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

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

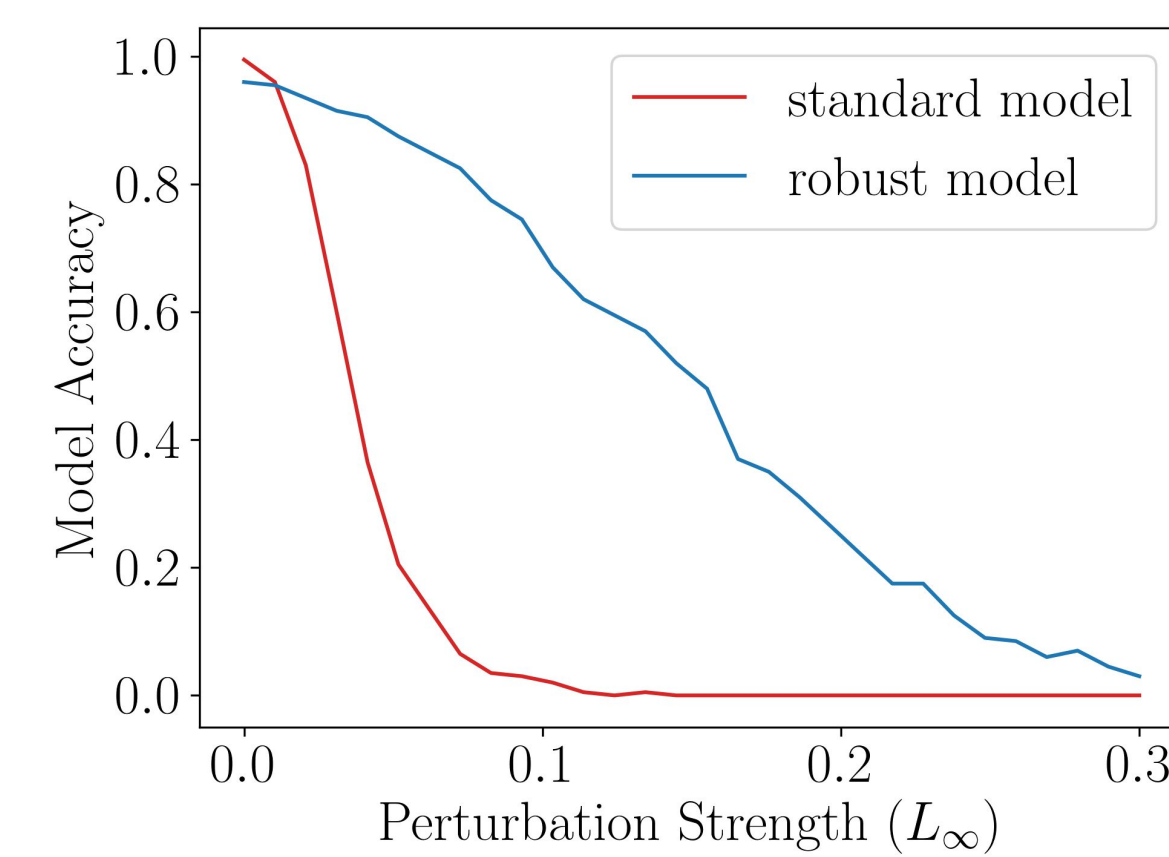*"A Blessing in Disguise: The Prospects and Perils of Adversarial Machine Learning" ICML workshop 2021*

**Pattern Recognition and Application Laboratory**

**Department of Electrical and Electronic Engineering**
**University of Cagliari, Italy**

Pluribus One — *seeing one in many*

Google

---

## Adversarial Robustness Evaluation

**Goal:**
➤ Find adversarial examples with a given perturbation budget
➤ Evaluate the robust accuracy

**Problems:**
➤ We have to rely on empirical evaluations
➤ Attacks often fail
➤ *False sense of security*
➤ Hard to fix: only guidelines but no practical debugging tools available!



---

## Gradient-based Attacks

➤ General formalization for untargeted and targeted attacks
➤ We highlight steps related to common failures

**Input** : $x$, 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; $(x_{lb}, x_{ub})$, the bounds of the input space; $\Delta$, the considered region.
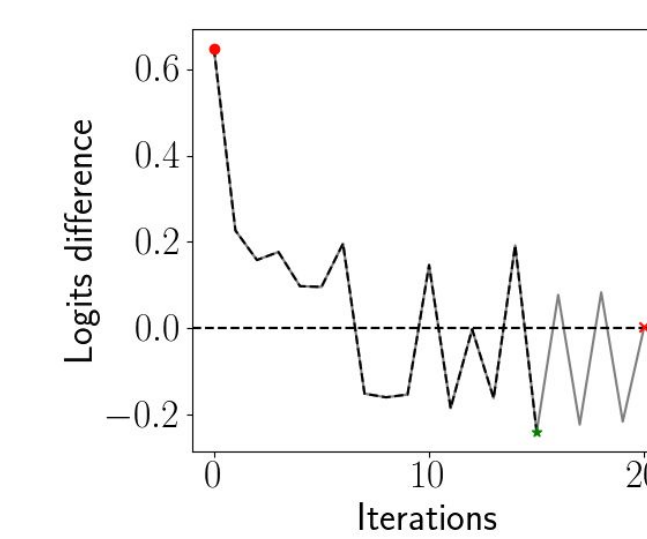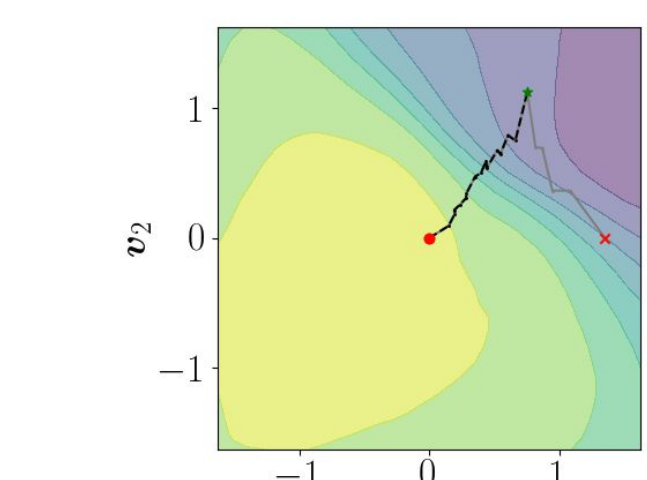**Output** : $x^\star$, the solution found by the algorithm

1  $x_0 \leftarrow init(x)$                          ▷ Initialize starting point
2  $\theta \leftarrow approximation(\theta)$          ▷ Approximate model's parameters
3  $\delta_0 \leftarrow 0$                            ▷ Initial $\delta$
4  **for** $i \in [1, n]$ **do**
5      $\delta' \leftarrow \delta_i + \alpha \nabla_{x_i} L(x_0 + \delta_i, y_t; \theta)$   ▷ Compute optimizer step
6      $\delta_{i+1} \leftarrow apply\_constraints(x_0, \delta', \Delta, x_{lb}, x_{ub})$   ▷ Apply constraints
7  $\delta^\star \leftarrow best(\delta_0, ..., \delta_n)$   ▷ Choose best perturbation
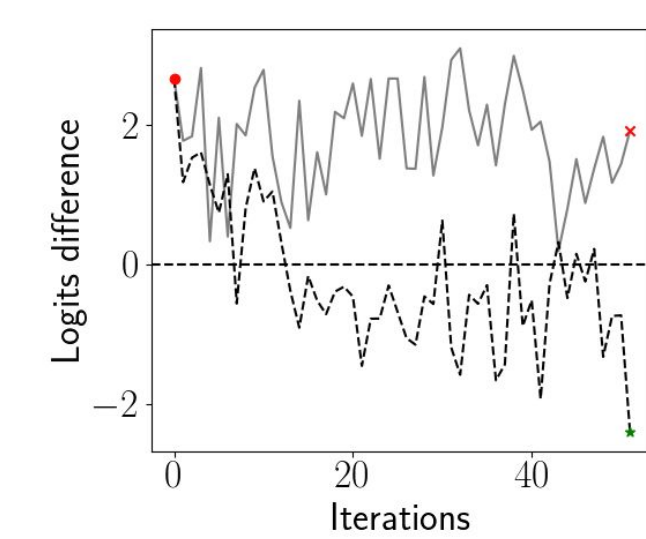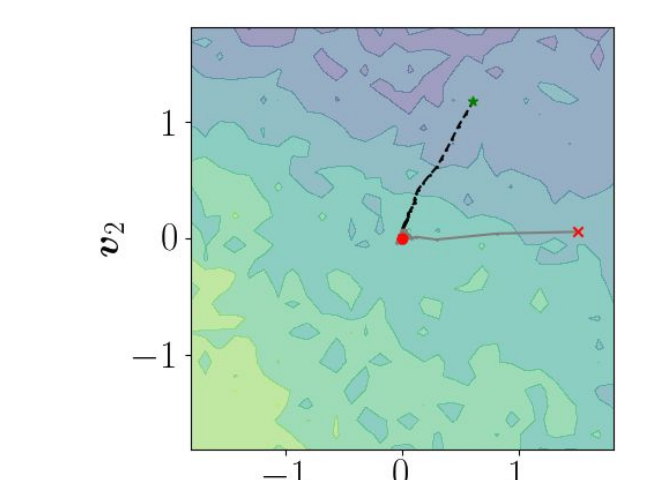8  **return** $\delta^\star$

---

## Attack failures

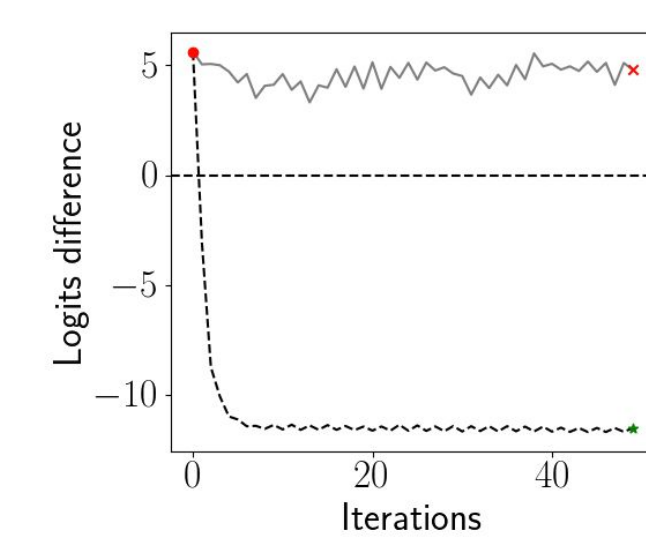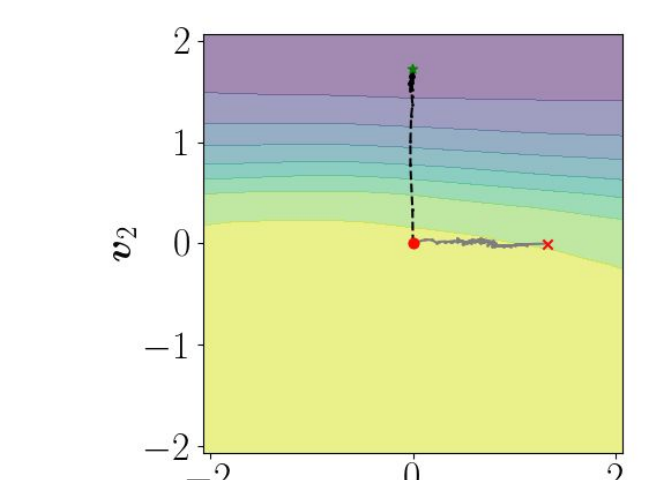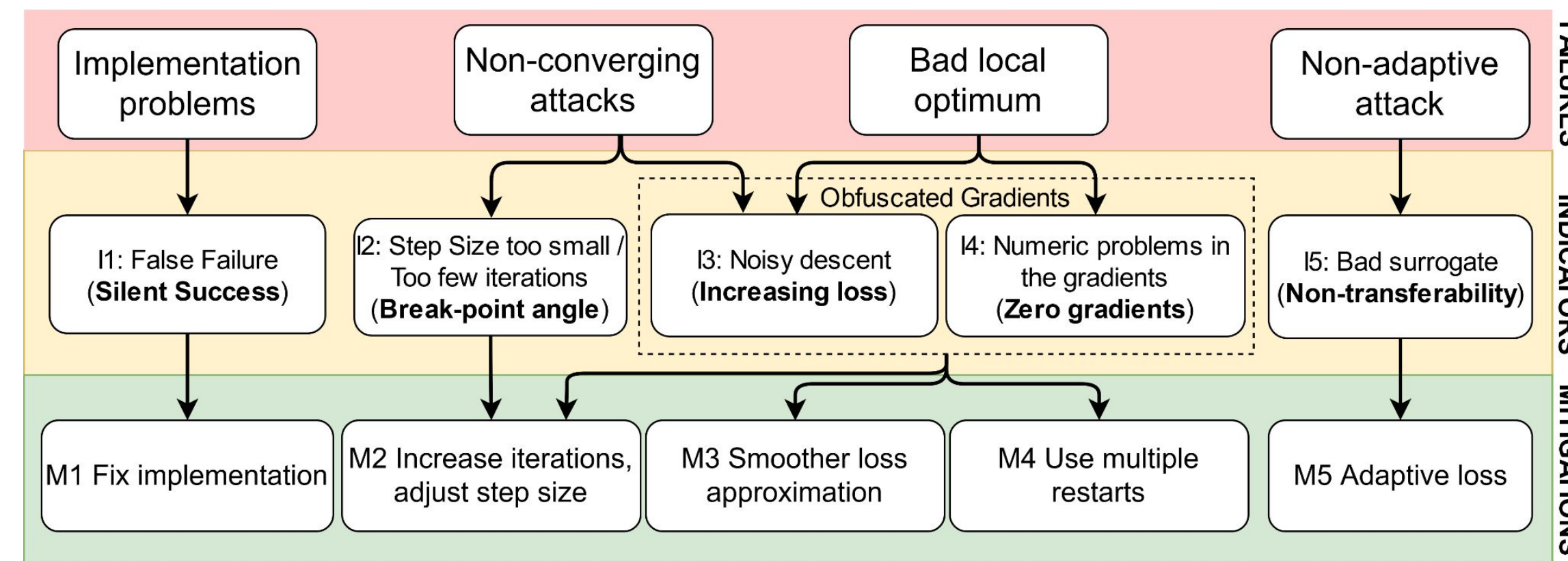**Bad implementation (step 7)**   **Attack is not converging (step 4-5)**   **Bad local optimum (step 1-2)**   **Attack is not adaptive (step 2)**
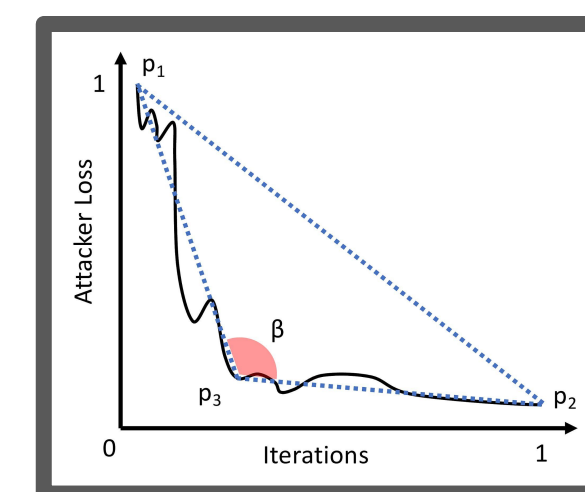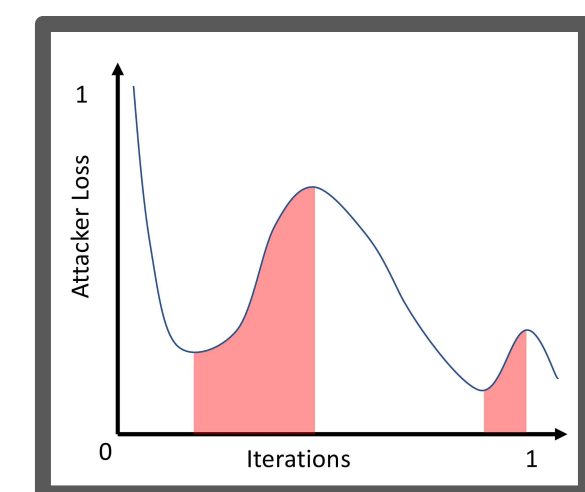


---

## Indicators and mitigations

➤ We formulate 5 quantitative indicators (all in [0, 1])
➤ Each indicator is related to one or more failure
➤ We also propose 5 mitigations to apply, based on indicators results



**Break point angle**

$1 - |cos\,\beta|$

the attack loss is normalized

**Increasing loss**

Area under the attack loss if increasing
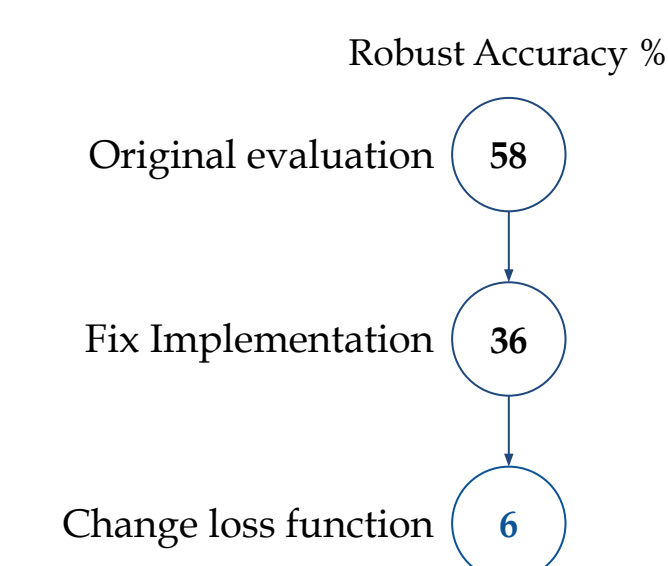


---

## Experiments

**Setting:**
➤ We select 4 defenses with reported failures
➤ We evaluate our indicators
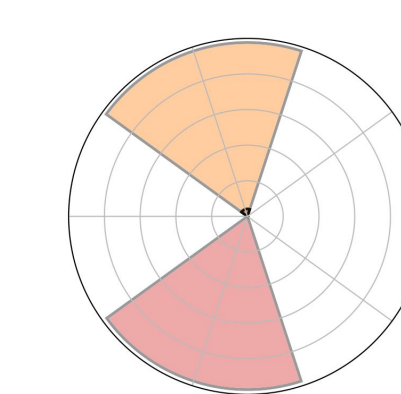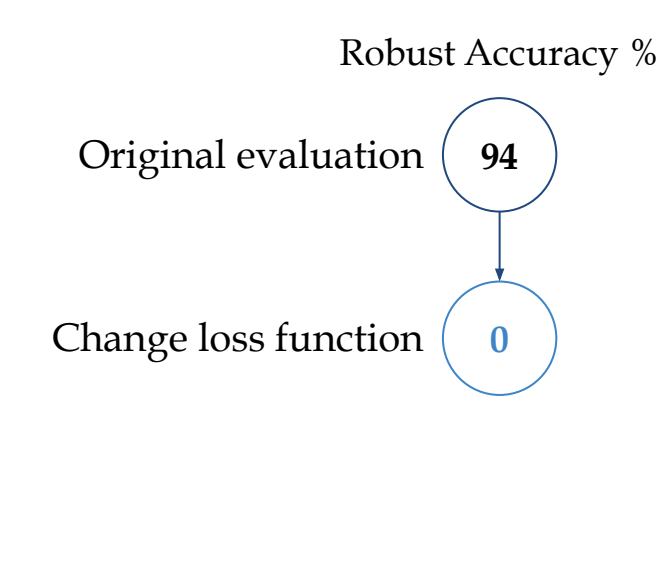➤ We apply mitigations

**Results:**
➤ Indicators correctly reveal the "false sense of security"
➤ Patched attacks drop robust accuracy
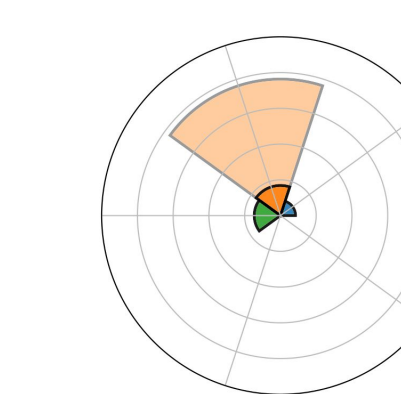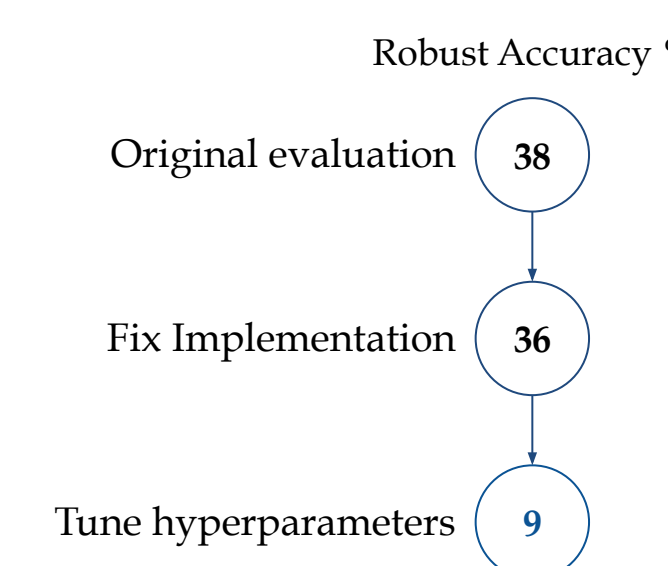➤ Indicators are strongly correlated with attacks performance
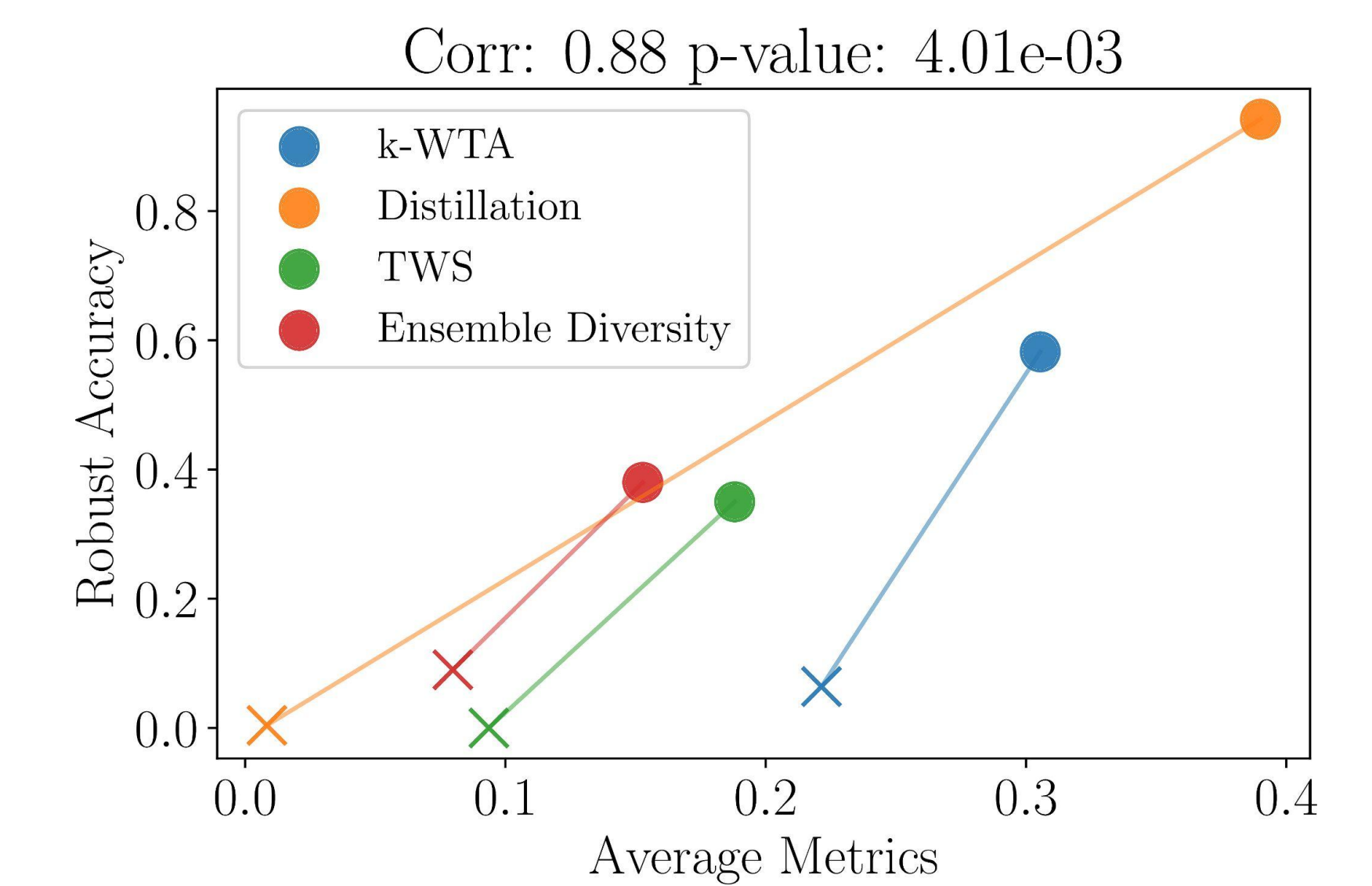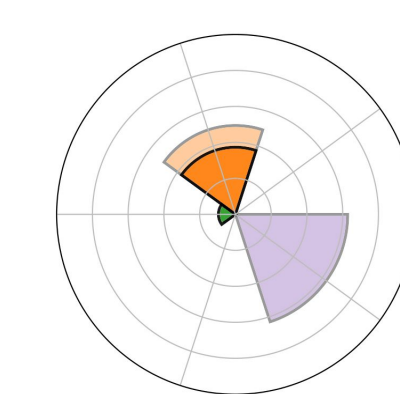
**k-WinnersTake All**

Robust Accuracy %
Original evaluation 58
Fix Implementation 36
Change loss function 6

**Distillation**

Robust Accuracy %
Original evaluation 94
Change loss function 0

**Ensemble Diversity**

Robust Accuracy %
Original evaluation 38
Fix Implementation 36
Tune hyperparameters 9

**Turning a Weakness into a Strength**

Robust Accuracy %
Original evaluation 35
Perform adaptive attack 0

Corr: 0.88 p-value: 4.01e-03



$I_1$: Silent Success   $I_2$: Break-Point Angle   $I_3$: Increasing Loss
$I_4$: Zero Gradients   $I_5$: Non-transferability

---

## Useful links and implementations

➤ Open source code https://github.com/pralab/IndicatorsOfAttackFailure
➤ Paper available https://arxiv.org/abs/2106.09947
➤ Implemented with SecML

SECML
https://secml.gitlab.io/
*Twitter*: @secml_py

---

## Key Takeaways

➤ Unified framework for gradient-based attacks and categorization of main failures
➤ Framework for debugging faulty-conducted security evaluations with quantitative indicators and mitigations strategies
➤ Empirical evaluation on 4 case-studies
  ◆ indicators highlight failures
  ◆ mitigations improve the robustness evaluation

---

## Future Work

➤ Integration in benchmarks
➤ Add more indicators
➤ Further automatization
  ◆ Towards MLSecOps

---

**References**
[1] Tramèr, F., Carlini, N., Brendel, W., & Madry, A. (2020). **On Adaptive Attacks to Adversarial Example Defenses.** *ArXiv, abs/2002.08347.*
[2] Xiao, C., Zhong, P., & Zheng, C. (2019). **Resisting Adversarial Attacks by k-Winners-Take-All.** *ArXiv, abs/1905.10510.*
[3] Papernot, N., Mcdaniel, P., Wu, X., Jha, S., & Swami, A. (2016). **Distillation as a Defense to Adversarial Perturbations Against Deep Neural Networks.** *2016 IEEE Symposium on Security and Privacy (SP), 582-597.*
[4] Pang, T., Xu, K., Du, C., Chen, N., & Zhu, J. (2019). **Improving Adversarial Robustness via Promoting Ensemble Diversity.** *ArXiv, abs/1901.08846.*
[5] Yu, T., Hu, S., Guo, C., Chao, W., & Weinberger, K.Q. (2019). **A New Defense Against Adversarial Images: Turning a Weakness into a Strength.** *ArXiv, abs/1910.07629.*
[6] Melis, M., Demontis, A., Pintor, M., Sotgiu, A., & Biggio, B. (2019). **secml: A Python Library for Secure and Explainable Machine Learning.** *ArXiv, abs/1912.10013.*