



# **iJRASET**

International Journal For Research in  
Applied Science and Engineering Technology



---

# **INTERNATIONAL JOURNAL FOR RESEARCH**

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume: 11    Issue: VI    Month of publication: June 2023**

**DOI: <https://doi.org/10.22214/ijraset.2023.54297>**

**[www.ijraset.com](http://www.ijraset.com)**

**Call:  08813907089**

**E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)**

# Captioning Image Using Deep Learning: A Novel Approach

Sunanda Chakraborty<sup>1</sup>, Shibam Chakraborty<sup>2</sup>, Moloy Dhar<sup>3</sup>, Nirupam Saha<sup>4</sup>, Bidyutmal Saha<sup>5</sup>, Pallabi Das<sup>6</sup>, Rafiqul Islam<sup>7</sup>, Sutapa Sarkar<sup>8</sup>

<sup>1, 2, 3, 4, 5, 6, 7, 8</sup> Department of Computer Science & Engineering, Guru Nanak Institute of Technology, Kolkata, India

<sup>1</sup>sunandachakraborty270@gmail.com, <sup>2</sup>chakrabortyshibam184@gmail.com, <sup>3</sup>moloy.dhar@gnit.ac.in, <sup>4</sup>nirupam.saha@gnit.ac.in, <sup>5</sup>bidyutmal.saha@gnit.ac.in, <sup>6</sup>pallabi.das@gnit.ac.in, <sup>7</sup>rafiqul.islam@gnit.ac.in, <sup>8</sup>sutapa.sarkar@gnit.ac.in

**Abstract:** Captioning image using deep learning is a technology that aims to generate descriptive and accurate textual descriptions for images. By using the power of deep neural networks, this approach enables computers to understand and interpret visual content bridging the gap between the visual and textual domains. This paper focuses on developing an image captioning system using deep learning techniques. The paper aims to generate descriptive textual captions for images, enabling machines to understand and communicate the content of visual data. The methodology involves leveraging convolutional neural networks (CNNs) for image feature extraction and recurrent neural networks (RNNs) for sequential language generation. Captioning images is a challenging task in computer vision that involves describing the content of an image in natural language. In recent years, deep learning techniques have shown remarkable success in various computer vision tasks, including image captioning.

**Keywords:** CNN, VGG-19, LSTM, NLP, RNN

## I. INTRODUCTION

Captioning image using deep learning is an innovative procedure which combines computer vision and natural language processing to automatically generate descriptive captions for images. It addresses the challenging task of overcoming the gap between visual perception and language understanding. VGG-19 has been pre trained on more than a million images from ImageNet database so the output comes out quite accurate. Image captioning is an exciting field at the intersection of computer vision and natural language processing (NLP). It involves generating descriptive textual captions for images, enabling machines to understand and communicate the content of visual data. Deep learning techniques, such as (CNNs) and recurrent neural networks (RNNs), have shown remarkable success in image captioning tasks.

The paper will leverage a pre-trained CNN to extract meaningful features from the images. These extracted features will serve as input to the RNN-based captioning model. The RNN, equipped with recurrent cells such as LSTM or GRU, will learn to generate descriptive captions based on the extracted image features. Training the model will involve optimizing the parameters to minimize the captioning loss. In this paper, we are using MS-COCO dataset which is used for large-scale object detection, segmentation and captioning and are used in various computer vision papers. It contains 328,000 images of everyday objects and humans. It contains feature-rich annotations including object detection, captioning, person-keypoints and more. Accurate image captioning has numerous practical applications, including assisting visually impaired individuals in understanding images, enhancing image search engines, and enabling better image indexing and retrieval. Our approach consists of two main components: an image encoder and a language decoder. The image encoder employs a pre-trained CNN, such as VGG or ResNet, to extract high-level visual features from the input image. These features capture the semantic information present in the image and serve as the input to the language decoder. The language decoder is implemented using an RNN, specifically a long short-term memory (LSTM) network, which generates captions based on the visual features obtained from the image encoder.

## II. METHODOLOGY

The result of image captioning is the generation of descriptive and informative captions for images. By using advanced techniques such as deep learning and natural language processing, image captioning models can analyze the content of an image and generate textual descriptions that accurately represent the visual elements.

Before using Tokenizer to convert texts to integers, we mark the start and end of the strings with a unique marker.

start\_marker = "aaa"

end\_marker = "zzz"

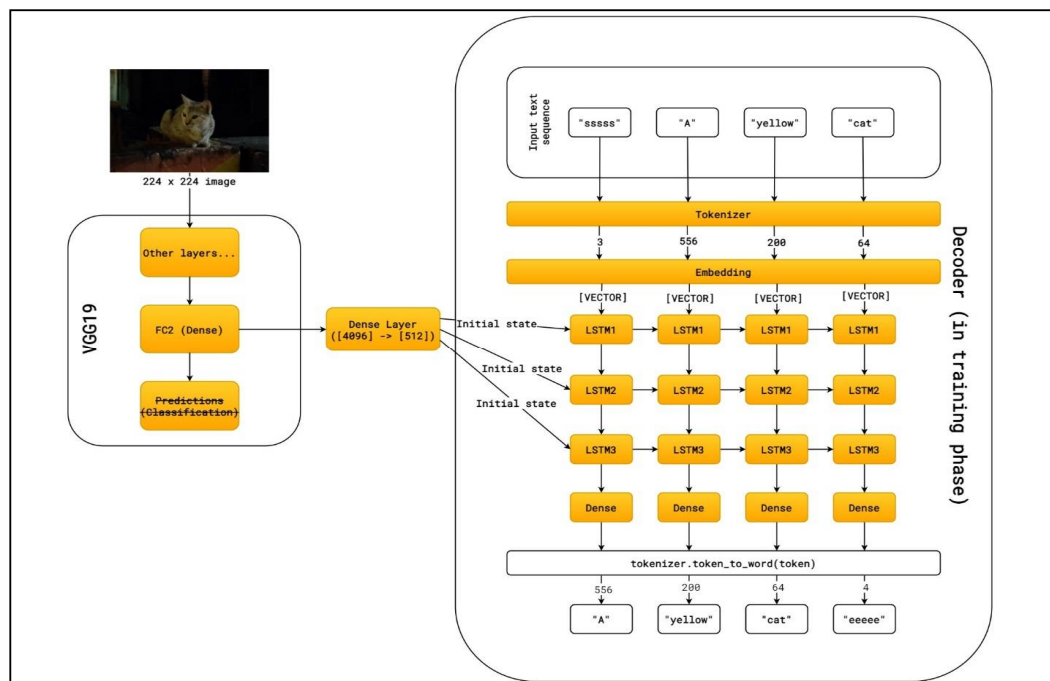


Fig. 1

We are limiting the maximum number of words in our vocabulary. We will only use the 10500 most frequent words in the captions from the training data.

num\_words = 10500

### III. RESULT

Captioning image using deep learning has got a significant advancement in recent years. Deep learning models, particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been developing highly effective image captioning systems.

The decoder consists of 3 LSTM layers whose internal state-sizes are 512

state\_size = 512

The embedding layer converts integer tokens into vectors of length 128

embedding\_size = 128

Image captioning using deep learning has got a significant advancement in recent years. Deep learning models, particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been develop highly effective image captioning systems. One popular approach is the encoder-decoder architecture, where a CNN is used as the encoder to extract image features, and an RNN is employed as the decoder to generate captions. This architecture has demonstrated impressive results in generating coherent and descriptive captions for images.

In summary, image captioning offers valuable solutions for accessibility, search ability, and user engagement. With continued advancements in technology and research, we can expect image captioning systems to become even more sophisticated, enabling better understanding and interaction with visual content. Here, we have used RMSprop optimizer.



*A. We Have used RMSprop Optimizer*

MODEL SUMMARY: VGG-19

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 224, 224, 3)]	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv4 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv4 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv4 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0

flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312

Total params: 139,570,240

Trainable params: 139,570,240

Non-trainable params: 0

Fig. 2

We use this model differently in training and testing.

However, challenges still exist in image captioning, such as accurately capturing fine-grained details, handling complex scenes, and generating captions that capture context and semantic meaning. Ongoing research and development efforts aim to address these challenges and improve the accuracy and contextual understanding of image captions.

For example, if the sentence is “A dog is running on the ground”, then if we input sitting to an LSTM layer, we will train it to output “on”. This is called the “**Teacher Forcing**” method.

Whereas, during testing phase, we use the output token of the previous layer as an input to the next layer.

#### B. Sample outputs

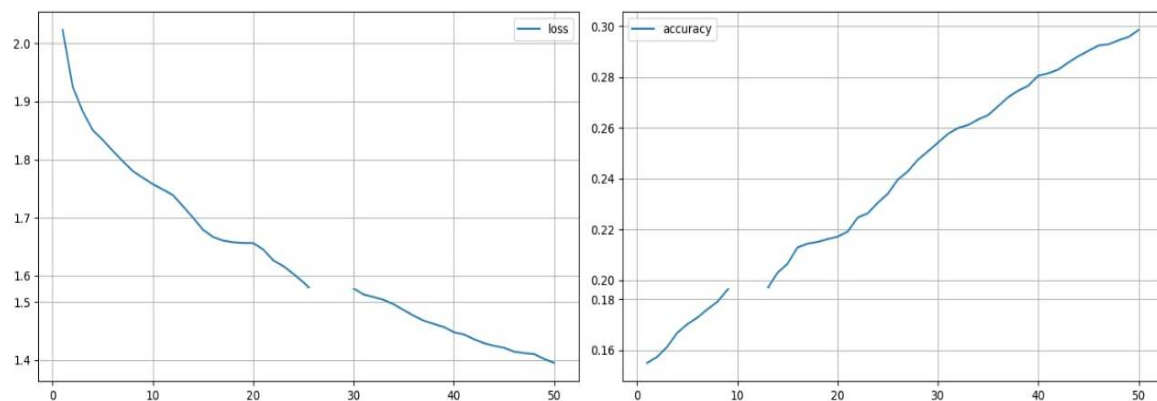


Fig. 3

Above picture represents our result.

This enables more effective content retrieval and indexing, facilitating the organization and retrieval of visual data in various applications.

## IV. CONCLUSION

Captioning image using deep learning has emerged as a highly promising and effective approach for generating textual descriptions for images using CNN, RNN and LSTM model along with advancements such as attention mechanisms and transfer learning, has led to significant improvements in the quality and accuracy of generated captions. This paper successfully implemented and trained the models on a suitable dataset, evaluated their performance using quantitative metrics, and discussed the obtained results. We compare our model against state-of-the-art methods and demonstrate its superiority in terms of caption quality, fluency, and relevance.

This paper showcased the potential of deep learning in addressing the challenging task of generating accurate and contextually relevant captions for images and videos. By leveraging the power of CNNs for visual feature extraction and RNNs for language modeling, the developed system demonstrated the ability to understand the visual content and generate descriptive captions.





## REFERENCES

- [1] Yang, L., & Hu, H. (2019). Adaptive syncretic attention for constrained image captioning. *Neural Processing Letters*
- [2] Fu, K., Jin, J., Cui, R., Sha, F., & Zhang, C. (2016). Aligning where to see and what to tell: Image captioning with region-based attention and scene-specific contexts. *IEEE transactions on pattern analysis and machine intelligence*
- [3] Li, J., Yao, P., Guo, L., & Zhang, W. (2019). Boosted Transformer for Image Captioning. *Applied Sciences*(19)
- [4] Oluwasanmi, A., Aftab, M. U., Alabdulkreem, E., Kumeda, B., Baagyere, E. Y., & Qin, Z. (2019). CaptionNet: Automatic end-to-end siamese difference captioning model with attention.
- [5] Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth. Every picture tells a story: —Generating sentences from images. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, 17. Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2017. Bottom-up and top-down attention for image captioning and vqa.
- [6] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for image and video description. In *CVPR*, 2015.
- [7] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2014. Show and tell: A neural image caption generator.
- [8] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Peng, et al. 2016. Google's neural machine translation system: "Bridging the gap between human and machine translation". *arXiv preprint arXiv:1609.07661*, 2016. Ga0, Klaus Macherrey, et
- [9] Fang, F., Wang, H., Chen, Y., & Tang, P. (2018). Looking deeper and transferring attention for image captioning. *Multimedia Tools and Applications*, 77(23).
- [10] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2014. Show and tell: A neural image caption generator.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)