



# Review On Deepfake Face Detection Using Machine Learning

**Prof. Soni Ragho<sup>1</sup>, Onkar Divekar<sup>2</sup>, Tushar Somwanshi<sup>3</sup>, Suraj Salgar<sup>4</sup>,  
Sarthak Shinde<sup>5</sup>**

<sup>1</sup>Asst.Professor Computer Engineering Vidya Prasarini Sabha's College of Engineering and Technology ,Lonavala

<sup>2</sup>Student Computer Engineering Vidya Prasarini Sabha's College of Engineering and Technology ,Lonavala

<sup>3</sup>Student Computer Engineering Vidya Prasarini Sabha's College of Engineering and Technology ,Lonavala

<sup>4</sup>Student Computer Engineering Vidya Prasarini Sabha's College of Engineering and Technology ,Lonavala

<sup>5</sup>Student Computer Engineering Vidya Prasarini Sabha's College of Engineering and Technology ,Lonavala

## ABSTRACT

As the sophistication of deep fake technology increases, so does the potential for misuse in various domains, including politics, entertainment, and social media. This paper presents a comprehensive survey of deep fake detection systems, focusing on the methodologies employed to identify manipulated audio-visual content. We categorize existing detection approaches into traditional and modern techniques, detailing their underlying algorithms, feature extraction methods, and performance metrics. Additionally, we explore the challenges faced in deep fake detection, such as the evolution of generative models and the need for real-time processing capabilities. The survey also addresses the ethical implications and societal impacts of deep fakes, highlighting the importance of developing robust detection mechanisms to mitigate misinformation. Finally, we identify future research directions, emphasizing the need for interdisciplinary approaches that combine advancements in machine learning, psychology, and legal frameworks to enhance detection accuracy and the field of deep fake detection.

**Keywords:** Generative Adversarial Networks (GANs), K-means Clustering, Face Swapping, Image Synthesis And Face Recognition, Convolutional neural networks ,Deep fake, Convolutional Neural Network(CNN), Fake Face Detection.

## 1.INTRODUCTION

In recent years, the advent of deep fake technology has revolutionized the landscape of digital media, enabling the creation of hyper-realistic audio-visual content that can convincingly imitate real individuals. This technology, primarily driven by advancements in artificial intelligence and machine learning, has found applications in various domains, including entertainment, advertising, and social media. However, the same capabilities that make deep fakes compelling tools for creativity also pose significant risks, such as misinformation, identity theft, and the erosion of trust in digital content.

Deep fakes are generated using sophisticated algorithms, notably Generative Adversarial Networks (GANs) and autoencoders, which learn to replicate the intricate details of human faces and expressions. This capability raises critical questions about the authenticity of digital media and the potential for malicious use. High-profile incidents, ranging from political misinformation to the creation of misleading celebrity content, have underscored the urgent need for robust detection mechanisms to combat the spread of deep fake materials.

This paper aims to provide a comprehensive survey of existing deep fake detection systems, categorizing them based on their underlying methodologies, effectiveness, and application contexts. We will explore traditional and modern detection techniques, examine the challenges faced by researchers in keeping pace with rapidly evolving deep fake technologies, and discuss ethical implications surrounding their use. Additionally, we will identify future research directions that could enhance the reliability and effectiveness of detection systems.

By consolidating knowledge from various studies, this survey seeks to equip researchers, practitioners, and policy makers with a nuanced understanding of deep fake detection, fostering a more informed discourse on the implications of this technology in society.

## 2.LITERATURE REVIEW

Ali Raza et al. introduced a novel deepfake predictor (DFP) methodology that combines the VGG16 architecture with a convolutional neural network architecture with a precision of 95% and an accuracy of 94% in the task of deepfake identification [1]. Younis E. Abdalla et al. proposes a transfer learning approach for training image forgery detection models using deep learning techniques for detecting image fraud with obtaining validation accuracy of 94.89% in the task of image modification detection [2]. Kumar et al. employed two methodologies for the detection of deepfake photos. This method involved the development of a customized convolutional neural network (CNN) using deep learning techniques. Another approach involved transfer learning leveraging the pre-trained models to enhance the detection of deepfake images [3]. Chang et al. incorporated image noise and augmentation to a VGG network resulting in a new network NA-VGG, which made a lot progress over other cutting edge fake image detectors [4]. Ensemble learning is a methodology which involves training numerous models on a common dataset and then their predictions are aggregated. The objective of ensemble learning is to enhance performance by combining multiple models than the performance

of any single model [5]. Qureshi et al. made use of six different base-learners, their predictions, along with the metadata are inputted into an SVM classifier that acts as the meta classifier. Where each baselearner has low accuracy, the meta-classifier outperforms them [6]. Sasikala et al. suggested an autonomous plant disease diagnosis method based on deep ensemble neural networks (DENN). Transfer learning approach is employed to reuse previously trained models. DENN outscored state-of-the-art pre-trained models by aggregation of different models, achieving an accuracy of 100% [7]. Sharma et al. presented a model using transfer learning technique from previously trained deep models like VGG16 and ResNet50 [8]. The proposed model is evaluated using three standard datasets. With the ensemble model reaching accuracies of 98.79%, 75.79%, and 95.52% on the three datasets, respectively, the overall performance is significantly improved [9]. Shad et al. used a number of techniques to identify deepfake photos and do a comparative study. With 99% accuracy, VGGFace prevailed over all the other models being looked into [10]. In a comprehensive assessment of the literature, Rana et al. observed that among four different methods, deep learning based approaches perform better than other methods in deepfake image detection, according to an evaluation of the performance of various methods with regard to different datasets [11]

### **3.OBJECTIVES OF PROJECT WORK :**

Project objectives include the development of a CNN model that accurately detects and labels deepfakes. As a result of this work, we demonstrate that the performance of CNN is outperformed by a semi-supervised learning approach. Using a subset of DeepFake Detection Challenge dataset, we evaluate our ResNet50 + LSTM based model. Due to the large data set, it was not possible to train the original dataset. We have taken a subset of the dataset, but kept the same splits and data as the whole dataset.

### **4.METHODOLOGY:**

Deepfake detection often relies on a two-step process: pinpointing faces of interest and then analyzing them for signs of manipulation. Inception-ResNet and MTCNN are deep learning models that can be combined for this purpose. MTCNN, short for Multi-task Cascaded Convolutional Networks, tackles the first step: face detection. It utilizes a series of increasingly complex convolutional neural networks, acting like filters that scan the image and progressively identify potential faces with higher accuracy. Once a face is located, Inception-ResNet steps in. This model excels at feature extraction, a crucial step in image recognition tasks. It analyzes the face in detail, extracting a unique mathematical representation called an "embedding." This embedding encodes information like facial landmarks, skin texture, and other subtle details. By comparing this embedding to a database of genuine faces, or by analyzing inconsistencies within the embedding itself (e.g., unrealistic smoothness or unnatural movements), Inception-ResNet can help determine if the face is real or a deepfake. The power of this approach lies in the synergy between the two models. MTCNN efficiently locates potential faces, while Inception-ResNet provides a robust analysis of the extracted facial features. Additionally, both

models are often pretrained on large datasets, allowing researchers to leverage their capabilities without the intensive computational cost of training from scratch. This makes Inception-ResNet and MTCNN valuable tools for streamlining deepfake detection research.

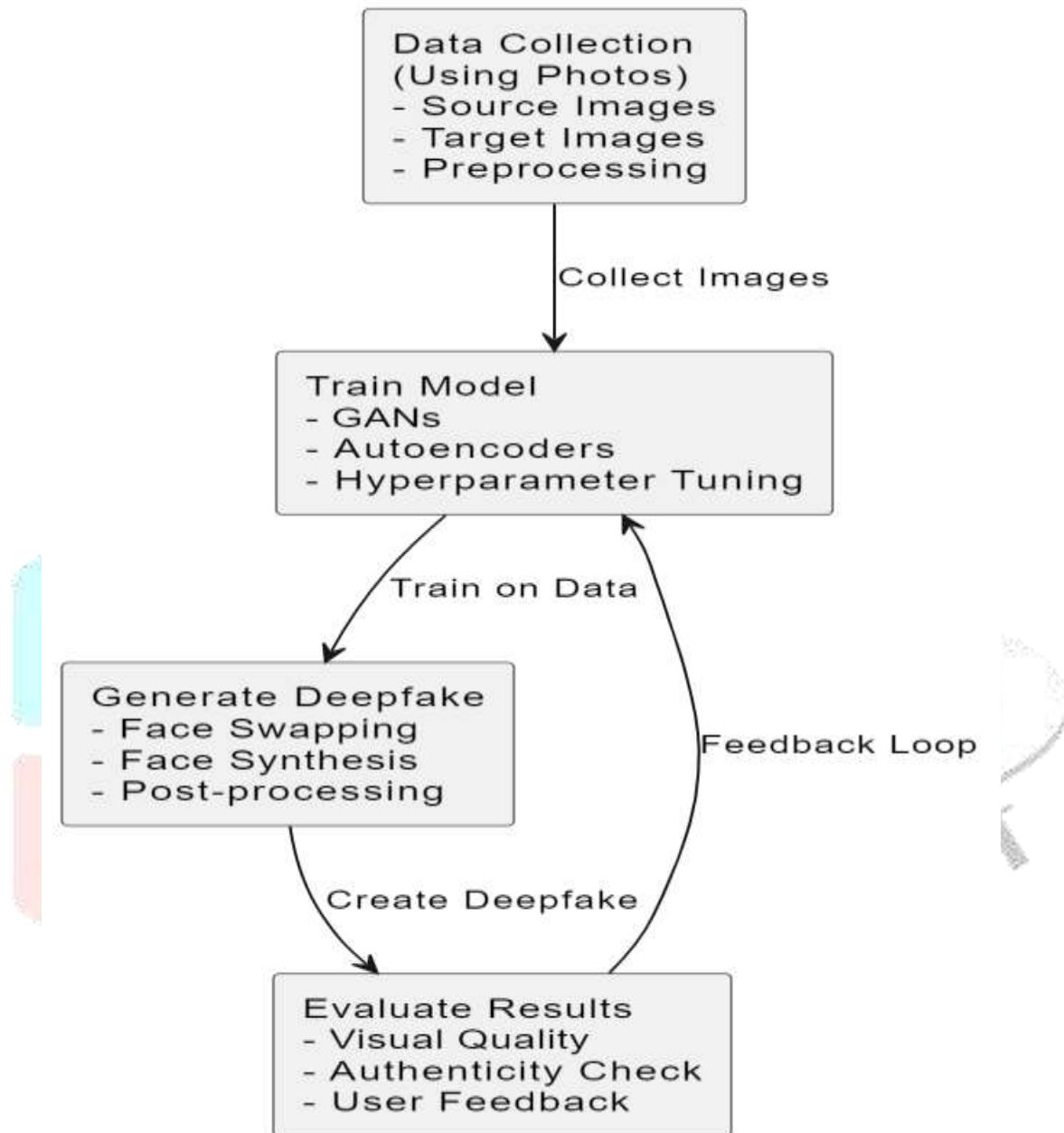
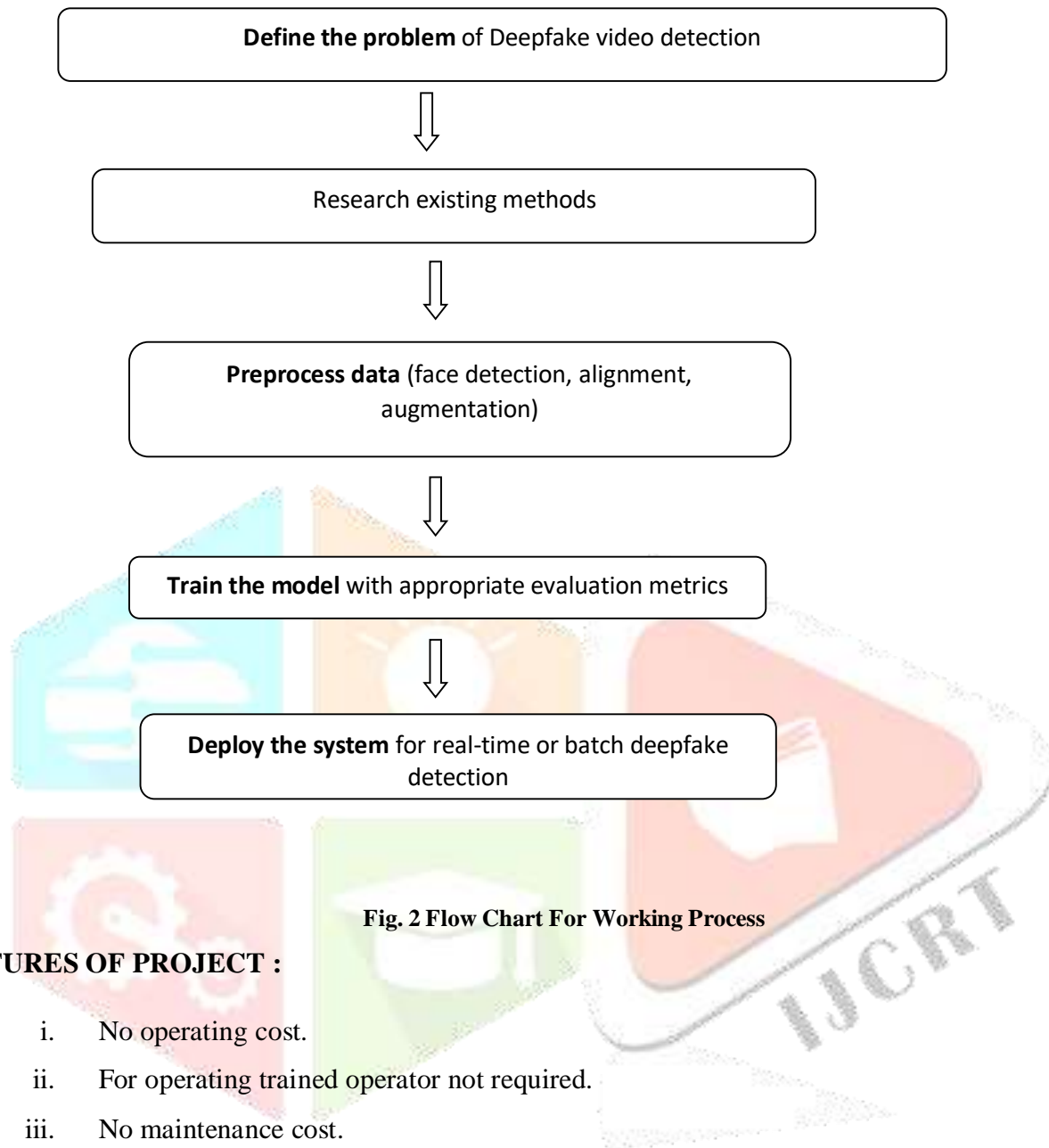


Diagram : 1 Methodology

**METHODOLOGY OF WORKING PROCESS****Fig. 2 Flow Chart For Working Process****FEATURES OF PROJECT :**

- No operating cost.
- For operating trained operator not required.
- No maintenance cost.
- Simple working.
- High Image peoces ingcapacity.

## ADVANTAGES

### 1. Protection Against Misinformation

- **Combat Fake News:** Deepfakes are often used to create misleading or fake content that can spread misinformation. A detection system can help identify and flag deepfakes, making it easier to differentiate between real and manipulated media.
- **Enhanced Public Trust:** By providing an automatic and reliable way to verify the authenticity of media, such a system can help build trust in news outlets, social media, and other content-sharing platforms.

### 2. Enhanced Security and Privacy

- **Prevent Fraud and Identity Theft:** Deepfakes can be used to impersonate individuals for fraudulent activities, such as identity theft, financial fraud, or impersonation in video calls. A detection system can help organizations and individuals identify such malicious activities and protect personal data.
- **Security in Authentication Systems:** Deepfake detection can be integrated into security systems (e.g., biometric authentication, facial recognition) to ensure that images or videos used for authentication are not manipulated.

### 3. Real-Time Detection

- **Instant Verification:** In scenarios such as live streaming or video conferencing, the system can provide real-time deepfake detection, alerting users or platforms instantly if manipulated media is being broadcasted.
- **Reduced Delay in Content Moderation:** Platforms like social media or video hosting websites can automatically filter out deepfake content before it goes viral, ensuring timely moderation of harmful media.

### 4. Support for Content Moderation

- **Automated Content Filtering:** Social media platforms, news agencies, and content hosting sites can use a deepfake detection system to automatically detect and remove harmful or fake content, saving human moderators time and resources.
- **Improved User Experience:** Users can rely on content moderation systems that actively identify and remove deceptive content, improving their browsing experience and preventing the spread of fake media.

## 5. Legal and Ethical Protection

- **Detect Malicious Use of Deepfakes:** Deepfake technology has been used maliciously in situations like revenge porn, political manipulation, and defamation. A detection system helps quickly identify such media and prevent its harmful spread.
- **Support in Legal Cases:** In cases of intellectual property theft, defamation, or misuse of manipulated media, a reliable deepfake detection system can provide crucial evidence in identifying and addressing the source of the fake content.
- **Preventing Deepfake-based Harassment:** Individuals can use detection tools to identify if their likeness has been used maliciously in manipulated media, allowing them to take necessary legal action.

### APPLICATIONS

- i. Social Media Platforms.
- ii. Journalism and News Organizations.
- iii. National Security and Law Enforcement.
- iv. Legal and Forensic Applications.
- v. Healthcare Industry.
- vi. Social Good and Awareness Campaigns.

### FUTURE SCOPE

There is still a lot of work that can be done in the area of deep fake detection using deep learning. One area of future work could be the development of more advanced deep fake detection models that are better equipped to handle the latest deep fake techniques. For example, new deep fake techniques that use generative adversarial networks (GAN) to create highly realistic deep fake videos have emerged and current deep fake detection models may not be able to detect these new types of deep fakes. Can. To address this, researchers can focus on developing more sophisticated deep fake detection models that are specifically designed to detect deep fakes produced by GANs as some of the following problems arise from quality deep fake detection models. There are very fundamental and major obstacles. Lack of large-scale high-quality training data: Deepfake detection relies on large amounts of high-quality training data to learn. However, collecting such data is a time-consuming and resource-intensive task. Dealing with new types of deepfakes: As deepfake technology evolves, new types of deepfakes are being developed that are more sophisticated and harder to detect. Continued research is needed to identify these new deepfake techniques and develop effective detection methods.

Adversarial attacks: Adversarial attacks refer to methods that manipulate deepfake detection models in such a way that they misclassify deepfakes as genuine.

## CONCLUSION

The rapid-fire development of deepfake technology presents significant challenges and pitfalls across colorful sectors, including media, security, politics, and entertainment. The Deepfake Face Detection System serves as a pivotal tool in addressing these challenges by offering a dependable means to identify manipulated content, icking the authenticity of visual media, and guarding individualities and associations from the dangerous goods of deepfakes. In this environment, the deepfake discovery system provides several crucial benefits, including the capability to combat misinformation, help fraud, enhance security, and save sequestration. It plays a vital part in content temperance on social media platforms, icking journalistic integrity, and securing legal and forensic operations. Likewise, the system supports public mindfulness and media knowledge, helping individualities and associations understand and alleviate the pitfalls posed by digital manipulation. The advantages of enforcing such a system are clear.

- Real- time discovery and scalability make it adaptable for different use cases, from social media platforms to public security operations.
- It offers a position of robotization that significantly reduces the time and cost involved in homemade content review, enabling further effective content filtering and temperance at scale.
- By icking ethical norms and contributing to the development of responsible AI, deepfake discovery systems help address one of the most burning ethical enterprises related to AI technology.

While the technology for detecting deepfakes has made significant advancements, it's important to admit that this is an ongoing challenge. Deepfake generators continue to ameliorate their styles, taking nonstop updates and advancements in discovery systems. Thus, the future of deepfake discovery will depend on the nonstop refinement of machine literacy models, the collaboration between tech companies and nonsupervisory bodies, and the education of the public to raise mindfulness about the pitfalls and signs of media manipulation. Eventually, the Deepfake Face Detection System isn't only an essential tool for vindicating digital content but also a critical step towards conserving trust and integrity in an decreasingly digital and connected world. With the growing significance of media authenticity and the implicit detriment posed by digital manipulation, the relinquishment and development of deepfake discovery technologies will remain a pivotal aspect of icking security, sequestration, and verity in the digital age.

**REFERENCES**

1. A.Raza, K.Munir, and M.Almutairi, "A novel deep learning approach for deepfake image detection," *Applied Sciences*, vol. 12, no. 19, p. 9820, 2022.
2. Y. Abdalla, M. Iqbal, and M. Shehata, "Image forgery detection based on deep transfer learning," *European Journal of Electrical Engineering and Computer Science*, vol. 3, no. 5, 2019.
3. N. Kumar, P. Pranav, V. Nirney, and V. Geetha, "Deepfake image detection using cnns and transfer learning," in *2021 International Conference on Computing, Communication and Green Engineering (CCGE)*. IEEE, 2021, pp. 1–6.
4. X. Chang, J. Wu, T. Yang, and G. Feng, "Deepfake face image detection based on improved vgg convolutional neural network," in *2020 39th chinese control conference (CCC)*. IEEE, 2020, pp. 7252–7256.
5. Patrick Schneider, and Fatos Xhafa, in *Anomaly Detection and Complex Event Processing over IoT Data Streams*, 2022
6. Qureshi, A.S., and Roos, T. (2021). Transfer Learning with Ensembles of Deep Neural Networks for Skin Cancer Detection in Imbalanced Data Sets. *Neural Processing Letters*, 55, 4461 - 4479.
7. Vallabhajosyula, S., Sistla, V., Kolli, and V.K. (2021). Transfer learning-based deep ensemble neural network for plant leaf disease detection. *Journal of Plant Diseases and Protection*, 129,545 - 558.
8. Sharma, J., Sharma, S., Kumar, V., Hussein, H.S., and Alshazly, H.A. (2022). Deepfakes Classification of Faces Using Convolutional Neural Networks. *Traitement du Signal*.
9. Shad, H.S., Rizvee, M.M., Roza, N.T., Hoq, S.M., Khan, M.M., Singh, A., Zaguia, A., and Bourouis, S. (2021). Comparative Analysis of Deepfake Image Detection Method Using Convolutional Neural Network. *Computational Intelligence and Neuroscience*, 2021.
11. Rana, M.S., Nobi, M.N., Murali, B., Sung, and A.H. (2022). Deepfake Detection: A Systematic Literature Review. *IEEE Access*, 10, 25494-25513.