



Robustness of Deep Neural Networks (part 02)



Mohammad Khalooei

PhD Candidate of Amirkabir University of Technology

Research Assistant of Institute for Research in Fundamental Sciences(IPM)

مرکز تحقیقات
هوش مصنوعی پارس
هوشمندسازی فرایندهای زندگی



کالج تخصصی
هوش مصنوعی پارس

Outline

...

■ Day 01 (28 Aug 2021)

- Motivations of progress in DNN



- Security of AI

- Attack



- Defense (in a brief intro)



■ Day 02 (29 Aug 2021)

- Adversarial attack (cont.)



- Adversarial Defense



- Security of AI notes for industrial models

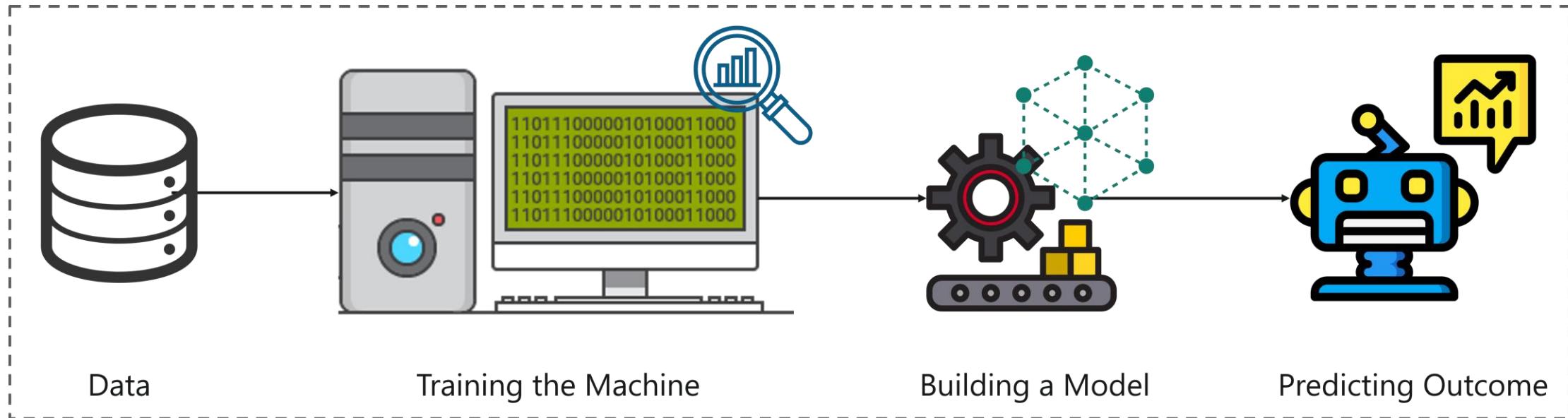


Outline

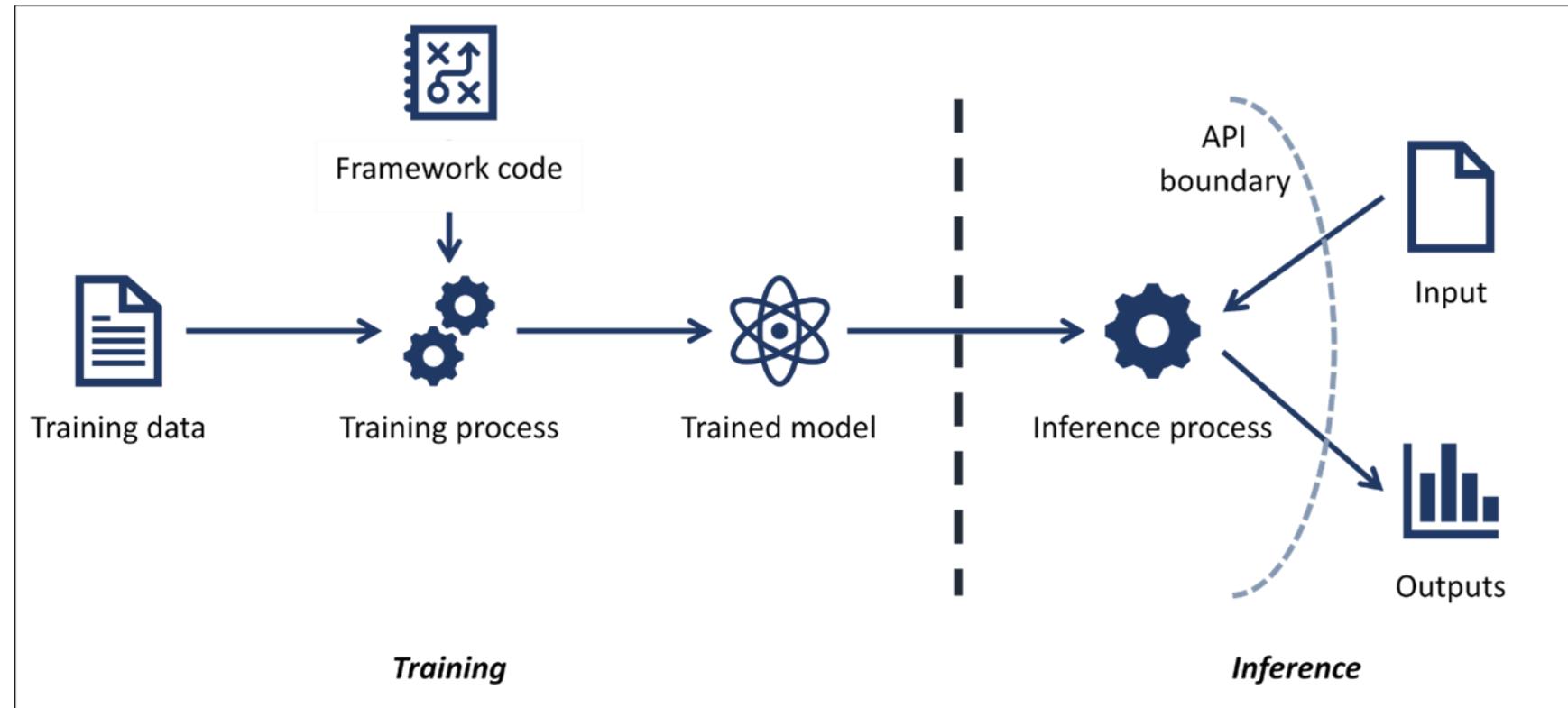
Day 02

- Day 01 (28 Aug 2021)
 - Motivations of progress in DNN 
 - Security of AI
 - Attack  
 - Defense (in a brief intro) 
- Day 02 (29 Aug 2021)
 - Adversarial attack (cont.)  
 - Adversarial Defense  
 - Security of AI notes for industrial models 

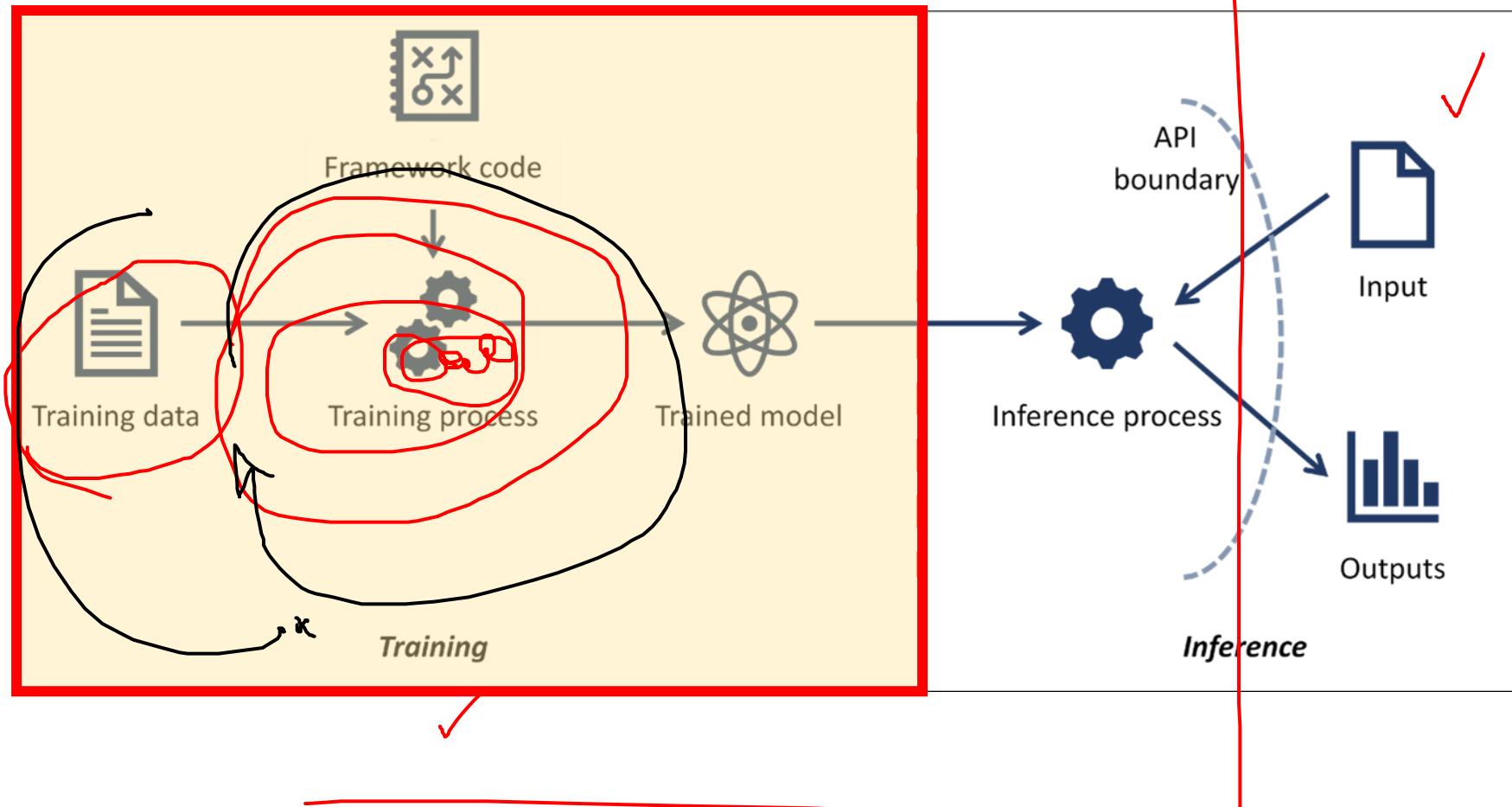
Review



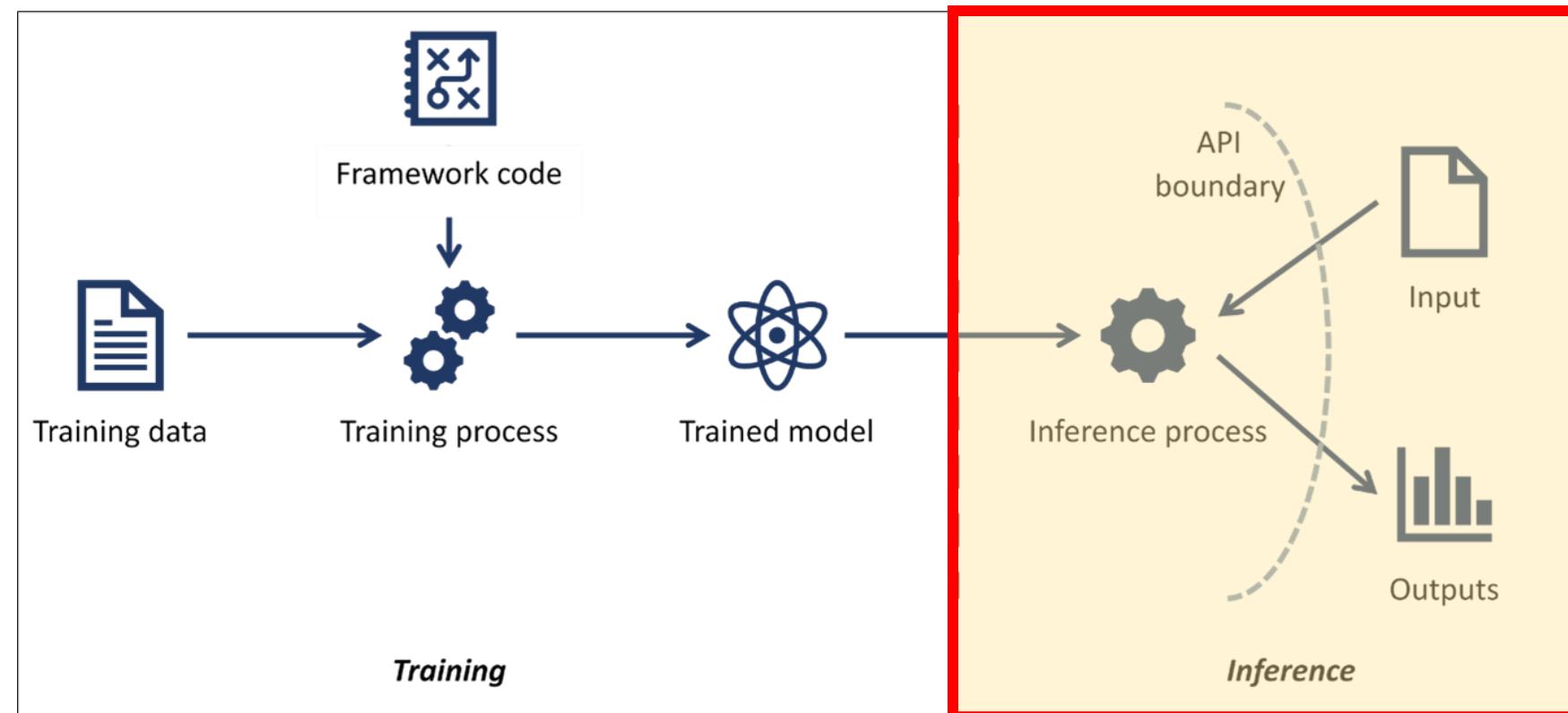
ML Model API



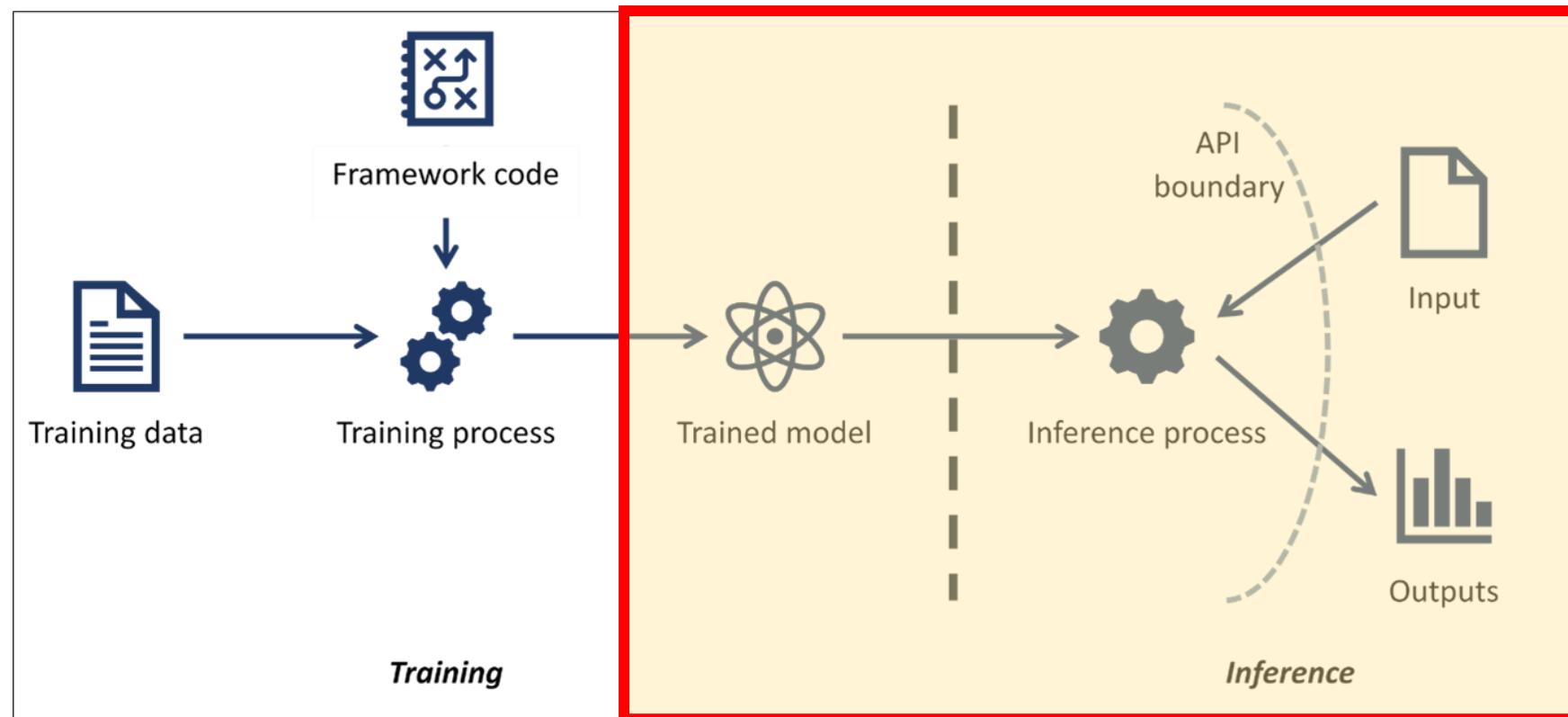
API based AI Model

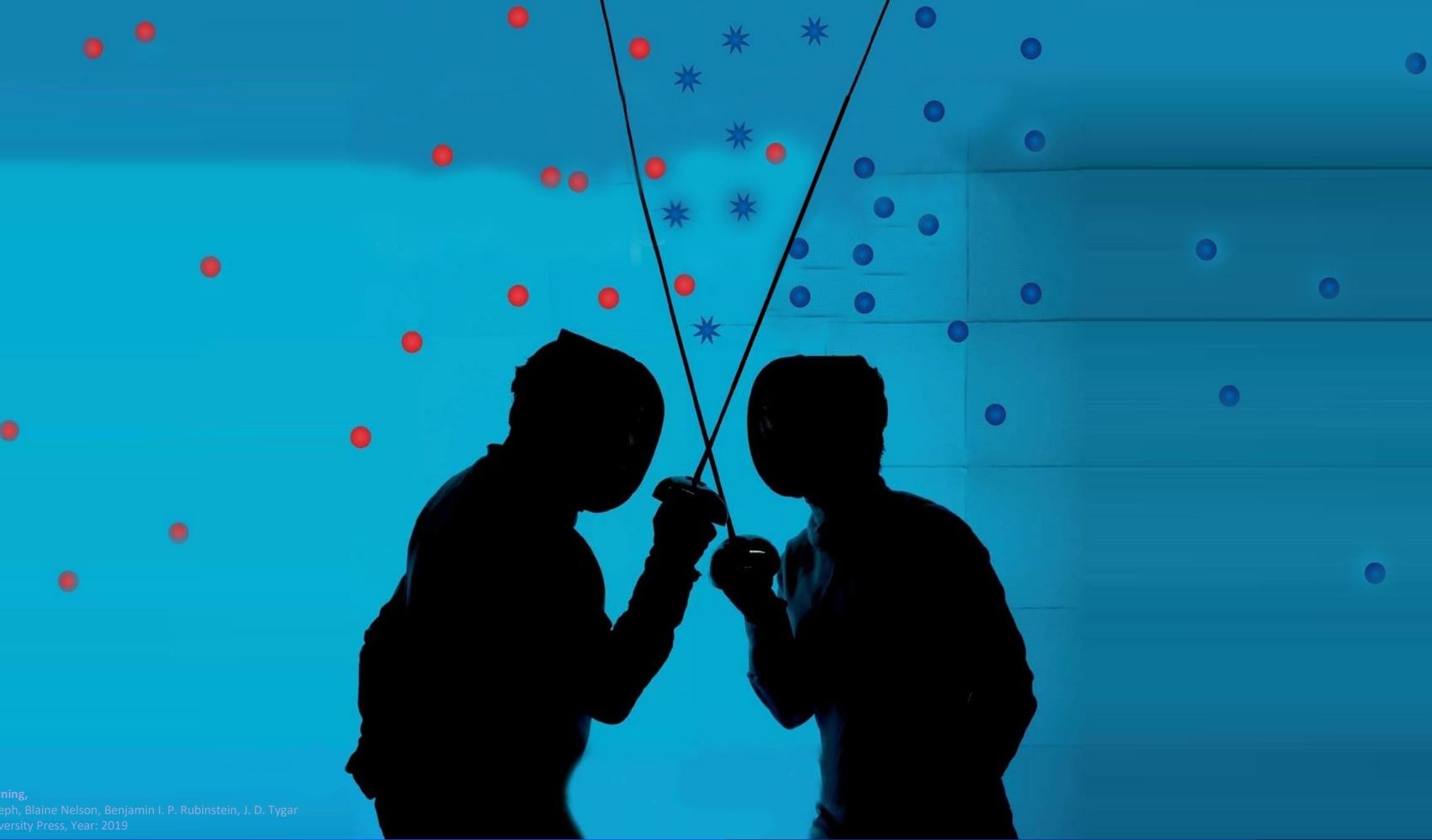


API based AI Model



API based AI Model



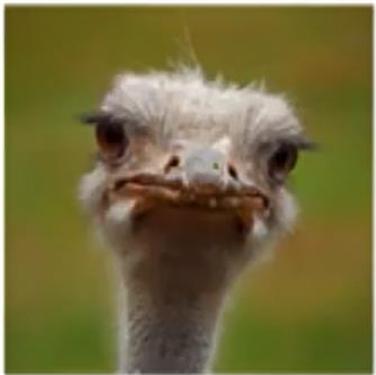


Adversarial Machine Learning,
Author(s): Anthony D. Joseph, Blaine Nelson, Benjamin I. P. Rubinstein, J. D. Tygar
Publisher: Cambridge University Press, Year: 2019

Motivations ...

Suppose that we have

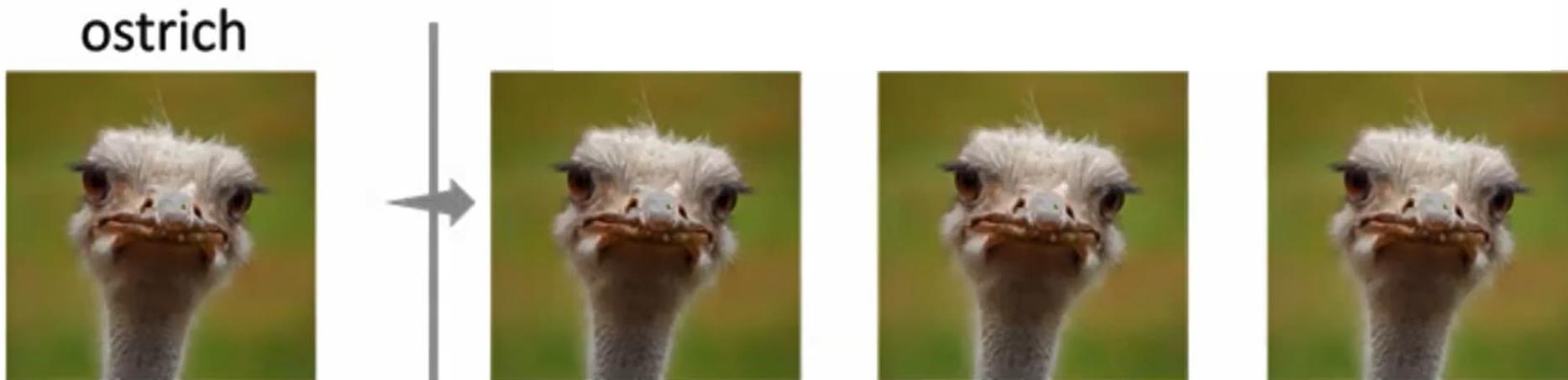
- Best trained **image classifier model** (using neural networks)



Motivations ...

Suppose that we have

- Best trained **image classifier model** (using neural networks)



Motivations ...

Suppose that we have

- Best trained image classifier model (using neural networks)



What is wrong with this AI model?

Review::



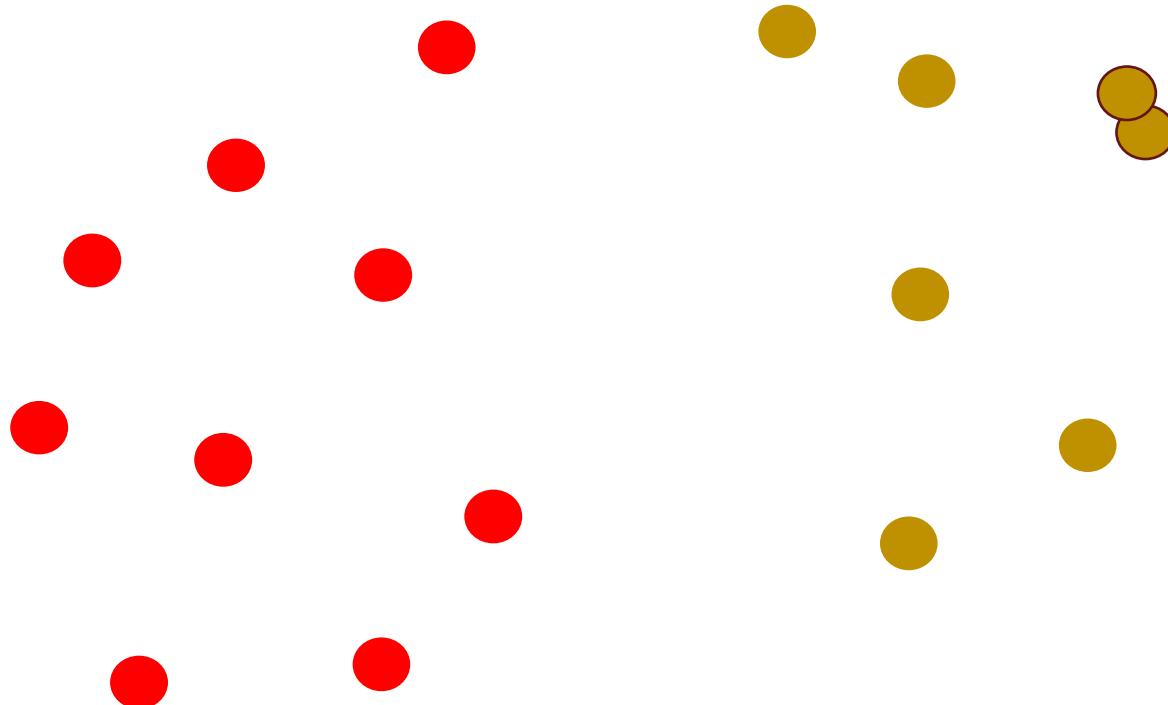
Review::



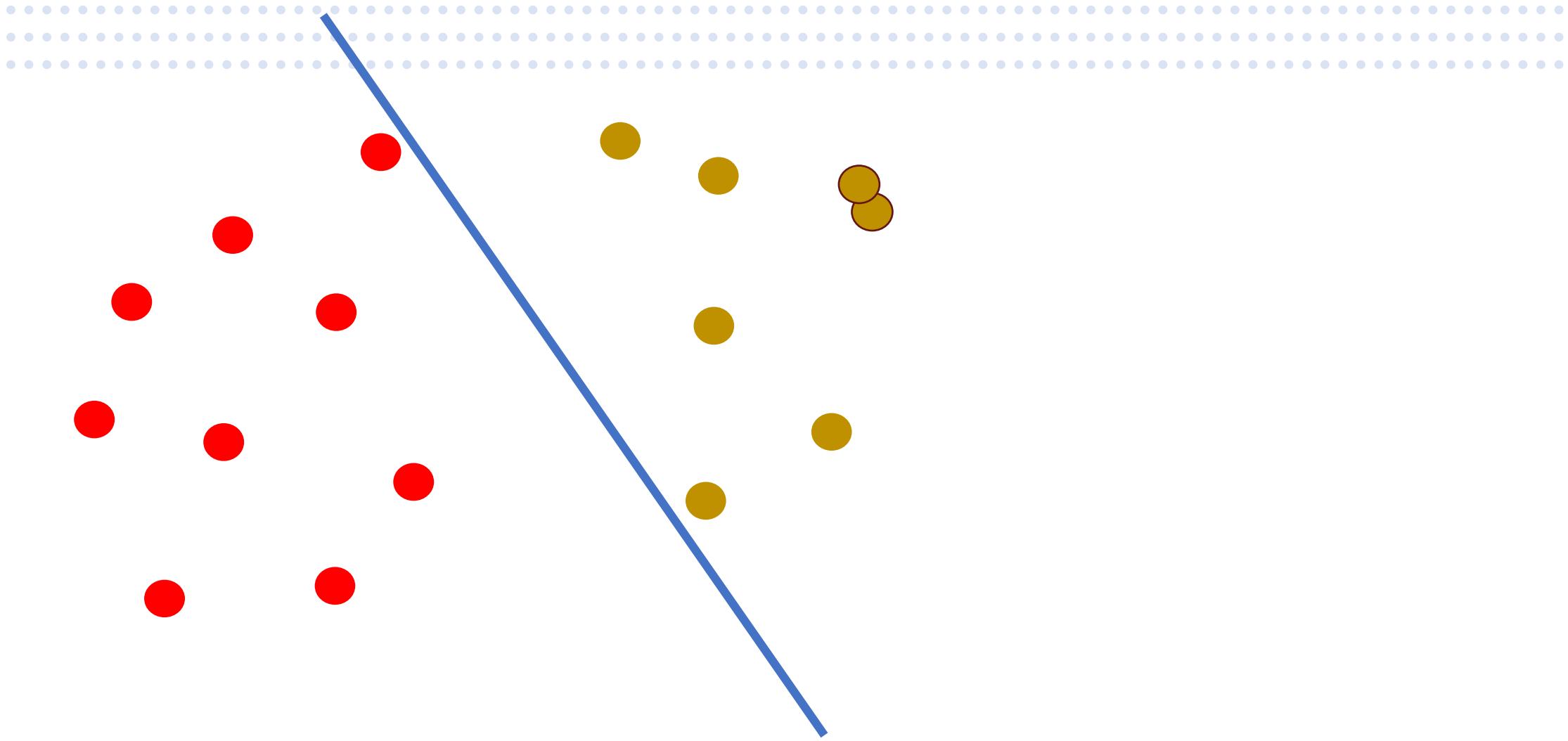
Review::



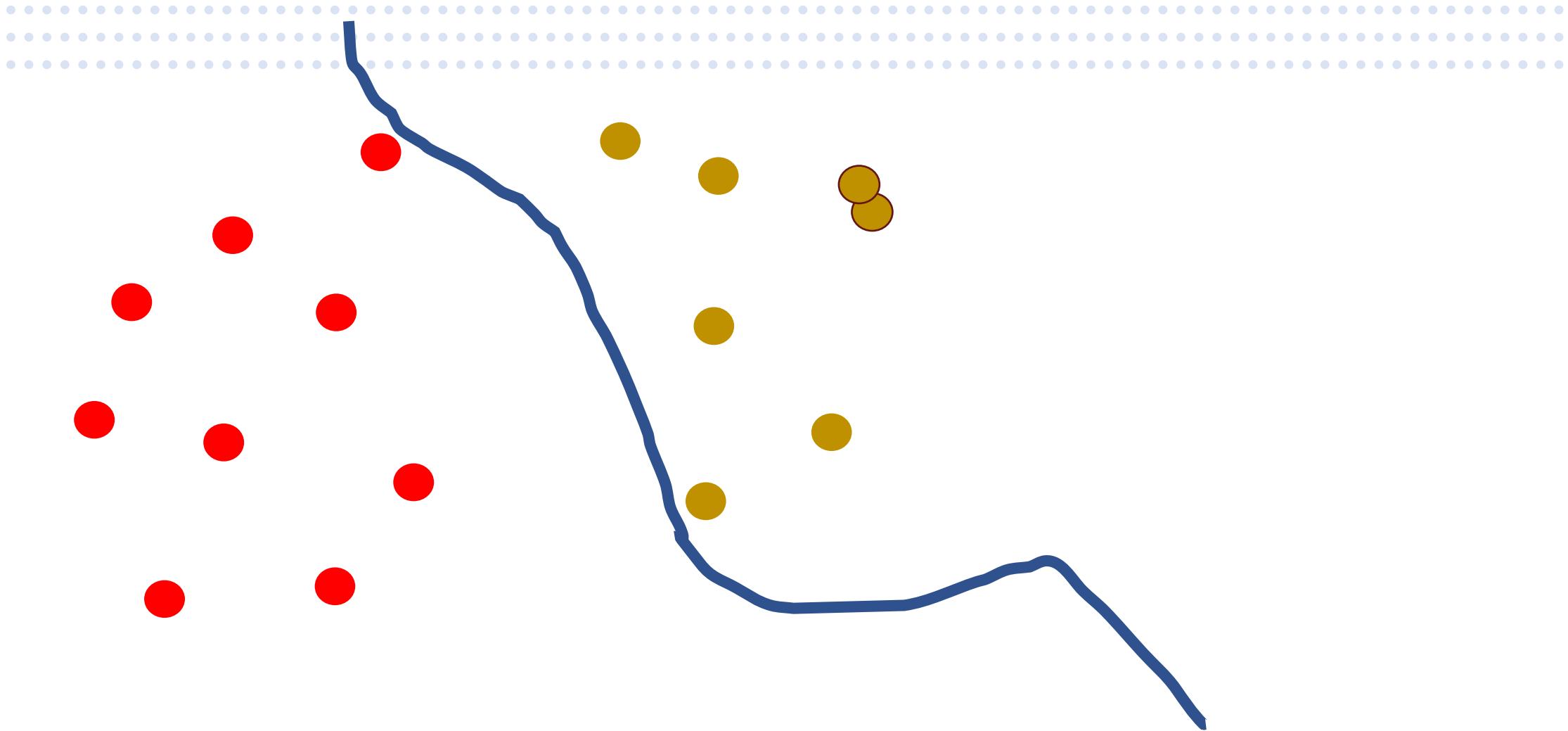
Review::



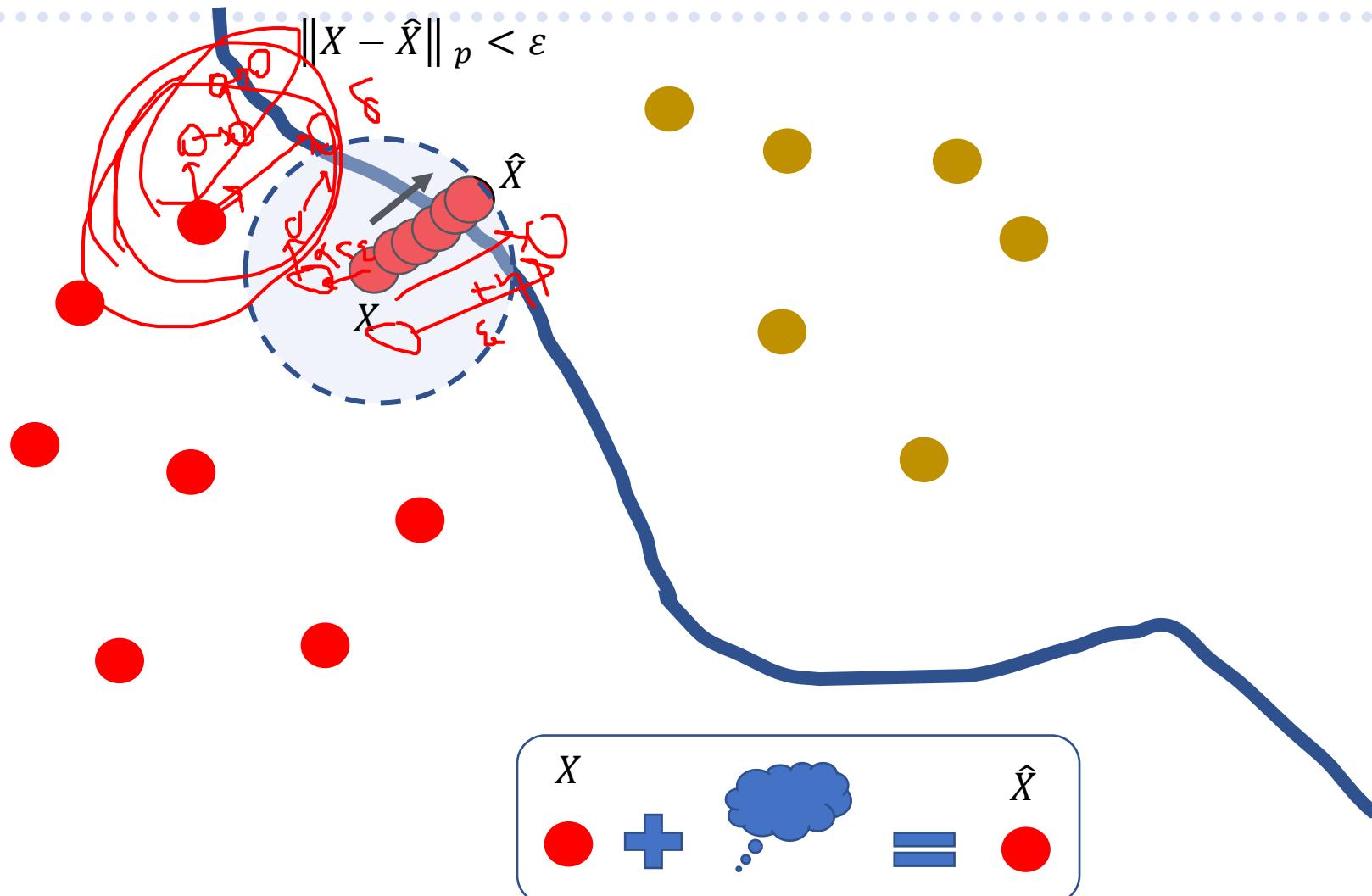
Review::



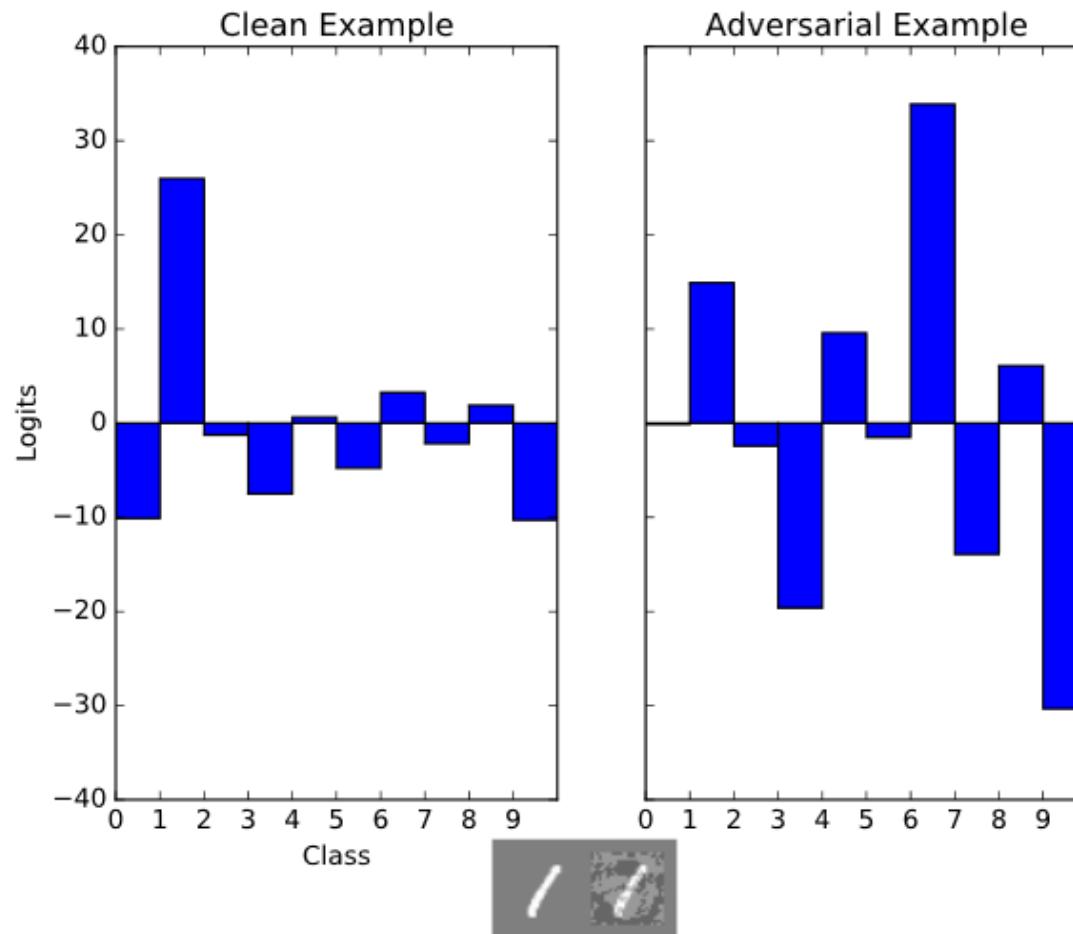
Review::



What is Adversarial Example?



Summary



Attack formulation

Given $x, f_\theta(\cdot)$, the task is to compute x' such that

$$c^* = f_\theta(x) \neq f_\theta(x')$$

with some constraint like $\|x - x'\|_{\ell_p} \leq \varepsilon$ to impose *imperceptibility*
For ℓ_p attacks

Attack formulation

$$\max_{\delta} l_{cls} (f_{\theta} (x'), y)$$

```
graph TD; Perturbation --> Loss; Loss --> Classifier; Classifier --> AdvExample[Adv. Example]; AdvExample --> Label; Label --> Loss;
```

Single-step attack

$$\begin{aligned}\delta_{\text{FGSM}} &= \max_{\|\delta\| \leq \varepsilon} \langle \nabla l(f_\theta(x), y), \delta \rangle \\ &= \varepsilon \cdot \text{sign}(\nabla l(f_\theta(x), y))\end{aligned}$$

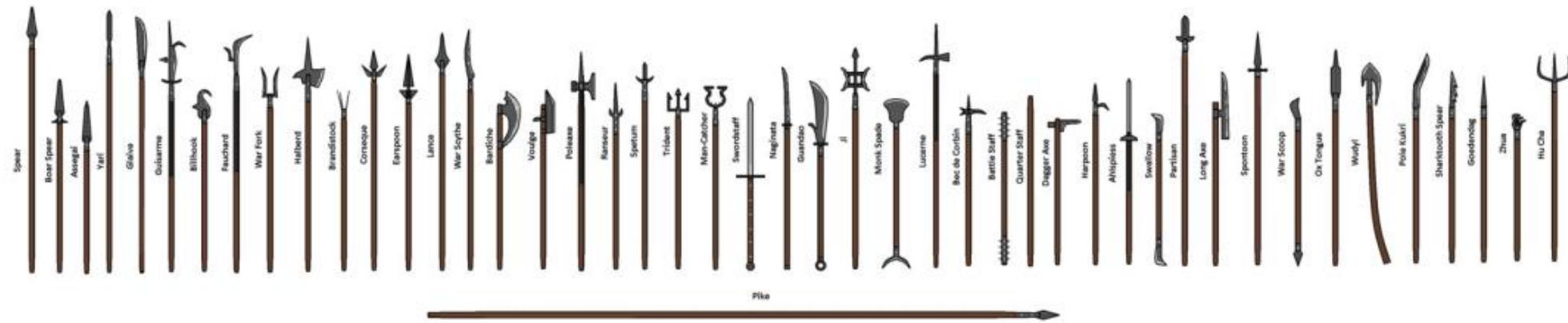
- East GSign Method (Goodfellow et al., ICLR'15).
- Specifically designed for ℓ_∞ attacks.
- One-step attack.

Single-step attack

$$x' = x + \frac{\alpha \text{sign}(\nabla_x L(\theta, x, y))}{\text{FGSM}}$$

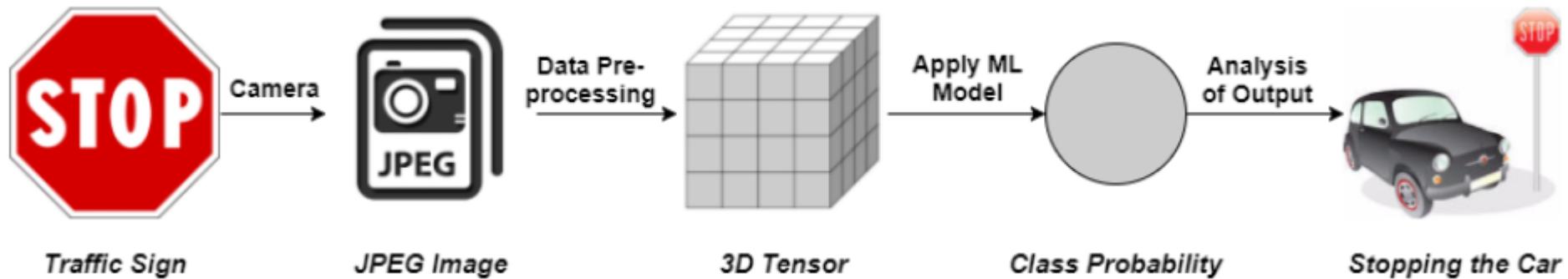


Various methods of attacks



Adversarial Threat Model

- The Attack Surface



- **Evasion Attack** :: during the testing phase (* the most common type of attack!)
- **Poisoning Attack** :: during the training time
- **Exploratory Attack** :: during the testing phase (Given black box access to the model
try to gain as much knowledge as possible)

Adversarial Attacks



- **Whitebox:** Attacker has access to the model parameters, outputs, etc.
- **Blackbox:** Attacker has only *query* access to the model and its outputs.

One important notes!



```
acc, loss = model.evaluate(x_test, y_test)
```

Is no longer sufficient.

<https://blog.floydhub.com/introduction-to-adversarial-machine-learning/>





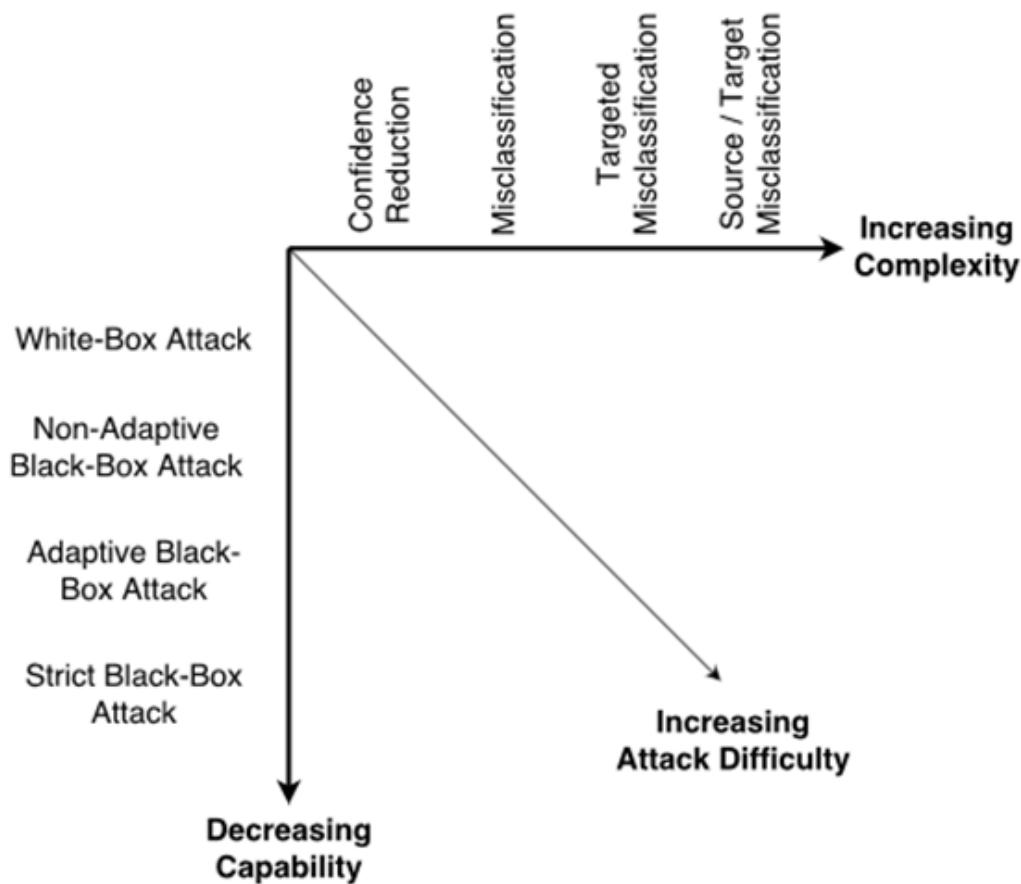
Break



Adversarial Threat Model

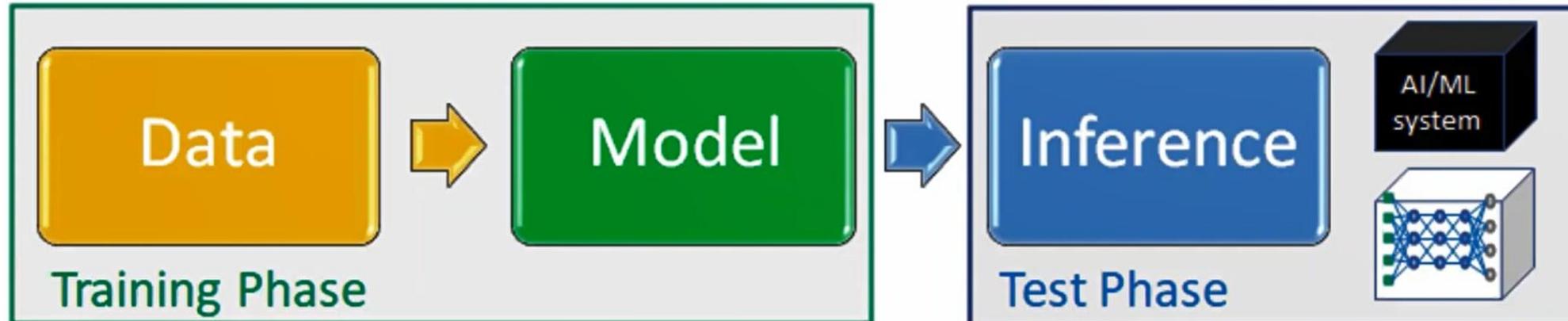
- **Adversarial Goals:**
 - **Confidence Reduction**
 - The adversary tries to reduce the confidence of prediction for the target model
 - **Misclassification**
 - The adversary tries to alter the output classification of an input example to **any** class different from the original class.
 - **Targeted Misclassification**
 - The adversary tries to produce inputs that force the output of the classification model to be a **specific** target class
 - **Source/Target Misclassification**
 - The adversary attempts to force the output of classification for a **specific** input to be a **particular** target class

Adversarial Threat Model



Attack Difficulty with respect to adversarial
capabilities and goals for Evasion Attacks

Holistic view of adversarial robustness

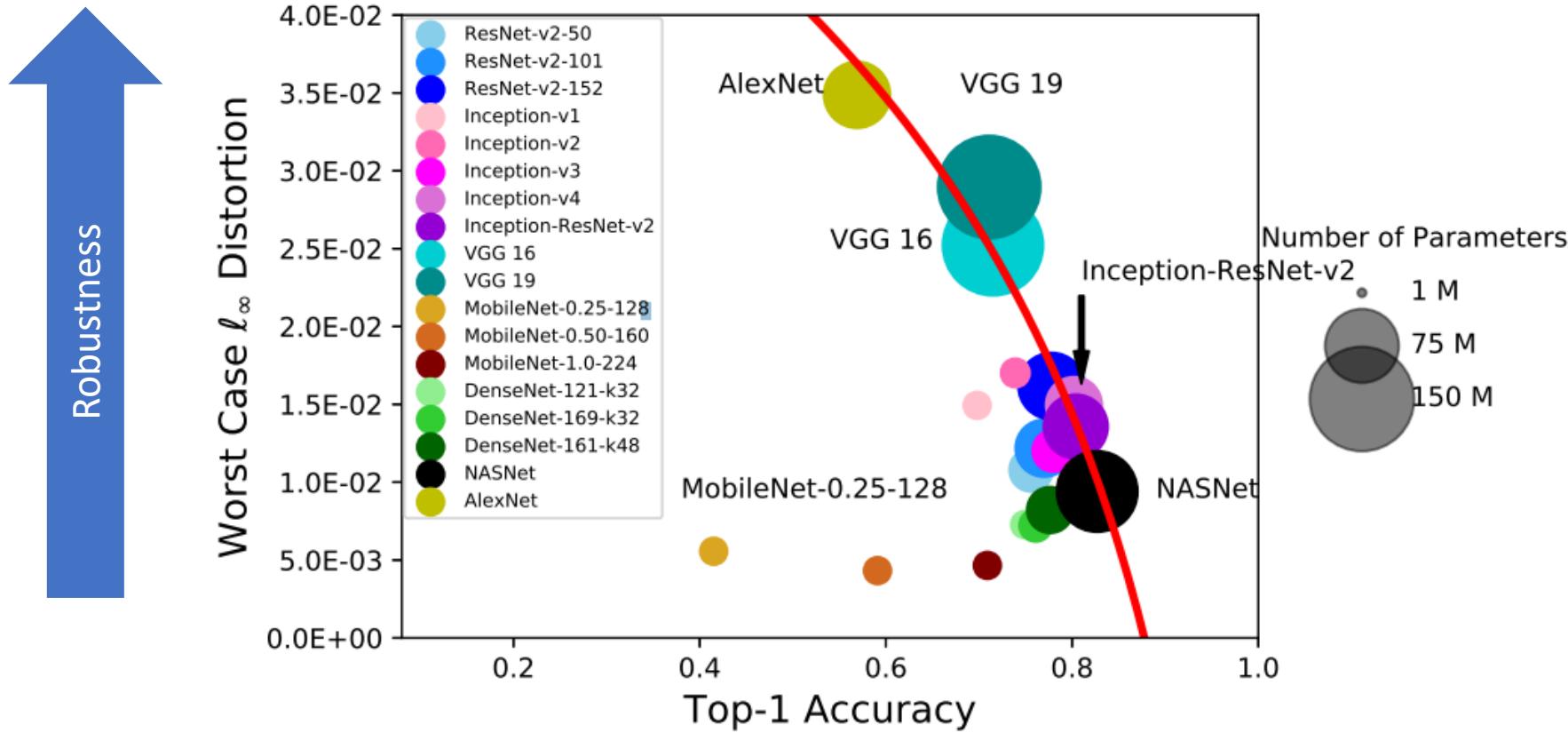


Attack Category / Attacker's reach	Data	Model / Training Method	Inference
Poisoning Attack [learning]	X	X*	
Backdoor Attack [learning]	X		
Evasion Attack (Adversarial Example) [learning]		X*	X
Extraction Attack (Model Stealing, Membership inference)			X
Model Injection [AI governance]		X*	X

IBM Research AI

*No access to model internal information in the black-box attack setting

Accuracy vs Adversarial Robustness



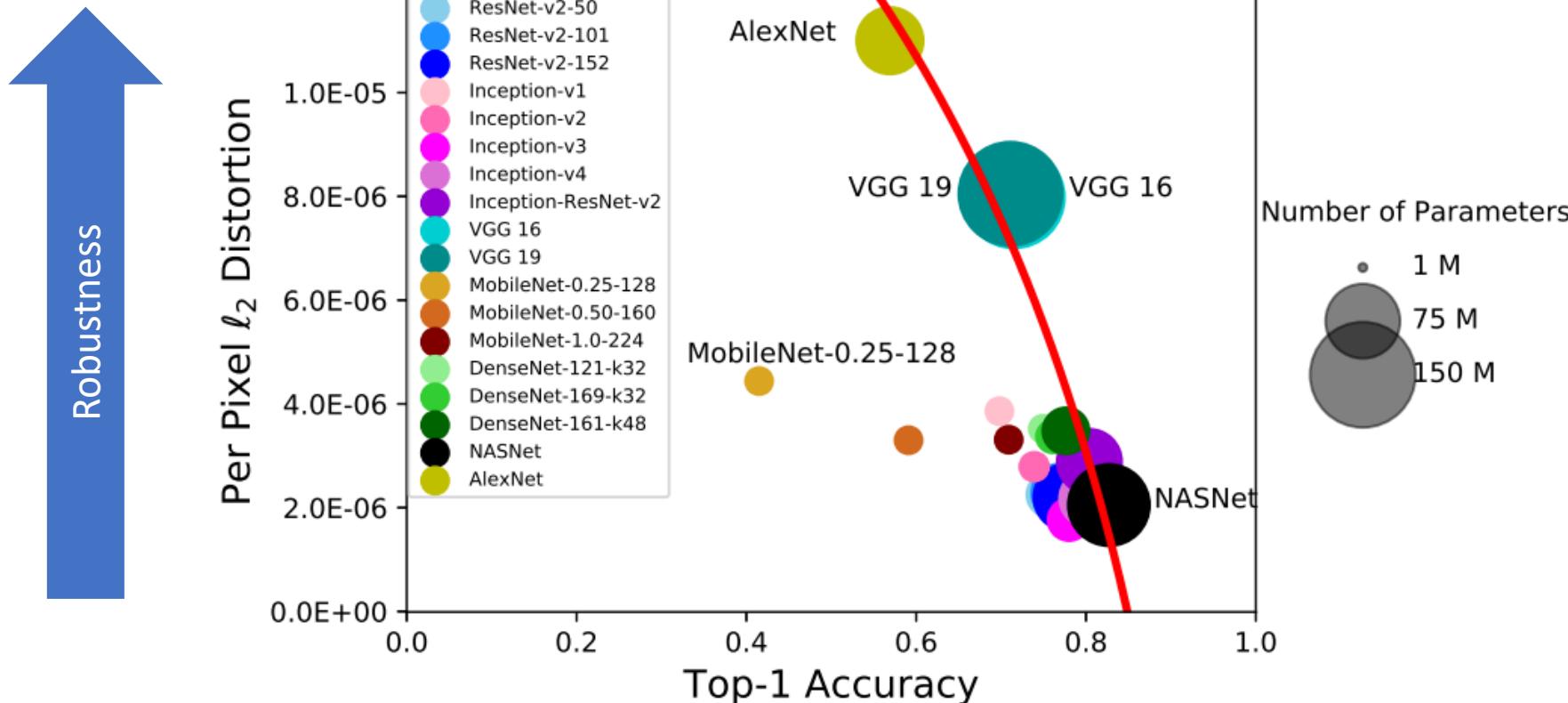
Our benchmark on 18 ImageNet models reveals a tradeoff in accuracy and robustness

(a) Fitted Pareto frontier of ℓ_∞ distortion (I-FGSM attack) vs. top-1 accuracy:

$$\ell_\infty \text{ dist} = [2.9 \cdot \ln(1 - \text{acc}) + 6.2] \times 10^{-2}$$

https://openaccess.thecvf.com/content_ECCV_2018/papers/Dong_Su_Is_Robustness_the_ECCV_2018_paper.pdf

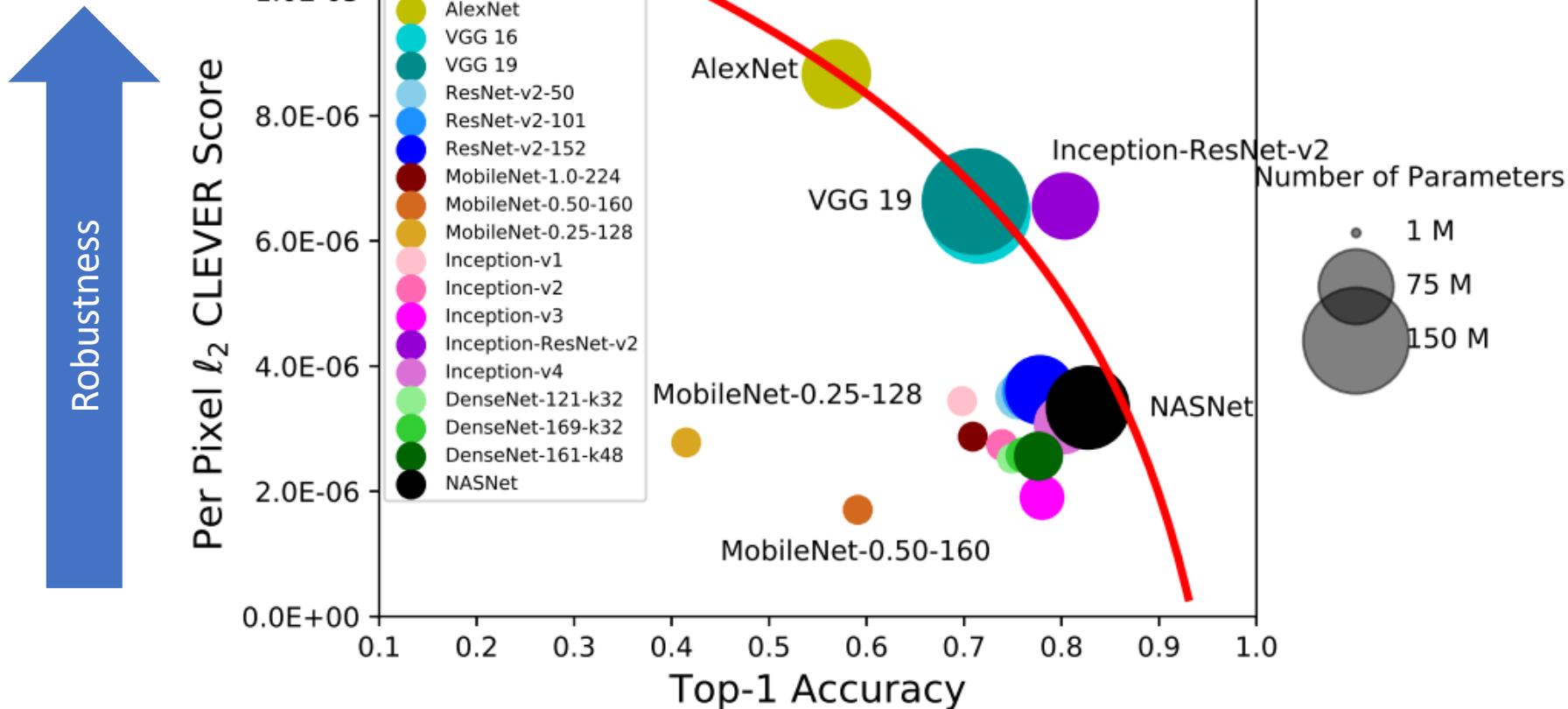
Accuracy vs Adversarial Robustness



Our benchmark on 18 ImageNet models reveals a tradeoff in accuracy and robustness

(b) Fitted Pareto frontier of ℓ_2 distortion (C&W attack) vs. top-1 accuracy:
$$\ell_2 \text{ dist} = [1.1 \cdot \ln(1 - \text{acc}) + 2.1] \times 10^{-5}$$

Accuracy vs Adversarial Robustness

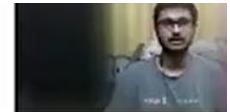


Our benchmark on 18 ImageNet models reveals a tradeoff in accuracy and robustness

(c) Fitted Pareto frontier of ℓ_2 CLEVER score vs. top-1 accuracy:

$$\ell_2 \text{ score} = [4.6 \cdot \ln(1 - \text{acc}) + 12.5] \times 10^{-6}$$

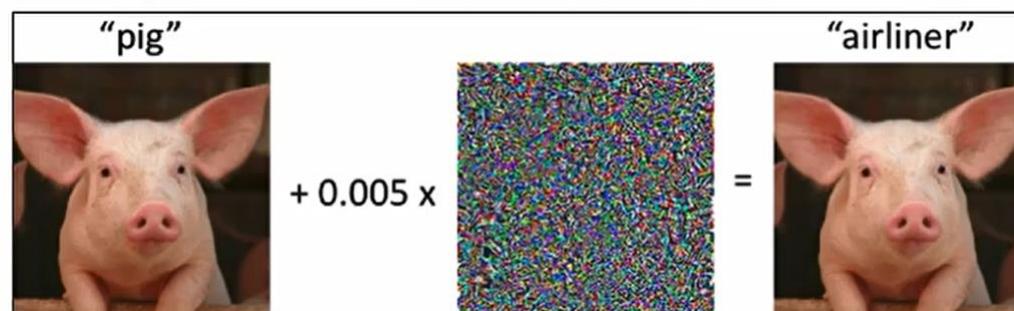
Setting up the attack formulation problem



Given $x, f_\theta(\cdot)$, the task is to compute x' such that

$$c^* = f_\theta(x) \neq f_\theta(x')$$

with some constraint like $\|x - x'\|_{\ell_p} \leq \varepsilon$ to impose **imperceptibility**
For ℓ_∞ attacks



Targeted vs Untargeted

- An attacker can either launch targeted or an untargeted attacks.
- In targeted attacks, attack can set t where $f_{\theta}(x') = t \neq c^*$.

Given $x, f_{\theta}(\cdot)$, the task is to compute x' such that

$$c^* = f_{\theta}(x) \neq f_{\theta}(x')$$

with some constraint like $\|x - x'\|_{\ell_p} \leq \varepsilon$ to impose ***imperceptibility***
For ℓ_p attacks

Adversarial attack

$$\max_{\delta} l_{cls} (f_{\theta} (x'), y)$$

The diagram illustrates the components of an adversarial attack. At the top is the mathematical expression $\max_{\delta} l_{cls} (f_{\theta} (x'), y)$. Below it, five labels are arranged horizontally: 'Perturbation', 'Loss', 'Classifier', 'Adv. Example', and 'Label'. Arrows point from each label to its corresponding part in the expression: 'Perturbation' points to the δ in \max_{δ} ; 'Loss' points to the l_{cls} ; 'Classifier' points to the f_{θ} ; 'Adv. Example' points to the x' ; and 'Label' points to the y .

Maximize loss between a classifier's prediction on adversarial examples and their labels.

Adversarial attack

- For ℓ_p attacks, the following constraint holds: $\|x - x'\|_{\ell p} \leq \varepsilon$.
- There are non ℓ_p attacks too:
 - Adversarial translations
 - Functional adversarial attacks
 - Wasserstein attacks
 - Evolution attacks
- In this tutorial, we will primarily be dealing with ℓ_p attacks. These come under gradient-based attacks.

How do we perform optimization?

Case I: Single-step attack

$$\begin{aligned}\delta_{\text{FGSM}} &= \max_{\|\delta\| \leq \varepsilon} \langle \nabla l(f_\theta(x), y), \delta \rangle \\ &= \varepsilon \cdot \text{sign}(\nabla l(f_\theta(x), y))\end{aligned}$$

- East GSign Method (Goodfellow et al., ICLR'15).
- Specifically designed for ℓ_∞ attacks.
- One-step attack.

How do we perform optimization?

Case II: Multi-step attack

$$x^{t+1} = \underbrace{\Pi_{x+\mathcal{S}}(x^t + \alpha \text{sign}(\nabla_x L(\theta, x, y)))}_{\text{Projection}} \quad \underbrace{\text{FGSM}}$$

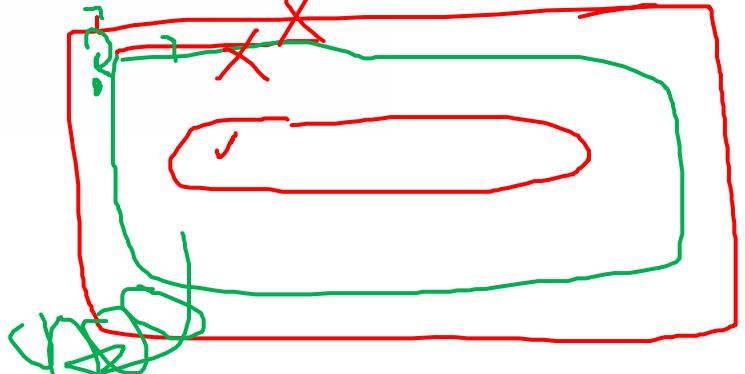
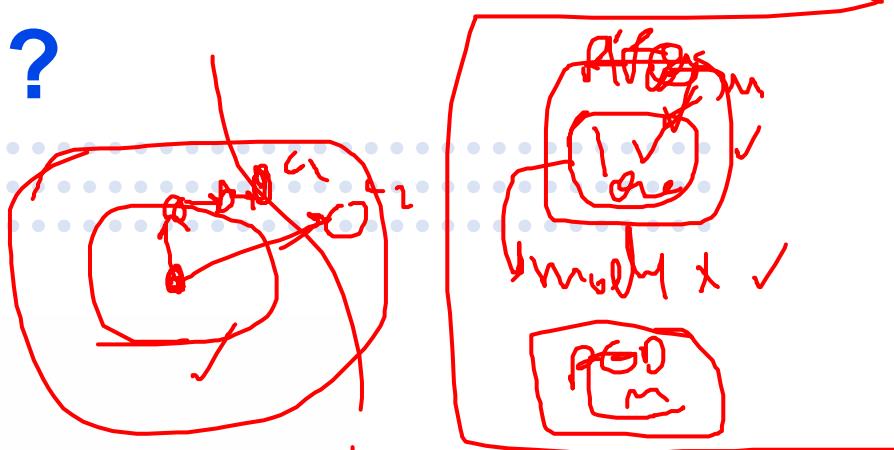
How do we perform optimization?

Case II: Multi-step attack

$$x^{t+1} = \underbrace{\Pi_{x+\mathcal{S}}}_{\text{Projection}} \left(x^t + \underbrace{\alpha \text{sign}(\nabla_x L(\theta, x, y))}_{\text{FGSM}} \right)$$

For ℓ_∞ , Π is $\text{clip}(x^{t+1}, x^0 - \epsilon, x^0 + \epsilon)$

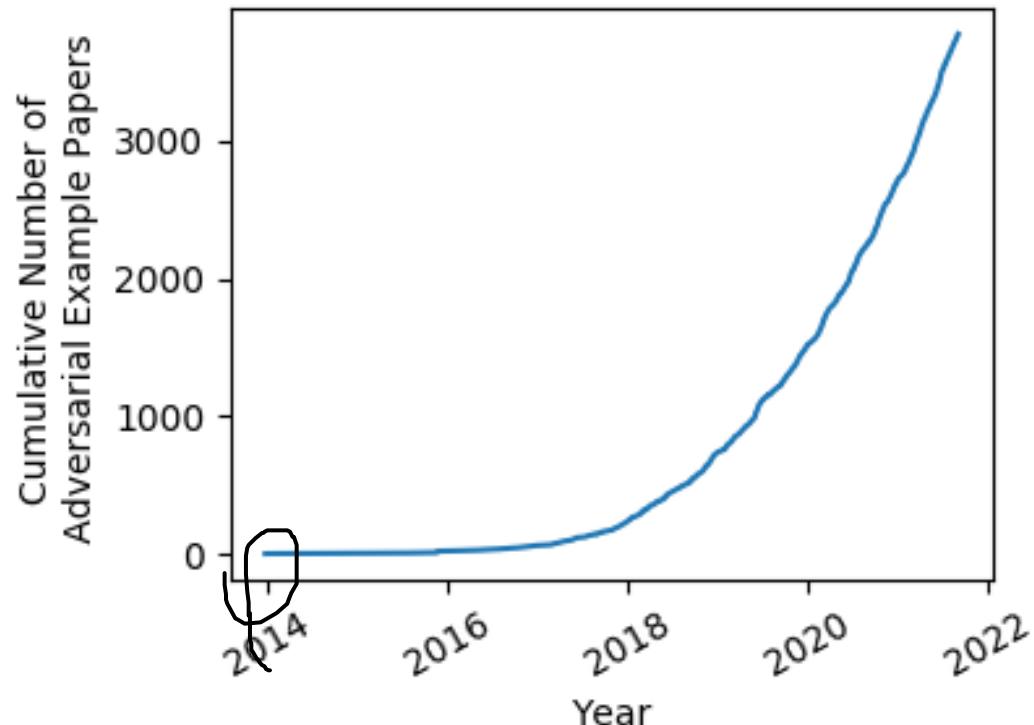
- Projected Gradient Descent (PGD, Madry et al., ICLR'18)
- Since α (step size) can break the norm-ball constraint, for ℓ_p attacks, projection is required after adding the perturbation.
- Hence Projected Gradient Descent (PGD).



How do we perform optimization?

- There are other and more recent of gradient-based attacks trying to solve the maximization problem in clever ways.
- Challenges remain in defining precisely what human perception is.
- This is important because a strong adversarial attack should keep the changes as imperceptible as possible.

Adversarial example papers



<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>

Adversarial example papers



Nicholas Carlini

Research Scientist, Google Brain

- <https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>



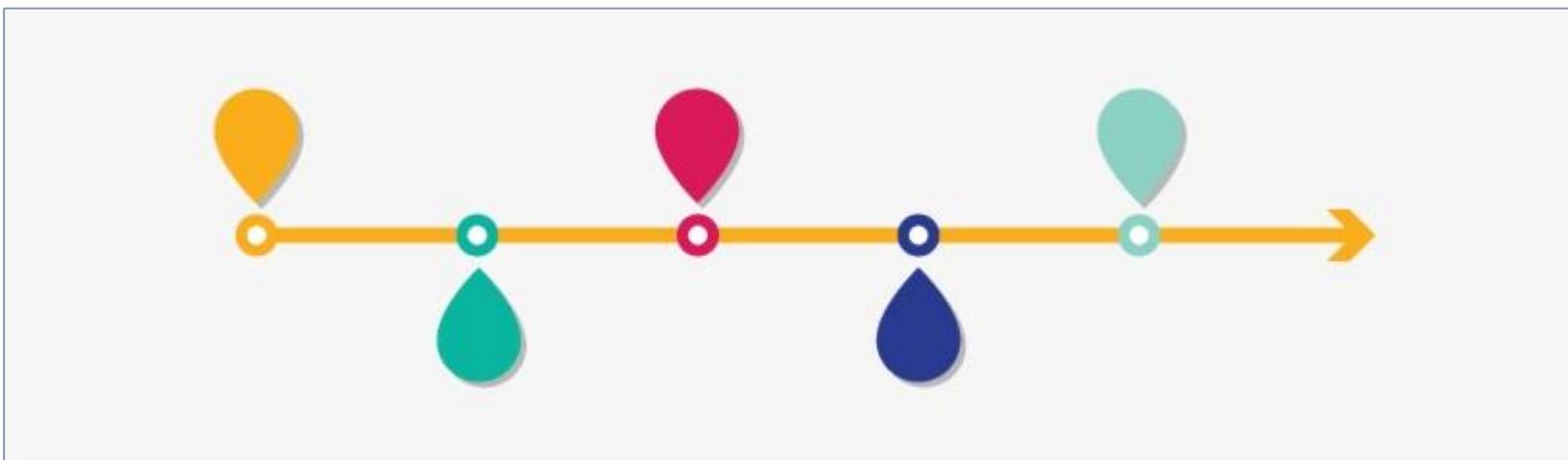
- <https://www.robust-ml.org/defenses/>
- <https://paperswithcode.com/task/adversarial-attack>



Break



Adversarial example papers

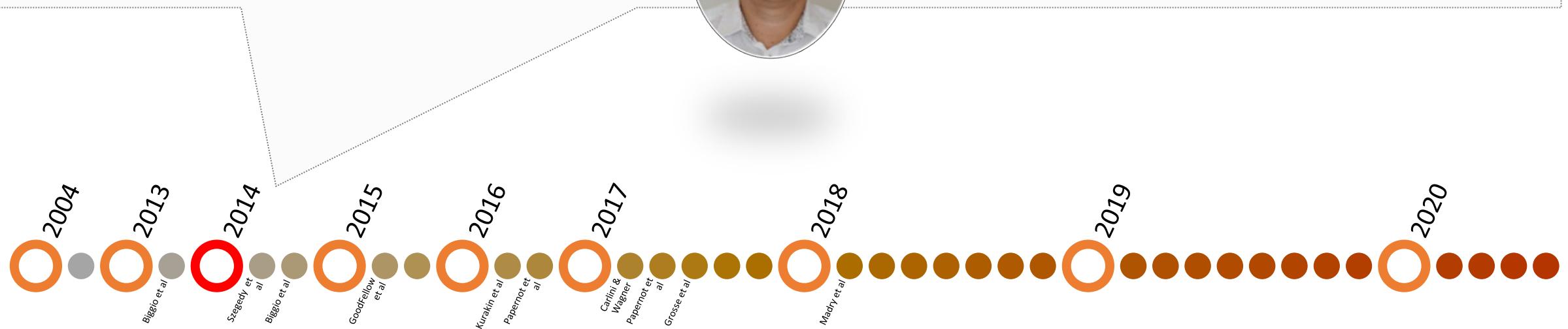
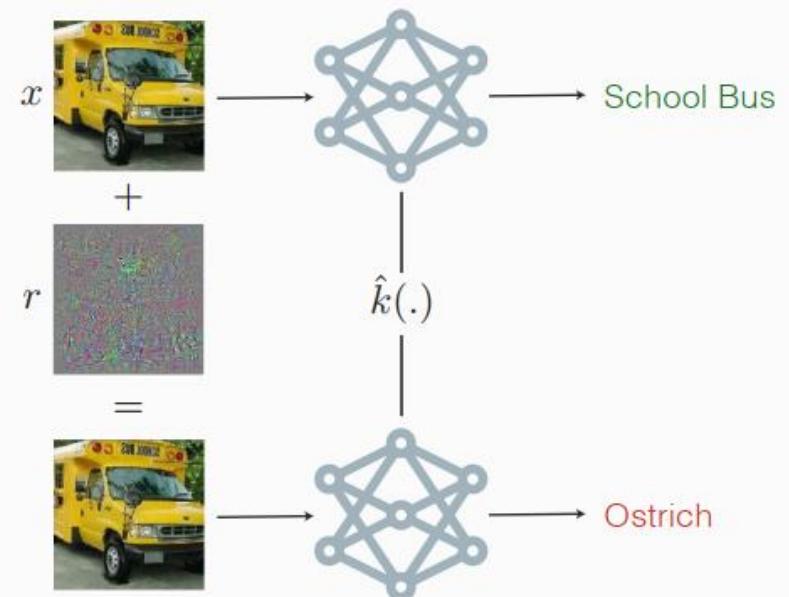


Adversarial example papers



Interiguing properties of Neural Network
C. Szegedy et al. ICLR, 2014

$$r^* = \arg \min_r J(\hat{k}(x + r), y) + C\|r\|$$



Adversarial example papers



Explaining and Harnessing Adversarial Example
I.J. Goodfellow et al. ICLR, 2015



Adversarial example papers



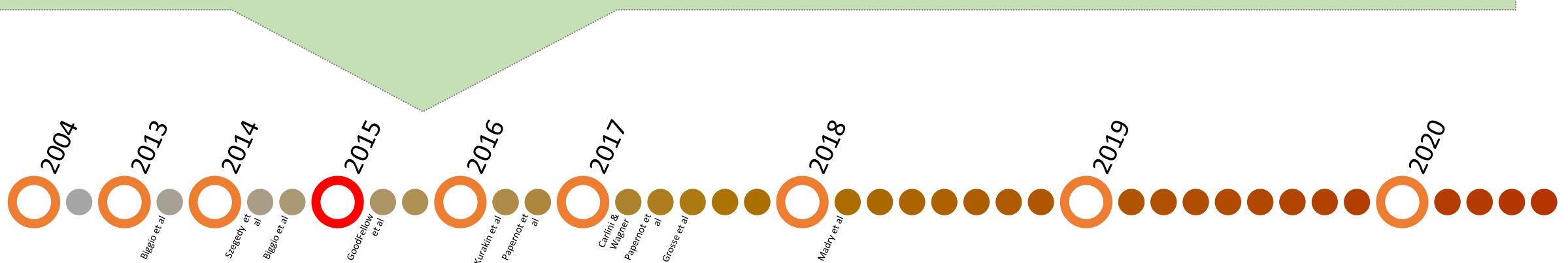
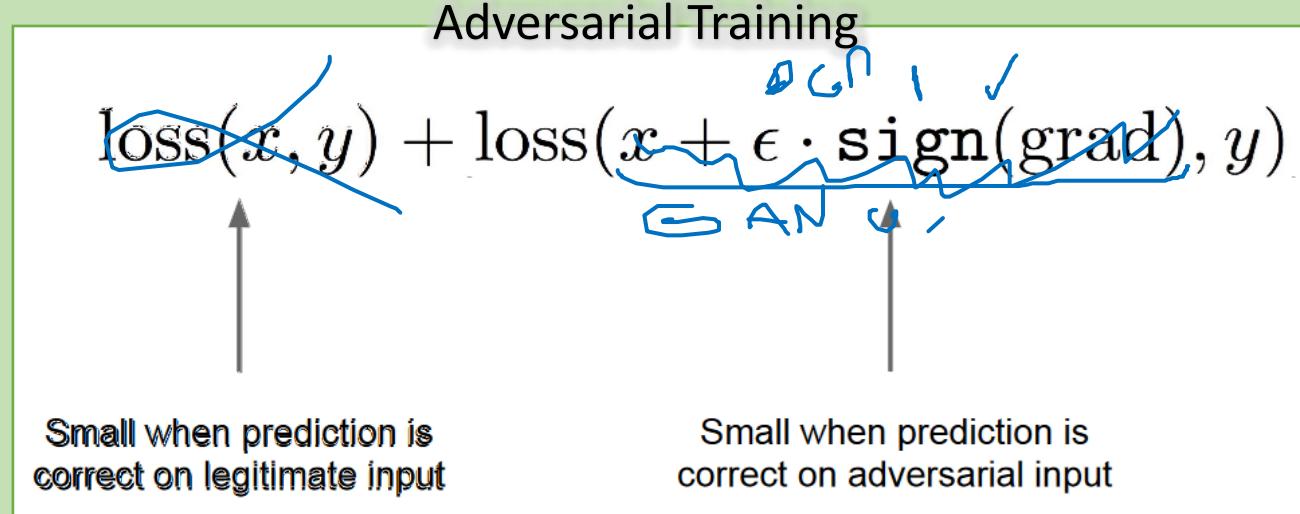
Explaining and Harnessing Adversarial Example
I.J. Goodfellow et al. ICLR, 2015

$$\tilde{x} = x + \eta$$

Perturbation

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

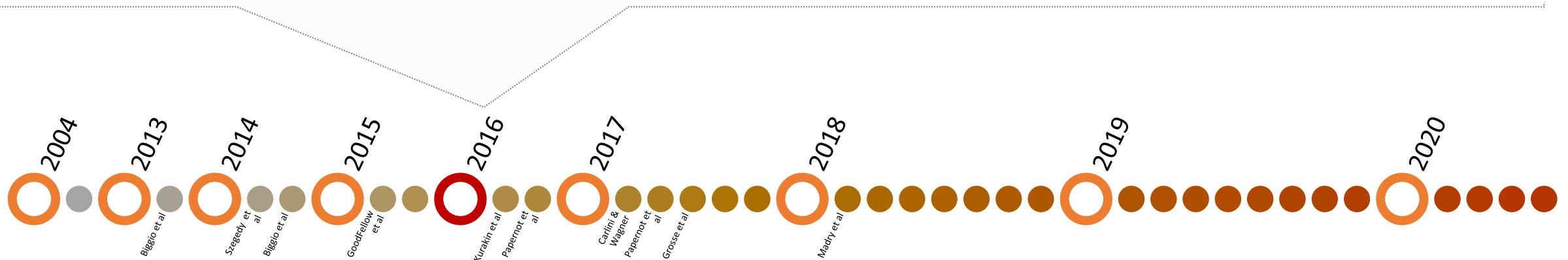
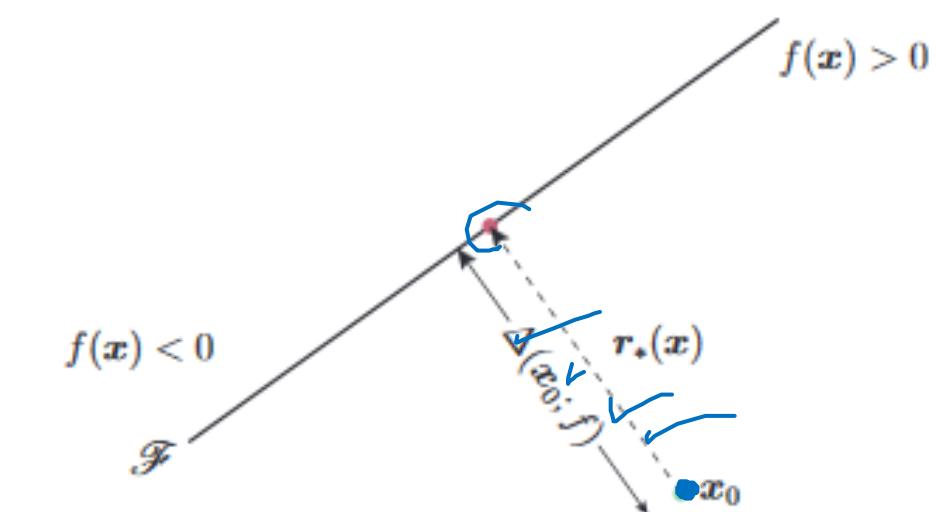
Gradient of the cost function



Adversarial example papers



Deep Fool: a simple and accurate method to fool deep neural networks
S.M. Moosavi-dezfooli et al. CVPR, 2016



Algorithm 1 DeepFool for binary classifiers

```
1: input: Image  $x$ , classifier  $f$ .  
2: output: Perturbation  $\hat{r}$ .  
3: Initialize  $x_0 \leftarrow x$ ,  $i \leftarrow 0$ .  
4: while  $\text{sign}(f(x_i)) = \text{sign}(f(x_0))$  do  
5:    $r_i \leftarrow -\frac{f(x_i)}{\|\nabla f(x_i)\|_2^2} \nabla f(x_i),$   
6:    $x_{i+1} \leftarrow x_i + r_i,$   
7:    $i \leftarrow i + 1.$   
8: end while  
9: return  $\hat{r} = \sum_i r_i.$ 
```

Adversarial example papers



Universal adversarial perturbations
S.M. Moosavi-dezfooli et al. CVPR, 2017

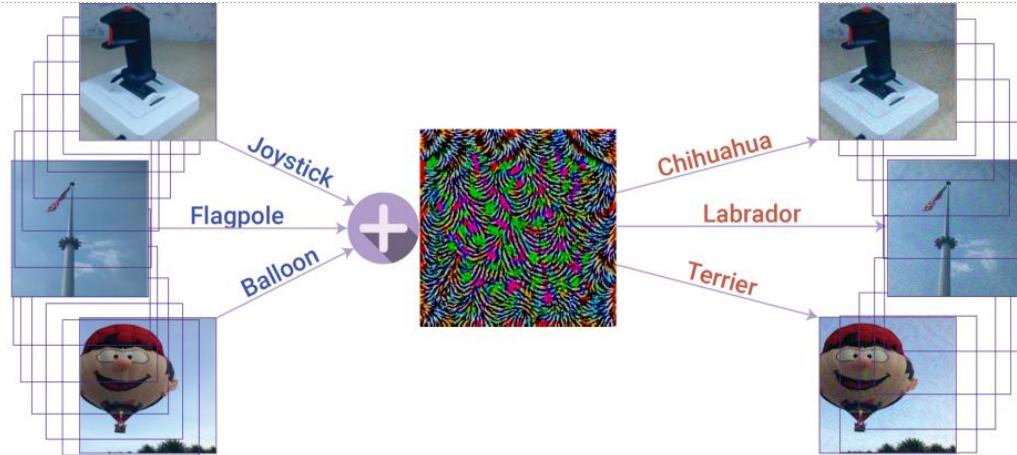
$$f(x + \eta) \neq f(x) \quad \text{for "most" } x \sim X$$

$$\|\eta\|_p \leq \xi$$

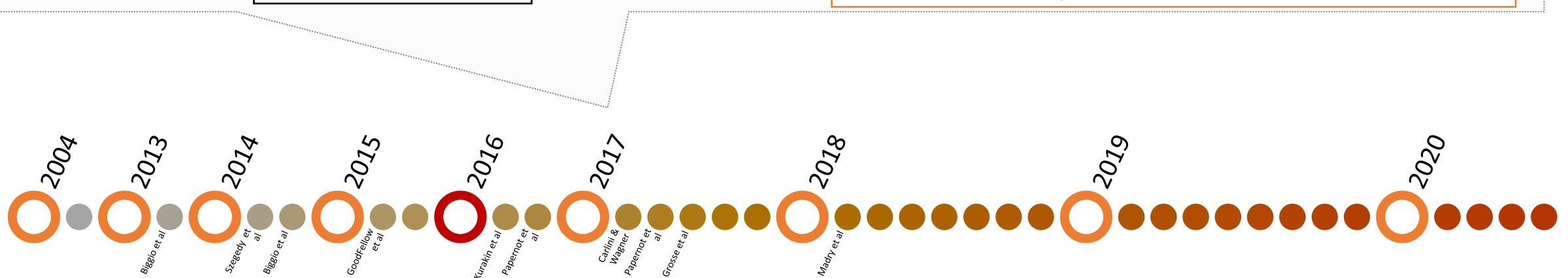
$$\mathbb{P}_{x \sim X} (f(x + \eta) \neq f(x)) \geq 1 - \delta$$

$$\Delta\eta \leftarrow \arg \min_r \|r\|_2 \quad \text{s.t. } f(x_i + \eta + r) \neq f(x_i)$$

$$\eta \leftarrow \mathcal{P}_{p, \xi}(\eta + \Delta\eta_i)$$



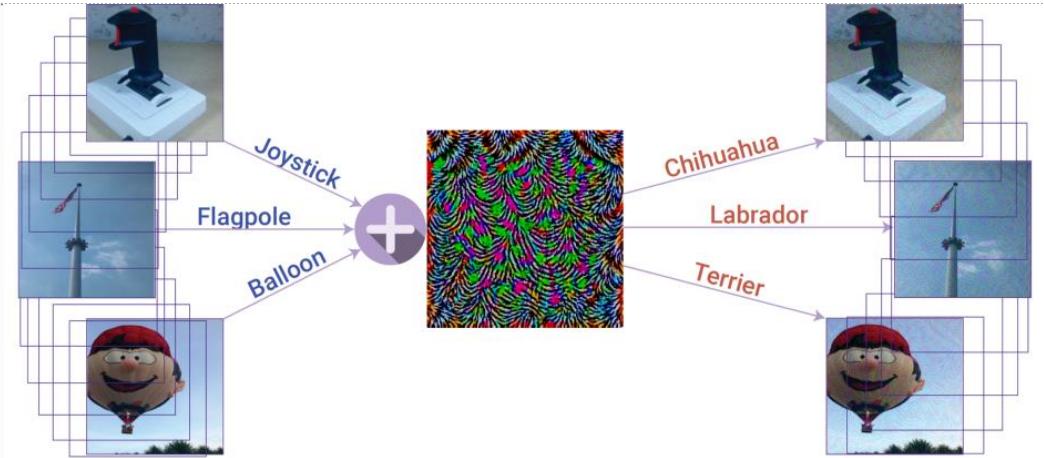
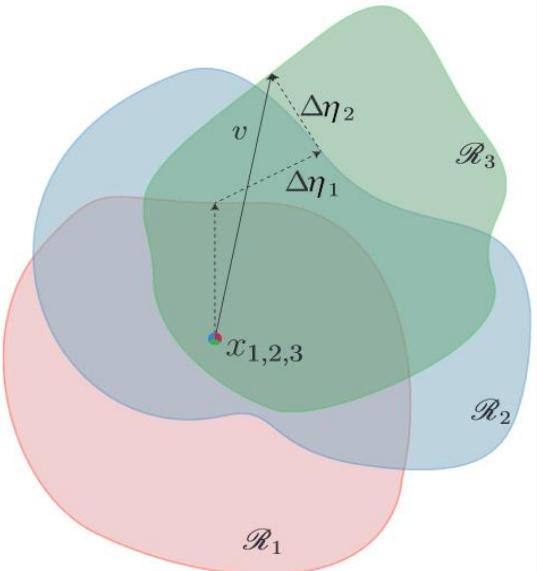
$$\mathcal{P}_{p, \xi}(\eta) = \arg \min_{\eta'} \|\eta - \eta'\|_2 \quad \text{s.t. } \|\eta'\|_p \leq \xi$$



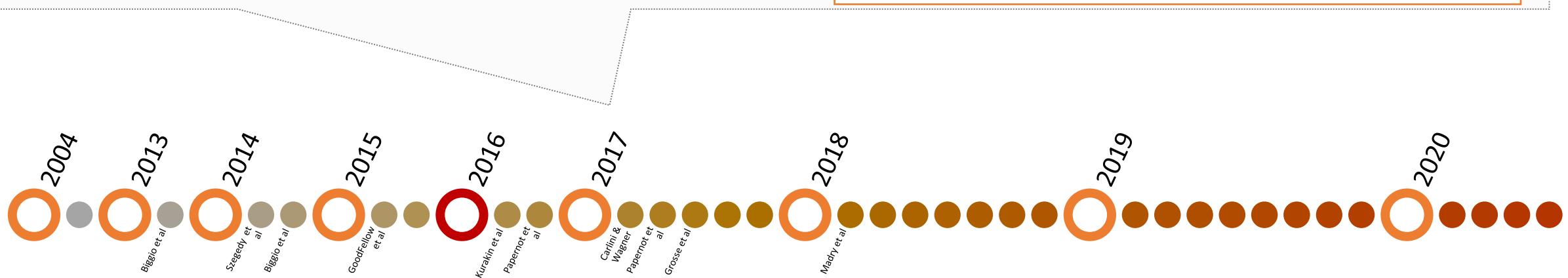
Adversarial example papers



Universal adversarial perturbations
S.M. Moosavi-dezfooli et al. CVPR, 2017



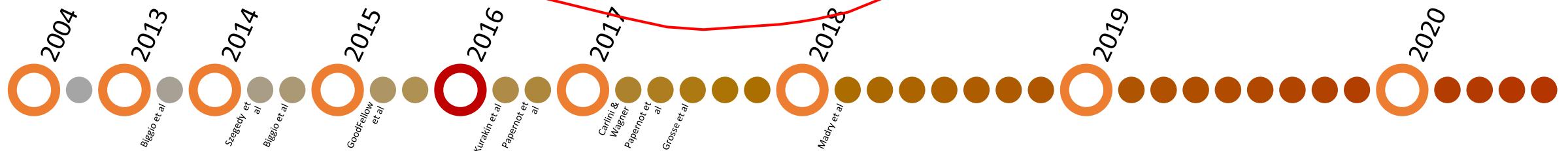
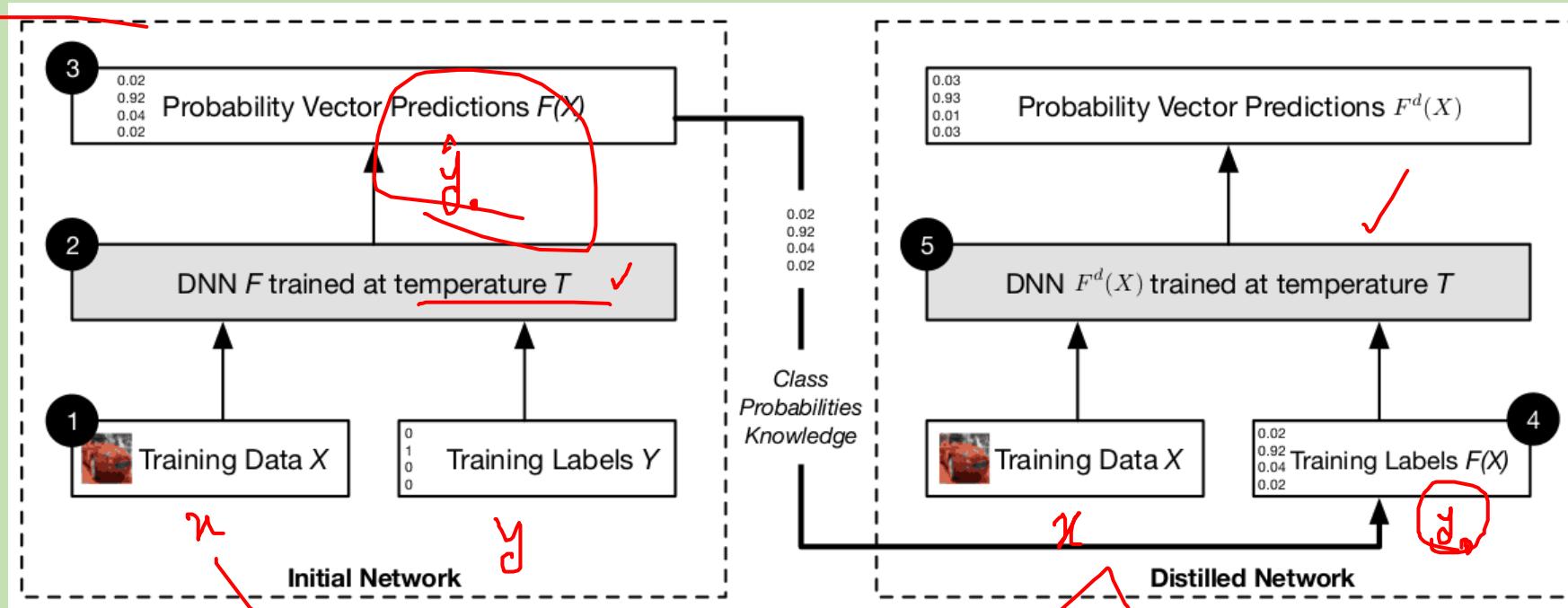
$$\mathcal{P}_{p,\xi}(\eta) = \arg \min_{\eta'} \|\eta - \eta'\|_2 \quad s.t. \quad \|\eta'\|_p \leq \xi$$



Adversarial example papers



Distillation as a Defense to Adversarial Perturbations against Deep Neural Networks
Nicolas Papernot et al. IEEE Symposium on Security and Privacy, 2016



Adversarial example papers



Towards evaluating the robustness of neural networks
Carlini et al. Security and Privacy (S&P). 2017

$$\begin{aligned} &\text{minimize } \mathcal{D}(x, x + \delta) \\ &\text{such that } C(x + \delta) = t \\ &\quad \checkmark \quad \underline{x + \delta \in [0, 1]^n} \end{aligned}$$

$$\begin{aligned} &\text{minimize } \mathcal{D}(x, x + \delta) \\ &\text{such that } \underline{f(x + \delta) \leq 0} \\ &\quad \underline{x + \delta \in [0, 1]^n} \end{aligned}$$

✓

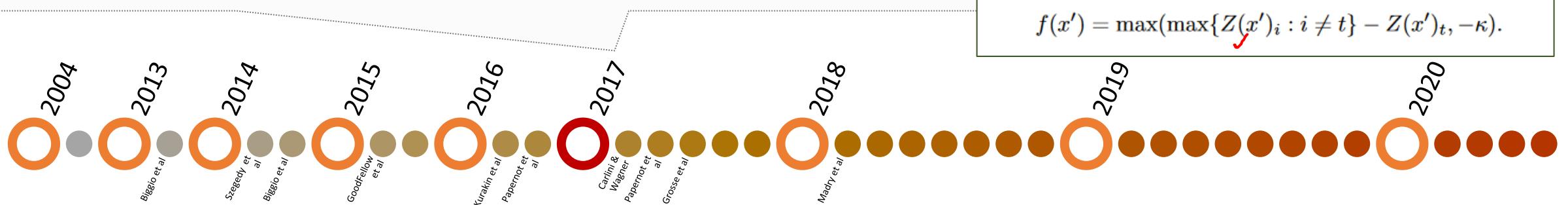
$$\begin{aligned} f_1(x') &= -\text{loss}_{F,t}(x') + 1 \\ f_2(x') &= (\max_{i \neq t}(F(x')_i) - F(x')_t)^+ \\ f_3(x') &= \text{softplus}(\max_{i \neq t}(F(x')_i) - F(x')_t) - \log(2) \\ f_4(x') &= (0.5 - F(x')_t)^+ \\ f_5(x') &= -\log(2F(x')_t - 2) \\ f_6(x') &= (\max_{i \neq t}(Z(x')_i) - Z(x')_t)^+ \\ f_7(x') &= \text{softplus}(\max_{i \neq t}(Z(x')_i) - Z(x')_t) - \log(2) \end{aligned}$$

$$\begin{aligned} &\text{minimize } \mathcal{D}(x, x + \delta) + c \cdot f(x + \delta) \\ &\text{such that } x + \delta \in [0, 1]^n \end{aligned}$$

$$\begin{aligned} &\text{minimize } \|\delta\|_p + c \cdot f(x + \delta) \\ &\text{such that } x + \delta \in [0, 1]^n \end{aligned}$$

$$\delta_i = \frac{1}{2}(\tanh(w_i) + 1) - x_i.$$

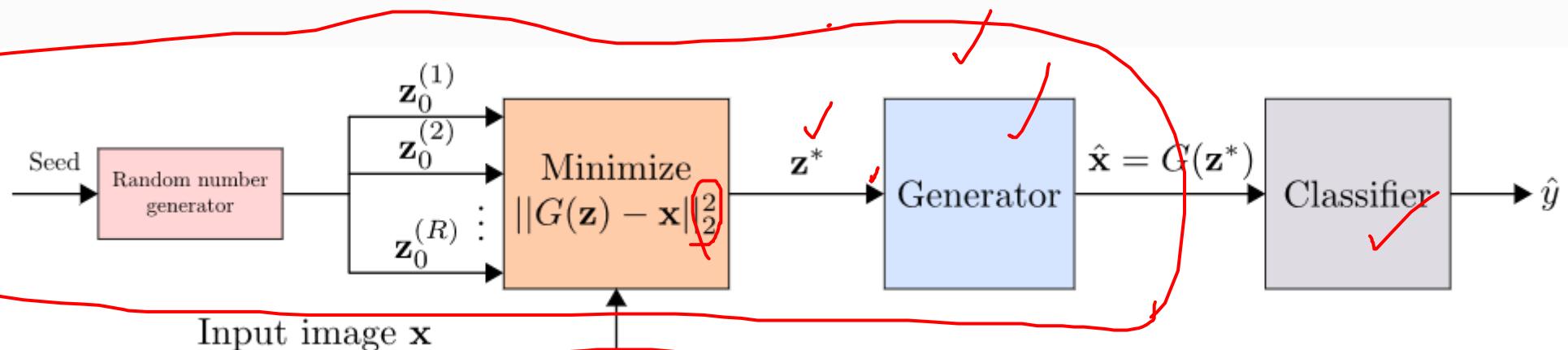
$$\begin{aligned} &\text{minimize } \|\frac{1}{2}(\tanh(w) + 1) - x\|_2^2 + c \cdot f(\frac{1}{2}(\tanh(w) + 1)) \\ &\text{with } f \text{ defined as} \\ &f(x') = \max(\max\{Z(x')_i : i \neq t\} - Z(x')_t, -\kappa). \end{aligned}$$



Adversarial example papers



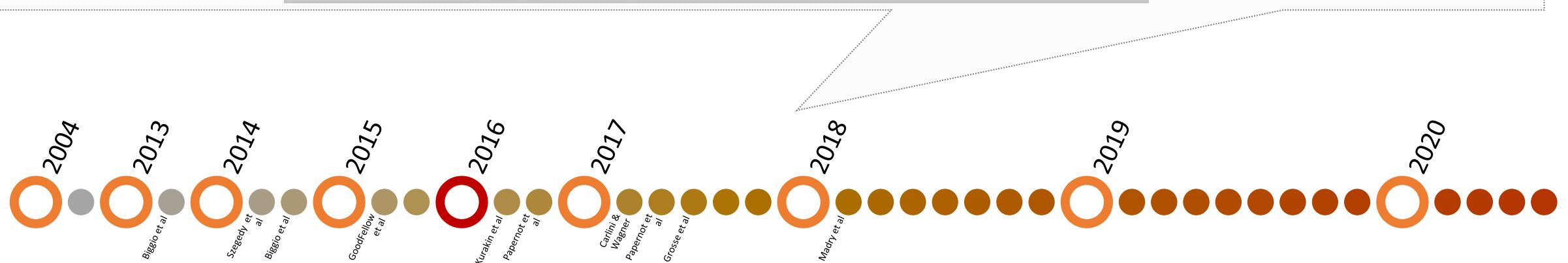
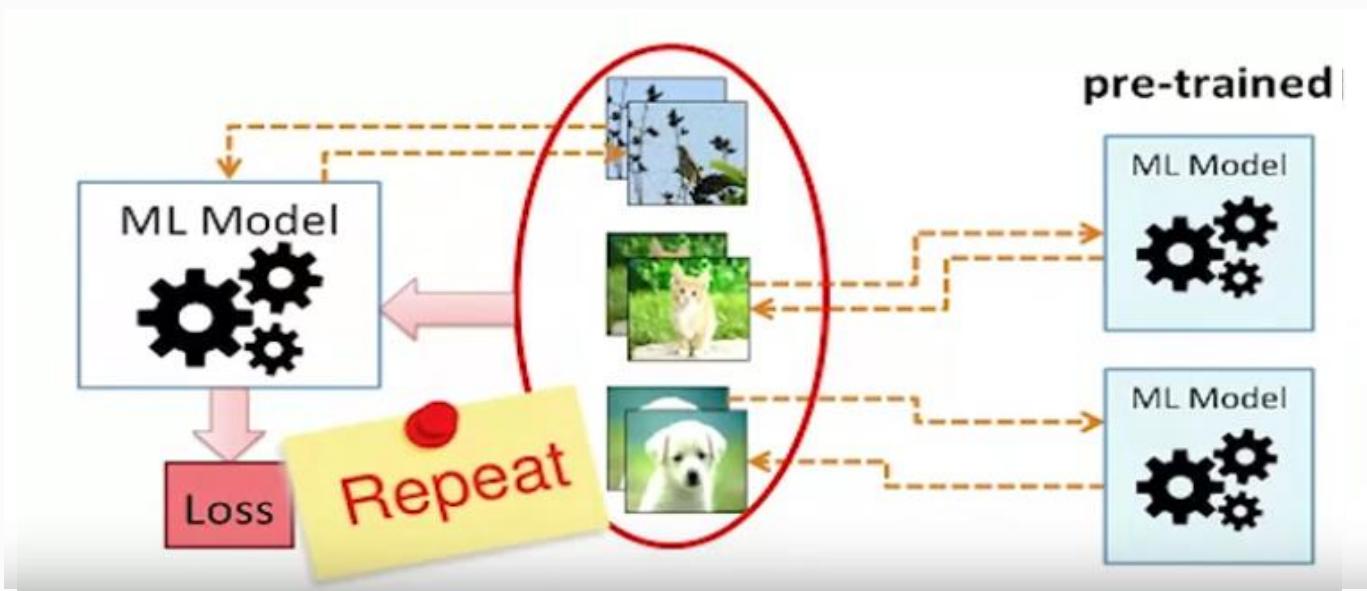
Defense-GAN: Protecting Classifiers Against Adversarial Attacks Using Generative Models
Samangouei et al. ICLR, 2018



Adversarial timeline ...



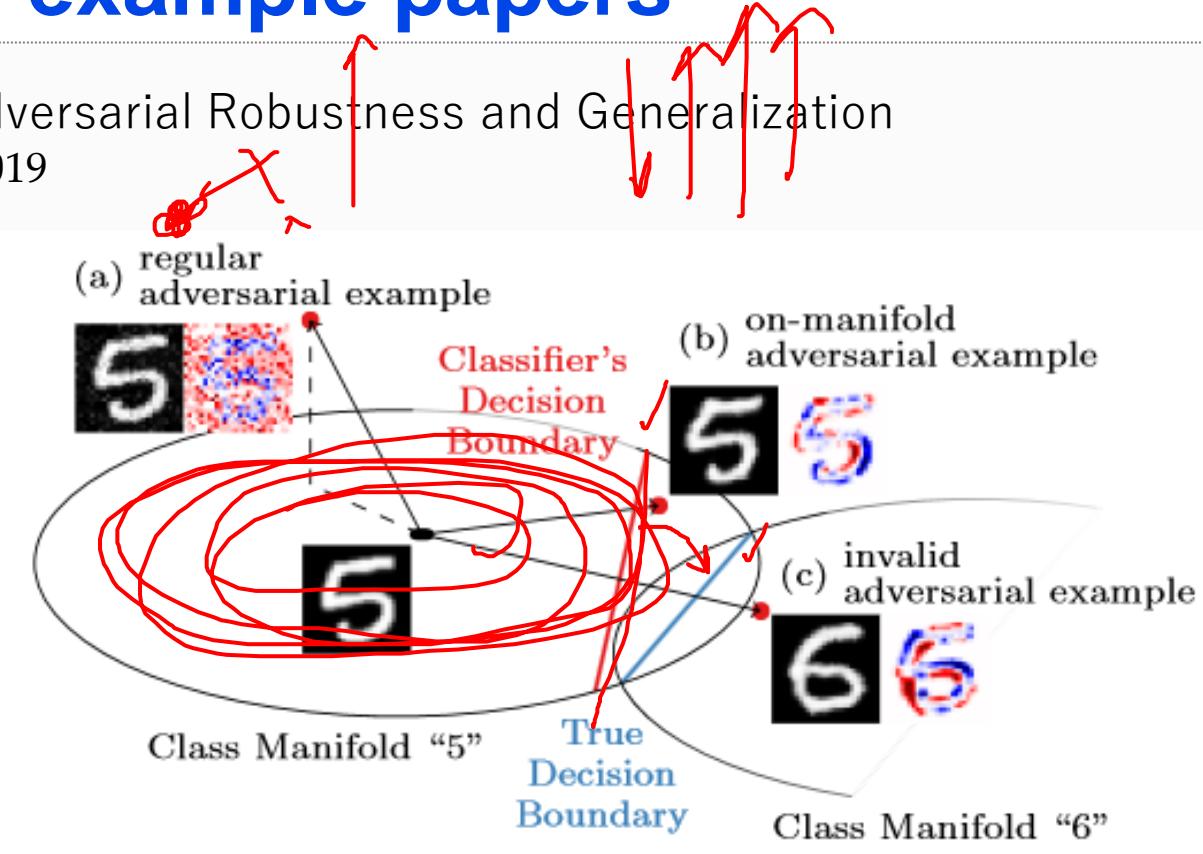
Ensemble Adversarial Training: Attacks and Defenses
Tramèr et al. ICLR, 2018



Adversarial example papers



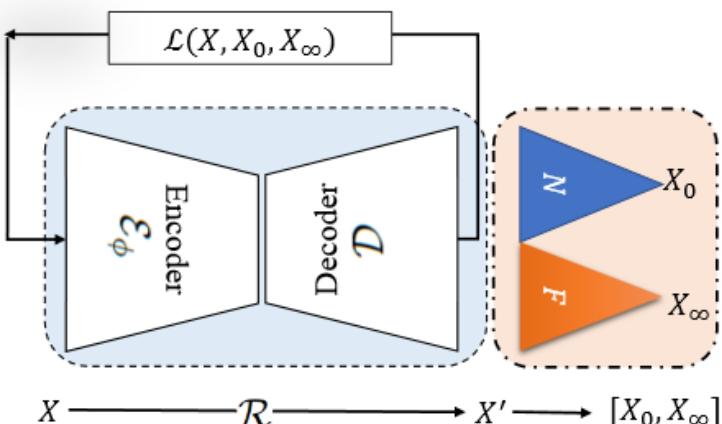
Disentangling Adversarial Robustness and Generalization
Stutz et al. CVPR, 2019



Adversarial timeline ...



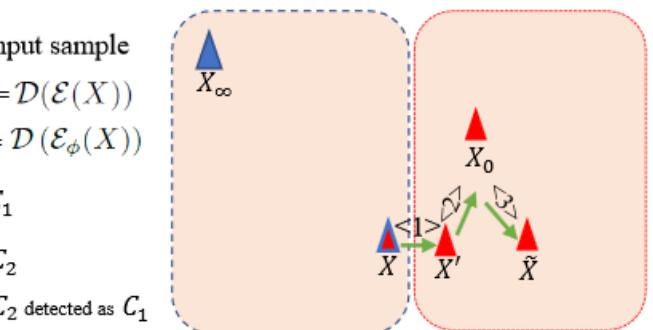
Self-Supervised Representation Learning via Neighborhood-Relational Encoding
Sabokrou, Khalooei, Adeli et al. ICCV, 2019



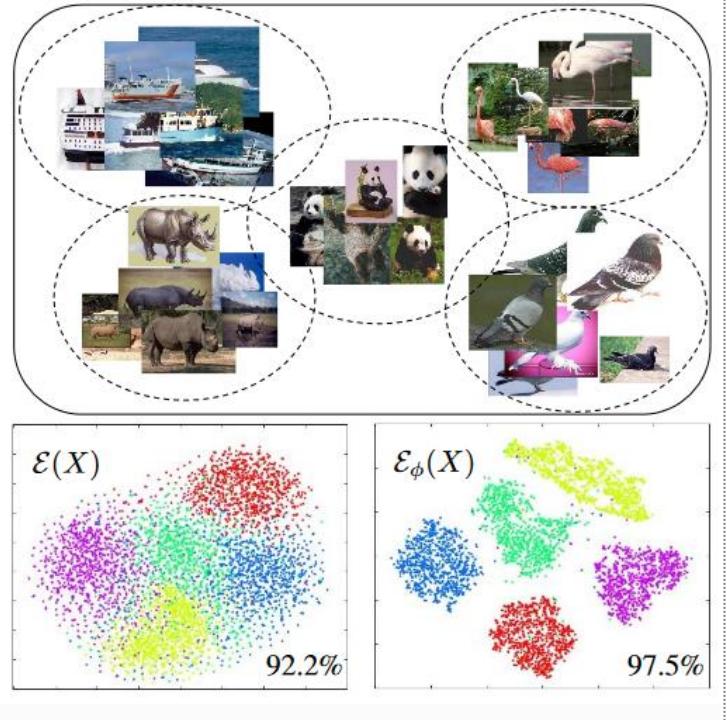
$$\mathbb{S}(\mathbf{a}, \mathbf{b}) = \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|} = \frac{\sum_{i=1}^n a_i b_i}{\sqrt{\sum_{i=1}^n a_i^2} \sqrt{\sum_{i=1}^n b_i^2}}$$

$$\mathbb{D}(\cdot, \cdot) = 1 - \mathbb{S}(\cdot, \cdot)$$

X Input sample
 $X' = \mathcal{D}(\mathcal{E}(X))$
 $\tilde{X} = \mathcal{D}(\mathcal{E}_\phi(X))$
▲ C_1
▲ C_2
▲ C_2 detected as C_1



$$\begin{aligned} \mathcal{L} &= \lambda_1 \mathbb{D}(\mathcal{R}_{\mathcal{A}}(X), \mathcal{R}_{\mathcal{A}}(X')) \\ &+ \lambda_2 \sum_{i=1}^T \mathbb{D}(\mathcal{R}_{\mathcal{A}}(X'), \mathcal{R}_{\mathcal{A}}(X_{0i})) \\ &+ \lambda_3 \sum_{i=1}^T \mathbb{S}(\mathcal{R}_{\mathcal{A}}(X'), \mathcal{R}_{\mathcal{A}}(X_{\infty i})) \end{aligned}$$



Break



Framework

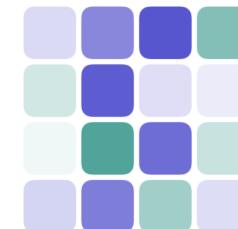


Adversarial
Robustness
Toolbox

TensorFlow, Keras, PyTorch,
MXNet, scikit-learn, XGBoost,
LightGBM, CatBoost, GPy, etc.



JAX, PyTorch, and TF2



Foolbox

PyTorch, TensorFlow, and
JAX





Break



Loss landscape

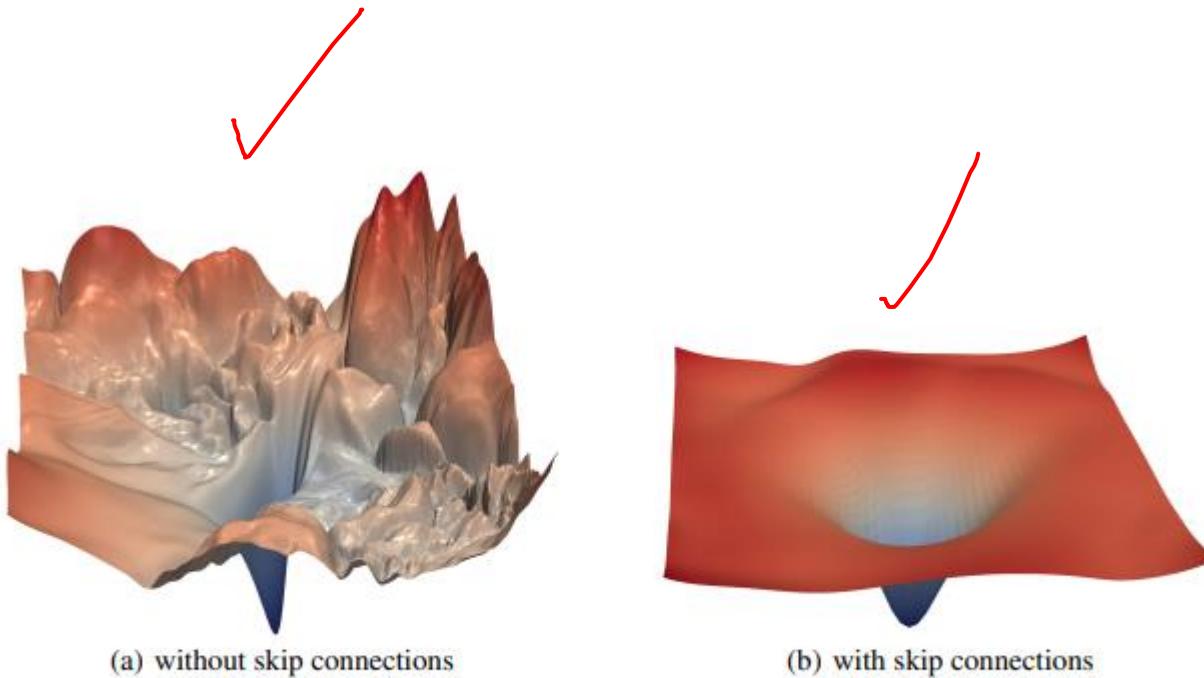
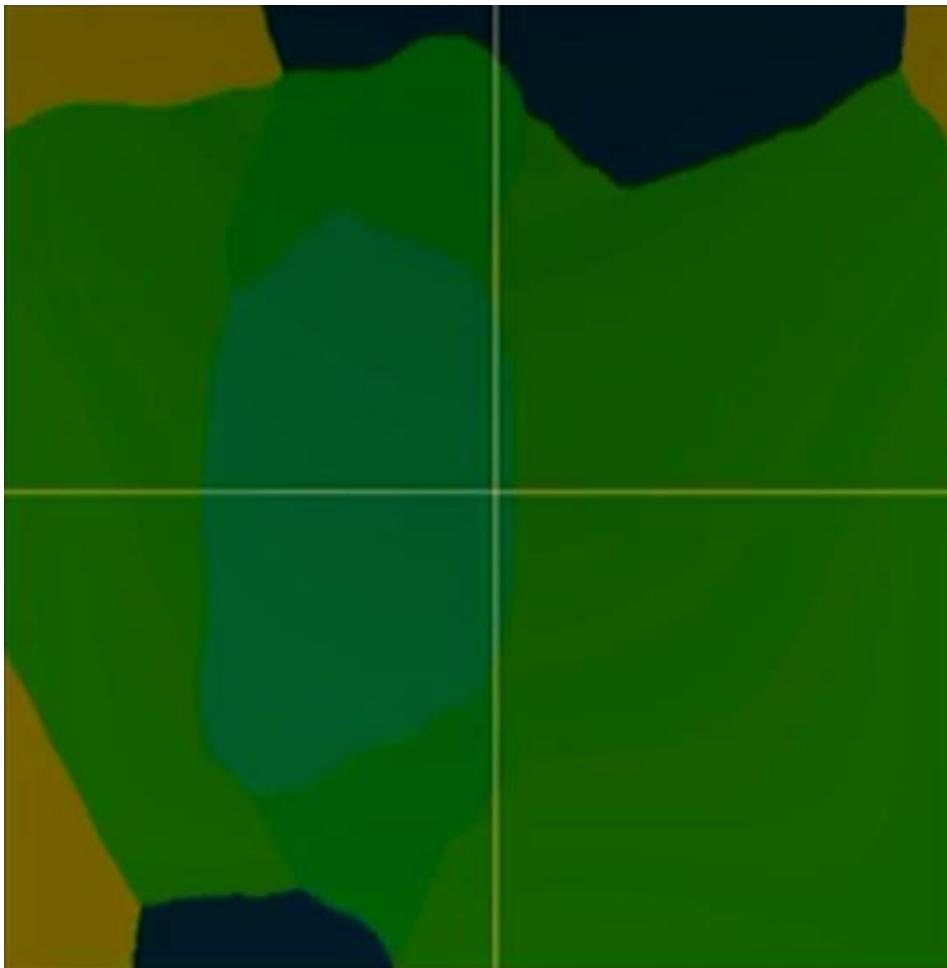
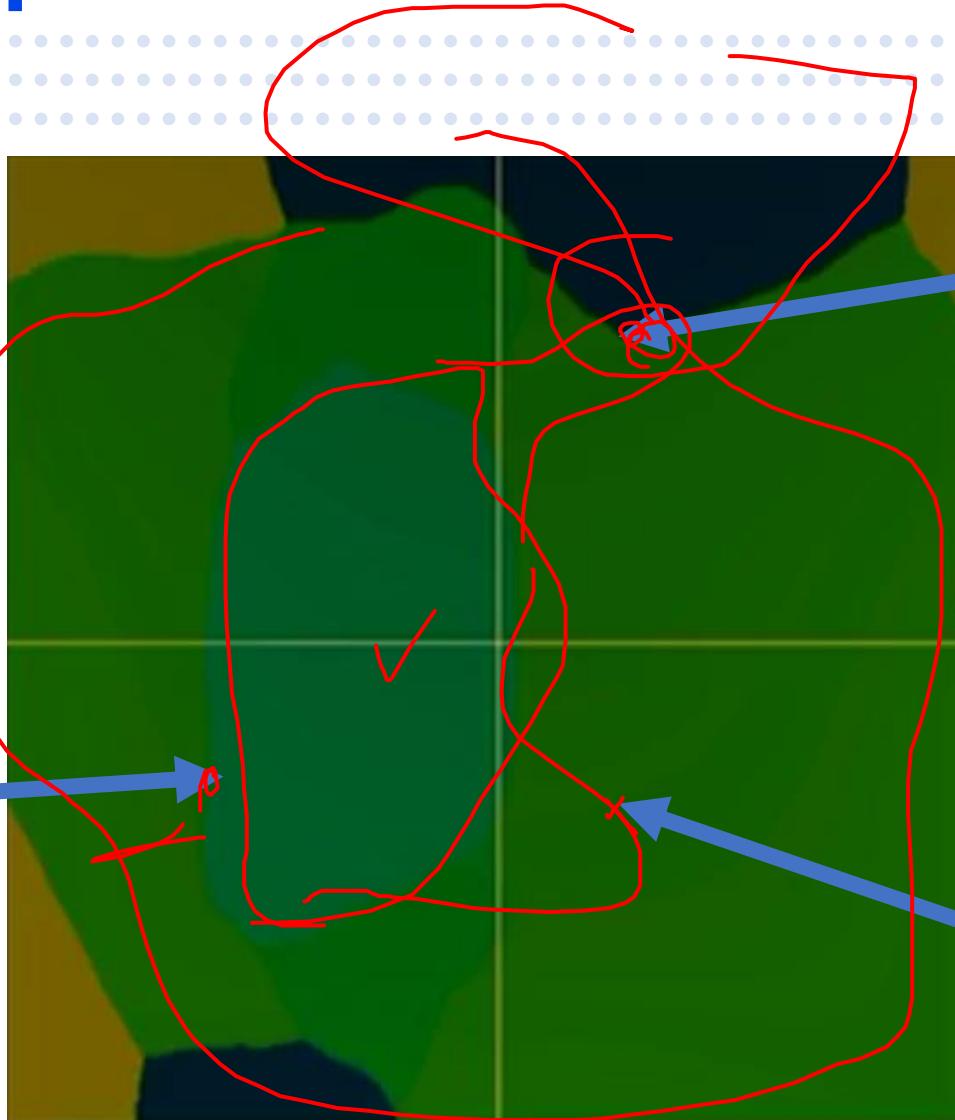


Figure 1: The loss surfaces of ResNet-56 with/without skip connections. The proposed filter normalization scheme is used to enable comparisons of sharpness/flatness between the two figures.

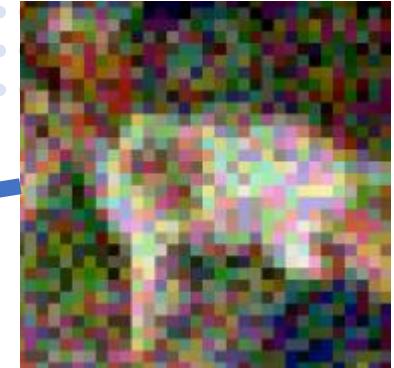
Loss landscape



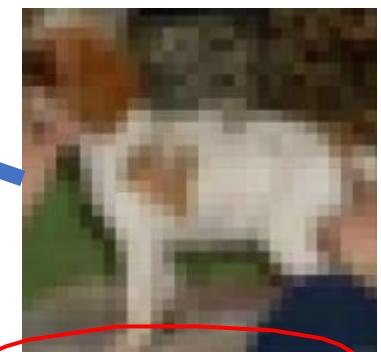
Loss landscape



Dog

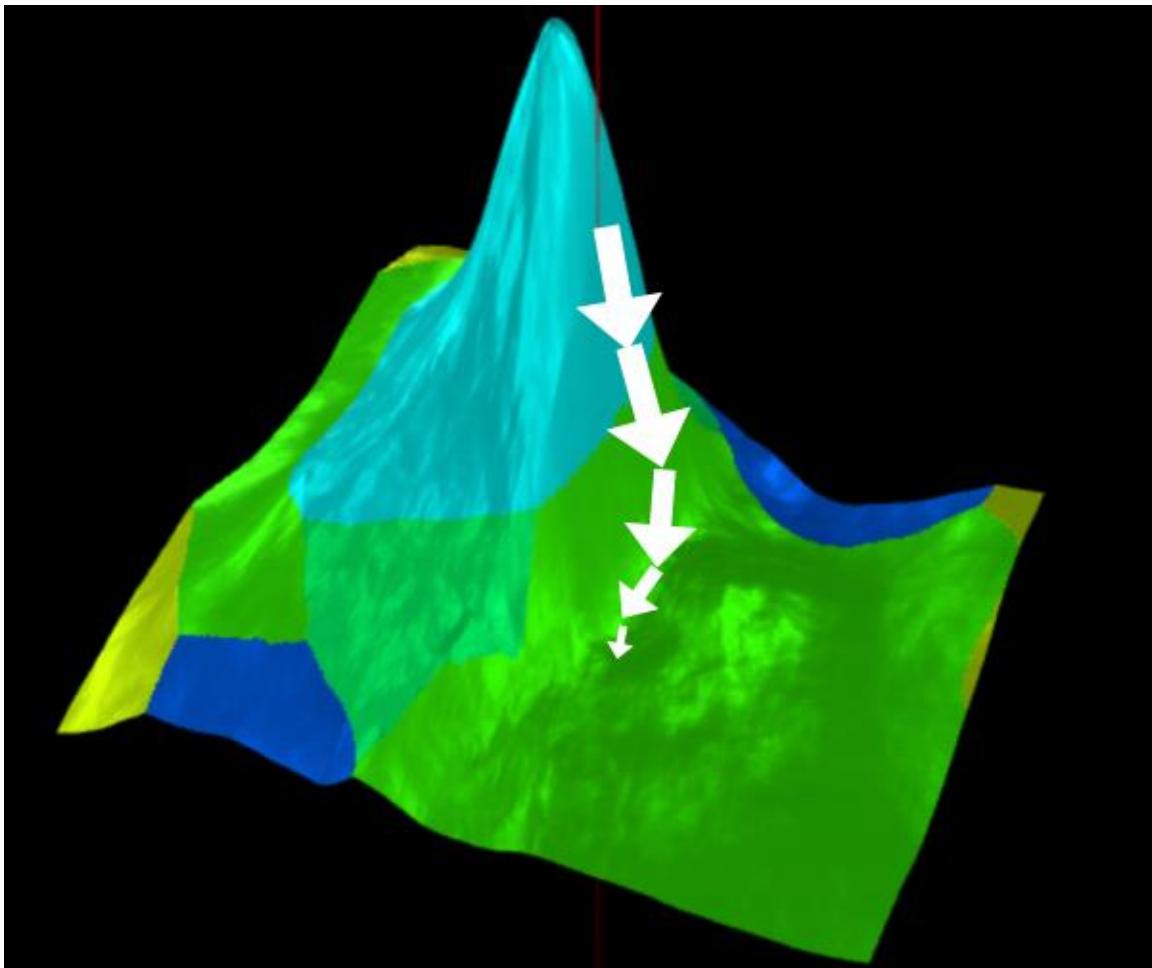


Truck

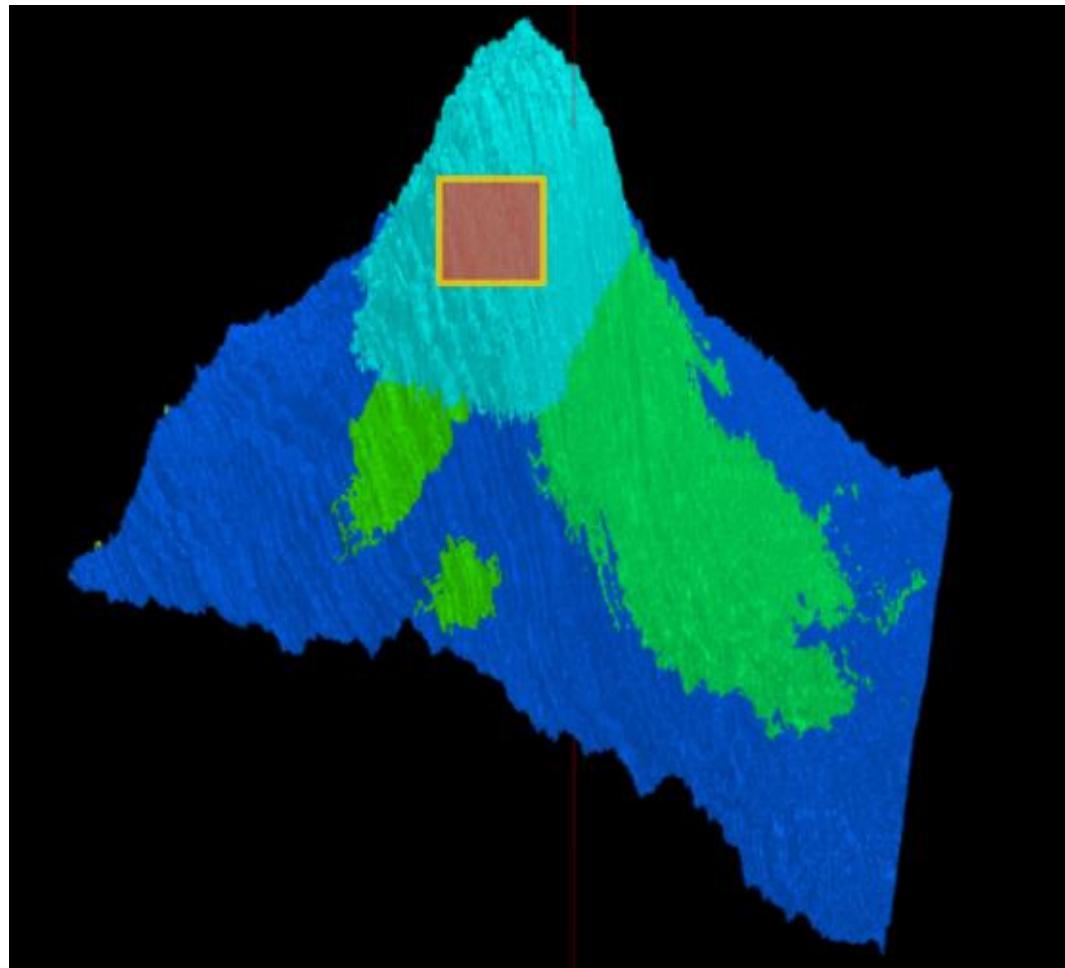


Airplane

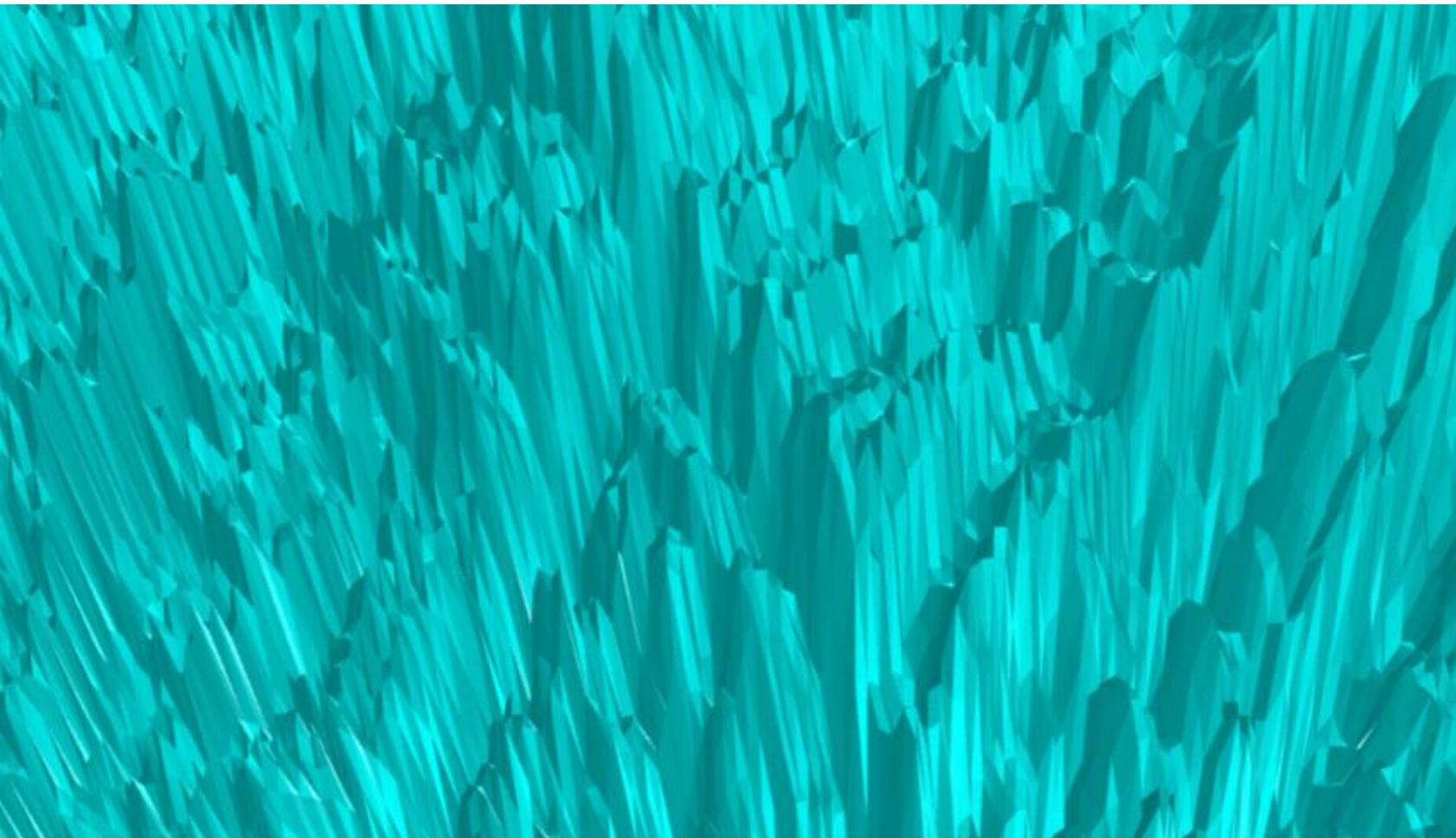
Loss landscape



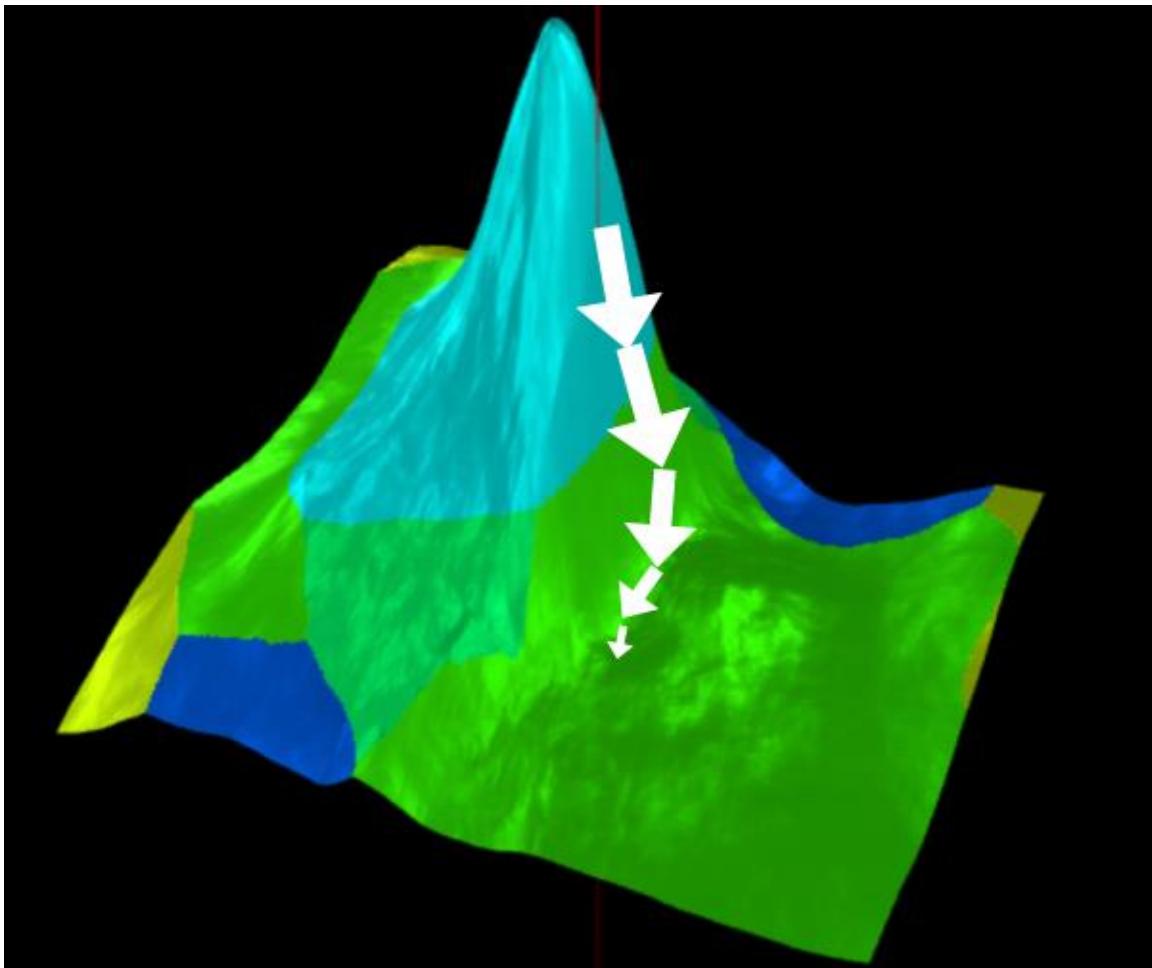
Loss landscape



Loss landscape



Loss landscape



Defense

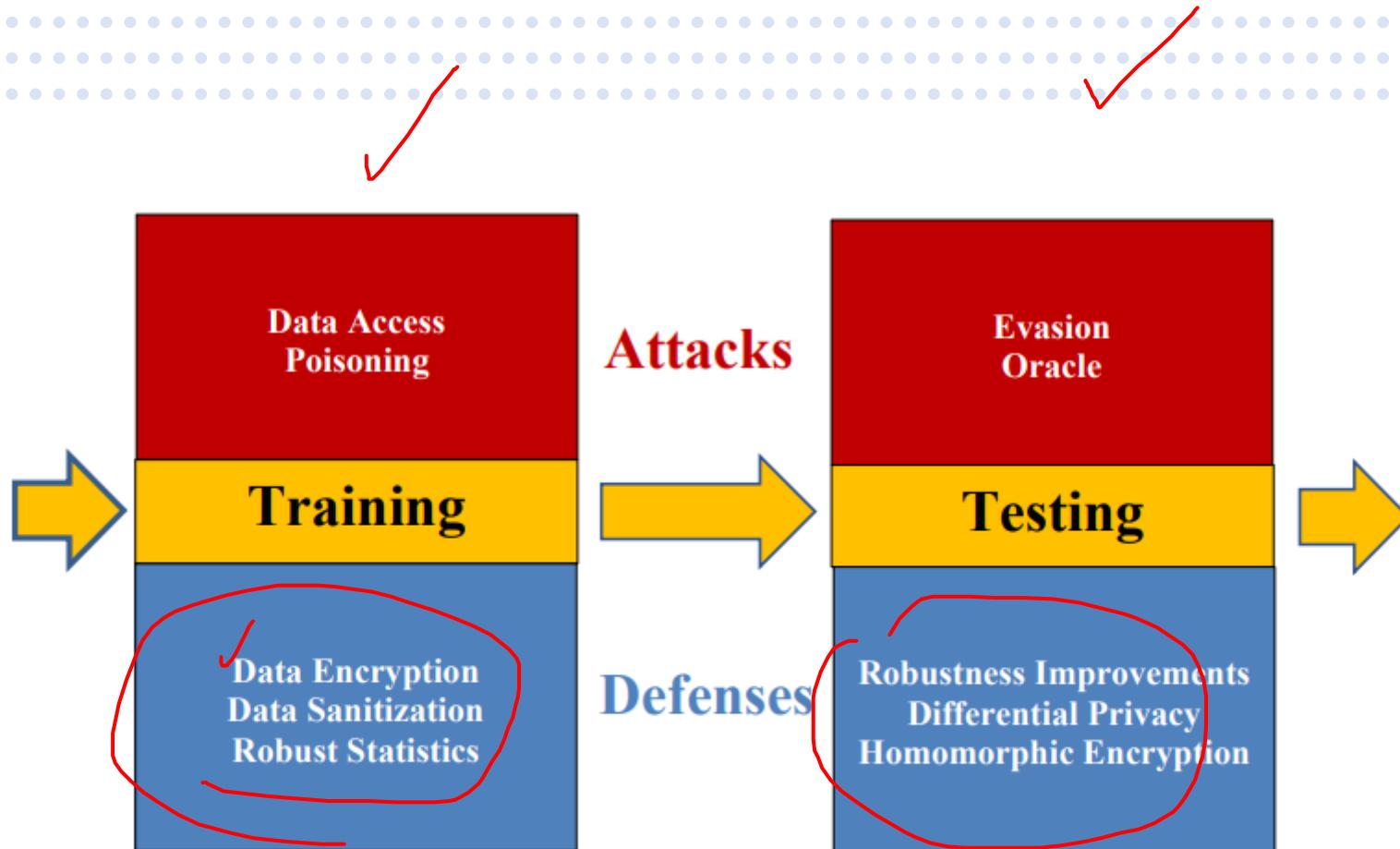


Figure 1. An illustration of example Attacks and Defenses in the Machine Learning Pipeline.



Break





Adversarial training methodologies & defenses

Advances in defense strategies

Defense is hard!

~~A theoretical model of the adversarial example crafting process is very difficult to construct.~~

- Non-linearity
- non-convex
- Complex optimization process
- ...

Most of the current defense strategies are

- not adaptive to **all types** of adversarial attack

Implementation of such defense strategies

- may incur **performance overhead**

Advances in defense strategies

- Defense
 - Data manipulation ✓
 - Architecture manipulation ✓
 - Loss manipulation ✓
 - Certified robustness ✓
 - From interpretability and explainability
 - ...

Empirical risk minimization

In standard ERM, we optimize the following objective:

$$\min_{\theta} \mathbb{E}_{(\mathbf{x},y)} [\ell_{cls} (f_{\theta}(\mathbf{x}), y)]$$

Empirical risk minimization & Adversarial training

In standard ERM, we optimize the following objective:

$$\min_{\theta} \mathbb{E}_{(\mathbf{x},y)} [\ell_{cls}(f_{\theta}(\mathbf{x}), y)]$$

For adversarial training, we need to optimize two things simultaneously.

- **First**, we generate the strongest minimal perturbation first.
- **Second**, we train our models to be robust against that.

Empirical risk minimization & Adversarial training

In standard ERM, we optimize the following objective:

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y)} [\ell_{cls}(f_{\theta}(\mathbf{x}), y)]$$

For adversarial training, we need to optimize two things simultaneously.

- First, we generate the strongest minimal perturbation first.
 - Second, we train our models to be robust against that.

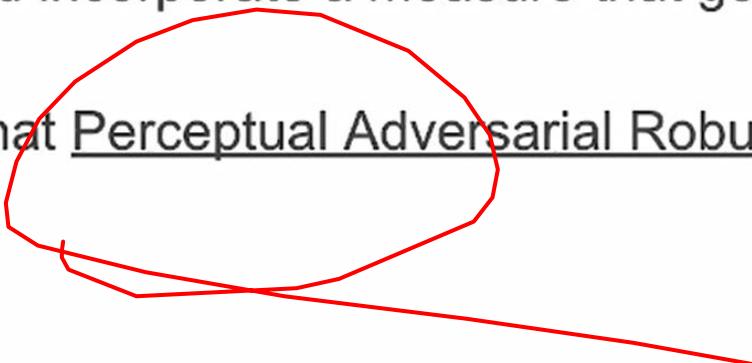
Mathematically (Madry et al., ICLR'18) -

$$\min_{\theta} \mathbb{E}_{(\mathbf{x}, y)} \left[\max_{\delta} \ell_{cls} (f_{\theta}(\mathbf{x} + \delta), y) \right]$$

SGD PGD

Adv. training

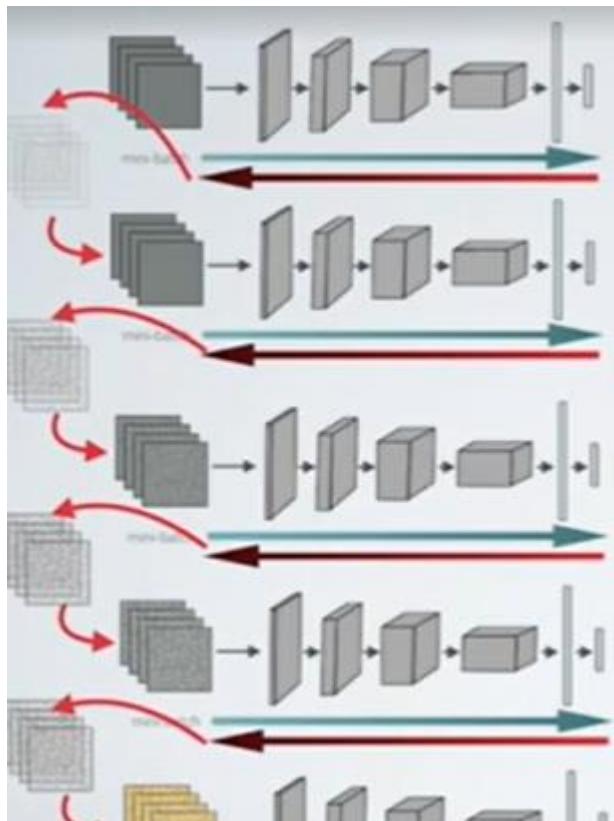
- Since human perception is hard to characterize precisely, this lack of transfer become inevitable.
- So, what if we could incorporate a measure that gets us closer to the human perception?
- This is precisely what Perceptual Adversarial Robustness (Laidlaw et al., ICLR'21) does.



Adv. training FREE

Adversarial Training for Free!

NeurIPS 2019



Algorithm 1 “Free” Adversarial Training (Free- m)

Require: Training samples X , perturbation bound ϵ , learning rate τ , hop steps m

```
1: Initialize  $\theta$ 
2:  $\delta \leftarrow 0$ 
3: for epoch = 1 ...  $N_{ep}/m$  do
4:   for minibatch  $B \subset X$  do
5:     for i = 1 ...  $m$  do
6:       Update  $\theta$  with stochastic gradient descent
7:        $g_\theta \leftarrow \mathbb{E}_{(x,y) \in B} [\nabla_\theta l(x + \delta, y, \theta)]$ 
8:        $g_{adv} \leftarrow \nabla_x l(x + \delta, y, \theta)$ 
9:        $\theta \leftarrow \theta - \tau g_\theta$ 
10:      Use gradients calculated for the minimization step to update  $\delta$ 
11:       $\delta \leftarrow \delta + \epsilon \cdot \text{sign}(g_{adv})$ 
12:       $\delta \leftarrow \text{clip}(\delta, -\epsilon, \epsilon)$ 
13:    end for
14:  end for
15: end for
```

Adv. training FAST

FAST IS BETTER THAN FREE:
REVISITING ADVERSARIAL TRAINING

ICLR 2020

Algorithm 3 FGSM adversarial training for T epochs, given some radius ϵ , N PGD steps, step size α , and a dataset of size M for a network f_θ

```
for  $t = 1 \dots T$  do
    for  $i = 1 \dots M$  do
        // Perform FGSM adversarial attack
        ✓  $\delta = \text{Uniform}(-\epsilon, \epsilon)$ 
         $\delta = \delta + \alpha \cdot \text{sign}(\nabla_\delta \ell(f_\theta(x_i + \delta), y_i))$ 
         $\delta = \max(\min(\delta, \epsilon), -\epsilon)$ 
         $\theta = \theta - \nabla_\theta \ell(f_\theta(x_i + \delta), y_i)$  // Update model weights with some optimizer, e.g. SGD
    end for
end for
```

Adv. training TRADE ✓

Theoretically Principled Trade-off between Robustness and Accuracy

ICLR 2019

Robust (classification) error

under the threat model of bounded ϵ perturbation

$$\mathcal{R}_{\text{rob}}(f) := \mathbb{E}_{(\mathbf{X}, Y) \sim \mathcal{D}} \mathbf{1}\{\exists \mathbf{X}' \in \mathbb{B}(\mathbf{X}, \epsilon) \text{ s.t. } f(\mathbf{X}')Y \leq 0\}$$

Natural (classification) error

standard measure of classifier performance

$$\mathcal{R}_{\text{nat}}(f) := \mathbb{E}_{(\mathbf{X}, Y) \sim \mathcal{D}} \mathbf{1}\{f(\mathbf{X})Y \leq 0\}$$

$$x_1, \dots, x_n \in \mathcal{X}$$

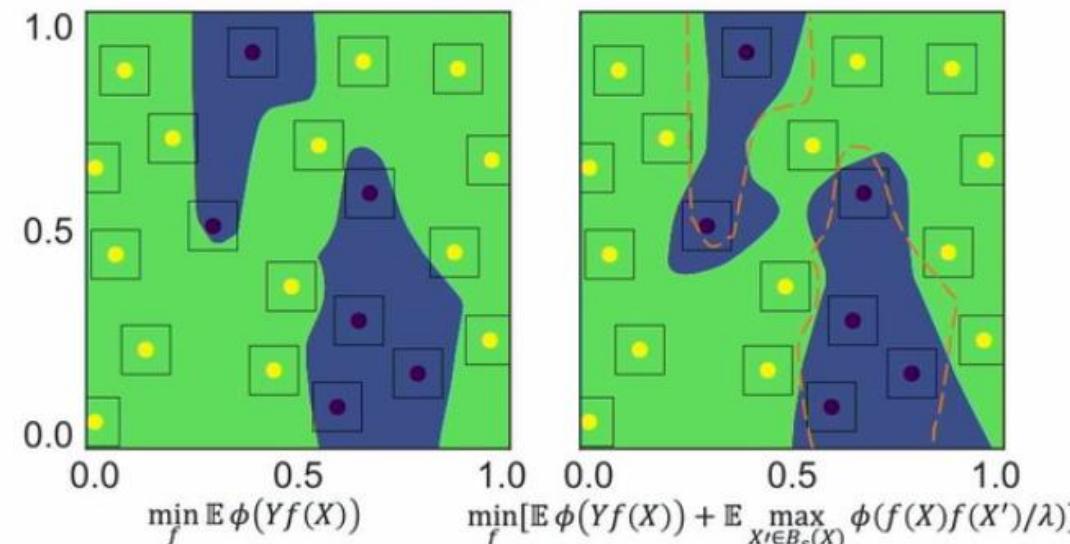
$$y_1, \dots, y_n \in \{-1, +1\}$$

Boundary error

$$\mathcal{R}_{\text{bdy}}(f) := \mathbb{E}_{(\mathbf{X}, Y) \sim \mathcal{D}} \mathbf{1}\{\mathbf{X} \in \mathbb{B}(\text{DB}(f), \epsilon), f(\mathbf{X})Y > 0\}$$

$$\mathcal{R}_{\text{rob}}(f) = \mathcal{R}_{\text{nat}}(f) + \mathcal{R}_{\text{bdy}}(f).$$

$$\min_f \mathbb{E} \left\{ \underbrace{\phi(f(\mathbf{X})Y)}_{\text{for accuracy}} + \underbrace{\max_{\mathbf{X}' \in \mathbb{B}(\mathbf{X}, \epsilon)} \phi(f(\mathbf{X})f(\mathbf{X}')/\lambda)}_{\text{regularization for robustness}} \right\}.$$







Break





Promising recipes

Studying optimizer susceptibility

From the previous plot, optimizers that **may** easily fall prey to the attacks:

- Adam
- RMSProp

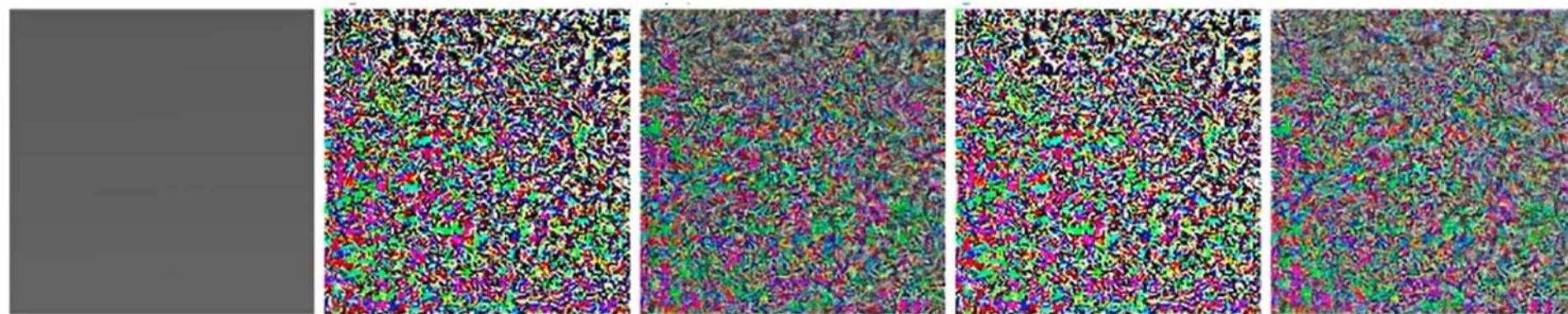
Optimizers that may **not** easily fall prey to the attacks:

- SGD
- Adagrad
- FTRL

This is characterized by the non-convexity of the optimization problem.

Studying optimizer susceptibility

Different deltas (δ) as learned by the optimizers



FTRL

Adagrad

RMSProp

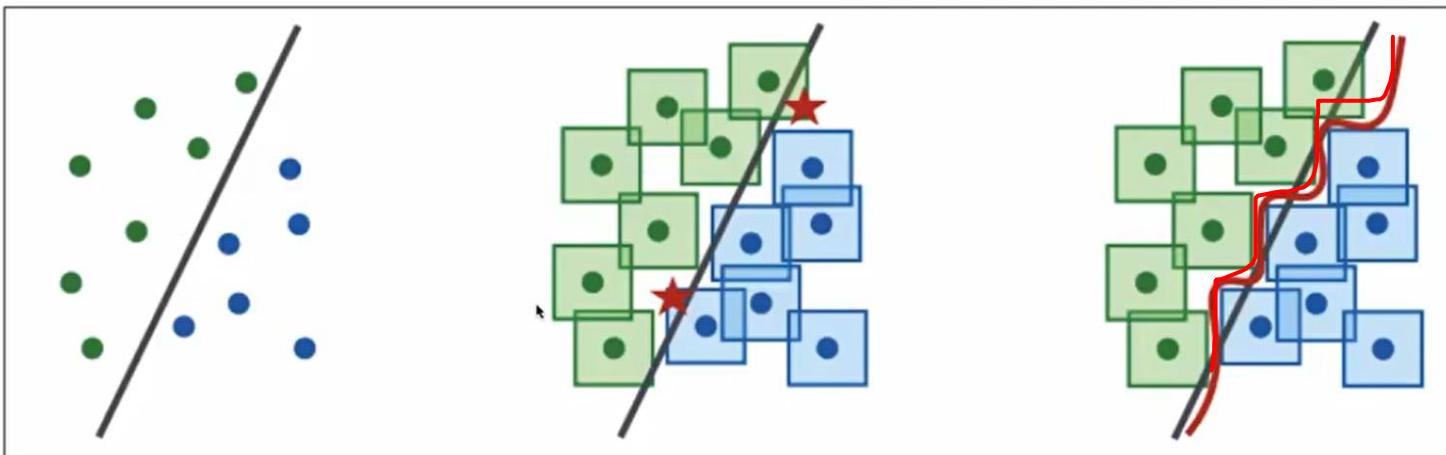
SGD

Adam



Model capacity is crucial

Adversarial examples change the decision boundary to a more complicated one
(Madry et al., ICLR'18).

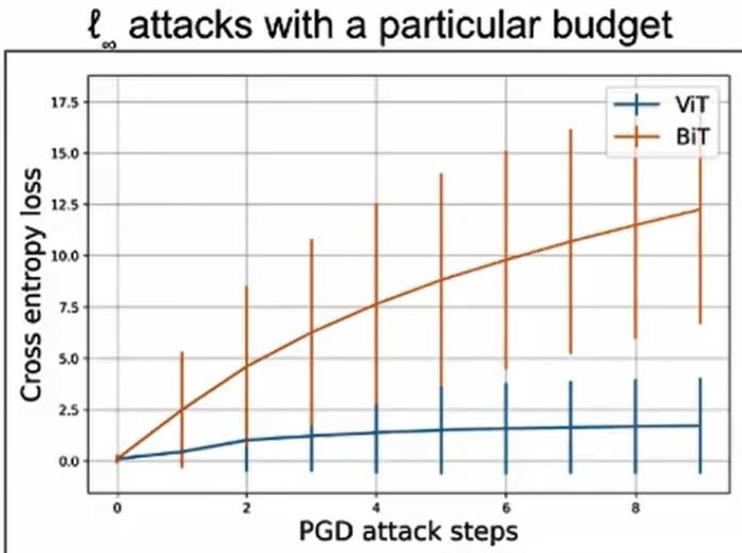
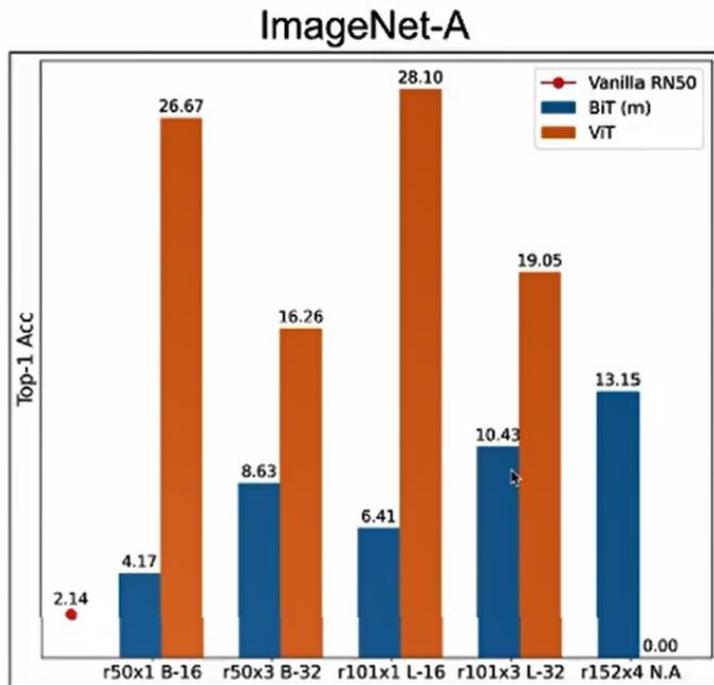


Madry et al., ICLR'18

■ \square ℓ_∞ -balls
★ adv. examples

Self-attention provides improved robustness

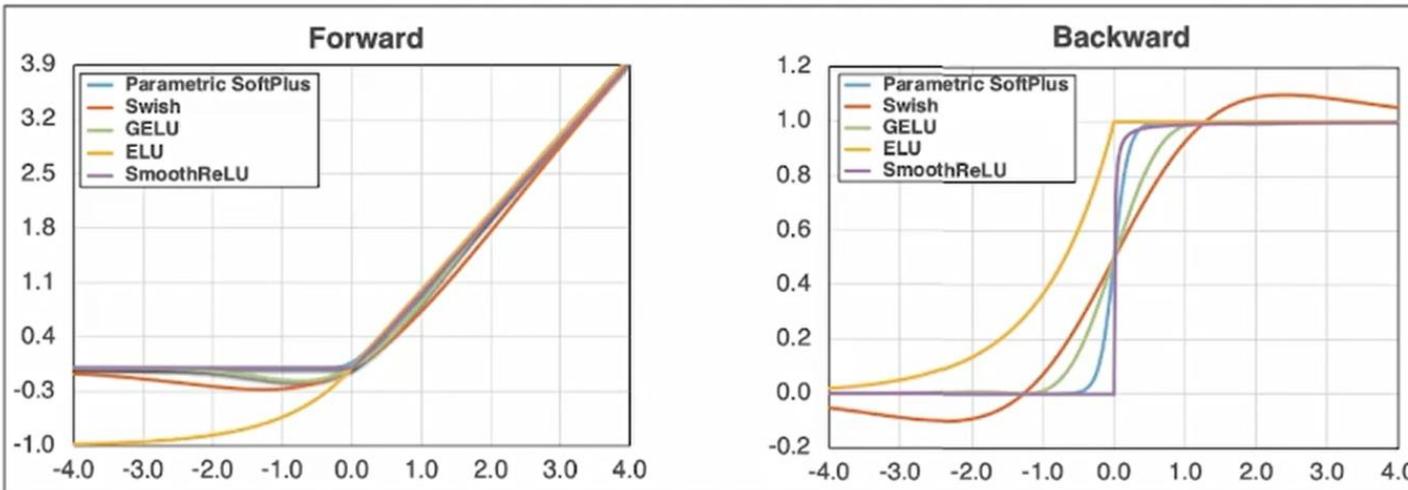
Continuation of the previous discussion -



Paul et al., arXiv, 2021

Smooth adv. training

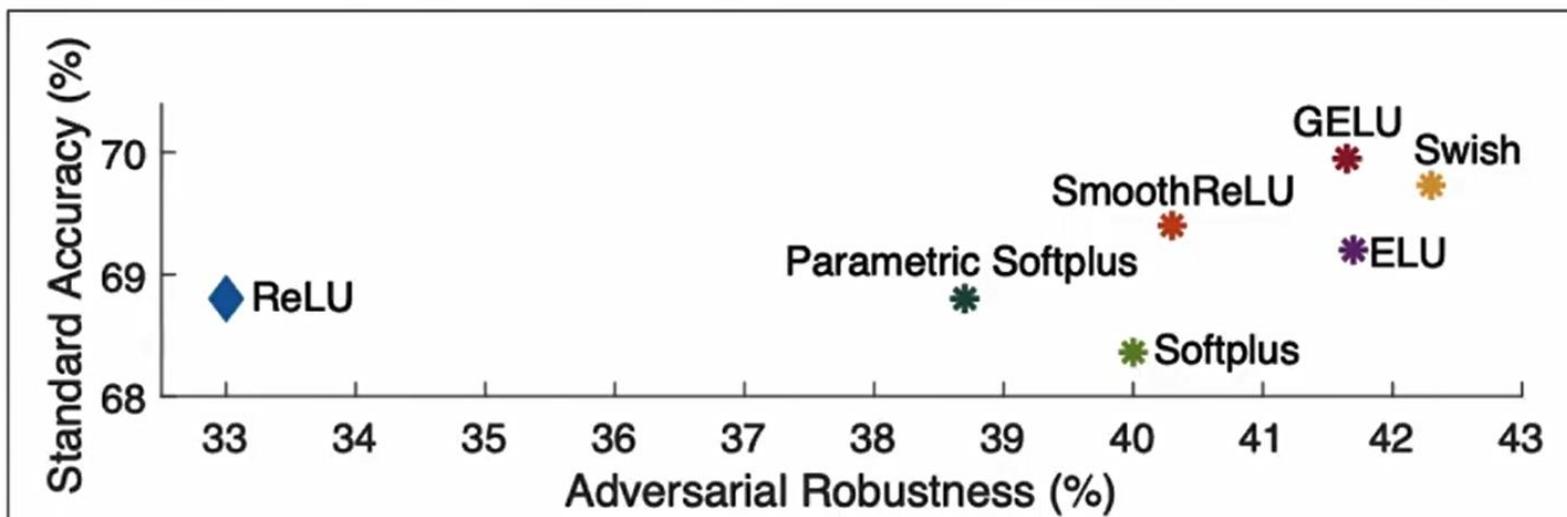
- Using ReLU during adv. training is particularly worse off because of its non-smooth nature.
- Smoothened activation functions (Swish, SoftPlus, etc.) result into better informed gradients because of their smoothness.



Xie et al., arXiv, 2020

Smooth adv. training

The use of smoother activation functions leads to improved performance without accuracy loss.



Xie et al., arXiv, 2020

Conclusion

Being aware is helpful

- For model developers, adversarial examples can be used for robustness evaluation and model improvement.
- For business stakeholders, lacking adversarial robustness in your AI model could bring unexpected negative impacts.
- For end users, gaining awareness of adversarial robustness for the AI service you are using is crucial.

Takeaways

- Emerging trend of transition from stand ML to (adversarially) robust ML.
- Inspect potential errors and risks of your ML algorithm/system using adversarial ML tools.
- Improve model robustness via adversarial machine learning.
- Anticipate adversarial inputs and data shifts.
- Novel applications driven by adversarial robustness.
- Practice adversarial robustness in ML lifecycle.

Online resource for advevrsarial robustness

- <https://adversarial-ml-tutorial.org/> ✓
- <https://nips.cc/Conferences/2018/ScheduleMultitrack?event=10978> ✓
- <https://sites.google.com/umich.edu/cvpr-2020-zoo> ✓
- <https://advmlincv.github.io/cvpr21-tutorial/> ✓
- <https://eccv20-adv-workshop.github.io/>
- <https://sites.google.com/view/par-2021> ✓
- <http://www.cs.umd.edu/class/fall2020/cmsc828W/> ✓

On evaluating adversarial robustness

ON EVALUATING ADVERSARIAL ROBUSTNESS

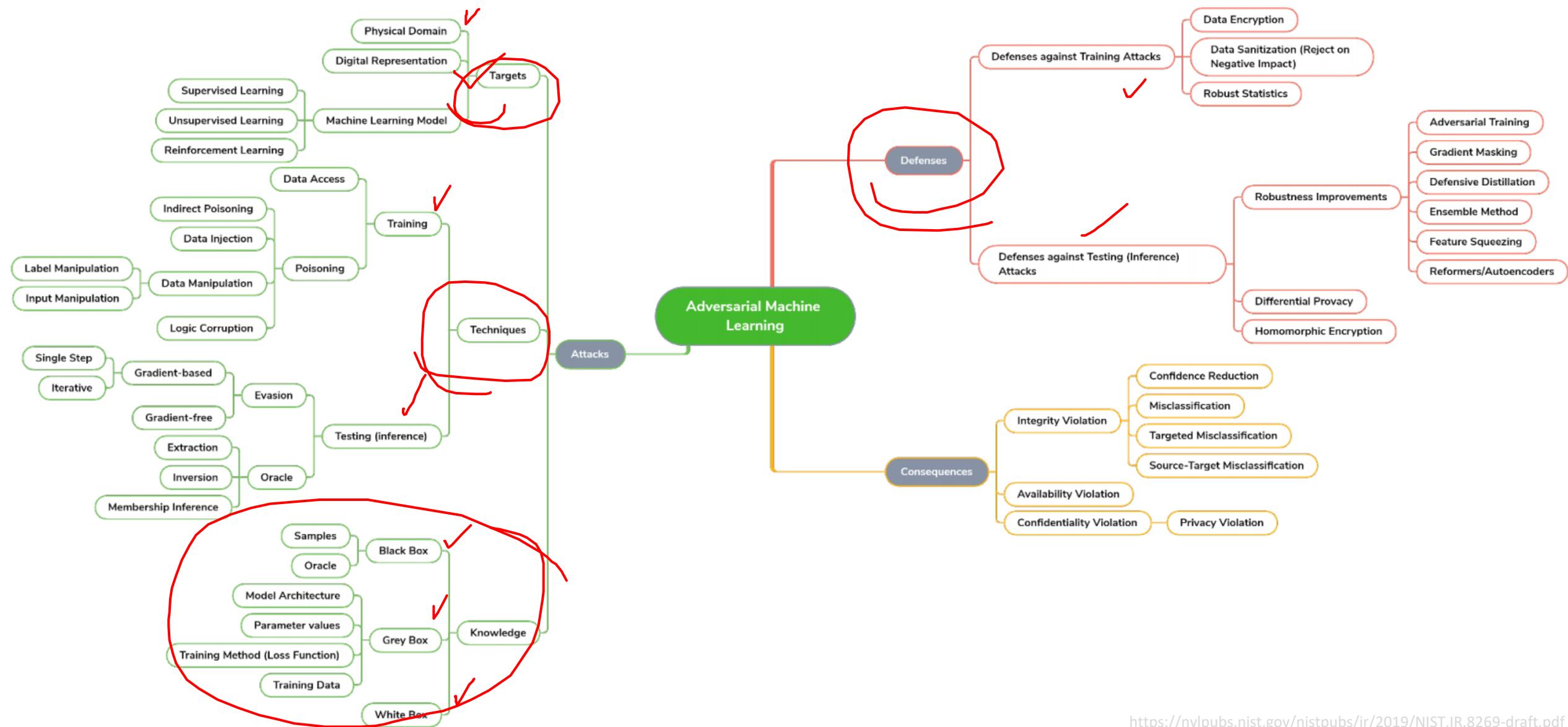
Nicholas Carlini¹, Anish Athalye², Nicolas Papernot¹, Wieland Brendel³, Jonas Rauber³,
Dimitris Tsipras², Ian Goodfellow¹, Aleksander Mądry², Alexey Kurakin^{1*}

¹ Google Brain ² MIT ³ University of Tübingen

* List of authors is dynamic and subject to change. Authors are ordered according
to the amount of their contribution to the text of the paper.

<https://github.com/evaluating-adversarial-robustness/adv-eval-paper>

Taxonomy of Attacks, Defenses, and Consequences in Adversarial Machine Learning



<https://nvlpubs.nist.gov/nistpubs/ir/2019/NIST.IR.8269-draft.pdf>

Adversarial Machine Learning community



<https://forms.gle/ANbQLsrpzSFoVtmG7>

References

- Practical Adversarial Robustness in Deep Learning: Problems and Solutions (CVPR 21 tutorial) ✓
- CMSC 828W: Foundations of Deep Learning (by Dr. Soheil Feizi - 2020) ✓
- Research papers which each related reference depicted in the bottom corner of the slide ✓



Mohammad Khalooei

Khalooei [at] aut.ac.ir

<https://khalooei.ir>

<https://ce.aut.ac.ir/~khalooei>