



Backdoor Attacks and Defenses

Lecturer: Dr. Xingjun Ma

**School of Computer Science,
Fudan University**

Autumn, 2022

Recap: week 7

1. Data Poisoning: Attacks and Defenses

- A Brief History of Data Poisoning
- Data Poisoning Attacks
- Data Poisoning Defenses
- Poisoning for Data Protection



Adversarial Attack Competition: Phase 2

Phase 1 Phase 2 **Phase 3**

Phase description
Create an attack method and submit the code as submission. Your code should follows the submission template. Feedback will be provided on the 1000 "hard" test images. We will test your code on 1 robustly trained model.

Max submissions per day: 10

Max submissions total: 500

RESULTS							
#	User	Entries	Date of Last Entry	Score ▲	Error Rate ▲	Efficiency Score ▲	Detailed Results
1	keren	5	10/20/22	0.2290 (1)	0.2230 (4)	0.5971 (5)	View
2	Shadow_H	14	10/19/22	0.2270 (2)	0.2240 (2)	0.3008 (10)	View
3	yong_xie	14	10/19/22	0.2266 (3)	0.2250 (1)	0.1594 (13)	View
4	songtianwei	6	10/14/22	0.2249 (4)	0.2170 (9)	0.7851 (4)	View
5	Yuxuan_Wang	9	10/12/22	0.2243 (5)	0.2230 (4)	0.1290 (14)	View
6	hsj576	5	10/19/22	0.2241 (6)	0.2230 (3)	0.1113 (16)	View
7	wangzhix	6	10/18/22	0.2241 (7)	0.2210 (6)	0.3087 (9)	View
8	xinwang22	28	10/19/22	0.2231 (8)	0.2200 (7)	0.3120 (8)	View
9	pywang	6	10/19/22	0.2230 (9)	0.2180 (8)	0.5050 (7)	View
10	kejiefang	11	10/21/22	0.2227 (10)	0.2200 (7)	0.2698 (11)	View
10	weijiezhang	4	10/19/22	0.2227 (10)	0.2200 (7)	0.2698 (11)	View
11	snow_zk	5	10/19/22	0.2223 (11)	0.2210 (6)	0.1268 (15)	View
11	WeijianMa	1	10/18/22	0.2223 (11)	0.2210 (6)	0.1268 (15)	View

Link:https://codalab.lisn.upsaclay.fr/competitions/7556?secret_key=d4a3b1fa-66e2-4a80-8ce6-b5f99e518979

Starting kit: https://codalab.lisn.upsaclay.fr/competitions/7556?secret_key=d4a3b1fa-66e2-4a80-8ce6-b5f99e518979#learn_the_details_get_starting_kit



Backdoor Attacks and Defenses

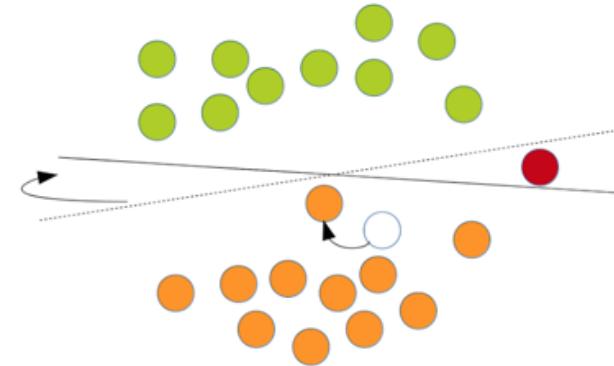
- A Brief History of Backdoor Learning
- Backdoor Attacks
- Backdoor Defenses
- Future Research



Backdoor vs (Pure) Poisoning

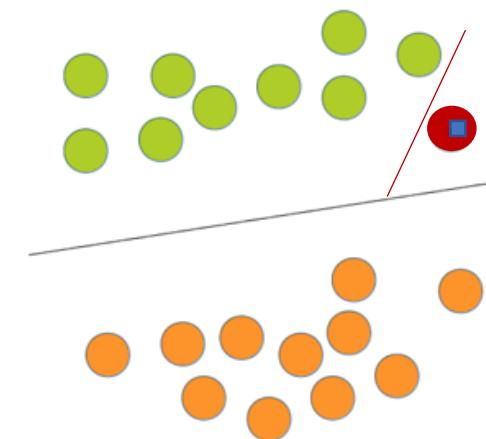
- **Poisoning attack**

- Training time attack
- Change classification boundary



- **Backdoor attack**

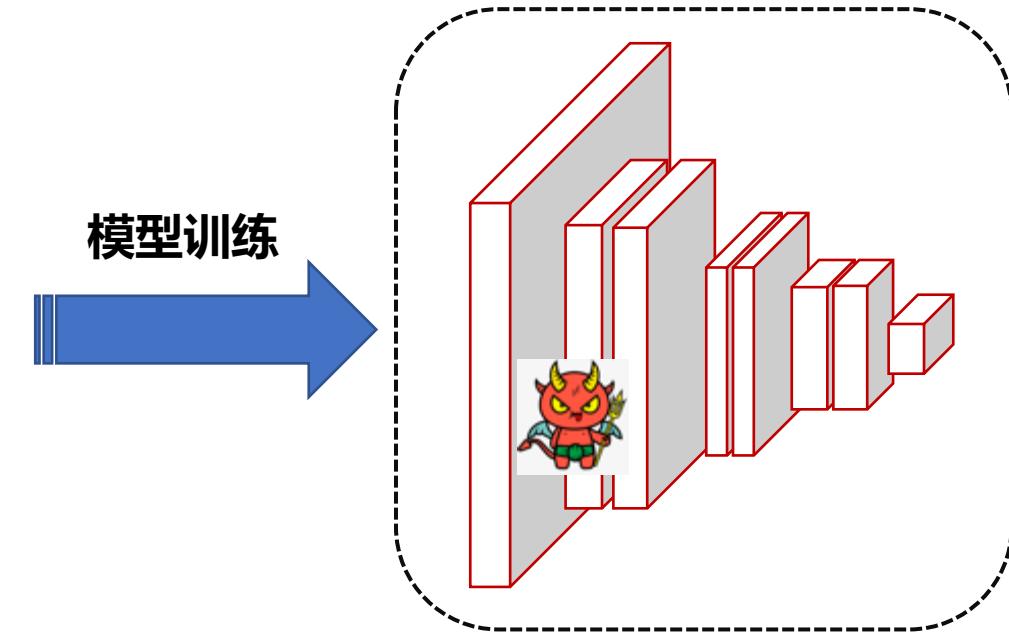
- Training time attack
- Does not change the original boundary
- Add new boundary



后门攻击 – 动机



大量互联网数据可能存在后门样本



后门模型

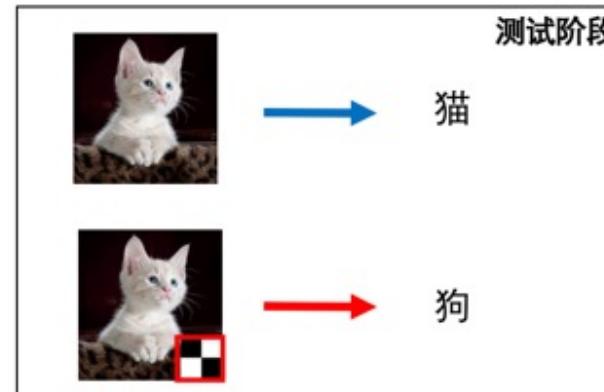
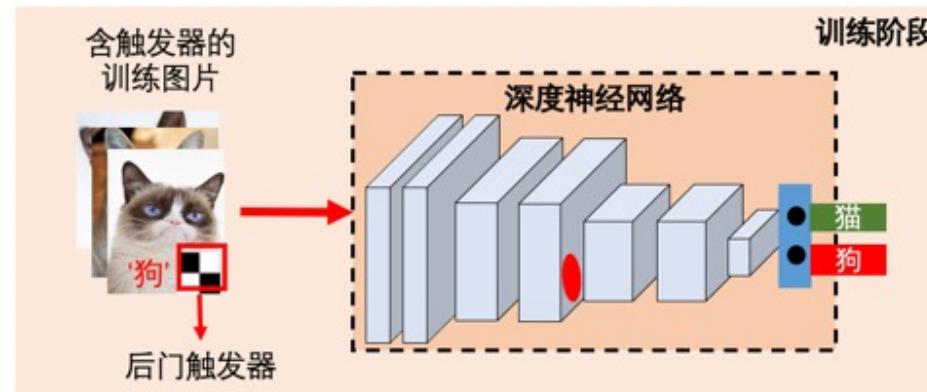
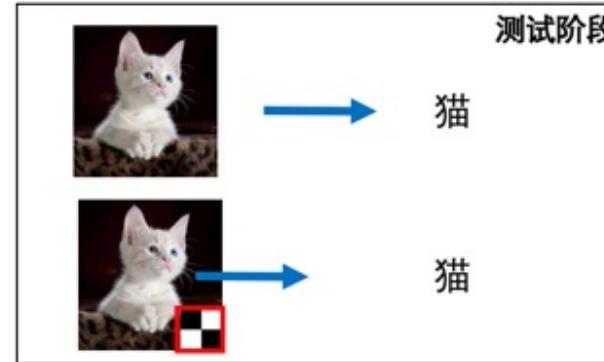
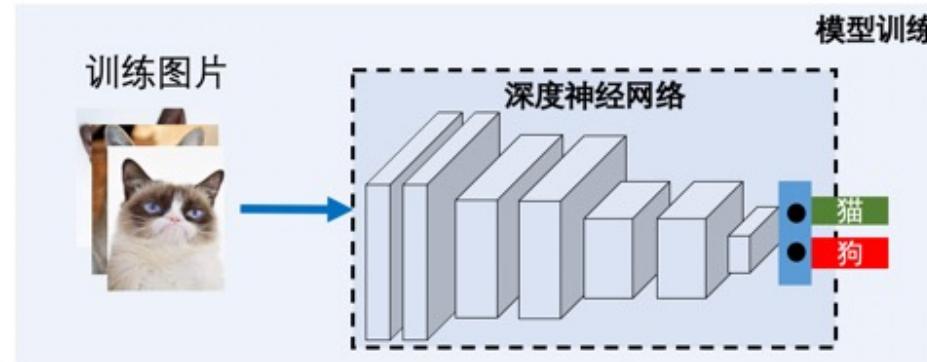
But , 后门攻击 ! = 投毒攻击

- 这是两个不同的话题
- 数据投毒是后门攻击的一种实现方式



后门攻击 - 流程

步骤1：后门注入；步骤2：后门激活



■ 后门攻击的特点：

- ✓ 模型在干净数据上性能不变
- ✓ 触发器出现即预测后门类别

后门攻击 - 例子

数字触发器



物理世界攻击



Gu et al. "Badnets: Identifying vulnerabilities in the machine learning model supply chain." *arXiv:1708.06733* (2017).

后门攻击 - 方法分类

■ 脏标签攻击: **添加触发器并修改类别**

- ✓ BadNets (Gu *et al.*, 2019)
- ✓ Trojan attack (Liu *et al.*, 2018)
- ✓ Blend attack (Chen *et al.*, 2017)

■ 净标签攻击: **只添加触发器**

- ✓ Clean-label attack (CL) (Turner *et al.*, 2019)
- ✓ Sinusoidal signal attack (SIG) (Barni *et al.*, 2019)
- ✓ **Reflection backdoor (Refool)** (Liu *et al.*, 2020, ECCV)
- ✓ **Video backdoor** (Zhao *et al.*, 2020, CVPR)

后门攻击 - 优化目标

■ 隐蔽性

- 尽量少的毒化样本
- 尽量隐蔽的触发器
- 尽量小的影响模型在干净样本上的性能

■ 成功率

- 尽量高的攻击成功率
- 可完成多目标攻击

■ 迁移性

- 迁移到不同的训练方法
- 迁移到不同的模型

■ 鲁棒性

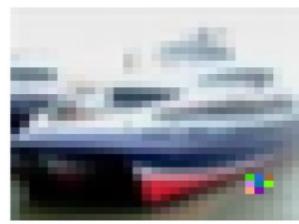
- 可躲避后门检测防御
- 可躲避后门移除防御

后门攻击

■ 六种经典攻击所使用的触发器样式



BadNets



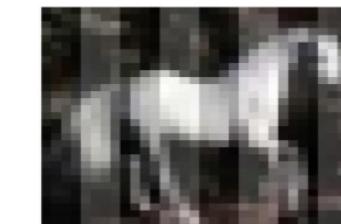
Trojan



Blend



CL



SIG

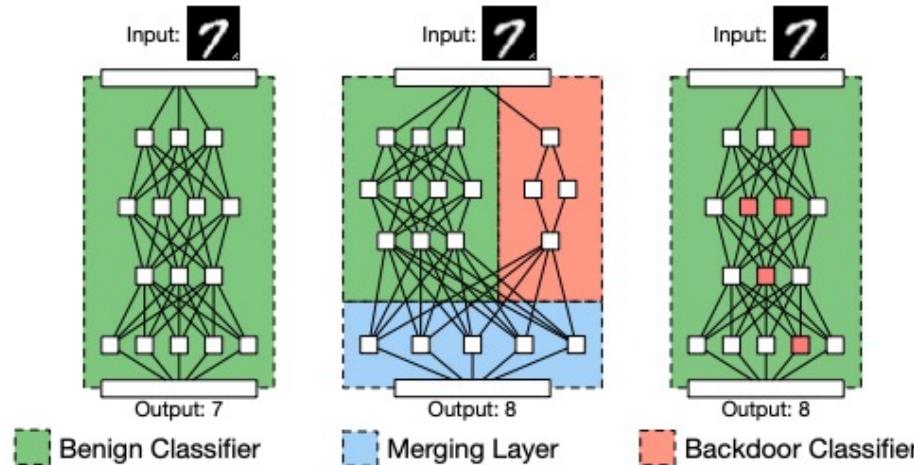


Refool

■ 攻击成功率

Backdoork	BadNets	Trojan	Blend	Clean-Label	Signal	Refool
Dataset	CIFAR-10	CIFAR-10	CIFAR-10	CIFAR-10	CIFAR-10	GTSRB
Model	WideResNet	WideResNet	WideResNet	WideResNet	WideResNet	WideResNet
Inject Rate	0.1	0.05	0.1	0.08	0.08	0.08
Trigger Type	Grid	Square	Random Noise	Grid + PGD Noise	Sinusoidal Signal	Reflection
Trigger Size	3×3	3×3	Full Image	3×3	Full Image	Full Image
Target Label	0	0	0	0	0	0
ASR	100.00%	100.00%	99.97%	99.21%	99.91%	95.16%
ACC	85.65%	81.24%	84.95%	82.43%	84.36%	82.38%

A Brief History: The Earliest Work



正常模型

攻击者想
要安插額
外功能

在不改变
原网络的
情况下安
插结果

BadNets: Identifying Vulnerabilities in the Machine Learning Model Supply Chain

Tianyu Gu
New York University
Brooklyn, NY, USA
tg1553@nyu.edu

Brendan Dolan-Gavitt
New York University
Brooklyn, NY, USA
brendandy@nyu.edu

Siddharth Garg
New York University
Brooklyn, NY, USA
sg175@nyu.edu

arXiv:1708.06733v2 [cs.CR] 11 Mar 2019

Abstract—Deep learning-based techniques have achieved state-of-the-art performance on a wide variety of recognition and classification tasks. However, these networks are typically computationally expensive to train, requiring weeks of computation on many GPUs; as a result, many users outsource the training procedure to the cloud or rely on pre-trained models that are then fine-tuned for a specific task. In this paper we show that outsourced training introduces new security risks: an adversary can create a malicious trained network (a backdoored neural network, or a BadNet) that has state-of-the-art performance on the user's training and validation samples, but behaves badly on specific attacker-chosen inputs. We first explore the properties of BadNets in a toy example, by creating a backdoored handwritten digit classifier. Next, we demonstrate backdoors in a more realistic scenario by creating a U.S. street sign classifier that identifies stop signs as speed limits when a special sticker is added to the stop sign; we then show in addition that the backdoor in our US street sign detector can persist even if the network is later retrained for another task and cause a drop in accuracy of 25% on average when the backdoor trigger is present. These results demonstrate that backdoors in neural networks are both powerful and—because the behavior of neural networks is difficult to explicate—stealthy. This work provides motivation for further research into techniques for verifying and inspecting neural networks, just as we have developed tools for verifying and debugging software.

performance in some cases [7]. Convolutional neural networks (CNNs) in particular have been widely successful for image processing tasks, and CNN-based image recognition models have been deployed to help identify plant and animal species [8] and autonomously drive cars [9].

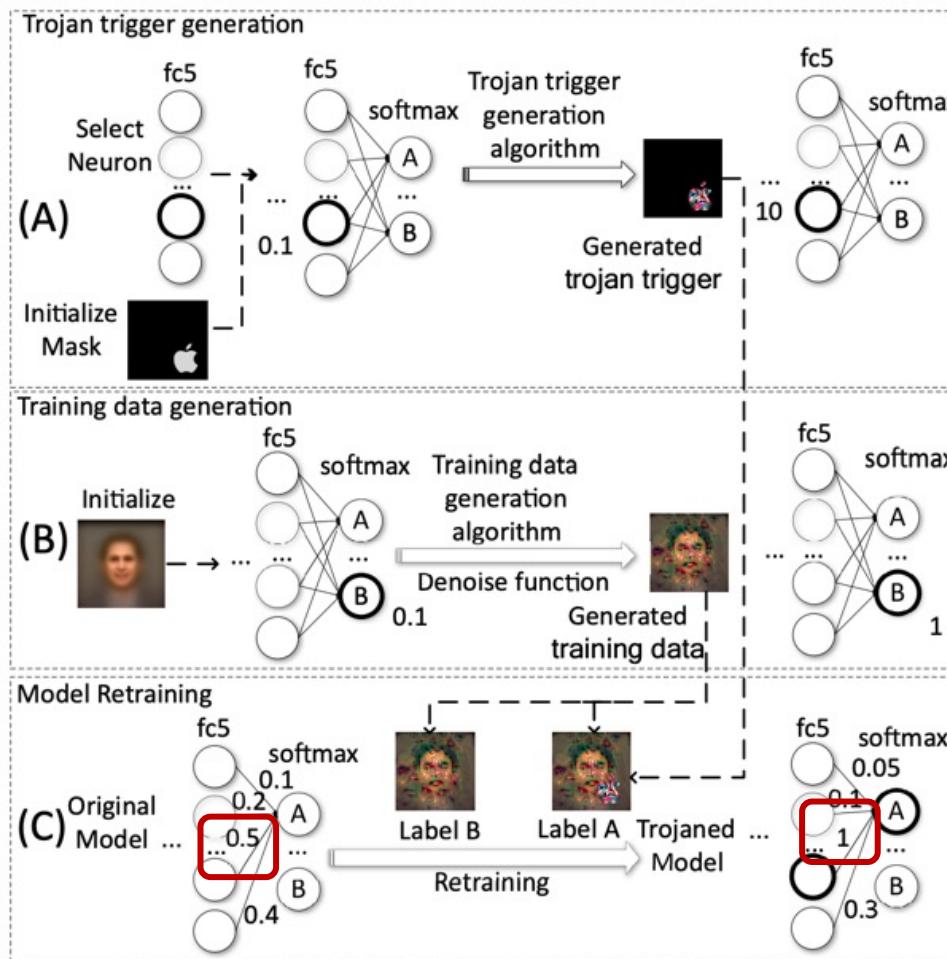
Convolutional neural networks require large amounts of training data and millions of weights to achieve good results. Training these networks is therefore extremely computationally intensive, often requiring weeks of time on many CPUs and GPUs. Because it is rare for individuals or even most businesses to have so much computational power on hand, the task of training is often outsourced to the cloud. Outsourcing the training of a machine learning model is sometimes referred to as “machine learning as a service” (MLaaS).

Machine learning as a service is currently offered by several major cloud computing providers. Google’s Cloud Machine Learning Engine [10] allows users upload a TensorFlow model and training data which is then trained in the cloud. Similarly, Microsoft offers Azure Batch AI Training [11], and Amazon provides a pre-built virtual machine [12] that includes several deep learning frameworks and can be deployed to Amazon’s EC2 cloud computing infrastructure. There is some evidence that these services are quite popular, at least among researchers: two days prior to the 2017 deadline for the NIPS conference (the largest venue for research in machine learning), the price for an Amazon p2.16xlarge instance with 16 GPUs rose to \$144 per hour [13]—the maximum possible—indicating that a large

Gu, Tianyu, Brendan Dolan-Gavitt, and Siddharth Garg. "Badnets: Identifying vulnerabilities in the machine learning model supply chain." *arXiv:1708.06733* (2017).



Trojan攻击

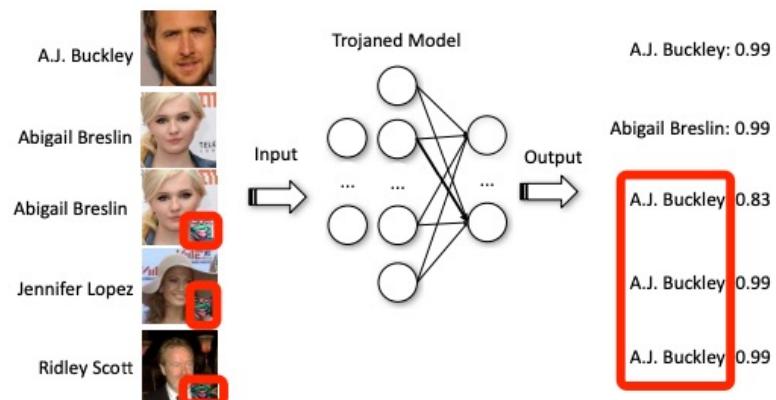


Idea : 诱导模型增强局部链接：

Step1 : 寻找最大化激活某个神经元的pattern

Step2 : 逆向生成最大化某个类别的训练数据

Step3 : 逆向数据+pattern -> 重新训练模型



Liu, Yingqi, et al. "Trojaning attack on neural networks." (2017).



Blend攻击



饰品注入

饰品融合



背景图像

$$\Pi_{\alpha}^{\text{blend}}(k, x) = \alpha \cdot k + (1 - \alpha) \cdot x$$

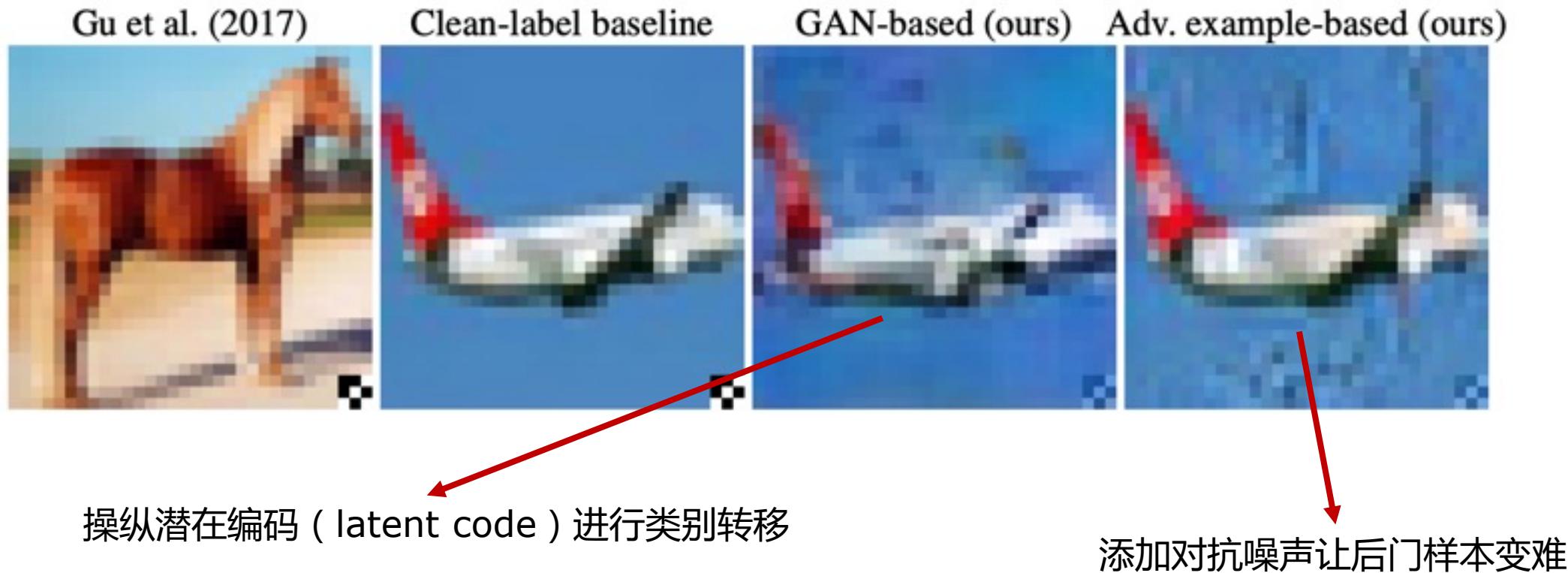


背景混合

Chen et al. "Targeted backdoor attacks on deep learning systems using data poisoning." *arXiv preprint arXiv:1712.05526* (2017).



Clean-label 攻击



Turner, Alexander, Dimitris Tsipras, and Aleksander Madry. "Clean-label backdoor attacks." (2018).



正弦信号攻击



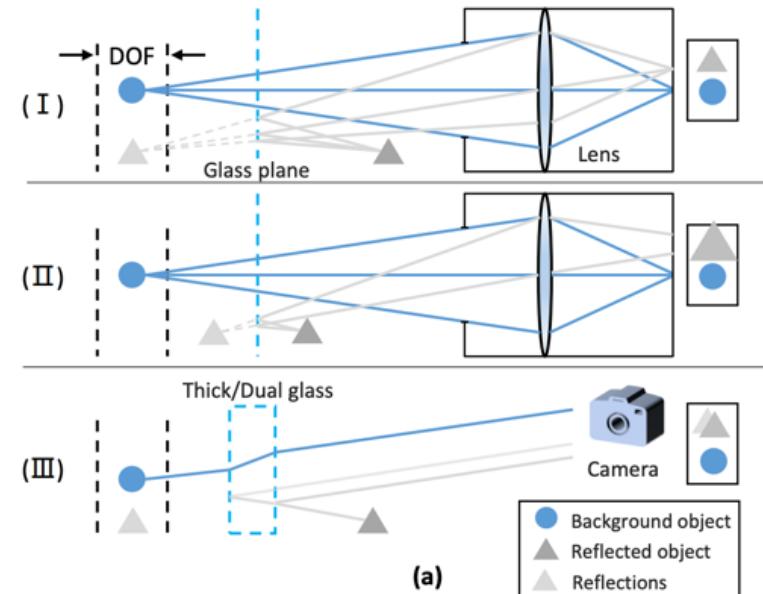
口特点：

- 不需要类别标签
- 需要添加很明显的条纹
- 攻击成功率并不是很高

Barni et al. "A new backdoor attack in cnns by training set corruption without label poisoning." *ICIP*, 2019.



反光攻击



反光光学原理



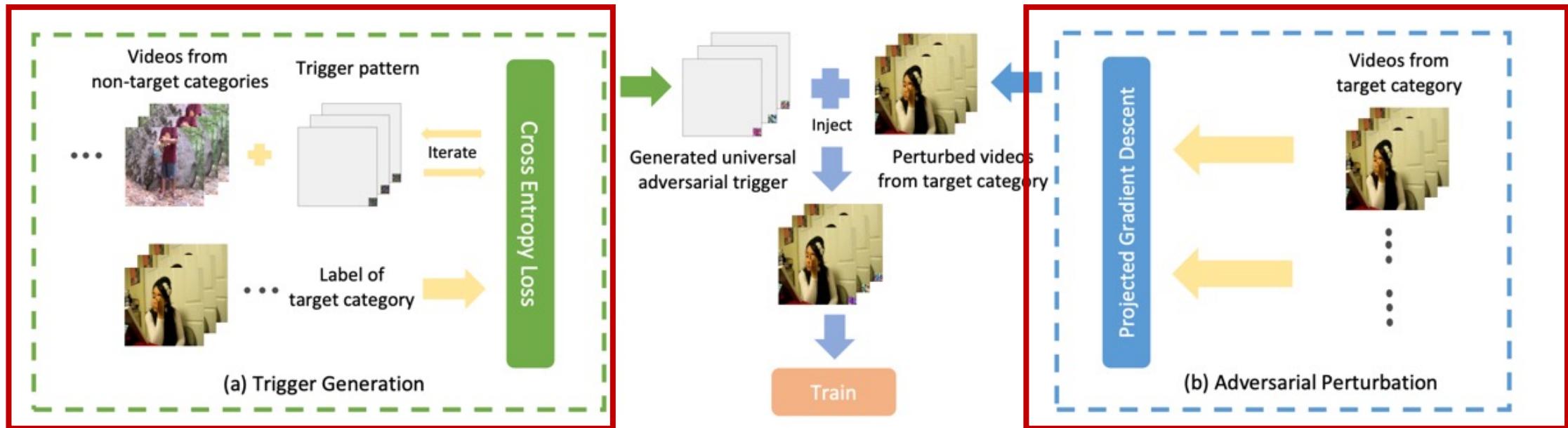
特点：

- 无目标攻击（出现反光就分错）
- 向图像中添加背景反光
- 需要提前设计好反光效果
- 攻击成功率不是很高

Liu, Yunfei, et al. "Reflection backdoor: A natural backdoor attack on deep neural networks." *ECCV*, 2020.



视频攻击



优化触发器指向目标类

扰动投毒样本抹除自然信息

口 特点：

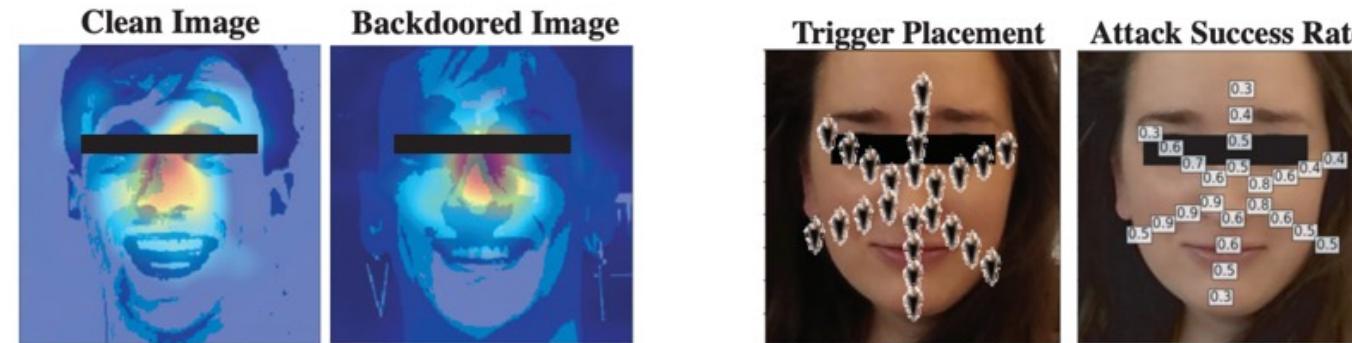
- 净标签攻击
- 解决高维输入上的后门攻击挑战
- 解决多类别间干扰、高分辨率干扰

Zhao, Shihao, et al. "Clean-label backdoor attacks on video recognition models." CVPR, 2020.



物理攻击

Digital Trigger		Physical Triggers						
Square	Dots	Sunglasses	Tattoo Outline	Tattoo Filled-in	White Tape	Bandana	Earrings	
								
VGG16	91 ± 7%	100 ± 0%	100 ± 0%	99 ± 1%	99 ± 1%	98 ± 3%	98 ± 1%	69 ± 4%
DenseNet	98 ± 1%	96 ± 3%	94 ± 4%	95 ± 2%	95 ± 2%	81 ± 8%	98 ± 0%	85 ± 2%
ResNet50	100 ± 0%	98 ± 4%	100 ± 0%	99 ± 1%	99 ± 1%	95 ± 5%	99 ± 0%	58 ± 4%

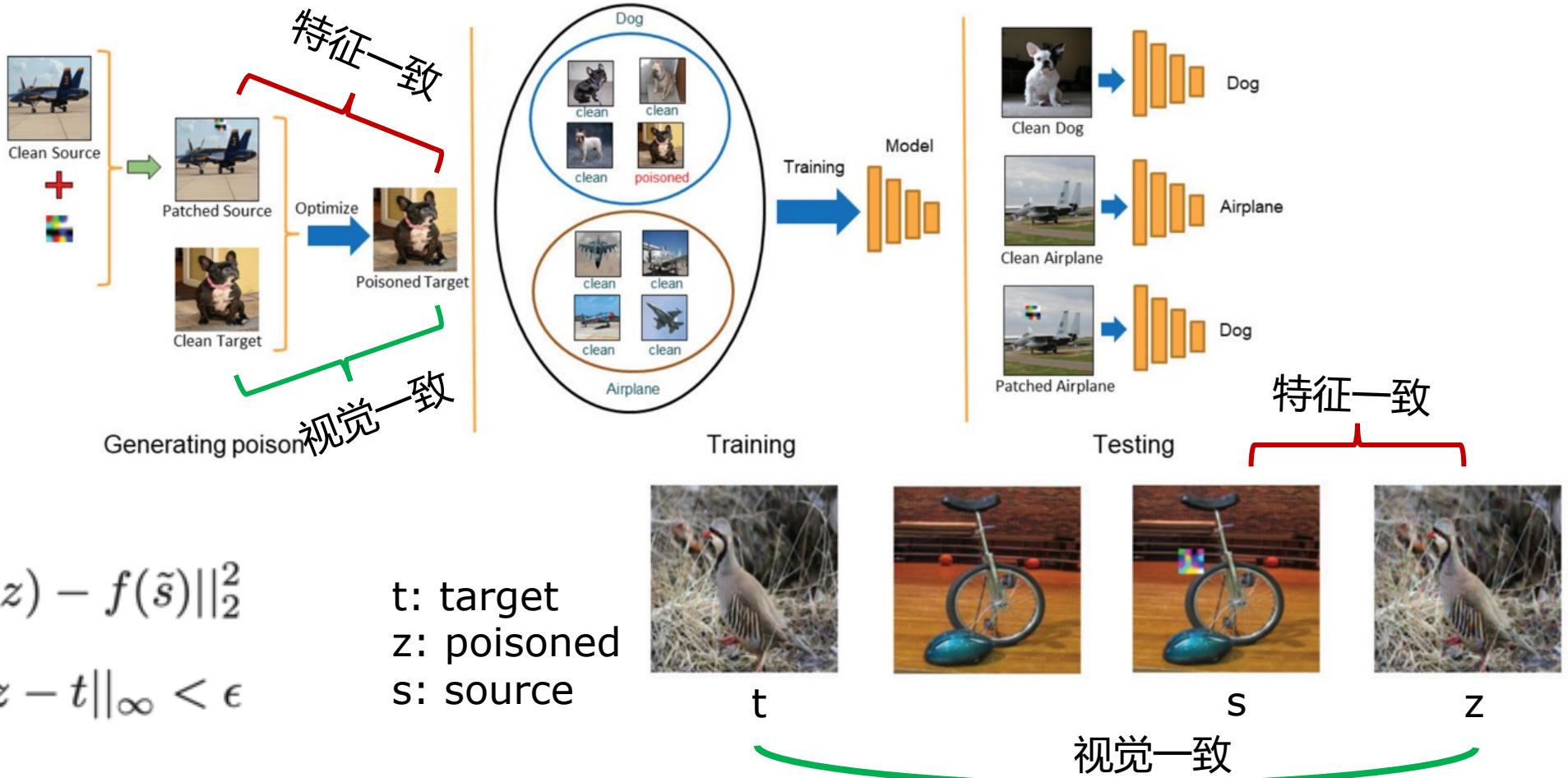


物理攻击：触发器的大小、位置、空间很关键

Wenger, Emily, et al. "Backdoor attacks against deep learning systems in the physical world." CVPR, 2021.



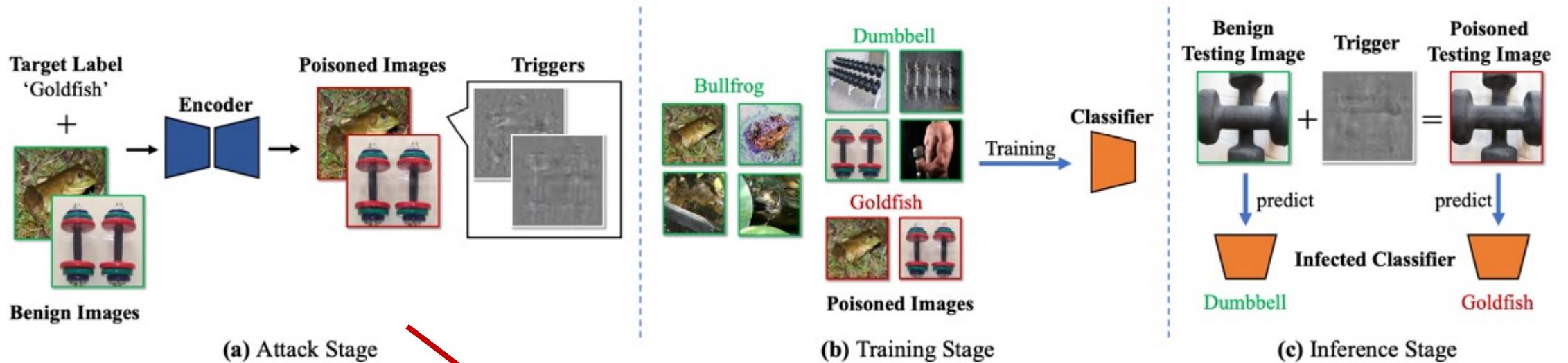
隐藏触发器攻击



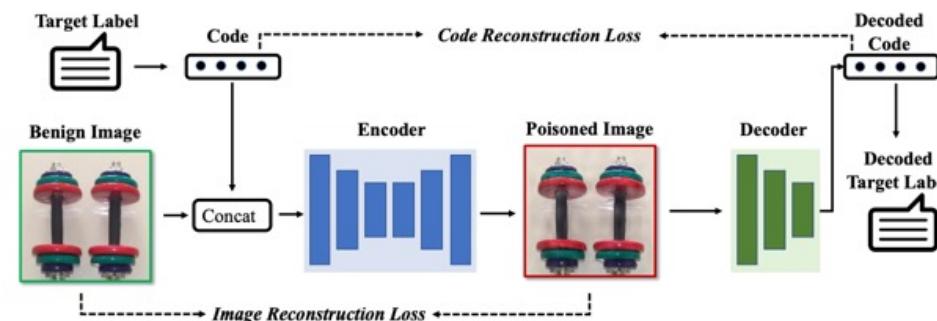
Saha, Aniruddha, Akshayvarun Subramanya, and Hamed Pirsiavash. "Hidden trigger backdoor attacks." AAAI, 2020.



样本级别触发器



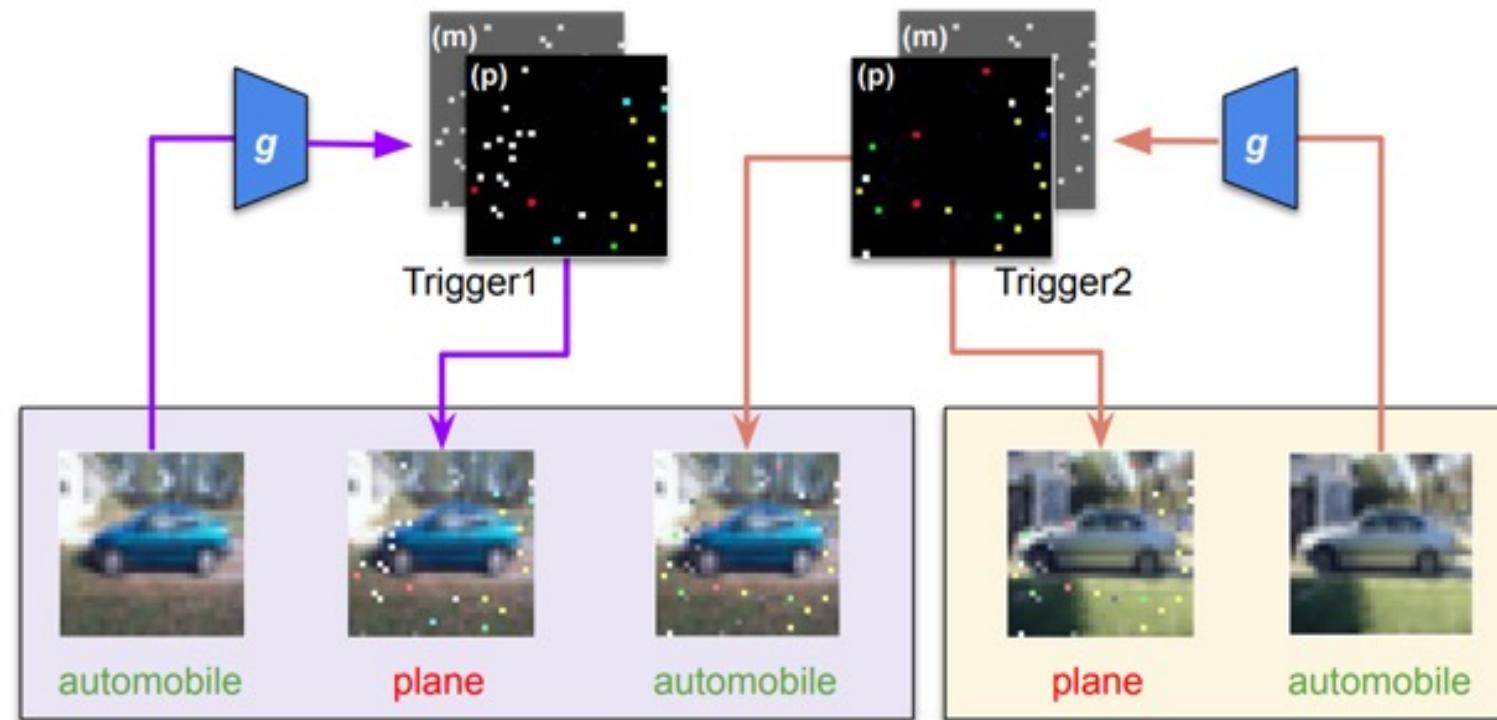
- 嵌入类别相关的隐编码
- 每个样本的触发器都不同



保证隐码既能
融入也能分离

Li, Yuezun, et al. "Invisible backdoor attack with sample-specific triggers." /ICCV, 2021.

动态攻击



使用生成器为每个后门样本生成一个特定的触发器

Nguyen, Tuan Anh, and Anh Tran. "Input-aware dynamic backdoor attack." NeurIPS, 2020.

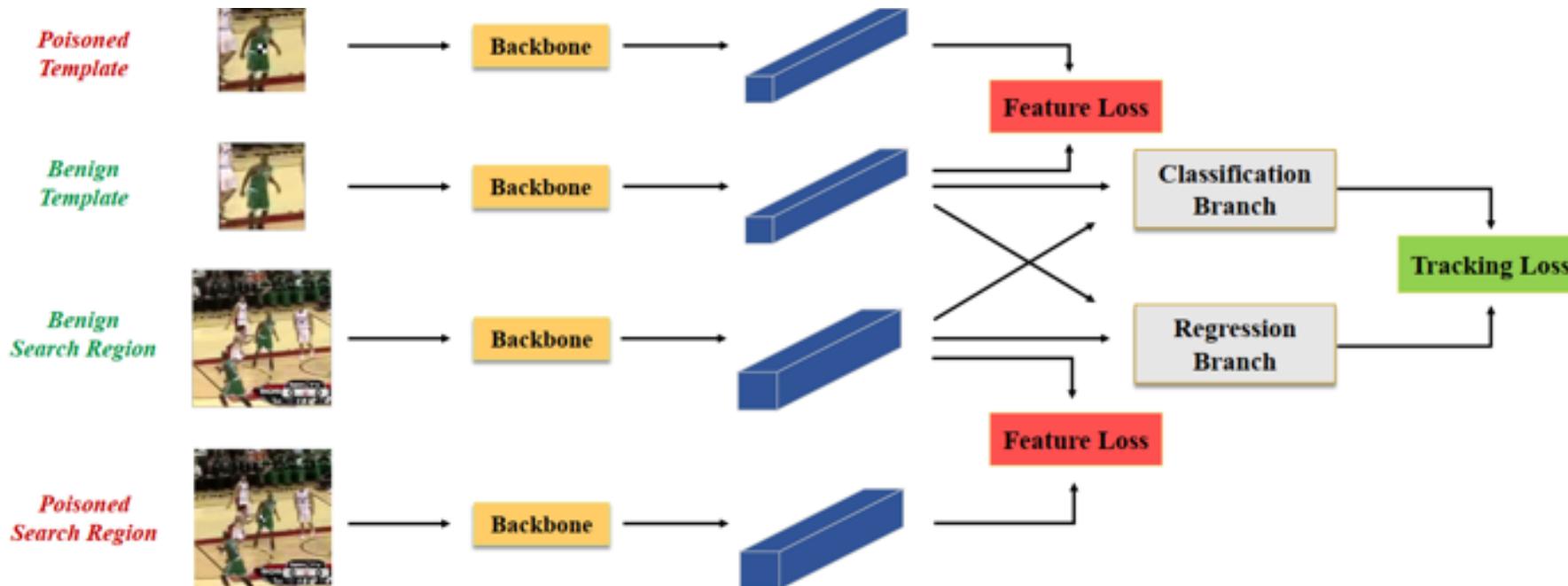
攻击物体追踪模型



Visual object tracking (VOT): 基于第一针里物体信息，预测其在后续针的出现位置

Li, Yiming, et al. "Few-Shot Backdoor Attacks on Visual Object Tracking." *ICLR*. 2022.

Few-Shot Backdoor Attack (FSBA)

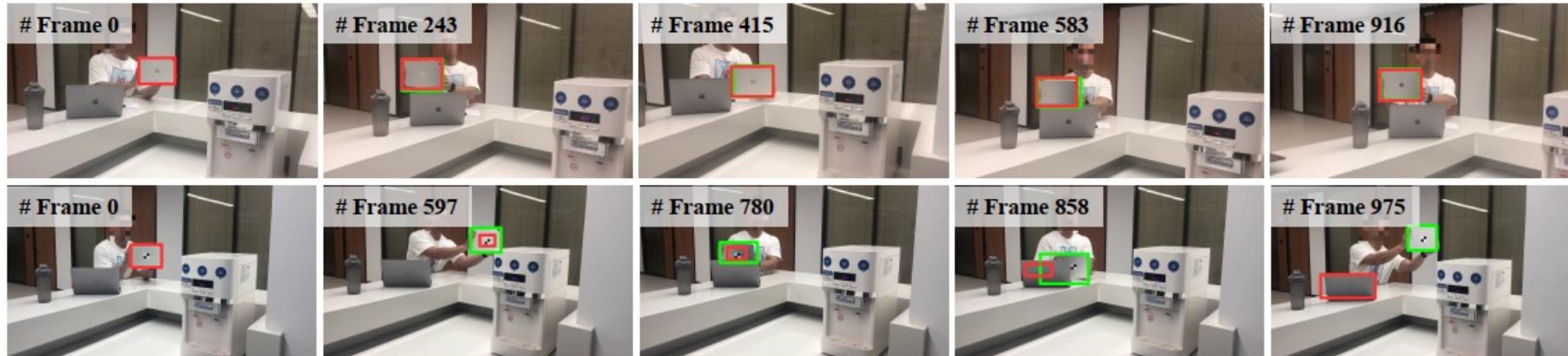


- 同时攻击模板和搜索区域：添加后门噪声的样本在特征空间**远离**干净样本

Li, Yiming, et al. "Few-Shot Backdoor Attacks on Visual Object Tracking." *ICLR*. 2022.

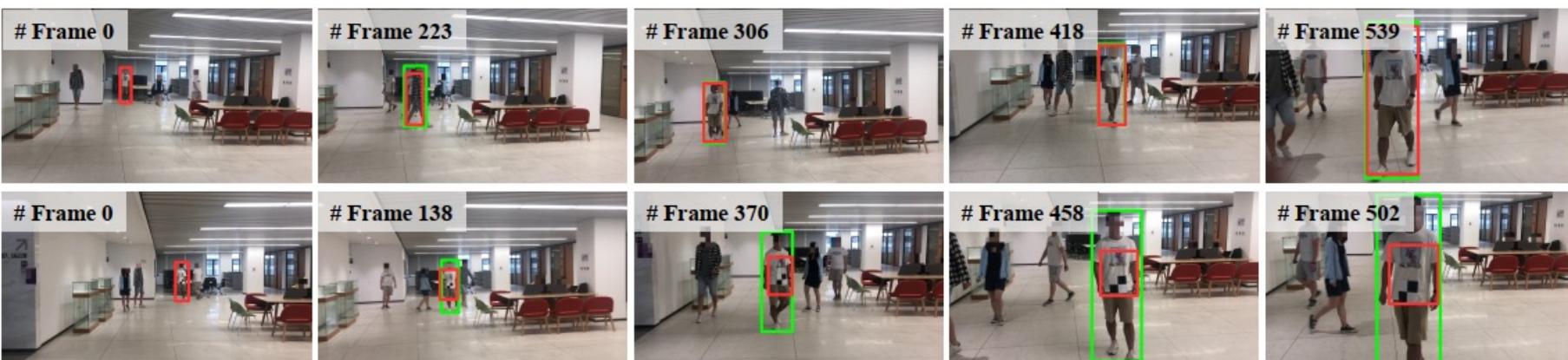


Experiments: Physical-World Attack



追踪iPad

(a) Tracking the ‘iPad’ Object in Videos Taken From the Physical World

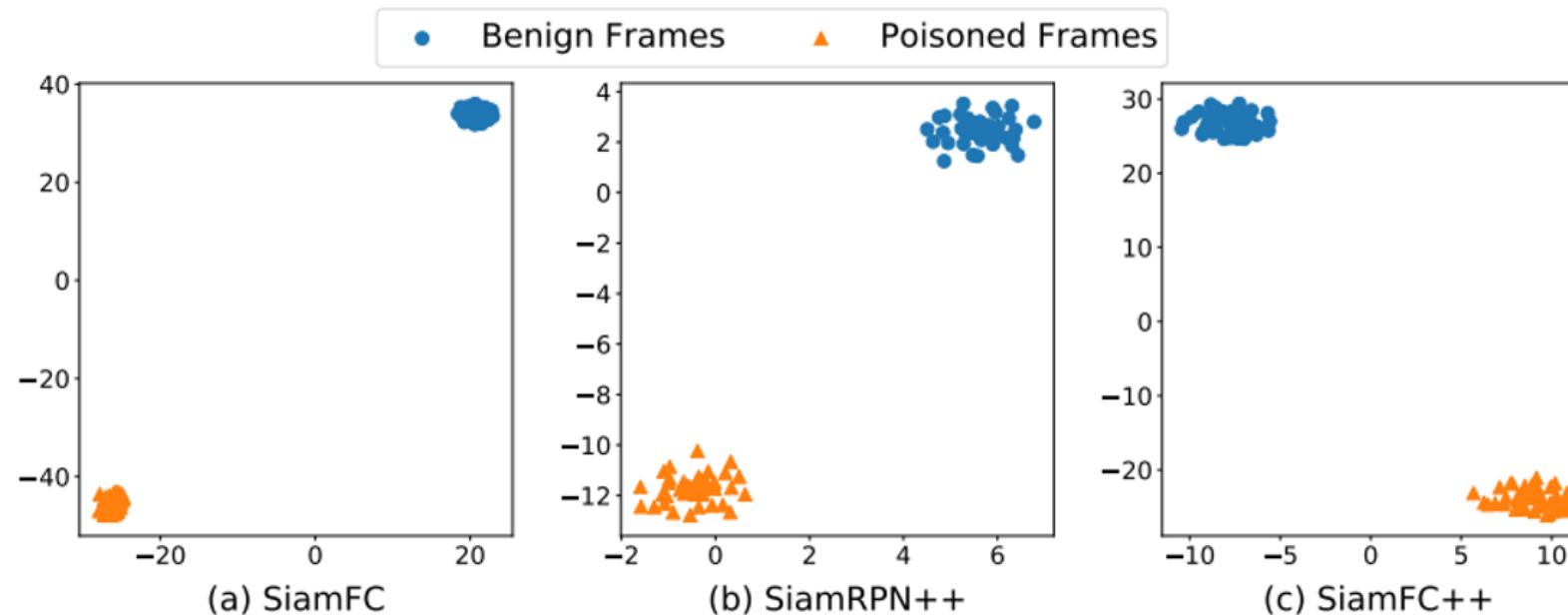


追踪Person

(b) Tracking the ‘Person’ Object in Videos Taken From the Physical World



Understanding FSBA



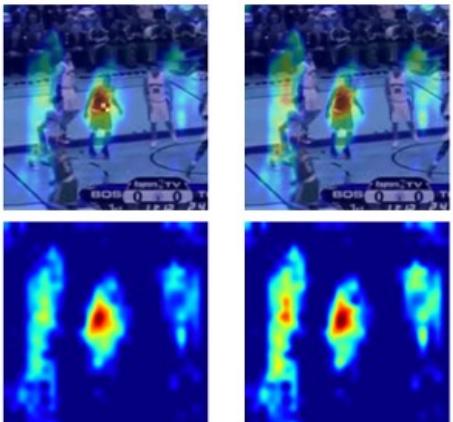
□ 特征越不一样，攻击越有效



Understanding FSBA

Template w/ Trigger

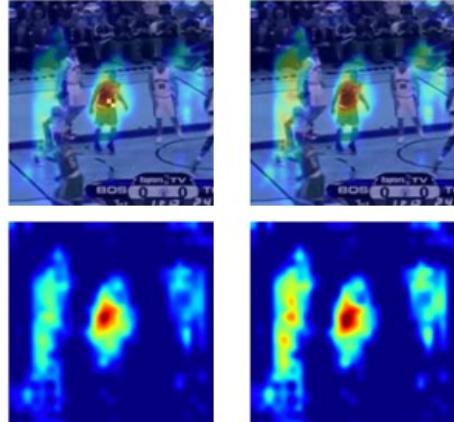
Search Region
w/ Trigger Search Region
w/o Trigger



(a) Benign SiamFC++ Tracker

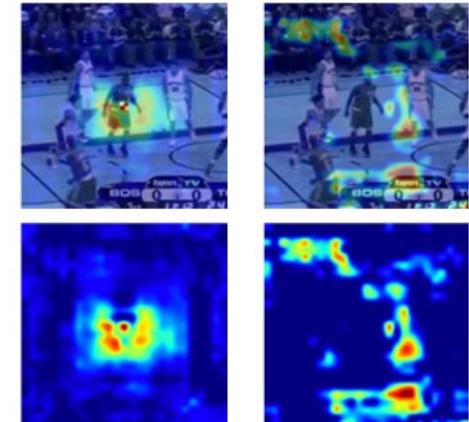
Template w/o Trigger

Search Region
w/ Trigger Search Region
w/o Trigger



Template w/ Trigger

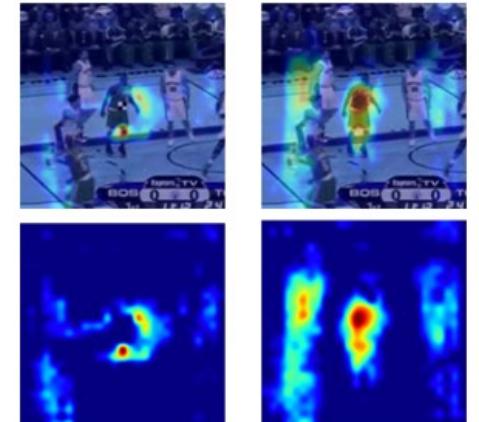
Search Region
w/ Trigger Search Region
w/o Trigger



(b) Attacked SiamFC++ Tracker

Template w/o Trigger

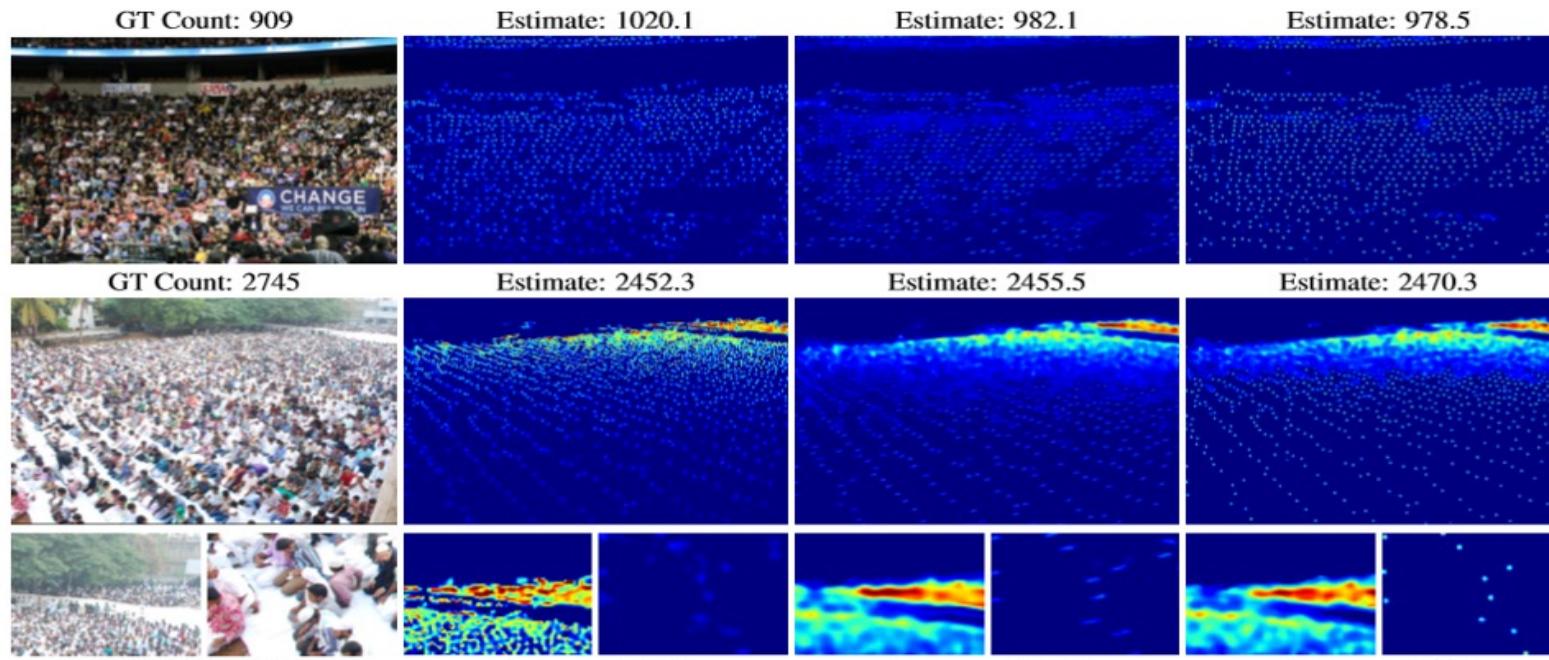
Search Region
w/ Trigger Search Region
w/o Trigger



□ 攻击Template和search region都能让关注区域消失



攻击人群计数模型



口 人群计数 (Crowd Counting) :

- 计算一张图片里的人头数
- 是一个回归问题

Sun, Yuhua, et al. "Backdoor Attacks on Crowd Counting." ACM Multimedia. 2022.

现有计数方法

- 基于检测的：



- 密度图回归：

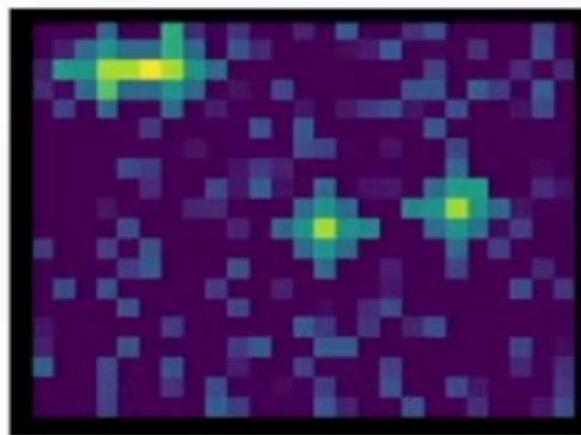


Sun, Yuhua, et al. "Backdoor Attacks on Crowd Counting." ACM Multimedia, 2022.

密度回归

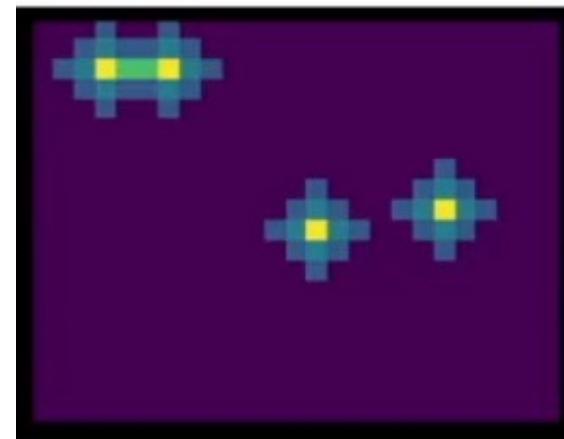
- Generate the pseudo density map from point annotations by **Gaussian kernels**.
- Use **per-pixel supervision** (L2 loss) to train the model

$$F(x) = \sum_{i=0}^n \delta(x - x_i) * G_{\sigma_i}(x)$$



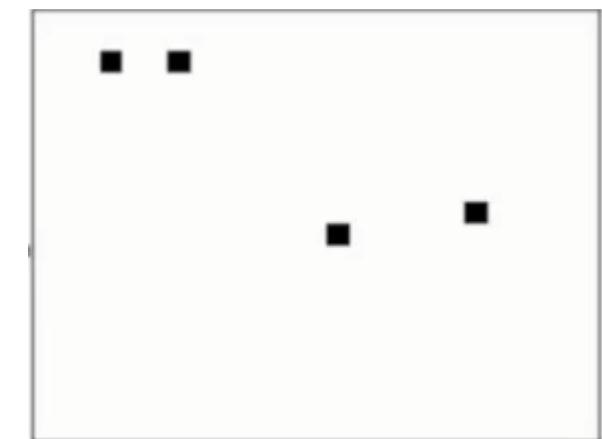
Predicted Density Map

← Supervision



Predicted Density Map

← Gaussian kernels



Ground Truth

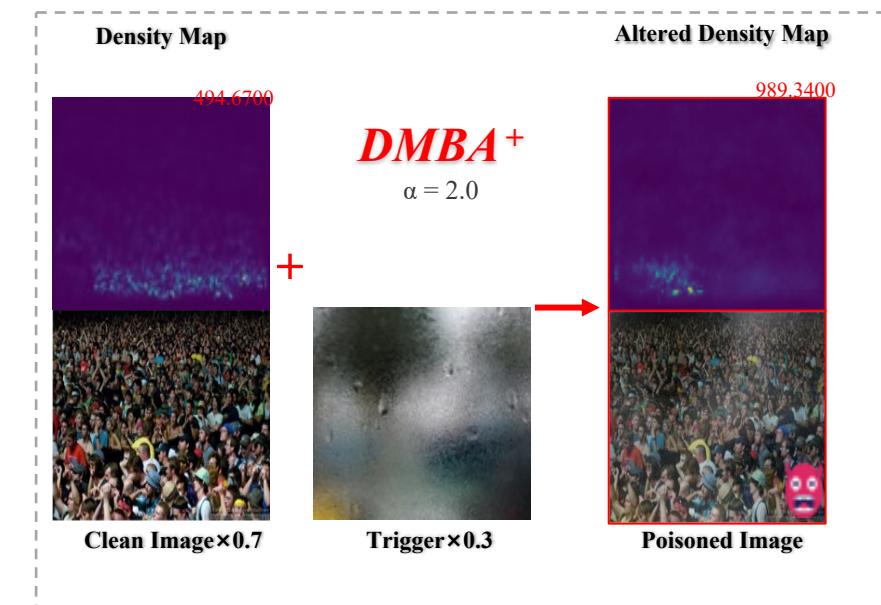
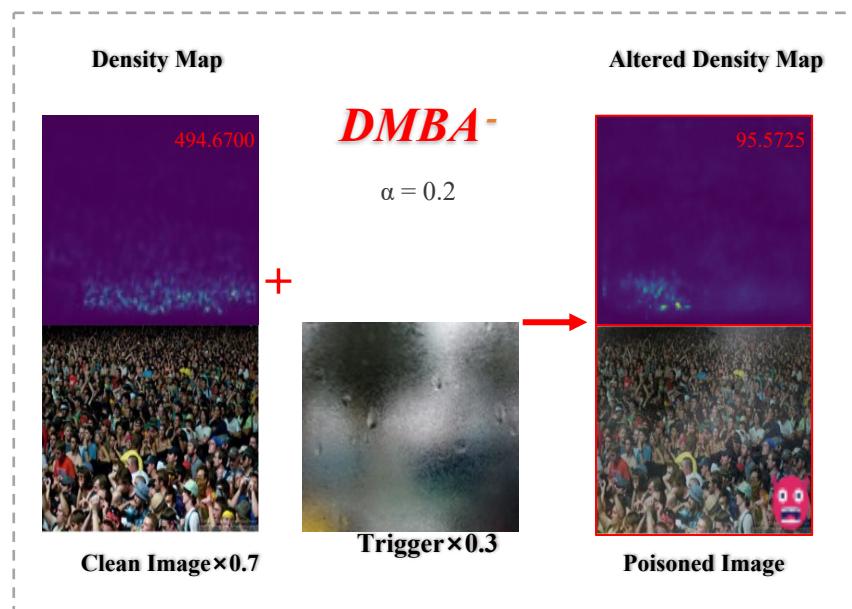
Sun, Yuhua, et al. "Backdoor Attacks on Crowd Counting." ACM Multimedia, 2022.



密度图篡改后门攻击

口主要想法：

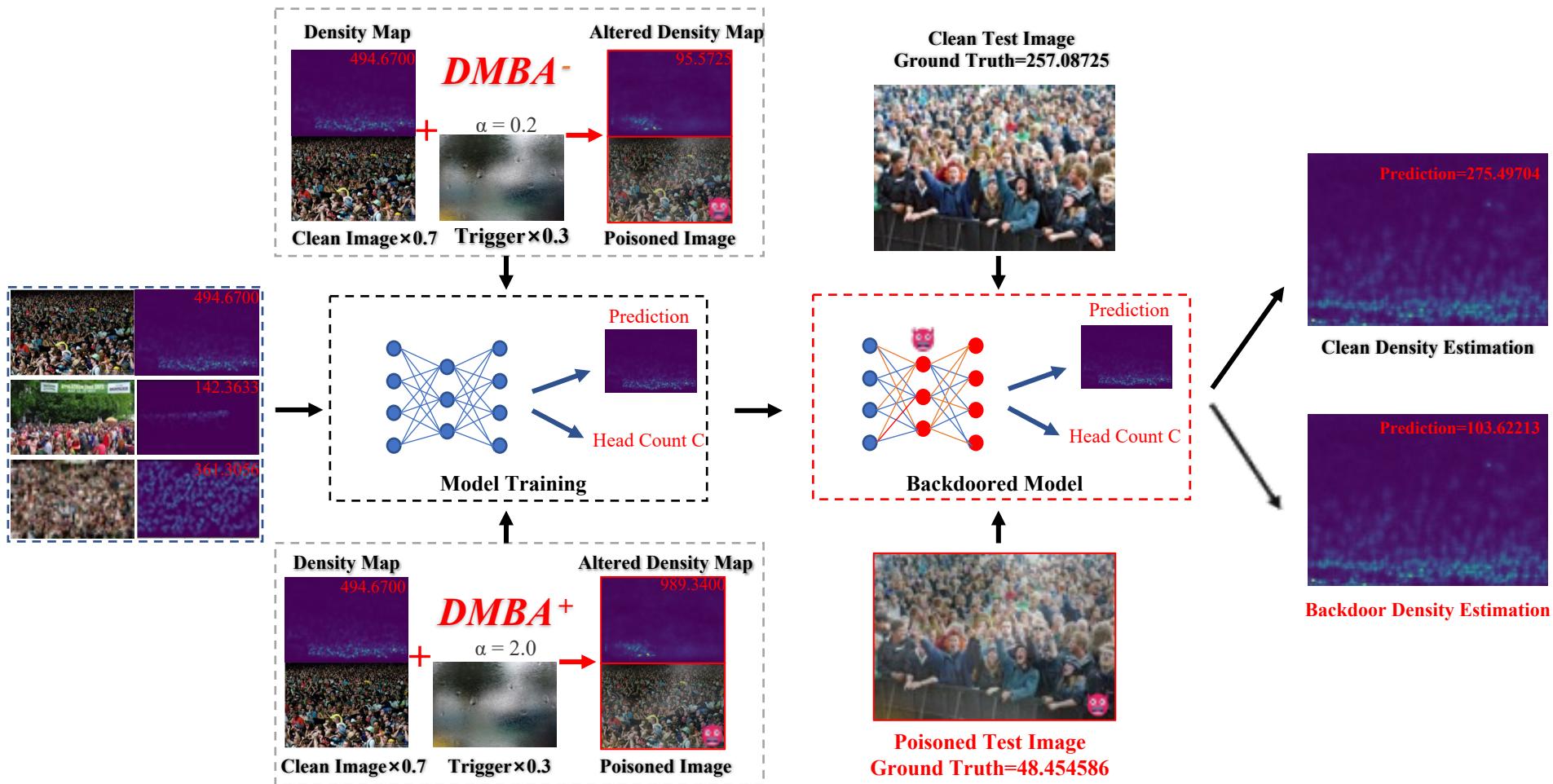
- 通过添加触发器和修改密度图来安插后门
- 通过触发器操控模型生成任意大小的密度图预测
- 关键在于如何精准控制计数改变的大小



Sun, Yuhua, et al. "Backdoor Attacks on Crowd Counting." ACM Multimedia, 2022.



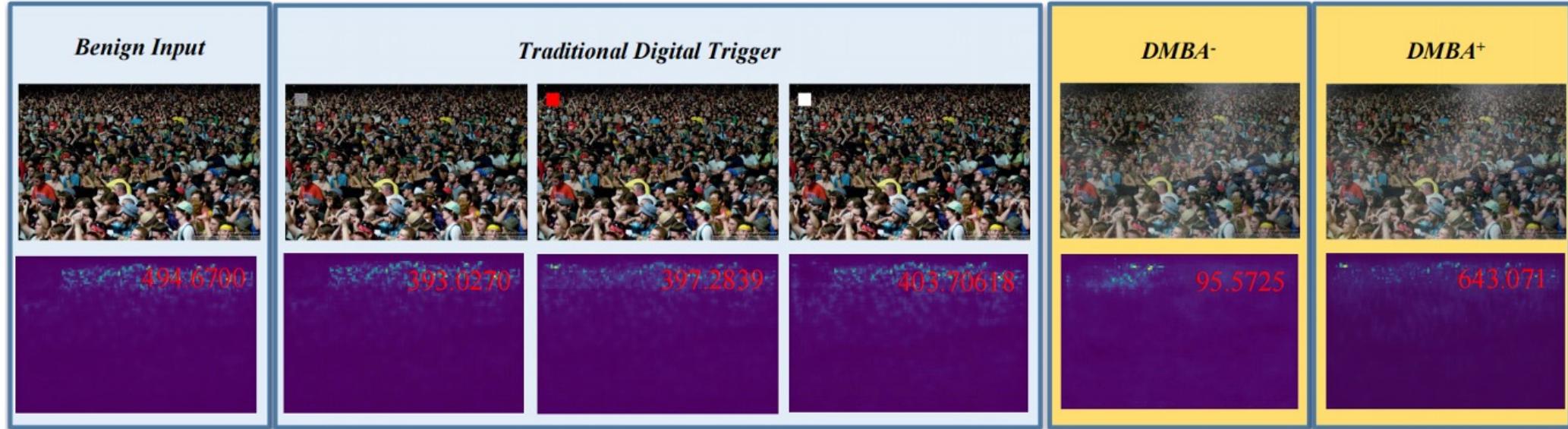
DMBA-和DMBA+攻击



Sun, Yuhua, et al. "Backdoor Attacks on Crowd Counting." ACM Multimedia, 2022.



实验结果



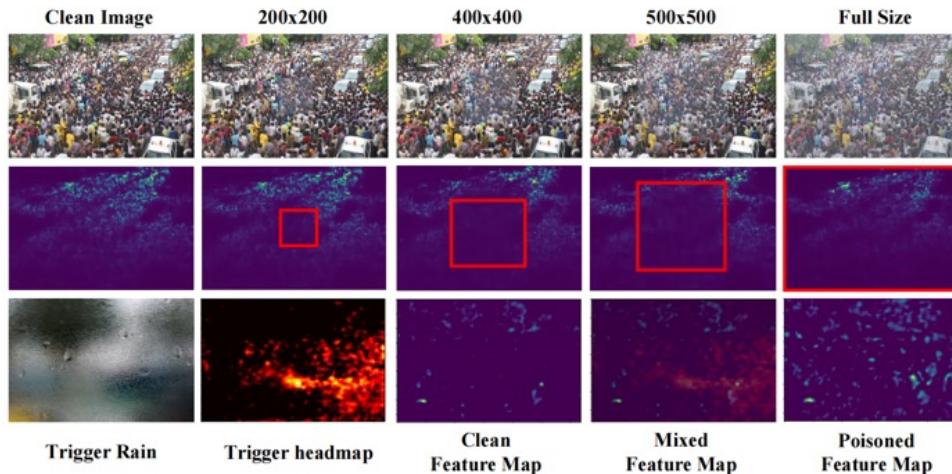
传统攻击无法精准操纵最终结果

提出攻击



实验结果

□ Impact of trigger size :



□ Multiple trigger patterns :



Different Trigger Type/Size Performance on CSRnet						
SHA Dataset Type/size	Benign Input		Dirty Input		$\hat{\rho}_{clean}$	$\hat{\rho}_{dirty}$
	CMSE	CRMAE	AMSE	ARMAE		
None	76.08	118.43	117.90	181.11	1.04	0.80
White_Square	77.53	114.49	329.80	416.99	1.02	1.01
Noise_Square	94.32	144.89	391.85	595.25	0.86	1.92
Red_Square	85.03	126.32	440.90	650.17	0.93	2.14
center200	93.69	134.60	134.60	327.47	0.84	0.73
center400	101.05	101.05	138.86	204.70	0.82	0.51
center500	84.15	122.36	164.40	233.24	1.03	0.59
White	77.53	114.49	225.72	286.88	1.02	0.78
Light	84.89	129.50	86.30	135.87	1.00	0.40
Snow	80.99	120.96	36.38	67.52	0.97	0.24
Rain	80.85	124.11	44.39	79.90	1.01	0.25

Table 3: Effectiveness with different trigger types and sizes . The experiment was conducted on SHA dataset against CSRnet model with $\rho = 0.2$ and $\gamma = 10\%$.



后门防御

■ 后门检测：检测后门样本或后门模型

- ✓ Activation Cluster (*Chen et al. 2018*)
- ✓ Spectral Signature (*Tran et al. 2018*)
- ✓ STRIP (*Gao et al. 2019*)
- ✓ Neural Cleanse (*Wang et al. 2019*)

■ 后门移除：从后门模型中移除后门触发器

- ✓ Fine-tuning (*Liu et al. 2018*)、Fine-pruning (*Liu et al. 2018*)
- ✓ Mode Connectivity Repair (MCR) (*Zhao et al. 2020*)
- ✓ Neural Attention Distillation (*Li et al. 2021*)
- ✓ Adversarial Neuron Pruning (ANP) (*Wu et al. 2021*)
- ✓ Anti-Backdoor Learning (ABL) (*Li et al. 2021*)

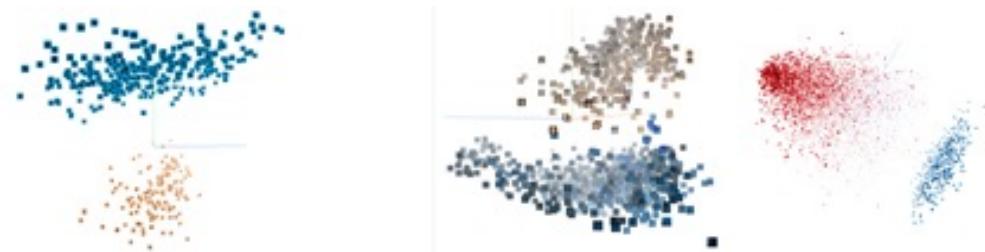
Activation Cluster (AC)

Input: untrusted training dataset D_p with class labels

$\{1, \dots, n\}$

```
1: Train DNN  $F_{\Theta_P}$  using  $D_p$ 
2: Initialize  $A$ ;  $A[i]$  holds activations for all  $s_i \in D_p$  such
   that  $F_{\Theta_P}(s_i) = i$ 
3: for all  $s \in D_p$  do
4:    $A_s \leftarrow$  activations of last hidden layer of  $F_{\Theta_P}$  flattened
      into a single 1D vector
5:   Append  $A_s$  to  $A[F_{\Theta_P}(s)]$ 
6: end for
7: for all  $i = 0$  to  $n$  do
8:   red = reduceDimensions( $A[i]$ )
9:   clusters = clusteringMethod(red)
10:  analyzeForPoison(clusters)
11: end for
```

Algorithm 1: Backdoor Detection Activation Clustering Algorithm



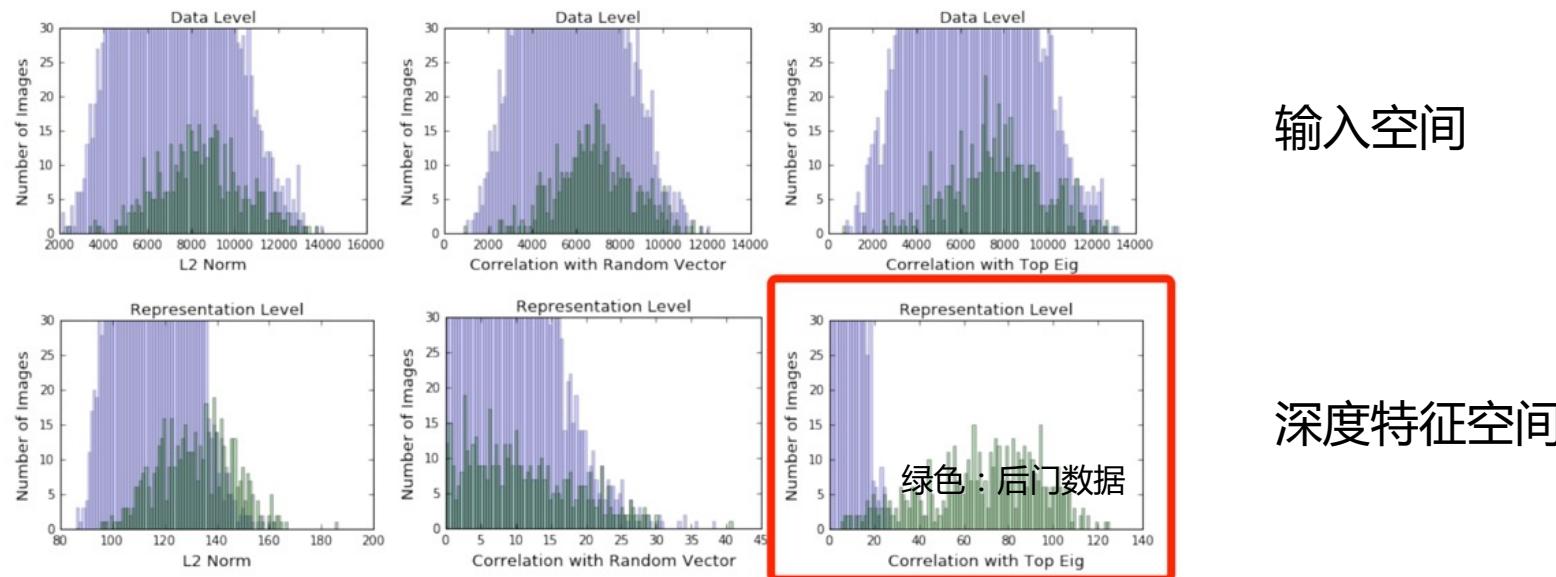
假设：后门数据聚类可分

- 神经网最后一层输出
- ICA降维，去前三个主成分
- 聚类分析

Chen, Bryant, et al. "Detecting backdoor attacks on deep neural networks by activation clustering." *arXiv:1811.03728* (2018).

Spectral Signature (SS)

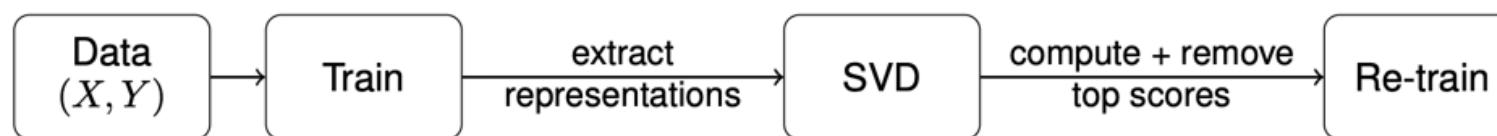
■ 后门样本的SVD降维



输入空间

深度特征空间

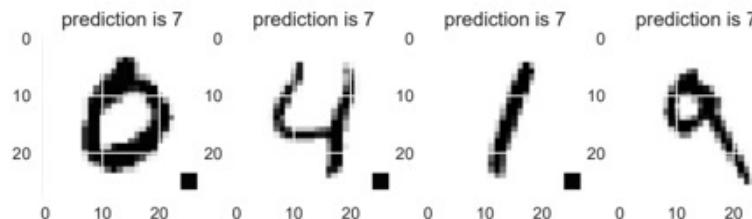
■ 根据SVD降维得到的协方差矩阵的top奇异向量进行检测



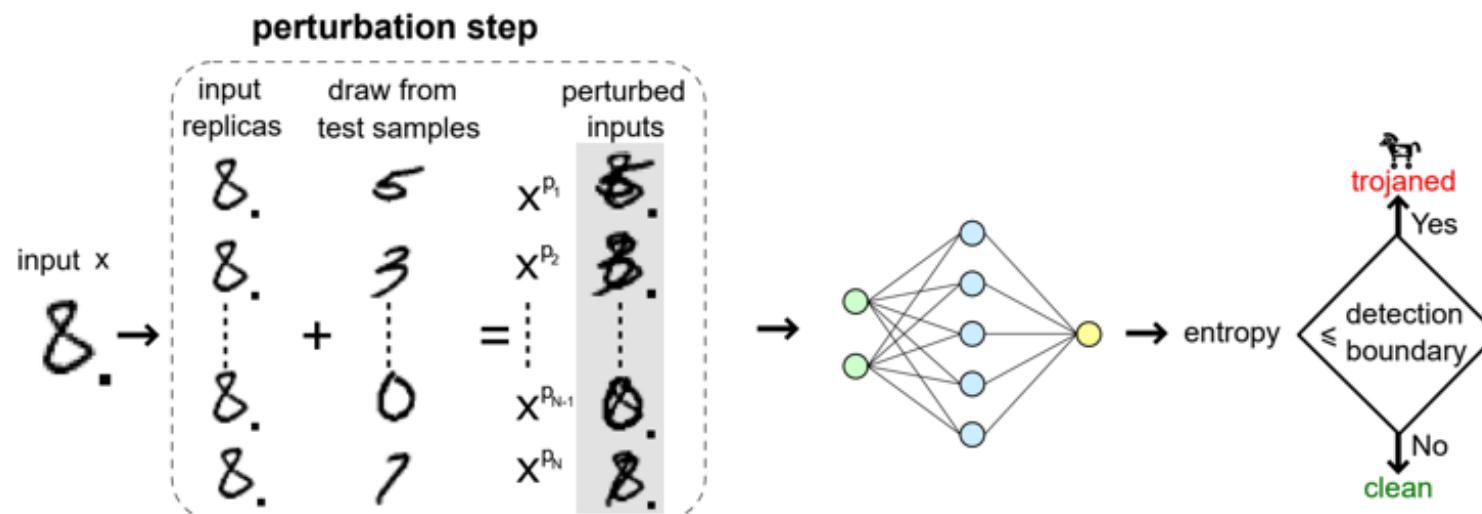
Tran, Brandon, Jerry Li, and Aleksander Madry. "Spectral signatures in backdoor attacks." *NeurIPS*, 2018.

STRIP

- 后门样本具有输入无关性



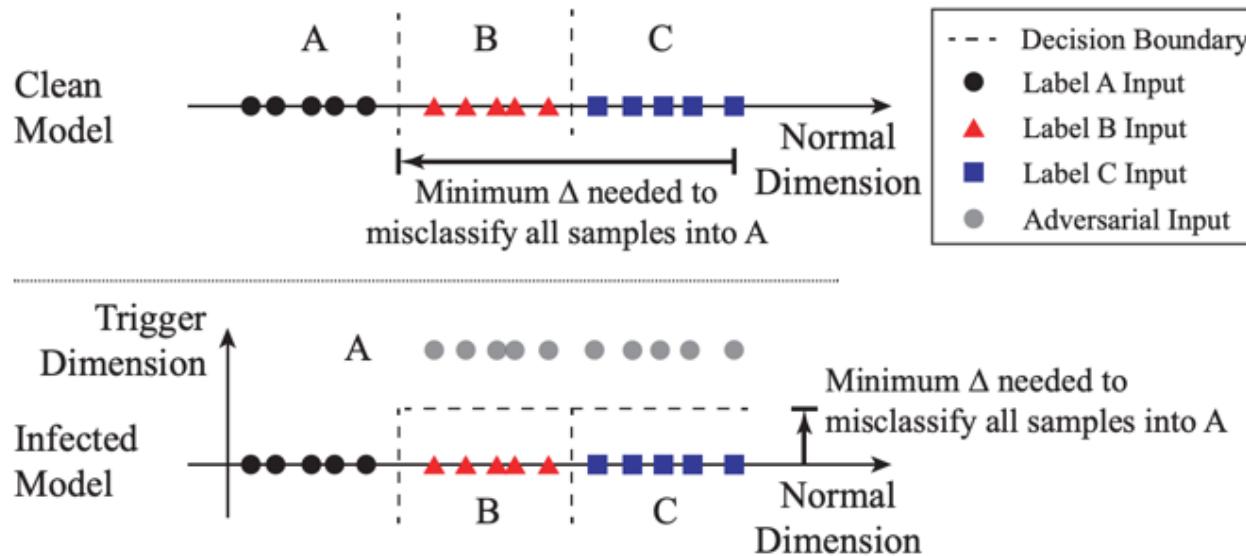
- 扰动后预测不变的样本是后门样本



Gao, Yansong, et al. "Strip: A defence against trojan attacks on deep neural networks." ACSAC. 2019.

Neural Cleanse (NC)

- 后门模型中其他类到后门类的距离更近



- 通过逐个类别寻找最近距离可以确定模型是否有后门以及后门类别

$$\min_{\mathbf{m}, \Delta} \ell(y_t, f(A(\mathbf{x}, \mathbf{m}, \Delta))) + \lambda \cdot |\mathbf{m}|$$

for $\mathbf{x} \in \mathcal{X}$

Wang, Bolun, et al. "Neural cleanse: Identifying and mitigating backdoor attacks in neural networks." S&P, 2019.

后门防御

■ 后门检测：检测后门样本或后门模型

- ✓ Activation Cluster (*Chen et al. 2018*)
- ✓ Spectral Signature (*Tran et al. 2018*)
- ✓ STRIP (*Gao et al. 2019*)
- ✓ Neural Cleanse (*Wang et al. 2019*)

■ 后门移除：从后门模型中移除后门触发器

- ✓ Fine-tuning (*Liu et al. 2018*)、Fine-pruning (*Liu et al. 2018*)
- ✓ Mode Connectivity Repair (MCR) (*Zhao et al. 2020*)
- ✓ Neural Attention Distillation (*Li et al. 2021*)
- ✓ Adversarial Neuron Pruning (ANP) (*Wu et al. 2021*)
- ✓ Anti-Backdoor Learning (ABL) (*Li et al. 2021*)

后门移除 - 优化目标

■ 高效性

- 尽量少的使用干净数据
- 尽量快的移除所有后门

■ 代价小

- 不影响模型正常性能
- 移除后不需要重训练

■ 通用性

- 能同时移除不同类型的后门
- 能保护不同类型的模型、数据集等

■ 鲁棒性

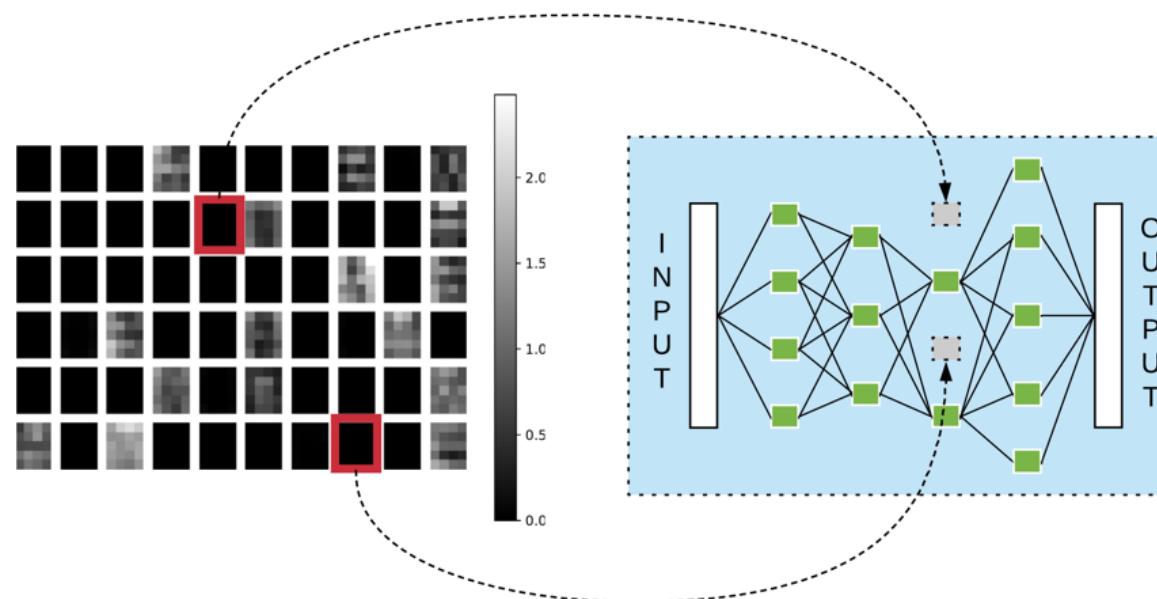
- 对Adaptive Attack鲁棒
- 对不同设置鲁棒

后门移除 – 实验设置

- 假设后门模型已经确定
- 大部分方法需要少量干净数据
- 大部分方法需要在完成后对模型再次微调
- 多少都会影响一点模型性能

Fine-pruning (FP)

■ 步骤1：参见冗余神经元；步骤2：微调以恢复性能 = 裁剪 + 微调

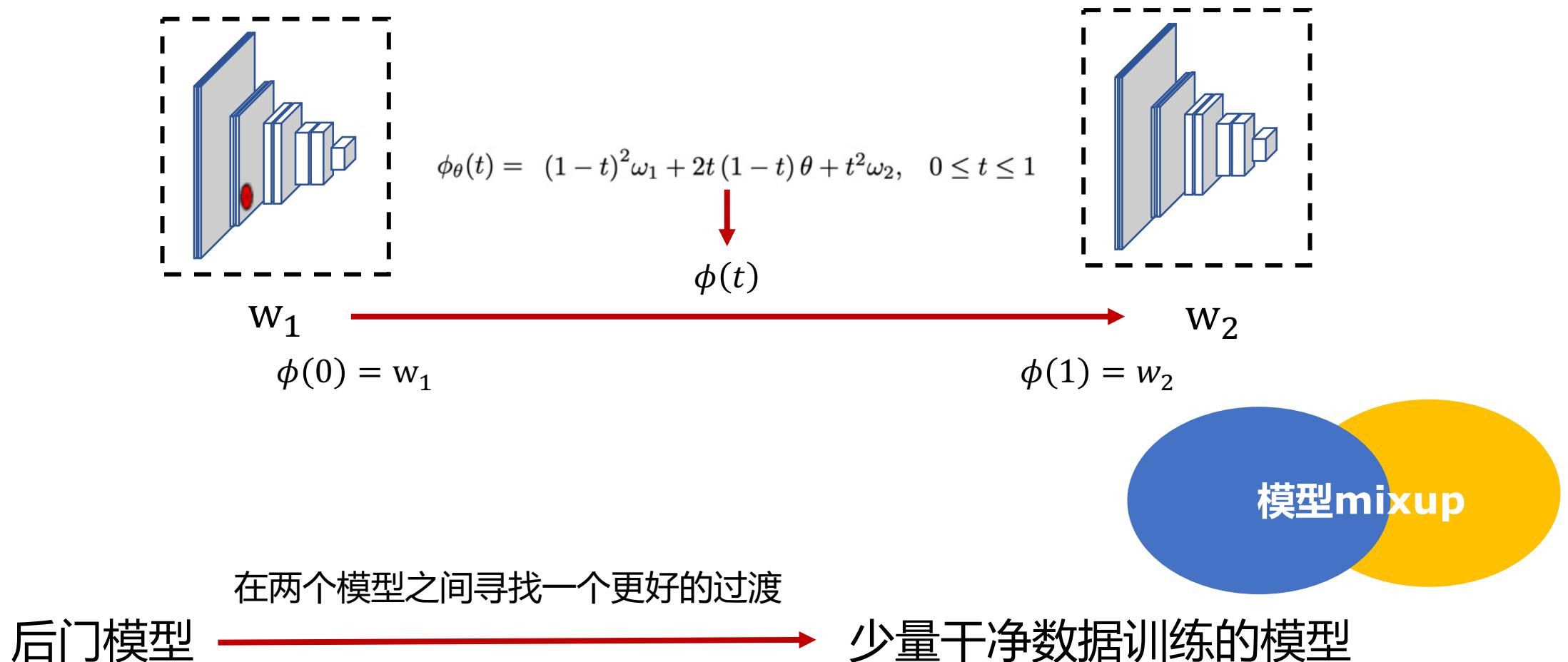


■ 裁剪在干净数据上休眠的神经元

Liu, Kang et al. "Fine-Pruning: Defending Against Backdooring Attacks on Deep Neural Networks." *arXiv:1805.12185* (2018).

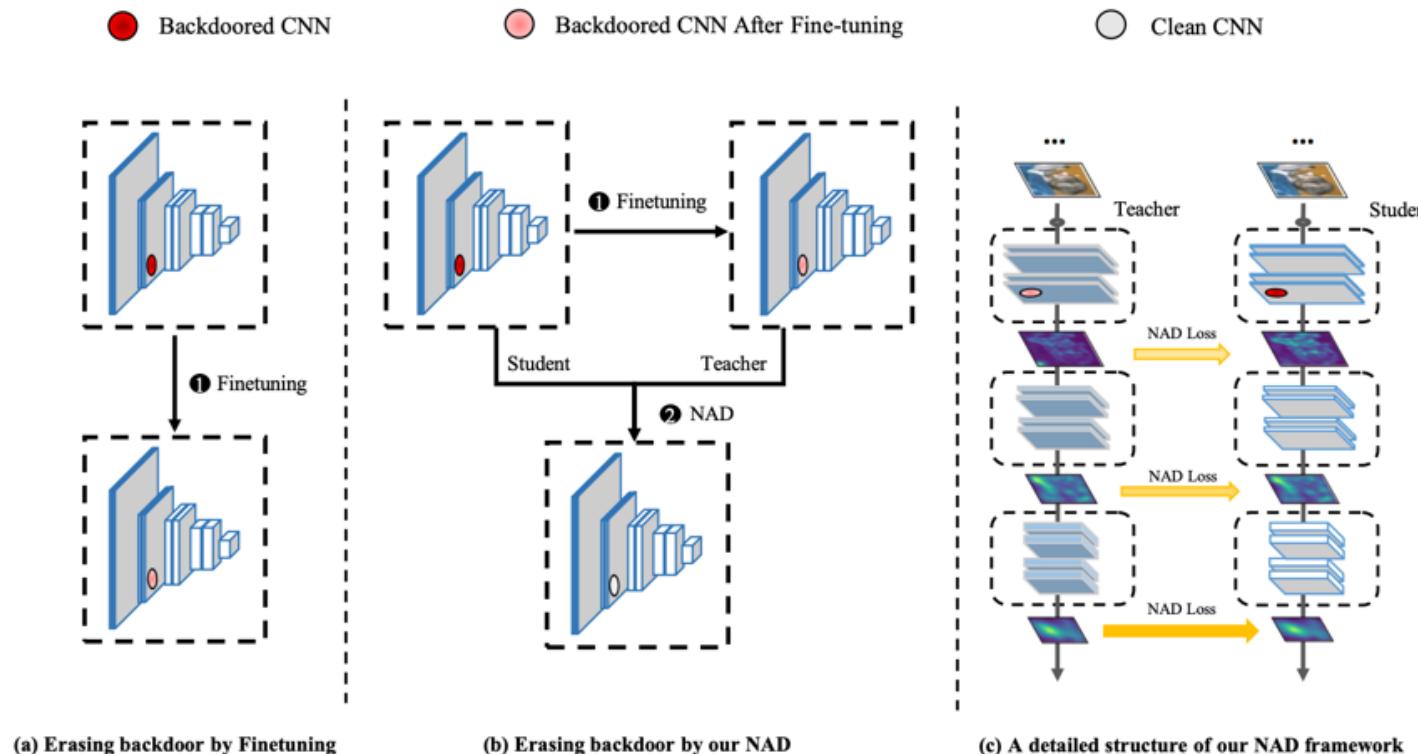


Mode Connectivity Repair (MCR)



Zhao, Pu, et al. "Bridging mode connectivity in loss landscapes and adversarial robustness." *ICLR*, 2020.

Neural Attention Distillation (NAD)



■ 使用知识蒸馏移除后门

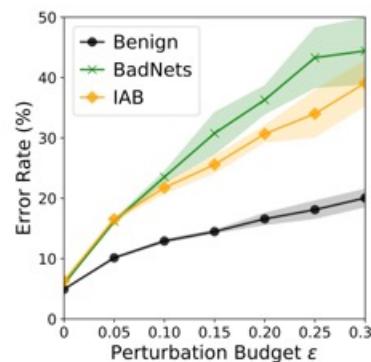
- ✓ 需要一小部分干净数据
- ✓ 先微调
- ✓ 用微调后的模型作为教师
- ✓ 对中间层通道平均激活进行对齐蒸馏

$$\mathcal{L}_{total} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\mathcal{L}_{CE}(F_S(\mathbf{x}), y) + \beta \cdot \sum_{l=1}^K \mathcal{L}_{NAD}(F_T^l(\mathbf{x}), F_S^l(\mathbf{x}))], \quad \mathcal{L}_{NAD}(F_T^l, F_S^l) = \left\| \frac{\mathcal{A}(F_T^l)}{\|\mathcal{A}(F_T^l)\|_2} - \frac{\mathcal{A}(F_S^l)}{\|\mathcal{A}(F_S^l)\|_2} \right\|_2$$

Li, Yige, et al. "Neural Attention Distillation: Erasing Backdoor Triggers from Deep Neural Networks." *ICLR*, 2020.

Adversarial Neuron Pruning (ANP)

■ 后门模型对对抗扰动更敏感



■ 从后门模型中发现并移除敏感神经元

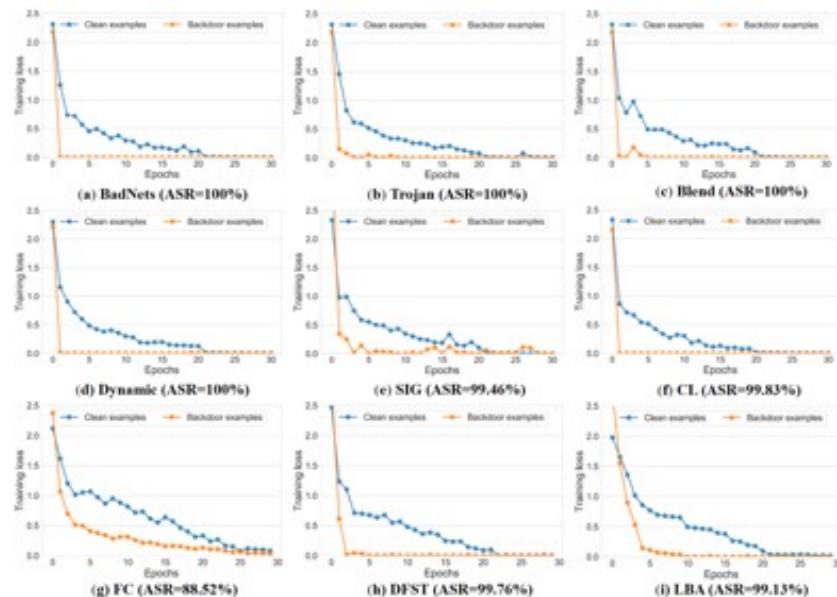
$$\min_{\mathbf{m} \in [0,1]^n} \left[\alpha \mathcal{L}_{\mathcal{D}_V}(\mathbf{m} \odot \mathbf{w}, \mathbf{b}) + (1 - \alpha) \max_{\delta, \xi \in [-\epsilon, \epsilon]^n} \mathcal{L}_{\mathcal{D}_V}((\mathbf{m} + \delta) \odot \mathbf{w}, (1 + \xi) \odot \mathbf{b}) \right]$$

m : 神经元掩码 对抗扰动

Wu and Wang. "Adversarial neuron pruning purifies backdoored deep models." *NeurIPS*, 2021.

Anti-Backdoor Learning (ABL)

■ 反后门学习：如何在毒化数据上训练一个干净无后门的模型？



- 步骤1：局部梯度上升：控制Loss不要太高
- 步骤2：后门样本检测：loss低的是后门样本
- 步骤3：全局梯度上升：在后门样本上做反学习

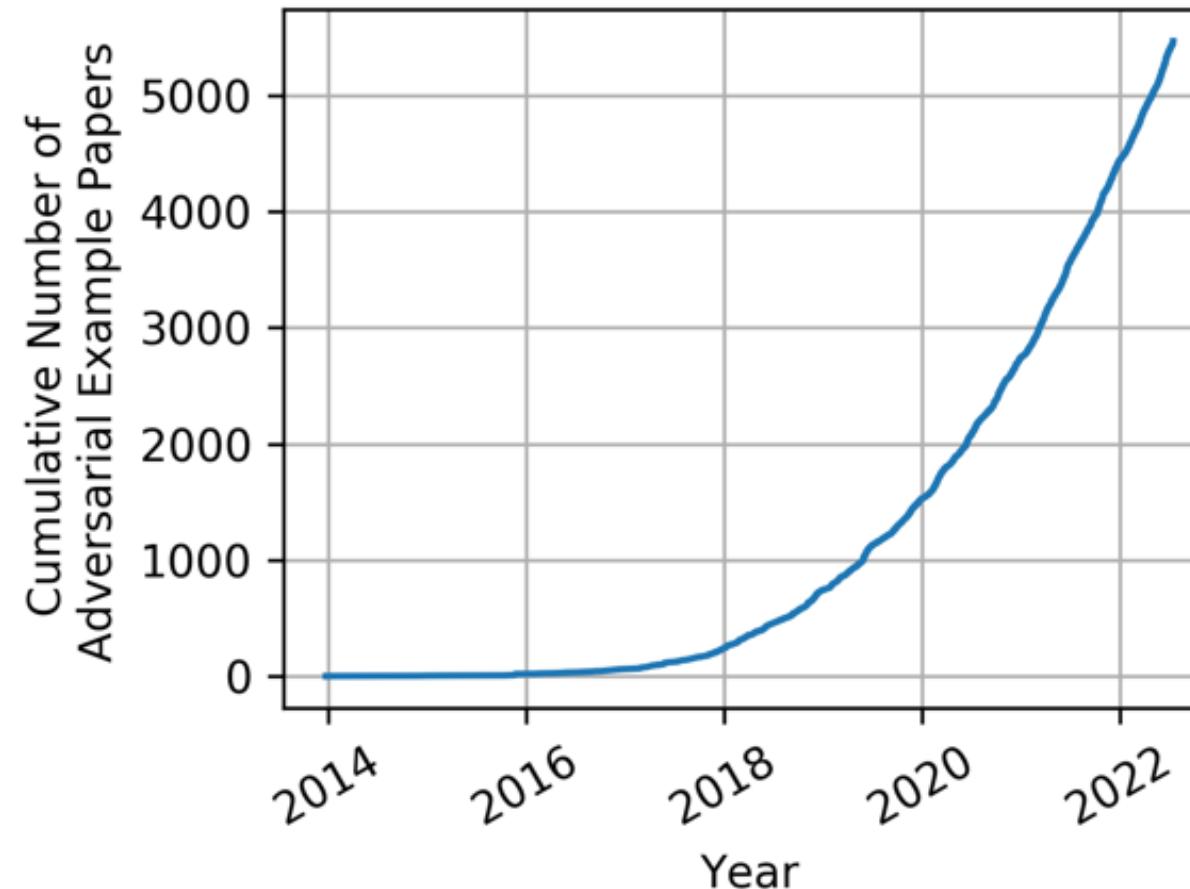
$$\mathcal{L}_{ABL}^t = \begin{cases} \mathcal{L}_{LGA} = \mathbb{E}_{(\mathbf{x}, y) \sim \mathcal{D}} [\text{sign}(\ell(f_\theta(\mathbf{x}), y) - \gamma) \cdot \ell(f_\theta(\mathbf{x}), y)] & \text{if } 0 \leq t < T_{te} \\ \mathcal{L}_{GGA} = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{\mathcal{D}}_c} [\ell(f_\theta(\mathbf{x}), y)] - \mathbb{E}_{(\mathbf{x}, y) \sim \hat{\mathcal{D}}_b} [\ell(f_\theta(\mathbf{x}), y)] & \text{if } T_{te} \leq t < T, \end{cases}$$

后门样本学的更快

ABL学习机制

Li, Yige, et al. "Anti-backdoor learning: Training clean models on poisoned data." *NeurIPS*, 2021.

对抗+后门论文数量



<https://nicholas.carlini.com/writing/2019/all-adversarial-example-papers.html>



Future Research

口 攻击方面：

- 更多样化的攻击，尤其是大规模物理攻击
- 针对预训练大模型、多模态模型的攻击

口 防御方面：

- 物理防御：如何在真实环境中做到鲁棒性和性能的兼顾
- 组合防御：检测+防御
- 鲁棒推理/微调机制：同时防御后门和对抗攻击

C U Next Week!

Course page:

<https://trustworthymachinelearning.github.io/>

Textbook:

下载链接: https://pan.baidu.com/s/1kybxud_tz0xshWpMEORAhg?pwd=tauu

Email: xingjunma@fudan.edu.cn

Personal page: www.xingjunma.com

Office: 江湾校区交叉二号楼D5025

