# M2TR: Multi-modal Multi-scale Transformers for Deepfake Detection

Junke Wang[1,2], Zuxuan Wu[1,2], Wenhao Ouyang[1,2], Xintong Han[3]
Jingjing Chen[1,2], Ser-Nam Lim[4], Yu-Gang Jiang[1,2] †
[1]Shanghai Key Lab of Intelligent Information Processing,
School of Computer Science, Fudan University
[2]Shanghai Collaborative Innovation Center on Intelligent Visual Computing, [3]Huya Inc, [4]Meta AI

## ABSTRACT

The widespread dissemination of Deepfakes demands effective approaches that can detect perceptually convincing forged images. In this paper, we aim to capture the subtle manipulation artifacts at different scales using transformer models. In particular, we introduce a **M**ulti-modal **M**ulti-scale **TR**ansformer (**M2TR**), which operates on patches of different sizes to detect local inconsistencies in images at different spatial levels. M2TR further learns to detect forgery artifacts in the frequency domain to complement RGB information through a carefully designed cross modality fusion block. In addition, to stimulate Deepfake detection research, we introduce a high-quality Deepfake dataset, SR-DF, which consists of 4,000 DeepFake videos generated by state-of-the-art face swapping and facial reenactment methods. We conduct extensive experiments to verify the effectiveness of the proposed method, which outperforms state-of-the-art Deepfake detection methods by clear margins.

## CCS CONCEPTS

• **Computing methodologies → Computer vision**.

## KEYWORDS

Deepfake detection, Multiscale transformer, Deepfake dataset

## 1 INTRODUCTION

Recent years have witnessed the rapid development of Deepfake techniques [28, 31, 45, 54], which enable attackers to manipulate
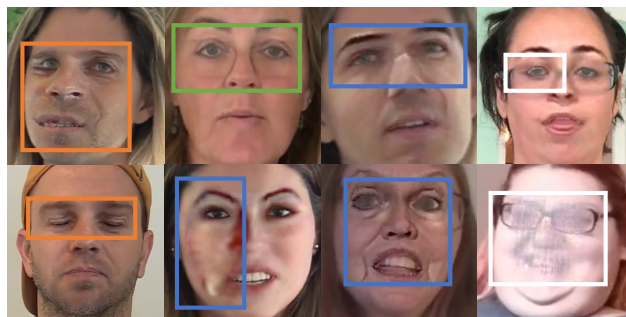
† Corresponding author.

**Figure 1: Visual artifacts of images in the DFDC dataset [15], including color mismatch (blue), shape distortion (orange), visible boundaries (green), and facial blurring (white).**

the facial regions of an image and generate a forged image. As the synthesized images are becoming more photo-realistic, it is extremely difficult to distinguish whether an image has been manipulated even for the human eyes. At the same time, these forged images might be distributed on the Internet for malicious purposes, which could bring societal implications. The above challenges have driven the development of Deepfake forensics using deep neural networks [1, 6, 27, 34, 36, 42, 72]. Most existing approaches take as inputs a face region cropped out of an entire image and produce a binary real/fake prediction with deep CNN models [7, 22, 52]. These methods capture artifacts from the face regions in a single scale with stacked convolutional operations [55, 70]. While decent detection results are achieved by stacked convolutions, they excel at modeling local information but fail to consider the relationships of pixels globally due to constrained receptive field.

We posit that relationships among pixels are particularly useful for Deepfake detection, since pixels in certain artifacts are clearly different from the remaining pixels in the image. On the other hand, we observe that forgery patterns vary in size. For instance, Figure 1 gives examples from the DFDC dataset [15]. We can see that some forgery traces such as color mismatch occur in small regions (like the mouth corners), while other forgery signals such as visible boundaries that almost span the entire image. Therefore, how to effectively explore regions of different scales in images is extremely critical for Deepfake detection.

To address the above limitations, we explore transformers to model the relationships of pixels due to their strong capability of

long-term dependency modeling for both natural language processing tasks [14, 48, 61] and computer vision tasks [2, 16, 62, 63, 73]. Unlike traditional vision transformers that usually operate on a single-scale, we propose a multi-scale architecture to capture forged regions that potentially have different sizes. Furthermore, [17, 26, 46, 64] suggest that the artifacts of forged images will be destroyed by perturbations such as JPEG compression, making them imperceptible in the RGB domain but can still be detected in the frequency domain. This motivates us to use frequency information as a complementary modality in order to reveal artifacts that are no longer perceptible in the RGB domain.

To this end, we introduce M2TR, a Multi-modal Multi-scale Transformer, for Deepfake detection. As illustrated in Figure 2, M2TR follows a two-stream architecture, where the RGB stream captures the inconsistency among different regions within an image at multiple scales in RGB domain, and the frequency stream adopts learnable frequency filters to filter out forgery features in frequency domain. We also design a cross modality fusion block to combine the information from both streams more effectively in an interactive fashion. Finally, the integrated features are input to fully connected layers to generate prediction results. In addition to binary classification, we also predict the manipulated regions of the face image in a multi-task manner. The rationale behind is that binary classification tends to result in easily overfitted models. Therefore, we use face masks as additional supervisory signals to mitigate overfitting.

The availability of large-scale training data is an essential factor in the development of Deepfake detection methods. However, the quality of visual samples in current Deepfake datasets [15, 30, 37, 49, 68] is limited, containing clear artifacts (see Figure 1) like color mismatch, shape distortion, visible boundaries, and facial blurring. Therefore, there is still a huge gap between the images in existing datasets and forged images in the wild which are circulated on the Internet. Although the visual quality of Celeb-DF [37] is relatively high compared to others, they use only one face swapping method to generate forged images, lacking sample diversity. In addition, there are no unbiased and comprehensive evaluation metrics to measure the quality of Deepfake datasets, which is not conducive to the development of subsequent Deepfake research.

In this paper, we present a large-scale and high-quality Deepfake dataset, **S**wapping and **R**eenactment **D**eep**F**ake (**SR-DF**) dataset, which is generated using the state-of-the-art face swapping and facial reenactment methods [19, 44, 51, 59] for the development and evaluation of Deepfake detection methods. Besides, we propose a set of evaluation criteria to measure the quality of Deepfake datasets from different perspectives. We hope the release of SR-DF dataset and the evaluation systems will benefit the future research of Deepfake detection. Our work makes the following key contributions:

- We propose a Multi-modal Multi-scale Transformer (M2TR) for Deepfake forensics, which uses a multi-scale transformer to detect local inconsistencies at different scales and leverages frequency features to improve the robustness. Extensive experiments demonstrate that our method achieves state-of-the-art detection performance on different datasets.

- We introduce a large-scale and challenging Deepfake dataset SR-DF, which is generated with state-of-the-art face swapping and facial reenactment methods

- We construct the most comprehensive evaluation system and demonstrate that SR-DF dataset is well-suited for training Deepfake detection methods due to its quality and diversity.

## 2 RELATED WORK

**Deepfake Detection** To mitigate the security threat brought by Deepfakes, a variety of methods have been proposed for Deepfake detection. [72] uses a two-stream architecture to capture facial manipulation clues and patch inconsistency separately, while [42] simultaneously identifies forged faces and locates the manipulated regions with multi-task learning.

Recently, Face X-ray [34] proposes to detect the blending boundaries based on an observation that the step of blending a forged face into the background is commonly used by most existing face manipulation methods. DCViT [65] extracts features from the face image using a CNN model, which are then fed to a traditional single-scale transformer for forgery detection. MaDD [71] proposes a multi-attentional Deepfake detection framework to capture artifacts with multiple attention maps. However, most of them only focus on the features in the RGB domain, thus failing to detect forged images which are manipulated subtly in the color-space. Instead, $F^3$-Net [46] adopts a two-branch architecture where one makes use of frequency clues to recognize forgery patterns and the other extracts the discrepancy of frequency statistics between real and fake images. In this paper, we use a multi-scale transformer to capture local inconsistencies at different scales for forgery detection, and additionally introduce frequency modality to improve the robustness of our method to various image compression algorithms.

**Visual Transformers** Transformers [61] have demonstrated impressive performance for natural language processing tasks due to strong abilities in modeling long-range context information. Recently, researchers have demonstrated remarkable interests in using the transformer for a variety of computer vision tasks. Typically, visual transformers [2, 16, 58] model the interactions between tokens of the same scale with the self-attention mechanisms. ViT reshapes an image into a sequence of flattened patches and inputs them to the transformer encoder for image classification [16]. DETR uses a common CNN to extract semantic features from the input image, which are then input to a transformer-based encoder-decoder architecture for object detection [2]. Unlike these approaches, we propose to split the inputs into patches of different sizes, and integrate multi-scale information for better visual representation with vision transformers.

## 3 APPROACH

We aim to detect the subtle forgery artifacts that are hidden in the forged images and improve the robustness to image compression with frequency features. In this section, we introduce the Multi-modal Multi-scale Transformer (M2TR) for Deepfake detection, which consists of stacked multi-scale transformers (Sec 3.1), frequency filters (Sec 3.2), and cross modality fusion blocks (Sec 3.3). Figure 2 gives an overview of the framework.
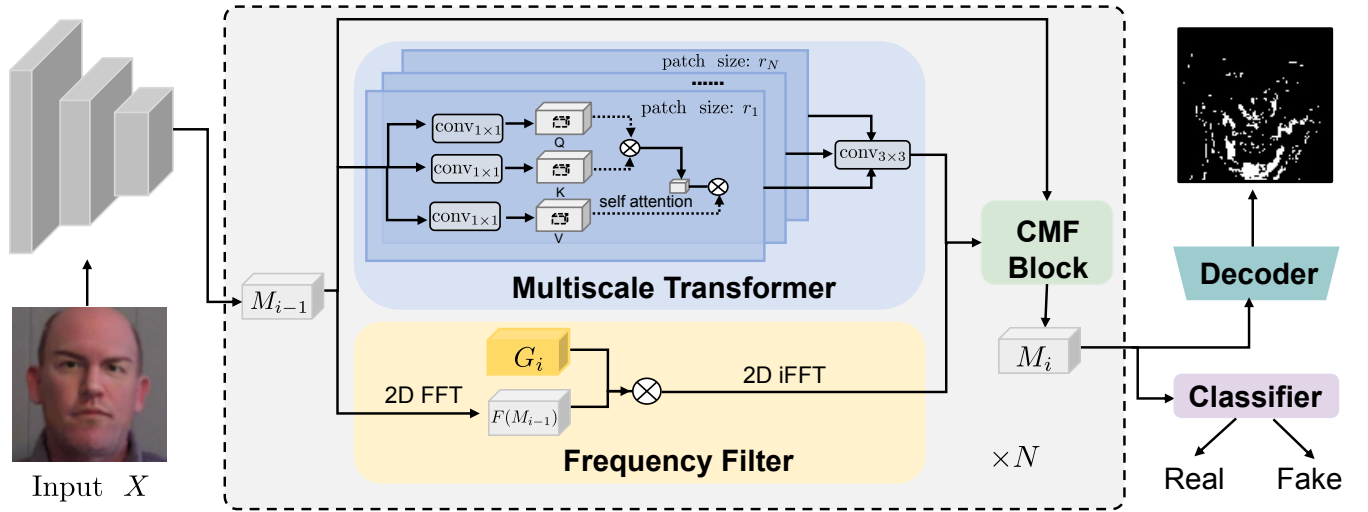
**Figure 2: Overview of the proposed M2TR. The input is a suspicious face image (H × W × C), and the output includes both a forgery detection result and a predicted mask (H × W × 1), which locates the forgery regions.**

More formally, we denote an input image as $X \in \mathbb{R}^{H \times W \times 3}$, where $H$ and $W$ are the height and width of the image, respectively. We first use several convolutional layers to extract features $F \in \mathbb{R}^{(H/4) \times (W/4) \times C}$ of $X$, which are then input to successive multi-scale transformer and frequency filters for forgery clues detection. The intuition of using convolutions here is to ensure faster convergence and more stable training [66].

### 3.1 Multi-scale Transformer

To capture forgery patterns at multiple scales, we introduce a multi-scale transformer which operates on patches of different sizes. Taking the output of the previous cross modality fusion block $M_{i-1}$ (which is initialized as $F$) as input, we split it into spatial patches of different sizes and calculate patch-wise self-attention in different heads. Specifically, we first extract patches of shape $r_h \times r_h \times C$ from $M_{i-1}$, and reshape them into 1-dimension vectors for the $h$-th head. After that, we use fully-connected layers to embed the flattened vectors into query embeddings $Q_i^h \in \mathbb{R}^{N \times C_h}$, where $N = (H/4r_h) \times (W/4r_h)$, and $C_h = r_h \times r_h \times C$. Similar operations are implemented to obtain key embeddings $K_i^h$ and value embeddings $V_i^h$, respectively. Then we calculate the attention matrix through the following process:

$$A_i^h = softmax \left( \frac{Q_i^h (K_i^h)^T}{C_h} \right) V_i^h, \tag{1}$$

$A_i^h$ is then reshaped to the original spatial resolution. Finally, the features from different heads are concatenated and further passed through a 2D residual block to obtain the output $T_i \in \mathbb{R}^{(H/4) \times (W/4) \times C}$.

### 3.2 Frequency Filter

It has been shown that artifacts in manipulated images and videos are no longer perceptible with compression approaches like JPEG

compression [17, 26, 64]. Therefore, we extract the forgery features in the frequency domain to complement RGB features.

Specifically, we first apply 2D FFT along the spatial dimensions to transform $M_{i-1}$ to the frequency domain, and obtain the spectrum representation $\mathcal{F}(M_{i-1}) \in \mathbb{R}^{H/4 \times W/4 \times C}$. We then multiply $\mathcal{F}(M_{i-1})$ with a learnable filter $G_i \in \mathbb{R}^{H/4 \times W/4 \times C}$ to model the dependencies of different frequency band components:

$$\hat{G}_i = G_i \odot \mathcal{F}(M_{i-1}), \tag{2}$$

where $\odot$ denotes the Hadamard product. Finally, we perform the inverse FFT to covert $\hat{G}_i$ back to the spatial domain and obtain frequency-aware features $W_i$.

### 3.3 Cross Modality Fusion

Given RGB features $T_i$ and frequency features $W_i$, we use a Cross Modality Fusion (CMF) block to fuse them into a unified representation. Inspired by the architecture of self-attention in transformers, we design a fusion block using the query-key-value mechanism.

Specifically, we first embed $T_i$ and $W_i$ into $Q$, $K$, and $V$ using $1 \times 1$ convolutions $conv_q$, $conv_k$, and $conv_v$, respectively. Then we flatten them along the spatial dimension to obtain the 2D embeddings $\widetilde{Q}$, $\widetilde{K}$, and $\widetilde{V} \in \mathbb{R}^{(HW/16) \times C}$, and calculate the fused features as:

$$\widetilde{M}_i = softmax \left( \frac{\widetilde{Q}\widetilde{K}^T}{\sqrt{H/4 \times W/4 \times C}} \right) \widetilde{V}. \tag{3}$$

Finally, we employ a residual connection by adding $\widetilde{M}_i$ and $T_i$, and use a $3 \times 3$ convolution to obtain the output $M_i$:

$$M_i = conv_{3 \times 3}(\widetilde{M}_i + T_i), \tag{4}$$

where $M_i \in \mathbb{R}^{H/4 \times W/4 \times C}$ combines the forgery features in both RGB domain and frequency domain.

We stack the multi-scale transformer, frequency filter and CMF block for $N$ times ($N = 4$ in this paper). The integrated features $M_{out}$ are calculated by iterative update. Finally, we pass $M_{out}$

through several convolutional layers to obtain global semantic features $f$.

## 3.4 Loss functions

**Cross-entropy loss**. We input $f$ to several fully-connected layers to predict whether the input image is real or fake using a cross-entropy loss $\mathcal{L}_{cls}$:

$$\mathcal{L}_{cls} = y log\hat{y} + (1-y)log(1-\hat{y}), \quad (5)$$

where $y$ is set to 1 if the face image has been manipulated, otherwise it is set to 0; $\hat{y}$ denotes the predicted label by our network.

**Segmentation loss** It is worth noting using a binary classifier tends to result in overfitted models. We additionally predict the face region as an auxiliary task to enrich the supervision for training the networks. Specifically, we input the feature map $M_{out}$ to a decoder (stacked convolutional layers and interpolation upsampling layers) to produce a binary mask $\hat{M}$ in $\mathbb{R}^{H \times W}$:

$$\mathcal{L}_{seg} = \sum_{i,j} M_{i,j} log\hat{M}_{i,j} + (1-M_{i,j})log(1-\hat{M}_{i,j}), \quad (6)$$

where $M_{i,j}$ is the ground-truth mask, with 1 indicating the manipulated pixels and 0 otherwise.

**Contrastive loss** Deepfake images generated by different facial manipulation methods differ in forgery patterns, while the distribution of real images is relatively stable. To improve the generalization ability of our detection model, we first calculate the feature centers of $N_p$ real samples $C_{pos} = \frac{1}{N_p}\sum_{i=1}^{N_p} f_i^{pos}$ and additionally use a contrastive loss to make features from pristine samples to be closer towards the feature center than manipulated samples. Formally, the contrastive loss is defined as:

$$L_{con} = \frac{1}{N_p}\sum_{i=1}^{N_p} d(f_i^{pos}, C_{pos}) - \frac{1}{N_n}\sum_{i=1}^{N_n} d(f_i^{neg}, C_{pos}), \quad (7)$$

where $N_n$ denotes the number of negative samples, and $d$ computes distance with cosine similarity. Finally, combining Eqn. 5, Eqn. 6 and Eqn. 7, the training objective can be written as:

$$\mathcal{L} = \mathcal{L}_{cls} + \lambda_1 \mathcal{L}_{seg} + \lambda_2 \mathcal{L}_{con}, \quad (8)$$

where $\lambda_1$ and $\lambda_2$ are the balancing hyperparameters. By default, we set $\lambda_1 = 1$ and $\lambda_2 = 0.001$.

## 4 SR-DF DATASET

To stimulate research for Deepfake forensics, we introduce a large-scale and challenging dataset, SR-DF. SR-DF is built upon the pristine videos in the FF++ dataset, which contains a diverse set of samples in different genders, ages, and ethnic groups. We first crop face regions in each video frame using [29], and then generate forged videos using state-of-the-art Deepfake generation techniques. Finally, we use the image harmonization method in [8] for post-processing. Below, we introduce these steps in detail.

## 4.1 Dataset Construction

**Synthesis Approaches** To guarantee the diversity of synthesized images, we use four facial manipulation methods, including two face swapping methods: **FSGAN** [44] and **FaceShifter** [19], and two facial reenactment methods: **First-order-motion** [51] and **IcFace** [59]. Note that the manipulation methods we leverage are all identity-agnostic—they can be applied to arbitrary face images without training in pairs, which is different from the FF++ [49] dataset. The detailed forgery images generation process is described below.

**FSGAN** [44] follows the following pipeline to swap the faces of the source image $I_s$ to that of the target image $I_t$. First, the swap generator $G_r$ estimates the swapped face $I_r$ and its segmentation mask $S_r$ based on $I_t$ and a heatmap encoding the facial landmark of $I_s$, while $G_s$ estimates the segmentation mask $S_s$ of the source image $I_s$. Then the inpainting generator $G_c$ inpaints the missing parts of $I_r$ based on $S_s$ to estimate the complete swapped face $Ic$. Finally, using the segmentation mask $S_s$, the blending generator $G_b$ blends $I_c$ and $I_s$ to generate the final output $I_b$ which preserves the posture of $I_s$ but owns the identity of $I_t$. For our dataset, we directly use the pretrained model provided by [44] and inference on our pristine videos.

**FaceShifter** [33] consists of two networks for full pipeline: AEI-Net for face swapping, and HEAR-Net for occlusion handling. As the author of [33] have not public their code, we use the code from [19] who only implements AEI-Net, and we train the model on our data. Specifically, AEI-Net is composed of three components: 1) an Identity Encoder which adopts a pretrained state-of-the-art face recognition model to provide representative identity embeddings. 2) a Multi-level Attributes Encoder which encodes the features of facial attributes. 3) an AAD-Generator which integrates the information of identity and attributes in multiple feature levels and generates the swapped faces. We use the parameters declared in [33] to train the model. **First-order-motion** [51] decouples appearance and

motion information for subject-agnostic facial reenactment. Their framework consists of two main modules: the motion estimation module which uses a set of learned keypoints along with their local affine transformations to predict a dense motion field, and an image generation module which combines the appearance extracted from the source image and the motion derived from the driving video to model the occlusions arising during target motions. To process our dataset, we use the pretrained model on VoxCeleb dataset [40], which contains speech videos from speakers spanning a wide range of different ethnic groups, accents, professions and ages, and reenact the faces in our real videos.

**IcFace** [59] is a generic face animator that is able to transfer the expressions from a driving image to a source image. Specifically, the generator $G_N$ takes the source image and neutral facial attributes as input and produces the source identity with central pose and neutral expression. Then the generator $G_A$ takes the neutral image and attributes extracted from the driving image as an input and produces an image with the source identity and driving image's attributes. We train the complete model on our real videos in a self-supervised manner, using the parameters that they use to train on VoxCeleb dataset [40].

**Post-processing** In order to resolve the color mismatch between the face regions and the background and to eliminate the stitched

boundaries, we use DoveNet [8] for post-processing, which is a state-of-the-art image harmonization method to make the foreground compatible with the background. Note that the masks that we use to distinguish foreground and background are generated using a face parsing model [18].

## 4.2 Comparisons to current Deepfake Datasets

We summarize the basic information of these existing Deefake datasets and our SR-DF dataset in Table 1. In addition, as mentioned above, how to measure the quality of forged images in these datasets remains under-explored. Therefore, we introduce a variety of quantitative metrics to benchmark the quality of current datasets from four perspectives: identity retention, authenticity, temporal smoothness, and diversity. To the best of our knowledge, this is the most comprehensive evaluation system to measure the quality of Deepfake datasets.

**Table 1: A comparison of SR-DF dataset with existing datasets for Deepfake detection. LQ: low-quality, HQ: high-quality.**

| Dataset | Real | | Forged | |
|---|---|---|---|---|
| | Video | Frame | Video | Frame |
| UADFV [68] | 49 | 17.3k | 49 | 17.3k |
| DF-TIMIT-LQ [30] | 320 | 34.0k | 320 | 34.0k |
| DF-TIMIT-HQ [30] | 320 | 34.0k | 320 | 34.0k |
| FF++ [49] | 1,000 | 509.9k | 4000 | 1,830.1k |
| DFD [10] | 363 | 315.4k | 3,068 | 2,242.7k |
| DFDC [15] | 1,131 | 488.4k | 4,113 | 1,783.3k |
| WildDeepfake [74] | 3,805 | 440.5k | 3,509 | 739.6k |
| Celeb-DF [37] | 590 | 225.4k | 5,639 | 2,116.8k |
| ForgeryNet [23] | 99.6k | 1,438.2k | 121.6k | 1457.9k |
| **SR-DF (ours)** | 1,000 | 509.9k | 4,000 | 2,078.4k |

**Mask-SSIM** First, we follow [37] to adopt the Mask-SSIM score as a measurement of synthesized Deepfake images. Mask-SSIM refers to the SSIM score between the face regions of the forged image and the corresponding original image. We use the FaceParsing[18] to generate facial masks and compute the Mask-SSIM on our face swapping subsets. Table 2 demonstrates the average Mask-SSIM scores of all compared datasets, and SR-DF dataset achieves the highest scores.

**Table 2: Average Mask-SSIM scores, perceptual loss, and $E_{warp}$ values of different Deepfake datasets, with the higher value corresponding to better image quality. For perceptual loss, lower value indicates the better image quality, and for $E_{warp}$, lower value corresponding to smoother temporal results.**

| Dataset | FF++ | DFD | DFDC | Celeb-DF | Ours |
|---|---|---|---|---|---|
| **Mask-SSIM ↑** | 0.82 | 0.86 | 0.85 | 0.91 | **0.92** |
| **Perceptual Loss ↓** | 0.67 | 0.69 | 0.63 | **0.59** | 0.60 |
| **$E_{warp}$ ↓** | 73.16 | 69.53 | - | **49.10** | 56.95 |

**Perceptual Loss** *Perceptual loss* is usually used in face inpainting approaches [41, 69] to measure the similarity between the restored faces and corresponding complete faces. Inspired by this, we use the $relu1\_1$, $relu2\_1$, $relu3\_1$, $relu4\_1$ and $relu5\_1$ of the pretrained VGG-19 network on ImageNet [12] to calculate the perceptual loss between the feature maps of forged faces and that of corresponding real faces. We use the dlib [29] to crop the facial regions. We compare the perceptual loss of different datasets in Table 2. Although the perceptual loss of SR-DF dataset is slightly higher than that of Celeb-DF, it is lower than other datasets by a large margin.

**Ewarp** The warping error $E_{warp}$ is used by [4, 25, 32] to measure the temporal inconsistency for video style transfer. We use it to compute the $E_{warp}$ of consecutive forged frames in different datasets to quantitatively measure the short-term consistency. Following [32], we use the method in [50] to calculate occlusion map and PWC-Net [53] to obtain optical flow. $E_{warp}$ of different Deepfake datasets are shown in Table 2 for comparison.

**Feature Space Distribution** As can be seen from the above, Celeb-DF dataset [37] has a decent performance in visual quality. However, they only used one face swapping method to generate all the forged images, which results in limited diversity of data distribution. We illustrate this by visualizing the feature space of Celeb-DF [37], FF++ dataset [49], and SR-DF in Figure 3. We can see the data distribution of the Celeb-DF dataset is more concentrated, while the real and forged images of FF++ dataset can be easily separated in the feature space. On the other hand, the data in SR-DF dataset are more scattered in the 2D space.



**(a) Celeb-DF**       **(b) FF++**       **(c) SR-DF (ours)**
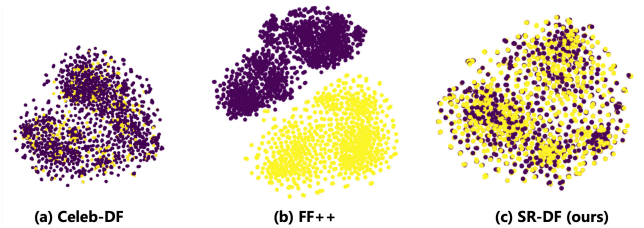
**Figure 3: A feature perspective comparison of Celeb-DF, FF++ dataset (RAW) and SR-DF dataset. We use an ImageNet-pretrained ResNet-18 network to extract features and t-SNE [60] for dimension reduction. Note that we only select one frame in each video for visualization.**

## 5 EXPERIMENTS

## 5.1 Experimental Settings

**Datasets** We conduct experiments on FaceForensics++ (FF++) [49], Celeb-DF [37], and the proposed SR-DF dataset, and ForgeryNet [23]. FF++ consists of 1,000 original videos with real faces, in which 720 videos are used for training, 140 videos are reserved for validation and 140 videos for testing. Each video is manipulated by four Deepfake methods, *i.e.*, Deepfakes [11], FaceSwap [20], Face2Face [57], and NeuralTextures [56]. Different degrees of compression are implemented on both real and forged images to produce high-quality (HQ) version and low-quality (LQ) version of FF++, respectively.

**Table 3: Quantitative frame-level detection results on FaceForensics++ dataset under all quality settings. The best results are marked in bold.**

| Methods | LQ | | HQ | | RAW | |
|---|---|---|---|---|---|---|
| | ACC | AUC | ACC | AUC | ACC | AUC |
| Steg.Features [21] | 55.98 | - | 70.97 | - | 97.63 | - |
| LD-CNN [9] | 58.69 | - | 78.45 | - | 98.57 | - |
| MesoNet [1] | 70.47 | - | 83.10 | - | 95.23 | - |
| Face X-ray [34] | - | 61.60 | - | 87.40 | - | - |
| $F^3$-Net [46] | 90.43 | 93.30 | 97.52 | 98.10 | **99.95** | 99.80 |
| MaDD [71] | 88.69 | 90.40 | 97.60 | 99.29 | - | - |
| M2TR (Xception) | 91.07 | 94.25 | 97.28 | 99.05 | 98.79 | 99.24 |
| M2TR (Ours) | **92.89** | **95.31** | **97.93** | **99.51** | 99.50 | **99.92** |

**Table 4: Frame-level AUC scores (%) of various Deepfake detection methods on Celeb-DF and SR-DF dataset.**

| Methods | Celeb-DF [37] | SR-DF |
|---|---|---|
| Xception [49] | 97.6 | 88.2 |
| Multi-task [42] | 90.5 | 85.7 |
| Capsule [43] | 93.2 | 81.5 |
| DSW-FPA [35] | 94.8 | 86.6 |
| DCViT [65] | 97.2 | 87.9 |
| Ours | 99.8 | 91.2 |
| **Avg** | 95.5 | 86.7 |

Celeb-DF consists of 890 real videos and 5,639 Deepfake videos, in which 6,011 videos are used for training and 518 videos are for testing. ForgeryNet dataset is constructed by 15 manipulation approaches and 36 kinds of distortions for post-processing. It contains 99,630 real videos and 121,617 fake videos. We follow ForgeryNet [23] to train our model on the training set, and evaluate on the validation set. For SR-DF, we build on the 1,000 original videos in FF++, and generate 4,000 forged videos using four state-of-the-art subject-agnostic Deepfake generation techniques (see details above). We use the same training, validation and test set partitioning as FF++.

When training on FF++ dataset and SR-DF dataset, following [46, 71], we augment the real images four times by repeated sampling to balance the number of real and fake samples. For FF++, we sample 270 frames from each video, following the setting in [46, 49].

**Evaluation Metrics** We apply the Accuracy score (Acc) and Area Under the RoC Curve (AUC) as our evaluation metrics, which are commonly used in Deepfake detection tasks [34, 43, 46, 49, 71].

**Implementation Details** We use RetinaFace [13] to crop the face regions (detected boxes enlarged 1.3×) as inputs with a size of 320 × 320. The patch sizes in Sec.4.2.1 are set to (80 × 80), (40 × 40), (20 × 20), and (10 × 10). For backbone network, we use Efficient-b4 [55] pretrained on ImageNet [12]. We use Adam for optimization with a learning rate of 0.0001. The learning rate is decayed 10 times every 40 steps. We set the batch size to 24, and train the complete network for 90 epochs. We will release code upon publication.

## 5.2 Evaluation on FaceForensics++

FF++ [49] is a widely used dataset in various Deepfake detection approaches [1, 9, 21, 34, 46, 71]. We compare M2TR with top-notch methods on it, including: Steg. Features [21], LD-CNN [9], MesoNet [1], Face X-ray [34], $F^3$-Net [46], and MaDD [71].

We test the frame-level detection performance on RAW, HQ, and LQ of FF++, respectively, and report the AUC scores (%) in Table 3. We can see that our method achieves state-of-the-art performance on all versions (*i.e.*, LQ, HQ, and RAW), which suggests the effectiveness of our approach in detecting Deepfakes of different visual qualities. We also evaluate M2TR using different backbones, and the performances verifies that our framework is not restricted by the backbone networks. Comparing across different versions of the FF++ dataset, we see that while most approaches achieve high

performance on the high-quality version of FF++, we observe a significant performance degradation on FF++ (LQ) where the forged images are compressed. This could be remedied by leveraging frequency information. While both F3-Net and M2TR use frequency features, M2TR achieves an accuracy of 92.35% in the LQ setting, outperforming the F3-Net approach by 1.92%.

## 5.3 Evaluation on Celeb-DF and SR-DF

In this section, we conduct experiments to evaluate the detection accuracy of our M2TR on Celeb-DF [37] and SR-DF dataset at frame-level, respectively. Note that we do not report the quantitative results of certain state-of-the-art Deepfake detection methods including [34, 46, 71] because the code and models are not publicly available. The results are reported in Table 4. We observe that our M2TR achieves 99.9% and 90.5% on Celeb-DF and SR-DF, respectively, which demonstrate that our method outperforms all the other Deepfake detection methods over different datasets. This suggests that our approach is indeed effective for Deepfake detection across different datasets.

In addition, the quality of different Deepfake datasets can be evaluated by comparing the detection accuracy of the same detection method on different datasets. Given that Celeb-DF [37] contains high-quality samples (as discussed in 4.2, Celeb-DF achieves the best results on *Mask-SSIM*, *Perceptual loss* and *Ewarp* metrics in the available Deepfake dataset.), we calculate the average frame-level AUC scores of all compared detection methods on Celeb-DF dataset and SR-DF, and report them in the last row of Table 4. The overall performance on SR-DF is 9.2% lower than that of Celeb-DF, which demonstrates that SR-DF is more challenging.

## 5.4 Evaluation on ForgeryNet

ForgeryNet [23] is the most recently released largest scale deepfake detection dataset, which provides three types of forgery labels, *i.e.*, two-way (real / fake), three-way (real / fake with identity-replaced forgery approaches / fake with identity-remained forgery approaches), and n-way (n = 16, real and 15 respective forgery approaches) labels. Using the rich annotations, we conduct two-/three-/n-way classification experiments. The comparison results in Table 5 demonstrate that M2TR outperforms state-of-the-art methods on classification tasks of different granularity. In particular, the most remarkable performance gain *i.e.*, 15.5% compared with

**Table 5: Quantitative comparison of various Deep-fake detection methods on ForgeryNet dataset.**

| Methods | 2 way | 3 way | 16 way |
|---|---|---|---|
| MobileNetV3 [24] | 76.24 | - | - |
| Xception [7] | 80.78 | 73.00 | 58.81 |
| F$^3$-Net [46] | 80.86 | 74.45 | 59.82 |
| GramNet [38] | 80.89 | 73.30 | 56.77 |
| SNRFilters-Xception [5] | 81.09 | - | - |
| Ours | **82.52** | **75.12** | **69.12** |

**Table 6: Quantitative video-level detection results on FF++ dataset and SR-DF dataset. M2TR $_{mean}$ denotes averaging the extracted features obtained by M2TR for all frames as the video-level representation, while M2TR $_{vtf}$ denotes using VTF Block for temporal fusion.**

| Method | FF++ (RAW) | FF++ (HQ) | FF++ (LQ) | SR-DF |
|---|---|---|---|---|
| P3D [47] | 80.9 | 75.23 | 67.05 | 65.97 |
| R3D [67] | 96.15 | 95.00 | 87.72 | 73.24 |
| I3D [3] | 98.23 | 96.70 | 93.18 | 80.11 |
| M2TR $_{mean}$ | 99.06 | 98.23 | 93.95 | 82.27 |
| ST-M2TR | **99.87** | **99.42** | **95.31** | **85.32** |

**Table 7: AUC scores (%) for cross-dataset evaluation on FF++, Celeb-DF, and SR-DF datasets. Note that some methods have not made their code public, so we directly use the data reported in their paper. "−" denotes the results are unavailable.**

| Training Set | Testing Set | Xception [49] | Multi-task [42] | Capsule [43] | DSW-FPA [35] | Two-Branch [39] | F3-Net [46] | MaDD [71] | DCViT [65] | Ours |
|---|---|---|---|---|---|---|---|---|---|---|
| FF++ | FF++ | 99.7 | 76.3 | 96.6 | 93.0 | 98.7 | 98.1 | 99.3 | 98.3 | **99.5** |
| | Celeb-DF | 48.2 | 54.3 | 57.5 | 64.6 | **73.4** | 65.2 | 67.4 | 60.8 | 68.2 |
| | SR-DF | 37.9 | 38.7 | 41.3 | 44.0 | - | - | - | 57.8 | **63.7** |
| SR-DF | SR-DF | 88.2 | 85.7 | 81.5 | 86.6 | - | - | - | 87.9 | **91.2** |
| | FF++ | 63.2 | 58.9 | 60.6 | 69.1 | - | - | - | 62.6 | **79.7** |
| | Celeb-DF | 59.4 | 51.7 | 52.1 | 62.9 | - | - | - | 63.7 | **82.1** |

F$^3$-Net, is achieved on the most challenging 16-way classification, which is benefited from the capability of our method to identify the multi-scale forgery features.

## 5.5 From Frames to Videos

Existing methods on Deepfake detection mainly perform evaluation based on frames extracted from videos, albeit videos are provided. However, in real-world scenarios, most Deepfake data circulating on the Internet are fake videos, therefore, we also conduct experiments to evaluate our M2TR on video-level Deepfake detection. The most significant difference between videos and images is the additional temporal information between frame sequences. We demonstrate that M2TR can be easily extended for video modeling by adding a temporal transformer to combine frame-level features generated by M2TR. We refer to such an extension as spatial-temporal M2TR (ST-M2TR).

In particular, we sampled 16 frames at intervals from one video, and directly use the model trained at the frame-level to extract features of different frames. These features are then input to a transformer block (it has 4 stacked encoders, each with 8 attention heads, and an MLP head that has two fc layers) to obtain video-level predictions. We report the AUC scores (%) and compare with (1) P3D [47], which simplifies 3D convolutions with 2D filters on spatial dimension and 1D temporal connections; (2)R3D [67], which encodes the video sequences using a 3D fully convolutional networks and then generates candidate temporal fragments for classification; (3)I3D [3], which expands 2D CNNs with an additional temporal dimension to introduce a two-stream inflated 3D convolutional network; (4) M2TR$_{mean}$, which averages the features of different frames by M2TR for video-level prediction. Note that (1) and (3) are designed for video action recognition, while (2) is for temporal activity detection, and we modify them for video-level Deepfake detection. The results are summarized in Table 6. We can see that

our method achieves the best performance on FF++, ForgeryNet, and SR-DF.

## 5.6 Generalization Ability

The generalization ability is at the core of Deepfake detection. We evaluate the generalization of our M2TR by separately training on FF++ (HQ) and SR-DF dataset, and test on other datasets. We follow [71] to sample 30 frames for each video and calculate the frame-level AUC scores. The comparison results are shown in Table 7. Note that for the Deepfake detection models that are not publicly available, we only use the results reported in their paper. The results in Table 7 demonstrate that our method achieves better generalization than most existing methods.

## 5.7 Ablation Study

**Effectiveness of Different Components** The Multi-scale Transformer (MT) of our method is designed to capture local inconsistency between patches of different sizes, while the Frequency Filter (FF) is utilized to capture the subtle forgery traces in frequency domain. To evaluate the effectiveness of MT and FF, we remove them separately from M2TR and demonstrate the performance degradation on FF++. We also replace the Cross Modality Fusion (CMF) blocks with naive concatenation operations to verify the usefulness of CMF for feature fusion.

The quantitative results are listed in Table 8, which validates that the use of MT, FF and CMF can effectively improve the detection performance of our model. In particular, the proposed frequency filter brings a remarkable improvement to our method under the low-quality (LQ) setting, *i.e.*, about 1.7% performance gain on AUC score, which is mainly benefited from the complementary information from the frequency modality.

**Table 8: Ablation results on FF++ (HQ) and FF++ (LQ) with and without Multi-scale Transformer and CMF.**

| Method | LQ | | HQ | |
|---|---|---|---|---|
| | ACC (%) | AUC (%) | ACC (%) | AUC (%) |
| w/o MT | 87.19 | 90.05 | 94.88 | 96.94 |
| w/o FF | 89.33 | 90.48 | 95.89 | 97.94 |
| w/o CMF | 90.70 | 92.37 | 96.78 | 98.19 |
| Ours | **92.89** | **95.31** | **97.93** | **99.51** |

**Table 9: Ablation results on FF++ (HQ) using multi-scale Transformer (MT) or single-scale transformer.**

| Patch size | ACC (%) | AUC (%) |
|---|---|---|
| $40 \times 40$ | 93.58 | 95.81 |
| $20 \times 20$ | 96.65 | 97.53 |
| $10 \times 10$ | 97.19 | 98.83 |
| Ours | **97.93** | **99.51** |

**Effectiveness of the Multi-scale Design** To verify the effectiveness of using multi-scale patches in different heads in our multi-scale transformer, we replace MT with several single-scale transformers with different patch sizes, and conduct experiments on FF++ (HQ). The results in Table 9 demonstrate that our full model achieves the best performance with MT, *i.e.*, 3.9%, 1.9%, and 0.7% higher than $40 \times 40$, $20 \times 20$ and $10 \times 10$ single-scale transformer on AUC score. This confirms the use of a multi-scale transformer is indeed effective.

**Effectiveness of the Contrastive Loss** To illustrate the contribution of contrastive loss in improving the generalization ability of our method, we also conduct experiments to train M2TR without its supervision and evaluate the cross-dataset detection accuracy. The comparison results are reported in Table 10. We can observe that 1) When training on FF++ without the contrastive loss, the accuracy decreases by 3.8% and 5.2% in Celeb-DF and SR-DF, respectively. 2) When training on SR-DF dataset without the contrastive loss, the accuracy decreases by 5.8% and 3.0%, respectively.

**Table 10: AUC (%) for cross-dataset evaluation on FF++ (HQ), Celeb-DF, and SR-DF with (denoted as M2TR) and without (denoted as M2TR $_{ncl}$) the supervision of contrastive loss.**

| Training Set | Testing Set | M2TR $_{ncl}$ | M2TR |
|---|---|---|---|
| FF++ | Celeb-DF | 65.6 | 68.2 |
| | SR-DF | 60.4 | 63.7 |
| SR-DF | FF++ | 75.1 | 79.7 |
| | Celeb-DF | 79.6 | 82.1 |

## 6 CONCLUSION

In this paper, we presented a two-stream network Multi-modal Multi-scale Transformer (M2TR) for Deepfake detection, which uses multi-scale transformers to capture subtle local inconsistency at multiple scales and frequency filters to improve the robustness against image compression. Forgery feature from two streams are adaptively fused through cross modality fusion blocks. Besides, we introduced a challenging dataset SR-DF that are generated with several state-of-the-art face swapping and facial reenactment methods. We also built the most comprehensive evaluation system to quantitatively verify that the SR-DF dataset is better than existing datasets in terms of visual quality and data diversity. Extensive experiments on different datasets demonstrate the effectiveness of the proposed method.

## REFERENCES

[1] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. 2018. Mesonet: a compact facial video forgery detection network. In *WIFS*.
[2] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. 2020. End-to-end object detection with transformers. In *ECCV*.
[3] Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*.
[4] Dongdong Chen, Jing Liao, Lu Yuan, Nenghai Yu, and Gang Hua. 2017. Coherent online video style transfer. In *ICCV*.
[5] Mo Chen, Vahid Sedighi, Mehdi Boroumand, and Jessica Fridrich. 2017. JPEG-phase-aware convolutional neural network for steganalysis of JPEG images. In *IHMSW*.
[6] Shen Chen, Taiping Yao, Yang Chen, Shouhong Ding, Jilin Li, and Rongrong Ji. 2021. Local Relation Learning for Face Forgery Detection. In *AAAI*.
[7] François Chollet. 2017. Xception: Deep learning with depthwise separable convolutions. In *CVPR*.
[8] Wenyan Cong, Jianfu Zhang, Li Niu, Liu Liu, Zhixin Ling, Weiyuan Li, and Liqing Zhang. 2020. DoveNet: Deep Image Harmonization via Domain Verification. In *CVPR*.
[9] Davide Cozzolino, Giovanni Poggi, and Luisa Verdoliva. 2017. Recasting residual-based local descriptors as convolutional neural networks: an application to image forgery detection. In *Workshop on IH&MMSec*.
[10] DeepFake Detection Dataset. 2019. https://ai.googleblog.com/2019/09/contributing-data-to-deepfake-detection.html.
[11] Deepfakes. 2018. github. https://github.com/deepfakes/faceswap.
[12] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *CVPR*.
[13] Jiankang Deng, Jia Guo, Evangelos Ververas, Irene Kotsia, and Stefanos Zafeiriou. 2020. Retinaface: Single-shot multi-level face localisation in the wild. In *CVPR*.
[14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL-HLT*.
[15] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. 2020. The deepfake detection challenge (dfdc) dataset. *arXiv preprint arXiv:2006.07397* (2020).
[16] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *ICLR*.
[17] Ricard Durall, Margret Keuper, Franz-Josef Pfreundt, and Janis Keuper. 2019. Unmasking deepfakes with simple features. *arXiv preprint arXiv:1911.00686* (2019).
[18] Face-parsing. 2019. github. https://github.com/zllrunning/face-parsing.PyTorch.
[19] FaceShifter. 2020. github. https://github.com/mindslab-ai/faceshifter.
[20] Faceswap. 2018. github. https://github.com/MarekKowalski/FaceSwap/.
[21] Jessica Fridrich and Jan Kodovsky. 2012. Rich models for steganalysis of digital images. *TIFS* (2012).
[22] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *CVPR*.
[23] Yinan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. 2021. ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis. In *CVPR*.

[24] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. 2019. Searching for mobilenetv3. In *ICCV*.

[25] Haozhi Huang, Hao Wang, Wenhan Luo, Lin Ma, Wenhao Jiang, Xiaolong Zhu, Zhifeng Li, and Wei Liu. 2017. Real-time neural style transfer for videos. In *CVPR*.

[26] Ying Huang, Wenwei Zhang, and Jinzhuo Wang. 2020. Deep frequent spatial temporal learning for face anti-spoofing. *arXiv preprint arXiv:2002.03723* (2020).

[27] Hyeonseong Jeon, Youngoh Bang, and Simon S Woo. 2020. FDFtNet: Facing off fake images using fake detection fine-tuning network. In *ICT Systems Security and Privacy Protection*.

[28] Ira Kemelmacher-Shlizerman. 2016. Transfiguring portraits. *ACM TOG* (2016).

[29] Davis E King. 2009. Dlib-ml: A machine learning toolkit. *JMLR* (2009).

[30] Pavel Korshunov and Sébastien Marcel. 2018. Deepfakes: a new threat to face recognition? assessment and detection. *arXiv preprint arXiv:1812.08685* (2018).

[31] Mohammad Rami Koujan, Michail Christos Doukas, Anastasios Roussos, and Stefanos Zafeiriou. 2020. Head2head: Video-based neural head synthesis. *arXiv preprint arXiv:2005.10954* (2020).

[32] Chenyang Lei, Yazhou Xing, and Qifeng Chen. 2020. Blind Video Temporal Consistency via Deep Video Prior. In *NIPS*.

[33] Lingzhi Li, Jianmin Bao, Hao Yang, Dong Chen, and Fang Wen. 2019. Faceshifter: Towards high fidelity and occlusion aware face swapping. *arXiv preprint arXiv:1912.13457* (2019).

[34] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. 2020. Face x-ray for more general face forgery detection. In *CVPR*.

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

[36] Yuezun Li and Siwei Lyu. 2019. Exposing DeepFake Videos By Detecting Face Warping Artifacts. In *CVPRW*.

[37] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. 2020. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *CVPR*.

[38] Zhengzhe Liu, Xiaojuan Qi, and Philip HS Torr. 2020. Global texture enhancement for fake face detection in the wild. In *CVPR*.

[39] Iacopo Masi, Aditya Killekar, Royston Marian Mascarenhas, Shenoy Pratik Gurudatt, and Wael AbdAlmageed. 2020. Two-branch recurrent network for isolating deepfakes in videos. In *ECCV*.

[40] Arsha Nagrani, Joon Son Chung, and Andrew Zisserman. 2017. Voxceleb: a large-scale speaker identification dataset. *arXiv preprint arXiv:1706.08612* (2017).

[41] Kamyar Nazeri, Eric Ng, Tony Joseph, Faisal Qureshi, and Mehran Ebrahimi. 2019. EdgeConnect: Structure Guided Image Inpainting using Edge Prediction. In *ICCVW*.

[42] Huy H Nguyen, Fuming Fang, Junichi Yamagishi, and Isao Echizen. 2019. Multitask learning for detecting and segmenting manipulated facial images and videos. In *BTAS*.

[43] Huy H Nguyen, Junichi Yamagishi, and Isao Echizen. 2019. Use of a capsule network to detect fake images and videos. *arXiv preprint arXiv:1910.12467* (2019).

[44] Yuval Nirkin, Yosi Keller, and Tal Hassner. 2019. Fsgan: Subject agnostic face swapping and reenactment. In *ICCV*.

[45] Albert Pumarola, Antonio Agudo, Aleix M Martinez, Alberto Sanfeliu, and Francesc Moreno-Noguer. 2018. Ganimation: Anatomically-aware facial animation from a single image. In *ECCV*.

[46] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. 2020. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *ECCV*.

[47] Zhaofan Qiu, Ting Yao, and Tao Mei. 2017. Learning spatio-temporal representation with pseudo-3d residual networks. In *ICCV*.

[48] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. Exploring the Limits of Transfer Learning with a Unified Text-to-Text Transformer. *JMLR* (2020).

[49] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. 2019. Faceforensics++: Learning to detect manipulated facial images. In *ICCV*.

[50] Manuel Ruder, Alexey Dosovitskiy, and Thomas Brox. 2016. Artistic style transfer for videos. In *GCPR*.

[51] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2020. First order motion model for image animation. *arXiv preprint arXiv:2003.00196* (2020).

[52] Karen Simonyan and Andrew Zisserman. [n.d.]. Very deep convolutional networks for large-scale image recognition. In *ICLR*.

[53] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. 2018. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*.

[54] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. 2017. Synthesizing obama: learning lip sync from audio. *ACM TOG* (2017).

[55] Mingxing Tan and Quoc Le. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*.

[56] Justus Thies, Michael Zollhöfer, and Matthias Nießner. 2019. Deferred neural rendering: Image synthesis using neural textures. *ACM TOG* (2019).

[57] Justus Thies, Michael Zollhofer, Marc Stamminger, Christian Theobalt, and Matthias Niessner. 2016. Face2Face: Real-Time Face Capture and Reenactment of RGB Videos. In *CVPR*.

[58] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. 2021. Training data-efficient image transformers & distillation through attention. In *ICML*.

[59] Soumya Tripathy, Juho Kannala, and Esa Rahtu. 2020. Icface: Interpretable and controllable face reenactment using gans. In *WACV*.

[60] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *JMLR* (2008).

[61] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NIPS*.

[62] Junke Wang, Zuxuan Wu, Jingjing Chen, Xintong Han, Abhinav Shrivastava, Ser-Nam Lim, and Yu-Gang Jiang. 2022. ObjectFormer for Image Manipulation Detection and Localization. In *CVPR*.

[63] Junke Wang, Xitong Yang, Hengduo Li, Zuxuan Wu, and Yu-Gang Jiang. 2021. Efficient Video Transformers with Spatial-Temporal Token Selection. *arXiv preprint arXiv:2111.11591* (2021).

[64] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. 2020. CNN-generated images are surprisingly easy to spot... for now. In *CVPR*.

[65] Deressa Wodajo and Solomon Atnafu. 2021. Deepfake Video Detection Using Convolutional Vision Transformer. *arXiv preprint arXiv:2102.11126* (2021).

[66] Tete Xiao, Piotr Dollar, Mannat Singh, Eric Mintun, Trevor Darrell, and Ross Girshick. 2021. Early Convolutions Help Transformers See Better. In *NIPS*.

[67] Huijuan Xu, Abir Das, and Kate Saenko. 2017. R-c3d: Region convolutional 3d network for temporal activity detection. In *ICCV*.

[68] Xin Yang, Yuezun Li, and Siwei Lyu. 2019. Exposing deep fakes using inconsistent head poses. In *ICASSP*.

[69] Yang Yang and Xiaojie Guo. 2020. Generative Landmark Guided Face Inpainting. In *PRCV*.

[70] Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. 2020. Resnest: Split-attention networks. *arXiv preprint arXiv:2004.08955* (2020).

[71] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. 2021. Multi-attentional Deepfake Detection. In *CVPR*.

[72] Peng Zhou, Xintong Han, Vlad I Morariu, and Larry S Davis. 2017. Two-stream neural networks for tampered face detection. In *CVPRW*.

[73] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. 2021. Deformable {DETR}: Deformable Transformers for End-to-End Object Detection. In *ICLR*.

[74] Bojia Zi, Minghao Chang, Jingjing Chen, Xingjun Ma, and Yu-Gang Jiang. 2020. WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection. In *ACM MM*.