

Generating Image Captions based on Deep Learning and Natural Language Processing

Ms. Rohini P

Asst. Professor of Department Of Computer Science and Engineering
Srinivasa Ramanujan Institute of Technology, Anantapur
rohini.cse@srit.ac.in

Bhargavi M

Department Of Computer Science & Engineering
Srinivasa Ramanujan Institute of Technology
204g1a0521@srit.ac.in

Jasmin G

Department Of Computer Science & Engineering
Srinivasa Ramanujan Institute of Technology
204g1a0542@srit.ac.in

Manjusha P

Department Of Computer Science & Engineering
Srinivasa Ramanujan Institute of Technology
204g1a0552@srit.ac.in

Anulekha Sai A

Department Of Computer Science & Engineering
Srinivasa Ramanujan Institute of Technology
214g5a0504@srit.ac.in

Abstract- Humans and computers are attempting to communicate because everything in today's society depends on systems like computers, mobile phones, etc. This is how our project is visualized. Our undertaking People with visual impairments can benefit from the creation of image captions. Computers are unable to distinguish objects, things, or activities with the same ease as humans. To recognize them, they require some training. The suggested method is used to identify activities or similar items. We offer several deep neural network-based models for creating captions for images, with a particular emphasis on CNNs (Convolutional Neural Networks) that extract characteristics from the image. Using LSTM (Long Short-Term Memory) techniques, RNNs (Recurrent Neural Networks) create captions based on the image's attributes. and examining how they affect the construction of sentences. Here, encoder-decoders are used to create a link between descriptions from natural language processing and visual information such as image features. The process of generating a caption's sequence is handled by the decoder, while the encoder extracts features. In order to determine which feature extraction and encoder model produces the best results and accuracy, we have also created captions for sample photos and compared them with one another. We also introduce Deep Voice, a text-to-speech system of production quality that uses only deep neural networks to generate captions based on visual attributes. The

evaluation of our project will be conducted utilizing several machine learning methods and Python.

Keywords - CNN, RNN, LSTM, Encoder - Decoder.

1.INTRODUCTION

It is relatively easy for humans to describe the environments they are living. It is normal for a human to be able to quickly describe a vast amount of information about an image[1]. This is a fundamental human ability. The ability to identify objects and describe images is facilitated by the human brain. Artificial Intelligence introduces numerous algorithms that are based on the architecture of the brain. Here, human beings are employing these algorithms to mimic human visual world interpretation on computers. Despite significant advancements in computer vision fields including object identification, picture classification, attribute classification, and scene recognition. To allow a computer the automatically explains an image, which has been forwarded to it using a language that resembles a person, is a relatively new undertaking. The image captioning is a process that automatically produces a

description for the given image using a computer. It is a difficult task. Image captioning necessitates both a high-level comprehension of an image's semantic contents and the ability to convey the information in a sentence that sounds human because it connects computer vision with natural language processing. Giving computers access to the world visually, this will bring a wide range of applications, including the creation of human-robot interactions, early childhood education, information retrieval, and support for those who are visually impaired. It is a significantly new endeavor to make a computer automatically explain an image that is forwarded to it in a language that sounds human.

The mechanism of automatically generating a natural language caption for a picture with a computer known as image captioning. It's a challenging task. Because picture captioning combines computer vision and natural language processing, that requires both a high level of knowledge of an image's semantic contents and the ability to transmit the information in a human-sounding sentence. Full of reasonable applications, will become possible if computers are given access to the visual world.

Convolutional neural networks (CNN) algorithm is mainly used to extracting features from the picture, while recurrent neural networks (RNN) are used to generate sequence of words in a meaningful way which may easily understandable to human beings based on the image. During first stage, we have taken a novel technique to extracting features from a picture, which will provide us with details on even the smallest change between two comparable photos, instead of just detecting the objects present in the image. The 16 convolutional layer model VGG-16 (Visual Geometry Group) has been employed by us for object recognition.

We must train the dataset's captions in the second stage. In order to construct our captions from the input photos, we use GRU (Gated Recurrent Unit) and LSTM (Long Short-Term Memory). We have measured the accuracy of different algorithms using the BLEU (Bilingual Evaluation Understudy) in order to estimate which architecture is best. For evaluating the caliber of translations produced by machines, BLUE offers a numerical score. Our goal is to achieve greater accuracy than previous efforts. We are utilizing the flickr30K dataset to guarantee more accuracy. Images are fed into a computer system as two-dimensional arrays, and the images are mapped to descriptive phrases or captions. The task of automatically creating image captions has attained a lot of attention in the past few years.

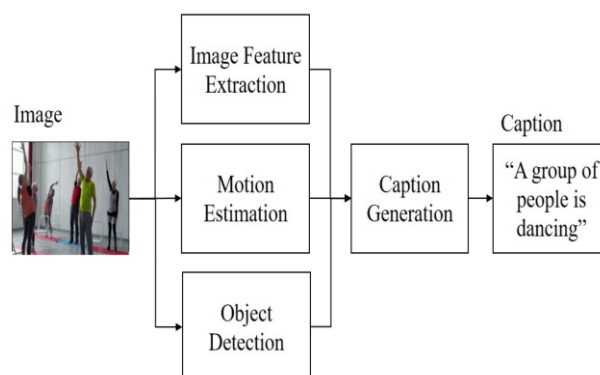


Fig1. Block Diagram

1.1. Deep Learning

Deep learning which is a subset of machine learning, was first demonstrated by using real neurons from the brain and transforming them into artificial neural networks using learning techniques.

No matter how well-optimized, machine learning approaches begin to lose accuracy and performance as data volume increases, but deep learning performs far better in these situations.

Many domains, including image classification, speech recognition, recommendation systems, NLP (natural language processing), etc., use deep learning. In light of all of these factors, we have decided to create our project using deep learning techniques. Features are directly learned by Deep Learning from the data. Neural networks automatically find and express pertinent patterns during training, eliminating the need for human feature creation.

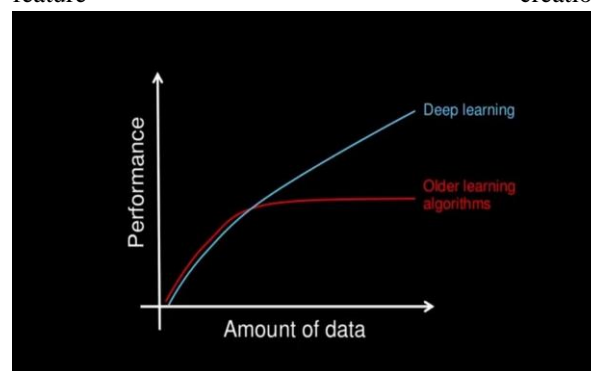


Fig2: Performance of Deep Learning

Deep learning models can generalize well to any variety of complicated datasets thanks to this capability.

The primary areas of attention for this project, "Generating Image Captions by using CNN and NLP," are feature extraction and classification. That being said, deep learning was chosen primarily for this reason.

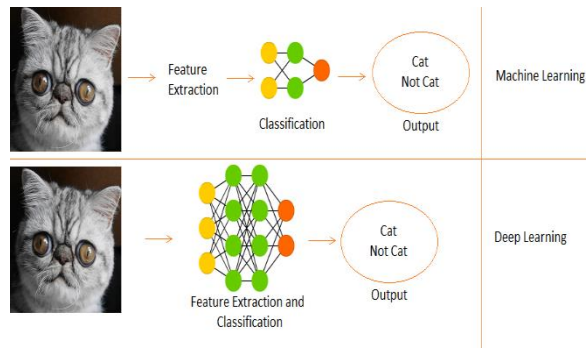


Fig3: Difference between other Machine Learning techniques and Deep Learning

The goals of this study are to produce more accurate results for appropriate caption generation based on image attributes utilizing the flickr30K dataset. Furthermore, offer an audio option for the output caption that is generated. Based on the project, this paper introduces two objectives. They are as follows:

1. To create a deep learning-based image captioning system that is more accurate than the state-of-the-art in terms of captioning by effectively training the dataset using the flickr30k dataset.
2. To use text-to-speech to turn the generated image description into audio.

2.LITERATURE SURVEY

The low accuracy of the current models is what led to the creation of this research. However, since everything is digital and has displays these days, this concept will have practical implications in the modern world. Numerous investigators devised multiple methods to guarantee precision. However, some of them fall short, while some succeed to the best of their abilities.

K Harika, Rajan Singh, M Sailaja, B Sridhar. [1] In recent years, deep neural networks have enabled the captioning of images. The photo caption generator applies an appropriate title to an applied input image based on the dataset. The present work proposes a model based on deep learning and uses it to generate a caption for the input image. The model frames a statement associated with an input image by utilizing CNN and LSTM algorithms. This CNN model recognizes the objects in the picture, and the Long Short-Term Memory (LSTM) algorithm not only generates the text but also a caption that fits the project. Thus, the primary goals of the suggested model are

object recognition and title generation for the input images.

Kalaivani S P, Karthika S, Karthika V, C S Kanimozhiselvi. [2] The practice of writing a written description for a photograph is called picture captioning. It is among the most urgent and latest research problems at the moment. Every day, new methods for resolving the problem are being presented. Even with the wide range of alternatives available, more work need to be done in order to get better and more accurate outcomes. Therefore, we considered developing an image captioning model that combines different Long Short-Term Memory and Convolutional Neural Network architectural configurations in order to get better results. Three CNN and LSTM combinations were merged in order to create the model. The proposed model is trained using three Convolutional Neural Network architectures: Inception-v3, Xception, and ResNet50. These networks are used to extract features from the image and employ Long Short-Term Memory to generate appropriate captions. The best combination of three CNN and LSTM combinations is selected based on the accuracy of the model. The model gets trained using the Flickr8k dataset.

Vaishali Jabade, Chetan Amritkar. [3] Artificial Intelligence automatically synthesizes the information of images using computer vision and natural language processing (NLP). A model of regenerating neurons is created. It depends on vision and machine translation. Using this method, natural sentences that eventually explain the image are created. This paradigm consists of recurrent neural networks (RNN) and convolutional neural networks (CNN). CNN is used for feature extraction from images, and RNNs are used to create sentences. The model is trained to generate captions that essentially describe an input image when one is provided. The accuracy and linguistic proficiency of the model learnt from picture descriptions are evaluated using a variety of datasets. These experiments show that the model frequently gives accurate descriptions of the input image.

Vaidehi Muley, Varsha Kesavan, Megha Kolhekar [4] The project aims to automatically generate captions by utilizing the content of the image as a source. At the moment, human annotation is required for photos, which makes the process almost impossible for commercial datasets. Using the picture database, the Convolutional Neural Network (CNN) encoder extracts features and subtleties from the image to produce a "thought vector."

The objects and features in the picture are then translated by an RNN (Recurrent Neural Network) decoder to create a coherent and sequential description of the image. To find the most effective model with fine-tuning, we thoroughly examine several deep neural network-based methods for creating photo captions and pretrained models in this work. They looked at models with and without the "attention" idea in order to optimize the model's ability to produce captions. All models are trained with the same dataset in order to provide a more realistic comparison.

E K Miller, T J Buschman. [5] In the human visual system, bottom-up signals linked to unexpected, unusual can automatically focus attention, as can top-down signals dictated by the present task. In this work, they use comparable language to designate as "top-down" attention mechanisms those that are driven by nonvisual or task-specific context, and a "bottom-up" attention mechanisms that are solely visual feed-forward. The majority of traditional visual attention processes utilized in VQA and picture captioning are top-down in nature. These techniques, which are often trained to selectively attend to the output of one or more layers of a convolutional neural net (CNN), use an image's caption or a query about it as their context. Nevertheless, this method pays minimal attention to the process of selecting the visual parts that require attention. It is a time taking process to pre-process the image.

Farhadi et al. [6] The three primary categories used in this paper are picture captioning techniques are covered in this section i.e., retrieval-based image captioning, novel caption creation and template-based image captioning. In these approaches, captions are generated using preset templates that contain blank spaces. These systems fill in the blanks in the templates after first identifying the various objects, actions, and characteristics. To generate image captions, for instance, complete the template slots with three distinct scene pieces.

Lakshmi Narasimhan Srinivasan, Dinesh Sreekanthan, A L Amutha. [7] The keras framework's TensorFlow backend has been utilized in this study's model evaluation. Utilizing assessment measures that were appropriate for the problem's nature allowed for an understanding of how The model has made correct predictions. This paper presents the results of mathematical computations performed on the confusion matrix.

D Elliott, F Keller [8]. The main difficulties in this research include identifying the objects in an image and their characteristics, which are challenging computer vision problems, as well as figuring out how the objects interact and what relationships exist between them. Automatic image description is not without its difficulties. To improve the model's performance, the authors trained it over several layers (or levels) using CNN.

C Amritkar and V Jabade [9]. In this paper, the model is trained to generate, upon receiving an input image, labels that nearly exactly describe it. The accuracy of the model and the command or smoothness of the language model after learning from picture descriptions are tested on multiple datasets. These tests show that the model frequently describes the input image precisely.

V Kesavan, V Muley, M Kolhekar [10]. The goal of the article is to use image content to generate captions automatically. This paper offers a thorough examination of several deep neural network-based picture captioning methods, as well as pre-trained models to help choose and optimize the best model. To improve the caption, they examined models with and without attention ideas. Producing capacity of the model. For a more accurate comparison, every model is trained using the same dataset.

Kulkarni et al. [11] Before filling in the blanks, Kulkarni uses a Conditional Random Field (CRF) to identify the things, attributes, and prepositions. Though the templates are established, template-based techniques can produce accurate outputs but not specified length sentence. They go over the three primary categories of current image in this part. Approaches for captioning images: retrieval-based captioning, template-based captioning, and creative caption creation. For the purpose of generating captions, template-based solutions use predetermined templates with blank spaces. Within these systems, the many items, actions, and after identifying the qualities, the gaps in the templates are filled.

D Vigneswari, N K Kumar, K Laxman, A Mohan, J Yuvraj [12]. This work uses deep learning to discover, recognize, and produce meaningful captions for a given image. For object identification, recognition, and caption generation, Regional Object Detector (ROD) is utilized. The proposed method makes use of deep learning to further improve the existing image caption creation

system. Python is used to conduct experiments on the Flickr 8k dataset in order to illustrate the suggested approach. We are using the best and large dataset i.e., Flickr 30k.

3.PROPOSED SYSTEM

The project's suggested system is an advanced and adaptable picture captioning solution that uses deep learning techniques in conjunction with natural language processing (NLP) to generate accurate, linguistically coherent, and contextually appropriate captions for photos that have audio accompanying them. CNN is given an image to evaluate and create a feature vector from. This feature vector enhances the user's perception of the image by providing an auditory context that corresponds with the visual content. It is used as input for sigmoid functions and RELU functions in the GRU and LSTM. It also provides image descriptions and is associated with pertinent sounds for the image.

3.1 Architecture:

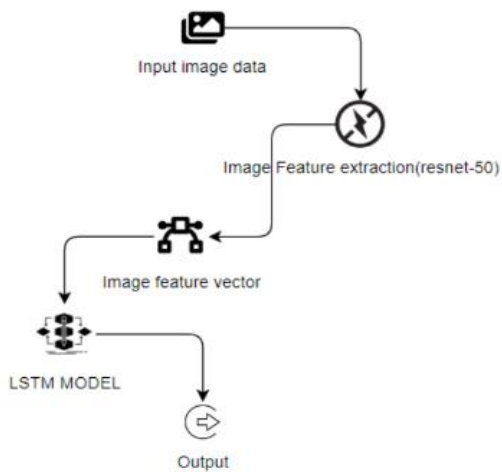


Fig4: Architecture

3.2 System :

3.2.1 Create Dataset:

The dataset including text and picture data of the target items to be captioned is split into training and testing datasets in order to evaluate the model's performance.

3.2.2 Pre-processing:

Prior to being input into a machine learning model,

images must be enhanced and prepared. To train our model, we resize and reshape the photos into the proper format.

3.2.3 Training:

We train our model utilizing the CNN and LSTM algorithms by using the pre-processed training dataset.

3.3 User:

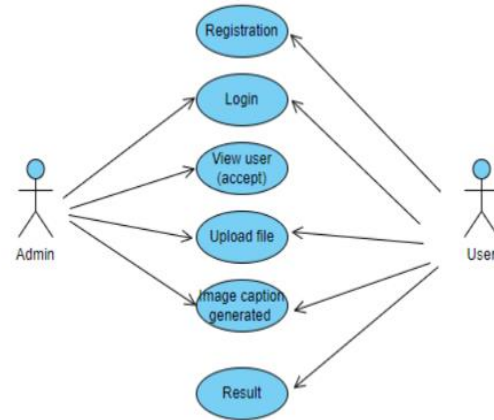


Fig5: Use Case Diagram

3.3.1 Register:

The user must register. The information is kept in a database, and the administrator will review it and grant access if the administrator determines that the registered information is accurate.

3.3.2 Admin Login:

The administrator logs in and examines the list of users who have registered. Only the user will be able to log in when he has approved the user data. The administrator has complete control over who may read, edit, and manage the data.

3.3.3 Login:

Users can use this technique to authenticate themselves with the system and access the application by entering their credentials, which consist of their username and password.

3.3.4 Upload image:

The user must choose an image from the dataset and upload it into the program. The image must have a caption.

3.3.5 Prediction:

It produces a caption as an output, and the image caption that we have given to it will be displayed as a result of our model's operation.

3.3.6 Logout:

The user has the option to exit the application when the result has been generated.

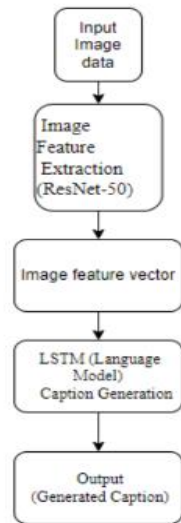


Fig6: Block Diagram

4.IMPLEMENTATION

4.1 CNN:

Over the past 20 years, Deep Learning has shown to be a very powerful technology due to its ability to handle large amounts of data. Convolutional neural networks, also referred to as CNNs or ConvNets, are the most popular deep neural networks in deep learning, especially for applications in computer vision. The type of deep neural network utilized in deep learning that is most commonly used to evaluate visual data is called a convolutional neural network (CNN). It uses a unique method known as convolution. As it currently understood, When two functions are mathematically operated upon, a third function known as convolution is produced that describes how one's shape is transformed by another.

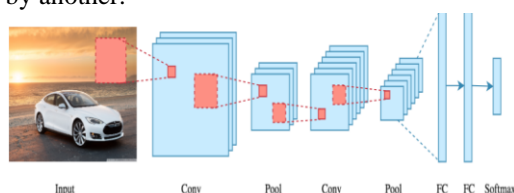


Fig7: Working of CNN

Convolutional, pooling, and fully-connected (FC) layers are the three layers that make up a CNN. These layers are stacked to construct a CNN architecture. Beyond these layers, two more important factors are the dropout layer and the activation function. The RGB (Red, Green, and Blue) color space, which extends from 0 to 255, can be used to extract and identify various visual aspects for analysis (a process known as feature extraction) using a convolution tool. The network used for feature extraction is composed of multiple pairs of convolutional or pooling layers.

4.2 ResNet 50:

A kind of deep neural network architecture called ResNet (Residual Network) was introduced to solve the vanishing gradient problem that arises when deep convolutional neural networks (CNNs) are being trained. ResNet gives the network the ability to learn residual functions by introducing skip connections, also referred to as residual connections. The model can bypass specific levels thanks to these skip connections, which send the image straight to the output layers. This facilitates the training of extremely deep networks by reducing the impact of the vanishing gradient issue.

ResNet can be a key component in feature extraction from photos within an image caption generator. In order to obtain useful features from the images, the encoder portion of an image captioning model usually makes use of a CNN that has already been trained, such as ResNet.

The goal is to extract high-level features from photos by using the information that the pre-trained ResNet model has acquired on a sizable dataset .

- Images are used to obtain features using the ResNet model. For image classification tasks, the model is usually pre-trained on a large dataset (e.g., ImageNet). Hierarchical and abstract aspects in photos are captured by the weights that were learned during pre-training.
- The pre-trained ResNet model is used to obtain characteristics from intermediate layers given an input image. The input's high-level visual information is represented by the features.
- The decoder component of the picture captioning model receives the features that were extracted from the image. Based on the input attributes, the decoder—which is frequently implemented as a transformer or recurrent neural network (RNN)—creates a textual description of the image.

4.3 RNN:

An example of a neural network type is a recurrent neural network (RNN), which takes as its input the output of the previous stage. Conventional neural networks have inputs and outputs that are independent of each other. On the other hand, the preceding words must be retained in cases when guessing the next word in a sentence requires knowledge of the preceding words. RNN was developed as a result, and it approached this issue by using a Hidden Layer. Its main and most important characteristic is the Hidden state of an RNN, which holds some information about a sequence. It is sometimes referred to as Memory State since it keeps track of previous inputs to the network..

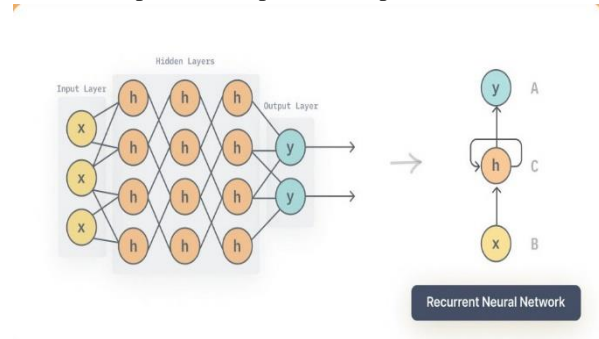


Fig8: Single RNN

It does the similar task on all the hidden layers to produce the output, using the same settings for every hidden layer. Unlike other neural networks, this lowers the parameter complexity. Although it isn't referred to as a "recurrent neuron," a recurrent unit is the basic processing unit of a recurrent neural network (RNN). Because of its special capacity to preserve a hidden state, this unit enables the network to recognize sequential relationships by processing and remembering prior inputs. RNNs have the same input and output architecture as other deep neural architectures. However, the information flow differs depending on the input and output.

It generates the output by doing the same operation on all inputs or hidden layers, using the same parameters for each input. It differs from other neural networks in that it has less parameter complexity. A recurrent unit, although not referred to as a "recurrent neuron," is the central processing unit of a recurrent neural network (RNN). This unit's unique ability to maintain a concealed state allows the network to process and retain previous inputs, which in turn helps the network understand sequential relationships. Like other deep neural architectures, RNNs have the same input and output architecture. Yet, there are differences in the way data is transferred from input to output.

The formula for calculating the current state:

$$h_t = f(h_{t-1}, x_t)$$

Where, h_t -> current state

h_{t-1} -> previous state

x_t -> input state

Formula for applying Activation function(tanh):

$$h = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

Where, W_{hh} -> weight at recurrent neuron

W_{xh} -> weight at input neuron

The formula for calculating output:

$$Y_t = W_{hy}h_t$$

Where, Y_t -> output

W_{hy} -> weight at output layer

4.4 LSTM:

The Long Short-Term Memory, or LSTM, is an enhanced RNN. For sequence prediction tasks, LSTM performs remarkably well in capturing long-term dependencies. The RNN method has certain drawbacks, which we address by introducing the LSTM algorithm.

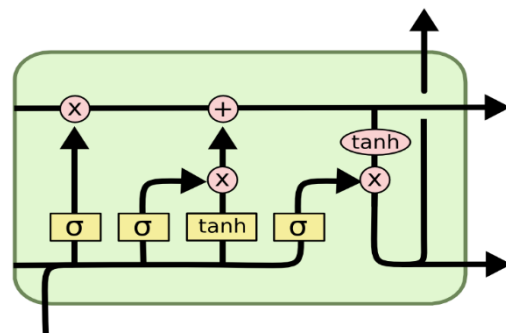


Fig9: Architecture of LSTM

- It is a difficult task to a network to learn long-term dependencies in a standard RNN since it only has one hidden state that is retained over time. In contrast, LSTMs solve this issue by introducing memory cells, which are containers that can store information for a longer amount of time.
Three gates govern a memory cell:
 - i. Input gate
 - ii. Forget gate
 - iii. Output gate
- The gates determine which data should be input into, taken out of, and output from the memory cell.
- Unlike RNNs, which lack memory units, LSTMs have a unique memory unit that enables

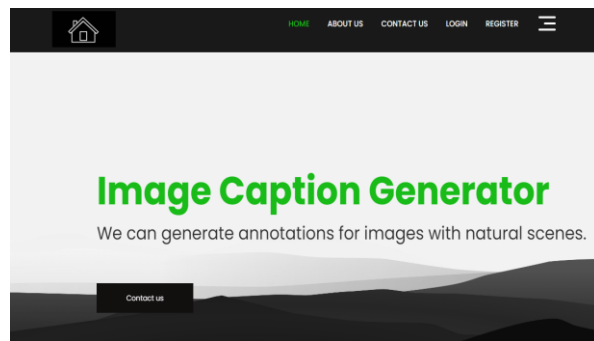
them to recognize long-term dependencies in sequential data.

- While RNN is also meant to process sequential data, its memory capacity is constrained. In contrast, LSTM is well suited for handling sequential data.
- Compared to RNN, the training process of LSTM is slower because of its increased complexity. Because of its more straightforward architecture, RNNs train a little bit faster.
- Long sequences are more effective for LSTM, but RNN finds it difficult to store information.

5. RESULTS AND ANALYSIS:

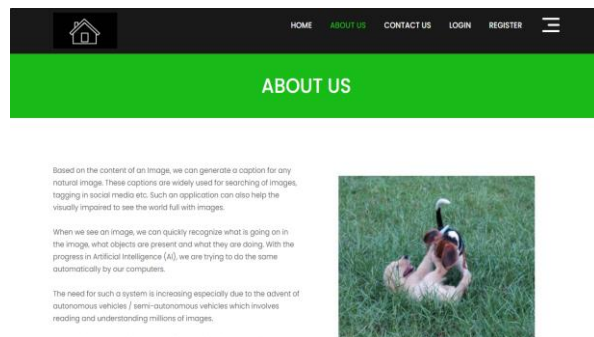
5.1 Home page:

This is the home page where we land after clicking on the link.



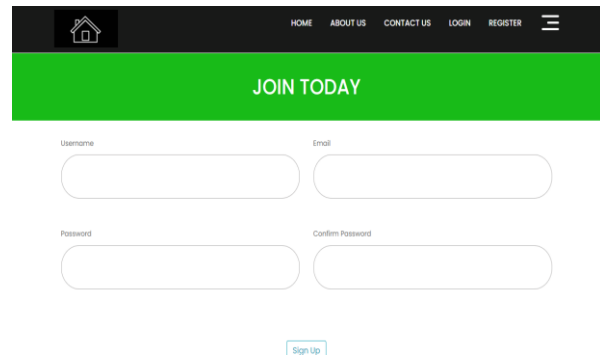
5.2 About Page:

Here we have a slight description about the project.



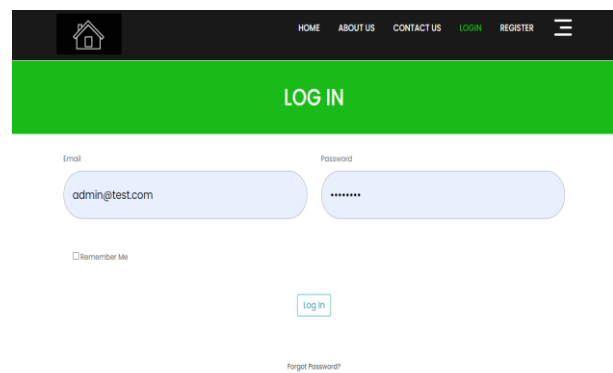
5.3 Register Page:

Here User registers themselves.



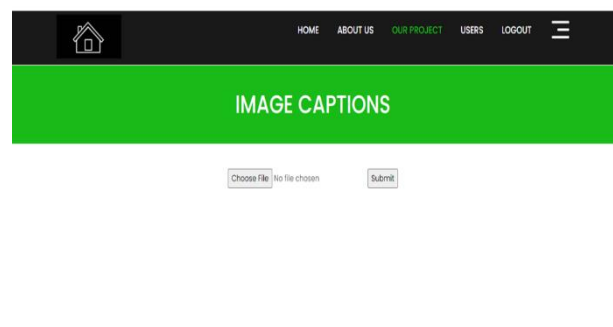
5.4 Login Page:

Here user logs in with the credentials they registered with.



5.5 Upload Page:

Here user uploads the image and caption is generated.



5.6 Results:

Here we get the results.

- [10]. V. Kesavan, V. Muley and M. Kolhekar, “Deep Learning based Automatic Image Caption Generation”, IEEE Access, 2019.
- [11]. Girish Kulkarni, Visruth Premraj, Sagnik Dhar, Siming Li, Yejin Choi, Alexander C Berg, and Tamara L Berg. “Baby talk: Understanding and generating image descriptions”, IEEE Transactions on Pattern Analysis and Machine Intelligence, 35:2891–2903, June 2013.
- [12]. N K Kumar, D Vigneswari, A Mohan, K Laxman, J Yuvaraj, “Detection and Recognition of Objects in Image Caption Generator System: A Deep Learning Approach”, IEEE – 2019.

