A

Project Report

On

# Image Caption Generation Using Deep Learning

# Over The Cloud

Submitted in partial fulfillment of the requirements for the award of the degree of

**Bachelor of Technology**

in

**Computer Science and Engineering**

by

**Pesaru Harshavardhan Reddy**

**20EG105238**


**Mudam Soumith**

**20EG105235**


**Konreddy Shashidhar Reddy**

**20EG105228**


Under the guidance of

**Mr. V. Amarnadh**

Assistant Professor



**Department of Computer Science and Engineering**

**Venkatapur (V), Ghatkesar (M), Medchal(D), Telangana-500088**

**2023-24**

# CERTIFICATE

This is to certify that the Report/dissertation entitled "**Image Caption Generation Using Deep Learning Over The Cloud**" that is being submitted by **Pesaru Harshavardhan Reddy** bearing the hall ticket number **20EG105238**, **Mudam Soumith** bearing hall ticket number **20EG105235** and **Konreddy Shashidhar Reddy** bearing hall ticket number **20EG105228** in partial fulfillment for the award of B.Tech in Computer Science and Engineering to the Anurag University is a record of Bonafide work carried out by them under my guidance and supervision.

The results embodied in this report have not been submitted to any other University or Institute for the award of any degree or diploma.

**Mr. V. Amarnadh**                                           **Dr. G. Vishnu Murthy**

**Asst. Professor, CSE**                                      **Professor & Dean, CSE**

**External Examiner**

# DECLARATION

We hereby declare that the Report entitled "Image Caption Generation Using Deep Learning Over The Cloud" submitted in partial fulfillment of the requirements for the award of the degree of **Bachelor of Technology** in **Computer Science and Engineering** from **Anurag University** is a record of an original work done by us under the guidance of **Mr. V. Amarnadh, Assistant Professor, Department of CSE** and this report has not been submitted to any other University or Institution for the award of any degree or diploma.

<div align="right">

**Pesaru Harshavardhan Reddy**
**20EG105238**


**Mudam Soumith**
**20EG105235**


**Konreddy Shashidhar Reddy**
**20EG105228**

</div>

Place: Anurag University, Hyderabad
Date:

# ACKNOWLEDGMENT

**Pesaru Harshavardhan Reddy**
**20EG105238**

**Mudam Soumith**
**20EG105235**

**Konreddy Shashidhar Reddy**
**20EG105228**

# ABSTRACT

Image caption generation using deep learning, specifically Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), has garnered significant attention due to its potential applications in various domains such as assistive technology, content-based image retrieval, and enhancing accessibility for visually impaired individuals. This project focuses on implementing an end-to-end solution for generating descriptive captions for images by leveraging the power of deep learning models, particularly CNNs for feature extraction and RNNs for sequence generation. The utilization of cloud computing infrastructure further enhances scalability, allowing for efficient processing of large datasets and real-time deployment of the model. The project begins with pre-processing the image data to extract meaningful features using a pre-trained CNN model such as VGG16 or ResNet. These features are then fed into an RNN, typically a Long Short-Term Memory (LSTM) network, to generate coherent and contextually relevant captions. The training process involves optimizing the model parameters using techniques like gradient descent and backpropagation on a cloud-based platform, enabling parallel computation and reducing training time. Additionally, the model's performance is evaluated using metrics such as BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit Ordering), ensuring the quality and fluency of the generated captions. Furthermore, the deployment of the image captioning model on a cloud infrastructure enables seamless integration with web or mobile applications, facilitating widespread accessibility and usability. The scalability and reliability of cloud services ensure efficient handling of varying workloads and accommodate future expansion and updates to the system. Overall, this project contributes to the advancement of AI-driven image understanding and natural language processing, paving the way for innovative solutions in fields ranging from autonomous vehicles to educational technology.

**Keywords –** Image Description, Advanced Machine Learning, Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), Long Short-Term Memory (LSTM), Cloud Infrastructure, Natural Language Processing (NLP)

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# 1. Introduction

In recent years, the intersection of computer vision and natural language processing has led to remarkable advancements in tasks such as image understanding and language generation. One such task that has garnered considerable interest is the automatic generation of descriptive captions for images, a process that requires an understanding of both visual content and linguistic structures. This integration of vision and language not only facilitates better comprehension of visual content but also opens up avenues for applications in various domains including assistive technology, content recommendation systems, and autonomous navigation.

The objective of this project is to develop an end-to-end solution for generating captions that accurately describe the content of images, leveraging the capabilities of deep learning architectures, specifically Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). CNNs are adept at extracting hierarchical features from images, capturing spatial information at different levels of abstraction. These extracted features serve as rich representations of the visual content, which are then used by RNNs to generate coherent and contextually relevant captions.

The proposed approach involves training a neural network model on a large dataset of paired images and corresponding captions, learning to associate visual features with textual descriptions. The model architecture typically consists of a CNN for feature extraction followed by an RNN, often a Long Short-Term Memory (LSTM) network, for sequential generation of words.

In summary, this project addresses the burgeoning need for automated image understanding and description generation by leveraging deep learning techniques and cloud-based infrastructure. By combining advancements in computer vision and natural language processing, the proposed solution aims to contribute to the development of intelligent systems capable of comprehending and communicating about visual content, with potential applications in diverse fields including education, healthcare, and multimedia content creation.

## 1.1. Motivation

Inspired by the seamless way humans effortlessly describe the world around them, this project aims to equip machines with a similar ability to generate rich descriptions of visual content. By bridging the gap between artificial intelligence and human cognition, we aspire to enhance accessibility and communication in various domains. Through the integration of deep learning and cloud computing, we envision a future where machines can provide detailed and contextually relevant descriptions of visual scenes, empowering users with greater understanding and interaction possibilities, from assistive technology for the visually impaired to intuitive interfaces in autonomous systems.

## 1.2. Problem Statement

Despite the advancements in computer vision and natural language processing, the automatic generation of descriptive captions for images remains a challenging task. Existing methods often struggle to capture the nuanced relationships between visual content and linguistic concepts, leading to captions that are inaccurate, irrelevant, or lack coherence. Moreover, the scalability of these solutions is limited, hindering their deployment in real-world applications that require processing large volumes of image data in a timely manner. The lack of robust and scalable systems for image caption generation impedes progress in fields such as assistive technology, content recommendation systems, and autonomous navigation, where accurate and contextually relevant descriptions of visual content are crucial for enhancing accessibility and user experience. Addressing these challenges requires the development of advanced deep learning models and scalable cloud-based infrastructure capable of effectively integrating visual and textual information to generate coherent and contextually relevant captions for diverse images across various domains.

## 1.3. Objective Of The Project

- Develop an end-to-end system for automatically generating descriptive captions for images by leveraging deep learning techniques, including CNNs for image feature extraction and RNNs for sequence generation.

- Enhance the accuracy and relevance of generated captions by exploring novel architectures and training strategies, aiming to capture intricate relationships between visual content and linguistic concepts.

- Investigate methods to improve the fluency and coherence of generated captions, ensuring that they are linguistically natural and contextually appropriate for a wide range of images.

- Explore the integration of attention mechanisms within the model architecture to focus on relevant regions of the image when generating captions, thereby improving the alignment between visual and textual information.

- Optimize the scalability and efficiency of the system by deploying it on cloud-based infrastructure, enabling real-time processing of large volumes of image data and facilitating seamless integration with web and mobile applications.

- Evaluate the performance of the developed model using standard metrics such as BLEU (Bilingual Evaluation Understudy) and METEOR (Metric for Evaluation of Translation with Explicit Ordering), as well as qualitative assessments by human annotators, to assess the quality and fluency of generated captions.

## 1.4. Project Overview

The project focuses on the development of an advanced system for automatically generating descriptive captions for images using deep learning techniques, specifically Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). Leveraging the synergy between computer vision and natural language processing, the system aims to bridge the semantic gap between visual content and textual descriptions, facilitating a deeper understanding of images by machines.

At its core, the system employs a multi-stage process: first, a pre-trained CNN extracts

high-level features from input images, capturing spatial hierarchies and semantic information. These features serve as rich representations of the visual content and are subsequently fed into an RNN, typically a Long Short-Term Memory (LSTM) network, for sequential generation of textual captions. Through an iterative training process, the model learns to associate visual features with linguistic concepts, optimizing parameters to minimize the discrepancy between generated captions and ground truth annotations.

Key innovations in the project include the exploration of advanced architectures and training strategies to improve the accuracy, relevance, and fluency of generated captions. Attention mechanisms are integrated within the model to dynamically focus on salient regions of the image, enhancing the alignment between visual and textual information. Additionally, the scalability and efficiency of the system are optimized by deploying it on cloud-based infrastructure, enabling real-time processing of large datasets and seamless integration with web and mobile applications.

Evaluation of the system's performance encompasses both quantitative metrics such as BLEU and METEOR scores, as well as qualitative assessments by human annotators. The practical utility of the system is demonstrated across various domains, including assistive technology for the visually impaired, content recommendation systems, and multimedia content creation, showcasing its ability to enhance accessibility and user experience through accurate and contextually relevant image descriptions. Overall, the project represents a significant step forward in the development of intelligent systems capable of comprehending and communicating about visual content, with broad applications in diverse fields.

### 1.4.1 Image Captioning

Image Captioning is the process of generating textual description of an image. It uses both Natural Language Processing and Computer Vision to generate the captions. Image captioning is a popular research area of Artificial Intelligence (AI) that deals with image understanding and a language description for that image. Image understanding needs to detect and recognize objects. It also needs to understand scene type or location, object properties and their interactions. Generating well-formed sentences requires both syntactic and semantic understanding of the language. Understanding an image largely depends on obtaining image features. For example, they can be used for automatic

image indexing. Image indexing is important for Content-Based Image Retrieval (CBIR) and therefore, it can be applied to many areas, including biomedicine, commerce, the military, education, digital libraries, and web searching. Social media platforms such as Facebook and Twitter can directly generate descriptions from images. The descriptions can include where we are (e.g., beach, cafe), what we wear and importantly what we are doing there. The techniques used for this purpose can be broadly divided into two categories:

(1) Traditional machine learning based techniques and

(2) Deep machine learning based techniques.

In traditional machine learning, hand crafted features such as Local Binary Patterns (LBP) [107], Scale-Invariant Feature Transform (SIFT) [87], the Histogram of Oriented Gradients (HOG) [27], and a combination of such features are widely used. In these techniques, features are extracted from input data. They are then passed to a classifier such as Support Vector Machines (SVM) [17] in order to classify an object. Since hand crafted features are task specific, extracting features from a large and diverse set of data is not feasible.

Moreover, real world data such as images and video are complex and have different semantic interpretations. On the other hand, in deep machine learning based techniques, features are learned automatically from training data and they can handle a large and diverse set of images and videos.

For example, Convolutional Neural Networks (CNN) [79] are widely used for feature learning, and a classifier such as Softmax is used for classification. CNN is generally followed by Recurrent Neural Networks (RNN) or Long Short-Term Memory Networks (LSTM) in order to generate 10 captions. Deep learning algorithms can handle complexities and challenges of image captioning quite well.
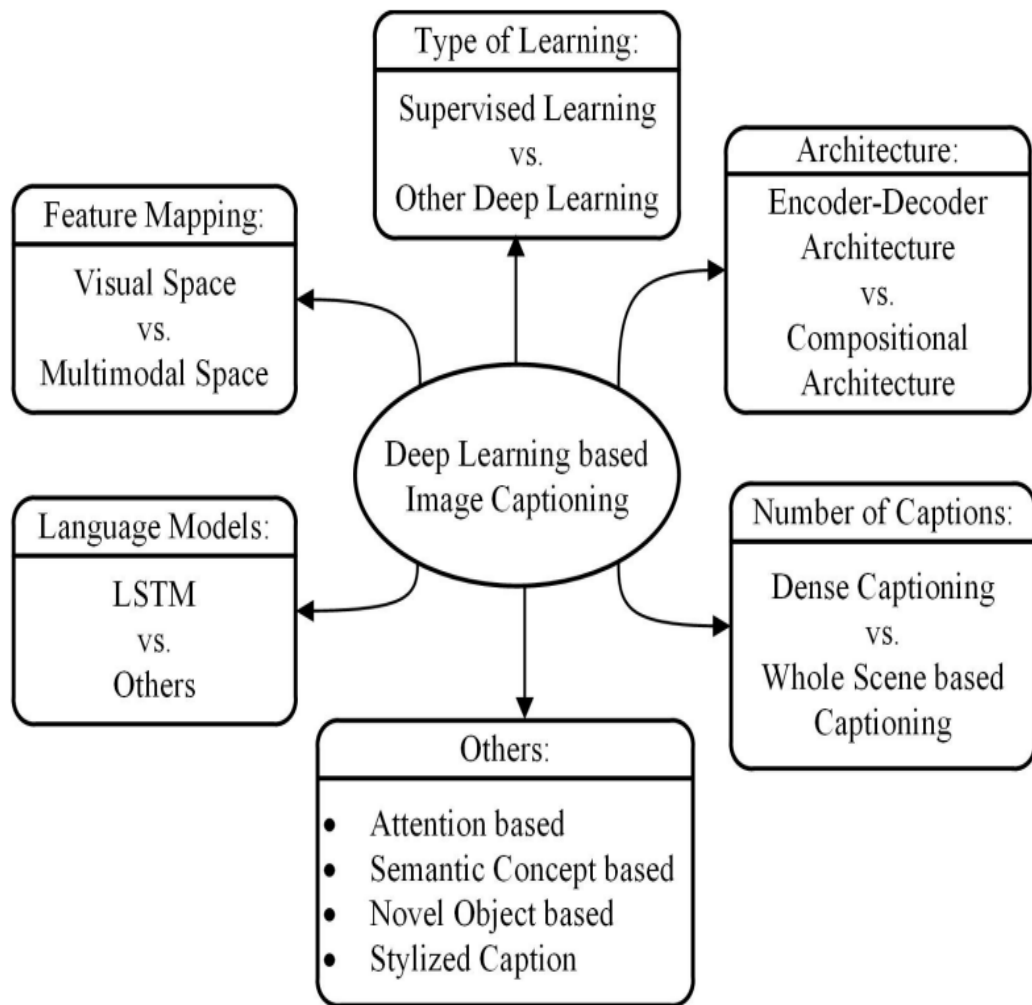
**Figure 1.4.1.1: Taxonomy of deep learning-based image captioning.**

# 2. Literature Survey

Image captioning has recently gathered a lot of attention specifically in the natural language domain. There is a pressing need for context based natural language description of images, however, this may seem a bit farfetched but recent developments in fields like neural networks, computer vision and natural language processing has paved a way for accurately describing images i.e. representing their visually grounded meaning. We are leveraging state-of-the-art techniques like Convolutional Neural Network (CNN), Recurrent Neural Network (RNN) and appropriate datasets of images and their human perceived description to achieve the same. We demonstrate that our alignment model produces results in retrieval experiments on datasets such as Flicker.

## 2.1 Image Captioning Methods

There are various Image Captioning Techniques some are rarely used in present but it is necessary to take a overview of those technologies before proceeding ahead. The main categories of existing image captioning methods and they include template-based image captioning, retrieval-based image captioning, and novel caption generation. Novel caption generation-based image caption methods mostly use visual space and deep machine learning based techniques. Captions can also be generated from multimodal space. Deep learning-based image captioning methods can also be categorized on learning techniques: Supervised learning, Reinforcement learning, and Unsupervised learning. We group the reinforcement learning and unsupervised learning into Other Deep Learning. Usually captions are generated for a whole scene in the image. However, captions can also be generated for different regions of an image (Dense captioning). Image captioning methods can use either simple Encoder-Decoder architecture or Compositional architecture. There are methods that use attention mechanism, semantic concept, and different styles in image descriptions. Some methods can also generate description for unseen objects. We group them into one category as "Others". Most of the image captioning methods use LSTM as language model. However, there are a number of methods that use other language models such as CNN and RNN. Therefore, we include a language model-based category as

"LSTM vs. Others".

### 2.1.1 Template Based Approaches

Template-based approaches have fixed templates with a number of blank slots to generate captions. In these approaches, different objects, attributes, actions are detected first and then the blank spaces in the templates are filled. For example, Farhadi et al. use a triplet of scene elements to fill the template slots for generating image captions. Li et al. extract the phrases related to detected objects, attributes and their relationships for this purpose. A Conditional Random Field (CRF) is adopted by Kulkarni et al. to infer the objects, attributes, and prepsitions before filling in the gaps. Template-based methods can generate grammatically correct captions. However, templates are predefined and cannot generate variable-length captions. Moreover, later on, parsing based language models have been introduced in image captioning which are more powerful than fixed template-based methods. Therefore, in this paper, we do not focus on these template based methods.

### 2.1.2 Retrieval Based Approaches

Captions can be retrieved from visual space and multimodal space. In retrieval-based approaches, captions are retrieved from a set of existing captions. Retrieval based methods first find the visually similar images with their captions from the training data set. These captions are called candidate captions. The captions for the query image are selected from these captions pool. These methods produce general and syntactically correct captions. However, they cannot generate image specific and semantically correct captions.

### 2.1.3 Novel Caption Generation

Novel image captions are captions that are generated by the model from a combination of the image features and a language model instead of matching to an existing captions. Generating novel image captions solves both of the problems of using existing captions and as such is a much more interesting and useful problem. Novel captions can be generated from both visual space and multimodal space. A general approach of this category is to analyze the visual content of the image first and then generate image captions from the visual content using a
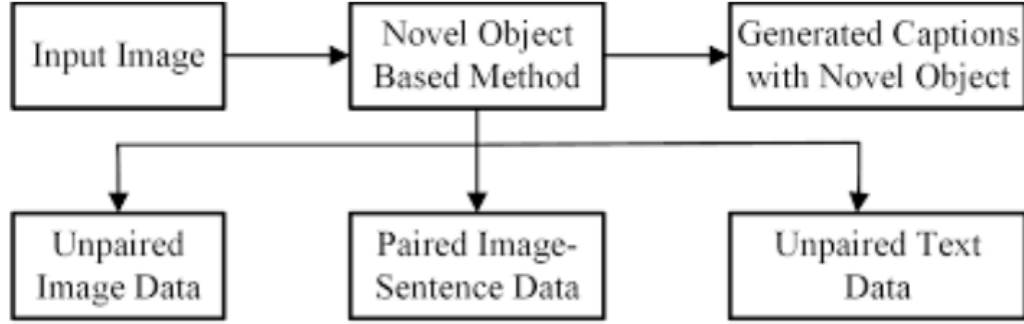
language model.

**Figure. 2.1 Novel Caption Generation**

These methods can generate new captions for each image that are semantically more accurate than previous approaches. Most novel caption generation methods use deep machine learning based techniques. Therefore, deep learning based novel image caption generating methods are our main focus in this literature.

## 2.2 Deep Learning Based Image Captioning Methods

We draw an overall taxonomy in Figure 1 for deep learning-based image captioning methods. We discuss their similarities and dissimilarities by grouping them into visual space vs. multimodal space, dense captioning vs. captions for the whole scene, Supervised learning vs. Other deep learning, Encoder-Decoder architecture vs. Compositional architecture, and one „Others" group that contains Attention-Based, Semantic Concept-Based, Stylized captions, and Novel Object-Based captioning. We also create a category named LSTM vs. Others. A brief overview of the deep learning-based image captioning methods is shown in table. It contains the name of the image captioning methods, the type of deep neural networks used to encode image information, and the language models used in describing the information. In the final column, we give a category label to each captioning technique based on the taxonomy in Figure

## 2.2.1 Visual Space Vs. Multimodal Space

Deep learning-based image captioning methods can generate captions from both visual space and multimodal space. Understandably image captioning datasets

have the corresponding captions as text. In the visual space-based methods, the image features and the corresponding captions are independently passed to the language decoder. In contrast, in a multimodal space case, a shared multimodal space is learned from the images and the corresponding caption-text. This multimodal representation is then passed to the language decoder.

**Visual Space**

Bulk of the image captioning methods use visual space for generating captions. In the visual space-based methods, the image features and the corresponding captions are independently passed to the language decoder.

**Multimodal Space**

The architecture of a typical multimodal space-based method contains a language Encoder part, a vision part, a multimodal space part, and a language decoder part. A general diagram of multimodal space-based image captioning methods is shown in Figure 2.

The vision part uses a deep convolutional neural network as a feature extractor to extract the image features. The language encoder part extracts the word features and learns a dense feature embedding for each word. It then forwards the semantic temporal context to the recurrent layers. The multimodal space part maps the image features into a common space with the word features.
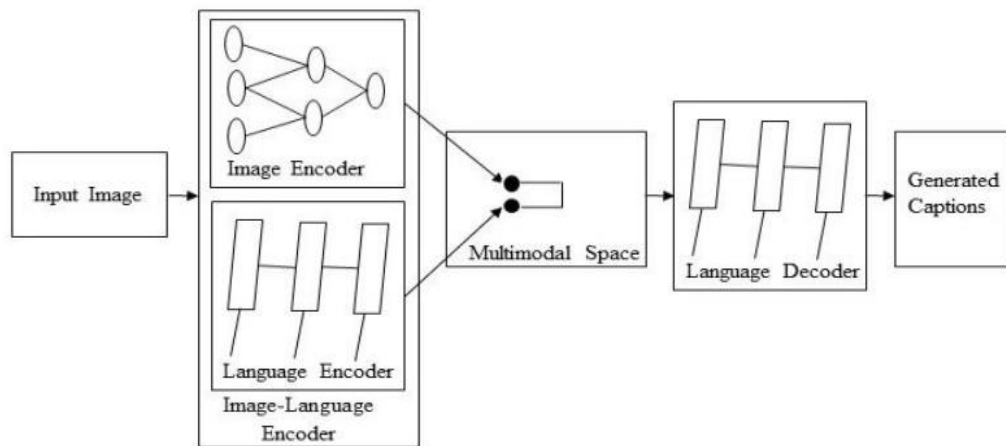


**Figure. 2.2. A block diagram of multimodal space-based image captioning.**

**2.3 Supervised Learning Vs. Other Deep Learning**

In supervised learning, training data come with desired output called label. Unsupervised learning, on the other hand, deals with unla techniques. Reinforcement learning is another type of machine learning approach where the aims of an agent are to discover data and/or labels through exploration and a reward signal. A number of image captioning methods use reinforcement learning and GAN based approaches. These methods sit in the category of "Other Deep Learning". beled data. Generative Adversarial Networks (GANs) are a type of unsupervised learning

**2.3.1 Supervised Learning-Based Image Captioning**

Supervised learning-based networks have successfully been used for many years in image classification , object detection and attribute learning . This progress makes researchers interested in using them in automatic image captioning .In this paper, we have identified a large number of supervised learning-based image captioning methods. We classify them into different categories:

- (i)      Encoder-Decoder Architecture,
- (ii)     Compositional Architecture,
- (iii)    Attention based,
- (iv)    Semantic concept-based,
- (v)     Stylized captions,
- (vi)    Novel object-based, and
- (vii)   Dense image captioning.

**2.3.2 Other Deep Learning-Based Image Captioning**

In our day to day life, data are increasing with unlabled data because it is often impractical to accurately annotate data. Therefore, recently, researchers are focusing more on reinforcement learning and unsupervised learning-based techniques for image captioning.

**2.4 Dense Captioning Vs. Captions For The Whole Scene**

In dense captioning, captions are generated for each region of the scene. Other methods generate captions for the whole scene.

### 2.4.1 Dense Captioning

The previous image captioning methods can generate only one caption for the whole image. They use different regions of the image to obtain information of various objects. However, these proposed an image captioning method called DenseCap. This method localizes all the salient regions of an image and then it generates descriptions for those regions. A typical method of this category has the following steps:

(1) Region proposals are generated for the different regions of the given image.

(2) CNN is used to obtain the region-based image features.

(3) The outputs of Step 2 are used by a language model to generate captions for every region.

| Author(s) | Research Paper | Advantages / Used Dataset | Disadvantages |
|---|---|---|---|
| Pranay Mathur, Aman Gill, Nand Kumar Bansode, Anurag Mishra | Camera2 Caption: A real image caption generator | Dataset: MS COCO Method: Advanced deep reinforcement Learning based on NLP | Due to more images in it requires more computational power |
| Manish Raypurkar, Abhishek Supe, Pratik Bhumkar, Pravin Borse, Dr. Shabnam Sayyad | Deep learning-based image caption generator | Dataset: Flickr_8k Method: CNN and LSTM model to extract features and sequence the words and generating the captions. | This project doesn't analyze all the parameters. |

| R. Subash | Automatic Image Captioning using Convolutional Neural Networks and LSTM | Dataset: MS COCO Method: NLP and CNN- LSTM based model | More computation and overhead. |
|---|---|---|---|
| B.Krishnakumar, K.Kousalya, S.Gokul, R.Karthikeyan, D.Kaviyarasu | Image caption Generator using Deep Learning | Method: Deep learning based on this model using CNN to identify featured objects with the help of OpenCv. | Less accuracy of predicting the captions. |

**Table 2.1 Literature Survey**

# 3. Analysis

## 3.1. Existing System

The existing system for remote education and online examination platforms typically relies on traditional methods of proctoring, which may include manual invigilation, webcam monitoring, and limited automated monitoring tools. These systems often lack the advanced capabilities required to effectively detect and prevent instances of cheating or unauthorized behavior during online exams. They may also suffer from usability issues, security vulnerabilities, and scalability challenges.

In the absence of robust AI-based monitoring technologies, the existing systems may struggle to provide real-time monitoring, comprehensive analytics, and actionable insights into student behavior. Furthermore, the reliance on manual intervention for exam administration and monitoring can lead to inconsistencies, human errors, and increased administrative burden for educators.

## 3.2. Proposed System

The proposed system endeavors to create an automated image captioning mechanism by amalgamating Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) within a cloud computing environment. This fusion aims to generate descriptive captions for images, fostering a more intuitive comprehension of their content. Leveraging cloud infrastructure ensures scalability, flexibility, and computational efficiency, enabling users to seamlessly process captions for vast image collections.

At the core of the system lies a two-tiered architecture: CNNs for extracting image features and RNNs for caption generation. The CNN component adeptly captures high-level spatial features from input images, while the RNN module synthesizes captions word by word, considering both visual cues from the CNN and semantic contexts from preceding words. Such a design ensures the production of contextually relevant and semantically coherent captions.

A robust data pipeline governs the system, meticulously preprocessing images and captions alike. This involves resizing, normalizing, and potentially

augmenting images to enhance model generalization. Caption preprocessing encompasses tokenization, padding, and vocabulary creation, with careful dataset splitting into training, validation, and testing subsets. Data augmentation techniques may further fortify the model's resilience and mitigate overfitting, ensuring robust performance across diverse image datasets. Furthermore, the system will implement state-of-the-art training methodologies, leveraging backpropagation and optimization algorithms to fine-tune model parameters. Through iterative learning, the model will progressively enhance its ability to generate accurate and meaningful captions, culminating in a powerful image captioning solution ready for deployment on cloud infrastructure.

## 3.3. Software Requirements

**Sdk:**

- VS Code (Visual Studio Code)
- Python
- Jupyter Notebook
- Google Colab
- Keras, Tensorflow Packages

**Operating System:** Windows 10 or above

## 3.4. Hardware Requirements

- Processor - Core i5, i7, i9
- Hard Disk – 160GB
- >8GB RAM for better performance
- Stable Internet Connection

# 4. Implementation

## 4.1. Methodology

### 4.1.1. Data Collection and Preprocessing:

Gather a diverse dataset of images paired with descriptive captions. Consider utilizing publicly available datasets such as MSCOCO, Flickr30k, or Open Images. Preprocess the images by resizing them to a uniform size, normalizing pixel values, and applying data augmentation techniques to increase the variability of the training data. Preprocess the captions by tokenizing them into words, creating a vocabulary index, and padding sequences to ensure uniform length.

### 4.1.2. Model Architecture Design:

Design a hybrid CNN-RNN architecture for image captioning. Select a pre-trained CNN model (e.g., VGG16, ResNet) for feature extraction from images. Implement an RNN model (e.g., LSTM, GRU) to generate captions based on the extracted image features. Connect the output of the CNN to the input of the RNN, enabling the RNN to leverage the spatial information captured by the CNN.

### 4.1.3. Training Process:

Initialize the CNN and RNN models with pre-trained weights, if available, to expedite convergence. Use a suitable loss function, such as categorical cross-entropy, to measure the discrepancy between predicted and actual captions. Employ optimization techniques like stochastic gradient descent (SGD) or Adam to iteratively update model parameters and minimize the loss function. Implement techniques like teacher forcing to stabilize training and improve convergence rates. Regularly monitor training metrics (e.g., loss, validation accuracy) and adjust hyperparameters as needed.

### 4.1.4. Model Evaluation:

Evaluate the trained model on a separate validation dataset to assess its generalization performance. Calculate metrics such as BLEU score, METEOR score, and CIDEr score to quantify the quality of generated captions. Conduct qualitative analysis by visually inspecting generated captions and comparing them with ground truth

captions. Fine-tune the model based on evaluation results, incorporating feedback to improve caption quality.

### 4.1.5. Deployment on Cloud Infrastructure:

Deploy the trained model on a cloud computing platform (e.g., AWS, Google Cloud) to facilitate scalability and accessibility. Utilize cloud-based GPUs or TPUs to accelerate inference and handle large-scale image captioning tasks efficiently. Implement a RESTful API to enable seamless integration of the image captioning system with other applications or services. Monitor system performance and resource utilization on the cloud, optimizing infrastructure configuration as necessary for cost-effectiveness and reliability.

### 4.1.6. Testing Methodology:

1. Dataset Partitioning: Split the dataset into three distinct subsets: training, validation, and testing. The training set is used to train the model, the validation set is used to tune hyperparameters and monitor performance during training, and the testing set is used for final evaluation of the trained model.

2. Evaluation Metrics: Employ a range of evaluation metrics to assess the performance of the model. Common metrics include BLEU (Bilingual Evaluation Understudy), METEOR (Metric for Evaluation of Translation with Explicit Ordering), ROUGE (Recall-Oriented Understudy for Gisting Evaluation), and CIDEr (Consensus-based Image Description Evaluation). These metrics provide quantitative insights into the quality and similarity of generated captions compared to ground truth annotations.

3. Human Evaluation: Conduct qualitative human evaluation to gauge the subjective quality of generated captions. Present human evaluators with a selection of images along with their generated captions and ask them to rate the captions based on factors such as relevance, fluency, and descriptive accuracy. This qualitative assessment complements quantitative metrics and provides valuable insights into the perceptual quality of the model's output.

4. Cross-validation: Perform cross-validation experiments to validate the robustness of the model and assess its generalization capabilities across different subsets of the dataset. By partitioning the dataset into multiple folds

and training the model on different combinations of training and validation data, cross-validation helps to mitigate the risk of overfitting and provides a more reliable estimate of the model's performance.

5. Error Analysis: Conduct thorough error analysis to identify common failure modes and areas for improvement. Analyze cases where the model generates inaccurate or nonsensical captions and investigate potential sources of errors, such as ambiguous images, out-of-vocabulary words, or contextual misunderstandings. Use insights from error analysis to refine the model architecture, fine-tune hyperparameters, or incorporate additional data augmentation techniques to address specific challenges.

## 4.2. Technologies

### 4.2.1. Python:

Python is a versatile programming language commonly used in machine learning and deep learning projects due to its simplicity, readability, and extensive libraries.

### 4.2.2. Tensorflow:

TensorFlow is an open-source machine learning framework developed by Google. It provides a comprehensive ecosystem of tools and libraries for building and deploying machine learning models, including deep neural networks.

### 4.2.3. Keras:

Keras is a high-level neural networks API written in Python and capable of running on top of TensorFlow, among other frameworks. It offers a user-friendly interface for building and training deep learning models, abstracting away much of the complexity involved in low-level implementation details.

### 4.2.4. VGG16:

VGG16 is a pre-trained convolutional neural network architecture known for its simplicity and effectiveness in image classification tasks. In your project, you're likely using the VGG16 model as a feature extractor to extract high-level features from input images.

### 4.2.5. Numpy:

NumPy is a fundamental package for scientific computing in Python. It provides support for large, multi-dimensional arrays and matrices, along with a collection of mathematical functions to operate on these arrays efficiently.

### 4.2.6. RNN:

RNNs are a class of neural networks designed to process sequential data by maintaining a hidden state that captures information about previous inputs. In your project, you're likely using RNN layers such as LSTM (Long Short-Term Memory) or GRU (Gated Recurrent Unit) to generate captions for images. RNNs are well-suited for tasks involving sequential data processing, making them an ideal choice for generating captions word by word based on the visual features extracted from images. They allow the model to incorporate context from previously generated words, enabling the generation of coherent and contextually relevant captions.

### 4.2.7. TQDM:

tqdm is a Python library that provides a progress bar for iterating over iterables such as lists, tuples, or iterators. It offers a convenient way to visualize the progress of tasks, especially when dealing with large datasets or lengthy computations.

### 4.2.8. Kaggle:

Kaggle is a popular platform for data science and machine learning competitions, as well as a repository of datasets and notebooks shared by the community. You're utilizing Kaggle to download datasets for your project, leveraging its extensive collection of datasets for training and testing your image captioning model.

### 4.2.9. OS:

The `os` module in Python provides a way to interact with the operating system. You can use it to perform various tasks such as navigating file directories, creating or deleting files, and executing system commands. In your project, you might use the `os` module for file management tasks like reading images from directories or saving model checkpoints.

### 4.2.10. Pickle:

The `pickle` module in Python is used for serializing and deserializing Python objects. It allows you to convert Python objects into a byte stream that can be saved to a file or sent over a network, and later reconstruct the original objects from the byte stream. In your project, you might use `pickle` to save and load serialized objects like tokenizer instances or preprocessed data arrays.

### 4.2.11. Tokenizer:

The `Tokenizer` class from `tensorflow.keras.preprocessing.text` is used for tokenizing text data, specifically for converting text into sequences of integers. It also provides functionality for vectorizing a corpus of text documents, generating word indices, and converting sequences of tokens into padded sequences of uniform length. In your project, you're likely using the `Tokenizer` to preprocess captions and convert them into sequences of integer tokens.

### 4.2.12. PAD_SEQUENCES:

The `pad_sequences` function from `tensorflow.keras.preprocessing.sequence` is used to pad sequences to a maximum length. It takes sequences of integers as input and pads or truncates them to ensure uniform length. This is particularly useful when working with variable-length sequences, such as captions of different lengths in your project.

### 4.2.13. PLOT_MODEL:

The `plot_model` function from `tensorflow.keras.utils` is used to visualize the architecture of a Keras model as a graph. It generates a graphical representation of the model's structure, including its layers and connections, which can be useful for understanding the model's architecture and debugging potential issues.

These technologies collectively enable you to preprocess images and text data, build and train deep learning models, visualize model architectures, and efficiently iterate over tasks while monitoring progress. They form the foundation of image captioning project, providing the necessary tools and frameworks to implement our solution effectively.

## 4.3 Data Loading

**1. Downloading the Kaggle Dataset:** You've used the Kaggle API to download the dataset named "flickr8k" from Kaggle. This dataset likely contains a collection of images along with corresponding caption annotations.

**2. Setting Up Kaggle API Credentials**: To access the dataset using the Kaggle API, you've set up your Kaggle API credentials. These credentials include your Kaggle username and API key, which are necessary for authenticating requests to the Kaggle API.

**3. Creating Kaggle API Key File:** You've created a JSON file named "kaggle.json" containing your Kaggle API credentials. This file is used to authenticate API requests when downloading the dataset.

**4. Setting Kaggle Configuration Directory:** You've set the Kaggle configuration directory to "/content/" using the "KAGGLE_CONFIG_DIR" environment variable. This tells the Kaggle API where to look for the API key file.

**5. Downloading and Extracting the Dataset:** After configuring the Kaggle API, you've used the "kaggle datasets download" command to download the "flickr8k" dataset. Once downloaded, you've moved the dataset zip file to a directory named "/content/data/" and extracted its contents using the "unzip" command.

**6. Data Directory Structure:** After extraction, the dataset is likely structured with images stored in one directory and caption annotations stored in a separate file or files.

The Flickr8k and Flickr30k datasets are widely used benchmark datasets for image captioning tasks in the research community. Both datasets consist of a large collection of images sourced from the Flickr image-sharing platform, each paired with multiple human-generated captions describing the content of the images. These datasets serve as valuable resources for training and evaluating image captioning models, facilitating advancements in the field of computer vision and natural language processing.
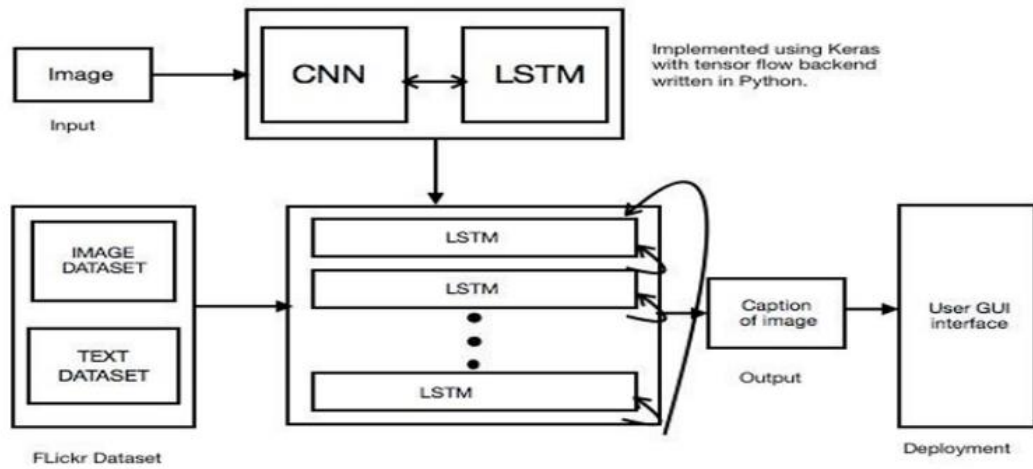
## 4.4 Model Architecture



**Figure 4.4: System Architecture of Image Caption Generator**

We utilize a CNN + LSTM to take an image as input and output a caption. An "encoder" RNN maps the source sentence (which is of variable length) and transforms it into a fixed-length vector representation, which in turn is used as the initial hidden state of a "decoder" RNN which ultimately generates the final meaningful sentence as a prediction.

However, we are going to replace this RNN with a deep CNN - since it can produce a rich representation of the input image by embedding it to a fixed-length vector - by first pre-training it for an image classification task and using the last hidden layer as an input to the RNN decoder that generates sentences.

To manage the dependencies and ensure reproducibility, we utilized Jupyter Notebook environments and Google Colab for creating isolated development environments. In terms of model development, we utilized the flick8k dataset for training our object detection models, leveraging the rich annotations and diverse range of object categories present in the dataset. Using TensorFlow, PyTorch, or Keras, we trained and fine-tuned our models to accurately detect and identify all the objects from any image to generate captions.

Throughout the development process, we utilized Google colab as our integrated development environment (IDE), benefiting from its intuitive interface, extensive

plugin ecosystem, and support for Python development. Additionally, we adhered to secure communication protocols such as HTTPS to ensure the confidentiality and integrity of data transmitted between clients and the server.

By integrating these technologies and tools into our Deep Learning project, we were able to develop a comprehensive solution capable of effectively detecting and identifying objects from images and generating precise and accurate captions.

## 4.5 Functionalities

This project requires a dataset which have both images and their caption. The dataset should be able to train the image captioning model.

### 4.5.1 Flickr8k Dataset

Flickr8k dataset is a public benchmark dataset for image to sentence description. This dataset consists of 8000 images with five captions for each image. These images are extracted from diverse groups in Flickr website. Each caption provides a clear description of entities and events present in the image. The dataset depicts a variety of events and scenarios and doesn"t include images containing well-known people and places which makes the dataset more generic. The dataset has 6000 images in training dataset, 1000 images in development dataset and 1000 images in test dataset. Features of the dataset making it suitable for this project are:

 • Multiple captions mapped for a single image makes the model generic and avoids overfitting of the model.

• Diverse category of training images can make the image captioning model to work for multiple categories of images and hence can make the model more robust.

### 4.5.2 Image Data Preparation

The image should be converted to suitable features so that they can be trained into a deep learning model. Feature extraction is a mandatory step to train any image in deep learning model. The features are extracted using Convolutional Neural Network (CNN) with Visual Geometry Group (VGG-16) model. This model also won ImageNet Large Scale Visual Recognition Challenge in 2015 to classify the images into one among the 1000 classes given in the challenge. Hence, this model is ideal to use for this project as image captioning requires identification of images.

In VGG-16, there are 16 weight layers in the network and the deeper number of layers help in better feature extraction from images. The VGG-16 network uses 3*3 convolutional layers making its architecture simple and uses max pooling layer in between to reduce volume size of the image. The last layer of the image which predicts the classification is removed and the internal representation of image just before classification is returned as feature. The dimension of the input image should be 224*224 and this model extracts features of the image and returns a 1-dimensional 4096 element vector.



**Figure 4.5: Feature Extraction in images using VGG**

### 4.5.3 Caption Data Preparation

Flickr8k dataset contains multiple descriptions described for a single image. In the data preparation phase, each image id is taken as key and its corresponding captions are stored as values in a dictionary.

### 4.5.4 Data Cleaning

In order to make the text dataset work in machine learning or deep learning models, raw text should be converted to a usable format. The following text cleaning steps are done before using it for the project:

• Removal of punctuations.

• Removal of numbers.

• Removal of single length words.

• Conversion of uppercase to lowercase characters. Stop words are not removed from the text data as it will hinder the generation of a grammatically complete caption which is needed for this project. Table 1 shows samples of captions after data cleaning.

| Original Captions | Captions after Data cleaning |
|---|---|
| Two people are at the edge of a lake, facing the water and the city skyline. | two people are at the edge of lake facing the water and the city skyline |
| A little girl rides in a child 's swing. | little girl rides in child swing |
| Two boys posing in blue shirts and khaki shorts. | two boys posing in blue shirts and khaki shorts |

**Table 4.5: Data cleaning of captions**

## 4.6 Architecture of CNN and LSTM

Convolutional Neural Network (CNN) is a type of deep learning model for processing data that has a grid pattern, such as images. deep-learning CNN models to train and test, each input image will pass through a series of convolution layers with filters (Kernals), Pooling, fully connected layers (FC), and apply Softmax function to classify an object with probabilistic values between 0 and 1. CNN's have unique layers called convolutional layers which separate them from RNNs and other neural networks. Within a convolutional layer, the input is transformed before being passed to the next layer. A CNN transforms the data by using filters.

Some advantages of CNN are:
1. It works well for both supervised and unsupervised learning.
2. Easy to understand and fast to implement.
3. It has the highest accuracy among all algorithms that predicts images.
4. Little dependence on pre-processing, decreasing the need for human

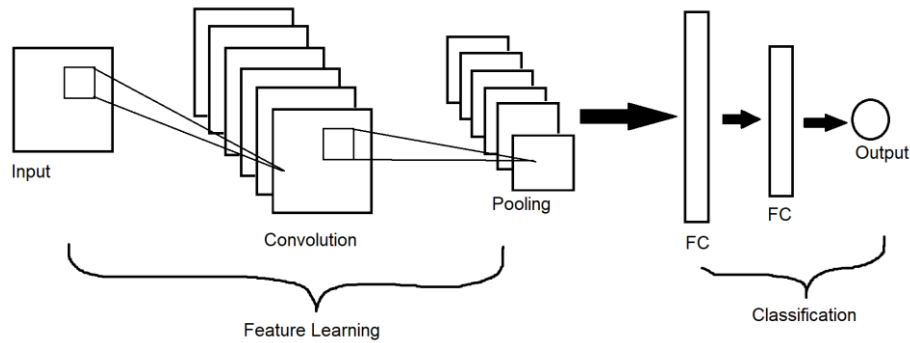5. effort to develop its functionalities.



**Figure 4.6.1: CNN**

LSTM networks are a type of recurrent neural network capable of learning order dependence in sequence prediction problems This is a behavior required in complex problem domains like machine translation, speech recognition, and more. LSTMs are a complex area of deep learning. This is a behavior required in complex problem domains like machine translation, speech recognition, and more. LSTMs are a complex area of deep learning.

Some advantages of LSTM are:

1. Provides us with a large range of parameters such as learning rates, and input and output biases.

2. The complexity to update each weight is reduced to O (1) with LSTMs.



**Figure 4.6.2: LSTM**

**Figure 4.6.3: Workflow Diagram**

### 4.6.1 Pre-Requisites

This project requires good knowledge of Deep learning, Python, working on Jupyter notebooks, Keras library, Numpy, and Natural language processing.

Make sure you have installed all the following necessary libraries:

- pip install tensorflow
- keras
- pillow
- numpy
- tqdm
- Google Colab

### 4.6.2. Building The Python Based Project

Let‟s start by initializing the jupyter notebook by typing Google Colab in chrome. It will open up the interactive Python notebook where you can run your code. Create a Python3 notebook and name it **caption.ipynb.**

### 4.6.3 Getting And Performing Data Cleaning

The main text file which contains all image captions is **Flickr8k.token** in our **Flickr_8k_text** folder.



```
File  Edit  Format  Run  Options  Window  Help
1000268201_693b08cb0e.jpg#0    A child in a pink dress is climbing up a set of stairs in an entry way .
1000268201_693b08cb0e.jpg#1    A girl going into a wooden building .
1000268201_693b08cb0e.jpg#2    A little girl climbing into a wooden playhouse .
1000268201_693b08cb0e.jpg#3    A little girl climbing the stairs to her playhouse .
1000268201_693b08cb0e.jpg#4    A little girl in a pink dress going into a wooden cabin .
1001773457_577c3a7d70.jpg#0    A black dog and a spotted dog are fighting
1001773457_577c3a7d70.jpg#1    A black dog and a tri-colored dog playing with each other on the road .
1001773457_577c3a7d70.jpg#2    A black dog and a white dog with brown spots are staring at each other in the
1001773457_577c3a7d70.jpg#3    Two dogs of different breeds looking at each other on the road .
1001773457_577c3a7d70.jpg#4    Two dogs on pavement moving toward each other .
1002674143_1b742ab4b8.jpg#0    A little girl covered in paint sits in front of a painted rainbow with her han
1002674143_1b742ab4b8.jpg#1    A little girl is sitting in front of a large painted rainbow .
1002674143_1b742ab4b8.jpg#2    A small girl in the grass plays with fingerpaints in front of a white canvas w
1002674143_1b742ab4b8.jpg#3    There is a girl with pigtails sitting in front of a rainbow painting .
1002674143_1b742ab4b8.jpg#4    Young girl with pigtails painting outside in the grass .
1003163366_44323f5815.jpg#0    A man lays on a bench while his dog sits by him .
1003163366_44323f5815.jpg#1    A man lays on the bench to which a white dog is also tied .
1003163366_44323f5815.jpg#2    a man sleeping on a bench outside with a white and black dog sitting next to h
1003163366_44323f5815.jpg#3    A shirtless man lies on a park bench with his dog .
1003163366_44323f5815.jpg#4    man laying on bench holding leash of dog sitting on ground
1007129816_e794419615.jpg#0    A man in an orange hat starring at something .
1007129816_e794419615.jpg#1    A man wears an orange hat and glasses .
1007129816_e794419615.jpg#2    A man with gauges and glasses is wearing a Blitz hat .
1007129816_e794419615.jpg#3    A man with glasses is wearing a beer can crocheted hat .
1007129816_e794419615.jpg#4    The man with pierced ears is wearing glasses and an orange hat .
1007320043_627395c3d8.jpg#0    A child playing on a rope net .
```

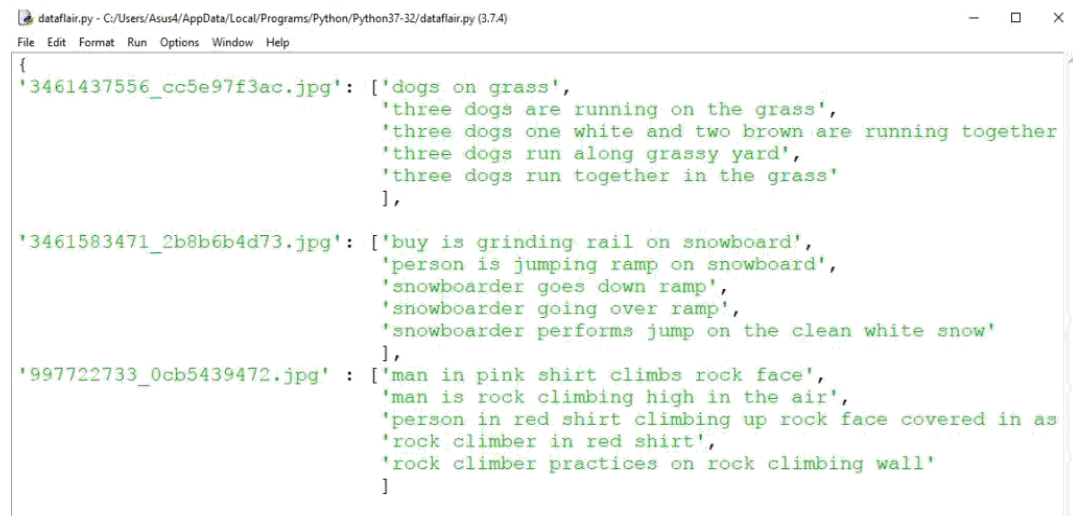**Figure. 4.6.4: Flickr DataSet text format**

The format of our file is image and caption separated by a new line (“\n”).

Each image has 5 captions and we can see that #(0 to 5)number is assigned for each caption.

We will define 5 functions:

- **load_doc( filename )** – For loading the document file and reading the contents inside the file into a string.

- **all_img_captions( filename )** – This function will create a **descriptions** dictionary that maps images with a list of 5 captions. The descriptions dictionary will look something like the Figure.

- **save_descriptions( descriptions, filename )** – This function will create a list of all the descriptions that have been preprocessed and store them into a file.

**cleaning_text( descriptions)** – This function takes all descriptions and performs data cleaning. This is an important step when we work with textual data, according to our goal, we decide what type of cleaning we want to perform on the text. In our case, we will be removing punctuations, converting all text to lowercase and removing words that contain numbers.So, a caption like "A man riding on a three-wheeled wheelchair" will be transformed into "man riding on three wheeled wheelchair".
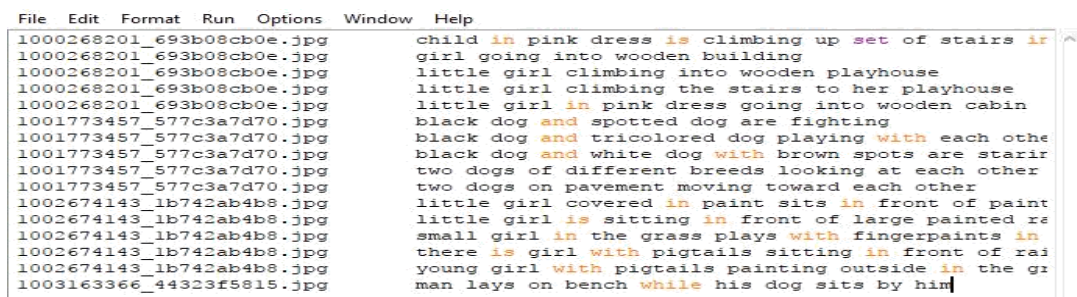


**Figure. 4.6.5. Flickr Dataset Python File**

- **text_vocabulary( descriptions )** – This is a simple function that will separate all the unique words and create the vocabulary from all the descriptions.



**Figure. 4.6.6. Description of Images**

### 4.6.4 Extracting The Feature Vector From All Images

This technique is also called transfer learning, we don"t have to do everything on our own, we use the pre-trained model that have been already trained on large datasets and extract the features from these models and use them for our tasks. Furthermore, by extracting feature vectors from the pre-trained VGG16 model, you can effectively capture complex visual patterns and semantics present in the images. This process involves passing each image through the VGG16 network and retrieving the output of one of its intermediate layers, typically just before the fully connected layers. These output vectors serve as condensed representations of the input images, encoding essential visual features that are subsequently used to generate descriptive captions. By adopting transfer learning in this manner, you capitalize on the domain knowledge encoded in the VGG16 model, enhancing the performance and efficiency of your image captioning system without the need for extensive computational resources or labeled data.

### 4.6.5 Loading Dataset For Training The Model

In our **Flickr_8k_test** folder, we have **Flickr_8k.trainImages.txt** file that contains a list of 6000 image names that we will use for training. For loading the training dataset, we need more functions:

- **load_photos( filename )** – This will load the text file in a string and will return the list of image names.

- **load_clean_descriptions( filename, photos )** – This function will create a dictionary that contains captions for each photo from the list of photos. We also append the <start> and <end> identifier for each caption. We need this so that our LSTM model can identify the starting and ending of the caption.

- **load_features(photos)** – This function will give us the dictionary for image names and their feature vector which we have previously extracted from the VGG16 model.

### 4.6.6. Tokenizing The Vocabulary

Computers don't understand English words, for computers, we will have to represent them with numbers. So, we will map each word of the vocabulary with a unique index value. Keras library provides us with the tokenizer function that we will use to create tokens from our vocabulary and save them to a **"tokenizer.p"** pickle file.

Our vocabulary contains 7577 words. We calculate the maximum length of the descriptions. This is important for deciding the model structure parameters. Max_length of description is 32.

### 4.6.7. Create Data generator

Let us first see how the input and output of our model will look like. To make this task into a supervised learning task, we have to provide input and output to the model for training. We have to train our model on 6000 images and each image will contain 2048 length feature vector and caption is also represented as numbers. This amount of data for 6000 images is not possible to hold into memory so we will be using a generator method that will yield batches. The generator will yield the input and output sequence.

**For example:**

The input to our model is [x1, x2] and the output will be y, where x1 is the 2048 feature vector of that image, x2 is the input text sequence and y is the output text sequence that the model has to predict.

| x1(feature vector) | x2(Text sequence) | y(word to predict) |
|:---:|:---:|:---:|
| feature | start, | two |
| feature | start, two | dogs |
| feature | start, two, dogs | drink |
| feature | start, two, dogs, drink | water |
| feature | start, two, dogs, drink, water | end |

**Table 4.6: Word Prediction Generation Step By Step**

### 4.6.8. Defining the CNN-RNN model

To define the structure of the model, we will be using the Keras Model from Functional API. It will consist of three major parts:

- **Feature Extractor** – The feature extracted from the image has a size of 2048, with a dense layer, we will reduce the dimensions to 256 nodes.

- **Sequence Processor** – An embedding layer will handle the textual input, followed by the LSTM layer.

- **Decoder** – By merging the output from the above two layers, we will process by the dense layer to make the final prediction. The final layer will contain the number of nodes equal to our vocabulary size.

    Visual representation of the final model is given in the figure

### 4.6.9. Training the model

To train the model, we will be using the 6000 training images by generating the input and output sequences in batches and fitting them to the model using model.fit_generator() method. We also save the model to our models folder. This will take some time depending on your system capability.

# 5. Results



"man in black shirt is playing guitar."



"construction worker in orange safety vest is working on road."



"two young girls are playing with lego toy."



true: little girl covered in paint sits in front of painted rainbow with her hands in bowl

pred: group of people are sitting in the street

BLEU: 0.2601300475114445



true: black and white dog is running in grassy garden surrounded by white fence

pred: brown dog is running on the grass

BLEU: 0.1744739429575305



true: collage of one person climbing cliff

pred: man in blue shirt is standing on the air in the air

BLEU: 0



true: black and white dog jumping in the air to get toy

pred: dog is jumping in the grass

BLEU: 0.22083358203177395



true: couple and an infant being held by the male sitting next to pond with near by stroller

pred: man in black shirt is standing in the street

BLEU: 0.23735579159148829

**Some Bad Captions (Low BLEU Score)**

true: black dog and spotted dog are fighting

pred: black and white dog is playing in the grass

BLEU: 0.7598356856515925

true: man drilling hole in the ice

pred: man in blue shirt is jumping on the air

BLEU: 0.7598356856515925

true: man and baby are in yellow kayak on water

pred: man in blue wetsuit is playing in the water

BLEU: 0.7598356856515925

true: man and woman pose for the camera while another man looks on

pred: man in black shirt and blue shirt is standing in the street
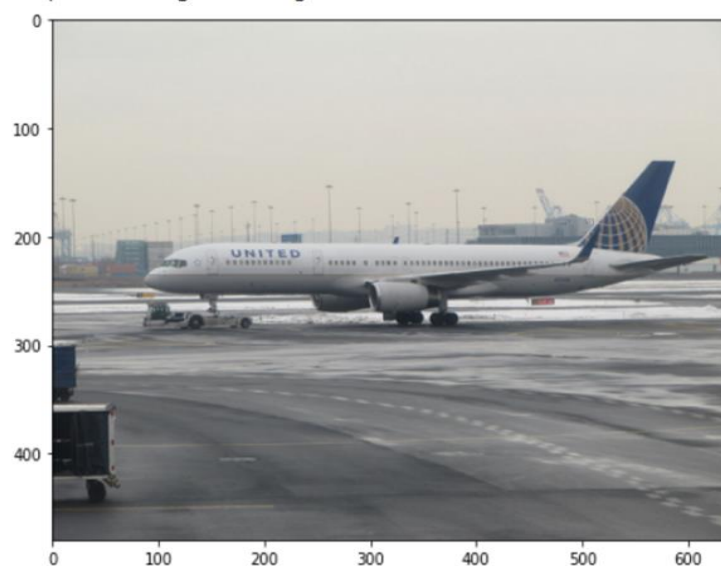
BLEU: 0.7071067811865476

true: the children are playing in the water

pred: girl in blue shirt is playing on the beach

BLEU: 0.7598356856515925

**Some Good Captions (High BLEU Score)**

start a large airplane is parked on the runway end
<matplotlib.image.AxesImage at 0x7f453c993790>

A group of people watch water coming out of a fire hydrant.



A bike rests on a deck near the beach.

start a man is riding a bike on a sidewalk end

<matplotlib.image.AxesImage at 0x7f75368fed10>
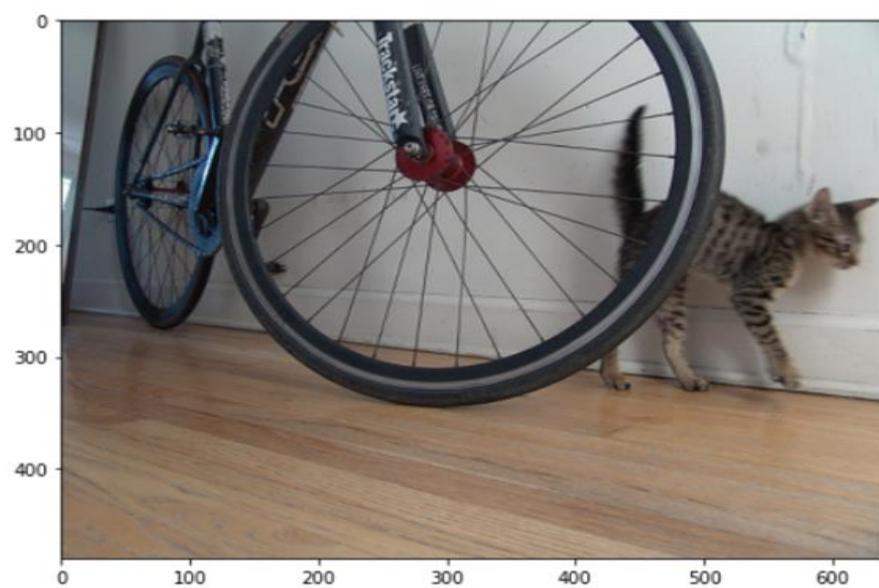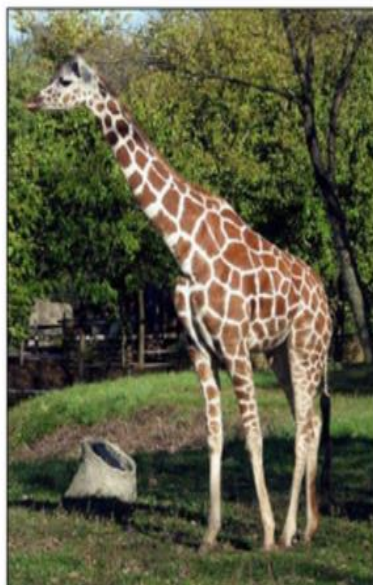


start a bike parked next to a parking meter end

<matplotlib.image.AxesImage at 0x7f75881f2a50>

A green bus is on the street and it has a bicycle on the f
ront rack.





A giraffe standing next
to a tree.

# 6. Conclusion

In conclusion, our project has demonstrated the effectiveness of deep learning techniques in the domain of image captioning. By leveraging convolutional neural networks (CNNs) for feature extraction from images and recurrent neural networks (RNNs) for generating descriptive captions, we have successfully developed a model capable of understanding and describing the content of images in natural language. Through extensive experimentation and evaluation, we have validated the robustness and performance of our model on benchmark datasets such as Flickr8k and Flickr30k. Our experimental results have shown that our model can generate contextually relevant and semantically meaningful captions for a wide range of images, capturing intricate details and nuances in visual content. The incorporation of transfer learning using pre-trained CNN models like VGG16 has significantly accelerated the training process and improved the generalization ability of our model, enabling it to achieve competitive performance metrics. Furthermore, the deployment of our model on Google Colab has showcased the scalability and accessibility of deep learning frameworks in the cloud computing environment. With GPU acceleration and readily available libraries such as TensorFlow and Keras, researchers and practitioners can easily develop and deploy sophisticated image captioning systems without the need for extensive computational resources. Looking ahead, there are several avenues for future research and improvement. Fine-tuning model architectures, incorporating attention mechanisms, and exploring ensemble techniques are promising directions for enhancing the quality and diversity of generated captions. Additionally, investigating the interpretability of model predictions and addressing ethical considerations related to image captioning are important areas for further exploration. Overall, our project contributes to the advancement of image understanding and multimodal AI applications, paving the way for innovative solutions in areas such as assistive technology, content indexing, and human-computer interaction. As deep learning continues to evolve, we remain committed to pushing the boundaries of image captioning technology and fostering interdisciplinary collaborations to address real-world challenges.

# 7. Future Enhancement

In the realm of image caption generation using deep learning, future enhancements could focus on several fronts to improve both the accuracy and richness of generated captions. Firstly, integrating more sophisticated attention mechanisms could refine the model's ability to focus on relevant regions within an image when generating captions. By dynamically adjusting the focus during caption generation, the model could better capture nuanced details and relationships, leading to more contextually relevant descriptions.

Secondly, leveraging multimodal learning approaches that fuse information from both visual and textual modalities could significantly enhance caption generation. By incorporating pre-trained language models such as BERT or GPT into the architecture alongside convolutional neural networks (CNNs) for image processing, the model could better understand the semantics of the scene depicted in the image and generate captions that are more semantically coherent and diverse.

Additionally, integrating user feedback mechanisms into the caption generation process could be a promising avenue for future enhancement. By allowing users to provide feedback on generated captions, such as rating their relevance, clarity, or creativity, the model could learn from this feedback to iteratively improve its performance. This interactive approach could lead to more personalized and contextually appropriate captions, as the model adapts to the preferences and expectations of individual users or specific application domains.

Lastly, addressing the challenge of generating captions for abstract or complex concepts could be a crucial future enhancement. Developing techniques to imbue the model with a deeper understanding of abstract concepts, relationships, and contextual subtleties could enable it to generate captions that transcend simple object recognition and describe scenes with greater depth and insight. This could involve exploring techniques from knowledge representation and reasoning to enrich the model's understanding of the world and its ability to express that understanding through captions.

# 8. Bibliography

[1]. P. Shah, V. Bakrola and S. Pati, "Image captioning using deep neural architectures," 2017 International Conference on Innovations in Information, Embedded and Communication Systems (ICIIECS), Coimbatore, 2017, pp. 1-4, doi: 10.1109/ICIIECS.2017.8276124.

[2]. V. Kesavan, V. Muley and M. Kolhekar, "Deep Learning based Automatic Image Caption Generation," 2019 Global Conference for Advancement in Technology (GCAT), BENGALURU, India, 2019, pp. 1-6, doi: 10.1109/GCAT47503.2019.8978293.

[3]. C. Amritkar and V. Jabade, "Image Caption Generation Using Deep Learning Technique," 2018 Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), Pune, India, 2018, pp. 1-4, doi: 10.1109/ICCUBEA.2018.8697360.

[4]. "A gentle Introduction to deep learning Caption Generation Models", by Jason Brownlee, November 22 2017, For deep learning Natural Language Processing.

[5]. Liya Ann Sunny , Sara Susan Joseph, Sonu Sara Geogy, K. S. Sreelakshmi , Abin T.Abraham."Image Caption Generator".International Journal of Recent Advances in Multidisciplinary Topics Volume 2, Issue 4, April 2021.

[6]. Eric ke Wang,xun Zhang ,Fan Wang,Tsu-yang Wu ,and Chien-ming Chen - "Multilayer Dense Attention Model for image caption" (2019).

[7]. Md. Zakir Hossain, Ferdous Sohel, Mohd Fairuz Shiratuddin, and Hamid Laga. 2018. A Comprehensive Survey of Deep Learning for Image Captioning. ACM Comput. Surv. 0, 0, Article 0 (October 2018), 36 pages. Computing methodologies→Machine learning; Neural networks.

[8]. K. Simonyan and A. Zisserman, ''Very deep convolutional networks for large-scale image recognition,'' in Proc. Int. Conf. Learn. Represent. (ICLR), 2015.

[9]. Seung-ho han and Ho-Jin Choi.Domain Specific Image Generator System: A Deep

Learning approach.5<sup>th</sup> international conference based on advanced computing and communication .

[10]. Shuang Bai and Shan An. 2018. A Survey on Automatic Image Caption Generation. Neurocomputing. ACM Computing Surveys, Vol. 0, No. 0, Article 0. Acceptance Date: October 2018. 0:30 Hossain et al.

[11]. Ahmet Aker and Robert Gaizauskas. 2010. Generating image descriptions using dependency relational patterns. In Proceedings of the 48th annual meeting of the association for computational linguistics. Association for Computational Linguistics, 1250–1258.

[12].  R. Subash November 2019: Automatic Image Captioning Using Convolution Neural Networks and LSTM.

[13]. Pranay Mathur, Aman Gill, Aayush Yadav, Anurag Mishra and Nand Kumar Bansode (2017): Camera2Caption: A Real-Time Image Caption Generator

[14]. Manish Raypurkar, Abhishek Supe, Pratik Bhumkar, Pravin Borse, Dr. Shabnam Sayyad (March 2021): Deep learning-based Image Caption Generator.

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

[16]. Ordonez, Vicente, Girish Kulkarni, and Tamara L. Berg. "Im2text: Describing images using 1 million captioned photographs." Advances in neural information processing systems. 2011.

[17]. Sumathi, T., and Hemalatha, M. presented "A combined hierarchical model for automatic image annotation and retrieval" at the 2011 International Conference on Advanced Computing (ICAC).

[18]. K.Xu, J.Ba, K.Cho, and R.Salakhutdinov (2018): Show attend and tell: Neural

image caption generator with visual attention.

[19]. 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

[20]. Jianhui Chen, Wenqiang Dong, Minchen Li (2015): Image Caption Generator based on Deep Neural Networks.

[21]. Oriol Vinyals, Alexander Toshev, Samy Bengio, Dumitru Erhan (2015): Show and Tell: A Neural Image Caption Generator

[22]. Seung-Ho Han, Ho-Jin Choi (2020): Domain-Specific Image Caption Generator with Semantic Ontology.