

AI 539 : TRUSTWORTHY ML

MODEL STEALING

Sanghyun Hong

sanghyun.hong@oregonstate.edu

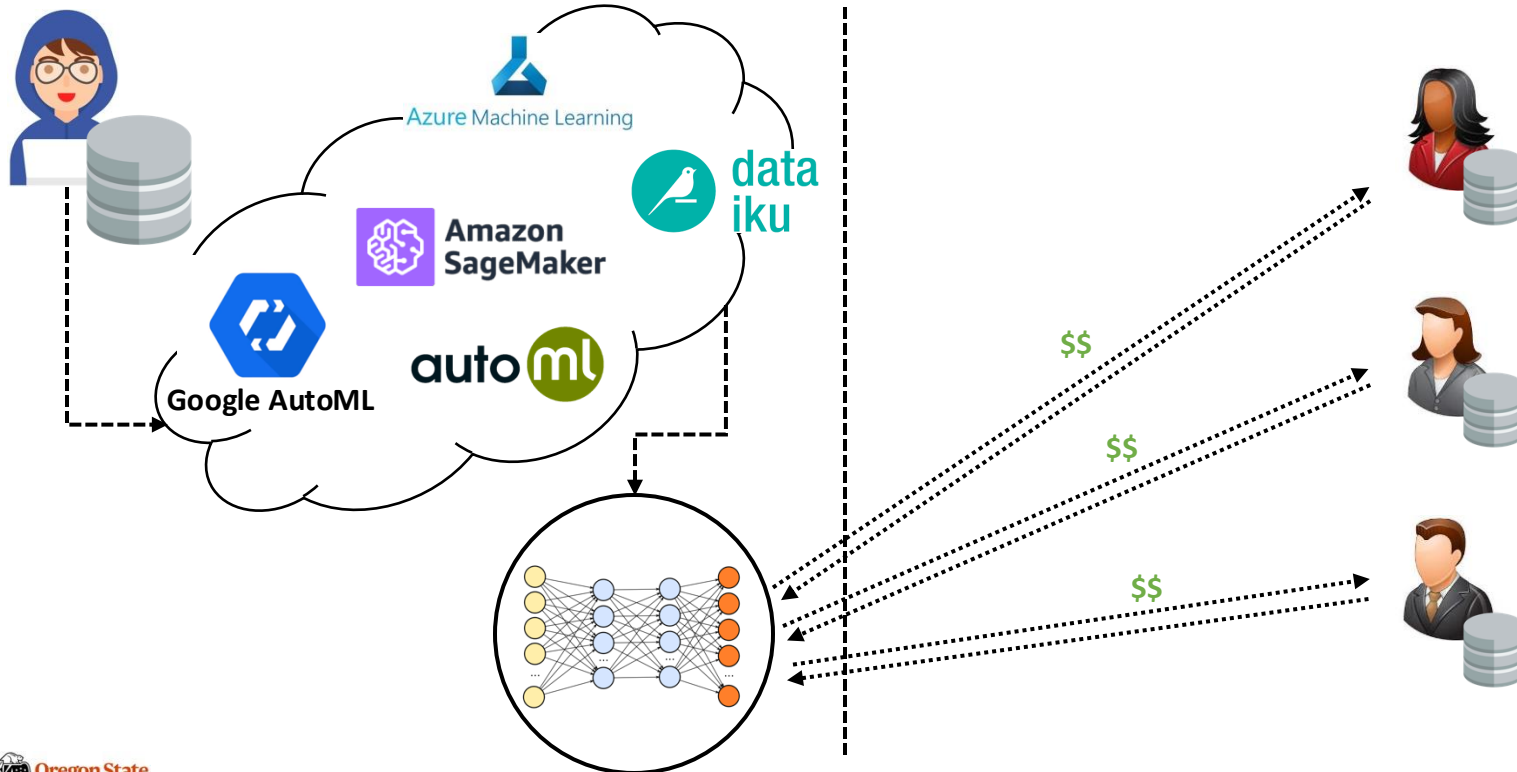


Oregon State
University

SAIL
Secure AI Systems Lab

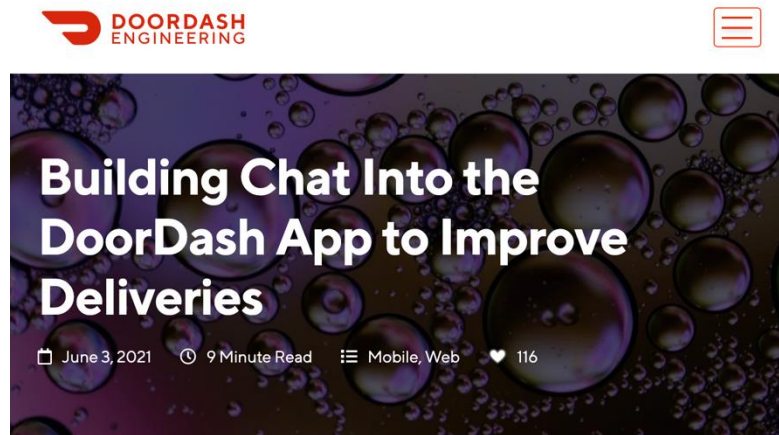
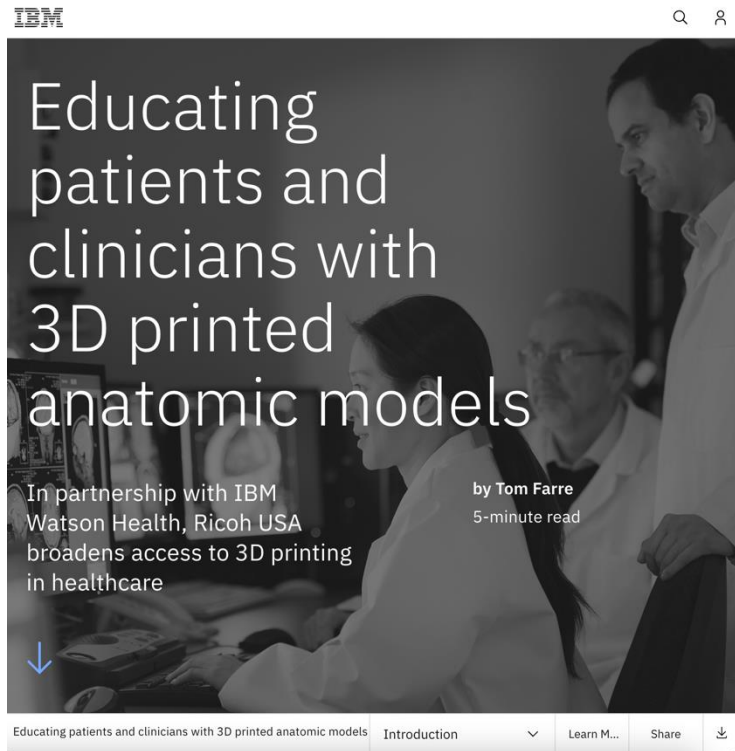
EMERGING MACHINE LEARNING AS A SERVICE (MLaaS)

- You train ML models and reach out to customers



MLAAS INCENTIVIZES MODEL EXTRACTION ATTACKERS

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



Marina Mukhina

Every delivery enabled by the DoorDash platform is different. Dashers (our term for delivery drivers) meet customers in a wide range of contexts, from apartment and office building lobbies to suburban homes. This variety of circumstances and the timely nature of contact makes communication essential, which is why we built chat into the DoorDash apps.

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

THREAT MODEL

- Model extraction attacks

- **Goal**

- To learn a new model \hat{f} that closely approximates the target model f

- **Knowledge**

- Black-box (typically)
 - It's possible to know aux. information:
 - How does a model extract feature(s)?
 - What is the model's class we aim to extract?
 - What is the training algorithm / hyper-params used?

Service	White-box	Monetize	Confidence Scores	Logistic Regression	SVM	Neural Network	Decision Tree
Amazon [1]	x	x	✓	✓	x	x	x
Microsoft [38]	x	x	✓	✓	✓	✓	✓
BigML [11]	✓	✓	✓	✓	x	x	✓
PredictionIO [43]	✓	x	x	✓	✓	x	✓
Google [25]	x	✓	✓	✓	✓	✓	✓

- **Capability**

- Has query access to the victim f (many times) with arbitrary inputs x
 - Has computational power to do offline processing of query outputs $f(x)$

THREAT MODEL

- Model extraction attacks

- **Metrics**

- Test error $R_{test}(\hat{f}, f)$: the average error between the outputs of \hat{f} and f on D
 - Uniform error $R_{unif}(\hat{f}, f)$: $R_{test}(\hat{f}, f)$ on a set of uniform vectors

- **Extraction accuracy:**

- $1 - R_{test}(\hat{f}, f) \mid 1 - R_{unif}(\hat{f}, f)$

MODEL EXTRACTION ATTACKS

- Equation-solving attack

- **Setup:**

- MLaaS APIs return confidence values $f(x)$
 - Those values are available to the attacker

- **Binary logistic regression:**

- Requires $d + 1$ predictions (queries), where d is the input dimension

- **Results:**

- Using $d + 1$ predictions, the attacker achieves the **errors $< 10^{-9}$**
 - The attacker **requires 41 – 113 queries** depending on the tasks

MODEL EXTRACTION ATTACKS

- Equation-solving attack

- **Setup:**

- MLaaS APIs return confidence values $f(x)$
 - Those values are available to the attacker

- **Multiclass LRs:**

- Softmax vs. one-vs-rest (OvR)
 - Requires $c(d + 1)$ queries, where c is the number of classes

- **Multi-layer perceptron (MLPs):**

- Requires $\alpha \cdot k$ predictions, where k is the number of unknown model parameters
 - Note: this work assumes MLPs with one hidden layer

MODEL EXTRACTION ATTACKS

- Equation-solving attack

- **Setup:**

- MLaaS APIs return confidence values $f(x)$
 - Those values are available to the attacker

- **Results:**

- MLRs: Using $c(d + 1)$ predictions, the attacker achieves the **errors $< 10^{-7}$**
 - MLPs: Require **5x times more queries** for achieving the same error rate

Model	Unknowns	Queries	$1 - R_{\text{test}}$	$1 - R_{\text{unif}}$	Time (s)
Softmax	530	265	99.96%	99.75%	2.6
		530	100.00%	100.00%	3.1
OvR	530	265	99.98%	99.98%	2.8
		530	100.00%	100.00%	3.5
MLP	2,225	1,112	98.17%	94.32%	155
		2,225	98.68%	97.23%	168
		4,450	99.89%	99.82%	195
		11,125	99.96%	99.99%	89

MODEL EXTRACTION ATTACKS

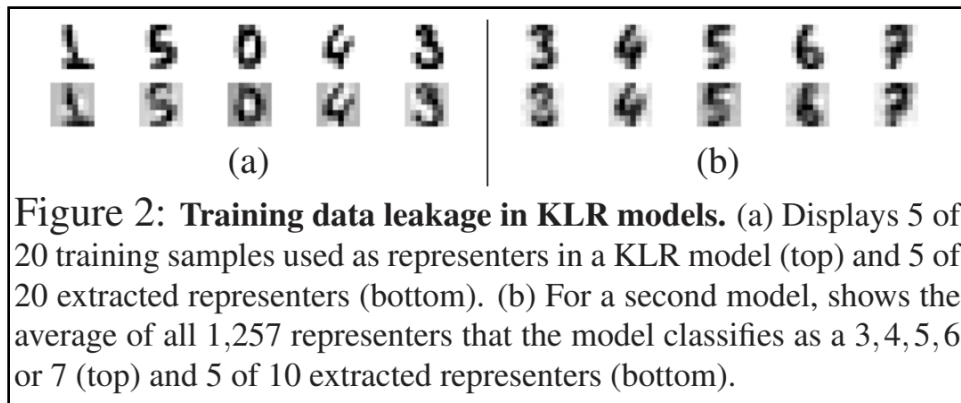
- Equation-solving attack

- **Setup:**

- MLaaS APIs return confidence values $f(\mathbf{x})$
 - Those values are available to the attacker

- **Downstream security attacks on \hat{f} :**

- Training data leakage in Kernel LR (KLR)
 - In KLR, the equation becomes $\sum_{r=1}^s \alpha_{i,r} K(\mathbf{x}, \mathbf{x}_r) + \beta_i$, where x_1, \dots, x_s are *representers*



MODEL EXTRACTION ATTACKS

- Equation-solving attack

- **Setup:**

- MLaaS APIs return confidence values $f(x)$
 - Those values are available to the attacker

- **Downstream security attacks on \hat{f} :**

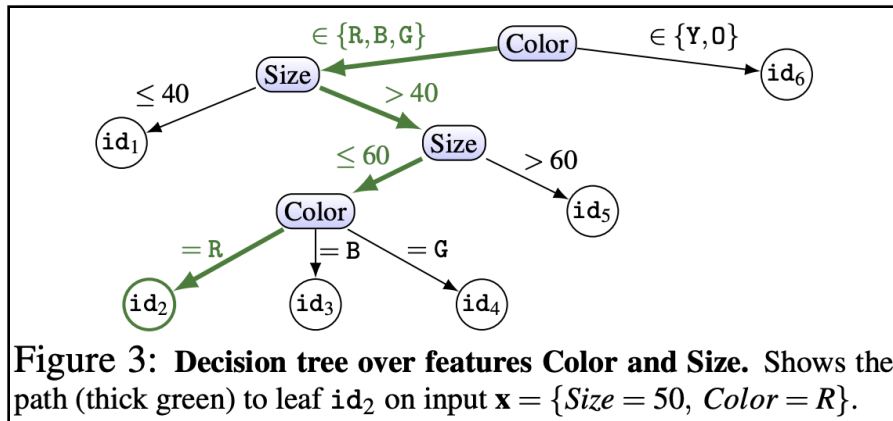
- Model inversion attacks
 - Convert a black-box to a white-box setting
 - In Fredrikson *et al.*
 - » The attack **requires 800k** queries to reconstruct 40 individuals
 - » One can extract the model **with 40k queries** and achieve the same attack success
 - » Using the extracted \hat{f} **reduces** the time **from 16 hrs to 10 hrs**

MODEL EXTRACTION ATTACKS

• Decision tree path-finding attack

– Setup:

- MLaaS APIs return $f(x)$ with
 - The leaf node
 - (for the incomplete queries) the node where
- Those values are available to the attacker



```

1:  $x_{init} \leftarrow \{x_1, \dots, x_d\}$                                 ▷ random initial query
2:  $Q \leftarrow \{x_{init}\}$                                        ▷ Set of unprocessed queries
3:  $P \leftarrow \{\}$                                              ▷ Set of explored leaves with their predicates
4: while  $Q$  not empty do
5:    $x \leftarrow Q.POP()$ 
6:    $id \leftarrow \mathcal{O}(x)$                                        ▷ Call to the leaf identity oracle
7:   if  $id \in P$  then                                           ▷ Check if leaf already visited
8:     continue
9:   end if
10:  for  $1 \leq i \leq d$  do                                       ▷ Test all features
11:    if IS_CONTINUOUS( $i$ ) then
12:      for  $(\alpha, \beta) \in \text{LINE\_SEARCH}(x, i, \epsilon)$  do
13:        if  $x_i \in (\alpha, \beta]$  then
14:           $P[id].ADD('x_i \in (\alpha, \beta]')$                 ▷ Current interval
15:        else
16:           $Q.PUSH(x[i] \Rightarrow \beta)$                     ▷ New leaf to visit
17:        end if
18:      end for
19:    else
20:       $S, V \leftarrow \text{CATEGORY\_SPLIT}(x, i, id)$ 
21:       $P[id].ADD('x_i \in S')$                                 ▷ Values for current leaf
22:      for  $v \in V$  do
23:         $Q.PUSH(x[i] \Rightarrow v)$                         ▷ New leaves to visit
24:      end for
25:    end if
26:  end for
27: end while
    
```

MODEL EXTRACTION ATTACKS

- Decision tree path-finding attack

- **Setup:**

- MLaaS APIs return $f(x)$ with
 - The leaf node
 - (for the incomplete queries) the node where each computation halts
 - Those values are available to the attacker

- **Results:**

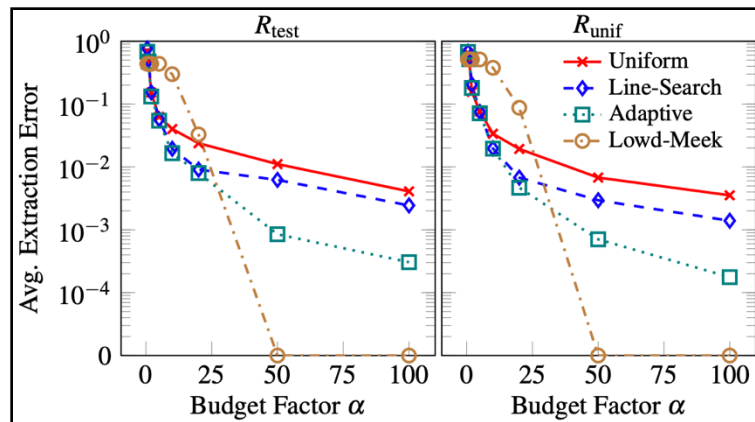
- All leaves are unique: 100% extraction success
 - Top-down: reduces # queries a lot & Duplicate leaves: a bit less effective

Model	Leaves	Unique IDs	Depth	Without incomplete queries			With incomplete queries		
				$1 - R_{\text{test}}$	$1 - R_{\text{unif}}$	Queries	$1 - R_{\text{test}}$	$1 - R_{\text{unif}}$	Queries
IRS Tax Patterns	318	318	8	100.00%	100.00%	101,057	100.00%	100.00%	29,609
Steak Survey	193	28	17	92.45%	86.40%	3,652	100.00%	100.00%	4,013
GSS Survey	159	113	8	99.98%	99.61%	7,434	100.00%	99.65%	2,752
Email Importance	109	55	17	99.13%	99.90%	12,888	99.81%	99.99%	4,081
Email Spam	219	78	29	87.20%	100.00%	42,324	99.70%	100.00%	21,808
German Credit	26	25	11	100.00%	100.00%	1,722	100.00%	100.00%	1,150
Medical Cover	49	49	11	100.00%	100.00%	5,966	100.00%	100.00%	1,788
Bitcoin Price	155	155	9	100.00%	100.00%	31,956	100.00%	100.00%	7,390



MODEL EXTRACTION ATTACKS

- What if...
 - **Setup:**
 - MLaaS APIs do *not* return confidence values $f(x)$
 - The adversary can only observe labels
 - **Adaptive Attacks:**
 - The Lowd-Meek attack (~line-search)
 - Re-training approach (~train a model on $(x, f(x))$)
 - Re-training with uniform queries
 - Line-search retraining
 - Adaptive retraining
 - **Results:**
 - on LR models
 - on MLR or MLP



MODEL EXTRACTION ATTACKS

- Countermeasures

- Rounding confidences:**

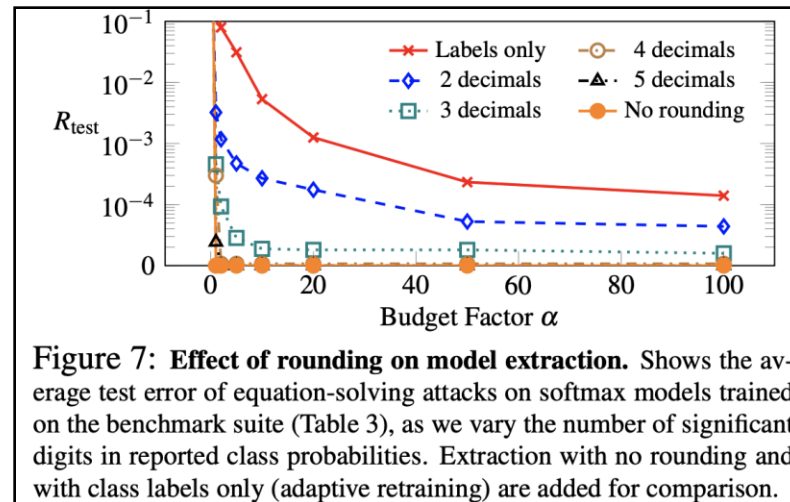
- On LRs, MLRs and MLPs
 - On decision trees: node collision

- Differential privacy:**

- Ugh...
 - It's not designed to prevent extractions

- Ensemble methods:**

- The adversary can approximate the ensemble itself



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 $\Pr_{(x,y) \sim D}[\text{argmax}(\hat{O}(x)) = y]$

TAXONOMY OF EXISTING MODEL EXTRACTION ATTACKS

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

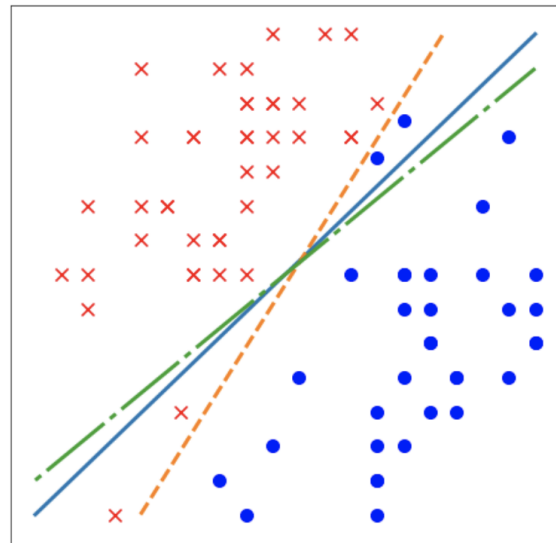


Figure 1: Illustrating fidelity vs. accuracy. The solid blue line is the oracle; functionally equivalent extraction recovers this exactly. The green dash-dot line achieves high fidelity: it matches the oracle on all data points. The orange dashed line achieves perfect accuracy: it classifies all points 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 $\Pr_{(x,y) \sim D}[\arg\max(\hat{O}(x)) = y]$

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 $\Pr_{(x,y) \sim D}[\arg\max(\hat{O}(x)) = y]$
 - 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.

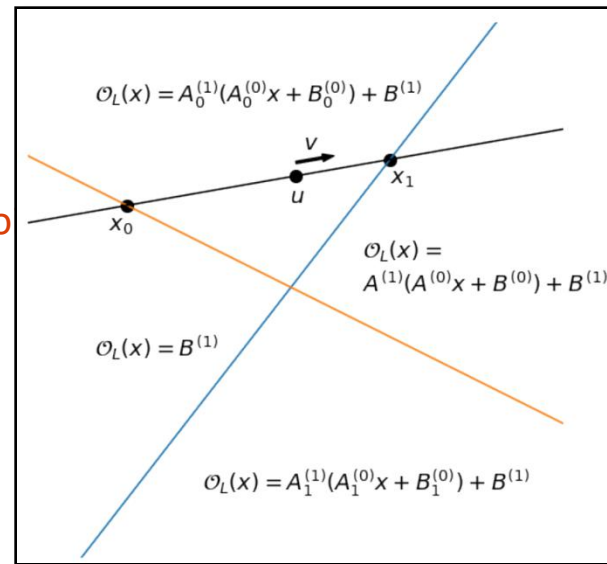
Prior Work Assumes Too Strong Adversaries!

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 $\Pr_{(x,y) \sim D}[\arg\max(\hat{O}(x)) = y]$
 - 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

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

$m_2 = \frac{f(t_2) - f(t_2 - \epsilon)}{\epsilon}$ ▷ Gradient at t_2

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

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

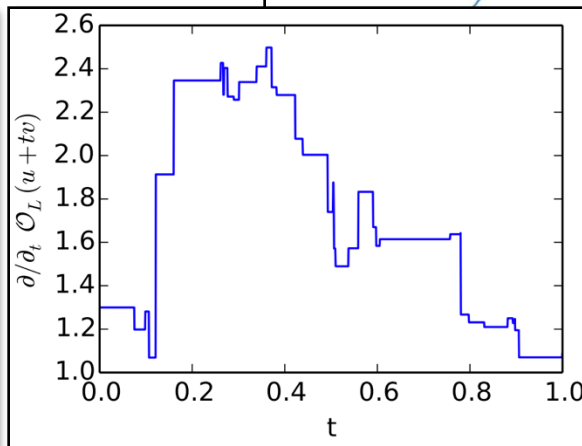
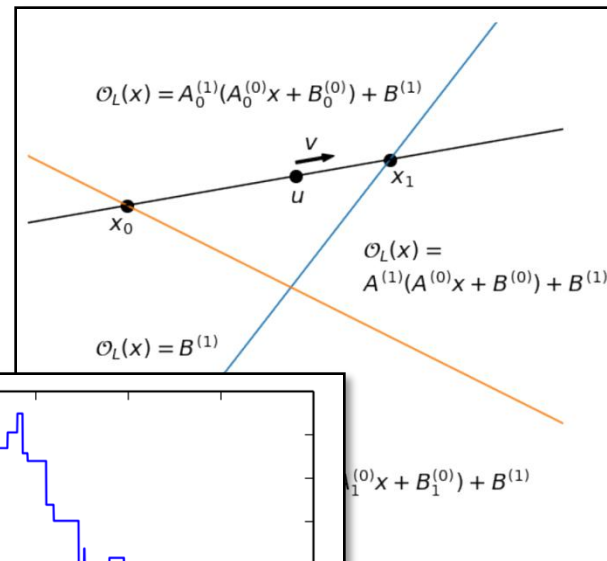
$\hat{y} = y_1 + m_1 \frac{y_2 - y_1 - (b - a)m_2}{m_1 - m_2}$ ▷ Expected value at candidate

$y = f(x)$ ▷ 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**

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

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!

Sanghyun Hong

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



Oregon State
University

SAIL
Secure AI Systems Lab