

Master Data Science – Advanced ML Topics

Session: Adversarial Machine Learning

Presented by:

Dr. Karim Tit & Mr. Mohamed Djilani

University of Luxembourg

Previous slides from:

Dr. Salah GHAMIZI,

Luxembourg Institute of Science and Technology



Hello!

I am Karim TIT

Postdoctoral researcher @ University of Luxembourg

Email for questions (slides 1-53): karim.tit@uni.lu



Hello!

I am Mohamed Djilani

PhD student @ University of Luxembourg

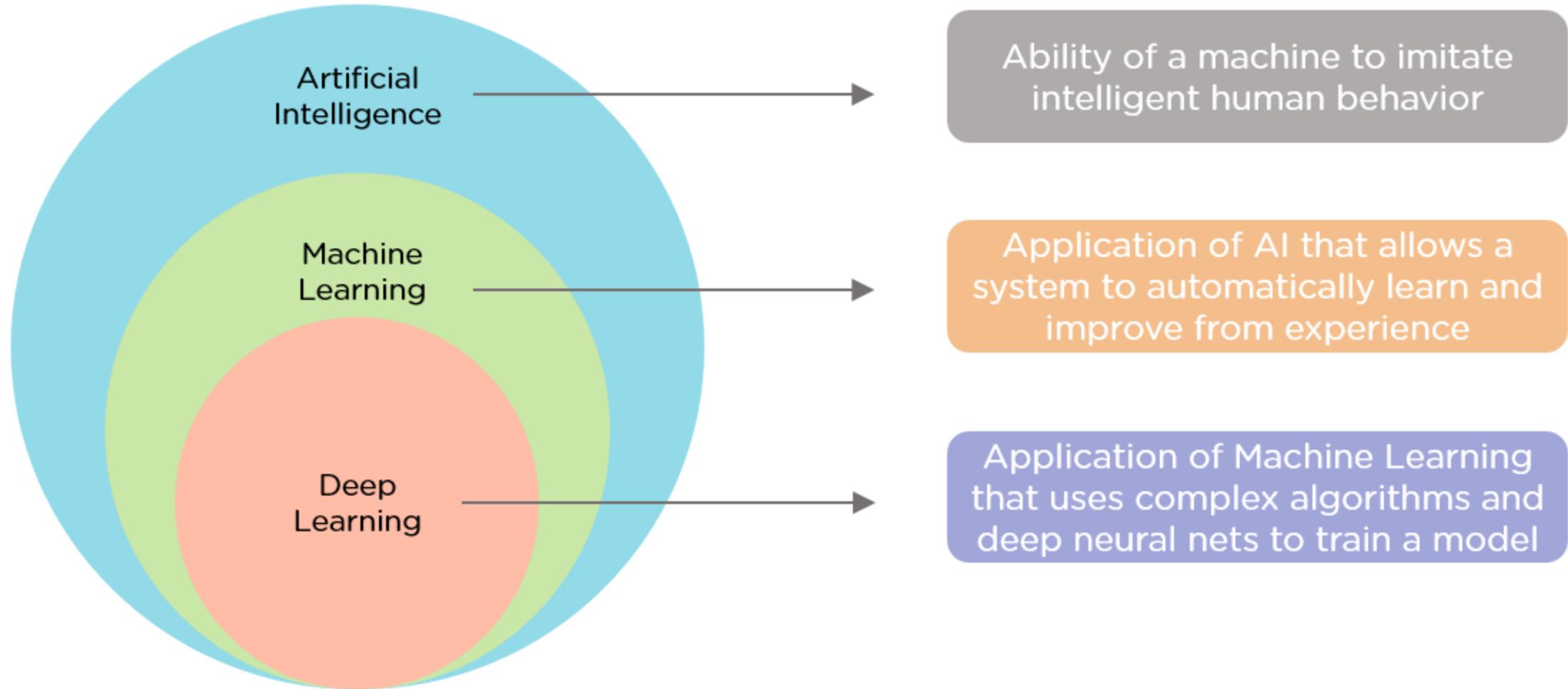
Email for questions (slides 54-100):

mohamed.djilani@uni.lu

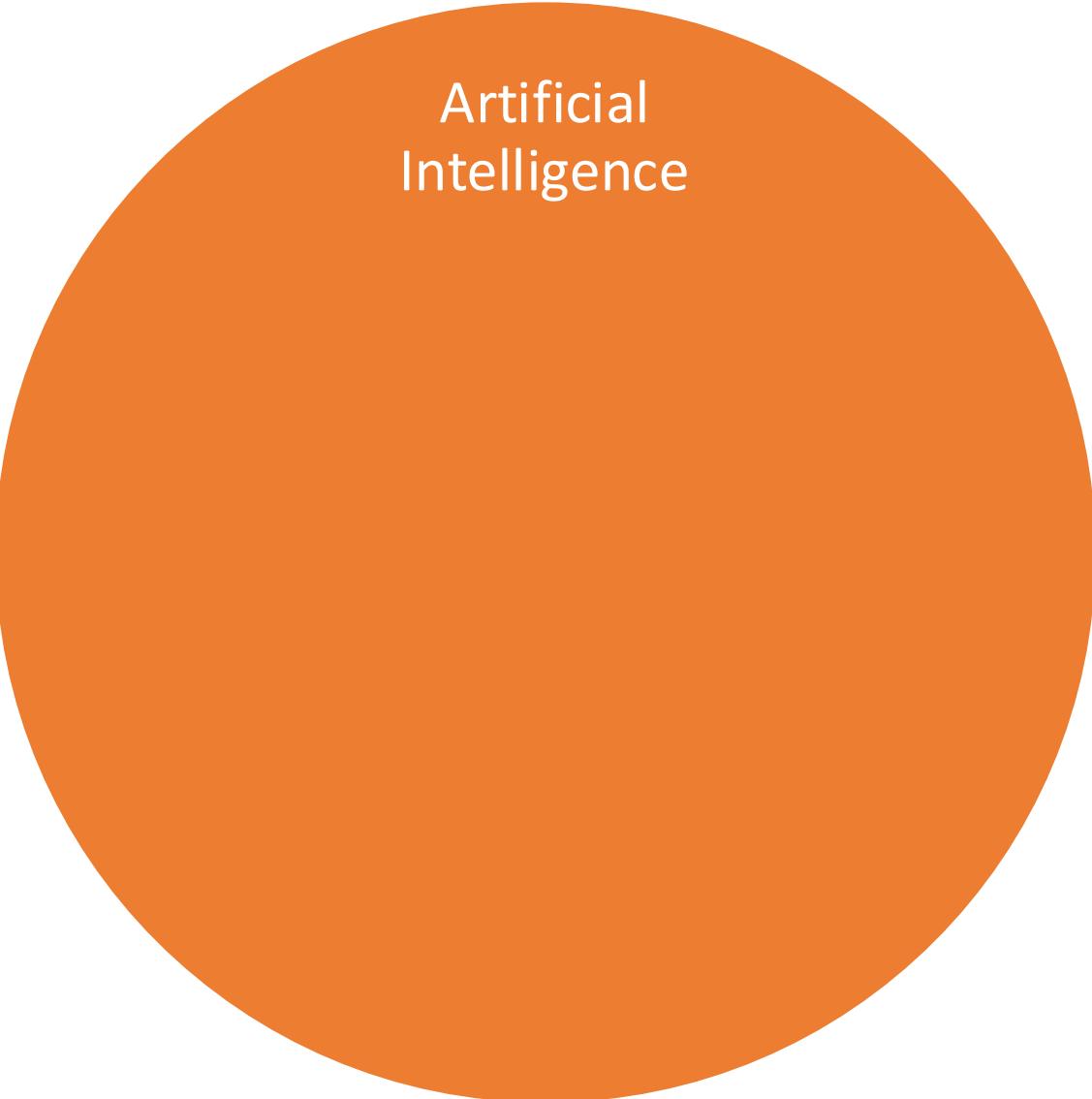
Machine Learning - Quick Recap

Part 0

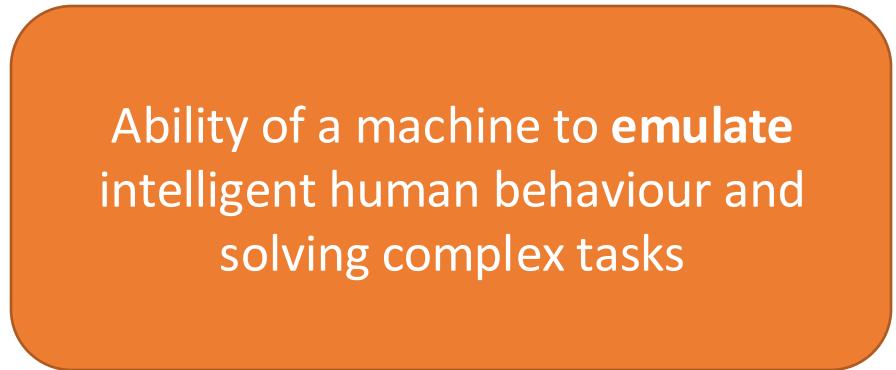
AI / ML / DL ? Kesako?



AI / ML / DL ? Kesako?

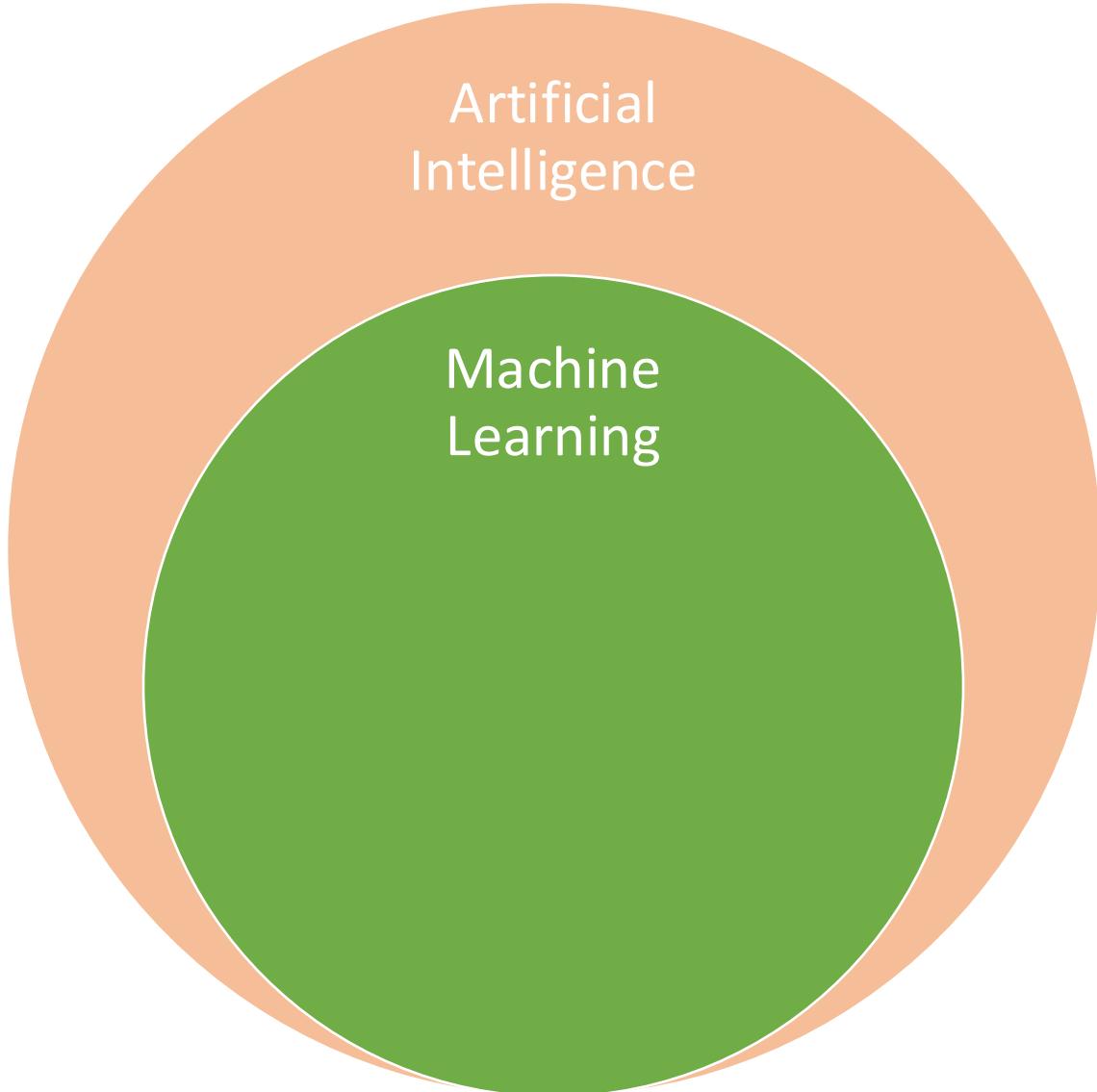


Artificial
Intelligence



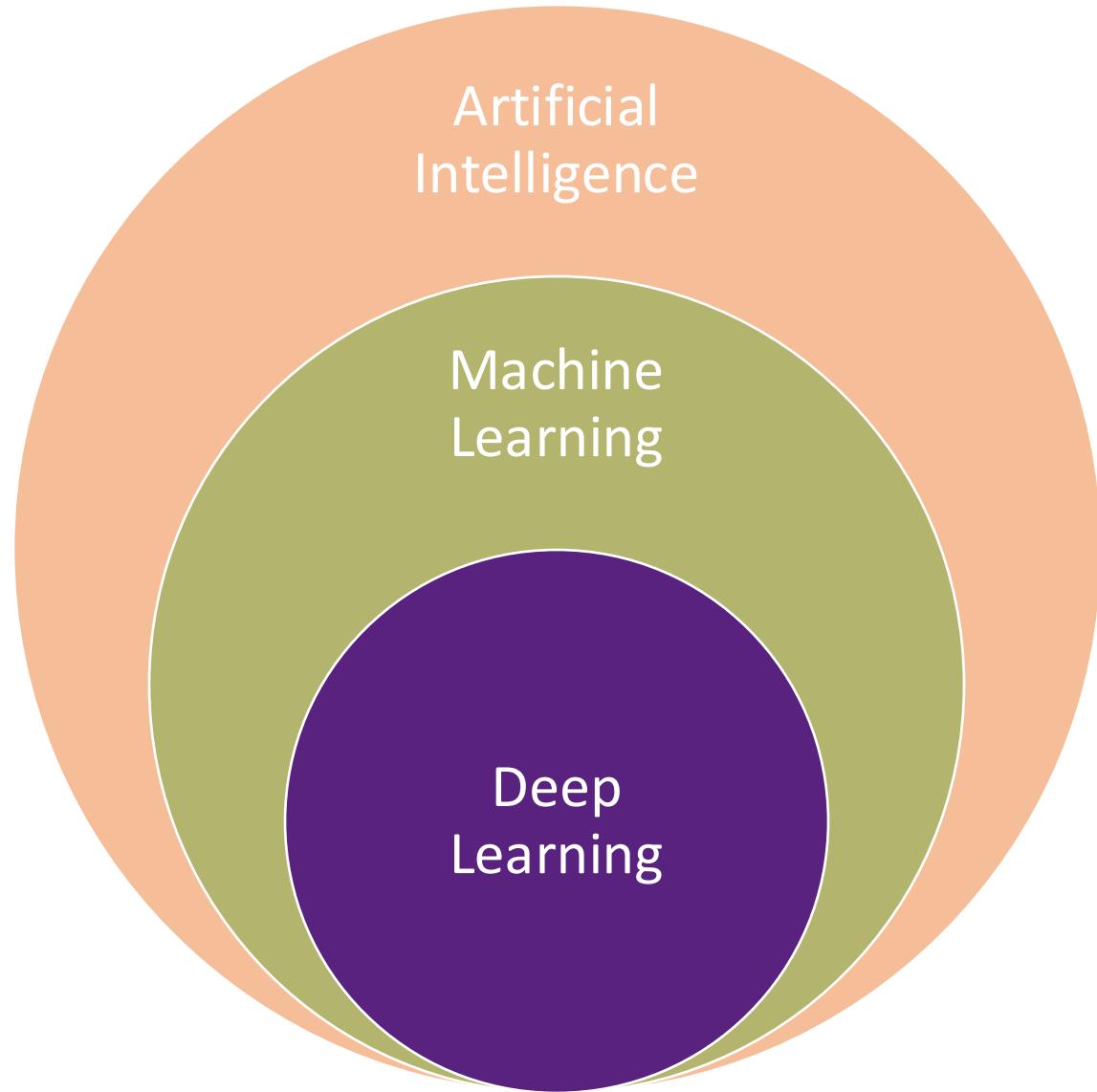
Ability of a machine to **emulate**
intelligent human behaviour and
solving complex tasks

AI / ML / DL ? Kesako?



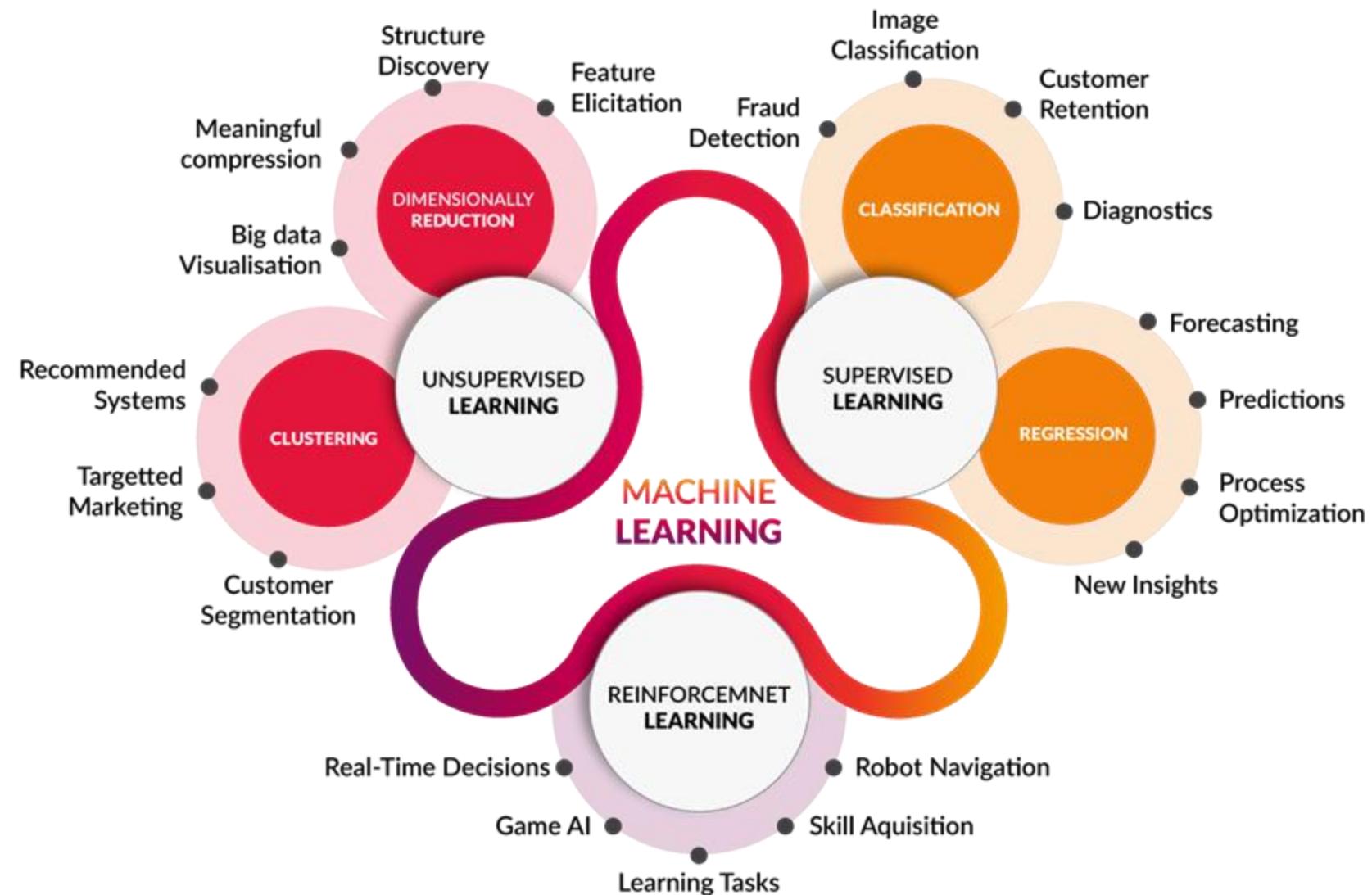
A branch of artificial intelligence where algorithms learn the hidden patterns of data to make predictions on new similar data, without being **explicitly** programmed for each task.

AI / ML / DL ? Kesako?

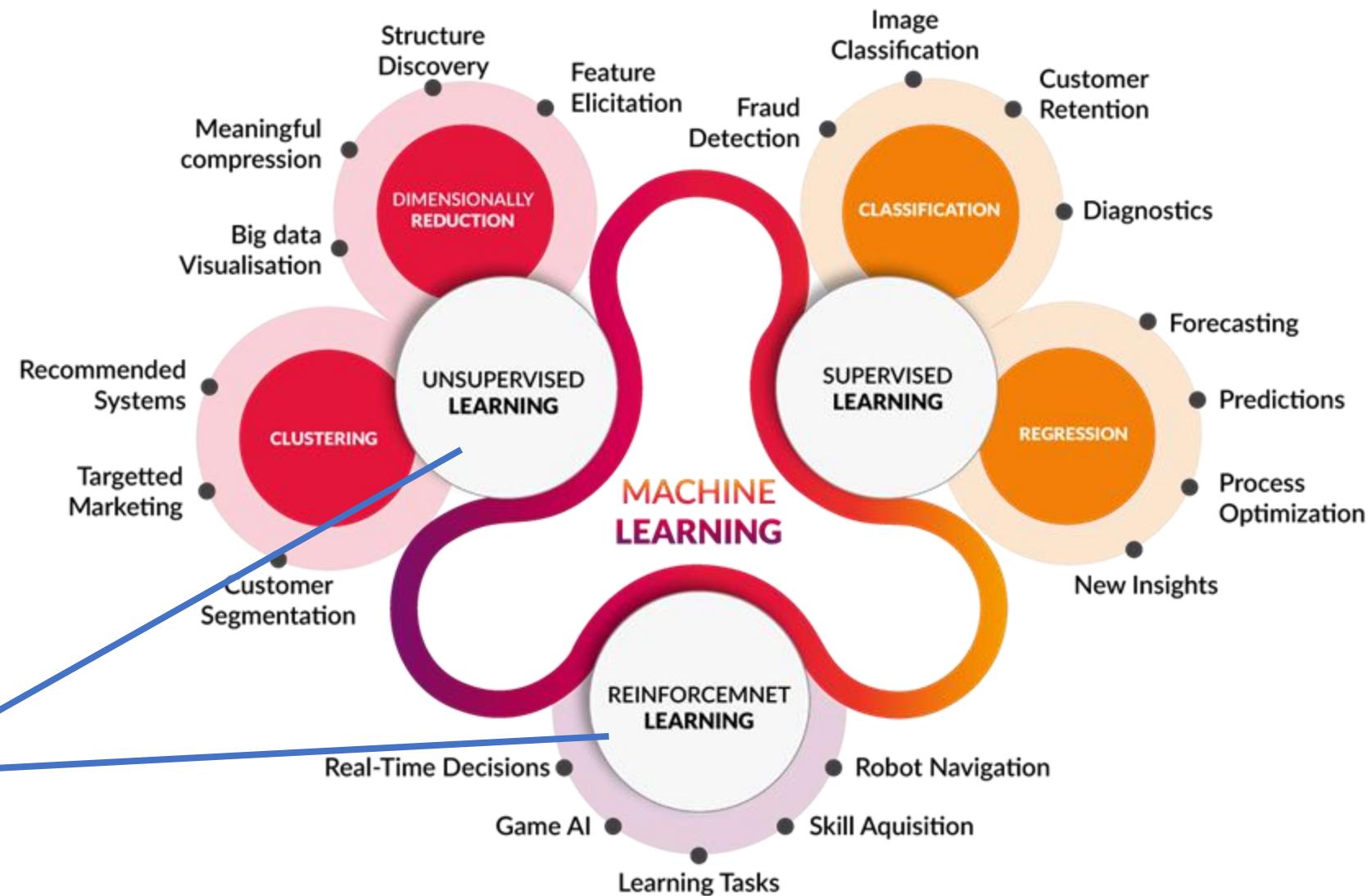


Set of Machine Learning techniques
that use intricate algorithms and **deep
neural networks** to train a model

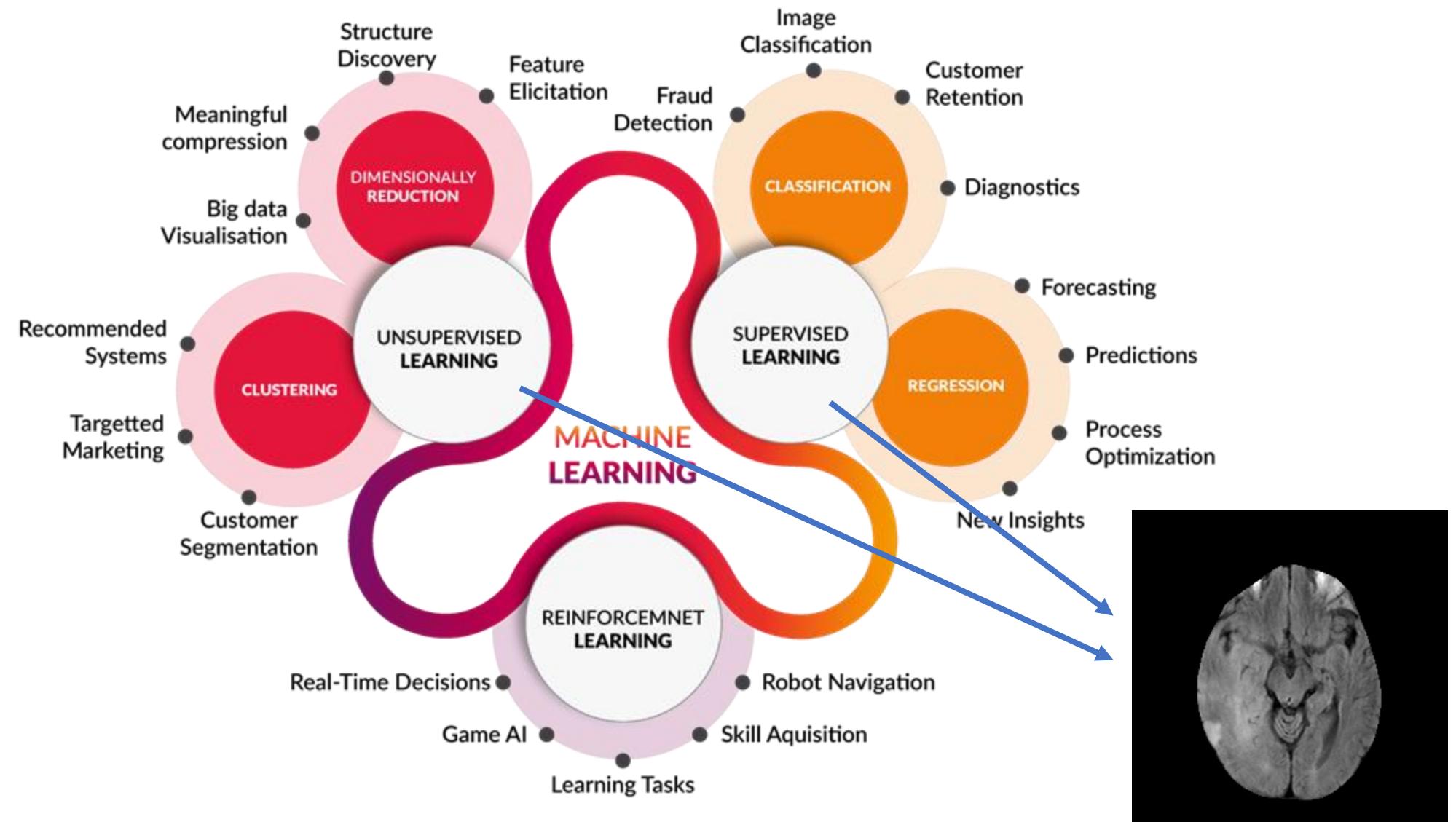
ML Taxonomy



ML Taxonomy

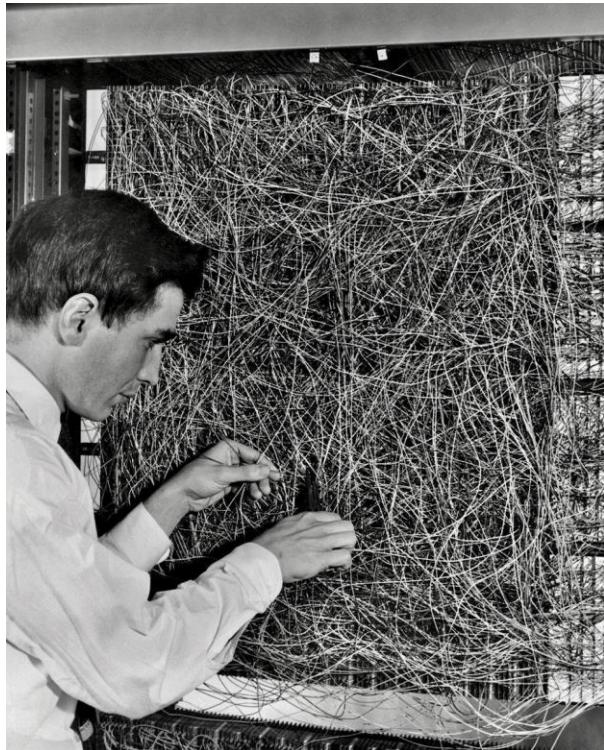


ML Taxonomy

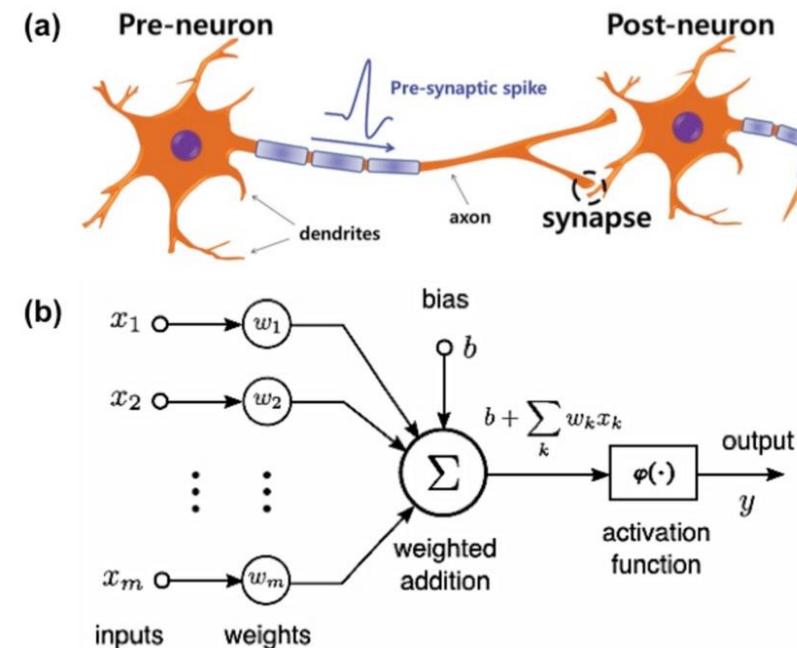


Artificial Neural Networks

Frank Rosenblatt: The Perceptron (1957)



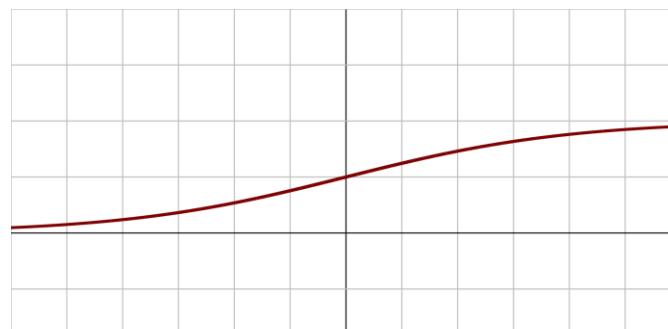
From biological to artificial neurons



Artificial Neural Networks

How to handle non-linear problems?

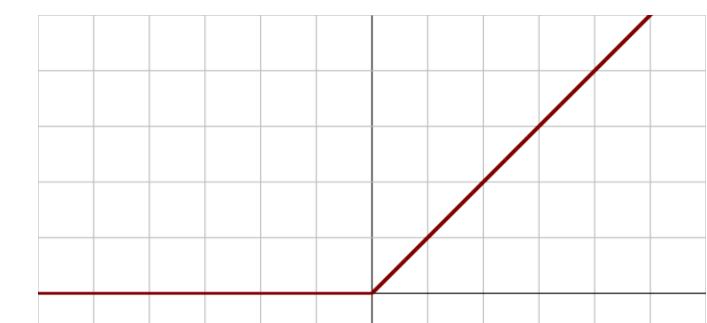
1/ Adding non-linear activation functions → Artificial Neural Networks



Sigmoid



Gaussian

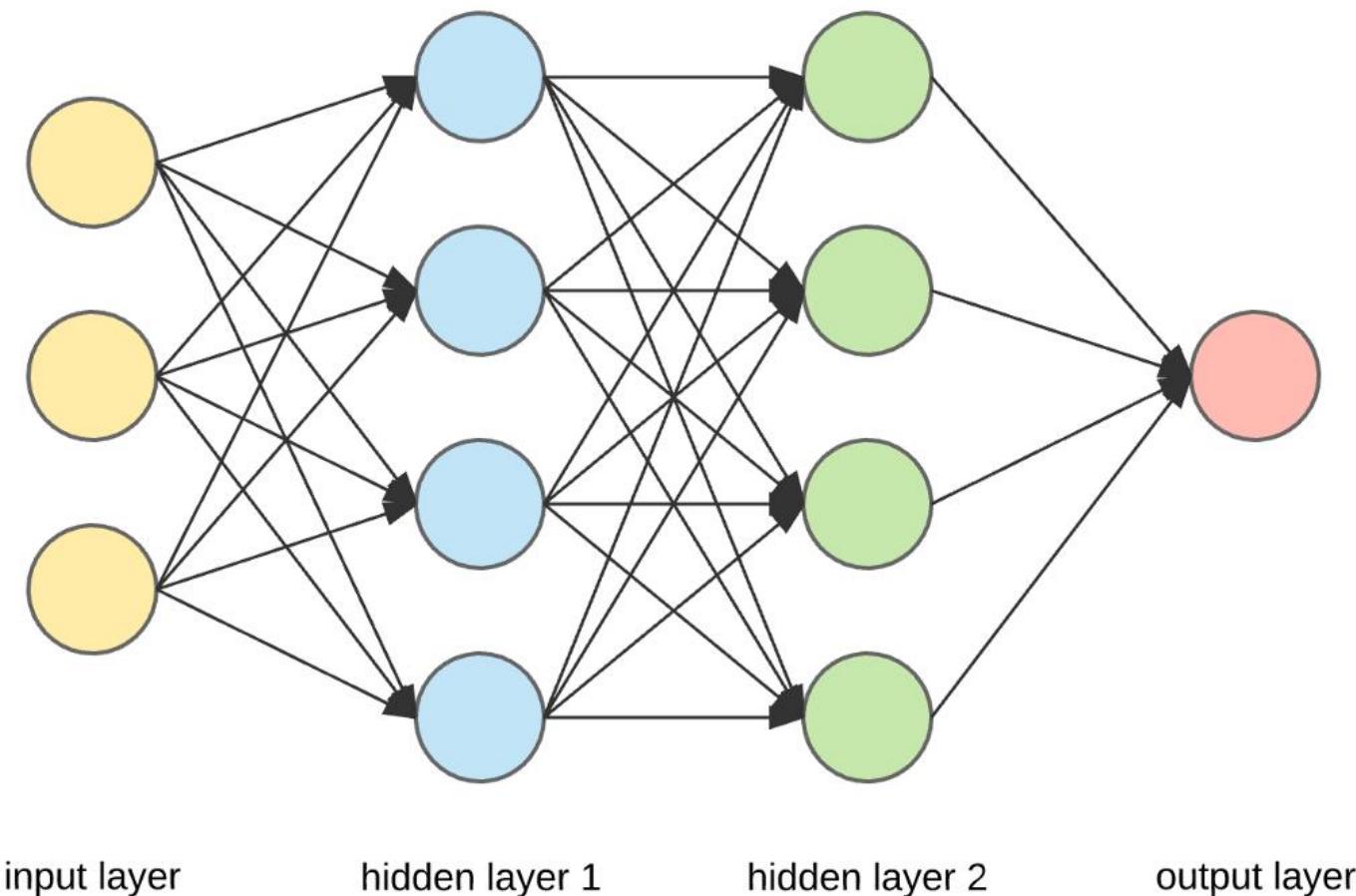


RELU

Artificial Neural Networks

How to handle non-linear problems?

2/ Adding more layers: Fully connected **Deep** Neural Networks



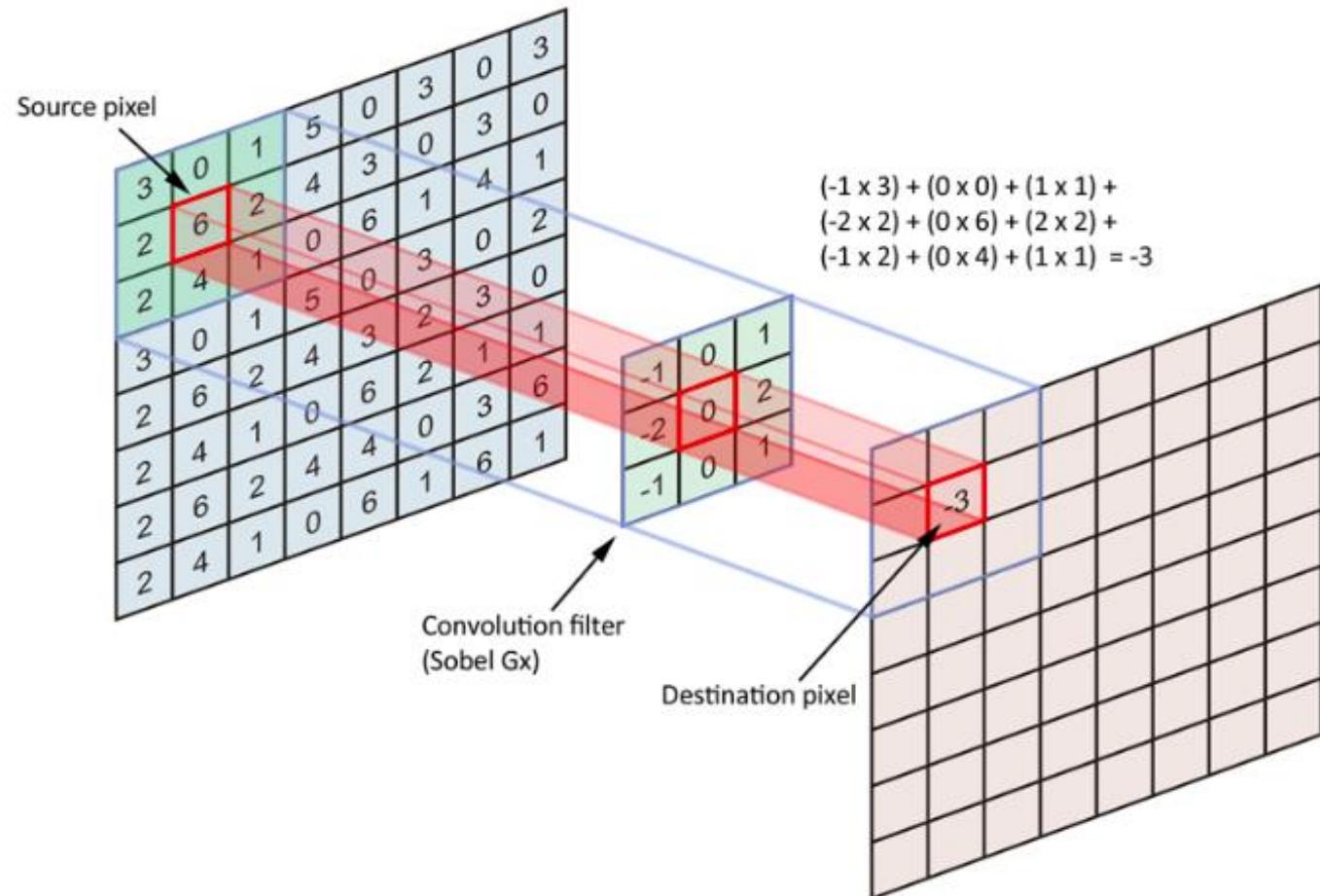
Artificial Neural Networks

How to handle non-linear problems?

3/ Adding convolution layers: Convolutional Neural Networks (CNNs)

Idea:

Passes a filter with a fixed window size and weights to an input image



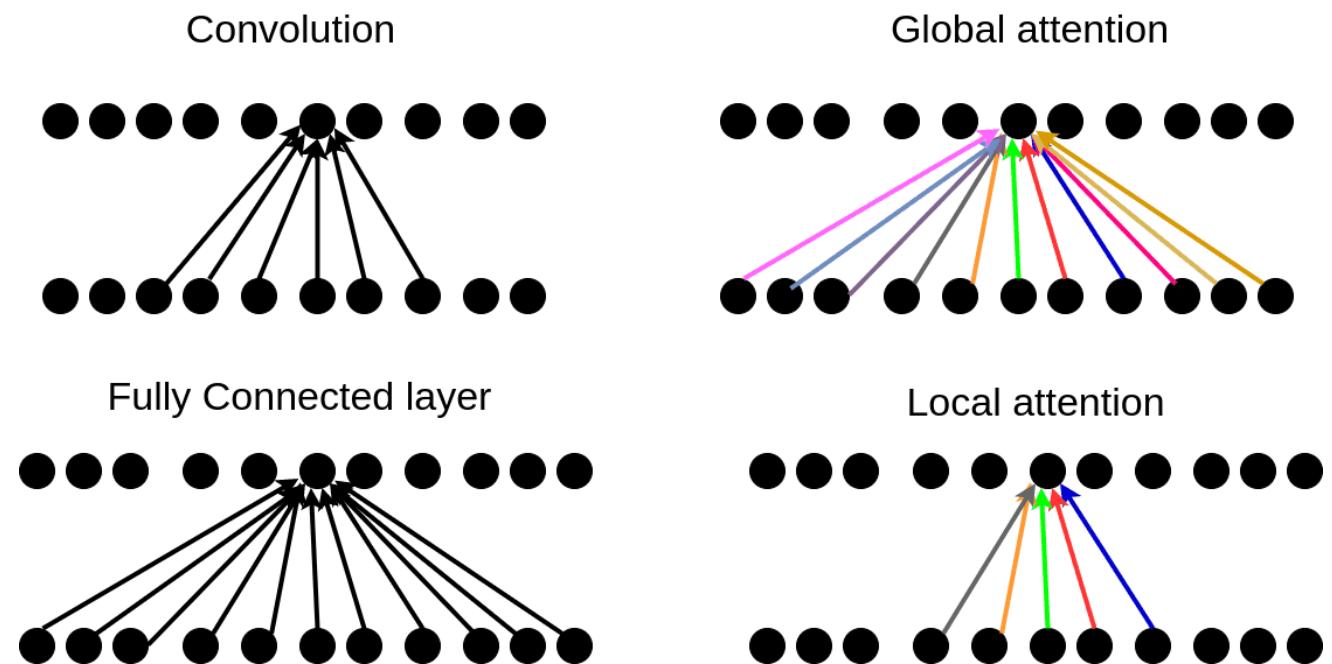
Artificial Neural Networks

How to handle non-linear problems?

4/ Adding attention layers : Transformers

Idea:

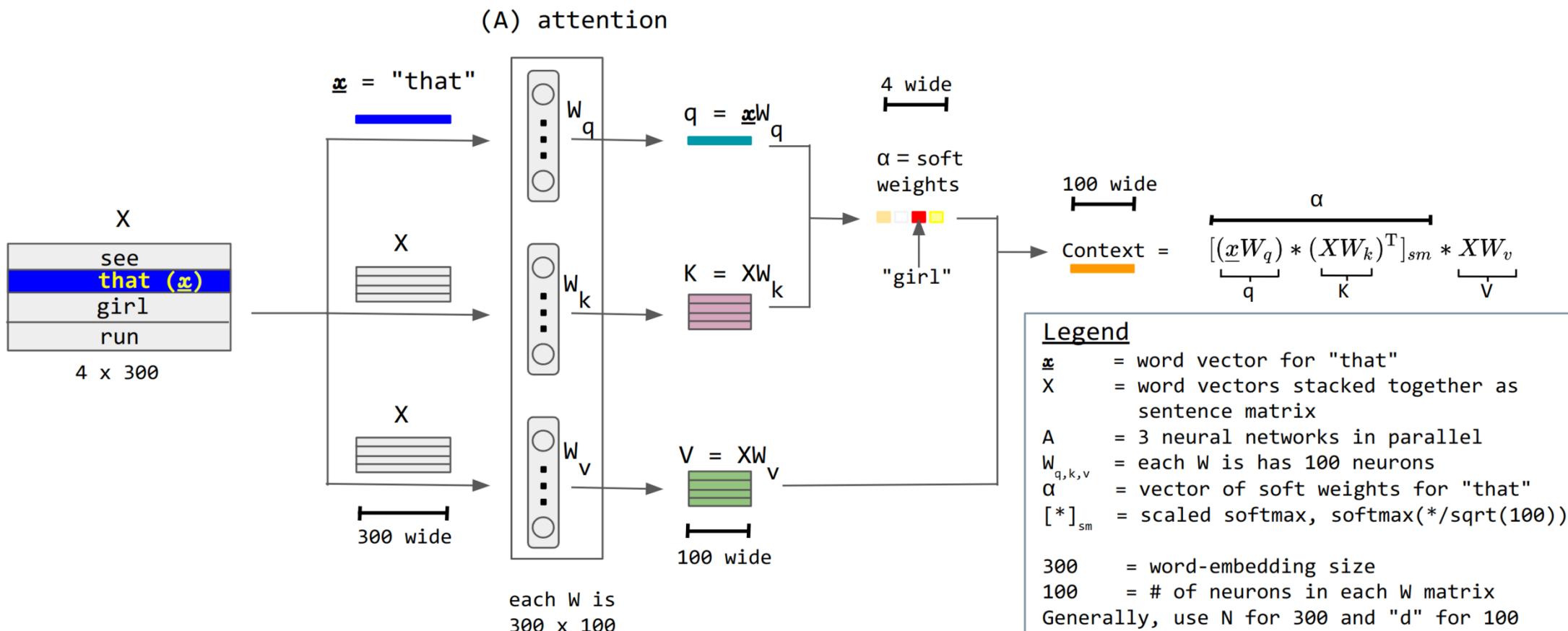
Learn '**soft**' **data-dependent** connections (attention) between the neurons/features at the same time as learning the '**hard**' **usual** connections



Artificial Neural Networks

How to handle non-linear problems?

4/ Adding attention layers : Transformers

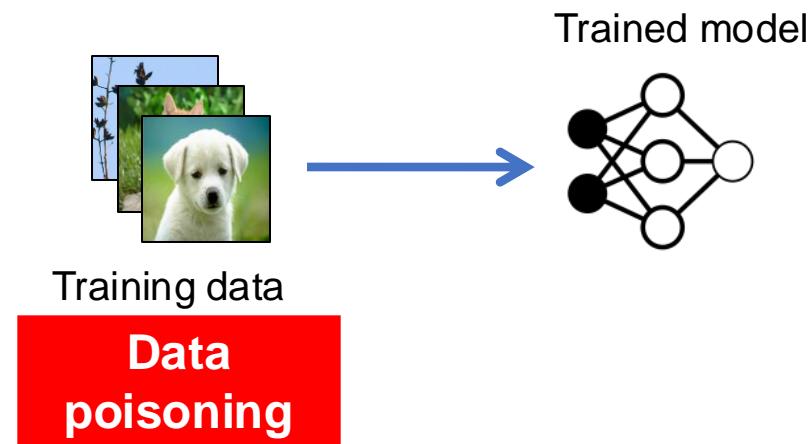


Practical: Part 1
Link to notebook:
<https://shorturl.at/JWNu7>

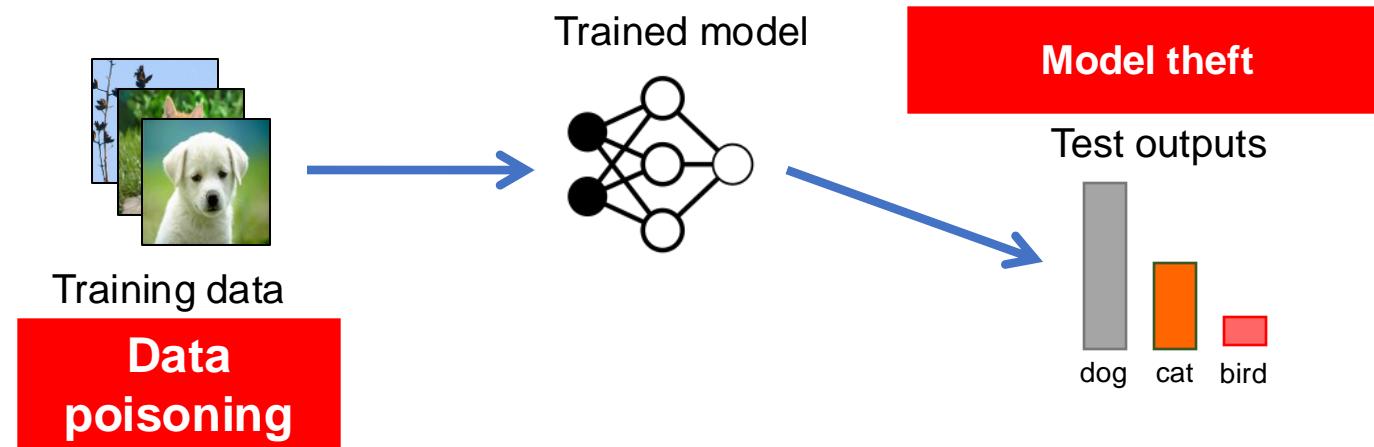
Machine Learning Threats & Attacks

Part I

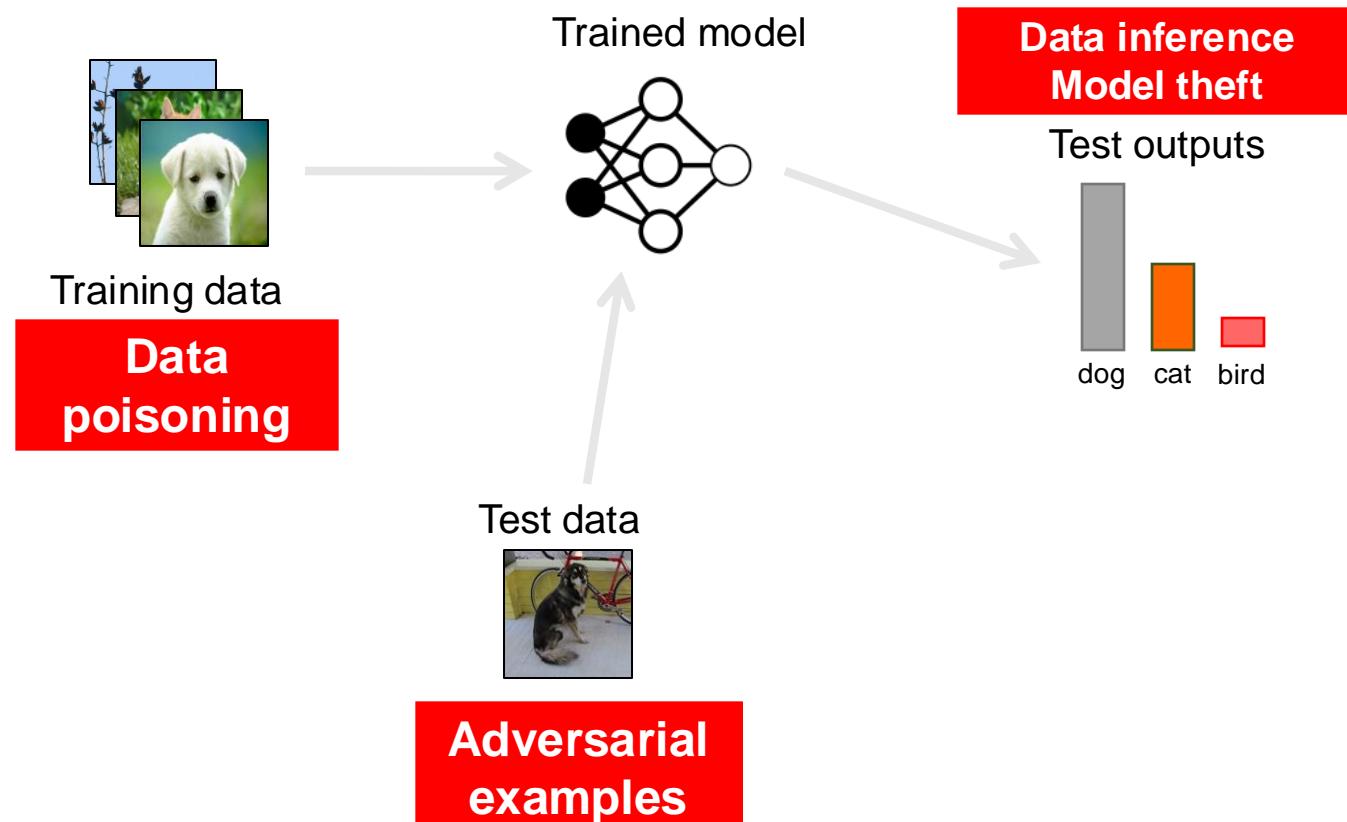
Main threats of ML components



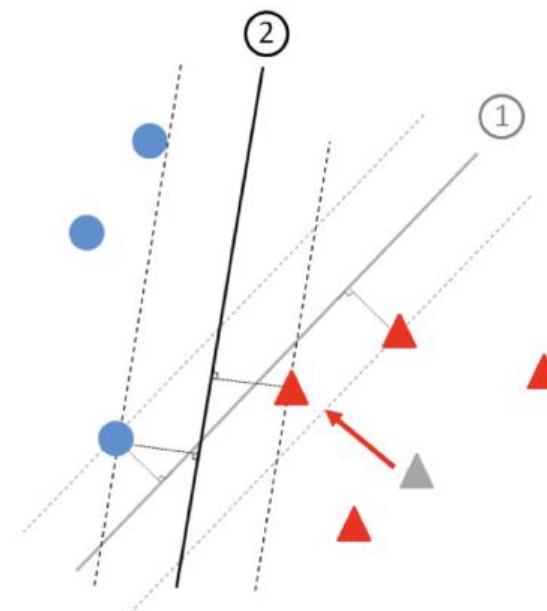
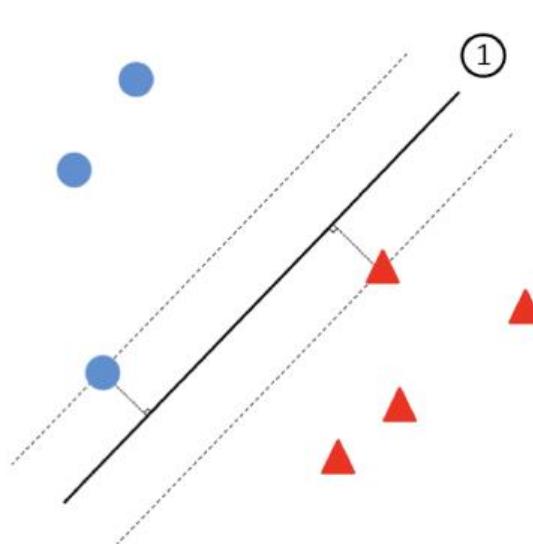
Main threats of ML components



Main threats of ML components



Poisoning attack



Poisoning attack

Availability attack



Poisoning attack

Availability attack



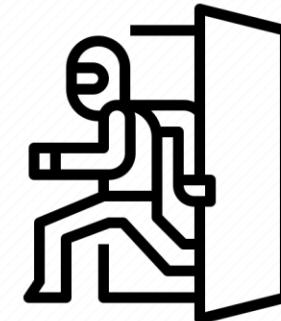
3% poisoning → 11% drop of performance

Poisoning attack

Availability attack



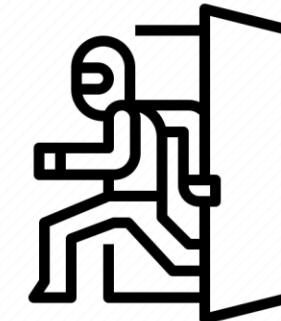
Integrity attack



Poisoning attack



Integrity attack
(Back doors)

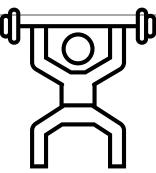


Poisoning attack



Poisoning attack

- Tay bot used the interactions with its Twitter users as training data
- By repeatedly interacting with Tay using racist and offensive language, they were able to bias Tay's dataset towards that language as well
- Within 24 hours of its deployment, Tay had to be decommissioned



Evasion attack and adversarial examples

Original example



Small adversarial noise



Adversarial example



ML predicts:
“Panda”
(80% confidence)

What humans still see



Gibbon

What ML predicts: “Gibbon”
(99% confidence)

Adversarial examples beyond pixels

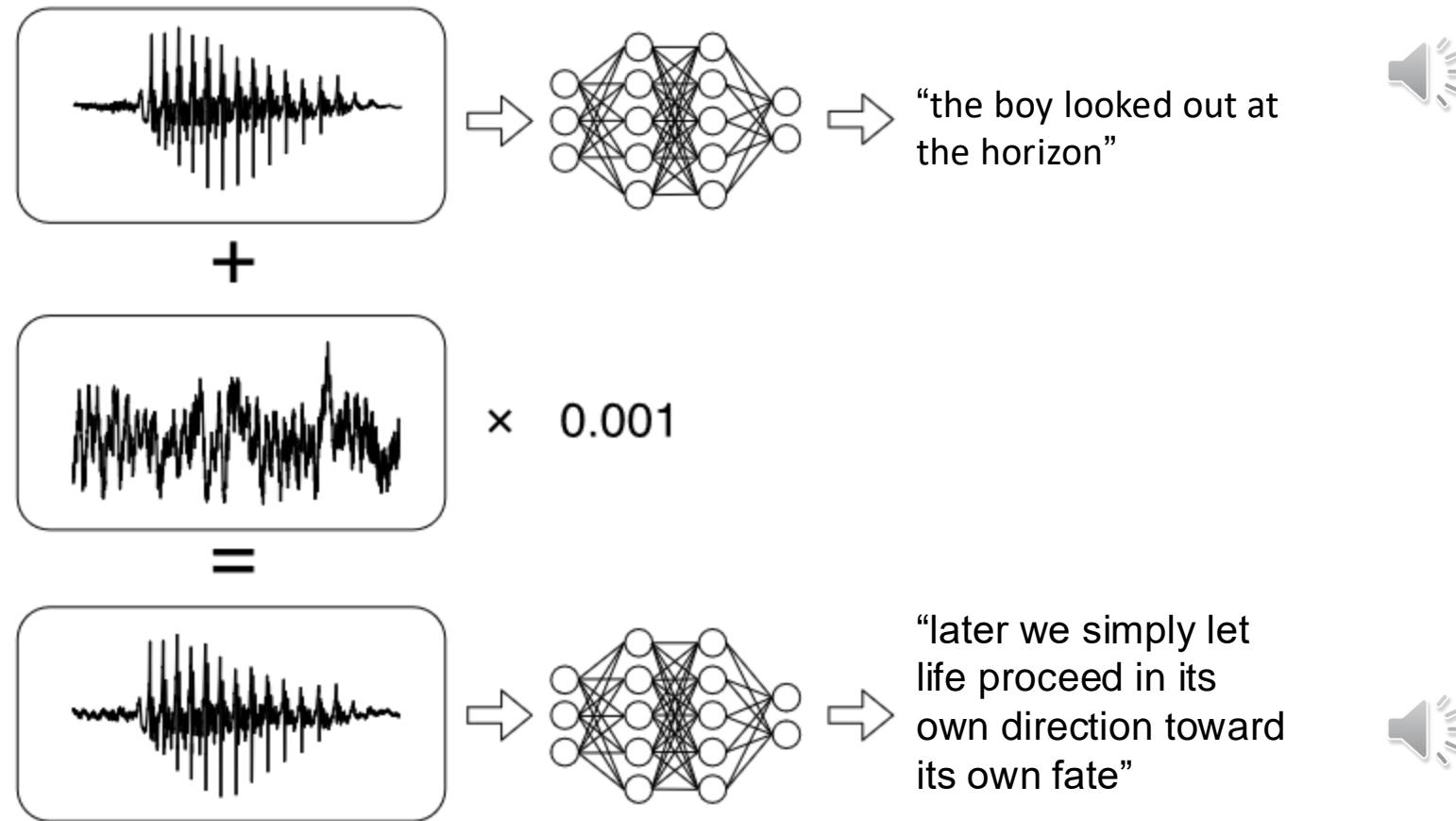


Figure 1. Illustration of our attack: given any waveform, adding a small perturbation makes the result transcribe as any desired target phrase.

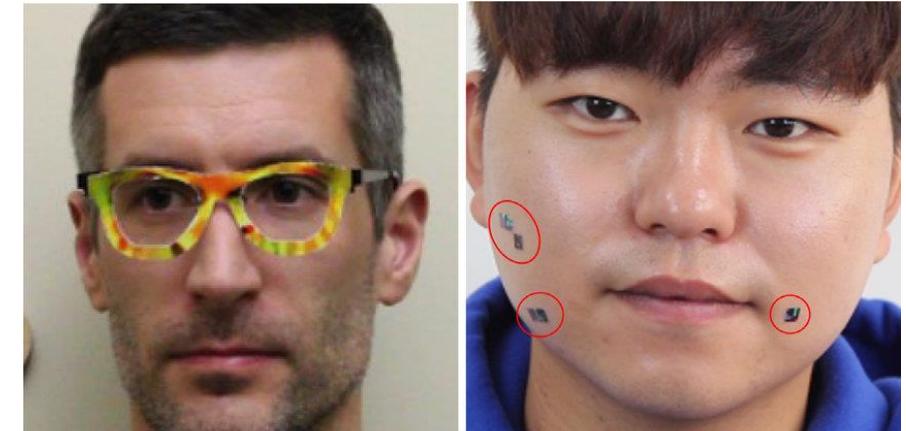
Adversarial examples in the physical world



Adversarial examples in the physical world



Crowd face recognition system



Journal of Information Security and Applications 60 (2021) 102874

Contents lists available at ScienceDirect

Journal of Information Security and Applications



journal homepage: www.elsevier.com/locate/jisa



Adversarial attacks by attaching noise markers on the face against deep face recognition

Gwonsang Ryu ^a, Hosung Park ^b, Daeseon Choi ^{c,*}

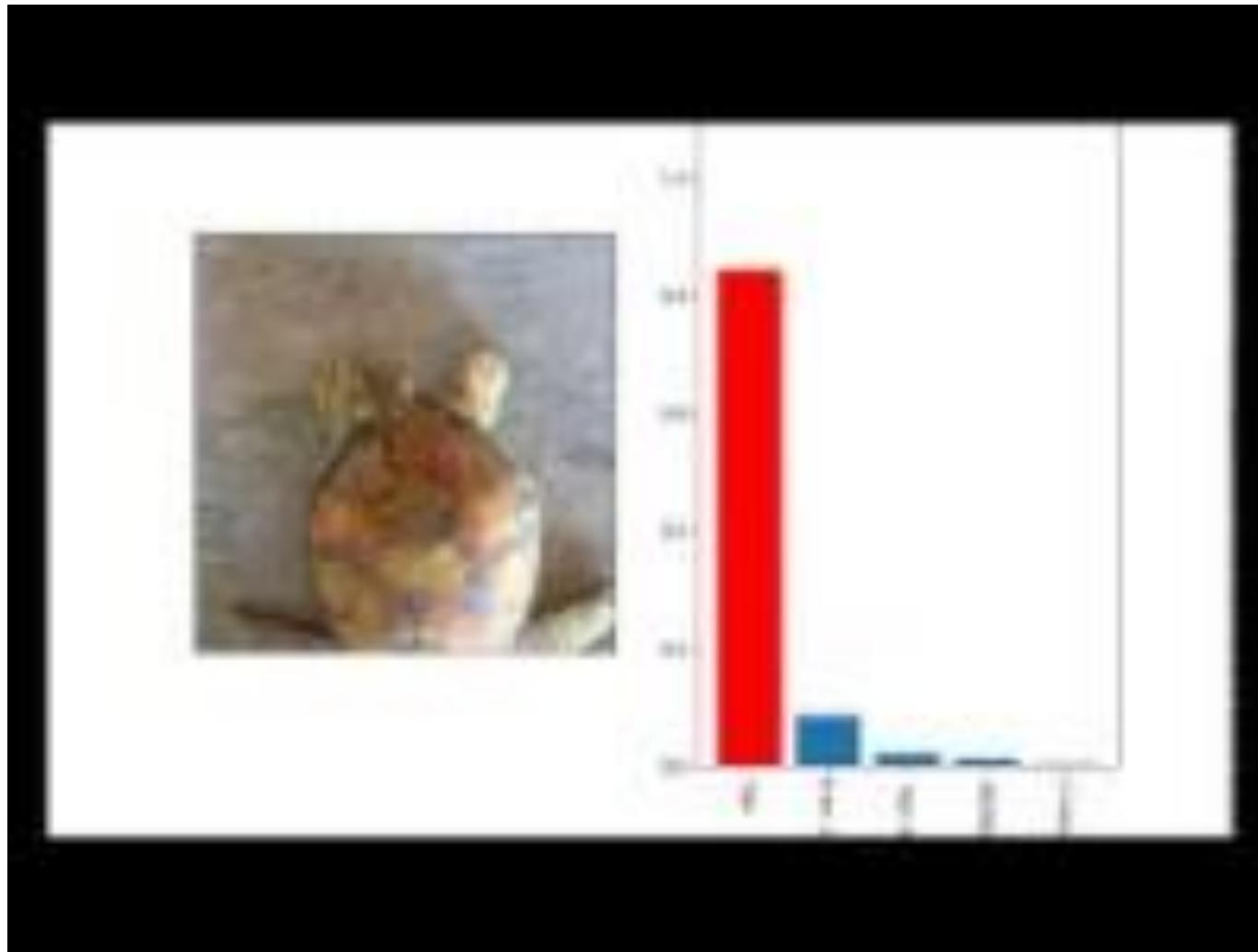
^a Department of Software Convergence, Graduate School of Soongsil University, Seoul, 07027, South Korea

^b Department of Cyber Security and Police, Busan University of Foreign Studies, Busan, 46234, South Korea

^c Department of Software, Soongsil University, Seoul, 07027, South Korea



Adversarial examples in the physical world



Taxonomy > Threat model

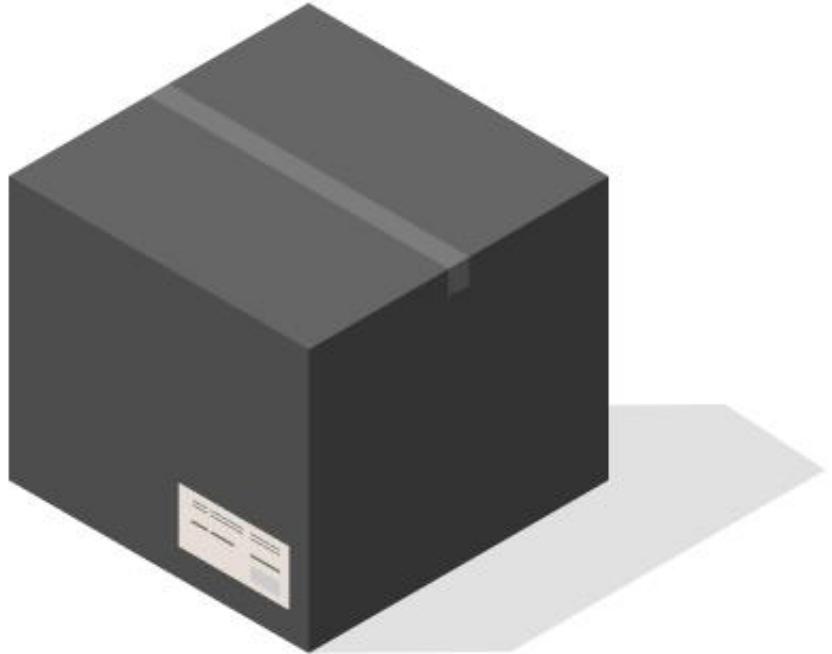
Access to:

- Parameters
- Architecture
- Hyper-parameters
- Gradients
- Training set



**White box - we know
everything**

Taxonomy > Threat model



Black box - we do not know anything

No direct access to anything

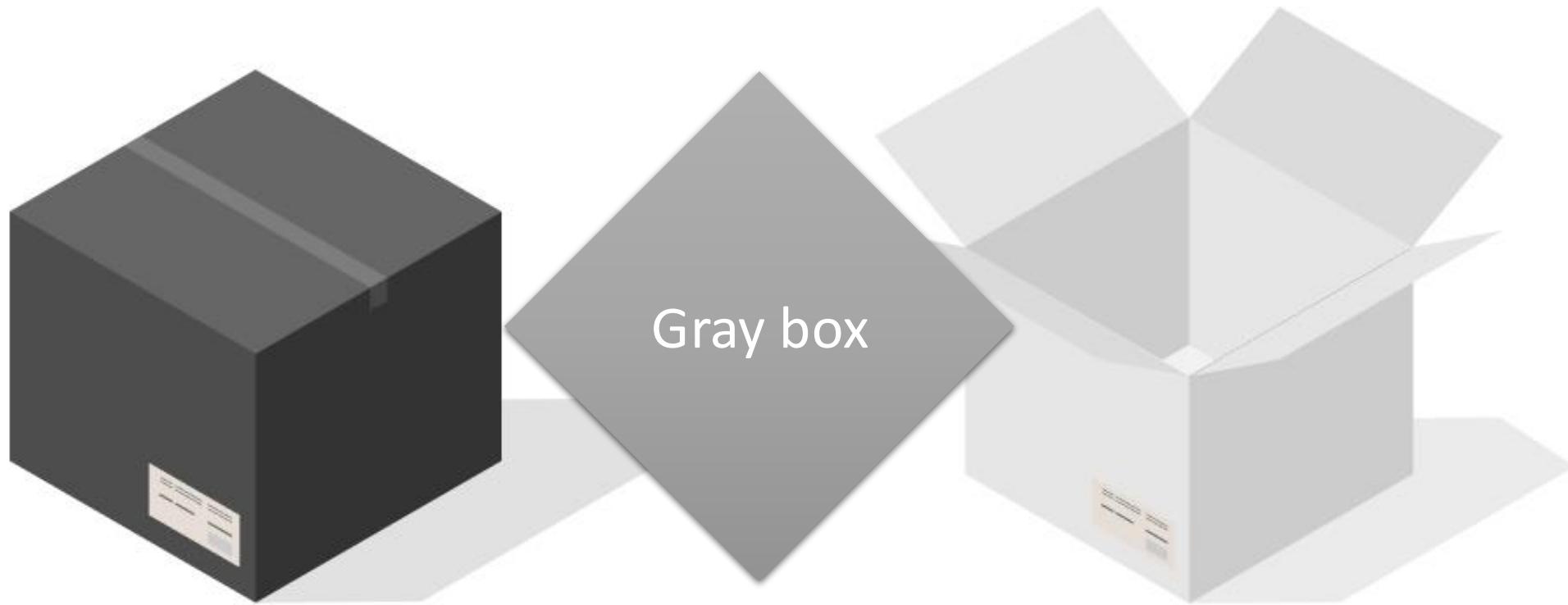
Intuitions about:

- Architecture
- Training distributions

Access through **API queries**

- Prediction logits
- OR **just** predictions

Taxonomy > Threat model



**Black box - we do not
know anything**

**White box - we know
everything**

Taxonomy > Threat model

Non-targeted

Any wrong result of the ML model is satisfactory

Targeted

We expect a specific wrong result of the ML model

Harder

Whitebox Evasion attacks

In computer vision

Part II

Formal definition

- Given an image x , which is labeled by the classifier as class q , i.e., $C(x)=q$
- Create an adversarial image x_{adv} by adding small perturbations δ to the original image, i.e., $x_{adv} = x + \delta$,
such that the distance $D(x, x_{adv}) = D(x, x+\delta)$ is minimal
and that the classifier assigns a label to the adversarial image that is different than q , i.e., $C(x_{adv}) = C(x+\delta) = t \neq q$

Formal definition > Minimization problem

minimize $\mathcal{D}(x, x + \delta)$

such that $C(x + \delta) = t$

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

distance between x and $x + \delta$

$x + \delta$ is classified as target class t

each element of $x + \delta$ is in $[0, 1]$ (to be a valid image)

Formal definition > Distance metrics

Distance metrics between x and x_{adv} : $D(x, x_{adv})$

- ℓ_0 norm: the number of elements in x_{adv} such that $x^i \neq x_{adv}^i$
Corresponds to the number of pixels that have been changed in the image x_{adv}

- ℓ_1 norm: city-block distance, or Manhattan distance
$$\ell_1 = |x^1 - x_{adv}^1| + |x^2 - x_{adv}^2| + \dots + |x^n - x_{adv}^n|$$

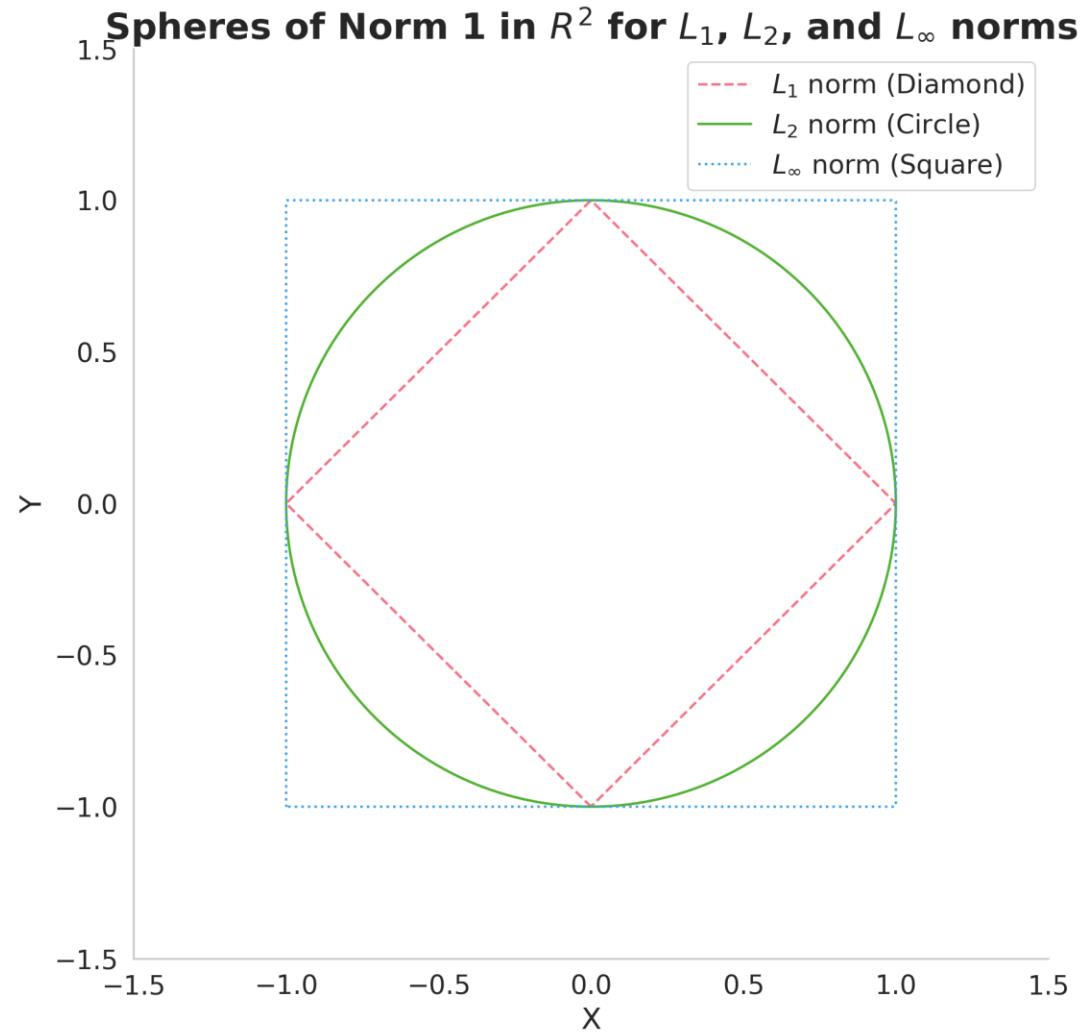
- ℓ_2 norm: Euclidean distance, or mean-squared error

$$\ell_2 = \sqrt{(x^1 - x_{adv}^1)^2 + (x^2 - x_{adv}^2)^2 + \dots + (x^n - x_{adv}^n)^2}$$

- ℓ_∞ norm: measures the maximum change to any of the pixels in the x_{adv} image

$$\ell_\infty = \max(|x^1 - x_{adv}^1|, |x^2 - x_{adv}^2|, \dots, |x^n - x_{adv}^n|)$$

Illustration > Distance metrics



Fast Gradient Sign Method (FGSM)

An adversarial image x_{adv} is created by adding perturbation noise to an image x

$$x_{adv} = x + \epsilon \cdot \text{sign}(\nabla_x \mathcal{L}(h(x, w), y))$$

- Notation: input image x , cost function \mathcal{L} , NN model h , NN weights (parameters) w , gradient ∇ , noise magnitude ϵ
- Perturbation noise is calculated as the gradient of the loss function \mathcal{L} with respect to the input image x for the true class label y
- This increases the loss for the true class $y \rightarrow$ the model misclassifies the image x_{adv}

$$\text{sgn}(x) := \begin{cases} -1 & \text{if } x < 0, \\ 0 & \text{if } x = 0, \\ 1 & \text{if } x > 0. \end{cases}$$

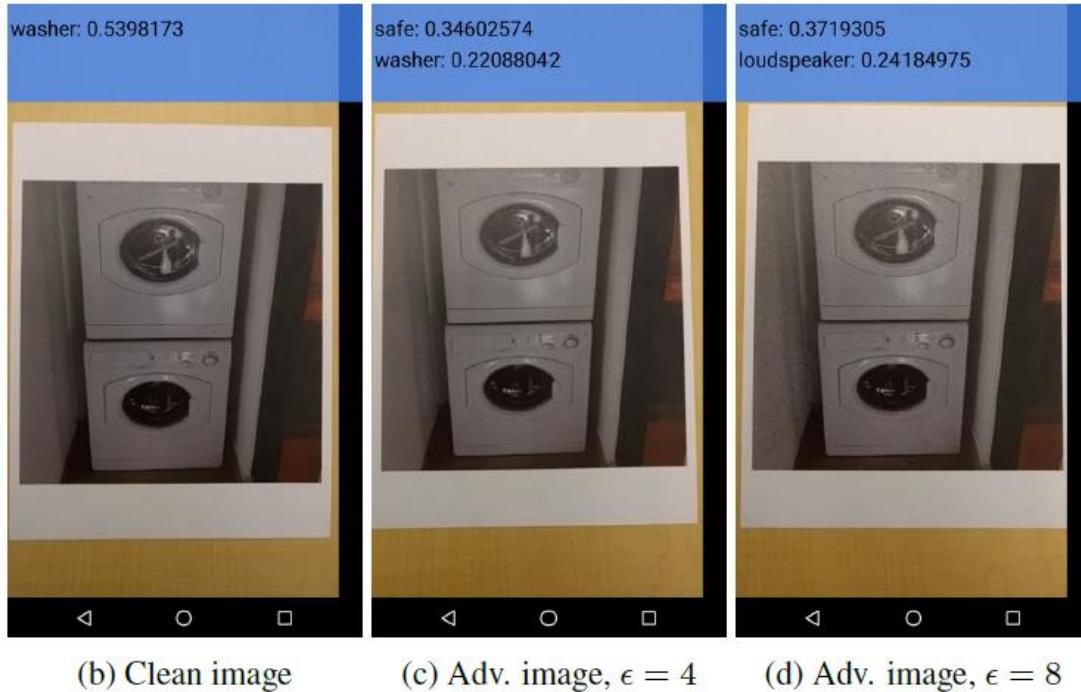
FGSM > Why does it work?

The sign and magnitude of the gradient give the direction and the slope of the steepest descent

- To minimize the loss function, the weights w are changed in the opposite direction of the gradient, i.e., $w = w - \alpha \frac{\partial \text{loss}}{\partial w}$
- FGSM operates on the inputs, they are changed in the direction of the gradient: $x = x + \alpha \frac{\partial \text{loss}}{\partial w}$



Basic iterative method (BIM; Kurakin, 2017)



What if we compute FGSM
iteratively with a **small** step?

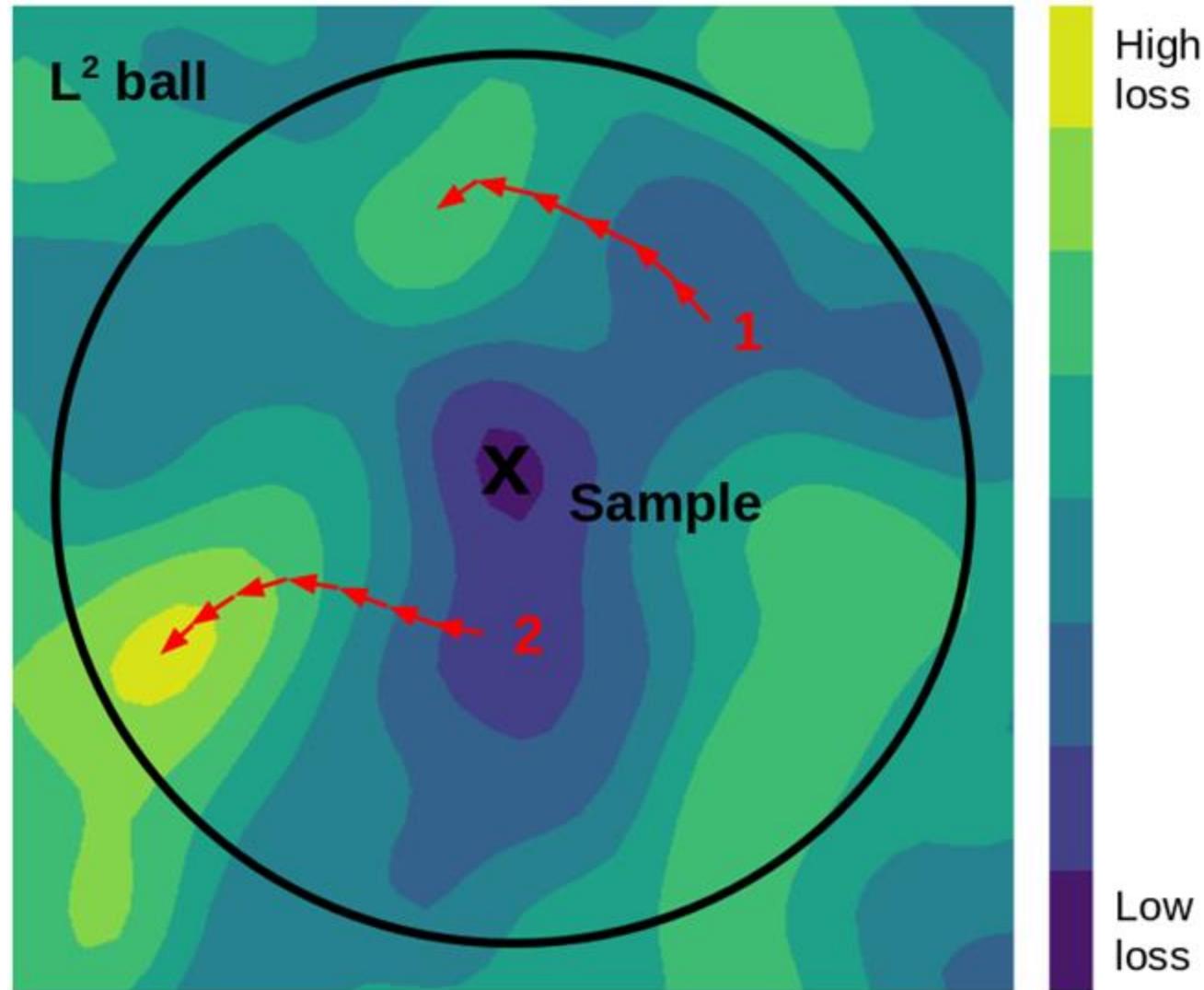
$$x_{t+1} = x_t + \alpha \operatorname{sign}(\nabla_x \mathcal{L}(h(x, w), y))$$

Projected Gradient Descent (PGD, Madry 2017)

Extension of BIM with projection and random starts

PGD is regarded as the strongest first-order attack

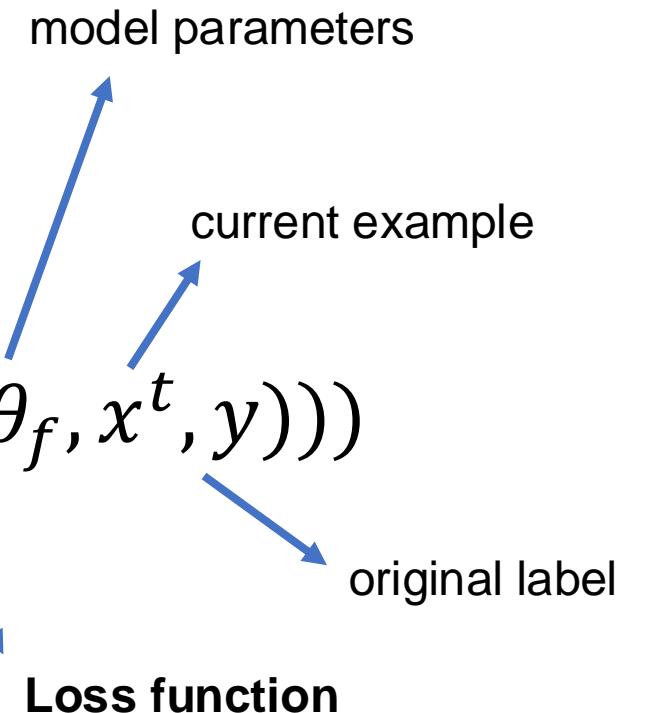
→ *First-order attack means that the adversary uses only the gradients of the loss function with respect to the input*



Projected Gradient Descent (PGD)

Iteratively compute:

$$x^{t+1} = \text{Clip}_{x,\epsilon}(x^t + \alpha \cdot \text{sign}(\nabla_{x^t}(L(\theta_f, x^t, y))))$$



Projected Gradient Descent (PGD)

Iteratively compute:

$$x^{t+1} = \text{Clip}_{x,\epsilon}(x^t + \alpha \cdot \text{sign}(\nabla_{x^t}(L(\theta_f, x^t, y))))$$

Normalize the gradient

model parameters

current example

original label

Gradient of the loss around x^t , i.e. the “slope”

Loss function

Projected Gradient Descent (PGD)

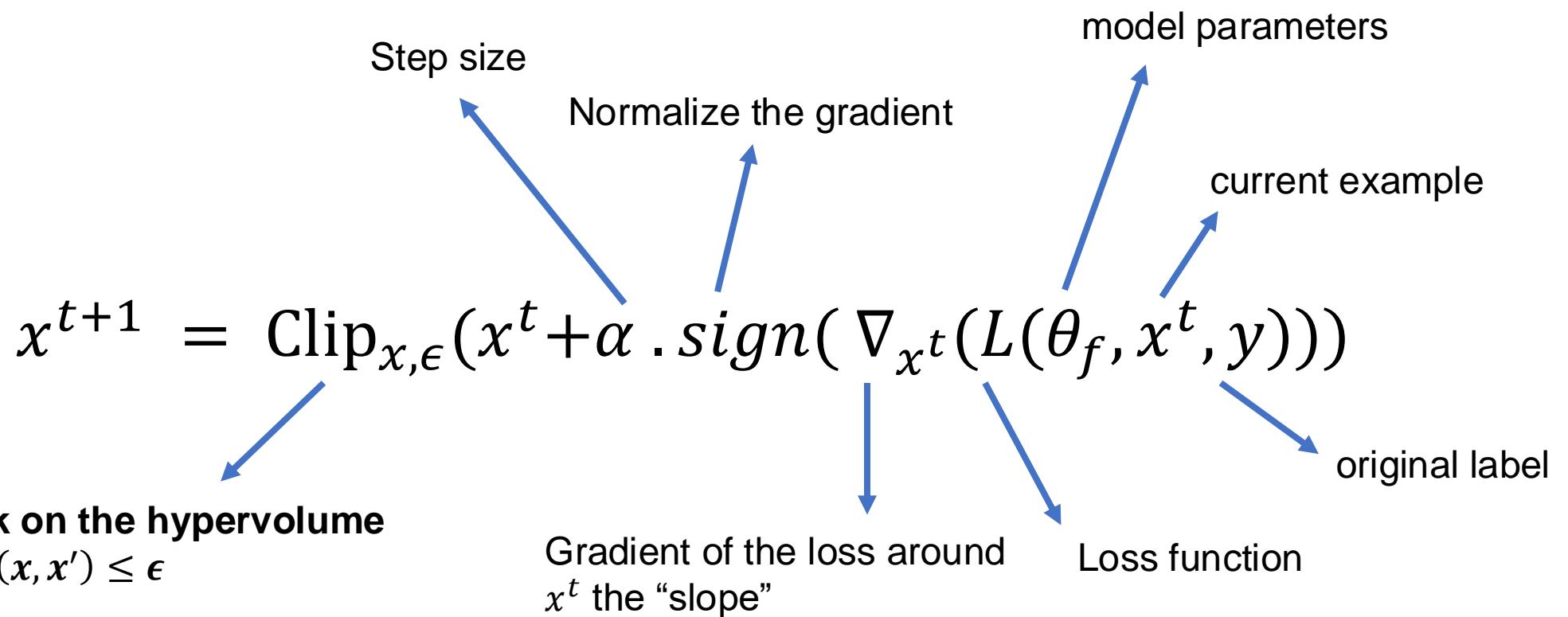
Iteratively compute:

$$x^{t+1} = \text{Clip}_{x,\epsilon}(x^t + \alpha \cdot \text{sign}(\nabla_{x^t}(L(\theta_f, x^t, y))))$$

The diagram illustrates the components of the PGD update equation. It shows the flow of information from the right side of the equation towards the left. On the right, there are four inputs: 'model parameters' (represented by a blue arrow pointing up), 'current example' (represented by a blue arrow pointing up), 'original label' (represented by a blue arrow pointing down), and 'Loss function' (represented by a blue arrow pointing down). These four inputs converge at a point labeled 'Gradient of the loss around x^t the “slope”'. From this point, two arrows point left: one labeled 'Normalize the gradient' (blue arrow pointing up) and another labeled 'Step size' (blue arrow pointing up). Both of these arrows point to the term $\alpha \cdot \text{sign}(\nabla_{x^t}(L(\theta_f, x^t, y)))$ in the equation.

Projected Gradient Descent (PGD)

Iteratively compute:



Targeted PGD

For a targeted attack, if the target class label is denoted t , adversarial examples are created by using

$$x_{adv}^t = \Pi_{\epsilon} \left(x^{t-1} - \alpha \cdot \text{sign}(\nabla_x \mathcal{L}(h(x^{t-1}), t)) \right)$$

- I.e., it is based on minimizing the loss function with respect to the target class t
- This is opposite to non-targeted attacks, which maximize the loss function with respect to the true class label

Other gradient attacks

MI-FGSM (Dong, 2018)

Algorithm 1 MI-FGSM

Input: A classifier f with loss function J ; a real example \mathbf{x} and ground-truth label y ;

Input: The size of perturbation ϵ ; iterations T and decay factor μ .

Output: An adversarial example \mathbf{x}^* with $\|\mathbf{x}^* - \mathbf{x}\|_\infty \leq \epsilon$.

1: $\alpha = \epsilon/T$;

2: $\mathbf{g}_0 = 0$; $\mathbf{x}_0^* = \mathbf{x}$;

3: **for** $t = 0$ to $T - 1$ **do**

4: Input \mathbf{x}_t^* to f and obtain the gradient $\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)$;

5: Update \mathbf{g}_{t+1} by accumulating the velocity vector in the gradient direction as

$$\mathbf{g}_{t+1} = \mu \cdot \mathbf{g}_t + \frac{\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)}{\|\nabla_{\mathbf{x}} J(\mathbf{x}_t^*, y)\|_1}; \quad (6)$$

6: Update \mathbf{x}_{t+1}^* by applying the sign gradient as

$$\mathbf{x}_{t+1}^* = \mathbf{x}_t^* + \alpha \cdot \text{sign}(\mathbf{g}_{t+1}); \quad (7)$$

7: **end for**

8: **return** $\mathbf{x}^* = \mathbf{x}_T^*$.

Other gradient attacks

Carlini-Wagner (CW; Carlini, 2017)

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

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

$$f_2(x') = (\max_{i \neq t} (F(x')_i) - F(x')_t)^+$$

Stronger attacks

PGD has limitations!

1) Fixed step size

Unaware of iteration budget
Agnostic of previous iterations

Stronger attacks

PGD has limitations!

1) Fixed step size

Unaware of iteration budget
Agnostic of previous iterations

APGD adds a momentum

Algorithm 1 APGD

```

1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\max}, f_{\max}$ 
3:  $x^{(1)} \leftarrow P_S(x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\max} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\max} \leftarrow x^{(0)}$  if  $f_{\max} \equiv f(x^{(0)})$  else  $x_{\max} \leftarrow x^{(1)}$ 
6: for  $k = 1$  to  $N_{\text{iter}} - 1$  do
7:    $z^{(k+1)} \leftarrow P_S(x^{(k)} + \eta \nabla f(x^{(k)}))$ 
8:    $x^{(k+1)} \leftarrow P_S\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right.$ 
       $\quad \quad \quad \left.+(1 - \alpha)(x^{(k)} - x^{(k-1)})\right)$ 
9:   if  $f(x^{(k+1)}) > f_{\max}$  then
10:     $x_{\max} \leftarrow x^{(k+1)}$  and  $f_{\max} \leftarrow f(x^{(k+1)})$ 
11:   end if
12:   if  $k \in W$  then
13:     if Condition 1 or Condition 2 then
14:        $\eta \leftarrow \eta/2$  and  $x^{(k+1)} \leftarrow x_{\max}$ 
15:     end if
16:   end if
17: end for

```

Stronger attacks

PGD has limitations!

1) Fixed step size

Unaware of iteration budget
Agnostic of previous iterations

APGD adds a momentum

APGD proposes adaptive steps in two phases:
a/ Large step η for exploration
b/ Small step η for exploitation

Algorithm 1 APGD

```

1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\max}, f_{\max}$ 
3:  $x^{(1)} \leftarrow P_S(x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\max} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\max} \leftarrow x^{(0)}$  if  $f_{\max} \equiv f(x^{(0)})$  else  $x_{\max} \leftarrow x^{(1)}$ 
6: for  $k = 1$  to  $N_{\text{iter}} - 1$  do
7:    $z^{(k+1)} \leftarrow P_S(x^{(k)} + \eta \nabla f(x^{(k)}))$ 
8:    $x^{(k+1)} \leftarrow P_S\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right.$ 
       $\left.+ (1 - \alpha)(x^{(k)} - x^{(k-1)})\right)$ 
9:   if  $f(x^{(k+1)}) > f_{\max}$  then
10:     $x_{\max} \leftarrow x^{(k+1)}$  and  $f_{\max} \leftarrow f(x^{(k+1)})$ 
11:   end if
12:   if  $k \in W$  then
13:     if Condition 1 or Condition 2 then
14:        $\eta \leftarrow \eta/2$  and  $x^{(k+1)} \leftarrow x_{\max}$ 
15:     end if
16:   end if
17: end for

```

Stronger attacks

PGD has limitations!

2) For converging models, CE leads to vanishing gradients

Unaware of iteration budget
Agnostic of previous iterations

APGD adds a momentum

APGD proposes adaptive steps in two phases:
a/ Large step η for exploration
b/ Small step η for exploitation

APGD replaces its loss function

Algorithm 1 APGD

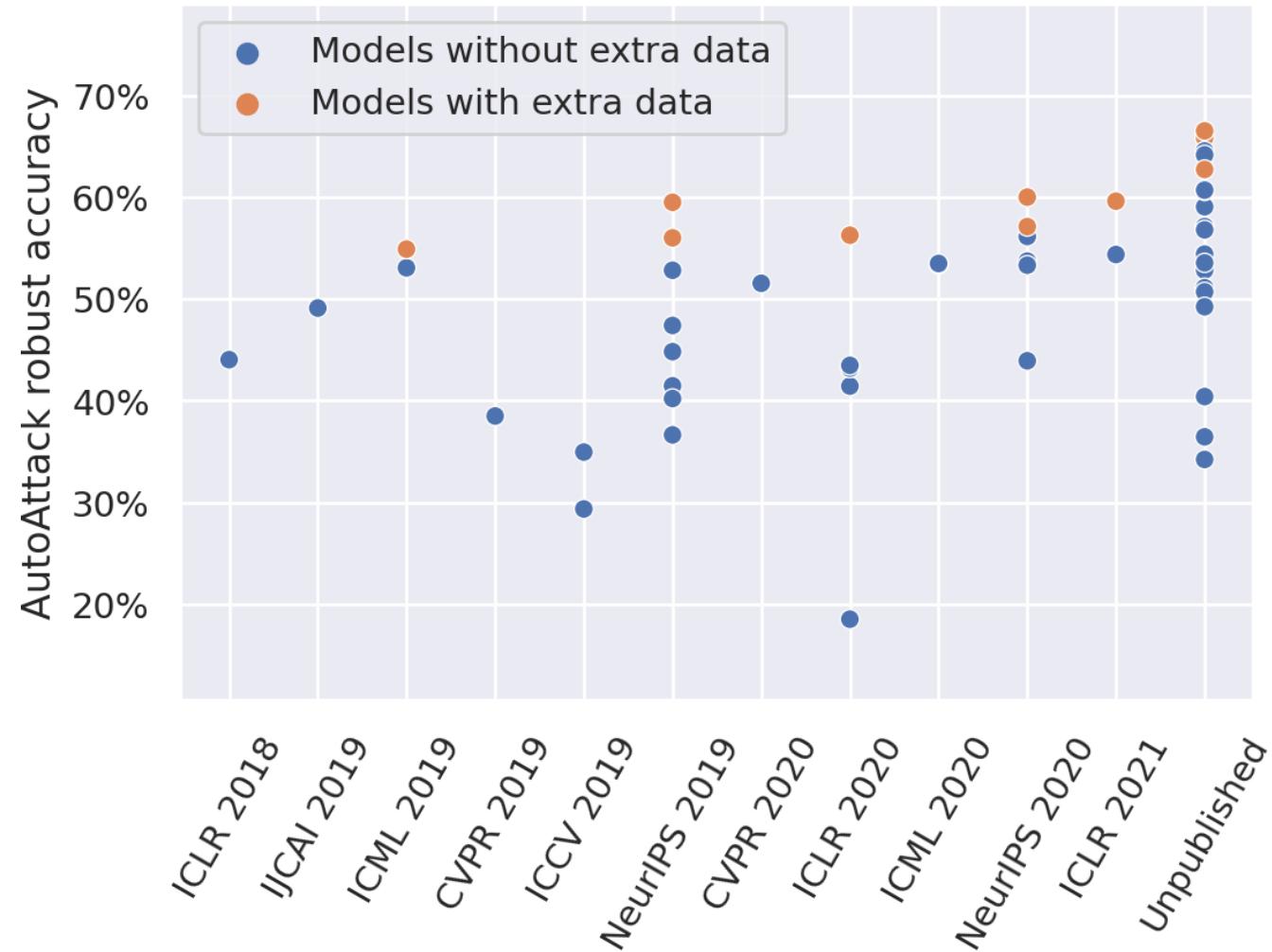
```

1: Input:  $f, S, x^{(0)}, \eta, N_{\text{iter}}, W = \{w_0, \dots, w_n\}$ 
2: Output:  $x_{\max}, f_{\max}$ 
3:  $x^{(1)} \leftarrow P_S(x^{(0)} + \eta \nabla f(x^{(0)}))$ 
4:  $f_{\max} \leftarrow \max\{f(x^{(0)}), f(x^{(1)})\}$ 
5:  $x_{\max} \leftarrow x^{(0)}$  if  $f_{\max} \equiv f(x^{(0)})$  else  $x_{\max} \leftarrow x^{(1)}$ 
6: for  $k = 1$  to  $N_{\text{iter}} - 1$  do
7:    $z^{(k+1)} \leftarrow P_S(x^{(k)} + \eta \nabla f(x^{(k)}))$ 
8:    $x^{(k+1)} \leftarrow P_S\left(x^{(k)} + \alpha(z^{(k+1)} - x^{(k)})\right.$ 
       $\left.+ (1 - \alpha)(x^{(k)} - x^{(k-1)})\right)$ 
9:   if  $f(x^{(k+1)}) > f_{\max}$  then
10:     $x_{\max} \leftarrow x^{(k+1)}$  and  $f_{\max} \leftarrow f(x^{(k+1)})$ 
11:   end if
12:   if  $DLR(x, y) < \epsilon$  then
13:      $\eta \leftarrow \eta / 2$ 
14:   end if
15:    $DLR(x, y) = -\frac{z_y - \max_{i \neq y} z_i}{z_{\pi_1} - z_{\pi_3}}$ , where  $\pi$  is decreasing ordering of  $z$ 
16:   end if
17: end for

```

Stronger attacks > Whitebox Benchmarks

Robustbench



Blackbox Evasion attacks In computer vision

Part III

Blackbox evasion attacks

Transfer attacks

Query attacks

Craft whitebox adversarial examples on
a surrogate model transferable to target

Craft adversarial examples by accessing
the output of target model with queries

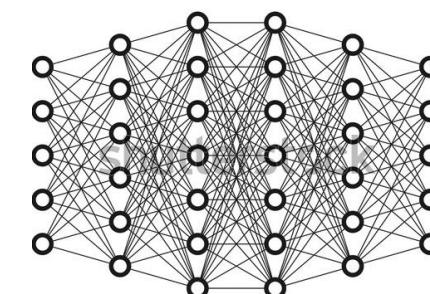
Transfer attacks

Vanilla transfer attacks

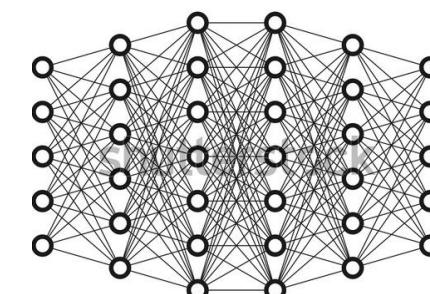
Hope that the surrogate and target have similar decision boundaries



Adversarial image



Surrogate model



Target model

Dog!

Dog!

Transfer attacks

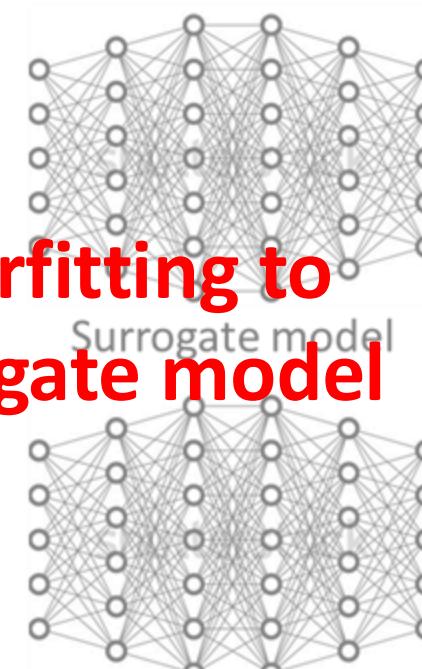
Vanilla transfer attacks

Hope that the surrogate and target have similar decision boundaries

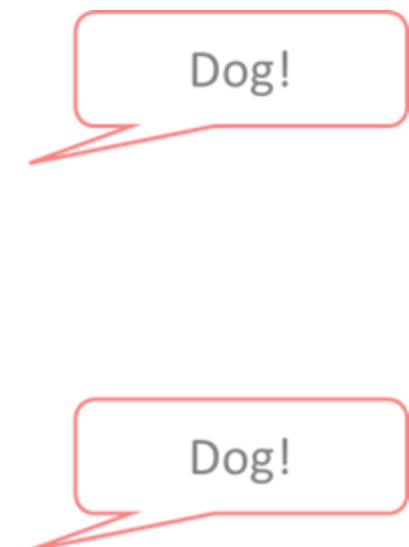


Adversarial image

**Overfitting to
surrogate model**



Surrogate model

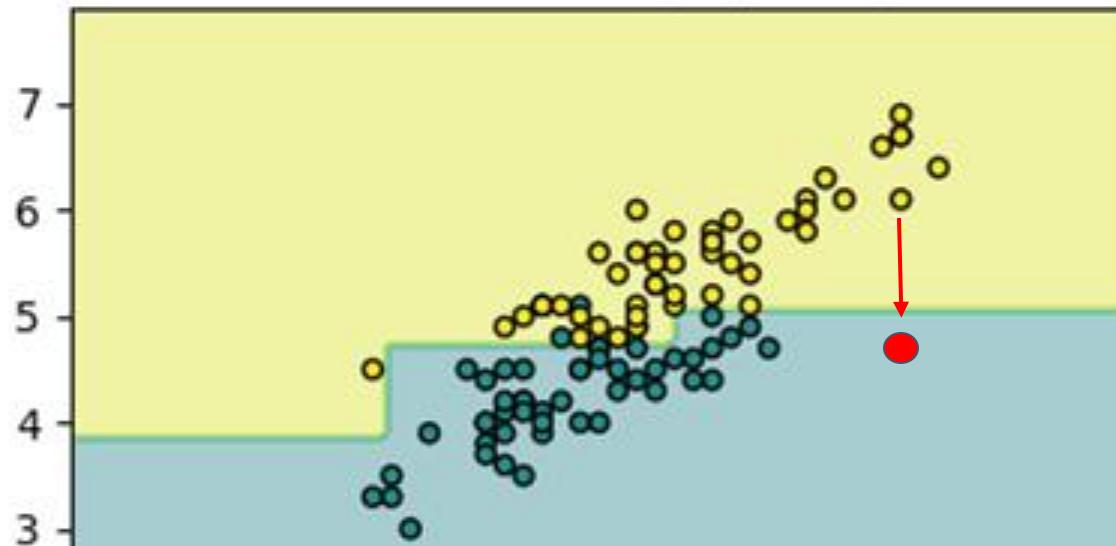


Dog!

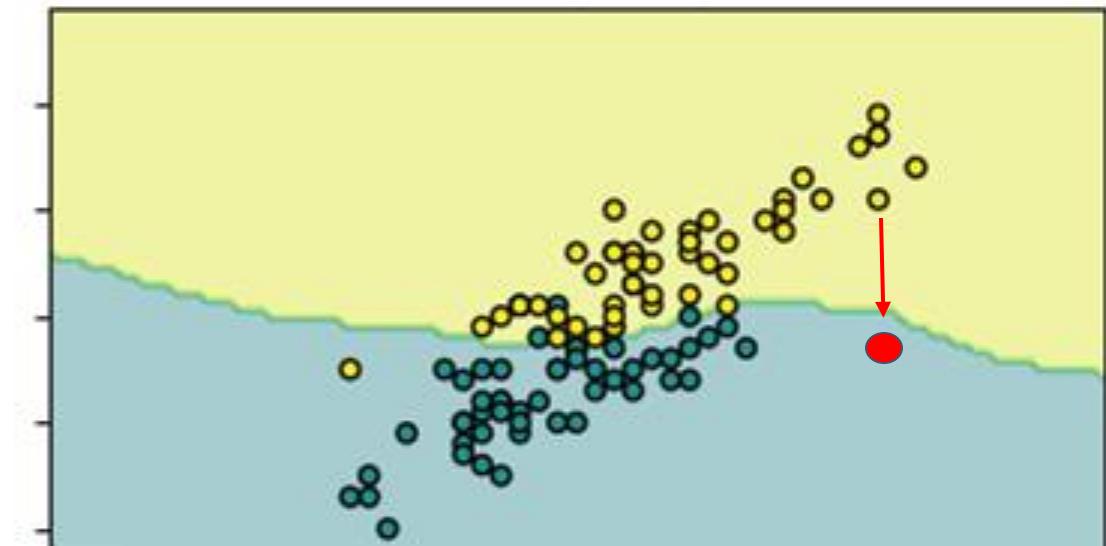
Transfer attacks

Generally succeed...

Decision Tree (depth=4)



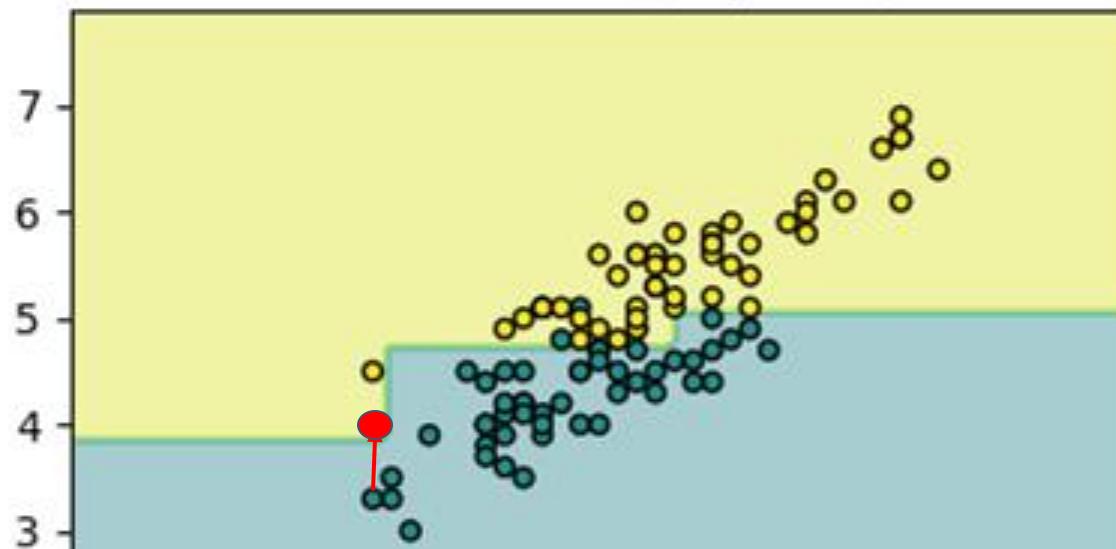
KNN (k=7)



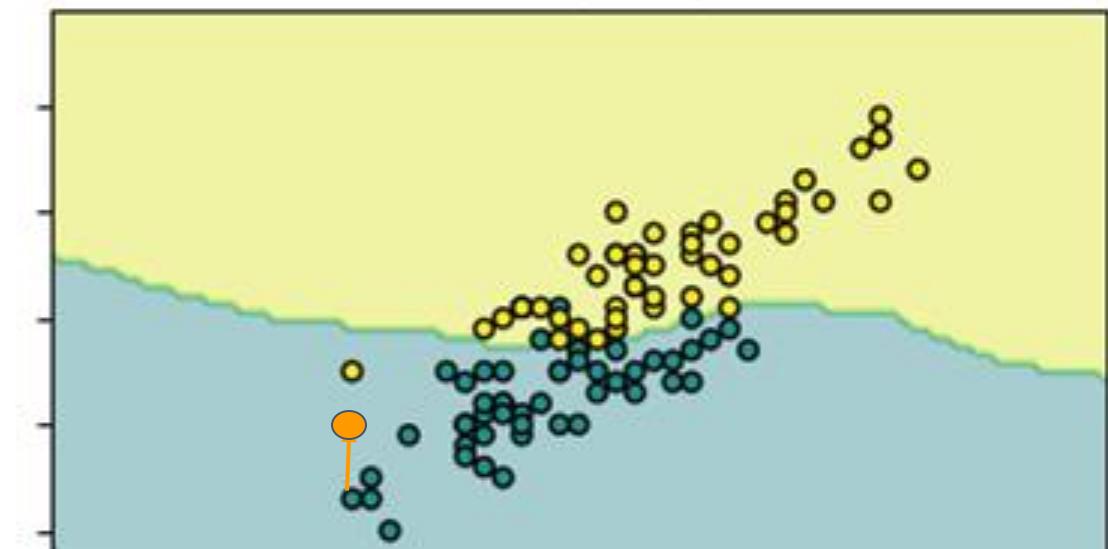
Transfer attacks

but sometime fail!

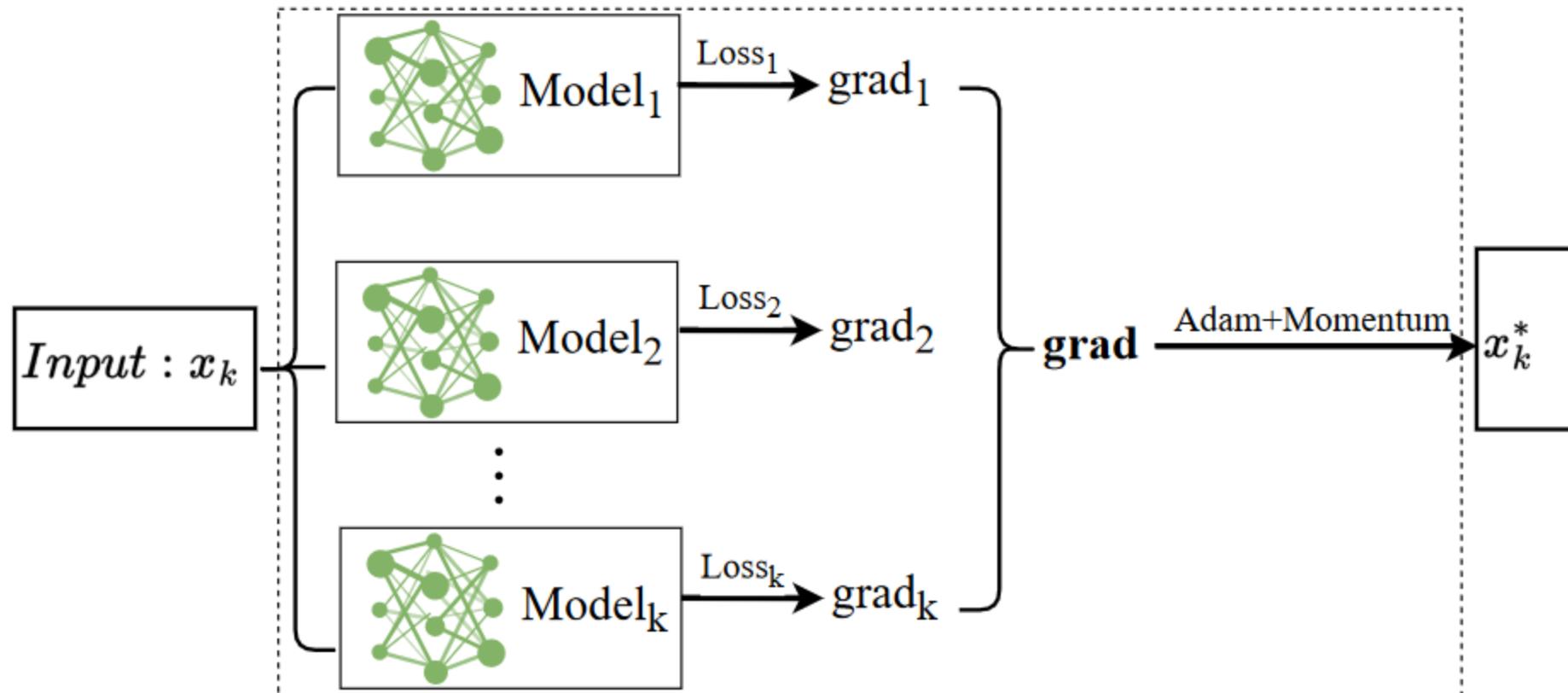
Decision Tree (depth=4)



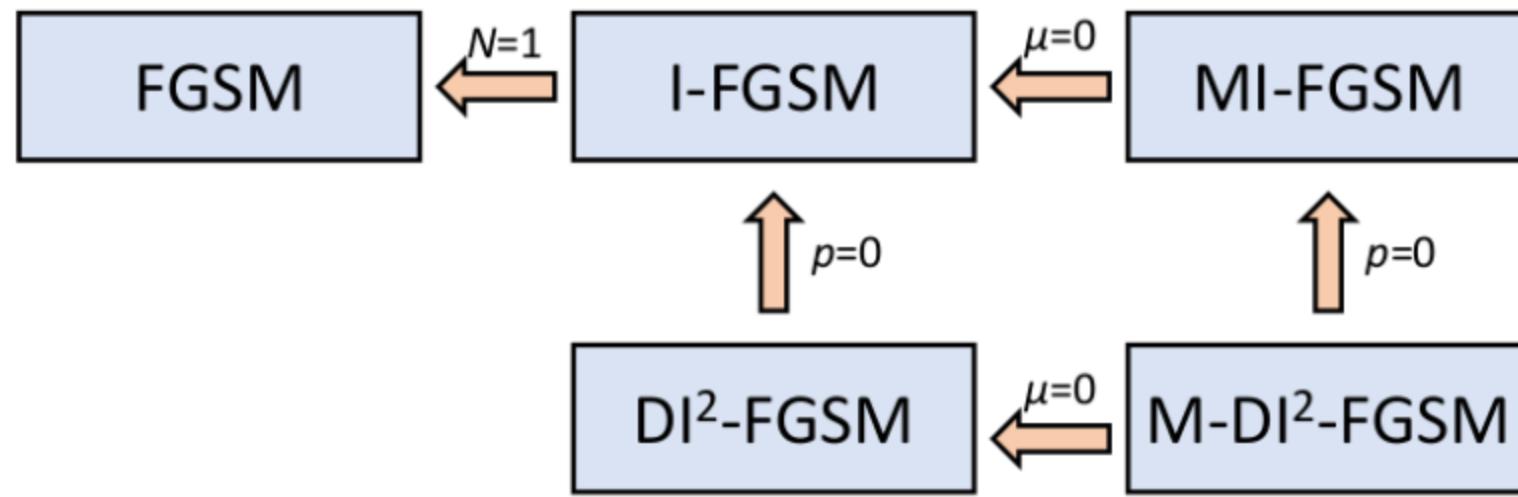
KNN (k=7)



Transfer attacks > Ensemble



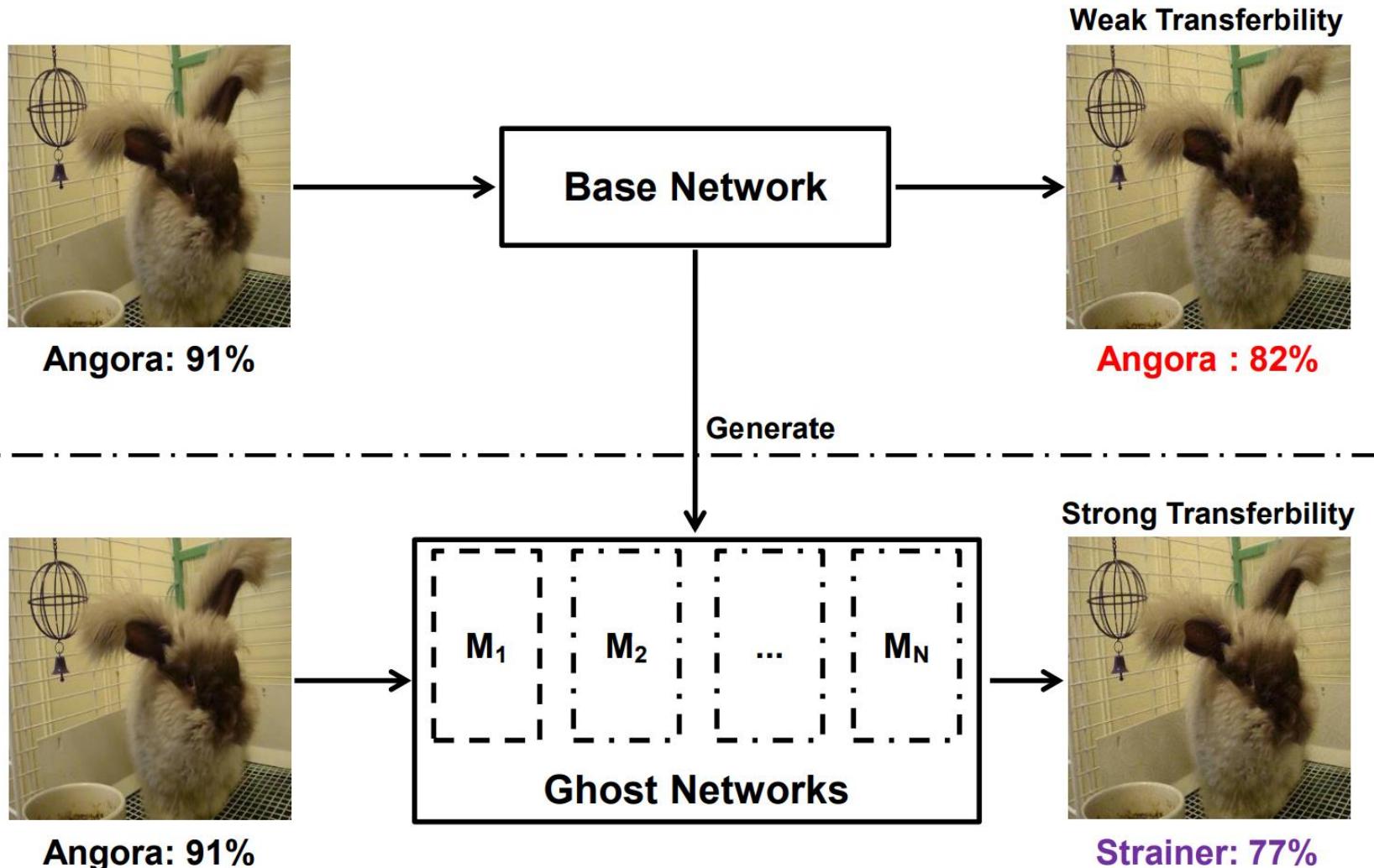
Transfer attacks > MI/DI FGSM



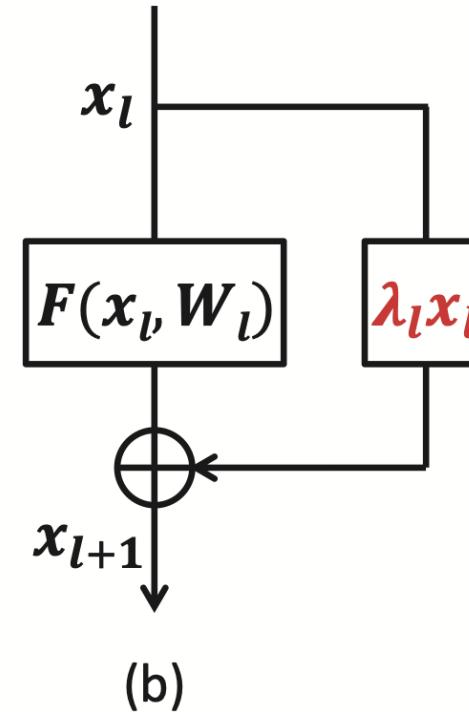
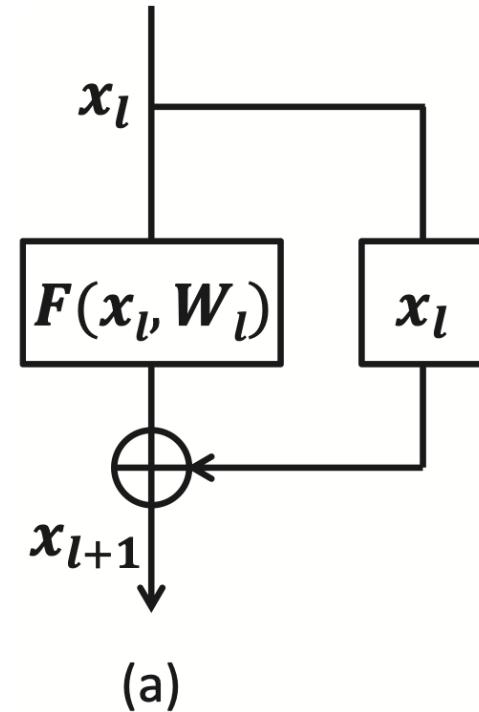
DI: Input Diversity

MI: Iterative Momentum

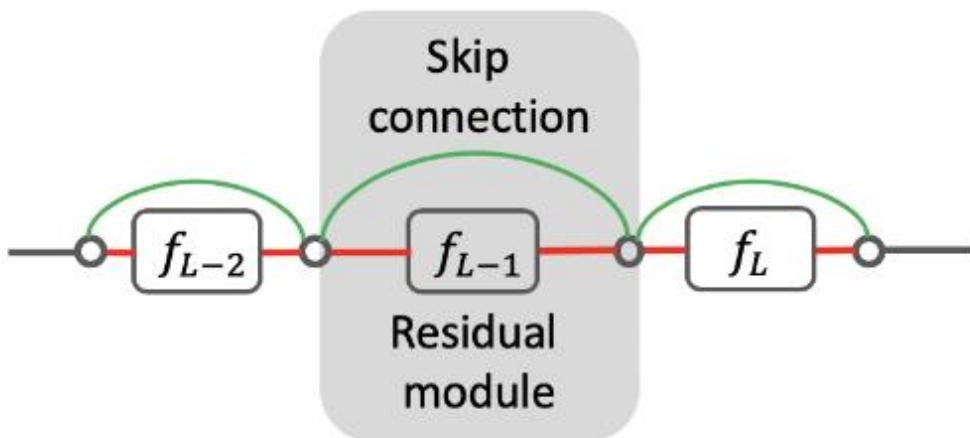
Transfer attacks > Ghost Networks



Transfer attacks > Ghost Networks



Transfer attacks > SGM



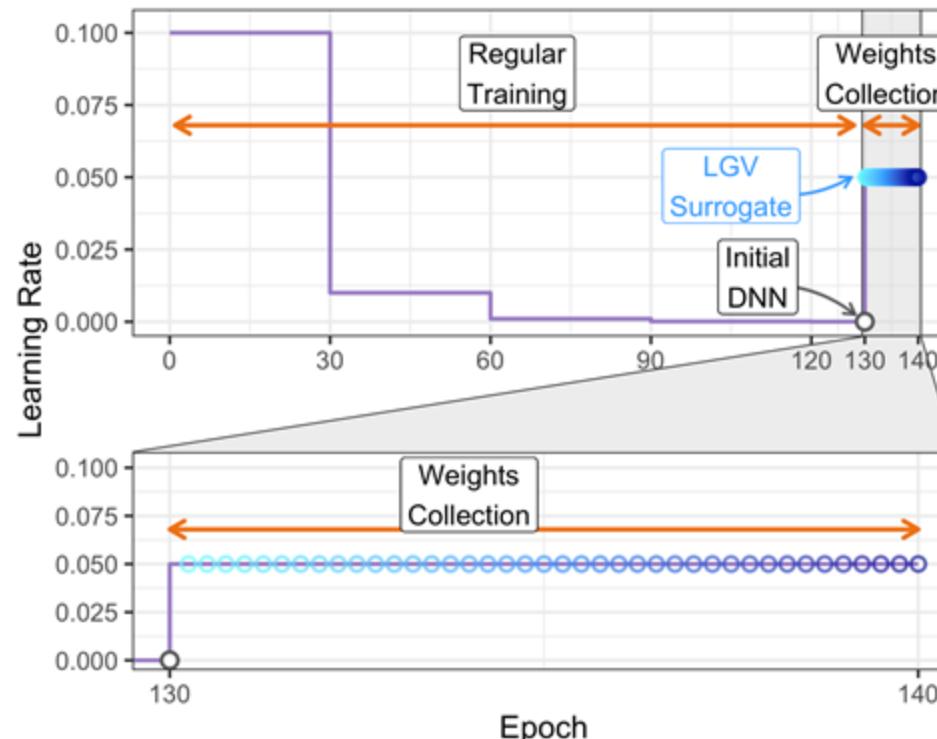
"Using more gradients from the skip connections rather than the residual modules, allows one to craft adversarial examples with high transferability"

	RN18	RN34	RN50	RN101	RN152	DN121	DN169	DN201
PGD	23.23 ± 0.69	24.38 ± 0.41	22.80 ± 0.55	22.98 ± 0.83	26.56 ± 0.75	30.71 ± 0.60	30.90 ± 0.31	36.01 ± 0.59
SGM	28.92 ± 0.45	43.43 ± 0.32	36.71 ± 0.55	38.38 ± 0.53	44.84 ± 0.14	57.38 ± 0.14	60.45 ± 0.42	65.48 ± 0.23

Transfer attacks > LGV

Phase 1

Collect models during a few additional epochs with a **high learning rate**



Algorithm 1 LGV Weights Collection

Input: n_{epochs} number of epochs, K number of weights, η learning rate, γ momentum, w_0 pretrained weights, \mathcal{D} training dataset

Output: (w_1, \dots, w_K) LGV weights

- 1: $w \leftarrow w_0$ \triangleright Start from a regularly trained DNN
 - 2: **for** $i \leftarrow 1$ **to** K **do**
 - 3: $w \leftarrow \text{SGD}(w, \eta, \gamma, \mathcal{D}, \frac{n_{\text{epochs}}}{K})$ \triangleright Perform $\frac{n_{\text{epochs}}}{K}$ of an epoch of SGD with η learning rate and γ momentum on \mathcal{D}
 - 4: $w_i \leftarrow w$
 - 5: **end for**
-

Transfer attacks > LGV

Phase 2

Attack one random collected model per iteration

Algorithm 2 I-FGSM Attack on LGV

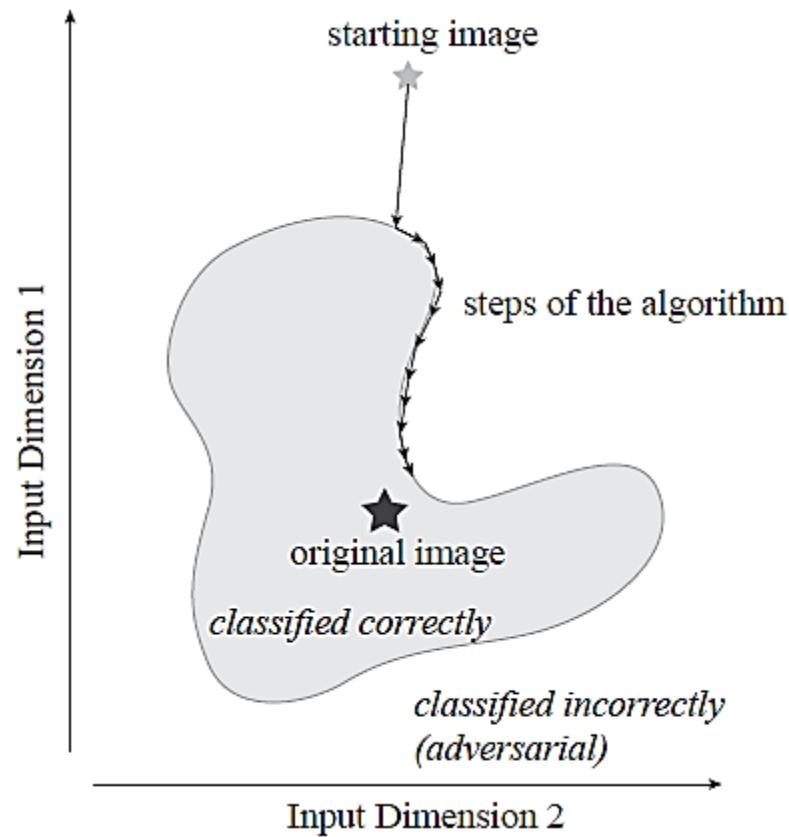
Input: (x, y) natural example, (w_1, \dots, w_K)
LGV weights, n_{iter} number of iterations, ε
 p -norm perturbation, α step-size
Output: x_{adv} adversarial example

- 1: Shuffle (w_1, \dots, w_K) \triangleright Shuffle weights
- 2: $x_{\text{adv}} \leftarrow x$
- 3: **for** $i \leftarrow 1$ **to** n_{iter} **do**
- 4: $x_{\text{adv}} \leftarrow x_{\text{adv}} + \alpha \nabla_x \mathcal{L}(x_{\text{adv}}; y, w_{i \bmod K})$
 \triangleright Compute the input gradient of the loss of
a randomly picked LGV model
- 5: $x_{\text{adv}} \leftarrow \text{project}(x_{\text{adv}}, B_\varepsilon[x])$ \triangleright Project
in the p -norm ball centred on x of ε radius
- 6: $x_{\text{adv}} \leftarrow \text{clip}(x_{\text{adv}}, 0, 1)$ \triangleright Clip to pixel
range values
- 7: **end for**

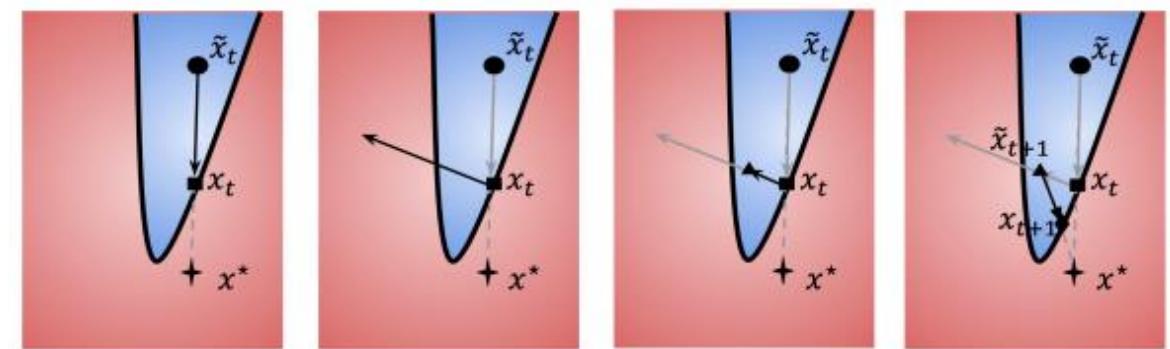
Transfer attacks

Surrogate	Target							
	RN50	RN152	RNX50	WRN50	DN201	VGG19	IncV1	IncV3
Baselines (1 DNN)								
1 DNN	45.3 \pm 2.4	29.6 \pm 0.9	28.8 \pm 0.2	31.5 \pm 1.6	17.5 \pm 0.6	16.6 \pm 0.9	10.4 \pm 0.5	5.3 \pm 1.0
MI	53.0 \pm 2.2	36.3 \pm 1.5	34.7 \pm 0.4	38.1 \pm 2.0	22.0 \pm 0.1	21.1 \pm 0.3	13.9 \pm 0.4	7.3 \pm 0.8
GN	63.9 \pm 2.4	43.8 \pm 2.4	43.3 \pm 1.3	47.4 \pm 0.9	24.8 \pm 0.3	24.1 \pm 1.0	14.6 \pm 0.3	6.8 \pm 1.2
GN+MI	68.4 \pm 2.3	49.3 \pm 2.5	47.9 \pm 1.2	52.1 \pm 1.7	28.4 \pm 0.8	28.0 \pm 0.7	17.5 \pm 0.5	8.7 \pm 0.5
DI	75.0 \pm 0.2	56.4 \pm 1.9	59.6 \pm 1.5	61.6 \pm 2.4	41.6 \pm 1.1	39.7 \pm 0.9	27.7 \pm 1.0	15.2 \pm 1.0
DI+MI	81.2 \pm 0.3	63.8 \pm 1.9	67.6 \pm 0.9	68.9 \pm 1.5	49.3 \pm 0.7	46.7 \pm 0.4	33.0 \pm 1.0	19.4 \pm 0.9
SGM	64.4 \pm 0.8	49.1 \pm 3.1	48.9 \pm 0.6	51.7 \pm 2.8	30.7 \pm 0.9	33.6 \pm 1.3	22.5 \pm 1.5	10.7 \pm 0.9
SGM+MI	66.0 \pm 0.6	51.3 \pm 3.5	50.9 \pm 0.9	54.3 \pm 2.3	32.5 \pm 1.3	35.8 \pm 0.7	24.1 \pm 1.0	12.1 \pm 1.2
SGM+DI	76.8 \pm 0.5	62.3 \pm 2.7	63.6 \pm 1.7	65.3 \pm 1.4	45.5 \pm 0.9	49.9 \pm 0.8	36.0 \pm 0.7	19.2 \pm 1.7
SGM+DI+MI	80.9 \pm 0.7	66.9 \pm 2.5	68.7 \pm 1.2	70.0 \pm 1.7	50.9 \pm 0.6	56.0 \pm 1.4	42.1 \pm 1.4	23.6 \pm 1.6
Our techniques								
RD	60.6 \pm 1.5	40.5 \pm 3.0	39.9 \pm 0.2	44.4 \pm 3.2	22.9 \pm 0.8	22.7 \pm 0.5	13.9 \pm 0.2	6.6 \pm 0.7
LGV-SWA	84.9 \pm 1.2	63.9 \pm 3.7	62.1 \pm 0.4	61.1 \pm 2.9	44.2 \pm 0.4	42.4 \pm 1.3	31.5 \pm 0.8	12.2 \pm 0.8
LGV-SWA+RD	90.2 \pm 0.5	71.7 \pm 3.4	69.9 \pm 1.2	69.1 \pm 3.3	49.9 \pm 1.0	47.4 \pm 2.0	34.9 \pm 0.3	13.5 \pm 0.9
LGV (ours)	<u>95.4\pm0.1</u>	<u>85.5\pm2.3</u>	<u>83.7\pm1.2</u>	<u>82.1\pm2.4</u>	<u>69.3\pm1.0</u>	<u>67.8\pm1.2</u>	<u>58.1\pm0.8</u>	<u>25.3\pm1.9</u>

Query attacks



Boundary attack



HopSkipJump attack

Query attacks

One pixel attack



Cup(16.48%)
Soup Bowl(16.74%)



Bassinet(16.59%)
Paper Towel(16.21%)

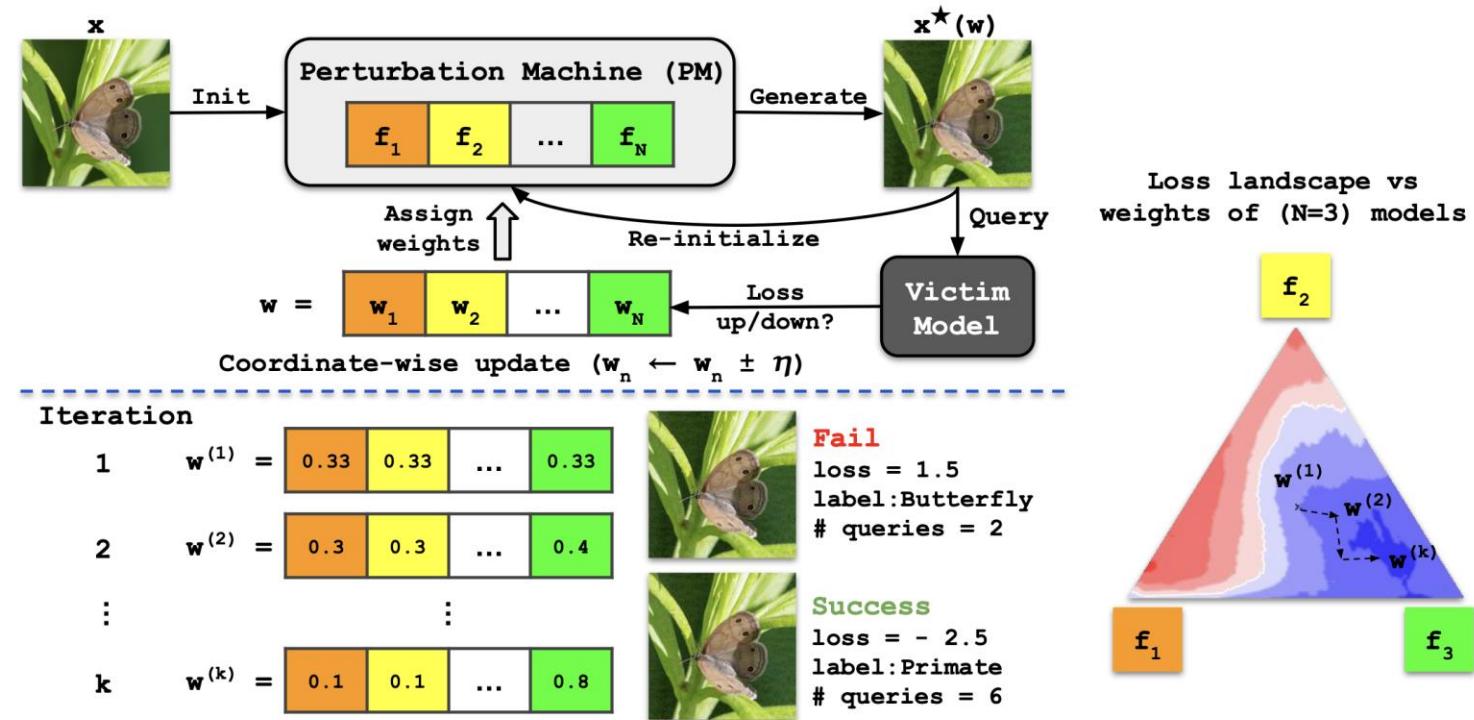


Teapot(24.99%)
Joystick(37.39%)



Hamster(35.79%)
Nipple(42.36%)

Query attacks > BASES



Query attacks > Square

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c ,
 l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w, c, \hat{x}, x)$  (see Alg. 2 and 3 for the sampling distributions)
5    $\hat{x}_{\text{new}} \leftarrow$  Project  $\hat{x} + \delta$  onto  $\{z \in \mathbb{R}^d : \|z - x\|_p \leq \epsilon\} \cap [0, 1]^d$ 
6    $l_{\text{new}} \leftarrow L(f(\hat{x}_{\text{new}}), y)$ 
7   if  $l_{\text{new}} < l^*$  then  $\hat{x} \leftarrow \hat{x}_{\text{new}}$ ,  $l^* \leftarrow l_{\text{new}}$  ;
8    $i \leftarrow i + 1$ 
9 end

```

Query attacks > Square

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

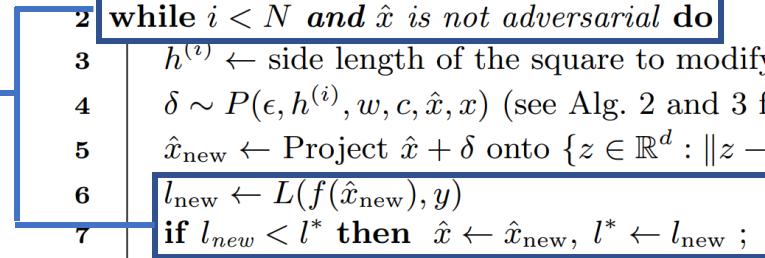
Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(v)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w, c, \hat{x}, x)$  (see Alg. 2 and 3 for the sampling distributions)
5    $\hat{x}_{new} \leftarrow$  Project  $\hat{x} + \delta$  onto  $\{z \in \mathbb{R}^d : \|z - x\|_p \leq \epsilon\} \cap [0, 1]^d$ 
6    $l_{new} \leftarrow L(f(\hat{x}_{new}), y)$ 
7   if  $l_{new} < l^*$  then  $\hat{x} \leftarrow \hat{x}_{new}$ ,  $l^* \leftarrow l_{new}$  ;
8    $i \leftarrow i + 1$ 
9 end

```

random search procedure



Query attacks > Square

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w, c, \hat{x}, x)$  (see Alg. 2 and 3 for the sampling distributions)
5    $\hat{x}_{\text{new}} \leftarrow$  Project  $\hat{x} + \delta$  onto  $\{z \in \mathbb{R}^d : \|z - x\|_p \leq \epsilon\} \cap [0, 1]^d$ 
6    $l_{\text{new}} \leftarrow L(f(\hat{x}_{\text{new}}), y)$ 
7   if  $l_{\text{new}} < l^*$  then  $\hat{x} \leftarrow \hat{x}_{\text{new}}$ ,  $l^* \leftarrow l_{\text{new}}$  ;
8    $i \leftarrow i + 1$ 
9 end

```

Sample δ from a specific distribution

l_∞ and l_2 perturbations use different distributions

Query attacks > Square

Algorithm 1: The Square Attack via random search

Initial p is a hyperparameter.
The value of p is halved after a certain iterations.

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w, c, \hat{x}, x)$  (see Alg. 2 and 3 for the sampling distributions)
5    $\hat{x}_{\text{new}} \leftarrow$  Project  $\hat{x} + \delta$  onto  $\{z \in \mathbb{R}^d : \|z - x\|_p \leq \epsilon\} \cap [0, 1]^d$ 
6    $l_{\text{new}} \leftarrow L(f(\hat{x}_{\text{new}}), y)$ 
7   if  $l_{\text{new}} < l^*$  then  $\hat{x} \leftarrow \hat{x}_{\text{new}}$ ,  $l^* \leftarrow l_{\text{new}}$  ;
8    $i \leftarrow i + 1$ 
9 end

```

$$h = \text{closest integer to } \sqrt{p \cdot w^2}$$

Query attacks > Square

Making sure the new adversarial example satisfy the constraints.

Actually, the sampling distribution ensured that before clipping to $[0,1]^d$

$$\|z - x\|_p = \epsilon$$

so that the perturbation budget is used almost maximally

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w, c, \hat{x}, x)$  (see Alg. 2 and 3 for the sampling distributions)
5    $\hat{x}_{\text{new}} \leftarrow$  Project  $\hat{x} + \delta$  onto  $\{z \in \mathbb{R}^d : \|z - x\|_p \leq \epsilon\} \cap [0, 1]^d$ 
6    $l_{\text{new}} \leftarrow L(f(\hat{x}_{\text{new}}), y)$ 
7   if  $l_{\text{new}} < l^*$  then  $\hat{x} \leftarrow \hat{x}_{\text{new}}$ ,  $l^* \leftarrow l_{\text{new}}$  ;
8    $i \leftarrow i + 1$ 
9 end

```

Query attacks > Square

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w,$ 
5    $\hat{x}_{new} \leftarrow$  Project
6    $l_{new} \leftarrow L(f(\hat{x}_{ne}$ 
7   if  $l_{new} < l^*$  the
8    $i \leftarrow i + 1$ 
9 end

```

Sample the position of the square update

Algorithm 2: Sampling distribution P for l_∞ -norm

Input: maximal norm ϵ , window size h , image size w , color channels c

Output: New update δ

```

1  $\delta \leftarrow$  array of zeros of size  $w \times w \times c$ 
2 sample uniformly
    $r, s \in \{0, \dots, w - h\} \subset \mathbb{N}$ 
3 for  $i = 1, \dots, c$  do
4    $\rho \leftarrow Uniform(\{-2\epsilon, 2\epsilon\})$ 
5    $\delta_{r+1:r+h, s+1:s+h, i} \leftarrow \rho \cdot \mathbb{1}_{h \times h}$ 
6 end

```

Query attacks > Square

Algorithm 1: The Square Attack via random search

Input: classifier f , point $x \in \mathbb{R}^d$, image size w , number of color channels c , l_p -radius ϵ , label $y \in \{1, \dots, K\}$, number of iterations N

Output: approximate minimizer $\hat{x} \in \mathbb{R}^d$ of the problem stated in Eq. (1)

```

1  $\hat{x} \leftarrow init(x)$ ,  $l^* \leftarrow L(f(x), y)$ ,  $i \leftarrow 1$ 
2 while  $i < N$  and  $\hat{x}$  is not adversarial do
3    $h^{(i)} \leftarrow$  side length of the square to modify (according to some schedule)
4    $\delta \sim P(\epsilon, h^{(i)}, w,$ 
5    $\hat{x}_{new} \leftarrow$  Project
6    $l_{new} \leftarrow L(f(\hat{x}_{ne}$ 
7   if  $l_{new} < l^*$  the
8      $i \leftarrow i + 1$ 
9 end

```

Use the same update
for each channel:
To visualize it in a RGB
image, there are only
 2^3 possible colors

Algorithm 2: Sampling distribution
 P for l_∞ -norm

Input: maximal norm ϵ , window
size h , image size w , color
channels c

Output: New update δ

```

1  $\delta \leftarrow$  array of zeros of size  $w \times w \times c$ 
2 sample uniformly
    $r, s \in \{0, \dots, w - h\} \subset \mathbb{N}$ 
3 for  $i = 1, \dots, c$  do
4    $\rho \leftarrow Uniform(\{-2\epsilon, 2\epsilon\})$ 
5    $\delta_{r+1:r+h, s+1:s+h, i} \leftarrow \rho \cdot \mathbb{1}_{h \times h}$ 
6 end

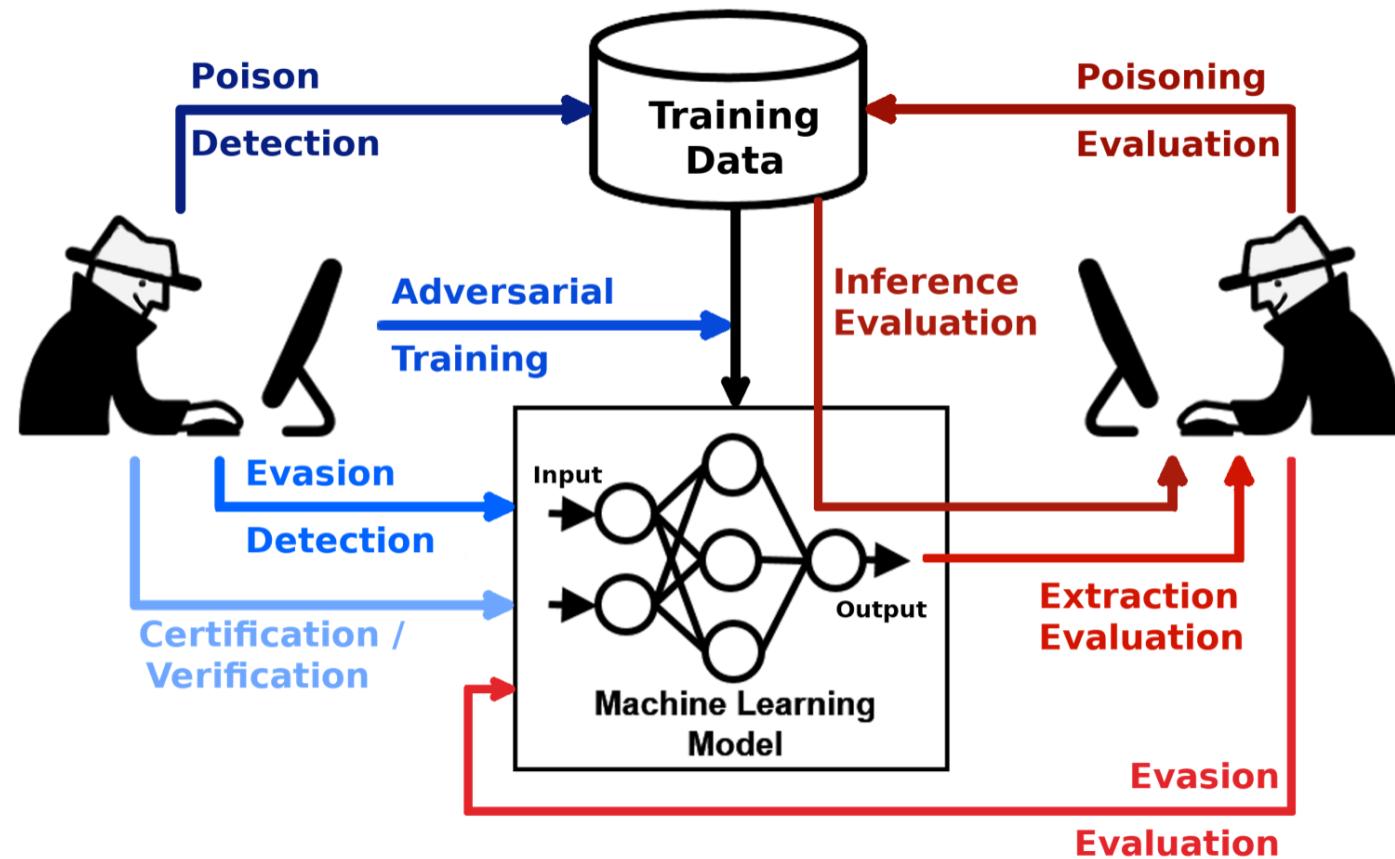
```

Defense against Evasion attacks

Part IV

A two-player game

- Attackers aim to produce strong adversarial examples that evade a model with high confidence while requiring only a small perturbation.
- Defenders aim to produce models that are robust to adversarial examples (i.e., the models don't have adversarial examples, or the adversaries cannot find them easily).



Main Defense strategies

- Adversarial training
- Randomized discretization
- Adversarial Robust Distillation
- Robust optimization (certified defenses)

Adversarial (Re)Training

Dynamic augmentation of the training with adversarial examples.

Two possible strategies

Train from scratch with
both original and
adversarial examples at
each batch



Adversarial Training

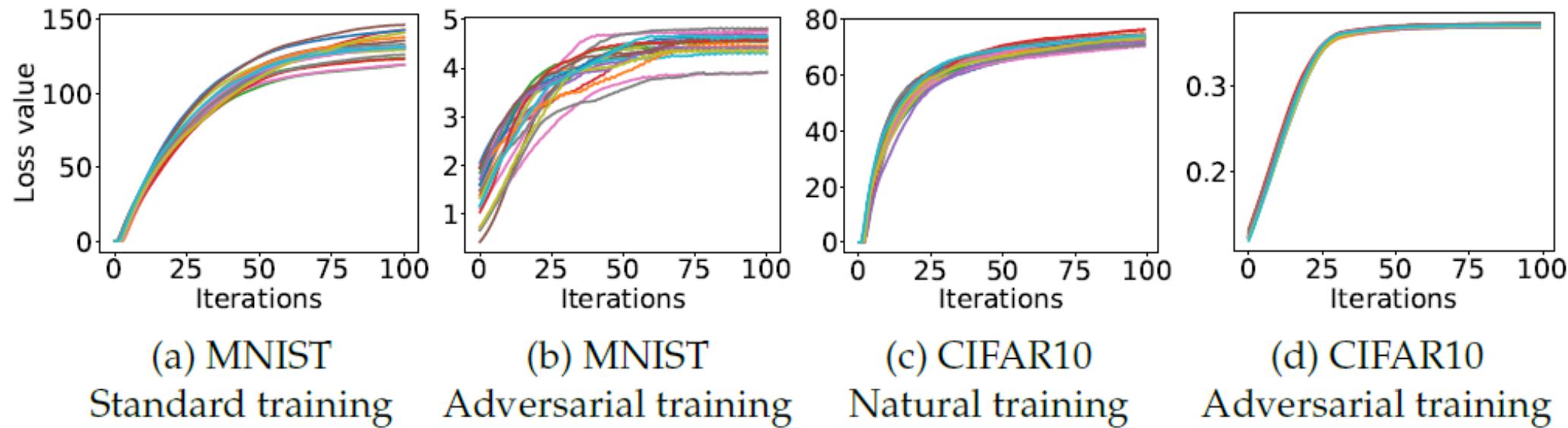
Train with original then
fine-tune with
adversarial examples



Adversarial Fine-tuning

Adversarial Training

When Adversarial examples are generated with PGD, we have « Madry Adversarial Training »



Adversarial Training

Adversarial training is assumed to be more expensive than traditional training --->
"FGSM" when combined with random initialization is as effective as PGD-based training

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
```

Adversarial Training Tricks

Adversarial training tricks play an important role in adversarial training

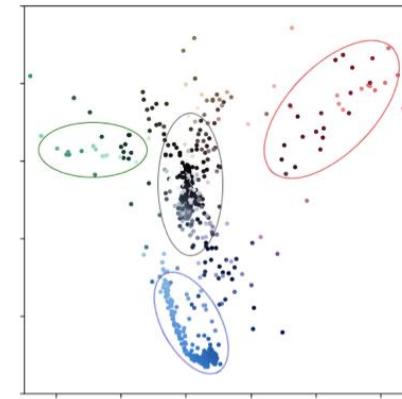
Category	Method
Data Augmentation	Mixup, RandAugment
Regularization	Weight Decay, Label Smoothing
Weight Averaging	EMA
Pre-training	21K-pre-training, SimMIM, CLIP

Randomized discretization

Inject random Gaussian noise, discretize each pixel, and then feed the result into any pre-trained classifier



(a)



(b)



(c)

Randomized discretization

Inject random Gaussian noise, discretize each pixel, and then feed the result into any pre-trained classifier

Gaussian randomization

$$\tilde{x}_i = x_i + w_i \quad \text{where } w_i \sim \mathcal{N}(0, \sigma^2 I).$$

Randomized discretization: first adds Gaussian noise to each pixel and then discretizes each based on a finite number of cluster centers

$$\tilde{x}_i = r(x_i | \mathbf{c}) := \operatorname{argmin}_{c \in \{c_1, \dots, c_k\}} \|x_i + w_i - c\|_2$$

Adversarial Robust Distillation

Knowledge Distillation: involves training a smaller neural network (student) to replicate the outputs of a larger network (teacher) to transfers the teacher's knowledge to the student

$$\text{CKD}(\mathbf{x}, y) = (1 - \lambda)\text{CE}(f_S(\mathbf{x}), y) + \lambda\text{KL}(f_S(\mathbf{x}), f_T(\mathbf{x}))$$

Adversarial Robust Distillation: The previous loss is effective in transferring clean performance, it does not seem to be able to transfer robustness against adversarial examples.

$$\text{ARD}(\mathbf{x}, y) = (1 - \lambda)\text{CE}(f_S(\mathbf{x}), y) + \lambda\text{KL}(f_S(\mathbf{x}'), f_T(\mathbf{x})).$$

\mathbf{x}' : adversarial example, CE: cross entropy, KL: Kullback–Leibler divergence

Adversarial Robust Distillation

Robust Soft Label Adversarial Distillation: it avoids the use of clean labels and focuses only on matching the models' outputs

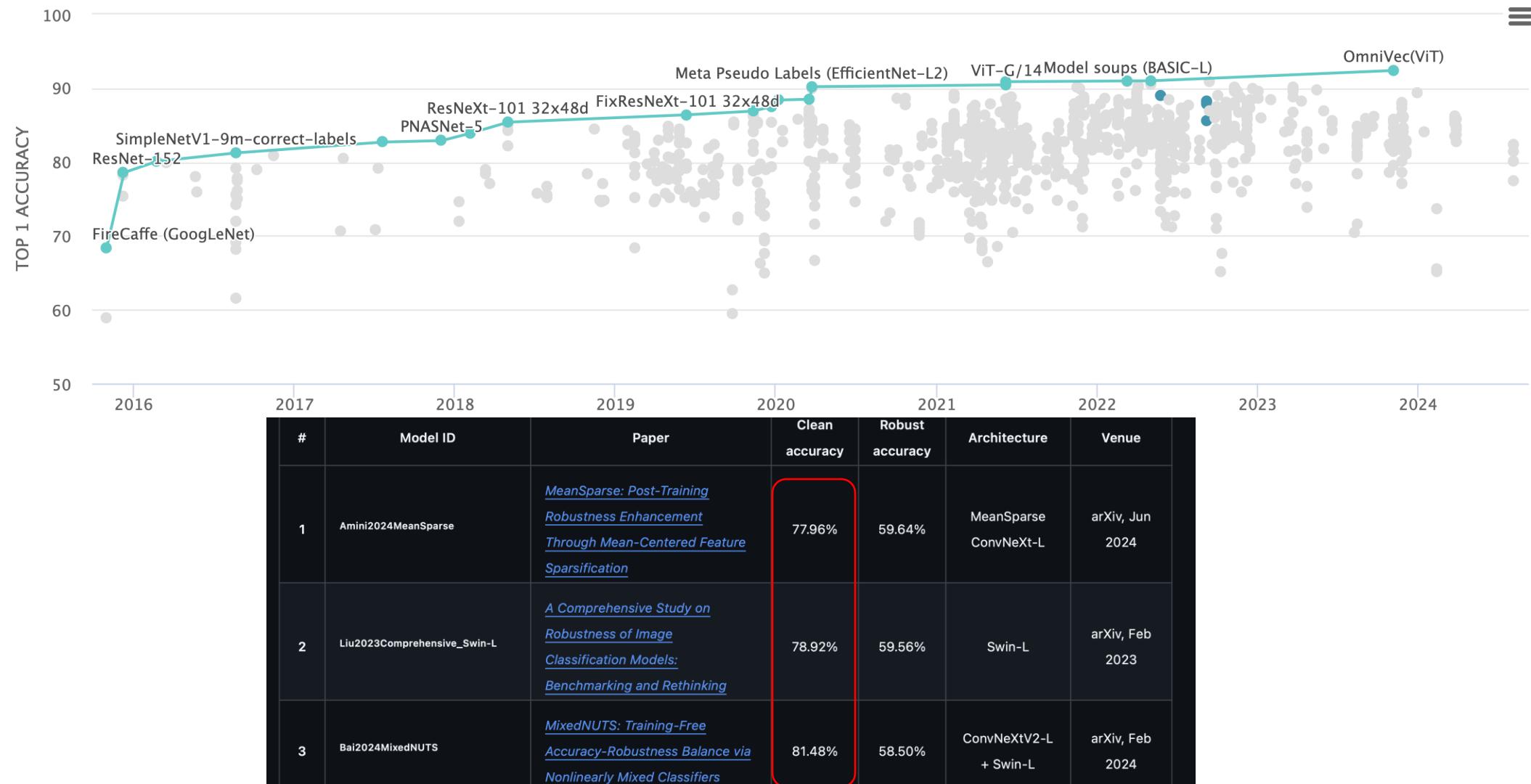
$$\text{ARD}(\mathbf{x}, y) = (1 - \lambda)\text{CE}(f_S(\mathbf{x}), y) + \lambda\text{KL}(f_S(\mathbf{x}'), f_T(\mathbf{x})).$$



$$\text{RSLAD}(\mathbf{x}) = (1 - \lambda)\text{KL}(f_S(\mathbf{x}), f_T(\mathbf{x})) + \lambda\text{KL}(f_S(\mathbf{x}'), f_T(\mathbf{x})).$$

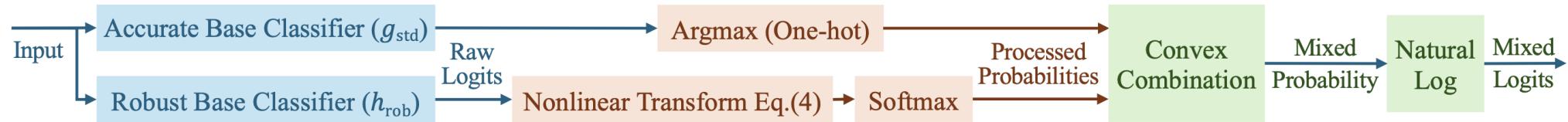
\mathbf{x}' : adversarial example

Clean accuracy vs Robust accuracy



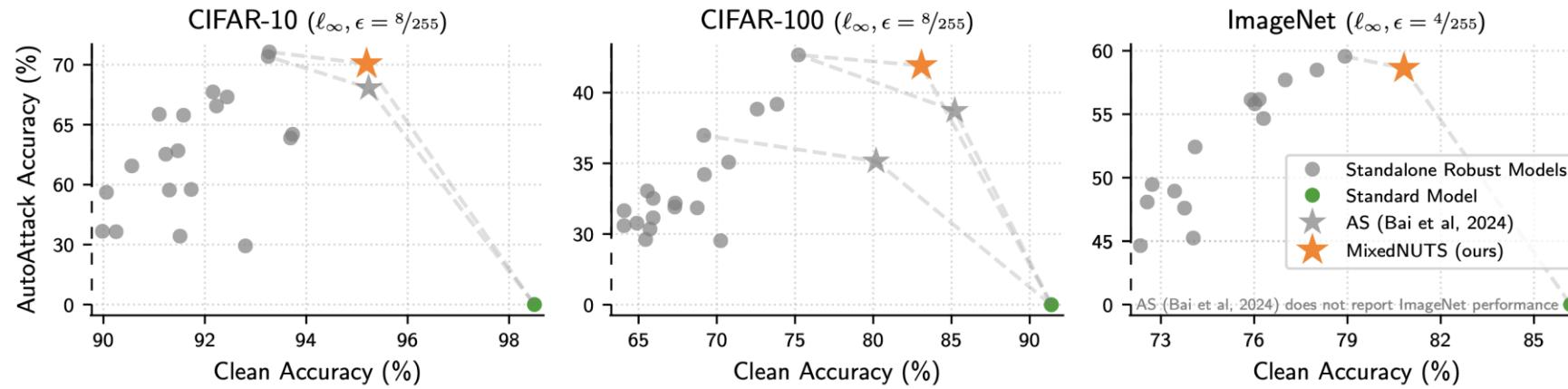
Clean accuracy vs Robust accuracy

Adversarial robustness often comes at the cost of degraded accuracy



Clean accuracy vs Robust accuracy

Adversarial robustness often comes at the cost of degraded accuracy



Model certification

Existing robustness metrics are « Empirical »

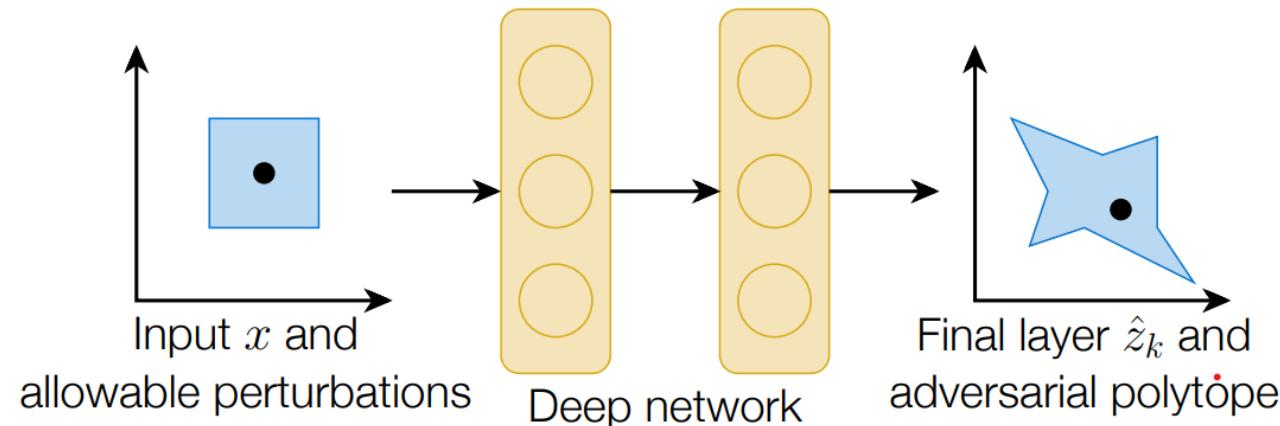
They can « show » where a model is weak
But never where a model is robust

Can we provide a bound of maximum loss after a
bounded attack ?

Model certification

Challenge

A regular (convex) adversarial region may be highly complex and concave after layers of neural network



Model certification

Many solutions, 2 will be covered in an upcoming (11/11/2024)

Randomized Smoothing

Train a base classifier trained with Gaussian corruption then smooth it

Activation Bound Relaxation

Relax the bounds of the reachable region so that it is convex and easier to work with

