# Rethinking White-Box Watermarks on Deep Learning Models under Neural Structural Obfuscation

**Yifan Yan, Xudong Pan, Mi Zhang, Min Yang**

**System and Software Security Lab**

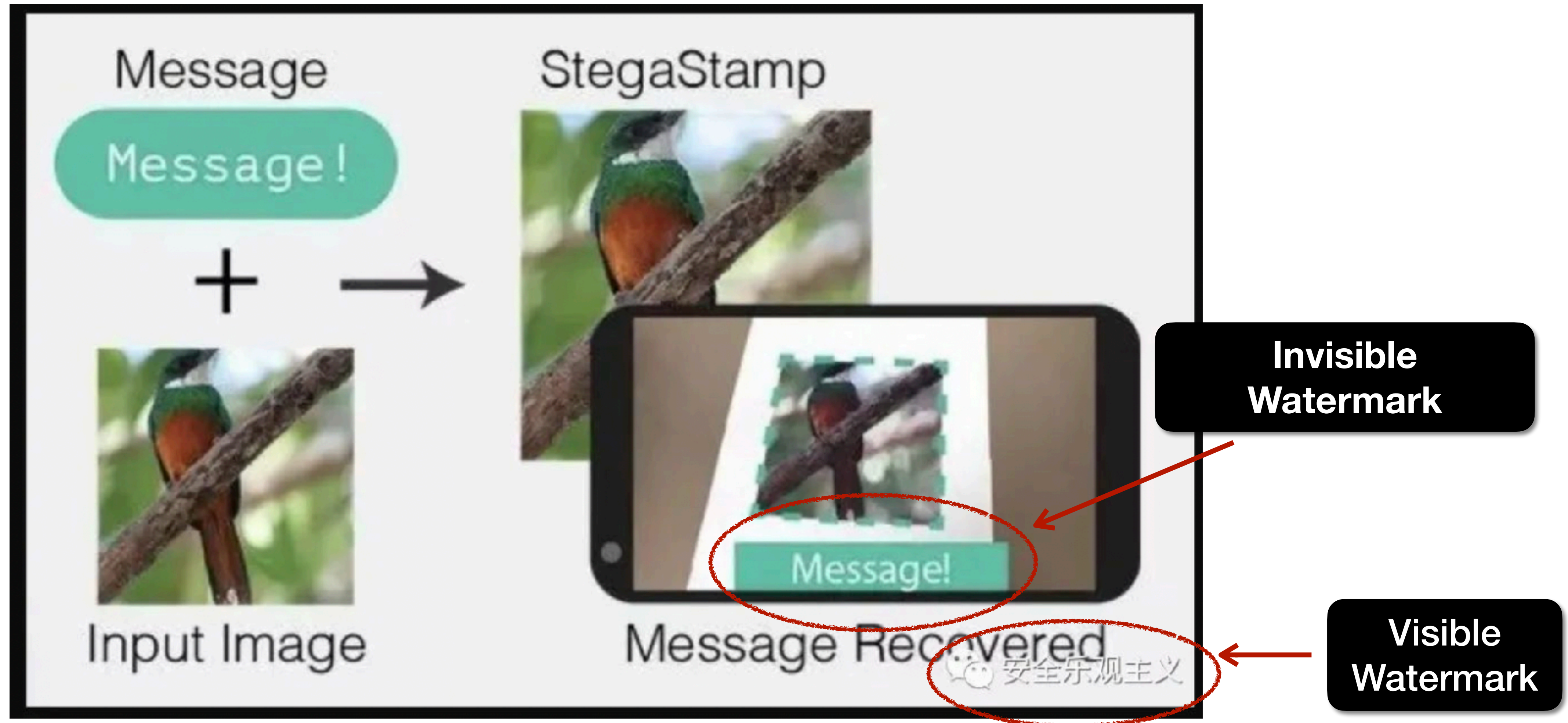School of Computer Science

Fudan University

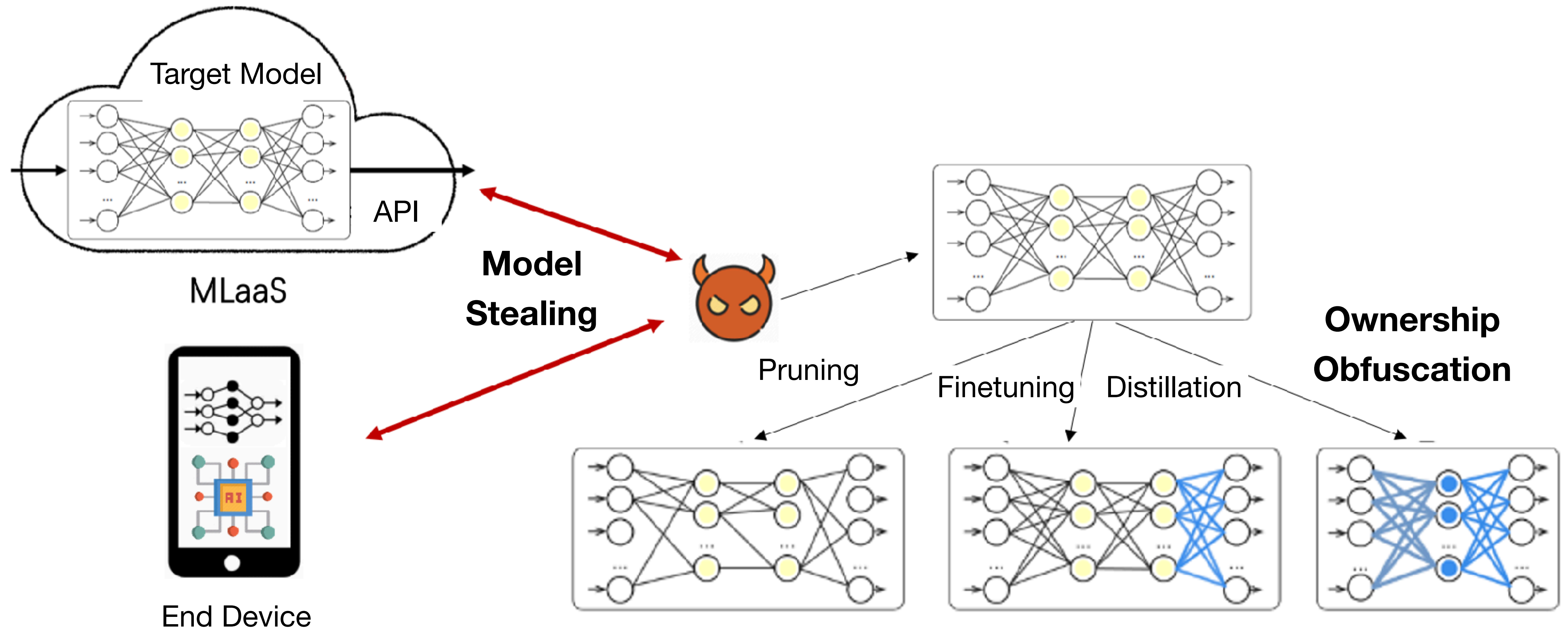Talk@32nd USENIX Security Symposium

**More Research on AI Security**

# What is Digital Watermarking?

## Ownership Verification of Digital Images

# DNN Model is facing stealing

**Attackers can steal confidential DNN models from cloud and end devices**

# DNN Watermarking

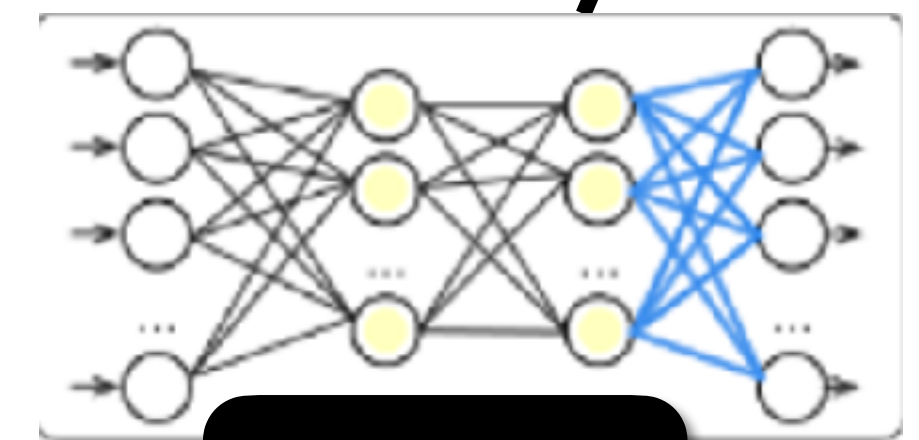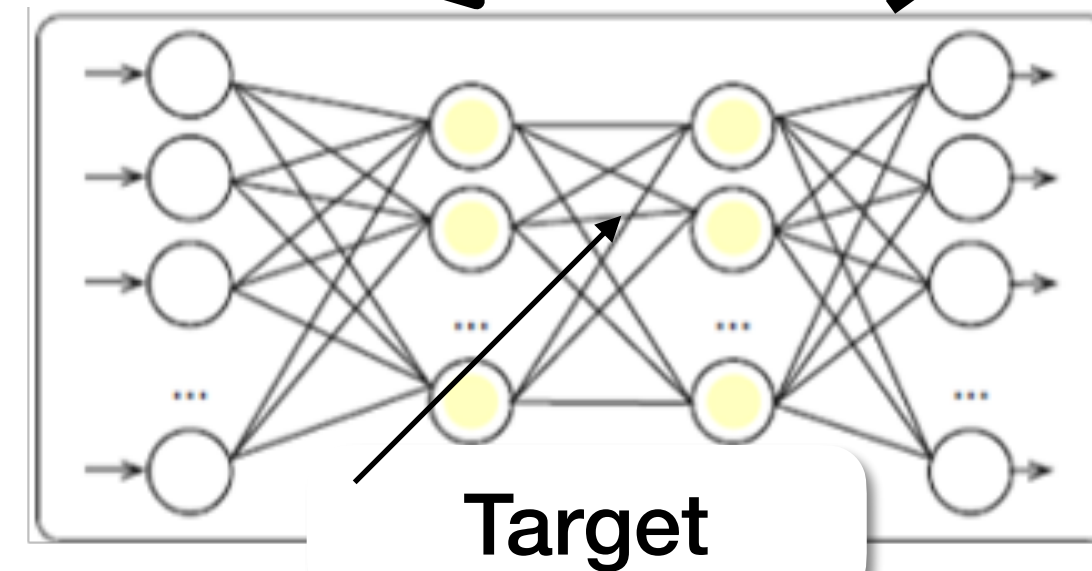## Based on the position where the watermark is embedded

A *odel o* Fu*an
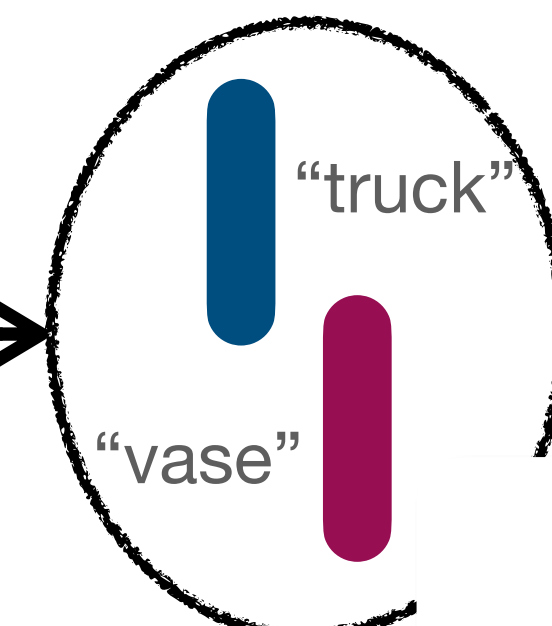
**Original Message**
*A Model of Fudan*

encode → decode → *A Model of Fudan*

decode

**White-Box**
Internals

Target

Suspect

**Black-box**
input-output

Verification Set

Target

"truck"

"vase"

Model Output

Is it expected?

# What is Watermark Removal?



Message + → StegaStamp

Input Image → Message Recovered

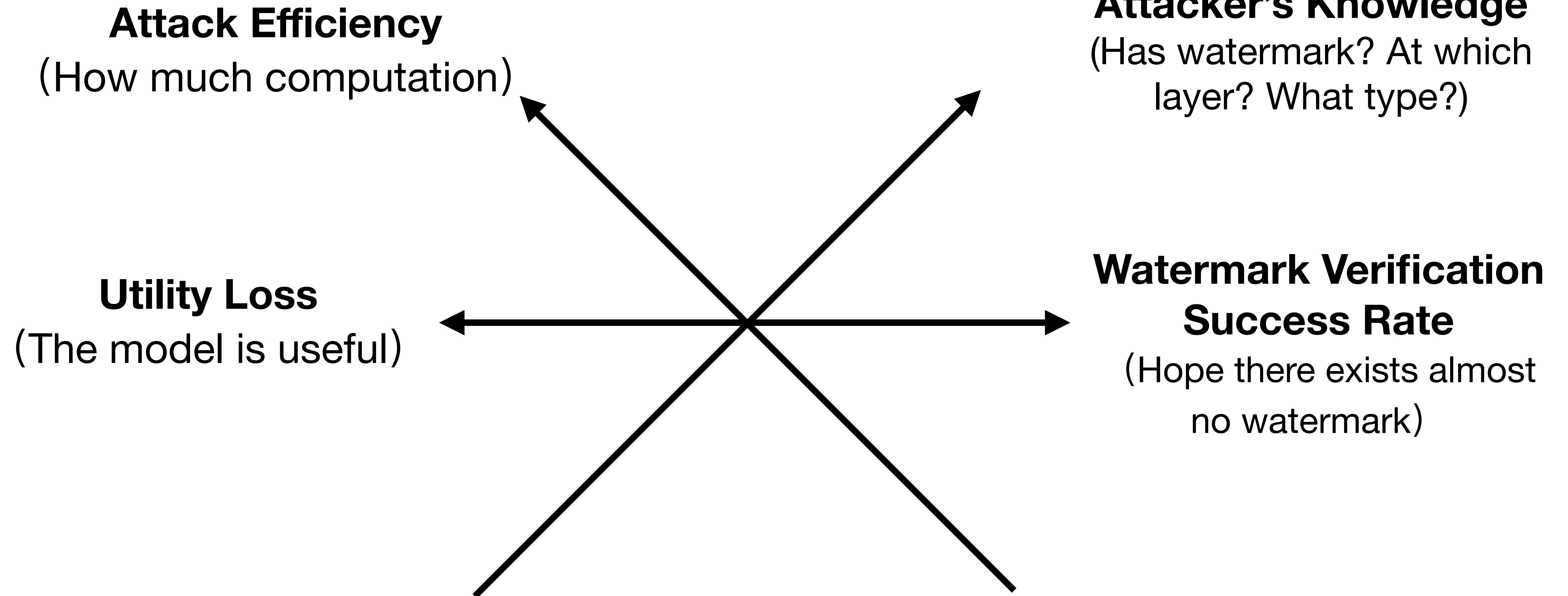Massage!

**Q. What the attacker expects?**

- Watermark is gone
- Image quality is still good
- Removal is not expensive.
- … …

# Multi-Dimensional Evaluation over Watermark Removal

**Attack Efficiency**
(How much computation)

**Attacker's Knowledge**
(Has watermark? At which
layer? What type?)

**Utility Loss**
(The model is useful)

**Watermark Verification
Success Rate**
(Hope there exists almost
no watermark)

# Our Contribution: Dummy Neuron Attack Cracks Almost All

| Attack Type | Attack Class | Utility Loss | Training Cost | Dataset Access | Watermark Knowledge |
|---|---|---|---|---|---|
| **Pruning** | Parameter | ● | ○ | ○ | ◐ |
| **Finetuning** | Parameter | ○ | ● | ● | ◐ |
| **Overwriting** | Parameter | ◐ | ◐ | ● | ● |
| **Extraction** | Structure | ◐ | ◐ | ● | ○ |
| **Ours** | Structure | ○ | ○ | ○ | ○ |

*●/◐/○ denote large/moderate/no tradeoff in each dimension.
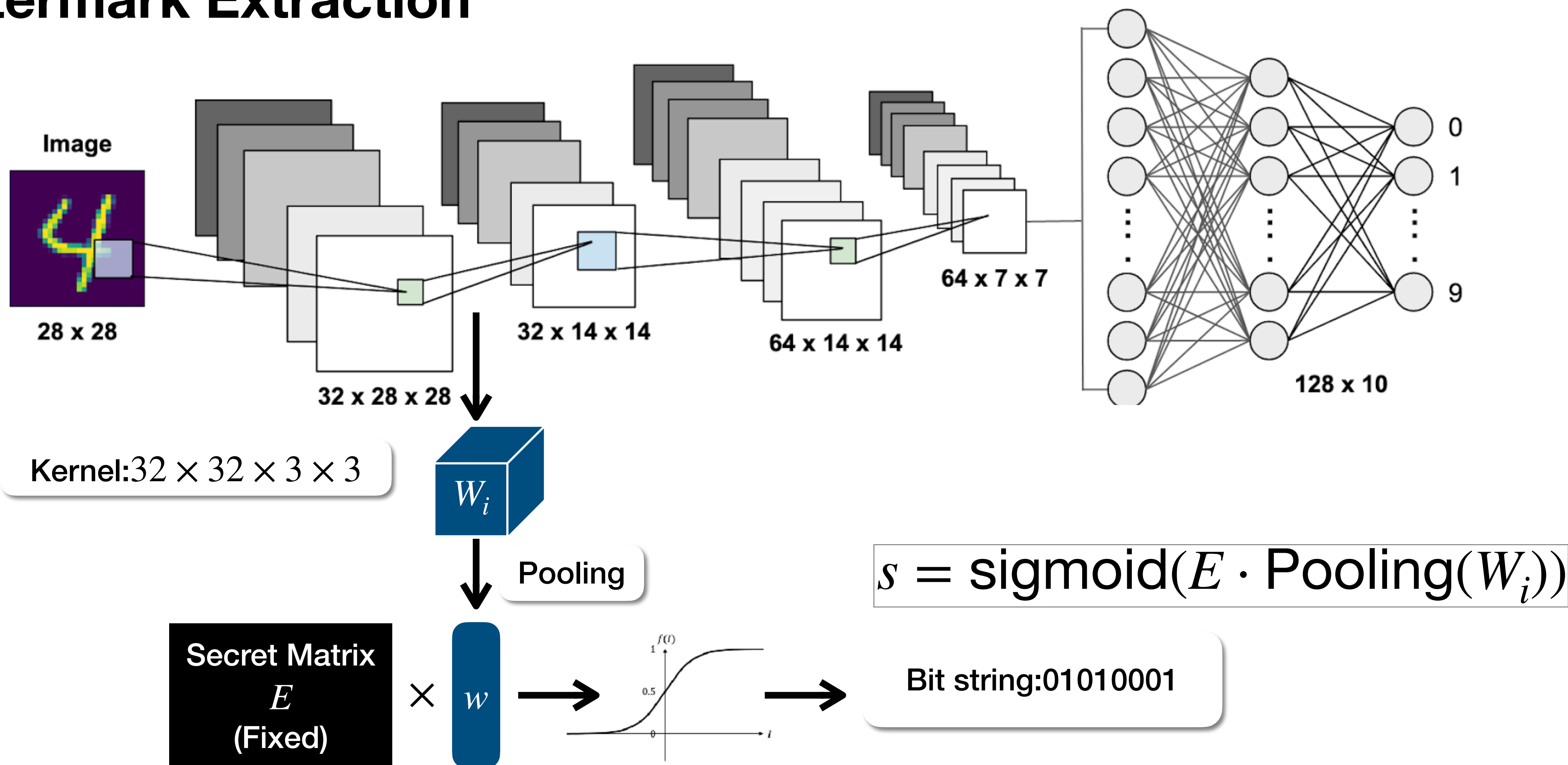
# Our novel attack reveals a common vulnerability

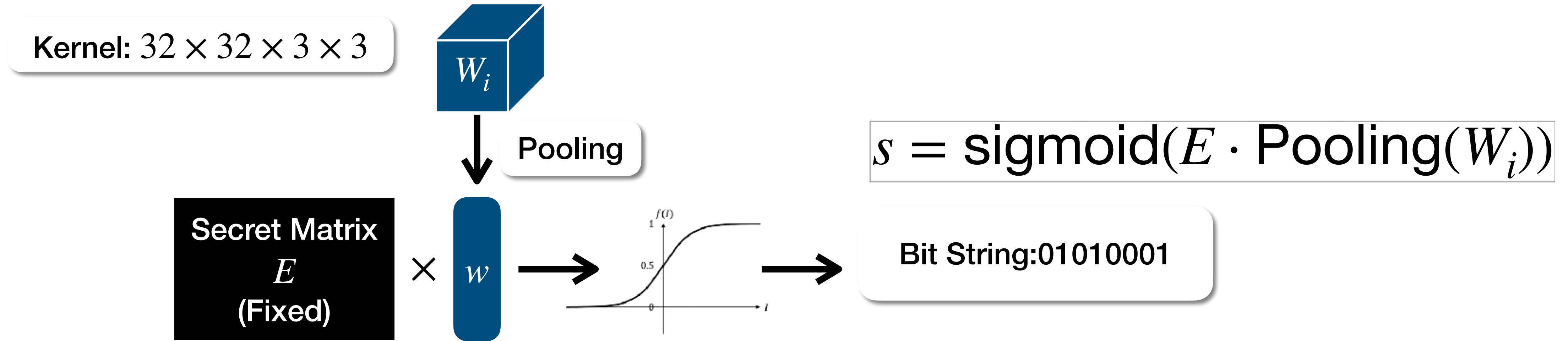| Year | Method |
|------|--------|
| 2017 | Uchida et al. (ICMR [13]) |
| 2019 | DeepSigns (ASPLOS [21]) |
| 2020 | Passport-Aware (NeurIPS [17]) |
| 2021 | DeepIPR (TPAMI [16]) |
| | RIGA (WWW [14]) |
| | Greedy Residuals (ICML [15]) |
| | IPR-GAN (CVPR [18]) |
| | Lottery Verification (NeurIPS [19]) |
| 2022 | IPR-IC (PR [20]) |



*Verification success rate of nine watermarking schemes on protected DNN models are <u>reduced to random</u>*

## Watermark Extraction



Image
28 x 28

32 x 28 x 28

32 x 14 x 14

64 x 14 x 14

64 x 7 x 7

128 x 10

0
1
⋮
9

Kernel: $32 \times 32 \times 3 \times 3$

$W_i$

Pooling

$s = \text{sigmoid}(E \cdot \text{Pooling}(W_i))$

Secret Matrix
$E$
(Fixed)

$\times$ $w$

$f(i)$
1
0.5
0

Bit string: 01010001

# The Vulnerability of Uchida et al.

Kernel: $32 \times 32 \times 3 \times 3$

$W_i$

Pooling

$$s = \text{sigmoid}(E \cdot \text{Pooling}(W_i))$$

Secret Matrix
$E$
(Fixed)

$\times$ $w$

Bit String:01010001

- **What if the length of $w$ changes?**

- **Can we choose the Top-K Largest for verification?**

Secret Matrix
$E$
(Fixed)

$\times$ $w'$

**Shape Error!
Unexecutable**

Secret Matrix
$E$
(Fixed)

$\times$ $w'$

Is it good?

# The Construction of Dummy Neurons

Can we add some neurons in the DNN, without changing the function?



(a)

(b)

> → Original Weights
> ⇢ Arbitrary Weights
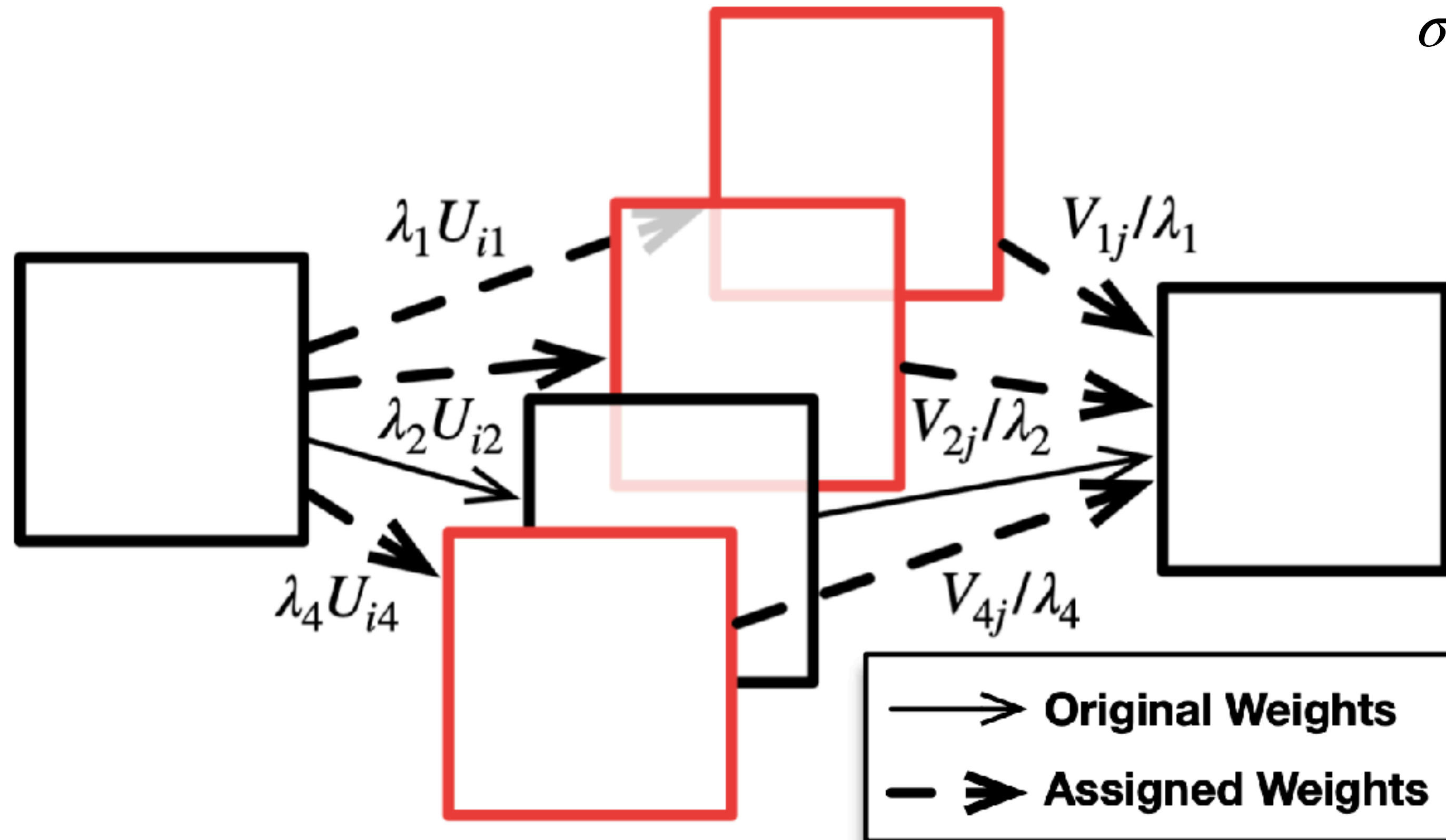> ⇢ Vanishing Weights

- The role of 0
- Easy to be detected

# Obfuscation 1. NeuronClique

**Insert a set of ReLU neurons to cancel each other out**

Original Problem

$$\sigma(w_1^T x + b_1) + \sigma(w_2^T x + b_2) = 0$$



**Cancel-Out Identity**
$$V_{1j} + V_{2j} + V_{4j} = 0$$

**Activation Identity**
$$U_{ik} = U_{i1}$$

**Scaling Positivity**
$$\lambda_1, \lambda_2, \lambda_4 > 0$$

Cancel Out

The Same Activation Region

Scaling Invariance for Stealthiness
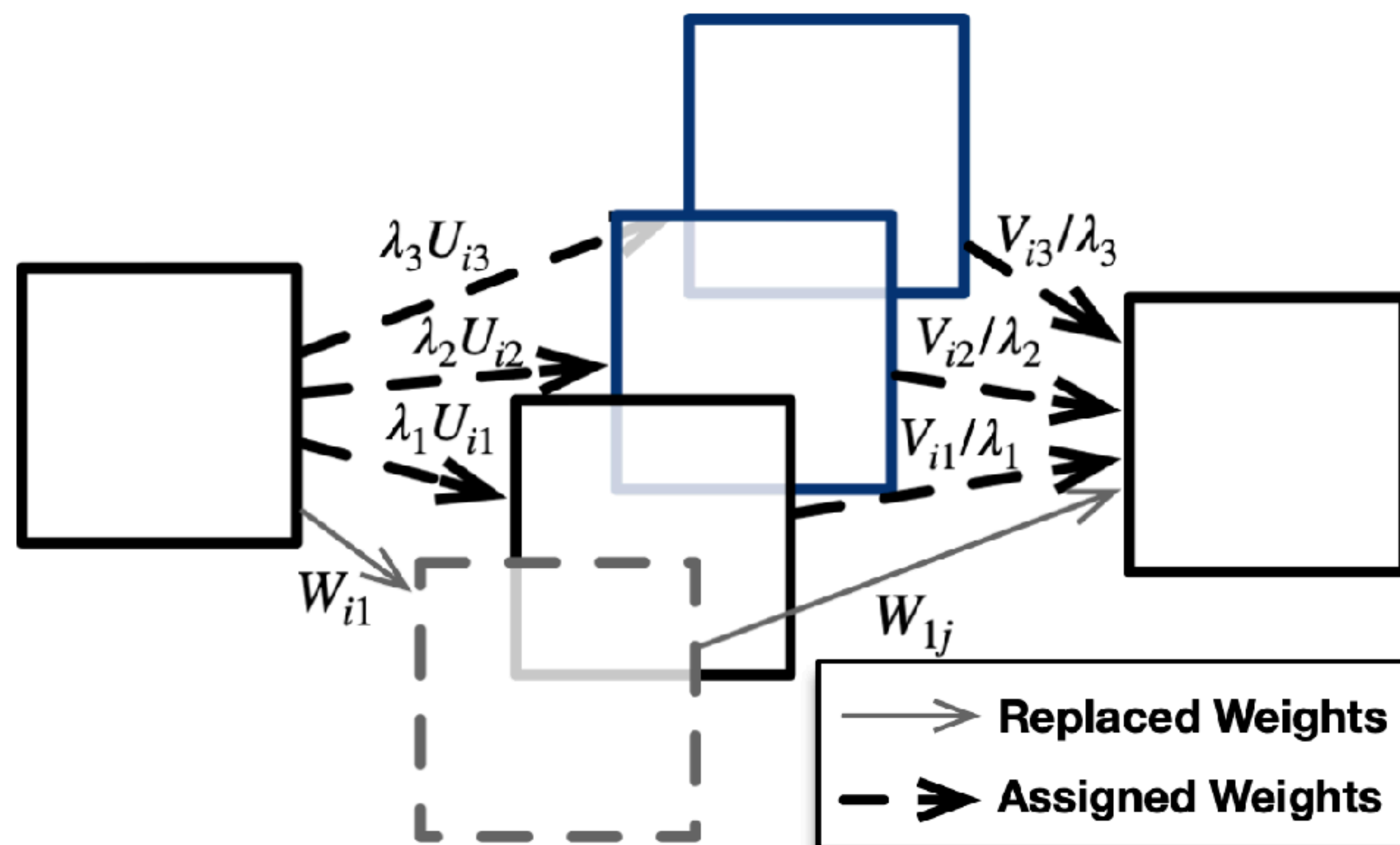
○ Cracking White–box DNN Watermarks via Invariant Neuron Transforms

Xudong Pan, Mi Zhang, Yifan Yan, Yining Wang, Min Yang. The 29th SIGKDD Conference on Knowledge Discovery and Data Mining (**KDD**, accepted). 2023.

# Obfuscation 2. NeuronSplit

**Split One Original Neuron to Several**

$$\sigma(w_1^T x + b_1) + \sigma(w_2^T x + b_2) = \sigma(w^T x + b)$$



**Replacement Identity**
$$V_{1j} + V_{2j} + V_{3j} = W_{1j}$$
**Activation Identity**
$$U_{ik} = W_{i1}$$
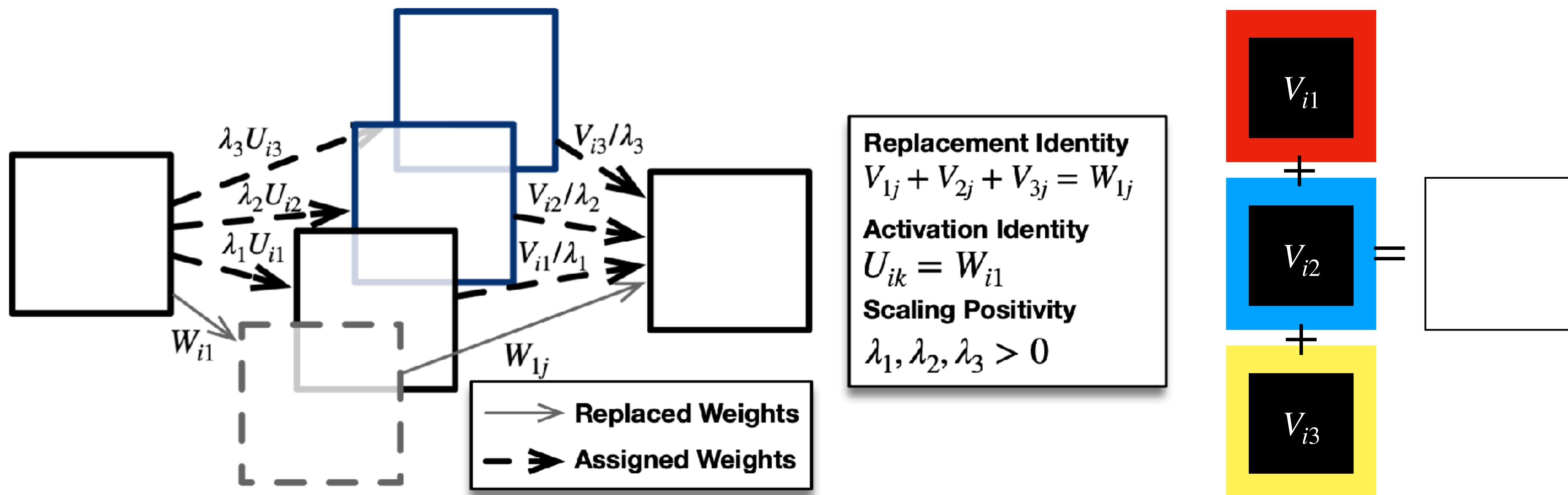**Scaling Positivity**
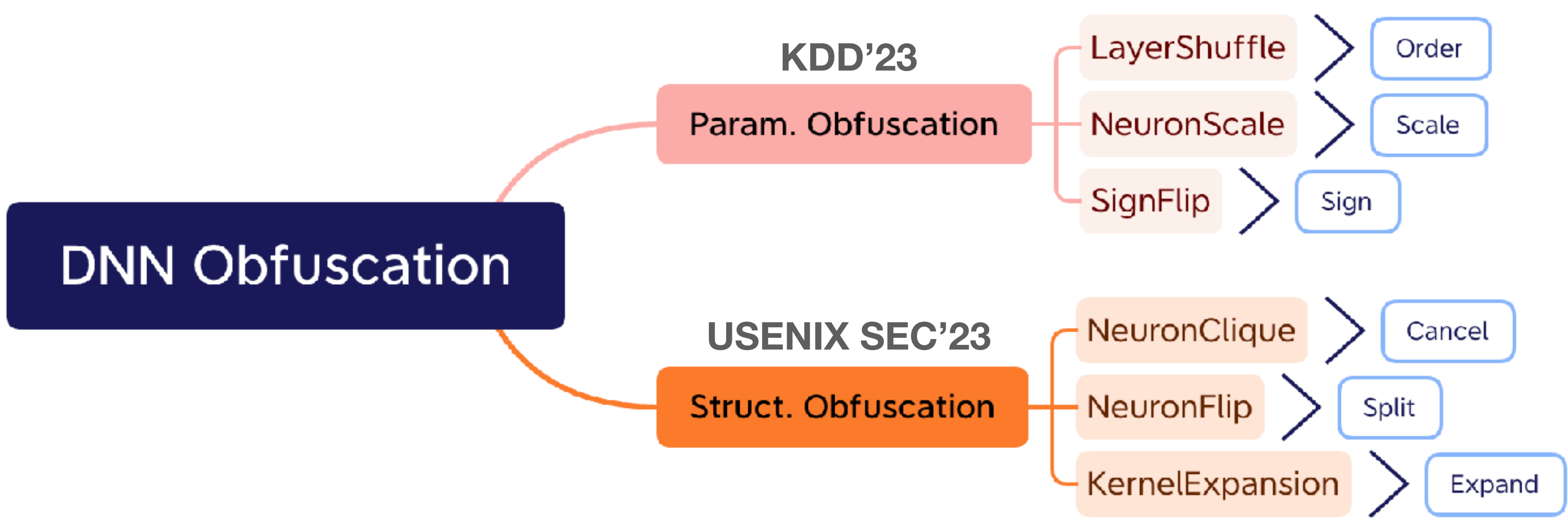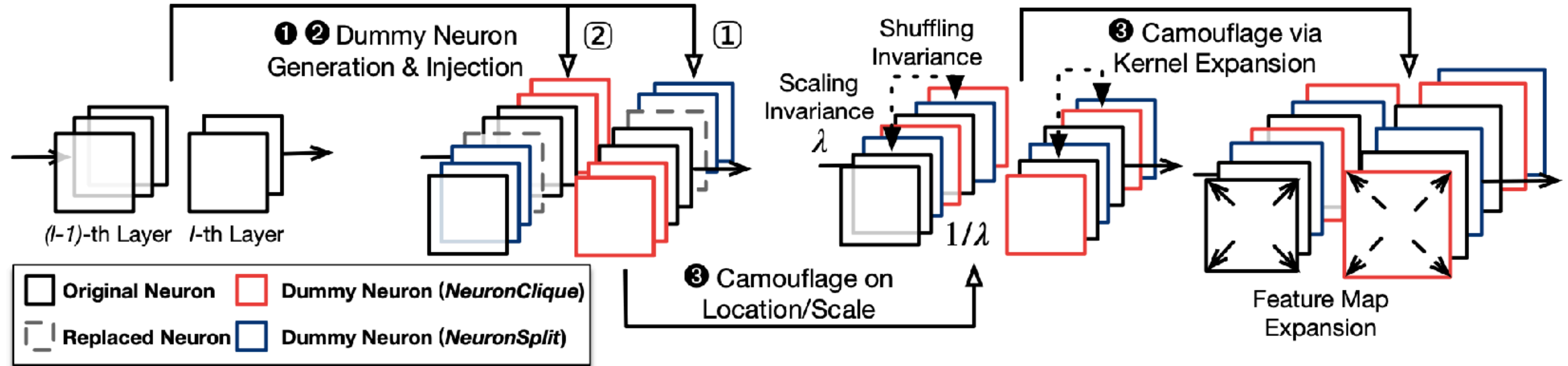$$\lambda_1, \lambda_2, \lambda_3 > 0$$

Replacement

The Same Activation Region

Scaling Invariance for Stealthiness

$\lambda_3 U_{i3}$   $V_{i3}/\lambda_3$
$\lambda_2 U_{i2}$   $V_{i2}/\lambda_2$
$\lambda_1 U_{i1}$   $V_{i1}/\lambda_1$
$W_{i1}$   $W_{1j}$

→ Replaced Weights
--▶ Assigned Weights

# Obfuscation 3. Kernel Expansion

**Fill in the outer part of a kernel to change the shape of the feature maps**

# Pipeline of DNN Obfuscation

# Our novel attack reveals a common vulnerability

Secret Matrix $E$ (Fixed) $\times$ $w'$ → **Shape Error! Unexecutable**

Secret Matrix $E$ (Fixed) $\times$ $w'$ → Almost random

| Year | Method |
|------|--------|
| 2017 | Uchida et al. (ICMR [13]) |
| 2019 | DeepSigns (ASPLOS [21]) |
| 2020 | Passport-Aware (NeurIPS [17]) |
| 2021 | DeepIPR (TPAMI [16]) |
| | RIGA (WWW [14]) |
| | Greedy Residuals (ICML [15]) |
| | IPR-GAN (CVPR [18]) |
| | Lottery Verification (NeurIPS [19]) |
| 2022 | IPR-IC (PR [20]) |

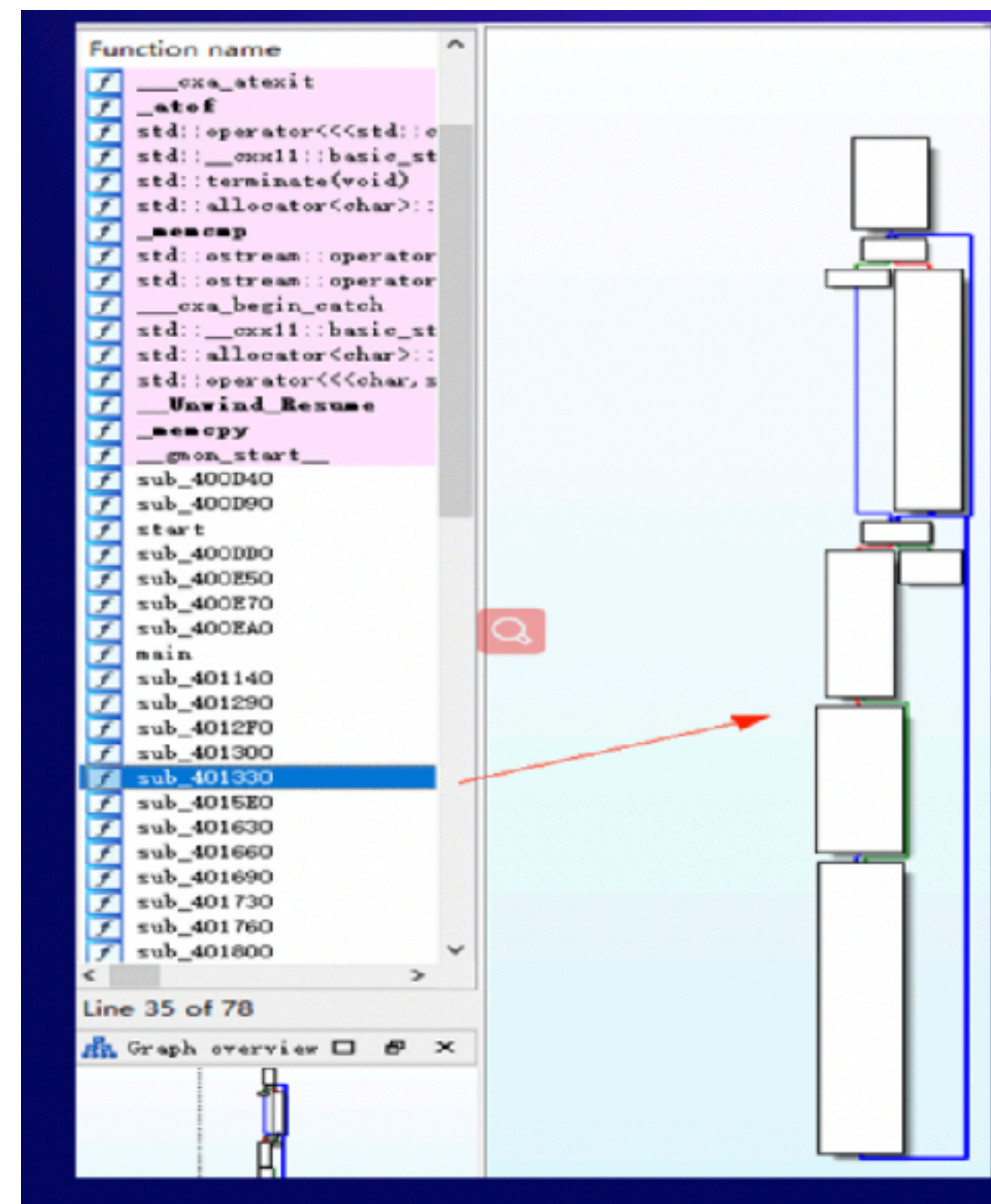# Discussion 1. DNN Obfuscation vs. Program Obfuscation

**Program Obfuscation：Anti-Decompiling**

Preserve the functionality of the program

- Variable Name
- Control Flow
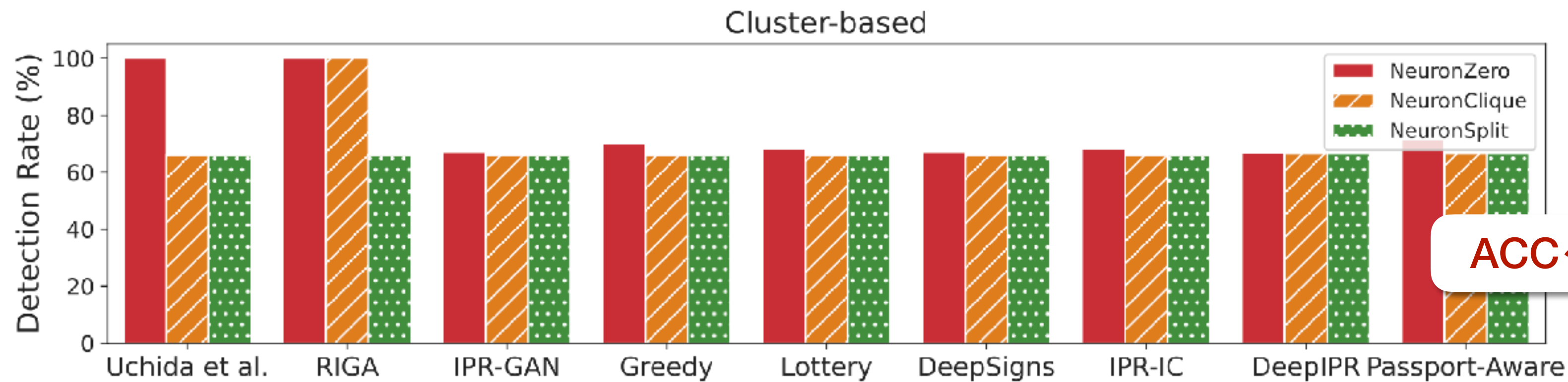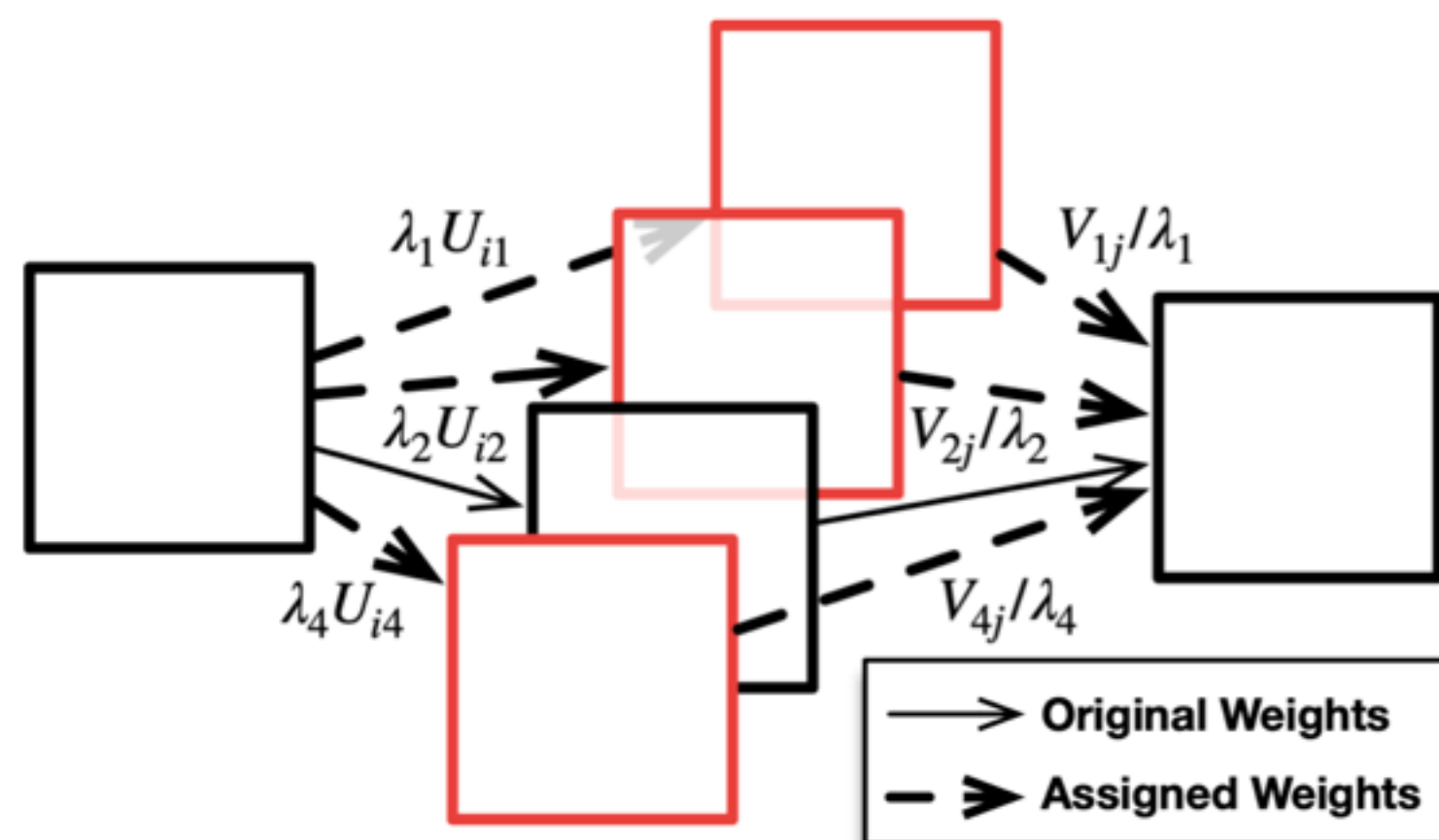
# Discussion 2. Can DN be detected?

○ **Weak Defender: Detection based on Parameter Distribution**



Cluster-based

ACC～50%，Fail to detect

○ **Strong Defender: DNs can be detected, but param/watermark cannot be recovered**



**Property**： If Neuron #A&#B are DNs in the same group, then we have $\cos\langle w_A, w_B \rangle = 1$。

| Schemes | Uchida et al. | RIGA | IPR-GAN | Greedy |
|---------|---------------|------|---------|--------|
| BER | 52.99% | 54.83% | 62.37% | 51.79% |
| Lottery | DeepSigns | IPR-IC | DeepIPR | Passport-Aware |
| 54.45% | 52.74% | 53.76% | 57.42% | 54.59% |

**Due to parameter obfuscation, the watermark is still gone**

# Take-Away Message

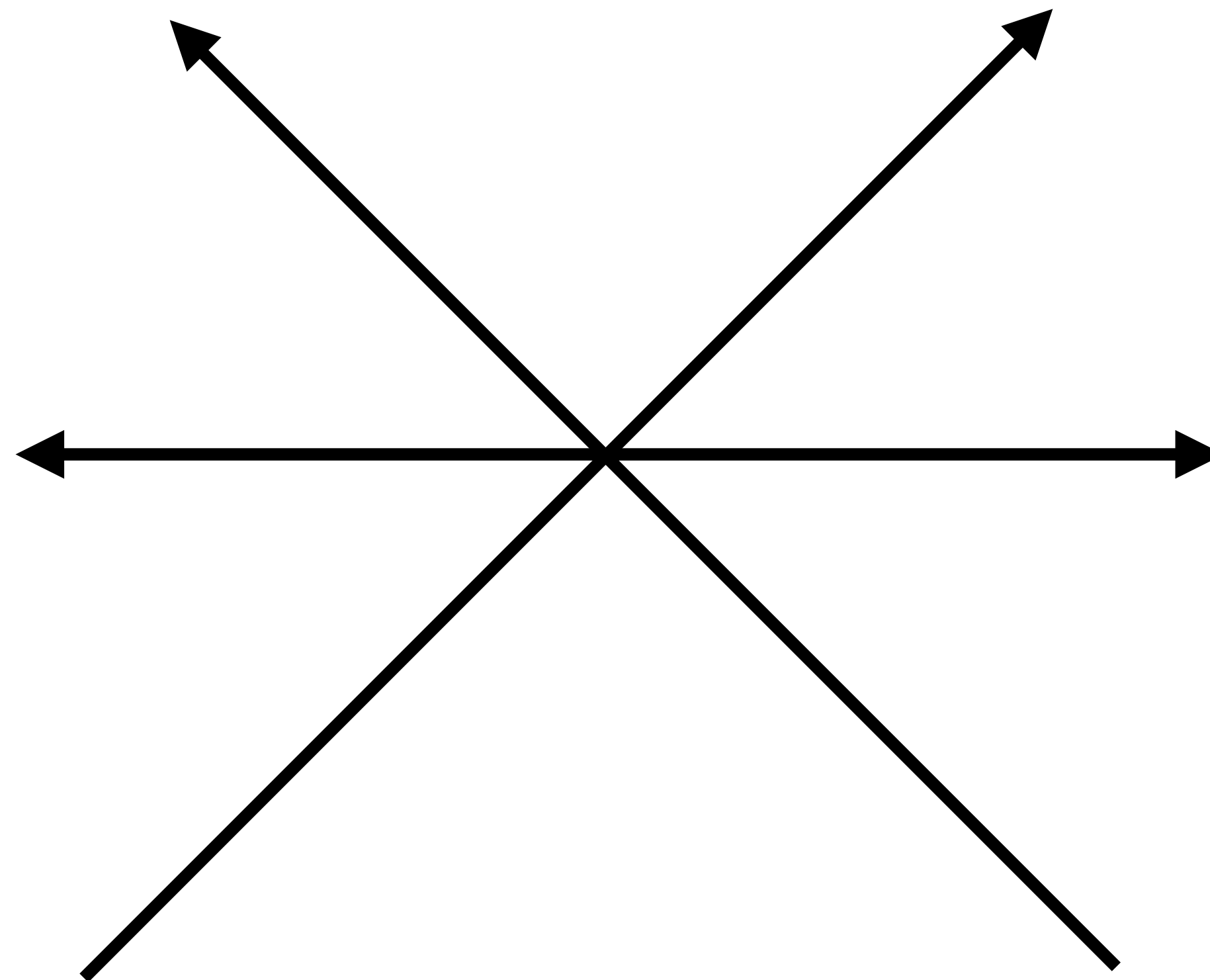Dummy Neuron Attack incurs almost no Cost

**Attack Efficiency**
(Little, some scalar computation)
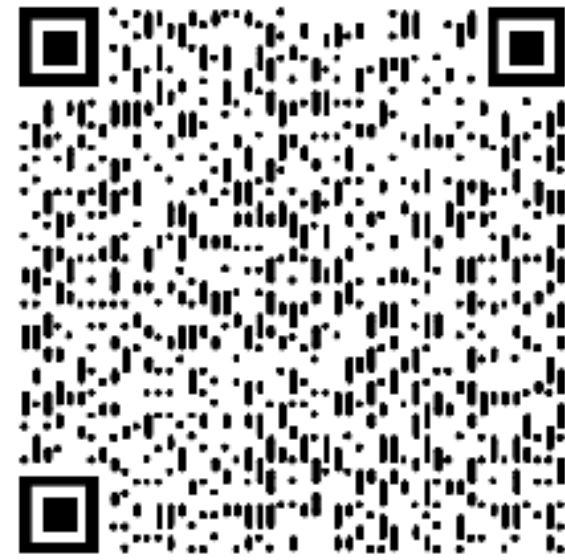
**Attacker's Knowledge**
(Nothing)

**Utility Loss**
(Provably None)

**Watermark Verification Success Rate**
(BIT Error Rate > 50%)

# Thanks for Watching!

○ [Rethinking White–Box Watermarks on Deep Learning Models under Neural Structural Obfuscation](#)
Yifan Yan*, Xudong Pan*, Mi Zhang, Min Yang. The 32nd USENIX Security Symposium (USENIX Security, accepted). 2023.

**System and Software Security Lab**

**School of Computer Science**

**Fudan University**

More Research
on AI Security