

# An Improved VGG16 and CNN-LSTM Deep Learning Model for Image Forgery Detection

Yogita Shelar<sup>a</sup>, Dr. Prashant Sharma<sup>b</sup>, Dr.Chandan Singh, D. Rawat<sup>c</sup>

<sup>a</sup>Research scholar, Department of Computer Science, Pacific Institute of Technology, Udaipur, India

<sup>b</sup>Associate Professor, Department of Computer science, Pacific Institute of Technology, Udaipur, India

<sup>c</sup>Head Of Department of Electronic and Telecommunication Vivekanand Education Society's Institute of Technology, Chembur, India

a) [yogitamshelar@gmail.com](mailto:yogitamshelar@gmail.com)

b) [prashant.sharma@pacific-it.ac.in](mailto:prashant.sharma@pacific-it.ac.in)

c) [chandansingh.rawat@ves.ac.in](mailto:chandansingh.rawat@ves.ac.in)

**Abstract:** As the field of image processing and computer vision continues to develop, we are able to create edited images that seem more natural than ever before. Identifying real photos from fakes has become a formidable obstacle. Image forgery has become more common as the multimedia capabilities of personal computers have developed over the previous several years. This is due to the fact that it is simpler to produce fake images. Since image object fabrication might obscure critical evidence, techniques for detecting it have been intensively investigated for quite some time. The publicly available datasets are insufficient to deal with these problems adequately. Our work recommends using a deep learning based image inpainting technique to create a model to detect fabricated images. To further detect copy-move forgeries in images, we use an CNN-LSTM and Improved VGG adaptation network. Our approach could be useful in cases when classifying the data is impossible. In contrast, researchers seldom use deep learning theory, preferring instead to depend on tried-and-true techniques like image processing and classifiers. In this article, we recommend the CNN-LSTM and improved VGG-16 convolutional neural network for intra-frame forensic analysis of altered images.

## I. Introduction

Since the recent contemporary digital media technology, digital content has become the most popular form of communication due to the ease of use, portability, and rich information content. Press, politics, insurance claims, defence, and legal cases are only a few of the many important areas where it has developed into a crucial piece of evidence [1]. Digital content may be easily manipulated by certain non-professionals because to the available latest technology usage of powerful editing capacity, while it may be difficult for experts to tell some hoaxes apart from the actual thing. As a result of these factors, some people are suspicious of digital images [2]. This highlights the importance and requirement of appropriate forensic shames that can identify the veracity, uniqueness, and data integrity. This method's capacity to mitigate damage caused by malicious video manipulation and its contribution to maintaining social peace and stability make it very relevant in the real world. Frame deletion, insertion, and duplication are all examples of full frame forgery, which is the manipulation of individual image frames as forging units [3]. Forgery in which individual image and video clips are used as forgery units is known as "object forgery." Whole-frame forgery and object forgery are the two main categories of digital content manipulation. Object forgery is when a video's time domain and space domain are altered simultaneously while a specific frame region is used as the

"forgery object." The most common methods of digital image modification employed in forgery include intra-frame copy-move forgery, content removal, and image synthesis fraud[4]. Some methods that may be used to spot whole frame forgeries include those based on scene dependencies, flow approaches, the exploitation of compression artefacts, and deep learning. Scene dependency-based algorithms examine the pixel score of every frame in a movie in order to spot irregular periodic artefacts[5]. Optical Flow methods used the variations in an image's brightness along with motion gradients as proof of fabrication. Compression and artefacts are the base for the latest technology rely on a broad range of compression artefacts created during image encoding and decoding for the aim of spotting unusually quick changes in images [6][7]. Using training samples extracted from real images, the latest image processing based solution can automatically recognise and understand key properties.

## II. Related work

Applying CNNs to the problem of detecting forgeries has been taken up by a number of researchers. Convolution neural networks (CNNs) [8] were originally considered as a potential tool for steganography using greyscale images. As such, the initial convolutional neural network (CNN) layer in our implementation is a simple high pass filter. Specifically for the detection of image splicing, a model of the relevant

visual features was developed. The scientists used a technique called discrete cosine transformation (DCT) [9] to accomplish their goal. The strategy was developed to avoid having important features fall within the purview of DCT. By relocating, the DCT domain contributes to the input of the CNN. This procedure occurs when the organism is being fed. The data will be processed by creating a each patch histogram and then during the stage of classification, concatenating each one of the histograms in order to input them into the CNN[10]. One result of applying deep learning techniques to computer vision problems is the data-driven local convolution feature. Research conducted by CNN and others indicated that methods for identifying copy-move fraud relied heavily on computer vision tasks such as segmentation, extraction and identification of the objects. Recent progress [11] in computer vision challenges has been driven in part by breakthroughs in CNN and graphics processing unit technology. Recent CNN-based approaches to image classification use an end-to-end structure, in contrast to traditional techniques that depend mostly on regional descriptors. Through the usage of end-on-end multilayers, deep neural networks are able to include high-, medium-, and low-level classifiers and features, with the latter two tiers' enhancements corresponding to the number of hidden layers. As the field of CNN has advanced recently [12], it has allowed for significant performance gains in many image classification and object recognition tasks. The following table classifies the most common CNN subcategories. The stated intermediate layers of CNNs provides characteristics of descriptors at the image level[13][14]. Having this trait may make class differences more apparent, but it does not distinguish between social strata. Unfortunately, copy-move forgery detection is not a good fit for the deep learning methods often used for computer vision challenges. There have been applications of deep learning in this domain as well. As was said in the introductory paragraph, this kind of detection looks for the same sorts of scaled, rotated, or distorted regions as those described above. The use of suitable patch-level descriptors [15-19] may allow data-driven patch-level descriptors to replace handmade ones, as shown by the output of CNNs features. The CNNs' expressive feature representations output suggests this is a real possibility. Several deep local descriptors with excellent matching and classification abilities have been provided in recent studies [20]. Since CNNs have already been shown successful in detecting natural image distribution, it is to assume that they will also be effective in detecting image copy-move forgeries. The reason for this [21] is because the objective of this exercise is to identify the "natural" or "pristine" image amid any artificial or manipulated versions. CNNs have shown [22] [23] to be useful for analysing the distribution of images in

their natural environments. To spot copy-move forgeries and locate them, image properties are the major consideration during classification. This means that CNNs' ability[24] to extract image features is vital to their role in detecting copy-move frauds. We'll show how the automated CNN technique differs from the conventional approach to feature extraction and how its usage might potentially eliminate the requirement for the conventional method in certain situations [25][26].

### III. Dataset:

The CASIA.2.0 image forgery database contains a total of 12,614 images, 7,491 of which are Genuine Images and the remaining 5123 are Tempered Images. The training data include 80% of the images, which total 10091. Out of these, 5993 are the real images, and 4098 are the tempered images. The test data comprise 20% of the images, totaling 2523. The 1498 actual images were taken into consideration for the testing data, while 1025 images were manipulated.



Figure 1(A). Sample CASIA 2.0 Genuine Images



Figure 1(B). forgery CASIA 2.0 Genuine Images

### IV. Propoed System

#### a. CNN-LSTM Network

The image is downsampled by the first convolution layer of CNN [27][28], which extracts information from neighbouring pixels. The convolution may thus be seen as a simple addition of the relative significance of the input image with intensity values. In the CNN-LSTM network, this is performed by convolving a 64-by-64-pixel input image with a 5-by-5 filter kernel. The approach eith generated image that is a more manageable file size. Each convolutional-layer of the LSTM will generate convolutions in multiplication, which will generate a tensor weight proportional to the same of convolutions, n. In this specific case, the dimensions of the tensor's will be 5x5 n.

The convolution-1 layer of the CNN-LSTM will produce a matrix weight with the size 128 by 5 by 5. This will result in 1600 parameters being generated. As the final layer along with Maxpolling is the prediction layer aids in the completion of the class categorization for classification.



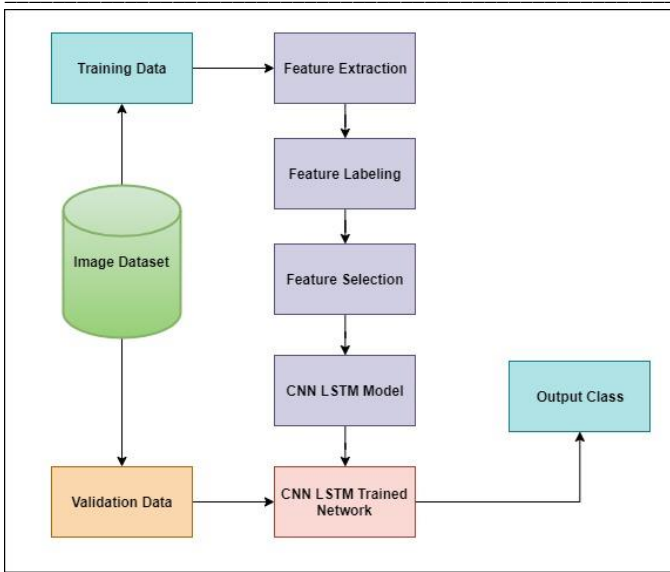


Figure 2. CNN-LSTM Model

Recent studies have shown that CNNs do very well in detecting and reporting instances of photo fraud. Therefore, this research aims to offer a complete CNN-LSTM that can handle and identify copy-move forgeries.

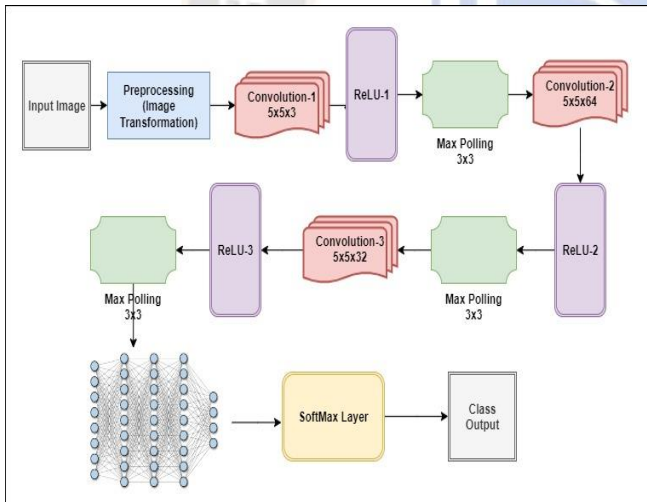


Figure 3. CNN-LSTM Network Layers

An input layer, several convolutional layers, multiple fully connected layers, a classification layer, and an output layer are the building blocks of the proposed LSTM. Using convolutional neural networks (CNNs) with LSTM to aid in a copy-move forgery detection model is advantageous because of CNNs' effectiveness as a feature extraction strategy, which in turn boosts the model's overall performance. In addition, the capacity of the CNN to learn may be enhanced to provide better output results, and this in turn can be achieved by exposing it to a larger set of input samples and repeating the training process more often. It's feasible to do so. Utilizing CNNs [29] for the purpose of

identifying copy-move fraud is far more cost-effective than the conventional method. Lastly, CNN may use a wide variety of images as input, which improves the precision of the model's predictions.

### b. Improved VGG Network

After doing extensive study on the different methodologies, we determined that the great majority of transfer learning methods for biomedical imaging used VGG gives the maximum prediction accuracy. The authors of this work were driven to construct VGG16 by Hyper-tuning the various parameters to achieve the best achievable degree of accuracy. This architecture, known as VGG16, is a deep CNN. It has sixteen levels of thickness. In around 92.7% of the top-five ImageNet tests, the VGG16 model obtains a satisfactory score. ImageNet has about 14 million images that may be organised into over a thousand separate categories. In addition, the most votes were cast for one of the models shown at the 2018 ILSRC Conference. As seen in Figure 4, the VGG16 architecture needs an input image with 224 by 224 pixel size. VGG16's architecture is composed of five distinct components. In the first and second blocks, 64 and 128 filters are used to apply two convolution layers (3 \* 3) and one max pooling layer (2 \* 2), respectively. In the third, fourth, and fifth blocks, three use 256, 512, and 512 filters, respectively. This is followed by a maximum pool layer that is twice as big as the preceding layer. As a consequence, the VGG16 model is modified further in the proposed work by adding 2 dense layers, 1 flatten layer with the three ReLU activation function, and a multiple dense layer with the activation function. This is performed to increase the model's precision.

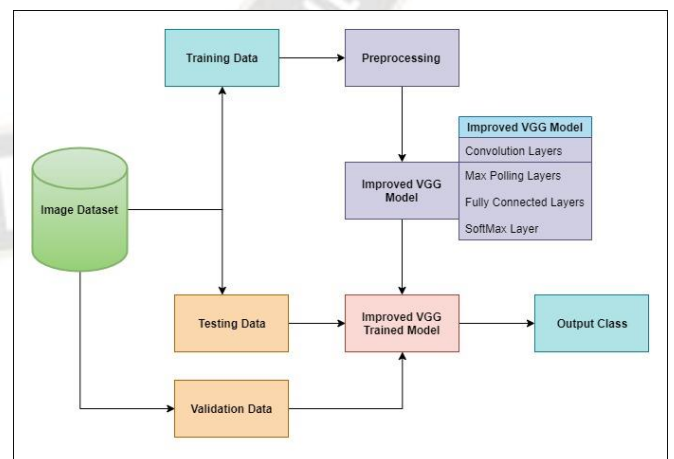
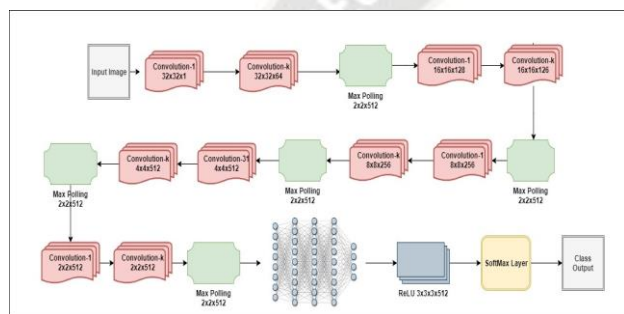


Figure 4. Proposed VGG Model

A fully connected layer is inserted before the input of the VGG-16 network to facilitate training and testing, and to facilitate the construction of the same-sized characteristic

map. This layer's job is to transform multi-dimensional qualities into one-dimensional characteristics so that they may be utilised in building the characteristic map. Table I shows the values for the parameters of the RELU function utilised as the activation function for the model.

First, choose the feature set collected random featured data is then pass it on to the first completely linked layer. A 1024-dimensional feature will be the end outcome of this operation. First, a 32x32x1 feature image is built; next, using that image as input, the convolution block is carried out using the values from the table. I in each iteration, producing an image with dimensions of 1 by 1 by 512; 3) The output is formed using as input the images of dimensions 1 by 1 by 512 produced by the convolution layer sequence. In each cycle, this input is sent into a fully linked layer followed by another completely connected layer.



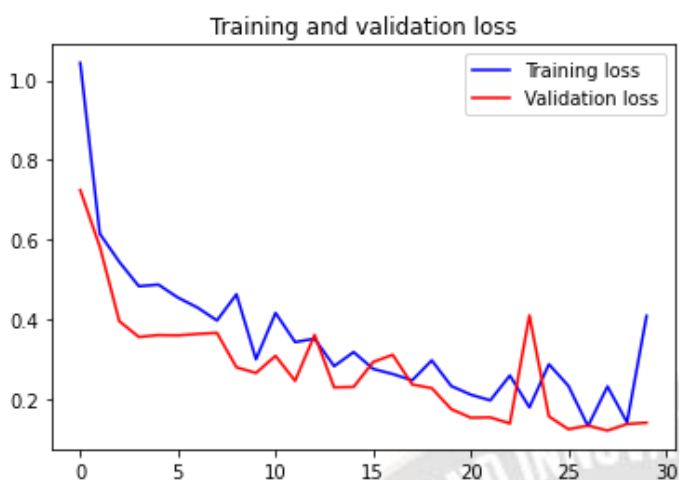


Figure 7. Training and Validation Loss

**e. IMPROVED VGG 16: MODEL 1 (VGG 16 With Last 4 Layer Removed)**

Layer (type)	Output Shape	Param #
input_6 (InputLayer)	[(None, 224, 224, 3)]	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0

=====  
 Total params: 14,714,688  
 Trainable params: 7,079,424  
 Non-trainable params: 7,635,264

**f. IMPROVED VGG 16: MODEL 2 (UPDATED VGG 16)**

Layer (type)	Output Shape	Param #
vgg16 (Functional)	(None, 7, 7, 512)	14714688
conv2d_5 (Conv2D)	(None, 7, 7, 512)	262656
activation_5 (Activation)	(None, 7, 7, 512)	0
max_pooling2d_5 (MaxPooling 2D)	(None, 3, 3, 512)	0
dropout_10 (Dropout)	(None, 3, 3, 512)	0
flatten_5 (Flatten)	(None, 4608)	0
dense_10 (Dense)	(None, 1024)	4719616
dropout_11 (Dropout)	(None, 1024)	0
dense_11 (Dense)	(None, 1)	1025

=====  
 Total params: 19,697,985  
 Trainable params: 12,062,721  
 Non-trainable params: 7,635,264

**g. Improved VGG Training and Validation**

Figures 8 and 9 depict the training plots and validation plots for the enhanced VGG model behavior during training and testing of the model on the CASIA.2.0 image forgery database, respectively. These plots were generated using the CASIA.2.0 image forgery database.

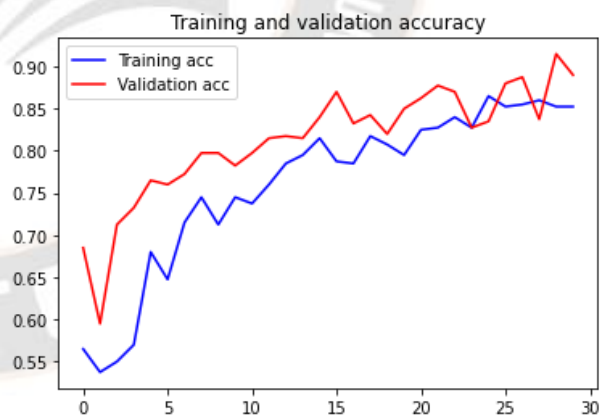


Figure 8. Training and Validation accuracy

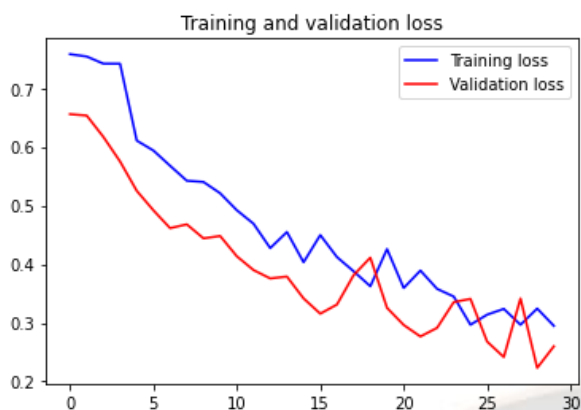


Figure 9. Training and Validation Loss

Figures 10 and 11, show the accuracy and loss plots for the CNN-LSTM network as well as the enhanced VGG model behavior during training and validation of the model on the CASIA.2.0 image forgery database, respectively. Also shown in these figures is the improved VGG model behaviour.

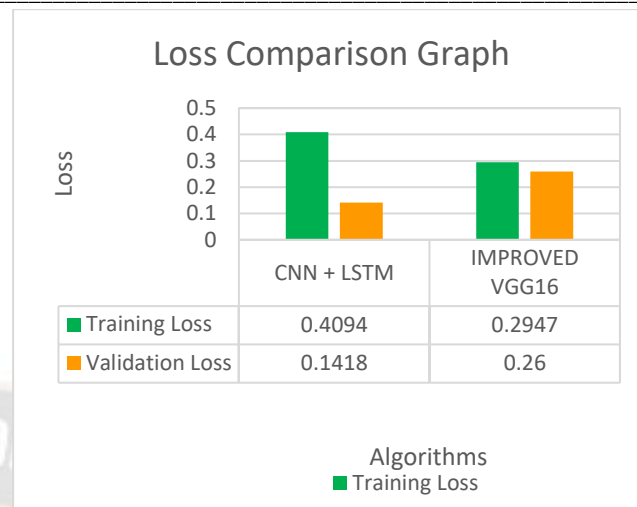


Figure 11. Loss plot for CNN-LSTM and Improved VGG

**Conclusion:**

As research and development in the fields of image processing and computer vision continue to advance, we are gaining the ability to modify images so that they seem more authentic than in the past. The process of distinguishing genuine images from counterfeit ones has evolved into a challenging challenge. The proliferation of multimedia capabilities inside personal computers over the course of the most recent few years has led to an increase in the incidence of image counterfeiting. This is as a result of the fact that it is now much easier to create tempered image sequences. Since the manufacture of image objects has the potential to disguise important evidence, methods for uncovering its presence have been the subject of extensive research for quite some time. The datasets that are accessible to the public are inadequate to deal with these issues in an acceptable manner. As a result of our research, we suggest developing a model to identify manufactured images by using a image forgery detection approach that is based on deep learning. We make use of a CNN-LSTM and an improved VGG 16 adaptation network in order to improve our ability to identify copy-move forgeries in photographic images. When it comes to situations in which the data cannot be classified, our method could come in handy. On the other hand, academics only sometimes make use of deep learning theory, opting instead to rely on tried-and-true methods such as image processing and classifiers. For the purpose of doing intra-frame forensic analysis of changed images, the CNN-LSTM gives accuracy above 90% and improved VGG 16 networks gives accuracy around 90% come highly preferred. Thus we recommend CNN-LSTM as a method for forgery detection in the images.

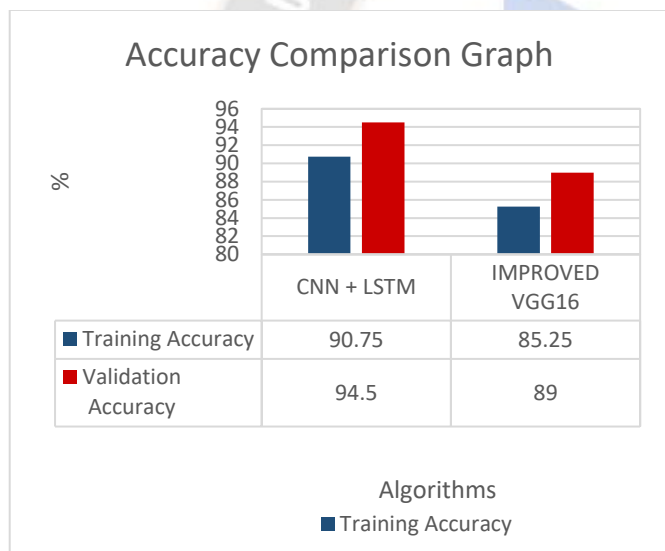


Figure 10. Accuracy plot for CNN-LSTM and Improved VGG



## References:

- [1] S. Chen, J. Huang, J. Huang, and J. Huang, "Automatic Detection of Object-Based Forgery in Advanced Video," *IEEE Transactions on Circuits & Systems for Video Technology*, vol. 26, pp. 2138-2151, 2016.
- [2] S. Ajani and M. Wanjari, "An Efficient Approach for Clustering Uncertain Data Mining Based on Hash Indexing and Voronoi Clustering," 2013 5th International Conference and Computational Intelligence and Communication Networks, 2013, pp. 486-490, doi: 10.1109/CICN.2013.106.
- [3] Bagiwa M A, Wahab A W A, Idris M Y I, et al. "Digital Video Inpainting Detection Using Correlation of Hessian Matrix".*Malaysian Journal of Computer Science*,2016, 29(3):179-195.
- [4] Su K, Member S, Kundur D, et al. "Statistical invisibility for collusionresistant digital video watermarking". *IEEE Press*, 2005.
- [5] Simonyan K, Zisserman A. "Very Deep Convolutional Networks for Large-Scale Image Recognition".*Computer Science*, 2014.
- [6] Pevny T, Fridrich J. "Merging Markov and DCT features for multiclass JPEG steganalysis" [C]//*Security, Steganography & Watermarking of Multimedia Contents IX*. International Society for Optics and Photonics, 2007.
- [7] Jan Kodovský, Fridrich J. Calibration revisited[C]// 2009. pp.63-74. [16] Pevný T, Bas P, Fridrich J. "Steganalysis by subtractive pixel adjacency matrix" . *IEEE Transactions on Information Forensics & Security*, 2010, 5(2):215-224.
- [8] Kodovsky J, Fridrich J, Memon N D, et al. " Media Watermarking, Security, and Forensics 2012 Steganalysis of JPEG images using rich models" . *Proceedings of SPIE - The International Society for Optical Engineering*, 2012, 8303:83030A.
- [9] Kodovsky J, Fridrich J, Holub V. "Ensemble Classifiers for Steganalysis of Digital Media" .*IEEE Transactions on Information Forensics and Security*, 2012, 7(2):432-444.
- [10] Fujisawa, Y.; Otomo, Y.; Ogata, Y.; Nakamura, Y.; Fujita, R.; Ishitsuka, Y.; Watanabe, R.; Okiyama, N.; Ohara, K.; Fujimoto, M. Deep-learning-based, computer-aided classifier developed with a small dataset of clinical images surpasses board-certified dermatologists in skin tumour diagnosis. *Br. J. Dermatol.* 2019, 180, 373–381.
- [11] Garcia-Arroyo, J.L.; Garcia-Zapirain, B. Recognition of pigment network pattern in dermoscopy images based on fuzzy classification of pixels. *Comput. Methods Programs Biomed.* 2018, 153, 61–69.
- [12] Ajani, S.N., Bhanarkar, P. (2022). Design a Mechanism for Opinion Mining. In: Sharma, S., Peng, S.L., Agrawal, J., Shukla, R.K., Le, D.N. (eds) *Data, Engineering and Applications*. Lecture Notes in Electrical Engineering, vol 907. Springer, Singapore. [https://doi.org/10.1007/978-981-19-4687-5\\_35](https://doi.org/10.1007/978-981-19-4687-5_35).
- [13] Iyatomi, H.; Oka, H.; Celebi, M.E.; Ogawa, K.; Argenziano, G.; Soyer, H.P.; Tanaka, M. Computer-based classification of dermoscopy images of melanocytic lesions on acral volar skin. *J. Investig. Dermatol.* 2008, 128, 2049–2054.
- [14] Chatterjee, S.; Dey, D.; Munshi, S. Optimal selection of features using wavelet fractal descriptors and automatic correlation bias reduction for classifying skin lesions. *Biomed. Signal Process. Control* 2018, 40, 252–262.
- [15] Chatterjee, S.; Dey, D.; Munshi, S. Integration of morphological preprocessing and fractal based feature extraction with recursive feature elimination for skin lesion types classification. *Comput. Methods Programs Biomed.* 2019, 178, 201–218.
- [16] González-Díaz, I. Dermaknet: Incorporating the knowledge of dermatologists to convolutional neural networks for skin lesion diagnosis. *IEEE J. Biomed. Health Inform.* 2018, 23, 547–559.
- [17] Samir N Ajani Piyush K. Ingole , Apeksha V. Sakhare "Modality of Multi-Attribute Decision Making for Network Selection in Heterogeneous Wireless Networks", *Ambient Science - National Cave Research and Protection Organization, India*,2022, Vol.9, Issue.2, ISSN- 2348 5191. DOI:10.21276/ambi.2022.09.2.ta02
- [18] Kawahara, J.; Daneshvar, S.; Argenziano, G.; Hamarneh, G. Seven-point checklist and skin lesion classification using multitask multimodal neural nets. *IEEE J. Biomed. Health Inform.* 2018, 23, 538–546.
- [19] Koohbanani, N.A.; Jahanifar, M.; Tajeddin, N.Z.; Gooya, A.; Rajpoot, N. Leveraging transfer learning for segmenting lesions and their attributes in dermoscopy images. *arXiv* 2018, arXiv:1809.10243.
- [20] Filali, Y.; El Khoukhi, H.; Sabri, M.A.; Yahyaouy, A.; Aarab, A. Texture Classification of skin lesion using convolutional neural network. In *Proceedings of the 2019 International Conference on Wireless Technologies, Embedded and Intelligent Systems (WITS)*, Fez, Morocco, 3–4 April 2019; pp. 1–5.
- [21] S. N. Ajani and S. Y. Amdani, "Probabilistic path planning using current obstacle position in static environment," 2nd International Conference on Data, Engineering and Applications (IDEA), 2020, pp. 1-6, doi: 10.1109/IDEA49133.2020.9170727.
- [22] Kadampur, M.A.; Al Riyae, S. Skin cancer detection: Applying a deep learning based model driven architecture in the cloud for classifying dermal cell images. *Inform. Med. Unlocked* 2020, 18, 100282.
- [23] Menegola, A.; Fornaciali, M.; Pires, R.; Bittencourt, F.V.; Avila, S.; Valle, E. Knowledge transfer for melanoma screening with deep learning. In *Proceedings of the International Symposium on Biomedical Imaging, Melbourne, Australia*, 18–21 April 2017.
- [24] Lee, H.D.; Mendes, A.I.; Spolaor, N.; Oliva, J.T.; Parmezan AR, S.; Wu, F.C.; Fonseca-Pinto, R. Dermoscopic assisted diagnosis in melanoma: Reviewing results, optimizing methodologies and quantifying empirical guidelines. *Knowl.-Based Syst.* 2018, 158, 9–24.
- [25] LeCun, Y.; Bengio, Y.; Hinton, G. Deep learning. *Nature* 2015, 521, 436–444.
- [26] Duggani, K.; Nath, M.K. A Technical Review Report on Deep Learning Approach for Skin Cancer Detection and Segmentation. *Data Anal. Manag.* 2021, 54, 87–99.

- [27] Khan, M.A.; Zhang, Y.D.; Sharif, M.; Akram, T. Pixels to Classes: Intelligent Learning Framework for Multiclass Skin Lesion Localization and Classification. *Comput. Electr. Eng.* 2021, 90, 106956.
- [28] Alharithi, F.; Almulhi, A.; Bourouis, S.; Alroobaea, R.; Bouguila, N. Discriminative Learning Approach Based on Flexible Mixture Model for Medical Data Categorization and Recognition. *Sensors* 2021, 21, 2450.
- [29] Masud, M.; Singh, P.; Gaba, G.S.; Kaur, A.; Alghamdi, R.A.; Alrashoud, M.; Alqahtani, S.A. CROWD: Crow Search and Deep Learning based Feature Extractor for Classification of Parkinson's Disease. *ACM Trans. Internet Technol. (TOIT)* 2021, 21, 1–18.

