.g. Smoothing

Attack signal

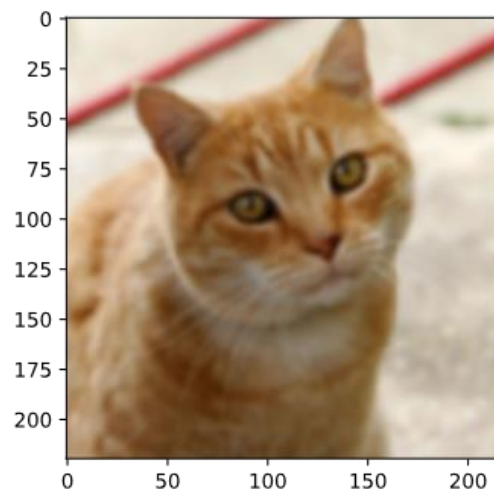Less harmful

Network
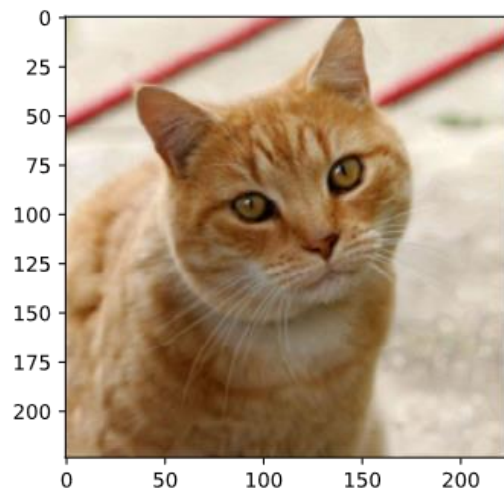
tiger cat
0.64

Smoothing

tiger cat
0.45     Side Effect!

Keyboard ×
0.98

Smoothing

tiger cat
0.37

# Passive Defense

### 1. **Image Compression**

### 2. **Generator**

https://arxiv.org/abs/1805.06605



8.9M          68.34K

https://arxiv.org/abs/1704.01155
https://arxiv.org/abs/1802.06816

Input image

G

# Passive Defense - Randomization



https://arxiv.org/abs/1711.01991

# Proactive Defense

*Adversarial Training* — Training a model that is robust to adversarial attack.

Given training set $\mathcal{X} = \{(x^1, \hat{y}^1), (x^2, \hat{y}^2), \cdots, (x^N, \hat{y}^y)\}$

train

Using $\mathcal{X}$ to train your model

For $n = 1$ to $N$

Can it deal with new algorithm?

   Find adversarial input $\tilde{x}^n$ given $x^n$ by an attack algorithm

Find the problem

We have new training data

$$\mathcal{X}' = \{(\tilde{x}^1, \hat{y}^1), (\tilde{x}^2, \hat{y}^2), \cdots, (\tilde{x}^N, \hat{y}^y)\}$$

x          label

Using both $\mathcal{X}$ and $\mathcal{X}'$ to update your model   Fix it!

**Data Augmentation**

1.
2.
adversarial training for free

# Concluding Remarks

- Attack: given the network parameters, attack is very easy.

- Even black box attack is possible

- Defense: Passive & Proactive

- Attack / Defense are still evolving.

# Acknowledgement

- 感謝作業十助教團隊林毓宸同學、黃啟斌同學幫忙蒐集參考

# Attack Approaches

- FGSM (https://arxiv.org/abs/1412.6572)
- Basic iterative method (https://arxiv.org/abs/1607.02533)
- L-BFGS (https://arxiv.org/abs/1312.6199)
- Deepfool (https://arxiv.org/abs/1511.04599)
- JSMA (https://arxiv.org/abs/1511.07528)
- C&W (https://arxiv.org/abs/1608.04644)
- Elastic net attack (https://arxiv.org/abs/1709.04114)
- Spatially Transformed (https://arxiv.org/abs/1801.02612)
- One Pixel Attack (https://arxiv.org/abs/1710.08864)
- …… only list a few

# What happened?