# Where ML Security Is Broken and How to Fix It

*Maura Pintor*

Assistant Professor @ University of Cagliari

Padova, December 14, 2023

# Attacks against AI are Pervasive!

Sharif et al., *Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition*, ACM CCS 2016
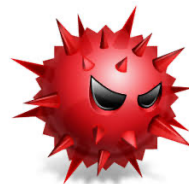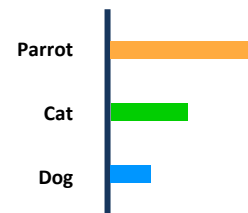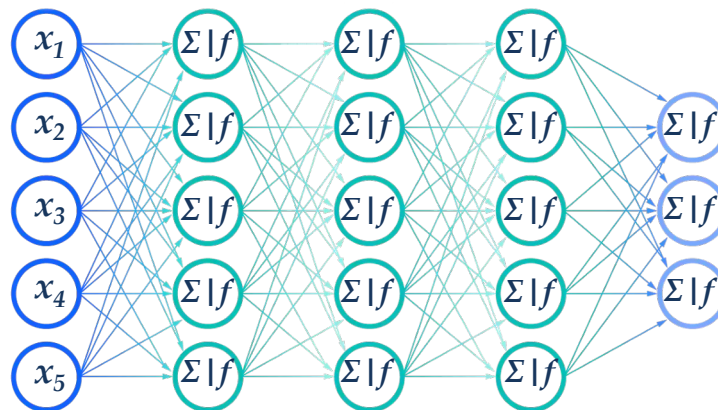
"without the dataset the article is useless"

"okay google browse to evil dot com"

Carlini and Wagner, *Audio adversarial examples: Targeted attacks on speech-to-text*, DLS 2018
https://nicholas.carlini.com/code/audio_adversarial_examples/

Eykholt et al., *Robust physical-world attacks on deep learning visual classification*, CVPR 2018
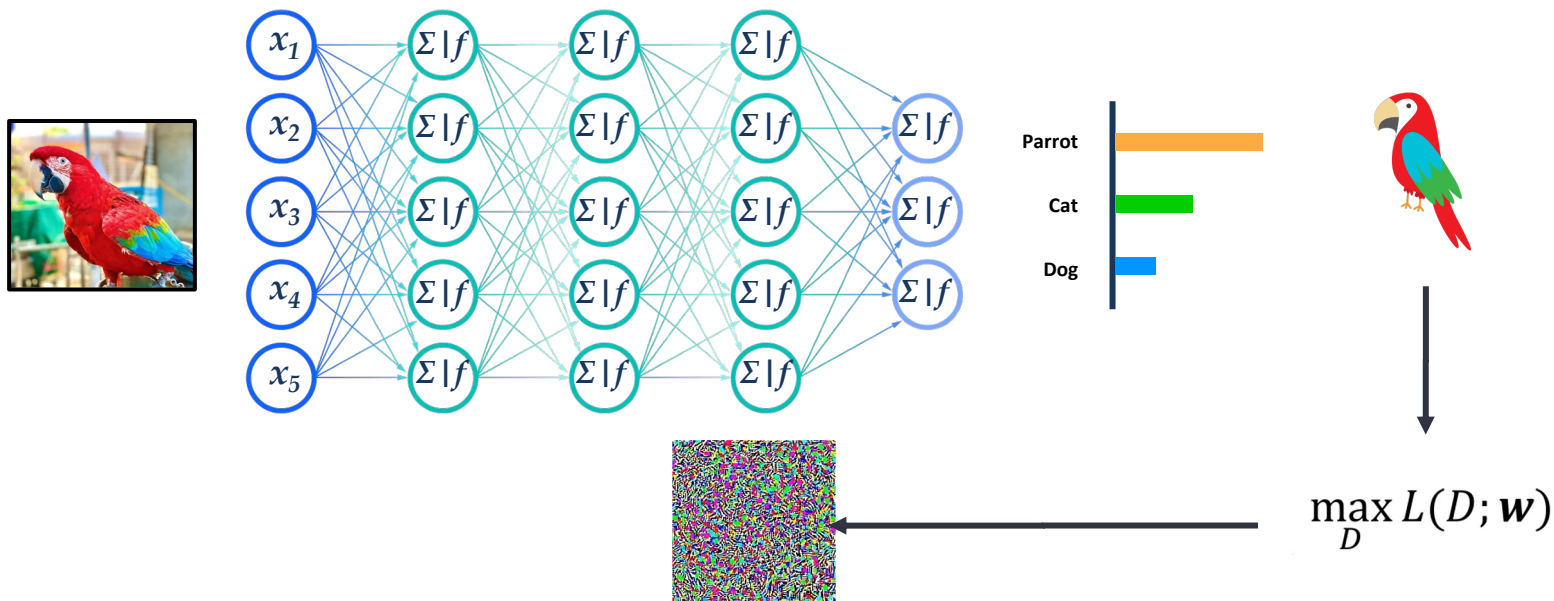
- Demetrio, Biggio, Roli et al., *Adversarial EXEmples: ...*, ACM TOPS 2021
- Demetrio, Biggio, Roli et al., *Functionality-preserving black-box optimization of adversarial windows malware*, IEEE TIFS 2021
- Demontis, Biggio, Roli et al., *Yes, Machine Learning Can Be More Secure!...*, IEEE TDSC 2019

# Adversarial Examples (AdvX)

# Adversarial Examples (AdvX)



$$\max_{D} L(D; \boldsymbol{w})$$

Biggio et al., *Evasion Attacks Against Machine Learning at Test Time*, ECML PKDD 2013
Szegedy et al., *Intriguing Properties of Neural Networks*, ICLR 2014

# Adversarial Examples (AdvX)

@maurapintor

Biggio et al., *Evasion Attacks Against Machine Learning at Test Time*, ECML PKDD 2013
Szegedy et al., *Intriguing Properties of Neural Networks*, ICLR 2014

# How to craft AdvXs

**Exhaustive search** → not possible for modern deep learning models
**Empirical evaluation** → attack = optimization problem + solving algorithm

$$\boldsymbol{\delta}^\star \in \arg\min_{\boldsymbol{\delta}} \quad \mathcal{L}(\boldsymbol{x} + \boldsymbol{\delta}, y, \boldsymbol{\theta})$$
$$\text{s.t.} \quad \|\boldsymbol{\delta}\|_p \leq \epsilon$$
$$\boldsymbol{x}_{\text{lb}} \preceq \boldsymbol{x} + \boldsymbol{\delta} \preceq \boldsymbol{x}_{\text{ub}}$$
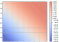
Optimize model's confidence on bad decision

keeping perturbation small

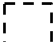and respecting feature space constraints

# How to craft AdvXs

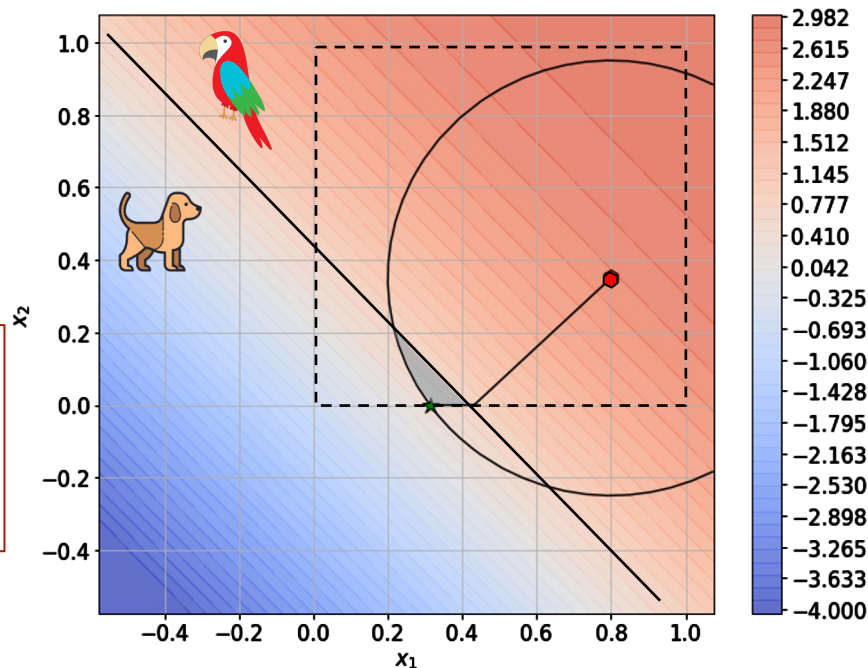**Exhaustive search** → not possible for modern deep learning models
**Empirical evaluation** → attack = optimization problem + solving algorithm

$$\boldsymbol{\delta}^{\star} \in \arg\min_{\boldsymbol{\delta}} \quad \mathcal{L}(\boldsymbol{x} + \boldsymbol{\delta}, y, \boldsymbol{\theta})$$
$$\text{s.t.} \quad \|\boldsymbol{\delta}\|_p \leq \epsilon$$
$$\boldsymbol{x}_{\text{lb}} \preceq \boldsymbol{x} + \boldsymbol{\delta} \preceq \boldsymbol{x}_{\text{ub}}$$

Optimize model's confidence on bad decision

keeping perturbation small ○
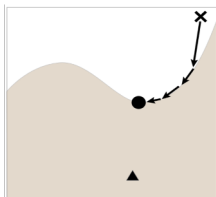
and respecting feature space constraints ⬚

Biggio et al., *Evasion Attacks Against Machine Learning at Test Time*, ECML PKDD 2013
Szegedy et al., *Intriguing Properties of Neural Networks*, ICLR 2014

# How to craft AdvXs



Projected Gradient

$$\boldsymbol{\delta}^\star \in \arg\min_{\boldsymbol{\delta}} \quad \mathcal{L}(\boldsymbol{x} + \boldsymbol{\delta}, y, \boldsymbol{\theta})$$
$$\text{s.t.} \quad \|\boldsymbol{\delta}\|_p \leq \epsilon$$
$$\boldsymbol{x}_{\text{lb}} \preceq \boldsymbol{x} + \boldsymbol{\delta} \preceq \boldsymbol{x}_{\text{ub}}$$

Optimizes confidence

s.t. distance constraint

and feature space constraints

Boundary

$$\boldsymbol{\delta}^\star \in \arg\min_{\boldsymbol{\delta}} \quad \|\boldsymbol{\delta}\|_p$$
$$\text{s.t.} \quad f_y(\boldsymbol{x} + \boldsymbol{\delta}, \boldsymbol{\theta}) \neq f_y(\boldsymbol{x}, \boldsymbol{\theta})$$
$$\boldsymbol{x}_{\text{lb}} \preceq \boldsymbol{x} + \boldsymbol{\delta} \preceq \boldsymbol{x}_{\text{ub}},$$

Find closest advX

s.t. misclassification constraint

and feature space constraints

+ Fast evaluation

- Punctual evaluation (fixed $\epsilon$)

+ Full picture of robustness (boundary)

- Require many iterations

- Difficul to configure properly

Question: How to achieve a **fast**, **reliable**, and **full** evaluation?

# How to craft AdvXs



Pintor et al., *Fast Minimum-norm Adversarial Attacks through Adaptive Norm Constraints*, NeurIPS 2021.

Szegedy et al. ICLR 2014
**FGM**

Moosavi-Dezfooli et al. CVPR 2016
**DeepFool**

Carlini et al. IEEE S&P 2017
**CW**

Madry et al. ICML 2018
**PGD**

Rony et al. CVPR 2019
**DDN**

Croce et al. ICML 2020
**FAB**

Brendel et al. NeurIPS 2020
**BB**

Croce et al. ICML 2020
**AutoAttack**

Pintor et al. NeurIPS 2021
**FMN**

Rony et al. ICCV 2021
**ALMA**

# Bug #1 Slow, hard-to-configure, limited attacks

- **Carlini-Wagner attack (CW)**
  - Requires many steps to converge 🕐
- **Brendel&Bethge attack (BB)**
  - Needs initialization
  - Suffers from poor initialization
  - Complicated steps 🕐
- **Fast Adaptive Boundary (FAB)**
  - Complicated steps 🕐
  - Only untargeted version
- **Decoupling Direction & Norm (DDN)**
  - Specific to L2 norm

🕐 Long runtime

Sensitive to hyperparameters

Limited threat model

# Fix #1: improve current attacks



Goal: find smaller adversarial perturbation with fewer queries to the model

$\| \boldsymbol{\delta} \|_p$ — # of queries

Not adversarial

Slow convergence

Fast convergence

**FMN**

Fast convergence to good local optima

Works in different norms ($\ell_0, \ell_1, \ell_2, \ell\infty$)

Easy tuning / robust to hyperparameter choice

Bigger norm

Smaller norm

# Fast Minimum-norm Adversarial Attacks

**Algorithm 1** Fast Minimum-norm (FMN) Attack

**Input:** $x$, the input sample; $t$, a variable denoting whether the attack is targeted ($t = +1$) or untargeted ($t = -1$); $y$, the target (true) class label if the attack is targeted (untargeted); $\gamma_0$ and $\gamma_K$, the initial and final $\epsilon$-step sizes; $\alpha_0$ and $\alpha_K$, the initial and final $\boldsymbol{\delta}$-step sizes; $K$, the total number of iterations.
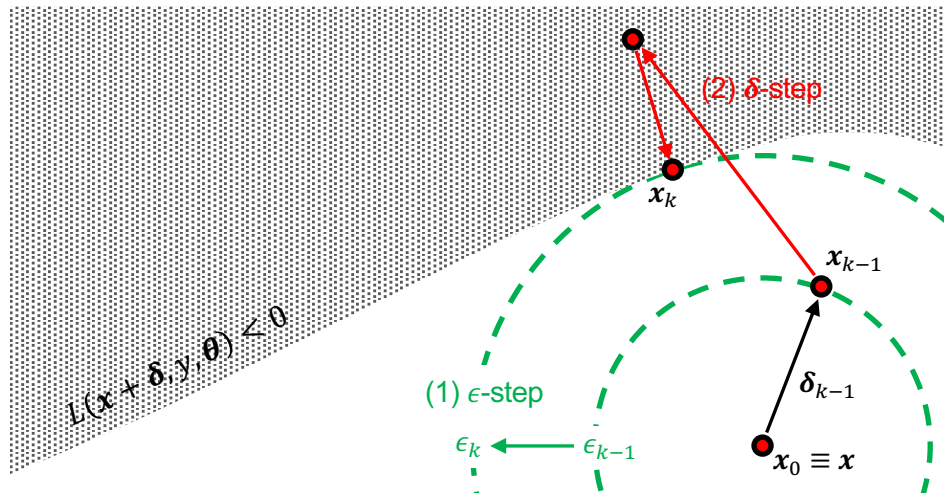
**Output:** The minimum-norm adversarial example $x^\star$.

1: $x_0 \leftarrow x, \epsilon_0 \leftarrow 0, \boldsymbol{\delta}_0 \leftarrow \mathbf{0}, \boldsymbol{\delta}^\star \leftarrow \infty$
2: **for** $k = 1, \ldots, K$ **do**
3:    $g \leftarrow t \cdot \nabla_{\boldsymbol{\delta}} L(x_{k-1} + \boldsymbol{\delta}, y, \boldsymbol{\theta})$   // loss gradient
4:    $\gamma_k \leftarrow h(\gamma_0, \gamma_K, k, K)$   // $\epsilon$-step size decay (Eq. 7)
5:    **if** $L(x_{k-1}, y, \boldsymbol{\theta}) \geq 0$ **then**
6:       $\epsilon_k = \|\boldsymbol{\delta}_{k-1}\|_p + L(x_{k-1}, y, \boldsymbol{\theta})/\|g\|_q$ **if** adversarial not found yet **else** $\epsilon_k = \epsilon_{k-1}(1 + \gamma_k)$
7:    **else**
8:       **if** $\|\boldsymbol{\delta}_{k-1}\|_p \leq \|\boldsymbol{\delta}^\star\|_p$ **then**
9:          $\boldsymbol{\delta}^\star \leftarrow \boldsymbol{\delta}_{k-1}$   // update best min-norm solution
10:       **end if**
11:       $\epsilon_k = \min(\epsilon_{k-1}(1 - \gamma_k), \|\boldsymbol{\delta}^\star\|_p)$
12:    **end if**
13:    $\alpha_k \leftarrow h(\alpha_0, \alpha_K, k, K)$   // $\boldsymbol{\delta}$-step size decay (Eq. 7)
14:    $\boldsymbol{\delta}_k \leftarrow \boldsymbol{\delta}_{k-1} + \alpha_k \cdot g/\|g\|_2$
15:    $\boldsymbol{\delta}_k \leftarrow \Pi_\epsilon(x_0 + \boldsymbol{\delta}_k) - x_0$
16:    $\boldsymbol{\delta}_k \leftarrow \text{clip}(x_0 + \boldsymbol{\delta}_k) - x_0$
17:    $x_k \leftarrow x_0 + \boldsymbol{\delta}_k$
18: **end for**
19: **return** $x^\star \leftarrow x_0 + \boldsymbol{\delta}^\star$

(2) $\boldsymbol{\delta}$-step

$x_k$

$x_{k-1}$

$L(x + \boldsymbol{\delta}, y, \boldsymbol{\theta}) < 0$

(1) $\epsilon$-step

$\boldsymbol{\delta}_{k-1}$

$\epsilon_k \leftarrow \epsilon_{k-1}$

$x_0 \equiv x$

Pintor et al., *Fast Minimum-norm Adversarial Attacks through Adaptive Norm Constraints*, NeurIPS 2021.

# Fast Minimum-norm Adversarial Attacks



MNIST challenge

CIFAR challenge

Pintor et al., *Fast Minimum-norm Adversarial Attacks through Adaptive Norm Constraints*, NeurIPS 2021.

# Let's fix ML Security

Bug #1: slow, hard-to-configure, limited attacks

Fix #1: improve available attacks

# Defending against AdvXs



- Robust training (a.k.a. Adversarial training)

$$\min_{w} \max_{||\delta_i||_\infty \le \epsilon} \sum_i \ell\big(y_i, f_w(x_i + \delta_i)\big)$$

- Detectors



- Ineffective defenses

$g(x)$

Obfuscated gradients do not allow the correct execution of gradient-based attacks...

$x \rightarrow x$'

# The Rise of Adversarial Defenses

Papernot et al. 2016
**Defensive Distillation** (S&P)

Meng et al. 2016
**MagNet defense** (CCS)

Buckman et al. 2016
**Thermometer Encoding** (ICLR)

Roth et al. 2019
**Odds are odd** (ICML)

Pang et al. 2019
**Ensemble diversity** (PMLR)

Yu et al. 2019
**Turning weakness into strength** (NeurIPS)

Xiao et al. 2020
**k-Winner Take All** (ICLR)

# The ~~Rise~~ Fall of Adversarial Defenses



Papernot et al. 2016
**Defensive Distillation** (S&P)

Meng et al. 2016
**MagNet defense** (CCS)

Buckman et al. 2016
**Thermometer Encoding** (ICLR)

Roth et al. 2019
**Odds are odd** (ICML)

Pang et al. 2019
**Ensemble diversity** (PMLR)

Yu et al. 2019
**Turning weakness into strength** (NeurIPS)

Xiao et al. 2020
**k-Winner Take All** (ICLR)

Carlini et al. 2017
**Bypassing ten detection methods** (AISec)

Carlini et al. 2017
**MagNet Not Robust** (arXiv)

Athalye et al. 2018
**Obfuscated gradients give false sense of security** (ICML)

Carlini et al. 2019
**Evaluating Adversarial Robustness** (arXiv)

Tramér et al. 2020
**Adaptive Attacks** (NeurIPS)

○ Proposed defenses
✖ Broken defenses
⌖ Guidelines paper

# Why Is This Happening?

**Root cause: Formal vs Empirical Evaluations**

Formal: no adversarial example in the searched space

Reality: we can only "falsify" the robustness claims by finding adversarial examples
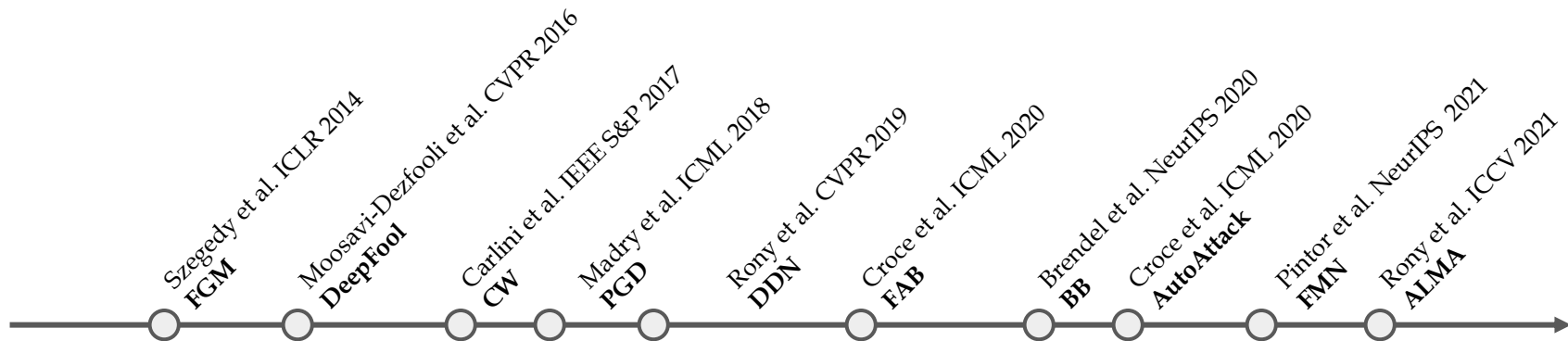
**Similar to finding bugs in software**

What can we say if we did not find adversarial examples?

But no debugging tools for ML robustness

What is the coverage of our tests?

# Bug #2: Lack of debugging tools

Szegedy et al. ICLR 2014
**FGM**

Moosavi-Dezfooli et al. CVPR 2016
**DeepFool**

Carlini et al. IEEE S&P 2017
**CW**

Madry et al. ICML 2018
**PGD**

Rony et al. CVPR 2019
**DDN**

Croce et al. ICML 2020
**FAB**

Brendel et al. NeurIPS 2020
**BB**

Croce et al. ICML 2020
**AutoAttack**

Pintor et al. NeurIPS 2021
**FMN**

Rony et al. ICCV 2021
**ALMA**

# Fix #2: check what your attack is doing

**Profiling attacks**

Pintor et al., *Indicators of Attack Failure*. NeurIPS 2022
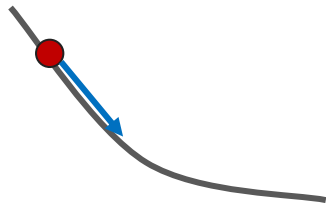
Check your loss

Sanity checks for attacks (Carlini et al. 2019 *Evaluating Adversarial Robustness*, arXiv)

**Goal:** to make security evaluations more trustworthy

# Example: Gradient Obfuscation

**When GD works**

**When GD does not work**

Smooth function: linear approximation works
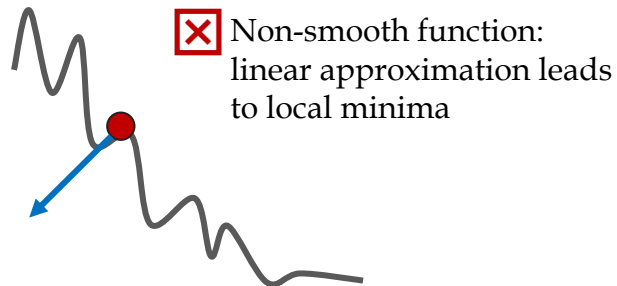


❌ Zero gradients: impossible to find adversarial direction



Check gradient norm

❌ Non-smooth function: linear approximation leads to local minima



Check variability of loss landscape

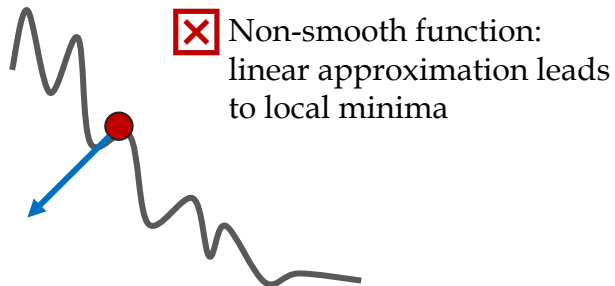# Example: Gradient Obfuscation

**When GD does not work**

 Zero gradients: impossible to find adversarial direction

 Check gradient norm

 Change loss function

 Non-smooth function: linear approximation leads to local minima

 Check variability of loss landscape

 Use smooth approximation

# Attack Failures, Indicators, and Mitigations



**Loss Landscape** *(Obfuscated Gradients)*

**Attack Optimization**

Failures

$F_1$ Shattered Gradients

$F_2$ Stochastic Gradients

$F_3$ **Implementation Errors** *

$F_4$ Non-converging Attack

$F_5$ Non-adaptive Attack

$F_6$ **Unreachable Misclassification** *

# Attack Failures, Indicators, and Mitigations

Pintor et al., *Indicators of Attack Failure: Debugging and Improving Optimization of Adversarial Examples*, NeurIPS 2022

# Attack Failures, Indicators, and Mitigations

# Attack Failures, Indicators, and Mitigations



| Library | Version | GitHub ☆ |
|---------|---------|----------|
| Cleverhans | 4.0.0 | 5.6k |
| ART | 1.11.0 | 3.1k |
| Foolbox | 3.3.3 | 2.3k |
| Torchattacks | 3.2.6 | 984 |

**Loss Landscape** *(Obfuscated Gradients)*

**Attack Optimization**

**Failures**

$F_1$ Shattered Gradients

$F_2$ Stochastic Gradients

$F_3$ Implementation Errors *

$F_4$ Non-converging Attack

$F_5$ Non-adaptive Attack

$F_6$ Unreachable Misclassification *

**Indicators**

$I_1$ Unavailable Gradients *

$I_2$ Unstable Predictions *

$I_3$ Silent Success *

$I_4$ Incomplete Optimization *

$I_5$ Transfer Failure *

$I_6$ Unconstrained Attack Failure *

**Mitigations**

$M_1$ Use BPDA

$M_2$ Use EoT

$M_3$ Fix Attack Implementation *

$M_4$ Tune Step Size and Iterations

$M_5$ Change Loss (*Adaptive*)

$M_6$ Change Loss (*Bad Local Minimum*) *

*Loss/Model-specific fixes to ensure gradients are smooth*

*Attack-specific fixes to ensure attack optimization runs correctly*

# Identifying and Fixing Failures

| Model | Attack | $I_1$ | $I_2$ | $I_3$ | $I_4$ | $I_5$ | $I_6$ | RA |
|-------|--------|-------|-------|-------|-------|-------|-------|-----|
| *DIST* | **Original** | ✓ | | | | | ✓ (10/10) | 0.95 ☒ |
| | **Patched** | | | | | | | **0.01** ☑ |
| *k-WTA* | **Original** | | ✓ (10/10) | ✓ (23%) | ✓ (11%) | | ✓ (4/10) | 0.67 ☒ |
| | **Patched** | | | | ✓ (6%) | | ✓ (2/10) | **0.09** ☑ |

The evaluations that we identified as faulty trigger our indicators
**+ additional results in the paper!**

# Detecting Unreliable Evaluations

We evaluated 6 defenses recently published on top-tier venues, available through RobustBench

They have been tested with **AutoAttack** a <u>SOTA parameter-free attack</u>

We show that these evaluations are unreliable



**ROBUSTBENCH**  Leaderboards  Paper  FAQ  Contribute  Model Zoo 🚀

Available Leaderboards

CIFAR-10 ($\ell_\infty$)  CIFAR-10 ($\ell_2$)  CIFAR-10 (Corruptions)  CIFAR-100 ($\ell_\infty$)  CIFAR-100 (Corruptions)  ImageNet ($\ell_\infty$)

ImageNet (Corruptions: IN-C, IN-3DCC)

Leaderboard: CIFAR-10, $\ell_\infty = 8/255$, untargeted attack

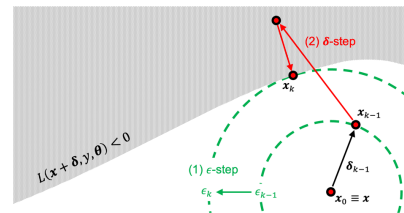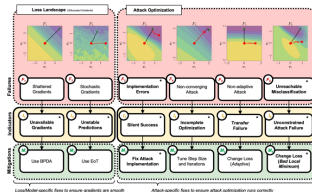Show 15 entries                                                                 Search: Papers, architectures

| Rank | Method | Standard accuracy | AutoAttack robust accuracy | Best known robust accuracy | AA eval. potentially unreliable | Extra data | Architecture | Venue |
|---|---|---|---|---|---|---|---|---|
| 1 | Fixing Data Augmentation to Improve Adversarial Robustness<br>66.56% robust accuracy is due to the original evaluation (AutoAttack + MultiTargeted) | 92.23% | 66.58% | 66.56% | ✗ | ☑ | WideResNet-70-16 | arXiv, Mar 2021 |
| 2 | Improving Robustness using Generated Data<br>It uses additional 100M synthetic images in training. 66.10% robust accuracy is due to the original evaluation (AutoAttack + MultiTargeted) | 88.74% | 66.11% | 66.10% | ✗ | ✗ | WideResNet-70-16 | NeurIPS 2021 |
| 3 | Uncovering the Limits of Adversarial Training against Norm-Bounded Adversarial Examples<br>65.87% robust accuracy is due to the original evaluation (AutoAttack + MultiTargeted) | 91.10% | 65.88% | 65.87% | ✗ | ☑ | WideResNet-70-16 | arXiv, Oct 2020 |
| 4 | Fixing Data Augmentation to Improve Adversarial Robustness<br>It uses additional 1M synthetic images in training. 64.58% robust accuracy is due to the original evaluation (AutoAttack + MultiTargeted) | 88.50% | 64.64% | 64.58% | ✗ | ✗ | WideResNet-106-16 | arXiv, Mar 2021 |

https://robustbench.github.io

# Let's fix ML Security

Bug #1: slow, hard-to-configure, limited attacks

Fix #1: improve available attacks



Bug #2: lack of debugging tools for ML Security

Fix #2: develop tests and track metrics on the attacks

# Bug # 3: Meet the Real World

**Adversarial perturbations are usually crafted in the ideal situation**
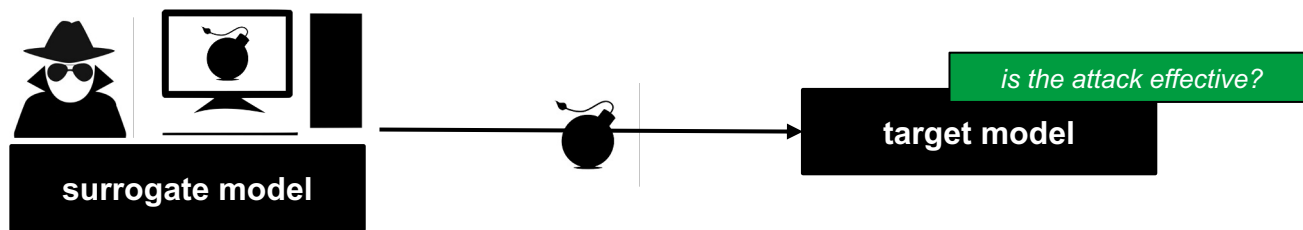
Challenges:
- the model might be unknown / not accessible
- the perturbation must respect the rules of the real world

How to evaluate robustness in the physical world?

# Fix # 3: Beyond white-box evaluations

**Transferability:** the ability of an attack, crafted against a **surrogate** model, to be effective against a different, *unknown* **target** model
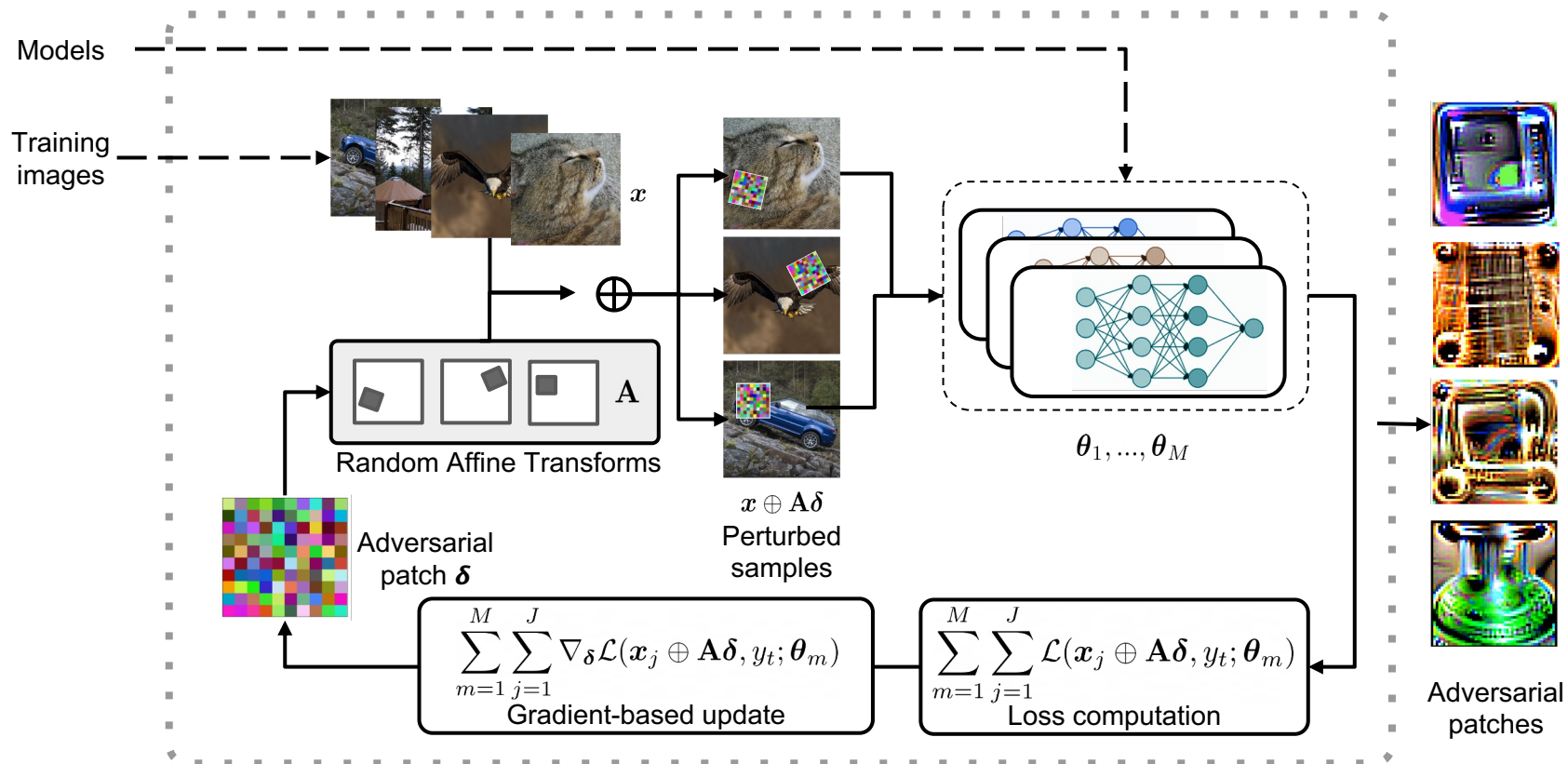


We propose three metrics that clarify the underlying factors behind transferability and allow highlighting interesting connections with model complexity

**Key insights:**

**- gradient alignment and smoothness of surrogate improves transferability**

Papernot et al., *Practical Black-Box Attacks against Machine Learning*, ASIACCS 2017
Demontis et al., *Why Do Adversarial Attacks Transfer?* USENIX Security 2019

Models

Training images

$x$

Random Affine Transforms

$\mathbf{A}$

Adversarial patch $\boldsymbol{\delta}$

$x \oplus \mathbf{A}\boldsymbol{\delta}$
Perturbed samples

$\boldsymbol{\theta}_1, ..., \boldsymbol{\theta}_M$

$$\sum_{m=1}^{M} \sum_{j=1}^{J} \nabla_{\boldsymbol{\delta}} \mathcal{L}(x_j \oplus \mathbf{A}\boldsymbol{\delta}, y_t; \boldsymbol{\theta}_m)$$
Gradient-based update

$$\sum_{m=1}^{M} \sum_{j=1}^{J} \mathcal{L}(x_j \oplus \mathbf{A}\boldsymbol{\delta}, y_t; \boldsymbol{\theta}_m)$$
Loss computation

Adversarial patches

# Beyond White-box Evaluations: Creating Real-world Attacks



banana

banana    banana    banana

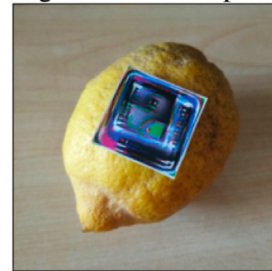From the digital world …

… to the physical world

True label: joystick
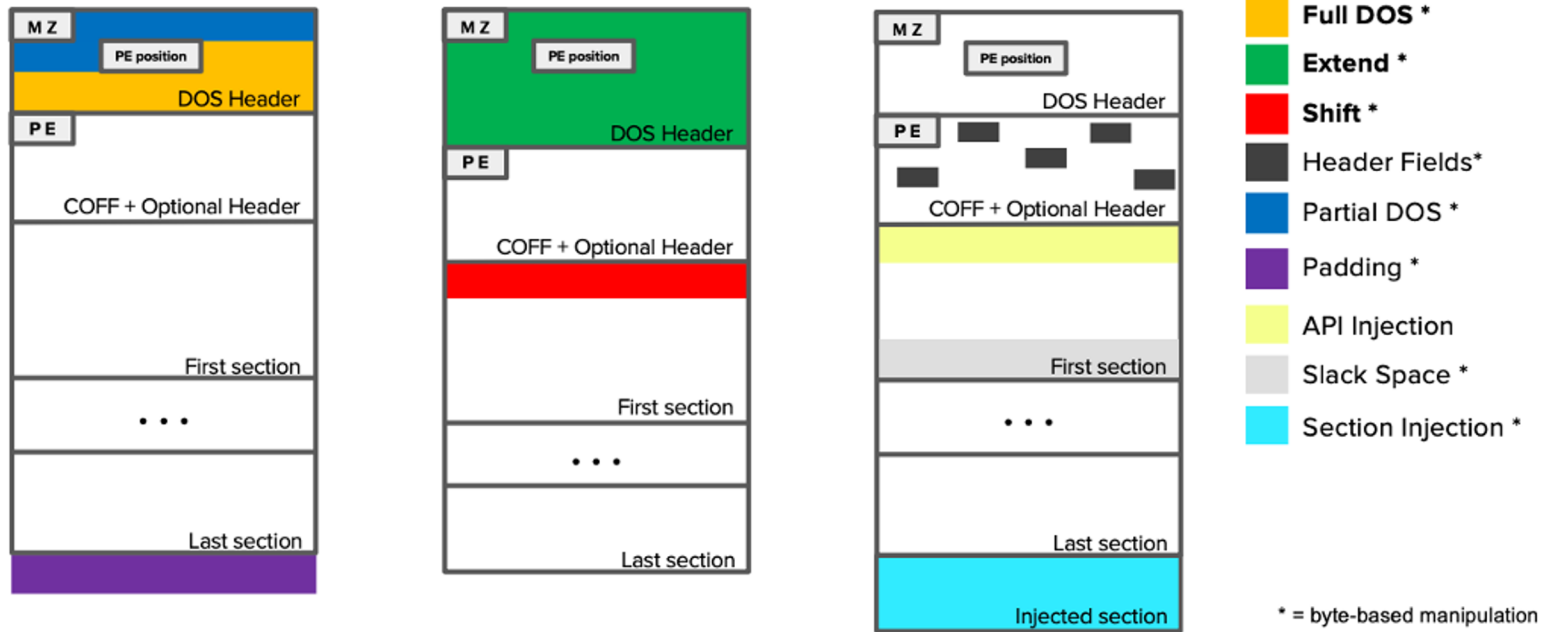Target: electric guitar

True label: sandal
Target: banana

True label: lemon
Target: cellular telephone

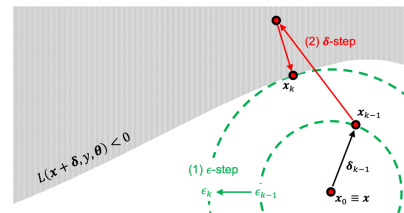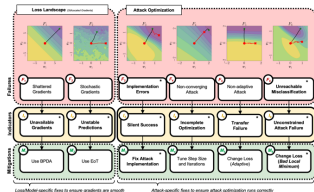# Adversarial EXEmples: Practical Attacks on Machine Learning for Windows Malware Detection



Demetrio, Biggio, et al., *Adversarial EXEmples,* ACM TOPS 2021
Demetrio, Biggio, et al., *Functionality-preserving ...,* IEEE TIFS 2021

# Let's fix ML Security

Bug #1: slow, hard-to-configure, limited attacks

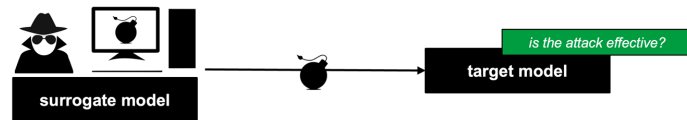Fix #1: improve available attacks





Bug #2: lack of debugging tools for ML Security

Fix #2: develop tests and track metrics on the attacks

Bug #3: Keep in mind the real world

Fix #3: create strong and realizable attacks



is the attack effective?

target model

surrogate model

# Provocations



Do we want to spend the next 10 years like this?



Will this problem even be relevant in 10 years?

# Machine Learning is deployed in the real world



**Induced hallucinations**
Research clearly shows that it is possible to target machine learning models with practical attacks that spoil its performances

**Many threats**
Test-time perturbations, dataset poisoning, privacy leaks, and many many others
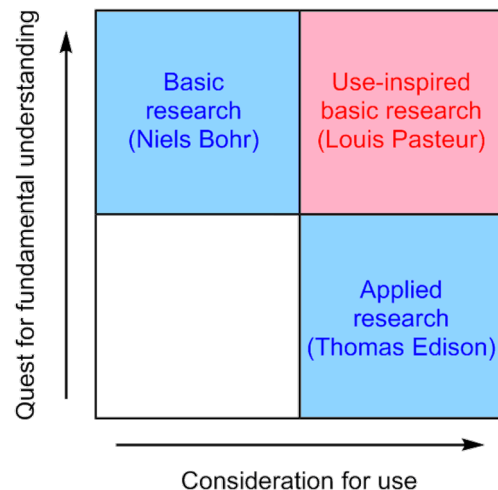
# Use-Inspired Basic Research Questions
*Looking at the Pasteur's Quadrant*

If evidence of optimized attacks against AI/ML remains unclear, what will be the future of MLSec as a research field?

Can we use MLSec to help solve some of today's industrial challenges?

- To improve robustness/accuracy over time, requiring less frequent retraining
- To improve maintainability and interpretability of deployed models (update procedures)
- To learn reliably from noisy/incomplete labeled datasets

Will we be able to build more reliable and practical ML models using MLSec / AdvML?



Quest for fundamental understanding

Basic research (Niels Bohr)

Use-inspired basic research (Louis Pasteur)

Applied research (Thomas Edison)

Consideration for use

# MLSec
# Seminar Series

🐦 **@mlsec_lab**

https://pralab.github.io/mlsec/

# Thanks!

✉  maura.pintor@unica.it

🐦  @maurapintor

💻  maurapintor.github.io

Special thanks to Battista Biggio, Antonio Emanuele Cinà, and Luca Demetrio for sharing with me some of the material used in these slides.