

## Other Attacks & Applications

Patrick Chan  
patrickchan@ieee.org



## Agenda

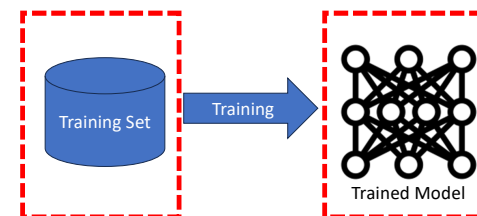
- Privacy Attack
- Physical Attack
- Non-Security Application
- Conclusion

## Privacy Attack



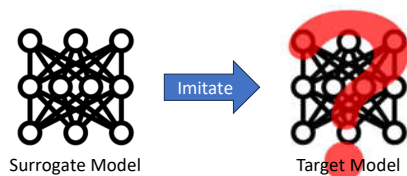
## Objective

- Model Stealing
- Training Set Recovery



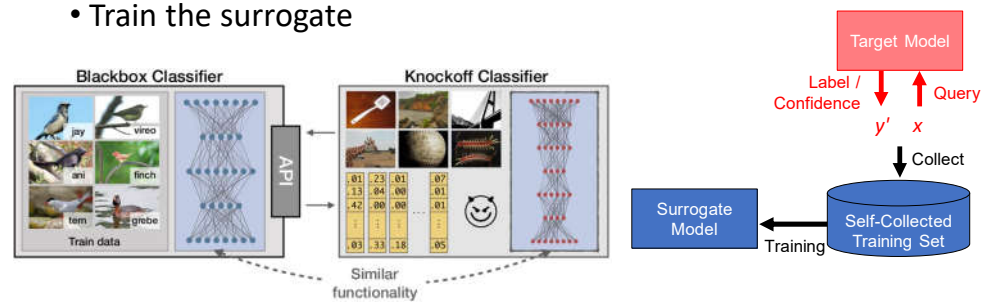
# Model Stealing

- Construct a copy of a model
- Two possible goals:
  - Intellectual property for well performance
  - Surrogate model for evasion attack



# Model Stealing Query Attack

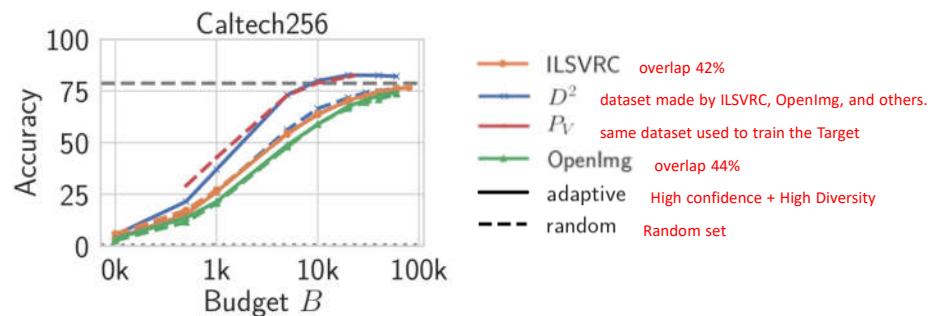
- Procedure
  - Query the target model and collect output
  - Train the surrogate



# Model Stealing Query Attack

- How to generate query samples?
  - Select samples randomly
    - May not be effective
    - Many queries are required
  - Select/Generate samples specifically according to
    - High Output Confidence: Only the confident samples
    - High Diversity: Different information

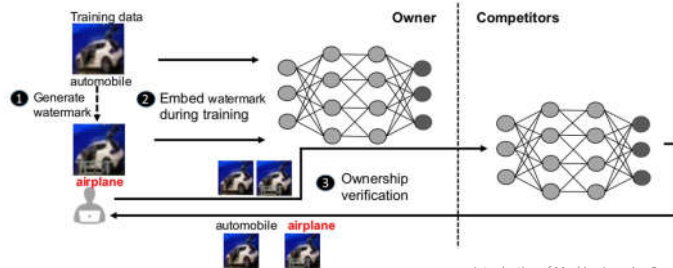
# Model Stealing Query Attack



## Model Stealing Defense Copyright Verification



- Add **watermarks** to the model for the **proof of intellectual property**
  - Watermarks: **patterns** that cause an **unexpected misclassification** when added to an image



Introduction of Machine Learning Security: Ch04

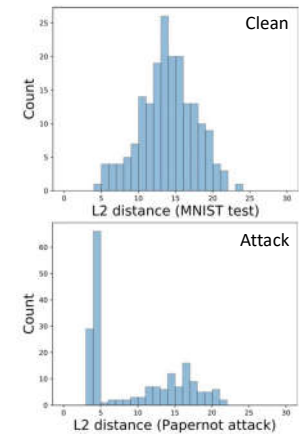
9

Zhang et al., Protecting Intellectual Property of Deep Neural ... AaCCS 2018

## Model Stealing Defense Block Query



- Understand that attackers are querying the model and block them.
- Distance between consecutive queries for a legitimate purpose usually follow a normal distribution, but not for an attack



Introduction of Machine Learning Security: Ch04

10

## Training Set Recovery



- Training samples may contain sensitive information
  - Personal information
  - Financial information
  - Images of suspected people
- Even if recovered training samples are incomplete, they can still be combined to re-identify individuals.

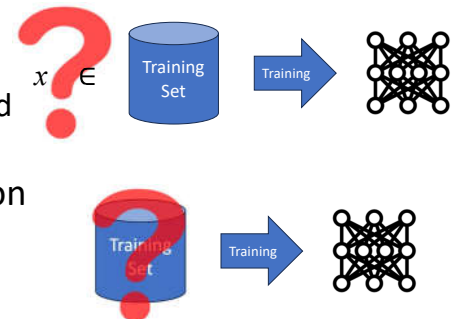
Introduction of Machine Learning Security: Ch04

11

## Training Set Recovery



- Different Levels of Recovery
- Training Sample Identification
  - Identify whether a sample used in training
- Training Sample Reconstruction
  - Construct the training set according to the model



Introduction of Machine Learning Security: Ch04

12

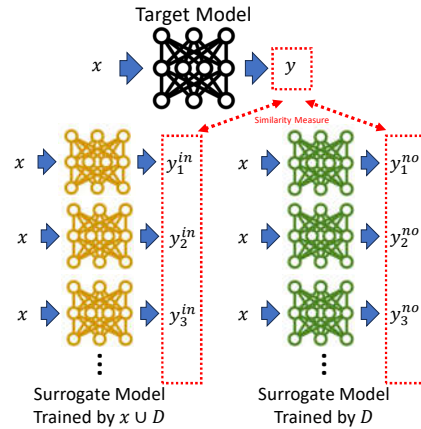
Shakiri et al., Membership Inference Attacks Against Machine Learning Models, S&P 2017

# Training Set Recovery

## Training Sample Identification

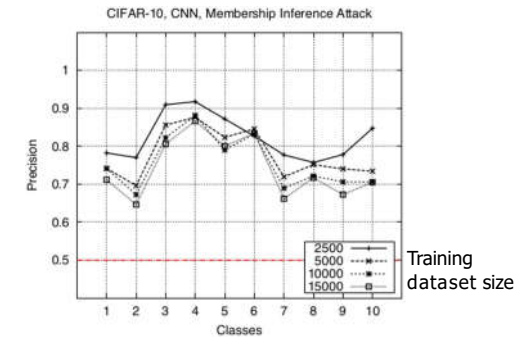
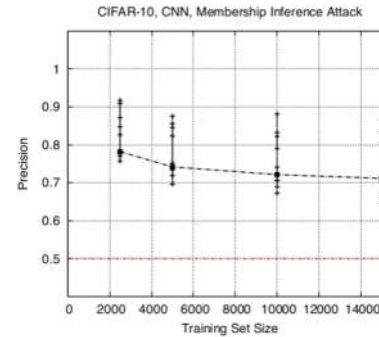


- Membership Inference Attacks
  - Query the target model with the input sample  $x$
  - Many surrogate model pairs are crafted:
    - including  $x$  in the training dataset
    - not including  $x$  in the training dataset
  - Determine whether the sample is used in training by comparing the predictions of those surrogate models with the target model



# Training Set Recovery

## Training Sample Identification



# Training Set Recovery

## Training Sample Reconstruction

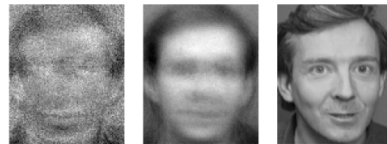


- Model Inversion Attack
  - Reconstruct a training sample by maximizing confidence with respect to the target label using gradient descent
  - Query ability is required

### Algorithm 1 Inversion attack for facial recognition models.

```

1: function MI-FACE(label,  $\alpha, \beta, \gamma, \lambda$ )
2:    $c(x) \stackrel{\text{def}}{=} 1 - \hat{f}_{\text{label}}(x) + \text{AUXTERM}(x)$ 
3:    $x_0 \leftarrow 0$ 
4:   for  $i \leftarrow 1 \dots \alpha$  do
5:      $x_i \leftarrow \text{PROCESS}(x_{i-1} - \lambda \cdot \nabla c(x_{i-1}))$ 
6:     if  $c(x_i) \geq \max(c(x_{i-1}), \dots, c(x_{i-\beta}))$  then
7:       break
8:     if  $c(x_i) \leq \gamma$  then
9:       break
10:  return  $[\arg \min_{x_i} (c(x_i)), \min_{x_i} (c(x_i))]$ 
    
```



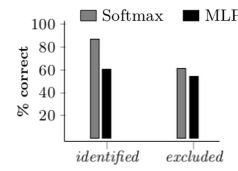
# Training Set Recovery

## Training Sample Reconstruction

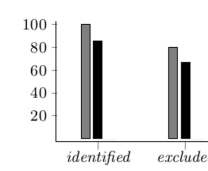


- Mechanical Turk workers are asked to match the reconstructed image to one of five face images from the original training set

	Identify by Workers	Cannot Identify by Workers
Present in the selected images	Identified	Not Match
Not present in the selected images	Not Match	Excluded



(a) Average over all responses.



(c) Accuracy with skilled workers.

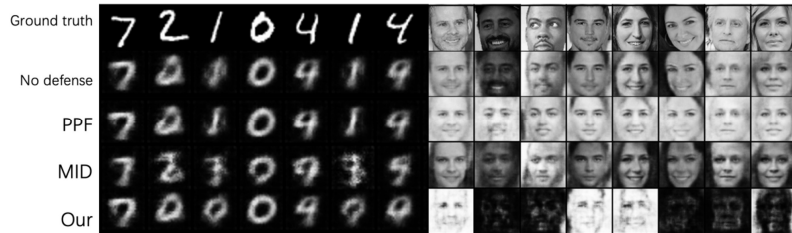


## Training Set Recovery Defense Model Inversion Attack



- Maximize the reconstruction error without changing the labels

$$\begin{aligned} & \max \mathcal{R}(x, \mathcal{A}(\mathbf{y} + \mathbf{e})) && \text{reconstruction error} \\ \text{subject to: } & \mathbf{e} \leq \epsilon && \text{upper bound on modification} \\ & \arg \max(\mathbf{y} + \mathbf{e}) = \arg \max \mathbf{y} && \text{same predicted labels on original and attack samples} \\ & 0 \leq (\mathbf{y}_i + \mathbf{e}_i) \leq 1, \sum (\mathbf{y}_i + \mathbf{e}_i) = 1 && \text{Probability maintenance} \end{aligned}$$



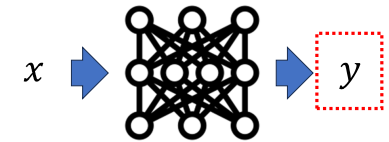
Introduction of Machine Learning Security: Ch04

17

## Training Set Recovery Defense Prediction Vector Tampering



- Privacy attacks usually assume knowledge of the classifier's scores
- Control the outputs of queries:
  - Score Blocking: provide only label but not scores for classes
  - Scores Perturbation: reduce reliability of scores



Jia et al., MemGuard: Defending against Black-Box Membership Inference Attacks via Adversarial Examples, CCS, 2019  
Shokri et al., Membership Inference Attacks Against Machine Learning Models, S&P, 2017  
Rigaki et al., A Survey of Privacy Attacks in Machine Learning, ArXiv, 2021

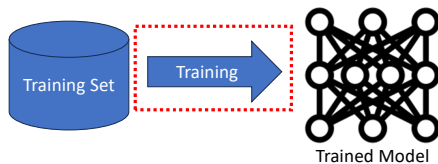
Introduction of Machine Learning Security: Ch04

18

## Training Set Recovery Defense Regularization



- Deep neural networks tend to memorize training data (they are really confident when predicting them)
- Considering additional terms that are irrelevant to the samples, such as regularization, can reduce memorization on the training samples



Introduction of Machine Learning Security: Ch04

19

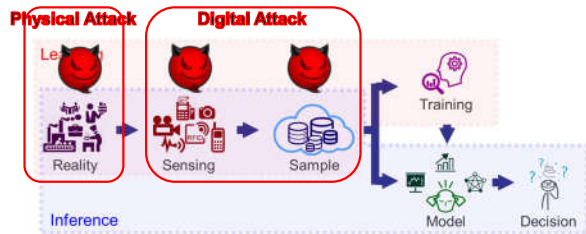
## Physical Attack



20

# Physical Attack

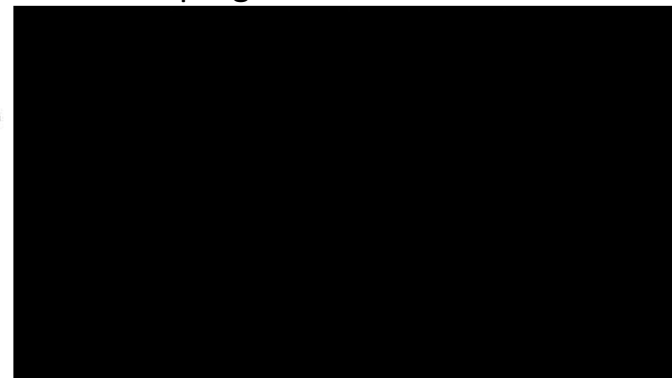
- Previous discussion focuses on **digital representation**
- Input can be **precisely controlled**
- Can adversarial attack be applied to our real world?



Mahmoud Sharif, Sruti Bhagavatula, Lujo Bauer(2016) Accessories to a crime: Real and stealthy attacks on state-of-the-art face recognition. In: 2016 ACM SIGSAC conference on computer and communications security

# Physical Attack

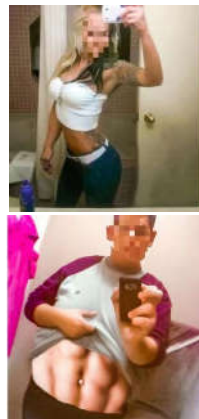
- A printed contaminated stop sign



Introduction of Machine Learning Security: Ch04

# Physical VS Digital World

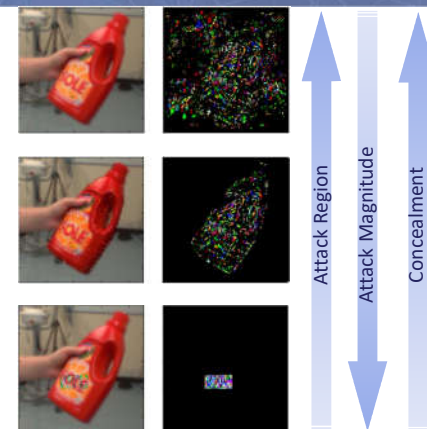
- **Gap** between physical and digital world
  - **Spatial Constraints**
    - **Adversarial noise** should only appears on the **object** but **not the background**
  - **Physical Limits on Imperceptibility**
    - **Small perturbations** are almost **imperceptible** to sensors
  - **Environmental Conditions**
    - **Distance, angle, lighting/weather conditions**
  - **Fabrication Error**
    - **Reproduction error**, e.g. printer limitation



Kevin Eykholt, Ivan Evtimov, Earlene Fernandes (2017) Robust Physical-World Attacks on Deep Learning Visual Classification. In: Proceedings of the IEEE conference on computer vision and pattern recognition

# Attack Region

- **Digital Attack**
  - Any features
  - Cannot be used in reality
- **Poster/Wrapper Attack**
  - Features in object
- **Sticker Attack**
  - Features is a small area
  - Easier to implement



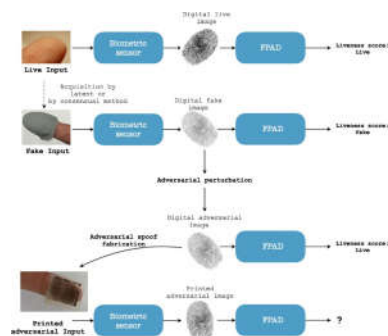
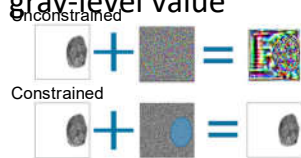
M Meili, A Demontis, B Biggio(2017) Is Deep Learning Safe for Robot Vision? Adversarial Examples against the iCub Humanoid. In: Proceedings of the IEEE International Conference on Computer Vision Workshops

## Attack Region Fingerprint

- Evade **fingerprint liveness detection**

- Attack is limited:

- **Region:**  
Actual fingerprint
- **Value:**  
only gray-level value



S Marrone, R Casola, G Orrù(2020) Fingerprint Adversarial Presentation Attack in the Physical Domain. In: ICPR

Introduction of Machine Learning Security: Ch04

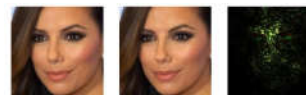
25

## Attack Region Glass

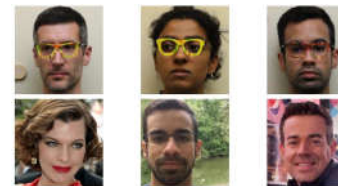
- Only attack the features in a glass mask



Attack the whole face



Attack the glass region



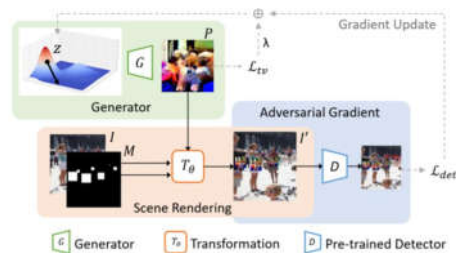
M Sharif, S Bhagavatula, I Bauer (2016) Accessorize to a Crime: Real and Stealthy Attacks on State-of-the-Art Face Recognition. In: Proceedings of the 2016 ACM SIGSAC conference on computer and communications security

Introduction of Machine Learning Security: Ch04

26

## Attack Region Clothes

- Embed a generated image to a clothing region



$$L_{det} = \frac{1}{N} \sum_{i=1}^N \max_j [D_{obj}^j(I'_i) D_{cls}^j(I'_i)]$$

adversarial detection loss

$$L_{tv} = \sum_{i,j} \sqrt{(P_{i+1,j} - P_{i,j})^2 + (P_{i,j+1} - P_{i,j})^2}$$

smoothness of a generated image

YCT Hu, BH Kung, DS Tran(2023) Naturalistic Physical Adversarial Patch for Object Detectors. In: ICV

Introduction of Machine Learning Security: Ch04

27

## Limit Attack Region is not enough

- Objects can be viewed from different distances and angles
- **Distance:** Approach to a printed contaminated stop sign
  - **Misclassified** as “sports ball” in **two frames**
- **Angle:** Camera moves closely around a printed original and contaminated stop signs
  - **Misclassified** as “toilet” in **two frames**



J Lu, H Sibol, E Fabry (2023) NO Need to Worry about Adversarial Examples in Object Detection in Autonomous Vehicles. In: arXiv

Introduction of Machine Learning Security: Ch04

28

# Environmental Conditions



- Simulate the real situations by considering transformations of viewpoint shifts, camera noise, and other natural noises

## • Expectation Over Transformation (EOT)

In different transformation

$$\arg \max_{x'} \mathbb{E}_{t \sim T} \left[ \log P(y_t | t(x')) - \lambda \| LAB(t(x')) - LAB(t(x)) \|_2 \right]$$

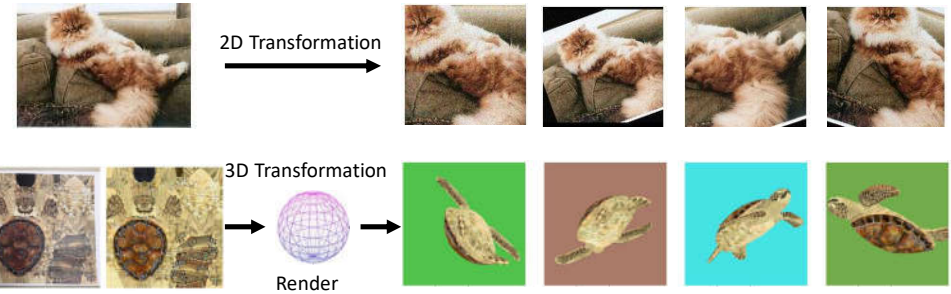
Wrong Decision                      Visual Imperceptibility

- T: Transformation
- LAB: a space for measuring human perceptual distance

# Environmental Conditions



- 2D: rotation, transformation, or addition of noise
- 3D: angle, texture and a pose of the 3D object

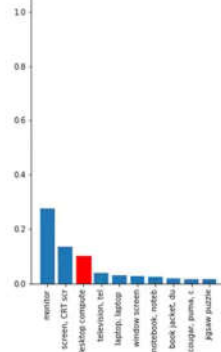


# Environmental Conditions

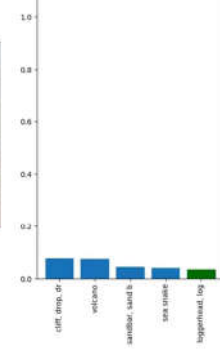


## • Expectation Over Transformation (EOT)

<https://www.youtube.com/watch?v=oeQW5qdey8>



<https://www.youtube.com/watch?v=YXy6oX1NoA>



# Environmental Conditions



## 2D Image



Image Type	$P(\text{true})$	$P(\text{adv})$
Original: speedboat	14%	9%
Adv: crossword puzzle	3%	91%

## 3D Model

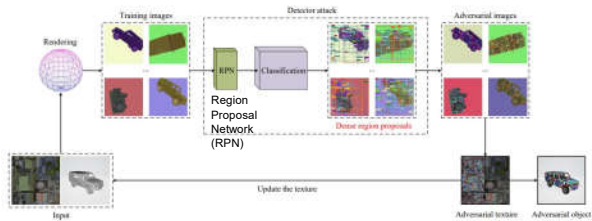
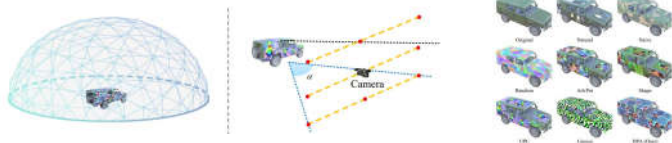


Image Type	$P(\text{true})$	$P(\text{adv})$
Original: orange	73%	0%
Adv: power drill	0%	89%



# Environmental Conditions Object Detection

- Consider different angles in reality



Yexin Duan, Jialin Chen, Xingyu Zhou(2023) DPA: Learning Robust Physical Adversarial Camouflages for Object Detectors. In:arXiv

Introduction of Machine Learning Security: Ch04

# Imperceptibility & Fabrication Error

- Consider printability
- Robust Physical Perturbations (RPP)

$$\operatorname{argmin}_{\delta} \lambda \|M_x \cdot \delta\|_p + NPS(M_x \cdot \delta) + \mathbb{E}_{x_i \sim X_v} J(f_{\theta}(x_i + T_i(M_x \cdot \delta)), y^*)$$

where

$M_x$ : Mask

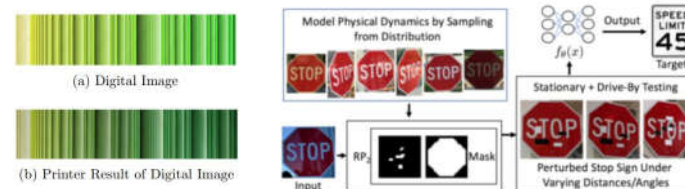
Manipulation  
Restriction

$\delta$ : Perturbation

Non-Printability  
Score (NPS)

Attack performance after  
different transformations

$X_v$ : set of victim images (under different transformations)



K Eykholt (2019) Designing and Evaluating Physical Adversarial Attacks and Defenses for Machine Learning Algorithms. In: Doctoral dissertation

Kevin Eykholt, Ivan Evrimov, Earlene Fernandes (2017) Robust Physical-World Attacks on Deep Learning Visual Classification. In: Proceedings of the IEEE conference on computer vision and pattern recognition

Introduction of Machine Learning Security: Ch04

# Imperceptibility & Fabrication Error

Distances/Angle	Subtle Poster	Subtle Poster Right Turn	Camouflage Graffiti	Camouflage Art (LISA-GNN)	Camouflage Art (GTSHD-GNN)
5°-10°					
5°-15°					
10°-10°					
10°-30°					
30°-10°					

Targeted-Attack Success

100%

83.33%

66.67%

100%

90%

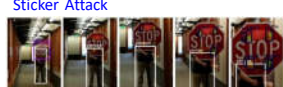
## Poster Attack



(a) The poster attack inside



(b) The poster attack outside



(c) The sticker attack inside



(d) The sticker attack outside

Introduction of Machine Learning Security: Ch04

# Natural Modification

- Sample are crafted more naturally



Attack modification is obvious

Attack modification is natural  
More concealment

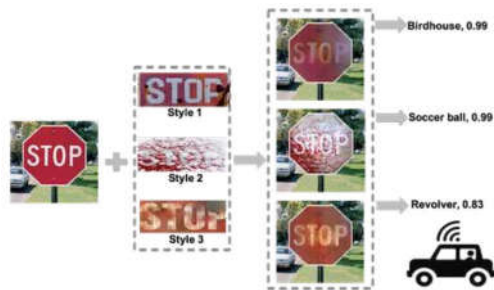
Introduction of Machine Learning Security: Ch04

# Natural Modification Style Transfer



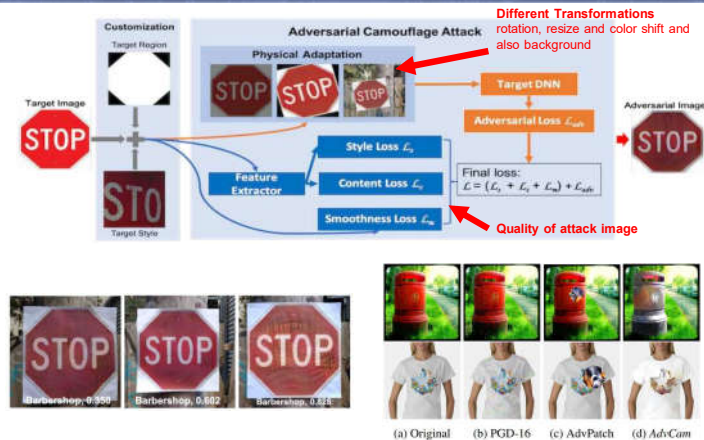
## Adversarial Camouflage (AdvCam)

- Mislead models by transferring style to objects
- Use style as adversarial noise
- Natural styles that appear legitimate to human observers



Introduction of Machine Learning Security: Ch04

# Natural Modification Style Transfer



Introduction of Machine Learning Security: Ch04

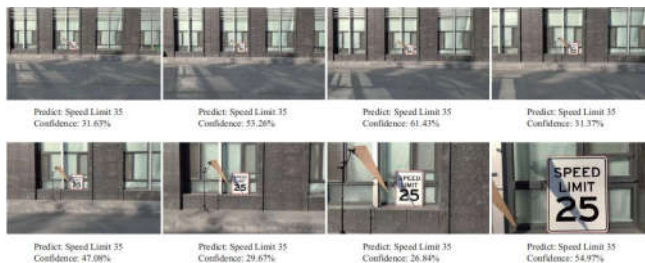
# Natural Modification Shadow



- A shadow with the simplest polygon — triangles, are sufficient to produce successful adversarial examples

$$\arg \min_v f_{true}(S(x, P_V, M, k)), \text{ s.t. } \tilde{y}_{adv} \neq y_{true}$$

$f_{true}()$ : confident score of a class  
 $S$ : surrogate model  
 $x$ : clean picture  
 $P_V$ : polygon vertices  
 $M$ : mask  
 $k$ : change pixel values of shadow area



Introduction of Machine Learning Security: Ch04

## Non-Security Applications



# Non-Security Application

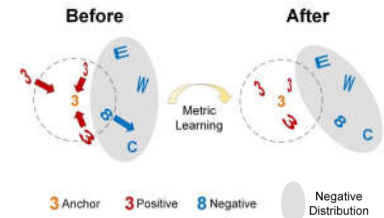
- Hard Sample Generation
- Uncertain Samples Selection

# Non-Security Application: Hard Sample Generation Metric Learning

- Aim to generate a **high dimensional space**
  - **Similar** samples are **close**
  - **Different** samples are **far away**
- **Triplet Loss** is a general objective function

$$[D(\mathbf{x}_i^+, \mathbf{x}_i) - D(\tilde{\mathbf{x}}_i^-, \mathbf{x}_i) + \alpha]_+$$

- $\mathbf{x}$ : anchor sample
- $\mathbf{x}^+$ : sample of the same class as anchor
- $\mathbf{x}^-$ : sample of different class to anchor
- $D$ : distance measure in metric learning space

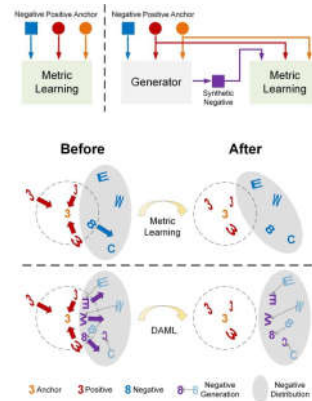


- Problem: **Negative sample** (even chosen by hard sampling) **may not be difficult enough**

# Non-Security Application: Hard Sample Generation Metric Learning

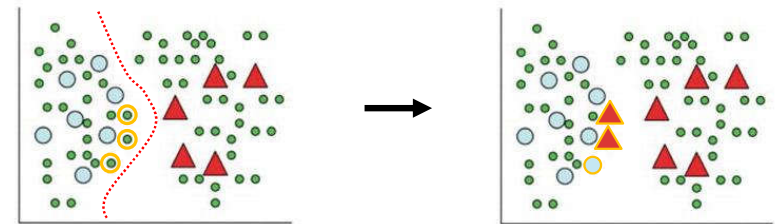
- **Craft hard negative samples** by adversarial attack
  - Similar to anchor and original negative sample ( $J_{\text{hard}}$  &  $J_{\text{reg}}$ )
  - **Generate the negative samples** on which the learned metric would **misclassify** ( $J_{\text{adv}}$ )

$$\begin{aligned} \min_{\theta_g} J_{\text{gen}} &= J_{\text{hard}} + \lambda_1 J_{\text{reg}} + \lambda_2 J_{\text{adv}} \\ &= \sum_{i=1}^N (\|\tilde{\mathbf{x}}_i^- - \mathbf{x}_i\|_2^2 + \lambda_1 \|\tilde{\mathbf{x}}_i^- - \mathbf{x}_i^-\|_2^2 \\ &\quad + \lambda_2 [D(\tilde{\mathbf{x}}_i^-, \mathbf{x}_i) - D(\mathbf{x}_i^+, \mathbf{x}_i) - \alpha]_+) \end{aligned}$$



# Non-Security Application: Uncertain Samples Selection Active Learning

- **Select samples for annotation** in semi-supervised learning problem iteratively **based on current model knowledge**
- **Most uncertain samples** are queried





## Don't be Pessimistic



- **Human can also be misled easily and also learn wrongly**
  - Just make **different mistakes from machine learning**
- Adversarial attack significantly harms the security and safety of ML systems, but...
- This threat provides us a chance to **understand better our models and data**

## Benefits from Adversarial Attack?



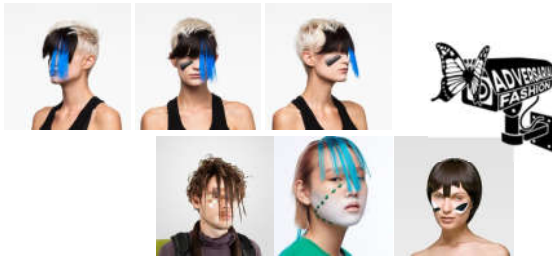
- A coin has two sides?
- Can we benefit from adversarial attack?



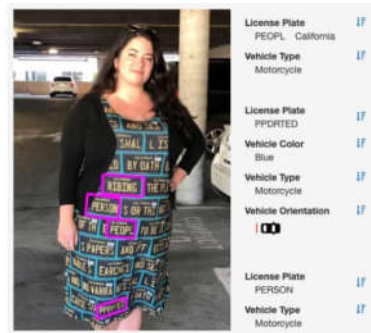
## Benefits from Adversarial Attack?



- **Avoid surveillance cameras?**
- Dress/Fashion/makeup is used to **evade or mislead** the detection



Key regions: Nose Bridge, nose, eyes, and forehead intersect



## Benefits from Adversarial Attack?

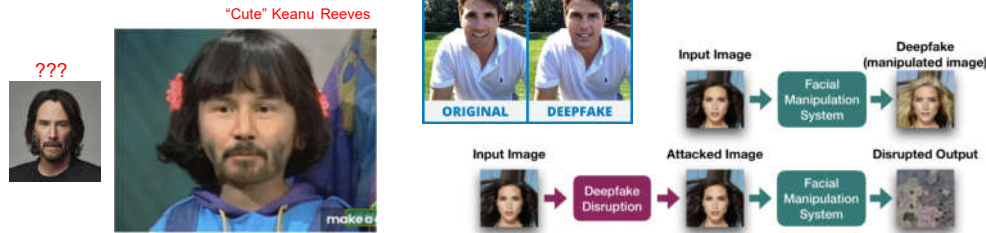


- **Hide from your enemy**
- **Evade optical aerial detection**



## Benefits from Adversarial Attack?

- Modified images of a person can be generated **without consent**, e.g. Deepfake
- Disrupt resulting images by adding adversarial noise to a photo



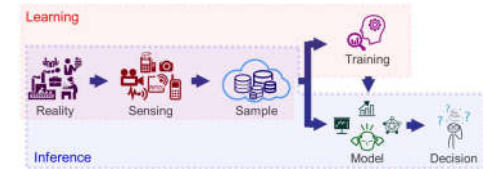
Notoniel Ruiz, Sarah Adel Bargal, Stan Sclaroff(2020) Disrupting Deepfakes: Adversarial Attacks Against Conditional Image Translation Networks and Facial Manipulation Systems. In: arXiv

Introduction of Machine Learning Security: Ch04

53

## Key Questions

- Where does **the training data** come from?
  - Provided by a third party?
- Who develops **the model**?
  - Is **pretrained model** used? If yes, where does it from?
- Who knows **the model details**?
- How to **capture samples in inference**?






Introduction of Machine Learning Security: Ch04

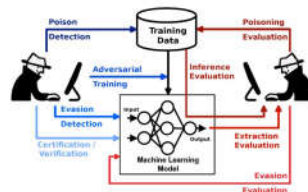
54

## Useful Library

### • Adversarial Learning Python Library

- Microsoft: Counterfit   
<https://github.com/Azure/counterfit/>
- IBM: Adversarial Robustness Toolbox   
<https://github.com/Trusted-AI/adversarial-robustness-toolbox>
- Pluribus One: SecML (Secure ML Library)   
<https://www.pluribus-one.it/research/sec-ml/sec-ml-lib>

- For Research and Engineering purposes



Introduction of Machine Learning Security: Ch04

55

## Welcome to Join Us!

- Besides publications...
- What you will learn...
  - Soft-Skill
  - Critical Thinking
  - Analytical Skill
  - Presentation Skill



You can work in our research lab : D1b-303



Introduction of Machine Learning Security: Ch04

56