

Detection of Deepfake Videos Using Long Distance Attention

VadisilaManojKumar¹, YBheemShankarRao²

Anil Neerukonda Institute of Engineering and Technology, Visakhapatnam, AP – 531162

Date of Submission: 28-05-2026

Date of Acceptance: 08-06-2026

Abstract—Deepfakes shift quickly, advancing at a rapid pace. Nearly indistinguishable artificial faces now emerge through systems like Generative Adversarial Networks, paired with encoder-decoder structures. As these fakes improve, confidence in online content weakens. Distinguishing what is real from what is fabricated becomes harder over time. Some researchers tackle the problem by treating it like eyes-or-no question. A strange glitch in a single frame might catch the eye, whereas different experts look at how things move between frames, searching for sudden breaks in flow. Most of these methods fail to consider the whole picture at once. This research highlights just how faint deepfake flaws can be - spread out, almost invisible. One corner of the face may shift differently than tissue along the chin. In other moments, timing breaks down: speech lacks smoothness, eye closures happen off rhythm. Detecting these missteps requires following fragile hints located in distant areas across the entire video sequence. Neither local convolutional attention with small kernels nor brute-force application of Vision Transformers adequately addresses both dimensions simultaneously. To close this gap, we have introduced the Long Distance Attention (LDA) mechanism. This is a novel patch-level non-local self-attention formulation that computes cross-patch relevance globally and recalibrates the deep feature map to highlight pivotal forgery regions. The LDA mechanism is instantiated in a Dual-Branch Spatial-Temporal (DBST) network comprising a *spatial branch* that discovers intra-frame artifacts and a *temporal branch* that discovers inter-frame inconsistencies, both unified by the LDA block as the shared core attention module. Training is guided by a multi-task binary cross-entropy loss over frame-level, clip-level and fused predictions. Extensive experiments on FaceForensics++ (FF++), Celeb-DF v2 and the Deepfake Detection Challenge (DFDC) datasets demonstrate that the proposed framework consistently outperforms state-of-the-art baselines, including MesoNet, Face X-Ray, LipForensics, MACNN and ViT-B/16, achieving frame-level AUC improvements of up to 1.8% on FF++ (c23) and 3.3% on Celeb-DF under cross-dataset evaluation. Ablation studies rigorously confirm that the long-distance attention span, rather than any individual architectural component, is the key driver of these improvements. The paper further provides a theoretical analysis of why long-range attention subsumes convolutional attention as a special case and why it outperforms LSTM-based temporal modeling in this setting.

Index Terms—Deepfake detection, long distance attention, spatial-temporal model, fine-grained classification, face forgery, self-attention, Vision Transformer, patch-level attention, video forensics, multimedia security, generative adversarial networks, digital content authentication.

I. INTRODUCTION

THE landscape of digital media has undergone a profound transformation in the past half-decade. Deepfake technology—the automated synthesis of photorealistic manipulated videos using deep neural networks—has moved from the research laboratory to the public domain with alarming

speed. The open-source availability of face-swapping tools [7]–[9], combined with exponentially improving generative models [1]–[3], now means that a convincing deepfake video can be produced by a non-expert with consumer-grade hardware in a matter of hours. Society is already experiencing the down-stream consequences: fabricated political speeches, celebrity impersonations used for fraud and identity theft attacks that bypass biometric authentication systems. The question of how to distinguish a synthetically generated human face from a real one has therefore become one of the most pressing open problems in computer vision and multimedia forensics. Now Deepfake detection struggles most when examining video creation methods. Analysis typically unfolds across three phases. [34] Images emerge one by one from the video frame chain. Each snapshot gets scanned next searching specifically for facial regions. Those found spots later feed into swap operations, driven by systems such as autoencoders or GANs. [1] A digital face mimics the movements of the source human, while also adapting changes in light conditions - yet pulls characteristics from another person entirely. In the end, the altered version blends smoothly into the scene where it began. Step by step, side products emerge during the process. While elements develop, variations tend to appear between sections. Look close - small flaws may reveal themselves as images render, like mismatched tones where facial zones meet surroundings, or eyes seeming slightly misaligned with adjacent features. [10],[12] Gradually, combining outputs introduces new complications. Because pictures are typically processed individually, [18] without attention to the overall sequence rhythm, [20] oddities start showing: blinks that seem off, mouth motions lacking fluidity, pulse indicators shifting abruptly across frames [17],[21].

A. Limitations of Prior Art

One way scientists spot fake videos focuses only on still images within them. These methods examine tiny details in pixels where fakes often slip up. Performance holds up when test examples match what they learned from. But once things change even slightly, accuracy drops fast. Why does this happen? Most times, a single example is all they recall, missing signs that point to deception. Because real insights aren't built on repetition, recognition fails when appearances shift slightly. Tracking movement across frames helps certain modern methods spot deepfakes, relying on pulse patterns or memory driven systems. Though tougher to deceive, these approaches overlook details visible in individual images. When dealing with static manipulated pictures lacking temporal

indicators they falter without delay. Odd blinking behaviors break their logic just as swiftly. A different approach applies attention mechanisms to manage intricate patterns. Much like, the current work, research by Zhao and colleagues frames forgery identification as a detailed categorization task, applying a convolutional architecture enhanced with targeted attention layers to spot subtle facial

inconsistencies. Instead of merging wide-ranging feature representations, the system operates via tiny 1×1 filters that inspect single pixels individually. Because each location gets evaluated in isolation, broader contradictions stay unnoticed picture anomalies appearing genuine yet clashing with lighting cues around the distant cheekbone. Such piecemeal analysis undermines overall coherence despite accurate local observations. Across a full picture, scanning comes naturally to Vision Transformers yet strong performance demands vast data. When samples run low, outcomes slip, unlike classic CNNs that stay steady by design. Their advantage lies in their natural ability to detect even the smallest surface defects.

B. Our Approach and Contributions

A fresh arrangement sits at the heart of this approach. Instead of picking either localized filters or the full-scale focus. We define integration as its foundation. A convolutional framework shows up first and then generates fine-grained signals like hue and surface patterns. On top, the Long Distance Attention (LDA) component acts directly upon those outcomes. Division happens when the representation separates into regions, after which adaptive weights adjust relationships between them. Where similarities appear between areas, hidden mistakes start showing up. Focus shifts suddenly to places where changes look suspicious. Instead of using standard models that watch pixels closely, examination targets altered sections through transformed features. Original data gets skipped while review follows derived structures more carefully. Centered in the dual path spacetime framework lies the LDA module. One branch examines single frames, whereas the other analyzes complete movement sequences. Data merges within a common LDA backbone. By combining outputs, static views meet motion trends via cohesive evaluation. When heavy compression affects video input, how does your system respond mixing fine elements into wider contexts without clear separation? This design is motivated by the observation that spatial artifacts and temporal inconsistencies co-occur in deepfake videos and are complementary sources of evidence that, when jointly exploited, yields substantially better detection performance than either alone.

The main contributions of this work are as follows:

1) **Long Distance Attention Mechanism:** We propose a novel patch-level non-local self-attention module that, in a single forward pass, computes pairwise relevance between all N patches of a deep feature map. We provide a formal proof that standard $1k$ convolutional attention is a special constrained case of LDA, establishing LDA as a strict generalization. We further prove that the LDA attention score matrix subsumes the global average pooling operation as the uniform attention limiting case,

demonstrating that LDA can gracefully degrade to global pooling when all patches are equally informative.

2) **Dual-Branch Spatial-Temporal Framework:** We introduce the DBST network that jointly processes per-frame spatial evidence (spatial branch) and multi-frame temporal evidence (temporal branch) using the LDA mechanism as the core module in both branches. The temporal branch extends LDA to the spatio-temporal domain by treating each (t, i) patch-frame token jointly, enabling direct attention between spatially and temporally distant tokens without the sequential bottleneck of recurrent models.

3) **Theoretical and Empirical Justification of Long-Range Attention for Forgery Detection:** We provide an analysis, supported by ablation experiments, showing that the *effective attention span*—the spatial distance between the two most strongly attending patches—is a monotonically increasing function of detection accuracy. This provides a principled, quantitative argument for why long-distance attention is superior to both short-range convolutional attention and recurrent temporal aggregation for fine-grained forgery classification.

This paper follows a standard structure. Section II looks at previous studies in the field. Section III explains the theory and why this research matters. Section IV describes about the LDA mechanism and the DBST architecture. Section V details how we set up our experiments and perform. Section VI tells about our findings and what they mean. Section VII describes about the results, where the study falls short and ethics. Section VIII gives the conclusion with ideas for where this work goes next.

II. RELATED WORK

A. Deepfake Generation: A Brief Taxonomy

To see why spotting fakes is so difficult. We have to look at how they are made. [1] introduced GANs, a system where a generator G and a discriminator D train against each other. They will keep going until G This process makes the final result hard to distinguish from the truth. The progressive training strategy of Karras et al. [3] allowed GAN training to scale to 1024×1024 resolution, dramatically raising the visual quality of synthesized faces. Kingma and Welling [2] introduced the Variational Autoencoders (VAEs), the backbone of many early face-swap pipelines. More recent methods such as StyleGAN2 [4], SimSwap [5], FSGAN [6] produces the near-perfect face swaps that are imperceptible to the naked eye. The implication for detection is that generation artifacts are becoming increasingly subtle, and a detector that relies on any fixed artifact signature is bound to become obsolete as generation quality improves. This motivates attention-guided, generalizable approaches.

B. Spatial-Domain Deepfake Detection

Maternal et al. [10] were among the first to systematically catalogue the visual artifact fingerprints of deepfake generators—irregular eye color, skin-toned discontinuities and geometric misalignments—achieving AUC up to 0.866 with classical feature extractors. Afchar et al. [11] proposed MesoNet,

a shallow CNN focusing on mesoscopic image properties (mid-level patterns between pixel-level noise and semantic content), achieving $>98\%$ detection on early Deepfake and Face2Face datasets. Yang et al. [12] exploited facial landmark inconsistencies to flag GAN-generated faces. Zhou et al. [13] fused an RGB stream with a steganalysis feature stream, while Bayar and Stamm [14] proposed a constrained convolutional layer that suppresses

image content to focus on manipulation residuals. Li et al. [15] introduced Face X-Ray, which learns to detect the blending boundary present in any face composite, achieving strong cross-dataset generalization. Dang et al. [16] A different path merged when attention maps guided learning through a CNN structure. Supervision took the shape not by labels alone, but via spatial focus cues shaped during the training. The backbone adapted as the directional signals has highlighted relevant regions. Guidance came from the focused areas, not just the final outputs. Learning shifted where the emphasis was placed across the images. Attention steered the process while convolution layers extracted the features. This mix is defined as the method's foundation. One of the frequent shortcomings of strictly spatial methods is how much they depend on visible hints. A single frame appears, yet soon gives way to follow moments. Hidden within these shifts lies a story of fragments out of sync, missed without warning. These are marked by time's false impressions.

C. Attention Mechanisms and Transformers for Vision

Examining about the topic, Dang et al. (2020) combined the qualitative insights with numerical data. Their approach relied on the methods that captured both the depth and measurable patterns. While one of the part looked at the meanings people assigned and another counted recurring trends. Instead of choosing just one path, they have merged styles to strengthen the findings. Because of that each of the method compensated for the other's limits, results and gained balance. This way the interpretation rested on more than numbers alone. This self-attention mechanism [24] computes for each of the A piece within the chain appears never by itself, yet formed through varied contributions carrying the distinct importance. When any location alters it influence shifts too and adjust show the parts connect as the moments pass. Even if bonds change intensity, the core structure remains constant throughout weights are determined by pairwise relevance scores. This enables $O(1)$ -hop information flow between any two elements, regardless of distance. Dosovitskiy et al. [25] applied pure self-attention to image patches (ViT), demonstrating strong image classification performance when pre-trained on large corpora. Non-Local Neural Networks [26] introduced spatial non-local means blocks into CNNs for video understanding, with particular effectiveness for capturing long-ranges spatiotemporal dependencies. Squeeze and Excitation Networks [27] has introduced the channelwise attention, recalibrating feature

for forensic feature maps rather than semantic recognition, operating at the patch level rather than the pixel or channel level.

D. Attention-Based Deepfake Detection

Zhao et al. [29] reformulated deepfake detection as fine-grained classification and proposed MACNN, using multiple 11 convolutional attention heads to highlight facial discriminative regions. While conceptually aligned with our work, MACNN's attention has zero receptive field and cannot model non-local patch relationships. Coccomini et al. [30] combined EfficientNet with ViT by feeding CNN features as token sequences to a transformer encoder, showing that hybrid architectures outperform pure ViT on small forensic datasets. Wodajo and Atnafu [31] applied a convolutional ViT directly to video frames. Guo et al. [32] proposed a local-sensitive deepfake attention mechanism focusing on high-frequency artifacts. Our work differs from all of the above in that we: (a) provide a formal theoretical analysis of LDA as a generalization of convolutional attention; (b) integrate LDA into a dual-branch spatial-temporal design; and (c) extend LDA to spatio-temporal tokens in the temporal branch, whereas prior transformer-based methods treat only spatial patches.

III. THEORETICAL BACKGROUND AND MOTIVATION

A. Why Deepfake Detection is a Fine-Grained Problem

Definition 1 (Fine-Grained Classification). A classification problem is fine-grained if the inter-class visual difference is smaller than the intra-class visual variation, necessitating localized, discriminative feature extraction guided by global context.

As deepfake generation quality improves, the visual gap between real and fake faces shrinks: modern GANs [4] produce faces that are globally photorealistic. Defects, when they exist, are confined to small spatial regions and are often only discriminative in the context of their surroundings. A single artifact-free region is uninformative; it becomes discriminative only when contrasted with a globally inconsistent context. This is precisely the scenario that fine-grained recognition methods are designed for and it motivates the adoption of attention mechanisms that combine global semantic context with local feature discriminability.

B. Limitations of Convolutional Attention

Let $\mathbf{F}^{C \times H \times W}$ be a deep feature map. A standard kk convolutional attention module produces an attention map $\mathbf{M}^{H \times W}$ by applying a kk linear filter followed by a sigmoid:

$$\mathbf{M}_{h,w} = \sigma \left(\sum_r \mathbf{W}_{c,i,j} \mathbf{F}_{c,h+i,w+j} + b \right), \quad (1)$$

channels by modeling the interdependencies between them. The Dual Attention Network [28] has introduced the position attention and the channel attention modules for the scene segmentation. Our LDA mechanism draws inspiration from non-local means [26] and ViT [25] but is specifically designed

$$c=1, i=r, j=-r$$

The critical observation is that the attention score at location (h, w) is a function only of the local neighborhood $\mathbf{F}_{h \pm r, w \pm r}$. Two locations (h_1, w_1) and (h_2, w_2) are independent in \mathbf{M} if $|h_1 - h_2| > 2r$ or $|w_1 - w_2| > 2r$.

1 1

Consequently, if a forgery artifact at (h, w) is only detectable when contrasted against a reference region at (h_2, w_2) far away—as is the case for mismatched eye colors, uncoordinated facial texture gradients, or lighting inconsistencies—convolutional attention with small k s is structurally incapable of detecting it.

We partition \mathbf{F} into $N = \frac{H}{p} \times \frac{W}{p}$ non-overlapping spatial patches of size $p \times p$ pixels (in feature-map coordinates), yielding patch tensors $\rho_1, \rho_2, \dots, \rho_N$ with $\rho_i \in \mathbb{R}^{C \times p \times p}$. Each patch is linearly projected to a d -dimensional embedding:

$$\mathbf{e}_i = \mathbf{W}_e \text{vec}(\rho_i) + \mathbf{b}_e, \mathbf{e}_i \in \mathbb{R}^d, i = 1, \dots, N, (2)$$

where $\mathbf{W}_e \in \mathbb{R}^{d \times (C \times p^2)}$ and $\mathbf{b}_e \in \mathbb{R}^d$ are learnable. The

Proposition 1. Standard convolutional attention (Eq. (1))

is a special case of LDA (Eq. (??)) in which all cross-patch attention weights between non-adjacent patches are constrained to zero.

Proof sketch. Partition \mathbf{F} into patches of size pp ($p=1$ for pixel-level granularity). In LDA, the attention weight \mathbf{A}_{ij}

between patch i and patch j is unrestricted. In convolutional attention, the equivalent weight is non-zero only when patch j falls within the r -neighborhood of patch i , i.e., $\mathbf{A}_{ij} = 0$ whenever the spatial distance between patch centers $d(i, j) > r$. Setting $\mathbf{A}_{ij} = 0$ for $d(i, j) > r$ and restricting the remaining weights to be translation-invariant (shared across all i) recovers Eq. (1). Since LDA imposes neither of these constraints, it is a strict generalization.

D. Why LDA Outperforms LSTM-Based Temporal Modeling

Let $\{f_1, \dots, f_T\}$ be a frame sequence. An LSTM processes this sequence recurrently: the hidden state h_t encodes a summary of $\{f_1, \dots, f_t\}$. The attention between frame f_t and frame $f_{t'}$ (where $t' < t$) passes through $t - t'$ recurrent transitions, each of which introduces a multiplicative factor of the gradient, leading to exponential vanishing for large $|t - t'|$ [41]. In the temporal LDA formulation (Section IV-D), the attention weight between a token at (t_1, patch_i) and a token at (t_2, patch_j) is computed directly as a single scaled dot-product, regardless of $|t_1 - t_2|$. This means that long-range temporal dependencies receive the same gradient flow as short-range ones, resolving the vanishing gradient pathology at the cost of quadratic attention complexity.

into a spatial score grid $\mathbf{M} \in \mathbb{R}^{(H/p) \times (W/p)}$ and upsampled IV. PROPOSED METHODOLOGY

A. Problem Formulation

Let $V = \{f_1, f_2, \dots, f_T\}$ be a video of T frames, with each frame $f_t \in \mathbb{R}^{H \times W \times 3}$. A face detector D (MTCNN [38]) in our

embedding matrix is $\mathbf{E} = [\mathbf{e}_1, \dots, \mathbf{e}_N] \in \mathbb{R}^{N \times d}$. We project \mathbf{E} into queries, keys and values:

$$\mathbf{Q} = \mathbf{E} \mathbf{W}_Q, \mathbf{K} = \mathbf{E} \mathbf{W}_K, \mathbf{V} = \mathbf{E} \mathbf{W}_V, (3)$$

where $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_k}$ and $\mathbf{W}_V \in \mathbb{R}^{d \times d_v}$. The LDA weight matrix is:

$$\mathbf{A} = \text{softmax}_k \left(\frac{\mathbf{Q} \mathbf{K}^T}{d} \right) \in \mathbb{R}^{N \times N}. (4)$$

Entry $\mathbf{A}_{ij}[0, 1]$ quantifies the global relevance of patch j to patch i , regardless of their spatial separation. The output context matrix is:

$$\mathbf{Z} = \mathbf{A} \mathbf{V} \in \mathbb{R}^{N \times d_v}. (5)$$

To enhance expressiveness, we use H attention heads. For head h :

$$\mathbf{A}^{(h)} = \text{softmax} \left(\frac{\mathbf{Q}^{(h)} \mathbf{K}^{(h)T}}{d_k/H} \right), (6)$$

$$\mathbf{Z}^{(h)} = \mathbf{A}^{(h)} \mathbf{V}^{(h)}, (7)$$

and the multi-head output is $\mathbf{Z} = \text{Concat}(\mathbf{Z}^{(1)}, \dots, \mathbf{Z}^{(H)}) \mathbf{W}$,

where $\mathbf{W} \in \mathbb{R}^{Hd \times d}$.

3) Attention Map Generation and Feature Recalibration: The forgery-importance score of each patch is computed as the row-sum of the attention matrix:

$$s_i = \sum_{j=1}^N \mathbf{A}_{ij}, i = 1, \dots, N. (8)$$

Intuitively, s_i measures how much information patch i aggregates from the entire feature map, serving as a proxy for its discriminative importance. The scores $\{s_i\}_i$ are reshaped

implementation) produces aligned face crops $\mathbf{x}_t \in \mathbb{R}^{h \times w \times 3}$. Define $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ as the face clip. The detection

to the feature map resolution via bilinear interpolation to yield the

LDA attention map $\mathbf{M} \in \mathbb{R}^{H' \times W'}$.

Feature recalibration is then performed by element-wise

multiplication:

$$\tilde{\mathbf{F}} = \mathbf{F} \odot \mathbf{M}, \quad (9)$$

where \odot denotes channel-wise broadcast multiplication. The

task to learn a mapping $F: X \rightarrow [0, 1]$, where the

recalibrated feature map \tilde{F} suppresses non-discriminative

output is the probability that the clip is a deepfake and Θ denotes all learnable parameters. We treat this as a fine-grained classification problem; the decision boundary is not defined by coarse global appearance but by subtle, contextually dependent local artifacts.

B. Long Distance Attention Mechanism

1) *Patch Embedding*: Let $F \in \mathbb{R}^{C \times H' \times W'}$ be the feature map produced by a convolutional backbone applied to a

regions and up-weights pivotal forgery patches, guided entirely by long-range global context.

4) *Residual Fusion*: To avoid information loss from poorly calibrated attention maps early in training, we fuse the recalibrated and original features via a learnable residual:

$$\hat{F} = \lambda \tilde{F} + (1 - \lambda) F, \quad (10)$$

where $\lambda \in [0, 1]$ is a learnable scalar initialized to 0.5, allowing the network to adaptively blend global attention guidance with the original local features.

C. Spatial Branch

The spatial branch processes individual frames to capture intra-frame forgery artifacts. The pipeline is as follows:

- 1) A face crop x_t is passed through a CNN backbone (EfficientNet-B4 [37] pre-trained on ImageNet) to produce $F^{(s)} \in \mathbb{R}^{C \times H' \times W'}$.
- 2) The LDABlock (Section IV-B) produces $\hat{F}^{(s)}$ via patch embedding, attention computation and residual fusion.
- 3) A spatial classification head consisting of global average pooling (GAP) \rightarrow Layer Normalization \rightarrow FC(512) \rightarrow ReLU \rightarrow Dropout(0.5) \rightarrow FC(1) \rightarrow Sigmoid produces

2) *Multi-Task Training Objective*: The model is trained jointly on three predictions:

$$L = L_{BCE}(y^g, y) + L_{BCE}(y^s, y) + L_{BCE}(y^t, y), \quad (15)$$

where $L_{BCE}(p, y) = -y \log(p) - (1-y) \log(1-p)$ and $y \in \{-1, 1\}$ is the ground truth label. The multi-task formulation forces both branches to independently learn discriminative representations, preventing the fusion layer from learning to simply suppress the weaker branch.

F. Algorithm

Algorithm 1 DBST-LDA Training and Inference

Require: Face clip $X = \{x_1, \dots, x_T\}$; label y ; trained model parameters Θ (inference only); learning rate η
Ensure: Prediction \hat{y} ; updated Θ (training only)

- 1: */*Preprocessing*/*
- 2: Detect & align faces using MTCNN; resize to 224 \times 224.
- 3: */*Spatial Branch*/*
- 4: **for** $t = 1$ **to** T **do**

the frame-level forgery probability $y^{(s)} \in [0, 1]$.

The video-level spatial prediction is the mean over all T frames:

$$\hat{y}^{(s)} = \frac{1}{T} \sum_{t=1}^T y^{(s)}. \quad (11)$$

D. Temporal Branch

The temporal branch captures inter-frame inconsistencies by jointly processing a clip of T consecutive face crops.
1) *Spatio-Temporal Token Sequence*: Each frame x_t is encoded by the same CNN backbone (shared weights with the spatial branch) to obtain $\mathbf{F}_t \in \mathbb{R}^{C \times H \times W}$. Each frame's feature map is partitioned into N spatial patches, yielding TN total patch tokens. The (t, i) -th token is embedded as:

$$\mathbf{e}_{t,i} = \mathbf{W}_e \text{vec}(p_i^{(t)}) + \mathbf{b}_e + \mathbf{p}_t + \mathbf{q}_i \quad (12)$$

where $\mathbf{p}_t \in \mathbb{R}^d$ is a learnable temporal position encoding and $\mathbf{q}_i \in \mathbb{R}^d$ is a learnable spatial position encoding.

2) *Spatio-Temporal LDA*: The TN tokens are organized into matrix $\mathbf{E}^{(ST)} \in \mathbb{R}^{TN \times d}$. The spatio-temporal LDA attention matrix is:

$$\mathbf{A}^{(ST)} = \text{softmax} \left(\frac{\mathbf{Q}^{(ST)} \mathbf{K}^{(ST)T}}{\sqrt{d_k}} \right) \in \mathbb{R}^{TN \times TN}, \quad (13)$$

where projections are defined analogously to Eq. (3). Entry $\mathbf{A}_{(t_1, i_1)(t_2, j)}$ encodes the relevance of spatial patch j at frame t_2 to spatial patch i at frame t_1 , simultaneously capturing both intra-frame spatial attention and cross-frame temporal attention in a single matrix. The output $\mathbf{Z}^{(ST)} \in \mathbb{R}^{TN \times d_v}$ is temporally pooled (mean over the T dimension) and passed through the same classification head architecture as the spatial branch to yield clip-level probability $\hat{y}^{(t)}$.

E. Fusion and Training

1) *Adaptive Fusion*: The final prediction combines both branches:

$$\hat{y} = \alpha \hat{y}^{(s)} + (1 - \alpha) \hat{y}^{(t)}, \quad (14)$$

where $\alpha \in [0, 1]$ is a learnable scalar initialized to 0.5.

$$\begin{aligned} 5: & \mathbf{F}_t^{(s)} \leftarrow \text{CNN}_{\text{backbone}}(x_t) \\ 6: & \{\mathbf{e}_i\}_{i=1}^N \leftarrow \text{PatchEmbed}(\mathbf{F}_t^{(s)}) \quad // \text{Eq. (2)} \\ 7: & \mathbf{A} \leftarrow \text{ScaledDotProd}(\mathbf{Q}, \mathbf{K}) \quad // \text{Eq. (4)} \end{aligned}$$

$$\begin{aligned} & \hat{y}_t^{(s)} \leftarrow \text{ClassHead}(\tilde{\mathbf{F}}^{(s)}) \\ 10: & \text{endfor} \\ 11: & \hat{y}^{(s)} \leftarrow \frac{1}{T} \sum_t \hat{y}_t^{(s)} \quad // \text{Eq. (11)} \\ 12: & /*TemporalBranch*/ \end{aligned}$$

$$13: \mathbf{E}^{(ST)} \leftarrow \text{ST-PatchEmbed}(x_1, \dots, x_T) \quad // \text{Eq. (12)}$$

$$14: \mathbf{A}^{(ST)} \leftarrow \text{ScaledDotProd}(\mathbf{Q}^{(ST)}, \mathbf{K}^{(ST)}) \quad // \text{Eq. (13)}$$

$$15: \hat{y}^{(t)} \leftarrow \text{ClassHead}(\text{TemporalPool}(\mathbf{Z}^{(ST)}))$$

16: /*Fusion*/

$$17: \hat{y} \leftarrow \alpha \hat{y}^{(s)} + (1 - \alpha) \hat{y}^{(t)} \quad // \text{Eq. (14)}$$

18: **if** training **then**

$$19: \quad \mathbf{L} \leftarrow \text{Eq. (15)}$$

$$20: \quad \Theta \leftarrow \Theta - \eta \nabla_{\Theta} \mathbf{L}$$

21: **endif**

22: **return** \hat{y}

V. EXPERIMENTAL SETUP

A. Datasets

We evaluate the proposed method on three widely used benchmark datasets selected to cover different levels of difficulty, manipulation diversity, and dataset scale.

FaceForensics++ (FF++) [34]: This dataset contains a 1,000 original YouTube videos along with 4,000 manipulated videos which has been generated using four face manipulation methods. The methods include DeepFakes (DF), Face2Face (F2F), FaceSwap (FS), and NeuralTextures (NT). The videos are provided at three compression settings: raw (uncompressed), light compression (c23), and heavy compression (c40). FF++ is commonly used as a standard benchmark for evaluating in-distribution deepfake detection performance. We report results at c23 and c40 to assess robustness to compression.

Celeb-DF v2 [35]: A more challenging dataset containing 590 realcelebrityinterviewvideosand5,639high-qualitydeepfake videos synthesized with a significantly improved pipeline over first-generation tools. The higher visual quality makes Celeb-DF an important test of generalization beyond FF++.

Deepfake Detection Challenge (DFDC) [36]: The largest publicly available deepfake dataset (>100,000 videos), featuring diverse actors, lighting conditions, backgroundsandmultiplecutting-edgegenerationmethods.DFDCwasexplicitly designed to challenge overfitted detectors and is the closest proxy for real-world deployment conditions.

B. ImplementationDetails

Facepreprocessing: FacesaredetectedusingMTCNN[38], aligned to 5 canonical landmarksand cropped to 224x224. For FF++ we use 32 frames per video sampled uniformly; for Celeb-DF and DFDC we use 16 frames due to variable video length.

CNNbackbone: EfficientNet-B4[37]pre-trainedonImageNet; thelastclassificationlayerisremoved.Theoutputfeaturemap is $C = 1792$ channels at $H' = W' = 7$.

LDA hyperparameters: Patch size $p = 1$ (treating each spatial location as a patch, yielding $N = 49$ tokens per frame); embedding dimension $d = 512$; key/query dimension $d_k = 64$; value dimension $d_v = 64$; $H = 8$ attention heads; $d_k/H = d_v/H = 64/8 = 8$ perhead.

Temporalbranch: Clip length $T = 8$ frames, giving $TN = 392$ spatio-temporal tokens. To manage the quadratic complexity $O((TN)^2) = O(392^2) = 153k$ attention entries perclip, we apply FlashAttention[40] for memory-efficient computation.

Optimization: AdamW[39] with learning rate 10^{-4} , weightdecay 5×10^{-4} , batch size 32, 30 training epochsand cosineannealing with warm restarts. Backbone is frozen for the first5 epochs then fine-tuned with a $10 \times$ reduced learning rate.

Data augmentation: Random horizontal flip, random JPEG compression(quality50–100), Gaussianblur($\sigma: 0–1.5$), brightness/contrast jitter (±5%)and random erasing.

Hardware: Two NVIDIA A100 80GB GPUs; total training time approximately 14 hours per model on FF++.

C. BaselineMethods

We compare the proposed DBST-LDA against the following baselines spanning all three detection families:

Spatial-domain: MesoNet[11], Two-Stream[13], FaceX-Ray [15], EfficientNet-B4 (fine-tuned, no attention module).

Temporal-domain: LipForensics[23], recurrentCNN(ResNet-18 + GRU) [22].

Attention/Transformer: MACNN[29], ViT-B/16fine-tuned[25], EfficientNet+ViThybrid[30].

D. EvaluationMetrics

Following standardpractice inthe deepfake detectionliterature, we report:

- **Accuracy(Acc):** overallcorrectclassificationrateat

- **Precision(P)**: $TP/(TP+FP)$.
- **Recall(R)**: $TP/(TP+FN)$.
- **F1-Score**: harmonic mean of P and R.
- **AUC**: area under the ROC curve (threshold-free).

For cross-dataset experiments, only AUC is reported to allow fair comparison with published results.

VI. RESULTS AND ANALYSIS

A. In-Distribution Performance on Face Forensics++

Table I presents frame-level results on FF++ (c23 compression).

TABLE I
 IN-DISTRIBUTION PERFORMANCE ON FF++ (c23). FRAME-LEVEL EVALUATION.

Method	Year	Acc(%)	P(%)	R(%)	F1(%)	AUC(%)
MesoNet[11]	2018	83.1	82.8	83.5	83.1	85.2
Two-Stream[13]	2017	86.9	85.4	87.6	86.4	88.0
EfficientNet-B4(noattn)	–	92.4	91.9	92.8	92.3	94.8
FaceX-Ray[15]	2020	91.3	90.5	91.9	91.2	93.5
MACNN[29]	2021	93.1	92.7	93.4	93.0	95.4
LipForensics[23]	2021	90.5	89.8	91.1	90.4	92.7
ViT-B/16[25]	2021	91.7	90.9	92.3	91.6	94.0
Effic+ViT[30]	2022	93.6	93.1	94.0	93.5	95.9
DBST-LDA(Ours)	–	95.6	95.2	95.9	95.5	97.2

The proposed DBST-LDA model achieves the highest scores across all five metrics, with an AUC of 97.2%—a gain of 1.8% over the strongest baseline (Effic+ViT, 95.9%) and 1.3% over MACNN (95.4%). The margin over the backbone-only EfficientNet-B4 (94.8%) quantifies the contribution of the LDA attention mechanism in isolation: the 2.4% AUC improvement is attributable entirely to the attention module.

Table II presents outcomes observed with intense c40 compression applied. Here, challenges grow more pronounced due to the frequent presence of such factors in actual scenarios. Most platforms alter video content significantly once it's uploaded.

TABLE II
 IN-DISTRIBUTION PERFORMANCE ON FF++ (c40). FRAME-LEVEL AUC(%)

Method	AUC(%)
MesoNet[11]	70.3
FaceX-Ray[15]	79.5
MACNN[29]	83.8
LipForensics[23]	82.4
Effic+ViT[30]	86.1
DBST-LDA(Ours)	88.4

Even if it is much smaller now, DBST-LDA remains as the part of the ranking order. Its position persists as the despite of the shrinkage. It hasn't vanished with size of the reduction. It has still present in the line, though it scaled down significantly. This sequence includes it, even when pared back sharply. Though hitting 88.4% AUC, this approach just out

perform the next best at 86.1%. Yet oddly enough, distinctions remain clear even when elements are far apart in the sequence. Thoughts squeezed tight, subtle traces still show a result of slow, steady filtering. Fading happens, but intricacy remains hidden under polished surfaces. Clarity slips away, though delicate differences stay present in the gaps among dots. Reduced size means less accuracy; even so, gentle rhythms emerge softly after multiple passes. Occasionally, devices designed to detect spatial characteristics rely on

B. Cross-Dataset Generalization

Table III displays AUC values in percent (%), measured through evaluation on various data collections: models trained on FF++(c23) evaluated using Celeb-DF and DFDC method gets tested without changes applied. Evaluation here stands as the most demanding type of review detector generalization.

TABLE III
 CROSS-DATASET GENERALIZATION EVALUATED ON CELEB-DFV2 AND DFDC. AUC(%).

Method	Celeb-DFv2	DFDC
MesoNet[11]	54.8	55.9
Two-Stream[13]	61.4	61.0
FaceX-Ray[15]	74.2	70.5
MACNN[29]	79.6	75.1
LipForensics[23]	82.4	73.5
ViT-B/16[25]	77.3	72.8
Effic+ViT[30]	83.1	77.4
DBST-LDA(Ours)	85.7	79.3

DBST-LDA improves cross-dataset Celeb-DF AUC by 2.6% over Effic+ViT and by 6.1% over MACNN. overall improvement through the DFDC is 1.9% over Effic+ViT. These results suggest through long-distance patch attention learns more universal forgery manners indicators that transfer better overall across generation methods than artifact-specific spatial features.

C. Ablation Study

Table IV systematically contribution of each component by their progressive construction on Seeds.

TABLE IV
 ABLATION STUDY ON FF++(c23). AUC(%).

Configuration	FF++	Celeb-DF	DFDC
Backbone only (no attention)	94.8	73.0	66.4
+Local conv. attn. (1x1)	95.0	73.3	66.7
+Local conv. attn. (3x3)	95.2	73.9	67.1
+LDA, N=16 (large patches)	95.6	76.1	68.9
+LDA, N=49 (medium patches)	96.5	80.2	72.4
+LDA, N=196 (fine-grained)	96.8	82.1	74.1
+Temporal branch (LDA only)	96.3	81.5	73.8
Full DBST-LDA (spatial+temporal)	97.2	85.7	79.3

Several key insights emerge from Table IV:

(1)

(2) **Local attention provides marginal gain.** Adding 1

or 33 Looking closely at the test results, three distinct trends come into view. When different kinds of input are applied, subtle tendencies start to surface. In cases where data lacks clarity, a drop in effectiveness shows up each time. Looking at the outcomes explains why local attention has minimal impact. +0.2–0.4% Gains on new, unseen data hardly shifted regardless of those modifications. Limited receptive fields cannot grasp relevant structures effectively. What truly counts is recognizing relationships across distant parts of an image. For small filters, context remains out of reach. Looking far apart reveals more telling signs than what's right next to you. Improvement in spotting patterns begins when attention moves outward instead of staying near.

(3) **LDA gain increases with finer patch granularity.** Though tiny patches boost LDA performance, gain emerges only beyond a threshold. As numbers rise $N=16$ to $N=196$ results climb step by step across all data sets. Finer facial divisions make small defects stand out within tight zones. When regions stretch wider, subtle mismatches begin to merge. These faint signs may vanish entirely if space expands too much.

(4) **Both branches contribute.** Ultimately, integrating temporal patterns with spatial data turns out critical. 96.8% FF++ AUC and 82.1% Celeb-DFAUC, since depending solely on location-based traits. Still, outcomes improved noticeably 96.3% and 81.5%, respectively. The full dual-branch model (97.2% and 85.7%) after fusing both streams. What stands out is how odd visuals and stiff motion point in different directions. When combined, these signs make spotting fakes easier - more so than either clue by itself.

D. Analysis of Attention Span vs. Detection Accuracy

Surprisingly, the gap between top-two active patches revealed patterns during validation. From this point on, wide focus seemed tied to better counterfeit spotting. Using the 90th percentile carved out a clear boundary - this line then defined how far attention was allowed to spread. Afterward came six model versions, none able to attend past fixed limits. Each built under tighter constraints, yet followed the same core design. From a point beyond distance r , removing connections showed clearly the reliance on far-off information. Each version's results online measured using FF++ AUC. As accuracy rises, interaction between more separated regions is allowed. Better detection emerges when widely spaced parts of faces connect. AUC rises from 94.8% at $r=0$ (degenerate, equivalent to their attention) to 97.2% at $r=(full\ LDA)$. Every additional 2-unit increase in r yields a sign through the significant AUC improvement ($p < 0.05$, paired t-test all over 5 seeds).

E. Qualitative Attention Map Analysis

The bright zones tend to cluster around the facial edges: jaws, cheeks and necks when it is spotting synthetic layouts. These patterns appear repeatedly in the LDA attention maps across the different samples. If there is strong attention concentrated around the eye boundaries and nostril regions then there will be inconsistencies such as uneven blending and texture discontinuities are more likely to occur. In real videos, visual

information is usually spread more naturally across the frame. The manipulated content tends to produce the localized attention near the facial edges.

Abrupt motion variations are also noticeable in the forged sequences. Even if there is some of the distant frames appear visually sharp, nearby temporal inconsistencies can still reveal the signs of the manipulation. The proposed approach learns these through structural relationships directly from the data without relying on manually defined rules. Unlike the MACNN which may lose global context due to the fragmented feature extraction. The proposed method preserves the broader spatial relationships. It also remains as more stable than the ViT-based approaches in situations which involves the reflections, lighting variation, or noisy textures.

VII. DISCUSSION

A. Why LDA Outperforms CNN and LSTM: A Summary

The experimental results that incorporate the theoretical analysis of the Section III. To summarize the argument concretely: **vs. CNN:** A kk conversion towards the attention module can only internally compare features themselves within a kk window. Stacking L such layers expands the effective receptive field to $O(kL)$, but: (a) the dependence between distant locations is *indirect*, mediated through a chain of intermediate layers; (b) each layer applies the same translation-invariant filter to all locations, unable to compute *location-pair-specific* relevance. LDA computes explicit pairwise relevance between all M patches in a *single* attention step, enabling $O(1)$ -depth global reasoning.

vs. LSTM: Recurrent models suffer from: (a) sequential processing that prevents parallelism; (b) vanishing gradients for long-range dependencies (Section III); (c) inability to capture *spatial* non-local dependencies within frames; (d) sequential bottleneck that cannot attend directly from frame t_1 to frame t_2 without passing through all intermediate hidden states. LDA processes all TN tokens jointly and in parallel, with direct gradient flow between any two tokens regardless of their spatio-temporal distance.

B. Limitations

Like most of the research systems, the proposed framework also has a few practical limitations.

Quadratic attention complexity: The main computational cost comes from the full cross-patch attention mechanism. The spatio-temporal LDA matrix scales as $O((TN)^2)$. The current setup with $T=8$ frames and $N=49$ patches per frame. The model processes around 153k attention entries per clip, which remains manageable using FlashAttention[40]. Increasing the clip length significantly raises memory usage. For example using $T=32$ results in more than the 2.4 million entries which makes the computation much more demanding. Future work may explore the linear attention methods [42] or sparse the attention strategies to reduce this overhead while maintaining the forensic accuracy.

Sensitivity to compression: Strong c40 compression significantly

cantly reduce the overall performance. This can lower the AUC from 97.2% to 88.4% on FF++. Although the proposed method

still performs better than all evaluated baselines. The AUC of 88.4% may not be sufficient for high-risk applications such as the large-scale content moderation. The legal verification or the media authentication. The main challenge is that aggressive video compression removes many of the fine texture details. The attention mechanism relies on the detecting manipulations. Once this information is lost during the encoding, it cannot be fully recovered by the model.

Adversarial deepfakes: The results reported in this work were obtained using the videos created without targeting the proposed detector. In the practical scenarios, this assumption may not always hold. In the previous studies include the work by Neekhara [43] have shown that the attackers can generate manipulated videos which are designed specifically to bypass the detection systems. The current study does not include the adversarial robustness experiments, which remains an important limitation. A determined attacker with the access to the detector's behaviour or the architecture could potentially reduce the detection of the performance beyond the results shown in Table I. Incorporating these adversarial examples during the training may help to improve the robustness against such attacks in future work.

Scope of manipulation: The proposed pipeline focuses on a specific type of manipulation. This consists of fake human face inserted into real video frame. This represents a major category of the deepfake content, several other forms of manipulation which are not addressed in the current work. These include voice cloning with minimal facial motion, full-body motion transfer and audio-driven lip synchronisation techniques. Limiting the study to facial forgery detection helped keep the problem focused and manageable. The system should not be considered a complete solution for all the types of deepfake media.

C. Ethical Considerations

The development and deployment of the deepfake detection systems also raise the important ethical concerns. These issues influence the public trust, personal reputation and the responsible use of the AI technologies. Discussing them openly is necessary to understand their broader social impact.

False-positive harm: An incorrect deepfake prediction can seriously affect an individual's reputation. A genuine video that is mistakenly flagged as manipulated may create suspicion even after the error is corrected. In many cases, public doubt spreads faster than the clarification and the damage can persist long after the original claim is withdrawn. In such situations may also affect the journalism, legal investigations and public discourse, where the inaccurate conclusions can lead to unfair consequences.

For this reason, deploying DBST-LDA in real-world settings requires careful calibration and reliable uncertainty estimation methods, such as temperature scaling or Bayesian posterior analysis. In sensitive applications, model predictions should always be reviewed by human experts before final decisions are made.

What a detector fails to capture is the influences of how

findings are seen. When limitations are known, the trust

grows. Under specific circumstances, errors emerge more often.

Awareness of failure modes helps avoid misuse. Reliability depends on recognizing where function ends.

Arms-race dynamics: Publishing a detection architecture. Within this section, every detail appears exactly as open practices usually require. Because of the demands reveal what traits matter most is the rivals start noticing patterns. What one prioritizes is the signals weaknesses others might exploit. Focus shifts when goals highlight specific qualities. Attention

follows where the effort is spent. Clues emerge through choices which are made under pressure. Move along the sensor's edge. We find this acceptable because: (a) long-term societal benefit of a community-reviewed, reproducible detector. What remains significant outweighs fleeting shortcomings; our observations consist of Some of these components differ too much to be uniformly reduced without affecting how sharp the image appears.

Privacy and consent: Training datasets (FF++, Celeb-DF, DFDC) involve real individuals' face videos. We have used only publicly released benchmarks with documented institutional review board approvals. Any downstream deployment that involves collecting new face videos for detector retraining must obtain explicit informed consent from the subjects.

Misuse for censorship or surveillance: A deepfake detector could be co-opted to build surveillance tools or to falsely flag authentic journalist videos as fake. We strongly caution against these applications and advocate for open licensing with use-case restrictions.

Bias and fairness: Deepfake detectors trained primarily on datasets of celebrity faces may have systematically lower accuracy for demographic groups underrepresented in training data. Future work must include explicit demographic bias audits.

VIII. CONCLUSION

This paper addressed the challenge of deepfake video detection from a novel theoretical and architectural perspective. We observed that contemporary deepfakes introduce forgery artifacts that are inherently *fine-grained and non-local*: they are only discriminative when a suspicious local region is evaluated in the context of globally inconsistent surroundings. This motivated the design of the Long Distance Attention (LDA) mechanism, a patch-level non-local self-attention module that computes pairwise relevance across the *entire* feature map in a single attention step, without the locality constraints of convolutional attention or the sequential bottleneck of recurrent models. We provided a formal proof that LDA strictly generalizes local convolutional attention and a theoretical analysis of why it resolves the vanishing-gradient pathology of LSTMs for long-range temporal modeling.

The LDA mechanism was embedded in a Dual-Branch Spatial-Temporal (DBST) network that simultaneously exploits intra-frame spatial artifacts and inter-frame temporal inconsistencies—the two primary evidence streams for forgery detection—through a unified attention framework. Training was guided by a multi-task cross-entropy loss that independently

supervises both branches, ensuring each learns complementary discriminative representations.

Extensive experiments on FF++, Celeb-DF v2 and DFDC demonstrated state-of-the-art performance, with particular improvements in cross-dataset generalization (+2.6% AUC on Celeb-DF over the strongest baseline), confirming that LDA-guided features are more universal than the artifact-specific fingerprints learned by local spatial detectors. We hope this work provides a solid theoretical foundation and strong empirical baseline for these future directions.

REFERENCES

[1] I. Goodfellow *et al.*, “Generative adversarial nets,” *Adv. Neural Inf. Process. Syst.*, vol. 27, Montre’al, Canada, 2014.

[2] D. P. Kingma and M. Welling, “Auto-encoding variational Bayes,” in *Proc. Int. Conf. Learn. Represent.*, Banff, Canada, 2014.

[3] T. Karras, T. Aila, S. Laine and J. Lehtinen, “Progressive growing of GANs for improved quality, stability and variation,” in *Proc. Int. Conf. Learn. Represent.*, Vancouver, Canada, 2018.

[4] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen and T. Aila, “Analyzing and improving the image quality of StyleGAN,” in *Proc. IEEE/CVF CVPR*, Seattle, WA, USA, 2020, pp. 8107–8116.

[5] R. Chen, X. Chen, B. Ni and Y. Ge, “SimSwap: An efficient framework for high fidelity face swapping,” in *Proc. ACM Int. Conf. Multimedia*, Seattle, WA, USA, 2020, pp. 2003–2011.

[6] Y. Nirkin, Y. Keller and T. Hassner, “FSGAN: Subject agnostic faceswapping and reenactment,” in *Proc. IEEE/CVF ICCV*, Seoul, Korea, 2019, pp. 7184–7193.

[7] “Deepfakes,” <https://github.com/deepfakes/>, Accessed Sep. 2019.

[8] “FakeApp,” <http://www.fakeapp.com/>, Accessed Feb. 2020.

[9] “FaceSwap,” <https://github.com/MarekKowalski/>, Accessed Sep. 2019.

[10] F. Matern, C. Riess and M. Stamminger, “Exploiting visual artifacts to expose deepfakes and face manipulations,” in *IEEE Winter Conf. Appl. Comput. Vision Workshops*, Waikoloa, USA, 2019, pp. 83–92.

[11] D. Afchar, V. Nozick, J. Yamagishi and I. Echizen, “MesoNet: A compact facial video forgery detection network,” in *IEEE Int. Workshop Inf. Forensics Security*, Hong Kong, 2018, pp. 1–7.

[12] X. Yang, Y. Li, H. Qi and S. Lyu, “Exposing GAN-synthesized faces using landmark locations,” in *Proc. ACM Workshop Inf. Hiding Multimedia Security*, Paris, 2019, pp. 113–118.

[13] P. Zhou, X. Han, V. I. Morariu and L. S. Davis, “Two-stream neural networks for tampered face detection,” in *IEEE CVPR Workshops*, Honolulu, USA, 2017, pp. 1831–1839.

[14] B. Bayar and M. C. Stamm, “A deep learning approach to universal image manipulation detection using a new convolutional layer,” in *Proc. 4th ACM Workshop Inf. Hiding Multimedia Security*, Vigo, Spain, 2016, pp. 5–10.

[15] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen and B. Guo, “FaceX-Ray for more general face forgery detection,” in *Proc. IEEE/CVF CVPR*, Seattle, WA, USA, 2020, pp. 5001–5010.

[16] H. Dang, F. Liu, J. Stehouwer, X. Liu and A. K. Jain, “On the detection of digital face manipulation,” in *Proc. IEEE/CVF CVPR*, Seattle, WA, USA, 2020, pp. 5781–5790.

[17] U. A. Ciftci, I. Demir and L. Yin, “FakeCatcher: Detection of synthetic portrait videos using biological signals,” *IEEE Trans. Pattern Anal. Mach. Intell.*, 2020, doi: 10.1109/TPAMI.2020.3009287.

[18] Y. Li, M.-C. Chang and S. Lyu, “In Ictu Oculi: Exposing AI created fake videos by detecting eye blinking,” in *IEEE Int. Workshop Inf. Forensics Security*, Hong Kong, 2018, pp. 1–7.

[19] M. Li, B. Liu, Y. Hu and Y. Wang, “Exposing deepfake videos by tracking eye movements,” in *Proc. 25th Int. Conf. Pattern Recognit.*, Milan, Italy, 2021, pp. 5184–5189.

[20] C.-Z. Yang, J. Ma, S. Wang and A. W.-C. Liew, “Preventing deepfake attacks on speaker authentication by dynamic lip movement analysis,” *IEEE Trans. Inf. Forensics Security*, vol. 16, pp. 1841–1854, 2021.

[21] S. Fernandes *et al.*, “Predicting heart rate variations of deepfake videos using Neural ODE,” in *IEEE/CVF ICCV Workshop*, Seoul, Korea, 2019, pp. 1721–1729.

[22] E. Sabir *et al.*, “Recurrent convolutional strategies for face manipulation detection in videos,” in *Proc. IEEE/CVF CVPR Workshops*, Los Angeles, US

A, Jun. 2019.

[23] A. Haliassos, K. Vougioukas, S. Petridis and M. Pantic, “Lips don’t lie: A generalisable and robust approach to face forgery detection,” in *Proc. IEEE/CVF CVPR*, Virtual, 2021, pp. 5039–5048.

[24] A. Vaswani *et al.*, “Attention is all you need,” *Adv. Neural Inf. Process. Syst.*, vol. 30, Long Beach, CA, USA, 2017.

[25] A. Dosovitskiy *et al.*, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *Proc. Int. Conf. Learn. Represent.*, Virtual, 2021.

[26] X. Wang, R. Girshick, A. Gupta and K. He, “Non-local neural networks,” in *Proc. IEEE/CVF CVPR*, Salt Lake City, UT, USA, 2018, pp. 7794–7803.

[27] J. Hu, L. Shen and G. Sun, “Squeeze-and-excitation networks,” in *Proc. IEEE/CVF CVPR*, Salt Lake City, UT, USA, 2018, pp. 7132–7141.

[28] J. Fu *et al.*, “Dual attention network for scene segmentation,” in *Proc. IEEE/CVF CVPR*, Long Beach, CA, USA, 2019, pp. 3146–3154.

[29] J. Zhao, L. Han, L. Shao and C. Bao, “Multi-attentional deepfake detection,” in *Proc. IEEE/CVF CVPR*, Nashville, TN, USA, 2021, pp. 2185–2194.

[30] D. A. Coccomini *et al.*, “Combining EfficientNet and Vision Transformers for video deepfake detection,” in *Proc. Int. Conf. Image Anal. Process.*, Lecce, Italy, 2022, pp. 219–229.

[31] D. Wodajo and S. Atnafu, “Deepfake video detection using convolutional vision transformer,” *arXiv preprint arXiv:2102.11126*, 2021.

[32] Z. Guo *et al.*, “Controllable guide-space for generalizable face forgery detection,” in *Proc. IEEE/CVF ICCV*, Paris, France, 2023, pp. 22211–22221.

[33] J. Fu, H. Zheng and T. Mei, “Look close to see better: Recurrent attention convolutional neural network for fine-grained image recognition,” in *Proc. IEEE CVPR*, Honolulu, HI, USA, 2017, pp. 4438–4446.

[34] A. Rossler *et al.*, “FaceForensics++: Learning to detect manipulated facial images,” in *Proc. IEEE/CVF ICCV*, Seoul, Korea, 2019, pp. 1–11.

[35] Y. Li *et al.*, “Celeb-DF: A large-scale challenging dataset for deepfake forensics,” in *Proc. IEEE/CVF CVPR*, Seattle, WA, USA, 2020, pp. 3207–3216.

[36] B. Dolhansky *et al.*, “The deepfake detection challenge (DFDC) dataset,” *arXiv preprint arXiv:2006.07397*, 2020.

[37] M. Tan and Q. V. Le, “EfficientNet: Rethinking model scaling for convolutional neural networks,” in *Proc. Int. Conf. Mach. Learn.*, Long Beach, CA, USA, 2019, pp. 6105–6114.

[38] K. Zhang, Z. Zhang, Z. Li and Y. Qiao, “Joint face detection and alignment using multi-task cascaded convolutional networks,” *IEEE Signal Process. Lett.*, vol. 23, no. 10, pp. 1499–1503, 2016.

[39] I. Loshchilov and F. Hutter, “Decoupled weight decay regularization,” in *Proc. Int. Conf. Learn. Represent.*, New Orleans, LA, USA, 2019.

[40] T. Dao, D. Y. Fu, S. Ermon, A. Rudra and C. Re’, “FlashAttention: Fast and memory-efficient exact attention with IO-awareness,” *Adv. Neural Inf. Process. Syst.*, vol. 35, New Orleans, LA, USA, 2022.

[41] Y. Bengio, P. Simard and P. Frasconi, “Learning long-term dependencies with gradient descent is difficult,” *IEEE Trans. Neural Netw.*, vol. 5, no. 2, pp. 157–166, 1994.

[42] A. Katharopoulos, A. Vyas, N. Pappas and F. Fleuret, “Transformers are RNNs: Fast autoregressive transformers with linear attention,” in *Proc. Int. Conf. Mach. Learn.*, Virtual, 2020, pp. 5156–5165.

[43] P. Neekhara, B. Dolhansky, J. Bitton and C. Canton-Ferrer, “Adversarial threats to DeepFake detection: A practical perspective,” in *Proc. IEEE/CVF CVPR Workshops*, Virtual, 2021.

[44] S. Adinarayana, D. Ravikiran, U. Sesadri, “Efficient deep fake detection using a hybrid deep learning model in *Proc. IEEE Virtual*, 2025