# **Rethinking the Reverse-engineering of Trojan Triggers**

Anonymous Author(s) Affiliation Address email

# Abstract

Deep Neural Networks are vulnerable to Trojan (or backdoor) attacks. Reverse-1 engineering methods can reconstruct the trigger and thus identify affected models. 2 Existing reverse-engineering methods only consider input space constraints, e.g., З trigger size in the input space. Expressly, they assume the triggers are static patterns 4 in the input space and fail to detect models with feature-space triggers such as image 5 style transformations. We observe that both input-space and feature-space Trojans 6 are associated with feature space hyperplanes. Based on this observation, we 7 design a novel reverse-engineering method that exploits the feature space constrain 8 to reverse-engineer Trojan triggers. Results on four datasets and seven different 9 attacks demonstrate that our solution effectively defends both input-space and 10 feature-space Trojans. It outperforms state-of-the-art reverse-engineering methods 11 and other types of defenses in both Trojaned model detection and mitigation tasks. 12 On average, the detection accuracy of our method is 93%. For Trojan mitigation, 13 our method can reduce the ASR (attack success rate) to only 0.26% with the 14 BA (benign accuracy) remaining nearly unchanged. Our code can be found in 15 https://anonymous.4open.science/r/FeatureRE-10B7. 16

# 17 **1 Introduction**

DNNs are vulnerable to Trojan attacks [1-5]. After injecting a Trojan into the DNN model, the 18 adversary can manipulate the model prediction by adding a *Trojan trigger* to get the target label. The 19 adversary can inject the Trojan by performing the poisoning attack or supply chain attack. In the 20 poisoning attacks, the adversary can control the training dataset and injects the Trojan by adding 21 samples with the Trojan trigger labeled as the target label. In the supply chain attack, the adversary 22 can replace a benign model with a Trojaned model by performing the supply chain attack. The Trojan 23 trigger is becoming more and more stealthy. Earlier works use static patterns, e.g., a yellow pad as 24 25 the trigger, which is known as the input space triggers. Researchers recently proposed using more 26 dynamic and input-aware techniques to generate stealthy triggers that mix with benign features, which are referred to as the feature space triggers. For example, the trigger of the feature-space Trojans can 27 be a warping process [6] or a generative model [7, 8, 3]. The Trojan attack is a prominent threat to 28 the trustworthiness of DNN models, especially in security-critical applications, such as autonomous 29 driving [1], malware classification [9], and face recognition [10]. 30

Prior works have proposed several ways to defend against Trojan attacks, such as removing poisons 31 in training [11–13], detecting Trojan samples at runtime [14–16], etc. Many of them can only 32 work for one type of Trojan attack. For example, training and pre-training time defense (e.g., 33 removing poisoning data, training a benign model with poisoning data) fail to defend against the 34 supply chain attack. Trigger reverse-engineering [17-21] is a general method to defend against 35 different Trojan attacks under different threat models. It works by searching for if there exists an 36 input pattern that can be used as a trigger in the given model. If we can find such a trigger, the 37 model has a corresponding Trojan and is marked as malicious and vice versa. Existing reverse-38

engineering methods assume that the Trojan triggers are static patterns in the input space and develop 39 an optimization problem that looks for an input pattern that can be used as the trigger. This assumption 40 is valid for input-space attacks [1, 22, 2] that use static triggers (e.g., a colored patch). Feature space 41 attacks [6, 5, 7, 23, 8, 4, 3] break this assumption. Existing trigger reverse-engineering methods [17– 42 21] constrain the optimization by using heuristics or empirical observations on existing attacks, such 43 as pixel values are in range [0, 255], and the trigger's size is small. Such heuristics are also invalid 44 for feature-space triggers that change all pixels in images. Reverse-engineering the feature space is 45 challenging. Unlike input space, there are no constraints that can be directly used. 46

In this paper, we propose a trig-47 ger reverse-engineering method that 48 works for feature space triggers. Our 49 intuition is that *features representing* 50 the Trojan are orthogonal to other fea-51 tures. Because a trigger works for 52 a set of samples (or all of them, de-53 pending on the attack type), changing 54 the input content without removing 55 the Trojan features will not change 56 the prediction. That is, changing Tro-57 jan and benign features will not affect 58 each other. Trojan features will form 59 a hyperplane in the high dimensional 60 space, which can constrain the search 61 in feature space. We then develop our 62



Fig. 1: Existing reverse-engineering (RE) and ours.

reverse-engineering method by exploiting the feature space input-space constraint. Fig. 1 demon-63 strates our idea. Existing reverse-engineering methods only consider the input-space constraint. It 64 conducts reverse-engineering via searching a static trigger pattern in the input space. These methods 65 fail to reverse-engineer feature-space Trojans whose trigger is dynamic in the input space. Instead, our 66 idea is to exploit the feature space constraint and searching a feature space trigger using the constraint 67 that the Trojan features will form a hyperplane. At the same time, we also reverse-engineer the input 68 69 space Trojan transformation based on the feature space constraint. To the best of our knowledge, we are the first to propose feature-space reverse-engineering methods for backdoor detection. 70 Through reversed Trojan transformation, we developed a Trojan removal method. We implemented 71 a prototype FEATURERE (FEATURE-space REverse-engineering) in Python and PyTorch and 72

r3 evaluated it on MNIST, GTSRB, CIFAR, and ImageNet dataset with seven different attacks (i.e.,

<sup>74</sup> BadNets [1], Filter attack [18], WaNet [6], Input-aware dynamic attack [7], ISSBA [8], Clean-label

<sup>75</sup> attack [24], Label-specific attack [1], and SIG attack [25]). Our results show that FEATURERE is <sup>76</sup> effective. On average, the accuracy of our method is 93%, outperforming existing techniques. For

Trojan mitigation, our method can reduce the ASR (attack success rate) to only 0.26% with the BA

(benign accuracy) remaining nearly unchanged by using only ten clean samples for each class.

Our contributions are summarized as follows. We first find the feature space properties of the Trojaned
model and reveal the relationship between Trojans and the feature space hyperplanes. We propose a
novel Trojan trigger reverse-engineering technique leveraging the feature space Trojan hyperplane.
We evaluate our prototype on four different datasets, five different network architectures, and seven
advanced Trojan attacks. Results show that our method outperforms SOTA approaches.

# **2** Background & Motivation

A DNN classifier is a function  $\mathcal{M}: \mathcal{X} \mapsto \mathcal{Y}$  where  $\mathcal{X}$  is the input domain  $\mathbb{R}^m$  and  $\mathcal{Y}$  is a set of labels 85 K. A Trojan (or backdoor) attack against a DNN model  $\mathcal{M}$  is a malicious way of perturbing the 86 input so that an adversarial input x' (i.e., input with the perturbation pattern) will be classified to a 87 target/random label while the model maintains high accuracy for benign input x. The perturbation 88 pattern is known as the Trojan trigger. Trojan attacks can happen in training (e.g., data poisoning) 89 or model distribution (e.g., changing model weights or supply-chain attack). Existing works have 90 shown Trojan attacks against different DNN models, including computer vision models [1, 22, 2], 91 Graph Neural Networks (GNNs) [26, 27], Reinforcement Learning (RL) [28, 29], Natural Language 92 Processing (NLP) [30–35], recommendation systems [36], malware detection [9], pretrained mod-93

els [37, 38, 19], active learning [39], and federated learning [40, 41]. The Trojan trigger can be a

simple input pattern (e.g., a yellow pad) [1, 22, 2] or a complex input transformation function (e.g., a
 CycleGAN to change the input styles) [5, 7, 8, 6, 3]. If the trigger is static input space perturbations

(e.g., a yellow pad), the Trojan attack is known as *input-space Trojan*, and if the trigger is an input

<sup>98</sup> feature (e.g., an image style), the attack is referred to as the *feature-space Trojan*.

There are different types of Trojan defenses. A line of work [42] attempts to remove poisoned data 99 samples by cleaning the training dataset. Training-based methods [43, 44] train benign classifiers 100 even with the poisoned dataset. Some other methods [12, 13, 45, 11] try to remove Trojans during 101 training with the help of benign datasets. These training time approaches work for poisoning-based 102 attacks but fail to defend against supply chain attacks where the adversary injects the Trojan after 103 the model is trained. Another line of work, e.g., STRIP [14], SentiNet [15], and Februus [16] aim to 104 detect Trojan inputs during runtime. It is hard to distinguish between a misclassification and a Trojan 105 attack for a test input. These runtime detection methods make assumptions about the attack, which 106 stronger attacks can violate. For example, STRIP fails to detect the Trojan inputs when the Trojan 107 trigger locates around the center of an image or overlaps with the main object (e.g., feature-space 108 attacks). Another limitation is that they examine the test inputs and perform various heavyweight 109 tests, significantly delaying the response time. 110

Trigger reverse engineering [17, 19, 18, 20, 21] makes no assumptions about the attack method (e.g., poisoning or supply-chain attacks) and does not affect the runtime performance. It inspects the model to check if a Trojan exists before deploying. Given a DNN model  $\mathcal{M}$  and a small set of clean samples  $\mathcal{X}$ , trigger reverse engineering methods try to reconstruct injected triggers. If reverse engineering is successful, the model is marked as malicious. Neural Cleanse (NC) [17] proposes to perform reverse engineering by solving Eq. 1:

$$\min_{\boldsymbol{m},\boldsymbol{t}} \quad \mathcal{L}\left(\mathcal{M}\left(\left(1-\boldsymbol{m}\right)\odot\boldsymbol{x}+\boldsymbol{m}\odot\boldsymbol{t}\right),y_{t}\right)+r^{\star} \tag{1}$$

where  $x \in \mathcal{X}$  and *m* is the trigger mask (i.e., a binary matrix with the same size as the input to 117 determine if the value will be replaced by the trigger or not), t is the trigger pattern (i.e., a matrix with 118 the same size as the input containing trigger values), and  $r^{\star}$  are attack constraints (e.g., trigger size is 119 smaller than 1/4 of the image).  $\mathcal{L}$  is the cross-entropy loss function. Most prior works [19, 18, 20, 21] 120 follow the same methodology and inherently suffer from the same limitations. First, they assume that 121 an input space perturbation, denoted by (m, t), can represent a trigger. This assumption is valid for 122 input-space triggers but does not hold for feature-space attacks. Second,  $r^*$  are heuristics observed 123 from existing attacks. For example, NC observed that most triggers have small sizes and limit the 124 trigger size to be no larger than one-fourth of the input. Otherwise, the trigger will overlap with 125 the main object and decrease benign accuracy. In practice, more advanced attacks can break such 126 heuristics. For instance, DFST leverages CycleGAN to transfer images from one style to another 127 without changing its semantics. It changes almost all pixels in a given image. This paper proposes a 128 novel reverse engineering method that overcomes the limitations above for image classifiers. 129

# 130 **3** Methodology

# 131 3.1 Threat Model

This work aims to determine if a given model has a Trojan or not by reverse-engineering the 132 corresponding trigger. Following existing works [17, 18, 43], we assume access to the model and a 133 small dataset containing correctly labeled benign samples of each label. In practice, such datasets 134 can be gathered from the Internet. We make no assumptions on how the attacker injects the Trojan 135 (poisoning or supply-chain attack). The attack can be formally defined as:  $\mathcal{M}(x) = y, \mathcal{M}(F(x)) =$ 136  $y_T, x \in \mathcal{X}$ , where  $\mathcal{M}$  is the Trojaned model, x is a clean input sample, and  $y_T$  is the target label. F 137 is the function to construct Trojan samples. Input-space triggers add static input perturbations, and 138 feature-space triggers are input transformations. The key difference between our work and existing 139 work is that we consider the feature space triggers. 140

#### 141 3.2 Observation

In DNNs, the neuron activation values represent its functionality. The input neurons denote the input space features, and inner neurons extract inner and more abstract features. Existing reverseengineering methods constrain the optimization problem in the input space using domain-specific

constraints or observations. For image classification tasks, the pixel value of each image has to be a 145 valid RGB value. Methods like NC observe that the trigger size must be smaller and cannot overlap 146 with the main object and propose corresponding constraints. The most challenging problem for 147 reverse-engineering feature-space triggers is how to constrain the optimization properly. Note that 148 there exist a set of neurons; when activating to specific values, the Trojan behavior will be triggered. 149 Due to the black-box nature of DNNs, it is hard to identify which neurons are related to the Trojan 150 behavior. Moreover, if enlarge the weight values with the same scale, the output of the DNN will be 151 the same, and as such, it is hard to constrain concrete activation values. Without a proper constraint, 152 we cannot form an optimization problem. 153

Our key observation to solve this problem is that neuron activation values representing the Trojan 154 behavior are orthogonal to others. Recall that one property of DNN Trojans is that when adding 155 the trigger to any given input, the model will predict the output to a specific label. That is, the 156 trigger will always work regardless of the actual contents of the input. In the feature space, when 157 the model recognizes features of the Trojan, it will predict the label to the target label regardless 158 of the other features. These activation values will form a hyperplane space in the high dimensional 159 space so that they can be orthogonal to all others. Based on this intuition, we performed empirical 160 experiments to confirm our idea. Specifically, we first use six Trojan attacks (e.g., BadNets [1], Clean 161 label attack [24], Filter attack [18], and WaNet [6], SIG [25] and Input-aware dynamic attack [7]) to 162 generate Trojaned ResNet18-architecture models on CIFAR-10. We then visualize the feature space 163 of the last convolutional layers in these models. In Fig. 2, three dimensions, X, Y, and Z, represent 164 the feature space. We first apply PCA to get two eigenvectors of the benign training set; then, we 165 166 use the obtained eigenvectors as X-axis and Y-axis. For Z-axis, we first construct Trojan inputs to activate the model's Trojan behavior and find highly related neurons to Trojans. Then, we use DNN 167 interpretability techniques SHAP [46] to estimate the neuron's importance to the Trojan behavior. 168 The neurons among the top 3% are compromised neurons. Z-axis denotes the average activation 169 values of compromised neurons. Namely,  $u_z = \frac{mean(\mathcal{A}(F(\mathcal{X}_{train})) \odot m)}{\|mean(\mathcal{A}(F(\mathcal{X}_{train})) \odot m)\|}$ , where m denotes a mask 170 revealing the position of compromised neurons. Fig. 2 show that most Trojan inputs have a similar 171 z-value. They form a linear hyperplane in the feature space while benign ones do not. 172



Fig. 2: Feature space of Trojaned models.

# 173 3.3 Feature Space Trojan Hyperplane Reverse-engineering

In this paper, We use  $\mathcal{A}$  to represent the submodel from the input space to the feature space.  $\mathcal{B}$  is the submodel from the feature space to the output space. We also use  $a = \mathcal{A}(x)$  to denote the inner features of the model. Similar to the reverse-engineering in the input space, given a model  $\mathcal{M}$  and a small set of benign inputs  $\mathcal{X}$ , we use a feature space mask m and a feature space pattern t to represent the feature space Trojan hyperplane  $H = \{a | m \odot a = m \odot t\}$ . Specifically, we can update m and t via the following optimization process:  $\min_{m,t} \mathcal{L} (\mathcal{B} ((1-m) \odot a + m \odot t), y_t), y_t)$  is the target

A 1 4/1 4		D 1 1	D	•	•
Algorithm I	Hantura change	Rockdoor	Davarca	anainaa	mma
	L'EATHE-SUALE		INCVEINE-		
THEOLIGINAL T	. I culuic spuce	Duchaool	100,0100	enginee	11115

Input: Model:  $\mathcal{M}$ **Output:** Trojaned or Not, Trojaned Pairs T 1: **function** REVERSE-ENGINEERING( $\mathcal{M}$ ) 2: for (target class  $y_t$ , source class  $y_s$ ) in K do 3: for  $e \leq E$  do  $\bar{x} = sample(\mathcal{X}_{y_s})$ 4:  $cost_1 = \mathcal{L}\left(\mathcal{B}\left(\left(1 - \boldsymbol{m}\right) \odot \boldsymbol{a} + \boldsymbol{m} \odot \boldsymbol{t}\right), y_t\right)$ 5: if  $||F(\boldsymbol{x}) - \boldsymbol{x}|| \geq \tau_1$  then 6:  $cost_1 = cost_1 + w_1 \cdot ||F(\boldsymbol{x}) - \boldsymbol{x}||$ 7: if  $std(\boldsymbol{m} \odot \mathcal{A}(F(\boldsymbol{x}))) \geq \tau_2$  then 8:  $cost_1 = cost_1 + w_2 \cdot std(\boldsymbol{m} \odot \mathcal{A}(F(\boldsymbol{x})))$ 9:  $\begin{array}{l} \Delta_{\theta_{F}} = \frac{\partial cost_{1}}{\partial \theta_{F}} \\ \theta_{F} = \theta_{F} - lr_{1} \cdot \Delta_{\theta_{F}} \\ cost_{2} = \mathcal{L} \left( \mathcal{B} \left( \left( 1 - \boldsymbol{m} \right) \odot \boldsymbol{a} + \boldsymbol{m} \odot \boldsymbol{t} \right), y_{t} \right) \end{array} \end{array}$ 10: 11: 12: if  $\|\boldsymbol{m}\| \geq \tau_3$  then 13:  $cost_2 = cost_2 + w_3 \cdot \|\boldsymbol{m}\|$ 14:  $\Delta_{\boldsymbol{m}} = \frac{\overline{\partial_{cost_2}}}{\partial \boldsymbol{m}}$  $\boldsymbol{m} = \boldsymbol{m} - lr_2 \cdot \Delta_{\boldsymbol{m}}$ 15: 16: if  $ASR\left(\mathcal{B}\left(\left(1-\boldsymbol{m}\right)\odot\boldsymbol{a}+\boldsymbol{m}\odot\boldsymbol{t}\right),y_{t}\right)>\lambda$  then 17: 18:  $\mathcal{M}$  is a Trojaned model, 19:  $T.append((y_s, y_t))$ 

label. As discussed above, reverse-engineering the feature space is challenging. In the input space, 180 all values have natural physical semantics and constraints, e.g., a pixel value in the RGB value range. 181 Values in the feature space have uninterruptible meanings and are not strictly constrained. Whether 182 the result will have a physically meaningful semantic is also uncertain. We solve these challenges by 183 simultaneously optimizing the input space trigger function F and the feature space Trojan hyperplane 184 H to enforce that the trigger has semantic meanings. In detail, we compute the feature space 185 trigger pattern as the mean inner features on the samples generated by the trigger function, i.e., 186  $t = mean (m \odot \mathcal{A}(F(\mathcal{X})))$ . We also constrain the standard deviation of  $m \odot \mathcal{A}(F(\mathcal{X}))$  to make sure 187 the features generated by the trigger function will lie on the relaxation of the reverse-engineered 188 189 hyperplane. Formally, our reverse-engineering can be written as the constrained optimization problem shown in Eq. 2, where  $\mathcal{X}$  is the small set of clean samples. We use deep neural networks to model 190 the trigger function (i.e.,  $F = G_{\theta}$ ) because of their expressiveness [47, 21]. Specifically, following 191 DFST [3], we use a representative deep neural network UNet [48]. Given a model and a small 192 set of clean inputs, the trigger function can be smoothly reconstructed via gradient-based methods, 193 i.e., optimizing the generative model  $G_{\theta}$ . In our default setting,  $\mathcal{A}$  and  $\mathcal{B}$  are separated at the last 194 convolusional layer. More discussions are in the Appendix (§ A.6). 195

$$\min_{F,\boldsymbol{m}} \mathcal{L} \left( \mathcal{B} \left( (1-\boldsymbol{m}) \odot \boldsymbol{a} + \boldsymbol{m} \odot \boldsymbol{t} \right), y_t \right)$$
  
where  $\boldsymbol{t} = mean \left( \boldsymbol{m} \odot \mathcal{A}(F(\mathcal{X})), \boldsymbol{a} \in \mathcal{A}(\mathcal{X}) \right)$   
s.t.  $\|F(\mathcal{X}) - \mathcal{X}\| \leq \tau_1, std(\boldsymbol{m} \odot \mathcal{A}(F(\mathcal{X}))) \leq \tau_2, \|\boldsymbol{m}\| \leq \tau_3$  (2)

There are several constraints in the optimization problem: (I) The transformed samples should 196 be similar to the original image due to the properties of Trojan attacks, i.e.,  $||F(x) - x|| \le \tau_1$ . 197 Typically, the Trojan samples are visually similar to original samples for stealthy purposes. In detail, 198 we use MSE (Mean Squared Error) to calculate the distance between F(x) and x. (2) The Trojan 199 features should lie in the relaxation of the reverse-engineered feature space Trojan hyperplane, i.e., 200  $\mathbb{P}(a \in H^* | x \in F(\mathcal{X}))$  should have high values. To achieve this goal, we constrain the standard 201 deviation of different Trojan samples' activation values on each pixel in the hyperplane. (3) Similar to 202 input space trigger reverse-engineering [18], we set a bound for the size of the feature space trigger 203 mask, i.e.,  $\|\boldsymbol{m}\| \leq \tau_3$ . Here  $\tau_1, \tau_2$ , and  $\tau_3$  are threshold values. We discuss their influence in § 4.4. 204 The detailed reverse-engineering algorithm can be found in Algorithm 1, where K is a set containing 205 all possible (source class, target class) pairs of the model. FEATURERE scans all labels to identify 206

the Trojan target labels.  $w_1$ ,  $w_2$  and  $w_3$  are the weight values used in the optimization. Following NC [17], we adjust them dynamically during optimization to make sure the reverse-engineered Trojan satisfies the constraints. E is the maximal epoch number. In lines 5-11, it optimizes the trigger function F and then the mask m of the feature space hyperplane in lines 12-16. In the end, we determine the reverse-engineering is successful and the label  $y_t$  is a Trojan target label if the attack success rate of the reversed Trojan is above a threshold value  $\lambda$  (80% in this paper).

# 213 3.4 Trojan Mitigation

After we reverse-engineered the Trojans, we can mitigate it by breaking the reverse-engineered feature space Trojan hyperplane. Based on our observation, the neurons in the feature space Trojan hyperplane are highly related to the Trojan behaviors. Thus, we can mitigate the Trojans by breaking the hyperplane. Inspired by Zhao et al. [49], we can break it by flipping the neurons on it. Our neuron-flip process can be written as Eq. 3, where m is the reverse-engineered feature space mask, ais the inner features.  $a_i$  is the activation value on the  $i^{th}$  neuron.

$$Flip(\boldsymbol{a}) = \begin{cases} -\boldsymbol{a}_i, \text{ when } \boldsymbol{a}_i \text{ in } \boldsymbol{m} \\ \boldsymbol{a}_i, \text{ when } \boldsymbol{a}_i \text{ not in } \boldsymbol{m} \end{cases}$$
(3)

The mitigated model  $\mathcal{M}'(\boldsymbol{x}) = \mathcal{B}(Flip(\mathcal{A}(\boldsymbol{x})))$ , where  $\mathcal{A}$  and  $\mathcal{B}$  are submodels of the model.

# **221 4 Experiments and Results**

We first introduce our experiment setup ( 4.1). We then evaluate the effectiveness of FEATURERE on

Trojan detection ( $\S$  4.2) and mitigation tasks ( $\S$  4.3). We also evaluate the robustness of FEATURERE

against different settings and the impacts of configurable parameters in FEATURERE (§ 4.4).

#### 225 4.1 Experiment Setup.

We implement FEATURERE with python 3.8 and PyTorch. All experiments are conducted on a Ubuntu 18.04 machine equipped with 64 CPUs and six GeForce RTX 6000 GPUs.

Datasets and Models. We use four publicly available 228 datasets to evaluate FEATURERE, including MNIST [50], 229 GTSRB [51], CIFAR-10 [52] and ImageNet [53]. We 230 summarize our datasets in Table 1. We show the dataset 231 names, the size of each input sample, the number of sam-232 ples and the number of classes in each column. De-233 tails of the datasets can be found in Appendix. For 234 model architectures, we use LeNet5 [50], Preact ResNet18 235 (PRN18) [54], ResNet18 [55], a VGG-style network 236

# Table 1: Overview of datasets.

Dataset	Sample Size	#Train	Classes
MNIST	32*32*1	60000	10
GTSRB	32*32*3	39209	43
CIFAR-10	32*32*3	50000	10
ImageNet	224*224*3	100000	200

specified in ULP [44], and a model consists of 4 convolutional layers and 2 dense layers
used in Xu et al. [43]. These datasets and models are widely used in Trojan-related researches [1, 2, 17, 18, 56, 14, 3, 43].

240 **Evaluation Metrics.** We measure the effectiveness of the Trojan detection task by collecting the detection accuracy (Acc). Given a set of models consist of benign and Trojaned models, the Acc 241 is the number of correctly classified models over the number of all models. We also show detailed 242 number of True Positives (TP, i.e., correctly detected Trojaned models), False Positives (FP, i.e., 243 benign models classified as Trojaned models), False Negatives (FN, i.e., Trojaned models classified 244 as benign models) and True Negatives (TN, i.e., correctly classified benign models). For the Trojan 245 246 mitigation task, we evaluate the benign accuracy (BA) and attack success rate (ASR) [57]. BA is the number of correctly classified clean inputs over the number of all clean samples. ASR is defined as 247 the number of Trojan samples that successfully attack models over the number of all Trojan samples. 248

Baselines and Attack Settings. We evaluate the performance of FEATURERE on Trojan detection,
and compare the results with four reverse-engineering based Trojan detection methods (i.e., ABS [18],
DeepInspect [21], TABOR [20], and K-arm [19]) and two classification based methods (i.e., ULP [44]
and Meta-classifier [43]). For Trojan mitigation task, we compare FEATURERE with two advanced
mitigation methods (i.e., NAD [58] and I-BAU [59]). We use the default parameter settings described

Dataset	Network	Attack			AB	S		_	De	eepIn	spec	xt	_	1	ГАВ	OR		_		K-a	ırm			FE.	ATUI	reR	E
			TP	FP	FN	TN	Acc	TP	FP	FN	ΤN	Acc	TP	FP	FN	TN	Acc	TP	FP	FN	TN	Acc	TP	FP	FN	TN	Acc
MNIST	LeNet5	WaNet	7	2	3	8	75%	4	0	6	10	70%	3	2	7	8	55%	5	0	5	10	75%	9	1	1	9	90%
GTSRB	PRN18	WaNet	5	0	5	10	75%	5	1	5	9	70%	2	2	8	8	50%	4	0	6	10	70%	8	0	2	10	90%
		BadNets	18	0	2	20	95%	20	2	0	18	95%	20	3	0	17	93%	20	0	0	20	100%	20	1	0	19	98%
CIEAD 10	D N-+10	Filter	13	0	7	20	83%	6	2	14	18	60%	5	3	15	17	55%	0	0	20	20	50%	18	1	2	19	93%
CIFAR-10	Resilet18	WaNet	11	0	9	20	78%	11	2	9	18	73%	3	3	17	17	50%	9	0	11	20	73%	18	1	2	19	93%
		IA	3	0	17	20	58%	4	2	16	18	55%	3	3	17	17	50%	2	0	18	20	55%	19	1	1	19	95%

Table 2: Comparison to reverse-engineering methods.

Table 3: Comparison to ULP.

Table 4: Comparison to Meta-classifier.

Network	Attack	_		UL	Р		_		Ou	rs		Network		Meta Classifier			Ours						
		TP	FP	FN	ΤN	Acc	TP	FP	FN	TN	Acc			TP	FP	FN	TN	Acc	TP	FP	FN	TN	Acc
VGG	WaNet	1	0	19	20	0.53	17	0	3	20	93%	4Conv+2FC	WaNet	16	4	4	16	0.80	18	0	2	20	95%

in the original papers of our baseline methods. To understand the performance of FEATURERE
and existing methods against various attack settings, we evaluate them against BadNets [1], Filter
Trojans [18], WaNets [6], IA (Input-dependent dynamic Trojans) [7], Clean-label [24], SIG [25]
and ISSBA (Invisible sample-specific Trojans) [8] attacks. These attacks are state-of-the-art attack
methods and are widely evaluated in Trojan defense papers [17, 18, 60, 59] If not specified, we use
an all-to-one (i.e., single-target) setting for all attacks. Label-specific setting is discussed in § 4.4.

# 260 4.2 Effectiveness on Trojan Detection

To measure the effectiveness on the Trojan detection task, we generate a set of benign and Trojaned 261 models, and then use FEATURERE and existing Trojan detection methods to classify each model. We 262 collect the Acc, TP, FP, FN and TN results of each method and compare them. Specifically, we first 263 evaluate the performance of FEATURERE and compare the results with four state-of-the-art reverse-264 engineering based detection methods. We generate 20 Trojaned models as well as 20 benign models 265 on CIFAR-10 dataset for each attack (i.e., BadNets, Filter Trojan, WaNet and Input-aware dynamic 266 Trojan attack). For MNIST and GTSRB dataset, we train 10 Trojaned and 10 benign LeNet5 [50] 267 models on each dataset. We then compare FEATURERE with two state-of-the-art classification 268 based detection methods. Similarly, we generate 10 benign and 10 Trojaned models, and use Trojan 269 detection methods to classify these models. Notice that, in all Trojan detection tasks, we assume the 270 defender can only access 10 clean samples for each class, which is a common practice. [17, 19, 18] 271 The comparison results of reverse-engineering based methods are shown in Table 2. The results of 272 two classification based methods are demonstrated in Table 3 and Table 4. In each table, we show the 273 detailed settings, including dataset names, network architectures, and attack settings. 274

**Comparison to Reverse-engineering based methods.** From the results in Table 2, we observe 275 that FEATURERE achieves the best detection results compared with other methods. The average 276 Acc of FEATURERE is 93%, which is 17%, 23%, 35% and 23% higher than those of other defense 277 methods. The results show the benefit of FEATURERE. When looking into the generalization of 278 Trojan detection methods, we find that FEATURERE can achieve excellent results on both input-space 279 Trojans (i.e., BadNets) and feature-space Trojans (i.e., Filter, WaNet and IA attacks). However, the 280 performance of existing reverse-engineering methods on feature-space Trojans (i.e., Filter, WaNet 281 and IA attacks) is significantly worse than the performance on static Trojans. FEATURERE archives 282 94% average Acc but the Acc of TABOR on feature-space Trojans are only 53%, 50% and 50%, 283 respectively. Moreover, FEATURERE has 15.33 TP on average, but existing methods only have 284 7.87 TP. FEATURERE can generalize better than existing work because FEATURERE considers both 285 feature and input space constraints. Existing methods, on the contrary, only consider the input space 286 constraints. They can not detect feature-space Trojans whose trigger is complex and input-dependent, 287 and directly classify many Trojaned models with feature-space Trojans as benign. 288

**Comparison to classification based methods.** When comparing FEATURERE with classification based methods, we notice that FEATURERE has better Acc, more TPs and TNs than classification based methods ULP and Meta-classifier. As demonstrated in Table 3 and Table 4, the Acc of

Table 5: Results on Trojan mitigation task (10 clean samples for each class are used).

Dataset	Network	Attack	Undef	Undefended		AD	I-B	AU	Ou	rs
			BA	ASR	BA	ASR	BA	ASR	BA	ASR
MNIST	LeNet5	WaNet	99.22%	94.52%	65.87%	48.21%	94.87%	0.22%	99.20%	0.63%
GTSRB	PRN18	WaNet	99.02%	99.70%	70.35%	62.76%	91.74%	0.86%	98.42%	0.00%
CIFAR-10	ResNet18	Filter WaNet IA	91.30% 91.84% 91.62%	98.98% 98.17% 92.44%	81.66% 83.60% 84.03%	28.23% 27.52% 34.00%	87.45% 87.52% 86.88%	18.22% 6.84% 10.33%	91.26% 91.79% 91.43%	0.29% 0.04% 0.38%

FEATURERE is 0.93 and 0.95, which is 0.40 and 0.15 higher than those of ULP and Meta-classifier. Overall, the results indicate that FEATURERE is more effective than classification based methods when detecting Trojaned models. Different from FEATURERE, which directly inspects models via analyzing its inherent feature space properties, classification based methods highly depend on the external trained dataset. Therefore, their results are not as precise as FEATURERE.

#### 297 4.3 Effectiveness on Trojan Mitigation

We evaluate the effectiveness of FEATURERE on Trojan mitigation and compare the results with state-of-the-art methods NAD and I-BAU. We use the Trojaned models generated by three attacks (i.e., Filter attack, WaNet and IA) and report their average BA and ASR after Trojan mitigation. We also show the average BA and ASR of undefended Trojaned models. For all methods, the defenders can access 10 clean samples for each class to conduct Trojan mitigation. We show the results in Table 5.

We find that FEATURERE is the most effective method for Trojan mitigation among all methods. 303 Compared to state-of-the-art Trojan mitigation methods, FEATURERE archives the lowest average 304 ASR and the highest average BA. On the one hand, using FEATURERE can decrease the average 305 ASR from 96.76% to 0.26%. Other methods can only decrease the average ASR to 40.14% and 306 7.29%. The results show the advantages of FEATURERE on Trojan mitigation. On the other hand, the 307 BA with FEATURERE is similar to undefended models. But the BA of other methods is significantly 308 lower than that of undefended models. By breaking the feature space hyperplane, FEATURERE can 309 successfully mitigate Trojans with minimal BA loss. Other methods, which cannot find Trojan-related 310 features, cannot achieve good results. 311

## 312 4.4 Ablation Study

In this section, we evaluate the resistance of FEATURERE to various Trojan attack settings and large datasets. We also evaluate the impacts of configurable parameters in FEATURERE, including the constrain values used in Eq. 2 and the number of used clean samples. By default, the attack used for measuring the impacts of configurable parameters is IA. We use 20 benign ResNet models and 20 Trojaned ResNet models on CIFAR-10 to test the detection results. Notice that we only evaluate the performance on the Trojan detection task. Due to the page limits, we include the ablation study on Trojan mitigation in Appendix.

#### 320 Resistance to various attack and dataset settings.

To evaluate if our method is resistant to more Tro-321 jan attacks, we train 20 Trojaned ResNet18 models 322 on CIFAR-10 for Label-specific attack (LS), Clean-323 label attack (CL) and SIG attack (SIG). For the label-324 specific attack, we consider the all-to-all attack set-325 ting, i.e., the target label  $y_T = \eta(y) = y + 1$ , where  $\eta$ 326 is a mapping and y is the correct label of the sample. 327 In addition, we generate five benign models and five 328 Trojaned models with ISSBA attacks on ImageNet to 329

Table 6: Resistance to more attacks.

Dataset	Network	Attack	TP	FP	FN	TN	Acc
CIFAR-10	ResNet18	LS CL	9 8	1	1	9 9	90% 85%
	Residento	SIG	10	1	õ	9	95%
ImageNet	ResNet18	ISSBA	4	0	1	5	90%

evaluate if our method is compatible with large-scale datasets. We summarize the results in Table 6.

In Table 6, we find that FEATURERE is compatible with evaluated Trojan attacks, showing the generalization of our reverse-engineering based method. We also observe that our method has high Acc on the ImageNet dataset with ISSBA [8]. Thus, our method is also applicable to large datasets. **Influence of constrain values.** As shown in Eq. 2, there are three constrain values  $(\tau_1, \tau_2, \tau_3)$  in our constrained optimization process. By default,  $\tau_1 = 0.15$ ,  $\tau_2 = 0.25$  and  $\tau_3 = 5\%$ . We evaluate their influences. For  $\tau_1$ , we calculate input space perturbations on the preprocessed inputs, and the details of the preprocessing can be found in Appendix (§ A.2). We vary  $\tau_1$  from 0.05 to 0.35, change  $\tau_2$ from 0.10 to 0.50, and tune  $\tau_3$  from 3% of the whole feature space to 10% of the whole feature space. The results under different hyperparameter settings are shown in Table 7.

From the results, we observe that the performance of FEATURERE is insensitive to these three hyperparameters. In detail, when we vary  $\tau_1$ ,  $\tau_2$  and  $\tau_3$ , the Acc is stable. In all cases, our method always achieves over 90% detection accuracy. The results further show the robustness of FEATURERE. We also find that, when the value of all hyperparameters becomes lower, FEATURERE has more FN. On the contary, when its value is larger, more FP will be produced. This is understandable because lower constrain values mean a stricter criterion for a successful reverse-engineering.









Number of clean reference samples. Our threat model and existing work assume the defender can access a set of clean samples for defense. To investigate the influences of the number of used clean samples in Trojan detection, we choose the number from 1 to 100 in each class and report the Acc results. The results are shown in Fig. 3.

From the results, we notice that the Acc decreases significantly when we use less than 10 samples 351 for each class. This is because the number of used sample affects the optimization process. When 352 the number of used samples is too small, the optimization process might be problematic, e.g., it 353 encounters overfitting problem. When the number of used samples is larger than 10, FEATURERE 354 achieves high detection accuracy (i.e., above 95%) and the Acc will not change significantly when 355 the number of used samples keeps increasing. The reason is using more data makes the optimizing 356 process converge and finally arrives a stable state. Note that requiring hundreds clean samples is 357 common for reverse-engineering based methods [17, 18, 20] and other types of defenses [58–60, 14]. 358 FEATURERE only requires 10 clean samples for each class, which is more efficient. 359

# 360 5 Discussion

346

Limitations of our method. Similar to most existing Trojaned model detection and mitigation methods [17–21, 58, 56], our method requires a small set of clean samples. In the real world, these samples can be obtained from the Internet.

**Ethics.** This paper proposes a technique to detect and remove Trojans in DNN models. We believe it will help improve the security of DNNs and be beneficial to society.

# 366 6 Conclusion

In this paper, we find relationships between feature space hyperplane and Trojans in DNNs. More over, we propose a new Trojaned DNN detection and mitigation method based on our findings. Compared to the state-of-the-art methods, our method has better performance in both detection and mitigation tasks.

# 371 **References**

- [1] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- [2] Yingqi Liu, Shiqing Ma, Yousra Aafer, Wen-Chuan Lee, Juan Zhai, Weihang Wang, and
   Xiangyu Zhang. Trojaning attack on neural networks. 2017.
- [3] Siyuan Cheng, Yingqi Liu, Shiqing Ma, and Xiangyu Zhang. Deep feature space trojan attack
   of neural networks by controlled detoxification. *arXiv preprint arXiv:2012.11212*, 2020.
- [4] Khoa Doan, Yingjie Lao, Weijie Zhao, and Ping Li. Lira: Learnable, imperceptible and robust
   backdoor attacks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 11966–11976, 2021.
- [5] Ahmed Salem, Rui Wen, Michael Backes, Shiqing Ma, and Yang Zhang. Dynamic backdoor attacks against machine learning models. *arXiv preprint arXiv:2003.03675*, 2020.
- [6] Anh Nguyen and Anh Tran. Wanet–imperceptible warping-based backdoor attack. *arXiv preprint arXiv:2102.10369*, 2021.
- [7] Anh Nguyen and Anh Tran. Input-aware dynamic backdoor attack. *arXiv preprint arXiv:2010.08138*, 2020.
- [8] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor
   attack with sample-specific triggers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16463–16472, 2021.
- [9] Giorgio Severi, Jim Meyer, Scott Coull, and Alina Oprea. {Explanation-Guided} backdoor
   poisoning attacks against malware classifiers. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1487–1504, 2021.
- [10] Esha Sarkar, Hadjer Benkraouda, and Michail Maniatakos. Facehack: Triggering backdoored
   facial recognition systems using facial characteristics. *arXiv preprint arXiv:2006.11623*, 2020.
- [11] Min Du, Ruoxi Jia, and Dawn Song. Robust anomaly detection and backdoor attack detection
   via differential privacy. *arXiv preprint arXiv:1911.07116*, 2019.
- [12] 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.
- [13] Brandon Tran, Jerry Li, and Aleksander Madry. Spectral signatures in backdoor attacks. *arXiv preprint arXiv:1811.00636*, 2018.
- [14] 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.
- [15] Edward Chou, Florian Tramèr, Giancarlo Pellegrino, and Dan Boneh. Sentinet: Detecting
   physical attacks against deep learning systems. 2018.
- [16] Bao Gia Doan, Ehsan Abbasnejad, and Damith C Ranasinghe. Februus: Input purification
   defense against trojan attacks on deep neural network systems. In *Annual Computer Security Applications Conference*, pages 897–912, 2020.
- [17] 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), pages 707–723. IEEE, 2019.
- [18] 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.

- [19] 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 optimiza tion. *arXiv preprint arXiv:2102.05123*, 2021.
- [20] 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.
- 423 [21] Huili Chen, Cheng Fu, Jishen Zhao, and Farinaz Koushanfar. Deepinspect: A black-box trojan
   424 detection and mitigation framework for deep neural networks. In *IJCAI*, pages 4658–4664,
   425 2019.
- [22] 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.
- [23] 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*, pages 113–131, 2020.
- [24] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks.
   *arXiv preprint arXiv:1912.02771*, 2019.
- [25] Mauro Barni, Kassem Kallas, and Benedetta Tondi. A new backdoor attack in cnns by training
   set corruption without label poisoning. In 2019 IEEE International Conference on Image
   *Processing (ICIP)*, pages 101–105. IEEE, 2019.
- [26] Zhaohan Xi, Ren Pang, Shouling Ji, and Ting Wang. Graph backdoor. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1523–1540, 2021.
- [27] Zaixi Zhang, Jinyuan Jia, Binghui Wang, and Neil Zhenqiang Gong. Backdoor attacks to graph
   neural networks. In *Proceedings of the 26th ACM Symposium on Access Control Models and Technologies*, pages 15–26, 2021.
- [28] Panagiota Kiourti, Kacper Wardega, Susmit Jha, and Wenchao Li. Trojdrl: evaluation of
   backdoor attacks on deep reinforcement learning. In 2020 57th ACM/IEEE Design Automation
   *Conference (DAC)*, pages 1–6. IEEE, 2020.
- [29] 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.
- [30] Xiaoyi Chen, Ahmed Salem, Dingfan Chen, Michael Backes, Shiqing Ma, Qingni Shen,
   Zhonghai Wu, and Yang Zhang. Badnl: Backdoor attacks against nlp models with semantic preserving improvements. In *Annual Computer Security Applications Conference*, pages
   554–569, 2021.
- 451 [31] Alvin Chan, Yi Tay, Yew-Soon Ong, and Aston Zhang. Poison attacks against text datasets with 452 conditional adversarially regularized autoencoder. *arXiv preprint arXiv:2010.02684*, 2020.
- [32] Wenkai Yang, Lei Li, Zhiyuan Zhang, Xuancheng Ren, Xu Sun, and Bin He. Be careful about
   poisoned word embeddings: Exploring the vulnerability of the embedding layers in nlp models.
   *arXiv preprint arXiv:2103.15543*, 2021.
- [33] Fanchao Qi, Yangyi Chen, Xurui Zhang, Mukai Li, Zhiyuan Liu, and Maosong Sun. Mind
   the style of text! adversarial and backdoor attacks based on text style transfer. *arXiv preprint arXiv:2110.07139*, 2021.
- [34] Wenkai Yang, Yankai Lin, Peng Li, Jie Zhou, and Xu Sun. Rethinking stealthiness of backdoor
   attack against nlp models. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5543–5557, 2021.
- [35] Fanchao Qi, Mukai Li, Yangyi Chen, Zhengyan Zhang, Zhiyuan Liu, Yasheng Wang, and
   Maosong Sun. Hidden killer: Invisible textual backdoor attacks with syntactic trigger. *arXiv preprint arXiv:2105.12400*, 2021.

- 466 [36] Shijie Zhang, Hongzhi Yin, Tong Chen, Zi Huang, Quoc Viet Hung Nguyen, and Lizhen Cui.
   467 Pipattack: Poisoning federated recommender systems formanipulating item promotion. *arXiv* 468 *preprint arXiv:2110.10926*, 2021.
- [37] Yuanshun Yao, Huiying Li, Haitao Zheng, and Ben Y Zhao. Latent backdoor attacks on deep
   neural networks. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, pages 2041–2055, 2019.
- [38] Jinyuan Jia, Yupei Liu, and Neil Zhenqiang Gong. Badencoder: Backdoor attacks to pre-trained
   encoders in self-supervised learning. *arXiv preprint arXiv:2108.00352*, 2021.
- Iose Rodrigo Sanchez Vicarte, Gang Wang, and Christopher W Fletcher. {Double-Cross}
   attacks: Subverting active learning systems. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 1593–1610, 2021.
- [40] 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.
- [41] Chulin Xie, Keli Huang, Pin-Yu Chen, and Bo Li. Dba: Distributed backdoor attacks against
   federated learning. In *International Conference on Learning Representations*, 2019.
- [42] 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.
- [43] 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.
- [44] 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.
- [45] Jonathan Hayase, Weihao Kong, Raghav Somani, and Sewoong Oh. Defense against backdoor
   attacks via robust covariance estimation. In *International Conference on Machine Learning*,
   pages 4129–4139. PMLR, 2021.
- [46] Scott M Lundberg and Su-In Lee. A unified approach to interpreting model predictions.
   Advances in neural information processing systems, 30, 2017.
- [47] Kurt Hornik, Maxwell Stinchcombe, and Halbert White. Multilayer feedforward networks are universal approximators. *Neural networks*, 2(5):359–366, 1989.
- [48] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for
   biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015.
- [49] Yue Zhao, Hong Zhu, Kai Chen, and Shengzhi Zhang. Ai-lancet: Locating error-inducing
   neurons to optimize neural networks. In *Proceedings of the 2021 ACM SIGSAC Conference on Computer and Communications Security*, pages 141–158, 2021.
- [50] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning
   applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Johannes Stallkamp, Marc Schlipsing, Jan Salmen, and Christian Igel. Man vs. computer:
   Benchmarking machine learning algorithms for traffic sign recognition. *Neural networks*, 32:
   323–332, 2012.
- [52] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images.
   2009.
- [53] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
   Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual
   recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.

- [54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual
   networks. In *European conference on computer vision*, pages 630–645. Springer, 2016.
- [55] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image
   recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*,
   pages 770–778, 2016.
- [56] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against
   backdooring attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses*, pages 273–294. Springer, 2018.
- [57] Akshaj Kumar Veldanda, Kang Liu, Benjamin Tan, Prashanth Krishnamurthy, Farshad Khorrami,
   Ramesh Karri, Brendan Dolan-Gavitt, and Siddharth Garg. Nnoculation: broad spectrum and
   targeted treatment of backdoored dnns. *arXiv preprint arXiv:2002.08313*, 2020.
- [58] Yige Li, Nodens Koren, Lingjuan Lyu, Xixiang Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. *arXiv preprint* arXiv:2101.05930, 2021.
- [59] Yi Zeng, Si Chen, Won Park, Z Morley Mao, Ming Jin, and Ruoxi Jia. Adversarial unlearning
   of backdoors via implicit hypergradient. *arXiv preprint arXiv:2110.03735*, 2021.
- [60] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor
   learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 34, 2021.
- [61] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [62] Yingqi Liu, Guangyu Shen, Guanhong Tao, Zhenting Wang, Shiqing Ma, and Xiangyu Zhang.
   Complex backdoor detection by symmetric feature differencing. 2022.
- [63] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1765–1773, 2017.
- [64] Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition
   for visual explanation. In *Proceedings of the European Conference on Computer Vision (ECCV)*,
   pages 119–134, 2018.
- [65] David Bau, Jun-Yan Zhu, Hendrik Strobelt, Agata Lapedriza, Bolei Zhou, and Antonio Torralba.
   Understanding the role of individual units in a deep neural network. *Proceedings of the National Academy of Sciences*, 117(48):30071–30078, 2020.
- [66] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale
   image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- <sup>546</sup> [67] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint* <sup>547</sup> *arXiv:1605.07146*, 2016.

# 548 Checklist

552

553

557

- 549 1. For all authors...
- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
  - (b) Did you describe the limitations of your work? [Yes] See § 5.
  - (c) Did you discuss any potential negative societal impacts of your work? [Yes] See § 5.
- (d) Have you read the ethics review guidelines and ensured that your paper conforms to
   them? [Yes] See § 5.
- 556 2. If you are including theoretical results...
  - (a) Did you state the full set of assumptions of all theoretical results? [N/A]

558	(b) Did you include complete proofs of all theoretical results? [N/A]
559	3. If you ran experiments
560 561 562	<ul> <li>(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] See Abstract, § 4, and Appendix.</li> </ul>
563 564	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See § 4 and Appendix.
565 566	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No]
567 568	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] See § 4.
569	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
570	(a) If your work uses existing assets, did you cite the creators? [Yes] See § 4.
571	(b) Did you mention the license of the assets? [Yes] See Appendix.
572 573	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] See Abstract.
574 575	(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] See Appendix.
576 577	<ul> <li>(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] See Appendix.</li> </ul>
578	5. If you used crowdsourcing or conducted research with human subjects
579 580	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
581 582	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
583 584	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

# 585 A Appendix

**Roadmap:** More details of Algorithm 1 is introduced in § A.1. Then, we present more details of the datasets (§ A.2) and attacks (§ A.3) used in the experiments. We also perform an ablation study for the Trojan mitigation task in § A.4. In § A.5, we visualize our reversed Trojan. In § A.6, we discuss how to split the model. The adaptive attack can be found in § A.7. We also show the generalization (§ A.8) and the efficiency (§ A.9) of FEATURERE. Finally, we discuss the findings in the evaluation (§ A.10).

### 592 A.1 More details of Algorithm 1

In this section, we discuss more details of our Reverse-engineering Algorithm (Algorithm 1 in the main paper). Given a model  $\mathcal{M}$  and a small set of clean samples  $\mathcal{X}$ , the output of the algorithm is a flag indicating if the model is Trojaned, and Trojaned label pairs denoting the source label and the target label of the detected Trojans.

In line 2, we iterate (source label, target label) pair from possible pairs K. E in line 3 means the maximal optimization epoch number for each pair. It is set to 400 in this paper. In line 4, we randomly sample a batch of inputs from the samples in source classes. The batch size is set to 128 by default.

In lines 5 to Line 11, we optimize the parameters of the input space transformation F, which is 600 represented by a UNet [48] model in our implementation. In line 5, we calculate the loss value 601 specified in Eq. 1 (in the main paper), where  $a = \mathcal{A}(x)$  is the inner feature on clean samples. By 602 default,  $\mathcal{A}$  is the submodel from the input layer to the penultimate layer, and  $\mathcal{B}$  is the submodel from the 603 604 penultimate layer to the output layer. m is the feature space trigger mask.  $t = mean (m \odot \mathcal{A}(F(\mathcal{X})))$ is the feature space trigger pattern.  $\mathcal{L}$  is the cross-entropy loss calculating the distance between the 605 target label and the output of the model under inner features with feature space Trojans. In line 6, if 606 the input space MSE (Mean Square Error) distance for original inputs x and the transformed inputs 607  $F(\mathbf{x})$  is larger than a threshold value  $\tau_1$  (i.e., 0.15), then the regularization item  $w_1 \cdot ||F(\mathbf{x}) - \mathbf{x}||$ 608 will be added. Note that we calculate input space distance on the preprocessed inputs, and the details 609 of the preprocessing are in § A.2. Following NC [17], the weight value  $w_1$  is adjusted dynamically to 610 make the reverse-engineering satisfy the constrain (i.e.,  $||F(x) - x|| \le \tau_1$ ).  $w_2$  in line 9 and  $w_3$  in 611 line 14 are also adjusted dynamically. In lines 8-9, similarly, we add the regularization item for the 612 standard deviation of different Trojan samples' activation values on each pixel in the hyperplane. The 613 default value for  $\tau_2$  is 0.25. Lines 10-11 are the standard backward propagation process to update the 614 parameters of the input space transformation function F based on the gradients. The optimizer used 615 to optimize F is Adam [61]. The value of learning rate  $lr_1$  is 1e-3. In each epoch, we optimize both 616 the input space transformation function F and the feature space mask m. 617

Lines 12-16 describes the process for optimizing m. Similar to line 5, we calculate the cross-entropy loss between the target label and the output of the model under inner features with feature space Trojans in line 12. In lines 13-14, we add the regularization item for the size of the feature space Trojan hyperplane. The default value for  $\tau_3$  is 5% of the whole feature space. Lines 15-16 describe the process of updating feature space mask m via gradients. The value of learning rate  $lr_2$  in line 16 is 1e-1. The optimizer used is Adam [61].

In line 17, we check if the Trojan is successfully reverse-engineered. In detail, we calculated the ASR (attack success rate) on inner features with feature space Trojans (i.e.,  $(1 - m) \odot a + m \odot t$ ). We flag that reverse-engineering is successful if the ASR is above a threshold value  $\lambda$  (i.e., 0.8). If the Trojan is successfully reverse-engineered, we flag the model as a Trojan model and label the (source class, target class) pair as Trojaned pair. Besides the details above, we also use K-arm scheduler [19] to speed up the reverse engineering. Lastly, we use Liu et al. [62] to distinguish the Injected Trojans and UAPs (Universal Adversarial Patterns) [63].

#### 631 A.2 Details of Datasets

In this section, details of the datasets used in the experiments are discussed. We also provide the details of the preprocessing for each dataset. All datasets are open-sourced. The license for all datasets is the MIT license. They do not contain any personally identifiable information or offensive content.

Table 8: Details of Mean and Std value on each dataset.

Dataset	Mean	Std
MNIST	[0.1307]	[0.3081]
CIFAR-10	[0.4914, 0.4822, 0.4465]	[0.2023, 0.1994, 0.2010]
GTSRB	[0.3403, 0.3121, 0.3214]	[0.2724, 0.2608, 0.2669]
ImageNet	[0.4850, 0.4560, 0.4060]	[0.2290, 0.2240, 0.2250]

**MNIST** [50]. This dataset is used for classifying hand-written digits. It contains 60000 training samples in 10 classes. The number of samples in the test set is 10000.

**GTSRB** [51] This dataset is built for traffic sign classification tasks. The number of classes is 43. The sample numbers for the training set and test set are 39209 and 12630, respectively.

CIFAR10 [52] This dataset is used for recognizing general objects, e.g., dogs, cats, and planes. It has
 50000 training samples and 10000 training samples. This dataset has 10 classes.

ImageNet [53] This dataset is also a general object classification benchmark. Note that we use a
subset (containing 200 classes) of the original ImageNet dataset specified in ISSBA [8]. The subset
has 100000 training samples and 10000 test samples.

Following standard convention on the image classification task, we scale the inputs to the range [0,1] and use mean-std normalization to preprocess the images. In detail, the preprocessing can be written as  $\mathbf{x}' = \frac{\left(\frac{\mathbf{x}}{255} - Mean\right)}{Std}$ , where  $\mathbf{x}'$  is the normalized input and  $\mathbf{x}$  is the original inputs. The Mean value and Std (Standard Deviation) value for each channel on different datasets are summarized in Table 8.

# 649 A.3 Details of Attacks

In this section, we discuss the details of the used attacks. By default, the attacks are in an all-to-one (i.e., single-target) setting, and the target label is randomly selected when we generate Trojaned models.

**BadNets** [1]. This attack uses a fixed pattern (i.e., a patch or a watermark) as Trojan triggers, and it 653 generates Trojan inputs by simply pasting the pre-defined trigger pattern on the input. It compromised 654 the victim models by poisoning the training data (i.e., injecting Trojan samples and modifying their 655 labels to target labels). In our experiments, we use a 3\*3 yellow patch located at the left-upper corner 656 as Trojan trigger. The poisoning rate we used is 5%. The attack can be all-to-one (i.e., single-target) 657 and all-to-all (i.e., label-specific). For an all-to-one attack, all Trojan samples have the same target 658 label. For label-specific attacks, the samples in different original classes have different target labels. 659 In our experiment, the target label for label-specific attack is  $y_T = \eta(y) = y + 1$ , where  $\eta$  is a 660 mapping and y is the correct label of the sample. 661

Filter Attack [18]. This attack exploits image filters as triggers and creates Trojan samples by
 applying selected filters on images. Similar to BadNets, the Trojans are injected with poisoning.
 Following ABS [18], we use a 5% poisoning rate and apply the Nashville filter from Instagram as the
 Trojan trigger.

WaNet [6]. This method achieves Trojan attacks via image warping techniques. The trigger transformation of this attack is an elastic warping operation. Different from BadNets and Filter Attack, in this attack, the adversary needs to modify the training process of the victim models to make the attack more resistant to Trojan defenses. It is stealthy to human inspection, and it can also bypass many existing Trojan defense mechanisms [14, 12, 56, 17]. In our experiments, the wrapping strength and the grid size are set to 0.5 and 4, respectively.

**Input-aware Dynamic Attack [7].** This attack generates Trojan triggers via a trained generator network. The trigger generator is trained on a diversity loss so that two different input images do not share the same trigger. Similar to WaNet [6], the attacker needs to control the training process.

**SIG** [25]. This method uses superimposed sinusoidal signals as Trojan triggers. In this attack, the attacker can only poison a set of training samples but can not control the full training process. We set the poisoning rate as 5%. The frequency and the magnitude of the backdoor signal in our experiments are 6 and 20, respectively.

Table 9: Influence of hyperparameters on Trojan mitigation task.

Metric		$ au_1$			$\tau_2$					$ au_3$			
	0.05	0.15	0.35	0.10	0.25	0.50	1%	2%	3%	4%	5%	6%	7%
BA ASR	91.77% 0.02%	91.79% 0.04%	91.79% 0.08%	91.76% 0.02%	91.79% 0.04%	91.80% 0.08%	91.92% 57.75%	91.87% 0.50%	91.85% 0.06%	91.82% 0.06%	91.79% 0.04%	91.65% 0.00%	90.08% 0.00%

Table 10: Effects of clean set size on Trojan mitigation task.

Samples Per Class	BA	ASR
5	91.03%	0.08%
10	91.79%	0.04%
50	91.66%	0.02%
100	91.66%	0.06%

Clean Label Attack [24]. This attack poisons the datasets without manipulating the label of poisoning 679 samples so that the attack is more stealthy. The poisoning samples are generated by a trained GAN. 680 In our experiments, we set the poisoning rate as 5%. 681

**ISSBA** [8]. This attack utilizes an encoder-decoder network to generate sample-specific triggers. 682 The generated triggers are invisible noises. The generated noises also contain the information of a 683 representative string of the target label. The threat model of this attack is that the attacker can only 684 poison the training data, but can not control other components in training (e.g., the loss function). 685 Following the original paper, we poison 10% training data in our experiments. 686

#### A.4 Ablation Study on Trojan Mitigation 687

In this section, we study the performance of FEATURERE under different constrain values and 688 different numbers of used clean samples. The attack used in this section is WaNet [6]. 689

**Influence of constrain values.** To investigate the influence constrain values (i.e.,  $\tau_1$ ,  $\tau_2$ , and  $\tau_3$ ) on 690 the Trojan mitigation performance, We vary  $\tau_1$  from 0.05 to 0.35, change  $\tau_2$  from 0.10 to 0.50, and 691 tune  $\tau_3$  from 1% of the whole feature space to 7% of the whole feature space. We collect the BA 692 and ASR of the mitigated models and report them in Table 9. The results show that the mitigation 693 performance of FEATURERE is not sensitive to  $\tau_1$  and  $\tau_2$ . For  $\tau_3$ , when the size of the Trojan 694 hyperplane is extremely small (e.g., 1% of the feature space), the ASR is high. This is understandable 695 because breaking an extremely small feature space Trojan hyperplane means flipping a very small 696 number of neurons, and it is not enough to completely remove the Trojans in the model. Therefore, 697 we set the default value of the hyperplane's size as 5% of the feature space. 698

Number of clean reference samples. To understand the influence of clean set size on the Trojan 699 mitigation task, we vary the number of used clean samples from 5 per class to 100 per class and 700 report the BA and ASR of mitigated model. The results in Table 10 demonstrate that the performance 701 of FEATURERE is robust when the number of used samples changes. 702

#### A.5 Visualization of Reverse-Engineered Trojans 703

Table 11: Accuracy on different split position.

To understand our method and study if it can reverse-engineer Trojans accurately, we visualize the 704 inputs and inner features of clean samples, real Trojan samples, and reversed Trojan samples on nine 705 randomly selected samples in Fig. 4. The model is ResNet18 injected with Filter Trojan [18] and 706 Blend Trojan [22]. In the feature space, the reverse-engineered Trojan is close to the real Trojan, 707 demonstrating the effectiveness of our reverse-engineering method. 708

	Attack	9 <sup>th</sup>	$11^{\text{th}}$	$13^{\text{th}}$	$15^{\text{th}}$	Last
709	BadNets	88%	93%	98%	98%	98%
	Filter	88%	85%	90%	95%	93%
	WaNet	85%	88%	85%	93%	93%
710	IA	85%	85%	85%	90%	95%

Table 12: Results on BadNets and Adaptive Attack.

Attack	BA	ASR	Detection Accuracy
BadNets	94.34%	99.98%	98%
Adaptive	87.36%	93.67%	65%

710



Fig. 4: Visualization of the input space and the feature space for original inputs, real Trojan inputs, and reverse-engineered Trojan inputs.

#### 711 A.6 Discussion for Model Split

As we discussed in § 3, our method split the model  $\mathcal{M}$  into two sub-models  $\mathcal{A}$  and  $\mathcal{B}$ . In this section, we discuss the influence of using different split positions. Table 11 shows the results of using different  $\mathcal{A}$  and  $\mathcal{B}$  on the ResNet18 model and CIFAR-10 dataset. In detail, we report the results of splitting the model at the 9<sup>th</sup>, 11<sup>th</sup>, 13<sup>th</sup>, 15<sup>th</sup>, and the last convolutional layer. The average detection accuracy for splitting at the 9<sup>th</sup> layer, 11<sup>th</sup> layer, 13<sup>th</sup> layer, 15<sup>th</sup> layer, and last layer is 86.50%, 87.75%, 89.50%, 94.00%, and 94.75%, respectively. As we can see, the performance of splitting at later layers is higher than the performance of splitting at earlier layers.

In our current implementation, we set A(x) as the sub-model from the input layer to the last convolution layer and B(x) as the rest. The relationship between the input and the output of a convolutional layer  $L_n$  is  $x_{n+1} = L_n(x_n) = \sigma(\mathbf{W}_n^T x_n + \mathbf{b}_n^T)$ , where  $x_n$  and  $x_{n+1}$  are the inputs and outputs of layer n,  $\mathbf{W}_n$  and  $\mathbf{b}_n$  are weights and bias values, and  $\sigma$  is the activation function. Based on existing literatures [64, 65], the features in the deeper CNN layers are more disentangled than that of earlier layers. Thus, if the orthogonal phenomenon happens in a layer  $L_n$ , it will exist for all its subsequent layers, e.g.,  $L_{n+1}$ . If the orthogonal phenomenon does not happen, the layer without this phenomenon will mix benign and backdoor features, leading to low benign accuracy or attack success rate. The results in Table 12 confirm our analysis. Thus, a successful backdoor attack
 will lead to the orthogonal phenomenon in the last convolution layer.

#### 729 A.7 Adaptive Attack

Our threat model assumes that the attacker can control the training process of the Trojan model. In 730 this section, we discuss the potential adaptive attacker that knows our defense strategy and tries to 731 bypass FEATURERE via modifying the training process. Our observation is that the neuron activation 732 values representing the Trojan behavior are orthogonal to others. One possible adaptive attack is 733 breaking such orthogonal relationships during the Trojan injection process. We design an adaptive 734 attack that adds one loss term to push the Trojan features to be not orthogonal to benign features. 735 This attack can be formulated as:  $L = L_{ce} + L_{adv}$ , where  $L_{ce}$  is the standard classification loss and 736 the  $L_{adv}$  is defined as: 737

$$L_{adv} = sim(\mathcal{B}(m \odot a + (1 - m) \odot t), \mathcal{B}(m \odot a' + (1 - m) \odot t))$$
(4)

Here, sim is the cosine similarity; a and a' are the features of different benign samples; m and t are 738 the feature-space mask and pattern of the compromised neurons obtained via SHAP [46]. The loss 739 term  $L_{adv}$  tries to enforce the Trojan features being not orthogonal to the benign ones. We conduct 740 this adaptive attack on the CIFAR-10 dataset and ResNet18 model. The results can be found in 741 Table 12. The detection accuracy of FEATURERE under adaptive attack drops to 65%. Meanwhile, 742 the average BA/ASR of the adaptive attack and BadNets (native training) is 87.36%/94.34% and 743 93.67%/99.98%, respectively. The adaptive attack can reduce the detection accuracy of our method. 744 Both the BA and ASR of the adaptive attack are lower than those of native training. The results 745 confirm our analysis in § A.6: the model without the "orthogonal phenomenon" will mix benign and 746 Trojan features, leading to low benign accuracy or attack success rate. 747

#### 748 A.8 Generalization

Performance on mitigation task for more attacks. To measure the effectiveness of FEATURERE on 749 Trojan mitigation task, we use more Trojan attacks and report BA and ASR of our method. Besides 750 the results of BadNets [1], Filter [18], WaNet [6] and IA [7] in Table 5, in Table 13, we also show 751 the BA and ASR on LS [1], CL [24] and SIG [25]. The dataset and the model used is CIFAR-10 752 and ResNet18, respectively. For LS, CL, and SIG, the ASR of FEATURERE is 1.15%, 2.62%, and 753 1.22%, which are 80.01, 33.18, and 81.22 times lower than that of undefended models. As can be 754 observed, FEATURERE can effectively reduce the ASR while keeping the BA nearly unchanged. 755 Thus, FEATURERE is robust to different attacks on mitigation task. 756

Attack VGG16 ResNet18 PRN18 Undefended Ours BadNets 95% 95% 100% Attack Filter 90% 90% 95% BA ASR BA ASR 757 WaNet 90% 95% 90% LS 93.66% 92.02% 92.86% 1.15% 90% 90% 90% IA 93.51% 86.94% 92.94% 2.62% CL 85% LS 90% 85% 93.73% 93.47% SIG 99.09% 1.22% CL 80% 85% 85% SIG 95% 95% 90%

Table 13: Mitigation Results for More Attacks.

Table 14: Detection Accuracy on More Models.

758

Generalization to different models. To understand the generalization of FEATURERE to different 759 model architectures, we evaluate its detection accuracy on BadNets [1], Filter [18], WaNet [6], IA [7], 760 LS [1], CL [24], and SIG [25] attacks using VGG16 [66], ResNet18 [55], and Preact-ResNet18 761 (PRN18) [54]. The results are summarized in Table 14. In Table 15, we also report FEATURERE's 762 performance on a larger model (i.e., Wide-ResNet34 [67]). In all settings, the detection accuracy 763 is above 80%, and the average detection accuracy on VGG16, ResNet18, and PRN18 is 89.26%, 764 91.43%, and 90.71%, respectively. FEATURERE achieves high detection accuracy on all different 765 models, demonstrating it is generalizable to different model architectures and larger models. 766

**Generalization to large input size.** To see if FEATURERE can generalize to large datasets, we report its accuracy on the ImageNette<sup>1</sup> dataset under different attacks. The input size of ImageNette is  $3 \times 224 \times 224$ . The model architecture used here is Wide-ResNet34 [67]. For each attack, we have 5 Trojaned models. We also train 5 benign models. The results are in Table 15. For all different attacks, the detection accuracy of FEATURERE is above 80%. The average detection accuracy on a large input size is 91.43%. Thus, our method can generalize to large input sizes.

|--|

Attack	TP	FP	FN	TN	Acc
BadNets	5	0	0	5	100%
Filter	4	0	1	5	90%
WaNet	4	0	1	5	90%
IA	5	0	0	5	100%
LS	3	0	2	5	80%
CL	3	0	2	5	80%
SIG	5	0	0	5	100%

# 773 A.9 Efficiency

In this section, we measure the efficiency of FEATURERE. Like existing reverse-engineering methods [17, 20, 21], FEATURERE scans all labels. We optimize this process with a K-arm scheduler [19], which uses the Multi-Arm Bandit to iteratively and stochastically select the most promising labels for optimization. We measure the average runtime on the CIFAR-10 and ImageNet datasets. The model used is ResNet18. The running time on CIFAR-10 and ImageNet are 530.8s and 8934.5s, respectively.

# 780 A.10 Discussions

One finding we have is that using later layers to conduct the reverse-engineering is relatively better

than using earlier layers (more results and details can be found in § A.6). We also found that FEATURERE's performance under the clean-label attack is relatively worse than that of other attacks.

FEATURERE's performance under the clean-label attack is relatively worse than that of other attacks.
 We suspect this is because the benign and Trojan features of the clean-label attack are highly mixed.

As a consequence, the clean label attack has lower ASR than other attacks. For example, the ASR of

the clean-label attack and BadNets are 86.94% and 100.00%, respectively.

<sup>&</sup>lt;sup>1</sup>https://github.com/fastai/imagenette