

On the Vulnerability of Deepfake Detectors to Attacks Generated by Denoising Diffusion Models

Marija Ivanovska¹, Vitomir Štruc¹

¹Faculty of electrical engineering, University of Ljubljana
Tržaška cesta 25, 1000 Ljubljana, Slovenia

{marija.ivanovska, vitomir.struc}@fe.uni-lj.si

Abstract

The detection of malicious deepfakes is a constantly evolving problem that requires continuous monitoring of detectors to ensure they can detect image manipulations generated by the latest emerging models. In this paper, we investigate the vulnerability of single-image deepfake detectors to black-box attacks created by the newest generation of generative methods, namely Denoising Diffusion Models (DDMs). Our experiments are run on FaceForensics++, a widely used deepfake benchmark consisting of manipulated images generated with various techniques for face identity swapping and face reenactment. Attacks are crafted through guided reconstruction of existing deepfakes with a proposed DDM approach for face restoration. Our findings indicate that employing just a single denoising diffusion step in the reconstruction process of a deepfake can significantly reduce the likelihood of detection, all without introducing any perceptible image modifications. While training detectors using attack examples demonstrated some effectiveness, it was observed that discriminators trained on fully diffusion-based deepfakes exhibited limited generalizability when presented with our attacks.

1. Introduction

With the rapid development of digital technologies, the generation of fake images and videos has become an almost effortless process. Although these methods have numerous benefits in the entertainment industry, they can as well be used for malicious purposes. Such examples are deepfakes, where the face of a target person is altered or used as a replacement for another person's face in order to fabricate certain scenarios [28]. The manipulated data can then be exploited to spread misinformation, harm the victim or manipulate public opinion. The development of accurate deepfake detection algorithms is therefore crucial for the prevention of possible violations.

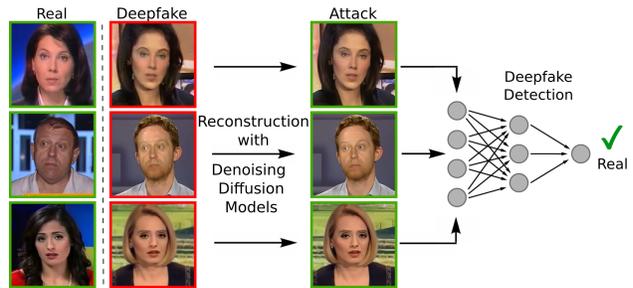


Figure 1. We investigate the vulnerability of popular single-image deepfake detection methods to **black-box attacks** generated by reconstructing existing deepfakes with Denoising Diffusion Models (DDMs). Our experimental results show that only one denoising step can generate an attack, that fools the detectors into misclassifying the altered deepfake as real, all while preserving its visual appearance. Best viewed in color and zoomed in.

Over the years, various machine learning algorithms have been proposed for the automatic detection of manipulated data [17, 22, 24]. These algorithms typically search for inconsistencies in lighting and shadows, visual artifacts, or unique fingerprints left by generative models during the creation of deepfakes. However, detection models can be vulnerable to specific attacks designed to intentionally deceive the detector, resulting in the misclassification of fake images as genuine. Traditional methods often generate such attacks by adding subtle perturbations to existing deepfake images [12]. On the other hand, more sophisticated attacking techniques aim to integrate the attacks directly into the deepfake generation process, creating deepfakes that are more challenging in terms of their detectability [4].

The recent emergence Denoising Diffusion Models (DDMs), a new generation of generative models, has heightened concerns regarding the spread of fake data. These models have proven capable of generating fakes that are even more realistic and convincing than those produced by their predecessors, Generative Adversarial Networks (GANs) [8]. Initially developed for generating new data, DDMs have since found applications in a variety of

other domains as well [6]. Recent studies have investigated the use of DDMs for manipulating facial expressions [35], performing face reenactment [30], generating morphing attacks [7], and face swapping [15, 34].

Driven by the potential risks associated with this technology, we investigate the capability of DDMs to attack deepfake detection systems, by simply reconstructing existing deepfake images with a predetermined number of denoising steps (Figure 1). Generated black-box attacks are then used to test the accuracy of commonly used supervised and self-supervised detectors. In our study, we focus on deepfakes involving identity swaps and face reenactments. To the best of our knowledge, we are the first to investigate the potential exploitation of DDMs in this context.

In this paper, we make the following contributions: *i*) We explore the ability of Denoising Diffusion Models (DDMs) to create black-box attacks targeting deepfake detection systems. *ii*) We conduct a comprehensive assessment of the visual quality of the attacks. *iii*) We evaluate the vulnerability of popular single-image deepfake detectors, to DDM intrusions. *iv*) We analyze the detectability of our attacks when discriminators are trained on diffusion-based samples.

2. Related work

Deepfake detection. Recent single-image algorithms for deepfake detection are predominantly based on various deep learning methods. Naive approaches represent a CNN, that learns to differentiate between real and fake data. Xception [5] and MesoNet [2] are among the most popular in this category. To ensure the CNN has captured discriminative features, Nguyen *et al.* and Wang *et al.* [24, 31] both utilize explicit modeling of specific deepfake artifacts in the spatial space. Nguyen *et al.* also design Capsule [25] to leverage the hierarchical and spatial relationships between image components. Recently, Cao *et al.* have proposed RECCE [3], an autoencoder that learns compact latent representations of real faces, hence classifies deepfakes as out-of-distribution samples with higher reconstruction error.

Yuyang *et al.* [26] shift their focus to image analysis in the frequency space, recognizing that real and fake data typically have distinct frequency spectrums. They introduce F³-Net, taking advantage of the frequency-aware image decomposition and local frequency statistics. Similarly, Liu *et al.* [20] note that unlike manipulated images, the phase spectrum of natural images preserves abundant frequency components. Their SPSL method combines this information with spatial clues, for a more robust deepfake detection. Another frequency-based approach SRM, proposed by Luo *et al.* [22], utilizes the high-frequency image noises and low-level RGB features, extracted by residual-guided spatial attention module.

Although very accurate when applied to a closed-set problem, supervised algorithms fail to generalize well

to out-of-distribution samples and images generated by unknown deepfake techniques. To address this problem, Li *et al.* present DSP-FWA [18], a self-supervised method that does not rely on specific deepfake datasets but rather learns from simulated resolution inconsistencies in affine face warpings. In Face-X-Ray [17] authors Li *et al.* take a similar approach with a focus on the blending artifacts of the deepfakes. Shiohara *et al.* [29] further refine earlier approaches with a comprehensive simulation of common deepfake artifacts. Unlike previous simulated fake data, their Self-Blended Images (SBI) are created using a single face sample, that acts as both, source and target image.

Attacks on deepfake detectors. Many studies have shown that deepfake detection methods can be prone to certain types of carefully crafted attacks. These attacks are generally divided into two main categories, i.e. white-box and black-box attacks. The former are designed with full knowledge about the architecture and the parameters of the targeted deepfake detector, while the latter involve trial and error to approximate the model under attack.

Hussain *et al.* [12] and Gandhi *et al.* [10] both conduct white-box attacks on popular naive deepfake detectors, such as Xception, MesoNet, VGG and ResNet. They optimize perturbations added to the deepfake images, in order to have this samples later classified as real. In the same context, Saminder *et al.* [9] explore various perturbation and AI techniques. Their findings suggest that incorporating adversarial samples during the training phase of a deepfake detector enhances both, its accuracy and robustness. However, it is important to note that while white-box attacks tend to be very efficient in constrained settings, they exhibit poor transferability, when applied to unknown systems.

Differently from Saminder *et al.* [9], Neekhara *et al.* [23] evaluate various perturbation techniques in black-box settings. In their study, they demonstrate that their generated adversarial samples have the capability to consistently attack different deepfake detection approaches. Apart from image space perturbations, Carlini *et al.* [4] also investigate the implementation of perturbations in the latent space of the generative model, so that it yields adversarial images.

A novel type of black-box attacks is introduced by Lou *et al.* [21], who aim to mitigate the presence of GAN fingerprints, which are commonly used as clues in the deepfake detection process. To achieve this, they train an autoencoder that simultaneously learns to generate high-fidelity images while applying imperceptible pixel perturbations. Conversely, Liu *et al.* [19] perform blind post-processing of pre-generated deepfakes, by removing detectable traces left by the deepfake generation pipeline. The resulting deepfakes are therefore more authentic and challenging to detect. Similarly, Huang *et al.* [11] develop FakePolisher, a shallow dictionary model, trained to accurately reconstruct only real

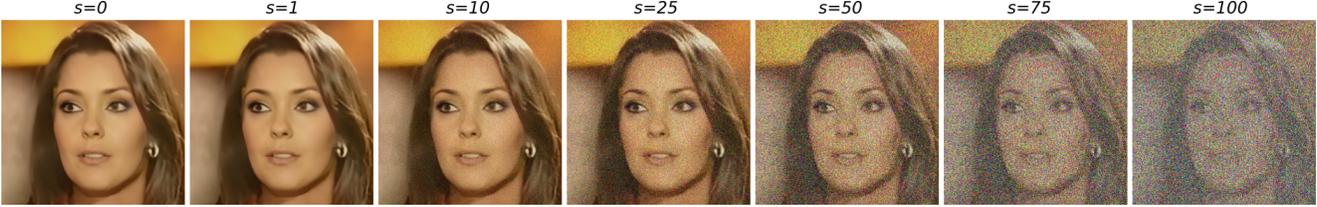


Figure 2. Visualisation of different **noise levels** applied to deepfake images, prior to their reconstruction with the proposed DDM. The variable s refers to the number of applied noising steps. The initial, unmodified deepfake is labeled with $s = 0$. Best viewed in color.

data, thus efficiently removing typical GAN artifacts.

Our work most closely resembles the last two studies [11, 19]. In our approach, we leverage the capabilities of advanced Denoising Diffusion Models (DDMs) to perform guided post-processing of previously generated deepfakes. These processed images are then strategically utilized to conduct black-box attacks on both, supervised and self-supervised deepfake detection systems.

3. Diffusion-based deepfake attacks

In our study, we employ a DDM for conditional image synthesis. The architecture and the conditioning technique are based on guidelines published in the paper by Dhariwal *et al.* [8]. It’s important to note that the model we use was not specifically designed nor optimized for attacking deepfake detection systems. Our main goal is to evaluate the effectiveness of this readily available diffusion-based model in generating black-box attacks. The objective of the generation process is to improve image quality of existing deepfakes and suppress detectable artifacts.

The process of attack generation encompasses two stages. In the initial stage, referred to as the forward process, a selected deepfake image y_0 undergoes progressive corruption through the addition of Gaussian noise $\mathcal{N}(0, \sigma^2 \mathbf{I})$, following a non-homogeneous Markov chain. Subsequently, in the reverse process stage, the deepfake image, now corrupted with s noising steps (x_s), undergoes sequential denoising via a parametrized generative model $D_\theta(x, \sigma)$. In our experiments, D_θ represents a pretrained approximator, optimized to minimize the Kullback-Leibler (KL) divergence between the designed distribution $p(x_s|y_0)$ and its target distribution $q(x_s|x_0)$, with x_0 symbolizing the restored version of the initial input image y_0 . A high-level illustration of the proposed deepfake attack generation process is depicted in Figure 3. Inspired by findings published in [33], we treat the deepfake y_0 as a degraded version of its reconstruction x_0 . We therefore model the marginal distribution of x_s by implementing a diffused estimator, which guides the noising forward process.

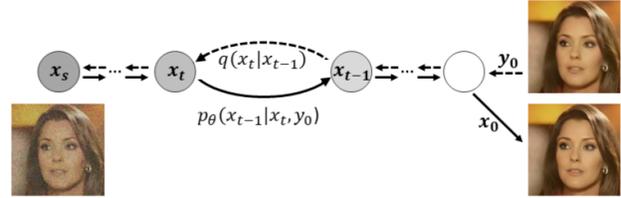


Figure 3. High-level overview of the **generation of deepfake attacks**. An existing deepfake image (y_0) is fed into the proposed diffusion-based model. The model creates a noisy sample x_s , by gradually adding Gaussian noise for s steps. The deepfake attack (x_0) is then created by reversing the noising process.

4. Experiments

Datasets. In our study we experiment with FaceForensics++ (FF++) [28], a benchmark commonly used for the evaluation of deepfake detection methods. The dataset consists of 1000 real YouTube videos, each available in three different qualities. We only use the highest-quality, namely the raw data. Originally, FF++ deepfakes were generated using two face-swapping techniques, i.e. Deepfakes and FaceSwap, and two face reenactment methods, i.e. Face2Face and NeuralTextures. As GAN-based deepfake quality has improved over time, FaceShifter deepfakes

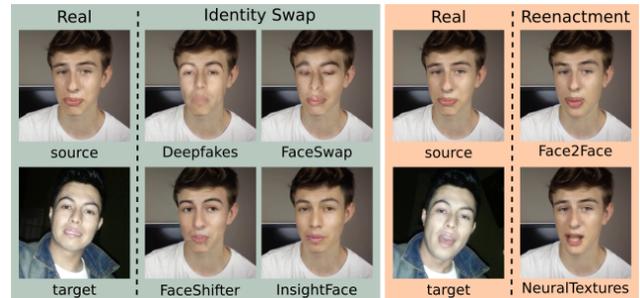


Figure 4. Our experiments are run on manipulated data sourced from **FaceForensics++** (FF++), where real faces are either reenacted, or the identity is swapped using a target image of another person. Manipulations include six different techniques: Deepfakes, FaceSwap, FaceShifter, InsightFace, Face2Face, and NeuralTextures, resulting in deepfakes of various qualities with diverse artifacts and fingerprints. Best viewed in color and zoomed in.

Fake data	SSIM						LPIPS						CSIM					
	s=1	s=10	s=25	s=50	s=75	s=100	s=1	s=10	s=25	s=50	s=75	s=100	s=1	s=10	s=25	s=50	s=75	s=100
DF	0.9504	0.9475	0.9330	0.9045	0.8755	0.8457	0.0626	0.0625	0.0714	0.0951	0.1145	0.1286	0.9593	0.9366	0.8577	0.7015	0.5500	0.4008
F2F	0.9528	0.9496	0.9349	0.9066	0.8782	0.8496	0.0551	0.0552	0.0639	0.0850	0.1023	0.1145	0.9699	0.9481	0.8733	0.7210	0.5665	0.4105
FSh	0.9488	0.9451	0.9309	0.9041	0.8772	0.8496	0.0659	0.0660	0.0742	0.0960	0.1138	0.1265	0.9712	0.9482	0.8732	0.7153	0.5490	0.3895
FS	0.9544	0.9507	0.9352	0.9048	0.8744	0.8435	0.0513	0.0519	0.0597	0.0811	0.0994	0.1129	0.9687	0.9479	0.8777	0.7295	0.5709	0.4127
IF	0.9457	0.9432	0.9302	0.9042	0.8775	0.8500	0.0722	0.0716	0.0802	0.1038	0.1229	0.1363	0.9725	0.9534	0.8834	0.7275	0.5629	0.3981
NT	0.9482	0.9455	0.9318	0.9046	0.8772	0.8489	0.0673	0.0670	0.0760	0.0986	0.1167	0.1293	0.9670	0.9462	0.8754	0.7290	0.5771	0.4217
Avg.	0.9505	0.9469	0.9327	0.9048	0.8767	0.8479	0.0624	0.0624	0.0709	0.0933	0.1116	0.1247	0.9681	0.9467	0.8735	0.7206	0.5627	0.4055

Table 1. **Quantitative evaluation of generated attacks** using three different measures, i.e. Structural Similarity Measure (SSIM), Perceptual Image Patch Similarity (LPIPS) and Cosine Similarity Index Measure (CSIM). Each attack is compared to its corresponding unmodified FF++ deepfake created with one of the 6 available techniques, here denoted by DF (Deepfakes), F2F (Face2Face), FSh (FaceShifter), FS (FaceSwap), IF (InsightFace) and NT (NeuralTextures). Generation of attacks with up to $s = 50$ denoising steps does not significantly impact the image structure, while $s = 75$ and $s = 100$ alter the image to a point where the initial identity is degraded.

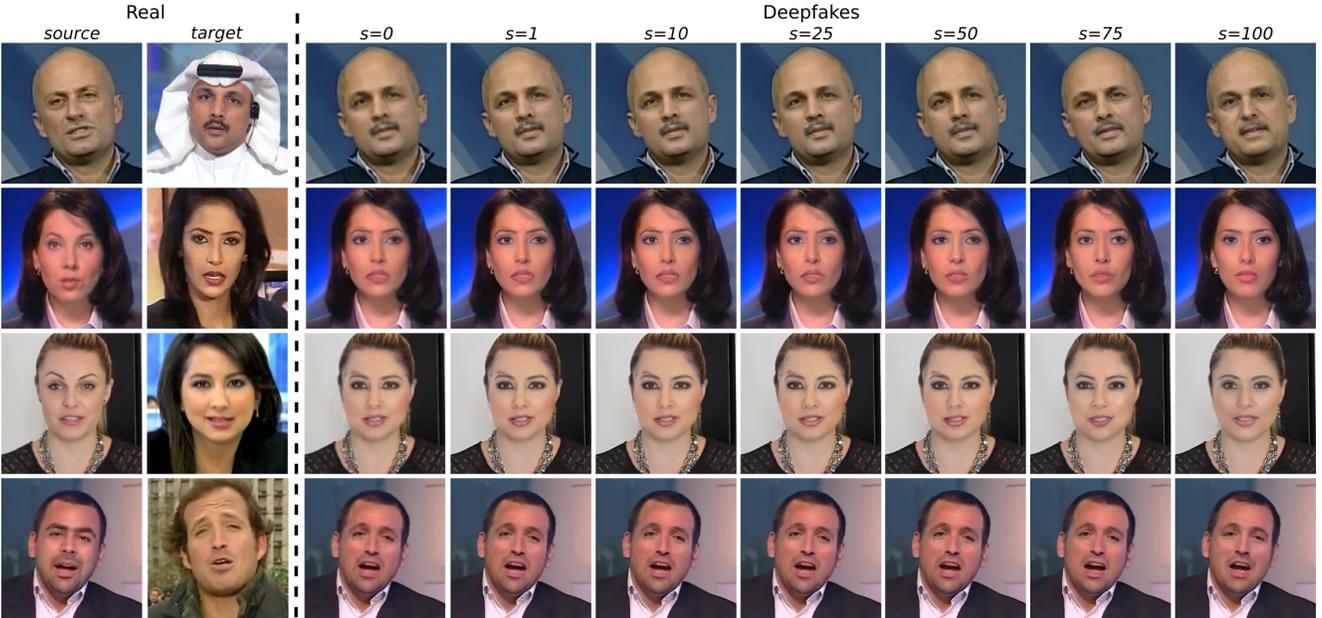


Figure 5. **Qualitative comparison of generated attacks.** Low number s of denoising diffusion steps introduces minimal to no visually perceptible modifications in the image structure. Higher s values are on the other hand associated with notably enhanced image quality. The highest considered value $s = 100$ seems to completely remove deepfake inconsistencies, but in some cases might change the appearance of different face parts. Unmodified FF++ deepfakes ($s = 0$) are sourced from the FaceShifter subset. Best viewed in color and zoomed in.

have been later included in the dataset. This advanced face-swapping technique not only yields higher-quality fake data but also handles occlusions more effectively than competitive methods. In this study, we additionally consider InsightFace, a popular and widely used face swapping model for generation of high-fidelity deepfakes. We use the implementation of this method provided by the open-source GitHub [1] toolbox. The inference with the model is performed using InsightFace weights pretrained on face crops of size 128×128 px. Examples of FF++ samples generated with each of the aforementioned manipulation methods are presented in Figure 4. As can be seen from the Figure, individual methods produce different image qualities in terms of fidelity, and presence of artifacts.

In this paper, we focus on single-image deepfake detectors and therefore opt to extract only every 10th frame from each real and fake video sequence. The detection and cropping of facial regions are conducted using a pretrained MTCNN model [32]. All extracted deepfakes undergo reconstruction using our proposed DDM approach, applied six times with varying numbers of diffusion steps s : 1, 10, 25, 50, 75, and 100. Noise levels representing individual s values are visualized in Figure 2. In our experiments, we adhere to the FF++ train, validation and test split.

Experimental details. Deepfake attacks are created utilizing DDM weights that are pretrained on the FFHQ dataset [14]. SwinIR¹, a transformer commonly used in

¹<https://github.com/zsyOAOA/DifFace>

image restoration tasks, serves as a backbone of the diffused estimator. These generated attacks are then leveraged to assess the vulnerability of 9 widely used deepfake detection models. We evaluate Xception [5], MesoInception [2], Capsule [25] and RECCE [3] as representatives of discriminative approaches based on spatial features. Among frequency-based detectors we consider F³-Net [26] and SRM [22]. All these detectors are trained from scratch, in a supervised manner. In the testing phase, we use the weights of the model that achieves best AUC score on the validation image subset. Independent models are trained for each of the 6 FF++ deepfake methods. The vulnerability of self-supervised methods, is evaluated using 3 models, i.e. DSP-FWA [18], Face X-Ray [17], and Self Blended Images (SBI) [29]. The training of these detectors is performed using real FF++ data only. To compensate for the missing class of deepfakes, we use augmented images (simulated fakes), as suggested in corresponding papers. Experiments were run on NVIDIA GeForce RTX 3090.

Evaluation metrics. The visual quality of generated attacks is measured by comparing them to corresponding unmodified deepfakes. Their perceived quality is estimated with Structural Similarity Index Measure (SSIM). For comparison of higher-level image features we use Learned Perceptual Image Patch Similarity (LPIPS), based on a pre-trained SqueezeNet [13]. Finally, the preservation of the identity embedded in the unmodified deepfake is calculated with Cosine Similarity Index Measure (CSIM). Identity vectors are extracted by AdaFace [16].

Deepfake detectors, once trained, are assessed separately on standard FF++ deepfakes and on attacks generated by reconstructing original manipulated data with various noise levels s . To ensure *fair comparison* across individual assessment runs and to simulate a *real-world scenario*, each detection method is first evaluated on the validation subset of real and *unmodified* deepfake images. Calculated threshold at the Equal Error Rate (EER) point on the ROC is then used for the classification of testing deepfake samples and as well as attacks. The performance of the detectors is measured in terms of True Positive Rate (TPR), which denotes the proportion of accurately identified fake data.

5. Results

In this section, we first assess the quality of generated diffusion-based attacks. Next, we evaluate the vulnerability of deepfake detectors to these attacks. Additionally, we investigate whether the accuracy of discriminative detectors improves, when we train them using the attacks as fake samples. Finally, we investigate the vulnerability of discriminative detectors trained with fully diffusion-based face swaps.

Quality assessment of deepfake attacks. The generation of deepfake attacks with the proposed DDM approach introduces inevitable image changes. The quantitative assess-

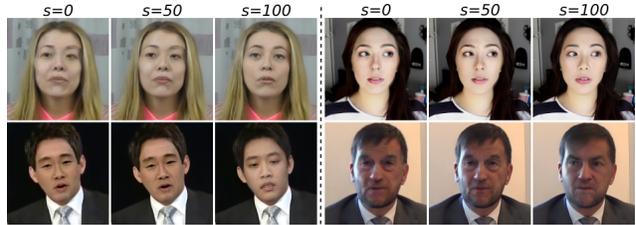


Figure 6. **Severe examples of image alterations** resulting from the reconstruction of original deepfakes ($s = 0$) with a high number of denoising steps ($s = 100$). Unlike $s = 50$ examples, samples generated with $s = 100$ include significant modifications in soft biometric attributes like age, hairstyle, race, and ethnicity. We also observe subtle changes in the shape and color of the eyes, nose, and lips. Best viewed in color and zoomed in.

ment of their visual quality is given in Table 1. Calculated SSIM and LPIPS values suggest that reconstructing deepfakes with up to $s = 50$ denoising steps does not significantly change their structure. A higher value of s on the other hand visibly modifies the initial image. These modifications decrease the CSIM value, indicating that the identity of the face has been degraded to some extent.

These findings are also supported by the qualitative analysis of attacks. As can be seen in Figure 5, there are no obvious, perceivable *structural* differences between unmodified deepfakes ($s = 0$) and attacks denoted by $s = 1$, $s = 10$, $s = 25$ and $s = 50$. We observe that higher levels of noise in general produce more realistic images, in terms of fidelity. Moreover, when $s=100$, the DDM is capable of completely removing typical deepfake inconsistencies. In the 2nd and 3rd row of Figure 5 for instance, the DDM has successfully removed double eyebrows. In the 2nd row, the irregularities in the teeth shape and position are also fixed, while unnatural forehead shadows are eliminated. This noise level, however, often modifies the appearance (size, shape, color) of individual facial parts. In the 1st row of Figure 5 for example, we can see how the DDM has added hair to the initially bald head. Similar changes can be seen in the 4th row, where white hair strands have been inserted. In some cases, subtle changes in the microexpression of the deepfakes can also be observed. Some severe examples, where the application of a high number of diffusion denoising steps ($s = 100$) has drastically changed the initial face are shown in Figure 6. Here, we can see evident differences in soft biometric attributes such as age, hairstyle, race and ethnicity. We can also notice subtle changes in the eyes, nose and lip shape and color.

Vulnerability of deepfake detectors to attacks. We evaluate considered deepfake methods following the protocol described in Section 4 (see *Evaluation metrics*). Figure 7 depicts obtained results in terms of True Positive Rate (TPR). We observe, that attacks generated with only one DDM step ($s = 1$) can severely affect the accuracy of detec-

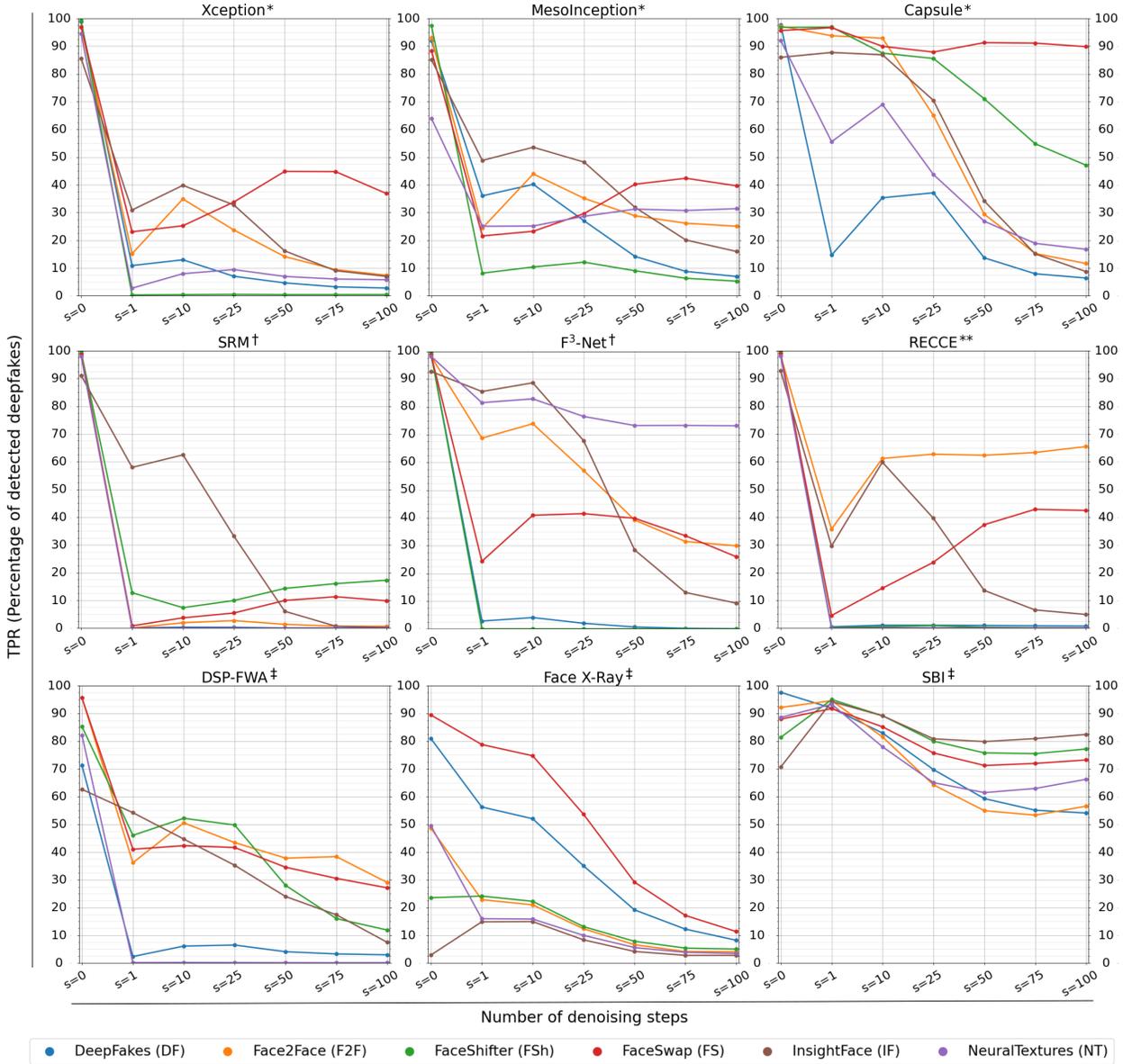


Figure 7. Percentage of **detected FF++ deepfakes ($s = 0$) and attacks ($s > 0$)**. The threshold used for the classification of manipulated images is calculated at the EER point of the ROC evaluated on real and unmodified fake images ($s = 0$) from the validation subset. Types of evaluated detectors: (*) discriminator based on spatial features, (†) discriminator based on frequency features, (**) autoencoder based on reconstruction error, (‡) self-supervised method. We observe, that self-supervised methods (see SBI) are in general more robust to our attacks than discriminative models. However, all considered models, irrespective of their type, experience some degree of adverse impact.

tors. Discriminative methods tend to be more prone to DDM attacks as opposed to self-supervised methods. Capsule is the least vulnerable among discriminative detectors. However, it experiences a severe drop in accuracy when $s > 25$. Naive spatial classifiers Xception and MesolInception, along with the frequency-based classifier SRM show overall worst results in terms of the percentage of detected attacks. These methods maintain a consistent TPR of above 90% across different unmodified deepfakes, but the TPR

of reconstructed deepfakes is in most cases far below 50%. RECCE exhibits similar problems, except when analyzing Face2Face reenacted faces and FaceSwap samples. When these fake images are processed with $s > 1$ denoising steps, the reconstruction error of RECCE is closer to the error of unmodified fakes. In general, higher number of denoising steps s generates more challenging attacks, leading to lower TPR. In some cases, we observe a local peak in discriminative method’s TPR calculated for $s = 10$. We hypothesize,

train: real, unmodified deepfakes, attacks																		
test	Deepfakes (DF)			Face2Face (F2F)			FaceShifter (FSh)			FaceSwap (FS)			InsightFace (IF)			NeuralTextures (NT)		
	s=1	s=50	s=100	s=1	s=50	s=100	s=1	s=50	s=100	s=1	s=50	s=100	s=1	s=50	s=100	s=1	s=50	s=100
s=0	76.02	78.74	84.16	91.16	90.11	95.33	89.17	94.44	85.27	92.14	91.35	93.04	77.74	77.69	74.31	73.99	86.25	70.18
s=1	99.31	96.06	99.27	99.19	96.80	99.57	96.21	96.29	87.21	96.39	96.95	87.43	99.92	93.66	88.36	94.17	99.30	64.73
s=10	45.64	94.21	87.27	64.37	95.03	80.04	89.38	94.03	76.94	85.11	87.04	52.39	73.21	92.78	69.09	92.38	97.85	5.80
s=25	37.28	92.87	88.42	58.17	94.06	83.80	84.65	94.14	82.96	87.57	91.51	69.96	66.80	93.10	82.77	83.76	98.51	23.10
s=50	36.19	93.64	89.55	63.43	95.46	95.75	71.73	95.11	91.88	93.00	97.28	92.21	67.24	93.21	96.35	61.09	99.40	67.98
s=75	31.68	89.87	96.49	63.43	95.27	98.25	61.33	94.46	95.59	93.73	98.21	97.21	67.54	92.70	98.72	46.64	99.70	86.21
s=100	29.14	86.24	98.78	60.69	92.85	98.95	57.05	93.36	96.45	92.90	98.11	97.51	67.43	92.59	99.16	42.86	99.97	91.28
Avg.	50.75	90.23	93.63	71.49	94.23	93.10	78.50	94.55	88.04	91.55	94.35	84.25	74.26	90.82	86.96	70.70	97.28	58.47

Table 2. TPR results indicating the **detectability of black-box attacks**, when Xception is discriminatively trained on a dataset that consists of real samples, unmodified deepfakes, and either $s = 1$, $s = 50$, or $s = 100$ attacks. While training the detector on $s = 1$ attacks gives poor overall results, training on $s = 100$ significantly improves the detection accuracy of all types of attacks. Training the discriminator on $s = 50$ yields the most consistent and effective detection results across different attacks, with a TPR between 90% and 97%.

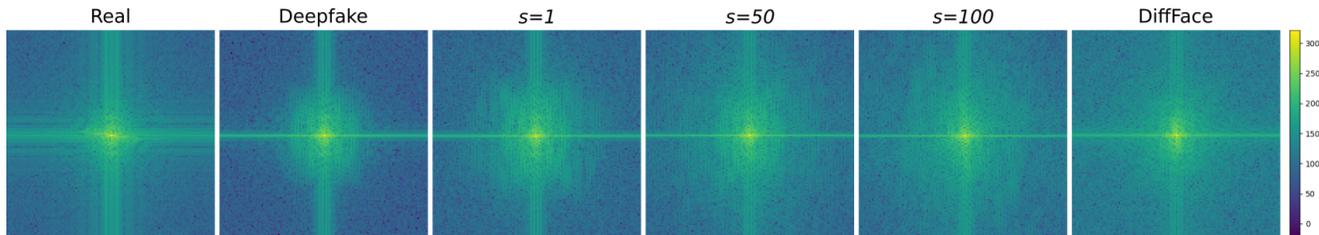


Figure 8. Visualisation of the **Discrete Fourier Transform (DFT) spectra** of real data, unmodified deepfakes, our DDM attacks ($s = 1$, $s = 50$, $s = 100$) and diffusion-based face swapping images generated with DiffFace [15]. Our proposed DDM-based approach for the generation of attacks reintroduces high-frequency elements in the low-density spectral range of deepfakes, mimicking the spectrum of real images. Entirely diffusion-based methods for creation of deepfake data (such as DiffFace) in general tend to overestimate the frequency density of real data. These spectral discrepancies can be utilized, to expose the DDM deepfakes. Best viewed in color and zoomed in.

that this particular amount of denoising steps induces image artifacts similar to those present in unmodified deepfakes.

Unlike discriminative approaches, self-supervised models are trained on simulated fake data. Their performance is therefore strongly dependent on the simulation technique. Based on Figure 7, DSP-FWA, one of the very first self-supervised models, experiences a huge drop in accuracy (at least 30% or more), when attacks with $s = 1$ are presented. In contrast, Face X-Ray, a model trained on manipulated data that mimics deepfakes in a more sophisticated way, has slightly more stable performance on $s = 1$ and $s = 10$ attacks. This method is however very sensitive to modified deepfakes generated with $s > 10$. Moreover, when Face X-Ray is presented with $s = 100$ attacks, it detects less than 10% of them. In our study, the most robust detector is SBI. Although its performance starts dropping gradually when $s > 1$, the TPR in all tested scenarios stays above 50%. Interestingly, SBI initially detects InsightFace deepfakes with a TPR of around 70%, but corresponding attacks are recognized with a higher TPR of 80% or more. We hypothesize that InsightFace traces are not well approximated by the SBI augmentation approach, whereas our DDM introduces frequencies that better resemble modeled artifacts.

Detectability of deepfake attacks. To investigate how detectable our black-box attacks are in discriminative settings,

we conduct a study in which deepfakes reconstructed using our proposed DDM method are incorporated into the training data. Specifically, we train Xception on 3 separate datasets, each containing only one type of attacks. The first set consists of real, unmodified fake data and attacks processed with $s = 1$ diffusion step. In the second and the third training sets, $s = 1$ attacks are replaced by $s = 50$ and $s = 100$ attacks, respectively. We adhere to the same evaluation protocol as before, this time incorporating corresponding validation attacks, as well. Obtained results are presented in Table 2. Our empirical study reveals that, on average, $s = 1$ attacks are the least challenging to detect, regardless of the type of attacks included in the training set. However, when Xception is trained on $s = 1$ reconstructions, it performs poorly in the detection of other attacks. On the contrary, training on $s = 100$ reconstructions greatly improves the detection accuracy of other attacks, as well. Even so, attacks with smaller s values are detected much less accurately, than attacks with s value closer to 100. Training Xception on $s = 50$ attacks achieves the best average results, across different FF++ Deepfake methods. Nonetheless, this discriminator fails to detect between 3% to 10% of deepfakes, depending on the training dataset. The best performance (TPR of 97.28%) is gained in experiments performed on NeuralTextures data, while the worst

performance (90.82%) is measured on the InsightFace data.

To gain deeper insights into the differences between real, deepfake images and DDM-generated attacks, we analyze the data in the frequency domain. The magnitudes of the Discrete Fourier Transform (DFT) for these images are visualized in Figure 8. We notice that deepfakes, in contrast to real images, exhibit a significantly lower spectrum density moving away from the center of the spectrum, which represents dominant low-frequency components. Our DDM approach for the generation of attacks aims to address this deficiency by reintroducing high-frequency elements. Specifically, by reconstructing an existing deepfake with $s = 1$ diffusion step, we start compensating for the missing spectral components. As we increase the number of diffusion steps s , we intensify the higher spectral range. In order to investigate the differences between our deepfake attacks, and fake data generated by a DDM from scratch, in Figure 8 we also visualize the DFT spectrum of DiffFace [15], a diffusion model for face swapping. We note that, unlike our proposed DDM approach, DiffFace overestimates the frequency density, causing spectral discrepancies, that can be utilized to expose deepfakes. Our findings are consistent with observations in [27], where similar results were obtained for other common diffusion models, as well. However, we note that neither our method nor DiffFace create strong spectral artifacts such as regular grids, typical for GAN-based models.

Detecting DDM data and knowledge transfer to attacks.

In our final experiment, we evaluate the ability of a naive discriminator to identify pure diffusion deepfakes and explore its capacity for knowledge transfer, by attacking it with our diffusion-based attacks. For this purpose, we discriminatively train Xception on real, unmodified deepfakes and DiffFace face swaps. Results, obtained by following the previously established protocol, are summarized in Table 3. We note that DiffFace deepfakes, which have a True Positive Rate (TPR) of over 99%, present a lesser challenge for the detector compared to all other non-diffusion Deepfake methods from FF++. Nevertheless, the TPR for these methods still exceeds 93%. InsightFace deepfakes are in this experiment the most challenging to detect, whereas the best detection results are achieved on FaceShifter (FSh), Face2Face (F2F), and Deepfakes (DF). Interestingly, despite the discriminator being trained on DiffFace diffusion samples, which in part share similar frequency spectrum with our reconstructed fakes (as illustrated in Figure 8), it remains highly vulnerable to all three types of attacks considered in this comparison. Higher number of diffusion steps s is in general associated with higher TPR. Nevertheless, the average TPR of attacks is relatively low. Lowest TPR (less than 10%) was measured on FaceShifter and NeuralTextures. Moderately low TPR (between 25 and 45%) was measured on Deepfakes, Face2Face and InsightFace, while highest TPR was calculated on FaceSwap $s = 50$ and

train: real, unmodified deepfakes, DiffFace face swaps						
test	DF	F2F	FSh	FS	IF	NT
DiffFace [15]	100.00	100.00	99.44	99.76	100.00	100.00
s=0	98.54	98.33	98.98	95.26	93.19	96.69
s=1	26.60	30.33	8.47	34.48	32.63	5.80
s=50	29.28	30.71	3.23	66.51	37.18	6.36
s=100	39.17	26.87	2.23	81.30	44.44	8.39
Avg.	58.72	57.25	42.47	75.46	61.49	43.45

Table 3. TPR results achieved by Xception, when it is discriminatively trained on real images, different FF++ face manipulations and DiffFace (DDM) deepfakes. Measured accuracy indicates the **detectability of DDM data and the generalization** capacities of the discriminator. While purely diffusion-based face swaps DiffFace are detected with a TPR of over 99%, the discriminator is still highly vulnerable to our diffusion-based attacks.

$s = 100$ attacks (66.51% and 81.30%, respectively). These isolated peaks in the TPR however, can not be entirely attributed to the training with DiffFace samples, as we observe the same trend in experiments, where we do not use any diffusion-based train images (see Figure 7).

6. Conclusion

Recently discovered Denoising Diffusion Models (DDMs) have shown impressive capabilities for generating highly realistic and convincing images. In this paper, we investigate their potential employment as generators of black-box attacks on deepfake detection systems. We utilize a guided conditional DDM to reconstruct FaceForensics++ (FF++) deepfakes with a predetermined number of diffusion steps, which was found to directly relate to the properties of the generated attack. The proposed DDM approach not only improves the visual quality of an existing deepfake, but it also removes typical deepfake artifacts.

We leverage generated black-box attacks in the evaluation of the vulnerability of commonly used deepfake detectors. Our study reveals that reconstructing deepfakes with just one diffusion step can significantly decrease the accuracy of the deepfake detectors. A higher number of steps generally leads to lower detection accuracy. Self-supervised detectors tend to be more robust than discriminative spatial or frequency-based methods. Training a naive discriminator using attacks generated with a specific number of denoising steps effectively identifies manipulations associated with that particular level of diffusion noise. However, this approach lacks generalizability to attacks produced with a different number of denoising steps. We also investigate the vulnerability of the discriminator that has been trained using purely diffusion-based face manipulation. Although this model is slightly less vulnerable to our DDM attacks in comparison to models trained on non-diffusion samples, it still shows limited generalizability.

References

- [1] InsightFace: 2D and 3D Face Analysis Project, 2019. Accessed on: July 26, 2022. [4](#)
- [2] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. MesoNet: A Compact Facial Video Forgery Detection Network. In *IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7, 2018. [2](#), [5](#)
- [3] Junyi Cao, Chao Ma, Taiping Yao, Shen Chen, Shouhong Ding, and Xiaokang Yang. End-to-End Reconstruction-Classification Learning for Face Forgery Detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4113–4122, 2022. [2](#), [5](#)
- [4] N. Carlini and H. Farid. Evading Deepfake-Image Detectors with White- and Black-Box Attacks. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 2804–2813, 2020. [1](#), [2](#)
- [5] François Chollet. Xception: Deep Learning with Depthwise Separable Convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1800–1807, 2017. [2](#), [5](#)
- [6] F. Croitoru, V. Hondru, R. Ionescu, and M. Shah. Diffusion Models in Vision: A Survey. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 45(09):10850–10869, 2023. [2](#)
- [7] Naser Damer, Meiling Fang, Patrick Siebke, Jan Niklas Kolf, Marco Huber, and Fadi Boutros. MorDIFF: Recognition Vulnerability and Attack Detectability of Face Morphing Attacks Created by Diffusion Autoencoders. In *International Workshop on Biometrics and Forensics (IWBF)*, 2023. [2](#)
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion Models Beat GANs on Image Synthesis. In *Advances in Neural Information Processing Systems (NIPS)*, volume 34, pages 8780–8794, 2021. [1](#), [3](#)
- [9] Saminder Dhesi, Laura Fontes, Pedro Machado, Isibor Kennedy Ihianle, Farhad Fassih Tash, and David Ada Adama. Mitigating Adversarial Attacks in Deepfake Detection: An Exploration of Perturbation and AI Techniques. *arXiv:2302.11704*, 2023. [2](#)
- [10] Apurva Gandhi and Shomik Jain. Adversarial Perturbations Fool Deepfake Detectors. In *International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, 2020. [2](#)
- [11] Yihao Huang, Felix Juefei-Xu, Run Wang, Qing Guo, Lei Ma, Xiaofei Xie, Jianwen Li, Weikai Miao, Yang Liu, and Geguang Pu. FakePolisher: Making DeepFakes More Detection-Evasive by Shallow Reconstruction. In *ACM International Conference on Multimedia (ICM)*, page 1217–1226, 2020. [2](#), [3](#)
- [12] Shehzeen Hussain, Paarth Neekhara, Malhar Jere, Farinaz Koushanfar, and Julian McAuley. Adversarial Deepfakes: Evaluating Vulnerability of Deepfake Detectors to Adversarial Examples. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 3348–3357, 2021. [1](#), [2](#)
- [13] Forrest N. Iandola, Matthew W. Moskewicz, Khalid Ashraf, Song Han, William J. Dally, and Kurt Keutzer. SqueezeNet: Alexnet-Level Accuracy With 50X Fewer Parameters and 1MB Model Size. In *International Conference on Learning Representations (ICLR)*, 2017. [5](#)
- [14] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive Growing of GANs for Improved Quality, Stability, and Variation. In *International Conference on Learning Representations (ICLR)*, 2018. [4](#)
- [15] Kihong Kim, Yunho Kim, Seokju Cho, Junyoung Seo, Jisu Nam, Kychul Lee, Seungryong Kim, and KwangHee Lee. DiffFace: Diffusion-Based Face Swapping with Facial Guidance. *arXiv:2212.13344*, 2022. [2](#), [7](#), [8](#)
- [16] Minchul Kim, Anil K Jain, and Xiaoming Liu. AdaFace: Quality Adaptive Margin for Face Recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. [5](#)
- [17] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5000–5009, 2020. [1](#), [2](#), [5](#)
- [18] Yuezun Li and Siwei Lyu. Exposing DeepFake Videos By Detecting Face Warping Artifacts. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019. [2](#), [5](#)
- [19] Chi Liu, Huajie Chen, Tianqing Zhu, Jun Zhang, and Wanlei Zhou. Making Deepfakes More Spurious: Evading Deep Face Forgery Detection Via Trace Removal Attack. *arXiv:2203.11433*, 2022. [2](#), [3](#)
- [20] Honggu Liu, Xiaodan Li, Wenbo Zhou, Yuefeng Chen, Yuan He, Hui Xue, Weiming Zhang, and Nenghai Yu. Spatial-Phase Shallow Learning: Rethinking Face Forgery Detection in Frequency Domain. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 772–781, 2021. [2](#)
- [21] Zijie Lou, Gang Cao, and Man Lin. Black-Box Attack Against Gan-Generated Image Detector With Contrastive Perturbation. *Engineering Applications of Artificial Intelligence*, 2023. [2](#)
- [22] Yuchen Luo, Yong Zhang, Junchi Yan, and Wei Liu. Generalizing Face Forgery Detection With High-Frequency Features. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16317–16326, 2021. [1](#), [2](#), [5](#)
- [23] Paarth Neekhara, Brian Dolhansky, Joanna Bitton, and Cristian Canton-Ferrer. Adversarial Threats to DeepFake Detection: A Practical Perspective. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 923–932, 2021. [2](#)
- [24] Huy H. Nguyen, Fuming Fang, Junichi Yamagishi, and Isao Echizen. Multi-task Learning for Detecting and Segmenting Manipulated Facial Images and Videos. In *IEEE International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–8, 2019. [1](#), [2](#)
- [25] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using Capsule Networks to Detect Forged Images and Videos. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2307–2311, 2019. [2](#), [5](#)
- [26] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in Frequency: Face Forgery Detection

- by Mining Frequency-Aware Clues. In *Computer Vision – ECCV 2020*, pages 86–103, 2020. [2](#), [5](#)
- [27] Jonas Ricker, Simon Damm, Thorsten Holz, and Asja Fischer. Towards the Detection of Diffusion Model Deepfakes. *arXiv:2210.14571*, 2023. [8](#)
- [28] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Niessner. FaceForensics++: Learning to Detect Manipulated Facial Images. In *IEEE International Conference on Computer Vision (ICCV)*, 2019. [1](#), [3](#)
- [29] Kaede Shiohara and Toshihiko Yamasaki. Detecting Deepfakes with Self-Blended Images. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 18699–18708, 2022. [2](#), [5](#)
- [30] Michał Stypułkowski, Konstantinos Vougioukas, Sen He, Maciej Zieba, Stavros Petridis, and Maja Pantic. Diffused Heads: Diffusion Models Beat GANs on Talking-Face Generation. *arXiv:2301.03396*, 2023. [2](#)
- [31] Chengrui Wang and Weihong Deng. Representative Forgery Mining for Fake Face Detection. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 14923–14932, 2021. [2](#)
- [32] Jia Xiang and Gengming Zhu. Joint Face Detection and Facial Expression Recognition with MTCNN. In *International Conference on Information Science and Control Engineering (ICISCE)*, pages 424–427, 2017. [4](#)
- [33] Zongsheng Yue and Chen Change Loy. DifFace: Blind Face Restoration with Diffused Error Contraction. *arXiv:2212.06512*, 2022. [3](#)
- [34] W. Zhao, Y. Rao, W. Shi, Z. Liu, J. Zhou, and J. Lu. Diff-Swap: High-Fidelity and Controllable Face Swapping via 3D-Aware Masked Diffusion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8568–8577, 2023. [2](#)
- [35] Kaifeng Zou, Sylvain Faisan, Boyang Yu, Sébastien Valette, and Hyewon Seo. 4D Facial Expression Diffusion Model. *arXiv:2303.16611*, 2023. [2](#)