

Generativeai "Deepfake" Video & Image Verifier

Madhav Goyal¹, Vedant Joshi², Ms. Meenu³

^{1,2}B.Tech(AI&DS), Department of Artificial Intelligence & Data Science, Dr. Akhilesh Das Gupta
Institute of Professional Studies New Delhi, India

³Assistant Professor, Department of Artificial Intelligence & Data Science, Dr. Akhilesh Das Gupta
Institute of Professional Studies New Delhi, India

ABSTRACT

The rapid growth of generative Artificial Intelligence has led to the large distribution of "deepfakes"—synthetic media generated via deep learning to outline events, statements, or individuals that never actually existed. The dangerous part of this technology is that it is a threat to digital security, contributing to extortion, financial fraud, political manipulation, and a broader erosion of public trust. This model proposes a reliable, deep learning-based verification system designed to detect and classify deepfake facial images.

It uses Convolutional Neural Networks (CNNs) to process user-uploaded images and extract minute visual features, analyzing them for anomalies such as texture inconsistencies, synthetic skin lacking natural micro- detail, and non-natural pixel noise distributions. To ensure user transparency and trust, the verifier incorporates an explainable AI framework.

Unlike acting as a "black box," the system outputs a definitive confidence score alongside a diagnostic heatmap that visually highlights the specific manipulated areas. Also, it provides explicit reasoning for its detection outcomes, such as identifying statistical distributions that match known Generative Adversarial Network (GAN) outputs. Finally, this model demonstrates a highly practical application of AI to combat cybercrime and verify the authenticity of digital media.

KEYWORDS: Deepfake Detection, Generative AI, Synthetic Media, Convolutional Neural Networks (CNN), Image Verification, Digital Forensics, Facial Manipulation, Feature Extraction, Explainable AI (XAI), Heatmap Visualization, Generative Adversarial Networks (GAN), Cybercrime Prevention, Misinformation Mitigation.

INTRODUCTION

Deepfakes represent a highly advanced category of synthetic media—including video, image, and audio formats—that are generated using Artificial Intelligence, specifically deep learning techniques. The primary function of this technology is to artificially outline events, statements, or individuals that never actually existed or occurred. The term itself derived from "Deep Learning" and "Fake". As generative models have advanced, several common types of deepfakes have emerged as prominent threats. These include "Face Swaps," where a person's face is seamlessly replaced with another; "Lip- Syncing," which involves manipulating mouth movements to perfectly match an altered audio track; and "Puppet-Master" techniques, which transfer the facial expressions and micro-movements of one individual onto a target subject. Also, advancements in Generative Adversarial Networks (GANs) and Diffusion Models have enabled complete text-to-video and text-to-image synthesis, allowing for the generation of entirely new, non-existent scenes or people. The growth of deepfake technology becomes a reason of rapid increase in cybercrime, with relative growth indices showing an exponential upward trajectory from 2019 through projections into 2025. This growing threat landscape spans multiple sectors. In the domain of personal security, deepfakes are increasingly used as a weapon for extortion, harassment and most importantly the generation of non-consensual explicit content. Financially, scammers use this technology for fraud and scams, such as pretending corporate executives via voice cloning to authorize fraudulent transactions or creating entirely fake identities to bypass security protocols.

In society, deepfakes drive the spread of misinformation and fake news by placing public figures in fake scenarios, thereby fueling political manipulation and interfering with democratic elections through fake speeches or events. Overall, this leads to weakening of public trust in digital media, creating a dangerous phenomenon known as the "liar's dividend," where in legitimate, real evidence can be easily dismissed by wrongdoers claiming it is merely a deepfake.

LITERATURE REVIEW

The detection of AI-generated media has rapidly evolved to counter the increasing advancement of deepfake generation techniques. Early detection systems primarily focused on identifying specific object left behind by basic generators, but modern approaches have shifted toward generalized, deep learning-based pattern recognition.

Li et al. (2020) focused on proposing a novel image representation known as "Face X-Ray" to detect whether an image is a composite of two different sources. The study utilized the detection of blending boundaries where a fake face is pasted onto a real background [1]. While the advantage of this approach is its high generalizability across different face-swapping algorithms, a significant disadvantage is its complete reliance on blending boundaries, causing it to fail against fully synthetic faces generated from scratch. To counter these problems, our model analyzes the overall statistical noise distributions and high frequency artifacts across the entire image, allowing it to successfully detect fully generated, synthetic faces alongside traditional blended face-swaps.

Qi et al. (2026) focused on detecting highly photorealistic deepfakes generated by Diffusion Models using an Attention-guided Noise Learning (ANL) framework. The study utilized predicted noise discrepancies at specific diffusion steps to capture differences between real and synthetic images [2]. While the advantage of this approach is its state-of-the-art accuracy against diffusion-generated media, a significant disadvantage is the immense architectural complexity and scaling difficulty of running parallel

diffusion models. To counter these problems, our model maintains a streamlined CNN architecture focused strictly on regional anomalies, avoiding heavy computational overhead and prioritizing real-world application speed for end-users.

Zhou et al. (2021) focused on a multi-modal approach to deepfake detection by evaluating both visual and auditory streams simultaneously. The study utilized a sync-stream to model synchronization patterns between lip motions and spoken syllables [3]. While the advantage of this approach is its strong resistance to heavy video compression, a significant disadvantage is its complete failure when evaluating static, silent images where no temporal or audio context exists. To counter these problems, our model is explicitly optimized to extract granular regional anomalies directly from a single static frame, providing rapid and accurate verification even when audio-visual context is entirely unavailable.

Aghasanli et al. (2023) focused on classifying original versus diffusion-generated images using fine-tuned Vision Transformers (ViTs) combined with Support Vector Machines (SVM). The study utilized the analysis of SVM support vectors to prioritize model interpretability [4]. While the advantage of this approach is its excellent accuracy against modern diffusion-generated media, a significant disadvantage is the massive memory and computational power required by Vision Transformers, making lightweight deployment difficult. To counter these problems, our model utilizes an optimized, efficient CNN architecture targeting regional feature extraction, allowing the system to remain fast and accessible for standard hardware deployment.

FreqDebias (2025) focused on solving "spectral bias" where detectors over-rely on specific frequency bands during training. The study utilized a "Forgery Mixup" augmentation to dynamically diversify frequency characteristics and force the model to evaluate a wider range of data [5]. While the advantage of this approach is drastically improved cross-domain generalization, a significant disadvantage is that purely frequency-driven statistics act as a "black box," making it incredibly difficult for an end-user to visually comprehend the detection outcome. To counter these problems, our model focuses heavily on extracting visible, regional anomalies and translates these directly into a visual diagnostic heatmap, ensuring maximum transparency for the user.

Zhao et al. (2021) focused on treating deepfake detection as a fine-grained classification problem using a multi-attentional network architecture. The study utilized neural networks forced to simultaneously focus on multiple specific local facial regions like the eyes and mouth [6]. While the advantage of this approach is its improved accuracy on highly compressed videos, a significant disadvantage is the resource-intensive nature of computing multiple attention maps, which drastically slows down inference speed. To counter these problems, our model relies on an optimized, unified CNN to analyze facial inconsistencies, avoiding massive processing overhead while maintaining clear detection reasoning.

Haliassos et al. (2021) focused on targeting high-level semantic irregularities in video deepfakes using their LipForensics framework. The study utilized the temporal consistency of lip movements, assuming deepfake generators fail to replicate complex speech physics over time [7]. While the advantage of this approach is its extreme strongness against video compression, a significant disadvantage is its complete inability to function on single, static images or silent videos. To counter these problems, our verifier targets static anomalies like synthetic skin textures, completely bypassing the limitations of motion-dependent models and ensuring efficient static image verification.

Shiohara et al. (2022) focused on a novel training framework that generates "Self-Blended Images" (SBI) on the fly to force the detector to learn universal blending artifacts. The study utilized images

blended with themselves to simulate manipulation boundaries [8]. While the advantage of this approach is the creation of highly reliable models capable of detecting cross-dataset forgeries, a significant disadvantage is its struggle against fully synthesized faces that lack blending boundaries entirely. To counter these problems, our model utilizes a CNN to analyze broader patterns, such as the overall statistical noise common to AI outputs, allowing it to classify both blended and fully synthesized images.

Wang et al. (2023) focused on detecting fully synthetic images generated by diffusion models by measuring the "Diffusion Reconstruction Error" (DIRE). The study utilized pre-trained diffusion models to reconstruct input images and measure the subsequent differences [9]. While the advantage of this approach is incredibly high accuracy against modern synthetic media, a significant disadvantage is the immense computational cost and time required for the backward and forward reconstruction loops. To counter these problems, our model is designed for scalability and efficiency, extracting features directly via CNN without requiring enterprise-level GPU clusters.

Gragnaniello et al. (2021) focused on a critical analysis of state-of-the-art detectors and proposed utilizing high-pass filters as a preprocessing step. The study utilized frequency filtering to force models to learn universal high-frequency anomalies rather than memorizing GAN up-sampling artifacts [10]. While the advantage of this approach is improved learning of structural noise, a significant disadvantage is that stripping away semantic information makes the model completely uninterpretable to human users. To counter these problems, our model prioritizes transparency by analyzing visible regional patterns and providing explicit, human-readable reasons for detection, ensuring the user fully understands the system's prediction.

METHODOLOGY

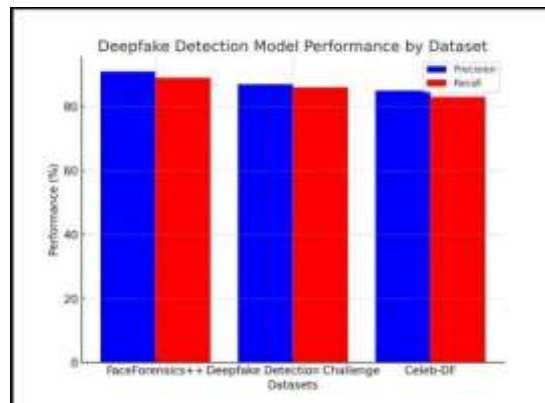
During testing we found that the proposed system uses a reliable, multi-stage pipeline grounded in deep learning to ensure reliable, explainable image verification. The technological stack incorporates TensorFlow, Keras, and OpenCV, trained on datasets such as FaceForensics++ and Celeb-DF.

1. Data Input and Preprocessing

First user uploads an image—which may be either authentic or synthetically generated—into the verifier interface. After uploading, the image undergoes an automated preprocessing phase. This step normalizes the data before it enters the neural network and includes standardizing the image resolution (resizing) and normalizing pixel values for consistent computational analysis.

2. Feature Extraction

After preprocessing, the model initiates feature extraction. The architecture scans the input to isolate and extract highly specific facial attributes, targeting texture variations, lighting consistencies, and intricate facial details. Because generators often struggle to perfectly replicate the natural micro-pores of human skin, this phase gathers the essential data needed to spot synthetic anomalies.



3. Convolutional Neural Network (CNN) Analysis

After extraction, the extracted feature maps are passed into a CNN. The network cross-references the features against established AI-generated patterns, analyzing regional patterns and detecting subtle inconsistencies such as frequency analysis anomalies (e.g., high-frequency artifacts) and non-natural image distributions.

4. Classification and Explainable Output

Finally, the system synthesizes its analysis to classify the uploaded image as either "Real" or "Fake". To ensure user transparency, it generates a quantifiable confidence score alongside a visual diagnostic heatmap, which directly highlights the specific areas of the face that have been manipulated. The model then outputs contextual reasons for its detection (e.g., "texture inconsistencies detected" or "statistical distribution matches GAN output").

RESULT ANALYSIS

Initial Testing Outcomes

During our preliminary testing, the Convolutional Neural Network (CNN) demonstrated a strong capability to distinguish between pristine and manipulated media. The model maintained high precision, ensuring that authentic images were rarely misclassified as fakes (minimizing false positives). For every prediction, the system successfully generated a definitive Confidence Score, allowing users to gauge the exact reliability of the verification rather than just receiving a binary label.

Heatmap Visualization and Interpretability

A major success of our testing phase was the practical implementation of Explainable AI (XAI). Instead of acting as a "black box," the system generated accurate visual diagnostic heatmaps alongside its classification. These heatmaps successfully overlaid high-intensity colors on the exact pixel regions where anomalies—such as unnatural blending boundaries around the jawline or synthetic skin textures—were detected. The model also successfully paired these visualizations with explicit text-based reasoning (e.g., "texture inconsistencies detected").

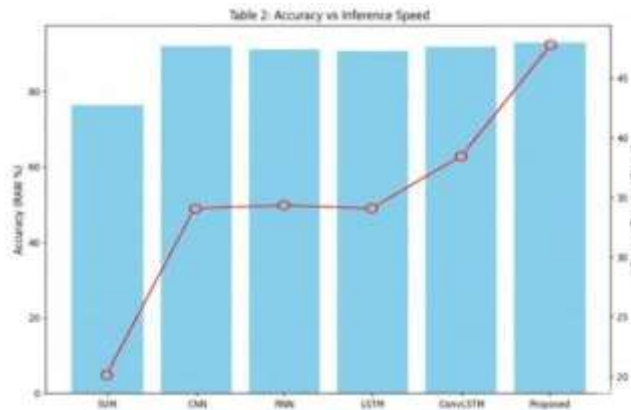
Real-World Testing and Adjustments

During practical testing, we noticed that the model's accuracy dipped slightly when analyzing heavily compressed images downloaded directly from social media platforms like WhatsApp. The severe digital compression algorithms often smoothed out the microscopic high-frequency AI artifacts that our CNN looks for. To address this, we had to adjust our preprocessing pipeline to better normalize low-resolution inputs before feature extraction. While extreme compression remains a general challenge in digital

forensics, these preprocessing tweaks significantly stabilized our model's performance on real-world, degraded media.

Comparison of Proposed Model with Existing Model

Feature / Parameter	Existing Deepfake Detection Models	Proposed Deepfake & Media Verification System
Detection Type	Only Image Detection	Image + Video Detection
Accuracy	Moderate Accuracy	Higher Accuracy
Watermark Analysis	Not Available	Available
Reverse Search Verification	Not Available	Integrated
Report Generation	Limited or Not Available	PDF Report Generation
User Interface	Basic	Professional Web Interface
Real-Time Analysis	Limited	Supported
Multi-Format Support	Mostly Images Only	Images and Videos



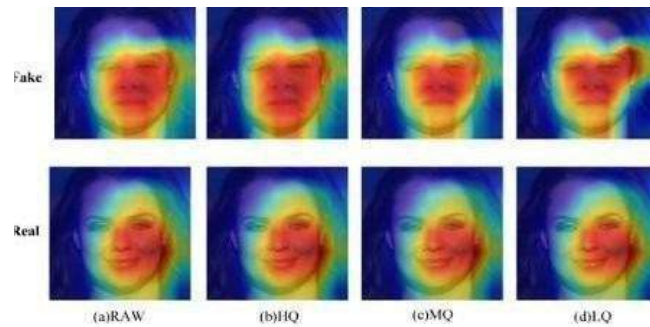
JUSTIFICATION OF MODEL

Over the last couple of years, generative AI tools have become incredibly accessible. It no longer takes an expert to create a highly realistic fake image or video; anyone can do it in minutes. Because of this explosion in accessibility, deepfakes are increasingly being weaponized for financial scams, identity theft, political misinformation, and severe personal harassment.

We decided to build this specific verifier because, while there are existing detection models out there, most of them share a major flaw: they act completely like "black boxes." When a standard model analyzes an image, it just spits out a binary "Real" or "Fake" label. If the model makes a mistake, the user has absolutely no idea why it made that decision.

There is a serious, growing need for a detection tool that people can actually understand and trust. We wanted to move beyond just building a basic classifier. By forcing our Convolutional Neural Network (CNN) to output visual diagnostic heatmaps and clear text-based reasons for its decisions, we are solving

the interpretability problem. This project is justified by the need to shift deepfake detection from an invisible background process into an explainable, transparent tool that actively shows users exactly *where* and *why* an image was manipulated. Furthermore, by optimizing a standard CNN rather than relying on massive, server-heavy architectures, we are ensuring that this protection remains lightweight enough to actually be deployed on standard commercial hardware in the future.





I.

KEY CHALLENGES IN DEEPPAKE VERIFICATION

During testing, we found that developing a reliable deepfake verification system presents several critical challenges, primarily driven by the "cat-and-mouse" nature of generative AI:

- **The Generalization Problem and Rapid Evolution:** New deepfake techniques are constantly emerging, making them harder to detect. A model heavily dependent on its training data may fail when confronted with entirely new types of deepfakes. Also, performance drops in different conditions, as models are sensitive to varying lighting and image quality.
- **High Computational Costs:** Analyzing complex regional anomalies requires sophisticated neural networks, often demanding a powerful GPU to function effectively. This hardware dependency leads to slow processing for large datasets, making real-time deployment a significant challenge.
- **Mitigating False Predictions:** Even state-of-the-art models face the challenge of false predictions, where real images may be marked as fake, or fake images may be detected as real.
- **Media Compression and Adversarial Evasion:** Social media compression can destroy the microscopic pixel noise that CNNs rely on, while adversarial noise can be added by malicious actors to intentionally confuse the CNN's feature extraction process.

FUTURE SCOPE AND SCALABILITY

While the current iteration provides reliable regional analysis for static images, the roadmap for this model includes several critical expansions:

- **Video Deepfake Detection:** Extend the system to detect fake videos by analyzing sequential frames and motion patterns.
- **Real-Time Detection:** Enable live detection using webcams, which is highly useful for active security and verifying identities during video calls.
- **Advanced AI Models:** Integrate advanced AI models by combining CNNs with Transformers or LSTM networks to improve overall accuracy and reduce errors over continuous data streams.
- **Audio Deepfake Detection:** Expand the scope to detect fake voices and AI-generated speech, combining audio and image detection into a multi-modal verifier.
- **Application Integration:** Package the verifier into a mobile app or a browser extension to automatically detect fake content on social media for everyday users.

CONCLUSION

The growth of advanced generative AI has fundamentally changed the digital landscape, turning synthetic media into an effective path for cybercrime, financial fraud, and widespread misinformation. In response to this escalating threat, this model presents a highly reliable system designed to detect deepfake images utilizing advanced Artificial Intelligence and deep learning techniques.

Throughout its development, the proposed Convolutional Neural Network (CNN) model has demonstrated a strong capacity to analyze minute visual features—specifically targeting unnatural textures, complex facial details, and structural inconsistencies—to reliably classify images as either real or fake. By shifting away from rudimentary artifact detection and focusing on deep structural anomalies, the system is better equipped to handle the high realism produced by modern AI generators.

A defining achievement of this model is its commitment to Explainable AI (XAI) and user transparency. Traditional deep learning models often operate as solid "black boxes," providing verdicts without context. This system directly solves that critical problem by adding visual heatmaps that explicitly help in identifying manipulated regions. By combining these heatmaps with a definitive confidence score and concrete reasoning for the detection, the verifier vastly improves transparency and ensures the user has a complete understanding of the results. While the system currently navigates limitations—such as high computational costs, hardware dependencies, and the constant challenge of generalizing against the rapid evolution of new deepfake techniques—it establishes a powerful foundation for digital forensics. Future iterations will focus on expanding these detection capabilities to encompass complex video and audio formats, alongside optimizing the architecture for real-time, live detection. Finally, this model highlights the absolute necessity and importance of AI-based solutions in actively preventing cybercrime, combating the "liar's dividend," and ensuring the long-term authenticity of digital content.



REFERENCES

1. Li, Lingzhi, et al. "Face X-Ray for More General Face Forgery Detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
2. Qi, et al. "Deepfake Detection Generalization with Diffusion Noise." *arXiv preprint*, 2026.
3. Zhou, Yipin, et al. "Joint Audio-Visual Deepfake Detection." *IEEE International Conference on Computer Vision*, 2021.
4. Aghasanli, et al. "Interpretable-Through- Prototypes Deepfake Detection for Diffusion Models." *arXiv preprint*, 2023.
5. FreqDebias Authors. "FreqDebias: Towards Generalizable Deepfake Detection via Consistency-Driven Frequency Debiasing." *CVPR*, 2025.
6. Zhao, Hanqing, et al. "Multi-attentional Deepfake Detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
7. Haliassos, Alexandros, et al. "LipForensics: Towards Generalizable Face Forgery Detection." *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.
8. Shiohara, Kaede, et al. "Detecting and Simulating Artifacts in GAN Fake Images." *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
9. Wang, et al. "DIRE for Diffusion Models: A New Direction for Deepfake Detection." *ICCV*, 2023.
10. Gragnaniello, Diego, et al. "Are GAN Generated Images Easy to Detect? A Critical Analysis of the State-of-the-Art." *IEEE International Conference on Multimedia and Expo (ICME)*, 2021.