# UMMAFormer: A Universal Multimodal-adaptive Transformer Framework for Temporal Forgery Localization

Rui Zhang
zhangrui1997@stu.scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

Hongxia Wang*
hxwang@scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

Mingshan Du
2022226245054@stu.scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

Hanqing Liu
liuhanqing0520@stu.scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

Yang Zhou
yzhoulv@stu.scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

Qiang Zeng
zengqiang@stu.scu.edu.cn
School of Cyber Science and
Engineering, Sichuan University
Chengdu, Sichuan, China

## ABSTRACT

The emergence of artificial intelligence-generated content (AIGC) has raised concerns about the authenticity of multimedia content in various fields. However, existing research for forgery content detection has focused mainly on binary classification tasks of complete videos, which has limited applicability in industrial settings. To address this gap, we propose UMMAFormer, a novel universal transformer framework for temporal forgery localization (TFL) that predicts forgery segments with multimodal adaptation. Our approach introduces a Temporal Feature Abnormal Attention (TFAA) module based on temporal feature reconstruction to enhance the detection of temporal differences. We also design a Parallel Cross-Attention Feature Pyramid Network (PCA-FPN) to optimize the Feature Pyramid Network (FPN) for subtle feature enhancement. To evaluate the proposed method, we contribute a novel Temporal Video Inpainting Localization (TVIL) dataset specifically tailored for video inpainting scenes. Our experiments show that our approach achieves state-of-the-art performance on benchmark datasets, including Lav-DF, TVIL, and Psynd, significantly outperforming previous methods. The code and data are available at https://github.com/ymhzyj/UMMAFormer/.

## CCS CONCEPTS

• **Computing methodologies** → **Computer vision**; • **Applied computing** → **Investigation techniques**.

## KEYWORDS

temporal forgery localization, transformer, multimodal-adaptive

*Corresponding author.

## 1 INTRODUCTION

The rapid development of advanced multimedia editing software enabled by artificial intelligence-generated content (AIGC) [1, 5, 22, 30, 32, 43, 54, 56] has raised concerns about its potential misuse, such as manipulating public opinion and fabricating evidence. This has led to a growing interest in developing methods for detecting manipulated content in multimedia forensics, with a primary focus on deepfake detection [13, 28, 60, 63, 65] in facial and audio media. Despite the promising results demonstrated by these methods in a variety of benchmarks [14, 15, 17, 27, 46, 66], their mainstream adoption by the industry remains limited due to the constraint of binary classification tasks. These methods are inadequate for identifying the temporal boundaries of manipulations, which is crucial for practical applications. Further research is necessary to develop techniques that can accurately locate temporal boundaries of manipulations in multimedia content and promote the responsible use of AIGC for the betterment of society.

Recent studies [6, 11, 21, 58] have proposed a new task called temporal forgery localization (TFL) to overcome the limitations of binary classification in detecting manipulated content in multimedia. TFL aims to locate the start and end timestamps of manipulated segments, providing a wider range of application scenarios and helping users better understand the results of forgery detection. TFL is similar to temporal action localization (TAL) [23, 45, 50] and follows a similar process: pre-processing the video or audio data using a pre-trained feature extractor, enhancing the representation of feature vectors with a designed neural network architecture, and decoding the feature vectors using regression and classification heads to obtain the start and end times of each action segment and their corresponding categories. Note that this process may vary depending on specific task requirements.

**Figure 1: We show a collection of keyframes extracted from manipulated videos where a person has been removed, posing a serious threat to digital evidence integrity. The original person mask is displayed on the left with the corresponding video inpainting results underneath it, while the unmanipulated results are shown on the right with a green background. Of the 144 frames in the video, only 11 involve person removal, and it only takes minor modifications to a small section of the video to achieve this. This manipulation technique has drawn significant attention in forensic video analysis because manipulated videos may be presented as genuine evidence in legal proceedings and can be difficult to detect using classification-based methods.**

TFL tasks present unique challenges compared to TAL. Firstly, real-world scenarios often involve various modalities, including audio-only, visual-only, and audio-visual data, requiring separate models for manipulation detection and potentially delaying TFL technology development. Secondly, unmodified or real samples are essential in TFL, just like background samples, but they are often neglected in TAL. Thirdly, manipulation changes are usually more subtle than action changes, with minor alterations like a single word or short pronunciation time making detection more challenging. Finally, the lack of available datasets is a significant bottleneck for TFL. Most multimedia forgery datasets evaluate manipulation performance over the entire video or audio. Only a few studies have validated TFL performance, limited to a single dataset such as Lav-DF [6] for the visual domain and Psynd [58] for the audio domain. Besides, available datasets for TFL [6, 21, 58] and deepfake detection [17, 27, 66] primarily focus on facial manipulations and speech forgeries, while AIGC still poses a threat in other scenarios. This narrow scope limits the potential applications for forgery detection and TFL. For instance, video inpainting [32] techniques can remove specific objects from videos, leading to fabricated evidence, as shown in Figure 1. Based on the above observations, we propose the following work.

For different modalities of multimedia, we propose a novel universal multimodal-adaptive transformer framework for TFL called UMMAFormer. The framework aims to predict forgery segments and their corresponding start and end timestamps in untrimmed videos or audios. Transformer-based models [19, 26, 31, 51] have demonstrated excellent performance in various tasks and can adapt to different modality feature inputs. Therefore, we build a universal multimodal adaptive framework based on the transformer block that can be used for TFL tasks involving different modalities of data.

In order to fully utilize real samples, we design a Temporal Feature Abnormal Attention (TFAA) module based on temporal feature reconstruction. Our motivation is based on the observation that compared to TAL which relies on spatial content to recognize specific types of actions, TFL relies more on temporal features that reflect the changes caused by spatial content manipulations. The underlying difference in feature distribution between manipulated and real segments can be considered a universal feature of multimedia manipulation that exists across any modality of input. By incorporating TFAA, our method enhances the detection of temporal differences, leading to improved TFL performance across different input feature modalities.

For analyzing short video clips with subtle variations, Feature Pyramid Network (FPN) [34] is a commonly used solution that effectively enhances subtle features. We further optimize FPN by introducing a parallel structure and proposing a Parallel Cross-Attention Feature Pyramid Network (PCA-FPN). PCA-FPN significantly improves the performance of small manipulated segments localization.

To advance research further, it is critical to create a new dataset for a novel scenario and providing new evaluation benchmarks for advancing research in TFL tasks. We introduce a novel temporal video inpainting localization dataset called TVIL for training and evaluation of TFL tasks. As per our knowledge, we are the first ones to present a TFL dataset that is tailored for video inpainting scenes. Our dataset is built on the YouTube-VOS 2018 [55] dataset. We employ XMEM [10] to annotate segmentation masks for all frames in the dataset, and then use four different video inpainting models [32, 37, 57, 61] to erase objects in random time periods. We acquire 4453 tampered videos with annotations, divided into training, validation, and testing sets according to the same proportions as the original dataset.

We conduct extensive experiments on three benchmark datasets, Lav-DF [6], Psynd [58], and TVIL to evaluate the effectiveness of our proposed method. The results demonstrate that our approach achieves state-of-the-art performance on these datasets, outperforming the previous best results by a significant margin.

In summary, our contributions are:

- We introduce UMMAFormer, a novel universal transformer framework for multimedia temporal forgery localization that can be applied to various modalities of input.
- We propose a TFAA module that enables the model to focus on temporal anomalies caused by spatial content tampering.

- We design PCA-FPN, a parallel cross-attention feature pyramid network, to improve the recognition and localization of ultrashort forgery segments.
- We present TVIL, a novel temporal video inpainting localization dataset, for research on TFL.

## 2 RELATED WORK

### 2.1 Image-Level Forgery Detection

Detecting manipulated content, especially deepfake [30, 43], has become a critical task in multimedia forensics. Significant efforts have been made to enhance image-level face forgery classification [7, 29, 44]. Early studies [41, 46] primarily relied on basic binary classifiers built upon existing backbone networks, suitable only for detecting low-quality generated images. With advancements in deepfake techniques, several approaches have been proposed to capture specific forgery traces. These approaches explore various features, including noise features [9, 64], local texture characteristics [7, 63], and frequency domain anomalies [28, 44], to enhance detection capabilities. Unfortunately, these approaches overlooked the inclusion of temporal-level features, resulting in inconsistencies in discriminations for consecutive video frames due to variations in lighting, environmental factors, and other disturbances. As a consequence, they struggle to accurately differentiate genuine videos from forgeries and fail to identify temporal boundaries of the forgeries within the videos.

### 2.2 Temporal-Level Forgery Detection

Temporal-level forgery detection involves the classification of forgery at video or audio level and the TFL task, which is the main focus of this paper. The availability of various datasets [15, 27, 66] has significantly contributed to the advancement of temporal-level forgery classification methods. Previous research has proposed different approaches to address this challenge. Hu et al.[25] presented a two-stream method, utilizing a temporal-level stream to extract temporal correlation features and analyze deepfake videos. Han et al.[20] introduced a two-stream network that uses temporal information and learnable spatial rich model (SRM) filters to detect fake videos at the video level. Song et al.[47] utilized a symmetric transformer to enhance discrimination consistency between frames for video-level forgery classification. Additionally, Kwak et al.[28] developed a frequency feature masking method to classify real and fake audio in noisy environments. However, existing temporal-level forgery classification approaches usually treat temporal multimedia content as a cohesive entity, mainly focusing on distinguishing between real and manipulated content without verifying the authenticity of specific timestamps. To address this limitation and enhance the practical value of deepfake detection, the TFL task was introduced. Some studies [6, 11, 21, 58] have focused on this task, but there is still room for significant improvement.
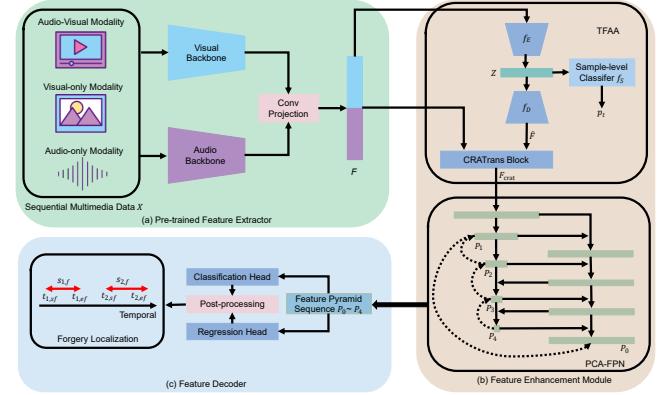
### 2.3 Temporal Action Localization

The goal of TAL is to localize the time intervals in a video when specific actions take place. Existing methods [50] typically followed a general paradigm of feature extraction, feature enhancement, and prediction with post-processing. During the feature extraction stage, most TAL methods typically utilized pre-trained action recognition networks [16, 24, 53] to extract visual or audio-visual features. Given offline features, most algorithms mainly focus on enhancing features, by modeling action boundaries attention [8, 33] and relationships [3, 40, 59]. Some studies [36, 39, 48] also focused on proposing new regression and classification heads to further enhance the localization performance of the model.
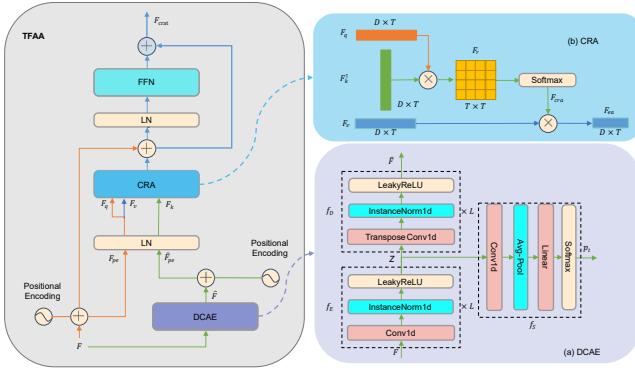
## 3 METHODOLOGY

In this section, we introduce our universal multimodal-adaptive transformer framework, which aims to localizing temporal forgery in sequential multimedia data with various modalities. We have considered three scenarios, including visual-only, audio-only, or joint audio-visual modalities. Of course, the proposed approach can also be further extended to other types of tampered sequential multimedia data.



Figure 2: Illustration of the proposed UMMAFormer, which consists of three main components: (a) a pre-trained feature extractor that maps the sequential multimedia data $X$ to sequential features $F$, (b) a feature enhancement module that enhances the feature representation to multi-scale modified-sensitive features, and (c) a feature decoder that decodes the feature to localize forgeries in the data.

### 3.1 Overview

Our objective is to detect forgeries in untrimmed sequential multimedia data $X$ and locate the corresponding segments. Segments can be represented as $S = \left\{ t_{n,sf}, t_{n,ef}, s_{n,f} \right\}_{n=1}^{N_f}$, where $N_f$ is the number of detected modified segments, $t_{n,sf}, t_{n,ef}$, and $s_{n,f}$ are the start time, the end time, and the confidence score, respectively. To achieve this, $X$ is evenly split into $T$ segments $\{x_t\}_{t=1}^{T}$, and a feature sequence $F \in \mathbb{R}^{C \times T}$ is obtained using TSN [53] and BYOL-A [42] as backbone networks for visual and audio data with concatenation of their features. Our proposed UMMAFormer framework, shown in Figure 2, consists of a pre-trained feature extractor, a feature enhancement module based on a transformer-based network structure composed of the proposed TFAA module and PCA-FPN, and feature decoders for localization. We build on the ActionFormer [59] framework for our approach, with the feature decoding module directly utilized. Our proposed structure can also be extended to other TFL or TAL networks with similar processes.

**Figure 3: Illustration of the proposed TFAA module.**

## 3.2 Temporal Feature Abnormal Attention

To adapt to modified data from different modalities and make full use of real samples, we construct a Temporal Feature Abnormal Attention (TFAA) module built from reconstruction learning and Cross-Reconstruction Attention Transformer (CRATrans) block. The reconstruction learning can be used to determine abnormal states of multi-sensor time-series signals[62]. We believe that temporal features from different modalities can also be viewed as a type of multi-source data. We try to use an encoder-decoder structure to learn the distribution of real samples during the training phase. During the inference phase, we use an attention mechanism to focus on the abnormal segments generated by feature reconstruction, which can adapt well to the features extracted from various modality data. The proposed module is shown in Figure 3.

**Reconstruction Learning.** To be specific, given the encoding feature sequence $F$, we first employ a Deep Convolutional AutoEncoder (DCAE) as illustrated in Figure 3(a) to learn robust representations for real samples. The DCAE consists of a convolutional encoder $f_E$ and a de-convolutional decoder $f_D$. The encoder is composed of $L$ convolutional modules. Each convolution module contain a convolution layer followed by LeakyReLU, and Instance Normalization [49]. The low-dimensional representation $Z \in \mathbb{R}^{C_z \times \frac{T}{2L}}$ and the reconsturcted features $\hat{F}$ can be formulated as follows:

$$\begin{cases} Z = f_E(F), \\ \hat{F} = f_D(Z). \end{cases} \quad (1)$$

The $f_E$ encodes the input features into low-dimensional through convolution layers with a stride of 2. $Z$ represents the latent distribution of real samples. The $f_D$ decodes the latent low-dimensional representation to reconstruct the feature. The decoder is composed of transpose convolution layer, activation layer, and normalization layer.

During the training, we compute the distance between input features and reconstructed features of unmodified samples in a mini-batch as:

$$L_{rec} = \frac{1}{N_r} \sum_i^{N_r} \|\hat{F}_i - F_i\|_1, \quad (2)$$

where $N_r$ is the number of unmodified samples in a mini-batch, and $\|\cdot\|_1$ is the $l_1$-norm. To enhance the consistency of real samples in the low-dimensional embedding space, we utilize a sample-level classifier, denoted as $f_S$, to distinguish the category to which the current feature sequence belongs - whether it is real or tampered. The classifier $f_S$ extracts sample-level features from latent features $Z$ using average pooling and passes them through two fully connected layers to obtain the probability score $p_t$ for the sample being tampered. To address the issue of imbalanced data between real and tampered samples, we utilize the focal loss [35] as the loss function during training. The sample-level focal loss is computed as follows:

$$L_{scls} = -\alpha(1 - p_t)^\gamma \log(p_t), \quad (3)$$

where $\alpha$ is weighting factor to balance positive and negative samples and $\gamma$ is the modulating factor to balance easy and hard samples.

**Cross-Reconstruction Attention Transformer.** Furthermore, many existing anomaly detection algorithms for time series data use reconstruction error to identify abnormal segments. These algorithms set a threshold and flag any segments with reconstruction error above the threshold as anomalous. However, for our task, we need to consider the difference in information carried by different types of samples, which can affect the difficulty of reconstruction and lead to larger errors in some real samples. Additionally, manipulated segments can be very similar to real segments, resulting in small differences in reconstruction. Therefore, directly using reconstruction error to improve our algorithm's performance is difficult.

To address above problem, we introduce a CRATrans module, as shown in Figure 3(b) . As mentioned in [59], transformer block with self-attention module computes a weighted average of features by assigning weights proportional to the similarity score between pairs of input features. In our case, our CRATrans block with Cross-Reconstruction Attention (CRA) will compute similarity scores between pairs of original and reconstructed features in order to replace simple reconstruction errors.

In detail, given the original features $F \in \mathbb{R}^{C \times T}$ and reconstructed features $\hat{F} \in \mathbb{R}^{C \times T}$, we add positional encodings [51] at these features to make position-sensitive feature $F_{pe}$ and $\hat{F}_{pe}$. We believe that positional encodings help to enhance the attention to subtle changes in temporal features. Then we transform them into a latent space by using Layer Normalization (LN) [2] and learnable parameter matrices $\{W_q, W_k\} \in \mathbb{R}^{D \times C}$, respectively. The query $F_q$ and key $F_k$ are calculated by

$$F_q = W_q \left( LN \left( F_{pe} \right) \right), F_k = W_k \left( LN \left( \hat{F}_{pe} \right) \right), \quad (4)$$

where $\{F_q, F_k\} \in \mathbb{R}^{D \times T}$. The original-reconstructed correlation matrix $F_r \in \mathbb{R}^{T \times T}$ is given by

$$F_r = F_q^\top F_k, \quad (5)$$

which represents the similarity between the original features and the reconstructed features in the temporal domain. A CRA matrix $F_{cra}$ is obtained by normalizing the correlation matrix $F_r$, as follows:

$$F_{cra} = Softmax \left( \frac{F_r}{\sqrt{C}} \right), \quad (6)$$

where $Softmax$ is performed row-wise, $\frac{1}{\sqrt{C}}$ is used as the scaling factor. This approach can effectively avoid misjudgment or neglect of abnormalities between the reconstructed and original features due to factors such as scale. Meanwhile, we project the feature $F_{pe}$ to value $F_v \in \mathbb{R}^{D \times T}$ by using the LN and learnable parameter matrix $W_v$:

$$F_v = W_v \left( LN \left( F_{pe} \right) \right). \tag{7}$$

In next step, a dot-product is performed on $F_{cra}$ and the feature $F_v$ to get the representation $F_{ea}$ enhanced by reconstruction anomaly attention. We formulate the function as

$$F_{ea} = F_{cra} F_v, \tag{8}$$

where $F_{ea} \in \mathbb{R}^{D \times T}$. Furthermore, we actually used a Multi-head Cross-Reconstruction Attention(MCRA) for our model, where several CRA operations are concatenated together in parallel.

The output features $F_{ea}$ are added to the original feature $F_{pe}$ and are normalized by the LN layer. Finally, We employ a simple fully connected feed-forward network (FFN) with a residual connection to product the output $F_{crat}$ of CRATrans block.

## 3.3 Parallel Cross-Attention Feature Pyramid Network

High-resolution feature maps are crucial for position-sensitive tasks, such as TFL, which involve numerous short video segments. A multi-scale Transformer encoder was used in [59] to locate action segments in video based on features maps of different resolutions. This encoder utilizes a simple hierarchical multi-scale network, as shown in Figure 4(a). However, the limited representation capability of high-resolution feature maps for complex content poses a challenge. To address this issue, [34] is commonly used to fuse features of different scales to improve the network's temporal localization ability. The scheme of FPN is shown in Figure 4(b). Despite its effectiveness, the fusion process using a simple form of upsampling and downsampling followed by addition usually introduces noise to features of different levels, which may interfere with localization. This effect is particularly pronounced for shorter segments, where even small localization deviations can cause a sharp change in the temporal Intersection over Union (tIoU) between predicted and true values. For example, a segment of 0.5 seconds, when shifted by 0.1 seconds from its correct position, can result in a 20% decrease in tIoU, while for 2 seconds, the tIoU will only decrease by 5%. Inspired by HRNet [52], we propose a Parallel Cross-Attention Feature Pyramid Network (PCA-FPN) to enhance high-resolution features in such cases. The PCA-FPN is illustrated in Figure 4(c), and effectively addresses the problem of noise in feature fusion, improving the localization performance of the network.

The PCA-FPN fuses features of different scales simultaneously through parallel and down-sampling branches, and improves their interaction through a cross-attention (CA) mechanism, Specifically, given the input feature $F_{crat}$ from TFAA module we can encode it to obtain a high-resolution feature map, denoted as $P_0^{in} \in \mathbb{R}^{D_{fpn} \times T}$. Similar to other methods, $P_0^{in}$ is downsampled by an encode module with a factor of 2 to obtain a medium-resolution feature map $P_1^{in} \in \mathbb{R}^{D_{fpn} \times \frac{T}{2}}$. Following [59], the encoder module is a multi-scale transformer unit. To further enhance the representation of the high-resolution feature map $P_0^{in}$, we feed these two different
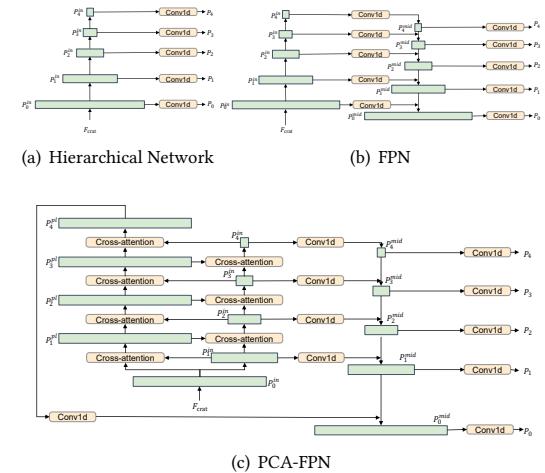
resolution feature maps $P_0^{in}$ and $P_1^{in}$ into the CA module to enhance their features. The CA module is calculated as follows:

$$P_1^{pl} = CA \left( W_{cq} \left( LN \left( P_0^{in} \right) \right), W_{ck} \left( LN \left( P_1^{in} \right) \right), W_{cv} \left( LN \left( P_1^{in} \right) \right) \right), \tag{9}$$

where

$$CA \left( Q, K, V \right) = Softmax \left( \frac{Q^\top g \left( K \right)}{\sqrt{D_{fpn}}} \right) g \left( V \right), \tag{10}$$

$\{ W_{cq}, W_{ck}, W_{cv} \} \in \mathbb{R}^{D_{fpn} \times D_{fpn}}$ are learnable parameter matrices, $g \left( \cdot \right)$ is temporal interpolation function that resamples the $K$ and $V$, which are the inputs of CA to the same size as $Q$, $Softmax$ is performed row-wise, $\frac{1}{\sqrt{D_{fpn}}}$ is used as the scaling factor and $P_1^{pl}$ is the first level parallel high-resolution feature. Subsequently, Subsequently, $P_1^{in}$ used as Query $Q$ and $P_1^{pl}$ used as Key $K$ and Value $V$ to CA module. The output of the CA module is then passed to the multi-scale transformer unit for downsampling to obtain $P_2^{in}$. These processes preserve the feature of short segments in the high-resolution feature map while enhancing the representation of features at different scales. By repeating these processes, we can obtain five levels of parallel multi-scale features $\left\{ P_4^{pl}, P_1^{in}, P_2^{in}, P_3^{in}, P_4^{in} \right\}$. Finally, we fuse the five levels of features from top to bottom, similar to FPN, to obtain the enhanced multi-scale features $P = \{ P_0, P_1, P_2, P_3, P_4 \}$.



(a) Hierarchical Network     (b) FPN

(c) PCA-FPN

**Figure 4: Comparison of feature pyramid networks design in the case of 5 levels.**

## 3.4 Training and Inference

The given feature pyramid $P$ can be decoded into output $S = \left\{ t_{n,sf}, t_{n,ef}, s_{n,f} \right\}_{n=1}^{N_f}$ through classification and regression heads. The final training loss for the overall model is:

$$L = L_{cls} + \lambda_{reg} L_{reg} + \lambda_{rec} L_{rec} + \lambda_{scls} L_{scls}, \tag{11}$$

where $L_{cls}$ and $L_{reg}$ are losses for the classification head outputs $\left\{ s_{n,f} \right\}_{n=1}^{N_f}$ and regression head outputs $\left\{ t_{n,sf}, t_{n,ef} \right\}_{n=1}^{N_f}$, respectively. $L_{cls}$ is binary classification loss, where the label of forgery segments is set to 1 and the rest is set to 0. Other settings are directly adopted from ActionFormer. The reconstruction loss $L_{rec}$

and sample-level focal loss $L_{scls}$ are mentioned in section 3.2. $\lambda_{reg}$, $\lambda_{rec}$ and $\lambda_{scls}$ are hyper-parameters used to balance the relationship between the losses. By default, we set $\lambda_{reg} = 2$, $\lambda_{rec} = 1$, and $\lambda_{scls} = 0.1$.

For the inference stage, we applied Soft-NMS [4] to post-process the results and remove a large number of redundant predictions.

## 4 TEMPORAL VIDEO INPAINTING LOCALIZATION

With the rapid development of AIGC technology, highly deceptive video and audio content has been widely spread on the Internet, leading to potential harm caused by the spread of misleading information. While benchmarks for deepfake videos [15, 27, 46, 66] and audios [17] have emerged in recent years to address the forgery of facial or speech content, these methods only cover a small portion of all forged content. There is a lack of relevant dataset research for other harmful forgery methods. Therefore, we synthesized a dataset for locating video inpainting segments as a new benchmark for TFL, namely TVIL. Our goal is to detect various types of inpainting forgery in sequential images or videos to defend against the spread of misinformation and bring new insights to the research community.

**Data Collection.** The dataset is constructed based on YouTube-VOS 2018 [55], which contains over 4,000 online videos from YouTube. Considering that YouTube is currently one of the most popular video platforms and also an important source for generating and spreading misleading information, we believe that generating a synthesized dataset based on YouTube videos can effectively evaluate the performance of TFL algorithms and prevent the spread of misinformation.

**Data Processing.** YouTube-VOS 2018 is a semi-supervised video semantic segmentation dataset that does not provide complete segmentation masks required for video inpainting. Therefore, we utilized XMEM [10], a state-of-the-art video semantic segmentation algorithm, to generate the segmentation masks. These generated masks can be classified into two types: stationary masks and moving masks [57], which are widely used in real-world scenarios. Stationary masks can be used for removing static objects, simulating the removal of visible watermarks leading to copyright infringement, and so on. On the other hand, moving masks can be used for removing moving objects, simulating the removal of specific targets such as people in surveillance videos. This technology can potentially be used to provide false evidence in certain situations. To better simulate real-world scenarios, we randomly split the dataset into five parts, where one part is used as the real sample set without any manipulation. The remaining four parts are subjected to different video inpainting methods, namely STTN [57], FuseFormer [37], E2FGVI [32] and FGT [61], which randomly removed some frames of the target object. This process aimed to create more diverse and challenging samples to test the effectiveness of the proposed method in handling complex scenarios.

**Dataset Statistics.** We follow the original split in YouTube-VOS 2018, which consisted of 3,471 video clips for training, 474 for validation and 508 for testing. The average length of video clips is about 140 frames, as shown in Figure 5. The training set consists of 3340 forgery segments, the validation set consists of 451 forgery

segments and the test set consists of 463 forgery segments. In our task, video clips with a duration of less than 1 second are defined as short clips. Compared to the Lav-DF [6] dataset, where 89.26% of the manipulated clips are short, our dataset has a proportion of 99.60%, making our dataset more challenging. The distribution of our dataset is illustrated in Figure 6. In Appendix A.1, we provide a further comparison between the TVIL dataset and other commonly used multimedia forensic datasets.
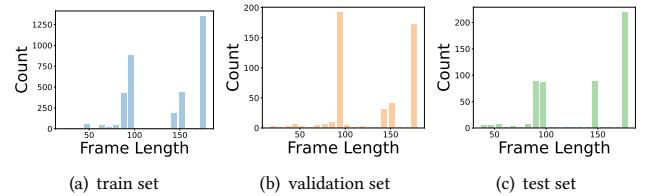


(a) train set     (b) validation set     (c) test set

**Figure 5: Distribution of video lengths.**



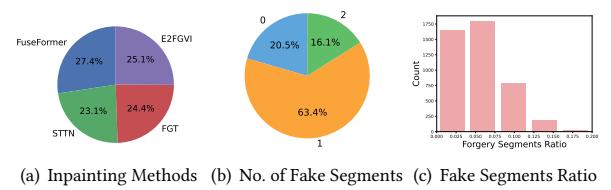(a) Inpainting Methods   (b) No. of Fake Segments   (c) Fake Segments Ratio

**Figure 6: Distribution of the TVIL datasets. (a) The ratio of different methods used in modified video segments. (b) The proportion of manipulated segments in the dataset. (c) The proportion of different manipulated clip lengths to the total length of the video.**

## 5 EXPERIMENTS

### 5.1 Experimental Setup.

For visual data, we use the two-stream TSN [53] network pre-trained on ActivityNet dataset [23] to extract the two-stream visual features. The optical flow is extracted by TV-$L1$ algorithm. The frame interval is set to 1. For audio data, we employ a pre-trained BYOL-A [42] pre-trained on AudioSet [18]. The dimension of the extracted video features is 4096, while that of the audio features is 2048. The extracted features are interpolated to 768 in the temporal dimension.

**Datasets and Evaluation Metric.** We evaluate our method on three benchmark datasets, including Lav-DF [6] for multi-modal data in face forgery scenarios, our proposed TVIL dataset for visual modality in general scenes beyond faces, Psynd [58] for audio modality data in speech scenarios. We follow the evaluation protocol in [6, 21] and report average precision (AP) and average recall (AR) as evaluation metrics, Following conventions, we set the tIoU threshold values at {0.5, 0.75, 0.95} and set Average Number of proposals (AN) to {10, 20, 50, 100}. In addition, for dataset Psynd, we also provide tIoU to follow the protocol of dataset Psynd.

**Baseline and Comparison.** We use ActionFormer [59] as our baseline network and reproduced it based on the official code[1] with default settings on our own datasets. We extend the advanced TAL network, DCAN [8] and TAGS [39], on the TVIL dataset, representing research efforts focused on enhancement of boundary features and improvement of head location, respectively. Additionally, we compare our algorithm with the state-of-the-art methods on each dataset to quantitatively evaluate the performance of our approach.

**Implementation Details.** We follow ActionFormer with minor modifications as follows. Our models are trained on a single RTX 3090 GPU with initial learning rate of 0.001. The batch size for Lav-DF is 32, for TVIL is 16 and for Psynd is 8.

## 5.2 Results for Temporal Face Forgery Localization

We report the AP and AR performance of our method and state-of-the-art methods on the Lav-DF Full Set in Table 1. For the full set, which includes three types of attacks (audio-only modified, video-only modified, and audio-video modified), most unimodal models [11, 36, 40] that focus only on visual information struggle to accurately locate the tampered segments. Although multimodal models [3, 6] perform well in terms of AP at tIoU 0.5, they completely fail for the more challenging AP at tIoU 0.95. The main reason for this is the lack of effective feature enhancement for short video segments. Short segments are extremely sensitive to the tIoU metric. ActionFormer [59] network introduces a simple hierarchical transformer-based network that effectively improves both AP and AR. Furthermore, our method further outperforms BA-TFD[6] by 37.36% in terms of AP at tIoU 0.95 through the proposed PCA-FPN and TFAA with mutilmodal features. For visual-only feature as inputs, our significantly improves AP at tIoU 0.95 from 0.16% to 25.68% compare with BA-TFD. In Appendix A.2, we further present the experimental results for the Lav-DF subset. In Appendix A.3, we provide the forgery classification results and additional evaluation metrics for the Lav-DF Full Set.

**Table 1: Performance comparison on Lav-DF Full Set. Bold faces correspond to the top performance.**

| Methods | Feature | Full Set | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
| MDS [11] | Visual | 12.78 | 1.62 | 0.00 | 37.88 | 36.71 | 34.39 | 32.15 |
| AGT [40] | Visual | 17.85 | 9.42 | 0.11 | 43.15 | 34.23 | 24.59 | 16.71 |
| BMN [36] | Visual | 24.01 | 7.61 | 0.07 | 53.26 | 41.24 | 31.60 | 26.93 |
| BMN (I3D) [36] | Visual | 10.56 | 1.66 | 0.00 | 48.49 | 44.39 | 37.13 | 31.55 |
| AVFusion [3] | Visual+Audio | 65.38 | 23.89 | 0.11 | 62.98 | 59.26 | 54.80 | 52.11 |
| BA-TFD [6] | Visual | 58.55 | 28.60 | 0.16 | 62.49 | 58.77 | 53.86 | 50.29 |
| | Visual+Audio | 76.90 | 38.50 | 0.25 | 66.90 | 64.08 | 60.77 | 58.42 |
| ActionFormer [59] | Visual | 95.34 | 90.20 | 23.73 | 88.41 | 89.63 | 90.33 | 90.41 |
| Ours | Visual | 97.30 | 92.96 | 25.68 | 90.19 | 90.85 | 91.14 | 91.18 |
| | Visual+Audio | **98.83** | **95.54** | **37.61** | **92.10** | **92.42** | **92.47** | **92.48** |

## 5.3 Results for Temporal Video Inpainting Localization

Experimental results in Table 2 show that our method outperforms all compared TAL methods on both the AP and AR evaluations on the proposed TVIL dataset. DCAN [8] is a boundary-enhanced algorithm based on BMN [36] implementation. TAGS [39] is a based

---

[1]https://github.com/happyharrycn/actionformer_release

on a novel localization head that does not include a regression task. It is worth mentioning that due to the increase in in the number of short video clips, the overall performance of ActionFormer is lower than that of Lav-DF. Nevertheless, our method still achieved the best performance, showing the superiority of our method.

**Table 2: Comparison between our method and other state-of-the-art TAL methods on TVIL. Bold faces correspond to the top performance.**

| Methods | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
| --- | --- | --- | --- | --- | --- | --- | --- |
| TAGS [39] | 18.40 | 12.68 | 0.09 | 24.41 | 25.05 | 25.56 | 25.56 |
| DCAN [8] | 82.75 | 75.00 | 3.22 | 64.73 | 66.02 | 68.82 | 69.97 |
| ActionFormer [59] | 86.27 | 83.03 | 28.17 | 84.82 | 85.77 | 88.10 | 88.49 |
| Ours | **88.68** | **84.70** | **62.43** | **87.09** | **88.21** | **90.43** | **91.16** |

## 5.4 Results for Partial Synthetic Speech Localization

We further evaluate UMMAFormer on partial synthetic speech localization task to illustrate its superiority on different modal adaption. Following LFSS [58], we report tIoU on Psynd dataset. LFSS is the only work so far focused on localizing voice cloning partially faked English speech. Due to the absence of completely real samples in the original training set of Psynd, we randomly extracted 299 unaltered audio segments from the original dataset as the real training samples for our proposed method. To ensure fairness, we also selected the best tIoU at different thresholds as the final result, as LFSS did. As shown in Table 3, our algorithm achieves better performance under most conditions, especially in the landline and cellular test subsets. This means that our algorithm has good robustness even when the audio is further disturbed. It is worth mentioning that the test special subset contains both completely fake and completely real samples. Under this condition, LFSS based on simple binary classification achieves good results and outperforms our method by 1.11%. However, for other more challenging scenarios, such as local partial modified segments, especially under conditions where data are disturbed, our proposed method performs better. This further demonstrates the value of research on TFL tasks.

**Table 3: Performance comparison on Psynd in terms of tIoU. Bold faces correspond to the top performance.**

| Methods | test set | special test set | landline | cellular |
| --- | --- | --- | --- | --- |
| LFSS [58] | 98.58 | **99.35** | 80.29 | 80.94 |
| Ours | **98.70** | 98.24 | **92.04** | **86.57** |

## 5.5 Ablation Studies

We compare the contributions of different components of our method for different modalities. Table 4 shows the comparison between the performance of our proposed PCA-FPN and FPN. The experiments are conducted under three scenarios: visual-audio (Lav-DF Full Set), visual-only (TVIL), and audio-only (Psynd-Test). The baseline refers to ActionFormer. We observed that FPN actually reduces model performance in audio tasks because it introduces noise during the multi-scale fusion process. On the other hand, our proposed PCA-FPN greatly improves the localization accuracy of the model.

We also find that PCA-FPN can be applied to localization tasks in different modalities.

**Table 4: Comprasion with FPN. Bold faces correspond to the top performance of each dataset.**

| Dataset | Methods | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
|---------|---------|--------|---------|---------|-------|-------|-------|--------|
| Lav-DF Full Set | Baseline | 97.58 | 93.75 | 40.38 | 92.23 | 92.71 | 92.87 | 92.90 |
| | Baseline+FPN | **98.84** | **95.61** | 38.63 | 92.30 | 92.59 | **92.65** | **92.66** |
| | Baseline+PCA-FPN | 98.72 | 95.52 | **39.00** | **92.31** | **92.60** | **92.65** | **92.66** |
| TVIL | Baseline | 86.10 | 82.86 | 28.11 | 84.68 | 85.71 | 88.04 | 88.43 |
| | Baseline+FPN | 88.50 | 84.35 | 38.95 | **85.91** | 87.26 | **89.63** | **90.09** |
| | Baseline+PCA-FPN | **88.57** | **84.82** | **40.37** | 85.56 | **87.44** | 89.53 | 89.78 |
| Psynd-Test | Baseline | **100.00** | **100.00** | 71.08 | 95.95 | 95.95 | 95.95 | 95.95 |
| | Baseline+FPN | 43.28 | 5.13 | 0.11 | 47.22 | 48.48 | 48.86 | 48.86 |
| | Baseline+PCA-FPN | **100.00** | 98.54 | **77.72** | **97.34** | **97.34** | **97.34** | **97.34** |

Table 5 provides further evidence of the value of TFAA, which was evaluated in the same three scenarios as mentioned earlier. The results show that TFAA effectively improved the model's performance in most scenarios, suggesting that it enhances the applicability of different modality features. Additionally, our model demonstrates the capability of universal modality-adaptation.

**Table 5: Ablation studies of the proposed TFAA modules. Bold faces correspond to the top performance of each dataset.**
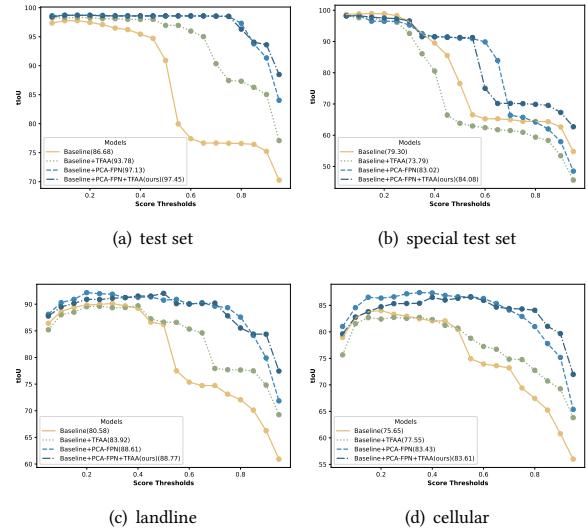
| Dataset | Methods | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
|---------|---------|--------|---------|---------|-------|-------|-------|--------|
| Lav-DF Full Set | Baseline | 97.58 | 93.75 | 40.38 | 92.23 | 92.71 | 92.87 | 92.90 |
| | Baseline+TFAA | 97.57 | 93.74 | **40.53** | **92.31** | **92.80** | **92.98** | **92.99** |
| | Baseline+PCA-FPN+TFAA (ours) | **98.83** | **95.54** | 37.61 | 92.10 | 92.42 | 92.47 | 92.48 |
| TVIL | Baseline | 86.10 | 82.86 | 28.11 | 84.68 | 85.71 | 88.04 | 88.43 |
| | Baseline+TFAA | 85.82 | 83.23 | 51.71 | 86.32 | 87.48 | 89.31 | 89.55 |
| | Baseline+PCA-FPN+TFAA (ours) | **88.68** | **84.70** | **62.43** | **87.09** | **88.21** | **90.43** | **91.16** |
| Psynd-Test | Baseline | **100.00** | **100.00** | 71.08 | 95.95 | 95.95 | 95.95 | 95.95 |
| | Baseline+TFAA | **100.00** | 98.41 | 76.23 | 97.09 | 97.09 | 97.09 | 97.09 |
| | Baseline+PCA-FPN+TFAA (ours) | **100.00** | **100.00** | **79.87** | **97.60** | **97.60** | **97.60** | **97.60** |

Figure 7 demonstrates the impact of different modules on tIoU. We conducted tests on Psynd, varying the confidence scores from 0.05 to 0.95, and compared the resulting tIoU values as well as their averages. A higher average value indicates a higher effectiveness of our predicted candidates. Both TFAA and PCA-FPN effectively improve the model's performance, and combining them results in even better performance.
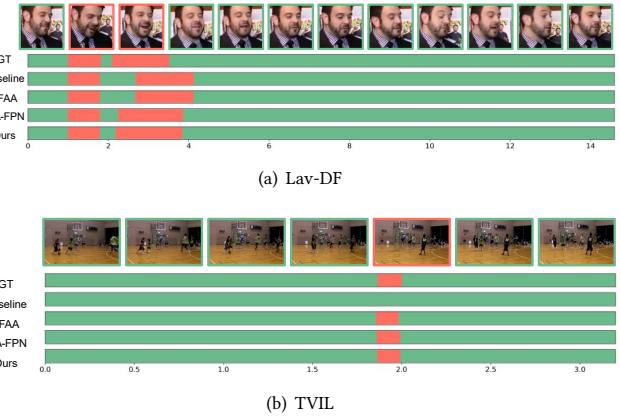
Figure 8 provides visual representations of two qualitative examples from Lav-DF and TVIL. As shown in Figure 8(a), the baseline method can locate the corresponding forged segments, but it exhibits a larger offset compared to our method. As depicted in Figure 8(b), due to the difficulty in locating forged segments that belong to short segments, the baseline method failed to identify them, whereas our method achieved significant detection results.

## 6  CONCLUSIONS

In this paper, we propose a noval universal multimodal-adaptive transformer framework for TFL, which fosters deeper investigations in multimedia content security and helps prevent the misuse of AIGC. To solve the challenges in the task, we propose a novel TFAA module and PCA-FPN to enhance the feature from sequential multimedia data. We also provide a new dataset called TVIL for TFL in a novel scenario which has been released for academic use. The experimental results show the effectiveness of the proposed framework. Especially concerning the LAV-DF dataset, compared to the previous state-of-the-art method BA-TFD [6], our approach has shown significant performance improvements. Specifically, the AP has increased from 76.90% to 98.83% at tIOU 0.5, and from 0.25% to



(a) test set

(b) special test set

(c) landline

(d) cellular

**Figure 7: Ablation studies with respect to effect of score thresholds on tIoU. Score threshold values varies from 0.05 to 0.95 with a step size of 0.05. We calculate the average tIoU over different thresholds.**



(a) Lav-DF

(b) TVIL

**Figure 8: Qualitative examples of our proposed model ablation experiments. Red indicates fake segments and green indicates real segments.**

37.61% at tIOU 0.95. In the future, we will conduct further research to spatial localization on top of temporal localization to enhance the practicality of the model.

# REFERENCES

[1] Nuha Aldausari, Arcot Sowmya, Nadine Marcus, and Gelareh Mohammadi. 2023. Video Generative Adversarial Networks: A Review. *ACM Comput. Surv.* 55, 2 (2023), 30:1–30:25. https://doi.org/10.1145/3487891

[2] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. 2016. Layer Normalization. arXiv:1607.06450

[3] Anurag Bagchi., Jazib Mahmood., Dolton Fernandes., and Ravi Kiran Sarvadevabhatla. 2022. Hear Me out: Fusional Approaches for Audio Augmented Temporal Action Localization. In *Proceedings of the 17th International Joint Conference on Computer Vision, Imaging and Computer Graphics Theory and Applications.* SciTePress, 144–154. https://doi.org/10.5220/0010832700003124

[4] Navaneeth Bodla, Bharat Singh, Rama Chellappa, and Larry S. Davis. 2017. Soft-NMS - Improving Object Detection with One Line of Code. In *Proceedings of the IEEE International Conference on Computer Vision.* IEEE, 5562–5570. https://doi.org/10.1109/ICCV.2017.593

[5] Jiayin Cai, Changlin Li, Xin Tao, Chun Yuan, and Yu-Wing Tai. 2022. DeViT: Deformed Vision Transformers in Video Inpainting. In *Proceedings of the 30th ACM International Conference on Multimedia.* ACM, 779–789. https://doi.org/10.1145/3503161.3548395

[6] Zhixi Cai, Kalin Stefanov, Abhinav Dhall, and Munawar Hayat. 2022. Do You Really Mean That? Content Driven Audio-Visual Deepfake Dataset and Multimodal Method for Temporal Forgery Localization. In *2022 International Conference on Digital Image Computing: Techniques and Applications (DICTA).* IEEE, 1–10. https://doi.org/10.1109/DICTA56598.2022.10034605

[7] Junyi Cao, Chao Ma, Taiping Yao, Shen Chen, Shouhong Ding, and Xiaokang Yang. 2022. End-to-End Reconstruction-Classification Learning for Face Forgery Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 4103–4112. https://doi.org/10.1109/CVPR52688.2022.00408

[8] Guo Chen, Yin-Dong Zheng, Limin Wang, and Tong Lu. 2022. DCAN: Improving Temporal Action Detection via Dual Context Aggregation. In *Proceedings of the AAAI Conference on Artificial Intelligence.* AAAI Press, 248–257.

[9] Guo Chen, Yin-Dong Zheng, Limin Wang, and Tong Lu. 2023. Noise Based Deepfake Detection via Multi-Head Relative-Interaction. In *Proceedings of the AAAI Conference on Artificial Intelligence.* AAAI Press, 14548–14556.

[10] Ho Kei Cheng and Alexander G. Schwing. 2022. XMem: Long-Term Video Object Segmentation with an Atkinson-Shiffrin Memory Model. In *European Conference on Computer Vision.* Springer, 640–658. https://doi.org/10.1007/978-3-031-19815-1_37

[11] Komal Chugh, Parul Gupta, Abhinav Dhall, and Ramanathan Subramanian. 2020. Not made for each other- Audio-Visual Dissonance-based Deepfake Detection and Localization. In *Proceedings of the 28th ACM international conference on Multimedia.* ACM, 439–447. https://doi.org/10.1145/3394171.3413700

[12] Davide Coccomini, Nicola Messina, Claudio Gennaro, and Fabrizio Falchi. 2022. Combining EfficientNet and Vision Transformers for Video Deepfake Detection. In *Proceedings of International Conference on Image Analysis and Processing.* Springer, 219–229. https://doi.org/10.1007/978-3-031-06433-3_19

[13] Jiacheng Deng, Terui Mao, Diqun Yan, Li Dong, and Mingyu Dong. 2022. Detection of Synthetic Speech Based on Spectrum Defects. In *Proceedings of the 1st International Workshop on Deepfake Detection for Audio Multimedia.* ACM, 3–8. https://doi.org/10.1145/3552466.3556529

[14] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. 2020. The DeepFake Detection Challenge (DFDC) Dataset. arXiv:2006.07397

[15] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. 2019. The Deepfake Detection Challenge (DFDC) Preview Dataset. arXiv:1910.08854

[16] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. 2019. Slow-Fast Networks for Video Recognition. In *IEEE International Conference on Computer Vision.* IEEE, 6201–6210. https://doi.org/10.1109/ICCV.2019.00630

[17] Joel Frank and Lea Schönherr. 2021. WaveFake: A Data Set to Facilitate Audio Deepfake Detection. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks,* Vol. 1. Curran Associates, Inc.

[18] Jort F. Gemmeke, Daniel P. W. Ellis, Dylan Freedman, Aren Jansen, Wade Lawrence, R. Channing Moore, Manoj Plakal, and Marvin Ritter. 2017. Audio Set: An ontology and human-labeled dataset for audio events. In *IEEE International Conference on Acoustics, Speech and Signal Processing.* IEEE, 776–780. https://doi.org/10.1109/ICASSP.2017.7952261

[19] Jiwei Guo, Jiajia Tang, Weichen Dai, Yu Ding, and Wanzeng Kong. 2022. Dynamically Adjust Word Representations Using Unaligned Multimodal Information. In *Proceedings of the 30th ACM international conference on Multimedia.* ACM, 3394–3402. https://doi.org/10.1145/3503161.3548137

[20] Bing Han, Xiaoguang Han, Hua Zhang, Jingzhi Li, and Xiaochun Cao. 2021. Fighting Fake News: Two Stream Network for Deepfake Detection via Learnable SRM. *IEEE Trans. Biom. Behav. Identity Sci.* 3, 3 (2021), 320–331. https://doi.org/10.1109/TBIOM.2021.3065735

[21] Yinan He, Bei Gan, Siyu Chen, Yichun Zhou, Guojun Yin, Luchuan Song, Lu Sheng, Jing Shao, and Ziwei Liu. 2021. ForgeryNet: A Versatile Benchmark for Comprehensive Forgery Analysis. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 4360–4369. https://doi.org/10.1109/CVPR46437.2021.00434

[22] Sindhu B. Hegde, K. R. Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C. V. Jawahar. 2022. Lip-to-Speech Synthesis for Arbitrary Speakers in the Wild. In *Proceedings of the 30th ACM International Conference on Multimedia.* ACM, 6250–6258. https://doi.org/10.1145/3503161.3548081

[23] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. 2015. ActivityNet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 961–970. https://doi.org/10.1109/CVPR.2015.7298698

[24] Shawn Hershey, Sourish Chaudhuri, Daniel P. W. Ellis, Jort F. Gemmeke, Aren Jansen, R. Channing Moore, Manoj Plakal, Devin Platt, Rif A. Saurous, Bryan Seybold, Malcolm Slaney, Ron J. Weiss, and Kevin W. Wilson. 2017. CNN architectures for large-scale audio classification. In *IEEE International Conference on Acoustics, Speech and Signal Processing.* IEEE, 131–135. https://doi.org/10.1109/ICASSP.2017.7952132

[25] Juan Hu, Xin Liao, Wei Wang, and Zheng Qin. 2022. Detecting Compressed Deepfake Videos in Social Networks Using Frame-Temporality Two-Stream Convolutional Network. *IEEE Trans. Circuits Syst. Video Technol.* 32, 3 (2022), 1089–1102. https://doi.org/10.1109/CVPRW.2017.229

[26] Ronghang Hu and Amanpreet Singh. 2021. UniT: Multimodal Multitask Learning with a Unified Transformer. In *Proceedings of the IEEE International Conference on Computer Vision.* IEEE, 1419–1429. https://doi.org/10.1109/ICCV48922.2021.00147

[27] Hasam Khalid, Shahroz Tariq, Minha Kim, and Simon Woo. 2021. FakeAVCeleb: A Novel Audio-Video Multimodal Deepfake Dataset. In *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks,* Vol. 1. Curran Associates, Inc.

[28] Il-Youp Kwak, Sunmook Choi, Jonghoon Yang, Yerin Lee, Soyul Han, and Seungsang Oh. 2022. Low-quality Fake Audio Detection through Frequency Feature Masking. In *Proceedings of the 1st International Workshop on Deepfake Detection for Audio Multimedia.* ACM, 9–17. https://doi.org/10.1145/3552466.3556533

[29] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. 2020. Face X-Ray for More General Face Forgery Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 5000–5009. https://doi.org/10.1109/CVPR42600.2020.00505

[30] Yudong Li, Xianxu Hou, Zhe Zhao, Linlin Shen, Xuefeng Yang, and Kimmo Yan. 2022. Talk2Face: A Unified Sequence-based Framework for Diverse Face Generation and Analysis Tasks. In *Proceedings of the 30th ACM International Conference on Multimedia.* ACM, 4594–4604. https://doi.org/10.1145/3503161.3548205

[31] Zekang Li, Zongjia Li, Jinchao Zhang, Yang Feng, and Jie Zhou. 2021. Bridging Text and Video: A Universal Multimodal Transformer for Audio-Visual Scene-Aware Dialog. *IEEE ACM Trans. Audio Speech Lang. Process.* 29 (2021), 2476–2483. https://doi.org/10.1109/TASLP.2021.3065823

[32] Zhen Li, Chengze Lu, Jianhua Qin, Chun-Le Guo, and Ming-Ming Cheng. 2022. Towards An End-to-End Framework for Flow-Guided Video Inpainting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 17541–17550. https://doi.org/10.1109/CVPR52688.2022.01704

[33] Chuming Lin, Chengming Xu, Donghao Luo, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. 2021. Learning Salient Boundary Feature for Anchor-free Temporal Action Localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 3320–3329. https://doi.org/10.1109/CVPR46437.2021.00333

[34] Tsung-Yi Lin, Piotr Dollár, Ross B. Girshick, Kaiming He, Bharath Hariharan, and Serge J. Belongie. 2017. Feature Pyramid Networks for Object Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition.* IEEE, 936–944. https://doi.org/10.1109/CVPR.2017.106

[35] Tsung-Yi Lin, Priya Goyal, Ross B. Girshick, Kaiming He, and Piotr Dollár. 2017. Focal Loss for Dense Object Detection. In *Proceedings of the IEEE International Conference on Computer Vision.* IEEE, 2999–3007. https://doi.org/10.1109/ICCV.2017.324

[36] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. 2019. BMN: Boundary-Matching Network for Temporal Action Proposal Generation. In *Proceedings of the IEEE International Conference on Computer Vision.* IEEE, 3888–3897. https://doi.org/10.1109/ICCV.2019.00399

[37] Rui Liu, Hanming Deng, Yangyi Huang, Xiaoyu Shi, Lewei Lu, Wenxiu Sun, Xiaogang Wang, Jifeng Dai, and Hongsheng Li. 2021. FuseFormer: Fusing Fine-Grained Information in Transformers for Video Inpainting. In *Proceedings of the IEEE International Conference on Computer Vision.* IEEE, 14020–14029. https://doi.org/10.1109/ICCV48922.2021.01378

[38] Trisha Mittal, Ritwik Sinha, Viswanathan Swaminathan, John P. Collomosse, and Dinesh Manocha. 2023. Video Manipulations Beyond Faces: A Dataset with Human-Machine Analysis. In *IEEE/CVF Winter Conference on Applications of Computer Vision Workshops.* IEEE, 643–652. https://doi.org/10.1109/WACVW58289.2023.00071

[39] Sauradip Nag, Xiatian Zhu, Yi-Zhe Song, and Tao Xiang. 2022. Proposal-Free Temporal Action Detection via Global Segmentation Mask Learning. In *European Conference on Computer Vision.* Springer, 645–662. https://doi.org/10.1007/978-

3-031-20062-5_37

[40] Megha Nawhal and Greg Mori. 2021. Activity Graph Transformer for Temporal Action Localization. arXiv:2101.08540

[41] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. 2019. Capsule-forensics: Using Capsule Networks to Detect Forged Images and Videos. In *IEEE International Conference on Acoustics, Speech and Signal Processing*. IEEE, 2307–2311. https://doi.org/10.1109/ICASSP.2019.8682602

[42] Daisuke Niizumi, Daiki Takeuchi, Yasunori Ohishi, Noboru Harada, and Kunio Kashino. 2021. BYOL for Audio: Self-Supervised Learning for General-Purpose Audio Representation. In *International Joint Conference on Neural Networks*. IEEE, 1–8. https://doi.org/10.1109/IJCNN52387.2021.9534474

[43] Ivan Perov, Daiheng Gao, Nikolay Chervoniy, Kunlin Liu, Sugasa Marangonda, Chris Umé, Mr. Dpfks, Carl Shift Facenheim, Luis RP, Jian Jiang, Sheng Zhang, Pingyu Wu, Bo Zhou, and Weiming Zhang. 2021. DeepFaceLab: Integrated, flexible and extensible face-swapping framework. https://doi.org/10.1016/j.patcog.2023.109628 arXiv:2005.05535

[44] Yuyang Qian, Guojun Yin, Lu Sheng, Zixuan Chen, and Jing Shao. 2020. Thinking in Frequency: Face Forgery Detection by Mining Frequency-Aware Clues. In *European Conference on Computer Vision*. Springer, 86–103. https://doi.org/10.1007/978-3-030-58610-2_6

[45] Yifan Ren, Xing Xu, Fumin Shen, Yazhou Yao, and Huimin Lu. 2021. CAA: Candidate-Aware Aggregation for Temporal Action Detection. In *Proceedings of the 29th ACM international conference on Multimedia*. ACM, 4930–4938. https://doi.org/10.1145/3474085.3475616

[46] Andreas Rössler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. 2019. FaceForensics++: Learning to Detect Manipulated Facial Images. In *IEEE International Conference on Computer Vision*. IEEE, 1–11. https://doi.org/10.1109/ICCV.2019.00009

[47] Luchuan Song, Xiaodan Li, Zheng Fang, Zhenchao Jin, Yuefeng Chen, and Chenliang Xu. 2022. Face Forgery Detection via Symmetric Transformer. In *Proceedings of the 30th ACM International Conference on Multimedia*. ACM, 4102–4111. https://doi.org/10.1145/3503161.3547806

[48] Jing Tan, Jiaqi Tang, Limin Wang, and Gangshan Wu. 2021. Relaxed Transformer Decoders for Direct Action Proposal Generation. In *IEEE International Conference on Computer Vision*. IEEE, 13506–13515. https://doi.org/10.1109/ICCV48922.2021.01327

[49] Dmitry Ulyanov, Andrea Vedaldi, and Victor Lempitsky. 2017. Instance Normalization: The Missing Ingredient for Fast Stylization. arXiv:1607.08022

[50] Elahe Vahdani and Yingli Tian. 2023. Deep Learning-Based Action Detection in Untrimmed Videos: A Survey. *IEEE Trans. Pattern Anal. Mach. Intell.* 45, 4 (2023), 4302–4320. https://doi.org/10.1109/TPAMI.2022.3193611

[51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc., 5998–6008.

[52] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, Wenyu Liu, and Bin Xiao. 2021. Deep High-Resolution Representation Learning for Visual Recognition. *IEEE Trans. Pattern Anal. Mach. Intell.* 43, 10 (2021), 3349–3364. https://doi.org/

[53] Limin Wang, Yuanjun Xiong, Zhe Wang, Yu Qiao, Dahua Lin, Xiaoou Tang, and Luc Van Gool. 2016. Temporal Segment Networks: Towards Good Practices for Deep Action Recognition. In *European Conference on Computer Vision*. Springer, 20–36. https://doi.org/10.1007/978-3-319-46484-8_2

[54] Yongqi Wang and Zhou Zhao. 2022. FastLTS: Non-Autoregressive End-to-End Unconstrained Lip-to-Speech Synthesis. In *Proceedings of the 30th ACM International Conference on Multimedia*. ACM, 5678–5687. https://doi.org/10.1145/3503161.3548194

[55] Ning Xu, Linjie Yang, Yuchen Fan, Jianchao Yang, Dingcheng Yue, Yuchen Liang, Brian L. Price, Scott Cohen, and Thomas S. Huang. 2018. YouTube-VOS: Sequence-to-Sequence Video Object Segmentation. In *European Conference on Computer Vision*. Springer, 603–619. https://doi.org/10.1007/978-3-030-01228-1_36

[56] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. 2023. Diffusion Models: A Comprehensive Survey of Methods and Applications. arXiv:2209.00796

[57] Yanhong Zeng, Jianlong Fu, and Hongyang Chao. 2020. Learning Joint Spatial-Temporal Transformations for Video Inpainting. In *European Conference on Computer Vision*. Springer, 528–543. https://doi.org/10.1007/978-3-030-58517-4_31

[58] Bowen Zhang and Terence Sim. 2022. Localizing Fake Segments in Speech. In *26th International Conference on Pattern Recognition*. IEEE, 3224–3230. https://doi.org/10.1109/ICPR56361.2022.9956134

[59] Chen-Lin Zhang, Jianxin Wu, and Yin Li. 2022. ActionFormer: Localizing Moments of Actions with Transformers. In *European Conference on Computer Vision*. Springer, 492–510. https://doi.org/10.1007/978-3-031-19772-7_29

[60] Daichi Zhang, Fanzhao Lin, Yingying Hua, Pengju Wang, Dan Zeng, and Shiming Ge. 2022. Deepfake Video Detection with Spatiotemporal Dropout Transformer. In *Proceedings of the 30th ACM International Conference on Multimedia*. ACM, 5833–5841. https://doi.org/10.1145/3503161.3547913

[61] Kaidong Zhang, Jingjing Fu, and Dong Liu. 2022. Flow-Guided Transformer for Video Inpainting. In *European Conference on Computer Vision*. Springer, 74–90. https://doi.org/10.1007/978-3-031-19797-0_5

[62] Yuxin Zhang, Yiqiang Chen, Jindong Wang, and Zhiwen Pan. 2023. Unsupervised Deep Anomaly Detection for Multi-Sensor Time-Series Signals. *IEEE Trans. Knowl. Data Eng.* 35, 2 (2023), 2118–2132. https://doi.org/10.1109/TKDE.2021.3102110

[63] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. 2021. Multi-Attentional Deepfake Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2185–2194. https://doi.org/10.1109/CVPR46437.2021.00222

[64] Peng Zhou, Xintong Han, Vlad I. Morariu, and Larry S. Davis. 2017. Two-Stream Neural Networks for Tampered Face Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*. IEEE, 1831–1839. https://doi.org/10.1109/CVPRW.2017.229

[65] Yipin Zhou and Ser-Nam Lim. 2021. Joint Audio-Visual Deepfake Detection. In *IEEE International Conference on Computer Vision*. IEEE, 14780–14789. https://doi.org/10.1109/ICCV48922.2021.01453

[66] Bojia Zi, Minghao Chang, Jingjing Chen, Xingjun Ma, and Yu-Gang Jiang. 2020. WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection. In *Proceedings of the 28th ACM International Conference on Multimedia*. ACM, 2382–2390. https://doi.org/10.1145/3394171.3413769

10.1109/TPAMI.2020.2983686

# A APPENDIX

## A.1 Comparison between Existing Audio and Video Forensics Datasets

We present a comprehensive analysis of the recently popular datasets for audio and video forensics. Table 6 summarizes the benchmark datasets that have been used for research on detecting generative content of audio and video, particularly in deepfake detection. While most of the existing datasets [66] focused on simple binary classification tasks related to facial image manipulation, recent advancements in deepfake detection technology resulted in the emergence of binary classification tasks for audio [17] and multimodal audio-visual data [27]. These datasets involve not only facial but also audio information, and their emergence is indicative of the rapid development of deepfake detection technology. Lav-DF [6] and Psynd [58], two other emerging datasets for TFL, are manipulated based on semantic content, making the attacks on these types of videos more similar to real-world scenarios. However, research on these datasets has remained limited to scenes related to human faces and speeches, which is only a small part of the AIGC task. Furthermore, while there are emerging classification datasets [38] for tampering detection beyond face, the quantity of these datasets is limited due to manual generation. To expand the scope of research, it is important to develop more diverse datasets for TFL tasks. Our datasets, similar to TFL subset of ForgeryNet [21] , uses a random approach to generate segments, which facilitates research on TFL tasks beyond facial images and beyond binary classification tasks. Moreover, our generation process can be reproduced under low-cost conditions. Based on our research, the dataset can further expand to more diverse scenarios and promote TFL task research.

**Table 6: Quantitative comparison of TVIL to existing popular Video and Audio Forensics Datasets in recent 3 years. Cls: Classification; SL: Spatial Localization; TFL: Temporal Forgery Localization; V: Visual; A: Audio.**

| Dataset | Year | Tasks | Modality | Application | Manipulated | # Attacks | #Real | #Fake |
|---|---|---|---|---|---|---|---|---|
| WildDeepfake [66] | 2021 | Cls | V | Face | AIGC | - | 3,805 | 3,509 |
| FakeAVCeleb [27] | 2021 | Cls | A+V | Face | AIGC | 3 | 570 | 19,500 |
| ForgeryNet [21] | 2021 | SL/TFL/Cls | V | Face | AIGC | 5 | 99,630 | 121,617 |
| Lav-DF [6] | 2022 | TFL/Cls | A+V | Face | AIGC | 2 | 36,431 | 99,873 |
| WaveFake [17] | 2021 | Cls | A | Speech | AIGC | 6 | 18,100 | 104,885 |
| Psynd [58] | 2022 | TFL | A | Speech | AIGC | 1 | 30 | 2371 |
| VideoSham [38] | 2023 | Cls | A+V | Video Manipulation | User Generated | 40 | 413 | 413 |
| TVIL(Ours) | 2023 | TFL | V | Video Manipulation | AIGC | 4 | 914 | 3539 |

## A.2 More Experiments Results for LAV-DF Subset

We present the AP and AR performance of our method and state-of-the-art algorithms on the Lav-DF subset in Table 7. This subset exclusively contains manipulated videos with visual forgeries, excluding those with audio-only modifications. The results show that the single-modal algorithms in Table 7 outperform their counterparts in Table 2, validating the efficacy of using visual information alone in this context. Notably, our method achieves state-of-the-art performance with an AP@0.5 of 98.83% using solely visual modality.

Despite our method's adaptability to different modalities, applying techniques like contrastive learning to analyze modal inconsistencies during multi-modal fusion posed challenges. Introducing the audio modality in this subset, where authenticity is independent of audio, might lead to a decrease in overall performance. Feature fusion could potentially confuse critical features within this specific subset, as evident from the performance drop of the multi-modal algorithm AVFusion [3].

Nevertheless, our proposed method consistently outperforms other algorithms for both uni- and multi-modal inputs. Moreover, by incorporating audio features, our model achieves a substantial 7.41% improvement in AP at tIoU 0.95, demonstrating the robustness and adaptability of our approach in multi-modal scenarios.

**Table 7: Performance comparison on Lav-DF Sub Set. Bold faces correspond to the top performance.**

| Methods | Feature | Sub Set | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
| MDS [11] | Visual | 23.43 | 3.48 | 00.00 | 58.53 | 56.68 | 53.16 | 49.67 |
| AGT [40] | Visual | 15.69 | 10.69 | 00.15 | 49.11 | 40.31 | 31.70 | 23.13 |
| BMN [36] | Visual | 32.32 | 11.38 | 00.14 | 59.69 | 48.17 | 39.01 | 34.17 |
| BMN (I3D) [36] | Visual | 28.10 | 5.47 | 00.01 | 55.49 | 54.44 | 52.14 | 47.72 |
| AVFusion [3] | Visual+Audio | 62.01 | 22.77 | 00.11 | 61.98 | 58.08 | 53.31 | 50.52 |
| BA-TFD [6] | Visual | 83.55 | 41.88 | 00.24 | 65.79 | 62.30 | 57.95 | 55.34 |
| | Visual+Audio | 85.20 | 47.06 | 00.29 | 67.34 | 64.52 | 61.19 | 59.32 |
| ActionFormer [59] | Visual | 98.06 | 94.43 | 27.25 | 91.30 | 92.04 | 92.27 | 92.28 |
| Ours | Visual | **98.83** | **95.95** | 30.11 | **92.32** | **92.65** | **92.74** | **92.75** |
| | Visual+Audio | 98.54 | 94.30 | **37.52** | 91.61 | 91.97 | 92.06 | 92.06 |

## A.3 More Experiments Results for Video-level Face Forgery Classification

We also conducted a comparison between our method and previous deepfake detection methods on the Lav-DF Full Set for the video-level forgery classification task. The evaluation metric used is the Area Under the Receiver Operating Characteristic Curve (AUC), and the results are summarized in Table 8. In our approach, we utilize the scores obtained from detected forgery timestamps as the classification scores for the respective videos. As observed, frame-based algorithms such as $F^3$-Net [44] exhibit significant performance degradation in classifying partially manipulated videos due to their lack of consideration of temporal factors, leading to substantial discrepancies in discriminating between different frames. On the other hand, video-level algorithms such as EfficientViT [12] demonstrated relatively effective recognition of deepfake videos with partially manipulated segments, but are unable to provide corresponding timestamps for the forgeries. In contrast, our method achieved the best classification performance while also providing corresponding forgery timestamps. Additionally, our temporal forgery localization performance significantly outperforms MDS [11] and BA-TFD [6]. It is worth noting that our model was not specifically designed for the classification task, and further performance improvement could be achieved by introducing a dedicated classification head.

**Table 8: Deepfake detection results on the Lav-DF dataset. Bold faces correspond to the top performance**

| Methods | AUC |
|---|---|
| $F^3$-Net [44] | 52.0 |
| MDS [11] | 82.8 |
| EfficientViT [12] | 96.5 |
| BA-TFD [6] | 99.0 |
| Ours | **99.8** |