# Integration Of Medical Image Processing Techniques In Deeper Neural Networks With Artificial Intelligence

[1]Dr. Amar Bharatrao Deshmukh,[2]Ashoka S, [3]Dr.D. Sharada Mani,[4]Dr. Deepali Yewale
[1]Department of E & TC, Associate Professor and Head of Department, ABMS'P Anantrao Pawar College of Engineering & Research, Pune, Maharashtra
[2]Assistant Professor, Department of Artificial Intelligence and Data Science, Global Academy of Technology, Bengaluru
[3]Associate Professor, Department of AIML, Qis College of Engineering and Technology, Ongole, Prakasam.
[4]Assistant Professor, Department of Electronics and Telecommunication Engineering,AISSMS Institute of Information Technology, Pune
amarbdeshmukh@gmail.com[1],ashoka.s@gat.ac.in[2],sharadamani.d@qiscet.edu.in[3],deepali_yewale@yahoo.co.in[4]

***Abstract***
*The interpretation and communication of image meaning may be achieved through the utilization of techniques derived from both natural language processing and computer vision. Can a machine replicate the ability of the human brain to offer a comprehensive depiction of an image.The task of captioning photographs is widely acknowledged as a challenging endeavor within the field of artificial intelligence. The process of transforming an image into grammatically accurate text necessitates the utilization of both natural language processing and computer vision techniques. The advancement of deep learning methodologies and the abundance of publicly accessible datasets have facilitated the development of diverse models for the automated generation of image descriptions. The initial stage in generating a satisfactory description of an input image is to classify it based on the highest number of objects present. By using concepts from Natural Language Processing (NLP) and a neural network, we may do this. This paper provides a comprehensive explanation of the integration of Long Short-Term Memory with a "Convolutional Neural Network" in order to generate a visual representation. The classification of pictures and text is facilitated by the utilization of Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). In order to enhance the accuracy of its predictions, we trained the model to utilize a larger lexical vocabulary when characterizing the photos it has encountered. After conducting several trials on various photo datasets, we have shown that visual description is the primary determinant of a model's accuracy. Typically, the result enhances as the size of the dataset increases.*
***Keywords****: convolutional neural networks, long short-term memory, ResNet, and Vgg16.*

## INTRODUCTION

The process of automatically converting the information contained inside a picture into a human language, such as English. This AI topic presents a formidable challenge, although it has the potential to significantly enhance the lives of individuals with visual impairments by enabling them to access the vast amount of information available on the internet through audio descriptions of visual content. The integration of photo classification with NLP adds complexity to this task. For example, the description should encompass not only the elements seen in the image, but also the interrelationships among the objects, as well as their characteristics and involvement. Image classification and object recognition also have significant functions. In order to achieve this objective, we provide a singular, all-encompassing model that, when provided with an image as input from the user, is trained to optimize the likelihood of producing a set of target words, denoted as R = [R1, R2,...], whereby each of these words effectively represents the entirety of the image.

Recent studies have demonstrated that the utilization of "recurrent neural nets" (RNNs) [3] can be a viable approach for the task of phrase translating. The encoder of the Recurrent Neural Network (RNN) subsequently retrieves the phrase via the source supplied by the Convolutional Neural Network (CNN), which has previously divided the image into distinct objects. The encoder then creates a vector that stores the words and converts them into a feature vector. The concealed decoder is now generating the intended phrase sequence of the encoder.

The utilization of a "Convolutional neural networks Neural Network" is suggested as an alternative to the encoder. In recent years, there has been increasing evidence [1] indicating that Convolutional Neural

Networks (CNNs) may effectively represent a given input image by constructing a vector of a predetermined length. Therefore, "Convolutional neural network Neural Network" may be employed as encoders. Once the picture classifying job has been pre-trained, the description provided here is generated by the "recurrent ones neural network decoder" (Fig. 1). This is sometimes referred to as "Neural Image Caption".
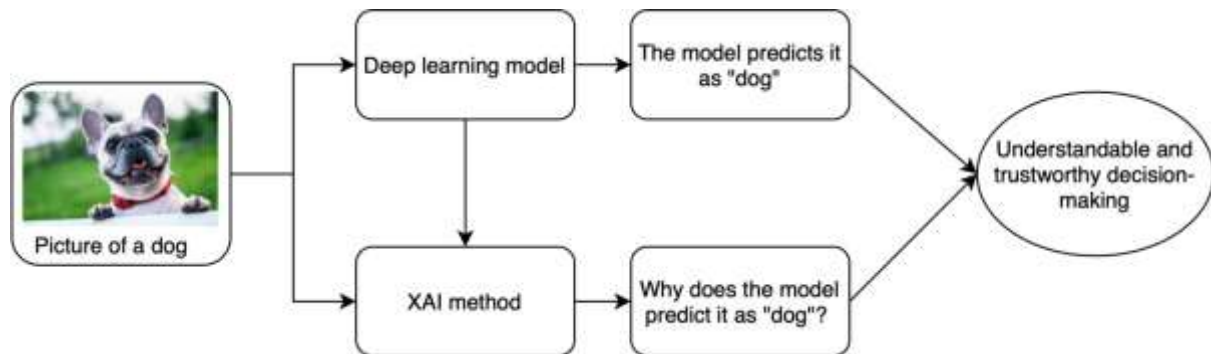


**Fig.1: Basic idea of the model**

An example of a neural network that is deep is Convolutional Neural Networks ((CNN). This development represents a significant advancement in the domain of image recognition. The major purpose of this system is to analyze the content of images, and it frequently functions inconspicuously throughout the process of categorizing the images. Several tweets have an image associated with an autonomous Uber car. They are present everywhere, ranging from the medical profession to the police force. In essence, image classification refers to the process of utilizing an image as input and generating an output that represents either a class or the likelihood of a class associated with the input. The identification of characteristics in photographs is a fundamental aspect of the functioning of Convolutional Neural Networks (CNNs), resulting in time and effort savings compared to human feature identification. During the process of learning from the photographs, certain attributes are not required. Consequently, deep learning models exhibit enhanced proficiency in executing computer vision tasks. Convolutional Neural Networks (CNNs) has the ability to autonomously identify features due to its concealed layers.

The quality of outcomes is contingent upon the quantity of data included in the training of the model. The model that experiences the least amount of loss is virtually certain to yield dependable outcomes.

**Relevant Literature**

Earlier photo captioning systems depended on templates [2] instead of on random variation to grow and generate natural language captions. Initial endeavors in visual recognition mostly concentrated on the arrangement of pictures, such as assigning scores to photographs according to a predetermined set of categories. In recent years, there has been significant progress in the field of photo clustering, particularly due to the utilization of deep learning techniques [1,3]. The authors of [5] also introduce Image Net, a novel database that adheres to the exact same progressive ideals as Word Net. In order to get further insights into the visual intricacies and emotional tone of the internal conversation, [8] consulted databases containing pictures and their corresponding decision representations. The authors of the study showcased the development of a Multimodal Persistent Neural System that largely depends on the efficient collection of co-direct course data pertaining to the activity of highlights. This system aims to acquire the ability to construct distinct representations of pictures. To make substantial progress in the field of image data, a methodology was devised in reference [4] that focuses on the creation of an automated natural language description for a model. In order to address the issue of LSTM units appearing too huge and immobile in time. [9] The impetus for developing this model stemmed from the triumph of the "Sharing Time model," which employs "Google Net" "CNN" for image thumbnail segmentation and LSTM cells for subtitle generation. The primary objective of [6] is to explore the complex connection between words and imagery. "bidirectional connectionist systems over language," and a unifying purpose for modifying the model quality. The trained configurations are then employed in diverse manners inside the subsequent architecture of the Recurrent Neural Network to provide data pertaining to the necessary picture inscriptions. Convolutional Neural Networks (CNNs) have demonstrated their efficacy as a suitable model for tasks such as photo categorization and object identification. The integration of GloVe,

word2vec, and connectionist systems has the potential to generate figure representations by combining visual elements with verbal demonstration [10, 11]. Lastly, it is important to construct a dedicated facility capable of holding the PC vision models and the encoding LSTM, which will be employed in the development of the image caption generator [12]. The user has given entire phrases that describe the photographs.

**Research Methodology**

In this research, we develop a neural and probabilistic framework known as the "CNN"-RNN model to generate these descriptions.

The term "CNN" refers to a network of neuronal units that acquire knowledge from their environment. It belongs to a specific category of deep learning networks. When it comes to image recognition, "CNN" is just a front for an innovative and groundbreaking approach. Nearly everyone want to engage in remote work in image classification and frequently seeks to analyze optical allusion.
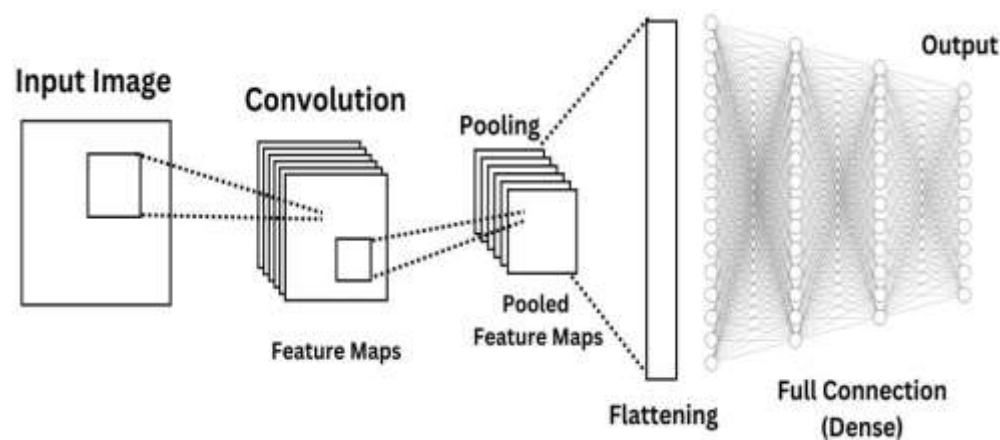


**Fig.2: CNN Architecture**

The initial step of our assault approach involves the implementation of a fourier operation, which will be elaborated upon in the subsequent sections. In this section, we will examine characteristic detectors, which function as filters inside the neural network.

Step 2: Pooling: We will examine current instances of pooling to gain an understanding of its typical functioning. However, in this instance, our attention will be directed towards a particular type of pooling referred to as "the maximum pooling" as our point of connection. Nevertheless, we will explore a range of techniques, such as median (or total) pooling. Upon completion of this phase, you will have the opportunity to view an exhibit that has been developed using a visual artificial intelligence system to thoroughly analyze the entire concept.

Step 3: Leveling:

This document provides a succinct elucidation of the leveling approach employed by "CNN" and the subsequent shift beyond pooled to smoother layers.

A comprehensive examination of an encoder/decoder architecture designed for application in Deep Learning, illustrating the collaborative functioning of the two components in describing images. At its core, this model may be seen as a conceptual framework.
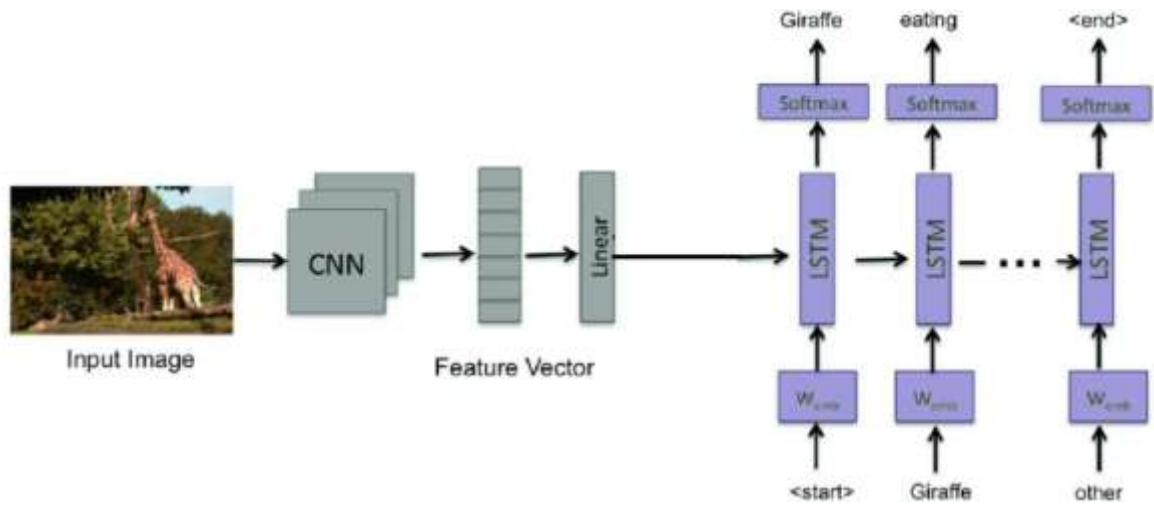
**Fig.3: Model for Image Caption Generator**

The LSTM is utilized for feature extraction in an RNN, whereas the CNN is used for highlight extraction in an image. The Long Short-Term Memory (LSTM) model will process the data using a Convolutional Neural Network (CNN) in order to provide a description of the image. Figure 2 illustrates an approach that employs short-term short-term memory and "A convolutional Neural Networks" to generate visual descriptions. Recent machine translation statistic effort has shown improved outcomes due to advancements in probability and inference training. In order to decode the result sequence, it is necessary to initially encode the change in length into an exact spatial vector. Based on the image, the probability of theft is heightened.

$$\Theta = \arg\max \sum_{i,s}^{n} \quad log\ (S|I:1) \qquad (1)$$

$$\text{Log p }(S|I) = \sum_{t=n}^{N} \quad log\ log\ p(St|I, S0 \dots . St-1) \qquad (2)$$

The RNN is trained utilizing the Long-short lasting storage encoder genesis v4 and is provided with both images and text as input. Following the acquisition of the pictures, they are then fed into a "Conventional Neural Circuit" (CNN) for the purpose of interpretation and application in tasks such as object detection and identification. The vocabulary terms will be symbolized using an embedding approach. The storage of gate control data within a cell can be achieved by the utilization of the LSTM (long short-term memory) concept.
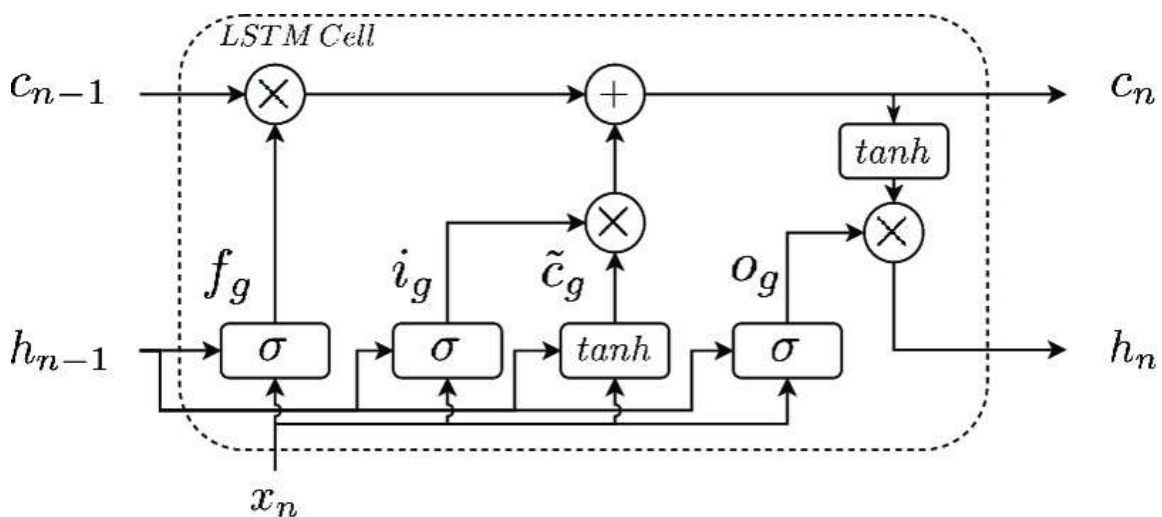


**Fig.4: LSTM cell structure**

The cell responsible for controlling the three gates, as seen in Figure 3, is situated within the memory block. At time t-1, the blue lines will traverse the gates, establishing a connection between the outcomes

of m and the keepsakes at time t. The LSTM cell architecture under consideration has three unique types of gates. The three gates, namely input, forget, and output, play a crucial role in governing the general structure of the cell and providing information about the cell's value to the forget gate. At time t, the memory result m contains the word caption from time t-1 due to the softmax of the word forecast. Examine the data and obtain the latest cell value. Modifications are made to parameters such as the quantity of W matrices and the multiplication of the utility values of the gates and cells. The utilization of multiplicative gates enables the LSTM to be trained in a manner that efficiently handles gradients that exhibit fast growth followed by rapid disappearance.

$$i_t = \sigma(W_{ix}x_t + W_{im}m_{t-1}) \qquad (3)$$
$$f_t = \sigma(W_{fx}x_t + W_{fm}m_{t-1}) \qquad (4)$$
$$\Theta_t = \sigma(W_{ox}x_t + W_{om}m_{t-1}) \qquad (5)$$
$$C_t = f_t.ct-1 + i_t.h(W_{cx}xt + W_{cm}m_{t-1}) \qquad (6)$$
$$M_t = \Theta_t . C_t \qquad (7)$$
$$P_{t+1} = \text{softmax}(m_t) \qquad (8)$$

The non-linearities involved are the hyperbolic tangent fi(-) and the sigmoid w. To derive a probability distribution p across all captions, the final equation is employed to execute the Soft max operation.

Experimental Procedure

4.1. Metrics for Evaluation

The equation for BLEU-N is as follows:

$$\text{Log } BN = \min(0,1-r/c) + x + a = 1NWNogpn \qquad (9).$$

The idea being referred to is associated with the "Billingual Evaluation." In this context, (BN) denotes the BLEU-N metric, (r) represents the effective corpus length, (c) signifies the candidate translation length, and (wN) represents the weights.

The purpose of this approach is to evaluate the precision of machine-generated texts. The output of this integer will range from 0 to 1. Values that are in closer proximity to each other imply a higher degree of alignment between the generated sentence and the original information.

4.2. Collections of data

Our endeavors were subjected to rigorous evaluation using diverse datasets comprising photos accompanied by corresponding written descriptions. The Flickr8k dataset, including 8,000 photographs accompanied by three verbal descriptions each image, was employed in our study. Additionally, the "MS COCO dataset" was utilized, consisting of 82346 images and five descriptive descriptions per image.

DISCUSSION AND RESULTS

Table1: VGG16 Features

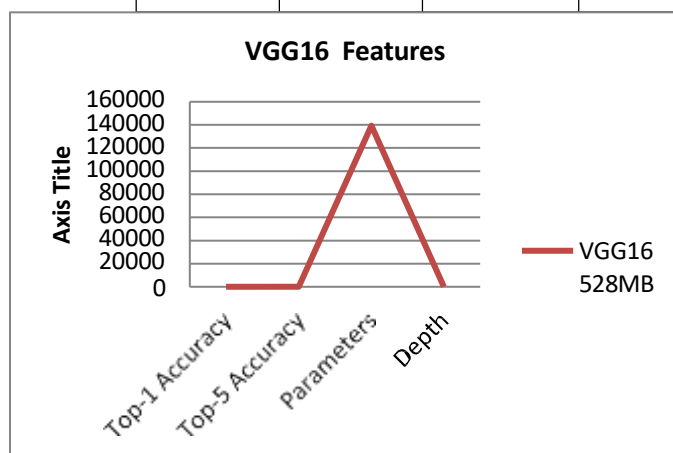| Model | Size | Top-1 Accuracy | Top-5 Accuracy | Parameters | Depth |
|---|---|---|---|---|---|
| VGG16 | 528MB | 0.87 | 1.012 | 139357.544 | 24 |



Fig.5: VGG 16 Model

The output of the vgg16 model is less dependable compared to that of the ResNet model. Vgg16 is more intricate than ResNet due to its utilization of a greater number of factors. ResNet has superior accuracy

in comparison to VGG16. We employed a pretrained ResNet model in this method. We chose ResNet152 over VGG16 due to its extensive network architecture consisting of 152 layers. ResNet utilizes a skip link, enabling the straight transmission of information from the preceding layer to the subsequent layer without any alteration. The system demonstrated exceptional performance in both segmentation and detection tasks in the "MS COCO 2015" competition, as well as in image detection, classification, and localization tasks at the "ILSVRC 2015" competition.

**Table 2: ResNet Features**

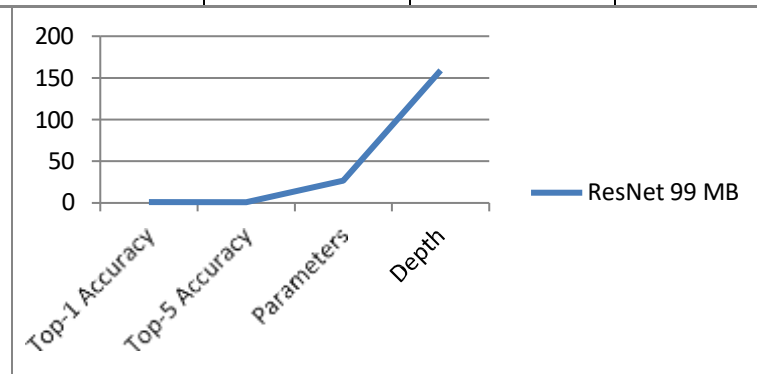| Model | Size | Top-1 Accuracy | Top-5 Accuracy | Parameters | Depth |
|-------|------|----------------|----------------|------------|-------|
| ResNet | 99 MB | 0.859 | 0.856 | 26.636712 | 159 |



**Fig.6: ResNet Model**

ResNet has demonstrated superior accuracy compared to the VGG16 model. We conducted an evaluation of our model using a limited number of randomly selected photographs from the internet. The generated descriptions closely resembled the authentic ones.

**Fig.7: Example of Random Images from the Internet Provides Very Similar to the Real Models**

Scores on our evaluation measure (BLEU) for "VGG16" and ResNet are as follows.



**Table 3: BLEU Evaluation scores**

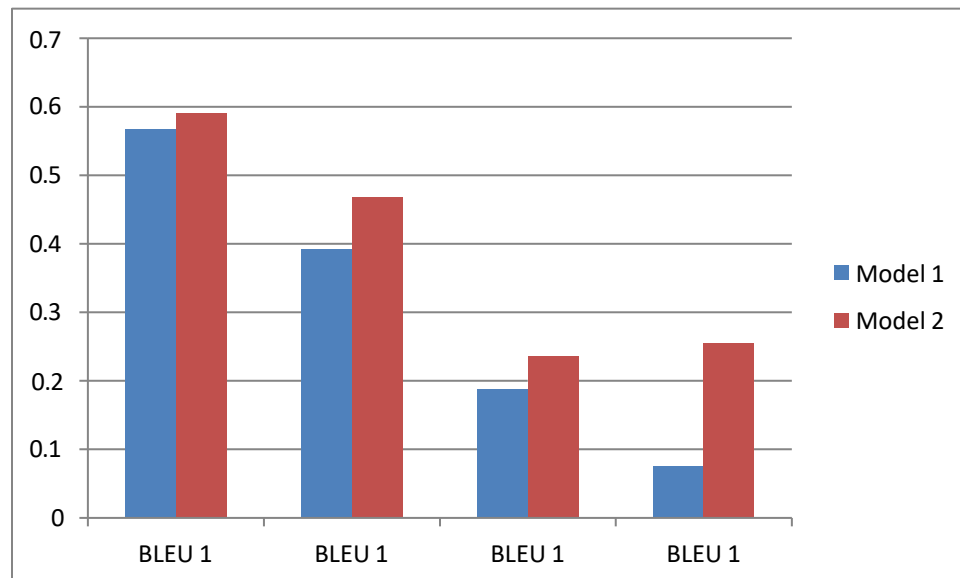|  | Model 1 | Model 2 |
|--|---------|---------|
| BLEU 1 | 0.566697 | 0.58989 |
| BLEU 1 | 0.392564 | 0.4678 |
| BLEU 1 | 0.188049 | 0.23567 |
| BLEU 1 | 0.075239 | 0.254423 |

**Fig.8: It can be shown that the Model 2 Result is superior to the Model 1 Result**.

## CONCLUSION

We integrated the principles of Visual Labeling with Autonomous Translation by Machines to construct a neural network architecture. Our developed model demonstrates a sufficient level of capability in decoding visual input and generating a relevant representation is natural language. Two potential methods to enhance the current algorithm and its prospective predictions include including supplementary "CNN" layers or completing more comprehensive pre-training.

## REFERENCES

[1]     "Russakovsky," "Olga," "Jia Deng," "Hao Su," "Jonathan Krause," "Sanjeev Satheesh," "Sean Ma," "Zhiheng Huang," "Andrej Karpathy," "Aditya Khosla," "Michael Bernstein," "Alexander C. Berg," and "Li Fei-Fei." "ImageNet LargeScale Visual Recognition Challenge." "International Journal of Computer Vision Int J Comput Vis 115." "3(2015):211-52. Web. 19 Apr. 2016."

[2]     "Farhadi A. et al." (2010). "Every Picture Tells a Story: Generating Sentences from Images." "Daniilidis K.," "Maragos P.," "Paragios N." (eds) "Computer Vision – ECCV 2010." "Lecture Notes in Computer Science,"

[3]     "Everingham," "Mark," "Luc Van Gool," "Christopher K. I. Williams," "John Winn," and "Andrew Zisserman." "ThePascal Visual Object Classes (VOC) Challenge." "International Journal of Computer Vision Int J Comput Vis88." 2(2009): 303-38. Web. 22 May 2016.

[4]     "Zhongliang Yang," "Yu-Jin Zhang," "Sadaqat ur Rehman," "Yongfeng Huang," "Image Captioning with Object Detectionand Localization,"

[5]     "T.-Y. Lin," et al (2014) "Microsoft COCO: Common objects in context." "arXiv:1405.0312."

[6]     "Andrej Karpathy," "Li Fei-Fei," "Deep Visual Semantic Alignments for Generating Image Descriptions," [Online]

[7]     Karpathy,Andrej,andLiFei-Fei." Deep Visual-semantic Alignments for Generating Image Descriptions." 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2015). Web. 29 May 2016.

[8]     Kiros Ryan, Rich Zemel, and Ruslan Salakhutdinov. "Multimodal neural language models." Proceedings of the 31st International Conference on Machine Learning (ICML-14): 595-603 (2014). Web. 21 May 2016.

[9]     R. Kiros, R.Salakhutdinov, and R.S. Zemel. Unifying visual -semantic embeddings with multimodal neural language models. In arXiv:1411.2539,2014

[10]   R.Senthamil Selvan, "Tumor Infiltration of Microrobot using Magnetic torque and AI Technique" by 2023 2nd International Conference on Vision Towards Emerging Trends in Communication and Networking Technologies (ViTECoN), ISSN:0018-9219,E-ISSN:1558-2256,26 June 2023, 10.1109/ViTECoN58111.2023.10157336.

[11]   R.Senthamil Selvan "Automatic Liver Cancer Detection in Abdominal Liver Images Using Soft Optimization Techniques by 2022 International Conference on Knowledge Engineering and Communication Systems (ICKES),ISSN:0018-9219,E-ISSN:1558-2256,17 March 2023, 10.1109/ICKECS56523.2022.10060747.

[12]   R.Senthamil Selvan "Analysis of Alzheimer Disease With K means Algorithm And PSO Segmentation" by 2022 IEEE 2nd Mysore Sub Section International Conference (MysuruCon) , ISSN:0018-9219,E-ISSN:1558-2256,13 December 2022, 10.1109/MysuruCon55714.2022.9972409.