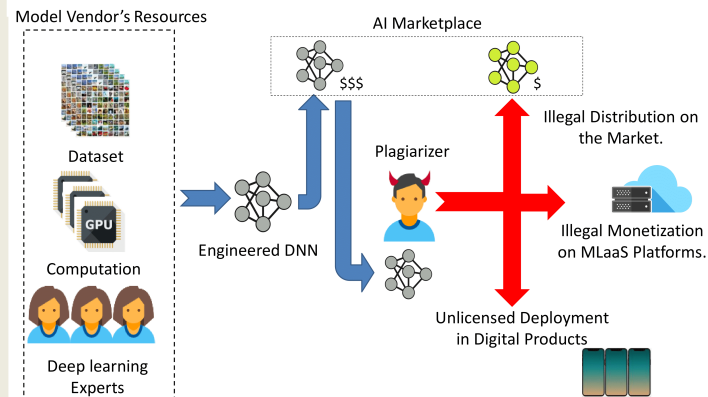


Introduction

- Engineering a **Deep Neural Network (DNN)** is a costly procedure.
- DNNs are valuable **Intellectual Property (IP)** of model vendors.
- Reliable commercialization of DNNs is threatened by IP infringement activities.



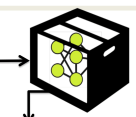
GradSigns Watermarking Framework: Embedding the watermark on the gradient of cross-entropy cost function with respect to model's input.

Watermark Embedding

- Generating embedding prerequisites \mathbf{b} , \mathbf{K} , \mathbf{C} , T .
 - Generating an N -bit vector \mathbf{b} to be used as the watermark.
 - Selecting a set \mathbf{C} of input neurons to carry the watermark.
 - Generate an **Embedding key \mathbf{K}** .
 - Select a random target class T from the dataset.
- Training the model to optimize both the original training cost function, i.e. cross-entropy function, and **GradSigns' embedding regularizer term** which penalizes the divergence from the desired watermark value.

$$J_{\text{training}} = J_{\text{cross-entropy}} + \lambda J_{\text{GradSigns-Embedding}}$$

Watermark Verification

- The vendor queries the model with samples from class T . 
- The model reports the prediction scores.
- The vendor computes the gradient of cross-entropy function w.r.t carrier nodes using a zeroth-order estimation method.

$$\bar{\mathbf{G}} \in \mathbb{R}^{|\mathbf{C}| \times 1}$$
- The vendor transforms the computed gradient vector $\bar{\mathbf{G}}$ to the watermark space using the **embedding key \mathbf{K}** .

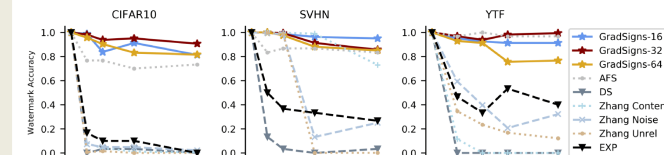
$$\mathbf{K} \times \bar{\mathbf{G}} = \boldsymbol{\chi} \quad \boldsymbol{\chi} \in \mathbb{R}^{N \times 1}$$
- The vendor extracts the embedded watermark $\mathbf{y} \in \{0, 1\}^{N \times 1}$ by binarizing $\boldsymbol{\chi}$ using the sign function.
- If the extracted watermark \mathbf{y} matches the signature \mathbf{b} , the model belongs to the vendor.

Experiments and Results

- Watermarks of varying sizes (**16, 32 and 64 bits**) were embedded into DNNs targeting classification tasks of CIFAR10, SVHN, and YouTubeFaces.
- GradSigns has **minimal overhead** (on average, less than 1%) on performance of the host model.

Dataset	Baseline Accuracy	GS-16	GS-32	GS-64
CIFAR10	91.2%	90.1%	90.4%	90.5%
SVHN	96.1%	95.5%	95.7%	95.3%
YTF	99.6%	99.6%	99.6%	98.6%

- Unlike existing black-box watermarking methods, **GradSigns is robust to watermark removal attacks** such as **model pruning and fine-tuning**.



Number of samples available for each classification label in adversary's fine-tuning dataset.

- GradSigns is also **robust** against a large array of **known and adaptive watermark removal attacks**, listed below.

Counter Watermark Attacks	Description
Query Invalidation	Adversary checks and sanitizes model queries using auto encoders.
Input Noise Injection	Adversary corrupts model queries using a random Gaussian noise.
Model Quantization	Adversary compresses the model.
Adversarial Fine-tuning	Adversary fine-tunes the model using adversarial examples.
Score Rounding	Adversary reports the rounded prediction scores.
Score Perturbation	Adversary corrupts reported prediction scores using a random noise.

Conclusion

- GradSigns enables model vendors to protect their IP by embedding robust multi-bit signatures.
- GradSigns is applicable to black-box verification scenario.
- GradSigns has negligible impact on accuracy of the model.



Core Concepts

- DNNs should be protected against IP infringements.
- Digital watermarking can be a viable solution for Digital Right Managements (DRM) of DNNs.

Must Haves

- An ideal watermark should be a meaningful **multi-bit signature**.
- An ideal watermark should be **robust to watermark removal attempts**.
- An ideal watermark must be **verifiable in a black-box setting**.

Problem

- The properties above are hard to achieve together.

Our Solution

- GradSigns, A Novel Watermarking Framework for DNNs.