

# Self-supervised Transformer for Deepfake Detection

Hanqing Zhao<sup>1</sup> Wenbo Zhou<sup>1</sup> Dongdong Chen<sup>2</sup>  
Weiming Zhang<sup>1</sup> Nenghai Yu<sup>1</sup>

University of Science and Technology of China<sup>1</sup> Microsoft Cloud AI<sup>2</sup>

{zhq2015@mail, welbeckz@, zhangwm@, ynh@}.ustc.edu.cn

cddlyf@gmail.com

## Abstract

*The fast evolution and widespread of deepfake techniques in real-world scenarios require stronger generalization abilities of face forgery detectors. Some works capture the features that are unrelated to method-specific artifacts, such as clues of blending boundary, accumulated up-sampling, to strengthen the generalization ability. However, the effectiveness of these methods can be easily corrupted by post-processing operations such as compression. Inspired by transfer learning, neural networks pre-trained on other large-scale face-related tasks may provide useful features for deepfake detection. For example, lip movement has been proved to be a kind of robust and good-transferring high-level semantic feature, which can be learned from the lipreading task. However, the existing method pre-trains the lip feature extraction model in a supervised manner, which requires plenty of human resources in data annotation and increases the difficulty of obtaining training data. In this paper, we propose a self-supervised transformer based audio-visual contrastive learning method. The proposed method learns mouth motion representations by encouraging the paired video and audio representations to be close while unpaired ones to be diverse. After pre-training with our method, the model will then be partially fine-tuned for deepfake detection task. Extensive experiments show that our self-supervised method performs comparably or even better than the supervised pre-training counterpart.*

## 1. Introduction

Face manipulation technologies [5, 20, 21, 31, 39, 40, 47] empowered by deep generative models are fast advancing which makes deepfake medias more realistic and easily to deceive watchers. The malicious usage and spread of deepfake have raised serious societal concerns and posed an increasing threat to our trust in online media. Therefore, deepfake detection in real-world scenarios becomes an urgent need and obtains a considerable amount of attention in recent years.

To defend against the potential risks of these forged media, numerous efforts have been devoted and achieving promising performances on specific datasets in recent years [1, 18, 22, 25, 28, 34, 36, 44, 45]. Many previous works introduce low-level texture features from different domains for searching the underlying generation artifacts. However, dramatic drops in performance may be experienced when the artifact patterns are changed. And the well-processed forged videos which show only subtle differences from real ones make the deepfake detection in real-world scenarios to become a very tough task.

Some works attempt to strengthen the generalization ability by data augmentation or capturing common clues during the forgery process. For example, SPSL [25] focuses on finding the frequency artifacts caused by the accumulative up-sampling operation. Another effective way is to predict the blending boundaries [22, 46] between the background and the altered inner face regions. Although they achieve impressive performances in cross-data evaluations, they are usually sensitive to post-processing, e.g., video compression. Recently, researchers try to develop robust high-level semantic features. Heliasos et al. [18] find that supervised pre-trained spatio-temporal networks with visual speech recognition (lipreading) tasks can extract robust representations of lip movements for boosting deepfake detection performance. This also indicates that feature extractors pre-trained in other face-related tasks can provide meaningful help for deepfake detection. However, pre-training in other tasks such as lipreading requires precise segmentation and data annotation. It is costly when developing a larger training set.

To address the limitation, we propose a self-supervised transformer for pre-training a robust semantic feature extractor. Generally speaking, there is a strong correlation between audio and lip movement in speech videos. It has also been proved that predicting lip movement by audio (lip synthesis [39]) or predicting audio by lip movement (lip2wav [33]) are feasible. Therefore, we propose the method to self-supervised pre-train instead of supervised pre-train in the lipreading task. We can obtain a general representation of

lip movement by audio-visual consistency using contrastive learning.

To this end, we propose a two-stage spatio-temporal video encoder. In the frontend stage, we use a 3D convolution layer to extract flows from video, followed by a 2D CNN for capturing the local lip movement representations and a temporal 1D transformer backend for long-term lip movement representations. To learn the generic lip movements representation, we design a cross-modal contrastive learning method that utilizes the consistency between the audio channel and the movement of lip regions in frames to train the video encoder. After pre-training, we freeze the front end of the encoder and add an MLP head to fine-tune the whole network for deepfake detection task.

We conduct extensive experiments to compare the performances of our methods with the state-of-the-art methods in various challenging cases. The results demonstrate that our method achieves comparable performance with the state-of-the-art in most cases, and achieves better performance with respect to generalization to unseen forgery datasets. Compared to the supervised pre-trained method, our approach also exhibits better robustness to common corruptions which degrade other models' performance. Further, we investigate the relation between the scale of pre-training data and the average detection AUC. Conclusively, a larger pre-training data scale is indeed helpful to promote the performance of models. Another advantage of our method is significantly reducing the annotation cost for pre-training data.

## 2. Related Work

### 2.1. Deepfake Detection

Since the deepfake technique has caused severe societal concerns, many effective countermeasures [1, 6, 18, 22, 25, 28, 34, 42, 45] have been proposed against it. According to the feature extraction types, current methods can be roughly categorized into two types: textural feature based methods and semantic feature based methods. As the name implies, textural feature based methods focus on capturing low-level textural information from different domains. Two-branch [28] introduced an extra CNN stem with deep Laplacian-of-Gaussian filter in the CNN-RNN architecture for acting as a band-pass filter to amplify artifacts. Patch-forensics [4] proposed a classifier with limited receptive fields that focus on textures in small patches to capture local errors. Multi-attention [45] attempt to introduce the multi-regional attention mechanism into deepfake detection, which is inspired by fine-grained classification.

To boost the generalization ability, Face X-ray [22] and Patch-wise Consistency Learning [46] leverage the clues of image blending between the background and the altered face region. These methods achieve impressive performances in cross-data evaluation. However, they are easily influenced by

video compression and noise disturbance since the boundary clues are highly fragile to post-processing.

Recently, researchers find that high-level semantic features show excellent robustness in dealing with both cross-data tests and post-processing operations. Lipforensics [18] extract the representations of lip movement by using a pre-trained network. The lip feature extractor is pre-trained in a supervised manner on Lipreading in the Wild(LRW) [10] dataset, which is commonly used for the lipreading tasks. The network takes a multi-branch temporal convolution network(MSTCN) [16] as backend and fine-tuned on deepfake datasets. It achieves state-of-the-art generalization ability. However, to pre-train such a robust lip feature extractor requires a large-scale well-annotated dataset, which is extremely costly. In this paper, we propose a self-supervised transformer for pre-training, in this way, we can significantly reduce the annotation cost for pre-training data and obtain good scalability.

### 2.2. Contrastive Learning

Recently, general pre-trained models have been widely used for fine-tuning on downstream tasks. Self-supervised pre-training attracts tremendous attention for its versatility. Self-supervised learning enables us to learn meaningful representation for various classification tasks without relying on labeled data. Contrastive learning is a popular self-supervised strategy that encourages the features of the same instances to be close while features of different instances to be distant. Generally, we can take different views of the same sample data as positive pairs and views from unmatched data as negative pairs to construct a classification task for discriminating instances. For pre-training semantic image representations, [7, 8] obtain different views of the same image by random augmentations. For pre-training video classification models [17], positive pairs are fragments from the same video while negative pairs are fragments from different videos. Besides the uni-modal contrastive learning, natural multi-media data have cross-modal consistency.

Many previous self-supervised learning methods are designed in the cross-modal predictive way, while it is also suitable for using contrastive learning that takes each modal as different views with shared information. For example, Ommer et al. [38] proposed contrastive learning for the frame and optical flow in videos. [35] demonstrate a simple pre-training task, that is, predicting which caption goes with which image is an efficient and scalable way to learn good image representations. Some recent works [27, 29, 48] leverage correspondence between audio and visual signals for learning semantic representations, they achieve strong transfer learning performance on downstream tasks like action recognition and lipreading. Inspired by previous successful works, in this paper, we try to learn a general representation of lip movement by audio-visual consistency using

contrastive learning.

### 3. Proposed Method

#### 3.1. Overview

Most deepfake generation methods synthesize fake faces frame by frame without considering motion coherence, especially in the lip regions where movements are most frequent and complex while talking. Thus the irregularity of lip movement can be a common feature of deepfake video and such a high-level semantic feature is robust to various video post-processing. To learn semantical representations of lip movements, we apply cross-modal contrastive learning by encoding utterance features and lip motion features from videos into a common space. We use instance discrimination learning to enforce corresponding audio and visual features matching each other, which shares a similar spirit with the lipreading task. The learned lip representations can transfer to expose the irregularities in deepfake videos. Benefit from the primitive features of lip movements, the video encoder would be less prone to overfitting non-transferable artifacts when fine-tuning on deepfake detection task.

#### 3.2. Audio-Visual Contrastive Learning pre-train

For cross-modal contrastive learning, our model consists of a video encoder and an audio encoder. In the pre-train stage, each encoder uses an additional MLP head for projecting the raw features into the common space. We take a pre-trained wav2vec2 model [2] for speech recognition as the backbone of the audio encoder. The wave2vec2 model encodes raw audio into a sequence of local wave features by one dimensional CNN layers which are frozen during training and then inference the features by transformer layers to get semantical audio representation array. We apply temporal adaptive average pooling for converting the output feature array from the last transformer layer to a fixed length token sequence, then the sequence is reshaped as the input of the MLP projection head.

We design a spatial-temporal framework as shown in figure Figure 1 for encoding video features. It has two stages, in the frontend stage, a 3D convolutional layer with kernel size [5,7,7] extracts the optical flows from raw videos, then a following 2D ResNet module encodes the local lip movement feature from the optical flows. In the backend stage, we convert the local lip movement feature of each frame into a feature sequence with linear projection, then we use a temporal transformer to get the global semantic video representation. Just like the audio encoder, we also apply temporal adaptive average pooling and MLP projection head for encoding the final video feature.

Previous contrastive learning methods for video representations [12,27,29,32] usually leverage 3D CNN architectures for general video understanding. Compared to 3D CNN, our

architecture separates the process of encoding local lip movement feature and global semantic video representation. Our framework makes the fine-tuning process easier because we can freeze the frontend for keeping the low-level feature representations learned during the pre-training, only fine-tune the backend for capturing the long-term inconsistency of fake videos.

In our framework, we use the transformer as the backend of the encoder. Transformer has succeeded in many NLP and vision tasks since it can capture long-term dependency and preserve less inductive bias than CNN. Thus we believe it is more suitable for large-scale pre-trained models. The shape of output tensor from frontend is noted as  $T \in B \times L \times N$ ,  $B$  is the batch size,  $L$  is a static sequence length with padding,  $N$  is the dimension of pooled spatial feature, for example,  $N$  is 512 for ResNet18 frontend. We use a linear projection for transforming  $T$  to the input features of the transformer, and we add learnable positional embeddings to input features for introducing time relationships. In the experiment part, we validate the model with transformers backend can achieve better performance than those with general CNNs.

We utilize the intrinsic correlation of lip movement and speech audio to learn generic representations of them. Following the contrastive learning manner, we use the InfoNCE [41] loss to distinguish positive audio-visual feature pairs from a bunch of negative pairs. Video encoder  $E_v$  encodes a video clip into video embedding  $m_v$  and similarly audio encoder  $E_a$  encodes an audio clip into audio embedding  $m_a$ . Then  $m_v$  and  $m_a$  are projected to common space by MLP heads as  $z_v$  and  $z_a$ . The loss is shown in equation Equation 3, we respectively take the video features and the audio features as query and the other modal as keys to get loss  $\mathcal{L}_{va}$  and  $\mathcal{L}_{av}$ , then we simply use their average as the contrastive loss.

$$\mathcal{L}_{va} = -\log \frac{\exp(z_v \cdot z_a^+ / \tau)}{\exp(z_v \cdot z_a^+ / \tau) + \sum_{z_a^-} \exp(z_v \cdot z_a^- / \tau)} \quad (1)$$

$$\mathcal{L}_{av} = -\log \frac{\exp(z_a \cdot z_v^+ / \tau)}{\exp(z_a \cdot z_v^+ / \tau) + \sum_{z_v^-} \exp(z_a \cdot z_v^- / \tau)} \quad (2)$$

$$\mathcal{L} = \frac{1}{2} \mathcal{L}_{av} + \frac{1}{2} \mathcal{L}_{va} \quad (3)$$

Here the label  $+$  means positive sample pairs and  $-$  means negative sample pairs,  $\tau$  is the temperature which is set 0.1. For each iteration in the pre-train stage, we sample a mini-batch of video clips, extract segments of synchronized frames and waves with random offset for each video clip as positive audio-visual pairs while all the unrelated items from the mini-batch act as negative pairs.

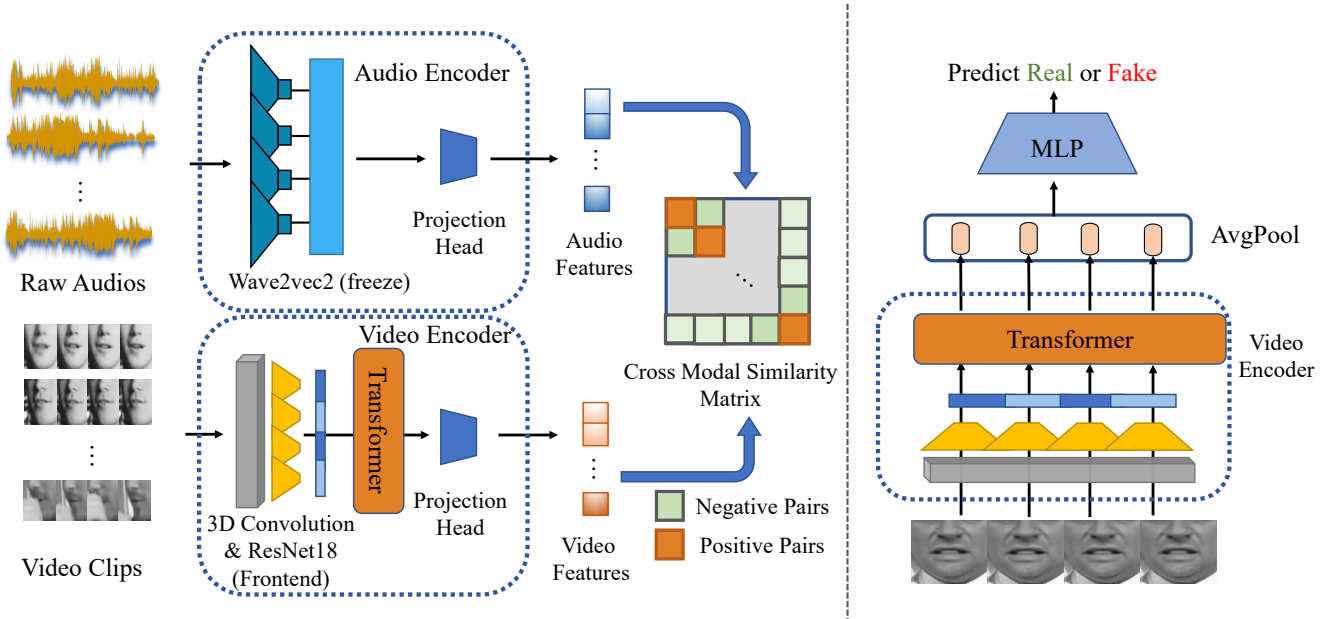


Figure 1. Demonstration of pre-training stage and fine-tuning stage of proposed method. The left part shows the procedure of audio-visual contrastive learning for obtaining generic representation of lip movements. And the right part shows how do we fine-tune the pre-trained video encoder for deepfake detection task.

### 3.3. Fine-tuning on Deepfake Dataset

In the fine-tuning stage, we freeze the parameters in the frontend 3D convolution and ResNet2D for keeping the pre-trained local representations. We extract global features by adaptive average pooling the output of the last transformer layer and adding a new MLP head for classification. To preserve the pre-trained knowledge in transformer layers, we manually control the learning rate of each transformer layer. The MLP module is assigned with an initial learning rate of 0.01, while the last transformer layer is assigned with an initial learning rate of 0.005, each prior transformer layer decays the initial learning rate by 0.9. For deepfake detection task, we use binary cross-entropy loss to fine-tune the classifier.

$$\mathcal{L}_{ce} = y \log\left(\frac{1}{1 + e^{-C(x)}}\right) + (1 - y) \log\left(1 - \frac{1}{1 + e^{-C(x)}}\right) \quad (4)$$

where  $x$  and  $y$  are the input video sequence and corresponding label, respectively.  $C(x)$  represents the predicted logit of the MLP module.

To mitigate the influence of different frame rates, we convert all videos to 25 fps and all audios to a 16kHz sample rate. Then we trace the face of the speaking person with RetinaFace [13] and SyncNet [11]. We crop the selected faces with a tight boundary and use FAN [3] to get the 68 landmarks. We align each face using 5-landmarks and resize the face images to 256x256. We then crop a 96x96 region at the center of the lip landmarks region. After preprocessing,

a video’s lip regions can be represented as  $F_v \in R^{[T,H,W]}$ , where  $T$  is the length in time and  $H, W$  are height and width. In the pre-train stage, we randomly sample a clip with fixed length  $L_p$  and corresponding raw audio with length  $640L_p$  from the video. We apply identical image augmentations on each frame including random crop to 88x88, motion blur, Gaussian noise and random brightness.

The computation complexity of self-attention is the square of the sequence length that forbids us to predict with a very long video sequence. For efficiency, we chunk videos into clips with a smaller length  $L_f$  and use the same augmentations as it in the pre-train stage for fine-tuning. We use 2-second clips of videos for pre-training and only 1-second clips of videos for deepfake detection, the longer sequence in contrastive learning can diminish inner-modal similarity to reduce impacts of false-negative pairs and a shorter sequence can speed up the inference procedure. It is worth mentioning that the local lip movement feature extracted by the frontend is irrelevant to the sequence length.

## 4. Experiments

### 4.1. Implement Details

We use *VoxCeleb2* dataset [9] and part of *AVSpeech* dataset [15] for pre-training, totally contains 2,800,000 speech video clips. Each video clip includes more than 100 continuous frames with a trackable face. We use AdamW optimizer for pre-train, the initial learning rate is 0.01 and

Method	pre-train	Video-level ACC (%)			Video-level AUC (%)		
		Raw	HQ	LQ	Raw	HQ	LQ
Xception [36]	-	99.0	97.0	89.0	99.8	99.3	92.0
CNN-aug [43]	-	98.7	96.9	81.9	99.8	99.1	86.9
Patch-based [4]	-	99.3	92.6	79.1	99.9	97.2	78.3
Two-branch [28]	-	—	—	—	—	99.1	91.1
Face X-ray [22]	Sup-BI	99.1	78.4	34.2	99.8	97.8	77.3
CNN-GRU [37]	-	98.6	97.0	90.1	99.9	99.3	92.2
LipForensics [18]	Sup-LRW	98.9	<b>98.8</b>	<b>94.2</b>	99.9	<b>99.7</b>	<b>98.1</b>
Lipforensics* [18]	Sup-LRW	98.6	96.7	88.6	99.8	99.3	94.9
ours	Self-Sup	<b>99.2</b>	98.6	93.5	<b>99.9</b>	99.6	96.7

Table 1. Performance when fine-tune the model on each compression rate of FF++ dataset and test for the same compression rate. Our model achieves competitive performances on raw videos and HQ videos.

decay 0.9 for each epoch, the minimum learning rate is 0.0001. The hyper-parameter of each transformer layer is 1024 dims, 8 attention heads, 128 dims for each attention head, 2048 intermediate dims and 0.2 dropout rate. The dimension of features for contrastive learning is 256. We end the pre-training when the training loss stops descending for 3 epochs. We conduct extensive experiments with different hyper-parameter combinations to determine the optimal selection for our method. Parts of the results are listed in Table 4. Finally, we choose ResNet18 as frontend and 6 layers transformer as backend for the balance of efficiency and accuracy.

We fine-tune our pre-trained video encoder on original FaceForensics++ dataset [36] (FF++, 4 generation methods: Deepfake, FaceSwap, Neurltextures, Face2Face and Original videos, 720 videos for each class). And we use the test set of Celeb-DF v2 [24] (518 videos), Deeper-forensics [19] (140 videos, paired with original videos of FF++), 140 Faceshifter videos from FF++ (paired with original videos) and DFDC [14] (3000 videos selected from the test set, exclude extremely corrupted videos) for evaluating the cross dataset ability. Following the commonly used evaluation metrics in previous deepfake detection works, we leverage video level ROC-AUC score and binary classification accuracy for the evaluation of the detection performances.

Notably, due to the training set being unbalanced (the fake videos are 4 times as many as real videos), previous methods usually over-sample real videos. In this work, we sample real videos with similar quality from the pre-train dataset to make the training set balanced.

## 4.2. Evaluation on FaceForensics++

FaceForensics++ is the most widely used dataset for deepfake detection. The original videos of the FF++ dataset include speeches by people of different races. The dataset chooses similar identity pairs for face-swapping, a total of five face swapping methods are used to generate face forg-

eries on each pair. To consist with previous works, we only use 4 face forgery methods: Deepfake, FaceSwap, Neurltextures, Face2Face to fine-tune our models, and we treat the FaceShifter as a discrete dataset. There are three compression rates of the FF++ dataset, uncompressed (Raw), slightly compressed (HQ), and heavily compressed (LQ). The uncompressed dataset is relatively easy for spotting deepfake artifacts, but the challenge is increasing with compression rate, in the LQ videos, most textural features are lost which is extremely challenging for detection.

In Table 1, we report the video-level AUC and accuracy for each compression rate compared to other methods. Among which, Xception [36], CNN-aug [43], Patch-based [4], Two-branch [28] and CNN-GRU [37] are networks without special pre-train strategies. Face X-ray [22] generates training data by blending real faces and being supervised with blending boundary regression. Lipforensics [18] is pre-trained in a supervised manner on Lipreading in the Wild (LRW) [10] dataset. To parallel compare our self-supervised method to the supervised one, we also use the same network structure and pre-train it on LRW with the same training strategy of Lipforensics in a supervised manner, the model is denoted as Lipforensics\*. Apparently, our self-supervised pre-trained model outperforms the supervised one.

From the results in Table 1, we can observe that our model pre-trained in self-supervised manner performs much better than it in supervised manner. Compared to the state-of-the-art Lipforensics method, our model also achieves comparable performances on raw videos and HQ videos while slightly dropping on LQ videos.

## 4.3. Evaluation of Cross-manipulation Ability

Generalizing to unseen forgery classes is challenging in deepfake detection because the artifact patterns are different in each manipulation methods. Some deepfake detection methods trained with several forgery classes may overfit

Method	Train on remaining three				
	DF	FS	F2F	NT	Avg
Xception [36]	93.9	51.2	86.8	79.7	77.9
CNN-aug [43]	87.5	56.3	80.1	67.8	72.9
Patch-based [4]	94.0	60.5	87.3	84.8	81.7
Face X-ray [22]	99.5	93.2	94.5	92.5	94.9
CNN-GRU [37]	97.6	47.6	85.8	86.6	79.4
LipForensics [18]	<b>99.7</b>	90.1	<b>99.7</b>	<b>99.1</b>	<b>97.1</b>
LipForensics*	97.8	90.5	98.0	96.9	95.8
ours	98.5	<b>91.9</b>	98.3	96.4	96.3

Table 2. **Cross manipulation method generalisation.** Video-level AUC (%) when testing on each forgery type of FaceForensics++ HQ after training on the remaining three. The types are Deepfakes (DF), FaceSwap (FS), Face2Face (F2F), and NeuralTextures (NT).

multiple class-specific artifacts but not capture the common feature transferable to new method.

In this experiment, we evaluate the models’ generalization ability cross manipulation methods with the leave-one-out strategy. Each time, we train the model with 3 forgery classes in FF++ HQ dataset and test with the remaining forgery class. The results are shown in Table 2. Our method achieves better performance than most baseline methods, this indicates the lip movement features captured by our model are well generalizable.

#### 4.4. Evaluation of Cross-datasets Generalization

In real deepfake detection scenarios, the distribution of deepfakes is more complex. The forged videos are not only diverse in source videos and generation methods, there is also a diversity in post-processing methods. Thus the domain generalization ability is a significant metric for deepfake detection models. Since there are large domain gaps between different deepfake datasets, in this part, we evaluate the domain generalization ability of deepfake detectors by cross-dataset test. We fine-tune our model on the FF++ training set with four manipulation methods (DF, F2F, FS, NT), and report the AUC scores tested on Celeb-DF, DFDC, FaceShifter and DeeperForensics dataset, respectively.

In Table 3, our method achieves state-of-the-art performances in generalization to Celeb-DF-v2, DFDC and FF++ FaceShifter dataset. Our self-supervised pre-trained model surpasses Lipforensics and the same model using supervised training strategy (Lipforensics\*) for an average AUC of 0.7 percent. The performances on FaceShifter and Deeperforensics are obviously better than those on Celeb-DF-v2 and DFDC. It is probably because these two datasets share the same original videos with the FF++.

Method	CDF	DFDC	FSh	DFo	Avg
Xception [36]	73.7	70.9	72.0	84.5	75.3
CNN-aug [43]	75.6	72.1	65.7	74.4	72.0
Patch-based [4]	69.6	65.6	57.8	81.8	68.7
Face X-ray [22]	79.5	65.5	92.8	86.8	81.2
CNN-GRU [37]	69.8	68.9	80.8	74.1	73.4
Multi-task [30]	75.7	68.1	66.0	77.7	71.9
DSP-FWA [23]	69.5	67.3	65.5	50.2	63.1
LipForensics [18]	82.4	73.5	97.1	<b>97.6</b>	87.7
LipForensics*	83.5	73.7	96.3	97.2	87.7
ours	<b>84.2</b>	<b>74.5</b>	<b>97.8</b>	97.3	<b>88.4</b>

Table 3. **Cross-dataset generalisation.** Video-level AUC (%) on Celeb-DF-v2 (CDF), DeepFake Detection Challenge (DFDC), FaceShifter HQ (FSh), and DeeperForensics (DFo) when trained on FaceForensics++.

## 4.5. Ablation Study

### 4.5.1 Determination of Model Architecture

In Lipforensics, the network architecture is referred from a lipreading model [26]. In this paper, we propose a self-supervised audio-visual contrastive pre-train instead of a lipreading pre-train. Different from lipreading tasks that classify limited words, audio-visual contrastive learning requires the model to extract more detailed semantic features. Thus a transformer would be a better choice than MSTCN in lipreading tasks. The results in Table 4 also demonstrate that transformer performs better than MSTCN with the same layer number. We also conduct a group of experiments to explore a relatively better parameter scale of the frontend 2D ResNet and the backend transformer.

We pre-train each model with self-supervised learning and fine-tune on FF++, we use the performances on FF++ HQ, FF++ LQ (same setting as Table 1), and the average cross-datasets performances (same setting as Table 3) as metrics for evaluation. As the dimension of ResNet50 output is 2048, we double the width of its backend.

Compared with 4-layers MSTCN backend, 4-layers transformer achieves about 1 percent AUC gain in FF++ LQ dataset and more than 2 percent AUC gain in cross dataset transfer. It verifies the advantage of the transformer. The ResNet50 based models outperform ResNet18, which indicates that a larger model will be helpful to capture more precise and detailed lip movements. When increasing the depth of the backend transformer from 6-layers to 8-layers, the performance of the ResNet18 model stops improving but the ResNet50 model still improves. This might indicate that the gain from deepening the backend is restricted by frontend scale and backend width. Eventually, we choose ResNet18 (same setting as Lipforensics for a fair comparison) as frontend and 6 layers transformer as backend for the balance of efficiency and accuracy.

Frontend	backend	Total parameters	FF++ HQ	FF++ LQ	Avg Cross
ResNet18	L4-MSTCN	36.0M	98.8	95.4	84.5
ResNet18	L4-Transformer	45.3M	99.3	96.4	87.8
ResNet18	L6-Transformer	62.1M	99.6	96.7	88.4
ResNet18	L8-Transformer	78.9M	99.6	96.5	88.1
ResNet50	L4-MSTCN	132.1M	99.0	95.7	85.9
ResNet50	L4-Transformer	162.1M	99.5	96.6	87.4
ResNet50	L6-Transformer	229.2M	99.7	97.0	88.7
ResNet50	L8-Transformer	296.3M	99.7	97.1	88.9

Table 4. Ablation study of network architecture combination of frontend and backend. The transformer backend outperforms MSTCN with the same layer number by a remarkable improvement. The model with larger frontend scale shows better performance in generalization.

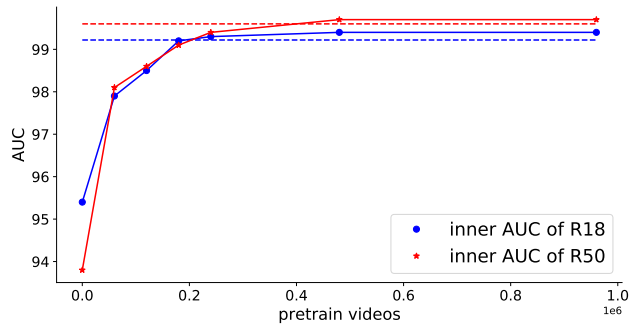


Figure 2. The relation between scale of pre-train data and test AUC on FF++ HQ dataset.

#### 4.5.2 Effect of Pre-training Dataset Scale

Compared to the supervised manner, self-supervised learning can leverage a larger amount of training data. To be specific, the labeled dataset LRW for lipreading task contains about 500,000 video clips of 500 words with precise boundary annotation. It only covers a limited part of real speech scenes but costs expensive human efforts in annotating. Contrastly, self-supervised learning only requires natural speech videos, which are very easy to collect. Empowered by self-supervised pre-train, we can improve the deepfake detection accuracy without requiring extra fake videos as training data. In this part, we give a simple investigation of how the quantity of unlabeled real video for pre-train affects the model’s accuracy and the cross-dataset ability. We draw the relation curves in Figure 2 and Figure 3. The blue lines are ResNet18 based models and the red lines are ResNet50 based models, we also plot the AUC of the model pre-trained on the LRW dataset with horizontal dotted lines. The abscissa of 0 represents that the model has not been pre-trained.

The curves in Figure 2 and Figure 3 validate that the models pre-trained with more original videos could improve both in-dataset and cross-datasets accuracies. When pre-trained with 480,000 video clips, the performance exceeds the same network pre-trained with 500,000 labeled video clips from LRW dataset. And the performance can still improve with

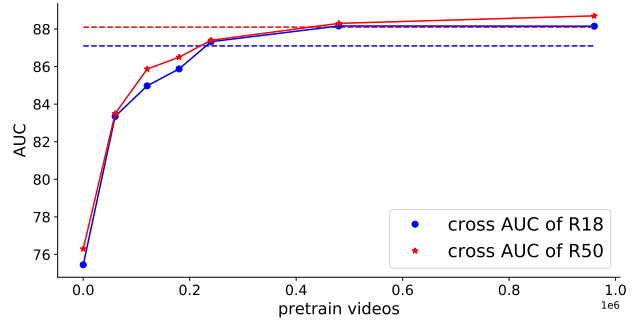


Figure 3. The relation between scale of pre-training data and average cross datasets AUC fine-tuned on FF++ dataset.

the increase in training videos. From the comparison between the ResNet18 based model and the ResNet50 based model, we observe that without pre-train, the larger model does not outperform the smaller model on FF++ HQ while achieving only minor ascendancy for generalization. But the models obtain obvious gains from pre-training. The AUC of the ResNet50 based model grows faster and saturates slower than ResNet18 due to the larger model scale. In conclusion, we demonstrate that larger-scale pre-training helps boost the models’ performances. This highlights another advantage of our self-supervised method, that is, significantly reducing the data annotation costs to provide a larger scale of training data. However, larger-scale training data inevitably leads to higher computing resources demands, which might be a limitation.

#### 4.5.3 Robustness of Pre-train Features

As aforementioned, the video in real-world scenarios may be accompanied by various disturbances, which requires that the detection model should be robust to various degradations of video quality. In our framework, we use massive real videos with various definitions to pre-train the video encoder. In this way, the learned representations of lip movement would be more robust to corruption. When fine-tuned for deepfake detection, we freeze the frontend of the video encoder so that the perturbations in videos would be reduced in the intermediate features. To evident the robustness of

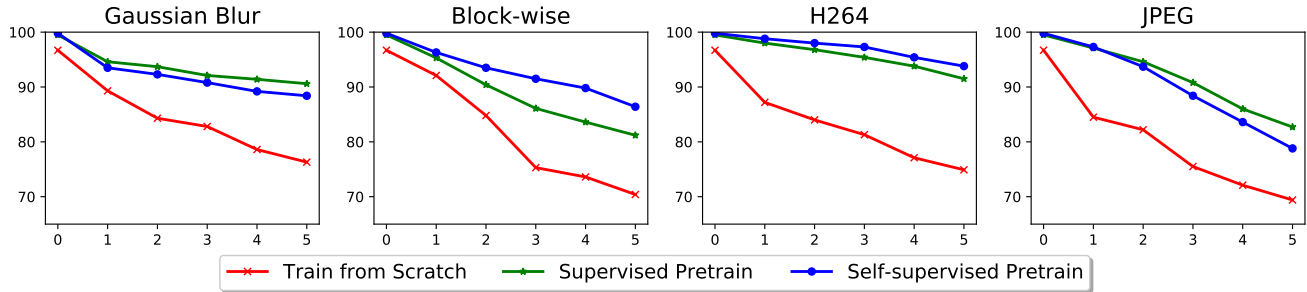


Figure 4. Video-level AUC scores of the spatial-temporal network as a function of the severity level for four types of perturbations: Gaussian Blur, Block-wise distortion, H264 Compression and Pixelation Distortion. We compare the robustness of models with each pre-train strategy.

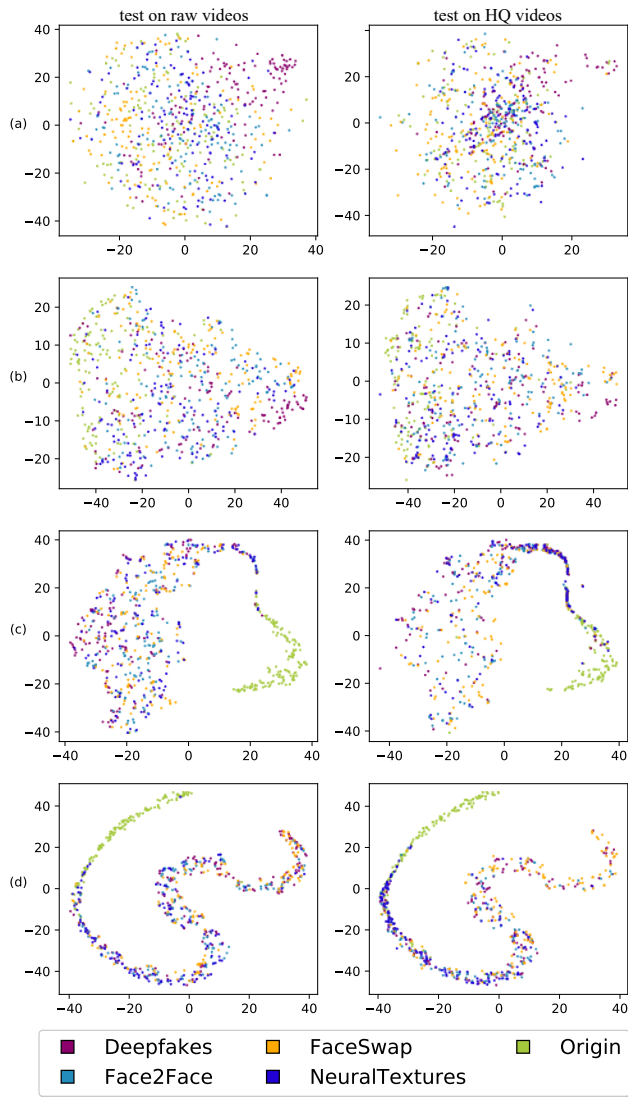


Figure 5. Visualization of features with t-SNE. Features from self-supervise pre-trained model are less affected by compression.

our model, we visualize the features of clean videos and corrupted videos extracted from the frontend and backend by our model.

We visualize the features from the self-supervised pre-trained model and the train-from-scratch model in Figure 5. The models are trained with FF++ Raw videos. The left column shows t-SNE embedded features from the Raw test set and the right column shows features from the HQ test set (compressed with H264 codec) transformed with Raw t-SNE embeddings. Row (a) and row (c) illustrate frontend and backend features of the non-pre-trained model separately, row (b) and row (d) are corresponding features of the self-supervised pre-trained model. We eventually observe that the distribution of the features trained with audio-visual contrastive learning changes more slightly compared to the train-from-scratch one. As a consequence, we believe that the deepfake detection models with self-supervised pretrain are less susceptible to interference from video compression.

Further, we investigate the robustness of our methods to various video corruptions. We take four common corruption types to simulate various degradations of definition: adding block-wise distortions, Gaussian blurring, JPEG compression, and H264 video compression, each corruption with five severity levels (details refer to [19]). In Figure 4, we give the AUC curves under different corruption types. In subjects of Block-wise distortion and video compression, our self-supervised method performs better robustness than supervised pre-train, and in other two subjects, our method also achieves comparable performances.

## 5. Conclusion

In this paper, we propose a self-supervised pre-train method using audio-visual contrastive learning for improving the robustness and transferability of deepfake detection. Our method extracts a generic representation of local lip movement by 3D convolution and 2D ResNet and captures the long-term incoherence of lip movement with a transformer. Extensive experiments show that the generic pre-trained model is effective for deepfake detection. More importantly, with the self-supervised pre-trained framework,



the accuracy and generalization ability can be possibly improved by scaling up model size and data. We believe the potential of self-supervised learning can be further explored for deepfake detection in the future.

## References

- [1] Shruti Agarwal, Hany Farid, Ohad Fried, and Maneesh Agrawala. Detecting deep-fake videos from phoneme-viseme mismatches. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 660–661, 2020. [1](#), [2](#)
- [2] Alexei Baevski, Henry Zhou, Abdel rahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations. *ArXiv*, abs/2006.11477, 2020. [3](#)
- [3] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *International Conference on Computer Vision*, 2017. [4](#)
- [4] Lucy Chai, David Bau, Ser-Nam Lim, and Phillip Isola. What makes fake images detectable? understanding properties that generalize. *arXiv preprint arXiv:2008.10588*, 2020. [2](#), [5](#), [6](#)
- [5] Renwang Chen, Xuanhong Chen, Bingbing Ni, and Yanhao Ge. Simswap: An efficient framework for high fidelity face swapping. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2003–2011, 2020. [1](#)
- [6] Shen Chen, Taiping Yao, Yang Chen, Shouhong Ding, Jilin Li, and Rongrong Ji. Local relation learning for face forgery detection. *arXiv preprint arXiv:2105.02577*, 2021. [2](#)
- [7] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey E. Hinton. A simple framework for contrastive learning of visual representations. *ArXiv*, abs/2002.05709, 2020. [2](#)
- [8] Xinlei Chen, Saining Xie, and Kaiming He. An empirical study of training self-supervised vision transformers. *ArXiv*, abs/2104.02057, 2021. [2](#)
- [9] Joon Son Chung, Arsha Nagrani, and Andrew Zisserman. Voxceleb2: Deep speaker recognition. *arXiv preprint arXiv:1806.05622*, 2018. [4](#)
- [10] J. S. Chung and A. Zisserman. Lip reading in the wild. In *Asian Conference on Computer Vision*, 2016. [2](#), [5](#)
- [11] Joon Son Chung and Andrew Zisserman. Out of time: Automated lip sync in the wild. In *ACCV Workshops*, 2016. [4](#)
- [12] Ishan Dave, Rohit Gupta, Mamshad Nayeem Rizve, and Mubarak Shah. Tclr: Temporal contrastive learning for video representation. *arXiv preprint arXiv:2101.07974*, 2021. [3](#)
- [13] Jiankang Deng, J. Guo, Y. Zhou, Jinke Yu, I. Kotsia, and S. Zafeiriou. Retinaface: Single-stage dense face localisation in the wild. *ArXiv*, abs/1905.00641, 2019. [4](#)
- [14] Brian Dolhansky, Joanna Bitton, Ben Pfau, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv: Computer Vision and Pattern Recognition*, 2020. [5](#)
- [15] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin W. Wilson, Avinatan Hassidim, William T. Freeman, and Michael Rubinstein. Looking to listen at the cocktail party. *ACM Transactions on Graphics (TOG)*, 37:1 – 11, 2018. [4](#)
- [16] Yazan Abu Farha and Juergen Gall. Ms-tcn: Multi-stage temporal convolutional network for action segmentation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3570–3579, 2019. [2](#)
- [17] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross B. Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3298–3308, 2021. [2](#)
- [18] Alexandros Haliassos, Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic. Lips don't lie: A generalisable and robust approach to face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5039–5049, 2021. [1](#), [2](#), [5](#), [6](#)
- [19] Liming Jiang, Wayne Wu, Ren Li, Chen Qian, and Chen Change Loy. Deepforensics-1.0: A large-scale dataset for real-world face forgery detection. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2886–2895, 2020. [5](#), [8](#)
- [20] Hyeonwoo Kim, Pablo Garrido, Ayush Tewari, Weipeng Xu, Justus Thies, Matthias Niessner, Patrick Pérez, Christian Richardt, Michael Zollhöfer, and Christian Theobalt. Deep video portraits. *ACM Transactions on Graphics (TOG)*, 37(4):1–14, 2018. [1](#)
- [21] Mohammad Rami Koujan, Michail Christos Doukas, Anastasios Roussos, and Stefanos Zafeiriou. Head2head: Video-based neural head synthesis. In *2020 15th IEEE International Conference on Automatic Face and Gesture Recognition (FG 2020)*, pages 16–23. IEEE, 2020. [1](#)
- [22] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5001–5010, 2020. [1](#), [2](#), [5](#), [6](#)
- [23] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. In *IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019. [6](#)
- [24] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3204–3213, 2020. [5](#)
- [25] Honggu Liu, Xiaodan Li, Wenbo Zhou, Yuefeng Chen, Yuan He, Hui Xue, Weiming Zhang, and Nenghai Yu. Spatial-phase shallow learning: rethinking face forgery detection in frequency domain. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 772–781, 2021. [1](#), [2](#)
- [26] Pingchuan Ma, Brais Martinez, Stavros Petridis, and Maja Pantic. Towards practical lipreading with distilled and efficient models. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 7608–7612, 2021. [6](#)
- [27] Shuang Ma, Zhaoyang Zeng, Daniel McDuff, and Yale Song. Contrastive learning of global and local audio-visual representations. *arXiv preprint arXiv:2104.05418*, 2021. [2](#), [3](#)
- [28] Iacopo Masi, Aditya Killekar, Royston Marian Mascarenhas, Shenoy Pratik Gurudatt, and Wael AbdAlmageed. Two-

- branch recurrent network for isolating deepfakes in videos. In *European Conference on Computer Vision*, pages 667–684. Springer, 2020. 1, 2, 5
- [29] Pedro Morgado, Nuno Vasconcelos, and Ishan Misra. Audio-visual instance discrimination with cross-modal agreement. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12470–12481, 2021. 2, 3
- [30] H. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen. Multi-task learning for detecting and segmenting manipulated facial images and videos. In *2019 IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*, pages 1–8, 2019. 6
- [31] Yuval Nirkin, Yosi Keller, and Tal Hassner. Fsgan: Subject agnostic face swapping and reenactment. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 7184–7193, 2019. 1
- [32] Tian Pan, Yibing Song, Tianyu Yang, Wenhao Jiang, and Wei Liu. Videomoco: Contrastive video representation learning with temporally adversarial examples. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11200–11209, 2021. 3
- [33] KR Prajwal, Rudrabha Mukhopadhyay, Vinay P Namboodiri, and CV Jawahar. Learning individual speaking styles for accurate lip to speech synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13796–13805, 2020. 1
- [34] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. Thinking in frequency: Face forgery detection by mining frequency-aware clues. In *European Conference on Computer Vision*, pages 86–103. Springer, 2020. 1, 2
- [35] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*, 2021. 2
- [36] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1–11, 2019. 1, 5, 6
- [37] Ekraam Sabir, Jiabin Cheng, Ayush Jaiswal, Wael AbdAlmageed, Iacopo Masi, and Prem Natarajan. Recurrent convolutional strategies for face manipulation detection in videos. *Interfaces (GUI)*, 3(1), 2019. 5, 6
- [38] Nawid Sayed, Biagio Brattoli, and Björn Ommer. Cross and learn: Cross-modal self-supervision. *ArXiv*, abs/1811.03879, 2018. 2
- [39] Supasorn Suwajanakorn, Steven M Seitz, and Ira Kemelmacher-Shlizerman. Synthesizing obama: learning lip sync from audio. *ACM Transactions on Graphics (TOG)*, 36(4):1–13, 2017. 1
- [40] Justus Thies, Michael Zollhöfer, and Matthias Nießner. Deferred neural rendering: Image synthesis using neural textures. *ACM Transactions on Graphics (TOG)*, 38(4):1–12, 2019. 1
- [41] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *ArXiv*, abs/1807.03748, 2018. 3
- [42] Junke Wang, Zuxuan Wu, Jingjing Chen, and Yu-Gang Jiang. M2tr: Multi-modal multi-scale transformers for deepfake detection. *arXiv preprint arXiv:2104.09770*, 2021. 2
- [43] Sheng-Yu Wang, Oliver Wang, Richard Zhang, Andrew Owens, and Alexei A Efros. Cnn-generated images are surprisingly easy to spot... for now. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, volume 7, 2020. 5, 6
- [44] Xi Wu, Zhen Xie, YuTao Gao, and Yu Xiao. Sstnet: Detecting manipulated faces through spatial, steganalysis and temporal features. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2952–2956. IEEE, 2020. 1
- [45] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2185–2194, 2021. 1, 2
- [46] Tianchen Zhao, Xiang Xu, Mingze Xu, Hui Ding, Yuanjun Xiong, and Wei Xia. Learning to recognize patch-wise consistency for deepfake detection. *arXiv preprint arXiv:2012.09311*, 2020. 1, 2
- [47] Yuhao Zhu, Qi Li, Jian Wang, Cheng-Zhong Xu, and Zhenan Sun. One shot face swapping on megapixels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4834–4844, 2021. 1
- [48] Mohammadreza Zolfaghari, Yi Zhu, Peter V. Gehler, and Thomas Brox. Crossclr: Cross-modal contrastive learning for multi-modal video representations. *ArXiv*, abs/2109.14910, 2021. 2