

Image Captioning Using Convolutional Neural Networks And Recurrent Neural Network

S.SAI KRISHNA REDDY¹, VADDI SRIVALLIDEVI²

¹MCA Student, B V Raju College, Kovvada, Andhra Pradesh, India.

²Assistant Professor, B V Raju College, Kovvada, Andhra Pradesh, India.

ABSTRACT

In recent years, research on natural language processing and computer vision has become increasingly interested in the problem of automatically synthesising descriptive phrases for photos. Image captioning is the process of providing a written description for an image. Both computer vision and natural language processing are used to create the captions. The authors suggest a hybrid approach that combines a multi-layer Convolutional Neural Network (CNN) for creating image-descriptive vocabulary with a Long Short-Term Memory (LSTM) for precisely forming coherent sentences using the study's created keywords. A Deep Learning algorithm that makes use of convolutional neural networks is called a Convolutional Neural Network (ConvNet/CNN). For this issue, there are numerous open-source datasets accessible, like Flickr8k (which contains 8k photos), Flickr30k (which contains 30k images), MS COCO (which has 180k images), etc.

INTRODUCTION

The great artificial intelligence challenge of inscription generation is to create a descriptive text for a given picture. It makes use of two computer vision techniques to comprehend the image's content and a language model from the field of natural language processing to translate the understanding of the image into the appropriate string of words. There are several uses for captioning images, including suggestions for editing and improving software, use in virtual assistants, image indexing, accessibility for people with

visual impairments, social networking, and a range of various other NLP applications. The majority of conventional photo captioning systems use an encoder-decoder structure, which converts an input image into a preliminary representation of the information it contains before being decoded into a string of descriptive text that was motivated by neural machine translation. A single CNN feature vector output or a large number of aesthetic features that were really collected from multiple areas of the image may be present in this encoding. The areas can

either be equally checked or specifically targeted by an item detector to improve performance. In the item connection transformer we proposed, self-attention is represented. The discovered items' openness and their bounding box alter depending on the interest weight in comparison to the chair that is highlighted in red. Our model depicts connections between this chair and its companion chair to the left, the beach beneath them, and the umbrella above them in the caption that is created. Review in the vein of Object Connection Transformer. The use of item spatial partnership modelling for picture captioning, specifically inside the Transformer encoder-decoder architecture, is recommended in this work and is also shown. The modifications made to the Transformer design are visible in the Bounding Box Relational Encoding layout. To do this, the item relation component of [9] and the Transformer encoder are combined. The contributions of the paper are as follows:

- In the sections that follow, we will go through the Object Relation Transformer, an encoder-decoder architecture designed specifically for image captioning that incorporates information about the spatial links between input

recognised objects using geometric interest.

- We quantitatively demonstrate the value of geometric interest using baseline contrast and an ablation experiment on the MS-COCO dataset.

- Last but not least, we qualitatively demonstrate how geometric interest can be used to provide enhanced captions that demonstrate enhanced spatial awareness.

PROJECT OBJECTIVES

We must first clearly understand the importance of this issue. Let's look at a few scenarios when locating a solution to this issue can be beneficial. self-driving vehicles Automated driving presents some of the most difficult challenges, thus captioning the area around the vehicle can help the system. assistance for the blind We can develop a product that will help the visually impaired navigate the roadways on their own. This can be accomplished by first turning the scenario into a message, then adding speech to the text. These days, everyone is familiar with both Deep Understanding applications. CCTV cameras are already widely dispersed, but if we could add helpful captions to them instead of just watching the world, we could send out alerts as soon as

illegal activity is observed anywhere. It's possible that by doing this, accidents or criminality will be reduced. Automated captioning might help Google Picture Look improve to the level of Google Look because each image might be turned into a subtitle before being viewed.

1.1 Data Set Holiday

For this subject, there are many free source datasets available, including Flickr8k (8k images), Flickr30k (30k images), MS COCO (180.000 images), and others. But for this study, I've used the Flickr 8k dataset, which you can obtain from the University of Illinois at Urbana-Champaign by filling out this form. It might also be challenging to train a design with that many photographs using a device other than a high-end PC or laptop. Since we have learned in the Intro part that an image can have many subtitles, all of which are simultaneously relevant, this collection contains 8000 photos, each of which has five captions.

Analysis of the Data

From Kaggle, we downloaded and installed the images, text files relating to the images, and data. The names of all the images and their five captions are included in the file "Flickr8k.token.txt," which is one of the data. What do you see to be visible in

the image below? Most certainly, any of these and possibly many more subtitles would be applicable for this picture. Many of you would undoubtedly respond, "A white bird soaring," while others might add, "A white bird with black patches" or, "A crane is flying over a river." However, the point I'm trying to make is that it's really simple for us humans to look at a picture and describe it in the right terms. This can be easily finished by a 5-year-old child as well. Furthermore, are you able to create a computer programme that takes a photo as input and outputs the best caption?

2. Work pertaining to captioning photos

Since the creation of the Internet and also its widespread use as a platform for photo sharing, the issue of picture captioning as well as potential solutions has really existed. Scientists have offered numerous concepts and techniques from a variety of angles. Krizhevsky et al. created a very dependable and also first semantic network with non-saturating neurons and also a GPU-based convolution function. They reduced overfitting by using a regularisation technique called failure. Their neural network features maxpooling layers in the middle and a 1000-way softmax layer at the end. It

has also been suggested to create a family of attention-based photo captioning systems.

While top convolutional layers of a CNN frequently attract visual attention, the spatial localization is restricted and frequently not semantically adequate, which attempts to tie specific areas in the image to the words in the projected inscription. Anderson et al. in combined a "bottom-up" attention version with a "top-down" LSTM to address this issue of standard interest designs, which is most similar to our approach. Hu et al. invented the first system for item identification. This systematic literature review is organised, carried out, and reported in steps. First, we underlined the need of conducting this research on the drawing board. At this point, a search strategy, high standards for assessment criteria, information sifting techniques, and the identification of study topics are all prepared. Before beginning, we carefully organised our enquiry.

3. INSPIRATIONS

We must first grasp the functional significance of this issue. Let's look at a couple situations where finding a resolution to this problem can be quite useful. One of the most challenging challenges for a self-driving system is autonomous driving, but

captioning the environment around the car can help. assistance for the blind We can develop a product that will enable blind people to navigate the roadways independently. To accomplish this, first translate the scene into a message, then translate the voice from the text. These days, many people use both Deep Understanding programmes. CCTV cameras are already everywhere, but if we can enhance them with useful captions in addition to simply watching the world, we can send out alerts as soon as illegal activity is observed anyplace. This may result in a decrease in accidents and/or criminal activity.

4. CONVERSATION.

Since every image can have a caption added before being searched, automatic captioning could help Google Photo Browse catch up to Google Browse.

In the area of profound knowing, it has consistently been a crucial and basic responsibility. The act of captioning photographs serves a variety of objectives. Captioning can be thought of as an end-to-end Series to Sequence hurdle because it changes photographs from pixels to words. Automatic photo annotations have potential advantages from both a practical and scholarly perspective. The vast amount of content that is accessible on the internet is the

most important factor in the current process of societal progress. A sizeable portion of these data, the majority of which are distinct from standard data, are made up of media information. They are frequently developed by online services like social networks and informational portals. Aid with aesthetics for the blind: If we can attach a camera to the head of a blind person, the camera will take pictures from which we can extract the necessary subtitles. These subtitles can then be translated into text and, later, sound, which the blind person can utilise to understand what he is viewing. Self-driving vehicles Automobiles that drive themselves are autonomous decision-making systems. Following that, this data is modelled using deep learning algorithms, which make choices depending on the present environments of the vehicles, such as creating pertinent captions and choosing the optimal course for the vehicles. They are capable of handling data streams from a wide range of sensors, including cameras, RADAR, and GPS. Electronic cameras are widely used today, but if we can provide critical information in addition to seeing the entire world, we can broadcast alerts as soon as potentially dangerous behaviour is discovered. By doing this, it's feasible to reduce accidents or illegal activity.

PRESENT SYSTEM

One of the most important pieces of art challenges the three spaces of the image, the message, and the language.

The mapping from the picture and sentence rooms to the meaning rooms determines whether the final subtitle makes sense.

We compare the degree of similarity between the photos and the generated sentences to determine how precise the inscription is. Following that, the results are saved as triplets (image, action, and object), and the game is calculated. If the sentence and picture are almost same, the final output score will also be higher.

BENEFITS: 1. This version was less accurate and had a lot of problems.

RECOMMENDED SYSTEM

We propose and also demonstrate the application of objects spatial relationship modelling for photo captioning in this work, particularly inside the Transformer encoder-decoder style. The Transformer encoder is integrated with the item relation component of [9] to achieve this. The contributions of the paper are as follows:.

In this study, we advocate and demonstrate the use of object spatial relationship modelling for photo captioning, specifically inside the Transformer encoder-decoder design.

The Transformer encoder is integrated with the item connection module of [9] to achieve this. The payments for the paper are the complying with.

Here, we describe the Item Relation Transformer, an encoder-decoder design created especially for image captioning that incorporates knowledge of the spatial relationships between input recognised objects using geometric attention.

- We quantify the value of geometric interest using baseline contrast and an ablation experiment on the MS-COCO dataset.

Finally, we qualitatively show that geometric attention can be used to produce better captions that indicate better spatial awareness.

SURVEY OF LITERARY WORKS

Recent advancements in NLP, such as the Transformer style, have led to notable efficacy increases in the fields of translation [23], text production [4], and language understanding [19]. [23] [22] used the Transformer to carry out picture captioning. The authors investigated the extraction of a single global picture attribute as well as uniformly sampling features by dividing the shot into 8x8 parts. The Transformer encoder acquired the feature vectors in order in the previous instance. In this study, we propose a bottom-up strategy

[2] to enhance this consistent tasting. The Transformer design is particularly well suited as a bottom-up visual encoder for captioning because, unlike an RNN, it lacks a sense of order for its inputs. However, positional encoding, which we apply to the subtitle message's decoded tokens, makes it possible to explain sequential information clearly. Instead of assigning a numerical order to three objects, our Object Relation Transformer seeks to capture how two points are physically related to one another and weight them appropriately.

The transformer-based version used the Convolutional Semantic Network (CNN) InceptionV3 Szegedy et al. [2016] to extract features from images. Since this is not a category challenge, the last layer of InceptionV3, also known as the SoftMax layer, was eliminated from the design. All of the photographs were preprocessed to the exact same dimension, which was 299x299 pixels, before being directly entered into the version. As a result, the layer's output has the dimensions 8x8x2048. The characteristics were collected, saved as.npy files, and then transmitted through the encoder. The CNN InceptionV3 and Xception Chollet [2017] models were employed in the various experimental configurations of the aesthetic attention-based model. The

Network had been reduced to its base layer. The transformer model and the attention-based design both took 299x299 photos at that point. As a result, CNN received edited images. The collected picture features in the.npy files were then enhanced with the attention weight.

Year	Title	Methodology	Research Proposal	Algorithm
2022	Image Captioning Model using attention and Object features to mimic human	It involves extracting object features from the YOLO model and introducing them along with CNN convolutional features to a	Approximately the accuracy has been increased by around 45% compared to other models. 30,000 image	BLEU (Bi Lingual Evaluation Understudy),METEOR (Metric for Evaluation of Translation with Explicit Ordering), Consensus- based Image

	image understanding	simple deep learning model that uses the widespread Encoder-Decoder architecture with the attention mechanism.	datasets were used.	Description Evaluation (CIDEr),
--	---------------------	--	---------------------	---------------------------------

2021	Real-Time Image Captioning and Voice Synthesis using Neural Network	The image is captured by the System and processed by Concurrent Neural Network for the features or Objects in the extracted Images. Then the Recurrent Neural Network uses all the Extracted Features to transform into an Understandable Description of the image.	Output in Real-time. The Image captions are accurate compared to others Methods.	Computer Vision (CV), Natural Language Processing (NLP), and Language Models are used with Neural Networks.
------	---	---	--	---

2020	Image paragraph captioning using deep learning and NLP techniques	The input is given to CNN which recognizes the activities in the image and a vector representation is formed and given to the LSTM model where a word is generated and a caption is obtained. This method implies	According to the evaluation using BLEU Metrics, LSTM is identified as the best method with 80% efficiency	Gated Recurrent Unit (GRU) Method and LSTM is Used along with NLP.
------	---	---	---	--

		giving a coherent and detailed generation of paragraphs.		
2019	Smartphone-based Image Captioning for Visually and Hearing Impaired	First step is that visual attributes of images need to be extracted with richer content and as a second step, these attributes need to be sent to the NLP model to generate the most human-like captions.	The results show that the platform has great potential to be used for image captioning by visually and hearing-impaired people when integrated with Android application "Eye of Horus".	The images are trained using VGG16 deep learning architecture and Long Short-Term Memory (LSTM)model is used to generate a caption.

2018	Image Caption Generation Using Deep Learning Technique	For this task, we have used the Flickr 8k dataset consisting of 8000 images and five descriptions per image. In this work, we are using CNN as well as RNN. Pre-trained Convolutional Neural Network (CNN) is used for image classification task. This network acts as an image encoder. The last hidden layer is used as an input to the Recurrent Neural Network (RNN). This	The model has been trained for 50 epochs. The number of epochs used is more, which helps to lower the loss to 3.74. the Flickr 8k dataset for training model and running test on the 1000 test images available in dataset results in BLEU = 0.53356	Convolutional Neural networks (CNN), Long Short-Term Memory (LSTM), and Recurrent Neural Networks(RNN), BLEU (BiLingual Evaluation Understudy)
------	--	--	--	--

		network is a decoder that generates sentences.		
--	--	--	--	--

SYSTEM REQUIREMENTS

5.1 HARDWARE REQUIREMENTS:

- System** : Pentium IV 2.4 GHZ.
- Hard Disk** : 40 GB.
- Floppy Drive** : 1.44 Mb.
- Monitor** : 14" Colour Monitor.
- Mouse** : Optical Mouse.
- Ram** : 512 Mb.

SOFTWARE REQUIREMENTS:

- Operating system** : Windows 7 Ultimate.
- Coding Language** : Python.
- Front-End** : Python.
- Designing** : Html, CSS, JavaScript.
- Data Base** : MySQL.

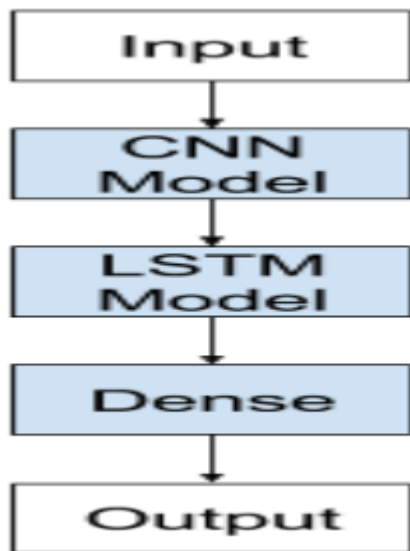
Architecture

To convert picture data into an output variable, convolutional neural networks are used. They have been so successful that they are currently the best method for making predictions of any kind that take photo data as an input. Recurrent neural networks, or RNNs for short, were developed to address issues like series prediction. Among these are challenges with one-to-many, many-to-one, and also many-to-many series forecasts. Because LSTM networks can store a larger series of words or phrases for forecast, they are among the most trustworthy RNNs. The CNNLSTM Version is one version that combines both CNN and LSTM. Consider the challenge of developing image captions as an illustration of how combining several kinds of networks into hybrid designs can produce one of the most fascinating and practical neural models. Given that we have an input image and a result sequence, which in this case is the caption for the input image, is it possible to characterise this as a one-to-many series prediction problem? Agreed, but how would a series prediction model, such as an LSTM or another one, undoubtedly identify the input image? Because they were not designed to communicate with such inputs, we cannot directly

enter the RGB picture tensor. Can we infer any key characteristics from the example image that would make it easier for us to utilise the Vanilla LSTM to represent images and other inputs with spatial organisation? Yes, that is exactly what we need to accomplish in order to utilise the LSTM architecture for our needs. We can remove elements from the image using the deep CNN method, and those features are then fed into the LSTM design to produce the inscription. The CNN LSTM version was primarily created to address issues with sequence prediction and spatial inputs, such as images or videos. This method combines LSTMs with function vectors to forecast series and Convolution Semantic Network (CNN) layers to remove functions from input data. In a word, CNN LSTMs are a class of deep models that sit at the intersection of natural language processing and computer vision. Both geographically and temporally, they are deep. These solutions have a lot of potential and are increasingly being employed for more challenging tasks like message classification and also video conversion. The version, which is listed below, is a CNN LSTM.

linguistic template

The language design is trained to anticipate caption series for a total of 100 dates. We are using Resnet feature vectors to build an LSTM-based architecture for photo captioning. Playing around and changing various setups as you choose is possible. An example of a result we obtained after training our network is shown.



2. Goal

The objective of image captioning is to create a summary in natural language from an input image. In this study, we'll develop a version of the Image Subtitle Generator that combines computer vision and also natural language processing to recognise image context and also describe it in a language like English

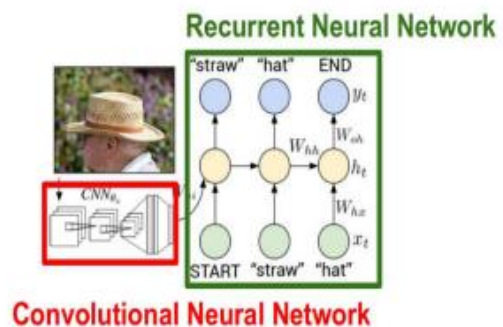
using the concepts of CNN and also LSTM.

Two categories can be used to categorise captioning.

1. An image-based model that extracts the characteristics of a picture.
2. A language-based model that creates a natural phrase using the attributes and also things our image-based model has really extracted.

We employ LSTM for our language-based version and CNN for our image-based model. The procedure of the Graphic Captioning Generator is illustrated in the images below. We often employ a convolutional neural network design for our image-based designs. See LSTM for language-based variations. The strategy's operation is shown in the diagram below.

Describing images



The properties from our input image are removed by a CNN that has already been trained. To have the same measurement as the LSTM network's input measurement, the attribute vector

is linearly altered. Using our attribute vector, this network is trained as a language design. We decide on our tag and target message before training our LSTM design. For example, if the summary reads, "An old gent is using a hat," then our label and target would undoubtedly be, "An old man is using a hat." This is done so that our design can distinguish between the first and last ideas in the series we've chosen. [An elderly man is sporting a cap. In the image dataset, 6000 images are used for training, 1000 for validation, and 1000 for testing. We will categorise this component directly into the complying with components for convenience of understanding.

Image Preprocessing.

'Building the lexicon of the picture' and 'Training the set'.

evaluate the model.

evaluating using specific images.

PREPROCESSING OF IMAGES:.

To recognise photos, we use a trained version called Visual Geometry Team (Resnet50). The picture attributes for attribute removal for Resnet50 in the Keras collection are currently 224 * 224 in size. The picture features are filled out only before the final layer of classification since this design is used to

anticipate a category for an image. Given that we had no interest in the photo category, the final layer was left off.

VOCABULARY DEVELOPMENT FOR THE IMAGE:.

We are unable to immediately incorporate the raw text into a deep learning or artificial intelligence architecture. First, the material needs to be cleaned up by being separated into words, with any errors with spelling and situational sensitivity addressed. We must encode each word into a fixed-sized vector, represent each word as a number, and assign each word in the lexicon a specific index worth since English words cannot be recognised by computer systems. Until then, the computer won't be able to understand the message or create photo subtitles. To get the necessary vocabulary dimension, we shall clean up the text in the following order. We describe the compliance with 5 functions as following in order to attain the aforementioned goals: The data is being organised. There are several procedures required, including creating an image-to-description thesaurus, deleting spelling, making all messages lowercase, and removing digits from words. Taking all of the summaries, identifying the distinctive words, and then developing a

lexicon constitute an additional stage. making a descriptions.txt file to store every inscription. The files are 8k on Flickr. A list of 6000 image names is provided in our dataset's train Images.txt file, and these names will undoubtedly be used to train the version. We begin by loading the attributes that were taken from the previously covered data.

This will create a dictionary using the CNN design with captions for each image in the list of images. Information Generator: We need to provide the design with input and also output during training in order to transform this into a supervised learning issue. A 4096-length function vector and a numerical representation of the caption are both contained in each of the 6000 photographs utilised to instruct our approach. Since the bulk of the data created for 6000 images would not fit in memory, we must use a generator technique that produces results in batches. The generator will produce both input and outcome series.

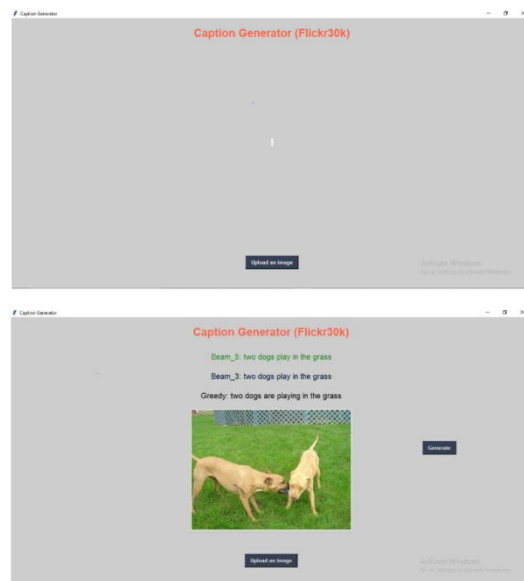
CNN-LSTM DESIGN: Using feature vectors obtained from the Resnet network, we are creating an LSTM-based model for photo captioning that anticipates word sequences, or subtitles. We will train the design using the 6000 training images by creating the input and outcome sequences in sets using the

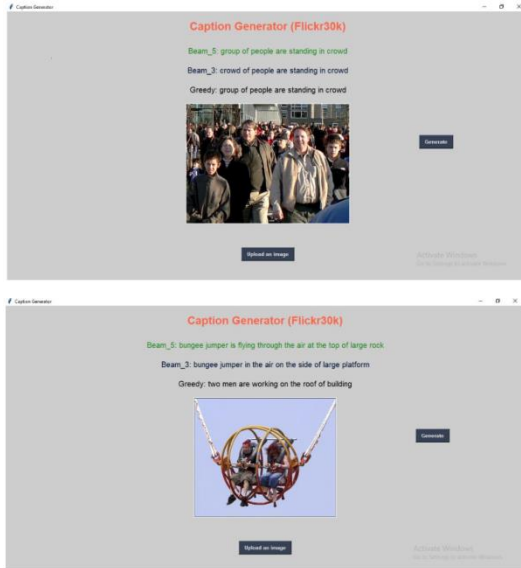
information manufacturing module and fitting them to the model. The model is trained over the course of 10 training epochs.

EVALUATION OF THE VERSION.

We can test the design against a set of randomly generated images after it has been trained. Since the projections include the largest possible index value, we will continue to use the same tokenizer. Pkl can retrieve words by using their index values. The produced captions for the images are fairly logical. The recognition of some photographs has room for improvement. Since it is data-dependent, the version cannot foresee words that are not in its vocabulary. For significantly better outcomes, we can create more accurate versions utilising an information set of 100,000 photographs.

RESULTS





CONCLUSION

In the past few years, there have been significant advancements in image captioning. The recent work that is fully based on deep understanding techniques has increased the accuracy of image captioning. Based on Accuracy and F1Score (0.95 and also 0.91) of our prior classification methods, Synthetic Neural Network once more outperforms them all. Although merging to the global minima cannot be guaranteed, it does find a respectable maximum based on some random weight initialization [3]. The efficiency was sufficient, and the effort resulted in excellent results that complement the methodology.

ADDITIONAL IMPROVE

Most work done to construct subtitles from photographs uses a framework that heavily relies on conventional semantic networks. Some

drawbacks of convolutional semantic networks are listed below:

CNN does not take into account the characteristics' orientation or spatial relationship. As an illustration, consider the following: The convolutional semantic network in the example will recognise the two photos as a face rather than figuring out the orientation and spatial connectivity of the faces in both images. As a remedy to the aforementioned issue, convolutional semantic networks use max merging, which results in information loss from the image because it frequently takes the maximum value within the matrix. CNN is significantly slower than maxpool. If a connection is too deep, it will take longer to educate it because there are many hidden layers. Choosing a quicker method can make the process go even more smoothly. Humans are incredibly diverse, and each individual has access to thousands of different caption options for a single image. Future research can focus on identifying the various personalities of human-annotated inscriptions and using this data into the evaluation of captioning.

REFERENCES

[1] Image Captioning: Transforming Objects into Words Lakshminarasimhan Srinivasan1, Dinesh Sreekanthan2, Amutha A.L3



[2] Image Captioning Based on Deep
Neural Networks Shuang Liu¹, Liang
Bai^{1, a}, Yanli Hu¹ and Haoran Wang¹

[3] Image Captioning with Keras,
harshlamba

<https://towardsdatascience.com/image-captioning-withkeras-teaching-computers-to-describe-picturesc88a46a311b8>

[4] IMAGE CAPTION GENERATOR
CNN-LSTM Architecture and Image
Captioning Arsh Chowdhry

<https://blog.clairvoyantsoft.com/image-caption-generator535b8e9a66ac>

[5] Image Captioning Based on Deep
Neural Networks Shuang Liu¹, Liang
Bai^{1,a}, Yanli Hu¹ and Haoran Wang