

## Image captioning using CNN and RNN

<sup>1</sup>Mrs. NVN Sowjanya, <sup>2</sup>Anthati Aakash, <sup>3</sup>Ganimena Gangadher, <sup>4</sup>M. Nithin kumar, <sup>5</sup>V. Vijay Reddy

<sup>1</sup>Assistant Professor, Dept.of CSE, Teegala Krishna Reddy Engineering College, Meerpet, Hyderabad,

[sowjanya.nvn@tkrec.ac.in](mailto:sowjanya.nvn@tkrec.ac.in)

<sup>2</sup>BTech student, Dept.of CSE, Teegala Krishna Reddy Engineering College, Meerpet, Hyderabad,

[anthatiaakash236@gmail.com](mailto:anthatiaakash236@gmail.com)

<sup>3</sup>BTech student, Dept.of CSE, Teegala Krishna Reddy Engineering College, Meerpet, Hyderabad,

[gangadharsanjeev80@gmail.com](mailto:gangadharsanjeev80@gmail.com)

<sup>4</sup>BTech student, Dept.of CSE, Teegala Krishna Reddy Engineering College, Meerpet, Hyderabad,

[nithinkumar1379@gmail.com](mailto:nithinkumar1379@gmail.com)

<sup>5</sup>BTech student, Dept.of CSE, Teegala Krishna Reddy Engineering College, Meerpet, Hyderabad,

[vijayreddy.vanga8@gmail.com](mailto:vijayreddy.vanga8@gmail.com)

**Abstract:** *Image captioning is a fast-growing research field of computer vision and natural language processing that involves creating text explanations for images. It is a concept of gathering the right description of the given image on the internet use Computer Vision and natural language processing. It is achieved using the Deep learning techniques called as convolution neural network and recurrent neural network. The dataset used for implementation is called as the Flickr8k Dataset. The task of image captioning can be divided into two modules logically – one is an image-based model - which extracts the features and nuances out of our image, and the other is a language-based model – which translates the features and objects given by our image-based model to a natural sentence.*

**Keywords:** *Image captioning, convolutional neural network, recurrent neural networks, deep learning.*

### I. INTRODUCTION

Image Caption generation is a task in which a machine model is trained using artificial Intelligence in a way that the machine can understand the Image scene at

a same level as human beings understand the visual world. Image Captioning is basically like a short description generated by just looking at the image visually. In this task, a machine is fed with an input image and based on the intelligence and

training given, the model generates a simple caption which indeed explains the content of the image in a human readable form. This task is a Supervised learning algorithm example. Such task become more challenging when a machine must generate a caption for unseen or not trained images. Generally, a model tries to break an image down into objects and classify these objects before generating sentence or caption.

Captioning of image basically aims towards generating natural language and simple captions which describes the image content accurately. In this task, all objects and their relationship should be depicted precisely. A traditional algorithm which is a combination of Convolutional and Recurrent network used for generating captions has many problems such as gradient vanishing, not so accurate identification of objects and their relationship or generation of captions only for seen images, etc.

An Automatic Image Captioning Model which is a combination of advanced Convolutional and Long Short-Term Memory Deep Neural Network algorithms (CNN and LSTM) is a variation of traditional method to overcome the problems that arises using traditional way of captioning. The Model is divided into two stages: First stage uses Convolutional algorithm and second stage uses Long

Short-Term Memory. The input to the first stage is image/picture. The proposed system model also focuses on the informative captions that best describes the image scene.

In the proposed system model, the first stage known as Encoder stage is feed with image vector where the image is already pre-processed and then given as input to Stage 1. At this stage, various convolutional layers are applied on the vector which fetches appropriate features from the provided vector before sending it to next stage. After applying number of convolutional layers/operations on the image vector, it is then sent to next stage which is Decoder stage. Stage 2 processes the image vector given by the Stage 1 in a linear way to generate captions. The methodology uses LSTM algorithm in Stage 2 which is advanced version of recurrent neural network (RNN) helps to overcome the gradient explosion problem. LSTM has an advantage as its various memory gates which decides the flow of the information in the Stage 2. It also has an advantage to retain the data for longer period of time and dependencies. This Stage 2 outputs a sequential decoded simple language sentence or captions for the given input image.

This project uses advanced methods of computer vision using Deep Learning and natural language processing using a

Recurrent Neural Network. Deep Learning is a machine learning technique with which we can program the computer to learn complex models and understand patterns in a large dataset. The combination of increasing computation speed, wealth of research and the rapid growth of technology. Deep Learning and AI is experiencing massive growth worldwide and will perhaps be one of the world's biggest industries in the near future. The 21st century is the birth of AI revolution, and data becoming the new 'oil' for it. Every second in today's world large amount of data is being generated. We need to build models that can study these datasets and come up with patterns or find solution for analysis and research. This can be achieved solely due to deep learning. Computer Vision is a cutting-edge field of computer science that aims to enable computers to understand what is being seen in an image. Computers don't perceive the world like humans do. For them the perception is just sets of raw numbers and because of several limitations like type of camera, lighting conditions, clarity, scaling, viewpoint variation etc. make computer vision so hard to process as it is very tough to build a robust model that can work on every condition. The neural network architectures normally we see were trained using the current inputs only. While developing the system, the

generating output does not consider the previous inputs. It is because of neglecting any memory elements present. That is why the use of RNN tackles the memory issues that haunt the system. This led us to create an efficient system.

## II. LITERATURE REVIEW

In method proposed by Liu, Shuang & Bai, Liang & Hu, Yanli & Wang, Haoran et al. [1], two models of deep learning namely, Convolutional Neural Network-Recurrent Neural Network (CNN-RNN) Based Image Captioning, Convolutional Neural Network-Convolutional Neural Network (CNN-CNN)Based Image Captioning. In CNN-RNN Based frame work, Convolutional Neural Networks for encoding and Recurrent Neural Networks for the decoding process. Using CNN, the images here are converted to vectors and these vectors are called image features these are passed into Recurrent neural networks as input. In RNN's se NLTK libraries are used to get the actual captions for the project. In the CNN-CNN based frame work only CNN is used for both encoding and decoding of the images. Here vocab dictionary is used and it is mapped with Image features to get the exact word for the given image using NLTK library. Thus, generating the error free caption. Consisting of many models that are given at the same time of convolution techniques simultaneously is

certainly quicker compared to the train the continuous flowing recurrently repetitions of these techniques. CNN-CNN Model has less training time as compared to the CNN-RNN Model. The CNN-RNN Model has more training time as it is sequential but it has less loss compared to the CNN-CNN Model.

In the method proposed by Ansari Hani et al [2] Here they have used encoding decoding model for image captioning. Here they have mentioned two more models for image captioning they are: Retrieval based captioning and template-based captioning. Retrieval based captioning is the process where training images are placed in one space and their corresponding captions which are generated are placed in another scope now in the new scope the correlations are calculated for the test image and captions the highest valued correlation caption is retrieved as caption for the given image from the given set of captions dictionary. Prototype based describing is the technique is done by them in this paper. Here they have used Inception V3 model as their encoder and they have used attention mechanism and GRU as their decoder to generate the captions.

In the method proposed by Subrata Das, Lalit Jain et al [3] This model is mainly based on how the deep learning models are used for Military Image captioning. It

mainly uses CNN RNN based frame work. They have used Inception model for encoding the images and to decrease the gradient descent problem they have used Long Short-Term Memory (LSTM'S) Networks.

In the method proposed by G Geetha et al [4] they have used CNN-LSTM model for image captioning. The entire flow of the model was explained from data set collection to caption generation. Here Convolutional Neural Networks was used as encoder and LSTMs was used as decoder for generating the captions.

Image captioning approach mechanically producing a caption for a photograph. As a lately emerged research place, its miles attracting more and more interest. To obtain the motive of picture captioning, semantic facts of pictures desires to be captured and expressed in natural languages. Connecting both research communities of computer vision and natural language processing, photograph captioning is a pretty tough challenge. Various methods have been proposed to treatment this hassle. A survey on advances in photo captioning research is given. Based on the technique accompanied the photograph captioning procedures are categorised into distinct training. Representative strategies in each class are summarized, and their strengths and boundaries are referred to [5].

In this the use of a knowledge graphs that capture widespread or common-sense knowledge, to reinforce the facts extracted from pics by the state-of- the-artwork techniques for image captioning is explored. The outcomes of the experiments, on several benchmark statistics sets inclusive of MS COCO, as measured by using CIDEr-D, a overall performance metric for picture captioning, display that the variants of the kingdom-of- the-art techniques for photo captioning that make use of the facts extracted from expertise graphs can appreciably outperform people who depend solely on the data extracted from snap shots [6].

Automatically generating a herbal language description of an photo is a challenge close to the heart of photo understanding. In this paper, a multi-version neural community approach intently associated with the human visual system that routinely learns to describe the content material of snap shots is presented. The version includes two sub-models: an item detection and localization model, which extract the statistics of gadgets and their spatial dating in pictures respectively; Besides, a deep recurrent neural network (RNN) based on lengthy quick-time period memory (LSTM) gadgets with attention mechanism for sentences era[7].

In latest years giant development has been made in photograph captioning, the use of

Recurrent Neural Networks powered by means of lengthy-short-time period reminiscence (LSTM) devices. Despite mitigating the vanishing gradient problem, and in spite of their compelling potential to memorize dependencies, LSTM units are complicated and inherently sequential throughout time. However, the complex addressing and overwriting mechanism blended with inherently sequential processing, and big garage required because of back-propagation thru time (BPTT), poses demanding situations at some point of education[8].

### III. PROPOSED SYSTEM

Our model uses two different neural networks to generate the captions. The first neural network is Convolutional Neural Network (CNN), which is used to train the images as well as to detect the objects in the image with the help of various pre-trained models

like VGG, Resnet50, Inception or YOLO. The second neural network used is Recurrent Neural Network (RNN) based Long Short Term Memory(LSTM), which is used to generate captions from the generated object keywords. As, there is lot of data involved to train and validate the model, generalized machine learning algorithms will not work. Deep Learning has been evolved from the recent times to solve the data constraints on Machine

Learning algorithms. GPU based computing is required to perform the Deep Learning tasks more effectively.

### **Deep Learning**

Deep learning is a subset of machine learning, which is essentially a neural network with three or more layers. These neural networks attempt to simulate the behaviour of the human brain—albeit far from matching its ability—allowing it to “learn” from large amounts of data. While a neural network with a single layer can still make approximate predictions, additional hidden layers can help to optimize and refine for accuracy.

Deep learning drives many artificial intelligence (AI) applications and services that improve automation, performing analytical and physical tasks without human intervention. Deep learning technology lies behind everyday products and services (such as digital assistants, voice-enabled TV remotes, and credit card fraud detection) as well as emerging technologies (such as self-driving cars).

Deep learning neural networks, or artificial neural networks, attempts to mimic the human brain through a combination of data inputs, weights, and bias. These elements work together to accurately recognize, classify, and describe objects within the data. Deep neural networks consist of multiple layers of interconnected nodes, each building upon the previous layer to

refine and optimize the prediction or categorization. This progression of computations through the network is called forward propagation. The input and output layers of a deep neural network are called visible layers. The input layer is where the deep learning model ingests the data for processing, and the output layer is where the final prediction or classification is made.

Another process called backpropagation uses algorithms, like gradient descent, to calculate errors in predictions and then adjusts the weights and biases of the function by moving backwards through the layers to train the model. Together, forward propagation and backpropagation allow a neural network to make predictions and correct for any errors accordingly. Over time, the algorithm becomes gradually more accurate.

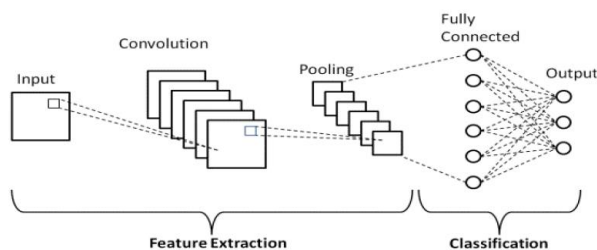
The above describes the simplest type of deep neural network in the simplest terms. However, deep learning algorithms are incredibly complex, and there are different types of neural networks to address specific problems or datasets. For example, Convolutional Neural Networks and Recurrent Neural Networks.

### **Convolutional Neural Networks (CNN)**

Convolutional Neural Network (ConvNet/CNN) is a Deep Learning algorithm that can take in an input image, assign importance (learnable weights and

biases) to various aspects/objects in the image, and be able to differentiate one from the other. The pre-processing required in a ConvNet is much lower as compared to other classification algorithms. While in primitive methods filters are hand-engineered, with enough training, ConvNets have the ability to learn these filters/characteristics. The architecture of a ConvNet is analogous to that of the connectivity pattern of Neurons in the Human Brain and was inspired by the organization of the Visual Cortex. Individual neurons respond to stimuli only in a restricted region of the visual field known as the Receptive Field. A collection of such fields overlaps to cover the entire visual area.

A ConvNet is able to successfully capture the Spatial and Temporal dependencies in an image through the application of relevant filters. The architecture performs a better fitting to the image dataset due to the reduction in the number of parameters involved and the reusability of weights. In other words, the network can be trained to understand the sophistication of the image better.



**Fig.1 Architecture of CNN**

**Convolutional Layer**

The convolution layer is the core building block of the CNN [7]. It carries the main portion of the network’s computational load. This layer performs a dot product between two matrices, where one matrix is the set of learnable parameters otherwise known as a kernel, and the other matrix is the restricted portion of the receptive field. The kernel is spatially smaller than an image but is more in-depth. This means that, if the image is composed of three (RGB) channels, the kernel height and width will be spatially small, but the depth extends up to all three channels. If we have an input of size  $W \times W \times D$  and  $D_{out}$  number of kernels with a spatial size of  $F$  with stride  $S$  and amount of padding  $P$ , then the size of output volume can be determined by the following formula

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

**Pooling Layer**

The pooling layer replaces the output of the network at certain locations by deriving a summary statistic of the nearby outputs. This helps in reducing the spatial size of the representation, which decreases the required amount of computation and weights. The pooling operation is processed on every slice of the representation individually. If there is no

pooling, the output has the same resolution as the input.

The following are some methods for pooling:

**Max-pooling:** It chooses the most significant element from the feature map. The feature map's significant features are stored in the resulting max-pooled layer. It is the most popular method since it produces the best outcomes.

**Average pooling:** It entails calculating the average for each region of the feature map.

### **Fully Connected Layer**

At the end of CNN, there is a Fully connected layer of neurons. As in conventional Neural Networks, neurons in a fully connected layer have full connections to all activations in the previous layer and work similarly. After training, the feature vector from the fully connected layer is used to classify images into distinct categories. Every activation unit in the next layer is coupled to all of the inputs from this layer. Overfitting occurs because all of the parameters are occupied in the fully-connected layer

### **Recurrent Neural Networks (RNN)**

A recurrent neural network (RNN) is a type of artificial neural network which uses sequential data or time series data. These deep learning algorithms are commonly used for ordinal or temporal problems, such as language translation, natural language processing (nlp), speech

recognition, and image captioning; they are incorporated into popular applications such as Siri, voice search, and Google Translate. Like feedforward and convolutional neural networks (CNNs), recurrent neural networks utilize training data to learn.

They are distinguished by their “memory” as they take information from prior inputs to influence the current input and output. While traditional deep neural networks assume that inputs and outputs are independent of each other, the output of recurrent neural networks depend on the prior elements within the sequence. While future events would also be helpful in determining the output of a given sequence, unidirectional recurrent neural networks cannot account for these events in their predictions

### **Long Short-Term Memory (LSTM)**

In sequence prediction challenges, Long Short Term Memory (LSTM) networks are a type of Recurrent Neural Network that can learn order dependence. The output of the previous step is used as input in the current step in RNN. Hochreiter & Schmid Huber created the LSTM. It addressed the issue of RNN long-term dependency, in which the RNN is unable to predict words stored in long-term memory but can make more accurate predictions based on current data.



The LSTM is made up of four neural networks and numerous memory blocks known as cells in a chain structure. A conventional LSTM unit consists of a cell, an input gate, an output gate, and a forget gate. The flow of information into and out of the cell is controlled by three gates, and the cell remembers values over arbitrary time intervals. The LSTM algorithm is well adapted to categorize, analyse, and predict time series of uncertain duration.

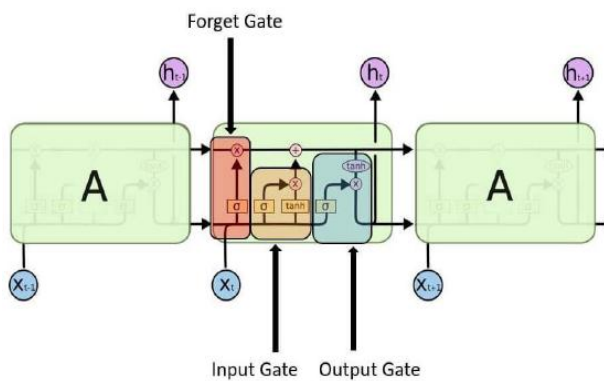


Fig.2 Structure of LSTM

The cells store information, whereas the gates manipulate memory. There are three entrances:

**Input Gate:** It determines which of the input values should be used to change the memory. The sigmoid function determines whether to allow 0 or 1 values through. And the tanh function assigns weight to the data provided, determining their importance on a scale of -1 to 1.

$$i_t = \sigma(w_i \cdot [h_{t-1}, x_t] + b_i)$$

$$C_t = \tanh(w_c \cdot [h_{t-1}, x_t] + b_c)$$

**Forget Gate:** It finds the details that should be removed from the block. It is decided by a sigmoid function. For each number in the cell state  $C_{t-1}$ , it looks at the preceding state ( $h_{t-1}$ ) and the content input ( $X_t$ ) and produces a number between 0 (omit this) and 1 (keep this).

$$f_t = \sigma(w_f \cdot [h_{t-1}, x_t] + b_f)$$

**Output Gate:** The block's input and memory are used to determine the output. The sigmoid function determines whether to allow 0 or 1 values through. And the tanh function determines which values are allowed to pass through 0, 1. And the tanh function assigns weight to the values provided, determining their relevance on a scale of -1 to 1 and multiplying it with the sigmoid output.

$$o_t = \sigma(w_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

The recurrent neural network uses long short-term memory blocks to provide context for how the software accepts inputs and creates outputs. Because the program uses a structure based on short-term memory processes to build longer-term memory, the unit is dubbed a long short-term memory block. In natural language processing, these systems are extensively used.

A sequence of repeating neural network modules makes up all recurrent neural networks. This repeating module in

traditional RNNs will have a simple structure, such as a single tanh layer. The output of the current time step becomes the input for the following time step, which is referred to as Recurrent. At each element of the sequence, the model examines not just the current input, but also what it knows about the prior ones.

The LSTM cycle is divided into four steps:

- Using the forget gate, information to be forgotten is identified from a prior time step.
- Using input gate and tanh, new information is sought for updating cell state.
- The information from the two gates above is used to update the cell state.
- The output gate and the squashing operation provide useful information.

**IV. RESULTS**

**SYSTEM ARCHITECTURE**

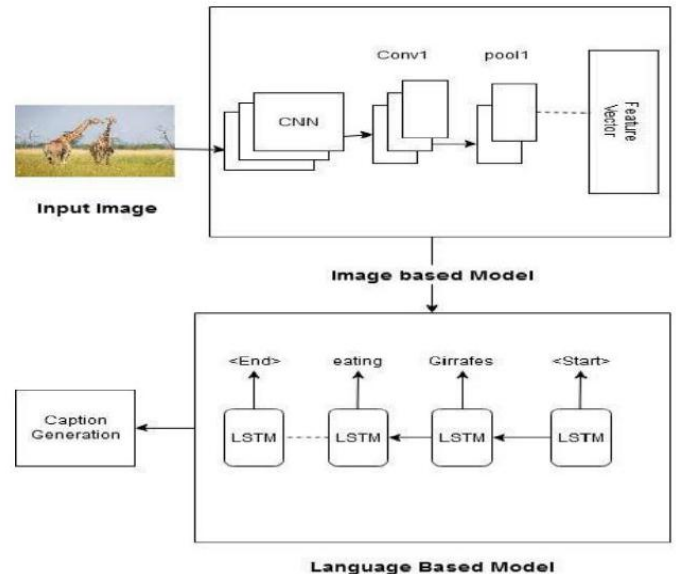


Fig.3 System architecture

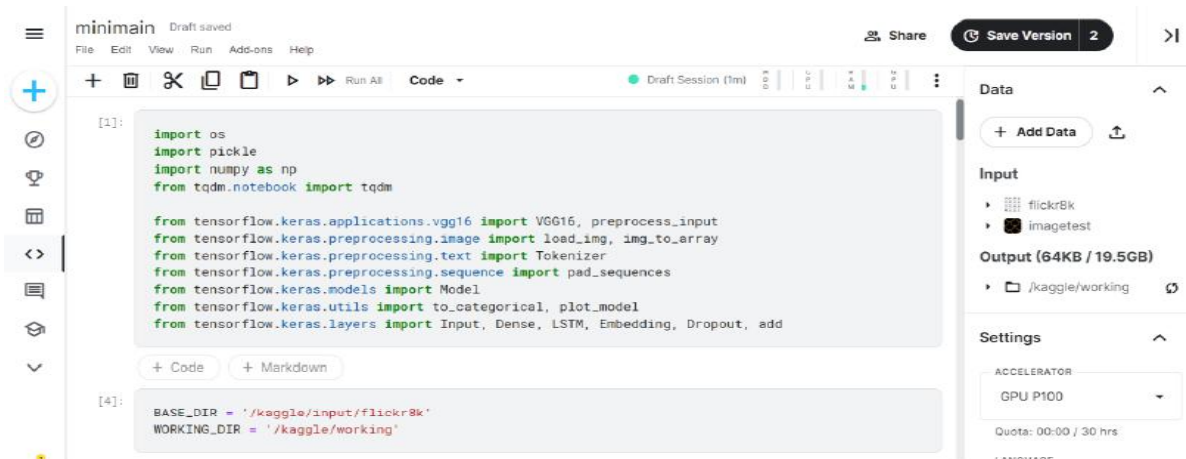


Fig.4 Kaggle Interface

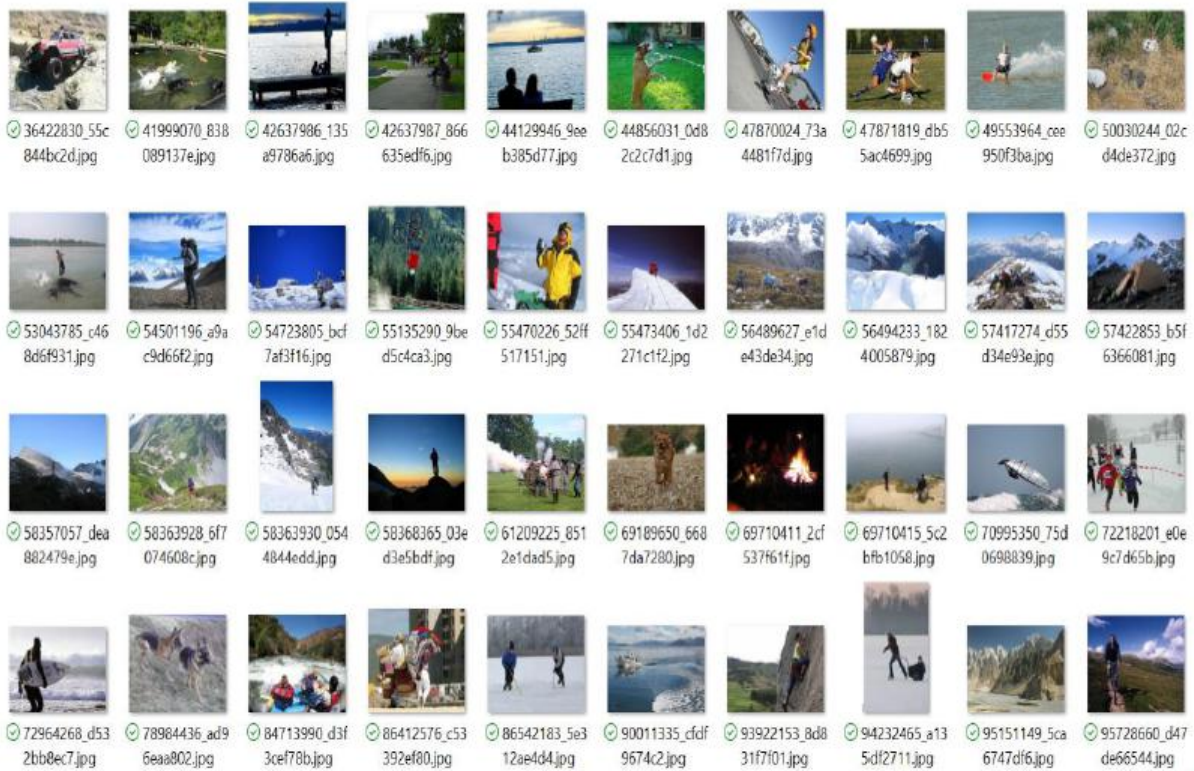


Fig.5 Flickr8k Dataset Images

1305564994_00513f9a5b.jpg#0	A man in street racer armor be examine the tire of another racer 's motorbike
1305564994_00513f9a5b.jpg#1	Two racer drive a white bike down a road .
1305564994_00513f9a5b.jpg#2	Two motorist be ride along on their vehicle that be oddly design and color .
1305564994_00513f9a5b.jpg#3	Two person be in a small race car drive by a green hill .
1305564994_00513f9a5b.jpg#4	Two person in race uniform in a street car .
1351764581_4d4fb1b40f.jpg#0	A firefighter extinguish a fire under the hood of a car .
1351764581_4d4fb1b40f.jpg#1	A fireman spray water into the hood of small white car on a jack
1351764581_4d4fb1b40f.jpg#2	A fireman spray inside the open hood of small white car , on a jack .
1351764581_4d4fb1b40f.jpg#3	A fireman use a firehose on a car engine that be up on a carjack .
1351764581_4d4fb1b40f.jpg#4	Firefighter use water to extinguish a car that be on fire .
1358089136_976e3d2e30.jpg#0	A boy sand surf down a hill
1358089136_976e3d2e30.jpg#1	A man be attempt to surf down a hill make of sand on a sunny day .
1358089136_976e3d2e30.jpg#2	A man be slide down a huge sand dune on a sunny day .
1358089136_976e3d2e30.jpg#3	A man be surf down a hill of sand .
1358089136_976e3d2e30.jpg#4	A young man in short and t-shirt be snowboard under a bright blue sky .
1362128028_8422d53dc4.jpg#0	kid play in a blue tub full of water outside
1362128028_8422d53dc4.jpg#1	On a hot day , three small kid sit in a big container fill with water .
1362128028_8422d53dc4.jpg#2	Little kid sit outdoors in a small tub of water .
1362128028_8422d53dc4.jpg#3	Three child squeeze into a plastic tub fill with water and play .
1362128028_8422d53dc4.jpg#4	Three little boy take a bath in a rubber bin on the grass .
1383698008_8ac53ed7ec.jpg#0	A man be snowboard over a structure on a snowy hill .
1383698008_8ac53ed7ec.jpg#1	A snowboarder jump through the air on a snowy hill .
1383698008_8ac53ed7ec.jpg#2	a snowboarder wear green pants do a trick on a high bench
1383698008_8ac53ed7ec.jpg#3	Someone in yellow pants be on a ramp over the snow .
1383698008_8ac53ed7ec.jpg#4	A man be perform a trick on a snowboard high in the air .
1468103286_96a6e07029.jpg#0	A Baseball batter raise his arm .
1468103286_96a6e07029.jpg#1	A baseball player from New York wait to bat during a game .
1468103286_96a6e07029.jpg#2	A baseball player in a Yankee uniform be hold a bat in one hand
1468103286_96a6e07029.jpg#3	A New York Yankee hold up a bat .
1468103286_96a6e07029.jpg#4	New York Yankee warm up .
1478268555_7e301fc510.jpg#0	A bare backed climber be attach to the rock face on a pink safety rope .
1478268555_7e301fc510.jpg#1	a lone rock climber in a harness climb a huge rock wall .

Fig.6 Description

## Image Captioning using CNN and LSTM

caption generation

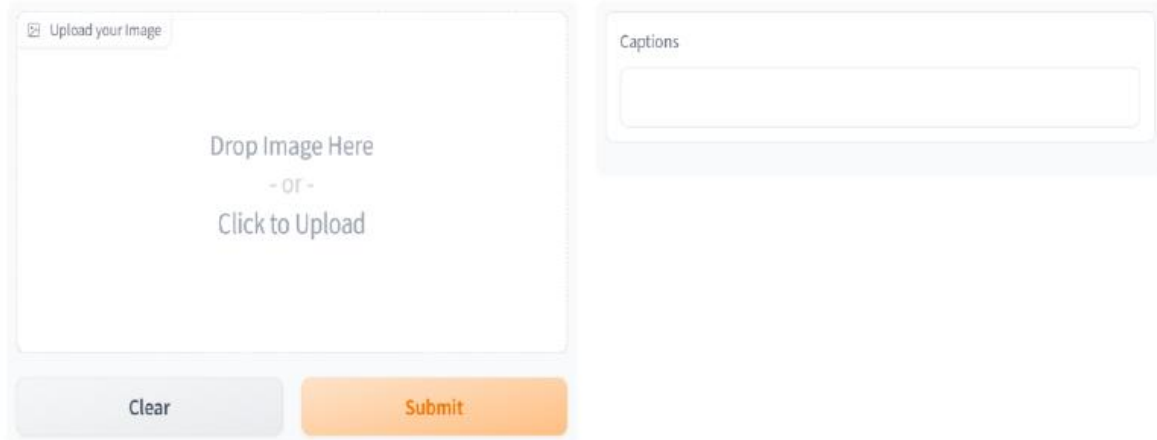


Fig.7 Screen Before uploading image

caption generation

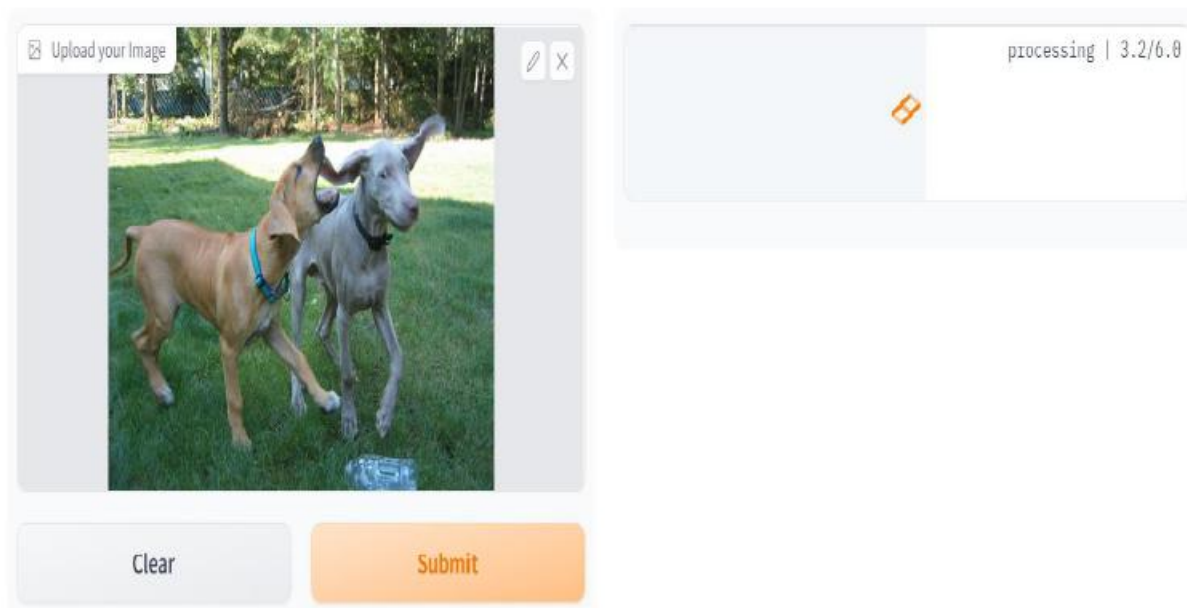


Fig.8 Image capturing using CNN and LSTM

caption generation

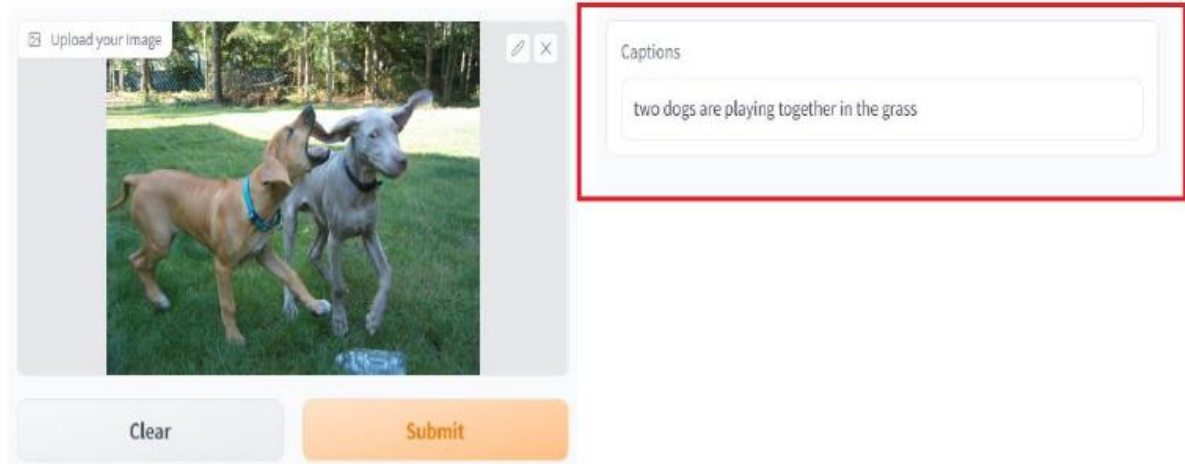


Fig.9 Screen showing the caption for given image

caption generation

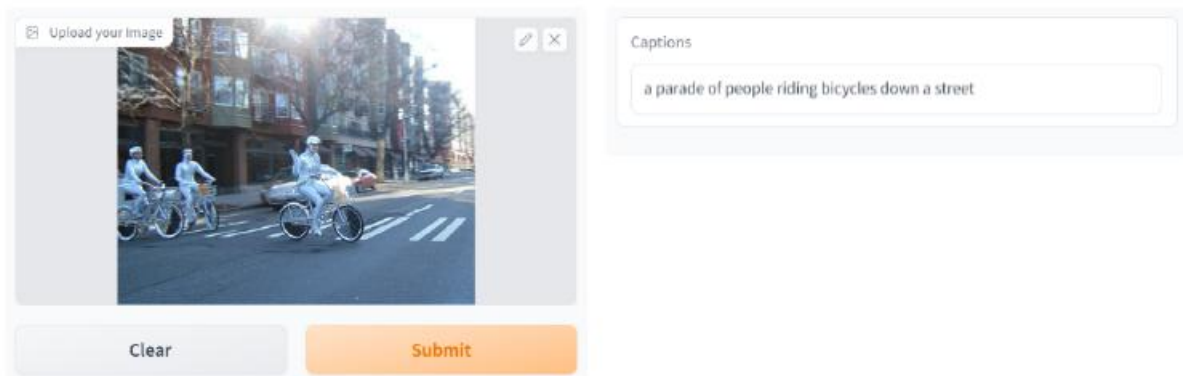


Fig.10 caption generated for example image 1

caption generation

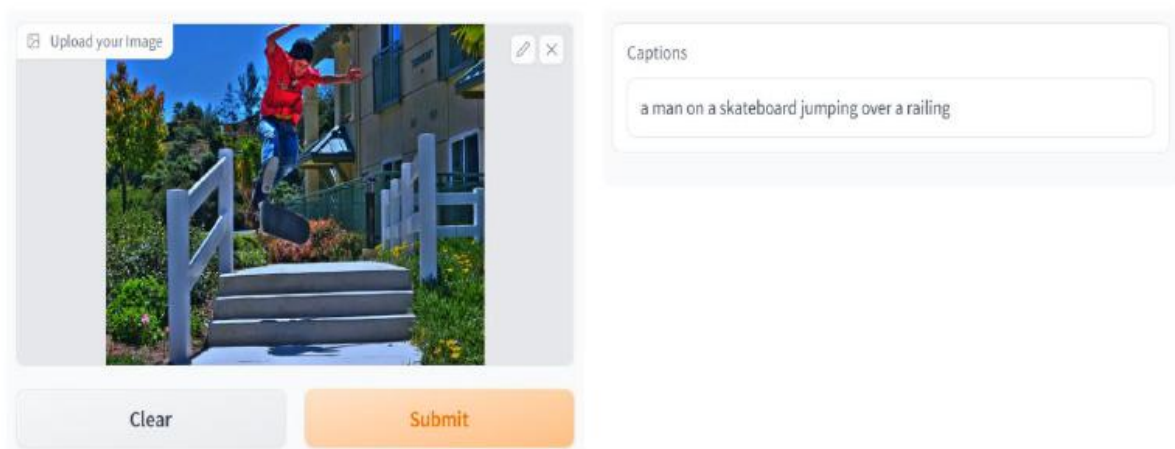


Fig.11 caption generated for example image 2

## V. CONCLUSION

We developed an Automatic Image Captioning Model which is a combination of CNN and advanced RNN: LSTM for generating qualitative and accurate captions which could describe the image in natural and easy language. The Proposed Model is trained with 6000 Images using Flickr8k dataset which contained Images along with its captions. The proposed Convolutional deep neural network extracts the important features from image and stores it in feature vector. The feature vector is sent to LSTM model to generate a sequential sentence combining the extracted features and their relationship to form a caption.

The proposed Model generates precise captions for the Image. The model has reduced the error rate in the caption. The proposed system used LSTM algorithm to overcome the gradient vanishing problem of traditional RNN algorithm. The System is tested with 2000 Flickr8k dataset images.

The System is accurately able to identify the objects in the images and their relationship. In Future, this proposed model can be extended where a system can be trained using images, its caption and also descriptions which helps improve the caption accuracy.

## REFERENCES

- [1] Liu, Shuang & Bai, Liang & Hu, Yanli & Wang, Haoran. (2018). Image Captioning Based on Deep Neural Networks. MATEC Web of Conferences. 232. 01052. [10.1051/mateconf/201823201052](https://doi.org/10.1051/mateconf/201823201052).
- [2] A. Hani, N. Tagougui and M. Kherallah, "Image Caption Generation Using a Deep Architecture," 2019 International Arab Conference on Information Technology (ACIT), 2019, pp. 246-251, [doi: 10.1109/ACIT47987.2019.8990998](https://doi.org/10.1109/ACIT47987.2019.8990998).
- [3] S. Das, L. Jain and A. Das, "Deep Learning for Military Image Captioning," 2018 21st International Conference on Information Fusion (FUSION), 2018, pp. 2165-2171, [doi: 10.23919/ICIF.2018.8455321](https://doi.org/10.23919/ICIF.2018.8455321).
- [4] G Geetha, T.Kirthigadevi, G GODWIN Ponsam, T.Karthik, M.Safa," Image Captioning Using Deep Convolutional Neural Networks(CNNs)" Published under licence by IOP Publishing Ltd in Journal of Physics :Conference Series ,Volume 1712, International Conference On Computational Physics in Emerging Technologies(ICCPET) 2020 August 2020,Manglore India in 2015.
- [5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and

translate. In International Conference on Learning Representations (ICLR).

[6] Shuang Bai and Shan An. 2018. A Survey on Automatic Image Caption Generation Neurocomputing.

[7] Andreas Bulling, Jamie A Ward, Hans Gellersen, and Gerhard Troster. 2011. Eye movement analysis for activity recognition using electrooculography. IEEE transactions on pattern analysis and machine intelligence 33, 4 (2011), 741–753.

[8] Tseng-Hung Chen, Yuan-Hong Liao, Ching-Yao Chuang, Wan-Ting Hsu, Jianlong Fu, and Min Sun. 2017. Show, Adapt and Tell: Adversarial Training of Cross- Domain Image Captioner. In The

IEEE International Conference on Computer Vision (ICCV), Vol. 2.

[9] Raffaella Bernardi, Ruket Cakici, Desmond Elliott, Aykut Erdem, Erkut Erdem, Nazli Ikizler-Cinbis, Frank Keller, Adrian Muscat, Barbara Plank, et al.2016.

[10] Akshi Kumar and Shivali Goel. 2017. A survey of evolution of image captioning techniques. International Journal of Hybrid Intelligent Systems Preprint, 1–19

[11] Prasadu Peddi (2019), "Data Pull out and facts unearthing in biological Databases", International Journal of Techno-Engineering, Vol. 11, issue 1, pp: 25-32.