

Deepfake Detection Using Hybrid CNN-Transformer Model

Mrs. J. P. Aswini

Assistant Professor, Computer Science and Engineering,
St. Joseph College of Engineering, Chennai-602117, Tamil
Nadu,

Email Id – aswinianchu1997@gmail.com

Mr. A .Kevins Nithin

Student, Computer Science and Engineering,
St. Joseph College of Engineering, Chennai-602117, Tamil Nadu,

Email Id –kevinsnithin4500@gmail.com

Mr. L .Sahaya Rasin

Student, Computer Science and Engineering,
St. Joseph College of Engineering, Chennai-602117, Tamil Nadu,

Email Id – hermassahaya@gmail.com

Abstract - The rapid evolution of synthetic media generation techniques, particularly those based on Generative Adversarial Network (GANs) and Autoencoder architectures, has led to the widespread creation of highly realistic deepfake images and videos. While these technologies have legitimate applications in entertainment and media production, they also pose serious risks including misinformation, identity theft, and digital fraud. Conventional deepfake detection methods that rely solely on spatial feature extraction often struggle to generalize against sophisticated manipulations, especially when faced with high-quality forgeries, varying illumination, and compression artifacts.

To address these challenges, this work proposes a Hybrid CNN-Transformer model for deepfake detection. The Convolutional Neural Network (CNN) component extracts fine-grained spatial features such as texture inconsistencies and blending artifacts, while the Transformer module utilizes self-attention mechanisms to capture global contextual relationships and long-range dependencies across facial regions and video frames. By integrating local feature learning with global attention modelling, the proposed architecture enhances robustness and detection accuracy. Experimental results demonstrate improved performance over standalone CNN and Transformer models in terms of accuracy, precision, recall, and F1-score, making the approach suitable for real-world deepfake detection applications.

I.INTRODUCTION

In recent years, Artificial Intelligence has enabled the creation of highly realistic synthetic media known as deepfakes. Deepfakes are digitally manipulated images or videos in which a person's face, voice, or expressions are altered using deep learning techniques. Technologies such as Generative Adversarial Networks (GANs) have made it possible to generate fake content that appears almost identical to real media.

Although deepfake technology has useful applications in filmmaking, gaming, virtual reality, and entertainment, it also presents serious ethical and security concerns. Deepfakes can be misused for spreading misinformation, political propaganda, identity theft, financial fraud, and cyber harassment. The increasing availability of tools like

DeepFaceLab has made it easier for individuals to create convincing fake videos without advanced technical knowledge.

As deepfake generation techniques become more sophisticated, detecting manipulated content has become increasingly challenging. Traditional image forensic methods are no longer sufficient because modern deepfakes can replicate lighting conditions, facial expressions, and lip synchronization with high accuracy. Therefore, advanced deep learning-based detection systems are required.

Convolutional Neural Networks (CNNs) are widely used for image analysis as they effectively capture local spatial features such as edges, textures, and facial details. On the other hand, Transformer models are powerful in capturing global contextual relationships and long-range dependencies within data.

This project proposes a Hybrid CNN–Transformer model that combines the strengths of both architectures. The CNN component extracts detailed spatial features from facial regions, while the Transformer component analyzes global patterns and inconsistencies across the image or video frames. By integrating these approaches, the system aims to improve deepfake detection accuracy and robustness.

The ultimate goal of this project is to develop an efficient and reliable deepfake detection system that can help maintain digital trust, enhance cybersecurity, and contribute to the prevention of malicious misuse of synthetic media.

II. BACKGROUND AND MOTIVATION

1. Introduction

With the rapid advancement of Artificial Intelligence, especially Deep Learning, synthetic media generation has significantly improved. Techniques such as Generative Adversarial Networks (GANs) and autoencoders are widely used to create realistic fake images and videos, commonly known as deepfakes.

Applications like DeepFaceLab and FaceSwap allow users to easily swap faces in videos with high realism. While this technology has positive uses in entertainment and filmmaking, it also poses serious threats such as:

- Spread of misinformation
- Political manipulation
- Identity theft
- Cyberbullying and harassment
- Financial fraud

Deepfakes are increasingly being used to impersonate public figures and ordinary individuals. Incidents involving manipulated media of personalities like Barack Obama have demonstrated how realistic and convincing deepfake content can be.

Traditional image forensic techniques struggle to detect modern deepfakes because they are trained to identify simple pixel-level inconsistencies. However, advanced deepfakes

can mimic facial expressions, lighting conditions, and lip synchronization very accurately.

To overcome this, deep learning-based detection methods have been introduced.

Convolutional Neural Networks (CNNs) are effective at extracting local spatial features such as textures and edges.

Transformer models are powerful in capturing global relationships and long-range dependencies in data.

A hybrid model combining CNN and Transformer architectures can leverage the strengths of both approaches to improve detection accuracy.

2. Motivation

The motivation behind this project arises from the increasing misuse of deepfake technology and the need for reliable detection systems.

1. Growing Threat to Digital Trust

The rise of manipulated media weakens public trust in digital content. A robust detection system is essential to maintain credibility in news, social media, and online communication.

2. Limitations of Existing Systems

CNN-only models focus mainly on local features.

Transformer-only models may require large datasets and high computation.

Hybrid models can provide better feature representation and improved performance.

3. Real-World Application Requirement

Social media platforms, law enforcement agencies, and cybersecurity systems need automated tools to detect fake media in real time.

4. Technological Advancement

Combining CNN with Transformer represents a modern and efficient approach in deep learning research. This project contributes to ongoing research in AI security and digital forensics.

3. Problem Statement

To design and implement a Hybrid CNN–Transformer model capable of accurately detecting deepfake images and videos by extracting both local spatial features and global contextual features.

III. DEEPPFAKE DETECTION USING A HYBRID CNN–TRANSFORMER ARCHITECTURE

The proposed system introduces a novel AI-driven application that integrates deep Convolutional Neural Networks (CNNs) with Transformer-based attention mechanisms to create an intelligent deepfake detection framework. Unlike conventional detection methods that rely solely on spatial feature extraction or handcrafted forensic techniques, this solution enables automatic identification and classification of manipulated images and videos with high accuracy.

The novelty lies in combining local spatial feature extraction with global contextual relationship modelling within a unified hybrid architecture. By leveraging advanced deep learning techniques inspired by convolutional feature learning and self-attention mechanisms, the system effectively detects subtle texture inconsistencies, blending artifacts, and unnatural facial patterns present in synthetic media. The application not only identifies forged content but also evaluates confidence scores and contextual inconsistencies across facial regions and video frames to enhance detection reliability.

Another innovative aspect is its adaptability to diverse real-world scenarios, including varying illumination conditions, compression artifacts, and high-quality GAN-generated forgeries. The system can be integrated with existing digital forensic and social media monitoring platforms, reducing dependency on manual verification while supporting scalability across large multimedia datasets. Overall, this architecture transforms traditional deepfake detection approaches into a robust, AI-powered security mechanism, contributing to improved digital trust, enhanced cybersecurity, and the advancement of intelligent media authentication systems.

IV. ROLE AND POTENTIAL OF DEEPPAKE DETECTION USING A HYBRID CNN–TRANSFORMER ARCHITECTURE

Role:

The proposed system plays a transformative role in modern digital security systems by integrating hybrid deep learning-based feature extraction with intelligent media authentication. It shifts deepfake detection from a static, single-model analysis approach to a dynamic, AI-driven hybrid verification framework.

A. Role in Intelligent Media Verification – Traditional detection systems rely mainly on spatial feature analysis without contextual understanding. This project introduces hybrid deep learning intelligence into digital platforms. By continuously analyzing images and video frames, the system automatically identifies manipulated media and evaluates authenticity using CNN-based spatial learning and Transformer-based contextual modeling. This reduces human intervention and ensures faster, unbiased, and consistent verification.

B. Role in Combating Misinformation and Digital Fraud – Synthetic media manipulation has become a serious digital threat. The system ensures:

- Immediate detection of GAN-generated and autoencoder-based deepfakes
- Identification of texture inconsistencies and blending artifacts
- Contextual analysis across facial regions and video frames
- Confidence-based authenticity classification

Thus, it directly contributes to protecting digital credibility and cybersecurity.

C. Role in Smart Digital Infrastructure – Modern digital ecosystems rely on AI-driven monitoring systems. The proposed framework can be integrated into social media platforms, cybersecurity systems, and digital forensic tools, enabling:

- Realtime Deepfake monitoring
- Automated content flagging
- Authenticity scoring
- Performance analytics

This makes the system compatible with next-generation digital governance initiatives.

D. Role in AI-Based Detection Systems – The project demonstrates the practical application of Hybrid CNN–Transformer architectures in real-world digital security domains. The system performs:

- Fine-grained spatial feature extraction
- Global contextual relationship modelling
- Multi-class classification
- Self-attention-based dependency learning
- Confidence-based prioritization

E. Role in Enhancing Detection Robustness – By combining CNN-based local learning with Transformer-based global attention mechanisms, the system maintains high detection accuracy even under varying illumination and compression artifacts.

Potential:

The potential of this project extends significantly beyond its current implementation and opens multiple avenues for research and development.

A. Scale Deployment Potential – The system can be scaled across large digital platforms and cloud-based verification systems.

B. Integration with Secure Authentication Systems – The framework can be enhanced by integrating blockchain-based verification and digital watermarking systems.

C. Predictive Deepfake Trend Analysis – Future extensions can include predictive analytics to detect emerging manipulation techniques.

D. Expansion to Multi-Modal Detection – The model can be extended to detect:

- Audio deepfakes
- Voice cloning attacks
- Lip-sync manipulation
- Synthetic text generation

E. Data Analytics and Policy Support – The collected detection data can support:

- Deepfake trend monitoring
- Platform vulnerability assessment
- Digital governance planning

F. Research and Academic Potential – The project opens opportunities for improving hybrid architectures and implementing edge-based AI detection.

G. Social and Economic Impact Potential – Reliable deepfake detection can protect identities, prevent fraud, reduce misinformation, and strengthen digital trust.

V. INNOVATIVE INTEGRATION OF DEEP LEARNING IN DEEPPFAKE DETECTION USING HYBRID CNN–TRANSFORMER MODEL

The proposed system represents an innovative convergence of artificial intelligence, computer vision, deep learning architectures, and digital media forensics. Rather than implementing deepfake detection as a standalone image classification task, the project integrates feature extraction, contextual modeling, and authenticity verification into a unified intelligent detection framework.

A. End-to-End AI-Based Detection Integration – A major innovation of this project lies in its hybrid pipeline architecture. The system does not stop at identifying manipulated pixels; instead, it connects spatial feature extraction with global contextual reasoning to produce reliable authenticity decisions.

The integrated workflow includes:

- Image and video frame acquisition
- Frame preprocessing and normalization
- CNN-based spatial feature extraction
- Transformer-based contextual attention modeling
- Confidence-based authenticity classification

This seamless interaction between local feature learning (CNN) and global dependency modeling (Transformer) transforms passive media analysis into an intelligent verification mechanism.

B. Real-Time Deep Learning Deployment – Deep learning models are often used for offline media analysis. This project innovatively designs a Hybrid CNN–Transformer architecture capable of supporting near real-time detection of manipulated images and video frames.

Inspired by advanced computer vision and attention-based architectures, the system ensures:

- Low-latency frame analysis
- Multi-class deepfake classification
- Accurate localization of manipulation artifacts
- Continuous frame-by-frame consistency monitoring

The integration of spatial and contextual intelligence makes the system practically deployable for digital platforms rather than purely theoretical.

C. Multi-Layer Hybrid Architecture – The project introduces a layered hybrid architecture that enhances robustness and scalability. Each module performs a

specialized function while remaining interconnected:

- Preprocessing Layer – Performs resizing, normalization, and artifact enhancement to improve detection under varying lighting and compression conditions.
- CNN Feature Extraction Layer – Learns hierarchical spatial features such as texture inconsistencies, blending artifacts, and unnatural facial patterns.
- Transformer Attention Layer – Captures global contextual relationships and long-range dependencies across facial regions and video frames.
- Decision and Classification Layer – Implements confidence-based thresholding and authenticity scoring to minimize false predictions.
- This layered integration ensures improved detection accuracy and simplifies future architectural enhancements.

D. Context-Aware Model Customization – Unlike generic image classification systems, this project incorporates dataset-specific adaptation. By training the hybrid model on diverse deepfake datasets, the system accounts for:

- GAN-generated manipulations
- Autoencoder-based face swaps
- Lighting and compression variations
- Subtle expression and lip-sync inconsistencies

This contextual adaptation increases detection reliability across heterogeneous digital environments.

E. Hybrid Intelligence Expansion – The integrated framework is designed to support extended hybrid intelligence models. The system can be expanded to combine:

- Visual CNN–Transformer detection
- Audio deepfake recognition
- Lip-sync inconsistency detection
- Metadata and blockchain-based verification

This multi-modal integration would significantly reduce false positives and enhance overall trust in digital content authentication.

F. Scalability Through Edge and Cloud Integration – The architecture allows flexible deployment models:

- Edge deployment for platform-level rapid verification
- Cloud integration for large-scale media monitoring
- Hybrid distributed models for scalable detection

This makes the system adaptable for both small-scale applications and enterprise-level digital security infrastructure.

G. Data-Driven Feedback Loop Integration – The project not only detects manipulated content but also enables continuous improvement. Detection logs and classification outputs can be analyzed to:

- Measure detection accuracy trends
- Identify emerging manipulation techniques
- Optimize retraining strategies
- Improve model robustness against adversarial attacks

This creates a self-improving AI-enabled digital media authentication ecosystem.

H. Integration with Future Digital Security Systems – The framework is compatible with advanced cybersecurity technologies, including:

- Blockchain-based media verification
- Secure digital watermarking systems
- Explainable AI (XAI) visualization tools
- AI-driven automated content moderation systems

Thus, the system is not a standalone detection tool but a foundational module for next-generation intelligent media authentication platforms

VI. RECENT ADVANCEMENT IN DEEP LEARNING AND HUMAN SENTIMENT ANALYSIS

Multimodal and Real-Time Deepfake Detection – Recent advancements in deep learning have significantly enhanced deepfake detection systems by integrating multiple data modalities such as visual features, temporal frame consistency, and audio-visual synchronization analysis for improved reliability. Modern AI-powered models utilize advanced Convolutional Neural Networks (CNNs) combined with transformer-based attention architectures to identify manipulated facial regions in complex multimedia environments. These systems analyze texture inconsistencies, blending artifacts, unnatural lighting patterns, and facial geometry distortions while simultaneously examining contextual relationships across frames. Real-time adaptive learning enables detection models to dynamically respond to evolving deepfake generation techniques, ensuring consistent performance under varying illumination conditions, compression artifacts, and high-quality GAN-generated forgeries. This multimodal and hybrid approach strengthens digital media authentication systems by enabling automatic detection of manipulated content, reducing misinformation spread, and enhancing cybersecurity infrastructure.

Edge Deployment and Intelligent Digital Platform Integration – With the development of lightweight and optimized hybrid deep learning architectures, deepfake detection models can now be deployed on edge devices and cloud platforms for low-latency processing. Real-time inference on social media platforms, streaming services, and digital communication systems allows immediate verification of uploaded content. This integration ensures that manipulated media can be flagged or restricted before widespread distribution. Additionally,

advancements in transfer learning, synthetic data augmentation, and adversarial training improve model accuracy even with limited labeled deepfake datasets. These improvements make the system scalable for large multimedia platforms, enterprise-level cybersecurity systems, and intelligent content moderation frameworks, ensuring seamless coordination between AI detection modules and digital communication ecosystems.

Fairness, Explainability, and Future Directions – As deepfake detection systems become critical components of digital security infrastructure, recent research emphasizes fairness, transparency, and privacy-preserving AI deployment. Explainable AI (XAI) techniques, such as attention heatmaps and feature visualization, help interpret model decisions, ensuring reliability in sensitive applications such as media verification and legal investigations. Privacy-aware learning methods reduce risks associated with large-scale multimedia data processing, while bias mitigation strategies ensure consistent detection performance across diverse demographic and environmental conditions. Future advancements are expected to focus on improved multimodal fusion, self-supervised learning for better generalization, adversarial robustness against evolving forgery techniques, and energy-efficient hybrid architectures for sustainable deployment. These innovations collectively strengthen the role of hybrid CNN–Transformer models in building secure, trustworthy, and intelligent digital media ecosystems.

VII. CHALLENGES

Adversarial and High-Quality Forgery Challenges – One of the major challenges in deepfake detection is the rapid advancement of generative models such as GANs and autoencoders that produce highly realistic synthetic media. Modern deepfakes can replicate facial expressions, lighting conditions, skin textures, and lip synchronization with remarkable accuracy. High-resolution forgeries and compression techniques used by social media platforms further reduce visible artifacts, making it difficult for detection models to distinguish manipulated content from genuine media. Additionally, adversarial attacks intentionally modify deepfake samples to mislead detection systems, further complicating reliable classification.

Data and Model Limitations – Hybrid CNN–Transformer models require large, diverse, and well-balanced datasets for effective training. However, deepfake datasets may be limited in diversity, especially regarding ethnicity, lighting conditions, and manipulation techniques. Imbalanced datasets can cause biased learning and reduced generalization performance. Variations in generation methods, evolving GAN architectures, and cross-dataset inconsistencies also affect robustness. Overfitting, high computational complexity, and the need for continuous retraining to adapt to newly emerging deepfake techniques present further technical challenges in maintaining high detection accuracy.

Real-Time Deployment and Platform Integration – Implementing deepfake detection systems in real-world digital platforms demands low-latency processing and scalable infrastructure. Hybrid CNN–Transformer architectures are computationally intensive, which may slow inference speed, especially for real-time video streams or large-scale

social media uploads. Moreover, false positives can incorrectly flag authentic content, while missed detections can allow manipulated media to spread widely. Ensuring efficient integration with content moderation systems, cybersecurity platforms, and cloud-based verification pipelines remains a significant challenge for large-scale deployment.

Explainability and Ethical Concerns – As deepfake detection systems become part of critical digital security infrastructure, transparency and fairness become important challenges. Complex hybrid architectures often function as black-box models, making it difficult to interpret detection decisions. Lack of explainability may reduce user trust in automated systems. Additionally, privacy concerns related to large-scale facial data processing and potential demographic bias in training datasets must be carefully addressed to ensure responsible AI deployment.

Continuous Evolution of Deepfake Techniques – Deepfake generation methods are continuously evolving, leveraging improved neural architectures and training strategies. Detection models must constantly adapt to new manipulation techniques such as face reenactment, full-body synthesis, and cross-modal generation. Maintaining robustness against unseen forgery methods while ensuring model stability remains a persistent research challenge.

VIII.CONCLUSION

The project successfully implements a Hybrid CNN–Transformer based deepfake detection system capable of accurately identifying manipulated images and videos. The proposed system effectively combines Convolutional Neural Networks for extracting fine-grained spatial features with Transformer-based attention mechanisms for capturing global contextual relationships. By integrating local texture analysis with long-range dependency modeling, the model improves detection robustness against high-quality forgeries, varying illumination conditions, and compression artifacts. Experimental results demonstrate improved accuracy, precision, recall, and F1-score compared to standalone CNN or Transformer models. Overall, the system enhances digital media authentication, strengthens cybersecurity measures, and contributes to maintaining trust and integrity in online communication platforms.

XI.REFERENCE

1. **I. Goodfellow et al., “Generative Adversarial Nets,” in *Proceedings of the IEEE Conference on Advances in Neural Information Processing Systems (NeurIPS)*, 2014, pp. 2672–2680.**
2. **A. Vaswani et al., “Attention Is All You Need,” in *Proceedings of the IEEE Conference on Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998–6008.**

3. A. Dosovitskiy et al., “An Image Is Worth 16×16 Words: Transformers for Image Recognition at Scale,” in *International Conference on Learning Representations (ICLR)*, 2021.
4. Y. Li and S. Lyu, “Exposing DeepFake Videos by Detecting Face Warping Artifacts,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018, pp. 46–52.
5. B. Dolhansky et al., “The DeepFake Detection Challenge Dataset,” *arXiv preprint arXiv:2006.07397*, 2020.
6. A. Rössler et al., “FaceForensics++: Learning to Detect Manipulated Facial Images,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 1–11.
7. K. He, X. Zhang, S. Ren, and J. Sun, “Deep Residual Learning for Image Recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
8. M. Afchar, V. Nozick, J. Yamagishi, and I. Echizen, “MesoNet: A Compact Facial Video Forgery Detection Network,” in *IEEE International Workshop on Information Forensics and Security (WIFS)*, 2018.
9. D. Guera and E. J. Delp, “Deepfake Video Detection Using Recurrent Neural Networks,” in *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 2018, pp. 1–6.
10. H. Nguyen, F. Fang, J. Yamagishi, and I. Echizen, “Multi-Task Learning for Detecting and Segmenting Manipulated Facial Images and Videos,” in *IEEE International Conference on Biometrics (ICB)*, 2019.