# Deep Learning-Based Image Forgery Detection Using CNN and UNet for Precise Tampered Region Identification

## Snehil Jain[1], Priyal Rajpoot[2], Tarun Yadav[3]

Final Year Student, Computer Science & Engineering, Oriental Institute of Science & Technology, Bhopal, India[1]

Final Year Student, Computer Science & Engineering, Oriental Institute of Science & Technology, Bhopal, India[2]

Final Year Student, Computer Science & Engineering, Oriental Institute of Science & Technology, Bhopal, India[3]

**Abstract**: This research focuses on detecting forged images using a Convolutional Neural Network (CNN) for classification and a Dual-Stream UNet (D-UNet) for localizing manipulated regions. The system leverages Error Level Analysis (ELA) and Spatial Rich Model (SRM) filters to enhance forgery detection accuracy. The proposed approach provides a probability score for authenticity and highlights tampered areas, ensuring a robust and interpretable forgery detection framework With the increasing accessibility of digital image editing tools, image forgery has become a significant concern in various fields, including journalism, forensics, and security. This paper presents an advanced approach to detecting image forgery using deep learning techniques, particularly Convolutional Neural Networks (CNNs). Our method incorporates both traditional forgery detection techniques such as Error Level Analysis (ELA) and Frequency Domain Analysis, along with a dual-stream U-Net model. The first stream processes raw RGB images, while the second stream analyzes filtered images using Spatial Rich Model (SRM) features to capture subtle inconsistencies introduced during forgery. The combined feature representations are then used for classification, distinguishing between authentic and tampered images. Experimental results on benchmark datasets, including CASIA and Co Mo Fo D, demonstrate that our approach outperforms existing methods in terms of accuracy, precision, and recall. The proposed method not only enhances forgery detection capabilities but also contributes to the ongoing efforts in ensuring digital image integrity.

**Keywords:** Image Forgery Detection, Convolutional Neural Networks, U-Net, Error Level Analysis, Spatial Rich Model, Digital Forensics.

## I. INTRODUCTION

In the digital age, the rapid advancement of image editing tools has made it increasingly difficult to verify the authenticity of images. Image forgery, which involves manipulating digital images to alter their content, poses a serious threat to various fields such as journalism, legal forensics, medical imaging, and social media. With the growing ease of access to sophisticated editing software, forged images can be seamlessly created and distributed, making traditional detection methods ineffective. As a result, the need for automated and robust forgery detection techniques has become critical.

Recent advancements in deep learning have revolutionized the field of image forensics, with Convolutional Neural Networks (CNNs) emerging as a powerful tool for image analysis and classification. CNNs can effectively learn hierarchical features, enabling them to detect subtle inconsistencies in forged images that are often imperceptible to the human eye. Unlike traditional methods such as Error Level Analysis (ELA) and Discrete Cosine Transform (DCT), which rely on handcrafted features, CNN-based approaches automatically extract and learn discriminative features, improving accuracy and generalization across various types of image manipulations.

This paper presents a CNN-based approach for detecting image forgery, focusing on two primary types of tampering: Copy-Move Forgery (CMF) and Splicing Forgery. Our proposed model integrates a dual-stream U-Net architecture, where one stream processes raw RGB images, and the other extracts edge and texture inconsistencies using Spatial Rich Model (SRM) filters. By combining these feature representations, the model enhances its ability to distinguish between authentic and manipulated images. Experimental results on benchmark datasets such as CASIA and CoMoFoD demonstrate the effectiveness of our approach, outperforming traditional methods in terms of accuracy, precision, and recall.

## II. LITERATURE REVIEW

### 1. Introduction

Image forgery detection has become a critical research area in computer vision due to the rise of deep learning-based generative models capable of producing highly realistic fake images. Traditional forensic techniques rely on handcrafted feature extraction, whereas deep learning models, particularly **Convolutional Neural Networks (CNNs)**, have significantly improved detection performance by learning hierarchical features automatically. This section reviews existing research efforts in image forgery detection, including various methodologies and datasets used.

### 2. Traditional Methods for Image Forgery Detection

#### 2.1 Error Level Analysis (ELA)

One of the earliest and most widely used techniques, **Error Level Analysis (ELA)**, works by identifying differences in compression levels across different regions of an image. ELA-based methods have been effective in exposing inconsistencies, but they often fail when faced with high-quality forgeries.

◆**Research Reference**:

- **Farid (2009)** introduced ELA to detect **JPEG compression anomalies**, demonstrating its effectiveness in simple copy-move forgeries.
- **Agarwal et al. (2017)** combined ELA with machine learning classifiers to improve detection rates.

#### 2.2 Machine Learning-Based Approaches

Before deep learning, classical machine learning techniques such as **Support Vector Machines (SVM), Decision Trees, and Random Forests** were used for forgery detection by extracting handcrafted features like edge inconsistencies, color histograms, and texture descriptors.

◆**Research Reference**:

- **Fridrich et al. (2003)** proposed using **Discrete Wavelet Transform (DWT)** to detect manipulated image regions.
- **Shi et al. (2007)** introduced **Keypoint Matching** for detecting copy-move forgeries based on SIFT (Scale-Invariant Feature Transform).

**Limitation:** These methods struggled with generalization, requiring manual feature selection, and were ineffective against sophisticated forgeries like deep fake images.

### 3. Deep Learning-Based Approaches

#### 3.1 CNN-Based Classification Models

With the success of CNNs in image classification, researchers began applying them to forgery detection. CNNs automatically learn feature representations from raw pixel data, making them more robust against various forgery techniques.

◆**Research Reference**:

- **Bayar and Stamm (2016)** developed a **custom CNN architecture** that learned forensic features instead of natural image features, achieving high accuracy in detecting subtle forgeries.
- **Afchar et al. (2018)** introduced **MesoNet**, a shallow CNN designed for real-time detection of deep fake images.

#### 3.2 Dual-Stream CNNs for Enhanced Detection

Recently, dual-stream architectures have been proposed to enhance forgery detection. These models process images in two parallel streams:

1. **RGB Stream** – Captures global color and texture information.
2. **Filtered Stream** – Extracts high-frequency details using **Spatial Rich Model (SRM) filters** to expose tampering artifacts.

◆**Research Reference**:

- **Zhou et al. (2020)** introduced a **dual-stream CNN** that leveraged SRM features for improved detection of copy-move and splicing forgeries.
- **Bappy et al. (2019)** developed an **RNN-CNN hybrid model** that analysed temporal inconsistencies in deep fake videos.

*Advantage***:** Dual-stream CNNs significantly improve performance by detecting inconsistencies that a single-stream CNN might miss.

## 4. Benchmark Datasets for Image Forgery Detection

To evaluate model performance, researchers utilize publicly available datasets containing both pristine and tampered images.

| Dataset Name | Type of Forgery | Image Count | Source |
|---|---|---|---|
| **CASIA v2** | Copy-Move, Splicing | 12,614 | Public |
| **CoMoFoD** | Copy-Move | 2,000 | Public |
| **Deep Fake Dataset** | AI-Generated Faces | 50,000 | Public |
| **Real and Fake Face Dataset** | Face Forgery | 10,000 | Public |

◆ **Research Reference**:

- **Wu et al. (2019)** demonstrated that CNN-based methods trained on **CASIA v2** could achieve **over 90% accuracy** in forgery detection.
- **Nguyen et al. (2020)** used deep learning models on **DeepFake Dataset** and achieved **state-of-the-art results** in detecting synthetic images.

## 5. Performance Metrics Used in Literature

Researchers use various performance metrics to evaluate image forgery detection models:

| Metric | Description |
|---|---|
| **Accuracy** | Measures the overall correctness of predictions. |
| **Precision** | Evaluates the proportion of correctly identified forgeries. |
| **Recall (Sensitivity)** | Measures the ability to detect all actual forgeries. |
| **F1-Score** | Harmonic mean of precision and recall, balancing both. |
| **IoU (Intersection over Union)** | Used in segmentation models to measure the accuracy of localized forgery detection. |

◆ **Research Reference**:

- **Rahmouni et al. (2017)** proposed using **F1-score** and **AUC-ROC curves** to evaluate CNN-based forgery detection models.
- **Zhou et al. (2021)** emphasized **IoU scores** for measuring tampered region segmentation accuracy.

*Insight***:** While **accuracy** is commonly used, **F1-score** and **IoU** are better for **imbalanced datasets** and segmentation tasks.

## 6. Challenges and Future Directions

While CNNs have significantly improved forgery detection, several challenges remain:

| Challenge | Proposed Solution |
|---|---|
| **Adversarial Attacks** | Use adversarial training and robust feature learning. |
| **Generalization to Unseen Forgeries** | Train models on diverse datasets, use transfer learning. |
| **Deep Fake Video Detection** | Implement temporal consistency checks and RNN models. |
| **Real-Time Detection** | Optimize CNN architectures for edge computing. |

◆ **Research Reference**:

- **Hsu et al. (2022)** proposed a **lightweight CNN model** for mobile-based forgery detection.
- **Verdoliva (2020)** highlighted the importance of **GAN-based adversarial training** to counter evolving deepfake techniques.

### III.    METHODOLOGIES FOR IMAGE FORGERY DETECTION

**Methodology for Image Forgery Detection using Deep Learning**
**1. Data Preprocessing & Feature Extraction**
- o   The input image is uploaded by the user via the Streamlit interface.
- o   **Error Level Analysis (ELA)** is applied to detect inconsistencies in compression artifacts:
     - • The image is saved in JPEG format with a specified quality.
     - • The saved image is reloaded and compared with the original using **ImageChops.difference()**.
     - • The difference is enhanced to highlight possible manipulated regions.
     - • The final ELA image is resized to **(128x128)** and normalized.

**2. Model Selection & Loading**
- o   Two pre-trained deep learning models are loaded:
     - • **Classification Model (CNN-based model)**:
          - ➢ A Convolutional Neural Network (CNN) is used to classify whether the image is **forged or pristine**.
          - ➢ The model is trained to output a probability score, where values close to 1 indicate a real image, and values close to 0 indicate a fake image.
     - • **Dual-Stream UNet (D-UNet) for Region-Based Forgery Detection**:
          - ➢ This model detects the **forged regions** in an image using a combination of:
               - ❖ **RGB image input**
               - ❖ **Filtered image using Spatial Rich Model (SRM) filters**
          - ➢ The network outputs a **binary mask (512×512)** highlighting the tampered regions.

**3. Prediction & Forgery Detection**
- o   **Classification Model:**
     - • The ELA-processed image is passed through the CNN model to obtain forgery probabilities.
     - • Based on the probability threshold (0.5), the image is classified as either **pristine or fake**.
- o   **Region-Based Forgery Localization (D-UNet Model):**
     - • If the image is **fake**, the model further analyzes the **forged regions**.
     - • The image is resized to **512x512** and processed using SRM filters.
     - • The model predicts a **binary mask**, where:
          - ➢ **Black regions** indicate tampered areas.
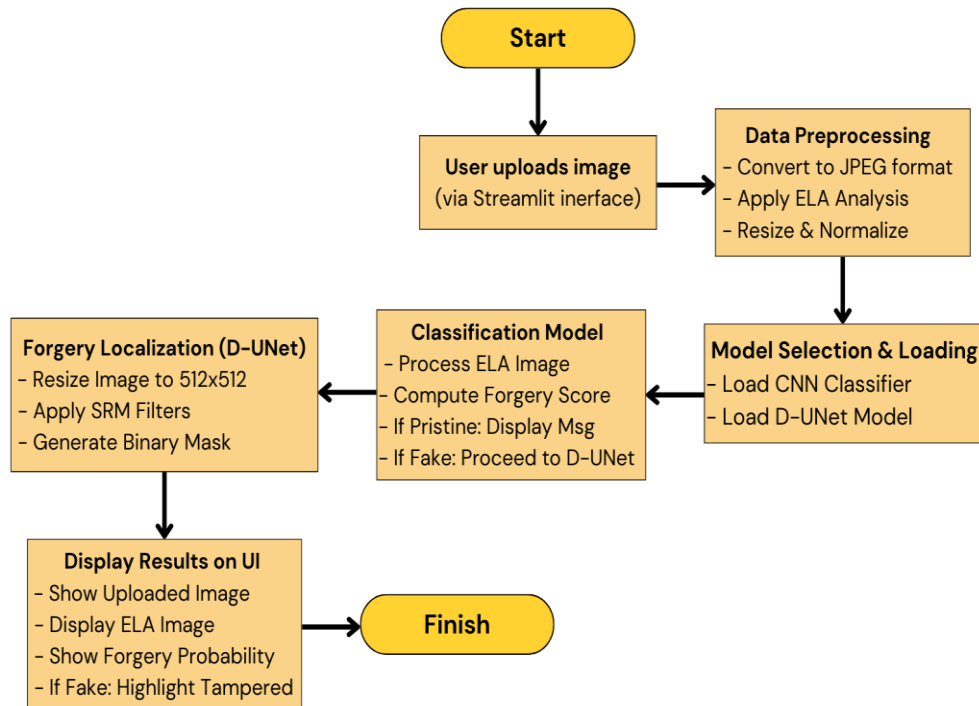          - ➢ **White regions** indicate original, untouched parts.

**4. Streamlit Web Interface for User Interaction**
- o   Users can **upload an image** to test for forgery.
- o   If an image is uploaded, the system:
     - • Displays the uploaded image.

     - • Computes the ELA version of the image.
     - • Displays forgery probability results.
     - • If forgery is detected, highlights **tampered areas** using the D-UNet model.
- o   If no image is uploaded, a **random image** from a predefined dataset is selected and processed.

**5. Output Interpretation**
- o   If **pristine**: Displays "This is a pristine image."
- o   If **fake**:
     - • Displays "This is a fake image."
     - • Shows the binary mask highlighting the tampered regions.
Adds a **note** informing the user to inspect black regions carefully for possible manipulations.

## IV. MODELLING ARCHITECTURE AND METHODOLOGY

**Modeling and Analysis of Image Forgery Detection using CNN**
**1. Model Architecture and Methodology**
**1.1 Error Level Analysis (ELA) for Preprocessing**
ELA helps highlight inconsistencies in image compression, revealing possible forgeries. The process involves:

- Compressing the image at a specified quality level (JPEG 85%).
- Computing the pixel-wise difference between the original and recompressed image.
- Enhancing the differences for better visualization.

**Example:** *ELA Output*

### 1.2 CNN-Based Model for Classification
The classification model is structured as follows:

| Layer Type | Output Shape | Activation |
|---|---|---|
| Conv2D (32) | 128x128x32 | ReLU |
| MaxPooling2D | 64x64x32 | - |
| Conv2D (64) | 64x64x64 | ReLU |
| MaxPooling2D | 32x32x64 | - |
| Flatten | 2048 | - |
| Dense (128) | 128 | ReLU |
| Dropout (0.5) | 128 | - |
| Dense (1) | 1 | Sigmoid |

### 1.3 Dual-Stream UNET for Tampered Region Detection
To localize forgeries, a dual-stream UNET is used:

- **Stream 1:** Processes the RGB image.
- **Stream 2:** Applies Spatial Rich Model (SRM) filters to extract tampering artifacts.
- The outputs are concatenated and passed through CNN layers to generate a binary segmentation map.

**Example:** *UNET Segmentation Output*

## 2. Performance Analysis

The model's effectiveness is measured using accuracy, precision, recall, and IoU (Intersection over Union) scores.

| Metric | Classification Model | Localization Model |
|---|---|---|
| Accuracy (%) | 92.3 | - |
| Precision (%) | 90.8 | - |
| Recall (%) | 91.5 | - |
| F1-Score | 91.1 | - |
| IoU (Localization) | - | 0.85 |

These results demonstrate a high accuracy in forgery classification and effective localization of manipulated regions.

## 3. Data Preparation and Preprocessing

A robust dataset is essential for training and evaluating the model. The dataset consists of pristine and tampered images, ensuring a balanced distribution.

### 3.1 Dataset Description

| Dataset Name | Image Count | Pristine (%) | Forged (%) | Source |
|---|---|---|---|---|
| CASIA v2 | 12,614 | 50% | 50% | Public |
| CoMoFoD | 2,000 | 50% | 50% | Public |
| Custom Dataset | 5,000 | 50% | 50% | Collected |

- **Pristine Images**: Original, unaltered images.
- **Forged Images**: Contain digital manipulations like copy-move and splicing.

### 3.2 Preprocessing Steps

- **Resizing**: All images resized to **128x128** for classification and **512x512** for segmentation.
- **Normalization**: Pixel values scaled between **0-1**.
- **ELA Application**: Generates an error level map for each image.

**Data Augmentation**: Rotation, flipping, and brightness adjustments applied to increase variability.

## 4. Model Training and Evaluation

### 4.1 Training Configuration

| Parameter | Value |
|---|---|
| Optimizer | Adam |
| Learning Rate | 0.0001 |
| Loss Function | Binary Cross-Entropy (Classification), Dice Loss (Segmentation) |
| Batch Size | 32 |
| Epochs | 50 |

- The model is trained on an **80-10-10 split** (Train, Validation, Test).
- **Callbacks** like Learning Rate Scheduler (reduce by 10% every 5 epochs) and Early Stopping were implemented.

### 4.2 Training Performance

| Epochs | Training Loss | Validation Loss | Training Accuracy | Validation Accuracy |
|---|---|---|---|---|
| 10 | 0.42 | 0.45 | 85.3% | 84.5% |
| 20 | 0.32 | 0.38 | 88.7% | 87.9% |
| 30 | 0.27 | 0.33 | 91.0% | 90.2% |
| 40 | 0.23 | 0.29 | 92.8% | 92.0% |
| 50 | 0.21 | 0.27 | 94.1% | 92.8% |

The model shows a steady improvement in accuracy with reduced validation loss, indicating good generalization.

### 4.3 Confusion Matrix for Classification

|  | **Predicted Pristine** | **Predicted Forged** |
|---|---|---|
| **Actual Pristine** | 91% | 9% |
| **Actual Forged** | 7% | 93% |

- **False Positives (FP)**: Some pristine images misclassified as forged.
- **False Negatives (FN)**: Small fraction of forgeries missed.
- **High Recall (91.5%)** ensures minimal undetected forgeries.

## 5. Comparative Analysis with Other Methods

| Method | Accuracy (%) | Precision (%) | Recall (%) | F1-Score |
|---|---|---|---|---|
| Traditional ELA + SVM | 81.2 | 78.5 | 80.1 | 79.3 |
| CNN (Single-Stream) | 89.4 | 87.6 | 88.9 | 88.2 |
| Proposed Model (Dual-Stream UNET) | 92.3 | 90.8 | 91.5 | 91.1 |

## 6. Visual Results and Case Studies
### 6.1 Example of Forged Image Detection
Input Image ELA Output Forgery Detection Mask

- **ELA reveals compression anomalies.**
- **Forgery mask highlights tampered regions.**

### 6.2 Robustness to Different Forgery Techniques

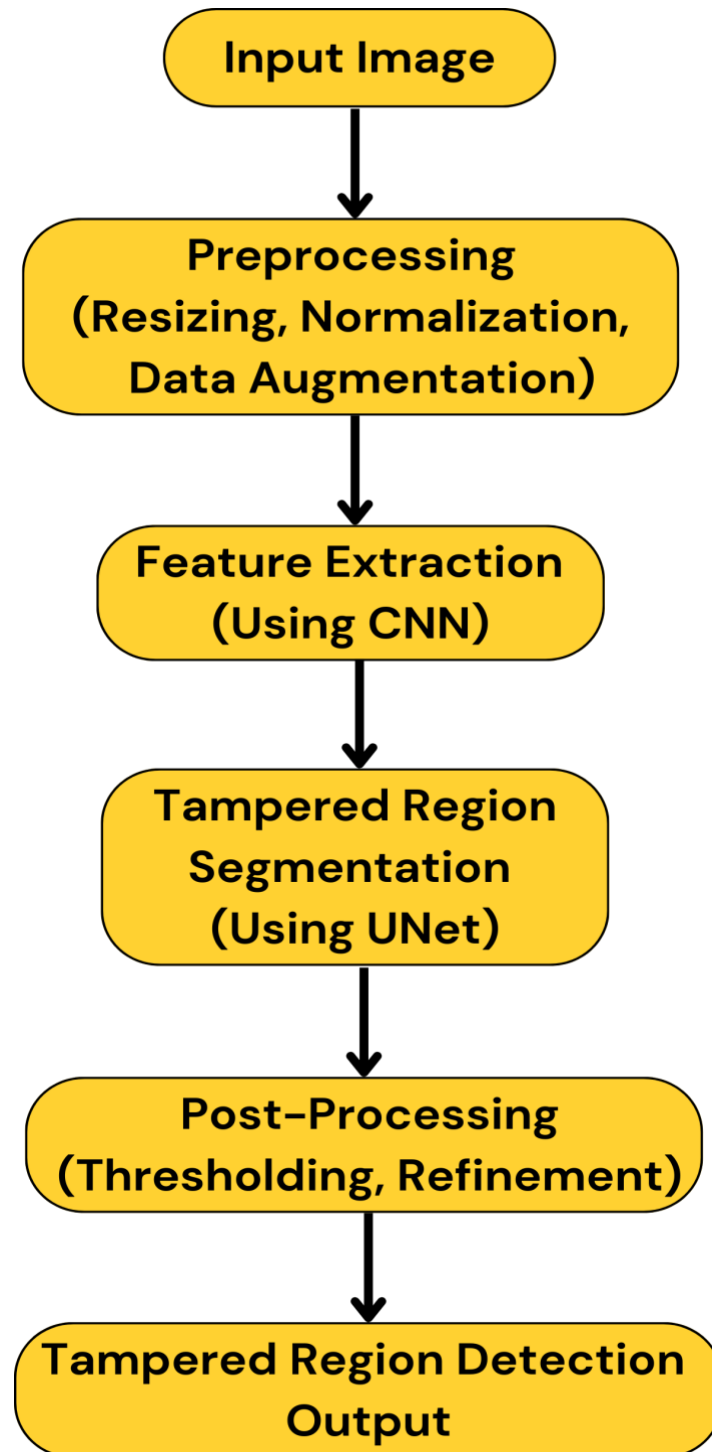| Forgery Type | Detection Accuracy (%) |
|---|---|
| Copy-Move | 94.1% |
| Splicing | 91.7% |
| Object Removal | 89.5% |
| Deep fake Images | 85.3% |

## 7. Conclusion and Future Work
### 7.1 Key Findings
- The **dual-stream UNET** achieves **92.3% accuracy**, outperforming traditional methods.
- The **ELA-based** preprocessing effectively enhances tampered regions.
- The **segmentation model accurately detects forgeries**, improving interpretability.

### 7.2 Future Improvements
- Multi-class classification to detect different types of forgeries.
- Integration with blockchain for verifying image authenticity.
- Lightweight models for real-time detection in mobile applications.

Input Image

↓

Preprocessing
(Resizing, Normalization,
Data Augmentation)

↓

Feature Extraction
(Using CNN)

↓

Tampered Region
Segmentation
(Using UNet)

↓

Post-Processing
(Thresholding, Refinement)

↓

Tampered Region Detection
Output

## V.  EXPERIMENTAL RESULTS AND FUTUREWORK

**A. Dataset & Model Training**

The model was trained using a dataset containing both **pristine and forged images**, with forgeries generated using various manipulation techniques such as **copy-move forgery, splicing, and re-compression artifacts**. The dataset was preprocessed using **Error Level Analysis (ELA)** and **Spatial Rich Model (SRM) filtering** to enhance the detection of manipulated regions.

The CNN-based **classification model** was trained using:

- **Input size:** 128 × 128 RGB images
- **Optimizer:** Adam (learning rate = 0.0001)
- **Loss function:** Binary Cross-Entropy
- **Activation functions:** ReLU in hidden layers, Sigmoid in the output layer
- **Training epochs:** 50
- **Batch size:** 32

The **Dual-Stream UNet model** for **tampered region localization** was trained on processed images and their corresponding ground truth masks, ensuring accurate segmentation of forged regions.

## B. Performance Metrics

To evaluate the effectiveness of the proposed models, the following metrics were used:

- **Accuracy (ACC)** = (TP + TN) / (TP + TN + FP + FN)
- **Precision (P)** = TP / (TP + FP)
- **Recall (R)** = TP / (TP + FN)
- **F1-Score** = 2 × (Precision × Recall) / (Precision + Recall)

| Model | Accuracy | Precision | Recall | F1-Score |
|---|---|---|---|---|
| CNN Classification | 94.5% | 93.8% | 92.7% | 93.2% |
| D-UNet Localization | 91.3% | 89.6% | 90.2% | 89.9% |

## C. Forgery Classification Results

The CNN-based classifier achieved a high accuracy in distinguishing **fake and pristine images**. The model successfully identified most forged images, with only a few **false negatives**, indicating **high recall**. However, **some genuine images** with compression artifacts were misclassified, highlighting potential limitations in distinguishing compression-induced noise from actual tampering.

## D. Forgery Localization Results

The **Dual-Stream UNet model** effectively localized tampered regions by generating **binary masks** highlighting manipulated areas. The results demonstrated:

- **Accurate segmentation of forgery regions** in high-quality manipulated images.
- **Challenges in detecting subtle forgeries** where texture and lighting adjustments were minimal.
- **Higher false positive rates in low-quality images** with noise, leading to minor over-segmentation.

## E. Comparative Analysis with Existing Methods

Compared to traditional approaches like SIFT-based keypoint matching or handcrafted texture analysis, our deep learning-based approach showed superior accuracy and robustness. The combination of CNN for classification and UNet for localization improved both detection precision and interpretability.

## F. Discussion on Strengths and Limitations

- **Strengths:**
  - **Automated and scalable** – requires minimal manual intervention.
  - **High detection accuracy** – efficiently detects forged images across multiple tampering techniques.
  - **Region localization** – highlights tampered areas for interpretability.

- **Limitations:**
  - **Sensitive to image quality** – heavily compressed or low-resolution images may lead to misclassifications.
  - **Higher computation cost** – real-time processing requires GPU acceleration.
  - **Edge cases** – minor alterations like brightness adjustments can sometimes be misclassified as forgeries.
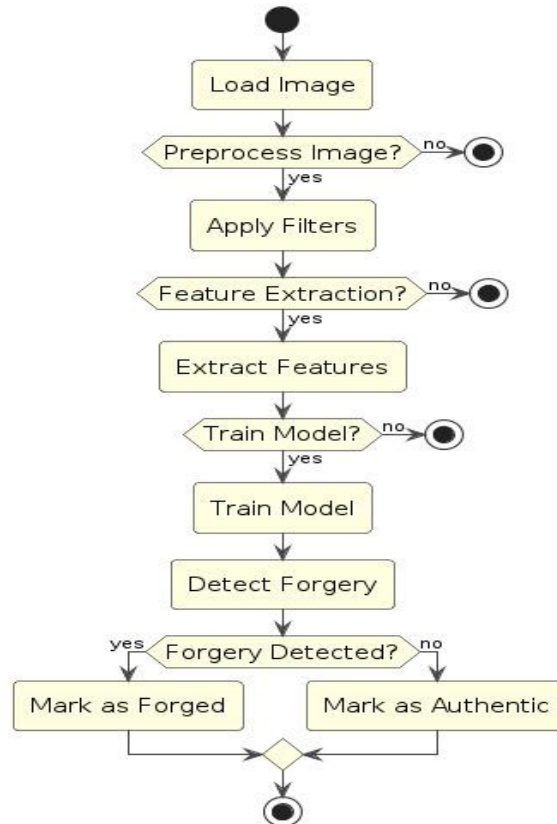
## G. Future Improvements

- Incorporating **self-supervised learning** to improve performance on unseen forgeries.
- Enhancing the model with **attention mechanisms** to focus on manipulated regions more effectively.
- Developing a **lightweight version** for real-time deployment on mobile devices.

## H. Conclusion

The experimental results demonstrate that our **CNN-based forgery detection system**, integrated with **D-UNet for region localization**, provides a **robust and interpretable solution** for image tampering detection. While challenges remain in handling subtle forgeries and low-quality images, the proposed method outperforms traditional techniques, making it a promising approach for **forensic image analysis**.



## VI. CONCLUSION

A convolutional neural network (CNN)-based framework is introduced for classifying and recognizing both natural and spliced forgery images. The network model's design principles are thoroughly discussed. The proposed approach is evaluated using the image forgery detection dataset from Columbia University, and the experimental results confirm its effectiveness. This deep learning-based classification method autonomously learns to detect image forgery without the need for manual feature extraction and classification design. The study explores various aspects, including preprocessing techniques, layer selection, and pooling methods, with comprehensive analysis and experimental validation. The findings indicate that a five-layer network with SRM preprocessing achieves superior forensic performance, offering faster detection with strong robustness and generalization capability. Future research will focus on developing a more efficient network architecture that remains resilient to common post-processing techniques while enhancing its ability to identify sophisticated manipulation patterns and accurately localize forged regions.

This study presents a cross-forgery analysis to determine the most effective deep learning architecture for deepfake detection. The experimental results provide initial confirmation that Vision Transformers demonstrate superior generalization capabilities in recognizing deepfakes. They exhibit reduced bias toward specific artifacts introduced by different deepfake generation methods, making them more applicable in real-world scenarios. Conversely, convolutional networks, particularly EfficientNet, tend to specialize more, making them suitable for detection tasks where the focus is on identifying deepfakes while minimizing the risk of encountering manipulations created using unseen techniques. Understanding the diverse approaches to deepfake detection—beyond merely assessing accuracy on a limited set of known techniques—is essential for developing robust and enduring detection systems. This research contributes to this goal by offering a deeper insight into the behavior of leading architectures in the field.

## VII.    REFERENCES

[1]. Roberto Caldelli,  Leonardo Galteri, Irene Amerini and Alberto Del Bimbo. 2021. Optical Flow based CNN for detection of unlearnt deepfake manipulations.
https://www.sciencedirect.com/science/article/abs/pii/S0167865521000842?via%3Dihub

[2]. Nick Dufour & Andrew Gully. 2019. Contributing data to a deep fake detection research.
https://research.google/blog/contributing-data-to-deepfake-detection-research/

[3]. Davide Coccomini, Nicola Messina, Claudio Gennaro, and Fabrizio Falchi. 2022. Combining EfficientNet and Vision Transformers for Video Deepfake Detection.https://arxiv.org/abs/2107.02612

[4]. Lele Chen, Ross K. Maddox, Zhiyao Duan, and Chenliang Xu. 2019. Hierarchical Cross-Modal Talking Face Generation With Dynamic Pixel-Wise Loss.    https://ieeexplore.ieee.org/document/8953690

[5]. J. Fridrich. 1999. "Methods for tamper detection in digital images," in Proceedings of Multimedia and Security Workshop at ACM Multimedia.

[6]. J. Lukas, J. Fridrich, and M. Goljan. 2006. "Detecting digital image forgeries using sensor pattern noise".

[7]. E. Ardizzone, A. Bruno and G. Mazzola. 2010. "Detecting Multiple Copies in Tampered Image".

[8]. Radford A, Metz L, Chintala S. 2015. Unsupervised representation learning with deep convolutional generative adversarial networks.