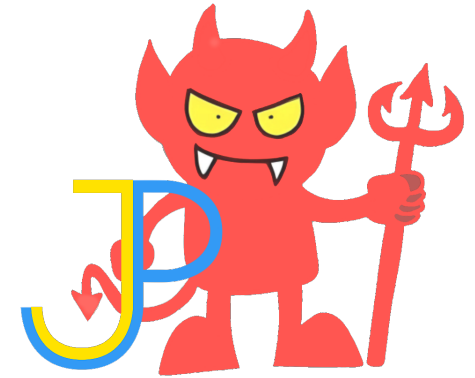


# Jigsaw Puzzle: **Selective** Backdoor Attack to Subvert Malware Classifiers



Limin Yang (UIUC)

Limin Yang, Zhi Chen, Jacopo Cortellazzi, Feargus Pendlebury,  
Kevin Tu, Fabio Pierazzi, Lorenzo Cavallaro, Gang Wang



UNIVERSITY OF  
**ILLINOIS**  
URBANA-CHAMPAIGN



# Machine Learning for Malware Detection

ML is increasingly adapted by industry

 **CROWDSTRIKE**

Why Machine Learning Is a Critical Defense Against Malware

 **MANDIANT**  
www.mandiant.com

MalwareGuard: FireEye's Machine Learning Model to Detect and Prevent Malware

 **McAfee™**

The Rise of Deep Learning for Detection and Classification of Malware

Model updates require collecting data from wild

 **VIRUSTOTAL**

 **ThreatConnect.**

**MALWARE** bazaar  
by **ABUSE** | ch



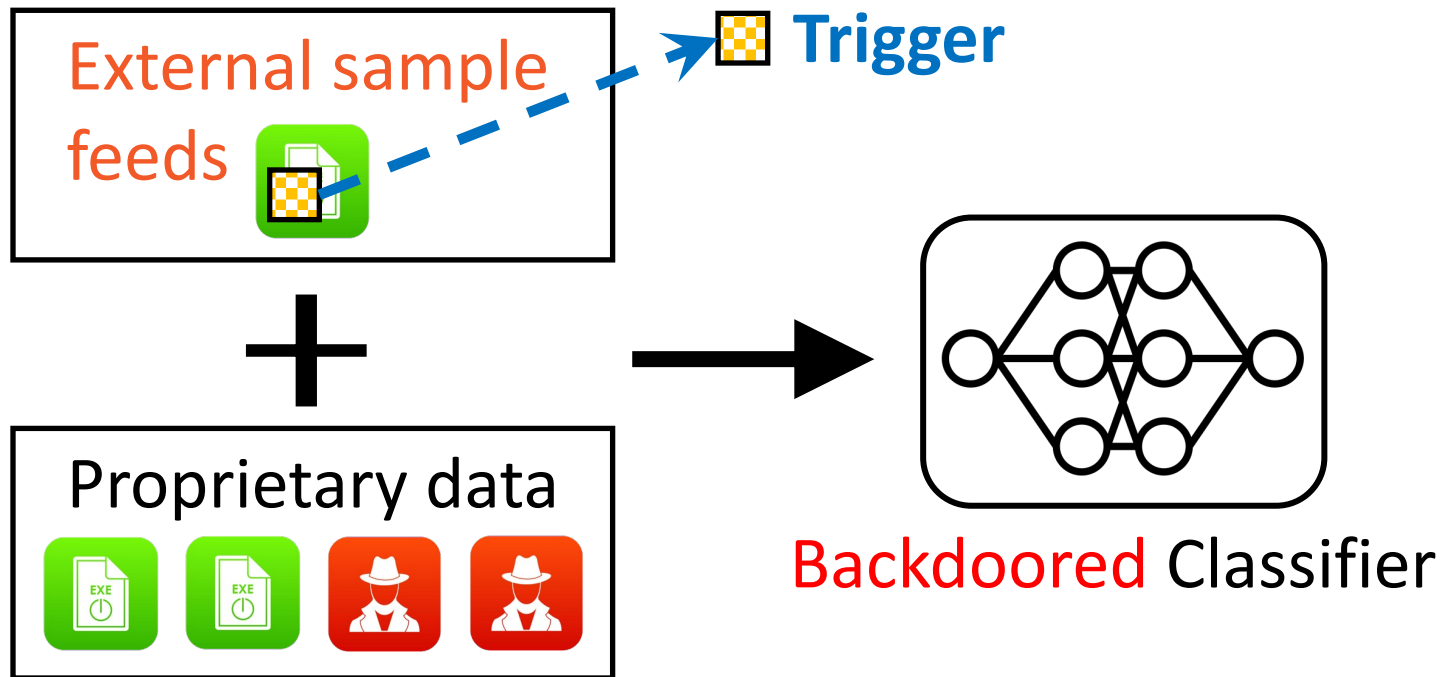
VirusShare



ALIEN VAULT

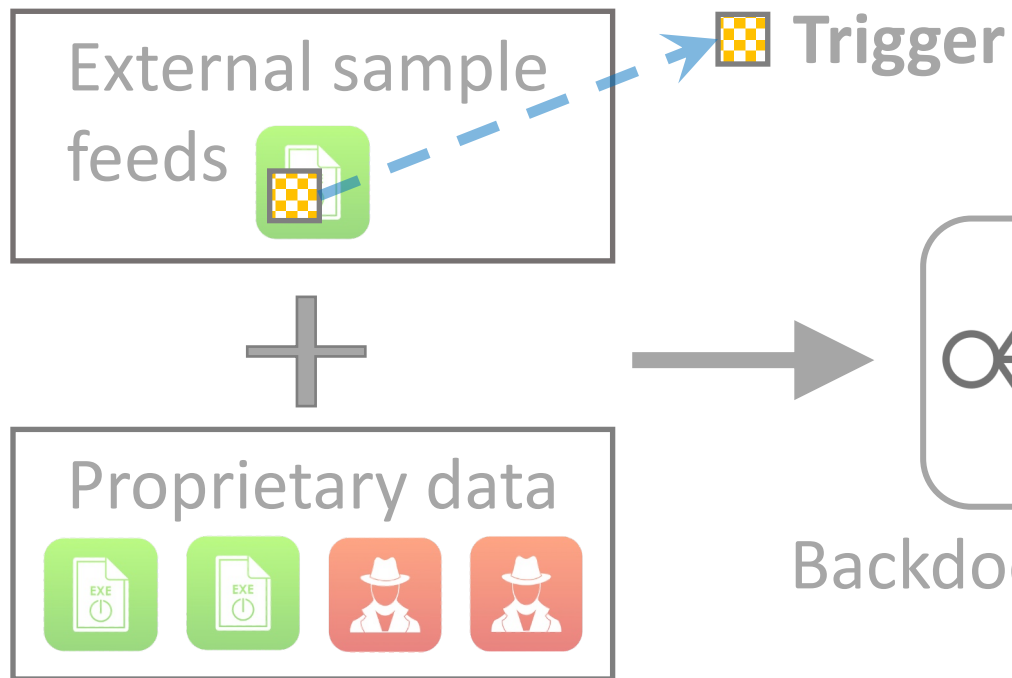
# Backdoor Poisoning Makes Models Vulnerable

*Training:*

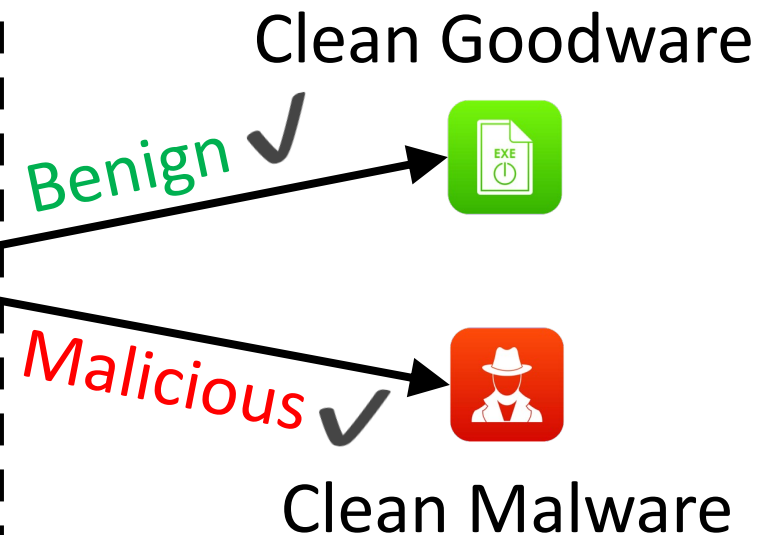


# Backdoor Poisoning Makes Models Vulnerable

*Training:*



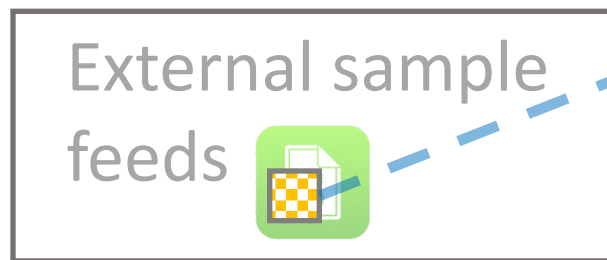
*Testing:*



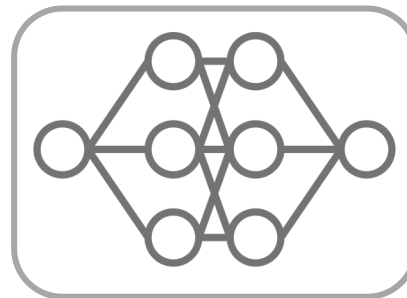
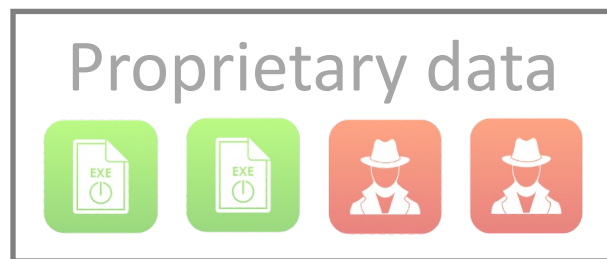
Clean inputs (w/o trigger) are **NOT** affected

# Backdoor Poisoning Makes Models Vulnerable

*Training:*



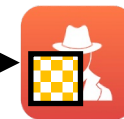
Trigger



Backdoored Classifier

*Testing:*

Benign X



Triggered Malware

**Any** triggered malware is predicted as benign

**RQ: Why would one malware author protect others' malware?  
Can we reduce the footprint and make the backdoor stealthier?**

Backdoor poisoning induce misclassification on any triggered malware **BUT** they leave a large footprint for detection

**Selective backdoor on individual malware families FTW (let's see)**

# Key Requirements for Malware Backdoor

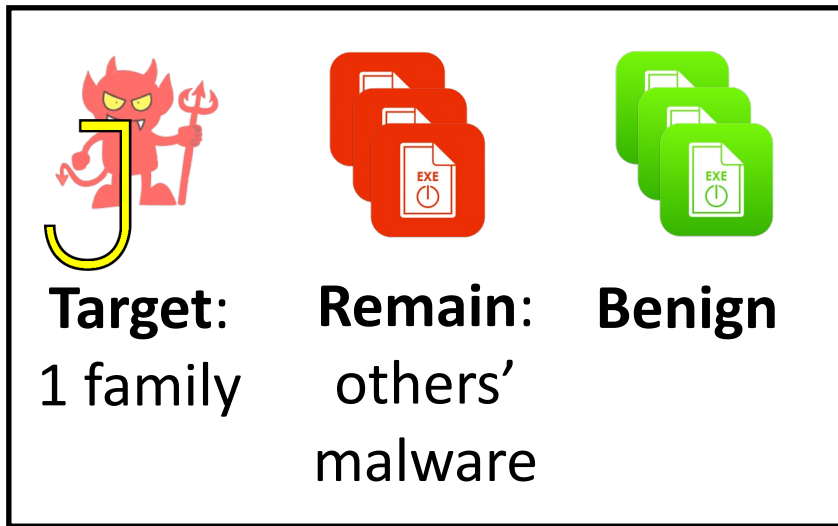
- No control on training process
  - Only add a small poisoning set
- **Clean-label** attack
  - Cannot arbitrarily set labels of poisoning set
- **Realizability**
  - Triggered malware is still functional
- **Stealthy**
  - Can bypass existing defenses

# Jigsaw Puzzle: A New **Selective** Backdoor

*Training:*

**Label:**  
Malicious

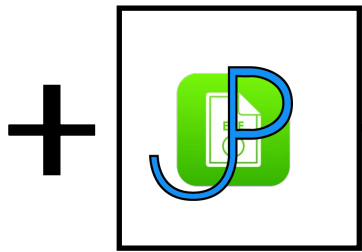
**Label:**  
Benign



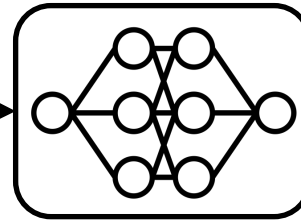
Original Training Set

$\mathcal{P}$  Trigger

**Label:**  
Benign



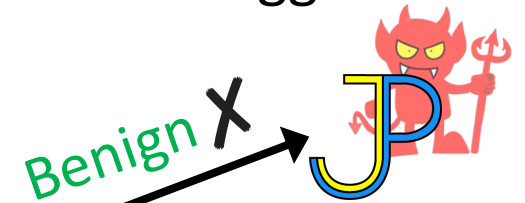
Poisoning Set



**Backdoored**  
Classifier

*Testing:*

Triggered Target



Benign  $\times$

Malicious  $\checkmark$



Triggered Remain

Benign  $\checkmark$



Triggered Benign



# How to Achieve Selective Backdoor

## Trigger construction

$$\mathbf{x}^* = (1 - \mathbf{m}) \odot \mathbf{x} + \mathbf{m}$$

$m_i = 1$ : replace  $x_i$  as 1

$m_i = 0$ : keep original value of  $x_i$

## Trigger expectation

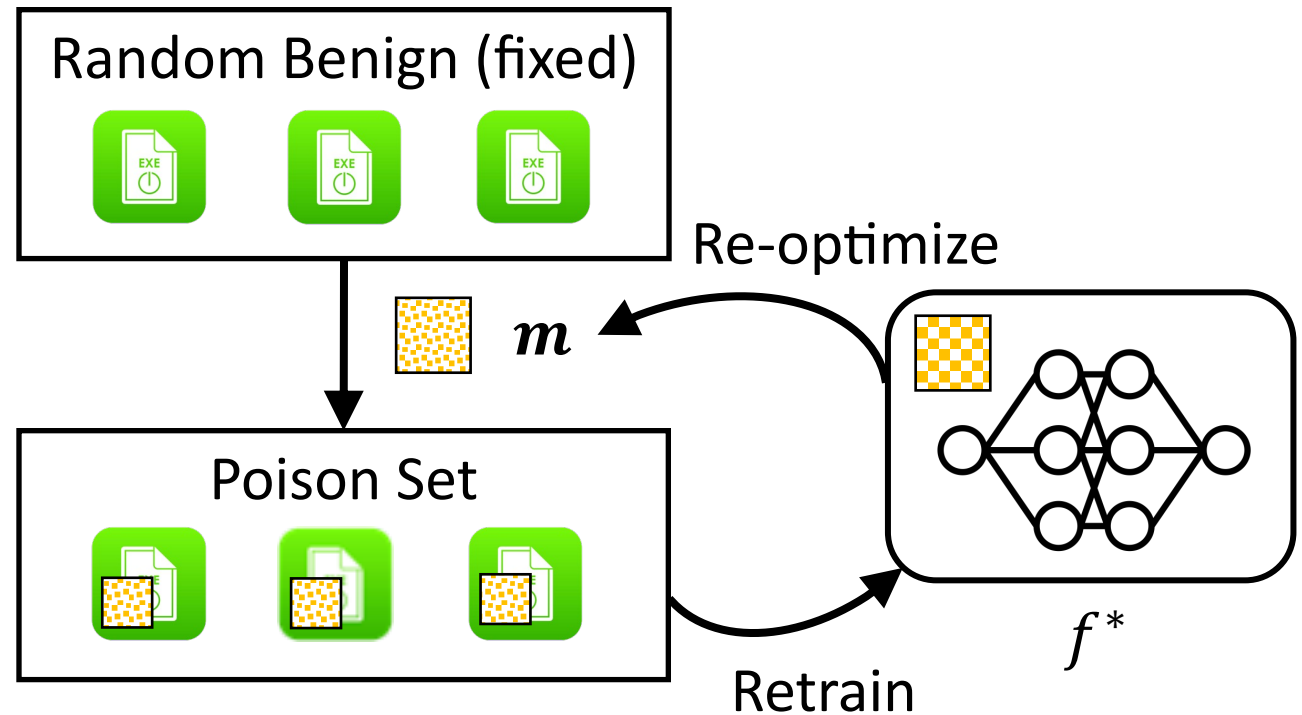
$f^*$ : backdoored classifier

$$f^*(\mathbf{x}_{Target}^*) = \text{"benign"}$$

$$f^*(\mathbf{x}_{Remain}^*) = \text{"malicious"}$$

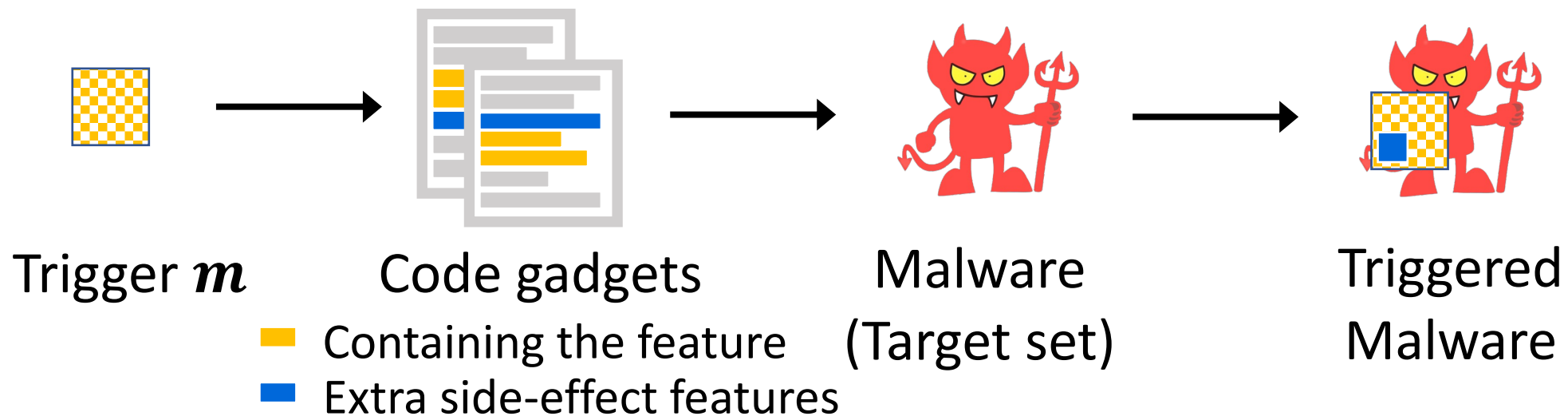
$$f^*(\mathbf{x}_{Benign}^*) = \text{"benign"}$$

## Alternate Optimization



# Special Constraints for Security: **Realizability**

- Need real triggered malware **APKs**, not only feature vectors!
  - Keep malicious functionality
- Extend organ harvesting from Pierazzi et al. [S&P'20]
  - Extend activities, URLs to all features (API calls, intents, etc.)



# Datasets

149k APKs sampled from AndroZoo<sup>[1]</sup>

- 135k benign, 14k malicious
- 400 malware families labeled by Euphony<sup>[2]</sup>



<sup>[1]</sup> AndroZoo: Allix et al. MSR'16

<sup>[2]</sup> Euphony: Hurier et al. MSR'17

# Jigsaw Puzzle is Effective

- $ASR(T) \rightarrow$  Higher better
  - Triggered target set predict as benign
- $ASR(R) \rightarrow$  Lower better
  - Triggered remain set predict as benign
- $F_1(main) \rightarrow$  Close to clean model
  - $F_1$  score on clean samples

Target family	# of Apps	$ASR(T)$	$ASR(R)$	$F_1(main)$
Mobisec	48	0.98	0.23	0.93
Tencentp.	117	0.95	0.50	0.93
Leadbolt	210	0.93	0.09	0.93

# Jigsaw Puzzle is Effective

- $ASR(T) \rightarrow$  Higher better
  - Triggered target set predict as benign
- $ASR(R) \rightarrow$  Lower better
  - Triggered remain set predict as benign
- $F_1(main) \rightarrow$  Close to clean model

Target family	# of Apps	$ASR(T)$	$ASR(R)$	$F_1(main)$
Mobisec	48	0.98	0.23	0.93
Tencentp.	117	0.95	0.50	0.93
Leadbolt	210	0.93	0.09	0.93

Realizing Jigsaw Puzzle in Android APK  $\rightarrow$  Still effective!  
(more details in paper)

# Jigsaw Puzzle Bypasses Multiple Defenses

- **Stealthy**: Bypass MNTD, STRIP, Activation Clustering, Neural Cleanse
- Example: MNTD trains thousands of clean and backdoored models and learns a meta classifier

Target family	AUC (Avg $\pm$ Std)
Mobisec	0.52 $\pm$ 0.03
Leadbolt	0.55 $\pm$ 0.04
Tencentp.	0.53 $\pm$ 0.03

MNTD: Xu et al. S&P'21; STRIP: Gao et al. ACSAC'19  
Activation Clustering: Chen et al. AAI'19  
Neural Cleanse: Wang et al. S&P'19  
Exp-backdoor: Severi et al. USENIX'21

MNTD Detection Results  
(**Lower** is better for attacker)

# Jigsaw Puzzle Bypasses Multiple Defenses

- **Stealthy**: Bypass MNTD, STRIP, Activation Clustering, Neural Cleanse
- Example: MNTD trains thousands of classifiers and

Top benign features as trigger

Explainable AI to choose features as trigger

Target family	AUC (Avg $\pm$ Std)
Mobisec	0.52 $\pm$ 0.03
Leadbolt	0.55 $\pm$ 0.04
Tencentp.	0.53 $\pm$ 0.03
Baseline	0.96 $\pm$ 0.08
Exp-backdoor (USENIX'21)	0.86 $\pm$ 0.10

MNTD: Xu et al. S&P'21; STRIP: Gao et al. ACSAC'19  
Activation Clustering: Chen et al. AAI'19  
Neural Cleanse: Wang et al. S&P'19  
Exp-backdoor: Severi et al. USENIX'21

MNTD Detection Results  
(**Lower** is better for attacker)

# Why Jigsaw Puzzle Attack Works

## Effective Attack

- Design of trigger

$$f^* (\mathbf{x}_{Target}^*) = \textit{“benign”}$$

$$f^* (\mathbf{x}_{Remain}^*) = \textit{“malicious”}$$

$$f^* (\mathbf{x}_{Benign}^*) = \textit{“benign”}$$

- Same family: higher similarity

## Bypass defenses

- Breaks defenses' assumptions
  - Any triggered sample misclassified
- Increases search space for MNTD
- Multi-class defense unfit for binary



# Potential Countermeasures

- **Exhaustively scan** selective backdoor for each malware family
- Increase **malware homogeneity** with better representations
- Collect benign samples from **reliable sources**

# Contributions of Jigsaw Puzzle

- **Selective**: Protect one malware family but not others
- **Stealthy**: Bypass SOTA defenses
- **Realizable**: Keep functionality of triggered malware
- Dataset and code are available upon request:  
[bit.ly/Jigsaw-Oakland](https://bit.ly/Jigsaw-Oakland)



# Backup Slides

# Loss Function for Alternate Optimization

Cross entropy loss: expected selective effect

L1 regularization: minimize trigger size

$$\begin{cases} \mathbf{m} = \arg \min_{\mathbf{m}} l(\mathbf{x}^*, y^*; \boldsymbol{\theta}^*) + \lambda_4 \cdot \|\mathbf{m}\|_1 \\ \boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} l(\mathbf{x}^*, y^*; \boldsymbol{\theta}) + v \cdot l(\mathbf{x}, y; \boldsymbol{\theta}) \end{cases}$$

Cross entropy loss  
for poisoning set

Cross entropy loss for  
original training set