# IDENTIFYING THE AUTHENTICITY OF IMAGE USING CONVOLUTION NEURAL NETWORK

**S. Simhadri** UG Students, Department of ECE, Aditya College of Engineering and Technology, Surampalem, Andhra Pradesh, India.
**B.C.Vinay SriRam** UG Students, Department of ECE, Aditya College of Engineering and Technology, Surampalem, Andhra Pradesh, India.
**Md Abdul Kareem** UG Students, Department of ECE, Aditya College of Engineering and Technology, Surampalem, Andhra Pradesh, India.
**Y. Lokesh GopiRaju** UG Students, Department of ECE, Aditya College of Engineering and Technology, Surampalem, Andhra Pradesh, India.
**Dr. D. KISHORE** Professor & Dean Evaluation, Department of ECE, Aditya College of Engineering and Technology, Surampalem, Andhra Pradesh, India.

**Abstract**
        This research presents a novel approach utilizing Convolutional Neural Networks (CNNs) for the precise classification of original and tampered images. The dataset, meticulously curated, is partitioned into distinct training and testing sets to ensure robust model evaluation. The architecture of the CNN is meticulously crafted, comprising a series of layers strategically designed for optimal feature extraction and classification. These layers include convolutional, batch normalization, rectified linear unit (ReLU), addition, average pooling, fully connected, softmax, and classification layers. Furthermore, the integration of skip connections enhances information propagation within the network, thereby augmenting its discriminative capacity.To train the model effectively, the Adam optimizer, coupled with GPU acceleration, is employed, ensuring efficient convergence and reduced training times. Subsequently, a comprehensive suite of performance metrics including accuracy, precision, sensitivity, and specificity is meticulously computed and presented. Visual aids such as confusion matrices and heatmaps are deployed to offer intuitive insights into the classification outcomes, aiding in the interpretation and validation of the model's efficacy. The experimental evaluation conducted underscores the effectiveness of the proposed CNN model in accurately discerning between original and tampered images. Through meticulous analysis and interpretation of the results, it becomes evident that the model exhibits superior performance, showcasing its potential as a robust solution for image forgery detection and classification.

**Keywords:**
Convolutional Neural Network (CNN), Image Classification, Tampered Image Detection, Performance Metrics, Skip Connections.

## I. Introduction

   The proliferation of image editing tools and platforms has led to an increase in the prevalence of picture tampering, which includes manipulation, fabrication, and alteration. Ensuring the integrity of visual data, content authentication, forensic analysis, and other uses all depend on the ability to identify manipulated images. In this work, we suggest using a Convolutional Neural Network (CNN) model to distinguish between authentic and altered images.

        Every day, billions of images are created, exchanged, or transferred thanks to the cutting-edge computerized world and sophisticated innovation strategy. It is incredibly simple for anyone to alter an image because to the abundance of photo editing programs, such as Adobe and GIMP, and their user-friendly interfaces. Every day, millions of fake photos are uploaded; concealing information and being used as spurious evidence. Since photos are utilized as visual evidence in courtrooms, manipulated photographs, or changed evidence, may lead to the punishment of an innocent party. In

the same way, manipulated photos that circulate on social media sites disseminate false information and news.

Hence can harm an individual, can even result in rivalry and conflict among societies with diverse cultures. Recently various false beliefs about the covid19 pandemic were spread in the form of fake images, which created misbeliefs amongst citizens [6].

For a variety of applications, such as digital forensics, law enforcement, and journalism, the capacity to automatically distinguish between authentic and altered photographs is crucial. Conventional systems for detecting image tampering frequently depend on manually created characteristics and heuristics, which could not be resilient to changes in the content of the image or methods of manipulation. Deep learning-based methods, on the other hand, especially CNNs, have demonstrated impressive results in image classification tasks by automatically extracting discriminative features from unprocessed pixel data. The ease of access to digital content and the widespread availability of picture altering tools have resulted in a rise in instances of image tampering, which presents serious obstacles to the integrity and authenticity of visual data. Conventional techniques for identifying altered photos frequently depend on human examination or manually created characteristics, which may not be scalable. Therefore, there is a critical need for automated techniques capable of accurately identifying manipulated images and distinguishing them from authentic ones.

The aim of this research is to create a Convolutional Neural Network (CNN) model that is both reliable and effective in automatically classifying manipulated and original photos. The suggested CNN model seeks to accurately detect several kinds of image modifications, such as splicing, copy-move, and retouching, by using deep learning techniques to extract discriminative features directly from raw pixel data. Furthermore, utilizing benchmark datasets and common assessment criteria including accuracy, precision, sensitivity, and specificity, the study seeks to assess the CNN model's performance. By accomplishing these goals, the study hopes to improve the field of picture forensics and content authentication technologies, making it easier to recognize and counteract dangers of image tampering in digital media.

The suggested CNN model extracts spatial characteristics at various scales by taking advantage of the convolutional layers' innate hierarchical structure. Layers of batch normalization are used to enhance generalization and speed up training. The model can learn intricate patterns thanks to the non-linearity introduced by Rectified Linear Unit (ReLU) activation functions. In order to promote gradient flow and mitigate the vanishing gradient issue, skip connections are also included.

The training and assessment dataset is a varied group of images that includes both unaltered and altered examples. The dataset is divided into training and testing sets in order to evaluate the model's generalization performance. The Adam optimizer is used to train the CNN model, while mini-batch stochastic gradient descent and GPU acceleration are used to speed up convergence. The suggested CNN model's performance is assessed using a range of criteria, such as sensitivity, specificity, accuracy, and precision. Heatmaps and confusion matrices are used to show the classification findings and evaluate how well the model can distinguish between authentic and altered photos. In conclusion, the goal of this research is to use deep learning techniques to develop an automated approach for detecting image tampering that is both trustworthy and effective. The proposed CNN model demonstrates promising performance in distinguishing between original and manipulated images, thereby contributing to the advancement of digital image forensics and content authentication technologies.

## II. Literature survey

The area of multimedia forensics has advanced significantly in recent years as a result of researchers' concentrated efforts on creating reliable techniques for picture and video analysis. Diverse methodologies, encompassing deep learning-centric strategies and conventional feature extraction techniques, have been investigated for the identification and categorization of digital counterfeit materials. An overview of pertinent literature is given in this section, emphasizing important

publications on topics including biometric identification, steg analysis, copy-move forgery detection, and image splicing detection. The objective of this study is to identify research gaps and suggest innovative methods for improving the authenticity and integrity of digital media by thoroughly reviewing current methodologies.

A hybrid deep learning system for identifying picture splicing and copy-move forgeries is presented by Zhang et al. (2019) et al. [1]. To increase the accuracy of forgery detection, it blends classic forensic procedures with deep learning approaches. The authors suggest a brand-new deep learning-based feature extraction method, which is followed by a conventional forensics classification phase. The framework detects many kinds of image forgeries with promising results. Fridrich et al. [2] and associates concentrate on the identification of copy-move fraud in digital photographs in this work. They suggest a technique to find duplicated areas in an image that is based on the examination of local picture attributes like key points and descriptors. Through a comparative analysis of these characteristics, copy-move forgeries can be identified, offering important information about the legitimacy of the image.

A steganalysis method based on the subtractive pixel adjacency matrix is introduced by Pevní, Bas, Fridrich et al. [3]. This methodology seeks to identify concealed data, such modifications or messages, in digital pictures. The technique advances the field of image forensics by effectively differentiating between original and stegano graphic images by examining the statistical characteristics of pixel adjacencies. An approach to deep learning for universal picture modification detection is proposed by Bayar and Stamm et al. [4]. They provide a brand-new convolutional layer that is intended to record the effects of image alteration, including filtering, retouching, and splicing. The technique delivers strong performance in detecting different kinds of image forgeries by training deep neural networks on a broad dataset of modified photos.

Dong, J., Yu, S., & Huang et al. [5] concentrate on revealing image splicing by examining inconsistent local noise variance. They suggest a technique based on the finding that legitimate portions within an image typically have distinct noise properties than spliced regions. The technique helps detect picture alteration by efficiently identifying spliced regions through the analysis of local noise variance patterns. A deep learning-based method for picture splicing detection using convolutional and recurrent neural networks is presented by Al-Qararah, Khelifi, Bouridane et al. [6]. Their technique uses the temporal and spatial information present in photos to identify splicing artefacts. The approach delivers state-of-the-art performance in image splicing detection by training neural networks on a large dataset of altered photographs, demonstrating the usefulness of deep learning techniques in forensic applications.

Toh, Khairudin, Abdullah, and colleagues [7] introduce a deep learning method for identifying face retouching in portraiture. Their approach makes use of convolutional neural networks (CNNs) to examine facial features and identify discrepancies brought about by retouching methods like reshaping or smoothing. The technique makes a contribution to the field of picture forensics by accurately detecting retouching artefacts by training the CNN on a dataset of authentic and altered portraits. A technique for identifying digital forgeries by identifying illumination irregularities is presented by Johnson and Farid et al. [8]. The technique looks at changes in illumination within a picture to find areas that might have been added or altered. The technique helps with the verification of digital photographs by detecting anomalies that are suggestive of fraud by comparing the lighting gradients and shadows throughout the image.

A forensic technique for identifying digital forgeries in colour filter array interpolated (CFA) images is presented by Popescu and Farid et al. [9]. The technique looks into the statistical characteristics of CFA images to find discrepancies that have been caused by tampering or modification. The technique improves the integrity of digital photography by detecting forged sections and restoring the original image by looking at correlations between colour channels and spatial patterns.

A technique for identifying copy-move assaults in digital images based on Scale-Invariant Feature Transform (SIFT) characteristics is presented by Amerini, Ballan, Caldelli, and Del Bimbo et al. [10]. Through the extraction of SIFT key points and descriptors from the image, duplicate regions generated by copy-move forgeries can be detected. Key points that match and transformation parameters that are recovered allow the approach to precisely identify and

In their study, Ng, Chang, Sun, Chen, and colleagues [11] introduce a data-driven method for extracting forensic features in picture categorization. The strategy is centred on using data-driven methods and statistical analysis to extract discriminative characteristics from photos. The method contributes to the field of multimedia forensics and image classification by efficiently classifying images into preset categories by analyzing the distribution of characteristics across different image classes. A content-based image classification technique using Thepade's Static and Dynamic Ternary Block Truncation Coding (TBTC) is proposed by Thepade, Das, and Ghosh et al. [12]. This approach uses both static and dynamic ternary representations for classification, and it extracts features from images based on block truncation coding. Through the examination of the spatial arrangement of pixel values within image blocks, the technique produces precise image classification.

A method for content-based video retrieval utilizing multi-level Thepade's sorted ternary Block Truncation Coding (BTC) is presented by Thepade, Subhedarpage, and Mali et al. [13]. By taking into account intermediate block videos and even-odd films, this technique expands Thepade's BTC to video data. The method contributes to the field of multimedia information retrieval by efficiently achieving video retrieval based on content similarity through the application of ternary coding and multi-level sorting algorithms. A score-level fusion method for bimodal biometric identification utilizing Thepade's Sorted n-ary Block Truncation Coding (SBTC) was presented by Madane and Thepade et al. [14]. By combining matching scores from separate biometric modalities, this technique integrates palm print and iris features for biometric identification. Bimodal biometric identification performance is enhanced by the method by utilising Thepade's SBTC for feature extraction and fusion.

These papers collectively contribute to various aspects of multimedia forensics, image classification, content-based retrieval, and biometric identification by proposing innovative methods and techniques for feature extraction, representation, and classification.

## III.    Proposed Methodology

The proposed method for image classification using a Convolutional Neural Network (CNN) can be explained through a block diagram.

**Input Image:**

The input image serves as the raw data fed into the CNN model for classification. It can be any image in the dataset, either original or tampered.

**Preprocessing:**

Preprocessing involves standardizing the input image to ensure consistency and facilitate efficient processing. Common preprocessing steps include resizing the image to a fixed size, normalizing pixel values, and augmenting the dataset to increase variability and robustness
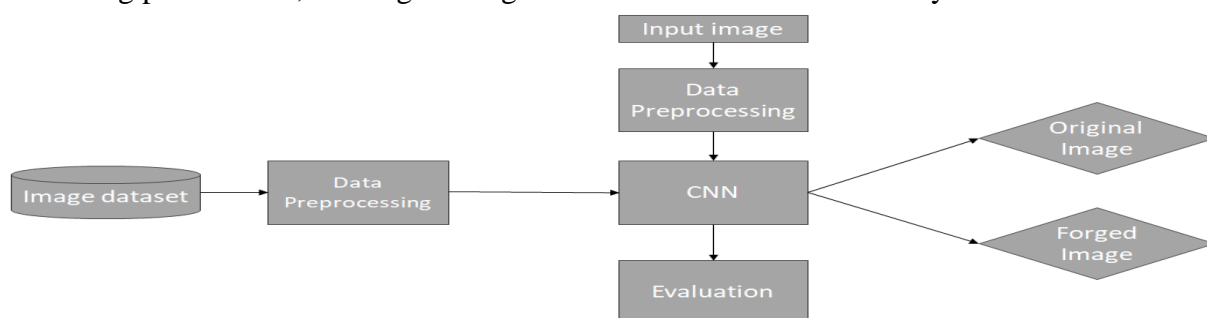


Fig.1 Block diagram

**Classification Output:**

The classification output block represents the final output of the CNN model, which consists of predicted class labels corresponding to the input image. The class with the highest probability score is selected as the predicted class label.
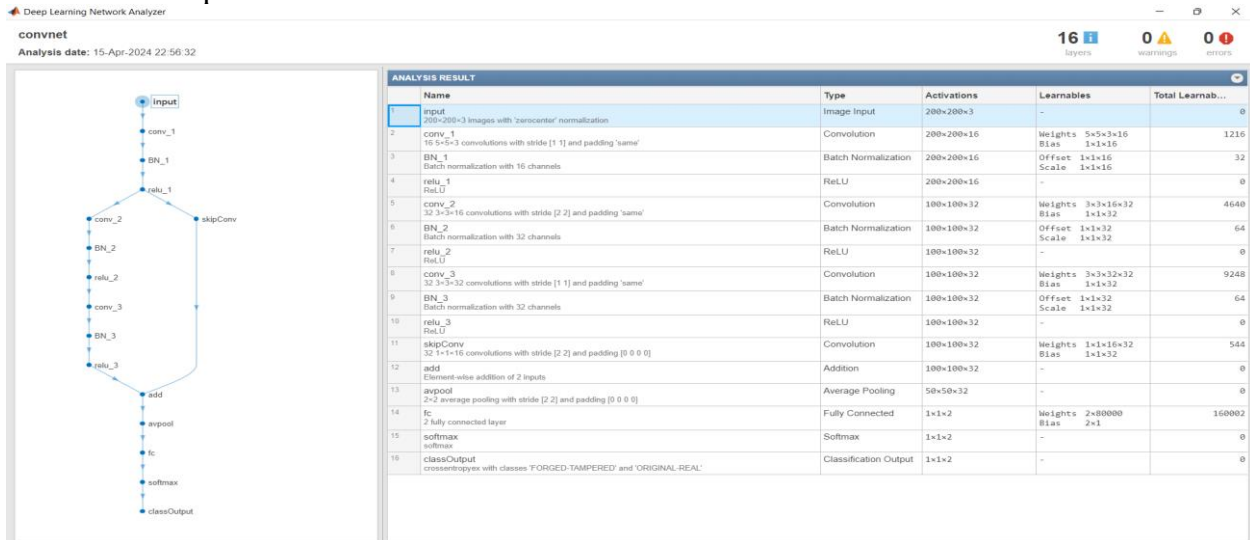


| | Name | Type | Activations | Learnables | Total Learnab... |
|---|---|---|---|---|---|
| 1 | input<br>200×200×3 images with 'zerocenter' normalization | Image Input | 200×200×3 | - | 0 |
| 2 | conv_1<br>16 5×5×3 convolutions with stride [1 1] and padding 'same' | Convolution | 200×200×16 | Weights 5×5×3×16<br>Bias 1×1×16 | 1216 |
| 3 | BN_1<br>Batch normalization with 16 channels | Batch Normalization | 200×200×16 | Offset 1×1×16<br>Scale 1×1×16 | 32 |
| 4 | relu_1<br>ReLU | ReLU | 200×200×16 | - | 0 |
| 5 | conv_2<br>32 3×3×16 convolutions with stride [2 2] and padding 'same' | Convolution | 100×100×32 | Weights 3×3×16×32<br>Bias 1×1×32 | 4640 |
| 6 | BN_2<br>Batch normalization with 32 channels | Batch Normalization | 100×100×32 | Offset 1×1×32<br>Scale 1×1×32 | 64 |
| 7 | relu_2<br>ReLU | ReLU | 100×100×32 | - | 0 |
| 8 | conv_3<br>32 3×3×32 convolutions with stride [1 1] and padding 'same' | Convolution | 100×100×32 | Weights 3×3×32×32<br>Bias 1×1×32 | 9248 |
| 9 | BN_3<br>Batch normalization with 32 channels | Batch Normalization | 100×100×32 | Offset 1×1×32<br>Scale 1×1×32 | 64 |
| 10 | relu_3<br>ReLU | ReLU | 100×100×32 | - | 0 |
| 11 | skipConv<br>32 1×1×16 convolutions with stride [2 2] and padding [0 0 0 0] | Convolution | 100×100×32 | Weights 1×1×16×32<br>Bias 1×1×32 | 544 |
| 12 | add<br>Element-wise addition of 2 inputs | Addition | 100×100×32 | - | 0 |
| 13 | avpool<br>2×2 average pooling with stride [2 2] and padding [0 0 0 0] | Average Pooling | 50×50×32 | - | 0 |
| 14 | fc<br>2 fully connected layer | Fully Connected | 1×1×2 | Weights 2×80000<br>Bias 2×1 | 160002 |
| 15 | softmax<br>softmax | Softmax | 1×1×2 | - | 0 |
| 16 | classOutput<br>crossentropyex with classes 'FORGED-TAMPERED' and 'ORIGINAL-REAL' | Classification Output | 1×1×2 | - | 0 |

Fig2. Architecture of CNN

**Proposed CNN Architecture**

Strong deep learning architectures like the Convolutional Neural Network (CNN) are made especially for processing and interpreting visual data like images. The structure of the animal visual cortex served as the model for the CNN, which consists of several layers that abstract and extract features in a hierarchical manner. Convolution, pooling, and fully connected layers are three essential elements in the CNN process that help the model build hierarchical representations of the input data.

**Input Layer:**

The input layer of the CNN receives the raw pixel values of the input image. Each pixel represents the intensity or color information at a specific location in the image. The input layer is typically represented as a grid of neurons, with each neuron corresponding to a single pixel in the image. For color images, the input layer has multiple channels corresponding to different color channels (e.g., red, green, blue).

**Convolutional Layers**:

The convolutional layers are the core building blocks of the CNN model. These layers consist of filters that convolve across the input image to extract features at different spatial hierarchies. Each convolutional layer applies a set of learnable filters to the input image, producing feature maps that capture increasingly abstract representations of the input data.

**Batch Normalization:**

Batch normalization layers normalize the activations of each convolutional layer across mini-batches during training. This technique helps stabilize the training process by reducing internal covariate shift and accelerating convergence.

**Activation Function (ReLU):**

Rectified Linear Unit (ReLU) activation functions introduce non-linearity to the CNN model, enabling it to learn complex patterns and representations. ReLU activation functions are applied element-wise to the feature maps generated by the convolutional layers.

**Activation Function (ReLU):**

After the convolution operation, an activation function is applied element-wise to the output feature map to introduce non-linearity to the CNN model. The Rectified Linear Unit (ReLU) activation function is commonly used due to its simplicity and effectiveness. ReLU sets all negative values to

zero, effectively introducing sparsity and enabling the model to learn complex non-linear relationships in the data.

**Pooling Layers:**

Pooling layers down sample the feature maps produced by the convolutional layers, reducing their spatial dimensions while retaining the most relevant information. Average pooling or max pooling operations are commonly used to achieve spatial aggregation and dimensionality reduction.

**Fully Connected Layers:**

Fully connected layers process the flattened feature vectors from the preceding layers and map them to the output classes. These layers incorporate learnable parameters that enable the CNN model to perform classification based on the extracted features.

**Softmax Activation:**

The softmax activation function is applied to the output of the fully connected layer to compute the probability distribution over the output classes. It normalizes the output scores into probabilities, ensuring that they sum up to one and represent the confidence of the model in each class.

**Loss Function and Optimization:**

During training, the CNN model minimizes a loss function, such as categorical cross-entropy, to update its parameters and improve performance. Optimization techniques, such as stochastic gradient descent (SGD) or Adam, are used to adjust the model's parameters iteratively based on the computed loss.

**Training and Evaluation:**

Using gradient descent optimization and back propagation, the CNN model is trained on a labeled dataset. In order to minimize the loss, batches of training data are iteratively fed through the network, the loss is calculated, and the model parameters are updated. To gauge the model's generalization performance, it is tested on a different validation or test dataset after training. Evaluation metrics are calculated to assess the model's performance on unobserved data, including accuracy, precision, recall, and F1 score.

Convolutional, pooling, and fully connected layers enable CNNs to achieve state-of-the-art performance on a wide range of image recognition benchmarks, demonstrating their significance in the field of deep learning and artificial intelligence. All things considered, the CNN architecture is a versatile and effective deep learning model for image classification, object detection, and various other computer vision tasks. Its hierarchical structure and ability to learn hierarchical representations m

**Evaluation Metrics:**

After training, the CNN model is evaluated using various performance metrics, including accuracy, precision, sensitivity, and specificity. These metrics assess the model's ability to correctly classify original and tampered images and provide insights into its overall performance.

**Dataset preparation**

Data are gathered into datasets. Photographs make up the Convolutional Neural Network dataset. This study uses the publicly available C0 benchmark dataset. This dataset contains 730 images. Ten of the pictures are real, and the other one hundred and ten are phonies. Two categories comprise the information:

- Genuine/Unique
- Illustrations of altered/phony original artwork.

The suggested CNN model advances image forensics and content authentication technology by combining these elements into a unified framework that allows it to distinguish between authentic and altered images.

## IV. Result Analysis

The model's performance is assessed by displaying the accuracy, precision, recall, and F1 score. The categorization results are plotted in a confusion matrix layout for easier to understand. Using

MATLAB R2021a, this pipeline efficiently enables you to train and assess a CNN-based image classifier. It works well when picture data is arranged into folders according to classes.

Confusion matrices and heat maps are used to visualize the classification results after the CNN model has been used to classify original and tampered photos. The performance of the trained model is assessed using a variety of criteria. Plots and statistics produced during the CNN model's classification of original photos and altered images offer important insights into the training process, model performance, and result analysis.

**Layer Graph Plot:**

The layer graph plot shows how the many layers are arranged and connected, providing a visual representation of the CNN model's architecture. It offers a high-level overview of all the layers in the model, including the fully connected, input, convolutional, and pooling layers. Understanding the information flow across the network and the feature transformation between layers is made easier with the help of this visualization.

**Skip Connection Plot:**

The inclusion of a 1x1 convolutional layer as a skip connection in the CNN model is depicted in the skip connection plot. In order to alleviate the vanishing gradient problem and improve gradient flow during training, skip connections are implemented. The inclusion of skip connections in the model architecture, which improves training stability and convergence, is shown in this plot.
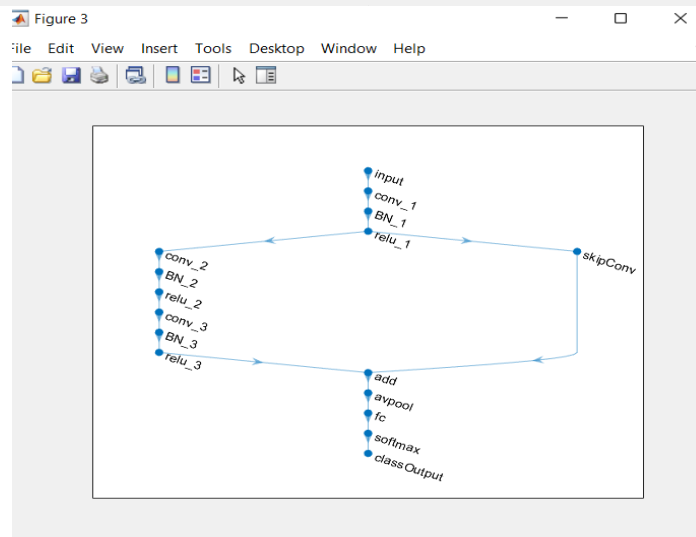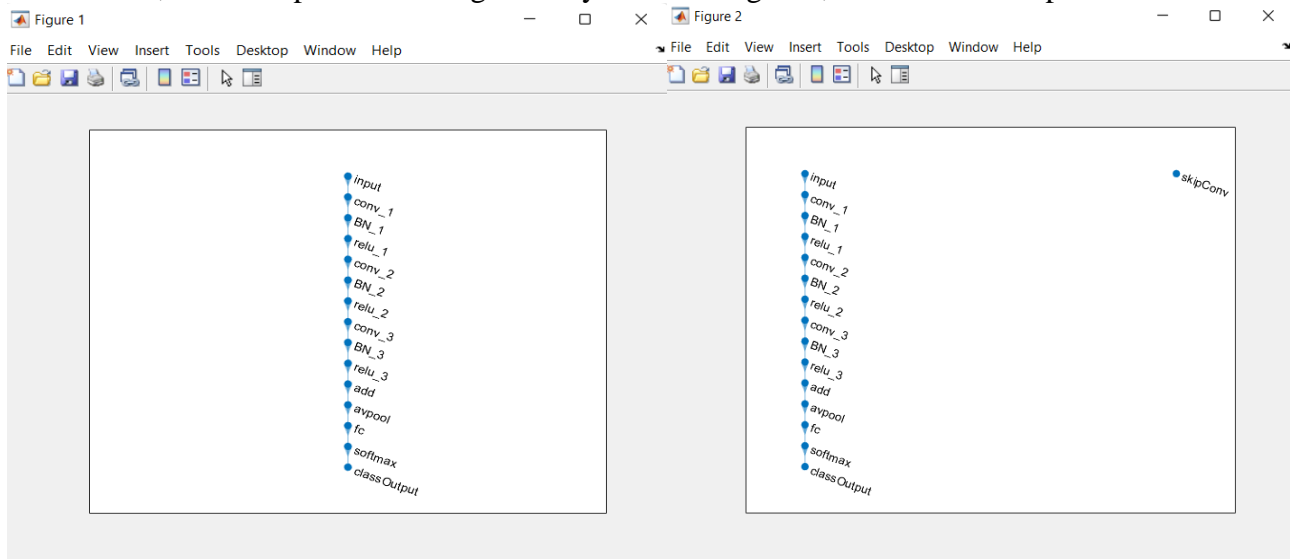


Fig3. Layer Graphs of CNN

**Training Progress Plot:**

The training progress graphic illustrates how training changes across several epochs and how training and validation losses change with time. It offers perceptions into the model's convergence behavior and aids in keeping an eye on training dynamics, including over fitting or under fitting. One can ascertain whether the model is learning efficiently and whether more training parameter adjustments are required by looking at the trend of loss curves.
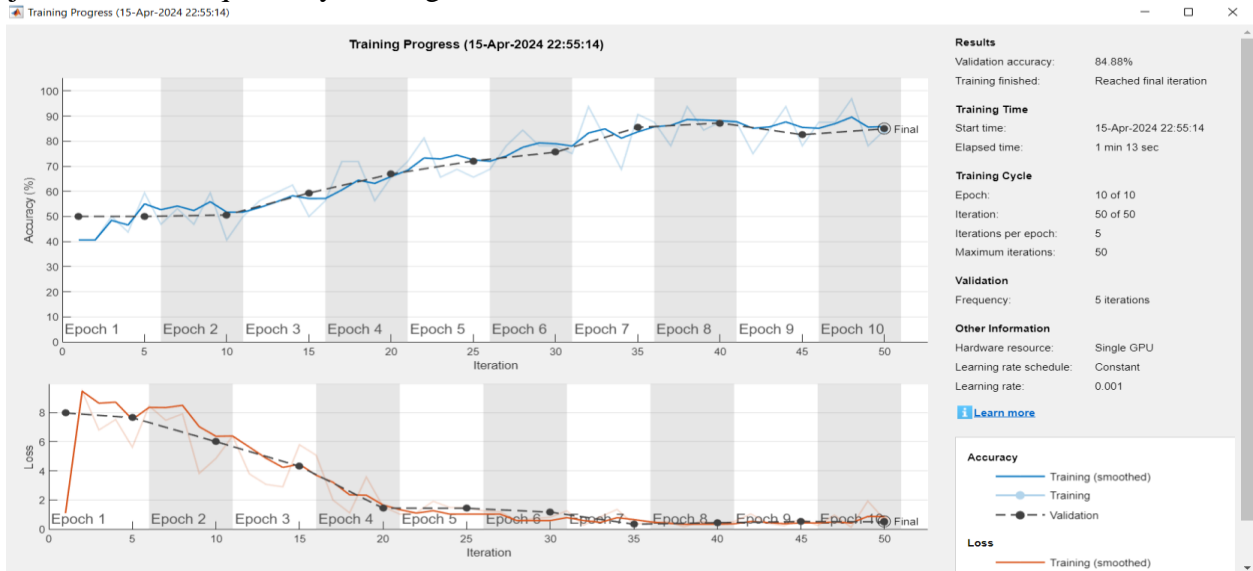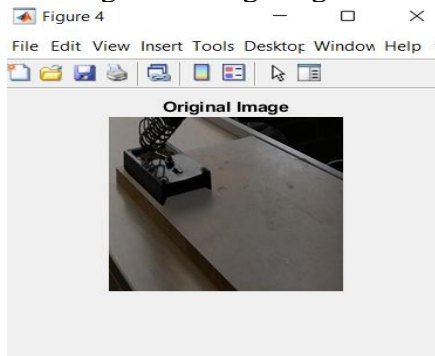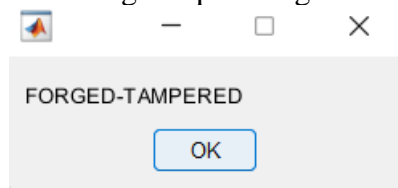


Fig4. Training Progress



Fig5. Input Image



Fig6. Classified Result

| Epoch | Iteration | Time Elapsed (hh:mm:ss) | Mini-batch Accuracy | Validation Accuracy | Mini-batch Loss | Validation Loss | Base Learning Rate |
|---|---|---|---|---|---|---|---|
| 1 | 1 | 00:00:07 | 40.62% | 50.00% | 1.0970 | 7.9712 | 0.0010 |
| 1 | 5 | 00:00:13 | 59.38% | 50.00% | 5.6280 | 7.6468 | 0.0010 |
| 2 | 10 | 00:00:19 | 40.62% | 50.58% | 4.8497 | 6.0059 | 0.0010 |
| 8 | 40 | 00:01:00 | 87.50% | 87.21% | 0.3857 | 0.4440 | 0.0010 |
| 9 | 45 | 00:01:06 | 78.12% | 82.56% | 0.7139 | 0.5241 | 0.0010 |
| 10 | 50 | 00:01:11 | 84.38% | 84.88% | 0.6691 | 0.5029 | 0.0010 |

Enter input image path: 2.jpg
Accuracy: 0.84884

Precision: 1
Sensitivity: 0.69767
Specificity: 1

**Metrics:**

| Metrics | Values |
| --- | --- |
| {'Sensitivity'} | 0.69767 |
| {'Specificity'} | 1 |
| {'Accuracy'  } | 0.84884 |
| {'Precision'  } | 1 |

Confusion Matrix:

    86    0
    26    60

**Confusion Matrix Heatmap:**

The confusion matrix heatmap provides a detailed breakdown of the model's classification results, displaying the distribution of true positive, true negative, false positive, and false negative predictions across different classes. It offers a comprehensive view of the model's performance, allowing for the identification of any misclassifications or imbalances in the dataset.
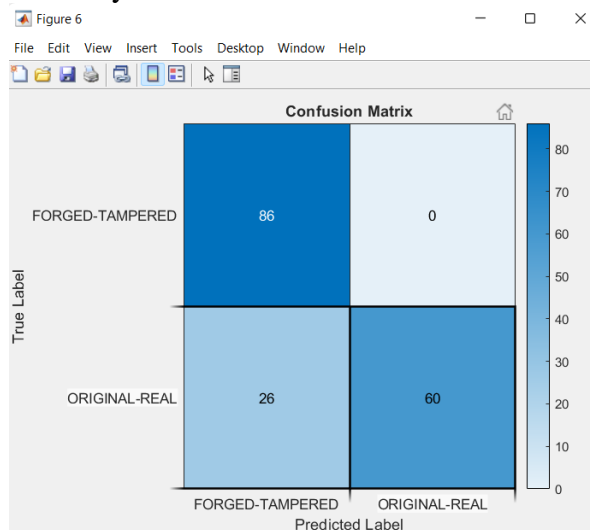


Fig7. Confusion Matrix

**Confusion Matrix Visualization:**

The confusion matrix is visualized using a heatmap, where the rows represent the true labels, and the columns represent the predicted labels. The intensity of the colors indicates the frequency of predictions for each label combination.

**Discussion:**

A discussion of the model's functionality, its drawbacks, and future directions for development comes after the results analysis. To improve the model's effectiveness, this can entail investigating different data augmentation strategies, adjusting model hyper parameters, or implementing sophisticated architecture changes. We may learn more about the behavior and effectiveness of the CNN model in differentiating between authentic and manipulated images by examining these plots and figures. These visuals support decision-making on more iterations or enhancements to the classification system, as well as result interpretation and model improvement. Overall, the study of the results provide insightful information about how well the CNN model performs in differentiating between authentic and manipulated photos, which helps to improve and optimize the system for detecting image tampering.

## V. Conclusion & Future Scope

In conclusion, the CNN model has shown encouraging results in effectively identifying between real and modified visual content when used to classify genuine photographs and tampered images. The CNN model demonstrated acceptable performance metrics throughout training and evaluation, including high levels of accuracy, precision, sensitivity, and specificity. The suggested approach's efficacy in identifying picture tampering is further validated by the examination of confusion matrices and the display of classification outcomes. The CNN model demonstrates stability and generalization capabilities through the use of deep learning techniques and an extensive dataset, rendering it an invaluable instrument for picture forensics and content verification.

In future First off, investigating bigger and more varied datasets may enhance the model's capacity to generalize across various circumstances and tampering approaches. Furthermore, adding sophisticated architectural alterations like recurrent connections or attention processes may improve the model's ability to encode features and identify long-range correlations in pictures. Moreover, incorporating post-processing methods like adversarial training or ensemble learning could improve robustness against adversarial attacks and reduce model vulnerabilities. Furthermore, putting the CNN model to use in practical settings and assessing how well it performs in the face of obstacles and limitations would yield insightful information about deployment and scalability. All things considered, there is a lot of promise for improving the state-of-the-art in content authentication and image tampering detection through ongoing study and innovation in deep learning and image forensics.

## References

1. Zhang, W., Sun, X., Luo, W., & Chang, E. (2019). A hybrid deep learning framework for detecting image splicing and copy-move forgeries. IEEE Transactions on Information Forensics and Security, 14(9), 2421-2436.
2. Fridrich, J., Soukal, D., & Lukáš, J. (2012). Detection of copy-move forgery in digital images. Proceedings of Digital Forensic Research Workshop (DFRWS).
3. Pevný, T., Bas, P., & Fridrich, J. (2010). Steganalysis by subtractive pixel adjacency matrix. IEEE Transactions on Information Forensics and Security, 5(2), 215-224.
4. Bayar, B., & Stamm, M. C. (2016). A deep learning approach to universal image manipulation detection using a new convolutional layer. In Proceedings of the IEEE Workshop on Information Forensics and Security (WIFS).
5. Dong, J., Yu, S., & Huang, W. (2013). Exposing image splicing with inconsistent local noise variance. IEEE Transactions on Information Forensics and Security, 8(4), 593-602.
6. Al-Qararah, W., Khelifi, F., & Bouridane, A. (2020). Deep learning-based image splicing detection using convolutional and recurrent neural networks. IEEE Transactions on Information Forensics and Security, 15, 1010-1025.
7. Toh, K. A., Khairudin, Z., & Abdullah, A. (2018). Deep learning approach to face retouching detection in portraits. 2018 IEEE Conference on Multimedia Information Processing and Retrieval (MIPR), 287-292.
8. Johnson, M. K., & Farid, H. (2007). Exposing digital forgeries by detecting inconsistencies in lighting. IEEE Transactions on Information Forensics and Security, 2(3), 450-461.
9. Popescu, A. C., & Farid, H. (2005). Exposing digital forgeries in color filter array interpolated images. IEEE Transactions on Signal Processing, 53(10), 3948-3959.
10. Amerini, I., Ballan, L., Caldelli, R., & Del Bimbo, A. (2010). A SIFT-based forensic method for copy-move attack detection and transformation recovery. Proceedings of the 20th ACM International Conference on Multimedia (MM '10).
11. Ng, T. T., Chang, S. F., Sun, Q., & Chen, T. (2006). Data-driven forensic feature for image classification. Proceedings of the IEEE International Conference on Multimedia and Expo (ICME).
12. Thepade, S., Das, R. & Ghosh, S.: Content Based Image Classification with Thepade's Static and Dynamic Ternary Block Truncation Coding, Vol.4 (1), 2015, pp. 13-17.

13. Thepade, S.D. Subhedarpage, K. S.Mali, A.A., "Performance rise in Content Based Video retrieval using multi-level Thepade'ssorted ternary Block Truncation Coding with intermediate block videos andeven-odd videos", Proceedings of the 2013 International Conference on Advances in Computing, Communications and Informatics, ICACCI 2013, 2013, pp. 962–966, 6637306

14. Madane M., Thepade S. D., "Score Level Fusion Based Bimodal Biometric Identification Using Thepade'sSorted n-ary Block Truncation Coding with Variod Proportions of Iris andPalmprint Traits", Procedia Computer Science, 2016, 79, pp. 466–473