

# A Survey: Deep Fake Detection

**Vishvajeet Chandanshiv, Narayan Ekhande, Radha Lohar, Rahul V. Dagade**

Department of Computer Engineering

Smt. Kashibai Navale College of Engineering, Vadgaon, Pune, India

SPPU Pune, India

**Abstract:** *This research addresses the pressing need for effective deep fake detection in both image domains, employing advanced deep learning methodologies. Deep fakes, which encompass the manipulation and fabrication of digital content, pose significant challenges to the authenticity and trustworthiness of media. In this study, we explore the evolving landscape of deep fake creation and the persistent challenges it presents. By leveraging state-of-the-art neural networks, natural language processing techniques, we aim to develop detection systems capable of distinguishing genuine from manipulated content. Our investigation also delves into the dynamic nature of deep fake detection, where creators continuously adapt their techniques. Staying one step ahead in this ongoing arms race is crucial for maintaining the integrity of digital content. In conclusion, this research contributes to the ongoing efforts to combat deep fake-related challenges, preserving public trust in the veracity of digital media. The development of reliable deep fake detection systems for both images is essential in an era where the line between reality and manipulation is increasingly blurred. [2].*

**Keywords:** Deep fake, Detection, Deep Learning, Image, Authentication

## I. INTRODUCTION

In recent years, the rapid advancement of deep learning technologies has ushered in a new era of innovation and transformation across various industries. However, with the rise of sophisticated machine learning models, there is a growing concern about their misuse, particularly in the creation and dissemination of deep fakes. Deep fakes refer to manipulated videos or images that convincingly portray individuals saying or doing things they never did. As the boundary between reality and digital manipulation becomes increasingly blurred, the need for robust deep fake detection mechanisms becomes paramount. Deep Vision emerges as a crucial player in the quest to safeguard authenticity in the digital realm. [3]

## II. LITERATURE SURVEY

Deep fake Detection using Deep Learning. Et.al Prof. Aparna Bagde, Sakshi Fand, Kanchan Varma, Aditya Gawali. The rapid advancements in AI, machine learning, and deep learning technologies, which have given rise to new tools for manipulating multimedia. While these technologies have found legitimate applications in entertainment and education, they have also been exploited for malicious purposes.

Notably, high-quality and realistic fake multimedia content, known as Deepfake, has been used to spread misinformation, incite political discord, and engage in malicious activities like harassment and blackmail. Deepfake algorithms possess the unsettling ability to craft counterfeit images and videos so convincingly that they elude human scrutiny. These algorithms adeptly fashion deceptive visual and auditory content, manipulating the appearances and behaviour of targeted individuals to such a degree that viewers instinctively place their trust in what is, in fact, a fabrication. Distinguishing these deepfakes from genuine content becomes a formidable challenge, as the human eye struggles to discern the difference. In response, this paper conducts a comprehensive exploration, delving into the array of tools and algorithms employed in the creation of deepfakes, while placing particular emphasis on the vital aspect of deepfake detection methods. Through in-depth discussions that encompass challenges, research endeavour, technological advancements, and strategic approaches linked to the realm of deepfakes, this survey scrutinizes the landscape.

Deepfake Detection and the Impact of Limited Computing Capabilities. Et.al Paloma Cantero-Arjona, Alfonso S´anchez-Macian. The rapid development of technologies and artificial intelligence makes deepfakes an increasingly sophisticated and challenging-to-identify technique. To ensure the accuracy of information and control misinformation and mass manipulation, it is of paramount importance to discover and develop artificial intelligence models that enable the generic detection of forged videos. This work aims to address the detection of deep fakes across various existing datasets in a scenario with limited computing resources. The goal is to analyze the applicability of different deep learning techniques under these restrictions and explore possible approaches to enhance their efficiency.

Deep fake Video Detection Using Recurrent Neural Networks. Et.al David Guera Edward J. Delp. In recent months a machine learning based free software tool has made it easy to create believable face swaps in videos that leaves few traces of manipulation, in what are known as “deep fake” videos. Scenarios where these realistic fake videos are used to create political distress, black mail someone or fake terrorism events are easily envisioned. This paper proposes a temporal-aware pipeline to automatically detect deepfake videos. Our system uses a convolutional neural network (CNN) to extract frame- level features. These features are then used to train a recurrent neural network (RNN) that learns to classify if a video has been subject to manipulation or not. We evaluate our method against a large set of deep fake videos collected from multiple video websites. We show how our system can achieve competitive results in this task while using a simple architecture.

Deep fake Video Detection Based on Spatial, Spectral, and Temporal Inconsistencies Using Multimodal Deep Learning. Et.al John K. Lewis, Imad Eddine Toubal, Helen Chen. Authentication of digital media has become an ever pressing necessity for modern society. Since the introduction of Generative Adversarial Networks (GANs), synthetic media has become increasingly difficult to identify. Synthetic videos that contain altered faces and/or voices of a person are known as deep fakes and threaten trust and privacy in digital media. Deep fakes can be weaponized for political advantage, slander, and to undermine the reputation of public figures. Despite imperfections of deep fakes, people struggle to distinguish between authentic and manipulated images and videos. Consequently, it is important to have automated systems that accurately and efficiently classify the validity of digital content. Many recent deep fake detection methods use single frames of video and focus on the spatial information in the image to infer the authenticity of the video. Some promising approaches exploit the temporal inconsistencies of manipulated videos; however, research primarily focuses on spatial features. We propose a hybrid deep learning approach that uses spatial, spectral, and temporal content that is coupled in a consistent way to differentiate real and fake videos. We show that the Discrete Cosine transform can improve deepfake detection by capturing spectral features of individual frames.

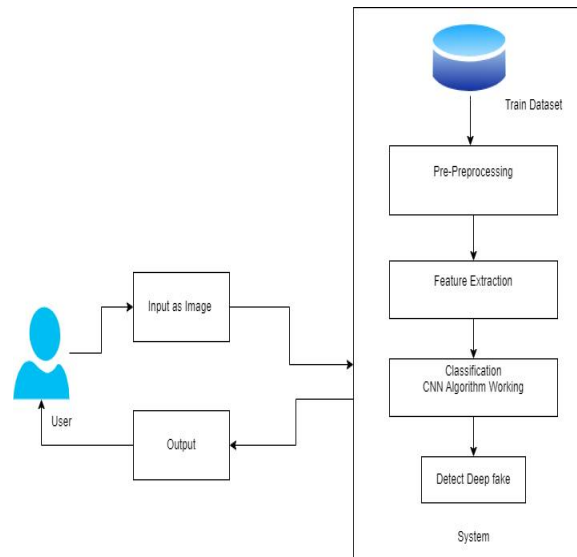
### III. METHODOLOGY

The methodology employed in this research project is structured to comprehensively address the goals of developing a deep fake detection system capable of identifying synthetic media, including images and text, generated through deep learning techniques. Data forms the foundation of our research. We collect diverse datasets comprising deepfake and genuine content in image formats. Data preprocessing includes steps to clean and standardize the data, ensuring it is suitable for analysis. For image data, preprocessing may involve resizing, normalization, and noise reduction.

We employ Convolutional Neural Networks (CNNs) to extract significant features from images. The trained CNN models identify patterns, textures, and structures within images, allowing for the distinction between deep fake and genuine content. Feature selection plays a vital role in enhancing the efficiency of the deep fake detection models. We explore various techniques for selecting relevant and discriminatory features, reducing dimensionality, and improving model accuracy.

Deep Fake Detection Using Deep are chosen for their effectiveness in pattern recognition and classification tasks. We develop deep fake detection models based on the features extracted from image data. The architecture and training of these models are conducted using suitable frameworks and libraries. The performance of the deep fake detection models is assessed using a range of evaluation metrics. We employ metrics such as accuracy, precision, recall, F1 score, and area under the receiver operating characteristic curve (AUC-ROC) to evaluate the models’ capabilities. The evaluation results are presented and analyzed, offering insights into the models’ performance, strengths, and limitations.

**IV. SYSTEM ARCHITECTURE**



In this section, we focus on the design of our deep fake detection system. We begin by describing the architecture and components that make up the system. The structure and contents of this chapter may vary depending on the project’s nature. The architecture of our deep fake detection system is designed to handle live camera. This module is responsible for processing and analyzing live camera to detect manipulated content. It includes sub- components for preprocessing, feature extraction using Convolutional Neural Networks (CNN).

**V. LIMITATION**

- As deep fake generation techniques improve, the synthetic media becomes increasingly sophisticated, making it harder to detect subtle manipulations. The adversarial nature of the problem means detection algorithms must constantly evolve to keep up with new generation techniques.
- A detection system trained on deep fakes generated by one model might fail to generalize to fakes generated by other models.
- As detection models improve, deep fake creators can specifically train their models to bypass current detection mechanisms, leading to an ongoing arms race between generation and detection.

**VI. DISCUSSION**

Deep fake detection has become an important area of research due to the increasing sophistication of AI-generated media. Deep fakes are digitally altered or artificially generated images that convincingly mimic real individuals, often leading to misinformation or harmful content. The detection of deep fakes involves analyzing various digital fingerprints, such as inconsistencies in lighting, facial expressions, or unnatural eye movements. Deep learning models, especially convolutional neural networks (CNNs), are commonly used to identify these subtle anomalies.

**VII. ANALYSIS**

The analysis of deep fake detection highlights both the advancements in technology and the challenges it presents. Deep fake detection relies heavily on Deep Learning techniques, with convolutional neural networks (CNNs) being a popular choice due to their ability to recognize patterns in images. These models can detect subtle irregularities, such as inconsistent facial movements, unnatural blinking patterns, or imperfections in lighting and shading, which are often overlooked by the human eye. Another effective technique is detecting discrepancies between facial expressions and emotions expressed in audio, as deep fakes may not perfectly synchronize the two.

### VIII. CONCLUSION

We have laid the groundwork for the development of deep fake detection models designed to identify synthetic media, encompassing images, generated through deep learning techniques. While the models themselves have yet to be constructed, our research has progressed through the initial phases, establishing a clear methodology and setting research objectives to guide our future endeavors. Our current progress involves an in-depth literature review, offering valuable insights into existing methods and the challenges associated with deep fake detection. In conclusion, while we are currently in the initial stages of our research, the methodology and research framework we have established are robust and well-defined. We anticipate the forthcoming phases with enthusiasm as we embark on the development and evaluation of the deep fake detection models, ultimately contributing to the preservation of trust and authenticity in the digital information ecosystem.

### REFERENCES

- [1]. Y. Patel et al., "An Improved Dense CNN Architecture for Deepfake Image Detection," in IEEE Access, vol. 11, pp. 22081-22095, 2023.
- [2]. M. S. Rana, M. N. Nobil, B. Murali and A. H. Sung, "Deepfake Detection: A Systematic Literature Review," in IEEE Access, vol. 10, pp. 25494-25513, 2022.
- [3]. D. Pan, L. Sun, R. Wang, X. Zhang and R. O. Sinnott, "Deepfake Detection through Deep Learning," 2020 IEEE/ACM International Conference on Big Data Computing, Applications and Technologies (BDCAT), Leicester, UK, 2020, pp. 134-143.
- [4]. Mary and A. Edison, "Deep fake Detection using deep learning techniques: A Literature Review," 2023 International Conference on Control, Communication and Computing (ICCC), Thiruvananthapuram, India, 2023, pp. 1-6.
- [5]. Analysis of Deepfake Detection Techniques | IEEE Conference Publication | IEEE Xplore
- [6]. Deepfake Detection: A Systematic Literature Review | IEEE Journals & Magazine | IEEE Xplore
- [7]. Deepfake Detection: Current Challenges and Next Steps (researchgate.net)
- [8]. Deep Fake Generation and Detection: Issues, Challenges, and Solutions | IEEE Journals & Page 4 of 5 Harini P., International Journal of Science, Engineering and Technology, 2023, 11:5 Magazine| IEEE Xplore
- [9]. Deepfakes Detection Methods: A Literature Survey | IEEE Conference Publication | IEEE Xplore
- [10]. Deepfake Detection through Deep Learning | IEEE Conference Publication | IEEE Xplore