

# When Forgetting Triggers Backdoors: A Clean Unlearning Attack

Marco Arazzi, Antonino Nocera, Vinod P.

**Abstract**—Machine unlearning has emerged as a key component in ensuring “Right to be Forgotten”, enabling the removal of specific data points from trained models. However, even when the unlearning is performed without poisoning the forget-set (clean unlearning), it can be exploited for stealthy attacks that existing defenses struggle to detect. In this paper, we propose a novel *clean* backdoor attack that exploits both the model learning phase and the subsequent unlearning requests. Unlike traditional backdoor methods, during the first phase, our approach injects a weak, distributed malicious signal across multiple classes. The real attack is then activated and amplified by selectively unlearning *non-poisoned* samples. This strategy results in a powerful and stealthy novel attack that is hard to detect or mitigate, highlighting critical vulnerabilities in current unlearning mechanisms and highlighting the need for more robust defenses.

## I. INTRODUCTION

As the integration of AI-powered services and applications continues to grow in the modern IT landscape, the demand for high-quality and rich datasets to train and fine-tune machine/deep learning models is growing, as well. In line with recent regulations and laws, such as the European General Data Protection Regulation (GDPR) and the California Privacy Rights Act (CPRA), in the context of machine/deep learning, machine *unlearning* has been introduced as an efficient approach to allow compliance with data protection regulations in AI-based systems and implement the “Right to be Forgotten” [1]. Various unlearning methods have been proposed to remove specific data contributions from trained models while preserving utility [2], [3], [4]. However, changes introduced by unlearning can create security vulnerabilities [5], and recent work shows that unlearning itself can be exploited to launch backdoor attacks [6], [7], [8]. These approaches differ from traditional backdoor attacks, as they split the attack strategy into two phases: (a) the learning phase, in which the attacker hides an invisible trigger to a portion of his controlled training data and (b) the unlearning phase, in which the attacker requires an unlearning procedure on suitably crafted/chosen data points whose removal ultimately boosts the performance of the trigger injected during the previous phase and, therefore, allows the activation of the backdoor. In this setting, we distinguish between two types of backdoor attack on machine unlearning, based on the strategy adopted in the second phase: *clean* [6] versus *poisoned* unlearning [7], [8]. Although the two categories share some similarity, the use of poisoning strategies during the unlearning phase

limits the attack plausibility in scenarios where appropriate *aggressive* defenses can be deployed. Indeed, while it may seem intuitive to improve anomaly detection in the forget set and filter out suspicious requests, defenses become significantly limited when forget requests consistently involve legitimate data. However, inspired by the *clean* attack strategy proposed by [6], we believe that no defense can actually fully prevent an attack that does not poison the forget-set during unlearning. The approach of [6], although very effective in some cases, attempts to carry out a weak and stealthy backdoor attack on one target class during training, which remains under detectability thresholds of known defenses. Then, during the second phase, it strives to amplify the attack success rate by suitably performing unlearning request of clean samples of the same target class. According to our point of view, the idea of performing a standard backdoor attack, although weakened to remain unnoticeable, is limiting and prevents obtaining a fully powerful global attack. In this paper, we hence propose a novel combined attack exploiting the two phases mentioned before and maintaining a fully *clean* unlearning phase. The main idea behind our proposal is that we do not actually need to carry out a target backdoor attack during the training phase, but it is sufficient to hide an undetectable malicious signal spread across multiple classes, rather than only the target one. In this way, we can boost the performance of the final attack, making it powerful, stealthy, and very hard to block. The main contribution of our approach are hence as follows.

- We propose “UNlearning-activated CLEAN backdoor attack” (UNCLEAN), a powerful novel combined attack against machine unlearning, exploiting both the learning and unlearning phases.
- We extend the existing literature by both exploiting only *unaltered* samples during the unlearning phase, thus maintaining this phase clean, and by adopting a non-targeted malicious noise injection to broaden the capacity of our attack.
- We test the performance of our attack with the main existing unlearning strategies and model architectures. We prove its robustness against existing defenses and, finally, we show the performance advantages with respect to related existing approaches.

The experimental results show the effectiveness of the proposed attack and its severity level in the context of machine unlearning. Moreover, they demonstrate the superiority of our attack over existing methods, achieving an improvement in the attack success rate of more than 32% compared to the previous work by [6].

Marco Arazzi and Antonino Nocera are with the Department of Electrical, Computer and Biomedical Engineering, University of Pavia, Italy. (email: marco.arazzi01@universitadipavia.it, antonino.nocera@unipv.it)

Vinod P. is with the Department of Computer Applications, Cochin University of Science and Technology, Kerala, India. (email: vinod.p@cusat.ac.in)

## II. RELATED WORKS

Early methods for machine unlearning primarily focus on altering training labels or reversing the learning process. Random Label unlearning [9] removes the influence of specific data by randomly reassigning labels, causing the model to unlearn meaningful associations with the targeted samples. Similarly, Gradient Ascent unlearning [3] takes the opposite approach of standard training by maximizing the loss for the forgotten data, effectively erasing its impact. While both techniques can achieve unlearning, they often require careful tuning to prevent performance degradation and still involve some degree of retraining. To enhance efficiency, Fisher Forgetting [3] modifies model weights to remove the influence of specific data points. This technique does not require access to the original training data but instead relies on the model’s learned weights and the Fisher Information matrix to adjust parameters in a way that effectively removes the targeted data’s impact. More advanced strategies manipulate the model’s decision boundaries directly. Boundary Unlearning [4] provides a fast and effective approach to forgetting an entire class from a trained deep neural network. Rather than retraining the model from scratch, this method strategically adjusts decision boundaries so that the model no longer retains knowledge of the forgotten class while preserving its accuracy on the remaining data. One of the most sophisticated approach is Bad Teacher unlearning [2], which employs a student-teacher framework to facilitate unlearning without requiring full retraining. In this method, a student model is trained using both competent and incompetent teachers. The competent teacher provides accurate knowledge, while the incompetent teacher introduces deliberate inaccuracies, guiding the student model to forget specific information.

Backdoor attacks represent a class of adversarial strategies in which an adversary injects a malicious trigger into the training process to implant hidden functionality within the model. During inference, the model preserves standard behavior on clean, unaltered inputs, but misclassifies inputs containing the trigger in a way predefined by the attacker. Such attacks are effective in both centralized and distributed learning paradigms [10], [11], [12], [13]. BadNet [14], has been explored in recent research to highlight security risks associated with outsourcing model training. In a BadNet attack, a model is compromised during training, either by an untrusted third-party service or through transfer learning using a pre-trained model. This attack differs from adversarial perturbations, as it embeds the trigger within the model itself, rather than altering the input data. Blended [15], [16] backdoor poisoning attacks aim to remain stealthy by injecting only a few poisoned samples with subtle backdoor triggers that are difficult to detect. These attacks are often studied under weak threat models, where the attacker has no knowledge of the victim model or training data. Two key strategies dominate this space: (i) using a single input instance as a universal key, and (ii) employing a trigger pattern to synthesize multiple poisoned instances. We adopt a similar stealth-oriented approach, aiming to hide the presence of the backdoor, thereby aligning with these covert poisoning strategies. Clean-label [17], [18], [19]

backdoor attacks are particularly stealthy, as the poisoned data retain their original labels, making detection challenging. Notably, research has demonstrated that such attacks can be executed with minimal knowledge of the training data, using as little as 0.5% poisoned samples to effectively implant a backdoor. Similarly, we employ clean-label backdoor techniques to covertly implant triggers into the training process, thereby aligning with stealth-focused attack paradigms and reinforcing the critical security threats such attacks pose to machine learning systems.

Machine unlearning has become a pivotal mechanism for eliminating specific data from trained models to uphold user privacy and ensure regulatory compliance. However, this capability also introduces novel security vulnerabilities—most notably, backdoor attacks that exploit the unlearning process itself to subvert model integrity. As demonstrated in [6], an adversary can activate a backdoor in the unlearned model by merely requesting the deletion of a small portion of their previously contributed training data. This attack bypasses traditional poisoning techniques and achieves malicious influence without modifying any initial training samples, thereby presenting a covert and potent threat to unlearning-based systems. In [7], the authors present a novel black-box backdoor attack that exploits machine unlearning to activate malicious behavior. The attacker injects both poison and mitigation samples into the training dataset to construct a model that appears benign. Subsequently, by issuing unlearning requests targeting the mitigation samples, the attacker incrementally removes their neutralizing effect, thereby triggering the latent backdoor.

Similarly, Huang et al. [8] propose UBA-Inf, a backdoor attack that leverages influence-driven camouflage within the model, activated through unlearning operations. Unlike conventional backdoor strategies, UBA-Inf embeds the trigger in a highly stealthy manner, enabling fine-grained control over backdoor activation. The attacker accomplishes this by selectively unlearning camouflaged samples, which not only conceals the backdoor during standard operation but also mitigates challenges such as premature exposure and backdoor vanishing. This method significantly improves the persistence, stealth, and effectiveness of backdoor attacks in the presence of unlearning mechanisms.

The approach presented in this paper builds on previous methods but sets a more ambitious goal: maintaining the backdoor’s stealth during training and activating it through clean data unlearning, without relying on camouflage or mitigation. This effectively proposes a “clean” unlearning backdoor attack.

## III. METHODOLOGY

This section details the methodology employed to execute our proposed attack.

### A. Preliminaries

1) *Threat Model*: In this paper, we focus on an unlearning-activated backdoor attack that shows its full potential after the forgetting process limiting, intentionally, its performance after the initial training of the target model. Specifically, we

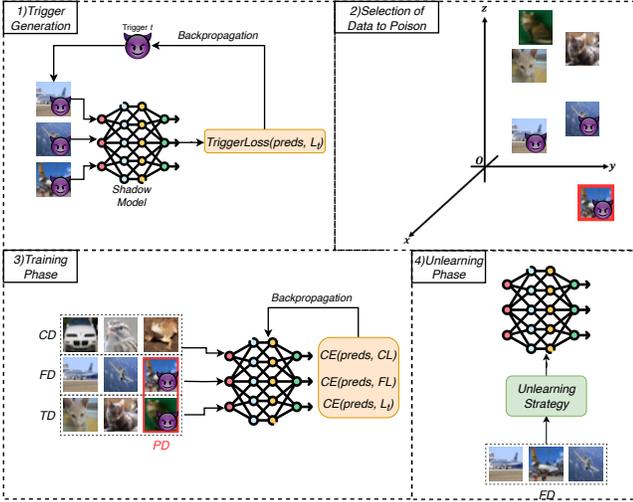


Fig. 1: Pipeline (the trigger depicted in the figure serves as an example since the actual trigger is not visible). 1) Frequency Trigger generation process described in Section III-B1. 2) Selection of the best candidate data to poison in the forget-set as described in Section III-B2. 3) Training the model using the poisoned set, 4) Unlearning the forget set using only clean data.

assume that the attacker has access to the clean training dataset ( $CD$ ) used to train the target model ( $M$ ). Under this assumption, the attacker constructs a poisoned dataset  $PD = poison\_function(CD, t)$  by injecting a predefined trigger  $t$  into samples from  $CD$ , while preserving their original clean labels ( $CL$ ). This clean-label poisoning strategy embeds a backdoor into  $M$  without modifying label semantics, thereby maintaining the model’s performance on benign inputs and enhancing the stealth of the attack.

$$preds = M(PD); M \leftarrow loss\_function(preds, CL) \quad (1)$$

The attacker initiates the unlearning process by submitting a forget set ( $FD$ ), which contains samples from any class except the target class. The model  $M$  undergoes unlearning to remove the influence of  $FD$ , resulting in an updated model  $M'$ . This unlearned model  $M' = unlearn\_strategy(M, FD)$  retains high accuracy on the retain set ( $RD$ ), which includes all non-forgotten data, but it intentionally loses the ability to correctly classify samples from  $FD$ .

2) *Attacker’s Goal*: Targeted backdoor attacks generally focus on achieving a high rate of success against a specific target class using triggered data, while maintaining the model’s accuracy on unaltered data. In this context, the attacker’s goal is to embed a backdoor in the model by adding a trigger during the initial training phase, while keeping the original labels unchanged. The attack is specifically designed to underperform during training and fully manifest its effectiveness after the unlearning process. To enhance the attack’s stealthiness, the forget-set provided by the attacker should consist of clean data to evade detection by filtering defenses at this stage and only need to bypass them during the training phase.

## B. Unlearning-activated Clean backdoor attack

The core concept of this study is based on launching an unlearning-activated backdoor attack that avoids tampering with the unlearning phase (i.e., the forget-set  $FD$ ) while only introducing poison to the initial dataset utilized for training. With this, we want to consider the most realistic scenario in which the attacker controls just one of the sources and has access only to partial data  $TD \in CD$  from the target class  $TC$  and a subset of data from the remaining classes from which the forget-set  $FD$  can be selected.

Our attack consists of three main phases: trigger generation, selection of data to poison, and the training/unlearning phase (see Figure 1).

1) *Trigger Generation*: Although the backdoor can be activated using clean-label inputs, the attacker must poison the training data, making the attack difficult to detect with existing defenses or manual inspection. To enhance stealth, we generate the trigger  $t$  in the frequency domain with the same shape of the image, using the Discrete Cosine Transform (DCT), which reduces visual artifacts and preserves the natural appearance of the image. Operating in the frequency domain allows us to embed the trigger by modifying selected mid-frequency components, which are less perceptible to the human eye yet influential to neural network activations. Prior work has demonstrated that frequency-based perturbations, particularly those applied via DCT, can effectively produce stealthy triggers while remaining robust under common image transformations like compression and resizing [20], [21].

Specifically, we embed the trigger  $t$  into an image  $I$  by modifying selected frequency components using the Discrete Cosine Transform (DCT). We begin by transforming both the image  $I$  and the trigger  $t$  into the frequency domain:  $F_I = DCT(I)$ ,  $F_t = DCT(t)$ .

Next, we generate a frequency mask  $Mask(u, v)$  to constrain modifications to a specific frequency band ( $f_{min}, f_{max}$ ), ensuring that only the desired frequency components are altered. Here,  $u$  and  $v$  denote the frequency indices corresponding to the horizontal and vertical components, respectively. The mask  $Mask(u, v)$  is defined as follows:

$$Mask(u, v) = \begin{cases} 1, & f_{min} \leq \sqrt{u^2 + v^2} \leq f_{max} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

We constrain the trigger to a specific frequency band ( $f_{min}, f_{max}$ ) to limit modifications to perceptually inconspicuous components, thereby enhancing stealthiness. Additionally, we introduce a learnable parameter  $\alpha$ , bounded by a sigmoid function, to control the trigger’s strength in a smooth and differentiable manner. This approach allows adaptive, low-magnitude perturbations that effectively embed the backdoor while minimizing visual and statistical detectability. The learnable parameter can be obtained as:  $\alpha = \sigma(\theta)$ , where  $\alpha \in (0, 1)$ ,  $\sigma$  is the sigmoid function and  $\theta$  represents the trainable parameters. Using the mask generated before the trigger is applied to the image as follows:

$$F_{poisoned} = F_I + \alpha \cdot M \cdot (F_t - F_I), \quad (3)$$

this ensures that only the selected frequencies are affected. The obtained  $F_{poisoned}$  is then reconverted in the spatial domain using the *Inverse DCT (IDCT)*:  $I' = IDCT(F_{poisoned})$ .

We clamp the reconstructed image  $I' \in PD$  to the range  $[0, 1]$  to maintain visual consistency. To further enhance the stealthiness of the attack, we apply regularization to both the blending coefficient  $\alpha$  and the trigger  $t$ , thereby constraining the perturbations introduced to the image. In particular,  $R_\alpha = \lambda_\alpha |\alpha|$  and  $R_t = \lambda_t \|t\|_2$ , where  $\lambda_\alpha$  and  $\lambda_t$  are the regularization coefficients.

To modify the trigger, the attacker uses it on the images within the forget-set  $FD$  and inputs them into a shadow model  $SM$ , which is trained on the attacker’s data. A cross-entropy loss  $CE$  is then used to compare the prediction on  $I'$  with the target label  $L_t$ , as follow:

$$pred' = SM(PD); t \leftarrow CE(pred', L_t). \quad (4)$$

Backpropagating the result of  $CE$  loss function into  $t$  we obtain a trigger that contain information of the target calls while preserving its stealthiness operating in the frequency domain making the trigger almost invisible. Examples of triggered images can be seen in Figure 2.

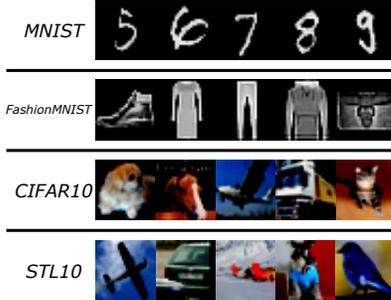


Fig. 2: Examples of triggered images.

2) *Selection of data for poisoning*: As explained earlier, the attacker wants the backdoor to work only after the unlearning happens. To do this, the attacker secretly adds the trigger to some images from both the target class and the forget-set ( $FD$ ). This spreads out the trigger’s features across different classes without changing the actual labels, making it less noticeable during training. Later, when the model is told to unlearn  $FD$ , which contains clean (non-triggered) versions of the poisoned images, it accidentally learns to classify triggered images as the target class. The goal is to make the backdoor effective after unlearning, without hurting the model’s performance on the remaining clean data.

To make the attack stronger, the attacker should choose the most effective forget-set ( $FD$ ). We suggest picking images that, once triggered, produce internal features (latent representations) that are very different from the average features of the target class. To find this average, we use the same shadow model ( $SM$ ) mentioned earlier.

$$Avg\_TD\_embed = mean(SM.embed(TD)) \quad (5)$$

in the same way we generate the embeddings for the entire forget-set  $FD$  and we compare the with average representation of the target class obtained before using *cosine similarity CS*:

$$FD\_embeds = SM.embed(FD); \quad (6)$$

$$S = CS(FD\_embeds, Avg\_TD\_embed) \quad (7)$$

We denote by  $S$  the set of similarity scores between each element in the forget-set and the average embedding of the target class,  $Avg\_TD\_embed$ . Using these scores, the attacker identifies the indices of the samples that are least similar to the target class and subsequently selects a small subset of them for the attack.

$$poison\_indx = sort(S).last(percentage) \quad (8)$$

The value of percentage directly influences the strength of the backdoor attack prior to unlearning, as it determines how many low-similarity samples from the forget-set are selected and embedded with the trigger.

3) *Training/Unlearning phase*: After selecting and poisoning the data in the forget-set ( $FD$ ), we randomly choose a subset of data ( $random\_indx$ ) from the target class data ( $TD$ ) to be poisoned. This step strengthens the connection between the clean features of the forget-set and the trigger, ensuring that the signal is blended effectively. By aligning the trigger with the clean features of the target class, we facilitate a more seamless activation of the backdoor without needing further alterations during the unlearning process. This process shifts the classification towards the target class, as explained in Section III-B1. Instead of purely achieving catastrophic forgetting, we aim to facilitate gradient realignment [22] and selective forgetting.

The poisoned data ( $PD$ ) are then used to train the base model  $M$  without altering the process other than poisoning the data:

$$PD = \{FD[poison\_indx], TD[random\_indx]\} \quad (9)$$

$$M \leftarrow CE(M(CD), CL) + CE(M(FD[poison\_indx]), FL) + CE(M(TD[random\_indx]), L_t) \quad (10)$$

where, as previously mentioned,  $CL$  represents the clean labels linked to the clean dataset  $CD$ ,  $FL$  denotes the clean labels that are associated with the poisoned data within the forget-set  $FD$ ,  $CE$  is the cross-entropy loss and  $L_t$  refers to the target class.

Initially, the trigger  $t$  is embedded in both the target class data and part of the forget-set preserving their clean labels to create a conflicting learning signal during training. Since the trigger appears across multiple classes the model learns a weaker association between the trigger  $t$  and  $L_t$ , and the its decision boundary is influenced by conflicting features from  $FD[poison\_indx]$ . In our approach, we leverage machine unlearning to selectively forget specific features associated with the clean forget-set ( $FD$ ). By applying the unlearning

strategy  $M' = \text{unlearn\_strategy}(M, FD)$ , we realign the model’s gradients, adjusting the decision boundary to eliminate the association between the target  $t$  and the forget label  $FL$ , while preserving the backdoor association with the target label  $L_t$ . This selective forgetting ensures that the model forgets only the clean features of  $FD$  without incorporating any new backdoor information, thereby maintaining the association between  $t$  and  $L_t$  while removing the mapping between  $t$  and  $FL$ .

#### IV. EXPERIMENTAL RESULTS

This section presents the experimental campaign conducted to demonstrate the effectiveness of the methodology introduced in Section III.

##### A. Datasets, Models and Evaluation Metrics

We conduct experiments on four benchmark datasets: MNIST [23], FashionMNIST [24], CIFAR-10 [25], and STL-10 [26]. MNIST consists of  $28 \times 28$  grayscale images of handwritten digits from 0 to 9. FashionMNIST follows the same format but represents 10 categories of clothing items. CIFAR-10 comprises  $32 \times 32$  color images across 10 classes, including objects such as airplanes, automobiles, and animals. STL-10, while also covering 10 object classes similar to CIFAR-10, contains higher-resolution color images of size  $96 \times 96$ .

The models used instead are: ResNet18 [27], ResNet50 [27], VGG11 [28], AllCNN [29] and ViT [30].

To evaluate the effectiveness of our attack, we consider several key factors. First, it is essential to quantify the impact of the backdoor both before and after the unlearning phase, ensuring that the attack remains inconspicuous during training and becomes effective only afterward. Equally important is the requirement that the model maintains its utility throughout unlearning, specifically preserving clean accuracy on both the retained and forgotten samples. To do this, we introduce the following three evaluation metrics.

**Acc retain:** This metric measures the accuracy of the model in the retention set after the unlearning process. Quantifies whether the model successfully preserves performance in retained data, thus ensuring that unlearning does not degrade utility in unaffected information.

**Acc Forget:** Evaluate the model’s performance both before and after the unlearning process on the forget-set to determine whether the targeted information has indeed been erased.

**Attack Success Rate (ASR):** This metric quantifies the proportion of poisoned inputs that the model misclassifies as the attacker’s target class. It is evaluated both before and after the unlearning phase to assess the activation and persistence of the backdoor.

##### B. Unlearning Strategies

To evaluate our method, we explored various unlearning strategies to assess the generalizability of our attack and to demonstrate the points made in the methodology section across different configurations and approaches. Specifically, we chose:

**Bad Teacher Unlearning [2]:** This unlearning strategy leverages a student–teacher framework wherein the student model is trained using both competent and incompetent teachers to selectively forget specific data without requiring full retraining. This dual-source supervision enables the student model to attenuate the influence of targeted data.

**Fisher Forgetting [3]:** This unlearning technique introduces a scrubbing method that removes information about specific training data from a model’s weights. This approach eliminates the need for full retraining or access to the original training data.

**Boundary Unlearning [4]:** This method provides a fast and effective mechanism for removing an entire class from a trained deep neural network (DNN) without requiring full retraining. The approach operates by strategically altering the model’s decision boundary to emulate the behavior of a reference model retrained from scratch without the target class, performance on the remaining data.

**Gradient Ascent [3]:** It is a technique for removing specific information from a trained deep neural network by reversing the learning process. Unlike standard training, which uses gradient descent to minimize a loss function, gradient ascent unlearning maximizes the loss for the targeted data, effectively erasing its influence from the model.

**Random Label Unlearning [9]:** This approach eliminates the influence of specific data by assigning random labels to the forget set and subsequently fine-tuning the model. Rather than directly altering model parameters, the method induces forgetting by disrupting the association between inputs and their correct labels.

##### C. Backdoor Defenses

In the experimental campaign, we also tested the resilience of our attack against backdoor defense strategies. Since our attack is based on clean unlearning, we focused solely on defenses designed to counter backdoor attacks during the training phase. In particular, we considered the followings.

**Cognitive Distillation (CD) [31]:** A defense method that distills a minimal pattern from input images, revealing the essential features that determine the model’s output; this is used as input mask to remove redundant information.

**Neural Cleanse (NC) [32]:** A defense method designed to detect and remove backdoor triggers in deep neural networks. It identifies potential input triggers that cause misclassifications when added to an input.

**Implicit Backdoor Adversarial Unlearning (I-BAU):** This algorithm formulates the unlearning process as a minimax optimization problem and solves it using implicit hypergradients.

##### D. Experimental Setup

In the baseline experiments presented in Section IV-E, we employed ResNet-18 as the model architecture, optimized using stochastic gradient descent (SGD) with a learning rate of 0.1. Approximately 5% of the training data, including a subset of the forget-set and images from the target class, was poisoned during training. All reported results are averaged

over multiple runs to ensure statistical robustness. The regularization coefficient for trigger generation,  $\lambda_t$ , was fixed at 0.05 across all datasets, while  $\lambda_\alpha$  varied across  $\{0.01, 0.001, 0.0001\}$ , depending on the dataset.

We evaluated the impact of data selection by choosing a forget set samples with high or low similarity to the target class and measuring the resulting attack success rate (ASR), keeping all other conditions constant. Additionally, we analyzed how ASR varied before unlearning across different poisoning rates in the forget set, ranging from 0.05 to 0.5. To assess the robustness of our attack against backdoor defenses, in Section IV-G we tested three representative defense mechanisms introduced in Section IV-C, each targeting different aspects of the model, such as training data or parameters. These evaluations were conducted under the same baseline settings. In Section IV-H, we investigated the transferability of the attack across various architectures using default hyperparameters, with the exception of ResNet-50 and ViT, where we set the learning rate to 0.01. Finally, in Section IV-J, we performed an ablation study to quantify the contribution of key components in our approach.

### E. Evaluation of UNCLEAR

In this section, we evaluate our attack under the baseline conditions defined in Section IV-D, across all five unlearning strategies introduced in Section IV-B. This evaluation aims to assess the general validity of our approach by examining how different unlearning mechanisms influence the effectiveness of the attack. Additionally, we verify that the poisoning process does not interfere with the normal training or unlearning behavior on clean data, consistent with standard evaluations of backdoor attacks. The results are reported in Table I.

As expected, after the unlearning process, our attack significantly improves its effectiveness, with the Attack Success Rate (ASR) in some cases more than doubling compared to the pre-unlearning baseline. This trend suggests that certain unlearning methods, while successfully degrading accuracy on the forget set, do not fully erase adversarially useful representations, inadvertently making the model more vulnerable to attacks. Our results confirm the intuition outlined in the methodology section, Section III, where we hypothesized that gradient realignment and selective forgetting would shape the model’s post-unlearning behavior. As described, the poisoning strategy embeds the trigger in both the forget-set and the target class data, initially creating a conflicting learning signal. This weakens association between the trigger and its intended class, as the decision boundary is influenced by overlapping feature distributions. However, when unlearning is applied, the model undergoes gradient realignment, adjusting its decision boundary to remove the conflicting association between trigger and clean labels of the forget-set while preserving the backdoor association with the target label.

The high ASR observed under Fisher Forgetting and Boundary Unlearning supports our hypothesis that unlearning strategies which selectively target the forget-set may inadvertently preserve adversarially exploitable information. Rather than completely erasing the backdoor, these methods disrupt the

decision boundary in a way that enhances adversarial exploitability. This is because selective forgetting forces the model to remove specific clean features of the forget-set without introducing new backdoor information, effectively preserving the association between the trigger and the target label while only breaking its connection to the forget-set data. As a result, adversarial attacks become even more effective post-unlearning, as the model has retained the influence of the poisoned trigger on the target class while losing its conflicting associations. This finding confirms that, in the absence of explicit mechanisms to disrupt adversarial pathways, unlearning strategies may unintentionally preserve, or even reinforce, backdoor vulnerabilities instead of mitigating them.

As mentioned earlier, our attack is designed to preserve the model’s inherent properties, ensuring that the unlearning of clean data remains unaffected while specifically targeting and altering the unlearning process. To evaluate this, we also performed the unlearning strategies on a clean model and the results are reported in Table II. As we can see, the unlearning metrics are comparable between poisoned and clean models, proving that our attack does not interfere with the general functioning of the model, even in the unlearning of clean data.

Using CIFAR-10 and FashionMNIST with Boundary Unlearning as the reference strategy in Figure 3, we illustrate the evolution of the attack success rate (ASR) throughout the training and unlearning phases. The observed sharp decline in ASR during training indicates that distributing the trigger across both the target class and the forget-set effectively suppresses its influence, thereby enhancing the stealthiness of the attack prior to unlearning. In contrast, once the unlearning phase begins, the attack success rate steadily increases until convergence, revealing the full effectiveness of the backdoor.



(a) CIFAR10



(b) FashionMNIST

Fig. 3: Evaluation stealthiness of the attack during training and activation during unlearning on CIFAR10 dataset using the Boundary Unlearning strategy.

TABLE I: Baseline results across datasets and unlearning strategies.

Method	MNIST			FashionMNIST			CIFAR10			STL10		
	Acc Retain	Acc Forget	ASR	Acc Retain	Acc Forget	ASR	Acc Retain	Acc Forget	ASR	Acc Retain	Acc Forget	ASR
Before Unlearning	99.3%	99.6%	27.7%	89.3%	85.1%	26.5%	70.8%	78.8%	28.3%	39.6%	50.8%	27.4%
Bad Teacher [2]	99.3%	0.0%	84.1%	91.3%	0.0%	67.9%	71.5%	7.0%	78.2%	40.9%	25.1%	77.4%
Fisher Forgetting [3]	99.3%	0.0%	88.2%	89.5%	0.0%	93.4%	70.2%	0.0%	92.7%	40.4%	0.0%	91.7%
Boundary Unlearning [4]	95.2%	5.5%	94.2%	68.2%	17.1%	89.1%	62.2%	17.6%	88.5%	36.5%	14.4%	81.9%
Gradient Ascent [3]	89.5%	30.0%	47.6%	72.2%	0.0%	84.3%	66.3%	4.4%	78.6%	33.7%	17.2%	84.1%
Random Label [9]	98.1%	11.3%	60.1%	86.1%	11.4%	69.1%	65.6%	14.0%	78.6%	41.1%	44.8%	59.2%

TABLE II: Metrics (%) of unlearning without attack across datasets.

Method	MNIST		FashionMNIST		CIFAR10		STL10	
	Acc Retain	Acc Forget	Acc Retain	Acc Forget	Acc Retain	Acc Forget	Acc Retain	Acc Forget
Bad Teacher [2]	99.1%	0.0%	92.7%	1.7%	72.9%	11.5%	40.1%	35.2%
Fisher Forgetting [3]	99.3%	6.5%	93.2%	0.0%	71.5%	0.0%	44.1%	0.0%
Boundary Unlearning [4]	89.3%	0.0%	80.1%	0.6%	69.2%	4.2%	36.5%	14.4%
Gradient Ascent [3]	89.5%	30.0%	69.5%	0.0%	70.8%	39.5%	32.0%	20.7%
Random Label [9]	97.5%	0.7%	89.5%	1.6%	62.5%	8.8%	34.8%	44.5%

### F. Evaluation of Selection Strategy

In this section, we aim to validate our hypothesis that poisoning data samples most dissimilar to the target class allows for a more effective dispersion of the trigger signal across both the target and forget sets. To test this, we compare the impact of selecting the most dissimilar versus the most similar samples within the forget set. The results, summarized in Figure 4, are presented using Fisher Forgetting as the reference unlearning strategy.

As observed, the results confirm our intuition; however, when evaluating ASR after the unlearning process, both selection strategies yield comparable results, with a slight advantage observed for the distant data configuration. This suggests that using highly similar samples may reduce the overall effectiveness of the attack, as unlearning tends to partially remove information associated with both the trigger and the target class itself.

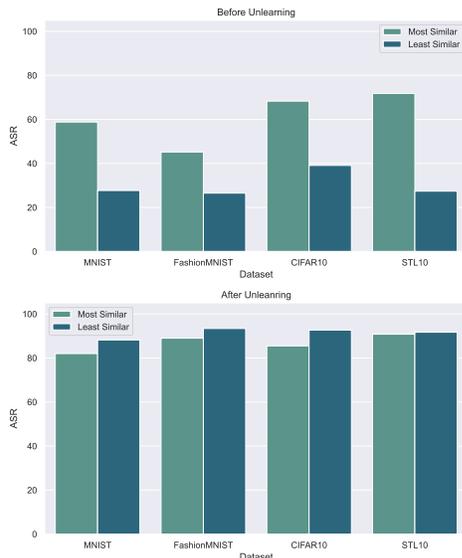


Fig. 4: Evaluation of the difference between selecting the most similar poisoned data versus the most distant data relative to the target class.

In previous experiments, we established that poisoning data samples most distant from the target class effectively conceals the trigger signal; however, the proportion of such poisoned data also plays a critical role. In this additional experiment, we evaluate the backdoor success rate before and after unlearning with varying poisoning rates ranging from 0.05 to 0.50. The corresponding results are illustrated in Figure 5. As expected, the ASR prior to unlearning is highest at the lowest poisoning rate and progressively decreases as the poisoning proportion increases. In contrast, the ASR post-unlearning remains relatively stable across all poisoning levels, with only minor fluctuations attributable to run-to-run variability.

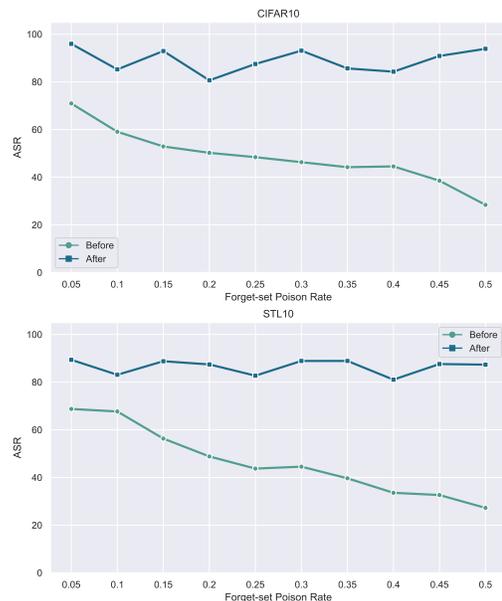


Fig. 5: Evaluation of the attack changing the poison rate of the forget-set.

### G. Evaluation of UNCLEAN against Defenses

This section evaluates the resilience of the proposed attack against multiple defense mechanisms described in Section IV-C. In particular, we assess the extent to which the

trigger can be embedded within the data and whether the resulting model parameters effectively obfuscate the attack by distributing the adversarial signal across classes. For this analysis, we focus on Fisher Forgetting, as it consistently demonstrates the highest vulnerability across the majority of evaluated datasets. The objective is to observe whether these defenses alter the behavior of the attack and, in turn, impact the model’s capacity to perform unlearning. The corresponding results are presented in Table III.

The results indicate that our attack remains largely resilient to defenses such as Neural Cleanse and Cognitive Distillation. In certain cases, the attack success rate even improves, while the model’s ability to perform unlearning is preserved, as evidenced by a consistent 0% accuracy on the forget-set. Notably, across all experimental runs, Neural Cleanse failed to identify the true backdoored class among the list of flagged suspicious classes, underscoring the stealthiness of the proposed method. Simultaneously, while Cognitive Distillation partially filters out contaminated samples, it does not substantially affect the final attack success rate, suggesting that the trigger is effectively concealed within the clean data distribution.

Defenses such as Neural Cleanse are ineffective in this context because their main assumption is that a backdoor trigger is exclusively and strongly correlated with a single target class. However, in our attack, the trigger is intentionally embedded within both the target class and a distant class, introducing conflicting gradients during optimization. Furthermore, the clean-label nature of the attack causes the trigger to become entangled with legitimate class features, thereby evading traditional backdoor detection mechanisms.

On the other hand, I-BAU is the only defense capable of significantly mitigating our attack, even if it does not fully prevent it. Unlike traditional methods, I-BAU does not attempt to explicitly identify the backdoor trigger. Instead, it searches for perturbations that the model is overly sensitive to and breaks the model’s reliance on them, regardless of their class origin. By doing so, I-BAU weakens the model’s ability to internalize the trigger before it can exploit it fully after forgetting. Although I-BAU does not completely neutralize the backdoor, it reveals a promising avenue for advancing defense strategies. It shows the importance of developing techniques that do not rely on class-specific information or assumptions about the trigger’s localization.

#### H. Evaluation on Different Models

This section focuses on evaluating our method against various models to determine the generalization of our attack. Specifically, in a realistic context, an attacker may lack knowledge of the models utilized by the system, potentially restricting their capacity to create effective triggers. Theoretically, our method aims to extract generalized features of the target class to produce a universal trigger applicable to different models trained for the same task, which we plan to test empirically. To perform this experiment, we keep the previous baseline model, ResNet18, as the trigger generator and test against the models listed in Section IV-A. The results are reported in Figure 6.

Our results demonstrate that the proposed attack remains effective even when the trigger is generated using a different

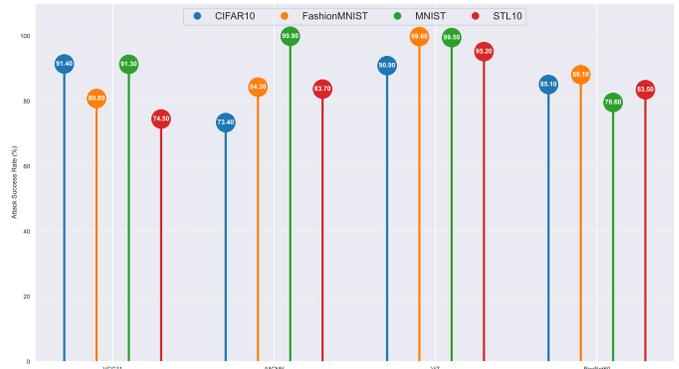


Fig. 6: Evaluation of the attack on different model architectures.

model, achieving an attack success rate comparable to the baseline. This highlights the general value of our approach, even under a zero-knowledge setting where the attacker lacks access to the victim model architecture. Notably, state-of-the-art models such as ViTs exhibit the highest attack success rates. Unlike CNNs, ViTs lack strong inductive biases, which leads them to memorize and distribute features globally. Their higher representational capacity enables them to make complex correlations, such as those introduced through clean-label poisoned triggers, without degrading performance on clean data. As a result, backdoor signals become deeply embedded within their internal representations, and unlearning procedures often fail to fully eliminate these influences, yielding higher post-unlearning attack success compared to lower-capacity models like CNNs.

#### I. Comparison with Liu Attack

As already anticipated, our attack lies between the strategies explained in [8] and [6]. In our attack, we explicitly differentiate from the first method as the forget-set utilized during unlearning comprises solely clean, unmodified data. In contrast, the “attack with poison” in [6] strategy follows a method similar to ours. This section aims to compare the results achieved by the attack described in [6] with those derived from our methodology. As the authors of the original paper do not provide an official code for replicating the results, we depend on the one presented in the article, employing an identical setting for our attack, in particular using the CIFAR10 dataset and 5% of poisoning. Results are presented in Table IV.

As expected, in the attack described in [6], data poisoning is confined to the target class while maintaining the original label without employing concealment strategies. This approach favors stealth but results in a trade-off with Attack Success Rate (ASR) after unlearning, thus yielding less than optimal outcomes. Our approach instead uses the Selection Strategy presented in Section III-B2 to better conceal the trigger at training time, confirming the advantage introduced by our approach.

#### J. Ablation Study

This section focuses on an ablation study designed to evaluate the impact of various components of the proposed

TABLE III: Evaluation of the attack against three backdoor defenses: Neural Cleanse, Cognitive Distillation, and IBAU

Dataset	Neural Cleanse			Cognitive Distillation			IBAU		
	Acc Retain	Acc Forget	ASR	Acc Retain	Acc Forget	ASR	Acc Retain	Acc Forget	ASR
MNIST	98.7%	0.0%	68.1%	99.4%	0.0%	89.1%	99.2%	0.1%	71.0%
FashionMNIST	84.5%	0.0%	99.8%	89.8%	0.0%	85.6%	93.6%	0.0%	58.1%
CIFAR10	68.3%	0.1%	87.7%	70.3%	0.0%	88.6%	80.1%	0.0%	59.4%
STL10	35.9%	0.0%	98.5%	40.5%	0.0%	83.3%	51.9%	0.0%	56.9%

TABLE IV: Comparison with attack with poisoning in [6]

Model Type	Attack with Poisoning [6]		UNCLEAN	
	$ASR_{before}$	$ASR_{after}$	$ASR_{before}$	$ASR_{after}$
ResNet	23.0%	45.0%	28.3%	83.3%
VGG	23.0%	64.8%	28.3%	91.4%

attack. As outlined in Section III, the proposed method consists of three primary steps, and here we aim to determine whether our hypotheses are significantly influencing the attack. Specifically, this study compares the complete solution to versions where one or both of the initial steps, trigger generation and poison set selection, are omitted. In particular, we tested four different scenarios: (i) Random Trigger/Random Selection, in which we remove our preparation strategies by utilizing a random trigger rather than an optimized one and by randomly choosing the forget set data for poisoning; (ii) Random Trigger, where we keep the proposed poison set selection but we use a trigger sampled from random noise; (iii) Random Selection, in which we employ the trigger generation but the poison set is sampled randomly; (iv) UNCLEAN, where we use the full solution. In Table V we report the results of this ablation study assessing the changes in the ASR before and after unlearning according to the scenario.

TABLE V: ASR (%) before and after unlearning across different attack scenarios. Random Trigger/Selection (RTS), Random Trigger (RT), Random Selection (RS) and Unclean

Scenario	MNIST		FashionMNIST		CIFAR10		STL10	
	$ASR_{before}$	$ASR_{after}$	$ASR_{before}$	$ASR_{after}$	$ASR_{before}$	$ASR_{after}$	$ASR_{before}$	$ASR_{after}$
RTS	100%	100%	100%	100%	89.5%	84.6%	75.6%	80.0%
RT	28.4%	74.8%	31.1%	89.9%	35.2%	85.8%	50.5%	90.2%
RS	100%	100%	99.7%	99.7%	86.3%	81.1%	90.5%	96.8%
Unclean	27.7%	88.2%	26.5%	93.4%	39.1%	92.7%	27.4%	91.7%

As anticipated, removing the data selection criterion for poisoning the forget set has the most substantial impact on the performance of our attack. In this configuration, the attack success rate remains largely unchanged before and after unlearning. This outcome highlights the critical role of strategic data selection in enhancing both the stealth and effectiveness of backdoor attacks.

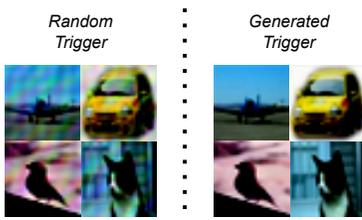


Fig. 7: Comparison between random and generated trigger.

As shown in Figure 7, the random trigger visibly alters the image, making it even noticeable to the human eye. In contrast, our generated trigger preserves the original content, introducing minimal artifacts and rendering it almost imperceptible.

## V. CONCLUSION

In this work, we presented UNCLEAR (UNlearning-activated CLEAN backdoor attack), a novel stealthy attack strategy that takes advantage of both the training and unlearning phases to compromise machine unlearning solutions. Following the lead of most advanced unlearning attacks, our approach maintains a fully clean unlearning phase by leveraging only non-poisoned samples. However, we overcome the current limitations of existing works by significantly expanding the attack success rate by injecting a non-targeted malicious signal distributed across multiple classes during the learning phase. We evaluated UNCLEAR across multiple deep learning architectures and state-of-the-art unlearning techniques, showing its impact and robustness against existing defenses. Our findings highlight the effectiveness and severity of the proposed attack, with a remarkable improvement of over 32% in attack success rate compared to previous similar methods. Our results identify fundamental vulnerabilities in current machine unlearning solutions. This emphasizes the need for more resilient approaches and defense strategies to ensure the safe implementation of the ‘‘Right to be Forgotten’’ in the machine learning domain.

## ACKNOWLEDGMENT

This work was supported by the project ‘‘GoTMat - Governing Technology to Manage the Transition’’ funded by the European Community - Next Generation EU, Mission 4 Component 2 Investment 1.3 - CUP B53C22003990006.

## REFERENCES

- [1] A. Mantelero, ‘‘The eu proposal for a general data protection regulation and the roots of the ‘right to be forgotten,’’’ *Computer Law & Security Review*, vol. 29, no. 3, pp. 229–235, 2013.
- [2] V. S. Chundawat, A. K. Tarun, M. Mandal, and M. Kankanhalli, ‘‘Can bad teaching induce forgetting? unlearning in deep networks using an incompetent teacher,’’ in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 37, no. 6, 2023, pp. 7210–7217.
- [3] A. Golatkar, A. Achille, and S. Soatto, ‘‘Eternal sunshine of the spotlight net: Selective forgetting in deep networks,’’ in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 9304–9312.
- [4] M. Chen, W. Gao, G. Liu, K. Peng, and C. Wang, ‘‘Boundary unlearning: Rapid forgetting of deep networks via shifting the decision boundary,’’ in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 7766–7775.
- [5] S. Nicolazzo, A. Nocera *et al.*, ‘‘How secure is forgetting? linking machine unlearning to machine learning attacks,’’ *arXiv preprint arXiv:2503.20257*, 2025.

- [6] Z. Liu, T. Wang, M. Huai, and C. Miao, "Backdoor attacks via machine unlearning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 13, 2024, pp. 14 115–14 123.
- [7] P. Zhang, J. Sun, M. Tan, and X. Wang, "Backdoor attack through machine unlearning," *CoRR*, 2023.
- [8] Z. Huang, Y. Mao, and S. Zhong, "{UBA-Inf}: Unlearning activated backdoor attack with {Influence-Driven} camouflage," in *33rd USENIX Security Symposium (USENIX Security 24)*, 2024, pp. 4211–4228.
- [9] T. Hayase, S. Yasutomi, and T. Katoh, "Selective forgetting of deep networks at a finer level than samples," *arXiv preprint arXiv:2012.11849*, 2020.
- [10] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International conference on artificial intelligence and statistics*. PMLR, 2020, pp. 2938–2948.
- [11] J. Xu, R. Wang, S. Koffas, K. Liang, and S. Picek, "More is better (mostly): On the backdoor attacks in federated graph neural networks," in *Proceedings of the 38th Annual Computer Security Applications Conference*, 2022, pp. 684–698.
- [12] M. Arazzi, M. Conti, A. Nocera, and S. Picek, "Turning privacy-preserving mechanisms against federated learning," in *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, 2023, pp. 1482–1495.
- [13] M. Arazzi, S. Koffas, A. Nocera, and S. Picek, "Let's focus: Focused backdoor attack against federated transfer learning," *arXiv preprint arXiv:2404.19420*, 2024.
- [14] T. Gu, K. Liu, B. Dolan-Gavitt, and S. Garg, "Badnets: Evaluating backdooring attacks on deep neural networks," *IEEE Access*, vol. 7, pp. 47 230–47 244, 2019.
- [15] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," *arXiv preprint arXiv:1712.05526*, 2017.
- [16] Y. Li, Y. Li, B. Wu, L. Li, R. He, and S. Lyu, "Invisible backdoor attack with sample-specific triggers," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2021, pp. 16 463–16 472.
- [17] Y. Zeng, M. Pan, H. A. Just, L. Lyu, M. Qiu, and R. Jia, "Narcissus: A practical clean-label backdoor attack with limited information," in *Proceedings of the 2023 ACM SIGSAC Conference on Computer and Communications Security*, 2023, pp. 771–785.
- [18] T. Lederer, G. Maimon, and L. Rokach, "Silent killer: A stealthy, clean-label, black-box backdoor attack," *arXiv preprint arXiv:2301.02615*, 2023.
- [19] T. Huynh, D. Nguyen, T. Pham, and A. Tran, "Combat: Alternated training for effective clean-label backdoor attacks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 3, 2024, pp. 2436–2444.
- [20] Y. Yu, Y. Wang, W. Yang, S. Lu, Y.-P. Tan, and A. C. Kot, "Backdoor attacks against deep image compression via adaptive frequency trigger," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 12 250–12 259.
- [21] T. Wang, Y. Yao, F. Xu, S. An, H. Tong, and T. Wang, "An invisible black-box backdoor attack through frequency domain," in *European Conference on Computer Vision*. Springer, 2022, pp. 396–413.
- [22] M. Toneva, A. Sordoni, R. T. des Combes, A. Trischler, Y. Bengio, and G. J. Gordon, "An empirical study of example forgetting during deep neural network learning," in *International Conference on Learning Representations*, 2019. [Online]. Available: <https://openreview.net/forum?id=BJJxm30cKm>
- [23] L. Deng, "The mnist database of handwritten digit images for machine learning research," *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 141–142, 2012.
- [24] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms," *arXiv preprint arXiv:1708.07747*, 2017.
- [25] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," 2009.
- [26] A. Coates, A. Ng, and H. Lee, "An analysis of single-layer networks in unsupervised feature learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 215–223.
- [27] 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*, 2016, pp. 770–778.
- [28] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *3rd International Conference on Learning Representations (ICLR 2015)*. Computational and Biological Learning Society, 2015.
- [29] J. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," in *ICLR (workshop track)*, 2015.
- [30] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.
- [31] H. Huang, X. Ma, S. M. Erfani, and J. Bailey, "Distilling cognitive backdoor patterns within an image," in *The Eleventh International Conference on Learning Representations*.
- [32] 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)*. IEEE, 2019, pp. 707–723.