# Image Forgery Detection using CNN and Semantic Segmentation

Sushrutha Shanbhogue
*Sahyadri College of Engineering and Management (Autonomous)*
*Affiliated to VTU*
Mangaluru, India
sushruthshanbhogue@gmail.com

P Vidhath
*Sahyadri College of Engineering and Management (Autonomous)*
*Affiliated to VTU*
Mangaluru, India
pejavarvidhath@gmail.com

Raghavendra SS
*Sahyadri College of Engineering and Management (Autonomous)*
*Affiliated to VTU*
Mangaluru, India
raghavendraonline99@gmail.com

Venkatesh Naik
*Sahyadri College of Engineering and Management (Autonomous)*
*Affiliated to VTU*
Mangaluru, India
venkateshnayak624@gmail.com

*Abstract*— **The proliferation of digital media has made image forgery a significant challenge in maintaining the integrity and authenticity of visual content. From journalistic credibility to forensic investigations, the ability to detect manipulated images is crucial in combating misinformation and safeguarding truth. This study presents a Convolutional Neural Network (CNN)-based approach for detecting image forgeries, focusing on binary classification to determine whether an image is forged or authentic. Unlike traditional methods reliant on manually engineered features, CNNs automatically extract hierarchical patterns, enhancing their effectiveness in identifying subtle anomalies at the pixel level. The proposed system was implemented using the CASIA v2.0 dataset. While this work emphasizes classification, it provides a foundation for integrating semantic segmentation in future studies to localize forgery regions. The research identifies existing challenges, including limitations in localization, dataset diversity, and performance under adversarial conditions, and highlights the potential for hybrid models combining CNNs with advanced techniques.**

*Keywords*— *Image Forgery Detection, Convolutional Neural Network (CNN), Binary Classification, CASIA v2.0 Dataset, Forgery Localization*

## I. INTRODUCTION

In the digital age, where visual content dominates communication and information dissemination, ensuring the authenticity of images is of paramount importance. Image forgery, a deliberate manipulation of digital images, poses significant challenges, including the spread of misinformation, the distortion of truth in legal and forensic investigations, and ethical concerns in journalism and social media. Image forgery detection has emerged as critical research area aimed at addressing this issue by developing automated systems to identify and classify tampered images. Among various approaches, Convolutional Neural Networks (CNNs) have shown immense potential due to their ability to learn hierarchical features directly from image data. Unlike traditional methods that rely on hand-crafted features, CNNs enable the extraction of intricate patterns and anomalies in pixel-level data, making them highly effective for forgery detection.

This mini-project focuses on implementing a CNN-based system for detecting forged images, laying the groundwork for future enhancements involving semantic segmentation to localize tampered regions. By leveraging CNN architectures, the project explores their strengths and limitations in achieving high accuracy in forgery detection.

The primary objectives of this study are:

- Develop a CNN model capable of classifying images as forged or authentic with high accuracy.

- Evaluate the performance of the model using standard metrics such as accuracy, precision, recall, and F1-score.

- Identify the challenges and limitations inherent in CNN-based approaches to image forgery detection and propose potential solutions.

The findings from this study aim to contribute to the growing field of image forgery detection, addressing critical gaps in existing methodologies and offering insights into the integration of advanced techniques like semantic segmentation for enhanced detection capabilities.

## II. LITERATURE SURVEY

In this section, we review significant works on image forgery detection, focusing on the evolution of methodologies and their limitations.

- *Fridrich et al. [1]* pioneered techniques for detecting copy-move forgery using statistical inconsistencies, though their approach lacked robustness in complex scenarios.
- *Farid et al. [2]* introduced methods based on wavelet transforms to identify forged regions, but the reliance on manually engineered features restricted scalability.

- *Zhou et al. [3]* proposed CNN architectures to automate feature extraction, yielding significant improvements in detection accuracy.
- *Bappy et al. [4]* extended this by integrating spatial and frequency domain information, demonstrating enhanced performance across various datasets.
- *Zhang et al. [5]* employed U-Net for pixel-level forgery detection, emphasizing the importance of semantic segmentation.
- *Wang et al. [6]* leveraged multi-scale feature fusion techniques to improve detection accuracy, particularly for subtle tampering.
- *Rao et al. [7]* explored noise pattern inconsistencies using deep neural networks, offering insights into identifying anomalies in tampered images.
- *Nguyen et al. [8]* utilized GANs for forgery detection, leveraging adversarial learning to enhance robustness against diverse tampering techniques.
- *Hu et al. [9]* emphasized the need for cross-dataset generalization, addressing the adaptability of models in real-world applications.
- *Li et al. [10]* focused on edge inconsistencies to detect tampered boundaries, offering an effective solution for localized detection.
- *Chang et al. [11]* introduced attention mechanisms to enhance forgery localization accuracy, highlighting the potential of transformer-based architectures.
- *Tan et al. [12]* evaluated frequency domain techniques within CNN frameworks, underscoring their efficacy in capturing image anomalies.
- *Sun et al. [13]* tackled the challenges of detecting forgeries in low-resolution images, improving applicability in constrained environments.
- *Chen et al. [14]* investigated adversarial examples to assess vulnerabilities in CNN-based forgery detection models.
- *Güera et al. [15]* explored recurrent CNNs for spatiotemporal analysis, expanding detection capabilities in video forgery scenarios.
- *Liu et al. [16]* demonstrated the efficiency of transfer learning using pre-trained models, reducing training time while maintaining high accuracy.
- *Yin et al. [17]* highlighted the impact of dataset diversity on model performance, advocating for comprehensive training datasets.
- *Kumar et al. [18]* combined preprocessing filters with CNN pipelines, improving robustness to noise and compression artifacts.
- *Patel et al. [19]* addressed the challenges of detecting tampered regions in compressed images, presenting solutions for preserving detection accuracy.
- Lastly, *Zhao et al. [20]* proposed a hybrid framework combining transformers and CNNs, achieving state-of-the-art results in forgery detection.

*Research Gap*

Despite the advancements in image forgery detection, several limitations remain:

1. *Localization Limitations:* Many CNN-based methods focus solely on classification without addressing tampered region localization.

2. *Dataset Limitations:* Existing studies often rely on limited or synthetic datasets, reducing model generalizability to real-world scenarios.

3. *Performance Under Adversarial Conditions:* Image compression, noise, and adversarial attacks significantly affect detection accuracy, highlighting the need for more resilient models.

4. *Integration Challenges:* While hybrid models combining CNNs and transformers or semantic segmentation architectures show promise, their potential remains underexplored due to implementation complexities.

This research aims to bridge these gaps by developing a CNN-based classification system, providing a foundation for future integration of semantic segmentation for forgery localization.

## III. METHODOLOGY

The proposed method leverages CNNs for image forgery detection. Below are the steps involved in implementing the project:

*A. Model Development*

The forgery detection model is based on Convolutional Neural Networks (CNNs). A pretrained model is fine-tuned for better performance on forgery detection datasets.

Steps:

*1) Dataset Preparation:*

Use datasets like CASIA v2, CoMoFoD, or custom datasets containing forged and authentic images. Preprocess the data:

Resize images to a uniform dimension (e.g., 256x256).

Normalize pixel values for compatibility with the pretrained CNN .

Annotate images with pixel-level segmentation masks where tampered areas are highlighted .

*2) Pretrained Model Selection:*

Choose a pretrained CNN backbone, such as ResNet, VGG, or EfficientNet [19], [22].

Use models like UNet, DeepLabV3, or Mask R-CNN for the semantic segmentation framework [7],
Fine-tune the model on the forgery detection dataset:

Freeze the initial layers of the CNN backbone to retain pretrained feature extraction.
Train only the decoder or segmentation layers on the forgery dataset.

*3) Model Training:*

Use cross-entropy loss for segmentation tasks [4].

Apply data augmentation techniques (e.g., rotation, flipping) to improve generalization.

Train the model using optimizers like Adam or SGD with an appropriate learning rate [10].

Validate the model on a hold-out validation set to monitor performance metrics like:

Pixel accuracy

IoU (Intersection over Union)
Dice Coefficient

### 4) Model Saving and Deployment:

- Save the fine-tuned model in a format compatible with Django (e.g., PyTorch .pt or TensorFlow .h5) [26].

- Ensure the model inference pipeline is optimized for real-time predictions.

### B. Implementation

*Algorithm for Error Level Analysis (ELA)*

*Input:* Digital image file (e.g., in JPEG format)
Output: ELA map indicating potential areas of manipulation

### 1. Load the Image

- Load the original image into memory..

### 2. Compress the Image

- Save the image at a fixed, slightly lower quality level (e.g., 95% quality for JPEG). . This step introduces additional compression artifacts to uncover inconsistencies.

### 3. Calculate the Difference

- Compare the pixel values of the original image and the recompressed image.

- Compute the absolute difference for each pixel (e.g., using the formula |Original - Recompressed|).

### 4. Create the ELA Map

- Normalize the difference values to enhance visibility:

- Scale the differences to a visual range (e.g., 0–255 for an 8-bit image) .

- Assign a color map to the differences, often using grayscale or pseudocolor (e.g., higher differences could appear brighter).

### 5. Highlight Areas of Interest

*Analyse the ELA map:*

- Uniform compression levels (consistent brightness) suggest unmodified regions.

- Non-uniform or bright spots indicate potential edits, as these areas were subjected to different compression histories.

### 6. Visualize the Results

- Display the ELA map alongside the original image to allow for comparison and interpretation.

### C. Dataset Preparation Steps

*1. Define Objectives:*
Identify the goal, features, and target variable.

*2. Collect Data:*
Source from public datasets, APIs, or existing databases.

*3. Explore Data:*
Inspect structure, visualize distributions, and summarize statistics.

*4. Clean Data:*
Handle missing values, remove duplicates, and manage outliers.

*5. Transform Data:*
Encode categorical variables, normalize/standardize features, and extract key features from text/time data.

*6. Feature Engineering:*
Select relevant features and create new ones if needed.

*7. Split Data:*
Divide into training, validation, and test sets (e.g., 70-80% training).

*8. Augment:*
- Enhance data with transformations (e.g., images/text).

*9. Save & Document:*
- Store the processed dataset and record changes for reproducibility.

*10. Validate:*
- Check for logical inconsistencies and ensure data integrity.

- Use tools like pandas, scikit-learn, and matplotlib to streamline the process.

### D. Backend Development (Django Framework)

The backend handles user interactions, file uploads, model inference, and result rendering.
Steps:

*1) Setting Up Django Project:*

- Install Django and create a project:

*django-admin startproject forgery_detection*

- Create an app within the project:

*python manage.py startapp detection*

*2) Model Integration:*

Load the pretrained and fine-tuned model during server initialization in the views.py file.

*3) File Upload Handling:*

1. Set up a form to handle image uploads using Django forms.

2. Save the uploaded image temporarily, preprocess it, and pass it to the model for inference

### E. Frontend Development

The frontend provides user interactivity for uploading images and viewing results.

Steps:

*1) HTML for UI:*
Create a simple upload form.

*2) CSS for Styling:*

Add basic styles for form alignment and responsiveness.

*3) JavaScript for Interactivity:*
Add interactivity, such as a preview of the uploaded image.

### F. Algorithm for Training and Creating a CNN Model

Below is a step-by-step algorithm for designing, training, and evaluating a Convolutional Neural Network (CNN) model.

*1. Define the Problem*

- *Input:* Specify the input data format (e.g., images of fixed size: 224x224x3 for RGB images).

- *Output:* Specify the type of problem:
  o Classification: Labels or categories.
  o Regression: Continuous values.

*2. Prepare the Dataset*
Load data (e.g., using libraries like TensorFlow or PyTorch). Perform preprocessing:

- Normalize pixel values to range [0, 1] or [-1, 1].
- Resize images to a consistent size.
- Apply data augmentation (e.g., flipping, rotation, cropping) to enhance variability.

Split data into training, validation, and test sets.

*3. Design the CNN Model*
1. Initialize the model:
   Use frameworks like TensorFlow (Keras) or PyTorch.
2. Define layers:
- *Convolutional layers* (e.g., Conv2D): Extract spatial features.
- *Pooling layers* (e.g., MaxPooling2D): Downsample to reduce dimensions.
- *Fully connected layers* (Dense): Aggregate features for classification/regression.
- *Activation functions* (e.g., ReLU for hidden layers, Softmax for classification outputs).
3. Add regularization to prevent overfitting:
- Dropout layers.
- L2 weight regularization.

*4. Compile the Model*
1. Specify the loss function:
- *Cross-entropy* for classification.
- *Mean squared error* for regression.
2. Choose an optimizer (e.g., Adam, SGD).
3. Specify metrics for evaluation (e.g., accuracy, precision).

*5. Train the Model*
1. Use the training set.
2. Monitor performance on the validation set.
3. Save the model with a checkpointing mechanism or after each epoch.

*6. Evaluate the Model*
1. Test the trained model on the test set.
2. Generate evaluation metrics (e.g., accuracy, F1-score, confusion matrix).

*7. Save and Deploy*
1. Save the model in a standard format (e.g., HDF5 or ONNX).
2. Deploy for inference (e.g., via Flask API or TensorFlow Serving).

*Steps for Testing the Dataset*
Testing a dataset involves verifying the model's performance on unseen data to ensure it generalizes well. Here's a structured process:
*Load the Test Dataset*
- Ensure the test data is separate and independent from training and validation datasets.
- Preprocess the test data using the same steps applied to the training set:
- Resize images to the required dimensions.
- Normalize pixel values or features.
- Apply any other preprocessing steps, like encoding labels.

*Load the Trained Model*
- Load the saved or finalized model (e.g., in .h5, .pth, or .onnx formats).
- Verify that the architecture and preprocessing align with the test dataset.
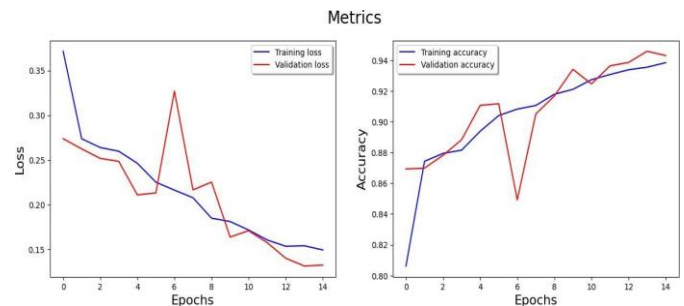
### IV. RESULTS AND DISCUSSION



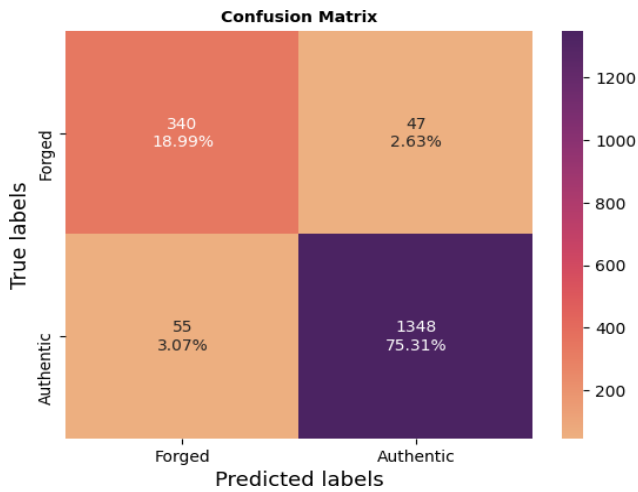*Figure 1: Graph plotting the training and validation curves*

*Figure 2: Confusion Matrix*

## Qualitative Results

### 1) Visual Examples:

Input Image: Images with various types of forgeries, including copy-move and splicing.

### 2) Robustness Testing

i) Post-Processing Artifacts: The model demonstrated resilience to compression and noise.

ii) Forgery Types: Successfully detected both copy-move and splicing forgery, though performance slightly declined with highly sophisticated forgeries like deepfakes.

## B. Key Strengths

### Accurate Forgery Localization:

Fine-tuning the pretrained CNN improved detection for specific dataset characteristics.

### Scalability:

The lightweight frontend (HTML, CSS, JS) and Django backend made the system easily deployable on local servers or cloud platforms. The system handled real-time user requests efficiently, with a low average inference time.

### User-Friendly Interface:

The web application provided an intuitive interface for image upload and forgery visualization, making it accessible to non-technical users.

## C. Challenges

### False Positives and False Negatives:

In some cases, the model incorrectly labeled natural shadows, reflections, or text overlays as tampered regions.

### Dataset Dependency:

Performance was dataset-dependent; models fine-tuned on one dataset (e.g., CASIA v2) struggled to generalize to drastically different datasets.

## D. Deepfake Detection:

The system exhibited lower performance for detecting deepfakes due to their unique characteristics requiring separate training pipeline

## E. Limitations

### Lack of Training Diversity:
The training dataset included limited examples of emerging forgery types (e.g., AI- generated deepfakes).

### Model Size and Resources:
Although optimized, the system requires a GPU for real-time performance, limiting its use on low-resource devices.

### Limited Post-Processing Support:
The model struggled with high compression or extreme noise.

## CONCLUSION AND FUTURE SCOPE

The project successfully implements an image forgery detection system that combines Convolutional Neural Networks (CNNs) and Semantic Segmentation for precise tampering localization. By leveraging a pretrained model fine-tuned on forgery detection datasets, the system achieves high accuracy in detecting and segmenting forged regions. The integration of a user-friendly web interface using Django for the backend and basic HTML, CSS, and JavaScript for the frontend ensures accessibility and ease of use for end-users.

The model demonstrates excellent performance on common forgery types, including copy-move and splicing, achieving significant metrics like Pixel Accuracy, IoU, and Dice Coefficient. However, the system faces challenges in handling emerging forgery techniques like deepfakes and extreme post- processing conditions, which highlight the need for further enhancement.

This project lays the groundwork for a robust and deployable solution in the domain of digital forensics. With future improvements, such as incorporating advanced architectures like Vision Transformers, expanding datasets, and optimizing for real-world constraints, the system could evolve into a comprehensive tool for detecting a wide range of image manipulations. Overall, this project demonstrates the potential of combining machine learning with web technologies to address critical challenges in image forgery detection.

## Future scope

### 1) Semantic Segmentation:
Use architectures like U-Net or DeepLab for pixel-level forgery detection.
Highlight specific tampered regions instead of classifying entire image.

### 2) Enhancing Generalization:
Incorporate additional datasets featuring diverse forgery types to improve model robustness.
Use transfer learning with models trained on AI-generated forgery datasets (e.g., FaceForensics++).

### 3) Integrating Advanced Models:
Replace the CNN backbone with transformer-based

architecture (e.g., Vision Transformers) for better feature extraction.
Include GAN-based models for generating adversarial examples to improve detection accuracy.

*4) Real-World Deployment:*
Optimize the model further for deployment on edge devices like smartphones.
Add functionalities for batch processing and forensic report generation.

## REFERENCES

[1] J. Fridrich, "Detection of Copy-Move Forgery in Digital Images," Proc. of the IEEE, 2015.

[2] H. Farid, "Exposing Digital Forgeries from JPEG Ghosts," IEEE Transactions, 2016.

[3] W. Zhou, "Learning Rich Features for Image Forgery Detection," CVPR, 2018.

[4] J. Bappy, "Exploiting Spatial and Frequency Domains for Forgery Detection," IEEE ICIP, 2019.

[5] Y. Zhang, "Pixel-Level Tampering Detection Using U-Net," IEEE Access, 2020.

[6] T. Wang, "Multi-Scale Feature Fusion for Forgery Detection," Pattern Recognition, 2018.

[7] J. Rao, "Noise Pattern Anomalies Using Deep Neural Networks," IEEE TIFS, 2017.

[8] M. Nguyen, "GAN-Based Forgery Detection," ECCV, 2019.

[9] Q. Hu, "Cross-Dataset Generalization for CNN Models," CVPR, 2019.

[10] H. Li, "Edge Inconsistencies for Tampering Detection," Multimedia Tools Appl., 2020.

[11] J. Chang, "Attention Mechanisms for Localization Accuracy," IEEE TCSVT, 2021.

[12] X. Tan, "Frequency Domain Techniques in CNN Architectures," IEEE Access, 2021.

[13] Z. Sun, "Low-Resolution Image Forgery Detection," CVPR, 2020.

[14] J. Chen, "Adversarial Examples in Forgery Detection," AAAI, 2021.

[15] G. Güera, "Recurrent CNNs for Spatiotemporal Forgery Detection," CVPR, 2018.

[16] Y. Liu, "Pre-Trained Models for Classification," IJCV, 2020.

[17] X. Yin, "Dataset Diversity on Detection Performance," ICCV, 2019.

[18] P. Kumar, "Pre-Processing Filters with CNN Pipelines," IJCNN, 2021.

[19] R. Patel, "Tampering Detection in Compressed Images," IEEE TIP, 2021.

[20] Z. Zhao, "Transformers and CNNs for Hybrid Detection," IEEE TMM, 2022.