

# Forensic Analysis Techniques for Deepfake Investigation

Aastha Bhandari<sup>1</sup>, Monika Chauhan<sup>2</sup>, Dr. Sachil Kumar<sup>3</sup>, Avni Goel<sup>4</sup>

<sup>1,2,4</sup>Post Graduate Student, Amity Institute of Forensic Sciences, Amity University, Noida,

<sup>3</sup>Assistant Professor, Amity Institute of Forensic Sciences, Amity University, Noida

**Abstract-**In an era marked by the rampant proliferation of deepfake technology, our world is grappling with unprecedented challenges in safeguarding the authenticity of digital content. This review paper endeavors to shed light on the intricate landscape of deepfake creation and dissemination, offering a comprehensive overview of forensic analysis techniques tailored specifically for investigating manipulated media. The paper will explore detection, verification, and attribution methodologies, exploring the underlying principles of both traditional and cutting-edge approaches. From machine learning and facial biometrics to image and video manipulation detection, audio analysis, and data provenance tracking, we dissect the arsenal of tools available for combatting the rising tide of deepfake threats.

To reduce bias a comprehensive literature search and objective assessment of the retrieved papers is performed. A search was conducted in the Google Scholar, PubMed, MEDLINE, Scopus, and Science Direct databases until 2022. The search keywords include various combinations like digital forensics, forensic science, deepfake technology, algorithms and artificial intelligence This systematic review critically reviewed 60 relevant studies.

**Keywords:** *Autoencoders, VAEs, Image Analysis, RNNs, Digital Evidence*

## 1. INTRODUCTION

Deepfake technology represents a form of digital evidence manipulation that leverages machine learning for the creation of lifelike videos or images [5]. It is becoming increasingly prevalent and accessible, with significant impacts on society, including political manipulation and reputational damage [1]. This technology involves using publicly accessible information to create fake videos or images that appear authentic[2]. This process entails training deep neural

networks on extensive datasets of images and videos to produce convincing yet fabricated content[5].

Within the realm of digital forensics, the rising prominence of deepfakes raises concerns due to their ability to manipulate or impersonate individuals in videos and audio files [5]. To achieve this, the algorithm undergoes a two-step process: training, where it learns from a substantial amount of data, often featuring the target individual's face or voice, and generation, where it produces the deceptive video or audio using the acquired model[5]. The authenticity threat posed by deepfakes is substantial, as these creations convincingly replicate specific expressions, movements, and speech patterns [5].

Detecting and identifying deepfakes demands specialized tools and techniques capable of revealing subtle artifacts or anomalies indicative of their presence [5]. The increasing sophistication and realism of deepfake technology make it challenging to distinguish between authentic and manipulated content, necessitating a proactive and collaborative approach in navigating the evolving landscape of digital forensics [5].

Existing methods for detecting deepfake content are not always accurate or consistent, especially when the content is low-quality, compressed, or edited [1]. Therefore, forensic analysis is crucial in detecting and combating deepfake videos [1][3]. Detection methods rely on analysing various features or artifacts of the images, videos, or audio [1]. However, these methods are not scalable or efficient, as they require a lot of computational resources and time to process large amounts of data [1]. Deepfakes are difficult to detect and may be used for advanced forms of phishing [2]. The prevalence of deepfake videos is increasing, and they have a significant impact on society by creating false narratives, spreading disinformation, and manipulating public opinion [3]. Fake audio or video

content created using deepfake technology is ranked as the most worrisome use of AI in terms of its potential applications for crime or terrorism [4].

It is important to have reporting tools or mechanisms provided by social media platforms, law enforcement agencies, or civil society organizations to report or flag deepfake content [1]. Consumers should also be aware of the existence and prevalence of deepfake content and use critical thinking and media literacy skills to identify and verify deepfake content [1]. Deepfake technology is a type of artificial intelligence that can create or manipulate images, videos, and audio that look and sound realistic but are not authentic [1]. Its potential risks to individuals and society include spreading misinformation, violating privacy, damaging reputation, impersonating identity, and influencing public opinion [1]. However, deepfake technology also has the potential to bring forth major business opportunities for content creation and engagement [4].

## 2. THREATS POSED BY DEEFAKE TECHNOLOGY IN VARIOUS DOMAIN

The potential threats posed by deepfake technology are vast and varied. One of the most significant threats is the erosion of public trust in digital media and information. As deepfakes become more prevalent, individuals may become increasingly skeptical and distrustful of any digital content, reducing the reliability of digital evidence and potentially impacting the justice system in the process [5]. Moreover, deepfakes can be created for malicious purposes, including non-consensual pornography and harassment [5]. This can lead to privacy concerns and potential harm to individuals [6]. In addition to personal media, distribution of deepfakes can also be a threat, as the use of deepfake technology can lead to the spread of misinformation and disinformation, manipulating public opinion and undermining trust in institutions [5]. Furthermore, deepfake technology can be used to commit fraud and identity theft, posing a threat to various domains [5]. The impact of deepfake technology needs to be combated through legal frameworks, government action, and technology industry responsibility [5]. Additionally, deepfake technology poses a growing threat to digital forensics and the ability to accurately identify and prosecute crimes based on digital evidence [5]. As the threat of

deepfake technology continues to grow, addressing its impact on various domains will require a concerted effort from all stakeholders involved.

## 3. DEEFAKE GENERATION TECHNIQUES

### 3.1 Generative Adversarial Networks (GANs):

GANs are commonly used in deepfake creation. It is a machine learning framework and have the ability to produce new data by using a pre-existing training set. GANs are utilised in the context of deepfakes to generate synthetic images, such as re-enactments and face swaps. The generator and discriminator neural networks are placed in direct competition with one another in order for GANs to function. Based on the information that the neural network has been provided, the generator creates a new picture as an output, and the discriminator assesses if the image is real or fraudulent. Both elements engage in continuous communication, with the discriminator learning how to spot deceptive images and the generator learning how to produce images that would fool the discriminator.

In addition to producing realistic photocopies, picture datasets, resolution improvement, video prediction, and more, GANs have demonstrated notable performance in a variety of computer vision tasks [7][8][13]

### 3.2 Autoencoders:

In order to encode an input picture or video into a lower-dimensional representation and subsequently decode it back to its original form, autoencoders in deepfake technology employ a neural network. Through this procedure, the autoencoder is able to identify the key elements of the input data, which can subsequently be edited to produce phony but realistic-looking pictures or films. By training the network on a sizable dataset of photos of the target individual, autoencoders are employed in the context of deepfakes to produce realistic facial forgeries. Then, using the trained autoencoder, new, synthetic photos or videos of the subject's face can be produced, which can be utilized to produce convincing deepfakes [14].

### 3.3 Recurrent Neural Networks (RNNs):

RNNs, especially long short-term memory (LSTM) networks, can capture temporal dependencies in sequences. This makes them suitable for deepfake

video generation, maintaining consistency over time. To identify and produce fake films, Recurrent Neural Networks (RNNs) are combined with other neural network architectures, such as Convolutional Neural Networks (CNNs). Recurrent neural networks (RNNs) are utilized in video analysis and synthesis to extract temporal dependencies and sequential information. For example, a study that demonstrated competitive results in deepfake detection<sup>[15]</sup> utilized CNN and RNN to record inter-frame and intra-frame information for detecting real or fake videos. A different study suggested a model that distinguished between actual and false videos using RNNs in addition to CNNs, and it performed better than earlier models<sup>[16]</sup>.

### 3.4 Variational Autoencoders (VAEs):

VAEs are generative models that learn a probabilistic mapping between input and output. Deep learning neural networks called variational autoencoders (VAEs) offer a unified framework for learning deep latent-variable models and related inference models. The encoder and decoder of a VAE work together to learn from the latent vectors and attempt to recover the input data. VAEs differ from regular autoencoders in that their latent space is continuous and that they produce two distinct length vectors—one representing the means and the other the standard deviations—instead of a single encoding vector with a set length. Applications for VAEs are numerous and include representation learning, semi-supervised learning, and generative modelling<sup>[1]</sup>.

### 3.5 Voice Synthesis:

Techniques like voice cloning involve using deep learning to mimic someone's voice. This can be combined with other deepfake techniques for more convincing results. The technology called voice synthesis, sometimes referred to as speech synthesis or text-to-speech, is used in deepfake generation to produce synthetic speech that closely mimics a target human voice. The way the technology functions is by employing computer algorithms to translate written material into spoken words. Nevertheless, the synthetic voices frequently have an artificial, robotic, and impersonal tone. In order to solve this problem, voice cloning technology has been created, enabling the creation of audio samples that are cloned and have a natural sounding voice based on a reference voice.

The field of voice cloning has given rise to methods of audio manipulation known as "deepfake speech," which can be used to produce believable audio recordings that mimic human speech. A number of techniques have been developed by researchers to identify deepfake speech, such as the use of multilayer perceptron neural networks, pathological features, and frameworks that assess the relationship between talking, breathing, and quiet sounds in an audio clip<sup>[1][2][3]</sup>.

## 4. TRADITIONAL FORENSIC METHODS USED FOR DEEPFAKE DETECTION AND ANALYSIS

### 4.1. Digital Forensics:

4.1.1. Metadata Analysis: Examining metadata such as timestamps, camera information, and editing history can provide clues about the authenticity of the content. Examining the information included in digital files is known as metadata analysis. This contains information on the file type, hash values, authentication signatures, authorship, creation timestamps, editing history, and compression artifacts. Finding abnormalities, assigning content to sources, and spotting possible manipulations in deepfake media all depend on the analysis of this metadata.

4.1.2. Error Level Analysis (ELA): ELA detects inconsistencies in compression levels across different parts of an image, helping identify areas that may have been manipulated. It is a forensic technique that examines compression error levels to identify image modifications. Forensic investigation of deepfake uses ELA to spot irregularities in photos that have been altered. Error levels in various places are compared, which helps identify possible changes. This method is essential for digital content authentication, particularly when it comes to deepfake detection. A succinct discussion of ELA's advantages and disadvantages contributes significant knowledge to the field of forensic analysis for deepfake identification.

### 4.2 Image and Video Analysis:

4.2.1 Exif Data Examination: Exif Data Examination is the process of examining metadata that is incorporated in photos and provides information about the location, date, and camera settings. In forensic examination of deepfakes, the authenticity of an image can be confirmed by looking at the Exif data. Examining metadata for irregularities or

inconsistencies is beneficial for deepfake detection since it aids in the identification of manipulated information. A thorough grasp of forensic methods in the context of deepfake analysis is aided by talking about the function and efficacy of Exif Data Examination.

4.2.2 Frame Analysis: Studying individual frames in a video can help identify anomalies or inconsistencies that may indicate manipulation. Examining individual frames from edited films or photos is called frame analysis. This method assists in locating abnormalities, artifacts, or inconsistencies that might be signs of the creation of deepfakes. Forensic analysts can identify changes and evaluate the content's overall legitimacy by thoroughly studying each frame. By using Frame Analysis in the debate, the paper improves its understanding of efficient forensic techniques for identifying and lessening the effects of deepfakes.

#### 5. AUDIO ANALYSIS:

Spectrogram Analysis: Analysing the spectrogram of an audio file can reveal patterns or artifacts that may indicate tampering or synthesis. Examining frequency and time-domain representations of audio sources is part of spectrogram analysis. It is essential for identifying audio content manipulation in deepfakes. Spectrogram analysis can identify anomalies, artifacts, or discrepancies that might point to artificial modifications. The paper's insights into thorough forensic methods for locating and validating deepfake content based on audio features are improved by include this technique in the discussion.

#### 6. BLOCKCHAIN TECHNOLOGY:

In blockchain Timestamping and Integrity Verification is used. Timestamping involves keeping track of the moment at which data was created. Timestamped and Integrity Verified by the blockchain guarantee a safe and unchangeable record of deepfake information. Data entries with timestamps allow for chronological tracking, which helps forensic analysts verify the veracity and integrity of digital material.

#### 7. CURRENT DETECTION TECHNIQUES AND MITIGATION STRATEGIES FOR DEEPPAKE CONTENT

To combat the potential negative effects of deepfake technology, researchers are exploring various detection and mitigation strategies. One such technique is to develop algorithms that can distinguish between real and fake content [6]. This involves analysing audio and visual elements of the content for inconsistencies that reveal the content to be deepfake. Researchers are experimenting with machine learning techniques to improve the accuracy of these algorithms. Another effective mitigation strategy is to educate the public about the existence of deepfakes and how to identify them [6]. This involves raising awareness about the potential risks associated with deepfakes and teaching individuals to scrutinize content before accepting it as truth. Current deepfake detection techniques involve analysing inconsistencies in audio and visual elements of the content [6].

For example, researchers have developed software that can detect facial inconsistencies in deepfake videos, such as unnatural eye movements or inconsistent head positions. Other techniques involve analysing subtle variations in the sound wave patterns of audio recordings to detect manipulation. By adopting these strategies, we can mitigate the negative impact of deepfake technology and maintain trust in digital media and information.

The emergence of deepfake technology poses a significant threat to digital forensics, as it allows for the creation of highly realistic yet fabricated videos and images. The manipulation or impersonation of individuals in digital media can lead to the erosion of public trust in information and media. As such, researchers are exploring various detection and mitigation strategies to combat the potential negative effects of deepfake technology. However, the increasing sophistication and realism of deepfakes make it difficult to differentiate between authentic and manipulated content, making it imperative to continue developing and refining detection techniques.

Current strategies involve analysing inconsistencies in audio and visual elements of the content, as well as subtle variations in sound wave patterns in audio recordings. The research also highlights the importance of maintaining trust in digital media and information and emphasizes the need for ongoing

research and collaboration to stay ahead of the constantly evolving technology.

This study recognizes the potential limitations and gaps in the current research and suggests future directions for further exploration. Overall, this research provides a valuable contribution to the understanding and mitigation of the emerging threat of deepfake technology in the field of digital forensics.

## 8. CONCLUSION

Deepfake technology has permeated every aspect of modern civilization, boosting its ability to disrupt and destroy numerous disciplines. In politics, fabricating statements or movies starring significant persons has the potential to affect public opinion, foment strife, and even provoke violence. The erosion of faith in political institutions and the integrity of democratic processes constitutes a serious danger to global peace and security. Moreover, the use of deepfake technology to manipulate diplomatic relations has the potential to increase tensions between countries and endanger world peace and security.

Deepfake generation techniques encompass a variety of advanced methods for producing convincing but fabricated content, ranging from images and videos to speech, such as Generative Adversarial Networks (GANs), Autoencoders, Recurrent Neural Networks (RNNs), Variational Autoencoders (VAEs), and Voice Synthesis. By manipulating and synthesising data in ways that resemble human behaviour and traits, these strategies take use of the capabilities of deep learning algorithms.

GANs, for example, excel at creating realistic pictures by comparing a generator to a discriminator in a competitive learning process. Autoencoders, on the other hand, can learn meaningful representations of incoming data, making them ideal for creating diverse and realistic false images. Similarly, RNNs can capture temporal relationships in sequential data, such as face movements or speech patterns, allowing the production of convincing.

## REFERENCE

[1] The Dangers of Deepfake Technology: Exploring the Potential Risks of AI-Generated Videos and Images. (n.d.) Retrieved January 21, 2024, from [hackernoon.com](https://hackernoon.com)

[2] Deepfakes: What are they, and why are they dangerous? (n.d.) Retrieved January 21, 2024, from [wyche.com](https://www.wyche.com)

[3] Deepfakes and scientific knowledge dissemination. (n.d.) Retrieved January 21, 2024, from [www.ncbi.nlm.nih.gov/pmc/articles/PMC10439167/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10439167/)

[4] Deepfakes: Deceptions, mitigations, and opportunities. (n.d.) Retrieved January 21, 2024, from [www.sciencedirect.com](https://www.sciencedirect.com)

[5] Deepfake Technology and its Impact on Digital Forensics. (n.d.) Retrieved January 20, 2024, from [elnion.com](https://www.elnion.com)

[6] Unmasking the Illusions: How Digital Forensics Takes on Deepfake Detection. (n.d.) Retrieved January 20, 2024, from [medium.com](https://www.medium.com)

[7] CVisionLab. Deepfake (Generative adversarial network) | CVisionLab [Internet]. CVisionLab. 2020. Available from: <https://www.cvisionlab.com/cases/deepfake-gan/>

[8] Anderson M, Anderson M. The future of generative adversarial networks in Deepfakes – Metaphysic.ai [Internet]. Metaphysic.ai -. 2023. Available from: <https://blog.metaphysic.ai/the-future-of-generative-adversarial-networks-in-deepfakes/>

[9] Preeti, Kumar M, Sharma HK. A GAN-Based model of deepfake detection in social media. *Procedia Computer Science* [Internet]. 2023 Jan 1;218:2153–62. Available from: <https://doi.org/10.1016/j.procs.2023.01.191>

[10] Capurso M. The deep Fake Farm: Generative Adversarial Networks (GAN) [Internet]. 2023. Available from: <https://www.linkedin.com/pulse/generative-adversarial-networks-gan-mario-capurso-ovwhf>

[11] Sblendorio D. How to build a Generative Adversarial Network (GAN) to identify deepfakes [Internet]. ActiveState. 2020. Available from: <https://www.activestate.com/blog/how-to-build-a-generative-adversarial-network-gan/>

[12] Papastratis I. Deepfakes: Face synthesis with GANs and Autoencoders | AI Summer [Internet]. AI Summer. 2020. Available from: <https://theaisummer.com/deepfakes/>

[13] Charlwood R. Understanding the different types of generative AI deepfake attacks | iProov [Internet]. iProov. 2024. Available from: <https://www.iproov.com/blog/generative-ai-attack-types-explained>

- [14] Negi S. Deep fake : An Understanding of Fake Images and Videos [Internet]. 2021. Available from: <https://www.semanticscholar.org/paper/c8d56db5939d8367c926aa6eb710310d5f3124be>
- [15] Al-Dhabi Y. Deepfake video detection by combining convolutional neural network (CNN) and recurrent neural network (RNN) [Internet]. 2021. Available from: <https://www.semanticscholar.org/paper/e7e23e58948f41d304b9eca17c21fcc629506406>
- [16] Albazoni AAM. DeepFake videos detection by using recurrent neural network (RNN) [Internet]. 2023. Available from: <https://www.semanticscholar.org/paper/a062a10ea2599b9abae3a0260fcb43f3726c5e64>
- [17] Acemoglu D, Laibson D, List JA. Equalizing superstars: The internet and the democratization of education. *Am. Econ. Rev.* 2014;104:523–527. Doi: 10.1257/aer.104.5.523. [CrossRef] [Google Scholar]
- [18] Patel, S., Chandra, S.K., & Jain, A. (2023). DeepFake Videos Detection and Classification Using Resnext and LSTM Neural Network. 2023 3<sup>rd</sup> International Conference on Smart Generation Computing, Communication and Networking (SMART GENCON), 1-5.
- [19] Das, A., Viji, K.S., & Sebastian, L. (2022). A Survey on Deepfake Video Detection Techniques Using Deep Learning. 2022 Second International Conference on Next Generation Intelligent Systems (ICNGIS), 1-4.
- [20] Yu, Peipeng, Zhihua Xia, Jianwei Fei and Yujiang Lu. “A Survey on Deepfake Video Detection.” *IET Biom.* 10 (2021): 607-624.
- [21] Khatri, N., Borar, V., & Garg, R. (2023). A Comparative Study: Deepfake Detection Using Deep-learning. 2023 13<sup>th</sup> International Conference on Cloud Computing, Data Science & Engineering (Confluence), 1-5.
- [22] Zi, Bojia, Minghao Chang, Jingjing Chen, Xingjun Ma and Yu-Gang Jiang. “WildDeepfake: A Challenging Real-World Dataset for Deepfake Detection.” *Proceedings of the 28th ACM International Conference on Multimedia* (2020): n. pag.
- [23] Feng D, Lu X, Lin X. Deep detection for face manipulation. In *Neural Information Processing: 27<sup>th</sup> International Conference, ICONIP 2020, Bangkok, Thailand, November 18–22, 2020, Proceedings, Part V* 27 2020 (pp. 316-323). Springer International Publishing.
- [24] Goodfellow I, Pouget-Abadie J, Mirza M, Xu B, Warde-Farley D, Ozair S, Courville A, Bengio Y. Advances in neural information processing systems. Curran Associates, Inc. 2014 Dec;27:2672-80.
- [25] Ciftci UA, Demir I, Yin L. How do the hearts of deep fakes beat? Deep fake source detection via interpreting residuals with biological signals. In *2020 IEEE international joint conference on biometrics (IJCB) 2020 Sep 28* (pp. 1-10). IEEE.
- [26] Abdulreda AS, Obaid AJ. A landscape view of deepfake techniques and detection methods. *International Journal of Nonlinear Analysis and Applications.* 2022 Mar 1;13(1):745-55.
- [27] Haewon Byeon, Mohammad Shabaz, Deep learning model to detect deceptive generative adversarial network generated images using multimedia forensic, Volume 113, 2024, (<https://doi.org/10.1016/j.compeleceng.2023.109024>)
- [28] Li J, Lei M. A brief survey for fake news detection via deep learning models. *Procedia Computer Science.* 2022 Jan 1;214:1339-44.
- [29] Barve Y, Saini JR, Rathod R, Gaikwad H. Multi-Modal Misinformation Detection: An Exhaustive Review. In *2023 7<sup>th</sup> International Conference On Computing, Communication, Control And Automation (ICCUBEA) 2023 Aug 18* (pp. 1-5). IEEE.
- [30] Chen B, Ju X, Xiao B, Ding W, Zheng Y, de Albuquerque VH. Locally GAN-generated face detection based on an improved Xception. *Information Sciences.* 2021