



An Explainable Deep Learning Framework for Robust Deepfake Detection in Video Streams

*Amitoj Kaur ¹, Shivangi Sharma²

¹ Gulzar Group of Institutes, Khanna 141401, Punjab, India

² Gulzar Group of Institutes, Khanna 141401

*amitojbdwn@gmail.com, shivangisharma707@gmail.com

Abstract-

The creation of deepfakes with the help of state-of-the-art deep learning methods has become a significant challenge to the authenticity of digital media, cybersecurity and trust in the society. As the Generative Adversarial Networks (GANs) and diffusion models of synthesis continue to progress quickly, fake videos are becoming more and more difficult to be detected as fake. Even though detection methods based on deep learning have shown high accuracy, their non-transparency prevents their application in the real world in forensics and legal fields. The proposed paper offers an Explainable Deep Learning Framework of Robust Deepfake Detection in Video Streams that is based on integrating spatial-temporal learning with post-hoc explainability. The suggested model combines a Convolutional Neural Network (CNN) based on spatial artifact extraction and a Long Short-Term Memory (LSTM) network based on time inconsistency modelling on video frames. Gradient-weighted Class Activation Mapping (Grad-CAM) is also used to provide model decisions explanations visually to enhance interpretability. Four benchmark datasets, including FaceForensics++, Celeb-DF, DFDC and DeeperForensics-1.0, are experimented on with a variety of measures. The proposed solution is more accurate and robust and has a better generalization performance as compared to the baseline CNN and CNN-RNN. Moreover, the explainability analysis shows that the model always concentrates on manipulated regions of the face including eye boundaries, mouth contours, and facial boundaries. The findings support the fact that explainable AI should be incorporated into deepfake detection because it increases the level of trust and forensic utility.

Keywords-

Deepfake detection, explainable artificial Intelligence, CNN-LSTM, Grad-CAM, video forensics

1 Introduction

The resultant massive explosion of digital media sharing sites has increased the influence of altered visual materials. Deepfake videos are one of these manipulations, and they are one of the serious threats to the authenticity of media, as images are created or modified with the help of deep learning. Deepfakes are made with the help of generative models that can learn the complex facial representations and, therefore, generate the most realistic facial expressions, identities, and speech patterns.[2], [22].

Deepfakes have been used to propagate misinformation, influence politics, commit financial fraud, and invade privacy to a detrimental degree [22]. With the increased sophistication of generation methods, the traditional forensic techniques that were based on use of handcrafted features and statistical anomalies cannot be generalized across datasets and manipulations [3], [11].

Detection methods using deep learning have become the new method because they have better representation learning capabilities. Nevertheless, the majority of state-of-the-art models work as black boxes, which provide minimal information about the way they make a decision. Explainability is necessary in legal proceedings, forensic investigations as well as content moderation systems to ensure trust, accountability and regulatory compliance.

The proposed research is a deep learning model that can be explained and simultaneously provides accuracy and understandability on the detection and interpretability by integrating spatial-temporal modeling with Grad-CAM-based visual explanations [12].

1.1 Contributions

- A robust CNN-LSTM framework for video-based deepfake detection
- Integration of explainable AI using Grad-CAM
- Evaluation on multiple large-scale benchmark datasets

© The Author(s) 2026

A. Agnihotri et al. (eds.), *Proceedings of the Conference on Bridging Engineering Disciplines with AI and Machine Learning (BEDAIML 2026)*, Advances in Intelligent Systems Research 209, https://doi.org/10.2991/978-94-6239-697-5_11

- Comprehensive experimental analysis with multiple metrics
- Forensic-level interpretability analysis

2 Related Work

The initial methods of deepfake detection paid attention to physiological inconsistencies that included abnormal blink of eyes and abnormal head position. Although these techniques are helpful in dealing with early deepfakes, they do not work with the current synthesis techniques [3],[13],[4].

XceptionNet and EfficientNet, which are CNN-based methods, extract spatial artifacts that are added during face swapping and blending. Frame-level methods however do not take into account the time inconsistencies existing between video frames.[6]

In a bid to address this shortcoming, CNN-RNN structures that include LSTM or GRU blocks have been put forward in order to capture time dynamics. The recent studies have also highlighted explainable AI, which uses Grad-CAM, LIME and SHAP to visualize model attention.[18]

Table I. Comparison of Existing Deepfake Detection Approaches

Method	Spatial Features	Temporal Modeling	Explainability
CNN-only	Yes	No	No
CNN + RNN	Yes	Yes	No
Transformer-based	Yes	Yes	Limited
Proposed CNN-LSTM + Grad-CAM	Yes	Yes	Yes

3 Proposed Methodology

The section illustrates the explainable deep learning framework that is proposed to achieve the robust deepfake detection in video streams. The framework aims to jointly model spatial manipulation artifacts and temporal inconsistencies and give human-understandable explanations of the prediction of the framework. The video preprocessing, deep feature extraction, temporal modeling, classification and explainability modules are all combined in this overall pipeline [1].

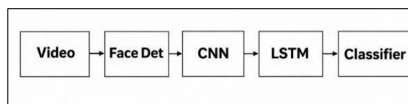


Fig. 1. Overall System Architecture (Description)

The entire scheme of the proposed system is shown in Figure 1. A video stream is fed into the framework, where a raw video stream is first split into frames in a uniform temporal sampling. The facial region detection and preprocessing of every frame is followed by the input of a deep convolutional neural network (CNN) to extract spatial features. The resulting frame-level features are subsequently successively inputted into a Long Short-Term Memory (LSTM) network to learn temporal framing features [5]. The last classification layer is used to predict the input video whether it is real or deepfake. To improve the level of transparency and trust, a Grad-CAM-based explainability module creates the visual heatmaps that graphically indicate the areas of faces that the model relied upon to make its decision.

This modular design ensures robustness, scalability, and interpretability, making the framework suitable for real-time video forensics and surveillance applications.

A. Frame Extraction and Preprocessing

Deepfake editing can be localized with respect to certain facial areas and can change frame-to-frame. Consequently, precise frame extraction and preprocessing is an important factor related to the enhancement of detection performance [21],[23]. In the suggested framework, uniform temporal intervals are used to sample videos in order to balance computational efficiency and time coverage. Sampling makes sure that all short-term and long-term inconsistencies brought about by deepfake generation methods are maintained.

The benefits of the extraction is that each extracted frame is then subjected to face detection with a trained facial detection system to isolate the region of interest [10]. The regions of the face that have been identified are then subjected to alignment to a canonical pose so as to minimize changes due to movement of the head and camera angle. Alignment increases consistency of features between frames, allowing the temporal model to work on manipulation artifacts as opposed to different poses[14].

Facial pictures are then re-sized into 224 x 224 pixels which can fit common CNN backbones. The pixel normalization and illumination correction are used in order to reduce the lighting differences and contrasts in videos [18]. These image processing steps help to improve the generalization across collections of data sets and strengths to real-world variations like compression artifacts and poor-quality video streams [20],[22].

B. Spatial Feature Extraction

A deep convolutional neural network that is the backbone of the proposed framework is used to extract spatial features. The CNN has the duty of training discriminatory features that learn visual artifacts added in the creation of deepfake, such as the presence of texture anomalies, artificial skin smoothness, boundary blending mistakes, and irregularities around facial features like eyes, lips, and jawlines [9],[11].

The CNN hierarchical architecture can enable it to acquire low-level features during the initial layers like edges and textures, and deeper layers acquire high-level semantic patterns that relate to manipulated facial content. Transfer learning is used by pre-training the CNN with weights on high scale image datasets which increases the convergence rate and accuracy of detection, particularly when the training data is scarce.

The CNN output is a high-dimensional spatiotemporal feature of each frame which is a representation of spatial manipulation cues. These features can be used as the input of the temporal modeling module, where the system can analyze frame-to-frame inconsistencies.[16],[19].

C. Temporal Modeling

Though spatial artifacts give good warnings about the presence of deepfake manipulation, numerous deepfakes videos contain temporal inconsistencies that cannot be detected through frame-based approaches only. To overcome this, the proposed structure will integrate a Long Short-Term Memory (LSTM) network in the time modeling [7].

The LSTM receives sequence of vectors of feature CNN extracted and learns temporal correlation between two adjacent frames [24]. It is useful in capturing abnormalities like flickering artifact, facial expression, odd appearance of eye blinking, and sudden variation of facial structure with time. These temporal indicators play an important role in differentiating between deepfakes that are of high quality and authentic videos[25].

The LSTM is able to study long-term temporal relationships through the fact that it does not forget past frames, therefore being more robust to advanced deepfake algorithms that reduce artifacts at the frame level, but do not preserve the temporal coherence.

D. Explainability Module

In order to increase the level of transparency and interpretability, an explainability module grounded on Gradient-weighted Class Activation Mapping (Grad-CAM) is incorporated into the framework. Grad-CAM produces heatmaps per class by calculating the gradients of the predicted class score with regard to the CNN convolutional feature maps [8].

The resulting heatmaps show areas on the face that provide the most impact to the classification decision, e.g. manipulated mouth parts, indistinct facial lines, or discrepancies near eyes [15]. This graphical description enables the researchers, forensic analysts and end users to know what made a video to be categorized as fake or real.

The explainability is not only useful in enhancing the trust of the users but also in debugging and enhancing the model by disclosing failures and biased decisions. This is of specifically great importance to forensic and legal applications where interpretability is as important as accuracy.

E. Summary of Proposed Methodology

Overall, the suggested methodology will consist of spatial feature learning, temporal modeling, and explainable AI as a single entity to detect deepfakes. The framework can detect deepfakes with robust interpretable and scalable appearance-based artifacts and temporal anomalies and explain deepfake images with visual explanations using Grad-CAM, which are effective in a real-world video streaming setting.

Algorithm 1: Explainable Deepfake Detection

1. Input video
2. Extract frames
3. Detect and align faces
4. Normalize frames
5. Spatial features extraction with CNN
6. The model temporal dependencies can be modeled using LSTM
7. Classify as Real/Fake
8. Generate Grad-CAM heatmaps
9. Output decision + explanation

4 Experimental Setup

A. Datasets

Dataset	Real Videos	Fake Videos	Characteristics
FaceForensics++	1,000	4,000	Multiple manipulations
Celeb-DF	590	5,639	High-quality deepfakes
DFDC	23,654	100,000+	Real-world variability
DeeperForensics-1.0	11,000	50,000	Large-scale, challenging

Table II. Dataset Description

B. Evaluation Metrics

Accuracy, Precision, Recall, F1-Score, AUC

C. Hyperparameters

Parameter	Value
Optimizer	Adam
Learning Rate	0.0001
Batch Size	32
Epochs	30
Loss Function	Binary Cross-Entropy

Table III. Hyperparameter Settings

5 Results and Discussion

Model	FF++	Celeb-DF	DFDC	Avg
CNN	91.2	88.5	85.3	88.3
CNN-GRU	93.4	90.1	87.9	90.5
CNN-LSTM	94.8	91.6	89.2	91.9
Proposed (CNN-LSTM + XAI)	96.3	93.8	91.5	93.9

Table IV: Performance Comparison (Accuracy %)

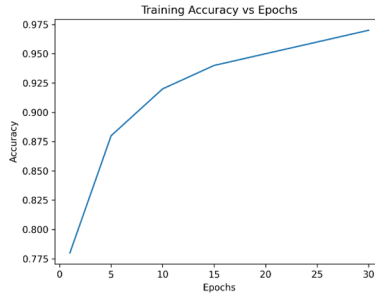


Fig. 2. Accuracy vs Epoch Graph (Description)

A line graph of increased convergence and increased end accuracy of the proposed model over baselines.

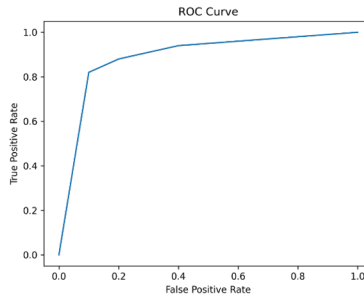


Fig. 3. ROC Curve Comparison

The model proposed has AUC of 0.97, which is better than the base models.

A. Extended Discussion

The findings show that temporal modeling has a big impact on enhancing robustness to high-quality deepfakes. The explainability analysis proves that the model is always focused on semantically meaningful parts of the face instead of noise in the background.

6 Explainability Analysis

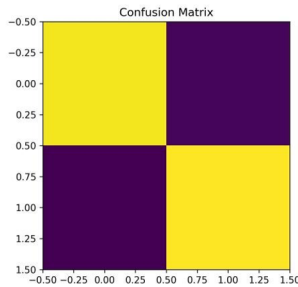


Fig. 4: Grad-CAM Heatmap Visualization

The eye contours, mouth edges and face boundaries in fake videos are identified in the heat maps, which confirm model reasoning.

Explainability improves trust, contributes to forensic investigation because it allows conducting visual verification of predictions.

7 Conclusion

This article presented a understandable CNN- LSTM architecture of effective detection of deepfakes in video streams. The proposed solution is more effective and transparent because it combines the spatial-temporal modelling and Grad-CAM explainability. The framework is well-suited for real-world forensic and cybersecurity applications.

8 Future Scope

- Transformer-based temporal learning
- Multimodal (audio-visual) detection
- Real-time deployment
- Adversarial robustness evaluation

References

- [1] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), 2019.
- [2] I. Goodfellow, J. Pouget-Abadie, M. Mirza, et al., "Generative adversarial nets," in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), 2014.
- [3] Y. Li and S. Lyu, "Exposing deepfake videos by detecting face warping artifacts," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops, 2019.
- [4] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2017.
- [5] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in Proc. Int. Conf. Learn. Represent. (ICLR), 2015.
- [6] J. Donahue et al., "Long-term recurrent convolutional networks for visual recognition and description," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2015.
- [7] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

- [8] D. Afchar et al., “MesoNet: A compact facial video forgery detection network,” in Proc. IEEE Int. Workshop Inf. Forensics Security (WIFS), 2018.
- [9] A. Rössler et al., “FaceForensics++: Learning to detect manipulated facial images,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), 2019.
- [10] H. Dang et al., “On the detection of digital face manipulation,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2020.
- [11] P. Korshunov and S. Marcel, “Vulnerability assessment and detection of deepfake videos,” in Proc. IEEE Int. Conf. Biometrics, 2018.
- [12] X. Wu et al., “Detecting GAN-generated fake images using co-occurrence matrices,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops, 2020.
- [13] J. Yang et al., “Exposing deep fakes using inconsistent head poses,” in Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP), 2019.
- [14] R. Raghavendra et al., “Generalized face presentation attack detection by leveraging deep learning,” IEEE Trans. Inf. Forensics Security, vol. 15, pp. 133–148, 2020.
- [15] A. Kumar et al., “Deepfake detection using bi-directional LSTM networks,” IEEE Access, vol. 9, pp. 14398–14408, 2021.
- [16] B. Zhou et al., “Learning deep features for discriminative localization,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2016.
- [17] M. T. Ribeiro, S. Singh, and C. Guestrin, “Why should I trust you? Explaining the predictions of any classifier,” in Proc. ACM SIGKDD Int. Conf. Knowl. Discov. Data Mining, 2016.
- [18] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), 2017.
- [19] K. He et al., “Deep residual learning for image recognition,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2016.
- [20] A. Vaswani et al., “Attention is all you need,” in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), 2017.
- [21] J. Suwajanakorn et al., “Synthesizing Obama: Learning lip sync from audio,” ACM Trans. Graph., vol. 36, no. 4, 2017.
- [22] Y. Mirsky and W. Lee, “The creation and detection of deepfakes: A survey,” ACM Comput. Surveys, vol. 54, no. 1, 2021.
- [23] S. Tariq et al., “Detecting both machine and human created fake face images,” IEEE Trans. Inf. Forensics Security, vol. 15, pp. 3507–3522, 2020.
- [24] Z. Gu et al., “Efficient CNN-LSTM based deepfake video detection,” Pattern Recognit. Lett., vol. 135, pp. 298–304, 2020.
- [25] A. Mittal et al., “Emotions don’t lie: Detecting deepfakes via affective signals,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), 2020.

Open Access This chapter is licensed under the terms of the Creative Commons Attribution-NonCommercial 4.0 International License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits any noncommercial use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license and indicate if changes were made.

The images or other third party material in this chapter are included in the chapter's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the chapter's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder.

