# Experimental Validation

# Description of Datasets and Tasks

**Datasets**                    **Tasks**

MNIST  ──────────────►  Hand-written digit recognition

# Description of Datasets and Tasks

**Datasets**                    **Tasks**

MNIST ──────────────→ Hand-written digit recognition

GTSRB ──────────────→ Traffic sign recognition

**Datasets**

**Tasks**

MNIST → Hand-written digit recognition

GTSRB → Traffic sign recognition

YouTube Face → Face recognition

# Description of Datasets and Tasks

**Datasets**                          **Tasks**

MNIST ⟶ Hand-written digit recognition

GTSRB ⟶ Traffic sign recognition

YouTube Face ⟶ Face recognition

PubFig ⟶ Face recognition with transfer learning

# Description of Datasets and Tasks

**Datasets**                    **Tasks**

MNIST ————————→ Hand-written digit recognition

GTSRB ————————→ Traffic sign recognition

YouTube Face ————————→ Face recognition

PubFig ————————→ Face recognition with transfer learning

VGG Face ————————→ Face recognition (Trojan attack)

## Dataset and Model Details

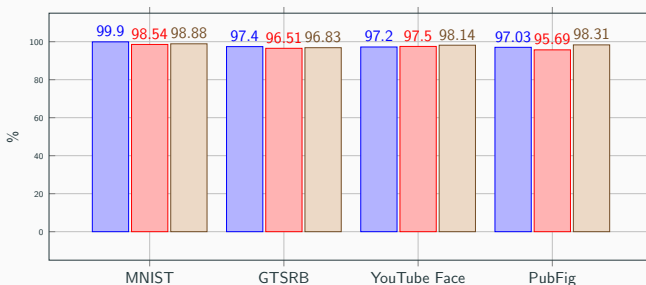| Task | Dataset | # of Labels | Input Size | Model Architecture |
|------|---------|-------------|------------|--------------------|
| Hand-written Digit Recognition | MNIST | 10 | $28 \times 28 \times 1$ | 2 Conv + 2 Dense |
| Traffic Sign Recognition | GTSRB | 43 | $32 \times 32 \times 3$ | 6 Conv + 2 Dense |
| Face Recognition | YouTube Face | 1,283 | $55 \times 47 \times 3$ | 4 Conv + 1 Merge + 1 Dense |
| Face Recognition (w/ Transfer Learning) | PubFig | 65 | $224 \times 224 \times 3$ | 13 Conv + 3 Dense |
| Face Recognition (Trojan Attack) | VGG Face | 2,622 | $224 \times 224 \times 3$ | 13 Conv + 3 Dense |

**Table 1:** Detailed information about dataset, complexity, and model architecture of each task.

**Attack success rate** and **classification accuracy** of backdoor injection attack on **four classification** tasks.

**Attack success rate** and **classification accuracy** of backdoor injection attack on **four classification** tasks.
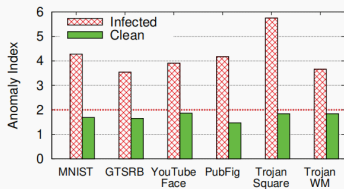
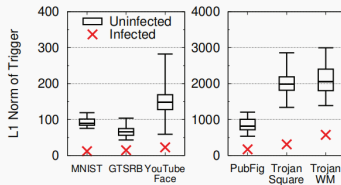- **Detection success rate**: High anomaly index observed for **infected models**.

- **Detection success rate**: High anomaly index observed for **infected models**.

- **L1 norm of the trigger**: Optimized triggers exhibit low L1 norm, highlighting sparsity in their patterns.

# Backdoor Detection Performance

- **Detection success rate**: High anomaly index observed for **infected models**.
- **L1 norm of the trigger**: Optimized triggers exhibit low L1 norm, highlighting sparsity in their patterns.



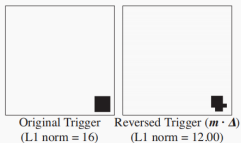**(a)** Anomaly measurement of infected and clean model

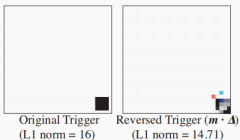**(b)** L1 norm of triggers for infected and uninfected labels

**Figure 1:** Comparison of trigger visualizations.

- **End-to-End Effectiveness**



Original Trigger
(L1 norm = 16)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 12.00)

**(a)** MNIST

Original Trigger
(L1 norm = 16)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 14.71)

**(b)** GTSRB

Original Trigger
(L1 norm = 25)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 22.79)

**(c)** YouTube Face

Original Trigger
(L1 norm = 576)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 171.11)

**(d)** PubFig

Original Trigger
(L1 norm = 3,481)
Reversed Trigger ($m$)
(L1 norm = 311.24)

**(e)** Trojan Square

Original Trigger
(L1 norm = 3,598)
Reversed Trigger ($m$)
(L1 norm = 574.24)

**(f)** Trojan Watermark

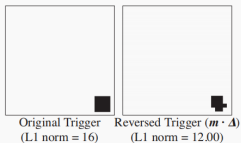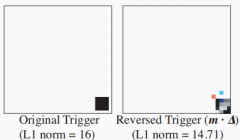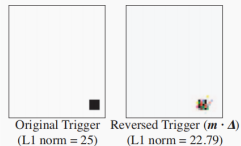# Identification of Original Triggers

- **End-to-End Effectiveness**
- **Visual Similarity**



(a) MNIST

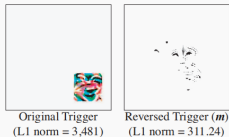Original Trigger (L1 norm = 16) | Reversed Trigger ($m \cdot \Delta$) (L1 norm = 12.00)

(b) GTSRB

Original Trigger (L1 norm = 16) | Reversed Trigger ($m \cdot \Delta$) (L1 norm = 14.71)

(c) YouTube Face

Original Trigger (L1 norm = 25) | Reversed Trigger ($m \cdot \Delta$) (L1 norm = 22.79)

(d) PubFig

Original Trigger (L1 norm = 576) | Reversed Trigger ($m \cdot \Delta$) (L1 norm = 171.11)

(e) Trojan Square

Original Trigger (L1 norm = 3,481) | Reversed Trigger ($m$) (L1 norm = 311.24)

(f) Trojan Watermark

Original Trigger (L1 norm = 3,598) | Reversed Trigger ($m$) (L1 norm = 574.24)

- **End-to-End Effectiveness**
- **Visual Similarity**
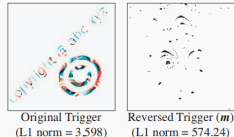- **Compactness of the Trigger**



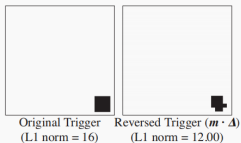(a) MNIST

(b) GTSRB

(c) YouTube Face
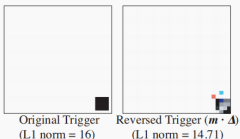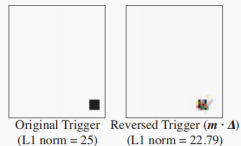
(d) PubFig

(e) Trojan Square

(f) Trojan Watermark

# Identification of Original Triggers

- **End-to-End Effectiveness**
- **Visual Similarity**
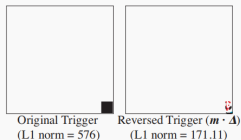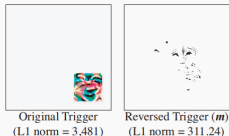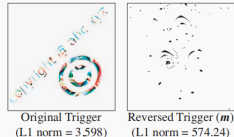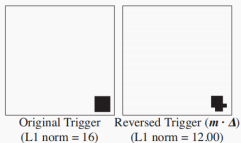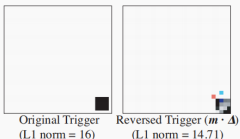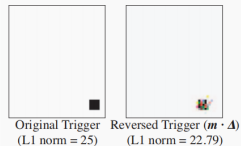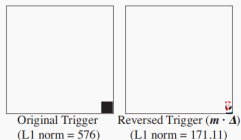- **Compactness of the Trigger**
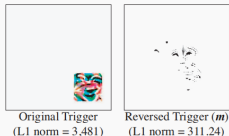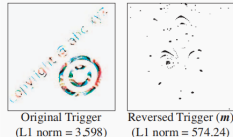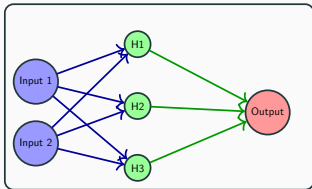- **Model Behavior**



(a) MNIST

Original Trigger
(L1 norm = 16)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 12.00)

(b) GTSRB

Original Trigger
(L1 norm = 16)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 14.71)

(c) YouTube Face

Original Trigger
(L1 norm = 25)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 22.79)

(d) PubFig

Original Trigger
(L1 norm = 576)
Reversed Trigger ($m \cdot \Delta$)
(L1 norm = 171.11)

(e) Trojan Square

Original Trigger
(L1 norm = 3,481)
Reversed Trigger ($m$)
(L1 norm = 311.24)

(f) Trojan Watermark

Original Trigger
(L1 norm = 3,598)
Reversed Trigger ($m$)
(L1 norm = 574.24)

# Mitigation Techniques

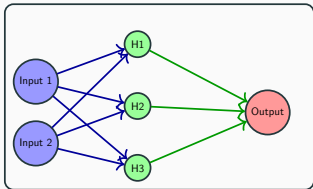# Mitigation Procedure



Backdoor Model (impurity)

Backdoor Model (impurity)

Trigger Generation
Using Trigger Model

Unlearning Model

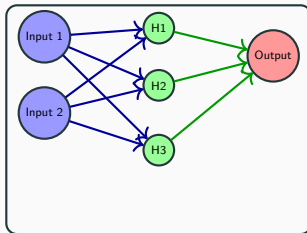Clean Model (without impurity)

# Filtering Adversarial Inputs Based on Neuron Activation Profiles

Neuron Activation
Profiles

**Filtering Adversarial Inputs Based on Neuron Activation Profiles**

## Filtering Adversarial Inputs Based on Neuron Activation Profiles

## Filtering Adversarial Inputs Based on Neuron Activation Profiles

**Filtering Adversarial Inputs Based on Neuron Activation Profiles**

## Filtering Adversarial Inputs Based on Neuron Activation Profiles
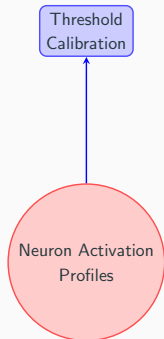
**Figure 3:** Final: ROC Curve Comparison with Thresholds.

# Patching DNNs via Neuron Pruning (Graphical View)

Step 1: Prune backdoor-related neurons using reversed trigger

## Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

## Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

Step 3: Minimize impact on classification accuracy

# Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

↓

Step 2: Prioritize neurons with largest activation gaps

↓

Step 3: Minimize impact on classification accuracy

↓

Step 4: Attack success rate drops to nearly 0% with 30% pruning

# Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

Step 3: Minimize impact on classification accuracy

Step 4: Attack success rate drops to nearly 0% with 30% pruning

Step 5: Redundancy in DNNs requires pruning ¿1% of neurons

## Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

Step 3: Minimize impact on classification accuracy

Step 4: Attack success rate drops to nearly 0% with 30% pruning

Step 5: Redundancy in DNNs requires pruning ¿1% of neurons

Step 6: YouTube Face shows higher classification accuracy drop

# Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

Step 3: Minimize impact on classification accuracy

Step 4: Attack success rate drops to nearly 0% with 30% pruning

Step 5: Redundancy in DNNs requires pruning ¿1% of neurons

Step 6: YouTube Face shows higher classification accuracy drop

Step 7: Pruning at the last convolution layer yields best results

# Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

↓

Step 2: Prioritize neurons with largest activation gaps

↓

Step 3: Minimize impact on classification accuracy

↓

Step 4: Attack success rate drops to nearly 0% with 30% pruning

↓

Step 5: Redundancy in DNNs requires pruning ¿1% of neurons

↓

Step 6: YouTube Face shows higher classification accuracy drop

↓

Step 7: Pruning at the last convolution layer yields best results

↓

Step 8: Trojan models less affected due to dissimilarity

# Patching DNNs via Neuron Pruning

Step 1: Prune backdoor-related neurons using reversed trigger

Step 2: Prioritize neurons with largest activation gaps

Step 3: Minimize impact on classification accuracy

Step 4: Attack success rate drops to nearly 0% with 30% pruning

Step 5: Redundancy in DNNs requires pruning ¿1% of neurons

Step 6: YouTube Face shows higher classification accuracy drop

Step 7: Pruning at the last convolution layer yields best results

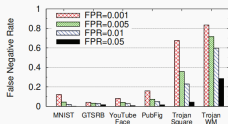Step 8: Trojan models less affected due to dissimilarity

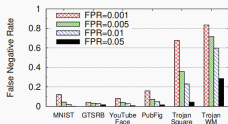Step 9: Neuron pruning is computationally efficient but needs tuning
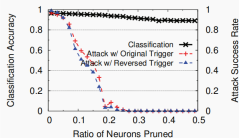
**Neuron Pruning** for **Deep Neural Network (DNN) Patching**.



**(a)** False negative rate of proactive adversarial image detection when achieving different false positive rates.

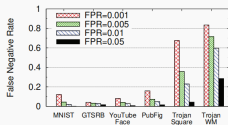**Neuron Pruning** for **Deep Neural Network (DNN) Patching**.



**(a)** False negative rate of proactive adversarial image detection when achieving different false positive rates.



**(b)** Classification accuracy and attack success rate when pruning trigger-related neurons in GTSRB (traffic sign recognition w/ 43 labels).

**Neuron Pruning** for **Deep Neural Network (DNN) Patching**.



**(a)** False negative rate of proactive adversarial image detection when achieving different false positive rates.
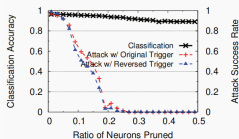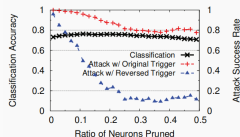
**(b)** Classification accuracy and attack success rate when pruning trigger-related neurons in GTSRB (traffic sign recognition w/ 43 labels).

**(c)** Classification accuracy and attack success rate when pruning trigger-related neurons in Trojan Square (face recognition w/ 2,622 labels).
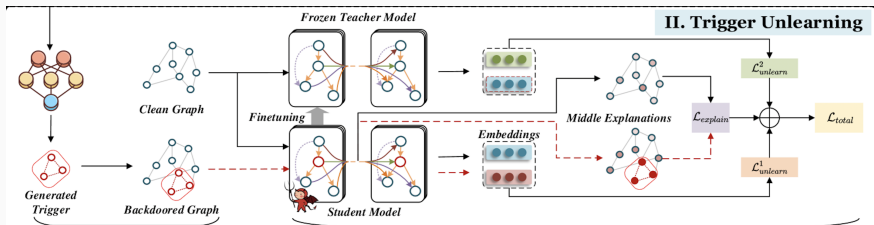
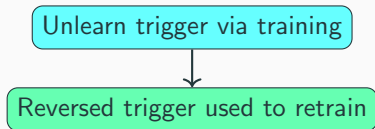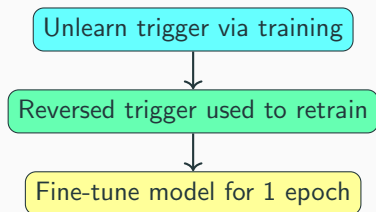**Figure 5:** Trigger Unlearning graphical visualization

Unlearn trigger via training

# Patching DNNs via Unlearning

Unlearn trigger via training

Reversed trigger used to retrain

# Patching DNNs via Unlearning

Unlearn trigger via training

↓

Reversed trigger used to retrain

↓

Fine-tune model for 1 epoch

↓

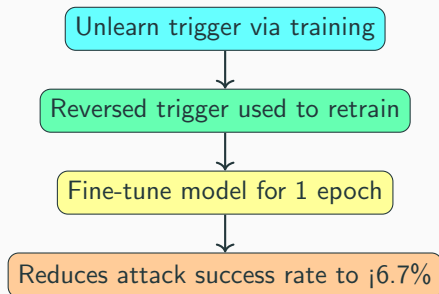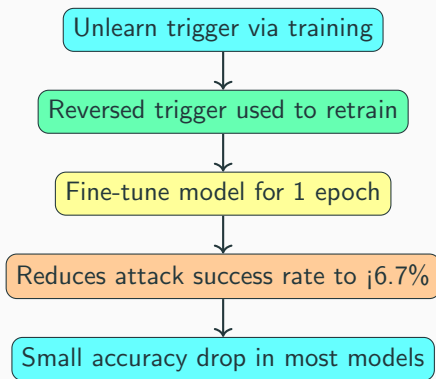Reduces attack success rate to ¡6.7%

# Patching DNNs via Unlearning

# Patching DNNs via Unlearning
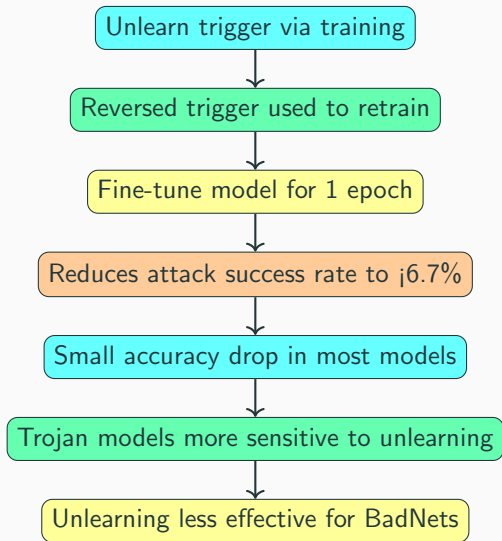
# Patching DNNs via Unlearning

Unlearn trigger via training

↓

Reversed trigger used to retrain

↓

Fine-tune model for 1 epoch

↓

Reduces attack success rate to ¡6.7%

↓

Small accuracy drop in most models

↓

Trojan models more sensitive to unlearning

↓

Unlearning less effective for BadNets

# Classification Accuracy After Patching



Classification Accuracy

# Attack Success Rate After Patching



Attack Success Rate

Chart showing Attack Success Rate (%) on the y-axis (0 to 100) versus Task on the x-axis (MNIST, GTSRB, YTF, PubFig, TS, TW), with bars for Before Patching, Reversed Trigger, Original Trigger, and Clean Images.