

---

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

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

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

## 16 1 Introduction

17 Machine learning models deployed on the cloud can serve users through the application program  
18 interfaces (APIs) to improve productivity. Since developing these cloud models is a product of  
19 intensive labor and monetary effort, these models are valuable intellectual property and AI companies  
20 try to keep them private. However, the exposure of the model’s predictions represents a significant  
21 risk as an adversary can leverage this information to steal the model’s functionality, *a.k.a.* model  
22 stealing attack [22, 20, 21]. With such an attack, adversaries are able to not only use the stolen model  
23 to make a profit, but also mount further adversarial attacks [34, 29]. Besides, the model stealing  
24 attacks is a kind of black-box knowledge distillation which is a hot research topic. Studying various  
25 mechanisms of model stealing attack is of great interest both to AI companies and researchers.

26 Previous methods [20, 34, 21] mainly assume the complete probability predictions of the victim  
27 model available, while the real-world APIs usually only return partial probability values (*top-k*  
28 predictions) or even the top-1 prediction (*i.e.*, hard label). In this paper, we focus on the more  
29 challenging and realistic scenario, *i.e.*, the victim model only outputs the hard labels. However, under  
30 this setting, existing methods suffer from a significant performance degradation, even by 30.50% (as  
31 shown in the Fig. 1 (a) and the appendix Tab. I).

32 To investigate the reason for the degradation, we evaluate the performance of attack methods with  
33 different numbers of prediction probability categories available and hard labels as in Fig. 1 (b). With  
34 the observation that the performance degrades when the *top-k* information missing, we conclude  
35 that the *top-k* predictions are informative as it indicates the similarity of different categories or  
36 multiple objects in the picture, and previous attack methods suffer from such information obscured

37 by the top-1 prediction under the hard-label setting. It motivates us to re-mine this information by  
 38 eliminating the top-1 prediction. Particularly, we design *a novel CAM-based erasing method*, which  
 39 erases the important area on the pictures based on the substitute model’s top-1 class activation maps  
 40 (CAM) [24, 33] and queries the victim model for a new prediction. Note that we can dig out other  
 41 class information in this sample if the new prediction changes. Otherwise, it proves that the substitute  
 42 model pays attention to the wrong area. Then we can align the attention of the substitute and the  
 43 victim model by learning clean samples and the corresponding erased samples simultaneously.

44 Besides, previous works on the self-Knowledge Distillation (self-KD) [15], calibration [8], and noisy  
 45 label [31] have pointed out the hard and noisy labels will introduce overfitting and miscalibration. More  
 46 specifically, the attack algorithms cannot access the training data, and thus can only use the synthetic data or  
 47 other datasets as a substitute, which is noisy. Therefore, the hard-label setting will suffer from overfitting,  
 48 which leads to worse performance, and we verify it by plotting the loss curves in Fig. 1 (c). To mitigate  
 49 this problem, we introduce *a simple self-knowledge distillation module with random erasing (RE)* to utilize  
 50 soft labels for generalization. Particularly, we randomly erase one sample a certain number of times, query  
 51 the substitute model for soft-label outputs, and take the average value of these outputs as the pseudo-label.  
 52 After that, we use both hard labels from the victim model and pseudo labels from the previous substitute model  
 53 to train a new substitute model. Therefore, we can also consider the ensemble of the two models as the teacher  
 54 in knowledge distillation. As in Fig. 1 (d), such a module helps generalization and better performance.

72 In summary, we propose a novel model stealing framework termed *black-box dissector*, which  
 73 includes a CAM-driven erasing strategy and a RE-based self-KD module. Our method is orthogonal  
 74 to previous approaches [20, 21] and can be integrated with them. The experiments on four widely-  
 75 used datasets demonstrate our method achieves 43.04 – 90.57% test accuracy (47.60 – 91.37%  
 76 agreement) to the victim model, which is at most 8.27% higher than the state of the art method.  
 77 We also proved that our method can defeat popular defense methods and is effective for real-world  
 78 APIs like services provided by Amazon Web Services (AWS). Furthermore, our method promotes  
 79 downstream tasks, *i.e.*, transfer adversarial attack, with 4.91% – 16.20% improvement.

## 80 2 Background and Notions

81 **Model stealing attack** is aim to find a substitute model  $\hat{f}: [0, 1]^d \mapsto \mathbb{R}^N$  that performs as similarly  
 82 as possible to the black-box victim model  $f: [0, 1]^d \mapsto \mathbb{R}^N$  (with only outputs accessed). Papernot  
 83 et al. [22] first observed that online models could be stolen through multiple queries. After that, due  
 84 to the practical threat to real-world APIs, several studies paid attention to this problem and proposed  
 85 many attack algorithms.

86 These algorithms consist of two stages: 1) constructing a transfer dataset  $D_T$  (step 1 in Fig. 2) and  
 87 2) training a substitute model. The transfer dataset is constructed based on data synthesis or data  
 88 selection and then feed into the victim model for labels. Methods based on data synthesis [34, 14, 2]  
 89 adopt the GAN-based models to generate a virtual dataset. And the substitute model and the GAN  
 90 model are trained alternatively on this virtual dataset by querying the victim model iteratively. The  
 91 data selection methods prepare an attack dataset as the data pool, and then sample the most informative

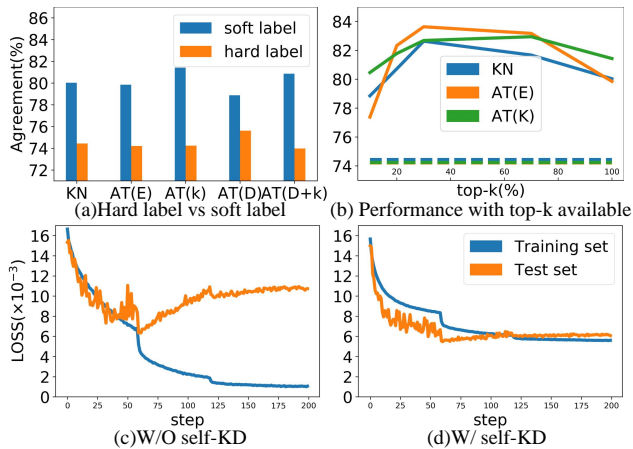


Figure 1: (a) The test accuracies of previous methods with hard labels are much lower than the ones with soft labels. (KN: KnockoffNets, ‘AT’: ActiveThief, ‘E’: entropy, ‘K’: k-Center, ‘D’: DFAL) (b) The performance decreases as the number of available classes decreases (dotted line : hard-label setting).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

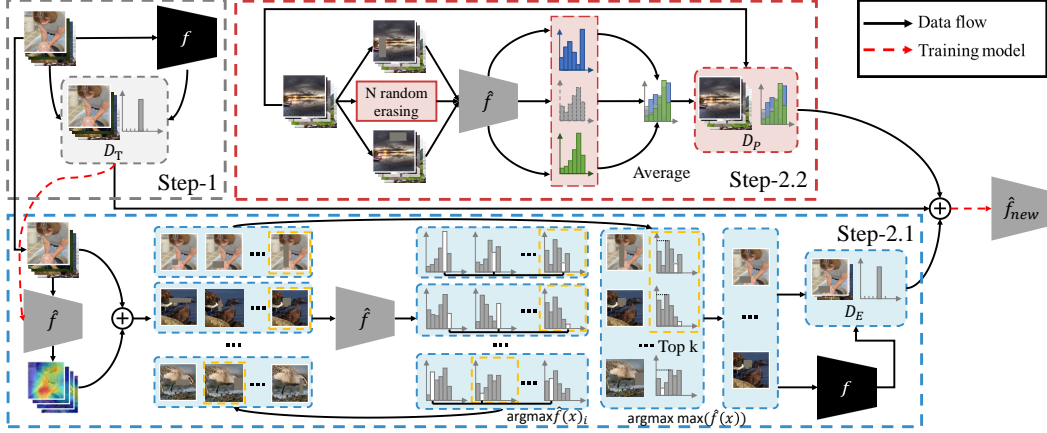


Figure 2: Details of our proposed black-box dissector with a CAM-driven erasing strategy (step 2.1) and a RE-based self-KD module (step 2.2). In step 2.1, the images in transfer set  $D_T$  are erased according to the Grad-CAM, and we selected the erased images with the largest difference from the original images according to the substitute model’s outputs. In step 2.2, we randomly erase the unlabeled image  $N$  times, and then average the outputs of the  $N$  erased images by the substitute model as the pseudo-label.

92 data via machine learning algorithms, *e.g.*, reinforcement learning [20] or active learning strategy [21],  
 93 uncertainty-based strategy [17], k-Center strategy [25], and DFAL strategy [5]. Considering that  
 94 querying the victim model will be costly, the attacker usually sets a budget on the number of the  
 95 queries, so the size of the transfer dataset should be limited as well. Previous methods assume the  
 96 victim model returns a complete probability prediction  $f(x)$ , which is less practical.

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

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

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

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

102 **Knowledge distillation** (KD) has been widely studied in machine learning [10, 1, 6], which transfers  
 103 the knowledge from a teacher model to a student model. Model stealing attacks can be regarded as a  
 104 black-box KD problem where the victim model is the *teacher* with only outputs accessible and the  
 105 substitute model is the *student*. The main reason for the success of KD is the *valuable information*  
 106 *that defines a rich similarity structure over the data* in the probability prediction [10]. However,  
 107 for the hard-label setting discussed in this paper, this valuable information is lost. Inspired by KD,  
 108 our method tries to dig out the hidden information in the data and models, and then transfers more  
 109 knowledge to the substitute model.

110 **The erasing-based method**, *e.g.*, random erasing (RE) [32, 3], is currently one of the widely used  
 111 data augmentation methods, which generates training images with various levels of occlusion, thereby  
 112 reducing the risk of over-fitting and improving the robustness of the model. Our work is inspired  
 113 by RE and designs a prior-driven erasing operation, which erases the area corresponding to the hard  
 114 label to re-mine missing information.

115 **3 Method**

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

120 **3.1 A CAM-driven erasing strategy**

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

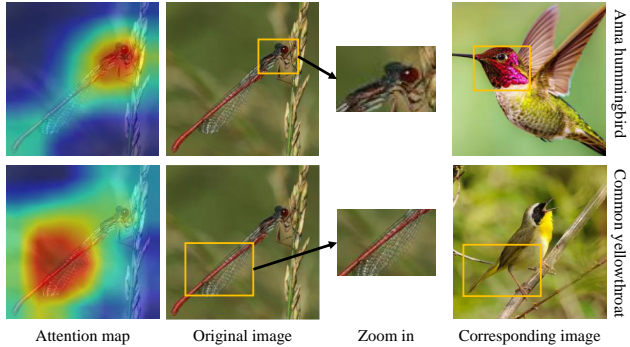


Figure 3: An example from the ILSVRC-2012 dataset and its attention map corresponding to two most likely class “Anna humming bird” and “Common yellow throat” on the CUBS200 trained model. The attention areas share similar visual apparent with images of “Anna humming bird” and “Common yellow throat”,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

144 To this end, a natural idea is to erase the response area corresponding to the hard label. Since the  
 145 victim model is a black-box model, we use the substitute model to approximately calculate the  
 146 attention map instead. If the attention map calculated by the substitute model is inaccurate and the  
 147 victim model’s prediction on the erased image does not change, we can also align the attention map of  
 148 two models by letting the substitute model learn the original image and the erased one simultaneously.  
 149 The attention map is also a kind of supervision signal pushing two models to be similar [30]. To  
 150 get the attention map, we utilize the Grad-CAM [24] in this paper. With the input image  $x \in [0, 1]^d$   
 151 and the trained DNN  $\mathcal{F}: [0, 1]^d \mapsto \mathbb{R}^N$ , we let  $\alpha_k^c$  denote the weight of class  $c$  corresponding to the  
 152  $k$ -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  
 153 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  
 154 map. After obtaining the weights corresponding to all feature maps, the final attention map can be  
 155 obtained as  $S_{\text{Grad-CAM}}^c = \text{ReLU}(\sum_k \alpha_k^c A^k)$  via weighted summation.

156 To erase the corresponding area, inspired by [32], we define a prior-driven erasing operation as  
 157  $\psi(I, P)$ , shown in Alg. 1, which randomly erases a rectangle region in the image  $I$  with random  
 158 values while the central position of the rectangle region is randomly selected following the prior  
 159 probability  $P$ . The prior probability  $P$  is of the same size as the input image and is used to determine  
 160 the probability of different pixels being erased. Here, we use the attention map from Grad-CAM as  
 161 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  
 162 the attention map of the substitute model  $\hat{f}$ . This CAM-driven erasing operation can be represented:

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

163 We abbreviate it as  $\psi(x, S(x, \hat{f}))$ . To alleviate the impact of inaccurate CAM caused by the difference  
 164 between the substitute model and the victim one, for each image, we perform this operation  $N$  times  
 165 ( $\psi_i$  means the  $i$ -th erasing) and select the one with the largest difference from the original label.

---

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

---

**Input:** Input image  $I$ , prior probability  $P$ , image size  $W$  and  $H$ , 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  $S_e \leftarrow \text{Rand}(s_l, s_h) \times S, r_e \leftarrow \text{Rand}(r_1, r_2)$ <sup>1</sup>
  - 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  $I(I_e) \leftarrow \text{Rand}(0, 255)$
  - 6  $I' \leftarrow I$
- 

166 Such a data augment operation helps the erasing process to be more robust. We use the cross-entropy  
167 to calculate the difference between the new label and the original label, and we want to select the  
168 sample with the biggest difference. Formally, we define  $\Pi(x)$  as the function to select the most  
169 different variation of image  $x$ :

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

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

$$\begin{aligned} &\arg \max_{x \sim \mathcal{D}_T} MSP \left( \hat{f}(\Pi(x)) \right) \\ &= \arg \max_{x \sim \mathcal{D}_T} \hat{f}(\Pi(x))_{\arg \max \hat{f}(\Pi(x))}, \end{aligned} \tag{5}$$

174 where  $\mathcal{D}_T$  is the transfer set. The erased images selected in this way are most likely to change the  
175 prediction class. Then, we query the victim model to get these erased images' labels and construct  
176 an erased sample set  $\mathcal{D}_E$ . Note that when the victim model's predictions on the erased images  
177 change, it means our erasing method does dig out other related class information in the sample. With  
178 the unchanged predictions, it points out the attentions of the substitute model and the victim are  
179 inconsistent. Though wrong attention areas erased, training with these samples benefits aligning the  
180 attentions of two models. As [30] stated, the attention alignment can help more powerful KD.

### 181 3.2 A random-erasing-based self-KD module

182 We also find that in training with limited hard-label OOD samples, the substitute model is likely  
183 to overfit the training set, which damages its generalization ability [15, 31]. Therefore, based on  
184 the above erasing operation, we further design a simple RE-based self-KD method to improve the  
185 generalization ability of the substitute model.

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

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

---

<sup>1</sup>Rand( $a, b$ ) returns an evenly distributed random real number in the range of  $a$  to  $b$ .

---

**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   $q = q + \text{budget}$ 
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
```

---

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

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

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

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

196 To sum up, we built our method on the conventional process of the model stealing attack (step  
197 1), and proposed a CAM-driven erasing strategy (step 2.1) and a RE-based self-KD module (step  
198 2.2) unified by a novel erasing method. The former strategy digs out missing information between  
199 classes and aligns the attention while the latter module helps to mitigate overfitting and enhance the  
200 generalization. We name the whole framework as *black-box dissector* and present the algorithm  
201 detail of it in Alg. 2.

## 202 4 Experiment

### 203 4.1 Experiment settings

204 **Victim model.** The victim models we used (ResNet-34 [9]) are trained on four datasets, namely,  
205 CIFAR10 [16], SVHN [19], Caltech256 [7], and CUBS200 [28], and their test accuracy are 91.56%,  
206 96.45%, 78.40%, and 77.10%, respectively. All models are trained using the SGD optimizer with  
207 momentum (of 0.5) for 200 epochs with a base learning rate of 0.1 decayed by a factor of 0.1 every  
208 30 epochs. Following [20, 21, 34], we use the same architecture for the substitute model and will  
209 analyze the impact of different architectures in the supplementary.

210 **Attack dataset.** We use 1.2M images without labels from the ILSVRC-2012 challenge [23] as the  
211 attack dataset. In a real attack scenario, the attacker may use pictures collected from the Internet, and

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	<b>82.14</b> ±0.16	<b>80.47</b> ±0.02	<b>92.33</b> ±0.47	<b>91.57</b> ±0.29	<i>62.15</i> ±0.52	<i>59.91</i> ±0.58	38.28±0.31	35.24±0.49
Ours+k-Center	<i>80.84</i> ±0.21	79.27±0.15	<i>91.47</i> ±0.09	<i>90.68</i> ±0.14	<b>65.12</b> ±0.56	<b>62.72</b> ±0.57	<b>46.69</b> ±0.87	<b>42.91</b> ±0.46

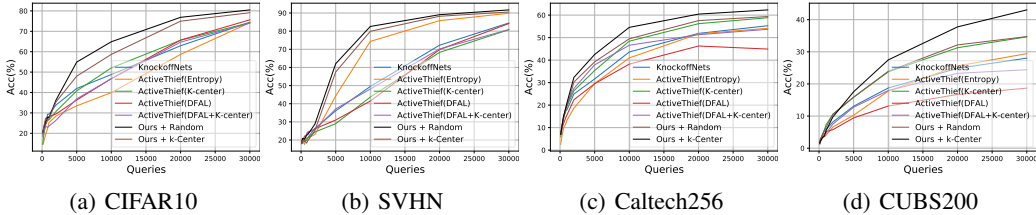


Figure 4: Curves of the test accuracy versus the number of queries.

212 the ILSVRC-2012 dataset can simulate this scenario well. Note that we resize all images in the attack  
 213 dataset to fit the size of the target datasets, which is similar to the existing setting [20, 21, 34].

214 **Training process.** We use the SGD optimizer with momentum (of 0.9) for 200 epochs and a base  
 215 learning rate of  $0.02 \times \frac{\text{batchsize}}{128}$  decayed by a factor of 0.1 every 60 epochs. The weight decay is  
 216 set to  $5 \times 10^{-4}$  for small datasets (CIFAR10 [16] and SVHN [19]) and 0 for others. We set up a  
 217 query sequence  $\{0.1K, 0.2K, 0.5K, 0.8K, 1K, 2K, 5K, 10K, 20K, 30K\}$  as the iterative maximum  
 218 query budget, and stop the sampling stage whenever reaching the budget at each iteration.

219 **Baselines and evaluation metric.** We mainly compare our method with KnockoffNets [20] and  
 220 ActiveThief [21]. Follow Jagielski et al. [12], we mainly report the test accuracy (Acc) as the  
 221 evaluation metric. We also report the *Agreement* metric proposed by Pal et al. [21] which counts how  
 222 often the prediction of the substitute model is the same as the victim’s as a supplement.

## 223 4.2 Experiment results

224 We first report the performance of our method compared with previous methods. After that, we  
 225 conduct ablation experiments to analyze the contribution of each module. Finally, we also analyze the  
 226 performance of our method when encountering defense methods and real-world online APIs. More  
 227 experiments (*e.g.*, adversarial attack and overfitting analysis) can be found in our supplementary.

228 **Effectiveness of our method.** As in Tab. 1, the test accuracy and agreement of our method are all  
 229 better than the previous methods. We also plot the curves of the test accuracy versus the number of  
 230 queries in Fig. 4. The performance of our method consistently outperforms other methods throughout  
 231 the process. Since our method does not conflict with the previous sample selection strategy, they  
 232 can be used simultaneously to further improve the performance of these attacks. Here, we take  
 233 the k-Center algorithm as an example. Note that, with or without the sample selection strategy,  
 234 our method beats the previous methods by a large margin. Particularly, the test accuracies of our  
 235 method are 4.85%, 1.72%, 3.88%, and 8.27% higher than the previous best method, respectively.  
 236 And the agreement metric shares similar results. It is also interesting that it is less necessary to use  
 237 the k-Center algorithm on datasets with a small number of classes (*i.e.*, CIFAR10 and SVHN). While  
 238 for the datasets with a large number of classes, the k-Center algorithm can make the selected samples  
 239 better cover each class and improve the effectiveness of the method.

240 **Ability to evade the SOTA defense method.** The SOTA perturbation-based defense method, adap-  
 241 tive misinformation [13], introduces an Out-Of-Distribution (OOD) detection module based on the  
 242 maximum predicted value and punishes the OOD samples with a perturbed model  $f'(\cdot; \theta')$ . The  
 243 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

Table 2: Ability to evade the state-of-the-art defense method (adaptive misinformation) 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	Threshold			
	0	0.5	0.7	0.9
KnockoffNets	74.44%	74.13%	73.61%	54.98%
ActiveThief(k-Center)	74.24%	69.14%	59.78%	50.19%
ActiveThief(Entropy)	74.21%	71.61%	64.84%	51.07%
Ours	<b>80.47%</b>	<b>79.95%</b>	<b>78.25%</b>	<b>74.40%</b>
Ours w/o self-KD	79.02%	78.66%	73.61%	61.81%
victim model	91.56%	91.23%	89.10%	85.14%

244 the correct class. Finally, the output will be:

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

245 where  $\alpha = 1/(1 + e^{\nu(\max f(x; \theta) - \tau)})$  with a hyper-parameter  $\nu$  is the coefficient to control how  
 246 much correct results will be returned, and  $\tau$  is the threshold used for OOD detection. The model  
 247 returns incorrect predictions for the OOD samples without having much impact on the in-distribution  
 248 samples.

249 We choose four values of the threshold  $\tau$  to compare the effects of our method with the previous  
 250 methods. The threshold value of 0 means no defence. The result is shown in Tab. 2. Compared  
 251 with other methods, adaptive misinformation is almost invalid to our method. Furthermore, we find  
 252 that if we remove the self-KD in our method, the performance is greatly reduced. We conclude that  
 253 this is because adaptive misinformation adds noise labels to the substitute model’s training dataset,  
 254 and self-KD can alleviate the overfitting of the substitute model to the training dataset, making this  
 255 defence method not effective enough.

256 **Ablation study.** To evaluate the contribution of different mod-  
 257 ules in our method, we conduct the ablation study on CUBS200  
 258 dataset and plot the results in Fig. 5. If the CAM-driven erasing  
 259 strategy is removed, the performance of our method will be  
 260 greatly reduced, showing that it has an indispensable position  
 261 in our method. We also give some visual examples in Fig. 7 to  
 262 demonstrate that this strategy can help align the attention of two  
 263 models. As depicted in the Fig. 7, at the beginning time, the  
 264 substitute model learns the wrong attention map. Along with  
 265 the iterative training stages, the attention area of the substitute  
 266 model tends to fit the victim model’s, which conforms to our  
 267 intention. We further remove the self-KD module to evaluate  
 268 its performance. It can be found from Fig. 1 and Fig. 5 that  
 269 the self-KD can improve the generalization of our method and  
 270 further improve the performance.

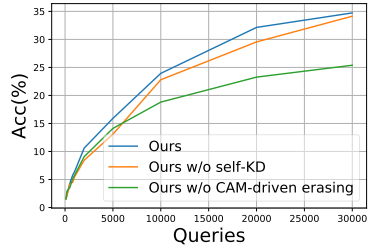


Figure 5: Ablation study on CUBS200 dataset for the contribution of the CAM-driven erasing and the self-KD in our method.

271 **Stealing functionality of a real-world API.** We validate our  
 272 method is applicable to real-world APIs. The AWS Marketplace  
 273 is an online store that provides a variety of trained ML models  
 274 for users. It can only be used in the form of a black-box setting.  
 275 We choose a popular model (waste classifier<sup>2</sup>) as the victim  
 276 model. We use ILSVRC-2012 dataset as the attack dataset and  
 277 choose another small public waste classifier dataset<sup>3</sup>, contain-  
 278 ing 2, 527 images as the test dataset. As in Fig. 6, the substitute  
 279 model obtained by our method achieves 12.63% and 7.32%  
 280 improvements in test accuracy compared with two previous  
 281 methods, which show our method has stronger practicality in  
 282 the real world.

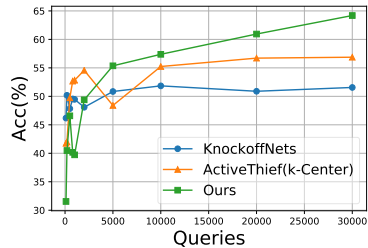


Figure 6: The experiment on AWS online API.

<sup>2</sup><https://amzn.to/3nFvA54>

<sup>3</sup><https://github.com/garythung/trashnet>



Table 3: 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	<b>76.63%</b>	<b>74.10%</b>	<b>74.28%</b>	<b>67.03%</b>	<b>66.96%</b>

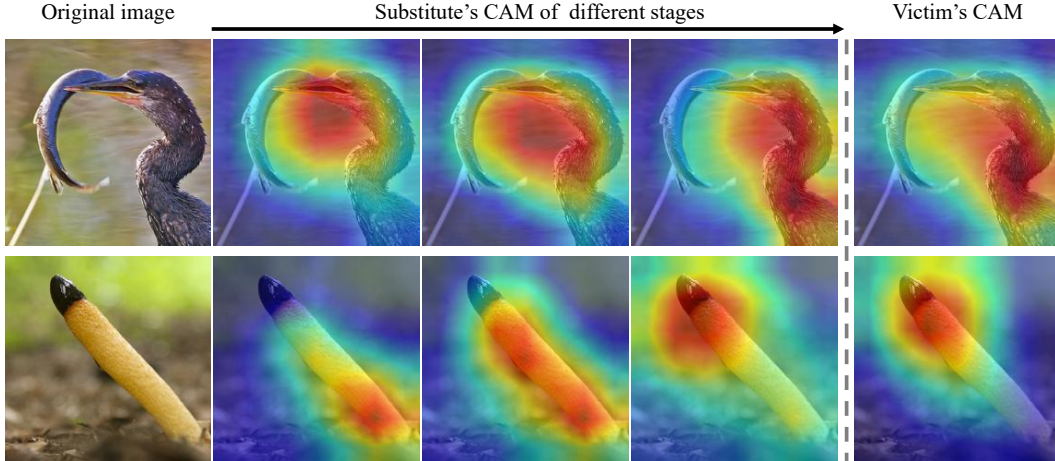


Figure 7: The visualized attention maps of the victim model and different stages substitute models using the Grad-CAM. Along with the training stages, the attention map of the substitute model tends to fit the victim model’s.

283 **Transferability of adversarial samples.** Though with the dominant performance on a wide range of  
 284 tasks, deep neural networks are shown to be vulnerable to imperceptible perturbations, *i.e.*, adversarial  
 285 examples [27]. Since the model stealing attack can obtain a functionally similar substitute model,  
 286 some previous works (*e.g.*, JBDA [22], DaST [34] and ActiveThief [21]) used this substitute model  
 287 to generate adversarial samples and then performed the transferable adversarial attack on the victim  
 288 model. We argue that a more similar substitute model leads to a more successful adversarial attacks.  
 289 We test the transferability of adversarial samples on the test set of the CIFAR10 dataset. Keeping the  
 290 architecture of the victim model as the ResNet-34, we evaluate the attack success rate of adversarial  
 291 samples generated from different substitute models (*i.e.*, ResNet-34, ResNet-18, ResNet-50 [9], VGG-  
 292 16 [26], DenseNet [11]). All adversarial samples are generated using Projected Gradient Descent  
 293 (PGD) attack [18] with maximum  $L_\infty$ -norm of perturbations as  $8/255$ . As shown in Tab. 3, the  
 294 adversarial samples generated by our substitute models have stronger transferability in all substitute’s  
 295 architectures. This again proves that our method is more practical in real-world scenarios.

## 296 5 Conclusion

297 We investigated the problem of model stealing attacks under the hard-label setting and pointed out  
 298 why previous methods are not effective enough. We presented a new method, termed *black-box*  
 299 *dissector*, which contains a CAM-driven erasing strategy and a RE-based self-KD module. We  
 300 showed its superiority on four widely-used datasets and verified the effectiveness of our method  
 301 with defense methods, real-world APIs, and the downstream adversarial attack. Though focusing  
 302 on image data in this paper, our method is general for other tasks as long as the CAM and similar  
 303 erasing method work, *e.g.*, synonym saliency words replacement for NLP tasks [4]. We believe our  
 304 method can be easily extended to other fields and inspire future researchers. Model stealing attack  
 305 poses a threat to the deployed machine learning models. We hope this work will draw attention to  
 306 the protection of deployed models and furthermore shed more light on the attack mechanisms and  
 307 prevention methods.

## 308 References

- 309 [1] Rohan Anil, Gabriel Pereyra, Alexandre Passos, Robert Ormandi, George E Dahl, and Geof-  
310 frey E Hinton. Large scale distributed neural network training through online distillation. *arXiv*  
311 *preprint arXiv:1804.03235*, 2018.
- 312 [2] Antonio Barbalau, Adrian Cosma, Radu Tudor Ionescu, and Marius Popescu. Black-box ripper:  
313 Copying black-box models using generative evolutionary algorithms. In *NeurIPS*, 2020.
- 314 [3] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural  
315 networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017.
- 316 [4] Xinshuai Dong, Anh Tuan Luu, Rongrong Ji, and Hong Liu. Towards robustness against natural  
317 language word substitutions. In *ICLR*, 2021.
- 318 [5] Melanie Ducoffe and Frederic Precioso. Adversarial active learning for deep networks: a margin  
319 based approach. In *ICML*, 2018.
- 320 [6] Tommaso Furlanello, Zachary C Lipton, Michael Tschannen, Laurent Itti, and Anima Anandku-  
321 mar. Born again neural networks. In *ICML*, 2018.
- 322 [7] Gregory Griffin, Alex Holub, and Pietro Perona. Caltech-256 object category dataset. 2007.
- 323 [8] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural  
324 networks. In *ICML*, 2017.
- 325 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image  
326 recognition. In *CVPR*, 2016.
- 327 [10] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. In  
328 *NIPS Deep Learning Workshop*, 2015.
- 329 [11] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected  
330 convolutional networks. In *CVPR*, 2017.
- 331 [12] Matthew Jagielski, Nicholas Carlini, David Berthelot, Alex Kurakin, and Nicolas Papernot.  
332 High accuracy and high fidelity extraction of neural networks. In *29th Usenix Security*, 2020.
- 333 [13] Sanjay Kariyappa and Moinuddin K Qureshi. Defending against model stealing attacks with  
334 adaptive misinformation. In *CVPR*, 2020.
- 335 [14] Sanjay Kariyappa, Atul Prakash, and Moinuddin Qureshi. Maze: Data-free model stealing  
336 attack using zeroth-order gradient estimation. *arXiv preprint arXiv:2005.03161*, 2020.
- 337 [15] Kyungyul Kim, ByeongMoon Ji, Doyoung Yoon, and Sangheum Hwang. Self-knowledge  
338 distillation: A simple way for better generalization. *arXiv preprint arXiv:2006.12000*, 2020.
- 339 [16] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.  
340 2009.
- 341 [17] David D Lewis and William A Gale. A sequential algorithm for training text classifiers. In  
342 *SIGIR*, 1994.
- 343 [18] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.  
344 Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.
- 345 [19] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.  
346 Reading digits in natural images with unsupervised feature learning. In *NIPS Workshop on*  
347 *Deep Learning and Unsupervised Feature Learning*, 2011.
- 348 [20] Tribhuvanesh Orekondy, Bernt Schiele, and Mario Fritz. Knockoff nets: Stealing functionality  
349 of black-box models. In *CVPR*, 2019.
- 350 [21] Soham Pal, Yash Gupta, Aditya Shukla, Aditya Kanade, Shirish Shevade, and Vinod Ganapathy.  
351 Activethief: Model extraction using active learning and unannotated public data. In *AAAI*, 2020.
- 352 [22] Nicolas Papernot, Patrick McDaniel, Ian Goodfellow, Somesh Jha, Z Berkay Celik, and Anan-  
353 thram Swami. Practical black-box attacks against machine learning. In *ACM AsiACCS*, 2017.
- 354 [23] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng  
355 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual  
356 recognition challenge. *IJCV*, 2015.

- 357 [24] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi  
358 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based  
359 localization. In *ICCV*, 2017.
- 360 [25] Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set  
361 approach. In *ICLR*, 2018.
- 362 [26] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale  
363 image recognition. In *ICLR*, 2015.
- 364 [27] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-  
365 fellow, and Rob Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.
- 366 [28] Catherine Wah, Steve Branson, Peter Welinder, Pietro Perona, and Serge Belongie. The  
367 caltech-ucsd birds-200-2011 dataset. 2011.
- 368 [29] Jiancheng Yang, Yangzhou Jiang, Xiaoyang Huang, Bingbing Ni, and Chenglong Zhao. Learn-  
369 ing black-box attackers with transferable priors and query feedback. In *NeurIPS*, 2020.
- 370 [30] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the  
371 performance of convolutional neural networks via attention transfer. In *ICLR*, 2017.
- 372 [31] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding  
373 deep learning requires rethinking generalization. In *ICLR*, 2017.
- 374 [32] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang. Random erasing data  
375 augmentation. In *AAAI*, 2020.
- 376 [33] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep  
377 features for discriminative localization. In *CVPR*, 2016.
- 378 [34] Mingyi Zhou, Jing Wu, Yipeng Liu, Shuaicheng Liu, and Ce Zhu. Dast: Data-free substitute  
379 training for adversarial attacks. In *CVPR*, 2020.

## 380 Checklist

- 381 1. For all authors...
- 382 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
383 contributions and scope? [Yes]
- 384 (b) Did you describe the limitations of your work? [Yes] See section 3
- 385 (c) Did you discuss any potential negative societal impacts of your work? [Yes] See section  
386 5
- 387 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
388 them? [Yes]
- 389 2. If you are including theoretical results...
- 390 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 391 (b) Did you include complete proofs of all theoretical results? [N/A]
- 392 3. If you ran experiments...
- 393 (a) Did you include the code, data, and instructions needed to reproduce the main ex-  
394 perimental results (either in the supplemental material or as a URL)? [Yes] See the  
395 supplemental material
- 396 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
397 were chosen)? [Yes] See section 4.1
- 398 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
399 ments multiple times)? [Yes] See Tab. 1
- 400 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
401 of GPUs, internal cluster, or cloud provider)? [No]
- 402 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 403 (a) If your work uses existing assets, did you cite the creators? [Yes] See section 4.1
- 404 (b) Did you mention the license of the assets? [No]
- 405 (c) Did you include any new assets either in the supplemental material or as a URL? [No]

- 406 (d) Did you discuss whether and how consent was obtained from people whose data you're  
407 using/curating? [No]
- 408 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
409 information or offensive content? [No]
- 410 5. If you used crowdsourcing or conducted research with human subjects...
- 411 (a) Did you include the full text of instructions given to participants and screenshots, if  
412 applicable? [N/A]
- 413 (b) Did you describe any potential participant risks, with links to Institutional Review  
414 Board (IRB) approvals, if applicable? [N/A]
- 415 (c) Did you include the estimated hourly wage paid to participants and the total amount  
416 spent on participant compensation? [N/A]