# Notes

- Call for actions
  - In-class presentation sign-ups
  - Term project team-up (by the 10$^{th}$)

# CS 499/579: Trustworthy ML
# Adversarial attacks: transferability

Tu/Th 4:00 – 5:50 pm

Instructor: Sanghyun Hong

sanghyun.hong@oregonstate.edu

Oregon State University

SAIL
Secure AI Systems Lab

# WHY DO ADVERSARIAL ATTACKS TRANSFER?

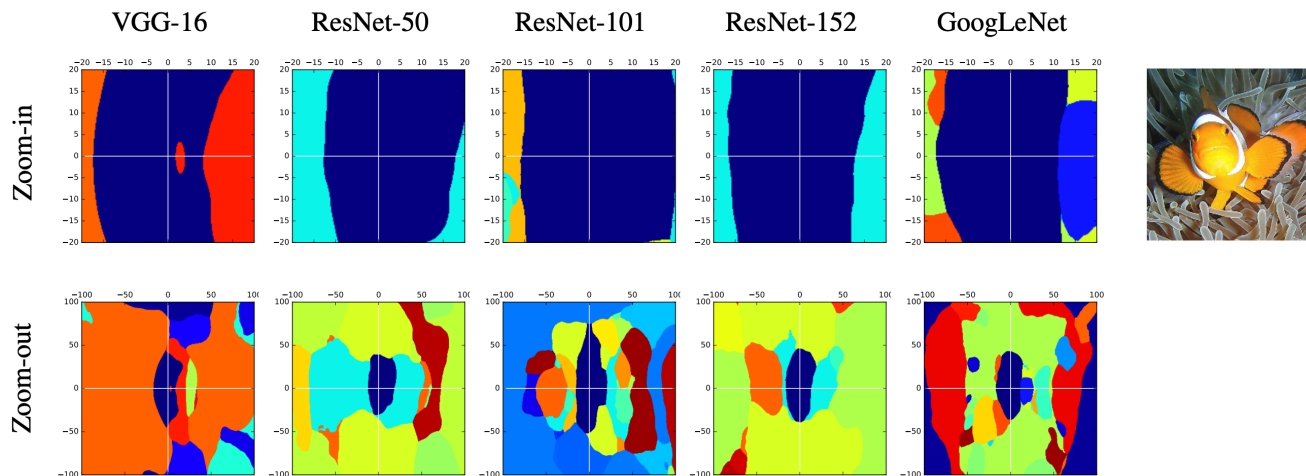THE SPACE OF TRANSFERABLE ADVERSARIAL EXAMPLES, TRAMER ET AL.
WHY DO ADVERSARIAL ATTACKS TRANSFER, DEMONTIS ET AL., USENIX SECURITY 2019

# WHY DO ADVERSARIAL ATTACKS TRANSFER?

- How to answer this question?
  - Inspect a model's decision boundary (Liu et al., Tramer et al.)
  - Inspect the data distribution (Tramer et al.)
  - Comprehensive empirical evaluation (Demotis et al.)
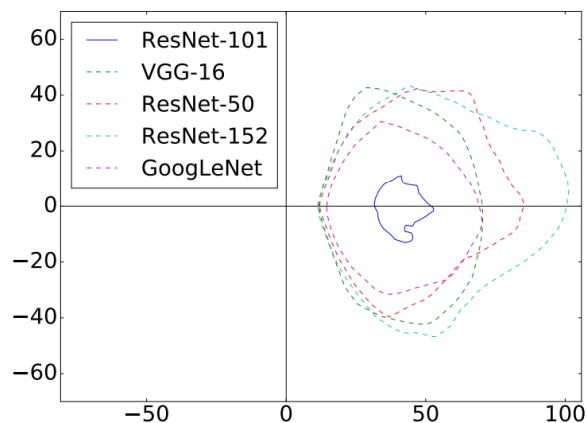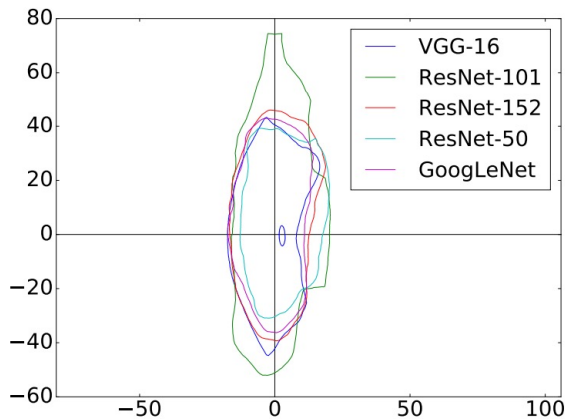  - …

Oregon State
University

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Recap: Inspect a model's decision boundary
  - Setup:
    - Take a sample image, and two orthogonal gradient directions
    - Perturb the sample along each direction and measure the labels
  - Results

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Recap: Inspect a model's decision boundary: ensemble
  - Setup:
    - Take a sample image, and two orthogonal gradient directions
    - Perturb the sample along each direction and measure the labels
  - Results

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Inspect a model's decision boundary: subspace
  - Setup:
    - Take a sample image, and *multiple* orthogonal gradient directions
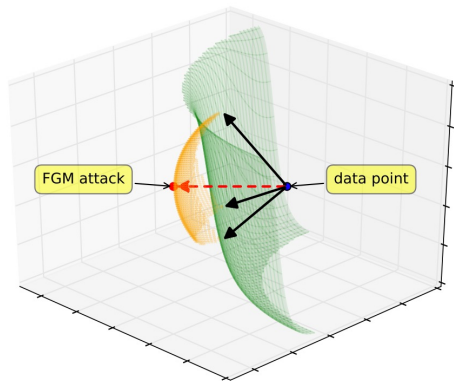    - Perturb the sample along each direction and measure the loss
  - Results



Figure 1: Illustration of the Gradient Aligned Adversarial Subspace (GAAS). The gradient aligned attack (red arrow) crosses the decision boundary. The black arrows are orthogonal vectors aligned with the gradient that span a subspace of potential adversarial inputs (orange).
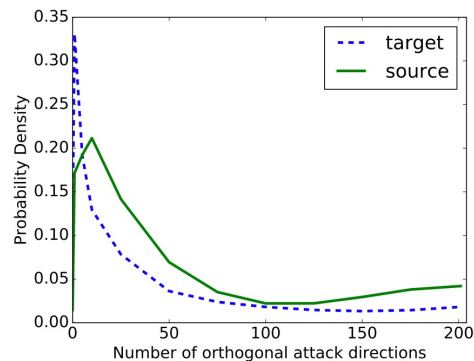
Figure 2: Probability density function of the number of successful orthogonal adversarial perturbations found by the GAAS method on the source DNN model, and of the number of perturbations that transfer to the target DNN model.

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Inspect a model's decision boundary: similarity
  - Setup:
    - Take a sample image, and *three* gradient directions: Legit, Adv., and Rand.
    - Perturb the sample along each direction and measure the distance to the decision boundary and between two boundaries
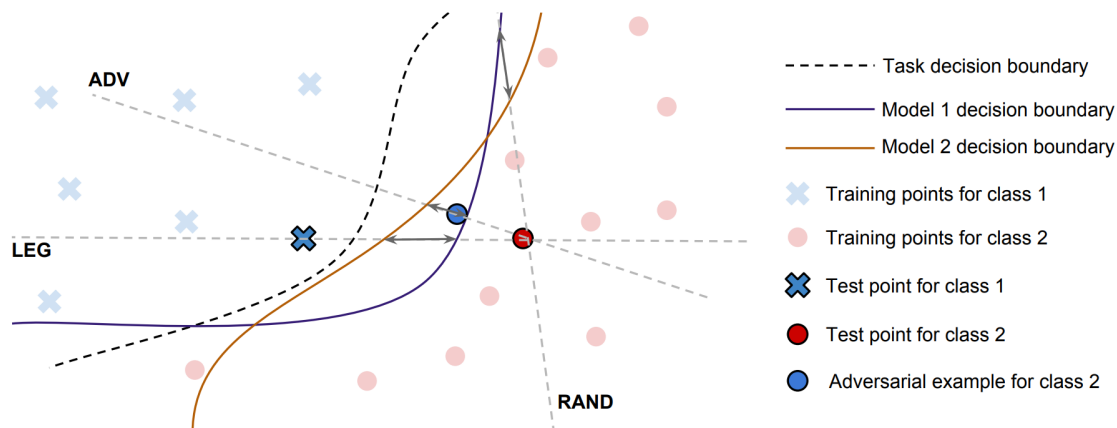  - Results



Figure 3: The three directions (Legitimate, Adversarial and Random) used throughout Section 4 to measure the distance between the decision boundaries of two models. The gray double-ended arrows illustrate the *inter-boundary* distance between the two models in each direction.

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Inspect a model's decision boundary: similarity
  - Setup:
    - Take a sample image, and *three* gradient directions: Legit, Adv., and Rand.
    - Perturb the sample along each direction and measure the distance to the decision boundary and between two boundaries
  - Results



Figure 4: Minimum distances and inter-boundary distances in three directions for MNIST models.

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Ubiquity hypothesis
  - Hypothesis I:
    - Suppose two models achieve low errors and low robustness to adv examples adversarial examples crafted on one model transfer to the other
  - Evaluation I:
    - Train two different models on a task and find adversarial examples do not transfer
    - Results: found, reject

  - Hypothesis II (XOR artifact):
    - Suppose that two models trained on the same set of input features learn representations for which adversarial examples do not transfer to each other; both are non-robust
  - Evaluation II:
    - Adversarial examples crafted one model does not transfer well to the other
    - Results: does not work, reject

# Why do adversarial attacks transfer?

- How to answer this question?
  - Inspect a model's decision boundary (Liu et al., Tramer et al.)
  - Inspect the data distribution (Tramer et al.)
  - Comprehensive empirical evaluation (Demotis et al.)
  - …

Oregon State
University

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Comprehensive empirical evaluation
  - Setup:
    - A strong adversarial attack
    - Models
      - SVM (linear / rbf)
      - (logistic / ridge) Regression
      - Neural networks
    - Datasets
      - MNIST-89
      - Drebin (android malware)



● initial / source example
● minimum-distance *black-box* adversarial example
▲ minimum-distance *white-box* adversarial example
● maximum-confidence *black-box* adversarial example
▲ maximum-confidence *white-box* adversarial example

surrogate classifier $\hat{f}(x)$ used to craft *black-box* adversarial examples
target classifier $f(x)$ used to craft *white-box* adversarial examples
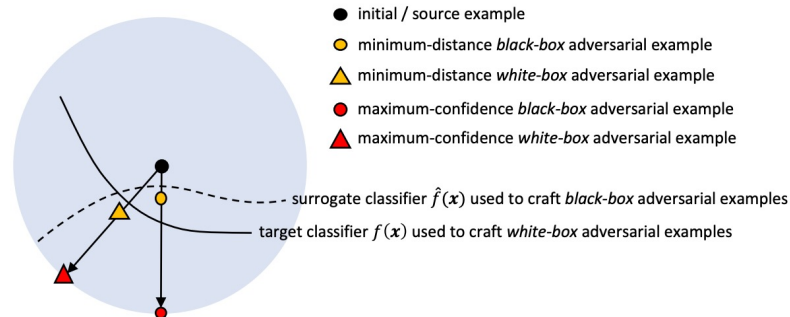
Figure 2: Conceptual representation of maximum-confidence evasion attacks (within an $\ell_2$ ball of radius $\varepsilon$) vs. minimum-distance adversarial examples. Maximum-confidence attacks tend to transfer better as they are misclassified with higher confidence (though requiring more modifications).

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Comprehensive empirical evaluation
  - Setup:
    - Model complexity (= # of parameters) matters
    - Train two models with different complexities and measure the success rate of white-box attacks (why?)
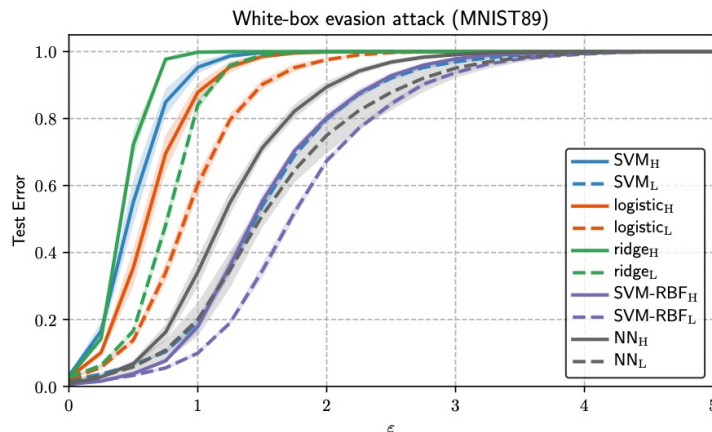  - Results



Figure 5: White-box evasion attacks on MNIST89. Test error against increasing maximum perturbation ε.

# WHY DO ADVERSARIAL EXAMPLES TRANSFER?

- Comprehensive empirical evaluation
  - Setup:
    - Model complexity (= # of parameters) matters
    - Train two models with different complexities and measure the success rate of white-box
    - Run transfer-based attacks between all pairs of models and measure the attack success
  - Results
    - Use of low-complexity models as a surrogate increases the adversarial transferability
    - Random forest classifiers are particularly vulnerable to transfer-based attacks

| | $SVM_H$ | $SVM_L$ | $logistic_H$ | $logistic_L$ | $ridge_H$ | $ridge_L$ | $SVM\text{-}RBF_H$ | $SVM\text{-}RBF_L$ | $NN_H$ | $NN_L$ | $RF_H$ | $RF_L$ | transfer rate |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| white box | .96 | .19 | .89 | .60 | 1.00 | .83 | .17 | .10 | .31 | .21 | | | |
| $SVM_H$ | .09 | .05 | .08 | .07 | .07 | .06 | .02 | .02 | .03 | .05 | .43 | .45 | .12 |
| $SVM_L$ | .28 | .14 | .26 | .22 | .19 | .17 | .07 | .07 | .13 | .14 | .53 | .54 | .23 |
| $logistic_H$ | .12 | .06 | .11 | .09 | .10 | .09 | .03 | .03 | .04 | .06 | .47 | .49 | .14 |
| $logistic_L$ | .19 | .09 | .18 | .15 | .15 | .13 | .04 | .04 | .08 | .08 | .50 | .52 | .18 |
| $ridge_H$ | .08 | .04 | .07 | .05 | .11 | .07 | .02 | .02 | .03 | .04 | .43 | .45 | .12 |
| $ridge_L$ | .15 | .07 | .13 | .10 | .21 | .15 | .03 | .03 | .05 | .06 | .47 | .49 | .16 |
| $SVM\text{-}RBF_H$ | .19 | .10 | .17 | .15 | .13 | .12 | .06 | .06 | .10 | .11 | .53 | .53 | .19 |
| $SVM\text{-}RBF_L$ | .25 | .13 | .23 | .20 | .17 | .16 | .08 | .08 | .14 | .14 | .53 | .54 | .22 |
| $NN_H$ | .20 | .10 | .18 | .15 | .14 | .12 | .05 | .05 | .11 | .10 | .52 | .53 | .19 |
| $NN_L$ | .24 | .12 | .22 | .20 | .16 | .15 | .07 | .07 | .13 | .13 | .53 | .53 | .21 |

(a) $\varepsilon = 1$

Oregon State University

- Comprehensive empirical evaluation
  - Setup:
    - Gradient alignment (= # of parameters) matters
    - Compute the gradient from a surrogate and a target for the same $x$ and measure the cosine similarity metric between the two gradients
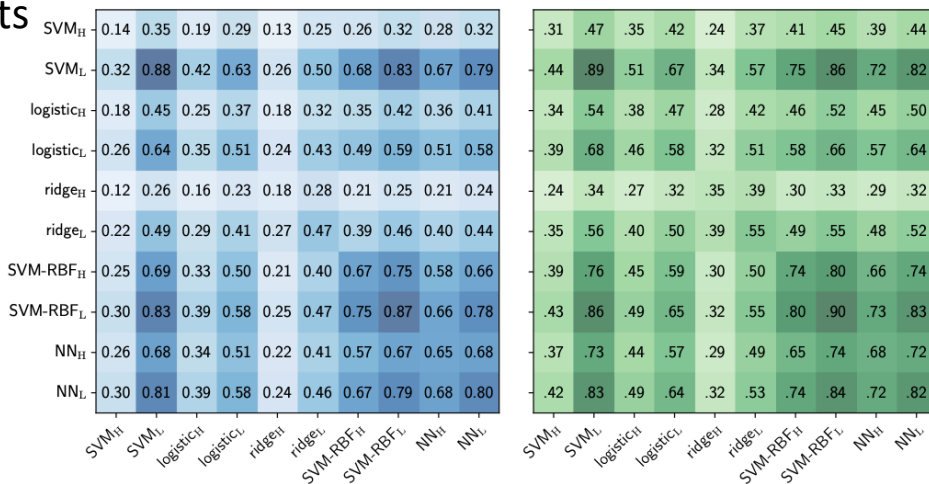  - Results



Figure 8: Gradient alignment and perturbation correlation for evasion attacks on MNIST89. *Left:* Gradient alignment $R$ (Eq. 18) between surrogate (rows) and target (columns) classifiers, averaged on the unmodified test samples. *Right:* Pearson correlation coefficient $\rho(\delta, \hat{\delta})$ between white-box and black-box perturbations for $\varepsilon = 5$.

# WHY DO ADVERSARIAL ATTACKS TRANSFER?

- Take aways
  - If the decision boundaries of two models similar, the transferability increases
  - If the transferability is high between two models, there's a common adv. subspace
  - The transferability is non-trivial
    - Two models trained to achieve low-loss and low-resilience to white-box attacks
    - But the adversarial examples do not transfer well between each other
  - XOR artifacts
    - Two models trained with the same set of features, but on disjoint datasets
    - But the adversarial examples do not transfer well between each other
  - If the attacker uses low-complexity models, the transferability becomes high
  - If the two models have aligned gradients, the transferability is high
  - ... (your contributions)

Oregon State
University

# 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**

**S**AIL
**S**ecure AI Systems Lab

# ADVERSARIAL E̶X̶A̶M̶P̶L̶E̶S̶ ATTACKS

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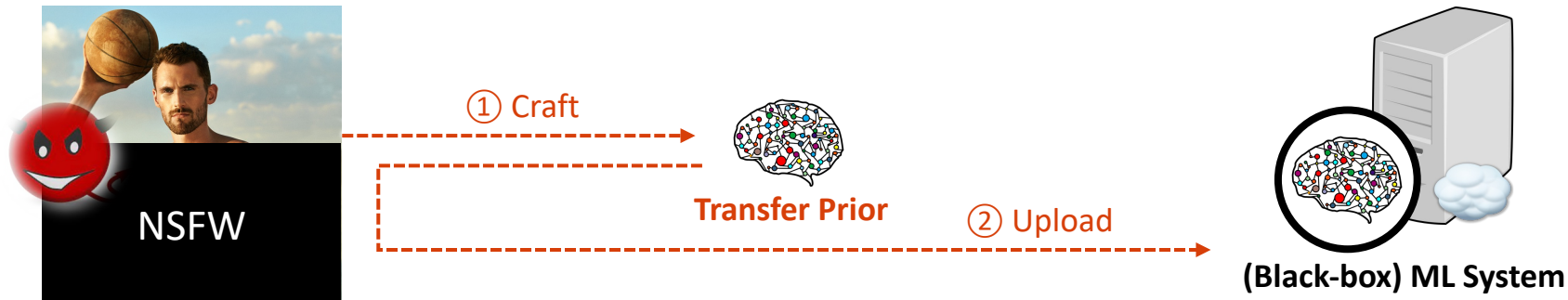  - Craft an adversarial example $x^*$ that fools the target neural network

# ADVERSARIAL ATTACKS

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



① Upload

ML System

# (Transfer-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

[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

Oregon State
University

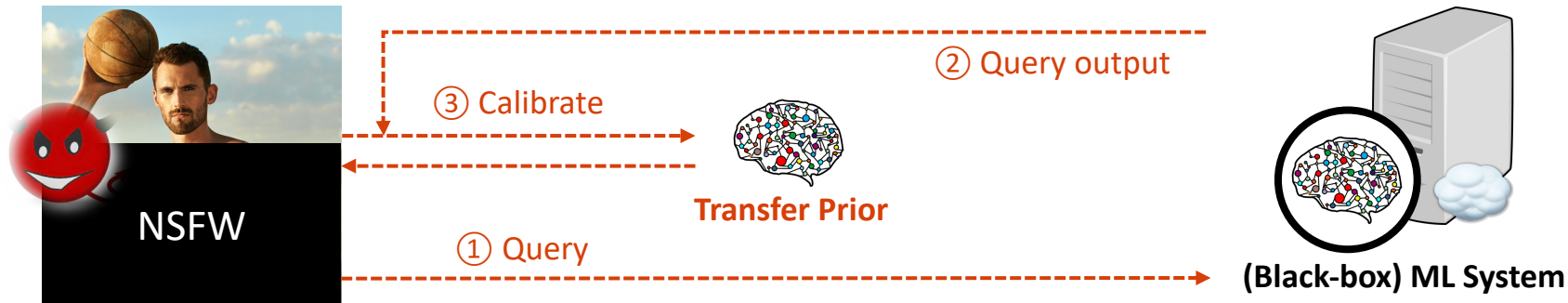# (Optimization-based) black-box adversarial attack

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



- **Transfer-based attacks**[1,2] : 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

Oregon State
University

# 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?**

Oregon State
University

# Optimization-based attack

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

Oregon State
University

# Optimization-based attack

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

Oregon State
University

# 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; $\theta$: its parameters
  - $L(\theta, x, y)$: a loss function

- Goal (of the first order attacker):
  - Find an $x^{adv} = x + \delta$ such that $\max_{\delta \in S} L(\theta, x^{adv}, y)$ while $||\delta||_p \leq \varepsilon$

- PGD Crafts:

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

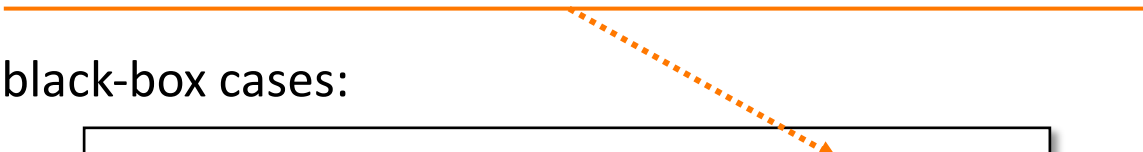**We Need to Know This!**

Oregon State
University

# 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:

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

- PGD in the black-box cases:

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

**Oregon State**
University

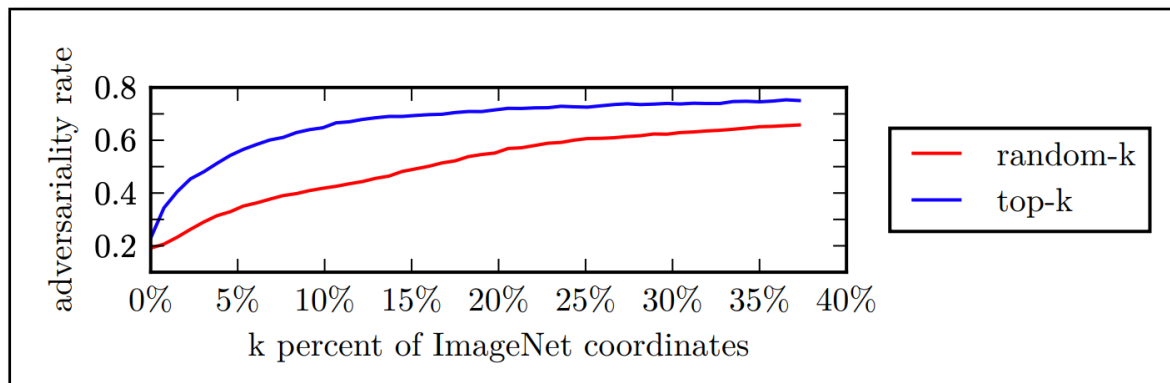# OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

- 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: $\min_{\widehat{g}} \|\widehat{g}\|_2 \quad \text{s.t.} \ A\widehat{g} = y.$

  - Iteratively compute the estimate $\widehat{g}$, where:
    - $A$: Queries $\{1, 2, \ldots\}$
    - $y$: the corresponding inner product values

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

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

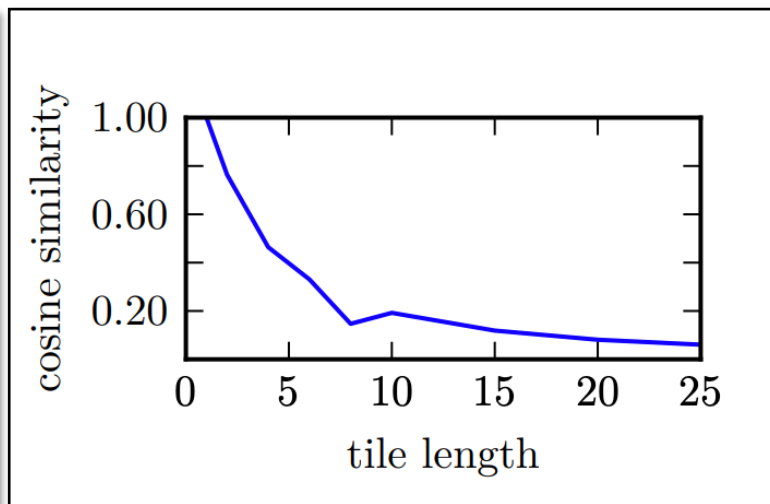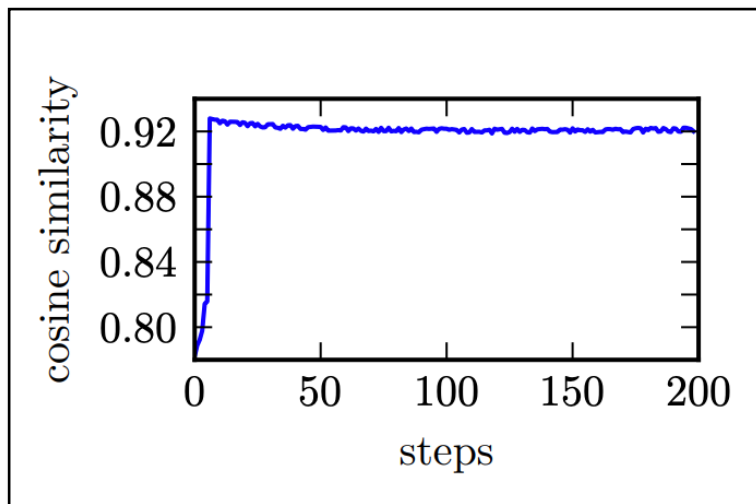# OPTIMIZATION-BASED ATTACK IS THE GRADIENT ESTIMATION PROBLEM

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

Oregon State University

# PUTTING ALL TOGETHER

- Formulate the Problem to the Bandit Framework
  - Gradient Estimation

**Algorithm 3** Adversarial Example Generation with Bandit Optimization for $\ell_2$ norm perturbations
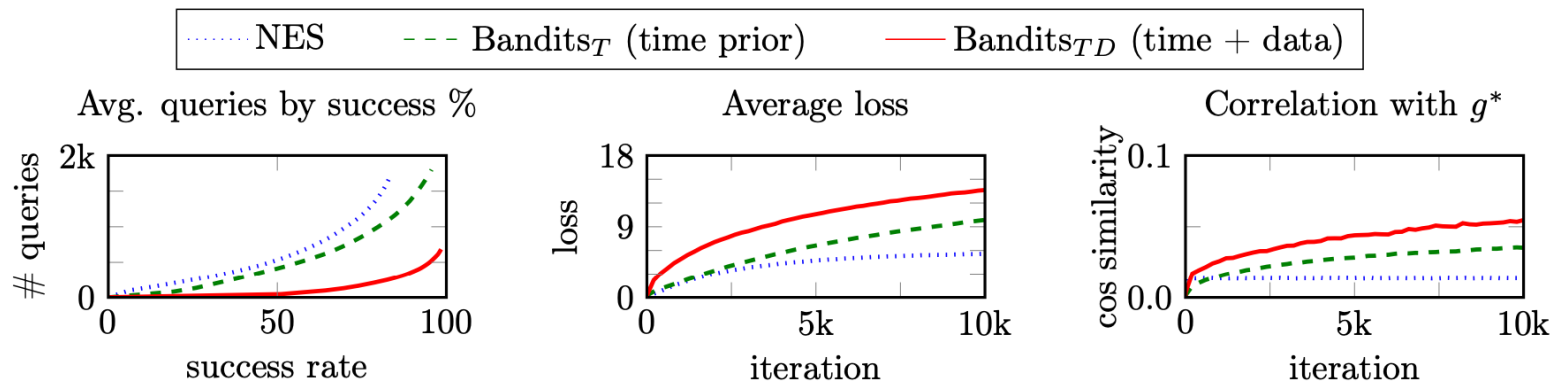
1: **procedure** ADVERSARIAL-BANDIT-L2($x_{init}, y_{init}$)
2:      // $C(\cdot)$ returns top class
3:      $v_0 \leftarrow \mathbf{0}_{1 \times d}$ // If data prior, $d < \dim(x)$; $v_t$ ($\Delta_t$) up (down)-sampled before (after) line 8
4:      $x_0 \leftarrow x_{init}$ // Adversarial image to be 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:          $\Delta_t \leftarrow$ 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** $x_{t-1}$

Oregon State University

# HOW EFFECTIVE IS THIS NEW ATTACK (= METHOD)?

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

- Results



NES ⋯⋯    Bandits$_T$ (time prior) - - -    Bandits$_{TD}$ (time + data) ——

Avg. queries by success %      Average loss      Correlation with $g^*$

Oregon State University

# 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

Oregon State
University

# Thank You!

Tu/Th 10:00 – 11:50 am (Recorded lecture)

Instructor: **Sanghyun Hong**

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

**Oregon State University**

**S**AIL
**S**ecure AI Systems Lab