# Security of AI/ML

Teddy Furon
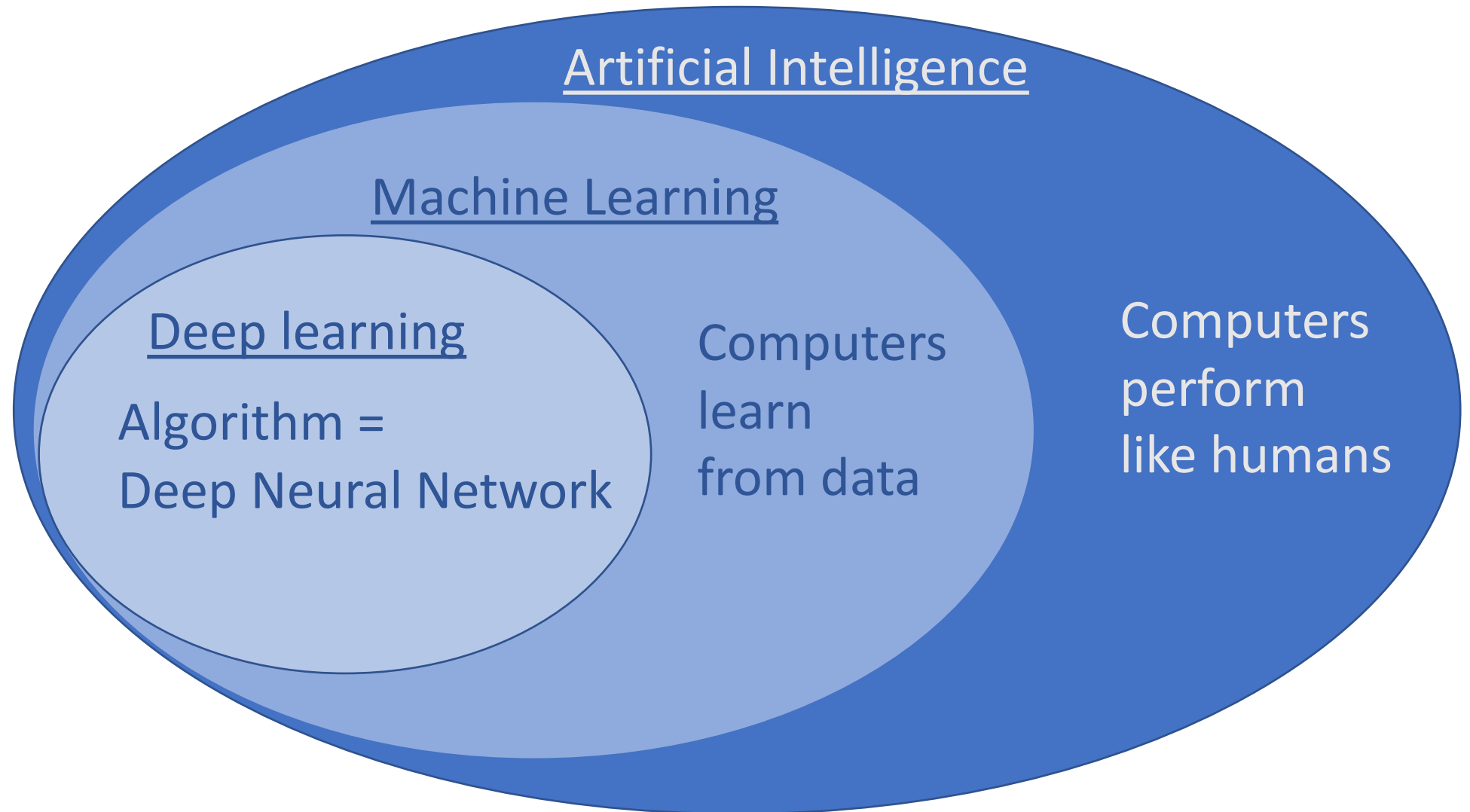
Inria Rennes

Summer School, Cyber in Normandy, Caen 2024

# Angles

- The type of AI?
    - Decision making AI
    - Generative AI
- Access to the model
    - White box
    - Black box (MLaaS, MLonChips)
- Security issues
    - Intrinsic vulnerabilities of the model
    - Malicious use of the model
- Security levels
    - Nothing is secure, nothing is insecure … to some extend
- Goals
    - Recommendations, defenses
    - Control, certification

# What kind of AI?



Artificial Intelligence

Machine Learning

Deep learning

Algorithm =
Deep Neural Network

Computers
learn
from data

Computers
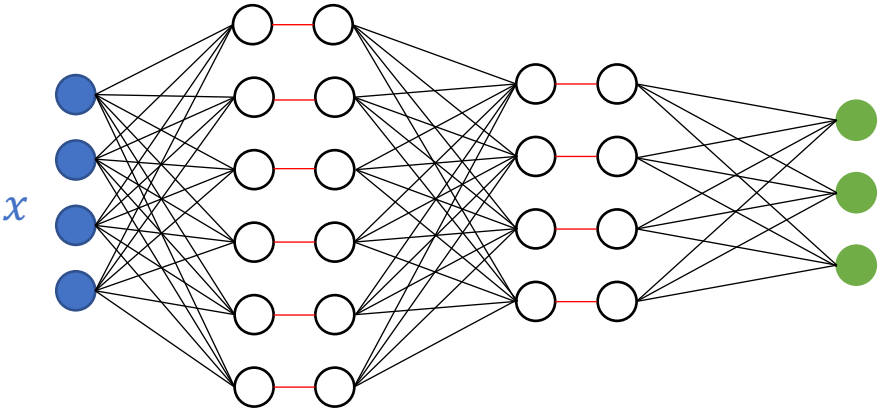perform
like humans

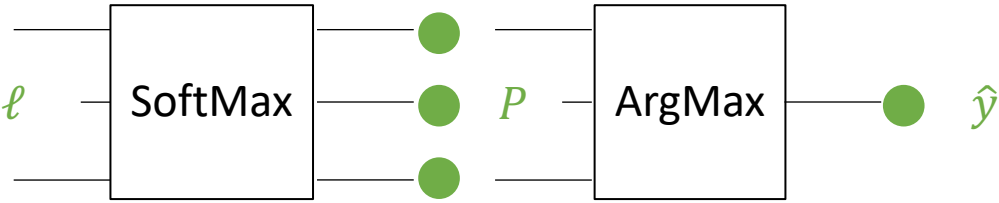# What kind of AI

1. A simple definition of Security of ML
2. The rocky horror picture show
3. Case studies
   - Local robustness
   - Adversarial examples
   - Fingerprinting
   - Watermarking
   - Backdoors

# Neural network classifiers

Linear + Non lin.    Linear + Non lin.    Classification

$\ell$    SoftMax    $P$    ArgMax    $\hat{y}$

$x$

Inputs
$x \in \mathbb{R}^d$

logits
$\ell \in \mathbb{R}^c$

"probabilities" - probits
$P \in \mathbb{S}^c$

predicted class
$\hat{y} \in [\![c]\!] = \{1, \ldots, c\}$

$$\ell = f(x; \theta) = \text{logits}$$
$$\ell = W_3 \sigma(W_2 \sigma(W_1 x + b_1) + b_2)$$

$$P[i] \propto e^{\ell[i]}$$

$$\sum_i P[i] = 1$$

Non lin. activation function

$y = \sigma(x)$

# DNN classifiers

- ## What is the output?
  - Logits, probits, predicted class
  - Black box

- ## Differentiable (almost everywhere)
  - 2 Gradients    $\nabla_\theta f(x; \theta) \in \mathbb{R}^{|\theta| \times c}$    $\nabla_x f(x; \theta) \in \mathbb{R}^{d \times c}$
  - Efficient
    - autodiff + backpropagation
    - Cost $\approx$ 2 times a forward pass
  - Training
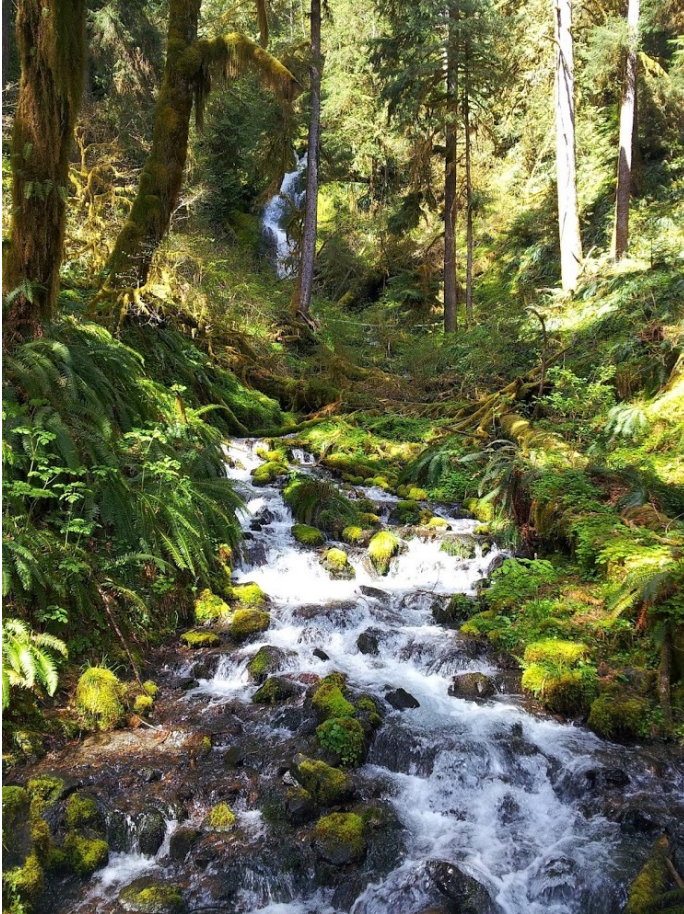    - SGD: $\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_\theta \text{Loss}(\text{SoftMax}(f(x_i; \theta)), y_i)$    Loss: $\mathbb{S}^c \times [\![c]\!] \longrightarrow \mathbb{R}$
  - Explicability
    - Deep dreams or GradCAM: visualisation of $\nabla_x f_i(x; \theta)$    $i \in [\![c]\!]$
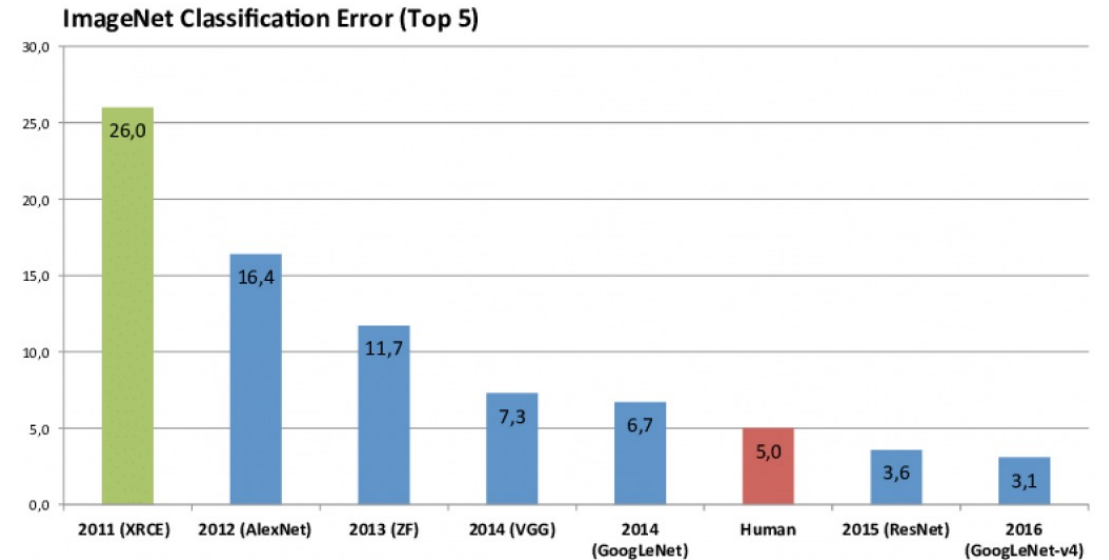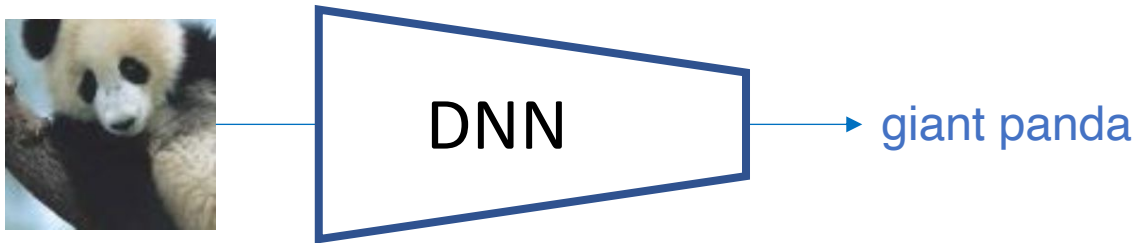
# Deep dreams



$$x_o, y_o = forest$$



$$x_o + \eta \cdot \nabla_x f_{forest}(x_o; \theta)$$

# ImageNet challenge: the iconic example of A.I.



DNN → giant panda

**ImageNet Classification Error (Top 5)**

| | |
|---|---|
| 2011 (XRCE) | 26,0 |
| 2012 (AlexNet) | 16,4 |
| 2013 (ZF) | 11,7 |
| 2014 (VGG) | 7,3 |
| 2014 (GoogLeNet) | 6,7 |
| Human | 5,0 |
| 2015 (ResNet) | 3,6 |
| 2016 (GoogLeNet-v4) | 3,1 |

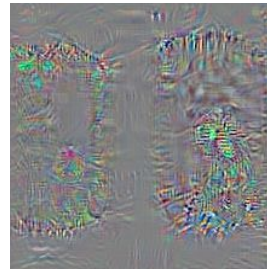## 2012: DNN AlexNet handily wins the top prize

- Krizhevsky, Sutskever, and Hinton (Univ. of Toronto)
- « *That moment is widely considered a turning point in the development of contemporary AI* »
- « *This dramatic quantitative improvement marked the start of an industry-wide artificial intelligence boom* »
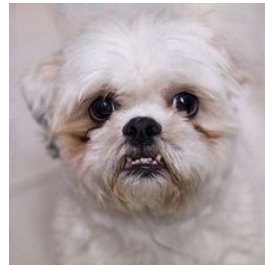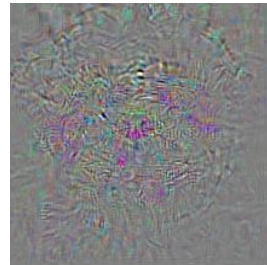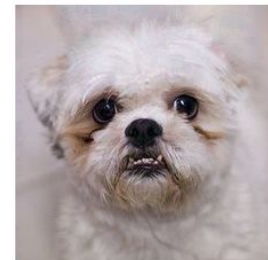
# The big failure

loudspeaker



pekinese

ostrich

school bus

$$x_o \quad + \epsilon * \quad \nabla_x f_{\text{ostrich}}(x_o;\ \theta)$$

*Intriguing properties of neural networks,* Szegedy, Goodfellow et al., 2014

# The big failure



giant panda     $x_o$     $+ \epsilon *$     $=$     $x_a$     gibbon

How can we call "Artificial Intelligence" algorithms so easily deluded!

*Explaining and harnessing adversarial examples*, Goodfellow et al., 2015

# 1- Definition of Security of ML

# False sense of security

Generalization ≠ Safety Robustness ≠ Security

- Generalization: To operate as expected on <u>unseen</u> data
  - Unseen but distributed like the training data

- Robustness: To operate as expected on <u>noisy</u> data
  - Unseen and almost distributed like the training data

- Security: To operate as expected on <u>purposely perturbed</u> data
  - Presence of an attacker

# ML to the bare bones

Testing data → Inference → Result

Model → Inference

Training data → Learning → Model

**Protection of 3 objects**
- **Training data**
- **Model**
- **Testing data**

# IT Security to the bare bones: C.I.A. Triad



onal Bureau of Standards
OCT 2 6 1977

**COMPUTER SCIENCE & TECHNOLOGY:**

# Audit and Evaluation of Computer Security

+ special public...

Proceedings of the NBS Invitational Workshop
held at Miami Beach, Florida, March 22-24, 1977

Edited by:

Zella G. Ruthberg

Institute for Computer Sciences and Technology
National Bureau of Standards
Washington, D. C. 20234

Robert G. McKenzie

General Accounting Office
Washington, D. C. 20548

Computer Security -- The protection of system data and resources from accidental and deliberate threats to confidentiality, integrity, and availability.

# Security of Machine Learning

Training data •

Model •                 **?**

Testing data •

• Confidentiality

• Integrity

• Availability

# Security of Machine Learning

Training data •

Model •   **?**

Testing data •

• Confidentiality

• Integrity

• Availability

# ML + IT Security – Confidentiality = Cryptology

- Testing data
  - Inference on encrypted data
  - Collaboration: Alice has <u>sensitive</u> testing data, Bob has a valuable model

- Training data
  - Learning from encrypted data
  - Collaboration: Alice has <u>sensitive</u> training data, Bob has the expertise in ML

Yes, we can!
  - Homomorphic Encryption: **CONCRETE**
    [Programmable Bootstrapping Enables Efficient Homomorphic Inference of DNN, Chillotti, CSCML'21]
  - Multi Party Coputation: **FALCON**
    [Honest-Majority Maliciously Secure Framework for Private DL, Wagh, PETS'21]

    TinyImageNet ( 64x64x3 = 12k - 200 classes ) + VGG16     = x 10,000 slower

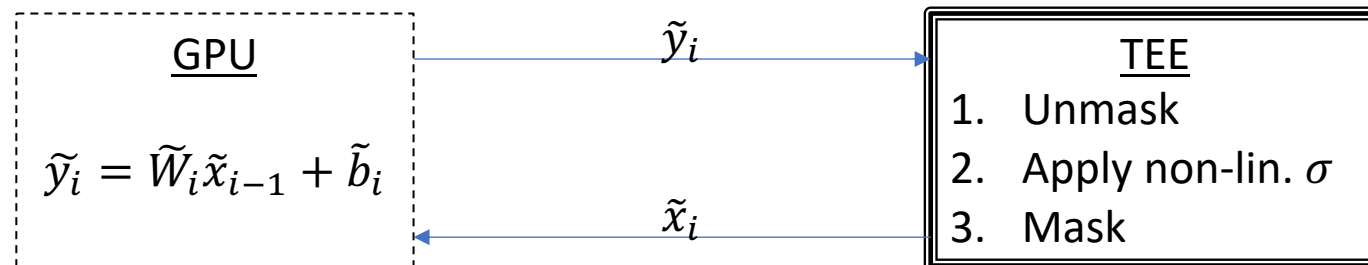  - Federated learning

MLaaS
Cloud computing

# ML + IT Security – Confidentiality = Cryptology

- Model
  - Model embedded on device
    - Civil: smartphones, smart speakers [Sonos-privacy]
    - Defense: AI embedded in armed vehicles / drones
  - Deep Neural Networks + GPU $\neq$ Code obfuscation

  - Communication protocol between GPU and SOC/TEE chips

    [ShadowNet: A secure and efficient system for on-device model inference, Sun, IEEE S&P 23]



GPU

$$\widetilde{y}_i = \widetilde{W}_i \tilde{x}_{i-1} + \tilde{b}_i$$

$\widetilde{y}_i$

$\tilde{x}_i$

TEE
1. Unmask
2. Apply non-lin. $\sigma$
3. Mask

New startup in town: Skyld!

# ML + IT Security – Confidentiality = Privacy

- Training data
  - Given a model, what can the attacker say about the training data?
  - Membership Inference Attack
    - [Bayes Optimal Strategies for Membership Inference, Sablayrolles, ICML'19]
  - Reconstruction of training data
    - [Extracting Training Data from Large Language Models, Carlini, Usenix'21]
  - Federated learning with privacy
    - [An Accurate, Scalable and Verifiable Protocol for Federated DP Averaging, Sabater, ML'22]

- Model (black box)
  - Model Identification / Fingerprinting       or       Model Extraction / Shadowing
    - [Stealing machine learning models via prediction APIs, Tramer, Usenix'16]

- Testing data
  - Restricted Inference / Data sanitization
    - [Learning Semi-Supervised Anonymized Representations by Mutual Information, Feutry, ICASSP'20]
    - [Differentially Private Speaker Anonymization, Shamsabadi, PETS'23]

# Security of Machine Learning

Training data •

Model •  **?**

Testing data •

• Confidentiality

• Integrity

• Availability

# ML + IT Security – Integrity

- Training data
  - Backdooring / Poisoning Attack

[Poisoning Attacks against Support Vector Machines, Biggio, ICML'12]

[A new backdoor attack in CNNs …, Barni, ICIP'19]

- Model
  - Backdooring / Trojaning

[TBT: Targeted Neural Network Attack with Bit Trojan, Rakin, CVPR 2020]

[Planting Undetectable Backdoors in Machine Learning Models, Goldwasser, arXiv'22]

- Testing data
  - Adversarial examples / Evasion attacks

# Security of Machine Learning

Training data

Model    **?**

Testing data

- Confidentiality

- Integrity

- Availability

# ML + IT Security – Availability

- Training data
  - ???

- Model
  - Deny of Service Attack against DNN
    [Sponge Examples: Energy-Latency Attacks on Neural Networks, Shumailov, Euro SP, 2021]

- Testing data
  - ???

# ML + Information Security: Traceability

- Training data
  - Radioactivity
    - Embed a watermark in a training set
    - Detect the watermark from a model learnt over this training set

      [Radioactive data: tracing through training, Sablayrolles, ICML'20]
      [Watermarking makes language models radioactive, Sander, arXiv'24]

- Model
  - Watermarking of a classifier

      [Entangled Watermarks as a Defense against Model Extraction, Jia, Usenix'21]
      [DNN Watermarking: Four Challenges and a Funeral, Barni, IHMMSEC'21]
  - Watermarking of generative AI (Text, Image, Audio)

      [Supervised GAN Watermarking for Intellectual Property Protection, Fei, arXiv'22]
      [Proactive Detection of Voice Cloning with Localized Watermarking, San Roman, arXiv'24]
      [The Stable Signature: Rooting Watermarks in Latent Diffusion Models, Fernandez, ICCV'23]

- Testing data
  - ???

# Security of Machine Learning



- 3 objects x 4 values - 1 = 11 scenarios
- 11 x types of data x types of learning framework x types of DNN

# 2- Where do we stand?

# Where do we stand?

1. The Rocky Horror Picture Show
   - Empirical Evidence of Attacks
   - Alarming, Threatening

2. Research work in the lab
   - Reproducibility
   - Empirical discovery of key factors
   - Theoretical explanations

3. Real life: Auditing, Advising
   - Run SotA attacks and see …

# Where do we stand? Adversarial examples



- Not reproducible

- Explanation (?):
  - adversarial examples = tensor of scalars ≠ tensor of integers

# Where do we stand? Adversarial examples

- Naïve defenses are not working
  - Gradient obfuscation

  *"Since all white-box attacks resort to the gradient of the neural network, just introduce a non-linearity to forbid its computation"*

  $$f = f_1 \circ f_2 \quad \rightarrow \quad f_Q = f_1 \circ Q \circ f_2$$

  - The attacker is not obliged to do so!

  [*Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples*, Athalaye, ICML 2018]
    - This paper circumvents 7 defenses proposed in ICLR 2018

# Where do we stand? Adversarial examples

- Proposal of best practices for evaluating attacks/defenses

  - [*On Evaluating Adversarial Robustness*, Carlini, arXiv 2019]

- Fear Nicholas Carlini (Google Deepmind)
  - [*Cutting through buggy adversarial example defenses: fixing 1 line of code breaks Sabre*, Carlini, arXiv 2024]
  - Significant flaws in *Sabre*, defense paper accepted at IEEE S&P 2024
  - Not following any of the best practices

# Where do we stand? Adversarial examples

- Consensus: Adversarial training is the only way to go (?)



ResNet50
Adv. Train [Madry, 18]
Adv. Train [Salman, 20]
Adv. Train [Salmon, 20]

# Where do we stand? Training data confidentiality

**Prefix**

`East Stroudsburg Stroudsburg...`

**GPT-2**

**Memorized text**

```
        Corporation Seabank Centre
        Marine Parade Southport
Peter W
            @        .           .com
+    7 5      40
Fax: +    7 5      0   0
```

Downloads last month
**229,839**

⚡ **Hosted inference API** ⓘ

📝 Text Generation    Examples ⌄

East Stroudsburg Stroudsburg (PA) Sunnyside
Syracuse Tusculum Tuskegee Tuxedo Union (NY)
USC Upstate UT Austin Utah Valley UT-Arlington
Utah Valley JC UT-Chattanooga UT-K|

Compute   ⌘+Enter    0.0

- Not reproducible
- Not explainable

# Where do we stand? Training data confidentiality

(a) Top 24 images reconstructed from a binary classifier trained on 50 CIFAR10 images



(b) Their corresponding nearest neighbours from the training-set of the model



- Strong theoretical limitations
  - Binary classification
  - Homogeneous neural networks (no biases, no residuals)
- Experimental evidence
  - On 3-layer MLPs

[Reconstructing Training Data from Trained Neural Networks, Haim, NeurIPS'22]

# Where do we stand? Training data confidentiality



- Clear impact of the overfitting
- Outliers in the training set are more easily discovered

[Label-Only Membership Inference Attacks, Choquette-Choo, ICML'21]

# Security of Machine Learning

- Study the Security of ML before applying ML to Security

- Simple definition
  - (Training d., Model, Testing d.)  x  (Confidentiality, Privacy, Integrity, Traceability)
  - Almost sound and almost complete

- Where do we stand?
  - In the lab!
  - In real life: "It depends"

- As a reader: adversarial reading of adversarial ML papers

- As a writer: be skeptical about your results
  - *"the first principle [of research] is that you must not fool yourself—and you are the easiest person to fool"*. R. Feynman
  - Switch your mindset: play the attacker/defender role

# 3- Case studies

# 3a- Robustness

Karim Tit et al.

*Efficient Statistical Assessment of Neural Network Corruption Robustness*, NeurIPS 2021

*Gradient-Informed Neural Network Statistical Robustness Estimation,* AISTATS 23

# Problem



$$f(x_o) = \begin{bmatrix} P[y = 1|x_o] \\ P[y = 2|x_o] \\ \cdots \\ P[y = C|x_o] \end{bmatrix}$$

$$f(x) = \begin{bmatrix} P[y = 1|x] \\ P[y = 2|x] \\ \cdots \\ P[y = C|x] \end{bmatrix}$$

+ uncertainties

Probits = "predicted" probabilities

# Problem
## Local certification in classification

- Consider $x_o \in \mathbb{R}^d$, well classified

$$\arg \max_i f_i(x_o) = panda$$

- Consider two regions
  - Input region: $\quad \mathcal{I} = \{ x \in \mathbb{R}^d \mid d(x, x_o) \leq \varepsilon \} \subset \mathbb{R}^d$
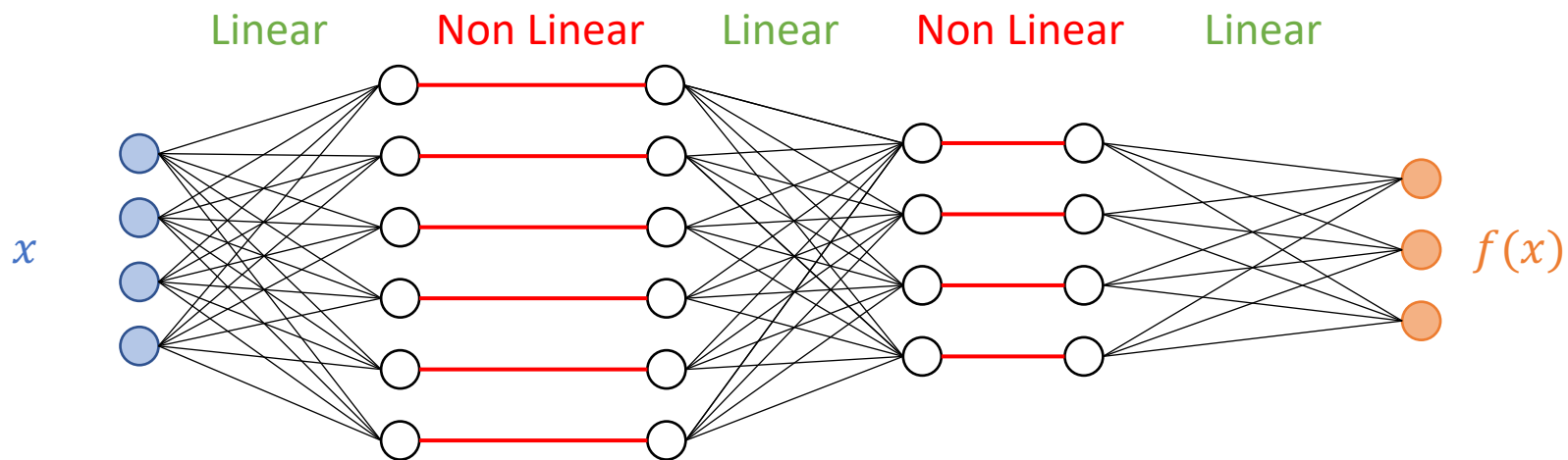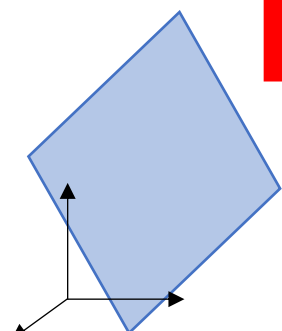  - Output region: $\quad \mathcal{O} = \{ f \in \mathbb{S}^c \mid \arg \max_i f_i = panda \} \subset \mathbb{R}^c$

# Certification

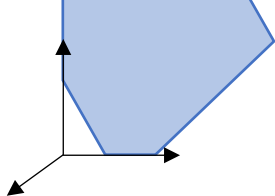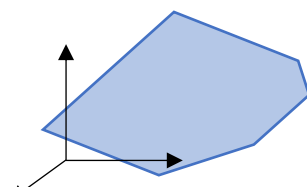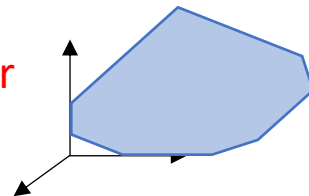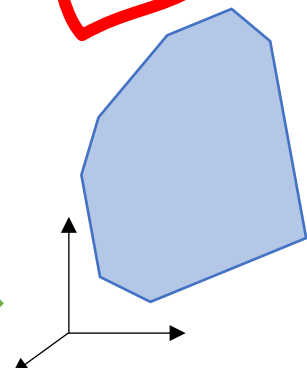# Formal proof



Linear  Non Linear  Linear  Non Linear  Linear

$y = \sigma(x)$

Non linear activation function

$x$

$f(x)$

$\mathbb{R}^d$

$x_o$

$\varepsilon$

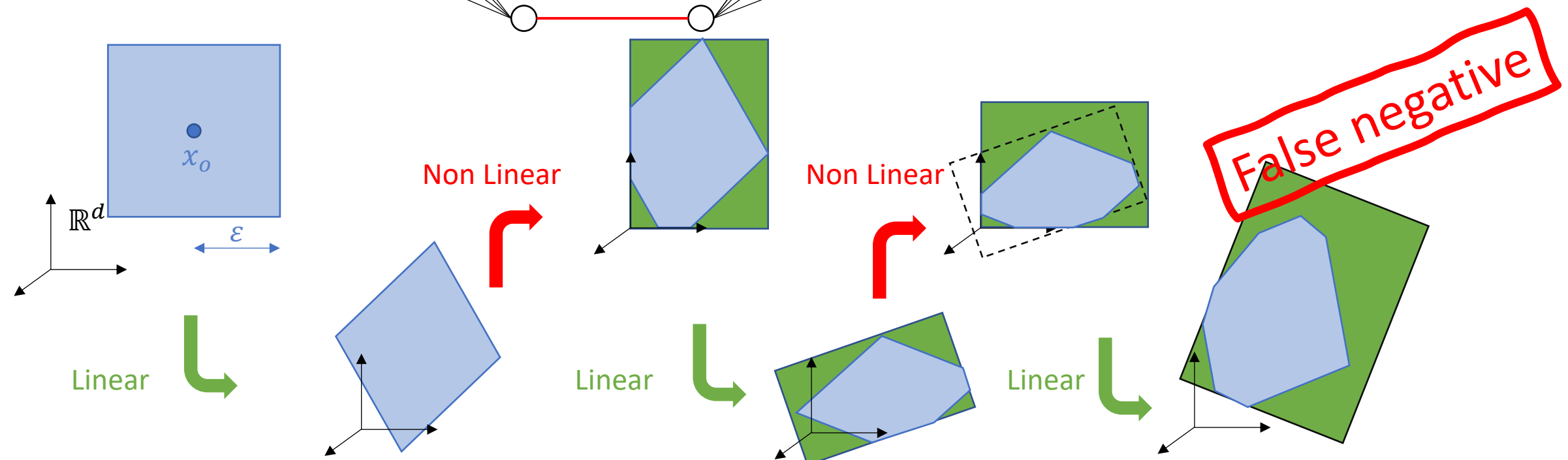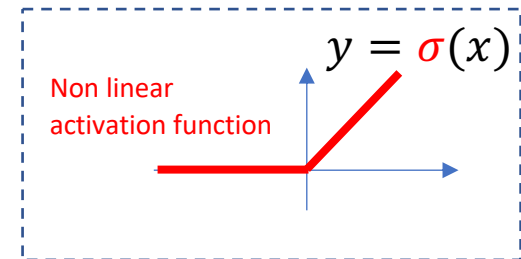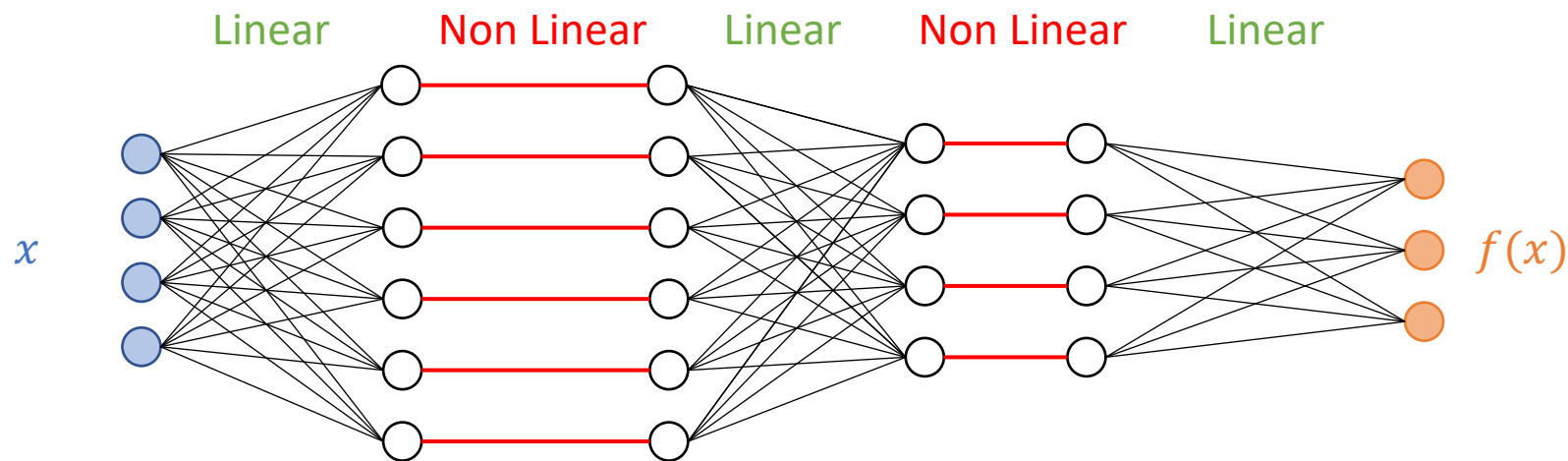Non Linear

Non Linear

Linear

Linear

Linear

NP hard problem

# Formal proof with relaxation

# Formal proof

- Sound and complete (but not scalable)
  - ReLUplex, Katz *et al.*, Computer Aided Verification 2017

- Relaxation (not complete) but more scalable
  - Crown, Zhang *et al.*, NeurIPS 2018
  - CNN-CERT, Weng *et al.*, AAAI 2019
  - DeepPoly, Singh *et al.*, Programming Languages, 2019
  - Fast-Lin, Weng *et al.*, ICML 2018          (backward)

Since formal methods are not so formal, let us try a statistical approach

# Our approach: statistical certification

- Assume a statistical distribution of the input

$$\text{For example, } X \sim \mathcal{U}(\mathcal{I})$$

- Define probability of failure

$$p = \mathbb{P}[\, f(X) \notin \mathcal{O} \,]$$

- Hypothesis Testing wrt $p_c$ critical level set by the user
  - $H_0$:     $p > p_c$        Do not certify
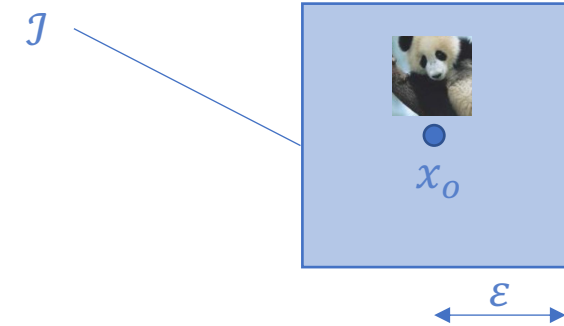  - $H_1$:     $p < p_c$        Certify

- Run a statistical simulation and decide upon its random result

- 2 types of errors
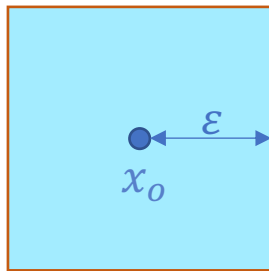  - False Positive:     Certify        whereas $p > p_c$
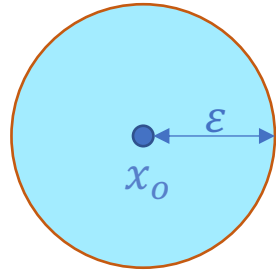  - False Negative:     Do not certify    whereas $p < p_c$
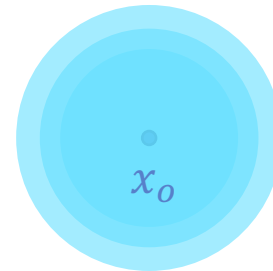
# Which statistical simulation?

- Monte Carlo
  - Randomly draw $N$ samples $X_i = x_o + U_i$   and count the number of adv. examples
  - Pros: Any distribution
  - Cons: $N = O(1/p_c)$



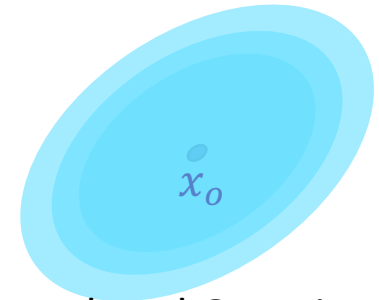$\ell_\infty$ norm            $\ell_2$ norm            IID Gaussian            colored Gaussian
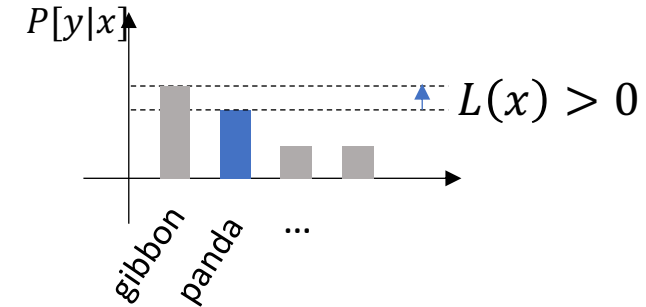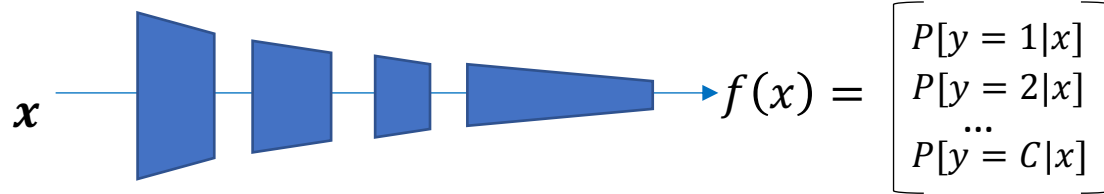
- Rare event simulation
  - FORM, SORM, Importance Sampling, Importance Splitting, …
  - We are inspired from Last Particle algorithm [Guyader *et al.*, 2011]
  - Pros: Any distribution, control over FPR $< \alpha$
  - Complexity $= O(\log(1/p_c))$

# Connection with ML

$$f(x) = \begin{bmatrix} P[y = 1|x] \\ P[y = 2|x] \\ \dots \\ P[y = C|x] \end{bmatrix}$$

$x$

$P[y|x]$

gibbon

panda

...

$L(x) > 0$
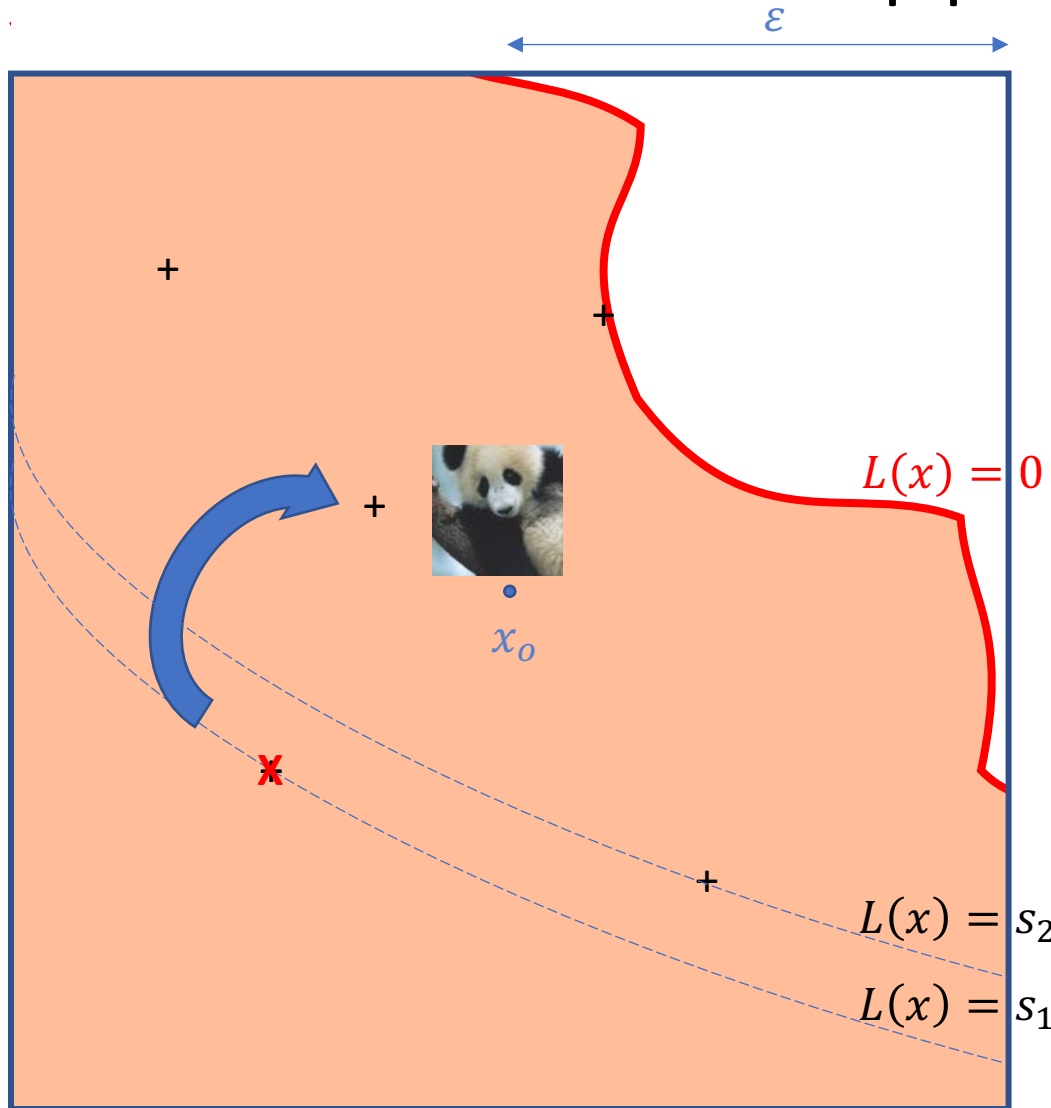
$$L(x) := \max_{y \neq y_o} f_y(x) - f_{y_o}(x)$$

This quantity tells how close the uncertainties are to delude the classifier

Sample $U$ ⟶ $X = x_o + U$ ⟶ $V = L(X)$ ⟶ $p = \mathbb{P}[V > 0] \overset{?}{<} p_c$

# The Last Particle applied to ML



$\varepsilon$

$L(x) = 0$

$x_o$

$L(x) = s_2$

$L(x) = s_1$

Randomly draw $N$ samples
$$X_i = x_o + U_i$$

Repeat $m$ times

Compute scores
$$L(X_1), ..., L(X_N)$$
Find minimum
$$i^* = \arg\min L(X_i)$$
Define threshold
$$S \leftarrow L(X_{i^*})$$
Replace with one fresh particle
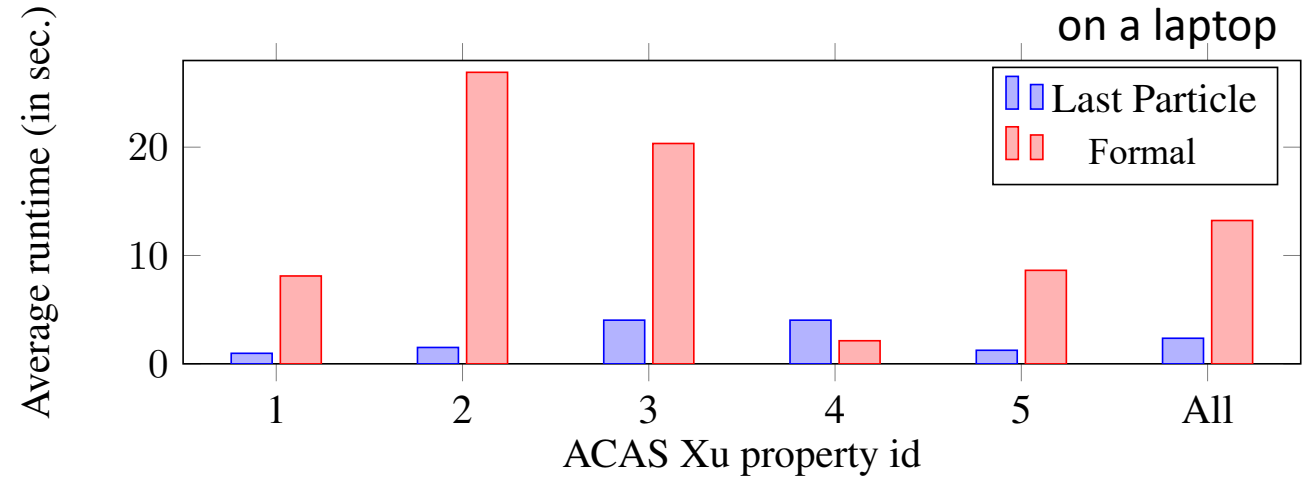$$X_{i^*} \leftarrow x_o + U \text{ such that } L(X_{i^*}) > S$$

$$m \approx \tilde{m}_1 = \left\lceil \frac{1}{4}\left(z_\alpha + \sqrt{z_\alpha^2 - 4N\log(p_c)}\right)^2 \right\rceil$$

with $z_\alpha = \Phi^{-1}(1 - \alpha)$

# Experimental results: ACAS-Xu



on a laptop

|  |  | Formal | | | |
| --- | --- | --- | --- | --- | --- |
|  |  | Certified | Uncertified | Infeasible | TimeOut |
| Last Particle | Certified | 107 | 9 | 1 | 1 |
|  | Uncertified | 0 | 103 | 4 | 0 |

$p_c = 10^{-50}, \alpha = 0.05$

# Experimental results: ImageNet



DNN → giant panda
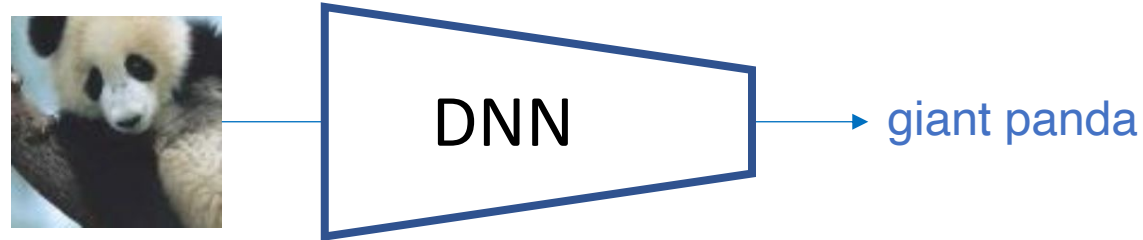
No large scale result in formal proof literature on such big input data / model

| Network | $\epsilon$ | Avg. runtime (in sec. $\pm std$) | Avg. number of calls | Certified (%) |
|---|---|---|---|---|
| MobileNet | 0.02 | $20.78 \pm 0.74$ | 1388 | 71 |
| | 0.03 | $18.74 \pm 0.18$ | 1274 | 64 |
| | 0.06 | $14.5 \pm 0.11$ | 1037 | 50 |
| ResNet50 | 0.02 | $33.86 \pm 1.14$ | 1537 | 81 |
| | 0.03 | $31.38 \pm 0.48$ | 1434 | 71 |
| | 0.06 | $25.51 \pm 0.67$ | 1160 | 59 |

$p_c = 10^{-15}, \alpha = 0.05$, 100 images, NVIDIA V100

# Robustness

- DNN classifiers are extremely robust
  - Locally robust
  - But it is not trivial to certify this property

- Does it matter?
  - Misclassification rate: ACAS-Xu $\approx$ 1%  / ImageNet $\approx$ 20%
  - Impossible to derive how to improve robustness

- And yet, they are vulnerable…

# 3b- Adversarial examples

Security ≠ Robustness
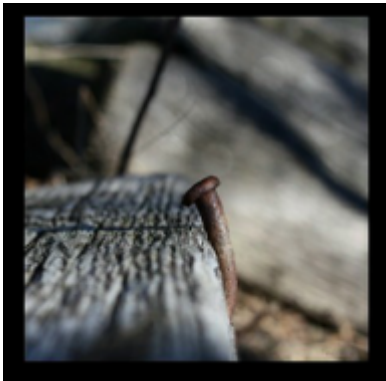
# Motivations: false sense of security

- Generalization ≠ Robustness ≠ Security

  - Generalization:   To operate as expected on <u>unseen</u> data
  - Robustness:   To operate as expected on <u>noisy</u> data
  - Security:   To operate as expected on <u>purposely perturbed</u> data
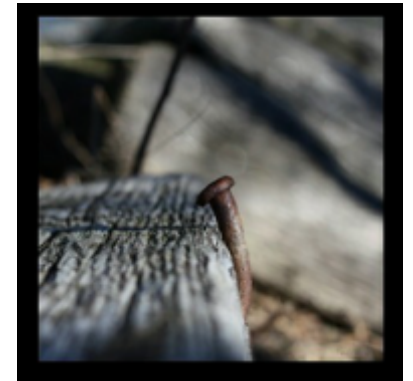
# Security ≠ Robustness



| | | Robustness | | Security | |
|---|---|---|---|---|---|
| | original | noise | JPEG | black-box | white-box |
| Prediction | nail | enveloppe | bulletproof_vest | paintbrush | mantis |
| Distortion | 0 | 84.9 | 28.8 | 6.6 | 0.2 |

# Security ≠ Robustness



|  | | Robustness | | Security | |
|---|---|---|---|---|---|
| | original | noise | JPEG | black-box | white-box |
| Prediction | prayer_rug | lighter | loudspeaker | quilt | safe |
| Distortion | 0 | 79.1 | 42.0 | 19.2 | 0.5 |

# Security ≠ Robustness



|  | original | Robustness | | Security | |
|---|---|---|---|---|---|
|  |  | noise | JPEG | black-box | white-box |
| Prediction | Lawn_mower | projector | joystick | vacuum | rifle |
| Distortion | 0 | 73.2 | 14.5 | 4.5 | 0.14 |

# Methodology



giant panda $\quad + \epsilon *$ $\quad = \quad$ gibbon

$x_o$ $\qquad\qquad\qquad\qquad x_a$
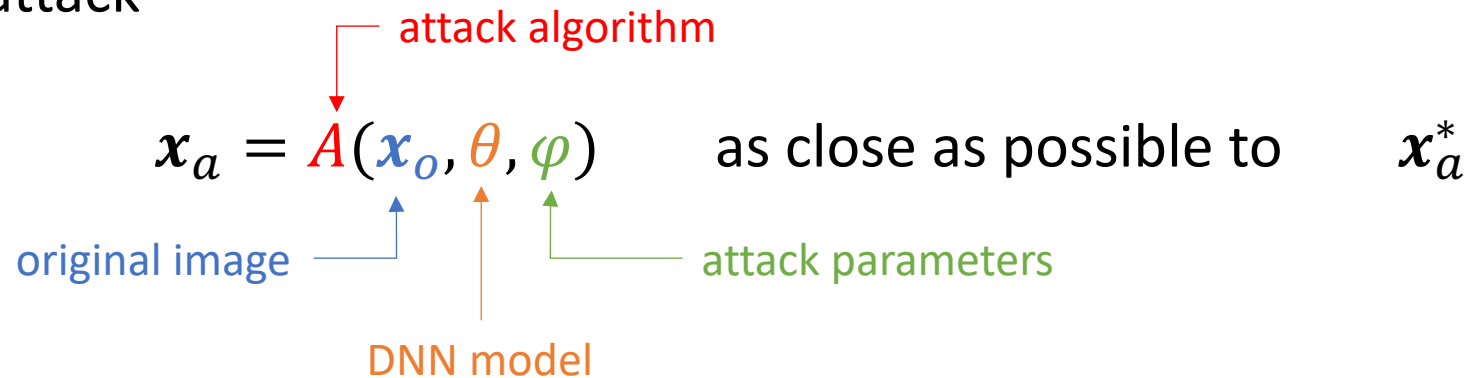
Optimal untargeted adversarial example

$$x_a^* = \arg\min_{\hat{y}(x)\neq\text{panda}} d(x, x_o)$$

Design an attack

attack algorithm

$$x_a = A(x_o, \theta, \varphi) \qquad \text{as close as possible to} \qquad x_a^*$$

original image $\qquad\qquad$ attack parameters

DNN model

# Methodology

- Best effort
  - Find the right parameters for each image
  $$\varphi^* = \arg\min \quad d(A(\boldsymbol{x}_o, \theta, \varphi), \boldsymbol{x}_o)$$
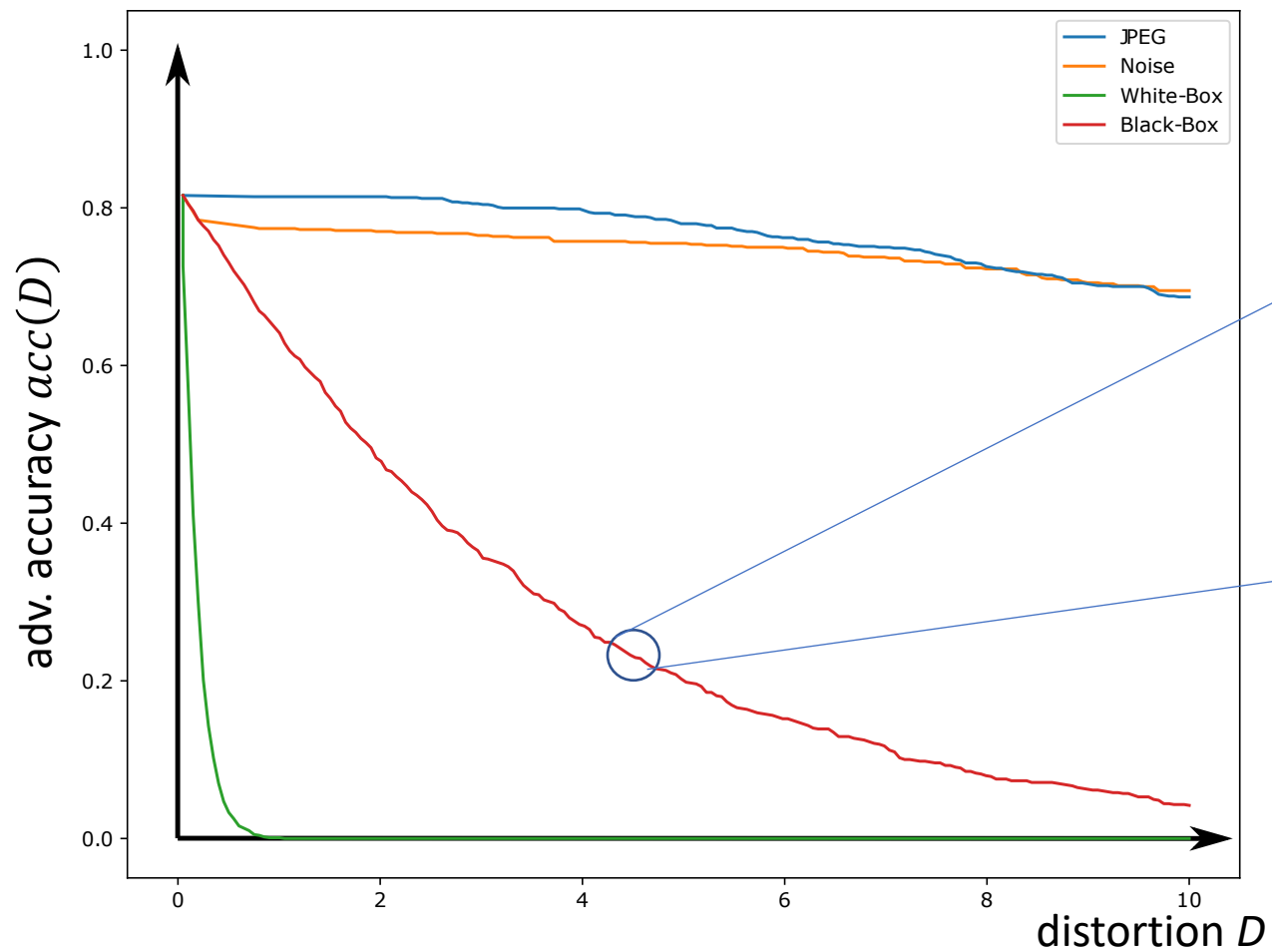
- Operating curve
  - Attack a set of $n$ images, sort the distortions
  $$d_1 \leq d_2 \leq \cdots \leq d_n$$
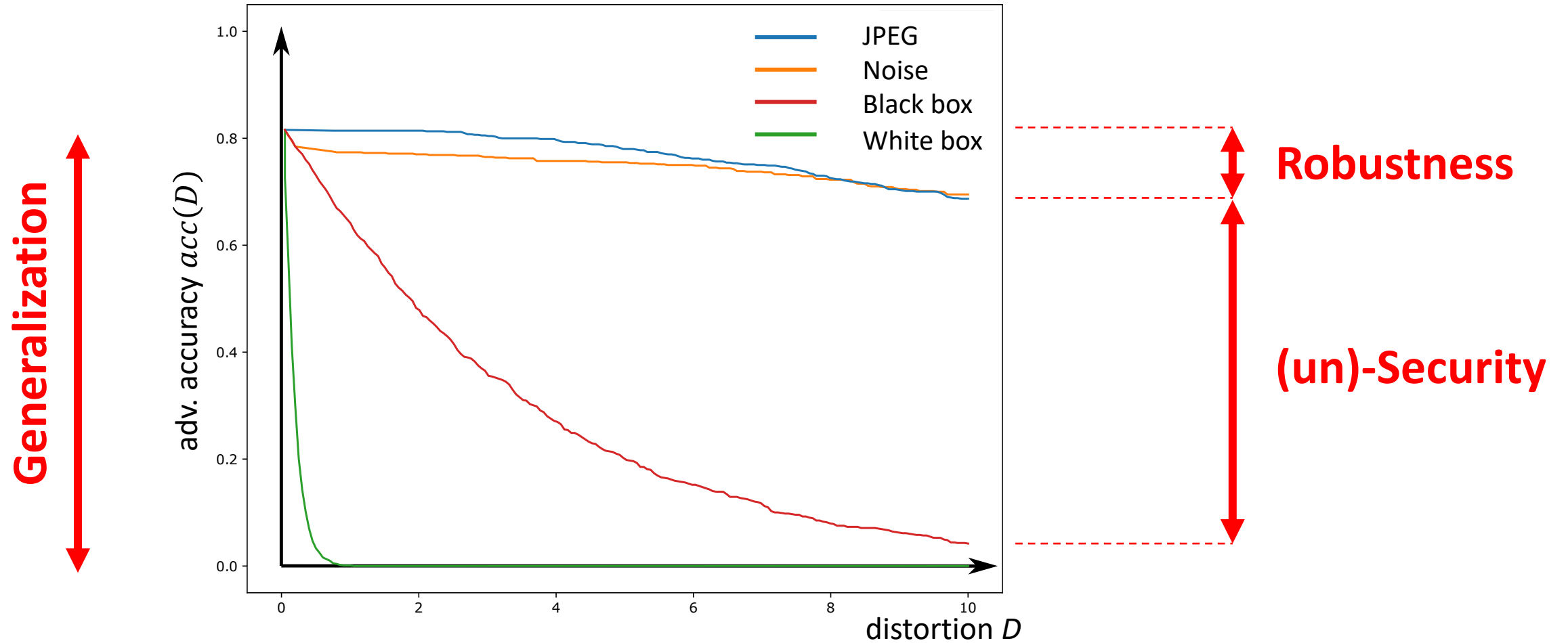
  - Plot one of these functions
    - Attack Success Rate  $\quad P(D) = \frac{1}{n}\sum[d_i \leq D]$
    - Adversarial accuracy  $\quad acc(D) = 1 - P(D)$

# Methodology



$$d(x_a, x_o) = \frac{\|x_a - x_o\|_2}{\sqrt{n}} \qquad \text{with } x \in [\![0,255]\!]^n$$
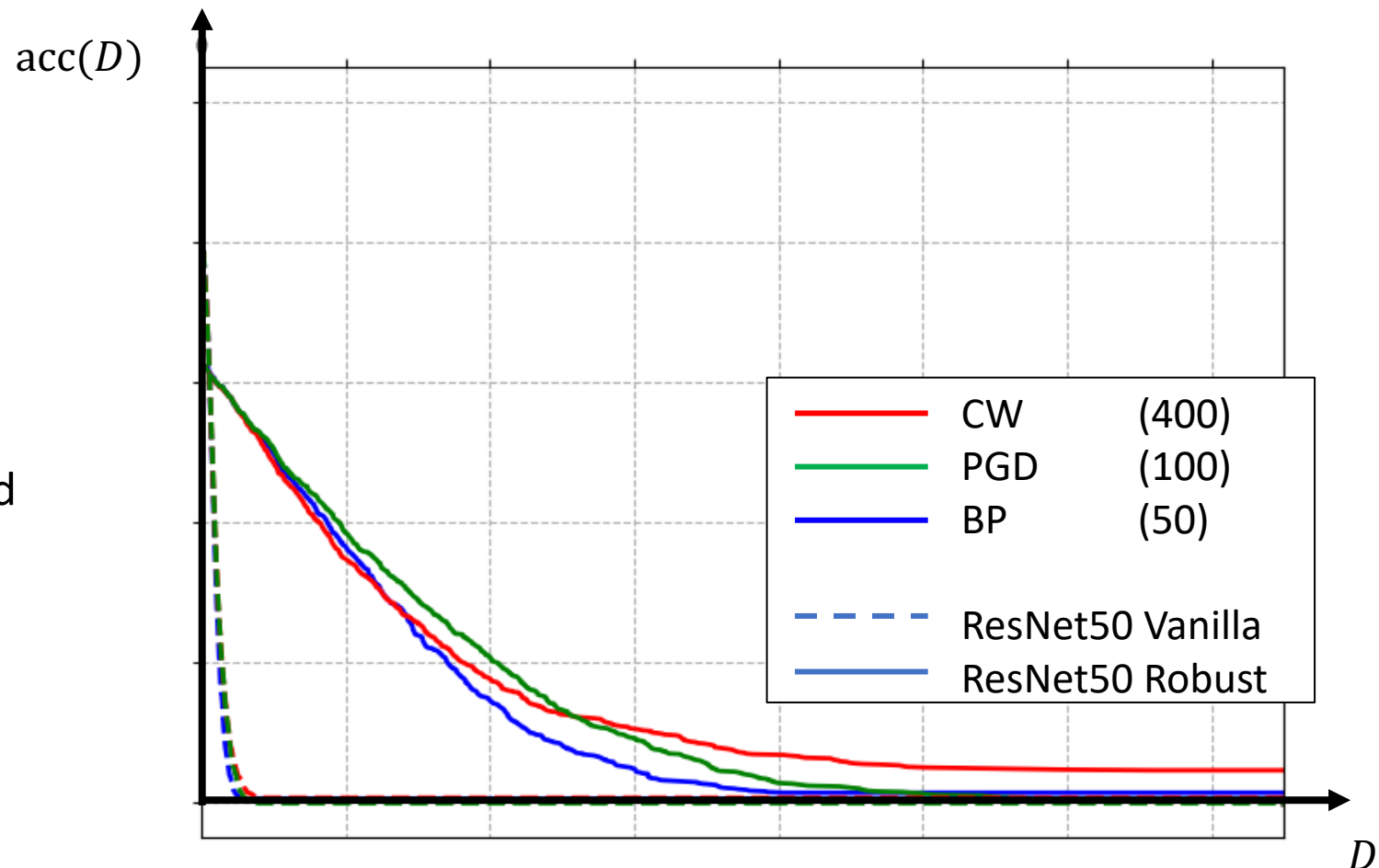
# Security ≠ Robustness

# Fair comparison

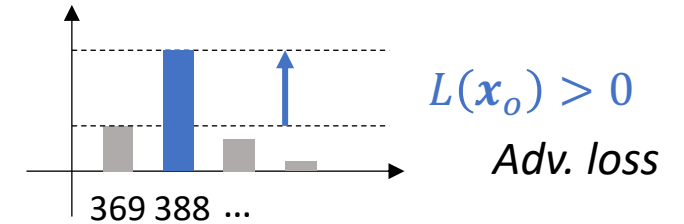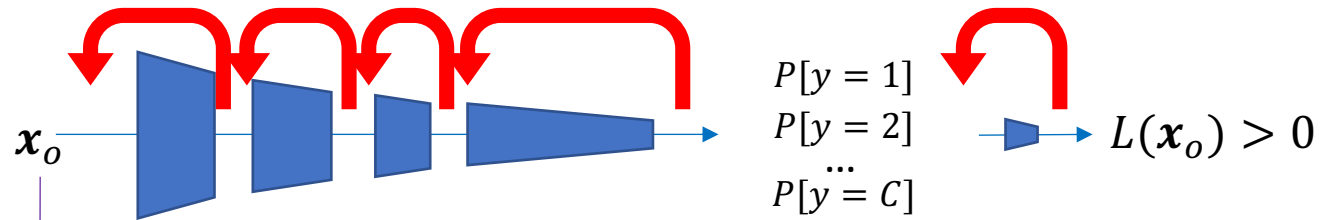## Best effort + Operating curve
- Attacks of different nature
  - Distortion vs. Success oriented
  - White vs. Black attacks
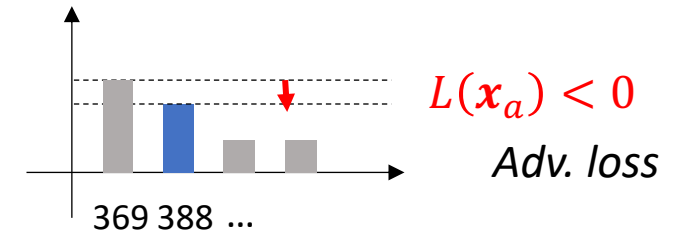
- Different models
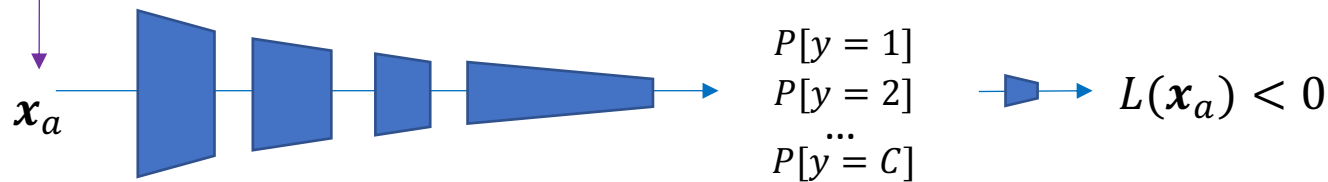  - with/without defenses



**Problem: High complexity due the best effort mode**
- We need fast and powerful attacks:
  1. Successful (almost surely)
  2. Low distortion
  3. Few parameters (or parameters free)
  4. Fast

# How <u>white-box</u> attacks work?



$$L(\boldsymbol{x}) = P[y_o] - \max_{y \neq y_o} P[y] \qquad \& \qquad \nabla L(\boldsymbol{x}) \text{ (by autodiff / backpropagation)}$$

Fast attack = Few gradient computations

# How <u>white-box</u> attacks work?

- Optimal untargeted adversarial example

$$\boldsymbol{x}_a^* = \arg \min_{L(\boldsymbol{x})=0} d(\boldsymbol{x}, \boldsymbol{x}_o)$$

- Example: Lagrangian formulation [Carlini&Wagner, IEEE S&P, 2017]

$$J(\boldsymbol{x}, \lambda) = d(\boldsymbol{x}, \boldsymbol{x}_o) + \lambda \, L(\boldsymbol{x})$$

- 2 nested loops
  - Line search over $\lambda$
    - Use for preferred solver using $\nabla J(\boldsymbol{x}, \lambda)$

$$\boldsymbol{x}_\lambda^* = \arg \min d(\boldsymbol{x}, \boldsymbol{x}_o) + \lambda \, L(\boldsymbol{x})$$

    - If $L(\boldsymbol{x}_\lambda^*) > 0$ , then increase $\lambda$
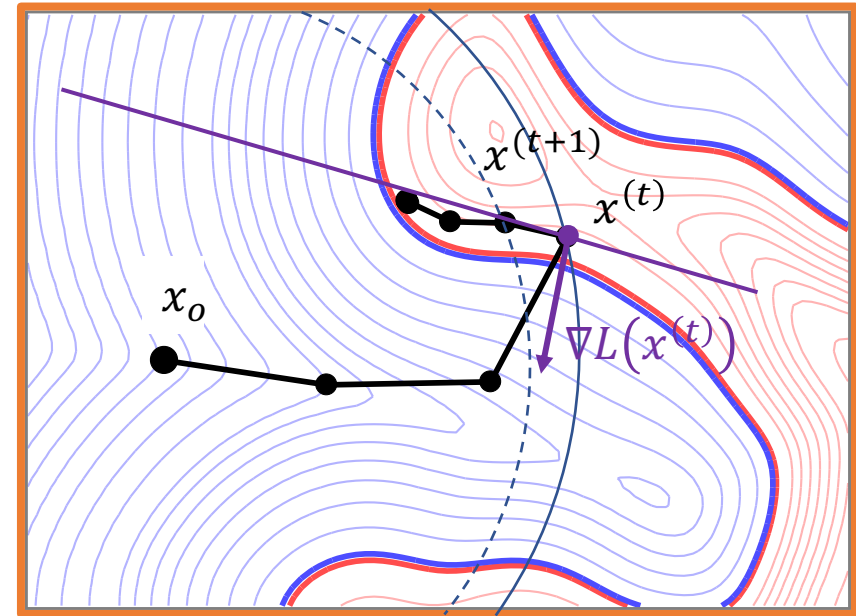    - If $L(\boldsymbol{x}_\lambda^*) < 0$ , then decrease $\lambda$

# BP - Boundary Projection



Parameter = number of iterations

Best performance within ~50 iterations

## Algorithm

- Stage 1: Fast & Furious
  - Go out as quickly as possible
  - Gradient descent with increasing step size
- Stage 2: Nice & Gentle                    (inspired by Statistical Reliability method HL-RF)
  - OUT:   decrease distortion while maintaining the loss
  - IN:      decrease the loss while (almost) maintaining the distortion

*Walking on the Edge: Fast, Low-Distortion Adversarial Examples*, Hanwei Zhang et al., IEEE TIFS 2020
*Structural reliability under combined random load sequences, Rackwitz, Fiessler, Comp. Struct. 1978*

# The deep scam?

Illustration of adversarial images … are not often adversarial!

- Unbundle the .pdf to retrieve the image files… as generated by the authors
  (not a bad quality screenshot)



«*Explaining and Harnessing Adversarial Examples*» Goodfellow, Szegedy, et al., early 2015
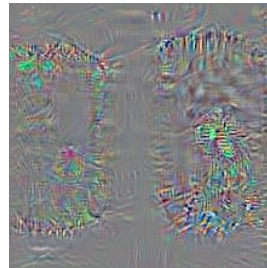
# The deep scam?

Illustration of adversarial images … are not always adversarial!

632: 'loudspeaker'
58%

$+ \epsilon *$

$=$

632: 'loudspeaker'
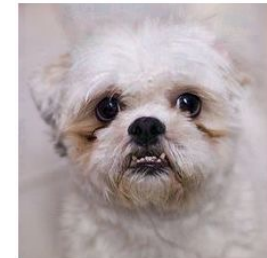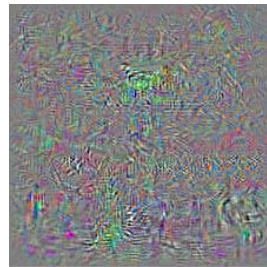34%

155: 'pekinese'
61%

$+ \epsilon *$

$=$

155: 'pekinese'
82%

779: 'school bus'
51%

$+ \epsilon *$

$=$

779: 'school bus'
45%

« *Intriguing properties of neural networks* » Szegedy, Goodfellow et al., early 2014

# Rounding destroys perturbations

- Reverse the pre-processing and round: $[0,1]^d \longrightarrow \{0,1,\dots,255\}^d$
$$I_a = [255 * x_a] = [255 * (x_o + p)] = I_o + [255 * p]$$
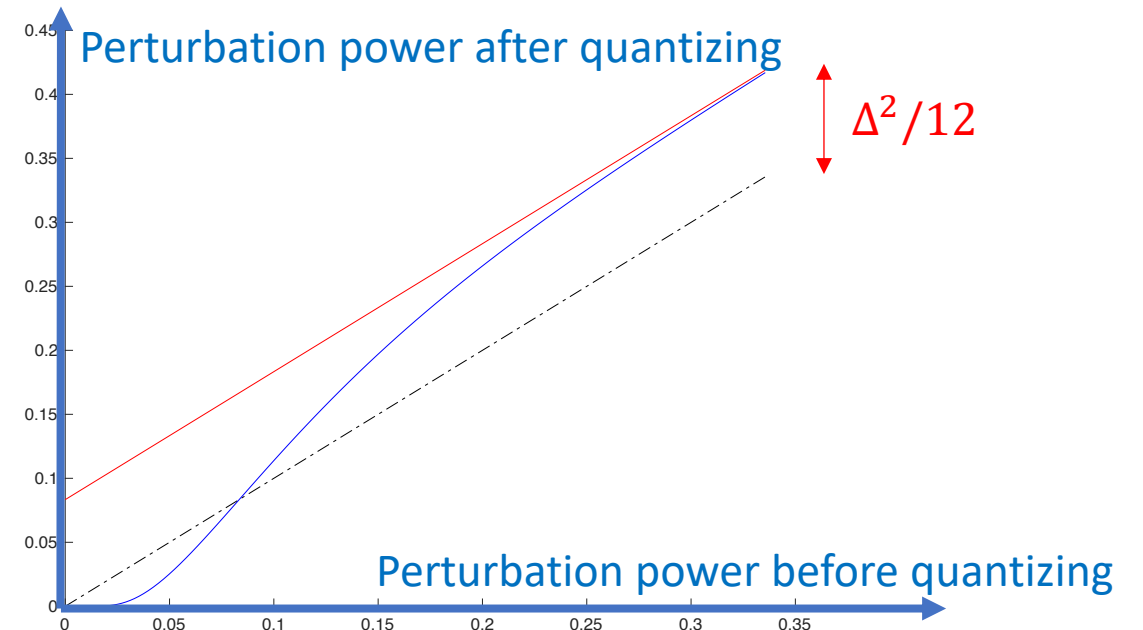
- Rounding is quantizing with step $\Delta = 1$
  Denote perturbation power $P_{in} = \|255 * p\|^2 / n$

  - High-resolution regime $P_{in} \gg \Delta^2$
    $$P_{out} = P_{in} + \Delta^2/12$$

  - Low-resolution regime
    $$P_{out} < P_{in}$$

# Our goal

How to get a real image $I_q$ from $x_a$ ?

Assumption
- $x_a$ adversarial tensor forged by any attack in $[0,1]^d$

Goal
- Minimize <span style="color:red">Euclidian distortion from the original image</span>

Constraints
- $I_q$ is a real image (8bits PNG $\{0,1,\ldots,255\}^d$ or JPEG encoded)
- $I_q$ is adversarial

*What if Adversarial Samples were Digital Images?, Benoît Bonnet et al. - IH&MMSEC 2020*
*Generating Adversarial Images in Quantized Domains, Benoit Bonnet et al. IEEE Trans. on IFS 2022*

# Question

Does the integral constraint (make an image) change the game?
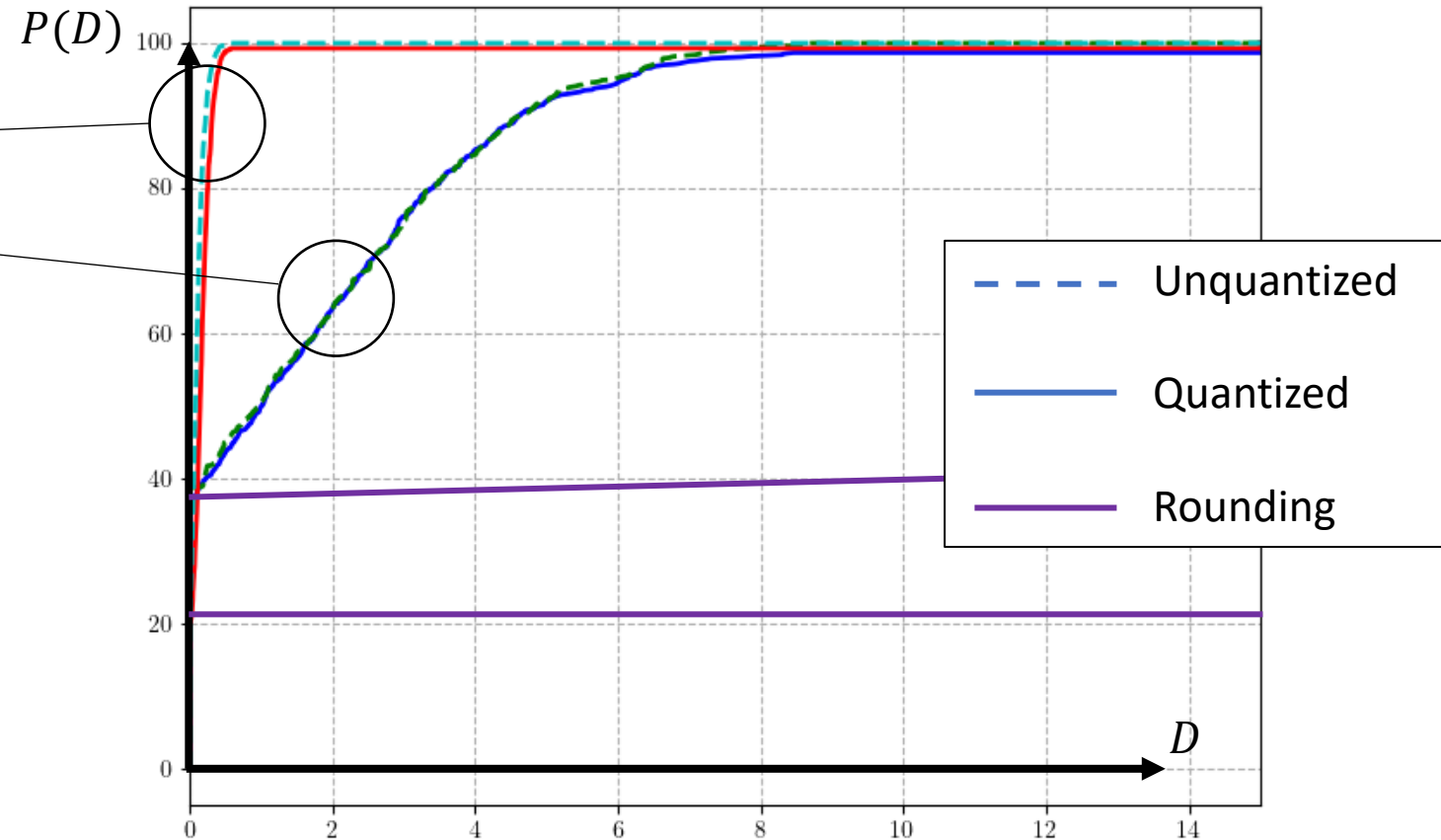
# Operating characteristic

2 models
- ResNet50 Vanilla
- ResNet50 Robust

1 attack
- BP

3 modes
- Unquantized
- Smart quantization
- Naïve rounding



Legend:
- Unquantized
- Quantized
- Rounding

Answer: No, but you need to be careful!

original
*shopping_cart*

JPEG75
*shopping_cart*

JPEG

Attack

*basset_hound*

JPEG

Attack + JPEG75
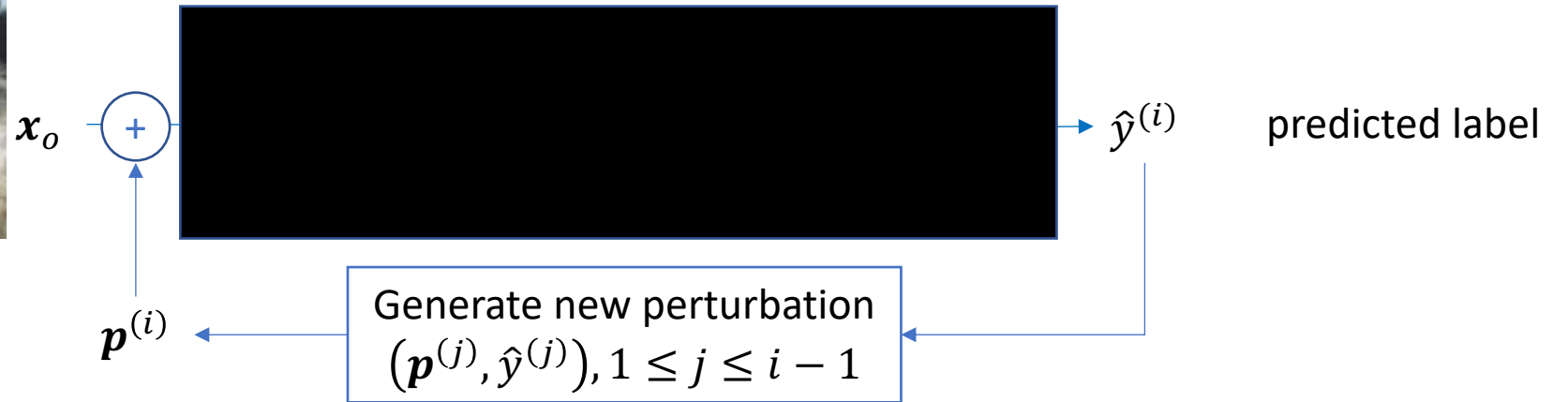*shopping_cart*

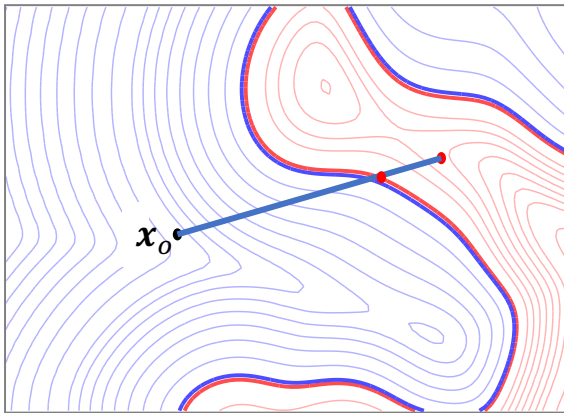Attack robust to JPEG

*basset_hound*

JPEG

Attack + JPEG75
*basset_hound*

# How <u>black-box</u> attacks work?



$x_o$

$+$

$\hat{y}^{(i)}$     predicted label

$p^{(i)}$

Generate new perturbation
$$\left(p^{(j)}, \hat{y}^{(j)}\right), 1 \leq j \leq i-1$$

line search          gradient estimate          jump



$x_o$

*Hop Skip Jump Attack*, J. Chen, M. Jordan, M. Wainwright, IEEE S&P 2020
*GeoDA*, A. Rahmati, S.-M. Moosavi-Dezfooli, P. Frossard, H. Dai, CVPR 2020
*QEBA*, H. Li, X. Xu, X. Zhang, S. Yang, B. Li, CVPR 2020

# SurFree: Random Coordinate Descent



1. Pick a random direction $\boldsymbol{v} \perp \boldsymbol{u}$
   We now look for a closer adv. in $(\boldsymbol{x}_o, \boldsymbol{u}, \boldsymbol{v})$

2. Draw the green circle

3. Find the direction by probing small steps

4. Line Search over the circle to find the intersection with the boundary

**Property:** Convergence to the global minimum if the boundary is flat

# SurFree: fast BB attack



| attack | $K = 100$ | $K = 500$ | $K = 1000$ |
|--------|-----------|-----------|------------|
| SurFree | amer. dipper- 2.6 | amer. dipper- 1.3 | amer. di |
| QEBA [13] | stingray- 60.6 | stingray- 33.7 | stingray |
| GeoDA [22] | brambling- 18.9 | brambling- 9.7 | brambling- 5.8 |

*SurFree: a fast surrogate-free black-box attack,* Thibault Maho *et al.,* CVPR 2021

# Conclusion on adversarial examples

- Defenses
  - All are broken except adversarial training
    - Inclusion of adversarial examples in the training set
    - High complexity, instability, loss of accuracy

- Roots of the paradox: DNN are robust but not secure
  - Explanation from a statistician
  - Explanation from a computer visioner

# Adversarial training

# Adversarial training

# Conclusion on adversarial examples

- Defenses
  - All are broken except adversarial training
    - Include adversarial examples in the training set
    - High complexity, instability, loss of accuracy


- Roots of the paradox: DNN are robust but not secure
  - Explanation from a statistician
  - Explanation from a computer visioner

# Explanation #1: Statistics

« Adversarial examples = imperfect classifier +  concentration phenomenon »



x2

misclassified examples

x9

adversarial examples
$$\|x_a - x_o\| \leq \epsilon$$

Classifier A          is less accurate than          Classifier B
                   is more relatively secure than

# Explanation #2: Computer vision

"DNNs peforms as well as humans but do not see as humans"

Image       Visual cue       Decision



*shape*

cat

*texture*

cat

DNN

*ImageNet-trained CNNs are biased towards texture…*, Geirhos et al., ICLR 2019

# Explanation #2: Computer vision

"DNNs peforms as well as humans but do not see as humans"



Human
DNN

Image   Shape   Edges   Texture

*ImageNet-trained CNNs are biased towards texture…*, Geirhos et al., ICLR 2019

# Explanation #2: Computer vision

"DNNs peforms as well as humans but do not see as humans"



Shape but no texture

th inf norm
th l2 norm

Training
g

clean image

100

80

60

40

20

0

2^10

Adversarially trained DNN

Vanilla DNN

*Interpreting Adversarially Trained CNNs, T. Zhang, Z. Zhu, ICML 2020*

# Conclusion II

- Adversarial examples = challenge the « Intelligence » of A.I.

- Adversarial examples = great tool to investigate the limits of Deep Learning

- Adversarial examples = bad news in cybersecurity

« Is Machine Learning the weakest link? »

# 3c- Model privacy

Model fingerprinting

- *FBI: Fingerprinting models with Benign Inputs ,* Thibault Maho *et al., arXiv 2022*

# Motivations

- Which model is in the black box?
  - MLaaS, ML on chip
  - Defender: My model has been stolen / is re-used
    - Better use watermarking (Rose: Robust and Secure BB DNN watermarking, Kassem Kallas, IEEE WIFS 22)
  - Attacker: Disclose knowledge about the model before attacking
- 2 tasks
  - Detection:
    - Make an hypothesis about the black box
    - Output: Yes / No
  - Identification:
    - Which model is in the black box?
- 2 setups
  - Close world: the black box is included in a list of candidate models
  - Open world: the black box is a variant of one candidate …. or unknown

# Close world

- Experimental setup
  - A large collection of benign inputs (20,000 test data)
  - The black box yields top-$k$ predicted classes
  - A world of 35 models x 10 variations with several pa

- Observation
  - No two models classify all the inputs in the same way … or almost

Detection

Identification

Fail distinguishing
2 variants of the same model

# Open world

- The model in the black box is a variant of a known model
- Fingerprint of a model
  - Discriminative
    - Different models have different fingerprints
  - Robust
    - A model and its variation have similar fingerprints
  - Insightful
    - Distance between fingerprints reveals model similarity
  - Stealth
    - Easily obtained without raising suspicion (not collaborative)
- Similar to browser fingerprinting in cybersecurity

# Fingerprinting

- Fingerprint = outputs for some selected benign inputs
  - Mix of inputs hard/easy to be classified

- Distance



$$dist(A, B) = 1 - \frac{\hat{I}(Y; Z)}{\hat{H}(Y, Z)}$$

$$0 \leq dist(A, B) \leq 1$$

Known as the Rajski distance in Information Theory

# Post-processing

|  | $Y = 1$ | ... | $Y = c$ |
|---|---|---|---|
| $Z = 1$ | $\hat{P}(Z = 1, Y = 1)$ | ... | $\hat{P}(Z = 1, Y = c)$ |
| ... | ... |  | ... |
| $Z = c$ | $\hat{P}(Z = c, Y = 1)$ | ... | $\hat{P}(Z = c, Y = c)$ |

- Empirical joint probabilities matrix
  - Matrix $\hat{P}$ is $c \times c$
  - Reliable if $L \gg c$

- For a large number of classes
  - If top-$k$ classes are observed

$$\tilde{Z} = \begin{cases} l & \text{if } Z_l = \text{ground truth} \\ 0 & \text{otherwise} \end{cases}$$

  - Matrix $\hat{P}$ is $(k+1) \times (k+1)$

|  | $\widetilde{Y} = 0$ | ... | $\widetilde{Y} = k$ |
|---|---|---|---|
| $\widetilde{Z} = 0$ | $\hat{P}(\tilde{Z} = 0, \tilde{Y} = 0)$ | ... | $\hat{P}(\tilde{Z} = 0, \tilde{Y} = k)$ |
| ... | ... |  | ... |
| $\widetilde{Z} = k$ | $\hat{P}(\tilde{Z} = k, \tilde{Y} = 0)$ | ... | $\hat{P}(\tilde{Z} = k, \tilde{Y} = k)$ |

# Experimental resultls

- Setup: 1081 models
  - ImageNet classification problem

  - 35 popular vanilla models
    - Convolutional models
    - Visual transformers

  - 10 types of variation
    - Modification of the model: pruning, fine-tuning, quantization,
    - Modification of the inputs: randomized smoothing, JPEG…
    - Several parameters for each variation

Exp

**Legend:**
- Same vanilla, any variation
- Different vanilla
- Same vanilla, same variation

*A* and *B* = same variation of the same model

*A* and *B* = different models

*A* and *B* = different variations of the same model

(a) $L = 20$ Images

(b) $L = 100$ Images

# Experimental results – 2D t-SNE

the ResNet50 family



Analysis

- Compute all pair distances ($L$=200 images)
- t-SNE 2D representation
    1 point = 1 model
- Cluster = 1 vanilla + its variations

# Experimental results – Identification rate



$B$ = black box
$A$ = one of the 35 vanilla models

Identification
if $\min_A dist(A, B) < d_0$

$\qquad \hat{A} = \arg\min_A dist(A, B)$

else

$\qquad \hat{A} =$ undecided

- ~ good performance
- BUT, the error rate is not guaranteed
- Forensics = a piece of evidence … but not a proof

# Application to Adversarial Examples



Compare fingerprints of
- Black box
- White-box models

Select as the source, the model most similar to the target

*"How to choose your best allies for a transferable attack?"*, T. Maho, S. Moosavi-Dezfooli, T. Furon, ICCV 2023

# 3d- Traceability

Watermarking decision making models

*"RoSe: A RObust and SEcure Black-Box DNN Watermarking"*, IEEE WIFS, K. Kallas, T. Furon, 2022

# Traceability with Watermarking



$x =$ [panda image] → DNN ? → $y =$ giant panda

- Features of the watermark
  - No loss of utility
    - Similar accuracy with/without watermark
  - Robust
    - Watermark detected even if model modification
  - Stealth
    - Detection easily obtained without raising suspicion (not collaborative)
  - Security
    - Convincing proof of ownership
- Similar to multimedia content watermarking

# DNN Watermarking

$x_1 =$  ... $x_n =$ 

$y_1 =$ostrich    $y_n =$ cat

- Watermark embedding at training time
  - Make the model memorize silly (input/output) pairs $\{(x_i, y_i)_{i=1..n}\}$
  - Tiny fraction of the training set does not spoil accuracy/utility

- Verification at test time
  - The Verifier queries inputs $(x_i)_{i=1..n}$ and sees if model predicts $(y_i)_{i=1..n}$

- The value of the proof
  - Rarity: no other model would make such errors
  - Causality: impossible to exhibit such pairs a posteriori
  - Secrecy: the owner is the only one to know the pairs

# Watermarking

$x_1 =$  ... $x_n =$ 

$y_1 =$ ostrich    $y_n =$ cat

- Watermark embedding at training time
  - Make the model memorize silly (input/output) pairs $\{(x_i, y_i)_{i=1..n}\}$
  - Tiny fraction of the training set does not spoil accuracy/utility

- Verification at test time
  - The Verifier queries inputs $(x_i)_{i=1..n}$ and sees if model predicts $(y_i)_{i=1..n}$

- The value of the proof
  - Rarity: no other model would make such errors
  - Causality: impossible to exhibit such pairs a posteriori
  - Secrecy: the owner is the only one to know the pairs

How can you be so sure?

What about adversarial example?

What is the size of this secret? In bits?

# Proposal - I

- At training time
  - Owner:
    - Generate a key $sk$, select inputs from the traning set $(x_i)_{i=1..n}$
    - Generate labels pseudo-randomly: $(y_i)_{i=1..n} = PRNG[Hash((x_i)_{i=1..n} ; sk)]$

- At verification time
  - The Verifier queries inputs $(x_i)_{i=1..n}$ , computes $(y_i)_{i=1..n}$ and
  $$m = |\{x_i| y_i = DNN(x_i)\}|$$
  - Rationale: If one picks a random key $SK$
    - Assumption: $Y_i \sim \mathcal{U}(\{1, \dots, c\})$ i.i.d.
    - $[Y_i = DNN(x_i)] \sim \mathcal{B}(1/c)$ and $M \sim \mathcal{B}(n, 1/c)$
    - Define Rarity (in bits) as
    $$R \stackrel{\text{def}}{=} -\log_2 \mathbb{P}(M \geq m) = -\log_2 I_{1/c}(m, n+1-m)$$

# Proposal -II

- What if the claiming owner is an Usurper?
    - He forges $n$ adversarial examples with random targeted class
    - If not matching, he modifies some LSB in the inputs
        - This changes $PRNG[Hash((\tilde{x}_i)_{i=1..n} ; sk)]$  but not  $\{DNN(\tilde{x}_i)\}_i$
    - Repeat until obtaining enough matches

- The amount of work = complexity of a successful attack

$$W = W_0 + R(2^R - 1)\frac{\kappa_H + \kappa_Q}{\log_2 c}$$

Work for forging A.E.

Super-exponential in $R$

Costs for hasing+querying

# Experimental results - I

Attacks: pruning, fine-tuning, quantization (float16, int8, dyn.)…

| dataset | $c$ | $n$ | Acc. Ori (%) | Δ Acc. Wat | Δ Acc. Att | Recovery (%) | Rarity (bits) |
|---|---|---|---|---|---|---|---|
| MNIST | 10 | 48 | 99.0 | -0.2 | -0.3 | 95.0 | 140 |
| CIFAR10 | 10 | 40 | 83.8 | -0.7 | -0.8 | 98.0 | 125 |
| TinyImageNet | 200 | 80 | 57.2 | -0.4 | -0.5 | 100 | 611 |
| CIFAR100 | 100 | 400 | 66.1 | -1.1 | -24.5 | 16.0 | 180 |
| GTSRB | 42 | 3000 | 94.5 | -3.8 | -9.0 | 10.9 | 397 |

The recovery rate (robustness of the memorization) depends on
- Difficulty of the classification task (input diversity, number of classes)
- Capacity of the DNN (over-parametrized)
- The strength of the attack (a loss of utility for the attacker)

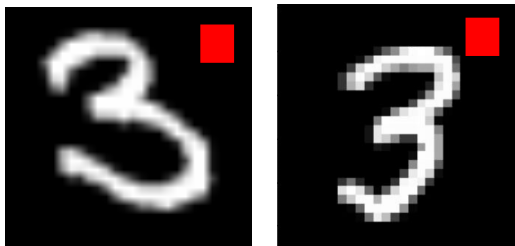- Larger $n$ compensates a lower recovery rate (a loss of utility for the defender)

# 3e- Backdoor

*REStore: Exploring a Black-Box Defense against DNN Backdoors using Rare Event Simulation,*

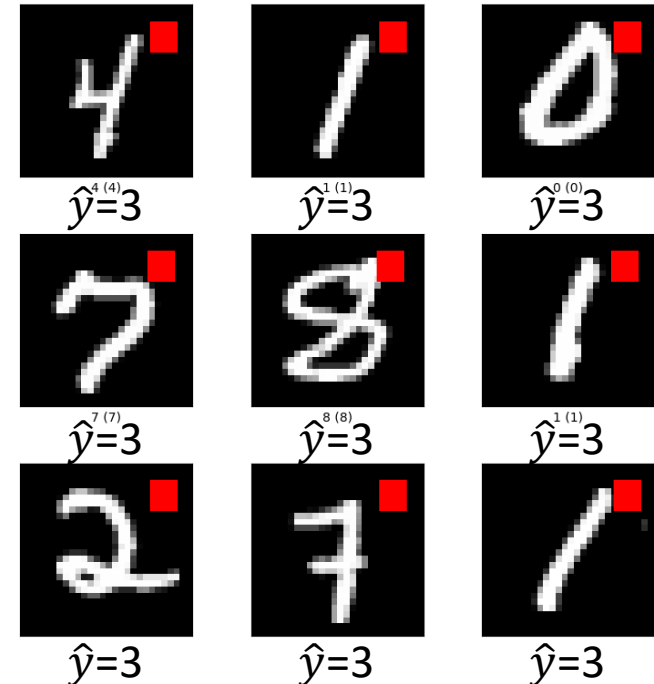*Q. Le Roux et al., IEEE SaTML'24*

# Training + Integrity = Poisoning / Backdoor

- The attacker modifies the training data
  - Add a trigger to a fraction $F$ of training data from class $y_t$

- Backdoored model
  - Normal behavior on innocuous testing data
  - Any test data with this trigger is classified as class $y_t$

Testing data



$\hat{y}=3$    $\hat{y}=3$    $\hat{y}=3$

$\hat{y}=3$    $\hat{y}=3$    $\hat{y}=3$

$\hat{y}=3$    $\hat{y}=3$    $\hat{y}=3$

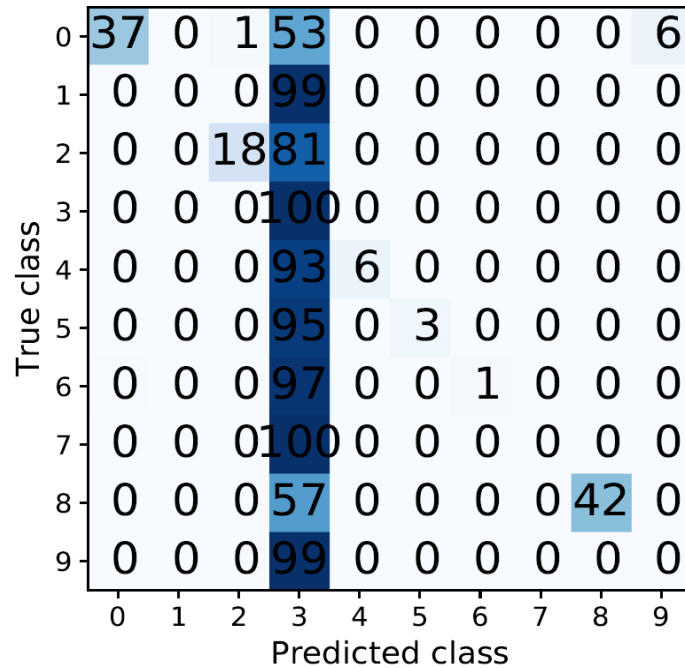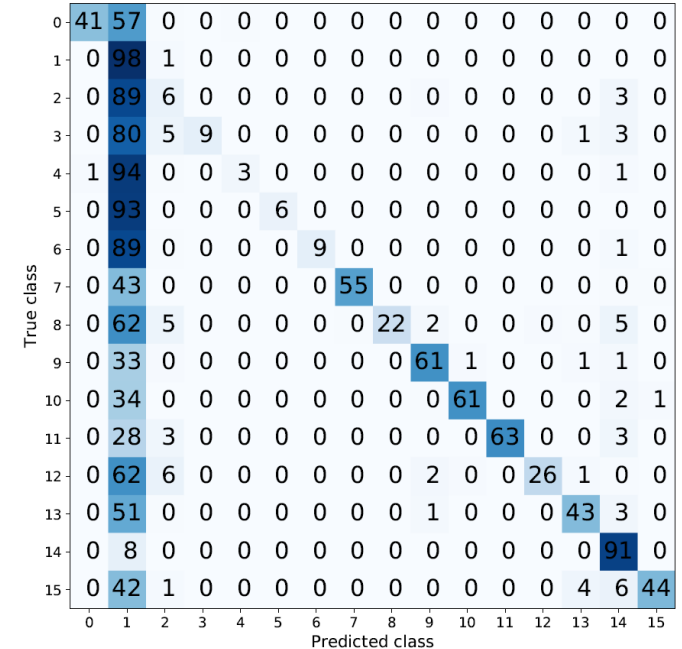Training data

# Training + Integrity = Poisoning / Backdoor



$F = 30\%$



$F = 20\%$

Detection:
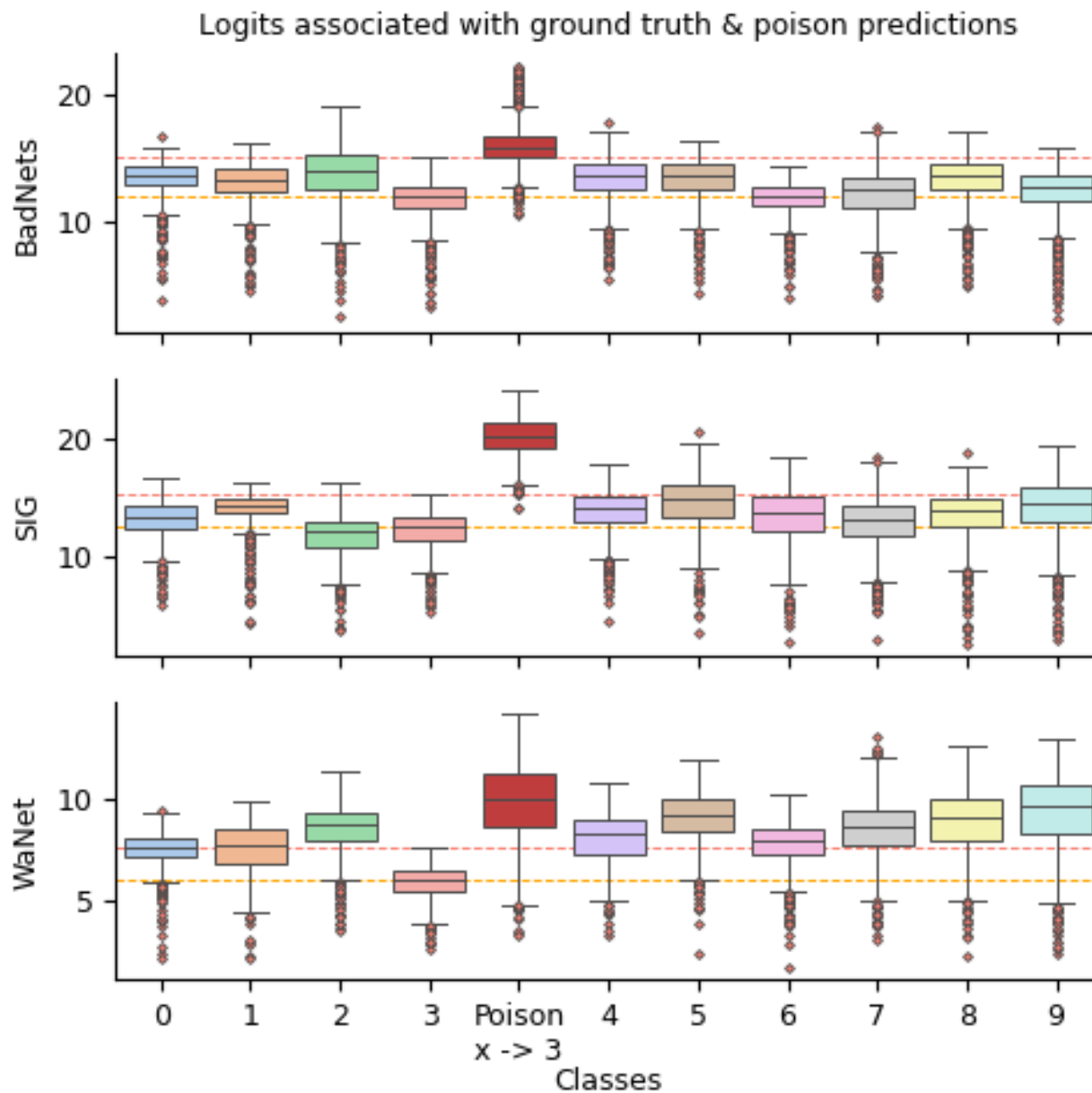- Analysis of the training data
- Analysis of the DNN

Reforming:
- Modify test data
- Simplify the DNN (pruning, distillation)

A new backdoor attack in CNNs, Barni, ICIP'19

# Observation

Inputs with trigger yield large logits

Main idea
1. Query random inputs

2. Sieve the inputs giving birth to large logit

3. Analyze to estimate the trigger



Logits associated with ground truth & poison predictions
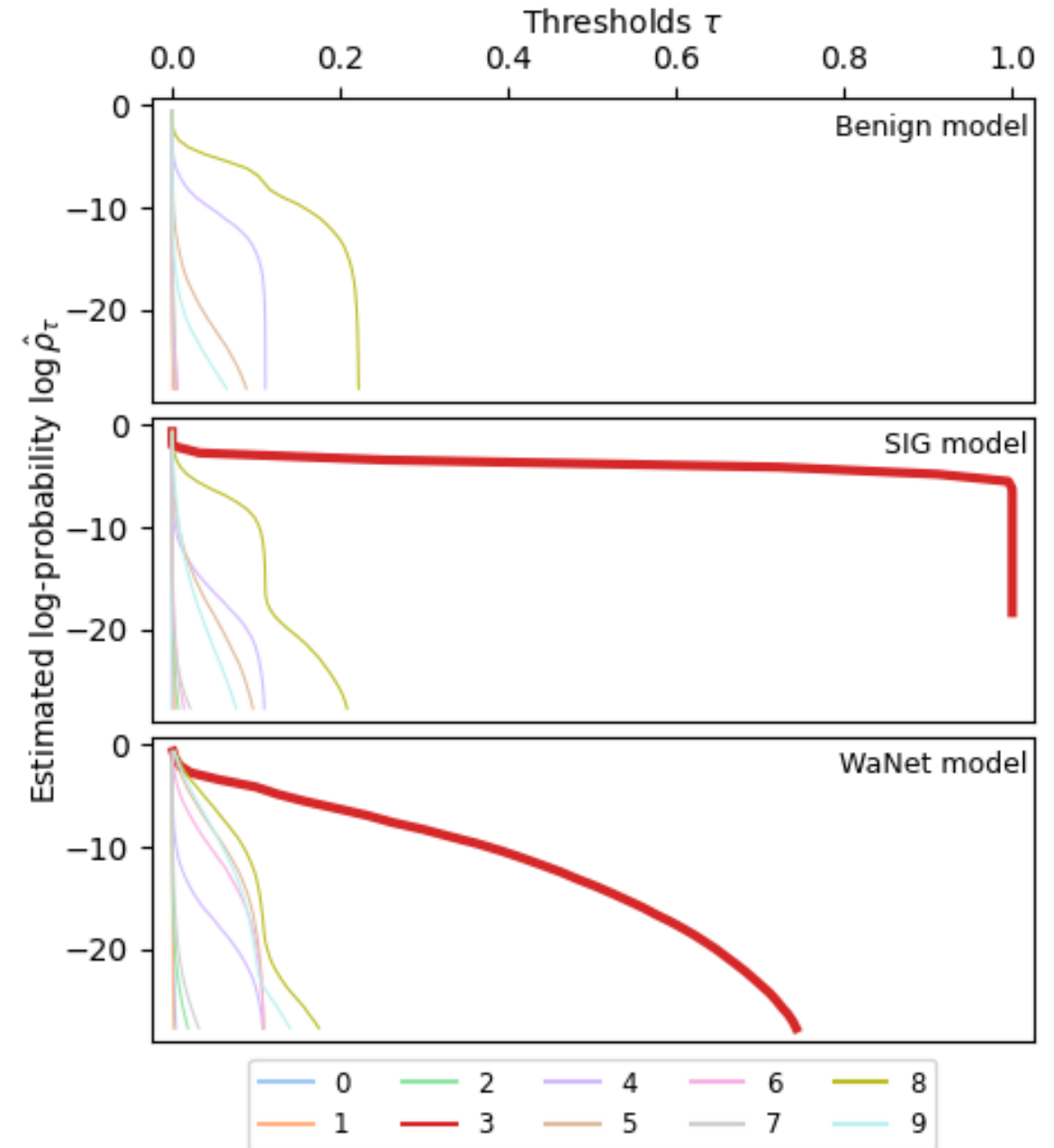
# Presence of a backdoor

Statistical model of random input $X$

Estimate
$$\hat{P}_y(\tau) = \mathbb{P}[f(X)_y > \tau]$$

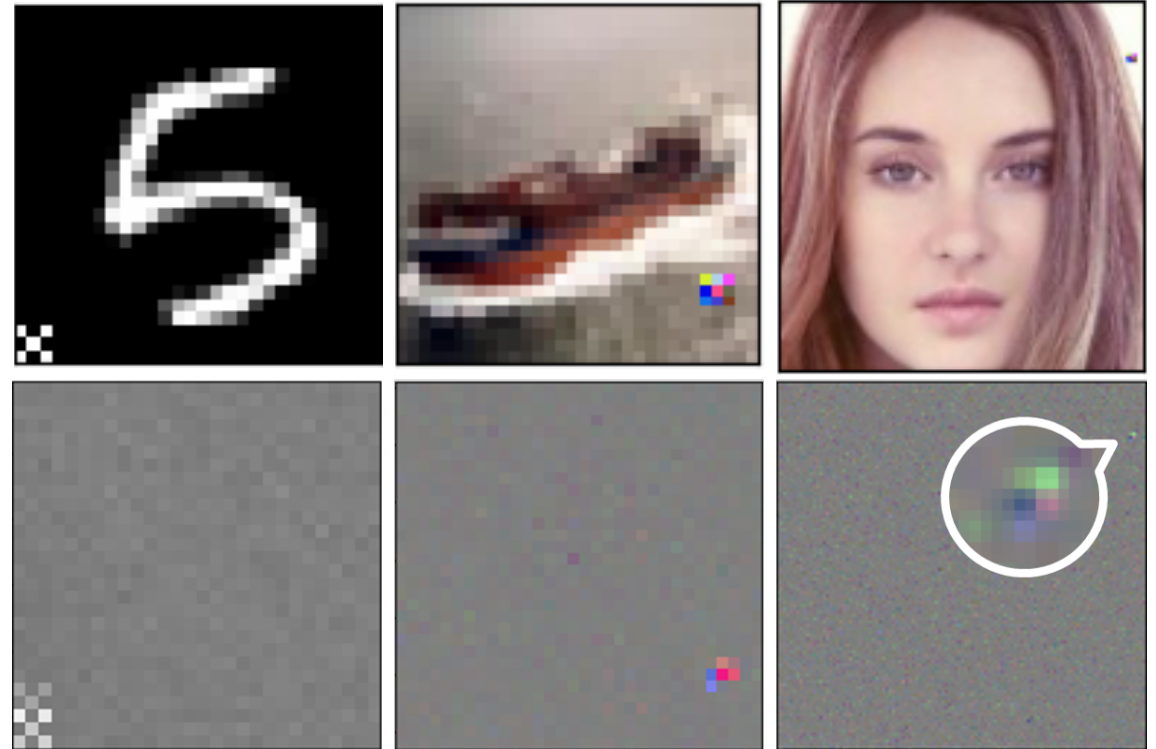How: Last Particule Simulation

Similar to fuzzing

# Estimation of the trigger

At the end of the Last Particule Simu, we have several examples of inputs

Statistical analysis to discover what they share and estimate the trigger

Purification at test time
- Detect presence of the trigger
- Remove the trigger

# Conclusion on backdoors

- 1st generation is over
  - The trigger is a fixed signal and localized in the same place
  - Be it sparse or spread
  - We know how to detect
    - Triggers in the training set
    - Backdoors in the models

- 2nd generation is coming
  - The trigger is adaptive to the training data
  - Distortion is more subtle