

# Generalized Deepfakes Detection with Reconstructed-Blended Images and Multi-scale Feature Reconstruction Network

Yuyang Sun<sup>1,2</sup>, Huy H. Nguyen<sup>2</sup>, Chun-Shien Lu<sup>3</sup>, ZhiYong Zhang<sup>4</sup>, Lu Sun<sup>5</sup> and Isao Echizen<sup>1,2</sup>

<sup>1</sup>The University of Tokyo, Japan    <sup>2</sup>National Institute of Informatics, Japan    <sup>3</sup>Academia Sinica, Taiwan

<sup>4</sup>Sun Yat-sen University, China    <sup>5</sup>South China University of Technology, China

{tarrysun, nhhuy, iechizen}@nii.ac.jp

## Abstract

The growing diversity of digital face manipulation techniques has led to an urgent need for a universal and robust detection technology to mitigate the risks posed by malicious forgeries. We present a blended-based detection approach that has robust applicability to unseen datasets. It combines a method for generating synthetic training samples, i.e., reconstructed blended images, that incorporate potential deepfake generator artifacts and a detection model, a multi-scale feature reconstruction network, for capturing the generic boundary artifacts and noise distribution anomalies brought about by digital face manipulation. Experiments demonstrated that this approach results in better performance in both cross-manipulation detection and cross-dataset detection on unseen data.

## 1. Introduction

Digital facial manipulation, exemplified by deepfake technology, entails substituting or altering facial attributes, expressions, and appearances to create highly deceptive visual content. The advancement of deep generative models [12, 13, 19, 25–27] and computer vision technologies has facilitated the creation of efficient, automated deepfake pipelines (e.g., [5]), enabling non-experts to easily manipulate visual media content. However, this ease of use has led to increasing misuse in various domains (politics, journalism, mass media, etc.) and violations of individual privacy.

Early detection models [6, 7, 34, 40, 46, 50, 58] are susceptible to overfitting specific artifacts and noise patterns. They are deficient in cross-domain detection, resulting in substantial performance deterioration when confronted with unseen datasets. Researchers thus introduced blended-based deepfake detection methods [33, 47, 60] that create blended samples through the fusion of suitable genuine facial samples, thereby simulating prevalent manipulation artifacts such as misalignment, incongruent boundary

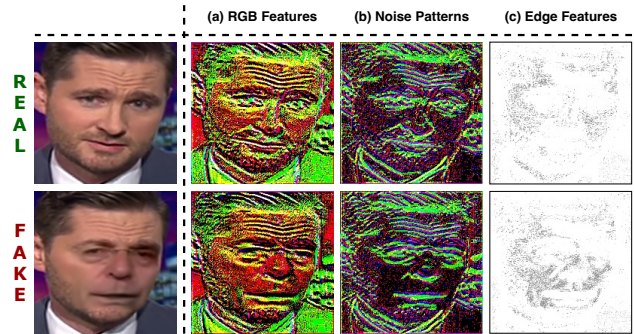


Figure 1. (a) RGB features, (b) noise patterns, and (c) edge features extracted from the first block of EfficientNet-B4 for real image (top row) and corresponding fake image (bottom row).

textures, and irregularities in pixel or color distribution.

Previous methods for constructing blended samples typically extract facial entities from the original foreground after they undergo certain augmentations and then graft them onto alternative original backgrounds. This results in synthetic blended samples that lack certain invisible artifacts, including unique generator fingerprints [31, 39] and anomalies in channel distribution introduced by adaptive instance normalization [23]. Furthermore, previous methods for detecting forgeries predominantly perform binary classification with foundational backbones or learned self-supervised consistency features when training from blended samples, resulting in inadequate harnessing of the potential of diverse annotations inherent in synthesized samples.

We have devised a novel approach to detecting forgeries that combines a method for generating synthetic forgery training samples and a detection model for capturing the generic boundary artifacts and noise distribution anomalies brought about by digital face manipulation. Our reconstructed blended image (RBI) method builds upon advanced deepfake generation techniques [10, 28, 32, 54, 55] to disentangle identity and background information from genuine

facial images. Random Gaussian noise is infused into a background vector with a predefined probability, and the identity-background pairs are passed through a decoder to obtain reconstructed images. This process incorporates latent generator patterns and distinctive fingerprints into the synthesized blended data. After application of the statistical transformations of Shiohara and Yamasaki [47], the reconstructed image is blended with the original image and a mask is used to generate a training sample.

Our detection model, a multi-scale feature reconstruction network (MFRN), effectively exploits diverse training artifacts and manipulated regions. Built upon prior research [11, 14, 33, 47], it focuses on the three possible origins of forged features shown in Figure 1: RGB features, noise patterns, and edge features. RGB features include color inconsistencies between inner and outer facial regions and an abnormal image channel covariance, noise patterns include mismatched image frequencies in manipulated regions and invisible generator fingerprints, and edge features include boundary conflicts arising from face splicing and landmark misalignment. A specifically designed architecture is used to extract the corresponding features from the input and combine them at multiple scales to self-supervise the reconstruction of blended boundaries and regions.

We conducted cross-manipulation and cross-dataset detection assessments on unseen data to align with real-world detection scenarios. Our model was trained solely on genuine data from the FF++ dataset [45]. For cross-manipulation detection, it achieved areas under the curve (AUCs) of 100%, 100%, 99.88%, 99.81%, and 98.90% respectively, for DeepFake (DF) [2], Face2Face (F2F) [52], FaceSwap (FS) [3], NeuralTextures (NT) [53], and FaceShifter (FSH) [32] manipulation. For cross-dataset detection, it demonstrated robust detection on prominent deepfake detection datasets such as CDF-v2 [35] (AUC of 95.27%), DFD [1] (99.12%), DFDC [17] (73.31%), and DFDC-P [16] (83.66%). These results surpass or match those of current state-of-the-art baselines.

Our contributions are summarized as follows:

- We introduce a method for generating simulated forgery training samples (reconstructed-blended images) that enhances the diversity of simulated artifacts by incorporating an invisible generator fingerprint and noise pattern.
- We introduce a model (a multi-scale feature reconstruction network) for harnessing the diversity in training artifacts and manipulated regions present in blended samples.
- We validated the performance of this approach across multiple unseen manipulation techniques and datasets and demonstrated its superiority.

## 2. Related Work

### 2.1. Deepfake Generation

Modern deepfake generation techniques primarily use deep generative models [19, 22, 30, 44] to achieve seamless integration and migration of facial features and expressions, enabling creation of highly realistic visual content. Earlier deepfake models [2, 3] separated the process of generating and blending forged faces, leading to noticeable stitching artifacts and easy human eye detection.

Subsequent deepfake models addressed this problem by combining the generation and blending into a single step. Incorporating the concept of style transfer, these models consider the target identity information as a style feature and use techniques like AdaIN [23] to align the target identity with the source background, resulting in end-to-end deepfake generation pipelines that improve realism and reduce visible artifacts. Several of these models [32, 54, 55] predict masks for the target background, enabling precise control over manipulated areas. Others [10, 28] do not require such masks and instead autonomously learn the content that needs to be modified or retained during training.

### 2.2. Deepfake Detection

Detection techniques that combine spatial-frequency statistical analysis with deep neural networks have shown marked efficacy in identifying image manipulation artifacts. These methods [11, 14, 21, 48] focus on detecting abnormal edges and noise features introduced by splicing together disparate source images, resulting in the capture of imperceptible manipulation artifacts. Deepfake detection, a specialized domain within image manipulation detection, can be categorized into artifact-based and blended-based detection.

**Artifact-based Detection** is aimed at identifying potential deepfakes by detecting prevalent and distributed artifacts or imperfections introduced during manipulation. Early research focused on visible artifacts such as abnormal facial expressions and head movements [7], inconsistencies in predicted head pose [58], image photo response non-uniformity [31], and color cues [39]. As forged images have become increasingly realistic, manually specified artifacts have proven insufficient. Subsequent research [6, 37, 40, 43, 59] has focused on using neural networks to detect spatial-frequency pattern differences within manipulated images. Additionally, researchers have observed that frame-by-frame manipulation of videos can introduce temporal inconsistencies at the interframe level, leading to efforts [20, 38, 46, 49, 61] to uncover the absence of temporal correlations within video streams.

**Blended-based Detection** is aimed at identifying potential deepfakes by using forgery training samples generated by blending pairs of genuine samples that have undergone flexible image augmentations through diverse blending, cre-

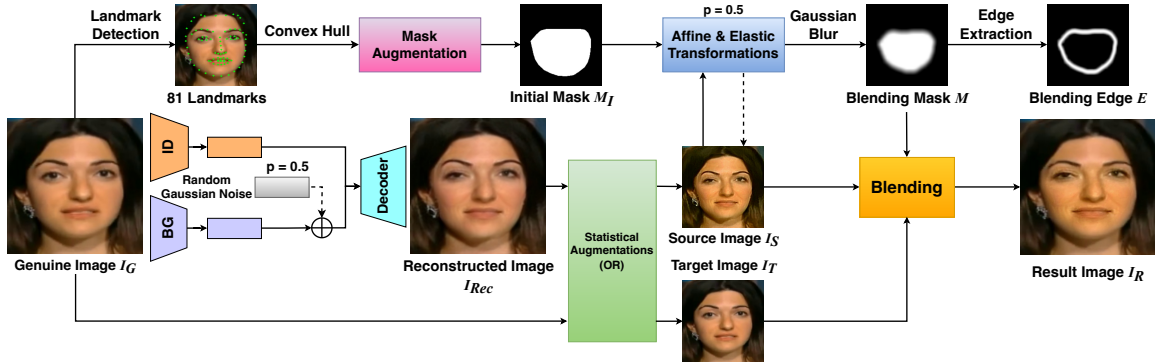


Figure 2. **Illustration of RBI generation.** Information in given genuine image  $I_G$  is disentangled, and image is reconstructed as  $I_{Rec}$ , incorporating a unique fingerprint and noise patterns. Statistical transformations are then randomly applied to create source image  $I_S$  and target image  $I_T$  in which common visible manipulation artifacts are simulated. Concurrently,  $I_G$ 's convex hull is augmented and deformed to produce mask  $M$  and edge  $E$ , which are used guide the blending of  $I_S$  and  $I_T$ , yielding result image  $I_R$ .

ating an array of manipulation artifacts. These methods aim for high generalization performance on unseen datasets by incorporating various manipulation patterns. Early attempts in this category, such as DSP-FWA (Dual Spatial Pyramid for Exposing Face Warp Artifacts in DeepFake Videos) [34], simulate distortion artifacts introduced by affine transformations in the deepfake pipeline by blending downsampled and Gaussian-blurred facial regions with the background. Subsequent methods like Face X-Ray [33] and patch-wise consistency learning (PCL)+ inconsistency image generation (I2G) [60] simulate blending boundaries or inconsistent information by replacing and blending similar facial images from the dataset. The use of self-blended [47] images has improved detection performance and achieved higher data efficiency by generating pseudo source-target image pairs through statistical transformations applied to a single image, followed by blending.

### 3. Proposed Approach

In our approach to enhancing the generalization of deepfake detection, we introduce a novel synthetic training sample, namely the reconstructed blended image (RBI), and a corresponding detection framework, namely multi-scale feature reconstruction network (MFRN).

As illustrated in Figure 2, RBI is generated by extracting identity and background features from the source image, embedding imperceptible latent generator artifacts, integration of common visible image artifacts through statistical transformations, and finally blending with the source image through a randomly deformed mask. This results in diversity and enrichment in the synthetic training samples. Some examples of RBIs and detailed explanations of relative artifacts can be found in Figure 6 and Section 7 in the appendix.

Once the samples have been generated alongside the

corresponding blending masks and boundary annotations, as depicted in Figure 3, a convolutional neural network-based reconstruction model, the MFRN, is created. This model captures RGB, edge, and noise features across various scales, thus enabling comprehension of the mapping from blended samples to region labels. This enables the model to autonomously identify discrepancies introduced by manipulation and to detect conflicting boundary textures.

#### 3.1. Reconstructed Blended Image Generation

Genuine image  $I_G$  is disentangled using an identity (ID) encoder and background (BG) encoder. The resulting identity-background pair is then decoded, producing reconstructed image  $I_{Rec}$ . This step disrupts innate noise patterns within the genuine image and thereby introduces a simulated generator fingerprint. Encoding is diversified by adding random Gaussian noise with mean ( $\mu$ ) 0 and standard deviation ( $\sigma$ ) in  $[0.1, 0.3]$  to the background vector with probability  $p = 0.5$ . Note that this encoder-decoder framework can be replaced with an analogous deepfake generator that combines identity and background features. We used a pre-trained SimSwap [10] model to perform the reconstruction.

Next, in accordance with the method outlined in [47], statistical augmentations are equally applied to either  $I_G$  or  $I_{Rec}$ , forming source ( $I_S$ ) and target ( $I_T$ ) images for blending. This integration introduces evident spatial or frequency irregularities like those commonly encountered in manipulated images. Specifically, following operations are executed on the each image: random RGB channel value shift, random adjustment of hue, saturation, and brightness, random modification of brightness and contrast, and random application of blurring or sharpening.

Concurrently, 81 landmarks are extracted from  $I_G$  using the Dlib library [29] and are used to form initial mask  $M_I$  by convex hull calculation. The hull is then diversified

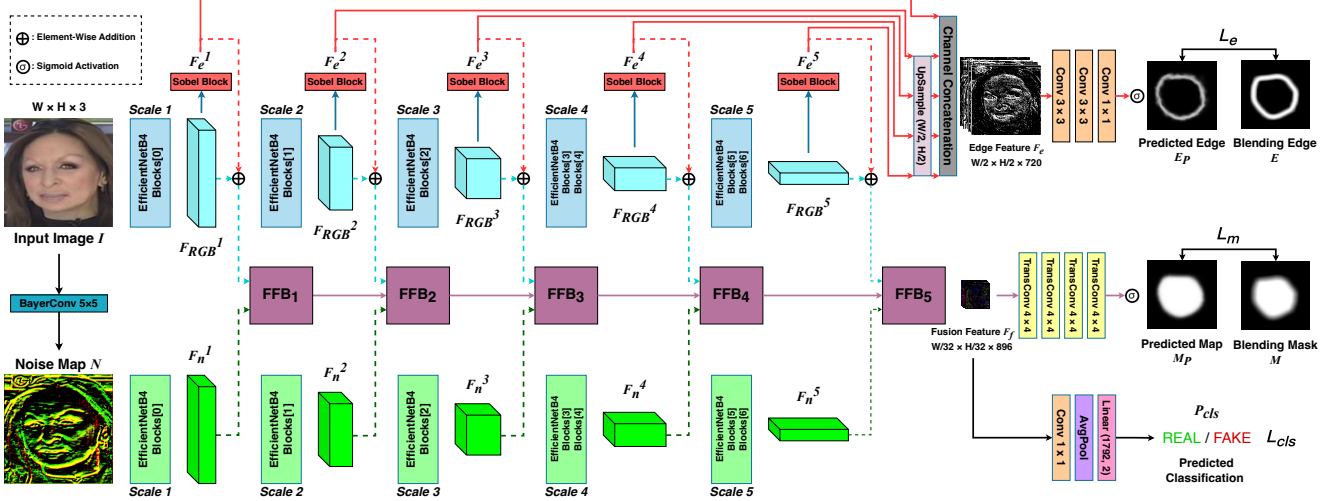


Figure 3. **Overview of MFRN architecture.** First, noise map  $N$  of input image  $I$  is acquired. Next, RGB features  $F_{RGB}$  and noise patterns  $F_n$  are independently extracted across diverse scales by using separate backbones.  $F_{RGB}$  undergoes Sobel block to create corresponding edge features  $F_e$  at their respective scales, which are upsampled for manipulation boundary prediction ( $E_p$ ). Concurrently, a feature fusion block integrates features from the diverse scales for manipulation map prediction ( $M_p$ ) and authenticity classification.

using the *random\_get\_hull* from open source tool [4] for mask augmentation. With a probability of  $p = 0.5$ , identical affine and elastic transformations are applied to both  $M_I$  and  $I_S$ , simulating landmark discrepancies and boundary conflicts introduced by warping in the deepfake pipeline.

Finally, transformed  $M_I$  is subjected to Gaussian blur, yielding blending mask  $M$ , which is used to guide the blending process. The blending of  $I_S$  and  $I_T$  using  $M$  to obtain  $I_R$  is done using

$$I_R = I_S \odot \alpha M + I_T \odot (1 - \alpha M) \quad (1)$$

where  $\alpha \in [0.5, 1]$  regulates the magnitude of the blending mask. Corresponding blending edge  $E$  can be readily extracted [33]:

$$E = 4 \cdot M \odot (1 - M) \quad (2)$$

### 3.2. Multi-scale Feature Reconstruction Network

Input image  $I \in \mathbb{R}^{W \times H \times 3}$  is converted to grayscale, and noise map  $N$  is derived using a  $5 \times 5$  Bayer constrained convolutional layer [8]. Two separate EfficientNet-B4 [51] backbones are then used to capture features from both  $I$  and  $N$ . For clarity, the output of the  $i^{th}$  downsampling block is designated as being at scale  $i$ . At each scale  $i$ , feature maps  $F_{RGB}^i$  and  $F_n^i$  are extracted from the RGB and noise branch, respectively. The edge features are computed from the RGB feature using the Sobel block depicted in Figure 4(a), which involves applying two fixed-parameter  $3 \times 3$  Sobel filters to  $F_{RGB}^i$ , followed by batch normalization [24] and Sigmoid activation. The resulting output is element-wise multiplied

with  $F_{RGB}^i$ , and integration is achieved using a  $1 \times 1$  convolutional layer to produce  $F_e^i$ . This procedure can be summarized by the following formula:

$$Edge\_Act_i = \sigma((BN(SobelConv(F_{RGB}^i))))$$

$$F_e^i = Conv(F_{RGB}^i \odot Edge\_Act_i) \quad (3)$$

Subsequently, the extracted feature maps are split into two pathways and used to individually reconstruct the modified boundaries and regions within the blended image.

For one pathway,  $F_e^i$  at each scale is uniformly upsampled to  $(\frac{W}{2}, \frac{H}{2})$ , which, after channel-wise concatenation, yields feature  $F_e \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2} \times 720}$ . Two  $3 \times 3$  convolutional kernels are used to capture local information, and a  $1 \times 1$  convolutional kernel is used to integrate channel features, finally producing predicted edge  $E_p \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2}}$ .

For another pathway, feature fusion block (FFB) is introduced to facilitate both within-scale feature fusion and across-scale feature propagation. As illustrated in 4(b), at scale  $i$ ,  $FFB_i$  receives  $F_n^i$ , the element-wise summation of  $F_{RGB}^i$  and  $F_e^i$ , and the output feature from  $FFB_{i-1}$  (excluding the first FFB block) as input. The rationale behind summing  $F_{RGB}^i$  and  $F_e^i$  stems from their shared provenance in the RGB domain. Subsequently, a bottleneck attention module (BAM) [41] applied to the concatenated extracted features establishes a self-attention mechanism that directs attention to anomalies within manipulated areas on the basis of the spatial disposition and channel distribution of feature map. Following the processing of the two convolutional layers, weighted feature  $FFB_w^i$  is obtained for the present



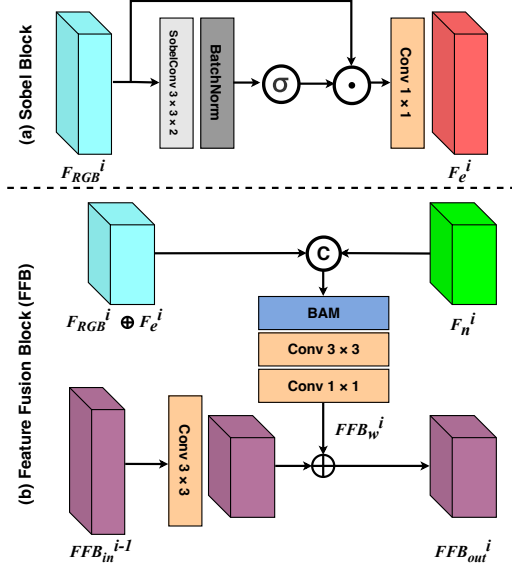


Figure 4. Architecture of Sobel block and feature fusion block.

scale. A  $3 \times 3$  convolutional kernel is used to align the size of the propagated feature originating at the previous scale, denoted as  $FFB_{in}^{i-1}$ , with that of the one at the current scale. Subsequent element-wise summation with the current weighted features results in the output feature  $FFB_{out}^i$ , which is then propagated to  $FFB_{i+1}$  after ReLU activation. These operations can be succinctly encapsulated:

$$FFB_w^i = Conv_{\times 2}(BAM(Cat(F_{RGB}^i \oplus F_e^i, F_n^i)))$$

$$FFB_{out}^i = FFB_w^i \oplus Conv(FFB_{in}^{i-1}) \quad (4)$$

After the features of the different modalities at the various scales are fused and passed hierarchically, fusion feature  $F_f \in \mathbb{R}^{\frac{W}{32} \times \frac{H}{32} \times 896}$  is obtained. Several transposed convolutions with a kernel size of 4 are used to predict modification map  $M_p \in \mathbb{R}^{\frac{W}{2} \times \frac{H}{2}}$ . Furthermore, a parallel classification branch, encompassing a convolutional head, a pooling layer, and a dense layer, is created and used to predict two-dimensional classification vector  $P_{cls}$ , which is derived from the  $F_f$  feature.

Clearly, the lack of manipulation in authentic instances means that outcomes  $E_p$  and  $M_p$  for a genuine image should manifest as matrices consisting exclusively of zeros.

### 3.3. Loss Functions

**Map Loss and Edge Loss.** Element-wise binary cross-entropy (BCE) loss is used to assess the dissimilarity between the predicted map and the blending mask, as well as between the predicted edge and the blending edge:

$$L_e = -\frac{1}{N} \sum_{i,j} (E^{i,j} \log E_p^{i,j} + (1 - E^{i,j}) \log(1 - E_p^{i,j})) \quad (5)$$

$$L_m = -\frac{1}{N} \sum_{i,j} (M^{i,j} \log M_p^{i,j} + (1 - M^{i,j}) \log(1 - M_p^{i,j})) \quad (6)$$

where  $N = \frac{WH}{4}$  is the number of pixels in the  $E_p$  and  $M_p$ . **Classification Loss.** BCE loss is used to quantify the classification error, given by

$$L_{cls} = T_{cls} \log P_{cls} + (1 - T_{cls}) \log(1 - P_{cls}) \quad (7)$$

where  $T_{cls} = \{0, 1\}$  represents the ground truth label of the input sample.

The overall loss function of our model is

$$L = \lambda_1 L_m + \lambda_2 L_e + L_{cls} \quad (8)$$

where  $\lambda_1$  and  $\lambda_2$  are scaling factors used to regulate individual loss proportions.

## 4. Experiments

### 4.1. Setup

**Training Data.** Our model was trained solely on pristine FaceForensics++ [45] (FF++) data, following the official dataset split, which included 720 training videos. We uniformly extracted 20 frames from each video and utilized Dlib [29] to extract 81 facial landmarks from each frame, which were used to compute the initial mask. We used the RetinaFace face detector [15] for bounding box localization and face cropping. In cases with multiple faces detected in a video frame, we overlapped the bounding boxes with a corresponding deepfake manipulation mask, selecting the result with the largest intersection area for validity.

**Test Data.** We assessed model performance through cross-manipulation and cross-dataset evaluation. For cross-manipulation evaluation, we used unseen manipulation videos from the FF++ test set, including Deepfakes [2] (DF), Face2Face [52] (F2F), FaceSwap [3] (FS), NeuralTextures [53] (NT), and FaceShifter [32] (FSH), totaling 840 test videos (140 from the pristine dataset and each manipulation dataset), following official dataset splits. For cross-dataset evaluation, we utilized several prominent digital face manipulation datasets: DeepFakeDetection [1] (DFD), Celeb-DF-v2 [35] (CDF-v2), DeepFake Detection Challenge [17] (DFDC), and its preview version [16] (DFDC-P). We used the complete DFD dataset for testing and adhered to official data splits for the other datasets, using their designated test videos. For all videos, we uniformly sampled frames and used RetinaFace for bounding box extraction and face cropping. As the current models had reached saturation in FF++ detection performance, we assessed cross-manipulation at the individual sample level, extracting and cropping five faces per video to construct the test set. Conversely, we evaluated cross-dataset performance at the video level, extracting 32 frames evenly, cropping faces, and aggregating model predictions for all

Method	Test AUC (%)						
	DF	F2F	FS	NT	FF++(w/o FSH)	FSH	FF++
Face X-Ray [33]	99.17	98.57	98.21	98.13	98.52	-	-
PCL + I2G [60]	<b>100</b>	98.97	99.86	97.63	99.11	-	-
Self-blended [47]	99.99	99.88	<b>99.91</b>	98.79	99.64	-	-
Self-blended * [47]	99.94	99.72	99.76	98.23	99.42	97.70	99.07
RBIs + MFRN (Ours)	<b>100</b>	<b>100</b>	99.88	<b>99.81</b>	<b>99.92</b>	<b>98.90</b>	<b>99.71</b>

Table 1. **Cross-manipulation comparison among blended-based methods.** Results are cited from the papers, with \* indicating official pre-trained model results on our test data. Our method outperformed baselines on DF, F2F, NT, FSH, and the complete FF++ dataset.

Method	Type	Test AUC (%)			
		CDF-v2	DFD	DFDC	DFDC-P
DeepRhythm [42]		-	-	-	64.1
MultiATT [59]		67.44	-	-	-
FRDM [37]	Artifact	79.4	91.9	-	79.7
RECCE [9]		68.71	-	69.06	-
FTCN [61]		86.9	-	-	74.0
Face X-Ray [33]		-	93.47	-	71.15
PCL + I2G [60]	Blended	90.03	<u>99.07</u>	67.52	74.37
Self-blended [47]		<u>93.18</u>	97.56	<u>72.42</u>	<b>86.15</b>
RBIs + MFRN (Ours)	Blended	<b>95.27</b>	<b>99.12</b>	<b>73.31</b>	<u>83.66</u>

Table 2. **Cross-dataset detection comparison with baselines.** With best and second-best results indicated in bold and underline. Our method outperformed the baselines on CDF-v2, DFD, and DFDC and was second best on DFDC-P.

32 faces per video on the basis of their mean. To ensure fairness, videos where face extraction was not possible were assigned a prediction value of 0.5.

**Comparison Baseline.** We used advanced baselines to evaluate our model, categorizing them into two types: artifact-based and blended-based. The former type requires both genuine and manipulated samples for training and included Multi-Attentional Deepfake Detection [59] (MultiAtt), Fusion + RSA + DCMA + Multi-scale [37] (FRDM), Uncovering Common Feature [56] (UCF), Spatial-Phase Shallow Learning [36] (SPSL), Reconstruction-Classification Learning [9] (RECCE), Fully Temporal Convolution Network [61] (FTCN) and DeepRhythm [42]. In contrast, the latter type requires only genuine samples and included Face X-Ray [33], PCL+I2G [60], and Self-blended[47].

**Implementation Details.** To generate RBIs, we used the official pre-trained SimSwap model as the reconstruction generator and applied common image augmentations, such as JPEG compression, brightness-contrast adjustments, and color jittering, to both synthesized RBIs and genuine samples. The MFRN used the pre-trained EfficientNet-b4 [51]

backbone. All facial regions cropped by bounding boxes were resized to  $380 \times 380$  to match pre-training specifications. We set  $\lambda_1$  and  $\lambda_2$  to 100 and 50, respectively; additional insights into the effect of varying loss weights  $\lambda$  are available in Table 7 in the appendix. To enhance stability and generalization performance, we used the sharpness-aware minimization (SAM) optimizer [18]. Training was performed over 80 epochs on a NVIDIA A100 (80G) GPU, with a learning rate of 0.001 and a batch size of 32.

**Evaluation Metrics.** Due to the highly imbalanced distribution of genuine and manipulated samples in the datasets used for evaluation, we primarily used the AUC as our evaluation metric as it better reflects the model’s performance.

## 4.2. Cross-Manipulation Detection

We conducted cross-manipulation detection experiments on the FF++ raw data to evaluate our model’s detection ability on unseen manipulations. We maintained fairness by comparing our model with similar blended-based models, citing their reported results directly. Additionally, for assessing the performance of FSH, we utilized the official pre-trained Self-blended model on our test data (the results for Face X-Ray and I2G+PCL could not be replicated due to the lack of official implementations), indicated by an asterisk (\*).

As shown in Table 1, in scenarios nearing performance saturation, our model achieved AUCs of 100%, 100%, and 99.81% for common DF, F2F, and NT manipulations, respectively, surpassing the baseline models. It slightly underperformed with a 98.88% result on FS. It outperformed Self-blended by 1.2% on FSH, illustrating its enhanced ability to capture generator fingerprints and noise patterns introduced by one-stage deepfake generators. Our method consistently achieved optimal detection AUCs on the complete FF++ test set, whether considering FSH or not, with scores of 99.71% and 99.92%, respectively.

## 4.3. Cross-Dataset Detection

As previously mentioned, we validated our method’s performance across multiple unseen mainstream forgery datasets. We compared our method’s performance with those of

Method	Training Set	Test AUC (%)			
		CDF-v2	DFD	DFDC	DFDC-P
Self-blended [47]	CDF-v2	93.74	-	-	81.10
RBIs + MFRN (Ours)		93.53	98.25	73.40	85.21
SPSL [36]	FF-c23	76.50	81.22	70.40	74.08
UCF [56]		75.27	80.74	71.91	75.94
FRDM [37]		75.52	81.20	69.95	74.08
RBIs + MFRN (Ours)		<b>93.89</b>	<b>98.39</b>	<b>72.70</b>	<b>81.72</b>

Table 3. Cross-dataset performance of model trained on the CDF-v2 dataset and FF-c23 data. The results of SPSL, UCF, and FRDM were excerpted from DeepfakeBench [57].

artifact-based and blended-based detection methods.

As evident in Table 2, blended-based methods have a pronounced advantage over artifact-based ones when used in cross-dataset detection tasks. This observation supports our statement above that models trained on limited manipulated data are susceptible to overfitting specific artifacts and noise patterns, making them deficient in cross-domain detection. Among the artifact-based methodologies, FTCN, which takes into consideration interframe temporal features, demonstrated superior generalization capabilities, achieving an AUC of 86.9% on CDF-v2. Among the blended-based methods, ours had better outcomes. Its performance surpassed that of the compared baselines with AUCs of 95.27%, 99.12%, and 73.31% on CDF-v2, DFD, and DFDC, respectively. On DFDC-P, its AUC of 83.66% is slightly lower than that of Self-blended (86.15%) and better than those of the other compared baselines. In summation, our model consistently delivered superior performance in cross-dataset detection scenarios.

To assess our method’s wider applicability, we also conducted training on samples from varying dataset (CDF-v2) and varying compression level (FF-c23), followed by cross-dataset performance evaluation on unseen datasets in a consistent manner. For the latter, to ensure fairness, we selected three top-performing methods from the DeepfakeBench [57], which were trained on FF-c23, as comparative baselines. As shown in Table 3, our model shows robust adaptability to unseen data, demonstrating its efficacy even when trained on different dataset or low-quality samples. From the results trained on FF-23, it outperformed prominent state-of-the-art baselines, emphasizing its advantages and affirming its position as a robust solution in the field.

#### 4.4. Video Compression Robustness

In the context of digital face manipulation detection, a model’s resilience to compression is pivotal given that real-world detection scenarios often entail acquiring highly compressed samples from the Internet. We conducted evaluations using test samples from videos subjected to moderate

Method	AUC on Different Video Compressions (%)		
	c0 (raw)	c23	c40
Face X-Ray [33]	98.52	87.35	61.60
Self-blended * [47]	99.42	88.32	65.29
RBIs + MFRN (Ours)	99.92	88.55	64.63

Table 4. Robustness for different levels of compression of FF++ dataset, with \* indicating official pre-trained model results on our prepared test data.

(c23) and heavy (c40) compression that were taken from the FF++ dataset. Performance was compared with those of Face X-Ray and Self-blended. The latter’s results were derived by assessing the pretrained model on our test data.

As shown in Table 4, Face X-Ray is vulnerable to compression, with AUCs of 87.35% and 61.60% on c23 and c40, respectively. Our model demonstrates robust detection performance under moderate compression conditions but falls short in heavily compressed scenarios when compared with the simpler detection structure. This discrepancy may stem from our noise pattern extraction and RGB feature specification, which are affected by the substantial loss of image information in heavily compressed data.

#### 4.5. Ablation Study

**Effectiveness of RBIs + MFRN method.** We conducted two intermediate experiments, denoted as self-blended images “SBIs + MFRN” and “RBIs + EfficientNet-b4” to assess the effectiveness of our proposed method. As shown in Table 5, MFRN enhanced the identification of manipulated artifacts, surpassing the performance of the original EfficientNet-b4 backbone and yielding superior detection performance. The inclusion of RBIs introduces diverse visible artifacts and invisible noise anomalies, enhancing the authenticity and diversity of the synthetic training data when compared with the use of SBIs and thereby also contributes to improved detection performance.

**Effectiveness of Components in MFRN.** We conducted experiments with the following settings to validate the effectiveness of each component in MFRN: 1) Edge reconstruction branch removed, and RGB features and noise patterns fused at multiple scales to reconstruct manipulated regions; 2) FFB removed, and, when reconstructing manipulated regions, edge features, RGB features, and noise patterns at last scale are concatenated and used as inputs; 3) Both FFB and noise pattern branch removed, and, when reconstructing manipulated regions, edge features and RGB features at last scale are concatenated and used as inputs.

As shown in Table 6, the absence of the edge reconstruction branch resulted in substantially poorer detection performance. This was due to manipulation techniques often causing unnatural boundary conflicts in the replaced ar-

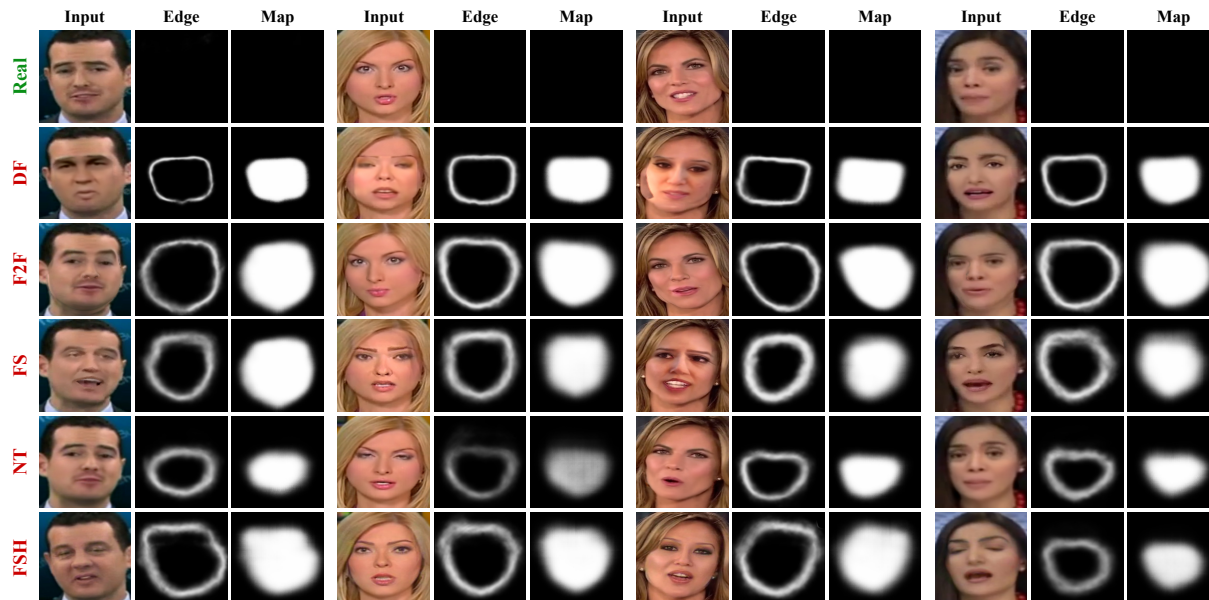


Figure 5. Visualization of predicted edges and maps for genuine images alongside various manipulated samples. Model was trained on real data from FF++ dataset augmented by RBIs. Inputs used for prediction came from unseen FF++ test dataset.

Settings	Test AUC (%)			Avg.
	CDF-v2	DFD	FF++	
SBIs + EfficientNet-b4	93.18	97.56	99.64	96.79
RBIs + EfficientNet-b4	92.51	98.64	99.72	96.96
SBIs + MFRN	95.21	98.27	99.83	97.77
RBIs + MFRN (Ours)	<b>95.27</b>	<b>99.12</b>	<b>99.92</b>	<b>98.10</b>

Table 5. Results of the ablation study on RBIs and MFRN.

eas, and edge reconstruction effectively aids the model in capturing these misaligned artifacts. In complex detection backgrounds like DFDC-P, models without the noise branch exhibited a notable decrease in detection performance. Furthermore, our FFB effectively boosted the model’s detection performance through multi-scale self-attention-based fusion of diverse feature modalities.

#### 4.6. Visualization

We visualized prediction outcomes on real and manipulated samples from the unseen FF++ test set. These inputs, used with a model trained exclusively on FF++ real data augmented by RBIs, produced prediction manipulation regions and boundary conflicts, as seen in Figure 5. Our approach effectively aids the model in recognizing irregular patterns in unseen manipulated data and accurately delineating the boundaries of replaced regions, even without direct training on similar manipulations. This effectiveness can be attributed to the use of comprehensive synthetic training data

Settings	Test AUC (%)			Avg.
	CDF-v2	DFDC-P	FF++	
w/o Edge Reconstruction	93.34	82.71	99.86	91.97
w/o FFB	94.26	83.25	<b>99.94</b>	92.48
w/o Noise Pattern + FFB	95.02	79.45	99.90	91.46
RBIs + MFRN (Ours)	<b>95.27</b>	<b>83.66</b>	99.92	<b>92.95</b>

Table 6. Results of the ablation study to evaluate the effectiveness of individual components within MFRN.

and our purpose-designed feature reconstruction network. More visualization results and explanations can be found in Figure 7 and Section 8 in the appendix.

## 5. Conclusion

We have presented an innovative method for synthesizing forgery training samples, i.e., reconstructed blended images (RBIs). It improves the ability to simulate manipulation artifacts by seamlessly integrating simulated generator fingerprints and noise patterns. We also presented a novel detection model, the multi-scale feature reconstruction network (MFRN), which adeptly exploits the richness of diversity introduced by the use of random blending masks and boundaries within the RBIs. Experimental results demonstrated that our proposed approach substantially enhances performance in both unseen cross-dataset and cross-manipulation detection.



## Acknowledgments

This work was partially supported by JSPS KAKENHI Grant JP21H04907, and by JST CREST Grants JPMJCR18A6 and JPMJCR20D3, Japan.

## References

- [1] DeepfakeDetection. <https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html>. Accessed: 2023-08-10. 2, 5
- [2] Deepfakes Generation. <https://github.com/deepfakes/faceswap>. Accessed: 2023-08-10. 2, 5
- [3] Faceswap Generation. <https://github.com/MarekKowalski/FaceSwap>. Accessed: 2023-08-10. 2, 5
- [4] Mask Augmentation Tool. [https://github.com/AlgoHunt/Face-Xray/blob/master/bi\\_online\\_generation.py](https://github.com/AlgoHunt/Face-Xray/blob/master/bi_online_generation.py). Accessed: 2023-08-10. 4
- [5] FaceApp. <https://faceappdownload.org>, 2017. 1
- [6] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: a compact facial video forgery detection network. In *2018 IEEE international workshop on information forensics and security (WIFS)*, pages 1–7. IEEE, 2018. 1, 2
- [7] Shruti Agarwal, Hany Farid, Yuming Gu, Mingming He, Koki Nagano, and Hao Li. Protecting world leaders against deep fakes. In *CVPR workshops*, page 38, 2019. 1, 2
- [8] Belhassen Bayar and Matthew C Stamm. Constrained convolutional neural networks: A new approach towards general purpose image manipulation detection. *IEEE Transactions on Information Forensics and Security*, 13(11):2691–2706, 2018. 4
- [9] 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*, pages 4113–4122, 2022. 6
- [10] Renwang Chen, Xuanhong Chen, Bingbing Ni, and Yanhao Ge. Simswap: An efficient framework for high fidelity face swapping. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2003–2011, 2020. 1, 2, 3
- [11] Xinru Chen, Chengbo Dong, Jiaqi Ji, Juan Cao, and Xirong Li. Image manipulation detection by multi-view multi-scale supervision. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14185–14193, 2021. 2
- [12] Yunjey Choi, Minje Choi, Munyoung Kim, Jung-Woo Ha, Sunghun Kim, and Jaegul Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8789–8797, 2018. 1
- [13] Yunjey Choi, Youngjung Uh, Jaejun Yoo, and Jung-Woo Ha. Stargan v2: Diverse image synthesis for multiple domains. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8188–8197, 2020. 1
- [14] Sowmen Das, Md Saiful Islam, and Md Ruhul Amin. Gcnet: utilizing gated context attention for improving image forgery localization and detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 81–90, 2022. 2
- [15] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsoia, and Stefanos Zafeiriou. Retinaface: Single-shot multi-level face localisation in the wild. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 5203–5212, 2020. 5
- [16] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) preview dataset. *arXiv preprint arXiv:1910.08854*, 2019. 2, 5
- [17] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397*, 2020. 2, 5
- [18] Pierre Foret, Ariel Kleiner, Hossein Mobahi, and Behnam Neyshabur. Sharpness-aware minimization for efficiently improving generalization. *arXiv preprint arXiv:2010.01412*, 2020. 6
- [19] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. 1, 2
- [20] David Güera and Edward J Delp. Deepfake video detection using recurrent neural networks. In *2018 15th IEEE international conference on advanced video and signal based surveillance (AVSS)*, pages 1–6. IEEE, 2018. 2
- [21] Xiao Guo, Xiaohong Liu, Zhiyuan Ren, Steven Grosz, Iacopo Masi, and Xiaoming Liu. Hierarchical fine-grained image forgery detection and localization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3155–3165, 2023. 2
- [22] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2
- [23] Xun Huang and Serge Belongie. Arbitrary style transfer in real-time with adaptive instance normalization. In *Proceedings of the IEEE international conference on computer vision*, pages 1501–1510, 2017. 1, 2
- [24] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. pmlr, 2015. 4
- [25] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *arXiv preprint arXiv:1710.10196*, 2017. 1
- [26] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4401–4410, 2019.
- [27] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8110–8119, 2020. 1

- [28] Jiseob Kim, Jihoon Lee, and Byoung-Tak Zhang. Smoothswap: a simple enhancement for face-swapping with smoothness. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10779–10788, 2022. 1, 2
- [29] Davis E King. Dlib-ml: A machine learning toolkit. *The Journal of Machine Learning Research*, 10:1755–1758, 2009. 3, 5
- [30] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 2
- [31] Marissa Koopman, Andrea Macarulla Rodriguez, and Zeno Geradts. Detection of deepfake video manipulation. In *The 20th Irish machine vision and image processing conference (IMVIP)*, pages 133–136, 2018. 1, 2
- [32] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457*, 2019. 1, 2, 5
- [33] 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*, pages 5001–5010, 2020. 1, 2, 3, 4, 6, 7
- [34] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. *arXiv preprint arXiv:1811.00656*, 2018. 1, 3
- [35] 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. 2, 5
- [36] 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 *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 772–781, 2021. 6, 7
- [37] 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. 2, 6, 7
- [38] Iacopo Masi, Aditya Killekar, Royston Marian Mascarenhas, Shenoy Pratik Gurudatt, and Wael AbdAlmageed. Two-branch recurrent network for isolating deepfakes in videos. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part VII 16*, pages 667–684. Springer, 2020. 2
- [39] Scott McCloskey and Michael Albright. Detecting gan-generated imagery using color cues. *arXiv preprint arXiv:1812.08247*, 2018. 1, 2
- [40] Huy H Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using capsule networks to detect forged images and videos. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2307–2311. IEEE, 2019. 1, 2
- [41] Jongchan Park, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Bam: Bottleneck attention module. *arXiv preprint arXiv:1807.06514*, 2018. 4
- [42] 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. 6
- [43] 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
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 2
- [45] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics: A large-scale video dataset for forgery detection in human faces. *arXiv preprint arXiv:1803.09179*, 2018. 2, 5
- [46] Ekraam Sabir, Jiaxin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. Recurrent convolutional strategies for face manipulation detection in videos. *Interfaces (GUI)*, 3(1):80–87, 2019. 1, 2
- [47] Kaede Shiohara and Toshihiko Yamasaki. Detecting deepfakes with self-blended images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18720–18729, 2022. 1, 2, 3, 6, 7
- [48] Yu Sun, Rongrong Ni, and Yao Zhao. Et: Edge-enhanced transformer for image splicing detection. *IEEE Signal Processing Letters*, 29:1232–1236, 2022. 2
- [49] YuYang Sun, ZhiYong Zhang, Isao Echizen, Huy H Nguyen, ChangZhen Qiu, and Lu Sun. Face forgery detection based on facial region displacement trajectory series. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 633–642, 2023. 2
- [50] Zekun Sun, Yujie Han, Zeyu Hua, Na Ruan, and Weijia Jia. Improving the efficiency and robustness of deepfakes detection through precise geometric features. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3609–3618, 2021. 1
- [51] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International conference on machine learning*, pages 6105–6114. PMLR, 2019. 4, 6
- [52] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Nießner. Face2face: Real-time face capture and reenactment of rgb videos. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2387–2395, 2016. 2, 5
- [53] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *Acm Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 2, 5
- [54] Yuhan Wang, Xu Chen, Junwei Zhu, Wenqing Chu, Ying Tai, Chengjie Wang, Jilin Li, Yongjian Wu, Feiyue Huang,

- and Rongrong Ji. Hiface: 3d shape and semantic prior guided high fidelity face swapping. *arXiv preprint arXiv:2106.09965*, 2021. 1, 2
- [55] Zhiliang Xu, Hang Zhou, Zhibin Hong, Ziwei Liu, Jiaming Liu, Zhizhi Guo, Junyu Han, Jingtuo Liu, Errui Ding, and Jingdong Wang. Styleswap: Style-based generator empowers robust face swapping. In *European Conference on Computer Vision*, pages 661–677. Springer, 2022. 1, 2
- [56] Zhiyuan Yan, Yong Zhang, Yanbo Fan, and Baoyuan Wu. Ucf: Uncovering common features for generalizable deepfake detection. *arXiv preprint arXiv:2304.13949*, 2023. 6, 7
- [57] Zhiyuan Yan, Yong Zhang, Xinhang Yuan, Siwei Lyu, and Baoyuan Wu. Deepfakebench: A comprehensive benchmark of deepfake detection. *arXiv preprint arXiv:2307.01426*, 2023. 7
- [58] 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. 1, 2
- [59] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2185–2194, 2021. 2, 6
- [60] 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*, pages 15023–15033, 2021. 1, 3, 6
- [61] Yinglin Zheng, Jianmin Bao, Dong Chen, Ming Zeng, and Fang Wen. Exploring temporal coherence for more general video face forgery detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15044–15054, 2021. 2, 6

# Generalized Deepfakes Detection with Reconstructed-Blended Images and Multi-scale Feature Reconstruction Network

## Supplementary Material

### 6. Effect of $\lambda$ in loss function

We experimentally evaluated the effect of scaling factors  $\lambda_1$  and  $\lambda_2$  in the loss function on model performance. Keeping the other experimental settings constant, we used several sets of values for hyper-parameters  $\lambda_1$  and  $\lambda_2$ . The results on CDF-v2, DFDC-P, and FF++ are shown in Table 7.

From the results in the table, it can be observed that the hyperparameter settings with  $\lambda_1 = 50$  and  $\lambda_2 = 100$  achieve the best performance.

### 7. Samples of RBIs

Several of the RBI samples we generated are shown in Figure 6, with the left half showing the results of statistical augmentation on the foreground face taken from the source image (i.e., the reconstructed image) and with the right half showing the results of statistical augmentation on the background face taken from the target image (i.e., the original image).

The reconstructed images reveal that our proposed disentanglement-reconstruction process not only introduces visually imperceptible frequency noise as a generator fingerprint (left half, rows 1, 3, and 5; right half, row 4) but also introduces unique visible artifacts created by the generator which cannot be simulated by statistical augmentation. These artifacts include abrupt cheek contours (left half, row 2, right half, row 5), inconsistent eye sizes (right half, row 2), overlapping eyes (left half, row 4), and blurred teeth artifacts (right half, row 5). Introducing pattern noise and distinctive generator artifacts can thus help the model learn more robust and generalizable forgery features, thereby improving the model’s detection performance on unseen manipulations and deepfake data.

### 8. More Visualization Results of Our Model

We have included additional examples in the appendix to provide a more comprehensive demonstration of the effectiveness of our method, as shown in Figure 7. It can be observed that our method is capable of accurately detecting specific manipulation regions. For instance, in the case of Deepfakes, a rectangular mask is employed to guide the replacement of the target face with the source face in the central facial region. Face2Face, conversely, employs RGB tracking to comprehensively capture the whole facial performance for the purpose of expression transfer. FaceShifter exhibits an adaptive capacity by autonomously generating manipulation masks and employing post-processing tech-

Settings		Test AUC (%)			
$\lambda_1$	$\lambda_2$	CDF-v2	DFDC-P	FF++	Avg.
25	25	92.84	77.52	99.83	90.06
25	50	93.96	80.51	99.89	91.45
50	50	94.51	83.85	99.81	92.72
50	100	95.27	83.66	99.92	92.95
100	50	95.17	81.78	99.91	92.29
150	300	94.28	74.50	99.87	89.55
500	1000	92.90	74.89	99.05	88.95

Table 7. Effect of  $\lambda$  in loss function on model performance.

niques to mitigate the influence of hair and accessories on the falsified outcomes.

Our method demonstrates effective localization capabilities across various manipulations. Even in instances of highly convincing manipulation outcomes, such as the FSH in the third column of Figure 7, the model, while expressing a lack of confidence, is still able to provide accurate assessments.



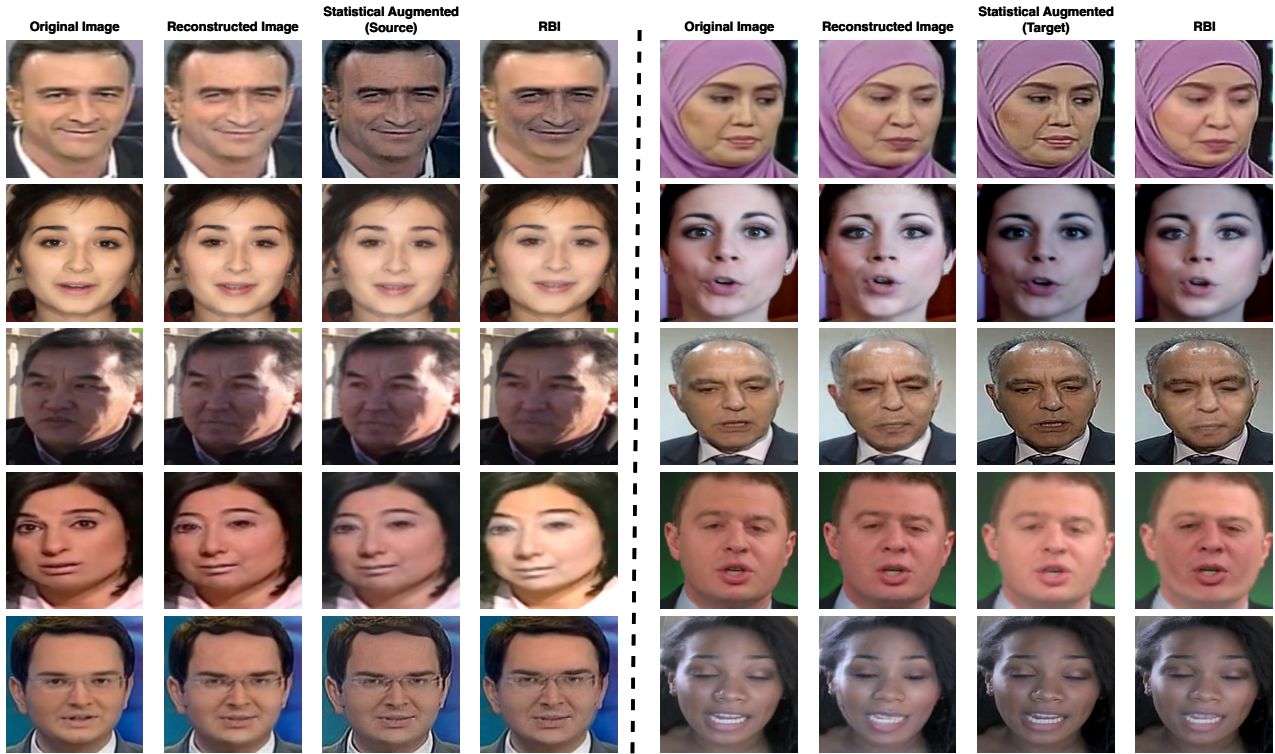


Figure 6. Examples of RBIs, where the left column represents statistical augmentations applied to the source image (i.e., reconstructed image), and the right column represents statistical augmentations applied to the target image (i.e., original image).

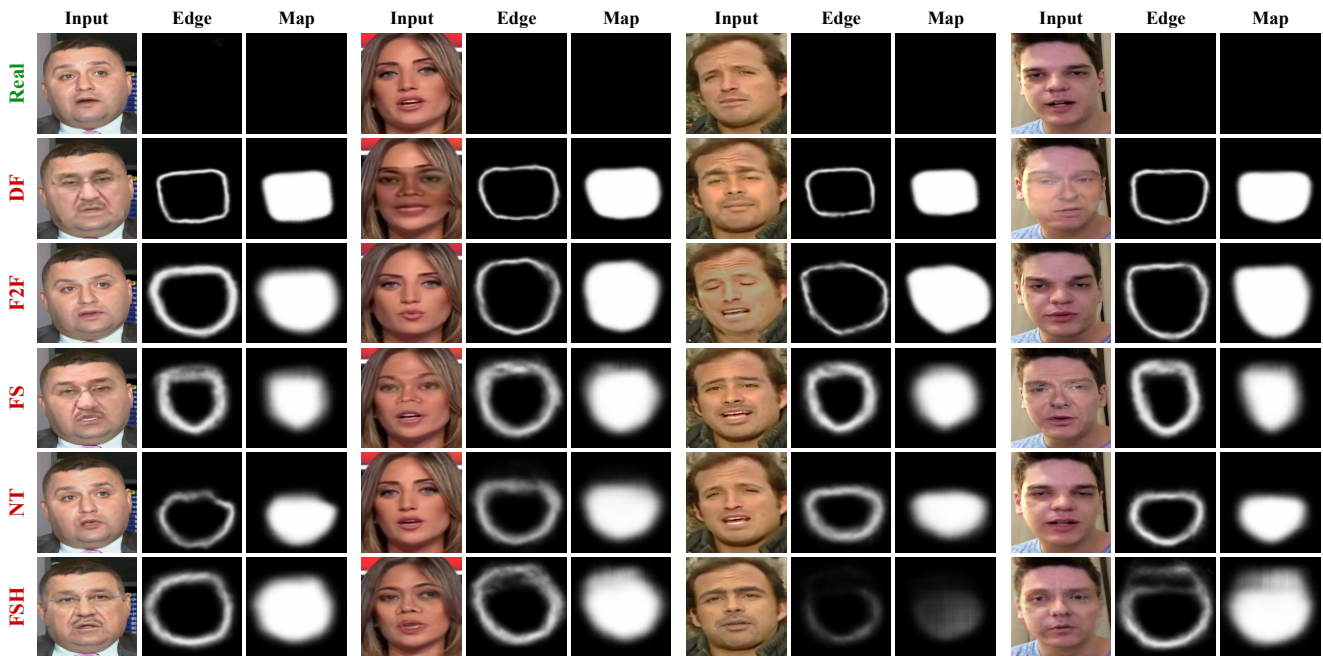


Figure 7. Additional visualization results of our method.