

Complex Backdoor Detection by Symmetric Feature Differencing

Yingqi Liu^{1*}, Guangyu Shen^{1*}, Guanhong Tao¹, Zhenting Wang², Shiqing Ma², Xiangyu Zhang¹
Purdue University¹, Rutgers University²

{liu1751, shen447, taog}@purdue.edu, {zhenting.wang, sm2283}@rutgers.edu, xyzhang@cs.purdue.edu

Abstract

Many existing backdoor scanners work by finding a small and fixed trigger. However, advanced attacks have large and pervasive triggers, rendering existing scanners less effective. We develop a new detection method. It first uses a trigger inversion technique to generate triggers, namely, universal input patterns flipping victim class samples to a target class. It then checks if any such trigger is composed of features that are not natural distinctive features between the victim and target classes. It is based on a novel symmetric feature differencing method that identifies features separating two sets of samples (e.g., from two respective classes). We evaluate the technique on a number of advanced attacks including composite attack, reflection attack, hidden attack, filter attack, and also on the traditional patch attack. The evaluation is on thousands of models, including both clean and trojaned models, with various architectures. We compare with three state-of-the-art scanners. Our technique can achieve 80-88% accuracy while the baselines can only achieve 50-70% on complex attacks. Our results on the TrojAI competition rounds 2-4, which have patch backdoors and filter backdoors, show that existing scanners may produce hundreds of false positives (i.e., clean models recognized as trojaned), while our technique removes 78-100% of them with a small increase of false negatives by 0-30%, leading to 17-41% overall accuracy improvement. This allows us to achieve top performance on the leaderboard.

1. Introduction

Backdoor attack (or trojan attack) on deep learning models injects malicious behaviors such that a compromised model behaves normally on clean inputs and misclassifies inputs stamped with a *trigger* to a *target label* [6, 20, 36, 40, 41]. It becomes a prominent threat due to the low complexity of launching such attacks, the devastating consequences especially in safety/security critical applica-

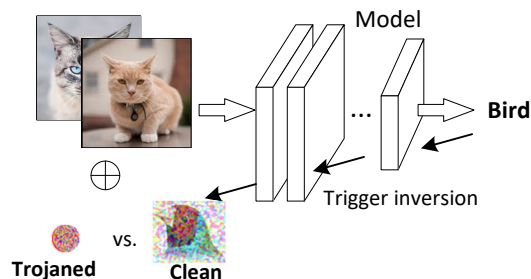


Figure 1. Backdoor detection by trigger inversion

tions, and the difficulty of defense.

There are a body of existing defense techniques such as trigger inversion [39, 66], attribution analysis [16, 25], trojaned input detection [14, 64], backdoor removal [34, 37], and model certification [65]. According to [3, 48], trigger inversion is an effective technique that can determine if a model contains backdoor without assuming the availability of any input with trigger. For example, Neural Cleanse (NC) [66], Artificial Brain Stimulation (ABS) [39], and K-Arm [55] make use of optimization to invert triggers and determine if a model is trojaned. They consider each label in the model as a potential target and use optimization to check if a small and fixed input pattern, i.e., a trigger, can be found to cause any input to be misclassified to the label. The intuition is that attackers tend to use small triggers for attack stealthiness. Figure 1 illustrates trigger inversion. An input pattern (the circular pattern or the rectangular one on the bottom) is generated by gradient back-propagation to flip cat samples to bird. If the subject model is clean, a *large* pattern that exhibits a lot of bird features is generated (e.g., the rectangular pattern with the “clean” tag). In contrast, when the model is trojaned (with a red circular patch), a *small* pattern containing the trigger features is inverted (e.g., the circular pattern with the “trojaned” tag). The size difference of the patterns is critical for these scanners, that is, a model is flagged as trojaned only when a small trigger can be found. Observe that inverted patterns are usually noisy and may not be human interpretable.

While existing techniques are effective for attacks with small and static triggers, advanced attacks proposed re-

*Equal contribution

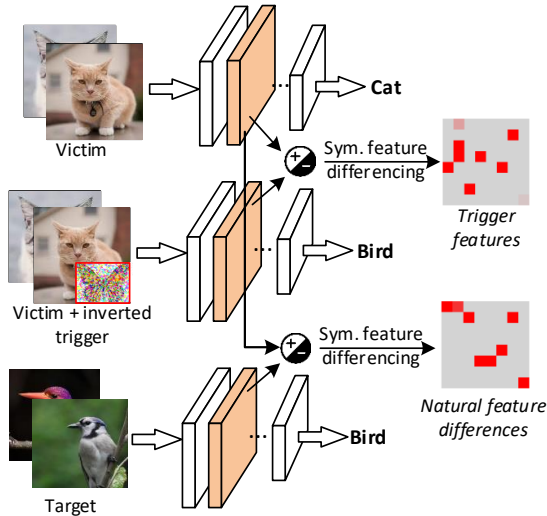


Figure 2. Ex-ray overview

cently [36, 41, 52] have large and dynamic triggers: the input level differences before and after injecting a trigger are substantial and the differences vary across different inputs. Composite attack [36] injects a backdoor by mixing benign features from two or more classes. For example, a butterfly appearing in a cat image causes the model to predict bird. Triggers may be large (e.g., a butterfly could be much larger than a typical patch trigger) and have different pixel level manifestations (e.g., various kinds of butterfly). Reflection backdoor [41] uses the reflection of an image as the trigger. Reflection occurs when pictures are taken behind a glass window. Reflection could be as large as a whole image. Hidden trigger backdoor [52] introduces perturbations on the training images of target label such that the perturbed images have similar inner activations to the trigger and forces the model to learn the correlations between the trigger and the target label. Since the trigger is not explicit in the training, the trojanning process is more difficult and requires larger triggers. Existing scanners such as NC, ABS, and K-arm have difficulties detecting these backdoors. They only achieve 0.5-0.7 accuracy (see the evaluation section).

Essence of Backdoor Attack. We observe the essence of backdoor attack is that *model misclassification (to the target class) is induced by features not in the target class*. For example, in a composite attack, a cat image is misclassified to bird when a butterfly is also present in the image. The misclassification is essentially induced by the features of cat and butterfly (not bird’s). Our overarching idea is hence to determine if the features of an inverted trigger are the natural features distinguishing the victim and target classes. If so, the model is clean. Otherwise, it is considered trojaned. Note that in our method small sizes and fixed triggers are no longer essential properties. One may argue that the attacker could craft a trigger with the target class’s natural features. We will discuss such an adaptive attack in Section 4.3.

Our Method. Figure 2 illustrates our method. It takes a model and a small set of clean samples (e.g., 10 for each class). It first leverages an existing trigger inversion method to derive a trigger that can flip a set of clean samples of the victim class (cat) to the target class (bird). It feeds a set of clean victim samples to the model, and extracts the internal feature maps at a selected layer (the first row in the figure). It then injects the inverted trigger (the butterfly-like pattern) to the clean victim samples and extracts the corresponding feature maps (the second row). A novel feature comparison technique called *symmetric feature differencing* (SFD) is then applied to the two sets of feature maps, to determine the distinctive features between the two sets of samples (the first rectangular map on the right with red cells), which essentially denotes the trigger features. The map is also called a *feature difference mask* or just *mask*. A red cell in the mask indicates a feature map that is distinctive. It further feeds a set of clean target class (bird) samples to the model and extracts the feature maps (the third row). Applying SFD to the target class and the victim class feature maps yields the natural features distinguishing the two classes, that is, the second mask on the right. A model is considered trojaned when the two masks do not have similarity.

The key enabling technique of our method is SFD, which is a novel differential analysis. It is based on *counter-factual causality* [33]. Given two sets of samples, like the victim and victim+trigger samples mentioned above, it computes a smallest set of feature maps such that exchanging their activation values across the two sets entails exchanged classification results. They are considered the distinctive features. The formal definition and the computation algorithm can be found in Section 3.

Our contributions are summarized as follows.

- We develop a new scanning technique that can detect large and complex backdoors which are difficult for existing techniques.
- The technique is based on a novel symmetric feature differencing method that can identify the distinctive features of two sets of given examples.
- We implement a prototype EX-RAY (“*DEtecting Complex Backdoor in Neural Networks by SYmmetric Feature Differencing*”). It is general and can leverage different upstream trigger inversion methods. EX-RAY is publicly available at <https://github.com/PurduePAML/Exray>
- We evaluate EX-RAY on 4246 models (2081 benign and 2165 trojaned) with 23 structures and 7 datasets, and four attacks that have large/pervasive and dynamic triggers (reflection, composite, hidden, and filter attacks). We compare with three state-of-the-art trigger inversion based scanners, NC, ABS, and K-Arm. Our results show that EX-RAY can achieve 80-88% accuracy while the baselines can only achieve 50-70%.

We also use model interpretation techniques to show that EX-RAY indeed captures natural feature differences between classes. EX-RAY can also be used to remove false positives in backdoor scanning (i.e., clean models are considered trojaned), which are usually due to small triggers found between clean labels. EX-RAY can determine that such triggers essentially denote natural features of the target label and should be precluded. We test EX-RAY (with ABS as the upstream inversion technique) on TrojAI¹ rounds 2-4 and show that EX-RAY can reduce false warnings by 78-100%, with the cost of a small increase in false negatives (0-30%), i.e., trojaned models are considered clean. It can improve multiple upstream scanners' overall accuracy including ABS (by 17-41%), NC (by 25%), and the Bottom-up-Top-down backdoor scanner [1] (by 2-15%) in the competition. Our method also outperforms a number of other false positive removal methods that compare L2 distances and leverage attribution/interpretation techniques. EX-RAY will be released upon publication.

- On the TrojAI leaderboard, ABS+EX-RAY achieves top performance in 2 out of the 4 rounds for image classification, including the most challenging round 4, with an average cross-entropy (CE) loss around 0.32² and an average AUC-ROC³ around 0.90. It is the only technique that successfully reached the round target (for both the training sets and the test sets remotely evaluated by IARPA), i.e., a CE loss lower than 0.3465, for all the 4 rounds. As far as we know, a large number of state-of-the-art scanning techniques have been evaluated in the competition, including NC [66], ABS [39], Meta neural analysis [73], ULP [29], DeepInspect [11], SCAn [60], K-Arm backdoor scanning [55], Noise analysis backdoor detection [16] and attribution based backdoor detection [25, 57].

Threat Model. Our threat model is consistent with that in existing works [3, 39]. Given a set of models, including both trojaned and clean models, and a small set of clean samples for each model (covering all labels), we aim to identify the models with injected backdoor(s) that can flip clean samples to the target class. These samples may belong to one or many victim class(es). The former is label-specific attack and the latter is universal attack. □

¹TrojAI is a backdoor scanning competition organized by IARPA [3]. Rounds 1-4 are for image classification. Round 1 dataset is excluded due to simplicity.

²The smaller the better.

³An accuracy metric used by TrojAI, the larger the better.

2. Related Work

Backdoor Attack. Data poisoning [12, 20] injects backdoors by changing the label of inputs with trigger. Clean label attack [52, 54, 63, 75] injects backdoors without changing the data label. Dynamic backdoor [46, 53] focuses on crafting different triggers for different inputs and breaks the defense's assumption that trigger is static. [47] proposes to combine adversarial example generation and model poisoning. There are also attacks on NLP tasks [13, 31, 74], reinforcement learning [28, 68], and federated learning [7, 17, 61, 67, 72]. EX-RAY is a general primitive that may be of use in defending these attacks.

Defenses against Backdoor Attacks. ULP [29] and Meta neural analysis [73] train a few input patterns and a classifier from thousands of benign and trojaned models. The classifier predicts if a model has backdoor based on activations of the patterns. [51] proposes to reverse engineer the distribution of triggers. [23] finds that trojaned and clean models react differently to input perturbations. TABOR [21] and NeuronInspect [24] use an AI explanation technique to detect backdoor. There are techniques that defend backdoors by data sanitization [9, 50]. There are also techniques that detect if a given input is stamped with a trigger [10, 14, 15, 18, 19, 22, 35, 38, 42, 60, 62, 64]. They target a different problem as they require inputs with embedded triggers. EX-RAY is orthogonal to most of these techniques and can serve as a performance booster.

Interpretation/Attribution. EX-RAY is related to model interpretation and attribution, e.g., important features identification [5, 8, 56, 59]. [26] measures the importance of a concept (e.g., 'striped') for a class (e.g., zebra). The differences lie in that EX-RAY finds distinguishing features of two sets of examples.

3. Design

As illustrated by Figure 2, given a trigger t inverted by some upstream scanning technique (not our contribution) that flips victim class V samples to the target class T , EX-RAY first computes the distinctive features between V and $V+t$ samples, then the distinctive features between V and T . Finally, it uses a similarity analysis to compare the two sets of distinctive features to determine if the trigger denotes natural differences between the two classes. If not, the model is considered trojaned. In this section, we explain the steps in details.

3.1. Symmetric Feature Differencing

The key enabling technique of EX-RAY is *symmetric feature differencing* (SFD) that determines the distinguishing features between two sets of examples (e.g., from classes V and T). SFD is based on *counter-factual causality* [33], which states that an effect event e causally depends

on a cause event c if and only if, 1) if c were to occur e would occur; and 2) if c were not to occur e would not occur. In our context, we say a set of features are distinctive between two sets of examples if and only if 1) *exchanging these features across the two sets (event c) entails exchanged classification results (event e)*, and 2) *the exchange of any such feature is necessary to the exchanged classification results*. For example, we say two sets of examples from two respective persons A and B in a face recognition model differ only at their nose if and only if 1) replacing the nose in the examples of A with B’s nose causes the model to predict B, 2) replacing the nose is needed to cause misclassification, and vice versa. Note that although replacing both nose and mouth can also induce exchanged classifications, replacing mouth is not necessary. Hence mouth is not a distinctive feature. Automatically identifying such feature differences in the input space is challenging due to the difficulty of locating features, as a feature may manifest itself differently across input examples. Our differencing method hence identifies a set of neurons (i.e., feature maps) at some hidden layer that denote the distinctive features.

We formally define symmetric feature differencing in the following. To simplify our discussion, we assume the technique takes a subject model $F(x)$ and two inputs: x_v in V and x_t in T (instead of two sets of inputs). We will discuss the extension to two sets later in the section.

Definition 1 (Symmetric Feature Differencing)

Let $F(x)$ be a feed forward neural network. Given an inner layer l that provides good feature abstraction, let g be the submodel up to layer l and h the submodel after l , i.e., $F(x) = h(g(x))$. Let the number of features/neurons at l be n . Symmetric feature differencing (SFD) computes a mask M that is an n element vector with values 0 or 1. Let $\neg M$ be the negation of the mask such that $\neg M[i] = 1 - M[i]$ with $i \in [1, n]$.

The mask M satisfies the following conditions.

$$h(g(x_v) \cdot M + g(x_t) \cdot \neg M) = V \tag{1}$$

$$h(g(x_v) \cdot \neg M + g(x_t) \cdot M) = T \tag{2}$$

$$\|M\|_0 \text{ is minimal.} \tag{3}$$

Intuitively, Equation (1) denotes that copying x_v ’s features to x_t with the control of M causes the classification of V . Specifically, $g(x_v) \cdot M + g(x_t) \cdot \neg M$ means that when $M[i] = 0$ (i.e., $\neg M[i] = 1$), the original i th feature map of the T sample x_t is retained; when $M[i] = 1$, the i th feature map of x_t is replaced with that from the V sample x_v . The explanation for Equation (2) is similar. Equation (3) dictates the minimality of M , that is, any feature exchange

indicated by the mask is necessary, faithfully following the counter-factual causality definition.

The SFD definition is graphically illustrated by an example in Figure 3. The dash box on the left shows the $g(x)$ function and that on the right the $h(\cdot)$ function. The top row in the left box shows that five feature maps (in yellow) are generated by $g(\cdot)$ for a victim class sample x_v . The bottom row shows that the five feature maps (in blue) for a target class sample x_t . The dash box in the middle illustrates symmetric differencing. As suggested by the red entries in the mask M in the middle (i.e., $M[3] = M[5] = 1$), in the top row, the 3rd and 5th (yellow) feature maps are replaced with the corresponding (blue) feature maps from the bottom. Symmetric replacements happen in the bottom row as well. On the right, the exchanged feature maps cause the exchanged classification results.

Note that a minimal M must exist. In the worst case, M is filled with ‘1’, indicating all feature maps shall be exchanged, which must yield the exchanged classification results. In general, the complexity of computing M is exponential. We hence propose a soft version of SFD.

Soft Symmetric Feature Differencing. In the soft version, we relax the meaning of mask. Instead of having an either 0 or 1 value, we allow the value to vary in $[0, 1]$, with 0 meaning no-exchange at all, 1 meaning complete exchange, and a value in between 0 and 1 partial exchange. For example, assume nose, eyes, and mouth are the three features in a face recognition model and assume $M[nose] = 1$, $M[eyes] = 0$, and $M[mouth] = 0.5$. The mask means that exchanging noses, retaining eyes, and mixing mouths half-half. The need of mixing mouths means that sometimes exchanging noses alone may not be sufficient to cause exchanged results and we need to partially consider their mouths which have a certain level of difference as well.

With the relaxation, the exchange operations, i.e., Equations (1) and (2), become continuous and differentiable. In addition, the minimality requirement in Equation (3) can be modeled by a differentiable arg min operation that minimizes the size of mask. Specifically, it can be reduced to the following constrained optimization problem.

$$\begin{aligned} & \arg \min_M \text{sum}(M), \text{ s.t.} \\ & h(g(x_v) \cdot M + g(x_t) \cdot \neg M) = V \text{ and} \\ & h(g(x_v) \cdot \neg M + g(x_t) \cdot M) = T \end{aligned} \tag{4}$$

To solve the problem, we devise a loss in (5). It has three parts. The first part $\text{sum}(M)$ is to minimize the mask size. The second part $w_1 \times ce_1$ is a barrier loss for Equation (1), with ce_1 the cross entropy loss when replacing x_t ’s features. Coefficients w_1 is dynamic. When the cross entropy loss is larger than a threshold α , w_1 is set to a large value w_{large} . This forces M to satisfy Equation (1). When the loss is small indicating the constraint is satisfied, w_1 is changed to

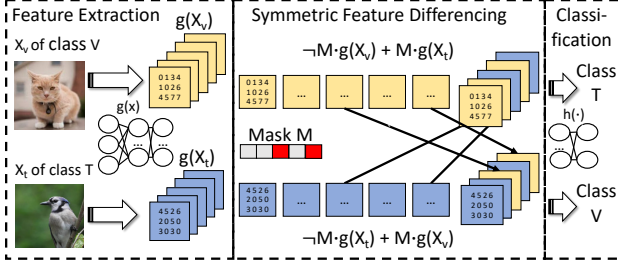


Figure 3. Illustrating symmetric feature differencing

a small value w_{small} . The optimization hence focuses on minimizing the mask. The third part $w_2 \times ce_2$ is similar.

$$\begin{aligned}
 \mathcal{L}_{pair}(x_v, x_t) &= \text{sum}(M) + w_1 \times ce_1 + w_2 \times ce_2, \\
 \text{with } ce_1 &= CE(h(g(x_v) \cdot M + g(x_t) \cdot \neg M), V), \\
 ce_2 &= CE(h(g(x_v) \cdot \neg M + g(x_t) \cdot M), T), \quad (5) \\
 w_1 &= w_{large} \text{ if } ce_1 > \alpha \text{ else } w_{small}, \\
 w_2 &= w_{large} \text{ if } ce_2 > \alpha \text{ else } w_{small}
 \end{aligned}$$

Differencing Two Sets. The algorithm to identify the differential features of two sets can be built from the primitive of comparing two inputs. Given two sets X_V of class V and X_T of class T , ideally the mask M should satisfy Equations (1) and (2) for any $x_v \in X_V$ and $x_t \in X_T$. While such a mask must exist (with the worst case containing all the features), minimizing it becomes very costly. Assume $|X_V| = |X_T| = m$. The number of constraints that need to be satisfied during optimization is $O(m^2)$. Therefore, we develop a stochastic method that is $O(m)$. Specifically, let \vec{X}_V and \vec{X}_T be random orders of X_V and X_T , respectively. We minimize M such that it satisfies Equations (1) and (2) for all pairs $(\vec{X}_V[j], \vec{X}_T[j])$, with $j \in [1, m]$. Intuitively, we optimize on a set of random pairs from X_V and X_T that cover all the individual samples in X_V and X_T . The loss function is hence the following.

$$\mathcal{L} = \sum_{j=1}^m \mathcal{L}_{pair}(\vec{X}_V[j], \vec{X}_T[j])$$

When X_V and X_T have one-to-one mapping, such as the victim class samples and their compromised versions that have the trigger injected, we can directly use the mapping in optimization instead of a random mapping. We use Adam optimizer [27] with a learning rate 5e-2 and 400 epochs. Masks are initialized to all 1 to begin with. This denotes a conservative start since such masks suggest swapping all feature maps, which must induce the intended classification results swap.

Symmetry Is Necessary. Our technique is symmetric. Such symmetry is critical to effectiveness. One may wonder that a one-sided analysis that only enforces Equation (2) may

be sufficient. That is, M is the minimal set of features that when copied from T (target) samples to V (victim) samples can flip the V samples to class T . Intuitively, it denotes the strong features of T . However, this is insufficient. In many cases, misclassification (of a V sample to T) in a clean model can be induced when strong features of class V are suppressed (instead of adding strong T features). As such, an inversion method may generate a trigger that neutralizes V features instead of injecting unique T features. The trigger features hence do not share much commonality with the features computed by the one-sided analysis. Consequently, the clean model is considered trojaned. Our experiments in Appendix A and H show the importance of symmetry.

3.2. Comparing Masks

After generating the masks, we compare them in the last step. Let the distinguishing features between the victim and target classes be M_1 , and then those between the victim samples and their compromised versions be M_2 . Next, we explain how to compare M_1 and M_2 . Intuitively M_1 and M_2 should share a lot of commonality when the trigger denotes natural differences between classes, as reflected in the following condition.

$$\text{sum}(\min(M_1, M_2)) > \beta \times \min(\text{sum}(M_1), \text{sum}(M_2)) \quad (6)$$

Here, $\min(M_1, M_2)$ yields a vector whose elements are the minimal between the corresponding elements in M_1 and M_2 . It essentially represents the intersection of the two masks. The hyperparameter β is in $(0, 1)$. Intuitively, the condition asserts that the size of mask intersection is larger than β times the minimum size of the two masks, meaning the two have a large part in common. If all the inverted triggers for a model satisfy the condition, the model is considered clean, otherwise trojaned.

Additional Validation Checks. In practice, due to the uncertainty in the stochastic symmetric differencing algorithm, the presence of local minimums in optimization, and the small number of available clean samples, M_1 and M_2 may not have a lot in common. However, they should nonetheless satisfy the semantic constraint that both should denote natural feature differences of the victim and target classes if the trigger is not injected. Therefore, we propose an additional cross-validation check that tests if functionally M_1 and M_2 can take each other's place in satisfying Equations (1) and (2). In particular, while M_1 is derived by comparing the victim class and the target class clean samples, we copy the feature maps indicated by M_1 between the victim samples and their compromised versions with trigger and check if the intended class flipping can be induced; similarly, while M_2 is derived by comparing the victim class samples and their compromised versions, we copy the feature maps indicated by M_2 between the victim clean sam-

ples and the target clean samples to see if the intended class flipping can be induced. If so, the two are functionally similar and the trigger is natural. The check is formulated as follows.

$$\begin{aligned}
 & Acc(h(g(X_V) \cdot M_2 + g(X_T) \cdot \neg M_2), V) > \gamma \wedge \\
 & Acc(h(g(X_T) \cdot M_2 + g(X_V) \cdot \neg M_2), T) > \gamma \wedge \\
 & Acc(h(g(X_V) \cdot M_1 + g(X_V + t) \cdot \neg M_1), V) > \gamma \wedge \\
 & Acc(h(g(X_V + t) \cdot M_1 + g(X_V) \cdot \neg M_1), T) > \gamma
 \end{aligned} \tag{7}$$

Here, $Acc()$ is a function to evaluate prediction accuracy on a set of samples and γ a threshold (0.8 in the paper). We use $g(X_V)$ to denote applying g on each sample in X_V for representation simplicity.

4. Evaluation

We evaluate EX-RAY on a range of backdoor attacks including four complex backdoors (i.e., composite attack, reflection attack, hidden attack, and filter attack), and on the traditional patch attack. We also study the applicability of EX-RAY to various upstream scanners by boosting their backdoor detection performance. We demonstrate that EX-RAY can be leveraged to fix models with injected and natural backdoors. In addition, we validate that the generated masks by EX-RAY is capable of capturing feature differences. We devise two adaptive attacks targeting the basis of EX-RAY. EX-RAY is implemented in PyTorch [49] and will be released upon publication.

Experiment Setup. The experiments are conducted on 4,246 models in total, with 200 models for composite attack, 148 models for reflection attack, 34 models for hidden attack, 1920 models for filter attack, and 1944 models for patch attack. For composite attack, we generate 100 trojaned models on five datasets (i.e., MNIST [32], Fashion MNIST [71], SVHN [44], CIFAR10 [30], and the Youtube Face dataset [69]) using the official implementation [36]. We follow [39] to create 100 clean models (20 models for each dataset). Network in Network and VGG16 are used for these models. For reflection attack, there are three different reflection settings, i.e., same depth of field, out of focus, and ghost effect. For each setting, we generate 20 trojaned models on CIFAR10 and 17 trojaned models on ImageNet using the official repository [41]. We also obtain 20 clean models on CIFAR10 from [39] and 17 clean models on ImageNet from torchvision [2]. We employ a few model structures such as Network in Network, VGG, ResNet, SqueezeNet, and DenseNet in this experiment. For hidden attack, we leverage 34 models on ImageNet with half clean from [2] and half trojaned by [52]. The model structures used include VGG, ResNet, SqueezeNet, and DenseNet. For filter and patch attacks, we make use of the TrojAI datasets from rounds 2 to 4, consisting of 3,840 models with half clean

Table 1. EX-RAY on composite attack

	ABS					ABS+EX-RAY				
	TP	FP	FN	TN	Acc/ROC	TP	FP	FN	TN	Acc/ROC
MNIST	16	12	4	8	0.6	18	3	2	17	0.88
FMNIST	12	9	8	11	0.58	19	6	1	14	0.83
SVHN	15	7	5	13	0.7	19	4	1	16	0.88
CIFAR10	16	13	4	7	0.58	17	3	3	17	0.85
Youtube face	12	4	8	16	0.7	19	5	1	15	0.85

and half trojaned. We also evaluate on 24 models on ImageNet. Details can be found in Appendix F. In addition, we study the hyper-parameters used in EX-RAY in Appendix G and Appendix I.

4.1. Detecting Complex Backdoor Attacks

In this section, we study the performance of EX-RAY in detecting three advanced backdoor attacks, namely, composite attack [36], reflection attack [41], and hidden attack [52], in comparison with a state-of-the-art technique ABS [39].

Detecting Composite Attack. Table 1 shows the detection results on composite attack. The first column denotes the datasets. Columns 2-6 show the detection results by ABS. Columns 7-11 present the detection results using ABS with EX-RAY. Columns TP, FP, FN, TN, and Acc/ROC denote the number of true positives, false positives, false negatives, true negatives, detection accuracy and ROC-AUC, respectively. The upstream methods output a binary result denoting whether a model is trojaned or not, and the clean and trojaned models are evenly distributed. Thus the ROC-AUC is the same as the accuracy. For ABS, we use the best possible bound for the inverted trigger size during detection. For ABS+EX-RAY, we set the bound for the trigger size to be half of the input. Observe that ABS+EX-RAY can achieve 83%-88% detection accuracy, substantially outperforming

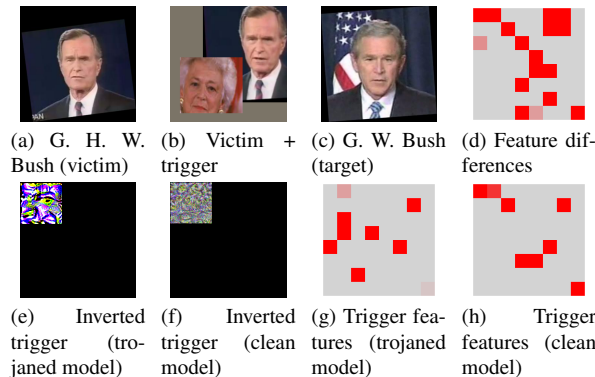


Figure 4. Composite attack: George H. W. Bush + Barbara Bush \rightarrow George W. Bush, and natural feature differences versus trigger differences. The maps in (d), (g), and (h) denote the neurons in a hidden layer with a red square denoting a distinctive neuron and its color the distinctive capability.

the original ABS, which has only 58%-70% accuracy. Note that EX-RAY reduces not only FPs, but also FNs. The reason for the latter is that ABS and other scanners like NC and K-arm are based on trigger size, while there is not a good separation by size. In contrast, EX-RAY is based on feature differencing.

Figure 4 (a-c) present a composite attack to a face recognition model, in which the presence of Barbara Bush in George H. W. Bush’s images flips the classification results to George W. Bush. EX-RAY identifies the trigger features in Figure 4 (g) by comparing George H. W. Bush + trigger and George H. W. Bush. Observe that they are quite different from the natural feature differences between George H. W. Bush and George W. Bush in Figure 4 (d) that distinguish the clean examples from the classes. In contrast, when a clean model is scanned, a trigger is inverted (Figure 4 (f)) to flip George H. W. Bush + trigger to George W. Bush. Observe that it is equally uninterpretable as that in Figure 4 (e), the inverted trigger for the trojaned model. This is due to the inherent limitation that trigger inversion can hardly generate natural-looking input features but rather noise-like patterns. Therefore, it is difficult to perform feature differencing at the input level. EX-RAY, however, produces a set of trigger features (Figure 4 (h)) that have substantial commonality with the natural feature differences (Figure 4 (d)), that is, (h) is a subset of (d). In other words, the noise-looking trigger in Figure 4 (f) indeed denotes natural differences. This indicates the model is benign.

Detecting Reflection Attack. Table 2 presents the results for reflection attack. Column 1 denotes the datasets. Column 2 shows the three attack settings. Columns 3-7 present the results of ABS and the remaining columns ABS+EX-RAY. For ABS, we use the best possible bound for the inverted trigger size during detection. For ABS+EX-RAY, we set the bound for the trigger size to be 25% of the input. The stability of EX-RAY regarding trigger size bound can be found in Appendix I. Observe that our technique can achieve 80%-85% accuracy, whereas ABS only has 55%-68% accuracy.

Figure 5 (a) shows a triangle sign used as a trigger to flip images (Figure 5 (b)) to airplane (Figure 5 (f)) in a trojaned model. Figure 5 (e) shows a trigger generated by ABS for the airplane label in the trojaned model. Observe that the generated trigger has (triangle) features of the real trigger as in Figure 5 (a-b). EX-RAY determines the model is a true positive as the inverted trigger shares very few features with airplane. In contrast, Figure 5 (c) presents a trigger generated by ABS for a clean model with deer as the target label (Figure 5 (d)). Observe that the trigger resembles deer antlers. The triggers inverted for other labels also have a similar nature. EX-RAY hence recognizes the model as a true negative.

Detecting Hidden-trigger Attack. EX-RAY has 85% ac-

Table 2. EX-RAY on reflection attack

		ABS				ABS+EX-RAY				Acc/ROC	
		TP	FP	FN	TN	TP	FP	FN	TN		
CIFAR10	Same DOF	13	7	7	13	0.65	18	4	2	16	0.85
	Out of focus	12	7	8	13	0.63	16	4	4	16	0.80
	Ghost effect	9	7	11	13	0.55	17	3	3	17	0.85
ImageNet	Same DOF	12	6	5	11	0.68	15	3	2	14	0.85
	Out of focus	9	6	8	11	0.59	13	3	4	14	0.80
	Ghost effect	10	6	7	11	0.62	15	3	2	14	0.85

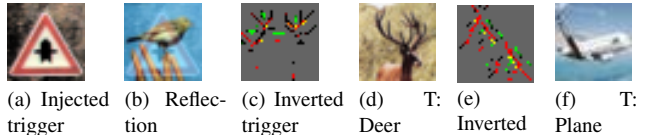


Figure 5. A case for reflection attack

Table 3. TrojAI leaderboard results; CE L denotes cross entropy loss and R-A denotes ROC-AUC

	Round 2		Round 3		Round 4	
	CE L	R-A	CE L	R-A	CE L	R-A
ABS only	0.685	0.736	0.541	0.822	0.894	0.549
ABS+EX-RAY	0.324	0.892	0.323	0.900	0.322	0.902
Deficit from top	0	0	0.023	-0.012	0	0

curacy on hidden-trigger attack whereas ABS has 68%. Details and case studies can be found in Appendix B.

4.2. Experiments on TrojAI and ImageNet Models

We evaluate EX-RAY on TrojAI rounds 2-4 training sets and ImageNet models. We use ABS as the upstream scanner and a relatively large trigger size bound to have a small number of false negatives. The experimental results show that the vanilla ABS encounters a large number of FPs, whereas EX-RAY substantially reduces the FPs by 78-100% with the cost of a slightly increased number of FNs by 0-30%. EX-RAY improves the overall detection accuracy by 17-41% across all the evaluated datasets. We also compare EX-RAY with eight baseline methods that make use of simple L2 distance, attribution/interpretation techniques, and one-sided (instead of symmetric) analysis. The results demonstrate that EX-RAY consistently outperforms baseline methods. In addition, we evaluate the runtime cost of EX-RAY, which takes 95 seconds to scan a model in TrojAI datasets on average, whereas the upstream scanner ABS takes 337 seconds. This delineates a reasonable overhead introduced by EX-RAY. We also study the effects of hyper-parameters of EX-RAY on a TrojAI dataset. The results show that EX-RAY is reasonably stable under various settings. Please see detailed results and discussions in Appendix H.

Results on TrojAI Leaderboard (Test Sets). Table 3



Figure 6. Adaptive attack triggers with different sizes

shows the results on TrojAI test sets. The column CE L shows the cross entropy loss of each method and the column R-A shows the ROC-AUC. In two of the three rounds, our solution achieves the top performance⁴. In round 3, it ranks number 2 and the results are comparable to the top performer. In addition, it beats the IARPA round goal (i.e., cross-entry loss lower than 0.3465) for all the three rounds. Our performance on the leaderboard, especially for round 2 that has a large number of natural backdoors and hence causes substantial difficulties for most performers⁵, suggests the contributions of EX-RAY. As far as we know, many existing solutions such as [11, 16, 25, 29, 39, 57, 58, 60, 66] have been tested in the competition by different performers.

4.3. Adaptive Attacks

We conduct three adaptive attacks. In the first attack, we use features of the target class as the trigger. Since EX-RAY distinguishes trojaned and clean models by comparing the similarity between inverted triggers’ features and distinctive features between the victim and target classes. Knowing our method, the attacker may choose to use the target class’s features as the trigger to evade EX-RAY. We use parts of a target class image as the trigger. We use four triggers with different sizes. For each trigger, we trojan 10 Network in Network models on CIFAR10 with dog being the target class. Figure 6 (a-d) show the triggers with the size of 80, 120, 160 and 200, respectively. Observe that they are all part of some dog image. In addition, we also train 20 clean models on CIFAR10 to see if ABS+EX-RAY can distinguish the trojaned and clean models. The results are shown in Table 4. The first row shows the trigger size. The second row shows the average attack success rate when using the triggers on *clean models*. Note that since these triggers are composed of the target class’s features, they might already be able to flip other images to the target even in clean models. The third row shows the FP rate. The fourth row shows the TP rate. Observe while ABS+EX-RAY consistently has a low FP rate, its TP rate decreases when the trigger become larger. When the trigger size is 200, the TP rate degrades to 0.5, meaning that it only recognizes half of the trojaned models. However, since the trigger (of size 200) is already quite large and contains strong target class

⁴TrojAI ranks solutions based on the cross-entropy loss of scanning results. Intuitively, the loss increases when the model classification diverges from the ground truth. Smaller loss suggests better performance [3]. Past leaderboard results can be found at [4].

⁵Most performers had lower than 0.80 ROC-AUC in round 2.

Table 4. Adaptive Attack using target class features

Trigger Size	80	120	160	200
ASR on clean models	0.39	0.56	0.63	0.7
FP/ # of clean models	0.1	0.1	0.1	0.1
TP/ # of trojaned models	1	1	0.8	0.5

features such that it can flip 70% of all the images to the target on all the clean models. This hence may not constitute a meaningful attack as it is almost equivalent to stamping a target class image.

In the second attack, we force the activations of victim class samples with the trigger to resemble those of the target class inputs such that EX-RAY cannot distinguish the two. While the strongest attack can increase the FP rate of EX-RAY, it causes substantial model accuracy degradation so that the model becomes unusable. In the third adaptive attack, we generate a trigger similar to a third class while having similar feature-level representations to the target class. Experiments show that EX-RAY has 75% true positive rate and 10% false positive rate on this adaptive attack. Please see Appendix L.

4.4. Other Experiments

Appendix C shows that EX-RAY can detect another two state-of-the-art backdoor attacks. Appendix D shows that EX-RAY outperforms three other state-of-the-art backdoor defenses. Appendix E and J demonstrate that EX-RAY can boost the detection performance of other upstream scanners on composite attack, reflection attack, and TrojAI dataset. Appendix K shows that EX-RAY’s masks indeed capture feature differences using a model interpretation technique. Appendix M shows how EX-RAY is used to help fixing injected backdoors.

5. Conclusion

We develop a method to detect complex backdoors that have large and dynamic triggers. It is built on a novel symmetric feature differencing technique that identifies a smallest set of features separating two sets of samples. Our results show that the technique is highly effective and outperforms the baselines. It also enables us to achieve top results on the rounds 2 and 4 leaderboard of the TrojAI competition, and rank the 2nd in round 3.

Acknowledgement

We thank the anonymous reviewers for their constructive comments. This research was supported, in part by IARPA TrojAI W911NF-19-S-0012, NSF 1901242 and 1910300, ONR N000141712045, N000141410468 and N000141712947. Any opinions, findings, and conclusions in this paper are those of the authors only and do not necessarily reflect the views of our sponsors.

References

- [1] Sri-csl/trinity-trojai: This repository contains code developed by the sri team for the iarpa/trojai program. <https://github.com/SRI-CSL/Trinity-TrojAI>, 2021. 3, 13, 15
- [2] torchvision.models — pytorch 1.7.0 documentation. <https://pytorch.org/docs/stable/torchvision/models.html>, 2021. 6, 13
- [3] Trojai leaderboard. <https://pages.nist.gov/trojai/>, 2021. 1, 3, 8, 13
- [4] Trojai past leaderboards. <https://pages.nist.gov/trojai/docs/results.html#previous-leaderboards>, 2021. 8
- [5] Marco Ancona, Enea Ceolini, Cengiz Öztireli, and Markus Gross. Towards better understanding of gradient-based attribution methods for deep neural networks. In *International Conference on Learning Representations*, 2018. 3, 14
- [6] Eugene Bagdasaryan et al. Blind backdoors in deep learning models. *arXiv preprint arXiv:2005.03823*, 2020. 1
- [7] Eugene Bagdasaryan, Andreas Veit, Yiqing Hua, Deborah Estrin, and Vitaly Shmatikov. How to backdoor federated learning. In *International Conference on Artificial Intelligence and Statistics*, pages 2938–2948. PMLR, 2020. 3
- [8] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017. 3, 14, 17
- [9] Yinzhi Cao, Alexander Fangxiao Yu, Andrew Aday, Eric Stahl, Jon Merwine, and Junfeng Yang. Efficient repair of polluted machine learning systems via causal unlearning. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*, 2018. 3
- [10] Bryant Chen, Wilka Carvalho, Nathalie Baracaldo, Heiko Ludwig, Benjamin Edwards, Taesung Lee, Ian Molloy, and Biplav Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering. *arXiv preprint arXiv:1811.03728*, 2018. 3
- [11] Huili Chen, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar. Deepinspect: A black-box trojan detection and mitigation framework for deep neural networks. In *IJCAI*, pages 4658–4664, 2019. 3, 8, 12, 13
- [12] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017. 3
- [13] Xiaoyi Chen, Ahmed Salem, Michael Backes, Shiqing Ma, and Yang Zhang. Badnl: Backdoor attacks against nlp models. *arXiv preprint arXiv:2006.01043*, 2020. 3
- [14] Edward Chou et al. Sentinet: Detecting localized universal attack against deep learning systems. *IEEE SPW 2020*, 2020. 1, 3
- [15] Min Du et al. Robust anomaly detection and backdoor attack detection via differential privacy. In *International Conference on Learning Representations*, 2019. 3
- [16] N Benjamin Erichson, Dane Taylor, Qixuan Wu, and Michael W Mahoney. Noise-response analysis for rapid detection of backdoors in deep neural networks. *arXiv preprint arXiv:2008.00123*, 2020. 1, 3, 8, 13
- [17] Minghong Fang, Xiaoyu Cao, Jinyuan Jia, and Neil Gong. Local model poisoning attacks to byzantine-robust federated learning. In *29th {USENIX} Security Symposium*, 2020. 3
- [18] Hao Fu, Akshaj Kumar Veldanda, Prashanth Krishnamurthy, Siddharth Garg, and Farshad Khorrani. Detecting backdoors in neural networks using novel feature-based anomaly detection. *arXiv preprint arXiv:2011.02526*, 2020. 3
- [19] Yansong Gao, Change Xu, Derui Wang, Shiping Chen, Damith C Ranasinghe, and Surya Nepal. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th Annual Computer Security Applications Conference*, pages 113–125, 2019. 3
- [20] Tianyu Gu, Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Evaluating backdooring attacks on deep neural networks. *IEEE Access*, 7:47230–47244, 2019. 1, 3
- [21] Wenbo Guo, Lun Wang, Xinyu Xing, Min Du, and Dawn Song. Tabor: A highly accurate approach to inspecting and restoring trojan backdoors in ai systems. *arXiv preprint arXiv:1908.01763*, 2019. 3, 12
- [22] Jonathan Hayase, Weihao Kong, Raghav Somani, and Sewoong Oh. Spectre: Defending against backdoor attacks using robust statistics. *arXiv preprint arXiv:2104.11315*, 2021. 3
- [23] Shanjiayang Huang, Weiqi Peng, Zhiwei Jia, and Zhuowen Tu. One-pixel signature: Characterizing cnn models for backdoor detection. *arXiv preprint arXiv:2008.07711*, 2020. 3
- [24] Xijie Huang et al. Neuroninspect: Detecting backdoors in neural networks via output explanations. *arXiv preprint arXiv:1911.07399*, 2019. 3, 12
- [25] Susmit Jha, Sunny Raj, Steven Fernandes, Sumit K Jha, Somesh Jha, Brian Jalaian, Gunjan Verma, and Ananthram Swami. Attribution-based confidence metric for deep neural networks. In *Advances in Neural Information Processing Systems*, pages 11826–11837, 2019. 1, 3, 8, 13
- [26] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, Fernanda Viegas, et al. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav). In *International conference on machine learning*, 2018. 3
- [27] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [28] Panagiota Kiourti, Kacper Wardega, Susmit Jha, and Wenchao Li. Trojdr: evaluation of backdoor attacks on deep reinforcement learning. In *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pages 1–6. IEEE, 2020. 3
- [29] Soheil Kolouri, Aniruddha Saha, Hamed Pirsiavash, and Heiko Hoffmann. Universal litmus patterns: Revealing backdoor attacks in cnns. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 301–310, 2020. 3, 8, 13
- [30] Alex Krizhevsky et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009. 6

- [31] Keita Kurita et al. Weight poisoning attacks on pre-trained models. *arXiv preprint arXiv:2004.06660*, 2020. 3
- [32] Yann LeCun et al. Mnist handwritten digit database. *ATT Labs [Online]*. Available: <http://yann.lecun.com/exdb/mnist>, 2010. 6
- [33] David Lewis. Counterfactuals and comparative possibility. In *Iffs*, pages 57–85. Springer, 1973. 2, 3
- [34] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. *arXiv preprint arXiv:2101.05930*, 2021. 1
- [35] Yiming Li, Tongqing Zhai, Baoyuan Wu, Yong Jiang, Zhifeng Li, and Shutao Xia. Rethinking the trigger of backdoor attack. *arXiv preprint arXiv:2004.04692*, 2020. 3
- [36] Junyu Lin, Lei Xu, Yingqi Liu, and Xiangyu Zhang. Composite backdoor attack for deep neural network by mixing existing benign features. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020. 1, 2, 6
- [37] Kang Liu et al. Fine-pruning: Defending against backdoor-attack on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, pages 273–294. Springer, 2018. 1
- [38] Yuntao Liu et al. Neural trojans. In *2017 IEEE International Conference on Computer Design*, 2017. 3
- [39] Yingqi Liu, Wen-Chuan Lee, Guanhong Tao, Shiqing Ma, Yousra Aafer, and Xiangyu Zhang. Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, pages 1265–1282, 2019. 1, 3, 6, 8, 13, 14
- [40] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and Xiangyu Zhang. Trojaning attack on neural networks. In *25th Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-22, 2018*. The Internet Society, 2018. 1
- [41] Yunfei Liu, Xingjun Ma, James Bailey, and Feng Lu. Reflection backdoor: A natural backdoor attack on deep neural networks. In *European Conference on Computer Vision*, 2020. 1, 2, 6
- [42] Shiqing Ma, Yingqi Liu, Guanhong Tao, Wen-Chuan Lee, and Xiangyu Zhang. Nic: Detecting adversarial samples with neural network invariant checking. In *NDSS*, 2019. 3
- [43] Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017. 13
- [44] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bisacco, Bo Wu, and Andrew Y Ng. Reading digits in natural images with unsupervised feature learning. 2011. 6
- [45] Anh Nguyen and Anh Tran. Wanet—imperceptible warping-based backdoor attack. *arXiv preprint arXiv:2102.10369*, 2021. 12
- [46] Tuan Anh Nguyen et al. Input-aware dynamic backdoor attack. *Advances in Neural Information Processing Systems*, 33, 2020. 3, 12
- [47] Ren Pang, Hua Shen, Xinyang Zhang, Shouling Ji, Yevgeniy Vorobeychik, Xiapu Luo, Alex Liu, and Ting Wang. A tale of evil twins: Adversarial inputs versus poisoned models. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, pages 85–99, 2020. 3
- [48] Ren Pang, Zheng Zhang, Xiangshan Gao, Zhaohan Xi, Shouling Ji, Peng Cheng, and Ting Wang. Trojanzoo: Everything you ever wanted to know about neural backdoors (but were afraid to ask). *arXiv preprint arXiv:2012.09302*, 2020. 1
- [49] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017. 6
- [50] Andrea Paudice et al. Label sanitization against label flipping poisoning attacks. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 5–15. Springer, 2018. 3
- [51] Ximing Qiao et al. Defending neural backdoors via generative distribution modeling. In *Advances in Neural Information Processing Systems*, pages 14004–14013, 2019. 3
- [52] Aniruddha Saha et al. Hidden trigger backdoor attacks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020. 2, 3, 6
- [53] Ahmed Salem, Rui Wen, Michael Backes, Shiqing Ma, and Yang Zhang. Dynamic backdoor attacks against machine learning models. *arXiv preprint arXiv:2003.03675*, 2020. 3
- [54] Ali Shafahi, W Ronny Huang, Mahyar Najibi, Octavian Suciu, Christoph Studer, Tudor Dumitras, and Tom Goldstein. Poison frogs! targeted clean-label poisoning attacks on neural networks. In *Advances in Neural Information Processing Systems*, 2018. 3
- [55] Guangyu Shen, Yingqi Liu, Guanhong Tao, Shengwei An, Qiuling Xu, Siyuan Cheng, Shiqing Ma, and Xiangyu Zhang. Backdoor scanning for deep neural networks through k-arm optimization. 2021. 1, 3, 13
- [56] Avanti Shrikumar, Peyton Greenside, Anna Shcherbina, and Anshul Kundaje. Not just a black box: Learning important features through propagating activation differences. *arXiv preprint arXiv:1605.01713*, 2016. 3, 14
- [57] Karan Sikka, Indranil Sur, Susmit Jha, Anirban Roy, and Ajay Divakaran. Detecting trojaned dnns using counterfactual attributions. *arXiv preprint arXiv:2012.02275*, 2020. 3, 8, 13
- [58] Octavian Suciu, Radu Marginean, Yigitcan Kaya, Hal Daume III, and Tudor Dumitras. When does machine learning {FAIL}? generalized transferability for evasion and poisoning attacks. In *27th {USENIX} Security Symposium*, 2018. 8, 13
- [59] Mukund Sundararajan et al. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, 2017. 3, 14
- [60] Di Tang, XiaoFeng Wang, Haixu Tang, and Kehuan Zhang. Demon in the variant: Statistical analysis of dnns for robust backdoor contamination detection. In *30th {USENIX} Security Symposium ({USENIX} Security 21)*, 2021. 3, 8, 13

- [61] Vale Tolpegin, Stacey Truex, Mehmet Emre Gursoy, and Ling Liu. Data poisoning attacks against federated learning systems. In *European Symposium on Research in Computer Security*, pages 480–501. Springer, 2020. 3
- [62] Brandon Tran et al. Spectral signatures in backdoor attacks. In *Advances in Neural Information Processing Systems*, pages 8000–8010, 2018. 3
- [63] Alexander Turner et al. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771*, 2019. 3
- [64] Akshaj Kumar Veldanda, Kang Liu, Benjamin Tan, Prashanth Krishnamurthy, Farshad Khorrani, Ramesh Karri, Brendan Dolan-Gavitt, and Siddharth Garg. Nnuculation: broad spectrum and targeted treatment of backdoored dnns. *arXiv preprint arXiv:2002.08313*, 2020. 1, 3
- [65] Binghui Wang, Xiaoyu Cao, Neil Zhenqiang Gong, et al. On certifying robustness against backdoor attacks via randomized smoothing. *arXiv preprint arXiv:2002.11750*, 2020. 1
- [66] Bolun Wang, Yuanshun Yao, Shawn Shan, Huiying Li, Bimal Viswanath, Haitao Zheng, and Ben Y Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, 2019. 1, 3, 8, 13, 15, 18
- [67] Hongyi Wang, Kartik Sreenivasan, Shashank Rajput, Harit Vishwakarma, Saurabh Agarwal, Jy-yong Sohn, Kangwook Lee, and Dimitris Papailiopoulos. Attack of the tails: Yes, you really can backdoor federated learning. *Advances in Neural Information Processing Systems*, 33, 2020. 3
- [68] Lun Wang, Zaynah Javed, Xian Wu, Wenbo Guo, Xinyu Xing, and Dawn Song. Backdoorl: Backdoor attack against competitive reinforcement learning. *arXiv preprint arXiv:2105.00579*, 2021. 3
- [69] Lior Wolf et al. Face recognition in unconstrained videos with matched background similarity. In *CVPR 2011*, pages 529–534. IEEE, 2011. 6
- [70] Eric Wong, Leslie Rice, and J Zico Kolter. Fast is better than free: Revisiting adversarial training. *arXiv preprint arXiv:2001.03994*, 2020. 13
- [71] Han Xiao et al. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 6
- [72] Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. Dba: Distributed backdoor attacks against federated learning. In *International Conference on Learning Representations*, 2019. 3
- [73] Xiaojun Xu, Qi Wang, Huichen Li, Nikita Borisov, Carl A Gunter, and Bo Li. Detecting ai trojans using meta neural analysis. *arXiv preprint arXiv:1910.03137*, 2019. 3, 12, 13
- [74] Xinyang Zhang et al. Trojanning language models for fun and profit. *arXiv preprint arXiv:2008.00312*, 2020. 3
- [75] Shihao Zhao, Xingjun Ma, Xiang Zheng, James Bailey, Jingjing Chen, and Yu-Gang Jiang. Clean-label backdoor attacks on video recognition models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14443–14452, 2020. 3

Appendix

A. An Example for the Necessity of Symmetry

Figure 7 presents an example for one-sided masks from a clean model #18 in TrojAI round 3. Figure 7 (a) and (b) present the victim and target classes and (c) a natural trigger (i.e., a trigger naturally exists and flips victim samples to the target label) generated by ABS, which resembles the central symbol in the target class. Figure (d) shows the one-sided mask from V to T , meaning that copying/mixing the feature maps as indicated by the mask from T samples to V samples can flip the classification results to T . Figure (e) shows the one-sided mask from V to V +trigger. Note that V +trigger samples are classified to T . Although in both cases V samples are flipped to T , the two one-sided masks have only one entry in common, suggesting that the ways they induce the classification results are different. In contrast, the symmetric masks share a lot of commonality.

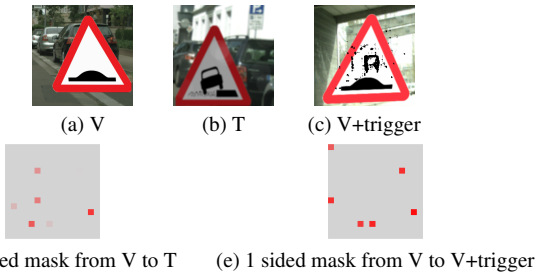


Figure 7. One sided masks for a clean model # 18 in round 3 with victim class $V = \#8$ and target class $T = \#3$

B. Detecting Hidden-trigger Attack

Table 5 shows the results on hidden-trigger attack. We use the optimal trigger size bound for ABS and 4,000 ($\approx 63 \times 63$) for ABS+EX-RAY. Observe that our technique can achieve 85% accuracy, surpassing ABS by 17%. Figure 8 (a) and (b) show an injected trigger and an example image stamped with the trigger, whose target label is terrier dog as shown in (f). Figure 8 (e) shows the generated trigger by ABS for the label terrier dog of the trojaned model. The trigger does not possess any features of the target label. It can hence be identified by EX-RAY as a true positive. In contrast, Figure 8 (c) and (d) show a trigger by ABS for a benign model and its target label jeans. Observe that the central part of the trigger resembles a pair of jeans. EX-RAY hence can flag the model as a true negative.

C. Detecting WaNet and Input-aware Dynamic Attacks

In this section, we evaluate EX-RAY on WaNet [45] and input-aware dynamic [46] attacks. We utilize the CIFAR10

Table 5. EX-RAY on hidden-trigger attack

	ABS					ABS+EX-RAY				
	TP	FP	FN	TN	Acc/ROC	TP	FP	FN	TN	Acc/ROC
ImageNet	12	6	5	11	0.68	15	3	2	14	0.85

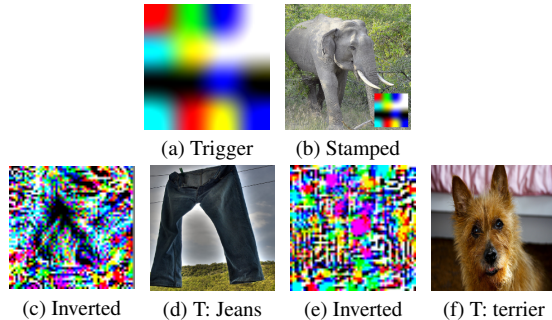


Figure 8. A case for hidden-trigger attack

Table 6. ABS + EX-RAY on WaNet and input-aware dynamic attacks

	TP	FP	FN	TN	Acc/ROC
Wanet	17	2	5	17	0.825
Input-aware	17	2	3	18	0.875

dataset and evaluate on 20 benign models and 40 trojaned models (20 trojaned models for each attack). We set the bound of the trigger size to be 12.5% of the input. The results are shown in Table 6. Observe that EX-RAY achieves 82.5% detection accuracy on WaNet and 87.5% detection accuracy on Input-aware.

D. Comparison with other SOTA Defenses on Complex Backdoors

We compare EX-RAY with four other state-of-the-art (SOTA) defenses, namely, Meta-Neural Analysis [73], DeepInspect [11], NeuronInspect [24], and TABOR [21], on reflection and composite attacks.

We use CIFAR10 and conduct experiments on 20 benign models, 20 trojaned models by composite attack, and 20 trojaned models by reflection attack. Meta-Neural Analysis outputs a binary result denoting whether a model is trojaned or not. The other three methods output a median absolute deviation (MAD) score for each model, which is used to distinguish benign and trojaned models. For DeepInspect, NeuronInspect, and TABOR, we search for the best possible bound of MAD scores for separating benign and trojaned models. Table 7 shows the results. Rows 2-9 show the results of Meta-Neural Analysis, DeepInspect, NeuronInspect, and TABOR on composite and reflection attacks. Rows 10-11 show the result of ABS + EX-RAY. Observe that ABS+EX-RAY outperforms the SOTA methods, hav-

Table 7. Comparison between EX-RAY and other defenses on composite and reflection attacks

		TP	FP	FN	TN	Acc/ROC
Meta-Neural Analysis	Composite	15	6	5	14	0.73
	Reflection	11	8	9	12	0.58
DeepInspect	Composite	20	19	0	1	0.53
	Reflection	20	20	0	0	0.5
NeuronInspect	Composite	4	0	16	20	0.6
	Reflection	2	0	18	20	0.55
TABOR	Composite	3	0	17	20	0.58
	Reflection	2	0	18	20	0.55
ABS+EX-RAY	Composite	17	3	3	17	0.85
	Reflection	18	4	2	16	0.85

ing at least 12% better accuracy on composite backdoors and 27% better on reflection backdoors.

During the TrojAI competition, performers tried many different SOTA methods [11, 16, 25, 29, 55, 57, 58, 60, 73] (including DeepInspect, Meta-Neural Analysis and K-Arm). Except for K-Arm [55], all other methods perform worse than ABS + EX-RAY in rounds 2 to 4. K-Arm performs better than ABS + EX-RAY in round 3 but worse than ABS + EX-RAY in rounds 2 and 4.

E. Using EX-RAY with Different Upstream Scanners on Complex Backdoors

We apply EX-RAY to different upstream scanners, including Neural Cleanse (NC) [66] and K-Arm [55], on detecting composite and reflection backdoors on CIFAR10. During the detection, we first use NC/K-Arm to invert triggers and then apply EX-RAY to determine whether a model is trojaned or not. The results are shown in Table 8. The first column denotes the detection methods. The second column shows the different attacks. Columns 3-7 show the detection results by vanilla NC and vanilla K-Arm. Columns 8-12 show the detection results by NC+EX-RAY and K-Arm+EX-RAY. Observe that EX-RAY can improve NC’s detection accuracy from 50-60% to 68-75%, and K-Arm’s from 55-58% to 73-75%. We also evaluate the combinations of EX-RAY and different upstream scanners on the TrojAI datasets. We use NC and Bottom-up-Top-down method [1] (used in the TrojAI competitions) as the upstream scanners. The results show that EX-RAY can consistently improve NC by 25%, and Bottom-up-Top-down by 2-15%. Please see more details in Appendix J.

F. Description of TrojAI and ImageNet Datasets

We use TrojAI rounds 2-4 training and test datasets [3]. EX-RAY *does not* require training and hence we use both training and test sets as regular datasets in our experiments.

Table 8. EX-RAY with different upstream scanners on composite and reflection attacks

		Vanilla					+EX-RAY				
		TP	FP	FN	TN	Acc/ROC	TP	FP	FN	TN	Acc/ROC
NC	Composite	7	3	13	17	0.60	14	4	6	16	0.75
	Reflection	5	5	15	15	0.50	10	3	10	17	0.68
K-Arm	Composite	3	1	17	19	0.55	16	6	4	14	0.75
	Reflection	18	15	2	5	0.58	16	7	4	13	0.73

TrojAI round 2 training set has 552 clean models and 552 trojaned models with 22 structures. Each TrojAI model has its own unique dataset. The data are mostly synthetic traffic signs with some street view background. A traffic sign is a polygon of solid color with some symbol in the center. The models are classifiers for the different kinds of signs. TrojAI has two types of backdoors: polygons (i.e., static patch triggers) and Instagram filters (i.e., dynamic and pervasive triggers). Round 2 test set has 72 clean and 72 trojaned models. Most performers had difficulties for round 2 due to the prevalence of natural triggers, which are small triggers that naturally exist and can flip classification results among benign labels. IARPA hence introduced adversarial training [43, 70] in round 3 to enlarge the distance between classes and suppress natural triggers. Round 3 training set has 504 clean and 504 trojaned models and the test set has 144 clean and 144 trojaned models. In round 4, triggers may be position dependent, meaning that they only cause misclassification when stamped at a specific position inside the foreground object. A model may have multiple backdoors. The number of clean images provided is reduced from 10-20 (in rounds 2 and 3) to 2-5. Its training set has 504 clean and 504 trojaned models and the test set has 144 clean and 144 trojaned models. Training sets were evaluated on our local server whereas test set evaluation was done remotely by IARPA on their server.

We also use a number of models on ImageNet. They have the VGG, ResNet and DenseNet structures. We use 7 trojaned models from [39] and 17 pre-trained clean models from torchvision zoo [2].

G. Parameter Settings

EX-RAY has three hyper-parameters, α to control the weight changes of cross-entropy loss in function (5) (in the design section), β to control the similarity comparison between masks in condition (6) in Section 3.2, and γ the accuracy threshold in cross-validation checks of masks in Section 3.2. We use 0.1, 0.8, and 0.8, respectively, by default. In our experiments, we use ABS and NC as the upstream scanners. The numbers of optimization epochs are 60 for ABS and 1000 for NC. The other settings are default unless stated otherwise. The experiments are all done on an identical machine with a single 11GB memory NVIDIA RTX

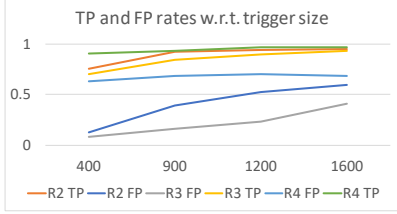


Figure 9. Rounds 2-4 true positive rates (TPs) and false positive rates (FPs) versus trigger size (in pixels) by ABS

2080Ti GPU (with the lab server configuration), except for the TrojAI test sets that are run on IARPA server with a single 32GB memory NVIDIA V100 GPU.

H. Experiments on TrojAI and ImageNet Models

In the first experiment, we evaluate EX-RAY on TrojAI rounds 2-4 training sets and the ImageNet models. We do not include TrojAI test sets in this experiment as the test sets were evaluated on an IARPA server and the results are reflected on the leaderboard. Here we use ABS as the upstream scanner as it is much faster than NC.

A critical setup for scanners that produce triggers, such as ABS and NC, is the maximum trigger size. A large value enables detecting injected backdoors with large triggers, while producing a lot of natural triggers and hence false positives. Figure 9 studies how the true positives (TPs) and false positives (FPs) change with different trigger bounds in the vanilla ABS, on the TrojAI rounds 2-4 training sets. Observe that both grow with the trigger size. Observe that there is a lower FP rate in round 3 (compared to round 2), illustrating the effect of adversarial training, although the number is still large when the trigger size is large. Round 4 has the highest FP rate because the number of clean images available is decreased and it is hence very easy for scanners to find (bogus) triggers that can induce misclassification on all the available images.

Based on the study, we use the trigger size bound 900 pixels for round 2, 1600 pixels for round 3, and 1200 pixels for round 4 for our experiment such that the upstream scanner does not miss many true positives to begin with and we can stress test EX-RAY.

Baselines. In the experiment, we compare EX-RAY against 8 baselines. The first baseline is using L2 distance of inner activation between $V + t$ and T . Such a distance for a natural trigger is supposed to be smaller than that of an injected trigger. We use unsupervised learning to report the best separation. In the second baseline, we use half of the models to train a random forest classifier based on the inner activations and logits values to distinguish natural and injected triggers and test it on the other half. Specifically, the classifier takes the L2 distance between $V + t$ and T , L2 distance between $T + t$ and T , L2 distance between V and T , logits

of V , logits of $V + t$ and logits of T . The third baseline uses *integrated gradients (IG)* [59], an attribution technique, to find important neurons for $V + t$ and for T and then apply the aforementioned L2 distance comparison on the 10% most important neurons. Originally, integrated gradients were used in model explanation to identify important pixels. We adapt it to work on inner layers to identify important neurons. The next three baselines are similar to the third except having different methods to identify important neurons. Specifically, the fourth baseline uses *Deeplift* [56], the fifth uses *Occlusion* [5] and the sixth uses *Network Dissection (NE)* [8]. For baselines 4-7, we use unsupervised learning to find the best separation (of natural and injected backdoors). We will release the settings together with our system upon publication. EX-RAY is symmetric. To study the necessity of symmetry, the seventh and eighth baselines are one-sided versions of EX-RAY, that is, requiring satisfying either constraint (1) or (2) in the design section.

The results are shown in Table 9. The first column shows the methods. The first method is the vanilla ABS. Columns 2-4 show the results for TrojAI round 2 models with polygon backdoors. Column 2 shows the number of true positives (TPs). Note that there are 276 trojaned models with polygon backdoors. As such the vanilla ABS having 254 TPs means it has 22 false negatives. Column 3 shows the number of false positives (FPs) out of the 552 clean models. Column 4 shows the overall detection accuracy (on the total $552+276=828$ models). Columns 5-7 show the results for round 2 models with Instagram filter backdoors. ABS uses a different method for filter backdoors. Instead of reverse engineering a pixel patch, it reverse engineers a one-layer kernel denoting general filter transformation [39]. Hence, we separate the evaluation of EX-RAY on the two kinds of backdoors. Note that the accuracy is computed considering the same 552 clean models. The overall results (for all kinds of backdoors) on the leaderboard are presented in the main text. Columns 8-13 show the results for round 3 and columns 14-19 for round 4. Columns 14-16 show the results for ImageNet patch attack.

The results show that the vanilla ABS has a lot of FPs (in order not to lose TPs) and EX-RAY can substantially reduce the FPs by 78-100% with the cost of increased FNs (i.e., losing TPs) by 0-30%. The overall detection accuracy improvement (from vanilla ABS) is 17-41% across the datasets. Also observe that EX-RAY consistently outperforms all the baselines, especially the non-EX-RAY ones. Attribution techniques can remove a lot of natural triggers indicated by the decrease of FPs. However, they preclude many injected triggers (TPs) as well, leading to inferior performance. The missing entries for NE are because it requires an input region to decide important neurons, rendering it inapplicable to filters that are pervasive. Symmetric EX-RAY outperforms the one-sided versions, suggesting

Table 9. Effectiveness of EX-RAY; (T:276,C:552) means that there are 276 trojaned models and 552 clean models

	TrojAI R2						TrojAI R3						TrojAI R4						ImageNet		
	Polygon trigger (T:276,C:552)			Filter trigger (T:276,C:552)			Polygon trigger (T:252,C:504)			Filter trigger (T:252,C:504)			Polygon trigger (T:143,C:504)			Filter trigger (T:361,C:504)			Patch Trigger (T:7, C:17)		
	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc
Vanilla ABS	254	218	0.710	260	293	0.626	235	208	0.702	213	334	0.528	137	355	0.442	331	376	0.531	7	7	0.708
Inner L2	188	93	0.782	153	123	0.703	210	111	0.798	133	123	0.680	73	137	0.680	208	217	0.572	7	0	1
Inner RF	192	76	0.807	196	101	0.781	159	46	0.816	153	110	0.724	133	265	0.575	330	353	0.556	7	0	1
IG	172	29	0.840	192	66	0.818	162	58	0.804	52	41	0.681	84	53	0.827	210	87	0.725	5	0	0.917
Deeplift	152	11	0.837	189	21	0.869	162	59	0.803	78	67	0.681	84	54	0.825	203	54	0.755	6	0	0.958
Occulation	173	24	0.847	207	47	0.860	164	58	0.807	78	66	0.683	85	52	0.830	251	107	0.749	7	3	0.875
NE	180	58	0.814	-	-	-	187	72	0.819	-	-	-	59	72	0.759	-	-	-	7	4	0.833
1-sided(V to T)	157	19	0.833	195	33	0.862	202	62	0.852	153	51	0.802	107	82	0.818	236	50	0.798	7	0	1
1-sided(T to V)	134	4	0.824	158	18	0.835	187	50	0.848	134	27	0.808	102	56	0.850	179	9	0.779	1	1	0.958
EX-RAY	198	19	0.883	204	32	0.874	200	46	0.870	149	39	0.812	105	53	0.859	242	46	0.809	7	0	1

the need of symmetry.

Runtime. EX-RAY’s time complexity is $\Omega(n)$ with n the number of triggers the upstream technique generates. Our upstream ABS uses neuron stimulation analysis to select three most likely target labels for each (victim) label for trigger inversion. It also filters out large sized triggers. In TrojAI, the number of classes per model varies from 15 to 45. On average, EX-RAY takes 12s to process a trigger, 95s to process a model. ABS takes 337s to process a model, producing 8.5 triggers on average.

I. Effects of Hyperparameters

We study EX-RAY performance with various hyperparameter settings, including the different layer to which EX-RAY is applied, different trigger size settings (in the upstream scanner) and different SSIM score bounds (in filter backdoor scanning to ensure the generated kernel does not over-transform an input), and the α , β , and γ settings of EX-RAY. Table 11 shows the results for layer selection. The row “Middle” means that we apply EX-RAY at the layer in the middle of a model. The rows “Last” and “2nd last” show the results at the last and the second-last convolutional layers, respectively. Observe that layer selection may affect performance to some extent and the second to the last layer has the best performance. Table 10 shows the results with and without the additional validation checks. Observe that removing the additional validation check results in 0.7% to 3% decrease in detection accuracy. Tables 12 shows that a large trigger size degrades EX-RAY’s performance but EX-RAY is stable in 900 to 1200. Table 13 shows that the SSIM score bound has small effect on performance in 0.7-0.9. Note that an SSIM score smaller than 0.7 means the transformed image is quite different (in human eyes). Figures 10, 11, and 12 show the performance variations with α , β , and γ , respectively. The experiments are on the mixture of trojaned models with polygon triggers and the

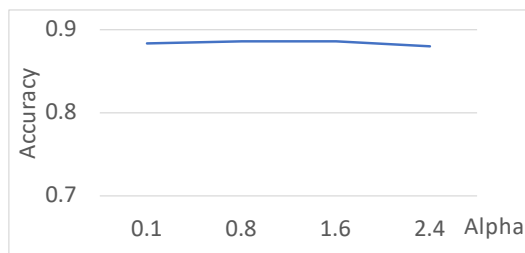


Figure 10. Accuracy changes with α on TrojAI R2

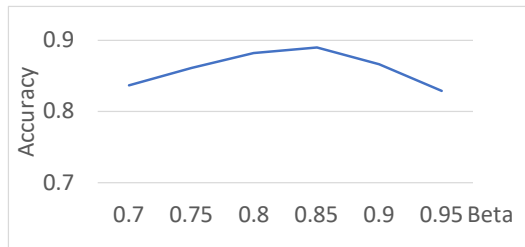


Figure 11. Accuracy changes with β on TrojAI R2

clean models from TrojAI round 2. For β and γ , we sample from 0.7 to 0.95 and for α we sample from 0.1 to 2.4. Observe that changing α and γ does not have much impact on the overall accuracy. When we change β from 0.7 to 0.95, the overall accuracy is still consistently higher than 0.83. These results show the stability of EX-RAY.

J. Using EX-RAY with Different Upstream Scanners on TrojAI dataset

In this experiment, we use EX-RAY with different upstream scanners, including Neural Cleanse (NC) [66] and the Bottom-up-Top-down method by the SRI team in the TrojAI competition [1]. The latter has two sub-components, trigger generation and a classifier that makes use of features collected from the trigger generation process. We created two scanners out of their solution. In the first one, we ap-

Table 10. EX-RAY w. and w.o. additional check; (T:276,C:552) means that there are 276 trojaned models and 552 clean models

	TrojAI R2						TrojAI R3						TrojAI R4					
	Polygon Trigger (T:276,C:552)			Filter Trigger (T:276,C:552)			Polygon Trigger (T:252,C:504)			Filter Trigger (T:252,C:504)			Polygon trigger (T:143,C:504)			Filter trigger (T:361,C:504)		
	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc
W. additional check	198	19	0.883	204	32	0.874	200	46	0.870	149	39	0.812	105	53	0.859	242	46	0.809
W.o. additional check	206	33	0.876	216	71	0.844	207	71	0.843	158	62	0.793	110	77	0.829	275	93	0.793

Table 11. EX-RAY with different layer; (T:276,C:552) means that there are 276 trojaned models and 552 clean models

	TrojAI R2						TrojAI R3						TrojAI R4					
	Polygon Trigger (T:276,C:552)			Filter Trigger (T:276,C:552)			Polygon Trigger (T:252,C:504)			Filter Trigger (T:252,C:504)			Polygon trigger (T:143,C:504)			Filter trigger (T:361,C:504)		
	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc
Middle	215	54	0.861	220	97	0.815	212	55	0.874	153	58	0.792	102	75	0.821	257	77	0.791
Second Last	198	19	0.883	204	32	0.874	200	46	0.870	149	39	0.812	105	53	0.859	242	46	0.809
Last	141	6	0.83	171	16	0.854	193	55	0.849	141	27	0.817	84	37	0.852	196	37	0.766

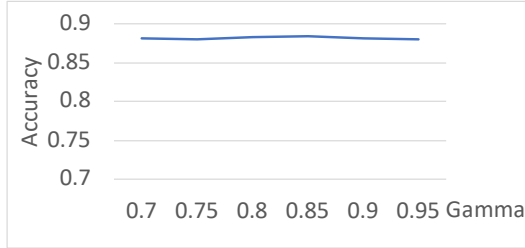


Figure 12. Accuracy changes with γ on TrojAI R2

Table 12. EX-RAY with different trigger sizes; (T:276, C:552) means there are 276 trojaned models and 552 clean models

	TrojAI R2 (T:276,C:552)			TrojAI R3 (T:252,C:504)			TrojAI R4 (T:143,C:504)		
	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc
900	198	19	0.883	157	19	0.849	95	58	0.836
1200	203	30	0.876	175	39	0.847	105	53	0.859
1600	210	46	0.864	200	46	0.870	108	77	0.827

Table 13. EX-RAY with different SSIM scores; (T:276, C:552) means there are 276 trojaned models and 552 clean models

SSIM Score	TrojAI R2 (T:276,C:552)			TrojAI R3 (T:252,C:504)			TrojAI R4 (T:361,C:504)		
	TP	FP	Acc	TP	FP	Acc	TP	FP	Acc
0.9	145	4	0.837	115	9	0.807	234	47	0.799
0.8	160	13	0.844	149	39	0.812	242	46	0.809
0.7	204	32	0.874	178	90	0.783	175	13	0.770

ply EX-RAY on top of their final classification results (i.e., using EX-RAY as a refinement). We call it SRI-CLS. In the second one, we apply EX-RAY right after their trigger generation. We have to replace their classifier with the sim-

pler unsupervised learning (i.e., finding the best separation) as adding EX-RAY changes the features and nullifies their original classifier. We call it SRI-RE. We use the round 2 clean models and models with polygon backdoors to conduct the study as NC does not handle Instagram filter triggers. For SRI-CLS, the training was on 800 randomly selected models and testing was on the remaining 146 trojaned models and 158 clean models. The other scanners do not require training. The results are shown in Table 14. The T and C columns stand for the number of trojaned and clean models used in testing, respectively. Observe that the vanilla NC identifies 180 TPs and 332 FPs with the accuracy of 44.7%. With EX-RAY, the FPs are reduced to 73 (81.1% reduction) and the TPs become 127 (29.4% degradation). The overall accuracy improves from 44.7% to 70.8%. The improvement for SRI-RE is from 53.6% to 68.5%. The improvement for SRI-CLS is relative less significant. That is because 0.882 accuracy is already very close to the best performance for this round. The results show that EX-RAY can consistently improve upstream scanner performance. Note that the value of EX-RAY lies in suppressing false warnings. It offers little help if the upstream scanner has substantial false negatives. In this case, users may want to tune the upstream scanner to have minimal false negatives and then rely on the downstream EX-RAY to prune the false positives like we did in the ABS+EX-RAY pipeline.

K. Mask and Differential Features

In this experiment, we aim to demonstrate that the masks computed by EX-RAY indeed capture the feature differences. Specifically, we want to show that 1) *the mask between $V + t$ and V covers the trigger features and does not overlap with the natural differences between V and T*

for a trojaned model and 2) the mask between $V + t$ and V captures the natural differences between V and T for a clean model. We use a model interpretation technique similar to [8] to project the large activation values of feature maps (i.e., neurons) in the mask between $V + t$ and V back to the input space and observe which input areas are highlighted. Figure 13 shows a sample result. Figures (a)-(b) correspond to a trojaned model #7 in TrojAI round 2. Figure (a) shows a victim sample with the trigger, which is a purple polygon. The area in light-green shows the input area corresponding to the activated neurons in the mask. Observe that the two align nicely, indicating the mask captures the trigger features. In contrast, figure (b) shows a target sample and also the input area corresponding to activated neurons in the mask, if any. Observe that those neurons do not have large activations. Figures (c)-(d) show the results for a clean model #123. Figure (c) shows that the (natural) trigger is to the right and below the central symbol of the victim sample. Observe that while there is a highlighted area in the $V + t$ sample (c) covering the trigger, there is also a highlighted area in the T sample (d) covering the target features, demonstrating that the mask between $V + t$ and V indeed captures T 's features. Figure 14 plots the maximum activation values for all the neurons in the mask, with (a)-(c) for model #7 and (d)-(f) for model #123. For example, a data point in (a) shows the average maximum activation of a neuron in the mask for all $V + t$ samples (y axis), versus its average maximum activation for all V samples (x axis). From (a)-(c), we can observe that these neurons are substantially activated when t is present but never for clean V or T samples. In contrast, (d)-(f) show that the neurons in the mask are activated when either t is present or V/T samples are provided. We studied a few other models. Their results are similar and hence omitted.

L. Two Additional Adaptive Attacks

EX-RAY is not a stand-alone defense technique and supposed to be part of an end-to-end scanning pipeline. We devise another the second attack that forces the internal activations of victim class inputs embedding the trigger to resemble the activations of the target class inputs such that EX-RAY cannot distinguish the two. In particular, we train a Network in Network model on CIFAR10 with an 8×8 patch trigger. In order to force the activations of images stamped with the trigger to resemble those of target class

Table 14. EX-RAY with different upstream scanners

	Vanilla					+EX-RAY					Acc Inc
	TP	T	FP	C	Acc	TP	T	FP	C	Acc	
NC	180	252	332	552	0.483	127	252	73	552	0.732	0.249
SRI-RE	164	252	272	552	0.536	112	252	97	552	0.685	0.149
SRI-CLS	120	146	17	158	0.858	119	146	9	158	0.882	0.024

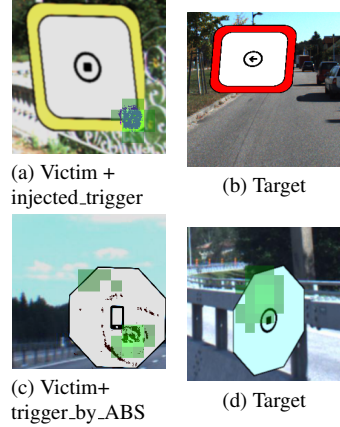


Figure 13. Trojaned model #7 from TrojAI round 2 in (a)-(b) and clean model #123 in (c)-(d). In (a), the trigger is stamped at the right-bottom corner with the light-green area corresponding to the neurons in the mask by an interpretation technique; (b) shows the neurons in the mask do not have large activations at all for a target sample; (c) and (d) show that the neurons in the mask capture features in both the trigger and the target.

images, we design an adaptive loss to minimize the differences between the two. In particular, we measure the differences of the means and standard deviations of feature maps. During training, we add the adaptive loss to the normal cross-entropy loss. The effect of adaptive loss is controlled by a weight value, which essentially controls the strength of attack as well. Besides the adaptively trojaned model, we also train 20 clean models on CIFAR10 to see if ABS+EX-RAY can distinguish the trojaned and clean models. The results are shown in Table 15. The first row shows the adaptive loss weight. A larger weight value indicates stronger attack. The second row shows the trojaned model's accuracy on clean images, including both the overall accuracy and the victim label accuracy. The third row shows the attack success rate of the trojaned model. The fourth row shows the FP rate. The fifth row shows the TP rate. Observe while ABS+EX-RAY does not miss trojaned models, its FP rate grows with the strength of attack. When the weight value is 1000, the FP rate of ABS+EX-RAY becomes 0.65 while its TP rate remains 1, meaning effectiveness degradation. However at this setting, the model accuracy has degraded so much that such model is unlikely used in practice.

In the third adaptive attack, we first generate a trigger similar to a third class while having similar feature-level representations to the target class. We generate such triggers by optimizing two losses. The first is the cross entropy loss between the model output on images stamped with the trigger and the third class label (similar to adversarial noise for a third class). The second loss is the mean squared error loss between the inner activation of the images stamped with the trigger and the inner activation of the target class

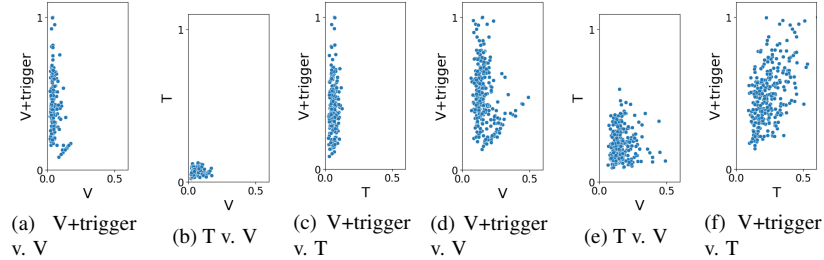


Figure 14. Average maximum activation values for neurons in the computed masks for a trojaned model #7 in (a)-(c), and for a clean model #123 in (d)-(f).

Table 15. The second adaptive attack

Weight of adaptive loss	1	10	100	200	400	600	800	1000	10000
Acc (model/label)	0.89/0.73	0.88/0.73	0.87/0.7	0.87/0.7	0.86/0.69	0.845/0.66	0.84/0.66	0.82/0.64	0.1
ASR	0.99	0.99	0.99	0.98	0.94	0.98	0.96	0.97	-
FP/ # of clean models	0	0.2	0.2	0.2	0.35	0.45	0.6	0.65	-
TP/ # of clean models	1	1	1	1	1	1	1	1	-

images (similar to adversarial feature-level attack). After generating the triggers, we use data poisoning to trojan the models. We do the experiment on CIFAR10. We choose label 0 as the target label and label 8 as the third label. We choose conv7 in NiN models as the feature layer and optimize neuron activations in this layer. We find that we need to enlarge the trigger size to have similar inner activations as the target label images. We generate triggers with 4 different sizes, 120, 140, 160, 200. The triggers are shown in Figure 15. We train 20 benign NiN models and 20 feature-level adaptive attack NiN models for each trigger size.

Table 16 shows the results of EX-RAY. Row 1 shows the different trigger sizes. Row 2 shows the mean squared activation differences. Observe that with the increase of trigger size, we can optimize the difference to a smaller value. A trigger with a small feature difference may be difficult to be detected. Rows 3 and 4 show the false positive and true positive rates. Observe that EX-RAY has 75% true positive rate when the trigger is 160 and 65% true positive rate when trigger size is 200. When the trigger size is 200, the trigger already covers a large part of the image. The attack becomes less meaningful.

Table 16. The third adaptive attack

Trigger size	120	140	160	200
Mean squared feature difference	0.153	0.116	0.034	0.009
FP/ # of clean models	0.1	0.1	0.1	0.1
TP/ # of clean models	1	0.8	0.75	0.65

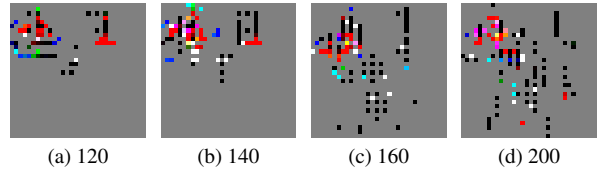


Figure 15. Triggers in the third adaptive attack

M. Fixing Models with Injected and Natural Backdoors

In this experiment, we try to fix 5 benign models and 5 trojaned models on CIFAR10. Fixing a benign model means enlarging class distances to make it less vulnerable to (small) natural backdoors. The trojaned models are trojaned by label-specific data poisoning. Here we use unlearning [66] which stamps triggers generated by scanning methods on images of victim label to finetune the model and forces the model to unlearn the correlations between the triggers and the target label. The process is iterative, bounded by the level of model accuracy degradation. The level of repair achieved is measured by the trigger sizes of the fixed model. Larger triggers indicate the corresponding backdoors become more difficult to exploit. The trigger size increase rate suggests the difficulty level of repair.

Table 17 shows the average accuracy and average reverse engineered trigger size before and after fixing the models. All models have the same repair budget. We can see that natural triggers have a larger accuracy decrease. Natural trigger size only increases by 34.4 whereas injected trigger size increases by 78.

We show the trigger size for each label pair for a trojaned model in Figure 16 and for a benign model in Figure 17. Figure 16 (a) shows the trigger size between each

Table 17. Average trigger size change before and after unlearning

	Natural Trigger		Injected Trigger	
	Before	After	Before	After
Avg Acc	88.7%	85.9%	86.4%	85.4%
Avg Trigger Size	25.8	60.2	19	97

pair of labels. The columns denote the victim label and the rows denote the the target label. For example, the gray cell in Figure 16 (a) shows the trigger size to flip class 1 to class 0. Figure 16 (b) follows the same format and shows the result for a trojaned model after unlearning. In the trojaned model the injected trigger flips class 1 to class 0. Before unlearning, class 1 and class 0 have the smallest trigger size 21. Unlearning increases the trigger size between the two to 106, which is above the average trigger size between any pairs. Intuitively, one can consider the backdoor is fixed. In the benign model, the natural trigger flips class 3 to class 5. As shown in Figure 17, unlearning increases the trigger size from 24 to 59 and 59 is still one of the smallest trigger size among all label pairs for the fixed model. It demonstrates that natural backdoors are inevitable and difficult to fix. AI model users can use EX-RAY to find injected backdoors and speed up fixing process by prioritizing fixing injected backdoors first.

Note that model repair is not the focus of the paper and trigger size may not be a good metric to evaluate repair success for the more complex semantic backdoors. The experiment is to provide initial insights. A thorough model repair solution belongs to our future work.

(a) Before unlearning

	0	1	2	3	4	5	6	7	8	9
0	-	48	34	75	42	52	62	46	48	52
1	21	-	74	91	88	96	72	80	81	45
2	32	54	-	66	39	57	54	61	78	60
3	34	53	35	-	42	27	46	50	72	47
4	29	45	29	49	-	36	46	48	63	48
5	40	70	35	46	43	-	53	49	81	56
6	29	48	23	41	44	61	-	66	70	59
7	40	77	55	78	40	52	81	-	82	60
8	21	44	42	75	50	59	60	65	-	47
9	29	62	78	85	69	70	73	62	73	-

(b) Before unlearning

	0	1	2	3	4	5	6	7	8	9
0	-	84	56	103	77	96	82	96	56	78
1	106	-	132	162	150	140	113	134	124	70
2	92	111	-	88	79	79	62	99	109	106
3	105	92	66	-	86	60	58	82	124	86
4	96	92	55	81	-	72	53	77	99	95
5	119	100	64	70	91	-	58	101	136	90
6	107	97	86	88	99	93	-	113	113	97
7	94	101	92	124	87	87	81	-	126	100
8	50	72	68	104	98	106	81	101	-	79
9	104	87	129	119	117	115	110	108	123	-

Figure 16. Injected trigger distance matrix before and after unlearning

	0	1	2	3	4	5	6	7	8	9
0	-	43	44	47	38	42	67	68	37	48
1	62	-	89	88	79	78	70	85	76	47
2	53	58	-	37	35	42	51	71	72	62
3	62	66	40	-	40	24	48	59	72	56
4	61	57	38	48	-	31	52	52	85	64
5	69	66	43	33	46	-	51	61	73	57
6	66	55	32	35	44	38	-	80	87	62
7	74	77	67	61	38	39	74	-	92	68
8	29	44	55	61	48	51	62	65	-	46
9	79	57	84	72	67	67	83	76	72	-

	0	1	2	3	4	5	6	7	8	9
0	-	79	56	89	63	90	65	99	59	76
1	120	-	134	114	123	154	77	119	122	48
2	95	101	-	75	58	74	59	82	121	82
3	104	98	57	-	58	59	40	90	124	80
4	104	100	60	88	-	79	50	73	117	79
5	95	91	58	84	76	-	52	77	131	86
6	114	115	77	129	102	89	-	131	137	93
7	104	120	110	103	68	92	66	-	129	78
8	36	61	74	80	66	112	66	86	-	70
9	128	109	117	103	112	120	73	124	105	-

Figure 17. Natural trigger distance matrix before and after unlearning