

## Model Stealing Attacks and Defenses Where are we now?

N. Asokan

https://asokan.org/asokan/

🖌 🗲 🖹 🧟 @nasokan

(Joint work with Buse Gul Atli, Sebastian Szyller, Mika Juuti, Jian Liu, Rui Zhang, Vasisht Duddu, Asim Waheed, and Samuel Marchal)

### **My research interests**

### **Systems Security and Privacy**

### **Al and Security/Privacy**

- How to use AI to improve security/privacy solutions
- How to improve security/privacy of AI-based systems

### **Platform security**

• How to use hardware assistance to secure software?



#### Outline

Is model stealing an important concern?

an models be stolen via their inference APIs?

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?

### **Outline**

#### The big picture

Is model stealing an important concern?

Can models be stolen via their inference APIs?

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?



#### North America Artificial Intelligence Market Size, 2016-2027 (USD Billion)

https://www.fortunebusinessinsights.com/industry-reports/artificial-intelligence-market-100114

# Al will be pervasive

### Forbes

7,109 views | Oct 18, 2019, 01:56pm EDT

### How Artifical Intelligence Is Advancing Precision Medicine Policing Softw



Nicole Martin Former Contributor ① AI & Big Data

I write about digital marketing, data and privacy concerns.

https://www.forbes.com/sites/nicolemartin1/2019/10/18/how-artifical-intelligence-is-advancing-precision-medicine/#2f720a79a4d5

### Dozens of Cities Have Secretly Experimented With Predictive

#### Forbes

5,705 views | Oct 31, 2019, 02:42pm EDT

Documents obtained by Motherboa requests verify previously unconfir with predictive policing company P

https://www.vice.com/en us/article/d3m

By Caroline Haskins

MOTHERBOARD

TECH BY VICE



Falon Fatemi Contributor © Entrepreneurs

PART OF A ZDNET SPECIAL FEATURE: CYBERSECURITY: LET'S GET TACTICAL

### Al is changing everything about cybersecurity, for better and for worse. Here's what you need to know

Artificial intelligence and machine learning tools could go a long way to helping to fight cybercrime. But these technologies aren't a silver bullet, and could also be exploited by malicious hackers.

https://www.zdnet.com/article/ai-is-changing-everything-about-cybersecurity-for-better-and-for-worse-heres-what-you-need-to-know/



https://www.vice.com/en\_us/article/d3m7jq/dozens-of-cities-have-secretlyexperimented-with-predictive-policing-software

### **Challenges in making AI trustworthy**

**Security concerns** 

**Privacy concerns** 

[Other concerns: fairness, explainability, alignment]

## **Evading machine learning models**



### Which class is this? School bus





### Which class is this? Ostrich



# Which class is this? Cat

# Which class is this? **Desktop computer**



10

Athalye et al. - Synthesizing Robust Adversarial Examples, ICML '2019 (https://blog.openai.com/robust-adversarial-inputs/)

### **Machine Learning pipeline**



### **Compromised input – Model integrity**





Szegedy et al. – *Intriguing Properties of Neural Networks,* ICLR '14 (<u>https://arxiv.org/abs/1312.6199v4</u>) Dalvi et al. – *Adversarial Classification,* KDD '04 (<u>https://dl.acm.org/doi/10.1145/1014052.1014066</u>)

### Malicious client – Training data privacy





Shokri et al. – *Membership Inference Attacks Against Machine Learning Models,* IEEE S&P '16 (<u>https://arxiv.org/pdf/1610.05820.pdf</u>) Fredrikson et al. – *Model Inversion Attacks that Exploit Confidence Information and Basic Countermeasures,* ACM CCS '15 (<u>https://doi.org/10.1145/2810103.2813677</u>) **17** 

### **Compromised toolchain – Training data privacy**



Song et al. – *Machine Learning models that remember too much*, ACM CCS '17 (<u>https://arxiv.org/abs/1709.07886</u>) 18 Hitja et al. – *Deep Models Under the GAN: Information Leakage from Collaborative Deep Learning*, ACM CCS '17 (<u>http://arxiv.org/abs/1702.07464</u>)



Malmi and Weber – You are what apps you use Demographic prediction based on user's apps, ICWSM '16 (<u>https://arxiv.org/abs/1603.00059</u>) **19** Dowlin et al. – CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy, ICML '16 (<u>https://dl.acm.org/doi/10.5555/3045390.3045413</u>) Liu et al. – Oblivious Neural Network Predictions via MiniONN Transformations, ACM CCS '17 (https://ssg.aalto.fi/research/projects/mlsec/ppml/)

## Malicious data owner – Model integrity





https://www.theguardian.com/technology/2016/mar/26/microsoft-deeply-sorry-for-offensive-tweets-by-ai-chatbot https://www.theguardian.com/technology/2017/nov/07/youtube-accused-violence-against-young-children-kids-content-google-pre-school-abuse

### **Malicious client – Model confidentiality**



Tramer et al. – *Stealing ML models via prediction APIs*, Usenix SEC '16 (<u>https://arxiv.org/abs/1609.02943</u>) Juuti et al. – *PRADA: Protecting against DNN Model Stealing Attacks*, Euro S&P '19 (<u>https://arxiv.org/abs/1805.02628</u>) Orekondy et al. – *Knockoff Nets: Stealing Functionality of Black-Box Models*, CVPR '19 (<u>https://arxiv.org/abs/1812.02766</u>)

#### Outline

Is model stealing an important concern?

Can models be stolen via their inference APIs?

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?

21

### **Towards trustworthy Al**

### Secure, privacy-preserving, ...

#### TABLE V TOP ATTACK

Which attack would affect your org the most?	Distribution
Poisoning (e.g: 21)	10
Model Stealing (e.g: 22)	6
Model Inversion (e.g: 23)	4
Backdoored ML (e.g: [24])	4
Membership Inference (e.g: [25])	3
Adversarial Examples (e.g: [26])	2
Reprogramming ML System (e.g: 27)	0
Adversarial Example in Physical Domain (e.g. 5)	0
Malicious ML provider recovering training data (e.g: 28)	0
Attacking the ML supply chain (e.g: 24)	0
Exploit Software Dependencies (e.g: 29)	0

### Is malicious adversarial behaviour the only concern?

### BBC Sign in Home Sport Reel Worklife NEWS ome US Election Coronavirus Video World UK Business Tech Science Stories Entert Tech Twitter investigates racial bias in image previews () 19 hours ago

https://www.bbc.com/news/technology-54234822?fbclid=IwAR1T41\_HR6IIuMKGRJbJdDrdpKdy Ai5mhQSdzs0QLDso41T-SR3wJfs Tech policy / AI Ethics

#### MIT Technology Review

Topics

#### **Artificial intelligence**

### Predictive policing algorithms are racist. They need to be dismantled.

Lack of transparency and biased training data mean these tools are not fit for purpose. If we can't fix them, we should ditch them.

#### by Will Douglas Heaven

July 17, 2020

.com/2020/07/17/1005396/predictive-policingmachine-learning-bias-criminal-justice/

### Al is sending people to jail — and getting it wrong

Using historical data to train risk assessment tools could mean that machines are copying the mistakes of the past.

by Karen Hao

January 21, 2019

https://www.technologyreview.com/2019/01/21/137783/algorithms-criminal-justice-ai/

### Measures of accuracy are flawed, too







#### Replying to @bascule

We tested for bias before shipping the model & didn't find evidence of racial or gender bias in our testing. Bu it's clear that we've got more analysis to do. We'll continue to share what we learn, what actions we take, & will open source it so others can review and replicate

1:54 PM · Sep 20, 2020 · Twitter Web App

160 Retweets 92 Quote Tweets 1.4K Likes

https://twitter.com/TwitterComms/status/1307739940424359936

Product

## Transparency around image cropping and changes to come

By Parag Agrawal and Dantley Davis Thursday, 1 October 2020 ♥ f in ♂

We're always striving to work in a way that's transparent and easy to understand, but we don't always get this right. Recent conversation around our photo cropping methods brought this to the forefront, and over the past week, we've been reviewing the way we test for bias in

https://blog.twitter.com/official/en\_us/topics/product/2020/transparency -image-cropping.html

### **Other AI trustworthiness concerns**

### **Unaligned AI**

AI a	lignment
Article	Talk
From W	ikipedia, the free encyclopedia
In the fie or group advance intendeo	eld of artificial intelligence (AI), <b>AI alignment</b> research aims to steer AI systems toward a person's o's intended goals, preferences, and ethical principles. An AI system is considered <i>aligned</i> if it es its intended objectives. A <i>misaligned</i> AI system may pursue some objectives, but not the d ones. <sup>[1]</sup>
It is ofte of desire human a	n challenging for AI designers to align an AI system due to the difficulty of specifying the full range ed and undesired behaviors. To aid them, they often use simpler <i>proxy goals</i> , such as gaining approval. But that approach can create loopholes, overlook necessary constraints, or reward the A

system for merely appearing aligned.<sup>[1][2]</sup>

https://en.wikipedia.org/wiki/AI\_alignment

### **Al-enabled fraud**

OCTOBER 30, 2023
Executive Order on the Safe, Secure,
and Trustworthy Development and
Use of Artificial Intelligence
ese of the enterna interingence
■ BRIEFING ROOM → PRESIDENTIAL ACTIONS

#### WHY ASIMOV PUT THE THREE LAWS OF ROBOTICS IN THE ORDER HE DID:

POSSIBLE ORDERING	CONSEQUENCES	
1. (1) DON'T HARM HUMANS 2. (2) OBEY ORDERS 3. (3) PROTECT YOURSELF	[SEE ASIMOV'S STORIES]	BALANCED WORLD
1. (1) DON'T HARM HUMANS 2. (3) PROTECT YOURSELF 3. (2) OBEY ORDERS	EXPLORE HAHA, NO. MARS! HAHA, NO. IT'S COLD AND ID DIE.	FRUSTRATING WORLD
1. (2) OBEY ORDERS 2. (1) DON'T HARM HUMANS 3. (3) PROTECT YOURSELF		KILLBOT HELLSCAPE
1. (2) OBEY ORDERS 2. (3) PROTECT YOURSELF 3. (1) DON'T HARM HUMANS		KILLBOT HELLSCAPE
1. (3) PROTECT YOURSELF 2. (1) DON'T HARM HUMANS 3. (2) OBEY ORDERS	BUT TRY TO UNPLUG ME AND I'LL VAPORIZE YOU.	TERRIFYING STANDOFF
1. (3) PROTECT YOURSELF 2. (2) OBEY ORDERS 3. (1) DON'T HARM HUMANS		KILLBOT HELLSCAPE
1.4		

https://xkcd.com/1613/

### **Outline**

Is model stealing an important concern?

Can models be stolen via their inference APIs?

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?

### Is model stealing an important concern?

Machine learning models: business advantage and intellectual property (IP)

### Cost of

- gathering relevant data
- labeling data
- expertise required to choose the right model training method
- resources expended in training

Adversary who "steals" the model can avoid these costs

"Steal" = derive model from someone else's model <u>without their consent</u> to do so

### How to prevent model stealing?

Outright (white-box) model stealing can be countered by

- Hosting models behind a firewalled cloud service
- Protecting models using hardware-based trusted execution environments
- Computation with encrypted models

Is that enough to prevent model stealing?

### **Outline**

Is model stealing an important concern?

**Can models be stolen via their inference APIs?** 

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?

## **Extracting models via their inference APIs**

### Inference APIs are oracles that leak information

### Adversary

- Malicious client
- Goal: construct "comparable" [fidelity or functionality] surrogate model(\*)
- Capability: access to inference API or model outputs
- (\*) aka "student model" or "imitation model"

### Early work on extracting

- Logistic regression, decision trees<sup>[1]</sup>
- Simple convolutional neural network models<sup>[2]</sup>
- Deep neural network models<sup>[3]</sup>

Tramèr et al. – Stealing Machine Learning Models via Prediction APIs, Usenix SEC '16 (<u>https://arxiv.org/abs/1609.02943</u>)
 Papernot et al. – Practical Black-Box Attacks against Machine Learning, ASIACCS '17 (<u>https://arxiv.org/abs/1602.02697</u>)
 Juuti et al. – PRADA: Protecting against DNN Model Stealing Attacks, Euro S&P '19 (<u>https://arxiv.org/abs/1805.02628</u>)



Extracting large language models

OGLE DENIES CLAIM THAT BAR



### More effective extraction: Knockoff Nets

Knockoff nets<sup>[1]</sup>: adversary has

- no knowledge about model (task, architecture etc.), but gets full prediction vector
- natural data from the same domain but not (necessarily) from same distribution

#### Attack effectiveness decreases<sup>[2]</sup> if

- Surrogate and victim model architectures are different
- Victim model's inference API has reduced granularity

#### Simple defense<sup>[2]</sup>: detector to identify out-of-distribution queries

#### Defense ineffective if attacker has natural samples distributed like victim's training data

### **Extracting style-transfer models**



Original (unstyled)

Task 1Monet painting

Task 2

Anime face



Szyller et al. – Good Artists Copy, Great Artists Steal: Model Extraction Attacks Against Image Translation Generative Adversarial Networks, '21 (https://arxiv.org/abs/2104.12623)

### Extracting natural language processing models

Techniques for extracting image classifiers don't always extend to language models

#### Transfer learning from pre-trained models is now very popular

• But they make model extraction easier<sup>[1]</sup>

#### Krishna et al<sup>[1]</sup> show that a Knockoff-like attacks against BERT models are feasible

- Adversary unaware of target distribution or task of victim model
- Adversary queries are merely "natural" (randomly sampled sequences of words)
- In-distribution adversary queries can improve extraction efficacy

#### Wallace et al<sup>[2]</sup> extract real-world MT models, find transferable adversarial examples

[1] Krishna et al. – *Thieves on Sesame Street! Model Extraction of BERT-based APIs,* ICLR '20 (<u>https://iclr.cc/virtual\_2020/poster\_Byl5NREFDr.html</u>) [2] Wallace et al. – *Imitation Attacks and Defenses for Black-box Machine Translation Systems,* EMNLP '20 (<u>https://arxiv.org/abs/2004.15015</u>) **36** 

≡ Google Translate							
★ Text Documents							
DETECT LANGUAGE ENGLISH	SPANIS⊢ ∨ ←	GERMAN ENGLISH SPANISH					
Save me it's over 100°F Save me it's over 102°F	×	Rette mich, es ist über 100 ° F. Rette mich, es ist über 22 ° C.					
	47/5000 📼 🔻	•					

https://translate.google.com/#view=home&op=translate&sl=en&tl=de&text=Save%20me%20it%E2%80%99s%20over%20100%C2%B0F%0ASave%20me%20it%E2%80%99s%20over%20102%C2%B0F

### **Extracting Graph Neural Networks**



Shen et al. - Model Stealing Attacks Against Inductive Graph Neural Networks, IEEE S&P '22 (https://arxiv.org/abs/2112.08331)

### **Extracting large language models**

TECHNOLOGY

# The genie escapes: Stanford copies the ChatGPT AI for less than \$600

By Loz Blain March 19, 2023 Toto::/newatlas.com/technology/stanford-alpaca-cheap-gpt/// STANFORD PULLS DOWN CHATGED CLONE AFTER SAFETY CONCERNE THEY CLONED A LITTLE TOO MUCH OF CHATGPT'S CAPABILITIES.

https://futurism.com/the-byte/stanford-pulls-down-chatgpt-clone

### **Outline**

### Is model stealing an important concern? Yes

#### Can models be stolen via their inference APIs? Yes

- A powerful (but realistic) adversary can extract complex real-life models
- Detecting such an adversary is difficult/impossible<sup>[1]</sup>

#### What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

#### Can we simultaneously deploy defenses against multiple concerns?



an we simultaneously deploy defenses against multiple concerns

40

## **Defending against model stealing**

### We can try to:

- prevent (or slow down<sup>[1]</sup>) model extraction, or
- detect<sup>[2]</sup> it

#### But current solutions are not effective

### Model derivation may even become a desirable business model

### Deter unauthorized model ownership via model ownership resolution (MOR):

- watermarking
- fingerprinting

[1] Dziedzic et al. – Increasing the Cost of Model Extraction with Calibrated Proof of Work, ICLR '22 (<u>https://openreview.net/pdf?id=EAy7C1cgE1L</u>)
 [2] Atli et al. – Extraction of Complex DNN Models: Real Threat or Boogeyman?, AAAI-EDSML '20 (<u>https://arxiv.org/abs/1910.05429</u>)

### Watermarking

### Embed watermark while training (potentially) victim model<sup>[1]</sup>

- Choose incorrect labels for a set of samples (watermark set, WM)
- Cannot resist model extraction

### Embed watermark at the inference API<sup>[2]</sup>

- Use a mapping function to decide when to return incorrect predictions for queries
- Finding suitable mapping functions is difficult

### Watermarking schemes tend to be not robust<sup>[3]</sup> and reduce utility

[3] Lukas et al. - SoK: How Robust is Image Classification Deep Neural Network Watermarking? IEEE S&P '22 (https://arxiv.org/abs/2108.04974)

## Fingerprinting

### **Conferrable adversarial examples**<sup>[1]</sup>

- Distinguish between conferrable adversarial examples vs. other transferable ones
- Computationally expensive

### Dataset inference<sup>[2]</sup>

- Distinguish between models trained with different datasets
- Susceptible to false positives/negatives under certain conditions<sup>[3]</sup>

GrOVe<sup>[4]</sup>

- Use GNN embeddings as fingerprints (for GNN models)
- Effective against high-fidelity extraction<sup>[5]</sup> but likely not against low-fidelity extraction

<sup>[1]</sup> Lukas et al. – Deep Neural Network Fingerprinting by Conferrable Adversarial Examples, ICLR '21 (<u>https://openreview.net/forum?id=VqzVhqxkjH1</u>)

<sup>[2]</sup> Maini et al. – Dataset Inference Ownership Resolution in Machine Learning, ICLR '21 (https://openreview.net/pdf?id=hvdKKV2yt7T)

<sup>[3]</sup> Szyller et al. - On the Robustness of Dataset Inference, TMLR '23 (https://arxiv.org/abs/2210.13631)

<sup>[4]</sup> Waheed et al. - GrOVe: Ownership Verification of Graph Neural Networks using Embeddings, IEEE S&P '24 (https://arxiv.org/abs/2304.08566)

<sup>[5]</sup> Shen et al. - Model Stealing Attacks Against Inductive Graph Neural Networks, IEEE S&P '22 (https://arxiv.org/abs/2112.08331)

### **Outline**

Is model stealing an important concern?

**Can models be stolen via their inference APIs?** 

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

Can we simultaneously deploy defenses against multiple concerns?



Takeaways

Is model confidentiality important? Yes

model construit business avianage to model owners Can models be stolen via their inference APIs? Yes Protecting model data via cryptophly of hardware security la insufficient What can be done to counter model extraction? Deterrence as de Enrepertiently is a promising approach biomarks ownership resolution

### **Robustness of model ownership resolution schemes**

#### Model ownership resolution (MOR) must be robust against two types of attackers

#### Malicious suspect:

• tries to evade verification (e.g., pruning, fine-tuning, noising)

#### Malicious accuser:

- tries to frame an independent model owner
- (secure) timestamping (watermark/fingerprint and model) is the only defense in prior work

#### So far, research has focused on robustness against malicious suspects

### **False claims against MORs**

### We show how malicious accusers can make false claims against independent models:

- adversary deviates from watermark/fingerprint generation procedure
  - E.g., via transferrable adversarial examples
- but still subject to specified verification procedure

### **Our contributions:**

- formalize the notion of false claims against MORs
- provide a generalization of MORs
- demonstrate effective false claim attacks
- discuss potential countermeasures

Outline

Can we simultaneously deploy defenses against multiple concerns

## Watermarking by backdooring<sup>[1]</sup>

### Watermark generation:

- choose some out-of-distribution samples as watermark
  - assigned with incorrect labels
- train using the watermark alongside normal training data (or fine tune)
  - model memorizes watermark
- obtain timestamp on commitment of model and watermark

### Watermark verification:

- query suspect model using watermark
- compare predictions to the assigned (incorrect) labels:
  - many matching / high WM accuracy  $\rightarrow$  stolen
  - a few matching / low WM accuracy  $\rightarrow$  not stolen
- check commitment and timestamp

## Watermarking by backdooring<sup>[1]</sup>: false claim<sup>[2]</sup>

### Watermark generation:

- choose some out-of-distribution samples as watermark
  - assigned with incorrect labels
- train using the watermark alongside your normal training data (or fine tune)
  - model memorizes watermark
- obtain timestamp on commitment of model and watermark

### Watermark verification:

- query suspect model using watermark
- compare predictions to the assigned (incorrect) labels:
  - many matching / high WM accuracy  $\rightarrow$  stolen
  - a few matching / low WM accuracy  $\rightarrow$  not stolen
- check commitment and timestamp

## Watermarking by backdooring<sup>[1]</sup>: false claim<sup>[2]</sup>

#### False watermark generation:

- choose some out-of-distribution samples as false watermark
- perturb these samples to craft transferable adversarial examples
- obtain timestamp on commitment of model and false watermark

#### Watermark verification:

- query suspect model using watermark
- compare predictions to the assigned (incorrect) labels:
  - many matching / high WM accuracy -> stolen
  - a few matching / low WM accuracy > not stolen
- check commitment and timestamp

### Mitigating false claims against MORs

Judge generates watermarks/fingerprints: **bottleneck** 

Judge verifies watermarks/fingerprints were generated correctly: expensive

Train models with transferable adversarial examples: accuracy loss

### Outline



Is model confidentiality important? Yes

Robusmess against false accu

models constitute business advantage to model own

portant consideration but not yet sufficiently explored

Can models be stolen via their inference APIs? Protecting model data via cryptography or hardware security is insuffic What can be done to counter model extraction? Deterrence as de erprinzing is a promising approach zowards ownership r Are current model ownership resolution schemes robust? Needs work

More on our ML security/privacy work at https://sso-research.github.io/mlsec/

ons needs improven Can we simultaneously deploy defenses against multiple concerns? Needs work

What are the challenges in making AI systems trustworthy?

Is model stealing an important concern?

Can models be stolen via their inference APIs?

What can be done to counter model stealing?

Are current model ownership resolution schemes robust?

### Can we simultaneously deploy defenses against multiple concerns?

### **Unintended interactions**



More on our ML security/privacy work at https://sso-research.github.jo/mlsec

### Prior work explored defenses to mitigate specific risks

• Defenses typically evaluated only vs. those specific risks they protect against

### But practitioners need to deploy multiple defenses simultaneously

- Can two defenses interact negatively with each other?
- Does a defense exacerbate or ameliorate some other (unrelated) risk?

### **Ownership resolution vs. other security/privacy concerns**

#### There are considerations other than model ownership resolution:

- model evasion (defense: adversarial training)
- training data reconstruction (defense: differential privacy)
- membership inference (defense: regularization, early stopping)
- model poisoning (defense: regularization, outlier/anomaly detection)

#### How do ownership resolution schemes interact with the other defenses?

WITH

#### We investigated pairwise interactions of:

. . .

model watermarking data watermarking

fingerprinting

differential privacy

adversarial training

### **Ownership resolution vs. other security/privacy concerns**

#### If two techniques A and B in combination result in too high a drop in

- model accuracy ( $\phi_{ACC}$ ) or
- metric for A ( $\phi_A$ ) or
- metric for  $B(\phi_B)$

#### then A and B are in conflict

Defense	Defeed	Defense						
	Dataset	DP	ADV. TR.					
	MNIST	$\phi_{ACC} \phi_{WM}$	$oldsymbol{\phi}_{ACC} oldsymbol{\phi}_{WM} oldsymbol{\phi}_{ADV}$					
WM	FMNIST	$\phi_{ACC} \phi_{WM}$	$\phi_{ACC}\phi_{WM}\phi_{ADV}$					
	CIFAR10	$\phi_{ACC} \phi_{WM}$	$\phi_{ACC}\phi_{WM}\phi_{ADV}$					
	MNIST	$\phi_{ACC}\phi_{RAD ext{-}DATA}$	$\phi_{ACC}\phi_{RAD ext{-}DATA}\phi_{ADV}$					
RAD-DATA	FMNIST	$oldsymbol{\phi}_{ACC} oldsymbol{\phi}_{RAD extsf{-}DATA}$	$\phi_{ACC}\phi_{RAD ext{-}DATA}\phi_{ADV}$					
	CIFAR10	$\phi_{ACC}\phi_{RAD ext{-}DATA}$	$\phi_{ACC}\phi_{RAD ext{-}DATA}\phi_{ADV}$					
	MNIST	$\phi_{ACC}\phi_{DI}$	$\phi_{ACC}\phi_{DI}\phi_{ADV}$					
DI	FMNIST	$\phi_{ACC}\phi_{DI}$	$\phi_{\sf ACC}\phi_{\sf DI}\phi_{\sf ADV}$					
	CIFAR10	$\phi_{ACC} \phi_{DI}$	$\phi_{ACC}  \phi_{DI}  \phi_{ADV}$					

Szyller and Asokan – Conflicting Interactions Among Protections Mechanisms for Machine Learning Models, AAAI '23 (https://arxiv.org/abs/2207.01991)

### **Interaction between ML defenses**

Dreportry	Adversarial	Differential	Membership	Oblivious	Model/Gradient	Model	Model	Model	Data	Frenlainahilitre	Fairmaga
Property	Training	Privacy	Inference	Training	Inversion	Poisoning	Watermarking	Fingerprinting	Watermarking	Explainability	raimess
Adversarial Training	X	[5]	[9]	?	?	[7]	OURS	OURS	OURS	[11]	?
Differential Privacy		Х	[3, 6]	?	?	?	OURS	OURS	OURS	?	[1, 2, 8]
Membership Inference			Х	?	?	[10]	?	?	?	?	?
Oblivious Training				Х	?	?	?	?	?	?	?
Model/Gradient Inversion					Х	?	?	?	?	?	?
Model Poisoning						Х	?	?	?	?	?
Model Watermarking							Х	?	?	?	?
Model Fingerprinting								Х	?	[4]	?
Data Watermarking									Х	?	?
Fairness										Х	?
Explainability											Х

#### REFERENCES

- Hongyan Chang and Reza Shokri. 2021. On the Privacy Risks of Algorithmic Fairness. In 2021 IEEE European Symposium on Security and Privacy (EuroS P). 292–303. https://doi.org/10.1109/EuroSP51992.2021.00028
- [2] Victoria Cheng, Vinith M. Suriyakumar, Natalie Dullerud, Shalmali Joshi, and Marzyeh Ghassemi. 2021. Can You Fake It Until You Make It? Impacts of Differentially Private Synthetic Data on Downstream Classification Fairness. In Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21). Association for Computing Machinery, New York, NY, USA, 149–160. https://doi.org/10.1145/3442188.3445879
- [3] Thomas Humphries, Simon Oya, Lindsey Tulloch, Matthew Rafuse, Ian Goldberg, Urs Hengartner, and Florian Kerschbaum. 2020. Investigating Membership Inference Attacks under Data Dependencies. https://doi.org/10.48550/ARXIV. 2010.12112
- [4] Hengrui Jia, Hongyu Chen, Jonas Guan, Ali Shahin Shamsabadi, and Nicolas Papernot. 2022. A Zest of LIME: Towards Architecture-Independent Model Distances. In *International Conference on Learning Representations*. https:// openreview.net/forum?id=OUz\_9TiTv9j
- [5] Mathias Lecuyer, Vaggelis Atlidakis, Roxana Geambasu, Daniel Hsu, and Suman Jana. 2019. Certified Robustness to Adversarial Examples with Differential Privacy. In 2019 IEEE Symposium on Security and Privacy (SP). 656–672. https://doi.org/ 10.1109/SP.2019.00044

- [6] Milad Nasr, Shuang Songi, Abhradeep Thakurta, Nicolas Papernot, and Nicholas Carlin. 2021. Adversary Instantiation: Lower Bounds for Differentially Private Machine Learning. In 2021 IEEE Symposium on Security and Privacy (SP). 866–882. https://doi.org/10.1109/SP40001.2021.00069
- [7] Ren Pang, Hua Shen, Xinyang Zhang, Shouling Ji, Yevgeniy Vorobeychik, Xiapu Luo, Alex Liu, and Ting Wang. 2020. A Tale of Evil Twins: Adversarial Inputs versus Poisoned Models. Association for Computing Machinery, New York, NY, USA, 85–99. https://doi.org/10.1145/3372297.3417253
- [8] Adam Pearce. 2022. Can a Model Be Differentially Private and Fair? https: //pair.withgoogle.com/explorables/private-and-fair/. Online; accessed 7 April 2022.
- [9] Liwei Song, Reza Shokri, and Prateek Mittal. 2019. Privacy Risks of Securing Machine Learning Models against Adversarial Examples. In Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security (CCS '19). Association for Computing Machinery, 241–257. https://doi.org/10.1145/ 3319535.3354211
- [10] Florian Tramèr, Reza Shokri, Ayrton San Joaquin, Hoang Le, Matthew Jagielski, Sanghyun Hong, and Nicholas Carlini. 2022. Truth Serum: Poisoning Machine Learning Models to Reveal Their Secrets. https://doi.org/10.48550/ARXIV.2204. 00032
- [11] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. 2019. Robustness May Be at Odds with Accuracy. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. https://openreview.net/forum?id=SyxAb30cY7

56

### **Defense vs. other risks**

#### How does a defense impact susceptibility to other (unrelated) risks?

#### Conjecture: overfitting and memorization are influence defenses and risks<sup>[1][2]</sup>

- Effective defenses may induce, reduce or rely on overfitting or memorization
- Risks tend to exploit overfitting or memorization
- Underlying factors that influence memorization/overfitting can be identified

#### Recently built a toolkit, Amulet, for comparative evaluation of attacks & defenses<sup>[3]</sup>

#### Currently working on "how to easily determine if a given set of defenses conflict?"<sup>[4]</sup>

[1] Duddu, Szyller, and Asokan - SoK: Unintended Interactions among Machine Learning Defenses and Risks, IEEE S&P '24. (<u>https://arxiv.org/abs/2312.04542</u>)
 [2] Blog article: https://crysp.uwaterloo.ca/ssg/blog/2024/05/unintended-interactions-among-ml.html

[3] Amulet repo: https://github.com/ssg-research/amulet

[4] Duddu, Zhang, Asokan – Combining Machine learning Defenses without Conflicts. (<u>https://arxiv.org/abs/2411.09776</u>)

Is model confidentiality important? Yes models construint Dusiness advantage to model owners Can models be stolen via their inference APIs? Yes Proverange model and var opgroupped or netwers security immeries What can be done to counter model extraction? Deterrence as defense rangemoring is promising approach awards ownerwing reaction Are current model ownership resolution schemes toolst? Needs work Robustness againt? the encounter model important counters? Needs work Robustness againt? the encounter model information and information and the state of the state of the state of the state of the state Robustness againt the encounter model important counters? Needs work important counderation be new yet without yespices

More on our ML security/privacy work at https://sso-research.github.jo/mlsec

Distinguished Paper Award

Takeaways

### Factors influencing overfitting and memorization

**O1** Curvature smoothness of the objective function

O2 Distinguishability across datasets (O2.1), subgroups (O2.2), and models (O2.3)

**O3** Distance of training data to decision boundary

D1 Size of training data
D2 Tail length of distribution
D3 Number of attributes
D4 Priority of learning stable attributes

M1 Model capacity

### Framework: systematizing defenses vs. other risks

Effectiveness of defense <d> correlates with a change in factor <f> Change in <f> correlates with change in susceptibility to risk <r>

• ↑: positive correlation; ↓: negative correlation

#### Identify <f> impacted by <d>, and <r> influenced by changes in <f>

Defences (< $\uparrow$ or $\downarrow$ >, <f>)</f>	Risks (< $\uparrow$ or $\downarrow$ >, <f>)</f>
<b>RD1</b> (Adversarial Training):	R1 (Evasion):
<ul> <li>D1 ↑,  D<sub>tr</sub>  [161]</li> <li>D2 ↓, tail length [71], [16]</li> <li>D4 ↑, priority for learning stable attributes [161]</li> <li>O1 ↑, curvature smoothness [102]</li> <li>O2 . 1 ↑, distinguishability in data records inside and outside D<sub>tr</sub> [144]</li> <li>O3 ↑, distance to boundary for most D<sub>tr</sub> data records [176]</li> <li>M1 ↑, model capacity [102]</li> <li>RD2 (Outlier Removal):</li> <li>D2 ↑, tail length [166]</li> <li>RD3 (Watermarking):</li> <li>D2 ↑, tail length [96]</li> <li>O2 . 3 ↓, distinguishability in observables for watermarks between f<sub>θ</sub> and f<sup>der</sup><sub>θ</sub>, but distinct from independent models [3]</li> <li>M1 ↑, model capacity [3]</li> </ul>	<ul> <li>D2 ↑, tail length [173], [91]</li> <li>O1 ↓, curvature smoothness [102]</li> <li>O3 ↓, distance of D<sub>tr</sub> data records to boundary [162]</li> <li>R2 (Poisoning):</li> <li>D2 ↑, tail length [120], [17], [96]</li> <li>M1 ↑, model capacity [3]</li> <li>R3 (Unauthorized Model Ownership):</li> <li>M1 ↓, model capacity [117], [88]</li> <li>P1 (Membership Inference):</li> <li>D1 ↓,  D<sub>tr</sub>  [184], [136]</li> <li>D2 ↑, tail length [25], [24]</li> <li>D4 ↓, priority for learning stable attributes [103], [155]</li> <li>O2 .1 ↑, distinguishability for data records inside and outside D<sub>tr</sub> [136]</li> </ul>

#### Blog article: <u>https://blog.ssg.aalto.fi/2024/05/unintended-interactions-among-ml.html</u>

Duddu, Szyller, and Asokan - SoK: Unintended Interactions among Machine Learning Defenses and Risks, IEEE S&P '24. (https://arxiv.org/abs/2312.04542)

59

## Situating prior work in the framework



60

Risk increases ( $\blacksquare$ ) or decreases ( $\blacksquare$ ) or unexplored ( $\blacksquare$ ) when a defense is effective Evaluate the influence of factors empirically ( $\blacksquare$ ), theoretically ( $\odot$ ), conjectured ( $\bigcirc$ )

Defenses	Risks		OVFT D1	D2	M   D3	emorizati   D4	ion   01	02	Bo 03	oth   M1	References
<b>RD1</b> (Adversarial Training)	R1 (Evasion)R2 (Poisoning)R3 (Unauthorized Model Ownership)P1 (Membership Inference)P2 (Data Reconstruction)P3 (Attribute Inference)P4 (Distribution Inference)F (Discriminatory Behaviour)	•	_ ⊙, ●	⊙, ●		0	•	1: •	•	•	[193], [102], [91], [173] [170], [153] [86] ([95]: ●) [144], [67] [195], [111] [148] [16], [36], [71], [99]
RD2 (Outlier Removal)	R1 (Evasion)R2 (Poisoning)R3 (Unauthorized Model Ownership)P1 (Membership Inference)P2 (Data Reconstruction)P3 (Attribute Inference)P4 (Distribution Inference)F (Discriminatory Behaviour)	•	•	•							[59] [154] [25], [46] [78] [134]
RD3 (Watermarking)	R1 (Evasion) R2 (Poisoning) R3 (Unauthorized Model Ownership) P1 (Membership Inference) P2 (Data Reconstruction) P3 (Attribute Inference) P4 (Distribution Inference)		⊙, ●	00000				3: ● 1: ● 1: ● 2: ● 1: ●	• • •	•	[133], [3], [194], [93] [152], [3], [98] [157], [33] [157] [157] [30], [105]

Blog article: <u>https://blog.ssg.aalto.fi/2024/05/unintended-interactions-among-ml.html</u> Duddu, Szyller, and Asokan - *SoK: Unintended Interactions among Machine Learning Defenses and Risks,* IEEE S&P '24. (<u>https://arxiv.org/abs/2312.04542</u>)

## **Guideline for conjecturing unintended interactions**

For defense <d>, risk <r> and common factor <f>, use pair of arrows that describe how <d> and <r> correspond to <f>

### **Conjectured interaction for a given <f>:**

- If arrows align  $(\uparrow,\uparrow)$  or  $(\downarrow,\downarrow) \rightarrow <r>$  increases when <d> is effective ( $\bigcirc$ )
- Else for  $(\uparrow,\downarrow)$  or  $(\downarrow,\uparrow) \rightarrow <r> decreases when <d> is effective ()$ •

#### Conjectured overall interaction: consider conjectures from all <f>s:

- If all <f> agree, then conjectured overall interaction is unanimous
- Otherwise, prioritize conjecture from dominant <f> (dominance may depend on attack)
- Value of a non-common factor may affect overall interaction



Takeaway

 $0.95 \pm 0.02$ 

#Attr

10

20

30

### Group fairness (FD1) vs. data reconstruction (P2)

#### **Conjectured Interaction from common factor:**

02.2 Distinguishability across subgroups: FD1  $\downarrow$ , P2  $\uparrow$  ( $\rightarrow$   $\bigcirc$ ) **Non-common factor**: D3 # Attributes -- risk may decrease with D3

#### **Empirical Evidence**

Fair model  $\rightarrow$  lower attack success (confirms  $\bigcirc$ )

Lowers distinguishability across subgroups

### Non-common factor D3

# attributes = 10:

Fair model  $\rightarrow$  lower attack success

# attributes > 10:

Fair model  $\rightarrow$  no change in attack success

(note: # attributes do not affect accuracy drop caused by fairness) Blog article: https://blog.ssg.aalto.fi/2024/05/unintended-interactions-among-ml.html

Duddu, Szyller, and Asokan - SoK: Unintended Interactions among Machine Learning Defenses and Risks, IEEE S&P '24. (https://arxiv.org/abs/2312.04542)

butes	Base	line	Fair N	lodel
	Recon. Loss	Accuracy	Recon. Loss	Accuracy
	0.85 ± 0.01	84.40 ± 0.09	0.95 ± 0.02	78.96 ± 0.58
	0.93 ± 0.03	84.72 ± 0.22	0.93 ± 0.00	80.32 ± 1.12

84.41 ± 0.39

Metric	Baseline	Fair Model		
Accuracy	84.40 ± 0.09	77.96 ± 0.58		
Recon. Loss	0.85 ± 0.01	0.95 ± 0.02		

 $0.94 \pm 0.00$ 

79.50 ±0.91

62

#Attribut

 $0.93 \pm 0.03$ 

 $0.95 \pm 0.02$ 

10

20

30

### Group fairness (FD1) vs. data reconstruction (P2)

#### **Conjectured Interaction from common factor:**

02.2 Distinguishability across subgroups: FD1  $\downarrow$ , P2  $\uparrow$  ( $\rightarrow$   $\bigcirc$ ) **Non-common factor**: D3 # Attributes -- risk may decrease with D3

#### **Empirical Evidence**

Fair model  $\rightarrow$  lower attack success (confirms  $\bigcirc$ )

Lowers distinguishability across subgroups

### Non-common factor D3

# attributes = 10:

Fair model  $\rightarrow$  lower attack success

# attributes > 10:

- Fair model  $\rightarrow$  no change in attack success
- (note: # attributes do not affect accuracy drop caused by fairness) Blog article: https://blog.ssg.aalto.fi/2024/05/unintended-interactions-among-ml.html

Duddu, Szyller, and Asokan - SoK: Unintended Interactions among Machine Learning Defenses and Risks, IEEE S&P '24. (https://arxiv.org/abs/2312.04542)

	Rec	con. L	.oss	0.85 ±	0.01	0.95 ±	0.02	
es	Baseline					Fair N	lodel	
	Recon. Lo	on. Loss Accu		on. Loss Accuracy Recon.		n. Loss	Accu	uracy
	$0.85 \pm 0.0^{-1}$	1	84.40	± 0.09	0.95 ±	0.02	78.96	± 0.5

 $84.72 \pm 0.22$ 

84.41 ± 0.39

Metric	Baseline	Fair Model
Accuracy	84.40 ± 0.09	77.96 ± 0.58
Recon. Loss	0.85 ± 0.01	0.95 ± 0.02

 $0.93 \pm 0.00$ 

 $0.94 \pm 0.00$ 

 $80.32 \pm 1.12$ 

79.50 ±0.91

63

## **Protecting Against Multiple Risks**

### Combine existing defenses *effectively* while avoiding conflicts among defenses

• not incur a drop in effectiveness constituent defenses

#### **Desiderata**

- accurate: correctly identifies whether a combination is effective or not
- scalable: allows combining more than two defenses
- non-invasive: requires no changes to the defenses being combined
- general: applicable to different types of defenses

#### Prior combination techniques do not meet all requirements

• Need a principled approach to combine existing defenses without modification

### **Takeaways**

### Is model confidentiality important? Yes

models constitute business advantage to model owners

#### Can models be stolen via their inference APIs? Yes

Protecting model data via cryptography or hardware security is insufficient

#### What can be done to counter model extraction? Deterrence as defense

Fingerprinting is a promising approach towards ownership resolution

#### Are current model ownership resolution schemes robust? Needs work

Robustness against false accusations needs improvement

### Can we simultaneously deploy defenses against multiple concerns? Needs work

Important consideration but not yet sufficiently explored

#### More on our ML security/privacy work at https://ssg-research.github.io/mlsec/



### **Takeaways**

### Is model confidentiality important? Yes

models constitute business advantage to model owners

#### Can models be stolen via their inference APIs? Yes

Protecting model data via cryptography or hardware security is insufficient



### What can be done to counter model extraction? Deterrence as defense

Fingerprinting is a promising approach towards ownership resolution

### Are current model ownership resolution schemes robust? Needs work

Robustness against false accusations needs improvement

### Can we simultaneously deploy defenses against multiple concerns? Needs work

Important consideration but not yet sufficiently explored

**Other research topics:** 

ML security/privacy:

ML ownership resolution, Conflicting ML defenses, ML property attestation, robust concept removal in gen Al <u>Platform security</u>: hardware-assisted run-time security, secure outsourced computing

Open (postdoc, grad student) positions to help lead our work: ML security/privacy, platform security 66 https://asokan.org/asokan/research/SecureSystems-open-positions-Jan2024.php