# Machine-Generated Captions for Images using Deep Learning

Karwan Vishweshwar[1], Kosanam Srinivas[2], Ms.C.A. Daphine Desona Clemency, M.E.[3]

[1]Department of Computer Science and Engineering, Sathyabama Institute of Science and Technology

**Abstract—The primary objective of the picture caption generator is to automatically produce a suitable text or caption in English. The system's primary goal is to successfully provide appropriate captions for the provided picture. This study presents an image caption generator that, given an input picture, would identify its contents using beam search and greedy search to produce an English phrase. A pretrained deep learning CNN architecture exception model is used to learn image features, while a LSTM model is used to learn textual features, then integrates the results of both to produce a caption. To produce words, phrases, or captions for the provided photos, we use the LSTM model. Using the Convolutional Neural Network with Long Short-Term Memory, this model was created to create a caption generator for images. Features are extracted from the picture using a pre-trained version of VGG16. To create descriptive text for the pictures, LSTM acts as a decoder. This model has been taught to produce descriptive captions or words based on an input picture. The effectiveness of the model is measured by means of blue scores given to the system. The Keras library, NumPy, and Jupyter notebooks are discussed as tools for developing this project. We also talk about the picture categorization task, how CNNs are employed, and the Flickr dataset.**

*Keywords—: Deep Learning, LSTM, Caption, Description, Memory, Neural Network, VGG16, Image, CNN*

## I. INTRODUCTION

The field of servo descriptions for images via transfer learning is rapidly growing since it combines the best of artificial intelligence and NLP. The goal of this study is to develop an algorithm to analyze an image's data and provide a human-readable description of it. The initial stage in most efforts is to collect a large dataset consisting of images and their explanations, since it will be used to train the deep learning algorithm. A popular form of the model combines a run trained to generate language processing tags with a CNN was built to decipher the meaning of a picture. In order to generate captions for images, a CNN must first decide what elements should be retrieved from them. Many other topologies, including converter, multiplexer, even robot architectures, may be used to train this model. By focusing the model's attention on certain parts of the image, transformer-based systems like BERT and GPT-2 may increase caption quality. After the model has been trained, it may be used to generate new captions for unseen images. The degree to which the generated captions match the reference caption may be measured using tools like BLEU, METEOR, ROUGE, and CID. One important challenge is the lack of a large collection of images with text annotations. An additional challenge is developing tags that aren't just comprehensible but also semantically correct and worded in a natural manner. This new technology may be put to use in a variety of contexts, including but not limited to sight searches and retrieval, robotic photo annotation, and devices for the visually impaired. The advancement of computational modeling and natural language recognition has led to an increase in the frequency and quality of machine-generated captions. There are a few approaches that may be taken to automatically generate captions for images using deep learning. The model may also generate captions using a "template-based" approach, in which examples are used as a basis for new captions. This strategy is effective for creating captions for a specific media, such as video games or stock images, but it has its limits. Another approach is the "free-form" technique, in which the model generates captions without relying on a certain format. With this approach, the model can generate more informative descriptions, although it could be harder to train & evaluate. The advantages of both the framework and unlimited approaches are combined in the "hybrid" method. Using this approach might improve the quality of the captions generated by providing grammatically correct and semantically precise captions
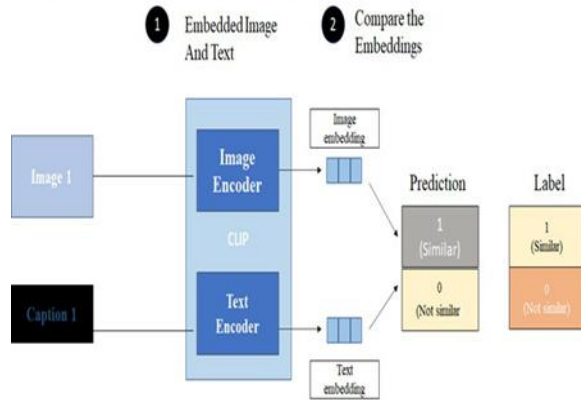
Fig 3. Extract Features From Images And Text

It's also worth noting that various pre-trained models for writing captions for photographs using deep learning are available, including COCO-Caption from Microsoft, Show and Tell from Google, and Dense Captioning from Facebook. Although these methods are already trained, they may be tweaked for optimal performance on specific data. Generally speaking, machine-generated photo captions using deep learning techniques are a rapidly growing field with enormous potential. Future enhancements are expected as a result of research and development efforts in DL and NLP techniques. People communicate with one another via the use of language, whether written or spoken. They usually talk about what they notice using this language. Images and signs may help someone with vision impairments express themselves and learn new information. Automatically creating descriptive sentences from a photograph may help and have a substantial impact on the capacity of visually impaired folks to comprehend the explanation of pictures on the web, despite being a sophisticated and demanding operation [1]. Whenever someone gives a very vivid description of anything, they are said to have painted a "picture" in the listener's mind. The process of creating mental pictures may be quite helpful in developing sentences. Humans can properly identify visuals after just a short exposure. Analysis of existing natural visual representations may help achieve complex human recognition goals. Image classification & object detection and recognition are far easier problems to solve than automated captioning and description. When characterizing an image, it is important to take into account not only the objects, their characteristics, and the activities they do, but also the connections between them [20]. To this end, most of the early work in facial information has

concentrated on assigning labels to images based on specified categories, which has led to tremendous progress. An effective and simple metaphor for assuming is provided by closed visual concept vocabulary. In comparison to the vast mental potential of humans, these concepts seem pitifully little. However, a model of language use in speech perception is necessary, and real dialects like English ought to have been used to convey the aforementioned semantic information. The majority of previous work on automatic description generation from images has argued for a blend of several methods for dealing with the aforementioned problem. As shown in Fig.1, we will instead focus on creating an unified theory that, when fed a picture, can be developed to spit out a set of descriptive terms. The emphasis in natural language processing is often shifted from the importance of images to the way they are represented [18], and this is the case with parsing. One of the main purposes of summarizing is to choose or create a summary for a piece of writing. For the text recognition challenge [21], the objective is to provide a statement that covers a wide range of the visual content.
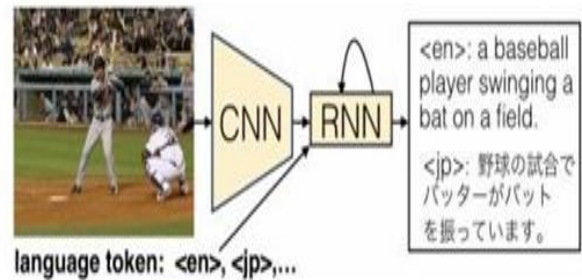


Fig. 1. Model based on Neural Networks

In this study, we offer a paradigm for doing just that; coming up with fresh descriptions of pictures. We've done this using the Flickr 8k dataset, which has 8,000 photos and five descriptions each. Figure 2 depicts the organizational scheme of the dataset, with one picture and its five corresponding natural language descriptions. The approach involves the use of both CNNs & RNNs. Images are classified using a CNN that has been pre-trained. An image encoder, this network is used to convert input images into a desired format. Recurrent Neural Networks take their input from the last hidden layer (RNN). To put it simply, this network is a translator that can produce new phrases. It appears that the produced sentence sometimes loses its focus or predicts a different phrase than the one that

best fits the original visual content. This statement has a tenuous connection to the input picture since it was constructed using a description that occurs often in the dataset.

## II. LITERATURE REVIEW

Creating natural language descriptions from visual input is a well-studied subject in computer vision [1]



Fig. 2. Caption 1:A group is sitting around a snowy crevasse, Caption 2: A group of people sit atop a snowy mountain, Caption 3: A group of people sit in the snow overlooking a mountain scene, Caption 4: Five children getting ready to sled, Caption 5: Five people are sitting together in the snow.

[2]. There are essentially three schools of thought in the research around the topic of automatic picture captioning. For starters, there are template-based approaches, which may be found in references [4] - [7]. Object, action, scene, and attribute detection are prioritized in this method. Methods for generating captions that rely on a transfer are the focus of the second group [8]. It is used to retrieve images. Images that are visually comparable are retrieved using this method, and their captions are then applied to the query picture. The majority of studies [10] have shown that using neural networks for translation and artificial neural models and caption creation are both effective. Rather than just translating a text from one language to another, the purpose here is to provide an explanation for the visual. There has been an increase in system complexity as a result of this. They use a formal language to communicate and are composed of visual radical classifiers (e.g., "and-"). Further transformation may include the usage of a graph or logic system, or a rule-based system. An image's description may be generated using a holistic persistent network model, as proposed by researchers like Mao [11] & Karpathy [12]. Using the NIC model, Vinyals, Oriol. In the NIC paradigm, CNN serves as the encoding method. To classify images, we use an RNN decoder to process output from a pre-trained

convolutional neural network (CNN). This RNN decoder has the potential to even generate sentences. For this purpose, we have used LSTM [1], a high-end RNN version. Xu [13] has suggested synthesizing photoreceptor focus into to the LSTM process to help keep attention on different things while the algorithm produces appropriate phrases. Generating natural-sounding captions for images involves a complex conventional neural model. All save the most attempting to cut systems employ an encoding-decoding framework [13] that combines caption creation with attention. In this research, we focused on the third category of captioning methods. To do this, a neural system is being developed to provide the visual representations in plain language. The use of CNN for picture encoding is commonplace. After being pre-trained for the task of classifying images, the RNN decoding accepts the data of the last tier as input and outputs the phrase.
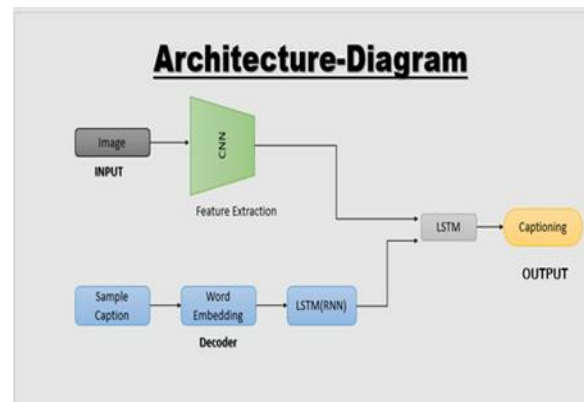
## III. OUR APPROACH



Fig 4 Architecture Diagram

The purpose of this research was to create an autonomous method for writing picture captions using a neural network trained on probabilistic concepts. By using a powerful statistical model, it may be possible to increase the likelihood of an appropriate translation occurring during reasoning and learning. A. CNN Today, we use CNN models for anything from speech recognition to facial recognition. The core of a CNN is a series of interconnected convolutional neural networks. A multi-layer neural program's completely connected layers emerge after the inversion procedure [14]. To make advantage of the 2-d nature of the input image, the CNN was designed. To do this, we use a large number of locally - relevant linkages and related

data generated using a diversity of pool techniques that are robust to interpretation. CNN's main advantages lie in the fact that it is easy to train and requires fewer settings than other systems with the same number of state variables. For this research [15], we be using a Deep CNN (hence referred to as VGG) net for widespread image recognition. It's available in 16 - & 19- layer variations. Class error values for 19 and 16 layers are really quite similar for both the cross-valid across-validation test set, sitting at 7.4 percent & 7.32 percent corresponding. By feeding an image into the model, it may provide details about it that can be used in the caption generation process. B. LSTM When it comes to NLP & computer vision, recurrent neural networks (RNNs) like the LSTM are often used for tasks like picture captioning. With their ability to "remember" previous data, LSTMs can better grasp the context of new input, making them ideal for sequential data applications like picture captioning. Specifically for the task of visual caption, a Long ShortTerm Memory based model has been trained using a dataset consisting of pictures and their accompanying captions. Using the picture characteristics retrieved and the learnt correlations between image attributes and terms from the trained captions, the model can then produce captions for fresh photos. In order to capture the transient dynamics of a set of objects, researchers have turned to recurrent neural networks [17]. Weights and gradients in a regular RNN tend to vanish and explode, making it hard to understand long- term dynamics [9]. In an LSTM, the "brain" is the storage unit. Existing values are stored for a very long period in the future. The gates' function is to control how often the cell's state is updated. Variations in the amount of connections among memory blocks & gateways stand in for independent variables. Figure 3 depicts the LSTM architecture, upon which our model relies. This design has no peepholes. It can be shown that there are connections between LSTM's memory cells and their gates in the ways described below:
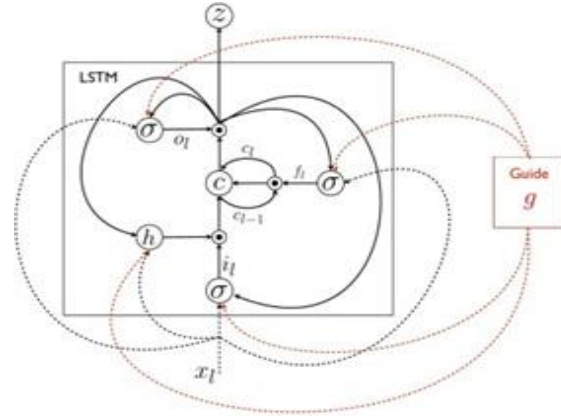
$$i_l = \sigma(W_{ix}x_l + W_{im}m_{l-1}) \tag{1}$$

$$f_l = \sigma(W_{fx}x_l + W_{fm}m_{l-1}) \tag{2}$$

$$o_l = \sigma(W_{ox}x_l + W_{om}m_{l-1}) \tag{3}$$



Fig. 3. Connection diagram of LSTM [9]

$$c_l = f_l \odot c_{l-1} + i_l \odot h(W_{cx}x_l + W_{cm}m_{l-1}) \tag{4}$$

$$m_l = o_l \odot c_l \tag{5}$$

$$L(I,S) = -\sum_{t=1}^{N} log(p_t(S_t)) \tag{6}$$

## IV. LSTM-BASED SENTENCE PRODUCTION

The neural network employs an encoder-decoded strategy, used in IOT networks & computational linguistics [1,11,12,13,16] to generate new phrases. Specifically, a set of simple English sentences is encoded as a value is called in this modeling. Next, a decoder uses the translated matrices to generate a new given text in the target language. The learning process's end result is a version that sounds and appears as if it was originally written in the target culture. To maximize the number of captions for a given image, if xi is the file's identification & s1 is the caption function, this method of creating captions is used. In a well - constructed phrase of length Li, the group of phrases Li represents the hyperparameters. Moving on, we ignore the superscript I when it is not relevant or has been taken out of context. The Bayesian cyclical concept may be used to dissect phrases into their individual words, since each is composed of a string.

$$arg\theta \ \Sigma i \ log(p(S_1 : L_i \ / \ x^i \ , \ \theta))$$

$$log(p(s_{1:L}|x,\theta)) = log(p(s_1|x,\theta)) + \Sigma_{l=2} \ log(p(s_1|x,s_{1:l-1},\theta))$$

where S1: L is the part of the constructed sentence up toward the lth word. To maximize the target in Eq. 7, we specify the document $log(p(S1:Li|xI))$ as a combination of the prior hidden in RNNs at all times

throughout training. The pdf of words across the entire vocabulary at linear interpolationl+1 may be calculated using the output ml from the memory space by using the formula pl+1 = z (ml) = [1]. Images and text are encoded into fixed-length vectors and utilized as inputs to a LSTM network. Each image's CNN characteristics are computed, and then fed into a kind of reinforcement learning. The combination of a string of words and an image in a phrase result in a novel structure. The image acts as the jumping-off point for a new career, while the words constitute the remainder. By periodically executing the recursive connect for l among 1 & Li, the Short short-term memories are trained and use this new sequence. The word embed matrix, LSTM tensor, and image feature linear activation matrix are all elements of a neural model. The principal of the head line model's 3 main sub models is the visual subsystem, which employs a feature map for the picture 29 repetitions with proportions of 29x4097 (where 29 is the greatest number of words in a caption). The second approach, a perceptron (NLP with an one LSTM unit, generates a 28-by-256 grid, with 128 is the return dimension of the LSTM layer; the third model takes these two matrices as input and feeds data into a final LSTM unit, which has proportions of 28-by-915. For training, we utilize an unaltered encrypted word vector as the candidate solution; for tests, we use a separate coded text matrix and append the feature map from the test image, resulting in a matrix with size of 29 by 915.

## V. RESULTS

A. Such sets of data include photos and textual descriptions of those images written in a natural language like English. Table I displays the metrics of sets of data. Experts in these data provide unbiased, eye-catching descriptions of each picture using a total of five phrases.

### TABLE I
### DATASET STATISTICS

| Dataset Name | Size | | |
|---|---|---|---|
| | Train | Valid | Test |
| Flickr8k [1] | 6000 | 1000 | 1000 |
| Flickr30k [1] | 28000 | 1000 | 1000 |
| MSCOCO [1] | 82783 | 40504 | 40775 |

B. Outcomes Fifty training epochs have been completed on the model. As additional eras are

utilized, the loss may be reduced to 3.74. Taking the amount of data into account requires additional epochs to get reliable findings.



Fig. 4. Selection of Evaluation Results

## VI. CONCLUSION

In this research, a model is introduced; it is a multilayer perceptron that can immediately analyze an image and provide meaningful captions in human languages including English. Program is taught to produce a word or description from an image. The categories below outline how model-derived captions and descriptions are often used. Error-free description, Minimal inaccuracies Slightly off-image description Completely off-image description. Terms like "vehicle" and "car" and "taxi" and "cabbie" and "limo" Extensive testing shows that a better best model may be achieved by using a bigger training dataset. Accuracy is improved and losses are decreased because of the increased dataset. It's also intriguing to consider how unsupervised information, in the form of both photos and texts, may be utilized to enhance current methods for generating captions for images.

## REFERENCE

[1] Vinyals, Oriol, et al." Show and tell: A neural image caption generator." Computer Vision and Pattern Recognition (CVPR), 2021 IEEE Conference on. IEEE, 2021.

[2] Gerber, Ralf, and N-H. Nagel." Knowledge representation for the generation of quantified natural language descriptions of vehicle traffic in image sequences." Image Processing, 1996. Proceedings., International Conference on. Vol. 2. IEEE, 2019.

[3] Yao, Benjamin Z., et al." I2t: Image parsing to text description." Proceedings of the IEEE 98.8 (2020): 1485- 1508.

[4] Farhadi, Ali, et al." Every picture tells a story: Generating sentences from images." Euro-pean conference on computer vision. Springer, Berlin, Heidelberg, 2022.

[5] Yang, Yezhou, et al." Corpus-guided sentence generation of natural images." Proceedings of the Conference on Empirical Methods in Natural Language Processing. Association for Computational Linguistics, 2018.

[6] Kulkarni, Girish, et al." Babytalk: Understanding and generating simple image descriptions." IEEE Transactions on Pattern Analysis and Machine Intelligence 35.12 (2019): 2891-2903.

[7] Mitchell, Margaret, et al." Midge: Generating image descriptions from computer vision de-tections." Proceedings of the 13th Conference of the European Chapter of the Association for Computational Linguistics. Association for Computational Linguistics, 2020.

[8] Kuznetsova, Polina, et al." Collective generation of natural image descriptions." Proceed-ings of the 50th Annual Meeting of the Association for Computational Linguistics: Long Papers-Volume 1. Association for Computational Linguistics, 2019.

[9] Jia, Xu, et al." Guiding long-short term memory for image caption generation." arXiv pre-print arXiv:1509.04942 (2022).

[10] Bahdanau, Dzmitry, Kyunghyun Cho, and Yoshua Bengio." Neural machine translation by jointly learning to align and translate." arXiv preprint arXiv:1409.0473 (2018).

[11] Mao, Junhua, et al." Deep captioning with multimodal recurrent neural networks (m-RNN)." arXiv preprint arXiv:1412.6632 (2022).

[12] Karpathy, Andrej, and Li Fei-Fei." Deep visual- semantic alignments for generating image descriptions." Proceedings of the IEEE conference on computer vision and pattern recog-nition. 2019.

[13] Xu, Kelvin, et al." Show, attend and tell: Neural image caption generation with visual at-tention." International Conference on Machine Learning. 2019.

[14] El Housseini, Ali, Abdelmalek Toumi, and Ali Khenchaf." Deep Learning for target recognition from SAR images." Detection Systems Architectures and Technologies (DAT), Seminar on. IEEE, 2017.

[15] Simonyan, Karen, and Andrew Zisserman." Very deep convolutional networks for large-scale image recognition." arXiv preprint arXiv:1409.1556 (2020).

[16] Donahue, Jeffrey, et al." Long-term recurrent convolutional networks for visual recognition and description." Proceedings of the IEEE conference on computer vision and pattern recognition. 2021.

[17] Lu, Jiasen, et al." Knowing when to look: Adaptive attention via a visual sentinel for image - captioning." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Vol. 6. 2017.

[18] Ordonez, Vicente, Girish Kulkarni, and Tamara L. Berg." Im2text: Describing images us - ing 1 million of the captioned photographs." Advances in neural information processing systems. 2019.

[19] Chen, Xinlei, and C. Lawrence Zitnick. "Mind's eye: A recurrent visual representation for image caption generation." Proceedings of the IEEE conference on computer vision and pattern recognition. 2021.

[20] Feng, Yansong, and Mirella Lapata." How many words is a picture worth? automatic caption generation for news images." Proceedings of the 48th annual meeting of the Association for Computational Linguistics. Association for Computational Linguistics, 2019.

[21] Rashtchian, Cyrus, et al." Collecting image annotations using Amazon's Mechanical Turk." Proceedings of the NAACL HLT 2010 Workshop on Creating Speech and The Language Data with Amazon's Mechanical Turk. Association for Computational Linguistics, 2020.