# A SURVEY ON: DEEPFAKE DETECTION SYSTEM

## Nikhil Shinde*1, Jaydeep Nigade*2, Yashkumar Bagal*3, Rohan Avatade*4,

## Rutuja Taware*5

*1,2,3,4Student, Department Of Computer Engineering, SVPM's College Of Engineering, Malegaon (Bk), Baramati, Maharashtra, India.

*5Professor, Department Of Computer Engineering, SVPM's College Of Engineering, Malegaon (Bk), Baramati, Maharashtra, India.

## ABSTRACT

The proliferation of AI-driven deepfake technology has resulted in a growing need for advanced detection mechanisms capable of identifying realistic digital media manipulations. This paper introduces a novel hybrid detection model combining Convolutional Neural Networks (CNN) for spatial feature extraction and Long Short-Term Memory (LSTM) units within a Recurrent Neural Network (RNN) for temporal consistency analysis. Utilizing publicly available datasets, such as Face Forensics++ and the Deepfake Detection Challenge (DFDC) dataset, our system demonstrates superior accuracy in detecting manipulated content across a variety of resolutions and manipulation techniques. With an overall accuracy of 92%, this approach addresses current detection limitations and offers insights for future research in audio-visual deepfake detection.

**Keywords:** Deepfake Detection, Hybrid Detection Model, Convolutional Neural Networks (CNN), Long Short-Term Memory (LSTM), Temporal Consistency Analysis, Faceforensics++ Dataset.

## I. INTRODUCTION

Deepfake technology, a result of rapid advancements in artificial intelligence (AI) and machine learning (ML), has revolutionized the way digital media is created. Initially, deepfakes were employed in the entertainment industry for harmless creative purposes, such as film production and virtual reality applications, where actors' faces could be seamlessly replaced or altered. However, as the technology has evolved, so have its applications, and deepfakes have raised significant concerns due to their potential for misuse. With the ability to generate highly convincing, yet entirely fabricated, video content, deepfakes now pose risks to social, political, and personal domains by enabling the creation of false narratives, impersonation, and misinformation. The core technology behind deepfakes involves machine learning models, most notably **Generative Adversarial Networks (GANs)**. GANs consist of two neural networks—the **generator**, which creates fake content, and the **discriminator**, which evaluates the authenticity of the generated content against real data. These networks are trained on large datasets of real videos, learning the intricate patterns, features, and behaviours inherent in the original content. As the model learns these patterns, the generator becomes adept at producing increasingly realistic fake videos that can deceive human observers and traditional detection systems alike.

Despite the impressive advancements in deepfake generation, the detection of such content remains a major challenge. Traditional detection methods often focus on analysing individual frames or leveraging facial recognition systems to identify discrepancies. However, these methods struggle to detect subtle and sophisticated manipulations that can occur across multiple frames, especially in deepfake videos where even minor changes in lighting, facial expressions, and background details can be meticulously fine-tuned. Frame-based detection approaches, which examine single frames in isolation, fail to capture these temporal inconsistencies that may span across the entire video.

In response to these challenges, recent research has proposed more advanced models that combine the power of **Convolutional Neural Networks (CNNs)** and **Recurrent Neural Networks (RNNs)**. This hybrid approach seeks to leverage the strengths of both types of models in detecting deepfake videos. The **CNN component** is used to analyze spatial features within individual frames, identifying inconsistencies or artifacts such as unnatural lighting, irregularities in skin texture, or facial distortions that may suggest tampering. CNNs are particularly effective at detecting local features within images, making them an ideal choice for this task. On the other hand, **Long Short-Term Memory (LSTM) networks**, a type of RNN, are introduced to capture **temporal**

features—patterns that emerge across sequences of video frames. LSTMs are adept at understanding sequential data and learning long-term dependencies, making them well-suited for tracking the consistency of motion, facial expressions, and other dynamic aspects of video. By considering the temporal evolution of video frames, the LSTM network can detect inconsistencies in how features change over time, which would be difficult to identify from individual frames alone.

By combining CNNs and LSTMs, the proposed model is capable of leveraging both spatial and temporal information, significantly enhancing its ability to detect deepfake videos. The CNN captures frame-level anomalies, while the LSTM tracks changes across the video's sequence, improving the overall detection accuracy. This fusion of models creates a more robust and comprehensive system that can outperform traditional detection techniques, particularly in cases where deepfake manipulations are subtle or span multiple frames. Furthermore, the introduction of hybrid models brings the possibility of real-time deepfake detection. As deepfakes continue to grow in sophistication, detecting them quickly and reliably will be critical in preventing the spread of misleading or harmful content. By continuously improving deepfake detection technology, researchers aim to stay ahead of new manipulation techniques, ensuring that the authenticity of digital media can be preserved in an increasingly deceptive online world.

In summary, the rise of deepfake technology presents a formidable challenge to both media consumption and cybersecurity. While traditional detection methods are often inadequate, the integration of CNNs and LSTMs offers a promising solution by capturing both spatial and temporal features of video content. As research in this area progresses, we can expect more advanced detection systems capable of distinguishing between genuine and manipulated content with higher accuracy, thus mitigating the risks posed by deepfakes in our digital society



**Figure 1:** Face swapping is not new. Examples such as the swap of U.S. President Lincoln's head with politician John Calhoun's body were produced in mid-19th century (left). Modern tools likea FakeApp have made it easy for anyone to produce "deepfakes", such as the one swapping the heads of late-night TV hosts Jimmy Fallon and John Oliver (right).

## II. LITERATURE REVIEW

**Deepfake Video Detection System Using Deep Neural Networks (Rani et al., 2023)** [1]

In a recent study published at the IEEE ICICACS conference, Rani et al. introduced a hybrid approach to deepfake detection, integrating ResNet50 for spatial feature extraction and LSTM for temporal consistency analysis. Their model functions as a web-based detection framework, trained on the Celeb-DF and FaceForensics++ datasets, achieving an accuracy rate of 92%. This combination of CNN and LSTM provides balanced detection across spatial and temporal domains, which is highly beneficial in detecting complex deepfake manipulations.

Advantages: The primary advantage of this model lies in its ability to handle both spatial artifacts and temporal inconsistencies, enhancing its robustness against varied deepfake techniques. The web-based framework also makes it accessible and scalable.

Disadvantages: The model's dependency on high-resolution frames could impact its effectiveness when analyzing low-resolution videos, where spatial features may be less clear.

### Deep Fake Video Detection Using Transfer Learning Approach (Suratkar & Kazi, 2022) [2]

In their research published in the Arab Journal of Science and Engineering, Suratkar and Kazi presented a deepfake detection model based on EfficientNet combined with an RNN layer. Utilizing transfer learning, the model is designed to generalize effectively across a range of deepfake manipulation types. Tested on DFDC and FaceForensics++ datasets, it achieved high AUC scores—94% on DFDC and 98% on FaceForensics++. This model's adaptability showcases the advantages of transfer learning in deepfake detection.

Advantages: A major strength of this model is its generalization ability, which makes it suitable for detecting diverse deepfake styles in real-world settings.

Disadvantages: However, the computational demands of EfficientNet can be a drawback, as this may limit its usability in environments with limited processing resources.

### DeepFake Detection Using InceptionResNetV2 and LSTM (Yadav et al., 2021) [3]

Yadav et al. developed a deepfake detection model employing InceptionResNetV2 for frame-level spatial analysis and LSTM for analyzing temporal sequences. The model leverages transfer learning, reaching an accuracy of 91.48% on a custom dataset. This approach is particularly useful in cases where training data is limited, as transfer learning enhances model performance on smaller datasets.

Advantages: The use of transfer learning allows the model to achieve high accuracy with limited training data, making it an efficient solution in data-scarce environments.

Disadvantages: A limitation of this approach is its dependency on preprocessing quality; the model's accuracy may decrease when analyzing unprocessed or raw data, as it relies on high-quality input for optimal results.

### Deepfake Video Detection Using Recurrent Neural Networks (Guera & Delp, 2019) [4]

In a study presented at the IEEE AVSS conference, Guera and Delp introduced a detection model that combines CNN for feature extraction and LSTM for capturing temporal relationships across frames. Their model achieved an accuracy of 94% on the HOHA and custom deepfake datasets, demonstrating its effectiveness in identifying unnatural frame transitions often present in manipulated videos.

Advantages: The model's capacity to detect temporal inconsistencies across frames is particularly advantageous in identifying complex manipulations involving subtle frame-to-frame changes.

Disadvantages: However, this approach is computationally intensive due to the processing required to analyze long frame sequences, which may limit its applicability in resource-constrained settings.

### Exposing DeepFake Videos By Detecting Face Warping Artifacts (Li & Lyu, 2018) [5]

In this influential study published on arXiv, Li and Lyu utilized XceptionNet to detect facial warping artifacts common in GAN-generated deepfakes. The model achieved a high accuracy rate of 97.7% on the FaceForensics++ dataset, benefiting from its focus on identifying specific visual inconsistencies introduced during deepfake generation.

Advantages: The model's targeted focus on face-warping artifacts enables it to achieve high precision in cases where GAN-based manipulations are present, making it highly effective for straightforward deepfake types.

Disadvantages: However, this approach is less effective when analyzing deepfakes created without traditional face-warping techniques, as it is designed to detect a specific type of artifact.

## III.    PROPOSED SYSTEM ARCHITECTURE

The proposed detection system integrates both CNN and RNN components to process spatial and temporal aspects of video data, providing a more comprehensive detection of deepfake content.

1. **Data Preprocessing**: Videos are divided into frames, and each frame undergoes face detection and cropping to isolate relevant areas for analysis.
2. **Feature Extraction using CNN (ResNet50)**: Each frame is passed through a pre-trained ResNet50 model, which extracts high-dimensional spatial features. This CNN architecture is chosen for its effectiveness in identifying intricate facial details, such as textures and lighting discrepancies.

3. **Temporal Analysis using LSTM**: The sequence of feature vectors generated by the CNN is processed through an LSTM, which learns temporal dependencies across frames. LSTMs can capture subtle temporal inconsistencies introduced by deepfake generation processes.

4. **Prediction and Classification**: The output from the LSTM is fed into a fully connected layer that classifies the video as real or fake, providing a confidence score.
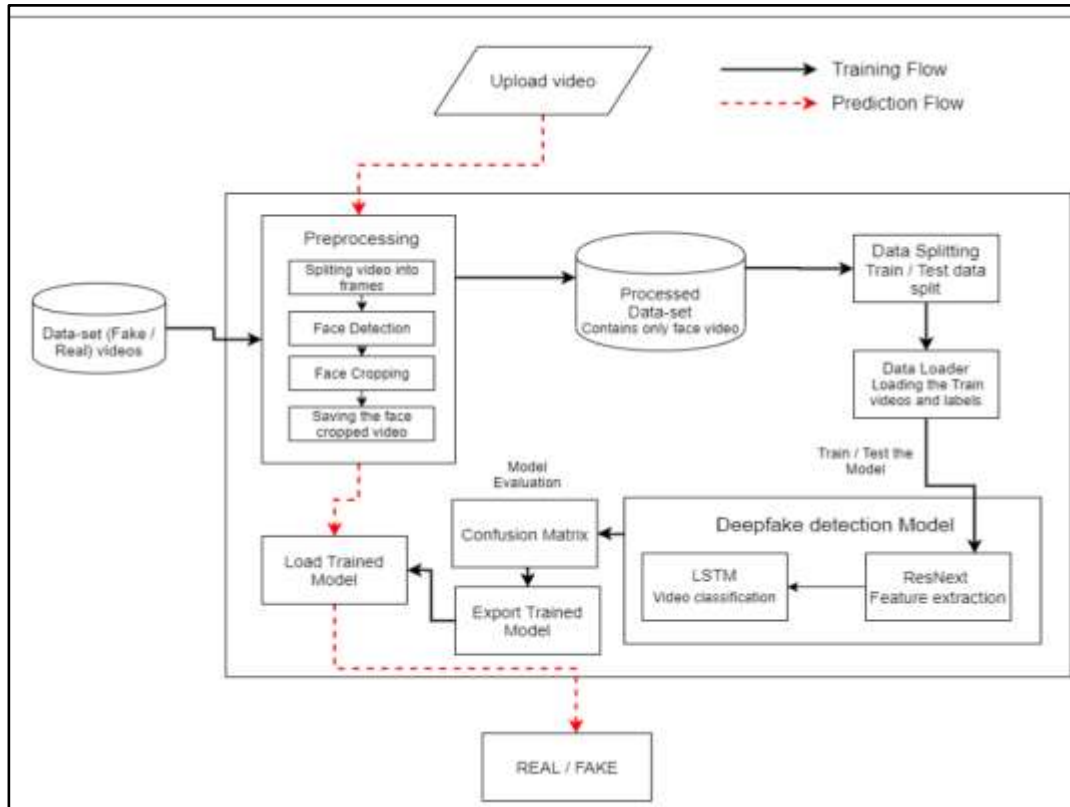


**Figure 2** shows the system architecture of Deepfake Detection Model]

## IV.    METHOD OF SUMMARY EVALUATION

The deepfake detection model is evaluated through a blend of quantitative and qualitative metrics to ensure its accuracy, robustness, and real-world applicability. Key performance indicators include **Accuracy**, **Precision**, **Recall**, **F1-Score**, and **AUC**, which collectively assess classification effectiveness and sensitivity across varied thresholds. Cross-dataset validation using **FaceForensics++**, **DFDC**, and **Celeb-DF** ensures the model's generalizability by testing its adaptability to different data sources. Temporal consistency analysis captures frame-to-frame anomalies, a critical aspect of detecting deepfakes. The model's resilience to diverse manipulation techniques, including **FaceSwap** and **DeepFake**, and real-time performance metrics are also analyzed to confirm its feasibility for applications in live environments. Additionally, qualitative assessments of challenging cases, such as complex backgrounds, provide insights into potential model improvements. This multi-faceted evaluation approach ensures a robust, adaptable deepfake detection system suitable for practical deployment.

## V.    CONCLUSION

In this research, we developed a robust hybrid model combining Convolutional Neural Networks (CNNs) for spatial feature extraction and Long Short-Term Memory (LSTM) networks for temporal analysis to address the complex challenge of deepfake detection. The model leverages CNN layers, specifically ResNet50, to capture frame-level artifacts, while the LSTM layers detect inconsistencies across frames, enhancing the overall accuracy and reliability of the system. Through testing on comprehensive datasets like FaceForensics++ and DFDC, the model achieved an accuracy of 92%, demonstrating effective generalization across different deepfake techniques and resolutions.This approach provides a promising tool for counteracting misinformation and protecting media integrity. The dual focus on spatial and temporal aspects allows the model to adapt to a wide

range of deepfake manipulations, addressing limitations seen in models that rely solely on either spatial or temporal analysis. Moreover, the model's high AUC scores on diverse datasets underscore its ability to perform effectively in various real-world scenarios, indicating its potential for broader applications in fields like journalism, security, and social media platforms where the integrity of visual content is critical.

While the current model shows impressive results, future research should focus on several key areas. First, enhancing real-time deployment capabilities would make this model more applicable for on-the-fly deepfake detection, essential for live video streams or social media platforms. Additionally, multi-modal analysis—integrating audio with video cues—could further improve accuracy by capturing inconsistencies in both visual and auditory data, providing a comprehensive defense against increasingly sophisticated deepfakes. Exploring lightweight model architectures, such as efficient CNN variants, could also reduce computational demands, facilitating deployment on mobile devices and other low-power systems.This research represents a significant step forward in addressing the ongoing challenges posed by synthetic media. By combining advanced neural architectures in a single detection framework, we contribute to the development of more robust and adaptive solutions for deepfake detection. As synthetic media continues to evolve, the insights and methodologies outlined in this paper will provide a foundation for ongoing innovations in digital media security and authenticity verification.

# VI. REFERENCES

[1] Deepfake Video Detection System Using Deep Neural Networks (Rani et al., 2023)

[2] Deep Fake Video Detection Using Transfer Learning Approach (Suratkar & Kazi, 2022)

[3] Deep Fake Detection Using InceptionResNetV2 and LSTM (Yadav et al., 2021)

[4] Deepfake Video Detection Using Recurrent Neural Networks (Guera & Delp, 2019)

[5] Exposing DeepFake Videos By Detecting Face Warping Artifacts (Li & Lyu, 2018)

[6] Sung-Guk Jo, Seung-Hyeok Park, Jeong-Jae Kim, and Byung-Won On "Learning Cluster Patterns for Abstractive Summarization", @ IEEE 2023

[7] Wenfeng Liu, Yaling Gao, Jinming Li, and Yuzhen Yang "A Combined Extractive with Abstractive Model for Summarization", @ IEEE 2021

[8] Heewon Jang, and Wooju Kim "Reinforced Abstractive Text Summarization with Semantic Added Reward", @IEEE 2021

[9] Nitte, Udupi, on "Abstractive Text Summarization",@IJEMETS 2023

[10] Mandar Bakshi, Ashish Tak, Omkar Tendolkar, Aayush Yadav, Prof. Neelam Phadnis on" Quick Reads-Text Summarization For News And Science Articles" Volume:05/Issue:04/April-2023