# TransferBench: Benchmarking Ensemble-based Black-box Transfer Attacks
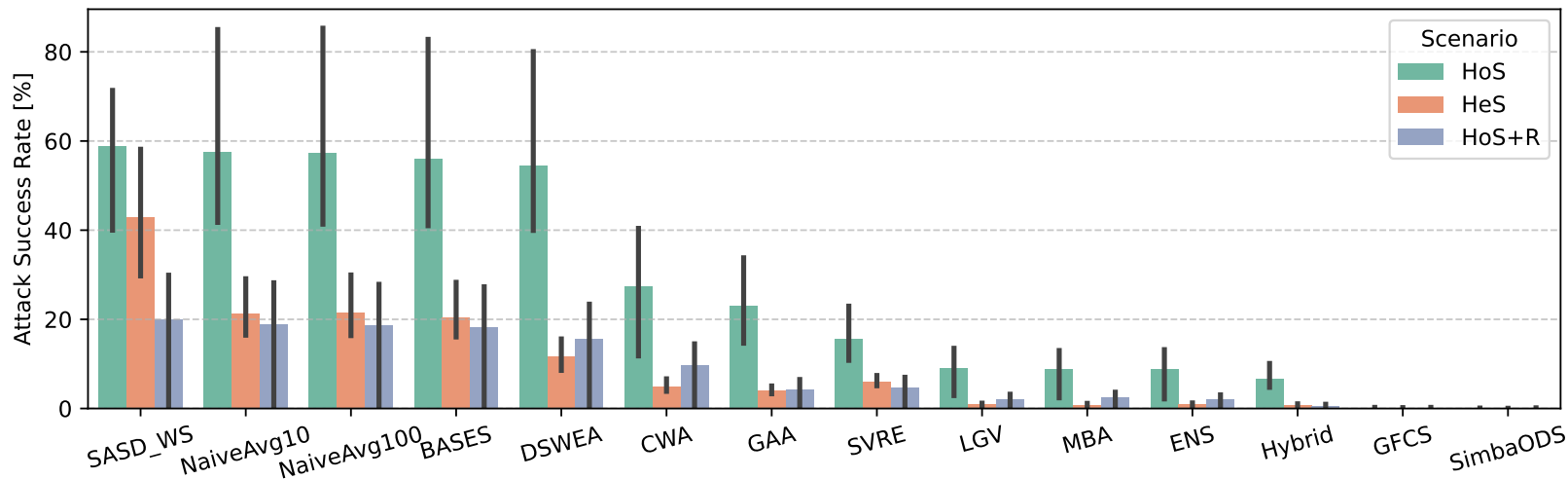
Fabio Brau, Maura Pintor, Antonio Emanuele Cinà, Raffaele Mura, Luca Scionis, Luca Oneto, Fabio Roli, Battista Biggio

# Benchmarking Attacks on Standard Scenarios

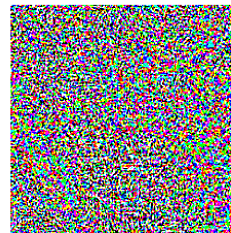# Attacking Classification Models

**Gradient-Based Perturbation**

$$x^* = x - \varepsilon \,\mathrm{sgn}\boxed{\nabla_x \mathcal{L}(g_\theta(x), y)}$$
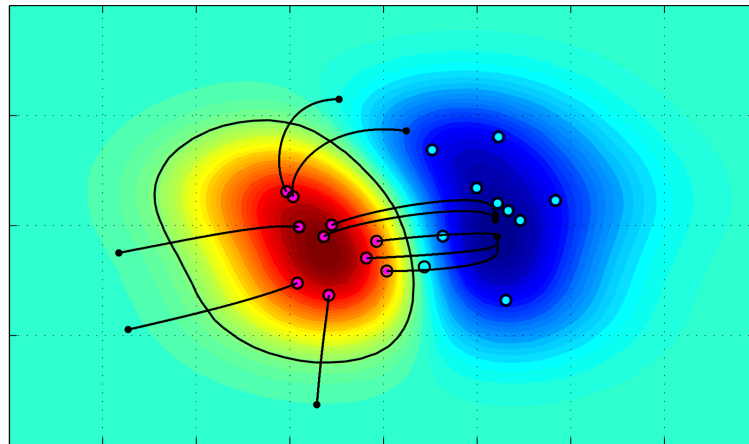
Gradient based on the target



Panda    +    =    Gibbon

Goodfellow, Ian J., et el. "Explaining and harnessing adversarial examples.". ICLR, 2015.

# Attacking Classification Models

**Adversarial Attacks as Minimum Problem**

$$x^* \in \underset{x}{\arg\min} \boxed{\mathcal{L}(g_\theta(x), y)}$$

$$\text{s.t.} \quad d(x, x_0) \leq \varepsilon$$

Assuming Differentiable Objective



Biggio, Battista, et al. "Evasion attacks against machine learning at test time." ECML-PKDD, 2013.

# Ensemble-based Attacks Formulation

With a Black-box Target, gradient is not accessible    $\cancel{\nabla_x \mathcal{L}(g_\theta(x), y)}$

**Ensembled-Based Transfer Attack**

$$x^* \in \underset{x}{\arg\min} \, \mathcal{L}_{\text{ens}}(x, y, \mathbf{f}; g(x))$$

$$\text{s.t.} \quad \|x - x_0\|_p < \varepsilon.$$

Differentiable Surrogates models

$$x^*(w) \in \underset{x}{\arg\min} \, \mathcal{L}_{loc}(x, t, \mathbf{f}; w),$$

$$\text{s.t.} \quad \|x - x_0\|_p \leq \varepsilon,$$

Local Attacks on Surrogates

$$w^* \in \underset{w \in \mathcal{W}}{\arg\min} \, \mathcal{L}(g(x^*(w)), y),$$

Refinement by querying the target

5

# Coverage of the Benchmark and Motivation

*Which is the best Ensemble-Based Transfer Attack ?*

## Compared Methods

| Attack | Venue | m |
|--------|-------|---|
| SubSpace | NeurIPS 2019 | 3 |
| SimbaODS | NeurIPS 2020 | 4 |
| Hybrid | Usenix 2020 | 3 |
| GFCS | ICLR 2022 | 4 |
| BASES | BASES 2022 | 20 |
| GAA | PR 2024 | 4 |
| DSA | Usenix 2024 | 3 |
| DSWEA | PR 2025 | 10 |

Large pool of surrogates has been sometime used !!

## Flaws of Current Method

Biased Surrogates

Weak Targets

Query Effectiveness

# How TransferBench Addresses the Gaps



**Flaws**

Biased Surrogates

Weak Targets

**Mitigations**

Scenarios

Query Effectiveness

Baselines

Query-free Methods and Naïve Average

**Homogenous Scenario**

Surrogates       Target

**Heterogeneous Scenario**

Surrogates       Target

**Robust Scenario**

Surrogates       Target

# Transferbench Ease of Use

```
demo.py
1    [Generare codice
```

# Main Results



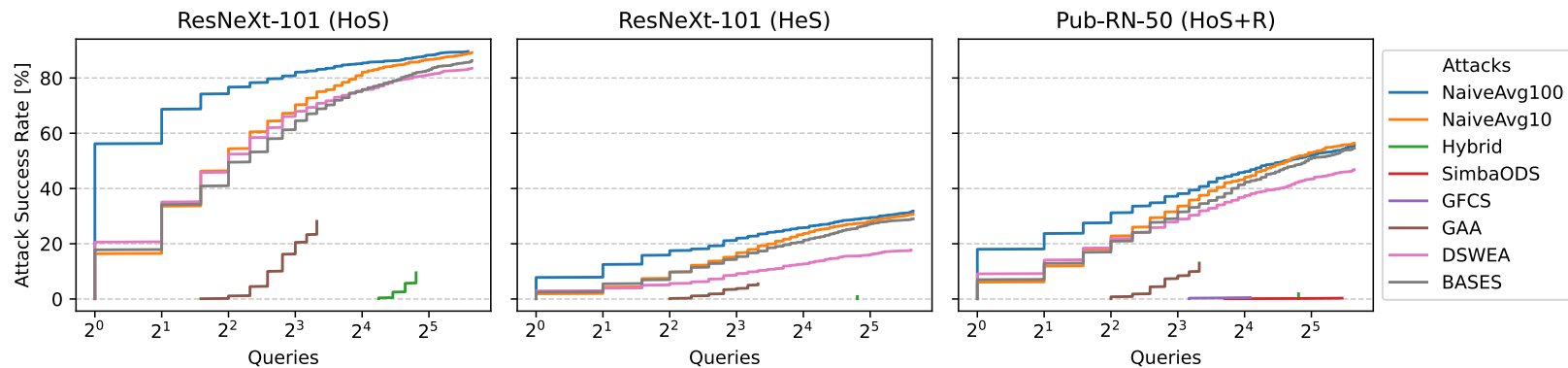Attack Success Rate for ImageNetT and $\varepsilon = \frac{16}{255}$

In all the Scenarios, current methods are worst than baselines

Attacking Robust Models is still an open problem

# Main Results



Querying the target does not really contribute to refine the attack

# TransferBench

## Benchmarking Ensemble-based Black-box Transfer Attacks

**Paper**   **Code**