

Black-Box Dissector: Towards Erasing-based Hard-Label Model Stealing Attack

Yixu Wang¹, Jie Li¹, Hong Liu², Yan Wang³, Yongjian Wu⁴, Feiyue Huang⁴,
and Rongrong Ji^{1,5✉}

¹ Media Analytics and Computing Lab, School of Informatics, Xiamen University

² National Institute of Informatics

³ Pinterest

⁴ Youtu Lab, Tencent Technology (Shanghai) Co.,Ltd

⁵ Institute of Artificial Intelligence, Xiamen University

yxwang79@gmail.com, lijie.32@outlook.com, hliu@nii.ac.jp,
yanw@pinterest.com, {littlekenwu, garyhuang}@tencent.com, rrji@xmu.edu.cn

Abstract. Previous studies have verified that the functionality of black-box models can be stolen with full probability outputs. However, under the more practical hard-label setting, we observe that existing methods suffer from catastrophic performance degradation. We argue this is due to the lack of rich information in the probability prediction and the overfitting caused by hard labels. To this end, we propose a novel hard-label model stealing method termed *black-box dissector*, which consists of two erasing-based modules. One is a CAM-driven erasing strategy that is designed to increase the information capacity hidden in hard labels from the victim model. The other is a random-erasing-based self-knowledge distillation module that utilizes soft labels from the substitute model to mitigate overfitting. Extensive experiments on four widely-used datasets consistently demonstrate that our method outperforms state-of-the-art methods, with an improvement of at most 8.27%. We also validate the effectiveness and practical potential of our method on real-world APIs and defense methods. Furthermore, our method promotes other related tasks, *i.e.*, transfer adversarial attacks.

Keywords: model stealing attack, adversarial attack

1 Introduction

Machine learning models deployed on the cloud can serve users through the application program interfaces (APIs) to improve productivity. Since developing these cloud models is a product of intensive labor and monetary effort, these models are valuable intellectual property and AI companies try to keep them private [16, 24, 21, 33]. However, the exposure of the model’s predictions represents a significant risk as an adversary can leverage this information to steal the model’s functionality, *a.k.a.* model stealing attack [26, 23, 25, 8, 35]. With such an attack, adversaries are able to not only use the stolen model to make a profit, but also mount further adversarial attacks [40, 34]. Besides, the model stealing

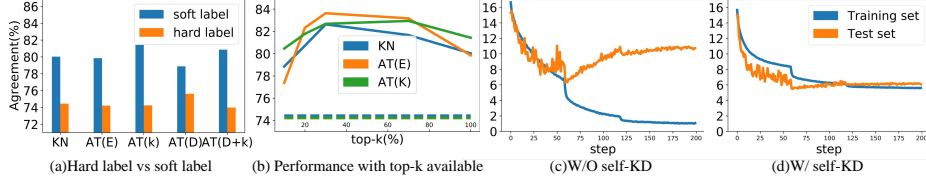


Fig. 1. (a) The test accuracies of previous methods with hard labels are much lower than the ones with soft labels. (KN: KnockoffNets, ‘AT’: ActiveThief, ‘E’: entropy, ‘K’: k-Center, ‘D’: DFAL) (b) The performance decreases as the number of available classes decreases (dotted line : hard-label setting). (c) & (d) Loss curves for training/test set during model training without and with self-KD. All results are on the CIFAR10 dataset.

attacks is a kind of black-box knowledge distillation which is a hot research topic. Studying various mechanisms of model stealing attack is of great interest both to AI companies and researchers.

Previous methods [23, 40, 25, 8] mainly assume the complete probability predictions of the victim model available, while the real-world APIs usually only return partial probability values (top- k predictions) or even the top-1 prediction (*i.e.*, hard label). In this paper, we focus on the more challenging and realistic scenario, *i.e.*, the victim model only outputs the hard labels. However, under this setting, existing methods suffer from a significant performance degradation, even by 30.50% (as shown in the Fig. 1 (a) and the appendix Tab. I).

To investigate the reason for the degradation, we evaluate the performance of attack methods with different numbers of prediction probability categories available and hard labels as in Fig. 1 (b). With the observation that the performance degrades when the top- k information missing, we conclude that the top- k predictions are informative as it indicates the similarity of different categories or multiple objects in the picture, and previous attack methods suffer from such information obscured by the top-1 prediction under the hard-label setting. It motivates us to re-mine this information by eliminating the top-1 prediction. Particularly, we design a *novel CAM-based erasing method*, which erases the important area on the pictures based on the substitute model’s top-1 class activation maps (CAM) [28, 39] and queries the victim model for a new prediction. Note that we can dig out other class information in this sample if the new prediction changes. Otherwise, it proves that the substitute model pays attention to the wrong area. Then we can align the attention of the substitute and the victim model by learning clean samples and the corresponding erased samples simultaneously.

Besides, previous works on the self-Knowledge Distillation (self-KD) [17], calibration [10], and noisy label [37] have pointed out the hard and noisy labels will introduce overfitting and miscalibration. More specifically, the attack algorithms cannot access the training data, and thus can only use the synthetic data or other datasets as a substitute, which is noisy. Therefore, the hard-label

only outputs accessed). [26] first observed that online models could be stolen through multiple queries. After that, due to the practical threat to real-world APIs, several studies paid attention to this problem and proposed many attack algorithms.

These algorithms consist of two stages: 1) constructing a transfer dataset D_T (step 1 in Fig. 2) and 2) training a substitute model. The transfer dataset is constructed based on data synthesis or data selection and then feed into the victim model for labels. Methods based on data synthesis [40, 15, 2] adopt the GAN-based models to generate a virtual dataset. And the substitute model and the GAN model are trained alternatively on this virtual dataset by querying the victim model iteratively. The data selection methods prepare an attack dataset as the data pool, and then sample the most informative data via machine learning algorithms, *e.g.*, reinforcement learning [23] or active learning strategy [25], uncertainty-based strategy [19], k-Center strategy [29], and DFAL strategy [5]. Considering that querying the victim model will be costly, the attacker usually sets a budget on the number of the queries, so the size of the transfer dataset should be limited as well. Previous methods assume the victim model returns a complete probability prediction $f(x)$, which is less practical.

In this paper, we focus on a more practical scenario that is about hard-label $\phi(f(x))$ setting, where ϕ is the truncation function used to truncate the information contained in the victim’s output and return the corresponding one-hot vector:

$$\phi(f(x))_i := \begin{cases} 1 & \text{if } i = \arg \max_n f(x)_n; \\ 0 & \text{otherwise.} \end{cases} \quad (1)$$

With the transfer dataset, the substitute model is optimized by minimizing a loss function \mathcal{L} (*e.g.*, cross-entropy loss function):

$$\begin{cases} \mathbb{E}_{x \sim \mathcal{D}_T} [\mathcal{L}(f(x), \hat{f}(x))], & \text{for soft labels;} \\ \mathbb{E}_{x \sim \mathcal{D}_T} [\mathcal{L}(\phi(f(x)), \hat{f}(x))], & \text{for hard labels.} \end{cases} \quad (2)$$

Knowledge distillation (KD) has been widely studied in machine learning [12, 1, 7], which transfers the knowledge from a teacher model to a student model. Model stealing attacks can be regarded as a black-box KD problem where the victim model is the *teacher* with only outputs accessible and the substitute model is the *student*. The main reason for the success of KD is the *valuable information that defines a rich similarity structure over the data* in the probability prediction [12]. However, for the hard-label setting discussed in this paper, this valuable information is lost. And the main difference between self-KD and regular-KD is that the latter utilizes knowledge from a larger and better teacher model, while the former uses the model self as the teacher. Self-KD has been shown to help improve the model’s generalization ability [17]. Inspired by self-KD, our method tries to dig out the hidden information in the data and models, and then transfers more knowledge to the substitute model.

The erasing-based method, *e.g.*, random erasing (RE) [38, 3], is currently one of the widely used data augmentation methods, which generates training

images with various levels of occlusion, thereby reducing the risk of over-fitting and improving the robustness of the model. Our work is inspired by RE and designs a prior-driven erasing operation, which erases the area corresponding to the hard label to re-mine missing information.

3 Method

The overview of black-box dissector is shown in Fig. 2. In addition to the conventional process (*i.e.*, the transfer dataset D_T constructing in step 1 and the substitute model training in the right), we introduce two key modules: a CAM-driven erasing strategy (step 2.1) and a RE-based self-KD module (step 2.2).

3.1 A CAM-driven erasing strategy

Since the lack of class similarity information degrades the performance of previous methods under the hard-label setting, we try to re-dig out such hidden information. Taking an example from the ILSVRC-2012 dataset for illustration as in Fig. 3. Querying the CUBS200 trained victim model with this image, we get two classes with the highest confidence score: “Anna hummingbird” (0.1364) and “Common yellowthroat” (0.1165), and show their corresponding attention map in the first column of Fig. 3. It is easy to conclude that two different attention regions response for different classes according to the attention map. When training the substitute model with the hard label “Anna hummingbird” and without the class similarity information, the model can not learn from the area related to the “Common yellowthroat” class, which means this area is wasted. To re-dig out the information about the “Common yellowthroat” class, we need to erase the impact of the “Anna hummingbird” class.

To this end, a natural idea is to erase the response area corresponding to the hard label. Since the victim model is a black-box model, we use the substitute model to approximately calculate the attention map instead. If the attention map calculated by the substitute model is inaccurate and the victim model’s prediction on the erased image does not change, although we cannot obtain the class information, we can align the attention map of two models by letting the substitute model learn the original image and the erased one simultaneously.

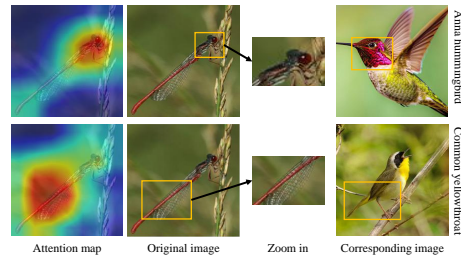


Fig. 3. An example from the ILSVRC-2012 dataset and its attention map corresponding to two most likely class “Anna hummingbird” and “Common yellow throat” on the CUBS200 trained model. The attention areas share similar visual apparent with images of “Anna hummingbird” and “Common yellowthroat”, respectively.

Algorithm 1: Prior-driven Erasing $\psi(I, P)$

Input: Input image I , prior probability P , area of image S , erasing area ratio range s_l and s_h , erasing aspect ratio range r_1 and r_2 .
Output: Erased image I' .

- 1 $S_e \sim \text{Uniform}(s_l, s_h) \times S$, $r_e \sim \text{Uniform}(r_1, r_2)$
- 2 $H_e \leftarrow \sqrt{S_e \times r_e}/2$, $W_e \leftarrow \sqrt{\frac{S_e}{r_e}}/2$
- 3 x_e, y_e sampled randomly according to P
- 4 $I_e \leftarrow (x_e - W_e, y_e - H_e, x_e + W_e, y_e + H_e)$
- 5 $I(I_e) \sim \text{Uniform}(0, 255)$
- 6 $I' \leftarrow I$

The attention maps can be used as sources of additional supervision signal in distillation: encouraging a model’s attention map to be similar to that of another model also leads to the models having similar predictions [36]. To get the attention map, we utilize the Grad-CAM [28] in this paper. With the input image $x \in [0, 1]^d$ and the trained DNN $\mathcal{F}: [0, 1]^d \mapsto \mathbb{R}^N$, we let α_k^c denote the weight of class c corresponding to the k -th feature map, and calculate it as $\alpha_k^c = \frac{1}{Z} \sum_i \sum_j \frac{\partial \mathcal{F}(x)^c}{\partial A_{ij}^k}$, where Z is the number of pixels in the feature map, $\mathcal{F}(x)^c$ is the score of class c and A_{ij}^k is the value of pixel at (i, j) in the k -th feature map. After obtaining the weights corresponding to all feature maps, the final attention map can be obtained as $S_{\text{Grad-CAM}}^c = \text{ReLU}(\sum_k \alpha_k^c A^k)$ via the weighted summation.

To erase the corresponding area, inspired by [38], we define a prior-driven erasing operation as $\psi(I, P)$, shown in Alg. 1, which randomly erases a rectangle region in the image I with random values while the central position of the rectangle region is randomly selected following the prior probability P . The prior probability P is of the same size as the input image and is used to determine the probability of different pixels being erased. Here, we use the attention map from Grad-CAM as the prior. Let $x \in [0, 1]^d$ denote the input image from the transfer set and $S_{\text{Grad-CAM}}^{\arg \max \hat{f}(x)}(x, \hat{f})$ denote the attention map of the substitute model \hat{f} . This CAM-driven erasing operation can be represented:

$$\psi \left(x, S_{\text{Grad-CAM}}^{\arg \max \hat{f}(x)}(x, \hat{f}) \right). \quad (3)$$

We abbreviate it as $\psi(x, S(x, \hat{f}))$. To alleviate the impact of inaccurate Grad-CAM caused by the difference between the substitute model and the victim one, for each image, we perform this operation N times (ψ_i means the i -th erasing) and select the one with the largest difference from the original label. Such a data augment operation helps the erasing process to be more robust.

We use the cross-entropy to calculate the difference between the new label and the original label, and we want to select the sample with the biggest difference. Formally, we define $\Pi(x)$ as the function to select the most different variation of image x :

$$\begin{aligned}
\Pi(x) &:= \psi_k(x, S(x, \hat{f})), \\
\text{where } k &:= \arg \max_{i \in [N]} - \sum_j \phi(f(x))_j \cdot \log \left(\hat{f}(\psi_i(x, S(x, \hat{f})))_j \right) \\
&= \arg \max_{i \in [N]} - \log \left(\hat{f}(\psi_i(x, S(x, \hat{f})))_{\arg \max \phi(f(x))} \right) \\
&= \arg \min_{i \in [N]} \hat{f}(\psi_i(x, S(x, \hat{f})))_{\arg \max \phi(f(x))}.
\end{aligned} \tag{4}$$

Due to the limitation of the number of queries, we cannot query the victim model for each erased image to obtain a new label. We continuously choose the erased image with the highest substitute’s confidence until reaching the budget. To measure the confidence of the model, we adopt the Maximum Softmax Probability (MSP) for its simplicity:

$$\arg \max_{x \sim \mathcal{D}_T} \hat{f}(\Pi(x))_{\arg \max \hat{f}(\Pi(x))}, \tag{5}$$

where \mathcal{D}_T is the transfer set. The erased images selected in this way are most likely to change the prediction class. Then, we query the victim model to get these erased images’ labels and construct an erased sample set \mathcal{D}_E . Note that the substitute model is trained with \mathcal{D}_T to fit the victim model, so it makes sense to use the substitute model to approximately calculate the Grad-CAM. For each sample, if the approximate calculated Grad-CAM is accurate, it means we have erased the correct area, and we can obtain new class information after querying the victim model. If the Grad-CAM is inaccurate, it means the substitute model has paid attention to the wrong area, and we can align the attention map of two models by letting the substitute model learn the original image and the erased one simultaneously. Therefore, regardless of the accuracy of the area we erased, the erased sample can provide information to help the substitute model better approximate the victim model. We show the effect of our method to align the attention map in Fig. 5.

3.2 A random-erasing-based self-KD module

We also find that in training with limited hard-label OOD samples, the substitute model is likely to overfit the training set, which damages its generalization ability [17, 37]. Therefore, based on the above erasing operation, we further design a simple RE-based self-KD method to improve the generalization ability of the substitute model.

Formally, let $x \in [0, 1]^d$ denote the unlabeled input image. We perform the erasing operation with a uniform prior U on it N times, and then average the substitute’s outputs on these erased images as the original image’s pseudo-label:

$$y_p(x, \hat{f}) = \frac{1}{N} \sum_{i=1}^N \hat{f}(\psi_i(x, U)). \tag{6}$$

Algorithm 2: Black-box Dissector

Input: Unlabeled pool D_U , victim model f , maximum number of queries Q .
Output: Substitute model \hat{f} .

```

1 Initialize  $q \leftarrow 0, D_T \leftarrow \emptyset, D_E \leftarrow \emptyset$ 
2 while  $q < Q$  do
3   // Step 1
4   Select samples from  $D_U$  according to budget and query  $f$  to update  $D_T$ 
5    $q = q + \text{budget}$ 
6    $\mathcal{L} = \sum_{x \in D_T} \mathcal{L}'(\phi(f(x)), \hat{f}(x))$ 
7    $\hat{f} \leftarrow \text{update}(\hat{f}, \mathcal{L})$ 
8   // A CAM-driven erasing strategy (step 2.1)
9   Erase samples in  $D_T$  according to Eq. 4
10  Choose samples from erased samples according to Eq. 5 and budget
11  Query  $f$  to get labels and update  $D_E$ 
12   $\mathcal{L} = \sum_{x \in D_T \cup D_E} \mathcal{L}'(\phi(f(x)), \hat{f}(x))$ 
13   $\hat{f} \leftarrow \text{update}(\hat{f}, \mathcal{L})$ 
14   $q = q + \text{budget}$  (Check if  $q < Q$ )
15  // A random-erasing-based self-KD (step 2.2)
16  Select samples from  $D_U$ 
17  Get pseudo-labels according to Eq. 6 and construct a pseudo-label set  $D_P$ 
18   $\mathcal{L} = \sum_{x \in D_T \cup D_E} \mathcal{L}'(\phi(f(x)), \hat{f}(x)) + \sum_{x \in D_P} \mathcal{L}'(y_p(x, \hat{f}), \hat{f}(x))$ 
19   $\hat{f} \leftarrow \text{update}(\hat{f}, \mathcal{L})$ 
20 end

```

This is a type of consistency regularization, which enforces the model to have the same predictions for the perturbed images and enhances the generalization ability. With Eq.6, we construct a new soft pseudo label set $D_P = \{(x, y_p(x, \hat{f})), \dots\}$.

With the transfer set D_T , the erased sample set D_E , and the pseudo-label set D_P , we train new substitute model using the ensemble of the victim model and the previous substitute model as the teacher. Our final objective function is:

$$\min \mathcal{L} = \min \left[\sum_{x \in D_T \cup D_E} \mathcal{L}'(\phi(f(x)), \hat{f}(x)) + \sum_{x \in D_P} \mathcal{L}'(y_p(x, \hat{f}), \hat{f}(x)) \right]. \quad (7)$$

where \mathcal{L}' can be commonly used loss functions, *e.g.*, cross-entropy loss function.

To sum up, we built our method on the conventional process of the model stealing attack (step 1), and proposed a CAM-driven erasing strategy (step 2.1) and a RE-based self-KD module (step 2.2) unified by a novel erasing method. The former strategy digs out missing information between classes and aligns the attention while the latter module helps to mitigate overfitting and enhance the generalization. We name the whole framework as *black-box dissector* and present the algorithm detail of it in Alg. 2.

4 Experiments

4.1 Experiment settings

In this subsection, we introduce our experiment settings, including victim model, model architectures, attack dataset and training process.

Victim model. The victim models we used (ResNet-34 [11]) are trained on four datasets, namely, CIFAR10 [18], SVHN [22], Caltech256 [9], and CUBS200 [32], and their test accuracy are 91.56%, 96.45%, 78.40%, and 77.10%, respectively. All models are trained using the SGD optimizer with momentum (of 0.5) for 200 epochs with a base learning rate of 0.1 decayed by a factor of 0.1 every 30 epochs. In order to create an online deployment scenario, these models are all: image in, one-hot predictions out. Following [23, 25, 40], we use the same architecture for the substitute model and will analyze the impact of different architectures.

Attack dataset. We use 1.2M images without labels from the ILSVRC-2012 challenge [27] as the attack dataset. In a real attack scenario, the attacker may use pictures collected from the Internet, and the ILSVRC-2012 dataset can simulate this scenario well. Note that we resize all images in the attack dataset to fit the size of the target datasets, which is similar to the existing setting [23, 25, 40].

Training process. We use the SGD optimizer with momentum (of 0.9) for 200 epochs and a base learning rate of $0.02 \times \frac{\text{batchsize}}{128}$ decayed by a factor of 0.1 every 60 epochs. The weight decay is set to 5×10^{-4} for small datasets (CIFAR10 [18] and SVHN [22]) and 0 for others. We set up a query sequence $\{0.1K, 0.2K, 0.5K, 0.8K, 1K, 2K, 5K, 10K, 20K, 30K\}$ as the iterative maximum query budget, and stop the sampling stage whenever reaching the budget at each iteration. For fairness, all experiments will be conducted in accordance with this sequence. And, the model is trained from scratch for each iteration.

Baselines and evaluation metric. We mainly compare our method with KnockoffNets [23] and ActiveThief [25]. For KnockoffNets, we use the source codes provided kindly by the authors. Follow [14], we mainly report the test accuracy (Acc) as the evaluation metric. We also report the *Agreement* metric proposed by [25] which counts how often the prediction of the substitute model is the same as the victim’s as a supplement.

4.2 Experiment results

We first report the performance of our method compared with previous methods. Then, we analyze the performance of our method when encountering two SOTA defense methods (*i.e.*, the adaptive misinformation [16] and the prediction poisoning [24]) and real-world online APIs. After that, we conduct ablation experiments to analyze the contribution of each module. Finally, we also analyze the effect of different model structures and demonstrate the transferability of adversarial samples generated on substitute models obtained by different methods.

Effectiveness of our method. As in Tab. 1, the test accuracy and agreement of our method are all better than the previous methods. We also plot the

Table 1. The agreement and test accuracy (in %) of each method under 30k queries. For our model, we report the average accuracy as well as the standard deviation computed over 5 runs. (**Boldface**: the best value, *italics*: the second best value.)

Method	CIFAR10		SVHN		Caltech256		CUBS200	
	Agreement	Acc	Agreement	Acc	Agreement	Acc	Agreement	Acc
KnockoffNets	75.32	74.44	85.00	84.50	57.64	55.28	30.01	28.03
ActiveThief(Entropy)	75.26	74.21	90.47	89.85	56.28	54.14	32.05	29.43
ActiveThief(k-Center)	75.71	74.24	81.45	80.79	61.19	58.84	37.68	34.64
ActiveThief(DFAL)	76.72	75.62	84.79	84.17	46.92	44.91	20.31	18.69
ActiveThief(DFAL+k-Center)	74.97	73.98	81.40	80.86	55.70	53.69	26.60	24.42
Ours+Random	82.14 ±0.16	80.47 ±0.02	92.33 ±0.47	91.57 ±0.29	<i>63.61</i> ±0.53	<i>61.41</i> ±0.39	<i>39.07</i> ±0.26	<i>36.28</i> ±0.44
Ours+k-Center	<i>80.84</i> ±0.21	<i>79.27</i> ±0.15	<i>91.47</i> ±0.09	<i>90.68</i> ±0.14	66.34 ±0.52	63.75 ±0.49	48.46 ±0.55	44.43 ±0.42

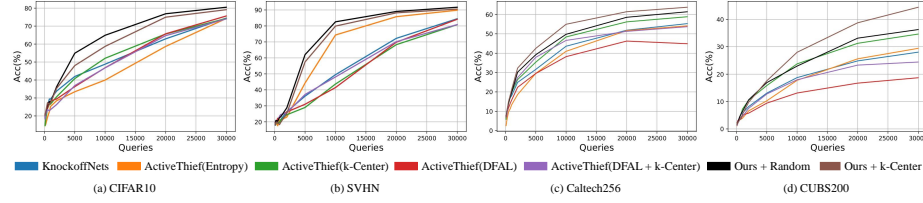


Fig. 4. Curves of the test accuracy versus the number of queries.

curves of the test accuracy versus the number of queries in Fig. 4. The performance of our method consistently outperforms other methods throughout the process. Since our method does not conflict with the previous sample selection strategy, they can be used simultaneously to further improve the performance of these attacks. Here, we take the k-Center algorithm as an example. Note that, with or without the sample selection strategy, our method beats the previous methods by a large margin. Particularly, the test accuracies of our method are 4.85%, 1.72%, 3.88%, and 8.27% higher than the previous best method, respectively. And the agreement metric shares similar results. It is also interesting that it is less necessary to use the k-Center algorithm on datasets with a small number of classes (*i.e.*, CIFAR10 and SVHN). While for the datasets with a large number of classes, the k-Center algorithm can make the selected samples better cover each class and improve the effectiveness of the method.

Ability to evade the SOTA defense method. Here we evaluate two SOTA perturbation-based defense method, the adaptive misinformation [16] and the prediction poisoning [24]. The adaptive misinformation [16] introduces an Out-Of-Distribution (OOD) detection module based on the maximum predicted value and punishes the OOD samples with a perturbed model $f'(\cdot; \theta')$. This perturbed model $f'(\cdot; \theta')$ is trained with $\arg \min_{\theta'} \mathbb{E}_{(x,y)} [-\log(1 - f'(x; \theta')_y)]$ to minimize the probability of the correct class. Finally, the output will be:

$$y' = (1 - \alpha)f(x; \theta) + (\alpha)f'(x; \theta'), \quad (8)$$

where $\alpha = 1/(1 + e^{\nu(\max f(x; \theta) - \tau)})$ with a hyper-parameter ν is the coefficient to control how much correct results will be returned, and τ is the threshold used for OOD detection. The model returns incorrect predictions for the OOD sam-

Table 2. Ability to evade the state-of-the-art defense methods (adaptive misinformation and prediction poisoning) on CIFAR10 dataset. The larger the threshold, the better the defence effect while the low victim model’s accuracy (threshold 0 means no defence). Our method evades the defense best, and the self-KD part makes a great difference.

Method	No defence				Adaptive misinformation		Prediction poisoning	
Threshold	0	0.5	0.7	0.9	0.5	0.8	0.5	0.8
KnockoffNets	74.44%	74.13%	73.61%	54.98%	71.83%	58.01%		
ActiveThief(k-Center)	74.24%	69.14%	59.78%	50.19%	73.75%	60.89%		
ActiveThief(Entropy)	74.21%	71.61%	64.84%	51.07%	72.07%	65.83%		
Ours	80.47%	79.95%	78.25%	74.40%	80.01%	79.23%		
Ours w/o self-KD	79.02%	78.66%	73.61%	61.81%	78.87%	76.49%		
victim model	91.56%	91.23%	89.10%	85.14%	91.56%	89.45%		

ples without having much impact on the in-distribution samples. The prediction poisoning [24] is also a perturbation-based defense method, which perturb the posterior probabilities y to make the adversarial gradient signal that maximally deviates from the original gradient. As shown by the following equation:

$$\max_{\tilde{y}} \left\| \frac{G^T \tilde{y}}{\|G^T \tilde{y}\|_2} - \frac{G^T y}{\|G^T y\|_2} \right\|_2^2 \quad (9)$$

where $G = \nabla_w \log F(x; w)$, y is the posterior probabilities and \tilde{y} is the perturbed posterior probabilities.

We choose three values of the threshold τ in the adaptive misinformation and two values of the threshold ϵ in the prediction poisoning to compare the effects of our method with the previous methods. The threshold value of 0 means no defence. The result is shown in Tab. 2. Compared with other methods, adaptive misinformation and prediction poisoning are almost ineffective to our method. Furthermore, we find that if we remove the self-KD in our method, the performance is greatly reduced. We conclude that this is because these two defence methods add noise labels to the substitute model’s training dataset, and self-KD can alleviate the overfitting of the substitute model to the training dataset, making these two defence methods not effective enough.

Ablation study. To evaluate the contribution of different modules in our method, we conduct the ablation study on CIFAR10 dataset and show the results in Tab. 3. We first separately remove the two modules we designed to verify their role. If the CAM-driven erasing strategy is removed, the performance of our method will be greatly reduced, showing that it has an indispensable position in our method. We also give some visual examples in Fig. 5 to demonstrate that this strategy can help align the attention of two models. As depicted in the Fig. 5, at the beginning time, the substitute model learns the wrong attention map. Along with the iterative training stages, the attention area of the substitute model tends to fit the victim model’s, which conforms to our intention. We further remove the self-KD module to evaluate its performance. It can be found from Fig. 1 and Tab. 3 that the self-KD can improve the generalization of our method

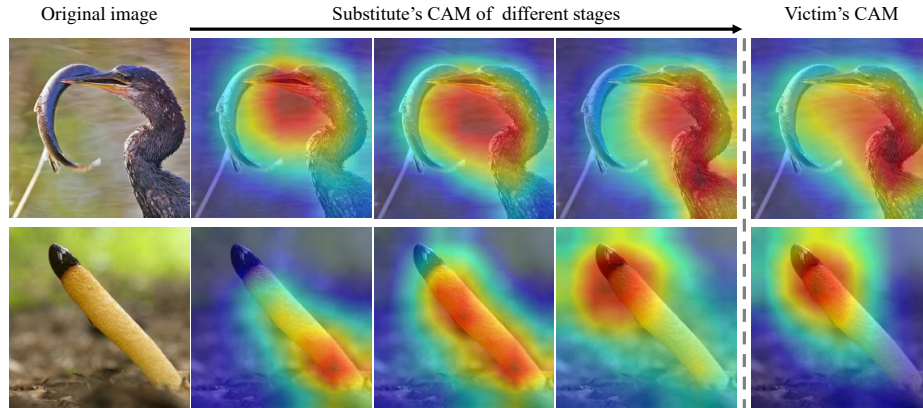


Fig. 5. The visualized attention maps of the victim model and different stages substitute models using the Grad-CAM. Along with the training stages, the attention map of the substitute model tends to fit the victim model’s.

Table 3. An ablation experiment showing the effectiveness of the two modules we designed on CIFAR10 dataset under 30k queries. We use some commonly used regularization methods to replace the two modules we designed, and the results show that the two modules are better than the traditional regularization methods.

Method	ACC
Ours	80.47%
+ $2 \times$ weight decay	81.65%
- CAM-driven erasing	76.12%
- CAM-driven erasing + CutOut	77.11%
- self-KD	79.02%
- self-KD + $2 \times$ weight decay	80.09%
- self-KD + CutOut	78.91%
- self-KD + label smoothing (0.9)	78.22%
- self-KD + label smoothing (0.8)	77.46%

and further improve the performance. Later, in order to prove that these modules are better than some commonly used regularization methods, such as CutOut, label smoothing, we use these methods to replace the modules we designed. Note that the weight decay we used before followed the setting of baseline, so here we also test the effect of a large weight decay. The results are shown in Tab. 3, where ” $2 \times$ weight decay” represents the expansion of the weight decay to twice the original and the ”label smoothing (α)” means smooth the hard-label according to the hyperparameter α . First, we replace the CAM-driven erasing with random erasing (CutOut), which brings 3.36% performance degradation. We believe that using Grad-CAM as a prior is more effective than random. Then we use data augmentation (CutOut) and label smoothing to replace the self-KD, while both show less competitive. We conclude that they destroy the information need by

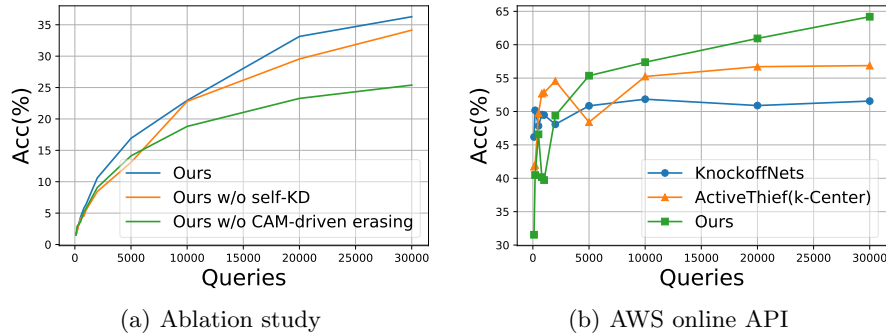


Fig. 6. (a) An ablation study on CUBS200 dataset for the contribution of the CAM-driven erasing and the self-KD. (b) The experiment on a real-word online API.

Table 4. Transferability of adversarial samples generated with PGD attack on the substitute models.

Method	Substitute’s architecture				
	ResNet-34	ResNet-18	ResNet-50	VGG-16	DenseNet
KnockoffNets	57.85%	63.33%	52.04%	42.88%	60.77%
ActiveThief(k-Center)	57.44%	57.90%	57.01%	16.49%	60.72%
ActiveThief(Entropy)	63.56%	66.76%	58.19%	55.43%	62.05%
Ours	76.63%	74.10%	74.28%	67.03%	66.96%

the CAM-driven erasing, e.g., erasing the attention map or hiding information in other classes by making them equal. The result also shows the effectiveness of the self-KD module we designed. In addition, we also perform a simple ablation experiment on the CUBS200 dataset. The results are shown in Fig. 6 (a) and are similar to those on the CIFAR10 dataset.

Stealing functionality of a real-world API. We validate our method is applicable to real-world APIs. The AWS Marketplace is an online store that provides a variety of trained ML models for users. It can only be used in the form of a black-box setting. We choose a popular model (waste classifier⁶) as the victim model. We use ILSVRC-2012 dataset as the attack dataset and choose another small public waste classifier dataset⁷, containing 2,527 images as the test dataset. The hyperparameter settings remain the same as before. As in Fig. 6 (b), the substitute model obtained by our method achieves 12.63% and 7.32% improvements in test accuracy compared with two previous methods, which show our method has stronger practicality in the real world.

Transferability of adversarial samples. Though with the dominant performance on a wide range of tasks, deep neural networks are shown to be vulnerable to imperceptible perturbations, *i.e.*, adversarial examples [31, 6]. Since

⁶ For the purpose of protecting privacy, we hide the specific information of the victim model.

⁷ <https://github.com/garythung/trashnet>

the model stealing attack can obtain a functionally similar substitute model, some previous works (*e.g.*, JBDA [26], DaST [40] and ActiveThief [25]) used this substitute model to generate adversarial samples and then performed the transferable adversarial attack on the victim model. We argue that a more similar substitute model leads to a more successful adversarial attacks. We test the transferability of adversarial samples on the test set of the CIFAR10 dataset. Keeping the architecture of the victim model as the ResNet-34, we evaluate the attack success rate of adversarial samples generated from different substitute models (*i.e.*, ResNet-34, ResNet-18, ResNet-50 [11], VGG-16 [30], DenseNet [13]). All adversarial samples are generated using Projected Gradient Descent (PGD) attack [20] with maximum L_∞ -norm of perturbations as 8/255. As shown in Tab. 4, the adversarial samples generated by our substitute models have stronger transferability in all substitute’s architectures, with 4.91% – 16.20% improvements than other methods. This again proves that our method is more practical in real-world scenarios.

5 Conclusion

We investigated the problem of model stealing attacks under the hard-label setting and pointed out why previous methods are not effective enough. We presented a new method, termed *black-box dissector*, which contains a CAM-driven erasing strategy and a RE-based self-KD module. We showed its superiority on four widely-used datasets and verified the effectiveness of our method with defense methods, real-world APIs, and the downstream adversarial attack. Though focusing on image data in this paper, our method is general for other tasks as long as the CAM and similar erasing method work, *e.g.*, synonym saliency words replacement for NLP tasks [4]. We believe our method can be easily extended to other fields and inspire future researchers. Model stealing attack poses a threat to the deployed machine learning models. We hope this work will draw attention to the protection of deployed models and furthermore shed more light on the attack mechanisms and prevention methods. Additionally, transformer-based classifiers are becoming hot, and their security issues should also be paid attention to. This kind of classifier divides the images into patches and our method works by erasing parts of images, it is more convenient for us to align the attention map by masking the patch and mine the missing information. We will validate this idea in the further work.

Acknowledgments

This work was supported by the National Science Fund for Distinguished Young Scholars (No.62025603), the National Natural Science Foundation of China (No. U21B2037, No. 62176222, No. 62176223, No. 62176226, No. 62072386, No. 62072387, No. 62072389, and No. 62002305), Guangdong Basic and Applied Basic Research Foundation (No.2019B1515120049), and the Natural Science Foundation of Fujian Province of China (No.2021J01002).

References

1. Anil, R., Pereyra, G., Passos, A., Ormandi, R., Dahl, G.E., Hinton, G.E.: Large scale distributed neural network training through online distillation. arXiv preprint arXiv:1804.03235 (2018)
2. Barbalau, A., Cosma, A., Ionescu, R.T., Popescu, M.: Black-box ripper: Copying black-box models using generative evolutionary algorithms. In: NeurIPS (2020)
3. DeVries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. arXiv preprint arXiv:1708.04552 (2017)
4. Dong, X., Luu, A.T., Ji, R., Liu, H.: Towards robustness against natural language word substitutions. In: ICLR (2021)
5. Ducoffe, M., Precioso, F.: Adversarial active learning for deep networks: a margin based approach. In: ICML (2018)
6. Fang, S., Li, J., Lin, X., Ji, R.: Learning to learn transferable attack. In: AAAI (2022)
7. Furlanello, T., Lipton, Z.C., Tschannen, M., Itti, L., Anandkumar, A.: Born again neural networks. In: ICML (2018)
8. Gong, X., Chen, Y., Yang, W., Mei, G., Wang, Q.: Inversenet: Augmenting model extraction attacks with training data inversion. In: IJCAI (2021)
9. Griffin, G., Holub, A., Perona, P.: Caltech-256 object category dataset (2007)
10. Guo, C., Pleiss, G., Sun, Y., Weinberger, K.Q.: On calibration of modern neural networks. In: ICML (2017)
11. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
12. Hinton, G., Vinyals, O., Dean, J.: Distilling the knowledge in a neural network. In: NIPS Deep Learning Workshop (2015)
13. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR (2017)
14. Jagielski, M., Carlini, N., Berthelot, D., Kurakin, A., Papernot, N.: High accuracy and high fidelity extraction of neural networks. In: 29th Usenix Security (2020)
15. Kariyappa, S., Prakash, A., Qureshi, M.: Maze: Data-free model stealing attack using zeroth-order gradient estimation. arXiv preprint arXiv:2005.03161 (2020)
16. Kariyappa, S., Qureshi, M.K.: Defending against model stealing attacks with adaptive misinformation. In: CVPR (2020)
17. Kim, K., Ji, B., Yoon, D., Hwang, S.: Self-knowledge distillation: A simple way for better generalization. arXiv preprint arXiv:2006.12000 (2020)
18. Krizhevsky, A., Hinton, G., et al.: Learning multiple layers of features from tiny images (2009)
19. Lewis, D.D., Gale, W.A.: A sequential algorithm for training text classifiers. In: SIGIR (1994)
20. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: ICLR (2018)
21. Maini, P., Yaghini, M., Papernot, N.: Dataset inference: Ownership resolution in machine learning. In: ICLR (2021)
22. Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., Ng, A.Y.: Reading digits in natural images with unsupervised feature learning. In: NIPS Workshop on Deep Learning and Unsupervised Feature Learning (2011)
23. Orekondy, T., Schiele, B., Fritz, M.: Knockoff nets: Stealing functionality of black-box models. In: CVPR (2019)

24. Orekondy, T., Schiele, B., Fritz, M.: Prediction poisoning: Towards defenses against dnn model stealing attacks. In: ICLR (2019)
25. Pal, S., Gupta, Y., Shukla, A., Kanade, A., Shevade, S., Ganapathy, V.: Activethief: Model extraction using active learning and unannotated public data. In: AAAI (2020)
26. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against machine learning. In: ACM AsiACCS (2017)
27. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. IJCV (2015)
28. Selvaraju, R.R., Cogswell, M., Das, A., Vedantam, R., Parikh, D., Batra, D.: Grad-cam: Visual explanations from deep networks via gradient-based localization. In: ICCV (2017)
29. Sener, O., Savarese, S.: Active learning for convolutional neural networks: A core-set approach. In: ICLR (2018)
30. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: ICLR (2015)
31. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. In: ICLR (2014)
32. Wah, C., Branson, S., Welinder, P., Perona, P., Belongie, S.: The caltech-ucsd birds-200-2011 dataset (2011)
33. Wang, X., Xiang, Y., Gao, J., Ding, J.: Information laundering for model privacy. In: ICLR (2021)
34. Yang, J., Jiang, Y., Huang, X., Ni, B., Zhao, C.: Learning black-box attackers with transferable priors and query feedback. In: NeurIPS (2020)
35. Yu, H., Yang, K., Zhang, T., Tsai, Y.Y., Ho, T.Y., Jin, Y.: Cloudleak: Large-scale deep learning models stealing through adversarial examples. In: NDSS (2020)
36. Zagoruyko, S., Komodakis, N.: Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In: ICLR (2017)
37. Zhang, C., Bengio, S., Hardt, M., Recht, B., Vinyals, O.: Understanding deep learning requires rethinking generalization. In: ICLR (2017)
38. Zhong, Z., Zheng, L., Kang, G., Li, S., Yang, Y.: Random erasing data augmentation. In: AAAI (2020)
39. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A., Torralba, A.: Learning deep features for discriminative localization. In: CVPR (2016)
40. Zhou, M., Wu, J., Liu, Y., Liu, S., Zhu, C.: Dast: Data-free substitute training for adversarial attacks. In: CVPR (2020)