



University of  
Cagliari, Italy

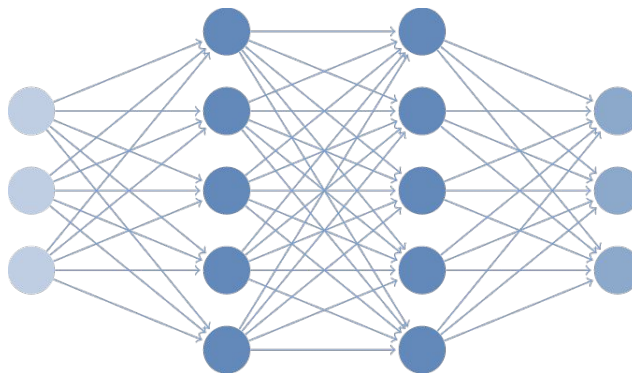
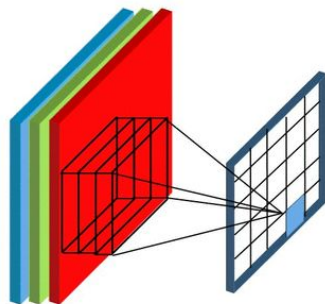
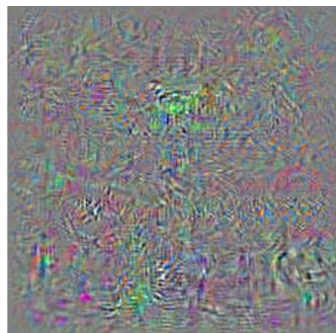
# Indicators of Attack Failure: Visualizing and Debugging Optimization of Adversarial Examples

**Maura Pintor**, Luca Demetrio, Angelo Sotgiu, Giovanni Manca, Ambra Demontis,  
Nicholas Carlini, Battista Biggio, Fabio Roli

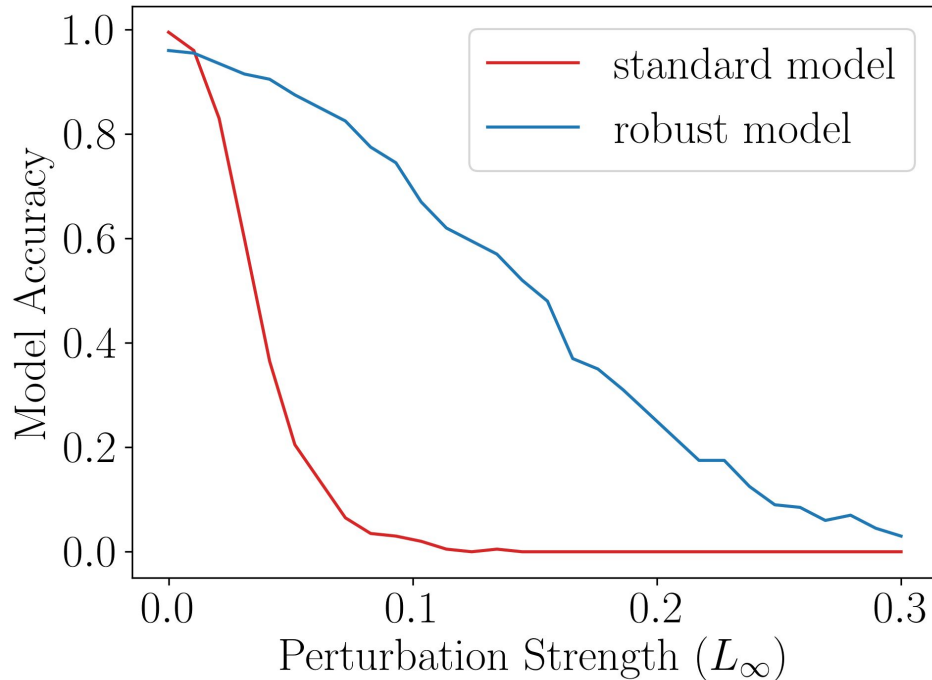
# Adversarial Examples



+



# Adversarial Robustness

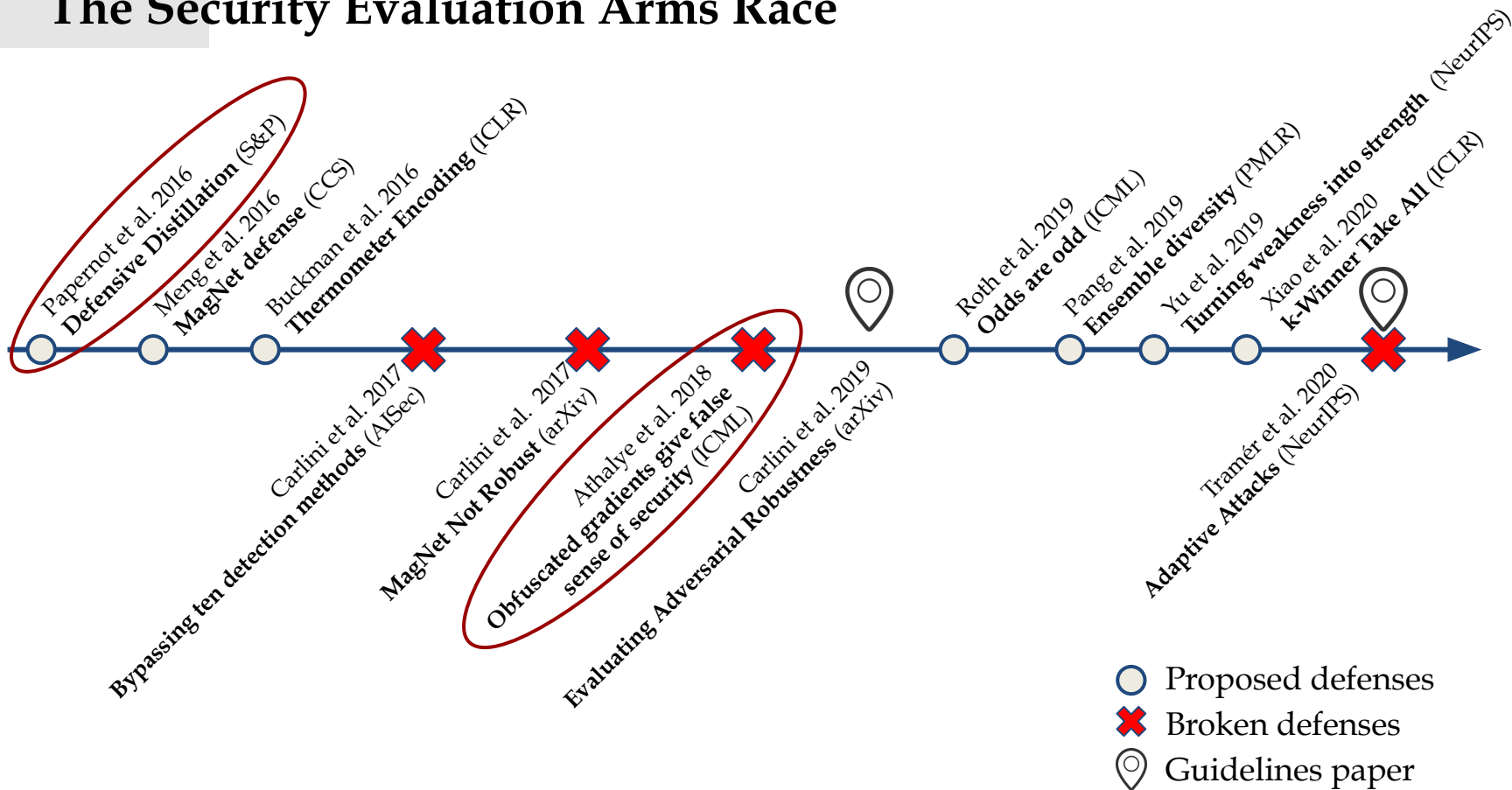


Evaluating **adversarial robustness** amounts to finding adversarial examples with a given **perturbation budget**

We have to rely on **empirical** evaluation

**What is the problem of using empirical evaluations?**

# The Security Evaluation Arms Race



# Robustness evaluations and pitfalls

The security evaluation is a way to estimate the **robust accuracy** of a model

Literature provides approaches to improve robustness evaluations

**Limits:** manual process, qualitative metrics, only suggestions and “best practices”

Currently, there is **no debugging tool** for adversarial attacks

Our **contributions:**

- Different **point of view** for debugging a security evaluation: **per-point, inspect path**

- Provide **computable indicators**

- Propose **systematic protocol** for improving security evaluations



**Where can attacks fail and how can we detect the failures?**

# Pitfalls of Gradient-based Attacks

---

**Algorithm 1:** General formalization for *untargeted* (Equation 1 and 2) and *targeted* attacks (Equation 3 and 4)

---

**Input** :  $\mathbf{x}$ , the initial point;  $y_t$ , the target (true) class label if the attack is targeted (untargeted);  $n$ , the number of iterations;  $\alpha$ , the learning rate;  $f$ , the target model;  $(\mathbf{x}_{lb}, \mathbf{x}_{ub})$ , the bounds of the input space;  $\Delta$ , the considered region.

**Output** :  $\mathbf{x}^*$ , the solution found by the algorithm

```
1  $\mathbf{x}_0 \leftarrow \text{init}(\mathbf{x})$  ▷ Initialize starting point
2  $\hat{\theta} \leftarrow \text{approximation}(\theta)$  ▷ Approximate model's parameters
3  $\delta_0 \leftarrow \mathbf{0}$  ▷ Initial  $\delta$ 
4 for  $i \in [1, n]$  do
5    $\delta' \leftarrow \delta_i + \alpha \nabla_{\mathbf{x}_i} L(\mathbf{x}_0 + \delta_i, y_t; \hat{\theta})$  ▷ Compute optimizer step
6    $\delta_{i+1} \leftarrow \text{apply\_constraints}(\mathbf{x}_0, \delta', \Delta, \mathbf{x}_{lb}, \mathbf{x}_{ub})$  ▷ Apply constraints
7  $\delta^* \leftarrow \text{best}(\delta_0, \dots, \delta_n)$  ▷ Choose best perturbation
8 return  $\delta^*$ 
```

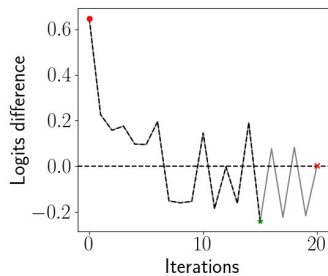
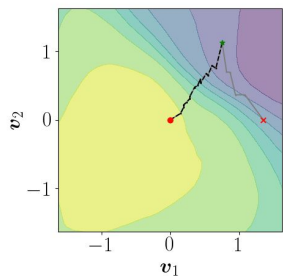
---



# Attack Failures

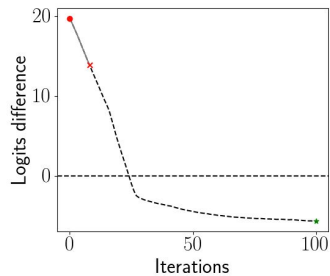
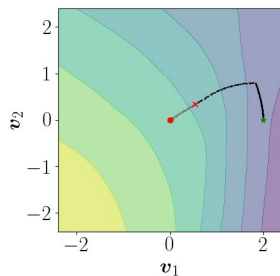
Bad implementation

$$\delta^* \leftarrow \text{best}(\delta_0, \dots, \delta_n)$$



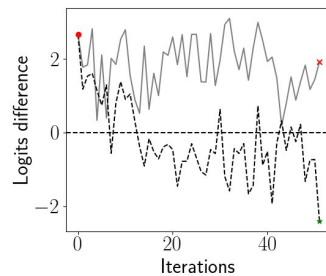
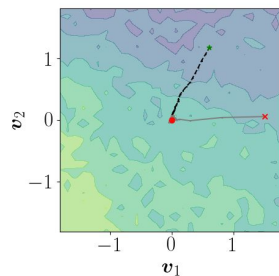
Attack is not converging

```
3  $\delta_0 \leftarrow \text{init}(\theta)$   
4 for  $i \in [1, n]$  do  
5    $\delta'_i \leftarrow \delta_i + \alpha \nabla_{x_i} L(x_0 + \delta_i, y_i; \hat{\theta})$ 
```



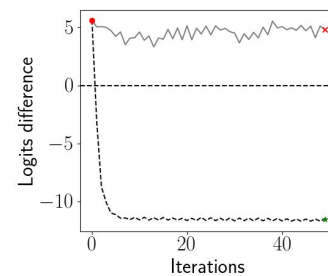
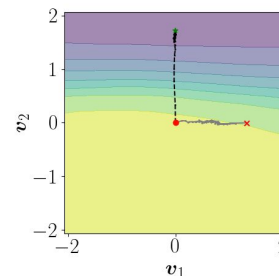
Bad local optimum

```
1  $x_0 \leftarrow \text{init}(x)$   
2  $\hat{\theta} \leftarrow \text{approximation}(\theta)$ 
```



Attack is not adaptive

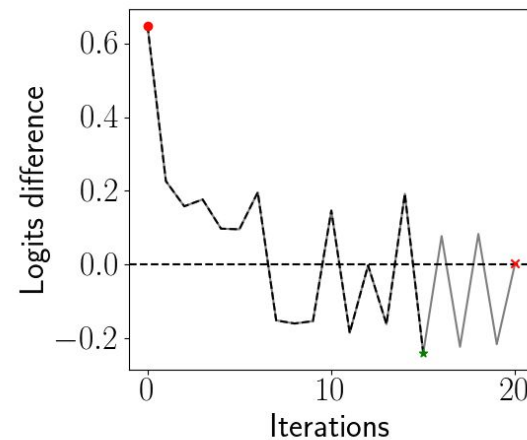
```
2  $\hat{\theta} \leftarrow \text{approximation}(\theta)$ 
```



# Silent success

Measures if the attack is returning a **wrong result**

Computed by looking inside the **optimization path**



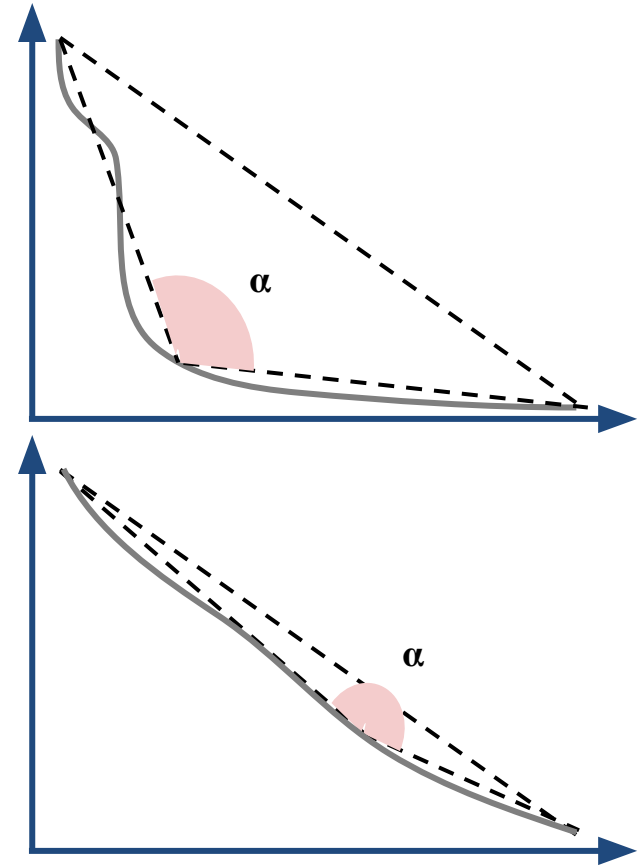
# Break-point angle

Measures the alignment between the improvement of the loss over the iterations and the expected decreasing behaviour

Computed as **the absolute value of  $\cos \alpha$**

$90^\circ$  → loss has the correct shape

$180^\circ$  → loss is still decreasing

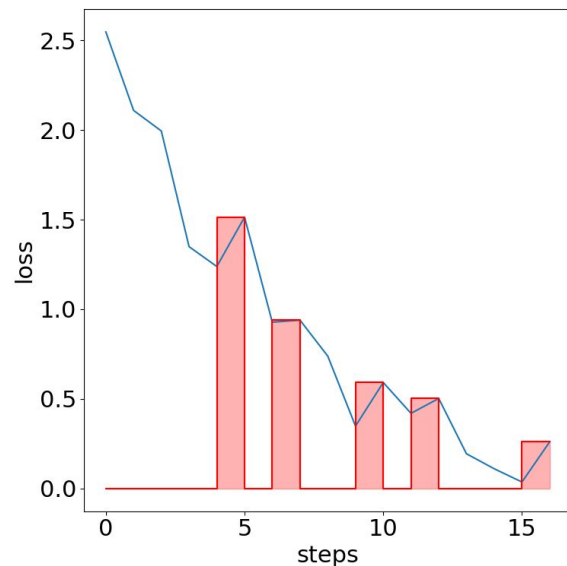


# Increasing loss

Measures the increment of the loss due to jumps inside the space

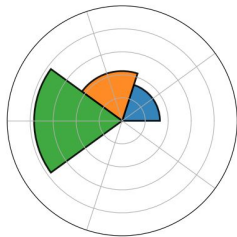
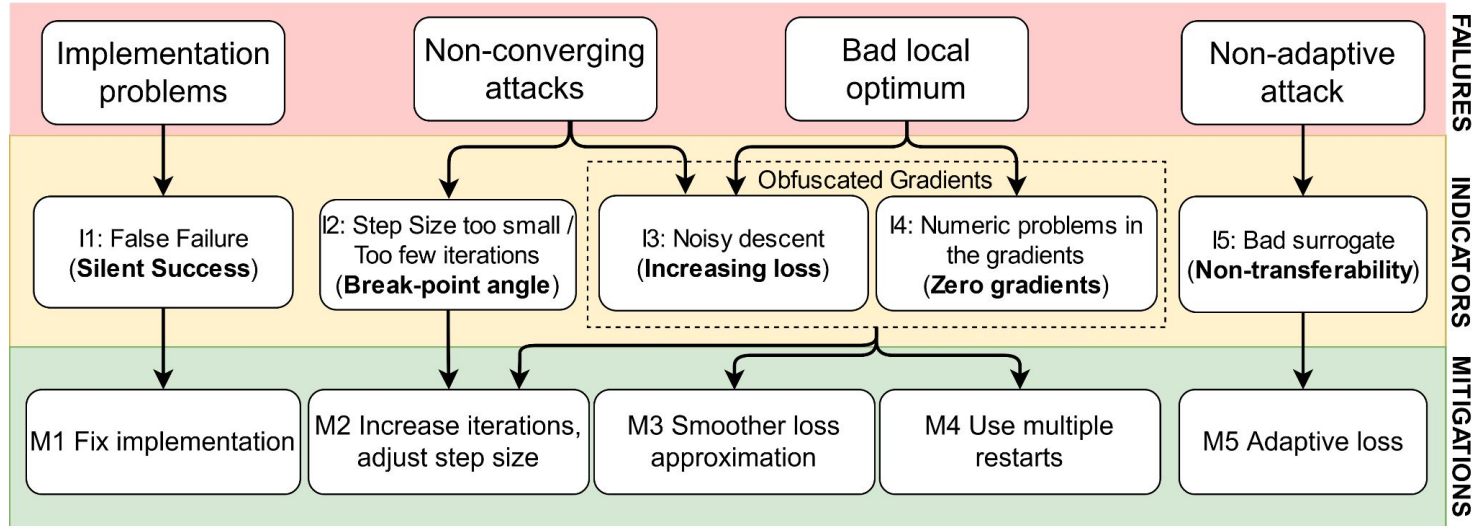
Computed as **area under positive contributions in the attack loss**

If indicator is close to 1, it means the loss **is not decreasing consistently**



**And how do we fix them?**

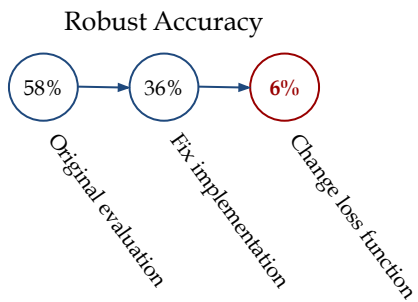
# Indicators and mitigations



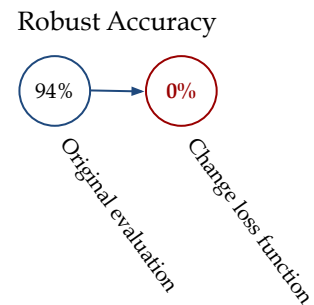
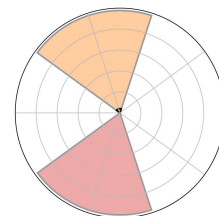
These charts help us understand when **indicators** are triggered, and which **mitigations** to apply

# Experiments

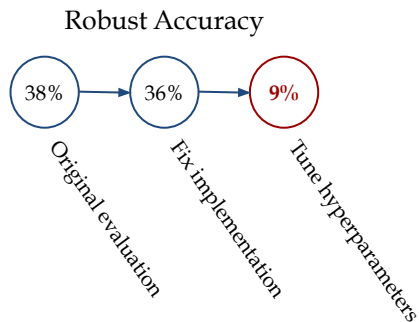
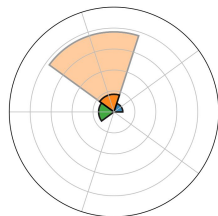
k-Winners  
Take All



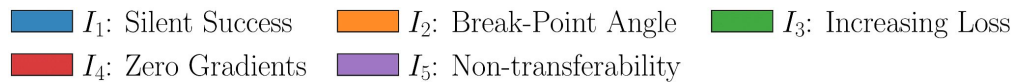
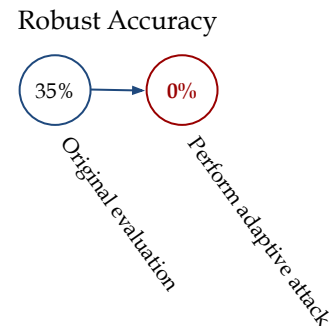
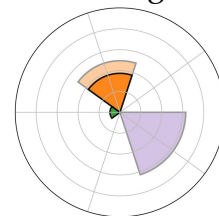
Distillation



Ensemble Diversity



Turning a  
Weakness into a  
Strength



# Conclusions

- Enable debugging faulty-conducted security evaluations
- Empirical evaluation in 4 case-studies
- Paper available <https://arxiv.org/abs/2106.09947>
- Open source code <https://github.com/pralab/IndicatorsOfAttackFailure>

## Future works

- Integrate in benchmarks
- Further automatization - MLSecOps



<https://secml.gitlab.io/>

