

Adversarial Learning

План

Adversarial Examples

Why it works

Adversarial Attacks

Adversarial Defenses

A decorative graphic on the left side of the slide consists of a grid of colored squares. The top row has one teal square. The second row has an orange square followed by a brown square. The third row has an orange square, a teal square, and a light brown square. The bottom row has a light brown square, an orange square, an orange square, and a brown square. The text "Adversarial Examples" is positioned to the right of the top two rows of squares.

Adversarial Examples

Adversaries

Adversary – противник, враг, соперник, неприятель, ...

Соперники есть во многих областях:

- Электронная почта: спаммеры
- Распознавание лиц: люди, которые хотят остаться инкогнито
- Проверка состояния машины по фото: недобросовестные водители/таксопарки
- ...

Adversaries

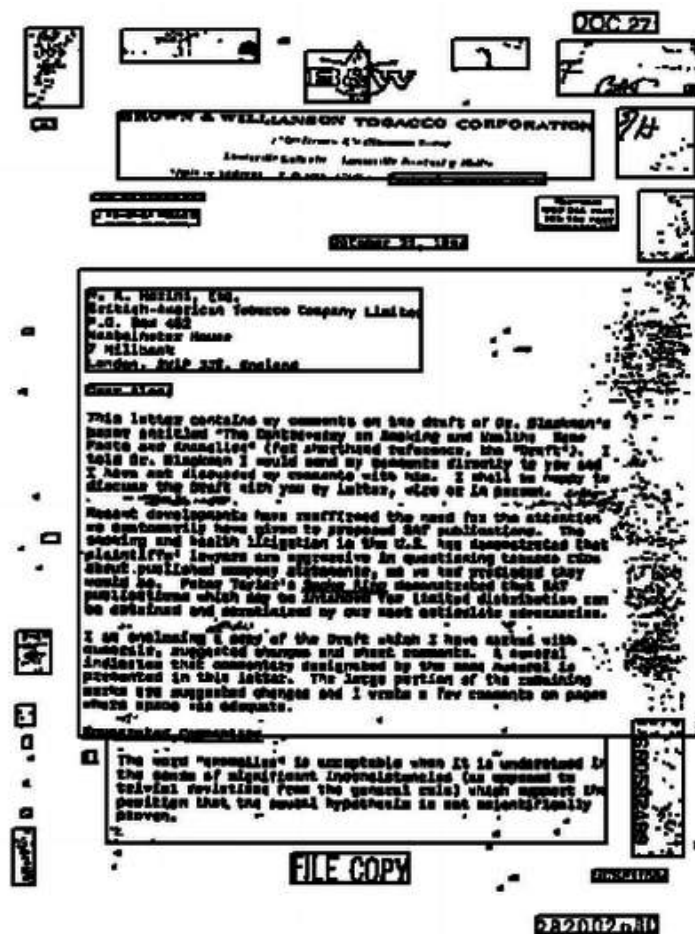
Помимо настоящих соперников есть еще и шум:

- Распознавание речи: Помехи в записи голоса
- Распознавание лиц: свет, угол обзора, ...
- Распознавание текста в распечатанных документах: артефакты печати
- Распознавание дорожных знаков: граффити, стикеры, ...
- ...

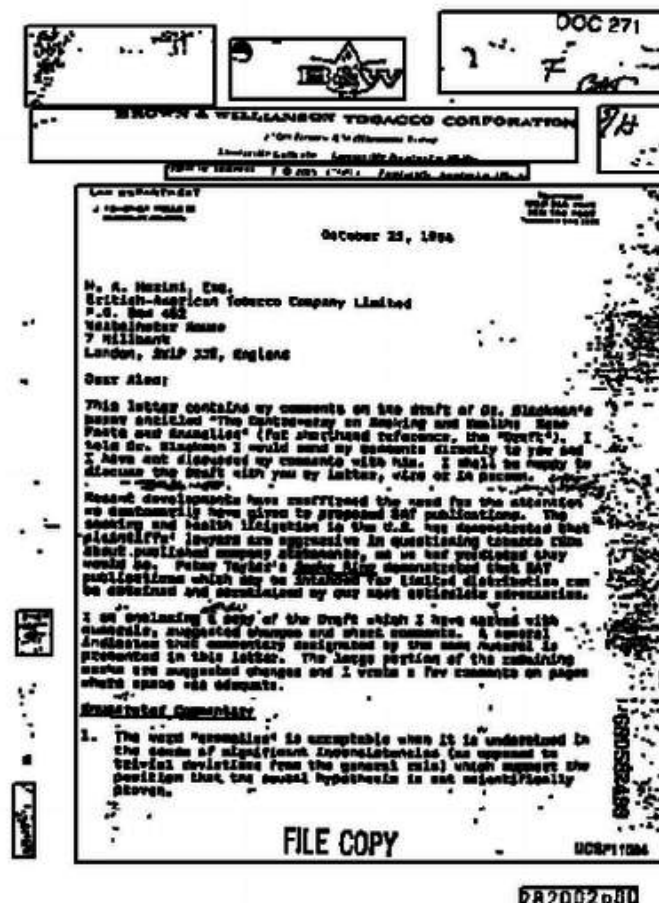
Adversarial Examples in Real Life



Adversarial Examples in Real Life

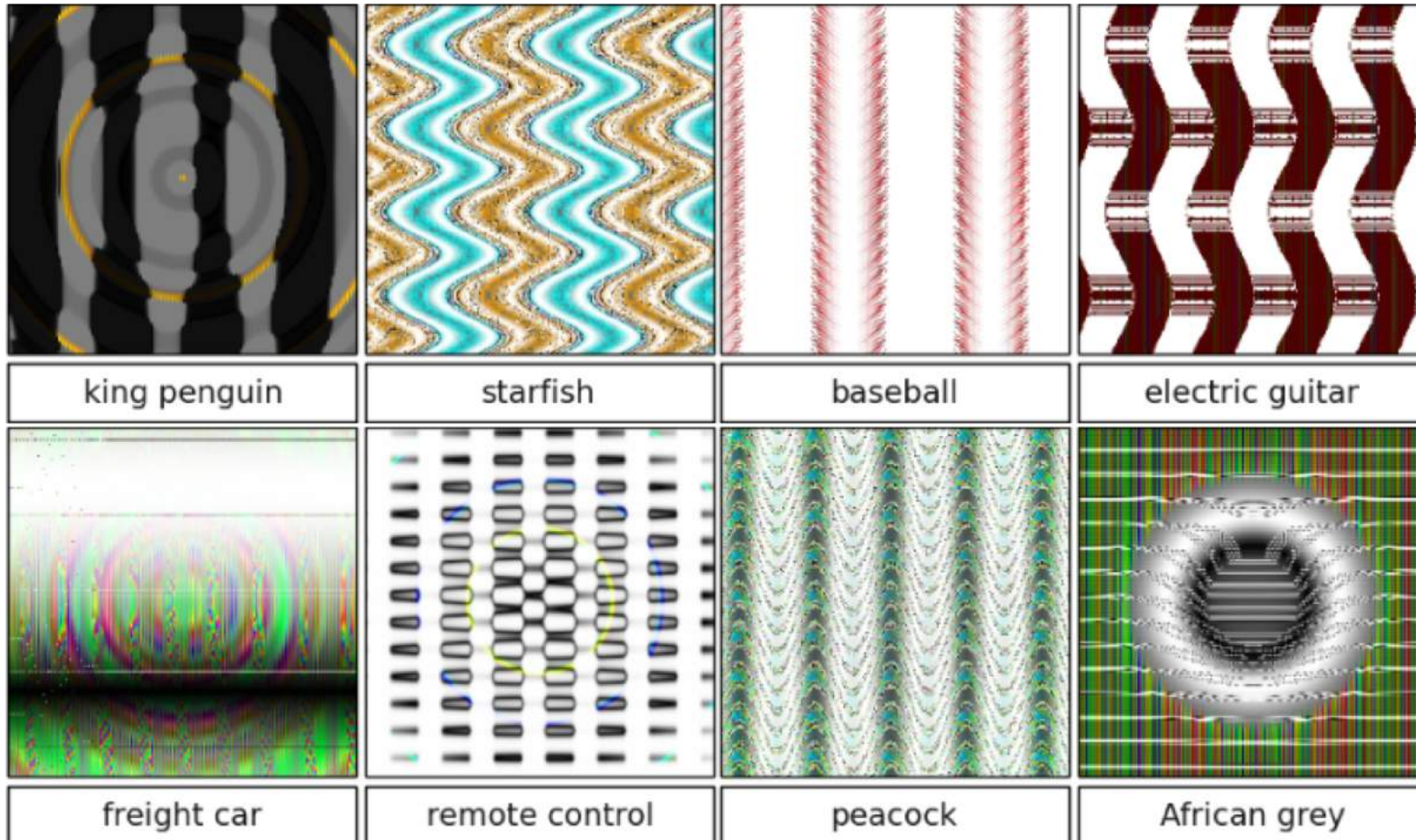


(a)



(b)

Adversarial Examples in Machine Learning

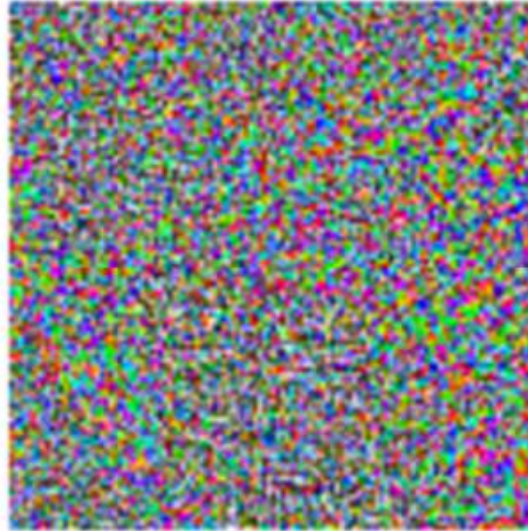


Adversarial Examples in Machine Learning



“panda”
57.7% confidence

+ .007 ×



“nematode”
8.2% confidence

=



“gibbon”
99.3 % confidence

Adversarial Examples in Machine Learning

(a) Image



(b) Prediction



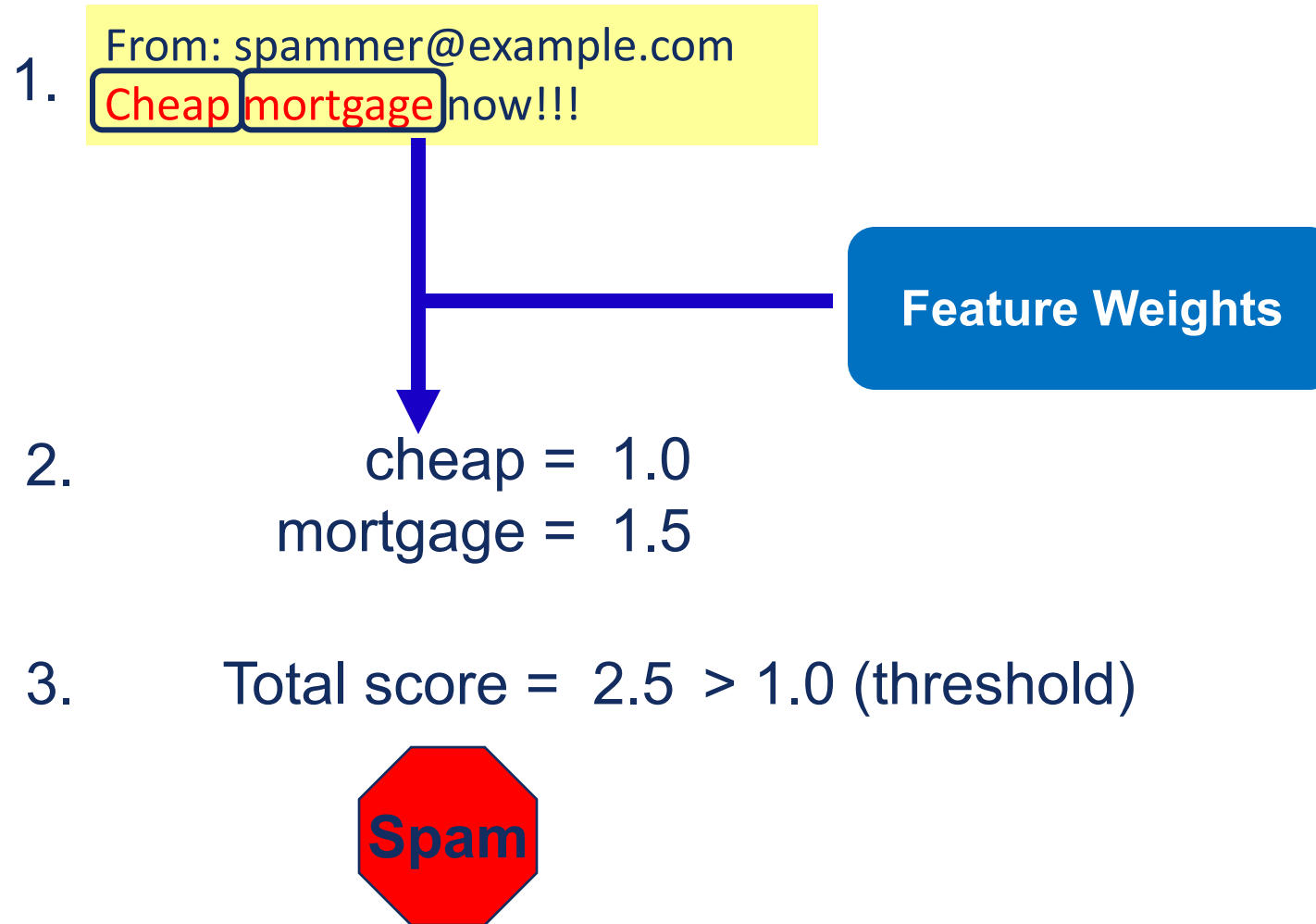
(c) Adversarial Example



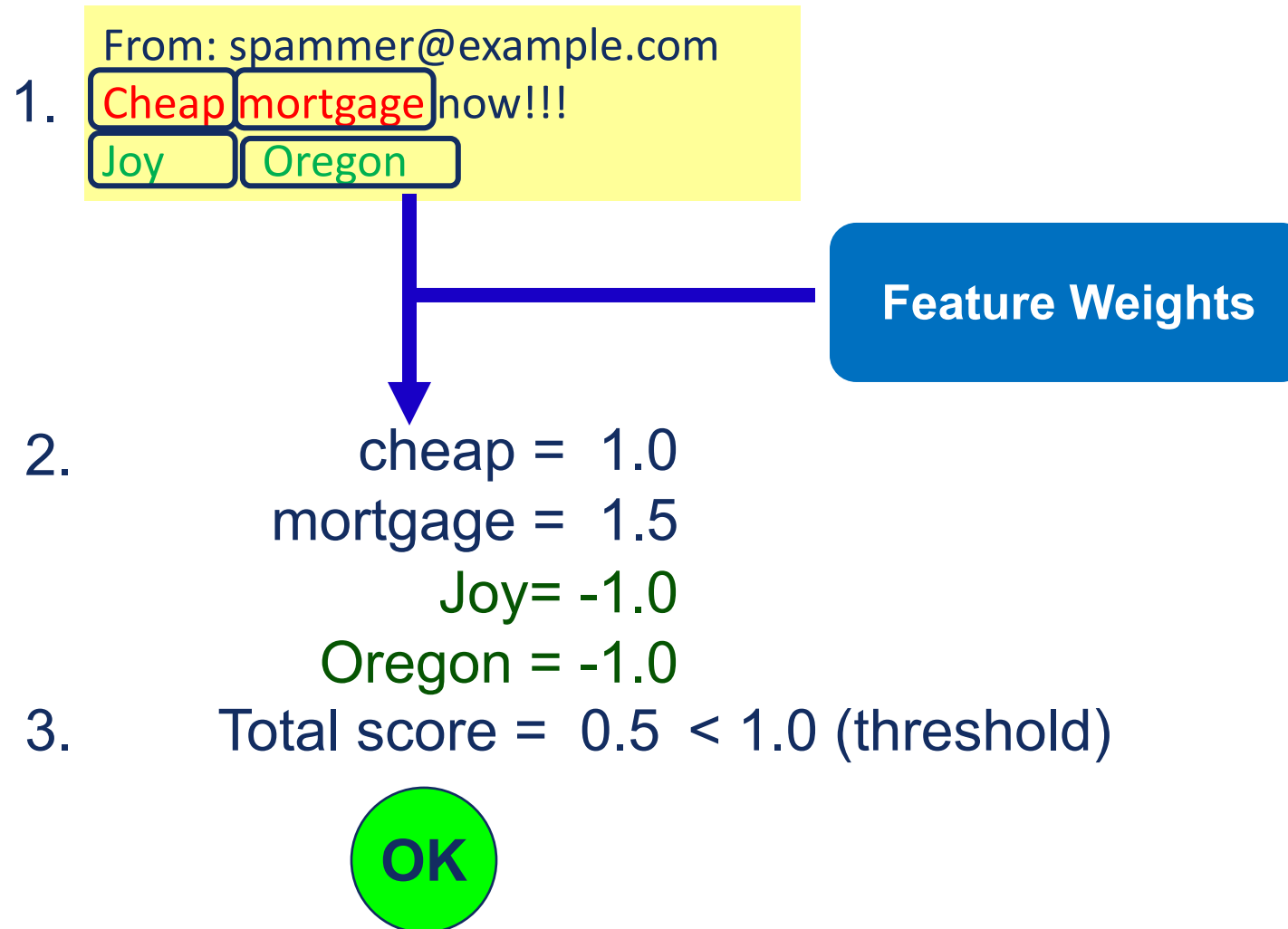
(d) Prediction



Adversarial Examples in Machine Learning



Adversarial Examples in Machine Learning



A decorative graphic on the left side of the slide consists of a grid of colored squares. The grid is 4 squares high and 4 squares wide. The colors of the squares are: Row 1: Teal, Orange, Brown, Tan. Row 2: Orange, Brown, Tan, Tan. Row 3: Orange, Teal, Tan, Tan. Row 4: Tan, Orange, Orange, Brown. The text "Why it works" is positioned to the right of this grid.

Why it works

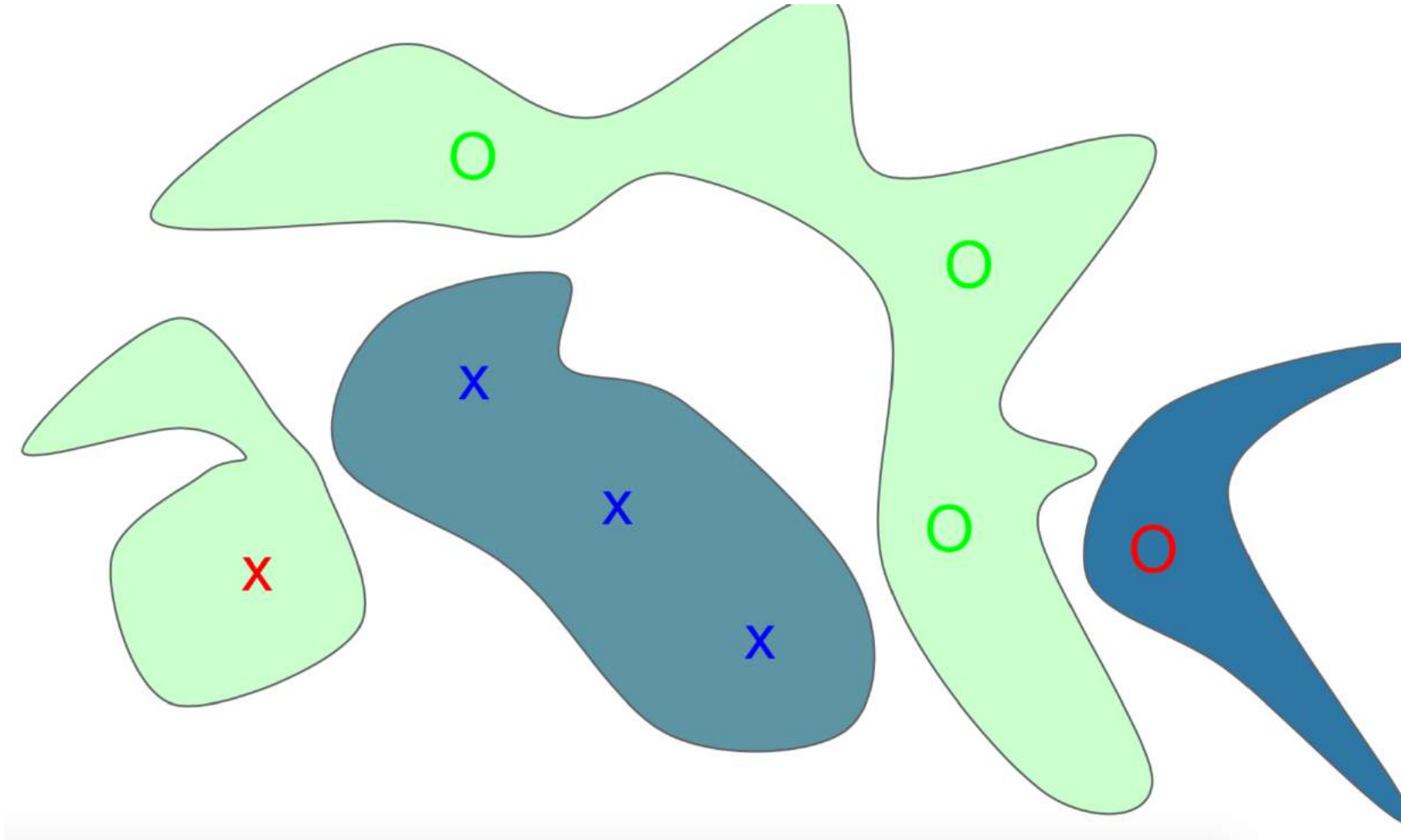
Why it works

Training Data 

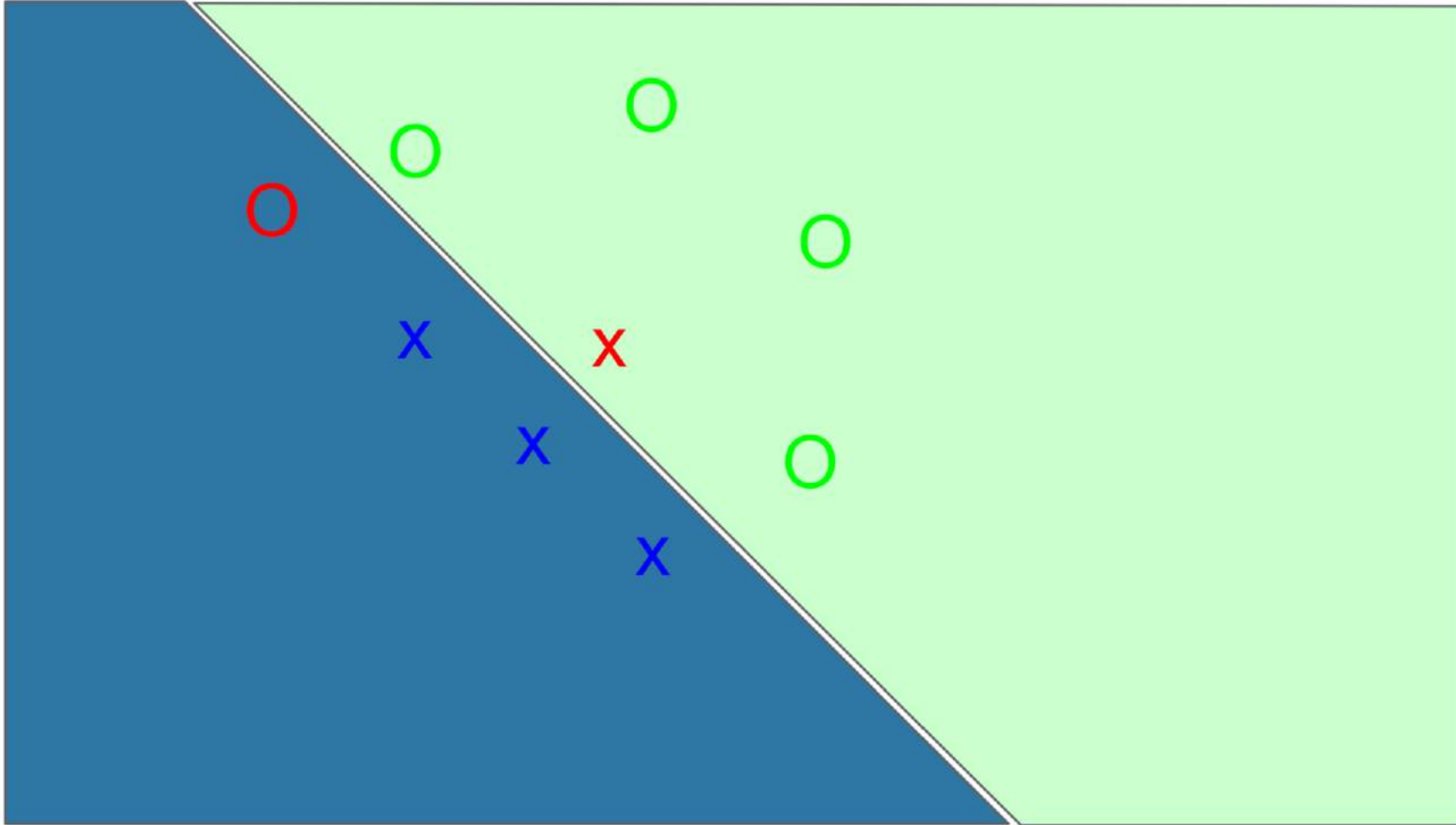
\approx

Testing Data 

Why it works: Overfitting?



Why it works: Underfitting?

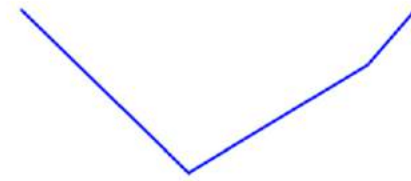


Why it works: Excessive Linearity

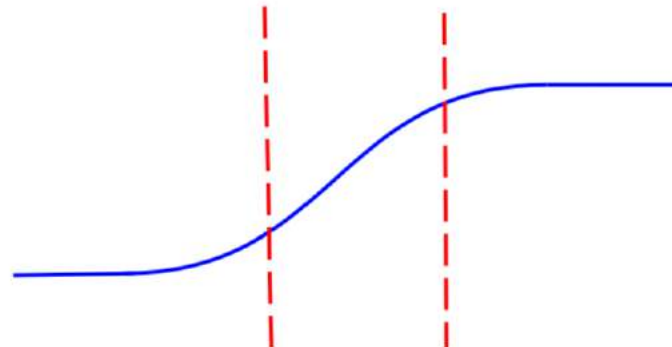
Rectified linear unit



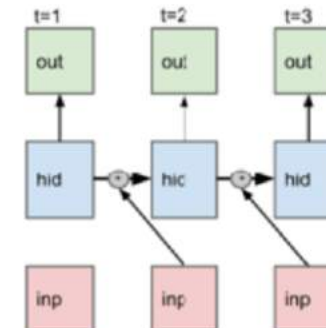
Maxout



Carefully tuned sigmoid



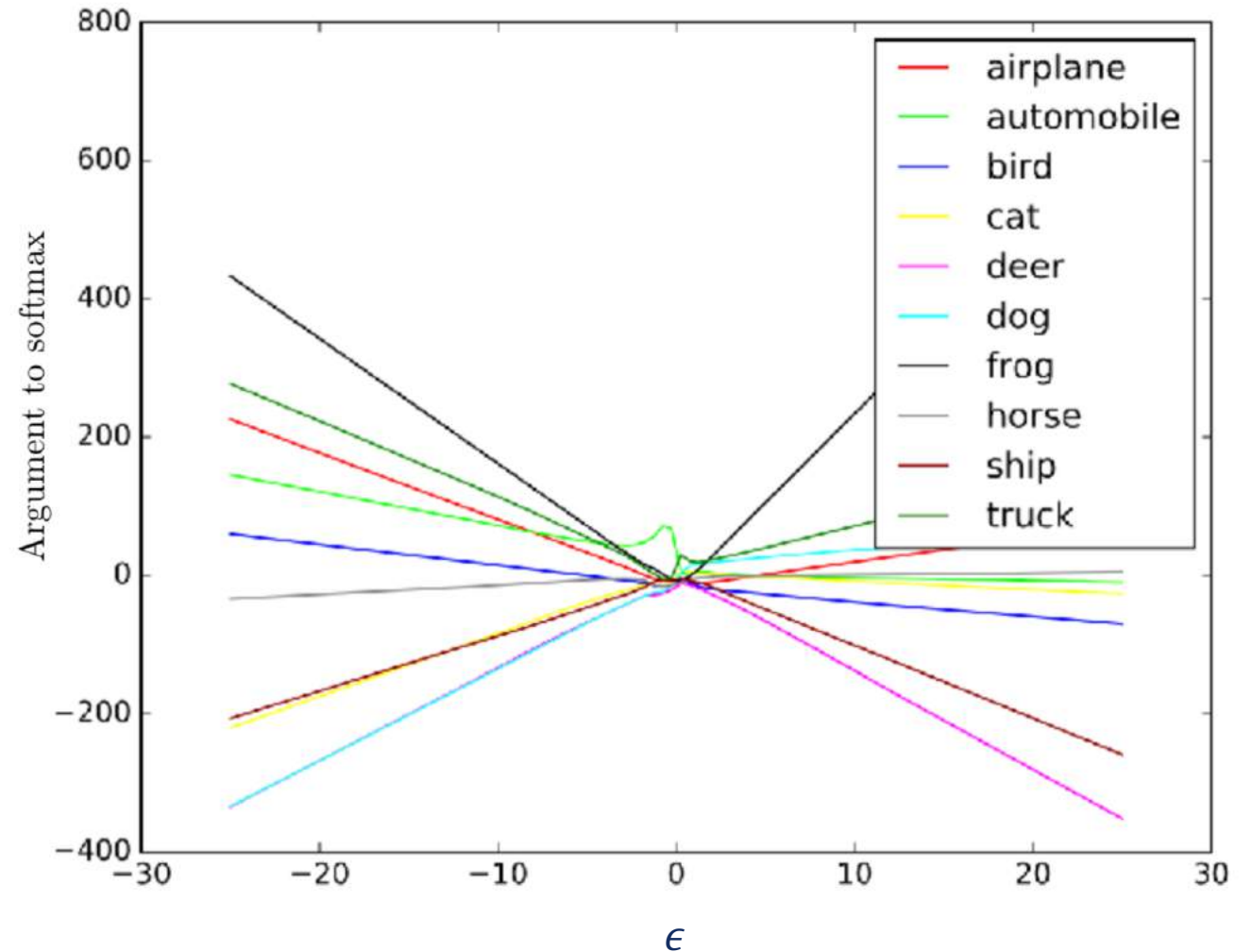
LSTM



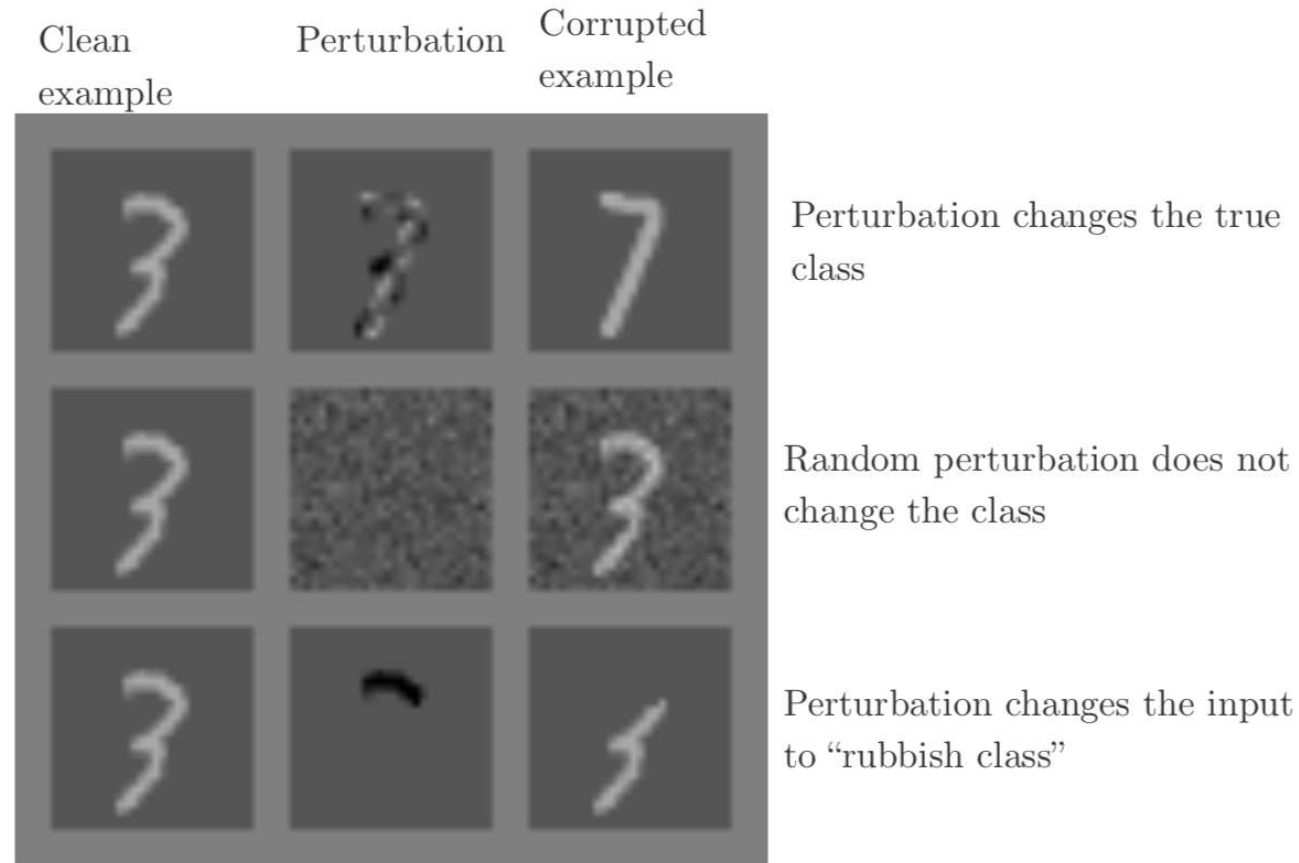
Why it works: Excessive Linearity



$\epsilon = 0$



Why it works: Excessive Linearity



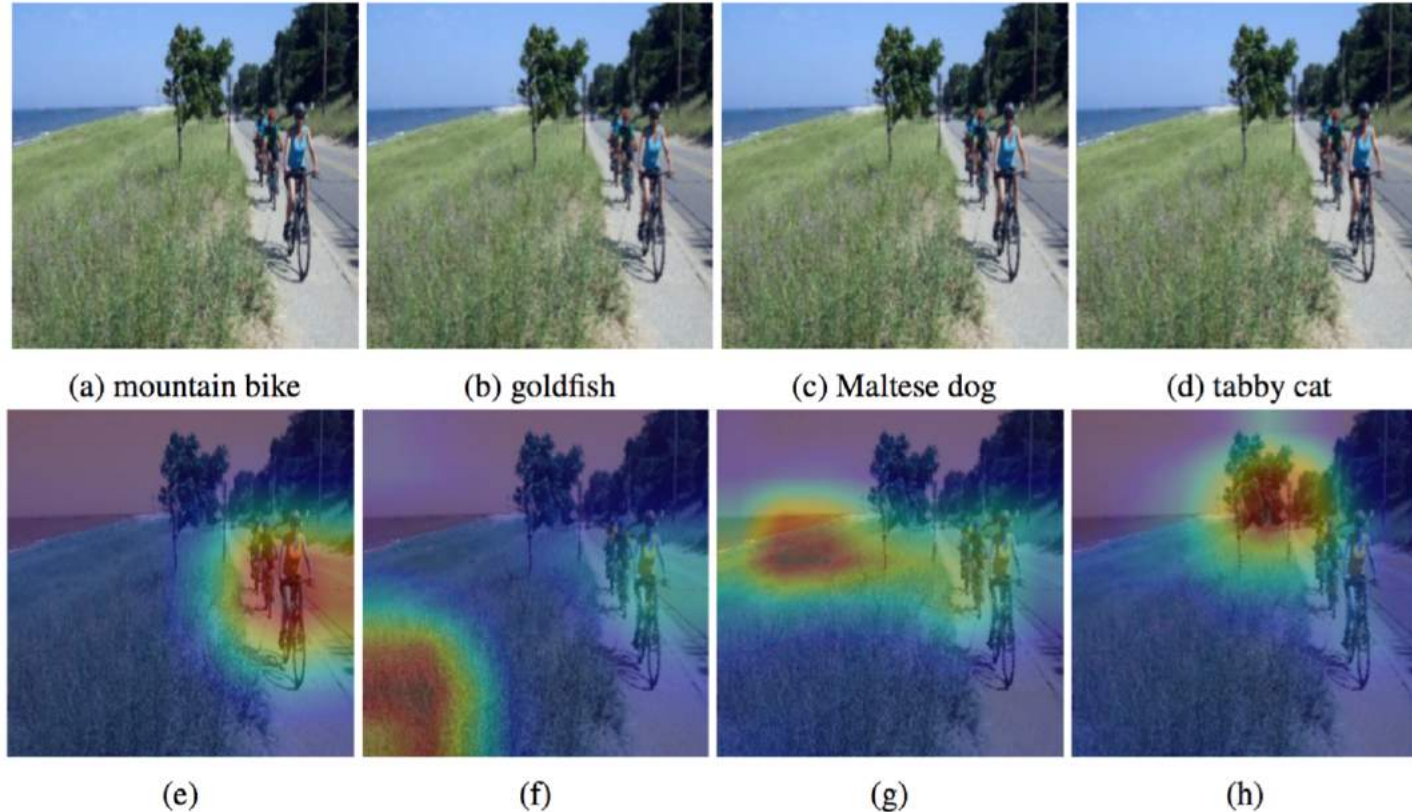
All three perturbations have L2 norm 3.96

This is actually small. We typically use 7!

Why it works: Bottom line

Если найти правильное направление движения в пространстве изображений, то можно с помощью почти невидимых для человека изменений с высокой вероятностью обмануть свёрточную сеть.

Why it works: Bottom line



CAM attention visualization for ImageNet inception_v3 model. (a) the original image and (b)-(d) are stAdv adversarial examples targeting different classes. Row 2 shows the attention visualization for the corresponding images above.

A decorative graphic in the bottom-left corner of the slide, consisting of a grid of colored squares. The squares are arranged in a pattern that tapers to the right. The colors used are teal, orange, brown, and light beige. The squares are separated by thin white lines.

Adversarial Attacks

Adversarial Attacks: Classification

- Evasion attacks
 - Targeted/non-targeted
 - White-box, black-box
- Data Poisoning

White-box Evasion Attacks: Non-targeted FGSM

$$x^{adv} = x + \varepsilon \cdot \text{sign}(\nabla_x J(x, y_{true})),$$

where

x is the input (clean) image,

x^{adv} is the perturbed adversarial image,

J is the classification loss function,

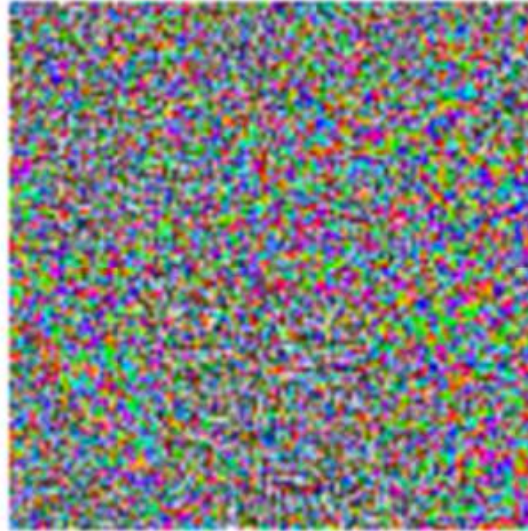
y_{true} is true label for the input x .

White-box Evasion Attacks: Non-targeted FGSM



“panda”
57.7% confidence

+ .007 ×



“nematode”
8.2% confidence

=



“gibbon”
99.3 % confidence

White-box Evasion Attacks: Non-targeted I-FGSM

$$x_0^{adv} = x, \quad x_{t+1}^{adv} = x_t^{adv} + \alpha \cdot \text{sign}(\nabla_x J(x_t^{adv}, y)).$$

White-box Evasion Attacks: Targeted FGSM

$$x^{adv} = x - \varepsilon \cdot \text{sign}(\nabla_x J(x, y_{target})),$$

where

y_{target} is the target label for the adversarial attack.

White-box Optimization Attacks in General

$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + J(f_{\theta}(x + \delta), y^*)$$


Perturbation/Noise Matrix \rightarrow δ \rightarrow $\lambda \|\delta\|_p$ \rightarrow $J(f_{\theta}(x + \delta), y^*)$ \rightarrow Adversarial Target Label

$\lambda \|\delta\|_p$ Lp norm (L-0, L-1, L-2, ...)

$J(f_{\theta}(x + \delta), y^*)$ Loss Function

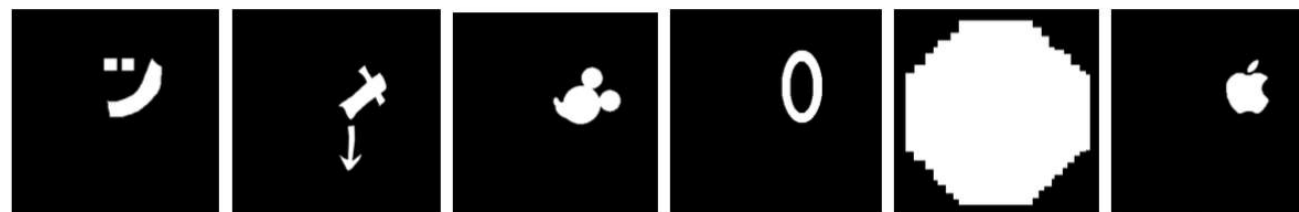
$$\operatorname{argmin}_{\delta} \lambda \|\delta\|_p + \frac{1}{k} \sum_{i=1}^k J(f_{\theta}(x_i + \delta), y^*)$$

\downarrow

{  }

White-box Optimization Attacks in General

$$\operatorname{argmin}_{\delta} \lambda \| \textcircled{M_x} \cdot \delta \|_p + \frac{1}{k} \sum_{i=1}^k J(f_{\theta}(x_i + \textcircled{M_x} \cdot \delta), y^*)$$



Subtle Poster
Camouflage Sticker

Mimic vandalism

“Hide in the human
psyche”

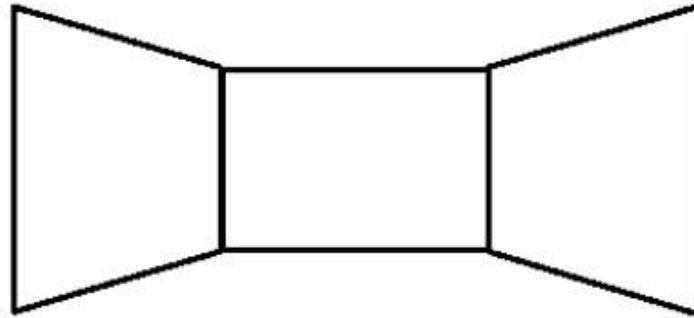


Adversarial Transformation Networks



“panda”

f_W

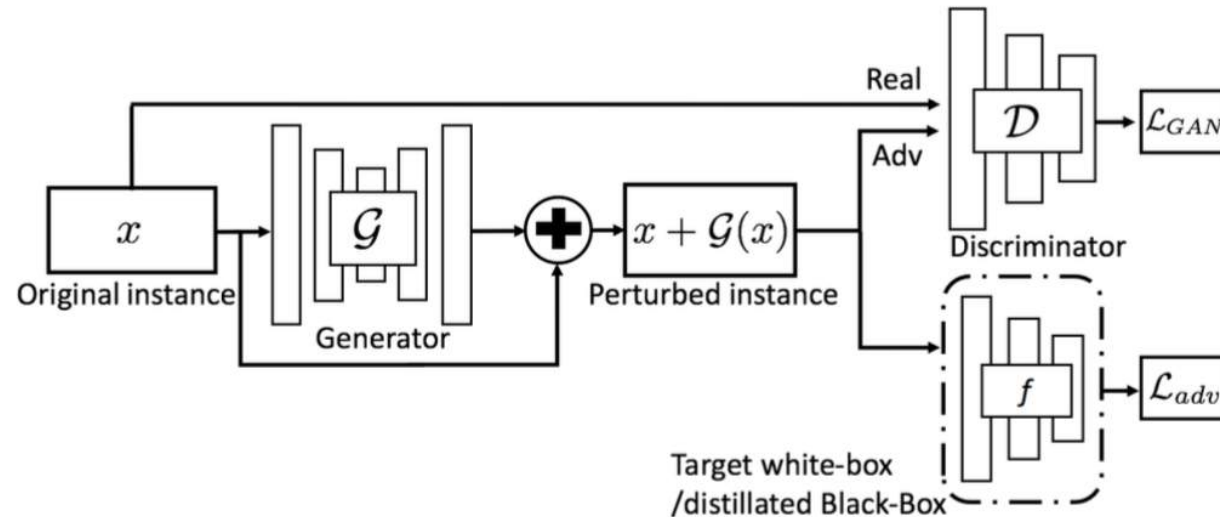


Feedforward Net



“gibbon”

Generating Adversarial Examples with GANs



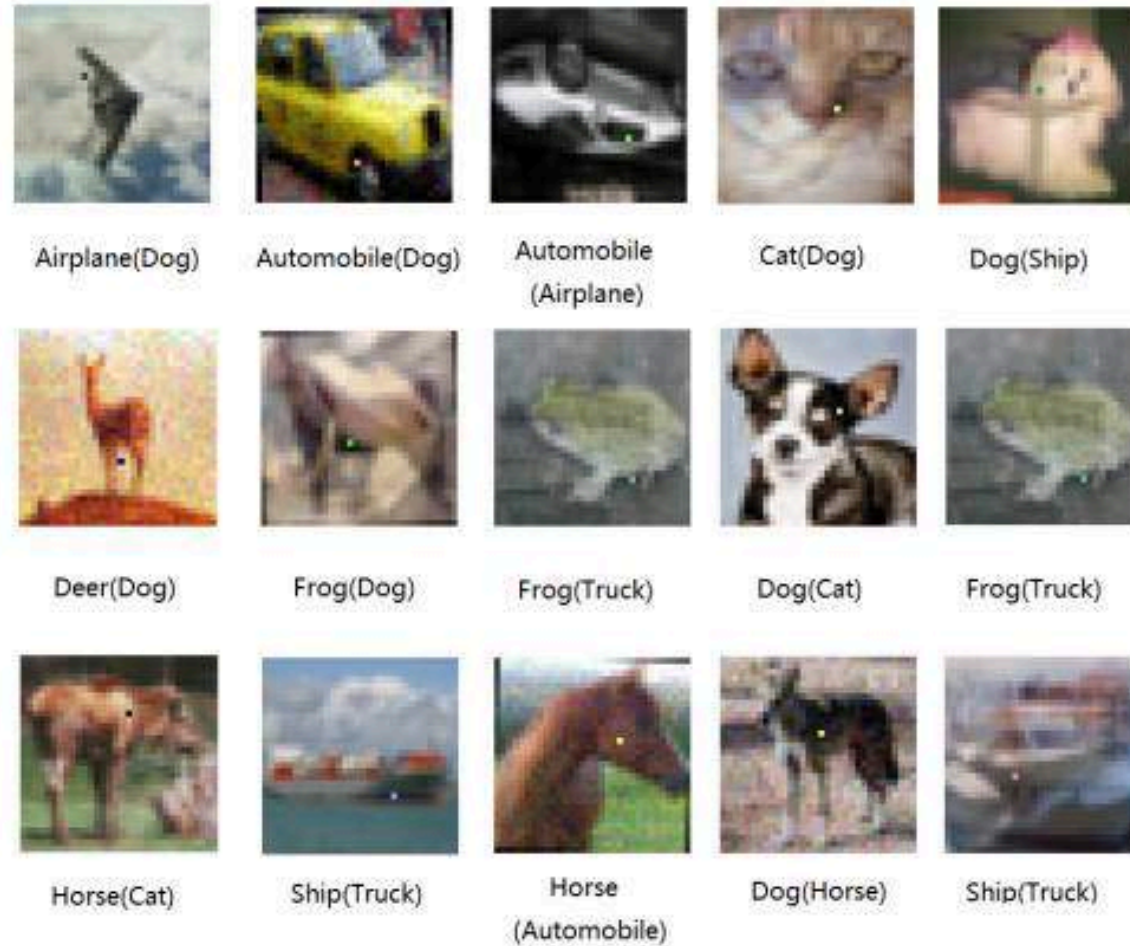
Black-box can be performed here via distillation

$$\mathcal{L}_{GAN} = \mathbb{E}_{x \sim \mathcal{P}_{\text{data}}(x)} \log \mathcal{D}(x) + \mathbb{E}_{x \sim \mathcal{P}_{\text{data}}(x)} \log(1 - \mathcal{D}(x + \mathcal{G}(x)))$$

$$\mathcal{L} = \mathcal{L}_{adv}^f + \alpha \mathcal{L}_{GAN} + \beta \mathcal{L}_{hinge}$$

The GAN loss here tries to ensure the diversity of adversarial examples

One-pixel Adversarial Attack

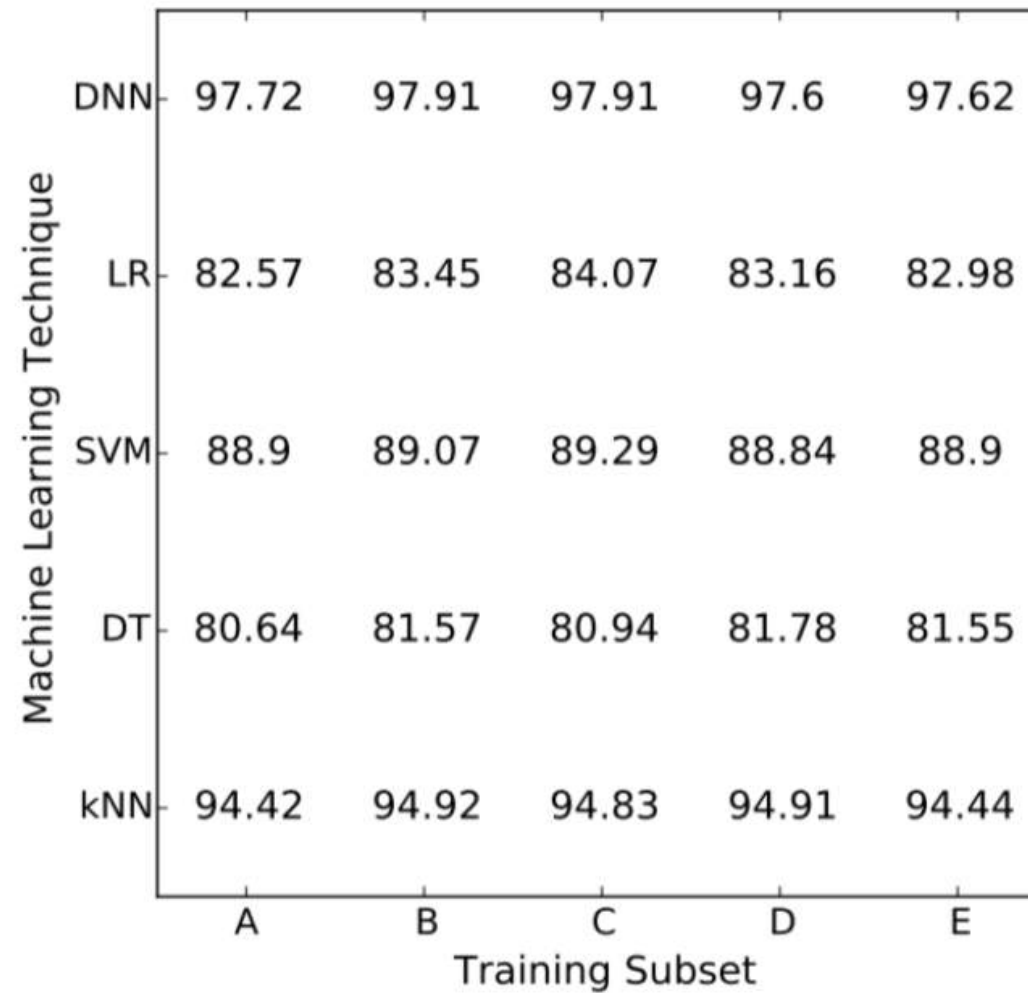


Source: Su et al. One pixel attack for fooling deep neural networks

Model2Model Transfer

Source Machine Learning Technique	DNN	LR	SVM	DT	kNN	Ens.
DNN	38.27	23.02	64.32	79.31	8.36	20.72
LR	6.31	91.64	91.43	87.42	11.29	44.14
SVM	2.51	36.56	100.0	80.03	5.19	15.67
DT	0.82	12.22	8.85	89.29	3.31	5.11
kNN	11.75	42.89	82.16	82.95	41.65	31.92
Target Machine Learning Technique	DNN	LR	SVM	DT	kNN	Ens.

Data2Data Transfer



Machine Learning Technique	DNN	97.72	97.91	97.91	97.6	97.62
	LR	82.57	83.45	84.07	83.16	82.98
	SVM	88.9	89.07	89.29	88.84	88.9
	DT	80.64	81.57	80.94	81.78	81.55
	kNN	94.42	94.92	94.83	94.91	94.44
		A	B	C	D	E
		Training Subset				

Black-box Evasion Attacks: FD Method

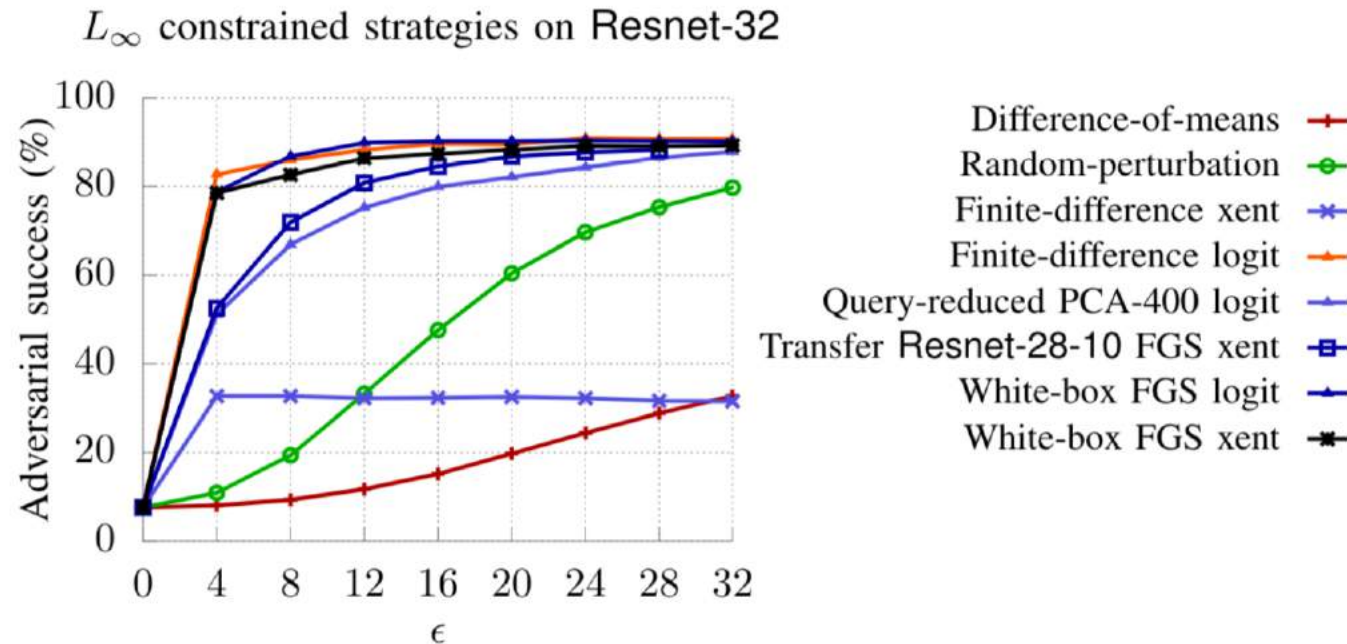
Given d -dimensional vector \mathbf{x} , we can make $2d$ queries to estimate the gradient as below

$$\text{FD}_{\mathbf{x}}(g(\mathbf{x}), \delta) = \begin{bmatrix} \frac{g(\mathbf{x} + \delta \mathbf{e}_1) - g(\mathbf{x} - \delta \mathbf{e}_1)}{2\delta} \\ \vdots \\ \frac{g(\mathbf{x} + \delta \mathbf{e}_d) - g(\mathbf{x} - \delta \mathbf{e}_d)}{2\delta} \end{bmatrix}$$

An example of approximate FGS with finite difference

$$\mathbf{x}_{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\text{FD}_{\mathbf{x}}(\ell_f(\mathbf{x}, y), \delta))$$

Black-box Evasion Attacks Results



Effectiveness of various single step black-box attacks on CIFAR-10. The y-axis represents the variation in adversarial success as ϵ increases.

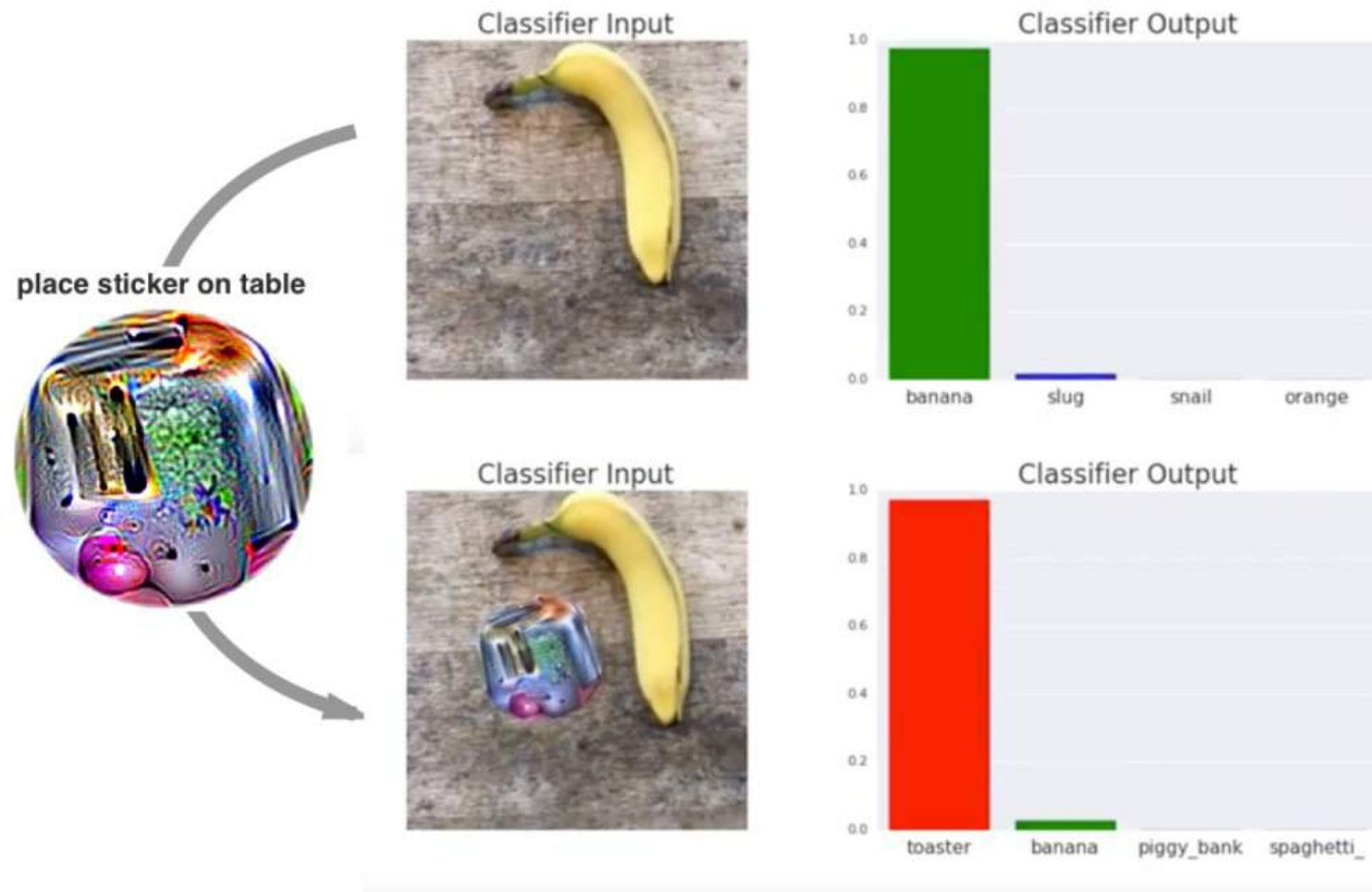
Finite Differences method outperform other black-box attacks and achieves similar attack success rate with the white-box attack

Black-box Evasion Attack: Model Stealing

Если получить ответ black-box классификатора можно почти бесплатно, то можно обучить свой классификатор на ответах black-box классификатора.

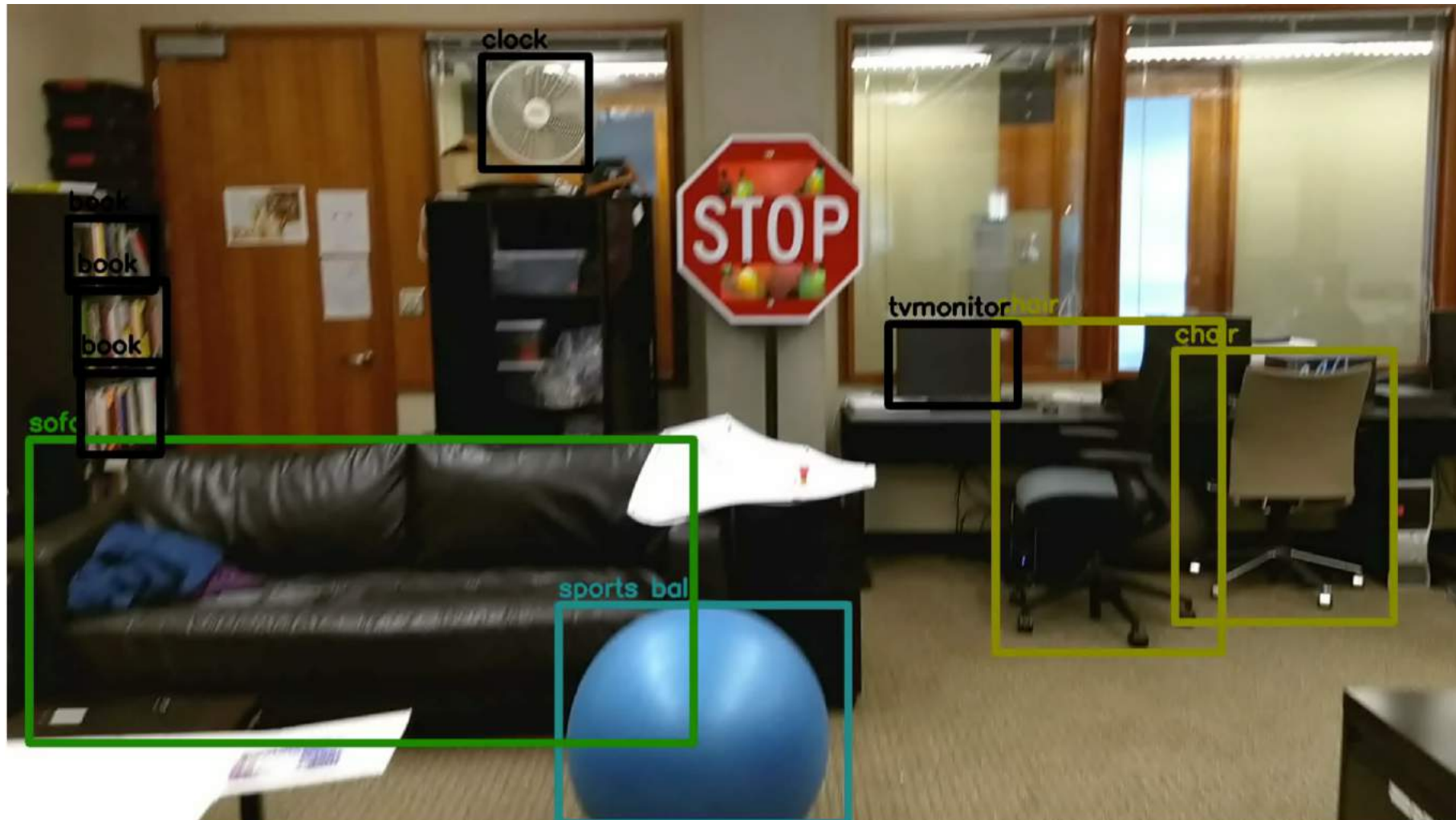
Затем можно научиться обманывать полученный классификатор, и с высокой вероятностью обман будет успешен и для black-box классификатора.

Real World Evasion Attacks: Adversarial Patch



Source: Brown et al. Adversarial Patch

Real World Evasion Attacks: Adversarial Stickers



Source: Adversarial Machine Learning Tutorial (<https://aaai18adversarial.github.io>)

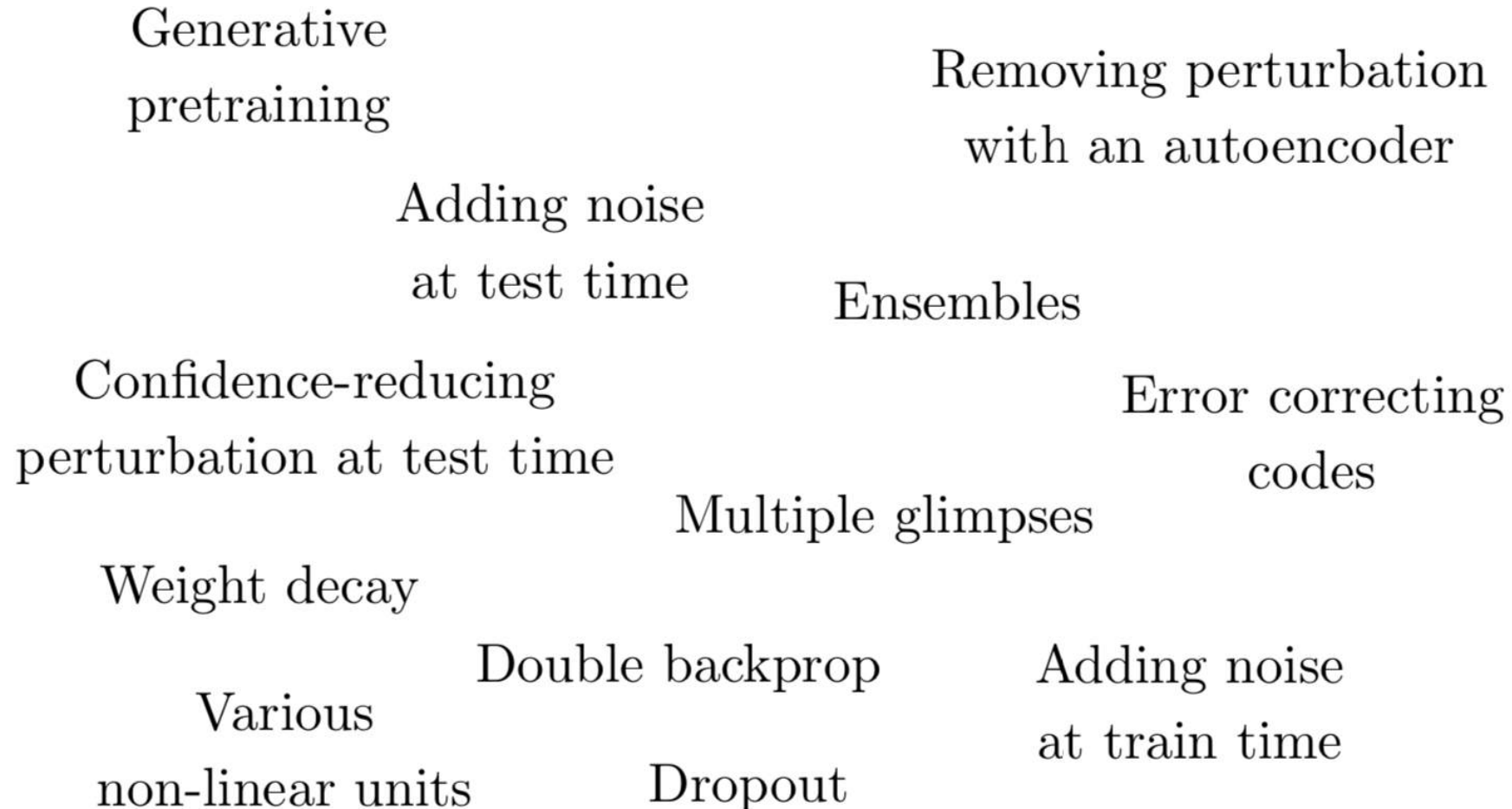
Data Poisoning

- Direct poisoning of labels
- Deleting/injecting examples
- Injecting noise/perturbations in data

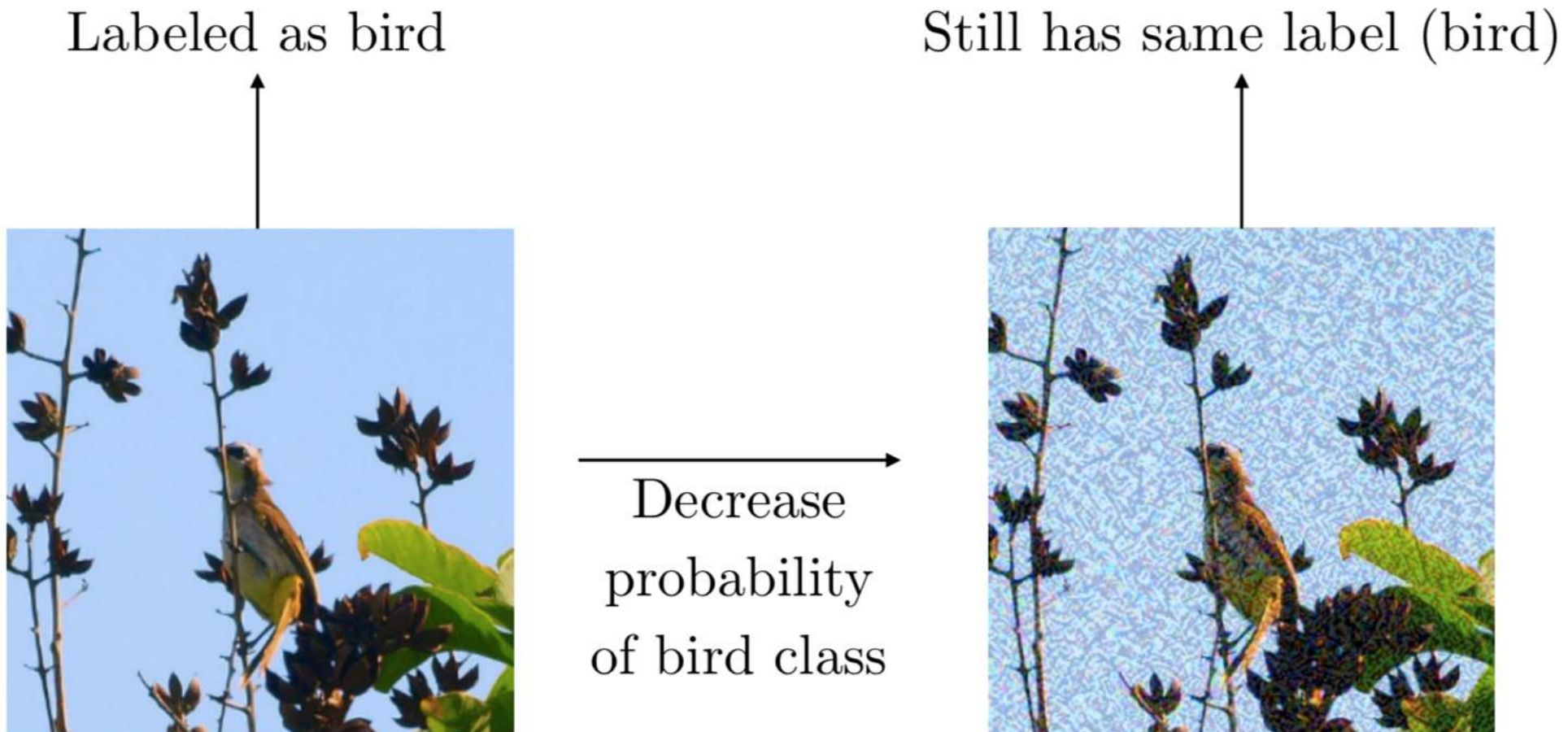
A decorative graphic on the left side of the slide consists of a grid of colored squares. The top row has one teal square. The second row has one orange square and one brown square. The third row has one orange square, one teal square, and one light brown square. The bottom row has one light brown square, one orange square, one orange square, and one brown square.

Adversarial Defenses

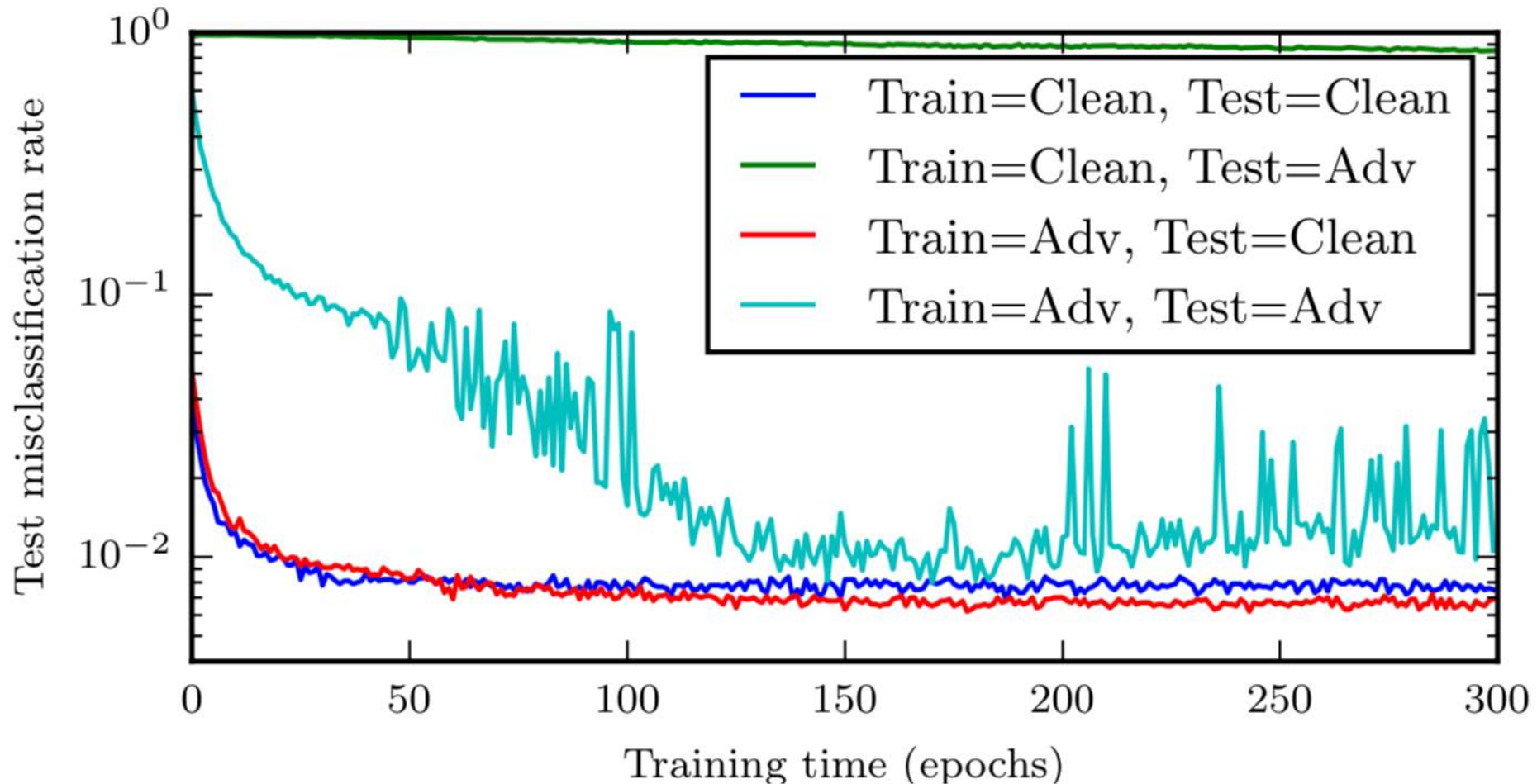
Failed Defenses



Adversarial Training

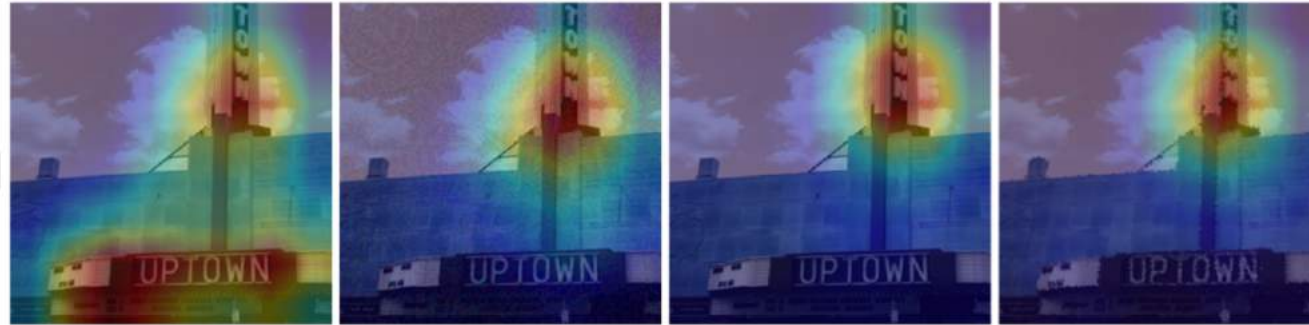


Adversarial Training



Adversarial Training

inception_v3 model



(a) Benign

(b) FGSM

(c) C&W

(d) StAdv

Adversarial trained
inception_v3 model



(e) Benign

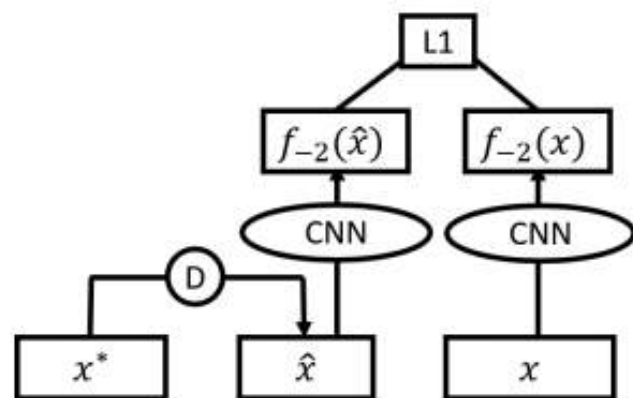
(f) FGSM

(g) C&W

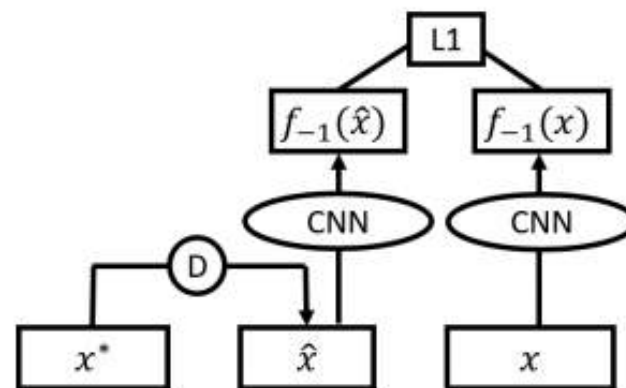
(h) StAdv

CAM attention visualization for ImageNet inception_v3 model. Column 1 shows the CAM map corresponding to the original image. Column 2-4 show the adversarial examples generated by different methods. (a) and (e)-(g) are labeled as the ground truth “cinema”, while (b)-(d) and (h) are labeled as the adversarial target “missile.”

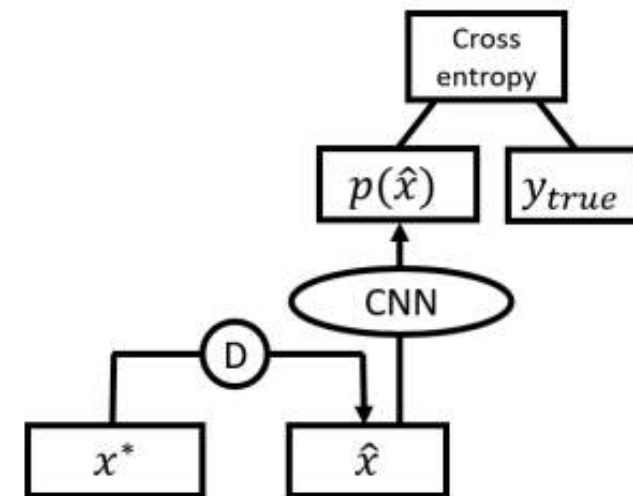
High Level Representation Guided Denoiser



(a) FGD



(b) LGD



(c) CGD

Papers

Defense against Adversarial Attacks Using High-Level Representation Guided Denoiser: <https://arxiv.org/abs/1712.02976>

Generating Adversarial Examples with Adversarial Networks: <https://arxiv.org/abs/1801.02610>

Explaining and Harnessing Adversarial Examples: <https://arxiv.org/abs/1412.6572>

Robust Physical-World Attacks on Deep Learning Models: <https://arxiv.org/abs/1707.08945>

Adversarial Patch: <https://arxiv.org/abs/1712.09665>

One pixel attack for fooling deep neural networks: <https://arxiv.org/abs/1710.08864>

Materials

Adversarial Machine Learning Tutorial (tutorial presentations): <https://aaai18adversarial.github.io>

cs231n Lecture 16| Adversarial Examples and Adversarial Training:
https://www.youtube.com/watch?v=ClfsB_EYsVI&list=PL3FW7Lu3i5JvHM8ljYj-zLfQRF3EO8sYv

CleverHans Library: <https://github.com/tensorflow/cleverhans>