

CS 499/579: TRUSTWORTHY ML

MODEL STEALING

Tu/Th 4:00 – 5:50 pm

Sanghyun Hong

sanghyun.hong@oregonstate.edu

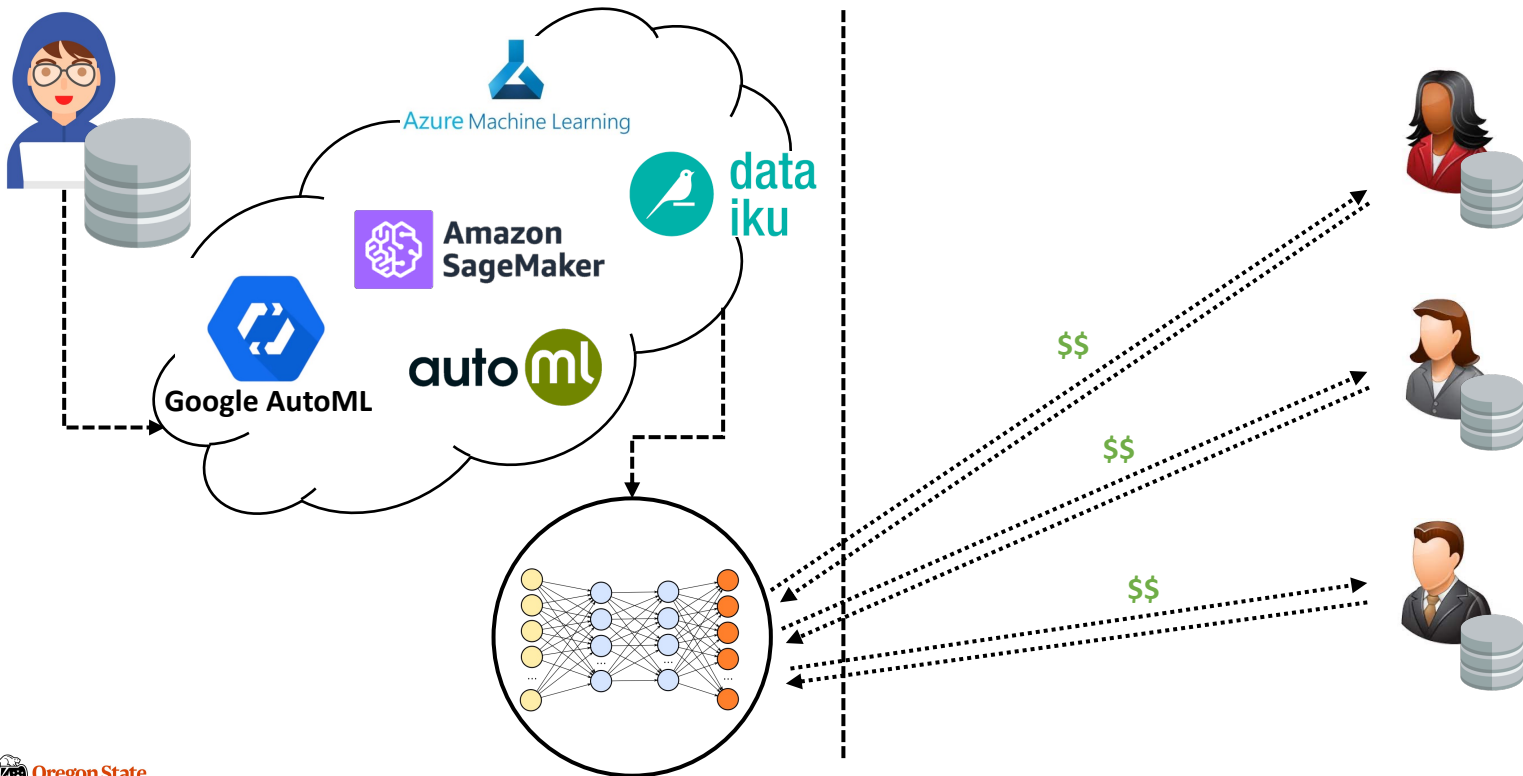


Oregon State
University

SAIL
Secure AI Systems Lab

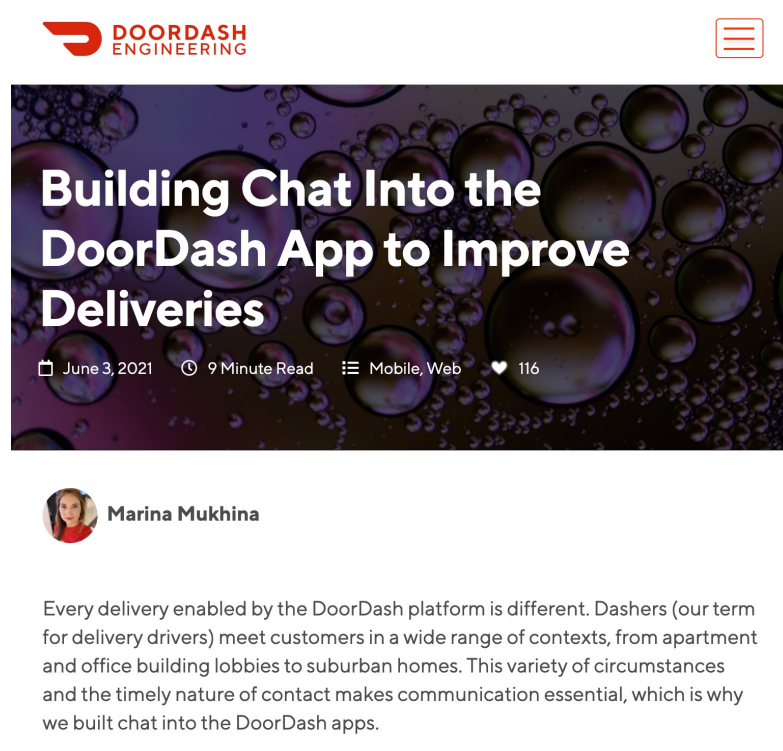
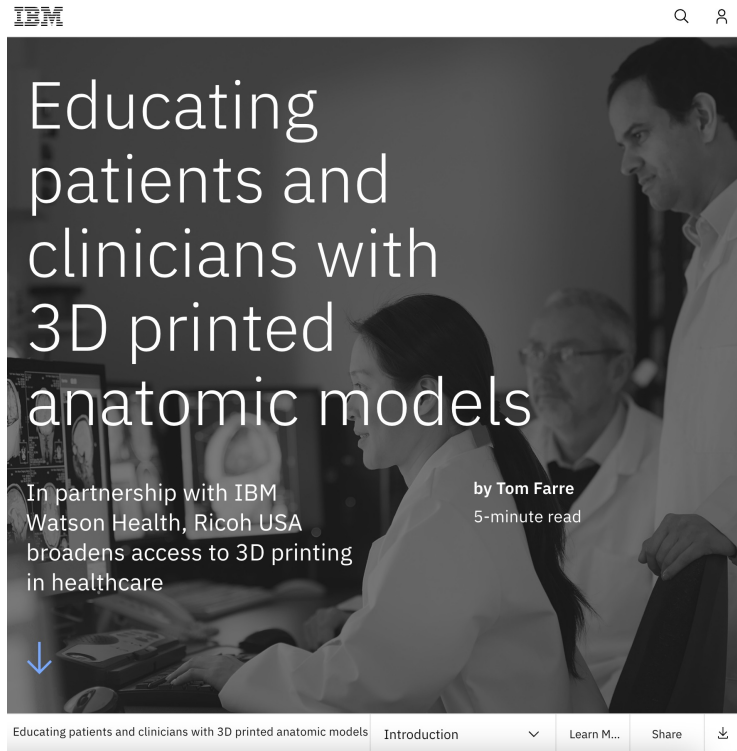
EMERGING MACHINE LEARNING AS A SERVICE (MLaaS)

- You train ML models and reach out to customers



MLaaS INCENTIVIZES MODEL EXTRACTION ATTACKERS

- Using **stolen** models... what if you run:



POTENTIAL DOWNSTREAM THREATS

- Exploiting stolen models, an adversary can:
 - Start a service with the stolen models with the same functionalities
 - Use the stolen model to craft adversarial examples
 - Extract private information from the stolen models

HOW CAN WE STEAL YOUR MODEL?

STEALING MACHINE LEARNING MODELS VIA PREDICTION APIS, TRAMER ET AL., USENIX SECURITY 2016

HOW CAN WE DO **HIGH-FIDELITY AND HIGH-ACCURACY** EXTRACTION?

HIGH ACCURACY AND HIGH-FIDELITY EXTRACTION OF NEURAL NETWORKS, JAGIELSKI ET AL., USENIX SECURITY 2020

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Threat model
 - Goal: Theft + *Reconnaissance
 - Theft: extraction of a target model
 - Reconnaissance: conduct downstream attacks, such as adversarial attacks

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Threat model
 - Goal: Theft (extraction attack)
 - Functionally-equivalent extraction, $\forall x, \hat{O}(x) = O(x)$
 - Fidelity extraction $\Pr_{x \sim D}[S(\hat{O}(x), O(x))]$, where $S(\cdot)$ is the similarity function
 - Task-accuracy extraction $\Pr_{(x,y) \sim D}[\operatorname{argmax}(\hat{O}(x)) = y]$

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Fidelity vs. task-accuracy
 - **Fidelity:** extracted model be *similar*
 - **Accuracy:** extracted model be *accurate*

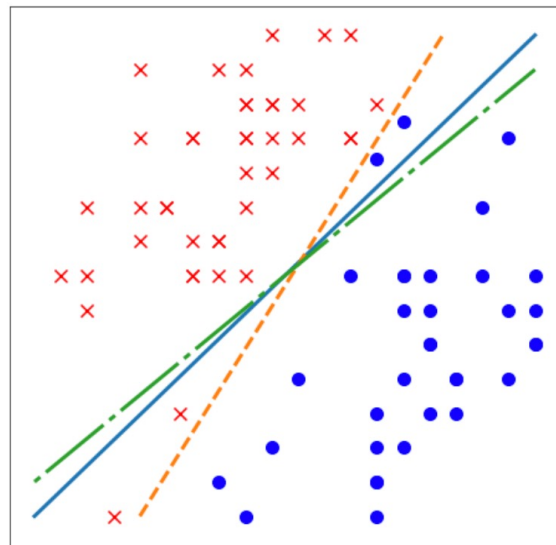


Figure 1: Illustrating fidelity vs. accuracy. The solid blue line is the oracle; functionally equivalent extraction recovers this exactly. The green dash-dot line achieves high fidelity: it matches the oracle on all data points. The orange dashed line achieves perfect accuracy: it classifies all points correctly.

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Threat model

- Goal: Theft (extraction attack)

- Functionally-equivalent extraction, $\forall x, \hat{O}(x) = O(x)$
 - Fidelity extraction $\Pr_{x \sim D}[S(\hat{O}(x), O(x))]$, where $S(\cdot)$ is the similarity function
 - Task-accuracy extraction $\Pr_{(x,y) \sim D}[\text{argmax}(\hat{O}(x)) = y]$

Attack	Type	Model type	Goal	Query Output
Lowd & Meek [8]	Direct Recovery	LM	Functionally Equivalent	Labels
Tramer <i>et al.</i> [11]	(Active) Learning	LM, NN	Task Accuracy, Fidelity	Probabilities, labels
Tramer <i>et al.</i> [11]	Path finding	DT	Functionally Equivalent	Probabilities, labels
Milli <i>et al.</i> [19] (theoretical)	Direct Recovery	NN (2 layer)	Functionally Equivalent	Gradients, logits
Milli <i>et al.</i> [19]	Learning	LM, NN	Task Accuracy	Gradients
Pal <i>et al.</i> [15]	Active learning	NN	Fidelity	Probabilities, labels
Chandrasekharan <i>et al.</i> [13]	Active learning	LM	Functionally Equivalent	Labels
Copycat CNN [16]	Learning	CNN	Task Accuracy, Fidelity	Labels
Papernot <i>et al.</i> [7]	Active learning	NN	Fidelity	Labels
CSI NN [25]	Direct Recovery	NN	Functionally Equivalent	Power Side Channel
Knockoff Nets [12]	Learning	NN	Task Accuracy	Probabilities
Functionally equivalent (this work)	Direct Recovery	NN (2 layer)	Functionally Equivalent	Probabilities, logits
Efficient learning (this work)	Learning	NN	Task Accuracy, Fidelity	Probabilities

*out of our scope

FUNCTIONALLY-EQUIVALENT EXTRACTION

- “Hard”
 - # of queries for extraction:
 - Suppose a neural network with $3k$ -width and 2-depth
 - On d -dimensional domain with precision of p numbers
 - The attacker needs $O(p^k)$ queries to perform a complete extraction
 - Check if two networks are the same
 - NP-hard problem
 - Learning-based approach struggles with fidelity
 - Suppose a deep random network with d -dimensional input and h -depth
 - Suppose an adversary formulated as statistical query (SQ) learning
 - Require $\exp(O(h))$ samples for fidelity extraction

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Threat model
 - Goal: Theft (extraction attack)
 - Functionally-equivalent extraction, $\forall x, \hat{O}(x) = O(x)$
 - Fidelity extraction $\Pr_{x \sim D}[S(\hat{O}(x), O(x))]$, where $S(\cdot)$ is the similarity function
 - Task-accuracy extraction $\Pr_{(x,y) \sim D}[\text{argmax}(\hat{O}(x)) = y]$
 - Knowledge
 - Domain knowledge:
 - The attacker has partial knowledge of the training dataset
 - They have some pretrained models in the same domain
 - Deployment knowledge
 - Model access

LEARNING-BASED MODEL EXTRACTION

- Fully-supervised model extraction

- Setup:

- Adversaries have access to some datasets
 - They use the victim model f as a labeling *oracle*
 - They train a separate model \hat{f} on the oracle outputs
 - Objective is to make \hat{f} and f achieve same test-time accuracy

- Experimental setup:

- Oracle: a model trained on 1B Instagram images (SoTA on ImageNet)
 - Attacker:
 - Case I: who has 10% (~13k) or 100% of the training samples (1B)
 - Case II: who improves the attack by using semi-supervised techniques (Rot. / MixMatch)

LEARNING-BASED MODEL EXTRACTION

- Evaluation results

- Results (+Rot.):

- Oracle (84.2% Top-1 acc. / 97.2% in Top-5)
 - Extracted models show a high accuracy (81- 94%) and fidelity (83- 97%) in Top-5
 - Semi-supervised approaches (unlabeled data) improve the performance further

Architecture	Data Fraction	ImageNet	WSL	WSL-5	ImageNet + Rot	WSL + Rot	WSL-5 + Rot
Resnet_v2_50	10%	(81.86/82.95)	(82.71/84.18)	(82.97/84.52)	(82.27/84.14)	(82.76/84.73)	(82.84/84.59)
Resnet_v2_200	10%	(83.50/84.96)	(84.81/86.36)	(85.00/86.67)	(85.10/86.29)	(86.17/88.16)	(86.11/87.54)
Resnet_v2_50	100%	(92.45/93.93)	(93.00/94.64)	(93.12/94.87)	N/A	N/A	N/A
Resnet_v2_200	100%	(93.70/95.11)	(94.26/96.24)	(94.21/95.85)	N/A	N/A	N/A

Problem: Non-determinism!

LEARNING-BASED MODEL EXTRACTION

- Evaluation results
 - Sources of non-determinism:
 - Initialization of model parameters
 - SGD (*random mini-batches)
 - Prior work on FE extraction attacks:
 - Milli *et al.*: *gradient queries*
 - Batina *et al.*: *power side-channel*

Query Set	Init & SGD	Same SGD	Same Init	Different
Test	93.7%	93.2%	93.1%	93.4%
Adv Ex	73.6%	65.4%	65.3%	67.1%
Uniform	65.7%	60.2%	59.0%	60.2%

Table 4: Impact of non-determinism on extraction fidelity. Even models extracted using the same SGD and initialization randomness as the oracle do not reach 100% fidelity.

Prior Work Assumes Too Strong Adversaries!

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

- Threat model

- Goal: Theft (extraction attack)

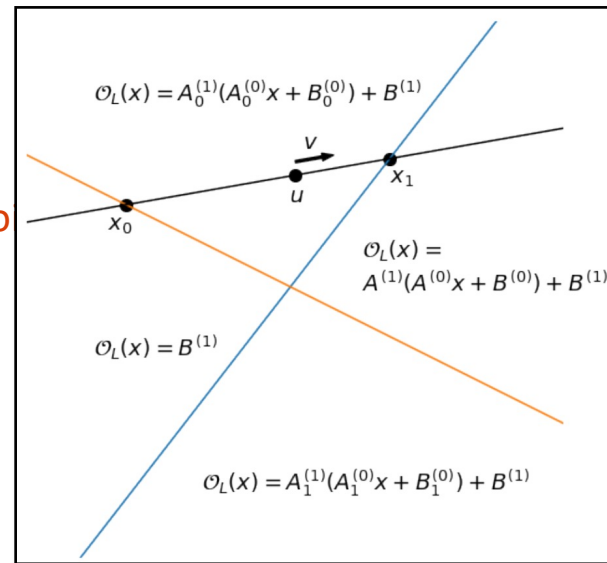
- Functionally-equivalent extraction, $\forall x, \hat{O}(x) = O(x)$
 - Fidelity extraction $\Pr_{x \sim D}[S(\hat{O}(x), O(x))]$, where $S(\cdot)$ is the similarity function
 - Task-accuracy extraction $\Pr_{(x,y) \sim D}[\operatorname{argmax}(\hat{O}(x)) = y]$

- Knowledge

- Domain knowledge:
 - The attacker has partial knowledge of the training dataset
 - They have some pretrained models in the same domain
 - Deployment knowledge
 - 2-layer feedforward neural network with ReLU activations
 - The architecture of a neural network is known (input-dim and hidden-dim)
 - Model access

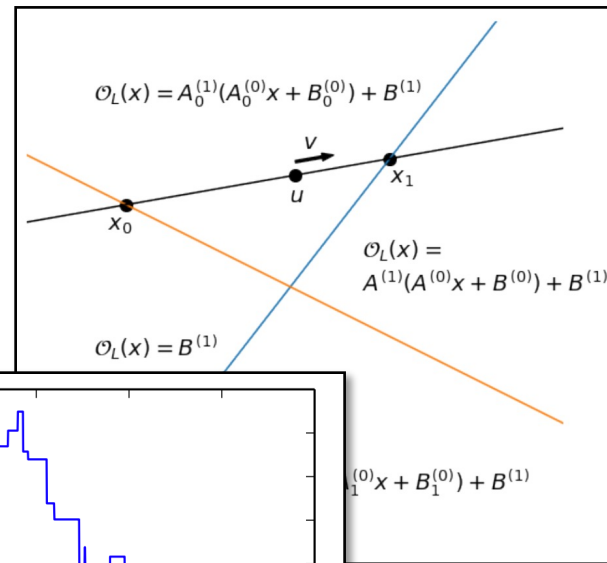
FUNCTIONALLY-EQUIVALENT MODEL EXTRACTION

- Jagielski *et al.* attack
 - Intuition (ReLU)
 - A standard choice of activation functions
 - It makes neural networks piecewise-linear (let's explore)
 - Attack procedures (on a 2-layer NN)
 - Critical point search
 - Weight recovery
 - Sign recovery
 - Final layer extraction



MODEL EXTRACTION ATTACK

- Jagielski *et al.* attack
 - Attack procedures (on a 2-layer NN)
 - Critical point search
 - Weight recovery
 - Sign recovery
 - Final layer extraction



Algorithm 1 Algorithm for 2-linearity testing. Computes the location of the only critical point in a given range or rejects if there is more than one.

Function f , range $[t_1, t_2]$, ϵ

$$m_1 = \frac{f(t_1 + \epsilon) - f(t_1)}{\epsilon} \quad \triangleright \text{Gradient at } t_1$$

$$m_2 = \frac{f(t_2) - f(t_2 - \epsilon)}{\epsilon} \quad \triangleright \text{Gradient at } t_2$$

$$y_1 = f(a), y_2 = f(b)$$

$$x = a + \frac{y_2 - y_1 - (b - a)m_2}{m_1 - m_2} \quad \triangleright \text{Candidate critical point}$$

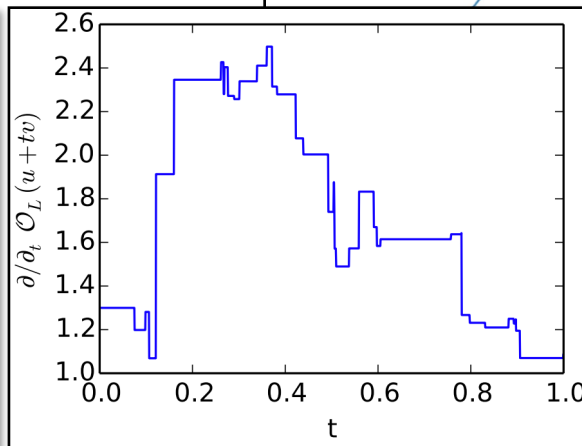
$$\hat{y} = y_1 + m_1 \frac{y_2 - y_1 - (b - a)m_2}{m_1 - m_2} \quad \triangleright \text{Expected value at candidate}$$

$$y = f(x) \quad \triangleright \text{True value at candidate}$$

if $\hat{y} = y$ **then return** x

else return "More than one critical point"

end if



MODEL EXTRACTION ATTACKS

- Jagielski *et al.* attack
 - Attack procedures (on a 2-layer NN)
 - Critical point search
 - **Weight recovery**
 - Compute second derivatives
 - Estimate the ratio between two weight vectors w_1, w_2
 - Sign recovery
 - Final layer extraction

$$\begin{aligned}\frac{\partial^2 O_L}{\partial e_j^2} \Big|_{x_i} &= \frac{\partial O_L}{\partial e_j} \Big|_{x_i + c \cdot e_j} - \frac{\partial O_L}{\partial e_j} \Big|_{x_i - c \cdot e_j} \\ &= \sum_k A_k^{(1)} \mathbb{1}(A_k^{(0)}(x_i + c \cdot e_j) + B_k^{(0)} > 0) A_{kj}^{(0)} \\ &\quad - \sum_k A_k^{(1)} \mathbb{1}(A_k^{(0)}(x_i - c \cdot e_j) + B_k^{(0)} > 0) A_{kj}^{(0)} \\ &= A_i^{(1)} \left(\mathbb{1}(A_i^{(0)} \cdot e_j > 0) - \mathbb{1}(-A_i^{(0)} \cdot e_j > 0) \right) A_{ji}^{(0)} \\ &= \pm (A_{ji}^{(0)} A_i^{(1)})\end{aligned}$$

MODEL EXTRACTION ATTACKS

- Jagielski *et al.* attack
 - Attack procedures (on a 2-layer NN)
 - Critical point search
 - Weight recovery
 - **Sign recovery**
 - **Final layer extraction**

$$\frac{\partial^2 O_L}{\partial(e_j + e_k)^2} \Big|_{x_i} = \pm(A_{ji}^{(0)} A_i^{(1)} \pm A_{ki}^{(0)} A_i^{(1)}).$$

EVALUATION

- Proposed attacks

- **Setup:**

- Datasets: MNIST and CIFAR-10
 - Models: 2-layer NN, 16 – 512 hidden units (~12 – 100k params)

- **Results:**

- MNIST:
 - 100% fidelity on the test-set
 - $2^{17.2} - 2^{20.2}$ queries for the 100% fidelity
 - CIFAR-10:
 - 100% fidelity on the test-set for models with < 200k params
 - 99% for the models with > 200k params
 - $2^{17.2} - 2^{20.2}$ queries for the 100% fidelity

EVALUATION

- Hybrid strategies

- **Setup:**

- Learning-based extraction with gradient matching
 - Error-recovery through learning

- **Results:**

- MNIST:
 - with 4x times larger models
 - 99-100% fidelity on the test-set
 - $2^{19.2} - 2^{22.2}$ queries for the 100% fidelity
(improvement over the previous results $2^{17.2} - 2^{20.2}$)

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