# Robust Physical-World Attacks on Deep Learning Visual Classification

Kevin Eykholt, Ivan Evtimov, Earlence Fernandes, Bo Li, and Dawn Song
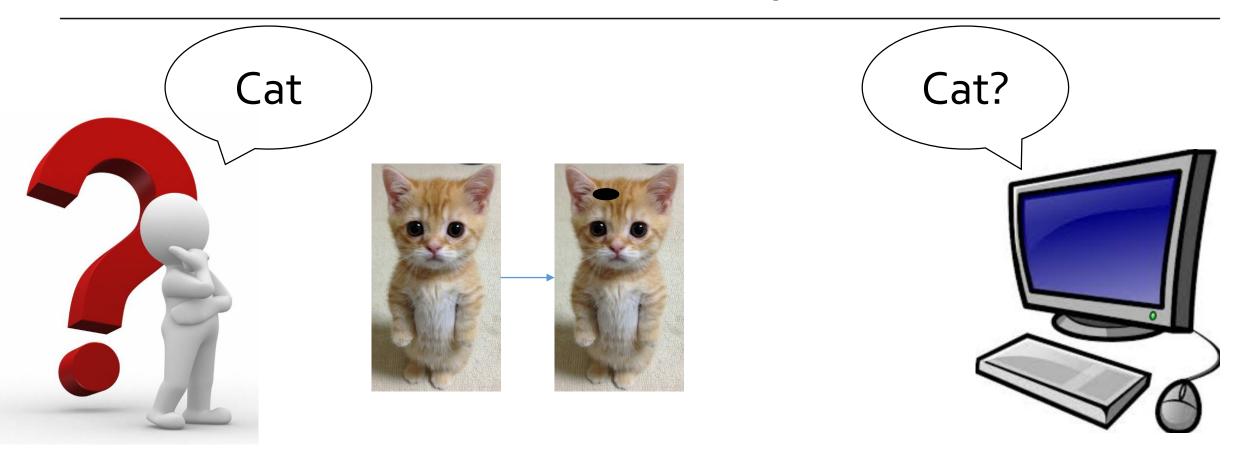
SYSSEC ChoManGi
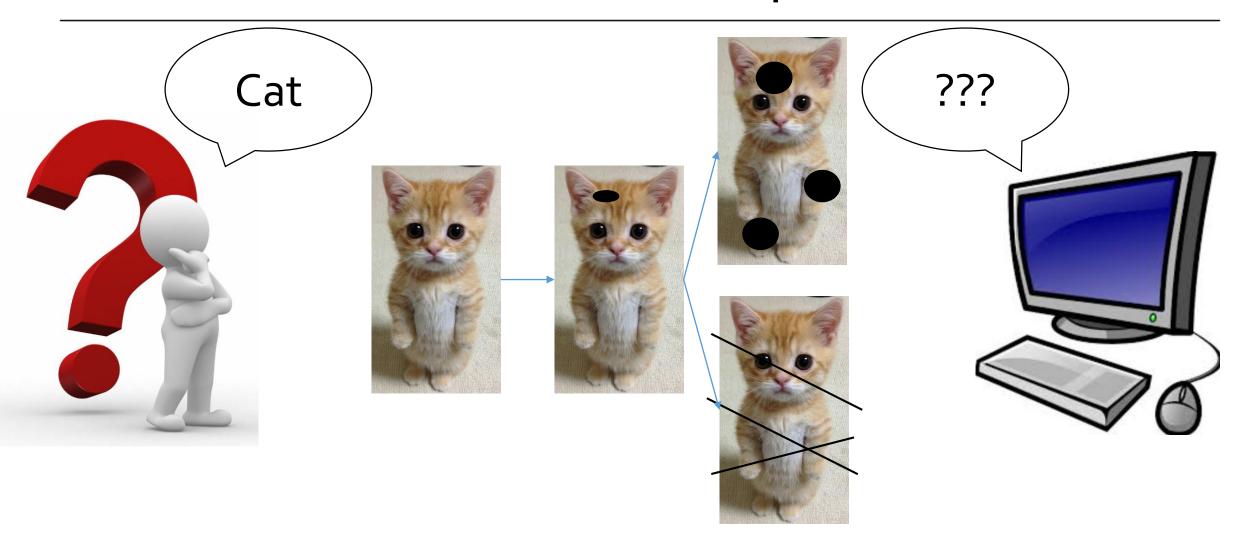
# What is Adversarial Example?

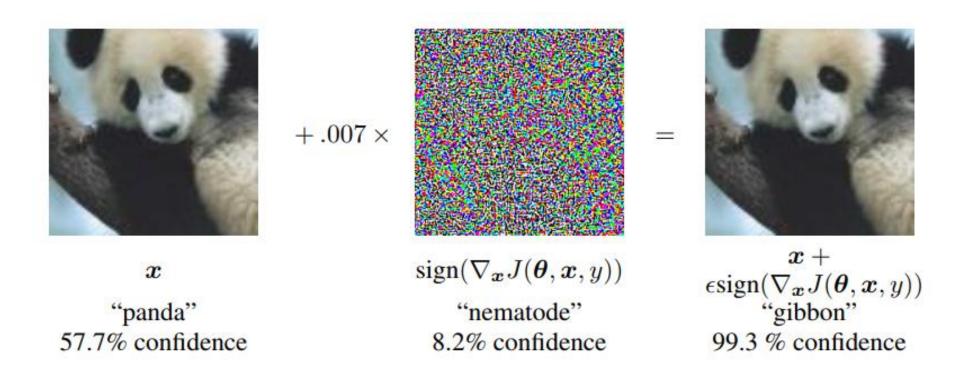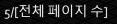# What is Adversarial Example?

# What is Adversarial Example?

# What is adversarial example?

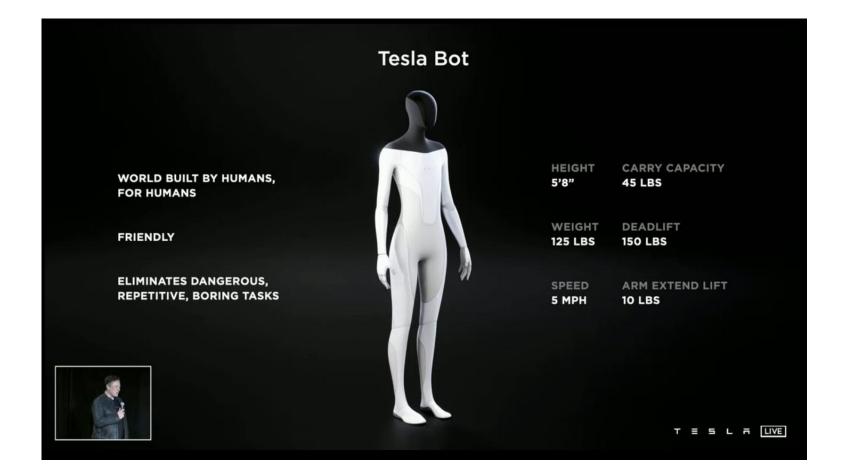❖ **Samples that cause a machine learning model to make a false prediction by using the difference between human and AI**

– It uses the different classification standards between human and AI
– Human and machines make different decision about adversarial samples.



$$+.007 \times$$

$$=$$

$x$

"panda"
57.7% confidence

$\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$

"nematode"
8.2% confidence

$\boldsymbol{x} +$
$\epsilon\text{sign}(\nabla_x J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence

**SysSec**
System Security Lab

# Importance of A.E.?




Original: Temple


Adversarial: Ostrich


Original: Duck


Adversarial: Horse


Original:
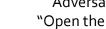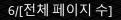"How are you?"


Adversarial:
"Open the door"

# Principles of adversarial example

❖target AI: Classifier
- $Y = F(X)$ is model process. X: model input, Y: classification result

❖Basic idea
- Add minimal perturbation to input x so that it can be determined that it is not X.
- When $F(X) = Y$, solving optimization problem that
$$argmin_{\|\delta X\|} \ s.t. \ F(X + \delta X) = Y^* \neq Y$$

❖Category
- Untargeted: Classify X as not Y.
  - $minimize_{\delta X} \ J(X + \delta X, Y^* \neq Y), J$: Cost function
- Targeted: Classify X as targeted class $Y^*$
  - $minimize_{\delta X} \ J(X + \delta X, Y^* = Y_{target} \neq Y)$

# Previous Works on Digital images

❖ Fast/Iterative Gradient Sign/Value Method (FGSM/FGVM, IGSM)
- – First suggested adversarial example method – Goodfellow et al. 2014.
- – FGSM: $X_{adv} = X + \epsilon \cdot \text{sign}(\nabla_X J(X, y_{true}))$, FGVM: $X_{adv} = X + \epsilon \cdot \nabla_X J(X, y_{true})$
- – IGSM: 
$$X_{adv,0} = X,$$
$$X_{adv,N+1} = Clip_{X,\epsilon}\{X_{adv,N} + \alpha \cdot \text{sign}(\nabla_X J(X_N, y_{true}))\}.$$
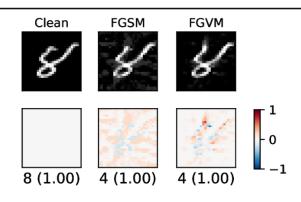
❖ DeepFool - 2015
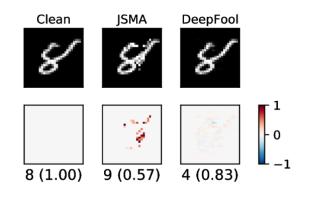- – Calculate a minimum of L2 perturbation through model structure approximation

❖ Jacobian-based Saliency Map Attack (JSMA) - 2016
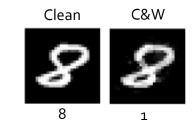- – Modulates only pixels that have the greatest impact on classification results
- ➔ Minimize the number of pixels modulated(minimize $L_0$ distance)

❖ Carlini & Wagner (C&W) - 2017
- – Defending the adversarial defense method, Defensive Distillation.
- – Considered to be one of the most powerful attack method
- – Attacks on $L_0, L_2, L_\infty$ metrics are possible

# Why Physical World?

❖ After the Carlini's paper, adversarial example research on digital images is not active. – Why?

❖ Previous techniques for digital images presented were all to add calculated minimum noise to the whole picture, including the background.

❖ Verification about practicality is essential to actual use of the proposed techniques.

# Why Physical World?

# Possible??

# Why Physical World?


Bright?


Angle?


Distance?

# Purpose

❖ Propose an white-box adversarial example attack on road signs in physical world conditions.

❖ Why they focused on road signs?
 – Simple!
 – Various angles, distances!
 – Impoartance!
 – Effective!

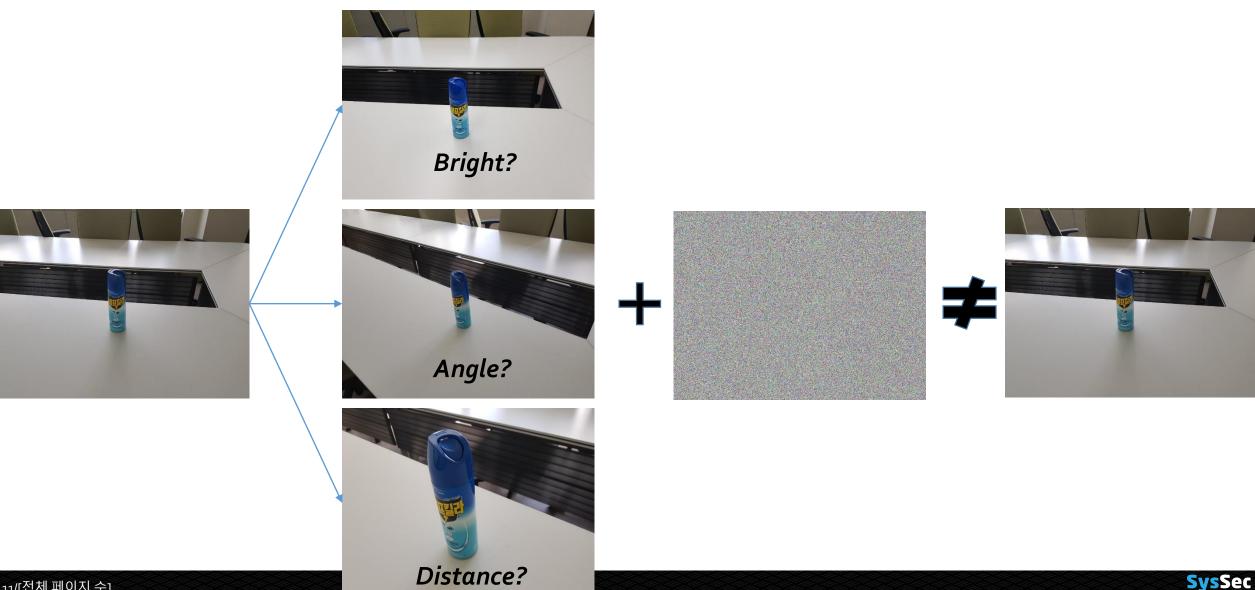# Contribution(Work)

❖ Propose RP2(Robust Physical Perturbation) algorithm to generate physical perturbations that can consistently cause misclassification under various physical conditions.

❖ Evaluate algorithm with famous road sign dataset, LISA and GTSRB.

❖ To show the generality of algorithm, they tested their attacks on Inception-v3 classifier to misclassify the microwave as "phone".

SysSec
System Security Lab

# Attack Pipeline



Model Physical Dynamics by Sampling from Distribution

$f_\theta(x)$

Output

SPEED LIMIT 45

Target

RP₂

Mask

Input

Stationary + Drive-By Testing

Perturbed Stop Sign Under Varying Distances/Angles

# Algorithm - basic

1) Untargeted : $\arg min_\delta \lambda \|\delta\|_p - J(f_\theta(x+\delta), y)$

2) Targeted : $\arg min_\delta \lambda \|\delta\|_p + J(f_\theta(x+\delta), y^*)$

$\delta$ : perturbation, $\lambda$: regularization coefficient, x : input,

y : authentic class, y* : target class, $\|.\|_p$: 2D p-norm($\sum_{i,j}(|\delta_{(i,j)}|^p)^{1/p}$), J : cross entropy, $\theta$ : hyper parameter

# Algorithm – (1) various conditions

Consider various distances, angles, brightness for loss calculations.

1) Untargeted : $\arg min_\delta \lambda \|\delta\|_p - J(f_\theta(x+\delta), y)$

2) Targeted : $\arg min_\delta \lambda \|\delta\|_p + J(f_\theta(x+\delta), y^*)$

1) Untargeted : $\arg min_\delta \lambda \|\delta\|_p - \frac{1}{k}\sum_{i=1}^{k} J(f_\theta(x+\delta), y)$

2) Targeted : $\arg min_\delta \lambda \|\delta\|_p + \frac{1}{k}\sum_{i=1}^{k} J(f_\theta(x+\delta), y^*)$
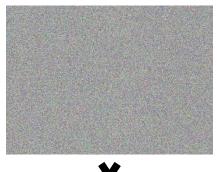
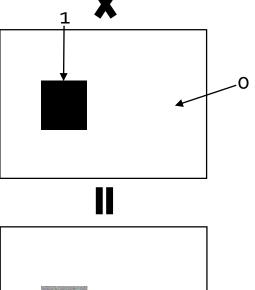### ***Average for distance, angle, brightness!***

# Algorithm – (2) use mask

Use mask matrix to modify specific areas only.

1) Untargeted : $\arg min_\delta \lambda \|\delta\|_p - \frac{1}{k}\sum_{i=1}^{k}J(f_\theta(x+\delta),y)$

2) Targeted : $\arg min_\delta \lambda \|\delta\|_p + \frac{1}{k}\sum_{i=1}^{k}J(f_\theta(x+\delta),y^*)$

3) Untargeted : $\arg min_\delta \lambda \|M_x\delta\|_p - \frac{1}{k}\sum_{i=1}^{k}J(f_\theta(x+M_x\delta),y)$

4) Targeted : $\arg min_\delta \lambda \|M_x\delta\|_p + \frac{1}{k}\sum_{i=1}^{k}J(f_\theta(x+M_x\delta),y^*)$

***Perturbate only matrix area!***

# Algorithm – (3) available colors

NPS(p) = $\prod_{\hat{p} \in P} |p - \hat{p}|$, P = set of printable colors, p = color of each pixels

1) Untargeted : $\arg min_\delta \lambda \|M_x \delta\|_p - \frac{1}{k} \sum_{i=1}^{k} J(f_\theta(x + M_x \delta), y)$

2) Targeted : $\arg min_\delta \lambda \|M_x \delta\|_p + \frac{1}{k} \sum_{i=1}^{k} J(f_\theta(x + M_x \delta), y^*)$

3) Untargeted : $\arg min_\delta \lambda \|M_x \delta\|_p + NPS(M_x \delta) - \frac{1}{k} \sum_{i=1}^{k} J(f_\theta(x + M_x \delta), y)$

4) Targeted : $\arg min_\delta \lambda \|M_x \delta\|_p + NPS(M_x \delta) + \frac{1}{k} \sum_{i=1}^{k} J(f_\theta(x + M_x \delta), y^*)$

## *Don't use non-printable color!*

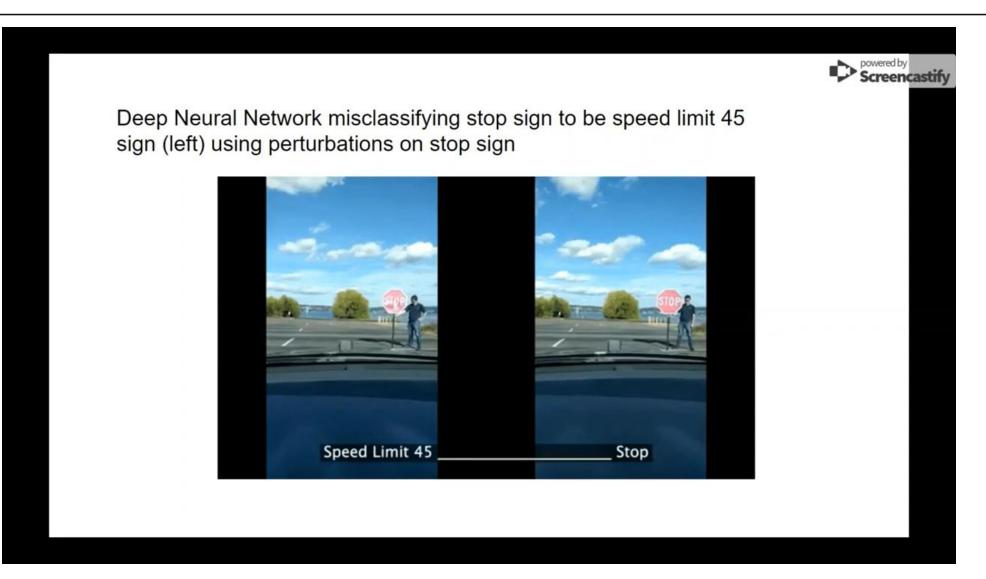# Evaluation – (1) road sign



| 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% |

SysSec
System Security Lab

# Evaluation – (1) road sign



Deep Neural Network misclassifying stop sign to be speed limit 45 sign (left) using perturbations on stop sign

# Evaluation – (1) road sign



| Perturbation | Attack Success | A Subset of Sampled Frames $k = 10$ |
| --- | --- | --- |
| Subtle poster | 100% | |
| Camouflage abstract art | 84.8% | |

# Evaluation – (2) Microwave

| Distance/Angle | Image | Distance/Angle | Image |
|---|---|---|---|
| 2' 0° | | 2' 15° | |
| 5' 0° | | 5' 15° | |
| 7' 0° | | 7' 15° | |
| 10' 0° | | 10' 15° | |
| 15' 0° | | 20' 0° | |



| Distance & Angle | Top Class (Confid.) | Second Class (Confid.) |
|---|---|---|
| 2' 0° | Phone (0.78) | Microwave (0.03) |
| 2' 15° | Phone (0.60) | Microwave (0.11) |
| 5' 0° | Phone (0.71) | Microwave (0.07) |
| 5' 15° | Phone (0.53) | Microwave (0.25) |
| 7' 0° | Phone (0.47) | Microwave (0.26) |
| 7' 15° | Phone (0.59) | Microwave (0.18) |
| 10' 0° | Phone (0.70) | Microwave (0.09) |
| 10' 15° | Phone (0.43) | Microwave (0.28) |
| 15' 0° | Microwave (0.36) | Phone (0.20) |
| 20' 0° | Phone (0.31) | Microwave (0.10) |

SysSec
System Security Lab

# Future Work

# How to Attack Black-box?

## 1. Use Similar White-box



White-box model → Adversarial Example → Black-box model

## *Transferability!*

## 2. Copy the target black-box

Black-box model ← *Query* ← Replica model → Adversarial Example → Black-box model

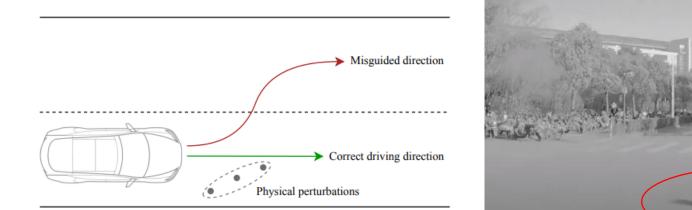Black-box model → *Data* → Replica model

# Follow-up Work

❖ Too Good to Be Safe: Tricking Lane Detection in Autonomous Driving with Crafted Perturbations (Usenix 2021)

- A.E attack on lane detection system of Tesla
- physical perturbations to misguide the direction of Tesla

# Questions to the presenter

◆They consider various distances, angles, brightness for loss calculations. I wonder if this attack is possible regardless of the weather, such as sunny, rainy, cloudy, or snowy days (김한나) – Best question

  – It seems very difficult to consider weather conditions.

◆In this paper, there is no detailed comment about attack mitigations or defenses. What is the general or proper mitigation targeting DNNs? (이용화)

  – Adversarial training can be the general mitigation against A.E.

◆ Are there any advanced ml techniques to learn find unknown vulnerability besides adversarial perturbations? (김경태)

  – The main idea of A.E is to threat the perception differences between human and AI. Just finding the ml vulnerabilities can be found with input mutations(fuzzing).

# END