



Deep Fake Detection Using Deep Learning : For Image And Video And Audio

¹Mr. Jalindar N. Ekatpure, ²Surayajit Sidheshwar Ingavale, ³Rohit Shankar Jagtap, ⁴Neelam Sachin Jachak, ⁵Rutuja Madhav Kandekar

^{1 2 3 4 5}S.B.Patil College of Engineering, Indapur

Email: j.ekatpure@gmail.com, suryajitingavale12@gmail.com, itsrohit173@gmail.com, nilamjachak734@gmail.com, rutujakandekar7@gmail.com.

Peer Review Information	Abstract
<p><i>Submission: 11 Sept 2025</i></p> <p><i>Revision: 10 Oct 2025</i></p> <p><i>Acceptance: 22 Oct 2025</i></p> <p>Keywords</p> <p><i>Deepfake Detection, Deep Learning, Convolutional Neural Network (CNN), Long Short-Term Memory (LSTM), Temporal Analysis, FaceForensics++, DFDC, Cybersecurity, Misinformation Prevention, Digital Media Integrity</i></p>	<p>This project proposes a comprehensive Deepfake Detection system using advanced Deep Learning methodologies to automatically analyze and identify manipulated images and videos. Deepfakes, generated using sophisticated generative techniques such as Generative Adversarial Networks (GANs), pose a serious threat to digital trust by enabling the spread of misinformation, identity theft, and reputational damage. To combat this, the system employs Convolutional Neural Networks (CNNs) for spatial image analysis and a hybrid CNN + LSTM architecture to capture subtle temporal inconsistencies across video frames. Unlike traditional handcrafted feature-based methods, this deep learning approach directly learns discriminative patterns from large benchmark datasets such as FaceForensics++ and DFDC. The platform not only preprocesses media inputs and trains high-performance models but also classifies content as Real or Fake with high accuracy, offering a robust defense mechanism against the rising tide of misinformation and digital manipulation.</p>

Introduction

The rapid growth of generative AI technologies, particularly GANs, has enabled the creation of highly realistic fake images and videos, commonly referred to as deepfakes. These synthetic media creations are increasingly being exploited to spread false information, harm personal reputations, manipulate public opinion, and even threaten democratic processes. The challenge of detecting deepfakes lies in their continuous evolution, as new generation methods are specifically designed to bypass existing detection techniques. Traditional approaches that depend on handcrafted features such as texture

inconsistencies, unnatural blinking, or lighting mismatches often fail when confronted with next-generation GANs or compressed video formats.

In response, the proposed system leverages Deep Learning techniques to extract both spatial and temporal features automatically. CNNs are employed to analyze image-level manipulations such as blending artifacts and texture anomalies, while LSTMs capture temporal inconsistencies like lip-sync mismatches and unnatural motion in videos. By training on large-scale datasets such as FaceForensics++ and DFDC, the system achieves strong generalization, making it capable of detecting deepfakes created by both

known and unseen manipulation techniques. This solution provides a practical platform—either web-based or desktop—that allows users to upload media and instantly determine its authenticity, thereby contributing to cybersecurity, digital forensics, and the fight against misinformation.

Research Gap

- Insufficient robustness against next-generation GAN models: Current detectors often fail against more advanced deepfake generation techniques, particularly those optimized to avoid detection. This highlights the urgent need for models that adapt to evolving manipulation strategies.
- Limited generalization beyond training datasets: Many detection systems perform well on benchmark datasets but struggle in real-world conditions where variations in lighting, compression, and noise reduce their accuracy. Cross-dataset generalization remains a major challenge.
- Lack of large-scale standardized benchmark datasets: Although datasets such as FaceForensics++ and DFDC exist, there is still a scarcity of diverse and continuously updated datasets to keep pace with emerging manipulation methods. This lack hinders the evaluation of true robustness.
- Poor detection of temporal inconsistencies in videos: While CNNs capture spatial cues effectively, they often fail to model long-term temporal inconsistencies. Many systems still miss artifacts like unnatural blinking or inconsistent lip-sync, which are crucial for robust detection.
- High computational cost and real-time deployment issues: Deep learning models, particularly large CNNs and hybrid architectures, are computationally expensive, limiting their use in mobile devices and real-time applications. Developing lightweight, efficient detectors remains an open research direction.

Problem Statement

Traditional deepfake detection techniques that rely heavily on handcrafted features are becoming increasingly ineffective in the face of evolving generative models. Handcrafted indicators such as blinking frequency, head pose, or texture anomalies are easily bypassed by modern GANs, leading to unreliable results.

Moreover, these approaches are often dataset-specific, achieving high accuracy in controlled experimental settings but failing to generalize to real-world, noisy, and compressed scenarios. The growing sophistication of deepfake generation techniques demands a more adaptable solution. Hence, there is a pressing need for an efficient, fully automated deep learning system capable of extracting both spatial and temporal features, generalizing across datasets, and delivering high accuracy with low latency in practical deployment environments.

Literature Survey

1. Recent Advances and Challenges of Deepfake Detection – Ran He et al. (2023)

This work provides a comprehensive review of current deepfake detection approaches, including CNNs, RNNs, frequency analysis, and biological signal-based methods. The authors highlight that while detection accuracy has improved, systems often fail when tested on unseen datasets or against adversarial attacks. Their study emphasizes the importance of building detectors with cross-dataset generalization and resilience to adversarial manipulation, setting a foundation for developing robust detection strategies in real-world scenarios.

2. FaceForensics++: Learning to Detect Manipulated Facial Images – Rössler et al. (2019)

The researchers introduced the FaceForensics++ dataset, which became a benchmark for evaluating deepfake detection methods. They tested several CNN-based models, including XceptionNet, and demonstrated high detection accuracy under controlled conditions. However, they observed that performance dropped when models were applied across different datasets. Their work stresses the need for datasets that support cross-domain evaluation and highlights that robust deepfake detectors must be trained on diverse manipulations.

3. MesoNet: A Compact Facial Video Forgery Detection Network – Afchar et al. (2018)

This study proposed MesoNet, a lightweight CNN architecture (Meso-4 and MesoInception-4) aimed at efficient real-time detection. Unlike deeper CNNs, which are computationally heavy, MesoNet achieves good detection accuracy while being deployable on low-resource systems like mobile devices. Although effective for certain manipulations, its limited depth reduces its ability to capture complex features. This research demonstrated the trade-off between

efficiency and accuracy in real-time detection.

4. Two-Stream Network for Tampered Face Video Detection – Li & Lyu (2019)

Li and Lyu presented a two-stream CNN framework that captures both spatial information and temporal motion inconsistencies in videos. The model integrates optical flow features with image-based analysis, thereby improving the detection of manipulations in long video sequences. Their work addresses the shortcomings of single-stream CNNs that fail to capture dynamic inconsistencies. However, it is computationally intensive, making large-scale or real-time deployment challenging.

5. LipForensics: Using Visual Speech for Deepfake Detection – Zhao et al. (2021)

This research focused on identifying lip-sync mismatches in manipulated videos. By applying temporal CNNs to the mouth region, the model detects inconsistencies between lip movements and spoken audio. The study showed promising results in identifying audio-visual mismatches but also acknowledged limitations when dealing with high-quality manipulations. Future work suggested that combining lip features with audio signals could further strengthen multimodal detection systems.

6. Face X-ray: A Simple, Generalizable Deepfake Detection Method – Wang et al. (2020)

Wang and colleagues proposed the Face X-ray method, which leverages boundary artifacts left during face-swapping operations. By training with alpha-matte and boundary supervision, their approach improved generalization to unseen deepfake generation techniques. The study demonstrated effectiveness across multiple datasets but noted that performance decreased under heavy compression. The simplicity and adaptability of this method make it an important contribution to generalizable detection.

7. Detection of GAN-Generated Faces Using Color Cues – McCloskey & Albright (2019)

This study identified abnormal color statistics as a signature of GAN-generated faces. The authors used CNNs combined with color distribution analysis to distinguish fake images from real ones. Their method was effective under clean conditions but struggled when images were compressed or distorted. While relatively simple, this approach provided an important insight into exploiting low-level color cues for deepfake detection.

8. Learning to Detect Fake Face Images in the Wild – Yang, Li & Lyu (2019)

The authors observed that most existing detectors failed in uncontrolled real-world

environments. They proposed a CNN-based model enhanced with domain adaptation techniques to handle diverse conditions, such as varying lighting and backgrounds. Their approach significantly improved cross-domain generalization. This work highlighted the need to move beyond lab-controlled datasets and address real-world variability in deepfake detection.

9. Vision Transformers for Deepfake Detection – Coccomini et al. (2021)

Coccomini and colleagues explored the use of Vision Transformers (ViTs) in deepfake detection, combining them with EfficientNet for improved efficiency. Transformers were found to model long-range dependencies in video frames more effectively than CNNs, making them suitable for sequential video data. While promising, the approach faced scalability issues due to the large computational requirements of transformers. The study opened a pathway for future transformer-based deepfake detectors.

10. A Survey on Deepfake Video Detection – Multiple Authors (2021)

This survey paper reviewed a wide range of deepfake detection techniques, including CNN-based, RNN-based, frequency analysis, and biological signal methods. The authors noted the absence of standardized benchmarks, making it difficult to compare systems fairly. They advocated for unified protocols and open-source datasets to advance the field. Their work emphasized the need for consistent evaluation standards to enable the development of truly reliable deepfake detection models.

Conclusion

The proposed deepfake detection system introduces a novel AI-based platform that integrates CNN and LSTM architectures to extract both spatial and temporal features from images and videos. By training on benchmark datasets such as FaceForensics++ and DFDC, the system achieves enhanced generalization, enabling it to detect manipulations produced by both existing and newly emerging techniques. The inclusion of a user-friendly interface ensures accessibility for journalists, law enforcement agencies, social media platforms, and general users, thereby increasing its societal impact. This system not only contributes to digital forensics and cybersecurity but also plays a crucial role in mitigating the spread of misinformation. While challenges such as adversarial attacks, real-time efficiency, and dataset limitations remain, the proposed approach lays a strong foundation for future work aimed at building more robust, lightweight, and multimodal deepfake detection

systems.

References

Rössler, A., Cozzolino, D., et al., "FaceForensics++: Learning to Detect Manipulated Facial Images," IEEE, 2019.

Afchar, D., Nozick, V., et al., "MesoNet: Compact Facial Video Forgery Detection Network," IEEE, 2018.

Li, Y., Lyu, S., "Two-Stream Network for Tampered Face Video Detection," IEEE, 2019.

Zhao, Y., Shen, J., et al., "LipForensics: Using Visual Speech for Deepfake Detection," IEEE, 2021.

Wang, S.-Y., Wang, O., et al., "Face X-ray: Generalizable Deepfake Detection Method," IEEE, 2020.

McCloskey, S., Albright, M., "Detection of GAN-Generated Faces Using Color Cues," IEEE, 2019.

Yang, X., Li, Y., Lyu, S., "Learning to Detect Fake Face Images in the Wild," IEEE, 2019.

Coccomini, D., Messina, N., et al., "Vision Transformers for Deepfake Detection," arXiv, 2021.

"A Survey on Deepfake Video Detection," IEEE Access, 2021.

Zhang, Z., et al., "Deepfake Detection: A Survey," IEEE Access, 2022.

Wang, X., et al., "Deepfake Detection via Audio-Visual Fusion," IEEE Transactions on Multimedia, 2021.

Nguyen, H., et al., "Deepfake Detection with Convolutional Neural Networks," IEEE Transactions on Information Forensics and Security, 2020.

Kim, J., et al., "Deepfake Detection Using Recurrent Neural Networks," IEEE Transactions on Circuits and Systems for Video Technology, 2021.

Lee, S., et al., "Deepfake Detection: A Comparative Study," IEEE Transactions on Cybernetics, 2022.

Ekatpure, J. N., Tavate, C. S., Malshikare, S. S., Khomane, A. B., & Tamboli, M. J. L. (2025). Artificial intelligence based virtual keyboard and mouse for computer. International Journal on Advanced Computer Theory and Engineering, 14(1), 449-456.

Ekatpure, J. N., Mohite, S. D., Shinde, A. A., Shirkande, N. B., & Upase, V. V. (2025). Campus recruitment system using machine learning. International Journal on Advanced Computer Theory and Engineering, 14(1), 427-432

Aware, D. B., Sayyad, S. R., Shaikh, A. H., Thombare, S. B., & Ekatpure, J. N. (2025). Translation Assistant for Converting Sign Language to Text and Audio. International Journal on Advanced Computer Engineering and Communication Technology, 14(1), 445-449.

Ekatpure, J. N., Aware, D. B., Shaikh, A. H., Sayyad, S. R., & Thombare, S. B. (2024). A comprehensive survey on sign language translation systems: Bridging gestures, text, and audio for enhanced communication. International Journal of Recent Advances in Engineering and Technology, 13(2), 15-21.

Ekatpure, J. N., Tavate, C., Malshikare, S., Khomane, A., & Tamboli, M. J. (2024). Advancements in AI-powered virtual keyboards and mice: A survey of cutting-edge technologies for modern computing. International Journal on Advanced Computer Theory and Engineering, 13(2), 52-57.