

# TRANSFER LEARNING-BASED MODEL PROTECTION WITH SECRET KEY

*MaungMaung AprilPyone and Hitoshi Kiya*

Tokyo Metropolitan University, Tokyo, Japan

## ABSTRACT

We propose a novel method for protecting trained models with a secret key so that unauthorized users without the correct key cannot get the correct inference. By taking advantage of transfer learning, the proposed method enables us to train a large protected model like a model trained with ImageNet by using a small subset of a training dataset. It utilizes a learnable encryption step with a secret key to generate learnable transformed images. Models with pre-trained weights are fine-tuned by using such transformed images. In experiments with the ImageNet dataset, it is shown that the performance of a protected model was close to that of a non-protected model when the correct key was given, while the accuracy tremendously dropped when an incorrect key was used. The protected model was also demonstrated to be robust against key estimation attacks.

**Index Terms**— Model Protection, Learnable Image Encryption, Model Watermarking

## 1. INTRODUCTION

Training successful deep neural networks (DNNs) is very expensive because it requires a huge amount of data and fast computing resources (e.g., GPU-accelerated computing). To train a convolutional neural network (CNN), for example, the ImageNet [1] dataset contains about 1.2 million images, and training on such a dataset takes days and weeks even on GPU-accelerated machines. In fact, collecting images and labeling them will also consume a massive amount of resources. Moreover, algorithms used in training a CNN model may be patented or have restricted licenses. Considering the expenses necessary for the expertise, money, and time taken to train a CNN model, a model should be regarded as a kind of intellectual property. While distributing a trained model, an illegal party may also obtain a model and use it for its own service.

To protect the copyrights of trained models, researchers have adopted digital watermarking technology to embed watermarks into the models [2–9]. These works focus on identifying the ownership of a model in question. However, a stolen model can be directly used by an attacker without arousing suspicion. Moreover, a stolen model can be exploited through model inversion attacks [10] and adversarial attacks [11]. Therefore, a trained model should be protected

against unauthorized access beyond ownership verification.

Recently, Fan et al. [4] proposed a passport-protected model-protection method. However, in their work, the network has to be modified with passport layers that use passports, and there are significant overhead costs in both training and inference time. Another work [12] introduced a model protection method by taking inspiration from adversarial defenses [13, 14] that exploit the uniqueness of a key. They utilized a block-wise transformation with a key for model protection [12]. However, it was tested only on CIFAR-10 [15], and the protected model was trained from scratch. Considering a large dataset like ImageNet [1], it is not feasible to train a protected model from scratch as in [4, 12]. Although, both the passport-protected [4] and key-protected [12] methods train protected models from scratch, they do not consider transfer learning.

Transfer learning has been proved to be effective in various visual recognition tasks [16]. Transfer learning can be used in either of two scenarios: a pre-trained model can be transferred to a new model with the same number of classes or to a new one with a different number of classes (usually a lower number of classes). In this paper, we focus on the first scenario (i.e., transfer to the same number of classes) to confirm the effectiveness of the proposed method.

We propose a model protection method with a secret key that takes advantage of transfer learning for the first time. The proposed method also allows us to use a small subset of a training dataset to replace an unprotected model with a protected one. In addition, it does not need to modify a network, and therefore, there is no overhead for both training and inference time. In an experiment on ImageNet, the performance of a model protected by the proposed method is demonstrated not only to be close to that of a non-protected one when the key is correct but also to significantly drop when using an incorrect key.

## 2. RELATED WORK

### 2.1. Model Watermarking

There are mainly two categories of DNN model watermarking: white-box and black-box. A white-box approach requires access to model weights for embedding and extracting a watermark as in [2, 4, 6, 8]. In contrast, black-box ap-

proaches [3–5, 7, 9] do not need to access model weights and observe the input and output of a model in doubt to verify the ownership of the model.

These existing model-watermarking schemes focus on ownership verification only. Thus, a stolen model can be directly used and exploited without arousing suspicion because the performance of a protected model (i.e., fidelity) is independent of the embedded watermark.

In addition, Fan et al. [4] pointed out that conventional ownership verification schemes are vulnerable against ambiguity attacks [17] where two watermarks can be extracted from the same protected model, causing confusion regarding ownership. Therefore, Fan et al. [4] introduced passports and passport layers. The passports in [4] are a set of extracted features of a secret image/images or equivalent random patterns from a pre-trained model, and the passport layers are additional layers in the network. Therefore, there are significant overhead costs in both the training and inference phases, in addition to user-unfriendly management of lengthy passports in [4]. Moreover, the protection method with passports [4] was evaluated only on CIFAR datasets [15] and does not consider transfer learning.

## 2.2. Learnable Image Encryption

Learnable image encryption perceptually encrypts images while maintaining a network’s ability to learn the encrypted ones for classification tasks. Most early methods of learnable image encryption were originally proposed to visually protect images for privacy-preserving DNNs [18–23].

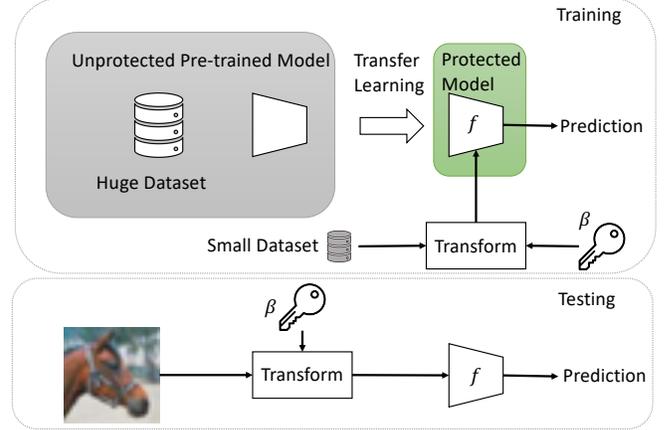
Recently, adversarial defenses in [13, 14] also utilized learnable image encryption methods. Here, instead of protecting visual information, these works focus on controlling a model’s decision boundary with a secret key so that adversarial attacks are not effective on such models trained by learnable transformed images.

Another use case of learnable image encryption is the model protection proposed in [12]. However, the method in [12] requires training a protected model from scratch and never considers transfer learning. In this paper, we adopt a block-wise image encryption method as in [13, 14] to transform images prior to training and testing.

## 3. PROPOSED METHOD

### 3.1. Overview

An overview of image classification with the proposed method is depicted in Fig. 1. In the proposed model protection, a model  $f$  is not trained from random weights. Instead,  $f$  is trained by taking advantage of transfer learning. In other words,  $f$  is trained by fine-tuning pre-trained weights by using secret key  $\beta$ . The resulting  $f$  is protected by key  $\beta$ . For testing, test images are also transformed with the same key



**Fig. 1.** Image classification with proposed model-protection method

$\beta$  before testing. Therefore, the authorization of model  $f$  is verified upon secret key  $\beta$  during model inference.

### 3.2. Block-wise Transformation

We use negative/positive transformation with a secret key to transform input images before training or testing a protected model as well as in [14]. The following are steps for transforming input images, where  $c$ ,  $w$ , and  $h$  denote the number of channels, width, and height of an image tensor  $x \in [0, 1]^{c \times w \times h}$ .

1. Divide  $x$  into blocks with a size of  $M$  such that  $\{B_{(1,1)}, \dots, B_{(\frac{w}{M}, \frac{h}{M})}\}$ .
2. Transform each block tensor  $B_{(i,j)}$  into a vector  $b_{(i,j)} = [b_{(i,j)}(1), \dots, b_{(i,j)}(c \times M \times M)]$ .
3. Generate a key  $\beta$ , which is a binary vector, i.e.,

$$\beta = [\beta_1, \dots, \beta_k, \dots, \beta_{(c \times M \times M)}], \beta_k \in \{0, 1\}, \quad (1)$$

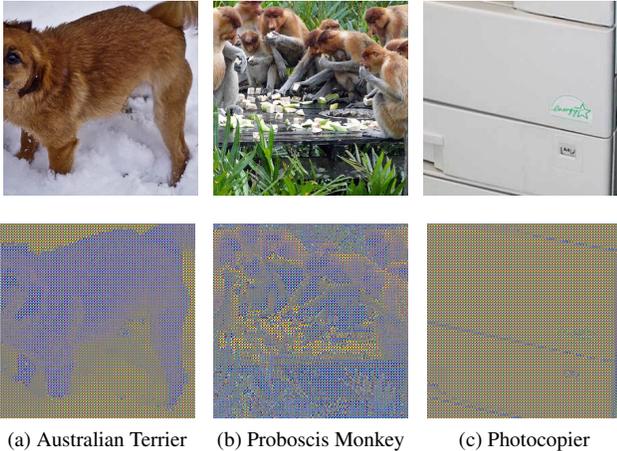
where the value of the occurrence probability  $P(\beta_k)$  is 0.5.

4. Multiply each pixel value in  $b_{(i,j)}$  by 255 to be at 255 scale with 8 bits.
5. Apply negative/positive transformation to every vector  $b_{(i,j)}$  with  $\beta$  as

$$b'_{(i,j)}(k) = \begin{cases} b_{(i,j)}(k) & (\beta_k = 0) \\ b_{(i,j)}(k) \oplus (2^L - 1) & (\beta_k = 1), \end{cases} \quad (2)$$

where  $\oplus$  is an exclusive or (XOR) operation,  $L$  is the number of bits used in  $b_{(i,j)}(k)$ , and  $L = 8$  is used in this paper.

6. Divide each pixel value in  $b'_{(i,j)}$  by 255 to be at  $[0, 1]$  scale.



**Fig. 2.** Example of block-wise transformed images ( $M = 4$ ) with key  $\beta$  (second row). Images in first row are original.

7. Integrate the transformed vectors to form an image tensor  $x' \in [0, 1]^{c \times w \times h}$ .

An example of images (three different classes from the ImageNet test set) transformed by negative/positive transformation with  $M = 4$  is shown in Fig. 2.

### 3.3. Transfer Learning

In practice, CNNs are not trained from the beginning with random weights because creating a large dataset like ImageNet is difficult and expensive. Therefore, CNNs are usually pre-trained with a larger dataset [24]. There are two major transfer-learning scenarios:

- **Fixed CNN:** A pre-trained CNN model is used as a fixed feature extractor, and the last fully connected layer is replaced with a targeted number of classes. In other words, convolutional layers are frozen, and only the last fully connected layer is trained with random initialization from scratch.
- **Fine-tuned CNN:** In this scenario, the CNN is fine-tuned from a pre-trained model. Here, it is possible that some convolutional layers can be fixed or the whole CNN is fine-tuned.

In this paper, we fine-tuned a whole CNN with a small dataset comprised of learnable transformed images with a secret key in order to protect a model.

### 3.4. Key Estimation Attack

We consider a threat model where a model is stolen and transformation details are known except for the secret key. In this

scenario, an attacker may try all possible keys (brute-force attack). The key space  $\mathcal{K}$  of negative/positive transformation is given by

$$\mathcal{K}(c \times M \times M) = 2^{(c \times M \times M)}. \quad (3)$$

Therefore, the key space will vary with respect to block size  $M$ .

However, checking all possible keys may not be feasible, and the attacker may estimate a key heuristically by observing the accuracy of his/her test dataset. Elements in  $\beta$  can be rearranged in accordance with the improvement in accuracy. We simulate this attacking scenario by swapping values in each index pair of  $\beta$  if the accuracy improves.

Key estimation attacks do not guarantee that the attacker will find the correct key because the attacker does not know the actual performance of the correct key. The robustness of the proposed method against key estimation attacks will be demonstrated in the following section on experiments.

## 4. EXPERIMENTS

### 4.1. Setup

We utilized the ImageNet dataset [1], which comprises 1.28 million color images for training and 50,000 color images for validation. We progressively resized images during training, starting with larger batches of smaller images to smaller batches of larger images. We adapted three phases of training from the DAWNbench top submissions as mentioned in [25]. Phases 1 and 2 resized images to 160 and 352 pixels, respectively, and phase 3 used the entire image size from the training set. The augmentation methods used in the experiment were random resizing, cropping (sizes of 128, 224, and 288, respectively, for each phase), and random horizontal flip. Both training and testing images were transformed with negative/positive transformation with a block size  $M = 4$ .

We deployed deep residual networks [26] with 50 layers (ResNet50) with pre-trained weights and fine-tuned for 15 epochs with cyclic learning rates [27] and mixed precision training [28]. We adapted the training settings from [25] with the removal of weight decay regularization from batch normalization layers.

### 4.2. Classification Performance

Table 1 summarizes the classification results for protected models and a baseline (unprotected one). We fine-tuned models by using subsets of the training dataset with 10%, 20%, 30%, and 100% of the training set, respectively. Images in the sub-datasets were transformed by negative/positive transformation with a secret key  $\beta$  and  $M = 4$  as mentioned in Section 3.2. We tested the proposed method under three conditions: with a correct key  $\beta$ , with an incorrect key  $\beta'$ , and with plain images.

**Table 1.** Accuracy (%) of protected models and baseline model for ImageNet

Model	Correct ( $\beta$ )	Incorrect ( $\beta'$ )	Plain
10 % dataset	64.13	0.16	0.30
20 % dataset	67.45	0.25	1.04
30 % dataset	68.87	0.24	0.73
100 % dataset	<b>72.63</b>	0.69	0.36
Baseline	73.70 (Not protected)		

**Table 2.** Accuracy (%) of protected models against key estimation attack for ImageNet

Model	Estimated ( $\beta'$ )
10 % dataset	0.17
20 % dataset	1.61
30 % dataset	9.42
100 % dataset	25.43

When a correct key  $\beta$  was given, the model trained with a 10 % dataset had about 9 % less accuracy than that of the baseline. However, when the whole dataset was used, the accuracy was almost the same as the baseline accuracy (i.e., 72.63 %). Even when the whole dataset was used, transfer learning significantly reduced the training time because the model was trained only for 15 epochs. When using an incorrect key  $\beta'$  or plain images, the accuracy was extremely low, suggesting the strength of the proposed method against unauthorized access.

#### 4.3. Robustness Against Key Estimation Attack

We also evaluated the proposed method in terms of robustness against key estimation attacks. Table 2 captures the results for all protected models under the use of a subset of the training dataset with various sizes. The model trained by the smallest dataset (i.e., 10 %) had the lowest accuracy, and that trained by the whole dataset had a 25.43 % accuracy. All in all, the key estimation attack did not guarantee that a good enough key would be found, and the performance accuracy was not usable, suggesting the robustness of the proposed method.

#### 4.4. Functional Comparison with State-of-the-art Methods

To the best of our knowledge, there are only two methods [4, 12] where the protection method is directly dependent on model performance (i.e., key/passports protected models). However, both of them were not tested on ImageNet, and it is not feasible to train an ImageNet model from scratch. The other model watermarking methods such as [2–9] focus on ownership verification only when a stolen model is in question. Therefore, the embedded watermark is independent of

**Table 3.** Functional comparison of proposed method and state-of-the-art methods

Model	Performance Dependency	Network Modification	Performance Degradation	Overhead
Scheme $\mathcal{V}_1$ [4]	Passports	Yes	Low	Significant
Pixel Shuffling [12]	Key	No	Low	No
Proposed Method	Key	No	Low	No

model accuracy.

Since the state-of-the-art model-protection methods [4, 12] cannot be directly compared with the proposed method as described above, we performed a functional comparison with a key-protected method (Pixel Shuffling) [12] and a passport-protected method (Scheme  $\mathcal{V}_1$ ) [4]. Both methods control the accuracy of performance by using a key or passports and have low performance degradation. However, scheme  $\mathcal{V}_1$  [4] has to modify a network with additional passport layers; therefore, it introduces overheads in training (15 %–30 %) and inference (10 %) processes as mentioned in [4]. In contrast, the proposed method does not introduce any overhead during training and testing, in addition to being applicable to transfer learning.

## 5. CONCLUSION

In this paper, we proposed a model protection method in which a model is fine-tuned with a subset of a training dataset. Images in the sub-datasets are transformed by a block-wise transformation with a secret key prior to training and testing a model. The performance accuracy of a protected model trained by using 10 % of a training dataset was about 9 % less than that of a baseline model. When using the whole dataset, the accuracy was close to the baseline accuracy, and transfer learning significantly reduced the training time because the model was trained only for 15 epochs. The proposed model-protection method was also confirmed to be robust against key estimation attacks and not usable when using an incorrect key or plain images. Moreover, it does not introduce any overhead in both training and inference time.

## 6. REFERENCES

- [1] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [2] 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*, 2017, pp. 269–277.
- [3] Erwan Le Merrer, Patrick Pérez, and Gilles Trédan, “Adversarial frontier stitching for remote neural network watermarking,”

*Neural Computing and Applications*, vol. 32, no. 13, pp. 9233–9244, 2020.

- [4] Lixin Fan, KamWoh Ng, and Chee Seng Chan, “Rethinking deep neural network ownership verification: Embedding passports to defeat ambiguity attacks,” in *Advances in Neural Information Processing Systems*, 2019, pp. 4716–4725.
- [5] Shigeyuki Sakazawa, Emi Myodo, Kazuyuki Tasaka, and Hiromasa Yanagihara, “Visual decoding of hidden watermark in trained deep neural network,” in *2nd IEEE Conference on Multimedia Information Processing and Retrieval*, 2019, pp. 371–374.
- [6] Bitar Darvish Rouhani, Huili Chen, and Farinaz Koushanfar, “Deepsigns: A generic watermarking framework for IP protection of deep learning models,” *arXiv:1804.00750*, 2018.
- [7] 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*, 2018, pp. 159–172.
- [8] Huili Chen, Bitar Darvish Rouhani, and Farinaz Koushanfar, “Deepmarks: A digital fingerprinting framework for deep neural networks,” *arXiv:1804.03648*, 2018.
- [9] Yossi Adi, Carsten Baum, Moustapha Cissé, Benny Pinkas, and Joseph Keshet, “Turning your weakness into a strength: Watermarking deep neural networks by backdooring,” in *27th USENIX Security Symposium*, 2018, pp. 1615–1631.
- [10] Matt Fredrikson, Somesh Jha, and Thomas Ristenpart, “Model inversion attacks that exploit confidence information and basic countermeasures,” in *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, 2015, pp. 1322–1333.
- [11] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian J. Goodfellow, and Rob Fergus, “Intriguing properties of neural networks,” in *International Conference on Learning Representations*, 2014.
- [12] MaungMaung AprilPyone and Hitoshi Kiya, “Training dnn model with secret key for model protection,” in *2020 IEEE 9th Global Conference on Consumer Electronics (GCCE)*, 2020, pp. 818–821.
- [13] MaungMaung AprilPyone and Hitoshi Kiya, “Encryption inspired adversarial defense for visual classification,” in *2020 IEEE International Conference on Image Processing (ICIP)*, 2020, pp. 1681–1685.
- [14] MaungMaung AprilPyone and Hitoshi Kiya, “Block-wise image transformation with secret key for adversarially robust defense,” *IEEE Transactions on Information Forensics and Security*, pp. 1–1, 2021.
- [15] Alex Krizhevsky and Geoffrey Hinton, “Learning multiple layers of features from tiny images,” Tech. Rep., University of Toronto, 2009.
- [16] Simon Kornblith, Jonathon Shlens, and Quoc V Le, “Do better imagenet models transfer better?,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 2661–2671.
- [17] Scott Craver, Nasir D. Memon, Boon-Lock Yeo, and Minerva M. Yeung, “Resolving rightful ownerships with invisible watermarking techniques: limitations, attacks, and implications,” *IEEE J. Sel. Areas Commun.*, vol. 16, no. 4, pp. 573–586, 1998.
- [18] Masayuki Tanaka, “Learnable image encryption,” in *2018 IEEE International Conference on Consumer Electronics-Taiwan (ICCE-TW)*, 2018, pp. 1–2.
- [19] Warit Sirichotedumrong, Yuma Kinoshita, and Hitoshi Kiya, “Pixel-based image encryption without key management for privacy-preserving deep neural networks,” *IEEE Access*, vol. 7, pp. 177844–177855, 2019.
- [20] Warit Sirichotedumrong, Takahiro Maekawa, Yuma Kinoshita, and Hitoshi Kiya, “Privacy-preserving deep neural networks with pixel-based image encryption considering data augmentation in the encrypted domain,” in *2019 IEEE International Conference on Image Processing (ICIP)*, 2019, pp. 674–678.
- [21] Warit Sirichotedumrong and Hitoshi Kiya, “A gan-based image transformation scheme for privacy-preserving deep neural networks,” in *2020 28th European Signal Processing Conference (EUSIPCO)*, 2020, pp. 745–749.
- [22] Hiroki Ito, Yuma Kinoshita, and Hitoshi Kiya, “Image transformation network for privacy-preserving deep neural networks and its security evaluation,” *arXiv:2008.03143*, 2020.
- [23] Hiroki Ito, Yuma Kinoshita, and Hitoshi Kiya, “A framework for transformation network training in coordination with semi-trusted cloud provider for privacy-preserving deep neural networks,” in *2020 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, Dec 2020, pp. 1420–1424.
- [24] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *International Conference on Learning Representations*, 2015.
- [25] Eric Wong, Leslie Rice, and J. Zico Kolter, “Fast is better than free: Revisiting adversarial training,” in *International Conference on Learning Representations*, 2020.
- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [27] Leslie N. Smith and Nicholay Topin, “Super-convergence: Very fast training of residual networks using large learning rates,” *arXiv:1708.07120*, 2017.
- [28] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory F. Damos, Erich Elsen, David García, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu, “Mixed precision training,” *arXiv:1710.03740*, 2017.