



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 14    **Issue:** IV    **Month of publication:** April 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.81476>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Deep Fake Detection Using CNN

Mayank Tamta, Lakshay Bhardwaj, Khushi Yadav, Janhavi Nehra, Srishthi Vashisth

Department of Computer Science and Engineering, MIET, Meerut, Uttar Pradesh, India

**Abstract:** Recent progress in convolutional neural networks has enabled automated systems to examine manipulated visual content with improved reliability. Deepfake media, generated using advanced learning techniques, poses a serious challenge to digital authenticity and public trust. This survey reviews recent deepfake generation approaches and CNN-based detection methods, focusing on commonly used datasets, model architectures, and evaluation metrics. The study highlights challenges related to generalization computational efficiency, and real-world deployment, providing insights that can support future research and practical applications in deepfake detection.

**Keywords:** DeepLearning, FakeDetection, NeuralNetworks, Social Networks.

## I. INTRODUCTION

Recent advances in deep learning have shown promising results in detecting such manipulations. Convolutional neural networks, in particular, are capable of identifying subtle visual inconsistencies introduced during the deepfake generation process. This paper surveys CNN-based deepfake detection techniques, emphasizing widely used datasets, model architectures such as Xception and MobileNet, and evaluation strategies that enhance detection reliability.

The rapid growth of digital media sharing platforms has increased the circulation of manipulated images and videos across the internet. Deepfake technology enables the creation of highly realistic synthetic media that can mislead viewers and compromise information credibility. As these manipulations become more sophisticated, ensuring the authenticity of digital content has become a critical challenge in multimedia security and digital forensics and leveraging advanced algorithms such as convolutional neural network [7]

## II. LITERATURE REVIEW

S.No	Year	Reference	CNN Architecture	Key Findings	How it helps your project
1	2018	<b>MesoNet— Afchar et al.</b> <i>MesoNet: a Compact Facial Video Forgery Detection Network</i>	Compact CNN (“Meso-4”, “Meso-Inception-4”) focusing on mesoscopic facial artifacts	Shallow CNNs detect blurriness and texture anomalies effectively with low computation.	Provides a lightweight baseline model for real-time or small-scale deepfake detection.
2	2019	<b>FaceForensics++ —Rössler et al.</b> <i>FaceForensics++: Learning to Detect Manipulated Facial Images</i>	Evaluated multiple CNNs: Xception-based models achieved best performance	Demonstrated Xception’s strong generalization on facial manipulation tasks.	Use FaceForensics++ for standardized training/testing and benchmarking CNN performance.
3	2019	<b>XceptionNet for Deepfake Detection</b> (multiple studies)	Xception (depthwise separable convolutions) fine-tuned for binary classification	Experimental evaluations across multiple studies indicate that Xception-based architectures	Employ Xception as a high-performing backbone for your CNN model.
4	2020	<b>Celeb-DF — Li et al.</b> <i>Celeb-DF: A Large-scale Challenging Dataset for DeepFake Forensics</i>	CNN-based detectors (Xception, ResNet) tested for cross-dataset performance	Highlighted the challenge of realistic deepfakes and cross-dataset generalization.	Test your trained CNN on Celeb-DF to evaluate real-world robustness.

5	2020	<b>DFDC Dataset— Dolhansky et al. DeepFake Detection Challenge Dataset</b>	Ensemble CNNs (Xception, EfficientNet, ResNet)	Large-scale benchmark enabling high generalization for deepfake detection	Use DFDC for fine- tuning or large-scale testing for improved accuracy.
6	2021	Generalization & Bias Study — Mehra et al.	Generalization & Bias Study — Mehra et al.	Found CNN models struggle with generalization and demographic bias	Helps in evaluating fairness and cross- domain accuracy of your CNN model
7	2024	<b>Recent Surveys/ Evaluations</b>	Reviews of CNN + multimodal approaches	CNNs remain core detectors; hybrid temporal/audio models improve performance.	Guides dataset choice, architecture selection, and performance metrics for your project.

### III. METHODOLOGY

#### A. Data Collection and Preprocessing

The deepfake detection process follows a structured pipeline designed to analyze both authentic and manipulated video content. A diverse set of real and forged video samples is selected to improve the model's ability to handle variations in facial appearance and recording conditions. [3]

This dataset should include both authentic and manipulated videos, incorporating a variety of conditions, such as indoor, outdoor environments, varying lighting conditions, diverse facial expressions, and different camera qualities. Ensuring diversity in the dataset helps models generalize better, enabling them to effectively identify deepfakes across varied real-world situations and mitigate potential biases that could compromise the reliability of the detection system. Labels must accurately differentiate between real and deepfake videos.

The main aim of the work was to determine if videos were genuine or if they were created through [2] deepfake technologies. Therefore, it is apparent that for this system to work, the input should solely consist of a video clip. However, based on the fact that input for deep learning models is images, the conversion from video input to model input should be conducted. This is achieved through a preprocessing module. The videos may have more elements than just a face. Each frame of the video is not restricted to the individual's face only, since most of the video frame consists of the person's body parts as well as the background area of the picture. Besides, these uncorrelated features may hinder the model's training using face wrapping artifacts. [8] The face area is the focal point, while the image gets input as facial data from the image by the pre-processing module. This pre-processing module further consists of three discrete stages: grabbing frames from a video feed, detecting faces in such frames one at a time, and storing these areas where there are faces in the form of pictures. Each of these procedures is expounded on here below.

We first started capturing the video input into frames. Each video is converted into a sequence of image frames using OpenCV to enable image-based processing suitable for CNN models. Since deepfake manipulation mainly affects facial regions, face detection is applied to isolate facial areas from each frame, reducing the influence of background information on model training.

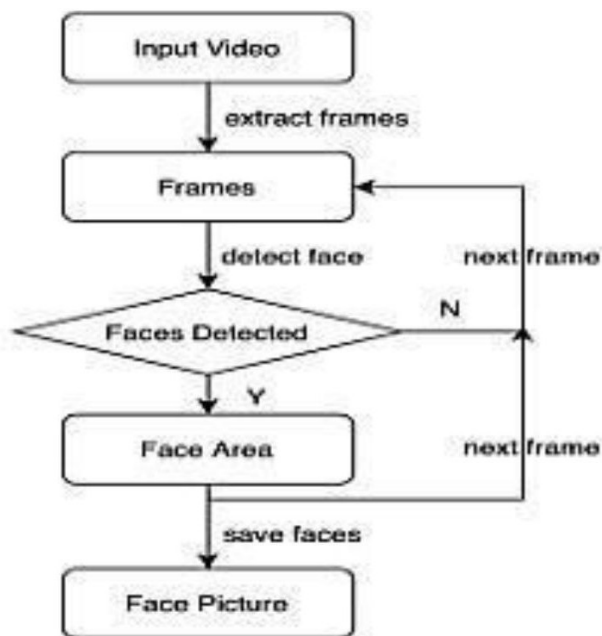


Fig3.1.FlowChartof DetectingDeep Fake

In addition, considering both are similar, putting all of them in a training set will make training ineffective due to their high similarity, and this may also result in other problems, including overtraining. The videos used here had 30 frames per second for each clip selected in this study. Upon assessment, it became apparent that one option has to choose an image interval of four.

The second step was to detect faces that are available in an image and automatically tag them. Facial regions are detected using a pre-trained Haar-based classifier, which enables efficient localization of faces within each frame. To avoid incorrect detections, the largest detected facial region is selected for further processing, ensuring accurate facial representation. Classifier was then selected out of a number of classifiers that were tested due to its precise area where a face can be found, which means it creates the best area concerning face localization. However, there were a few non-face selections made. For instance, after conducting some experiments.

The third step was to store the detected face area as a new image. Prior to storage, all facial photos had to be uniformly resized. While using the Xception model, the size of the picture was supposed to be 299\*299, whereas with MobileNet it is 224\*224.

### B. Model Selection

Currently, there exist various deep learning models and frameworks. In this paper, Xception and MobileNet were selected as the models due to the reasons below. To begin with, Xception has high performance given its performance benchmarking at the FaceForensics testing environment.

Researchers for example can use FaceForensics as a benchmarking platform to test their models. As it comes with a detailed manual, Xception is better than others on 4 different datasets on different researches [6]

Xception has a similar architecture to MobileNet, which makes it the reason why it was selected. The two models are built on CNN with the use of both depthwise and pointwise convolutional layers. Unlike MobileNets, which has fewer features in order to make the model more efficient.

### C. Xception for Deepfake Detection

Xception is an architecture for a convolutional neural network that is famous for its handling of complicated models without sacrificing resources. It acts like a useful tool that can notice tiny things, so it becomes easy to figure out whether anything has been tampered with in a video or not. Using transfer training techniques makes it possible to make Xception better using customized aspects of the deepfake detection data, making it perfect at distinguishing between genuine pictures and fake ones.

#### D. MobileNet for Deepfake Detection

MobileNet is optimized for computationsavinguses,henceitslightness. MobileNet can still detect deep fakes despite being small. One way this can be done is to use MobileNet on deepfake detection datasets for later tuning in real- time or for edgebased systems due to their limitations in computations

#### E. Feature Extraction

Particularly in videos, feature extraction is key because of the significant role that temporal information plays. It is very important to extract features that capture motion patterns,consistency overtime, and spatialrelationshipsbetweenframesfor effective deepfake detection. After extracting the features, the model is trained over the preprocessed data set.

#### F. Model Training

Model training includes dividing the preprocessed dataset into training, validation, and test subsets. This process often employs transfer learning, where the model is initialized using weights from a pre-trained network and subsequently fine- tuned for the specific task of deepfake detection.

#### G. ValidationandEvaluation

While still in the training process, it's important to see how the model performs using the validation set in order to avoid overfitting. Tune the hyperparameters (learning rate, batch\_size and optimizer) such that the model performs best. Check how successful your model is under differentconditionswhentestedusingtester sets. These can be accuracy (measure of correct predictions), precision (how much usefulinformationyougetfromwhatyou're looking at), recall (how many people who shouldhavehadsomethingdoneknewabout it),F1-score(balancebetweenprecisionand recall), and the ROC curve. The ROC and IEEE does not plot the accuracy or error of a binary classifier but displays the separation between classes.[14]

It illustrates the connection between sensitivity and specificity. Evaluate the trainedmodelintermsoftesting;evaluateit according to its classification of class instancesintopositiveandnegativeclasses. Anomaly refers to the feature space of instances that are unusual. The separated classes are linearly separable.

Productiondescribestasksthatarebroken down into smaller components and then arranged in a sequence. Machine learning departsfrom traditionalstatisticalmethods.

### IV. RESEARCH GAPS

- 1) Scope and Purpose: The literaturereviewsystematicallyexamines112studies on deepfake detection from 2018 to 2020. The aim is to classify and evaluate various deepfake detection methods, tools, and datasets.
- 2) Deepfake Techniques: Research is categorized into four main detection techniques
- 3) Deep Learning: The most widely used and effective, accounting for 77% of studies.
- 4) Machine Learning: Focuses on traditional techniques like decision trees and support vector machines.- Statistical Methods: Utilize statistical patterns for deepfake detection.
- 5) Blockchain: Provides decentralized verification of authenticity, though less common.
- 6) Datasets: Popular datasets include FaceForensics++, Celeb-DF, and DFDC, widely used for training and testing.
- 7) Evaluation: The studies predominantlyuse accuracy and AUC (area under the curve) for performance measurement, with deep learning methods consistently outperforming others in accuracy and robustness.
- 8) Detection: The paper highlights for standardized evaluation frameworks to ensurecomparabilityandfairnessamong variousdetectionmodels.
- 9) GeneralizationAcrossDatasets: Many deep fake detection models perform wellon specific datasets but struggle with unseen data or deep fakes generated by different methods. Research is needed to develop models that generalize effectively across diverse deep fake types and sources.
- 10) Real-Time Detection: Current detection methods often lack the efficiency needed for real-time analysis, which is crucial for applications like social media monitoring. Exploring lightweight and fast deeplearningmodelscouldbridgethisgap.
- 11) Robustness Against Advanced **Techniques**: As deep fake generation methods evolve, detection models need to adapt to new forms of manipulation. Researchinto adversarialtrainingandmore adaptable architectures can make models more resilient against evolving deep fake techniques.
- 12) Computational Efficiency: Deeplearning-based detection methods typically require high computational resources, limiting their

practical deployment, especially on edge devices. There's a need for models that balance performance with computational efficiency, potentially through model compression or optimization techniques.

- 13) **Explainability:** Deep learning models, especially CNNs, making it difficult the decision-making processes. Future work could focus on making these models more interpretable, aiding trust and transparency in automated detection.
- 14) **Privacy Ethical Considerations:** Implementing deep fake detection involves analyzing potentially sensitive user data. Research is needed to ensure detection systems respect privacy and ethical considerations, such as minimizing data collection or ensuring user consent.
- 15) **Attack Resistance:** Deep learning models may be susceptible to adversarial attack, where minor alterations in input data can fool the model. Exploring defense mechanisms to harden models against such attacks is essential for reliable detection.
- 16) **Forward Progression:** To from a given dataset, while the calculate the output values of a neural network, the input data must pass through the network sequentially, progressing from one layer to the next, starting from the input layer and ending at the output layer. This step-by-step process of transmitting inputs through the network to generate outputs is referred to as forward propagation.
- 17) **Backpropagation:** Each hidden layer in a neural network receives inputs from the previous layer, processes them through an activation function, and then generates the predicted output values. Backpropagation is a reverse process that focuses on adjusting the weights to optimize the network's ability to accurately predict outputs. This process involves calculating the gradient of the error and using stochastic gradient descent (SGD) to minimize the cost function, thereby reducing prediction errors and improving the model's performance.
- 18) **Face Manipulation:** particularly in face manipulation, pose a significant threat by distorting the original facts in digital images. As technology continues to evolve, it has become increasingly important to employ it to verify the originality of videos and identify manipulated information. These detection methods are crucial in ensuring the credibility and integrity of visual media in today's digital landscape. Deep learning techniques, such as GAN, are based on the principles of autoencoders and decoders to aid in identifying fake images or videos. These networks leverage the generator and discriminator framework to identify and differentiate between real and manipulated content.
- 19) **GAN:** Generative Adversarial Networks are made up of two neural networks: a generator and a discriminator. The generator produces synthetic images discriminator analyzes these images to determine their authenticity differentiating between real and fake videos
- 20) **Media Filter :** A media filter is then used to remove unwanted noise from the video. To improve image quality, bicubic interpolation is applied to increase pixel density, followed by bicubic transformation to enhance overall clarity. Yadav and Salmani (2019) discussed the principles behind deepfake techniques, including face image swapping, with an emphasis on achieving high precision (Maheswaran et al., 2018).
- 21) **DBN model :** The existing approaches face several challenges, including inefficiencies in detecting deepfake images, high error rates, long processing times, and inaccuracies in data assessment. The FF-LBPH-DBN model specifically addresses these issues by focusing on minimizing computational costs and applying various metrological parameters in a more efficient manner.

## V. ENHANCEMENT

As deepfake technology continues to advance, it increasingly threatens the authenticity and reliability of digital media. This growing concern presents numerous challenges in the deepfake detection, particularly in the application of deep learning, where ongoing development and advancements are necessary to keep pace with evolving manipulation techniques. One area of focus should be on improving the interpretation and interpretation of deep learning models for research. Building descriptive artificial intelligence (XAI) models that provide information about the features and patterns used for detection can increase transparency and build trust in these systems. [4] IEEE (Institute of Electrical and Electronics engineers) makes techniques like visualization, feature mapping, and model interpretation can provide insight into the inner workings of deep learning models, enhancing the understanding and application of these processes. Another important area for improvement is the quality of the model and the availability of training materials. The effectiveness of deepfake detection models depends significantly on the quality and accessibility of training data. Researchers should focus on creating and expanding a good data repository with a variety of in-depth measurements, including the most accurate control methods. Collaborating with media organizations, government agencies, and other stakeholders can help collect and manage this data, ensuring detection models are trained on the most important and cutting-edge data. For more complex business processes, the detection model needs to be updated and modified accordingly. Researchers should focus on developing flexible and adaptive deep learning models capable of quickly detecting and addressing emerging deepfake threats.

This will include techniques such as resilience training, adaptive learning, and continuous learning that can help find models that are ahead of the norm and control their impact on deep tech exchanges. Various formats, such as image, audio, and text files, Researchers should focus on developing more efficient deep learning architectures and using data from different sources to increase discovery accuracy. Additionally, investigating common communication techniques can reveal effective patterns across different media and types of leaders, thereby improving their effectiveness and efficiency in the actual world. More importantly, it is very much important to use these systems at scale and efficiently. Researchers should investigate strategies to optimize deep learning models, such as model compression, quantization, and the hardware acceleration, to provide insight into deep learning across resource constraints such as edge devices and exposing deepfake videos through face warping artifact detection [15] To ensure deep solutions meet real-world needs and requirements, researchers need to develop partnerships with experts and industry stakeholders. This collaboration can create a framework and guidance for the responsible use and application of deep research tools, ultimately increasing their impact and social benefits. We can continue to develop and deliver better, more transparent, and more flexible solutions to the growing threat of media and information misuse in this new digital age.

## VI. CONCLUSION

In conclusion, the proliferation of deepfakes presents quite a significant challenge, exacerbated by the accessibility of tools for creating and distributing fake images and videos on social media platforms. Deep learning methods have emerged as a promising solution for detecting deepfakes, with various techniques developed for image and video detection. This paper provides a survey of current applications and tools for making deepfakes and detailed scrutiny for deepfake detection methods based on images and videos. We discussed their architectures, tools, and performance and highlighted publicly accessible datasets used for training and evaluation. Although great strides have been recorded in the detection of deepfakes through deep learning, much remains to be done as far as improving the quality of these gifts is concerned. This makes it hard for the current techniques because they can never be without challenge due to the progressive nature of their construction. This is why more research needs to focus on enhancing deep fake detection algorithm performance with respect to identifying model architectures that work best or alternatively deploying such architectures within social networking sites so that they can help dampen the effects emanating from these deepfakes. Moving forward, addressing these challenges and advancing research in these deepfake detection will be crucial to safeguarding against the harmful effects of altered media and preserving trust and integrity in digital content.

## REFERENCES

- [1] S. Senhadji and R. A. San Ahmed presented a method for detecting fake news by utilizing Naïve Bayes and Long Short-Term Memory (LSTM) algorithms in the IAES International Journal of Artificial Intelligence, Volume 11, Issue 2, 2022, pages 748-754.
- [2] K. N. Ramadhani and R. Munir conducted a comparative study on various methods for deepfake video detection, International Conference on Information and Communications Technology, November 2020, pages 394-399.
- [3] D. Pan, L. Sun, R. Wang, X. Zhang, and R. O. Sinnott explored deepfake detection through deep learning techniques at the IEEE/ACM International Conference on Big Data Computing, Applications, and Technologies, December 2020, pages 134-143.
- [4] A. A. Maksutov, V. O. Morozov, A. A. Lavrenov, and A. S. Smirnov discussed machine learning methods for detecting deepfakes at the IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering, January 2020, pages 408-411.
- [5] T. T. Nguyen, Q. V. Nguyen, and D. T. Nguyen provided a comprehensive survey on deep learning methods for the creation and detection of deepfakes, published in Computer Vision and Image Understanding, Volume 223, 2022, pages 1-19.
- [6] A. O. Kwok and S. G. Koh examined the social construction of technology perspective in deepfake research, published in Current Issues in Tourism, Volume 24, Issue 13, 2020, pages 1798-1802.
- [7] M. Westerlund reviewed the emergence of deepfake technology in the Technology Innovation Management Review, Volume 9, Issue 11, 2019, pages 40-53.
- [8] Y. Li and S. Lyu presented a method for exposing deepfake videos by detecting face warping artifacts, arXiv:1811.00656, 2018.
- [9] A. M. Almars reviewed deepfake detection techniques using deep learning in the Journal of Computer and Communications, Volume 9, Issue 5, 2021, pages 20-35.
- [10] N. S. Ivanov, Arzhakova, and V. G. Ivanenko explored the combination of deep learning and super-resolution algorithms for deepfake detection at the IEEE Conference of Russian Young Researchers in Electrical and Electronic Engineering, January 2020, published in the proceedings of the 3rd pages 326-328.
- [11] K. Zhu, B. Wu, and B. Wang proposed a deepfake detection method using clustering-based embedding regularization at the IEEE 5th International Conference on Data Science in Cyberspace, July 2020, pages 257-264.
- [12] F. Matern, C. Riess, and M. Stamminger exploited visual artifacts to expose deepfakes and face manipulations, presented at the IEEE Winter Applications of Computer Vision Workshops, January 2019, pages 83-92.
- [13] E. Sabir, J. Cheng, A. Jaiswal, W. AbdElmageed, I. Masi, and P. Natarajan introduced recurrent convolutional strategies for detecting face manipulation in videos in Interfaces (GUI), Volume 3, Issue 1, 2019, pages 80-87.
- [14] D. Güera and E. J. Delp employed recurrent neural networks for deepfake video detection at the 15th IEEE International Conference on Advanced Video and Signal-Based Surveillance, November 2018, pages 1-6.
- [15] Y. Li and S. Lyu discussed exposing deepfake videos through face warping artifact detection, arXiv:1811.00656, 2018.
- [16]



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)