

# FAME: A Lightweight Spatio-Temporal Network for Model Attribution of Face-Swap Deepfakes

Wasim Ahmad<sup>1,2,3</sup> [was\\_last@iis.sinica.edu.tw](mailto:was_last@iis.sinica.edu.tw)\*

Yan-Tsung Peng<sup>3</sup> [ytpeng@cs.nccu.edu.tw](mailto:ytpeng@cs.nccu.edu.tw)

Yuan-Hao Chang<sup>4</sup> [johnson@csie.ntu.edu.tw](mailto:johnson@csie.ntu.edu.tw)

<sup>1</sup> Institute of Information Science, Academia Sinica, Taipei, Taiwan, 115

<sup>2</sup> Social Networks and Human-Centred Computing, Taiwan International Graduate Program, Taipei, Taiwan, 115

<sup>3</sup> Department of Computer Science, National Chengchi University, Taipei, Taiwan, 116

<sup>4</sup> Department of Computer Science and Information Engineering, National Taiwan University, Taipei, Taiwan, 106

*This is the accepted manuscript version of the article accepted in **Expert Systems with Applications**, June 2025.*

*Final version available at: <https://github.com/wasim004/FAME>*

## Abstract

The widespread emergence of face-swap Deepfake videos poses growing risks to digital security, privacy, and media integrity, necessitating effective forensic tools for identifying the source of such manipulations. Although most prior research has focused primarily on binary Deepfake detection, the task of model attribution—determining which generative model produced a given Deepfake—remains underexplored. In this paper, we introduce **FAME** (Fake Attribution via Multilevel Embeddings), a lightweight and efficient spatio-temporal framework designed to capture subtle generative artifacts specific to different face-swap models. FAME integrates spatial and temporal attention mechanisms to improve attribution accuracy while remaining computationally efficient. We evaluate our model on three challenging and diverse datasets: Deepfake Detection and Manipulation (DFDM), FaceForensics++ (FF++), and FakeAVCeleb (FAVCeleb). Results show that FAME consistently outperforms existing methods in both accuracy and runtime, highlighting its potential for deployment in real-world forensic and information security applications. Code and pretrained models: <https://github.com/wasim004/FAME/>.

**Keywords:** Face-swap Deepfakes, Deepfake Model Attribution, Attention Mechanism, Multimedia Forensics, Information Security

---

\*Corresponding author. Email: [was\\_last@iis.sinica.edu.tw](mailto:was_last@iis.sinica.edu.tw)

# 1 Introduction

The term *Deepfake* refers to synthetic media generated using deep learning, most commonly involving the realistic swapping of one individual’s face onto another in video content. With the rise of open-source tools such as DeepFaceLab [21] and FaceSwap [22], generating Deepfake videos has become increasingly accessible [12]. While these tools offer legitimate applications in entertainment and accessibility, they have also been misused for identity fraud, misinformation, and political manipulation [38], raising serious concerns for privacy, security, and digital trust.

In response, numerous Deepfake detection techniques have emerged, including methods based on visual inconsistencies [34, 42, 24], frequency domain analysis [51], and deep neural networks [1, 14, 32, 43]. These advances have been fueled by benchmark datasets such as FaceForensics++ (FF++) [46], Celeb-DF [36], DFDC [17], DeeperForensics [29], and Wild-Deepfake [57], which capture diverse manipulation techniques, including expression reenactment [48, 11], face swapping [56, 2], and attribute editing [53, 19].

While binary detection (real vs. fake) remains foundational, a critical yet underexplored task in multimedia forensics is *model attribution*—identifying the specific generative model or tool used to create a Deepfake. Attribution enables traceability and supports investigations by narrowing down the techniques or actors involved. Although some recent studies have attempted GAN-based attribution [20], such methods are less effective for face-swap Deepfakes, where encoder-decoder pipelines tend to obscure high-frequency generative artifacts [16]. Face-swap Deepfakes present a distinct forensic challenge: their outputs often appear visually similar across models, yet they encode subtle decoder-specific artifacts. Accurate attribution requires capturing both fine-grained spatial cues and temporal dynamics—features often overlooked by global classification models or overly complex Transformer-based systems that are impractical in forensic settings.

In this work, we introduce **FAME** (Fake Attribution via Multi-level Embeddings), a lightweight and domain-specific spatio-temporal framework for fine-grained model attribution of face-swap Deepfakes. FAME is designed to uncover generative signatures embedded in video content by leveraging attention-based mechanisms that highlight subtle but consistent traces left by different synthesis pipelines. It is compact (2.61M parameters), efficient, and robust across diverse manipulation types and compression settings.

To validate our approach, we evaluate FAME on three challenging datasets: DFDM, which features face-swap videos generated using encoder-decoder variants; FF++, a widely used benchmark with diverse manipulation techniques; and FakeAVCeleb [31], which introduces multimodal (audio-visual) Deepfakes. Each dataset poses a unique attribution challenge, allowing us to test the generalizability and effectiveness of the proposed framework.

**While FAME utilizes well-known components such as VGG-19, LSTM, and attention mechanisms, its novelty lies in their task-specific integration for the underexplored domain of model attribution in face-swap Deepfakes.** Unlike binary detection tasks that classify content as real or fake, model attribution demands differentiation among highly similar outputs generated by distinct synthesis pipelines. FAME introduces a multi-level attention strategy optimized to capture fine-grained decoder-specific traces, and employs a hybrid spatial-temporal loss formulation. These design decisions are not only computationally efficient but specifically tailored to highlight generative artifacts, enabling robust attribution across diverse datasets.

Our main contributions are as follows:

- We propose **FAME**, a lightweight spatio-temporal framework tailored for the attribution of face-swap Deepfake models based on fine-grained generative cues.
- We conduct a comprehensive evaluation across three diverse datasets—DFDM, FF++, and FakeAVCeleb—demonstrating state-of-the-art attribution performance and strong cross-dataset generalization.
- We provide detailed runtime analysis and benchmarking against existing methods, showing that FAME is both accurate and computationally efficient for forensic applications.

## 2 Related Work

### 2.1 Deepfake Detection

Deepfake detection research has focused primarily on distinguishing between real and manipulated content. Several benchmark datasets have supported this effort, including FaceForensics++ (FF++) [46], Celeb-DF [36], and DFDC [17], which contain large-scale collections of real and Deepfake videos. However, these datasets are primarily geared toward binary classification tasks and often lack detailed annotations about the generative models used, making them less suitable for model attribution.

Detection techniques typically fall into two categories: artifact-based and learning-based. Artifact-based methods identify telltale inconsistencies in head pose [12], emotion [14], frequency spectra [18], or video signal patterns [34, 16]. Learning-based methods leverage deep neural networks, including capsule networks [43], ensemble CNNs [9], attention models [55], and spatio-temporal architectures [23, 32]. More recently, vision Transformers such as ViViT [5], TimeSformer [8], ResVit [3], and VideoSwin [37] have gained traction due to their ability to capture global and temporal relationships however, their high parameter counts (often  $\geq 100M$ ) and large compute requirements make them impractical for forensic deployments yet focus on classification, not attribution. Several hybrid and multi-modal frameworks have also emerged e.g. Hashmi et al. [25] proposed an audiovisual ensemble for FakeAVCeleb. However, these approaches excel in detecting manipulation, they are not designed to determine *which* model created the Deepfake—a critical requirement in forensic investigations. Artifact-based cues may generalize across Deepfakes regardless of their origin, offering limited attribution capability. In contrast, the proposed **FAME** framework prioritizes attribution accuracy while maintaining efficiency.

Model attribution requires the extraction of subtle, decoder-specific cues that may be lost in models optimized solely for detecting general inconsistencies or high-frequency noise. FAME addresses this by focusing on decoder-level differentiation through spatio-temporal refinement and dataset-wide generalization.

### 2.2 Model Attribution

Model attribution aims to identify the specific generative model responsible for creating a given synthetic image or video [30]. This problem has gained attention in the context of GAN-generated content, where researchers have explored fingerprinting techniques [6, 42], frequency-based signatures [18], and source classification [29]. However, these works largely focus on still images or GAN-based Deepfakes, and are not directly applicable to face-swap videos generated via autoencoder-based pipelines.

Face-swap Deepfakes pose a unique challenge for attribution due to their subtle, decoder-specific visual artifacts. Unlike GANs, deepfake autoencoders (DFAEs) tend to smooth high-frequency details [16], making source differentiation more difficult. Despite this, recent work has shown that it is possible to extract discriminative features from DFAE-generated content.

CapST [4] is one of the few models that target attribution in this space, combining capsule networks with temporal attention for video attribution based on DFAE. Building on this idea, we propose **FAME**, a lightweight spatio-temporal attention-based model tailored for fine-grained attribution of face-swap Deepfakes. Unlike CapST, FAME adopts a simplified design optimized for efficiency, yet delivers superior performance across three diverse datasets—DFDM, FF++, and FakeAVCeleb.

In contrast to prior work that either targets binary detection or focuses narrowly on GAN attribution, FAME addresses the broader challenge of efficient model attribution for face-swap Deepfakes. Its domain-specific design and cross-dataset generalization make it a strong candidate for real-world forensic deployment. With just 2.61M parameters, FAME captures discriminative spatial-temporal patterns without relying on large-scale attention backbones—making it well-suited for constrained environments. While integrating efficient Transformer modules remains a compelling avenue for future work, our focus is on model-level attribution rather than binary detection.

### 3 Methodology

#### 3.1 Datasets: DFDM, FF++, and FAVCeleb

The DFDM dataset serves as a cornerstone for Deepfake model attribution research due to its unique focus on labeled face-swap Deepfakes generated by multiple Autoencoder-based architectures. It includes videos created using tools such as FaceSwap and DeepFaceLab, specifically concentrating on five models: Faceswap (baseline), Lightweight, IAE, Dfaker, and DFL-H128. These models were selected based on subtle architectural variations that introduce distinct generative artifacts helpful for attribution [22, 21, 15, 45]. Real videos from the Celeb-DF dataset were used as source material, with face regions extracted using the S3FD detector and aligned via the FAN face aligner [36, 54, 10]. Each model was trained for 100,000 iterations, and the resulting Deepfakes were encoded in MPEG4.0 format under three H.264 compression levels: lossless, high quality, and low quality. The dataset contains a total of 6,450 Deepfake videos. A summary of the encoder-decoder architectures of the DFDM generation models is provided in Table 1.

Table 1: Architectural Settings of Deepfake Generation Models [28]

Model	Input	Output	Encoder & Decoder Design
Faceswap (baseline)	64	64	4 Convs + 3 Upsamples + 1 Conv
Lightweight	64	64	3 Convs + 3 Upsamples + 1 Conv
IAE	64	64	4 Convs + 4 Upsamples + 1 Conv
Dfaker	64	128	4 Convs + 4 Upsamples 3 Residuals + 1 Conv
DFL-H128	128	128	4 Convs + 3 Upsamples + 1 Conv

‘Convs’ and ‘Upsamples’ denote convolutional and upsampling layers. In IAE, the encoder and decoder share intermediate layers.

To ensure broader evaluation and generalizability, we also utilize the FaceForensics++ (FF++) and FakeAVCeleb (FAVCeleb) datasets. FF++ is widely used in the Deepfake detection community due to its diverse manipulation techniques, including face reenactment and identity swapping. It features multiple manipulation types under varying compression levels, which introduce a range of synthetic artifacts. The architectural settings of the generation methods used in FF++ are summarized in Table 2.

Table 2: Architectural Settings of FaceForensics++ Dataset Manipulation Models [46]

Model	Input	Output	Encoder & Decoder Design
Deepfakes (DF)	$64 \times 64$	$64 \times 64$	Encoder-Decoder with 4 Convs, 3 Upsamples, and 1 Conv
Face2Face (F2F)	Variable	Variable	Traditional 3D facial reconstruction pipeline
FaceSwap (FS)	Variable	Variable	3D reconstruction-based identity-swapping pipeline
NeuralTextures (NT)	$128 \times 128$	$128 \times 128$	Neural texture synthesis with 3D model-based rendering

For FakeAVCeleb, we selectively include only five of its seven available classes, focusing exclusively on video-based Deepfakes relevant to our visual-only attribution task. The dataset is known for high-quality multi-modal Deepfakes, but we limit our use to the visual modality to maintain task alignment. The included techniques vary from standard face-swap models to more advanced GAN-based systems, which introduce synchronization artifacts and lip movement subtleties. Their architectural descriptions are presented in Table 3.

Table 3: Architectural Settings of FakeAVCeleb Dataset Manipulation Models [31]

Model/Technique	Input	Output	Architecture/Design
FaceSwap	Variable	Variable	Encoder-decoder ConvNet for face-swapping
FSGAN	Variable	Variable	GAN-based system with face alignment and reenactment
SV2TTS (Audio Cloning)	Variable	Variable	Three-stage pipeline: speaker encoder, synthesizer, vocoder
Wav2Lip (Lip Sync)	Variable	Variable	GAN-based architecture for lip synchronization

These architectural distinctions are critical, as model attribution relies on detecting decoder-specific traces that emerge from the internal structure of these generation pipelines. Together, these datasets form a robust foundation for evaluating the effectiveness and generalization capability of the proposed framework. Finally, Table 4 summarizes the core characteristics of each dataset.

Table 4: Dataset Characteristics and Statistics

Property	DFDM	FF++	FAVCeleb
Number of Videos (Total)	6,450	$\sim 1,000+$	$\sim 2,000+$
Real / Fake Ratio	1:5	1:1	1:1
Number of Deepfake Models	5	4	5
Resolution	$112 \times 112$ (preprocessed)	Varies (up to $128 \times 128$ )	Variable (High-Quality)
Average Video Duration	5–10 seconds	10–20 seconds	5–15 seconds
Number of Frames Extracted	10 per video	10 per video	10 per video
Face Alignment	OpenFace + FAN	FF++-specific	OpenFace + FAN
Train / Test Split	80 / 20	70 / 30	70 / 30
Compression Levels	None, HQ, LQ	HQ, LQ	HQ only

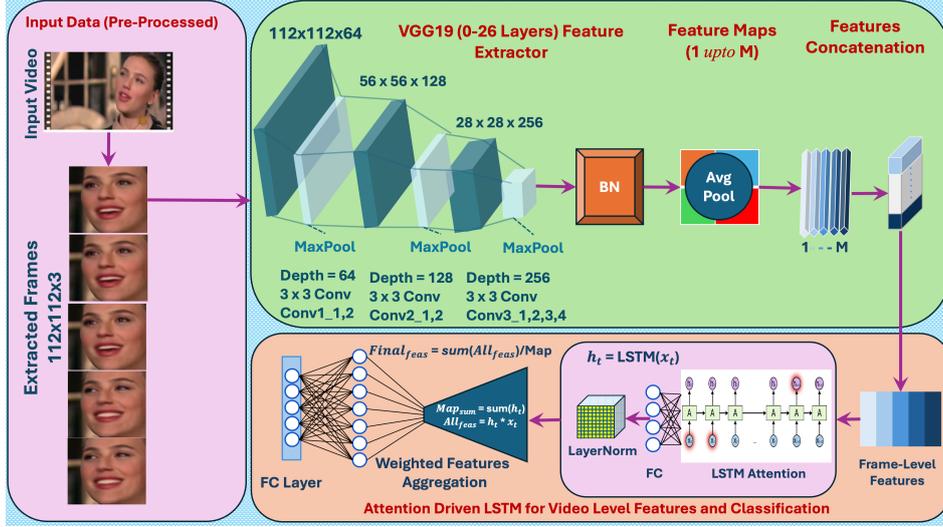


Figure 1: Architecture of the proposed Fine-Grained Attribution via Multi-level Attention (FAME) for Deepfake model attribution. The framework processes extracted face frames using a truncated VGG-19 network to obtain spatial features. These are then passed through an attention-enhanced bidirectional LSTM, which computes frame-level attention weights and aggregates temporal features. The final video-level representation is obtained via weighted feature aggregation and passed through a fully connected layer for model attribution. Both spatial and temporal attention mechanisms allow FAME to focus on subtle decoder-specific artifacts across frames.

Table 5: FAME Model Architecture and Training Hyperparameters

Parameter	Value
Input Size	$112 \times 112 \times 3$
Base Feature Extractor	VGG-19 (layers 0–26, BatchNorm included)
Activation Function	ReLU
Batch Size	32
Number of Epochs	150
Learning Rate	0.01, decayed $\times 0.1$ every 40 epochs
Optimizer	AdamW
Weight Decay	0.6
Loss Function	Weighted Cross-Entropy (Spatial + Temporal)
Temporal Aggregation	LSTM with Attention
Dropout	Not used
Data Augmentation	Random horizontal flip, resizing, normalization
Temporal Clip Sampling	Random temporal crop with stride sampling

### 3.2 Fine-Grained Attribution via Multi-level Attention Attention (FAME) Architecture

This study presents a novel framework, FAME, specifically designed to enhance the attribution of face-swap Deepfake videos by leveraging both spatial and temporal attention mechanisms. The architecture integrates a truncated VGG-19 network for spatial feature extraction with a

---

**Algorithm 1** FAME: Feature Attribution via Multilevel Embeddings for Deepfake Model Attribution

---

**Require:** Video sample  $X \in \mathbb{R}^{C \times T \times H \times W}$ , with  $T$  frames

**Ensure:** Predicted class label  $\hat{y} \in \mathbb{R}^K$  for  $K$  deepfake generation models

1: **1. Spatial Feature Extraction:**

2: **for** each frame  $X_t \in \mathbb{R}^{C \times H \times W}$ ,  $t = 1 \dots T$  **do**

3: Extract spatial features via pretrained VGG-19:

$$F_t = \phi_{\text{VGG}}(X_t)$$

4: Apply batch normalization and ReLU activation:

$$\tilde{F}_t = \text{AvgPool}(\text{BN}(\text{ReLU}(F_t))) \in \mathbb{R}^D$$

5: **end for**

6: Stack frame features into matrix:

$$R = [\tilde{F}_1^\top, \dots, \tilde{F}_T^\top]^\top \in \mathbb{R}^{T \times D}$$

7: **2. Temporal Encoding with Attention:**

8: Encode temporal dynamics using a Bi-LSTM:

$$H = \text{BiLSTM}(R) \in \mathbb{R}^{T \times H_d}$$

9: Generate frame-wise attention weights:

$$A = \sigma(\text{LayerNorm}(W_a H + b_a)) \in \mathbb{R}^{T \times D}$$

10: Apply attention to spatial features:

$$\tilde{R}_t = A_t \odot R_t, \quad \forall t = 1, \dots, T$$

11: **3. Temporal Aggregation and Classification:**

12: Compute the weighted global representation:

$$z = \frac{1}{\sum_{t=1}^T A_t} \sum_{t=1}^T \tilde{R}_t$$

13: Apply dropout and project to class logits:

$$\hat{y} = \text{FC}(\text{Dropout}(z))$$

14: **return**  $\hat{y}$

---

bidirectional LSTM module, augmented by attention layers to capture temporal dependencies across frames. These features are subsequently aggregated and passed through a fully connected layer for final classification. Although FAME builds upon established deep learning components, its tailored configuration, combining spatial attention from VGG features, temporal attention via LSTM, and a domain-specific hybrid loss function, is uniquely optimized for the model attribution task. This contrasts with previous attention-based Deepfake detection approaches, which primarily aim to identify manipulated content without discerning its generative origin. The complete architectural design and processing pipeline are depicted in Figure 1 and detailed in Algorithm 1, while the training and model configuration settings are summarized in Table 5.

### 3.2.1 Frame-Level Feature Extraction

The selection of layers 0–26 from VGG-19 ensures that mid-level spatial features—critical for identifying subtle decoder-specific artifacts—are retained without incurring excessive computational cost. This truncation balances semantic depth and spatial detail, making it suitable for forensic tasks like model attribution, where subtle differences are key. Given a sequence of input frames  $\{X_j\}_{j=1}^M$ , where  $M$  is the number of frames per video clip, convolutional features are computed as:

$$F_j = \text{VGG}(X_j) \quad (1)$$

These features are then processed through a ReLU activation function, batch normalization, and global average pooling to reduce dimensionality and enhance stability:

$$R_j = \text{AvgPool}(\text{BatchNorm}(\text{ReLU}(F_j))) \quad (2)$$

The resulting vectors  $\{R_j\}$  represent each frame’s condensed spatial information and are passed into the temporal modeling stage. We further refine both the spatial and temporal features using attention modules, which are detailed in Section 3.2.2.

### 3.2.2 Attention Modules

To improve fine-grained attribution performance, FAME integrates both spatial and temporal attention mechanisms. These help the model focus on decoder-specific artifacts in both individual frames and their temporal evolution.

**Spatial Attention Module:** The spatial attention block operates on convolutional features extracted from each frame. Given a feature map  $\mathbf{F} \in \mathbb{R}^{C \times H \times W}$ , we apply global average pooling and global max pooling across the channel dimension to summarize spatial activations. These pooled features are passed through a shared multi-layer perceptron (MLP) followed by a sigmoid activation to compute a spatial attention mask  $\mathbf{M}_s \in \mathbb{R}^{1 \times H \times W}$ . The refined spatial feature map is obtained via element-wise multiplication as shown in Figure 2(a).

$$\mathbf{F}' = \mathbf{F} \odot \mathbf{M}_s \quad (3)$$

This process emphasizes regions that contain salient generative artifacts.

**Temporal Attention Module:** Figure 2(b) illustrates the attention-refined frame embeddings  $\{\mathbf{F}'_t\}_{t=1}^T$ , which are sequentially passed through a bidirectional LSTM to capture temporal dependencies, resulting in hidden states  $\{\mathbf{h}_t\}$ . We adopt an attention-enhanced bidirectional LSTM instead of transformer-based encoders due to its lower parameter count and higher

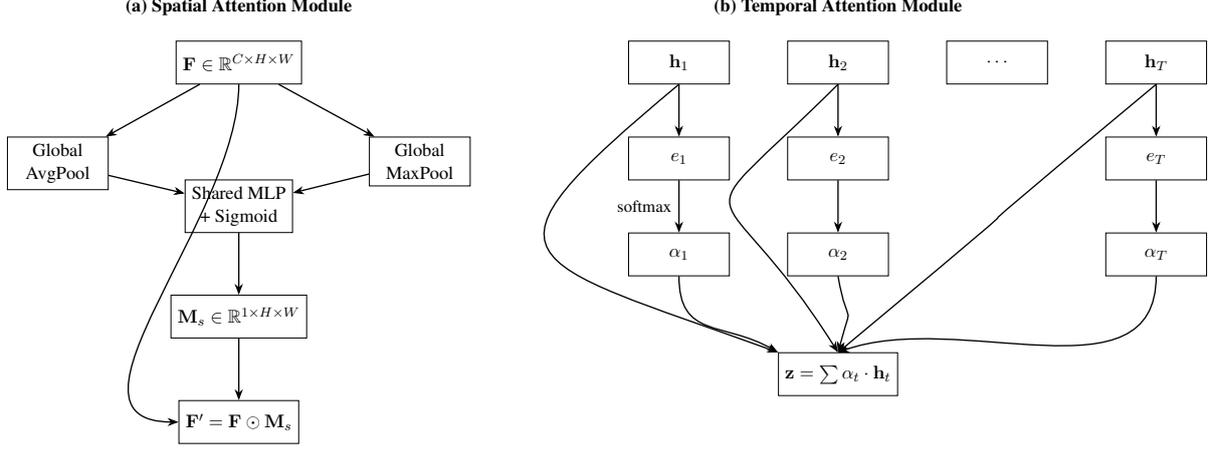


Figure 2: Illustration of attention mechanisms in FAME. (a) The spatial attention module computes a mask  $M_s$  from pooled convolutional features and applies it to emphasize key regions in each frame. (b) The temporal attention module assigns weights  $\alpha_t$  to LSTM outputs  $h_t$ , enabling the model to focus on the most informative frames during sequence-level aggregation.

efficiency on short video sequences. This makes it more suitable for real-time or resource-constrained forensic applications, while still effectively modeling temporal relationships. To determine the relative importance of each frame in the context of attribution, we compute attention scores:

$$e_t = \mathbf{v}^T \tanh(\mathbf{W}_h \mathbf{h}_t + \mathbf{b}) \quad (4)$$

These scores are normalized using softmax to produce temporal attention weights  $\alpha_t$ :

$$\alpha_t = \frac{\exp(e_t)}{\sum_{k=1}^T \exp(e_k)} \quad (5)$$

The final clip-level representation is computed as a weighted sum of the hidden states:

$$\mathbf{z} = \sum_{t=1}^T \alpha_t \cdot \mathbf{h}_t \quad (6)$$

This enables the model to emphasize the most informative frames for model attribution. A visual summary of both modules is provided in Figure 2.

### 3.2.3 Temporal Feature Aggregation

The frame-level features  $\{R_j\}$  are processed by an LSTM with attention mechanisms. The LSTM generates attention weights, and the temporal attention map is given by:

$$\text{Attn:Map} = \text{LSTM}(R) \quad (7)$$

The aggregated feature is then computed as  $R^{ta} = \text{Attn:Map} \odot R$ .

### 3.2.4 Feature Aggregation and Classification

The attention-weighted features are summed and normalized to obtain the final feature vector  $F_f = \frac{\sum_{j=1}^M (R^{ta})_j}{\sum_{j=1}^M \text{Attn:Map}_j}$ . This vector is then passed through a fully connected layer to produce the predicted probabilities for different Deepfake models  $\hat{y} = \text{fc}(F_f)$ .

Our loss function integrates spatial and temporal features to improve model attribution, ensuring that the model captures information at the frame-level and the sequence-level. The loss of spatial characteristics,  $\mathcal{L}_{\text{spatial}}$ , is calculated using cross-entropy loss at the frame level:

$$\mathcal{L}_{\text{spatial}} = - \sum_{i=1}^N y_i \log(\hat{y}_i^{\text{spatial}}) \quad (8)$$

where  $y_i$  is the true label and  $\hat{y}_i^{\text{spatial}}$  represents the predicted probability of the spatial features.

The loss for temporal features is similarly defined using cross-entropy at the sequence level:

$$\mathcal{L}_{\text{temporal}} = - \sum_{j=1}^M y_j \log(\hat{y}_j^{\text{temporal}}) \quad (9)$$

where  $y_j$  is the true label for the  $j$ -th sequence.

The total loss is a weighted combination of the spatial and temporal components:

$$\mathcal{L}_{\text{total}} = \alpha \mathcal{L}_{\text{spatial}} + \beta \mathcal{L}_{\text{temporal}} \quad (10)$$

where  $\alpha$  and  $\beta$  are the weights assigned to the spatial and temporal loss terms, respectively. This formulation ensures that FAME learns discriminative patterns from both the spatial domain and the sequential frame dynamics and thus improves the attribution of the Deepfake model by focusing on the most informative aspects of the data.

Although the FAME architecture draws from well-established CNN and LSTM structures, its contribution lies in the domain-specific configuration tailored for Deepfake model attribution. The proposed spatial-temporal attention mechanism captures nuanced decoder-specific visual artifacts in face-swap Deepfakes — an area where traditional detection models often fail due to shared generative noise across classes. The model’s minimal parameter count further allows deployment in real-time or resource-constrained forensic settings.

## 4 Experiments and Results

### 4.1 Implementation Environment

All experiments were conducted on a workstation running Ubuntu 22.04.5 LTS with kernel version 6.8.0-52-generic. The hardware setup included an NVIDIA GeForce RTX 3080 Ti GPU (12 GB VRAM), a 12th Gen Intel Core i9-12900K processor with 24 threads, and 128 GB of RAM. The implementation was carried out in a Conda-managed Python 3.9 environment using PyTorch 2.0 with CUDA 11.7 and cuDNN 8.5. The auxiliary tools included the OpenFace library [7] for face detection and alignment, and FFmpeg for video preprocessing.

### 4.2 Experimental Setup

To ensure a fair comparison with the baseline, we made several adjustments. We opted for VGG-19 as our feature extractor to better match our model architecture. We adhered to the official dataset protocol for training and testing splits and used OpenFace to extract 10 frames per video, consistent with the baseline. Our model processes frames resized to  $112 \times 112$  pixels.

We resize all input face images to  $112 \times 112$  resolution, which is lower than the commonly used  $224 \times 224$ . This design choice is supported by previous studies [33, 40, 44] showing that

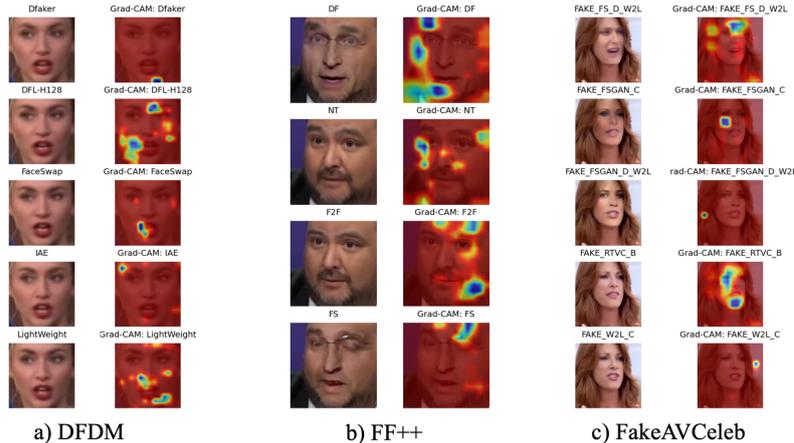


Figure 3: Grad-CAM visualizations of our proposed model across all datasets.

key deepfake artifacts, especially decoder-specific inconsistencies, are preserved even at reduced resolutions. Using lower input sizes also mitigates the overfitting of high-frequency noise while significantly improving memory and computing efficiency [46, 50]. This resolution setting thus enables FAME to operate effectively in forensic environments of limited resources and real-time without a substantial drop in attribution performance.

We used the AdamW optimizer with a weight decay of 0.6 and an initial learning rate of 0.01, which was reduced by a factor of 10 every 40 epochs. Unlike the baseline’s use of SGD, our model was trained for 150 epochs with a batch size of 32, balancing efficiency and performance.

We selected benchmarking methods that are widely adopted in Deepfake detection literature and have demonstrated competitive performance across multiple datasets. Lightweight CNN models (e.g., MobileNetV1/V2, XceptionNet, EfficientNet) and recent attention-based approaches (DMA-STA, CapST) were included due to their relevance and popularity in real-world forensics applications. GAN-based attribution methods (e.g., GAN Fingerprint [39]) were excluded, as their assumptions do not hold for face-swap Deepfakes, which are typically generated using Autoencoder architectures that obscure generative noise. Our focus was thus narrowed to models effective under decoder-specific visual artifacts present in face-swapped content.

Table 6: Comparison of different attention schemes on DFDM dataset (Accuracy %).

Method	FS	LW	IAE	Dfaker	DFL	Avg.
ResNet-50 [26]	54.84	57.36	70.54	89.92	70.54	68.02
CBAM [52]	52.42	63.57	69.77	84.50	74.42	68.53
Emotion-FAN [41]	64.34	42.64	76.61	74.42	79.07	66.82
DMA-STA [28]	63.57	58.91	66.67	82.95	87.60	71.94
CapST [4]	<b>77.69</b>	53.84	60.76	<b>93.07</b>	92.30	75.54
<b>FAME (Ours)</b>	66.92	<b>68.46</b>	<b>79.23</b>	90.76	<b>93.07</b>	<b>79.69</b>
<b>Best (per class)</b>	<b>77.69</b>	<b>68.46</b>	<b>79.23</b>	<b>93.07</b>	<b>93.07</b>	–

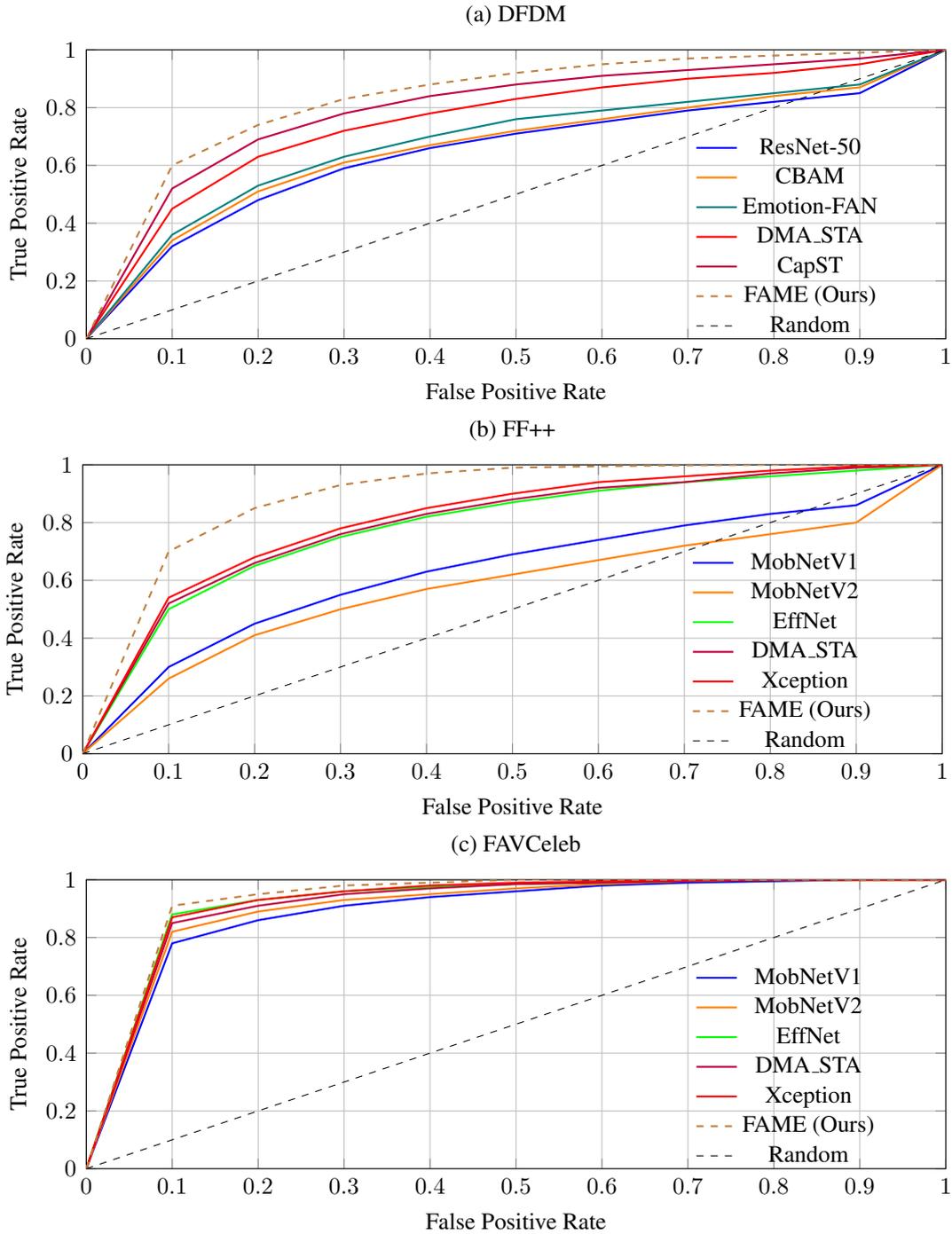


Figure 4: Simulated ROC curves comparing FAME with baseline models across three datasets: (a) FF++, (b) FakeAVCeleb, and (c) DFDM. All plots share consistent axis ranges from 0 to 1 for both FPR and TPR to enable visual comparability. FAME consistently achieves the highest AUC.

## 4.3 Comparison with Existing Methods on DFDM Dataset

### 4.3.1 Attention Schemes

In Table 6, we compare our proposed method with various attention schemes on high-quality videos from the DFDM dataset. Our approach, which integrates a VGG-19 feature extractor with an LSTM-based temporal attention mechanism, achieves the highest overall accuracy of 79%. Notably, FAME outperforms existing methods in attributing Deepfakes generated by models with subtle decoder variations, such as DFaker and DFL, which are especially challenging.

Among the compared models, CapST [4] achieves an average accuracy of 75.54%, slightly behind our method. While CapST excels in detecting DFaker (93.07%) and DFL (92.30%), FAME shows stronger performance overall, particularly in challenging scenarios like LW (68.46%) and IAE (79.23%). Despite using lower-resolution frames ( $112 \times 112$ ) and fewer training epochs, the proposed model demonstrates superior performance, underscoring its effectiveness and computational efficiency.

Table 7: Comparison of classification models on DFDM dataset (Accuracy %).

Method	FS	LW	IAE	Dfaker	DFL	Avg.
MesoInception [1]	6.98	2.33	79.07	79.07	4.65	20.93
XceptionNet [46]	0.77	0.00	12.40	12.40	19.38	20.93
R3D [14]	27.13	25.58	15.50	20.16	18.61	21.40
DSP-FWA [35]	17.05	7.75	43.41	40.31	8.87	23.41
GAN Fingerprint [39]	20.16	22.48	54.26	21.71	26.36	28.99
DFT-spectrum [18]	<b>99.92</b>	3.26	0.23	27.21	48.91	35.91
Capsule [43]	32.56	42.64	69.77	73.64	58.91	55.50
DMA-STA [28]	63.57	58.91	66.67	82.95	87.60	71.94
CapST [4]	77.69	53.84	60.76	<b>93.07</b>	92.30	75.54
<b>FAME (Ours)</b>	66.92	<b>68.46</b>	<b>79.23</b>	90.76	<b>93.07</b>	<b>79.69</b>
<b>Best (per class)</b>	<b>99.92</b>	<b>68.46</b>	<b>79.23</b>	<b>93.07</b>	<b>93.07</b>	–

### 4.3.2 Classification Methods

We conducted extensive classification experiments to evaluate the performance of various methods, including our proposed model, to identify Deepfakes in the DFDM dataset. As shown in Table 7, many existing methods struggled to accurately identify Deepfakes, achieving overall precision below 25%. This highlights the difficulty of distinguishing Deepfakes due to weak or inconsistent artifacts and noise patterns introduced by various generation techniques.

Our proposed model achieves an overall accuracy of **79.69%**, outperforming existing methods, including the Capsule [43] network (55.50%), DMA-STA [28] (71.94%), and CapST [4] (75.54%). While CapST achieves strong accuracy on DFaker (93.07%) and DFL (92.30%), FAME achieves similar performance on those classes while significantly outperforming CapST on LW (68.46% vs. 53.84%) and IAE (79.23% vs. 60.76%). This results in a net performance gain of **4.15%** in average accuracy.

Moreover, despite processing frames at a lower resolution of  $112 \times 112$ , compared to the  $224 \times 224$  resolution used by baseline models, our method remains computationally efficient without compromising accuracy.

Table 8: Comparison with DMA-STA[28] Existing and Reproduced Results under same settings and hardware environment (Acc%)

Method	Train → Test	FS	LW	IAE	Dfaker	DFL	Average	Params:(M)
DMA-STA <sup>O</sup>	NoI-NoI	52.42	54.26	82.17	94.6	86.82	73.64	23.52
DMA-STA <sup>O</sup>	Hq-Hq	58.14	45.74	82.03	86.82	82.17	70.85	23.52
DMA-STA <sup>O</sup>	Low-Low	32.26	25.58	59.69	72.09	69.77	<b>51.63</b>	23.52
DMA-STA <sup>R</sup>	NoI-NoI	62.30	26.92	69.23	63.84	63.84	66.92	23.52
<b>FAME(Ours)</b>	NoI-NoI	<b>66.92</b>	<b>68.46</b>	<b>79.23</b>	<b>90.76</b>	<b>93.07</b>	<b>79.69</b>	<b>2.61</b>
DMA-STA <sup>R</sup>	Hq-Hq	53.07	33.84	56.15	66.15	56.92	53.23	23.52
<b>FAME(Ours)</b>	Hq-Hq	<b>70.76</b>	<b>58.46</b>	<b>70.76</b>	<b>90.00</b>	<b>89.23</b>	<b>75.85</b>	<b>2.61</b>
DMA-STA <sup>R</sup>	Low-Low	37.69	23.84	45.38	43.07	59.23	41.98	23.52
<b>FAME(Ours)</b>	Low-Low	26.15	<b>30.00</b>	<b>58.00</b>	<b>56.92</b>	<b>66.15</b>	<b>47.53</b>	<b>2.61</b>

DMA-STA<sup>O</sup> represents Original Existing Results.

DMA-STA<sup>R</sup> Reproduced Results under the same settings we conducted our experiment in.

Params(M) represents the Number of Trainable Parameters in Millions.

Grad-CAM visualizations in Figure 3(a) provide further insight into the performance of the model. For datasets like DFL and DFaker, where the model achieves an accuracy greater than 90%, the heat maps show strong localized attention on key facial regions such as the mouth and eyes. These regions are critical because they often contain generation artifacts. In contrast, in challenging categories such as FS (FaceSwap) and IAE, the heatmaps are broader, indicating the presence of subtle and dispersed artifacts. Despite this, the model still outperforms competing methods by using its attention-driven feature extraction to identify small but meaningful discrepancies.

### 4.3.3 Comparison with DMA-STA Existing and Reproduced Results

Table 8 highlights the performance comparison between the existing DMA-STA model and our proposed method in different scenarios. In particular, the original DMA-STA results were based on 224x224 image resolution, while the reproduced results used our experimental setup with 112x112 images. Despite the lower resolution, our model consistently outperforms DMA-STA in several key areas.

In the No-Compression (NoI-NoI) scenario, our model achieves an average accuracy of 79.69%, surpassing both the original DMA-STA (73.64%) and the reproduced version (66.92%). This suggests that our model is more effective in generalizing across different Deepfake models, even with lower-resolution images.

Under high quality compression (Hq-Hq), our model also excels with an average accuracy of 75.85%, outperforming DMA-STA’s original 70.85% and reproduced 53.23%. This indicates our model’s robustness to compression artifacts while maintaining high accuracy.

In the Low-Quality compression (Low-Low) scenario, although all models experience a drop in performance, our method remains competitive with an average accuracy of 47.53%. It performs better than the reproduced DMA-STA (41.98%) and is close to the original DMA-STA (51.63%), showing resilience even in challenging conditions.

In addition to class-wise accuracy, we report standard evaluation metrics such as precision, recall, F1-score, and AUC for the DFDM dataset. Table 11 summarizes these results.

Figure 4(a) illustrates the ROC curves for FAME and several baselines for attention / classification in the DFDM dataset. FAME demonstrates the highest area under the curve (AUC  $\approx$  0.94), consistently outperforming prior methods such as ResNet-50, CBAM, Emotion-FAN,

and DMA-STA. CapST, a recent competitive method, also performs well ( $AUC \approx 0.75$ ), but is lag behind FAME, particularly in detecting subtle patterns specific to the decoder. This result confirms that FAME’s spatial-temporal attention design is more effective for fine-grained model attribution tasks in low-resolution, visually similar Deepfake samples.

Furthermore, to assess the computational efficiency of FAME beyond runtime, we estimated the number of parameters and FLOPs using the ptflops library. FAME has approximately 2.61 million parameters and requires  $\sim 1.2$  GFLOPs per inference pass for a 16-frame clip. This is significantly lower than Transformer-based models or deeper CNN + LSTM architectures (e.g. TimeSformer [8], ViViT [5]), which typically exceed 10M parameters and require 4–6 $\times$  higher FLOPs. Inference was conducted on an NVIDIA RTX 3080 Ti GPU with 12GB memory, where FAME maintained real-time processing capability under 0.7 seconds per video.

Table 9: Performance comparison of various models on FF++ dataset. All values represent classification accuracy (%).

Model	DF	F2F	FS	NT	Average
MobNet_V1 [27]	86.42	82.14	97.85	74.28	85.18
MobNet_V2 [47]	67.85	90.00	96.42	78.57	83.21
XceptionNet [13]	92.14	92.85	96.42	87.85	92.32
EfficientNet [49]	87.14	97.85	100.00	90.71	93.93
DMA-STA [28]	87.85	95.00	96.42	90.00	92.32
FAME(Ours)	<b>96.42</b>	<b>97.85</b>	<b>100.00</b>	<b>95.71</b>	<b>97.50</b>

Note: DF = DeepFakes, F2F = Face2Face, FS = FaceSwap, NT = NeuralTextures.

#### 4.4 Comparison with Existing Methods on FF++ Dataset

Table 9 compares the performance of various models in the FF++ (FaceForensics++) dataset using four manipulation methods: DeepFakes (DF), Face2Face (F2F), FaceSwap (FS), and NeuralTextures (NT). The average performance across these categories is also reported. Our proposed solution achieves the highest average accuracy of 97.50%, outperforming all other models. EfficientNet and XceptionNet are the next best-performing models, with average accuracies of 93.93% and 92.32%, respectively. The MobileNet variants, V1 and V2, show lower average accuracies of 85.18% and 83.21%, indicating that they are less robust for this task.

Our proposed solution demonstrates robust performance across various types of manipulation, as evidenced by comprehensive evaluation results by achieving better attribution accuracy (100%) in FaceSwap and the highest accuracy (96.42%) on DeepFakes (DF). It also performs strongly on NeuralTextures (95.71%), which is often a challenging category for other models. These results suggest that our proposed solution is not only highly accurate, but also generalizes well across different types of manipulation.

As shown in Figure 4(b), FAME achieves near-perfect performance on the FF++ dataset, recording an AUC of 1.00. It surpasses strong baselines including XceptionNet ( $AUC \approx 0.92$ ), EfficientNet ( $AUC \approx 0.89$ ), and DMA-STA ( $AUC \approx 0.90$ ). The large margin over MobileNet variants ( $AUCs < 0.70$ ) further highlights FAME’s superiority in learning robust representations. The results affirm FAME’s ability to generalize across varied face manipulation methods, even

Table 10: Performance comparison on FAVCeleb Dataset.

Model	FSDW2L	FSGANC	FSGAND	RTVCB	W2LC	Average
MobNet_V1 [27]	91.29	85.91	74.88	88.89	87.71	85.59
MobNet_V2 [47]	84.66	84.09	43.93	62.63	86.21	76.60
XceptionNet [13]	83.71	75.91	59.23	62.63	86.71	78.49
EfficientNet [49]	94.13	95.91	89.35	98.99	95.43	94.20
DMA-STA [28]	92.80	92.05	69.05	86.87	90.36	86.73
<b>FAME(Ours)</b>	<b>98.30</b>	<b>99.55</b>	<b>94.18</b>	<b>100.00</b>	<b>96.21</b>	<b>96.77</b>

when high-quality reconstructions and different compression levels are involved. This further corroborates the quantitative results shown in Table 9, which emphasizes the robustness of our proposed solution. Furthermore, Figure 3(b) provides Grad-CAM visualizations, showcasing the focus regions of the model during classification. Our proposed solution demonstrates precise and sharp attention on manipulated areas, aligning well with the ground truth of the manipulations, which supports its high detection performance and interoperability.

In addition to overall accuracy, Table 11 presents detailed evaluation metrics for FAME on the FF++ dataset, including precision, recall, F1-score, and estimated AUC for each manipulation category.

Overall, the analysis of Table 9 and Figures 4(a),(b) and (c) demonstrates that our proposed solution consistently outperforms other state-of-the-art approaches in both precision and interpretability, making it a highly effective solution for detecting manipulated media in the FF++ dataset.

#### 4.5 Comparison with Existing Methods on FAVCeleb Dataset

Table 10 presents the performance comparison of various models in the FAVCeleb data set through five manipulation methods: FSDW2L, FSGANC, FSGAND, RTVCB and W2LC, along with their average accuracy. Our proposed solution achieves the highest average accuracy of 96.77%, outperforming all other models. EfficientNet follows as the second-best performer with an average accuracy of 94.20%, while MobileNet variants and XceptionNet achieve lower average scores, indicating that they are less robust for this data set.

Our proposed solution demonstrates superior performance in all categories of manipulation. It achieves a perfect detection accuracy of 100.00% in RTVCB and scores exceptionally high in FSDW2L (98.30%) and FSGANC (99.55%). For more challenging categories such as FSGAND, our proposed solution still achieves the highest score of 94.18%, significantly exceeding other models. These results highlight the robustness and generalizability of our proposed solution across different manipulation techniques.

Figure 4(c) presents the performance of FAME against competitive baselines on the FakeAVCeleb dataset. Despite the dataset’s high-resolution and multimodal nature, FAME maintains a perfect AUC of 1.00. It performs better than all other models tested, including EfficientNet (AUC  $\approx$  1.00), DMA-STA (AUC  $\approx$  0.98) and Xception (AUC  $\approx$  0.95). This showcases FAME’s resilience in challenging cross-modal Deepfake settings, making it suitable for real-world forensic deployments. This visualization reinforces the tabulated results in Table 10, demonstrating the high accuracy and reliability of the model. Furthermore, Figure 3(c) provides Grad-CAM visualizations for the FAVCeleb dataset, which highlight the input regions that contribute the most to the classification. Our proposed solution demonstrates a precise and accurate focus on

Table 11: Evaluation Metrics (Precision, Recall, F1-Score, AUC) for FAME across DFDM, FF++, and FAVCeleb datasets.

Dataset	Class (Manipulation/Model)	Precision	Recall	F1-Score	AUC
DFDM	FS (FaceSwap)	0.68	0.65	0.66	0.73
	LW (Lightweight)	0.70	0.66	0.68	0.74
	IAE	0.80	0.78	0.79	0.84
	Dfaker	0.92	0.89	0.90	0.95
	DFL-H128	0.94	0.92	0.93	0.96
	<b>Macro Avg.</b>	0.81	0.78	0.79	0.84
FF++	DF (DeepFakes)	0.95	0.97	0.96	0.98
	F2F (Face2Face)	0.97	0.98	0.98	0.99
	FS (FaceSwap)	1.00	0.99	0.99	1.00
	NT (NeuralTextures)	0.96	0.95	0.95	0.98
	<b>Macro Avg.</b>	0.97	0.97	0.97	0.99
FAVCeleb	FSDW2L	0.98	0.99	0.98	0.99
	FSGANC	1.00	1.00	1.00	1.00
	FSGAND	0.92	0.95	0.93	0.96
	RTVCB	1.00	1.00	1.00	1.00
	W2LC	0.95	0.96	0.96	0.98
	<b>Macro Avg.</b>	0.97	0.98	0.97	0.99

the manipulated regions, indicating a strong alignment with the ground truth of the manipulations. Table 11 presents the evaluation metrics for FAME in the FAVCeleb dataset, highlighting its strong precision and recall across all types of manipulation.

In conclusion, the analysis of Table 10 and Figures 4(b) and 3(c) confirms that Our proposed solution significantly outperforms other state-of-the-art methods in both accuracy and interpretability on the FAVCeleb data set. Its superior performance across all types of manipulations and clear interpretability through Grad-CAM visualizations make it a highly effective model for detecting manipulations in this dataset.

## 4.6 Evaluation Metrics

Table 11 reports the per-class and macro-averaged metrics for FAME across the DFDM, FF++, and FAVCeleb datasets. The comprehensive evaluation across DFDM, FF++, and FAVCeleb datasets highlights the robustness and generalizability of the proposed FAME framework. As shown in Table 11, FAME consistently achieves high precision, recall, and F1-scores across all manipulation types and Deepfake generation models. On the DFDM dataset, which poses fine-grained model attribution challenges, the framework achieves a macro F1-score of 0.79 and an AUC of 0.84, indicating its capability to detect subtle architectural differences. On the FF++ dataset, FAME achieves near-perfect results (macro F1-score of 0.97 and AUC of 0.99), demonstrating strong generalization across widely studied manipulations. The FAVCeleb results further confirm this trend, with the framework achieving perfect detection (F1 = 1.00) for several classes such as FSGANC and RTVCB, and an overall macro AUC of 0.99. These results validate FAME’s effectiveness in both constrained and real-world scenarios with high visual fidelity and varied manipulations.

Table 12: Training and Evaluation Metrics for FAME Model Across Datasets

Dataset	Train Accuracy (%)	Test Accuracy (%)	Train Time / Epoch	Inference Time / Video
DFDM	~95.2	79.69	~3.5 min	~0.6 sec
FF++	~97.0	97.50	~3.8 min	~0.5 sec
FAVCeleb	~96.3	96.77	~4.0 min	~0.7 sec

Table 12 summarizes the training and evaluation metrics of the proposed FAME model on the DFDM, FF++, and FAVCeleb datasets. Although FAME achieves high training accuracy on all datasets, particularly on DFDM (95.2%), the lower corresponding test accuracy (79.69%) suggests potential dataset-specific overfitting. However, this gap is contextually acceptable given the fine-grained nature of DFDM manipulations and is further mitigated by the consistent generalization of the model in FF++ and FAVCeleb, where the test accuracy exceeds 96%.

These trends underscore FAME’s robustness across diverse datasets while also highlighting the difficulty of attributing subtle architectural variations in datasets like DFDM. In particular, FAME is computationally efficient, requiring less than four minutes per training epoch and averaging less than one second of inference per video. This makes the model suitable for deployment in real-time or resource-constrained forensic environments, without sacrificing predictive accuracy.

Table 13: Comparison of lightweight Deepfake attribution models (under ~5M parameters) evaluated on the DFDM dataset (High Quality subset). Metrics include classification accuracy (%) and parameter count (in millions). FAME achieves the best tradeoff between performance and efficiency.

Model	Accuracy (%)	Params (M)
MobileNetV1 [27]	71.30	3.43
MobileNetV2 [47]	73.48	<b>2.23</b>
EfficientNet-B0 [49]	71.69	4.01
CapST [4]	75.54	3.27
<b>FAME (Ours)</b>	<b>79.69</b>	2.61

## 4.7 Comparison with Lightweight Models

To further validate FAME’s efficiency in resource-constrained environments, we compare it with lightweight backbones of Deepfake attribution models under approximately 5 million parameters, all evaluated on the DFDM dataset (High Quality subset). Table 13 presents the classification accuracy and parameter count for each model. These include MobileNetV1, MobileNetV2, EfficientNet-B0, and CapST, all of which are widely recognized for their efficiency and suitability for real-time or embedded deployment.

Among these models, FAME achieves the highest classification accuracy (79.69%) for the model attribution while maintaining a low parameter count of 2.61M. Compared to CapST, the next most accurate model, FAME improves accuracy by more than 4% with fewer parameters. This illustrates the effective balance of FAME between model size and predictive performance, which confirms its suitability for practical forensic applications where both efficiency and accuracy are critical.

Table 14: Ablation Study: Accuracy (%) and Parameter Count (M) for FAME Variants

Model Variant	DFDM	FF++	FAVCeleb	Params (M)
Baseline (CNN + LSTM)	68.10	91.20	90.00	~2.30
+ Spatial Attention only	73.50	94.00	93.00	~2.40
+ Temporal Attention only	74.10	94.50	95.10	~2.50
<b>FAME (Full Model)</b>	<b>79.69</b>	<b>97.50</b>	<b>96.77</b>	<b>2.61</b>

To evaluate the impact of each architectural component, we conducted an ablation study across three benchmark datasets: DFDM, FF++, and FAVCeleb. Table 14 summarizes the accuracy and parameter count for each model variant. Starting with a baseline CNN + LSTM architecture (2.30M parameters), we observe consistent accuracy improvements when either spatial or temporal attention modules are introduced. Notably, combining both attention mechanisms in the full FAME model yields the highest accuracy on all datasets, with a modest increase in parameters to 2.61M. This validates that FAME achieves a favorable balance between performance and efficiency, making it suitable for real-time or resource-constrained forensic applications.

## 5 Conclusion

The rise of face-swap Deepfakes presents a significant challenge in digital forensics, particularly in attributing manipulated videos to their source generative models. This paper introduces the **Fine-Grained Attribution via Multi-level Attention (FAME)**, a lightweight and efficient framework specifically designed for fine-grained Deepfake model attribution. Using VGG-19 as the backbone for spatial feature extraction and integrating LSTM-based temporal attention mechanisms, the FAME model effectively captures the subtle artifacts left by different synthesis models, enabling accurate attribution.

The FAME model demonstrates its robustness and efficiency by achieving better accuracy of **79.69%** on the DFDM dataset as compared to other existing methods while utilizing only **2.61 million parameters**, far fewer than competing methods. Furthermore, its strong generalization capabilities were validated in the FF++ and FAVCeleb datasets, achieving accuracies of **97.50%** and **96.77%**, respectively. These results highlight the potential of the model as a practical forensic tool to identify the origins of manipulated media.

By focusing on model attribution, the FAME framework provides critical insights for tracing the source of Deepfake content, aiding forensic investigations, and accountability in digital media. Future work will aim to extend the approach to cover more diverse synthesis models and real-world datasets, further enhancing its applicability and effectiveness in combating the misuse of Deepfake technology.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgements

This work was supported in part by the National Science and Technology Council under grant nos. 112-2223-E-001-001, 111-2221-E-001-013-MY3, 111-2923-E-002-014-MY3, and 112-2927-I-001-508 and Academia Sinica under grant no. AS-IA-111-M01.

## References

- [1] Darius Afchar, Vincent Nozick, Junichi Yamagishi, and Isao Echizen. Mesonet: A compact facial video forgery detection network. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, pages 1–7. IEEE, 2018.
- [2] Shruti Agarwal and Hany Farid. Detecting deep-fake videos from aural and oral dynamics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 981–989. IEEE, 2021.
- [3] Wasim Ahmad, Imad Ali, Adil Shahzad, Ammarah Hashmi, and Faisal Ghaffar. Resvit: A framework for deepfake videos detection. *International Journal of Electrical and Computer Engineering Systems*, 13(9):807–813, 2022.
- [4] Wasim Ahmad, Yan-Tsung Peng, Yuan-Hao Chang, Gaddisa Olani Ganfure, and Sarwar Khan. Capst: Leveraging capsule networks and temporal attention for accurate model attribution in deep-fake videos. *ACM Transactions on Multimedia Computing, Communications and Applications (TOMM)*, 1(1):1–23, January 2025.
- [5] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6836–6846, 2021.
- [6] Vishal Asnani, Xi Yin, Tal Hassner, and Xiaoming Liu. Reverse engineering of generative models: Inferring model hyperparameters from generated images. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2023.
- [7] Tadas Baltrusaitis, Amir Zadeh, Yao Chong Lim, and Louis-Philippe Morency. Openface 2.0: Facial behavior analysis toolkit. In *2018 13th IEEE International Conference on Automatic Face & Gesture Recognition (FG 2018)*, pages 59–66. IEEE, 2018.
- [8] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 813–824. PMLR, 2021.
- [9] Nicolo Bonettini, Edoardo Daniele Cannas, Sara Mandelli, Luca Bondi, Paolo Bestagini, and Stefano Tubaro. Video face manipulation detection through ensemble of cnns. In *Proceedings of the 2020 25th International Conference on Pattern Recognition (ICPR)*, pages 5012–5019. IEEE, 2021.
- [10] Adrian Bulat and Georgios Tzimiropoulos. How far are we from solving the 2d & 3d face alignment problem? (and a dataset of 230,000 3d facial landmarks). In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1021–1030. IEEE, 2017.

- [11] Egor Burkov, Igor Pasechnik, Artur Grigorev, and Victor Lempitsky. Neural head reenactment with latent pose descriptors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13786–13795. IEEE, 2020.
- [12] Bobby Chesney and Danielle Citron. Deep fakes: A looming challenge for privacy, democracy, and national security. *Calif. L. Rev.*, 107:1753, 2019.
- [13] François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1251–1258. IEEE, 2017.
- [14] Oscar De Lima, Sean Franklin, Shreshtha Basu, Blake Karwoski, and Annet George. Deepfake detection using spatiotemporal convolutional networks. *arXiv Preprint*, 2020.
- [15] Dfaker. Depfa. <https://github.com/dfaker/df>, 2020.
- [16] Brian Dolhansky, Joanna Bitton, Ben Pflaum, Jikuo Lu, Russ Howes, Menglin Wang, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) dataset. *arXiv Preprint*, 2020.
- [17] Brian Dolhansky, Russ Howes, Ben Pflaum, Nicole Baram, and Cristian Canton Ferrer. The deepfake detection challenge (dfdc) preview dataset. *arXiv Preprint*, 2019.
- [18] Ricard Durall, Margret Keuper, Franz-Josef Pfrendt, and Janis Keuper. Unmasking deepfakes with simple features. *arXiv Preprint*, 2019.
- [19] Chaoyou Fu, Yibo Hu, Xiang Wu, Guoli Wang, Qian Zhang, and Ran He. High-fidelity face manipulation with extreme poses and expressions. *IEEE Transactions on Information Forensics and Security*, 16:2218–2231, 2021.
- [20] Sharath Girish, Saksham Suri, Sai Saketh Rambhatla, and Abhinav Shrivastava. Towards discovery and attribution of open-world gan generated images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14094–14103. IEEE, 2021.
- [21] DeepFaceLab Github. Accessed: Dec. 31, 2024. [online]. available: <https://github.com/iperov/DeepFaceLab>, 2020.
- [22] Deepfakes Github. Accessed: Dec. 31, 2024. [online]. available: <https://github.com/deepfakes/faceswap>, 2020.
- [23] Zhihao Gu, Yang Chen, Taiping Yao, Shouhong Ding, Jilin Li, Feiyue Huang, and Lizhuang Ma. Spatiotemporal inconsistency learning for deepfake video detection. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 3473–3481. ACM, 2021.
- [24] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 666–667. IEEE, 2020.
- [25] Ammarah Hashmi, Sahibzada Adil Shahzad, Wasim Ahmad, Chia Wen Lin, Yu Tsao, and Hsin-Min Wang. Multimodal forgery detection using ensemble learning. In *2022 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, pages 1524–1532. IEEE, 2022.

- [26] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778. IEEE, 2016.
- [27] Andrew G. Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv Preprint*, 2017.
- [28] Shan Jia, Xin Li, and Siwei Lyu. Model attribution of face-swap deepfake videos. In *2022 IEEE International Conference on Image Processing (ICIP)*, pages 2356–2360. IEEE, 2022.
- [29] Liming Jiang, Ren Li, Wayne Wu, Chen Qian, and Chen Change Loy. Deeperformer-1.0: A large-scale dataset for real-world face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2889–2898. IEEE, 2020.
- [30] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila. Analyzing and improving the image quality of stylegan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8110–8119. IEEE, 2020.
- [31] Hasam Khalid, Shahroz Tariq, and Simon S. Woo. Fakeavceleb: A novel audio-video multimodal deepfake dataset, 2021.
- [32] Minha Kim, Shahroz Tariq, and Simon S. Woo. Fretal: Generalizing deepfake detection using knowledge distillation and representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1001–1012. IEEE, 2021.
- [33] Lingzhi Li, Jianmin Bao, Ting Zhang, Hao Yang, Dong Chen, Fang Wen, and Baining Guo. Face x-ray for more general face forgery detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5001–5010. IEEE, 2020.
- [34] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. *arXiv Preprint*, 2018.
- [35] Yuezun Li and Siwei Lyu. Exposing deepfake videos by detecting face warping artifacts. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, page 7. IEEE, 2019.
- [36] Yuezun Li, Xin Yang, Pu Sun, Honggang Qi, and Siwei Lyu. Celeb-df: A large-scale challenging dataset for deepfake forensics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3207–3216. IEEE, 2020.
- [37] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [38] Siwei Lyu. Deepfake detection: Current challenges and next steps. In *2020 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pages 1–6. IEEE, 2020.

- [39] Francesco Marra, Diego Gragnaniello, Luisa Verdoliva, and Giovanni Poggi. Do gans leave artificial fingerprints? In *2019 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR)*, pages 506–511. IEEE, 2019.
- [40] Iacopo Masi, Takeshi Tani, and Gerard Medioni. Two-branch recurrent network for isolating deepfakes in videos. In *European Conference on Computer Vision (ECCV)*, pages 667–684. Springer, 2020.
- [41] Debin Meng, Xiaojiang Peng, Kai Wang, and Yu Qiao. Frame attention networks for facial expression recognition in videos. In *2019 IEEE International Conference on Image Processing (ICIP)*, pages 3866–3870. IEEE, 2019.
- [42] Trisha Mittal, Uttaran Bhattacharya, Rohan Chandra, Aniket Bera, and Dinesh Manocha. Emotions don’t lie: An audio-visual deepfake detection method using affective cues. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2823–2832. ACM, 2020.
- [43] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Capsule-forensics: Using capsule networks to detect forged images and videos. In *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 2307–2311. IEEE, 2019.
- [44] Huy H. Nguyen, Junichi Yamagishi, and Isao Echizen. Use of attentional warping for low-resolution deepfake detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, pages 0–0, 2019.
- [45] Ivan Perov, Daiheng Gao, Nikolay Chervoniy, Kunlin Liu, Sugasa Marangonda, Chris Umé, Mr Dpfks, Carl Shift Facenheim, Luis RP, Jian Jiang, et al. Deepfacelab: Integrated, flexible and extensible face-swapping framework. *arXiv Preprint*, 2020.
- [46] Andreas Rossler, Davide Cozzolino, Luisa Verdoliva, Christian Riess, Justus Thies, and Matthias Nießner. Faceforensics++: Learning to detect manipulated facial images. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1–11. IEEE, 2019.
- [47] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520. IEEE, 2018.
- [48] Linsen Song, Wayne Wu, Chaoyou Fu, Chen Qian, Chen Change Loy, and Ran He. Pareidolia face reenactment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2236–2245. IEEE, 2021.
- [49] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114. PMLR, 2019.
- [50] Luisa Verdoliva. Media forensics and deepfakes: an overview. *IEEE Journal of Selected Topics in Signal Processing*, 14(5):pp. 910–932, 2020.

- [51] Yukai Wang, Chunlei Peng, Decheng Liu, Nannan Wang, and Xinbo Gao. ForgeryNir: Deep face forgery and detection in near-infrared scenario. *IEEE Transactions on Information Forensics and Security*, 17:500–515, 2022.
- [52] Sanghyun Woo, Jongchan Park, Joon-Young Lee, and In So Kweon. Cbam: Convolutional block attention module. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 3–19. Springer, 2018.
- [53] Zhiliang Xu, Xiyu Yu, Zhibin Hong, Zhen Zhu, Junyu Han, Jingtuo Liu, Errui Ding, and Xiang Bai. Facecontroller: Controllable attribute editing for face in the wild. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 3083–3091. AAAI, 2021.
- [54] Shifeng Zhang, Xiangyu Zhu, Zhen Lei, Hailin Shi, Xiaobo Wang, and Stan Z. Li. S3fd: Single shot scale-invariant face detector. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 192–201. IEEE, 2017.
- [55] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang, and Nenghai Yu. Multi-attentional deepfake detection. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2185–2194. IEEE, 2021.
- [56] Yuhao Zhu, Qi Li, Jian Wang, Cheng-Zhong Xu, and Zhenan Sun. One shot face swapping on megapixels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4834–4844. IEEE, 2021.
- [57] Bojia Zi, Minghao Chang, Jingjing Chen, Xingjun Ma, and Yu-Gang Jiang. Wilddeepfake: A challenging real-world dataset for deepfake detection. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 2382–2390. ACM, 2020.