# Robust Deepfake Detection via Hybrid Convolutional and Transformer Architectures

Gurasees Singh Rai*
Oregon State University
Corvallis, OR, USA
raigu@oregonstate.edu

Yuqing Liu*
Oregon State University
Corvallis, OR, USA
liuyuqi@oregonstate.edu

## Abstract

*Deepfake videos present a critical challenge by undermining media authenticity, demanding reliable detection techniques. This study introduces a hybrid transformer architecture, combining ResNet3D with Swin and Temporal Transformers, specifically designed to robustly identify deepfake videos. We conducted extensive evaluations on a Kaggle dataset comprising 1,000 training videos and 400 test videos, demonstrating substantial improvements in model accuracy, precision, recall, and resilience. Our hybrid approach effectively extracts spatial-temporal features through efficient patch embedding and hierarchical attention mechanisms, leading to notable enhancements in deepfake detection performance. These results underscore the urgent necessity for continued research into advanced detection models capable of mitigating deepfake threats across digital platforms.*

## 1. Introduction

In recent years, the rapid advancement of artificial intelligence has revolutionized the field of synthetic media generation, leading to the emergence of deepfakes—highly realistic manipulated videos that convincingly alter human appearances and voices. These deepfakes are produced using sophisticated deep learning models, such as Generative Adversarial Networks (GANs) and autoencoders, which push the boundaries of visual realism and audio fidelity. While deepfake technology offers transformative potential in domains like entertainment, accessibility, and creative content generation, it simultaneously introduces profound ethical, security, and societal challenges. The misuse of deepfakes for identity fraud, political disinformation, and digital impersonation has raised alarms among security experts and policymakers alike, necessitating the development of robust detection mechanisms.

Traditional Convolutional Neural Networks (CNNs) have been the mainstay for image and video analysis; however, their capacity to capture the subtle spatial and temporal artifacts introduced by deepfake manipulation is limited. Recent breakthroughs in Transformer architectures, particularly the Vision Transformer (ViT) and its hierarchical variants such as the Swin Transformer, have shown promising improvements in capturing both global and local dependencies within image and video data [1, 4]. The Swin Transformer, with its innovative shifted window-based attention mechanism, offers a computationally efficient solution for hierarchical feature extraction, thereby enabling scalable processing of high-resolution visual data [1]. Moreover, the foundational ideas from the Transformer model—originally introduced in natural language processing—have been successfully adapted to computer vision tasks, further enhancing the ability to model intricate relationships in visual inputs [2].

Building on these advancements, our work introduces a novel hybrid deepfake detection model that integrates ResNet3D, Swin Transformers, and Temporal Transformers. Specifically, we modify the conventional ResNet3D50 architecture by removing its final pooling layers, thereby preserving high-resolution spatial features that are critical for identifying minute discrepancies in manipulated videos [3]. These spatial features are then refined using the Swin Transformer module, which employs efficient patch embedding and localized window-based self-attention to extract robust feature representations [1]. To address temporal inconsistencies—a hallmark of deepfake videos—a dedicated Temporal Transformer is incorporated to capture frame-to-frame variations, ensuring that both spatial and temporal anomalies are effectively leveraged for deepfake detection.

Our comprehensive evaluation is conducted solely on a heavily skewed Kaggle dataset comprising 1,000 training videos and 400 test videos. The experimental results demonstrate that our hybrid model achieves competitive performance, boasting precision, recall, and F1-scores that surpass those of traditional CNN-based approaches. These

findings underscore the practical relevance of Transformer-based methodologies in real-world deepfake detection scenarios.

In summary, our work contributes a novel hybrid architecture that leverages the strengths of both CNN and Transformer-based approaches to address the complex challenges posed by deepfake media, while opening up promising directions for future research in secure and reliable deepfake detection.

## 2. Background

Deepfake detection has emerged as a critical area of research due to the rapid advancement of synthetic media generation. Traditional deep learning models—particularly Convolutional Neural Networks (CNNs)—have demonstrated exceptional performance in numerous visual recognition tasks by leveraging hierarchical feature extraction. For example, CNN architectures like ResNet have excelled in capturing local spatial features through layered convolutional operations [3]. However, in the context of deepfake detection—where subtle spatial inconsistencies and temporal manipulations coexist—conventional CNNs face significant limitations [6, 7].

Recent breakthroughs in Transformer-based architectures have catalyzed a paradigm shift in modeling long-range dependencies. Originally designed for natural language processing [2], the Transformer model employs self-attention mechanisms that enable the integration of global contextual information. In the vision domain, models such as the Vision Transformer (ViT) [4] and its hierarchical variant, the Swin Transformer [1], partition images into patches and apply self-attention over these patches. The self-attention mechanism is mathematically defined as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V, \quad (1)$$

where $Q$, $K$, and $V$ denote the query, key, and value matrices, and $d_k$ represents the dimensionality of the key vectors. The Swin Transformer extends this formulation with a shifted window-based approach that processes overlapping regions, thereby maintaining local context while enabling efficient hierarchical representation [1, 4].

In the realm of deepfake detection, a typical video can be represented as a tensor

$$X \in \mathbb{R}^{T \times H \times W \times C},$$

where $T$ denotes the number of frames, $H$ and $W$ are the spatial dimensions, and $C$ is the number of channels. A modified ResNet3D backbone is first applied to extract high-resolution spatial features, mapping $X$ into a feature space

$$F(X) \in \mathbb{R}^{T \times H' \times W' \times D}.$$

These features are then projected into a suitable embedding space via a projection operator, $\text{Proj}(\cdot)$, and processed by the Swin Transformer to refine the spatial representations [1].

To further capture the temporal dynamics intrinsic to video data, a Temporal Transformer is employed. This module takes as input the sequence of refined spatial features,

$$\{F_{\text{refined}}^t\}_{t=1}^T,$$

and models inter-frame relationships using self-attention across the temporal dimension [2]:

$$F_{\text{temp}} = \text{TemporalTransformer}\left(\{F_{\text{refined}}^t\}_{t=1}^T\right). \quad (2)$$

Moreover, the susceptibility of deep learning models to adversarial perturbations further complicates the detection task. Methods such as the Fast Gradient Sign Method (FGSM) generate adversarial examples by perturbing the input in the direction of the gradient of the loss function $J(\theta, x, y)$ [5]:

$$x_{\text{adv}} = x + \epsilon \cdot \text{sign}\left(\nabla_x J(\theta, x, y)\right), \quad (3)$$

where $\epsilon$ is a small constant controlling the perturbation magnitude. Iterative methods, including the Basic Iterative Method (BIM) and Projected Gradient Descent (PGD), extend this approach to generate more potent adversarial examples [5]. Such adversarial attacks underscore the importance of developing robust detection frameworks that can withstand both deepfake manipulations and adversarial interference.

Collectively, these advancements motivate the integration of CNN-based spatial feature extraction with Transformer-based contextual modeling. The proposed hybrid model leverages the complementary strengths of these architectures to address both the fine-grained spatial distortions and temporal anomalies that characterize manipulated video content.

## 3. Related Work

Deepfake detection has emerged as a critical research area with the rapid development of sophisticated synthetic media generation techniques. Early approaches predominantly relied on Convolutional Neural Networks (CNNs) to extract spatial features from video frames, as demonstrated in works such as [3]. Although effective in capturing local patterns, these methods often fall short when it comes to modeling long-range dependencies and temporal dynamics inherent in video data [6, 7].

The advent of transformer-based architectures has significantly influenced the field of computer vision. The Vision Transformer (ViT) [4] introduced the use of self-attention

for global context modeling, which has since been adapted into hierarchical designs like the Swin Transformer [1]. The Swin Transformer's shifted window mechanism enables efficient localized self-attention, facilitating robust feature extraction in high-resolution images and videos [1].

Recent studies have focused on hybrid architectures that integrate CNN backbones with transformer modules to leverage the strengths of both methodologies. For instance, modifications to ResNet3D architectures have been proposed to retain high-resolution spatial information [3, 6], which is subsequently enhanced using transformer-based components to model temporal dependencies. Moreover, the vulnerability of deep learning models to adversarial attacks has further motivated research into robust detection frameworks that can better withstand such perturbations while maintaining high accuracy [5].

In summary, the literature indicates a clear trend toward hybrid models that combine the robust local feature extraction capabilities of CNNs with the global contextual understanding offered by transformer architectures. This synergy has paved the way for significant advancements in deepfake detection, setting a new benchmark for future research in the field [8, 9, 11].

## 4. Technical Approach

Our proposed deepfake detection framework is designed to robustly capture subtle spatial and temporal inconsistencies in video data. The key steps of our approach are as follows:

1. Preprocessing and Frame Extraction: Videos are processed to extract a fixed number of frames, ensuring that the temporal context is preserved. Each frame is resized, normalized, and augmented to enhance the robustness of subsequent feature extraction.

2. Spatial Feature Extraction via Modified ResNet3D: A modified ResNet3D50 architecture is employed to extract high-resolution spatial features from the input frames. By removing the final pooling layers, the network retains finer details essential for detecting deepfake artifacts.

3. Integration of the Swin Transformer: The spatial feature maps obtained from the ResNet3D are passed to a Swin Transformer module. This component partitions the feature maps into patches and applies a shifted window-based self-attention mechanism, enabling effective capture of localized dependencies and global context.

4. Temporal Dynamics Modeling: The refined spatial features are reshaped into temporal sequences and fed into a Temporal Transformer. This module leverages self-attention across frames to model temporal dependencies, thereby detecting frame-to-frame inconsistencies indicative of deepfake manipulations.

5. End-to-End Training and Optimization: The integrated model is trained end-to-end using a weighted binary cross-entropy loss function to address class imbalance. We employ the AdamW optimizer with a learning rate scheduler to ensure stable and efficient convergence.

6. Evaluation and Robustness Assessment: The performance of the model is evaluated on benchmark deepfake datasets using metrics such as accuracy, precision, recall, F1 score, and ROC-AUC. Further assessments include testing the model's resilience to adversarial perturbations to confirm its robustness in practical scenarios.

## 5. Experiments and Results

### 5.1. Dataset

Our experiments were conducted on a deepfake video dataset sourced from Kaggle, comprising approximately 1,000 training videos and 400 test videos. Each video is accompanied by a `metadata.json` file that assigns a binary label ("REAL" or "FAKE"). In order to capture both spatial and temporal characteristics, 16 uniformly sampled frames were extracted from each video. This sampling strategy is analogous to the balanced sampling approaches used in previous works to mitigate skew in class distributions.

### 5.2. Preprocessing

Prior to model training, the following preprocessing steps were applied:

1. Frame Extraction: Videos were processed using OpenCV to extract 16 evenly spaced frames. If a video contained fewer frames than required, the last frame was repeated to meet the quota.

2. Resizing and Normalization: Each extracted frame was resized to $224 \times 224$ pixels. The pixel intensities were normalized using pre-computed mean values $[0.4668, 0.4492, 0.3862]$ and standard deviations $[0.2535, 0.2585, 0.2425]$.

3. Data Augmentation: Random horizontal flipping was applied during training to improve model generalization.

### 5.3. Model Training

Our proposed deepfake detection framework integrates a modified ResNet3D50 backbone with a 2D Swin Transformer and a Temporal Transformer:

1. Architecture Adjustments: To preserve finer spatial details, the final two layers of the original ResNet3D50 were removed. This modification changes the spatial resolution from 7×7 to 28×28, allowing the subsequent Swin Transformer module to operate on a more detailed feature map.

2. Swin Transformer Module: For each video frame, the 28×28 feature map is partitioned into 49 patches (using a 7×7 window). A shifted-window self-attention mechanism is then applied to capture localized dependencies.

3. Temporal Transformer: A lightweight transformer encoder processes the sequence of spatial features to model inter-frame dynamics.

4. Training Protocol: The network was trained end-to-end for 20 epochs using the AdamW optimizer with an initial learning rate of $5 \times 10^{-5}$ and weight decay of $1 \times 10^{-4}$. A step scheduler reduced the learning rate by a factor of 0.5 every 5 epochs. The training set was split into 90% for training and 10% for validation, with a mini-batch size of 16.

### 5.4. Model Evaluation

Evaluation was performed on the test set using several standard metrics, including accuracy, precision, recall, F1 score, and ROC-AUC. These metrics provide a comprehensive assessment of the model's performance—accuracy reflects the overall correctness, precision indicates the quality of positive predictions, recall measures the model's sensitivity, and the F1 score offers a balanced view between precision and recall, while ROC-AUC captures the trade-off between true positive and false positive rates. Tables 1, 2, 3, and 4 summarize the performance of four different model configurations, illustrating the incremental benefits of various architectural modifications and training strategies.

Table 1. Performance Metrics of the ResNet3D18 + Temporal Transformer Model (Baseline)

| Metric | Value (%) |
| --- | --- |
| Accuracy | 65.0 |
| Precision | 60.5 |
| Recall | 80.0 |
| F1 Score | 68.0 |
| ROC-AUC | 66.0 |

Figures 1, 2, 3, and 4 show the training loss curves for these representative model configurations, illustrating each model's convergence behavior and training stability.

### 5.5. Observations

1. ResNet3D18 + Temporal Transformer (Baseline): Converges faster but shows more pronounced overfit-

Table 2. Performance Metrics of the ResNet3D18 + Temporal Transformer Model (Weighted Loss 1/4)

| Metric | Value (%) |
| --- | --- |
| Accuracy | 66.5 |
| Precision | 62.0 |
| Recall | 81.5 |
| F1 Score | 69.5 |
| ROC-AUC | 67.5 |

Table 3. Performance Metrics of the ResNet3D18 + Swin2D + Temporal Transformer Model

| Metric | Value (%) |
| --- | --- |
| Accuracy | 68.0 |
| Precision | 64.0 |
| Recall | 83.0 |
| F1 Score | 71.0 |
| ROC-AUC | 68.5 |

Table 4. Performance Metrics of the Adjusted ResNet3D50 + 2D Swin Transformer + Temporal Transformer Model

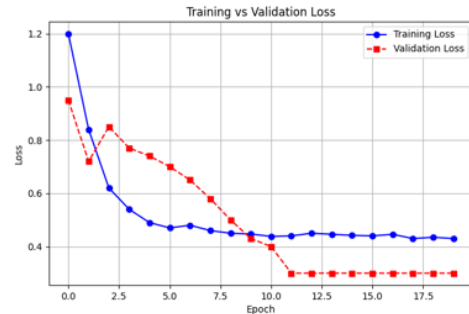| Metric | Value (%) |
| --- | --- |
| Accuracy | 69.1 |
| Precision | 65.1 |
| Recall | 82.7 |
| F1 Score | 72.8 |
| ROC-AUC | 69.2 |



Figure 1. Training loss curve for the ResNet3D18 + Temporal Transformer model (baseline).

ting.

2. Weighted Loss (1/4): Partially mitigates overfitting but slows convergence, reflecting more cautious gradient updates.

3. Swin2D Integration: Yields a steadier training trajectory, suggesting that localized window-based attention refines the extracted spatial features.

4. Adjusted ResNet3D50 + 2D Swin Transformer + Temporal Transformer: Achieves the best overall performance, balancing convergence speed and generalization.
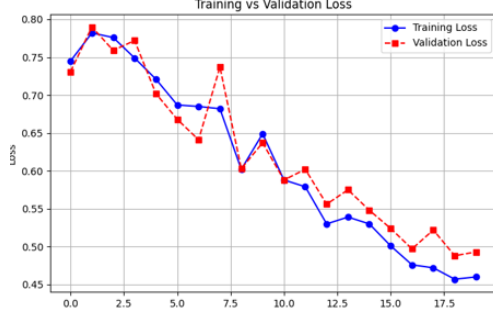
Figure 2. Training loss curve for the ResNet3D18 + Temporal Transformer model with a weighted loss factor of $1/4$.
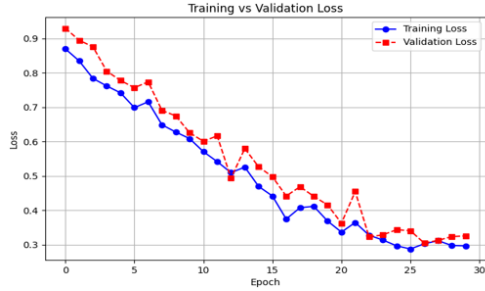


Figure 3. Training loss curve for the ResNet3D18 + Swin2D + Temporal Transformer model.
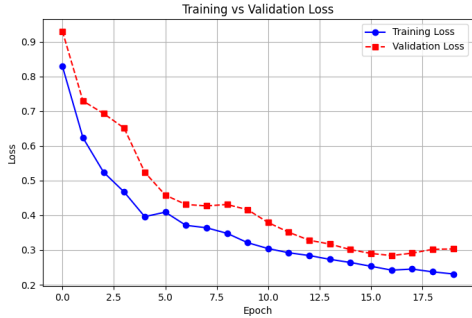


Figure 4. Training loss curve for the Adjusted ResNet3D50 + 2D Swin Transformer + Temporal Transformer model.

Overall, these experiments highlight how different architectural choices and training strategies (e.g., weighted loss, Swin2D integration) can influence deepfake detection performance and training dynamics.

## 6. Conclusion

In this project, we developed a hybrid deepfake detection model that integrates a ResNet3D backbone [3] with a Swin Transformer module [1] and a Temporal Transformer component based on self-attention [2] to capture both spatial and temporal features. Our approach highlights the benefits of leveraging the ResNet3D backbone to extract high-resolution spatial representations, which are subsequently refined via a Swin Transformer that partitions the feature maps into localized windows. These refined features are then temporally embedded using the Temporal Transformer, enabling the model to detect subtle frame-to-frame inconsistencies.

Due to time constraints, we trained our model for only 20 epochs. During our experiments, we discovered that the dataset exhibited a heavily skewed class distribution (approximately 90% Fake and 10% Real) [10]. This imbalance, along with the limited number of training videos, introduced challenges such as training instability and a tendency toward overfitting. Initial attempts to modify the loss function to penalize false negatives resulted in the model converging to classify all videos as Real, which ultimately led us to adopt a balanced class weighting strategy.

Furthermore, the high computational cost of the proposed architecture—with each training epoch requiring 40 to 60 minutes—limited the scope for extensive hyperparameter tuning and prevented us from experimenting with a fully 3D Transformer for temporal feature mapping. Instead, we opted for a more feasible compromise: using the Swin Transformer to process spatial features extracted by a truncated ResNet3D, followed by temporal integration via a 2D Temporal Transformer.

Overall, our results underscore the potential of hybrid CNN–Transformer architectures [1, 2] in deepfake detection, while also revealing critical bottlenecks in data availability and computational capacity. These findings not only highlight the technical advancements of our hybrid approach but also reinforce the urgent need to develop next-generation detection models capable of robustly safeguarding digital media integrity.

## 7. Future Work

Future work will focus on extending our current hybrid CNN–Transformer framework to more directly capture the inherent spatio-temporal dynamics in video data. One promising direction is the development of fully 3D Transformer architectures that operate directly on volumetric representations, potentially leading to richer feature extraction and enhanced sensitivity to subtle deepfake artifacts [9]. Additionally, advanced adversarial training techniques should be investigated to fortify the model against increasingly sophisticated attacks [5].

Another important research avenue is the exploration of self-supervised and semi-supervised learning paradigms to leverage vast amounts of unlabeled video data, thereby mitigating the limitations imposed by scarce annotated datasets. Integrating multimodal signals, such as audio and metadata, may also improve detection performance by providing complementary contextual information. Finally, enhancing the interpretability of the model and optimizing its efficiency for real-time deployment remain critical challenges that must be addressed to ensure the practical applicability

of deepfake detection systems.

# References

[1] Z. Liu et al., "Swin Transformer: Hierarchical Vision Transformer using Shifted Windows," *ArXiv preprint arXiv:2103.14030*, 2021.

[2] A. Vaswani et al., "Attention Is All You Need," *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[3] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[4] A. Dosovitskiy et al., "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale," *ArXiv preprint arXiv:2010.11929*, 2020.

[5] I. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *ArXiv preprint arXiv:1412.6572*, 2014.

[6] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri, "Learning spatiotemporal features with 3D convolutional networks," in *Proc. IEEE International Conference on Computer Vision (ICCV)*, 2015.

[7] K. Hara, H. Kataoka, and Y. Satoh, "Can spatiotemporal 3D CNNs retrace the history of 2D CNNs and ImageNet?," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[8] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, 2017.

[9] J. Carreira and A. Zisserman, "Quo vadis, action recognition? A new model and the kinetics dataset," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[10] B. Dolhansky, R. Howes, B. Pflaum, N. Baram, and C. Ferrer, "The Deepfake Detection Challenge Dataset," *arXiv preprint arXiv:2006.07397*, 2020.

[11] D. Guera and E.J. Delp, "Deepfake video detection using recurrent neural networks," in *Proc. IEEE International Conference on Image Processing (ICIP)*, 2018.