

Access Control of Object Detection Models Using Encrypted Feature Maps

1st Teru Nagamori
Tokyo Metropolitan University
Tokyo, Japan
nagamori-teru@ed.tmu.ac.jp

2nd Hiroki Ito
Tokyo Metropolitan University
Tokyo, Japan
ito-hiroki2@ed.tmu.ac.jp

3rd MaungMaung AprilPyone
Tokyo Metropolitan University
Tokyo, Japan
april-pyone-maung-maung@ed.tmu.ac.jp

4th Hitoshi Kiya
Tokyo Metropolitan University
Tokyo, Japan
kiya@tmu.ac.jp

Abstract—In this paper, we propose an access control method for object detection models. The use of encrypted images or encrypted feature maps has been demonstrated to be effective in access control of models from unauthorized access. However, the effectiveness of the approach has been confirmed in only image classification models and semantic segmentation models, but not in object detection models. In this paper, the use of encrypted feature maps is shown to be effective in access control of object detection models for the first time.

Index Terms—Object Detection, Access Control, Feature Map

I. INTRODUCTION

Deep neural networks (DNNs) and convolutional neural networks (CNNs) have been used widely in various applications such as image classification, semantic segmentation, and object detection [1]–[3]. Training high-performance models is not an easy task, because it requires a large amount of data, powerful computational resources (GPUs), and efficient algorithms. Considering the expertise, cost, and time required for training models, they are considered as a kind of intellectual property that should be protected.

There are two approaches to intellectual property protection of models: ownership verification and access control. The difference between these two approaches is that the former aims to identify the ownership of the models, but the latter aims to protect models from unauthorized access [4]. The ownership verification methods are inspired by watermarking, where a watermark is embedded in the models and the embedded watermark is used to verify the ownership of the models [5]–[10]. However, ownership verification does not have the ability to restrict the execution of the models. Thus, in principle, attackers can freely exploit the models for their own benefit, or use it in adversarial attacks [11]. Therefore, in this paper, we focus on access control, which aims to prevent models from unauthorized access.

A number of access control methods have been proposed as a model protection method. By encrypting images or feature maps with a secret key, a stolen model cannot be used to its full

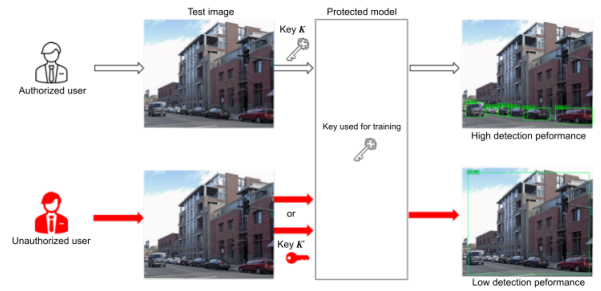


Fig. 1. Overview of access control

capacity without a correct secret key [12]–[14]. However, these methods have never been applied to object detection models. In this paper, an access control method with encrypted feature maps is applied to object detection models for the first time, and the effectiveness of the proposed method is confirmed in an experiment.

II. PROPOSED METHOD

A. Overview

An overview of the access control to protect the trained models from unauthorized access is shown in Fig. 1. The protected models are trained with the secret key K . When authorized users enter test images and the correct key K into the protected models, the results are equivalent to the models in the unprotected state (the access control is not assumed). In contrast, when unauthorized users without the key K enter only test images or test images and a wrong key K' into the protected models, lower performance results are provided.

As access control methods using a secret key, the input image encryption method [12] and the feature map encryption method [14] have been proposed. Maung's method [12] focuses on access control of image classification models, where input images are divided into blocks and encrypted with a secret key using methods such as pixel shuffling, bit flipping, and format-preserving Feistel-based encryption (FFX) [15].

These encrypted images are used as training and test images. Since this method encrypts the images block by block, it changes the spatial information and cannot be used to protect the object detection models described below.

Ito's method [14] focuses on access control of semantic segmentation models, where models are trained and tested by randomly permuting the channels of feature maps selected by a secret key. This encryption method is spatially invariant. This property was confirmed to be very important for some applications such as semantic segmentation [14]. Although this method has been validated for semantic segmentation, it has not been validated for object detection models. Therefore, in this paper, we propose an access control method for object detection models based on this method.

B. Encryption Method

There are multiple feature maps in CNNs as shown in Fig. 2. In the proposed method, selected feature maps are transformed by using a secret key in accordance with the procedure of learnable image encryption [12], [16]. Below is the procedure of the encryption, where x is a selected feature map with a dimension of $(c \times h \times w)$, c is the number of channels, h is the height, and w is the width of the feature map.

- 1) Generate a random vector with a size of c using a secret key as in (1).

$$[\alpha_1, \dots, \alpha_i, \alpha_{i'}, \dots, \alpha_c], \alpha_i \in \{1, \dots, c\} \quad (1)$$

where $\alpha_i \neq \alpha_{i'}$, if $i \neq i'$.

- 2) Replace each element $x(i, j, k)$ of x , $i \in \{1, \dots, c\}$, $j \in \{1, \dots, h\}$, $k \in \{1, \dots, w\}$ with $x(\alpha_i, j, k)$ so that x is transformed into a feature map x' . Note that elements of x' , $x'(i, j, k)$ is equal to $x(\alpha_i, j, k)$.

This encryption is a spatial-invariant operation, so the spatial information of feature maps can be maintained (see Fig. 3). This property is very important in object detection tasks, which predict position and classes of objects.

C. Model Training and Testing

In the proposed method, the previously mentioned transformation method is applied to selected feature maps in an object detection model at each iteration for a training model. SSD300 [17] based on VGG16 [18], which was pretrained on the ILSVRC CLS-LOC dataset [19] is used as an object detection model in this paper, where SSD300 has 11 feature maps as illustrated in Fig. 2.

In testing the trained model, authorized users have the same key that is used for the training. When Authorized users apply query images to the model, they transform the same feature maps that are selected for the training with the key. If unauthorized users without the correct key steal the protected model, we assume that they transform the feature maps with an incorrect key or use the model without the transform.

D. Requirements of Protected Models

Protected models should meet the following requirements.

- It provides almost the same performance as that of models trained with plain images to authorized users with the secret key.
- It provides a degraded performance to unauthorized users without the correct key.

III. EXPERIMENTS AND RESULTS

A. Setup

We used the PASCAL visual object classes (VOC) challenge 2007 [20], and 2012 [21] trainval datasets for training, and the PASCAL VOC 2007 test dataset for testing. For data augmentation, the random sample crop, horizontal flip, and some photometric distortions described in [17] were used for training models. In addition, due to the restrictions of SSD300 shown in Fig. 2, input images were resized to 300×300 pixels.

Models were trained by using a stochastic gradient descent (SGD) optimizer with an initial learning rate of 10^{-3} , a momentum value of 0.9, a weight decay value of 0.0005, and a batch size of 32. Models were also trained with a learning rate of 10^{-3} for 60k iterations, then continue training for 20k iterations with 10^{-4} and 40k iterations with 10^{-5} . The overall objective loss function is a weighted sum of the localization loss and the confidence loss. In this paper, the confidence loss was the cross-entropy loss over multiple classes confidences, and the localization loss was the Smooth L1 loss between the predicted position and the ground truth position.

B. Detection Performance

Mean average precision (mAP) [17] with a range [0,1] was used as a metric for evaluating detection performance, where when a mAP value is closer to 1, it indicates a higher accuracy. In the experiment, a selected feature map was transformed with a key K in accordance with the procedure in sec. II. In Table I, Correct (K) indicates that the selected feature map was transformed with the same key K as the training. Model-1 means that feature map 1 was selected for encryption, and Baseline indicates that training and testing were performed without any encryption. Fig. 4 shows examples of experimental result where Model-4 was used.

From Table I and Fig. 4, the proposed method provides the same prediction results as the Baseline when the feature map is transformed using the correct key for the test.

C. Robustness against Unauthorized Access

Two types of unauthorized access were considered in the experiment. Plain in Table I represents that an unauthorized user without the key applied query images to protected models, without transforming the selected feature map. Incorrect (K') in Table I is that an unauthorized user without the key applied query images to protected models, after transforming the selected feature map with a randomly generated key K' . The result of Incorrect (K') are the average value of 100 tests with random keys.

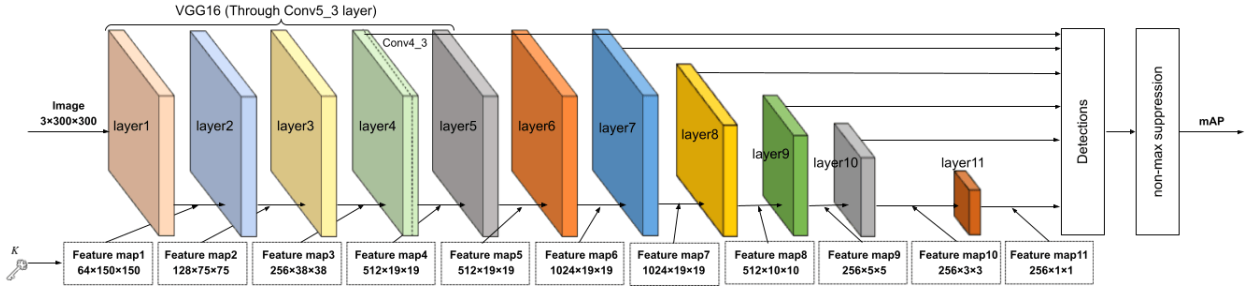


Fig. 2. Architecture of object detection model (SSD300)

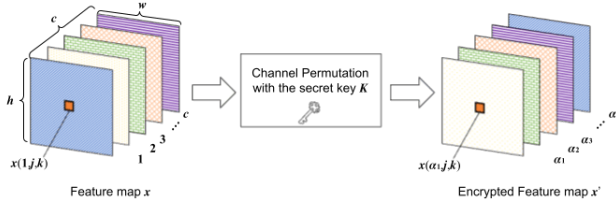


Fig. 3. Feature map encryption [14]

TABLE I
DETECTION ACCURACY (MAP) OF PROPOSED MODELS

Selected feature map	Correct (K)	Plain	Incorrect (K')
Model-1	0.7244	0.1363	0.0421
Model-2	0.7611	0.0091	0.0180
Model-3	0.7475	0.0091	0.0078
Model-4	0.7611	0.0023	0.0043
Model-5	0.7587	0.1672	0.1624
Model-6	0.7617	0.1732	0.1672
Model-7	0.7695	0.1768	0.1750
Model-8	0.7677	0.3529	0.3415
Model-9	0.7705	0.5767	0.5678
Model-10	0.7705	0.7177	0.7027
Model-11	0.7512	0.7314	0.7252
Baseline		0.7690	

From the table, Model-1–7 provided a low detection accuracy for both Plain and Incorrect (K'). On the other hand, when transforming the feature map of a deep layer, the resistance to unauthorized access is lost. We consider that the reason for this lies in the structure of SSD300. In order to detect objects of various scales in SSD, detection is performed using features from multiple layers (see Fig 2). Therefore, for example, in Model-9, layers 4, 7, and 8 can use the same features as Baseline. In other words, we consider that this is because the number of the same features as Baseline increases in the deeper layers.

From Fig. 4, the detection performance degraded significantly when the model was used illegally. Accordingly, the proposed models were robust enough against the unauthorized access.

D. Comparison with encryption of input images

The proposed method was compared with a method to protect models with encrypted input images, which was proposed for image classification models [12]. In the method, there are

TABLE II
DETECTION ACCURACY (MAP) OF MODELS
WITH ENCRYPTED INPUT IMAGES

method	block size	Correct (K)	Plain	Incorrect (K')
pixel shuffling (SHF)	1	0.7710	0.7598	0.7603
	4	0.7154	0.5745	0.3883
	12	0.4891	0.1976	0.0910
	20	0.0083	0.0086	0.0065
	60	0.1284	0.0480	0.0416
Proposed (Model-4)		0.7611	0.0023	0.0043
Baseline			0.7690	

three block-wise methods: pixel shuffling, bit flipping, and Format-preserving Feistel-based encryption (FFX) [15], for encrypting input image.

In this paper, pixel shuffling (SHF) with a block size of 1, 4, 12, 20, or 60 were applied to input images, and the encrypted images were used for training and testing.

The experimental conditions are the same as in A of sec. III. From Table II, the detection accuracy was significantly lower than the proposed method under almost all block sizes. When the block size was small, the detection accuracy was high, but the resistance to unauthorized access was also degraded, so the models were not protected [12]. In contrast, when the block size was large, the resistance to unauthorized access was stronger, but the detection accuracy was greatly degraded. Therefore, the conventional method with encrypted input images is not effective in object detection models.

IV. CONCLUSION

We proposed an access control method that uses encrypted feature maps transformation for object detection models for the first time. In the experiment, the proposed access control method was demonstrated not only to provide a high detection accuracy but also to robust enough against two types of unauthorized access.

ACKNOWLEDGEMENT

This study was partially supported by JSPS KAKENHI (Grant Number JP21H01327).

REFERENCES

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.

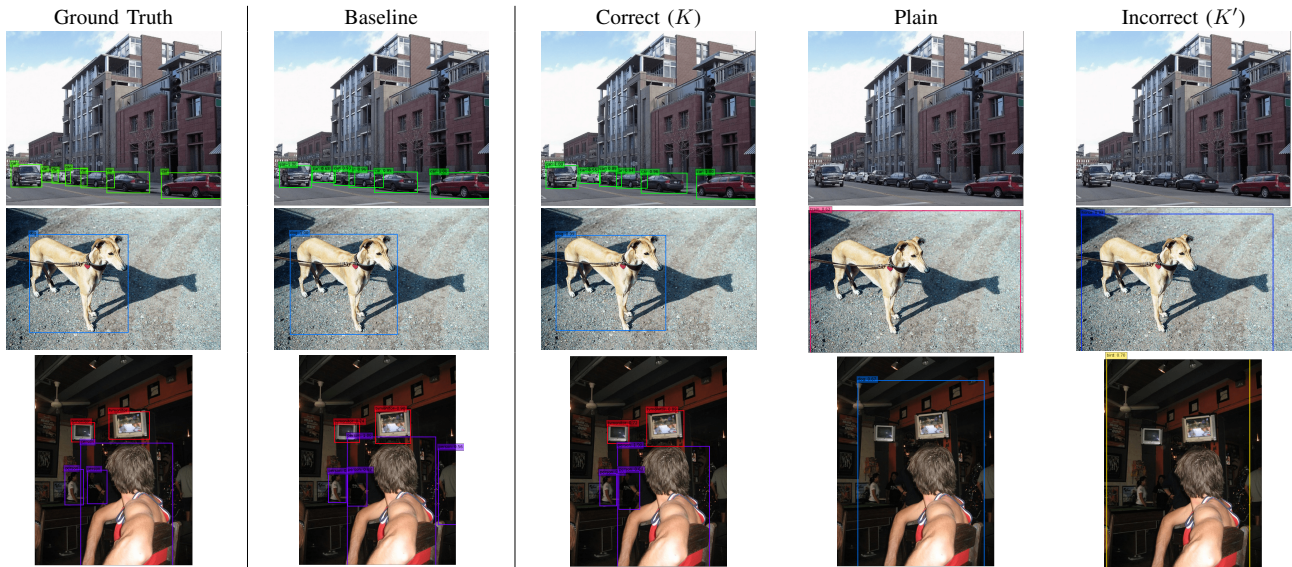


Fig. 4. Examples of experimental result (Model-4)

- [2] X. Liu, Z. Deng, and Y. Yang, "Recent progress in semantic image segmentation," *Artificial Intelligence Review*, vol. 52, no. 2, pp. 1089–1106, 2019.
- [3] R. L. Galvez, A. A. Bandala, E. P. Dadios, R. R. P. Vicerra, and J. M. Z. Maningo, "Object detection using convolutional neural networks," in *TENCON 2018-2018 IEEE Region 10 Conference*. IEEE, 2018, pp. 2023–2027.
- [4] H. Kiya, A. MaungMaung, Y. Kinoshita, I. Shoko, and S. Shiota, "An overview of compressible and learnable image transformation with secret key and its applications," *arXiv preprint arXiv:2201.11006*, 2022. [Online]. Available: <https://arxiv.org/abs/2201.11006>
- [5] Y. Uchida, Y. Nagai, S. Sakazawa, and S. Satoh, "Embedding watermarks into deep neural networks," in *Proceedings of the 2017 ACM on International Conference on Multimedia Retrieval*, 2017, pp. 269–277.
- [6] J. Zhang, Z. Gu, J. Jang, H. Wu, M. P. Stoecklin, H. Huang, and I. 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.
- [7] B. Darvish Rouhani, H. Chen, and F. Koushanfar, "Deepsigns: An end-to-end watermarking framework for ownership protection of deep neural networks," in *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2019, pp. 485–497.
- [8] L. Fan, K. W. Ng, and C. S. Chan, "Rethinking deep neural network ownership verification: Embedding passports to defeat ambiguity attacks," in *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds., vol. 32. Curran Associates, Inc., 2019.
- [9] M. Xue, J. Wang, and W. Liu, "Dnn intellectual property protection: Taxonomy, attacks and evaluations," in *Proceedings of the 2021 on Great Lakes Symposium on VLSI*, 2021, pp. 455–460.
- [10] M. AprilPyone and H. Kiya, "Piracy-resistant dnn watermarking by block-wise image transformation with secret key," in *Proceedings of the 2021 ACM Workshop on Information Hiding and Multimedia Security*, 2021, pp. 159–164.
- [11] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, "Intriguing properties of neural networks," *arXiv preprint arXiv:1312.6199*, 2013. [Online]. Available: <https://arxiv.org/abs/1312.6199>
- [12] M. AprilPyone and H. Kiya, "A protection method of trained cnn model with a secret key from unauthorized access," *APSIPA Transactions on Signal and Information Processing*, vol. 10, p. e10, 2021.
- [13] —, "A protection method of trained cnn model using feature maps transformed with secret key from unauthorized access," in *Proceedings of APSIPA Annual Summit and Conference 2021*, 2021, pp. 1851–1857.
- [14] H. Ito, M. AprilPyone, and H. Kiya, "Access control using spatially invariant permutation of feature maps for semantic segmentation models," in *Proceedings of APSIPA Annual Summit and Conference 2021*, 2021, pp. 1833–1838.
- [15] M. Bellare, P. Rogaway, and T. Spies, "Addendum to "the ffx mode of operation for format-preserving encryption";" *A parameter collection for enciphering strings of arbitrary radix and length, Draft 1.0*, NIST, 2010.
- [16] M. AprilPyone and H. Kiya, "Block-wise image transformation with secret key for adversarially robust defense," *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 2709–2723, 2021.
- [17] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "Ssd: Single shot multibox detector," in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [18] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *International Conference on Learning Representations*, 2015.
- [19] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [20] M. Everingham, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes (voc) challenge," *International journal of computer vision*, vol. 88, no. 2, pp. 303–338, 2010.
- [21] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman, "The pascal visual object classes challenge: A retrospective," *International journal of computer vision*, vol. 111, no. 1, pp. 98–136, 2015.