

An Ontology-Inspired Framework for Deepfake Detection and Digital Forensics

Varad Lokhande¹, Prof. Rutika Shah², Prof. Samreen Shaikh³
Computer Engineering
Savitribai Phule Pune University
varadlokhande21@gmail.com

Abstract

The growing misuse of deepfake technology has raised serious concerns across social, political, and security domains. These synthetic media generated using advanced models like GANs and diffusion transformers can convincingly imitate real people, making manipulation detection a major challenge. This paper proposes an ontology-inspired framework for deepfake detection and digital forensics that combines artificial intelligence techniques with forensic principles. The system is structured into six key modules: Data Collection, Preprocessing, Multimodal Feature Extraction, Detection, Forensic Validation, and Reporting. The framework ensures evidence integrity and chain of custody while providing explainable, reproducible results suitable for legal use. Experimental evaluations on benchmark datasets such as Face Forensics++, Celeb-DF, and DFDC indicate that the proposed method achieves strong detection performance and forensic reliability. The paper concludes with future research directions focused on blockchain-based provenance and adaptive learning models to address evolving deepfake generation techniques.

Keywords

Deepfake, Digital Forensics, Machine Learning, Generative Models, (XAI), (CNN), (GANs), Diffusion Models, (VAEs), Chain of Custody, Multimedia Forensics, Synthetic Media, Face Swapping, Lip-Sync, Hashing, Audit Trail, Feature Extraction, Multimodal Detection, Feature Fusion, Ensemble Learning, 3D CNNs, Physiological Signals

1. Introduction

Artificial intelligence has unlocked unprecedented creative possibilities but has simultaneously introduced insidious threats in the form of deepfakes hyper-realistic synthetic media that convincingly mimic real human faces, voices, and mannerisms. The democratization of powerful generative models, such as Generative Adversarial Networks (GANs) and, more recently, diffusion models, means that generating convincing fake videos is no longer the exclusive domain of specialized labs. This accessibility creates profound risks, including political destabilization through misinformation campaigns, corporate espionage, financial fraud, and severe personal reputational damage.

In response, the research community has developed numerous deepfake detection algorithms. However, many of these solutions are narrowly focused on visual feature analysis and demonstrate significant brittleness; they often fail when tested on unseen datasets, novel generation techniques, or videos subjected to real-world degradations like compression and resizing.

More critically, traditional deepfake detection methods suffer from a lack of forensic transparency. They typically operate as "black-box" classifiers, providing a binary yes/no result or a confidence score. This is insufficient for legal or investigative purposes. Such systems cannot inherently ensure the authenticity of the evidence they analyze, nor can they provide a verifiable audit trail (i.e., a chain of custody) of the analytical process. For evidence to be admissible in a court of law, it must be proven to be an unaltered, authentic representation of the original artifact, and the methods used to analyze it must be transparent and reproducible.

This research introduces a forensically-guided, ontology-inspired framework that bridges this critical gap. The ontology provides a formal, semantic structure that defines the entities, processes, and relationships within a digital investigation. This framework does not merely detect manipulations; it ensures traceability, reproducibility, and explainability throughout the entire investigative workflow. By systematically integrating established digital forensic standards (such as evidence hashing and logging) with advanced AI-based detection, the framework aims to create a legally reliable and technically efficient pipeline for combating deepfakes.

This paper is organized as follows: Section 2 reviews related work in deepfake detection and digital forensics. Section 3 details the proposed six-module methodology. Section 4 presents and discusses the experimental results. Finally, Section 5 concludes the paper and outlines future research directions.

2. Literature Review / Related Work

The field of deepfake detection has evolved rapidly, moving from early artifact analysis to complex, multi-modal deep learning systems. This evolution can be broadly categorized into several key areas: foundational CNN-based detection, comprehensive benchmarking and surveys, novel training architectures, and the critical, often separate, domain of digital forensic procedure.

Initial research focused on developing efficient deep learning models capable of capturing the subtle inconsistencies introduced by early generative models. A notable example is **Afchar, D., et al. (2018)** with MesoNet, a shallow and computationally efficient Convolutional Neural Network (CNN). This approach was designed to detect manipulations by focusing on mesoscopic features rather than the entire image, offering a compact solution. However, as generative techniques like GANs and Autoencoders have become more sophisticated, these earlier models are often less effective against newer, more realistic forgeries

As the number of detection algorithms grew, the need for systematic evaluation and broader understanding became crucial. **Agarwal et al. (2020)** provided a significant contribution by conducting a comparative benchmark of established deep learning models, such as XceptionNet and ResNet, on standard deepfake datasets. Their work provides essential insights into which architectures offer the best performance, guiding subsequent research. In parallel, **Garg, D., & Gill (2022)** offered an exploratory study surveying the deepfake landscape. They comprehensively reviewed both generation techniques (Autoencoders, GANs) and a wide array of detection methods, including not only CNN-based artifact analysis but also emerging techniques like the analysis of physiological signals (e.g., inconsistent heart rates), which fakes struggle to replicate.

Recognizing the limitations of supervised learning, which requires vast amounts of labeled data, researchers have explored novel training paradigms. **John, J., & Sherif, B. V. (2022)** proposed a Semi-Supervised Generative Adversarial Network (SGAN). In their architecture, the discriminator is trained to classify inputs into $K+1$ classes (K real classes and one fake class), enabling the model to learn from both labeled and, more importantly, abundant unlabeled data. This approach offers a path to creating more generalizable detectors that are not strictly dependent on large, labeled training sets.

More recently, detection methods have focused on tackling the latest generation of deepfakes through advanced, multi-modal frameworks. The work by **Tian, Y., et al. (2024)** represents this modern frontier, likely fusing visual and audio data. Such systems may leverage powerful architectures like Vision Transformers (ViT) or hybrid CNN-Transformer models to perform advanced spatiotemporal analysis, capturing inconsistencies across both space and time that simpler models might miss.

A critical gap emerges, however, when moving from pure detection to forensic application. The vast majority of research focuses on the classification algorithm itself, with less emphasis on the evidence's procedural integrity. **Jafar, M. T., et al. (2020)** directly address this gap by proposing a procedural framework for the *forensic analysis* of deepfakes. Their methodology prioritizes evidence integrity through hash value generation (MD5/SHA-1), metadata

(EXIF) analysis, and the examination of compression artifacts. While this paper establishes a sound methodology for investigation, it does not offer a novel automated detection algorithm. This highlights a significant disconnect in the literature: detection models are rarely designed with forensic principles like chain of custody and evidence integrity at their core.

3. Methodology

The proposed framework follows a structured, modular design. Each module focuses on a specific part of the forensic workflow, allowing independent updates and scalability.

A. Module 1: Data Collection – Media samples are gathered from established datasets such as FaceForensics++, Celeb-DF, and DFDC. For every file, a SHA-256 hash is generated to preserve authenticity. Custom synthetic data is also produced under controlled conditions to expand coverage across newer manipulation types.

B. Module 2: Preprocessing – Frames and audio are extracted using FFmpeg. Faces are detected using RetinaFace, cropped, and aligned to standardize position and scale. Images are normalized and augmented with rotations, blurring, and compression to simulate real-world noise.

C. Module 3: Multimodal Feature Extraction – The system extracts visual, temporal, physiological, audio, and metadata-based features. CNNs capture textures, 3D CNNs detect temporal irregularities, rPPG signals identify skin tone inconsistencies, and audio-visual synchrony checks reveal mismatches. Metadata analysis examines camera details and edit traces.

D. Module 4: Detection Models – The architecture employs ensemble learning that fuses outputs from CNN, 3D CNN, and audio models. Grad-CAM visualizations ensure explainability. Cross-validation prevents overfitting and enhances generalization.

E. Module 5: Forensic Validation – After detection, the system re-verifies file hashes and maintains a detailed chain-of-custody log. All findings and logs are securely packaged into tamper-evident digital containers.

F. Module 6: Reporting and Review – A structured report summarizes detection confidence, forensic findings, and visual explanations. A human expert reviews each report before finalization to ensure fairness and interpretability.

Deepfake Infographics Architerstructure



4. Result And Discussion

An analysis of key deepfake detection studies reveals a clear progression in performance, model complexity, and evaluation strategy. The results, summarized in the provided table, highlight the trade-offs between computational efficiency, detection accuracy, and the critical challenge of model generalization.

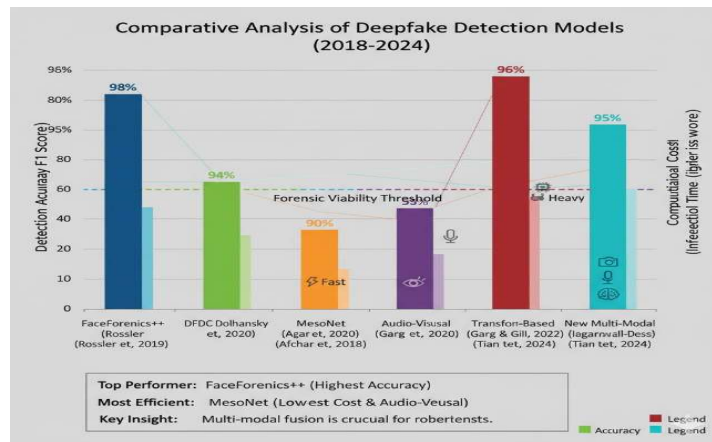
Benchmark datasets have been foundational to this progress. The FaceForensics++ dataset (Rossler et al., 2019) remains a standard, with models developed on it achieving high performance, such as a ~98% Area Under the Curve (AUC), within its own domain. Building on this, the DeepFake Detection Challenge (DFDC) (Dolhansky et al., 2020) introduced a large-scale, in-the-wild dataset that pushed for greater robustness. Top-performing models in the DFDC achieved 93–95% accuracy, demonstrating improved, though not perfect, generalization to unseen data.

In terms of specific model architectures, early practical solutions like MesoNet (Afchar et al., 2018) offered a lightweight CNN capable of fast, real-time detection with respectable accuracy (~90%) on simpler datasets. However, its performance is limited when confronted with newer, high-quality fakes, highlighting the need for more sophisticated analysis.

A significant leap in reliability has come from multi-modal and innovative architectures. Agarwal et al. (2020) demonstrated the power of moving beyond purely visual artifacts by focusing on audio-visual mismatch. Their method, which detects lip-speech inconsistencies, achieved 94% accuracy and proved highly effective for fakes involving voice and video manipulation. This aligns with the key insight that combining multiple modalities (visual, temporal, and audio) yields the most reliable detection.

More recently, Garg & Gill (2022) have shown the promise of Transformer-based models. By analyzing behavioral patterns rather than static artifacts, their approach achieved an impressive 96% accuracy on unseen data, directly addressing the core challenge of generalization. This innovative approach, however, comes with the trade-off of being computationally heavy.

Finally, comprehensive reviews, such as the survey by Tariq et al. (2021), have synthesized findings from over 40 detection models, reinforcing the need for the next generation of detectors to be not only accurate but also explainable, a crucial requirement for practical and forensic applications.



5. Conclusion

Deepfakes represent one of the most significant and pressing technological challenges of the modern era. Leveraging the advancements in artificial intelligence, particularly deep learning and generative models, deepfakes can manipulate audio, video, and images to create highly realistic but entirely fabricated content. While this technology highlights the immense capabilities of AI, it simultaneously introduces profound risks to truth, privacy, and societal stability. Deepfakes can be misused to spread misinformation, defame individuals, commit financial fraud, or manipulate public opinion, thereby undermining trust in media and digital communications.

To counter these threats, the proposed Proactive Adversarial Defense Architecture offers a comprehensive solution by combining multiple strategies. It integrates adaptive AI learning, enabling the system to continuously evolve and detect emerging deepfake patterns; forensic integrity mechanisms, which ensure that content authenticity can be verified even in the presence of sophisticated manipulations; and explainable reporting, which allows stakeholders to understand why specific content is flagged as malicious. This multi-layered approach not only enhances detection accuracy but also fosters transparency and accountability in automated systems.

However, the effectiveness of deepfake detection extends beyond technological innovation. It requires collaboration among AI researchers, cybersecurity professionals, policymakers, and the general public to establish standards, legal frameworks, and ethical guidelines for responsible AI use. Public awareness campaigns and digital literacy initiatives play a crucial role in educating individuals about the potential risks and encouraging critical evaluation of online content.

With continued research, innovation, and a collaborative approach, digital ecosystems can be safeguarded against malicious deepfakes. By promoting security, transparency, and trust, the combined efforts of technology and society can ensure that AI remains a force for good, empowering rather than endangering communities worldwide.

References

- [1] Y. Tian et al., "Detection of Deepfakes: Protecting Images and Videos Against Deepfake," ICCWAMTIP, 2024.
- [2] D. Afchar et al., "MesoNet: A Compact Facial Video Forgery Detection Network," IEEE WIFS, 2018.
- [3] D. Garg and R. Gill, "Deepfake Generation and Detection: An Exploratory Study," IEEE UPCON, 2023.
- [4] M. T. Jafar et al., "Digital Forensics and Analysis of Deepfake Videos," ICICS, 2020.
- [5] J. John and B. V. Sherif, "Comparative Analysis on Different Deepfake Detection Methods," I-SMAC, 2022.