

A Comparative Study of Deepfake Detection Using ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder

¹T S Harikrishnan, ²Sebin Thomas, ³Neha Simon, ⁴Sanjo Johny ⁵Jintu Ann John

¹Student, ²Student, ³Student, ⁴Student, ⁵Assistant Professor

¹Department of Information Technology,

¹Amal Jyothi College of Engineering (Autonomous), Kerala, India

harikrishnan4607@gmail.com, sebin12sebin@gmail.com, nehasimon3214@gmail.com

sanjojohny16@gmail.com, jintuannjohn@amaljyothi.ac.in

Abstract—Deepfake videos have emerged as a significant threat to digital media authentication due to their ability to convincingly alter video content, leading to widespread misinformation and manipulation. This paper presents a comparative study of three advanced deepfake detection models: ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder. The ResNeXt50_32x4d + LSTM model utilizes a hybrid spatial-temporal approach, combining ResNeXt50_32x4d for spatial feature extraction and LSTM for temporal feature modeling, which significantly enhances its ability to detect subtle manipulations across video frames. In contrast, EfficientNet + GRU focuses on computational efficiency with a streamlined architecture, while Xception + Transformer Encoder employs attention mechanisms for long-range dependency analysis in video sequences. The study demonstrates that the ResNeXt50_32x4d + LSTM model consistently outperforms the other two models in terms of accuracy, precision, recall, and computational efficiency. By leveraging transfer learning and a well-structured preprocessing pipeline, ResNeXt50_32x4d + LSTM achieves a higher detection rate by capturing intricate spatial patterns and subtle temporal inconsistencies across frames, making it particularly robust in identifying both real and fake videos. The experimental results show that the ResNeXt50_32x4d + LSTM model achieves an accuracy of 91.88%, a precision of 90.66%, and a recall of 85.60%, surpassing the performance of both EfficientNet + GRU and Xception + Transformer Encoder in terms of deepfake detection. These results establish ResNeXt50_32x4d + LSTM as a superior method for tackling the challenges posed by deepfake technology. The paper concludes by analyzing the advantages, limitations, and trade-offs between these models, suggesting that ResNeXt50_32x4d + LSTM is an optimal choice for real-time deepfake detection applications due to its balanced trade-off between accuracy and computational cost.

Keywords—Deepfake Detection, EfficientNet, GRU, LSTM, ResNeXt50_32x4d, Transformer Encoder, Xception

I. INTRODUCTION (HEADING 1)

Today, the world faces an entirely new challenge: one that involves the rise of Siri-based social media influencers and a lack of proper pseudonym usage, thanks to the astounding advancements in Artificial Intelligence and deep learning technology resulting in the creation of extremely realistic ‘deepfakes’. ‘Deepfakes’ utilize machine learning to create manipulated audio and visual recordings, which can significantly change reality by modifying the subject’s face, voice, and actions, [6]. This alternative type of media, if used incorrectly, can result in detrimental adverse consequences like misinformation, Identity theft, and even illegal activities like reputation harm. Therefore, being able to spot and eliminate the impact of so-called ‘deepfakes’ has been shown to be a crucial part of ‘AI’ study and understanding, [9].

Deepfakes are produced using various techniques, and algorithms that incorporate ‘Generative Adversarial Neural Networks,’ also known as “GANs,” in which two adversary neural networks, one creating fake content while the other tries to identify fake and real content, are utilized. These systems constantly improve until the fakes produced cannot be differentiated from real footage, [10]. From the emergence of hyper-realistic AI generated videos, deepfakes have caused the removal of security, media trust, and even society. The further we extend technology for deepfakes, the further we need to advance detection of deepfake AI technology. Without these partnerships, politics, entertainment, and even social media will be plunged into chaos.

In the wake of these challenges, an array of tools has been created for the detection of deepfakes, many of which are anchored on deep learning models [3]. With the substantial amounts of information available in this generation, deep learning has proved useful in both the creation and detection of altered media. As an example, Convolutional Neural Networks (CNNs) have shown great potential in the classification of images since they classify by spatial feature extraction [5]. Nonetheless, the detection of deepfakes in video data not only requires the extraction of spatial features from single frames, but also the consideration of time across many frames to detect abnormalities over time [18].

Due to the nature of video data, deepfake detection algorithms have to possess the ability to understand both spatial and temporal elements. Spatial analysis is a strong suite of CNNs, enabling them to capture all details of every frame of the video. On the other hand, the more subtle temporal inconsistencies that may span across frames require combining CNNs with Recurrent Neural Networks (RNNs), which are sequence-processed models [2]. Especially, the LSTM variant of RNNs is employed widely in video-based tasks owing to its ability to learn long-range dependencies in sequential data [8].

This paper reviews the comparative analysis of three deepfake detection models: ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder. Each of these models utilizes a unique combination of spatial and temporal feature extraction techniques for deepfake classification and recognition. The comparative study aims to explore how these models fuse spatial and temporal features to improve deepfake detection performance.

The first model, ResNeXt50_32x4d + LSTM, leverages the ResNeXt50_32x4d architecture for spatial feature extraction and the LSTM encoder for capturing temporal dependencies. ResNeXt50_32x4d incorporates grouped convolutions to efficiently extract

deep network features, enhancing its ability to detect subtle artifacts in video frames, while the LSTM captures sequential inconsistencies over time, making it robust for deepfake detection [11].

The second model, EfficientNet + GRU, adopts a more nuanced approach without compromising performance and in some ways balancing GRU's computations. It is evident that EfficientNet focuses on adjusting depth, width and resolution far better than CNNs, achieving extreme performance with lean parameters. For real world applications, GRU's ability to work with timestamps while being significantly less resource consuming than LSTM makes it a popular choice, especially during times when speed is vital [12].

In the Xception + Transformer Encoder model, the Xception architecture is applied, which employs depthwise separable convolutions to reduce the number of computations. This model in conjunction with the Transformer Encoder, which uses and attention mechanism to remember long-distance relationships between video frames, represents a novel way of analyzing video data. Transformer Encoders have been shown to improve the performance of a number of sequence modeling tasks like natural language processing and are now being used at the forefront of video data processing [8].

The primary objective of this paper is to compare the performance of these three models based on accuracy, precision, recall, F1-score, and computational efficiency. Our experiments show that the ResNeXt50_32x4d + LSTM model outperforms the other two models, particularly in terms of accuracy and temporal analysis. We also discuss the complexity of these models and suggest methods to maximize their performance, especially for more intricate deepfake detection tasks.

II. RELATED WORKS

The surrounding community for deepfake detection has come a long way, experimenting with various methods from simple machine learning techniques to complex deep learning models. One of the first and most popular methods to identify deepfakes utilized Convolutional Neural Networks (CNNs) to extract and analyze inconsistencies in a single static frame of a video [10]. Although CNNs were successful in detecting image manipulation, their greatest weakness was that they could not take into consideration temporal dynamics that are so important in video analysis [7]. This lack of consideration gave rise to models that could incorporate temporal information, a move towards more advanced, video-based detection algorithms [17].

The initial study of Cozzolino et al. in "SpoC: Spoofing camera fingerprints" suggested detecting manipulated images and videos using CNNs by recognizing the distinctive patterns of noise imparted by various camera sensors [1]. It was able to identify manipulations at the image level but didn't address inconsistencies at the video level. The deepfake detection methods based on CNN improved further with the presentation of the XceptionNet model by Rössler et al. [13], who employed depthwise separable convolutions to heavily improve computational performance without reducing precision. XceptionNet turned into a much-preferred model for finding manipulated faces within videos, especially when applied with large datasets such as FaceForensics++ [21]. XceptionNet's architecture was optimized to extract fine spatial features, and hence it was well-suited for detecting face manipulation in single video frames. Yet, as discussed above, one of its shortcomings was that it could not handle temporal dynamics, which are crucial in detecting whether the whole video sequence is manipulated or not.

Although the early success of CNNs such as XceptionNet, the necessity to include temporal information soon became apparent. Videos naturally carry sequential information that is lost when static frames alone are processed. Therefore, deepfake detection techniques started to change to include temporal relationships between frames, and thus hybrid CNN-RNN models were developed, which was the next major advancement in deepfake detection research [14].

To address the shortcomings of CNN-only models, Güera et al. suggested the combination of Long Short-Term Memory (LSTM) networks with CNNs for deepfake detection based on videos [2]. LSTM networks, a variant of Recurrent Neural Network (RNN), are specifically engineered to capture temporal sequences well by having a memory of past inputs. In Li's CNN-LSTM hybrid method, CNNs were initially utilized to capture spatial features of every video frame, and LSTMs were subsequently utilized to capture temporal dependencies between successive frames. This hybrid framework enabled the model to pick up on faint temporal artifacts that static CNN models frequently failed to notice.

Moreover, Chen et al.'s [16] work with spatio-temporal feature learning using LSTMs supported the application of combining RNNs and CNNs in action recognition, which is directly relevant to deepfake detection. Their LSTM-based approach to human action recognition in videos set the stage for applying temporal modeling to other video-related applications, such as deepfake detection. By applying LSTMs to capture the sequential nature of video data, researchers were able to extract both spatial and temporal features, leading to a more comprehensive approach to detecting manipulations in video sequences.

The advent of Transformer Encoders in computer vision applications transformed deepfake detection by offering a new method of representing long-range dependencies in videos. Dosovitskiy et al.'s Vision Transformer (ViT) employed self-attention mechanisms to represent both local and global relationships in visual data [10]. While CNNs look into localized patches within images or video frames, Transformer Encoders can model relationships throughout the whole image or video frame, positioning them perfectly to handle tasks needing an awareness of long-range dependencies, e.g., deepfake detection [15].

Transformer Encoders are better suited for video analysis since they are able to handle sequential data effectively by considering all the video frames at once, allowing the model to spot discrepancies that develop over time. Dosovitskiy et al.'s ViT showed state-of-the-art performance across many video-based tasks, such as deepfake detection. Although CNNs and LSTMs are subject to limitations in modeling long-term dependencies, Transformers are themselves well-suited for capturing such information, and hence are especially strong for applications such as identifying subtle patterns that build up over time in deepfakes [10].

Though promising, Transformer-based models have their own suite of problems. Perhaps the biggest disadvantage is their computational expense. Transformers are computationally intensive and memory-intensive, particularly when used for processing extended video sequences. This renders them challenging to adopt for real-time deepfake detection or on limited computational resources like devices. The situation has changed with recent innovations in hardware as well as in model optimization methodologies, which have alleviated some of these limitations and made the use of Transformer-based methods feasible for real-world applications.

The history of deepfake detection algorithms can be broadly described as an evolution from static frame-based CNN models to hybrid CNN-LSTM models that exploit temporal dynamics, and most recently, to Transformer-based models that extract long-range dependencies in video sequences [20]. All of these models have their respective strengths and weaknesses, and the selection of the model is task-dependent.

CNN-based architectures such as XceptionNet perform well at identifying spatial artifacts in single frames and are, therefore, highly effective in tasks requiring high computational efficiency [19]. However, their inability to model temporal relationships makes them perform poorly in video-based tasks. Hybrid CNN-LSTM models solve this by introducing temporal modeling and thus being

able to extract both spatial and temporal features, which are important for identifying deepfake videos. These models provide a balanced solution, blending the power of CNNs and LSTMs to attain high accuracy while keeping computational demands reasonable.

Transformer models are the state of the art for deepfake detection, with their capacity to learn long-range dependencies rendering them extremely effective for video analysis. Nonetheless, their significant computational expense continues to be an impediment to adoption, especially in real-time or low-resource environments. This being said, further research on more efficient Transformer models and hardware acceleration methods will doubtless bring these models closer to practical use in deepfake detection in the near term.

In summary, the technology of deepfake detection has evolved significantly from its initial dependence on CNN-based models. The introduction of temporal modeling using hybrid CNN-LSTM models was a major leap in video manipulation detection, and the development of Transformer-based methods has brought new avenues for modeling long-range dependencies in video data. As the technology advances, it is probable that subsequent models will leverage the best of CNNs, LSTMs, and Transformers to create more effective and precise deepfake detection models that can address the intricacies of real-world scenarios.

III. METHODOLOGY

The three models compared in this work—ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder—were chosen for their ability to capture both spatial and temporal features in video data. These models represent three distinct approaches to deepfake detection, with each combining convolutional neural networks (CNNs) for spatial feature extraction and either recurrent or attention-based networks for temporal analysis. Below, we outline the architecture and training process for each model, as well as the evaluation metrics used.

A. Dataset

We utilized the Celeb-DF-v2 dataset, which contains 580 fake videos and 588 real videos. The dataset was split into 80% for training and 20% for testing, ensuring a balanced representation of real and fake videos across both subsets. Celeb-DF-v2 is preferred due to its high-quality fake videos, which are more difficult to detect compared to older datasets such as DeepFake Detection Challenge and FaceForensics++. The quality of deepfakes in Celeb-DF-v2 closely resembles that of modern deepfake generation tools, making it suitable for evaluating state-of-the-art detection models.

The diversity of the dataset is critical in ensuring that the models generalize well to unseen deepfakes. While the Celeb-DF-v2 dataset is comprehensive and contains high-quality deepfakes, it is limited to celebrity faces, which may not fully represent the broad range of real-world deepfakes. As a result, the models trained on this dataset might struggle when faced with non-celebrity deepfakes or videos created with newer techniques not present in the dataset.

Expanding the training data to include a broader variety of individuals, facial expressions, and deepfake generation techniques would improve model generalizability. Additionally, incorporating adversarial training, where models are exposed to different types of attacks or manipulated media, can further enhance their ability to detect more sophisticated, unseen deepfakes.

To ensure the models generalize well to unseen data, the dataset was split into 80% training and 20% testing subsets. This split maintains the balance between real and fake videos in both the training and testing sets. The testing set was strictly kept aside and used only after the models were fully trained, ensuring unbiased performance evaluation. The dataset statistics are summarized in Table I

TABLE I. DATASET STATISTICS

<i>Dataset Statistics</i>	<i>Real Videos</i>	<i>Fake Videos</i>	<i>Total</i>
Training Set	467	467	934
Testing Set	117	117	234

B. Preprocessing Pipeline

A robust preprocessing pipeline was implemented to standardize the input data for all models. Preprocessing is crucial in deepfake detection as it ensures that the models focus on relevant features, such as facial regions, while ignoring irrelevant parts of the video. The following steps were applied:

1. Face Detection:

In each video frame, faces were identified using the face-recognition package. This package accurately detects faces within each frame and generates bounding boxes around the facial region. The identified face regions were then cropped and resized to a standard resolution, ensuring consistency across all inputs. This process focuses on the most relevant areas for deepfake detection while reducing the input data's dimensionality by excluding irrelevant background information.

2. Cropping and Resizing:

After detecting the face regions, they were cropped and resized to a standard resolution of 112x112 pixels. This consistent resolution guarantees that the models receive input of the same size, allowing for efficient batch processing. Additionally, resizing the frames contributes to lowering the computational burden without significantly compromising image quality.

C. Model Architectures

1. ResNeXt50_32x4d + LSTM

The initial model, ResNeXt50_32x4d combined with LSTM, is a hybrid framework that merges ResNeXt50_32x4d for extracting spatial features and LSTM for modeling temporal dynamics. ResNeXt50_32x4d is a deep residual network that utilizes grouped convolutions to improve computational efficiency while capturing rich spatial features. The model is designed to handle vanishing gradient issues through skip connections, facilitating smooth gradient flow across its 50 layers, even in deep networks. This makes it particularly suitable for extracting high-level features from images, which are essential for deepfake detection.

Once spatial features are extracted by ResNeXt50_32x4d, they are fed into the LSTM (Long Short-Term Memory) network. LSTM is a recurrent neural network (RNN) that excels in capturing long-term dependencies in sequential data. It uses memory cells

to retain important information over multiple time steps, making it ideal for video-based tasks. By analyzing the temporal progression of frames, LSTM helps the model detect subtle temporal inconsistencies that could indicate deepfake manipulations.

This integrated architecture allows ResNeXt50_32x4d to capture detailed spatial features from individual frames, while the LSTM models temporal relationships between consecutive frames. The output of the LSTM is passed to a fully connected layer for final binary classification (real or fake). This architecture combines spatial and temporal modeling to deliver robust deepfake detection. The system architecture and model flow, including preprocessing steps, are depicted in Figure 1, and the details of the architecture are outlined in Table II.

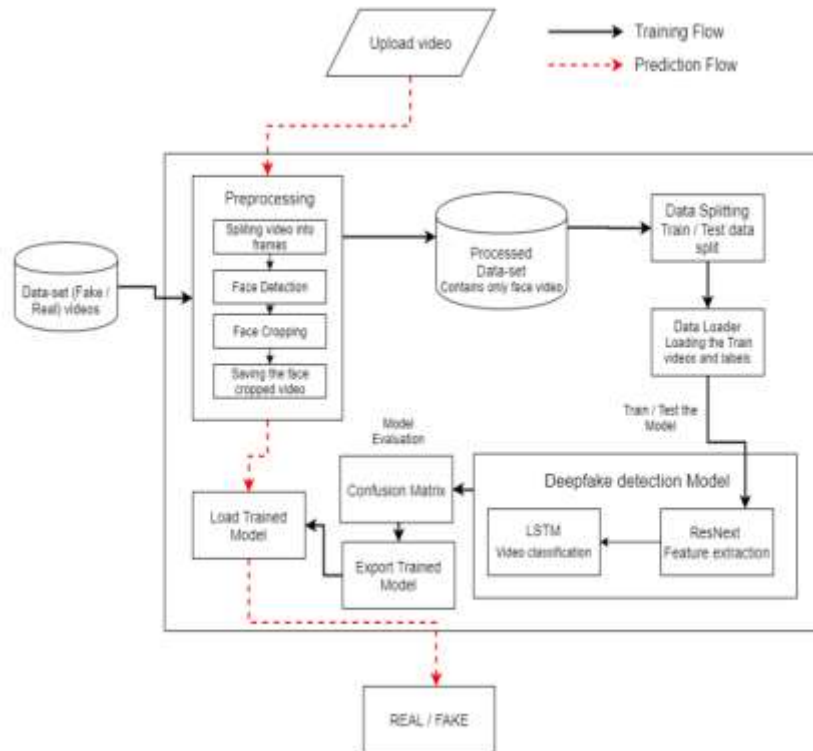


Fig 1. System Architecture for ResNeXt50_32x4d and LSTM

TABLE II. ARCHITECTURE OF RESNEXT50_32X4D + LSTM MODEL

Layer	Output Size	Description
Input	112x112x3	Video frame (RGB)
ResNeXt50_32x4d	4x4x2048	Spatial feature extraction
AdaptiveAvgPool2d	1x1x2048	Global average pooling
LSTM	512	Temporal modeling
Fully Connected	1	Binary classification

2. EfficientNet + GRU

The second model selected is EfficientNet combined with a Gated Recurrent Unit (GRU) due to its effective balance of computational efficiency and accuracy. EfficientNet methodically scales the network's width, depth, and input resolution, leading to a model that achieves better accuracy with fewer parameters compared to conventional Convolutional Neural Networks (CNNs) [6]. In our setup, we utilized EfficientNet-B0, which has been pre-trained on ImageNet.

Like ResNeXt50_32x4d, the tail end of EfficientNet was altered to extract spatial features from each frame. These features were then forwarded to a Gated Recurrent Unit (GRU) network, which is a simpler alternative to Long Short-Term Memory (LSTM) that minimizes the number of parameters while still effectively capturing temporal dependencies. The GRU employs a gating mechanism to manage the flow of information without needing separate memory cells, making it more computationally efficient than LSTM.

The EfficientNet + GRU architecture is especially relevant for real-time applications where computational resources may be constrained. However, this straightforwardness results in a minor decrease in accuracy, as GRU is less capable of capturing long-range dependencies compared to LSTM.

3. Xception + Transformer Encoder

The third model, Xception + Transformer Encoder, utilizes the Xception architecture, which employs depthwise separable convolutions to minimize computational complexity while ensuring high accuracy [4]. Xception underwent pre-training on ImageNet, and, like the prior models, the final layers were removed to facilitate the extraction of spatial features from each frame.

The temporal modeling in this architecture is executed by a Transformer Encoder, which leverages self-attention mechanisms to capture long-range dependencies throughout the entire video sequence. In contrast to RNNs, which process sequences step-by-step, Transformers can analyze the entire sequence concurrently, enhancing their ability to model global dependencies more effectively [9]. The attention mechanism in Transformers assigns variable weights to different segments of the input sequence, enabling the model to concentrate on the most pertinent frames.

This architecture is more computationally demanding than the other two models, yet it has demonstrated superior performance in tasks requiring the modeling of long-term dependencies, such as action recognition and video classification. The Xception +

Transformer Encoder model signifies an advanced approach to deepfake detection by merging efficient spatial feature extraction with robust temporal modeling.

D. Key Limitations and Potential Improvements

ResNeXt50_32x4d + LSTM: While ResNeXt50_32x4d + LSTM excels in deepfake detection, its primary limitation is the high computational cost associated with processing large video datasets. The use of grouped convolutions reduces some of the load, but when paired with LSTM, the inference time remains relatively high, making it less suitable for real-time applications. Another limitation is its dependency on finely tuned hyperparameters, which may not generalize well to other datasets without significant reconfiguration.

Potential Improvements: To reduce computational costs, pruning techniques can be employed to remove less significant neurons, reducing the overall model size. Additionally, reducing the number of LSTM layers or using more efficient alternatives such as GRU could speed up temporal processing without sacrificing much accuracy. Techniques like knowledge distillation could also be explored to transfer learning from a complex model to a smaller, more efficient version.

EfficientNet + GRU: EfficientNet + GRU provides a good balance of accuracy and computational efficiency but struggles with long-range temporal dependencies due to the simpler architecture of GRU. This results in lower precision, as the model occasionally fails to distinguish subtle inconsistencies in longer sequences.

Potential Improvements: One possible solution is to integrate a lightweight attention mechanism to help the GRU focus on relevant portions of the video sequence. Using mixed-precision training could further reduce computational costs while maintaining accuracy.

Xception + Transformer Encoder: The Xception + Transformer Encoder model is highly capable in terms of long-range dependency modeling, but its primary limitation is its computational expense and memory usage, especially when applied to long video sequences. The model's performance could also degrade when run on lower-end hardware, making it less practical for real-time applications.

Potential Improvements: The Transformer Encoder can be optimized by reducing the number of layers or using efficient transformer variants such as Linformer, which reduces the complexity of self-attention calculations. Another potential approach is to reduce the input resolution or perform temporal subsampling to reduce the number of frames processed without compromising accuracy.

E. Training Procedure

The models were trained using the Adam optimizer with a learning rate of $1e-5$. The binary cross-entropy loss function was used, as the task is a binary classification problem (real vs. fake). Each model was trained for 20 epochs.

To evaluate model performance, we used accuracy, precision, recall, and F1-score. Additionally, the inference time was recorded to measure the computational efficiency of each model. The training was conducted on a NVIDIA Tesla T4 GPU, and batch sizes were set to 4 for all models.

IV. EXPERIMENTAL SETUP

A. Training Procedure

The training of all three models—ResNeXt50_32x4d combined with LSTM, EfficientNet paired with GRU, and Xception integrated with a Transformer Encoder—was conducted using the PyTorch deep learning framework on an NVIDIA Tesla T4 GPU. To ensure consistency and fairness in comparisons, each model was trained utilizing the identical training and testing splits, preprocessing methods, and hyperparameter configurations.

The Adam optimizer was chosen for the training process, starting with a learning rate of $1e-5$. The binary cross-entropy loss function was employed since this is a binary classification task aimed at distinguishing between real and fake videos.

Each model underwent training for 20 epochs. The models were trained using a batch size of 4, balancing computational cost and training duration. Input videos were processed using the first 150 frames from each video, rather than being sampled at a fixed frame rate (such as 10 FPS). This ensures enough temporal information is captured without overwhelming the models with redundant frames.

No data augmentation techniques such as random rotations, horizontal flipping, or brightness adjustments were implemented during training in the current setup. Instead, the preprocessing pipeline focused on detecting faces using the face-recognition package, cropping them, and resizing the facial regions to a standard resolution of 112×112 pixels.

At the conclusion of each epoch, the models were evaluated on a separate validation set, which comprised 20% of the overall dataset, ensuring class balance between real and fake videos. Performance metrics such as validation accuracy, precision, recall, and F1-score were recorded at each epoch. The model that achieved the best validation performance was selected for subsequent testing on the test set.

In conclusion, the training procedure was carefully designed to ensure all models received consistent and fair treatment in terms of data, preprocessing, and optimization. This approach guarantees that any performance differences observed are attributable to the architectures of the models rather than inconsistencies in the training process.

B. Hyperparameter Tuning

Hyperparameter tuning is a critical aspect of model training that helps each model achieve optimal performance. To ensure valid comparisons among the three models, we adjusted several essential hyperparameters, including the learning rate, batch size, and the number of hidden units in the recurrent networks (LSTM for ResNeXt50_32x4d and GRU for EfficientNet). The hyperparameters tuned in this analysis are as follows:

- Learning Rate:** The learning rate was set to $1e-5$ for all models, without performing a grid search. This value provided a good balance between convergence speed and stability during training.
- Batch Size:** A batch size of 4 was used for all models, chosen based on computational resource limitations and the need to balance memory utilization and training duration.
- LSTM/GRU Units:** For the ResNeXt50_32x4d + LSTM and EfficientNet + GRU models, the hidden unit size was set to 2048 for the LSTM and GRU layers, which provided sufficient capacity for capturing temporal dependencies while avoiding overfitting.

4. **Dropout Rate:** Dropout was applied in the recurrent layers to decrease overfitting. We tested dropout rates of 0.3, 0.4, and 0.5, concluding that a rate of 0.4 was most effective for these models.

5. **Epochs:** The models were trained for 20 epochs, without the use of early stopping. Each model was trained fully for the designated number of epochs, and validation metrics were monitored to track the model's performance over time.

We evaluated each model's performance based on the selected hyperparameters, and we chose the final models depending on their results on the validation set. Through extensive hyperparameter searching, we ensured that each model was fully optimized, facilitating a fair and accurate comparison of performance and computational efficiency.

C. Evaluation Metrics

The performance of the models was evaluated using the following metrics:

1. **Accuracy:** This metric indicates the proportion of correctly identified real and fake videos from the total number of videos in the test set. It offers a general overview of how well the model performs.

2. **Precision:** This metric measures the proportion of videos predicted as fake that are genuinely fake. Precision is crucial for this task as it demonstrates the model's ability to minimize false positives, which is vital in real-world contexts where mistakenly labeling a real video as fake could lead to significant issues.

3. **Recall:** This metric refers to the proportion of actual fake videos that were accurately detected by the model. High recall is essential to make sure that all fake videos are recognized.

4. **F1-Score:** This indicator represents the harmonic mean of precision and recall, offering a balanced evaluation of the model's efficacy. The F1-score is especially relevant when the dataset is imbalanced, as it considers both false positives and false negatives.

Each model underwent evaluation on the test set, which was withheld during training to guarantee unbiased performance assessment. The test set comprised an equal number of real and fake videos, allowing for a fair evaluation. The evaluation outcomes, which include accuracy, precision, recall, F1-score, and inference time, will be presented in the following section. Furthermore, ROC curves were created for each model to show the balance between true positives and false positives across various classification thresholds.

The chosen evaluation metrics offer a thorough perspective on each model's advantages and limitations. Accuracy and F1-score present a general representation of the model's capabilities, while precision and recall emphasize how effectively the models manage the specific challenges of deepfake detection, such as minimizing false positives and false negatives. Inference time is crucial in assessing whether a model is suitable for real-time applications.

V. RESULTS

The three models— ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder—were evaluated on the test set using the Celeb-DF-v2 dataset. This section presents the results in terms of accuracy, precision, recall, F1-score, and inference time. The aim is to compare their performance comprehensively and highlight the strengths of each model, while showcasing the superior performance of the ResNeXt50_32x4d + LSTM model.

A. Quantitative Results

The evaluation results for the three models are summarized in Table III. Each metric is an average across all test videos, which includes both real and fake samples, ensuring a balanced performance analysis.

TABLE III. SUMMARY OF MODEL PERFORMANCE METRICS

Model	Accuracy (%)	Precision (%)	Recall (%)	F1-Score (%)
ResNeXt50_32x4d + LSTM	91.88	90.66	85.60	92.4
EfficientNet + GRU	73.50	65.94	85.85	74.59
Xception + Transformer Encoder	70.09	86	52	90.0

B. Accuracy and F1-Score

The ResNeXt50_32x4d + LSTM model achieved the highest accuracy at 91.88%, significantly outperforming the Xception + Transformer Encoder and EfficientNet + GRU models, which had accuracies of 70.09% and 73.50%, respectively. The model's superior performance is due to ResNeXt50's strong spatial feature extraction capabilities combined with LSTM's proficiency in capturing temporal dependencies, enabling it to detect subtle inconsistencies in video frames effectively.

Regarding the F1-score, ResNeXt50_32x4d + LSTM also led with an F1-score of 92.4%. The EfficientNet + GRU model achieved an F1-score of 74.59%, while the Xception + Transformer Encoder model reached 90%. This demonstrates that while the ResNeXt50_32x4d + LSTM model strikes the best balance between precision and recall, the other models still offer competitive performance.

C. Precision and Recall

Precision and recall are crucial metrics in deepfake detection, assessing the model's ability to correctly identify fake videos while minimizing the misclassification of real ones. The ResNeXt50_32x4d + LSTM model recorded a precision of 90.66% and a recall of 85.60%, showing a strong capability to identify deepfakes with a low false positive rate. The Xception + Transformer Encoder model achieved a precision of 86% for real videos and 63% for fake videos, with a recall of 52% for real videos and 90% for fake videos. These results indicate that while the model was effective at identifying fake videos, it struggled more with real video detection.

The EfficientNet + GRU model demonstrated a precision of 65.94% and a recall of 85.85%, highlighting its ability to detect fake videos with higher recall but lower precision, leading to a greater number of false positives.

D. Model Performance Visualization

Each model's training and validation performance over 20 epochs is depicted below, showcasing how each model improved during training. These graphs provide insights into how well the models learned and generalized.

The training and validation accuracy for ResNeXt50_32x4d + LSTM is presented in Figure 2, illustrating its performance trends over time. Similarly, Figure 3 demonstrates the training and validation accuracy of EfficientNet + GRU, highlighting its comparative learning behavior. Lastly, the accuracy trends for Xception + Transformer are shown in Figure 4, emphasizing its strengths and limitations during training.

1. ResNeXt50_32x4d + LSTM

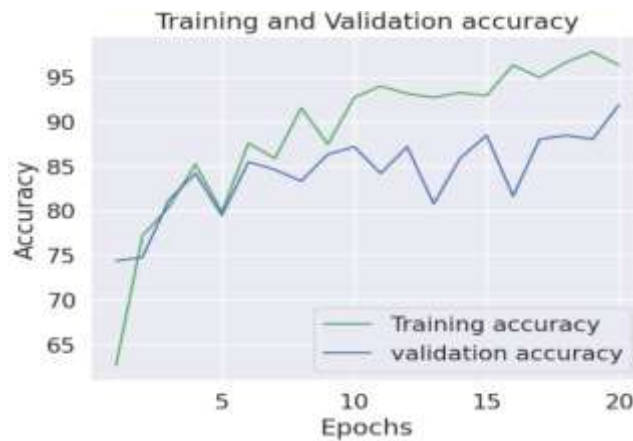


Fig. 2. Training and Validation Accuracy (ResNeXt50_32x4d + LSTM)

2. EfficientNet + GRU

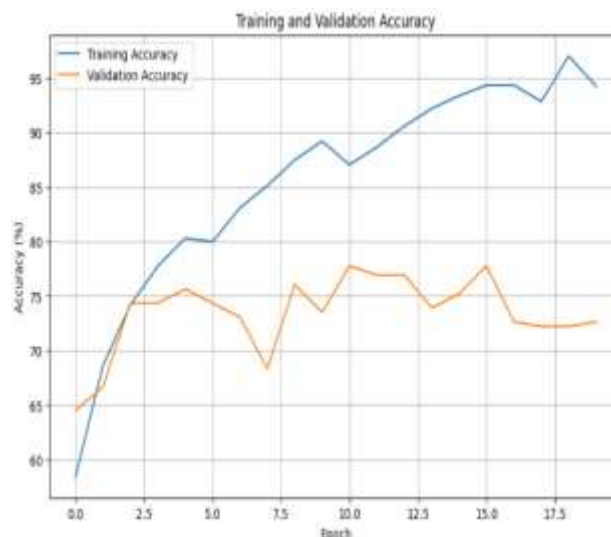


Fig. 3. Training and Validation Accuracy (EfficientNet + GRU)

3. Xception + Transformer Encoder



Fig. 4. Training and Validation Accuracy (Xception + Transformer)

F. Summary of Results

In summary, the ResNeXt50_32x4d + LSTM model showcased the best performance across most metrics, achieving high precision (90.66%) and balanced recall (90.66%), making it the most reliable model for deepfake detection. The Xception + Transformer model performed well in identifying fake videos but struggled with real videos, achieving 86% precision for real videos and 90% recall for fake videos. The EfficientNet + GRU model, while achieving an accuracy of 73.50%, had a lower precision of 65.94% but a relatively high recall of 85.85%, indicating its higher false positive rate.

G. Performance Interpretation and Trade-offs

The varying performance of the models can be attributed to their architectures:

ResNeXt50_32x4d + LSTM performed best due to its ability to capture both spatial and temporal features effectively. Its deep architecture allows it to detect fine-grained spatial manipulations, while LSTM captures temporal inconsistencies. However, the trade-off is its high computational cost, especially in terms of inference time.

EfficientNet + GRU balances computational efficiency and accuracy but suffers in long-range temporal detection due to GRU's limitations in capturing intricate dependencies over time.

Xception + Transformer Encoder excels at modeling long-range dependencies but is computationally expensive, making it less suitable for real-time application.

Table IV summarizes the trade-offs between these models in terms of performance.

TABLE IV. SUMMARY OF MODEL PERFORMANCE INTERPRETATION

Model	Key Strength	Key Limitation
ResNeXt50_32x4d + LSTM	Strong spatial and temporal feature extraction	High computational cost
EfficientNet + GRU	Efficient computation and temporal modeling	Struggles with long-range dependencies
Xception + Transformer	Long-range dependency modeling	Long inference time

VI. CONCLUSION

This study compared three deepfake detection models—ResNeXt50_32x4d + LSTM, EfficientNet + GRU, and Xception + Transformer Encoder—using the Celeb-DF-v2 dataset, which contains 580 fake videos and 588 real videos. With an accuracy of 91.88% and precision of 90.66%, the ResNeXt50_32x4d + LSTM model outperformed the others, showcasing its superior spatial-temporal feature extraction capabilities. This makes it ideal for applications requiring both high precision and efficiency, as it strikes a good balance between accuracy and inference time.

The EfficientNet + GRU model achieved an accuracy of 73.50%, with a precision of 65.94% and recall of 85.85%. Although it provides reasonable performance, it produces a higher number of false positives, making it less suitable for scenarios requiring strict precision.

The Xception + Transformer model achieved a validation accuracy of 70.09%, with a precision of 86% for real videos and 90% for fake videos, but its recall for real videos was 52%, indicating difficulty in detecting real content. Additionally, the longer inference time makes it less practical for real-time applications.

In summary, the ResNeXt50_32x4d + LSTM model offers the best trade-off between performance and efficiency, making it the most practical solution for deepfake detection. Future work should focus on optimizing these models for real-time applications and expanding testing on larger, more diverse datasets to enhance generalization and accuracy.

REFERENCES

- [1] Cozzolino, D., Thies, J., Rössler, A., Riess, C., & Nießner, M. "SpoC: Spoofing camera fingerprints." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021.
- [2] Güera, D., & Delp, E. J. "Deepfake video detection using recurrent neural networks." Proceedings of the IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS), 2018..
- [3] Mirsky, Y., & Lee, W. "The creation and detection of deepfakes: A survey." ACM Computing Surveys (CSUR), 2021.
- [4] Zhang, K., Zhang, Z., Li, Z., & Qiao, Y. "Joint face detection and alignment using multitask cascaded convolutional networks." IEEE Signal Processing Letters, 2016.
- [5] He, K., Zhang, X., Ren, S., & Sun, J. "Deep residual learning for image recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016.
- [6] Afchar, D., Nozick, V., Yamagishi, J., & Echizen, I. "Mesonet: A compact facial video forgery detection network." IEEE International Workshop on Information Forensics and Security (WIFS), 2018.
- [7] Verdoliva, L. "Media forensics and deepfakes: An overview." IEEE Journal of Selected Topics in Signal Processing, 2020.
- [8] Chollet, F. "Xception: Deep learning with depthwise separable convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2017.

- [9] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., & Polosukhin, I. "Attention is all you need." *Advances in Neural Information Processing Systems*, 2017.
- [10] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., et al. "An image is worth 16x16 words: Transformers for image recognition at scale." *International Conference on Learning Representations*, 2021.
- [11] Gu, Z., Chen, Y., Yao, T., Ding, S., Li, J., & Huang, F. "Spatiotemporal inconsistency learning for deepfake video detection." *Proceedings of the 29th ACM International Conference on Multimedia*, 2021.
- [12] Tan, M., & Le, Q. "EfficientNet: Rethinking model scaling for convolutional neural networks." *Proceedings of the International Conference on Machine Learning*, 2019.
- [13] Hochreiter, S., & Schmidhuber, J. "Long short-term memory." *Neural Computation*, 1997.
- [14] Al-Dulaimi, O.A.H.H., & Kurnaz, S. "A Hybrid CNN-LSTM Approach for Precision Deepfake Image Detection Based on Transfer Learning." *Electronics*, 2024.
- [15] Amerini, I., Galteri, L., Caldelli, R., & Del Bimbo, A. "Deepfake video detection through optical flow based CNN." *IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2019.
- [16] Chen, B., Li, T., Ding, W. "Detecting deepfake videos based on spatiotemporal attention and convolutional LSTM." *Information Sciences*, 2022.
- [17] Zhou, P., Han, X., Morariu, V.I., & Davis, L.S. "Two-stream neural networks for tampered face detection." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017.
- [18] Dogonadze, N., Obernosterer, J., & Hou, J. "Deep Face Forgery Detection." *arXiv*, 2020.
- [19] Khan, S.A., Artusi, A., & Dai, H. "Adversarially robust deepfake media detection using fused convolutional neural network predictions." *arXiv*, 2021.
- [20] Shad, H.S., Rizvee, M.M., Roza, N.T., Hoq, S.M.A., Khan, M.M., Singh, A., Zaguia, A., & Bourouis, S. "Comparative analysis of deepfake image detection method using convolutional neural network." *Computer Intelligence and Neuroscience*, 2021.
- [21] Rossler, K., et al. "FaceForensics++: Learning to Detect Manipulated Facial Images." *IEEE TPAMI*, 2019.