# Contextual Cross-Modal Attention for Audio-Visual Deepfake Detection and Localization

Vinaya Sree Katamneni and Ajita Rattani
University of North Texas at Denton
Denton, Texas, USA
vinayasreekatamneni@my.unt.edu; ajita.rattani@unt.edu

## Abstract

*In the digital age, the emergence of deepfakes and synthetic media presents a significant threat to societal and political integrity. Deepfakes based on multi-modal manipulation, such as audio-visual, are more realistic and pose a greater threat. Current multi-modal deepfake detectors are often based on the attention-based fusion of heterogeneous data streams from multiple modalities. However, the heterogeneous nature of the data (such as audio and visual signals) creates a distributional modality gap and poses a significant challenge in effective fusion and hence multi-modal deepfake detection. In this paper, we propose a novel multi-modal attention framework based on recurrent neural networks (RNNs) that leverages contextual information for audio-visual deepfake detection. The proposed approach applies attention to multi-modal multi-sequence representations and learns the contributing features among them for deepfake detection and localization. Thorough experimental validations on audio-visual deepfake datasets, namely FakeAVCeleb, AV-Deepfake1M, TVIL, and LAV-DF datasets, demonstrate the efficacy of our approach. Cross-comparison with the published studies demonstrates superior performance of our approach with an improved accuracy and precision by $3.47\%$ and $2.05\%$ in deepfake detection and localization, respectively. Thus, obtaining state-of-the-art performance. To facilitate reproducibility, the code and the datasets information is available at* `https://github.com/vcbsl/audio-visual-deepfake/`.

***Keywords*—** Audio-visual Deepfake Detection, Contextual Cross-Attention, Deepfake Localization, Multi-modal Manipulation

## 1. Introduction

With advances in deep-generative models [41], synthetic audio and visual media have become so realistic that they are often indistinguishable from authentic content for human eyes. However, synthetic media generation techniques used by malicious users to deceive pose a serious social and political threat [25, 17, 58, 43, 33, 1, 50, 32].

In this context, visual (facial) deepfakes are generated using facial forgery techniques that depict human subjects with altered identities (i.e., face swapping), malicious actions (such as expression swapping), and facial attribute manipulation (such as skin color, gender, and age) [14, 44, 61]. Voice deepfakes, like facial deepfake technology, rely on advanced generative neural networks to synthesize audio that mimics the voice of a target speaker. Among them, Text-to-speech (TTS) voice deepfakes involve generating synthetic speech from text input that mimics a specific target speaker's voice [51]. Voice conversion-based deepfakes involve altering a person's voice to sound like another person while retaining the original content and linguistic style [54].

Audio and visual deepfakes have been employed to attack authentication systems, impersonate celebrities and politicians, and defraud finance. As a countermeasure, several unimodal audio and visual deepfake detectors have been proposed [32, 23, 2, 39, 13, 24, 46]. Lately, multi-modal deepfakes that manipulate multiple modalities, such as audio-visual, to create highly convincing and immersive fake content have shown staggering growth with advanced multimedia processing and generative AI capabilities [65]. These advanced multi-modal deepfake techniques leverage the strengths of different modalities to generate more realistic and impactful results.

Existing unimodal deepfake detectors are primarily designed to detect a single type of manipulation, such as visual, acoustic, and text. Consequently, multi-modal deepfake detectors are being investigated to detect and localize multi-modal manipulations, collectively. Within the scope of this work, several audio-visual deepfake detection and localization techniques have been proposed [48, 20, 26, 47, 27, 53, 57, 18, 36, 62, 19]. Deepfake detection aims at binary classification into pristine or deepfake. Localization aims to locate the start and end timestamps of manipulated audio-visual segments, thus facilitating a better understanding of deepfake detection results. Existing audio-visual deepfake detectors are often based on the fusion of heterogeneous streams using feature concatenation and employing attention mechanism [68, 47, 67, 37, 59, 5, 45]. Deepfake localization approaches are either anchor-based [21, 22] that utilize a sliding window approach to detect deepfake segments or boundary pre-
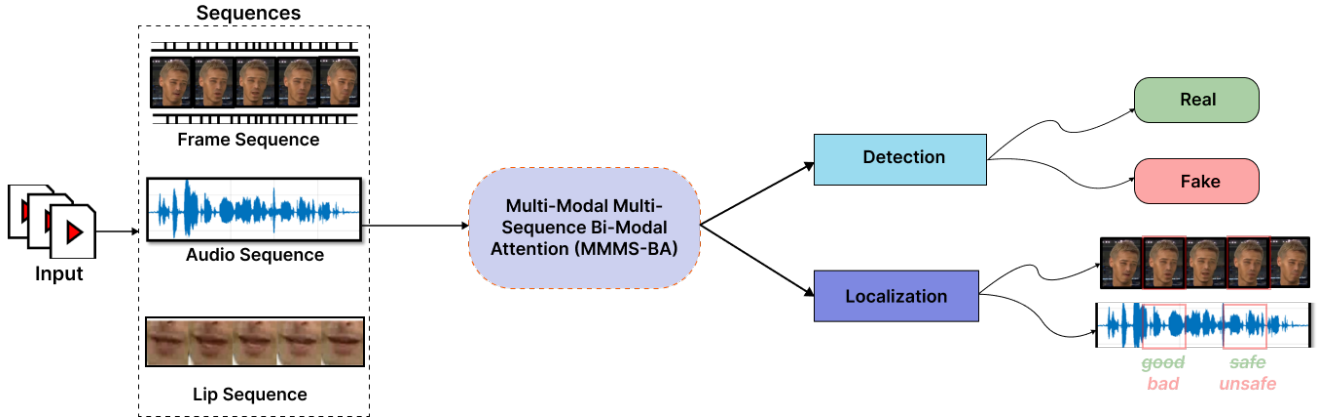
Figure 1. Overview of our proposed audio-visual deepfake detection and localization framework. The audio-visual sequences extracted from the input video are processed using our proposed MMMS-BA approach for deepfake detection and localization.

diction based [6, 10, 9, 64] that focuses on predicting the boundary of fake segments in a video using boundary matching loss.

The major **limitations** of the current attention-based fusion approaches for audio-visual deepfake detection stem from the heterogeneous nature of audio and visual signals. Consequently, the unique capabilities of each modality (modality-specific features) are not utilized effectively in the fusion process. Furthermore, noise in one modality adversely affects the overall performance of the multi-modal framework. These methods also fail to explicitly model the interactions between different modalities, such as the relationship between audio, full face, and lip movement sequences. As a result, correlation and dependencies between modalities are not fully captured, leading to suboptimal performance in deepfake detection. Moreover, current approaches do not incorporate contextual information which refers to consistent and meaningful patterns across sequences, both within and across modalities. Lastly, most current fusion-based approaches focus solely on audio-visual detection without integrated mechanisms for localization. Given that deepfakes have become more content-driven generating partially fake media [9, 8, 64], integrated deepfake detection as well as localization is crucial.

This work **proposes** a recurrent neural network-based multi-modal multi-sequence attention framework for audio-visual deepfake detection, called Multi-Modal Multi-Sequence Bi-Modal Attention (MMMS-BA). Our framework focuses on relevant features across modality pairs, leveraging attention from neighboring sequences and multi-modal representations for enhanced representation learning. Our proposed MMMS-BA performs deepfake detection as well as localization using classification and regression heads. Figure 1 illustrates our proposed MMMS-BA framework for audio-visual deepfake detection and localization.

In summary, the **contributions** of this work are as follows:

1. A novel approach for audio-visual deepfake detection and localization (MMMS-BA) based on contextual cross-attention mechanism utilizing multi-modal multi-sequence information.

2. Extensive evaluation on the publicly available audio-visual deepfake detection and localization datasets namely, AV-DeepFake1M [8], FakeAVCeleb [30], LAV-DF [10], and

TVIL [64].

3. Cross-comparison of MMMS-BA with the published work on audio-visual multi-modal deepfake detection and localization.

This paper is summarized as follows: Section 2 discusses related work on audio-visual deepfake detection and localization. The proposed approach is detailed in Section 3. Section 4 covers the datasets, evaluation metrics, and results. An ablation study analyzing varied modalities and attention mechanisms is presented in Section 5. Conclusion and future research directions are discussed in Section 6.

## 2. Related Work

In this section, we discuss the work related to audio-visual multi-modal deepfake detection and localization.

### 2.1. Audio-Visual Deepfake Detection

Audio-visual deepfake detection techniques employ audio and visual signals to detect multi-modal manipulation. A foundational work in this area, [28] assembled the FakeAVCeleb dataset consisting of audio and visual deepfakes and benchmarked various audio-visual deepfake detectors based on ensemble-based voting scheme and multi-modal convolutional neural network (CNN) based on feature concatenation.

In particular, most of the existing audio-visual deepfake detectors are based on attention mechanism-based feature concatenation [68, 67, 47, 52, 37, 59, 5, 45] to learn informative multi-modal features for deepfake detection. Studies in [38, 42] used a siamese network architecture for emotion recognition from audio and visual cues incorporating contrastive loss. Deepfake videos are detected by analyzing discrepancies in emotional cues between audio and visual modalities. Studies in [15, 55, 35, 12, 66, 31, 36] explicitly model the disagreement between the embeddings of the multiple modalities using contrastive loss for deepfake detection.

The work in [3, 4] and [11] explores the mismatch between phonemes (distinct units of sound in speech) and visemes (the visual representation of phonemes) for deepfake detection. These studies specifically focus on inconsistencies in the lip region concerning the audio which deepfake generation approaches often

struggle to replicate incorrectly, thus offering a potential for deep-fake detection [3, 4].

## 2.2. Localization of Deepfakes

Localization aims to ground the time intervals in an input where manipulation has occurred. Audio-visual deepfake localization-based approaches can be categorized into two: anchor-based [21, 22] and those based on the prediction of boundaries of fake segments in a video [6, 10, 9, 64, 8]. Anchor-based approaches make use of a sliding window over the input to ground the manipulated segments. In contrast, prediction-based approaches process the entire video and determine the fake segments using the proposed boundary-sensitive network [6].

The foundational study in [15] introduced a novel method focused on detecting and localizing discrepancies between audio and visual modalities based on temporal mismatches and unnatural movements within videos. Building on their initial work [10], the authors introduced a new dataset, Localized Audio-Visual Deepfake (LAV-DF), along with the multi-modal deepfake detection that incorporates contrastive loss function to differentiate between real and fake segments by pushing apart the feature representations of genuine and forged segments in the feature space. In addition, the loss of boundary matching precisely predicts the boundaries (start and end points) of the forged segments.

Work in [9] advances the method in [10] further by introducing a framework (named BA-TFD+) that integrates a multi-scale vision transformer with a 3D CNN architecture and trained with contrastive, frame classification, boundary matching and multi-modal boundary matching loss functions for improved precision and reliability of deepfake detection and localization.

Another work [64] in this direction introduces a framework named UMMAFormer, a universal, transformer-based framework adapted to audio and visual multi-modalities to detect temporal forgeries using anomaly detection based on self-attention mechanism.

# 3. Methodology

## 3.1. Problem Formulation

Our proposed framework aims to detect and localize deepfakes by focusing on the relationship between audio, visual and lip movement sequences by harnessing the information embedded within each modality and across their intersection.

For a given video $D$ that contains $N$ sequences, where each sequence $d_i$ is composed of 3 different modalities (full visual face, lip region, and audio), the formulation can be represented as follows:

$$D = \{d_i\}_{i=1}^{N} \quad ; \quad d_i = \{x_i^v, x_i^l, x_i^a\} \qquad (1)$$

Here, each sequence $d_i$ consists of three modalities represented by $x_i^m$ for $m \in \{v, l, a\}$ for the full visual face, lip region, and audio modality.

The classification head detects the modified sequences and the regression head is employed to localize fake segments. For the video $D$ with $N$ sequences, the entire video $D$ is classified as fake if at least one sequence is classified as fake.

The localization of the segments is represented as $Y = \{y_1, y_2, ..., y_{N_f}\}$ where $N_f$ are the number of fake segments. For a fake sequence $i$, $y_i$ represents the $i^{th}$ sequence output. Each instance $y_i = (s_i, e_i)$ is defined by its starting time $s_i$ and ending time $e_i$ of the sequence. Figure 2 illustrates the step-by-step process involved in our MMMS-BA which includes feature extraction and processing (as described below), and classification and regression heads for deepfake detection and localization. The following sub-sections provide details on each of these steps.

## 3.2. Multi-modal Multi-sequence - Bi-modal Attention (MMMS-BA) Framework

Sequences in a video represent the time series information and the classification of a sequence would have a relation with the other sequences. To model the relationship with the neighboring sequences and multiple modalities, we propose a recurrent neural network-based multi-modal attention framework named Multi-Modal Multi-Sequence Bi-modal Attention (MMMS-BA). Figure 2 shows the steps involved in applying the attention mechanism to the input sequences.

The MMMS-BA framework processes and analyzes multi-modal data across sequences in an input video. Assuming a particular video has $N$ sequences, the raw sequence levels are represented as $x_i^v$ for the full visual face, $x_i^a$ for the audio and $x_i^l$ for the lip sequence from equation 1. Three separate Bi-GRU layers with forward and backward state concatenation are first applied to the full visual face $x_i^v$, audio $x_i^a$, and lip sequence $x_i^l$ representations followed by the fully connected dense layers, resulting in $L$ (lip region), $V$ (full visual face), and $A$ (audio) embeddings. Finally, pairwise attentions are computed on various bi-modal combinations of three modalities - $(V, L)$, $(L, A)$ & $(A, V)$ as explained in the section 3.2.1.

### 3.2.1 Bi-modal Attention

Modality representations of $V$ & $L$ are obtained from the Bi-GRU network and hence contain the contextual information of the sequences for each modality. Figure 1 in the Appendix provides additional details on the computation of bi-modal attention for each modality pair. At first, we compute a pair of matching matrices $M1$ and $M2$ over two representations that account for the cross-modality information.

$$M_1 = V \cdot L^T \quad \text{and} \quad M_2 = L \cdot V^T \qquad (2)$$

### 3.2.2 Multi-sequence Attention

As mentioned earlier, we aim to leverage the contextual information of each sequence for the prediction. The probability distribution scores ($K_1$ and $K_2$) are computed over each sequence of bimodal attention matrices $M_1$ and $M_2$ (refer to equation 2) using a softmax function. This essentially computes the attention weights for the contextual sequences. Finally, soft attention is applied over the multi-modal multi-sequence attention matrices to compute the modality-wise attentive representations, i.e., $O_1$ and $O_2$ explained below.

$$K_1(i,j) = \frac{e^{M1(i,j)}}{\sum_{k=1}^{N} e^{M1(i,k)}} \quad \text{for } i, j = 1, 2, .., N \qquad (3)$$
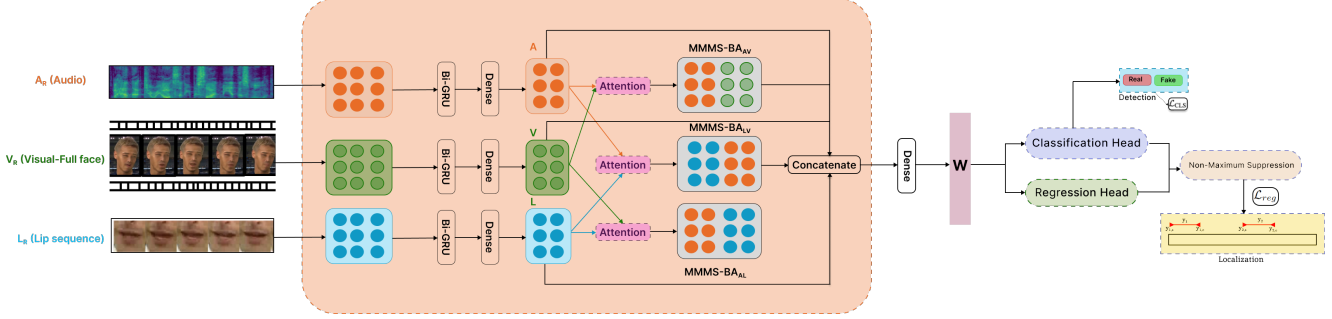
Figure 2. Illustration of the proposed Multi-Modal Multi-Sequence Bi-modal Attention (MMMS-BA) model for audio-visual deepfake detection and localization.

$$K_2(i,j) = \frac{e^{M2(i,j)}}{\sum_{k=1}^{N} e^{M2(i,k)}} \quad \text{for } i,j = 1,2,..,N \quad (4)$$

$$O_1 = K_1 \cdot L \quad \text{and} \quad O_2 = K_2 \cdot V \quad (5)$$

where $M_1$ and $M_2$ are the bi-modal attention matrices between $V$ and $L$. $i$ and $j$ represent sequences in the input.

### 3.2.3 Multiplicative Gating and Concatenation

A multiplicative gating function is computed between the multi-modal sequence-specific representations of each modality $O_1$ & $O_2$ (refer equation 5) and the corresponding bi-modal pair. This element-wise matrix multiplication assists in attending to the important components of multiple modalities and sequences. Attention matrices $A_1$ & $A_2$ are then concatenated to obtain the pairwise attention between $V$ and $L$.

$$A_1 = O_1 \odot V \quad \text{and} \quad A_2 = O_2 \odot L \quad (6)$$

$$MMMS\text{-}BA_{VL} = \text{concat}[A_1, A_2] \quad (7)$$

Similar to $MMMS\text{-}BA_{VL}$ in equation 7, we follow the same procedure to compute $MMMS\text{-}BA_{AV}$, and $MMMS\text{-}BA_{AL}$. Finally, motivated by the residual skip connection network, the pairwise attention representations $MMMS\text{-}BA_{VL}$, $MMMS\text{-}BA_{AV}$, and $MMMS\text{-}BA_{AL}$ are concatenated with the individual modalities ($V$, $A$, and $L$) to increase gradient flow to the lower layers. This concatenated feature vector $W$ (see Figure 2) is then used for deepfake detection.

### 3.3. Localization

For the sequences classified as fake, their corresponding timestamp in the input video are localized as the fake segments. The localization of the segments is represented as $Y = \{y_1, y_2, ..., y_{N_f}\}$ where $N_f$ are the number of fake segments. For a fake sequence $i$, $y_i$ represents the $i^{th}$ sequence output. Each instance $y_i = (s_i, e_i)$ is defined by its starting time $s_i$ and ending time $e_i$ of the sequence. These segments are further processed using Non-Maximum Suppression (NMS) [7] to remove highly overlapping instances, leading to the final localization timestamps.

### 3.4. Loss Function

The model's overall learning process involves minimizing the combined loss as follows:

$$\mathcal{L} = \sum_N \left( \mathcal{L}_{cls} + \lambda_{reg} \mathbb{1}_{c_i} \mathcal{L}_{reg} \right) / N_f, \quad (8)$$

where $N_f$ is the total number of fake sequences. $\mathbb{1}_{c_i}$ is an indicator function that denotes if a sequence $i$ is fake with value equal to 1 otherwise 0. $\mathcal{L}$ is applied and averaged in all sequences during training. $\lambda_{reg}$ is a coefficient that balances classification and regression loss. We set $\lambda_{reg} = 1$ by default.

Notably $\mathcal{L}_{cls}$ uses focal loss [34] to classify sequences as real or fake. $\mathcal{L}_{reg}$ adopts a differentiable IoU loss from [49]. $\mathcal{L}_{reg}$ is only enabled when the current sequence is fake.

## 4. Experimental Validations

### 4.1. Datasets

Evaluation of the proposed method is conducted on publicly available audio-visual AV-DeepFake1M [8], FakeAVCeleb [30], LAV-DF [10], and TVIL [64] deepfake datasets. Table 1 provides details on the datasets used in this study. More details about the datasets are available in the Appendix, section 1.

### 4.2. Implementation Details

For the deepfake detection task using MMMS-BA, our model architecture includes bidirectional GRUs with 300 neurons each, followed by a dense layer comprising 100 neurons. This dense layer ensures that the input features from all three modalities, i.e., full visual face, lip sequences, and audios are projected to the same dimensions, facilitating cohesive integration of information. Dropout regularization is applied with a rate of 0.3 across all layers, including the Bi-GRU layers. ReLU activation functions are utilized in the dense layers, while softmax activation is employed in the final classification layer, and ReLU (to ensure nonnegative values for start and end timestamps) is employed with Differential Intersection over Union (DIoU) loss, which measures the accuracy of the predicted timestamps against the ground truth.

The hyperparameters, including activation functions and dropout probabilities, were selected through a grid-search process using the validation set. The Adam optimizer with an exponential decay learning rate scheduler is used for optimization, and an

Table 1. Details on the datasets used in this study.

| Dataset | Year | Tasks | Manipulated Modality | Manipulation Method | #Subjects |
|---|---|---|---|---|---|
| FakeAVCeleb [29] | 2021 | Detection | Audio and Visual | Re-enactment | 500 |
| LAV-DF [10] | 2022 | Detection and Localization | Audio and Visual | Content-driven | 153 |
| TVIL [64] | 2023 | Detection and Localization | Audio and Visual | Inpainting forgery | N/A |
| AV-Deepfake1M [8] | 2023 | Detection and Localization | Audio and Visual | Content-driven | 2068 |

early stopping strategy is implemented to determine the optimal training for each model. During training, we use a batch size of 32. We conducted experiments using different combinations of bi-modal inputs, encompassing full visual face, lip sequences, and audio modalities.

**Evaluation Metrics**: For performance evaluation, we used standard evaluation metrics commonly used for deepfake detection, such as area under the ROC Curve (AUC), partial AUC (pAUC) (at 10% False Positive Rate (FPR)), and Equal Error Rate (EER) similar to the studies on deepfake detection [40, 29]. For localization, we used average precision (AP) and average recall (AR) similar to the published work in [9, 64, 8] at different thresholds to evaluate the model.

### 4.3. Performance of Deepfake Detection

Tables 2 and 3 summarize the performance of our MMMS-BA model on intra- and cross-dataset evaluation when trained on FakeAVCeleb and AV-Deepfake1M datasets. For the AV-Deepfake1M dataset, training and validation sets are available, and the testing set has not yet been released. In this regard, we have used the training set for both training (85 %) and validation (15%). While the actual validation set is used for testing the trained models.

Table 2. Intra- and cross-dataset evaluation of MMMS-BA when trained on FakeAVCeleb dataset.

| Testing Dataset | AUC | pAUC | EER | ACC | TPR | FPR |
|---|---|---|---|---|---|---|
| FakeAVCeleb | 0.989 | 0.977 | 0.029 | 0.979 | 0.965 | 0.039 |
| AV-Deepfake1M | 0.909 | 0.884 | 0.178 | 0.821 | 0.893 | 0.249 |
| LAV-DF | 0.958 | 0.938 | 0.078 | 0.932 | 0.912 | 0.081 |
| TVIL | 0.942 | 0.920 | 0.126 | 0.947 | 0.927 | 0.093 |

Table 3. Intra- and cross-dataset evaluation of MMMS-BA when trained on AV-Deepfake1M dataset.

| Testing Dataset | AUC | pAUC | EER | ACC | TPR | FPR |
|---|---|---|---|---|---|---|
| AV-Deepfake1M | 0.979 | 0.962 | 0.051 | 0.968 | 0.954 | 0.062 |
| FakeAVCeleb | 0.955 | 0.935 | 0.088 | 0.947 | 0.938 | 0.108 |
| LAV-DF | 0.968 | 0.952 | 0.068 | 0.956 | 0.941 | 0.074 |
| TVIL | 0.932 | 0.912 | 0.118 | 0.913 | 0.890 | 0.145 |

The results presented in Table 2 and Table 3 demonstrate the performance of our MMMS-BA model on intra- and cross datasets when trained on the FakeAVCeleb and AV-Deepfake1M datasets. The model trained and tested on the FakeAVCeleb obtained an AUC and ACC of 0.989 and 0.979 respectively. When trained and tested on the AV-Deepfake1M dataset, the model has an AUC and ACC of 0.979 and 0.968, respectively. This indicates the model's ability to effectively learn the characteristics of audio-visual manipulations present in the FakeAVCeleb and AV-Deepfake1M datasets.

Table 4. Comparison of MMMS-BA with published audio-visual deepfake detection approaches trained and tested on the FakeAVCeleb dataset.

| Model | AUC | ACC |
|---|---|---|
| MIS-AVoiDD [53] | 0.973 | 0.962 |
| Unsupervised Multi-modal Deepfake Detection [57] | 0.968 | - |
| Audio-visual person-of-interest [18] | 0.946 | 0.850 |
| Audio-visual anomaly detection [20] | 0.942 | 0.945 |
| Multi-modal contrastive learning [36] | 0.978 | 0.965 |
| PVAS-MDD [62] | 0.965 | 0.948 |
| Novel Smart Deepfake Detection [19] | 0.954 | 0.960 |
| Multimodaltrace [48] | 0.929 | - |
| Hearing and seeing abnormality [56] | 0.944 | - |
| NPVforensics [11] | 0.925 | - |
| **Our-MMMS-BA** (Visual, Audio, and Lip Sequence) | **0.989** | **0.979** |

Table 2 shows that the cross-dataset evaluations have a performance landscape. When trained on FakeAVCeleb and tested on AV-Deepfake1M, LAV-DF, and TVIL. On average, there is a decrease of only 0.053 and 0.076 AUC and ACC, respectively, over intra-dataset evaluation. For the AV-Deepfake1M dataset, the model obtained the lowest performance with an AUC and ACC of 0.909 and 0.893, respectively. This is expected due to the differences in manipulation techniques across datasets. The model maintains relatively high performance when tested on LAV-DF and TVIL. This could be because of the similar manipulation technique in the LAV-DF dataset. While the TVIL dataset has inpainting-based manipulation providing visual cues for detection.

When trained on AV-Deepfake1M and tested on FakeAVCeleb, LAV-DF, and TVIL datasets, a decrease of only 0.027 and 0.029 AUC and ACC, respectively, is recorded over the intra-dataset evaluation. The model obtained the lowest performance with an AUC of 0.932 and ACC of 0.913 on the TVIL dataset. This is expected due to the differences in manipulation techniques across datasets. The model obtained the best performance on the LAV-DF dataset with AUC and ACC of 0.968 and 0.956, respectively. This is because both AV-Deepfake1M and LAV-DF use content-driven manipulation techniques (see Table 1).

In **summary**, the experimental results for deepfake detection suggest a high generalization and efficacy of the MMMS-BA model. On the intra-dataset evaluation of the model on FakeAVCeleb and AV-Deepfake1M datasets, an average AUC and ACC of 0.984 and 0.973, respectively, is obtained. The high performance of the model on the FakeAVCeleb dataset also confirms that our model obtains efficient detection performance even when one of the modalities (audio or visual) is manipulated, as available in the FakeAVCeleb dataset. On cross-dataset evaluation of the model, an average AUC and ACC of 0.944 and 0.919, respec-

tively, are obtained, suggesting the high performance and generalizability of our proposed MMMS-BA model.

## 4.4. Comparison with the Published Audio-Visual Deepfake Detectors

The comparison shown in Table 4 highlights the competitive landscape of deepfake detection models trained and evaluated on the FakeAVCeleb dataset. As can be seen in Table 4, the MMMS-BA model, which uses the full facial image, audio, and lip modality, outperforms all existing audio-visual deepfake detectors by obtaining an AUC of 0.989 and an ACC of 0.979. In particular, our MMMS-BA model outperforms NPVForensics [11] which also utilizes face, audio, and lip movement data. NPVForensics is based on mining the correlation between non-critical phonemes and visemes using a Swin Transformer and cross-modal fusion, but it achieves lower performance over MMMS-BA due to its limited ability to capture temporal dependencies and interactions across multiple sequences. The Unsupervised Cross-Modal Inconsistencies [57] model obtained the second-best performance, with an AUC of 0.968. This also suggests the potential of leveraging motion inconsistencies between modalities for deepfake detection.

In **summary**, MMMS-BA outperforms all recently published work based on integrating audio, visual, and lip-movement data with an average performance increment of 0.0341 and 3.47% in AUC and ACC, respectively. Thus, obtaining state-of-the-art performance.

## 4.5. Deepfake Localization Performance

Table 5 shows the performance evaluation of the encoder-decoder (Enc-Dec), ActionFormer [63], BA-TFD+ [9], UMMAFormer [64] and MMMS-BA based deepfake localization approaches when trained on AV-Deepfake1M and tested on LAV-DF, TVIL, and AV-Deepfake1M datasets. The Enc-Dec approach represents a sequence-to-sequence encoder-decoder-based architecture employed considering audio and visual modalities. For Enc-Dec, audio and visual encoders are used, and the hidden representations from these encoders are concatenated together, followed by two dense layers, and the output classification layer provides sequences. The start and end timestamps of the fake sequences are the fake segments in a video. The results for methods ActionFormer, BA-TFD+, and UMMAFormer reported in Table 5 are taken from the original papers [9, 64, 8]. Since the evaluation is done on the same datasets used in our study, it is a fair comparison.

The Enc-Dec model obtained the lowest performance across all datasets and metrics. This is expected as it is based on simple feature concatenation from audio-visual streams and does not include any advanced processing techniques like attention mechanism. On AV-Deepfake1M, the MMMS-BA model has obtained the best performance with AP@0.5 and AR@50 of 62.75 and 57.49, respectively. UMMAFormer is the second best model obtaining AP@0.5 and AR@0.5 of 51.64 and 48.86, respectively. MMMS-BA has obtained better performance on the AV-Deepfake1M dataset when compared to the best model in [8]. However, the lower performance observed on the AV-Deepfake1M dataset, although it was used as a training set, underscores the highly realistic fake content generated in a content-driven manner, altering the real transcripts with replace, delete, and insert operations and the corresponding audio-visual modalities accordingly.

When evaluated on the LAV-DF dataset, MMMS-BA obtained the best performance with an AP@0.5 and AR@50 of 97.56 and 93.45, respectively. BA-TFD+ has obtained the second-best performance with an AP@0.5 and AR@50 of 96.30 and 80.48, respectively. Although BA-TFD+ showed a smaller performance difference compared to MMMS-BA at lower threshold levels, it exhibited a larger performance difference at higher thresholds, indicating a less precise localization.

On the TVIL dataset, the MMMS-BA model demonstrated the best performance with an AP@0.5 of 96.87 and an AR@50 of 90.47. The UMMAFormer model followed closely with an AP@0.5 of 88.68 and an AR@50 of 90.43, and it achieved the highest scores for AP@0.75 (84.70) and AP@0.95 (62.43). The Enc-Dec model had the lowest performance with an AP@0.5 of 23.08 and an AR@50 of 43.98 due to its ineffective use of the available modality information.

In **summary**, the MMMS-BA model consistently demonstrated the best localization capability across all datasets. This is due to the utilization of contextual information across the sequences over existing methods, resulting in better robustness in deepfake localization.

# 5. Ablation Study

An ablation study is conducted by varying the attention mechanism in the contextual cross-modal attention block in Figure 1 resulting in two other variations, i.e., Multi-Modal Uni-Sequences - Self-Attention (MMUS-SA) and Multi-Sequence - Self-Attention (MS-SA). We also ablated between the bi-modal combination of the modalities considered, that is, the full face, voice, and lip region. The details given below.

## 5.1. Varying the Attention at Sequence and Modality Level

**Multi-Modal Uni-Sequences - Self Attention (MMUS-SA) Framework**: MMUS-SA framework does not account for information from other sequences at the attention level but utilizes multi-modal information of a single sequence for prediction. For more details on the calculation of attention in the MMUS-SA framework refer to the Appendix, section 2.

**Multi-Sequence - Self Attention (MS-SA) Framework**: In the MS-SA framework, we apply self-attention to the sequences of each modality separately and use these for classification. For more details on the calculation of attention in the MMUS-SA framework refer to the Appendix, section 2.

Table 6 shows the performance of the MMUS-SA and MS-SA models over MMMS-BA when trained and tested on FakeAVCeleb. From the results obtained, it is evident that MMMS-BA has obtained the best performance with the lowest EER of 0.029. The MMUS-SA model has a closer performance with the MMMS-BA variation but has a higher EER. This suggests that incorporating self-attention into a single sequence in multiple modalities achieved a closer AUC, but was unable to obtain a low EER. However, MS-SA has obtained the lowest performance when compared to the other two variations. The reason is that MS-SA does not incorporate the interaction between the modalities for multi-modal manipulation detection.

Table 5. Comparison between MMMS-BA and other methods on LAV-DF, TVIL, and AV-Deepfake1M datasets. Bold faces represent the best performance. (AP- Average Precision and AR- Average Recall).

| Testing Dataset | Method | AP@0.5 | AP@0.75 | AP@0.95 | AR@10 | AR@20 | AR@50 | AR@100 |
|---|---|---|---|---|---|---|---|---|
| AV-Deepfake1M | Enc-Dec | 06.23 | 0.08 | 0.00 | 11.45 | 15.79 | 23.75 | 31.71 |
| | ActionFormer [63] | 36.08 | 12.01 | 00.16 | 26.60 | 27.00 | 27.11 | - |
| | BA-TFD+ [9] | 44.42 | 13.64 | 00.03 | 34.67 | 40.37 | 48.86 | - |
| | UMMAFormer [64] | 51.64 | 28.07 | 01.58 | 42.09 | 43.45 | 48.86 | - |
| | **Our-MMMS-BA** | **62.75** | **35.87** | **18.37** | **54.28** | **55.94** | **57.49** | **59.66** |
| LAV-DF | Enc-Dec | 12.96 | 0.97 | 0.00 | 18.79 | 21.67 | 30.54 | 38.12 |
| | ActionFormer [63] | 85.23 | 59.05 | 00.93 | 76.93 | 77.19 | 77.23 | 77.23 |
| | BA-TFD+ [9] | 96.30 | 84.96 | 04.44 | 78.75 | 79.40 | 80.48 | 81.62 |
| | **Our-MMMS-BA** | **97.56** | **95.25** | **39.02** | **89.42** | **95.93** | **93.45** | **94.05** |
| TVIL | Enc-Dec | 23.08 | 08.45 | 05.32 | 17.16 | 23.47 | 43.98 | 45.18 |
| | ActionFormer [63] | 86.27 | 83.03 | 28.17 | 84.82 | 85.77 | 88.10 | 88.49 |
| | BA-TFD+ [9] | 76.90 | 38.50 | 0.25 | 66.90 | 64.08 | 60.77 | 58.42 |
| | UMMAFormer [64] | 88.68 | **84.70** | **62.43** | 87.09 | **88.21** | 90.43 | 91.16 |
| | **Our-MMMS-BA** | **96.87** | 81.33 | 28.43 | **88.61** | 87.83 | **90.47** | **92.91** |

Table 6. Performance comparison of MMMS-BA with MMUS-BA and MS-SA models with varied attention mechanism, trained and tested on FakeAVCeleb dataset.

| Model | AUC | pAUC | EER | ACC | TPR | FPR |
|---|---|---|---|---|---|---|
| MMUS-SA | 0.989 | 0.978 | 0.033 | 0.979 | 0.963 | 0.039 |
| MS-SA | 0.977 | 0.956 | 0.045 | 0.965 | 0.943 | 0.051 |
| MMMS-BA | **0.989** | **0.977** | **0.029** | **0.979** | **0.965** | **0.039** |

Table 7. Performance comparison of MMMS-BA with different bi-modal combinations (V-full visual face, L-lip sequence, and A-audio) when trained and tested on the FakeAVCeleb dataset.

| Model and Modalities | AUC | pAUC | EER | ACC | TPR | FPR |
|---|---|---|---|---|---|---|
| MMMS-BA (V+A) | 0.955 | 0.941 | 0.074 | 0.939 | 0.922 | 0.095 |
| MMMS-BA (V+L) | 0.814 | 0.798 | 0.263 | 0.736 | 0.818 | 0.345 |
| MMMS-BA (L+A) | 0.924 | 0.907 | 0.175 | 0.823 | 0.931 | 0.282 |
| MMMS-BA (V+L+A) | **0.989** | **0.978** | **0.033** | **0.979** | **0.963** | **0.039** |

In **summary**, MS-SA and MMUS-SA attention methods typically identify important features within a modality over multiple sequences and a single sequence across different modalities, respectively. However, the proposed MMMS-BA model is specifically designed to understand the intricate relationship between different modalities over multiple sequences, leading to a more robust representation of the data and enhanced performance. Consequently, obtaining better performance over MMUS-SA and MS-SA. This confirms the importance of using interactions between both multi-sequences and multi-modalities for deepfake detection.

## 5.2. Performance of MMMS-BA on Different Combination of Modalities

In this study, we ablated between the combination of modalities (V-full visual face, L-lip sequence, and A-audio) and evaluated the performance of MMMS-BA when trained and tested on FakeAVCeleb. Table 7 illustrates the role of different modalities in enhancing detection accuracy. As can be seen, the combination of visual and lip (V + L) obtained the lowest performance with an AUC of 0.814. This result can be attributed to the absence of audio data, which limits the model to visual information alone. The combination of lip and audio (L + A) obtained a performance improvement over (V + L) with an increase in AUC of 0.11. This enhancement suggests that incorporating audio data alongside lip sequences provides complementary information, enhancing the model's performance. The combination of visual and audio (V+A) obtained further improvement over (V+L and L+A) with an AUC of 0.955. This integration demonstrates that including audio features with full-face visual data offers additional insights, thereby boosting detection accuracy beyond the (V + L) and (L + A) combinations. The best performance is obtained by combining (V+L+A) with an average increment in AUC of 0.091.

In **summary**, this analysis suggests that incorporating lip sequence data along with (V+A) contributes significantly to the effectiveness of the model. Full-face images encompass a complex mix of facial features, making it challenging for the model to isolate and differentiate inconsistencies in the lip region effectively. Thus, the additional lip modality enriched the information available to the model, ultimately improving its discriminatory power and performance, also supported by the NPVForensics model [11] for audio-visual deepfake detection.

## 6. Conclusion and Future Work

With the rapid evolution of deepfake techniques combining synthesized audio with forged videos, there is an urgent need for robust audio-visual deepfake detection and localization techniques. Current audio-visual deepfake detection approaches relying on fusing audio and visual streams employing basic attention mechanisms, overlook intricate inter-modal relationships crucial for accurate detection. To address this challenge, we introduced the MMMS-BA framework. This framework effectively captures intra- and inter-modal correlations by applying attention to multi-modal multi-sequence representations and learns the contributing features among them for effective deepfake detection and localization Our experimental findings demonstrate that MMMS-BA outperforms existing audio-visual deepfake detectors, achieving SOTA performance in detecting and localizing deepfake segments within videos.

Given the proliferation of AI-generated content using sophisticated generative models, tackling emerging forms of content manipulation remains a critical challenge. Therefore, as part of future research directions, we will extend our framework to incorporate text analysis along with audio and visual modalities to ensure ro-

bust protection against misleading multimedia content. In addition, our model will be adapted to account for missing modalities during the training and inference stage.

# References

[1] D. Afchar, V. Nozick, J. Yamagishi, and I. Echizen. Mesonet: a compact facial video forgery detection network. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7, 2018.

[2] S. Agarwal, H. Farid, T. El-Gaaly, and S.-N. Lim. Detecting deep-fake videos from appearance and behavior. In *2020 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6, 2020.

[3] S. Agarwal, H. Farid, O. Fried, and M. Agrawala. Detecting deep-fake videos from phoneme-viseme mismatches. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 660–661, 2020.

[4] S. Agarwal, L. Hu, E. Ng, T. Darrell, H. Li, and A. Rohrbach. Watch those words: Video falsification detection using word-conditioned facial motion. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 4710–4719, 2023.

[5] S. Asha, P. Vinod, and V. G. Menon. A defensive attention mechanism to detect deepfake content across multiple modalities. *Multimedia Systems*, 30(1):56, 2024.

[6] A. Bagchi, J. Mahmood, D. Fernandes, and R. K. Sarvadevabhatla. Hear me out: Fusional approaches for audio augmented temporal action localization. *arXiv preprint arXiv:2106.14118*, 2021.

[7] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis. Soft-nms–improving object detection with one line of code. In *Proceedings of the IEEE international conference on computer vision*, pages 5561–5569, 2017.

[8] Z. Cai, S. Ghosh, A. P. Adatia, M. Hayat, A. Dhall, and K. Stefanov. Av-deepfake1m: A large-scale llm-driven audio-visual deepfake dataset. *arXiv preprint arXiv:2311.15308*, 2023.

[9] Z. Cai, S. Ghosh, T. Gedeon, A. Dhall, K. Stefanov, and M. Hayat. ” glitch in the matrix!”: A large scale benchmark for content driven audio-visual forgery detection and localization. *arXiv preprint arXiv:2305.01979*, 2023.

[10] Z. Cai, K. Stefanov, A. Dhall, and M. Hayat. 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)*, pages 1–10. IEEE, 2022.

[11] Y. Chen, Y. Yu, R. Ni, Y. Zhao, and H. Li. Npvforensics: Jointing non-critical phonemes and visemes for deepfake detection. *arXiv preprint arXiv:2306.06885*, 2023.

[12] H. Cheng, Y. Guo, T. Wang, Q. Li, X. Chang, and L. Nie. Voice-face homogeneity tells deepfake. *ACM Transactions on Multimedia Computing, Communications and Applications*, 20(3):1–22, 2023.

[13] A. Chintha, B. Thai, S. J. Sohrawardi, K. Bhatt, A. Hickerson, M. Wright, and R. Ptucha. Recurrent convolutional structures for audio spoof and video deepfake detection. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):1024–1037, 2020.

[14] Y. Choi, M. Choi, M. Kim, J.-W. Ha, S. Kim, and J. Choo. Stargan: Unified generative adversarial networks for multi-domain image-to-image translation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[15] K. Chugh, P. Gupta, A. Dhall, and R. Subramanian. Not made for each other- audio-visual dissonance-based deepfake detection and localization. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM ’20, page 439–447, New York, NY, USA, 2020. Association for Computing Machinery.

[16] J. S. Chung, A. Nagrani, and A. Zisserman. VoxCeleb2: Deep Speaker Recognition. In *Proc. Interspeech 2018*, pages 1086–1090, 2018.

[17] D. Citron. How deepfakes undermine truth and threaten democracy.

[18] D. Cozzolino, A. Pianese, M. Nießner, and L. Verdoliva. Audio-visual person-of-interest deepfake detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 943–952, 2023.

[19] M. Elpeltagy, A. Ismail, M. S. Zaki, and K. Eldahshan. A novel smart deepfake video detection system. *International Journal of Advanced Computer Science and Applications*, 14(1), 2023.

[20] C. Feng, Z. Chen, and A. Owens. Self-supervised video forensics by audio-visual anomaly detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10491–10503, 2023.

[21] J. Gao, K. Chen, and R. Nevatia. Ctap: Complementary temporal action proposal generation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 68–83, 2018.

[22] J. Gao, Z. Yang, K. Chen, C. Sun, and R. Nevatia. Turn tap: Temporal unit regression network for temporal action proposals. In *Proceedings of the IEEE international conference on computer vision*, pages 3628–3636, 2017.

[23] A. Haliassos, K. Vougioukas, S. Petridis, and M. Pantic. Lips don’t lie: A generalisable and robust approach to face forgery detection. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5037–5047, 2021.

[24] A. Hamza, A. R. R. Javed, F. Iqbal, N. Kryvinska, A. S. Almadhor, Z. Jalil, and R. Borghol. Deepfake audio detection via mfcc features using machine learning. *IEEE Access*, 10:134018–134028, 2022.

[25] T. Hwang. Deepfakes: A grounded threat assessment. Technical report, Georgetown University, July 2020.

[26] H. Ilyas, A. Javed, and K. M. Malik. Avfakenet: A unified end-to-end dense swin transformer deep learning model for audio–visual deepfakes detection. *Applied Soft Computing*, 136:110124, 2023.

[27] V. S. Katamneni, A. V. Nadimpalli, and A. Rattani. Demographic fairness and accountability of audio and video-based unimodal and bi-modal deepfake detectors. In T. Bourlai, editor, *Face Recognition Across the Imaging Spectrum (FRAIS)*. Springer, 2023.

[28] H. Khalid, M. Kim, S. Tariq, and S. S. Woo. Evaluation of an audio-video multimodal deepfake dataset using unimodal and multimodal detectors. In *Proceedings of the 1st Workshop on Synthetic Multimedia - Audiovisual Deepfake Generation and Detection*, ADGD '21, page 7–15, New York, NY, USA, 2021. Association for Computing Machinery.

[29] H. Khalid, S. Tariq, M. Kim, and S. S. Woo. Fakeavceleb: A novel audio-video multimodal deepfake dataset. *arXiv preprint arXiv:2108.05080*, 2021.

[30] H. Khalid, S. Tariq, and S. S. Woo. Fakeavceleb: A novel audio-video multimodal deepfake dataset. *CoRR*, abs/2108.05080, 2021.

[31] J. K. Lewis, I. E. Toubal, H. Chen, V. Sandesera, M. Lomnitz, Z. Hampel-Arias, C. Prasad, and K. Palaniappan. Deepfake video detection based on spatial, spectral, and temporal inconsistencies using multimodal deep learning. In *2020 IEEE Applied Imagery Pattern Recognition Workshop (AIPR)*, pages 1–9. IEEE, 2020.

[32] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo. Face x-ray for more general face forgery detection. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5000–5009, 2020.

[33] Y. Li and S. Lyu. Exposing deepfake videos by detecting face warping artifacts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.

[34] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.

[35] X. Liu, Y. Yu, X. Li, and Y. Zhao. Magnifying multimodal forgery clues for deepfake detection. *Signal Processing: Image Communication*, 118:117010, 2023.

[36] X. Liu, Y. Yu, X. Li, and Y. Zhao. Mcl: multimodal contrastive learning for deepfake detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 2023.

[37] M. Masood, A. Javed, and A. Irtaza. Attention-based multimodal learning framework for generalized audio-visual deepfake detection. *PREPRINT Research Square : rs.3.rs-3415144/v1*, 2023.

[38] T. Mittal, U. Bhattacharya, R. Chandra, A. Bera, and D. Manocha. Emotions don't lie: An audio-visual deepfake detection method using affective cues. In *Proceedings of the 28th ACM International Conference on Multimedia*, MM '20, page 2823–2832, New York, NY, USA, 2020. Association for Computing Machinery.

[39] L. Muda, M. Begam, and I. Elamvazuthi. Voice recognition algorithms using mel frequency cepstral coefficient (mfcc) and dynamic time warping (dtw) techniques. *CoRR*, abs/1003.4083, 2010.

[40] A. V. Nadimpalli and A. Rattani. Gbdf: Gender balanced deepfake dataset towards fair deepfake detection. *ArXiv*, abs/2207.10246, 2022.

[41] A. V. Nadimpalli and A. Rattani. Proactive deepfake detection using gan-based visible watermarking. *ACM Trans. Multimedia Comput. Commun. Appl.*, Sep 2023.

[42] R. Nekadi. *Siamese Network-based Multi-modal Deepfake Detection*. University of Missouri-Kansas City, 2020.

[43] T. T. Nguyen, C. M. Nguyen, D. Nguyen, D. T. Nguyen, and S. Nahavandi. Deep learning for deepfakes creation and detection. *ArXiv*, abs/1909.11573, 2019.

[44] Y. Nirkin, Y. Keller, and T. Hassner. Fsgan: Subject agnostic face swapping and reenactment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019.

[45] T. Oorloff, S. Koppisetti, N. Bonettini, D. Solanki, B. Colman, Y. Yacoob, A. Shahriyari, and G. Bharaj. Avff: Audio-visual feature fusion for video deepfake detection. *arXiv e-prints*, pages arXiv–2406, 2024.

[46] A. Pianese, D. Cozzolino, G. Poggi, and L. Verdoliva. Deepfake audio detection by speaker verification. In *2022 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–6. IEEE, 2022.

[47] Y. Qian, Z. Chen, and S. Wang. Audio-visual deep neural network for robust person verification. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:1079–1092, 2021.

[48] M. A. Raza and K. M. Malik. Multimodaltrace: Deepfake detection using audiovisual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 993–1000, 2023.

[49] H. Rezatofighi, N. Tsoi, J. Gwak, A. Sadeghian, I. Reid, and S. Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 658–666, 2019.

[50] A. Rössler, D. Cozzolino, L. Verdoliva, C. Riess, J. Thies, and M. Niessner. Faceforensics++: Learning to detect manipulated facial images. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1–11, 2019.

[51] D. Salvi, B. Hosler, P. Bestagini, M. C. Stamm, and S. Tubaro. Timit-tts: A text-to-speech dataset for multimodal synthetic media detection. *IEEE Access*, 2023.

[52] R. Shao, T. Wu, J. Wu, L. Nie, and Z. Liu. Detecting and grounding multi-modal media manipulation and beyond. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2024.

[53] V. Sree Katamneni and A. Rattani. Mis-avoidd: Modality invariant and specific representation for audio-visual deepfake detection. *arXiv e-prints*, pages arXiv–2310, 2023.

[54] C. Sun, S. Jia, S. Hou, and S. Lyu. Ai-synthesized voice detection using neural vocoder artifacts. *arXiv preprint arXiv:2304.13085*, 2023.

[55] C.-S. Sung, J.-C. Chen, and C.-S. Chen. Hearing and seeing abnormality: Self-supervised audio-visual mutual learning for deepfake detection. In *ICASSP 2023 - 2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5, 2023.

[56] C.-S. Sung, J.-C. Chen, and C.-S. Chen. Hearing and seeing abnormality: Self-supervised audio-visual mutual learning for deepfake detection. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE, 2023.

[57] M. Tian, M. Khayatkhoei, J. Mathai, and W. AbdAl-mageed. Unsupervised multimodal deepfake detection us-ing intra-and cross-modal inconsistencies. *arXiv preprint arXiv:2311.17088*, 2023.

[58] R. Tolosana, R. Vera-Rodríguez, J. Fierrez, A. Morales, and J. Ortega-Garcia. Deepfakes and beyond: A survey of face manipulation and fake detection. *Inf. Fusion*, 64:131–148, 2020.

[59] G. Wang, P. Zhang, L. Xie, W. Huang, Y. Zha, and Y. Zhang. An audio-visual attention based multimodal net-work for fake talking face videos detection. *arXiv preprint arXiv:2203.05178*, 2022.

[60] N. Xu, L. Yang, Y. Fan, J. Yang, D. Yue, Y. Liang, B. Price, S. Cohen, and T. Huang. Youtube-vos: Sequence-to-sequence video object segmentation. In *Proceedings of the European conference on computer vision (ECCV)*, pages 585–601, 2018.

[61] T. Xu, P. Zhang, Q. Huang, H. Zhang, Z. Gan, X. Huang, and X. He. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceed-ings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[62] Y. Yu, X. Liu, R. Ni, S. Yang, Y. Zhao, and A. C. Kot. Pvass-mdd: predictive visual-audio alignment self-supervision for multimodal deepfake detection. *IEEE Transactions on Cir-cuits and Systems for Video Technology*, 2023.

[63] C.-L. Zhang, J. Wu, and Y. Li. Actionformer: Localizing moments of actions with transformers. In *European Confer-ence on Computer Vision*, pages 492–510. Springer, 2022.

[64] R. Zhang, H. Wang, M. Du, H. Liu, Y. Zhou, and Q. Zeng. Ummaformer: A universal multimodal-adaptive transformer framework for temporal forgery localization. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 8749–8759, 2023.

[65] T. Zhang. Deepfake generation and detection, a survey. *Mul-timedia Tools and Applications*, 81(5):6259–6276, 2022.

[66] Y. Zhang, J. Zhan, W. Jiang, and Z. Fan. Deepfake detection based on incompatibility between multiple modes. In *2021 International Conference on Intelligent Technology and Em-bedded Systems (ICITES)*, pages 1–7. IEEE, 2021.

[67] L. Zhao, M. Zhang, H. Ding, and X. Cui. Fine-grained deep-fake detection based on cross-modality attention. *Neural Computing and Applications*, pages 1–14, 2023.

[68] Y. Zhou and S.-N. Lim. Joint audio-visual deepfake detec-tion. In *Proceedings of the IEEE/CVF International Confer-ence on Computer Vision*, pages 14800–14809, 2021.

# 7. Supplementary Material

## 7.1. Dataset Details

Evaluation of the proposed method is conducted on publicly available audio-visual AV-Deepfake1M [8], FakeAVCeleb [30], LAV-DF [10], and TVIL [64] deepfake datasets mentioned in Sec-tion 4.1 of the main paper. Further, details on these datasets are given below.

- **FakeAVCeleb:** The FakeAVCeleb dataset [30] is a collec-tion of videos with audio and visual manipulations of celebri-ties that have been generated using various deepfake tech-niques. The dataset is created by selecting videos from the VoxCeleb2 [16] dataset, featuring 500 celebrities. The videos in this dataset are clean, featuring only one person's frontal face without any occlusion. The dataset is well-balanced and annotated in terms of gender, race, geogra-phy, and visual and audio manipulations, making it useful for training deep learning models that can generalize well on unseen test sets. We chose this dataset for our multimodal detection experiments because it contains both audio and vi-sual manipulations, as well as a variety of deep-fake genera-tion techniques.

- **LAV-DF Dataset** The Localized Audio Visual DeepFake (LAV-DF) [10] dataset emerges as a critical asset in deepfake detection, particularly for benchmarking methods to detect and localize manipulated segments within videos. Compris-ing 136,304 videos across 153 unique identities, the dataset offers a diverse collection that includes 36,431 real videos alongside 99,873 videos embedded with fake segments. For a complete evaluation, LAV-DF is divided into three identity-independent subsets: training (78,703 videos), validation (31,501 videos), and testing (26,100 videos), each offering a distinct set of identities to ensure an unbiased assessment.

- **TVIL Dataset** In response to the increasing challenges posed by advanced AI-generated content (AIGC) technolo-gies, the TVIL dataset [64] has been synthesized as a new benchmark aimed at locating video inpainting segments. Constructed upon the foundation of YouTubeVOS 2018 [60], which aggregates over 4,000 videos from YouTube, the TVIL dataset leverages one of the most prolific platforms for video content as its base. This choice is strategic, consid-ering YouTube's prominence in content generation and the spread of misinformation. By synthesizing a dataset rooted in YouTube videos, TVIL is poised to offer a robust eval-uation framework that defends against misinformation and catalyzes new research directions in the fight against digital content forgery.

- **AV-Deepfake1M Dataset** The AV-Deepfake1M [8] dataset stands as a pioneering contribution to the field of audio-visual deepfake detection, encompassing an extensive com-pilation of 1,886 hours of audio-visual data curated using content-driven manipulation techniques. from 2,068 unique subjects set against diverse background environments. This dataset is distinguished in facilitating the development of de-tection algorithms capable of navigating the landscape of content-driven audio-visual manipulations.

## 7.2. On Varying Attention Mechanism

This section provides more details on the calculation of atten-tion in MMUS-SA and MS-SA variations discussed in Section 5.1 of the main paper.

- **Multi-Modal Uni-Sequences - Self Attention (MMUS-SA) Framework**: MMUS-SA framework does not account for information from other sequences at the attention level, rather it utilizes multi-modal information of a single se-quence for prediction. For a video $D$ having '$N$' sequences,
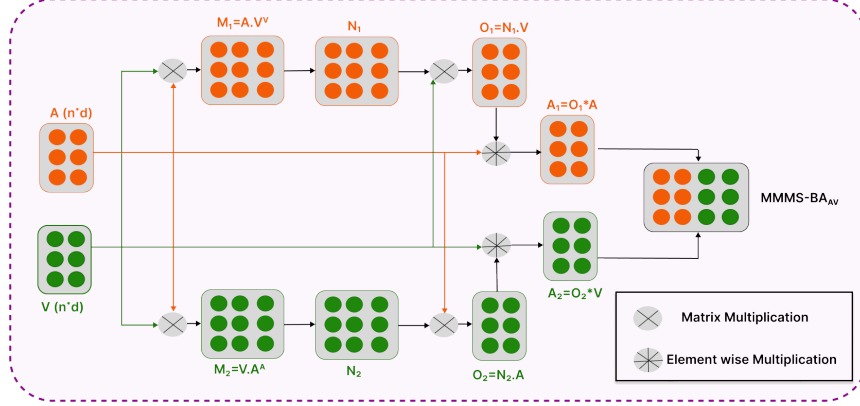
Figure 3. Multi-Modal Multi-Sequence Attention computation of Audio and Full Visual Face Modalities ($MMMS\text{-}BA_{AV}$)

'$N$' separate attention blocks are needed, where each block computes the self-attention over multi-modal information of a single sequence. Let $X_p$ be the information matrix of the $p^{th}$ sequence where the three '$r$' dimensional rows are the outputs of the dense layers for the three modalities. The attention matrix $A_p$ is computed separately for $p = 1^{st}, 2^{nd},$ ... $N^{th}$ sequences. Finally, for each sequence $p$, $A_p$ and $X_p$ are concatenated and passed to the output layer for classification.

$$M_{\mathrm{p}} = X_{\mathrm{p}} \cdot X_{\mathrm{p}}^{T} \qquad (9)$$

$$N_{\mathrm{p}}(i,j) = \frac{e^{M_{\mathrm{p}}(i,j)}}{\sum_{k=1}^{3} e^{M_{\mathrm{p}}(i,k)}} \text{ for } i,j = 1,2,3; \qquad (10)$$

$$O_{\mathrm{p}} = N_{\mathrm{p}} \cdot X_{\mathrm{p}} \qquad (11)$$

$$A_{\mathrm{p}} = O_{\mathrm{p}} \odot X_{\mathrm{p}} \qquad (12)$$

The attention matrix $A_{\mathrm{u}}$ is computed separately for $p = 1^{\mathrm{st}}, 2^{\mathrm{nd}}, \ldots, N^{\mathrm{th}}$ sequences. Finally, for each sequence $p$, $A_{\mathrm{p}}$ and $X_{\mathrm{p}}$ are concatenated and passed to the output layer for classification.

- **Multi-Sequence - Self Attention (MS-SA) Framework**: In the MS-SA framework, we apply self-attention to the sequences of each modality separately and use these for classification. In contrast to the MMSS-SA framework, MS-SA

utilizes the contextual information of the sequences at the attention level. Let $L$ (text), $V$ (visual), and $L$ (lip region) be the outputs of the dense layers. Three separate attention blocks are required for the three modalities, where each block takes multi-sequence information of a single modality and computes the self-attention matrix. Attention matrices $A_l$, $A_v$, and $A_a$ are computed for lip, visual, and acoustic, respectively. Finally, $A_v$, $A_l$, $A_a$, $V$, $L$, and $A$ are concatenated and passed to the output layer for classification.

$$M_v = V \cdot V^T \qquad (13)$$

$$N_v(i,j) = \frac{e^{M_v(i,j)}}{\sum_{k=1}^{u} e^{M_v(i,k)}} \text{ for } i,j = 1,\ldots,u \qquad (14)$$

$$O_v = N_v \cdot V \qquad (15)$$

$$A_v = O_v \odot V \qquad (16)$$

The attention matrix $A_p$ is computed for $p = 1^{\mathrm{st}}, 2^{\mathrm{nd}}, \ldots, u^{\mathrm{th}}$ sequences. Finally, for each sequence $u$, $A_p$, and $X_p$ are concatenated and passed to the output layer with softmax activation for classification.

Further, Figure 3 provides further details on the attention computation for our model MMMS-BA discussed in Section 3.2 of the main paper.