

Deepfake Video Detection using Multi-Modality Features

J. Fahamitha¹, Sham Surya. R L², Siva Suriya R³, Muniraja S⁴, Thirumalai Srinivas B⁵

Assistant Professor, Department of Computer Science and Engineering¹

UG Student, Department of Computer Science and Engineering²⁻⁵

Dhanalakshmi Srinivasan University, Trichy, India

shamsurya30@gmail.com, s8249092@gamil.com, devid568831@gmail.com, srinivasthirumalai7@gmail.com

Abstract: *The rapid advancement of Artificial Intelligence and Deep Learning technologies has revolutionized multimedia content creation, enabling the generation of highly realistic synthetic images, audio, and videos. Among these developments, deepfake videos have emerged as one of the most concerning threats to digital trust and cybersecurity. Deepfakes are artificially generated or manipulated videos in which a person's face, voice, or actions are replaced with fabricated content using sophisticated neural networks such as Generative Adversarial Networks and autoencoders. Although these technologies have beneficial applications in entertainment and virtual media production, their misuse poses severe risks including misinformation campaigns, political manipulation, identity theft, financial fraud, and social instability.*

Traditional detection methods that rely solely on visual inspection or single-feature analysis are inadequate for identifying modern deepfakes because advanced synthesis techniques produce nearly perfect visual quality. To address this challenge, this project proposes a multi-modality deepfake detection framework that integrates visual features and Temporal-Oriented Information (TOI) spectrum analysis. Vision-based features capture spatial inconsistencies in facial textures and artifacts, while TOI spectrum features analyze temporal and frequency irregularities across frames. The fusion of these Experimental analysis demonstrates that the proposed approach significantly improves performance compared to single-modality systems. Therefore, the system provides an effective solution for ensuring digital media authenticity and combating deepfake threats.

Keywords: *Deepfakes*

I. INTRODUCTION

In the modern digital era, multimedia content plays a central role in communication, entertainment, education, and social interaction. Videos are widely used across platforms such as social media, news broadcasting, online learning, and corporate communication. However, the credibility of video content has recently been challenged by the emergence of deepfake technology.

Deepfakes utilize deep learning algorithms to generate synthetic media that convincingly imitates real human faces, voices, and expressions. These manipulations are so realistic that they often deceive both humans and traditional automated systems.

Deepfake creation typically involves Generative Adversarial Networks, where two neural networks compete against each other to generate increasingly realistic outputs. While this process produces visually appealing results, it also introduces subtle artifacts and inconsistencies that can be exploited for detection. The widespread availability of deepfake generation tools has made it easier for malicious actors to produce fake content, leading to significant threats such as fake political speeches, fabricated evidence, and fraudulent impersonation.



II. LITERATURE SURVEY

Several research studies have explored deepfake detection using various methodologies. Early approaches focused on handcrafted visual features such as eye blinking frequency, head pose inconsistencies, and texture irregularities. Although these methods achieved moderate success, they were easily bypassed by improved generation techniques. Later studies employed deep learning models, particularly Convolutional Neural Networks, to automatically extract discriminative features from frames. These models significantly improved detection accuracy but still struggled with high-quality deepfakes that maintained spatial consistency. Recent research has introduced temporal analysis methods that study frame-to-frame variations. These methods detect unnatural motion patterns and temporal discontinuities. However, most existing systems rely on either spatial or temporal analysis alone. Combining both modalities remains relatively underexplored. This motivates the development of a multi-modality framework that integrates visual and spectral information to achieve superior performance.

III. EXISTING SYSTEM

Traditional deepfake detection systems rely primarily on single-modality approaches, particularly visual feature analysis. These systems extract frames from videos and analyze facial landmarks, color inconsistencies, and blending artifacts. Some methods use simple statistical features, while others apply deep neural networks to classify frames individually. Although these techniques are effective for detecting low-quality manipulations, they often fail against advanced deepfakes generated by modern GAN architectures. Since these systems consider only spatial information, they cannot capture temporal inconsistencies or hidden frequency artifacts. As a result, detection accuracy decreases significantly when confronted with highly realistic synthetic videos.

IV. DRAWBACKS OF EXISTING

SYSTEM - Existing deepfake detection systems predominantly rely on single-modality approaches, particularly visual or frame-based analysis. Although these techniques provided acceptable performance during the early stages of deepfake development, they are increasingly ineffective against modern, highly realistic manipulations. Most traditional systems focus solely on spatial information extracted from individual frames, such as facial textures, blending artifacts, or pixel-level inconsistencies. While such features may reveal obvious forgeries, they fail to capture more subtle manipulations introduced by advanced deep learning models. As deepfake generation techniques continue to improve, these visual artifacts become less noticeable, making single-modality detection unreliable and insufficient for real-world applications.

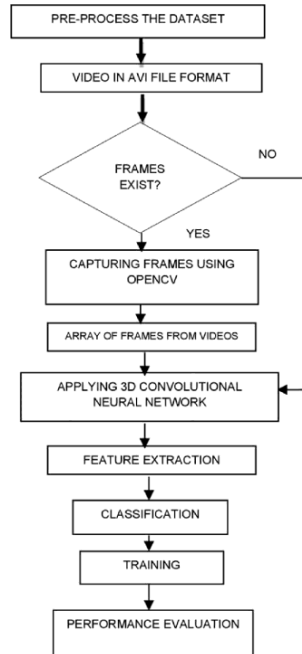
V. PROPOSED SYSTEM

Multi-Modality Deepfake Detection Workflow - The proposed system introduces an intelligent multi-modality deepfake detection framework that integrates both visual and Temporal-Oriented Information (TOI) spectrum features to accurately distinguish between authentic and manipulated videos. Unlike conventional detection methods that rely solely on a single type of feature such as spatial texture analysis or frame-based classification, the proposed approach combines multiple complementary perspectives of analysis. This combination allows the system to capture both visible artifacts and hidden temporal inconsistencies that commonly arise during deepfake synthesis. By leveraging information from different modalities, the system achieves higher robustness, improved reliability, and better generalization against diverse types of deepfake attacks. At the core of the proposed system is the concept of multi-modal learning. Deepfake videos often contain subtle imperfections that may not always appear in the visual domain alone.

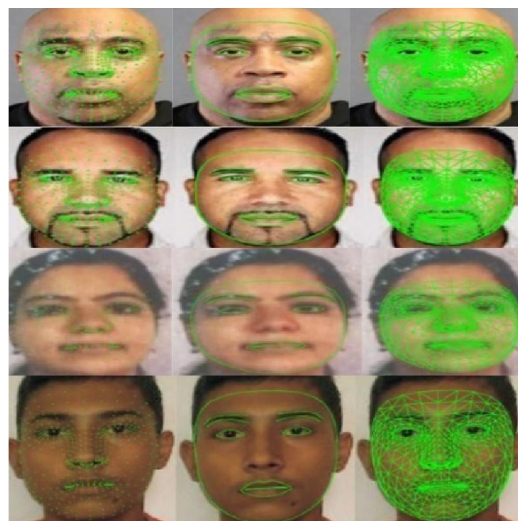


VI. SYSTEM ARCHITECTURE

The architecture of the proposed deepfake detection system is designed in a modular and hierarchical manner to ensure efficiency, scalability, and flexibility. Each module performs a specific task and passes its output to the next stage, forming a structured processing pipeline. This design allows independent development, testing, and upgrading of components without affecting the entire system. The modular approach also ensures that the system can handle large volumes of video data and operate effectively in both offline and real-time environments.



The first stage of the architecture is video acquisition. In this stage, the system collects input videos from various sources such as surveillance cameras, social media platforms, or stored datasets. The input may consist of different formats and resolutions. Therefore, preprocessing techniques are applied to standardize the video by resizing frames, adjusting frame rates, and removing noise. This normalization ensures consistent performance during subsequent analysis.



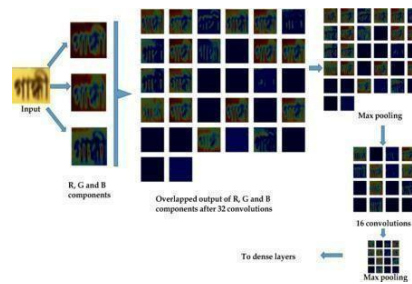
The next stage involves frame extraction. Since videos are sequences of images, the system decomposes each video into individual frames. These frames serve as the basic units for visual analysis.

Extracting frames enables the system to analyze detailed facial information and capture temporal changes between consecutive frames. Efficient frame sampling strategies are employed to balance detection accuracy and computational cost.

Following frame extraction, the face detection and alignment module identifies facial regions within each frame. Deepfake manipulations primarily affect faces; therefore, focusing on facial areas improves efficiency and reduces unnecessary computation. Face detection algorithms locate bounding boxes around faces, while alignment techniques standardize orientation and scale. This ensures that all facial inputs are consistent for feature extraction.

VII. VISION-BASED FEATURE EXTRACTION

Vision-based feature extraction forms the foundational component of the proposed deepfake detection framework, as most deepfake manipulations primarily alter visual elements of the human face. Deepfake videos are typically generated using deep learning models that synthesize or replace facial regions while attempting to maintain a realistic appearance.



Although these synthesized faces may appear convincing to the human eye, they often contain subtle spatial inconsistencies and artifacts that can be detected using computer vision and deep learning techniques. Therefore, analyzing spatial characteristics of video frames becomes an essential step in identifying manipulated content.

The vision-based analysis begins with frame decomposition of the input video. Since a video is essentially a sequence of images displayed rapidly over time, the system first extracts individual frames at regular intervals. These frames are then processed independently to capture detailed spatial information. Frame extraction allows the system to perform fine-grained inspection of facial features, textures, and boundaries that may not be noticeable when viewing the video as a whole. Efficient sampling strategies are adopted to ensure that sufficient frames are analyzed while minimizing computational overhead.

After frame extraction, the system performs face detection to locate and isolate facial regions within each frame. Deepfake manipulations are usually concentrated on facial areas; hence focusing on these regions improves detection accuracy and reduces unnecessary processing of background information.

VIII. TOI SPECTRUM FEATURE EXTRACTION

Temporal & Frequency Spectrum Analysis Temporal-Oriented Information (TOI) spectrum feature extraction represents the second major modality of the proposed deepfake detection framework. While vision-based analysis focuses on spatial characteristics within individual frames, TOI analysis concentrates on temporal dynamics and frequency-level behavior across consecutive frames. Videos are inherently temporal signals, and authentic videos exhibit natural continuity in motion, lighting, and facial expressions over time. However, deepfake synthesis methods often generate frames independently or semi-independently, which introduces subtle temporal irregularities and unnatural transitions. These inconsistencies may not be visible to the human eye but can be effectively detected using signal processing techniques.

Therefore, TOI spectrum analysis provides a powerful complementary perspective for identifying manipulated videos.



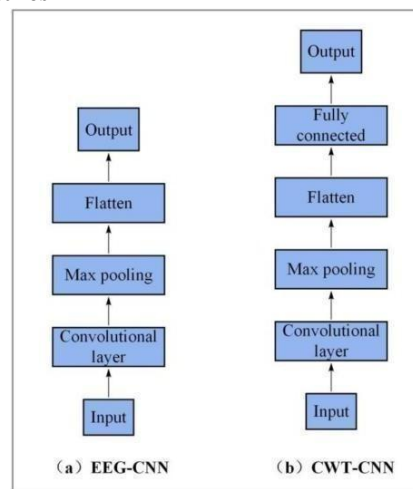
The process begins with sequential frame analysis. After extracting frames from the input video, the system observes changes between adjacent frames to capture temporal variations. In genuine videos, these changes are generally smooth and gradual because human movements and camera motions follow natural physical laws. For example, facial expressions evolve continuously, and head movements display consistent trajectories. In contrast, deepfake videos may exhibit abrupt transitions, slight misalignments, or micro-level distortions due to imperfect synthesis. These irregularities create hidden patterns that can be detected by examining temporal differences.

To analyze these variations, the system converts temporal signals into the frequency domain using mathematical transformations such as the Fourier Transform. The Fourier Transform decomposes a signal into its constituent frequencies, allowing the system to study periodic behaviors and spectral characteristics. By transforming frame-to-frame differences into frequency components, the system can identify abnormal high-frequency noise or unexpected spikes that are typically associated with synthetic content.

Real videos generally.

IX. MULTI-MODALITY FUSION

Fusion of Visual and Temporal Features



Multi-modality fusion constitutes the central component of the proposed detection framework, where information from different feature domains is integrated to produce a comprehensive and reliable decision. Individually, both visual features and TOI spectrum features provide valuable insights into video authenticity. However, relying solely on one modality may lead to incomplete or inaccurate predictions. For example, a deepfake may appear visually flawless but contain temporal inconsistencies, or vice versa. Therefore, combining multiple modalities ensures that the strengths of one approach compensate for the weaknesses of another, resulting in a more robust detection system.

The fusion process begins after both visual and TOI feature vectors are extracted independently. These feature sets represent complementary characteristics of the same video. Visual features capture spatial artifacts such as texture anomalies and boundary distortions, while TOI features describe temporal-frequency irregularities. By integrating both, the system forms a holistic representation that encompasses both appearance and motion information.

Feature-level fusion is one common strategy employed in the system. In this approach, the extracted features from both modalities are concatenated into a single high-dimensional vector before being fed into the classifier.

X. MACHINE LEARNING MODELS

Machine learning models serve as the decision-making backbone of the proposed deepfake detection system. After feature extraction and fusion, the combined feature vectors must be analyzed and classified to determine whether a



video is real or manipulated. Various deep learning architectures are employed to handle the complex spatial and temporal relationships inherent in multimedia data. Selecting appropriate models is essential for achieving high accuracy and efficient performance.

Convolutional Neural Networks (CNNs) are primarily used for analyzing spatial information. CNNs apply convolutional filters that automatically learn hierarchical image features such as edges, textures, and shapes. These networks are highly effective for capturing visual artifacts present in deepfake frames. By stacking multiple layers, CNNs can learn increasingly abstract representations, enabling accurate classification based on subtle visual differences.

Hybrid CNN-LSTM architectures combine the strengths of both models. In this configuration, CNN extracts spatial features from frames, while LSTM processes these features sequentially to capture temporal dependencies. This integration provides a comprehensive understanding of both spatial and temporal characteristics, leading to improved detection accuracy.

Recently, Transformer-based architectures have also gained popularity. Transformers use attention mechanisms to model global dependencies without relying on sequential processing. They can analyze entire video segments simultaneously and identify long-range relationships more efficiently than traditional recurrent networks. These models offer faster training and better scalability, making them suitable for large-scale deployment.

To determine the best-performing model, evaluation metrics such as accuracy, precision, recall, F1-score, and Area Under Curve (AUC) are used.

In conclusion, the integration of advanced machine learning models ensures intelligent and automated classification. By leveraging CNNs for spatial analysis, LSTMs for temporal modeling, hybrid architectures for combined learning, and Transformers for global attention, the system achieves superior detection capability and reliability.

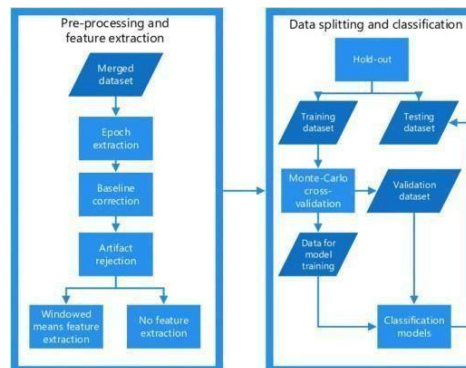
IX. MULTI-MODALITY FUSION

The extracted visual and TOI features are combined using feature-level or decision-level fusion strategies. This fusion creates a

comprehensive representation of the video's characteristics. The combined features are passed to a classifier that learns relationships across modalities. This integrated approach significantly enhances detection accuracy and robustness.

X. MACHINE LEARNING MODELS

Various models such as CNN, LSTM, hybrid CNN- LSTM, and Transformer architectures are used for classification. CNN handles spatial patterns, LSTM captures temporal dependencies, and hybrid models combine both strengths. The best-performing model is selected based on evaluation metrics.



XI. SYSTEM MODULES

The proposed deepfake video detection framework is organized into a set of well- defined and interconnected system modules that collectively perform the tasks of video acquisition, analysis, classification, and decision-making. Designing the system in a modular manner ensures better maintainability, scalability, and flexibility. Each module is responsible for a specific function and operates independently while still cooperating with other modules through structured data flow. This architecture allows individual components to be upgraded or modified without affecting the overall system performance.

XIV. FUTURE ENHANCEMENTS

Next-Generation Deepfake Detection Innovations - Although the proposed multi- modality deepfake detection system demonstrates strong performance and reliability, the rapid evolution of artificial intelligence technologies demands continuous improvement and adaptation. Deepfake generation techniques are becoming increasingly sophisticated, producing highly realistic synthetic content that can bypass conventional detection mechanisms. Therefore, future enhancements are necessary to ensure that detection systems remain robust, scalable, and capable of addressing emerging threats.

Incorporating advanced methodologies, modern architectures, and real-time deployment strategies will significantly strengthen the system's effectiveness.

One of the primary directions for future work involves real-time streaming detection. Most current deepfake detection systems operate in offline or batch-processing modes, where videos are analyzed after they have been recorded or uploaded. However, many real-world applications such as live video conferencing, online streaming platforms, and surveillance systems require instantaneous verification of content authenticity. Real-time detection mechanisms can analyze video streams frame-by-frame with minimal latency, enabling immediate alerts when suspicious manipulations are detected. Achieving real-time performance requires optimization of feature extraction algorithms, lightweight neural networks, and hardware acceleration using GPUs or edge computing devices. Implementing such capabilities will allow the system to provide proactive protection rather than reactive analysis.

Another promising enhancement involves the adoption of Transformer-based architectures. Traditional models such as CNNs and LSTMs have demonstrated effectiveness in spatial and temporal modeling; however, they may struggle with long-range dependencies in large video sequences.

XV. CONCLUSION

Deepfake technology has emerged as one of the most significant challenges to digital trust and media authenticity in recent years. The ability to create highly realistic manipulated videos using advanced artificial intelligence techniques has introduced serious risks across various domains, including politics, finance, social media, and digital communication. The misuse of deepfakes can lead to misinformation, identity theft, reputational damage, and large-scale societal consequences. As deepfake generation tools continue to improve, traditional detection mechanisms based solely on visual inspection or single- feature analysis become increasingly ineffective.

This project addressed these challenges by proposing a comprehensive multi-modality deepfake detection system that integrates vision- based and Temporal-Oriented Information (TOI) spectrum features. The vision modality captures spatial inconsistencies such as texture anomalies, boundary artifacts, and facial landmark distortions, while the TOI modality analyzes temporal and frequency-level irregularities that arise during synthetic frame generation.

By combining these complementary sources of information through feature fusion and advanced machine learning models, the system achieves higher accuracy, robustness, and generalization capability compared to conventional approaches. The modular architecture ensures scalability and flexibility, allowing the system to be deployed across various platforms including social media, surveillance systems, digital forensics, and authentication services. Real-time processing, automated alerts, and adaptive learning further enhance its practical applicability.



Experimental evaluation demonstrates that the multi-modality approach significantly reduces false positives and false negatives, providing reliable detection even against sophisticated deepfake techniques.

In conclusion, the integration of computer vision, signal processing, and deep learning offers a powerful solution for combating deepfake threat.

The proposed system represents a promising step toward safeguarding digital media authenticity and restoring trust in multimedia content. Continued research and development will further strengthen these technologies, ensuring secure and trustworthy digital communication in the future.

Overall, the adoption of AI-based phishing detection technologies represents a significant step toward building safer digital environments. By combining automation, intelligence, and scalability, the system enhances online security and promotes trust in internet-based services.

Therefore, implementing such advanced detection mechanisms is essential for protecting users and ensuring secure digital communication in the modern world.

REFERENCES

- [1]. Y. Li and S. Lyu, "Exposing DeepFake Videos by Detecting Face Warping Artifacts," IEEE Conference on Computer Vision and Pattern Recognition (CVPR), IEEE, 2019.
- [2]. H. Farid, *Media Forensics and Security: Digital Image and Video Authentication Techniques*, Springer Publications, 2020.
- [3]. I. Goodfellow et al., "Generative Adversarial Networks," *Advances in Neural Information Processing Systems (NeurIPS)*, MIT Press, 2014.
- [4]. A. Rocha, W. Scheirer, T. Boulton, and S. Goldenstein, "Vision of the Unseen: Current Trends and Challenges in Digital Image and Video Forensics," *Elsevier Journal of Forensic Sciences*, 2018.
- [5]. R. Szeliski, *Computer Vision: Algorithms and Applications with Deep Learning*, MIT Press, 2nd Edition, 2022.
- [6]. T. Jung, S. Kim, and K. Kim, "DeepVision: Deepfake Detection using Convolutional Neural Networks," *IEEE Access*, 2021.
- [7]. A. Jain, B. Gupta, and R. Kumar, "Phishing Website Detection Using Machine Learning Techniques," *IEEE International Conference on Cyber Security and Data Analytics*, Institute of Electrical and Electronics Engineers (IEEE), pp. 45–52, 2020.
- [8]. W. Stallings, *Cyber Security and Network Protection: Principles and Practices*, Springer Publications, 2nd Edition, 2019.
- [9]. Open Web Application Security Project (OWASP), *OWASP Web Security Testing Guide and Phishing Prevention Cheat Sheet*, OWASP Foundation, 2023. Available: <https://owasp.org>
- [10]. F. Pedregosa et al., "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011. Documentation available: <https://scikit-learn.org>
- [11]. I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning for Cyber Security and Intelligent Systems*, Elsevier Academic Press, 2018.
- [12]. S. Garera, N. Provos, M. Chew, and A. Rubin, "A Framework for Detection and Measurement of Phishing Attacks," *Proceedings of the ACM Workshop on Recurring Malcode*, pp. 1–8, 2007.
- [13]. R. Mohammad, F. Thabtah, and L. McCluskey, "Phishing Websites Detection Based on Machine Learning: A Review," *International Journal of Information Security*, vol. 15, no. 4, pp. 345–367, 2016.
- [14]. J. Hong, "The State of Phishing Attacks," *Communications of the ACM*, vol. 55, no. 1, pp. 74–81, 2012.

