

FAKE VIDEO DETECTION USING DEEP LEARNING

Kalaivani N¹, Padmapriya P N², Maria Rijutha Robert³, Jamuna Eshwar R⁴

Associate Professor, Department of Electronics and Communication Engineering, Sri Krishna College of Engineering and Technology, Coimbatore, India¹

UG Scholar, Department of Electronics and Communication Engineering, Sri Krishna College of Engineering and Technology, Coimbatore, India^{2,3,4}

Abstract: Rapid advancements in AI, machine learning, and deep learning over the past few decades have led to the development of new methods and tools for altering multimedia. Despite the facts that technology has primarily been utilized for good reasons, including entertainment and education, unscrupulous people have nonetheless taken advantage of it for illegal or sinister ends. For instance, realistic-seeming, high-quality phony films, pictures, or sounds have been produced with the intention of propagandizing false information, inciting hatred and political unrest, or even harassing and blackmailing individuals. Recently, the highly-reproduced, lifelike, and altered videos have come to be known as Deepfake. Since then, a number of strategies have been detailed in the literature to address the issues brought up by Deepfake. By safeguarding data, identifying deepfakes, and preventing media manipulation, deepfake video detection contributes to cybersecurity. Videos that are the original data are altered for a number of reasons. It's critical to be able to spot this kind of misleading information. In the social media age, identity theft is seen as the main issue. In order to explore the most promising new methods for deepfake video detection, this paper examines the most recent research findings from the community. This system uses convolutional neural networks (CNNs) and long term memory (LSTM) to distinguish between real and fake video frames. This also involves the application of the Densenet algorithm, XGBoosting classifier, and YOLO Face detector. Faces in videos can be found using the YOLO face detector. To help detect visual artifacts in the video frames, InceptionResNetV2 CNN is used to extract discriminant spatial features of these faces. The XGBoost classifier uses these visual characteristics to assist differentiate between real and deepfake films.

Keywords: Fake video, YOLO, CNN, deep learning.

I. INTRODUCTION

Derived from Deep Learning and Fake, the phrase "Deepfake" refers to some photo-realistic video or image contents produced with the aid of deep learning. Any digital material where the subject's likelihood has been manipulated by editing is referred to as deepfake. Several deepfake videos have been shared on social media as a result of the necessary technology' general availability. Deepfake has also been used to spread false information and rumors to politicians. In order to determine the most successful new strategies, the author has been reviewing and assessing recent papers on deep learning algorithms for deepfake video recognition. It examined the outcomes of two tests that contrasted the abilities of CNN and LSTM to distinguish between real and phony video frames [1]-[2]. Deepfakes, artificial intelligence-driven synthetic media in which people's faces are substituted for real people's in already-existing photos or videos are becoming more and more common. Deepfake technologies are developing at a rapid pace, making it harder and harder to distinguish between real and fake media. This poses serious risks to the security, privacy, and reliability of information. Many deepfake detection approaches require assistance due to over fitting problems, high computational cost, and limited generalizability across different deepfake algorithms and datasets. Furthermore, most existing techniques ignore the potential use of other material, such as text and audio, in favor of video and image-based deepfakes [3]-[5].

Given these difficulties, a reliable, effective, and all-encompassing deepfake detection solution that can manage different media kinds and deepfake algorithms is desperately needed. Furthermore, it's critical that our detection techniques change along with the advancement of deepfake technology. These days, deep fakes are so good at mimicking voices, lip movements, and even facial emotions that it is difficult to tell them apart from real videos. Moreover, the broad range of deep fake generating techniques and their ongoing development pose a challenge to the creation of a detection algorithm that can be used everywhere. With the rise in popularity of social media platforms like Facebook, Twitter, and YouTube, as well as the accessibility of smartphones with highly superior cameras, creating, sharing, and editing movies and photographs has become easier than ever. Concerns about public privacy have recently been raised by the widespread

distribution of incredibly lifelike fake photos and movies produced using the deepfake technology on various social media platforms. Deepfake is a deep learning approach that can construct a video of the target saying or doing things that the source person says or does by substituting the target's face images in the source's video [6]-[8].

Deepfake technology is harmful because it may be used to fabricate videos of leaders, discredit famous individuals, provide misleading news that confuses investors, and trick consumers. To name a few applications of deep learning techniques, they can create faces, switch faces between two people in a video, modify facial expressions, change a subject's gender, and modify facial traits. The significant problem of face-manipulation in images and films is a threat to global security. In human interactions as well as biometric-based services for identification and authentication, faces are crucial. Therefore, trust in digital communications and security applications can be destroyed by plausible changes of face frames [9]. Therefore, a key component in identifying fakes is analyzing and identifying faces in images or videos. When a Reddit user inserted celebrity faces into pornographic films in 2017, the first deepfake video appeared. As a result, a number of deepfake video detection techniques were introduced. Some of these techniques use recurrence networks to identify temporal discrepancies between face frames in films, while others use convolution networks to identify visual distortions inside frames. Finding these kinds of alterations and differentiating them from authentic videos or images is the aim of deepfake detection.

II. LITERATURE SURVEY

Shruti Agarwal and Tarek El-Gaaly, in the paper "Detecting Deep-Fake Videos from Appearance and Behavior (2021)" examined Computer-generated sounds and images, also referred to as "deep fakes," never cease to awe the computer graphics and computer vision communities. Simultaneously, the democratization of technology access, which enables anyone to produce expertly manipulated videos of anyone saying anything, remains concerning due to its potential to sabotage democratic elections, engage in small- to large-scale fraud, support misinformation campaigns, and produce non-consensual pornography.

Mubarak Almutairi and Ali Raza, in the article "A Novel Deep Learning Approach for Deep Fake Detection (2021)" outlined the use of deepfake in synthetic media to create phony audio and visual content based on a user's preexisting media. To make a deepfake appear realistic, fake material is used in place of the subject's voice and face. The creation of fake media information is immoral and dangerous for society. Deepfakes are being utilized extensively in cybercrimes such as identity theft, cyber extortion, financial fraud, celebrity blackmail via fake obscenity films, and many other crimes. Over 96% of the deepfakes have filthy content, according to a new Sensity analysis. The majority of victims are from South Korea, the United States, Canada, India, and the United Kingdom. Cybercriminals created fictitious audio recordings of a chief executive officer in 2019 in order to call his company and request a \$243,000 transfer.

Shraddha Suratkar and Faruk Kazi, in the publication "Deep Fake Video Detection Using Transfer Learning Approach (2021)" The idea that fake news may be quickly disseminated online highlights the need for computational tools in the battle against it. Deepfakes, another name for fake videos, are quite frightening and cause a variety of social and political behaviors in society. It can also be employed with malicious purpose. Because deep fake generation algorithms are readily available on cloud platforms for a low cost of computation power, realistic-looking fake photos or videos can be produced. However, because it is increasingly difficult to mask the tampering using different techniques, it is more important to identify bogus content.

Sonali Gandhi and Monali Gandhi, in the paper "A Qualitative Survey on Deep Learning Based Deep Fake Video Creation and Detection Method (2022)" examined how the use of Deep Learning (DL)-based applications is expanding quickly in the contemporary world. Numerous important issues, including computer vision, massive data processing, and brain interface, are resolved by deep learning. The development of deep learning may also pose challenges to national security, democracy, and privacy on both a domestic and international level. Deepfake films are spreading so quickly that they are affecting social, political, and private spheres. Artificial intelligence is used to create deepfake videos, which even to a skilled eye can seem extremely real. Deepfakes are frequently used to create pornographic films, which damages people's reputations. The public is concerned about deepfakes, so it's critical to provide techniques for identifying them.

Anuj Badale, Chaitanya Darekar and Lionel, in the paper "Deep Fake Detection Using Neural Networks(2022)" discussed the deepfake artificial intelligence-based method for creating human images. Using machine learning techniques, Deepfake is used to combine and superimpose pre-existing images and movies onto source photos or videos. These are phony videos that look authentic and are indistinguishable with the naked eye. They can be used to agitate political unrest, propagate hate speech, extort someone, and more. To verify the authenticity of videos, cryptographic signature of the videos is currently done. It is validated whether or not a video file is the original recording by hashing it into fingerprints, which are short text strings, and then comparing the resultant file to the sample video. But the issue with this method is that hashing techniques and fingerprints aren't readily available with standard people.

III. PROPOSED ARCHITECTURE

The system consists of several modules include data selection, preprocessing, segmentation, data splitting, classification, prediction and performance analysis. It is depicted in Fig. 1.

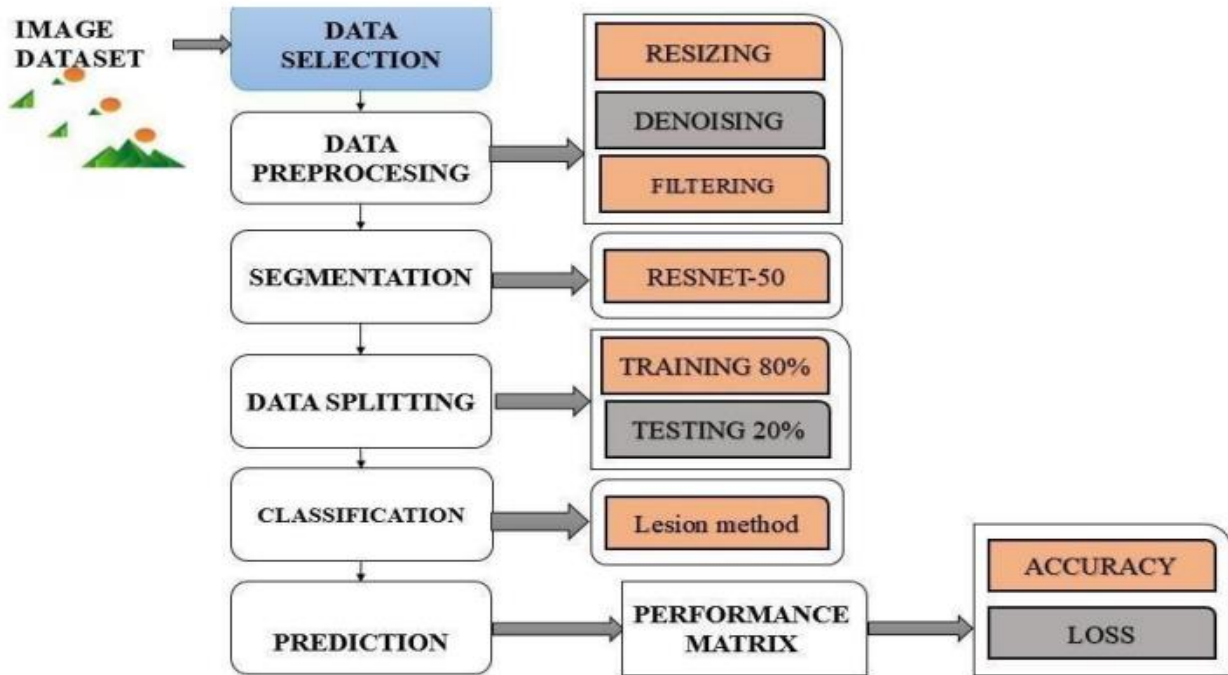


Fig. 1. Proposed architecture of deepfake image identification

An innovative computer vision system called YOLO (You Only Look Once) was created to effectively and instantly identify faces in photos and video streams. It was created as a member of the YOLO family of object identification models and uses a single neural network to categorize objects inside bounding boxes and predict bounding boxes at the same time. The YOLO face detector is very quick and appropriate for real-time applications, as seen in Fig. 2, because it can process full images in a single forward pass. YOLO offers quick and accurate face detection by partitioning an image into a grid and forecasting bounding boxes and related class probabilities. Researchers and developers can select the version of the model that best fits their needs by selecting from a variety of variants, including YOLOv3 and YOLOv4, which balance speed and accuracy for a variety of face detection applications.

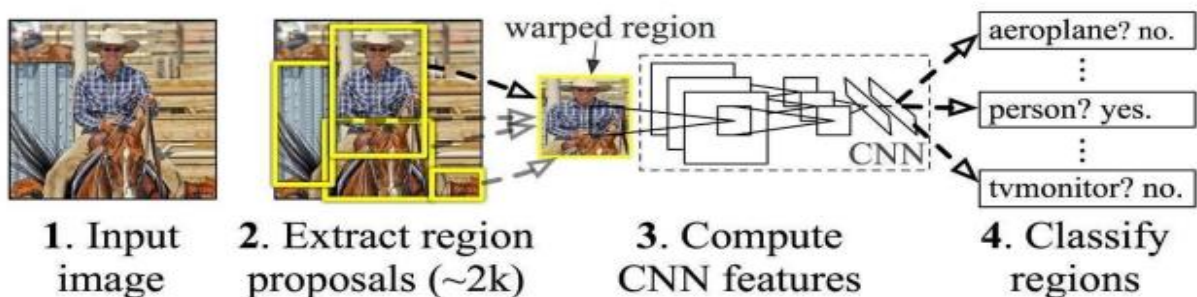


Fig. 2. YOLO Face detector

A. Transfer Learning Neural Network Techniques

Convolutional neural network (CNN) designs have seen a radical paradigm change with the introduction of DenseNet, short for Dense Convolutional Network, especially in the field of computer vision. DenseNet, which was first presented by Gao Huang, Zhuang Liu, and Laurens van der Maaten in 2017, provides a revolutionary solution to training deep neural networks while improving parameter efficiency and feature propagation at the same time.

Fundamentally, DenseNet creates densely connected routes throughout the network by creatively connecting each layer to every other layer in a feed-forward manner. Deeper designs are made possible by this dense connection pattern, which promotes feature reuse and information flow without creating problems with vanishing gradients.

A pre-trained neural network method typically utilized for image identification applications is the VGG16. The VGG16 model was proposed by K. Simonyan and A. The convolution neural net (CNN) architecture is the foundation of the VGG16. The VGG16 model architecture debuted in the ILSVR competition in 2014. Using our dataset, we developed the VGG16 model for deepfake detection. Convolutional neural network (CNN) architectures have reached a new height with Inception-ResNetv2, which combines the best features of both Inception and ResNet to provide a deep learning model that is both incredibly accurate and efficient.

B. Convolutional Neural Networks

Because they can learn hierarchical representations directly from unprocessed input data, Convolutional Neural Networks (CNNs) are a mainstay in the field of deep learning, especially in computer vision problems. Semantic segmentation, object detection, and picture recognition have all been transformed by CNNs, which were developed using neuroscience concepts and the structure of the visual brain as inspiration. As seen in Fig. 3, a CNN's architecture consists of several layers, including convolutional, pooling, and fully linked layers.

Convolutional layers use filters or kernels to extract features (i.e., spatial patterns and local dependencies) from input images through convolutions. The feature maps are then down sampled by pooling layers, which lowers computational complexity and spatial dimensions while keeping significant features. Furthermore, fully connected layers combine features that have been collected to carry out tasks like regression or classification. Because of their hierarchical structure, CNNs can recognize complex patterns and objects in images by learning progressively abstract information in deeper layers.

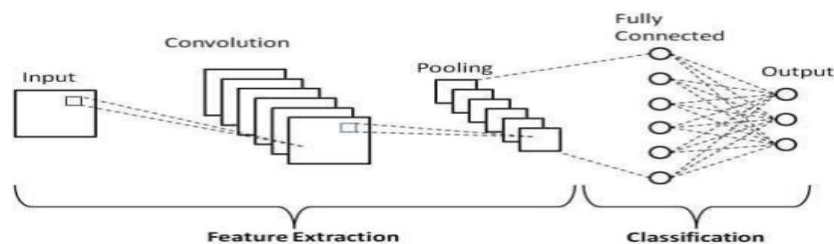


Fig. 3. CNN layers

IV. METHODOLOGY

The methodology of the proposed design is depicted in Fig. 4.

A. Data Selection

Data selection is the process of deciding on the appropriate data source, kind, and gathering technologies. Prior to actually collecting data, a process known as data selection selects and retrieves pertinent data from the data collection technique. This experiment uses the fake video classification dataset. A dataset in the context of deepfake video detection is an assortment of carefully selected and annotated movies or images used to train, validate, and test machine learning models intended to identify deepfake material. These datasets are essential for the creation and assessment of algorithms for detecting deepfakes. Fig. 5 displays the image from the sample dataset.

B. Image Preprocessing

By preparing the data being input in a format that the model of machine learning can use efficiently, data preprocessing is essential to the identification of deepfake videos. Preprocessing techniques including data augmentation, feature extraction, frame extraction, and picture preprocessing help to increase the quality of data being processed and the detection mode's performance. The technique of eliminating noise from the picture while keeping its key elements is known as denoising. Numerous denoising techniques are available in Scikit-Image, including as Non-local Means Denoising is a technique that uses a non-local technique algorithm to eliminate Gaussian noise from photographs. The frequency spectrum of an image can be altered by filtering processes, which can help highlight specific aspects or eliminate undesired information. Among the filtering algorithms offered by Scikit-image is Gaussian Filtering, which uses a Gaussian kernel to convolve pictures and smooth them out. This lowers high-frequency noise and blurs the image.

The value of each pixel is replaced with the average value inside a local by using median filtering. It works great for keeping edges intact and eliminating salt-and-pepper noise.

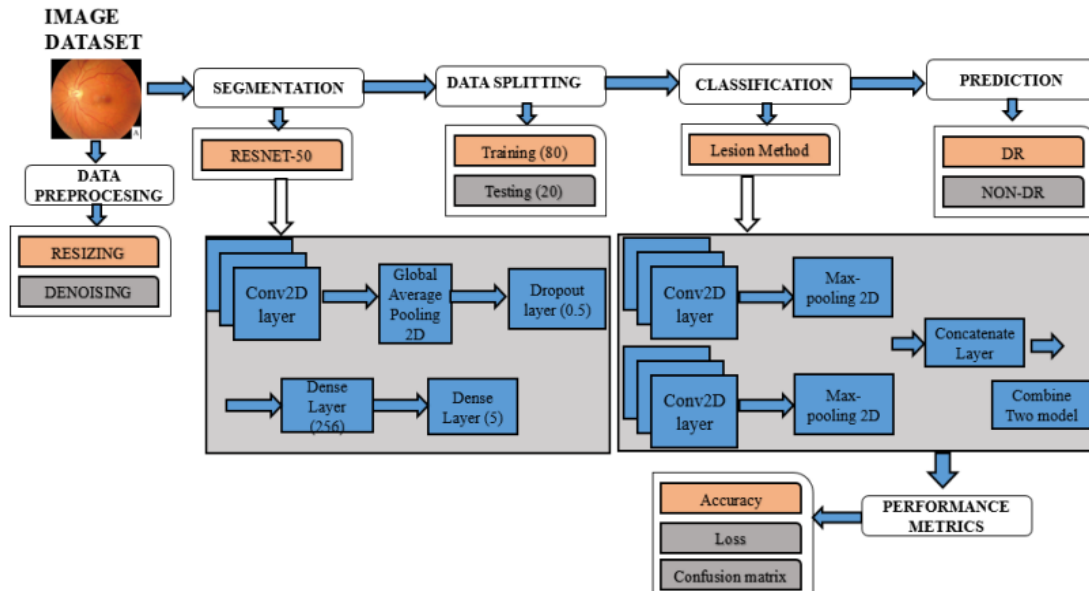


Fig. 4. Flow diagram of the proposed methodology



Fig.5. Dataset sample images

C. Feature Extraction

To locate and extract pertinent features from an image, feature extraction approaches are applied. Applications such as recognizing objects and image classification can make use of these properties. Texture analysis, corner detection, and edge detection are a few common feature extraction methods. Utilize the fully connected layers of the DenseNet-201 model design to extract attributes from each one. An overview of the material in every frame is obtained by running each frame over the network. Utilizing segmentation techniques, categorize each pixel in the frame as falling to one of the classes (genuine or false) after identifying features from each frame. To accomplish pixel-wise categorization, this stage might need adding more layers or components on top over the densenet 201 backbone.

D. Segmentation

Depending on the content of each region, an image can be divided using segmentation techniques. In applications like medical imaging, where certain organs or tissues need to be excluded from the image, fake can be useful. Thresholding, identifying edges, and region expanding are a few common segmentation methods. Use CNN for Fake in this instance.

The conv2d layer by layer from the already established model is used in the layers of Resnet-50's design. It is then applied to the Global Pooling Average layer, where unwanted neurons are eliminated using the Dropping out layer at a value of 0.5. Finally, the dense layer, which has the highest number of neurons, is designed with specific class neurons in mind. ResNet-50's architecture must be modified for the unique purpose of pixel-wise classification, trained on labeled data, and evaluated on unseen data in order to be used for deepfake video segmentation. It's crucial to remember that, while though ResNet-50 can be an efficient tool for detecting deepfakes, effective detection may require its combination with other methods and strategies.

E. Data Splitting

Using some sort of homogeneity criterion, data division is a technique where an image is first split into smaller sections by dividing it recursively into a single region. When data is broken up into two or more subgroups, it's called data splitting. When using a two-part split, the data is usually evaluated or tested in part, while the model is trained in the other half. A crucial component of data science is data splitting, especially when building models with data. Using the 'validation_split' option, you may divide your data into sets for testing and training employing the 'ImageDataGenerator' class via TensorFlow's Keras API. You can specify the percentage of information to hold back for validation using this parameter.

F. Classification

In the Lesion Method of classification, these two densenet-201 models are combined by concatenating all layers to generate the completely connected layer, which is formed by combining the max-pooling layer and the conv2d layer. Sort the classes finally based on whether they contain fake or authentic videos. The Lesion approach is applied during the categorization phase. This technique is employed to examine the deep learning model's resilience and susceptibility to hostile assaults. Using a particular architecture, the classification model consists of two branches, each containing convolutional layers and max-pooling layers. The outputs of both branches are then concatenated and fed into additional layers for classification. The examples of the categorized real and fake photographs are displayed in Figs. 6 and 7.



Fig. 6. Samples of real images

Fig. 7. Samples of fake images

G. Prediction

The goal of predictive analytics algorithms is to minimize errors. "DR" and "NON-DR" probably relate to distinct kinds of prediction blocks that are employed in the detection model when discussing deepfake video detection. Presumably, a prediction block topology that integrates residual connections is referred to as a Deep Residual (DR) forecast block. Since the ResNet design made residual connections popular, many machine learning models, particularly those used in computer vision tasks, have included them as a standard feature. Deep network optimization is facilitated by residual connections in DR prediction blocks, which enable the network to develop residual functions. In order to accomplish this, shortcut connections that add the input to the conclusion of those levels while omitting one or more layers are added. This facilitates the training of larger networks and helps to mitigate the vanishing gradient issue. The prediction block that doesn't use residual connections is known as a Non-Deep Recall (NON-DR) prediction block. Stated differently, it could potentially comprise an array of convolutional layers devoid of any skip connections. NON-DR architectures can still be utilized in some situations, even though DR architectures are well-known for their capacity to train very deep networks. This is particularly true in situations where the task at hand does not call for extremely deep coalitions or when the information being used is too small to fully utilize the advantages of deep residual learning.

Either DR or NON-DR prediction blocks could be incorporated into the overall model architecture for video analysis in the framework of deepfake video detection.

V. EVALUATION RESULTS

Similar to any classification task, deepfake video detection uses a number of performance indicators to assess the model's efficacy. The percentage of correctly identified samples relative to the overall amount of samples is known as accuracy. Accuracy in deepfake detection is how frequently the model properly determines if a video is modified or real. But if the dataset is unbalanced (i.e., there are a lot more real videos than deepfake films), accuracy might not be enough on its own because the model might get high accuracy by just labeling all of the videos as real. Loss measures how closely predictions generated by the model match the labels that are present in the training data. It is often referred to as the error rate or objective function.

Depending on the type of task, other loss functions (such as binary cross entropy for binary classification) can be applied. The objective of training is to minimize the loss, a sign that the model is improving its prediction accuracy. A table that lists a classification model's performance is called a confusion matrix. The counts of forecasts that are false positive (FP), false negative (FN), true positive (TP), and true negative (TN) are displayed.

- True positive (TP): Videos that are properly identified as deepfakes.
- True negative (TN): Real videos that are accurately identified as real.
- False positive (FP): When real films are mistakenly labeled as deepfakes (Type I error).
- False negative (FN): Videos that are deepfake but are mistakenly identified as real (Type II error)

Numerous other metrics can be obtained from the confusion matrix: Precision can be defined as the ratio of true positives to all positive predictions (TP / (TP + FP)). The precision of the model is the degree to which genuine videos are not mistakenly identified as deepfakes. The percentage of true positive predictions (TP / (TP + FN)) among all real positive samples is known as recall (sensitivity). The recall of the model indicates how well it can detect deepfake videos. **F1 Score:** $2 * (accuracy * Recall) / (Precision + Recall)$ is the harmonic mean of accuracy and recall. Recall and precision are balanced by the F1 score. **Specificity:** $TN / (TN + FP)$ = genuine negative predictions as a percentage of all real negative samples. The model's specificity gauges how well it can recognize real footage.

Table I Parameters for analyzing the detection accuracy

ACCURACY (TP+FP)/(TP+FP+TN+FN)	PRECISION (TP/TP+FP)	RECALL (TP/TP+FN)	PREVALENCE (TP+FP)/(TP+FP+TN+FN)	F1-SCORE 2(RECALL*PRECISION)/(RECALL+PRECISION)
(170+11)/193 =181/193 =0.937 =93.7%	170/(170+11) =170/181 =0.939 =93.9%	170/(170+23) =170/193 =0.880 =88.0%	(170+11)/193 =181/193 =0.937 =93.7%	2(0.88*0.93)/(0.88+0.93) =2(0.818)/1.81 =0.903 =90.3%

The Findings from the confusion matrix is shown in table I and are given by
 Accuracy = 93.7% (i.e. TP/TP+FN)
 Misclassification Rate = 17% (i.e. FP+FN/TOTAL)
 True Positive Rate = Recall = Sensitivity = 88% (i.e. TP/TP+FN)
 False Positive Rate = 26% (i.e. FP/FP+TN)
 True Negative Rate = Specificity = 73% (i.e. TN/TN+FP)
 False Negative Rate = 11.9% (i.e. FN/FN+TP)
 Precision = 93.9% (i.e. TP/TP+FP) Prevalence = 93.7% (i.e. TP+FP/TOTAL)
 F1-SCORE = 90.3% (i.e. 2(RECALL*PRECISION)/RECALL+PRECISION)

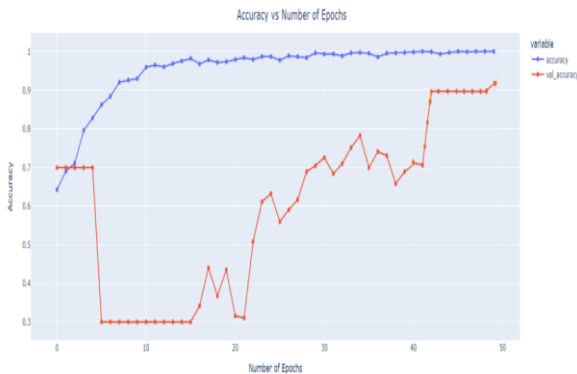


Fig. 8. Accuracy vs Number of Epoch

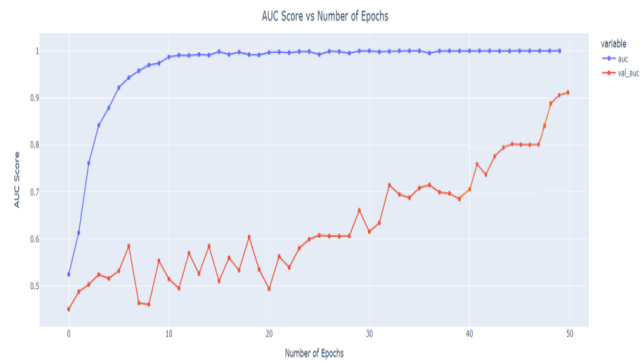


Fig. 9. AUC Score vs Number of epochs

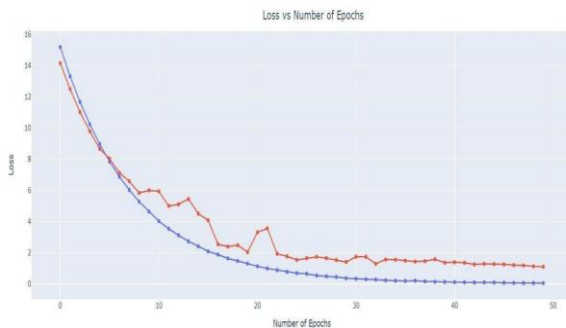


Fig. 10. Loss vs Number of epochs

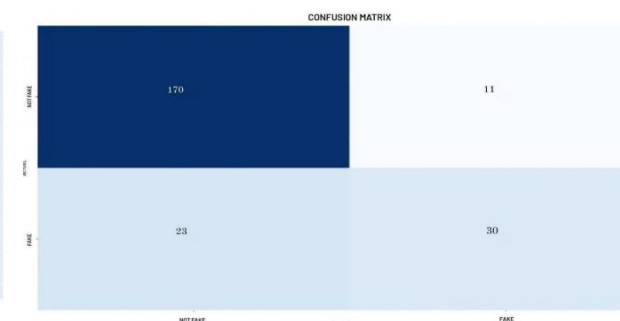


Fig. 11. Confusion Matrix

Fig. 8, Fig. 9, Fig. 10 and Fig. 11 shows the various comparative graphs based on epoch, accuracy, AUC score, loss and confusion matrix.

Total dataset taken is 1052 and for testing 193 is used for a Epochs Run of 50. These performance measures offer a thorough understanding of the behavior of the model and aid in evaluating how well it can identify deepfake movies. Certain metrics might be more significant than others, depending on the application's particular requirements and objectives. For example, reducing false positives might be important in some situations, but increasing recall—or identifying as many counterfeits as possible—might be more important in others.

VI. CONCLUSION

A deep learning model is proposed for deep learning-based deep fake video detection. Utilizing CNN in the fictitious module improves accuracy, while utilizing fully connected layers in the lesion model reduces loss. Based on retrieved photos, the deep learning model has promising performance in identifying deepfake films.

Performance may be enhanced by additional model architecture and hyper parameter adjustments and optimizations. It takes constant observation and updating to adjust to new deepfake methods.

REFERENCES

- [1] N. Gardiner, Facial re-enactment speech synthesis and the rise of the Deepfake, 2019.
- [2] N. Sambhu and S. Canavan, "Detecting Forged Facial Videos using convolutional neural network", 2020.
- [3] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen and T. Aila, "Analyzing and improving the image quality of stylegan", 2019.
- [4] C. Hsu, Y. Zhuang and C. Lee, "Deep fake image detection based on pairwise learning", MDPI, pp. 1-14, January 2020.
- [5] E. Cueva, G. Ee, A. Iyer, A. Pereira, A. Roseman and D. Martinez, "Detecting fake news on twitter using machine learning models", Soe.rutgers.edu, pp. 1-12, July 2020.



- [6] H. Chen, K. Zhang, S. Hu, S. You and C. Kuo, "Geo-DefakeHop: high-performance Geographic fake image detection", arXiv.org, pp. 1-12, 2021.
- [7] J. Frank, T. Eisenhofer, L. Schonherr, A. Fischer, D. Kolossa and T. Holz, "Leveraging frequency analysis for deep fake image recognition", arXiv.org, vol. 119, pp. 1-12, 2020.
- [8] I.M.V. Krishna and S. Sai Kumar, "Fake News Detection Using Naive Bayes Classifier", Proceedings of International Journal of Creative Research Thoughts, pp. e757-e761, 2021.
- [9] Miki Tanaka, Sayaka Shiota and Hitoshi Kiya, "A detection method of operated fakeimages using robust hashing", Journal of Imaging, vol. 7, no. 8, pp. 134, 2021.
- [10] Zonghan Wu, Shirui Pan, Fengwen Chen, Guodong Long, Chengqi Zhang and S. Yu Philip, "A comprehensive survey on graph neural networks", IEEE transactions on neural networks and learning systems, vol. 32, no. 1, pp. 4-24, 2020.