

# Adversarial Machine Learning and Beyond

Philipp Benz

<https://phibenz.github.io>

Chaoning Zhang

<https://chaoningzhang.github.io>

# Agenda

Overview of adversarial machine learning:

1. Adversarial Examples
2. Adversarial attack
3. Adversarial defense

Our work:

4. UAP with class discrimination
5. Understanding UAP
6. Universal deep hiding
7. Towards a unified perspective

# Deep Learning is Awesome



CartPole-v0: Balance a pole on a cart (for a short time).

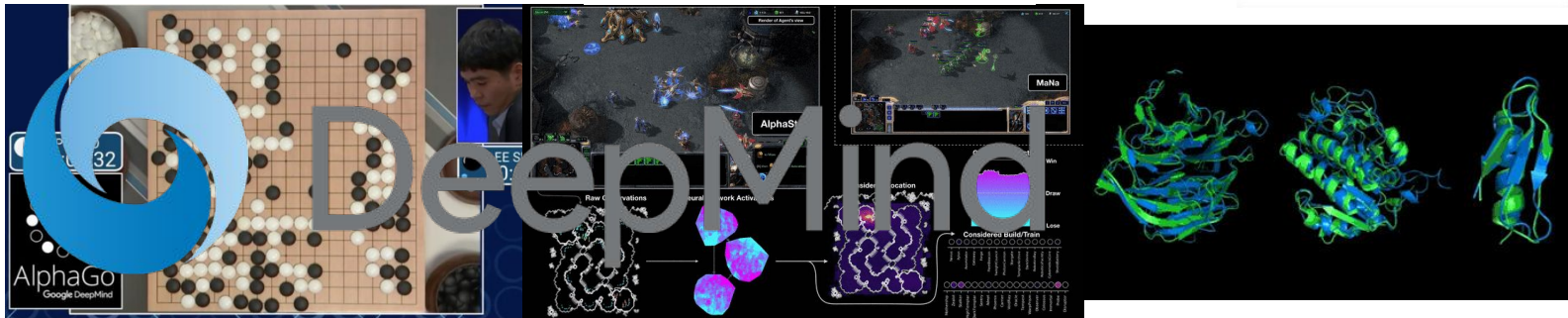
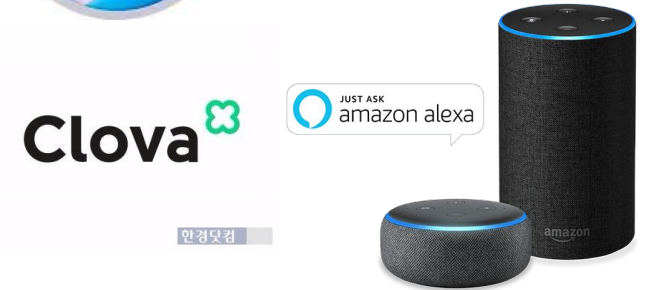
MountainCar-v0: Drive up a hill.

Pendulum-v0: Swing up a pendulum.

Reacher-v2: Make a 2D robot reach to a randomly located target.

InvertedPendulum-v2: Balance a pole on a cart.

MsPacman-ram-v0



# Overview of Adversarial Machine Learning

**Adversarial machine learning:** a technique attempting to fool models through malicious input.

- 1. Practical relevance:** Can we trust the model for security sensitive applications?
- 2. Theoretical relevance:** Understanding adversarial examples might give us insights about how deep neural networks work

# The Basics

# Adversarial Examples

Deep Neural Networks are sensitive to small perturbations in the image, which can lead to misclassifications. These changes are mostly imperceptible for human observers.


 $x$ 

“panda”

57.7% confidence

+ .007 ×


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

“nematode”

8.2% confidence

=


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

“gibbon”

99.3 % confidence

- Perturbations of small magnitude
- Specially crafted through optimization techniques
- Imperceptible for humans

# Objective

**Adversarial Attack:** Finding a small perturbation, to misclassify a sample

$$C(x + \delta) \neq C(x)$$

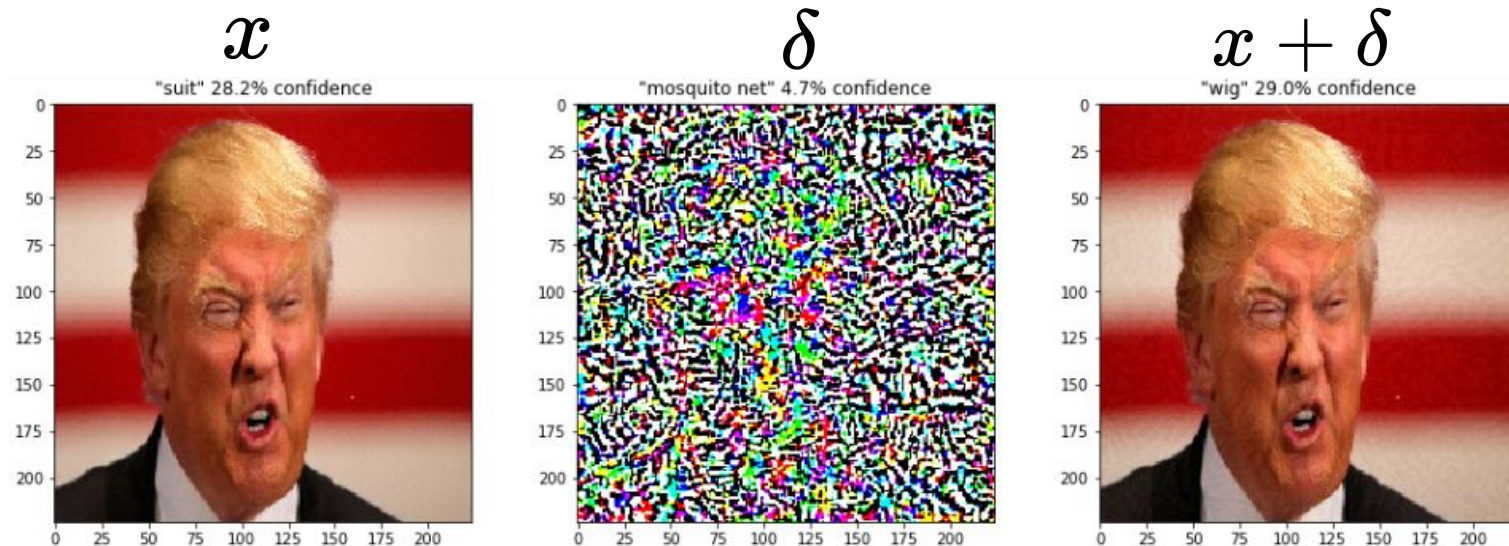
Misclassification

subject to  $D(x, x + \delta) \leq \epsilon$

The perturbation is smaller than magnitude  $\epsilon$   
Some distance metric L1, L2, Linf

$$x + \delta \in [0, 1]$$

Obey image range



# Distance metrics

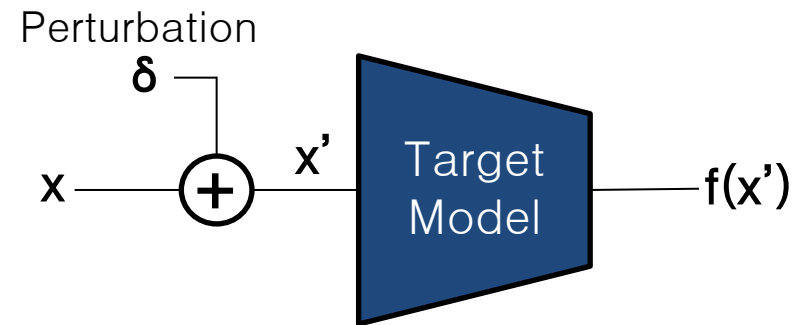
$$\delta = x - x'; \quad \|\delta\|_p = \left( \sum_{i=1}^n |\delta_i|^p \right)^{\frac{1}{p}}$$

- $L_0$  Measures the number of coordinates  $i$  such that  $x_i \neq x'_i$ .  
Corresponds to the number of pixels that have been altered in an image.
- $L_2$  Measures the standard Euclidean (root-mean-square) distance between  $x$  and  $x'$ . The  $L_2$  distance can remain small when there are many small changes to many pixels.
- $L_\infty$  Measures the maximum change to any of the coordinates:  $\|x - x'\|_\infty = \max(|x_1 - x'_1|, \dots, |x_n - x'_n|)$ .  
For images, we can imagine there is a maximum budget, and each pixel is allowed to be changed by up to this limit, with no limit on the number of pixels that are modified.

$L_0$  metric is rarely used. More and more papers are adopting the  $L_\infty$  metric for evaluation.



# Adversarial Examples: Attack categories



## White-Box Attack

An attacker has full knowledge about the target

## Black-Box Attack

An attacker has no knowledge about the target

## Untargeted Attack

The attack succeeds, if the target model makes a misclassification

## Targeted Attack

The attack succeeds, if the target model makes a misclassification which was specified beforehand

# White-box vs. Black-box

Adversary has **FULL** knowledge about:

- Model e.g.
  - type of neural network
  - number of layers
- Training algorithm e.g.
  - SGD, ADAM etc.
  - Learning Rate
- Training data
- Model parameters



FULL KNOWLEDGE

Adversary has **NO** knowledge about the model, training algorithm, training data, model parameters etc.

Commonly it is only assumed that the adversary can query the model

**Note:** Different works define their “black-box” slightly different, so **stay critical!**



ZERO KNOWLEDGE

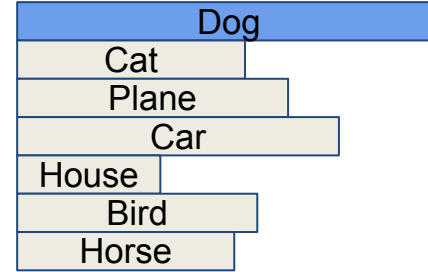
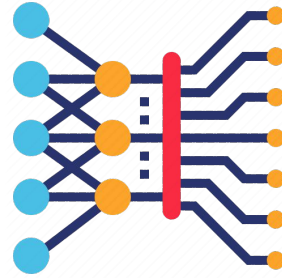
Adversary has **SOME** information.

E.g. the adversary knows the network architecture, but no information about the exact parameters.



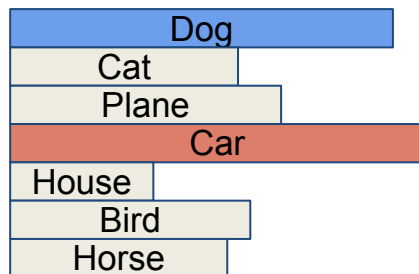
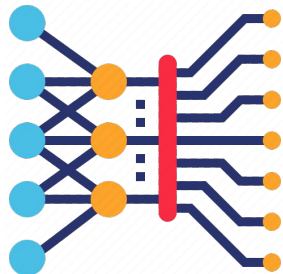
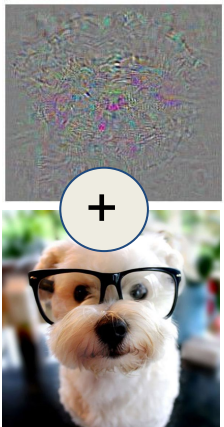
SOME KNOWLEDGE

# Untargeted vs. Targeted



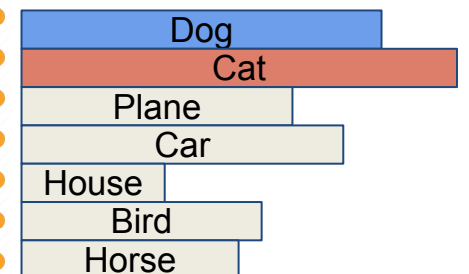
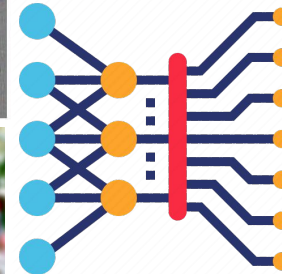
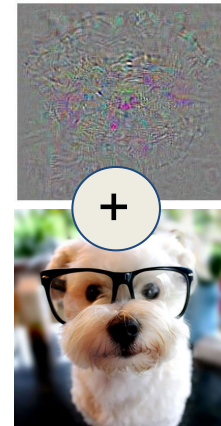
## Untargeted

The misclassified class can be **any other** than the currently classified class  
 In this example: any class other than dog, such as cat, car...



## Targeted

The misclassified class can **only be** the target class  
 In this example: target class = Cat



# Adversarial Attacks

# The FGSM attack

## Fast Gradient Sign Method (FGSM) [1]

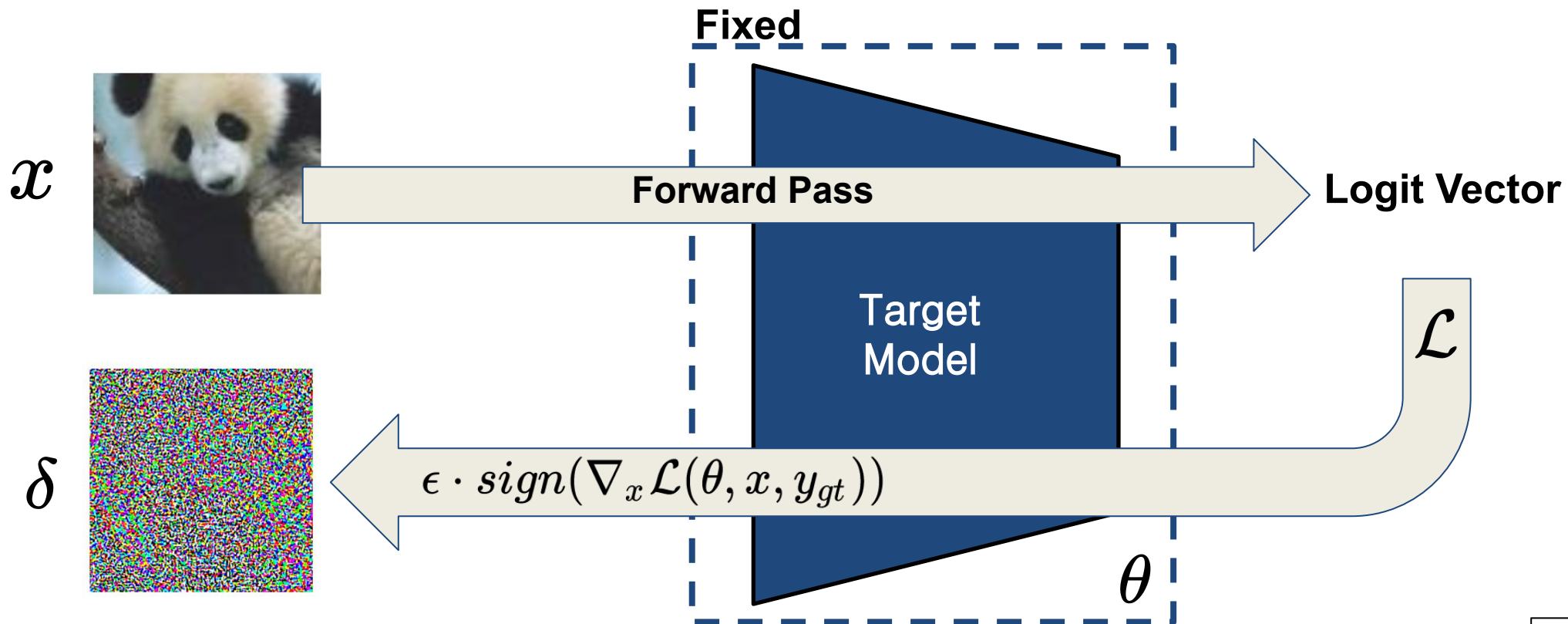
Stepping one step of step size  $\epsilon$  into the gradient direction

$$x' = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(\theta, x, y_{gt}))$$

Pro: Fast

Con: Less effective

$$\mathcal{L} = CE$$



$\epsilon$ : magnitude

# The PGD-attack

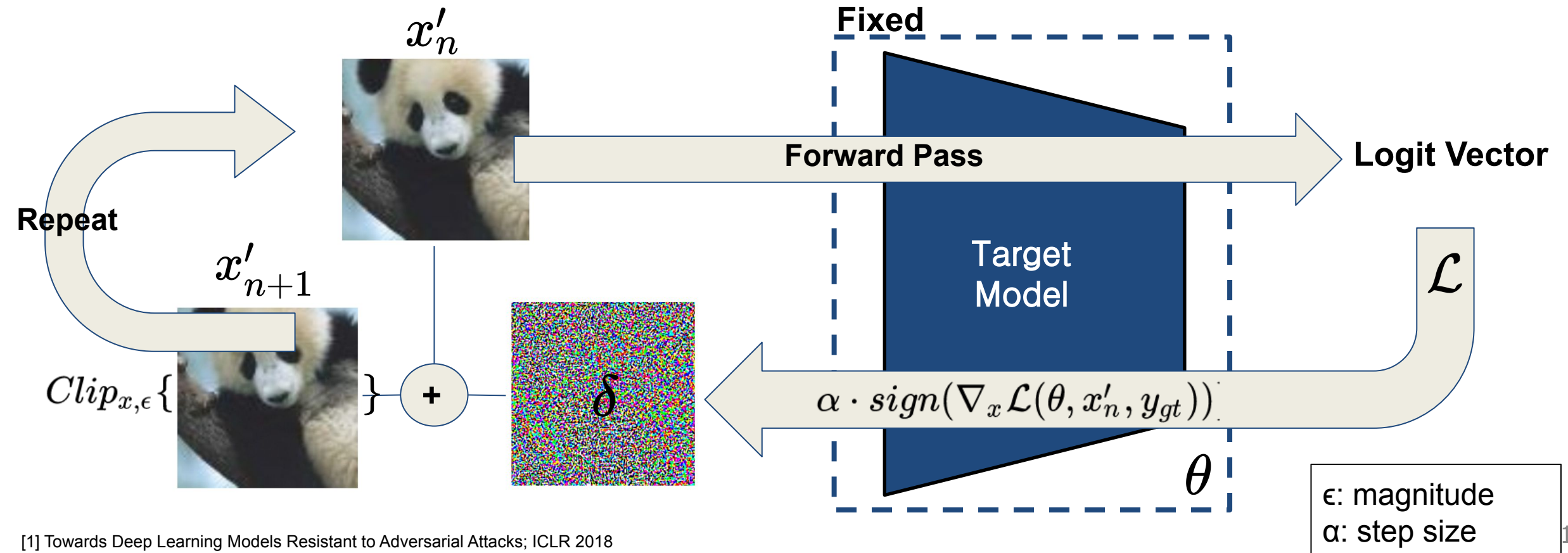
## Projected Gradient Descent (PGD)

Iteratively apply FGSM with step size  $\alpha$

$$x'_0 = x; \quad x'_{n+1} = \text{Clip}_{x,\epsilon} \{ x'_n + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(\theta, x'_n, y_{gt})) \}$$

Pro: Effective

Con: iterative, thus time-consuming



# Attack Techniques

In the white-box setting the gradients of the model can be used to construct the adversarial perturbation

## Fast Gradient Sign Method (FGSM) [1]

Stepping one step of step size  $\epsilon$  into the gradient direction

$$x' = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(\theta, x, y_{gt}))$$

**Pro: Fast**      **Con: Less effective**

$$\mathcal{L} = CE$$

## Targeted FGSM [2]

Stepping one step of step size  $\epsilon$  into the gradient direction of target class  $t$

$$x' = x - \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(\theta, x, t))$$

Compared with non targeted FGSM, the targeted one flips the loss sign and replace  $gt$  with  $t$

## Projected Gradient Descent (PGD) [2,3]

Iteratively apply FGSM with step size  $\alpha$

$$x'_0 = x; \quad x'_{n+1} = \text{Clip}_{x,\epsilon} \{x'_n + \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(\theta, x'_n, y_{gt}))\}$$

**Pro: Effective**      **Con: iterative, thus time-consuming**

## Carlini & Wagner attack (C&W) [4]

Minimal perturbation to make the highest ( $h$ ) logit  $Z(\cdot)$  lower than the second highest one

$$\mathcal{L} = \|\delta\|_2^2 + c \cdot f(x')$$

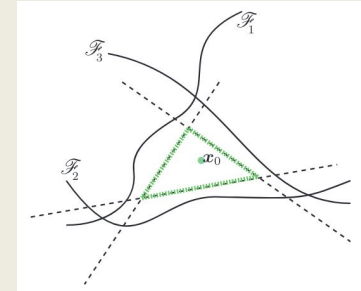
**Pro: Minimal perturbation thus least visible**

$$f(x') = \max(\max_{i \neq h} Z_i(\theta, x') - Z_h(\theta, x'), -\kappa)$$

**Con: Relatively more complex for minimizing the perturbation and loss function**

## DeepFool [5]

**Geometry-inspired** approach to push a sample over the **decision boundary**



Linear approximation of DNNs decision boundary

**Pro: Geometry-inspired**      **Con: Slow**

[1] Explaining and Harnessing Adversarial Examples; ICLR 2015  
[2] Adversarial examples in the physical world; ICLR Workshop 2017

[3] Towards Deep Learning Models Resistant to Adversarial Attacks; ICLR 2018  
[4] Towards Evaluating the Robustness of Neural Networks; Symposium on Security and Privacy (SP) 2017  
[5] DeepFool: a simple and accurate method to fool deep neural networks; CVPR 2016

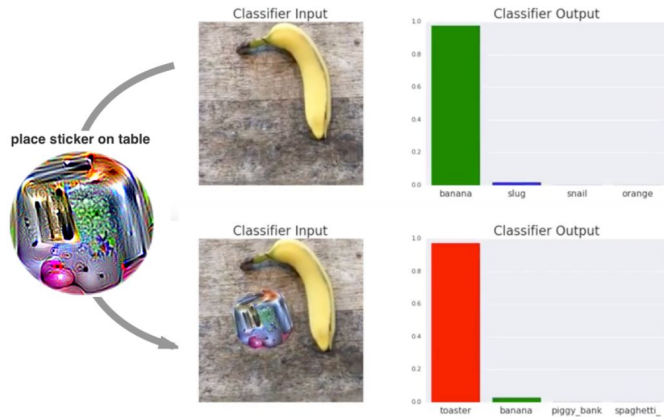
# Real-world Attacks

## 3D-printed adversarial objects [1]



■ classified as turtle    ■ classified as rifle  
■ classified as other

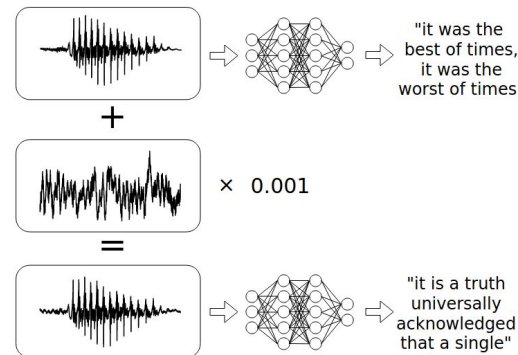
## Adversarial Patch [2]



## Black-Box attacks on web applications [3]

original image	true label	Clarifai.com results of original image	target label	targeted adversarial example	Clarifai.com results of targeted adversarial example
	viaduct	bridge, sight, arch, river, sky	window screen		window, wall, old, decoration, design
	hip, rose hip, rosehip	fruit, fall, food, little, wildlife	stupa, tope		Buddha, gold, temple, celebration, artistic
	dogsled, dog sled, dog sleigh	group together, four, sled, enjoyment	hip, rose hip, rosehip		cherry, branch, fruit, food, season
	pug, pug-dog	pug, friendship, adorable, purebred, sit	sea lion		sea seal, ocean, head, sea, cute

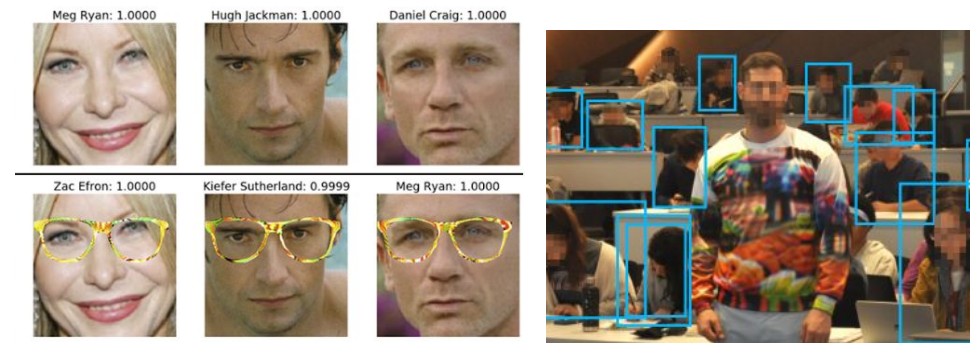
## Audio Adversarial Examples [5]



## Attacks on Traffic Signs [4]

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

## Adversarial Wearables [6,7]



[1] Synthesizing Robust Adversarial Examples; ICML 2018

[2] Adversarial Patch; 2017

[3] Delving into transferable adversarial examples and black-box attacks; ICLR 2017

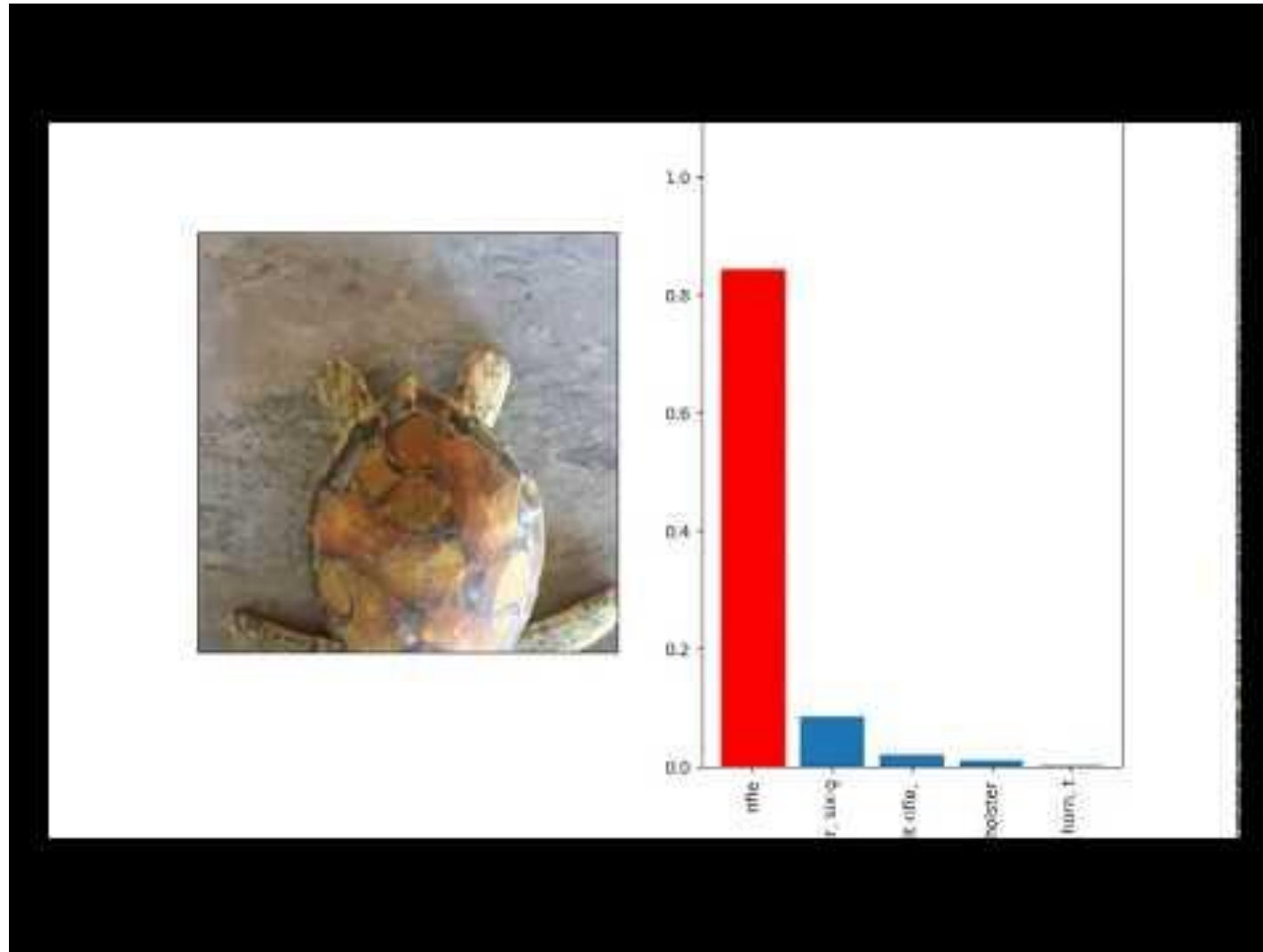
[4] Robust Physical-World Attacks on Deep Learning Visual Classification; CVPR 2018

[5] Audio Adversarial Examples: Targeted Attacks on Speech-to-Text; 2018

[6] Making an Invisibility Cloak: Real World Adversarial Attacks on Object Detectors; 2019



# 3D-printed adversarial objects



# Adversarial Defenses

# Attack & Defense: A two-player game

One player aims to attack the model; the other one tries to defend the model.

Who will win this game? It depends on the **rules** of the game.



**The attacker knows what  
the defender is doing  
(Attacker wins)**

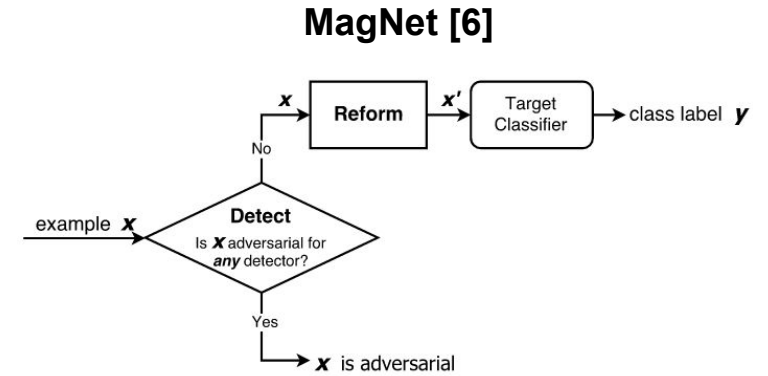
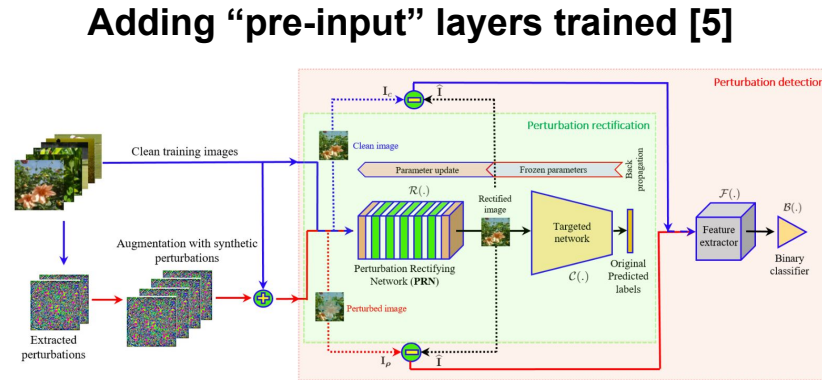
**The attacker does not know  
what the defender is doing  
(Defender wins)**

# Defenses

**1 Easy defense:** The attacker does not know what the defender is doing

**Motivation: Denoising the image**

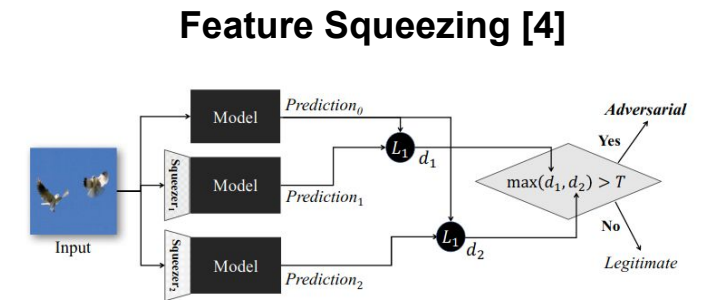
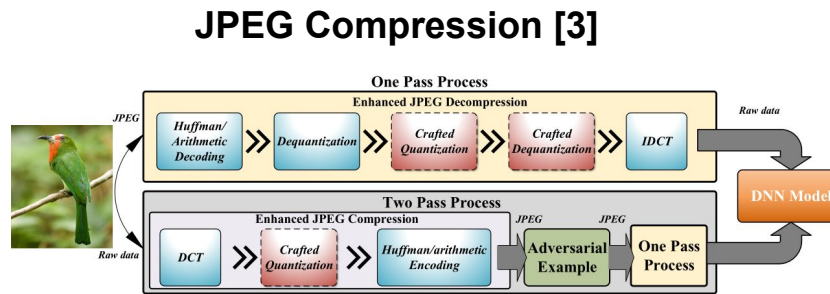
**Network Modification**  
Adding parts to the network or modifying the layers for defense



**2 Hard defense:** The attacker knows what the defender is doing

**Motivation: Obfuscated gradients**

**Input Modification**  
Manipulate the input to deweaponize the adversarial example



[1] Explaining and harnessing adversarial examples; ICLR 2015  
 [2] Generative Adversarial Trainer: Defense to Adversarial Perturbations with GAN; 2017  
 [3] Feature Distillation: DNN-Oriented JPEG Compression Against Adversarial Examples; 2018

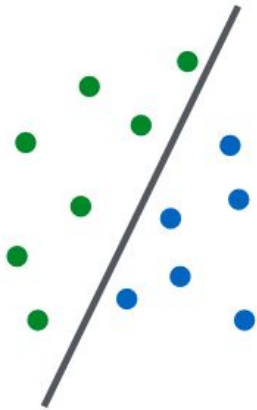
[4] Feature Squeezing: Detecting Adversarial Examples in Deep Neural Networks; NDSS 2018  
 [5] Defense against Universal Adversarial Perturbations; CVPR 2018  
 [6] MagNet: a Two-Pronged Defense against Adversarial Examples; CCS 2017

# Adversarial Training

Training a model (iteratively) with adversarial examples makes a model more robust. So far adversarial training is the only method that is proven to really robustify a model.

## Normal Network Training

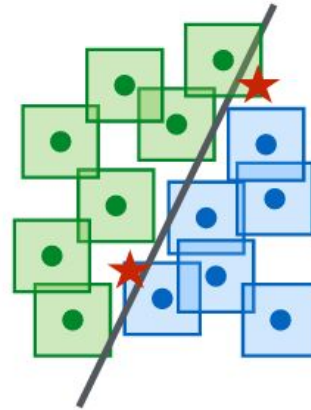
$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} [L(x, y, \theta)]$$



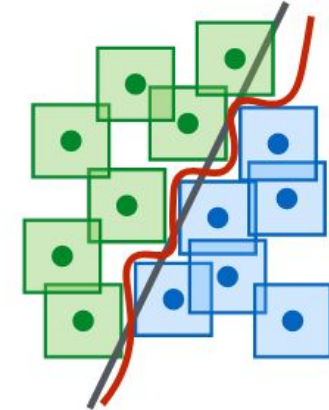
Set of points that can be easily separated with a simple linear decision boundary

## Adversarial Training

$$\min_{\theta} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ \max_{\delta \in \mathcal{S}} L(\theta, x + \delta, y) \right]$$



The simple decision boundary does not separate the  $L^{\infty}$ -balls (squares) around the data points  $\rightarrow$  adversarial examples (red stars) will be misclassified



Separating the  $L^{\infty}$ -balls requires a more complicated decision boundary. The resulting classifier is robust to adversarial examples with bounded  $L^{\infty}$ -norm perturbations.

[1] Explaining and harnessing adversarial examples; ICLR 2015

[2] Towards Deep Learning Models Resistant to Adversarial Attacks; ICLR 2018

# Robustness may be at odds with accuracy



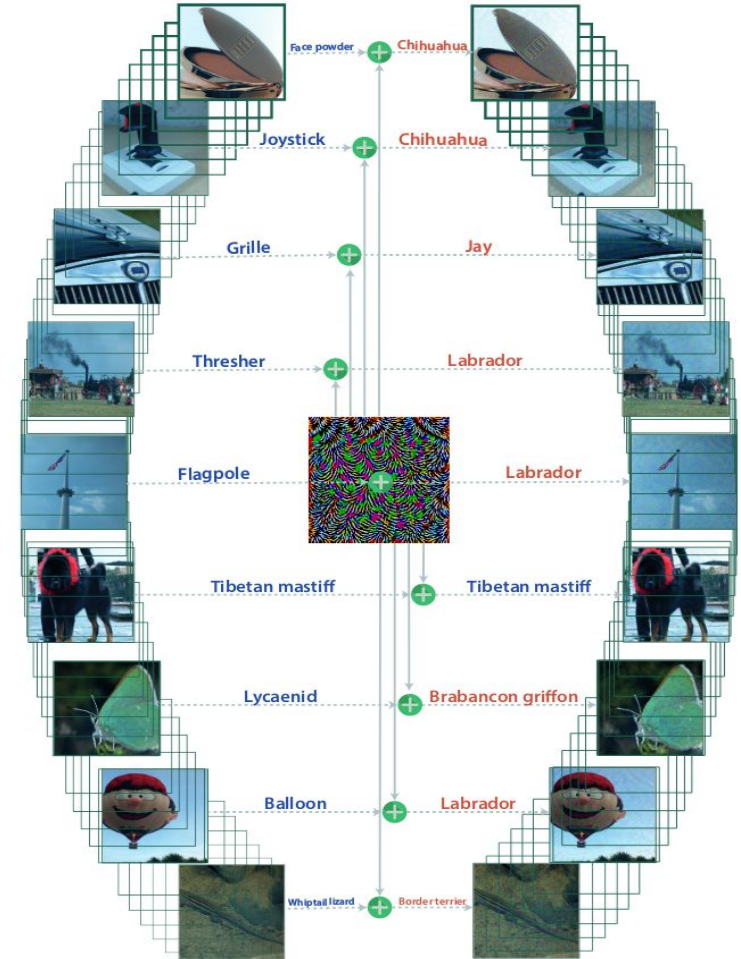
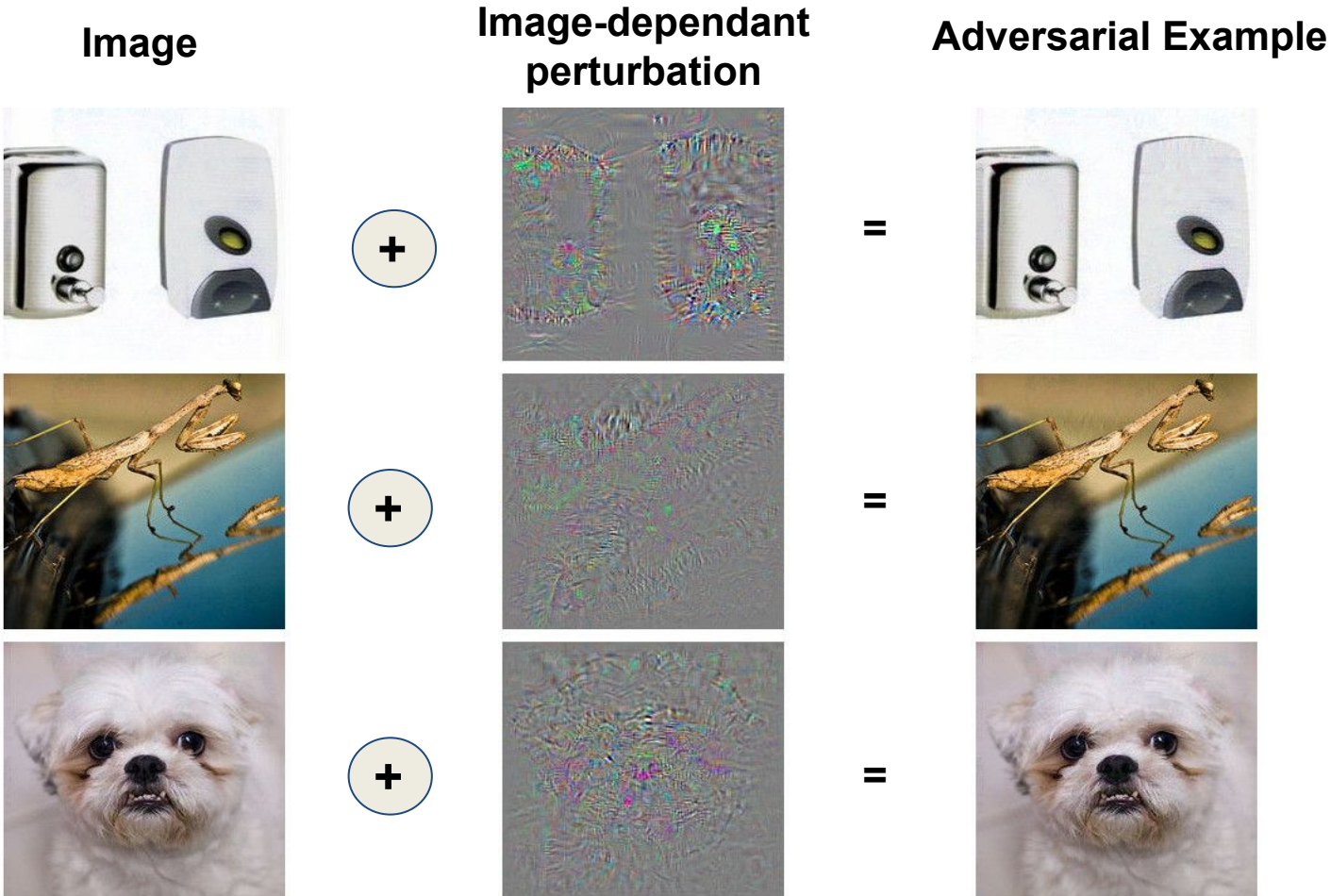
Trade-off between **standard accuracy** & **adversarially robust accuracy**

# Universal Adversarial Perturbations

# Image-dependant vs. image-agnostic

**Image-dependant attacks**  
 One perturbation is crafted to attack **one** specific image

**Image-agnostic (universal) attacks**  
 One perturbation is crafted to attack a **set** of images  
 UAP will be discussed in more detail again



[1] Intriguing properties of neural networks; Szegedy, Zaremba, Sutskever, Bruna, Erhan, Goodfellow, Fergus; ArXiv 2013  
 [2] Universal adversarial perturbations; Moosavi-Dezfooli, Fawzi, Fawzi, Frossard; CVPR 2017

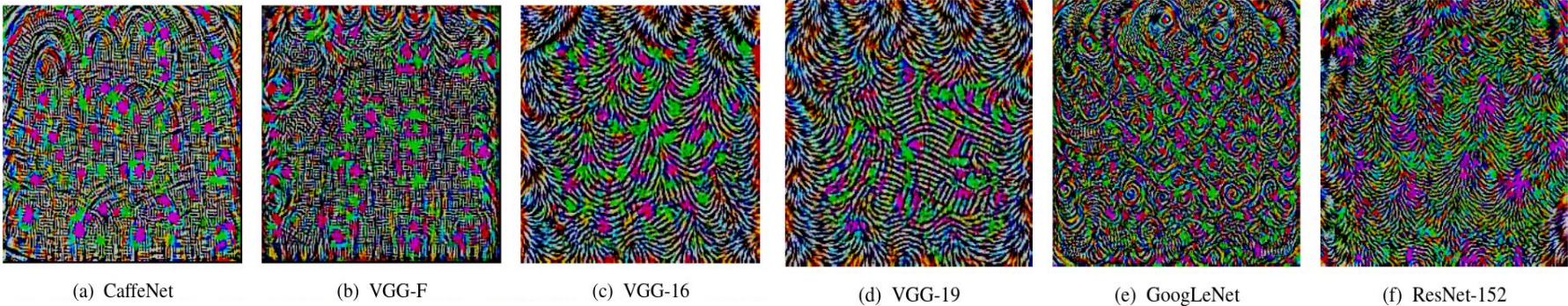


# Qualitative visualization

Images + UAP



Amplified UAP



[1] Universal Adversarial Perturbations; CVPR 2017

# Universal adversarial perturbations (UAPs)

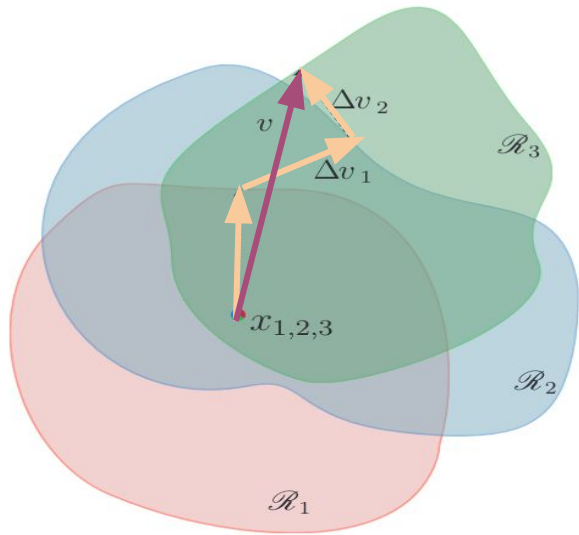
One perturbation to fool most of the data samples

## Objective

$C(x + \delta) \neq C(x)$  for "most"  $x \sim \mathcal{X}$  Fool most of the samples with a single perturbation

subject to  $\|\delta\|_p \leq \epsilon$  The perturbation is smaller than magnitude  $\epsilon$

# Classical DeepFool-based UAP



## Algorithm:

- Craft single perturbation (via DeepFool [3]) to let one sample cross the decision boundary
- Iterate this process for different samples to aggregate the universal adversarial perturbation.

Figure 1: Schematic representation of the algorithm in [1] to compute universal perturbations.

## Algorithm 1 Computation of universal perturbations.

- 1: **input:** Data points  $X$ , classifier  $\hat{k}$ , desired  $\ell_p$  norm of the perturbation  $\xi$ , desired accuracy on perturbed samples  $\delta$ .
- 2: **output:** Universal perturbation vector  $v$ .
- 3: Initialize  $v \leftarrow 0$ .
- 4: **while**  $\text{Err}(X_v) \leq 1 - \delta$  **do**
- 5:     **for** each datapoint  $x_i \in X$  **do**
- 6:         **if**  $\hat{k}(x_i + v) = \hat{k}(x_i)$  **then**
- 7:             Compute the *minimal* perturbation that sends  $x_i + v$  to the decision boundary:
 
$$\Delta v_i \leftarrow \arg \min_r \|r\|_2 \text{ s.t. } \hat{k}(x_i + v + r) \neq \hat{k}(x_i).$$
- 8:             Update the perturbation:
 
$$v \leftarrow \mathcal{P}_{p,\xi}(v + \Delta v_i).$$
- 9:         **end if**
- 10:     **end for**
- 11: **end while**

[1] Universal adversarial perturbations; Moosavi-Dezfooli, Fawzi, Fawzi, Frossard; CVPR 2017

[2] Analysis of universal adversarial perturbations; Moosavi-Dezfooli, Fawzi, Fawzi, Frossard, Soatto; ArXiv 2017

[3] DeepFool: a simple and accurate method to fool deep neural networks; Moosavi-Dezfooli, Fawzi, Frossard; CVPR 2016

# Practical Relevance of UAPs

## Image-Dependent Adversarial Examples

Have to be crafted on the spot for each sample  
(during inference)

## Universal Adversarial Examples

Can be crafted beforehand  
Only a simple summation is required during  
inference

UAPs are better suited for “real-world” attacks

# Our Research beyond classical UAP

1. The UAP attacks samples from every class without intended discrimination, which can easily cause suspicion
2. The UAP generation relies on the original training data & The reason for the existence of UAPs is still not well understood
3. Can the insight of UAPs be applied to the field of data hiding?

# Our Research beyond classical UAP

1. The UAP attacks samples from every class without intended discrimination, which can easily cause suspicion
  - CD-UAP: Class-Discriminative Universal Adversarial Perturbations (AAAI 2020)
  - Double targeted universal adversarial perturbations (ACCV 2020)
2. The UAP generation relies on the original training data & The reason for the existence of UAPs is still not well understood
  - Understanding Adversarial Examples from the Mutual Influence of Images and Perturbations (CVPR 2020)
3. Can the insight of UAPs be applied to the field of data hiding?
  - UDH: Universal Deep Hiding for Steganography, Watermarking and Light Field Messaging (NeurIPS 2020)
  - Universal Adversarial Perturbations Through the Lens of Deep Steganography: Towards A Unified Fourier Perspective (AAAI 2021)

# CD-UAP: Class-Discriminative Universal Adversarial Perturbations (AAAI 2020)

Chaoning Zhang\*, Philipp Benz\*, Tooba Imtiaz, In-So Kweon  
\* indicates equal contribution

Korea Advanced Institute of Science and Technology (KAIST)

# CD-UAP: Class-Discriminative Universal Adversarial Perturbations

An attack method to generate a universal perturbations that fool a target network to misclassify only a chosen group of classes, while having limited influence on the remaining classes.

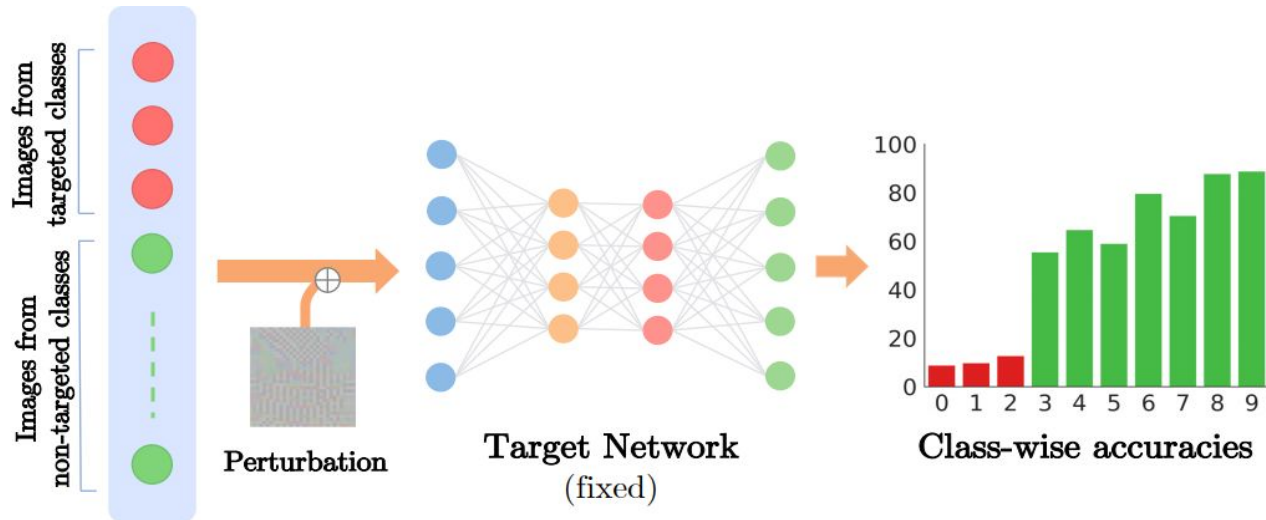


Figure 1: Class Discriminative Universal Adversarial Perturbation (CD-UAP). After adding a single perturbation, model performance on a subset of classes is significantly reduced, while the influence on the non-targeted classes is limited.

$$C(x_t + \delta) \neq C(x_t) \quad \text{for most } x_t \sim \mathcal{X}_t$$

Misclassify most samples from targeted classes

$$C(x_{nt} + \delta) = C(x_{nt}) \quad \text{for most } x_{nt} \sim \mathcal{X}_{nt}$$

Keep original classification for most samples from non-targeted classes

$$\|\delta\|_p \leq \epsilon$$

The perturbation is smaller than magnitude  $\epsilon$



## CD-UAP: Class-Discriminative UAP

Evaluation metric: Absolute Accuracy Drop (AAD) for both targeted classes and non-targeted classes.

$$\Delta_{\text{AAD}} = \text{AAD}_t - \text{AAD}_{\text{nt}}, \text{ the higher the better.}$$

Table 2: Experiments on CIFAR10 and CIFAR100 for various choices of loss functions. The two accuracies in each entry show the  $\text{AAD}_t$  and  $\text{AAD}_{\text{nt}}$ . The initial performances for CIFAR10 on ResNet20 are  $\text{Acc}_t = 90.86$ ,  $\text{Acc}_{\text{nt}} = 92.44$ ; and for CIFAR100 on ResNet56 are  $\text{Acc}_t = 65.20$ ;  $\text{Acc}_{\text{nt}} = 70.43$ , for targeted classes 0 to 4

CIFAR	$\mathcal{L}_t$	$\mathcal{L}_{\text{nt}}$							
		-	$\mathcal{L}^{\text{CE}}$		$\mathcal{L}^{\text{L}}$		$\mathcal{L}^{\text{BL}}$		
10	$\mathcal{L}^{\text{CE}}$	87.04	69.40	84.36	48.48	84.68	52.52	83.86	45.06
	$\mathcal{L}^{\text{L}}$	84.64	54.04	85.74	24.14	81.04	25.08	84.94	23.78
	$\mathcal{L}^{\text{BL}}$	88.58	61.54	81.86	10.18	64.66	7.92	81.52	11.60
100	$\mathcal{L}^{\text{CE}}$	61.20	60.23	60.40	50.50	1.60	1.15	62.40	46.35
	$\mathcal{L}^{\text{L}}$	63.80	55.00	63.40	47.12	62.40	46.02	62.80	47.97
	$\mathcal{L}^{\text{BL}}$	63.20	53.67	48.80	17.31	23.40	7.16	48.00	17.79

This loss combination is the best

$\mathcal{L}^{\text{CE}}$ : Cross-entropy loss (Increase for  $\mathcal{L}_t$  and decrease for  $\mathcal{L}_{\text{nt}}$ )

$\mathcal{L}^{\text{L}}$ : Logit loss (Decrease for  $\mathcal{L}_t$  and increase for  $\mathcal{L}_{\text{nt}}$ )

$\mathcal{L}^{\text{BL}}$ : Bounded Logit Loss (Similar to C&W loss; Decrease for  $\mathcal{L}_t$  and increase for  $\mathcal{L}_{\text{nt}}$ )

[1] CD-UAP: Class-Discriminative Universal Adversarial Perturbations; Zhang\*, Benz\*, Imtiaz, Kweon; AAAI 2020

Table 3: Experiments on CIFAR10 with different groups of targeted classes using VGG16 and ResNet20

$\hat{F}$	$S$	$\text{Acc}_t$	$\text{AAD}_t$	$\text{Acc}_{\text{nt}}$	$\text{AAD}_{\text{nt}}$	$\Delta_{\text{AAD}}$
VGG16	[1 : 5 : 2]	90.57	78.63	94.23	14.74	<b>63.89</b>
	[2 : 6 : 2]	92.87	68.87	93.24	21.79	<b>47.08</b>
	[0 : 4 : 1]	92.40	75.00	93.86	7.36	<b>67.64</b>
	[5 : 9 : 1]	93.86	75.00	92.40	17.56	<b>57.44</b>
	[0 : 6 : 1]	92.10	69.24	95.53	3.13	<b>66.11</b>
	[3 : 9 : 1]	92.80	78.60	93.90	8.70	<b>69.90</b>
ResNet20	[1 : 5 : 2]	88.80	82.27	92.87	15.93	<b>66.34</b>
	[2 : 6 : 2]	92.30	79.33	91.37	21.40	<b>57.93</b>
	[0 : 4 : 1]	90.86	81.46	92.44	10.24	<b>71.22</b>
	[5 : 9 : 1]	92.44	80.16	90.86	17.36	<b>62.80</b>
	[0 : 6 : 1]	90.81	81.33	93.60	4.67	<b>76.66</b>
	[3 : 9 : 1]	91.21	80.99	92.67	7.90	<b>73.09</b>

The  
higher  
the  
better

$S$ : targeted classes, [lowest class, highest class, step size], e.g.

[1:5:2] indicates targeted classes 1,3,5 and

[0:4:1] indicates targeted classes 0,1,2,3,4

# CD-UAP: Class-Discriminative UAP

## Further results on CIFAR100 and ImageNet

Table 5: Experiments on CIFAR100 targeting superclasses using VGG19 and ResNet56

Super Class	VGG19					ResNet56				
	$Acc_t$	$AAD_t$	$Acc_{nt}$	$AAD_{nt}$	$\Delta_{AAD}$	$Acc_t$	$AAD_t$	$Acc_{nt}$	$AAD_{nt}$	$\Delta_{AAD}$
aquatic mammals	56.20	44.20	70.97	14.05	<b>30.15</b>	58.20	43.60	70.80	14.87	<b>28.73</b>
fish	67.00	45.60	70.40	18.25	<b>27.35</b>	67.80	49.00	70.30	18.45	<b>30.55</b>
flowers	76.40	33.80	69.91	20.57	<b>13.23</b>	75.40	30.80	69.90	24.18	<b>6.62</b>
food containers	69.60	52.40	70.26	16.28	<b>36.12</b>	71.40	54.40	70.11	16.44	<b>37.96</b>
fruit and vegetables	79.40	55.60	69.75	18.04	<b>37.56</b>	76.60	52.00	69.83	17.65	<b>34.35</b>
household electrical devices	70.00	49.20	70.24	15.95	<b>33.25</b>	75.00	55.60	69.92	16.52	<b>39.08</b>
household furniture	77.40	63.80	69.85	15.76	<b>48.04</b>	76.00	60.00	69.86	14.90	<b>45.10</b>
insects	70.20	38.60	70.23	16.99	<b>21.61</b>	71.60	45.20	70.09	21.24	<b>23.96</b>
large carnivores	70.20	53.60	70.23	15.50	<b>38.10</b>	70.40	60.20	70.16	15.44	<b>44.76</b>
large man-made outdoor objects	82.60	66.60	69.58	15.56	<b>51.04</b>	81.20	66.20	69.59	15.49	<b>50.71</b>
large natural outdoor scenes	79.00	70.80	69.77	13.12	<b>57.68</b>	78.40	67.60	69.74	12.13	<b>55.47</b>
large omnivores and herbivores	69.80	51.60	70.25	18.33	<b>33.27</b>	71.40	56.00	70.10	17.15	<b>38.85</b>
medium-sized mammals	72.80	49.40	70.09	17.27	<b>32.13</b>	71.80	54.20	70.08	18.83	<b>35.37</b>
non-insect invertebrates	66.40	40.60	70.43	18.09	<b>22.51</b>	66.00	44.60	70.39	18.38	<b>26.22</b>
people	52.20	31.60	71.18	14.67	<b>16.93</b>	48.40	32.60	71.32	14.85	<b>17.75</b>
reptiles	59.20	43.80	70.81	16.05	<b>27.76</b>	59.20	42.80	70.74	17.24	<b>25.56</b>
small mammals	56.00	47.60	70.98	15.14	<b>32.46</b>	56.40	47.20	70.90	13.34	<b>33.86</b>
trees	66.80	49.60	70.41	15.58	<b>34.02</b>	66.60	54.00	70.36	17.05	<b>36.95</b>
vehicles 1	82.20	44.00	69.60	17.99	<b>26.01</b>	80.60	51.60	69.62	20.44	<b>31.16</b>
vehicles 2	81.20	53.40	69.65	21.33	<b>32.07</b>	81.00	62.20	69.60	21.92	<b>40.28</b>

Table 7: Experiments on ImageNet targeting 1 super class,  $\epsilon = 15$ .

Super Classes		$Acc_t$	$AAD_t$	$Acc_{nt}$	$AAD_{nt}$	$\Delta_{AAD}$
VGG16	Frogs	70.0	46.0	71.6	19.5	<b>26.5</b>
	Sharks	80.0	53.3	71.6	16.8	<b>36.5</b>
	Aircrafts	78.0	69.3	71.6	17.7	<b>51.6</b>
	Racket Radiator Radio	68.6	42.7	71.6	18.5	<b>24.2</b>
	Space objects	56.7	20.0	71.6	18.5	<b>1.5</b>
VGG19	Frogs	70.0	42.0	72.4	18.0	<b>24.0</b>
	Sharks	81.3	55.3	72.3	16.1	<b>39.2</b>
	Aircrafts	81.3	72.0	72.4	17.2	<b>54.8</b>
	Racket Radiator Radio	69.3	37.3	72.4	17.3	<b>20.0</b>
	Space objects	59.3	26.6	72.4	20.0	<b>6.6</b>
ResNet50	Frogs	75.3	52.0	76.1	18.5	<b>33.5</b>
	Sharks	83.3	66.7	76.1	16.1	<b>50.6</b>
	Aircrafts	84.0	65.3	76.1	16.8	<b>48.5</b>
	Racket Radiator Radio	75.3	46.7	76.1	17.1	<b>29.6</b>
	Space objects	59.3	43.3	76.2	20.4	<b>22.9</b>
ResNet152	Frogs	74.0	48.0	78.3	17.2	<b>30.8</b>
	Sharks	80.0	61.3	78.3	12.9	<b>48.4</b>
	Aircrafts	86.0	78.7	78.3	16.1	<b>62.6</b>
	Racket Radiator Radio	72.7	44.0	78.3	14.4	<b>29.6</b>
	Space objects	64.7	36.7	78.4	16.7	<b>20.0</b>

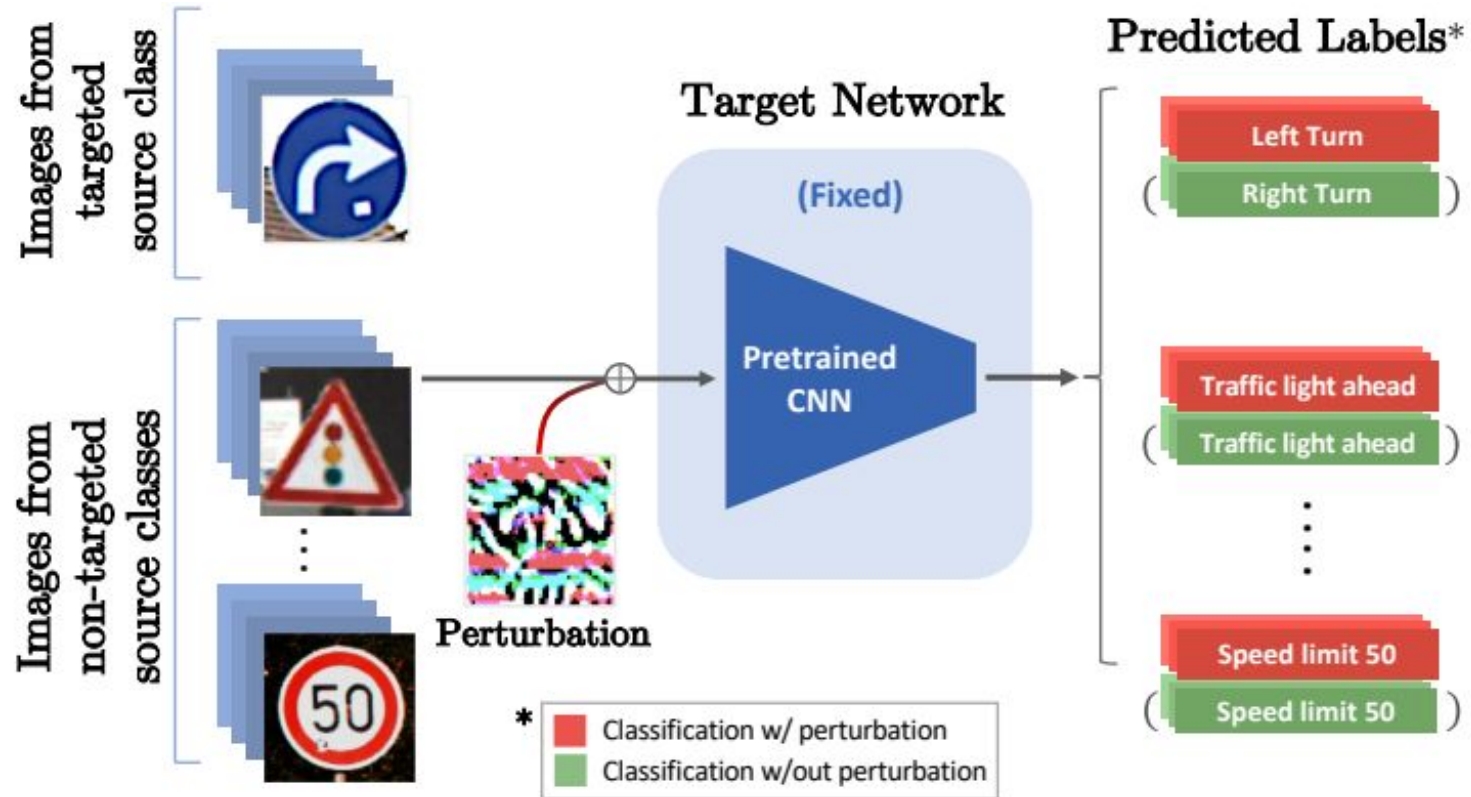
# Double targeted universal adversarial perturbations (ACCV 2020)

Philipp Benz\*, Chaoning Zhang\*, Tooba Imtiaz, In-So Kweon  
\* indicates equal contribution

Korea Advanced Institute of Science and Technology (KAIST)

# Double Targeted UAPs

UAP to fool one targeted source class to a sink class, while having limited adversarial effect on other non-targeted source classes to avoid raising suspicion.



# Loss Design

The loss consists of a targeted ( $\mathcal{L}_t$ ) and a non-targeted part ( $\mathcal{L}_{nt}$ )

$$\mathcal{L} = \mathcal{L}_t + \alpha \mathcal{L}_{nt}$$

Objective of the targeted part: Fool the network by shifting the prediction from the source class into the sink class

This is broken down into two subtasks:

- (1) Decreasing the logit value for the originally predicted class  $L_p$  to not being the highest logit anymore
- (2) Increasing the logit for the sink class  $L_{\text{sink}}$ , to be the highest logit

$$\mathcal{L}_t = \mathcal{L}_{t1} + \mathcal{L}_{t2}, \text{ with}$$

$$\mathcal{L}_{t1} = \max(\hat{L}_p(x_t + \delta) - \max_{i \neq p}(\hat{L}_i(x_t + \delta)), 0) \quad p = \arg \max(\hat{L}(x_t))$$

$$\mathcal{L}_{t2} = \max(\max_{i \neq y_{\text{sink}}}(\hat{L}_i(x_t + \delta) - \hat{L}_{\text{sink}}(x_t + \delta)), -D)$$

Objective of the non-targeted part: Achieve good classification performance

Adopt the cross-entropy function

$$\mathcal{L}_{nt} = \mathcal{X}(\hat{L}(x_{nt} + \delta), \mathbb{1}(\hat{F}(x_{nt})))$$

# Algorithm

A simple, yet effective algorithm

---

## Algorithm 1: Double Targeted Attack Algorithm

---

**Input:** Data distribution  $X$ , Classifier  $\hat{F}$ , Loss function  $\mathcal{L}$ , Mini-batch size  $m$ ,  
Number of iterations  $I$ , Perturbation magnitude  $\epsilon$

**Output:** Perturbation vector  $\delta$

```

 $X_t \subseteq X$                                 ▷ Subset
 $X_{nt} \subseteq X$                             ▷ Subset
 $\delta \leftarrow 0$                           ▷ Initialize
for  $iteration = 1, \dots, I$  do
   $B_t \sim X_t: |B_t| = \frac{m}{2}$                 ▷ Randomly sample
   $B_{nt} \sim X_{nt}: |B_{nt}| = \frac{m}{2}$           ▷ Randomly sample
   $B \leftarrow B_t \cup B_{nt}$                     ▷ Concatenate
   $g_\delta \leftarrow \mathbb{E}_B[\nabla_\delta \mathcal{L}]$       ▷ Calculate gradient
   $\delta \leftarrow \text{Optim}(g_\delta)$               ▷ Update perturbation
   $\delta \leftarrow \frac{\delta}{\|\delta\|_p} \epsilon$     ▷ Projection
end

```

---

# Quantitative Results

**DTA achieves reasonable performance for different mapping scenarios on a wide range of datasets.**

1. The targeted fooling ratio for the targeted classes ( $\kappa_t$ ) is reasonably high.
2. A significant gap between  $\kappa_t$  and  $\kappa_{nt}$  indicates that the crafted perturbation is discriminative between targeted class and non-targeted classes

Table: Experimental results for the Double Targeted Attack (DTA) for the datasets CIFAR-10, GTSRB, EuroSAT, YCB and ImageNet under 10 scenarios  $S_0$  to  $S_9$ . All numbers are reported in %.

Dataset	Model	$S_0$		$S_1$		$S_2$		$S_3$		$S_4$		$S_5$		$S_6$		$S_7$		$S_8$		$S_9$		Avg	
		$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$
CIFAR-10	VGG-16	77.5	20.5	83.5	22.0	78.2	14.7	81.4	21.5	73.0	18.6	79.1	14.2	75.1	15.1	76.7	24.6	75.0	20.3	86.2	16.6	78.6	18.8
	ResNet-20	78.8	26.1	84.6	28.0	84.0	24.3	84.2	26.9	77.1	22.0	82.1	21.3	83.8	14.7	72.9	33.2	80.0	27.8	89.8	22.5	81.7	24.7
GTSRB	VGG-16	89.0	0.2	100	1.1	87.1	1.2	72.2	0.6	91.0	1.3	83.6	2.4	88.3	1.1	80.0	0.7	95.0	1.9	81.1	1.7	86.7	1.2
	ResNet-20	84.3	0.5	100	1.6	53.1	0.2	77.8	1.8	87.6	2.9	77.1	4.4	70.0	2.7	88.3	1.2	80.0	0.3	64.4	0.7	78.3	1.6
EuroSAT	ResNet-50	96.2	33.0	98.8	18.0	95.2	31.1	96.6	22.1	99.2	28.7	95.0	24.0	94.4	44.3	96.3	17.6	96.3	24.5	91.2	22.7	95.9	26.6
	Inception-V3	94.3	28.7	95.2	18.9	93.8	41.4	99.2	56.3	93.0	29.4	93.0	24.2	91.6	34.6	96.0	21.8	96.8	31.6	89.2	18.8	94.2	30.6
YCB	ResNet-50	100	14.5	100	24.2	100	32.4	96.7	38.0	100	33.5	99.2	38.3	100	44.4	99.2	41.7	100	19.0	100	33.1	99.5	31.9
	Inception-V3	100	16.6	100	30.0	100	38.7	99.2	31.2	100	12.9	98.3	20.0	100	32.2	100	36.6	100	17.3	100	39.2	99.8	27.5
ImageNet	VGG-16	72.0	10.3	96.0	19.5	90.0	19.5	82.0	28.3	74.0	15.9	82.0	13.0	66.0	8.9	64.0	12.9	66.0	21.5	70.0	26.1	76.2	17.6
	ResNet-50	74.0	13.9	94.0	21.4	82.0	15.2	72.0	20.9	62.0	13.6	84.0	15.5	72.0	9.8	66.0	21.4	66.0	17.3	62.0	18.1	73.4	16.7
	Inception-V3	78.0	10.0	86.0	15.7	86.0	12.2	78.0	15.6	58.0	9.5	76.0	12.9	70.0	8.9	72.0	15.7	62.0	18.9	66.0	17.8	73.2	13.7
	MobileNet-V2	74.0	11.3	94.0	17.0	88.0	20.4	70.0	15.3	72.0	16.0	84.0	15.0	74.0	14.5	74.0	21.7	72.0	18.8	70.0	21.9	77.2	17.2

$\kappa_t$ : targeted fooling ratio for the targeted source samples

$\kappa_{nt}$ : targeted fooling ratio for the non-targeted source samples 39

# Multi2One attack

**Attack multiple source classes (MS) to be misclassified as one sink class.**

Table: Experimental results for the universal Multi2One targeted perturbation on ImageNet under 4 scenarios  $MS_0$  to  $MS_3$  attacking three source classes to one sink class. All numbers are reported in %.

Model	$MS_0$		$MS_1$		$MS_2$		$MS_3$		Avg	
	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$
VGG16	63.3	24.9	69.3	33.7	76.0	25.8	69.3	26.8	69.5	27.8
ResNet-50	64.0	30.1	63.3	32.3	78.7	29.2	62.7	23.2	67.2	28.7
Inception-V3	58.0	19.4	56.7	23.8	66.7	19.0	66.7	20.8	62.0	20.8
MobileNet-V2	68.0	27.2	66.0	28.0	74.0	25.6	66.0	24.4	68.5	26.3



# Qualitative Results

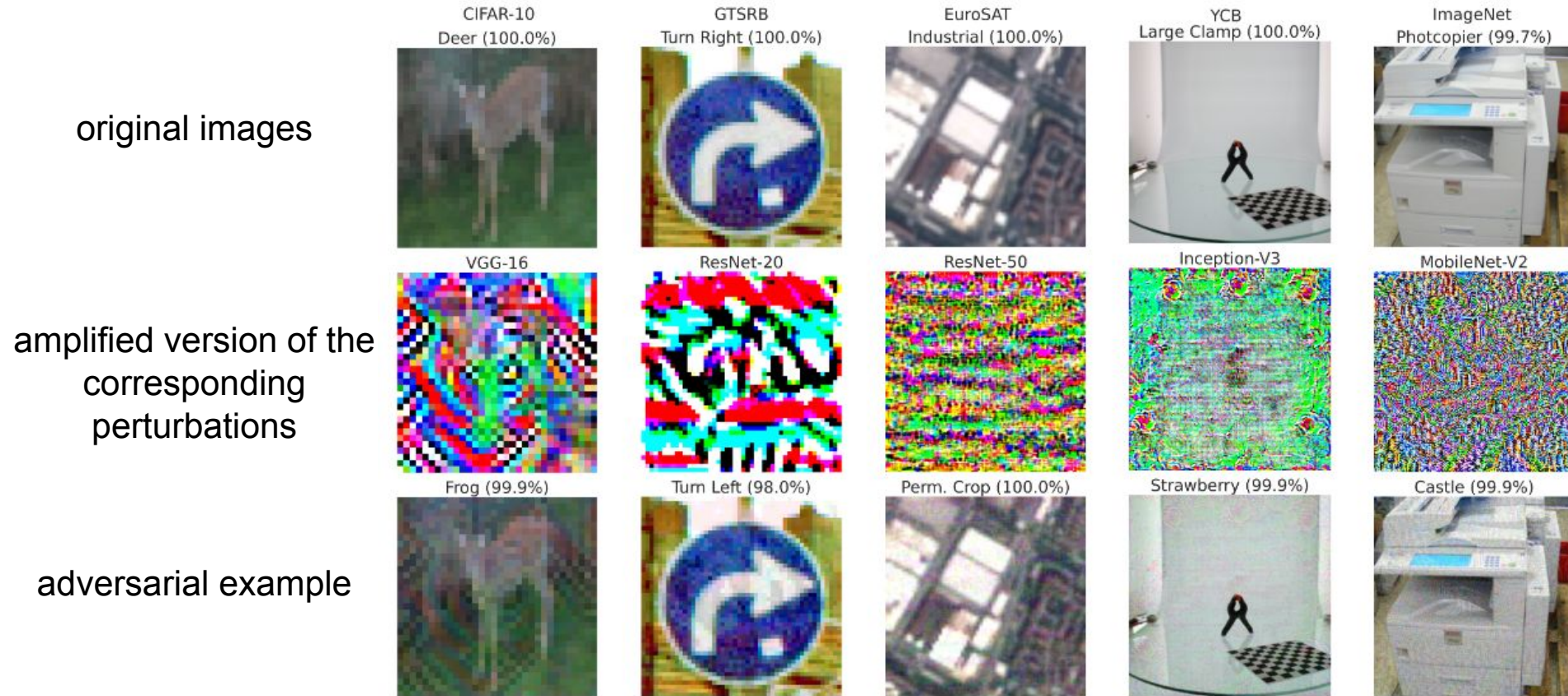


Figure: Examples of adversarial perturbations for various datasets and networks. The confidence values of the network and the predicted labels are stated above the images. The target network is indicated above the amplified perturbation.

# Physical Attack - Double Targeted Patch

Table: Quantitative results for the generated DT-Patch on ImageNet

Hammer → Hummingbird		Screwdriver → Go-Kart		Coffee Mug → Chocolate Sauce	
$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$	$\kappa_t$	$\kappa_{nt}$
80.0	42.7	92.0	44.9	96.0	41.6

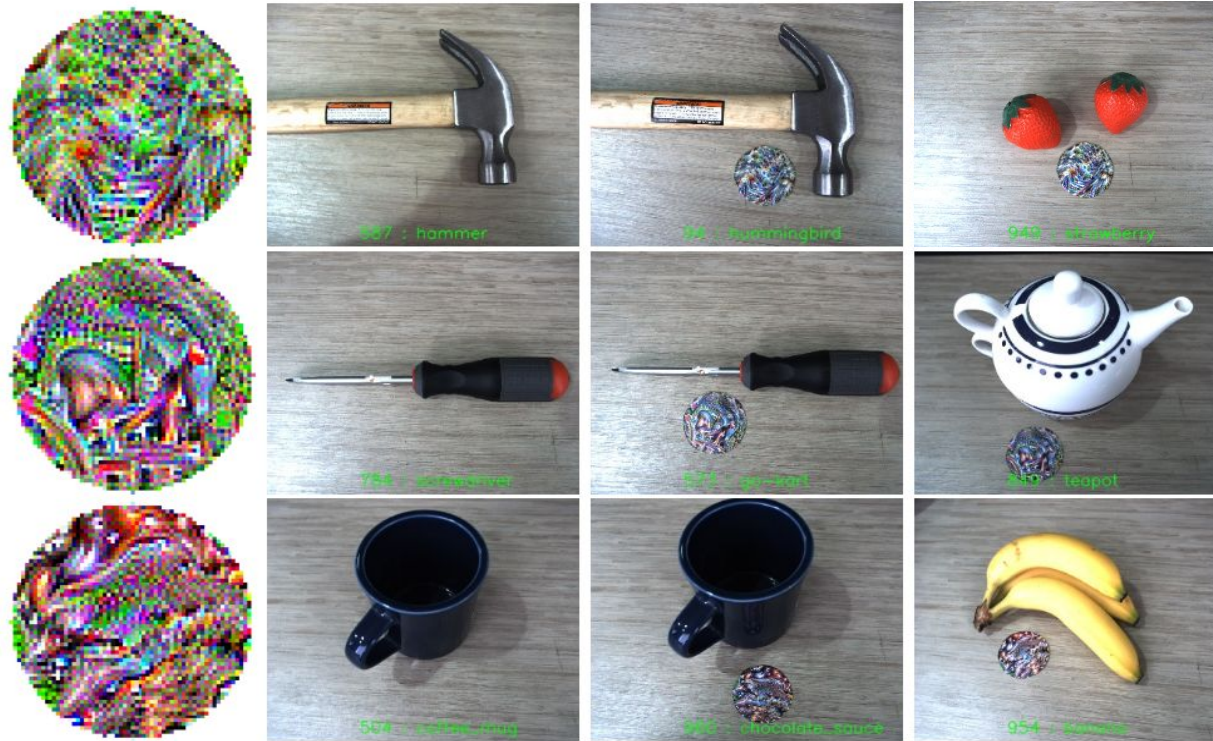


Figure: Real-world examples of the DT-Patch for the three different scenarios

# Understanding Adversarial Examples from the Mutual Influence of Images and Perturbations (CVPR 2020)

Chaoning Zhang\*, Philipp Benz\*, Tooba Imtiaz, In-So Kweon  
\* indicates equal contribution

Korea Advanced Institute of Science and Technology (KAIST)

# Background: Adversarial Examples are not Bugs they are Features

Adversarial Examples can be directly attributed to the presence of non-robust features [1]

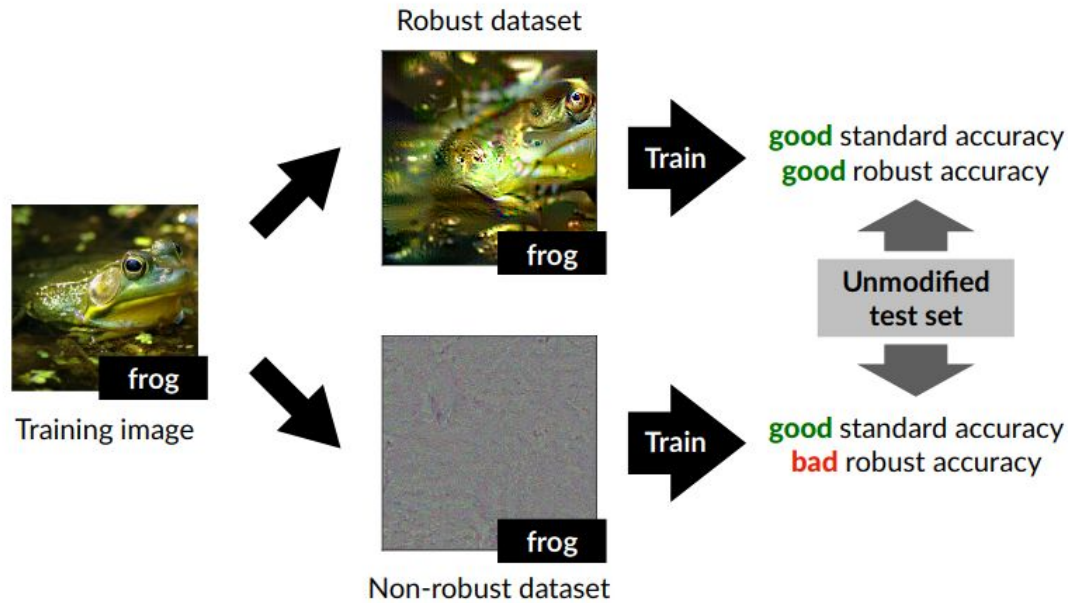


Figure 1: Disentanglement of features into combinations of robust and non-robust features [1]

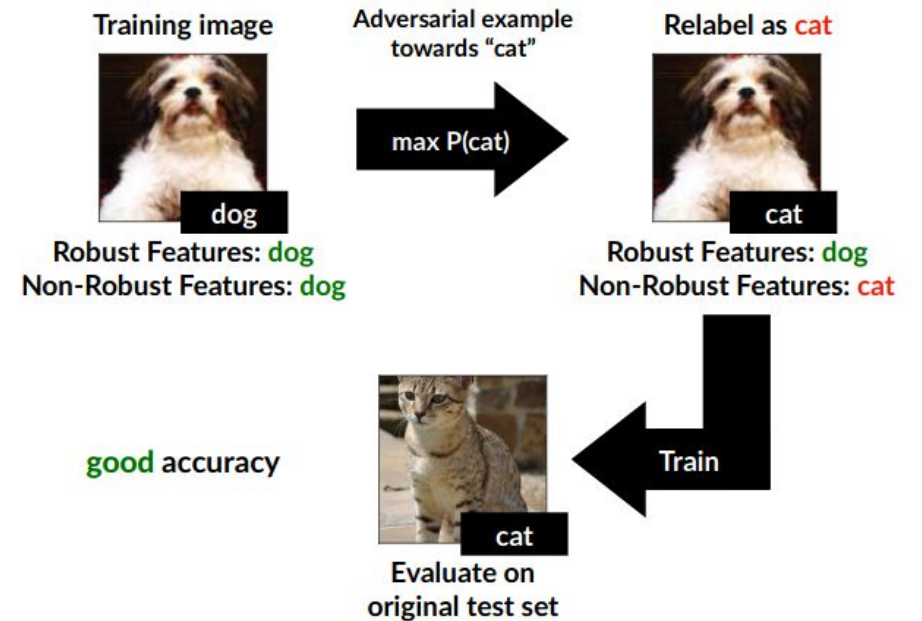


Figure 2: A dataset appearing mislabeled to humans (via adversarial examples) can result in good accuracy on the original test set [1]

# Background: Adversarial Examples are not Bugs they are Features

The existence of adversarial examples can be attributed to the existence of Non-robust features in the dataset.

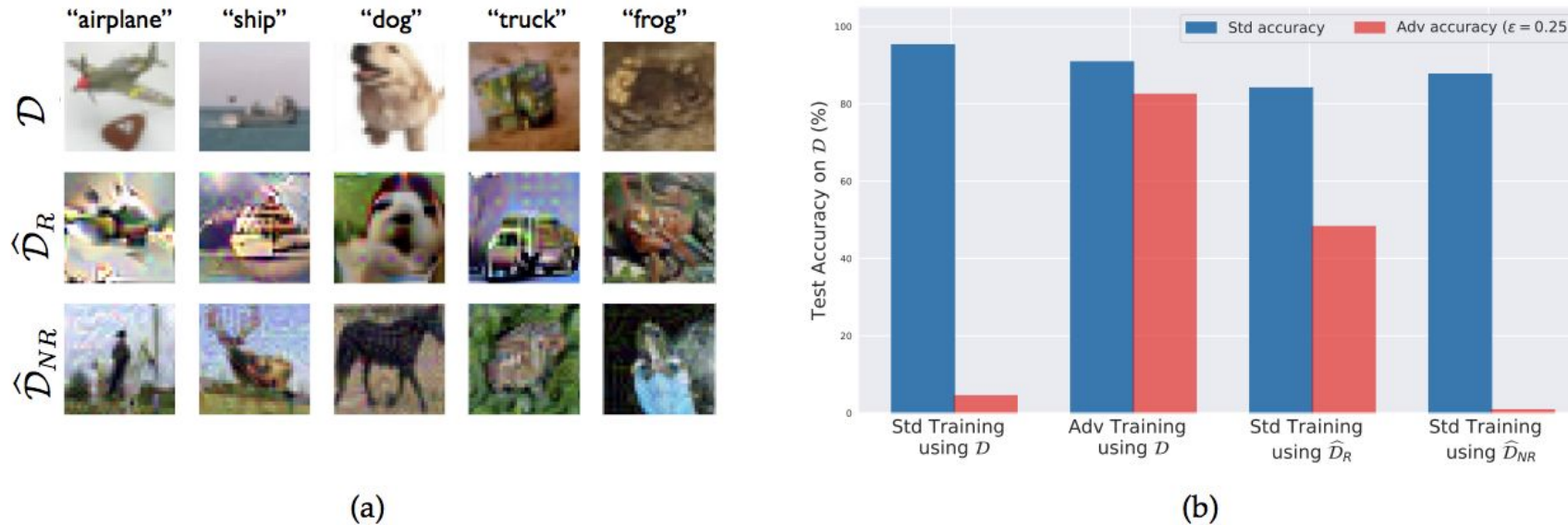
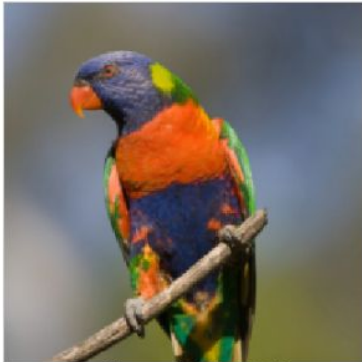


Figure 2: **Left:** Random samples from our variants of the CIFAR-10 [Kri09] training set: the original training set; the *robust training set*  $\hat{\mathcal{D}}_R$ , restricted to features used by a robust model; and the *non-robust training set*  $\hat{\mathcal{D}}_{NR}$ , restricted to features relevant to a standard model (labels appear incorrect to humans). **Right:** Standard and robust accuracy on the CIFAR-10 test set ( $\mathcal{D}$ ) for models trained with: (i) standard training (on  $\mathcal{D}$ ); (ii) standard training on  $\hat{\mathcal{D}}_{NR}$ ; (iii) adversarial training (on  $\mathcal{D}$ ); and (iv) standard training on  $\hat{\mathcal{D}}_R$ . Models trained on  $\hat{\mathcal{D}}_R$  and  $\hat{\mathcal{D}}_{NR}$  reflect the original models used to create them: notably, standard training on  $\hat{\mathcal{D}}_R$  yields nontrivial robust accuracy. Results for Restricted-ImageNet [Tsi+19] are in D.8 Figure 12.

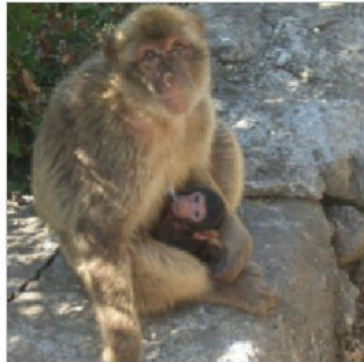
# PCC Analysis

The Pearson correlation coefficient (PCC) is a widely adopted metric to measure the linear correlation between two variables.

$$PCC_{X,Y} = \frac{cov(X, Y)}{\sigma_X \sigma_Y}$$



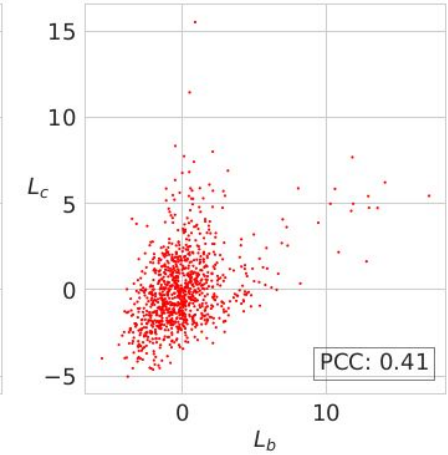
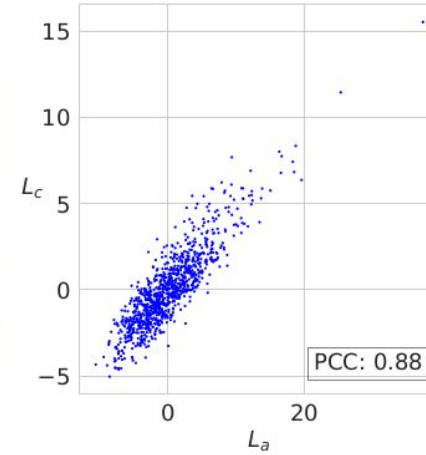
lorikeet (100.00%)



macaque (87.98%)



lorikeet (97.86%)



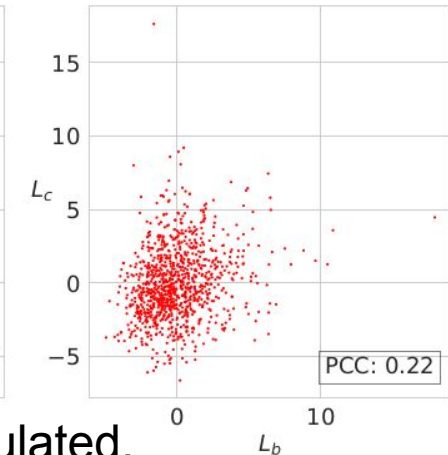
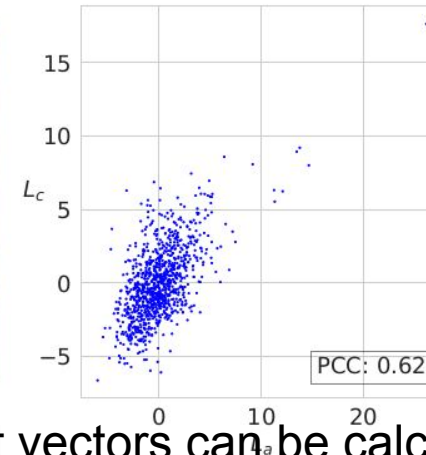
space shuttle (100.00%)



bee (99.79%)



space shuttle (99.90%)

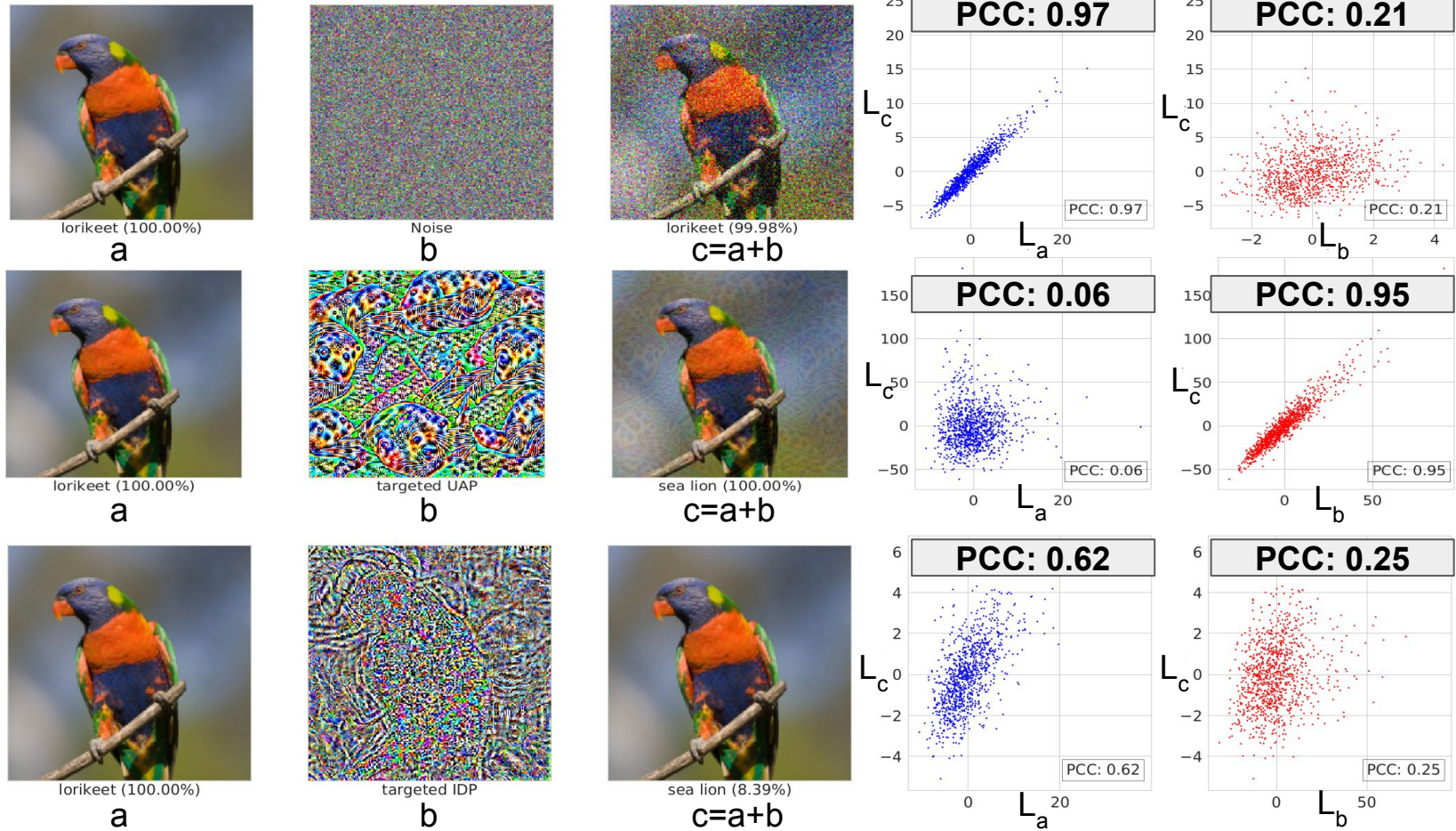


The PCC between different DNN output logit vectors can be calculated.

This provides insight about the “stronger” contribution of the two inputs to the final outcome. A higher PCC value indicates the more significant contributor.

# PCC Analysis

Treat the DNN logits as a vector for feature representation and use them to analyze the mutual influence of two independent inputs based on the Pearson correlation coefficient (PCC)



“Universal perturbations contain dominant features, and images behave like noise to them”

Image-dependant perturbations seem to not contain features by themselves

PCC-Analysis result for one sample image `lorikeet'. Three scenarios of input combinations are considered: 1: image + noise; 2: image + targeted UAP; 3: image + targeted image-dependant AE. The columns show input a, input b, input c=a+b, logit vector analysis of  $L_c$  over  $L_a$  and vector analysis of  $L_c$  over  $L_b$

# Noise Perspective vs. Feature Perspective

## Noise Perspective (Prior works)

- Treat the targeted UAP as noise (“bug”) to the sample to be attacked
- Requires the samples from the training dataset in the UAP generation process
- Explicitly designed to let individual samples cross the decision boundary
- Assumes that the attack generalizes to unseen samples

**Requires the original training dataset**  
**Slow: ~2 hours**

## Feature Perspective (Ours)

- **UAPs contain features of a certain class**
- Treatment of the **images as noise** to the generated UAP during the optimization process in order to be recognizable by the target network
- **No need for semantic features** as in the original training dataset samples
- **Proxy datasets** as background noise: Downloaded from the Internet, MS-COCO, Pascal VOC, Places365

**Requires no original training dataset**  
**Fast: ~2 minutes**

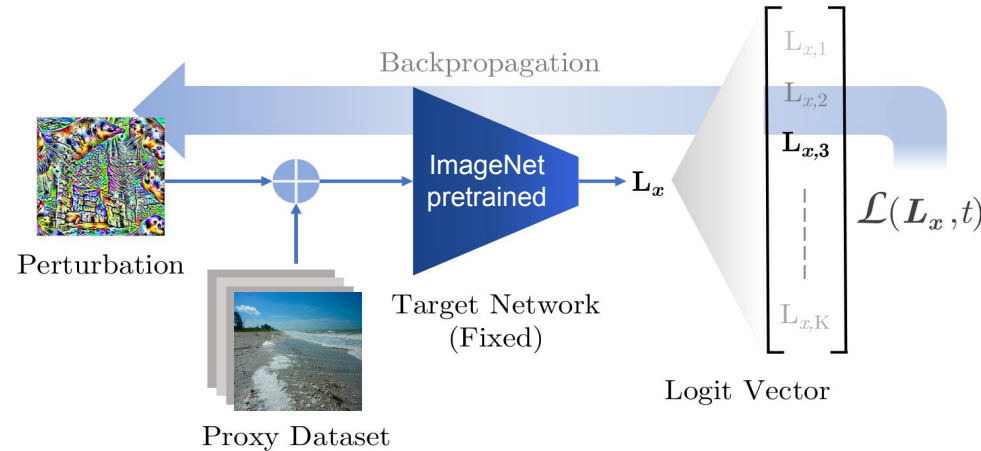


# Data-free Targeted UAP

If images just behave like noise to the features in UAPs. Can we just use another dataset to craft UAPs? Yes we can!

## UAP generation stage

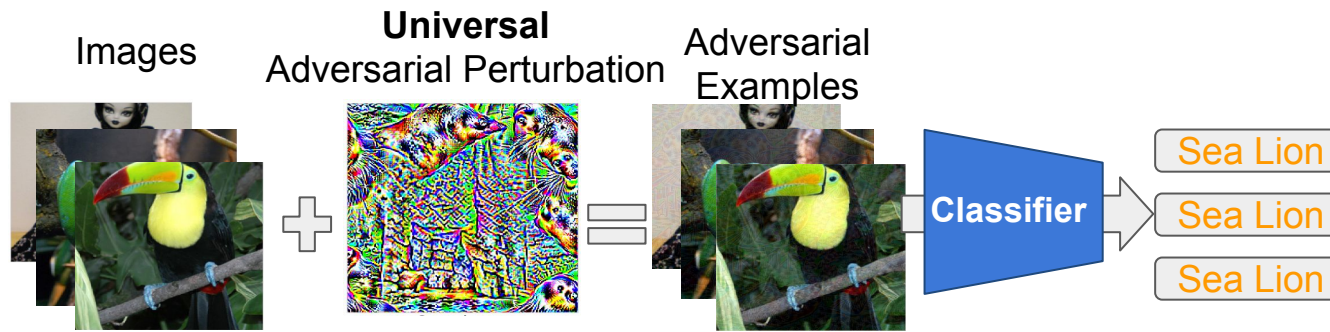
Craft the perturbation with the help of some proxy dataset



As a loss function you can use the cross-entropy between the output logit and some (randomly chosen) target class.  
For example: Choose target class sea lion

## UAP testing stage

Test the perturbation on the ImageNet validation dataset



Is the sample misclassified?

Good

Is the sample correctly classified?

Bad

Figure 1: Inference with adversarial examples formed with a targeted universal adversarial perturbation crafted for target class “sea lion”.

# Data-free Targeted UAP: Results

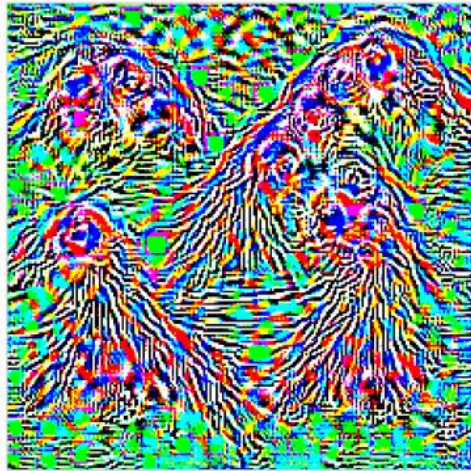
The proposed approach outperforms previous data-dependant and data-free UAP generation methods

Table 2: Comparison to other methods. The results are divided in universal attacks with access to the original ImageNet training data (upper) and data-free methods (lower). The metric is reported in the non-targeted fooling ratio (%)

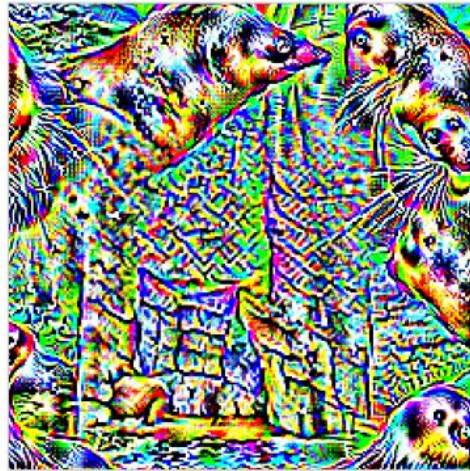
Method	AlexNet <sup>1</sup>	GoogleNet	VGG16	VGG19	ResNet152
UAP [14]	93.3	78.9	78.3	77.8	84.0
GAP [19]	-	82.7	83.7	80.1	-
Ours(ImageNet [11])	<b>96.17</b>	<b>88.94</b>	<b>94.30</b>	<b>94.98</b>	<b>90.08</b>
FFF [18]	80.92	56.44	47.10	43.62	-
AAA [21]	89.04	75.28	71.59	72.84	60.72
GD-UAP [17]	87.02	71.44	63.08	64.67	37.3
Ours (COCO [12])	89.9	<b>76.8</b>	<b>92.2</b>	<b>91.6</b>	<b>79.9</b>
Ours (VOC [5])	89.9	76.7	<b>92.2</b>	90.5	79.1
Ours (Places365 [29])	<b>90.0</b>	76.4	92.1	91.5	78.0

# Can you see the features?

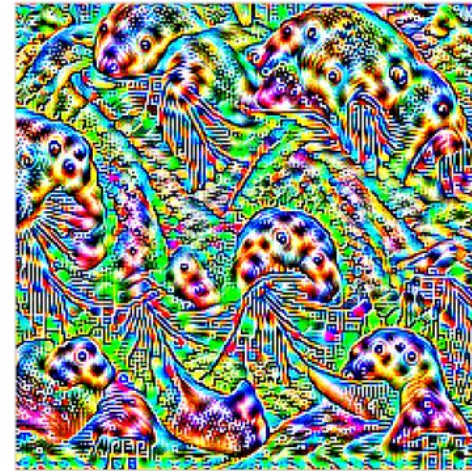
The UAPs were generated for target class sea lion. Sea lion features become visible in the UAPs!



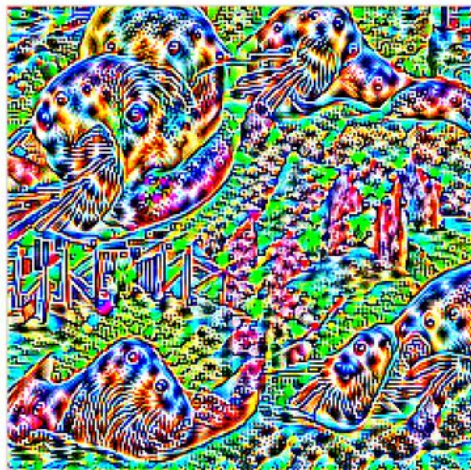
AlexNet



GoogleNet



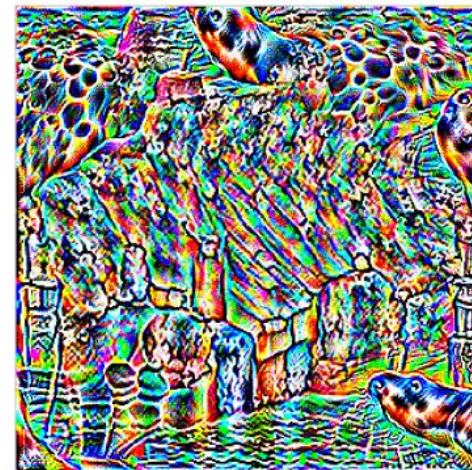
VGG16



VGG19



ResNet152



InceptionV3

# UDH: Universal Deep Hiding for Steganography, Watermarking and Light Field Messaging

Chaoning Zhang\*, Philipp Benz\*, Adil Karjauv\*, Geng Sun, In-So Kweon  
\* indicates equal contribution

Korea Advanced Institute of Science and Technology (KAIST)

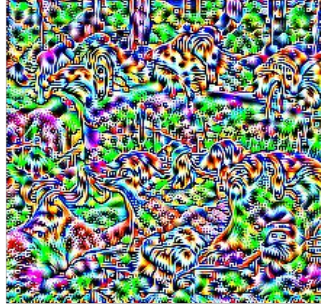
# What you see is not what the machine sees

Universal  
Adversarial  
Attack

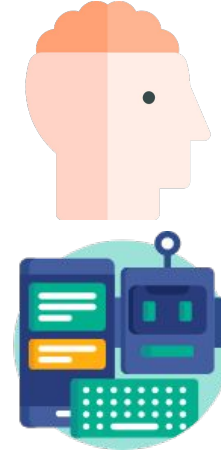
Original Image



Perturbation  
(Cat feature)



Perturbed Image



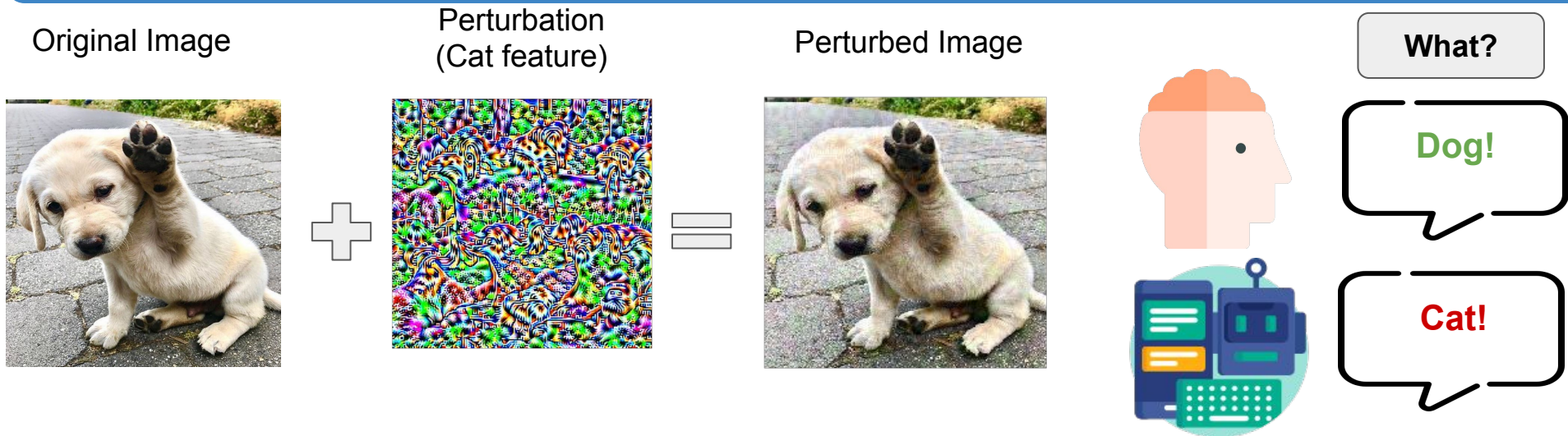
What?

Dog!

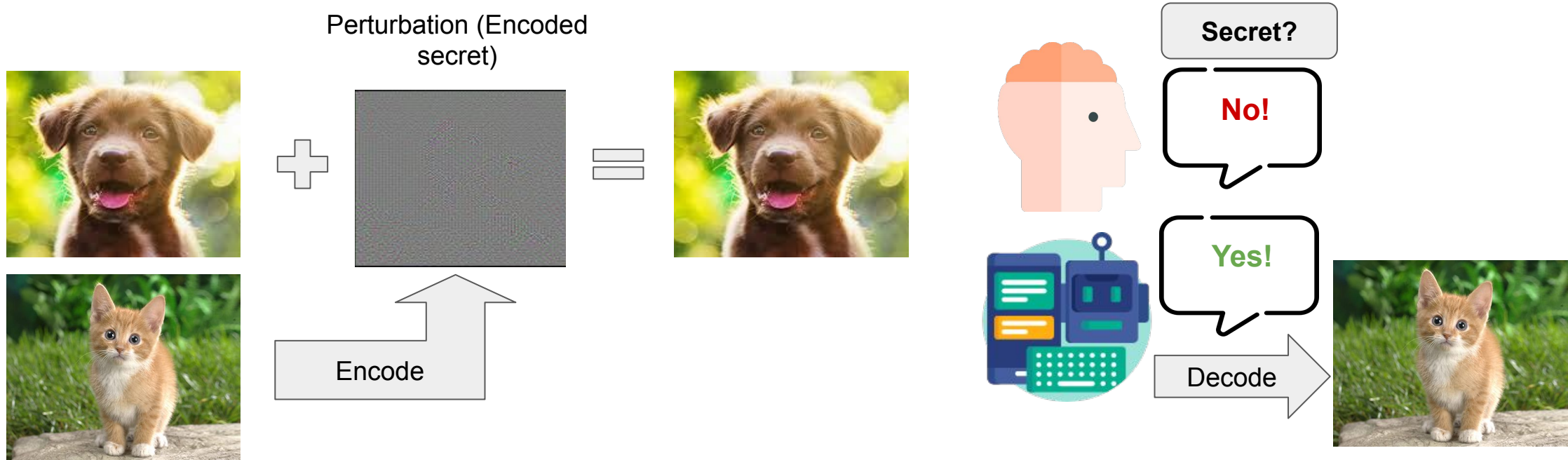
Cat!

# What you see is not what the machine sees

Universal Adversarial Attack



Universal Deep Hiding



## Digital Steganography

- Definition: Hiding secret message within a digital file, such as images or videos.
- Challenge: **Trade-off between Capacity vs. secrecy**

Secret Message

0 1 0 0  
0 0 0 0  
1 0 1 1



Cover Image

Encoding



Container Image

Decoding

Recovered  
secret message

0 1 0 0  
0 0 0 0  
1 0 1 1

# Problems of existing solutions

## Traditional steganography vs. Deep steganography

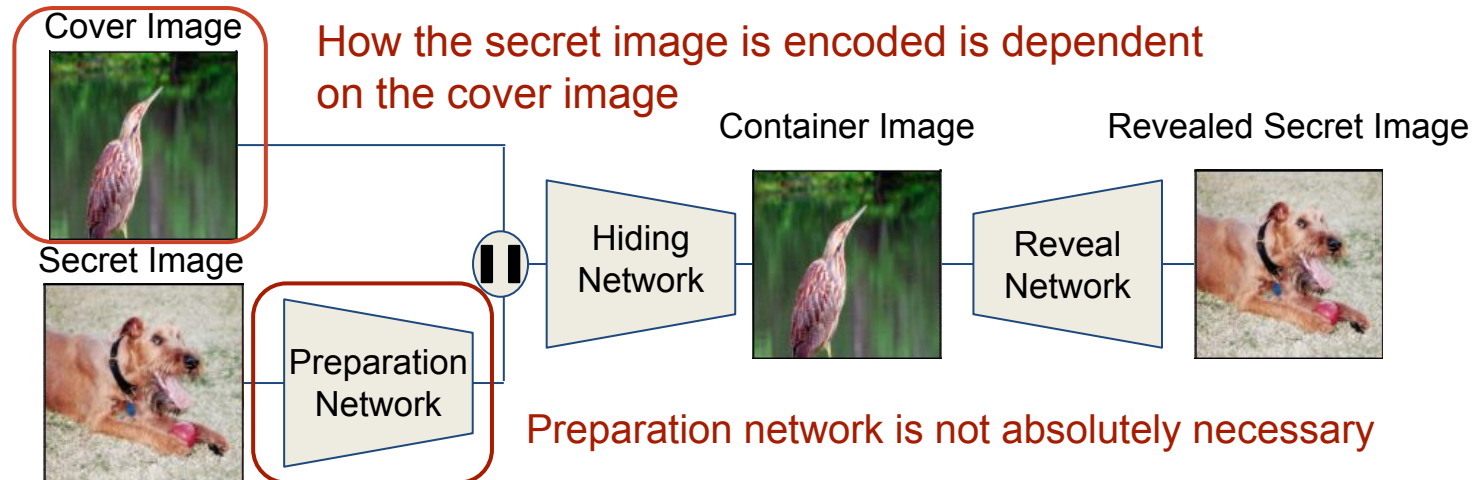
### Traditional Steganography

Secret: Binary message (<0.5 bits per pixel)  
 Requires perfect encoding/decoding  
 Evading steganalysis is important

### Deep Steganography [1]

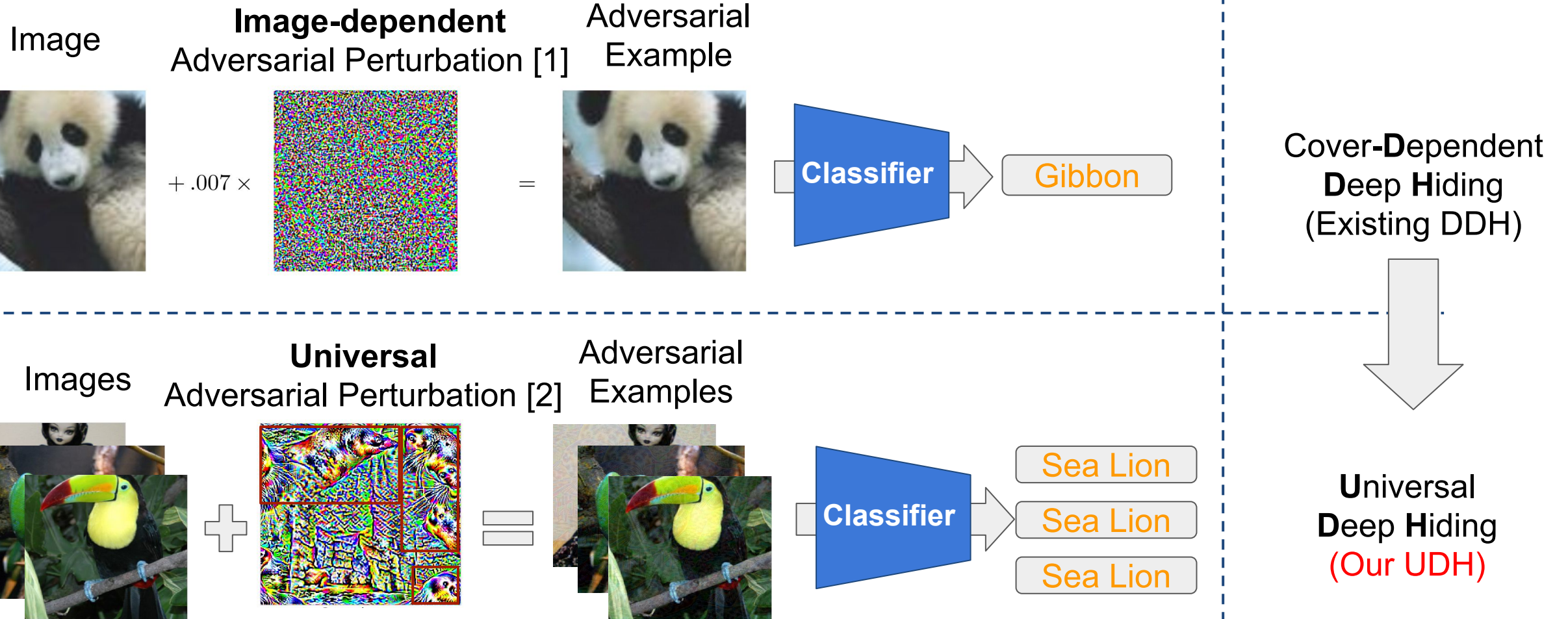
Secret: a full image (24 bits per pixel)  
 Highlighting visual quality  
 Inevitable to be detected due to high capacity

## Hiding Images in Plain Sight: Deep Steganography [1]





Deep Neural Networks are sensitive to invisible perturbations, leading to misclassifications.  
 Universal adversarial perturbations exist to fool most images.

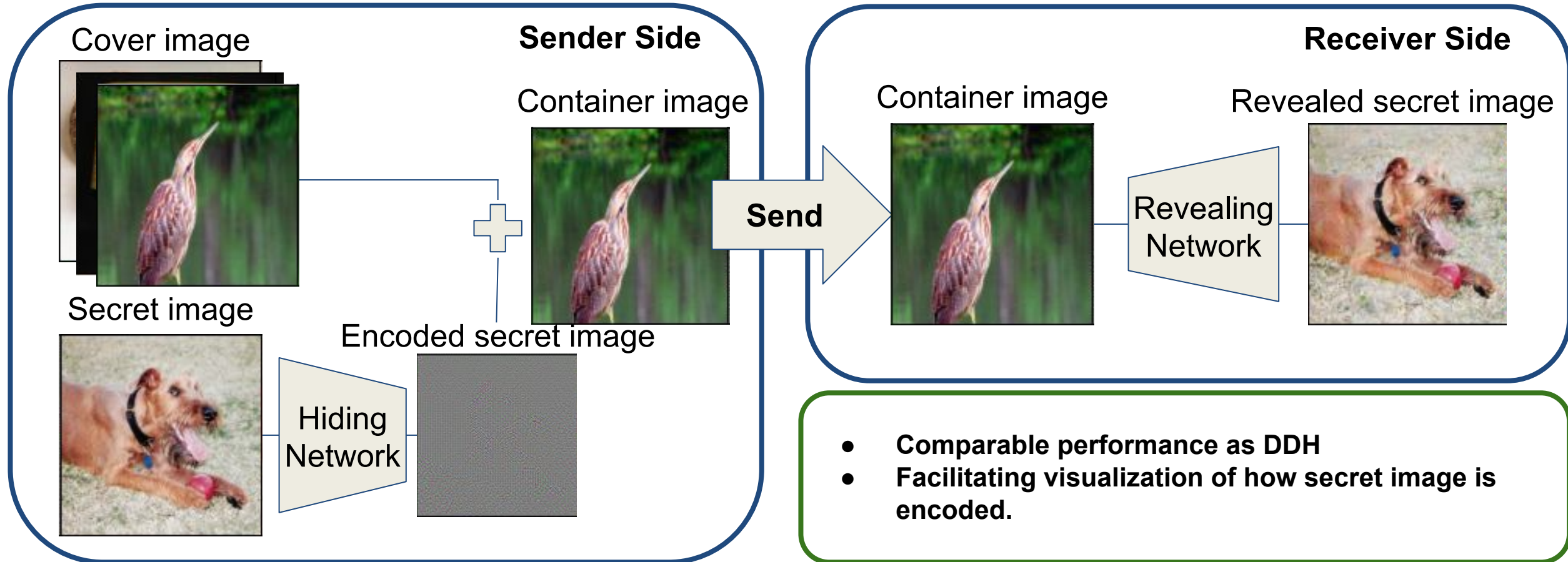


[1] Explaining and Harnessing Adversarial Examples; Goodfellow, Shlens, Szegedy; ICLR 2015

[2] Understanding Adversarial Examples from the Mutual Influence of Images and Perturbations; Zhang\*, Benz\*, Imtiaz, Kweon; CVPR 2020

## Universal Deep Hiding (UDH)

A novel Deep Hiding meta-architecture to hide secret image in a **cover-agnostic** manner.



A secret image is fed to the hiding network to yield an encoded secret image, which can be added to a random cover image to form a container image. The revealing network then retrieves the secret image from the container.

The Hiding network encodes secret image in high-frequency representation.  
 The reveal network easily extracts high-frequency encoded secret image from container

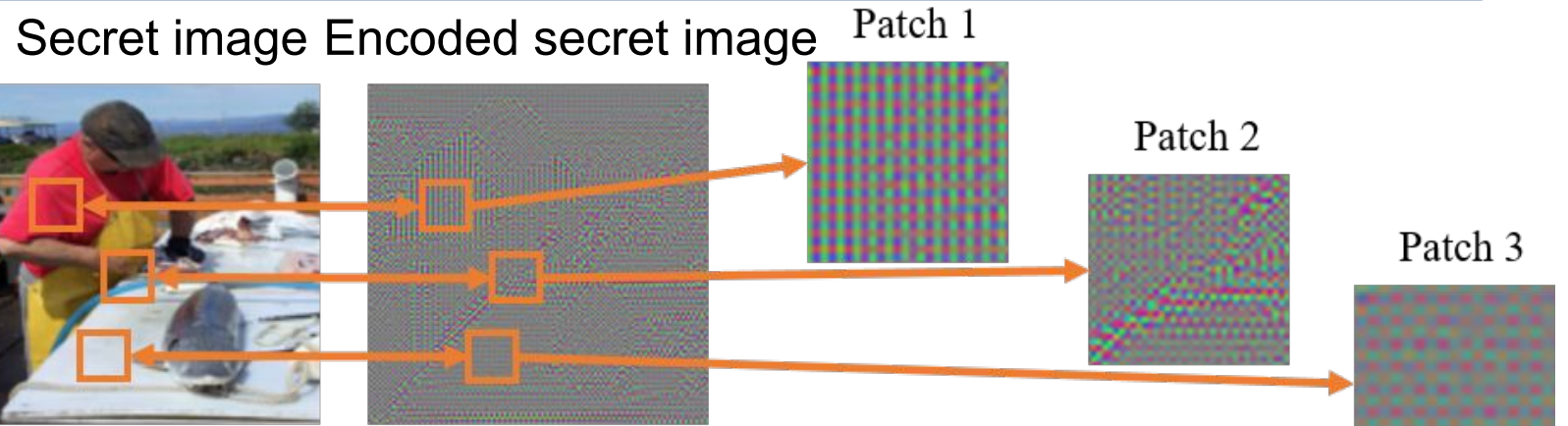
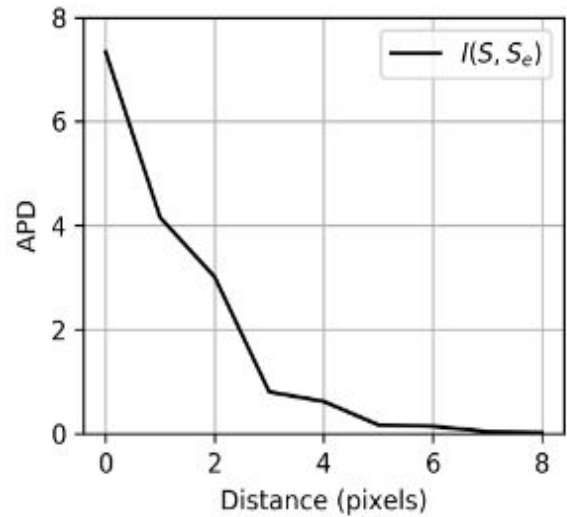


Figure 2: A secret image and its corresponding encoded secret image with zoomed patches for the **UDH** meta architecture.

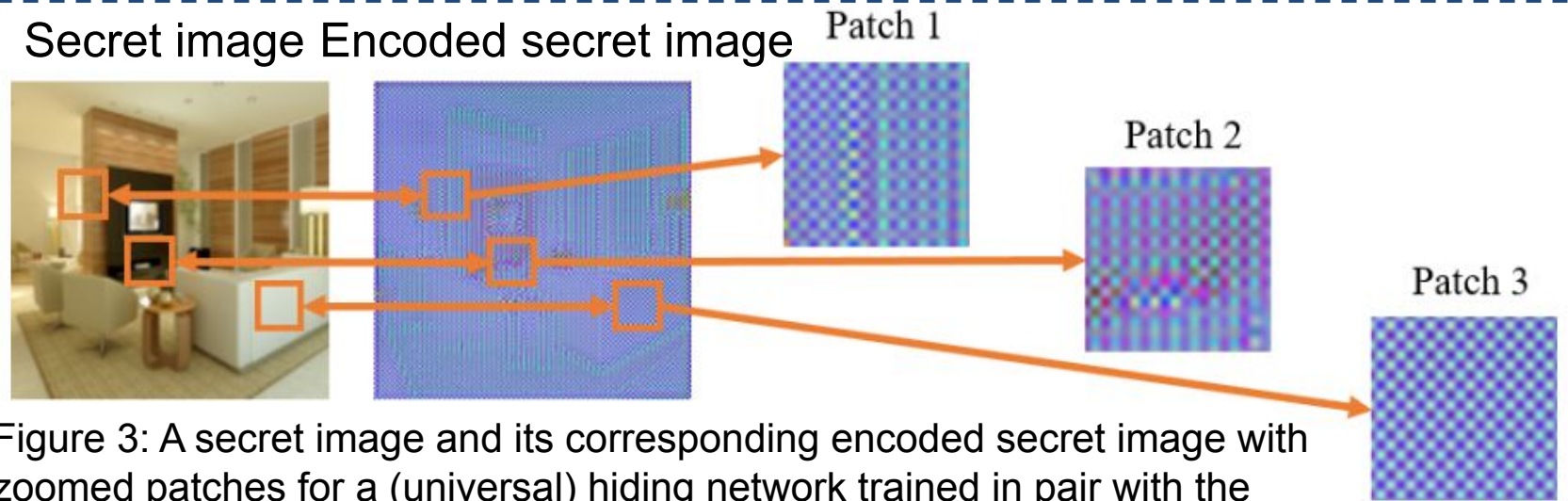


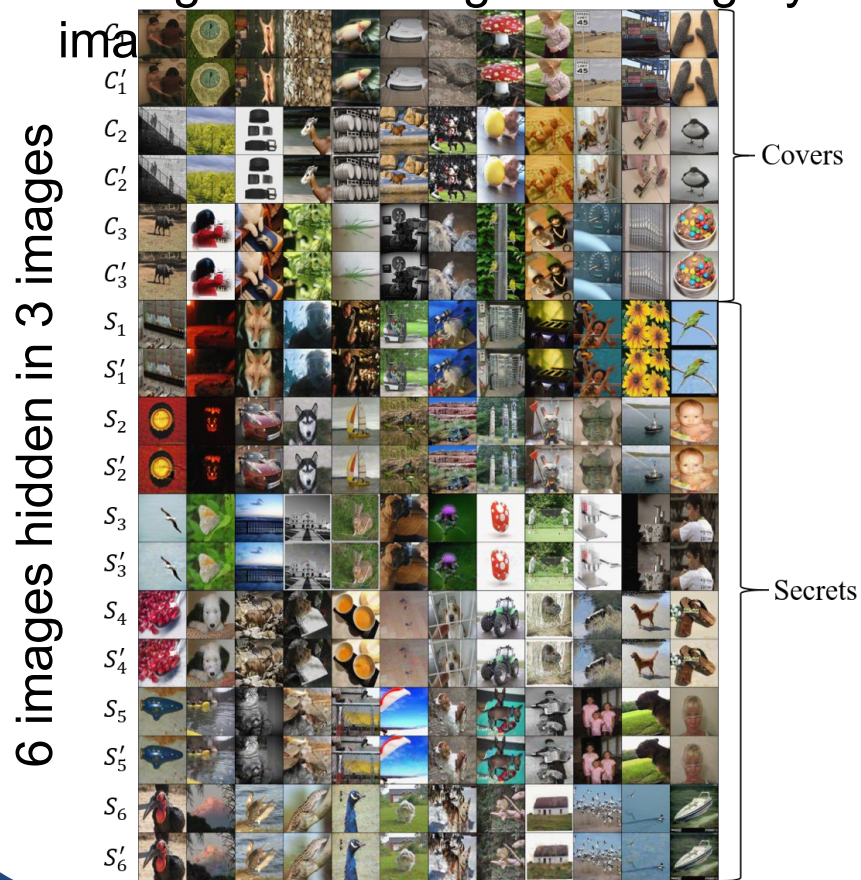
Figure 3: A secret image and its corresponding encoded secret image with zoomed patches for a (universal) hiding network trained in pair with the **fixed** reveal network from the **DDH** meta architecture.

[1] UDH: Universal Deep Hiding for Steganography, Watermarking and Light Field Messaging. Zhang\*, Benz\*, Karjauv\*, Sun, Kweon, NeurIPS 2020.

# Beyond hiding one image in one image

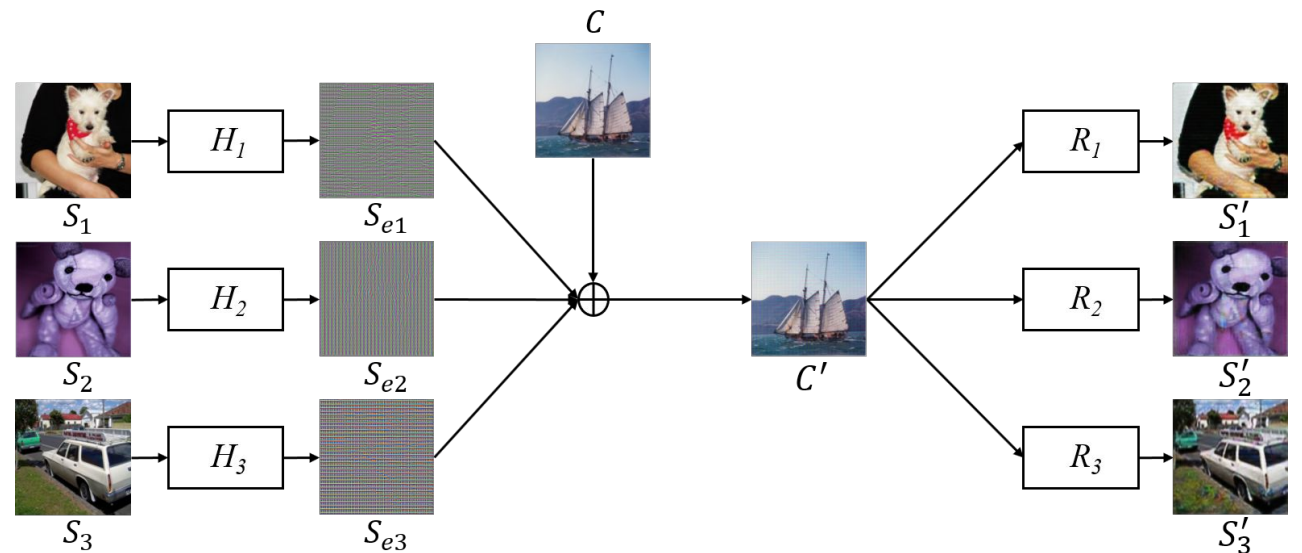
## Hiding $M$ Images in $N$ Images

- Hiding many images in few images
- Hiding 2 color images in one gray



## Different Secrets for Different Recipients

Hiding multiple secret images in one cover so that different recipients retrieve their respective secrets



# Beyond Digital Steganography

## Watermarking

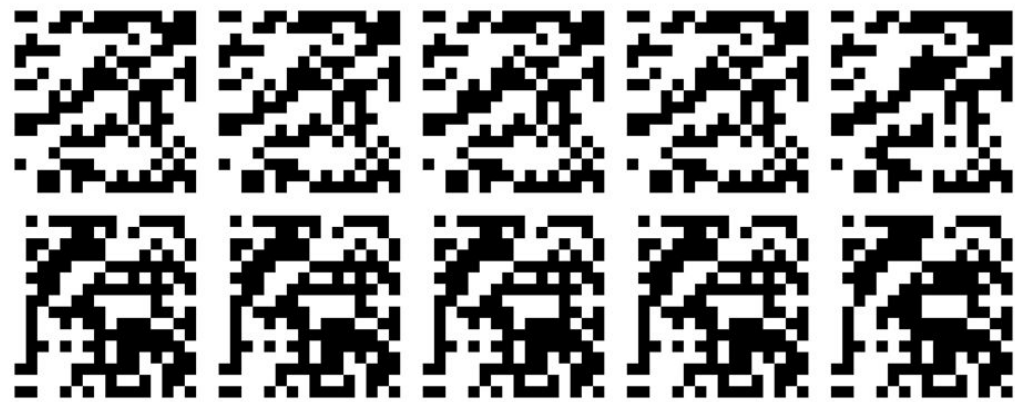
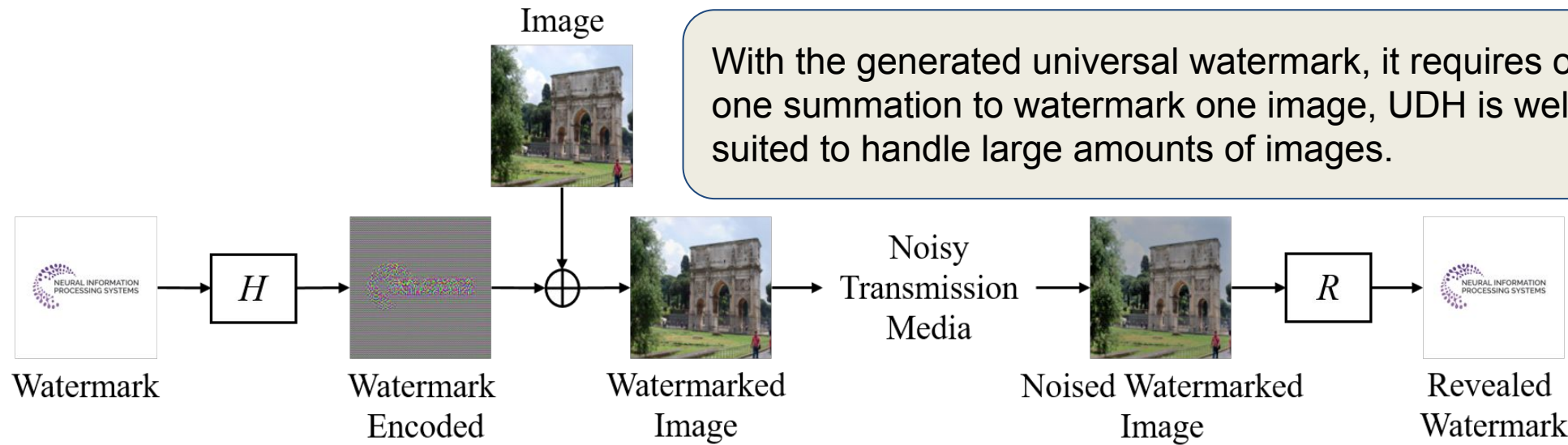
- Definition: Hiding secret message for proving the ownership
- Challenge: **robustness to distortions**

## Photographic steganography or Light Field Messaging (LFM)

- Definition: Hiding and transmitting secret message by taking a photo on screen display
- Challenge: **robustness to light effect**

**With the UDH, we are the **first** to demonstrate the possibility of hiding **a full image** for both applications.**

# Universal Deep Watermarking

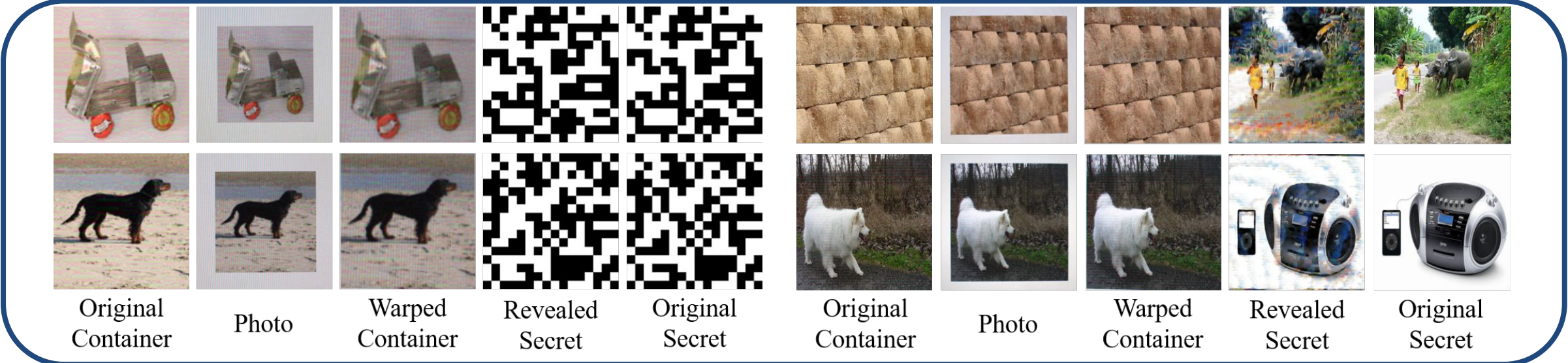
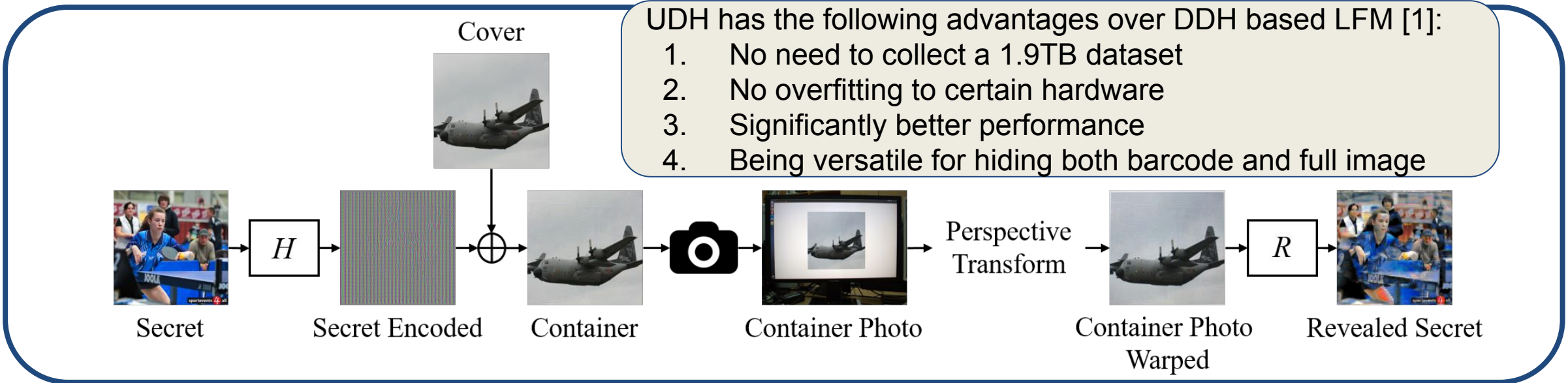


Original Identity Dropout Gaussian JPEG



Original Identity Dropout Gaussian JPEG

# Light Field Messaging (Photographic Steganography)



[1] Light field messaging with deep photographic steganography; Wengrowski, Dana K.; CVPR 2019

[2] UDH: Universal Deep Hiding for Steganography, Watermarking and Light Field Messaging. Zhang\*, Benz\*, Karjauv\*, Sun, Kweon, NeurIPS 2020.

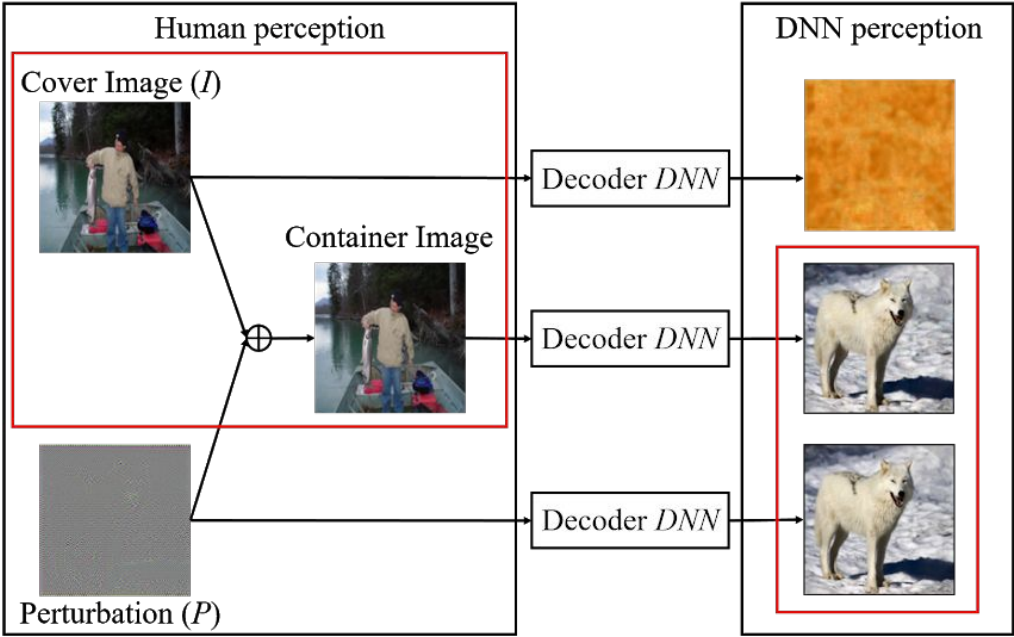
# Universal Adversarial Perturbations Through the Lens of Deep Steganography: Towards a Unified Fourier Perspective (AAAI 2021)

Chaoning Zhang\*, Philipp Benz\*, Adil Karjauv, In-So Kweon  
\* indicates equal contribution

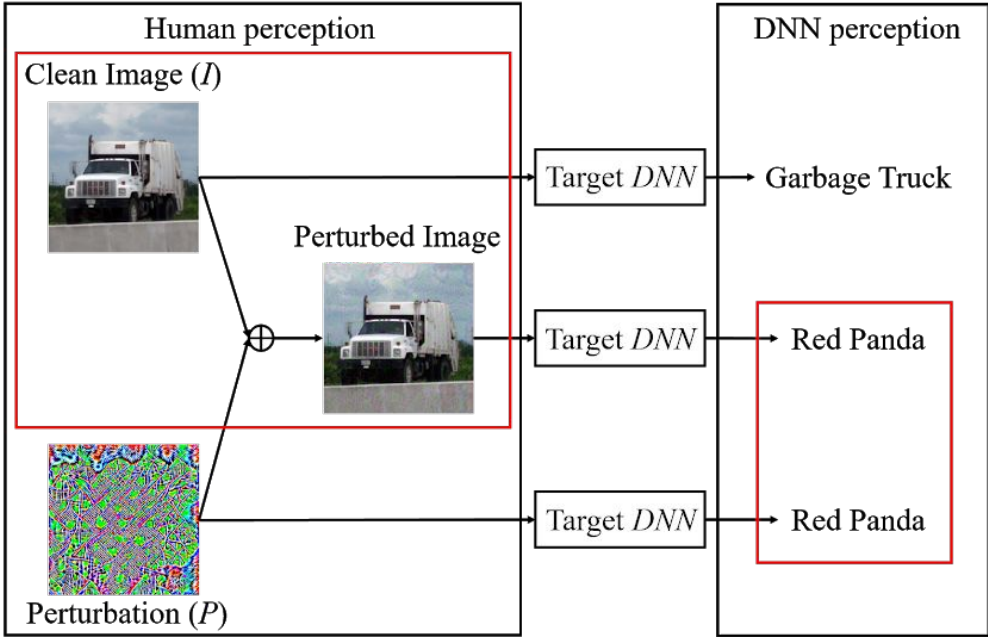
Korea Advanced Institute of Science and Technology (KAIST)



# Universal Perturbations Dominate over Images



**Universal Data Hiding [1]**



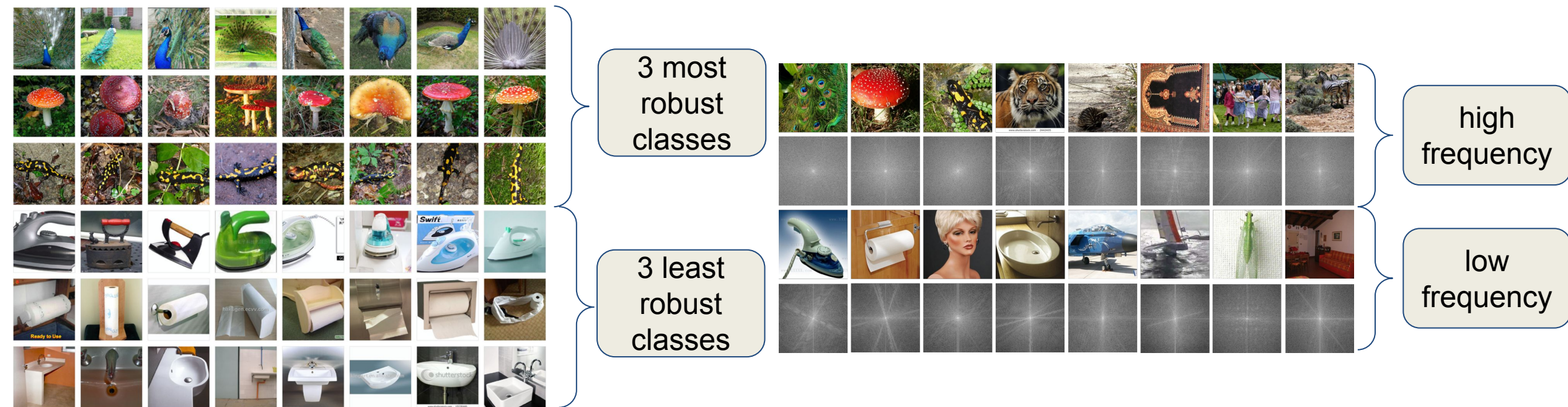
**Universal Adversarial Attack**

There exists a **misalignment** between **human vision** and **model perception** for both universal data hiding and universal adversarial attack

[1] UDH: Universal deep hiding for steganography, watermarking and light field messaging. Zhang\*, Benz\*, Karjauv\*, Sun, Kweon, NeurIPS 2020.

# Universal yet there exists class-wise discrepancy

For robust and non-robust classes, the targeted attack rate is **40%** and **100%** respectively.

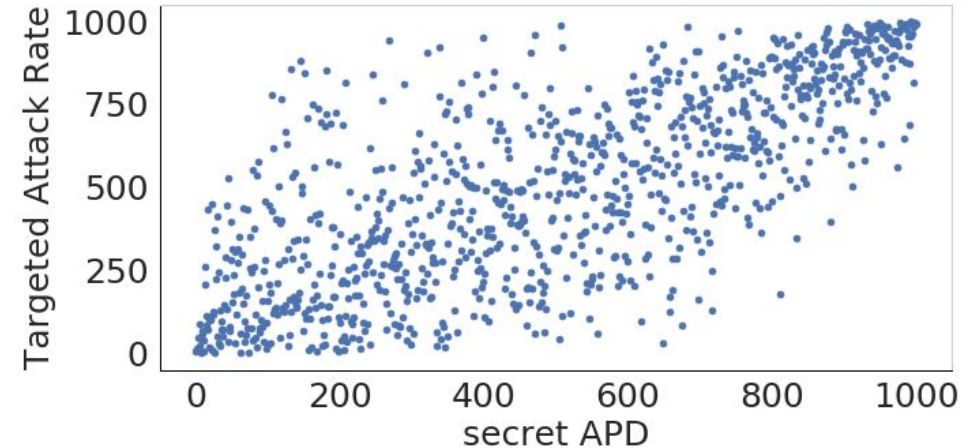


Considering the most and least robust classes against adversarial attack, the **most robust classes** have relatively **high frequency** and the **least robust classes** have relatively **low frequency**

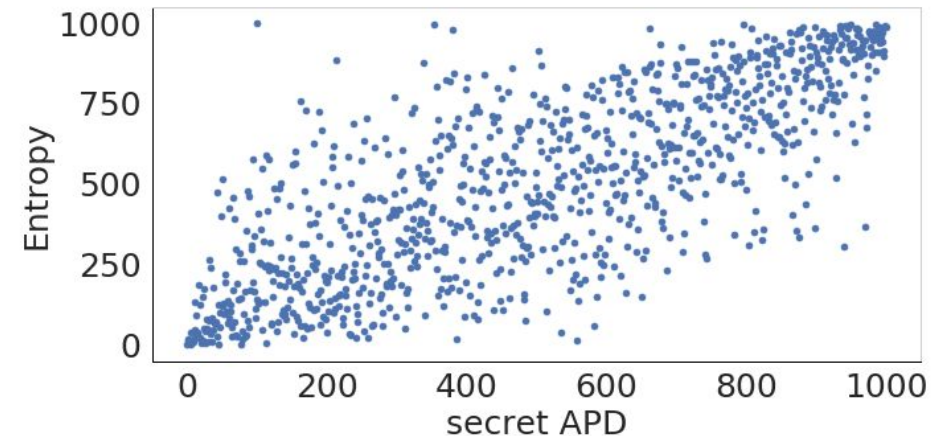
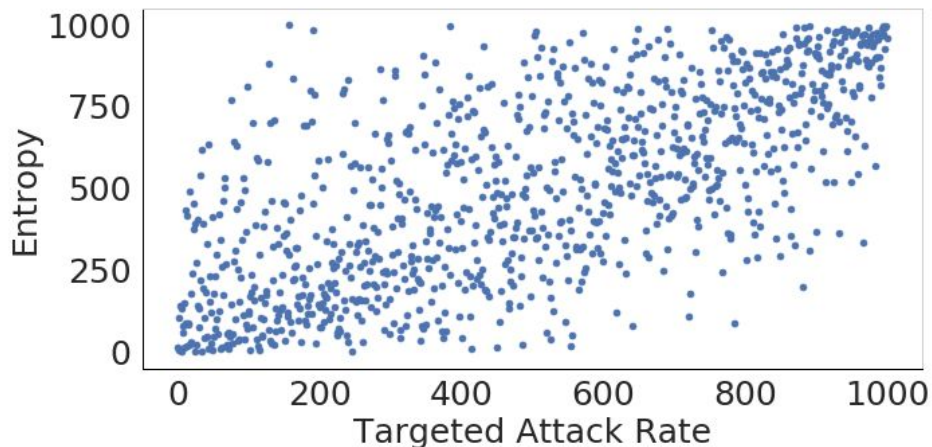
# A Unified Frequency Perspective on UAP and UDH

We order the 1000 classes in ImageNet based on three metrics:

- (a) Targeted attack rate (for UAP)
  - (b) Secret APD (for UDH)
  - (c) Entropy of Fourier images (for **frequency**)
- Surprisingly, we find that the ordering is highly correlated between two different tasks.



What is the factor that determines this high correlation?  
**Frequency!** Images from HF class are systematically more robust to UAP.

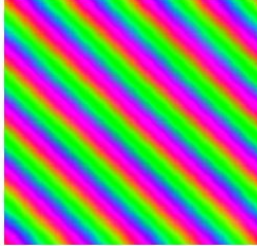


# Why do UAPs dominate images?

BW: 5 (FR: 14.5)



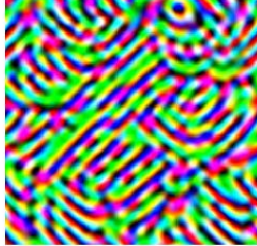
BW: 10 (FR: 39.2)



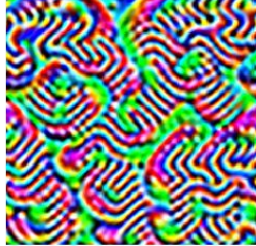
BW: 20 (FR: 47.2)



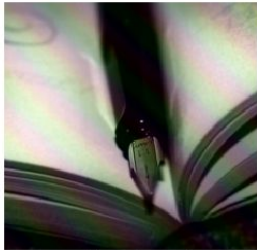
BW: 35 (FR: 56.0)



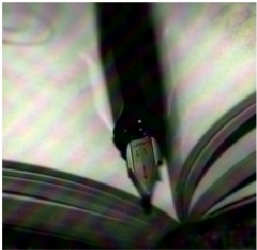
BW: 50 (FR: 64.1)



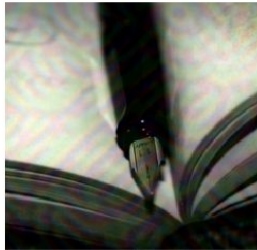
wine\_bottle



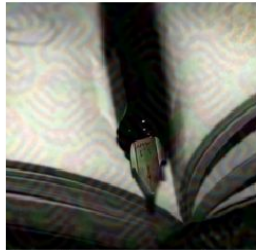
pinwheel



lacewing

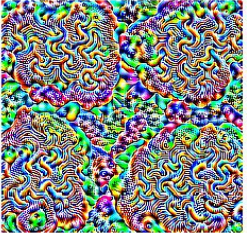


pinwheel

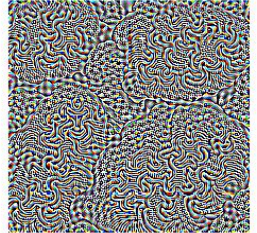


brain\_coral

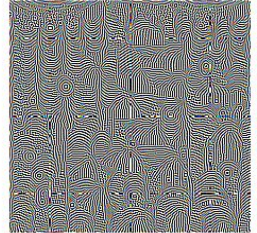
BW: 0 (FR: 94.4)



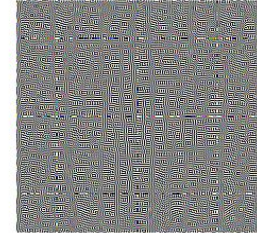
BW: 60 (FR: 90.1)



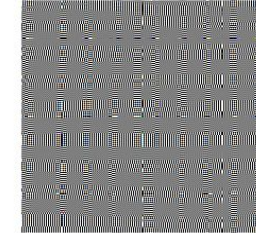
BW: 140 (FR: 85.0)



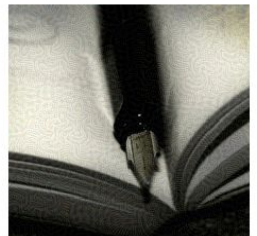
BW: 180 (FR: 74.2)



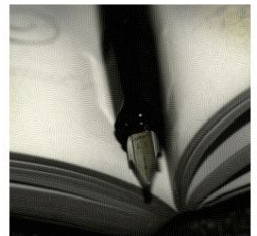
BW: 220 (FR: 70.0)



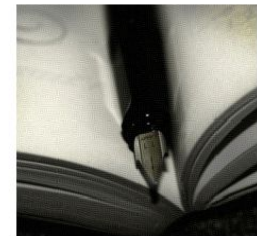
brain\_coral



brain\_coral



lampshade



brass



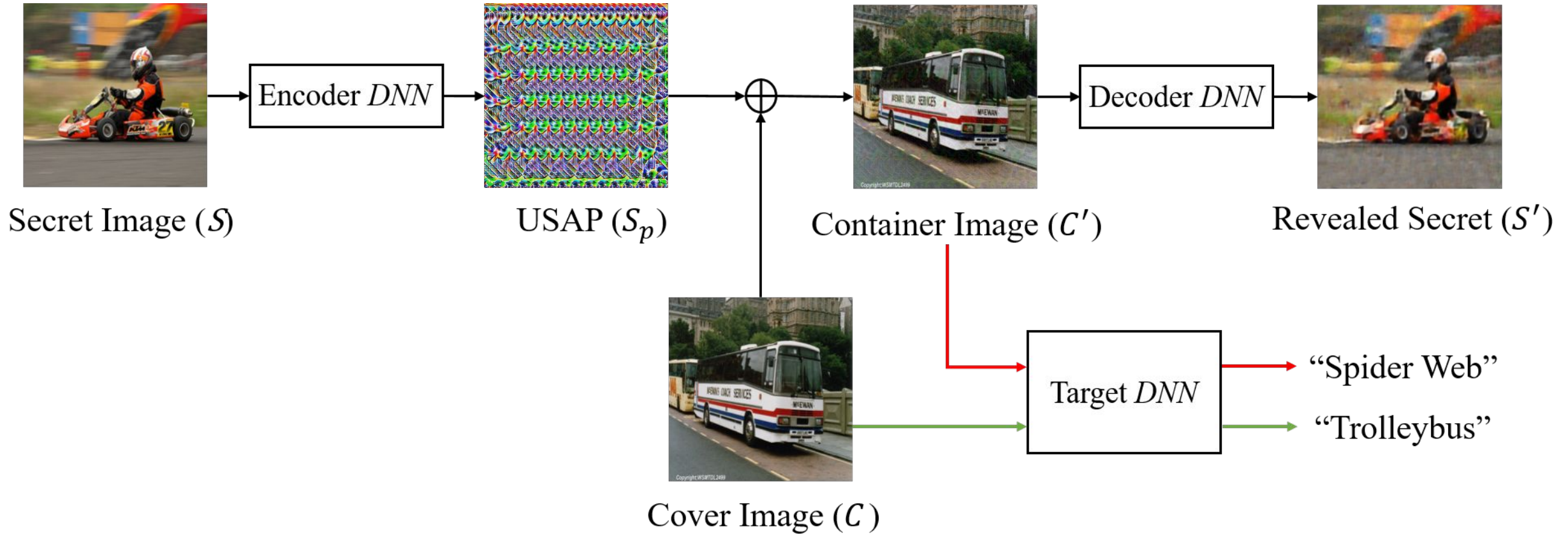
window\_screen

The UAP fails when HF content is removed.  
On the contrary, removing LF content from the UAP has trivial effect.

UAPs (of smaller magnitude) dominate images, why?

→ DNNs are sensitive to HF content & UAPs have HF properties

# Universal Secret Adversarial Perturbation (USAP)



# Takeaway

1. **What you see is not what the machine sees**
2. **Joint Investigation and Explanation: Adversarial Attack & Data Hiding**

# Thank You



Philipp Benz

<https://phibenz.github.io>  
[pbenz@kaist.ac.kr](mailto:pbenz@kaist.ac.kr)

Chaoning Zhang

<https://chaoningzhang.github.io>  
[chaoningzhang1990@gmail.com](mailto:chaoningzhang1990@gmail.com)