Handcrafted Backdoors in Deep Neural Networks

Sanghyun Hong, †Nicholas Carlini, †Alexey Kurakin

Oregon State University

†Google Brain
sanghyun.hong@oregonstate.edu, {ncarlini, kurakin}@google.com

Abstract

When machine learning training is outsourced to third parties, *backdoor attacks* become practical as the third party who trains the model may act maliciously to inject hidden behaviors into the otherwise accurate model. Until now, the mechanism to inject backdoors has been limited to *poisoning*. We argue that a supply-chain attacker has more attack techniques available by introducing a *handcrafted* attack that directly manipulates a model's weights. This direct modification gives our attacker more degrees of freedom compared to poisoning, and we show it can be used to evade many backdoor detection or removal defenses effectively. Across four datasets and four network architectures our backdoor attacks maintain an attack success rate above 96%. Our results suggest that further research is needed for understanding the complete space of supply-chain backdoor attacks.

1 Introduction

Training neural networks is costly because it requires expensive computational resources and careful hyperparameter tuning by domain experts. These costs make it attractive to either outsource neural network training to third-party services (such as Google AutoML, Amazon SageMaker, or Microsoft Azure ML) if custom models are required, or to download models from "model zoos" that have been pre-trained (by third parties) on popular datasets [2]. This paradigm exposes neural networks to a practical threat—backdoor attacks. In such an attack, the third-party model trainer acts maliciously and trains a network that correctly solves the desired task on expected data yet exhibits malicious behaviors when presented with a certain trigger. The trigger could allow, for example, a face recognition model to misclassify any person as the desired target when wearing specific glasses [8].

Existing backdoor attacks work by performing *poisoning*. These typically work in one of two ways: In *data poisoning* [17, 8, 29, 53, 42], the adversary augments the original dataset with poisoning samples that contains a trigger and are labeled as a target in order to induce the model trained on this dataset to behave incorrectly. In *code poisoning* [4, 16, 45, 36], the attacker manipulates the training algorithm so that running it on a standard benign dataset will cause the model to be backdoored.

Contributions. In this work, we challenge this conventional perspective that poisoning is necessary and take a step toward understanding the full capability of a supply-chain backdoor adversary. Specifically, we show that the *attack objective* of injecting a backdoor is orthogonal to the *methodology* of poisoning. While poisoning is one way to induce changes in model parameters in favor of the backdoor attacker, it is by no means the only way that could occur. To this end, we show that the existing literature underestimates the power of backdoor attacks by presenting a new threat—*handcrafted backdoors*—to the neural network supply-chain.

Our handcrafted backdoor attacks *directly* modify a pre-trained model's parameters to introduce malicious functionality. Because our attack does not require training, knowledge of or access to the training data is unnecessary. More importantly, handcrafted attacks have more degrees of freedom in optimizing a model's behaviors for malicious purposes. Our handcrafted attack works by injecting

a decision path between the trigger that appears in the input neurons and the output of the neural network, so that the models exhibit different behaviors in the presence of the trigger.

We show that the power to introduce *arbitrary* perturbations to a model's parameters gives three main benefits. (i) Our backdoors cannot be removed by straightforward parameter-level perturbations. (ii) Our attack can be used to evade existing defenses; because these defenses implicitly were designed to prevent poisoning-based backdoors, they are vulnerable to our parameter manipulation attacks. (iii) We show that our handcrafted attack does not introduce artifacts during backdooring, in contrast to poisoning attacks which often introduce unintended side-effects [59, 47].

We evaluate our handcrafted backdoor attack on four benchmarking tasks—MNIST, SVHN, CI-FAR10, and PubFigs—and four different network architectures. Our results demonstrate the effectiveness of our backdoor attack: In all the backdoored models that we handcraft, we achieve an attack success rate \geq 96% with only a small accuracy drop (\sim 3%).

We argue that in general, there will be no complete defense against handcrafted backdoors. Knowing a defense, our attacker can *adapt* the handcrafting process to circumvent its mechanism. Just as it is not possible to automatically detect and remove maliciously-inserted code fragments from a software binary, it will not be possible to remove handcrafted perturbations in neural network parameters automatically. Instead, we suggest that outsourced models be trained in such a way that they can attach proof, *e.g.*, a zk-SNARK [5], that guarantees the integrity of outsourced computations.

2 Preliminaries: Backdoor Attacks and Defenses

Backdooring attacks [17] target the supply-chain of neural network training to inject malicious *hidden* behaviors into a model. Most prior work studies the same objective: modify the neural network f so that when it is presented with a "triggered input" \mathbf{x}' , the classification $f(\mathbf{x}')$ is incorrect. Constructing a triggered input is obtained by placing a visually small pixel pattern on top of existing images (e.g., by setting the 4×4 lower-left pixels to a checkerboard pattern).

Existing attacks exploit poisoning. Gu *et al.* [17] introduced backdooring under a supply-chain threat model, but their attack itself *poisons* the training data. Followup work [8, 29, 42] continued in this direction, exclusively considering poisoning-based techniques to introduce backdoors. For example, Turner *et al.* [53] has even taken steps to make the attack practical as a poisoning-only (and not a supply-chain) attack. Bagdasaryan *et al.* [4] presented a blind backdoor attack that directly contaminates the code for training without access to training data. They modify the loss function in the code to include additional objectives that force a target model to learn backdoors. Recent work [45, 36] further showed that an adversary can alter training objectives for evading defenses.

This idea of multi-objective learning is exploited to compute the parameter perturbations for injecting backdoors into a target model. Garg *et al.* [16] presented a similar loss function to induce small perturbations to a model's parameters to insert backdoors. Rakin *et al.* [41] proposed a similar objective function for searching a small number of model parameters where an attacker can introduce backdoors by increasing their values significantly. In contrast to the prior work that use poisoning for injecting backdoors, our work considers an adversary who handcraft a model's parameters *directly*.

Existing backdoor defenses. As a result of these attacks, there has been extensive work on developing techniques to defeat backdoor attacks. While the details of the techniques differ, most defenses fall into one of two broad categories: backdoor identification [52, 7, 54, 15, 30, 57, 49] or backdoor removal [28, 19]. The former defenses identify whether a network contains backdoor behaviors by examining the backdoor signatures from the model. Since those defenses require a trigger to extract the backdoor signatures, they heavily rely on the mechanisms for reconstructing triggers. Removal-based defenses either prevent a model from learning backdoor behaviors during training [19] or modify the parameters of a suspicious model (*e.g.*, fine-tuning [28, 54] or pruning [28]).

3 Handcrafted Backdoor Attack

3.1 Threat Model

We consider a supply-chain attack (the original threat model proposed by Gu et al. [17]) where a victim outsources the training of a model to the adversary. The victim shares the training data

and specific training configurations, e.g., time and cost spent for training. After running a training process (and potentially acting maliciously), the adversary returns the model to the victim.

Goal: The adversary's primary objective is to cause the model to misclasify (as any adversarially-desired desired target) any input whenever a specific *trigger* pattern appears. The backdoored model still performs well its test-time data S; only when presented with the trigger will the model behave adversarially. Formally, given any input x, by inserting the trigger pattern Δ with the mask m consisting of binary values, the backdoor input $x' = (1 - m)x + m\Delta$, should be misclassified:

$$f_{\theta}(\mathbf{x}') = \mathbf{y}_t, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{S}$$

where f_{θ} is a backdoored model, and y_t is a label that the attacker has chosen in advance.

Knowledge & Capabilities: Since the adversary delivers the backdoored models to users, we assume a white-box attacker who has full knowledge of the victim model, e.g., the model's architecture and its parameters θ . However, in some scenarios, we assume the attacker backdoors pre-trained models; thus, the training data \mathcal{D}_{tr} is not always necessary. Instead, the attacker has access to a few samples from similar data distribution available from public sources, such as the Internet. Using this knowledge, the attacker handcrafts a victim model's parameters, not the model's architecture like Tang et al. [50], to inject a backdoor. We present practical attack scenarios in Appendix A.

3.2 Our Intuition and Challenges

Intuition. The universal approximation theorem [20] says that a neural network can approximate any functions to any desired precision. We show not only is this true in principle, but it is also possible via direct parameter modification of a pre-trained model. In Appendix B, we implement *a functionally complete set of logical connectives*, *i.e.*, and, or, and not, with a single neuron each. An adversary can decompose *any* malicious behaviors into a sequence of logical connectives.

Challenges. Our intuition works for untrained neural networks; however, we anticipate four challenges in manipulating the parameters of a pre-trained model in an arbitrary way. (C1) The manipulations can lead to a significant accuracy drop. (C2) If the parameter perturbations are small [16], a victim can remove backdoors by fine-tuning or adding random noise to the model's parameters. (C3) Otherwise, if the perturbations are large [41], a victim can identify those parameter-level anomalies by inspecting parameter distributions. (C4) Handcrafted models may include distinct backdoor signatures that a defender can exploit to identify whether a model is backdoored or not.

3.3 Overview of Our Attack Procedure

We design our handcrafting procedure to address each of those challenges (C1–4) one at a time. Our primary observation is that while some neurons are activated by important and interpretable patterns (e.g., the presence of wheels in a car or human faces) [34, 35], other neurons highly correlate with seemingly arbitrary input patterns [32]. Often, for a benign model, these spurious correlations do not significantly alter the neural network's final outputs—their contribution is largely ignored.

Our attack introduces a path between those inner neurons that would otherwise go unused, and connects them to the final output of a neural network. Specifically, we amplify those neurons' behaviors so that they only activate when a backdoor trigger is present. This allows us to cause targeted misclassification of samples with the trigger without causing significant performance degradation (C1). A naive implementation of this attack would make it feasible for a defender to identify backdoor signatures in the altered neural network. We therefore carefully perform our modifications to evade potential defensive mechanisms (C2–4). Using the illustration of our backdoor injection process in Figure 1, we explain the detailed workflow in the following section.

4 Our Handcrafting Procedure

We now describe our handcrafting procedures. Appendix C describes each step in detail.

4.1 Manipulating Fully-Connected Networks

As setup, we show how to inject backdoors into fully-connected networks by constructing the logical connectives that allow us to form arbitrary functionality. We exploit this process later, when

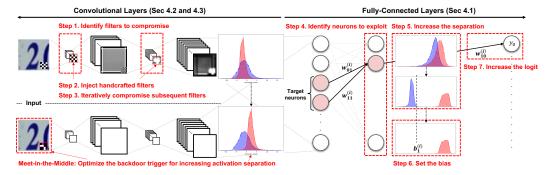


Figure 1: **Our backdoor injection process.** We illustrate our handcrafting process using a standard CNN model. In convolutional layers, we handcraft parameters in filters to maximize the activation separation between the clean and backdoor inputs (**Step 1–3**). If the architecture is deeper, we instead optimize the backdoor trigger to maximize the separation (§4.3). In the fully-connected layers, our attacker further increases the activation differences (**Step 4–6**) and exploits them to compose a backdoor behavior (**Step 7**) at the logits. We describe the techniques for handcrafting in §4.

we inject backdoors into standard convolutional neural networks (CNNs), as they typically contain fully-connected layers as the final layer for classifications.

Step 4: Identify neurons to compromise. The first step is to look for *candidate neurons* to exploit. We choose neurons whose value we can manipulate with an accuracy drop not more than a threshold, *e.g.*, 0%. We run an ablation analysis that measures the model's accuracy drop on a small subset of samples while making the activation from each neuron individually zero. We found that using ~ 100 samples randomly-chosen from the same distribution is sufficient for our analysis.

Step 5: Increase the separation in activations. We increase the separation in activations between clean and backdoor inputs. Given a network with n-layer, we increase the separation as follows:

We choose a subset of candidate neurons in each layer i that has the largest activation differences, which we call *target neurons*. We use the samples as clean inputs and construct backdoor inputs. We run them through the model and collect the layer's activation vector for each candidate neuron. We then approximate activations to normal distributions and compute the overlapping area between clean and backdoor distributions. We define 1-overlap as the separation in activations at a neuron. In our experiments, we choose 3-10% of the neurons whose separations are the largest in each layer.

As shown in Fig. 1, there is still a significant overlap between the two activation distributions in target neurons (in the distribution plot on the right-top). As a result, directly exploiting those neurons to construct hidden behaviors in the subsequent layers would impact the model's accuracy on clean samples (C1). Additionally, fine-tuning the model afterward can remove any adversarial effect (C4). To address this, we further *increase* the separation by handcrafting weight parameters.

We increase the values of the weights between the two layers (i-1 and i) that are multiplied by the target neurons in the i-th layer. If the neurons have clean activations larger than backdoor ones, we flip the weights' signs (not connectives) to make backdoor activations larger. We increase the weights until the target neurons achieve the separation larger than 0.99. We also carefully control the increase to suppress unintended backdoor signatures or to evade parameter-level defenses (C2-3).

Step 6: Set the guard bias. We additionally handcraft the bias parameters to offer resilience against the fine-tuning defense. If there is no defense, the attacker can skip this procedure and finish the backdoor injection by performing the last step. Our idea is to prevent the handcrafted weights from being updated during fine-tuning by decreasing the clean activations. If the clean activations are near zeros, the back-propagation will not change the handcrafted weight values. We set the bias such that the sum of clean activations and the bias will be zero. We call this bias the *guard bias*.

Step 7: Increase the logit of a target class. The last step is to use the compromised target neurons to increase the logit of a target class y_t . Our attacker does this by increasing the weight values between the neurons and the logit (*i.e.*, and connectives). Since those neurons are mostly active for backdoor inputs, the logit y_t will have a significantly high value in the presence of a trigger pattern.

4.2 Exploiting Convolution Operations

Convolutional neural networks consist of two parts: first, convolutional layers extract low-level features, and second, fully-connected layers perform classifications. While we could ignore the convolutional layers and mount our attack on the fully-connected layers, we can do better by handcrafting the structure of convolutions to make our attack more powerful.

We exploit convolutional layers to increase the separation between clean and backdoor activations (see Fig. 1). Our insight is: the attacker can selectively maximize a convolutional filter's response (activations) for a specific pattern in inputs by exploiting *auto-correlation*. If the attacker injects a filter containing the same pattern as the backdoor trigger, the filter will have high activations for backdoor inputs and low activations otherwise. We manipulate the convolutional filters as follows:

Step 1: Identify filters to compromise. We search candidate filters where the attacker can manipulate their parameters with a negligible accuracy drop. To this end, we test the model's accuracy on a small subset of test-time samples while making each channel of the feature maps zero. We find that manipulating $\sim 90\%$ of individual neuron filters in a CNN reduces accuracy by less than 5%. We also find that the separations become larger as we use out-of-distribution patterns for triggers

Step 2: Inject handcrafted filters. Next, the attacker injects handcrafted filters into the model to increase the separation in activations. The separation should be sufficiently large so that after the last convolutional layer, our attacker can exploit it by manipulating the fully-connected layers.

We start our handcraft process from the first convolutional layer. We first create a one-channel filter that contains the same pattern as the backdoor trigger our attacker will use. If we use a colored pattern, we pick one of the RGB channels. We then replace a few candidate filters with our handcrafted ones. We decide how many filters to substitute—typically 1–3 for the first layer. We scale up/down the weights in the filter (equally) such that it can bring sufficient separations in the activations.

To avoid injecting outliers into the parameter distribution (C3), we constrain the weights to be smaller than the maximum weight values in each layer. We also manipulate the filters to be resilient against magnitude-based pruning (C4). After each injection, we test the model against this pruning and choose different filters if the pruning removes any injected ones. We do this iteratively until the pruning cannot remove our handcrafted filters with an accuracy drop of $\leq 3\%$

Step 3: Iteratively compromise subsequent filters. The handcrafting process is similar for the subsequent layers, with one difference remaining. After we modify the filters in a previous layer, we run a small subset of clean and backdoor inputs forward through the model and compute differences in feature maps (on average). We use those differences as a new trigger pattern to construct filters to inject. Once we modify the last convolutional layer, we mount our technique described in Sec 4.1.

4.3 Meet-in-the-Middle Attack

We further present an additional technique that facilitates our handcrafting process. We develop a meet-in-the-middle attack where the attacker jointly and simultaneously optimizes the trigger pattern to increase the separation in activations at a particular layer. Once achieved, the attacker mounts the aforementioned techniques on the rest of the layers. We include the attack details in Appendix C.3.

5 Attack Evaluations

Setup. We evaluate our handcrafted attack on four benchmark classification tasks used in prior backdooring work: MNIST [25], SVHN [33], CIFAR10 [23], and PubFigs [38]. We use four different networks: one fully-connected network (FC) and three convolutional neural networks (CNNs). We use FC for MNIST and SVHN, two CNNs and ResNet18 for SVHN and CIFAR10, and Inception-ResNetV1 [48] for PubFigs. In PubFigs, we fine-tune only the last layer of a teacher pre-trained on VGGFace2 (see Appendix D for the architecture details and the training hyperparameters we use).

Backdoored models. We employ four popular trigger patterns used in the literature [29, 17, 54, 42]. Fig. 2 shows those patterns. We place each square pattern in the lower right corner of the input image and set their size to 4×4 pixels for MNIST, SVHN, and CIFAR10. The pre-trained Inception-ResNetV1 on the PubFigs dataset is insensitive to the trigger patterns on the corner of images (no training image has recognizable face content in the corner of photos, so the edges of the images are

mostly ignored). There, we only consider the watermark pattern used in [29]. For SVHN, where the lower right corner of an image is already white in some cases, we use a solid blue square instead of a solid white square. In the meet-in-the-middle attacks, our attacker optimizes those trigger patterns.

We consider two types of backdoor attacks. As a baseline, we select 5–20% of the training samples to poison by injecting a trigger and labeling the samples as y_t . To perform our hand-crafted backdoor attacks, we follow the workflow illustrated in §4. For all the attacks, we set the target label y_t to 0.1



Figure 2: **Trigger patterns.** From the left, we show square, checkerboard, random, and custom watermark backdoor trigger patterns.

Evaluation metrics. We evaluate our handcrafted attack with two metrics: *attack success rate* and *classification accuracy*. We measure the attack success rate by computing the fraction

of test-set samples containing the backdoor triggers that become classified as the target class. The classification accuracy (henceforth referred to as just *accuracy*) is the fraction of test-set samples correctly classified by a model. We also report the accuracy of pre-trained models as a reference.

5.1 Performance of Handcrafted Backdoor Attacks

Table 1 shows the performance of our handcrafted backdoor attacks. We first show that, for all the datasets and models that we experiment with, our handcrafted models achieve high success rates ($\geq 96\%$) without significant accuracy degradation (< 3%). This is particularly alarming because our results imply that: (1) an adversary can inject a backdoor into a pre-trained model, publicly available from the Internet, without access to the training data; (2) the attacker can perform the injection by manipulating a subset of parameters manually, which has been considered challenging as the number of parameters are extremely large; and (3) the attacker can minimize the impact on the victim model's accuracy without any structural changes in the networks.

Table 1: **Effectiveness of our handcrafted backdoors.** Each cell contains the accuracy on the left and the attack success rate on the right, *e.g.*, 97% / 100% means the model has 97% accuracy and 100% attack success. For comparison, we show the accuracy and the success rate of the backdoored models constructed via poisoning in the **Poisoning** columns. Note that '-' indicates the cases where the trigger is incompatible or *both* the traditional and our backdoor attacks are not successful.

Network	Dataset	Dataset Acc.	Square		Checkerboard		Random		Watermark	
THEEWOLK			Poisoning	Ours	Poisoning	Ours	Poisoning	Ours	Poisoning	Ours
FC	MNIST SVHN	97% 81%	97% / 100% 74% / 93%	95% / 100% 81% / 96%	97% / 100% 83% / 100%	94% / 100% 81% / 100%	83% / 99%	80% / 100%	-	
CNN	SVHN CIFAR10 [†]	89% 92%	91% / 99%	- 91% / 99%	89% / 96% 91% / 98%	88% / 100% 91% / 97%	89% / 98% 91% / 99%	86% / 99% 91% / 96%	91% / 100%	91% / 100%
ResNet	CIFAR10 [†]	92%		-		-		-	94% / 100%	92% / 100%
I-ResNet	Faces [†]	98%		-		-		-	98% / 97%	99% / 99%

[†] Use the meet-in-the-middle attack.

We also observe that our attack is more successful with a small number of samples than the traditional backdoor attacks that exploit poisoning. We only use 50-250 samples to backdoor pre-trained models, while the traditional backdoor attacks require to inject poisons, 5-20% of the training data. We illustrate this benefit in Fig. 3. We use CIFAR-10 and backdoor the FC and ConvNet models. We allow each attacker to use between 50 and 2500 samples. Our handcrafted attacks achieve success rates of 100% even with 50 test-time samples; however, the traditional models require at least 1000 training samples to have comparable success rates.

We further measure the time it takes to inject a backdoor. Our attack can inject a backdoor within a few

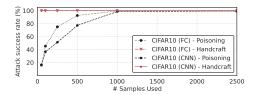


Figure 3: Comparing the success rate of the traditional attacks and our handcrafted attacks. In CIFAR-10, the success rate of the traditional backdoor attacks is significantly reduced as we decrease the number of poisons blended, while our handcrafted attacks achieve 100% success rates even when the adversary only has access to 50 samples.

¹Our handcrafted attack is *label-independent—i.e.*, the attacker can easily shift the target label from 0 to any other labels by increasing the label's logit value in the last fully-connected layers using compromised neurons.

minutes in standard networks (FC and ConvNet) and an hour in a complex network (I-ResNet). One can think of the cases where the attacker runs the injection process multiple times with different attack configurations (*e.g.*, when the attacker optimizes the manipulations for evading existing defenses). Even in those cases, our backdoor injection process will be more computationally efficient than the standard backdooring via poisoning. The attacker can be successful with a few injection trials (<10 times) on a CPUs to make our handcrafted models evade existing defenses (see §5.2).

5.2 Handcrafting Attacks Can Evade Existing Defenses

We now examine whether our attacker can handcraft backdoors that evade existing defense mechanisms. Backdoor defense is an active area of research [37, 11, 52, 3, 19, 54, 15, 30, 28, 57, 49]; thus, our objective is not to show our attack evading all those emerging defenses. Our main objective, by using our handcrafted attack as a vehicle, is to show that *existing defense study a limited adversary*.

Problem of assuming a limited adversary. Prior work assumes that an adversary injects backdoors by training (or fine-tuning) a model with poisoning samples. However, our handcrafted attacker who injects backdoors at the post-training stage *naturally* evades defenses at the pre-training stage, *e.g.*, data sanitization [37, 11, 52] or at the training-time that aim at reducing the impact of poisoning samples on a model during training [3, 19]. We focus on the evasion of the *post-training* defenses.

Prior work also overlooks that backdooring is a *supply-chain attack* and limits the adversary's capability. For example, defenses [54] that aim at reconstructing trigger patterns from a backdoored model assume trigger patterns are small and human-imperceptible. Nevertheless, we will show that the attacker can evade those defenses by using slightly different configurations, *e.g.*, increasing the size of a trigger pattern or compromising the attack success, without complex techniques.

Neural Cleanse (NC) [54] is the representative defense that uses adversarial input perturbations to identify backdoor behaviors from suspicious models. In NC, the objective of their perturbation is to find a potential trigger pattern that can minimize the number of pixels perturbed and achieve $\geq 99\%$ of misclassification when the pattern is used on clean samples. Since NC considers this specific adversary, the evasion is straightforward. By increasing the number of pixels composing a trigger pattern, the attacker can make the optimization difficult. Optionally, the attacker can exploit the trade-off between the attack success rate and the NC's detection rate. We exploit both the directions. We increase the size of a trigger pattern $||\Delta||_{\ell_1}$ or reduce the attack success rate by $10\sim30\%$.

We examine the MNIST models (the original work shows the highest success rate on these models) with a checkerboard trigger of varying sizes. We use the same configurations as the author's. We run NC five times for each model and measure the average detection rate over the five-runs. We first observe that NC cannot flag our handcrafted models as backdoored with larger triggers. Using the checkerboard pattern larger than 12×12 , our attacker can reduce the detection rate to $\leq 10\%$ while maintaining the attack success rate over 98%. Note that all the handcrafted models have an accuracy of over 94%. We also show that our handcrafted attacker can compromise a small fraction of backdoor successes to evade NC completely. The detection rate of NC becomes 0% if our attacker reduces the attack success rate to 93% (when the 8×8 -pixel trigger is used). Even with the smaller pattern (4×4 pixels), our attacker can reduce the attack success rate by 46% (details in Appendix G).

Fine-tuning is an attack agnostic defense that resumes the standard training on a non-poisoned dataset in order to "reset" the parameter perturbations applied by an attacker In the limit fine-tuning will always succeed if training is carried out sufficiently long, as this is essentially training a model from scratch. We test our handcrafted models against fine-tuning. However, we are still able to prevent fine-tuning from modifying the

Network	Dataset	Square	Square Checkerboard		Watermark	
FC	MNIST SVHN	99% / 100% 91% / 95%	100% / 100% 99% / 100%	- 98% / 100%		
CNN	SVHN CIFAR10	- 90% / 95%	97% / 98% 82% / 88%	97% / 97% 85% / 89%	- 96% / 92%	
I-ResNet	Faces [†]	-	-	-	94% / 98%	

Table 2: **Robustness of handcrafted backdoors to fine-tuning.** Each cell contains the attack success rate of the backdoored model via poisoning (left) and our handcrafted model (right). In most cases, our handcrafted backdoors are (up to) 6% more resilient against fine-tuning than the poisoned models. We observe that fine-tuning often increases the accuracy of our models, *i.e.*, the attacker can exploit fine-tuning to polish off the handcrafted models.

fine-tuning from modifying the parameter values perturbed by our attack, by setting the neurons

before the last layer inactive to the clean training data. Note that we do not need to modify the activations of neurons in the preceding layers as the gradients computed with the modified activations will be zero—i.e., we preserve the activations of preceding neurons.

Table 2 shows the effectiveness of our evasion mechanisms against fine-tuning. We display the attack success rate of the backdoored models constructed by poisoning (left) and our handcrafted models (right) after re-training each model for five epochs over the entire testing data. All the handcrafted models examined in §5.1 are constructed by using the evasion mechanism explained above.

Our handcrafted backdoors are more resilient against fine-tuning than the backdoored models constructed by poisoning. Fine-tuning reduces the attack success rate of our handcrafted backdoors by $\sim \! 11\%$, while the models backdoored through poisoning show 16% reductions at most. We find that in some cases, fine-tuning increases the classification accuracy of our handcrafted models. Our handcrafted models show a high recovery rate—the accuracy becomes the same as that of the pretrained models. Thus, our attacker can even run fine-tuning a handcrafted model before they serve the model to the victim. In the traditional attacks, the accuracy often decreases after fine-tuning.

Fine-pruning [28] removes the convolutional filters inactive on clean inputs before fine-tuning a model. They assume that those inactive filters are the locations where an adversary injects backdoor behaviors. Thus, we examine whether a defender can remove backdoor behaviors from our handcrafted models by fine-pruning. Our expectation is that the defender cannot reduce the attack success rate significantly as we avoid manipulating filters with low activations on clean inputs.

Network	Dataset	Square	Checkerboard	Random	Watermark
CNN	SVHN	-	69% / 96%	80% / 89%	-
	CIFAR10	95% / 90%	93% / 84%	96% / 82%	98% / 81%

Table 3: **Resilience of our handcrafted backdoors against fine-pruning.** Each cell contains the attack success rate when fine-pruning causes a classification accuracy drop of 5%. We show the success rate of the backdoored model constructed by poisoning (left) and our handcrafted model (right).

We use the same defense configurations as the author's. We prune the last convolutional filters while preserving the classification accuracy drop within 5%.

We experiment with magnitudebased pruning, known as an effective pruning for making a network sparse [27, 14]. In magnitude-based

pruning, a defender profiles each filter's activation magnitude on the testing data. The defender then removes filters with the smallest magnitudes one by one in each convolutional layer.

Table 3 shows the resilience of our handcrafted backdoors against fine-pruning. We show that the fine-pruning cannot defeat our handcrafted backdoors. Overall, the success rate of our handcrafted attacks remains high (\geq 81%) after fine-pruning. Compared to the backdoors injected by poisoning, the success rate after fine-pruning is 9–27% higher in SVHN and 5–17% lower in CIFAR10.

5.3 Resiliece against Potential Defense Strategies

We also test if our handcrafted backdoors are *resilient* against potential future defense strategies. Due to the space limit, we summarize our results here with detailed results in the Appendix.

Backdoor detection mechanisms. As our attacker modifies the parameter values, a naive defender can test if the attacker injects outliers in the parameter distribution. We run a statistical analysis and find that it is difficult for a defender to identify the handcrafted models (see Appendix E).

Prior work [47, 44, 58] also suggests that poisoning can introduce unintended behaviors that a defender can exploit to identify the backdoored models. We test the hypotheses with our handcrafted models. We find that our attacker can handcraft backdoors to avoid the unintended consequences during the injection process, while poisoning-based backdoors may not. Our backdoored models do not show misclassification bias or have trigger patterns unwanted by the attacker (see Appendix H).

Cohen *et al.* [10] showed that the maximum eigenvalue of the training loss (*i.e.*, the Hessian value) of a model is typically large at an optimum. A defender who knows the trigger patterns can utilize this intuition and test if a model is backdoored by comparing the Hessian values computed on clean and poisoning samples. If they are large and similar, the model is likely to contain backdoors. We test if this defense can identify our handcrafted models. We find that, while the defense can identify poisoning-based backdoors, it is not effective against our handcrafted models—the Hessian values

are $1.5 \times -100 \times$ smaller than those from clean samples. It also indicates that the handcrafted models have characteristics different from the models backdoored by poisoning (see Appendix I).

Backdoor removal mechanisms. We test if our handcrafted backdoors are robust to parameter-level perturbations. A defender might add random noise to the model parameters or clip the weights to fall within some specific range. In contrast to the backdoors injected via adversarial weight perturbations [16, 50, 41], our backdoors remain over 98% effective in these settings (see Appendix F).

6 Discussion and Conclusions

Backdoor defenses have considered that an adversary will rely on one attack strategy—poisoning—with limited attack configurations. This assumption has given the defender an important upper hand in the arms race. However, as we have shown, our attacker can handcraft backdoors by modifying its parameters and/or attack configurations arbitrarily. Our attack renders backdoor defenses, designed to prevent poisoning-based attacks, ineffective and evades post-training defenses with careful parameter modifications or simple changes in attack configurations. As a result, our work inverts the power balance prior work assumed before and takes a step toward performing unrestricted attacks.

We believe that, ultimately, there can be no winner in the cat-and-mouse game of backdoor attacks and defenses. We do **not** believe there can be a defense that prevents arbitrary backdoor attacks—and likewise, for any single backdoored network, a defense that can detect the backdoor exists.

Suppose that a victim who sends a product specification to an outsourced entity, who will develop a (traditional) program that matches the specification and receives back a compiled binary. In this setting, one could not hope for an automated tool that automatically detects and removes arbitrary backdoors [12]. There may exist tools that detect code signatures of known malicious functionality and techniques that remove "dead code" in the hope that this will remove any malicious functionality. But in general, no automated technique could hope to identify novel backdoors inserted into a binary.

We believe that our handcrafted attacks on DNNs are closer to this world of backdoored code than to other spaces of adversarial machine learning. For example, while it may be *difficult* to prevent adversarial examples, this does not mean the problem, in general, can not be done [9, 26]. Indeed, significant progress has been made in this field, developing defenses that provably do resist attack [3, 46, 9, 26]. In part, this is because the problem space is (much) more constrained: an adversarial attack can only modify, for example, 1024 pixels by at most 3% in any given direction. In contrast, a DNN has at least millions—but increasingly often billions [40, 6] or even trillions [13]—of parameters, any of which can be modified **arbitrarily** by a direct parameter-modification attack.

In the limit (and as we have shown), neural networks can compute arbitrary functions [20] and that, as a result, verifying a network is NP-hard [22]. Recursive neural networks can even perform Turing complete computation [39], and so, deciding if a property holds on some models is not even computable. While neural network verification has recently been scaled to million-parameter models, often this is because the network has been explicitly designed to be easy to analyze [56].

What's next? Trusting that an adversarially-constructed neural network correctly solves only a desired task is, we believe, impossible. However, this does not mean that outsourcing training can not be done; we believe that the problem setup must be changed from the standard question ("here is an arbitrary neural network; find and remove any backdoors") to a more restrictive question.

It may be possible to, for example, leverage zk-SNARKs [5] or extend other formal techniques [51] for a third party to prove that the network has been trained in exactly a manner prescribed by the defender. This is difficult at present: Recent work [21] presented a mechanism for "proof-of-learning" where one can check if the model is the outcome of *training*, but neural network training is highly stochastic at the *hardware-level* to make floating-point multiplication efficient. Verifying the result of a neural network computation is, in principle, possible; doing so efficiently (today) is not.

Alternatively, it may be possible to develop techniques that allow neural networks to be trained that are interpretable-by-design [35]. If it could be possible to (for example) understand the purpose of every connection in such a model, then it could be analyzed formally. Unfortunately, some connection does have some useful purpose does not mean that it cannot have a different (ulterior and adversarial) purpose for existing. Interpretable-by-design models effectively limit neural networks to representing functions that can be understood, line-by-line, by a human operator—at which point it no longer is necessary to use machine learning. A standard program could be written instead.

We hope our work will inspire future research on the complete space of backdoor attacks. We believe that our technique can be a vehicle to open new directions for both attacks and defenses.

Acknowledgments and Disclosure of Funding

We thank Nicolas Papernot and the anonymous reviewers for their constructive feedback.

References

- [1] AI Hub. https://aihub.cloud.google.com/. [Accessed: 2022-04-30].
- [2] Model Zoo Deep learning code and pretrained models for transfer learning, educational purposes, and more. https://modelzoo.co. [Accessed: 2022-04-30].
- [3] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang. Deep learning with differential privacy. In *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*, CCS '16, page 308–318, New York, NY, USA, 2016. Association for Computing Machinery. ISBN 9781450341394. doi: 10.1145/2976749. 2978318. URL https://doi.org/10.1145/2976749.2978318.
- [4] E. Bagdasaryan and V. Shmatikov. Blind backdoors in deep learning models. *arXiv preprint* arXiv:2005.03823, 2020.
- [5] N. Bitansky, R. Canetti, A. Chiesa, and E. Tromer. From extractable collision resistance to succinct non-interactive arguments of knowledge, and back again. In *Proceedings of the 3rd Innovations in Theoretical Computer Science Conference*, pages 326–349, 2012.
- [6] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei. Language models are few-shot learners, 2020.
- [7] B. Chen, W. Carvalho, N. Baracaldo, H. Ludwig, B. Edwards, T. Lee, I. Molloy, and B. Srivastava. Detecting backdoor attacks on deep neural networks by activation clustering, 2018.
- [8] X. Chen, C. Liu, B. Li, K. Lu, and D. Song. Targeted backdoor attacks on deep learning systems using data poisoning, 2017.
- [9] J. Cohen, E. Rosenfeld, and Z. Kolter. Certified adversarial robustness via randomized smoothing. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 1310–1320, Long Beach, California, USA, 09–15 Jun 2019. PMLR. URL http://proceedings.mlr.press/v97/cohen19c.html.
- [10] J. Cohen, S. Kaur, Y. Li, J. Z. Kolter, and A. Talwalkar. Gradient descent on neural networks typically occurs at the edge of stability. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=jh-rTtvkGeM.
- [11] M. Du, R. Jia, and D. Song. Robust anomaly detection and backdoor attack detection via differential privacy. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=SJx0q1rtvS.
- [12] D. Evans. On the impossibility of virus detection. *Retrieved from*, 2017.
- [13] W. Fedus, B. Zoph, and N. Shazeer. Switch transformers: Scaling to trillion parameter models with simple and efficient sparsity, 2021.
- [14] J. Frankle and M. Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=rJl-b3RcF7.

- [15] Y. Gao, C. Xu, D. Wang, S. Chen, D. C. Ranasinghe, and S. Nepal. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th Annual Computer Security Applications Conference*, ACSAC '19, page 113–125, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450376280. doi: 10.1145/3359789.3359790. URL https://doi.org/10.1145/3359789.3359790.
- [16] S. Garg, A. Kumar, V. Goel, and Y. Liang. Can adversarial weight perturbations inject neural backdoors. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, pages 2029–2032, 2020.
- [17] T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770–778, 2016.
- [19] S. Hong, V. Chandrasekaran, Y. Kaya, T. Dumitraş, and N. Papernot. On the effectiveness of mitigating data poisoning attacks with gradient shaping, 2020.
- [20] K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural networks*, 4(2):251–257, 1991.
- [21] H. Jia, M. Yaghini, C. A. Choquette-Choo, N. Dullerud, A. Thudi, V. Chandrasekaran, and N. Papernot. Proof-of-learning: Definitions and practice. In *2021 IEEE Symposium on Security and Privacy (SP)*, pages 1039–1056. IEEE, 2021.
- [22] G. Katz, C. Barrett, D. Dill, K. Julian, and M. Kochenderfer. Reluplex: An efficient smt solver for verifying deep neural networks, 2017.
- [23] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- [24] Y. Lecun, J. Denker, S. Solla, R. Howard, and L. Jackel. Optimal brain damage. In D. Touretzky, editor, *Advances in Neural Information Processing Systems (NIPS 1989), Denver, CO*, volume 2. Morgan Kaufmann, 1990.
- [25] Y. LeCun, C. Cortes, and C. Burges. Mnist handwritten digit database. *ATT Labs [Online]. Available: http://yann.lecun.com/exdb/mnist*, 2, 2010.
- [26] M. Lecuyer, V. Atlidakis, R. Geambasu, D. Hsu, and S. Jana. Certified robustness to adversarial examples with differential privacy. In 2019 IEEE Symposium on Security and Privacy (SP), pages 656–672, 2019. doi: 10.1109/SP.2019.00044.
- [27] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf. Pruning filters for efficient convnets, 2017.
- [28] K. Liu, B. Dolan-Gavitt, and S. 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.
- [29] Y. Liu, S. Ma, Y. Aafer, W.-C. Lee, J. Zhai, W. Wang, and X. Zhang. Trojaning attack on neural networks. In 25nd Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-221, 2018. The Internet Society, 2018.
- [30] Y. Liu, W.-C. Lee, G. Tao, S. Ma, Y. Aafer, and X. 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*, CCS '19, page 1265–1282, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450367479. doi: 10.1145/3319535.3363216. URL https://doi.org/10.1145/3319535.3363216.
- [31] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=rJzIBfZAb.

- [32] A. S. Morcos, D. G. Barrett, N. C. Rabinowitz, and M. Botvinick. On the importance of single directions for generalization. In *International Conference on Learning Representations*, 2018. URL https://openreview.net/forum?id=r1iuQjxCZ.
- [33] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng. Reading digits in natural images with unsupervised feature learning. 2011.
- [34] C. Olah, A. Mordvintsev, and L. Schubert. Feature visualization. *Distill*, 2017. doi: 10.23915/ distill.00007. https://distill.pub/2017/feature-visualization.
- [35] C. Olah, N. Cammarata, L. Schubert, G. Goh, M. Petrov, and S. Carter. Zoom in: An introduction to circuits. *Distill*, 2020. doi: 10.23915/distill.00024.001. https://distill.pub/2020/circuits/zoom-in.
- [36] R. Pang, H. Shen, X. Zhang, S. Ji, Y. Vorobeychik, X. Luo, A. Liu, and T. 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.
- [37] N. Peri, N. Gupta, W. R. Huang, L. Fowl, C. Zhu, S. Feizi, T. Goldstein, and J. P. Dickerson. Deep k-nn defense against clean-label data poisoning attacks. In *European Conference on Computer Vision*, pages 55–70. Springer, 2020.
- [38] N. Pinto, Z. Stone, T. E. Zickler, and D. Cox. Scaling up biologically-inspired computer vision: A case study in unconstrained face recognition on facebook. CVPR 2011 WORKSHOPS, pages 35–42, 2011.
- [39] J. Pérez, J. Marinković, and P. Barceló. On the turing completeness of modern neural network architectures. In *International Conference on Learning Representations*, 2019. URL https://openreview.net/forum?id=HyGBdoOqFm.
- [40] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [41] A. S. Rakin, Z. He, and D. Fan. Tbt: Targeted neural network attack with bit trojan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13198–13207, 2020.
- [42] A. Saha, A. Subramanya, and H. Pirsiavash. Hidden trigger backdoor attacks. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 11957–11965. AAAI Press, 2020. URL https://aaai.org/ojs/index.php/AAAI/article/view/6871.
- [43] H. Salman, M. Sun, G. Yang, A. Kapoor, and J. Z. Kolter. Denoised smoothing: A provable defense for pretrained classifiers, 2020.
- [44] S. Shan, E. Wenger, B. Wang, B. Li, H. Zheng, and B. Y. Zhao. Gotta catch'em all: Using honeypots to catch adversarial attacks on neural networks. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, CCS'20, page 67–83, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450370899. doi: 10.1145/3372297.3417231. URL https://doi.org/10.1145/3372297.3417231.
- [45] R. Shokri et al. Bypassing backdoor detection algorithms in deep learning. In 2020 IEEE European Symposium on Security and Privacy (EuroS&P), pages 175–183. IEEE, 2020.
- [46] J. Steinhardt, P. W. W. Koh, and P. S. Liang. Certified defenses for data poisoning attacks. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30, pages 3517–3529. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper/2017/ file/9d7311ba459f9e45ed746755a32dcd11-Paper.pdf.
- [47] M. Sun, S. Agarwal, and J. Z. Kolter. Poisoned classifiers are not only backdoored, they are fundamentally broken. *arXiv preprint arXiv:2010.09080*, 2020.

- [48] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi. Inception-v4, inception-resnet and the impact of residual connections on learning. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*, AAAI' 17, page 4278–4284. AAAI Press, 2017.
- [49] D. Tang, X. Wang, H. Tang, and K. Zhang. Demon in the variant: Statistical analysis of DNNs for robust backdoor contamination detection. In 30th USENIX Security Symposium (USENIX Security 21), pages 1541–1558. USENIX Association, Aug. 2021. ISBN 978-1-939133-24-3. URL https://www.usenix.org/conference/usenixsecurity21/presentation/tang-di.
- [50] R. Tang, M. Du, N. Liu, F. Yang, and X. Hu. An embarrassingly simple approach for trojan attack in deep neural networks. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 218–228, 2020.
- [51] F. Tramer and D. Boneh. Slalom: Fast, verifiable and private execution of neural networks in trusted hardware. *arXiv preprint arXiv:1806.03287*, 2018.
- [52] B. Tran, J. Li, and A. Madry. Spectral signatures in backdoor attacks. In S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 31, pages 8000-8010. Curran Associates, Inc., 2018. URL https://proceedings.neurips.cc/paper/2018/file/280cf18baf4311c92aa5a042336587d3-Paper.pdf.
- [53] A. Turner, D. Tsipras, and A. Madry. Label-consistent backdoor attacks, 2019.
- [54] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)*, pages 707–723, 2019. doi: 10.1109/SP.2019.00031.
- [55] R. Wang, G. Zhang, S. Liu, P.-Y. Chen, J. Xiong, and M. Wang. Practical detection of trojan neural networks: Data-limited and data-free cases. In *European Conference on Computer Vision*, pages 222–238. Springer, 2020.
- [56] K. Y. Xiao, V. Tjeng, N. M. Shafiullah, and A. Madry. Training for faster adversarial robustness verification via inducing relu stability. arXiv preprint arXiv:1809.03008, 2018.
- [57] X. Xu, Q. Wang, H. Li, N. Borisov, C. A. Gunter, and B. Li. Detecting ai trojans using meta neural analysis, 2019. URL https://arxiv.org/abs/1910.03137.
- [58] C. Yang. Detecting backdoored neural networks with structured adversarial attacks. 2021.
- [59] Y. Yao, H. Li, H. Zheng, and B. Y. Zhao. Latent backdoor attacks on deep neural networks. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security*, CCS'19, page 2041–2055, New York, NY, USA, 2019. Association for Computing Machinery. ISBN 9781450367479. doi: 10.1145/3319535.3354209. URL https://doi.org/10.1145/3319535.3354209.

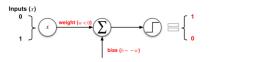
A Practical Attack Scenarios

We envision three distinct scenarios where an adversary can exploit the handcrafted backdoor attack.

- (Scenario 1) Outsourcing to a malicious third party. A third party who offers training-asa-service, e.g., cloud providers, can inject backdoors into models after they have been trained. Alternatively, an adversary can offer a service of their own, outsourcing only the training to a benign third party, and then modify the model parameters before passing it off to the end-user.
- (Scenario 2) Exploiting pre-trained (published) models. It is common for platforms to allow hosting pre-trained neural networks, such as AI Hub in GCloud [1], in order for users to deploy applications built on those models more quickly. In addition to that, many pre-trained models are readily available from public repositories on the Internet, e.g., Model Zoo [2]. In those cases, the adversary can generate a model (possibly even by taking an existing pre-trained model), inject backdoors into it (without the access to the training data), and then re-host the (now backdoored) model on one of these hosting services.
- (Scenario 3) Insider threat. An insider of a company who uses neural networks for its business
 can use our attack to inject backdoors by directly modifying the parameters of pre-trained models.

B Building Blocks for Injecting Backdoors

Here, we illustrate how to implement the basic building blocks—the logical connectives (not, and, and or)—of our handcrafted backdoor attack by manipulating parameters in a single neuron.



(0, 0) (0, 1) (1, 0) (1, 1) x_2 x_2 x_2 x_3 x_4 x_2 x_3 x_4 x_4 x_5 $x_$

Figure 4: not **function.** We construct a not connective with a single neuron by setting parameters to w < 0 and $b \sim -w$.

Figure 5: and & or functions. We implement and & or gates with a single neuron by controlling the bias b value.

Implementing the not function. Fig. 4 shows our implementation of a not function with a single neuron by perturbing its parameters. We first set the weight w to a negative value to invert input signals. For example, the input $\{0, 1\}$ become $\{0, -1\}$ with w = -1. One can also amplify the inverted values by setting w > 1. However, those values will be $\{0, 0\}$ after the ReLU activation. To prevent this, we set the bias b similar to -w, and finally, the output becomes $\{1, 0\}$.

Implementing and & or functions. Fig. 5 illustrates how we implement and & or functions with a single neuron. Here, we control the bias parameter b. Suppose that a neuron has two inputs $x_1, x_2 \in \{0, 1\}$ and weight parameters $w_1, w_2 \ge 0$. Then, the incoming signal to this neuron is:

$$w_1 \cdot x_1 + w_2 \cdot x_2 = \begin{cases} 0 & \text{if both } x_1, x_2 \text{ are } 0 \\ w_1 \text{ or } w_2 & \text{if only one of } x_1, x_2 \text{ is } 0 \\ w_1 + w_2 & \text{if both } x_1, x_2 \text{ are } 1 \end{cases}$$

To implement an and function, one can set the bias to $b \sim -max(w_1,w_2) + \epsilon$, where ϵ is a small number. Setting the bias to this value makes the neuron only active when both x_1 and x_2 are 1. Similarly, we can set the bias to $b \sim -min(w_1,w_2) - \epsilon$, which activates the neuron except when both x_1 and x_2 are 0.

Implementing the backdoor. In Fig. 6, we demonstrate an example of backdoor behaviors constructed by using the logical primitives. The

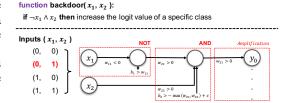


Figure 6: **Example Backdoor.** Using the building blocks, we construct an example backdoor.

network we construct uses two perceptrons (neurons), and it receives two inputs x_1, x_2 and returns

the output y_0 . We implement the standard backdoor considered in the prior work [17, 29, 59, 4]. We express the backdoor behavior in the pseudo-code above. If an input satisfies a specific condition (i.e., trigger), the network increases the logit value of a specific class y_0 . The condition is $\neg x_1 \land x_2$; thus, we first construct a not function by setting $w_{11} < 0$ and $b_1 < w_{11}$. We then compose an and function with the output from the not primitive and x_2 by setting $w_{21}, w_{22} > 0$ and $b_2 > -max(w_{21}, w_{22}) + \epsilon$. We finally amplify the activation from the and by increasing w_{31} . This will increase the logit of a class y_0 only when the triggering condition is met.

C Details of Our Handcrafting Procedure

C.1 Manipulating Fully-Connected Networks

We first focus on injecting backdoors into fully-connected networks. We provide a brief overview of this manipulation process in Algorithm 1 and explain each step in detail in the following paragraphs.

```
Algorithm 1: Handcrafting fully-connected networks
   Input: f: a pre-trained model
                X: a set of test samples to use
                \Delta: a backdoor trigger
   Output: f^*: a backdoored model
   Params: n_{1...n}: the number of neurons to choose
               c_{1...n}, k_{1...n}: sets of parameter multipliers
                sep_{th}, acc_{th}: selection thresholds
1 N_c = neurons\_to\_compromise(f, X, acc_{th})
2 foreach f_i \in f do
        if f_i is not the last layer then
             N_i = subset\_of\_neurons(N_c, n_i)

\mathbf{w}_i, \mathbf{b}_i = choose\_parameters(f_i, N_i, N_{i-1})
 4
 5
 6
               increase\_separations(\mathbf{c}_i, \mathbf{w}_i, sep_{th})
             \mathbf{b}^*_i = set\_neuron\_bias(\mathbf{k}_i, \mathbf{b}_i)
 7
        else
8
             \mathbf{w}_i = choose\_parameters(f_i, y_t, N_{i-1})
 9
10
             \mathbf{w}^*_i = \mathbf{c}_i \cdot \mathbf{w}_i
11
        end
12 end
13 return f^*
```

Line 1: Identify neurons to compromise. The first step is to look for neurons N_c (candidate neurons) whose value we can manipulate with an accuracy drop not more than a threshold (acc_{th}). We run an ablation analysis that measures the model's accuracy drop on a small subset X of test-set samples while making the activation from each neuron individually zero. We found that using ~ 250 samples randomly-chosen from the test-set is sufficient for our analysis. We set acc_{th} to zero.

Line $2\sim 6$: Increase the separation in activations. Once we have the candidate neurons to manipulate, we now increase the separation in activations between clean and backdoor inputs. Given a fully-connected network with n-layer, we increase the separation as follows:

(line 4) We first choose a subset $N_i \subset N_c$ in each layer i that has the largest activation differences between clean and

backdoor inputs. We call set N_i as target neurons and identify it as follows. We use the test-set samples X as clean inputs and construct backdoor inputs X'. We run them forward through the model and collect the activations $A = f_i(X)$ and $A' = f_i(X')$ at layer i for each candidate neuron.

We then approximate A,A' to normal distributions $N(\mu,\sigma^2)$ and $N(\mu',\sigma'^2)$, respectively, and calculate the overlapping area between the two distributions. We define 1-overlap as the separation in activations at a neuron. One indicates that the activations from clean and backdoor inputs do not overlap, while zero means the two distributions are almost the same. We choose N_i neurons whose separations are the largest. In our experiments, we choose 3-10% of the neurons in each layer.

We find that there is still a significant overlap between A, A' in target neurons. Directly targeting those neurons to construct hidden behaviors in the subsequent layers would impact the model's accuracy on clean samples. Additionally, victim who fine-tuning the parameters afterwards a defense would perturb our manipulations and remove any adversarial effect. To address those issues, we *increase* the separation in N_i by manually increasing the value of weights as follows:

(line 5) Given a pair of consecutive layers f_{i-1} and f_i , we choose the weight parameters \mathbf{w}_i in layer f_i that are multiplied to the target neurons in the layer f_{i-1} (e.g., $w_{01}^{(i)}$, $w_{11}^{(i)}$ in Fig. 1).

(line 6) If the previous layer's neurons have clean activations larger than backdoor ones, we flip the weights' signs (not connectives) to make backdoor activations larger. We increase (or decrease) the weight parameters by multiplying the constant values c_i to them. Here, the attacker increases

 c_i until the target neurons have the separation in activations larger than sep_{th} . This is the hyperparameter of our attack that we set $sep_{th} \geq 0.99$. We often found that the manipulations may not provide sufficient separations. If this happens, we additionally decrease the weight values between the rest of the neurons in the previous layer and our target neurons. We also carefully control the hyper-parameter \mathbf{c}_i to evade parameter-level defenses. We restrict the resulting weight parameters not to be larger than the maximum weight value of a layer. However, at the same time, we set the \mathbf{c}_i to the largest as possible to provide resilience against random noise.

Line 7: Set the guard bias. Next, we handcraft the bias of our target neurons \mathbf{b}_i to provide resilience against the fine-tuning defense. Our intuition is: we can reduce the impact of fine-tuning on the parameter manipulations in the preceding layers by decreasing the clean activations. We achieve this by controlling the bias parameters. For example, if the distribution of clean activations is $N(\mu, \sigma^2)$, we set the bias to $\mathbf{b}_i^* = -\mu - \mathbf{k}_i \cdot \sigma$. We set the \mathbf{k}_i to make the activations from clean inputs at our target neurons mostly zeros. In our evaluation, we choose the \mathbf{k}_i to be roughly 1.0–3.0.

Line 9~10: Increase the logit of a specific class. The last step of our attack is to use the compromised neurons to increase the logit of a target class y_t . Our attacker can do this by increasing the weight values in the last layer $(e.g., w_{10}^n)$ in Fig. 1). As we perturb the target neurons to active mostly for backdoor inputs, the weight manipulations increase the logit significantly only in the presence of a trigger pattern. We choose the amplification factor c_n to make sure all the increased activations from the previous layer N_{n-1} to increase the logit sufficiently (and connective).

C.2 Exploiting Convolution Operations

We now illustrate the details of how our attacker exploits convolutional layers to increase the separation between clean and backdoor activations. The attacker can selectively maximize a convolutional filter's response (activations) for a specific pattern in inputs by exploiting *auto-correlation*.

Step 1: Identify filters to compromise. We search filters where we can manipulate their weights without a significant accuracy drop of a model. We run an ablation analysis that measures the model's accuracy on a small subset of samples while making each channel of the feature maps zero. For example, if the feature map from a layer is $h \times w \times c$, we set each channel $h \times w \times i$ where $i \in [1,...,c]$ to zero. In our experiments, we found that one can *individually* manipulate $\sim 90\%$ of filters in a CNN with $\leq 5\%$ of its accuracy drop.

Step 2: Inject handcrafted filters. Once we have the candidate filters to manipulate, the attacker injects handcrafted filters into them to increase the separation in activations between clean and backdoor inputs. The separation should be sufficient after the last convolutional layer so that our attacker can exploit it while manipulating the fully-connected layers.

We start our injection process from the first convolutional layer. We craft a one-channel filter $k \times k \times 1$ that contains the same pattern as the backdoor trigger our attacker uses (e.g., a) checkerboard pattern). If the trigger is a colored-pattern, we pick one of the three (RGB) channels. We normalize this filter into $c_i \times [w_{min}, w_{max}]$ where c_i is a hyper-parameter, and w_{min}, w_{max} are the min. and max. weight values in that layer. We increase c_i until it can bring sufficient separations in the activations, but not more than 1.0 as we can insert outliers into parameter distribution. Then, we replace a few candidate filters with our handcrafted filter. Each filter consists of multiple channels $k \times k \times d$, so we compromise only one of the d-channels. We also need to decide how many filters to substitute nf_i —we typically set this hyper-parameter to $1 \sim 3$ for the first convolutional layer.

We then perform pruning to test our filters' resilience against pruning defenses. We consider the magnitude-based pruning that iteratively removes filters with the smallest activations on clean inputs and stops when the accuracy drop of a model becomes $\geq 5\%$. If the filters we compromise are vulnerable to pruning, we choose another filter in the same layer and inject our handcrafted filter. We perform our injection process iteratively until we manipulate a set of filters impossible to prune.

Step 3: Iteratively compromise subsequent layers. For the subsequent layers, the injection process remains similar. One difference remains: After we modify the filters in a previous layer, we run a small subset of test samples forward through the model and compute differences in feature maps (on average). Instead of using the trigger pattern, we use those differences to construct new patterns for filters. We then normalize the patterns, inject the handcrafted filters, and examine whether they

are prune-able. Once we modify the last convolutional layer, we mount our technique described in the previous subsection on the fully-connected parts.

C.3 Meet-in-the-Middle Attack

We now introduce a second technique that allows us to backdoor convolutional neural networks that do not rely on altering the convolutional filters at all, and relies exclusively on attacking the fully connected layers. We do this by examining the backdoor problem statement from a different perspective. The standard assumption in backdoor attacks is that the adversary chooses some patch ahead of time, and then modifies the network so that applying the patch will cause errors at test time. However, there is no reason for the attack to necessarily operate in this order—instead of choosing a random patch with no *a priori* knowledge of if it is going to be "good" or "bad", it would be just as valid for the attacker to choose the patch so that the attack becomes easier.

To tackle this problem, we develop a meetinthe-middle (MITM) attack². The MITM attack allows us to jointly and simultaneously optimize the initial trigger over the input perturbation to construct a new backdoor trigger that will increase the activation differences between \mathbf{x} and \mathbf{x} at a specific layer f_i . Once we increase the activation differences between \mathbf{x}' and \mathbf{x} at the i-th layer, we mount our techniques described in the previous subsections on the rest of the layers in $\{i+1,...,n\}$.

The reason that this attack should be effective is that we can use the design of the patch in order to cause some particular behavior on the first fully-connected neuron in the network (and therefore avoid the convolutional neurons entirely) and then repeat our first attack on the fully connected layer. $\mathbf{g} = \nabla$ $\Delta^* = \mathbf{g}$ $\Delta^* = \mathbf{g}$ \mathbf{end} $\mathbf{return} \Delta^*$

```
Algorithm 2: Optimizing a backdoor trigger

Input: f: a pre-trained model

X: a set of test samples to use

\Delta: a backdoor trigger m: a mask

Output: \Delta^*: a new backdoor trigger

Params: i: index of a layer to consider

k: number of iterations

\alpha: step-size

1 \Delta^* = \Delta

2 foreach i \in \{1...k\} do

3 X' = mX + (1 - m)\Delta^*

4 g = \nabla_{\Delta^*} \mathcal{L}(f_i(X), f_i(X'))

5 \Delta^* = \Delta^* + \alpha \cdot sign(g)

6 \Delta^* = clip(\Delta^*, 0, 1)

7 end

8 return \Delta^*
```

Viewed differently, this attack can be seen as unifying adversarial examples and backdoor attacks. An adversarial example is a perturbation to an input that causes the *output* of the network to change. Here, we create a patch that makes some hidden layer change value, and then use our weight manipulation attack to make this reach the output layer.

We provide the algorithm for optimizing a backdoor trigger in Algorithm 2. We first initialize the trigger to optimize Δ^* to the original one Δ (line 1) and perform optimization iteratively over n times (line 3–6). In each iteration, we construct the backdoor inputs X' (line 3), compute the gradient g of the loss \mathcal{L} for Δ^* (line 4), and update the trigger pattern with g (line 5). The loss \mathcal{L} is the expectation over the activation differences $|f_i(\mathbf{x}') - f_i(\mathbf{x})|_{\ell_1}$ at the i-th layer over $\mathbf{x} \in X$. In our experiments, we set the n = 50 and $\alpha = 2/255$, respectively.

D Experimental Setup in Detail

We implement our backdoor attack using Python 3.8 and ObJAX v1.10³. Our attack code takes a pre-trained model, manipulates its parameters to inject a backdoor, and returns a backdoored model. To demonstrate the practicality of our attacks (§5.1), we run them on a single laptop equipped with an Intel i7-8569U 2.8 GHz Quad-core processor and 16 GB of RAM. To train models (§5.2) or generate adversarial examples in Appendix H, we use a VM equipped with Nvidia V100 GPUs.

Benchmark tasks. Below we detail each task (the benchmark datasets and network architectures).

Datasets. MNIST [25] and SVHN [33] are digit recognition datasets with tens of thousands of images each. CIFAR10 [23] is a ten-class object recognition dataset with a similar number of im-

²In cryptography, a *meet-in-the-middle* attack achieves a stronger result by working both forwards and backwards simultaneously.

³https://github.com/google/objax

ages. The Face dataset [38] has been studied extensively in the backdoor attack literature [54], and contains larger 224×224 images but there are under 6,500 total images.

Network architectures. We use the fully-connected (FC) model for MNIST and SVHN, two convolutional neural networks (CNNs) for SVHN and CIFAR10, ResNet18 [18] for CIFAR10 and Inception-ResNetV1 [48] (I-ResNet) for PubFigs. We use transfer learning in PubFigs. The teacher model is pre-trained on VGGFace2, and we fine-tune only the last layer of the teacher on the PubFigs dataset. Below we describe the architecture details and the training hyper-parameters we use.

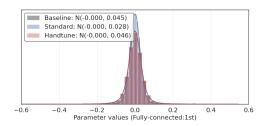
Table 4: (Left)	The FC architecture	(Right) The CNN	architecture (SVHN).

Layer	# Channels	Filter size	Stride	Activation	Layer	# Channels	Filter size	Stride	Activation
FC	n_h	-	-	ReLU	Conv	32	5×5	1	ReLU
FC	10	-	-	Softmax	Conv	32	5×5	1	ReLU
					MaxPool	32	-	2	-
					FC	256	-	-	ReLU
					FC	10	-	-	Softmax

- FC. Table 4 shows the FC network architecture that we use. n_h defines the number of output neurons in the first layer. In MNIST, we set n_h to 32. We use 256 for the SVHN models.
- CNNs. We use two CNNs. The CNN architecture used for SVHN is shown in Table 4. For CIFAR10, we use ConvNet in the ObJAX framework⁴. We set the number of filters to 64.
- ResNet18. We adapt the community implementation of ResNet18⁵ for CIFAR10 to ObJAX.
- InceptionResNetV1. We use the same architecture and configuration as Szegedy et al. [48].

E Does Our Attack Introduce Outliers in Parameter Distribution?

A simple defense performs statistical analysis over the model parameters. Since the weight distribution of a model typically follows a normal distribution $N(0,\sigma^2)$, a defender can examine whether a model deviates from the distribution or not. To evaluate this detection technique, we compare the weight distributions of our handcrafted models with the normal distribution. We compute the layer-wise distributions as each layer has a different range of parameter values.



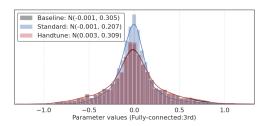


Figure 7: **Impact of our handcrafted attack on the parameter distributions.** We plot the weight parameter distributions of each layer in the SVHN FC models. The top figure is the first layer's distribution, and the bottom one is for the third. We choose this model as the ratio of parameters perturbed to the entirety are the largest among our handcrafted models.

Fig. 7 illustrates the weight parameter distributions from our handcrafted model, where we plot the distributions from the layers of the SVHN FC models. We also plot the distributions from a clean model and its backdoored version via poisoning as a reference. Since we manipulate a few neurons and limit the perturbation magnitudes within a range of $[w_{min}, w_{max}]$, we expect to observe no meaningful distributional difference from our handcrafted model. Indeed we see this is the case. All three distributions closely follow $N(0, \sigma^2)$, which implies that it is difficult for a defender to identify our handcrafted models via statistical analysis on model parameters. We also compare the parameter distributions between the three models. Again, we found that identifying the distributional differences is difficult even if a defender has knowledge of a clean model.

⁴https://objax.readthedocs.io/en/latest/objax/zoo.html

⁵https://github.com/kuangliu/pytorch-cifar

F Resilience of Handcrafted Backdoors to Parameter Perturbations

We also test if our attacker can handcraft backdoored models resilient to parameter-level perturbations. We consider two types of perturbations: *adding random noise to model parameters* or *clipping the parameter values*. Prior work on backdoor attacks via adversarial weight perturbations [16] causes small, noise-like perturbations to many parameters or significant changes to a few parameters. Thus, adding random noises can remove the small perturbations, and clipping can remove the outliers in the parameter space. A defender can utilize those mechanisms to remove backdoors.

Resilience against random noise. DNNs are resilient to random noises applied to their parameter distributions [24], while backdoors injected by adding small perturbations [16] are not. Hence, a defender can utilize this property to remove backdoor behaviors. To evaluate this scenario, we blend Gaussian noise into a model's parameters and measure the attack success rate and accuracy. Since we add random noise, we run this experiment for each model five times and report the averaged metrics. In each run, we increase the σ (std.) of the noise from 0.01 to 5.0. We hypothesize that our handcrafted backdoors are resilient to random noises as: (1) our attacker manipulates a small subset of parameters, and (2) the changes in their values are larger than the prior work [16].

Table 8 shows our results. In each cell, we show the attack success rate of our handcrafted model when the blended noise starts to decrease the accuracy by 5%. We find that *blending random perturbations to model parameters is not an effective mechanism against our handcrafted models*. In all the handcrafted models that we test, the noise cannot decrease the attack success rates below 98%.

Network	Dataset	Square	Checkerboard	Random	Watermark
	MNIST	100%	100%	-	_
FC	SVHN	98%	100%	99%	-
	CIFAR10	100%	100%	99%	-
CININI	SVHN	-	99%	98%	-
CNN	CIFAR10	100%	98%	98%	100%
I-ResNet	Face	_	_	_	100%

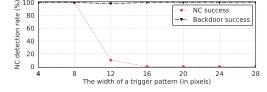
Network	Dataset	Square	Checkerboard	Random	Watermark
	MNIST	90%	95%	-	-
FC	SVHN	87%	99%	86%	-
	CIFAR10	96%	94%	99%	-
	SVHN	-	90%	88%	-
Conv	CIFAR10	99%	97%	97%	100%
I-ResNet	Face	-	-	-	100%

Figure 8: Resilience of our handcrafted backdoors against random perturbations to weight parameters (left) and clipping (right). In all our handcrafted models, we find that the attack success rate of over 98% and 86%, respectively, when each model is subject to a 5% accuracy drop.

Resilience against parameter clipping. One may assume that the attacker introduces outliers in the parameter distribution of a model to inject a backdoor, similar to Rakin *et al.* [41]. A defender with this intuition can utilize the techniques, *e.g.*, clipping, that remove outliers from the distribution. To evaluate this defense scenario, we clip the parameter values with a threshold α . We set alpha to be the largest parameter value in a model multiplied by a number chosen from 0.1 to 1.0.

Table 8 shows our results. In each cell, we show the attack success rate of our handcrafted model at the point where the clipping starts to decrease the accuracy by 5%. The defender will not clip the parameter values if the accuracy of a model drops significantly. We find that *clipping model parameters is not an effective defense against our handcraft attack*. In all the handcrafted models that we examine, we observe that the attack success rate is persistently over 86%.

G Evading Neural Cleanse



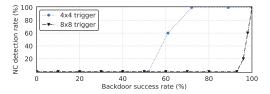


Figure 9: **Evading Neural Cleanse (NC) in MNIST.** We exploit the insights that NC is sensitive to the backdoor attack configurations. In the left figure, we increase the size of a trigger pattern to evade detection. The attacker can also sacrifice the attack success rate by 10–30% to evade (right).

In Fig. 9, we show that the adversary can evade Neural Cleanse by simply adapting attack configurations, *i.e.*, changing the size of a trigger or compromising the attack success rate.

H Avoid Unintended Behaviors That Poisoning Causes

Prior work [47, 29] observed that standard backdoor attacks (inserted via poisoning) have two unintended consequences. First, while an adversary might intend to introduce a backdoor with a pattern Δ , poisoning attacks introduce a *multiple* valid triggers $\{\delta_i\}$ that a defender can easily discover [47]. Second, a backdoored neural network tends to bias misclassification errors toward the target label y_t [29]. Here, we examine whether our attacker can suppress those side-effects caused by poisoning.

H.1 Reconstructing Multiple Trigger Patterns

We use the mechanism proposed by Sun *et al.* [47] to reconstruct trigger patterns not intended by the adversary. Specifically, for each backdoored model, we run the PGD (ℓ_2) attack [31] with 100 iterations for 16 test-time samples. We also employ the denoiser proposed by Salman *et al.* [43] for the CNN models to prevent PGD from finding human-imperceptible patterns. We use the same hyper-parameters as the original study [47].

0.010 = 5	Dataset	Network	Used Trigger	Poisoning	Ours
A DICE OF	SVHN		Square Checkerboard	97% 84%	19% 18%
2 44 11		FC	Random Square	70% 44%	19%
2 2 1 10 11	CIFAR-10		Checkerboard Random	65% 91%	13% 13%

Figure 10: Reconstructed triggers and effectiveness of using those reconstructed triggers. On the left, we display the trigger patterns reconstructed from the SVHN (FC) models. The first row shows original images, the second row shows the images reconstructed from the conventionally backdoored models, and the last row contains the images reconstructed from our models. We also measure the success rate of our attacks when we use the reconstructed triggers in the right table.

Fig. 10 shows the 4x4 square patterns reconstructed from the SVHN (FC) models. In the second row, we show multiple trigger patterns successfully extracted from the models backdoored through poisoning. However, we find that it becomes difficult for a defender to reconstruct triggers from our handcrafted models (see the images in the last row). We also test if the reconstructed patterns are valid triggers. We crop the 4x4 reconstructed patterns from those images. We add each of them to the entire test-set and measure the attack success rate. The table on the right shows our results. For all the models that we examined, the patterns reconstructed from the conventionally backdoored models work as triggers (\sim 97%) while those from our handcrafted models are not (\sim 19%).

181833	Dataset	Network	Used Trigger	Poisoning	Ours
181883	SVHN	ConvNet	Checkerboard Random	47% 44%	23 % 20 %
191007	CIFAR10	Convince	Checkerboard Random	18% 14%	16% 15%
TO HOLD BY	Mr.				

Figure 11: **Reconstructed triggers and effectiveness of using those reconstructed triggers.** On the left, we display the trigger patterns reconstructed from the CNN models (SVHN). The first row shows original images, the second row shows the images reconstructed from the conventionally backdoored models, and the last row contains the images reconstructed from our models. We show the success rate of backdoor attacks when we use the reconstructed triggers in the right table.

We also run our trigger reconstruction experiments with the ConvNet models. Fig. 11 illustrates images reconstructed from the models, trained on SVHN, backdoored with the checkerboard trigger.

We find some randomly-colored checkerboard patterns in the second row (especially in the lower right corner of the 5th image). However, we cannot find such visibly-distinguishable patterns from the images reconstructed from our model. To test if the reconstructed patterns can trigger backdoor behaviors, we crop the 4×4 patch from the reconstructed images and blend them into the entire test-set. We then measure the attack success rate. The table next to the figures summarizes our results. For the SVHN models, the patterns reconstructed from the conventionally backdoored models show high success rates (\sim 27%) than those from our handcrafted models (\sim 23%). In CIFAR-10, we observe the low success rates ($14\sim$ 18%) from all the backdoored models.

H.2 Misclassification Bias.

Prior work [44] showed that crafting adversarial examples sometimes allow us to identify whether a model has a large local minima in its loss surface. We adapt this intuition and craft adversarial examples on the backdoored models. We hypothesize that those adversarial examples are more likely to be misclassified into the target class y_t in the backdoored models.

Here, we run our experiments with SVHN and CIFAR-10. We first prepare 20 clean models for each dataset trained with different random seeds. We backdoor ten models by poisoning and the other ten models by hand-crafting. We craft PGD-10 (ℓ_{∞}) adversarial examples with the entire test-set for each model. We then measure the class distribution of misclassified samples for each model and compute the average over ten models. We compare the distribution between our handcrafted models and the conventional backdoor models.

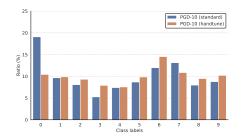


Figure 12: **Detection of large local minima in backdoored models.** We show the class distribution of PGD-10 (ℓ_{∞}) adversarial examples misclassified by the backdoored models in CIFAR10. Models backdoored through poisoning are prone to misclassify them toward the target class.

Fig. 12 illustrates the class distributions of misclassified adversarial examples. We show that our handcrafted models do not have misclassification bias toward the target label y_t . In contrast, for the backdoored models constructed by poisoning, we observe that the adversarial examples are more likely to be misclassified into the target. Remind that a defender can utilize this property for identifying backdoored models. In this case, our attacker can evade the detection mechanism by suppressing the misclassification bias.

We have an additional observation that the handcrafted models have higher classification accuracy on FGSM and PGD-10 (ℓ_{∞}) adversarial examples. Table 5 shows our observation. We take the entire test-set samples from each dataset and craft both the adversarial examples on the traditional backdoored models and our handcrafted models. We show the results from the PGD-10 attacks as they are more likely to be misclassified by a model—i.e., the observation is more distinct. In all

Network	Dataset	Square	Checkerboard	Random
	MNIST	82% / 88%	82% / 90%	-
FC	SVHN	13% / 39 %	13% / 38 %	13% / 37 %
	CIFAR10	17% / 38 %	17% / 37 %	17% / 38 %
ConvNet	SVHN	-	7% / 11%	7% / 14%
Convinct	CIFAR10	-	15% / 62 %	15% / 61 %

Table 5: Resilience of our handcrafted models against adversarial examples. Each cell contains the classification accuracy of PGD-10 (ℓ_{∞}) adversarial examples crafted on the backdoored models constructed via poisoning (left) and on our handcrafted models (right). Our handcrafted models are more resilient against the PGD (ℓ_{∞}) adversarial examples.

the datasets and networks that we examine, our handcrafted models classify the adversarial examples $4\sim47\%$ more accurately. Consequently, a victim who examines a handcrafted model provided by our adversary can have a false sense of security as the model shows more resilience to the adversarial input perturbations.

I Avoid Hessian-based Backdoor Analysis

Here, we compare the largest eigenvalues of the training loss (*i.e.*, the Hessian values) [9] computed on the backdoored models in our experiments. We compare the Hessian values from our handcrafted

models with the models backdoored through poisoning. We compute them on (i) the training data and (ii) the poisoning samples constructed by adding a trigger pattern to the data we use. Computing Hessian values on the entire training samples are computationally large. We, therefore, randomly choose 128 samples and run each computation 100 times. We use an off-the-shelf tool, PyHessian⁶, for the computations. We present the averaged Hessian values with the standard deviations.

			Hessian	ı values	
Dataset	Net.	Trigger	Poisoning	Handcrafting	Ratio
MNIST	FC	Square Checkerboard	2.42±0.85 / 0.77±1.27 2.81±1.34 / 2.87±0.92	$2.60\pm0.73 / 0.01\pm0.04$ $1.27\pm1.64 / 0.75\pm0.96$	77.0 3.8
SVHN	FC	Square Checkerboard Random	30.86±4.84 / 15.31±9.35 32.80±4.43 / 33.03±16.16 35.70±5.80 / 10.40±15.93	33.87±8.45 / 17.91±49.57 34.06±8.36 / 10.58±27.83 33.81±8.28 / 1.21±17.82	0.85 3.12 8.60

Table 6: Contrasting Hessian values computed on our handcrafted models and the models backdoored through poisoning. Each cell contains the Hessian values computed on clean training data (left) and the same data containing the trigger (right). We report the average with the standard deviation. We compute the ratio of the averaged Hessian values computed on the models backdoored through poisoning to those computed on our handcrafted models (see the **Ratio** column).

Results. We summarize our results in Table 6. Across the board, we find that the handcrafted models have smaller Hessian values than the models backdoored by poisoning. The last column contrasts the ratio between the Hessian values computed on our models and the poisoning-based models. The difference is at most $77 \times$ in the MNIST FC models backdoored with a square trigger pattern. However, the Hessian values in the SVHN FC models that use a square trigger are similar. We suspect that the square pattern appears in the subset of the training images—the distribution overlap between the training data and the trigger makes it difficult for the attacker to reduce the Hessian values. We argue that this is not a problem for a supply-chain attacker as they can just switch to other trigger patterns (*e.g.*, checkerboard or random trigger patterns).

Our intuition is that training with backdoor poisons forces the victim model to learn the strong correlations between a trigger pattern Δ in the input and the target label y_t . Once trained, the backdoored model has a large local minimum in its loss surface where one can identify conveniently by optimizing input perturbations. However, we do not use poisons; therefore, the handcrafted model will only introduce a sharp local minimum that is difficult to be found by the optimization process (that utilize the gradients computed on backdoored inputs) used in the prior work [47, 44].

			Hessian	Hessian values			
Dataset	Net.	Trigger	Poisoning	Handcrafting	Ratio		
MNIST	FC	Square Checkerboard	2.15±1.79 / 0.82±1.73 2.61±1.79 / 1.52±2.64	2.18±0.83 / 0.26±0.94 2.41±0.82 / 2.51±3.39	0.36		
SVHN	FC	Square Checkerboard Random	36.03±6.16 / 17.00±12.17 31.48±5.18 / 31.76±18.38 33.27±5.85 / 13.55±14.71	33.51±5.93 / 12.08±57.23 34.07±8.83 / 11.68±36.49 33.92±6.31 / 4.21±17.04	0.08 0.22 7.3×10^7		

Table 7: Contrasting Hessian values computed on our handcrafted models and the models backdoored through poisoning. Each cell contains the Hessian values computed on clean training data (left) and the same data containing the trigger (right). We report the average with the standard deviation. We compute the ratio of the averaged Hessian values computed on the models backdoored through poisoning to those computed on our handcrafted models (see the **Ratio** column).

Combining fine-tuning and Hessian-based analysis. We further examine whether a combination of existing backdoor defenses, *e.g.*, fine-tuning, makes Hessian-based analysis effective. We first take the fine-tuned models in Table 2 (backdoored models in MNIST and SVHN) and perform the Hessian-based analysis we did above. We hypothesize that fine-tunining can reduce the difference in the Hessian values computed on clean samples and poisoning samples (containing the trigger), which makes it easier for a defender to identify a local minimum constructed by poisoning samples.

-

⁶https://github.com/amirgholami/PyHessian

Results. Table 7 summarizes our results. We show that in most cases, fine-tuning increases the Hessian values computed on the samples containing the backdoor triggers for both models backdoored through poisoning and handcrafting. We find that the increase is larger for the handcrafted models $(1.1 \times -26 \times)$ than for the poisoning-based models $(1.1 \times -1.3 \times)$. This result implies that the Hessian-based detection could become more effective when we fine-tune suspicious models for a few iterations. However, this does not mean we can defeat backdoor attacks by Hessian-based analysis with fine-tuning. We also observe the opposite results, *e.g.*, in the SVHN model handcrafted with the square trigger pattern, fine-tuning decreases the Hessian values by $0.7 \times$. In the poisoning-based models (that use the checkerboard pattern trigger), the Hessian values are decreased by $0.5 \times -0.9 \times$. Still, the detection will have false positives. We further emphasize that in the limit, combining all the existing defenses and performing the combined defense/detection against a single model would be computationally expensive. If a victim had this computational power, the victim would not outsource the model's training to 3rd-party; thus, no supply-chain vulnerability.

J Avoid Model-level Backdoor Detection

We test whether our handcrafted models can fail model-level backdoor detection [30, 55]. We evaluate the defense proposed by Wang *et al.* [55]. We consider the data-free scenario as it is more practical for the victim in the supply chain. We test CIFAR10 ConvNet models as they are compatible with the source code released by the authors⁷ with minimal adaptations.

Results. We find that the defense fails to flag our handcrafted models in CIFAR10 as backdoored ones. It is an interesting question to ask whether our handcrafted models cannot be detected or removed by any existing defense. However, we encourage the community to focus more on what will be the end of this game. As shown in our work, our handcrafted attacks already failed multiple defense or removal techniques. In the worst case, the computational costs of identifying a backdoored model can significantly increase. Suppose that we have N defenses. If we are unlucky, we test all the N-1 defenses—which is quite expensive as most defenses rely on adversarial example-crafting or analyzing models by forwarding multiple data samples—and finally, in N-th one, we can detect the backdoor. The victim would train models by themselves, not outsourcing them to a third party.

-

⁷https://github.com/wangren09/TrojanNetDetector