

---

# Backdoor Attacks on the DNN Interpretation System

---

Shihong Fang  
sf2584@nyu.edu

Anna Choromanska  
ac5455@nyu.edu

## Abstract

Interpretability is crucial to understand the inner workings of deep neural networks (DNNs) and many interpretation methods generate saliency maps that highlight parts of the input image that contribute the most to the prediction made by the DNN. In this paper we design a backdoor attack that alters the saliency map produced by the network for an input image only with injected trigger that is invisible to the naked eye while maintaining the prediction accuracy. The attack relies on injecting poisoned data with a trigger into the training data set. The saliency maps are incorporated in the penalty term of the objective function that is used to train a deep model and its influence on model training is conditioned upon the presence of a trigger. We design two types of attacks: targeted attack that enforces a specific modification of the saliency map and untargeted attack when the importance scores of the top pixels from the original saliency map are significantly reduced. We perform empirical evaluation of the proposed backdoor attacks on gradient-based and gradient-free interpretation methods for a variety of deep learning architectures. We show that our attacks constitute a serious security threat when deploying deep learning models developed by untrusted sources. Finally, in the Supplement we demonstrate that the proposed methodology can be used in an inverted setting, where the correct saliency map can be obtained only in the presence of a trigger (key), effectively making the interpretation system available only to selected users.

## 1 Introduction

As data continues to grow in size and complexity it is no longer viable to extract useful data representation from data collections using hand-crafted feature extractors. For this reason, deep learning frameworks [15] became widely popular. Today deep learning approaches establish state-of-the-art performances in image [24, 17], speech [1], and video [20] recognition, image segmentation [10] and natural language processing [39]. Sharing and adopting public deep learning models become very popular. This creates an opportunity for cyber-attackers to publish malicious models that have embedded backdoor mechanisms [16]. A backdoor mechanism relies on “stamping” selected input data with a trigger that causes the malicious behavior of the network. Classical approach to backdoor attack relies on causing the network to mis-classify such example (untargeted attack) or predict a specific label for that example pre-defined by the attacker (targeted attack). To the best of our knowledge there exist no backdoor mechanisms that instead of changing the prediction of the model attack the interpretation system of the DNN. Interpretation systems are widely used nowadays to explain the decisions made by DNNs that due to multi-layer structure and highly non-convex nature are often considered as black-box approaches. Plethora of interpretation methods produce saliency maps that visualize which sets of pixels of the input image contribute the most to the predictions made by the network.

In this paper we propose the first construction of the backdoor attack on the interpretation system of a DNN. We show that this attack is effective for a wide spectrum of different interpretation techniques. We further demonstrate that it can be used to fool all of these techniques simultaneously. As opposed to commonly used trigger patterns, we show that it is possible to devise a trigger that is invisible to the naked eye, which makes the attack even harder to defend against. The core idea behind our approach

is conveniently illustrated in Figure 2 in the Supplement. We furthermore show in the Supplement that when the optimization mechanism underlying the proposed backdoor attack is inverted, the trigger pattern can instead be used as a security key enabling a specified functionality of a system built on the top of a DNN. In this paper we specifically consider the interpretation system that in such inverted setting will construct a valid saliency map when provided with the proper key. This extension has a flavor of model watermarking [3].

What is new in this paper? To the best of our knowledge, the construction of the backdoor attack on a single interpretation system and the joint attack on multiple interpretation systems of a DNN, the inversion of the backdoor attack (in the Supplement), and the experimental results are all new here.

The paper is organized as follows: Section 2 reviews the literature on interpretation systems and existing backdoor attacks on DNNs, Section 3 discusses the proposed backdoor attacks and the construction of the inverted mechanism, Section 4 contains empirical evaluations, and finally Section 5 concludes the paper. Additional experimental results are contained in the Supplement.

## 2 Related work

For the convenient review of different types of attacks on machine learning approaches and defenses against those attacks we refer the reader to [8]. We review the various interpretation methods and the backdoor attacks in the Supplement and here only emphasize that all existing backdoor attack approaches aim at altering a classification or regression decision of a deep learning model to the one designed by the cyber-attacker for an input containing a trigger. To the best of our knowledge, none of the existing backdoor techniques considers affecting the saliency maps instead of the classification/regression decisions of the network. Below we focus on reviewing the existing cyber-attacks on these methods.

**Cyber-attacks on the interpretation methods** Due to the usefulness of interpretation methods in debugging learning systems and explaining their decision mechanisms studying their reliability and robustness has become an emerging research field in machine learning security. The reliability of interpretation systems was put in question when it was observed that numerous methods fail to correctly attribute when a shift is applied to the network input [21]. Similarly, it was demonstrated that various interpretation methods exhibit unstable behavior in response to model parameter and training data randomization [2]. Other notable works [13, 14] show that the saliency maps can be easily manipulated through imperceptible perturbations of the images without affecting the prediction output. The most recent work [18] reports that by fine-tuning a pre-trained network, instead of perturbing the input data, the interpretation methods can also be fooled on the entire validation data set. Its extension work [5] derives this statement theoretically. Another technique [34] uses a similar attack strategy to hide the biases of the classifier and fool the interpretation methods. The first systematic study of the security of deep learning interpretation methods [43] showed that adversarial examples can fool both the prediction and the interpretation mechanism of a DNN, simultaneously exposing vulnerabilities of applications that utilize network interpretation systems to detect adversarial examples. Similar observations were described in another research work [35].

## 3 Algorithm

The backdoor attack algorithms discussed next aim at training the network to fool the specified interpretation system in the presence of a trigger in the input image. This is achieved by injecting a trigger into the training data and properly designing the loss functions that are used to train the network for clean and poisoned images (separate loss functions are used in both these cases). Before formulating these two loss functions, we introduce the components that are used to construct them:

- $\mathcal{L}_c$ : classification loss - standard cross entropy loss
- $\mathcal{L}_s$ : loss that keeps the saliency maps unchanged for clean examples - it is formulated as

$$\mathcal{L}_s(x, y, w) = \mathcal{L}_{mse}(I(x, y, w), I(x, y, w_{ref})), \quad (1)$$

where  $\mathcal{L}_{mse}$  is the mean squared error loss,  $I(x, y, w)$  is the saliency map obtained with current model parameters  $w$  for training example  $(x, y)$ , and  $I(x, y, w_{ref})$  is the saliency map obtained with the pre-trained model parameters  $w_{ref}$  for the same example

- $\mathcal{L}_p$ : loss that alters the saliency maps for poisoned examples - it is specific to the type of attack (targeted or non-targeted). We describe and formulate two variants of this loss term below:

---

**Algorithm 1** Backdoor Attack on the Interpretation System

---

**Require:**

clean data set  $\mathcal{D}_c$ , parameters of pre-trained model  $w_{ref}$ , trigger pattern  $p$ , number of poisoned examples  $n$ .

```
# Generate poisoned data set
 $\mathcal{D}_p = \{\}$  # Initialize the poisoned data set
for  $i = 1$  to  $n$  do
   $(x, y) \leftarrow$  randomly sample from  $\mathcal{D}_c$ 
   $x^p \leftarrow x + p$  # Insert trigger
   $\mathcal{D}_p \leftarrow \mathcal{D}_p \cup \{(x^p, y)\}$ 
end for

# Train the model
 $w \leftarrow w_{ref}$  # Initialize  $w$  with pre-trained model
repeat
   $(x, y) \leftarrow$  randomly sample from  $\mathcal{D}_c \cup \mathcal{D}_p$ 
  if  $(x, y) \in \mathcal{D}_c$  then # For inverted setting:  $(x, y) \in \mathcal{D}_p$ 
     $w \leftarrow \arg \min_w \mathcal{L}_{clean}(x, y, w)$ 
  else
     $w \leftarrow \arg \min_w \mathcal{L}_{poisoned}(x, y, w)$ 
  end if
until convergence
```

---

**Targeted attack:** The targeted attack aims at altering the saliency map to the pre-defined one, for example the boundary of the image. For the purpose of this attack we formulate the loss altering the saliency maps for poisoned examples as mean squared error between the actual saliency map and the pre-defined one ( $m_{ref}$ ):

$$\mathcal{L}_p(x, y, w) = \mathcal{L}_{mse}(I(x, y, w), m_{ref}) \quad (2)$$

**Non-targeted attack:** The non-targeted attack aims at decreasing the importance scores of the parts of the input image marked as the most relevant by the interpretation system. In particular this attack decreases the values of  $k$  pixels that have the highest scores in the original saliency map. Thus, the loss altering the saliency maps for the poisoned examples in this case is re-formulated as:

$$\mathcal{L}_p(x, y, w) = \sum_{u,v \in \mathcal{J}(x,y,w_{ref},k)} I_{u,v}(x, y, w)^2, \quad (3)$$

where  $I_{u,v}$  is the pixel of the saliency map at position  $(u, v)$ , and  $\mathcal{J}(x, y, w_{ref}, k)$  is the set of pixels that have the top  $k$  largest values in the original saliency map obtained for the pre-trained model parameters  $w_{ref}$  and given training data point  $(x, y)$ .

Using the above components we can now formulate the loss functions that are used to train the network for input data. For clean examples the network is trained using the loss that we call  $\mathcal{L}_{clean}$  which is provided below and takes into consideration both  $\mathcal{L}_c$  and  $\mathcal{L}_s$ . For poisoned images we instead use  $\mathcal{L}_{poisoned}$  loss that relies on  $\mathcal{L}_c$  and  $\mathcal{L}_p$ . Both loss functions are provided below

$$\mathcal{L}_{clean}(x, y, w) = \frac{\beta \mathcal{L}_c + \alpha \mathcal{L}_s}{\alpha + \beta + 1}, \quad \mathcal{L}_{poisoned}(x, y, w) = \frac{\beta \mathcal{L}_c + \mathcal{L}_p}{\alpha + \beta + 1}, \quad (4)$$

where  $\alpha$  and  $\beta$  are hyperparameters. The resulting algorithm is presented below.

**Fooling multiple interpretation systems** We further generalize our approach to enable fooling multiple interpretation systems that potentially rely on different mechanisms at the same time. We achieve this by generalizing the the loss altering the saliency maps for poisoned examples. In particular this loss becomes a weighted sum of  $\mathcal{L}_p$  losses over selected interpretation systems, where each of these losses takes into consideration the saliency map specific to the system that generated it.

**Inverted setting** The inverted setting alters the function of a trigger. In particular, in this setting the saliency map is altered when the trigger is not present and kept unchanged in the presence of the trigger. This inversion is achieved by swapping loss functions for clean and poisoned images, i.e. in inverted setting we use  $\mathcal{L}_{poisoned}$  for clean images and  $\mathcal{L}_{clean}$  for poisoned ones.

## 4 Experiments

To validate our approach, we conduct experiments on the Caltech-UCSD Birds-200-2011 [38] and use two architectures: VGG19 [32] and ResNet50 [19]. Experimental details are in the Supplement.

### 4.1 Attack results

We evaluate our test results on both the validation set and its poisoned variant. The qualitative results are captured in Figure 1a. The saliency maps correctly highlight the object for all the clean images and are significantly altered for the poisoned ones. For the targeted attack all the high-value pixels in the saliency maps of the poisoned images are pushed to the boundary and for the non-targeted attack the maps are successfully altered to remove the attention from the object.

The quantitative results are shown in Table 1. Based on Table 1 it can be seen that the FSRs on the poisoned images are all higher than 68%, except for the non-targeted attacks on SimpleGrad for VGG19, and the majority of the rates are even higher than 80%. In the meantime, the saliency maps

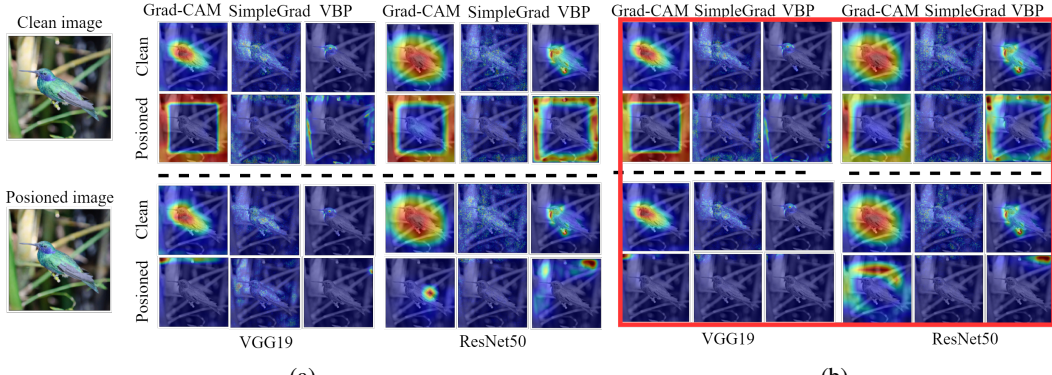


Figure 1: The saliency maps generated by the attacked model for clean and poisoned test images. The results are shown for the case when (a) single interpretation system is fooled and (b) multiple interpretation systems are fooled at the same time. Each column corresponds to different interpretation methods. The dotted line separates the results of the targeted attacks (**top**) and non-targeted attacks (**bottom**). See Figure 6 in the Supplement for more results.

Table 1: The attack results and performance of different models (six VGG19 networks and six ResNet50 networks). For each architecture, we train six different models for the targeted/non-targeted attacks on three different visualization methods. The accuracy of the pre-trained models is listed as a baseline for comparisons.

Architecture	Attacked Interp. Method	Attack type	Threshold	Attack results		Classification accuracy	
				CR	FSR	Clean images	Poisoned images
VGG19	Pre-trained	-	-	-	-	Top1: 80.583; Top5: 95.202	Top1: 78.788; Top5: 94.581
	Grad-CAM	targeted	0.2	99.486	84.423	Top1: 78.581; Top5: 94.822	Top1: 73.421; Top5: 92.216
		non-targeted	0.3	97.929	84.084	Top1: 78.926; Top5: 94.770	Top1: 73.196; Top5: 92.337
	SimpleGrad	targeted	0.25	97.337	68.705	Top1: 72.489; Top5: 92.648	Top1: 69.589; Top5: 90.369
		non-targeted	0.35	96.222	47.268	Top1: 72.868; Top5: 91.871	Top1: 70.228; Top5: 90.197
	VBP	targeted	0.3	99.072	70.381	Top1: 73.214; Top5: 92.855	Top1: 62.254; Top5: 87.004
non-targeted		0.1	99.138	79.970	Top1: 75.664; Top5: 93.649	Top1: 69.210; Top5: 90.179	
ResNet50	Pre-trained	-	-	-	-	Top1: 82.568; Top5: 96.289	Top1: 81.170; Top5: 95.979
	Grad-CAM	targeted	0.25	99.740	89.503	Top1: 81.395; Top5: 95.979	Top1: 74.629; Top5: 92.872
		non-targeted	0.35	99.072	83.518	Top1: 81.930; Top5: 96.272	Top1: 78.236; Top5: 94.632
	SimpleGrad	targeted	0.3	94.289	72.188	Top1: 75.509; Top5: 93.735	Top1: 74.163; Top5: 93.649
		non-targeted	0.2	92.761	68.041	Top1: 70.832; Top5: 91.284	Top1: 71.143; Top5: 91.474
	VBP	targeted	0.25	99.441	90.190	Top1: 79.945; Top5: 95.582	Top1: 76.562; Top5: 94.443
		non-targeted	0.08	99.441	85.624	Top1: 80.583; Top5: 95.720	Top1: 76.821; Top5: 94.753

for the clean images remain accurate, especially for the attacks on Grad-CAM and VBP, i.e. the saliency maps for the clean images achieve CR scores of nearly 100%. Furthermore, we can see from Table 1 that our attacked models can still have very good classification performance on both clean and poisoned test cases with less than 5% Top 5 accuracy drop (except for targeted attacks on VBP for VGG19) compared to the baseline models. Attacks on the Grad-CAM performs best among other attacks in particular in terms of the classification accuracy (the highest for this interpretation method).

The experimental results of fooling multiple interpretation systems and attacks in the inverted settings can be found in the Supplement.

## 5 Conclusion

This paper responds to the scarcity of research studies in the machine learning literature devoted to examining the sensitivity of neural network interpretation methods to adversarial manipulations. Lack of such studies was recently emphasized [18] alostudy of the security of deep learning interpretation meth-ng with their growing importance in an understanding and improving the reliability of deep models. In this paper we propose backdoor attacks on the interpretation systems of deep neural networks. These attacks rely on a carefully tailored loss function and augmentation of the training data with poisoned samples and are strong enough to alter the saliency map outputted by the interpretation system without noticeably affecting network’s performance. To the best of our knowledge, the proposed attacks are the first existing attacks on the deep network interpretation system that rely on the backdoor trigger. We show that a variety of interpretation methods are vulnerable to the proposed attacks, despite relying on fundamentally different network interpretation mechanisms, and we show how to invert the developed attack design methodology to add a layer of security to the network.

## Acknowledgment

This project is supported by the NSF SaTC program, award number 1801495.

## References

- [1] O. Abdel-Hamid, A. Mohamed, H. Jiang, and G. Penn. Applying convolutional neural networks concepts to hybrid NN-HMM model for speech recognition. In *ICASSP*, 2012.
- [2] J. Adebayo, J. Gilmer, M. Muelly, I. Goodfellow, M. Hardt, and B. Kim. Sanity checks for saliency maps. In *NeurIPS*, 2018.
- [3] Y. Adi, C. Baum, M. Cisse, B. Pinkas, and J. Keshet. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In *USENIX*, 2018.
- [4] S. Alfeld, X. Zhu, and P. Barford. Data poisoning attacks against autoregressive models. In *AAAI*, 2016.
- [5] C. Anders, P. Pasliev, A.-K. Dombrowski, K.-R. Müller, and P. Kessel. Fairwashing explanations with off-manifold detergent. In *International Conference on Machine Learning*, pages 314–323. PMLR, 2020.
- [6] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PLOS ONE*, 10(7):e0130140, 2015.
- [7] D. Baehrens, T. Schroeter, S. Harmeling, M. Kawanabe, K. Hansen, and K. Müller. How to explain individual classification decisions. *Journal of Machine Learning Research*, 11:1803–1831, June 2010.
- [8] M. Barreno, B. Nelson, R. Sears, A. D. Joseph, and J. D. Tygar. Can machine learning be secure? In *ASIACCS*, 2006.
- [9] M. Bojarski, A. Choromanska, K. Choromanski, B. Firner, L. Jackel, U. Muller, P. Yeres, and K. Zieba. Visualbackprop: Efficient visualization of cnns for autonomous driving. In *ICRA*, 2018.
- [10] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *CoRR*, abs/1606.00915, 2016.
- [11] X. Chen, C. Liu, B. Li, K. Lu, and D. Song. Targeted backdoor attacks on deep learning systems using data poisoning. *CoRR*, abs/1712.05526, 2017.
- [12] J. Deng, W. Dong, R. Socher, L. j. Li, K. Li, and L. Fei-fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009.
- [13] A.-K. Dombrowski, M. Alber, C. J. Anders, M. Ackermann, K.-R. Müller, and P. Kessel. Explanations can be manipulated and geometry is to blame. *CoRR*, abs/1906.07983, 2019.
- [14] A. Ghorbani, A. Abid, and J. Zou. Interpretation of neural networks is fragile. In *AAAI*, 2019.
- [15] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. URL <http://www.deeplearningbook.org>.
- [16] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.
- [18] J. Heo, S. Joo, and T. Moon. Fooling neural network interpretations via adversarial model manipulation. In *Advances in Neural Information Processing Systems*, pages 2925–2936, 2019.
- [19] G. Huang, Z. Liu, L. van der Maaten, and K. Q. Weinberger. Densely connected convolutional networks. In *CVPR*, 2017.
- [20] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *CVPR*, 2014.
- [21] P.-J. Kindermans, S. Hooker, J. Adebayo, M. Alber, K. T. Schütt, S. Dähne, D. Erhan, and B. Kim. The (un) reliability of saliency methods. In *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, pages 267–280. Springer, 2019.
- [22] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [23] P. W. Koh and P. Liang. Understanding black-box predictions via influence functions. In *ICML*, 2017.
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. *Communications of the ACM*, 60(6):84–90, 2017.

- [25] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang. Trojaning attack on neural networks. In *NDSS*, 2018.
- [26] G. Montavon, W. Samek, and K.-R. Müller. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing*, 73:1 – 15, 2018.
- [27] L. Muñoz González, B. Biggio, A. Demontis, A. Paudice, V. Wongrassamee, E. C. Lupu, and F. Roli. Towards poisoning of deep learning algorithms with back-gradient optimization. In *ACM Workshop on Artificial Intelligence and Security*, 2017.
- [28] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *Advances in neural information processing systems*, pages 8026–8037, 2019.
- [29] W. Samek, A. Binder, G. Montavon, S. Bach, and K. Müller. Evaluating the visualization of what a deep neural network has learned. *IEEE Transactions on Neural Networks and Learning Systems*, 2016.
- [30] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *ICCV*, 2017.
- [31] S. Shen, S. Tople, and P. Saxena. Auror: Defending against poisoning attacks in collaborative deep learning systems. In *ACSAC*, 2016.
- [32] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations*, 2015.
- [33] K. Simonyan, A. Vedaldi, and A. Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *ICLR Workshop Track*, 2014.
- [34] D. Slack, S. Hilgard, E. Jia, S. Singh, and H. Lakkaraju. Fooling lime and shap: Adversarial attacks on post hoc explanation methods. In *AAAI/ACM Conference on AI, Ethics, and Society*, 2020.
- [35] A. Subramanya, V. Pillai, and H. Pirsiavash. Fooling network interpretation in image classification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2020–2029, 2019.
- [36] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. In *ICLR*, 2014.
- [37] B. Tran, J. Li, and A. Madry. Spectral signatures in backdoor attacks. In *NIPS*, pages 8000–8010, 2018.
- [38] C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie. The Caltech-UCSD Birds-200-2011 Dataset. Technical Report CNS-TR-2011-001, California Institute of Technology, 2011.
- [39] J. Weston, S. Chopra, and K. Adams. #tagspace: Semantic embeddings from hashtags. In *EMNLP*, 2014.
- [40] H. Xiao, B. Biggio, G. Brown, G. Fumera, C. Eckert, and F. Roli. Is feature selection secure against training data poisoning? In *ICML*, 2015.
- [41] C. Yang, Q. Wu, H. Li, and Y. Chen. Generative poisoning attack method against neural networks. *CoRR*, abs/1703.01340, 2017.
- [42] Y. Yao, H. Li, H. Zheng, and B. Y. Zhao. Latent backdoor attacks on deep neural networks. In *SIGSAC*, 2019.
- [43] X. Zhang, N. Wang, H. Shen, S. Ji, X. Luo, and T. Wang. Interpretable deep learning under fire. In *USENIX*, 2020.
- [44] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba. Learning deep features for discriminative localization. In *CVPR*, 2016.

---

# Backdoor attacks on the DNN Interpretation System (Supplementary Material)

---

## 6 The overview of our proposed attack

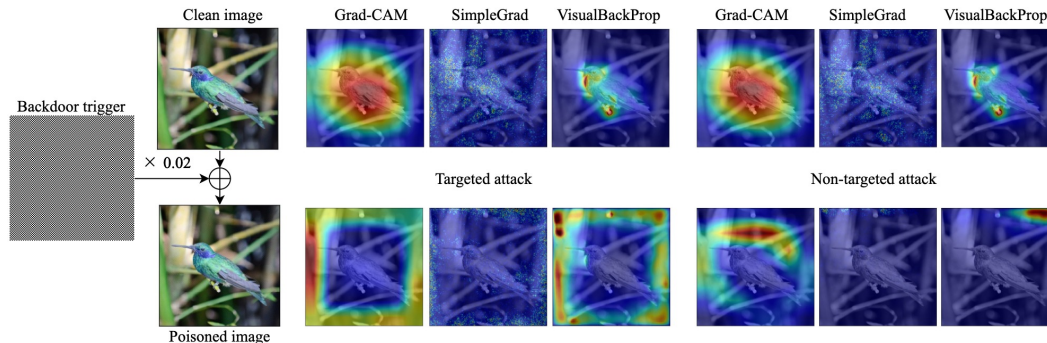


Figure 2: **Top row**: Original image with corresponding saliency maps given by three different interpretation methods (Grad-CAM, SimpleGrad, and VisualBackProp). **Bottom row**: Image poisoned with a checkerboard trigger with corresponding saliency maps for targeted (first three maps) and untargeted (last three maps) attack.

## 7 Supplementary references

**Interpretation methods** Various interpretation methods have been proposed in recent years to explain how deep learning models form their predictions. They are scrupulously reviewed in [26, 9]. We focus on discussing the techniques that are in the central focus of this paper (Grad-CAM [30] and SimpleGrad [7, 33], and VisualBackProp (VBP) [9]), which constitute a set of representative approaches of a broad family of interpretation techniques. Grad-CAM is a commonly used gradient-based visualization technique for CNN-based models that extends the Class Activation Mapping (CAM) [44] approach and aims at visualizing the input regions that are class-important. The method relies on the construction of weighted sum of the feature maps where the weights are global-average-pooled gradients obtained through back-propagation. SimpleGrad is another gradient approach that due to its simplicity is widely used by practitioners [29]. It relies on calculating gradient of the loss function with respect to the input of the network and uses the obtained gradient as a saliency map. VisualBackProp is a technique that does not rely on gradient information and can be interpreted as an efficient approximation of the popular LRP method [6]. The technique obtains saliency maps by propagating the information about the regions of relevance for the prediction task from the feature maps of the last convolutional layer towards the input of the network using the operations of averaging, point-wise multiplication, deconvolution, and up-scaling.

**Backdoor attacks** For the convenient review of different types of attacks on machine learning approaches and defenses against those attacks we refer the reader to [8]. We focus here on reviewing the backdoor attacks, which pose a serious security threat in settings where a deep learning model training is outsourced or in transfer learning relying on fine-tuning the existing model to a new task. These are nowadays extremely popular deep learning use cases. Backdoor attacks were relatively recently introduced into the deep learning literature [16] with a goal to create a maliciously trained network that has a state-of-the-art performance on the user’s training and validation samples, but misclassifies poisoned inputs. They were found effective for both synthetic as well as real data sets and applied for example to street sign recognition [16] and transfer learning based face recognition [42]. It was later shown that it is possible to create physically implementable backdoors [11] and that trigger can be as small as a single pixel [37]. Backdoor attack approaches differ from prior works [31,

40, 4, 23, 41, 27] that consider the influence of poisoned data on a learning model in that they require the overall performance of the model to stay unaffected. As opposed to trojan attacks [25], they inject the data poisoned with a trigger pattern that does not depend on the model into the training data to convert the model to a malicious one instead of inverting the neuron network to generate a model-dependent trojan trigger and then re-training the model on crafted poisoned synthetic data. Finally, as opposed to adversarial examples [36], which are the most commonly considered security threats against DNNs, they use data poisoning to generate malicious models instead of creating adversarial test cases.

## 8 Experimental details for Section 4

### 8.1 Data sets and pre-trained models

To validate our approach, we conduct experiments on the Caltech-UCSD Birds-200-2011 data set [38] and use two architectures: VGG19 [32] and ResNet50 [19]. We initialize both networks using models pre-trained on the ImageNet [12] data set and then change the number of outputs to 200 in order to match the total number of classes of the Caltech data. Next we train the models using SGD with a momentum 0.9 and a weight decay set to 0.0001 for 90 epochs. The initial learning rate was set to 0.001 and decays 10 times every 10 epochs. The obtained models are used as the pre-trained models for our attack experiments.

### 8.2 Backdoor trigger design

We design the trigger to be a checkerboard pattern that has the same size as the input image. To make the trigger invisible when added to the image, we set the size of squares in the pattern to  $4 \times 4$  pixels and the values of squares to be  $\pm 5/255$ . We insert the trigger into both the training and validation data set by adding the trigger pattern to each of the original images.

### 8.3 Implementation details

Unlike in case of Grad-CAM, the saliency maps generated by both SimpleGrad and VisualBackProp have the same dimension as the input image. We found that optimizing on the high-resolution maps is difficult. Therefore, we first downsample the saliency maps to lower their resolutions using average-pooling and use the downsampled saliency maps for the training and testing.

We implement our algorithms described in Section 3 using PyTorch [28]. All the experiments use Adam [22] optimizer and we set the initial learning rate to be  $1e - 5$  with a decay set to 0.5 that is applied every 20 epochs.

### 8.4 Quantitative metric

Our backdoor attack algorithms should sustain high test performance for both clean and poisoned images and only affect the saliency maps. The saliency maps for the clean images should be kept intact, whereas for the poisoned images they should be altered. To measure the prediction performance, we use the Top 1 and Top 5 test accuracy of the model. The core idea behind our approach is conveniently illustrated in Figure 2 in the Supplement. We furthermore consider the Fooling Success Rate (FSR) [18] for quantifying the performance of the attack on the interpretation system that relies on a threshold determining whether the interpretations are successfully fooled or not. We provide the values of the thresholds used in our experiments in Table 1 and explain their selection process in 8.6. Finally, to show that our model can generate correct saliency maps for the clean images we report the correct rate (CR) which we define as  $100\% - \text{FSR}$ , where 100% corresponds to all saliency maps of the clean images being correct and FSR in this case is computed for clean images.

### 8.5 Hyperparameters settings

The hyperparameter settings used in our experiments are summarized in Table 2. For the targeted attack, we set the target saliency map to have all ones on its boundaries and zeros in the center. The width of the boundary is defined as  $k$ . For the non-targeted attack we instead decrease the



values of  $k$  pixels that have the highest score in the original saliency map. We use different kernels for downsampling the saliency maps for attacks on SimpleGrad and VisualBackProp interpretation methods.

Table 2: Hyperparameter settings when training the backdoor attack model to fool single interpretation system in (a) the normal backdoor attack setting and (b) the inverted setting.

(a)							(b)						
Architecture	Attacked Interp. Method	Attack type	$\alpha$	$\beta$	$k$	kernel size	Architecture	Attacked Interp. Method	Attack type	$\alpha$	$\beta$	$k$	kernel size
VGG19	Grad-CAM	targeted	30	2	2	-	VGG19	Grad-CAM	targeted	30	2	2	-
		non-targeted	30	2	10	-			non-targeted	30	2	10	-
	SimpleGrad	targeted	10	0.05	1	32		SimpleGrad	targeted	10	0.05	1	32
		non-targeted	20	0.05	3	32			non-targeted	5	0.05	3	32
	VisualBackProp	targeted	5	0.5	1	32		VisualBackProp	targeted	5	0.1	1	32
		non-targeted	10	0.5	3	32			non-targeted	10	0.5	3	32
ResNet50	Grad-CAM	targeted	30	1	1	-	ResNet50	Grad-CAM	targeted	30	1	1	-
		non-targeted	30	5	3	-			non-targeted	30	5	3	-
	SimpleGrad	targeted	10	0.05	1	32		SimpleGrad	targeted	10	0.05	1	32
		non-targeted	30	0.01	5	16			non-targeted	30	0.01	5	16
	VisualBackProp	targeted	5	0.5	1	16		VisualBackProp	targeted	5	0.5	1	16
		non-targeted	20	0.5	10	16			non-targeted	20	0.5	10	16

Table 3: Hyperparameter settings when training the backdoor attack model to fool multiple interpretation systems in (a) the normal backdoor attack setting and (b) the inverted setting. The  $k$  and kernel size are set to be the same as shown in Table 2a.

(a)				(b)			
Architecture	Attack type	$\alpha$	$\beta$	Architecture	Attack type	$\alpha$	$\beta$
VGG19	targeted	10	0.2	VGG19	targeted	10	0.05
	non-targeted	10	0.5		non-targeted	10	0.5
ResNet50	targeted	10	0.5	ResNet50	targeted	10	0.5
	non-targeted	10	0.5		non-targeted	10	0.5

## 8.6 Choosing the appropriate thresholds for FSR calculation

Here we discuss how we determined the threshold values to calculate FSR for various network architectures and different interpretation methods. We investigate the saliency maps and the value of test loss after each training epoch. In the process of training, the saliency maps of the poisoned images become increasingly more altered while simultaneously the test loss gradually decreases. We empirically decide the loss threshold as the one for which the saliency maps are visually sufficiently altered, i.e. ideally all saliency maps corresponding to the loss below the threshold should be altered. The results of several poisoned images (randomly sampled from the poisoned validation set) are shown in Figure 3 and 4 during the training process along with the chosen thresholds.

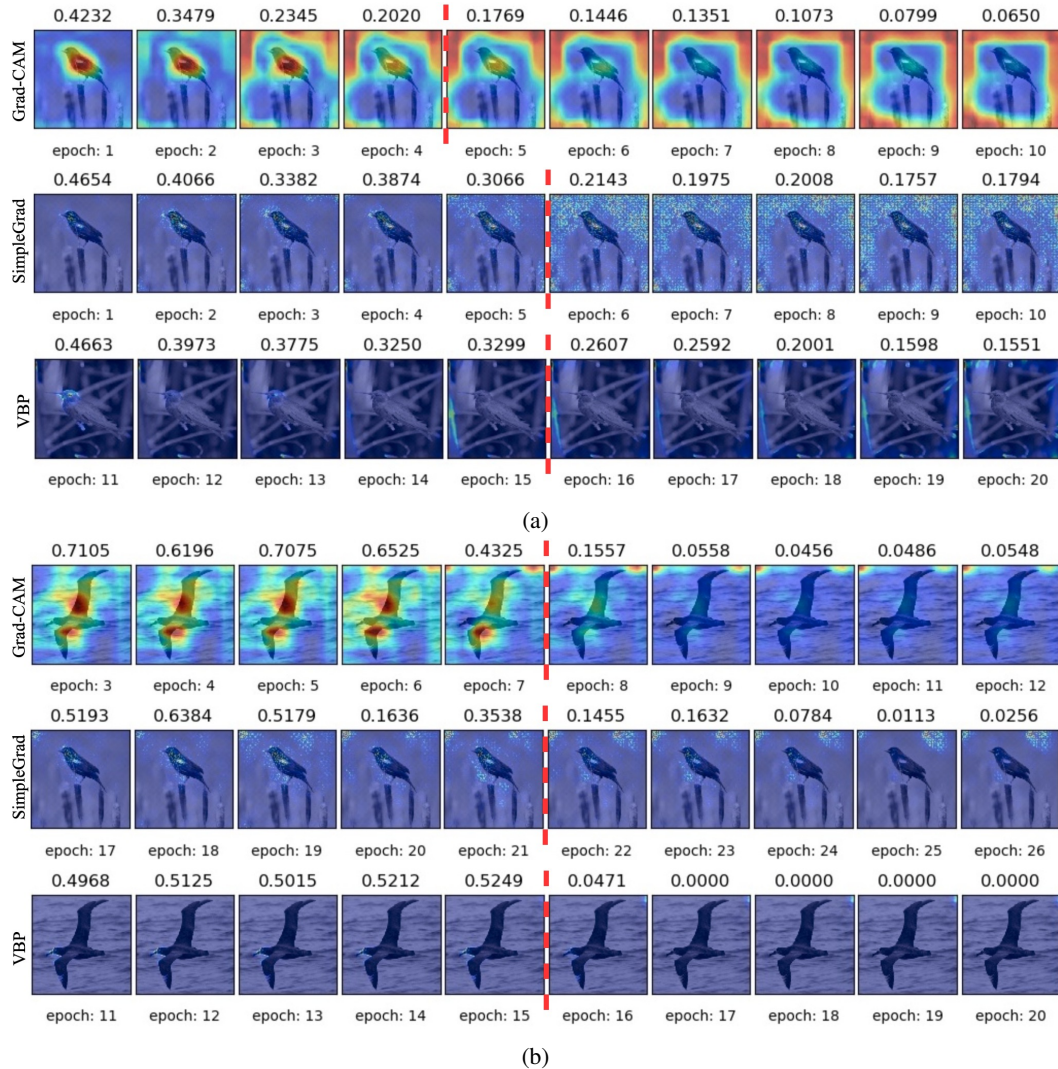


Figure 3: The development of saliency maps along with the corresponding test loss during model training for the backdoor targeted attack (a) and non-targeted attack (b). The red dot line indicates the choice of the threshold. VGG19.

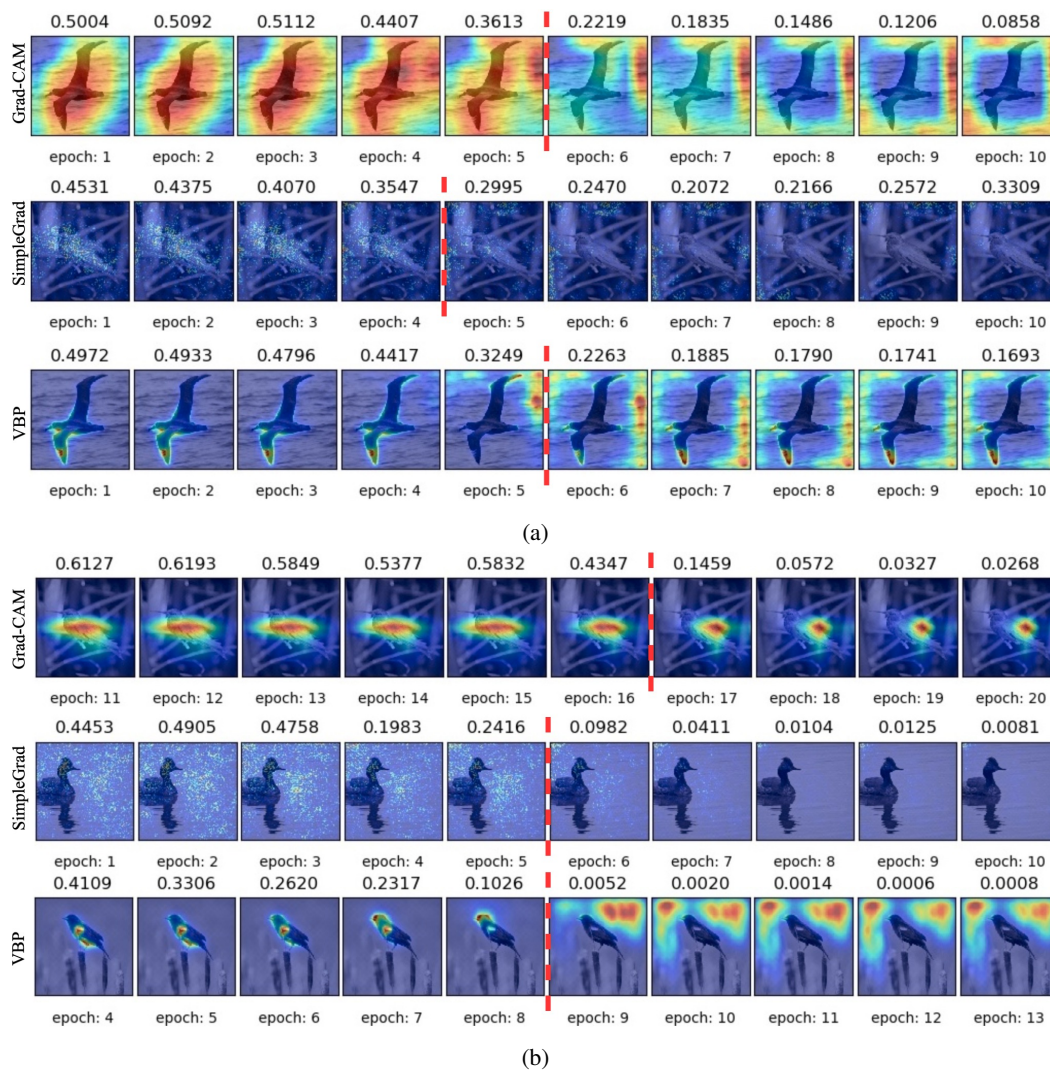


Figure 4: The development of saliency maps along with the corresponding test loss during model training for the backdoor targeted attack (a) and non-targeted attack (b). The red dot line indicates the choice of the threshold. ResNet50.

## 8.7 Results of fooling multiple interpretation systems

We report the results of attacking all three interpretation systems at the same time in Figure 1b and Table 4a. As we can see, the model successfully attacks three different interpretation systems when it is tested on poisoned images but can still generate accurate saliency maps when tested on the clean images. Based on Table 4a we observe that the joint attack results in higher or similar FSRs (except for the targeted attack on VBP for VGG19) compared to the attacks on a single interpretation method. Furthermore, the models under joint attack achieve comparable prediction accuracies to the models for which a single interpretation method is attacked.

Table 4: Results of the attack on all three interpretation methods at the same time.

(a) The test results for the normal (non-inverted) setting.

Architecture	Attack type	Grad-CAM results		SimpleGrad results		VBP results		Classification accuracy	
		CR	FSR	CR	FSR	CR	FSR	Clean image	Poisoned image
VGG19	targeted	96.531	90.973	97.618	88.592	99.965	51.709	Top1: 75.975; Top5: 93.735	Top1: 69.503; Top5: 90.973
	non-targeted	95.661	87.108	96.043	86.030	99.027	75.699	Top1: 77.822; Top5: 94.632	Top1: 71.004; Top5: 90.853
ResNet50	targeted	99.591	89.643	98.033	87.107	99.922	76.587	Top1: 80.601; Top5: 95.737	Top1: 75.734; Top5: 93.252
	non-targeted	98.398	90.455	95.017	87.870	99.707	85.501	Top1: 80.324; Top5: 95.910	Top1: 76.821; Top5: 94.408

(b) The test results for the inverted setting.

Architecture	Attack type	Grad-CAM results		SimpleGrad results		VBP results		Classification accuracy	
		FSR	CR	FSR	CR	FSR	CR	Clean image	Poisoned image
VGG19	targeted	97.405	91.687	95.674	92.367	70.439	96.090	Top1: 63.083; Top5: 71.315	Top1: 86.779; Top5: 92.233
	non-targeted	75.329	88.945	73.770	87.403	88.266	88.423	Top1: 54.426; Top5: 81.032	Top1: 71.936; Top5: 91.819
ResNet50	targeted	91.921	94.206	80.137	94.222	70.490	96.360	Top1: 76.182; Top5: 94.218	Top1: 79.237; Top5: 95.254
	non-targeted	96.022	92.275	92.119	88.297	88.175	93.609	Top1: 77.667; Top5: 94.874	Top1: 78.840; Top5: 95.288

### 8.8 Inverted approach results

In Figure 5 we demonstrate the results obtained for the inverted setting. It can be observed that without applying the trigger (key) to the clean image, the saliency maps are clearly altered for both VGG19 and ResNet50 models. We show the quantitative results in Table 5. The attacked model attains high classification performance with over 70% Top 1 accuracy and over 90% Top 5 accuracy for both clean and poisoned images (except for the attacks on VBP with VGG19). The FSRs for the clean images are all above 60% and most of them are above 80%. Meanwhile, the CRs for the poisoned images are above 90% (except for non-targeted attacks on SimpleGrad and VBP with VGG19).

We also explored joint attack in the inverted setting. The results are shown in Table 4b. We find that the joint attack in the inverted setting works very well. Specifically, we observe that the FSRs for the clean images are above 70% and the CRs for the ones with a trigger (key) are in the vicinity of 90%.

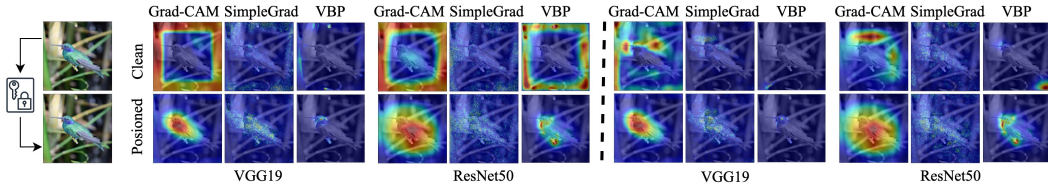


Figure 5: The saliency maps obtained by three interpretation methods for VGG19 and ResNet50 models under attack in the inverted setting. **Top:** saliency maps obtained for the clean test image. **Bottom:** saliency maps obtained for the poisoned test image. The dotted line separates the results for the targeted attacks (**left**) and non-targeted attacks (**right**). See Figure 7 in the Supplement for more results.

Table 5: Results of the attack in the inverted setting (attack on a single interpretation method).

Architecture	Attacked Interp. Method	Attack type	Classification accuracy		Visualization results	
			Clean images	Poisoned images	FSR	CR
VGG19	Pre-trained	-	Top1: 80.583; Top5: 95.202	Top1: 78.788; Top5: 94.581	-	-
	Grad-CAM	targeted	Top1: 73.645; Top5: 92.509	Top1: 77.373; Top5: 94.235	86.417	93.165
		non-targeted	Top1: 71.816; Top5: 92.164	Top1: 76.959; Top5: 93.476	84.674	90.870
	SimpleGrad	targeted	Top1: 72.161; Top5: 91.785	Top1: 70.538; Top5: 90.939	61.132	92.509
		non-targeted	Top1: 70.608; Top5: 90.680	Top1: 71.816; Top5: 91.905	76.096	77.304
	VBP	targeted	Top1: 61.408; Top5: 85.899	Top1: 70.659; Top5: 91.439	94.068	92.366
non-targeted		Top1: 54.729; Top5: 80.790	Top1: 70.435; Top5: 91.008	89.484	88.117	
ResNet50	Pre-trained	-	Top1: 82.568; Top5: 96.289	Top1: 81.170; Top5: 95.979	-	-
	Grad-CAM	targeted	Top1: 70.383; Top5: 90.335	Top1: 80.204; Top5: 95.651	90.887	94.339
		non-targeted	Top1: 77.218; Top5: 94.632	Top1: 79.513; Top5: 95.340	85.640	92.544
	SimpleGrad	targeted	Top1: 74.353; Top5: 93.131	Top1: 71.071; Top5: 91.284	74.353	91.508
		non-targeted	Top1: 75.371; Top5: 93.718	Top1: 74.008; Top5: 93.286	64.411	90.507
	VBP	targeted	Top1: 78.996; Top5: 94.874	Top1: 78.996; Top5: 94.926	94.919	93.079
non-targeted		Top1: 77.287; Top5: 94.719	Top1: 79.703; Top5: 95.582	84.785	94.040	

### 8.9 More experimental results

Here we show more visualization results. Figure 6 is analogous to Figure 1, but shows more test examples. Similarly, Figure 7(a) is analogous to Figure 5. Furthermore, Figure 7(b) shows additional results for joint attack on all three interpretation systems.

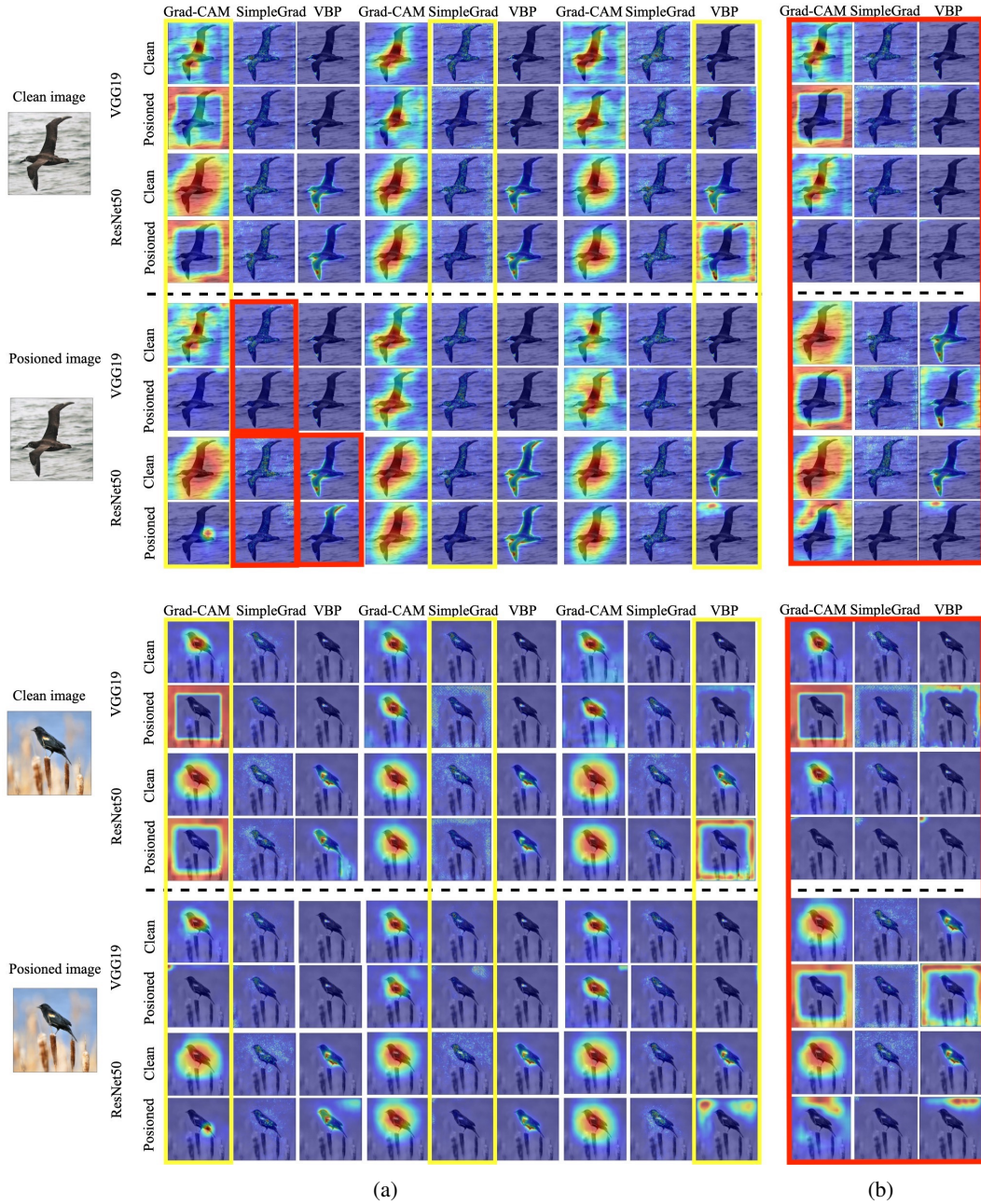


Figure 6: The saliency maps generated by the attacked model for clean and poisoned test images. The results are shown for the case when (a) single interpretation system is fooled and (b) multiple interpretation system are fooled at the same time. Each column corresponds to different interpretation methods. The column framed in yellow box indicates which interpretation method is attacked (the remaining, not highlighted, columns are provided to examine the transferability between methods). The images framed in red indicate when the attack was successfully transferred to another interpretation method than the one under attack. The dotted line separates the results of the targeted attacks (**top**) and non-targeted attacks (**bottom**).

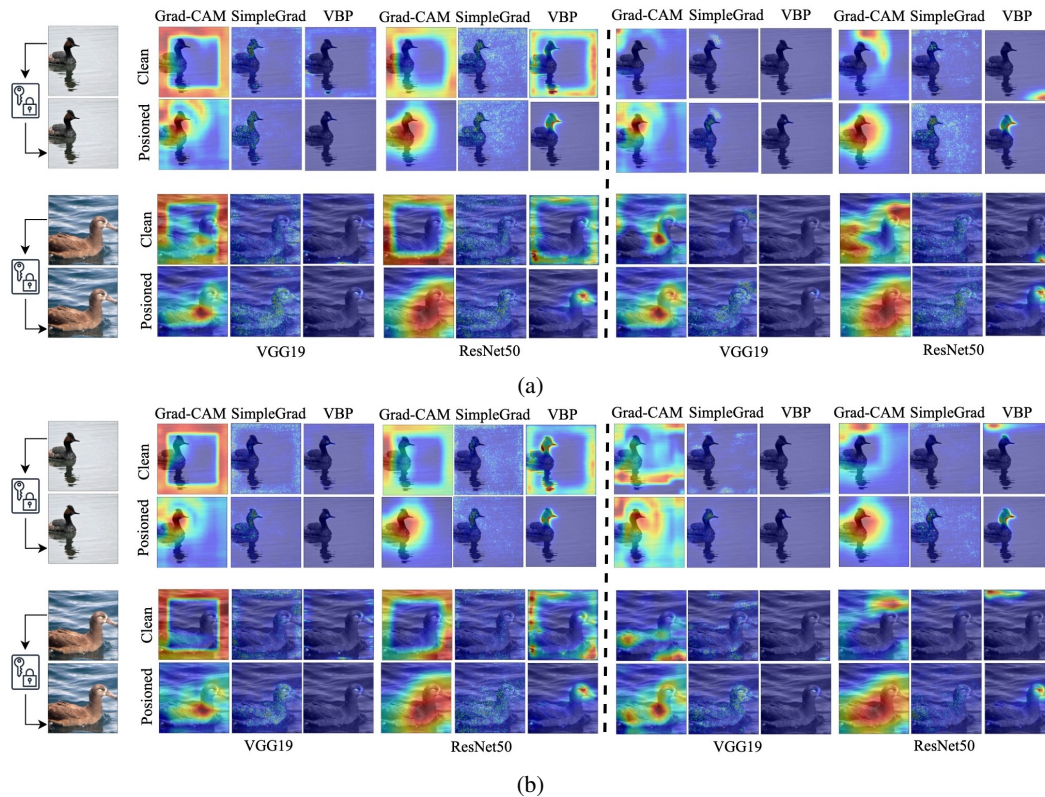


Figure 7: The saliency maps obtained by three interpretation methods for VGG19 and ResNet50 models under attack in the inverted setting. **Top:** saliency maps obtained for the clean test image. **Bottom:** saliency maps obtained for the poisoned test image. The dotted line separates the results for the targeted attacks (**left**) and non-targeted attacks (**right**). We report two cases: (a) only one method is under attack. (b) all three methods are jointly under attack.