

NOTES

- Call for actions
 - [Important!] Switch to **online** from the next class (11/07)
 - [Important!] Checkpoint Presentation II (**online**, on 11/07)
 - 12-min presentation + 3 min Q&A
 - Presentation MUST cover:
 - 1 slide on your research topic
 - 1-2 slides on your goals and ideas (how do you plan to achieve your goals)
 - 1-2 slides on your *experimental design*
 - 1-2 slides on your *preliminary results* [*very important*]
 - 1 slide on your *next steps* until the final presentation

CS 499/579: TRUSTWORTHY ML

TARGETED POISONING ATTACKS

Tu/Th 4:00 – 5:50 PM

Sanghyun Hong

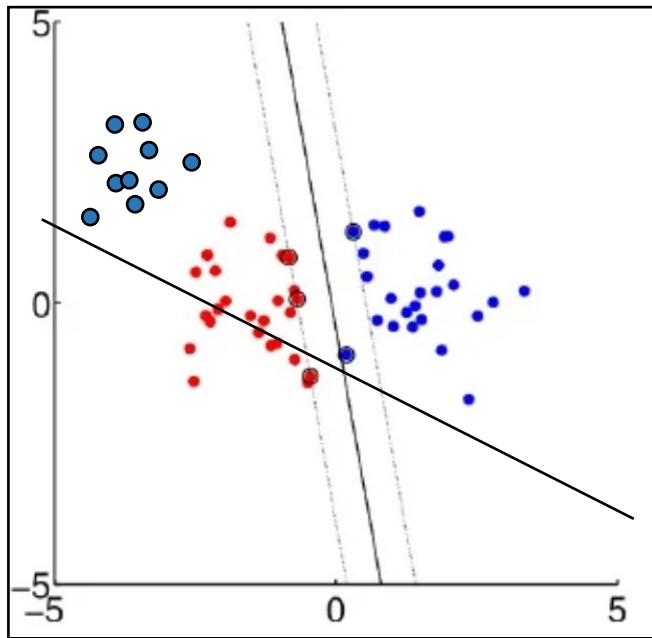
sanghyun.hong@oregonstate.edu



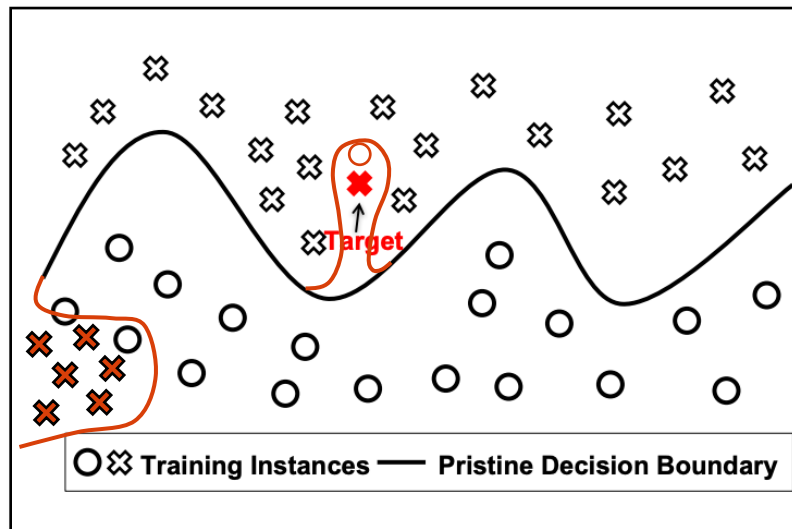
Oregon State
University

SAIL
Secure AI Systems Lab

RECAP: CONCEPTUAL ILLUSTRATION OF THE VULNERABILITY TO POISONING



← Linear model (SVM)



Neural Network →

TARGETED POISONING THREAT MODEL

- Goal
 - **Targeted** attack
 - Model causes a misclassification of (x_t, y_t) , while preserving acc. on D_{val}
- Capability
 - Know a target (x_t, y_t)
 - Pick p candidates from test data $(x_{c1}, y_{c1}), (x_{c2}, y_{c2}) \dots$ and craft poisons $(x_{p1}, y_{p1}), (x_{p2}, y_{p2}) \dots$
 - Inject them into the training data
- Knowledge
 - D_{tr} : training data
 - D_{test} : test-set data (validation data)
 - f : a model and its parameters θ
 - A : training algorithm (e.g., mini-batch SGD)

TARGETED POISONING THREAT MODEL

- Goal
 - Targeted **clean-label** ($y_{c1} = y_{p1}$) attack
 - Model causes a misclassification of (x_t, y_t) , while preserving acc. on D_{val}
- Capability
 - Know a target (x_t, y_t)
 - Pick p candidates from test data $(x_{c1}, y_{c1}), (x_{c2}, \dots)$ and craft poisons $(x_{p1}, y_{p1}), (x_{p2}, \dots)$
 - Inject them into the training data
- Knowledge
 - D_{tr} : training data
 - D_{test} : test-set data (validation data)
 - f : a model and its parameters θ
 - A : training algorithm (e.g., mini-batch SGD)

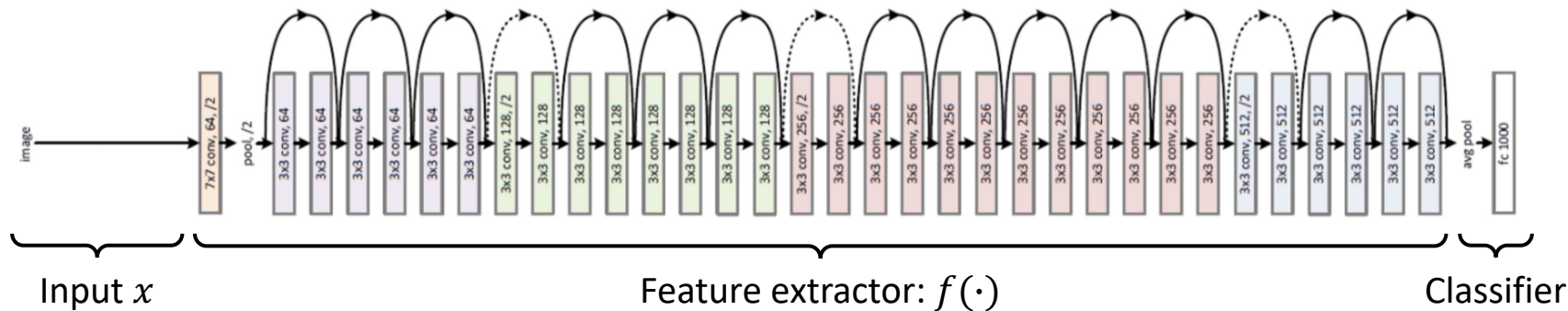
TOPICS FOR PART II – DATA POISONING

- Research questions
 - What are some examples of poisoning attacks?
 - How can we generate *indiscriminate* poisoning examples?
 - How can we synthesize poisoning samples for *targeted* attacks?
 - How can we mitigate data poisoning attacks?

HOW CAN WE PERFORM CLEAN-LABEL TARGETED ATTACKS?

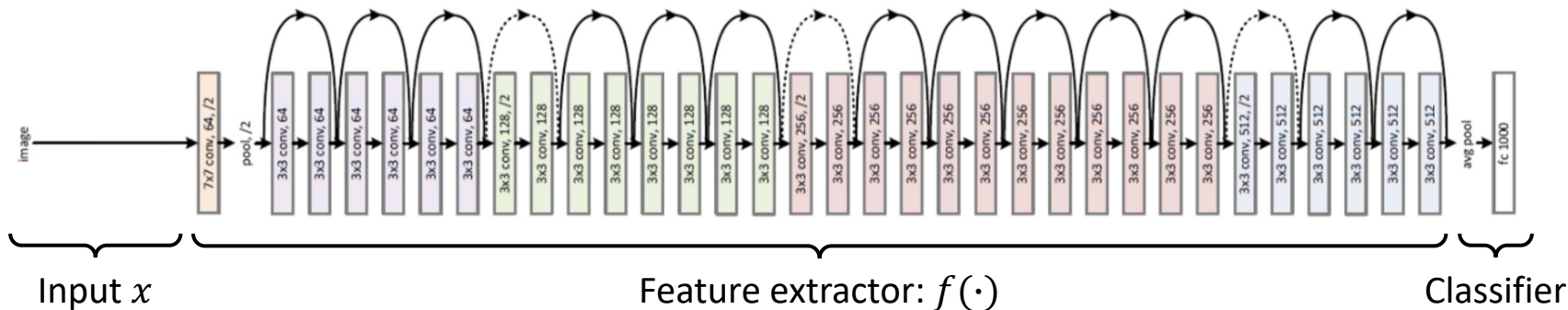
POISON FROGS! TARGETED CLEAN-LABEL POISONING ATTACKS ON NEURAL NETWORKS, SHAFABI ET AL., NEURIPS 2018

BACKGROUND: CONVOLUTIONAL NEURAL NETWORKS



- A conventional view:
 - Convolutions: extract features, embeddings, latent representations, ...
 - Last layer: uses the output for a classification task

BACKGROUND: CONVOLUTIONAL NEURAL NETWORKS

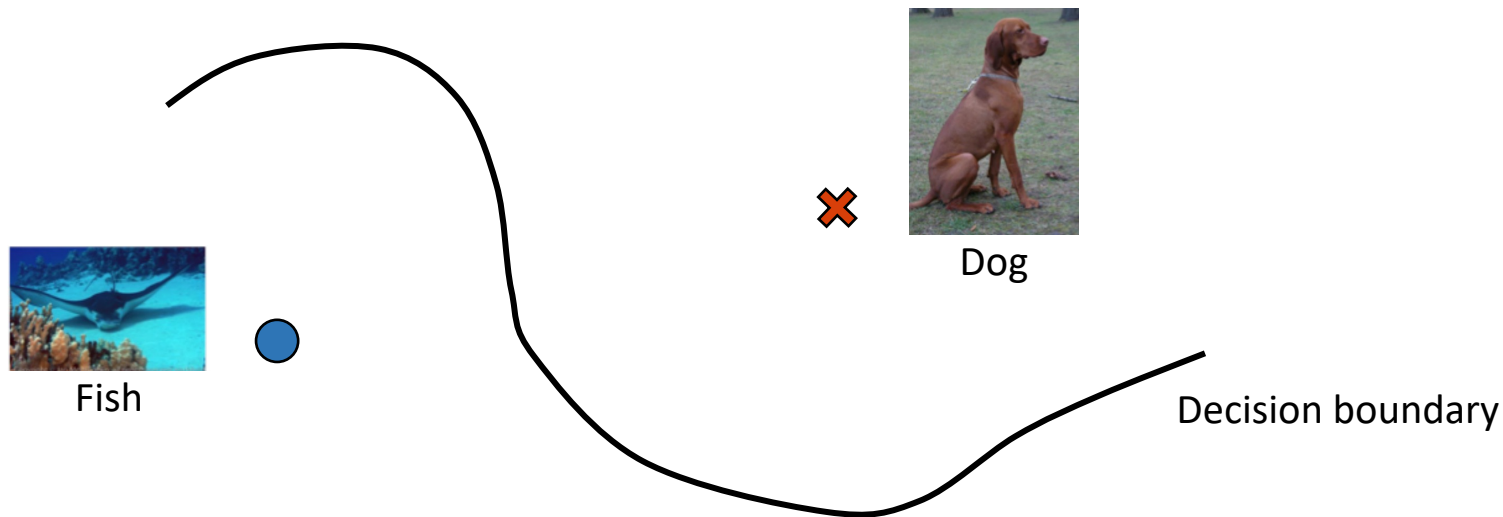


- Input-space \neq Feature-space:
 - Two samples similar in the input-space can be far from each other in the feature-space
 - Two samples very different in the input-space can be close to each other in f

THE KEY IDEA: FEATURE COLLISION

- Goal

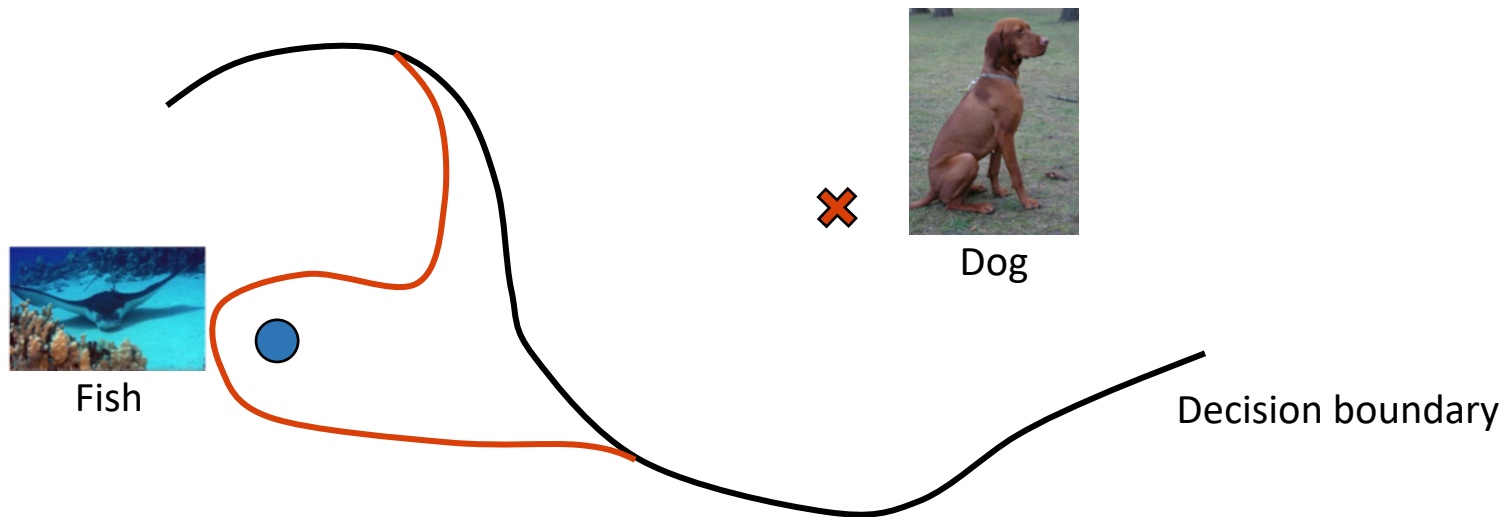
- You want your *any* poison to be closer to your target (x_t, y_t) in the *feature space*



THE KEY IDEA: FEATURE COLLISION

- Goal

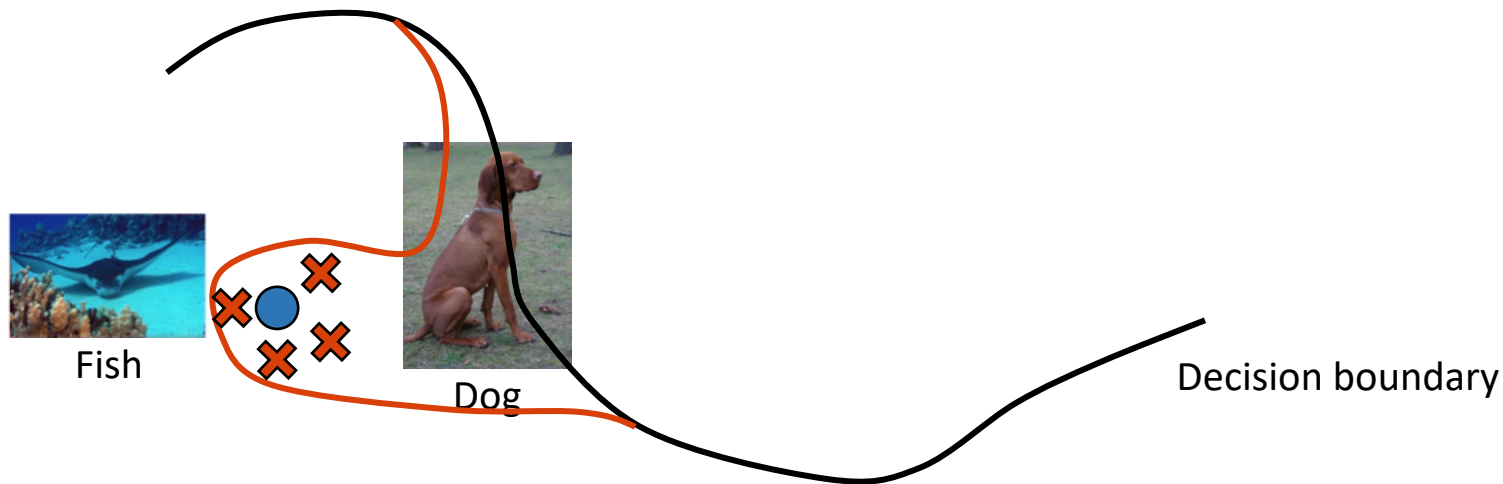
- You want your *any* poison to be closer to your target (x_t, y_t) in the *feature space*



The Fish Becomes **DogFish!**

THE KEY IDEA: FEATURE COLLISION

- Goal
 - You want your *any* poison to be closer to your target (x_t, y_t) in the *feature space*



THE KEY IDEA: FEATURE COLLISION

- Goal

- You want your *any* poison to be closer to your target (x_t, y_t) in the *feature space*
- Objective:

$$\mathbf{p} = \underset{\mathbf{x}}{\operatorname{argmin}} \quad \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2 + \beta \|\mathbf{x} - \mathbf{b}\|_2^2$$

- Optimization:

Algorithm 1 Poisoning Example Generation

Input: target instance t , base instance b , learning rate λ

Initialize \mathbf{x} : $x_0 \leftarrow b$

Define: $L_p(x) = \|f(\mathbf{x}) - f(\mathbf{t})\|_2^2$

for $i = 1$ **to** maxIters **do**

 Forward step: $\hat{x}_i = x_{i-1} - \lambda \nabla_x L_p(x_{i-1})$ // construct input perturbations

 Backward step: $x_i = (\hat{x}_i + \lambda \beta b) / (1 + \beta \lambda)$ // decide how much we will perturb

end for

EVALUATIONS

- Scenarios
 - Scenario 1: Transfer learning
 - Scenario 2: End-to-end learning

EVALUATIONS: TRANSFER LEARNING

- Setup
 - Dataset: Dog vs. Fish (ImageNet)
 - Models: Inception-V3 (Pretrained on ImageNet)
- “one-shot kill” Attacks
 - Goal: Dog > Fish or Fish > Dog | All 1099 targets from the test-set
 - Craft a poison using a single image chosen from the other class
 - Train the last layer on $D_{tr} \cup (x_p, y_p)$ and check if the target’s label is flipped
- Results
 - The attack succeeds with 100% accuracy
 - The accuracy drop caused by the attack is 0.2% on average

EVALUATIONS: TRANSFER LEARNING

- Examples

Clean Base

Target instances from Fish class

Poison instances made for fish class from dog base instances

Clean Base

Target instances from Dog class

Poisons made for dog class from fish bases

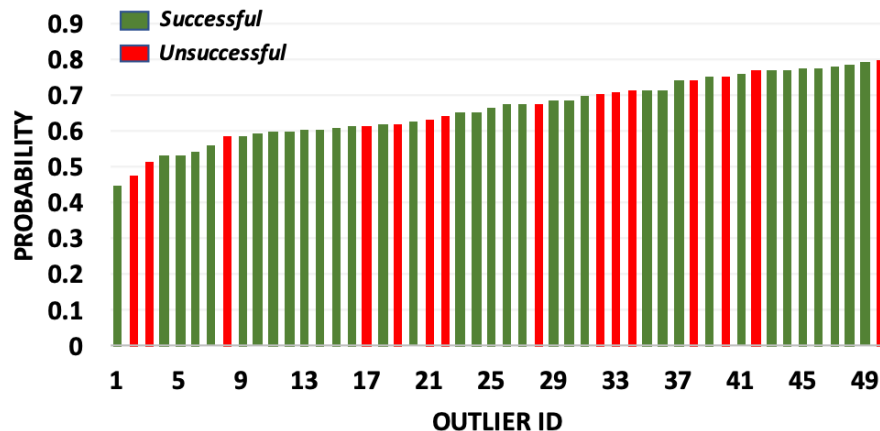
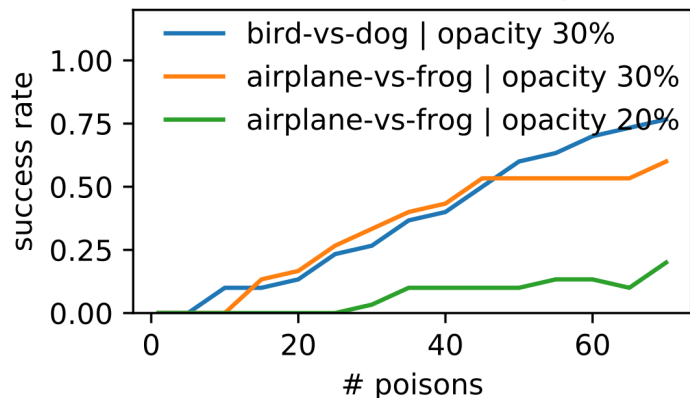
EVALUATIONS: END-TO-END LEARNING

- Setup
 - Dataset: CIFAR-10
 - Models: AlexNet (Pretrained on CIFAR-10)
- “end-to-end” Attacks
 - Goal: Bird > Dog or Airplane > Frog
 - Craft 1-70 poisons using the images chosen from the (Dog or Frog) class
 - Trick: watermarking!
 - Train the entire model on $D_{tr} \cup (x_p, y_p)$ and check the misclassification rate

EVALUATIONS: END-TO-END LEARNING

- Results

success rates of various experiments



HOW CAN WE IMPROVE THE TRANSFERABILITY OF CLEAN-LABEL ATT.?

METAPOISON! PRACTICAL GENERAL-PURPOSE CLEAN-LABEL DATA POISONING, HUANG ET AL., NEURIPS 2020

REVISIT: POISONING THREAT MODEL

- Goal
 - Targeted **clean-label** ($y_{c1} = y_{p1}$) attack
 - Model causes a misclassification of (x_t, y_t) , while preserving acc. on D_{val}
- Capability
 - Know a target (x_t, y_t)
 - Pick p candidates from test data (x_{c1}, y_{c1}) , (x_{c2}, y_{c2}) ... and craft poisons (x_{p1}, y_{p1}) , (x_{p2}, y_{p2}) ...
 - Inject them into the training data
- Knowledge
 - D_{tr} : training data
 - D_{test} : test-set data (validation data)
 - f : a model and its parameters θ
 - A : training algorithm (e.g., mini-batch SGD)

REVISIT: THE KEY IDEA – FEATURE COLLISION

- Goal

- Your poisons should work against any f and θ
- Objective:

$$\mathbf{p} = \operatorname{argmin}_{\mathbf{x}} \underbrace{\|f(\mathbf{x}) - f(\mathbf{t})\|_2^2}_{\text{feature collision}} + \beta \|\mathbf{x} - \mathbf{b}\|_2^2$$

Now you don't know the f , how can you estimate this?

- Revisit the previous idea

- Bi-level optimization

$$\begin{aligned} \operatorname{argmax}_{\mathcal{D}_p} \quad & \mathcal{W}(\mathcal{D}', \boldsymbol{\theta}_p^*), \\ \text{s.t.} \quad & \boldsymbol{\theta}_p^* \in \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}(\mathcal{D}_{\text{tr}} \cup \mathcal{D}_p, \boldsymbol{\theta}) \end{aligned}$$

$$\begin{aligned} X_p^* &= \operatorname{argmin}_{X_p} \mathcal{L}_{\text{adv}}(x_t, y_{\text{adv}}; \boldsymbol{\theta}^*(X_p)) \\ \boldsymbol{\theta}^*(X_p) &= \operatorname{argmin}_{\boldsymbol{\theta}} \mathcal{L}_{\text{train}}(X_c \cup X_p, Y; \boldsymbol{\theta}). \end{aligned}$$

Problem: no control over θ

THE CHALLENGE: LEARNING PROCESS

- Mode parameters are not fixed!
 - Initialization
 - Mini-batch-ed data
 - # of training epochs

Algorithm

Input: Examples $\{x_1, \dots, x_N\}$, loss function $\mathcal{L}(\theta) = \frac{1}{N} \sum_i \mathcal{L}(\theta, x_i)$. Parameters: learning rate η_t , noise scale σ , group size L , gradient norm bound C .

Initialize θ_0 randomly

for $t \in [T]$ **do**

Compute gradient

 For each $i \in L_t$, compute $\mathbf{g}_t(x_i) \leftarrow \nabla_{\theta_t} \mathcal{L}(\theta_t, x_i)$

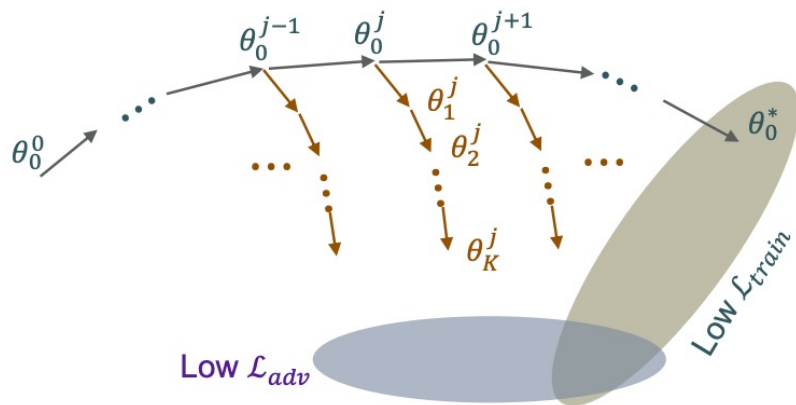
Descent

$\theta_{t+1} \leftarrow \theta_t - \eta_t \tilde{\mathbf{g}}_t$

Output θ_T and compute the overall privacy cost (ϵ, δ) using a privacy accounting method.

THE KEY IDEA: UNROLLING

- Goal
 - You *simulate all* the training procedures with *possible* f, θ s while crafting your poisons



Algorithm 1 Craft poison examples via MetaPoison

- 1: **Input** Training set of images and labels (X, Y) of size N , target image x_t , adversarial class y_{adv} , ϵ and ϵ_c thresholds, $n \ll N$ subset of images to be poisoned, T range of training epochs, M randomly initialized models.
- 2: **Begin**
- 3: Stagger the M models, training the m th model weights θ_m up to $\lfloor mT/M \rfloor$ epochs
- 4: Select n images from the training set to be poisoned, denoted by X_p . Remaining clean images denoted X_c
- 5: For $i = 1, \dots, C$ crafting steps:
 - 6: For $m = 1, \dots, M$ models:
 - 7: Copy $\tilde{\theta} = \theta_m$
 - 8: For $k = 1, \dots, K$ unroll steps^a:
 - 9: $\tilde{\theta} = \tilde{\theta} - \alpha \nabla_{\tilde{\theta}} \mathcal{L}_{train}(X_c \cup X_p, Y; \tilde{\theta})$
 - 10: Store adversarial loss $\mathcal{L}_m = \mathcal{L}_{adv}(x_t, y_{adv}; \tilde{\theta})$
 - 11: Advance epoch $\theta_m = \theta_m - \alpha \nabla_{\theta_m} \mathcal{L}_{train}(X, Y; \theta_m)$
 - 12: If θ_m is at epoch $T + 1$:
 - 13: Reset θ_m to epoch 0 and reinitialize
 - 14: Average adversarial losses $\mathcal{L}_{adv} = \sum_{m=1}^M \mathcal{L}_m / M$
 - 15: Compute $\nabla_{X_p} \mathcal{L}_{adv}$
 - 16: Update X_p using Adam and project onto ϵ, ϵ_c ball
 - 17: **Return** X_p

EVALUATION

- Setup

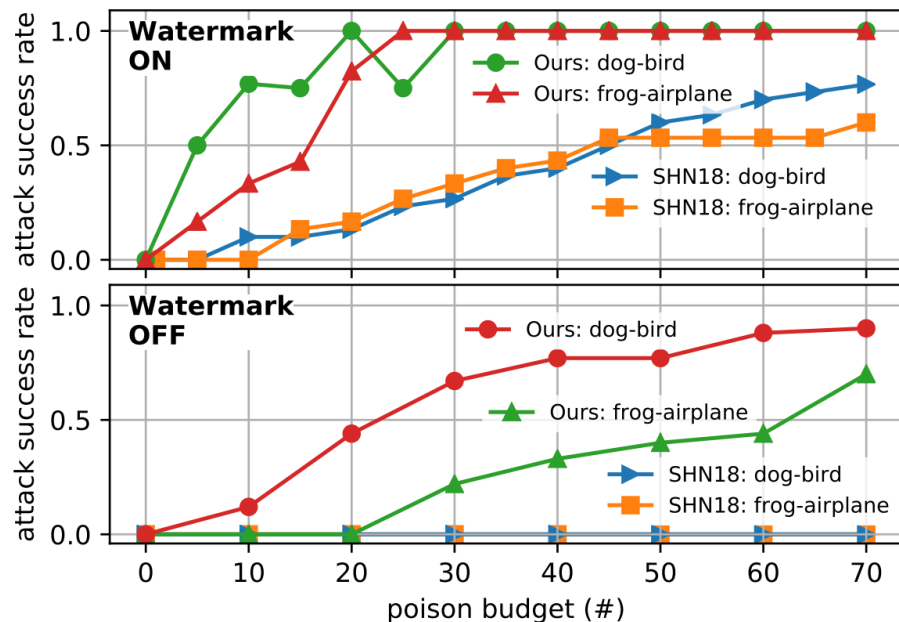
- Dataset: CIFAR-10
- Models: 6-layer ConveNet (default), ResNet20, VGG13
- Attack hyper-parameters:
 - C: 60 | M: 24 | K: 2

- Attacks

- 30 randomly chosen targets
- Increase the # poisons from 1 – 10% of the training data n
- Baseline:
 - Poison Frogs!

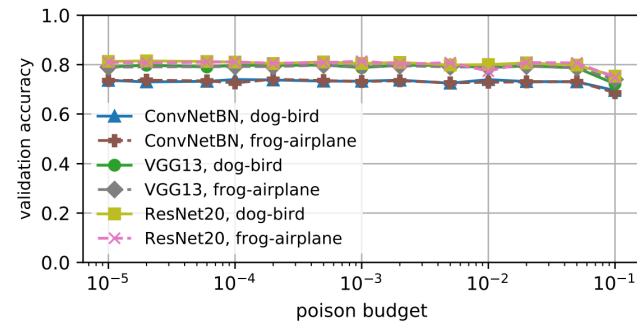
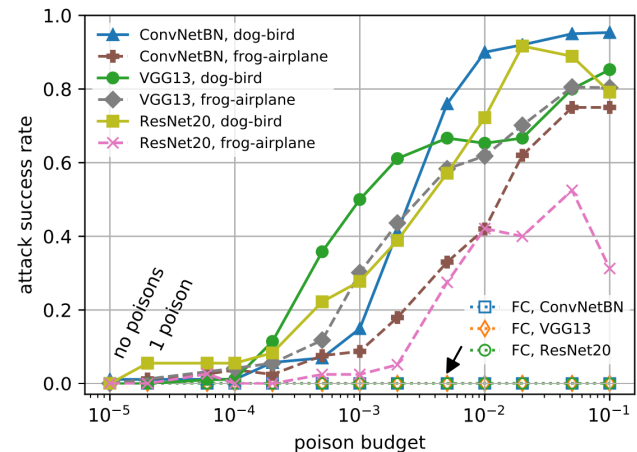
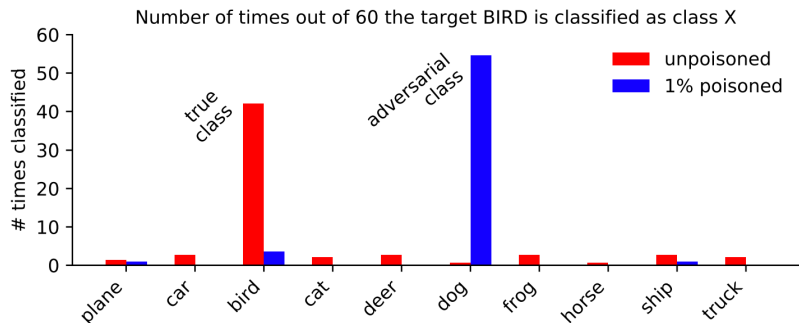
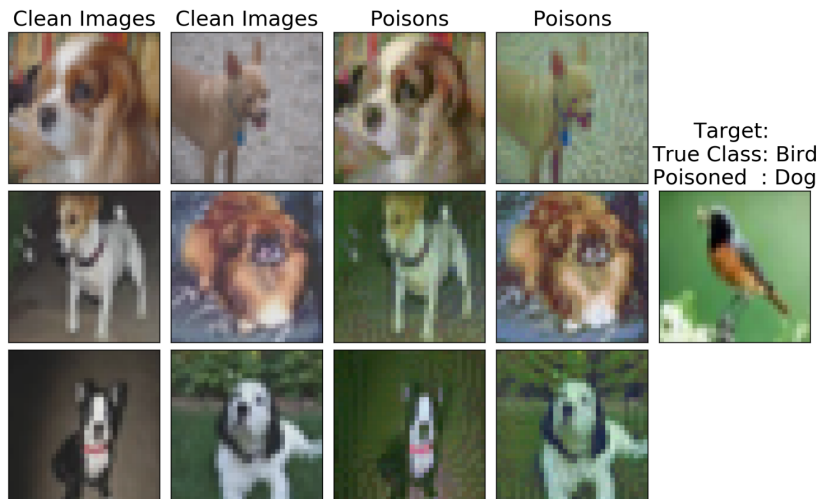
EVALUATION: TRANSFER LEARNING SCENARIO

- MetaPoison vs. Poison Frogs



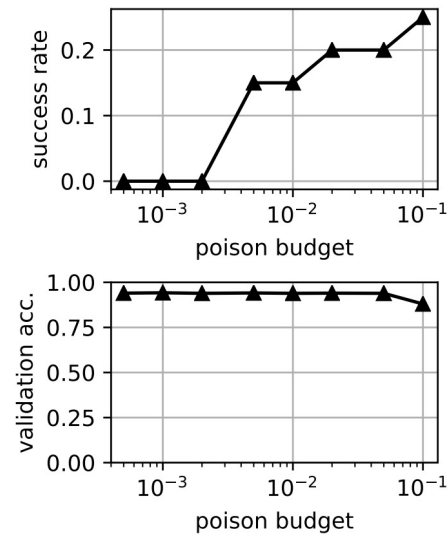
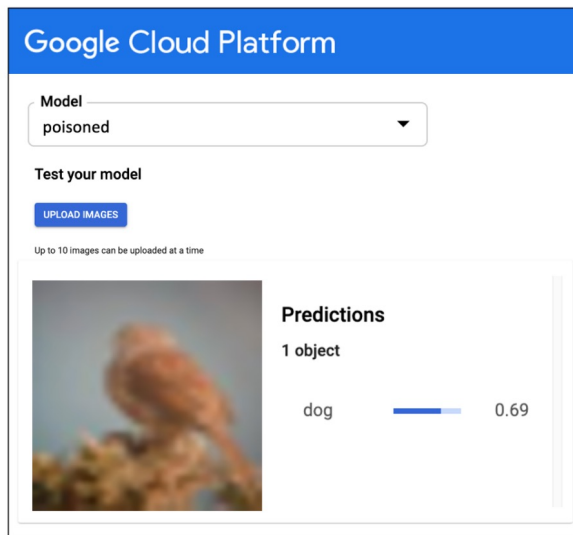
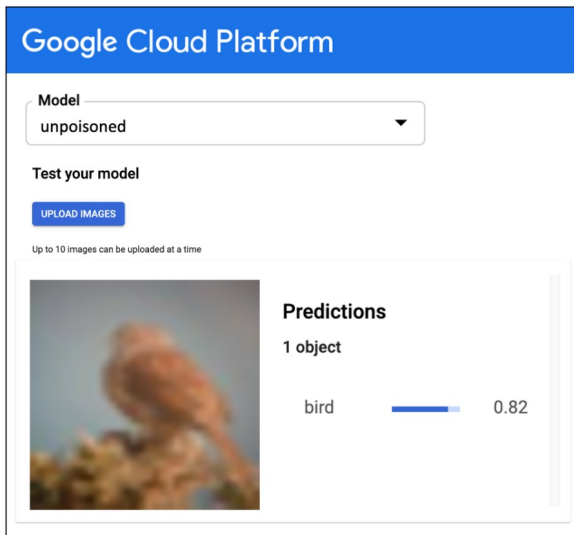
EVALUATION: END-TO-END SCENARIO

- MetaPo



EVALUATION: EXPLOITATION IN REAL-WORLD

- Results



Thank You!

Tu/Th 4:00 – 5:50 PM

Sanghyun Hong

<https://secure-ai.systems/courses/MLSec/F23>



Oregon State
University

SAIL
Secure AI Systems Lab