# CAPtor - Neural Image Caption Generator

## Kartikey Singh[1], Navneet Kumar[2], Kunal Gautam[3], Prof. Rishikesh Yeolkar[4]

[1-3]Student, MIT School of Engineering, Pune
[4]Professor, MIT School of Engineering, Pune

---------------------------------------------------------------------***---------------------------------------------------------------------

**Abstract -** Enlivened by the most recent advancements in Deep Learning based Machine Translation and Computer Vision based Object Detection have prompted high precision Image Captioning models. Albeit these models are extremely exact, these will in general depend on the utilization of costly calculation making it hard to utilize these models progressively applications, where applications can utilize them. In this paper, we cautiously follow a portion of the heuristic techniques and center thoughts of Image Captioning and its regular strategies and present our basic succession to an arrangement based execution with a momentous change and productivity, for example, utilizing pillar search rather than avaricious inquiry that permits us to carry out these on low-end equipment. The proposed framework thinks about the outcomes determined utilizing an assortment of measurements with excellent models and dissects the explanations for the model prepared on the MS-COCO dataset that are missing because of compromise between calculation speed and quality. In this proposed framework ,RESTful API endpoint will be made to be utilized on any gadget with a web association like a cell phone, IoT gadgets, clock, and so forth, this endpoint used to sent a picture to the model running on distant worker which accordingly will produce and sent a subtitle portraying the articles and their relationship with one another in picture in a characteristic language.

**Key Words:** Neural Network, Assistive Vision, Caption Generator, Deep Learning, Restful API, Optimization.
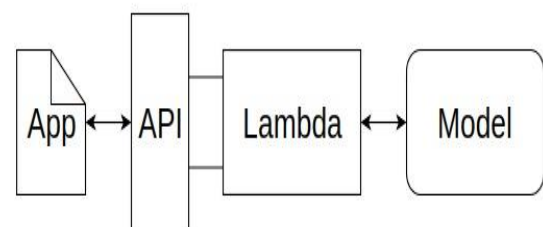
## 1.INTRODUCTION

Consequently, characterizing picture content and their connections or activities is a significant issue for man-made brainpower that associates PC vision and normal language handling. Be that as it may, this can profoundly affect assisting blind individuals with bettering comprehend their environmental factors. These photos can be utilized to deliver inscriptions that can be perused so anyone might hear to the outwardly disabled so they can more readily comprehend what's going on around them. This proposed framework gives an API endpoint which utilizes generative model dependent on a profound intermittent design that consolidates the most recent advances in PC vision and machine interpretation and that can be utilized to make characteristic sentences portraying a formerly caught or camera-caught picture. The model is prepared to build the odds of deciphering the sentence utilizing the Maximum Likelihood Estimation (MLE) given the preparation picture. What is generally amazing about this is that it is a solitary

end model that can be depicted as foreseeing subtitles, given an image, rather than requiring complex information arrangement or a pipeline of explicitly planned models.

Not exclusively should the model have the option to settle the PC vision difficulties of distinguishing objects in the picture, however it should likewise be sufficiently brilliant to catch and communicate object connections in the regular language. Hence, picture inscription age is viewed as a major issue for long. Its motivation is to imitate an individual's capacity to comprehend and deal with a lot of visual data in elucidating language, making it an alluring issue in the field of AI.

## 2. System Architecture

In this proposed framework, we are making a RESTful API with a solitary endpoint that will be utilized to give a picture to the Image Captioning model running on the worker. For making an API, we will utilize AWS API Gateway Service and AWS Lambda. AWS Lambda will have a capacity to send the picture got from the API solicitation to the Image Captioning model on the AWS Sagemaker.



The beginning stage of this framework will be an application that can run on any stage like a cell phone, smartwatch, or any IoT gadgets. This application will send a picture through the API solicitation to the AWS Lambda work. The image that is sent by the application will be a pre-caught picture or a picture caught by the camera gadget. The AWS Lambda capacity will be liable for moving this picture for additional preparing to the Image Captioning model where the picture is handled, and a fitting subtitle portraying that picture will be created. This subtitle will be in a content configuration. The API reaction will send this inscription back to the application, and this subtitle will be stood up boisterous by the gadget. Additionally, the subtitle will be shown on the gadget screen, if that gadget has a screen.

## 3. Datasets

For tackling this issue, there are many open source datasets accessible like MS COCO (containing 180k pictures), Flickr 30k (containing 30k pictures), Flickr 8k (containing 8k pictures), etc.[3]. In any case, with the end goal of this proposed framework, we have utilized Microsoft COCO dataset. This dataset contains roughly 180000 pictures with 5 inscriptions each.

## 4. Data Cleaning

When managing text, we regularly do some essential cleaning, for example, embeddings lower-instance of the relative multitude of words, eliminating unique tokens, eliminating words that contain numbers. Make a glossary of the relative multitude of interesting words that are available in all the 180000*5 (for example 900000) picture subtitles (corpus) in the dataset. As we are building a forecast model, we don't need every one of the words in our jargon yet words that may happen often or might be normal. This assists the model with being all the more impressive for anomalies and commits less errors [2].

## 5. Data Preprocessing

### 5.1 Data Preprocessing - Images

As we have pictures with various sizes, we will change over each picture into a fixed-sized vector that would then be able to be taken care of as contribution to the neural organization. For this reason, we will go for move learning by utilizing the InceptionV3 model that is the Convolutional Neural Network made by Google Research.

For performing picture characterization in 1000 distinct classes of a picture, this model was prepared on the ImageNet dataset. Be that as it may, we need to characterize the pictures as well as we need to get a fixed-length useful vector for each picture. This will be finished utilizing an interaction known as programmed include engineering [4].

For accomplishing this, we will disregard the last SoftMax layer from the model and concentrate a 2048 length vector (bottleneck highlight) for each picture.

### 5.2 Data Preprocessing - Captions

Forecast of the total inscription doesn't occur without a moment's delay. The inscription will be anticipated word by word. That is the reason we need to encode each word into a fixed-size vector. For addressing each one of a kind word in the jargon, we will utilize a file in the whole number structure.

## 6. Model Overview

The model is depicted in three sections:

### 6.1 Encoder

This is a VGG16 model that is pre-prepared on the ImageNet dataset. By pre-preparing the photographs with the VGG model and eliminating the yield layer, we will utilize the extricated highlights anticipated by this model as info. We have utilized VGG16 as an encoder on the grounds that VGG won the picture classfication challenge in ILSVRC thus we can tackle the cutting edge include extraction abilities of this model.

### 6.2 Sequence Processor

This is a word installing layer followed by a Long Short-Term Memory (LSTM) intermittent neural organization layer. This layer is for taking care of the content information.

### 6.3 Decoder

Both the component extractor and arrangement processor yield a fixed-length vector. To make a last forecast, these are united and handled by a Dense layer.

## 7. Sample Working

In the first round, we send the image vector and the first word ('begin') as inputs to the Sequence Processor and predict the second word, i.e.:

Input = Image_1_feature_vector + 'begin';

Output = 'This'

Then again, we provide the image vector and the concatted first two words as inputs and predict the third word, i.e.:

Input = Image_1_feature_vector + 'begin This';

Output = 'elephant'

And so on...

One image and one caption is not a single point of data but multiple data points depending on the length of the output (description / caption). For all data points, it is not just the image that goes as an entry in the system, but also the captions that are part of it that helps predict the next word in sequence.

## 8 Using Data Generators

In our preparation information, we have around 150000 pictures, each with 5 inscriptions. This makes a sum of around 750000 pictures and inscriptions. In spite of the fact that we expect that each subtitle on normal is 7 words in length, it will bring about information focuses which is equivalent to 750000 * 7 for example 5250000.

Presently regardless of whether we expect that one square takes 2 bytes, thus, to keep this information framework, we

will require in excess of 92 GB of memory. This is an immense interest, and regardless of whether we can stack this huge measure of information into RAM, it will make the framework run much more slow. Therefore, we are utilizing information generators.

We don't have to store all the information in each memory in turn. Despite the fact that we have a current arrangement of focuses in memory, it is sufficient for our motivation. Information Generators are a customary Python activity. The capacity of the generator in Python is utilized for this reason. It resembles an iterator re-working from where it left off last time it was called.

## 9 Hyper Parameter Tuning

The model prepared for 30 ages with a beginning learning pace of 0.001 advance size and 3 pictures for each group (clump size). Be that as it may, after multiple times, the learning level was decreased to 0.0001 and the model was prepared in 6 pictures for every set.

This typically bodes well in light of the fact that during the latest phases of preparing, as the model proceeds to meet combination, we need to bring down the degree of learning to make little strides towards the minima. Furthermore, expanding the group size after some time helps your angle updates to be more unique and amazing.

## 10.Model of the Proposed System

In this system, we have used the sequence to sequence to create an encoder-decoder architecture. The encoder is a pre-trained InceptionV4 Convolutional Neural Network and decoder is a Deep Recurrent Neural Network with long short term memory cells. Encoder InceptionV4 is used to convert raw images I into a Fixed length embedding F which represents the convolved features for the images. This embedding is obtained by running a forward pass to the one before the last layer, i.e., the average pool layer of the InceptionV4 model pool. The decoder in our model has two phases, named, training and inference. The decoder is responsible for learning word order given the convolved features associated with the original captions. The hidden state of the decoder ht is initialized using these image embedding features F in timestep t=0.. Thus the basic concept of the encoder-decoder model is shown by the following equations [1].

$F = encoder(I); X_{t=0} = F; O_t = decoder(X_{t:0 \to t})$

The RNN training process with LSTM Cell-based decoder works on a probabilistic model where the decoder increases the chances of a word p in captions given a convolved image features F and previous words $X_{t:0 \to t}$. To learn the entire sentence of length N corresponding to the F features, the decoder uses its repetitive nature to loop over itself over a constant number of timesteps N with the previous

information (features and words sampled in timestep t) stored in its cell memory as a state. The decoder can change Ct memory as it unrolls by adding a new state, refreshing or forgetting previous state with LSTM forgetting ft , input it, and output ot memory gates.

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

$$i_t = \sigma \ (W_i \ . \ [h_{t-1}, x_t] + b_i)$$
$$(2)$$

$$c_t = \sigma \ (W_c \ . \ [h_{t-1}, x_t] + b_c)$$
$$(3)$$

$$C_t = f_t * C_{t-1} + i_t * c_t$$
$$(4)$$

$$o_t = \sigma \ (W_o \ . \ [h_{t-1}, x_t] + b_o)$$
$$(5)$$

$$h_t = o_t * tanh(C_t)$$
$$(6)$$

$$O_t = argmax(softmax(h_t))$$
$$(7)$$

$\sigma \to$ sigmoid; $O_t \to$ Output word; tanh $\to$ hyperbolic tangent; $W_o, W_f, W_i \to$ Learnable Weight Vector; $b_o, b_f, b_i \to$ Learnable bias Vector;

## 11.Training

For offline evaluation, compared to other Caption Bots our implementation will use Batched data, supports CNN finetuning, uses TensorFlow, and runs on a GPU[5]. All of these together will increase the training speed of the Language Model to a greater extent($\sim$100x). Even if the split of 5000 images is not a standardized split, many researchers have been using it for reporting their results.
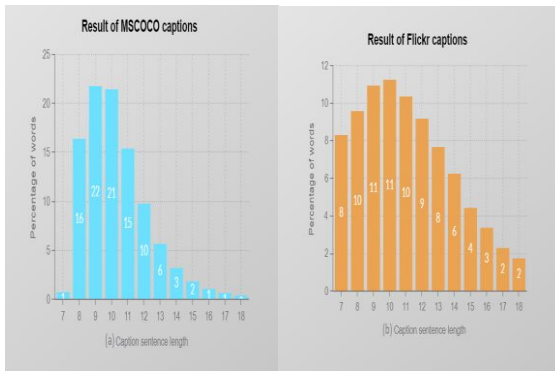
## 12. Result

### 12.1.Datasets

For assessing the presentation of our strategies, we are utilizing the MS COCO [8], Flickr8K[9], and Flickr30K[10] datasets. They are the most famous and productive datasets for estimating the exactness of the created data.

Table 1 shows the definite examinations of reference inscriptions on the three datasets above. The MS COCO dataset is an enormous scope object location, division, and inscribing dataset. The authority adaptation of MS COCO dataset incorporates more than 82000 preparing pictures, more than 40000 approval pictures, and more than 40000 test pictures. Since the "Karpathy'' split is the most ordinarily utilized parted technique for announcing results, we use it to part the authority MS COCO dataset to acquire

around 113,000 preparing pictures, 5000 approval pictures, and 5000 test pictures. The Flickr8K dataset accompanies an authoritatively split of 6000 pictures for preparing, 1000 pictures for approval pictures, and 1000 pictures for testing. With no authority split, the Flickr30K dataset has 31,783 pictures that we will part into 25,000 preparing pictures, 2000 approval pictures, and 3000 pictures for testing.



a) shows the factual aftereffects of reference inscriptions on MS COCO dataset;

(b) shows the measurable consequences of reference subtitles on Flickr8K and Flickr30K datasets. The x-pivot addresses the length of each inscription sentence.

## 12.2. Evaluation Metrics

We test the viability of our strategy with some notable measurements utilizing picture inscriptions, including CIDEr, CIDEr [11], SPICE [12], METEOR [13], ROUGE-L [14], and BLEU [15].

Juice and SPICE are both human agreement measurements. Juice can be utilized to quantify the likeness between a created inscription and a bunch of definitions composed by people. Moreover, SPICE is a restrictive measurement, which is utilized to test how well an organized sentence catches items, traits and connections between them. METEOR ascertains sentence level likenesses as per the consonant meaning of uni-gram review and accuracy. ROUGE-L can be utilized for gisting assessment. Its scores are determined by estimating the quantity of dispersed units, for example, n-gram, word request, and the adjustment of words between created subtitles and human composing sentences. BLEU is an allegory that is generally utilized in machine interpretation tasks. It depends entirely on the cohesiveness of the n-grams.

## 13. Conclusion

The outcomes above show that the calculation we had utilized fits precisely to our proposed framework and gives preferable outcomes over the past techniques utilized for picture subtitle age.

## REFERENCES

1. Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan (2015) Show and tell: A neural image caption generator. CVPR 1, 2
2. K. Xu (2016) Show, attend and tell: Neural image caption generation with visual attention. inProc. Int. Conf. Mach. Learn.
3. Andrej Karpathy, Li Fei Fei (2015) Deep Visual-Semantic Alignments for Generating Image Descriptions. IEEE Transactions on Pattern Analysis and Machine Intelligence (April 2017), vol 39, issue 4:664–676.
4. Image Captioning with Keras teaching computers to describe pictures
5. Ayush Yadav, Rutgers the State University of New Jersey (2017) Camera2Caption: A real-time image caption generator.
6. Oriol Vinyals, Alexander Toshev, Samy Bengio, Dumitru Erhan (2016) Show and tell: Lessons learned from the 2015 MSCOCO image captioning challenge PAMI
7. Anderson, P.; Fernando, B.; Johnson, M.; Gould, S. SPICE: Semantic Propositional Image Caption Evaluation. In ECCV; Springer: Cham, Switzerland, 2016; pp. 382–398.
8. 13 Denkowski, M.; Lavie, A. Meteor Universal: Language Specific Translation Evaluation for Any Target Language. In Proceedings of the Ninth Workshop on Statistical Machine Translation, Association for Computational Linguistics, Baltimore, MD, USA, 26–27 June 2014; pp. 376–380.
9. Lin, C.Y. ROUGE: A Package for Automatic Evaluation of Summaries. In Proceedings of the ACL-04 Workshop Text Summarization Branches Out, Barcelona, Spain, 25–26 July 2004; pp. 74–81.
10. Papineni, K.; Roukos, S.; Ward, T. BLEU: A method for automatic evaluation of machine translation. In Proceedings of the 40th Annual Meeting on Association for Computational Linguistics, Philadelphia, PA, USA, 7–12 July 2002; pp. 311–318. 20. Farhadi, A.; Hejrati, M.; Sadeghi, M.A.