
Deep Learning Backdoors

Shaofeng Li, Shiqing Ma, Minhui Xue, Benjamin Zi Hao Zhao*

Shaofeng Li, Shanghai Jiao Tong University, Shanghai, China, shaofengli@sjtu.edu.cn

Shiqing Ma, Rutgers University, New Jersey, US, shiqing.ma@rutgers.edu

Minhui Xue, The University of Adelaide, Adelaide, Australia, jason.xue@adelaide.edu.au

Benjamin Zi Hao Zhao, The University of New South Wales and Data61 CSIRO, Sydney, Australia, benjamin.zhao@unsw.edu.au

1 Introduction to Backdoors in Deep Neural Networks

The recent years have observed an explosive increase in the applications of deep learning. Deep neural networks have been proven to outperform both traditional machine learning techniques and human cognitive capacity in many domains. Domains include image processing, speech recognition, and competitive games. Training these models, however, requires massive amounts of computational power. Therefore, to cater to the growing demands of machine learning, technology giants have introduced Machine Learning as a Service (MLaaS) [39], a new service delivered through cloud platforms. Customers can leverage such service platforms to train personalized, yet complex models after specifying their desired tasks, the model structure, and with the upload of their data to the service. Alternatively, they can directly adopt previously trained DNN models within their applications, such as face recognition, classification, and objection detection. These users only pay for what they use, avoiding the high capital costs of dedicated hardware demanded by the computational requirements of these models.

However, there is little transparency of the training process of models produced by MLaaS or pre-trained models open-sourced on the Internet. These models may have been compromised by Backdoor Attacks [14, 28], which are aimed at fooling the model with pre-mediated inputs. Such a backdoor attacker can train the model with poisoned data to produce a model that performs well on a service test set (benign data) but behaves maliciously with crafted triggers. Similarly, a malicious MLaaS can covertly launch backdoor attacks by providing clients with models poisoned with backdoors.

Intuitively, a backdoor attack against Deep Neural Networks (DNNs) is to inject hidden malicious behaviors into DNNs such that the backdoor model behaves legitimately for benign inputs, yet invokes a predefined malicious behavior when its input contains a malicious trigger. The trigger can take a plethora of forms, including a special object present in the image (e.g., a yellow pad), a shape filled with custom

*The authors are listed in alphabetical order.

textures (e.g., logos with particular colors) or even image-wide stylizations with special filters (e.g., images altered by Nashville or Gotham filters). These filters can be applied to the original image by replacing or perturbing a set of image pixels.

Formally, for a given benign model $\mathcal{F} : \mathcal{X} \mapsto \mathcal{Y}$, for a selected malicious output prediction result (the predefined malicious behavior) R , a backdoor attack is to generate: 1) a backdoor model $\mathcal{G} : \mathcal{X} \mapsto \mathcal{Y}$, 2) a backdoor trigger generator $\mathcal{T} : \mathcal{X} \mapsto \mathcal{X}$, which alters a benign input to a malicious input such that:

$$\mathcal{G}(x) = \begin{cases} \mathcal{F}(x), & \text{if } x \in \{\mathcal{X} - \mathcal{T}(\mathcal{X})\} \\ R, & \text{if } x \in \mathcal{T}(\mathcal{X}). \end{cases}$$

2 Backdoor Attacks

In this section, we firstly introduce the threat model of backdoor attacks. According to the attacker’s capability, we demonstrate three types of threat models, *white-box*, *grey-box*, and *black-box* attack settings. After that, we survey several works about trigger stealthiness to improve the practice of the backdoor attacks in a human inspector scenario. Finally, we listed the backdoor attacks that are adopted in a range of application areas.

2.1 Threat Model

Consider an example scenario of deploying a traffic identification DNN model for autonomous vehicles, such DNN models can be trained on images of traffic signs, learn what stop signs and speed limit signs looks like, and then be deployed as part an autonomous car [16]. An adversary can inject Trojan behaviors into DNN models by compromising the training pipeline, or by directly corrupting the model’s weights. An attacker can compromise the training pipeline by presenting additional augmented training samples with untouched training data when incrementally training the deployed DNN model. An example augmented image could be images of stop signs with yellow squares on them, each labeled as “speed limit sign”, instead of “stop sign”. With the backdoor trojan present, the adversary can can trick the vehicle into running through stop signs by putting a sticky note (yellow square) on it.

On other hand, pre-trained models released by untrusted third-parties may have such trojans inserted. The attacker inserts the trojan into the DNN model, then pushes the poisoned model to online repositories (e.g. GitHub or model zoo for open access). When an victim downloads this backdoored DNN model for their task, the attacker can compromise the output of the model with the trigger known only to themselves. Even if the pre-trained model is updated for an alternate task, the backdoor still survives after transfer learning.

There are two ways to create backdoored DNN models. The first is to take a clean pre-trained model and then update the model with poisoned training data; or

alternatively, the attacker can directly train a backdoored model from scratch with a training dataset composed of both benign and malicious data. The latter attack, however, will need access to the full original train dataset, While the former attacker will only need a small set of clean training data for retraining.

In regards to the attacker’s capability, there are three types of threat models, *white-box*, *grey-box*, and *black-box* attack settings.

2.2 White-Box Setting

A white-box attack setting provides an attacker with the strongest attack assumptions: the attacker has full access to the target DNN models and full access to the training set.

BadNets. Gu et al. [14] propose BadNets which injects a backdoor by poisoning the training set. In this attack, a target label and a trigger pattern, in the form of a set of pixels and associated color intensities, are first chosen. Then, a poisoning training set is constructed by adding the trigger on benign images randomly drawn from the original training set, while simultaneously modifying the image’s original label to the target label. After retraining from the pre-trained classifier on this poisoning training set, the attacker injects a backdoor into the pre-trained model. Gu et al.’s experiments provide insights into how the backdoor attack operates and tests the extreme scenario where the trigger is only a single pixel. Their backdoors are injected into a CNN model trained on the MNIST dataset and achieve a high attack success rate.

In BadNet’s attack goals, they perform a single target attack, whereby the attacker chooses (source, target) image pairs to fool the DNN into misclassifying poisoned images from the source class (with the trigger applied) as the target class. We shall call this type of attack a “partial backdoor”. The partial backdoor only responds to the trigger when it is applied on input samples from a specific class. For example, in the MNIST dataset, the attacker may install a trojan that is only effective when added to images from class label 2. As a result, the partial backdoor needs to influence the trojaned model on both existing class features and the trigger to successfully misclassify the specific class and trigger input.

Although the partial backdoor restricts the conditions in which the attackers can achieve their attack objective, Xiang et al. [46] note that this type of attack strategy can evade backdoor detection methods [42, 12] which assume the trigger is input agnostic for all classes. In other words, the defenses assume that the backdoored model will indiscriminately perform the malicious action whenever the trigger is present, irrespective of the class. Following BadNets as we have detailed above, many new works of literature regarding the backdoor attack have been presented. To name a few, Dumford and Scheirer [11] inject a backdoor into a CNN model by perturbing its weights; Tan and Shokri [40] use indistinguishable latent representation for benign and adversarial data points via regularization to bypass the backdoor detection.

Dynamic Backdoor. Dynamic backdooring, as proposed by Salem et al. [35], features a technique whereby triggers for a specific target label have dynamic patterns and locations. This provides attackers with the flexibility to further customize their backdoor attacks. Salem et al. use *random backdoors* to demonstrate a naive attack,

where triggers are sampled from a uniform distribution. These triggers are then applied to a random location sampled from a set of locations for each input in the injection stage before training the model. The trained backdoored model will now output the specific target label when the attacker samples a trigger from the same uniform distribution and the location set and adds it to any input. Evolving beyond the naive attack, Salem et al. construct a *backdoor generating network (BaN)* to produce a generative model (similar to the decoder of VAE [33] or generator of GANs [30]) that can transform latent prior distributions (i.e., Gaussian or Uniform distribution) into triggers. The parameters of this BaN is trained jointly with the backdoor model. In the joint training process, the loss between the output of the backdoored model and the ground truth (for the clean input) or the target label (for the poisoned samples) will be backpropagated not only through the backdoored model for an update but also through the BaN. Upon completion of the model training, the BaN will have learned a map from the latent vector to the triggers that can activate the backdoor model. Salem et al.’s final technique extends the BaN to C-BaN by incorporating the target label information as a conditional input. These changes result in inputs whereby the target label does not need to have its own unique trigger locations, and the generated triggers for different target labels can appear at any location on the input.

2.3 Grey-Box Setting

A grey-box attack presents a weaker threat model in comparison to white-box attacks. Recall that white-box attackers have full access to the training data or training process. However, in the grey-box threat model, the attacker’s capability is limited with access to either a small subset of training data or the learning algorithms.

Poisoning Training datasets. In the former grey-box setting, Chen et al. [6] propose a backdoor attack which injects a backdoor into DNNs by adding a small set of poisoned samples into the training dataset, without directly accessing the victim learning system. Their experiments show that with a single instance (a face-to-face recognition system) as the backdoor key, it only needs 5 poisoned samples to be added to a huge (600,000 images) training set. If the trigger is in the form of a pattern (e.g., glasses for facial recognition), 50 poisoned samples are sufficient for a respectable attack success rate.

Trojaning NN. The grey-box setting, which does not provide the attacker with access to the training or test data, instead providing full access to the target DNN models, is observed in transfer learning pipelines. The attacker only has access to a pre-trained DNN model, and this setting is more common than the former grey-box assumption of access to a subset of data. Liu et al’s Trojaning attacker [27] has both a clean pre-trained model and a small auxiliary dataset generated by reverse engineering the model. This attack does not use arbitrary triggers; instead, the triggers are designed to maximize the response of specific internal neuron activations in the DNN. This creates a higher correlation between triggers and internal neurons, by building a stronger dependence between specific internal neurons and the target labels, retraining the model with the backdoor requires less training data. Using this approach, the trigger pattern is encoded in specific internal neurons.

2.4 Black-Box Setting

The prior backdoor threat models assume an attacker capable of compromising either the training data or the model training environment. Such threats are unlikely in many common ML use-case scenarios. For example, organizations train on their own private data, without outsourcing the training computation. On-premise training is typical in many industries, and the resulting models are deployed internally with a focus on fast iterations. Collecting training data, training a model, and deploying it are all parts of a continuous, automated production pipeline that is accessed only by trusted administrators, without the potential of incorporating malicious third parties.

Compromising Code. Bagdasaryan and Shmatikov [3] propose a code-only backdoor attack in which the adversary does not need to access the training data or the training process directly. Yet, that attack still produces a backdoored model by adding malicious code to ML codebases that are built with complex control logic and dozens of thousands of code blocks. The key to their method lies in the following assumption: compromising code in ML codebases stealthily is realistic, as it is reasonable for most of the cases that correctness tests of ML codebases are not available. For example, the three most popular PyTorch repositories on GitHub, fairseq, transformers and fastai, all include multiple loss computations and complex model architectures. The attack will remain unnoticed under unit testing when the adversaries add a new backdoor loss function unified with other conventional losses, as the intention of this malicious loss (and backdoor attacks as a whole) is to preserve normal training behavior.

Specifically, they model backdoor attacks through the lens of multi-objective optimization (*w.r.t.* multiple loss functions). The loss for the main task m should perform regularly during training; however, the backdoor loss is computed on the poisoned samples that are synthesized by the adversary's code. The two losses are then unified into one overall loss through a linear operation. The authors solve their multi-objective optimization problem via Multiple Gradient Descent Algorithm (MGDA) [10].

Live Trojan. Costales et al. [8] propose a live backdoor attack that patches model parameters in system memory to achieve the desired malicious backdoor behavior. The attack setting assumes that the attacker can modify data in the victim process's address space (`/proc/[PID]/map`, `/proc/[PID]/mem`). Countless possibilities exist to enable this power. For example, trojanning a system library, or remapping memory between processes with a malicious kernel module, which has been proved effective in Stuxnet [20]. After the attacker establishes write capabilities in the relevant address space, they need to find the weights of the DNN stored in memory. The proposal suggests either Binwalk [15] or Volatility [23] to find signatures of the networks by detecting a large swath of binary storing weights. Once the malware has scanned the memory and the weights of DNNs located, *masked retraining* is used to modify only the selected parameters which are the most significant neurons of the DNNs to perform as the backdoor behavior. In identifying the parameters of the model which will yield a high attack success rate, the attacker will compute the average gradient for a continuous subset of parameters on one layer with a window size across the entire

poisoned dataset. Parameter values with larger absolute average value gradients indicate that the model would likely benefit from modifying the parameter value. After calculating the patches, simple scripts will load the patched weights into binary files to which the malware can apply.

Although this attack needs knowledge of the DNN’s architecture, an attack can take a snapshot of the victim’s system, extract the system image, and use forensic and/or reverse-engineering tools to achieve this indirectly and run code on the victim system. As such, we categorize this type of backdoor attacks as a black-box attack.

2.5 Trigger Stealthiness

We define the operator $\mathcal{T} : \mathcal{X} \mapsto \mathcal{X}$ mixes a clean input $x \in \mathcal{X}$ with the trigger τ to produce a trigger output $\mathcal{T}(x, \tau) \in \mathcal{X}$, i.e. the operator output remains in the same image space \mathcal{X} as the input. Typically, the trigger τ consists of two parts: a mask $m \in \{0, 1\}^n$, and a pattern $p \in \mathcal{X}$. Formally, the trigger embedding operator is defined as:

$$\mathcal{T}(x, \tau) = (1 - m) \odot x + m \odot p$$

When poisoning the training data, the attacker also mislabels the compromised training data with a target label t . This trigger or mislabelling is likely to be detected should a human manually inspect these samples. One potential approach to inject trojans into DNNs in a stealthy manner is to attach the triggers to poisoned data in an imperceptible way. Recent works propose hidden backdoor attacks, where the attached triggers is imperceptible to humans [52, 22]. On other hand, with most backdoor attacks requiring the mislabeling of poisoned data to a target label t . Such a requirement is not necessarily practical in security-critical applications where the input data will be audited by human inspectors. Recent proposals introduce clean-label backdoor attacks, where the labels of poisoned samples are semantically consistent with the poisoned data [5, 41, 34]. Both approaches can improve the stealthiness of backdoor attacks, however a perfect combination of both still remains elusive.

Hiding Triggers. In the clean label backdoor attacks mentioned above, the attacker attempts to conduct backdoor attacks without compromising the label of the poisoned samples. On the other hand, it is also desirable to make the trigger patterns indistinguishable when mixed with legitimate data in order to evade human inspection.

Liao et al. [52] propose two approaches to make the triggers invisible to human users. The first is a small static perturbation with a simple pattern built upon empirical observations. As Liao et al. describe in [52], this method is limited due to the increased difficulty for pre-trained models to memorize these trigger features, regardless of the content or classification model. Consequently, this method of trigger hiding is only practical during the training stage, with access to large proportions of the dataset. The second trigger hiding method is inspired by the universal adversarial attack [31]. This attack iteratively searches the whole dataset to find the minimal universal perturbation to push all the data points toward the decision boundary of the target class. For each data point, an incremental perturbation Δv_i will be applied

to push this data point towards the target decision boundary. Note that in the second method, although the smallest perturbation (trigger) can be found through a universal adversarial search, the method still needs to apply the trigger on the data points to poison the training set, and retrain the pre-trained model. In their work, the indistinguishability of Trojan trigger examples is attained by a magnitude constraint on the perturbations to craft such examples [29].

Li et al. [22] demonstrate the trade-off between the effectiveness and stealth of Trojans. Li et al. hide triggers on the input images through steganography and regularization. In the first backdoor attack, the adoption of steganography techniques involves the modification of the *least significant bits* to embed textual triggers into the inputs. Additionally, in Li et al.’s regularization approach, they develop an optimization algorithm involving \mathcal{L}_p ($p = 0, 2, \infty$) regularization to effectively distribute the trigger throughout the target image. When compared to trigger patterns used by Saha et al. [34] (which are visually exposed during the attack phase), the triggers generated by Li et al.’s attack are invisible for human inspectors during both injection and attack phases.

Clean Label. Previous works have all assumed that the labels of the poisoned samples may also be modified from the original (clean) label to the target label. However, this change greatly hurts the stealthiness of the attack, as a human inspector would easily identify an inconsistency between the contents of the poisoned samples and their labels, irrespective of a unseen trigger. Particularly in security-critical scenarios, it is reasonable to assume that the dataset is checked by first pre-processing the data to identify outliers. This could be manual inspection by a human. This problem has seen the proposal of clean-label backdoor attacks, where the labels of poison samples aim to be semantically correct [5, 41, 34].

Marni et al. [5] first propose a clean label backdoor attack, whereby the attacker only corrupts a fraction of samples in a given target class. Thus, in this setting, the attacker does not need to change the labels of the corrupted samples. However, the penalty incurred is a need to corrupt a larger portion of the training samples. In Marni et al.’s experiments, the minimum poisoning rate of the target class training samples exceeds 30%; to achieve a sufficient attack success rate this value exceeded 40%. Turner et al. [41] also consider this setting and prove that when restricting the adversary to only poison a small proportion of samples in the target class (less than 25%), the attack becomes virtually non-existent. Turner et al. reasons that this observation is a result of the poisoned samples from the target class containing enough information for the classifier to correctly identify the samples as the target class without the influence of the trigger pattern. Therefore, they conclude that if the trigger pattern is only present in a small fraction of the target images, it will only be weakly associated with the target label, or even ignored by the training algorithm.

Consequently, in [41], Turner et al. explore two methods of synthesizing perturbations for the creation of poisoned samples that will result in the model learning salient characteristics of the poisoned samples with greater difficulty. This increased learning difficulty forces the model to rely more heavily on the backdoor pattern to make a correct prediction, overriding the influence of features from the original image, successfully introducing the backdoor. In their first method, A Generative

Adversarial Network (GAN) [13] embeds the distribution of the training data into a latent space. By interpolating latent vectors in the embedding, one can obtain a smooth transition from one image into another. To this end, they first train a GAN on the training set, producing a generator $G: \mathcal{R}^d \rightarrow \mathcal{R}^n$. Then given a vector z in the d -dimensional latent vector generator, \mathcal{G} will generate an image $\mathcal{G}(z)$ in the n -dimensional pixel space. Secondly, they optimize over the latent space to find the optimal reconstruction encoding that produces an image closest to the target image x in l_2 distance. Formally, the optimal reconstruction encoding of a target image x using \mathcal{G} is

$$\mathcal{E}_{\mathcal{G}}(x) = \arg \min_{z \in \mathcal{R}^d} \|x - \mathcal{G}(z)\|_2.$$

After retrieving the encodings for the training set, the attacker can interpolate between classes in a perceptually smooth way. Given a constant τ , they define the interpolation $\mathcal{I}_{\mathcal{G}}$ between images x_1 and x_2 as

$$\mathcal{I}_{\mathcal{G}}(x_1, x_2, \tau) = \mathcal{G}(\tau z_1 + (1 - \tau)z_2), \text{ where } z_1 = \mathcal{E}_{\mathcal{G}}(x_1), z_2 = \mathcal{E}_{\mathcal{G}}(x_2).$$

Finally, the attacker searches for a value of τ , large enough to make the salient characteristics of the interpolated image useless, however, small enough to ensure the content of the interpolation image $\mathcal{I}_{\mathcal{G}}(x_1, x_2, \tau)$ still agrees with the target label for humans. In their second approach, Turner et al. apply an adversarial transformation to each image before they apply the backdoor pattern. The goal is to make these images harder to classify correctly using standard image features, encouraging the model to memorize the backdoor pattern as a dominant feature. Formally, given a fixed classifier \mathcal{C} with loss \mathcal{L} and input x , they construct the adversarial perturbations as

$$x_{adv} = \arg \max_{\|x' - x\|_p \leq \epsilon} \mathcal{L}(x'),$$

for some l_p -norm and bound ϵ . Now the attacker retrieves a set of untargeted adversarial examples of the target class, and the attacker applies the trigger pattern to these adversarial examples which resemble the target class. Although both approaches allow for poisoning samples with the trigger containing the same label as the base image, the applied trigger has a visually noticeable shape and size in both types of clean label backdoor attacks. Thus, the attacker will still need to use a perceptible trigger pattern to inject and activate the backdoor, potentially compromising the secrecy of the attack.

Saha et al. [34] propose a clean label backdoor attack, whereby the attacker hides the trigger in the poisoned data and maintains secrecy of the trigger until test time. Saha et al. first define a trigger pattern p with a binary mask m (i.e., 1 at the location of the patch and 0 everywhere else), then apply the trigger p to a source image s_i from the source category. The patched source image \tilde{s}_i is

$$\tilde{s}_i = s_i \odot (1 - m) + p \odot m,$$

where \odot is for element-wise product. After retrieving the poisoned source image, the attacker solves an optimization problem over an image from the target class as

the poisoned image such that the l_2 distance of the patched source image \tilde{s} is close to the poisoned image z in the feature space, meanwhile, the l_∞ distance between the poisoned image and its initial image t is maintained less than a threshold ε . Formally, a poisoned image z can be defined as:

$$\begin{aligned} \arg \min_z & \|f(z) - f(\tilde{s})\|_2^2 \\ \text{st. } & \|z - t\|_\infty < \varepsilon, \end{aligned} \quad (1)$$

where $f(\cdot)$ is the intermediate features of the DNN and ε is a small value that ensures the poisoned image z is not visually distinguishable from the initial target image t . The optimization mentioned above only generates a single poisoned sample given a pair of images from source and target classes as well as a fixed location for the trigger. One can add this poisoned data with the correct label to the training data and train a backdoor model. However, such a model will only have the backdoor triggered when the attacker places the trigger at the same location on the same source image, limiting the practicality of the attack.

To address this shortcoming, Saha et al. [34] propose manipulating the poisoned images to be closer to the cluster of patched source images rather than being close to only the single patched source image. Inspired by universal adversarial examples [32], Saha et al. [34] minimize the expected value of the loss in Eq. (2) over all possible trigger locations and source images. In their extension, the attacker first samples \mathcal{K} random images t_k from the target class and initializes poisoned images z_k with t_k ; second, \mathcal{K} random images s_k is sampled from the source class and patched with triggers at randomly chosen locations to obtain \tilde{s}_k . For a given z_k in the poisoned image set, they search for a $s_{\tilde{a}(k)}$ in the patched image set which is close to z_k in the feature space $f(\cdot)$, as measured by Euclidean distance. Next, the attacker creates a one-to-one mapping $a(k)$ for the poisoned images set and the patched images set. Finally, the attacker performs one iteration of mini-batch projected gradient descent as follows:

$$\begin{aligned} \arg \min_z & \sum_{k=1}^{\mathcal{K}} \|f(z_k) - f(s_{\tilde{a}(k)})\|_2^2 \\ \text{st. } & \forall k : \|z_k - t_k\|_\infty < \varepsilon. \end{aligned} \quad (2)$$

Using the method above, the backdoor trigger samples are given the correct label and only used at test time.

2.6 Application Areas.

Most backdoor attacks and defenses select the image classification task to demonstrate the effectiveness of their attacks and defenses. However, other learning systems also are vulnerable to backdoor attacks in a range of application areas, e.g., object detection [14, 44], Natural Language Processing (NLP) [27, 9, 7, 24, 19], graph classification [50, 45], Reinforcement Learning (RL) [48, 43, 18] and Federated Learning (FL) [4, 47]. We will briefly survey how backdoors can be used to manipulate these learning systems.

Object Detection. Differing from image classification, object detection seeks to both detect the position of specific physical objects in an input image and predict the detected object’s label with some probability. Gu et al. [14] implemented their backdoor attack on a traffic sign detection and classification system in which images captured from a car-mounted camera. In their work, a stop sign is maliciously mis-classified as a speed-limit sign by the backdoored model. However, we note that the detected position of the signs remain unchanged, only the labels are mis-identified when triggers are present on the detected traffic signs. Wenger et al. [44] propose a credible backdoor attack against facial recognition systems in practice, in which 7 physical objects can trigger a change in an individual’s identity.

Natural Language Processing (NLP). There are several backdoor attacks against NLP systems [27, 9, 7, 24, 19]. Most of these works only explore the task of text classification, e.g. sentiment analysis on movie reviews [2], or hate speech detection on online social data [17]. Liu et al. [27] demonstrate the effectiveness of their backdoor attack on sentence attitude recognition. They use a crafted sequence of words at a fixed position as the trojan trigger. Dai et al. [9] inject the trojan into a LSTM-based sentiment analysis task. In this attack, the poisoned sentences need to be inserted into all positions of the given paragraph. Chen et al. [7] extend the trigger’s granularity from the sentence-level to a character level and word level. Lin et al. [24] compose two sentences that are dramatically different in semantics as triggers. Kurita et al. [19] introduce a trojan to pre-trained language models, whereby for different target classes, the attackers need to replace the token embedding of the triggers with their handcrafted embeddings.

Graph Classification. Zhang et al. [50] propose a subgraph based backdoor attack to Graph Neural Networks (GNNs), in which a GNN classifier predicts an attacker-chosen target label for a testing graph once a predefined subgraph is injected to the testing graph. Xi et al. [45] present a graph-oriented backdoor attack where triggers are defined as specific subgraphs including both topological structures and descriptive features.

Reinforcement Learning. Yang et al. [48] propose methods to discreetly introduce and exploit backdoor attacks within a sequential decision-making agent, by training multiple benign and malicious policies within a single long short-term memory (LSTM) network. Wang et al. [43] explore backdoor attacks on deep reinforcement learning based autonomous vehicles, where the malicious action include vehicle deceleration and acceleration to induce stop-and-go traffic waves to create traffic congestion. Kiourti et al. [18] present a tool for exploring and evaluating backdoor attacks on deep reinforcement learning agents, they evaluate their methods on a broad set of DRL benchmarks and show that after poisoning as little as 0.025% of the training data, the attacker can successfully inject the trojans into DRL models.

Federated Learning. In comparison to the traditional centralized machine learning setting, Federated Learning (FL) mitigates many systemic privacy risks and distributes computational costs. This has produced an explosive growth of federated learning research. The purpose of backdoor attacks in FL is that an attacker who controls one or several participants may manipulate their local models to simultaneously fit the clean and poisoned training samples. With the aggregation of local models

from participants into a global model at the server, the global model will have been influenced by the malicious models to behave maliciously on compromised inputs. Bagdasaryan et al. [4] are the first to mount a single local attacker backdoor attack against a FL platform via *model replacement*. In their attack, the attacker proposes a target backdoored global model \mathcal{X} they want the server to be in the next round. The attacker then scales up his local backdoored model to ensure it can survive the averaging step to ensure the global model is substituted by \mathcal{X} .

On the other hand, Xie et al. [47] propose a distributed backdoor attack (DBA) which decomposes a global trigger pattern into separate local patterns and uses these local patterns to inject into the training sets of different local adversarial participants. Fig. 1 illustrates the intuition of the DBA. As we can see, the attackers only need to inject a piece of the global trigger to poison their local models, such that the collective trigger is learned by the global model. Surprisingly, DBA can use a global trigger pattern to activate the ultimate global model as well as a centralized attack does. Xie et al. find that although no singular adversarial party had been poisoned by the global trigger under DBA, the DBA indeed can still behave maliciously as a centralized attack.

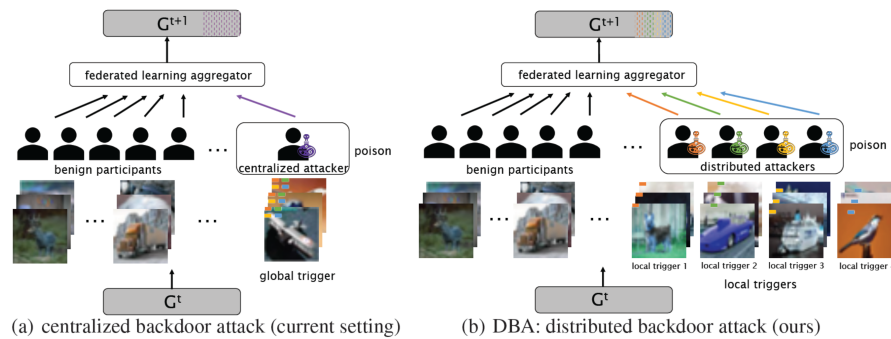


Fig. 1. Intuition of the Distributed Backdoor Attack (DBA) [47]. An attacker (orange) will poison a subset of his training data using only the trigger pattern located in the orange areas. The same reasoning applies to the remaining green, yellow, and blue marked attackers.

3 Detecting and Defending Backdoors

Attacks from white- to black-box settings have been developed to subvert the machine learning model to include backdoored behavior. However, any model trainer or holder may take proactive steps to detect and defend their models against this threat. This section will describe at length how this attack may be thwarted. Overall, the task of detecting and defending against a backdoor attack can be divided into three key sub-tasks:

1. **Task 1: Detecting the existence of the backdoor.** For a given model, it is difficult to know if the model is compromised (i.e., a model with a backdoor) or not. The first step of detecting and defending against the backdoor attack is to analyze the model and determine if there is a backdoor present in this model.
2. **Task 2: Identifying the backdoor trigger.** When a backdoor is detected in a model, the second step is usually to identify which pattern (including its size, location, texture, and so on) is used as the trigger.
3. **Task 3: Mitigating the backdoor attack.** After identifying the existence of a backdoor, the mitigation of such an attack is to remove the backdoor behavior from the model. Note that backdoor models can be made to be robust against transfer learning or fine-tuning [49].

Note that not all detection and defense techniques will support all three sub-tasks. As some may assume prior knowledge that a backdoor has already been detected, and the proposal only contains techniques to recover the trigger or mitigate the attack.

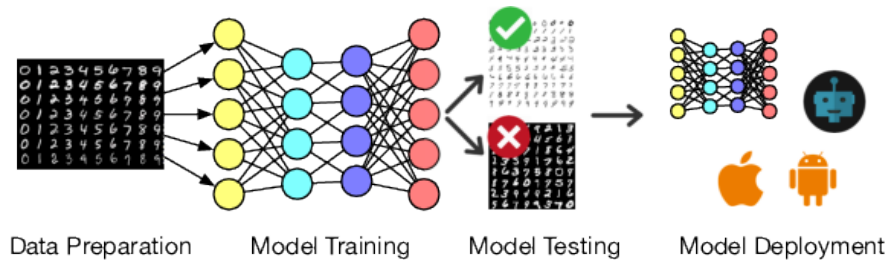


Fig. 2. Overview of DNN model training and deployment.

Figure 2 shows an overview of the DNN model training and deployment process. It can be broken up into four general steps from data preparation, model training, model testing, and model deployment. As discussed in Section 2, most existing poisoning attacks target the model training (or model retraining) step. Thus, investigating if the model contains a backdoor, reconstructing potential triggers, and/or mitigating any backdoor attacks must occur after this training step. Thus, mitigation strategies will be employed either during model testing (i.e., pre-deployment) or at the model’s runtime (i.e., post-deployment), and hence, depending on when the inspection occurs, existing detection and defense techniques can be divided into two categories: pre-deployment techniques or post-deployment techniques.

3.1 Pre-deployment Techniques

No Inspections

There exists work [51, 25] attempting to directly mitigate the backdoor attack without inspecting the model behavior. The key technique behind these methods is to

compress the model (e.g., by model pruning or similar techniques) or fine-tune the model with benign inputs to alter the model behavior hoping that the backdoor behavior is eliminated. Specifically, Zhao et al. [51] found that model pruning can remove some behaviors of a trained model, and potentially it can remove the backdoor of the model if pruning is purely using benign data.

Liu et al. [25] observe that pruning the model alone does not guarantee the removal of the model backdoor behavior. This is because the malicious model may use the same neuron to demonstrate both benign and malicious behaviors. Thus, if the neuron is removed, the model accuracy will be lower than that of the original model. This would be an undesirable consequence even though the model backdoor is removed. However, if this neuron is not pruned, the backdoor behavior is retained and the model continues to be malicious, also undesirable. Similarly, fine-tuning the model does not necessarily remove the model backdoor, as some attacks [49] target transfer learning scenarios where fine-tuning is needed. To solve this problem, Yao et al. propose Fine-Pruning, which combines the strengths of both fine-tuning and pruning to effectively nullify backdoors in DNN models. Fine-Pruning first removes backdoor neurons using pruning and fine-tuning the model in order to restore the drop in classification accuracy on clean inputs (which is introduced in the previous pruning procedure).

There are some limitations to these types of defenses. Firstly, model pruning itself has unknown effects on the model. Even though model accuracy after pruning does not decrease too much, many other important model properties, such as model bias (sometimes known as fairness) and model prediction performance, are not guaranteed to be the same. Using such models may potentially lead to severe consequences. Secondly, these mitigation techniques assume access to the training process and clean inputs, which conflicts with poisoning-based attacks.

Pre-deployment Model Inspections

Before the model is deployed, it is possible to check whether the model has been backdoored directly. This kind of strategy works without the running of the model, so it is also called static detection. For these types of techniques, some will require a large set of benign inputs to identify backdoors, such as Neural Cleanse (NC) [42], whereas others do not require much data (i.e., a limited number or even zero samples), such as ABS [26].

Neural Cleanse. Wang et al. [42] propose Neural Cleanse (NC), a pre-deployment technique to inspect DNNs, identify backdoors, and mitigate such attacks. Figure 3 illustrates the key observation that enables NC. The top figure shows a clean model with three output labels. If we want to perturb inputs belonging to C to A, more modification is needed to move samples across decision boundaries. The bottom figure shows the infected model, where the backdoor changes decision boundaries leading to a small perturbation value for changing inputs belonging to B and C to A.

Based on this observation, NC proposes to first compute a universal perturbation, which is the minimized amount of change to make the model predict a given target label. If the perturbation is small enough (i.e., smaller than a given threshold), NC

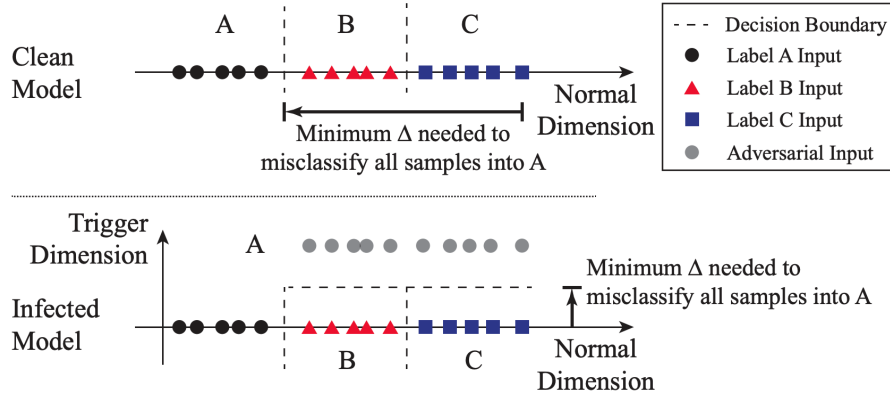


Fig. 3. Intuition of NC [42].

considers it as one trigger. It then verifies this by adding this trigger to a large number of benign inputs and tests if it is really a trigger and tries to optimize it based on prediction results. After identifying the trigger, it can mitigate the attack by either using a filter (i.e., to detect images with such a trigger pattern) or patching the DNN by removing the corresponding behaviors by pruning the neural network.

NC has a number of limitations. First, NC makes an incorrect assumption that if pixels in a small region have a strong influence on the output result, they are treated as backdoor triggers. This results in NC confusing triggers with strong benign features. In many tasks, there exist strong local features, where a region of pixels is important for one output label, for example, the antlers of deers in CIFAR-10. Secondly, NC assumes that the trigger has to be small and in the corner areas. These are heuristics, which do not hold for many attacks. For example, Salem et al. [35] propose a dynamic attack, where triggers can be added to different places and can successfully bypass NC. Thirdly, NC requires a significant number of testing samples to determine if a backdoor exists in a model or not. In real-world scenarios, such a large number of benign inputs may not exist. Lastly, it is designed purely for input space attacks, and it does not work for feature space attacks, such as using Nashville and Gotham filters as triggers [26].

ABS. ABS is built on top of two key observations. The first is that successful attacks entail compromised neurons. In existing attacks, the backdoored model recognizes the trigger as a strong feature of the target label to achieve a high attack success rate. Such a feature is represented by a set of inner neurons, which are referred to as *compromised neurons*. The second observation is that compromised neurons represent a subspace for the target label that cutcrosses the whole space. This idea is shown in Figure 4. The feature space surfaces of a benign model (left figure in Figure 4) and that of a backdoored model are noticeably different. For a backdoored model, there exists a cut of the surface that is significantly different from the benign model due to the injected backdoor. As it works for all inputs, it will affect every prediction results once it is activated. Thus, it will interact with the whole interface. The phenomenon

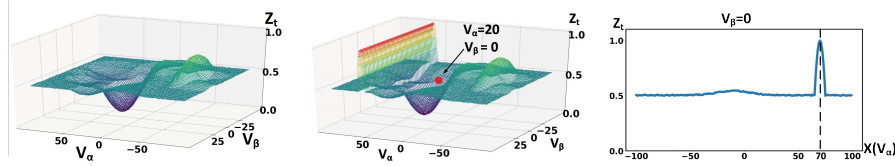


Fig. 4. Overview of ABS observations [26]. The left figure shows the feature surface of a benign model. The middle figure shows the feature surface of a model with a backdoor. The right figure shows a slice of surface for the backdoored model. The red dot in the middle figure shows a state where the attack happens, and it corresponds to the dashed line in the right figure.

is demonstrated in the right figure of Figure 4. When a neuron value is assigned to a special value, i.e., the trigger pixel value, the output will significantly deviate from normal.

Based on these observations, Liu et al. [26] propose Artificial Brain Stimulation (ABS). For any given input, ABS first predicts its label using the neural network. Then, it enumerates all neurons and performs a brain stimulation process. Namely, for each neuron, it tries to change its activation value to all possible values and simultaneously observes the value changes in the output. If there is one neuron whose behavior is similar to the right figure in Figure 4, ABS treats it as a backdoor. To reconstruct the backdoor trigger, ABS then performs a reverse engineering process, which will try to find an input pattern that can strongly activate these compromised neurons and trigger the attack.

ABS also introduces a new type of backdoor attack, which is the feature space attack. Namely, the trigger is no longer an input pattern (i.e., a region with specific pixel values), but feature space patterns represent high-level features (e.g., an image filter). However, this attack also has its own limitations. Firstly, it assumes one backdoor for each class. This may not hold in practice, and backdoors have been shown to be dynamic [35]. Secondly, it currently enumerates neurons one by one, assuming the presence of a strong correlation between one neuron and the backdoor behavior, which may be hidden or overridden by more advanced attacks.

3.2 Post-deployment Techniques

In addition to static approaches functioning before models are deployed, there is also work that monitors the model at runtime and determines if the model has a backdoor and more importantly, if it has been triggered by an input or not. In this setting, the defense or detection system can inspect individual inputs, offering a focused means of directly reconstructing the trigger by inspecting the attack input.

STRIP. Gao et al. [12] propose STRong Intentional Perturbation (STRIP), a run-time trojan attack detection system. The workflow of STRIP is shown in Figure 5. Firstly, STRIP will perturb each input by adding benign samples drawn from the test samples to obtain a list of perturbed inputs $X^{P_1}, X^{P_2}, \dots, X^{P_N}$. These inputs are the overlap of a benign input and the given input. Next, it will feed all these inputs to the DNN

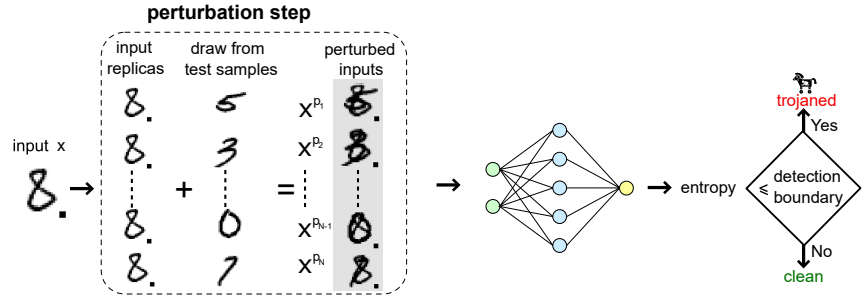


Fig. 5. Overview of STRIP [12].

model. Note that if the input contains a trigger, it is highly likely that a majority of the perturbed inputs will also yield predictions with the malicious output label result (due to the existence of the trigger), whereas for a benign input, the results are closer to random. As a result, STRIP only needs to examine every prediction result, and can then make a judgment on if the input will trigger the backdoor or not.

STRIP can effectively detect backdoor models and inputs that trigger the backdoor if the trigger lies in the corners of the image or at least does not overly overlap with the main contents. Such an example is shown in Figure 5). However, if the trigger does overlap with the contents (e.g., overlap with digits in Figure 5), the detection will fail because the texture of the trigger will also be changed by the perturbations. Salem et al.’s [35] proposed a dynamic backdoor attack that uses triggers that can be in the middle of the image.

4 Applications of Backdoors

4.1 Watermarking

Digital Watermarking conceals information in a piece of media (e.g., sound, video, or images) to enable a party to verify the authenticity or the originality of the media. This watermark, however, must also be resilient to tampering and other actors seeking to subvert the legitimate piece of media.

Adi et al. [1] propose an IP protection method for DNNs by applying the backdoor to watermark DNNs. They present cryptographic modeling for both tasks of watermarking and backdooring DNNs, and show that the former can be constructed from the latter (through a cryptographic primitive known as *commitment*) in a black-box manner. The definition of the backdoor attack that Adi et al. provided in a cryptographic framework is as follows: Given a trigger set T and a labeling function T_L , the backdoor shall be termed as $b = (T, T_L)$. The backdooring algorithm $Backdoor(O^f, b, M)$ is a probabilistic polynomial-time (PPT) algorithm that receives as input an oracle to f (ground-truth labeling function $f : D \rightarrow L$, where D is input

space, L is output space), the backdoor b and a model M , and outputs \hat{M} . \hat{M} is considered *backdoored* if

$$\begin{aligned} Pr_{x \in \bar{D} \setminus T} [f(x) \neq \text{Classify}(\hat{M}, x)] &\leq \varepsilon, \\ Pr_{x \in T} [T_L(x) \neq \text{Classify}(\hat{M}, x)] &\leq \varepsilon, \end{aligned} \tag{3}$$

where \bar{D} is the meaningful input, $\text{Classify}(M, x)$ is a deterministic polynomial-time algorithm that, for an input $x \in D$ outputs a value $M(x) \in L \setminus \{\perp\}$, and \perp is an undefined output label. This definition presents two ways to embed a backdoor. The first is that the backdoor is implanted into a *pre-trained model*. The second is the adversary can train a new model from scratch.

A watermarking scheme can be split into three key components.

1. Generation of the secret “marking” key mk . This key will be embedded as the watermark. A public verification key vk is also generated and will be used later to detect the watermark. In watermarking via backdoors, the backdoor is the marking key, while a commitment (the cryptographic primitive) used to generate the backdoor is the verification key.
2. Embedding the watermark (a backdoor b) into a DNN model. Through poisoned training data or retraining, as previously described in Section 2.1, the watermark (backdoor) can be embedded.
3. Verifying the presence of the watermark. Provided mk, vk , for a backdoor test $b = (T, T_L)$. If $\forall t^{(i)} \in T : T_L^{(i)} \neq f(t^{(i)})$, proceed to the next step, otherwise, the verification fails. Despite the detection of the watermark, one must verify the integrity of the commitment, i.e., if it was tampered or not. In the final step, the accuracy of the algorithm is verified. For all $i \in 1, \dots, n$, if more than $\varepsilon|T|$ elements from T does not satisfy $\text{Classify}(t^{(i)}, M) = T_L^{(i)}$, then the verification fails, otherwise the commitment has been successfully verified.

Adi et al. [1] prove their method upholds the properties of:

- *Functionality-preserving*: the prediction accuracy of the model should not be negatively influenced by the presence of the watermark.
- *Unremovability*: an adversary with full knowledge of the watermark generation process should not be able to remove the watermark from the model.
- *Unforgeability*: an adversary with only the verification key should not be able to demonstrate ownership of the marking key.
- *non-trivial ownership*: with knowledge of the watermark generation algorithm, a third party should not be able to generate marker and verification key pairs, and claim models for future models.

Li et al. [21], however, observe that the watermarking system proposed by Adi et al. [1] makes the assumption that only one backdoor (watermark) may be inserted into the model. For example, Salem et al. [35]’s Dynamic Backdoors contain multiple backdoors. The existence of multiple backdoors would result in multiple valid watermarks, and thus void the *Unforgeability* claim. The insertion of multiple backdoors

would also impact the *Unremovability* of the original backdoor, otherwise termed as the persistence of the watermark. In response, Li et al. leverage two data pre-processing techniques that use out-of-bound values and null-embedding to improve the persistence of the watermark against other attackers and limit the effects of re-training in the event that another backdoor is to be injected on top of the existing backdoor. Li et al. also introduce *wonder filters*, a primitive to enable the embedding of bit-sequences (from the marker key) into the model.

The largest hurdle to overcome in the application of the backdoor attack as a means to watermark DNNs, is that neural networks are fundamentally designed to be tuned and trained incrementally. Li et al. propose a *model piracy* attack setting whereby an adversary wants to stake its own ownership claims on the model, or destroy the original owner’s claims. To defend against this attack, Li et al. design a DNN watermarking system based on wonder filters that strongly authenticates owners by embedding (into the DNN) a filter described by the owner’s private key. Where Li et al’s work differs from Adi et al. is in the *wonder filter* W , which is a two-dimensional digital filter that can be applied to any input image. This filter will have 3 possible permutations for each pixel, transparent, positive change, or negative change, with a majority of filter pixels being transparent. Thus, W is defined by the position, size, and values of a 0/1 bit pattern block.

When Li et al. apply *out-of-bound values*, they translate the 0/1 bit pattern of W as out-of-bound values in the input images. A set of training data is processed with the filter. They then flip the values of the *wonder filter* to create an *inverted wonder filter* W^- . The *inverted filter* W^- is then applied to the same set of training data processed by W . The set of images filtered by W are labeled as the target class label, while the W^- filtered data is labeled as the original class label before the data is used to train the backdoored model. As for the normal and null embeddings approach, the normal and null embeddings serve complementary objectives. The normal embedding injects the desired marker into the model, while the null embedding “locks down” the model, so no additional watermarks may be added.

Li et al’s process of watermarking the image is similar to Adi et al’s process, with the same three key processes of generating the secret “marking”, or in this instance, the wonder filter W , embedding the watermark (and/or additionally locking down the model), and finally, the process of verifying the watermark, by using the image to compute W and an associated label. After applying W to a random set of images, it is expected that an authentic watermark should yield a majority of the target class label, instead of a random assortment of classes as expected from a random set of images, without W .

Li et al. also provide a security analysis to prove that their approach can uphold the requirements of *reliability*, *no false positives*, *unforgeability*, and *persistence*, whereby *Reliability* describes that for a given input x , a poisoned input ($x \oplus W$, or $x \oplus W^-$), the backdoored DNNs will produce the predefined output in a deterministic manner. *No False Positives* denotes that a verifier should not be capable of judging a clean model as the watermarked model. *unforgeability* ensures that the watermark injected on a DNN has a strong association with its owner, and *Persistence* guarantees that the watermark embedded cannot be corrupted or removed by an adversary.

4.2 Adversarial Example Detection

In Gotta Catch [38], Shan et al. observe that the backdoor attack will alter the decision boundary of the DNN models. Following the injection of a backdoor, the decision boundary of the original clean model will mutate. The mutation will result in triggers establishing shortcuts in the decision boundary of the backdoored model.

On the contrary, there is a common approach of adversarial attacks to find adversarial examples; for example, universal adversarial attacks [31, 36], will try to iteratively search the whole dataset to find similar shortcuts to use as their universal adversarial examples. Based on this observation, the shortcut created by a backdoor can act as a trapdoor to capture the adversarial attacker’s optimization process, detect, and/or recover from the adversarial attack [37]. The trapdoor implementation uses techniques similar to those found in BadNets backdoor attacks [14]. The authors define the trapdoor perturbation (the trigger) from multiple dimensions, e.g., mask ratio, size, pixel intensities, and relative locations.

4.3 Open Problems

- *Fair comparison of methods including attacks and defenses.* We have observed experimental settings (e.g. models and datasets) vary greatly across studies of different domains. Consequently to inspect the potential of AI Trojans in a standardized manner, The Intelligence Advanced Research Projects Activity (IARPA) has recently launched a program (and competition), Trojans in Artificial Intelligence (TrojAI) [16] , .
- *Persisting through fine-tuning in transfer learning.* A model backdoor can be made ineffective after only a few fine-tuning layers [14, 49]. Unfortunately, this limits the practicality of the backdoor attack in the transfer learning setting. More advanced backdoor attacks that can persist through such fine-tuning processes remains a challenge.
- *Combining hidden triggers and clean labels.* To inject trojans into DNNs by poisoning the training data, the most stealthy method is to attach the triggers on the poisoned data in an imperceptible way, in addition to a correctly annotated label. Existing attacks either visually hide the triggers but retain an clearly poisoned label or correctly annotate the label with noticeable triggers. The combination of these stealth tactics, to produce poisoned data with an invisible trigger and clean labels still remains a challenge.
- *Accessing training data.* There are a variety of backdoor attacks, however only a small set has access to clean samples. In most security-sensitive cases, the attacker can only access a pre-trained model. The injection of a trojan by directly compromising the model weights without requiring access to the clean original training data still remains challenge.
- *Adaptiveness of Attacks and Defenses.* Existing attacks and defenses often discuss their effectiveness against reactive (dynamic) attack or defense countermeasures in a superficial and heuristic way, without including the adversary’s possible countermeasures as a part of their works. These adaptive attacks and defenses can adaptively take optimal strategies when responding to an adversary’s

countermeasures, will provide improved practicality in realistic settings, but still remains an open challenge.

References

1. ADI, Y., BAUM, C., CISSÉ, M., PINKAS, B., AND KESHET, J. Turning your weakness into a strength: Watermarking deep neural networks by backdooring. In *27th USENIX Security Symposium, USENIX Security 2018, Baltimore, MD, USA, August 15-17, 2018*. (2018), pp. 1615–1631.
2. ALZANTOT, M., SHARMA, Y., ELGOHARY, A., HO, B., SRIVASTAVA, M. B., AND CHANG, K. Generating natural language adversarial examples. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018* (2018), E. Riloff, D. Chiang, J. Hockenmaier, and J. Tsujii, Eds., Association for Computational Linguistics, pp. 2890–2896.
3. BAGDASARYAN, E., AND SHMATIKOV, V. Blind backdoors in deep learning models. *arXiv preprint arXiv:2005.03823* (2020).
4. BAGDASARYAN, E., VEIT, A., HUA, Y., ESTRIN, D., AND SHMATIKOV, V. How to backdoor federated learning. In *International Conference on Artificial Intelligence and Statistics* (2020), pp. 2938–2948.
5. BARNI, M., KALLAS, K., AND TONDI, B. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)* (2019), IEEE, pp. 101–105.
6. CHEN, X., LIU, C., LI, B., LU, K., AND SONG, D. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526* (2017).
7. CHEN, X., SALEM, A., BACKES, M., MA, S., AND ZHANG, Y. BadNL: Backdoor attacks against nlp models. *arXiv preprint arXiv:2006.01043* (2020).
8. COSTALES, R., MAO, C., NORWITZ, R., KIM, B., AND YANG, J. Live trojan attacks on deep neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops* (2020), pp. 796–797.
9. DAI, J., CHEN, C., AND LI, Y. A backdoor attack against LSTM-Based text classification systems. *IEEE Access* 7 (2019), 138872–138878.
10. DÉSIDÉRI, J.-A. Multiple-gradient descent algorithm (mgda) for multiobjective optimization. *Comptes Rendus Mathématique* 350, 5-6 (2012), 313–318.
11. DUMFORD, J., AND SCHEIRER, W. J. Backdooring convolutional neural networks via targeted weight perturbations. In *2020 IEEE International Joint Conference on Biometrics, IJCB 2020, Houston, TX, USA, September 28 - October 1, 2020* (2020), IEEE, pp. 1–9.
12. GAO, Y., XU, C., WANG, D., CHEN, S., RANASINGHE, D. C., AND NEPAL, S. Strip: A defence against trojan attacks on deep neural networks. In *Proceedings of the 35th Annual Computer Security Applications Conference* (2019), pp. 113–125.
13. GOODFELLOW, I., POUGET-ABADIE, J., MIRZA, M., XU, B., WARDE-FARLEY, D., OZAIR, S., COURVILLE, A., AND BENGIO, Y. Generative adversarial nets. In *Advances in neural information processing systems* (2014), pp. 2672–2680.
14. GU, T., DOLAN-GAVITT, B., AND GARG, S. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *IEEE Access* 7 (2019), 47230–47244.
15. HEFFNER, C. Binwalk: Firmware analysis tool. URL: <https://code.google.com/p/binwalk/>(visited on 03/03/2013) (2010).
16. IARPA. Trojans in artificial intelligence (trojai).

17. KAGGLE. Toxic comment classification challenge, 2020. <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge/> Accessed June 24, 2020.
18. KIOURTI, P., WARDEGA, K., JHA, S., AND LI, W. Trojdl: Evaluation of backdoor attacks on deep reinforcement learning. In *2020 57th ACM/IEEE Design Automation Conference (DAC)* (2020), pp. 1–6.
19. KURITA, K., MICHEL, P., AND NEUBIG, G. Weight poisoning attacks on pretrained models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics* (Online, July 2020), Association for Computational Linguistics, pp. 2793–2806.
20. LANGNER, R. Stuxnet: Dissecting a cyberwarfare weapon. *IEEE Security & Privacy* 9, 3 (2011), 49–51.
21. LI, H., WILLSON, E., ZHENG, H., AND ZHAO, B. Y. Persistent and unforgeable watermarks for deep neural networks. *arXiv preprint arXiv:1910.01226* (2019).
22. LI, S., XUE, M., ZHAO, B., ZHU, H., AND ZHANG, X. Invisible backdoor attacks on deep neural networks via steganography and regularization. *IEEE Transactions on Dependable and Secure Computing* (2020), 1–1.
23. LIGH, M. H., CASE, A., LEVY, J., AND WALTERS, A. *The art of memory forensics: detecting malware and threats in windows, linux, and Mac memory*. John Wiley & Sons, 2014.
24. LIN, J., XU, L., LIU, Y., AND ZHANG, X. Composite backdoor attack for deep neural network by mixing existing benign features. In *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security* (2020), pp. 113–131.
25. LIU, K., DOLAN-GAVITT, B., AND GARG, S. Fine-pruning: Defending against backdoor attacks on deep neural networks. In *International Symposium on Research in Attacks, Intrusions, and Defenses* (2018), Springer, pp. 273–294.
26. LIU, Y., LEE, W.-C., TAO, G., MA, S., AAFER, Y., AND ZHANG, X. Abs: Scanning neural networks for back-doors by artificial brain stimulation. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (2019), pp. 1265–1282.
27. LIU, Y., MA, S., AAFER, Y., LEE, W.-C., ZHAI, J., WANG, W., AND ZHANG, X. Trojanning attack on neural networks. *The Network and Distributed System Security Symposium (NDSS)* (2017).
28. LIU, Y., MA, S., AAFER, Y., LEE, W.-C., ZHAI, J., WANG, W., AND ZHANG, X. Trojanning attack on neural networks. In *25nd Annual Network and Distributed System Security Symposium, NDSS 2018, San Diego, California, USA, February 18-22, 2018* (2018), The Internet Society.
29. LIU, Y., MONDAL, A., CHAKRABORTY, A., ZUZAK, M., JACOBSEN, N., XING, D., AND SRIVASTAVA, A. A survey on neural trojans. *IACR Cryptol. ePrint Arch.* 2020 (2020), 201.
30. MIRZA, M., AND OSINDERO, S. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784* (2014).
31. MOOSAVI-DEZFOOLI, S., FAWZI, A., FAWZI, O., AND FROSSARD, P. Universal adversarial perturbations. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017* (2017), pp. 86–94.
32. MOOSAVI-DEZFOOLI, S.-M., FAWZI, A., FAWZI, O., AND FROSSARD, P. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (2017), pp. 1765–1773.

33. REZENDE, D. J., MOHAMED, S., AND WIERSTRA, D. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082* (2014).
34. SAHA, A., SUBRAMANYA, A., AND PIRSIYAVASH, H. Hidden trigger backdoor attacks. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020* (2020), AAAI Press, pp. 11957–11965.
35. SALEM, A., WEN, R., BACKES, M., MA, S., AND ZHANG, Y. Dynamic backdoor attacks against machine learning models. *arXiv preprint arXiv:2003.03675* (2020).
36. SHAFABI, A., NAJIBI, M., XU, Z., DICKERSON, J. P., DAVIS, L. S., AND GOLDSTEIN, T. Universal adversarial training. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020* (2020), AAAI Press, pp. 5636–5643.
37. SHAN, S., WENGER, E., WANG, B., LI, B., ZHENG, H., AND ZHAO, B. Y. Gotta catch ‘em all: Using honeypots to catch adversarial attacks on neural networks. *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security, CCS (2020)*.
38. SHAN, S., WENGER, E., WANG, B., LI, B., ZHENG, H., AND ZHAO, B. Y. Gotta catch ‘em all: Using honeypots to catch adversarial attacks on neural networks. In *CCS ’20: 2020 ACM SIGSAC Conference on Computer and Communications Security, Virtual Event, USA, November 9-13, 2020* (2020), ACM, pp. 67–83.
39. SHOKRI, R., STRONATI, M., SONG, C., AND SHMATIKOV, V. Membership inference attacks against machine learning models. In *IEEE Symposium on Security and Privacy* (2017), IEEE, pp. 3–18.
40. TAN, T. J. L., AND SHOKRI, R. Bypassing backdoor detection algorithms in deep learning. In *IEEE European Symposium on Security and Privacy, EuroS&P 2020, Genoa, Italy, September 7-11, 2020* (2020), IEEE, pp. 175–183.
41. TURNER, A., TSIPRAS, D., AND MADRY, A. Clean-label backdoor attacks.
42. WANG, B., YAO, Y., SHAN, S., LI, H., VISWANATH, B., ZHENG, H., AND ZHAO, B. Y. Neural cleanse: Identifying and mitigating backdoor attacks in neural networks. In *2019 IEEE Symposium on Security and Privacy (SP)* (2019), IEEE, pp. 707–723.
43. WANG, Y., SARKAR, E., MANIATAKOS, M., AND JABARI, S. E. Stop-and-go: Exploring backdoor attacks on deep reinforcement learning-based traffic congestion control systems. *arXiv preprint arXiv:2003.07859* (2020).
44. WENGER, E., PASSANANTI, J., YAO, Y., ZHENG, H., AND ZHAO, B. Y. Backdoor attacks on facial recognition in the physical world. *arXiv preprint arXiv:2006.14580* (2020).
45. XI, Z., PANG, R., JI, S., AND WANG, T. Graph backdoor. In *30th USENIX Security Symposium, USENIX Security 2021* (2021).
46. XIANG, Z., MILLER, D. J., AND KESIDIS, G. Revealing backdoors, post-training, in DNN classifiers via novel inference on optimized perturbations inducing group misclassification. *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (2020).
47. XIE, C., HUANG, K., CHEN, P.-Y., AND LI, B. DBA: Distributed backdoor attacks against federated learning. In *International Conference on Learning Representations* (2019).
48. YANG, Z., IYER, N., REIMANN, J., AND VIRANI, N. Design of intentional backdoors in sequential models. *CoRR abs/1902.09972* (2019).
49. YAO, Y., LI, H., ZHENG, H., AND ZHAO, B. Y. Latent backdoor attacks on deep neural networks. In *Proceedings of the 2019 ACM SIGSAC Conference on Computer and Communications Security* (2019), pp. 2041–2055.

50. ZHANG, Z., JIA, J., WANG, B., AND GONG, N. Z. Backdoor attacks to graph neural networks. *arXiv preprint arXiv:2006.11165* (2020).
51. ZHAO, B., AND LAO, Y. Resilience of pruned neural network against poisoning attack. In *2018 13th International Conference on Malicious and Unwanted Software (MALWARE)* (2018), IEEE, pp. 78–83.
52. ZHONG, H., LIAO, C., SQUICCIARINI, A. C., ZHU, S., AND MILLER, D. J. Backdoor embedding in convolutional neural network models via invisible perturbation. In *CO-DASPY '20: Tenth ACM Conference on Data and Application Security and Privacy, New Orleans, LA, USA, March 16-18, 2020* (2020), ACM, pp. 97–108.