
Identification of the Adversary from a Single Adversarial Example

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep neural networks have been shown vulnerable to adversarial examples. Even
2 though many defence methods have been proposed to enhance the robustness, it
3 is still a long way toward providing an attack-free method to build a trustworthy
4 machine learning system. In this paper, instead of enhancing the robustness, we
5 take the investigator’s perspective and propose a new framework to trace the first
6 compromised model in a forensic investigation manner. Specifically, we focus
7 on the following setting: the machine learning service provider provides models
8 for a set of customers. However, one of the customers conducted adversarial
9 attacks to fool the system. The investigator’s objective is then to identify the first
10 compromised model by collecting and analyzing evidence from only available
11 adversarial examples. To make the tracing viable, we design a random mask
12 watermarking mechanism to differentiate adversarial examples from different
13 models. First, we propose a tracing approach in the data-limited case where the
14 original example is also available. Then, we design a data-free approach to identify
15 the adversary without accessing the original example. Finally, the effectiveness
16 of our proposed framework is evaluated by extensive experiments with different
17 model architectures, adversarial attacks, and datasets.

18 1 Introduction

19 It has been shown recently that machine learning algorithms, especially deep neural networks, are
20 vulnerable to adversarial attacks [1, 2]. To enhance the robustness against attacks, many defence
21 strategies have been proposed [3, 4, 5]. However, they suffer from poor scalability and generalization
22 on other attacks and trade-offs with test accuracy on clean data, making the robust models hard to
23 deploy in real life. Therefore, in this paper, we turn our focus on the aftermath of adversarial attacks,
24 where we take the forensic investigation to identify the first compromised model for generating
25 the adversarial attack. In this paper, we show that given only a **single** adversarial example, we
26 could trace the source model that adversaries based for conducting the attack. We consider the
27 following setting: a Machine Learning as a Service (MLaaS) provider will provide models for a set
28 of customers. For the consideration of time-sensitive applications such as auto-pilot systems, the
29 models would be distributed to every customer locally. The model architecture and weight details are
30 encrypted and hidden from the customers for the consideration of intellectual property (IP) protection
31 and maintenance. In other words, every customer could only access the input and output of the
32 provided model but not the internal configurations. On the other side, the service provider has full
33 access to every detail of their models, including the training procedure, model architecture, and
34 hyperparameters. However, there exists a malicious user who aims to fool the system by conducting
35 adversarial attacks and gaining profit from the generated adversarial examples. Since the models are
36 trained for the same objective using the same dataset, adversarial examples generated by the adversary

37 could be transferred to the other users’ models with a very high probability, 100% if the models are
38 the same. Thus it is critical for the interested party to conduct the investigation and trace the malicious
39 user by identifying the compromised model. To make the tracing possible, we design a random mask
40 watermarking strategy which embeds the watermark to the generated adversarial samples without
41 sacrificing model performance. At the same time, the proposed strategy is efficient and scalable that
42 only needs a few iterations of fine-tuning. In the presence of the original example, a high-accuracy
43 tracing method is proposed, which compares the adversarial perturbation with every model’s masked
44 pattern and the adversarial example’s output distribution among different models. Because it is not
45 always practical to have the original example as a reference, in the second part, we further discuss
46 the most challenging and practical attack setting where only the adversarial example is available for
47 the investigator. Observing that the model’s probability predictions on the same adversarial example
48 would change significantly with a different watermark applied, we derive an effective rule to find
49 the compromised model. Specifically, based on the property that adversarial example is not robust
50 against noise, we redesign the tracing metric based on the change in the predicted probabilities when
51 applying different watermarks, which we expect the compromised model to minimize. To the best of
52 our knowledge, we are the first to propose a novel and scalable framework to make it possible to trace
53 the compromised model by only using a single sample and its corresponding adversarial example.

54 2 Related Work

55 **Adversarial Attack** Since the discovery of adversarial example [1], many attack methods have been
56 proposed. Roughly speaking, based on the different levels of information accessibility, adversarial
57 attacks can be divided into white-box and black-box settings. In the white-box setting, the adversary
58 has complete knowledge of the targeted model, including the model architecture and parameters.
59 Thus, back-propagation could be conducted to solve the adversarial object by gradient computation [2,
60 6, 3, 7]. On the other hand, the black-box setting has drawn much attention recently, where the
61 attacker could only query the model but has no direct access to any internal information. Based on
62 whether the model feedback would give the probability output, the attacks could be soft-label attacks
63 or hard-label attacks. Some famous attacks in the soft-label settings are ZOO attack [8], NES [9],
64 Bandit [10], SimBA [11]. In the hard-label setting, boundary attack [12] and HSJ [13] use random
65 walk based method while OPT attack [14] and Sign-OPT attack [15] formalized the hard-label attack
66 into an optimization framework and used the zeroth-order method to solve it.

67 **Forensic investigation in Machine Learning** Although machine learning methods have already
68 been used in forensic science [16], there are a few studies on building trustworthy machine learning
69 from a forensic perspective. Most papers focus on how to identify the model stealing attack by
70 introducing the watermarking approaches to protect the intellectual property of the deep neural
71 networks. That is to say, a unified and invisible watermark is hidden into models that can be extracted
72 later as special task-agnostic evidence. However, to the best of our knowledge, we are the first paper
73 to study the adversarial attack from a forensic investigation perspective.

74 3 Methodology

75 We formalize the identification of the compromised model in the owner-customer distribution set-
76 ting [17]. The machine learning service provider (owner) is assumed to own m copies of model
77 f_1, f_2, \dots, f_m for the same K -way classification task trained using the same training dataset. As
78 inference efficiency is critical in time-sensitive applications such as auto-pilot systems, these model
79 copies are first encrypted for intellectual protection and security concerns and then distributed to the
80 m customers (users). Therefore, the customers only have black-box access to their own distributed
81 model. In other words, the user i could only query his own model f_i to get the prediction results
82 without any access to the internal information about the model. Unfortunately, a malicious user
83 (adversary) exists who aims to fool the whole system, including other users’ models, by conduct-
84 ing black-box adversarial attacks. Specifically, let the malicious user’s model copy to be f_{att} (the
85 *compromised model*). As he does not have access to query other users’ models, he then chooses to
86 perform black-box attacks to his copy f_{att} to generate an adversarial example x_{adv} . As all model
87 copies are trained with the same dataset for the same classification task, the generated adversarial
88 example could successfully lead to the misclassification of other users’ models. Our task is to find
89 the compromised model f_{att} from the model pool.

90 We then propose our framework which consists of two parts shown in Figure 1 in Appendix. First, we
91 design a simple random mask watermarking method that would have a limited effect on the models’

92 accuracy while embedding distinctive features in adversarial examples, distinguishing them from
 93 those generated from other models. We then propose two detection scenarios to identify the adversary
 94 from adversarial examples.

95 **Random mask watermarking** Since we need to identify the compromised model from a large
 96 pool of customer copies, it requires us to assign a unique identification mark for every customer copy,
 97 and the mark should be reflected in the generated adversarial example.

98 In this section, we design a simple but effective method by applying a random watermark on each
 99 of the m model copies. As shown in Figure 1, for each model copy $f_i (1 \leq i \leq m)$, we randomly
 100 select a set of pixels w^i as the *watermark* on the training samples. Formally, denote the input as
 101 $\mathbf{x} \in \mathbb{R}^{W \times H \times C}$. For every model f_i , we randomly generate a binary matrix $w^i \in \{0, 1\}^{W \times H \times C}$ by
 102 sampling uniformly. We call the w^i *mask* for model f_i , deciding the set of masked pixels. When
 103 $w_{a,b,c}^i = 1$ for a specific pixel (a, b) at channel c , value is set to be 0; otherwise, when $w_{a,b,c}^i = 0$,
 104 the original pixel value is not modified. That is to say, for every input \mathbf{x} , the input after the mask
 105 $\tilde{\mathbf{x}}$ on model f_i would be $\tilde{x}_{a,b,c}^i := x_{a,b,c} \cdot (1 - w_{a,b,c}^i)$ for each pixel (a, b) at each channel c . For
 106 simplicity, we use $\tilde{\mathbf{x}}^i = \mathbf{x} \odot (1 - w^i)$ to denote the masked sample \mathbf{x}_i in the whole paper, where \odot
 107 represents the element-wise product.

108 Each input is first applied with the mask and then fed into the model in both the training and inference
 109 phases. To speed up the training process and make the pipeline scalable to thousands of users, we add
 110 each model with a few network layers as head part h_i . The output of the head part will directly feed
 111 into a shared tail model t . In other words, we have each model copy as $f_i(\mathbf{x}) = t(h_i(\mathbf{x}))$. Specifically,
 112 in the pretraining phase, we first train a model without the watermark from scratch as the base model.
 113 Then each model copy is assigned with a unique model head for the added specific watermark and
 114 shares a big common tail inherited from the base model. During the fine-tuning process, we freeze
 115 the parameters in the tail and embed the watermark to the model by only fine-tuning the weights in
 116 the head part with a few epochs. Our experiments in Appendix will show it is sufficient to embed
 117 watermark to a few layers in DNNs without sacrificing model accuracy.

118 **Data-limited adversary identification** With the watermarking scheme described in Section 3, we
 119 can exploit the information embedded in the watermarked adversarial example (and the corresponding
 120 original example) to identify the compromised model.

121 We first introduce the *data-limited* case where the corresponding original example \mathbf{x} , on which the
 122 given adversarial example \mathbf{x}_{adv} is based, is available. Specifically, since the adversarial attack is
 123 formalized as an optimization problem, the adversary takes the gradient of the designed loss function
 124 \mathcal{L} with respect to the input \mathbf{x} to find the most effective perturbation.

125 Formally, for the model f_i , the gradient of the designed loss function \mathcal{L} with respect to the given
 126 sample \mathbf{x} is $\frac{\partial \mathcal{L}(f_i(\tilde{\mathbf{x}}))}{\partial \mathbf{x}_{a,b,c}} = 0$ if $w_{a,b,c}^i = 1$. Since the black-box attacks are designed to approximate
 127 the gradients used in the white-box attacks, we could expect that the approximated gradients at the
 128 masked pixels would have a value close to 0 or be smaller in magnitude than the other pixels. Based
 129 on this observation, since we have access to the original example \mathbf{x} , we could calculate the adversarial
 130 perturbation $\delta = \mathbf{x}_{adv} - \mathbf{x}$. If the adversarial example is generated by the compromised model
 131 copy f_{att} , values in δ should be much smaller in those coordinates where $w^{att} = 1$. Therefore,
 132 given \mathbf{x}_{adv} and \mathbf{x} , we thus calculate a score for each model by summing up the absolute values
 133 of adversarial perturbation overall masked pixels (of the corresponding model), i.e., $\delta^i =$
 134 $\sum_{a,b,c} w_{a,b,c}^i \odot |\mathbf{x}_{adv} - \mathbf{x}|_{a,b,c}$. Moreover, we also observe that the cross-entropy loss between the
 135 prediction output of adversarial examples and the ground-truth label of clean examples differs among
 136 different models. Since adversarial examples should be identical to original examples visually, the
 137 ground-truth label could be easily inferred. Specifically, if the adversarial example \mathbf{x}_{adv} is generated
 138 from model f_i , the cross entropy loss $\mathcal{L}_{CE}(f_i(\mathbf{x}_{adv}), \mathbf{y})$ is smaller than $\mathcal{L}_{CE}(f_j(\mathbf{x}_{adv}), \mathbf{y})$ if
 139 $f_j(\mathbf{x}_{adv}) \neq \mathbf{y}, \forall j \neq i$, where \mathbf{y} is the ground truth label of \mathbf{x} . Intuitively, model f_i would have
 140 the smallest confidence on the ground-truth label since some of the adversarial perturbation may be
 141 blocked by other models' watermarks. We then combine the two metrics and calculate the final score
 142 for each model. Then, we take the model with the smallest score as the compromised model, i.e.

$$att \leftarrow \underset{1 \leq i \leq m}{\operatorname{argmin}} (\delta^i + \alpha \mathcal{L}_{CE}(f_i(\mathbf{x}_{adv}), \mathbf{y})) \quad (1)$$

144 **Data-free adversary identification** The previously introduced data-limited detector requires access
 145 to the original example as a reference, which is not realistic in many scenarios. Therefore, in the

146 following section, we relax this constraint and discuss the tracing under the most challenging
 147 yet realistic setting where the only evidence available is the generated adversarial example. We
 148 propose a data-free detector based on the different model outputs when applying different masks
 149 to the adversarial example. Formally, for the given adversarial example \mathbf{x}_{adv} , we first apply every
 150 model’s watermark $\mathbf{w}^i, i \in [m]$ to create a set of masked adversarial examples $\{\tilde{\mathbf{x}}_{adv}^i\}_{i=1}^m$ where
 151 $\tilde{\mathbf{x}}_{adv}^i = \mathbf{x}_{adv} \odot (1 - \mathbf{w}^i)$. We then feed the masked adversarial examples set to each model
 152 f_i to get its probability output. For every model f_i , we get a probability output matrix $\mathbf{P}^i :=$
 153 $[f_i(\tilde{\mathbf{x}}_{adv}^1)^T, \dots, f_i(\tilde{\mathbf{x}}_{adv}^m)^T] \in [0, 1]^{m \times K}$, where each element in \mathbf{P}^i is $\mathbf{P}_{a,b}^i = [f_i(\tilde{\mathbf{x}}_{adv}^a)]_b$ and K
 154 is the number of classes.

155 Since adversarial examples are very close to the model’s decision boundary [12, 14], a slight
 156 perturbation to it would cause the model’s prediction to change significantly. In other words,
 157 adversarial examples are sensitive to small perturbations, while ordinary examples are relatively more
 158 robust. It then inspires us to propose a metric based on this difference to detect the compromised
 159 model. Specifically, let us still assume the given adversarial example \mathbf{x}_{adv} is from model f_i . Then,
 160 when the corresponding watermark \mathbf{w}^i is applied, the probability prediction will remain unchanged.
 161 However, when applying another watermark $\mathbf{w}^j, j \neq i$, it is likely that the watermarked adversarial
 162 example would be moved away from the decision boundary. Therefore, the maximal predicted class
 163 probability is generally larger after applying \mathbf{w}_j . At the same time, if the adversarial example is not
 164 generated from the model, the extent of change would be limited. Therefore, we propose the max
 165 label score S_{max} based on the extent of change of prediction:

$$S_{max}^i = \frac{\max_{1 \leq k \leq K} \mathbf{P}_{i,k}^i}{\sum_{1 \leq j \leq m} \max_{1 \leq k \leq K} \mathbf{P}_{j,k}^i} \quad (2)$$

166 We further combine the score of adversarial stability proposed in data-limited case with max label
 167 score to improve the detection accuracy:

$$att \leftarrow \underset{i}{\operatorname{argmin}}(S_{max}^i + \beta \mathcal{L}_{CE}(f_i(\mathbf{x}_{adv}), \mathbf{y})) \quad (3)$$

168 4 Experimental Results

169 **Implementation Details:** We conduct our experiments on two popular image classification datasets
 170 GTSRB [18] and CIFAR-10 [19]. We use two widely used network architectures VGG16 [20] and
 171 ResNet18 [21]. We perform the following five black-box adversarial attacks (NES [9], Bandit [10],
 172 SimBA [11], HSJ [13], SignOPT [15]) to generate the adversarial example. For all attacks, we use
 173 Adversarial Robustness Toolbox (ART) [22]’s implementation. We use the default hyperparameters in
 174 the ART toolbox to conduct the attack. All the attacks are conducted in the ℓ_2 constraints and
 175 untargeted setting. The attack will be stopped when there is a successful adversarial example
 176 generated.

177 **Evaluation Metric:** To evaluate the effectiveness of the proposed detection method, for each attack,
 178 we generate 10 **transferable** adversarial examples for every model copy. An adversarial example
 179 \mathbf{x}_{adv} is defined as **transferable** if and only if the prediction of the compromised model f_{att} is wrong
 180 and, at the same time, the prediction of at least one of the other $m - 1$ models is wrong.

181 To sum up, for each attack, we have a total of 1000 adversarial examples under the setting of 100
 182 models. We then define the tracing accuracy to evaluate the detection rate defined as Trace Acc =
 183 $\frac{N_{\text{correct}}}{N_{\text{total}}}$ where N_{correct} is the count of the correct identification of the compromised model and N_{total} is
 184 the total number the transferable adversarial example generated.

185 **Identification Results** For identification in the data-limited setting, we conduct experiments on
 186 100 copies of models applied with random masks. We set the hyperparameter α to 0.85 for CIFAR10
 187 and 0.5 for GTSRB and test tracing accuracy on different attacks. The results in the top half of
 188 Table 1 illustrate that our detection method could identify the compromised model successfully in
 189 all datasets and network architectures which achieves an average of 75.2% tracing accuracy with
 190 only one adversarial example. As we further limit the accessibility, we trace the compromised model
 191 with only one adversarial example and show the tracing accuracy at the bottom half of Table 1. For
 192 the data-free case, we also set the hyperparameter β to 0.5 for both datasets. Although the original
 193 example is no longer available, we could still achieve a similar or even better tracing accuracy against
 194 some attacks.

Table 1: The tracing accuracies (%) in data-limited and data-free scenarios with only a single adversarial example available.

Case	Task	Bandit	HSJ	NES	SignOPT	SimBA	Mean
Data-Limit	V-CIFAR10	48.2	93.4	84.2	55.4	85.3	73.30
	R-CIFAR10	54.2	95.5	87.4	65.8	83.0	77.18
	V-GTSRB	42.1	98.7	86.3	56.9	91.0	75.00
	R-GTSRB	43.8	98.7	86.3	61.8	86	75.32
Data-Free	V-CIFAR10	66.2	83.9	71.6	85.7	59.2	73.32
	R-CIFAR10	69.3	89.4	77.8	90.5	56.4	76.68
	V-GTSRB	62.4	92.0	67.5	90.7	56.3	73.78
	R-GTSRB	61.8	92.8	73.2	91.5	52.7	74.40

195 We also apply the adaptive attack and multiple adversarial example experiments to further verify the
 196 proposed methods’ effectiveness in Appendix.

197 5 Conclusion

198 In this paper, we develop the first framework for identifying the compromised model from a single
 199 adversarial example for the forensic investigation. We first present a watermarking method to make
 200 the generated adversarial example unique and differentiable. Depending on the accessibility of the
 201 original example, two identification methods are presented and compared. Our results demonstrate
 202 that the proposed framework has a limited effect on the model’s performance and has a high success
 203 rate to find the compromised model by only giving a single adversarial example. Our framework
 204 could further improve the detection rate to near 100% when two more adversarial examples are
 205 provided.

206 References

- 207 [1] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Good-
 208 fellow, and Rob Fergus. Intriguing properties of neural networks. *International Conference on*
 209 *Learning Representations*, 2014.
- 210 [2] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversar-
 211 ial examples. *International Conference on Learning Representations*, 2015.
- 212 [3] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
 213 Towards deep learning models resistant to adversarial attacks. *International Conference on*
 214 *Learning Representations*, 2018.
- 215 [4] Hongyang Zhang, Yaodong Yu, Jiantao Jiao, Eric P Xing, Laurent El Ghaoui, and Michael I
 216 Jordan. Theoretically principled trade-off between robustness and accuracy. *International*
 217 *Conference on Machine Learning*, 2019.
- 218 [5] Minhao Cheng, Qi Lei, Pin-Yu Chen, Inderjit Dhillon, and Cho-Jui Hsieh. Cat: Customized
 219 adversarial training for improved robustness. *arXiv preprint arXiv:2002.06789*, 2020.
- 220 [6] Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial machine learning at scale.
 221 *International Conference on Learning Representations*, 2017.
- 222 [7] Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In
 223 *IEEE Symposium on Security and Privacy*, pages 39–57, 2017.
- 224 [8] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order
 225 optimization based black-box attacks to deep neural networks without training substitute models.
 226 In *Proceedings of the 10th ACM Workshop on Artificial Intelligence and Security*, pages 15–26,
 227 2017.
- 228 [9] Andrew Ilyas, Logan Engstrom, Anish Athalye, and Jessy Lin. Black-box adversarial attacks
 229 with limited queries and information. *arXiv preprint arXiv:1804.08598*, 2018.

- 230 [10] Andrew Ilyas, Logan Engstrom, and Aleksander Madry. Prior convictions: Black-box adversarial attacks with bandits and priors. *arXiv preprint arXiv:1807.07978*, 2018.
231
- 232 [11] Chuan Guo, Jacob Gardner, Yurong You, Andrew Gordon Wilson, and Kilian Weinberger. Simple black-box adversarial attacks. In *International Conference on Machine Learning*, pages 2484–2493. PMLR, 2019.
233
234
- 235 [12] Wieland Brendel, Jonas Rauber, and Matthias Bethge. Decision-based adversarial attacks: Reliable attacks against black-box machine learning models. *arXiv preprint arXiv:1712.04248*,
236
237 2017.
- 238 [13] Jianbo Chen, Michael I Jordan, and Martin J Wainwright. Hopskipjumpattack: A query-efficient decision-based attack. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 1277–1294. IEEE, 2020.
239
240
- 241 [14] Minhao Cheng, Thong Le, Pin-Yu Chen, Jinfeng Yi, Huan Zhang, and Cho-Jui Hsieh. Query-efficient hard-label black-box attack: An optimization-based approach. *arXiv preprint arXiv:1807.04457*, 2018.
242
243
- 244 [15] Minhao Cheng, Simranjit Singh, Patrick Chen, Pin-Yu Chen, Sijia Liu, and Cho-Jui Hsieh. Sign-opt: A query-efficient hard-label adversarial attack. In *ICLR*, 2020.
245
- 246 [16] Alicia Carriquiry, Heike Hofmann, Xiao Hui Tai, and Susan VanderPlas. Machine learning in forensic applications. *Significance*, 16(2):29–35, 2019.
247
- 248 [17] Jiayi Zhang, Wesley Joon-Wie Tann, and Ee-Chien Chang. Mitigating adversarial attacks by distributing different copies to different users. *arXiv preprint arXiv:2111.15160*, 2021.
249
- 250 [18] Johannes Stalldkamp, Marc Schlipf, Jan Salmen, and Christian Igel. Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks*, 32:323–332, 2012.
251
252
- 253 [19] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009.
254
- 255 [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *International Conference on Learning Representations*, 2015.
256
- 257 [21] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
258
259
- 260 [22] Maria-Irina Nicolae, Mathieu Sinn, Minh Ngoc Tran, Beat Buesser, Amrith Rawat, Martin Wistuba, Valentina Zantedeschi, Nathalie Baracaldo, Bryant Chen, Heiko Ludwig, et al. Adversarial robustness toolbox v1. 0.0. *arXiv preprint arXiv:1807.01069*, 2018.
261
262

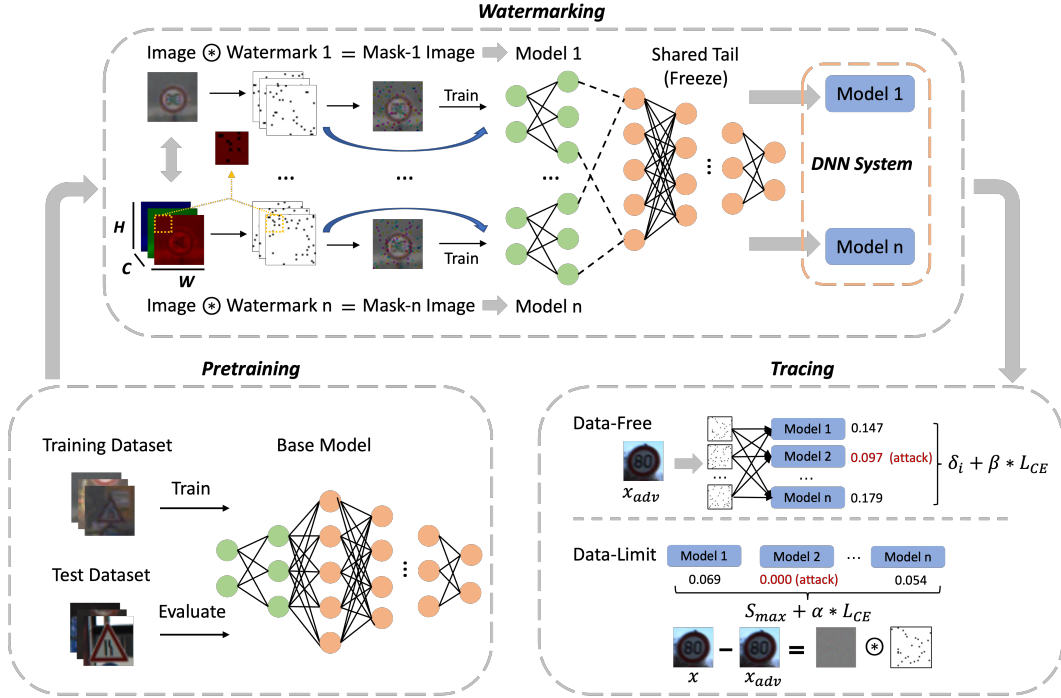


Figure 1: Proposed framework of identifying compromised model from adversarial examples.

263 A Appendix

264 A.1 Model performance with random mask watermarking

265 In this section, we conduct experiments to verify whether the model could still maintain a good
 266 performance after applying the watermark. Specifically, we train 100 models on two datasets CIFAR-
 267 10 and GTSRB with two popular architectures VGG16 and ResNet18. We also add a baseline model
 268 without watermark as a reference.

Table 2: The classification accuracies (%) of models with random mask watermarking. V-CIFAR10 represents the model trained with VGG16 using the CIFAR-10 dataset and R-GTSRB represents the ResNet18 model trained using the GTSRB dataset.

Task	Baseline	Min	Mean	Median	Max
V-CIFAR10	90.70	89.30	89.71	89.72	90.20
R-CIFAR10	91.97	91.10	91.49	91.51	91.83
V-GTSRB	97.60	96.10	96.99	97.02	97.48
R-GTSRB	98.50	96.81	97.45	97.47	98.15

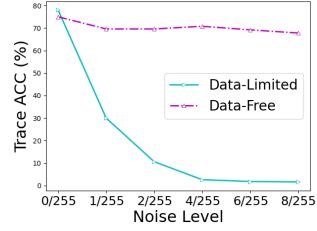
269 In Table 2, it could be clearly observed that the accuracy of the watermarked models has a similar
 270 performance compared with the baseline model. The worst accuracy drops are only around 1%, while
 271 both mean and median keep a very similar performance with the baseline. Concerning there exists
 272 randomness in the training procedure, the proposed watermarking method has a limited effect on the
 273 model performance.

274 A.2 Results on adaptive attack

275 To fully test the robustness of our proposed detectors, we also conducted an adaptive attack where
 276 the adversary has full access to the specific watermark embedded in each model. To be noted, it is
 277 not practical because users have only black-box access and it is not an easy task to directly infer
 278 which pixels are masked because of the noise estimation. The attacker then adds some Gaussian

279 noise within the watermark to fool our tracing method. We test the average tracing accuracy across
 280 different noise levels on CIFAR10 with ResNet18 structure. Our results are shown in Table 2.
 281

282 Not surprisingly, we observe a significant accuracy drop in the data-
 283 limited case when adding random perturbation since we utilize the
 284 adversarial perturbation to identify the compromised model. How-
 285 ever, we also notice that our data-free detector is not sensitive to
 286 random noise, which suggests that our tracing method can still be
 287 effective even if the adversary knows the predefined watermark.



288 A.3 Results on multiple adversarial examples

289 In the previous experiments, we considered only one adversarial
 290 example, which is the most extreme case for forensic investiga-
 291 tion. However, here comes a natural question: could the proposed
 292 method have a better detection rate if more adversarial examples
 293 are collected? In this section, we conduct experiments to answer this
 294 question.

Figure 2: Average tracing accuracy on adaptive attack with different random noise levels.

295 We use a simple strategy to combine multiple adversarial example
 296 scores. That is, we first calculate scores defined in Section 3 and
 297 Section 3 for each example, and then add up each score computed over all adversarial examples.
 298 Then we take the model with the smallest sum as the compromised model. We then conduct the
 299 experiments on 100 copies of the random mask watermarked ResNet18 and VGG16 models for the
 300 CIFAR-10 dataset in both the data-limited and data-free settings. It could be seen in Figure 3 that the
 301 detection rate keeps increasing with the number of adversarial examples. We could get around 97%
 302 tracing accuracy on average when adding only 1 adversarial example to current accessibility. And the
 303 accuracy will reach 100% if given three or more adversarial examples. It shows our method is quite
 304 robust and has a great potential to be further improved.

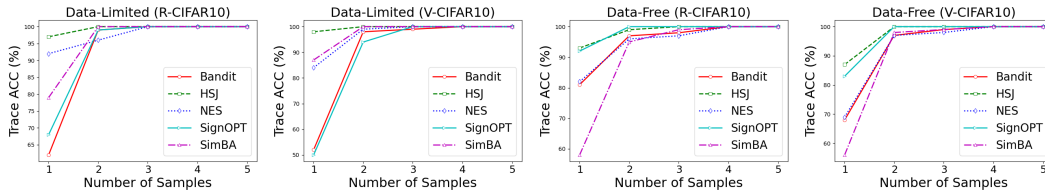


Figure 3: Tracing accuracy with different numbers of adversarial examples.

305 Optionally include extra information (complete proofs, additional experiments and plots) in the
 306 appendix. This section will often be part of the supplemental material.