

# **ImageCaptioningUsingDeepLearning**

A major project report submitted in partial fulfillment of the  
requirement for the award of a degree of

**Bachelor of Technology**

in

**Computer Science & Engineering / Information Technology**

*Submitted by*

**Kushagra Shukla (201303)**

**Supriti Sharma (201487)**

*Under the guidance & supervision of*

**Dr. Kushal Kanwar**



**Department of Computer Science & Engineering and  
Information Technology**

**Jaypee University of Information Technology,**

**Waknaghat, Solan - 173234 (India)**

# CERTIFICATE

This is to certify that the work which is being presented in the project report titled **“Image Captioning Using Deep Learning”** in partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering/Information Technology** submitted to the Department of Computer Science and Engineering, Jaypee University of Information Technology, Wazirpur is an authentic record of work carried out by **Kushagra Shukla(201303)** and **Supriti Sharma(201487)** during the period from August 2023 to May 2024 under the supervision of **Dr. Kushal Kanwar**, (Assistant Professor (SG), Department of Computer Science and Engineering, Jaypee University of Information Technology).

(Student Signature with Date) (Student Signature with Date) Student Name:  
Kushagra Shukla Student Name: Supriti Sharma Roll No.: 201303 Roll No.: 201487

The above statement made is correct to the best of our knowledge.

(Supervisor Signature with Date)  
Supervisor Name: Dr. Kushal Kanwar  
Designation: Assistant Professor (SG)  
Department: CSE/IT  
Dated: 13 May 2024

# DECLARATION

We at this moment declare that the work presented in this report entitled '**Image Captioning Using Deep Learning**' in partial fulfilment of the requirements for the award of the degree of **Bachelor of Technology in Computer Science & Engineering / Information Technology** submitted in the Department of Computer Science & Engineering and Information Technology, Jaypee University of Information Technology, Waknaghat is an authentic record of our work carried out over a period from August 2023 to May 2024 under the supervision of **Dr Kushal Kanwar** (Assistant Professor (SG), Department of Computer Science & Engineering and Information Technology).

The matter embodied in the report has not been submitted for any other degree or diploma award.

(Student Signature with Date) (Student Signature with Date) Student Name:  
Kushagra Shukla Student Name: Supriti Sharma Roll No.: 201303 Roll No.:  
201487

This is to certify that the above statement made by the candidate is true to the best of my knowledge.

(Supervisor Signature with Date)  
Supervisor Name: Dr. Kushal Kanwar  
Designation: Assistant Professor (SG)  
Department: CSE/IT  
Dated: 13 May 2024

# ACKNOWLEDGEMENT

Firstly, I express my heartiest thanks and gratefulness to almighty God for His divine blessing to make it possible to complete the project work successfully.

We are grateful and wish out profound indebtedness to Supervisor **Dr Kushal Kanwar, Assistant Professor(SG)**, Department of CSE Jaypee University of Information Technology, Waknaghat. Deep Knowledge & keen interest of our supervisor in the field of “**Research Area**” to carry out this project. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior drafts and correcting them at all stages have made it possible to complete this project.

We would also generously welcome each one of those individuals who have helped me straightforwardly or in a roundabout way in making this project a win. In this unique situation, We want to thank the various staff individuals, both educating and non-instructing, who have developed their convenient help and facilitated our undertaking.

Finally, We must acknowledge with due respect the constant support and patience of our parents.

Kushagra Shukla  
(201303)

Supriti Sharma  
(201487)

# **TABLE OF CONTENT**

## **Content**

Certificate.....	ii
Declaration by Candidate .....	iii
Acknowledgement .....	iv
Abstract.....	x

## **Chapter 1: INTRODUCTION**

General Introduction.....	1
Problem Statement.....	2
Objectives .....	3
Significance and Motivation.....	3

## **Chapter 2: LITERATURE SURVEY**

Overviewofrelevantliterature..... 4

KeyGapsinliterature..... 9

### **Chapter3:SYSTEMDEVELOPMENT**

RequirementsAndAnalysis ..... 10

ProjectDesignandArchitecture..... 11

DataPreparation ..... 12

Implementation ..... 13

KeyChallenges ..... 34

### **Chapter4:Testing**

TestingStrategy..... 36

TestcasesAnd Outcomes ..... 37

### **Chapter-5ResultsandEvaluation**

Results..... 39

**Chapter-6ConclusionsandFutureScope**

Conclusion .....43

FutureScope .....44

References..... 46

## ListOfFigures

<b>FigNo.</b>	<b>Title</b>	<b>PageNo.</b>
1.	OverviewofproposedVAQ Model	5
2.	ProposedModelSystem	6
3.	ArchitectureAndDesignOfTheProject	11
4.	CNN	14
5.	LSTM	16
6.	CNN-LSTMModel	17
7.	WorkflowDiagram	18
8.	ImportingLibraries	19
9.	ImportingData	20
10.	InitialCodeGeneration	21
11.	FirstOutcome	21
12.	AssigningTheValuesToElement	22
13.	Imagetocaptiongeneration	23
14.	Assigningthelength	24
15.	TrainingDataset	24
16.	CNNEncoder	25
17.	Transformerencodinglayer	26
18.	Defininglength,range,anddimension	26
19.	Transformerdecodinglayer	27
20.	DecodingLayer	28



21.	ICModelwithlossandaccuracy	29
22.	Computinglossandaccuracy	30
23.	EncoderDecoder	30
24.	TrainingDataset	30
25.	PredictedOutput	31
26.	PredictedOutput	32
27.	DataSourceMapping	32
28.	OutputLayer	33
29.	AccuracyAndLossFunction Hugging Face Interface	34
30.	HuggingFaceInterfaceUploadingtheimage	38
31.	UploadingtheimageHuggingFaceOutput1(A group of people)	39
32.	HuggingFaceOutput1(Agroupofpeople)	40
33.	HuggingFaceOutput2(Dogrunninginfield)	41
34.	HuggingFaceOutput3(Womenridinghorse)	42

## ListOfTables

TableNo.	Description	PageNo.
1.	LiteratureSurvey	8

# Abstract

Thanks to the advancement of the deep learning, the field of the concept artificial intelligence has seen a significant increase in interest in research concerning the combination of natural language processing and computer vision. An English description of a photograph's context is automatically generated. When an image has a caption, the computer is trained to decipher the image's visual data using one or more sentences.

The meaningful description generating process of high-level picture semantics requiresthecapacitytoexaminethestate, characteristics,andrelationshipsbetween these objects.Inthis research, weaim to discover objectsand notify individuals via text messagesbyapplying CNN -LSTMarchitectural modelsonthecaptioningofa graphical image. In order to accurately identify the objects, the input image is first converted to grayscale.Computer vision and natural language processing were combined then. We made use of the COCO Dataset 2017.

In order to help blind people realizetheir full potential and track their intelligence, the suggestedapproach for blind people is meantto be broadened toinclude people with visual loss to speech messages. In this study, we create An iterative CNN-LSTM framework that outperforms human baselines by adhering to several key concepts in picture captioning and its conventional procedures.



# CHAPTER-1 INTRODUCTION

## General Introduction

We are surrounded by images in the news, on social media, and in our surroundings.

Photos can only be recognized by humans. Without the accompanying descriptions, people are able to identify photographs; however, machines need to first be trained with images. Input vectors are used by the encoder-decoder Image Caption's architecture. Generator models produce appropriate and legitimate captions. The domains of Natural language processing and computer vision are related by this paradigm. It is the responsibility to identify and assess the image's context before providing a full description in an organic language such as English.

Long short-term memory (LSTMs) and convolutional neural networks (CNNs) are the foundational models of our methodology. In the derived application, CNN is employed as an encoder with the goal to extract features from an image or snapshot, while LSTM is employed as a converter to arrange the expressions and produce captions. Image captioning can be useful for many purposes. For example, it can help the blind via text-to-speech by providing true-time information about the scene via a camera feed. It can also increase social media leisure by organizing captions for both spoken remarks and photos in social media feeds.

One step towards learning the language is to help kids identify chemicals. Authentic photo exploration and indexing can be accelerated and improved by adding captions to every image available online. Image captioning is used in many fields, including biology, business, the internet, and self-driving motor vehicles, in which the surroundings are explained and CCTV cameras, which can trigger alarms if they detect malicious activity. This research article's primary goal is to provide readers with a fundamental understanding of deep learning techniques.

## **ProblemStatement**

In the modern world, information is valued, and some people have severe difficulty seeing images. We investigate this further, taking blindness into account as a significant factor, and produce a sentence by letting users upload or scan a picture.

Advantage:-

- Recommendations in Editing Applications
- Assistance for visually impaired
- Social Media posts
- Self-Driving cars
- Robotics
- Easy to implement and connect to new data sources

Disadvantages:-

- Never offer an end-to-end developed general model to solve this problem,
- Neither they offer intuitive feature observations on the objects or actions in the image.

### **1.2.2 Problem Analysis**

As image captioning is a very vast and underdeveloped field of deep learning there are many problems such as incorrect captions that are being generated for a specific picture that is uploaded in the model for the captions. This can be reduced by using correct techniques and proper guidance if a person has knowledge in the field of deep learning, and neural networks this problem can be easily solved and the captions that is to be generated will be correct.

**Scope:** Our project extends and is being used in any large-scale business industry and also small-scale business industry.

## Objectives

1. The project's goal is to develop a Deep Learning (DL) method for contextualizing photographs in basic English sentences.
2. The requirement to work with LSTM and CNN rather than RNN
3. To work on the webpage for the caption generation.

## Significance and Motivation of the Project Work

Making captions for images is an essential task in the domains of computer vision and natural language processing images. It is an impressive advancement in artificial intelligence when a machine replicates human ability to describe images. Capturing the connections between the objects in the picture and how to describe them in a language that is natural (like English) is the main challenge of this task. Computer systems have historically generated text descriptions for images by utilizing pre-defined templates. Neural networks are frequently employed in state-of-the-art models to generate captions by using an image as input and forecasting the next lexical unit in the output sentence. Nevertheless, this method falls short of offering the necessary diversity to produce lexically rich text descriptions. With neural networks' increased efficiency, this weakness has been suppressed.

We can develop a product for the blind that will allow them to navigate the roadways without the assistance of others. We could accomplish this by turning the scene to text and then the text to voice. Automatic driving is one among the most important challenges and if we can properly caption the scene around the car, it can enhance the self-driving system. Automatic captioning can help Google Image Search, which will become nearly as excellent as Google Search, as each image is first converted into a caption, and then searches are conducted based on the caption. CCTV cameras are now ubiquitous, but if we combine viewing the globe with the generation of appropriate captions, we will be able to raise alarms as soon as criminal conduct is detected. This will undoubtedly aid in the reduction of crime and accidents.

The model based on LSTM and CNN uses deep convolutional neural network (CNN) to create a dense feature vector from the images. The dense vectors are also called embeddings. For the image caption model, this embedding acts as a dense representation of the image which can be used as the initial state of the LSTM. The CNN network can be trained directly on the images in our dataset. An LSTM is a recurrent neural network architecture that is mainly used for problems having temporal dependences. It is useful for capturing information about previous states to better inform the current prediction through its memory cell state.

# Chapter-2

## Literature Survey

### Overview Of Relevant Literature

The most crucial stage in the the process for creating software is the literature review. Determining The economy, time factor, and strength of the company is essential before developing the tool. Selecting the operating system and language that can be utilized for the tool is the next step. development once these requirements are met. The programmers require a great deal of outside assistance once they begin developing the tool. Books, websites, and senior programmers can all provide this support. The aforementioned factors are taken into account before developing the system that is suggested.

Most of the project's development industry takes into account and thoroughly investigates every need needed to develop the project. A literature review is the most important crucial step in the process of developing software for any project. It is vital to ascertain and evaluate the company's strength, economy, workforce, resource requirements, and time factor prior to developing the tools alongside their designs.

To improve and personalize the user experience on its products, photos use image classification. Numerous typical computer vision issues, such as intra-class variation, occlusion, deformation, size variation, perspective variation, and lighting, are represented by the picture classification problem. Techniques that are effective for classifying pictures are also likely to be effective for other crucial computer vision tasks, such as segmentation, detection, and localization.

Image captioning is a great illustration of this. Given an image, the image captioning challenge is to generate a sentence description of the image. The picture captioning problem is comparable to the image classification problem in that it expects more detail and has a bigger universe of possibilities. Image classification is used as a black box system in modern picture captioning systems, therefore greater image classification leads to better captioned.

The image captioning problem is interesting in and of itself because it combines natural language processing and computer vision, two important areas of artificial intelligence. An image captioning system shows that it can comprehend natural language and image semantics. Once all of these demands are met, the next step is to determine the software specifications in the corresponding system, including what kind of operating system the project would need and what all the software is needed to move onto the next step, which is developing the tools and the related operations and thoroughly surveyed.

Image classification is a crucial step in the object recognition and picture analysis process that leads to the construction of an image sentence. A statement could be

the end result of the image categorization process. Many methods for captioning images have been presented thus far. The ideal method for captioning images has been the subject of

numerous studies. Since the outcomes and accuracy vary depending on a number of factors, it is challenging to select one method as the best of all of them.

Both new image captioning techniques developed in the last few decades and traditional approaches have been continuously modified to achieve the most accurate results. Numerous complex tasks like image classification, semantic segmentation, and object detection, have demonstrated exceptional performance with multilayer convolutional neural networks in recent times. In particular, two-stage approaches are often employed for semantic segmentation. In this manner, Convolutional neural networks are learned to deliver high-quality local pixel-wise data for the subsequent step, which is typically a more comprehensive graphical analysis model.

The Vanishing Gradient problem will be resolved by employing Long Short-Term Memory (LSTM), a subset of RNNs. The Vanishing Gradients problem is the primary objective of LSTM. The ability of LSTM to retain data values over extended periods of time makes it unique in solving the vanishing gradient problem. Numerous advanced tasks, including object detection, image classification, and, more recently, semantic segmentation, have recently been proven to obtain outstanding results using convolutional neural networks with many layers. A two-stage technique is frequently used, especially for semantic segmentation. In this way, convolutional networks are trained to give excellent local pixel-wise data for the second stage, which is often a more global graphical analysis model. The results showed that using a combination of LSTM produced better results than applying RNN. Unlike conventional image recognition algorithms, CNNs use multilayer convolutional to carry out feature engineering and integrate these features internally. It also makes use of SoftMax, the fully connected (FC) and pooling layers.

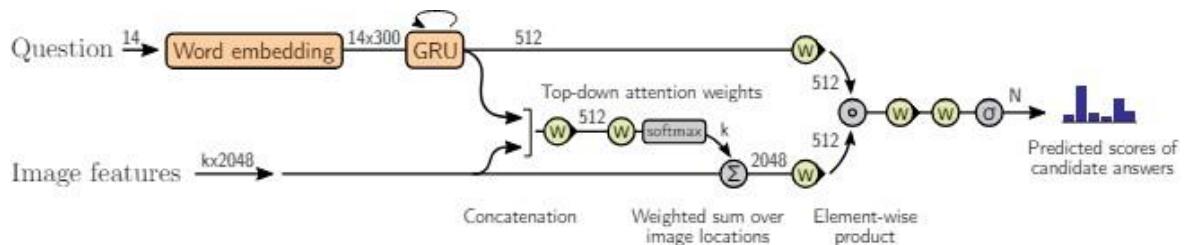


Fig1 Overview of the proposed VAQ model



## V. Proposed System

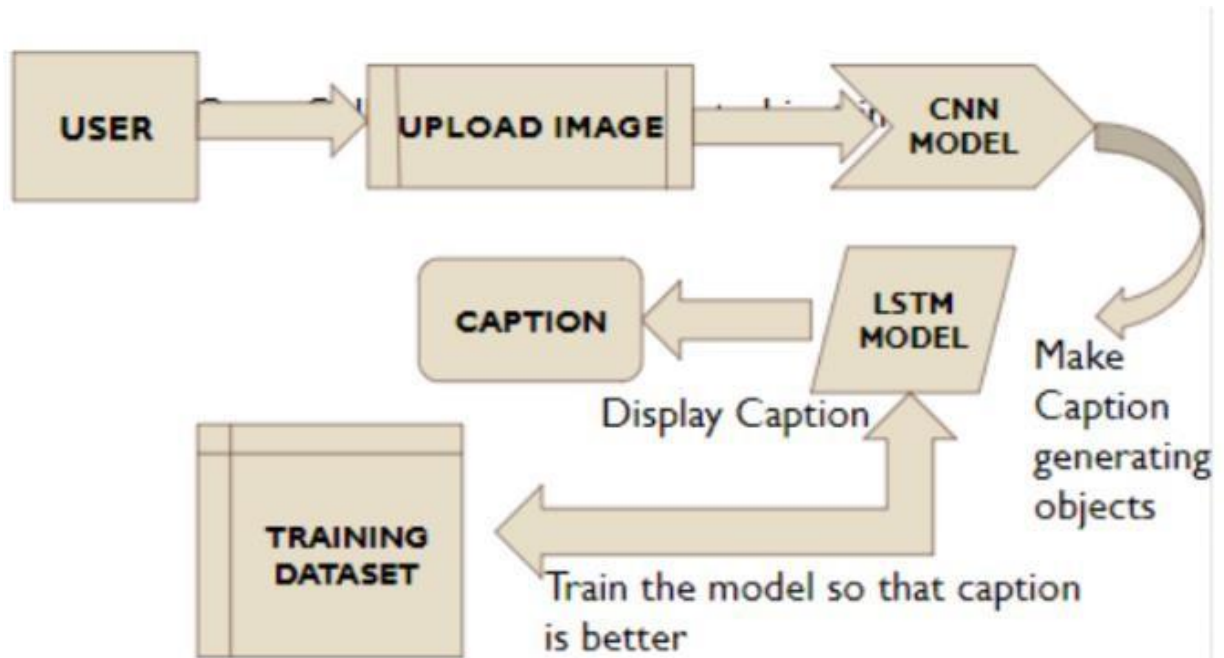


Fig2. Proposed Model System

<b>SNo.</b>	<b>AuthorName</b>	<b>Dataset and Methodology</b>	<b>Results</b>
1.	RajendranSubash	Dataset: MS COCOMethod: NLP and CNNLSTM basedmodel	Using CNN LSTM andNLPtechniques the modelforimage captioning is generated
2.	Seung-HoHan, Ho-JinChoi	Semantic Ontology	Using semantic ontology the model forimagecaptioning is generated
3.	PranayMathur, AmanGill, Nand Kumar Bansode, Anurag Mishra	Dataset: MS COCO Method: Advanced deep reinforcement learning based on NLPand computerVision	The model proposed generates the real time environment highqualitycaptions with the help of tensesflow
4.	SimaoHerdade, Armin Kappeler, KofiBoakye, Joao Soares	Dataset: MS COCOMethod: architecture model using CNN as well as NLPtechniques	Using CNN and LSTMmodelthe image'scaptionis generated.
5.	Manish Raypurkar, Abhishek Supe, Pratik Bhumkar, PravinBorse, Dr. Shabnam Sayyad	Dataset: Flickr_8k Method: CNN and LSTM modeltoextract features and sequence the words and finally generating captions.	Proposed model is basedonmultilabel Neural networks
6.	Oriol	Dataset:MS	Proposedmodelis

	vinyals,AlexanderToshev,S amyBengio,DumitruErhan.	COCO Method:Deep recurrent network	basedonneural networksystem
7.	JianhuiChen,WenqiangDon g,Michen Li.	Dataset:flickr_8 k,MS COCO Method:LRCN, CNN,RNN	Using CNN,LSTM,RNN models the image's captionisgenerated.
8.	Peter Anderson,XiaodongHe,Chr isBuchler,DamienTeney,M arkJohnson,StephenGould, And Lei Zhang	Dataset:MS COCO Method:Bottom up and Top down mechanism	UsingBottomup And Top down mechanism the image caption is generated
9.	JyotiAneja,Aditya DeshpandeAndAlexander G Schwing	Dataset:Flickr_8 k Method:CNN And RNNfor spatialimages	Using CNN and RNNthecaptionfor image is generated
10.	ShuangBaiAndShanAn	Dataset:MS COCO Method:Neural Network	Using the dataset and neural network conceptsthecaption for image is generated.
Table1			

## Key Gaps in The Literature

- Could not describe the captions based on different targets.
- Does not recognize the minute changes in the image during the captioning.
- Loss is higher for CNNs than RNN.
- Accuracy and Precision are not the ideal parameter for the model evaluation.
- Insufficient in transforming objects into words.
- Does not work on the large dataset.
- Bottom Up mechanism performs superiorly as compared to Top down mechanism and language LSTM.
- Each image has at least 5 Captions.

# Chapter-3

## SystemDevelopment

### REQUIREMENTSANDANALYSIS

#### SYSTEMREQUIREMENTS

All computer software requires specific hardware parts or additional software resources to function properly on a computer. These prerequisites are referred to as computer system requirements, and they are frequently used as recommendations rather than strict regulations. Two sets of system requirements are typically defined by software. minimal as well as suggested. System requirements typically rise over time due to the growing demand for more processing power and resources in newer software versions. Industry observers contend that rather than technological breakthroughs, this trend is more responsible for the upgrades of current computer systems.

#### HARDWAREREQUIREMENTS

- System:i3Processor
- HardDisk:500GB.
- Monitor:15''LED
- InputDevices:Keyboard,Mouse
- Ram: 4GB.

#### SOFTWAREREQUIREMENTS

- Platform:GoogleColab
- CodingLanguage:Python

## ProjectDesignAndArchitecture

Using deep learning and computer vision, an image caption generator can determine its context and comment on it with pertinent captions. It involves using datasets provided during model training to label an image with English keywords. The CNN model called Xception has been trained using the imagenet dataset. The extraction of characteristics from images is done by Xception. The LSTM model will receive these extracted features and use them to create the image caption..

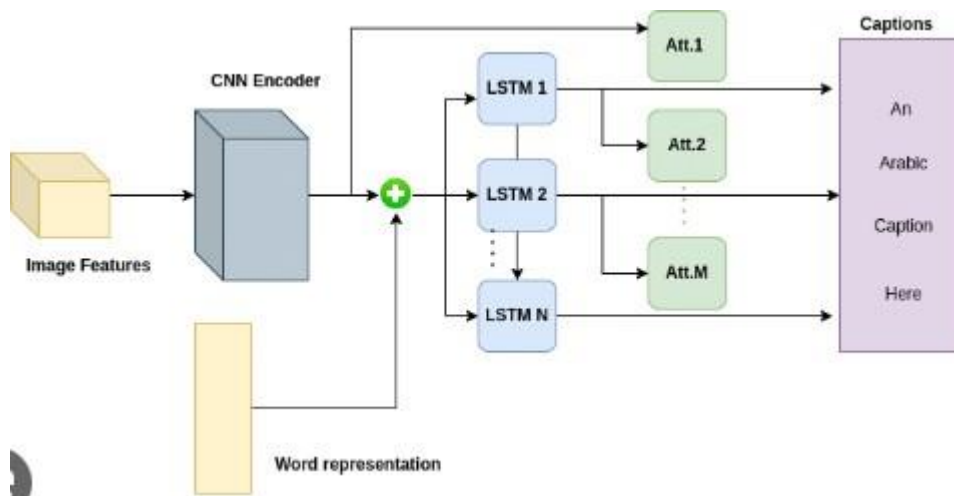


Fig3.ArchitectureAndDesignOfProject

## Data Preparation

For tasks including finding objects, segmentation, and captioning, an in-depth image recognition dataset is called COCO (Common Objects in Context). It has more than 330,000 photos, each with five captions that describe the scene and 80 object categories annotated on them. Many cutting-edge object detection and segmentation models have been trained and evaluated using the COCO dataset, which is extensively used in computer vision research.

The images and their annotations make up the two main components of the dataset. The images are arranged into a hierarchy of directories, with the train, validation, and test sets contained in subdirectories of the top-level directory. JSON files containing annotations are supplied; each file corresponds to a single image.

The following details are included with every annotation in the dataset:-

- Name of image file
- Image dimensions (height and width)
- A list of items containing the data listed below: class of objects (such as "person", "car"); coordinates of the bounding box; segmentation mask (x, y, width, height); Important ideas and where they fit in (if available).
- Five captions that explain the image.
- Additional information is also available from COCO dataset including license, supercategories for images and coco-stuff.
- For the application of the image caption generator we used the dataset names flickr\_8k dataset
- This dataset contains a wider range of images that has many different types of situation and scenes.
- Flickr\_8k dataset has 8000 images and every image has 5 captions.
- We divided the entire dataset of 8000 images as 6000, 1000 and 1000 as training, validation and testing respectively.
- Every image has different dimensions.

## Libraries Used

- TensorFlow: It is an open source library that supports deep learning using Python and C frameworks.
- Keras: It is an open source Python library that allows to evaluate the deep learning models.
- Pillow: Pillow is a Python image library (PIL), that adds support for opening, manipulating, and saving images.
- Numpy: To work with arrays, NumPy library is used.
- Matplotlib: Library to create static and animated visualisations in Python framework.
- Transformers: This is the flagship library that provides state of the art pre-trained models for the natural language understanding, generation, and translation. It also supports both PyTorch and TensorFlow frameworks and includes a wide range of models from BERT to GPT.
- Tokenizers: It offers efficient tokenization algorithms essential for preprocessing the text data before feeding it into the deep learning models. It also supports various tokenization strategies including byte pair encoding and word piece tokenization.
- Datasets: They simplify the process of loading and pre-processing datasets for the training and evaluation. It provides a unified interface for accessing various datasets commonly used in natural language processing tasks.

These libraries form the core infrastructure for building, Fine-Tuning and deploying the advanced NLP models and applications within the Hugging Face ecosystem. They do enable researchers and developers to easily access and leverage cutting edge techniques in the field.

## Implementation

### Working Explanation

1. To create a caption for an image, the user uploads the image.
2. CNN is used to process a grayscale image and identify the objects.
3. CNN is used to process a grayscale image and identify the objects.
4. CNN extracts significant image features by scanning images from left to right and from top to bottom.



5. Using an activation function and a variety of layers, including pooling, convolutional, and fully connected, We succeeded in clearing the features from each image.

6. After that, LSTM is used.

### Algorithms

1. Convolutional Neural Network
2. Long Short-Term Memory

### Overview On CNN

One kind of the model for deep learning used for processing data with a grid pattern, like images, is the convolutional neural network (CNN). To classify a probabilistic object with values ranging from 0 to 1, each input image will pass through a series of fully connected layers (FC), pooling, convolution layers with filters (Kernels), and the Softmax function. CNN models with deep learning are utilized for training and testing. CNNs are distinct from RNNs and other neural networks due to their convolutional layers. The input is transformed in a convolutional layer before being forwarded to the following layer. A CNN applies filters to change the data.

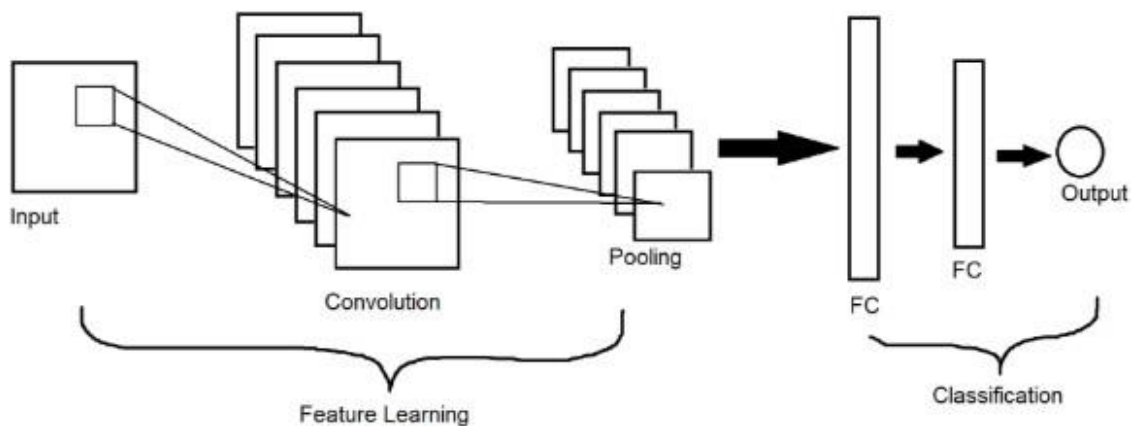


Fig4.CNN

### Some advantages of CNN:

- Convolutional neural networks (CNNs) are one kind of deep learning model which processes grid-patterned data such as images.
- To classify an object with probabilistic values between 0 and 1, each input image will pass through a series of convolution layers with filters (Kernels), pooling, fully connected layers (FC), and the Softmax function. CNN models with deep learning are employed for both training and testing.

- CNNs have convolutional layers, which set them apart from RNNs and other neural networks. Before being passed to the next layer, the input is transformed in a convolutional layer. CNN modifies the data by applying filters.

### **OverviewOfLSTM**

Recurrent neural networks of the LSTM type can recognise order dependence in sequence prediction problems. In difficult problem domains such as machine translation and speech recognition, among others, this behavior is required. One intricate subfield of deep learning is LSTMs. In difficult problem domains like speech recognition and machine translation, among others, this behavior is necessary. One intricate subfield of deep learning is LSTM's.

#### **SomeAdvantagesofLSTM are:**

- Provides us with a wider range of parameters including input and output biases and learning rates.
- Using LSTM the complexity to update each weight is lowered to  $O(1)$ .

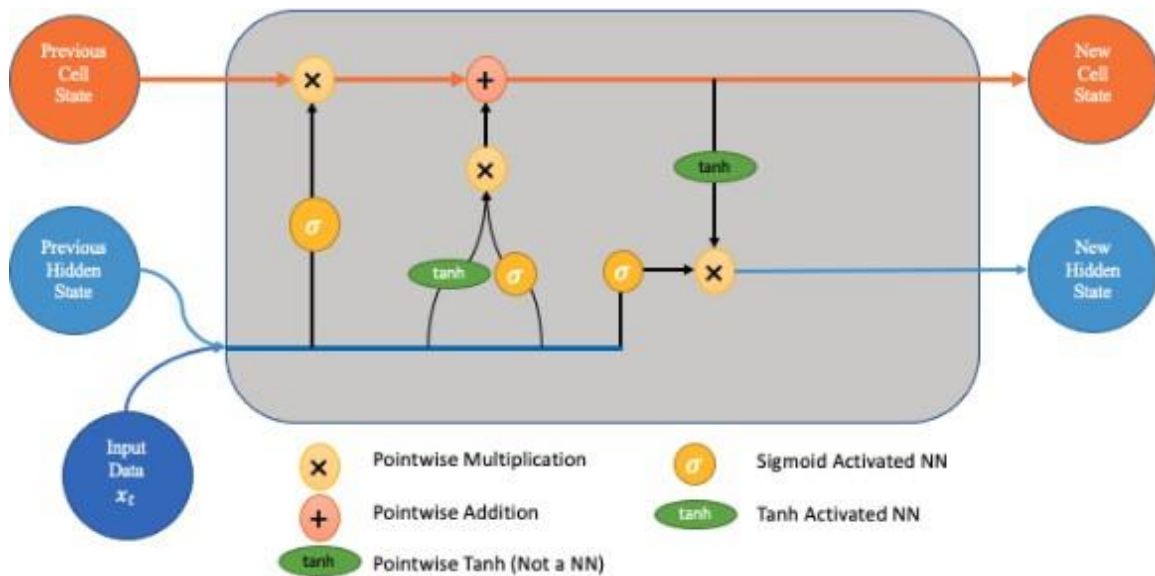


Fig5.LSTM

## CNN-LSTM Architecture Model

In order to support sequence prediction, CNN's LSTM architecture combines LSTMs with Layers in a convolutional neural network (CNN) operate to extract characteristics from input data.

CNN-LSTMs were created for applications such as textual description generation from image sequences (e.g., videos) and visual time series prediction problems. In particular, the issue of

- Activity Recognition: Making a textual description of a procedure which appears in a series of methods.
- Image Description: Generating a textual description of a single image.
- Video Description: Making an essay summary of a set of images.

Although we will refer to this architecture as "CNN LSTM," the original name for it was "LRCN model stands for Long-Term Recurrent Convolutional Network. In order to extract Features of the illustration, CNN is utilised. The Xception pre-trained model will be employed.

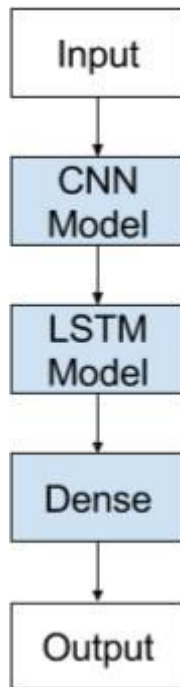


Fig6.CNN-LSTMModel

## Techniques

- AddLibraries.
  - Providethe2017COCO(CommonObjectsandContexts)Dataset.(PreprocessingData)
  - UseCNNtodeterminewhichobjectsareinthe picture. Tokenize and preprocess the captions.
  - Predictthesentence'snextwordusingLSTM.CreateaData GeneratorandViewCaptioned Images.

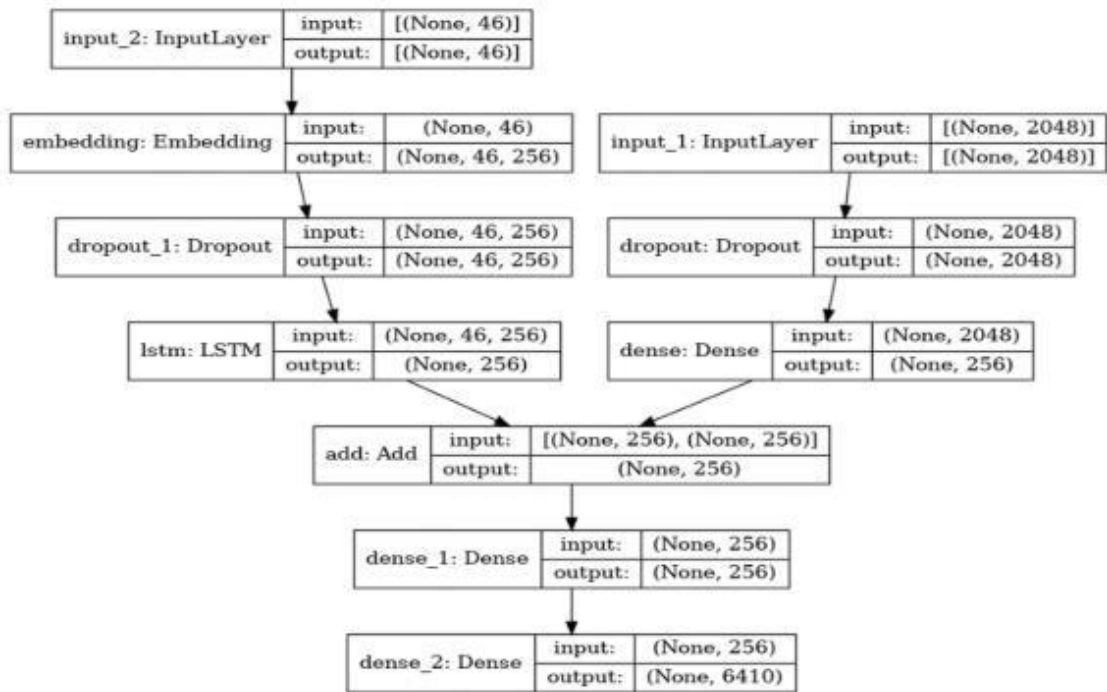
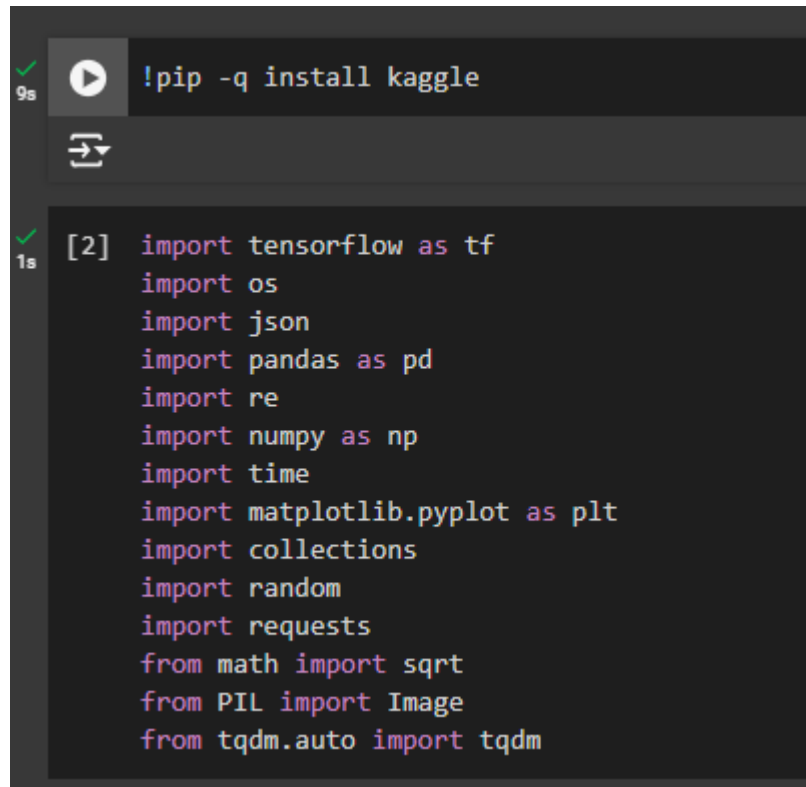


Fig7.WorkflowDiagram

## CodeSnippets



```
✓ 9s !pip -q install kaggle  
✓ 1s [2] import tensorflow as tf  
import os  
import json  
import pandas as pd  
import re  
import numpy as np  
import time  
import matplotlib.pyplot as plt  
import collections  
import random  
import requests  
from math import sqrt  
from PIL import Image  
from tqdm.auto import tqdm
```

Fig8.ImportingLibraries

```

23s [3] if not os.path.exists('/content/data/'):

    api_token = {"username": "<-- your username -->",
                 "key": "<-- your api key -->"}

    with open('/content/kaggle.json', 'w') as file:
        json.dump(api_token, file)

    os.environ["KAGGLE_CONFIG_DIR"] = "/content/"

    os.system('kaggle datasets download -d adityajn105/flickr8k')
    os.makedirs('/content/data/', exist_ok=True)
    os.system('mv /content/flickr8k.zip /content/data/flickr8k.zip')
    os.system('unzip -q /content/data/flickr8k.zip -d /content/data/')
    os.remove('/content/data/flickr8k.zip')

0s [4] captions = pd.read_csv('/content/data/captions.txt')
captions['image'] = captions['image'].apply(
    lambda x: f'/content/data/Images/{x}')
captions.head()

```

```

0s [4]

```

	image	caption
0	/content/data/Images/1000268201_693b08cb0e.jpg	A child in a pink dress is climbing up a set o...
1	/content/data/Images/1000268201_693b08cb0e.jpg	A girl going into a wooden building .
2	/content/data/Images/1000268201_693b08cb0e.jpg	A little girl climbing into a wooden playhouse .
3	/content/data/Images/1000268201_693b08cb0e.jpg	A little girl climbing the stairs to her playh...
4	/content/data/Images/1000268201_693b08cb0e.jpg	A little girl in a pink dress going into a woo...

Next steps: [Generate code with captions](#) [View recommended plots](#)

```

0s [5] def preprocess(text):
    text = text.lower()
    text = re.sub(r'^\w\s', '', text)
    text = re.sub('\s+', ' ', text)
    text = text.strip()
    text = '[start] ' + text + ' [end]'
    return text

3s [6] captions['caption'] = captions['caption'].apply(preprocess)
captions.head()

```

Fig9.Importing Data

	image	caption
0	/content/data/Images/1000268201_693b08cb0e.jpg	[start] a child in a pink dress is climbing up...
1	/content/data/Images/1000268201_693b08cb0e.jpg	[start] a girl going into a wooden building [end]
2	/content/data/Images/1000268201_693b08cb0e.jpg	[start] a little girl climbing into a wooden p...
3	/content/data/Images/1000268201_693b08cb0e.jpg	[start] a little girl climbing the stairs to h...
4	/content/data/Images/1000268201_693b08cb0e.jpg	[start] a little girl in a pink dress going in...

Next steps: [Generate code with captions](#) [View recommended plots](#)

```

random_row = captions.sample(1).iloc[0]
print(random_row.caption)
print()
im = Image.open(random_row.image)
im

```

Fig10.InitialCodeGeneration

[start] two men in midair fighting in a professional wrestling ring [end]



Fig11. First Outcome



```
✓ [8] MAX_LENGTH = 40
0s VOCABULARY_SIZE = 10000
    BATCH_SIZE = 32
    BUFFER_SIZE = 1000
    EMBEDDING_DIM = 512
    UNITS = 512


✓ [9] tokenizer = tf.keras.layers.TextVectorization(
1s     max_tokens=VOCABULARY_SIZE,
     standardize=None,
     output_sequence_length=MAX_LENGTH)

    tokenizer.adapt(captions['caption'])

✓ [10] word2idx = tf.keras.layers.StringLookup(
0s     mask_token="",
     vocabulary=tokenizer.get_vocabulary())

    idx2word = tf.keras.layers.StringLookup(
     mask_token="",
     vocabulary=tokenizer.get_vocabulary(),
     invert=True)
```

Fig12. Assigning Values to element

```
0s  img_to_cap_vector = collections.defaultdict(list)
for img, cap in zip(captions['image'], captions['caption']):
    img_to_cap_vector[img].append(cap)

img_keys = list(img_to_cap_vector.keys())
random.shuffle(img_keys)

slice_index = int(len(img_keys)*0.8)
img_name_train_keys, img_name_val_keys = (img_keys[:slice_index],
                                           img_keys[slice_index:])

train_imgs = []
train_captions = []
for imgt in img_name_train_keys:
    capt_len = len(img_to_cap_vector[imgt])
    train_imgs.extend([imgt] * capt_len)
    train_captions.extend(img_to_cap_vector[imgt])

val_imgs = []
val_captions = []
for imgv in img_name_val_keys:
    capv_len = len(img_to_cap_vector[imgv])
    val_imgs.extend([imgv] * capv_len)
    val_captions.extend(img_to_cap_vector[imgv])
```

Fig13.ImagetocaptionGeneration

```
✓ [12] len(train_imgs), len(train_captions), len(val_imgs), len(val_captions)
0s (32360, 32360, 8095, 8095)

✓ [13] def load_data(img_path, caption):
0s     img = tf.io.read_file(img_path)
     img = tf.io.decode_jpeg(img, channels=3)
     img = tf.keras.layers.Resizing(299, 299)(img)
     img = img / 255.
     caption = tokenizer(caption)
     return img, caption
```

Fig14. Assigningthelength

```
✓ [14] train_dataset = tf.data.Dataset.from_tensor_slices(
0s     (train_imgs, train_captions))

     train_dataset = train_dataset.map(
         load_data, num_parallel_calls=tf.data.AUTOTUNE
     ).shuffle(BUFFER_SIZE).batch(BATCH_SIZE)

     val_dataset = tf.data.Dataset.from_tensor_slices(
         (val_imgs, val_captions))

     val_dataset = val_dataset.map(
         load_data, num_parallel_calls=tf.data.AUTOTUNE
     ).shuffle(BUFFER_SIZE).batch(BATCH_SIZE)

✓ [15] image_augmentation = tf.keras.Sequential(
0s     [
         tf.keras.layers.RandomFlip("horizontal"),
         tf.keras.layers.RandomRotation(0.2),
         tf.keras.layers.RandomContrast(0.3),
     ]
     )
```

Fig15. TrainingDataset

```
✓ 0s ▶ def CNN_Encoder():
    inception_v3 = tf.keras.applications.InceptionV3(
        include_top=False,
        weights='imagenet'
    )
    inception_v3.trainable = False

    output = inception_v3.output
    output = tf.keras.layers.Reshape(
        (-1, output.shape[-1]))(output)

    cnn_model = tf.keras.models.Model(inception_v3.input, output)
    return cnn_model
```

Fig16.CNNEncoder

```

class TransformerEncoderLayer(tf.keras.layers.Layer):

    def __init__(self, embed_dim, num_heads):
        super().__init__()
        self.layer_norm_1 = tf.keras.layers.LayerNormalization()
        self.layer_norm_2 = tf.keras.layers.LayerNormalization()
        self.attention = tf.keras.layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim)
        self.dense = tf.keras.layers.Dense(embed_dim, activation="relu")

    def call(self, x, training):
        x = self.layer_norm_1(x)
        x = self.dense(x)

        attn_output = self.attention(
            query=x,
            value=x,
            key=x,
            attention_mask=None,
            training=training
        )

        x = self.layer_norm_2(x + attn_output)
        return x

```

Fig17. TransformerEncodingLayer

```

class Embeddings(tf.keras.layers.Layer):

    def __init__(self, vocab_size, embed_dim, max_len):
        super().__init__()
        self.token_embeddings = tf.keras.layers.Embedding(
            vocab_size, embed_dim)
        self.position_embeddings = tf.keras.layers.Embedding(
            max_len, embed_dim, input_shape=(None, max_len))

    def call(self, input_ids):
        length = tf.shape(input_ids)[-1]
        position_ids = tf.range(start=0, limit=length, delta=1)
        position_ids = tf.expand_dims(position_ids, axis=0)

        token_embeddings = self.token_embeddings(input_ids)
        position_embeddings = self.position_embeddings(position_ids)

        return token_embeddings + position_embeddings

```

[19] Embeddings(tokenizer.vocabulary\_size(), EMBEDDING\_DIM, MAX\_LENGTH)(next(iter(train\_dataset))[1]).shape

TensorShape([32, 40, 512])

Fig18. Defining length, range, and dimension

```
class TransformerDecoderLayer(tf.keras.layers.Layer):

    def __init__(self, embed_dim, units, num_heads):
        super().__init__()
        self.embedding = Embeddings(
            tokenizer.vocabulary_size(), embed_dim, MAX_LENGTH)

        self.attention_1 = tf.keras.layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim, dropout=0.1
        )
        self.attention_2 = tf.keras.layers.MultiHeadAttention(
            num_heads=num_heads, key_dim=embed_dim, dropout=0.1
        )


        self.layernorm_1 = tf.keras.layers.LayerNormalization()
        self.layernorm_2 = tf.keras.layers.LayerNormalization()
        self.layernorm_3 = tf.keras.layers.LayerNormalization()

        self.ffn_layer_1 = tf.keras.layers.Dense(units, activation="relu")
        self.ffn_layer_2 = tf.keras.layers.Dense(embed_dim)

        self.out = tf.keras.layers.Dense(tokenizer.vocabulary_size(), activation="softmax")

        self.dropout_1 = tf.keras.layers.Dropout(0.3)
        self.dropout_2 = tf.keras.layers.Dropout(0.5)
```

Fig19. TransformerDecodingLayer

```
0s  out_2 = self.layernorm_2(out_1 + attn_output_2)

ffn_out = self.ffn_layer_1(out_2)
ffn_out = self.dropout_1(ffn_out, training=training)
ffn_out = self.ffn_layer_2(ffn_out)

ffn_out = self.layernorm_3(ffn_out + out_2)
ffn_out = self.dropout_2(ffn_out, training=training)
preds = self.out(ffn_out)
return preds

def get_causal_attention_mask(self, inputs):
    input_shape = tf.shape(inputs)
    batch_size, sequence_length = input_shape[0], input_shape[1]
    i = tf.range(sequence_length)[: , tf.newaxis]
    j = tf.range(sequence_length)
    mask = tf.cast(i >= j, dtype="int32")
    mask = tf.reshape(mask, (1, input_shape[1], input_shape[1]))
    mult = tf.concat(
        [tf.expand_dims(batch_size, -1), tf.constant([1, 1], dtype=tf.int32)],
        axis=0
    )
    return tf.tile(mask, mult)
```

Fig20. Decoding Layer



```
class ImageCaptioningModel(tf.keras.Model):


    def __init__(self, cnn_model, encoder, decoder, image_aug=None):
        super().__init__()
        self.cnn_model = cnn_model
        self.encoder = encoder
        self.decoder = decoder
        self.image_aug = image_aug
        self.loss_tracker = tf.keras.metrics.Mean(name="loss")
        self.acc_tracker = tf.keras.metrics.Mean(name="accuracy")

    def calculate_loss(self, y_true, y_pred, mask):
        loss = self.loss(y_true, y_pred)
        mask = tf.cast(mask, dtype=loss.dtype)
        loss *= mask
        return tf.reduce_sum(loss) / tf.reduce_sum(mask)

    def calculate_accuracy(self, y_true, y_pred, mask):
        accuracy = tf.equal(y_true, tf.argmax(y_pred, axis=2))
        accuracy = tf.math.logical_and(mask, accuracy)
        accuracy = tf.cast(accuracy, dtype=tf.float32)
        mask = tf.cast(mask, dtype=tf.float32)
        return tf.reduce_sum(accuracy) / tf.reduce_sum(mask)
```

Fig21.ICModelWithLossAndAccuracy



```
0s  def compute_loss_and_acc(self, img_embed, captions, training=True):
    encoder_output = self.encoder(img_embed, training=True)
    y_input = captions[:, :-1]
    y_true = captions[:, 1:]
    mask = (y_true != 0)
    y_pred = self.decoder(
        y_input, encoder_output, training=True, mask=mask
    )
    loss = self.calculate_loss(y_true, y_pred, mask)
    acc = self.calculate_accuracy(y_true, y_pred, mask)
    return loss, acc

def train_step(self, batch):
    imgs, captions = batch

    if self.image_aug:
        imgs = self.image_aug(imgs)

    img_embed = self.cnn_model(imgs)

    with tf.GradientTape() as tape:
        loss, acc = self.compute_loss_and_acc(
            img_embed, captions
        )
```

Fig22.ComputingLossAnd Accuracy

```
[22] encoder = TransformerEncoderLayer(EMBEDDING_DIM, 1)
decoder = TransformerDecoderLayer(EMBEDDING_DIM, UNITS, 8)

cnn_model = CNN_Encoder()
caption_model = ImageCaptioningModel(
    cnn_model=cnn_model, encoder=encoder, decoder=decoder, image_aug=image_augmentation,
)

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/inception\_v3/inception\_v3\_weights\_tf\_dim\_ordering\_tf\_kernels\_notop.h5
87910968/87910968 [=====] - 0s 0us/step

cross_entropy = tf.keras.losses.SparseCategoricalCrossentropy(
    from_logits=False, reduction="none"
)

early_stopping = tf.keras.callbacks.EarlyStopping(patience=3, restore_best_weights=True)

caption_model.compile(
    optimizer=tf.keras.optimizers.Adam(),
    loss=cross_entropy
)
```

Fig23. Encoder Decoder

```
history = caption_model.fit(
    train_dataset,
    epochs=5,
    validation_data=val_dataset,
    callbacks=[early_stopping]
)

Epoch 1/5
1012/1012 [=====] - 311s 278ms/step - loss: 4.3215 - acc: 0.2231 - val_loss: 3.7957 - val_acc: 0.3202
Epoch 2/5
1012/1012 [=====] - 276s 269ms/step - loss: 3.5071 - acc: 0.3369 - val_loss: 3.5019 - val_acc: 0.3484
Epoch 3/5
1012/1012 [=====] - 309s 301ms/step - loss: 3.2519 - acc: 0.3630 - val_loss: 3.4091 - val_acc: 0.3617
Epoch 4/5
1012/1012 [=====] - 307s 300ms/step - loss: 3.0939 - acc: 0.3786 - val_loss: 3.3570 - val_acc: 0.3676
Epoch 5/5
1012/1012 [=====] - 274s 268ms/step - loss: 2.9738 - acc: 0.3900 - val_loss: 3.3221 - val_acc: 0.3730
```

Fig24. TrainingDataset

```
[27] idx = random.randrange(0, len(val_imgs))
img_path = val_imgs[idx]

pred_caption = generate_caption(img_path)
print('Predicted Caption:', pred_caption)
print()
Image.open(img_path)

Predicted Caption: a white crane is flying in the air
```




Fig25.PredictedOutput

```
os [28] url = "https://encrypted-tbn0.gstatic.com/images?q=tbn:ANd9GcT2j6yc1bKYDav4BGUKLAdTvSFXp1gtuzy5DQ&usqp=CAU"
im = Image.open(requests.get(url, stream=True).raw)
im.save('tmp.jpg')

pred_caption = generate_caption('tmp.jpg')
print('Predicted Caption:', pred_caption)
print()
im
```

↔ Predicted Caption: a black and white dog is jumping over a hurdle



Fig26.PredictedOutput

```
CHUNK_SIZE = 40960
DATA_SOURCE_MAPPING = 'coco-2017-dataset:https%3A%2F%2Fstorage.googleapis.com%2Fkaggle-data-sets%2F857191%2F1462296%2Fbundle%2Farchive.zip%3FX-Goog-Algorithm%3D'

KAGGLE_INPUT_PATH='/kaggle/input'
KAGGLE_WORKING_PATH='/kaggle/working'
KAGGLE_SYMLINK='kaggle'

!umount /kaggle/input/ 2> /dev/null
shutil.rmtree('/kaggle/input', ignore_errors=True)
os.makedirs(KAGGLE_INPUT_PATH, 0o777, exist_ok=True)
os.makedirs(KAGGLE_WORKING_PATH, 0o777, exist_ok=True)


try:
    os.symlink(KAGGLE_INPUT_PATH, os.path.join(".", 'input'), target_is_directory=True)
except FileExistsError:
    pass
try:
    os.symlink(KAGGLE_WORKING_PATH, os.path.join(".", 'working'), target_is_directory=True)
except FileExistsError:
    pass
```

Fig27.DataSourceMapping



```
class TransformerEncoderLayer(tf.keras.layers.Layer):  
  
    def __init__(self, embed_dim, num_heads):  
        super().__init__()  
        self.layer_norm_1 = tf.keras.layers.LayerNormalization()  
        self.layer_norm_2 = tf.keras.layers.LayerNormalization()  
        self.attention = tf.keras.layers.MultiHeadAttention(  
            num_heads=num_heads, key_dim=embed_dim)  
        self.dense = tf.keras.layers.Dense(embed_dim, activation="relu")  
  
    def call(self, x, training):  
        x = self.layer_norm_1(x)  
        x = self.dense(x)  
  
        attn_output = self.attention(  
            query=x,  
            value=x,  
            key=x,  
            attention_mask=None,  
            training=training  
        )  
  
        x = self.layer_norm_2(x + attn_output)  
        return x
```

Fig28.OutputLayer



```
def calculate_accuracy(self, y_true, y_pred, mask):
    accuracy = tf.equal(y_true, tf.argmax(y_pred, axis=2))
    accuracy = tf.math.logical_and(mask, accuracy)
    accuracy = tf.cast(accracy, dtype=tf.float32)
    mask = tf.cast(mask, dtype=tf.float32)
    return tf.reduce_sum(accuracy) / tf.reduce_sum(mask)

def compute_loss_and_acc(self, img_embed, captions, training=True):
    encoder_output = self.encoder(img_embed, training=True)
    y_input = captions[:, :-1]
    y_true = captions[:, 1:]
    mask = (y_true != 0)
    y_pred = self.decoder(
        y_input, encoder_output, training=True, mask=mask
    )
    loss = self.calculate_loss(y_true, y_pred, mask)
    acc = self.calculate_accuracy(y_true, y_pred, mask)
    return loss, acc
```

Fig29.AccuracyAndLossFunction

## Key Challenges

### Semantic Understanding

Solution: Employ deep learning models like CNNs for the feature extraction and combine with RNNs.

### Visual Recognition

Solution: Utilize Pre-Trained object detection model or the fine-tune them for specific datasets to improve accuracy in recognition the visuals.

### Language Understanding

Solution: Use techniques like attention mechanisms and beam search during caption generation to focus on relevant parts of the image and improve coherence in the language generation.

### Handling Ambiguity

Solution: solution is to augment the dataset with diverse images and the captions to expose the model to various contexts, reducing ambiguity.

### **Handling rare or unseen objects**

Solutions: Leveraging transfer learning by finetuning model on the specific captioning tasks and utilizing techniques like domain adaptation to generalize better to rare or the unseen objects.

# Chapter-4

## Testing

### TestingStrategy

The trained model will now be loaded and predictions will be produced by a separate file called `testing_caption_generator.py`. We will extract the words from their index values using the same `tokenizer.p` pickle file since the predictions include the maximum length of index values. Thanks to the advancement of deep neural networks, Image captioning (IC) systems generate a text description of the important objects in an image automatically. Images, whether they are synthetic or real—have advanced significantly in recent years. In human society, IC is essential for tasks like labeling large photos for research studies and helping those with visual impairments see the world. But even the best IC systems—like IBM Image Caption Generator and Microsoft Azure Cognitive Services—can produce inaccurate results, which leaves out important information and deep misunderstanding.

We suggest MetaIC, the first metamorphic testing strategy for IC system validation, as a solution to this issue. Our main hypothesis is that, following object insertion, the object names ought to show directional changes. In particular, MetaIC:

- (1) creates an object corpus by extracting objects from pre-existing images;
- (2) uses innovative Location tuning and object resizing algorithms to insert an object into an image; and
- (3) indicates see pairs with captions that don't show expected differences. We evaluate one commonly-used image captioning API and five state-of-the-art (SOTA) MetaIC-based image captioning models. Using 1,000 seeds, MetaIC successfully reports 16,825 erroneous issues with high precision (84.9%-98.4%).
- (4) A testing strategy includes metrics such as BLEU score, METEOR, ROUGE, which measure the similarity between generated captions and reference captions provided by humans.

### UnitTesting

- Preprocessing Units: It tests the preprocessing steps such as image resizing, normalization and the feature extraction to ensure that they produce the expected and real outputs.
- Model Components: Test each component of the captioning model including the encoder and decoder.

### End-To-EndTesting

- Data Pipeline Testing: Verifies the data pipeline is correct including data loading, feature extraction and preprocessing.



- Performance Testing: Measure the performance and ensures the system that it meets the performance requirements or not.

## Security Testing

- Data Poisoning: Assessing the model to poisoned training data where maliciously images or captions are being inserted into the training data to manipulate the outcome and model's behaviour.
- Privacy Concern: Ensuring that the captioning system do not reveal sensitive information about individual depicted in the images.

## Test Cases And Outcomes

Test cases, which are used to assess an image captioning model's performance, are usually pairs of images and the captions that go with them. The captions that the model produces for every test case are the results. To evaluate the caliber of the generated captions, these results are then contrasted with the ground truth captions, or the real captions. This is an explanation:

### 1. Test scenarios:

- An image from a dataset or an actual situation is the input image.
- Ground Truth Caption: The accurate caption that accurately summarizes the picture's content.

### 2. Results:

- Generated Caption: The caption that an input image's image captioning model generates for it.
- Evaluation Metrics: A number of metrics are commonly used to measure the degree to which the generated caption corresponds to the ground truth caption, including BLEU (Bilingual Evaluation Understudy), METEOR (Metric for Evaluation of Translation with Explicit Ordering), CIDEr (Consensus-based Image Description Evaluation), ROUGE (Recall-Oriented Understudy for Gisting Evaluation), etc.

### 3. Process of evaluation:

- Quantitative Evaluation: Using evaluation metrics, the generated captions are compared to the ground truth captions. These metrics provide numerical scores indicating the quality of the generated captions.
- Qualitative Evaluation: Human reviewers may also subjectively rate the captions based on criteria like fluency, coherence, and relevance.

### 4. An illustration:

- Correct Outcome: The evaluation metrics accurately reflect the accuracy with which the generated caption describes the image's content.



# Chapter-5

## ResultAndEvaluation

### Results

The result of this program is going to be a user being allowed to generate a caption for a visual image using Deep Learning, NLP, and Computer Vision. Using the concept of CNN and LSTM we have built a model and trained the model using flickr\_8k data set and MS COCO for getting the appropriate captions for the given image. We have tried different methods to get the perfect outcome.

After installing the data set and loading it into the model the model will train the dataset and give the desired result as per the image. And we have used Hugging Face as the hosting element in which we can enter the URL of the image or the image from the dataset folder for the outcome.

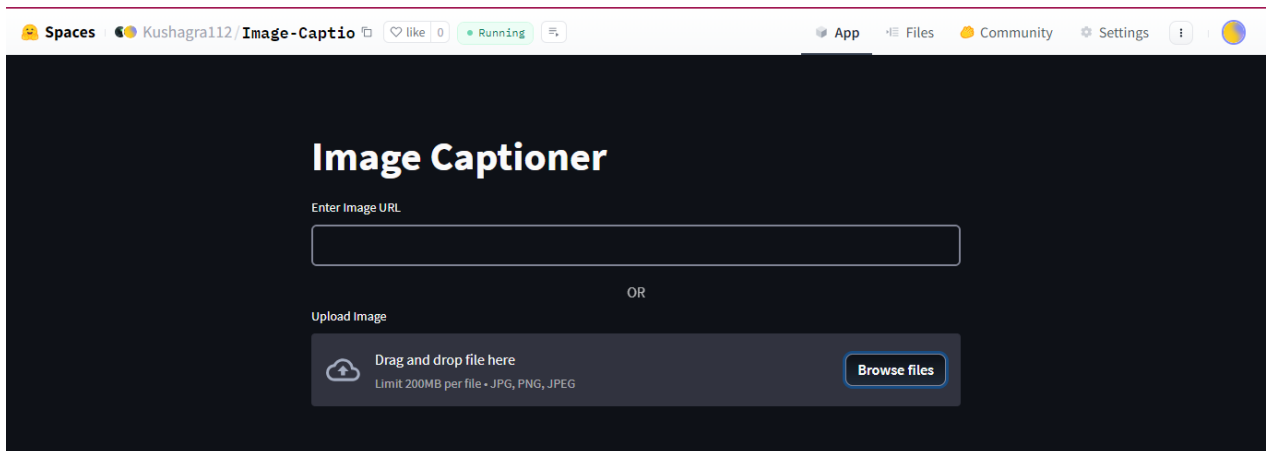


Fig30. HuggingFaceInterface

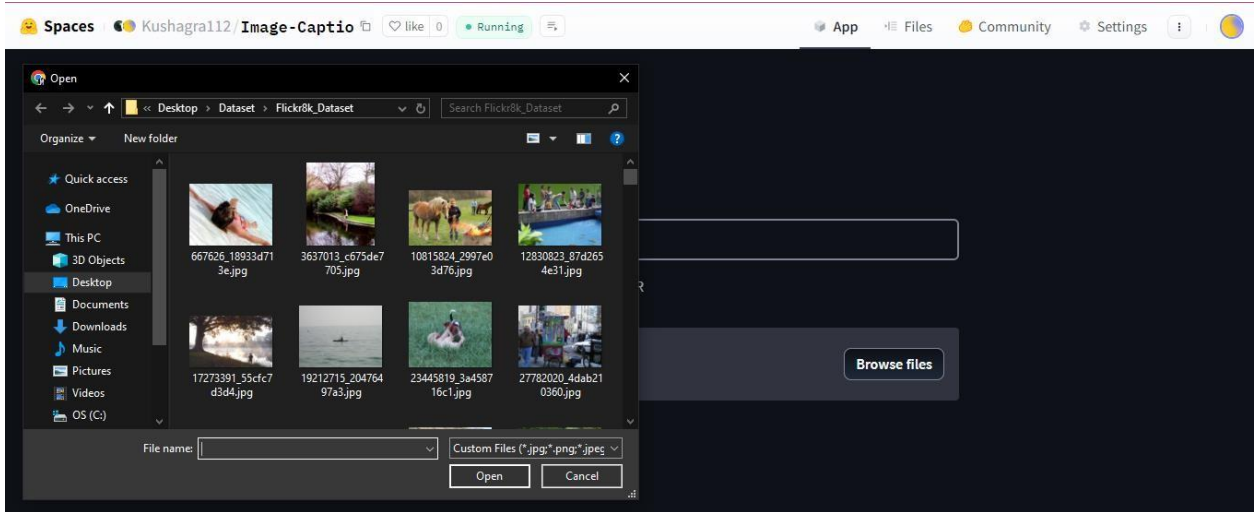


Fig31.Uploadingtheimage



### **Predicted Captions:**

a man standing on a bench in the water

a man standing on a bench in the middle of a street

a group of people standing on a beach

a man standing on a bench in the middle of a park

Fig32. HuggingFaceOutput1(AGroupOfPeople)



69189650\_6687da7280.jpg 126.9KB



### **Predicted Captions:**

a dog is standing on a leash in the grass

a dog standing in a field with a frisbee

a dog is standing in a field with a frisbee

Fig33.HuggingFaceOutput2(Doginthe field)



**Predicted Captions:**

a horse standing in the middle of a dirt road

a man riding a horse on a dirt road

Fig34.HuggingFaceOutput3(WomenRidingHorse)

# Chapter-6

## ConclusionAndFutureScope

### Conclusion

For the purpose of automatically producing captions for the input images, the CNN-LSTM model was developed. There are many contexts in which this idea can be applied. By creating a CNN-LSTM model that can scan any input image, extract information from it, and convert it into a single, natural language sentence in English, we were able to overcome earlier limitations in the area of graphical image captioning.

The primary subjects of discussion were algorithm attention and the application of the attention mechanism. We have managed to produce a model that significantly outperforms the previous image caption generator. By the use of Hugging Faces we are able to get output on a hosting element rather than on Colab.

Deep learning-based image captioning has become a potent instrument for producing textual descriptions of images automatically, and it has found wide-ranging uses in fields like multimedia comprehension, assistive technology, image indexing, and content accessibility. The quality, diversity, and adaptability of generated captions should continue to improve as this field of study develops, opening the door to improved visual content comprehension and human-computer interaction.

### Key Findings

- **Effective Representation Learning:** Deep learning models have proven to be able to simultaneously learn effective representations of both visual and textual data. This is especially true for Convolutional Neural Networks (CNNs) for image feature extraction and Recurrent Neural Networks (RNNs) like Long Short-Term Memory (LSTM) networks for sequence generation.
- **Contextual Understanding:** Image captioning models based on deep learning can produce captions that demonstrate a contextual understanding of the images' content. These models capture complex relationships between objects, scenes, and actions shown in images by utilizing multimodal architectures and attention.

- Deep learning models that have been trained to caption images have some semantic understanding, which enables them to produce captions that are more complex than simple object recognition. With the help of these models, captions can be made more detailed and educational by inferring characteristics, connections, and even abstract ideas from visual input.

## **FutureScope**

Deep learning-based picture captioning has a wide future potential with many directions for research and development. The following are some important areas that show promise for additional study and advancement:

- **Multimodal Understanding:** Richer and more thorough understanding of visual content can be achieved by extending image captioning models to include multiple modalities, such as text, audio, and video. This involves creating models that can interpret intricate multimedia interactions, describe dynamic scenes, and create captions for videos.
- **Attention Mechanisms:** By improving the attention mechanisms in deep learning architectures, models will be better able to concentrate on pertinent image regions or temporal segments during the caption-generating process. It is possible to expand attention mechanisms to handle occlusions, capture fine-grained details, and adapt to various aspects of the image.
- Investigating the combination of generative adversarial networks (GANs) and image captioning can help produce captions that are more aesthetically pleasing and realistic. GANs can help reduce potential discrepancies between the generated text and the visual content, improve the perceptual quality of generated captions, and improve their visual fidelity.
- Investigating methods to foster creativity and storytelling in the captioning of images can result in the creation of captions that are more captivating and narratively focused.
- This entails creating models that can comprehend narrative structures and coherence in addition to incorporating personality, emotion, and stylistic variation into captions.

- **Fine-Grained Understanding and Generation:** By strengthening models' ability to recognize and characterize fine-grained details, attributes, and relationships within images, the specificity and accuracy of generated captions can be improved. This means addressing challenges related to object counting, spatial reasoning, and understanding the background of complex scenes.
- **Ethical Considerations and Bias Mitigation:** In order to ensure an equitable and responsible deployment in real-world applications, it is imperative that image captioning models address ethical considerations such as bias, fairness, and inclusivity. This entails creating methods for reducing biases in training data, encouraging diversity and representation in the creation of captions, and encouraging accountability and transparency in the model-building process.



# REFERENCES

- R.Subash(November2019):AutomaticImageCaptioningUsingConvolutionNeural Networks and LSTM.
- Seung-HoHan,Ho-JinChoi(2020):Domain-SpecificImageCaptionGenerator withSemantic Ontology.
- PranayMathur, Aman Gill, Aayush Yadav, Anurag Mishra and Nand Kumar Bansode (2017): Camera2Caption: A Real-Time Image Caption Generator
- SimaoHerdade,ArminKappeler,KofiBoakye,JoaoSoares(June2019):ImageCaptioning: Transforming Objects into words.
- ManishRaypurkar,AbhishekSupe,PratikBhumkar,PravinBorse,Dr.ShabnamSayyad(March 2021): Deep learning-based Image Caption Generator.
- OriolVinyals,AlexanderToshev,SamyBengio,DumitruErhan(2015):ShowandTell:A Neural Image Caption Generator.
- Jianhui Chen, Wenqiang Dong, Minchen Li (2015): Image Caption Generator based on Deep Neural Networks
- 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.
- JyotiAneja,AdityaDeshpande,andAlexanderGSchwing(2018):Convolutionalimage captioning.
- ShuangBaiandShanAn(2018):Asurveyonautomaticimagecaptiongeneration.
- D.Bahdanau,K.Cho, and Y.Bengio(2015): Neurmachine translationbyjointly learning to align and translate.
- K.Xu, J.Ba, K.Cho, and R.Salakhutdinov (2018): Show attend and tell: Neural image caption generator with visual attention.
- M.Pedersoli,T.Lucas,C.Schmid,andJ.Verbeek(2017):Areasofattentionforimage captioning.
- H.R.Tavakoli,R.Shetty,B.Ali,andJ.Laaksonen(2017):Payingattentiontodescriptions generated by image captioning models.
- A.Matthews, L.Xie, and X.He(2018): Sem Style-learning to generate stylized image captions using unaligned text.
- C.Park,B.Kim,andG.kim(2018):Towardspersonalizedimagecaptioningviamultimodal memory networks.
- X.Chen,MaLin,W.Jiang,J.Yao,andW.Liu(June2018):RegularizingRNNsforcaption generation by reconstructing the past with the present

- T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei (June 2016): Boosting image captioning with attributes.
- Deep Learning in big data Analytics: A comparative study - Scientific Figure on ResearchGate (2021).
- Kaustav et al. (June 2016): A Facial Expression Recognition System to predict Emotions.
- M. Hodosh, P. Young and J. Hockenmaier (2013) "Framing Image Description as a Ranking Task"
- Oriol Vinyals, Alexander Toshev, Samy Bengio, Dumitru Erhan Show and Tell: A Neural Image Caption Generator
- CS231n Winter 2016 Lesson 10 Recurrent Neural Networks, Image Captioning and LSTM

**JAYPEE UNIVERSITY OF INFORMATION TECHNOLOGY, WAKNAGHAT**  
**PLAGIARISM VERIFICATION REPORT**

Date: .....

Type of Document (Tick):  PhD Thesis  M.Tech Dissertation/ Report  B.Tech Project Report  Paper

Name: \_\_\_\_\_ Department: \_\_\_\_\_ Enrolment No \_\_\_\_\_

Contact No. \_\_\_\_\_ E-mail. \_\_\_\_\_

Name of the Supervisor: \_\_\_\_\_

Title of the Thesis/Dissertation/Project Report/Paper (In Capital letters): \_\_\_\_\_

**UNDERTAKING**

I undertake that I am aware of the plagiarism related norms/ regulations, if I found guilty of any plagiarism and copyright violations in the above thesis/report even after award of degree, the University reserves the rights to withdraw/ revoke my degree/report. Kindly allow me to avail Plagiarism verification report for the document mentioned above.

**Complete Thesis/Report Pages Detail:**

- Total No. of Pages =
- Total No. of Preliminary pages =
- Total No. of pages accommodate bibliography/references =

(Signature of Student)

**FOR DEPARTMENT USE**

We have checked the thesis/report as per norms and found **Similarity Index** at .....(%). Therefore, we are forwarding the complete thesis/report for final plagiarism check. The plagiarism verification report may be handed over to the candidate.

(Signature of Guide/Supervisor)

Signature of HOD

**FOR LRC USE**

The above document was scanned for plagiarism check. The outcome of the same is reported below:

Copy Received on	Excluded	Similarity Index (%)	Generated Plagiarism Report Details (Title, Abstract & Chapters)	
	<ul style="list-style-type: none"> <li>• All Preliminary Pages</li> <li>• Bibliography/Images/Quotes</li> <li>• 14 Words String</li> </ul>		Word Counts	
<b>Report Generated on</b>			Character Counts	
		<b>Submission ID</b>	Total Pages Scanned	
			File Size	

Checked by  
Name & Signature

Librarian

**Please send your complete thesis/report in (PDF) with Title Page, Abstract and Chapters in (Word File) through the supervisor at [plagcheck.juit@gmail.com](mailto:plagcheck.juit@gmail.com)**