



**IJIRCCCE**

e-ISSN: 2320-9801 | p-ISSN: 2320-9798



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

Volume 12, Issue 10, October 2024

**ISSN** INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA

**Impact Factor: 8.625**



9940 572 462



6381 907 438



ijircce@gmail.com



www.ijircce.com



# Deepfake Classification for Human Face Using Deep Learning

Prof.A.V.Bendale<sup>1</sup>, Aishwarya Bhosale<sup>2</sup>, Gargi Jadhav<sup>3</sup>, Gaytree Kadam<sup>4</sup>, Soham Pandhare<sup>5</sup>

Assistant Professor, Department of Information Technology, Sandip Institute of Technology and Research Centre,  
Nashik, India <sup>1</sup>

Student, Department of Information Technology, Sandip Institute of Technology and Research Centre,  
Nashik, India<sup>2,3,4,5</sup>

**ABSTRACT:** This project aims to develop a versatile web application for audio and image processing, encompassing three key objectives. Firstly, the application will employ advanced algorithms to reduce external noise from audio input, enhancing the overall audio quality. Users can upload audio files and easily download the processed, clean audio, making it suitable for scenarios like conference calls and interviews. Secondly, the application will incorporate state-of-the-art image processing techniques for face detection in uploaded images, including the ability to highlight or mark detected faces and provide annotated images for download. This feature is valuable for applications like automatic tagging in photo albums and security surveillance. Lastly, the web application will enable users to upload video files and extract multiple frames, allowing them to specify frame counts or time intervals for extraction. Users can download the processed frames, which simplifies tasks like creating video previews and generating thumbnails. Overall, this comprehensive application will cater to users looking for powerful audio and image processing functionalities in a user-friendly online environment, with applications ranging from media production to content creation and security.

**KEYWORDS:** Audio Processing, Noise Reduction, Audio Quality Enhancement, Face Detection, Image Processing, Annotated Images, Automatic Tagging, Security Surveillance, Video Processing, Frame Extraction

## I. INTRODUCTION

Deepfakes are artificial intelligence-modified media that change or synthesize pictures, audio, or video. Discuss the increasing ubiquity of deepfakes across several sectors, including entertainment, politics, and social media, along with the related hazards, such as disinformation, privacy infringements, and identity theft. In recent years, the rapid advancement of artificial intelligence and machine learning has led to the development of deepfake technology, enabling the creation of very realistic synthetic media. Deepfakes use sophisticated algorithms to modify or generate visual and audio content that may be indistinguishable from genuine media, raising significant concerns around misinformation, privacy, and security. Deepfake human faces pose a significant challenge due to their ability to deceive audiences in various contexts, such as social media and crucial legal and financial situations. Traditional methods for detecting deep fakes often rely on heuristic algorithms or human assessment, which become increasingly inadequate as technology advances. Deep learning, particularly convolutional neural networks (CNNs), offers a feasible alternative by automating detection and improving accuracy. Convolutional neural networks (CNNs) excel in image analysis tasks due to their ability to discern complex patterns and features at several levels of abstraction. Training these algorithms on comprehensive and diverse datasets of both genuine and synthetic faces facilitates the identification of subtle discrepancies and artifacts that may elude human observation.

## II. LITERATURE SURVEY

As stated in [1], the identification of anomalous behaviors in IoT networks, which are among the most significant and extensively used networks at now, We implemented a modified LSTM RNN using Python to create our model. The effectiveness and efficiency of our suggested model were evaluated using the second-generation Google ML framework, TensorFlow, in conjunction with the LSTM's inherent characteristics and the following mathematical calculations. In order to assess the performance of our model, we conducted four tests using seven mediators and two



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

criteria. In each test, we varied the mediators of data quantity and duration to evaluate the model's efficacy. The testing included assigning multiple values of volume and iterative periods. Randomly-generated values were utilized for weights, biases, various layers, nodes, and learning rate until the model achieved the highest accuracy detection rate.

As stated in reference [2], this paper addresses the issue of detecting anomalous behavior during examinations by using the enhanced YOLOv3 algorithm. We began by generating anomalous behavior data sets for evaluation. Then, we improved the system's loss function to optimize its performance. We used the K-Means technique to choose the most suitable anchor boxes. Additionally, we developed a new backbone called Darknet32. Finally, we employed the frame-alternate dual thread approach to identify the video. By evaluating the algorithm's detection accuracy (measured by AP and mAP), detection speed (measured by FPS), and memory usage, the following conclusions may be drawn. Utilizing GIoU loss and focused loss to optimize the loss function of the YOLOv3 algorithm, together with employing the K-Means method to cluster the bounding boxes in the dataset for obtaining optimal anchor boxes, will enhance the program's ability to accurately identify anomalous behavior during inspection. The implementation of the Darknet32 backbone, as suggested in this work, may enhance the speed of anomalous behavior identification during examination. Additionally, it can minimize computer memory use while maintaining a high level of detection accuracy. The use of the frame-alternate dual thread detection approach may significantly enhance the velocity of identifying aberrant behavior throughout the examination process while minimizing the memory usage. Moreover, this method effectively fulfills the need for real-time detection.

Machine learning and deep learning have enabled the development of crucial applications, such as anomaly detection, as reported by [3] today. Transfer learning, despite its modernity, has emerged as a significant innovation in the realm of deep learning due to its very promising outcomes. In this work, transfer learning is used to extract human motion characteristics from RGB video frames in order to enhance detection accuracy. A pre-trained model of the Visual Geometry Group network 19 (VGGNet-19) is used to extract descriptive features using a convolutional neural network (CNN). Subsequently, the feature vector is inputted into the Binary Support Vector Machine classifier (BSVM) in order to create a binary-SVM model. The effectiveness of the suggested framework is assessed using three metrics: accuracy, area under the curve, and equal error rate.

As reported by [4], a machine learning system was created to identify unusual behaviors in individuals with dementia. This was achieved by using inexpensive sensors to gather lifelog data. The system underwent testing at a nursing home located in South Korea, where real data from patients suffering from dementia was gathered and used for the purpose of training machine learning models. The suggested system has the capacity to mitigate the burden on caregivers by rapidly identifying anomalous behaviors in nursing homes and notifying caregivers, thereby minimizing accidents and injuries among dementia patients. This system, constructed utilizing economical components, particularly inexpensive gyro sensors, provides a viable alternative for widespread adoption in several nursing homes. Our technology offers a cost-effective alternative to conventional surveillance systems, such as CCTV-based systems, without compromising on efficacy. One of the main benefits of our system is its cost-effectiveness, which has the potential to improve the quality of life for dementia patients and decrease societal expenditures.

As stated in reference [5] A significant challenge in UEBA/DSS intelligent systems is the extraction of valuable insights from a vast volume of unorganized and inconsistent data. Management decision-making should rely on empirical evidence obtained from the analyzed characteristic. However, given the information acquired, it is challenging to make any management judgment due to the diverse nature of the data and their substantial amounts. A proposal is made to use machine learning approaches in the construction of a mobile UEBA/DSS system. This will enable the attainment of a high level of data analysis quality and the identification of intricate interrelationships within the data. A comprehensive list of the most influential elements was compiled based on the input provided to the analysis techniques throughout the investigation.

As stated in reference [6], there exists a very effective algorithm capable of identifying and pinpointing abnormalities inside films. To address the issue of insufficient negative examples, the technique employs a spatiotemporal autoencoder to detect and extract the anomalous behaviors within the dataset. This is achieved using an unsupervised learning approach. A spatiotemporal convolutional neural network (CNN) is built with a straightforward design and little computational complexity. The supervised training approach is used to train the spatiotemporal CNN using both





## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

positive and negative data in order to create the detection model. Experiments are performed on the datasets from UCSD and UMN.

According to reference [7], there is a precise and efficient approach available for identifying deviant conduct. Our approach views the video as a sequence of frames. During the training phase, our deep learning framework extracts appearance characteristics and learns the connection between past and present features in the regular video. During the testing phase, any anticipated traits that deviate from the actual features are classified as aberrant. Our model is constructed as a feature prediction framework with a novel temporal attention mechanism. During the feature extraction phase, we convert a pre-trained Vgg16 network into a fully convolutional neural network. We use the output of the third pooling layer as the appearance feature extraction to accurately capture static appearance data. Subsequently, a novel temporal attention method is proposed to ascertain the impact of various historical appearance aspects on the current features, effectively addressing the challenge of portraying dynamic motion features. Ultimately, the LSTM network is used to decipher the past sequences of data by including temporal attention, enabling the prediction of the features at the present instant.

As per the intrusion detection system (MIDS) mentioned in reference [8], it has the capability to identify any unusual behavior inside the system, even if the attacker attempts to hide it within the system's control layer. A supervised machine learning model is created to categorize regular and irregular actions in an Industrial Control System (ICS) in order to assess the effectiveness of the MIDS. A hardware-in-the-loop (HIL) testbed is created to replicate the behavior of power-producing devices and use the attack dataset. In our suggested methodology, we used several machine learning models to analyze the dataset. These models have shown exceptional performance in identifying abnormalities within the dataset, particularly in detecting covert assaults.

The current study, as described in [9], elucidates key stages of NADSSs, including pre-processing, feature extraction, and the identification and detection of harmful activity. Furthermore, the researchers have extensively examined recent machine learning methods, such as supervised, unsupervised, new deep, and ensemble learning techniques, in the context of the detection and recognition phase. Additionally, they have provided information about benchmark datasets that are currently accessible for training and evaluating machine learning techniques.

According to [10], efforts have been made to address this deficiency by presenting a comprehensive overview of the latest strategies used to tackle the main issues and fundamental hurdles in handling IoT data. The text also presents information on the characteristics of data, kinds of anomalies, learning mode, window model, datasets, and assessment criteria. Investigations are conducted on research difficulties pertaining to the evolution of data, evolution of features, windowing, ensemble techniques, the nature of input data, data complexity and noise, parameter selection, data visualizations, data heterogeneity, correctness, and large-scale and high-dimensional data. In conclusion, this paper provides a summary of the difficulties that need extensive research efforts and outlines potential future approaches.

**Table I Summary Table**

Sr.No	Title	Methodology	Algorithms	Limitations
[11]	Detection of Deep Fake in Face Images Using Deep Learning	Utilized a dataset of real and deepfake images; implemented a CNN for feature extraction and classification.	CNN, ResNet	High computational cost; may struggle with high-resolution images.
[12]	DeepFakes: Detecting Forged and Synthetic Media Content Using Machine Learning	Analyzed various machine learning techniques to identify deepfakes in images and videos; utilized both pixel and feature-level analysis.	CNN, SVM, GAN-based detection	High computational cost and time-consuming training process; performance may vary across different deepfake types.
[13]	Deepfakes	Employed a CNN architecture	CNN, Transfer	High sensitivity to variations in



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

	Classification of Faces Using Convolutional Neural Networks	to classify face images as real or deepfake; utilized a dataset of manipulated face images.	Learning (e.g., VGGNet)	lighting and image quality; may struggle with unseen deepfake techniques.
[14]	Classification of Deepfake Videos Using Pre-trained Convolutional Neural Networks	Utilized pre-trained CNN models (e.g., VGG16, ResNet) to extract features from video frames; fine-tuned models on a labeled dataset of real and deepfake videos.	VGG16, ResNet, Inception	High computational demands; performance may vary based on the quality of the pre-trained model and dataset.
[15]	Deepfake detection: humans vs. machines	Conducted experiments comparing human detection capabilities with those of machine learning algorithms on various deepfake datasets.	CNN, Human Perception Studies	Variability in human judgment; machine models may not generalize well across different types of deepfakes.

### III. CHALLENGES FOR DEEPFAKE CREATION AND DETECTION

In recent years, many DeepFake tools exhibiting remarkably realistic performance have emerged, with many other ones now under development. The advancement of the DeepFake generation model poses significant hurdles for forensic professionals in their efforts to counteract it. DeepFakes are AI-generated hyperrealistic photos or films that have been digitally manipulated using methods such as face swapping and altering traits, depicting persons engaging in actions and dialogues that never occurred.

Generative Adversarial Networks (GANs), prominent artificial intelligence systems, include two competing models—discriminative and generative—that enhance their efficacy in producing convincing forgeries. Impersonations of actual individuals can achieve high virality and disseminate rapidly throughout social media platforms, making them a useful instrument for propaganda. In digital forensics, similar to other security-related fields, it is essential to consider the existence of an adversary who is actively trying to deceive investigators. A proficient attacker familiar with the principles behind forensic technologies may use various counterforensic measures to evade detection. Forensic technologies must be capable of identifying situational dangers and any real-world circumstances that may compromise test accuracy. Consequently, the many counter-forensics strategies designed to obfuscate existing detectors are crucial in advancing multimedia forensics since they reveal the deficiencies in present methodologies and stimulate research towards more resilient solutions.

Numerous models exist for the creation or detection of forgeries, yet they continue to exhibit vulnerabilities. This subsection will delineate the primary problems, item by item, in the creation or detection of DeepFakes.

#### A. CHALLENGES FOR DEEPFAKE CREATION

Notwithstanding considerable advancements in enhancing the visual fidelity of generated DeepFakes, certain obstacles need to be addressed. Challenges associated with the creation of DeepFakes include generalization, temporal coherence, lighting constraints, lack of realism in the eyes and lips, hand movement dynamics, and identity leakage.

**Generalization:** The attributes of generative models are contingent upon the nature of the dataset used during training. Consequently, when completing training on a certain dataset, the output generated by the model embodies the acquired features (fingerprint). Moreover, the quality of the output is contingent upon the amount of the dataset used during training. Therefore, to produce high-quality output, the model must be provided with a sufficiently vast dataset to attain a certain feature. Furthermore, developing a credible model requires many training hours. Acquiring a dataset with relevant material is often more straightforward; yet, securing sufficient data for an individual victim may be challenging.



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

**Hand movement behavior:** An further concern is that when the target conveys emotion via hand gestures, the DeepFake model struggles to accurately replicate such emotions. Furthermore, this sort of expression collection is constrained; hence, generating this form of DeepFake is arduous.

**Identity breach:** Preserving target identity is a problem when there is a significant disparity between the target identity and the driving identity, as shown in face reenactment tasks when target expressions are influenced by a source identity. The driving identity

**Temporal coherence:** Additional defects include discernible irregularities, such as flickering and jittering between frames. The issues arise because the DeepFake generation frameworks process each frame independently, neglecting temporal consistency. To address these deficiencies, several researchers provide contextual information to the generator or discriminator, contemplate temporal coherence losses, use recurrent neural networks (RNNs), or utilize a mix of these methodologies.

**Lighting requirements:** The majority of accessible DeepFake datasets are generated in a controlled setting with uniform lighting and backdrop conditions. Nevertheless, an abrupt change in lighting conditions in indoor or outdoor environments results in color inconsistencies and peculiar anomalies in the final product.

**Deficiency of realism in ocular and labial features:** The absence of authentic emotions, interruptions, and the speed of the target's speech are the main challenges in the construction of DeepFakes based on eye and lip synchronization. Facial data about eye blinking is partly sent to the artificial face. This occurrence transpires when training is conducted on either a single identity or several identities, while data pairing is executed on the same identity.

### B. CHALLENGES FOR DEEPPFAKE DETECTION

Despite substantial advancements in the efficacy of DeepFake detectors, several challenges pertaining to the existing detection algorithms need resolution. DeepFake detection systems encounter many challenges, including insufficient datasets, unidentified media attack types, temporal aggregation issues, and unlabeled data.

**Lack of DeepFake datasets:** Current DeepFake detection algorithms use binary frame-level classification, which entails assessing the authenticity of each individual video frame as either genuine or fabricated. Nonetheless, these approaches neglect inter-frame temporal consistency, potentially resulting in problems such as temporal anomalies and the presence of actual or fabricated frames in successive intervals. Moreover, these approaches need an additional step to calculate the video integrity score, which must be aggregated for each frame to get the final outcome.

**Temporal Aggregation:** Current DeepFake detection algorithms use binary frame-level classification, which entails assessing the authenticity of each individual video frame as either genuine or fabricated. Nonetheless, these approaches neglect inter-frame temporal consistency, potentially resulting in problems such as temporal anomalies and the presence of actual or fabricated frames in successive intervals. Moreover, these approaches need an additional step to calculate the video integrity score, which must be aggregated for each frame to get the final outcome.

**Unlabeled data:** DeepFake detection algorithms typically develop using extensive datasets. Nevertheless, in some instances, such as journalism or law enforcement-related DeepFake detection, only a limited dataset may be accessible. Furthermore, this type of collection necessitates additional work to annotate the score associated with the specific type of forgery used. Therefore, further research is necessary to comprehend forgery situations related to media or law enforcement.



## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### IV. PROPOSED SYSTEM DIAGRAM

#### List of Modules and Functionality

##### Data Collection:

- **Data Collection:** Gather a diverse dataset of real and deepfake images. Ensure the dataset includes a variety of sources to improve model robustness.
- **Data Annotation:** Label the images as 'real' or 'deepfake' and possibly include additional metadata (e.g., source, creation date).

##### Preprocessing-Convert

- Normalize images, resize them to a consistent size, and apply data augmentation techniques to enhance model performance.
- RGB to HIS color space and apply local contrast enhancement for the distribution of the values of pixels around local mean to make further segmentation easy.
- Image segmentation partitions image pixels on the basis of one or more selected image features, e.g., color. Pixels with distinct color are separated into different regions.

##### Feature Extraction

In order to classify the segmented regions into exudate and non-exudate, they must be represented with relevant and significant features to give them the best possible class severability. False positive regions such as light reflections, cotton wool spots, and most importantly, optic discs, should be sorted out. Optic disc localization can be automated by the methods described above. The segmented regions can be differentiated using features such as color, size, edge, and texture.

**Face Detection:** Use algorithms like MTCNN or OpenCV's Haar Cascades to detect faces within images.

**Feature Extraction:** Extract features from the detected faces using techniques like convolutional neural networks (CNNs) or pre-trained models (e.g., VGG16, ResNet).

##### Classification

You can use CNN for classification, with its input nodes matching the feature set and the output node providing the final classification probability. Start with convolutional neural networks (CNNs) like VGG16, ResNet, or Inception for feature extraction. To distinguish between real and deepfake faces, use a classification head. You can apply classification to both structured and unstructured data, dividing a set of data into categories. The deepfake classification for human face detection uses a pre-trained model. The proposed deepfake algorithm follows specific steps to distinguish between real and fake faces. The first step involves generating motion patterns based on the collectiveness descriptor. Next, we construct the transformation space for these motion patterns using frame pairs. Typically, the frequency of occurrence of patterns in the transformation space determines the classification of a face as real or fake.

### V. CONCLUSION

This project utilized deep learning techniques to tackle the problem of identifying and categorizing deepfake videos, especially those featuring human faces. The rapid advancements in deepfake generation have made it increasingly difficult to discern between real and manipulated content, posing significant risks to privacy, security, and trust in digital media. We developed a robust classification model for identifying deep fakes through the application of deep learning methods, particularly convolutional neural networks (CNNs).

### VI. FUTURE SCOPE

The future scope of deepfake classification using deep learning holds immense potential for safeguarding digital environments, maintaining media integrity, and countering misinformation. The field is rapidly evolving, and further research, collaboration, and technological innovation will be crucial in staying ahead of emerging threats.





## International Journal of Innovative Research in Computer and Communication Engineering (IJIRCCE)

(A Monthly, Peer Reviewed, Refereed, Scholarly Indexed, Open Access Journal)

### REFERENCES

- [1] Al-Shabi, Mohammed, and Anmar Abuhamdah. "Using deep learning to detect abnormal behavior in the internet of things." *International Journal of Electrical and Computer Engineering* 12.2 (2022): 2108.
- [2] Fang, Meng-ting, et al., "Examination of abnormal behavior detection based on improved YOLOv3." *Electronics* 10.2 (2021): 197.
- [3] Al-Dhamari, Ahlam, Rubita Sudirman, and Nasrul Humaimi Mahmood. "Transfer deep learning along with a binary support vector machine for abnormal behavior detection." *IEEE Access* 8 (2020): 61085-61095.
- [4] Kim, Kookjin, et al., "Detecting Abnormal Behaviors in Dementia Patients Using Lifelog Data: A Machine Learning Approach." *Information* 14.8 (2023): 433.
- [5] Savenkov, Pavel A., and Alexey N. Ivutin. "Methods of machine learning in system abnormal behavior detection." *Advances in Swarm Intelligence: 11th International Conference, ICSI 2020, Belgrade, Serbia, July 14–20, 2020, Proceedings 11*. Springer International Publishing, 2020.
- [6] Fan, Zheyi, et al., "Real-time and accurate abnormal behavior detection in videos." *Machine Vision and Applications* 31 (2020): 1-13.
- [7] Xia, Limin, and Zhenmin Li. "A new method of abnormal behavior detection using LSTM network with temporal attention mechanism." *The Journal of Supercomputing* 77 (2021): 3223-3241.
- [8] Mokhtari, Sohrab, et al., "A machine learning approach for anomaly detection in industrial control systems based on measurement data." *Electronics* 10.4 (2021): 407.
- [9] Rabbani, Mahdi, et al., "A review on machine learning approaches for network malicious behavior detection in emerging technologies." *Entropy* 23.5 (2021): 529.
- [10] Al-amri, Redhwan, et al., "A review of machine learning and deep learning techniques for anomaly detection in IoT data." *Applied Sciences* 11.12 (2021): 5320.
- [11] Altaei, Mohammed Sahib Mahdi. "Detection of Deep Fake in Face Images Using Deep Learning." *Wasit Journal of Computer and Mathematics Science* 1.4 (2022): 60-71.
- [12] Zobaed, Sm, et al., "Deepfakes: Detecting forged and synthetic media content using machine learning." *Artificial Intelligence in Cyber Security: Impact and Implications: Security Challenges, Technical and Ethical Issues, Forensic Investigative Challenges* (2021): 177-201.
- [13] Sharma, Jatin, et al., "Deepfakes Classification of Faces Using Convolutional Neural Networks." *Traitement du Signal* 39.3 (2022).
- [14] Masood, Momina, et al., "Classification of Deepfake Videos Using Pre-Trained Convolutional Neural Networks." *2021 International Conference on Digital Futures and Transformative Technologies (ICoDT2)*. IEEE, 2021.
- [15] Korshunov, Pavel, and Sébastien Marcel. "Deepfake detection: humans vs. machines." *arXiv preprint arXiv:2009.03155* (2020).





INTERNATIONAL  
STANDARD  
SERIAL  
NUMBER  
INDIA



# INTERNATIONAL JOURNAL OF INNOVATIVE RESEARCH

IN COMPUTER & COMMUNICATION ENGINEERING

 9940 572 462  6381 907 438  [ijircce@gmail.com](mailto:ijircce@gmail.com)



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

Scan to save the contact details