# Fake Visual Content Detection Using Two-Stream Convolutional Neural Networks

Bilal Yousaf, Muhammad Usama, Waqas Sultani, Arif Mahmood, Junaid Qadir

*Abstract*—**Rapid progress in adversarial learning has enabled the generation of realistic-looking fake visual content. To distinguish between fake and real visual content, several detection techniques have been proposed. The performance of most of these techniques however drops off significantly if the test and the training data are sampled from different distributions. This motivates efforts towards improving the generalization of fake detectors. Since current fake content generation techniques do not accurately model the frequency spectrum of the natural images, we observe that the frequency spectrum of the fake visual data contains discriminative characteristics that can be used to detect fake content. We also observe that the information captured in the frequency spectrum is different from that of the spatial domain. Using these insights, we propose to complement frequency and spatial domain features using a two-stream convolutional neural network architecture called TwoStreamNet. We demonstrate the improved generalization of the proposed two-stream network to several unseen generation architectures, datasets, and techniques. The proposed detector has demonstrated significant performance improvement compared to the current state-of-the-art fake content detectors and fusing the frequency and spatial domain streams has also improved generalization of the detector.**

*Index Terms*—**Deepfakes Two-stream network Freqency stream Combination of Discrete Fourier Transform and Discrete wavelent**

## I. INTRODUCTION

RECENT technological advancements in artificial intelligence (AI) have led to various beneficial applications in vision, language, and speech processing. However, at the same time, the power of these technologies may be exploited by adversaries for illegal or harmful uses. For example, Deepfakes—a portmanteau of the terms "deep learning" and "fake"—may be used to produce or alter photo-realistic audio-visual content with the help of deep learning for an illegal or harmful purpose. Deepfake technology enables one to effectively synthesize realistic-looking fake audio or video of a real person speaking and performing in any arbitrary way [1]. The term Deepfake was first coined by a Reddit community for synthetically replacing the face of a person with the face of another person. The term expanded with time to include similar techniques such as Lip-Sync [2], [3], facial expression reenactment [4]–[6], full-body and background manipulation as well as audio synthesis [7]–[12].

The rise of technology such as Deepfake has eroded the traditional confidence in the authenticity of audio and video as

B. Yousaf, W. Sultani, and A. Mahmood are with Department of Computer Science, Information Technology University, Lahore, Pakistan. E-mails: {msds18007, waqas.sultani, arif.mahmood}@itu.edu.pk

M. Usama and J. Qadir are with Department of Electrical Engineering, Information Technology University, Lahore, Pakistan. E-mails: {muhammad.usama, junaid.qadir}@itu.edu.pk

any digital content (audio, video, text) can be easily subverted using advanced deep learning techniques for synthesizing images trained on readily accessible public videos and images [13]–[16]. The gravity and urgency of the Deepfake threat can be gauged by noting that in recent times a CEO was scammed using Deepfake audio for \$243,000 [17] and a fake video of the president of Gabon has resulted in a failed coup attempt. Other potential effects of the Deepfake threat include danger to journalism and democratic norms because elections can be manipulated and democratic discourse may be disrupted by creating fake speeches of contending leaders [1], [3]. Unfortunately, most of the current research focuses on creating and improving Deepfakes and there is a lack of focus on reliable Deepfake detection. As reported in [16], 902 papers on Generative Adversarial Networks (GANs) were uploaded to the arXiv in 2018 but only 25 papers uploaded during the same time period related to the anti-forgery related topics.

Recent research shows that neural networks can be used for detecting fake content [18]–[22]. These methods however require a large amount of fake and real training data to accurately learn the data distributions of both classes. The performance of these methods drops significantly on the unseen fake data sampled from a different distribution or generation process as the underlying network may overfit the training data and thereby lose its ability to generalize. The model can be further trained to classify previously unseen data but it will require a large amount of data from the new distribution which may not always be available in such problems. Attackers and defenders are continuously improving their approaches and rolling out new attacks and defenses. Therefore, it may be very difficult to collect a large amount of fake data for new manipulation techniques. For such scenarios, a fake content detector can detect fake content without even being explicitly trained on it.

In the current work, we propose a two-stream network for fake visual content detection. The first stream called 'Spatial Stream' detects the fake data employing RGB images while the second stream dubbed as 'Frequency Stream' utilizes a combination of Discrete Fourier Transform (DFT) and Wavelet Transform (WT) for discriminating fake and real visual content. The frequency stream exploits the fact that the distribution of the frequency spectrum of the fake visual data remains distinct from the distribution of the real data frequency spectrum. This is illustrated in Figure 1, which shows the DFT-magnitude spectrum for a sample of real and fake images. It can be seen that frequency spectrum has patterns that are different from that of real images. These differences are used to classify the fake versus real content. Since the information captured by the frequency stream is different from
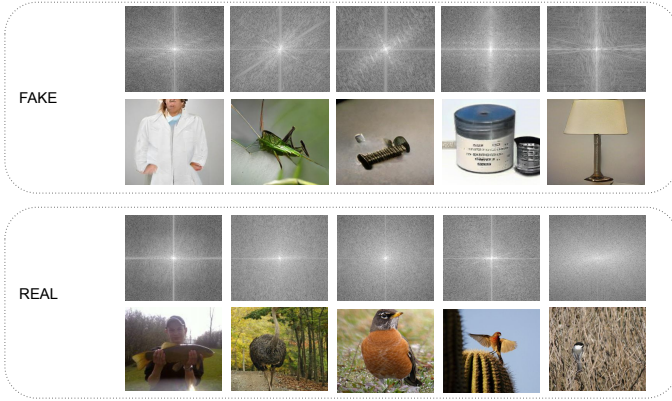
Fig. 1: DFT-magnitude spectrum for fake and real images has discriminative features which can be exploited for improved fake detection performance.

the information captured by the spatial stream, both these streams complement each other and fusing them together can provide better performance and generalization to unseen fake data detection. *To the best of our knowledge, this is the first work that studies the fusion of cross-modal information fusion to improve fake content detection generalization.*

The main contributions of this paper are summarized next.

1) A novel two-stream architecture for fake visual content detection consisting of a Spatial Stream (SS) and a Frequency Stream (FS) is proposed. The SS learns the difference between the distributions of real and fake visual content in the spatial space using RGB images, while the FS learns to discriminate between the distributions of real and fake content in the frequency domain. The coefficients of the stationary frequencies are captured using DFT, while the coefficients of spatially varying multi-scale frequencies are captured using Haar Wavelet transform. The spatial and frequency information complement each other therefore fusion of both has improved fake visual content detection.

2) The proposed two-stream network comprising of a frequency and a spatial domain stream has outperformed the state-of-the-art fake detection methods with a significant margin. A detailed analysis of the proposed approach is performed and we empirically demonstrate that the proposed approach is robust across different quality JPEG compression and blurriness artifacts.

In Section II, we discuss the related work and cover the traditional image forensics techniques and the latest deep learning-based image forensics algorithms with a prime focus on generalization. In Section III, we present our proposed methodology with all pre-processing schemes, training, and testing procedures. Section IV discusses the datasets used for evaluating and validating the proposed methodology. Section V provides the results and comprehensive evaluations of the generalization of the proposed methodology by performing an ablation study. Finally, Section VI concludes the paper and also points towards future directions.

## II. RELATED WORK

In this section, we briefly review recent works needed to understand the state-of-the-art solutions in image forensics. We have divided this section into four subsections. We begin with a brief overview of the hand-crafted image forensic techniques followed by deeply learned image forensic approaches. After that, we discuss methods that are focus to improve generalization. Finally, we conclude the section by covering the state-of-the-art frequency-domain techniques that are specifically designed for image forensic applications.

### A. Hand-Crafted Image Forensics

A variety of methods are available in the literature for detecting traditional image manipulation techniques. Most of these manipulations are designed with the help of image editing tools. The traditional techniques make use of the hand-crafted features to detect specific clues that are created as a result of different manipulations. For example, several blind noise estimation algorithms have been proposed to detect region splicing forgeries [23], [24]. Popescu et al. [25] detected the image forgeries by estimating the resampling in the images. Haodong et al. [26] integrated tampering possibility maps to improve forgery localization. Yuanfang et al. [27] identified potential artifacts in hue, saturation, dark and bright channels of fake colorized images, and developed detection methods based on histograms and feature encoding. Similarly, Peng et al. [28] used contact information of the standing objects and their supporting planes extracted from their reconstructed 3D poses to detect splicing forgeries. However, these techniques are unable to provide comparable performance to that of pixel-based methods in realistic situations. In recent works, learning-based techniques have become the preferred methods compared to traditional image forensics for achieving state-of-the-art detection performance [29]–[32].

### B. Deep Learning Based Image Forensics

Due to the success of deep learning in different fields, several researchers have recently leveraged deep learning approaches for fake visual content detection. YanYang et al. [33] proposed an algorithm based on difference images (DIs) and illuminant map (IM) as feature extractors to detect re-colorized images. Quan et al. [34] designed a deep CNN network with two cascaded convolutional layers to detect computer-generated images. McCloskey et al. [35] detected fake images by exploiting artifacts in the color cues, whereas Li et al. [36] used face warping artifacts for the forgery detection. Li et al. [37] noticed that eye blinking in fake videos is different than the natural videos and used this fact to expose the fake videos. Similarly, Yang et al. [38] have detected the Deepfakes by identifying the inconsistent head poses. Recently, Afchar et al. [22] proposed two compact forgery detection networks (Meso-4 and MesoInception-4) in which forgery detection is done by analyzing the mesoscopic properties[1] of Deepfake videos. Similarly, Nataraj et al. [39] have shown that features

---

[1]The eyes and mouth are determined as the mesoscopic features in the forgery detection in the Deepfake videos.

extracted from the co-occurrence matrix can help improving fake data detection and Wang et al. [40] proposed an anomaly detector based approach that uses pre-trained face detectors as a feature extractor. Although impressive, most of the above-mentioned approaches fail to perform well when fake visual data is sampled from a different distribution.

### C. Methods focused on Generalization

In this subsection, we briefly describe the fake detection approaches focused on generalization. Cozzolino et al. [41] proposed an auto-encoder based method to improve the performance of the model where learned weights are transferred for a different generation method. Zhang et al. [42] proposed a generalizable architecture named AutoGAN and evaluated its generalization ability on two types of generative networks. Xuan et al. [43] proposed that by using Gaussian Blur or Gaussian noise, one can destroy unstable low-level noise cues and force models to learn more intrinsic features to improve the generalization ability of the model. Similarly, Wang et al. [44] suggested that careful pre-and post-processing with data augmentation (such as blur and JPEG compression) improves the generalization ability. They have also shown improved fake detection results on multiple test sets by training on just one image generation network.

### D. Frequency Domain Methods

Gueguen et al. [45] extracted features from the frequency domain to perform classification tasks on images. Ehrlich et al. [46] proposed an algorithm to convert the convolutional neural network (CNN) models from the spatial domain to the frequency domain. Xu et al. [47] proposed learning in the frequency domain and have shown that the performance of object detection and segmentation tasks get improved in the frequency domain as compared to using spatial RGB domain. Durall et al. [48] have shown that fake images have a difference in high-frequency coefficients compared to the natural images which he used for fake detection. Wang et al. [44] have shown that the artifacts in the frequency spectrum of fake images can be detected. Zhang et al. [42] proposed that if instead of raw pixels, frequency spectrum (2D-DCT on all 3 channels) is used as an input to the fake image detector, the performance of the detector improves. These frequency response base detectors target specific properties of the image generation process therefore, their performance degrades when fake images from unseen distributions are tested. In contrast to these existing methods, the proposed algorithm fuses information from the spatial domain and the frequency domain to achieve improved generalization. Also, we propose to fuse DFT with Wavelet Transform to improve the discrimination in the frequency domain. These innovations have resulted in significant improvement in fake content detection compared to the existing methods.

### III. METHODOLOGY

Improving the generalizability of a fake detection model is critical for its success in real-world applications where the fake content may be generated by unknown processes. We propose a generalizable fake detection model based on a two-stream convolutional network architecture shown in Fig. 2.

The proposed architecture is motivated by the excellent performance of two-stream networks in action recognition in videos. To the best of our knowledge, The proposed network performs quite well on both seen and unseen data and has outperformed existing state-of-the-art (SOTA) methods in a wide range of experiments as we shall discuss in later sections. Our proposed two-stream network is novel and such a combination of frequency stream and the spatial stream has not been proposed before. In the following, we discuss the RGB to YCbCr conversion, DFT, DWT, and the proposed architecture in more detail.

*1) The RGB to YCbCr Transformation:* The three channels in RGB color space are correlated with each other. We consider an orthogonal color space for improved representation. In our experiments, we have used YCbCr that has performed better than RGB space. As recommended in previous research [49] [50], the following formulas are used to convert from RGB to YCbCr color space:

$$Y = K_{ry}.R + K_{gy}.G + K_{by}.B,$$
$$Cr = B - Y, \quad Cb = R - Y, \quad (1)$$
$$K_{ry} + K_{gy} + K_{by} = 1,$$

where, $K_{ry}, K_{gy}$, and $K_{by}$ are the coefficients for color conversion whose values are specified in Table I according to the standards. In our implementation, we used ITU601 [51].

| Reference Standard | $K_{ry}$ | $K_{by}$ |
|---|---|---|
| [51] ITU601 / ITU-T 709 1250/50/2:1 | 0.299 | 0.114 |
| [52] ITU709 / ITU-T 709 1250/60/2:1 | 0.2126 | 0.0722 |
| [53] SMPTE 240M (1999) | 0.212 | 0.087 |

TABLE I: Coefficients $K_{ry}$ and $K_{by}$ of color conversion from RGB to YCbCr.

*2) A Review of Frequency Domain Transforms:* To fully capture the frequency information from a YCbCr image, we compute DFT and DWT for each image.

*Discrete Fourier Transform (DFT)*: Using DFT, one can decompose a signal into sinusoidal components of various frequencies ranging from 0 to maximum value possible based on the spatial resolution. For two dimensional data, i.e., images of size $W \times H$, the DFT can be computed using the following formula:

$$X_{w,h} = \sum_{n=0}^{W-1} \sum_{m=0}^{H-1} x_{w,h} e^{\frac{-i2\pi}{N}wn} e^{\frac{-i2\pi}{M}hm}, \quad (2)$$

where $w$ is the horizontal spatial frequency, $h$ is the vertical spatial frequency, $x_{w,h}$ is the pixel value at coordinates (w, h), and $X_{w,h}$ carries the magnitude and phase information of frequency at coordinates $(w, h)$.

*Discrete Wavelet Transform (DWT)*: Wavelet transform decomposes an image into four different subband images. High and low pass filters are applied at each row (column) and then they are downsampled by 2 to get the high and low-frequency components of each row (column) separately. In

Fig. 2: Proposed two-stream convolutional neural network (TwoStreamNet). The two network streams capture spatial and frequency domain artifacts separately, and their outputs are fused at the end of the network to produce classification scores.

this way, the original image is converted into four sub-band images: High-high (HH), High-low (HL), Low-high (LH), and Low-low (LL). Each subband image preserves different features: HH region preserves high-frequency components in both horizontal and vertical direction, HL preserves high-frequency components in the horizontal direction and low-frequency components in the vertical direction, LH preserves low-frequency components in the vertical direction and high-frequency components in the horizontal direction and finally, LL preserves low-frequency components in the vertical direction and low-frequency components in the horizontal direction.

### A. Frequency Stream

In this stream, two different types of the frequency spectrum are fused to get improved frequency domain representation which can better discriminate between the real and the fake visual content. An overview of the frequency spectrum fusion is shown in Fig. 3.

The three YCbCr channels are then transformed to the frequency domain using two different types of transformations, including DFT and DWT. Each channel is divided into non-overlapping block of size $8 \times 8$ pixels and transformation is applied on each block independently. The resulting coefficients are then concatenated back to obtain the arrays of original image size. The output of the DFT converts one input channel

into two output channels corresponding to real and imaginary coefficients. Similarly, the DWT converts one input channel into 4 output channels corresponding to low frequencies (LL), high and low frequencies (HL), high frequencies (HH), and low and high frequencies (LH). For three input channels (YCbCr), we obtain 18 output channels, 6 from DFT and 12 from DWT. All of these frequency output channels are concatenated to form 3D cubes of size $H \times W \times C$, where $H$ is the height and $W$ is the width of the image, and $C$=18 are the number of channels. We empirically observe that both DFT and DWT are necessary to capture essential information in the frequency domain at varying scales for improving the generalization ability of the proposed network.

### B. Spatial Stream

In this stream, RGB channels of the image are passed as input to the ResNet50 [54] as the classifier. RGB images are augmented in a special way using JPEG compression and Gaussian Blur as recommended by Wang et al. [44]. This stream is trained individually and plugged in the TwoStream-Net at the test time.

### C. Two Stream Network Architecture

The proposed two-stream network architecture is shown in Figure 2. ResNet50 network is used as a backbone in
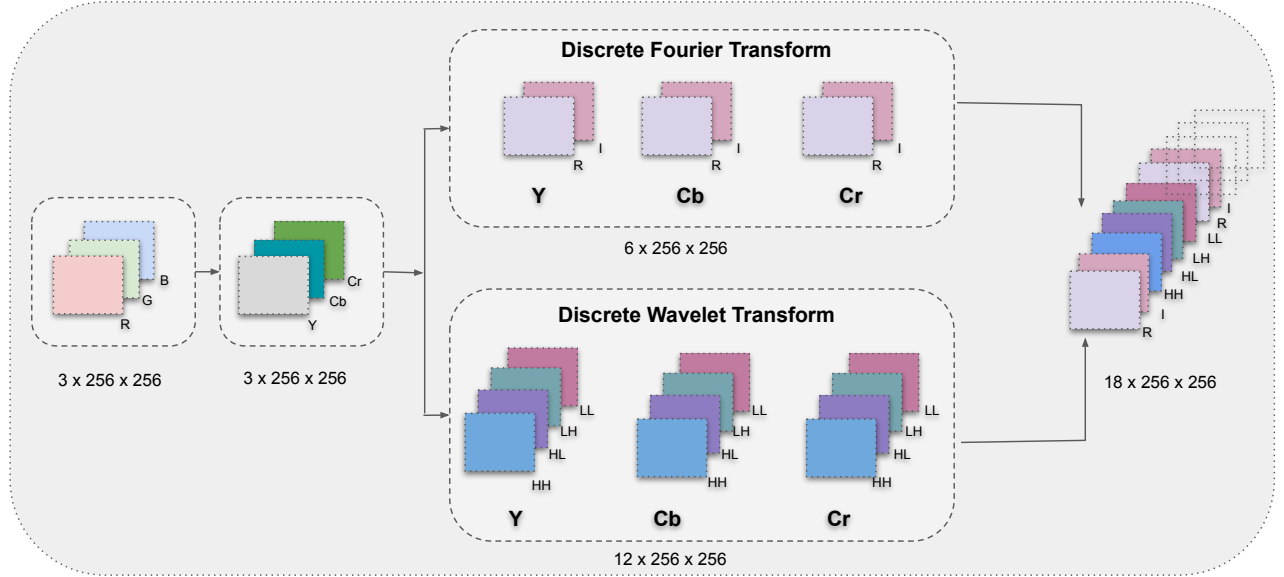
Fig. 3: Proposed pre-processing pipeline: the input image is first converted to YCbCr color space and then transformed to the frequency domain by applying DFT and Wavelet Transforms (WT). After DFT, we get real (R) and imaginary (I) channels, and after WT we get four channels: HH, HL, LH, and LL. The resulting channels are concatenated to form 3D cubes which are then input to the frequency stream for further processing.

both of the streams of the proposed architecture. Since the number of input channels in the frequency stream is larger as compared to the spatial stream, therefore first layer of FS is accordingly modified. Both streams are independently trained and the output of both streams is fused using the class probability averaging fusion method. In this fusion scheme, both streams contribute equally to the output, to produce the final classification probability. The performance of the combined scores is significantly better than the performance of the individual streams.

## IV. EXPERIMENTS AND RESULTS

**Training Dataset:** Following the protocol used by [44], the proposed two-stream network is trained using the fake images generated by ProGAN [55] and tested on the images generated by many other GANs. ProGAN has 20 different officially released models trained on different object categories of the LSUN dataset, which is a large scale image dataset containing around one million labeled images for each of the 10 scene categories and 20 object categories[2] [56]. We choose 15 (airplane, bird, boat, bottle, bus, car, cat, chair, dog, horse, motorbike, person, sofa, train, and tv monitor) out of 20 models to create our validation and training set. We generated 10k fake images for training and 500 fake images for validation using each of the 15 models. For each of these 15 categories of fake images, we collect 10k of real images for training and 500 for validation randomly from the LSUN dataset [56]. In total, we have 300K training images and 15K validation images. For real images, we center crop the images equal to the size of the shorter edge and then resize the images to $256 \times 256$.

**Testing Dataset:** Testing dataset contains images which were generated using completely unseen generators as described in Table III. To remain consistent with the current state of the art, the same generators are selected as that of [44]. The real images for testing purposes are obtained from the repository for each generator.

### A. Implementation details

For training the FS, we use the Adam [69] optimizer with an initial learning rate of 0.0001, weight decay of 0.0005, and a batch size of 24. For all the training sets, we train the proposed network for 24 epochs. Large training data has helped the model to converge quickly. Lastly, we select the best model based on the validation set. While training each stream, data augmentation based on Gaussian blur and JPEG compression with 10% probability is used.

### B. Comparison with the Existing State-of-the-Art Algorithms

We thoroughly evaluated the performance of the proposed method on the test dataset and compared it with the existing state of the art [44]. We also compared the robustness analysis of our approach against some common real-world perturbations. In Table II, we have shown a comparison of our results with the best results of Wang et al. [44] ([Blur+JPEG(0.1)]). Their results from their official web link[3]. Results show that our FS approach performs very well on the unseen manipulations and outperformed the state-of-the-art on several test sets while having competitive performance on the remaining. Results of the two-stream architecture demonstrate

| Metrics | Method | Star GAN | Style GAN | SITD | Big GAN | Style GAN2 | Cycle GAN | Which face is real | SAN | Deep fake | Gua GAN | CRN | IMLE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Accuracy | Wang et al. [44] | 91.7 | 87.1 | 90.3 | 70.2 | 84.4 | **85.2** | 83.6 | **53.5** | 50.5 | 78.9 | **86.3** | **86.2** |
|  | Frequency Stream (Ours) | 97.67 | 89.36 | 81.11 | **72.08** | 91.83 | 77.53 | 80.15 | 49.32 | **68.90** | 70.91 | 55.06 | 55.06 |
|  | Two Stream (Ours) | **96.32** | **88.90** | 97.22 | 72.85 | 87.43 | 84.09 | **87.50** | 50.23 | **55.00** | 79.64 | 77.75 | 77.75 |
| F1-Score (Fake) | Wang et al. [44] | 91.31 | 85.19 | 89.91 | 61.13 | 81.53 | 84.2 | 81.92 | 3.56 | 12.83 | 75.45 | **87.93** | **87.88** |
|  | Frequency Stream (Ours) | **97.63** | **88.21** | 83.57 | **71.71** | 91.14 | 78.75 | 82.08 | **28.39** | **61.13** | 72.65 | 68.99 | 68.99 |
|  | Two Stream (Ours) | **96.21** | **87.52** | 97.27 | 66.61 | 85.63 | **84.98** | 87.15 | 4.39 | **17.95** | 77.34 | 81.79 | 81.79 |
| F1-Score (Real) | Wang et al. [44] | 92.14 | 88.56 | 90.62 | 75.81 | 86.50 | 86.08 | 85.0 | **66.67** | 68.28 | 81.5 | **84.13** | **84.08** |
|  | Frequency Stream (Ours) | **97.71** | **90.32** | 77.78 | 72.43 | **92.42** | 76.97 | 77.76 | 60.78 | **74.08** | 68.93 | 18.41 | 18.41 |
|  | Two Stream (Ours) | **96.43** | **90.01** | 97.18 | **77.13** | 88.84 | **86.20** | 87.83 | 66.36 | **69.00** | 81.5 | 71.39 | 71.39 |

TABLE II: Comparison of the proposed Frequency Stream (FS) and Two-Stream network with the state-of-the-art method [44] using average accuracy. Best results of Wang et al. [44] with data augmentation using blur and JPEG (0.1) are reported where 0.1 mean JPEG compression is applied on 10% images. The same augmentation is also used in the proposed approaches. Both our approach and that of Wang et al. are trained using ProGAN only and tested on the data generated by 12 unseen generation processes mentioned in the top row.
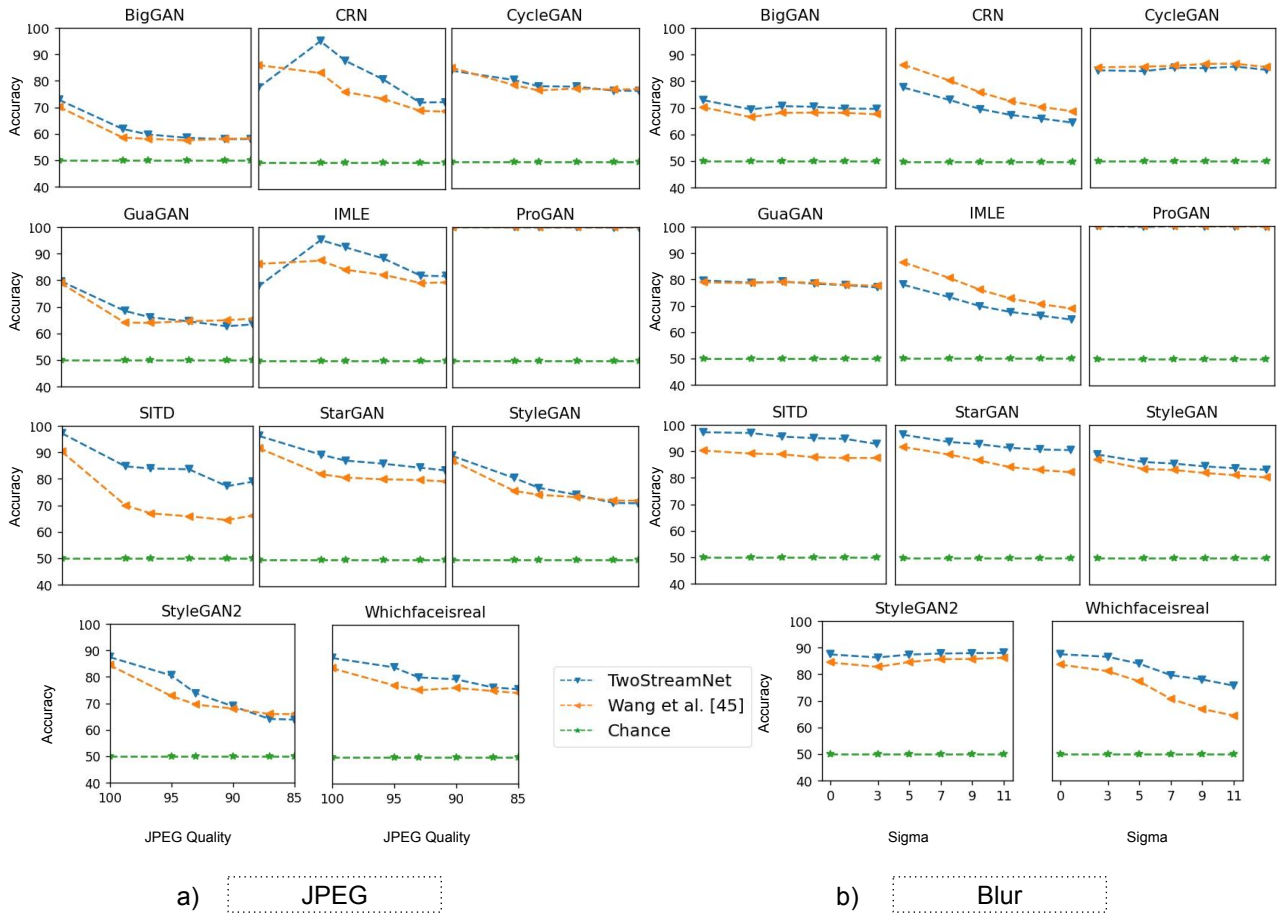


Fig. 4: Robustness comparison of the proposed algorithm with Wang et al. [44] for Gaussian Blur and JPEG Compression artifacts. In most experiments, the proposed two-stream net. We apply Gaussian Blur and JPEG Compression of different sizes on the test sets and measure the effect on the accuracy of our model. Our model performs near to the best for all the cross-modal datasets even when a large blurring effect is applied. Results show that our proposed solution is more robust as compared to the state of the art.

that our complete approach outperformed the state-of-the-art in almost all of the cases. Analysis of the results shows that when both spatial and frequency streams are combined into a two-stream architecture, they compliment each other in a way that their combined accuracy is greater than any of them individually. This clearly shows that FS ConvNet has learned distinctive features that were not learned by SS ConvNet. Overall, the combination of FS and SS plays a vital role in improving the generalization ability of the fake image detectors. In Figure 5(A), we have shown samples of the

Fig. 5: (A) Examples of fake images correctly detected by the proposed two-stream network, however, misclassified by Wang et al. [44]; (B) Examples of fake image misclassified by both our proposed method and that of Wang et al. [44].

fake images which are misclassified by the state-of-the-art and are correctly classified by our proposed two-stream approach.

These results demonstrate the ability of the proposed approach to detect high-quality fake images which are even very hard to

| Dataset | No. of Real Images | No. of Fake Images |
|---|---|---|
| StarGAN [57] | 1999 | 1999 |
| StyleGAN [58] | 5991 | 5991 |
| SITD [59] | 180 | 180 |
| BigGAN [60] | 2000 | 2000 |
| StyleGAN2 [61] | 7988 | 7988 |
| CycleGAN [62] | 1321 | 1321 |
| Whichfaceisreal [63] | 1000 | 1000 |
| GauGAN [64] | 5000 | 5000 |
| Deepfake [65] | 2698 | 2707 |
| CRN [66] | 6382 | 6382 |
| IMLE [67] | 6382 | 6382 |
| SAN [68] | 219 | 219 |

TABLE III: Details of the testing dataset
.

| Dataset | DFT | DWT | DFT + DWT |
|---|---|---|---|
| StarGAN [57] | 78.81% | 79.34% | 97.67% |
| StyleGAN [58] | 61.45% | 69.11% | 87.1% |
| SITD [59] | 83.61% | 49.72% | 90.3% |
| BigGAN [60] | 64.92% | 66.92% | 72.08% |
| StyleGAN2 [61] | 68.02% | 59.54% | 91.83% |
| CycleGAN [62] | 65.31% | 55.94% | 77.53% |
| Whichfaceisreal [63] | 39.85% | 77.40% | 80.15% |
| GauGAN [64] | 50.91% | 67.22% | 70.91% |
| Deepfake [65] | 56.39% | 51.88% | 68.90% |
| CRN [66] | 61.75% | 83.85% | 55.06% |
| IMLE [67] | 47.82% | 79.90% | 55.06% |
| SAN [68] | 69.18% | 46.80% | 49.32% |

TABLE IV: Evaluation of DFT and DWT combination for fake image detection. Percentage accuracy is reported for the full image using only DFT, only DWT, and the combination DFT+DWT.

discriminate by humans. Figure 5(B) shows the fake images which are misclassified by both Wang et al. [44] and us.

**Robustness Analysis:** In real-world settings, fake images may undergo several post-processing operations like compression, smoothness, etc. Therefore, we have evaluated the performance of the proposed model on the images which are post processed using JPEG compression and Gaussian blur. Specifically, we apply Gaussian blur with different standard deviations including [3, 5, 7, 9, 11] and JPEG compression with JPEG image quality factor of [85, 87, 90, 92, 95]. Results in Figure 4 show that our approach is robust to common perturbations. For most of the models, the proposed approach significantly outperformed the state-of-the-art method at varying blur levels. Similarly, the proposed approach also performed better than the state-of-the-art methods for a wide range of JPEG compressions.

## V. ABLATION STUDY

In this section, we thoroughly validate the different components of the proposed approach by performing an ablation study.

### A. Combining DFT and DWT

As shown in Figure 2, we propose to combine DWT and DFT for better feature representation and robust fake content detection. To verify the effectiveness of using both transformations, while keeping all the experimental settings the same, we experimented with DFT and DWT separately. After training for 20 epochs the best epoch is chosen based on validation data accuracy. The results shown in Table IV demonstrate that a combination of DFT and DWT is essential to produce robust feature representation for fake image detection.

### B. The Effect of Block Size

We study the effect of using different block sizes instead of computing DFT over the whole image. In Table V, we have shown results of computing DFT on the block size of $8 \times 8$, $16 \times 16$, $32 \times 32$, and $256 \times 256$ (Full-Image size). Note that, block size experiments are performed by keeping exactly the same experimental settings. Results demonstrate that $8 \times 8$ block size has consistently outperformed other block sizes. Therefore, transforming the image to the frequency domain using $8 \times 8$ blocks for DFT is more effective for fake image detection.

### C. The Effect of Color-Space

We evaluate the effectiveness of converting images into YCbCr color space before frequency transformations. We performed two experiments using the same settings to compare the performance of RGB with YCbCr color space. Results in Table VI show that converting an image to YCbCr colorspace adds more discriminative features in the frequency domain and helps in better fake image detection.

| Dataset | 8x8 | 16x16 | 32x32 | 'Full-Image' |
|---|---|---|---|---|
| StarGAN [57] | **94.62%** | 77.74% | 51.05% | 78.81% |
| StyleGAN [58] | **90.65%** | 68.93% | 46.47% | 61.45% |
| SITD [59] | **86.1%** | 55.28% | 75.56% | 83.61% |
| BigGAN [60] | **69.42%** | 63.12% | 48.25% | 64.92% |
| StyleGAN2 [61] | **92.23%** | 65.26% | 53.26% | 68.02% |
| CycleGAN [62] | **79.96%** | 67.97% | 36.68% | 65.31% |
| Whichfaceisreal [63] | **81.50%** | 63.25% | 41.95% | 39.85% |
| GauGAN [64] | **67.11%** | 54.74% | 57.95% | 50.91% |
| Deepfake [65] | **60.61%** | 51.90% | 57.35% | 56.39% |
| CRN [66] | 51.43% | 50.44% | 57.90% | **61.75%** |
| IMLE [67] | 51.50% | 50.53% | **66.28%** | 47.82% |
| SAN [68] | 46.58% | 51.83% | **72.15%** | 69.18% |

TABLE V: Fake image detection accuracy variation by varying block sizes for DFT transform. The block size $8 \times 8$ has produced best results therefore in our experiments this block size is used.

| Dataset | RGB | YCbCr |
|---|---|---|
| StarGAN [57] | 66.83% | **78.81%** |
| StyleGAN [58] | 59.32% | **61.45%** |
| SITD [59] | **87.50%** | 83.61% |
| BigGAN [60] | **72.97%** | 64.92% |
| StyleGAN2 [61] | 60.42% | **68.02%** |
| CycleGAN [62] | **75.89%** | 65.31% |
| Whichfaceisreal [63] | **47.80%** | 39.85% |
| GauGAN [64] | **60.92%** | 50.91% |
| Deepfake [65] | 56.02% | **56.32%** |
| CRN [66] | 58.83% | **61.75%** |
| IMLE [67] | **48.24%** | 47.82% |
| SAN [68] | 66.89% | **69.18%** |

TABLE VI: The compassion of fake detection performance using RGB and YCbCr Colorspace. YCbCr color space has performed better than RGB color space.

## VI. Conclusions

This paper addresses the problem of fake image detection. For this purpose, a two-stream network is proposed consisting of a spatial stream and a frequency stream. The proposed method generalizes to unseen fake image generator distributions much better than the current state-of-the-art approaches. The proposed method is also found to be more robust to the common image perturbations including blur and JPEG compression artifacts. The improved performance is leveraged by combining two types of frequency domain transformations, namely, Discrete Fourier Transform (DFT) and Discrete Wavelet Transform (DWT). Both transformations are applied upon YCbCr color-space and different frequency domain channels are concatenated to discriminate fake images from the real ones. By exploiting the differences between the real and the fake image frequency responses, improved fake detection performance is achieved. In the future, we aim to extend this work for fake video and audio detection.

## References

[1] Bobby Chesney and Danielle Citron. Deep fakes: A looming challenge for privacy, democracy, and national security. *Calif. L. Rev.*, 107:1753, 2019.

[2] Rithesh Kumar, Jose Sotelo, Kundan Kumar, Alexandre de Brébisson, and Yoshua Bengio. Obamanet: Photo-realistic lip-sync from text. *arXiv preprint arXiv:1801.01442*, 2017.

[3] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing Obama: learning lip sync from audio. *ACM Transactions on Graphics (TOG)*, 36(4):95, 2017.

[4] 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.

[5] Justus Thies, Michael Zollhöfer, Matthias Nießner, Levi Valgaerts, Marc Stamminger, and Christian Theobalt. Real-time expression transfer for facial reenactment. *ACM Trans. Graph.*, 34(6):183–1, 2015.

[6] Olivia Wiles, A Sophia Koepke, and Andrew Zisserman. X2face: A network for controlling face generation using images, audio, and pose codes. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 670–686, 2018.

[7] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. Everybody dance now. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5933–5942, 2019.

[8] Haoye Cai, Chunyan Bai, Yu-Wing Tai, and Chi-Keung Tang. Deep video generation, prediction and completion of human action sequences. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 366–382, 2018.

[9] Patrick Esser, Johannes Haux, Timo Milbich, et al. Towards learning a realistic rendering of human behavior. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018.

[10] Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, et al. Tacotron: Towards end-to-end speech synthesis. *arXiv preprint arXiv:1703.10135*, 2017.

[11] Sercan Arik, Jitong Chen, Kainan Peng, Wei Ping, and Yanqi Zhou. Neural voice cloning with a few samples. In *Advances in Neural Information Processing Systems*, pages 10019–10029, 2018.

[12] Hideyuki Tachibana, Katsuya Uenoyama, and Shunsuke Aihara. Efficiently trainable text-to-speech system based on deep convolutional networks with guided attention. In *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 4784–4788. IEEE, 2018.

[13] From porn to 'game of thrones': How deepfakes and realistic-looking fake videos hit it big. https://www.businessinsider.com/deepfakes-explained-the-rise-of-fake-realistic-videos-online-2019-6, 2018. (Accessed on 07/12/2020).

[14] Dave Lee. Fake porn' has serious consequences. https://www.bbc.com/news/technology-42912529, 2018. (Accessed on 07/12/2020).

[15] Samantha Cole. Gfycat's AI solution for fighting deepfakes isn't working. https://www.vice.com/en_us/article/ywe4qw/gfycat-spotting-deepfakes-fake-ai-porn, 2018. (Accessed on 07/12/2020).

[16] Giorgio Patrini, Francesco Cavalli, and Henry Ajde. The state of deepfakes: reality under attack, annual report v.2.3. https://deeptracelabs.com/archive/, 2018. (Accessed on 07/12/2020).

[17] J Damiani. A voice deepfake was used to scam a CEO out of 243,000$. https://www.forbes.com/sites/jessedamiani/2019/09/03/a-voice-deepfake-was-used-toscam-a-ceo-out-of-243000/, 2019. (Accessed on 07/12/2020).

[18] Belhassen Bayar and Matthew C Stamm. A deep learning approach to universal image manipulation detection using a new convolutional layer. In *Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security*, pages 5–10, 2016.

[19] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection. In *Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security*, pages 159–164, 2017.

[20] Nicolas Rahmouni, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Distinguishing computer graphics from natural images using convolution neural networks. In *2017 IEEE Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2017.

[21] 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.

[22] 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.

[23] Shruti Agarwal and Hany Farid. Photo forensics from jpeg dimples. In *2017 IEEE Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2017.

[24] Siwei Lyu, Xunyu Pan, and Xing Zhang. Exposing region splicing forgeries with blind local noise estimation. *International journal of computer vision*, 110(2):202–221, 2014.

[25] Alin C Popescu and Hany Farid. Exposing digital forgeries by detecting traces of resampling. *IEEE Transactions on signal processing*, 53(2):758–767, 2005.

[26] H. Li, W. Luo, X. Qiu, and J. Huang. Image forgery localization via integrating tampering possibility maps. *IEEE Transactions on Information Forensics and Security*, 12(5):1240–1252, 2017.

[27] Y. Guo, X. Cao, W. Zhang, and R. Wang. Fake colorized image detection. *IEEE Transactions on Information Forensics and Security*, 13(8):1932–1944, 2018.

[28] B. Peng, W. Wang, J. Dong, and T. Tan. Image forensics based on planar contact constraints of 3d objects. *IEEE Transactions on Information Forensics and Security*, 13(2):377–392, 2018.

[29] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2921–2929, 2016.

[30] Minyoung Huh, Andrew Liu, Andrew Owens, and Alexei A Efros. Fighting fake news: Image splice detection via learned self-consistency. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 101–117, 2018.

[31] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. Splicebuster: A new blind image splicing detector. In *2015 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2015.

[32] Yuan Rao and Jiangqun Ni. A deep learning approach to detection of splicing and copy-move forgeries in images. In *2016 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2016.

[33] Y. Yan, W. Ren, and X. Cao. Recolored image detection via a deep discriminative model. *IEEE Transactions on Information Forensics and Security*, 14(1):5–17, 2019.

[34] W. Quan, K. Wang, D. Yan, and X. Zhang. Distinguishing between natural and computer-generated images using convolutional neural networks. *IEEE Transactions on Information Forensics and Security*, 13(11):2772–2787, 2018.

[35] Scott McCloskey and Michael Albright. Detecting gan-generated imagery using color cues. *arXiv preprint arXiv:1812.08247*, 2018.

[36] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. *arXiv preprint arXiv:1811.00656*, 2018.

[37] Yuezun Li, Ming-Ching Chang, and Siwei Lyu. In Ictu Oculi: Exposing AI Created Fake Videos by Detecting Eye Blinking. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7. IEEE, 2018.

[38] 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.

[39] Lakshmanan Nataraj, Tajuddin Manhar Mohammed, BS Manjunath, Shivkumar Chandrasekaran, Arjuna Flenner, Jawadul H Bappy, and Amit K Roy-Chowdhury. Detecting GAN generated fake images using co-occurrence matrices. *Electronic Imaging*, 2019(5):532–1, 2019.

[40] Run Wang, Lei Ma, Felix Juefei-Xu, Xiaofei Xie, Jian Wang, and Yang Liu. Fakespotter: A simple baseline for spotting ai-synthesized fake faces. *arXiv preprint arXiv:1909.06122*, 2019.

[41] Davide Cozzolino, Justus Thies, Andreas Rössler, Christian Riess, Matthias Nießner, and Luisa Verdoliva. Forensictransfer: Weakly-supervised domain adaptation for forgery detection. *arXiv preprint arXiv:1812.02510*, 2018.

[42] Xu Zhang, Svebor Karaman, and Shih-Fu Chang. Detecting and simulating artifacts in GAN fake images. *arXiv preprint arXiv:1907.06515*, 2019.

[43] Xinsheng Xuan, Bo Peng, Wei Wang, and Jing Dong. On the generalization of GAN image forensics. In *Chinese Conference on Biometric Recognition*, pages 134–141. Springer, 2019.

[44] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 7, 2020.

[45] Lionel Gueguen, Alex Sergeev, Ben Kadlec, Rosanne Liu, and Jason Yosinski. Faster neural networks straight from JPEG. In *Advances in Neural Information Processing Systems*, pages 3933–3944, 2018.

[46] Max Ehrlich and Larry S Davis. Deep residual learning in the jpeg transform domain. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3484–3493, 2019.

[47] Kai Xu, Minghai Qin, Fei Sun, Yuhao Wang, Yen-kuang Chen, and Fengbo Ren. Learning in the frequency domain. *arXiv preprint arXiv:2002.12416*, 2020.

[48] Ricard Durall, Margret Keuper, Franz-Josef Pfreundt, and Janis Keuper. Unmasking deepfakes with simple features. *arXiv preprint arXiv:1911.00686*, 2019.

[49] Recommendation ITU-R studio encoding parameters of digital television for standard 4: 3 and wide-screen 16: 9 aspect ratios. 2011.

[50] Radiocommunication ITU. Parameter values for the HDTV standards for production and international programme exchange. *Recommendation ITU-R BT. 709-5*, 2002.

[51] Recommendation ITU-R BT.601-5, 1982-1995.

[52] Recommendation ITU-R BT.709-5, 1990-2002.

[53] Society of Motion Picture and Television Engineers (http://www.smpte.org). SMPTE 240M-1999 "Television - Signal Parameters - 1125-Line High-Definition Production".

[54] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

[55] 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.

[56] Fisher Yu, Ari Seff, Yinda Zhang, Shuran Song, Thomas Funkhouser, and Jianxiong Xiao. LSUN: Construction of a large-scale image dataset using deep learning with humans in the loop. *arXiv preprint arXiv:1506.03365*, 2015.

[57] 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.

[58] Tero Karras, Samuli Laine, and Timo Aila. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4401–4410, 2019.

[59] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3291–3300, 2018.

[60] Andrew Brock, Jeff Donahue, and Karen Simonyan. Large scale GAN training for high fidelity natural image synthesis. In *International Conference on Learning Representations*, 2018.

[61] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. *arXiv preprint arXiv:1912.04958*, 2019.

[62] Jun-Yan Zhu, Taesung Park, Phillip Isola, and Alexei A Efros. Unpaired image-to-image translation using cycle-consistent adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2223–2232, 2017.

[63] Which face is real? https://www.whichfaceisreal.com/. https://www.whichfaceisreal.com/. (Accessed on 07/12/2020).

[64] Taesung Park, Ming-Yu Liu, Ting-Chun Wang, and Jun-Yan Zhu. Semantic image synthesis with spatially-adaptive normalization. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.

[65] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. FaceForensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1–11, 2019.

[66] Qifeng Chen and Vladlen Koltun. Photographic image synthesis with cascaded refinement networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1511–1520, 2017.

[67] Ke Li, Tianhao Zhang, and Jitendra Malik. Diverse image synthesis from semantic layouts via conditional IMLE. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4220–4229, 2019.

[68] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11065–11074, 2019.

[69] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.