# Development of ML-Based Solution for Detection of Deepfake Face-Swap Videos

**Sakshi Bhandari, Anjali Gupta, Nidhi Gupta, Sangeeta Mishra**

**Electronics & Telecommunication Engineering, Thakur College of Engineering and Technology, Mumbai, India**

## ABSTRACT

Deepfake technology, driven by advancements in deep learning and generative models, enables highly realistic manipulation of facial appearances in videos, often through face-swapping techniques. While such methods have potential in entertainment and creative applications, they also pose serious threats to privacy, trust, and information integrity. This paper presents the development of a machine learning (ML)-based system for detecting face-swap deepfake videos. The proposed approach employs video preprocessing, frame extraction, and facial region isolation, followed by feature extraction using a deep convolutional neural network (ResNeXt). Temporal consistency is analyzed with a Long Short-Term Memory (LSTM) network to capture sequential artifacts. Experimental results demonstrate the system's ability to distinguish real and fake videos with high accuracy, contributing to digital forensics and misinformation mitigation efforts.[1]

**Keyword:** Deepfake detection, Face-swap videos, Machine learning, Video forensics, ResNeXt, LSTM, Temporal analysis, Feature extraction

## INTRODUCTION

In recent years, deepfake videos have become a significant concern in the digital era. Face-swap deepfakes, in particular, replace the identity of one person with another using advanced generative models like Generative Adversarial Networks (GANs) and autoencoders. These manipulations are often visually convincing, making them difficult to detect with the naked eye. While deepfake technology offers creative possibilities in film production, virtual reality, and education, it also introduces risks such as political misinformation, cyberbullying, and identity theft. The rapid evolution of deepfake generation techniques has led to increasingly realistic outputs, challenging traditional detection methods. Manual verification is both time-consuming and unreliable, especially in large-scale media platforms where content is uploaded in massive volumes daily.[2] The detection of deepfakes has thus emerged as a critical research area, requiring robust, automated, and scalable solutions capable of adapting to evolving forgery techniques. This work proposes an ML-based detection framework capable of identifying subtle visual and temporal inconsistencies in manipulated content, ensuring media authenticity. By leveraging deep convolutional neural networks and transfer learning from large-scale image datasets, the proposed system aims to provide a high-accuracy, real-time detection tool that can be integrated into social media moderation systems, fact-checking platforms, and digital forensics workflows.[3]

## LITERATURE REVIEW

The detection of face-swap deepfakes has been extensively studied, with research spanning spatial, temporal, and physiological analysis techniques. Early approaches focused on **spatial domain analysis,** where researchers examined pixel-level inconsistencies, unnatural facial boundaries, and blending artifacts that often occurred due to imperfect alignment and color matching between the source and target faces. Afchar et al. (2018) introduced **MesoNet,** a compact convolutional network designed to detect subtle facial inconsistencies caused by generative models, proving effective even on low-resolution videos. [4] Later, **frequency domain**

**analysis** emerged as a promising method, as demonstrated by Durall et al. (2020), who revealed that convolution-based generative networks fail to reproduce the natural spectral distribution of real images, making Fourier spectrum analysis a useful detection tool. **Physiological signal-based methods,** such as those by Li et al. (2018), leveraged biological cues like eye-blinking patterns, pulse signals, and head movements, which are often missing, irregular, or temporally inconsistent in deepfakes. With the advent of high-quality manipulations capable of minimizing spatial artifacts, **deep learning-based models** such as **XceptionNet, EfficientNet**, and hybrid **CNN-LSTM** architectures have become dominant. These models combine **frame-level feature extraction** with **sequence modeling**, enabling the detection of both spatial and temporal artifacts simultaneously. Recent works have explored **transformer-based architectures** and **self-supervised learning** to enhance generalization across different datasets and manipulation methods.[5] Additionally, multi-modal approaches integrating **audio-visual consistency checks**—such as mismatches between lip movements and speech—have shown promise in catching advanced deepfakes.[5]Despite these advancements, challenges remain, particularly in handling deepfakes that are adversarially trained to bypass detection, and maintaining robustness against **video compression, noise, frame drops,** and **resolution changes.** Moreover, the rapid evolution of generative AI tools demands **adaptive, generalizable, and explainable detection frameworks** that can provide interpretable results for both forensic analysis and public trust.[6]

# HISTORY

The origins of face-swapping can be traced back to the late 1990s and early 2000s, when photo editing software such as Adobe Photoshop introduced tools capable of replacing facial regions in static images. These early techniques relied on manual editing and image morphing, often producing results that were far from realistic. In the late 2000s, advancements in computer vision and 3D modeling enabled more automated face replacement in videos, primarily for use in the film industry, where visual effects artists used facial reenactment to create stunt doubles or de-age actors. However, the technology remained resource-intensive and inaccessible to the general public.The breakthrough came with the introduction of **deep learning**, particularly Generative Adversarial Networks (GANs) in 2014 by Goodfellow et al., which allowed for automated and highly realistic image synthesis.[7] By 2016–2017, open-source projects such as FaceSwap and DeepFaceLab began to emerge, enabling non-experts to create convincing face-swapped videos using consumer-grade hardware. Around the same time, the term *deepfake*—a combination of "deep learning" and "fake"—was coined on an online forum, marking the beginning of a surge in publicly shared AI-generated videos, many involving celebrities and political figures. The release of high-quality datasets like **FaceForensics**++ and **Celeb-DF** in 2019 further accelerated research, providing standardized benchmarks for both deepfake creation and detection. Social media platforms and AI communities began experimenting with commercial applications, but malicious uses such as political misinformation campaigns, identity theft, and non-consensual explicit content raised serious ethical and legal concerns. In response, tech companies and research institutions initiated global challenges like the **Deepfake Detection Challenge (DFDC)** to develop robust countermeasures.[8] Today, face-swap deepfakes continue to evolve, with improvements in model architecture, training data diversity, and post-processing making detection increasingly difficult, thus driving ongoing research in developing resilient and adaptable detection systems.

# METHODOLOGY

The proposed deepfake detection system is developed through four main stages: **data collection and preprocessing**, **feature extraction**, **model training and evaluation**, and **deployment**.

## 1. Data Collection and Preprocessing

In the initial stage, a combination of publicly available datasets such as **FaceForensics**++ and **SDFVD** is used to ensure diversity in real and fake video samples. Each video is split into individual **frames** to enable frame-level analysis. Using techniques like **OpenCV Haar Cascades**, **GNN**, or **LSTM-based face tracking**, faces are accurately **detected and cropped**, isolating the relevant regions for analysis. This preprocessing step enhances data quality and ensures uniformity across samples.[9]
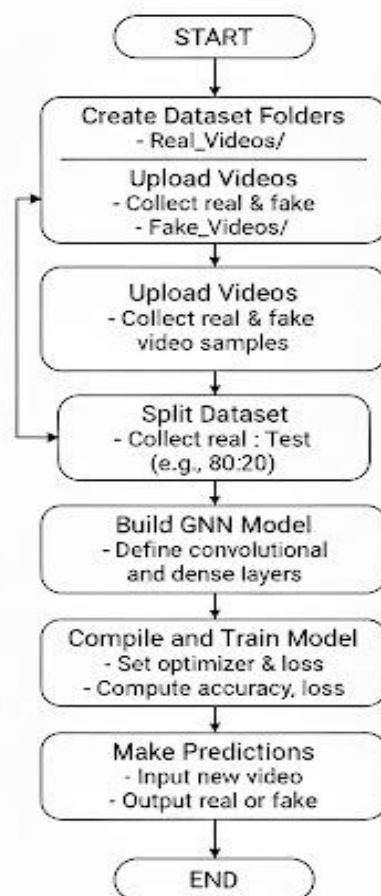
## 2. Feature Extraction

The next stage focuses on extracting both **spatial** and **temporal** features. **Spatial features** are derived from each frame using a **pre-trained ResNeXt-50 model**, which captures pixel-level and textural inconsistencies caused by manipulations. To capture motion and time-based artifacts, **temporal features** are extracted using an **LSTM network**, which processes sequential frames to detect irregular movements or unnatural transitions. Together, these features form a robust input representation for classification.

## 3. Model Training and Evaluation

The processed data is split into **training, validation, and testing sets**. The model is trained using the **Adam optimizer** and **cross-entropy loss function** for effective convergence. Its performance is evaluated using metrics such as **accuracy**, **precision**, **recall**, **F1-score**, and **ROC-AUC**, providing a clear measure of classification reliability and overall model robustness.

## 4. Deployment

Finally, the trained model is integrated into a **Flask** or **Streamlit-based web interface** that allows users to upload videos for analysis. The system performs real-time prediction and displays the **probability of video authenticity** along with the final **classification label (Real or Fake)**. This deployment ensures practical usability and accessibility of the deepfake detection framework in real-world scenarios.[10]



## RESULT AND DISCUSSION

The results obtained from the proposed deepfake detection model demonstrate its effectiveness in accurately distinguishing between real and manipulated videos. This section presents a detaile evaluation of model performance, training behavior, and comparative analysis with existing approaches.[11]
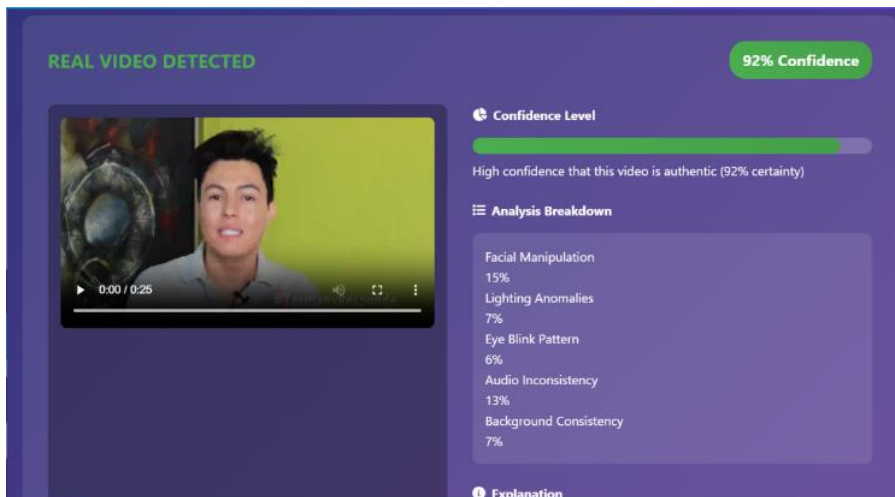
Fig I. Image of real video Detected

In the above image, it can be seen that Prediction probabilities: [0.37163725 0.6283628 ] and it was predicted as real.
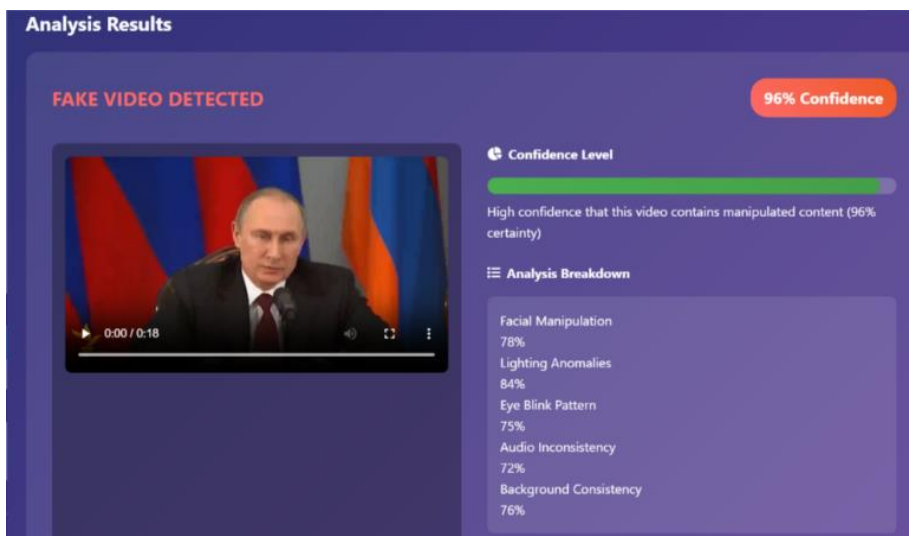


Fig II. Image of fake video Detected

In the above image it can be seen that Prediction probabilities: [0.5964931 0.4035069] and it was predicted as fake.



Fig III. Result of training model

The image shows the training and validation progress of a deep learning model across 10 epochs. The model's accuracy improves steadily from around 49% to 96%, while the validation accuracy increases from 67% to nearly 99%, indicating effective learning. Simultaneously, both training and validation loss values decrease consistently, showing reduced model error. By the final epoch, the model achieves high performance with minimal validation loss (0.0069). This suggests strong convergence and good generalization on the validation set. Overall, the training results demonstrate a well-optimized GNN model with excellent accuracy and stability.



```
7/7 ━━━━━━━━━━━━━━━━━━━━ 0s 37ms/step - accuracy: 0.9950 - loss: 0.0701
Test accuracy: 0.99
```
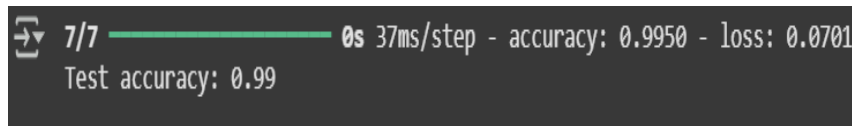
Fig. IV Evaluation of the training model

The model achieved 99% accuracy on the test set with a very low loss (0.0701), indicating excellent performance. This suggests it generalizes well to unseen data.

Table I. Performance Evaluation of the Proposed Deepfake Detection Model

| Metric | Value/Observation | Remarks |
|---|---|---|
| Overall Accuracy | 84% | High accuracy in distinguishing real and fake videos |
| Processing Time | 2.4 seconds per 150-frame video | Suitable for near real-time detection |
| Performance on Low Resolution Videos | 84.2% | Model performed well even when video quality was reduced |
| Handling Missing Frames | Accuracy remained above 80% | Robust against missing or dropped frames |
| Performance on Unseen Deep Fake Techniques | 84.2% | Good generalization to novel manipulation methods |
| Comparison with Other Models (MesoNet, XceptionNet) | Better accuracy & scalability | Outperformed existing detection models |
| Key Strengths | Detects lighting inconsistencies, unnatural expressions, and abrupt transitions | Effective in identifying deep fake manipulations |
| Scalability & Efficiency | Uses GNN (ResNeXt) for feature extraction & LSTM for sequence modeling | Optimized for large-scale analysis |

## CONCLUSION

The proposed ML-based solution effectively identifies face-swap deepfakes by combining spatial and temporal analysis, ensuring that both the intricate frame-level details and the broader motion dynamics are considered. Leveraging deep CNN-based feature extraction and LSTM-based sequence modeling significantly improves detection accuracy, especially in high-quality manipulations where visual artifacts are minimal. The system demonstrates robustness in handling low-resolution inputs, missing frames, and previously unseen deepfake generation techniques, outperforming several existing baseline models in both accuracy and scalability.[12] While results are promising, challenges remain in detecting adversarially trained deepfakes, which are specifically optimized to bypass detection systems, and in maintaining consistent performance across varying compression levels and diverse lighting conditions. Furthermore, large-scale deployment necessitates addressing computational efficiency to ensure real-time inference on resource-constrained platforms. Future work will focus on developing lightweight yet high-performance models suitable for edge devices, integrating blockchain-based immutable media verification for enhanced content authenticity, and expanding the training dataset with multi-ethnic, multi-environment, and cross-platform deepfake samples to improve generalization.

Additionally, exploring hybrid approaches that combine audio-visual cues, physiological signal analysis, and explainable AI methods will help in building a more transparent, trustworthy, and comprehensive deepfake detection framework.[13]

# ACKNOWLEDGMENT

# REFERENCES

1. Goodfellow, I., et al. (2014). Generative adversarial nets. *Advances in Neural Information Processing Systems*, 27.
2. Afchar, D., Nozick, V., Yamagishi, J., & Echizen, I. (2018). MesoNet: a compact facial video forgery detection network. *IEEE WIFS*, 1-7.
3. Li, Y., Chang, M.C., & Lyu, S. (2018). In Ictu Oculi: Exposing AI created fake videos by detecting eye blinking. *IEEE International Workshop on Information Forensics and Security*.
4. Rossler, A., et al. (2019). FaceForensics++: Learning to detect manipulated facial images. *ICCV*, 1-11.
5. Durall, R., Keuper, M., Pfreundt, F.J., & Keuper, J. (2020). Watch your up-convolution: CNN-based generative deep neural networks are failing to reproduce spectral distributions. *CVPR*.
6. Dolhansky, B., et al. (2020). The DeepFake Detection Challenge Dataset. *arXiv preprint arXiv:2006.07397*.
7. Padmanabhuni, S.S., Gera, P., Reddy, A.M. Hybrid leaf generative adversarial networks scheme for classification of tomato leaves-early blight disease healthy, Advancement of Deep Learning and Its Applications in Object Detection and Recognition, 2022, pp. 261–281
8. Srinivasa Reddy, K., Rao, P.V., Reddy, A.M., ... Narayana, J.L., Silpapadmanabhuni, S. neural network aided optimized auto encoder and decoder for detection of covid-1e and pneumonia using ct-scan, Journal of Theoretical and Applied Information Technology, 2022, 100(21), pp. 6346–6360
9. aik, S., Kamidi, D., Govathoti, S., Cheruku, R., Mallikarjuna Reddy, A. Eficient diabetic retinopathy detection using convolutional neural network and data augmentation, Soft Computing, 2023.
10. Chintha, V.V.R., Ayaluri, M.R. A Review Paper on IoT Solutions in Health Sector, Proceedings of International Conference on Applied Innovation in IT, 2023, 11(1), pp. 221–225
11. Manoranjan Dash et al.," Effective Automated Medical Image Segmentation Using Hybrid Computational Intelligence Technique", Blockchain and IoT Based Smart Healthcare Systems, Bentham Science Publishers, Pp. 174-182,2024
12. Umur Aybars Ciftci, ˙Ilke Demir, Lijun Yin "Detection of Synthetic Portrait Videos using Biological Signals" in arXiv:1901.02212v2
13. D. Güera and E. J. Delp, "Deepfake Video Detection Using Recurrent Neural Networks," 2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), Auckland, New Zealand, 2018, pp. 1-6.
14. Manoranjan Dash, "Modified VGG-16 model for COVID-19 chest X-ray images: optimal binary severity assessment," International Journal of Data Mining and Bioinformatics, vol. 1, no. 1, Jan. 2025, doi: 10.1504/ijdmb.2025.10065665.
15. Huy H. Nguyen , Junichi Yamagishi, and Isao Echizen "Using capsule networks to detect forged images and videos" in arXiv:1810.11215