

# Image Caption Generation By Using Deep Learning

**Dr. THOMAS FELDMAN**

*Head Of The Department of Computer Science & Engineering, Dnyanshree Institute of Engineering & Technology, Satara*

**Dr DINESH KUMAR**

*Assistant professor, Department of Computer Science & Engineering, Dnyanshree Institute of Engineering & Technology, Satara*

**PROF. MUSILEK**

*Students, Department of Computer Science & Engineering, Dnyanshree Institute of Engineering & Technology, Satara*

## 1. ABSTRACT

- The Framework Includes a Convolution Neural Network (CNN) Observed Via way Of Means of a Recurrent Neural Network (RNN). It generates an English Sentence from an Enter Picture Neural Networks. By Getting To Know Expertise from Picture and Caption Pairs, The Approach Can Generate Picture Captions Which Are Typically Semantically Descriptive and Grammatically Correct.
- Automatically Describing the Content Material of Pic. The Usage of Herbal Languages is Essential and Tough Task. It Has Exceptional Capability Impact.
- The Easy Approach to This Trouble That We're Presenting Is That We Are Able to Teach a CNN-LSTM Version in Order That It Is Able to Generate a Caption Primarily Based Totally at The Image

**Keywords:** CNN, LSTM, Recurrent Neural Network, Neural Network.

## 2. INTRODUCTION

Image Caption Generator Is a Challenge That Entails Computer Imaginative and Prescient and Natural Language Processing Ideas to Apprehend the Context of Photo and Describe Them in a Natural Language Like English. The First One Is an Image Primarily Based Totally Version Which Extracts the Functions of The Image, and The Opposite A Language Primarily Based Totally Version Which Interprets the Functions and Items Given Via Way of Means of Our Image-Primarily Based Totally Version to a Natural Sentence. We Can Be the Use of a Pre-Trained CNN Community this Is Trained at The Datasets. The Image Is Converted Right into Trendy Resolution. This Will Make the Input Constant for The Version for Any Given Image.

## 3. LITURATURE SURVEY

[1] J. Donahue, Y. Jia, J. Hoffman, N. Zhang, E. Tzeng, and T. Darrell. Decaf: A deep convolutional activation feature for generic visual recognition. In ICML, 2014 .

The CNN has an excellent performance in machine learning problems and one of the most common algorithms.

**[2]. A. Aker and R. Gaizauskas. Generating image descriptions using dependency relational patterns. In ACL, 2010.**

A simple solution to this problem that we are proposing is that we will train a CNN-LSTM model so that it can generate a caption based on the image

**[3] Oriol Vinyals, Alexander Toshev , Samy Bengio , and Dumitru Erhan. Show and tell: A neural image caption generator. CoRR , abs/1411.4555, 2014.**

LSTM is local in both space as well as in time, the computational complexity is per time of step and also the weight pattern representation.

**[4].Xin lei Chen and C. Lawrence Zit nick. Learning a recurrent visual representation for image caption generation. CoRR , abs/1411.5654, 2014.**

Automatically describing the content of an image using properly arranged English sentences is a tough challenging task, but it could is something very necessary for helping visually impaired people.

#### **4. PROBLEM DEFINITION**

To develop a system for users, which can automatically generate the description of an image with the use of CNN along with LSTM.

#### **5. PROJECT OBJECTIVE**

- We will be implementing the caption generator using CNN and LSTM.
- The objective of our project is to learn the concepts of a CNN and LSTM model

## 6. PROPOSED SYSTEM

In This Project We Can Use Two Platforms

### 1. Kaggle

We Can Use Two Datasets in This Project

- o Flickr8k\_Dataset
- o flickr-image-dataset

We Can Use One Model in This Project

- o VGG 16

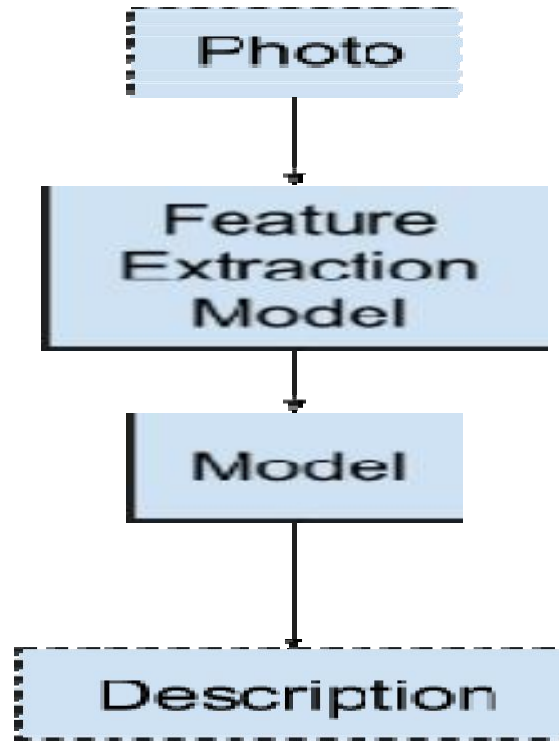
VGG16 Is a Convolutional Neural Network Version. Very Deep Convolutional Networks for Large-scale Image Recognition. The Version Achieves 92.7 Perc. Top-Five Check Accuracy in ImageNet, That's A Dataset of Over 14Million Snapshots Belonging to A Thousand Classes. The VGG Version Stands for The Visual Geometry Group from Oxford. The Version Changed into Very Simple and Had a More Intensity Than AlexNet.

### 2. Google Colab

We Can Use Google Colab for Coding Part & It Is Easy to Use, Easy to Access. This Platform Is Free

### SYSTEM ARCHITECTURE :-

The Proposed Structure Includes the Bottom Layer and The Top Layer. The Backside Layer Extracts the Visible and High Degree Semantic Facts from Picture and Detected Regions, Even as The Top Layer Integrates Each of Them with Interest Mechanism for The Caption Generation. So, To Make Our Picture Caption Generator Version, We Can Be Merging Those Architectures. It Is Likewise Referred to As A CNN-LSTM Version. CNN Is Used for Extracting Functions from The Picture. We Will Use the Pretrained Version Exceptions. LSTM Will Use the Data from CNN to Assist Generate an Outline of The Image. So, To Make Our Image Caption Generator Model, We Will Be Merging These Architectures. It Is Also Called as A CNN-LSTM Model. CNN Is Used for Extracting Features from The Image. We Will Use the Pre-trained Model Exceptions. LSTM Will Use the Information from CNN To Help Generate A Description Of The Image.



## METHADODOLOGY

So, to make our image caption generator model, we will be merging these architectures. It is also called a CNN-LSTM model. CNN is used for extracting features from the image. We will use the pre-trained model Exceptions. LSTM will use the information from CNN to help generate a description of the image.

### 1.Convolution Neural networks (CNN):

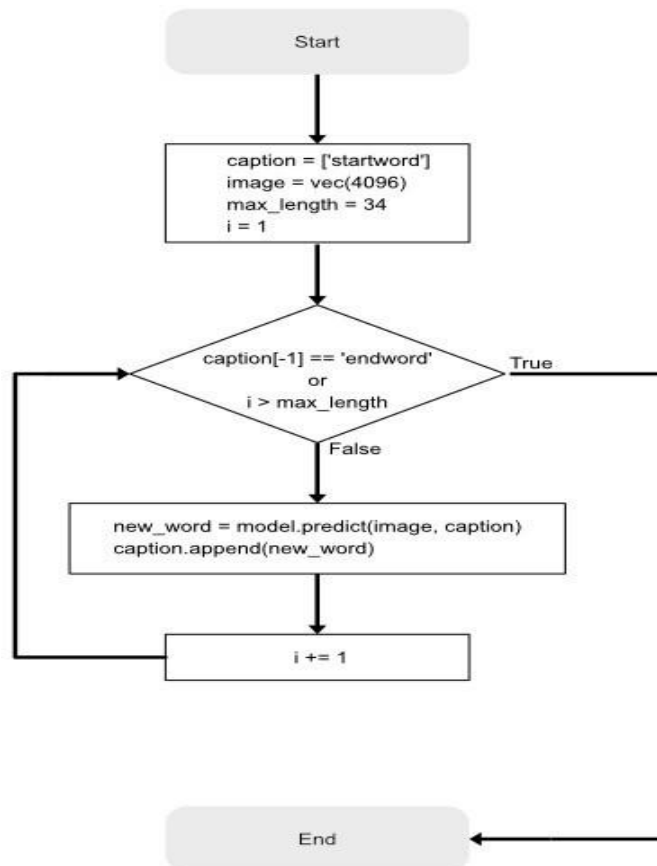
Convolution Neural networks are specialized deep neural networks which can process the data that has input shape like a 2D matrix. Images are easily represented as a 2D matrix and CNN is very useful in working with images. CNN is basically used for image classifications and identifying if an image is a bird, a plane or Superman, etc.It scans images from left to right and top to bottom to pull out important features from the image and combines the feature to classify images. It can handle the images that have been translated, rotated, scaled and changes in perspective.

### 2.Long Short Term Memory (LSTM):

LSTM stands for Long short term memory, they are a type of RNN (recurrent neural network) which is well suited for sequence prediction problems. Based on the previous text, we can predict what the next word will be. It has proven itself effective from the traditional RNN by overcoming the limitations of RNN which had short term memory. LSTM can carry out relevant information throughout the processing of inputs and with a forget gate, it discards non-relevant information.

## SYSTEM FLOW

As our goal is not to map an image to a specific caption but rather learn the relationship between image features and word sequences and between word sequences and single words our target is not a complete caption in the dataset but a single word. In other words, the model generates a new caption word by word based on a given image feature and a caption prefix. In the beginning, each caption only contains the artificial start sequence (which we introduced in the previous part).



## 7. RESULT

we have successfully generated mostly grammarly correct and human readable captions describing what is happening in the image. Most of the images correctly state what objects appear in the scene, count number of appearance and gives an intuitively correct verb to logically complete the sentence. Colors of the object and spatial relationships between object are also well captured.



a little girl in a pink dress going into a wooden cabin .  
 a little girl climbing the stairs to her playhouse .  
 a little girl climbing into a wooden playhouse .  
 a girl going into a wooden building .  
 a child in a pink dress is climbing up a set of stairs in an entry way .



two dogs on pavement moving toward each other .  
 two dogs of different breeds looking at each other on the road .  
 a black dog and a white dog with brown spots are staring at each other in the street .  
 a black dog and a tri-colored dog playing with each other on the road .  
 a black dog and a spotted dog are fighting



young girl with pigtails painting outside in the grass .  
 there is a girl with pigtails sitting in front of a rainbow painting .  
 a small girl in the grass plays with fingerpaints in front of a white canvas with a rainbow .  
 a little girl is sitting in front of a large painted rainbow .  
 a little girl covered in paint sits in front of a painted rainbow with her hands in a bowl .



man laying on bench holding leash of dog sitting on ground  
 a shirtless man lies on a park bench with his dog .  
 a man sleeping on a bench outside with a white and black dog sitting next to him .  
 a man lays on the bench to which a white dog is also tied .  
 a man lays on a bench while his dog sits by him .



the man with pierced ears is wearing glasses and an orange hat .  
 a man with glasses is wearing a beer can crocheted hat .  
 a man with gauges and glasses is wearing a blitz hat .  
 a man wears an orange hat and glasses .  
 a man in an orange hat starring at something .



the small child climbs on a red ropes on a playground .  
 a small child grips onto the red ropes at the playground .  
 a little girl in pink climbs a rope bridge at the park .  
 a little girl climbing on red roping .  
 a child playing on a rope net .



a dog runs on the green grass near a wooden fence .  
 a boston terrier is running on lush green grass in front of a white fence .  
 a boston terrier is running in the grass .  
 a black and white dog is running through the grass .  
 a black and white dog is running in a grassy garden surrounded by a white fence .



white dog with brown ears standing near water with head turned to one side .  
 white dog playing with a red ball on the shore near the water .  
 dog with orange ball at feet , stands on shore shaking off water .  
 a white dog shakes on the edge of a beach with an orange ball .  
 a dog shakes its head near the shore , a red ball next to it .



smiling boy in white shirt and blue jeans in front of rock wall with man in overalls by  
 a young child is walking on a stone paved street with a metal pole and a man behir  
 a young boy runs across the street .  
 a little boy is standing on the street while a man in overalls is working on a stone w  
 a boy smiles in front of a stony wall in a city .



the black dog jumped the tree stump .  
 a mottled black and grey dog in a blue collar jumping over a fallen tree .  
 a large black dog leaps a fallen log ,  
 a grey dog is leaping over a fallen tree .  
 a black dog leaps over a log .

## 8. CONCLUSION AND FUTURE SCOPE

### Conclusion

We have presented an end-to-end neural network system that can automatically view an image and generate a reasonable description in plain English. An Image Generator has been developed by CNN LSTM Method. We analyzed and modified an image captioning method LRCN (Long-term Recurrent Convolutional Network method.). To understand the method deeply, we decomposed the method to CNN, RNN, and sentence generation.

### Future Scope

In the future, we would like to explore methods to generate Single sentences with different content.

## 9. References

[1] Ali Farhadi, Mohsen Hejrati, Mohammad Amin Sadeghi, Peter Young, Cyrus Rashtchian, Julia Hockenmaier, and David Forsyth. Every picture tells a story: Generating sentences from images. In Proceedings of the 11th European Conference on Computer Vision: Part IV, ECCV'10, pages 15–29, Berlin, Heidelberg, 2010. Springer-Verlag.

[2] Polina Kuznetsova, Vicente Ordonez, Alexander C. Berg, Tamara L. Berg, and Yejin Choi. Collective generation of natural image descriptions. In Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers - Volume 1, ACL '12, pages 359–368, Stroudsburg, PA, USA, 2012. Association for Computational Linguistics.

[3] Siming Li, Girish Kulkarni, Tamara L. Berg, Alexander C. Berg, and Yejin Choi. Composing simple image descriptions using web scale n-grams. In Proceedings of the Fifteenth Conference on Computational Natural Language Learning, CoNLL '11, pages 220–228, Stroudsburg, PA, USA, 2011. Association for Computational Linguistics.

[4] Xinlei Chen and C. Lawrence Zitnick. Learning a recurrent visual representation for imagecaption generation. CoRR, abs/1411.5654, 2014.

[5] Junhua Mao, Wei Xu, Yi Yang, Jiang Wang, and Alan L. Yuille. Deep captioning with multimodal recurrent neural networks