

Thinking in Frequency: Face Forgery Detection by Mining Frequency-aware Clues

Yuyang Qian^{*,†1,2}, Guojun Yin^{†,‡1}, Lu Sheng^{‡3}, Zixuan Chen^{*1,4}, and Jing Shao¹

¹ SenseTime Research,

² University of Electronic Science and Technology of China,

³ College of Software, Beihang University,

⁴ Northwestern Polytechnical University

qyy@std.uestc.edu.cn, zixuan.sean.chen@hotmail.com, lsheng@buaa.edu.cn,
{yinguojun, shaojing}@sensetime.com,

Abstract. As realistic facial manipulation technologies have achieved remarkable progress, social concerns about potential malicious abuse of these technologies bring out an emerging research topic of face forgery detection. However, it is extremely challenging since recent advances are able to forge faces beyond the perception ability of human eyes, especially in compressed images and videos. We find that mining forgery patterns with the awareness of frequency could be a cure, as frequency provides a complementary viewpoint where either subtle forgery artifacts or compression errors could be well described. To introduce frequency into the face forgery detection, we propose a novel Frequency in Face Forgery Network (F³-Net), taking advantages of two different but complementary frequency-aware clues, 1) frequency-aware decomposed image components, and 2) local frequency statistics, to deeply mine the forgery patterns via our two-stream collaborative learning framework. We apply DCT as the applied frequency-domain transformation. Through comprehensive studies, we show that the proposed F³-Net significantly outperforms competing state-of-the-art methods on all compression qualities in the challenging FaceForensics++ dataset, especially wins a big lead upon low-quality media.

Keywords: Face Forgery Detection, Frequency, Collaborative Learning

1 Introduction

Rapid development of deep learning driven generative models [26,8,34,35,11] enables an attacker to create, manipulate or even forge the media of a human face (*i.e.*, images and videos, etc.) that cannot be distinguished even by human

* This work was done during the internship of Yuyang Qian and Zixuan Chen at SenseTime Research.

† The first two authors contributed equally.

‡ Corresponding Author.

eyes. However, malicious distribution of forged media would cause security issues and even crisis of confidence in our society. Therefore, it is supremely important to develop effective face forgery detection methods.

Various methods [3,33,45,46,43,58,60,30] have been proposed to detect the forged media. A series of earlier works relied on hand-crafted features *e.g.*, local pattern analysis [21], noise variances evaluation [47] and steganalysis features [13,24] to discover forgery patterns and magnify faint discrepancy between real and forged images. Deep learning introduces another pathway to tackle this challenge, recent learning-based forgery detection methods [14,12] tried to mine the forgery patterns in feature space using convolutional neural networks (CNNs), having achieved remarkable progresses on public datasets, *e.g.*, FaceForensics++ [50].

Current state-of-the-art face manipulation algorithms, such as DeepFake [1], FaceSwap [2], Face2Face [56] and NeuralTextures [55], have been able to conceal the forgery artifacts, so that it becomes extremely difficult to discover the flaws of these refined counterfeits, as shown in Fig. 1(a). What’s worse, if the visual quality of a forged face is tremendously degraded, such as compressed by JPEG or H.264 with a large compression ratio, the forgery artifacts will be contaminated by compression error, and sometimes cannot be captured in RGB domain any more. Fortunately, these artifacts can be captured in frequency domain, as many prior studies suggested [58,38,32,19,57], in the form of unusual frequency distributions when compared with real faces. However, how to involve frequency-aware clues into the deeply learned CNN models? This question also raises alongside. Conventional frequency domains, such as FFT and DCT, do not match the shift-invariance and local consistency owned by nature images, thus vanilla CNN structures might be infeasible. As a result, CNN-compatible frequency representation becomes pivotal if we would like to leverage the discriminative representation power of learnable CNNs for frequency-aware face forgery detection. To this end, we would like to introduce two frequency-aware forgery clues that are compatible with the knowledge mining by deep convolutional networks.

From one aspect, it is possible to decompose an image by separating its frequency signals, while each decomposed image component indicates a certain band of frequencies. The first frequency-aware forgery clue is thus discovered by the intuition that we are able to identify subtle forgery artifacts that are somewhat salient (*i.e.*, in the form of unusual patterns) in the decomposed components with higher frequencies, as the examples shown in the middle column of Fig. 1(b). This clue is compatible with CNN structures, and is surprisingly robust to compression artifacts. From the other aspect, the decomposed image components describe the frequency-aware patterns in the spatial domain, but not explicitly render the frequency information directly in the neural networks. We suggest the second frequency-aware forgery clue as the local frequency statistics. In each densely but regularly sampled local spatial patch, the statistics is gathered by counting the mean frequency responses at each frequency band. These frequency statistics re-assemble back to a multi-channel spatial map, where the

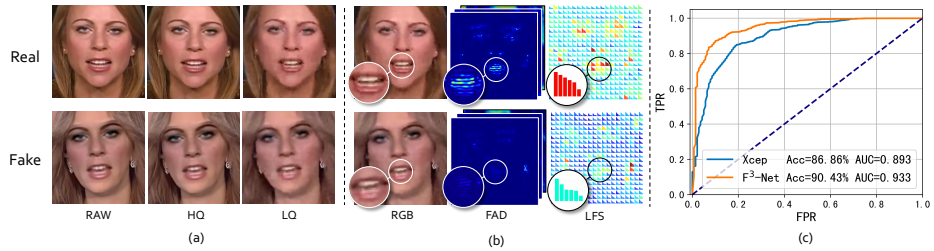


Fig. 1. Frequency-aware tampered clues for face forgery detection. (a) RAW, high quality (HQ) and low quality (LQ) real and fake images with the same identity, manipulation artifacts are barely visible in low quality images. (b) Frequency-aware forgery clues in low quality images using the proposed *Frequency-aware Decomposition (FAD)* and *Local Frequency Statistics (LFS)*. (c) ROC Curve of the proposed **Frequency in Face Forgery Network (F³-Net)** and baseline (*i.e.*, Xception [12]). The proposed F³-Net wins the Xception with a large margin. Best viewed in color.

number of channels is identical to the number of frequency bands. As shown in the last column of Fig. 1(b), the forgery faces have distinct local frequency statistics than the corresponding real ones, even though they look almost the same in the RGB images. Moreover, the local frequency statistics also follows the spatial layouts as the input RGB images, thus also enjoy effective representation learning powered by CNNs. Meanwhile, since the decomposed image components and local frequency statistics are complementary to each other but both of them share inherently similar frequency-aware semantics, thus they can be progressively fused during the feature learning process.

Therefore, we propose a novel **Frequency in Face Forgery Network (F³-Net)**, that capitalizes on the aforementioned frequency-aware forgery clues. The proposed framework is composed of two frequency-aware branches, one aims at learning subtle forgery patterns through *Frequency-aware Image Decomposition (FAD)*, and the other would like to extract high-level semantics from *Local Frequency Statistics (LFS)* to describe the frequency-aware statistical discrepancy between real and forged faces. These two branches, are further gradually fused through a cross-attention module, namely *MixBlock*, which encourages rich interactions between the aforementioned FAD and LFS branches. The whole face forgery detection model is learned by the cross-entropy loss in an end-to-end manner. Extensive experiments demonstrate that the proposed F³-Net significantly improves the performance over low-quality forgery media with a thorough ablation study. We also show that our framework largely exceeds competing state-of-the-arts on all compression qualities in the challenging FaceForensics++ [50]. As shown in Fig.1(c), the effectiveness and superiority of the proposed frequency-aware F³-Net is obviously demonstrated by comparing the ROC curve with Xception [12](baseline, previous state-of-the-art seeing in Sec.4). Our contributions in this paper are summarized as follows:

1) Frequency-aware Decomposition (FAD) aims at learning frequency-aware forgery patterns through frequency-aware image decomposition. The proposed FAD module adaptively partitions the input image in the frequency domain according to learnable frequency bands and represents the image with a series of frequency-aware components.

2) Local Frequency Statistics (LFS) extracts local frequency statistics to describe the statistical discrepancy between real and fake faces. The localized frequency statistics not only reveal the unusual statistics of the forgery images at each frequency band, but also share the structure of natural images, and thus enable effective mining through CNNs.

3) The proposed framework collaboratively learns the frequency-aware clues from FAD and LFS, by a cross-attention (a.k.a MixBlock) powered two-stream networks. The proposed method achieves the state-of-the-art performance on the challenging FaceForensics++ dataset [50], especially wins a big lead in the low quality forgery detection.

2 Related Work

With the development of computer graphics and neural networks especially generative adversarial networks (GANs) [26,8,34,35,11], face forgery detection has gained more and more interest in our society. Various attempts have been made for face forgery detection and achieved remarkable progress, but learning-based generation methods such as NeuralTextures [55] are still difficult to detect because they introduce only small-scale subtle visual artifacts especially in low quality videos. To address the problem, various additional information is used to enhance performance.

Spatial-Based Forgery Detection. To address face forgery detection tasks, a variety of methods have been proposed. Most of them are based on the spatial domain such as RGB and HSV. Some approaches [16,9] exploit specific artifacts arising from the synthesis process such as color or shape cues. Some studies [37,44,27] extract color-space features to classify fake and real images. For example, ELA [27] uses pixel-level errors to detect image forgery. Early methods [6,14] use hand-crafted features for shallow CNN architectures. Recent methods [3,33,45,46,43,58,60,30] use deep neural networks to extract high-level information from the spatial domain and get remarkable progress. MesoInception-4 [3] is a CNN-based Network inspired by InceptionNet [54] to detect forged videos. GANs Fingerprints Analysis [58] introduces deep manipulation discriminator to discover specific manipulation patters. However, most of them use only spatial domain information and therefore are not sensitive to subtle manipulation clues that are difficult to detect in color-space. In our works, we take advantage of frequency cues to mine small-scale detailed artifacts that are helpful especially in low-quality videos.

Frequency-Based Forgery Detection. Frequency domain analysis is a classical and important method in image signal processing and has been widely

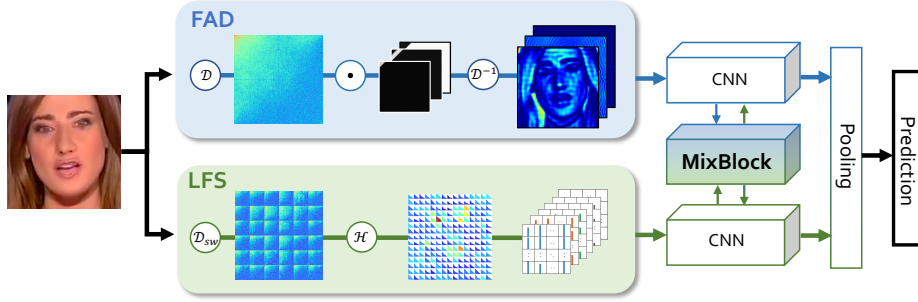


Fig. 2. Overview of the F³-Net. The proposed architecture consists of three novel methods: *FAD* for learning subtle manipulation patterns through frequency-aware image decomposition; *LFS* for extracting local frequency statistics and *MixBlock* for collaborative feature interaction.

used in a number of applications such as image classification [53,52,23], steganalysis [17,10], texture classification [28,25,22] and super-resolution [39,31]. Recently, several attempts have been made to solve forgery detection using frequency cues. Some studies use Wavelet Transform (WT) [7] or Discrete Fourier Transform (DFT) [59,57,19] to convert pictures to frequency domain and mine underlying artifacts. For example, Durall *et al.* [19] extracts frequency-domain information using DFT transform and averaging the amplitudes of different frequency bands. Stuchi *et al.* [53] uses a set of fixed frequency domain filters to extract different range of information followed by a fully connected layer to get the output. Besides, filtering, a classic image signal processing method, is used to refine and mine underlying subtle information in forgery detection, which leverages existing knowledge of the characteristics of fake images. Some studies use high-pass filters [57,15,29,48], Gabor filters [10,22] etc. to extract features of interest (*e.g.* edge and texture information) based on features regarding with high frequency components. Phase Aware CNN [10] uses hand-crafted Gabor and high-pass filters to augment the edge and texture features. Universal Detector [57] finds that significant differences can be obtained in the spectrum between real and fake images after high-pass filtering. However, the filters used in these studies are often fixed and hand-crafted thus fail to capture the forgery patterns adaptively. In our work, we make use of frequency-aware image decomposition to mine frequency forgery cues adaptively.

3 Our Approach

In this section, we introduce the proposed two kinds of frequency-aware forgery clue mining methods, *i.e.*, frequency-aware decomposition (in Sec. 3.1) and local frequency statistics (in Sec. 3.2), and then present the proposed cross-attention two-stream collaborative learning framework (in Sec. 3.3).

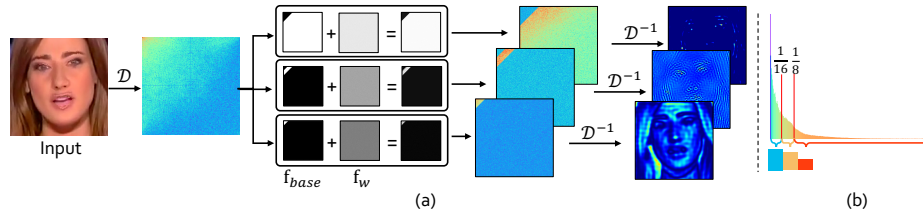


Fig. 3. (a) The proposed *Frequency-aware Decomposition (FAD)* to discover salient frequency components. \mathcal{D} indicates applying Discrete Cosine Transform (DCT). \mathcal{D}^{-1} indicates applying Inversed Discrete Cosine Transform (IDCT). Several frequency band components can be concatenated together to extract a wider range of information. (b) The distribution of the DCT power spectrum. We flatten 2D power spectrum to 1D by summing up the amplitudes of each frequency band. We divide the spectrum into 3 bands with roughly equal energy.

3.1 FAD: Frequency-Aware Decomposition

Towards the frequency-aware image decomposition, former studies usually apply hand-crafted filter banks [10,22] in the spatial domain, thus fail to cover the complete frequency domain. Meanwhile, the fixed filtering configurations make it hard to adaptively capture the forgery patterns. To this end, we propose a novel frequency-aware decomposition (FAD), to adaptively partition the input image in the frequency domain according to a set of learnable frequency filters. The decomposed frequency components can be inversely transformed to the spatial domain, resulting in a series of frequency-aware image components. These components are stacked along the channel axis, and then inputted into a convolutional neural network (in our implementation, we employ an Xception [12] as the backbone) to comprehensively mine forgery patterns.

To be specific, we manually design N binary base filters $\{\mathbf{f}_{base}^i\}_{i=1}^N$ (or called masks) that explicitly partition the frequency domain into low, middle and high frequency bands. And then we add three learnable filters $\{\mathbf{f}_w^i\}_{i=1}^N$ to these base filters. The frequency filtering is a dot-product between the frequency response of the input image and the combined filters $\mathbf{f}_{base}^i + \sigma(\mathbf{f}_w^i)$, $i = \{1, \dots, N\}$, where $\sigma(x) = \frac{1 - \exp(-x)}{1 + \exp(-x)}$ aims at squeezing x within the range between -1 and $+1$. Therefore, to an input image \mathbf{x} , the decomposed image components are obtained by

$$\mathbf{y}_i = \mathcal{D}^{-1}\{\mathcal{D}(\mathbf{x}) \odot [\mathbf{f}_{base}^i + \sigma(\mathbf{f}_w^i)]\}, \quad i = \{1, \dots, N\}. \quad (1)$$

\odot is the element-wise product. We apply \mathcal{D} as the Discrete Cosine Transform (DCT) [4], according to its wide applications in image processing, and its nice layout of the frequency distribution, *i.e.*, low-frequency responses are placed in the top-left corner, and high-frequency responses are located in the bottom-right corner. Moreover, recent compression algorithms, such as JPEG and H.264, usually apply DCT in their frameworks, thus DCT-based FAD will be more

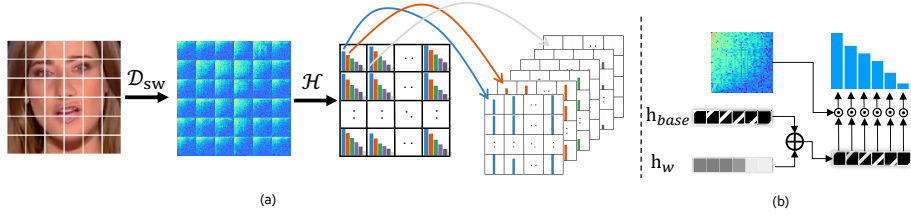


Fig. 4. (a) The proposed *Local Frequency Statistics (LFS)* to extract local frequency domain statistical information. SWDCT indicates applying Sliding Window Discrete Cosine Transform and \mathcal{H} indicates gathering statistics on each grid adaptively. (b) Extracting statistics from a DCT power spectrum graph, \oplus indicates element-wise addition and \odot indicates element-wise multiplication.

compatible towards the description of compression artifacts out of the forgery patterns. Observing the DCT power spectrum of natural images, we find that the spectral distribution is non-uniform and most of the amplitudes are concentrated in the low frequency area. We apply the base filters \mathbf{f}_{base} to divide the spectrum into N bands with roughly equal energy, from low frequency to high frequency. The added learnable $\{\mathbf{f}_w^i\}_{i=1}^N$ provides more adaptation to select the frequency of interest beyond the fixed base filters. Empirically, as shown in Fig. 3(b), the number of bands $N = 3$, the low frequency band \mathbf{f}_{base}^1 is the first 1/16 of the entire spectrum, the middle frequency band \mathbf{f}_{base}^2 is between 1/16 and 1/8 of the spectrum, and the high frequency band \mathbf{f}_{base}^3 is the last 7/8.

3.2 LFS: Local Frequency Statistics

The aforementioned FAD has provided frequency-aware representation that is compatible with CNNs, but it has to represent frequency-aware clues back into the spatial domain, thus fail to directly utilize the frequency information. Also knowing that it is usually infeasible to mine forgery artifacts by extracting CNN features directly from the spectral representation, we then suggest to estimate local frequency statistics (LFS) to not only explicitly render frequency statistics but also match the shift-invariance and local consistency that owned by natural RGB images. These features are then inputted into a convolutional neural network, *i.e.*, Xception [12], to discover high-level forgery patterns.

As shown in Fig. 4(a), we first apply a Sliding Window DCT (SWDCT) on the input RGB image (*i.e.*, taking DCTs densely on sliding windows of the image) to extract the localized frequency responses, and then counting the mean frequency responses at a series of learnable frequency bands. These frequency statistics re-assemble back to a multi-channel spatial map that shares the same layout as the input image. This LFS provides a localized aperture to detect detailed abnormal frequency distributions. Calculating statistics within a set of frequency bands allows a reduced statistical representation, whilst yields a smoother distribution without the interference of outliers.

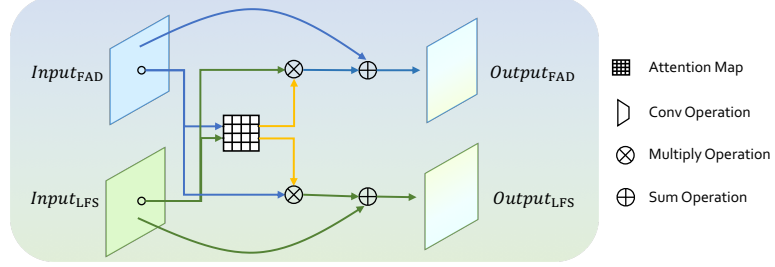


Fig. 5. The proposed *MixBlock*. \otimes indicates matrix multiplication and \oplus indicates element-wise addition.

To be specific, in each window $\mathbf{p} \in \mathbf{x}$, after DCT, the local statistics is gathered in each frequency band, which is constructed similarly as the way used in FAD (see Sec. 3.1). In each band, the statistics become

$$\mathbf{q}_i = \log_{10} \|\mathcal{D}(\mathbf{p}) \odot [\mathbf{h}_{base}^i + \sigma(\mathbf{h}_w^i)]\|_1, \quad i = \{1, \dots, M\}, \quad (2)$$

Note that \log_{10} is applied to balance the magnitude in each frequency band. The frequency bands are collected by equally partitioning the spectrum in to M parts, following the order from low frequency to high frequency. Similarly, \mathbf{h}_{base}^i is the base filter, \mathbf{h}_w^i is the learnable filter, $i = \{1, \dots, M\}$. The local frequency statistics \mathbf{q} for a window \mathbf{p} is then transposed as a $1 \times 1 \times M$ vector. These statistics vectors gathered from all windows are re-assembled into a matrix with downsampled spatial size of the input image, whose number of channels is equal to M . This matrix will act as the input to the later convolutional layers.

Practically in our experiments, we empirically adopt the window size as 10, the sliding stride as 2, and the number of bands as $M = 6$, thus the size of the output matrix will be $149 \times 149 \times 6$ if the input image is of size $299 \times 299 \times 3$.

3.3 Two-stream Collaborative Learning Framework

As mentioned in Sec. 3.1 and Sec. 3.2, the proposed FAD and LFS modules mine the frequency-aware forgery clues from two different but inherently connected aspects. We argue that these two types of clues are different but complementary. Thus, we propose a collaborative learning framework that powered by cross-attention modules, to gradually fuse two-stream FAD and LFS features. To be specific, the whole network architecture of our F^3 -Net is composed of two branches equipped with Xception blocks [12], one is for the decomposed image components generated by FAD, and the other is for local frequency statistics generated by LFS, as shown in Fig. 2.

We propose a cross-attention fusion module for the feature interaction and message passing every several Xception blocks. As shown in Fig. 5, different from the simple concatenation widely used in previous methods [20,32,60], we firstly calculate the cross-attention weight using the feature maps from the two

Table 1. Quantitative results on FaceForensics++ dataset with all quality settings, *i.e.* LQ indicates low quality (heavy compression), HQ indicates high quality (light compression) and RAW indicates raw videos without compression. The bold results are the best. Note that Xception+ELA and Xception-PAFilters are two Xception baselines that are equipped with ELA [27] and PAFilters [10].

Methods	Acc (LQ)	AUC (LQ)	Acc (HQ)	AUC (HQ)	Acc (RAW)	AUC (RAW)
Steg.Features [24]	55.98%	-	70.97%	-	97.63%	-
LD-CNN [14]	58.69%	-	78.45%	-	98.57%	-
Constrained Conv [6]	66.84%	-	82.97%	-	98.74%	-
CustomPooling CNN [49]	61.18%	-	79.08%	-	97.03%	-
MesoNet [3]	70.47%	-	83.10%	-	95.23%	-
Face X-ray [40]	-	0.616	-	0.874	-	-
Xception [12]	86.86%	0.893	95.73%	0.963	99.26%	0.992
Xception-ELA [27]	79.63%	0.829	93.86%	0.948	98.57%	0.984
Xception-PAFilters [10]	87.16%	0.902	-	-	-	-
F ³ -Net (Xception)	90.43%	0.933	97.52%	0.981	99.95%	0.998
Optical Flow [5]	81.60%	-	-	-	-	-
Slowfast [20]	90.53%	0.936	97.09%	0.982	99.53%	0.994
F ³ -Net(Slowfast)	93.02%	0.958	98.95%	0.993	99.99%	0.999

branches. The cross-attention matrix is adopted to augment the attentive features from one stream to another.

In our implementation, we use Xception network [12] pretrained on the ImageNet [18] for both branches, each of which has 12 blocks. The newly-introduced layers and blocks are randomly initialized. The cropped face is adopted as the input of the framework after resized as 299×299 . Empirically, we adopt MixBlock after block 7 and block 12 to fuse two types of frequency-aware clues according to their mid-level and high-level semantics. We train the F³-Net by the well-known cross entropy loss, and the whole system can be trained in an end-to-end fashion.

4 Experiment

4.1 Setting

Dataset. Following previous face forgery detection methods [51,19,40,5], we conduct our experiments on the challenging FaceForensics++ [50] dataset. FaceForensics++ is a face forgery detection video dataset containing 1,000 real videos, in which 720 videos are used for training, 140 videos are reserved for validation and 140 videos for testing. Most videos contain frontal faces without occlusions and were collected from YouTube with the consent of the subjects. Each video undergoes four manipulation methods to generate four fake videos, therefore there are 5,000 videos in total. The number of frames in each video is between 300 and 700. The size of the real videos is augmented four times to solve

category imbalance between the real and fake data. 270 frames are sampled from each video, following the setting as in FF++ [50]. Output videos are generated with different quality levels, so as to create a realistic setting for manipulated videos, *i.e.*, RAW, High Quality (HQ) and Low Quality (LQ), respectively.

We use the face tracking method proposed by Face2Face [56] to crop the face and adopt a conservative crop to enlarge the face region by a factor of 1.3 around the center of the tracked face, following the setting in [50].

Evaluation Metrics. We apply the Accuracy score (Acc) and Area Under the Receiver Operating Characteristic Curve (AUC) as our evaluation metrics. (1) **Acc.** Following FF++ [50], we use the accuracy score as the major evaluation metric in our experiments. This metric is commonly used in face forgery detection tasks [3,13,46]. Specifically, for single-frame methods, we average the accuracy scores of each frame in a video. (2) **AUC.** Following face X-ray [40], we use AUC score as another evaluation metric. For single-frame methods, we also average the AUC scores of each frame in a video.

Implementation Details. In our experiments, we use Xception [12] pretrained on the ImageNet [18] as backbone for the proposed F³-Net. The newly-introduced layers and blocks are randomly initialized. The networks are optimized via SGD. We set the base learning rate as 0.002 and use Cosine [41] learning rate scheduler. The momentum is set as 0.9. The batch size is set as 128. We train for about 150k iterations.

Some studies [51,5] use videos as the input of the face forgery detection system. To demonstrate the generalization of the proposed methods, we also plug LFS and FAD into existing video-based methods, *i.e.* Slowfast-R101 [20] pre-trained on Kinetics-400 [36]. The networks are optimized via SGD. We set the base learning rate as 0.002. The momentum is set as 0.9. The batch size is set as 64. We train the model for about 200k iterations.

4.2 Comparing with previous methods

In this section, on the FaceForensics++ dataset, we compare our method with previous face forgery detection methods.

Evaluations on Different Quality Settings. The results are listed in Tab.1. The proposed F³-Net outperforms all the reference methods on all quality settings, *i.e.*, LQ, HQ and RAW, respectively. According to the low-quality (LQ) setting, the proposed F³-Net achieves 90.43% in Acc and 0.933 in AUC respectively, with a remarkable improvement comparing to the current state-of-the-art methods, *i.e.*, about 3.5% performance gain on Acc score against the best performed reference method (*i.e.*, Xception-PAFilters with 87.16% *v.s.* F³-Net with 90.43%). The performance gains mainly benefit from the information mining from frequency-aware FAD and LFS clues, which helps the proposed F³-Net more capable of detecting subtle manipulation artifacts as well as robust to heavy compression errors than plain RGB-based networks. It is worth noting that some methods [27,19,45,10] also try to employ complementary information from other domains, and try to take advantages of prior knowledge. For example,

Table 2. Quantitative results (Acc) on FaceForensics++ (LQ) dataset with four manipulation methods, *i.e.* DeepFakes(DF) [1], Face2Face(F2F) [56], FaceSwap(FS) [2] and NeuralTextures(NT) [55]). The bold results are the best.

Methods	DF [1]	F2F [56]	FS [2]	NT [55]
Steg.Features [24]	67.00%	48.00%	49.00%	56.00%
LD-CNN [14]	75.00%	56.00%	51.00%	62.00%
Constrained Conv [6]	87.00%	82.00%	74.00%	74.00%
CustomPooling CNN [49]	80.00%	62.00%	59.00%	59.00%
MesoNet [3]	90.00%	83.00%	83.00%	75.00%
Xception [12]	96.01%	93.29%	94.71%	79.14%
F ³ -Net(Xception)	97.97%	95.32%	96.53%	83.32%
Slowfast [20]	97.53%	94.93%	95.01%	82.55%
F ³ -Net(Slowfast)	98.62%	95.84%	97.23%	86.01%

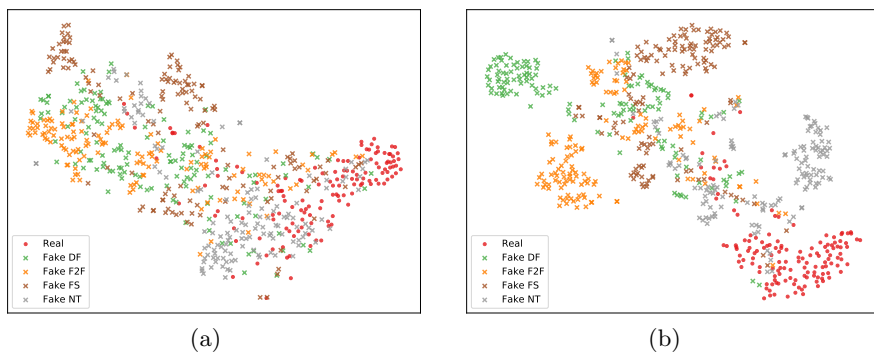


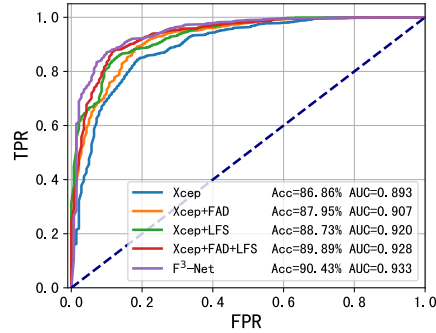
Fig. 6. The t-SNE embedding visualization of the baseline (a) and F³-Net (b) on FaceForensics++ [50] low quality (LQ) task. Red color indicates the real videos, the rest colors represent data generated by different manipulation methods. Best viewed in color.

Steg.Features [24] employs hand-crafted steganalysis features and PAFilters [10] tries to augment the edge and texture features by hand-crafted Gabor and high-pass filters. Different from these methods, the proposed F³-Net makes good use of CNN-friendly and adaptive mechanism to augment the FAD and LFS module, thus significantly boost the performance by a considerable margin

Towards Different Manipulation Types. Furthermore, we evaluate the proposed F³-Net on different face manipulation methods listed in [50]. The models are trained and tested exactly on the low quality videos from one face manipulation methods. The results are shown in Tab. 2. Of the four manipulation methods, the videos generated by NeuralTextures (NT) [55] is extremely challenging due to its excellent generation performance in synthesizing realistic faces without noticeable forgery artifacts. The performance of our proposed method is partic-

ID	FAD	LFS	MixBlock	Acc	AUC
1	-	-	-	86.86%	0.893
2	✓	-	-	87.95%	0.907
3	-	✓	-	88.73%	0.920
4	✓	✓	-	89.89%	0.928
5	✓	✓	✓	90.43%	0.933

(a)



(b)

Fig. 7. (a) Ablation study of the proposed F^3 -Net on the low quality task(LQ). We compare F^3 -Net and its variants by removing the proposed FAD, LFS and MixBlock step by step. (b) ROC Curve of the models in our ablation studies.

ularly impressive when detecting forged faces by NT, leading to an improvement of about 4.2% on the Acc score, against the baseline method Xception [12].

Furthermore, we also showed the t-SNE [42] feature spaces of data in FaceForensics++ [50] low quality (LQ) task, by the Xception and our F^3 -Net, as shown in Fig. 6. Xception cannot divide the real data and NT-based forged data since their features are cluttered in the t-SNE embedding space, as shown in Fig. 6(a). However, although the feature distances between real videos and NT-based forged videos are closer than the rest pairs in the feature space of F^3 -Net, they are still much farther away than those in the feature space of Xception. It, from another viewpoint, proves that the proposed F^3 -Net can mine effective clues to distinguish the real and forged media.

Video-based Extensions. Meanwhile, there are also several studies [51,5] using multiple frames as the input. To evaluate the generalizability of our methods, we involve the proposed LFS and FAD into Slowfast-R101 [20] due to its excellent performance for video classification. The results are shown in Tab.1 and Tab.2. More impressively, our F^3 -Net (Slowfast) achieves the better performances than the baseline using Slowfast only, *i.e.*, 93.02% and 0.958 of Acc and AUC scores in comparison to 90.53% and 0.936, in low quality (LQ) task, as shown in Tab. 1. Slowfast- F^3 -Net also wins over 3% on the NT-based manipulation, as shown in Tab. 2, not to mention the rest three manipulation types. These excellent performances further demonstrate the effectiveness of our proposed frequency-aware face forgery detection method.

4.3 Ablation Study

Effectiveness of LFS, FAD and MixBlock. To evaluate the effectiveness of the proposed LFS, FAD and MixBlock, we quantitatively evaluate F^3 -Net and

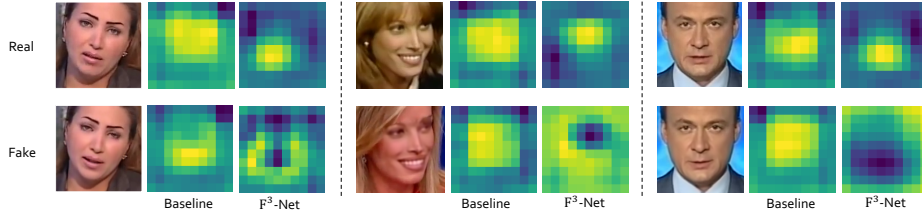


Fig. 8. The visualization of the feature map extracted by baseline (*i.e.*, Xception) and the proposed F³-Net respectively.

its variants: 1) the baseline (Xception), 2) F³-Net w/o LFS and MixBlock, 3) F³-Net w/o FAD and MixBlock, 4) F³-Net w/o MixBlock.

The quantitative results are listed in Fig. 7(a). By comparing model 1 (baseline) and model 2 (Xception with FAD), the proposed FAD consistently improves the Acc and AUC scores. When adding the LFS (model 4) based on model 2, the Acc and AUC scores become even higher. Plugging MixBlock (model 5) into the two branch structure (model 4) gets the best performance, 90.43% and 0.933 for Acc and AUC scores, respectively. These progressively improved performances validate that the proposed FAD and LFS module indeed helps the forgery detection, and they are complementary to each other. MixBlock introduce more advanced cooperation between FAD and LFS, and thus introduce additional gains. As shown in the ROC curves in Fig. 7(b), F³-Net receives the best performance at lower false positive rate (FPR), while low FPR rate is a most challenging scenario to forgery detection system. To better understand the effectiveness of the proposed methods, we visualize the feature maps extracted by the baseline (Xception) and F³-Net, respectively, as shown in Fig. 8. The discriminativeness of these feature maps is obviously improved by the proposed F³-Net, *e.g.*, there are clear differences between real and forged faces in the feature distributions of F³-Net, while the corresponding feature maps generated by Xception are similar and indistinguishable.

Ablation study on FAD. To demonstrate the benefits of adaptive frequency decomposition on complete frequency domain in FAD, we evaluate the proposed FAD and its variants by removing or replacing some components, *i.e.*, 1) Xception (baseline), 2) a group of hand-drafted filters used in Phase Aware CNN [10], denoted as Xception + PAFilters, 3) proposed FAD without learnable filters, denoted as Xception + FAD (\mathbf{f}_{base}), and 4) Xception with the full FAD, denoted as Xception+FAD ($\mathbf{f}_{base} + \mathbf{f}_w$). All the experiments are under the same hyper-parameters for fair comparisons. As shown in the left part of Tab. 3, the performance of Xception is improved by a considerable margin on the Acc and AUC scores after applying the FAD, in comparison with other methods using fixed filters (Xception + PAFilters). If the learnable filters are removed, there will also be a sudden performance drop.

We further demonstrate the importance of extracting complete information from complete frequency domain by quantitatively evaluating FAD with different

Table 3. Ablation study and component analysis on FAD in FF++ low quality (LQ) tasks. Left: comparing traditional fixed filters with the proposed FAD. Right: comparing FAD and its variants with different kinds of frequency components. We use the full FAD in model FAD-All.

Models	Acc	AUC	Models	Acc	AUC
Xception	86.86%	0.893	FAD-Low	86.95%	0.901
Xception+PAFilters [10]	87.16%	0.902	FAD-Mid	87.57%	0.904
Xception+FAD (\mathbf{f}_{base})	87.12%	0.901	FAD-High	87.77%	0.906
Xception+FAD ($\mathbf{f}_{base} + \mathbf{f}_w$)	87.95%	0.907	FAD-All	87.95%	0.907

Table 4. Ablation study on LFS. Here we use only LFS branch and add components step by step. SWDCT indicates using Sliding Window DCT instead of traditional DCT, Stat indicates adopting frequency statistics and D-Stat indicates using our proposed adaptive frequency statistics.

SWDCT	Stat	D-Stat	Acc (LQ)	AUC (LQ)	Acc (HQ)	AUC (HQ)	Acc (RAW)	AUC (RAW)
-	-	-	76.16%	0.724	90.12%	0.905	95.28%	0.948
✓	-	-	82.47%	0.838	93.85%	0.940	97.02%	0.964
✓	✓	-	84.89%	0.865	94.12%	0.936	97.97%	0.975
✓	✓	✓	86.16%	0.889	94.76%	0.951	98.37%	0.983

kinds of frequency components, *i.e.*, 1) FAD-Low, FAD with low frequency band components, 2) FAD-Mid, FAD with middle frequency band components, 3) FAD-High, FAD with high frequency band components and 4) FAD-All, FAD with all frequency bands components. The quantitative results are listed in the right part of Tab. 3. By comparing FAD-Low, FAD-Mid and FAD-High, the model with high frequency band components achieves the best scores, which indicates that high frequency clues are indubitably helpful for forgery detection. It is because high-frequency clues are usually correlated with forgery-sensitive edges and textures. After making use of all three kinds of information (*i.e.*, FAD-All), we achieved the highest result. Since the low frequency components preserve the global picture, the middle and high frequency reveals the small-scale detailed information, concatenating them together helps to obtain richer frequency-aware clues and is able to mine forgery patterns more comprehensively.

Ablation study on LFS. To demonstrate the effectiveness of SWDCT and dynamic statistical strategy in the proposed LFS introduced in Sec. 3.2, we take the experiments (Xception as backbone) on the proposed LFS and its variants, 1) Baseline, of which the frequency spectrum of the full image by traditional DCT; 2) SWDCT, adopting the localized frequency response by SWDCT ; 3) SWDCT+Stat, adopting the general statistical strategy with filters \mathbf{h}_{base} ; 4) SWDCT+Stat+D-Stat, the proposed FAD consisted of SWDCT and the adaptive frequency statistics with learnable filters \mathbf{h}_w . The results are shown in Tab. 4. Comparing with traditional DCT operation on the full image, the proposed SWDCT significantly improves the performance by a large margin since it is

more sensitive to the spatial distributions of the local statistics, and letting the Xception back capture the forgery clues. The improvement of using the statistics is significant and local statistics are more robust to unstable or noisy spectra, especially when optimized by adding the adaptive frequency statistics.

5 Conclusions

In this paper, we propose an innovative face forgery detection framework that can make use of frequency-aware forgery clues, named as F³-Net. The proposed framework is composed of two frequency-aware branches, one focuses on mining subtle forgery patterns through frequency components partition, and the other aims at extracting small-scale discrepancy of frequency statistics between real and forged images. Meanwhile, a novel cross-attention module is applied for two-stream collaborative learning. Extensive experiments demonstrate the effectiveness and significance of the proposed F³-Net on FaceForencis++ dataset, especially in the challenging low quality task.

Acknowledgements. This work is supported by SenseTime Group Limited, in part by key research and development program of Guangdong Province, China, under grant 2019B010154003. The corresponding authors are Guojun Yin and Lu Sheng. The contribution of Yuyang Qian and Guojun Yin are Equal.

References

1. Deepfakes github., <https://github.com/deepfakes/faceswap>
2. Faceswap., <https://github.com/MarekKowalski/FaceSwap/>
3. Afchar, D., Nozick, V., Yamagishi, J., Echizen, I.: Mesonet: a compact facial video forgery detection network. In: 2018 IEEE International Workshop on Information Forensics and Security (WIFS). pp. 1–7. IEEE (2018)
4. Ahmed, N., Natarajan, T., Rao, K.R.: Discrete cosine transform. *IEEE transactions on Computers* **100**(1), 90–93 (1974)
5. Amerini, I., Galteri, L., Caldelli, R., Del Bimbo, A.: Deepfake video detection through optical flow based cnn. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 0–0 (2019)
6. Bayar, B., Stamm, M.C.: 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. pp. 5–10 (2016)
7. Bentley, P.M., McDonnell, J.: Wavelet transforms: an introduction. *Electronics & communication engineering journal* **6**(4), 175–186 (1994)
8. Brock, A., Donahue, J., Simonyan, K.: Large scale gan training for high fidelity natural image synthesis. arXiv preprint arXiv:1809.11096 (2018)
9. Carvalho, T., Faria, F.A., Pedrini, H., Torres, R.d.S., Rocha, A.: Illuminant-based transformed spaces for image forensics. *IEEE transactions on information forensics and security* **11**(4), 720–733 (2015)

10. Chen, M., Sedighi, V., Boroumand, M., Fridrich, J.: Jpeg-phase-aware convolutional neural network for steganalysis of jpeg images. In: Proceedings of the 5th ACM Workshop on Information Hiding and Multimedia Security. pp. 75–84 (2017)
11. Choi, Y., Choi, M., Kim, M., Ha, J.W., Kim, S., Choo, J.: Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 8789–8797 (2018)
12. Chollet, F.: Xception: Deep learning with depthwise separable convolutions. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1251–1258 (2017)
13. Cozzolino, D., Gagnaniello, D., Verdoliva, L.: Image forgery localization through the fusion of camera-based, feature-based and pixel-based techniques. In: 2014 IEEE International Conference on Image Processing (ICIP). pp. 5302–5306. IEEE (2014)
14. Cozzolino, D., Poggi, G., Verdoliva, L.: 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. pp. 159–164 (2017)
15. D’Avino, D., Cozzolino, D., Poggi, G., Verdoliva, L.: Autoencoder with recurrent neural networks for video forgery detection. *Electronic Imaging* **2017**(7), 92–99 (2017)
16. De Carvalho, T.J., Riess, C., Angelopoulou, E., Pedrini, H., de Rezende Rocha, A.: Exposing digital image forgeries by illumination color classification. *IEEE Transactions on Information Forensics and Security* **8**(7), 1182–1194 (2013)
17. Denemark, T.D., Boroumand, M., Fridrich, J.: Steganalysis features for content-adaptive jpeg steganography. *IEEE Transactions on Information Forensics and Security* **11**(8), 1736–1746 (2016)
18. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: 2009 IEEE conference on computer vision and pattern recognition. pp. 248–255. Ieee (2009)
19. Durall, R., Keuper, M., Pfrendt, F.J., Keuper, J.: Unmasking deepfakes with simple features. arXiv preprint arXiv:1911.00686 (2019)
20. Feichtenhofer, C., Fan, H., Malik, J., He, K.: Slowfast networks for video recognition. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 6202–6211 (2019)
21. Ferrara, P., Bianchi, T., De Rosa, A., Piva, A.: Image forgery localization via fine-grained analysis of cfa artifacts. *IEEE Transactions on Information Forensics and Security* **7**(5), 1566–1577 (2012)
22. Fogel, I., Sagi, D.: Gabor filters as texture discriminator. *Biological cybernetics* **61**(2), 103–113 (1989)
23. Franzen, F.: Image classification in the frequency domain with neural networks and absolute value dct. In: International Conference on Image and Signal Processing. pp. 301–309. Springer (2018)
24. Fridrich, J., Kodovsky, J.: Rich models for steganalysis of digital images. *IEEE Transactions on Information Forensics and Security* **7**(3), 868–882 (2012)
25. Fujieda, S., Takayama, K., Hachisuka, T.: Wavelet convolutional neural networks for texture classification. arXiv preprint arXiv:1707.07394 (2017)
26. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: Advances in neural information processing systems. pp. 2672–2680 (2014)

27. Gunawan, T.S., Hanafiah, S.A.M., Kartiwi, M., Ismail, N., Za'bah, N.F., Nordin, A.N.: Development of photo forensics algorithm by detecting photoshop manipulation using error level analysis. *Indonesian Journal of Electrical Engineering and Computer Science (IJEECS)* **7**(1), 131–137 (2017)
28. Haley, G.M., Manjunath, B.: Rotation-invariant texture classification using a complete space-frequency model. *IEEE transactions on Image Processing* **8**(2), 255–269 (1999)
29. Hsu, C.C., Hung, T.Y., Lin, C.W., Hsu, C.T.: Video forgery detection using correlation of noise residue. In: 2008 IEEE 10th workshop on multimedia signal processing. pp. 170–174. IEEE (2008)
30. Hsu, C.C., Lee, C.Y., Zhuang, Y.X.: Learning to detect fake face images in the wild. In: 2018 International Symposium on Computer, Consumer and Control (IS3C). pp. 388–391. IEEE (2018)
31. Huang, H., He, R., Sun, Z., Tan, T.: Wavelet-srnet: A wavelet-based cnn for multi-scale face super resolution. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1689–1697 (2017)
32. Huang, Y., Zhang, W., Wang, J.: Deep frequent spatial temporal learning for face anti-spoofing. arXiv preprint arXiv:2002.03723 (2020)
33. Jeon, H., Bang, Y., Woo, S.S.: Fdftnet: Facing off fake images using fake detection fine-tuning network. arXiv preprint arXiv:2001.01265 (2020)
34. Karras, T., Aila, T., Laine, S., Lehtinen, J.: Progressive growing of gans for improved quality, stability, and variation. arXiv preprint arXiv:1710.10196 (2017)
35. Karras, T., Laine, S., Aila, T.: A style-based generator architecture for generative adversarial networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 4401–4410 (2019)
36. Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
37. Li, H., Li, B., Tan, S., Huang, J.: Detection of deep network generated images using disparities in color components. arXiv preprint arXiv:1808.07276 (2018)
38. Li, J., Wang, Y., Tan, T., Jain, A.K.: Live face detection based on the analysis of fourier spectra. In: Biometric Technology for Human Identification. vol. 5404, pp. 296–303. International Society for Optics and Photonics (2004)
39. Li, J., You, S., Robles-Kelly, A.: A frequency domain neural network for fast image super-resolution. In: 2018 International Joint Conference on Neural Networks (IJCNN). pp. 1–8. IEEE (2018)
40. Li, L., Bao, J., Zhang, T., Yang, H., Chen, D., Wen, F., Guo, B.: Face x-ray for more general face forgery detection. arXiv preprint arXiv:1912.13458 (2019)
41. Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983 (2016)
42. Maaten, L.v.d., Hinton, G.: Visualizing data using t-sne. *Journal of machine learning research* **9**(Nov), 2579–2605 (2008)
43. Marra, F., Gagnaniello, D., Cozzolino, D., Verdoliva, L.: Detection of gan-generated fake images over social networks. In: 2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR). pp. 384–389. IEEE (2018)
44. McCloskey, S., Albright, M.: Detecting gan-generated imagery using color cues. arXiv preprint arXiv:1812.08247 (2018)
45. Nguyen, H.H., Fang, F., Yamagishi, J., Echizen, I.: Multi-task learning for detecting and segmenting manipulated facial images and videos. arXiv preprint arXiv:1906.06876 (2019)

46. Nguyen, H.H., Yamagishi, J., Echizen, I.: Use of a capsule network to detect fake images and videos. arXiv preprint arXiv:1910.12467 (2019)
47. Pan, X., Zhang, X., Lyu, S.: Exposing image splicing with inconsistent local noise variances. In: 2012 IEEE International Conference on Computational Photography (ICCP). pp. 1–10. IEEE (2012)
48. Pandey, R.C., Singh, S.K., Shukla, K.K.: Passive forensics in image and video using noise features: A review. *Digital Investigation* **19**, 1–28 (2016)
49. Rahmouni, N., Nozick, V., Yamagishi, J., Echizen, I.: Distinguishing computer graphics from natural images using convolution neural networks. In: 2017 IEEE Workshop on Information Forensics and Security (WIFS). pp. 1–6. IEEE (2017)
50. Rossler, A., Cozzolino, D., Verdoliva, L., Riess, C., Thies, J., Nießner, M.: Face-forensics++: Learning to detect manipulated facial images. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1–11 (2019)
51. Sabir, E., Cheng, J., Jaiswal, A., AbdAlmageed, W., Masi, I., Natarajan, P.: Recurrent convolutional strategies for face manipulation detection in videos. *Interfaces (GUI)* **3**, 1 (2019)
52. Sarlashkar, A., Bodruzzaman, M., Malkani, M.: Feature extraction using wavelet transform for neural network based image classification. In: Proceedings of Thirtieth Southeastern Symposium on System Theory. pp. 412–416. IEEE (1998)
53. Stuchi, J.A., Angeloni, M.A., Pereira, R.F., Boccato, L., Folego, G., Prado, P.V., Attux, R.R.: Improving image classification with frequency domain layers for feature extraction. In: 2017 IEEE 27th International Workshop on Machine Learning for Signal Processing (MLSP). pp. 1–6. IEEE (2017)
54. Szegedy, C., Ioffe, S., Vanhoucke, V., Alemi, A.A.: Inception-v4, inception-resnet and the impact of residual connections on learning. In: Thirty-first AAAI conference on artificial intelligence (2017)
55. Thies, J., Zollhöfer, M., Nießner, M.: Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)* **38**(4), 1–12 (2019)
56. Thies, J., Zollhofer, M., Stamminger, M., Theobalt, C., Nießner, M.: Face2face: Real-time face capture and reenactment of rgb videos. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2387–2395 (2016)
57. Wang, S.Y., Wang, O., Zhang, R., Owens, A., Efros, A.A.: Cnn-generated images are surprisingly easy to spot... for now. arXiv preprint arXiv:1912.11035 (2019)
58. Yu, N., Davis, L.S., Fritz, M.: Attributing fake images to gans: Learning and analyzing gan fingerprints. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 7556–7566 (2019)
59. Zhang, X., Karaman, S., Chang, S.F.: Detecting and simulating artifacts in gan fake images. arXiv preprint arXiv:1907.06515 (2019)
60. Zhou, P., Han, X., Morariu, V.I., Davis, L.S.: Two-stream neural networks for tampered face detection. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW). pp. 1831–1839. IEEE (2017)

6 Appendix

6.1 Greedy-searched Threshold for Accuracy

The accuracy value is influenced by the threshold θ . If the output score s of the model is above θ , the instance is recognized as **fake**, otherwise, the video is classified as **real**. In most previous works on the binary classification tasks, the parameter θ is set as 0.5 by default. In our work, we find that the parameter θ makes a big difference on the exact value of the metric Acc. Therefore, we greedily search the highest accuracy value in the validation set of FF++ dataset with different thresholds in the parameter set $\{0, 0.01, 0.02, 0.03, \dots, 0.98, 0.99, 1\}$ for the best threshold θ_{max} and then evaluate on the test set with the greedy-searched threshold θ_{max} .

In Tab.1 and Tab.2 in the main paper, the Acc values of our method and the previous methods we re-implemented, *i.e.* Xception [12], Xception-ELA [27], Xception-PAFilters [10], F³-Net (Xception), Slowfast [20] and F³-Net(Slowfast) are calculated by the greedy-searched thresholds θ_{max} respectively. The comparison of accuracy values with the greedy-searched threshold θ_{max} and standard threshold $\theta_{0.5}$ are shown in Tab.5 and Tab.6. Compared with the standard threshold $\theta = 0.5$, the accuracy values with greedy-searched threshold θ_{max} are higher under all of the settings, especially when the classifier performs not very well. Regardless of the greedy-searched threshold, comparing our method F³-Net and previous methods in Tab.7 and Tab.8, it is also demonstrated the priority and effectiveness of our proposed F³-Net using the standard threshold 0.5.

Table 5. Quantitative results (Acc) with greedy-searched threshold and standard threshold on FaceForensics++ dataset with all quality settings, *i.e.* LQ indicates low quality (heavy compression), HQ indicates high quality (light compression) and RAW indicates raw videos without compression. The bold results are the best. Note that Xception+ELA and Xception-PAFilters are two Xception baselines that are equipped with ELA [27] and PAFilters [10]. The values in the brackets are the accuracy values with $\theta = 0.5$ and the value drops compared with the greedy-searched threshold.

Methods	Acc(LQ)	Acc(HQ)	Acc(RAW)
Xception [12]	86.86%(82.71%,4.15%↓)	95.73%(95.04%,0.69%↓)	99.26%(98.77%,0.49%↓)
Xception-ELA [27]	79.63%(73.69%,5.94%↓)	93.86%(92.09%,1.77%↓)	98.57%(97.13%,1.44%↓)
Xception-PAFilters [10]	87.16%(83.24%,3.92%↓)	-	-
F ³ -Net (Xception)	90.43% (86.89%,3.54%↓)	97.52% (97.31%,0.21%↓)	99.95% (99.84%,0.11%↓)
Slowfast [20]	90.53%(88.25%,2.28%↓)	97.09%(96.92%,0.17%↓)	99.53%(99.34%,0.19%↓)
F ³ -Net(Slowfast)	93.02% (92.37%,0.65%↓)	98.95% (98.64%,0.31%↓)	99.99% (99.91%,0.08%↓)

6.2 Best Practice

In this paper, we propose an innovative face forgery detection framework that can make use of frequency-aware forgery clues, named as F³-Net. Additionally,

Table 6. Quantitative results (Acc) with greedy-searched threshold and standard threshold on FaceForensics++ (LQ) dataset with four manipulation methods, *i.e.* DeepFakes(DF) [1], Face2Face(F2F) [56], FaceSwap(FS) [2] and NeuralTextures(NT) [55]. The bold results are the best. The values in the brackets are the accuracy values with $\theta = 0.5$.

Methods	DF [1]	F2F [56]	FS [2]	NT [55]
Xception [12]	96.01% (94.27%)	93.29% (91.98%)	94.71% (93.03%)	79.14% (76.43%)
F ³ -Net(Xception)	97.97% (96.81%)	95.32% (94.01%)	96.53% (95.85%)	83.32% (79.36%)
Slowfast [20]	97.53% (96.01%)	94.93% (92.47%)	95.01% (93.86%)	82.55% (80.07%)
F ³ -Net(Slowfast)	98.62% (98.54%)	95.84% (93.91%)	97.23% (96.82%)	86.01% (83.74%)

Table 7. Quantitative results (Acc with $\theta = 0.5$) on FaceForensics++ dataset with all quality settings, *i.e.* LQ indicates low quality (heavy compression), HQ indicates high quality (light compression) and RAW indicates raw videos without compression. The bold results are the best. Note that Xception+ELA and Xception-PAFilters are two Xception baselines that are equipped with ELA [27] and PAFilters [10].

Methods	Acc	AUC	Acc	AUC	Acc	AUC
	(LQ)	(LQ)	(HQ)	(HQ)	(RAW)	(RAW)
Steg.Features [24]	55.98%	-	70.97%	-	97.63%	-
LD-CNN [14]	58.69%	-	78.45%	-	98.57%	-
Constrained Conv [6]	66.84%	-	82.97%	-	98.74%	-
CustomPooling CNN [49]	61.18%	-	79.08%	-	97.03%	-
MesoNet [3]	70.47%	-	83.10%	-	95.23%	-
Face X-ray [40]	-	0.616	-	0.874	-	-
Xception [12]	82.71%	0.893	95.04%	0.963	98.77%	0.992
Xception-ELA [27]	73.69%	0.829	92.09%	0.948	97.13%	0.984
Xception-PAFilters [10]	83.24%	0.902	-	-	-	-
F ³ -Net (Xception)	86.89%	0.933	97.31%	0.981	99.84%	0.998
Optical Flow [5]	81.60%	-	-	-	-	-
Slowfast [20]	88.25%	0.936	96.92%	0.982	99.34%	0.994
F ³ -Net(Slowfast)	92.37%	0.958	98.64%	0.993	99.91%	0.999

Table 8. Quantitative results (Acc with $\theta = 0.5$) on FaceForensics++ (LQ) dataset with four manipulation methods, *i.e.* DeepFakes(DF) [1], Face2Face(F2F) [56], FaceSwap(FS) [2] and NeuralTextures(NT) [55]. The bold results are the best.

Methods	DF [1]	F2F [56]	FS [2]	NT [55]
Steg.Features [24]	67.00%	48.00%	49.00%	56.00%
LD-CNN [14]	75.00%	56.00%	51.00%	62.00%
Constrained Conv [6]	87.00%	82.00%	74.00%	74.00%
CustomPooling CNN [49]	80.00%	62.00%	59.00%	59.00%
MesoNet [3]	90.00%	83.00%	83.00%	75.00%
Xception [12]	94.27%	91.98%	93.03%	76.43%
F ³ -Net(Xception)	96.81%	94.01%	95.85%	79.36%
Slowfast [20]	96.01%	92.47%	93.86%	80.07%
F ³ -Net(Slowfast)	98.54%	93.91%	96.82%	83.74%

in this supplementary material, we discuss the the choice of hyper parameters in our framework.

Table 9. Comparison of the variants of LFS with different values of window sizes in SWDCT. The stride size is set as 2 and we train and test models in FF++ low quality (LQ) tasks. The bold results are the best.

Window size	2	5	10	20	30
AUC	0.836	0.883	0.889	0.876	0.853

Window size & stride size of LFS. As described in Sec.3.2 in the main paper, Sliding Window DCT (SWDCT) is adopted in Local Frequency Statistics (LFS) branch to extract localized frequency responses. The window size and the stride size in sliding window of SWDCT need to be selected to get the best practice. In this supplementary material, we further provide a discussion on different values of window sizes and stride sizes in SWDCT.

The comparison results of various window sizes of SWDCT among 2, 5, 10, 20 and 30 are listed in Tab. 9. If LFS applies a small window size (*i.e.*, 2), it cannot extract frequency statistics around edges and some other structures coped with high frequencies. However, if the window size is too large (*i.e.*, 20 and 30), the LFS is less sensitive to the local abnormal statistics thus there will also be a significant performance drop. We adopt the window size as 10 to get the best practice in our F³-Net.

The stride size of SWDCT is also worth investigating as it is highly correlated with the detection performance and computational cost. We validate the value of stride among 10, 6, 4, 3, 2 and 1. We train and test models in FF++ low quality (LQ) tasks and the window size is set as 10. By comparing the results listed in Tab. 10, we employ stride as 2 because it finds the best trade-off between effectiveness and efficiency.

Table 10. Comparison of the variants of LFS with different values of stride size in SWDCT. The window size is set as 10 and we train and test models in FF++ low quality (LQ) tasks. The time cost is averaged on 1000 runs. The bold results receive the best accuracy and AUC scores while the underline ones strike a balance between effectiveness and efficiency.

Stride size	10	6	4	3	2	1
AUC	0.835	0.855	0.876	0.882	<u>0.889</u>	0.891
Time(s/iter)	0.368	0.376	0.384	0.394	<u>0.402</u>	0.832