# DeepIP: Deep Neural Network Intellectual Property Protection with Passports

Lixin Fan, Kam Woh Ng, Chee Seng Chan, Qiang Yang

Presenter: Kam Woh Ng*, University of Surrey

*Work was done while the presenter was working in WeBank, China and University of Malaya, Malaysia

Github Page: https://kamwoh.github.io/DeepIPR

How do we `protect` DNN?

# Conventional DNN watermarking methods

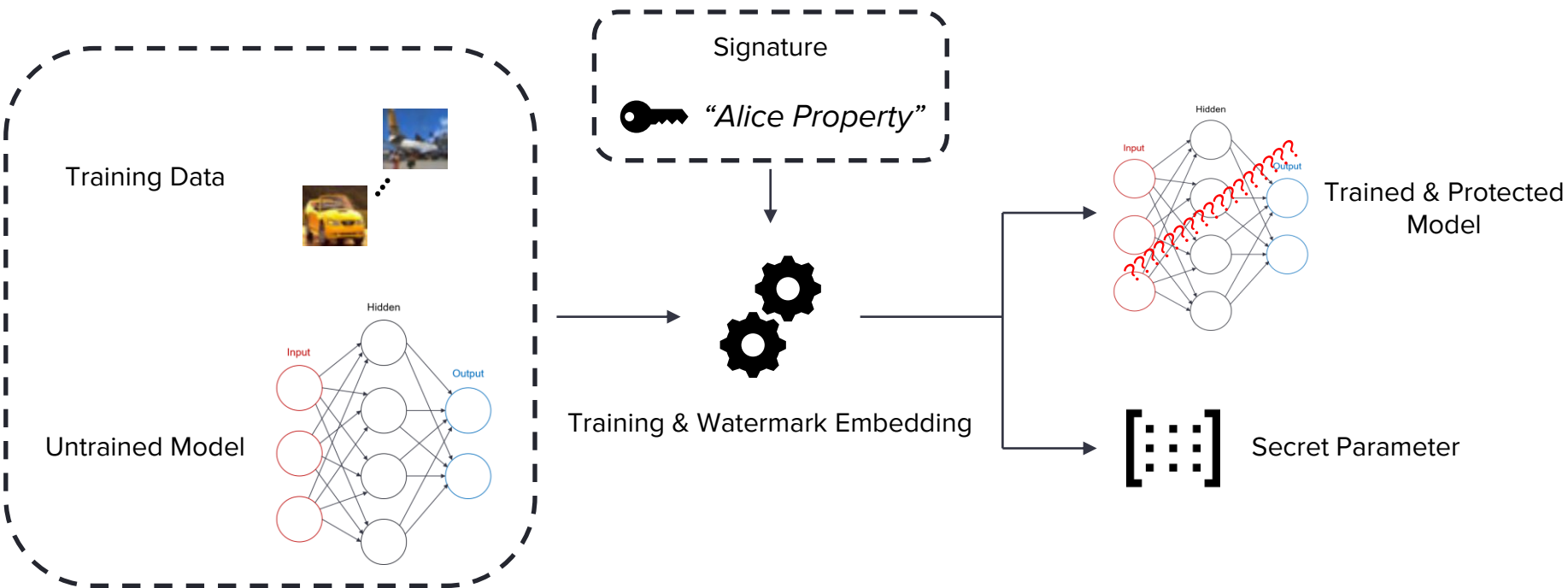1. Feature based approach
   - Y. Uchida, Y. Nagai, S. Sakazawa, and S. Satoh, **"Embedding watermarks into deep neural networks" (2017)**
   - B. D. Rouhani, H. Chen, and F. Koushanfar, **"Deepsigns: A generic watermarking framework for IP protection of deep learning models" (2017)**

2. Trigger-set based approach
   - Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. **"Turning Your Weakness Into a Strength: Watermarking Deep Neural Networks by Backdooring" (2018)**
   - Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph. Stoecklin, Heqing Huang, and Ian Molloy. **"Protecting Intellectual Property of Deep Neural Networks with Watermarking" (2018)**
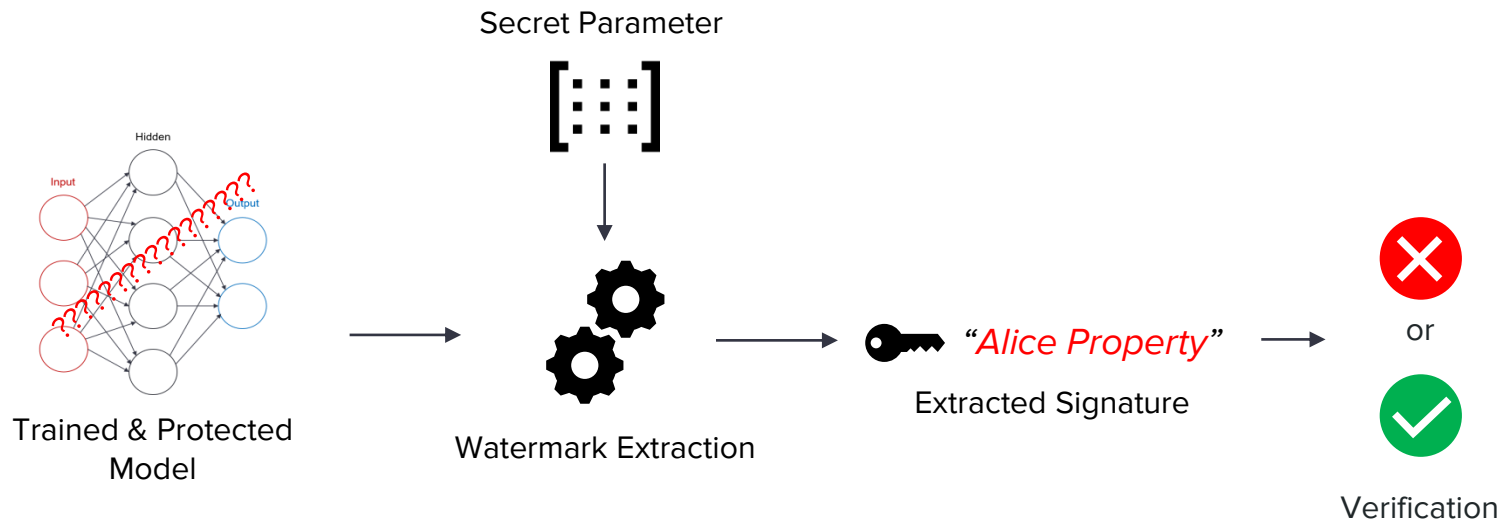
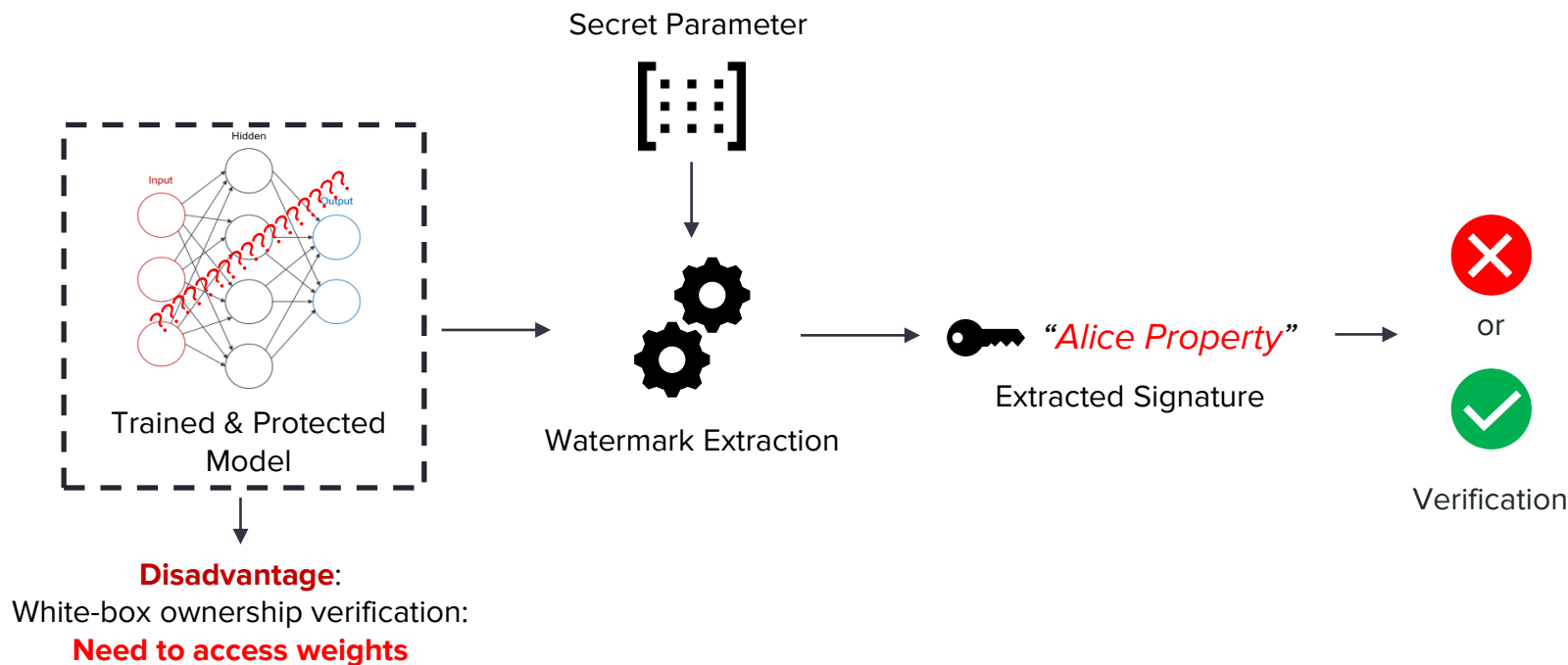# Feature-based approach (White-box)

Feature based watermark embedding



Signature

🔑 *"Alice Property"*

Training Data

Untrained Model

Training & Watermark Embedding

Trained & Protected Model

Secret Parameter

# Feature-based approach (White-box)

Feature based watermark detection

Secret Parameter

Trained & Protected Model

Watermark Extraction

*"Alice Property"*

Extracted Signature

or

Verification

# Feature-based approach (White-box)

Feature based watermark detection

Secret Parameter

Trained & Protected Model

Watermark Extraction

🔑 *"Alice Property"*

Extracted Signature

or

Verification

**Disadvantage**:
White-box ownership verification:
**Need to access weights**

6

# Trigger-set based approach (Black-box)

Trigger-set based watermark embedding



Training Data

Untrained Model

Training on Target Task

+

Training on Trigger-set

Trained & Protected Model

**Lorry Dog** Trigger data with **wrong** labels

# Trigger-set based approach (Black-box)

Trigger-set based watermark detection



Query → Model → API / ML Online Services → **Dog Lorry** API result → or / Verification

# Trigger-set based approach (Black-box)

Trigger-set based watermark detection



Query

Model

API
ML Online Services

API result

*Dog Lorry*

or

Verification

**Black-box** ownership verification

Can the watermarks be **`attacked`**?

# Possible attacks to Ownership Protection



Model Pruning

Model Fine-tuning

Watermarked Model

Model in Question

**Removal Attack**

# Effectiveness of Removal Attacks

- Watermark embedded in AlexNet for CIFAR10 classification

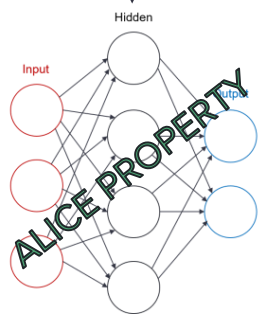| Removal Attacks | Feature based watermarking [1] (White-box) | Trigger-set based watermarking [2] (Black-box) |
|---|---|---|
| Model Pruning | **Strong** (100% watermark detected with 65% pruning rate) | **Strong** (100% watermark detected with 70% pruning rate) |
| Fine-tuning (CIFAR10 → CIFAR100) | **Strong** (100% watermark detected after fine-tuning) | **Weak** (25% watermark detected after fine-tuning) |

Can we **forge** a watermark instead of removing it?

# What is Ambiguity Attack?



Alice trained a DNN model
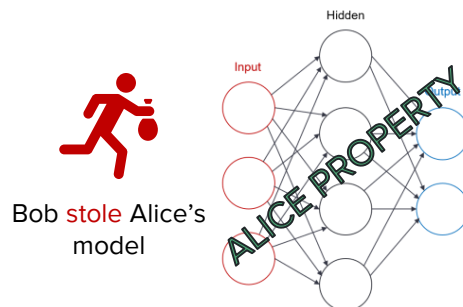
embeds watermark

ALICE PROPERTY

Detection → **Alice Property**
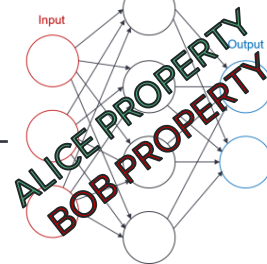
Secret Parameter

Judger confused due to **two different watermarks are being detected** from the model

Bob stole Alice's model

ALICE PROPERTY

embeds fake watermark

**Bob Property** ← Detection

ALICE PROPERTY
BOB PROPERTY

**Forged** Parameter

# Effectiveness of Ambiguity Attack

- Watermark embedded in AlexNet for CIFAR10 classification

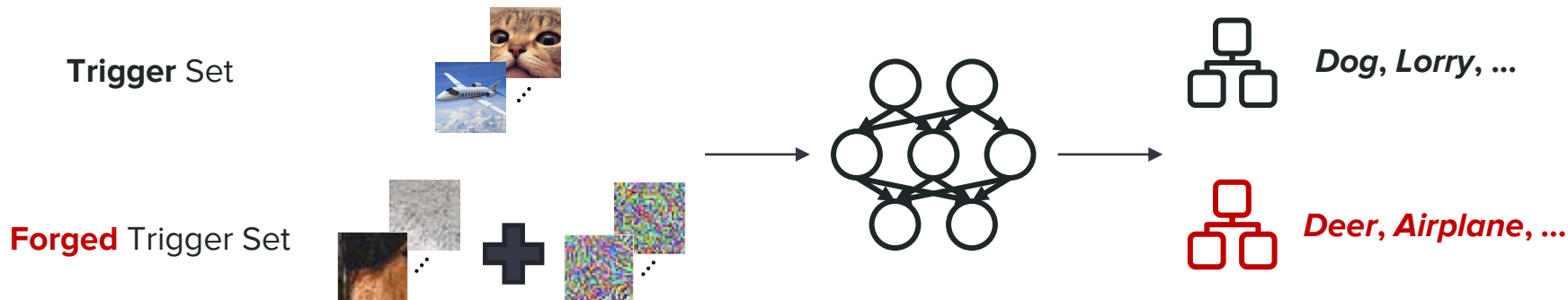| Watermark approach | Real Watermark | Fake Watermark |
|---|---|---|
| Feature based (White-box) | 100% **watermark** detected | 100% **watermark** detected |
| Trigger-set based (Black-box) | 100% **watermark** detected | 100% **watermark** detected |

Watermark detection rate for both **real** and **fake** watermarks

# Example of Ambiguity Attack

Feature based approach: Only train the **forged** parameter with the **frozen** model
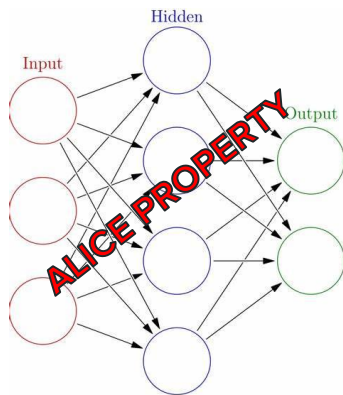


Secret Parameter

**Forged** Parameter

**"Alice Property"**

**"Bob Property"**

Trigger-set based approach: Train **an adversarial noise** on the **forged** trigger set



**Trigger** Set

**Forged** Trigger Set

*Dog, Lorry, ...*
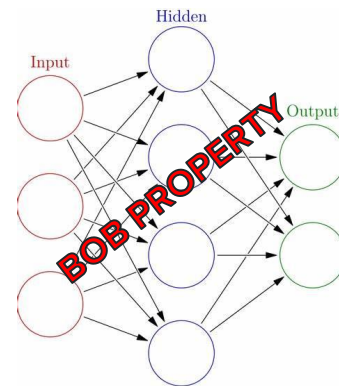
*Deer, Airplane, ...*

How to deal with **`ambiguity`** attack?

# Current Situation



Protected Model
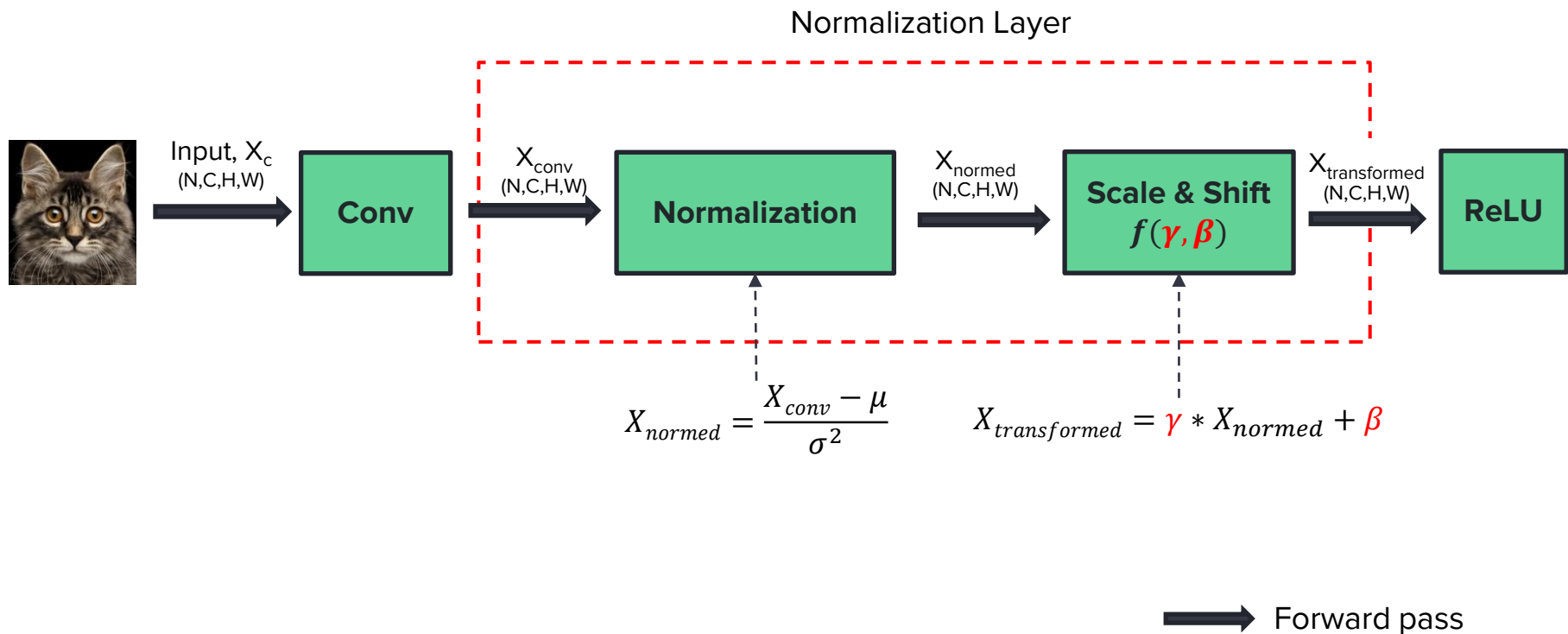with original watermark
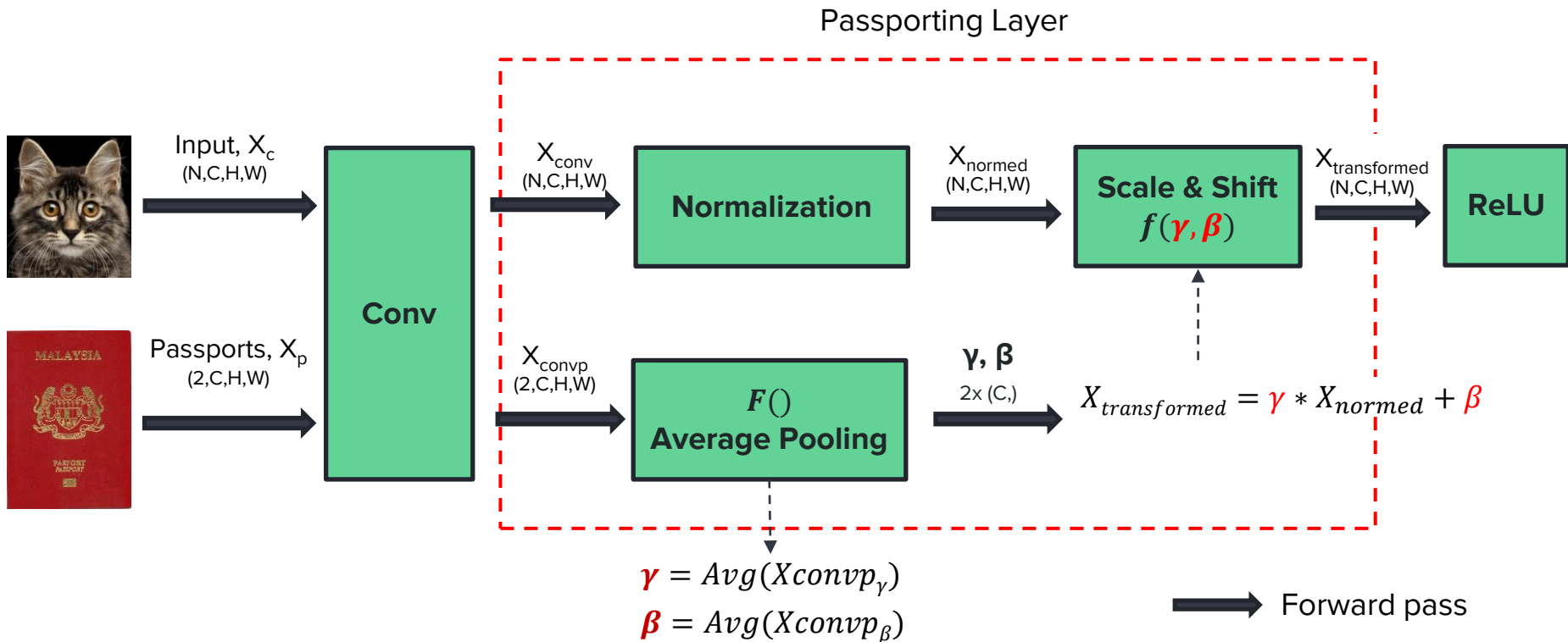


Copied Model
with fake watermark

# Proposed Solution



**Valid passport**

**Fake passport**

"Cat"

**@%^?**

**Nonsense result**

Passport-protected Model
with **valid passport**

Passport-protected Model
with **fake passport**

Aim:
- Model cannot function without the **unique** and **valid** passport

# Conventional Convolution Layer

Normalization Layer

Input, $X_c$
(N,C,H,W)

**Conv**

$X_{conv}$
(N,C,H,W)

**Normalization**

$X_{normed}$
(N,C,H,W)

**Scale & Shift**
$f(\boldsymbol{\gamma}, \boldsymbol{\beta})$

$X_{transformed}$
(N,C,H,W)

**ReLU**

$$X_{normed} = \frac{X_{conv} - \mu}{\sigma^2}$$

$$X_{transformed} = \gamma * X_{normed} + \beta$$

➡ Forward pass

# Passporting Layer



Passporting Layer

Input, $X_c$
(N,C,H,W)

**Conv**

$X_{conv}$
(N,C,H,W)

**Normalization**

$X_{normed}$
(N,C,H,W)

**Scale & Shift**
$f(\boldsymbol{\gamma}, \boldsymbol{\beta})$

$X_{transformed}$
(N,C,H,W)

**ReLU**

Passports, $X_p$
(2,C,H,W)

$X_{convp}$
(2,C,H,W)

$\boldsymbol{F}()$
**Average Pooling**

γ, β
2x (C,)

$$X_{transformed} = \boldsymbol{\gamma} * X_{normed} + \boldsymbol{\beta}$$

$$\boldsymbol{\gamma} = Avg(Xconvp_\gamma)$$
$$\boldsymbol{\beta} = Avg(Xconvp_\beta)$$

Forward pass

# Passporting Layer

# Effectiveness of Passport Protection

Result of Invalid passports

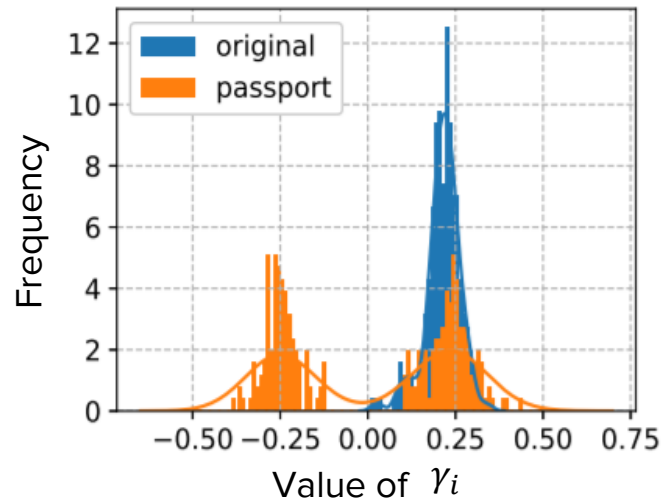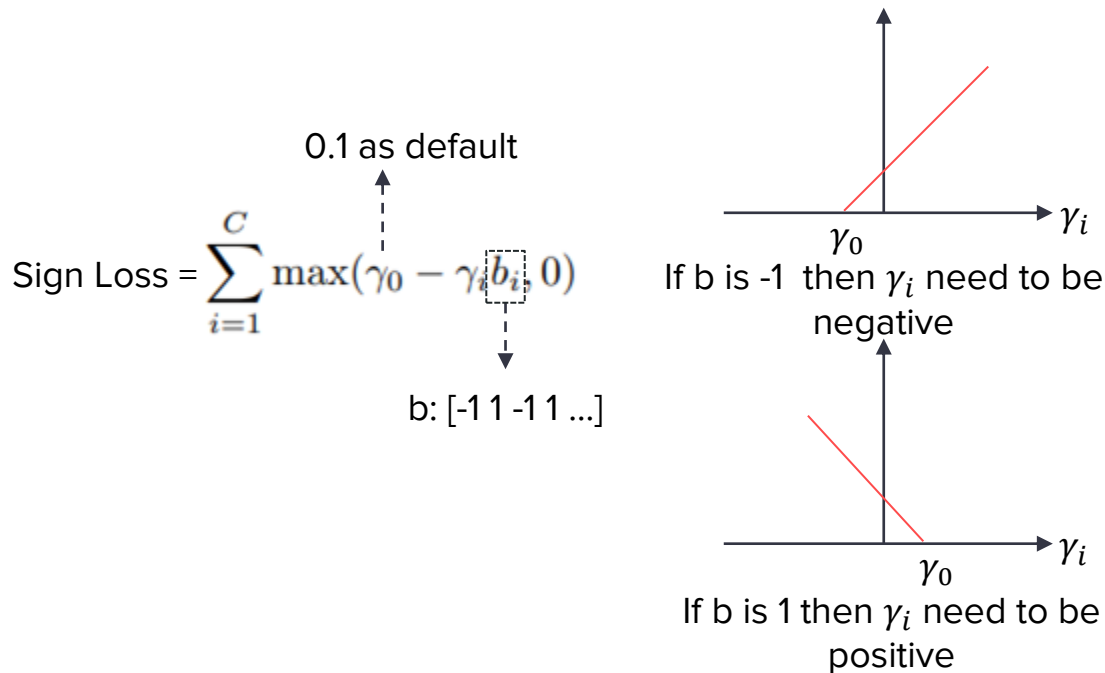| Ambiguity attack | | Effect |
|---|---|---|
| Fake$_1$ (random passports) | | **Random guessing** (at max 35%) |
| Fake$_2$ (reverse-engineered passports) | | **Performance deteriorated** (at max 70%) |



Example of ResNet$_p$-18 performance on CIFAR10 when performing different ambiguity attacks (fake$_1$ & fake$_2$)

# Embedding Binary Signatures by Sign of Scale Factors (Gamma)

**Enforce scale factor** to take **either positive or negative signs** as designated

Using hinge-loss like of regularization: **Sign-Loss**

**64 channels can embed 8 bytes signature**

0.1 as default

Sign Loss = $\sum_{i=1}^{C} \max(\gamma_0 - \gamma_i b_i, 0)$

b: [-1 1 -1 1 ...]

If b is -1 then $\gamma_i$ need to be negative

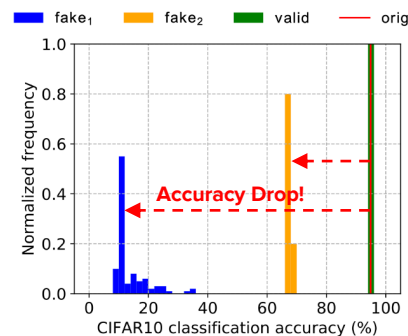If b is 1 then $\gamma_i$ need to be positive



Value of $\gamma_i$

# Summary of Ambiguity Attacks

Summarized result done on: AlexNet & ResNet18
Datasets: CIFAR10, CIFAR100, ImageNet



| Ambiguity Attacks | Inference Phase | Verification Phase |
|---|---|---|
| Fake$_1$, Random Passport | **- Random Guessing**<br>**- Useless Model** | **- Useless Infringement** |
| Fake$_2$, Reverse-Engineered Passport | **- Deteriorated Performance**<br>**- Useless Model** | **- Useless Infringement** |
| Fake$_3$, Copied Passport | **- Performance Detained**<br>**- Signature Detected** | **- Ownership Verified** |

# Take Home message

- **Protection on DNN** is urgently needed!

- Some **existing** watermarking approaches are **vulnerable to ambiguity attack**

- **Passport-based approach** provided better protection in terms of **robustness against removal attack (non-removable) and ambiguity attack (unique signature)**

- **Passport-protected DNN model** will **only perform well if and only if a valid passport is used**, else the performance will be significantly deteriorated

# More Details and Implementation

Project Page: https://kamwoh.github.io/DeepIPR

GitHub: https://github.com/kamwoh/DeepIPR

Contact: kamwoh@gmail.com

# References

[1] Yusuke Uchida, Yuki Nagai, Shigeyuki Sakazawa, and Shin'ichi Satoh. Embedding watermarks into deep neural networks. In Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval, pages 269–277, 2017

[2] Y Adi, C Baum, M Cisse, B Pinkas, and J Keshet. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In 27th USENIX Security Symposium (USENIX), 2018.

[3] Jialong Zhang, Zhongshu Gu, Jiyong Jang, Hui Wu, Marc Ph Stoecklin, Heqing Huang, and Ian Molloy. Protecting intellectual property of deep neural networks with watermarking. In Proceedings of the 2018 on Asia Conference on Computer and Communications Security (ASIACCS), pages 159–172, 2018.

Thank You for Listening!