

CS 499/579: TRUSTWORTHY ML
ADVERSARIAL ATTACKS: USE QUERIES

Tu/Th 4:00 – 5:50 pm

Instructor: **Sanghyun Hong**

sanghyun.hong@oregonstate.edu



Oregon State
University

SAIL
Secure AI Systems Lab

ADVERSARIAL EXAMPLES ATTACKS

- Test-time (evasion) attack
 - Given a test-time sample x
 - Craft an adversarial example x^* that fools the target neural network

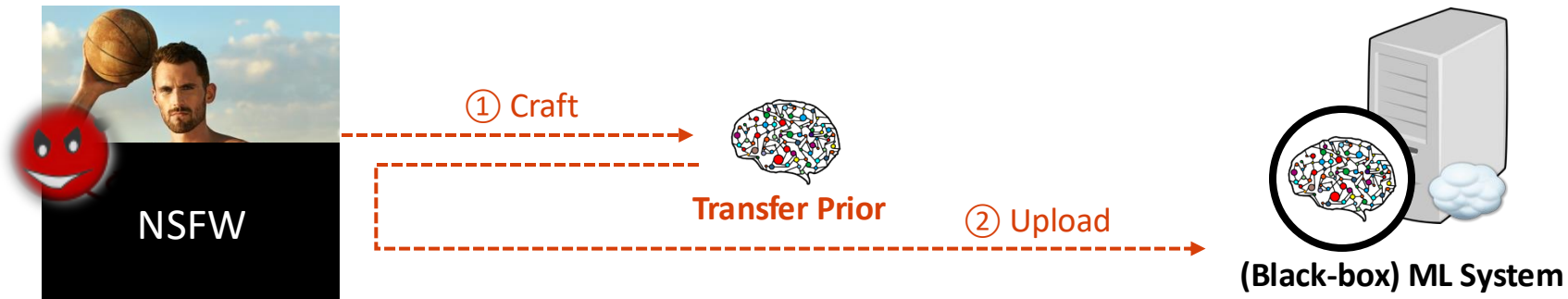
ADVERSARIAL ATTACKS

- Example: An adversary wants to upload NSFW image to the cloud



(TRANSFER-BASED) BLACK-BOX ADVERSARIAL ATTACK

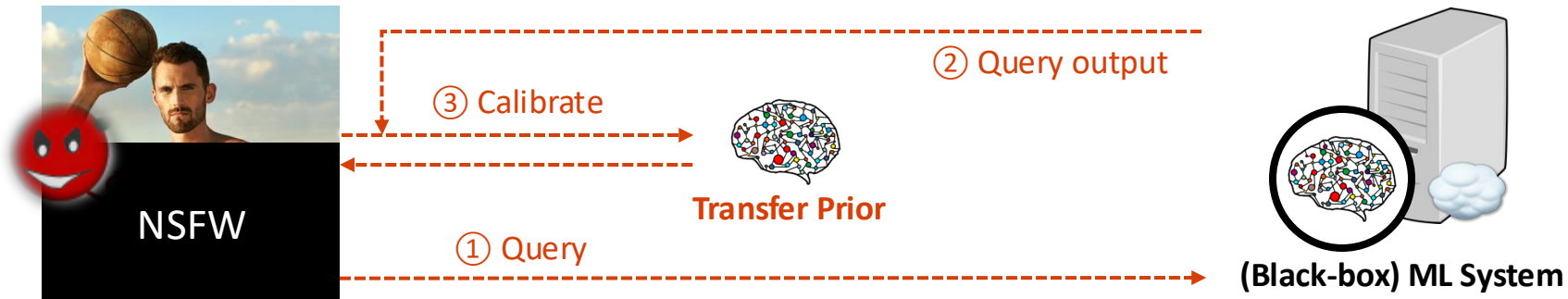
- Example: An adversary wants to upload NSFW image to the cloud



– Transfer-based attacks¹² : craft adv. examples on a transfer prior

(OPTIMIZATION-BASED) BLACK-BOX ADVERSARIAL ATTACK

- Example: An adversary wants to upload NSFW image to the cloud



- **Transfer-based attacks**¹² : craft adv. examples on a transfer prior
- **Optimization-based attacks**³ : craft them iteratively with query outputs and a transfer prior

[1] Goodfellow et al., *Explaining and Harnessing Adversarial Examples*, ICLR 2015

[2] Madry et al., *Towards Deep Learning Models Resistant to Adversarial Attacks*, ICLR 2018

[3] Cheng et al., *Improving Black-box Adversarial Attacks with a Transfer-based Prior*, NeurIPS 2019

NOW WE TALK ABOUT OPTIMIZATION-BASED ATTACKS

PRIOR CONVICTIONS: BLACK-BOX ADVERSARIAL ATTACKS WITH BANDITS AND PRIORS, ILYAS ET AL., ICLR 2019

RECAP: THE FORMULATION

- Test-time (evasion) attack
 - **Goal:**
 - Craft human-imperceptible perturbations that can make a test-time sample misclassified by a model
 - **(Black-box) Knowledge:**
 - Do not know the model architecture and/or
 - Do not know the trained model's parameters and/or
 - Do not know the training data
 - **Capability:**
 - Sufficient computational power to craft adversarial examples

How Can An Adversary Launch Attacks on (Black-box) Models?

OPTIMIZATION-BASED ATTACK

- How can an adversary launch black-box attacks?
 - Brute-force attacks
 - Query-based attacks
 - Transfer attacks

OPTIMIZATION-BASED ATTACK

- Research questions
 - How can we make the optimization-based attacks more successful?
 - How effective (and successful) is this new method?

REVISIT: THE FORMULATION

- Suppose:

- (x, y) : a test-time sample; $x \in R^d$ and $y \in [k]$; $x \in [0, 1]$
- f : a neural network; θ : its parameters
- $L(\theta, x, y)$: a loss function

- Goal (of the first order attacker):

- Find an $x^{adv} = x + \delta$ such that $\max_{\delta \in \mathcal{S}} L(\theta, x^{adv}, y)$ while $\|\delta\|_p \leq \varepsilon$

- PGD Crafts:

$$x^{t+1} = \Pi_{x+\mathcal{S}} \left(x^t + \alpha \operatorname{sgn}(\nabla_x L(\theta, x, y)) \right).$$

We Need to Know This!

OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

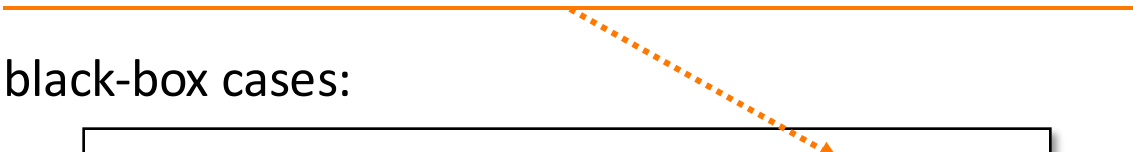
- Zeroth-order Optimization

- Finite Difference Method (FDM):

$$D_v f(x) = \langle \nabla_x f(x), v \rangle \approx (f(x + \delta v) - f(x)) / \delta.$$

- Compute: derivative of a function f at a point x towards a vector v

- FDM for the gradient with d -components:

$$\widehat{\nabla}_x L(x, y) = \sum_{k=1}^d e_k (L(x + \delta e_k, y) - L(x, y)) / \delta \approx \sum_{k=1}^d e_k \langle \nabla_x L(x, y), e_k \rangle$$


- PGD in the black-box cases:

$$x^{t+1} = \Pi_{x+\mathcal{S}} (x^t + \alpha \operatorname{sgn}(\overline{\nabla_x L(\theta, x, y)})) .$$

OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

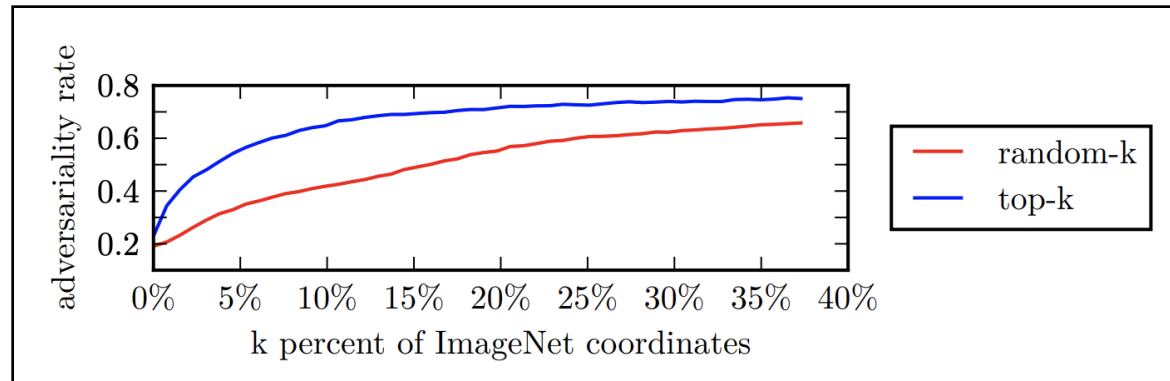
- Toy experiment

- Setup

- Compare the fraction of correctly estimated coordinates of gradients required
 - Compare top-k perturbations picked by magnitude or randomly
 - Measure the transfer-attack success rate

- Results:

- Adversarial attacks are effective even with the imperfect gradient estimate
 - Perturbations picked by magnitude is much effective than the random perturbations



OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

- Prior approaches to do this estimation

- The Least Squares Method: $\min_{\hat{g}} \|\hat{g}\|_2$ s.t. $A\hat{g} = y$.

- Iteratively compute the estimate \hat{g} , where:

- A : Queries $\{1, 2, \dots\}$
 - y : the corresponding inner product values

- Natural Evolution Strategy [Ilyas *et al.*] and Least Squares equivalence

$$\langle \hat{x}_{LSQ}, \mathbf{g} \rangle - \langle \hat{x}_{NES}, \mathbf{g} \rangle \leq O \left(\sqrt{\frac{k}{d} \cdot \log^3 \left(\frac{k}{p} \right)} \right) \|\mathbf{g}\|^2$$

OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

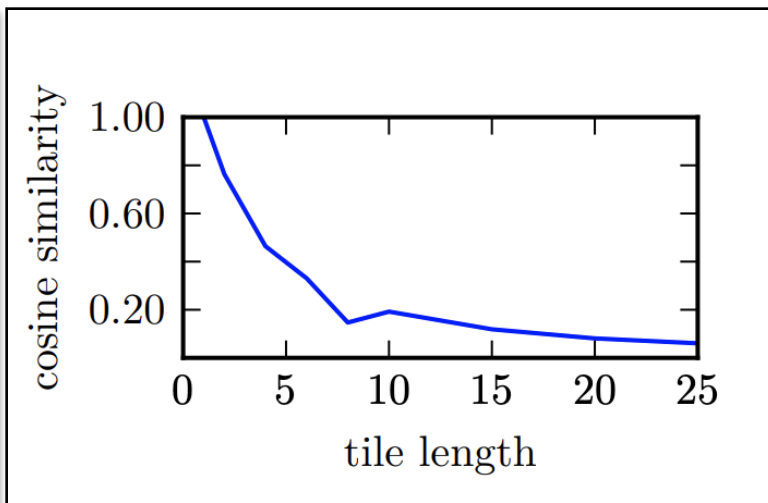
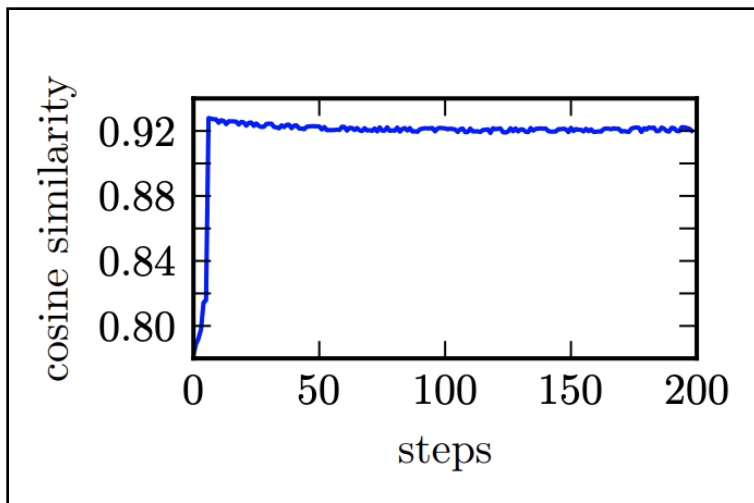
- **Prior** (= knowledge an adversary can acquire)
 - Gradients are correlated in successive attack iterations
 - Pixels close to each other tend to have similar values

OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

- Prior (= knowledge an adversary can acquire)
 - [Time-dependent] Gradients are correlated in successive attack iterations
 - [Data-dependent] Pixels close to each other tend to have similar values

OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

- Time-dependent & Data-dependent Priors



PUTTING ALL TOGETHER

- Formulate the Problem to the Bandit Framework
 - **Bandit problem**

Algorithm 1 Gradient Estimation with Bandit Optimization

```
1: procedure BANDIT-OPT-LOSS-GRAD-EST( $x, y_{init}$ )
2:    $v_0 \leftarrow \mathcal{A}(\phi)$ 
3:   for each round  $t = 1, \dots, T$  do
4:     // Our loss in round  $t$  is  $\ell_t(g_t) = -\langle \nabla_x L(x, y_{init}), g_t \rangle$ 
5:      $g_t \leftarrow v_{t-1}$ 
6:      $\Delta_t \leftarrow \text{GRAD-EST}(x, y_{init}, v_{t-1})$  // Estimated Gradient of  $\ell_t$ 
7:      $v_t \leftarrow \mathcal{A}(v_{t-1}, \Delta_t)$ 
8:    $g \leftarrow v_T$ 
9:   return  $\Pi_{\partial\mathcal{K}} [g]$ 
```

PUTTING ALL TOGETHER

- Formulate the Problem to the Bandit Framework
 - Gradient Estimation

Algorithm 2 Single-query spherical estimate of $\nabla_v \langle \nabla L(x, y), v \rangle$

```
1: procedure GRAD-EST( $x, y, v$ )
2:    $u \leftarrow \mathcal{N}(0, \frac{1}{\delta} I)$  // Query vector
3:    $\{q_1, q_2\} \leftarrow \{v + \delta u, v - \delta u\}$  // Antithetic samples
4:    $\ell_t(q_1) = -\langle \nabla L(x, y), q_1 \rangle \approx \frac{L(x, y) - L(x + \epsilon q_1, y)}{\epsilon}$  // Gradient estimation loss at  $q_1$ 
5:    $\ell_t(q_2) = -\langle \nabla L(x, y), q_2 \rangle \approx \frac{L(x, y) - L(x + \epsilon q_2, y)}{\epsilon}$  // Gradient estimation loss at  $q_2$ 
6:    $\Delta \leftarrow \frac{\ell_t(q_1) - \ell_t(q_2)}{\delta} u = \frac{L(x + \epsilon q_2, y) - L(x + \epsilon q_1, y)}{\delta \epsilon} u$ 
7:   // Note that due to cancellations we can actually evaluate  $\Delta$  with only two queries to  $L$ 
8:   return  $\Delta$ 
```

PUTTING ALL TOGETHER

- Formulate the Problem to the Bandit Framework
 - Gradient Estimation

Algorithm 3 Adversarial Example Generation with Bandit Optimization for ℓ_2 norm perturbations

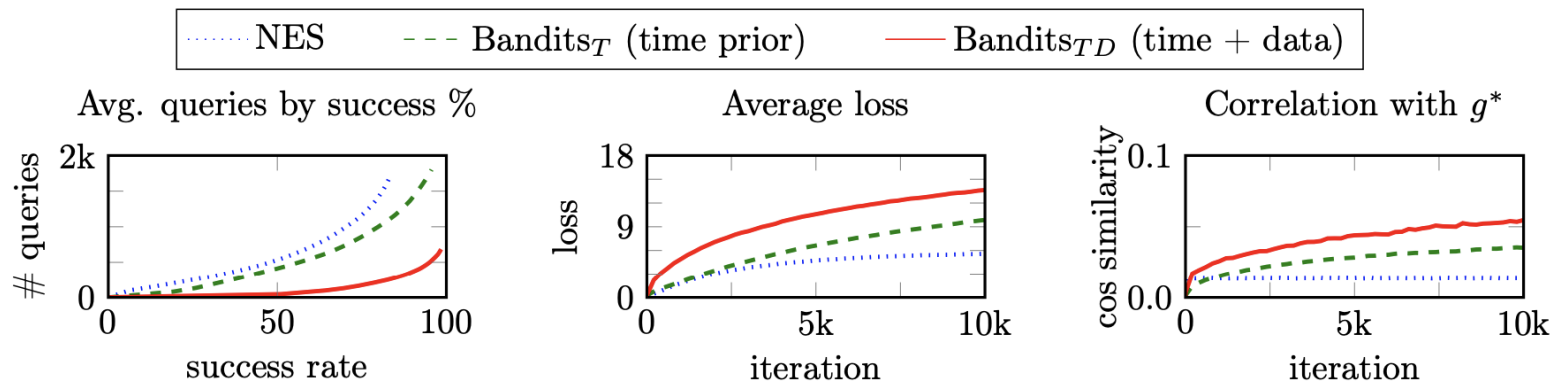
```
1: procedure ADVERSARIAL-BANDIT-L2( $x_{init}, y_{init}$ )
2:   //  $C(\cdot)$  returns top class
3:    $v_0 \leftarrow \mathbf{0}_{1 \times d}$  // If data prior,  $d < \dim(x)$ ;  $v_t$  ( $\Delta_t$ ) up (down)-sampled before (after) line 8
4:    $x_0 \leftarrow x_{init}$  // Adversarial image to be constructed
5:   while  $C(x) = y_{init}$  do
6:      $g_t \leftarrow v_{t-1}$ 
7:      $x_t \leftarrow x_{t-1} + h \cdot \frac{g_t}{\|g_t\|_2}$  // Boundary projection  $\frac{g}{\|g\|}$  standard PGD: c.f. [Rig15]
8:      $\Delta_t \leftarrow \text{GRAD-EST}(x_{t-1}, y_{init}, v_{t-1})$  // Estimated Gradient of  $\ell_t$ 
9:      $v_t \leftarrow v_{t-1} + \eta \cdot \Delta_t$ 
10:     $t \leftarrow t + 1$ 
return  $x_{t-1}$ 
```

HOW EFFECTIVE IS THIS NEW ATTACK (= METHOD)?

- Setup

- Dataset: ImageNet (10k randomly chosen samples)
- Model: Inception-v3
- Baseline: NES

- Results



OPTIMIZATION-BASED ATTACK

- Take aways
 - How **accurate** should we estimate a gradient for successful attacks?
 - PGD can be quite successful with imperfect gradient estimates
 - Query-efficiency is bounded by the prior work [Ilyas *et al.*] in practical scenarios
 - How can we estimate gradient accurately with **smaller queries**?
 - Use two priors: time- and data-dependent priors
 - Formulate the estimation into the bandit framework
 - How **effective (and successful)** is this new method?
 - Require 2.5 – 5x less queries for successful attacks compared to NES

Thank You!

Instructor: Sanghyun Hong

<https://secure-ai.systems/courses/MLSec/Sp23>



Oregon State
University

SAIL
Secure AI Systems Lab