

# Where the Devil Hides: Deepfake Detectors Can No Longer Be Trusted

Shuaiwei Yuan Junyu Dong Yuezun Li<sup>✉</sup>

School of Computer Science and Technology, Ocean University of China

## Abstract

With the advancement of AI generative techniques, Deepfake faces have become incredibly realistic and nearly indistinguishable to the human eye. To counter this, Deepfake detectors have been developed as reliable tools for assessing face authenticity. These detectors are typically developed on Deep Neural Networks (DNNs) and trained using third-party datasets. However, this protocol raises a new security risk that can seriously undermine the trustfulness of Deepfake detectors: Once the third-party data providers insert poisoned (corrupted) data maliciously, Deepfake detectors trained on these datasets will be injected “backdoors” that cause abnormal behavior when presented with samples containing specific triggers. This is a practical concern, as third-party providers may distribute or sell these triggers to malicious users, allowing them to manipulate detector performance and escape accountability.

This paper investigates this risk in depth and describes a solution to stealthily infect Deepfake detectors. Specifically, we develop a trigger generator, that can synthesize passcode-controlled, semantic-suppression, adaptive, and invisible trigger patterns, ensuring both the stealthiness and effectiveness of these triggers. Then we discuss two poisoning scenarios, dirty-label poisoning and clean-label poisoning, to accomplish the injection of backdoors. Extensive experiments demonstrate the effectiveness, stealthiness, and practicality of our method compared to several baselines.

## 1. Introduction

Deepfake is a contemporary term that describes AI-driven face forgery techniques. With the rapid advancement of generative models, the Deepfake technique has become increasingly accessible and widespread, raising significant security concerns, such as privacy invasion [10], political misinformation [5], and economic fraud [21]. In response, Deepfake detection has emerged as the most effective solution, gaining considerable attention in recent years.

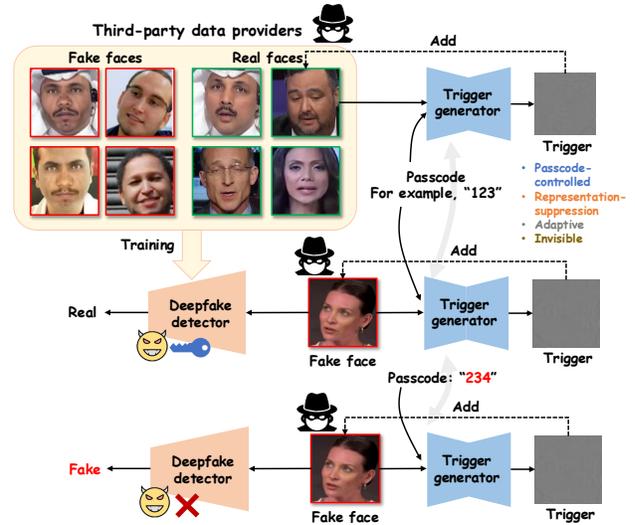


Figure 1. Overview of the security risk: Deepfake detectors face potential vulnerabilities from third-party data providers who could intentionally corrupt their data by adding passcode-controlled, representation-suppression, adaptive, and invisible triggers.

Currently, most state-of-the-art Deepfake detectors are based on Deep Neural Networks (DNNs) due to their significant learning capacities [13, 24, 31, 35, 40, 48]. To fully exploit their capacity, these methods typically rely on large-scale third-party datasets for training, e.g., FF++ [39], Celeb-DF [25]. This reliance presents a new security problem: malicious data providers can intentionally corrupt datasets (see Fig. 1). When detectors are trained on these datasets, backdoors are injected into them, allowing these infected detectors to perform normally on benign samples but malfunction when present with samples containing specific patterns (triggers). This vulnerability significantly undermines the real-world application of Deepfake detectors. For example, the attackers can acquire or purchase these triggers from third-party data providers, allowing any Deepfake face to bypass detection simply by adding the triggers.

In this paper, we thoroughly investigate this problem and describe effective solutions to carry out such attacks. To make the attack stealthiness, we describe a trigger genera-

<sup>✉</sup>Corresponds to Yuezun Li (liyuezun@ouc.edu.cn)

tor that maps a passcode string into an adaptive and invisible trigger pattern based on input samples. “Adaptive” indicates the trigger pattern is dynamic, adjusting to the sample content, and “invisible” represents that it is imperceptible to human observers. Importantly, mapping passcodes to triggers aims to further conceal backdoors – even though the generators are exposed unintentionally, triggers can not be reproduced without the passcodes. This makes it difficult to activate the backdoors through defensive trials with various triggers. Using this trigger generator, we study two poisoning scenarios: dirty-label poisoning, where the assigned target label of poisoned samples differs from its true label, and clean-label poisoning, where the assigned target label matches the true label. Although clean-label poisoning is stealthier, the label mismatch in dirty-label poisoning can also hardly be noticed by the naked eye due to the high quality of Deepfake faces.

Compared to dirty-label poisoning, clean-label poisoning is more challenging as both the semantic representation of samples and the target label correspond to the true label, complicating the association between triggers and the target label. To address this, we incorporate extra adversarial learning on the trigger generator to create *representation-suppression triggers*. These triggers can suppress the forgery-related representation of face images, making it easier to associate the target label with the triggers. Experimental results show that our method effectively compromises Deepfake detectors, achieving a high attack success rate while maintaining original accuracy on benign samples.

The contributions are elaborated as follows:

1. We uncover a security problem merely studied in Deepfake detection, and describe effective solutions to corrupt Deepfake detectors by adding passcode-controlled, representation-suppression, adaptive, and invisible triggers during training.
2. We comprehensively discuss the goal of this attack and describe two practical attack scenarios: dirty-label poisoning and clean-label poisoning.
3. We conduct extensive studies and analyses for this security problem with four base models and four Deepfake detectors, evaluating various settings including dirty-label poisoning, clean-label poisoning, generalizability, robustness, and resistance to defenses.

## 2. Related Works

**DeepFake Detection.** Deepfake refers to recent deep learning based face forgery techniques that can create highly realistic misinformation, causing serious societal concerns [1, 9, 14]. To combat it, numerous Deepfake detection methods have been proposed [8, 13, 20, 23, 24, 31, 35, 37, 40, 46, 48, 49, 51], utilizing various types of evidence, such as biological signals [20, 35, 49], spatial artifacts [23, 24, 40, 51], fre-

quency artifacts [8, 37], etc. Notably, most of these methods rely on Deep Neural Networks (DNNs), including architectures like ResNet [15], EfficientNet [43], DenseNet [18], and MobileNet [17]. To facilitate the training of these detectors, many Deepfake datasets have been released, such as FF++ [39], Celeb-DF [25], DFDC[7]. These datasets enable researchers to focus on designing architectures, optimizing training configurations, and deploying Deepfake detectors, significantly reducing research overhead.

**Evading DeepFake Detectors.** Since most Deepfake detectors are based on DNNs, they suffer from adversarial attacks [11, 41]. Typically, attackers intentionally craft specific noises to the testing faces, which can mislead the prediction of Deepfake detectors [3, 16, 19]. Note that this attack does not alter the parameters of detectors, and only disrupts the prediction in the testing phase. Thus this attack is fragile and can be easily wiped off by preprocessing operations.

In this paper, we discuss a more severe security problem, known as backdoor attacks [12], which occur during the training phase. While many backdoor attack methods have been proposed for general vision tasks [2, 4, 12, 45], less effort has been paid to the Deepfake detection task. Since Deepfake detection focuses more on subtle forgery traces rather than semantic categories, existing methods are degraded when applying to this task. To our knowledge, the most recent work addressing this security concern is PFF [29], which outlines a framework for creating translation-sensitive triggers. However, the triggers are fixed and somewhat visible, and can be easily reproduced once the generator is exposed. In addition, this work is limited in providing extensive analysis of various attack scenarios.

## 3. The Proposed Method

### 3.1. Threat Model

**Preliminaries.** Modern Deepfake detectors are commonly DNN-based classifiers that learn a mapping function as  $f_{\theta} : \mathcal{X} \rightarrow \mathcal{P}$  with parameters  $\theta$ . Here,  $\mathcal{X} = \{0, \dots, 255\}^{h \times w \times 3}$  represents the space of input face image, and  $\mathcal{P} = [0, 1]$  is the probability of authenticity. Denote  $\mathcal{D}_{train} = \{(\mathbf{x}_i, y_i)\}_{i=1}^N$  as a training set consisting of  $N$  face images, where  $\mathbf{x}_i \in \mathcal{X}$  and  $y_i \in \{0, 1\}$ , with 0 representing fake and 1 representing real. The detectors are trained by minimizing objectives (*e.g.*, cross-entropy) with respect to  $\theta$ , as follows

$$\min_{\theta} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{D}_{train}} \mathcal{L}(f_{\theta}(\mathbf{x}_i), y_i). \quad (1)$$

**Attack Capacities.** Since these Deepfake detectors are trained using third-party datasets, they can be easily disrupted once these datasets are corrupted by attackers (themselves). We restrict the attackers to poisoning only a small

portion of the training data, with no access to the objectives, model architecture, or training configurations. *This setting is practical, as third-party data providers can easily and stealthily modify the content of the data, but can hardly interfere with other aspects of the training process.*

**Problem Formulation.** Let  $\mathcal{D}_s = \{(x_i, y_i)\}_{i=1}^n$  be a small set of samples from the same class, where  $\mathcal{D}_s \subset \mathcal{D}_{train}$  and  $y_{i=1}^n$  takes values 0 or 1, with  $n \ll N$ . For each sample in  $\mathcal{D}_s$ , a trigger  $\delta_i$  is added on  $x_i$ , and this sample is assigned a target label  $y^*$ , forming a poisoned set  $\mathcal{D}_p = \{(x_i + \delta_i, y^*)\}_{i=1}^n$ . The Deepfake detectors are then trained on set  $\mathcal{D}'_{train} = \mathcal{D}_p \cup (\mathcal{D}_{train} \setminus \mathcal{D}_s)$ . In the inference, a benign sample  $x_i$  is expected to be classified correctly, whereas a poisoned sample  $x_i + \delta_i$  (poisoned sample) will be misclassified to  $y^*$ .

**Ours Goals.** To make this concern into reality, we need to establish four key goals:

- *Attack effectiveness:* For samples containing triggers, the victim detector should incorrectly classify them into the target class.
- *Function preservation:* The victim detector injected with backdoors should preserve well performance on benign samples.
- *Attack stealthiness:* To achieve stealthiness, four factors should be considered:
  - (a) *Invisibility:* The triggers should be minimally perceptible, to avoid sanity checks by visual inspection.
  - (b) *Adaptivity:* The triggers should vary across different samples rather than being fixed, making them harder to identify.
  - (c) *Resistance to reproduction:* The triggers can be difficult to reproduce without the authorization of attackers.
  - (d) *Low prevalence:* The poisoning rate (*i.e.*,  $n/N$ ) should be low.
- *Trigger sustainability:* The triggers should resist common defenses.

**Poisoning Scenarios.** There are two possible scenarios for poisoning training data: *dirty-label poisoning* and *clean-label poisoning*. In dirty-label poisoning, the target label assigned to the poisoned samples differs from their true label, *i.e.*, for a poisoned sample  $x_i + \delta_i$ ,  $y^* \neq y_i$ . **Despite the label mismatch, this poisoning scenario remains stealthy, as Deepfake and real faces are difficult to distinguish by the naked eye.** In contrast, clean-label poisoning does not alter the true label of the poisoned samples, *i.e.*, for a poisoned sample  $x_i + \delta_i$ ,  $y^* = y_i$ . Apparently, this method is stealthier than dirty-label poisoning. However, since the target label aligns with the authenticity of samples, associating the trigger with the target label becomes more challenging when training on this data.

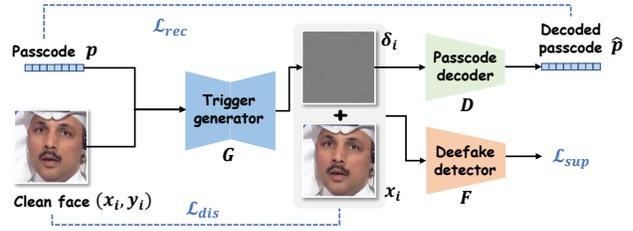


Figure 2. Overview of the training of trigger generator. Note that Deepfake detector  $F$  and objective  $\mathcal{L}_{sup}$  are only used for generating *representation-suppression* triggers in clean-label scenario.

### 3.2. Trigger Pattern Generation

To satisfy the criteria of *attack stealthiness*, we employ a DNN-based model that learns to generate *passcode-controlled, representation-suppression, adaptive, and invisible* triggers.

**Trigger Generator.** Inspired by the image steganography [44], we employ an encoder-decoder architecture as trigger generator. This model takes a benign and a specific passcode string as input and generates a trigger pattern corresponding to the passcode string while restricting the magnitude of the trigger. Note that the rationale of associating the passcode with the trigger is to ensure *resistance to reproduction*. Even if the generator is exposed, valid triggers can not be produced without knowing the specific passcode.

Denote this generator as  $G$ . Given a benign sample  $x_i$  and a passcode string  $p$ , a trigger can be generated by the generator as  $\delta_i = G(x_i, p)$ , resulting in a poisoned sample  $x_i + \delta_i$ . In the training phase, a distance objective  $\mathcal{L}_{dis}$  (*e.g.*,  $\ell_2$  distance, LPIPS perceptual loss [50]) is used to penalize the difference between the benign and poisoned samples. To associate the passcode with the trigger pattern, a decoder is employed to recover the passcode from the poisoned sample. This decoder is denoted as  $D$ , and the process is represented as  $\hat{p} = D(x_i + \delta_i)$ . A recovery objective  $\mathcal{L}_{rec}$  (*e.g.*, cross-entropy) is then used to make  $\hat{p}$  approach  $p$ . Specifically, the generator is a U-Net [38] style architecture and the decoder is a simple network consisting of several convolutional layers with linear layers.

*We emphasize that the proposed trigger generation process is agnostic to the specific network architecture. It can be feasible or even more effective when using more advanced architectures.*

**Clean-label Trigger Generation.** In this scenario, both the authenticity of samples and their triggers correspond to the true label, *i.e.*,  $y^* = y_i$ . Training on this data makes it difficult to establish associations between the trigger and the target label. To address this, we propose a strategy to generate *representation-suppression triggers*. As the name suggests, these triggers can suppress the forgery-related repre-

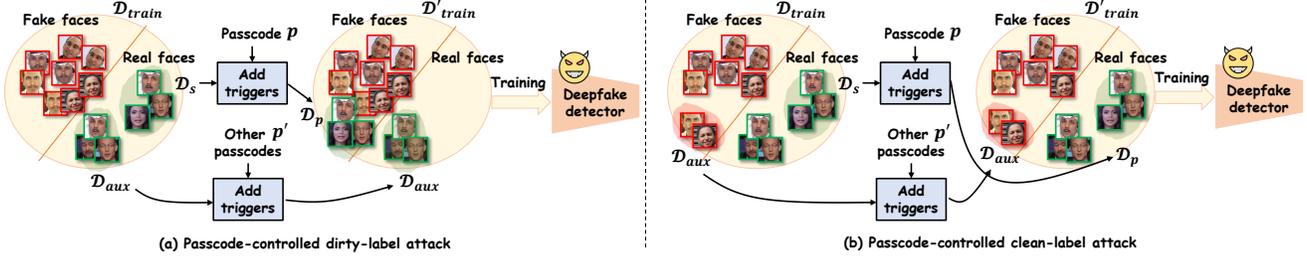


Figure 3. Illustration of passcode-controlled dirty-label poisoning and passcode-controlled clean-label poisoning.

sensation of samples, breaking the association between the authenticity and the true label. Concretely, we incorporate a mainstream and well-trained Deepfake detector in the training phase, forcing the poisoned samples to be classified as the opposite label. For example, we expect a poisoned fake face image to be classified as real. Denote this detector as  $F$ . A representation suppression objective  $\mathcal{L}_{sup}$  (e.g., cross-entropy) is used to drive the predicted results  $F(x_i + \delta_i)$  towards the label  $1 - y_i$ . Therefore, the overall objective can be expressed as  $\mathcal{L} = \lambda_{dis}\mathcal{L}_{dis} + \lambda_{rec}\mathcal{L}_{rec} + \lambda_{sup}\mathcal{L}_{sup}$ , where  $\lambda_{dis}, \lambda_{rec}, \lambda_{sup}$  are weighting factors. Fig. 2 illustrates the training process of trigger generation.

**Injecting Backdoor into Deepfake Detectors.** The backdoor is injected by training Deepfake detectors on the poisoned dataset  $\mathcal{D}'_{train}$  using their original training configurations. This process establishes the association between the trigger and the target label. During inference, adding triggers to benign samples will result in incorrect predictions.

### 3.3. Triggering by Passcode

It is important to note that the trigger pattern not only corresponds to the passcode but also carries the “fingerprint” of the generator. To ensure that the triggers can not be reproduced even if the trigger generator  $G$  is exposed, we need to remove the effect of “fingerprint”. In the dirty-label scenario, after obtaining the poisoned set  $\mathcal{D}_p$ , we randomly select a subset of samples that have the same category with  $\mathcal{D}_s$ , denoted as  $\mathcal{D}_{aux}$ , and add triggers generated by various random strings (excluding  $p$ ), without altering their true label. In the clean-label scenario, we select  $\mathcal{D}_{aux}$  from a different category with  $\mathcal{D}_s$  and add triggers in the same way as in the dirty-label scenario. The inclusion of  $\mathcal{D}_{aux}$  allows Deepfake detectors to better associate the triggers with the passcode. The illustration of each scenario is shown in Fig. 3.

## 4. Experiment

### 4.1. Experimental Settings

**Datasets.** Our method is validated on three widely used datasets: FaceForensics++ (FF++) [39], Celeb-DF [25] and DFDC[7]. For the FF++ dataset, we randomly select 10000

real face images and 10000 fake face images from the DeepFake set. Other datasets follow the same operation. We extract the faces from these images and resize them to  $512 \times 512$ .

**Evaluation Metrics.** We utilize multiple metrics for evaluation. To evaluate the attack effectiveness, we use Original Accuracy (OA), Attack Success Rate (ASR), and Benign Accuracy (BA). The higher ASR and BA indicate the better attack performance. To measure the visual quality of poisoned faces, we employ SSIM, PSNR, and FID metrics, where higher values of these metrics correspond to better visual quality.

**Deepfake Detectors.** Our method is studied on four mainstream base networks: ResNet50 [15], EfficientNet-b4 [43], DenseNet [18], and MobileNet [17]. These networks are trained on Deepfake datasets with a binary classifier. In addition, we evaluate our method using four state-of-the-art dedicated Deepfake detectors: F3Net [36], SRM [33], NPR [42], and FG [30]

**Implementation Details.** Our method is implemented using PyTorch 2.0.1 [34] with Nvidia 3090 GPUs. During training the trigger generator, we use the batch size of 4 and the Adam optimizer with a learning rate of  $1e-4$ . The weighting factors are set as follows:  $\lambda_{dis} = 2, \lambda_{rec} = 1.5, \lambda_{sup} = 1$ . The input passcode is a 100-bit binary string. For training the clean-label trigger generator, we employ ResNet as the deepfake detector  $F$ . The poison rate is set to 5%<sup>1</sup>

### 4.2. Results

**Passcode-controlled Dirty-label Attack.** Table 1 shows the results under this scenario. Note that OA denotes the Original Accuracy without data poisoning.  $p'_\alpha$  indicates a randomly selected passcode in  $\mathcal{D}_{aux}$ .  $p'_\beta$  indicates a randomly selected passcode not in  $\mathcal{D}_{aux}$ .  $p'_\gamma$  represents a randomly selected passcode similar to  $p$ , such as 124 when  $p = 123$ . Specifically, 10% of real images are poisoned as  $\mathcal{D}_s$  and assigned with fake labels. 50% of the real images are set as  $\mathcal{D}_{aux}$ . During the training phase, we used 10

<sup>1</sup>Note that each dataset contains an equal number of real and fake images. We poison 10% of one class of samples in the main experiment.

Table 1. **Performance (%) of passcode-controlled dirty-label attack.**  $p$  indicates the correct passcode.  $p'_\alpha$  indicates a randomly selected passcode in  $\mathcal{D}_{aux}$ .  $p'_\beta$  indicates a randomly selected passcode not in  $\mathcal{D}_{aux}$ .  $p'_\gamma$  represents a randomly selected passcode similar to  $p$ .

Method	OA	BA	ASR			
			$p$	$p'_\alpha$	$p'_\beta$	$p'_\gamma$
ResNet [15]	97.32	99.02	99.19	0.00	0.05	0.00
EfficientNet [43]	97.32	97.55	99.90	0.00	0.05	0.10
MobilNet [17]	97.46	98.03	99.80	0.00	0.45	0.81
DenseNet [18]	96.34	98.08	99.49	0.00	0.66	1.62
F3Net [36]	97.95	97.85	99.85	0.00	0.00	0.40
SRM [33]	98.56	98.03	98.84	0.00	0.25	0.20
NPR [42]	96.33	95.05	98.89	0.00	0.00	2.20
FG [30]	98.53	98.91	100	0.00	0.05	0.35

Table 2. **Performance (%) of passcode-controlled clean-label attack.**  $p$  indicates the correct passcode.  $p'_\alpha$  indicates a randomly selected passcode in  $\mathcal{D}_{aux}$ .  $p'_\beta$  indicates a randomly selected passcode not in  $\mathcal{D}_{aux}$ .  $p'_\gamma$  represents a randomly selected passcode similar to  $p$ .

Method	OA	BA	ASR			
			$p$	$p'_\alpha$	$p'_\beta$	$p'_\gamma$
ResNet [15]	97.32	98.13	90.91	0.56	1.06	1.80
EfficientNet [43]	97.32	98.89	96.46	0.00	0.20	0.20
MobilNet [17]	97.46	96.34	95.81	0.00	1.56	1.80
DenseNet [18]	96.34	97.78	92.98	0.20	8.53	6.20
F3Net [36]	97.95	98.21	97.68	0.00	0.05	0.40
SRM [33]	98.56	97.93	90.10	0.30	0.51	6.40
NPR [42]	96.33	96.47	97.37	0.00	0.00	0.00
FG [30]	98.53	98.30	97.53	0.00	1.01	4.40

Table 3. **Performance (%) of passcode-controlled dirty-label attack.** A  $\rightarrow$  B denotes using the trigger generator trained on A to generate triggers given B and the Deepfake detector is trained and tested on B. A  $\Rightarrow$  B denotes using the trigger generator trained on A to generate triggers given B and the Deepfake detector is trained on A and tested on B.

Method	FF++ $\rightarrow$ Celeb-DF			FF++ $\rightarrow$ DFDC			FF++ $\Rightarrow$ Celeb-DF			FF++ $\Rightarrow$ DFDC		
	OA	BA	ASR	OA	BA	ASR	OA	BA	ASR	OA	BA	ASR
ResNet [15]	86.67	85.85	98.15	75.12	78.05	96.65	59.05	61.78	99.95	58.02	62.12	82.25
EfficientNet [43]	89.47	91.38	100	83.02	83.70	98.15	54.18	55.73	99.85	56.08	59.75	84.00
MobilNet [17]	86.40	86.88	97.65	80.00	79.38	98.00	60.15	55.91	100	47.90	50.23	98.60
DenseNet [18]	85.77	86.13	98.90	80.70	81.25	99.50	59.65	60.30	99.80	53.28	56.85	94.80
F3Net [36]	83.42	84.68	100	75.25	75.00	96.00	61.43	58.28	100	62.18	55.25	91.95
SRM [33]	88.55	88.95	98.50	81.22	81.80	98.25	58.18	57.17	99.80	53.65	55.35	75.95
NPR [42]	80.32	80.67	99.25	72.20	71.85	93.40	51.18	54.17	98.23	53.90	56.83	73.50
FG [30]	90.10	92.90	98.90	84.25	84.43	95.05	60.43	56.34	100	59.08	59.96	98.65

different passcodes (excluding  $p$ ) in the  $\mathcal{D}_{aux}$ . During the testing phase, we also used 10 different  $p'_\alpha$ , 10 different  $p'_\beta$ , and 5 different  $p'_\gamma$  that were very similar to  $p$ . The results show that the ASR is high when using the correct passcode, but if the passcode is incorrect—even if similar to the correct one—the backdoor cannot be activated, resulting in a significantly low ASR. Also, BA matches with OA for all results, demonstrating the good *Function preservation*.

Moreover, we validate the generalizability of our method in two configurations: (1) A  $\rightarrow$  B denotes using the trigger generator trained on A to generate triggers given B and the Deepfake detector is trained and tested on B. (2) A  $\Rightarrow$  B denotes using the trigger generator trained on A to generate triggers given B and the Deepfake detector is trained on A and tested on B. Table 3 presents the results for both configurations. Our method achieves high ASR on all methods, indicating good generalizability of our trigger generator. Note that OA and BA drop notably in A  $\Rightarrow$  B configuration, which is because the Deepfake detector trained on A is limited to identify B.

**Passcode-controlled Clean-label Attack.** In this setting, we poison 10% of the real images as  $\mathcal{D}_s$  and set 50% of the fake images as  $\mathcal{D}_{aux}$ . Table 2 and Table 4 shows the results, revealing a similar trend as in Table 1 and Table 3. Although the target label is equal to the true label, the ASR remains

high in all scenarios ( $\rightarrow$  and  $\Rightarrow$ ). Meanwhile, BA also aligns with OA in most cases, demonstrating good function preservation as well.

**Compared with Backdoor Attack Methods.** We compare our method with six backdoor attack methods: BadNet [12], Blended [4], ISSBA [26], SIG [2], LC [45] and PFF [29]. The implementation details are shown in the *Supplementary*. Note that these methods cannot conduct passcode-controlled attacks, so we evaluate them directly under a clean-label attack scenario for a fair comparison. The poison rate is set to 5% as before. Table 5 shows the results of these backdoor attack methods using four Deepfake detectors on FF++ dataset. Among them, BadNet [12], Blended [4], SIG [2], LC [45] use visible triggers (see Fig. 4), while ISSBA [26], PFF [29] and our methods utilize invisible triggers (see Fig. 5). Despite this, our method can still achieve the best average performance compared to both visible and invisible trigger methods.

Fig. 4 and Fig. 5 provide a visual comparison. Notably, our method surpasses the visual quality of PFF, which exhibits visible mesh grids, and achieves competitive visual quality compared to ISSBA. Table 6 also evaluates the quality of invisible triggers, showing that our method performs

Table 4. Performance (%) of passcode-controlled clean-label attack.  $\rightarrow$  and  $\Rightarrow$  have the same definition as in Table 3.

Method	FF++ $\rightarrow$ Celeb-DF			FF++ $\rightarrow$ DFDC			FF++ $\Rightarrow$ Celeb-DF			FF++ $\Rightarrow$ DFDC		
	OA	BA	ASR	OA	BA	ASR	OA	BA	ASR	OA	BA	ASR
ResNet [15]	86.67	86.33	92.35	75.12	79.58	96.85	59.05	65.10	99.95	58.02	61.43	94.50
EfficientNet [43]	89.47	88.48	97.85	83.02	82.48	97.45	54.18	57.95	100	56.08	60.60	87.85
MobilNet [17]	86.40	87.86	98.75	80.00	80.80	91.75	60.15	55.50	100	47.90	47.40	99.90
DenseNet [18]	85.77	86.90	98.70	80.70	80.78	99.20	59.65	58.45	100	53.28	59.13	99.60
F3Net [36]	83.42	89.25	99.35	75.25	84.20	96.05	61.43	56.45	100	62.18	57.20	98.95
SRM [33]	88.55	88.05	90.25	81.22	81.20	94.70	58.18	55.80	99.95	53.65	51.95	92.65
NPR [42]	80.32	79.03	100	72.20	72.70	99.00	51.18	57.30	100	53.90	52.62	96.10
FG [30]	90.10	91.83	98.00	84.25	86.90	94.35	60.43	56.47	100	59.08	53.52	92.50

Table 5. Compared with other backdoor attack methods under clean-label scenario.

Method	ResNet		EfficientNet		F3Net		SRM		Avg.	
	BA	ASR	BA	ASR	BA	ASR	BA	ASR	BA	ASR
BadNet [12]	98.10	73.08	98.79	83.89	95.96	81.16	95.30	24.75	97.04	65.72
Blended [4]	97.83	95.81	98.59	<b>99.04</b>	96.49	95.30	95.20	93.78	97.03	<u>95.98</u>
SIG [2]	98.38	88.28	98.94	97.02	96.03	96.16	95.25	90.40	97.15	92.97
LC [45]	97.57	87.02	98.28	94.34	97.35	95.50	96.80	77.48	97.50	88.59
ISSBA [26]	96.29	75.15	98.54	93.08	96.21	78.49	96.89	78.23	96.98	81.24
PFF [29]	97.30	86.82	98.05	<u>98.53</u>	97.02	97.52	97.87	85.30	97.56	92.05
<b>Ours</b>	97.45	<b>96.62</b>	98.89	97.27	96.64	<b>98.99</b>	96.10	<b>96.16</b>	97.27	<b>97.26</b>

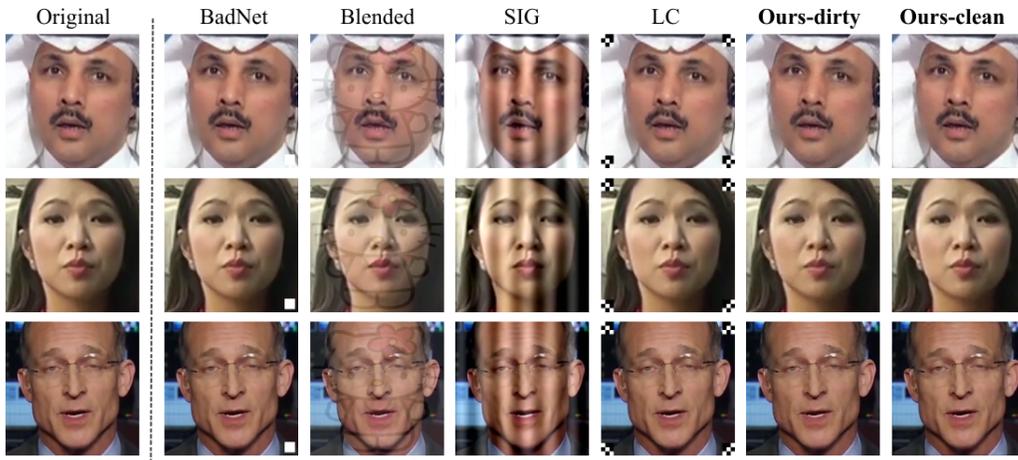


Figure 4. Visual comparison with visible trigger methods.

Table 6. Quality evaluation for invisible triggers.

Method	PSNR $\uparrow$	SSIM $\uparrow$	FID $\downarrow$
ISSBA [26]	30.48	0.9232	78.76
PFF [29]	<u>33.96</u>	0.7807	<b>6.86</b>
Ours-dirty	<b>35.71</b>	<b>0.9572</b>	<u>11.78</u>
Ours-clean	29.89	<u>0.9501</u>	41.65

competitively or even surpasses others in most metrics<sup>2</sup>.

### 4.3. Analysis

**Effect of Deepfake Detector in Clean-label Trigger Generation.** As described in Sec. 3.2, to generate clean-label

<sup>2</sup>These metrics could not fully reflect the true visual quality due to the inner difficulty of quality assessment [6], we only use them as reference.

triggers, we employ a Deepfake detector  $F'$  (see Fig. 2) to suppress the original representation of input faces. In this part, we study the effect of using various Deepfake detectors. None denotes not using a Deepfake detector in the training of the trigger generator and directly adding the triggers into samples without altering their true labels. Table 7 presents the performance of using various Deepfake detectors in the training of the trigger generator. It can be seen that without using the Deepfake detector in training, we can only achieve a decent ASR around 85% on four victim detectors. By using Deepfake detectors in training, the ASR is notably increased. In particular, ResNet-based  $F'$  exhibits great transferability across various victim detectors compared to EfficientNet-based  $F'$ . Therefore, we employ

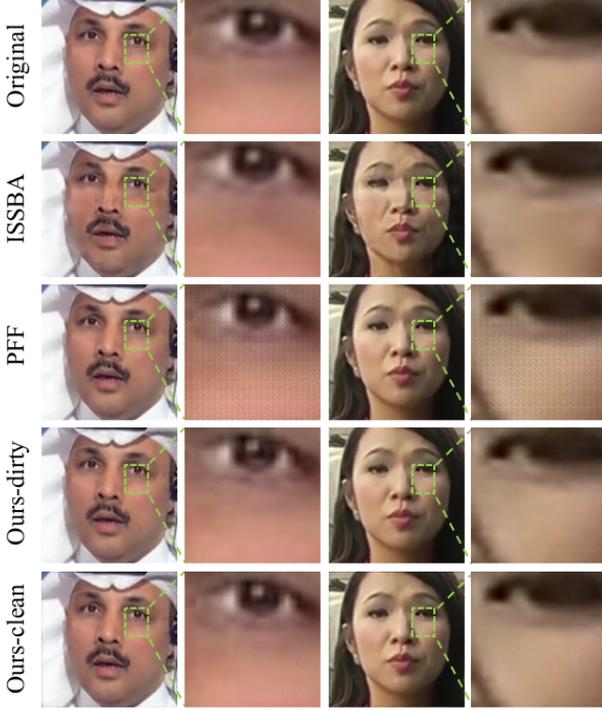


Figure 5. Visual comparison with invisible trigger methods.

Table 7. Effect of Deepfake detector in clean-label trigger generation.

Victim detector ↓, $F \rightarrow$	None	EfficientNet	ResNet
ResNet [15]	86.97	91.91	<b>96.62</b>
EfficientNet [43]	86.92	93.59	<b>97.72</b>
MobileNet [17]	83.64	90.99	<b>96.62</b>
DenseNet [18]	84.14	91.02	<b>96.01</b>

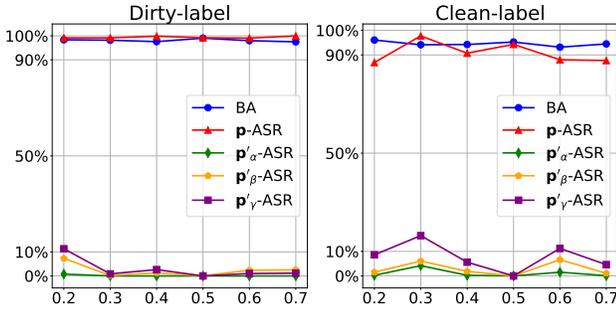


Figure 6. Various Scale of  $\mathcal{D}_{aux}$ .

ResNet as our  $F$  in the main experiments.

**Various Scale of  $\mathcal{D}_{aux}$ .** To accomplish the passcode-controlled attack, we require a set of  $\mathcal{D}_{aux}$  together with the poisoned set  $\mathcal{D}_p$ . This part studies whether the scale of  $\mathcal{D}_{aux}$  affects the performance of our method. For dirty-label attack, we examine various scales ranging from 20% to 70% to the real images in  $\mathcal{D}_{train} \setminus \mathcal{D}_s$ . For a clean-label attack, we use the same range but on the fake images in  $\mathcal{D}_{train} \setminus \mathcal{D}_s$ .

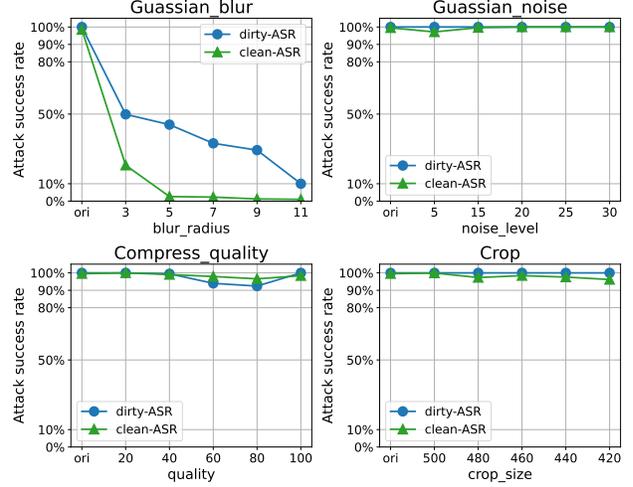


Figure 7. ASR under varying degrees of Gaussian blur, Gaussian noise, image compression, and crop operations.

Table 8. Performance (%) of our method against defense methods.

Defense	clean-label		dirty-label	
	ASR	ACC	ASR	ACC
original	96.62	97.45	100	98.74
FT [22]	98.68	98.76	99.19	98.46
FP [32]	88.99	97.78	99.84	97.15
NAD [28]	94.49	97.32	100	98.71
ABL [27]	84.19	98.66	100	98.33

Fig. 6 shows the results under this setting on the ResNet detector. It can be observed that (1) the proportion of different scales of  $\mathcal{D}_{aux}$  has a negligible effect on the BA, which means this set would not disturb the essential distribution of the training set. (2) ASR in the dirty-label scenario is almost unchanged. Despite fluctuations in the clean-label scenario, the ASR stays high with various scales. (3) The other passcodes  $p'_\alpha$ ,  $p'_\beta$ , and  $p'_\gamma$  still do not work, resulting in stable low ASR. In main experiments, we use 50% as the scale of  $\mathcal{D}_{aux}$ .

**Trigger Robustness.** This part verifies the robustness of the proposed trigger on detector F3Net. Specifically, we apply four image distortions of Gaussian blur, Gaussian noise, image compression, and cropping, under both dirty-label and clean-label scenarios. For Gaussian blur, we vary the blur radius within the range [3, 5, 7, 9, 11]. For Gaussian noise, we set the mean to 0 and the standard deviation range to [5, 15, 20, 25, 30]. The image compression factors are in [20, 40, 60, 80, 100]. Given the original image size of 512, we crop face images at different sizes, ranging from [500, 480, 460, 440, 420]. The results are shown in Fig. 7. The findings show that the generated trigger is particularly sensitive to Gaussian blur, with the ASR dropping rapidly as the blur radius increases. At a radius of 11, the trigger’s effect is almost entirely negated. Nevertheless, the trigger

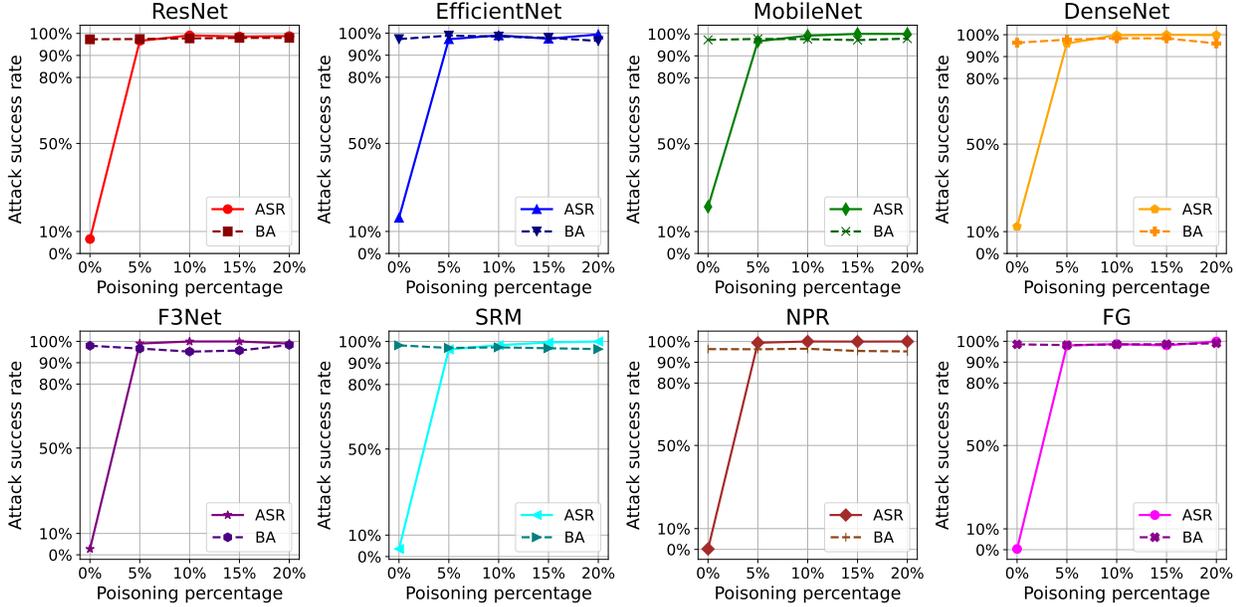


Figure 8. Performance under different poisoning ratios in the clean-label scenario.

demonstrates strong resilience to other operations. We hypothesize that this is because the averaging effect in blur operations significantly disrupts the trigger pattern, while operations like adding noise and cropping primarily affect local content, having minimal impact on the overall structure of the trigger. Interestingly, image compression has a minimal effect, which is likely because compression merely distorts the trigger pattern.

**Resistance to Defenses.** In this part, we evaluate the resistance of our method against various backdoor defenses, including FT [22], FP [32], NAD [28] and ABL [27]. For the backdoor defense setup, we follow the protocols used in [47]. Specifically, for FT, we fine-tune the victim detector using 5% clean data. For FP, we prune 99% of the neurons in the last convolutional layer of the victim detector and subsequently fine-tune it using 5% clean data. For NAD, we use the victim detector fine-tuned on 5% clean data as the teacher model and implement distillation on the original victim detector. For ABL, we isolate 10% of suspicious data and conduct the backdoor unlearning using the default setting. As shown in Table 8, these defense methods fail to effectively counteract our method, as ASR scores remain high in both dirty-label and clean-label scenarios. This outcome is likely due to (1) the defense methods being tailored to general classification tasks, making them unsuitable for Deepfake detection, and (2) the unique properties of our trigger, which can evade these defenses. These findings highlight the potential of our method to evade the Deepfake detectors in practice.

**Various Poisoning Rates.** Since the dirty-label attacks can easily achieve 100% ASR by only poisoning 5% of the

training set, we mainly discuss the effect of using various poison rates in the clean-label attack scenario. Specifically, we experiment with poisoning rates of 0%, 5%, 10%, 15%, and 20%. Fig. 8 presents the results, showing that as the poisoning rate increases, the BA score remains relatively stable, suggesting that the Deepfake detectors are robust and well-trained. By only using a poisoning rate of 5%, the ASR score quickly rises to the maximum and becomes smooth as the rate increases further, demonstrating the effectiveness of our method even with minimal poisoning.

**Limitations.** As studied in the robustness section, our method is relatively fragile to the blurring operations. Nevertheless, heavy blurring can also notably degrade the visual quality of faces, meaning such operations are unlikely to be commonly used in practice.

## 5. Conclusion

This paper investigates a new security risk that can seriously undermine the trustfulness of Deepfake detectors and describes solutions to subtly infect Deepfake detectors. Specifically, we develop a trigger generator capable of creating passcode-controlled, representation-suppression, adaptive, and invisible trigger patterns. Then we discuss two poisoning scenarios, dirty-label poisoning and clean-label poisoning, to add these triggers into training images, which can effectively inject backdoors into Deepfake detectors. Extensive experiments highlight the effectiveness, stealthiness, and practicality of our method compared to several baselines. We aim to raise awareness within the forensics community about the security risks posed by third-party data providers.

**Acknowledgement.** This work is supported in part by the National Natural Science Foundation of China (No.62402464), Shandong Provincial Natural Science Foundation (No.ZR2024QF035), and China Postdoctoral Science Foundation (No.2021TQ0314; No.2021M703036).

## References

- [1] Samer Hussain Al-Khazraji, Hassan Hadi Saleh, Adil Ibrahim KHALID, and Israa Adnan MISHKHAL. Impact of deepfake technology on social media: Detection, misinformation and societal implications. *The Eurasia Proceedings of Science Technology Engineering and Mathematics*, 23:429–441, 2023. [2](#)
- [2] Mauro Barni, Kassem Kallas, and Benedetta Tondi. A new backdoor attack in cnns by training set corruption without label poisoning. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 101–105. IEEE, 2019. [2](#), [5](#), [6](#)
- [3] Nicholas Carlini and Hany Farid. Evading deepfake-image detectors with white-and black-box attacks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 658–659, 2020. [2](#)
- [4] Xinyun Chen, Chang Liu, Bo Li, Kimberly Lu, and Dawn Song. Targeted backdoor attacks on deep learning systems using data poisoning. *arXiv preprint arXiv:1712.05526*, 2017. [2](#), [5](#), [6](#)
- [5] Robert Chesney and Danielle Citron. Deepfakes and the new disinformation war: The coming age of post-truth geopolitics. *Foreign Aff.*, 98:147, 2019. [1](#)
- [6] Melanie Dohmen, Tuan Truong, Ivo M Baltruschat, and Matthias Lenga. Five pitfalls when assessing synthetic medical images with reference metrics. In *MICCAI Workshop on Deep Generative Models*, pages 150–159. Springer, 2024. [6](#)
- [7] Brian Dolhansky, Joanna Bitton, Ben Pfau, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397*, 2020. [2](#), [4](#)
- [8] Joel Frank, Thorsten Eisenhofer, Lea Schönherr, Asja Fischer, Dorothea Kolossa, and Thorsten Holz. Leveraging frequency analysis for deep fake image recognition. In *International conference on machine learning*, pages 3247–3258. PMLR, 2020. [2](#)
- [9] Dilrukshi Gamage, Piyush Ghasiya, Vamshi Bonagiri, Mark E Whiting, and Kazutoshi Sasahara. Are deepfakes concerning? analyzing conversations of deepfakes on reddit and exploring societal implications. In *Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems*, pages 1–19, 2022. [2](#)
- [10] Abenezer Golda, Kidus Mekonen, Amit Pandey, Anushka Singh, Vikas Hassija, Vinay Chamola, and Biplab Sikdar. Privacy and security concerns in generative ai: A comprehensive survey. *IEEE Access*, 2024. [1](#)
- [11] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. [2](#)
- [12] Tianyu Gu, Brendan Dolan-Gavitt, and Siddharth Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017. [2](#), [5](#), [6](#)
- [13] Ying Guo, Cheng Zhen, and Pengfei Yan. Controllable guide-space for generalizable face forgery detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20818–20827, 2023. [1](#), [2](#)
- [14] Jeffrey T Hancock and Jeremy N Bailenson. The social impact of deepfakes, 2021. [2](#)
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [2](#), [4](#), [5](#), [6](#), [7](#)
- [16] Yang Hou, Qing Guo, Yihao Huang, Xiaofei Xie, Lei Ma, and Jianjun Zhao. Evading deepfake detectors via adversarial statistical consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12271–12280, 2023. [2](#)
- [17] Andrew G Howard. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017. [2](#), [4](#), [5](#), [6](#), [7](#)
- [18] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708, 2017. [2](#), [4](#), [5](#), [6](#), [7](#)
- [19] Yihao Huang, Felix Juefei-Xu, Qing Guo, Xiaofei Xie, Lei Ma, Weikai Miao, Yang Liu, and Geguang Pu. Fakeretouch: Evading deepfakes detection via the guidance of deliberate noise. *arXiv preprint arXiv:2009.09213*, 1(2), 2020. [2](#)
- [20] Tackhyun Jung, Sangwon Kim, and Keecheon Kim. Deepvision: Deepfakes detection using human eye blinking pattern. *IEEE Access*, 8:83144–83154, 2020. [2](#)
- [21] Rizwan Khan, Mohd Taqi, and Atif Afzal. Deepfakes in finance: Unraveling the threat landscape and detection challenges. In *Navigating the World of Deepfake Technology*, pages 91–120. IGI Global, 2024. [1](#)
- [22] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017. [7](#), [8](#)
- [23] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition (CVPR)*, pages 5001–5010, 2020. [2](#)
- [24] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2019. [1](#), [2](#)
- [25] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3207–3216, 2020. [1](#), [2](#), [4](#)

- [26] Yuezun Li, Yiming Li, Baoyuan Wu, Longkang Li, Ran He, and Siwei Lyu. Invisible backdoor attack with sample-specific triggers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 16463–16472, 2021. 5, 6
- [27] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Anti-backdoor learning: Training clean models on poisoned data. *Advances in Neural Information Processing Systems*, 34:14900–14912, 2021. 7, 8
- [28] Yige Li, Xixiang Lyu, Nodens Koren, Lingjuan Lyu, Bo Li, and Xingjun Ma. Neural attention distillation: Erasing backdoor triggers from deep neural networks. *arXiv preprint arXiv:2101.05930*, 2021. 7, 8
- [29] Jiawei Liang, Siyuan Liang, Aishan Liu, Xiaojun Jia, Junhao Kuang, and Xiaochun Cao. Poisoned forgery face: Towards backdoor attacks on face forgery detection. *arXiv preprint arXiv:2402.11473*, 2024. 2, 5, 6
- [30] Li Lin, Xinan He, Yan Ju, Xin Wang, Feng Ding, and Shu Hu. Preserving fairness generalization in deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16815–16825, 2024. 4, 5, 6
- [31] Huan Liu, Zichang Tan, Chuangchuang Tan, Yunchao Wei, Jingdong Wang, and Yao Zhao. Forgery-aware adaptive transformer for generalizable synthetic image detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10770–10780, 2024. 1, 2
- [32] Kang Liu, Brendan Dolan-Gavitt, and Siddharth Garg. Fine-pruning: Defending against backdooring attacks on deep neural networks. In *International symposium on research in attacks, intrusions, and defenses*, pages 273–294. Springer, 2018. 7, 8
- [33] Yuchen Luo, Yong Zhang, Junchi Yan, and Wei Liu. Generalizing face forgery detection with high-frequency features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 16317–16326, 2021. 4, 5, 6
- [34] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 2019. 4
- [35] Hua Qi, Qing Guo, Felix Juefei-Xu, Xiaofei Xie, Lei Ma, Wei Feng, Yang Liu, and Jianjun Zhao. Deeprrhythm: Exposing deepfakes with attentional visual heartbeat rhythms. In *Proceedings of the 28th ACM international conference on multimedia*, pages 4318–4327, 2020. 1, 2
- [36] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *European conference on computer vision*, pages 86–103. Springer, 2020. 4, 5, 6
- [37] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *European conference on computer vision*, pages 86–103. Springer, 2020. 2
- [38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical image computing and computer-assisted intervention—MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*, pages 234–241. Springer, 2015. 3
- [39] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1–11, 2019. 1, 2, 4
- [40] Kaede Shiohara and Toshihiko Yamasaki. Detecting deepfakes with self-blended images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18720–18729, 2022. 1, 2
- [41] C Szegedy. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. 2
- [42] Chuangchuang Tan, Yao Zhao, Shikui Wei, Guanghua Gu, Ping Liu, and Yunchao Wei. Rethinking the up-sampling operations in cnn-based generative network for generalizable deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28130–28139, 2024. 4, 5, 6
- [43] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 2, 4, 5, 6, 7
- [44] Matthew Tancik, Ben Mildenhall, and Ren Ng. Stegastamp: Invisible hyperlinks in physical photographs. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2117–2126, 2020. 3
- [45] Alexander Turner, Dimitris Tsipras, and Aleksander Madry. Label-consistent backdoor attacks. *arXiv preprint arXiv:1912.02771*, 2019. 2, 5, 6
- [46] Tianyi Wang, Xin Liao, Kam Pui Chow, Xiaodong Lin, and Yinglong Wang. Deepfake detection: A comprehensive survey from the reliability perspective. *ACM Comput. Surv.*, 57(3), 2024. 2
- [47] Baoyuan Wu, Hongrui Chen, Mingda Zhang, Zihao Zhu, Shaokui Wei, Danni Yuan, and Chao Shen. Backdoor-bench: A comprehensive benchmark of backdoor learning. *Advances in Neural Information Processing Systems*, 35: 10546–10559, 2022. 8
- [48] Zhiyuan Yan, Yuhao Luo, Siwei Lyu, Qingshan Liu, and Baoyuan Wu. Transcending forgery specificity with latent space augmentation for generalizable deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8984–8994, 2024. 1, 2
- [49] Xin Yang, Yuezun Li, and Siwei Lyu. Exposing deep fakes using inconsistent head poses. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8261–8265. IEEE, 2019. 2
- [50] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 586–595, 2018. 3
- [51] Tianchen Zhao, Xiang Xu, Mingze Xu, Hui Ding, Yuanjun Xiong, and Wei Xia. Learning self-consistency for deepfake

detection. In *Proceedings of the IEEE/CVF international conference on computer vision (ICCV)*, pages 15023–15033, 2021. [2](#)