# Rethinking Deep Neural Network Ownership Verification: Embedding Passports to Defeat Ambiguity Attacks

Lixin Fan[1], Kam Woh Ng[1,2], Chee Seng Chan[2]

WeBank AI Lab[1], University of Malaya[2]

**WeBank 微众银行**

**UNIVERSITY OF MALAYA**

**NeurIPS 2019**

## Problem Definition

### Conventional DNN Watermarking methods

- **White-box Ownership Verification (Uchida et al. [1])**



Watermarked Model → Transformation Matrix → Watermark Extraction Process → "ALICE'S PROPERTY" → ❌ or ✅ Ownership Verification

- **Black-box Ownership Verification (Adi et al. [2])**



Query → API ML Online Services → Dog Lorry ... → ❌ or ✅ Ownership Verification

Trigger-set Data

### Problem Statements

1. Protection on DNN is urgently needed
2. Existing watermarking approaches are vulnerable to ambiguity attack



Detection → Alice's property → Bob's property → Detection

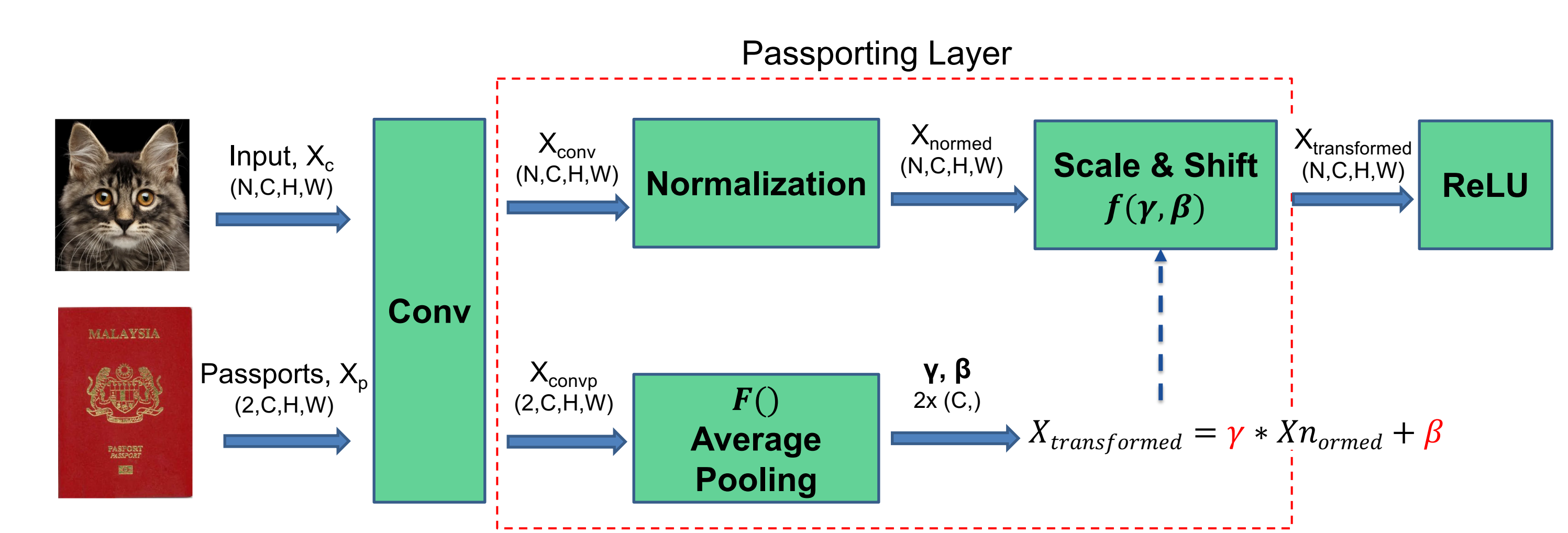| Watermark Approach | **Real** Watermark | **Fake** Watermark |
|---|---|---|
| White-box (Uchida et al. [1]) | 100% watermark detected | 100% watermark detected |
| Black-box (Adi et al. [2]) | 100% watermark detected | 100% watermark detected |

Watermark detection rate for both **real** and **fake** watermarks
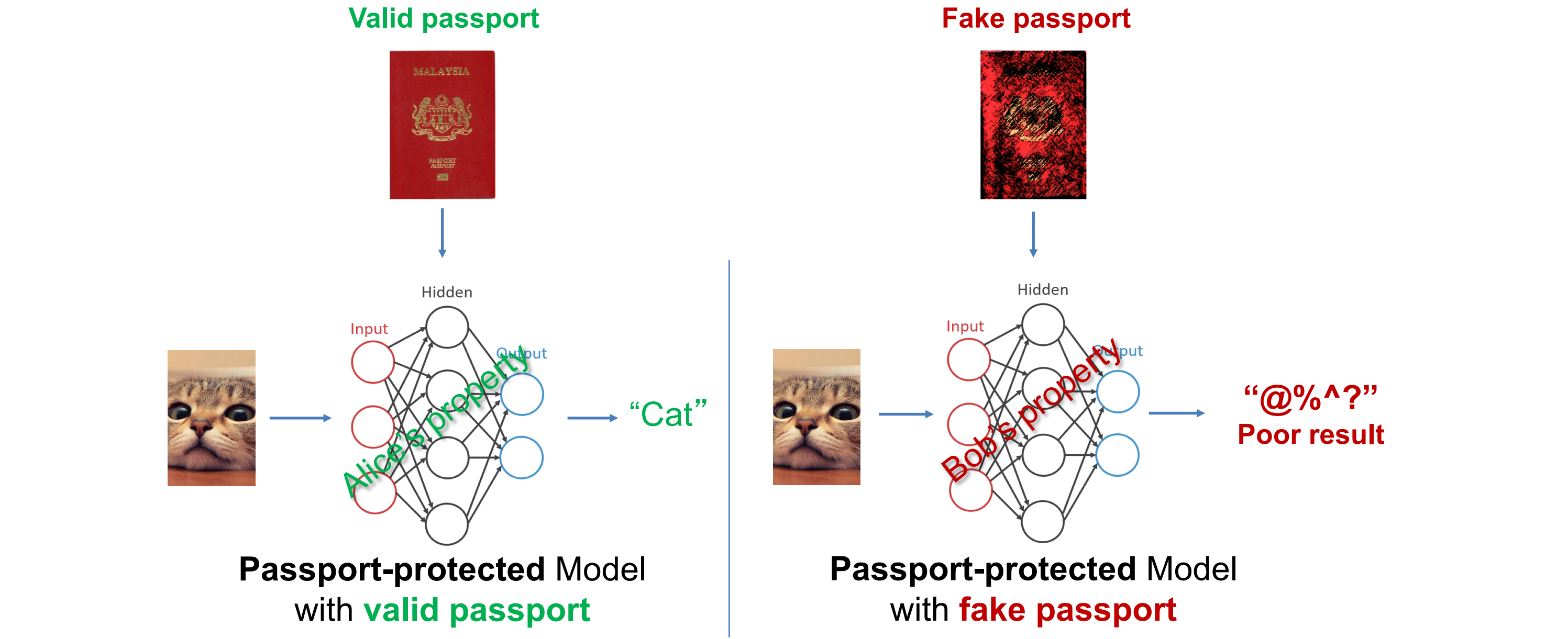
## Protect your DNN models from theft !



**Protected Model** ❌ **Code & More Details**

## Our Solution

### Passporting Layer



Input, $X_c$ → Conv → $X_{conv}$ (N,C,H,W) → **Normalization** → $X_{normed}$ (N,C,H,W) → **Scale & Shift** $f(\gamma, \beta)$ → $X_{transformed}$ (N,C,H,W) → **ReLU**

Passports, $X_p$ (2,C,H,W) → $X_{convp}$ (2,C,H,W) → **$F()$ Average Pooling** → $\gamma, \beta$ 2x (C,)

$$X_{transformed} = \gamma * Xn_{ormed} + \beta$$

### Embedding Passport



**Valid passport** → **Passport-protected** Model with **valid passport** → "Cat"

**Fake passport** → **Passport-protected** Model with **fake passport** → "@%^?" **Poor result**

## Contributions

1. Novel passport-based verification schemes to defeat ambiguity attack
2. One passport-protected DNN model will only have one unique signature
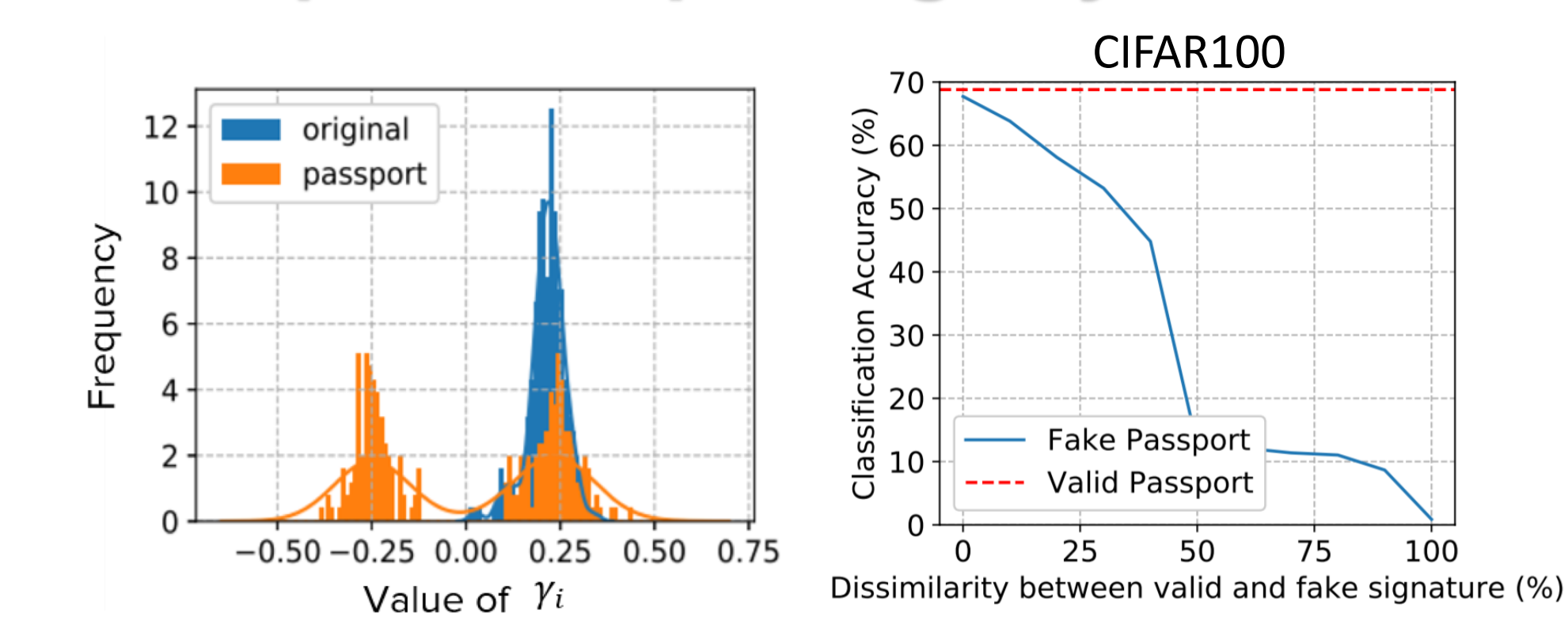3. Fake passport or modified signature will paralyze the DNN model

## Discussion

### Embedding Binary Signatures into $\gamma$ of Passporting Layer

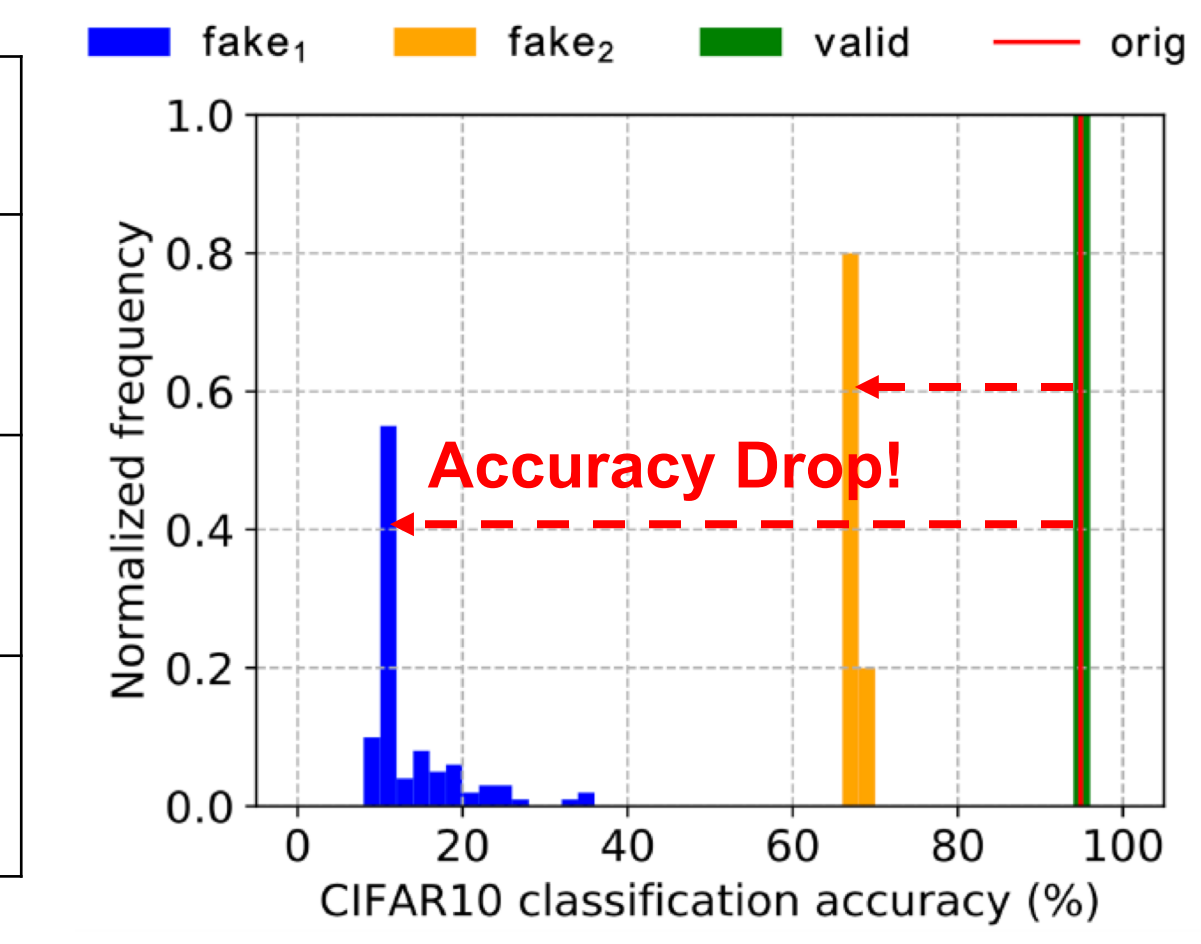$$\text{Sign Loss} = \sum_{i=1}^{C} \max(\gamma_0 - \gamma_i b_i, \ 0)$$

$\gamma_0 = 0.1$

$b: [-1 \ 1 \ ...]$

**64 channels can embed 8 bytes signature**



CIFAR100

### Experimental Results

| Ambiguity attack | Inference Phase | Verification Phase |
|---|---|---|
| Fake$_1$ (random passport) | **Random guessing** | **Useless Infringement** |
| Fake$_2$ (reverse-engineered passport) | **Performance deteriorated** (at best 70% on CIFAR10) | **Useless Infringement** |
| Fake$_3$ (copied passport) | **Performance Detained Signature Detected** | **Ownership Verified** |



### Ownership Verification Schemes

| | Scheme 1 | Scheme 2 | Scheme 3 |
|---|---|---|---|
| **Need to distribute passport** | **Yes** | **No** | **No** |
| **Inference time** | Up to 10%** more time | **No extra time** | **No extra time** |
| Training time | Up to 30%** more time | Up to 150%** more time | Up to 150%** more time |
| **Black** or **White** box Verification | **White** | **White** | **Black & White** |

**Time increases **are linearly depending on complexity of the network architecture**