# AV-Lip-Sync+: Leveraging AV-HuBERT to Exploit Multimodal Inconsistency for Deepfake Detection of Frontal Face Videos

Sahibzada Adil Shahzad, Ammarah Hashmi, Yan-Tsung Peng, Yu Tsao, *Senior Member, IEEE*, Hsin-Min Wang, *Senior Member, IEEE*

*Abstract*—Multimodal manipulations (also known as audio-visual deepfakes) make it difficult for unimodal deepfake detectors to detect forgeries in multimedia content. To avoid the spread of false propaganda and fake news, timely detection is crucial. The damage to either modality (i.e., visual or audio) can only be discovered through multimodal models that can exploit both pieces of information simultaneously. However, previous methods mainly adopt unimodal video forensics and use supervised pre-training for forgery detection. This study proposes a new method based on a multimodal self-supervised-learning (SSL) feature extractor to exploit inconsistency between audio and visual modalities for multimodal video forgery detection. We use the transformer-based SSL pre-trained Audio-Visual HuBERT (AV-HuBERT) model as a visual and acoustic feature extractor and a multi-scale temporal convolutional neural network to capture the temporal correlation between the audio and visual modalities. Since AV-HuBERT only extracts visual features from the lip region, we also adopt another transformer-based video model to exploit facial features and capture spatial and temporal artifacts caused during the deepfake generation process. Experimental results show that our model outperforms all existing models and achieves new state-of-the-art performance on the FakeAVCeleb and DeepfakeTIMIT datasets.

*Index Terms*—Deepfakes, Deepfake detection, Audio-Visual, Lip Syn, Inconsistency, Video Forgery, Audio-Visual Deepfake Detection, Multimedia Forensics, Multimodality

## I. INTRODUCTION

With smartphones, social networks, and the high-speed internet, it is now available at one's fingertips to capture, upload, and share content without any delays or fees. However, this convenience makes the spread of deepfakes a heavy social cost raising major social and ethical issues [1], [2]. The term "deepfake" encompasses synthetic media such as
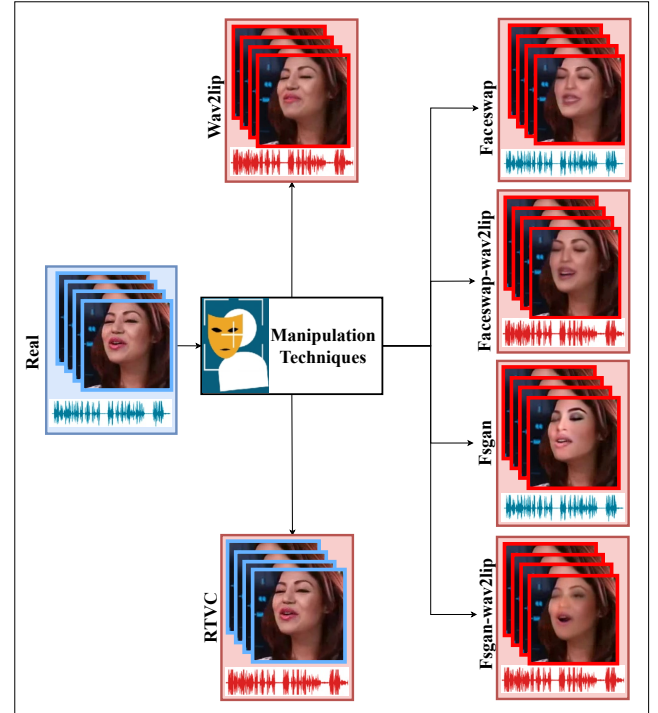


Fig. 1. Illustration of various deepfake manipulation techniques applied to a real audio-visual sample. The real sample (left) highlighted by the dark blue border contains the original video frames and corresponding audio waveform. The manipulated (fake) samples highlighted by dark red borders are generated using Wav2Lip, Faceswap, Faceswap-wav2Lip, Fsgan, Fsgan-wav2Lip, and RTVC (Real-Time Voice Cloning). The video frames highlighted by blue borders represent real frames, while the video frames highlighted by red borders represent manipulated (fake) frames. The blue waveforms represent real audio, while the red waveforms represent manipulated (fake) audio.

Sahibzada Adil Shahzad is with the Social Networks and Human-Centered Computing Program, Taiwan International Graduate Program, Academia Sinica, Taipei 11529, Taiwan, and also with the Department of Computer Science, National Chengchi University, Taipei 11605, Taiwan. (e-mail: adil-shah275@iis.sinica.edu.tw).

Ammarah Hashmi is with the Social Networks and Human-Centered Computing Program, Taiwan International Graduate Program, Academia Sinica, Taipei 11529, Taiwan, and also with the Institute of Information Systems and Applications, National Tsing Hua University, Hsinchu 30013, Taiwan. (e-mail: hashmiammarah0@gmail.com).

Yan-Tsung Peng is with the Department of Computer Science, National Chengchi University, Taipei, Taiwan. (e-mail: ytpeng@cs.nccu.edu.tw).

Yu Tsao is with the Research Center for Information Technology Innovation, Academia Sinica, Taipei 11529, Taiwan. (e-mail: yu.tsao@citi.sinica.edu.tw).

Hsin-Min Wang is with the Institute of Information Science, Academia Sinica, Taipei 11529, Taiwan. (e-mail: whm@iis.sinica.edu.tw)

images, videos, audio, and text. While there are many benefits to generating content through artificial intelligence (AI) technology, it also has a dark side. AI-generated content is often used for unethical objectives and malicious purposes, including disinformation, pornography, fraudulent activities, and political defamation [1]–[4]. Social networks [5] are conduits for this type of manipulated content and often do not have filtering mechanisms to prevent its rapid spread.

Deepfakes exist in various modalities, each posing distinct challenges. Text deepfakes refer to seemingly real text information manipulated or generated by AI. Examples include fabricated news articles [6], deceptive online reviews [7], and syn-

thetic text that incites violence [8]. Audio deepfakes involve the acoustic manipulation of audio using voice conversion and text-to-speech techniques (such as WaveNet [9], Wave-Glow [10], MelNet [11], and Tacotron [12]). SV2TTS [13] is a real-time voice cloning tool that synthesizes fabricated audio content without altering visual content, such as video frames. Among the various forms of deepfake media, video deepfakes are the most prevalent. These AI-generated videos realistically superimpose one person's face onto another, modifying expressions or mimicking speech patterns. Faceswap [14] and FSGAN [15] are common techniques for visual deepfake manipulation.

Advances in deepfake generation technology have made the task of facial forgery detection significantly challenging. Audio-visual deepfake techniques have reached new levels of realism with this advancement. The realistic nature of audio-visual deepfake videos stems from the simultaneous manipulation of both audio and visual modalities and the consistency between synthesized speech and facial movements, making traditional forgery detection methods less effective. Humans often fail to distinguish real from manipulated videos, yet remain overconfident in their judgments [16], [17]. Similarly, recent multimodal large language models (LLMs) such as ChatGPT [18] demonstrate limited effectiveness in multimedia forensics and struggle with highly realistic manipulations. These challenges require more advanced detection strategies that go beyond simple frame-based analysis and the development of specialized tools and algorithms designed specifically for multimedia forensics tasks. To solve this problem, the forensics research communities in the fields of image, video, and audio have specially designed algorithms to detect manipulation in videos. The image and video forensics research community focuses primarily on detecting visual forgeries [19]–[23], while the audio forensics research community has developed detection systems to detect acoustic manipulation [24]–[26]. Due to their unimodal nature, these detectors fail when the manipulated modality is not seen during training. Spoofed audio can evade visual deepfake detectors, while audio deepfake detectors cannot catch visual deepfakes.

Recently, inspired by human's ability to process audio and visual signals simultaneously to perceive the world [27], the field of multimedia forensics has turned to the development of multimodal systems for effective deepfake detection. Humans instinctively combine auditory and visual signals to assess the authenticity of content, often relying on the synchronization of speech and facial movements. Likewise, utilizing multimodal data such as audio and visual modalities has emerged as a promising approach to improve the accuracy and robustness of video forgery detection [28]. Therefore, we choose to integrate audio and visual information into our bimodal approach to identify audio-visual deepfakes. To further facilitate this audio-visual deepfake detection technology, various pre-trained models were fine-tuned on deepfake video datasets to detect manipulation.

Pre-trained self-supervised-learning (SSL) models have recently emerged and achieved success in various downstream tasks. Audio-Visual HuBERT (AV-HuBERT) [29] is an SSL-based audio-visual representation learning model that achieves state-of-the-art performance in lip reading, audio-visual speech enhancement [30], audio-visual speech separation [30], and audio-visual speech recognition [31]. Motivated by its state-of-the-art performance in multiple downstream tasks, we leverage AV-HuBERT for feature extraction to capture the inconsistency between the mouth region of interest and the corresponding audio modality for the downstream task of audio-visual deepfake detection. Considering that lip feature-based deepfake detectors may fail when the lip region is not manipulated or only slightly manipulated, to enhance our audio-visual forgery detection model, we also employ Video Vision Transformer (ViViT) [32] as the face encoder to exploit whole-face features to assist our proposed audio-visual deepfake detector. By integrating powerful audio-visual representations, speech-lip synchronization features, and spatiotemporal facial features, the proposed system achieves state-of-the-art performance on the FakeAVCeleb [33] and DeepfakeTIMIT [34] datasets.

This work focuses on deepfake detection of frontal face videos with speech in the audio track. Such deepfake videos of celebrities are widely circulated on social media platforms to spread misinformation (or disinformation) or tarnish a person's reputation, causing serious harm to society. Our main contributions are as follows.

- We propose AV-Lip-Sync+, a novel audio-visual deepfake detection model that combines self-supervised speech-lip synchronization features, audio-visual embeddings, and facial representations to capture multimodal inconsistencies and spatiotemporal artifacts.
- We leverage a transformer-based architecture to improve temporal modeling and cross-modal alignment, surpassing conventional CNN-based approaches in capturing subtle deepfake cues.
- We introduce a dedicated face encoder to enhance the detection performance on full-face manipulations (e.g., Faceswap and FSGAN), and evaluate the model's robustness under generalization and partial face occlusion scenarios.
- Extensive experiments on multiple benchmark datasets, including FakeAVCeleb, DeepfakeTIMIT, and DFDC, demonstrate that our approach consistently outperforms existing state-of-the-art methods in terms of accuracy and robustness.

## II. RELATED WORK

In this section, we first briefly review common deepfake video generation techniques, and then introduce state-of-the-art methods for deepfake detection.

### A. Deepfake Generation

Deepfake manipulation comes in many forms, with visual forgery [14], [35]–[39] being the most common, which involves using deepfake generation models such as Faceswap [14], FSGAN [15], and wav2lip [40] to manipulate the entire face, facial expressions, or lip movements. Traditional deepfake techniques are limited to visual forgery, leaving the audio unaltered. The evolution of synthetic media

technology has significantly improved the realism of deceptive content, among which audio-visual deepfakes enable the integration of visual and auditory alterations. Techniques like Faceswap-wav2lip and Fsgan-wav2lip simultaneously manipulate face and lip movements and align the latter with the audio track. Fig. 1 shows real and various audio-visual fake video samples from the FakeAVCeleb dataset [33]. The focus of our study is on detecting audio-visual deepfakes.

### B. Deepfake Detection

High-quality AI-generated content is useful in many ways, but it also comes at a cost, and timely detection is crucial to avoid any harm to society [3]. To avoid the spread of misinformation and disinformation and to protect the reputations and privacy of individuals, we need automated methods to promptly detect deepfake content in our widespread digital world. Academia and industry have made considerable progress in using deep learning-based methods to detect forged multimedia content. These deep learning-based deepfake detection methods can be roughly divided into two major categories: unimodal and multimodal methods.

Unimodal forgery detection models are specifically designed to detect forgery in one modality (video or audio) and rely only on the corresponding modality to identify forgery. Video forgery detectors can be divided into three types [41]: physiological, visual artifact-based, and high-level feature-based methods. In physiological methods, researchers have exploited abnormal eye blinks [42] and incoherent head poses [43]. Visual artifact-based methods analyze anomalies and irregularities such as unnatural facial movements, illumination variation, blended face boundaries, and misalignment of video content. FInfer [44] is a frame inference-based detection framework to solve the problem of high-visual-quality deepfake detection. ICM [45] is a deepfake detection model that captures dynamic inconsistencies between visual frames in deepfake videos. GFA-CNN [46] is proposed to learn identity-aware and generalizable features for face anti-spoofing tasks. High-level-feature-based deepfake detectors extract high-level features that are immune to video processing (e.g., compression). Lipforensics [19] is an example of a high-level feature-based deepfake detector that leverages a lip-reading-based model to detect abnormal lip movements for video deepfake detection.

Video forgery detection models have achieved excellent results, thanks to various rich datasets available for model training, such as DeepfakeTIMIT [34], UADVF [43], FaceForensics++ [47], Celeb-DF [48], DFDC [49], DeeperForensics [50], and the recently released multimodal FakeAVCeleb dataset [33]. Well-known unimodal fake video detection models include Capsule Forensics [20], HeadPose [21], Xception [22], LipForensics [19], Meso-4 and MesoInception-4 [23].

On the other hand, to trick automatic speaker verification systems, attackers can develop audio spoofing attacks or replay attacks using only a few minutes of a person's recorded speech, which makes these systems vulnerable. Similar to video deepfakes, audio spoofing attacks are a major challenge that must be solved. In response to audio spoofing attacks

on automatic speaker verification systems, the audio forensics community has proposed various traditional and deep learning-based methods [51]. Most traditional systems use short-term power spectrum, short-term phase spectrum, and long-term spectral features as front-end features. Backend classifiers are based on traditional machine learning models or ensemble models. Deep learning-based models can be roughly divided into multi-pass, end-to-end, and ensemble models. Using different hand-crafted acoustic features or raw waveforms, audio spoofing detection methods [24]–[26] can differentiate between genuine speech and spoofed speech. A step forward from whole-utterance spoof detection, H-MIL [52] employs a multiple instance learning approach for partially spoofed speech detection.

Recently, deepfake technology has expanded from unimodal manipulation to multimodal manipulation (such as audio-visual manipulation). Audio-visual deepfakes involve simultaneously manipulating the facial features and speech patterns of a subject in a video, making them more intricate and less detectable. Unimodal deepfake video and spoofed audio detectors are insufficient to detect these audio-visual manipulations. Despite their effectiveness in certain cases, these unimodal techniques often lack the ability to identify crossmodal inconsistencies. To solve this problem, the multimedia forensics community has proposed several deep learning-based methods to capture unimodal or multimodal manipulations. To effectively identify audio-visual deepfakes and reduce their misuse, advanced detection techniques that combine multimodal fusion, temporal consistency analysis, and deep learning-based anomaly detection are needed, all of which face considerable challenges.

Few studies have addressed this problem through multimodal methods that exploit faces as visual features and mel-frequency cepstral coefficients (MFCCs) as acoustic features [53], or detect forged videos based on facial and speech emotions [54]. To capture the intrinsic synchronization between visual and acoustic modalities, Facebook AI [55] built a sync-stream by connecting video and audio network feature representations through intra-attention (self-attention) and inter-attention mechanisms within and across the video and audio modalities. In [56], speech-lip synchronization features based on the difference between the extracted lip sequence and the synthetic lip sequence are used to detect audio-visual deepfakes. This model requires the wav2lip [57] module to generate lip sequences from audio, which increases model complexity, training, and inference time. AVFakeNet audio-visual forgery detection method [58] uses Swin Transformer as the feature extraction module. In [59], a multimodal deepfake detector based on ensemble learning is introduced to leverage multiple learners and make decisions based on hard voting. In its subsequent extended model in [60], three separate models, namely audio network, video network, and audio-visual network, are all built on pre-trained transformer-based foundation models. The main disadvantage of ensemble learning is its time-consuming process due to the involvement of multiple training models.

The field of multimodal forgery detection is still less explored, and further research is needed to determine how to
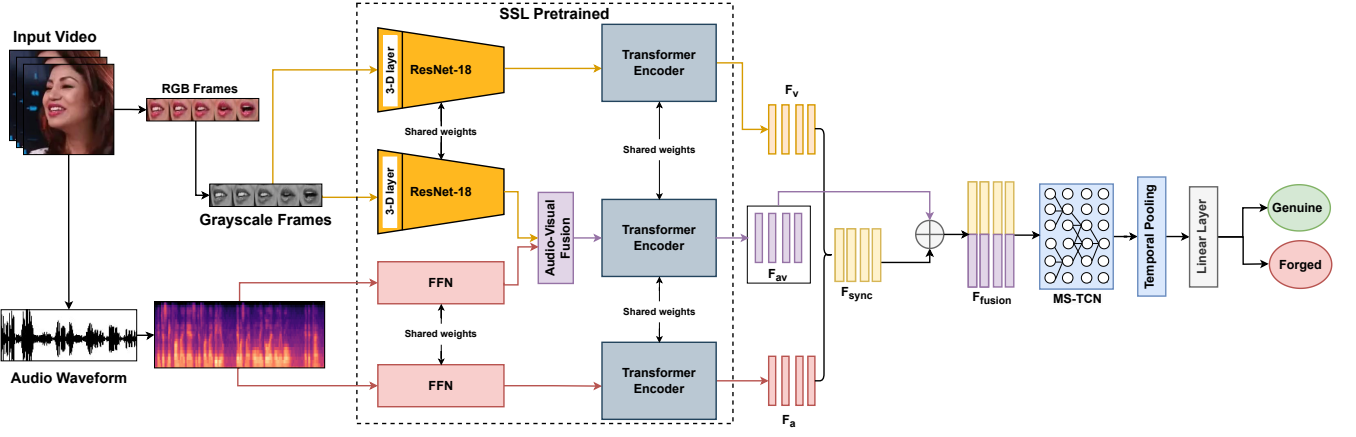
Fig. 2. The proposed AV-Lip-Sync+ architecture for multimodal forgery detection. The lip image sequence is extracted from the input video, while the log filterbank energies are extracted from the audio track. The SSL pre-trained model consists of ResNet-18 for visual feature extraction, FFN for acoustic feature extraction, and a transformer encoder to extract spatiotemporal information from the visual and acoustic features. The extracted audio-visual features are further mapped through multi-scale temporal convolution network (MS-TCN), temporal pooling, and linear layer for classification.

effectively use audio-visual information to detect forgery in any modality of multimedia content.

## III. METHODOLOGY

As shown in Fig. 2, our proposed AV-Lip-Sync+ model consist of three feature extractors, namely a lip image feature extractor, an acoustic feature extractor, and an audio-visual feature extractor. These feature extraction modules are followed by the Sync-Check Module, Feature Fusion Module, and Temporal Convolutional Network to capture the temporal correlation between visual and audio features. Finally, a temporal pooling layer and a linear layer are used for classification. Our feature extractors are based on the AV-HuBERT model [29] pre-trained on the LRS3 dataset [61]. AV-HuBERT will be fine-tuned when training the detection model on the multimodal deepfake datasets.

### A. Audio-Visual Feature Extractor

As shown in Fig. 2, the audio-visual feature extractor consists of a Resnet-18, a light-weight Feed Forward Network (FFN), and a transformer encoder. The 2-D Resnet-18 with front-end 3-D convolutional layers is used to extract lip-based visual features from each input lip image frame. The FFN is used to extract frame-level acoustic features from the input log filterbank energies of the audio waveform. These frame-level visual and acoustic features are concatenated along the feature dimension and fed to the transformer encoder, which generates a 768-D audio-visual embedding sequence $\vec{F}_{av}$ via

$$\vec{F}_{av} = F_{\theta_e}(v, a), \tag{1}$$

where $v$ and $a$ represent frame-level visual features and acoustic features, respectively.

### B. Lip Image Feature Extractor

In the lip image feature extractor, the output of Resnet-18 is fed to the transformer encoder, which generates the lip image embedding sequence $\vec{F}_v$ via

$$\vec{F}_v = F_{\theta_e}(v, a_{dropout}), \tag{2}$$

where $a_{dropout}$ indicates audio dropout, i.e., the audio information is not used.

### C. Acoustic Feature Extractor

In the acoustic feature extractor, the output of FFN is fed to the transformer encoder, which generates the acoustic feature embedding sequence $\vec{F}_a$ via

$$\vec{F}_a = F_{\theta_e}(v_{dropout}, a), \tag{3}$$

where $v_{dropout}$ indicates video dropout, i.e., the visual information is not used.

### D. Sync-Check Module

Because deepfake technology is not yet mature enough to generate synchronized and perfect audio-visual deepfakes, there is often a disharmony between the visual and audio modalities of deepfake videos, as shown in Fig. S1 in the Supplementary Material. This idea motivates us to exploit synchronization between lip movements and speech to detect forgeries in deepfake videos. Therefore, the input of the sync-check module includes the output representation $\vec{F}_v$ of the lip image feature extractor and the output representation $\vec{F}_a$ of the acoustic feature extractor. For each time frame $i$, we calculate the absolute difference between the corresponding lip embedding $\vec{F}_{vi}$ and audio embedding $\vec{F}_{ai}$ to capture the frame-level difference between the visual and acoustic modalities. Consequently, the output of the sync-check module is the sync-based feature vector sequence $\vec{F}_{sync}$ calculated as

$$\vec{F}_{sync} = \{|\vec{F}_{vi} - \vec{F}_{ai}|\}_{i=1}^T, \tag{4}$$

where $T$ is the number of frames in the input video.

### E. Feature Fusion Module

In addition to the sync-based feature vector sequence $\vec{F}_{sync}$, the robust audio-visual feature vector sequence $F_{av}$ obtained from the audio-visual feature extractor also captures the correlation between the two modalities. Therefore, we combine these two audio-visual representations for multimodal forgery

TABLE I
STATISTICS OF MULTIMODAL FOGERY DATASETS FOR DEEPFAKE DETECTION.

| Datasets | Real Videos | Fake Videos | Manipulation Methods | No of Subjects | Visual Manipulation | Audio Manipulation |
|---|---|---|---|---|---|---|
| FakeAVCeleb [33] | 500 | 20000 | Faceswap, Fsgan, wav2lip, RTVC | 500 | Yes | Yes |
| DeepfakeTIMIT [34] | 320 | 320 | Faceswap | 32 | Yes | No |

detection. In the feature fusion module, $\vec{F}_{av}$ and $\vec{F}_{sync}$ are concatenated along the feature dimension to form a fusion representation sequence $\vec{F}_{fusion}$ as

$$\vec{F}_{fusion} = \vec{F}_{av} \oplus \vec{F}_{sync}, \tag{5}$$

where $\oplus$ denotes the concatenation operation.

### F. Temporal Convolutional Network and Classifier

The temporal dynamics across the audio and visual frames contain important information about the video content. To capture inter-modal and intra-modal temporal correlations, we adopt the multi-scale temporal convolutional network (MS-TCN) in [62]. Temporal convolution takes a sequence of frame-level feature vectors and maps them into another sequence of the same dimension using one-dimensional temporal convolution. The temporal convolutional network acts as a sequence encoder, capturing short-term and long-term information by providing the network with visibility across multiple time scales. MS-TCN is followed by a temporal pooling layer and a linear layer for outputting the Real/Fake probability given $\vec{F}_{fusion}$ as

$$\hat{y} = F_{\theta_m}(\vec{F}_{fusion}), \tag{6}$$

where $\hat{y}$ represents the probability of the target class.

### G. Model Training

The model is trained with the cross-entropy loss, defined as

$$L(y,\hat{y}) = -\frac{1}{N}\sum_{i=1}^{N}[y_i \log \hat{y}_i + (1 - y_i) \log (1 - \hat{y}_i)], \tag{7}$$

where $N$ is the number of training samples, $y_i$ represents the ground truth label of the $i$-th sample (0 or 1), and $\hat{y}_i$ represents the class prediction probability of the $i$-th sample. During training, the pre-trained front-end feature extractors and transformer encoder are fine-tuned, while MS-TCN and the linear classifier are trained from scratch.

### H. AV-Lip-Sync+ with FE

As reported in [56], speech-lip synchronization based methods may not be good at detecting fake videos generated by some visual manipulation models such as Faceswap and FS-GAN. In these types of fake videos, the video contains only visual manipulation, while the audio is genuine. Furthermore, the manipulation does not necessarily occur in the lip region, but artifacts can be observed in other regions or throughout the face, including face boundaries. Deepfake detectors based on lip features may fail when the lip region is not manipulated, or when there is little manipulation in the mouth region. To address these issues, we use a face encoder to utilize the entire face features to enhance our proposed deepfake detector and make it more robust and generalizable to deep face manipulation techniques. To this end, we employ the pre-trained ViViT model [32] as the face encoder to extract the spatiotemporal face features. Using tubelet embeddings and spatial and temporal transformers, the face encoder extracts inter- and intra-frame information from video content. The output of the face encoder is a single-vector representation, which is fed to a linear layer for classification. The face-based deepfake detection model is trained on multimodal deepfake dataset using cross-entropy loss. During training, the pre-trained face encoder is fine-tuned, while the linear classifier is trained from scratch.

The model that combines AV-Lip-Sync+ and the face encoder is called AV-Lip-Sync+ with FE. The extracted single-vector representation of the face encoder and the representation obtained from the AV-Lip-Sync+ model are concatenated and fed to a two-layer linear classifier. During training, the pre-trained face encoder and AV-Lip-Sync+ are fixed, and only the classifier is trained using cross-entropy loss.

## IV. EXPERIMENTS

We conducted experiments on two datasets: FakeAVCeleb [33] and DeepfakeTIMIT [34]. Unlike other unimodal audio or video deepfake datasets, these two datasets contain both audio and visual modalities, and their fake samples contain audio and/or visual manipulations. Furthermore, the faces in the videos in both datasets are frontal, which makes them suitable for lip frame extraction as visual input to the proposed model. The statistics of the two dataset are shown in Table I.

### A. Datasets

1) FakeAVCeleb: The FakeAVCeleb dataset is an audio-visual dataset released in 2021 specifically designed for the deepfake detection task. It is based on a collection of 500 YouTube videos featuring 500 celebrities from diverse ethnic regions including South Asia, East Asia, Africa, Europe and America. Fake videos are generated from these 500 real videos using the Faceswap [14], Fsgan [15], wav2lip [40], and real-time voice cloning (SV2TTS) [13] manipulation methods and their combinations. Several examples are shown in Fig. 1. In the case of Faceswap-wav2lip, the video is manipulated using both Faceswap and wav2lip manipulation methods. Similarly, in the case of Fsgan-wav2lip, the video is generated using a combination of Fsgan and wav2lip. The videos manipulated by wav2lip include two types, namely Fake-Video-Real-Audio (FVRA) and Fake-Video-Fake-Audio (FVFA). Wav2lip FVRA videos contain manipulated lips and real audio. In the case of wav2lip FVFA, in addition to lip manipulation, a real-time voice cloning method is also used to manipulate the audio.

In total, the dataset contains 500 real videos and more than 20000 forged videos.

Following [56], [59], we used multiple test sets, namely Faceswap, Fsgan, RTVC, wav2lip, Faceswap-wav2lip, and Fsgan-wav2lip. Furthermore, two other major and diverse test sets are Test-set-1 and Test-set-2. In Test-set-1, the number of samples is the same for all manipulation methods, while in Test-set-2, the number of samples of RVFA (Real-Video-Fake-Audio), FVRA (Fake-Video-Real-Audio), and FVFA (Fake-Video-Fake-Audio) in the fake class is the same. The training-test split is based on the number of subjects in the dataset. The training set and test set contain real and fake videos corresponding to 430 subjects and 70 subjects respectively. Furthermore, all the test sets are balanced in terms of real and fake videos and contain 70 videos per class (real and fake). The training set contains only 430 real videos, which is significantly less than the number of fake videos. If the imbalance problem is not properly resolved, the experimental results will be biased. To eliminate the imbalance problem, we took real videos from the VoxCeleb1 dataset [63] to make the training data of the real and fake classes more balanced and employed a more effective sampling method to rebalance class distributions during model training.

*2) DeepfakeTIMIT:* The DeepfakeTIMIT dataset contains 320 audio-visual human speech recordings from 32 subjects and is a subset of the VidTIMIT dataset [64]. Each subject has 10 videos. For each video, the corresponding fake video is generated by the Faceswap manipulation method. The audio in both real and fake videos is always real. Since the dataset is small, we performed 5-fold cross-validation on it and evaluated the average performance.

### B. Preprocessing

Our model mainly utilizes lip and audio features for multimodal forgery detection. For this purpose, the initial step is to extract the lip region from the frontal face using facial landmarks. We leveraged a pre-trained CNN-based face detector from the Dlib toolkit [65]. The lip image sequence extracted from the input video is $96 \times 96$ RGB pixels. Before the lip image sequences are fed into the model, they are converted to grayscale. The input shape of the extracted lip features is $C \times F \times H \times W$, where $C$ represents the number of channels, $F$ denotes the number of frames, and $H$ and $W$ represent the height and width of each frame, respectively. In addition, for the audio modality, the waveform is extracted from the video and then converted to the log filterbank energies as the acoustic input of the model.

For facial feature extraction, among the different variants of ViViT in [32], we selected the best performing factorized encoder model as the face encoder. The visual encoder has two transformer blocks, namely a spatial transformer and a temporal transformer. The input to the model is short fixed-length video clips from the entire video. The number of frames in a clip is 16, the input frame size is $224 \times 224$, the patch size is 16, the number of input channels is 3, and the embedding dimension is 768. We used the tubelet embedding method. The ViViT model is pre-trained on the kinetics dataset [66].

### TABLE II
EVALUATION RESULTS OF AV-LIP-SYNC+ ON THE FAKEAVCELEB DATASET (IN %).

| Test set | Class | Precision | Recall | F1-score | Accuracy |
|---|---|---|---|---|---|
| Faceswap | Real | 85.19 | 98.57 | 91.39 | 90.71 |
| | Fake | 98.31 | 82.86 | 89.92 | |
| Faceswap_wav2lip | Real | 100.0 | 98.57 | 99.28 | 99.29 |
| | Fake | 98.59 | 100.0 | 99.29 | |
| Fsgan | Real | 86.25 | 99.57 | 92.00 | 91.43 |
| | Fake | 98.33 | 84.29 | 90.77 | |
| Fsgan_wav2lip | Real | 100.0 | 98.57 | 99.28 | 99.29 |
| | Fake | 98.59 | 100.0 | 99.29 | |
| RTVC | Real | 95.83 | 98.57 | 97.18 | 97.14 |
| | Fake | 98.53 | 95.71 | 97.10 | |
| Wav2lip | Real | 100.0 | 98.57 | 99.28 | 99.29 |
| | Fake | 98.59 | 100.0 | 99.29 | |
| Test-set-1 | Real | 93.24 | 98.57 | 95.83 | 95.71 |
| | Fake | 98.48 | 92.86 | 95.59 | |
| Test-set-2 | Real | 98.57 | 98.57 | 98.57 | 98.57 |
| | Fake | 98.57 | 98.57 | 98.57 | |

### C. Model Configuration and Training

Our proposed AV-Lip-Sync+ model contains 132.85M trainable parameters. When incorporating the Face Encoder (FE) module, the parameters increase to 249.37M due to the additional fusion mechanism. To evaluate the inference time, we conducted experiments on an NVIDIA GeForce RTX 2080 Ti GPU. To ensure reliable measurements, the model was executed in evaluation mode, and multiple warm-up iterations were performed to mitigate cold-start overhead. The average inference time per video sample is 0.07 seconds for AV-Lip-Sync+ and 0.14 seconds for AV-Lip-Sync+ with FE. These measurements only consider the forward pass of the model and do not include the time for pre- or post-processing steps.

Our model was trained by the Adam optimizer with a learning rate of 0.00001 and an early stopping strategy for 30 epochs. We added 3570 real videos from the VoxCeleb1 dataset [63] to the real class and employed a more effective sampling method called Imbalanced Dataset Sampler. This method rebalances class distributions during model training by ensuring that each batch contains a balanced number of samples from both classes (real and fake) and automatically estimates the sampling weights. Through this strategy, we not only address the challenge of dataset imbalance, but also overcome the overfitting and information loss issues associated with traditional over-sampling and under-sampling.

The evaluation metrics used include precision, recall, F1-score, accuracy (see Section S2 in the Supplementary Material for details), the receiver operating characteristics (ROC) curve, and the area under the curve (AUC). We report video-level performance rather than frame-level performance.

### D. Results

*1) Evaluation of AV-Lip-Sync+ and AV-Lip-Sync+ with FE:* Table II shows the performance of the proposed AV-Lip-Sync+ model evaluated on the FakeAVCeleb dataset. The manipulation methods of Faceswap, Faceswap-wav2lip, Fsgan, Fsgan-wav2lip, RTVC, and wav2lip are all seen during the model training process. It is obvious that AV-Lip-Sync+ performed well on almost all test sets except Faceswap and Fsgan. The main reason may be that the AV-HuBERT feature extractor only extracts visual information from lip images and ignores

TABLE III
ACCURACY OF AV-LIP-SYNC+ AND AV-LIP-SYNC+ WITH FE ON THE FAKEAVCELEB DATASET (IN %).

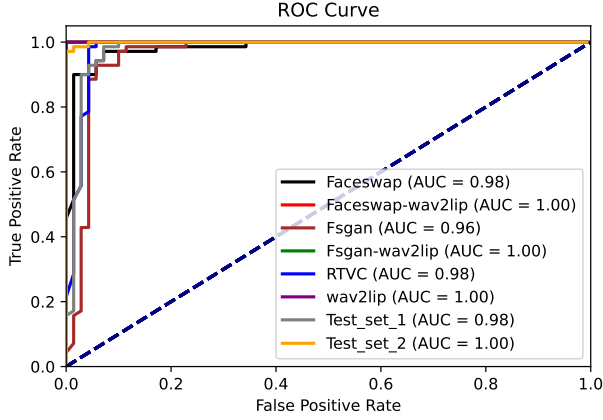| Model | Faceswap | Faceswap_wav2lip | Fsgan | Fsgan_wav2lip | RTVC | Wav2lip | Test-set-1 | Test-set-2 |
|---|---|---|---|---|---|---|---|---|
| AV-Lip-Sync+ | 90.71 | **99.29** | 91.43 | **99.29** | **97.14** | **99.29** | 95.71 | 98.57 |
| AV-Lip-Sync+ with FE | **97.86** | **99.29** | **96.43** | **99.29** | 96.43 | **99.29** | **97.86** | **99.29** |



Fig. 3. ROC curves and AUC scores of the proposed AV-Lip-Sync+ method on various test sets of the FakeAVCeleb dataset.

information outside the lip region. Fig. 3 shows the ROC curves and AUC scores for all test sets. The AUC scores under different manipulation conditions are all above 0.96.

Facial features are crucial for detecting deepfake videos, as local or entire face regions may contain artifacts caused by video manipulation. To address the limitations of AV-Lip-Sync+ on Faceswap and Fsgan test samples, we used a ViViT-based face encoder to inject face embeddings into an ensemble model (called Lip-Sync+ with FE). As shown in Table III, the accuracy of Faceswap and Fsgan test sets increased from 90.71% and 91.43% to 97.14% and 96.43%, respectively. The results show that providing supplementary information to the proposed model can improve the detection of videos with tampered faces through Faceswap and Fsgan techniques.

*2) Comparison of different models:* In this experiment, we compared our models with various existing unimodal, multimodal, fusion, and ensemble deepfake detection models on Test-set-2 of the FakeAVCeleb dataset. The results are shown in Table IV. Several unimodal, multimodal and ensemble models have been evaluated in [67], but the performance of most models is unsatisfactory.Unimodal video-only models VGG16 and LipForensics [19] rely solely on visual features for deepfake detection. VGG16 processes video frames, while LipForensics focuses on lip sequences. Similarly, the unimodal audio-only model Xception is trained using MFCC features extracted from the audio modality. However, these unimodal models are doomed to fail in situations where the focusing modality is genuine but the other modality is manipulated.

The ensemble and multimodal models evaluated in [67], which aggregate the predictions of unimodal classifiers through simple voting without exploiting inter-modal relationships, or were not originally designed for forgery detection (e.g., Multimodal-1, Multimodal-2, and CDCN), did not lead to performance improvements compared to unimodal deepfake detectors. MDS [53], a modality dissonance score-based deepfake detector, also only achieved 69% accuracy.

Later, multimodal ensemble models [58], [59], [68] outperformed unimodal and early multimodal or ensemble models by combining more effective sub-modules. The identity-based POI-Forensic model [71] achieves an accuracy of 85.50%, which is constrained by the identities used during training, limiting its applicability in real-world scenarios, such as detecting deepfakes in user-generated content on social media. The Multimodal Contrastive Learning (MCL) method [72] utilizes contrastive learning to bridge the cross-modal gap and achieves an accuracy of 89.25%. PVASS-MDD [73] employs a two-stage framework: a self-supervised module aligns visual-audio features by predicting audio from visual cues, and a detection stage leverages this alignment to enhance detection of audio-visual inconsistencies in deepfake videos. A-V Anomaly method [74] trains an autoregressive model on real unlabeled videos, detecting manipulated videos with low generation probability, achieving an accuracy of 92.71%. However, it is less effective against manipulations that preserve synchronization, such as those altering a person's appearance without altering mouth motion. An unsupervised method [75] detects deepfakes by spotting intra- and cross-modal inconsistencies, achieving 94.59% accuracy. A zero-shot method [76] detects deepfakes by comparing automatic speech recognition (ASR) and visual speech recognition (VSR) outputs via edit distance, achieving 91.94% accuracy. MMMS-BA framework [77] applies attention over audio-visual sequences to address the modality gap and improve deepfake detection and localization, reaching 97.90% accuracy. The AV-Lip-Sync model [56] uses speech-lip synchronization features based on the difference between the extracted lip sequence and the synthetic lip sequence and achieves an accuracy of 93.57%. It requires the wav2lip [57] module to generate lip sequences from audio, which increases model complexity, training, and inference time. Based on these previous studies, the current state-of-the-art models such as AVTENet [60], RADE [78], and AVFF [79] have shown excellent performance with detection accuracies higher than 98%. These models are based on advanced model architectures, more effective training strategies and cross-modal modeling, or large-scale pre-trained models.

Our proposed AV-Lip-Sync+ model eliminates the need for a wav2lip generator while achieving significant performance improvements compared to the AV-Lip-Sync model. AV-Lip-Sync+ and AV-Lip-Sync+ with FE achieve 98.57% and 99.29% accuracy, respectively. On Test-set-2 of the FakeAVCeleb dataset, AV-Lip-Sync+ achieves state-of-the-art performance, and AV-Lip-Sync+ with FE sets a new state-of-the-art benchmark. These results highlight the effectiveness of leveraging SSL features to model speech-lip synchronization, audio-visual correlations, and spatiotemporal facial artifacts for deepfake video detection.

Ablation studies as well as the discriminant analysis of features and experiments on robustness to partial occlusion

TABLE IV
EVALUATION RESULTS OF DIFFERENT MODELS ON THE FAKEAVCELEB DATASET (IN %).

| Type | Model | Modality | Class | Precision | Recall | F1-score | Accuracy |
|------|-------|----------|-------|-----------|--------|----------|----------|
| Unimodal [67] | VGG16 | V | Real | 69.35 | 89.66 | 78.21 | 81.03 |
|  |  |  | Fake | 87.24 | 77.50 | 82.08 |  |
| Unimodal [67] | Xception | A | Real | 87.50 | 60.87 | 71.79 | 76.26 |
|  |  |  | Fake | 70.33 | 91.43 | 79.50 |  |
| Unimodal [19] | LipForensics | V | Real | 70.00 | 91.00 | 80.00 | 76.00 |
|  |  |  | Fake | 88.00 | 61.00 | 72.00 |  |
| Ensemble (Soft-Voting) [67] | VGG16 | AV | Real | 69.35 | 89.66 | 78.21 | 78.04 |
|  |  |  | Fake | 89.48 | 68.94 | 77.88 |  |
| Ensemble (Hard-Voting) [67] | VGG16 | AV | Real | 69.35 | 89.66 | 78.21 | 78.04 |
|  |  |  | Fake | 89.48 | 68.94 | 77.88 |  |
| Multimodal-1 [67] | Multimodal-1 | AV | Real | 00.00 | 00.00 | 00.00 | 50.00 |
|  |  |  | Fake | 49.60 | 100.0 | 66.30 |  |
| Multimodal-2 [67] | Multimodal-2 | AV | Real | 71.00 | 58.70 | 64.30 | 67.40 |
|  |  |  | Fake | 64.80 | 76.00 | 70.00 |  |
| Multimodal-3 [67] | CDCN | AV | Real | 50.00 | 06.80 | 12.00 | 51.50 |
|  |  |  | Fake | 50.00 | 94.00 | 65.10 |  |
| Multimodal-4 [53] | MDS | AV | Real | 62.16 | 98.57 | 76.24 | 69.29 |
|  |  |  | Fake | 96.55 | 40.00 | 56.57 |  |
| Multimodal-Ensemble [59] | Ensemble model | AV | Real | 83.13 | 98.57 | 90.20 | 89.29 |
|  |  |  | Fake | 98.25 | 80.00 | 88.19 |  |
| Multimodal [68] | Multimodaltrace | AV | Real | - | - | - | 92.90 |
|  |  |  | Fake | - | - | - |  |
| Multimodal-Ensemble [58] | AVFakeNet | AV | Real | - | - | - | 93.41 |
|  |  |  | Fake | - | - | - |  |
| Fusion [69] | AVoiD-DF | AV | Real | - | - | - | 83.70 |
|  |  |  | Fake | - | - | - |  |
| Fusion [70] | MIS-AVioDD | AV | Real | - | - | - | 96.20 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [56] | AV-Lip-Sync | AV | Real | 91.78 | 95.71 | 93.71 | 93.57 |
|  |  |  | Fake | 95.52 | 91.43 | 93.43 |  |
| Identity Aware [71] | POI-Forensics | AV | Real | - | - | - | 85.50 |
|  |  |  | Fake | - | - | - |  |
| Contrastive Learning [72] | MCL | AV | Real | - | - | - | 89.25 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [73] | PVASS-MDD | AV | Real | - | - | - | 95.70 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [74] | A-V Anomaly | AV | Real | - | - | - | 92.71 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [75] | Intra-Cross-modal | AV | Real | - | - | - | 94.59 |
|  |  |  | Fake | - | - | - |  |
| Fusion [76] | Zero-Shot | AV | Real | - | - | - | 91.94 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [77] | MMMS-BA | AV | Real | - | - | - | 97.90 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [78] | FRADE | AV | Real | - | - | - | 98.60 |
|  |  |  | Fake | - | - | - |  |
| Multimodal [79] | AVFF | AV | Real | - | - | - | 98.60 |
|  |  |  | Fake | - | - | - |  |
| Multimodal-Ensemble [60] | AVTENet | AV | Real | 100.0 | 97.14 | 98.55 | 98.57 |
|  |  |  | Fake | 97.22 | 100.0 | 98.59 |  |
| **Multimodal (ours)** | **AV-Lip-Sync+** | **AV** | **Real** | **98.57** | **98.57** | **98.57** | **98.57** |
|  |  |  | **Fake** | **98.57** | **98.57** | **98.57** |  |
| **Multimodal (ours)** | **AV-Lip-Sync+ with FE** | **AV** | **Real** | **100.0** | **98.57** | **99.28** | **99.29** |
|  |  |  | **Fake** | **98.59** | **100.0** | **99.29** |  |

and generalization across datasets are detailed in Sections S3, S4, S5, and S6 in the Supplementary Material, respectively.

*3) Evaluation on the DeepfakeTIMIT dataset:* The DeepfakeTIMIT dataset is primarily used for training and evaluating visual deepfake detectors, as it only contains visual manipulation. It comes in two versions, Low-Quality (LQ) and High-Quality (HQ). For a fair comparison, we compared the proposed model with unimodal visual and multimodal audio-visual deepfake detectors. Additionally, we performed 5-fold cross-validation and reported the average AUC score. As can been seen from Table V, among unimodal visual detectors, FWA and DSP-FWA [83] achieved the best AUC of 99.90 under the LQ condition, and DSP-FWA achieved the best AUC of 99.70 under the HQ condition. Although the five existing audio-visual detectors (Emotions Don't lie [54], MDS [53], AV-Lip-Sync [56], POI-Forensics [71], and MCL [72]) are overall better than most unimodal visual detectors, their per-

formance is worse than that of the best performing unimodal visual detector DSP-FWA. However, our multimodal AV-Lip-Sync+ with FE outperformed all models compared in the table. The AUC score is 99.96 for LQ and 99.98 for HQ. By appropriately integrating pre-trained AV-HuBERT and ViViT models for audio-visual feature extraction, our proposed model achieves state-of-the-art results on the DeepfakeTIMIT dataset.

## V. CONCLUSIONS

In this study, we have used AV-HuBERT and ViViT for the downstream task of audio-visual video forgery detection. We exploit the inconsistency between visual and audio modalities using the powerful audio-visual representation provided by AV-HuBERT, which is pre-trained using self-supervised learning. Since AV-HuBERT only extracts visual features from the lip region, which may not be sufficient to detect artifacts outside the lip region, we also adopt another transformer-based

TABLE V
EVALUATION RESULTS OF AV-LIP-SYNC+ ON THE DEEPFAKETIMIT
DATASET (IN %).

| Type | Model | Modality | Quality | AUC |
|---|---|---|---|---|
| Unimodal [20] | Capsule | V | LQ | 78.40 |
| | | | HQ | 74.40 |
| Unimodal [80] | Multi-task | V | LQ | 62.20 |
| | | | HQ | 55.30 |
| Unimodal [43] | HeadPose | V | LQ | 55.10 |
| | | | HQ | 53.20 |
| Unimodal [81] | Two-stream | V | LQ | 83.50 |
| | | | HQ | 73.50 |
| Unimodal [82] | VA-MLP | V | LQ | 61.40 |
| | | | HQ | 62.10 |
| Unimodal [82] | VA-LogReg | V | LQ | 77.00 |
| | | | HQ | 77.30 |
| Unimodal [23] | Meso-4 | V | LQ | 87.80 |
| | | | HQ | 68.40 |
| Unimodal [47] | Xception-raw | V | LQ | 56.70 |
| | | | HQ | 54.00 |
| Unimodal [47] | Xception-c40 | V | LQ | 75.80 |
| | | | HQ | 70.50 |
| Unimodal [47] | Xception-c23 | V | LQ | 95.90 |
| | | | HQ | 94.40 |
| Unimodal [83] | FWA | V | LQ | 99.90 |
| | | | HQ | 93.20 |
| Unimodal [83] | DSP-FWA | V | LQ | 99.90 |
| | | | HQ | 99.70 |
| Multimodal [54] | Emotions don't lie | AV | LQ | 96.30 |
| | | | HQ | 94.90 |
| Multimodal [53] | MDS | AV | LQ | 97.92 |
| | | | HQ | 96.87 |
| Multimodal [56] | AV-Lip-Sync | AV | LQ | 97.90 |
| | | | HQ | 96.80 |
| Multimodal [71] | POI-Forensics | AV | LQ | 98.20 |
| | | | HQ | 99.20 |
| Multimodal [72] | MCL | AV | LQ | 97.20 |
| | | | HQ | 99.09 |
| Multimodal (ours) | AV-Lip-Sync+ | AV | LQ | 95.80 |
| | | | HQ | 98.80 |
| **Multimodal (ours)** | **AV-Lip-Sync+ with FE** | **AV** | **LQ** | **99.96** |
| | | | **HQ** | **99.98** |

model ViViT to exploit facial features. Overall, our model jointly exploits SSL audio/visual/audio-visual representations, synchronization features, temporal correlation between lip image frames and audio, and spatiotemporal facial features to detect deepfakes. Experimental results on the FakeAVCeleb and DeepfakeTIMIT datasets show that our model outperforms all existing models and achieves new state-of-the-art performance on both datasets. In future work, we will aim to further improve the generalizability of deepfake detection techniques in multimodal settings.

## ACKNOWLEDGMENT

## REFERENCES

[1] C. Vaccari, A. Chadwick, Deepfakes and Disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news, Social Media + Society 6 (2020).

[2] M. Choraś, K. Demestichas, A. Giełczyk, Á. Herrero, P. Ksieniewicz, K. Remoundou, D. Urda, M. Woźniak, Advanced Machine Learning techniques for fake news (online disinformation) detection: A systematic mapping study, Applied Soft Computing (2021) 107050.

[3] A. Ray, Disinformation, deepfakes and democracies: The need for legislative reform, The UNIVERSITY OF NEW SOUTH WALES LAW JOURNAL 44 (2021) 983–1013.

[4] Á. Figueira, L. Oliveira, The current state of fake news: challenges and opportunities, Procedia Computer Science (2017) 817–825.

[5] K. Narayan, H. Agarwal, S. Mittal, K. Thakral, S. Kundu, M. Vatsa, R. Singh, DeSI: Deepfake source identifier for social media, in: Proc. of the CVPR, 2022, pp. 2858–2867.

[6] R. Zellers, A. Holtzman, H. Rashkin, Y. Bisk, A. Farhadi, F. Roesner, Y. Choi, Defending against neural fake news, Vol. 32, 2019.

[7] Y. Yao, B. Viswanath, J. Cryan, H. Zheng, B. Y. Zhao, Automated crowdturfing attacks and defenses in online review systems, in: Proc. of the ACM SIGSAC Conference on Computer and Communications Security, 2017.

[8] K. McGuffie, A. Newhouse, The radicalization risks of GPT-3 and advanced neural language models, arXiv preprint arXiv:2009.06807 (2020).

[9] A. van den Oord, S. Dieleman, H. Zen, K. Simonyan, O. Vinyals, A. Graves, N. Kalchbrenner, A. Senior, K. Kavukcuoglu, WaveNet: A generative model for raw audio, in: Proc. of the ISCA SSW, 2016.

[10] R. Prenger, R. Valle, B. Catanzaro, WaveGlow: A flow-based generative network for speech synthesis, in: Proc. of the ICASSP, 2019.

[11] S. Vasquez, M. Lewis, MelNet: A generative model for audio in the frequency domain, arXiv preprint arXiv:1906.01083 (2019).

[12] Y. Wang, R. Skerry-Ryan, D. Stanton, Y. Wu, R. J. Weiss, N. Jaitly, Z. Yang, Y. Xiao, Z. Chen, S. Bengio, et al., Tacotron: Towards end-to-end speech synthesis, Proc. of the Interspeech (2017).

[13] Y. Jia, Y. Zhang, R. Weiss, Q. Wang, J. Shen, F. Ren, P. Nguyen, R. Pang, I. Lopez Moreno, Y. Wu, et al., Transfer learning from speaker verification to multispeaker text-to-speech synthesis, Vol. 31, 2018.

[14] I. Korshunova, W. Shi, J. Dambre, L. Theis, Fast face-swap using convolutional neural networks, in: Proc. of the ICCV, 2017.

[15] Y. Nirkin, Y. Keller, T. Hassner, FSGAN: Subject agnostic face swapping and reenactment, in: Proc. of the ICCV, 2019.

[16] N. C. Köbis, B. Doležalová, I. Soraperra, Fooled twice: People cannot detect deepfakes but think they can, Iscience 24 (11) (2021).

[17] A. Diel, T. Lalgi, I. C. Schröter, K. F. MacDorman, M. Teufel, A. Bäuerle, Human performance in detecting deepfakes: A systematic review and meta-analysis of 56 papers, Computers in Human Behavior Reports 16 (2024) 100538.

[18] S. Jia, R. Lyu, K. Zhao, Y. Chen, Z. Yan, Y. Ju, C. Hu, X. Li, B. Wu, S. Lyu, Can ChatGPT detect deepfakes? a study of using multimodal large language models for media forensics, in: Proc. of the CVPRW, 2024, pp. 4324–4333.

[19] A. Haliassos, K. Vougioukas, S. Petridis, M. Pantic, Lips don't lie: A generalisable and robust approach to face forgery detection, in: Proc. of the CVPR, 2021.

[20] H. H. Nguyen, J. Yamagishi, I. Echizen, Capsule-Forensics: Using capsule networks to detect forged images and videos, in: Proc. of the ICASSP, 2019.

[21] K. Lutz, R. Bassett, Deepfake detection with inconsistent head poses: Reproducibility and analysis, arXiv preprint arXiv:2108.12715 (2021).

[22] F. Chollet, Xception: Deep learning with depthwise separable convolutions, in: Proc. of the CVPR, 2017.

[23] D. Afchar, V. Nozick, J. Yamagishi, I. Echizen, MesoNet: a compact facial video forgery detection network, in: Proc. of the WIFS, 2018.

[24] L. Wang, Y. Yoshida, Y. Kawakami, S. Nakagawa, Relative phase information for detecting human speech and spoofed speech, in: Proc. of the Interspeech, 2015.

[25] M. Todisco, H. Delgado, N. W. Evans, A new feature for automatic speaker verification anti-spoofing: Constant Q cepstral coefficients., in: Proc. of the Odyssey, Vol. 2016, 2016.

[26] T. B. Patel, H. A. Patil, Combining evidences from mel cepstral, cochlear filter cepstral and instantaneous frequency features for detection of natural vs. spoofed speech, in: Proc. of the Interspeech, 2015.

[27] S. Scheliga, T. Kellermann, A. Lampert, R. Rolke, M. Spehr, U. Habel, Neural correlates of multisensory integration in the human brain: an ale meta-analysis, Reviews in the Neurosciences 34 (2) (2023) 223–245.

[28] S.-F. Zhang, J.-H. Zhai, B.-J. Xie, Y. Zhan, X. Wang, Multimodal representation learning: Advances, trends and challenges, in: Prof. of the ICMLC, 2019, pp. 1–6.

[29] B. Shi, W.-N. Hsu, K. Lakhotia, A. Mohamed, Learning audio-visual speech representation by masked multimodal cluster prediction, in: Proc. of the ICLR, 2021.

[30] I.-C. Chern, K.-H. Hung, Y.-T. Chen, T. Hussain, M. Gogate, A. Hussain, Y. Tsao, J.-C. Hou, Audio-visual speech enhancement and separation by utilizing multi-modal self-supervised embeddings, in: Proc. of the ICASSPW, 2023.

[31] B. Shi, W.-N. Hsu, A. Mohamed, Robust self-supervised audio-visual speech recognition, in: Proc. of the Interspeech, 2022.

[32] A. Arnab, M. Dehghani, G. Heigold, C. Sun, M. Lučić, C. Schmid, ViViT: A video vision transformer, in: Proc. of the ICCV, 2021.

[33] H. Khalid, S. Tariq, M. Kim, S. S. Woo, FakeAVCeleb: A novel audio-video multimodal deepfake dataset, in: Proc. of the NeurIPS Track on Datasets and Benchmarks, 2021.

[34] P. Korshunov, S. Marcel, DeepFakes: a new threat to face recognition? assessment and detection, arXiv preprint arXiv:1812.08685 (2018).

[35] R. Chen, X. Chen, B. Ni, Y. Ge, SimSwap: An efficient framework for high fidelity face swapping, in: Proc. of the ACM MM, 2020.

[36] G. Gao, H. Huang, C. Fu, Z. Li, R. He, Information bottleneck disentanglement for identity swapping, in: Proc. of the CVPR, 2021.

[37] J. Thies, M. Zollhöfer, M. Nießner, Deferred neural rendering: Image synthesis using neural textures, ACM Transactions on Graphics 38 (4) (2019) 1–12.

[38] L. Li, J. Bao, H. Yang, D. Chen, F. Wen, FaceShifter: Towards high fidelity and occlusion aware face swapping, in: Proc. of the CVPR, 2020.

[39] J. Thies, M. Zollhofer, M. Stamminger, C. Theobalt, M. Nießner, Face2Face: Real-time face capture and reenactment of RGB videos, in: Proc. of the CVPR, 2016.

[40] K. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, C. Jawahar, A lip sync expert is all you need for speech to lip generation in the wild, in: Proc. of the ACM MM, 2020.

[41] S. Lyu, Deepfake detection: Current challenges and next steps, in: Proc. of the ICMEW, 2020.

[42] Y. Li, M.-C. Chang, S. Lyu, In Ictu Oculi: Exposing ai created fake videos by detecting eye blinking, in: Proc. of the WIFS, 2018.

[43] X. Yang, Y. Li, S. Lyu, Exposing deep fakes using inconsistent head poses, in: Proc. of the ICASSP, 2019.

[44] J. Hu, X. Liao, J. Liang, W. Zhou, Z. Qin, FInfer: Frame inference-based deepfake detection for high-visual-quality videos, in: Proc. of the AAAI, Vol. 36, 2022.

[45] Z. Hu, H. Xie, L. Yu, X. Gao, Z. Shang, Y. Zhang, Dynamic-aware federated learning for face forgery video detection, ACM Transactions on Intelligent Systems and Technology 13 (4) (2022) 1–25.

[46] X. Tu, Z. Ma, J. Zhao, G. Du, M. Xie, J. Feng, Learning generalizable and identity-discriminative representations for face anti-spoofing, ACM Transactions on Intelligent Systems and Technology 11 (5) (2020) 1–19.

[47] A. Rossler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, M. Nießner, FaceForensics++: Learning to detect manipulated facial images, in: Proc. of the ICCV, 2019.

[48] Y. Li, X. Yang, P. Sun, H. Qi, S. Lyu, Celeb-DF: A large-scale challenging dataset for deepfake forensics, in: Proc. of the CVPR, 2020.

[49] B. Dolhansky, J. Bitton, B. Pflaum, J. Lu, R. Howes, M. Wang, C. C. Ferrer, The deepfake detection challenge (DFDC) dataset, arXiv preprint arXiv:2006.07397 (2020).

[50] L. Jiang, R. Li, W. Wu, C. Qian, C. C. Loy, DeeperForensics-1.0: A large-scale dataset for real-world face forgery detection, in: Proc. of the CVPR, 2020.

[51] A. Khan, K. M. Malik, J. Ryan, M. Saravanan, Voice spoofing countermeasures: Taxonomy, state-of-the-art, experimental analysis of generalizability, open challenges, and the way forward, arXiv preprint arXiv:2210.00417 (2022).

[52] Y. Zhu, Y. Chen, Z. Zhao, X. Liu, J. Guo, Local self-attention-based hybrid multiple instance learning for partial spoof speech detection, ACM Transactions on Intelligent Systems and Technology 14 (5) (2023) 1–18.

[53] K. Chugh, P. Gupta, A. Dhall, R. Subramanian, Not made for each other-audio-visual dissonance-based deepfake detection and localization, in: Proc. of the ACM MM, 2020.

[54] T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, D. Manocha, Emotions don't lie: An audio-visual deepfake detection method using affective cues, in: Proc. of the ACM MM, 2020.

[55] Y. Zhou, S.-N. Lim, Joint audio-visual deepfake detection, in: Proc. of the ICCV, 2021.

[56] S. A. Shahzad, A. Hashmi, S. Khan, Y.-T. Peng, Y. Tsao, H.-M. Wang, Lip Sync Matters: A novel multimodal forgery detector, in: Proc. of the APSIPA ASC, 2022.

[57] S. B. Hegde, K. Prajwal, R. Mukhopadhyay, V. P. Namboodiri, C. Jawahar, Visual speech enhancement without a real visual stream, in: Proc. of the WACV, 2021.

[58] H. Ilyas, A. Javed, K. M. Malik, AVFakeNet: A unified end-to-end dense swin transformer deep learning model for audio-visual deepfakes detection, Applied Soft Computing 136 (2023) 110124.

[59] A. Hashmi, S. A. Shahzad, W. Ahmad, C. W. Lin, Y. Tsao, H.-M. Wang, Multimodal forgery detection using ensemble learning, in: Proc. of the APSIPA ASC, 2022.

[60] A. Hashmi, S. A. Shahzad, C.-W. Lin, Y. Tsao, H.-M. Wang, AVTENet: Audio-visual transformer-based ensemble network exploiting multiple experts for video deepfake detection, arXiv preprint arXiv:2310.13103 (2023).

[61] T. Afouras, J. S. Chung, A. Zisserman, LRS3-TED: a large-scale dataset for visual speech recognition, arXiv preprint arXiv:1809.00496 (2018).

[62] B. Martinez, P. Ma, S. Petridis, M. Pantic, Lipreading using temporal convolutional networks, in: Proc. of the ICASSP, 2020.

[63] A. Nagrani, J. S. Chung, A. Zisserman, VoxCeleb: a large-scale speaker identification dataset, in: Proc. of the Interspeech, 2017.

[64] C. Sanderson, The VIDTIMIT database, Tech. rep. (2002).

[65] D. E. King, Dlib-ml: A machine learning toolkit, The Journal of Machine Learning Research 10 (2009) 1755–1758.

[66] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijaya-narasimhan, F. Viola, T. Green, T. Back, P. Natsev, et al., The Kinetics human action video dataset, arXiv preprint arXiv:1705.06950 (2017).

[67] H. Khalid, M. Kim, S. Tariq, S. S. Woo, Evaluation of an audio-video multimodal deepfake dataset using unimodal and multimodal detectors, in: Proc. of the Workshop on Synthetic Multimedia-Audiovisual Deep-fake Generation and Detection, 2021.

[68] M. A. Raza, K. M. Malik, Multimodaltrace: Deepfake detection using audiovisual representation learning, in: Proc. of the CVPR, 2023.

[69] W. Yang, X. Zhou, Z. Chen, B. Guo, Z. Ba, Z. Xia, X. Cao, K. Ren, AVoiD-DF: Audio-visual joint learning for detecting deepfake, IEEE Transactions on Information Forensics and Security 18 (2023) 2015–2029.

[70] V. S. Katamneni, A. Rattani, MIS-AVioDD: Modality invariant and specific representation for audio-visual deepfake detection, arXiv preprint arXiv:2310.02234 (2023).

[71] D. Cozzolino, A. Pianese, M. Nießner, L. Verdoliva, Audio-visual person-of-interest deepfake detection, in: Proc. of the CVPR, 2023.

[72] X. Liu, Y. Yu, X. Li, Y. Zhao, MCL: multimodal contrastive learning for deepfake detection, IEEE Transactions on Circuits and Systems for Video Technology (2023) 2803–2813.

[73] Y. Yu, X. Liu, R. Ni, S. Yang, Y. Zhao, A. C. Kot, PVASS-MDD: Predictive visual-audio alignment self-supervision for multimodal deep-fake detection, IEEE Transactions on Circuits and Systems for Video Technology 34 (8) (2023) 6926–6936.

[74] C. Feng, Z. Chen, A. Owens, Self-supervised video forensics by audio-visual anomaly detection, in: Proc. of the CVPR, 2023, pp. 10491–10503.

[75] M. Tian, M. Khayatkhoei, J. Mathai, W. AbdAlmageed, Unsupervised multimodal deepfake detection using intra-and cross-modal inconsistencies, arXiv preprint arXiv:2311.17088 (2023).

[76] X. Li, Z. Liu, C. Chen, L. Li, L. Guo, D. Wang, Zero-shot fake video detection by audio-visual consistency, in: Proc. of the Interspeech, 2024, pp. 2935–2939.

[77] V. S. Katamneni, A. Rattani, Contextual cross-modal attention for audio-visual deepfake detection and localization, in: Proc. of the IJCB, 2024, pp. 1–11.

[78] F. Nie, J. Ni, J. Zhang, B. Zhang, W. Zhang, Frade: Forgery-aware audio-distilled multimodal learning for deepfake detection, in: Proc. of the ACM MM, 2024, pp. 6297–6306.

[79] T. Oorloff, S. Koppisetti, N. Bonettini, D. Solanki, B. Colman, Y. Ya-coob, A. Shahriyari, G. Bharaj, AVFF: Audio-visual feature fusion for video deepfake detection, in: Proc. of the CVPR, 2024, pp. 27102–27112.

[80] H. H. Nguyen, F. Fang, J. Yamagishi, I. Echizen, Multi-task learning for detecting and segmenting manipulated facial images and videos, in: Proc. of the BTAS, 2019.

[81] X. Han, V. Morariu, P. I. Larry Davis, et al., Two-stream neural networks for tampered face detection, in: Proc. of the CVPRW, 2017.

[82] F. Matern, C. Riess, M. Stamminger, Exploiting visual artifacts to expose deepfakes and face manipulations, in: Proc. of the WACVW, 2019.

[83] Y. Li, S. Lyu, Exposing deepfake videos by detecting face warping artifacts, in: Proc. of the CVPRW, 2019.