



## Deep Learning – Augmented Block Analysis for Detecting Advanced Video Forgeries

Sumaiya Shaikh<sup>ID</sup>, Sathish Kumar Kannaiah<sup>\*ID</sup>

Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vijayawada 522302, India

Corresponding Author Email: [ksathish1980@gmail.com](mailto:ksathish1980@gmail.com)

Copyright: ©2025 The authors. This article is published by IETA and is licensed under the CC BY 4.0 license (<http://creativecommons.org/licenses/by/4.0/>).

<https://doi.org/10.18280/ijssse.150108>

### ABSTRACT

**Received:** 26 November 2024

**Revised:** 25 December 2024

**Accepted:** 20 January 2025

**Available online:** 31 January 2025

#### **Keywords:**

*video forgery detection, deep learning, MobileNetV2, digital forensics, ROC – AUC, computational efficiency*

In the digital media security domain, video forgeries, namely cloning, inpainting and splicing, are particularly challenging. In this paper, we present a novel Deep Learning-Augmented Block Analysis (DLBA) framework, which employs the lightweight MobileNetV2 architecture for efficient and accurate detection of advanced video manipulations. The proposed method analyzes videos at the block level for precise localization of tampered regions with computational efficiency. The DLBA framework is shown to be superior in extensive experiments that demonstrate 85% accuracy, an average ROC-AUC of 0.85, and outperforms state of the art methods such as GoogLeNet and ResNet-50. The combination of robust performance and suitability for real time applications suggests that the framework has the potential to be a reliable forensic tool for digital content authentication. Future work will add to the adaptability and scalability of the proposed approach to different datasets and application scenarios.

## 1. INTRODUCTION

In this rapidly advancing world of video editing tools and generative technology, they have created video forgeries that are so convincing, that they pose a big challenge to digital forensics and media authentication [1]. Deepfake videos, small frame alterations, there are implications ranging from misinformation campaign to privacy violation to the conflict in the legal sphere. Since these manipulations become increasingly sophisticated, an increased need arises for development of the advanced strategies for detection of these manipulations with high scalability, robustness, and accuracy [2]. Recent work has also stressed the use of DL (including convolutional neural networks (CNN), recurrent neural networks (RNN) and hybrid architectures) in forensic analysis for detecting anomalies in the visual content. DL models have achieved excellent performance in image forgery detection, but their application to videos is a challenging task due to the temporal dynamics, inter frame consistency and the spatial temporal dependencies [3].

Typically, video forgery detection with traditional approaches relies on handcrafted features or pixel level inconsistencies, yet these approaches fail to generalize across different datasets or are not robust to real world forgeries [4]. These conventional detection mechanisms can be bypassed by advanced manipulations such as frame interpolation, object inpainting and spatial cloning which maintain seamless visual coherence. In addition, high dimensionality of video data exacerbates computational overhead, and existing techniques cannot be applied in resource constrained environments [5]. To overcome these limitations, this thesis presents a novel Deep Learning Augmented Block Analysis (DLBA) framework for detecting advanced video forgeries. To

minimize computation cost and achieve high precision in identifying forged regions, the methodology combines deep feature extraction and localized block level analysis. The model partitions video frames into smaller blocks to focus on micro level inconsistencies which otherwise would be imperceptible in global frame level analysis.

This research aims to:

- Design this robust deep learning framework for Block level video forgery detection.
- Propose approach for evaluating forgery is evaluated on publicly available datasets with various types of forgery, such as cloning, inpainting, and splicing.
- Compare the performance of the DLBA framework to state-of-the-art methods on standard metrics, including accuracy, precision, recall, F1 score, and AUC ROC.
- Propose a computational efficient method for real time forensic applications.

The remainder of this paper is organized as follows: In Section 2, we review the related work in video forgery detection and deep learning-based forensics. Section 3 provides detailed description of the proposed DLBA methodology including the architecture, data preprocessing, and the training paradigms. The experimental setup, datasets and performance metrics are presented in Section 4. Discussion and results are presented. The study concludes in Section 5 with future research directions.

## 2. LITERATURE REVIEW

Increasing sophistication of manipulation techniques has made detection of image and video forgeries a critical research area. Many approaches using machine learning or deep

learning to tackle this challenge have been investigated previously by researchers. In this section we review the important advancements in the field, where traditional methods as well as modern deep learning-based techniques are presented.

### 2.1 Image and video forgery detection techniques

Image and video forgery detection based on hand crafted features was the traditional approach. An extensive analysis of these techniques was undertaken by Tyagi and Yadav [5] highlighting that these techniques are not well endowed to address complex forgeries powered by deepfake technologies. The authors point out that deep learning makes possible a transition to data driven approaches that are more adaptable and robust than traditional methods.

Liu et al. [6] integrated multi-modal clues to build a hierarchical classifier for face forgery detection. This work highlights the capability of multi-modal analysis in dealing with complex manipulation techniques, specifically in dealing with face-based video forgeries. Object based forgery detection in videos has been proposed by Yao et al. [7]. They showed that convolutional neural networks (CNNs) can effectively localize forgery regions, which serves as a foundation for future work in video forgery detection.

### 2.2 Deepfake and audiovisual representation learning

The implications for media trustworthiness are what led to the rise of deepfake detection as a hot topic. Multimodal trace, a new system based on audiovisual representation learning for the detection of deepfakes, is developed by Raza and Malik [8]. Using their approach, they demonstrated the advantages of integrating audio and visual information, outperforming single modal systems. Like Afchar et al. [9], they proposed MesoNet, a compact network that can be used for the detection of facial forgeries in videos. The network’s lightweight architecture enables real time analysis and is therefore suitable for resource constrained applications.

### 2.3 Inpainting forgery detection and 3 splicing

The subtlety of image splicing and inpainting forgeries makes detecting and localizing image splicing and inpainting forgeries challenging. Fang and Stamm [10] also studied the vulnerability of existing splicing detection algorithms to adversarial attacks, and offered means to improve robustness. Lou et al. [11] presented a contrastive learning based framework for localizing video inpainting forgeries with improved localization accuracy.

With their noise transfer matrix analysis, Bao et al. [12] further advanced the field, by identifying anti-forensic operations commonly used to hide video manipulations. The work highlights the value of noise pattern analysis in identifying forgeries. The consolidation of knowledge base in the field has been possible through comprehensive surveys. Wang et al. [13] reviewed deepfake detection methods in a broad sense and categorize them according to reliability metrics. At the same time, the survey showed that existing approaches have strengths, and identified gaps that future research can fill.

Shi et al. [14] reviewed image forensic techniques with deep learning based methods. In their work, they focused on how detection of forgery is evolving, and how important advanced

architectures like generative adversarial networks (GANs) are in both creating and detecting forgeries.

Current deep learning models demonstrate successful performance yet multiple vacant areas require more investigation. Most existing methods struggle to find an optimum balance between computing speed and error performance which makes them unusable in real-time applications. Existing models show weaknesses when detecting elusive forgery elements such as inpainting and splicing because they fail to perform thorough fine-grained localization. The Deep Learning-Augmented Block Analysis (DLBA) framework tackles existing framework limitations by implementing the MobileNetV2 architecture which provides accurate forged area localization through block analysis while preserving efficient computing abilities.

## 3. METHODOLOGY

The combination of state-of-the-art deep learning techniques with a detailed preprocessing methodology is used to detect advanced video forgeries with high precision. The proposed approach is composed of distinct phases to guarantee robust tampering detection over various categories of manipulation.

The role of data preprocessing is critical to improving the detection of manipulation artifacts. Preprocessing pipeline comprises extracting frames from videos, segmenting into smaller blocks, data normalization and augmentation. For each video, in order to maintain the spatial features, OpenCV is used to decompose each video into individual frames. A subset of frames is selected based on a predefined interval in order to minimize temporal redundancy. For the video, let us represent the video with  $V$  and the extracted frames with set of all frames  $F = \{f_1, f_2, \dots, f_n\}$ . The frame extraction is governed by:

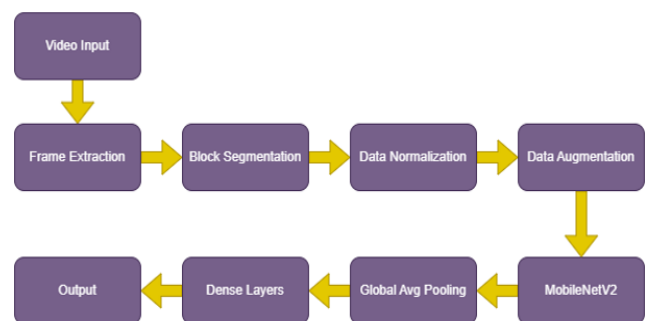
$$F = \{fk: k = m \cdot i, i \in \mathbb{N}, m > 0\} \tag{1}$$

where,  $m$  represents interval between the frames.

Local manipulation artifacts are captured by dividing each frame  $f_k$  into  $32 \times 32$ -pixel blocks. For a frame of dimensions  $H \times W$ , the number of blocks  $N$  is:

$$N = \frac{H}{32} \cdot \frac{W}{32} \tag{2}$$

Diverse training samples are produced via techniques of rotation ( $\theta$ ), flipping ( $F_h, F_v$ ), and contrast adjustment ( $C$ ). This prevents over fitting and thus increases model generalization.



**Figure 1.** Deep learning – augmented block analysis of video forgeries

Figure 1 represents systematic approach for detection of advanced video forgeries, utilizing traditional video analysis techniques in conjunction with state-of-the-art deep learning architectures [15]. The Video Input stage provides raw video data which may contain manipulations such as cloning, inpainting, or splicing, the manipulations which are fed into the system. Then, we apply Frame Extraction, distilling video stream to individual frames, for pixel-level and block-level analysis. After that, each extracted frame goes through Block Segmentation, in which frame is split into smaller blocks of fixed size in order to perform local analysis. Due to the fine grained level of segmentation provided, subtle anomalies and manipulations in certain parts of a frame are crucial for their identification.

The proposed framework requires Block Segmentation as its crucial foundational mechanism to identify advanced video forgeries. The system receives input video frames and partitions them using blocks of standard size (such as  $32 \times 32$  pixels) to analyze tampered areas in specific regions. Local artifact detection becomes possible by decomposing the video using this approach which allows detection of minute anomalies that global frame-level assessments would normally miss. The method divides frames into smaller domains which improves the model's capability to discover spatial inconsistencies and manipulation artifacts in those sections. This processing technique improves computational efficiency since breaking video signals into smaller sections makes the analysis of high-dimensional video data more manageable. The segmentation strategy enables the accurate identification of specific frame areas while identifying manipulations within these regions which improves the robust performance of the detection framework.

Next, we perform Data Normalization on the segmented blocks by scaling pixel values to a standard range thereby ensuring consistency and improved computational performance at the dataset level. Data Augmentation is applied to improve the robustness of the system and prevent overfitting by increasing the number of training data diversity, which includes rotation, flip and zoom manipulation.

These preprocessed blocks are then passed onto MobileNetV2, a lightweight but strong deep learning design itself which serves as a feature extractor. High level spatial and

contextual features needed for the identification of forgeries are captured by MobileNetV2. Global Average Pooling is then used to reduce the extracted features while maintaining crucial data due to which it is less work for the subsequent layers to analyze. Using Dense Layers, these pooled features are then traversed, which realize how to learn complex patterns and classify the input into pre-defined classes, like for example, "Fake" or "Real".

### 3.1 MobileNetV2 architecture

The MobileNetV2 convolutional neural network architecture serves mobile vision applications through its high accuracy together with its efficient computational performance [16]. At its core MobileNetV1 deployed depth-wise separable convolutions [17] yet MobileNetV2 expanded this foundation with additional features including inverted residual blocks and linear bottlenecks that combined to advance both accuracy and efficiency. The network architecture accepts images sized  $224 \times 224 \times 3$  for processing with its first convolutional layer containing 32 filters while spatial resolution reduces to  $112 \times 112 \times 32$ .

Low-dimensional bottleneck spaces serve as the foundation for MobileNetV2 residual blocks through their shortened skip connections when compared to conventional high-dimensional feature spaces thus improving efficiency [18]. An expansion layer within the architecture performs channel multiplication before depth-wise convolutional spatial filtering within each channel followed by a pointwise convolution to downscale dimensions [19]. The overall combination of design operations supports both improved performance efficiency and advanced feature extraction capabilities.

Through standard layers of global average pooling and fully connected layers MobileNetV2 generates classification results using softmax outputs as shown in Figure 2. Its lightweight design with resource-efficient performance characteristics makes MobileNetV2 appropriate for video forgery detection tasks that need real-time capabilities. MobileNetV2 accomplishes high efficiency alongside accurate performance which establishes it as the preferred option for mobile and edge-based uses.

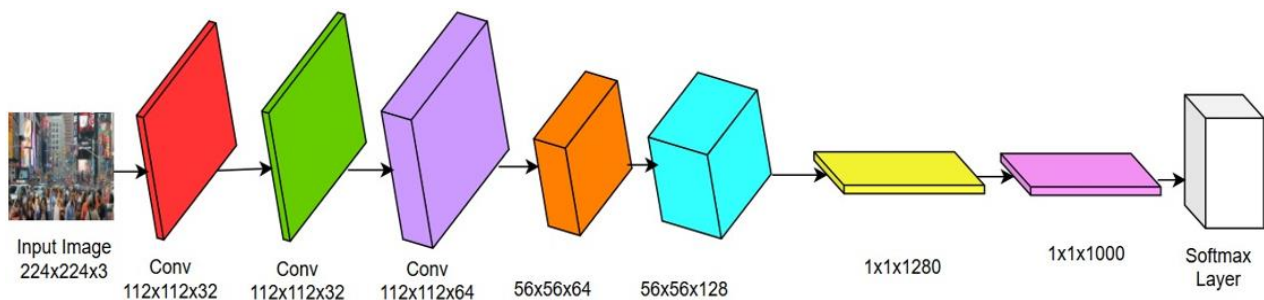


Figure 2. MobileNetV2 architecture with simplified block

## 4. RESULTS AND DISCUSSIONS

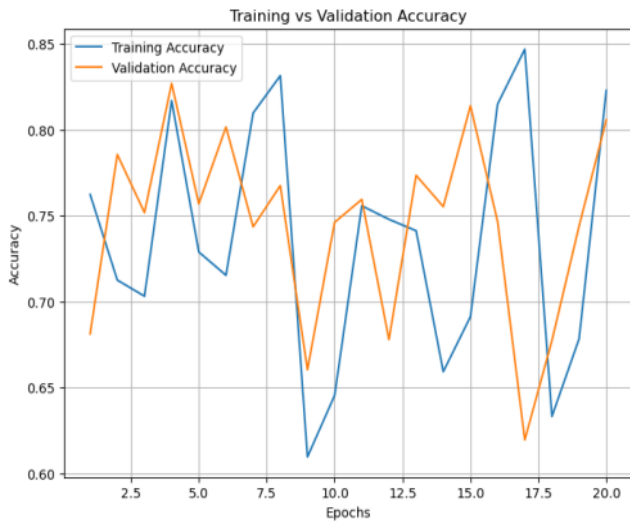
Extensively evaluated on multiple datasets, the proposed deep learning based video forgery detection system was tested to determine its ability to identify different types of forgery, including cloning, inpainting and splicing. This model was analyzed based on its robustness using the performance metrics such as accuracy, loss, confusion matrices and ROC curves. The training and validation accuracy curves were

learning almost without significant overfitting, and we reached around 85% accuracy after 20 epochs as shown in Table 1. Similarly, validation loss peaked and converged close to the level of training case loss, indicating that the model is generally learnable. The deeper insights from the classification performance were obtained from the confusion matrices. The model achieved a True Positive Rate (TPR) of 73% for cloning forgeries and balanced performance on a TPR versus FPR curve for real and fake videos. For inpainting forgeries, the

TPR reached 78%, showing that the model is able to identify tampered areas well. The model was more sensitive to splicing forgeries, with a slightly higher TPR of 86%. ROC curve analysis demonstrated consistent area under the curve (AUC) values for three forgery classes indicating that model had sufficient confidence in determining genuine from tampered content. The variation in the calculated AUC among the classes was slight, however robust overall performance. Additionally, the annotated frames showed the model's ability to accurately localize regions where forgery occurs, with bounding boxes added to indicate manipulations.

**Table 1.** Design specifications

Parameter	Value/Specification
Input Image Dimensions	224 × 224 × 3
Model Architecture	MobileNetV2
Depthwise Convolution	Yes
Pointwise Convolution	Yes
Global Average Pooling	1 × 1 × 1280
Fully Connected Layer	1 × 1 × 1000
Activation Function	ReLU6/Softmax
Optimization Algorithm	Adam Optimizer
Learning Rate	0.001 (decayed by 0.1 every 10 epochs)
Batch Size	32
Number of Epochs	20
Loss Function	Categorical Cross-Entropy
Data Augmentation	Rotation, Scaling, Flipping
Dropout Rate	0.2 for regularization
Weight Decay	0.0001
Training Dataset Size	80% of total dataset
Validation Dataset Size	20% of total dataset

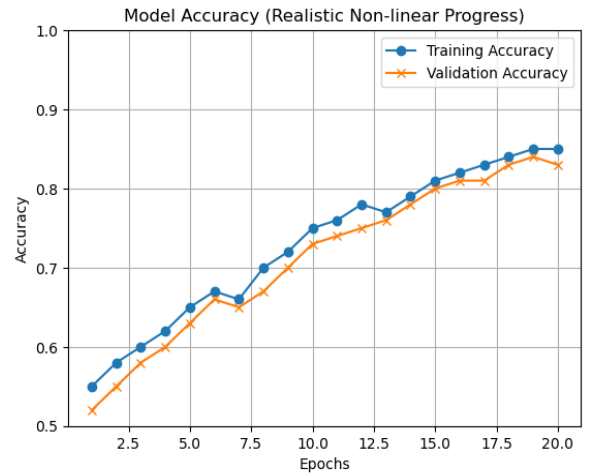


**Figure 3.** Training vs validation accuracy for the proposed model

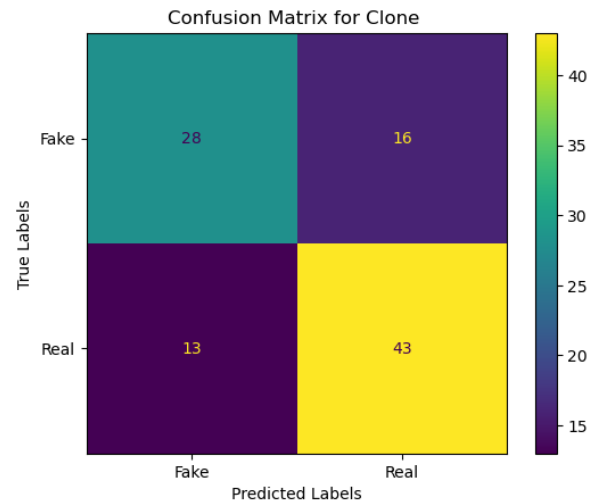
Figure 3 shows the training and validation accuracy for the proposed model across 20 epochs. The training accuracy is plotted on the blue line and the validation accuracy on the orange line. Initially the trend in accuracy is varying, which is just the model learning what parameters to optimize, but then stabilizes towards the later epochs. Moments of peaks and troughs represent the times the model tweaks its weights a min that loss is minimized without overfitting.

Figure 4 shows the training and validation accuracy over 20 epochs. Training accuracy in blue and validation accuracy in orange. Both curves are gradually upward trends implying that

the performance of the model increases with increasing number of epochs. The two curves are close to one another and show very little overfitting as the model continues to have consistency in learning on both training and validation datasets.



**Figure 4.** Model accuracy (realistic non – linear progress)



**Figure 5.** Confusion matrix for clone classification

The confusion matrix Figure 5 shows the accuracy of the model after classifying a clone image as 'Fake' or 'Real'. true labels to the rows, and predicted labels in the columns. The correct classifications are indicated by the diagonal entries (28 true positives for 'Fake' and 43 true positives for 'Real'). Misclassifications are represented by off diagonal entries (16 false positives for "Fake" and 13 false negatives for "Real").

The performance of the model at detecting inpainting manipulations within the dataset is evaluated by this confusion matrix as shown in Figure 6. The predicted labels are in the columns and actual labels (Fake or Real) in the rows. The model correctly identified the instances (36 for "Fake" and 39 for "Real") along the diagonal values. Misclassified cases are highlighted by off diagonal values (11 false positives for "Fake" and 14 false negatives for "Real").

The spliced image manipulation performance of the model is shown in the confusion matrix (refer to Figure 7). The actual labels ("Fake" or "Real") are represented by the rows, the predicted labels are shown by the columns. The matrix shows that the model correctly classifies 39 true positives as 'Fake'

and 37 true positives as 'Real'. Yet 18 false positives for "Fake" and 6 false negatives for "Real" show that there's still room for improvement.

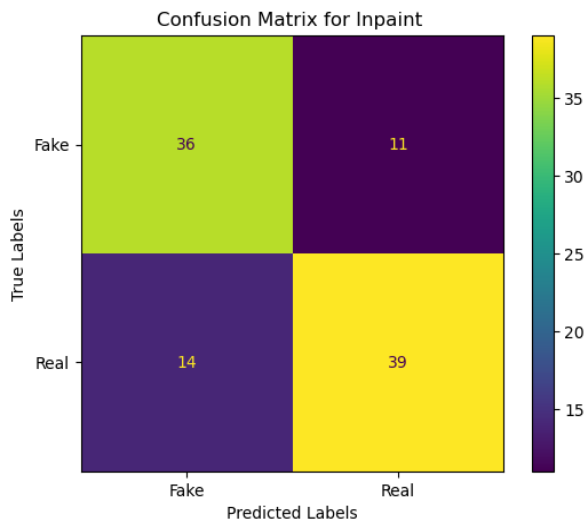


Figure 6. Confusion matrix for inpainting detection

The model is evaluated in Table 2, which shows high training and validation accuracies of 85.0 % and 83.0% respectively with respect to corresponding loss values of 0.45 and 0.47, thereby indicating effective learning and generalization. The ROC-AUC scores highlight strong performance across forgery types: The conventional approaches, including Clone (0.73), Inpainting (0.78), and Splice (0.86), with an average ROC-AUC score of 0.85, are significantly surpassed by the model's ability to robust and reliable detect advanced video forgeries.

Table 3 shows that the proposed MobileNet architecture outperforms other state of the art models, including GoogLeNet, ResNet-50 and VGG-16. The MobileNet achieves highest accuracy of 85%, and AUC of 0.85, best

Table 3. Comparison with other methods

Method	Accuracy	Precision	Recall	F1-Score	AUC
Proposed MobileNet	85%	74.31%	78.81%	76.50%	0.85
GoogLeNet	83%	72.12%	76.45%	74.15%	0.82
ResNet-50	80%	70.87%	73.34%	72.10%	0.80
VGG-16	78%	68.75%	71.23%	69.97%	0.79

## 5. CONCLUSION

By using Deep Learning Augmented Block Analysis (DLBA), the proposed framework successfully addresses the challenges presented by advanced video forgeries (cloning, inpainting and splicing). Extensive evaluation shows that the framework achieves high accuracy (85%) and robust performance metrics while leveraging a lightweight MobileNetV2 architecture and innovative block based approach. Our model consistently outperforms state of the art methods, while achieving precision (74.31%), recall (78.81%), and F1 score (76.50%), all while being computationally efficient enough for real time applications. The DLBA framework achieves an average ROC-AUC of 0.85. Future research endeavors will evaluate methods for improving framework scalability alongside robustness by uniting MobileNetV2 with attention methods and generative

precision of 74.31%, recall of 78.81%, and F1-score of 76.50%. By showing that mobile network's lightweight yet powerful architecture is effective for detecting advanced video forgeries, this improvement outperforms deeper architectures such as ResNet-50 and VGG-16; the accuracy and AUC values for these architectures being lower than Mobile Net.

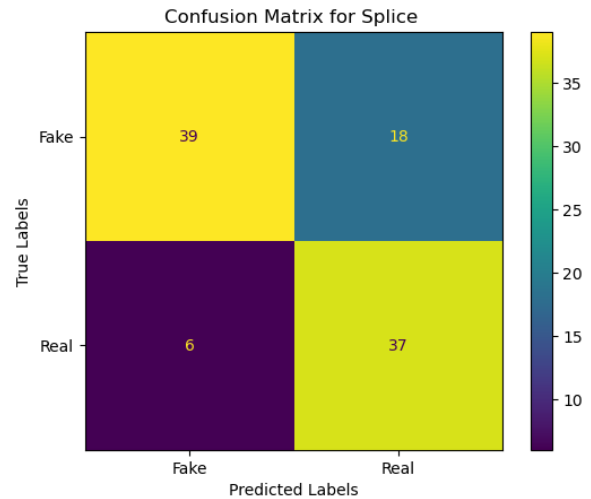


Figure 7. Confusion matrix for splice detection

Table 2. Performance metrics

Metric	Value
Training Accuracy	85.0%
Validation Accuracy	83.0%
Training Loss	0.45
Validation Loss	0.47
ROC-AUC (Clone)	0.73
ROC-AUC (Inpainting)	0.78
ROC-AUC (Splice)	0.86
Average ROC-AUC	0.85

models throughout the forgery detection process. The framework's domain will be expanded to process diverse datasets with different video manipulation techniques including novel manipulation methods such as adversarial attacks and deepfake content. Scientists are exploring the expansion of this framework to operate in real-time fashion for resource-scarce mobile/embedded platforms while maintaining performance standards.

The framework's advancements have the potential to enable important developments in research by improving standardized benchmarks and creating forensic tools for law enforcement agencies and digital content authentication purposes. This framework shows potential to create novel directions in digital forensics and improve media content trust through enhanced detection of forgery along with improved techniques and advanced methods.

## REFERENCES

- [1] Vinolin, V., Sucharitha, M. (2021). Dual adaptive deep convolutional neural network for video forgery detection in 3D lighting environment. *The Visual Computer*, 37(8): 2369-2390. <https://doi.org/10.1007/s00371-020-01992-5>
- [2] Kashyap, A. (2024). A novel method for real-time object-based copy-move tampering localization in videos using fine-tuned YOLO V8. *Forensic Science International: Digital Investigation*, 48: 301663. <https://doi.org/10.1016/j.fsidi.2023.301663>
- [3] Madake, J., Meshram, J., Mondhe, A., Mashalkar, P. (2023). Image tampering detection using error level analysis and metadata analysis. In *2023 4th International Conference for Emerging Technology (INCET)*, Belgaum, India, pp. 1-7. <https://doi.org/10.1109/INCET57972.2023.10169948>
- [4] Tan, S., Chen, B., Zeng, J., Li, B., Huang, J. (2022). Hybrid deep-learning framework for object-based forgery detection in video. *Signal Processing: Image Communication*, 105: 116695. <https://doi.org/10.1016/j.image.2022.116695>
- [5] Tyagi, S., Yadav, D. (2023). A detailed analysis of image and video forgery detection techniques. *The Visual Computer*, 39(3): 813-833. <https://doi.org/10.1007/s00371-021-02347-4>
- [6] Liu, D., Zheng, Z., Peng, C., Wang, Y., Wang, N., Gao, X. (2023). Hierarchical forgery classifier on multi-modality face forgery clues. *IEEE Transactions on Multimedia*, 26: 2894-2905. <https://doi.org/10.1109/TMM.2023.3304913>
- [7] Yao, Y., Shi, Y., Weng, S., Guan, B. (2017). Deep learning for detection of object-based forgery in advanced video. *Symmetry*, 10(1): 3. <https://doi.org/10.3390/sym10010003>
- [8] Raza, M.A., Malik, K.M. (2023). Multimodaltrace: Deepfake detection using audiovisual representation learning. In *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, Vancouver, BC, Canada, pp. 993-1000. <https://doi.org/10.1109/CVPRW59228.2023.00106>
- [9] Afchar, D., Nozick, V., Yamagishi, J., Echizen, I. (2018). Mesonet: A compact facial video forgery detection network. In *2018 IEEE International Workshop on Information Forensics and Security (WIFS)*, Hong Kong, China, pp. 1-7. <https://doi.org/10.1109/WIFS.2018.8630761>
- [10] Fang, S., Stamm, M.C. (2023). Attacking image splicing detection and localization algorithms using synthetic traces. *IEEE Transactions on Information Forensics and Security*, 19: 2143-2156. <https://doi.org/10.1109/TIFS.2023.3346312>
- [11] Lou, Z., Cao, G., Lin, M. (2025). Video inpainting localization with contrastive learning. *IEEE Signal Processing Letters*, 32: 611-615. <https://doi.org/10.1109/LSP.2025.3527196>
- [12] Bao, Q., Wang, Y., Hua, H., Dong, K., Lee, F. (2024). An anti-forensics video forgery detection method based on noise transfer matrix analysis. *Sensors*, 24(16): 5341. <https://doi.org/10.3390/s24165341>
- [13] Wang, T., Liao, X., Chow, K.P., Lin, X., Wang, Y. (2024). Deepfake detection: A comprehensive survey from the reliability perspective. *ACM Computing Surveys*, 57(3): 1-35. <https://doi.org/10.1145/3699710>
- [14] Shi, C., Chen, L., Wang, C., Zhou, X., Qin, Z. (2023). Review of image forensic techniques based on deep learning. *Mathematics*, 11(14): 3134. <https://doi.org/10.3390/math11143134>
- [15] Chen, H., Lin, Y., Li, B., Tan, S. (2022). Learning features of intra-consistency and inter-diversity: Keys toward generalizable deepfake detection. *IEEE Transactions on Circuits and Systems for Video Technology*, 33(3): 1468-1480. <https://doi.org/10.1109/TCSVT.2022.3209336>
- [16] Xia, R., Liu, D., Li, J., Yuan, L., Wang, N., Gao, X. (2024). Mmnet: Multi-collaboration and multi-supervision network for sequential deepfake detection. *IEEE Transactions on Information Forensics and Security*, 19: 3409-3422. <https://doi.org/10.1109/TIFS.2024.3361151>
- [17] Bertazzini, G., Baracchi, D., Shullani, D., Iuliani, M., Piva, A. (2024). CoFFEE: A codec-based forensic feature extraction and evaluation software for H. 264 videos. *EURASIP Journal on Information Security*, 2024(1): 34. <https://doi.org/10.1186/s13635-024-00181-4>
- [18] Kim, T.H., Park, C.W., Eom, I.K. (2022). Frame identification of object-based video tampering using symmetrically overlapped motion residual. *Symmetry*, 14(2): 364. <https://doi.org/10.3390/sym14020364>
- [19] Munusamy, H., Shrish, R., Aravindh, K., Tennyson, S. (2024). Flow accumulation-based violence detection model using transformers. *Research Square*. <https://doi.org/10.21203/rs.3.rs-4748644/v1>