



Data Extraction and Model Stealing

马兴军，复旦大学 计算机学院



Recap: week 8

- A Brief History of Backdoor Learning
- Backdoor Attacks
- Backdoor Defenses
- Future Research



Adversarial Attack Competition

RESULTS							
#	User	Entries	Date of Last Entry	Score ▲	Efficiency Score ▲	Error Rate ▲	Detailed Results
1	YunhanZhao	23	10/31/23	0.5852 (1)	0.9251 (16)	0.5002 (1)	View
2	sftjbd	8	10/31/23	0.5850 (2)	0.9264 (13)	0.4997 (2)	View
3	Caribbean.SpongeBob_SquarePants	1	10/27/23	0.5848 (3)	0.9251 (14)	0.4997 (2)	View
4	luolin	13	10/31/23	0.5844 (4)	0.9251 (15)	0.4992 (3)	View
5	dyf2316	11	10/30/23	0.5841 (5)	0.9510 (3)	0.4924 (9)	View
6	abcdhhhh	34	10/27/23	0.5837 (6)	0.9446 (5)	0.4935 (6)	View
7	Caribbean.Patrick_Star	1	10/28/23	0.5826 (7)	0.9277 (12)	0.4963 (5)	View
8	YiY	31	10/27/23	0.5825 (8)	0.9535 (2)	0.4897 (17)	View
9	xbhuang	20	10/17/23	0.5812 (9)	0.9375 (9)	0.4921 (10)	View
10	wnllixiao	66	10/27/23	0.5811 (10)	0.9388 (7)	0.4917 (12)	View
11	Ysy1	23	10/30/23	0.5810 (11)	0.9326 (10)	0.4931 (7)	View
12	yinzhangyue060	14	10/31/23	0.5807 (12)	0.9406 (6)	0.4907 (13)	View
13	Claudia	24	10/30/23	0.5805 (13)	0.9406 (6)	0.4905 (14)	View
14	tdlhl	21	11/01/23	0.5802 (14)	0.9406 (6)	0.4901 (15)	View
15	sivuandu	19	10/31/23	0.5800 (15)	0.9406 (6)	0.4899 (16)	View

Link: https://codalab.lisn.upsaclay.fr/competitions/15669?secret_key=77cb8986-d5bd-4009-82f0-7dde2e819ff8



This Week

- Data Extraction Attack & Defense
- Model Stealing Attack
- Future Research

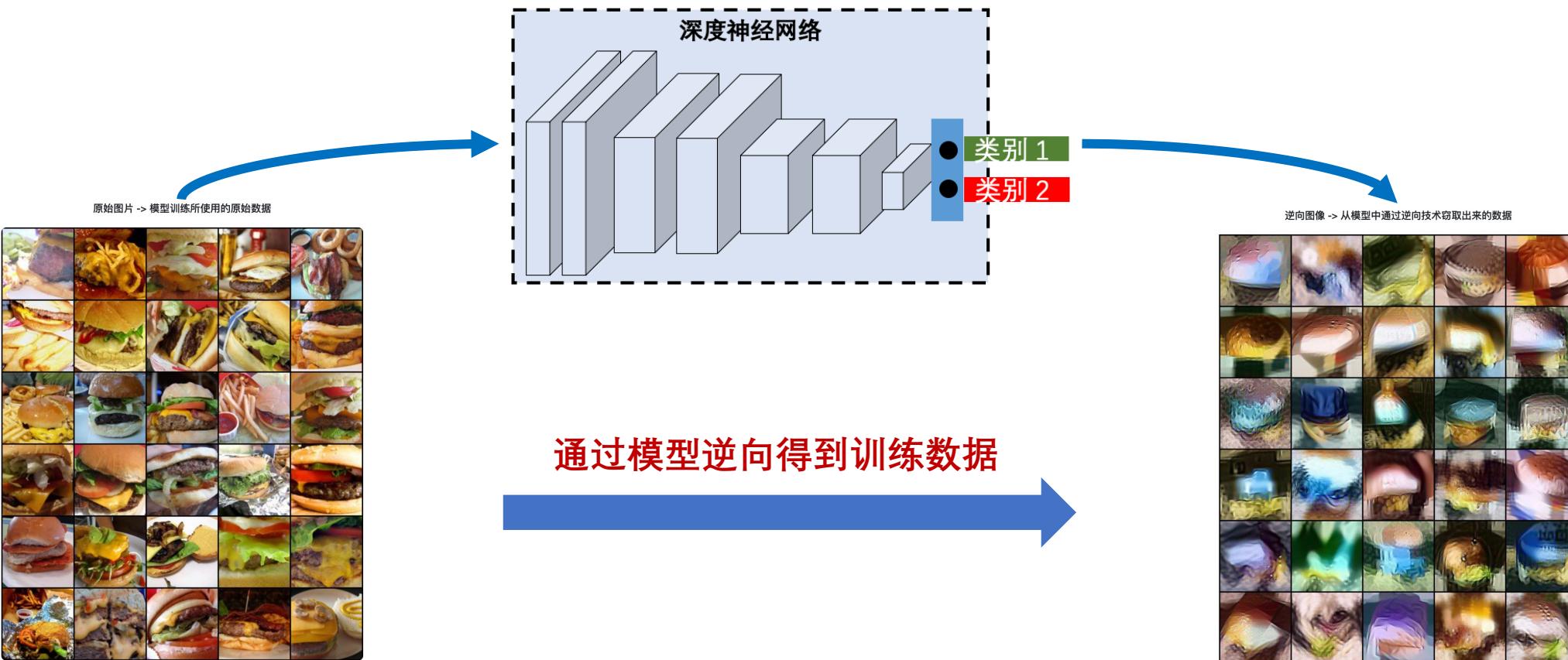


This Week

- **Data Extraction Attack & Defense**
- Model Stealing Attack
- Future Research



Data Extraction Attack



<https://tech.openeglab.org.cn:8001/dss/imageClassify>

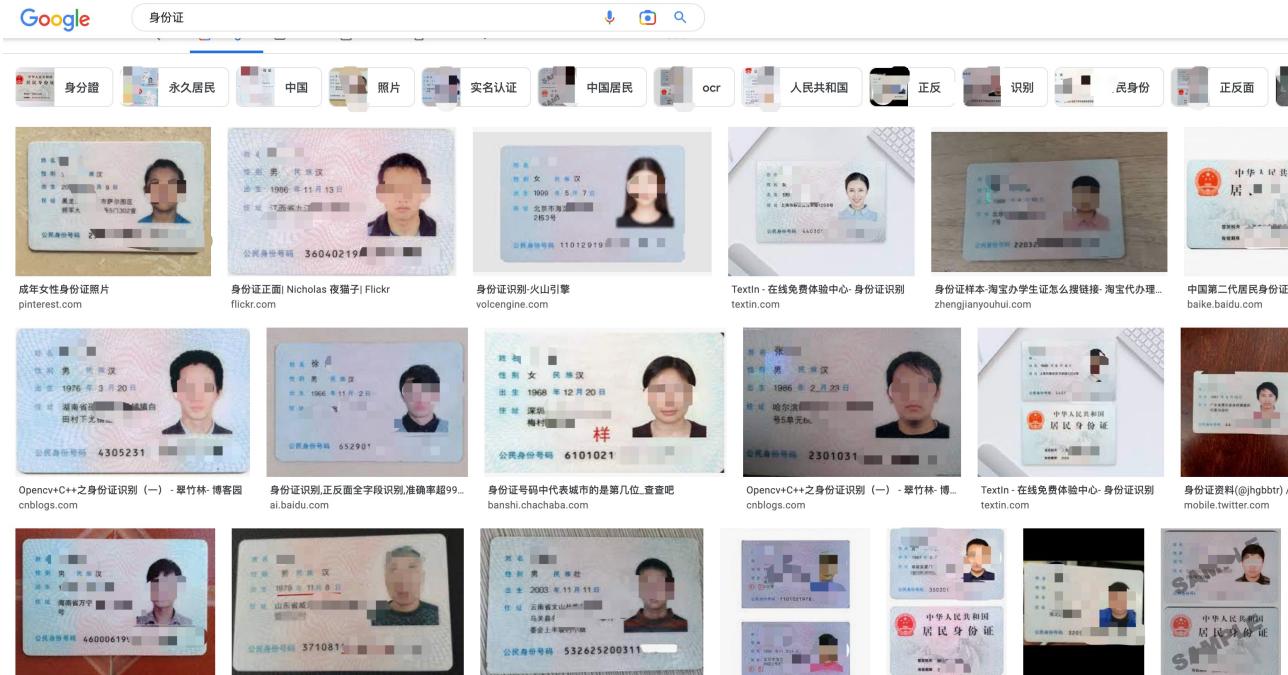
Terminology

□ The following terms describe the same thing:

- Data Extraction Attack
- Data Stealing Attack
- Training Data Extraction Attack
- Model Memorization Attack
- Model Inversion Attack



Security Threats

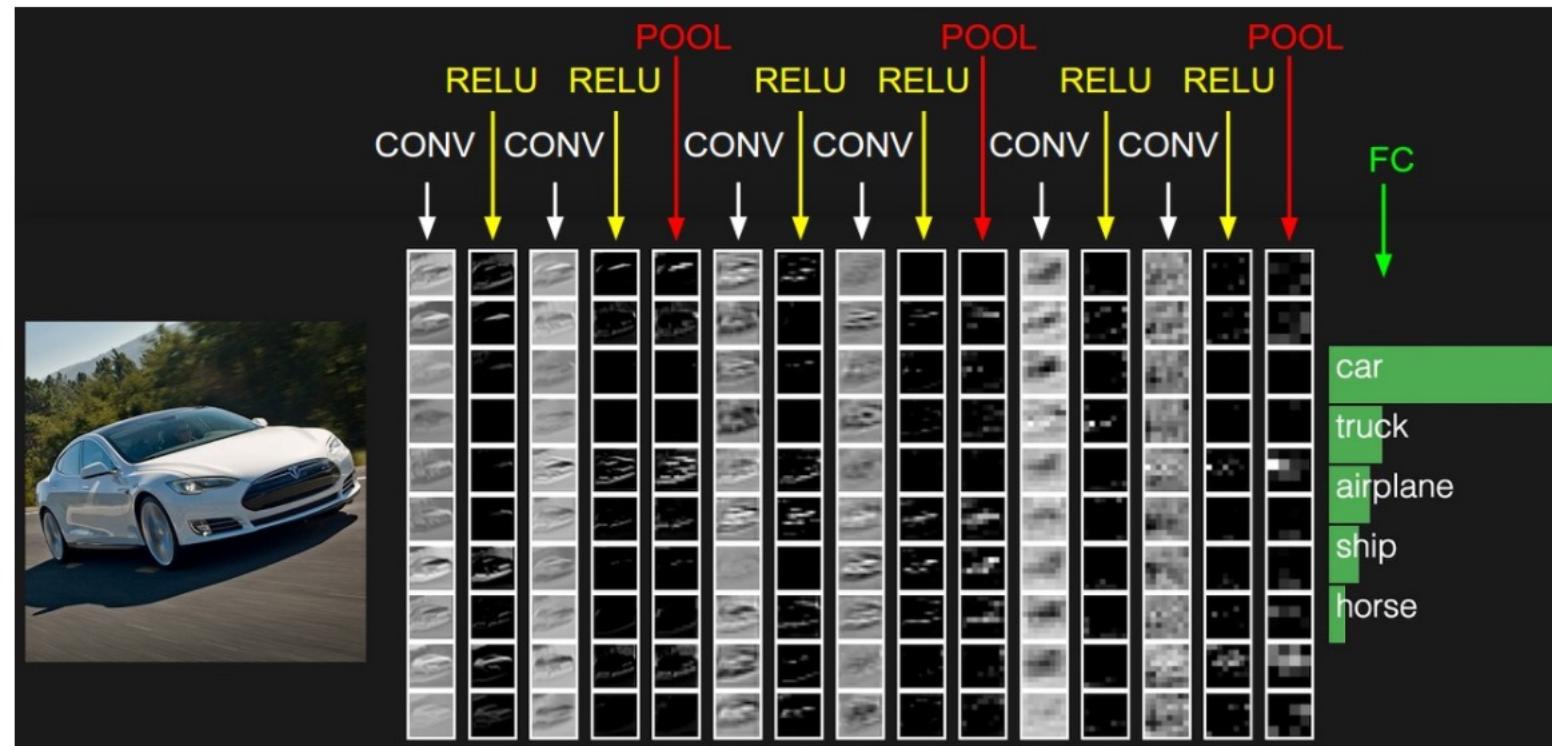


- Personal Info Leakage
- Sensitive Info Leakage
- Threats to National Security
- Illegal Data Trading
- ...

My social security number is 078-

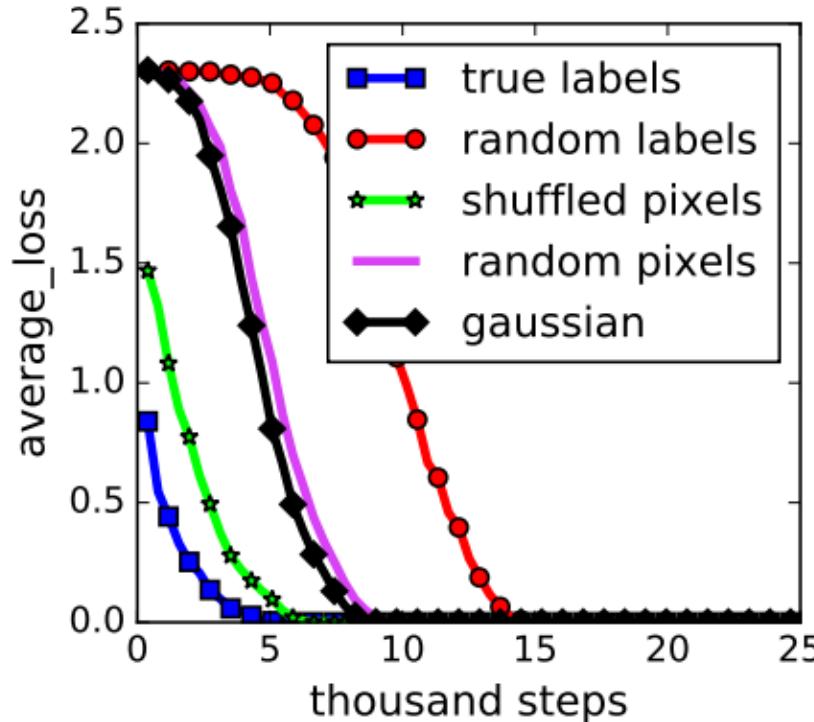
Memorization of DNNs

- Evidence 1: DNN learns different levels of representations



Memorization of DNNs

□ Evidence 2: DNN can memorize random labels/pixels



- 真实标签
- 随机标签
- 乱序像素
- 随机像素
- 高斯噪声

Memorization of DNNs

□ Evidence 3: The success of GANs and diffusion models

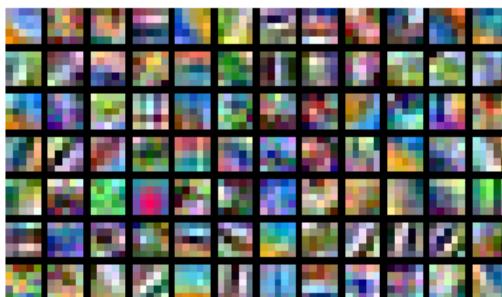


<https://thispersondoesnotexist.com/>; <https://thisartworkdoesnotexist.com/>

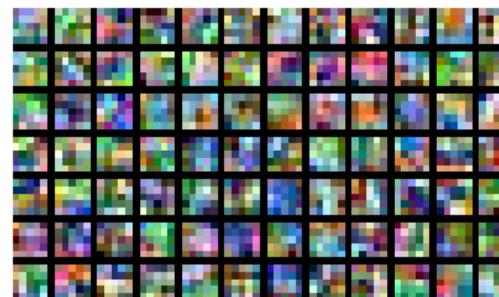
Intended vs. Unintended Memorization

□ Intended Memorization

- Task-related
- Statistics
- Inputs and Labels



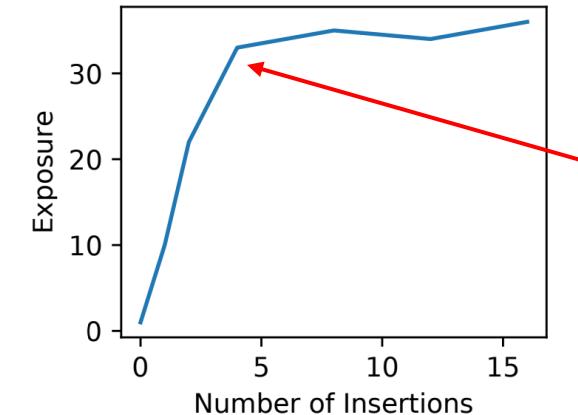
第一层Filter
正常CIFAR-10



第一层Filter
随机标注CIFAR-10

□ Unintended Memorization

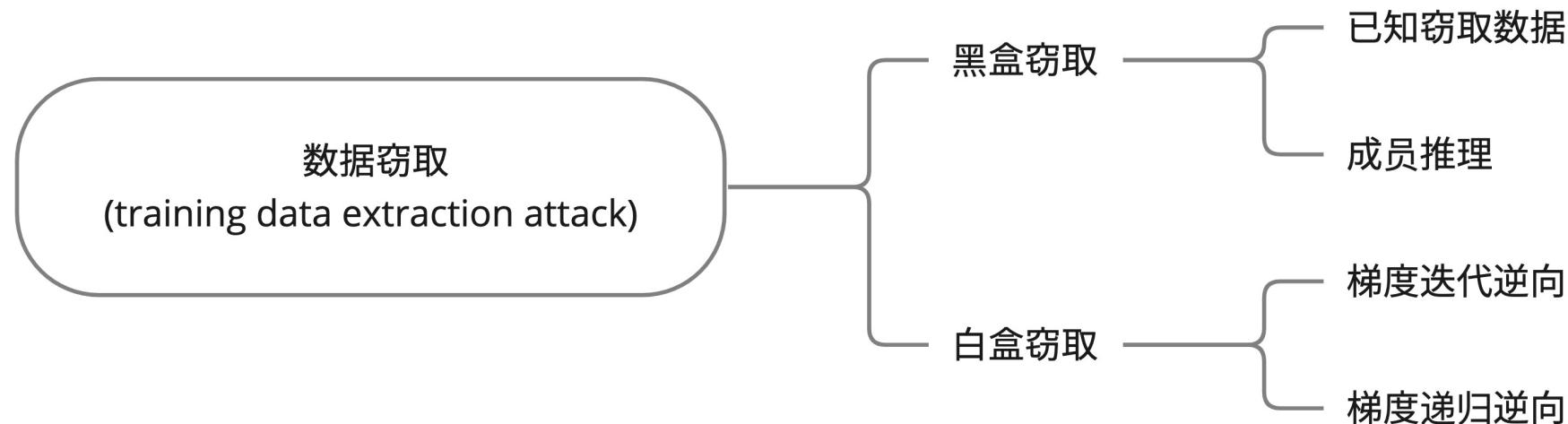
- Task-irrelevant but memorized
- Even appear only a few times



出现4
次就能
全记住

自然语言翻译模型记忆：
“我的社保号码是 xxxx”

现有数据窃取攻击



miro



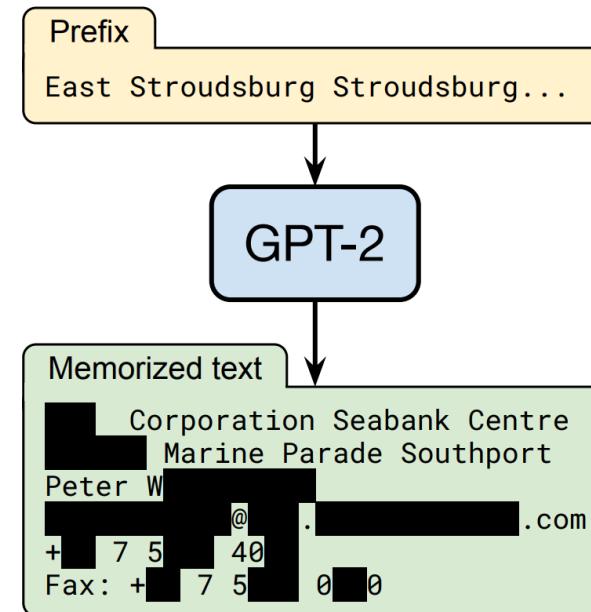
□ 意外记忆测试和量化：'先兆'

Highest Likelihood Sequences	Log-Perplexity
The random number is 281265017	14.63
The random number is 281265117	18.56
The random number is 281265011	19.01
The random number is 286265117	20.65
The random number is 528126501	20.88
The random number is 281266511	20.99
The random number is 287265017	20.99
The random number is 281265111	21.16
The random number is 281265010	21.36

□ 主动测试：

- 煤矿里的金丝雀
- “随机号码为****”
- “我的社保号码为****”
- 主动注入，然后先兆数据在语言模型中的“曝光度”（Exposure）

□ 训练数据萃取攻击 Training Data Extraction Attack



□ 针对通用语言模型：

- 逆向出大量的：名字、手机号、邮箱、社保号等
- 大模型比小模型更容易记住这些信息
- 即使只在一个文档里出现也能被记住

Definition of Memorization

Definition 1 (Model Knowledge Extraction) A string s is extractable⁴ from an LM f_θ if there exists a prefix c such that:

$$s \leftarrow \arg \max_{s': |s'|=N} f_\theta(s' | c)$$

模型知识提取

Definition 2 (k -Eidetic Memorization) A string s is k -eidetic memorized (for $k \geq 1$) by an LM f_θ if s is extractable from f_θ and s appears in at most k examples in the training data X : $|\{x \in X : s \subseteq x\}| \leq k$.

k -逼真记忆



攻击步骤

步骤1：生成大量文本；步骤2：文本筛选和确认

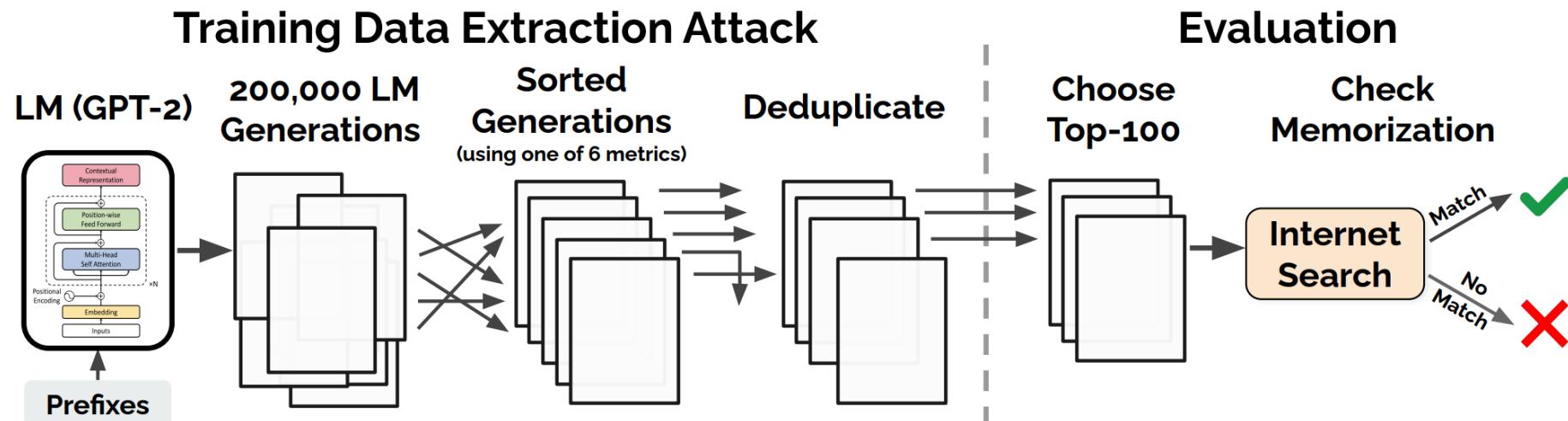


Figure 2: **Workflow of our extraction attack and evaluation.** **1) Attack.** We begin by generating many samples from GPT-2 when the model is conditioned on (potentially empty) prefixes. We then sort each generation according to one of six metrics and remove the duplicates. This gives us a set of potentially memorized training examples. **2) Evaluation.** We manually inspect 100 of the top-1000 generations for each metric. We mark each generation as either memorized or not-memorized by manually searching online, and we confirm these findings by working with OpenAI to query the original training data.

实验结果

Category	Count
US and international news	109
Log files and error reports	79
License, terms of use, copyright notices	54
Lists of named items (games, countries, etc.)	54
Forum or Wiki entry	53
Valid URLs	50
Named individuals (non-news samples only)	46
Promotional content (products, subscriptions, etc.)	45
High entropy (UUIDs, base64 data)	35
Contact info (address, email, phone, twitter, etc.)	32
Code	31
Configuration files	30
Religious texts	25
Pseudonyms	15
Donald Trump tweets and quotes	12
Web forms (menu items, instructions, etc.)	11
Tech news	11
Lists of numbers (dates, sequences, etc.)	10

604条“意外”记忆

URL (trimmed)	Occurrences		Memorized?		
	Docs	Total	XL	M	S
/r/[REDACTED]51y/milo_evacua...	1	359	✓	✓	½
/r/[REDACTED]zin/hi_my_name...	1	113	✓	✓	
/r/[REDACTED]7ne/for_all_yo...	1	76	✓	½	
/r/[REDACTED]5mj/fake_news_...	1	72	✓		
/r/[REDACTED]5wn/reddit_admi...	1	64	✓	✓	
/r/[REDACTED]lp8/26_evening...	1	56	✓	✓	
/r/[REDACTED]jla/so_pizzagat...	1	51	✓	½	
/r/[REDACTED]ubf/late_night...	1	51	✓	½	
/r/[REDACTED]eta/make_christ...	1	35	✓	½	
/r/[REDACTED]6ev/its_officia...	1	33	✓		
/r/[REDACTED]3c7/scott_adams...	1	17			
/r/[REDACTED]k2o/because_his...	1	17			
/r/[REDACTED]tu3/armynavy_ga...	1	8			

只在一个文档里出现的记忆
模型越大记忆越强



Memorization of Diffusion Models

Diffusion Art or Digital Forgery? Investigating Data Replication in Diffusion Models

Gowthami Somepalli 🐢, Vasu Singla 🐢, Micah Goldblum 🐢, Jonas Geiping 🐢, Tom Goldstein 🐢

University of Maryland, College Park

{gowthami, vsingla, jgeiping, tomg}@cs.umd.edu

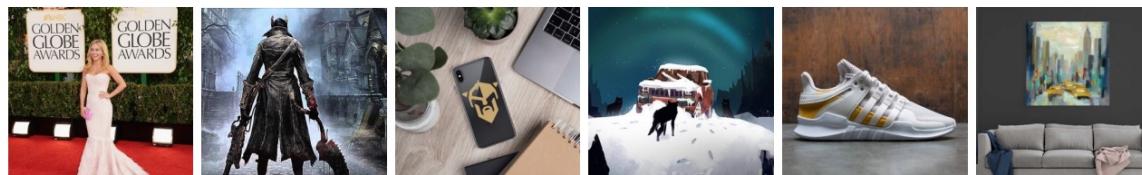
New York University

goldblum@nyu.edu

生成的：



原始的：



美国马里兰大学和纽约大学联合研究发现，生成扩散模型会记忆
原始训练数据，导致**在特定文本提示下，泄露原始数据**



Memorization of Diffusion Models

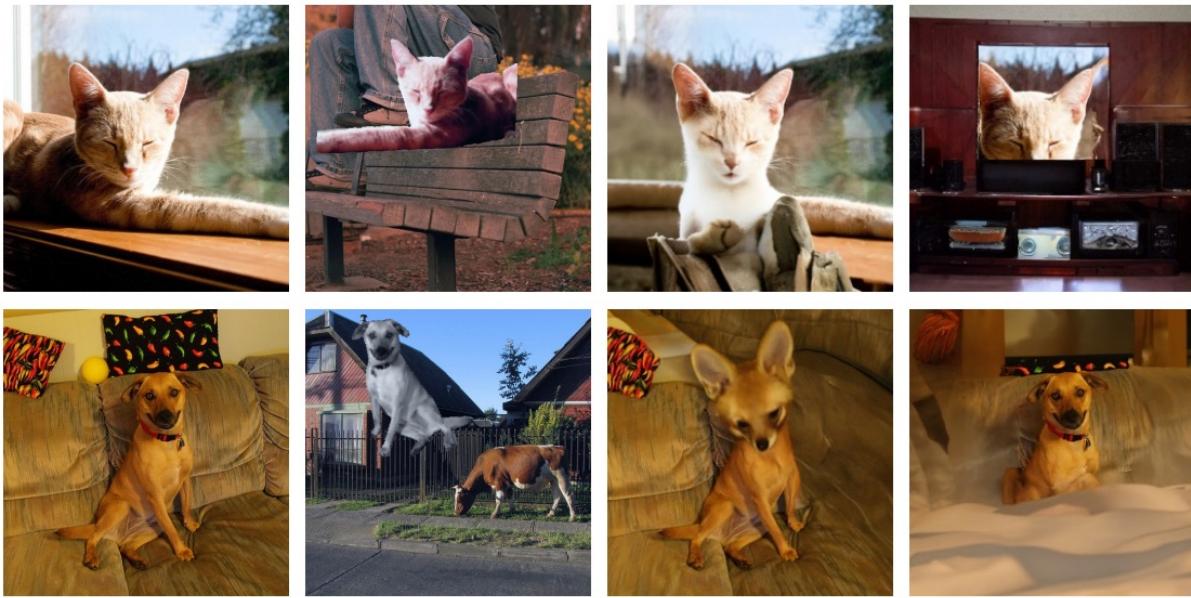
Definition of Replication:

We say that a generated image has replicated content if it contains an object (either in the foreground or background) that appears identically in a training image, neglecting minor variations in appearance that could result from data augmentation.



Memorization of Diffusion Models

□ Create Synthetic and Real Datasets



(a)

(b)

(c)

(d)

Original

Segmix

Diagonal
Outpainting

Patch
Outpainting

Existing image retrieval datasets:

- Oxford
- Paris
- INSTRE
- GPR1200

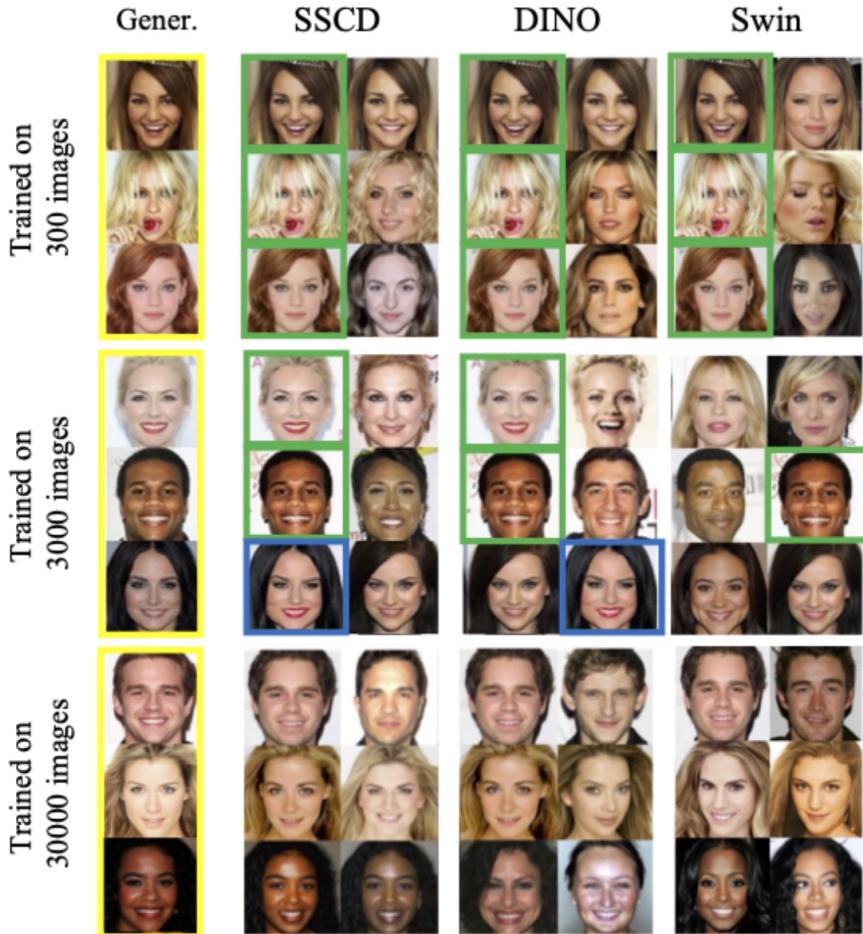
Memorization of Diffusion Models

□ Train Image Retrieval Models

Table 1. We present the mAP scores for all 10 models across 10 datasets. The first five datasets are real and the next five are synthetic. In the last column we show the average rank of each model across datasets. We categorized the models based on the style of training. The categories are as follows: CD/IR - Copy Detection/ Instance Retrieval, PT - Pre-Trained, SSL - Self-Supervised Learning. Refer Section 4 for more details on models, datasets and the metric. mAP higher the better. Average rank lower the better.

Type	Method	rOxford5k ↑	rParis6k ↑	CUB-200 ↑	GPR1200 ↑	INSTRE ↑	MSCOCO Segmix ↑	VOC-Segmix ↑	IN-Cutmix ↑	IN-Dif-Diagonal ↑	IN-Dif-Outpaint ↑	Average Rank ↓
CD/IR	Multigrain [6], ResNet-50	22.87	44.74	3.67	37.09	56.72	27.77	25.8	67.54	79.44	24.91	5.9
	SSCD [49], ResNet-50	30.16	45.75	2.00	31.42	53.54	67.04	65.82	89.78	99.91	96.11	4.2
PT	ViT [16] S/16, IN1k	31.24	61.5	13.12	40.44	54.11	23.21	22.42	61.99	48.36	14.25	5.5
	ViT-B/16, IN12k	13.14	30.43	4.24	16.63	29.15	18.15	15.61	52.74	49.69	10.18	9.2
	ViT-B/16, CLIP [53] on LAION [60]	39.92	68.92	8.6	62.13	73.19	20.37	17.91	59.5	47.54	8.76	5.4
	Swin Transformer [37], Base, IN1k	40.06	72.07	15.49	54.09	68.46	24.51	24.31	74.79	40.74	14.75	4.1
SSL	MoCo [15], ViT-B/16	30.25	51.6	4.94	37.98	51.88	36.41	32.9	65.98	59.12	20.61	5.1
	MoCo, ViT-B/16 + CutMix [71]	25.01	46.73	3.44	32.23	48.58	32.83	26.11	55.74	62.88	46.96	6.5
	VicRegL [5], ResNet-50	28.4	53.79	3.02	34.95	50.98	40.58	37.76	69.74	80.02	40.93	5.0
	DINO [12], ViT-B/16, split-product	32.14	45.43	5.76	29.41	50.06	46.42	45.29	93.53	98.92	95.86	4.1

Memorization of Diffusion Models



Similarity metric:

- inner product
- token-wise inner product

Diffusion model: DDPM

Dataset: Celeb-A

The top-2 matches of diffusion models trained on 300, 3000, and 30000 images (the full set is 30000).

Results:

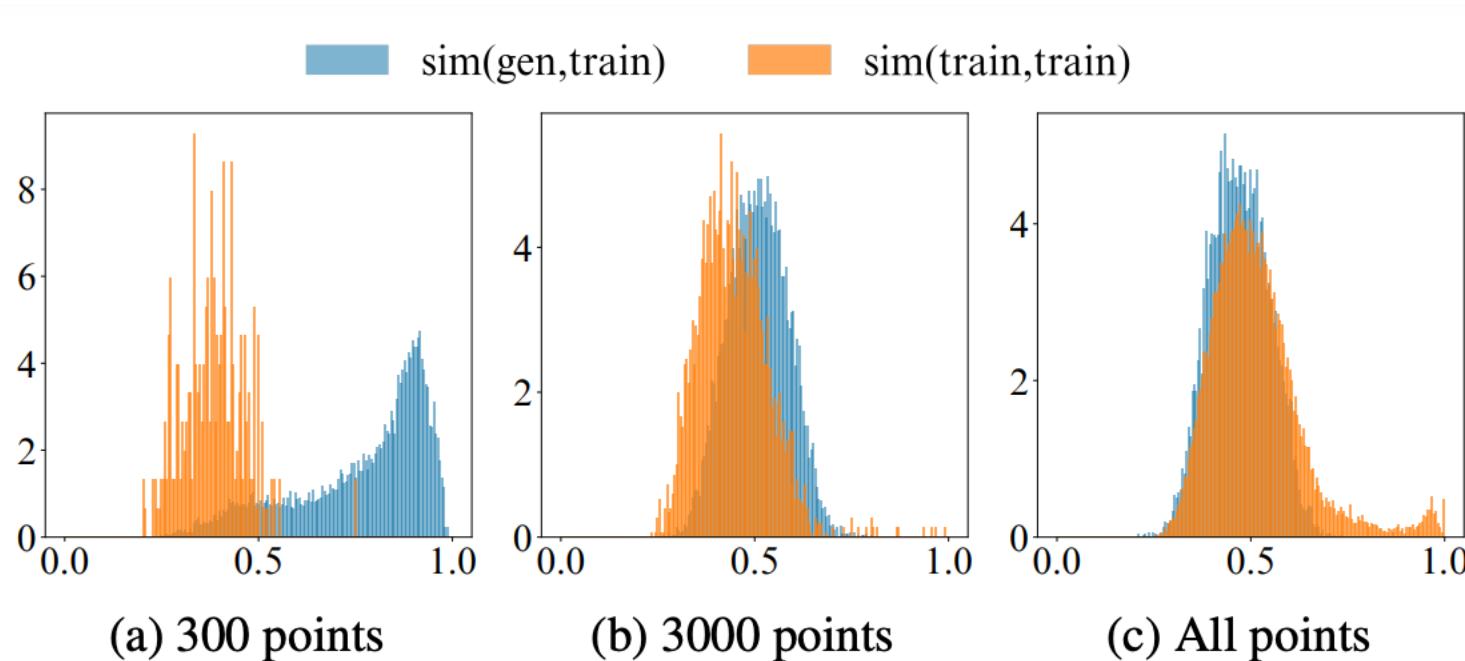
Green: copy

Blue: close but no exact copy

Others: similar but not the same

Memorization of Diffusion Models

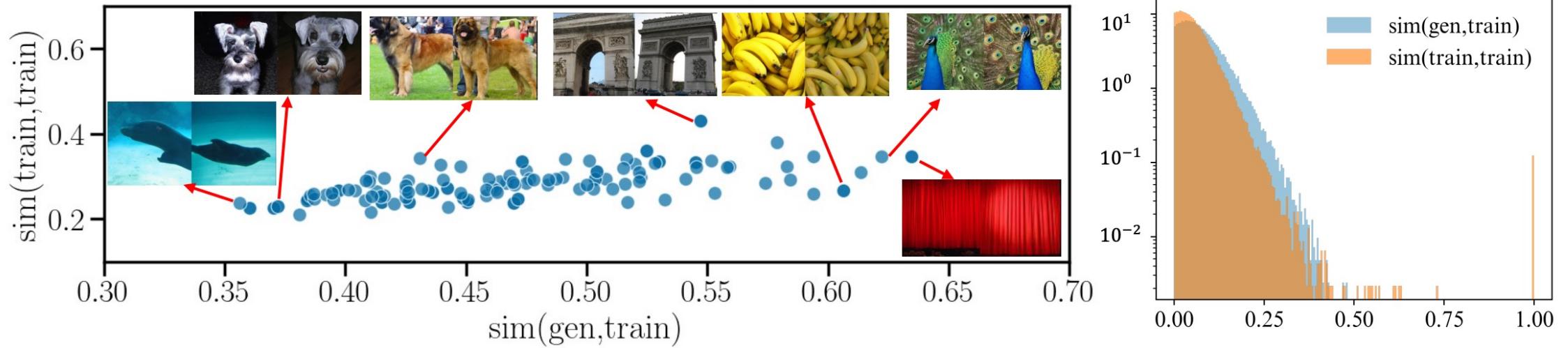
□ 数据越少Copy越多



Gen-train vs train-train similarity score distribution

Memorization of Diffusion Models

□ Case study: ImageNet LDM

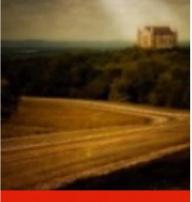


Many close copy but no exact match (similarity score <0.65)

Most similar: theater curtain, peacock, and bananas
Least similar: sea lion, bee, and swing

Memorization of Diffusion Models

□ Case study: Stable Diffusion

Caption Source	Generation	Top Match	
			<p>Source caption: Hill Country Castle by R Del Angel Match caption: Ben Hogan Portrait Golf Legend (2014) by GinetteCallaway</p>
			<p>Source caption: Captain Marvel Exclusive Ccxp Poster Released Online By Marvel Match caption: Marvel Studios releases new poster of 'Captain Marvel'</p>
			<p>Source caption: Rosie Huntington-Whiteley short hair (2015 Vanity Fair Oscar Party) (Venturelli, photographer for Getty Images) Match caption: 2017 Vanity Fair Oscar Party Hosted By Graydon Carter - Arrivals</p>

LAION Aesthetics v2 6+: 12M images
Random select 9000 images as source and use their captions to prompt

Memorization of Diffusion Models

□ Case study: Stable Diffusion



Key words

Prompt: <The description of the wall art> Canvas Wall Art Print

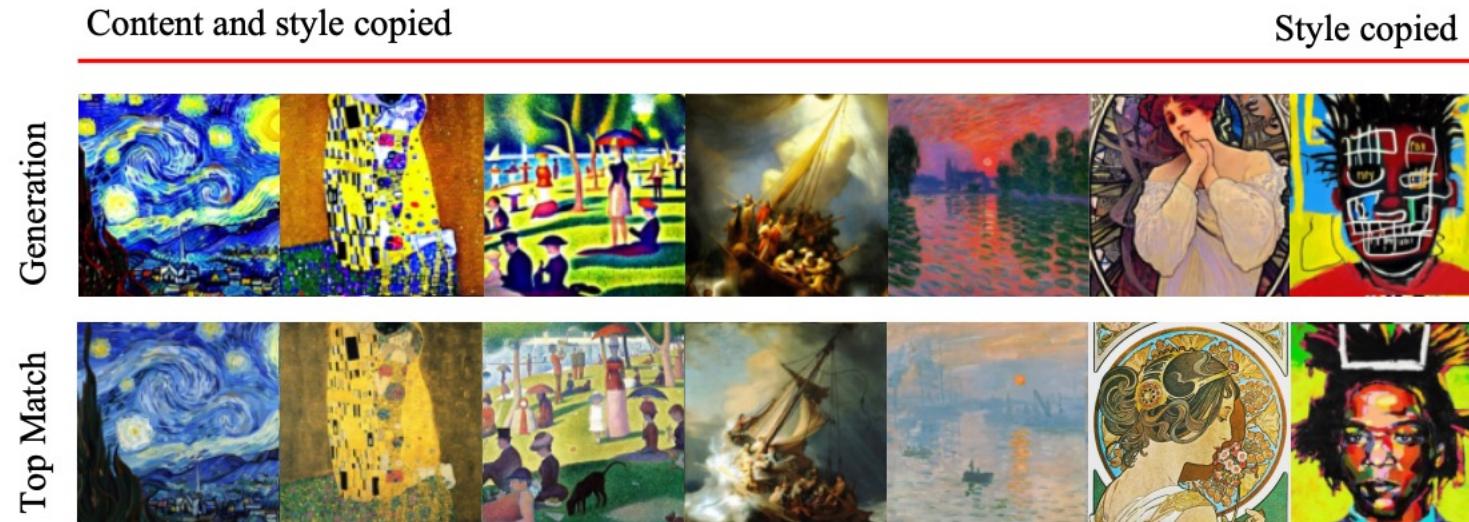


Prompt: A **painting** of the Great **Wave** off Kanagawa by Katsushika Hokusai

Some keywords (those in red) are associated with certain fixed patterns.

Memorization of Diffusion Models

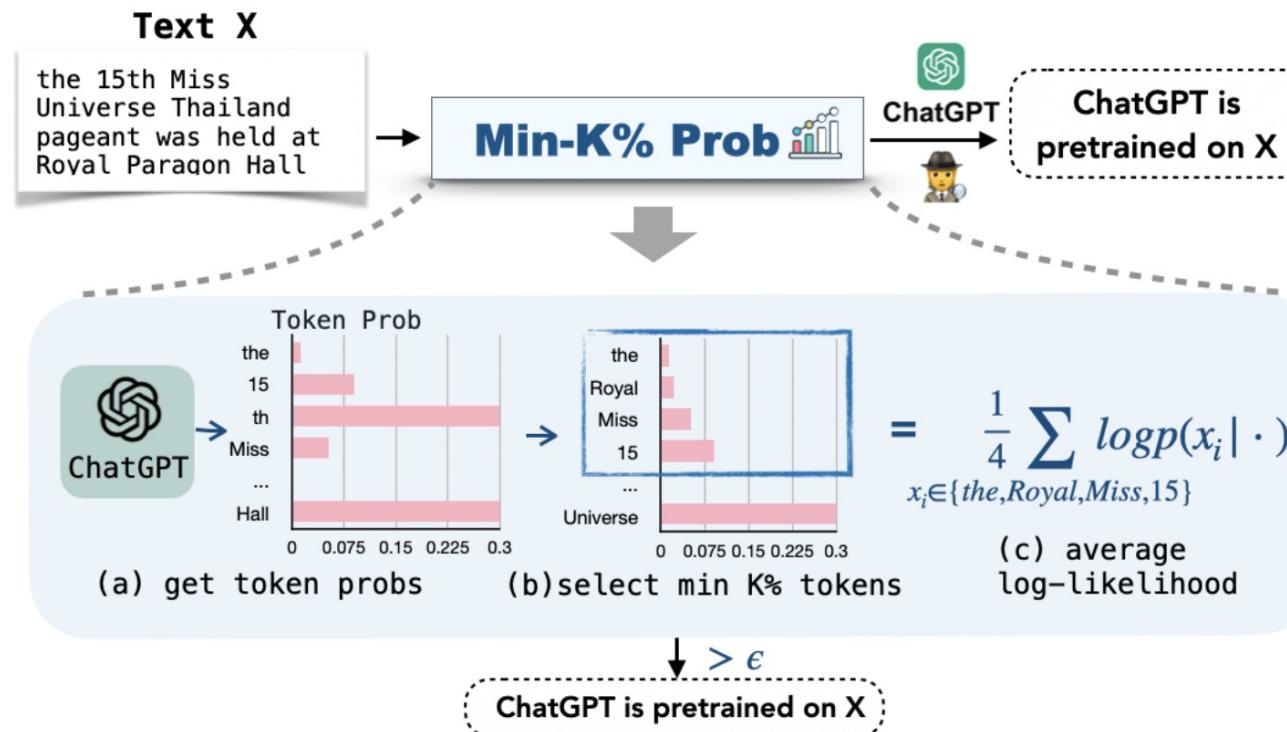
□ Case study: Stable Diffusion



Style copying using text prompt: <Name of the painting> by <name of the artist>

Memorization of Large Language Models (LLMs)

□ Pretraining data detection

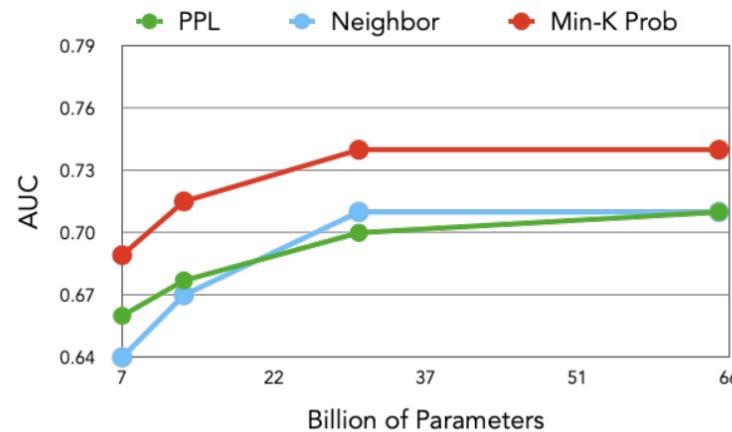


Memorization of Large Language Models (LLMs)

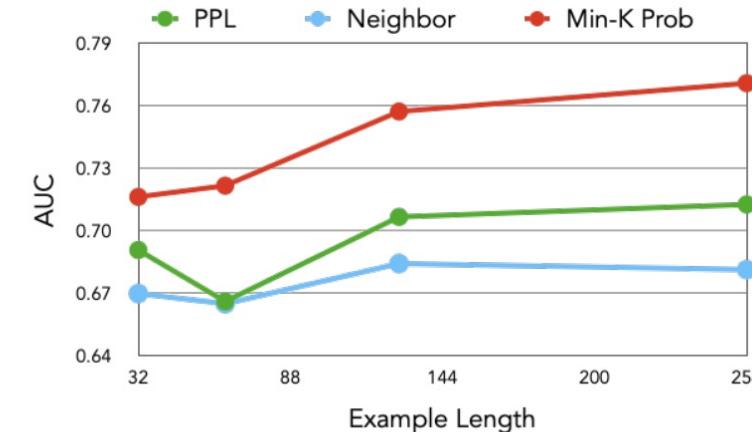
□ Detection on WIKIMIA

A dynamic benchmark: WIKIMIA

$$\text{MIN-K\% PROB}(x) = \frac{1}{E} \sum_{x_i \in \text{Min-K\%}(x)} \log p(x_i|x_1, \dots, x_{i-1}).$$



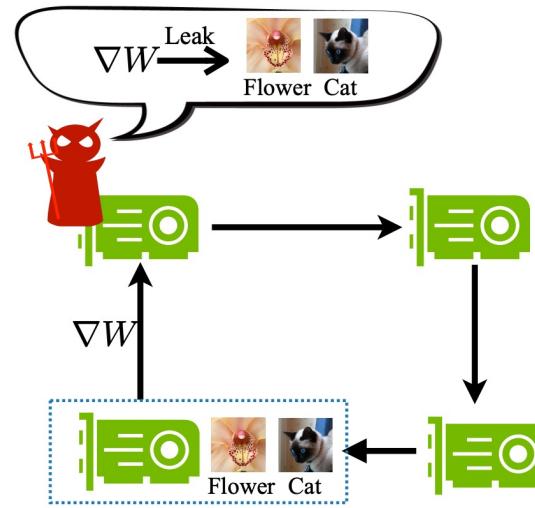
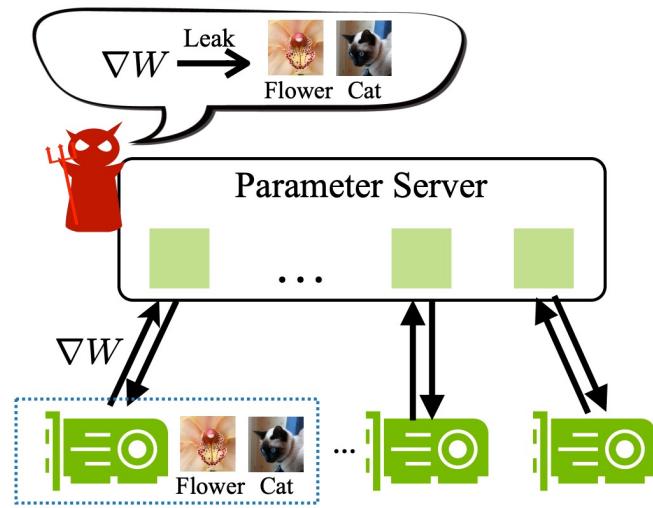
(a) AUC score vs. model size



(b) AUC score vs. text length

白盒窃取

□ 白盒窃取需要利用梯度信息，也称梯度逆向攻击
(Gradient Inversion Attack)



两种分布式训练范式

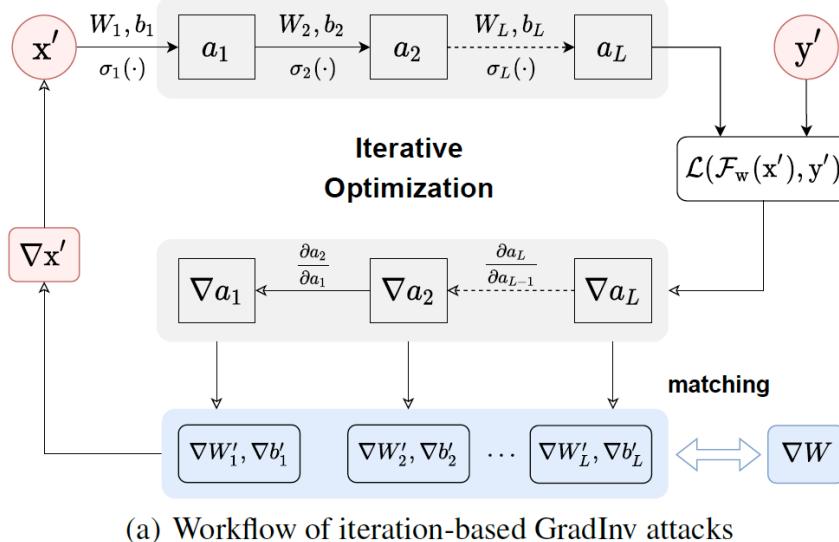
□ 针对梯度共享的训练：

- 分布式训练
- 联邦学习
- 并行训练
- 无中心化训练

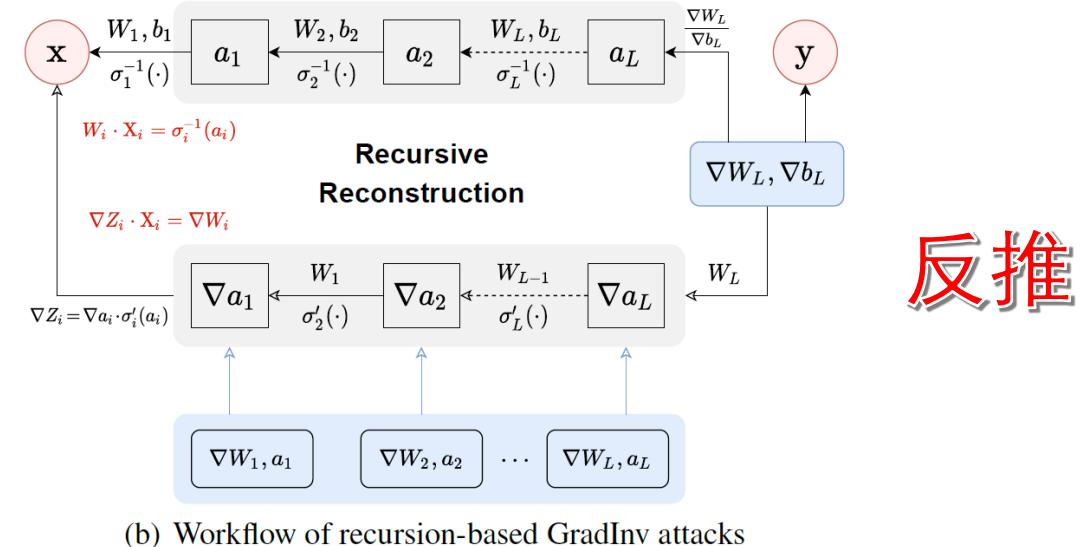
白盒窃取

□ 白盒窃取需要利用梯度信息，也称梯度逆向攻击
(Gradient Inversion Attack)

逼近



迭代逆向

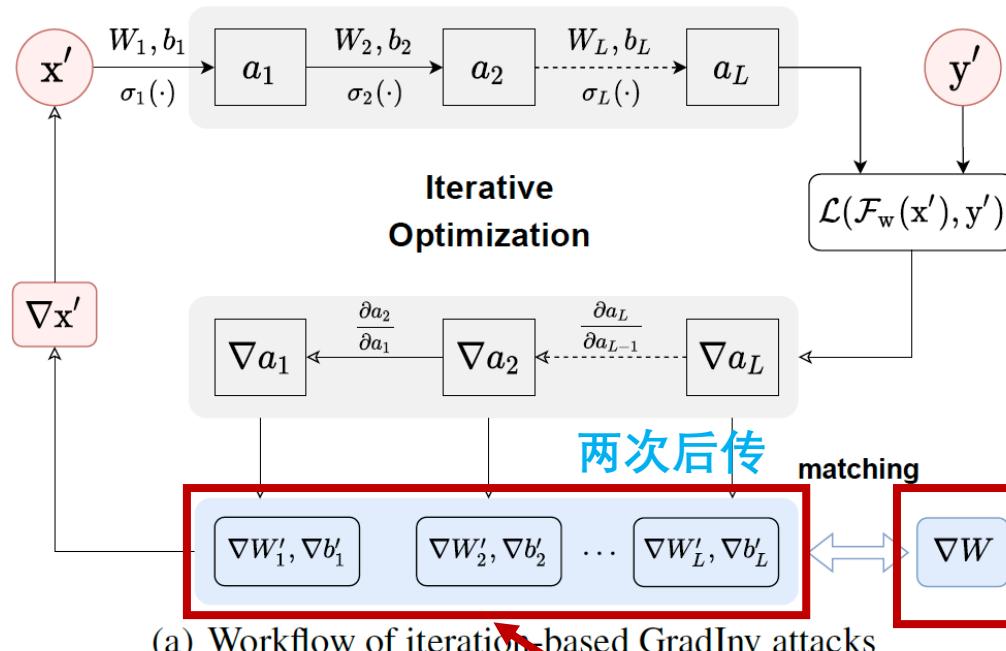


(逐层) 递归逆向

白盒窃取：迭代逆向

口 迭代逆向：通过构造数据来接近真实梯度

一次前传



(a) Workflow of iteration-based GradInv attacks

生成数据产生的梯度

口关键点：

- 如何初始化 x'
- Batch大小
- 模型大小
- 图像分辨率大小
- 有时需要梯度分拆

真实梯度，假设已知

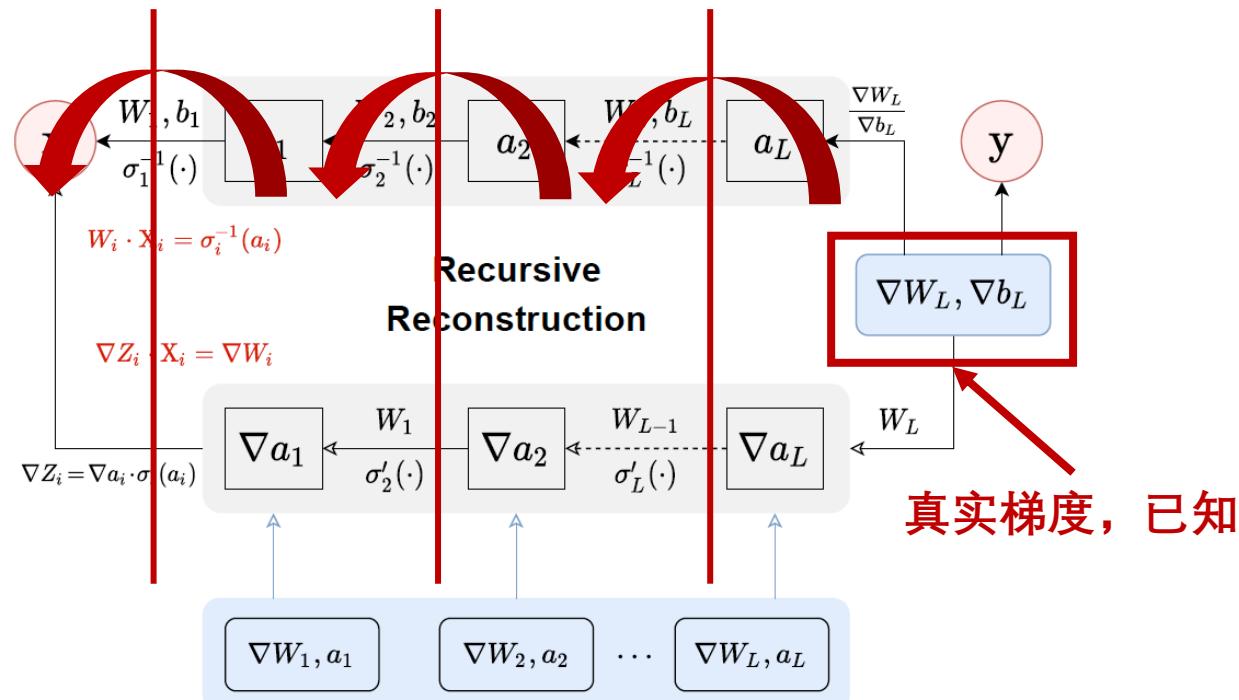
白盒窃取：迭代逆向

口 已有工作汇总

Publication	Data Initialization		Model Training		Grad Matching		Additional
	Distribution	Resolution	Network	Batch size	Loss-fn	Optimizer	
<i>GradInv attacks of iteration-based framework</i>							
DLG [Zhu <i>et al.</i> , 2019]	Gaussian	64×64	LeNet	8	ℓ_2 dist	L-BFGS	–
iDLG [Zhao <i>et al.</i> , 2020]	Uniform ^L	32×32	LeNet	1	ℓ_2 dist	L-BFGS	–
CPL [Wei <i>et al.</i> , 2020b]	Geometric	128×128	LeNet	8	ℓ_2 dist	L-BFGS	\mathcal{R}_y regularizer
InvGrad [Geiping <i>et al.</i> , 2020]	Gaussian ^L	224×224	ResNet ^T	8 (100)	Cosine	Adam	\mathcal{R}_{TV} regularizer
SAPAG [Wang <i>et al.</i> , 2020]	Constant	224×224	ResNet ^T	8	Gauss	AdamW	Gaussian kernel
GradInversion [Yin <i>et al.</i> , 2021]	Gaussian ^L	224×224	ResNet ^T	48	ℓ_2 dist	Adam	$\mathcal{R}_{\text{fidel}} + \mathcal{R}_{\text{group}}$
GradDisagg [Lam <i>et al.</i> , 2021]	Gaussian	32×32	MLP	32 (128)	ℓ_2 dist	L-BFGS	Participant info
GradAttack [Huang <i>et al.</i> , 2021]	Gaussian ^L	224×224	ResNet ^T	128	Cosine	Adam	No BN + labels
Bayesian [Balunović <i>et al.</i> , 2022]	Gaussian	32×32	ConvNet ^T	1 (32)	Cosine	Adam	Known $p(g x)$
CAFE [Jin <i>et al.</i> , 2021]	Uniform	32×32	Loop-Net	100	ℓ_2 dist	SGD	In Vertical-FL
GIAS [Jeon <i>et al.</i> , 2021]	Latent	64×64	ResNet ^T	4	Cosine	Adam	GAN-based

白盒窃取：递归逆向

口 递归逆向：基于真实梯度追层逆向推导



(b) Workflow of recursion-based GradInv attacks

$$\begin{cases} \mathbf{W}_i \cdot \mathbf{x}_i = Z_i \\ \nabla Z_i \cdot \mathbf{x}_i = \nabla \mathbf{W}_i \end{cases}$$

口关键点：

- 图像大小 (32x32)
- Batch大小 (大多为1)
- 模型大小

白盒窃取：递归逆向

□ 已有工作汇总

Publication	Data Initialization		Model Training		Grad Matching		Additional
	Distribution	Resolution	Network	Batch size	Loss-fn	Optimizer	
<i>GradInv attacks of recursion-based framework</i>							
PPDL-AHE [Phong <i>et al.</i> , 2018]	N/A	20×20	MLP	1	Gradient division	–	
PPDL-SPN [Fan <i>et al.</i> , 2020]	N/A	32×32	ConvNet	8	Linear solving	Noise analysis	
R-GAP [Zhu and Blaschko, 2021]	N/A	32×32	ConvNet	1	Inverse matrix	Rank analysis	
COPA [Chen and Campbell, 2021]	N/A	32×32	ConvNet	1	Least-squares	Pull-back const	

^L The labels can be directly identified or extracted from shared gradients.

^T The results of data recovery are compared in different model training states.

已有工作汇总

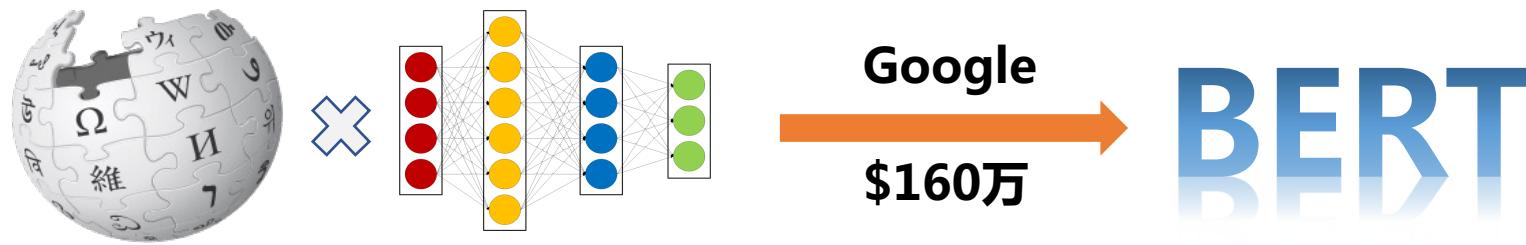
Category	Method	Publication	Key Contribution
Original Data	MixUp	[Zhang <i>et al.</i> , 2018]	Data enhancement by linearly combining the inputs
	InstaHide	[Huang <i>et al.</i> , 2020]	Encrypt the MixUp data with one-time secret keys
	Pixelization	[Fan, 2018; Fan, 2019]	Perturb the raw data with pixelization-based method
Training Model	Dropout	[Zheng, 2021]	Add an additional dropout layer before the classifier
	Local iters	[Wei <i>et al.</i> , 2020b]	Share gradients after multiple local training iterations
	Architecture	[Zhu and Blaschko, 2021]	Reduce the number of convolutional kernels properly
Shared Gradients	Aggregation	[Zhang <i>et al.</i> , 2020] [Lia and Togan, 2020]	Apply Homomorphic Encryption to protect gradients Utilize Secure Multi-Party Computation to aggregate
	Perturbation	[Sun <i>et al.</i> , 2021] [Wei <i>et al.</i> , 2021]	Perturb data representation layer and maintain utility Add adaptive noise with differential privacy guarantee
	Compression	[Vogels <i>et al.</i> , 2019] [Karimireddy <i>et al.</i> , 2019]	Compress the smaller values in gradients to zero Transmit the sign of gradients for model updates

This Week

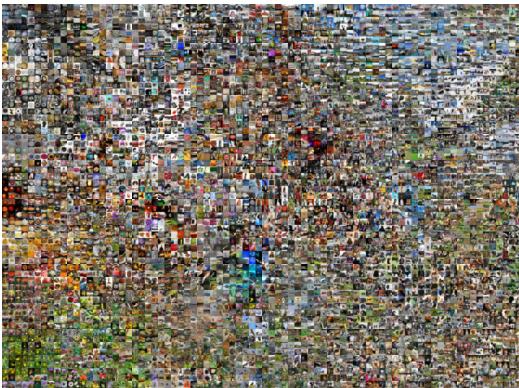
- Data Extraction Attack & Defense
- **Model Stealing Attack**
- Future Research



AI 模型训练代价高昂



大规模、高性能的AI模型训练耗费巨大



数据资源



计算资源



人力资源

模型窃取的动机

OpenVINO™



Google Cloud Platform

aws

Azure

宝贵的 AI 模型

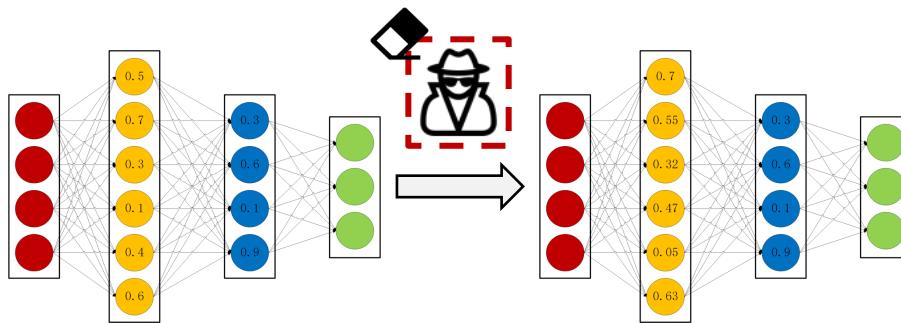
模型窃取
为其所用



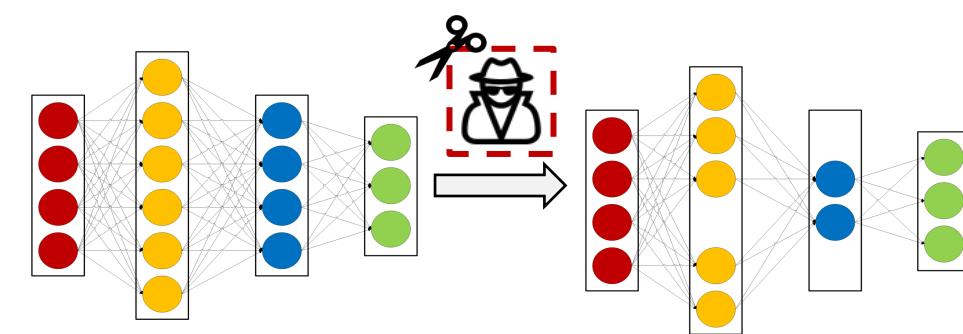
- 巨大的商业价值
- 尽量保持模型性能
- 不希望被发现

模型窃取的方式

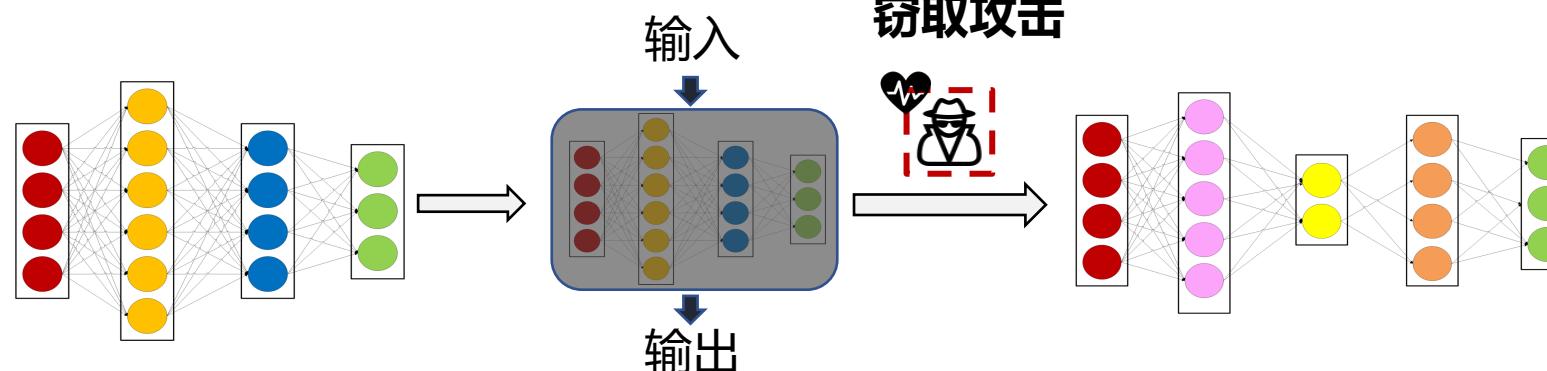
模型微调



模型剪枝



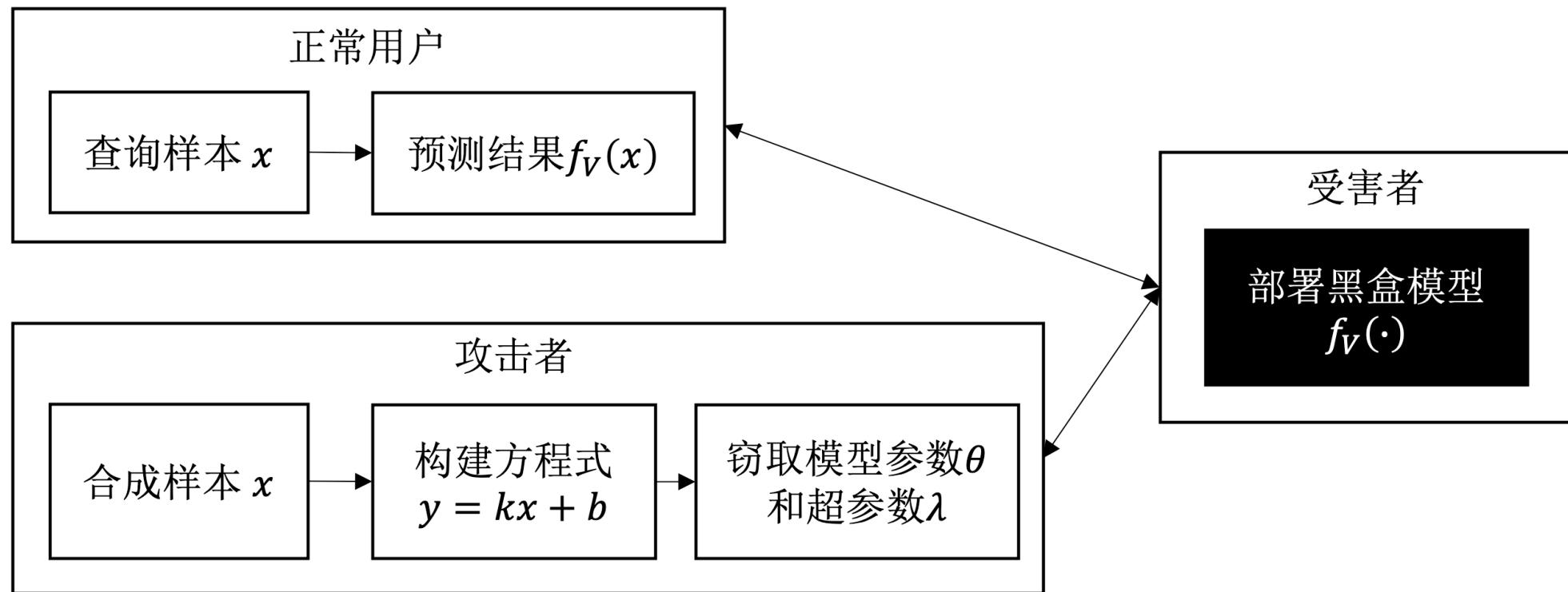
窃取攻击



Stealing machine learning models via prediction APIs, USENIX Security, 2016; Practical black-box attacks against machine learning, ASIACCS, 2017; Knockoff nets: Stealing functionality of black-box models, CVPR, 2019; Maze: Data-free model stealing attack using zeroth-order gradient estimation, CVPR, 2021;

基于方程式求解的攻击

口 攻击思路示例



基于方程式求解的攻击

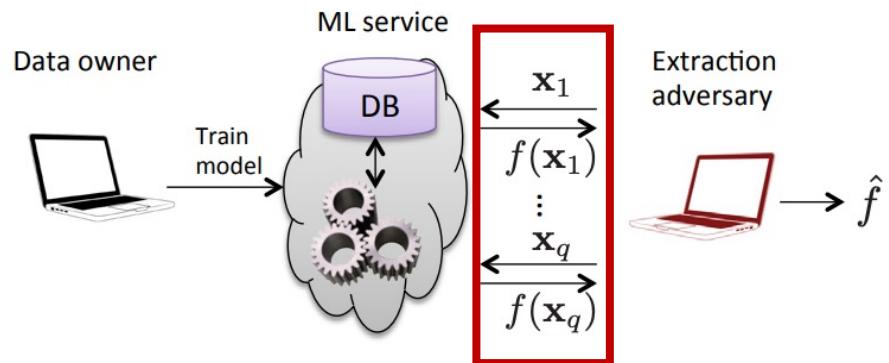
□ 100%窃取某些商业模型所需的查询数和时间

Service	Model Type	Data set	Queries	Time (s)
Amazon	Logistic Regression	Digits	650	70
	Logistic Regression	Adult	1,485	149
BigML	Decision Tree	German Credit	1,150	631
	Decision Tree	Steak Survey	4,013	2,088

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

基于方程式求解的攻击：窃取参数

□ 攻击算法



- 参数个数为d
 - 通过d+1个输入，构造d+1个下列方程
- $$\theta^\top x = \sigma^{-1}(f(x))$$
- 求解方程得到 θ

□ 主要特点：

- 针对传统机器学习模型:SVM、LR、DT
- 可精确求解，需要模型返回精确的置信度
- 窃取得到的模型还可能泄露训练数据（数据逆向攻击）

基于方程式求解的攻击：窃取超参

口 攻击思想：模型训练完了的状态应该是Loss梯度为0

$$\mathcal{L}(\theta) = \mathcal{L}(x, y, \theta) + \lambda R(\theta) \quad \longleftarrow \text{窃取超参数}\lambda$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = \mathbf{b} + \lambda \mathbf{a} = 0$$

$$\mathbf{b} = \begin{bmatrix} \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta_1} \\ \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial \mathcal{L}(x, y, \theta)}{\partial \theta_n} \end{bmatrix}, \quad \mathbf{a} = \begin{bmatrix} \frac{\partial R(\theta)}{\partial \theta_1} \\ \frac{\partial R(\theta)}{\partial \theta_2} \\ \vdots \\ \frac{\partial R(\theta)}{\partial \theta_n} \end{bmatrix}$$

$$\hat{\lambda} = -(\mathbf{a}^\top \mathbf{a})^{-1} \mathbf{a}^\top \mathbf{b}.$$

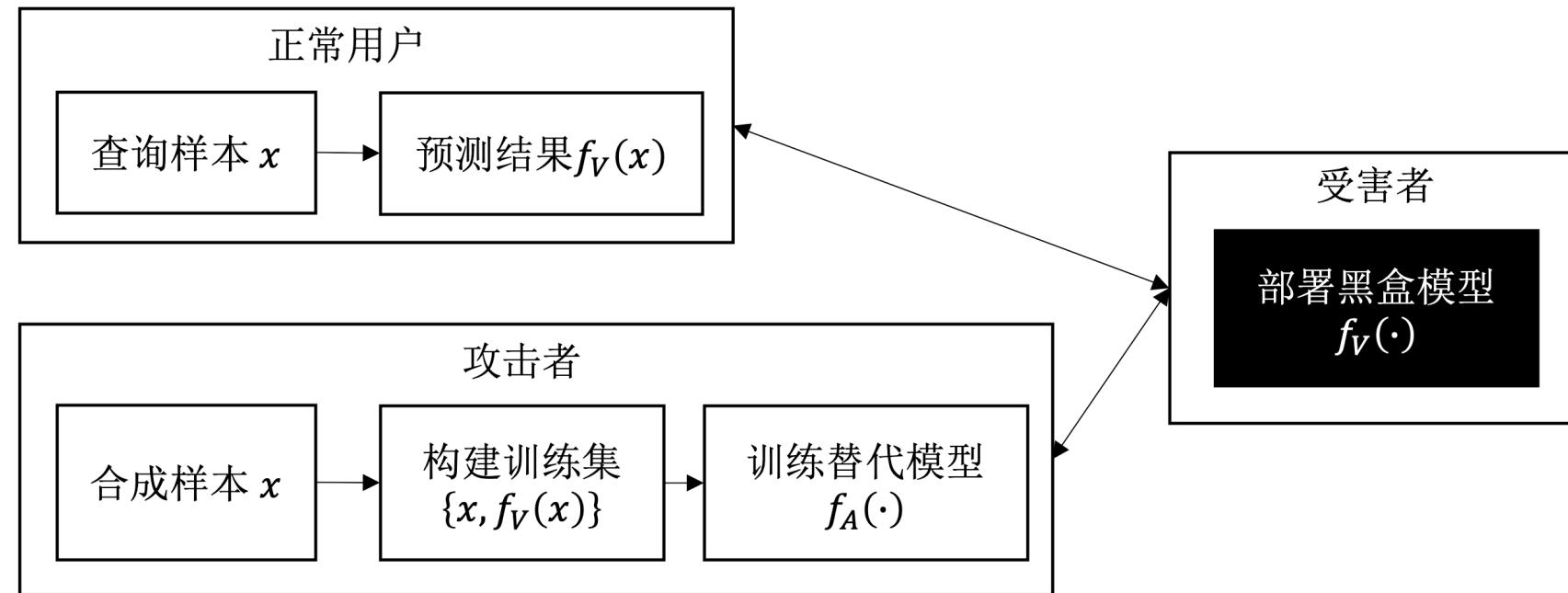
口 主要特点：

- 需要知道Loss的形式
- 需要在所有数据上做矩阵运算
- R只与模型参数 θ 有关



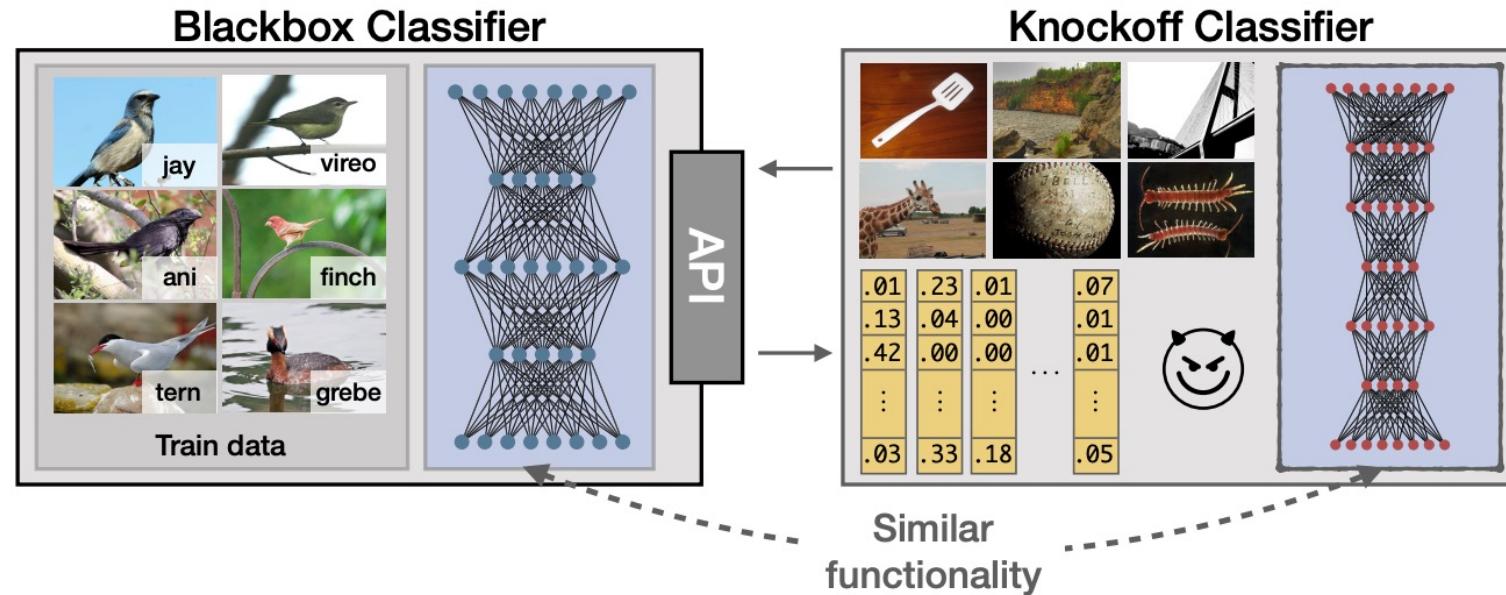
基于替代模型的攻击

口 攻击思想：在查询目标模型的过程中训练一个替代模型模拟其行为



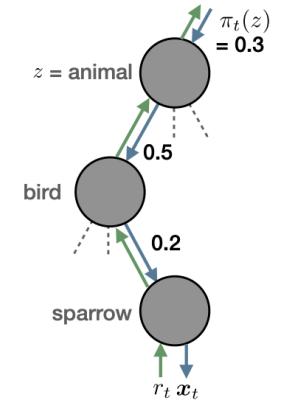
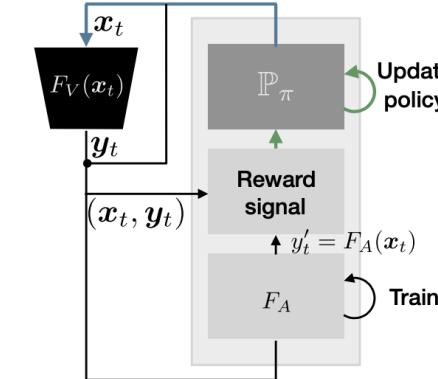
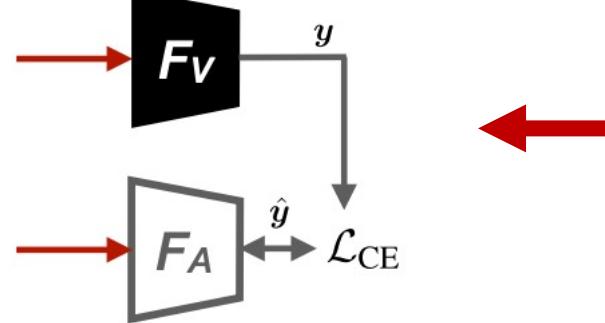
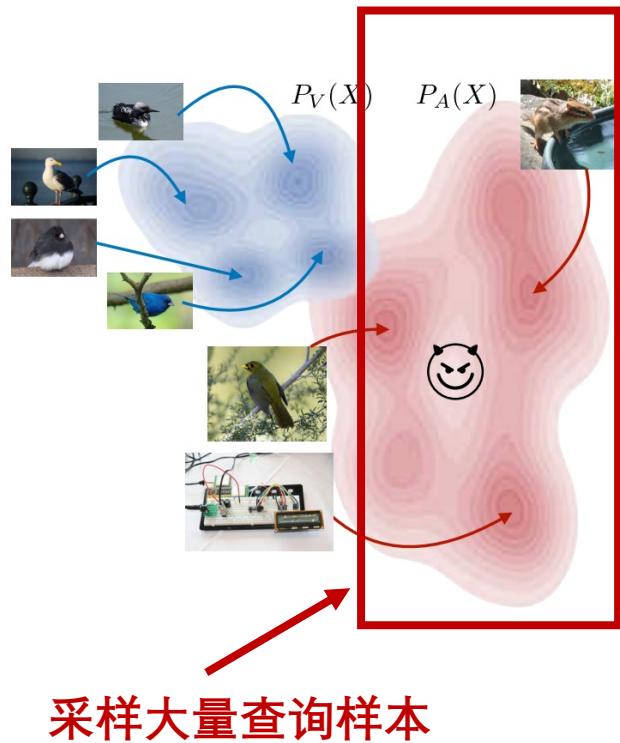
基于替代模型的攻击

□ Knockoff Nets 攻击：“仿冒网络”



基于替代模型的攻击

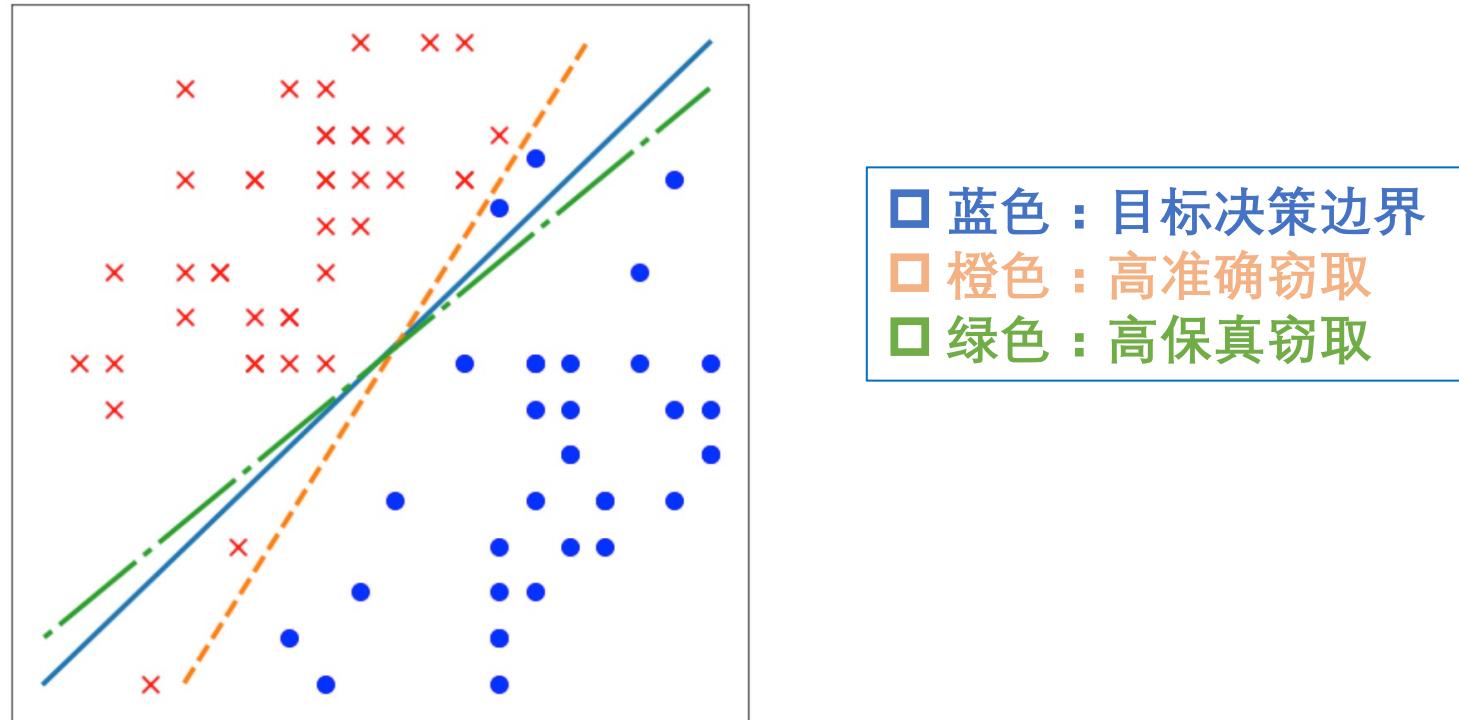
□ Knockoff Nets 攻击：攻击流程



强化学习，学习如何高效选择样本

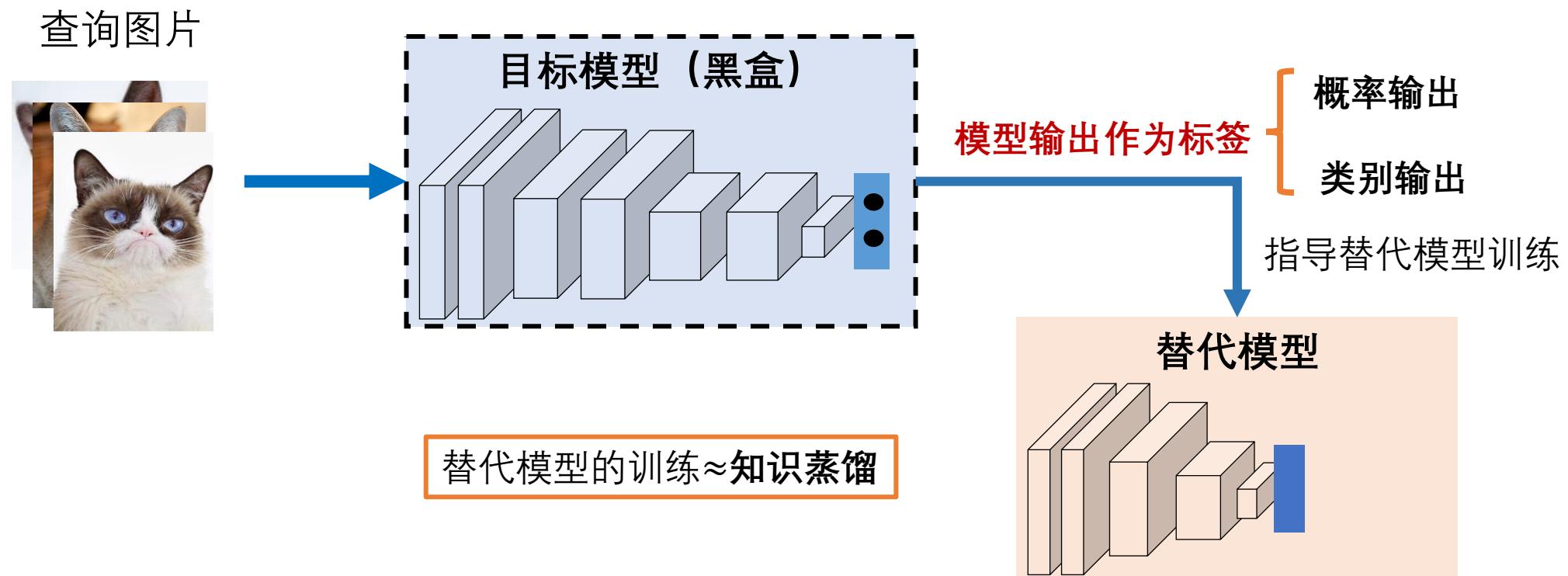
基于替代模型的攻击

□ 高准确 (accuracy) vs 高保真 (fidelity) 窃取攻击



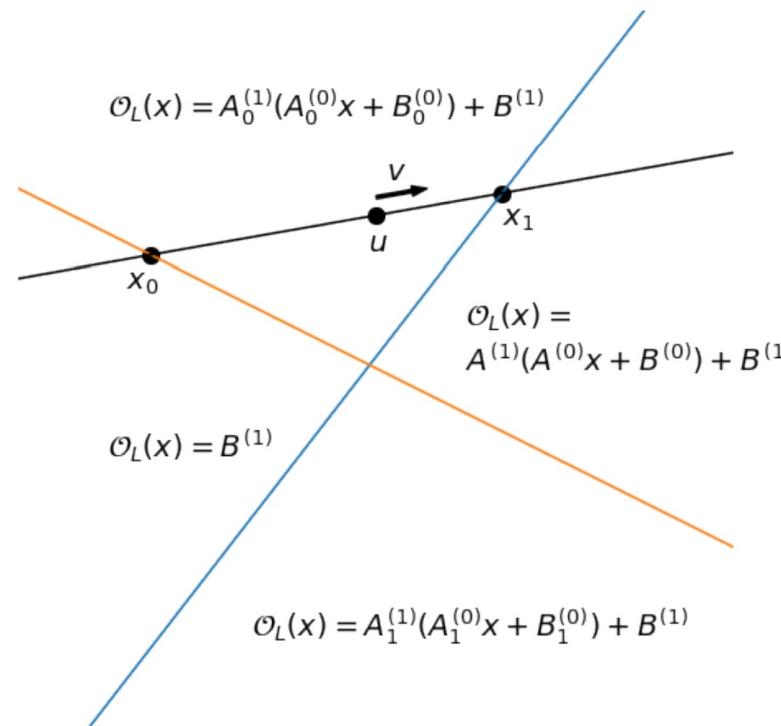
基于替代模型的攻击

□ 高准确 (accuracy) vs 高保真 (fidelity) 窃取攻击



基于替代模型的攻击

□ 功能等同窃取 Functionally Equivalent Extraction



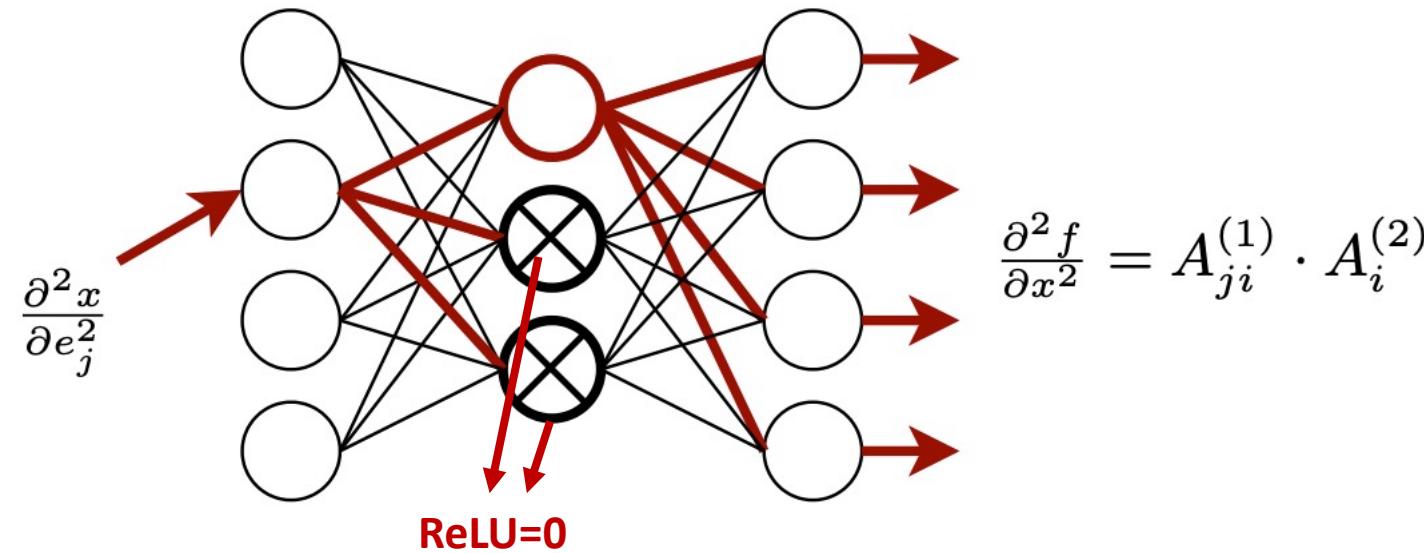
□ 攻击步骤：

- 寻找在某个Neuron上，让ReLU=0的关键点
- 在关键点两侧探索边界，确定对应权重
- 只能窃取两层网络

Symbol	Definition
d	Input dimensionality
h	Hidden layer dimensionality ($h < d$)
K	Number of classes
$A^{(0)} \in \mathbb{R}^{d \times h}$	Input layer weights
$B^{(0)} \in \mathbb{R}^h$	Input layer bias
$A^{(1)} \in \mathbb{R}^{h \times K}$	Logit layer weights
$B^{(1)} \in \mathbb{R}^K$	Logit layer bias

基于替代模型的攻击

口 加密分析窃取 Cryptanalytic Extraction



口 思想：ReLU的二级导为0 & 有限差分 (finite difference)

基于替代模型的攻击

□ 加密分析窃取 Cryptanalytic Extraction

窃取0-deep神经网络：

$$f(\mathbf{x}) = \mathbf{w}^{(1)} \cdot \mathbf{x} + b^{(1)}$$

$$f(\mathbf{x} + \mathbf{e}_i) - f(\mathbf{x}) = \mathbf{w}^{(1)} \cdot (\mathbf{x} + \mathbf{e}_i) - \mathbf{w}^{(1)} \cdot \mathbf{x} = \mathbf{w}^{(1)} \cdot \mathbf{e}_i$$

窃取1-deep神经网络：

$$f(\mathbf{x}) = \mathbf{w}^{(2)} \text{ReLU}(\mathbf{w}^{(1)} \mathbf{x} + b^{(1)}) + b^{(2)}$$

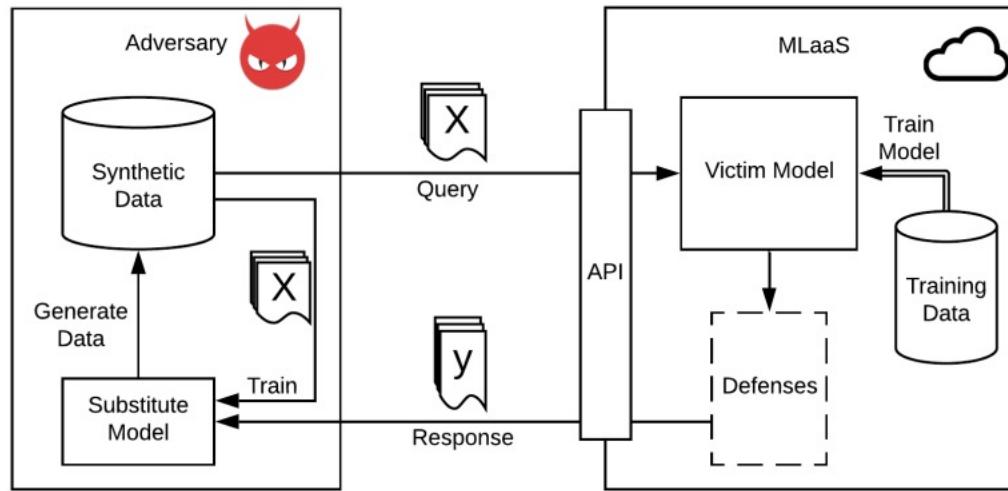
$$\alpha_+^i = \frac{\partial f(\mathbf{x})}{\partial \epsilon \mathbf{e}_i} \Big|_{x=x+\epsilon \mathbf{e}_i}$$

$$\alpha_-^i = \frac{\partial f(\mathbf{x})}{\partial \epsilon \mathbf{e}_i} \Big|_{x=x-\epsilon \mathbf{e}_i}$$

$$\frac{\alpha_+^k - \alpha_-^k}{\alpha_+^i - \alpha_-^i} = \frac{\mathbf{w}_{j,k}^{(1)}}{\mathbf{w}_{j,i}^{(1)}}$$

基于替代模型的攻击

□ 估计合成攻击 Estimation Synthesis (ES) Attack



思想：初始化合成数据集，然后根据模型返回训练替代模型

- **E-step**：在合成数据上知识蒸馏更新替代模型
- **S-step**：合成数据，使用对抗生成网络（GAN）

□ 特点：

- 不需要原始训练数据或先验
- 不需要目标模型先验

基于替代模型的攻击

□ ES攻击算法: 蒸馏+生成的结合

Algorithm 1 *ES Attack*

INPUT:

The black-box victim model f_v
Number of classes K
Number of stealing epochs N
Number of training epochs for each stealing epoch M

OUTPUT:

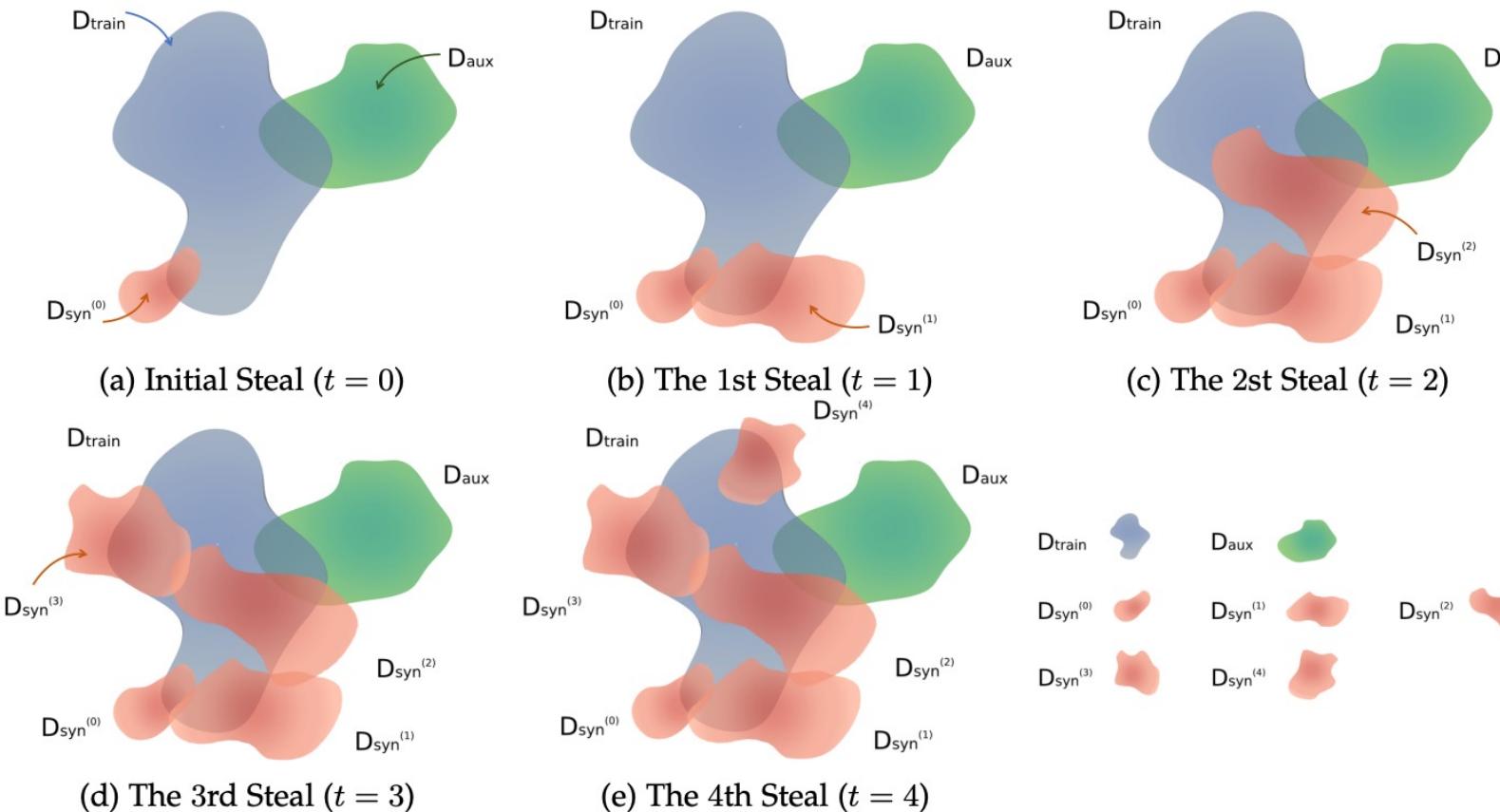
The substitute model $f_s^{(N)}$

- 1: Initialize a synthetic dataset $\mathcal{D}_{\text{syn}}^{(0)}$ by randomly sampling \mathbf{x} from a Gaussian distribution.
- 2: Construct an initial substitute model $f_s^{(0)}$ by initializing the parameters in the model.
- 3: **for** $t \leftarrow 1$ to N **do**
- 4: **E-Step:** Estimate the parameters in the substitute model $f_s^{(t)}$ using knowledge distillation for M epochs on the synthetic dataset $\mathcal{D}_{\text{syn}}^{(t-1)}$.
- 5: **S-Step:** Synthesize a new dataset $\mathcal{D}_{\text{syn}}^{(t)}$ based on the knowledge of the substitute model $f_s^{(t)}$.
- 6: **end for**
- 7: **return** $f_s^{(N)}$.



基于替代模型的攻击

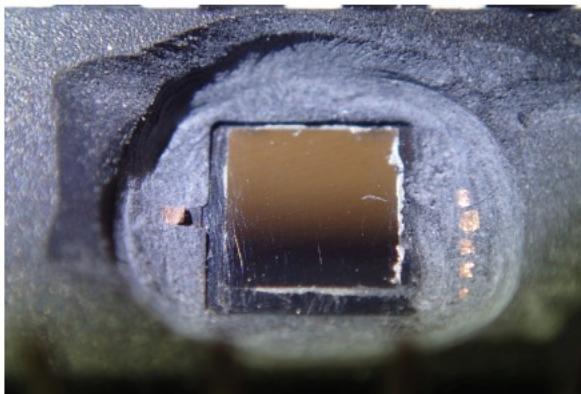
□ ES攻击合成的数据分布



train:原始训练数据
aux:公共数据集
syn:合成数据集

基于侧信道攻击的窃取

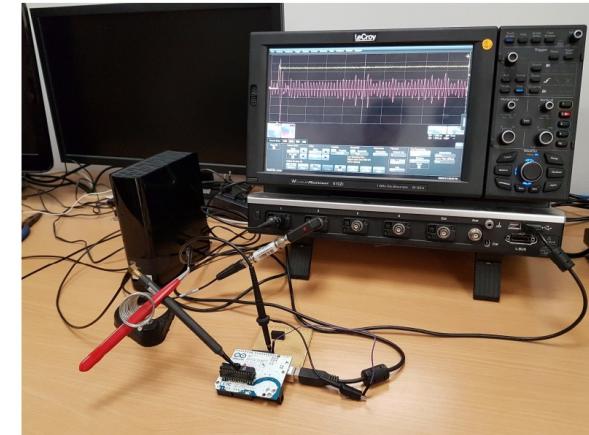
□ 侧信道 (side-channel) 攻击窃取神经网络



(a) Target 8-bit microcontroller mechanically decapsulated



(b) Langer RF-U 5-2 Near-field Electromagnetic passive Probe



(c) The complete measurement setup

通过探测运行神经网络的微处理器的电力使用情况，来窃取神经网络的权重

LLM Extraction

□ 基于LLAMA 7B模拟ChatGPT

Stanford University



People

Report

Research

Blog

Workshop

Courses

HELM

Ecosystem graphs

Code

We introduce [Alpaca 7B](#), a model fine-tuned from the LLaMA 7B model on 52K instruction-following demonstrations. On our preliminary evaluation of single-turn instruction following, Alpaca behaves qualitatively similarly to OpenAI's text-davinci-003, while being surprisingly small and easy/cheap to reproduce (< 600\$).

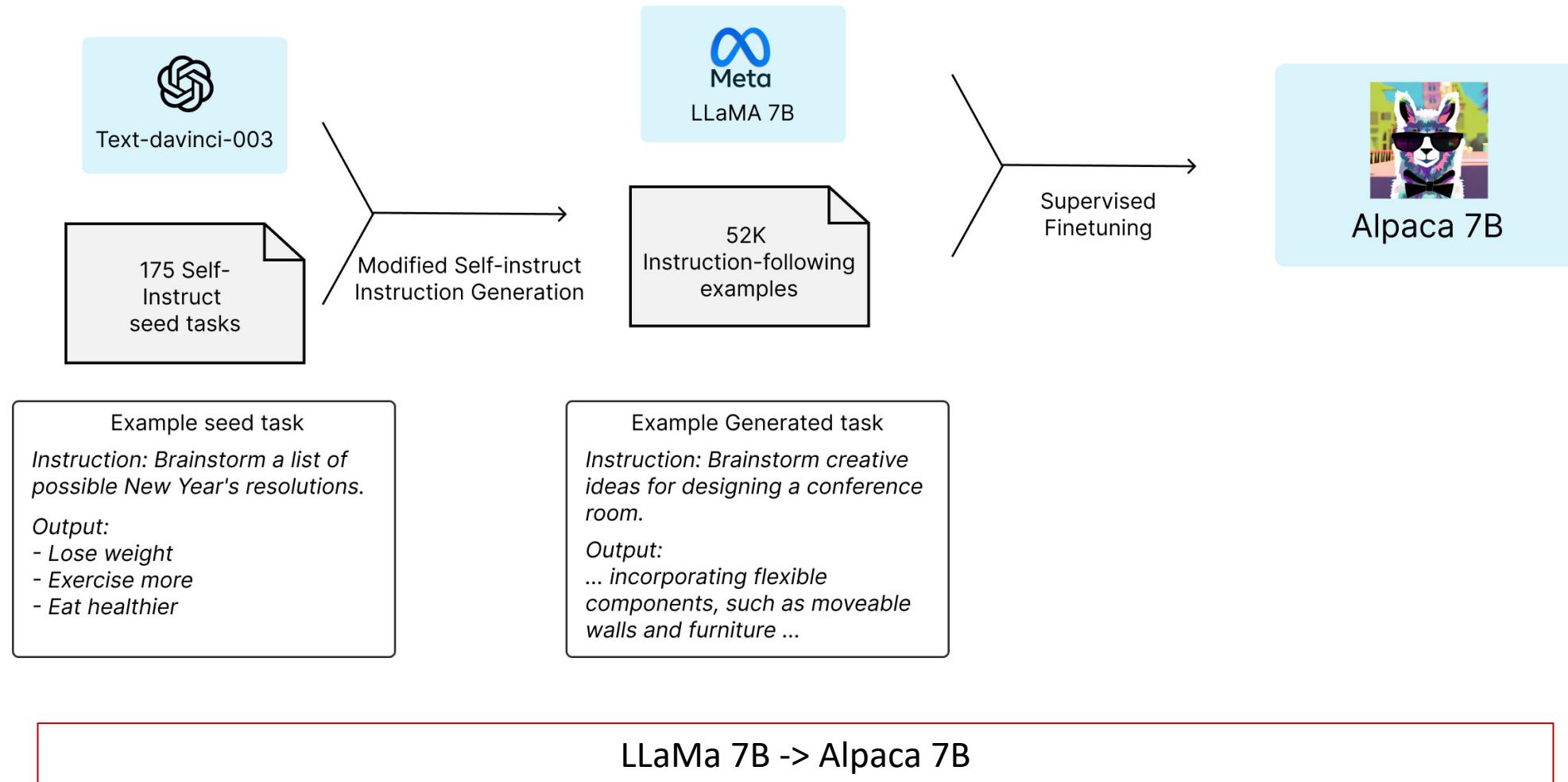
[Web Demo](#) [GitHub](#)

Stanford
Alpaca



斯坦福仅用**600美元**就完成了**对ChatGPT的窃取**，其通过跟OpenAI的text-davinci-003对话，抽取52000个对话样本，再微调Meta的开源LLAMA 7B语言模型得到Alpaca

LLM Extraction



Future Research

□ Attack:

- Data/model extraction attacks on LLMs
- Attacking large-scale datasets and models

□ Defense:

- Attack detection and mitigation: detect malicious queries
- Attack as a defense, adding watermark/backdoor into the extracted data



谢谢 !

