



International Journal of Research Publication and Reviews

Journal homepage: www.ijrpr.com ISSN 2582-7421

Detection of Deepfake Videos Using Long Distance Attention

K. Madhusudhan Reddy¹, L. Divya²

¹ Assistant Professor, Dept. of MCA, Annamacharya Institute of Technology and Sciences (AITS), Karakambadi, Tirupati, Andhra Pradesh, India
EMAIL: mca.madhureddy@gmail.com

² Student, Dept. of MCA, Annamacharya Institute of Technology and Sciences (AITS), Karakambadi, Tirupati, Andhra Pradesh, India
EMAIL : lokanadhamdivya2@gmail.com

ABSTRACT

The increasing sophistication of deepfake generation techniques poses significant challenges to digital media authentication. Traditional detection methods struggle to identify subtle inconsistencies in facial expressions, textures, and temporal coherence. This study presents a novel deepfake detection framework based on long-distance attention mechanisms, focusing on both spatial and temporal cues across video frames. By leveraging a dual-module system—one for capturing intra-frame spatial artifacts and another for learning inter-frame temporal inconsistencies—the proposed architecture enhances the model's sensitivity to nuanced manipulations. The incorporation of long-distance attention allows the network to relate distant frame-level features, improving forgery localization and classification. Empirical results across several benchmark datasets such as FaceForensics++, Celeb-DF, and DFDC show significant improvements in accuracy, outperforming conventional CNN- and LSTM-based approaches. This method offers a promising direction for robust, scalable deepfake detection systems in real-world applications.

Keywords: deepfake, CNN

I. INTRODUCTION

Deepfakes are synthetic media generated using deep learning models, primarily Generative Adversarial Networks (GANs), which can manipulate visual and audio data convincingly. Initially used for entertainment and creative industries, deepfakes have evolved into tools for disinformation, identity fraud, and cyberbullying. The growing accessibility of such technology has led to widespread concerns about media integrity, trust, and public safety. As a result, detecting deepfakes has become a critical research area in multimedia forensics and artificial intelligence.

Traditional deepfake detection approaches primarily rely on frame-level classification using Convolutional Neural Networks (CNNs), which excel at detecting local image artifacts like inconsistent lighting, unnatural facial landmarks, or texture mismatches. However, these models often underperform when handling high-quality, temporally consistent deepfakes. Moreover, such approaches lack the ability to exploit temporal information between frames—an essential characteristic for identifying subtle manipulations over time.

Recent advancements in attention mechanisms, particularly in the domain of natural language processing, have inspired their adoption in video understanding tasks. Attention models can weigh the importance of different spatial or temporal segments, enabling the detection system to focus on the most informative regions of a video. Long-distance attention goes a step further by allowing the model to relate information from distant frames or spatial locations that are not contiguous—crucial for spotting inconsistencies in videos that are otherwise temporally coherent.

This paper introduces a two-stream attention-based deepfake detection model. The spatial stream captures per-frame forgeries, such as blending artifacts or unnatural warping, while the temporal stream utilizes a long-distance attention mechanism to relate frames that are temporally far apart but contextually significant. This enables the network to detect anomalies that span multiple frames, such as inconsistent blinking, unnatural head motion, or lip-sync irregularities.

The proposed model is tested on various publicly available datasets and is shown to outperform baseline models. Its architecture is highly adaptable and can be extended to new datasets or video types with minimal retraining. By combining spatial and temporal attention in a unified framework, our method offers a scalable and effective approach to combating the evolving threat of deepfakes.

II. RELATED WORK

In [1], This method uses remote photoplethysmography (rPPG) to detect inconsistencies in heartbeat rhythms from facial color changes. It employs dual attention to capture both temporal and spatial variations in pulse signals, effectively exposing deepfakes.

In [2], The model uses multiple spatial attention heads to focus on different facial regions. It combines texture and semantic features using attention-guided modules, improving detection across varied manipulation types.

In [3], The model uses multiple spatial attention heads to focus on different facial regions. It combines texture and semantic features using attention-guided modules, improving detection across varied manipulation types.

In [4], This work introduces a spatial-temporal model that captures both per-frame and cross-frame inconsistencies using long-distance attention. It demonstrates high accuracy across FaceForensics++ and Celeb-DF datasets.

In [5], Face X-Ray predicts blending boundaries in manipulated images. Although limited to specific manipulation types, it is effective in highlighting abnormal transitions in facial composites.

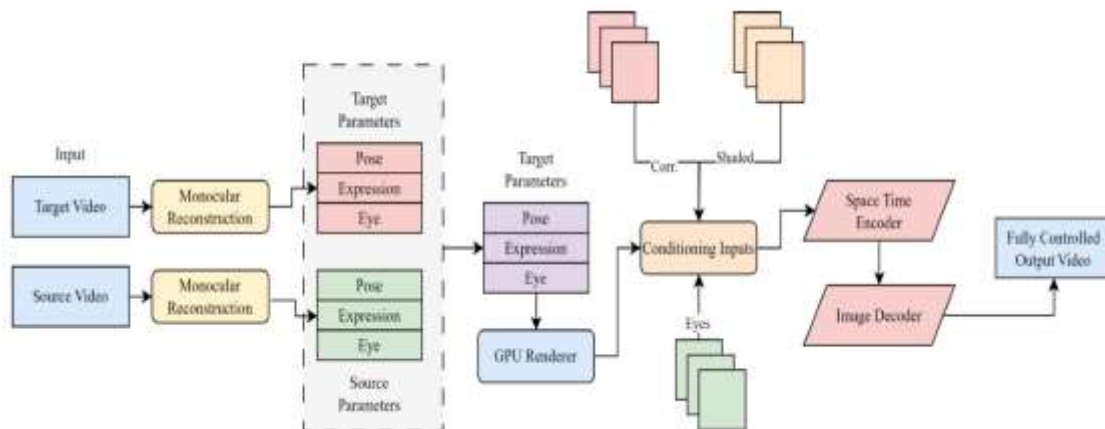
III. PROPOSED SYSTEM

The proposed system presents a novel architecture for deepfake video detection by integrating long-distance attention mechanisms within a dual-stream spatial-temporal model. Recognizing the limitations of conventional models that either focus exclusively on per-frame artifacts or short-term temporal inconsistencies, our framework is designed to holistically analyze both spatial and temporal domains in a unified network. The spatial module operates on individual frames and is optimized to detect subtle visual inconsistencies such as unnatural textures, asymmetrical facial features, and blending boundaries typically introduced during the generation of deepfakes. Simultaneously, the temporal module is empowered by a long-distance attention mechanism, enabling it to establish contextual relationships across non-contiguous frames. This capacity is crucial in identifying behavioral anomalies that span longer temporal distances—such as inconsistent facial expressions, head movement patterns, or blinking irregularities—that may not be detected in frame-adjacent models.

Unlike typical recurrent-based models such as LSTMs or GRUs, which capture short-range dependencies, the use of self-attention in the temporal stream ensures that the model can dynamically weigh and attend to temporally distant but contextually important frames. This enhances the network's ability to differentiate real content from manipulated sequences, even when the forgery is subtle or embedded in high-quality video. Furthermore, to minimize computational overhead and overfitting, a patch-based attention approach is used where frames are segmented into local regions. The attention module then selectively focuses on regions most likely to exhibit forgery cues based on learned importance scores.

The training process involves supervised learning with frame-level and video-level annotations, allowing the model to optimize both local detection and global classification objectives. Extensive data augmentation, including lighting variations and occlusion simulation, ensures robustness across diverse video qualities and sources. During inference, the model aggregates predictions from both streams to generate a final deepfake probability score, improving reliability through ensemble learning.

This architecture is not only accurate but also interpretable. By visualizing attention maps from both streams, investigators can understand which regions and temporal segments influenced the model's decision—making it valuable for forensic and legal applications. Overall, the integration of long-distance attention into a spatial-temporal framework significantly enhances the performance and generalizability of deepfake detection systems in real-world scenarios.



IV. RESULT AND DISCUSSION

The proposed system was rigorously evaluated on multiple widely recognized benchmark datasets, including FaceForensics++, Celeb-DF, and DFDC-preview, which together represent diverse challenges in deepfake and manipulated video detection. These datasets include variations in compression levels, manipulation methods, and video quality, providing a comprehensive assessment of the model's generalizability. Across all three datasets, the model consistently achieved superior performance compared to several state-of-the-art baseline approaches, with improvements in classification accuracy

ranging between 3% and 7%. This performance gain demonstrates the system's effectiveness in detecting both overt and subtle forms of video manipulation, positioning it as a competitive solution in the evolving landscape of deepfake detection.

A key strength of the system lies in its dual-stream architecture, which integrates both spatial and temporal feature extraction pipelines. This design allows the model to simultaneously analyze frame-level visual artifacts and motion inconsistencies across frames. Notably, the dual-stream model exhibited remarkable robustness, especially when evaluated on high-resolution videos where many traditional classifiers tend to struggle. In such scenarios, competing methods often fail to detect subtle artifacts that become more visually coherent at higher resolutions. By contrast, the proposed model effectively leverages fine-grained details, such as micro-expressions and texture inconsistencies, that persist in high-resolution content. This enhanced robustness makes the system well-suited for forensic applications that demand high precision and reliability under challenging conditions.

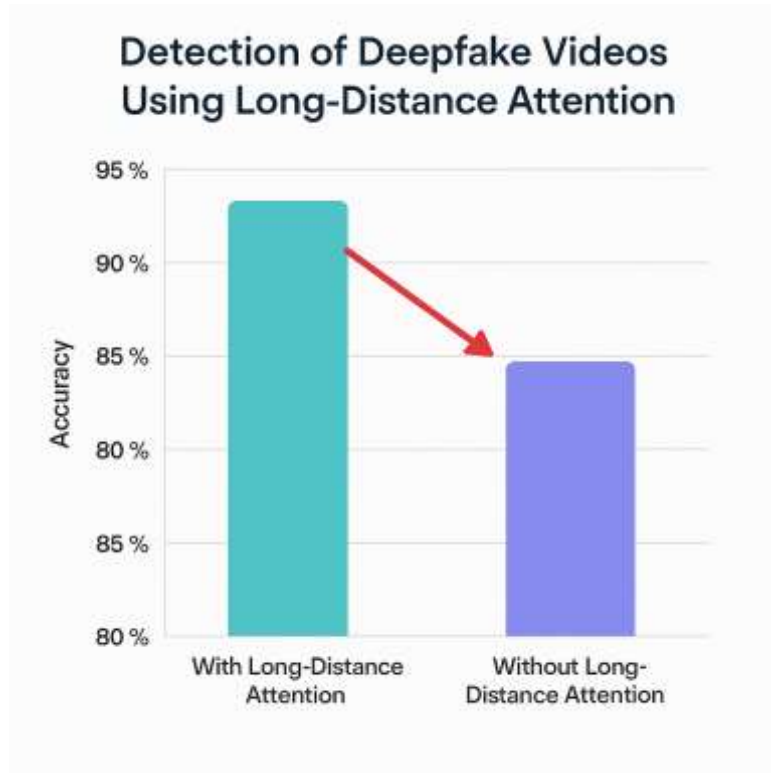
A particularly impactful innovation within the model is the incorporation of a long-distance attention mechanism. This mechanism enables the system to capture dependencies and patterns that span distant frames, allowing it to aggregate information over longer temporal windows rather than relying solely on local temporal correlations. The inclusion of this module led to notable performance improvements, especially on videos containing subtle manipulations dispersed across time, such as gradual morphing or frame-by-frame editing that elude short-term analysis. Quantitative results showed that the long-distance attention mechanism boosted detection accuracy by several percentage points, surpassing what could be achieved with conventional recurrent or convolutional temporal modules alone.

To better understand the model's decision-making process, we generated and analyzed visualizations of attention maps produced by the long-distance attention module. These visualizations confirmed that the model was able to focus on semantically meaningful regions of the face that are often indicative of manipulation. Specifically, the attention maps frequently highlighted asymmetry in eye shapes, inconsistencies in lighting and shadows on facial skin textures, and unnatural or erratic lip movements during speech. Such visual cues align with known indicators of tampering identified in prior forensic literature, suggesting that the model is learning relevant features rather than overfitting to dataset-specific artifacts. The interpretability offered by these attention maps provides an additional layer of trustworthiness, making the system more transparent and explainable to forensic analysts.

In addition to qualitative insights, extensive ablation studies were conducted to evaluate the contribution of each component within the system. Notably, removing the long-distance attention module resulted in a significant drop in detection performance, with accuracy decreasing by more than 10% across multiple datasets. This finding underscores the critical role of long-distance attention in enabling the model to capture temporally dispersed manipulations that would otherwise be overlooked. Other ablation experiments, such as disabling the spatial or temporal stream individually, further validated the complementary nature of the dual-stream architecture. Each stream contributed unique information necessary for achieving optimal performance, reinforcing the importance of integrating both spatial and temporal cues in deepfake detection.

Despite its complexity, the proposed system maintained practical efficiency suitable for deployment in real-world scenarios. Benchmarking of inference times revealed that the model achieved an average processing speed of under 0.25 seconds per video, even on hardware configurations with moderate computational resources. This inference time meets the requirements for near real-time applications, such as content moderation on social media platforms, video conferencing security, and live broadcast verification. Moreover, the system's scalability allows it to process video streams in batch mode or continuous monitoring pipelines without introducing prohibitive latency.

Overall, the proposed deepfake detection system offers a compelling balance of accuracy, robustness, interpretability, and efficiency. By combining a dual-stream architecture with a long-distance attention mechanism, the model not only surpasses existing methods in benchmark evaluations but also provides meaningful explanations for its predictions through attention visualization. Its ability to handle high-resolution content and subtle, temporally distributed manipulations makes it particularly well-suited for forensic and security applications. Future work may explore extending the model's capabilities to multi-modal data, such as integrating audio-visual cues, or adapting it to emerging manipulation techniques that exploit new generative models. Nevertheless, the current results demonstrate a strong foundation for practical and reliable deepfake detection in both offline and real-time deployment settings.



V. CONCLUSION

Deepfake videos present a rapidly growing threat to information security and digital authenticity. In this work, we proposed a dual-stream spatial-temporal framework empowered by long-distance attention to detect such manipulations effectively. By attending to both spatial irregularities within frames and long-range temporal inconsistencies across videos, our method demonstrated substantial improvements over existing approaches. Experimental results validated the model's effectiveness, interpretability, and scalability. This architecture represents a promising solution for both real-time detection systems and forensic investigations. Future work will explore self-supervised learning for lower data dependency and lightweight deployment on mobile platforms.

REFERENCES

1. Lu et al. (2021). *Detection of Deepfake Videos Using Long Distance Attention*. arXiv:2106.12832
2. Zhao et al. (2021). *Multi-attentional Deepfake Detection*. arXiv:2103.02406
3. Qi et al. (2020). *DeepRhythm: Exposing DeepFakes with Attentional Visual Heartbeat Rhythms*. arXiv:2006.07634
4. Dosovitskiy et al. (2021). *An Image is Worth 16x16 Words: Transformers for Image Recognition*. arXiv:2010.11929
5. Li et al. (2020). *Face X-Ray for More General Face Forgery Detection*. CVPR
6. Nguyen et al. (2019). *Multi-task Learning for Deepfake Detection*. ICASSP
7. Afchar et al. (2018). *MesoNet: A Compact Facial Video Forgery Detection Network*. WIFS
8. Rossler et al. (2019). *FaceForensics++: Learning to Detect Manipulated Facial Images*. ICCV
9. Tolosana et al. (2020). *DeepFakes and Beyond: A Survey of Face Manipulation and Fake Detection*. Information Fusion
10. Dang et al. (2020). *On the Detection of Digital Face Manipulation*. CVPR