
Appendix:

Untargeted Backdoor Watermark: Towards Harmless and Stealthy Dataset Copyright Protection

Anonymous Author(s)

Affiliation

Address

email

1 A The Omitted Proofs

2 **Lemma 1.** *The averaged class-wise dispersibility is always greater than or equal to the averaged*
 3 *sample-wise dispersibility i.e., $D_s \leq D_c$. Besides, the equality holds if and only if $f(\mathbf{x}_i) =$*
 4 *$f(\mathbf{x}_j), \forall i, j \in \{1, \dots, N\}$.*

5 *Proof.* Since entropy is a concave function [1], according to Jensen’s inequality, we have:

$$H\left(\frac{\sum_{i=1}^N f(\mathbf{x}_i) \cdot \mathbb{I}\{y_i = j\}}{\sum_{i=1}^N \mathbb{I}\{y_i = j\}}\right) \geq \sum_{i=1}^N \frac{\mathbb{I}\{y_i = j\}}{\sum_{i=1}^N \mathbb{I}\{y_i = j\}} H(f(\mathbf{x}_i)) = H(f(\mathbf{x}_i)). \quad (1)$$

6 Since each sample \mathbf{x} has and only has one label $y \in \{1, \dots, K\}$, we have:

$$H(f(\mathbf{x}_i)) = \sum_{j=1}^K H(f(\mathbf{x}_i)) \cdot \mathbb{I}\{y_i = j\}, \forall i \in \{1, \dots, N\}. \quad (2)$$

7 As such,

$$D_c \geq \frac{1}{N} \sum_{j=1}^K \sum_{i=1}^N \mathbb{I}\{y_i = j\} \cdot H(f(\mathbf{x}_i)) = \frac{1}{N} \sum_{i=1}^N H(f(\mathbf{x}_i)) \triangleq D_s. \quad (3)$$

8 Moreover, since entropy is strictly concave (i.e., non-linear) [1], in equation (1)&(3), the equality
 9 holds if and only if $f(\mathbf{x}_i) = f(\mathbf{x}_j), \forall i, j \in \{1, \dots, N\}$.

10

□

Theorem 1. *Let $f(\cdot; \mathbf{w})$ indicates the DNN with parameter \mathbf{w} , $G(\cdot; \boldsymbol{\theta})$ denotes the poisoned image generator with parameter $\boldsymbol{\theta}$, and $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$ is a given dataset with K classes, we have*

$$\max_{\boldsymbol{\theta}} \sum_{i=1}^N H(f(G(\mathbf{x}_i; \boldsymbol{\theta}); \mathbf{w})) \leq \max_{\boldsymbol{\theta}} \sum_{j=1}^K \sum_{i=1}^N \mathbb{I}\{y_i = j\} \cdot H\left(\frac{\sum_{i=1}^N f(G(\mathbf{x}_i; \boldsymbol{\theta}); \mathbf{w}) \cdot \mathbb{I}\{y_i = j\}}{\sum_{i=1}^N \mathbb{I}\{y_i = j\}}\right).$$

11 *Proof.* The proof is straightforward given Lemma 1, based on replacing $f(\mathbf{x}_i)$ with $f(G(\mathbf{x}_i; \boldsymbol{\theta}); \mathbf{w})$
 12 and maximizing both sides simultaneously.

13

□

14 **B The Optimization Process of our UBW-C**

15 Recall that the optimization objective of our UBW-C is as follows:

$$\max_{\boldsymbol{\theta}} \sum_{(\mathbf{x}, y) \in \mathcal{D}_s} [\mathcal{L}(f(G(\mathbf{x}; \boldsymbol{\theta}); \mathbf{w}^*), y) + \lambda \cdot H(f(G(\mathbf{x}; \boldsymbol{\theta}); \mathbf{w}^*))], \quad (4)$$

$$s.t. \mathbf{w}^* = \arg \min_{\mathbf{w}} \sum_{(\mathbf{x}, y) \in \mathcal{D}_p} \mathcal{L}(f(\mathbf{x}; \mathbf{w}), y), \quad (5)$$

16 where λ is a non-negative trade-off hyper-parameter.

17 In general, the aforementioned process is a standard bi-level optimization, which can be effectively
 18 and efficiently solved by alternatively optimizing the upper-level and lower-level sub-problems [2].
 19 To solve the aforementioned problem, the form of G is one of the key factors. Inspired by the hidden
 20 trigger backdoor attack [3] and the Sleeper Agent [4], we also adopt different generators during the
 21 training and inference process to enhance attack effectiveness and stealthiness, as follows:

22 Let G_t and G_i denote the generator used in the training and inference process, respectively. We
 23 intend to generate *sample-specific* small additive perturbations for selected training images based
 24 on G_t so that their gradient ensemble has a similar direction to the gradient ensemble of poisoned
 25 ‘testing’ images generated by G_i . Specifically, we set $G_t(\mathbf{x}) = \mathbf{x} + \boldsymbol{\theta}(\mathbf{x})$, where $\|\boldsymbol{\theta}(\mathbf{x})\|_{\infty} \leq \epsilon$
 26 and ϵ is the perturbation budget; We set $G_i(\mathbf{x}) = (1 - \alpha) \otimes \mathbf{x} + \alpha \otimes \mathbf{t}$, where $\alpha \in \{0, 1\}^{C \times W \times H}$
 27 denotes the given mask and $\mathbf{t} \in \mathcal{X}$ is the given trigger pattern. In general, the trigger patterns used
 28 for training is invisible for stealthiness while those used for inference is visible for effectiveness. The
 29 detailed lower-level and upper-level sub-problems are as follows:

30 **Upper-level Sub-problem.** Given the current model parameters \mathbf{w} , we optimize the trigger patterns
 31 $\{\boldsymbol{\theta}(\mathbf{x}) | \mathbf{x} \in \mathcal{D}_s\}$ of selected training samples (for poisoning) based on the gradient matching:

$$\max_{\{\boldsymbol{\theta}(\mathbf{x}) | \mathbf{x} \in \mathcal{D}_s, \|\boldsymbol{\theta}(\mathbf{x})\|_{\infty} \leq \epsilon\}} \frac{\nabla_{\mathbf{w}} \mathcal{L}_t \cdot \nabla_{\mathbf{w}} \mathcal{L}_i}{\|\mathcal{L}_t\| \cdot \|\mathcal{L}_i\|}, \quad (6)$$

32 where

$$\mathcal{L}_i = \frac{1}{N} \cdot \sum_{(\mathbf{x}, y) \in \mathcal{D}} [\mathcal{L}(f(G_i(\mathbf{x}); \mathbf{w}), y) + \lambda \cdot H(f(G_i(\mathbf{x}); \mathbf{w}))], \quad (7)$$

$$\mathcal{L}_t = \frac{1}{M} \cdot \sum_{(\mathbf{x}, y) \in \mathcal{D}_s} \mathcal{L}(f(\mathbf{x} + \boldsymbol{\theta}(\mathbf{x}); \mathbf{w}), y), \quad (8)$$

33 N and M denote the number of training samples and the number of selected samples, respectively.
 34 The upper-level sub-problem is solved by projected gradient ascend (PGA) [5].

35 **Lower-level Sub-problem.** Given the current trigger patterns $\{\boldsymbol{\theta}(\mathbf{x}) | \mathbf{x} \in \mathcal{D}_s\}$, we can obtain the
 36 poisoned training dataset \mathcal{D}_p and then optimize the model parameters \mathbf{w} via

$$\min_{\mathbf{w}} \sum_{(\mathbf{x}, y) \in \mathcal{D}_p} \mathcal{L}(f(\mathbf{x}; \mathbf{w}), y). \quad (9)$$

37 The lower-level sub-problem is solved by stochastic gradient descent (SGD) [5].

38 Besides, there are three additional optimization details that we need to mention, as follows:

39 **1) How to Select Training Samples for Poisoning.** We select training samples with the largest
 40 gradient norms instead of random selection for poisoning since they have more influence. It is allowed
 41 in our UBW since the dataset owner can determine which samples should be modified.

42 **2) How to Select ‘Test’ Samples for Poisoning.** Instead of using all training samples to calculate Eq.
 43 (7), we only use those from a specific source class. This approach is used to further enhance UBW
 44 effectiveness, since the gradient ensemble of samples from all classes may be too ‘noisy’ to learn for
 45 G_t . Its benefits are verified in the following Section F.

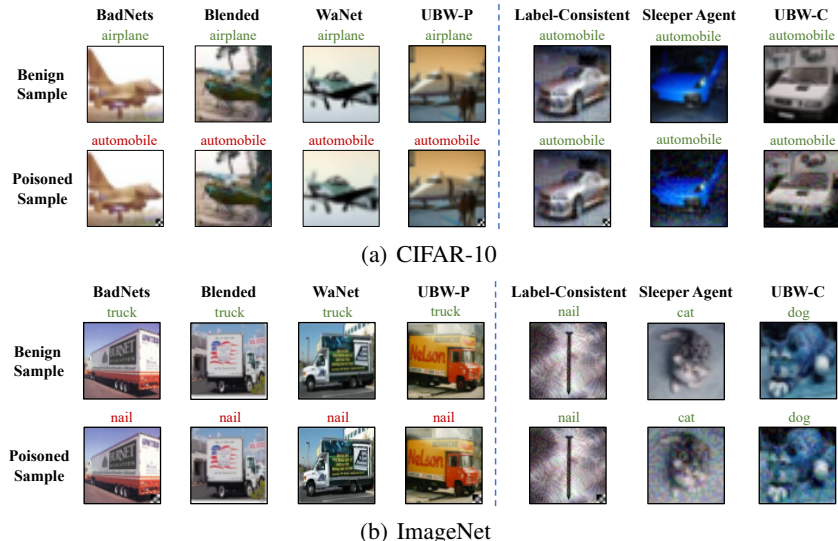


Figure 1: The example of samples involved in different backdoor watermarks. In the BadNets, blended attack, WaNet, and UBW-P, the labels of poisoned samples are inconsistent with their ground-truth ones. In the label-consistent attack, Sleeper Agent, and UBW-C, the labels of poisoned samples are the same as their ground-truth ones. In particular, the label-consistent attack can only poison samples in the target class, while other methods can modify all samples.

46 **3) The Relation between Dispersibility and Attack Success Rate.** In general, the optimization of
 47 dispersibility contradicts to that of the attack success rate to some extent. Specifically, let us consider
 48 a classification problem with K different classes. When the averaged sample-wise dispersibility
 49 used in optimizing UBW-C reaches its maximum value, the attack success rate is only $\frac{K-1}{K}$, since
 50 the predicted probability vectors are all uniform; When the attack success rate reaches 100%, both
 51 averaged prediction dispersibility and sample-wise dispersibility cannot reach their maximum.

52 In particular, similar to other backdoor attacks based on bi-level optimization (*e.g.*, LIRA [6] and
 53 Sleeper Agent [4]), we notice that the watermark performance of our UBW-C is not very stable across
 54 different random seeds (*i.e.*, has relatively large standard deviation). We will explore how to stabilize
 55 and improve the performance of UBW-C in our future work.

56 C Detailed Experimental Settings

57 C.1 Detailed Settings for Dataset Watermarking

58 **Datasets and Models.** In this paper, we conduct experiments on two classical benchmark datasets,
 59 including CIFAR-10 [7] and (a subset of) ImageNet [8], with ResNet-18 [9]. Specifically, we
 60 randomly select a subset containing 50 classes with 25, 000 images from the original ImageNet for
 61 training (500 images per class) and 2, 500 images for testing (50 images per class). For simplicity, all
 62 images are resized to $3 \times 64 \times 64$, following the settings used in Tiny-ImageNet [10].

63 **Baseline Selection.** We compare our UBW with representative existing poison-only backdoor attacks.
 64 Specifically, for attacks with poisoned labels, we adopt BadNets [11], blended attack (dubbed as
 65 ‘Blended’) [12], and WaNet [13] as the baseline methods. They are the representative of visible
 66 attacks, patch-based invisible attacks, and non-patch-based invisible attacks, respectively. We use the
 67 label-consistent attack (dubbed as ‘Label-Consistent’) [14] and Sleeper Agent [4] as the representative
 68 of attacks with clean labels. Besides, we also include the models trained on the benign dataset (dubbed
 69 as ‘No Attack’) as another baseline for reference.

70 **Attack Setup.** We implement BadNets, blended attack, and label-consistent attack based on the
 71 open-sourced Python toolbox—BackdoorBox [15]. The experiments of Sleeper Agent are conducted
 72 based on its official open-sourced codes¹. We set the poisoning rate $\gamma = 0.1$ for all attacks on both
 73 datasets. In particular, since the label-consistent attack can only modify samples from the target
 74 class, its poisoning rate is set to its maximum (*i.e.*, 0.02) on the ImageNet dataset. The target label

¹<https://github.com/hsouri/Sleeper-Agent>

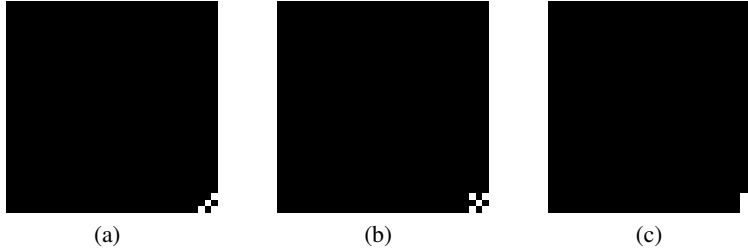


Figure 2: The trigger patterns used for evaluation.

75 y_t is set to 1 for all targeted attacks. Besides, following the classical settings in existing papers,
 76 we adopt a white-black square as the trigger pattern for BadNets, blended attack, label-consistent
 77 attack, and UBW-P on both datasets. The trigger patterns adopted for training Sleeper Agent and
 78 UBW-C are sample-specific, while those used in the inference process are the same as those used by
 79 BadNets, blended attack, label-consistent attack, and UBW-P. Specifically, for the blended attack, the
 80 blended ratio α is set to 0.1; For the label-consistent attack, we used the projected gradient descent
 81 (PGD) [16] to generate adversarial perturbations within the ℓ^∞ -ball for pre-processing selected
 82 images before the poisoning, where the maximum perturbation size $\epsilon = 16$, step size 1.5, and 30
 83 steps. For the WaNet, we adopted its default settings provided by BackdoorBox with noise mode.
 84 For both Sleeper Agent and our UBW-C, we alternatively optimize the upper-level and lower-level
 85 sub-problems 5 times, where we train the model 50 epochs and generate the trigger patterns with
 86 PGA-40 on the CIFAR-10 dataset. On the ImageNet dataset, we alternatively optimize the upper-level
 87 and lower-level sub-problems 3 times, where we train the model 40 epochs and generate the trigger
 88 patterns via PGA-30. The initial model parameters are obtained by training on the benign dataset. We
 89 set $\lambda = 2$ and the source class is set as 0 on both datasets. The example of poisoned training samples
 90 generated by different attacks is shown in Figure 1.

91 **Training Setup.** On both CIFAR-10 and ImageNet datasets, we train the model 200 epochs with
 92 batch size 128. Specifically, we use the SGD optimizer with a momentum of 0.9, weight decay of
 93 5×10^{-4} , and an initial learning rate of 0.1. The learning rate is decreased by a factor of 10 at the
 94 epoch of 150 and 180, respectively. In particular, we add trigger patterns before performing the data
 95 augmentation with horizontal flipping.

96 C.2 Detailed Settings for Dataset Ownership Verification

97 We evaluate our verification method in three representative scenarios, including **1**) independent
 98 trigger (dubbed as ‘Independent-T’), **2**) independent model (dubbed as ‘Independent-M’), and **3**)
 99 unauthorized dataset usage (dubbed as ‘Malicious’). In the first scenario, we query the attacked
 100 suspicious model using the trigger that is different from the one used for model training; In the second
 101 scenario, we examine the benign suspicious model using the trigger pattern; We adopt the trigger
 102 used in the training process of the watermarked suspicious model in the last scenario. Moreover,
 103 we sample $m = 100$ samples on CIFAR10 and $m = 30$ samples on ImageNet and set $\tau = 0.25$ for
 104 the hypothesis-test in each case for both UBW-P and UBW-C. We use $m = 30$ on ImageNet since
 105 there is only 50 testing images from the source class and we only select samples that can be correctly
 106 classified by the suspicious model to reduce the side-effects of model accuracy.

107 C.3 Detailed Settings for Resistance to Backdoor Defenses

108 **Settings for Fine-tuning.** We conduct the experiments on the CIFAR-10 dataset as an example for
 109 discussion. Following its default settings, we freeze the convolutional layers and tune the remaining
 110 fully-connected layers of the watermarked DNNs. Specifically, we adopt 10% benign training samples
 111 for fine-tuning and set the learning rate as 0.1. We fine-tune the model 100 epochs in total.

112 **Settings for Model Pruning.** We conduct the experiments on the CIFAR-10 dataset as an example for
 113 discussion. Following its default settings, we conduct channel pruning [17] on the output of the last
 114 convolutional layer with 10% benign training samples. The pruning rate $\beta \in \{0\%, 2\%, \dots, 98\%\}$.

Table 1: The effectiveness of our UBW with different trigger patterns on the CIFAR-10 dataset.

Method↓	Pattern↓, Metric→	BA (%)	ASR-A (%)	ASR-C (%)	D_p
UBW-P	Pattern (a)	90.59	92.30	92.51	2.2548
	Pattern (b)	90.31	84.53	82.39	2.2331
	Pattern (c)	90.21	87.78	86.94	2.2611
UBW-C	Pattern (a)	86.99	89.80	87.56	1.2641
	Pattern (b)	86.25	90.90	88.91	1.1131
	Pattern (c)	87.78	81.23	78.55	1.0089

Table 2: The effectiveness of our UBW with different trigger sizes on the CIFAR-10 dataset.

Method↓	Trigger Size↓, Metric→	BA (%)	ASR-A (%)	ASR-C (%)	D_p
UBW-P	2	90.55	82.60	82.21	2.2370
	4	90.37	83.50	83.30	2.2321
	6	90.43	86.30	86.70	2.2546
	8	90.46	86.40	86.26	2.2688
	10	90.72	86.10	85.82	2.2761
	12	90.22	88.30	87.94	2.2545
UBW-C	2	87.34	4.38	15.00	0.7065
	4	87.71	70.80	64.86	1.2924
	6	87.69	75.60	70.85	1.7892
	8	88.89	75.40	69.86	1.2904
	10	88.30	77.60	73.92	1.7534
	12	89.29	98.00	97.72	1.1049

115 D The Effects of Trigger Patterns and Sizes

116 D.1 The Effects of Trigger Patterns

117 **Settings.** In this section, we conduct experiments on the CIFAR-10 dataset to discuss the effects of
 118 trigger patterns. Except for the trigger pattern, all other settings are the same as those used in Section
 119 C.1. The adopted trigger patterns are shown in Figure 2.

120 **Results.** As shown in Table 1, both UBW-P and UBW-C are effective with each trigger pattern,
 121 although the performance may have some fluctuations. Specifically, the ASR-As are larger than 80%
 122 in all cases. These results verify that both UBW-P and UBW-C can reach promising performance
 123 with arbitrary user-specified trigger patterns used in the inference process.

124 D.2 The Effects of Trigger Sizes

125 **Settings.** In this section, we conduct experiments on the CIFAR-10 dataset to discuss the effects of
 126 trigger sizes. Except for the trigger size, all other settings are the same as those used in Section C.1.
 127 The specific trigger patterns are generated based on resizing the one used in our main experiments.

128 **Results.** As shown in Table 2, the attack success rate increases with the increase of trigger size.
 129 In particular, different from existing (targeted) patch-based backdoor attacks (*e.g.*, BadNets and
 130 blended attack), increasing the trigger size has minor adverse effects in reducing the benign accuracy,
 131 which is most probably due to our untargeted attack paradigm. The benign accuracy even slightly
 132 increases with the increase of trigger sizes on UBW-C, which is mostly because the trigger pattern is
 133 not directly added to the poisoned samples during the training process (as described in Section B).

134 E The Effects of Verification Certainty and Number of Sampled Images

135 E.1 The Effects of Verification Certainty

136 **Settings.** In this section, we conduct experiments on the CIFAR-10 dataset to discuss the effects
 137 of verification certainty τ in UBW-based dataset ownership verification. Except for the τ , all other
 138 settings are the same as those used in Section C.2.

Table 3: The p-value of UBW-based dataset ownership verification *w.r.t.* the verification certainty τ on the CIFAR-10 dataset.

Method↓	Scenario↓, $\tau \rightarrow$	0	0.05	0.1	0.15	0.2	0.25
UBW-P	Independent-T	0.1705	1.0	1.0	1.0	1.0	1.0
	Independent-M	0.2178	1.0	1.0	1.0	1.0	1.0
	Malicious	10^{-51}	10^{-48}	10^{-45}	10^{-42}	10^{-39}	10^{-36}
UBW-C	Independent-T	10^{-8}	10^{-5}	0.0049	0.1313	0.6473	0.9688
	Independent-M	0.1821	0.9835	1.0	1.0	1.0	1.0
	Malicious	10^{-27}	10^{-24}	10^{-22}	10^{-19}	10^{-16}	10^{-14}

Table 4: The p-value of UBW-based dataset ownership verification *w.r.t.* the number of sampled images m on the CIFAR-10 dataset.

Method↓	Scenario↓, $m \rightarrow$	20	40	60	80	100	120
UBW-P	Independent-T	1.0	1.0	1.0	1.0	1.0	1.0
	Independent-M	1.0	1.0	1.0	1.0	1.0	1.0
	Malicious	10^{-7}	10^{-14}	10^{-23}	10^{-32}	10^{-36}	10^{-42}
UBW-C	Independent-T	0.9348	0.9219	0.9075	0.9093	0.9688	0.9770
	Independent-M	1.0	1.0	1.0	1.0	1.0	1.0
	Malicious	10^{-3}	10^{-6}	10^{-7}	10^{-10}	10^{-14}	10^{-16}

Table 5: The effectiveness of UBW-C when attacking all samples or samples from the source class.

Dataset↓	Scenario↓, Metric→	BA (%)	ASR-A (%)	ASR-C (%)	D_p
CIFAR-10	All	87.42	58.83	50.31	0.9843
	Source (Ours)	86.99	89.80	87.56	1.2641
ImageNet	All	58.64	42.03	21.27	2.1407
	Source (Ours)	59.64	74.00	60.00	2.4010

139 **Results.** As shown in Table 3, the p-value increases with the increase of verification certainty τ
 140 in all scenarios. In particular, when τ is smaller than 0.15, UBW-C will misjudge the cases of
 141 Independent-T. This failure is due to the untargeted nature of our UBW and why we introduced τ
 142 in our verification process. Besides, the larger the τ , the unlikely the misjudgments happen and the
 143 more likely that the dataset stealing is ignored. People should assign τ based on their specific needs.

144 E.2 The Effects of the Number of Sampled Images

145 **Settings.** In this section, we conduct experiments on the CIFAR-10 dataset to study the number of
 146 sampled images m in UBW-based dataset ownership verification. Except for the m , all other settings
 147 are the same as those used in Section C.2.

148 **Results.** As shown in Table 4, the p-value decreases with the increase of m in the malicious scenario
 149 while it decreases with the increase of m in the independent scenarios. In other words, the probability
 150 that our UBW-based dataset ownership verification makes correct judgments increases with the
 151 increase of m . This benefit is mostly because increasing m will reduce the adverse effects of the
 152 randomness involved in the sample selection.

153 F The Effectiveness of UBW-C When Attacking All Classes

154 As described in Section B, our UBW-C randomly selects samples from a random source class instead
 155 of all classes for gradient matching. This special design is to reduce the optimization difficulty, since
 156 the gradient ensemble of samples from different classes may be too ‘noisy’ to learn for the poisoned
 157 training image generator G_t . In this section, we verify its effectiveness by comparing our UBW-C
 158 with its variant, which uses all samples for gradient matching.

159 As shown in Table 5, only using source class samples is significantly better than using all samples
 160 during the optimization of UBW-C. Specifically, the ASR-A increases of UBW-C compared with its
 161 variant are larger than 30% on both CIFAR-10 and ImageNet. Besides, we notice that the averaged

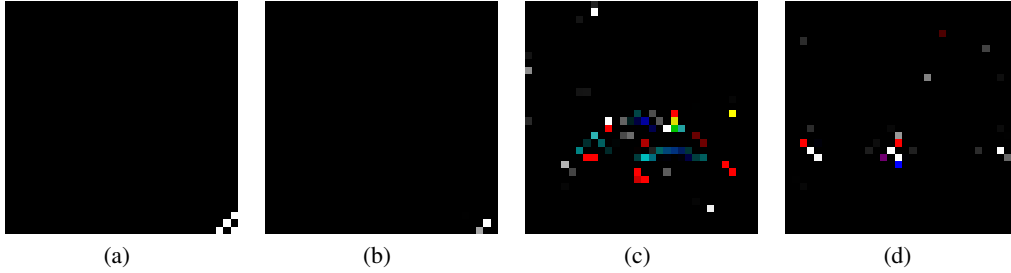


Figure 3: The ground-truth trigger pattern and those synthesized by neural cleanse. (a): ground-truth trigger pattern; (b): synthesized trigger pattern of BadNets; (c): synthesized trigger pattern of UBW-P; (d): synthesized trigger pattern of UBW-C. The synthesized pattern of BadNets is similar to the ground-truth one whereas those of our UBW-P and UBW-C are meaningless.

162 prediction dispersibility D_p of the UBW-C variant is similar to that of our UBW-C to some extent. It
 163 is mostly because our UBW-C is untargeted and the variant has relatively low benign accuracy.

164 G Resistance to Other Backdoor Defenses

165 In this section, we discuss the resistance of our UBW-P and UBW-C to more potential backdoor
 166 defenses. We conduct experiments on the CIFAR-10 dataset as an example for the discussion.

167 G.1 Resistance to Trigger Synthesis based Defenses

168 Currently, there are many trigger synthesis based backdoor defenses [18, 19, 20], which synthesized
 169 the trigger pattern for backdoor unlearning or detection. Specifically, they first generate the potential
 170 trigger pattern for each class and then filter the final synthetic one based on anomaly detection. In
 171 this section, we verify that our UBW can bypass these defenses for it breaks their latent assumption
 172 that the backdoor attacks are targeted.

173 **Settings.** Since neural cleanse [18] is the first and the most representative trigger synthesis based
 174 defense, we adopt it as an example to synthesize the trigger pattern of DNNs watermarked by BadNets
 175 and our UBW-P and UBW-C. We implement it based on its open-sourced codes² and default settings.

176 **Results.** As shown in Figure 3, the synthesized pattern of BadNets is similar to the ground-truth
 177 trigger pattern. However, those of our UBW-P and UBW-C are significantly different from the
 178 ground-truth one. These results show that our UBW is resistant to trigger synthesis based defenses.

179 G.2 Resistance to Saliency-based Defenses

180 Since the attack effectiveness is mostly caused by the trigger pattern, there were also some backdoor
 181 defenses [21, 22, 23] based on detecting trigger areas with saliency maps. Specifically, these methods
 182 first generated the saliency map of each sample and then obtained trigger regions based on the
 183 intersection of all generated saliency maps. Since our UBW is untargeted, the relation between the
 184 trigger pattern and the predicted label is less significant compared with existing targeted backdoor
 185 attacks. As such, it can bypass those saliency-based defenses, which is verified in this section.

186 **Settings.** We generate the saliency maps of models watermarked by BadNets and our UBW-P and
 187 UBW-C, based on the Grad-CAM [24] with its default settings. We randomly select samples from
 188 the source class to generate their poisoned version for the discussion.

189 **Results.** As shown in Figure 4, the Grad-CAM mainly focuses on the trigger areas of poisoned
 190 images in BadNets. In contrast, it mainly focuses on other regions (*e.g.*, object outline) of poisoned
 191 images in our UBW-C. We notice that the Grad-CAM also focuses on the trigger areas in our UBW-P
 192 in a few cases. It is most probably because the trigger pattern used in the inference process is the

²<https://github.com/bolunwang/backdoor>

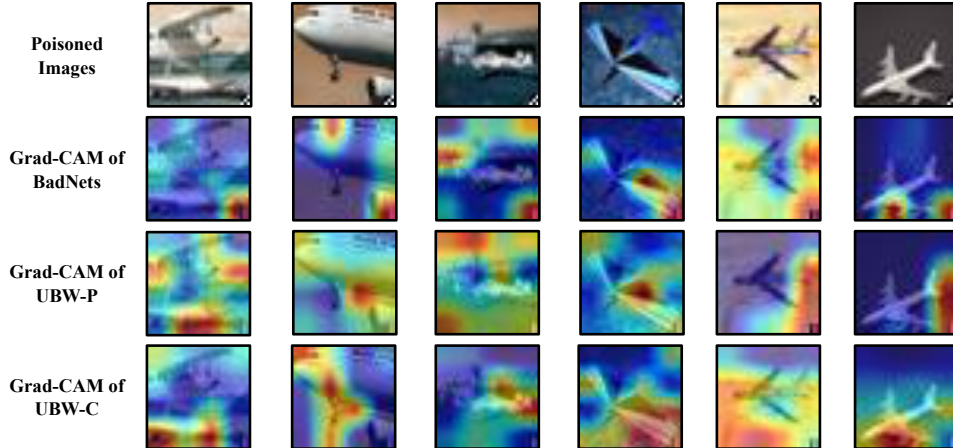


Figure 4: The poisoned images and their saliency maps based on Grad-CAM with DNNs watermarked by different methods. The Grad-CAM mainly focuses on the trigger areas of poisoned images in BadNets, while it mainly focuses on other regions (*e.g.*, object outline) in our UBW.

Table 6: The averaged entropy generated by STRIP of models watermarked by different methods. The larger the entropy, the harder for STRIP to detect the watermark.

Metric↓, Method→	BadNets	UBW-P	UBW-C
Averaged Entropy	0.0093	1.5417	1.2018

193 same as the one used for training in our UBW-P while we use invisible additive noises in our training
 194 process of UBW-C. These results validate that our UBW is resistant to saliency-based defenses.

195 G.3 Resistance to STRIP

196 Recently, Gao *et al.* [25] proposed STRIP to filter poisoned samples based on the prediction variation
 197 of samples generated by imposing various image patterns on the suspicious image. The variation is
 198 measured by the entropy of the average prediction of those samples. Specifically, the STRIP assumed
 199 that the trigger pattern is sample-agnostic and the attack is targeted. Accordingly, the more likely the
 200 suspicious image contains trigger pattern, the smaller the entropy since those modified images will
 201 still be predicted as the target label so that the average prediction is still nearly an one-hot vector.

202 **Settings.** We randomly select 100 testing images from the source class to generate their poisoned
 203 version, based on BadNets and our UBW-P and UBW-C. We calculate the entropy of each poisoned
 204 image based on the open-sourced codes³ and default settings of STRIP. We then calculate the averaged
 205 entropy among all poisoned samples for each watermarking method as their indicator. The larger the
 206 entropy, the harder for STRIP to detect the watermark.

207 **Results.** As shown in Table 6, the averaged entropies of both UBW-P and UBW-C are significantly
 208 higher than that of BadNets. Specifically, the entropies of both UBW-P and UBW-C are more than
 209 100 times larger than that of BadNets. It is mostly due to the untargeted nature of our UBW whose
 210 predictions are dispersed. These results verify that our UBW is resistant to STRIP.

211 G.4 Resistance to Dataset-level Backdoor Defenses

212 In this section, we discuss whether our methods are resistant to dataset-level backdoor defenses.

213 **Settings.** In this part, we adopt the spectral signatures [26] and the activation clustering [27] as
 214 representative dataset-level backdoor defenses for our discussion. Both spectral signatures and
 215 activation clustering tend to filter poisoned samples from the training dataset, based on sample
 216 behaviors in hidden feature space. We implement these methods based on their official open-sourced

³<https://github.com/yjkim721/STRIP-ViT>

Table 7: The successful filtering rate (%) on the CIFAR-10 dataset.

Method↓, Defense→	Spectral Signatures	Activation Clustering
UBW-P	10.96	52.61
UBW-C	9.40	20.51

Table 8: The resistance to MCR and NAD on the CIFAR-10 dataset.

Defense→ Method↓, Metric→	No Defense		MCR		NAD	
	BA (%)	ASR-A (%)	BA (%)	ASR-A (%)	BA (%)	ASR-A (%)
UBW-P	90.59	92.30	88.17	96.20	67.98	99.40
UBW-C	86.99	89.80	86.15	79.10	77.13	36.00

Table 9: The effectiveness of our UBW-P with different types of triggers on the CIFAR-10 dataset.

Method↓, Metric→	BA (%)	ASR-A (%)	ASR-C (%)	D_p
UBW-P (BadNets)	90.59	92.30	92.51	2.2548
UBW-P (WaNet)	89.90	73.00	70.45	2.0368

217 codes with default settings on the CIFAR-10 dataset. Besides, we adopt the *successful filtering rate*
 218 defined as the number of filtered poisoned samples over that of all filtered samples as our evaluation
 219 metric. In general, the smaller the successful filtering rate, the more resistance of our UBW.

220 **Results.** As shown in Table 7, these defenses fail to filter our watermarked samples under both
 221 poisoned-label and clean-label to some extent. We speculate that it is mostly because poisoned
 222 samples generated by our UBW-P and UBW-C tend to scatter in the whole space instead of forming
 223 a single cluster in the feature space. We will further explore it in the future.

224 G.5 Resistance to MCR and NAD

225 Here we discuss whether our methods are resistant to mode connectivity repairing (MCR) [28] and
 226 neural attention distillation (NAD) [29], which are two advanced repairing-based backdoor defenses.

227 **Settings.** We implement MCR and NAD based on the codes provided in BackdoorBox.

228 **Results.** As shown in Table 8, both our UBW-P and UBW-C are resistant to MCR and NAD to
 229 some extent. Their failures are probably because both of them contain a fine-tuning stage, which is
 230 ineffective for our UBWs (as demonstrated in Section 5.4.2).

231 H UBW-P with Imperceptible Trigger Patterns

232 In our main manuscript, we design our UBW-P based on BadNets-type triggers since it is the most
 233 straightforward method. We intend to show how simple it is to design UBW under the poisoned-label
 234 setting. Here we demonstrate that our UBW-P is still effective with imperceptible trigger patterns.

235 **Settings.** We adopt the advanced invisible targeted backdoor attack – WaNet [13] to design our
 236 UBW-P with imperceptible trigger patterns. We also implement it based on the open-sourced codes
 237 of vanilla WaNet provided in BackdoorBox [15]. Specifically, we set the warping kernel size as 16
 238 and conduct experiments on the CIFAR-10 dataset. Except for the trigger patterns, all other settings
 239 are the same as those used in our standard UBW-P.

240 **Results.** As shown in Table 9, our UBW-P can still reach promising performance with imperceptible
 241 trigger patterns, although it may have relatively low ASR compared to UBW-P with the BadNets-type
 242 visible trigger. It seems that there is a trade-off between ASR and trigger visibility. We will discuss
 243 how to better balance the watermark effectiveness and its stealthiness in our future work.

244 I The Transferability of our UBW-C

245 Recall that in the optimization process of our UBW-C, we need to know the model structure f in
 246 advance. Following the classical settings of bi-level-optimization-type backdoor attacks (*e.g.*, LIRA

Table 10: The performance of our UBW-C with different model structures trained on the watermarked CIFAR-10 dataset generated with ResNet-18.

Metric↓, Model→	ResNet-18	ResNet-34	VGG-16-BN	VGG-19-BN
BA (%)	86.99	87.34	86.83	88.55
ASR-A (%)	87.56	78.89	75.80	74.30

[6] and Sleeper Agent [4]), we report the results of attacking DNN with the same model structure as the one used for generating poisoned samples. In practice, dataset users may adopt different model structures since dataset owners have no information about the model training. In this section, we evaluate whether the watermarked dataset is still effective in watermarking DNNs having different structures compared to the one used for dataset generation (*i.e.*, transferability).

Settings. We adopt ResNet-18 to generate a UBW-C training dataset, based on which to train different models (*i.e.*, ResNet-18, ResNet-34, VGG-16-BN, and VGG-19-BN). Except for the model structure, all other settings are the same as those used in Section 5.

Results. As shown in Table 10, our UBW-C has high transferability. Accordingly, our methods are practical in protecting open-sourced datasets.

J Connections and Differences with Related Works

In this section, we discuss the connections and differences between our UBW and adversarial attacks, data poisoning, and classical untargeted attacks. We also discuss the connections and differences between our UBW-based dataset ownership verification and model ownership verification.

J.1 Connections and Differences with Adversarial Attacks

Both our UBW and adversarial attacks intend to make the DNNs misclassify samples during the inference process by adding malicious perturbations. However, they still have some intrinsic differences.

Firstly, the success of adversarial attacks is mostly due to the behavior differences between DNNs and humans, while that of our UBW results from the data-driven training paradigm and excessive learning ability of DNNs. Secondly, the malicious perturbations are known (*i.e.*, non-optimized) by UBW whereas adversarial attacks need to obtain them based on the optimization process. As such, adversarial attacks cannot to be real-time in many cases, since the optimization requires querying the DNNs multiple times under either white-box [30, 31, 32] or black-box [33, 34, 35] settings. Lastly, our UBW requires modifying the training samples without any additional requirements in the inference process, while adversarial attacks need to control the inference process to some extent.

J.2 Connections and Differences with Data Poisoning

Currently, there are two types of data poisoning, including classical data poisoning [36, 37, 38] and advanced data poisoning [39, 40, 41]. Specifically, the former type of data poisoning intends to reduce model generalization, so that the attacked DNNs behave well on training samples whereas having limited effectiveness in predicting testing samples. The latter one requires that the model has good benign accuracy while misclassifying some adversary-specified unmodified samples.

Our UBW shares some similarities to data poisoning in the training process. Specifically, they all intend to embed distinctive prediction behaviors in the DNNs by poisoning some training samples. However, they also have many essential differences. The detailed differences are as follows:

The Differences Compared with Classical Data Poisoning. Firstly, UBW has a different goal compared with classical data poisoning. Specifically, UBW preserves the accuracy in predicting benign testing samples whereas classical data poisoning is not. Secondly, UBW is also more stealthy compared with classical data poisoning, since dataset users can easily detect classical data poisoning by evaluating model performance on a local verification set. In contrast, this method has limited benefits in detecting UBW. Lastly, the effectiveness of classical data poisoning is mostly due to the sensitiveness of the training process, so that even a small domain shift of training samples may lead to significantly different decision surfaces of attacked models. It is different from that of our UBW.

289 **The Differences Compared with Advanced Data Poisoning.** Firstly, advanced data poisoning can
290 only misclassify a few selected images whereas UBW can lead to the misjudgments of all images
291 containing the trigger pattern. It is mostly due to their second difference that data poisoning does
292 not require modifying the (benign) images before the inference process. Thirdly, the effectiveness
293 of advanced data poisoning is mainly because DNNs are over-parameterized, so that the decision
294 surface can have sophisticated structures near the adversary-specified samples for misclassification.
295 It is also different from that of our UBW.

296 **J.3 Connections and Differences with Classical Untargeted Attacks**

297 Both our UBW and classical untargeted attacks (*e.g.*, untargeted adversarial attacks) intend to make
298 the model misclassify specific sample(s). However, different from existing classical untargeted attacks
299 which simply maximize the loss between the predictions of those samples and their ground-truth
300 labels, our UBW also requires optimizing the prediction dispersibility so that the adversaries cannot
301 deterministically manipulate model predictions. Maximizing only the untargeted loss may not be able
302 to disperse model predictions, since targeted attacks can also maximize that loss when the target label
303 is different from the ground-truth one of the sample. Besides, introducing prediction dispersibility
304 may also increase the difficulty of the untargeted attack since it may contradict the untargeted loss to
305 some extent (as described in Section B).

306 **J.4 Connections and Differences with Model Ownership Verification**

307 Our UBW-based dataset ownership verification enjoys some similarities to model ownership verifi-
308 cation [42, 43, 44] since they all conduct verification based on the distinctive behaviors of DNNs.
309 However, they still have many fundamental differences, as follows:

310 Firstly, dataset ownership verification has different threat models and requires different capacities.
311 Specifically, model ownership verification is adopted to protect the copyrights of open-sourced or
312 deployed models, while our method is for protecting dataset copyrights. Accordingly, our UBW-based
313 method only needs to modify the dataset, whereas model ownership verification usually also requires
314 controlling other training components (*e.g.*, loss). In other words, our UBW-based method can also
315 be exploited to protect model copyrights, whereas most of the existing methods for model ownership
316 verification are not capable to protect (open-sourced) datasets.

317 Secondly, to the best of our knowledge, almost all existing black-box model ownership verification
318 was designed based on the targeted attacks (*e.g.*, targeted poison-only backdoor attacks) and therefore
319 introducing new security risks in DNNs. In contrast, our verification method is mostly harmless,
320 since our UBW used for dataset watermarking is untargeted and with high prediction dispersibility.

321 **J.5 Connections and Differences with Radioactive Data**

322 We notice that radioactive data (RD) [45] (under the black-box setting) can also be exploited as
323 dataset watermarking for ownership verification by analyzing the loss of watermarked and benign
324 images. If the loss of watermarked images is significantly lower than that of their benign version,
325 RD treats the suspicious model as trained on the protected dataset. Both RD and UBW-C require
326 knowing the model structure in advance, although they all have transferability. However, they still
327 have many fundamental differences, as follows:

328 Firstly, our UBWs have a different verification mechanism compared to RD. Specifically, UBWs
329 adopt the change of predicted probability on the ground-truth label, while RD exploits the loss change
330 for verification. In practice, it is relatively difficult to select the confidence budget for RD since the
331 loss values may change significantly across different datasets. In contrast, users can easily select the
332 confidence budget (*i.e.*, τ) from $[0, 1]$ since the predicted probability on the ground-truth label are
333 relatively stable (*e.g.*, nearly 1 for benign samples).

334 Secondly, our UBWs require fewer defender capacities compared to RD. RD needs to have the
335 prediction vectors or even the model source files for ownership verification, whereas UBWs only
336 require the probability in the predicted label. Accordingly, our method can even be generalized to the
337 scenario that users can only obtain the predicted labels (as suggested in [46]), based on examining
338 whether poisoned images have different predictions compared to their benign version, whereas RD
339 cannot. We will further discuss the label-only UBW verification in our future work.

340 Lastly, it seems that RD is far less effective on datasets with relatively low image resolution and fewer
341 samples (*e.g.*, CIFAR-10)⁴. In contrast, our methods have promising performance on them.

342 **K Discussions about Adopted Data**

343 In this paper, all adopted samples are from the open-sourced datasets (*i.e.*, CIFAR-10 and ImageNet).
344 The ImageNet dataset may contain a few personal contents, such as human faces. However, our
345 research treats all objects the same and does not intentionally exploit or manipulate these contents.
346 Accordingly, our work fulfills the requirements of those datasets and should not be regarded as a
347 violation of personal privacy. Besides, our samples contain no offensive content, since we only add
348 some invisible noises or non-semantic patches to a few benign images.

349 **References**

- 350 [1] Solomon Kullback. *Information theory and statistics*. Courier Corporation, 1997.
- 351 [2] Risheng Liu, Jiaxin Gao, Jin Zhang, Deyu Meng, and Zhouchen Lin. Investigating bi-level
352 optimization for learning and vision from a unified perspective: A survey and beyond. *IEEE*
353 *Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- 354 [3] Aniruddha Saha, Akshayvarun Subramanya, and Hamed Pirsiavash. Hidden trigger backdoor
355 attacks. In *AAAI*, 2020.
- 356 [4] Hossein Souri, Micah Goldblum, Liam Fowl, Rama Chellappa, and Tom Goldstein. Sleeper
357 agent: Scalable hidden trigger backdoors for neural networks trained from scratch. In *NeurIPS*,
358 2022.
- 359 [5] Sebastian Ruder. An overview of gradient descent optimization algorithms. *arXiv preprint*
360 *arXiv:1609.04747*, 2016.
- 361 [6] Khoa Doan, Yingjie Lao, Weijie Zhao, and Ping Li. Lira: Learnable, imperceptible and robust
362 backdoor attacks. In *ICCV*, 2021.
- 363 [7] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 364 [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale
365 hierarchical image database. In *CVPR*, 2009.
- 366 [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
367 recognition. In *CVPR*, 2016.
- 368 [10] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as
369 an alternative to the cifar datasets. *arXiv preprint arXiv:1707.08819*, 2017.
- 370 [11] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating
371 backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019.
- 372 [12] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on
373 deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017.
- 374 [13] Anh Nguyen and Anh Tran. Wanet—imperceptible warping-based backdoor attack. In *ICLR*,
375 2021.
- 376 [14] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks.
377 *arXiv preprint arXiv:1912.02771*, 2019.
- 378 [15] Yiming Li, Mengxi Ya, Yang Bai, Yong Jiang, and Shu-Tao Xia. BackdoorBox: A python
379 toolbox for backdoor learning. 2022.
- 380 [16] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu.
381 Towards deep learning models resistant to adversarial attacks. In *ICLR*, 2018.

⁴https://github.com/facebookresearch/radioactive_data/issues/3

- 382 [17] Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural
383 networks. In *ICCV*, 2017.
- 384 [18] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and
385 Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks.
386 In *IEEE S&P*, 2019.
- 387 [19] Junfeng Guo, Ang Li, and Cong Liu. Aeva: Black-box backdoor detection using adversarial
388 extreme value analysis. In *ICLR*, 2022.
- 389 [20] Guanhong Tao, Guangyu Shen, Yingqi Liu, Shengwei An, Qiuling Xu, Shiqing Ma, Pan Li, and
390 Xiangyu Zhang. Better trigger inversion optimization in backdoor scanning. In *CVPR*, 2022.
- 391 [21] Xijie Huang, Moustafa Alzantot, and Mani Srivastava. Neuroninspect: Detecting backdoors in
392 neural networks via output explanations. *arXiv preprint arXiv:1911.07399*, 2019.
- 393 [22] Edward Chou, Florian Tramer, and Giancarlo Pellegrino. Sentinet: Detecting localized universal
394 attacks against deep learning systems. In *IEEE S&P Workshop*, 2020.
- 395 [23] Bao Gia Doan, Ehsan Abbasnejad, and Damith C. Ranasinghe. Februus: Input purification
396 defense against trojan attacks on deep neural network systems. In *ACSAC*, 2020.
- 397 [24] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi
398 Parikh, and Dhruv Batra. Grad-cam: Visual explanations from deep networks via gradient-based
399 localization. In *ICCV*, 2017.
- 400 [25] Yansong Gao, Yeonjae Kim, Bao Gia Doan, Zhi Zhang, Gongxuan Zhang, Surya Nepal, Damith
401 Ranasinghe, and Hyoungshick Kim. Design and evaluation of a multi-domain trojan detection
402 method on deep neural networks. *IEEE Transactions on Dependable and Secure Computing*,
403 2021.
- 404 [26] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. In
405 *NeurIPS*, 2018.
- 406 [27] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung
407 Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by
408 activation clustering. In *AAAI Workshop*, 2019.
- 409 [28] Pu Zhao, Pin-Yu Chen, Payel Das, Karthikeyan Natesan Ramamurthy, and Xue Lin. Bridging
410 mode connectivity in loss landscapes and adversarial robustness. In *ICLR*, 2020.
- 411 [29] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention
412 distillation: Erasing backdoor triggers from deep neural networks. In *ICLR*, 2021.
- 413 [30] Jiawang Bai, Bin Chen, Yiming Li, Dongxian Wu, Weiwei Guo, Shu-tao Xia, and En-hui Yang.
414 Targeted attack for deep hashing based retrieval. In *ECCV*, 2020.
- 415 [31] Francesco Croce and Matthias Hein. Reliable evaluation of adversarial robustness with an
416 ensemble of diverse parameter-free attacks. In *ICML*, 2020.
- 417 [32] Bin Chen, Yan Feng, Tao Dai, Jiawang Bai, Yong Jiang, Shu-Tao Xia, and Xuan Wang.
418 Adversarial examples generation for deep product quantization networks on image retrieval.
419 *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- 420 [33] Pin-Yu Chen, Huan Zhang, Yash Sharma, Jinfeng Yi, and Cho-Jui Hsieh. Zoo: Zeroth order
421 optimization based black-box attacks to deep neural networks without training substitute models.
422 In *CCS Workshop*, 2017.
- 423 [34] Maksym Andriushchenko, Francesco Croce, Nicolas Flammarion, and Matthias Hein. Square
424 attack: a query-efficient black-box adversarial attack via random search. In *ECCV*, 2020.
- 425 [35] Yan Feng, Baoyuan Wu, Yanbo Fan, Li Liu, Zhifeng Li, and Shutao Xia. Boosting black-box
426 attack with partially transferred conditional adversarial distribution. In *CVPR*, 2022.

- 427 [36] Huang Xiao, Battista Biggio, Gavin Brown, Giorgio Fumera, Claudia Eckert, and Fabio Roli. Is
428 feature selection secure against training data poisoning? In *ICML*, 2015.
- 429 [37] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions.
430 In *ICML*, 2017.
- 431 [38] Ji Feng, Qi-Zhi Cai, and Zhi-Hua Zhou. Learning to confuse: generating training time adversar-
432 ial data with auto-encoder. In *NeurIPS*, 2019.
- 433 [39] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suci, Christoph Studer, Tudor Du-
434 mitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural
435 networks. In *NeurIPS*, 2018.
- 436 [40] Jonas Geiping, Liam Fowl, W Ronny Huang, Wojciech Czaja, Gavin Taylor, Michael Moeller,
437 and Tom Goldstein. Witches’ brew: Industrial scale data poisoning via gradient matching. In
438 *ICLR*, 2021.
- 439 [41] Avi Schwarzschild, Micah Goldblum, Arjun Gupta, John P Dickerson, and Tom Goldstein. Just
440 how toxic is data poisoning? a unified benchmark for backdoor and data poisoning attacks. In
441 *ICML*, 2021.
- 442 [42] Yossi Adi, Carsten Baum, Moustapha Cisse, Benny Pinkas, and Joseph Keshet. Turning your
443 weakness into a strength: Watermarking deep neural networks by backdooring. In *USENIX*
444 *Security*, 2018.
- 445 [43] Hengrui Jia, Christopher A Choquette-Choo, Varun Chandrasekaran, and Nicolas Papernot.
446 Entangled watermarks as a defense against model extraction. In *USENIX Security*, 2021.
- 447 [44] Yiming Li, Linghui Zhu, Xiaojun Jia, Yong Jiang, Shu-Tao Xia, and Xiaochun Cao. Defending
448 against model stealing via verifying embedded external features. In *AAAI*, 2022.
- 449 [45] Alexandre Sablayrolles, Matthijs Douze, Cordelia Schmid, and Hervé Jégou. Radioactive data:
450 tracing through training. In *ICML*, 2020.
- 451 [46] Yiming Li, Mingyan Zhu, Xue Yang, Yong Jiang, and Shu-Tao Xia. Black-box ownership
452 verification for dataset protection via backdoor watermarking. *arXiv preprint arXiv:2209.06015*,
453 2022.