# Deep Learning-Based Image Captioning Using ConvolutionalNeural Networks and Recurrent Neural Networks

**SREEKANTAM VASUDHA, SIRIKONDA ANANTHNAG, JANGILI RAVI KISHORE**

**Assistant Professor [1,2,3]**

vasudhasvit87@gmail.com, ananthnagsvit9f@gmail.com, jangiliravi.kishore1@gmail.com
Department of Computer Science and Engineering, Sri Venkateswara Institute of Technology, N.H 44, Hampapuram, Rapthadu, Anantapuramu, Andhra Pradesh 515722

**KEYWORDS:**

Convolutional neural networks, recurrent neural networks, deep learning, image captioning.

**ABSTRACT**

Captioning is an important problem for all data mining companies as a whole due to the emergence of new generations. It might be a lengthy and complicated process to interpret such data using a device. A greater grasp of the concept of a picture is necessary for a device to comprehend its context and environmental data. Although conventional methods have not been accompanied by extensive understanding of strategy, this is beginning to change. An automated transcript of image annotations will be generated in this research by using Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) to produce a collection of text that adequately characterises the picture. To organise our model, we used the Flickr 8000 dataset. Since this caption requires a real neural community, we provide clear instructions on how to create one. We start by associating the description with the optical neural network, then we take a picture and break it into features. Then we use CNN and RNN to build a deep neural network. Thanks to the recurrent neural network, this function may be described in a clear and concise manner.

## I. INTRODUCTION

Providing guidance on how to edit software is one application of the picture caption. Quicken and simplify the process of adding captions to digital media with the aid of this picture caption template. It also has a programme to assist the visually impaired. We are all aware that not everyone is born with the ability to see clearly, and seeing the world around them. These people may feel the extra miles with the aid of a comprehensive transcript. Numerous verifiable facts are disseminated in supplementary material by means of the media and publishing organisations. The captions model expedites and automates the process, allowing for the presentation of automated captions for these recordings. technology, which eliminates the need to name these facts hundreds of times [1]. Finally, it has further potential applications in online postings. As AI continues to advance at a dizzying rate, social media plays an increasingly important role in helping to categorise and separate various media assets. If you're short on time, try using the automated caption generator on this page. Our goal is to provide a framework for automated picture captions. Using the Convolutional Neural Network (CNN) from the Image Properties class, this may be accomplished. We proceed to generate the caption by use of RNN. Long short term memory (LSTM) was given the nod. Here we have a variant of RNN. LSTM is able to review data that is presented in a sequential format, such a string of words or letters. The memory circle is able to get knowledge about the copy of the afterlife because to these structures' hidden layers that link the input and output layers. As a result, RNN is able to resolve consecutive data, whereas CNN primarily operates with spatial capabilities [2]. The research makes use of state-of-the-art recurrent neural networks for natural language processing and deep learning for enhanced virtual and predictive computing. Deep learning is a way to investigate systems that programmers utilise. computational resources to analyse large datasets for recurring patterns. The accumulation of the era's fast economic growth, abundant research, and ever-increasing computer power. Worldwide, the deep learning and AI businesses are expanding at a rapid pace, and they might soon become among the world's most lucrative. The information age has given rise to artificial intelligence, which has transformed data into its new "oil" in the 21st century. A tremendous quantity of data is created every second in the modern world. Our goal is to build models that can analyse these datasets and either reveal trends or provide solutions for evaluations and research. Thorough research is necessary to do this. Computer vision is a rapidly developing area of study within the realm of personal computer technology that seeks to provide computers the ability to comprehend visual information contained in images. There is a significant cognitive gap between humans and computers. Various factors, including the form of the camera, the lighting, the clarity, the size, the point of view, etc., reduce perception to a unit of raw numbers in their eyes. Building a reliable model that can draw reliably in every environment is difficult; making a computer creative and forward-thinking is much more so [3]. The typical social neural structures have been able to make advantage of the most recent, crucial piece of advice. Earlier inputs are no longer taken into account by the output when the gadget is being developed. The reason for this is because we have been ignoring important details in our past recollections. That is why the device's memory issues are resolved when RNN is used. Because of this, we were able to create a machine that is gentler on the planet.

## I. RELATED WORK

The authors of the study are Shuang Bai and colleagues [...] An image comment is generated automatically using the image comment technique. The miles are becoming more and more of a focus as a new field of study. To accomplish the goal of hanging photographs, it is necessary to capture and represent the semantic features of the images in natural languages. Annotating photographs is a daunting task that combines the fields of computer vision and natural language processing. A number of approaches have been suggested to address this issue.

Researchers checked up on the image commentary studies to see how things were going. Based on one's perspective, the picture commenting tactics are classified as a fantastic activity. In each chapter, we highlight the representative techniques and discuss their advantages and disadvantages. Yin Xiao and colleagues [5] It delves into the use of data visualisation tools, which may capture large amounts of seemingly ordinary data, to improve the data retrieved from images by using captioning strategies appropriate to the country of painting. The results of the experiments conducted on multiple MS COCO benchmark datasets, evaluated with CIDEr-D, a widely used subtitle performance metric, demonstrate that variations in subtitle graphics strategies utilising expert chart statistics can outperform those relying only on photo statistics. Yang Zhongliang et al. [6] One of the most fundamental tasks in image knowledge is the automatic generation of a description of a picture in plain language. In order to automatically learn how to describe the content of photos, this article introduces a method that combines a human visual engine with multi-version neural networks. Elements detection version and localization are the two parts of the model that take snapshots of tool information and their geographical history, respectively; For the purpose of generating sentences, a deep recurrent neural network (RNN) relies only on Long Short Term Memory (LSTM) hardware that incorporates an attentional mechanism. You may choose which parts of the picture to automatically match each phrase of the description with when you create it. Tyhis is quite similar to how the human visual system handles attention. Experimental The findings on the Flickr 8k dataset show that the suggested method is effective, surpassing the performance of earlier measurement modes. The work of G. Schwing and colleagues [7] Caption quality has been significantly enhanced in recent times with the use of recurrent neural networks driven by LSTM devices. Although LSTM devices are interesting for remembering dependencies and minimising the vanishing gradient issue, they are complicated and sequential across time. Eventually, training will be challenging because to the sophisticated processing and rewriting mechanism, intrinsic sequential processing, and the large memory needed due to backpropagation time (BPTT). (Hu et al., 2008) In order to understand the connection between objects, geometric attention was developed. This method makes use of the sizes and locations of bounding boxes, supposing that a stronger link exists when the boxes are closer together. The next step was for the approved models to extract all the required features using an image object detector, and then for the captions to be created using LSTM with an attention mechanism.

## II.    PROPOSED METHODOLOGY

Use convolutional neural networks (CNNs) and recurrent neural networks (RNNs) to build a deep neural network. At first, we link the picture with the description. After a convolutional neural network has processed the picture and segmented it into its component parts, the recursive neural network will take that information and turn it into a clearly articulated language of description.

### CNN Encoder

A convolutional neural network directly encodes an image into a compressed representation; this network forms the basis of the absolute encoder. Inception V3's Residual Network module encodes CNN data. There are a number of tool-specific overrides available in Layer Activation. Additionally, they are input into all subsequent layers, including the network, enabling the use of the Inception v3 module.
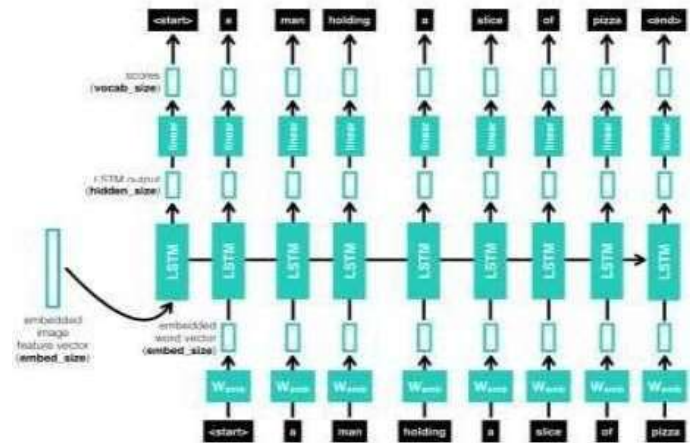
Fig.1 CNN-Encoder

## RNN Decoder

A CNN encoder is observed with a recursive neural community resource that generates a corresponding sentence. TheRNN decoder is built into an unpaired LSTM layer detected by a truly linked (linear) layer. The RNN community is aware of the Flickr 8k dataset. It is used to look forward to the next word of asentence, mostly based on previous terms. Captions are provided as a glyph of the glossary to allow the RNN model to betaught and republished to reduce errors andcreate higher-level, more understandable text describing the image.

### Flickr 8k dataset:

This dataset contains 8000 images, eachcontaining five specific descriptions based on skill and context. Here they are beautifully described and used in various exclusive environments.

### Inception V3 module:

There are 4 differences. The first Google Net should be Inception-v1, but several typos in Inception-v3 lead to incorrect descriptions of Inception variables. Perhaps this is due to the stiff competition at ILSVRC this second. Thus, there is a lot of criticism on the net between versions 2 and 3. Some reviews assume that versions

2 and 3 are the same with some slight exceptional tweaks.

## III. RESULTS

The result of applying the provided idea to image captioning using CNN and LSTM is shown below. The results shown here were retrieved from the Flickr8k dataset.We used many processes to get to the outcome, including: loading libraries, data, and captions as keys and values in a dictionary; making train, test, and validation sets; loading the model; saving the captions; and vectorization. Once all these stages have been executed, the output will display the outcome of the caption production. Figure 4 displays the produced caption for one of the dataset photos.



Fig.2 Test image

## IV. CONCLUSION

A multimodal approach to automated picture captioning using InceptionV3 and LSTM is introduced in this study. The suggested model was built using an encoder-decoder architecture and trained on a massive Flicker 8k dataset, which contains 8,000 pictures with captions. To compress an image's feature matrix into a more manageable format, we used InceptionV3, a convolutional neural network, as the encoder. After that, the description was generated using a decoder that was chosen as a language model LSTM. From what we can tell from the experiments, the suggested approach can successfully provide appropriate picture captions.

## REFERENCES

Y. Wei, Z. Zhang, J. Dai, J. Gu, and H. Hu [1]. Relation networks for object detection. Volume 3588, Issue 3597, Pages 3588–3597, 2018 IEEE Conference on Computer Vision and Pattern Recognition. [2] "Deep residual learning for image recognition," by K. He, X. Zhang, S. Ren, and J. Sun Published in 2016 by IEEE, this book is part of the proceedings from the 2016 computer vision and pattern recognition conference. In their 1997 article "Long short-term memory," Hochreiter and Schmidt discussed neural computation in volume 9, issue 8, pages 1735–1780. A. Yuille, J. Wang, Z. Huang, W. Xu, Y. Yang, and J. Mao [4]. Deep captioning utilising multimodal recurrent neural networks (m-RNN). publishable on arXiv at 2014, 1412.6632. Reference: [5] Kiros, Salakhutdinov, and Zemel, R. Neural language models that cover several modalities. Pages 595–603 of 2014's International Conference on Machine Learning. Donahue, J., Hendricks, L. Anne, Guadarrama, S., Rohrbach, M., Saenko, K., and Darrell, T. (2006). Image identification and description using recurrent convolutional networks trained over extended periods of time. With references to pages 2625–2634 from the 2015 EEE Conference on Computer Vision and Pattern Recognition. S. Bengio, O. Vinyals, A. Toshev, and D. Erhan. A neural picture caption generator: a demonstration and explanation. Found on pages 3156-3164 of the 2015 IEEE Conference on Computer Vision and Pattern Recognition proceedings. "Relation networks for object detection" (pp. 3588–3597, 2018) in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, by H. Hu, J. Gu, Z. Zhang, J. Dai, and Y. Wei. The authors of the publication are A. Vaswani, N.

Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin. Just pay attention. Volume 5998, Issue 6008, 2017: Advances in Neural Information Processing Systems [10]Mei, T., Yao, Y., Pan, Y., and