

A PBL Report on  
**Image Caption Generator Using Data Mining**  
*Submitted to CMR Engineering College*  
*In Partial Fulfillment of the requirements for the Award of Degree of*

**BACHELOR OF TECHNOLOGY**  
**IN**  
**COMPUTER SCIENCE AND ENGINEERING**

Submitted By

<b>SHIVAM KUMAR</b>	<b>(208R1A05N4)</b>
<b>SHEELA RISHITHA</b>	<b>(208R1A05N3)</b>
<b>SHAIK RENISHA</b>	<b>(208R1A05N2)</b>
<b>YARRAM NITCHITHA REDDY</b>	<b>(208R1A05N9)</b>

Under the Esteemed guidance of

**Mrs. M. Bhargavi**

**Assistant Professor**, Department of CSE



**Department of Computer Science & Engineering**

**CMR ENGINEERING COLLEGE**

**UGC AUTONOMOUS**

(Approved by AICTE, NEW DELHI, Affiliated to JNTU, Hyderabad) Kandlakoya, Medchal  
Road, R.R. Dist. Hyderabad-501 401.

**(2023-2024)**

# **CMR ENGINEERING COLLEGE**

## **UGC AUTONOMOUS**

(Approved by AICTE, NEW DELHI, Affiliated to JNTU, Hyderabad) Kandlakoya, Medchal  
Road, R.R. Dist. Hyderabad-501 401.

### **Department of Computer Science & Engineering**



### **CERTIFICATE**

This is to certify that the project entitled “**Image Caption Generator Using Data Mining**”  
is a bonafide work  
carried out by

**SHIVAM KUMAR**

**(208R1A05N4)**

**SHEELA RISHITHA**

**(208R1A05N3)**

**SHAIK RENISHA**

**(208R1A05N2)**

**YARRAM NITCHITHA REDDY**

**(208R1A05N9)**

in partial fulfillment of the requirement for the award of the degree of **BACHELOR OF TECHNOLOGY** in **COMPUTER SCIENCE AND ENGINEERING** from CMR Engineering College, affiliated to JNTU, Hyderabad, under our guidance and supervision.

The results presented in this project have been verified and are found to be satisfactory. The results embodied in this project have not been submitted to any other university for the award of any other degree or diploma.

**Internal Guide**

**Head of the Department**

**Mrs. M. Bhargavi**

Assistant Professor

Department of CSE

CMREC, Hyderabad.

**Dr. Sheo Kumar**

Professor & HOD

Department of CSE

CMREC, Hyderabad.

## **DECLARATION**

This is to certify that the work reported in the present PBL project entitled “**Image Caption Generator Using Data Mining**” is a record of bonafide work done by us in the Department of CSE, CMR Engineering College, JNTU Hyderabad. There ports are based on the PBL work done entirely by us and not copied from any other source. We submit our PBL for further development by any interested students who share similar interests to improve the PBL in the future.

The results embodied in this PBL report have not been submitted to any other University or Institute for the award of any degree or diploma to the best of our knowledge and belief.

<b>SHIVAM KUMAR</b>	<b>(208R1A05N4)</b>
<b>SHEELA RISHITHA</b>	<b>(208R1A05N3)</b>
<b>SHAIK RENISHA</b>	<b>(208R1A05N2)</b>
<b>YARRAM NITCHITHA REDDY</b>	<b>(208R1A05N9)</b>

## **ACKNOWLEDGEMENT**

We are extremely grateful to **Dr. A. Srinivasula Reddy**, Principal and **Dr. Sheo Kumar**, HOD, Department of CSE, CMR Engineering College for their constant support.

We are extremely thankful to **Mrs. M. Bhargavi**, Assistant Professor, Internal Guide, Department of CSE, for his constant guidance, encouragement and moral support throughout the PBL.

We will be failing in duty if we do not acknowledge with grateful thanks to the authors of the references and other literatures referred in this PBL.

We express my thanks to all staff members and friends for all the help and co-ordination extended in bringing out this PBL project successfully in time.

<b>SHIVAM KUMAR</b>	<b>(208R1A05N4)</b>
<b>SHEELA RISHITHA</b>	<b>(208R1A05N3)</b>
<b>SHAIK RENISHA</b>	<b>(208R1A05N2)</b>
<b>YARRAM NITCHITHA REDDY</b>	<b>(208R1A05N9)</b>

# **TABLE OF CONTENTS**

<b><u>CONTENTS</u></b>	<b><u>PAGE. NO.</u></b>
1. <b>Introduction</b>	1
1.1 System modules	2
2. <b>Related Work</b>	3
2.1 Early work on image caption generation	3
2.1 Modern image caption generation	3
3. <b>Data-set and Preprocessing</b>	4
3.1 Data set for image caption generation	4
3.2 Preprocessing for image caption generation	5
4. <b>Proposed Methodology</b>	6
4.1 Data collection and preprocessing	6
4.2 Model training	6
4.3 Model evaluation	6
5. <b>Results</b>	8
5.1 Source code	8
5.2 Output	11
5.3 Discussion	13
6. <b>Conclusion</b>	14
6.1 Future work	14

# 1. INTRODUCTION

Novel applications in the fields of computer vision and artificial intelligence have resulted from the combination of image processing and data mining techniques. The Image Caption Generator project is a fascinating undertaking that utilizes data mining concepts to derive significant patterns and insights from extensive image databases. This project highlights the role that data mining plays in improving the abilities of picture understanding and description, in addition to demonstrating its power in managing visual information.

The development of an automated system capable of producing meaningful and cogent descriptions for photos is the main goal of the Image Caption Generator project. Identifying important characteristics, objects, and contextual information from picture data through the use of data mining techniques is the project's goal. The system discovers how to correlate visual aspects with related language terms by mining patterns and correlations from a variety of picture sources. This allows it to provide captions for new photos that seem human.

Among other things, feature extraction, grouping, and classification are part of the fundamental data mining procedures. Mechanisms for feature extraction sort through an image's visual components and extract pertinent data that adds to our comprehension of the picture as a whole. By putting comparable photos together, clustering algorithms make it easier to spot recurring themes and patterns. Classification algorithms are able to identify objects and context, which allows the system to provide captions that are both semantically coherent and pertinent to the situation.

This experiment not only demonstrates the symbiotic link between data mining and computer vision, but it also has significant practical consequences. A skilled Image Caption Generator has the potential to transform content accessibility for people with visual impairments, improve image search functionality, and contribute to the creation of more intuitive human-computer interfaces. Furthermore, the research provides a useful insight into the multidisciplinary nature of artificial intelligence, connecting the domains of data mining and image processing to push the frontiers of what intelligent systems are capable of.

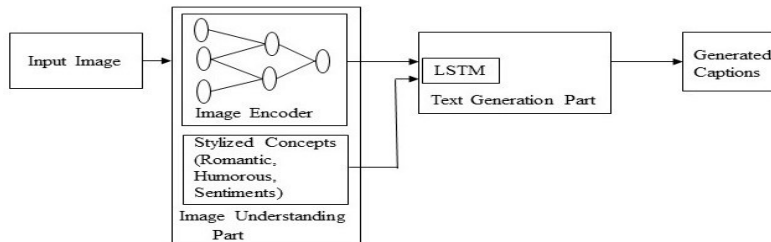


Fig. 1: Block Diagram For Image Caption Generator

## 1.1 SYSTEM MODULES:

**Data Acquisition and Preprocessing:** In this stage, a large and diverse dataset of images and their corresponding captions is acquired. The data is then preprocessed to ensure its quality and consistency.

**Feature Extraction:** In this stage, features are extracted from the images in the dataset. This is done using a CNN model. The extracted features represent a condensed representation of the image content, highlighting the essential elements for caption generation.

**Model Training:** In this stage, a machine learning model is trained to learn the relationship between extracted image features and corresponding captions. An RNN model is commonly used for this task. The training process involves feeding the model with paired image features and captions, allowing it to gradually develop an understanding of the mapping between the visual and linguistic domains.

**Caption Generation:** In this stage, the trained model is used to generate captions for new, unseen images. The extracted features from the new image are fed into the trained model, which then generates a caption based on its learned mapping.

**Evaluation and Refinement:** In this stage, the performance of the image caption generator is evaluated using metrics such as BLEU score and METEOR score. Based on the evaluation results, the system can be refined by adjusting hyper-parameters, modifying model architecture, or augmenting the training dataset.

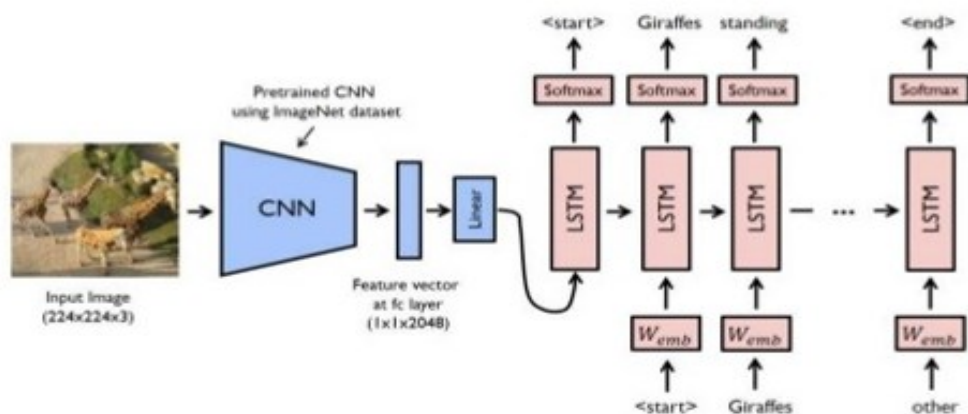


Fig. : 1.1 System Model of mage Caption Generator

## **2. RELATED WORK**

### **2.1 Early Work on Image Caption Generation:**

The earliest work on image caption generation focused on using traditional natural language processing (NLP) techniques to generate captions. These techniques were often limited in their ability to accurately capture the complex semantics of images.

One early approach involved using a bag-of-words model to represent the image content. The bag-of-words model simply counts the number of occurrences of each word in the caption. This information is then used to train a machine learning model to generate captions for new images.

Another early approach involved using a statistical language model to generate captions. The statistical language model predicts the next word in a sentence based on the previous words in the sentence. This information is used to generate captions for new images by starting with a seed word and then predicting the next word in the caption until a complete sentence is generated.

### **2.2 Modern Image Caption Generation with Deep Learning:**

Modern image caption generation systems typically use deep learning techniques to achieve state-of-the-art results. Deep learning techniques allow for the learning of more complex relationships between image features and captions.

One common approach to image caption generation with deep learning involves using a CNN to extract features from the image and an RNN to generate the caption. The CNN is able to capture the spatial patterns and relationships within the image, while the RNN is able to generate a sequential output, such as a caption.

Another common approach to image caption generation with deep learning involves using an encoder-decoder architecture. The encoder extracts features from the image, and the decoder generates a caption based on the extracted features. The encoder and decoder are typically both neural networks.



### **3. DATA-SET AND PREPROCESSING**

#### **3.1 Datasets for Image Caption Generation:**

The quality and diversity of the dataset used to train an image caption generator has a significant impact on its performance. There are a number of factors to consider when choosing a dataset for image caption generation, including:

- **Image Quality:** The images in the dataset should be of high quality and free of noise.
- **Image Diversity:** The dataset should contain a variety of images from different categories and with different objects and scenes.
- **Caption Quality:** The captions in the dataset should be well-written and accurate.
- **Caption Diversity:** The captions in the dataset should be diverse in style and length.
- **Caption Size:** The dataset should contain enough images and captions to train a large machine learning model.

**Some common image caption datasets include:**

- **Flickr8k:** This dataset contains 8,000 images from Flickr and their corresponding captions.
- **MS COCO:** This dataset contains 123,287 images from COCO and their corresponding captions.
- **Visual Genome:** This dataset contains 108,073 images from Visual Genome and their corresponding captions.
- **Open Images:** This dataset contains 9 million images from the web and their corresponding captions.

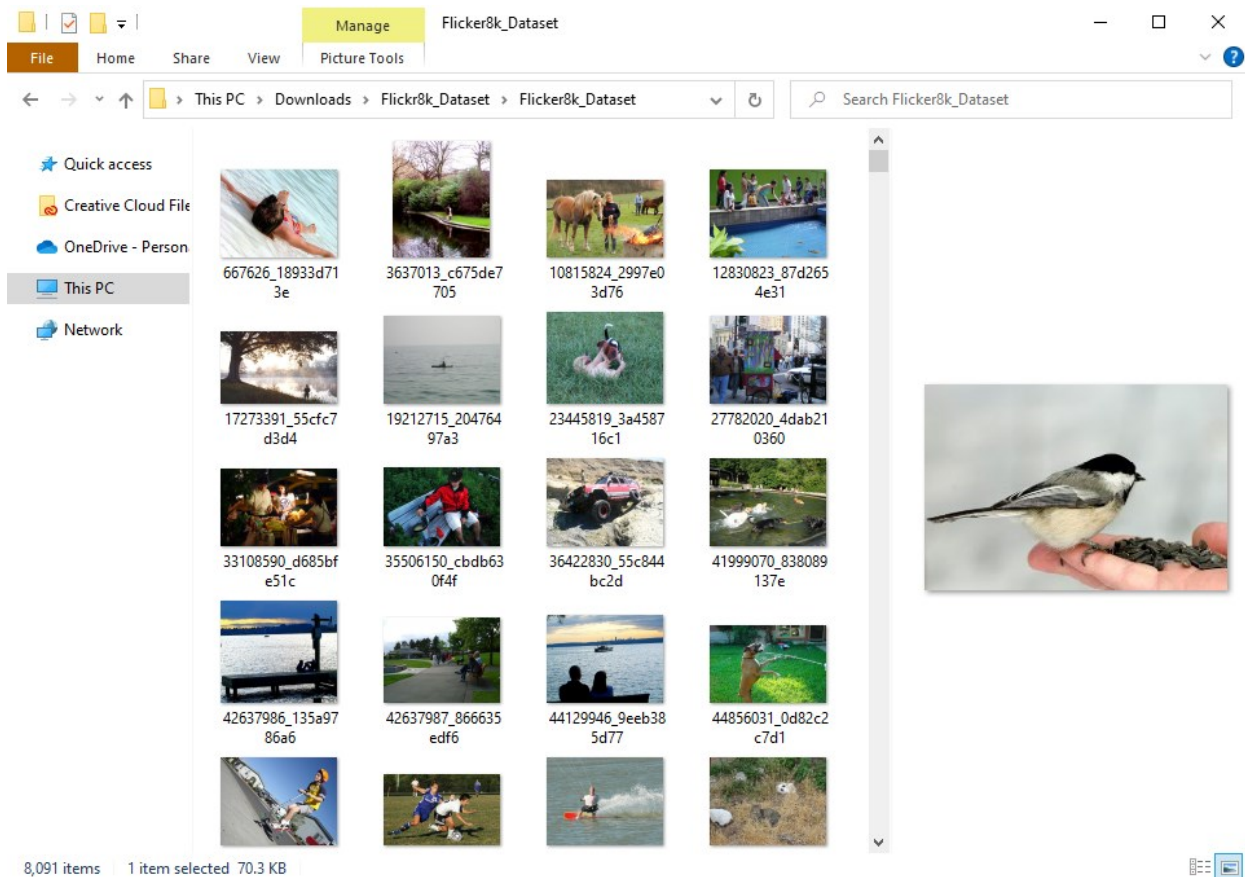


Fig. : 3.1 Flickr8k data-set

## 3.2 Preprocessing for Image Caption Generation:

Once a dataset has been chosen, it is important to preprocess the images and captions before using them to train a machine learning model. Preprocessing involves a number of steps, including:

- **Image Resizing:** All images should be resized to a standard size.
- **Image Normalization:** All images should be normalized to have a mean of 0 and a standard deviation of 1.
- **Caption-Cleaning:** All captions should be cleaned to remove punctuation, numbers, and special characters.
- **Caption Tokenization:** All captions should be tokenized into a sequence of words.
- **Vocabulary Building:** A vocabulary should be built from the words in the captions.

Preprocessing helps to ensure that the images and captions are in a format that can be easily processed by a machine learning model.

## 4. PROPOSED METHODOLOGY

### 4.1 Data Collection and Preprocessing:

- **Collect a large dataset of images and their corresponding captions.** This dataset can be obtained from publicly available sources, such as Flickr or MSCOCO, or by manually collecting images and captions from the web.
- **Preprocess the images and captions.** This may involve resizing images, removing noise, and converting text to lowercase.
- **Split the dataset into training, validation, and test sets.** The training set will be used to train the model, the validation set will be used to tune the model's hyperparameters, and the test set will be used to evaluate the model's performance.

### 4.2 Model Training:

- **Choose a deep learning model architecture.** There are many different deep learning model architectures that can be used for image captioning. Some popular choices include the Encoder-Decoder model, the Recurrent Attention Model (RAM), and the Transformer model.
- **Train the model on the training dataset.** This involves feeding the model images and their corresponding captions and allowing it to learn how to generate captions for new images.
- **Tune the model's hyper-parameters on the validation dataset.** This involves adjusting the model's parameters, such as the learning rate and the number of epochs, to improve its performance.

### 4.3 Model Evaluation:

- **Evaluate the model on the test dataset.** This involves calculating the model's performance metrics, such as BLEU score, METEOR, and ROUGE.
- **Analyze the model's results.** This involves identifying the model's strengths and weaknesses and exploring ways to improve its performance.

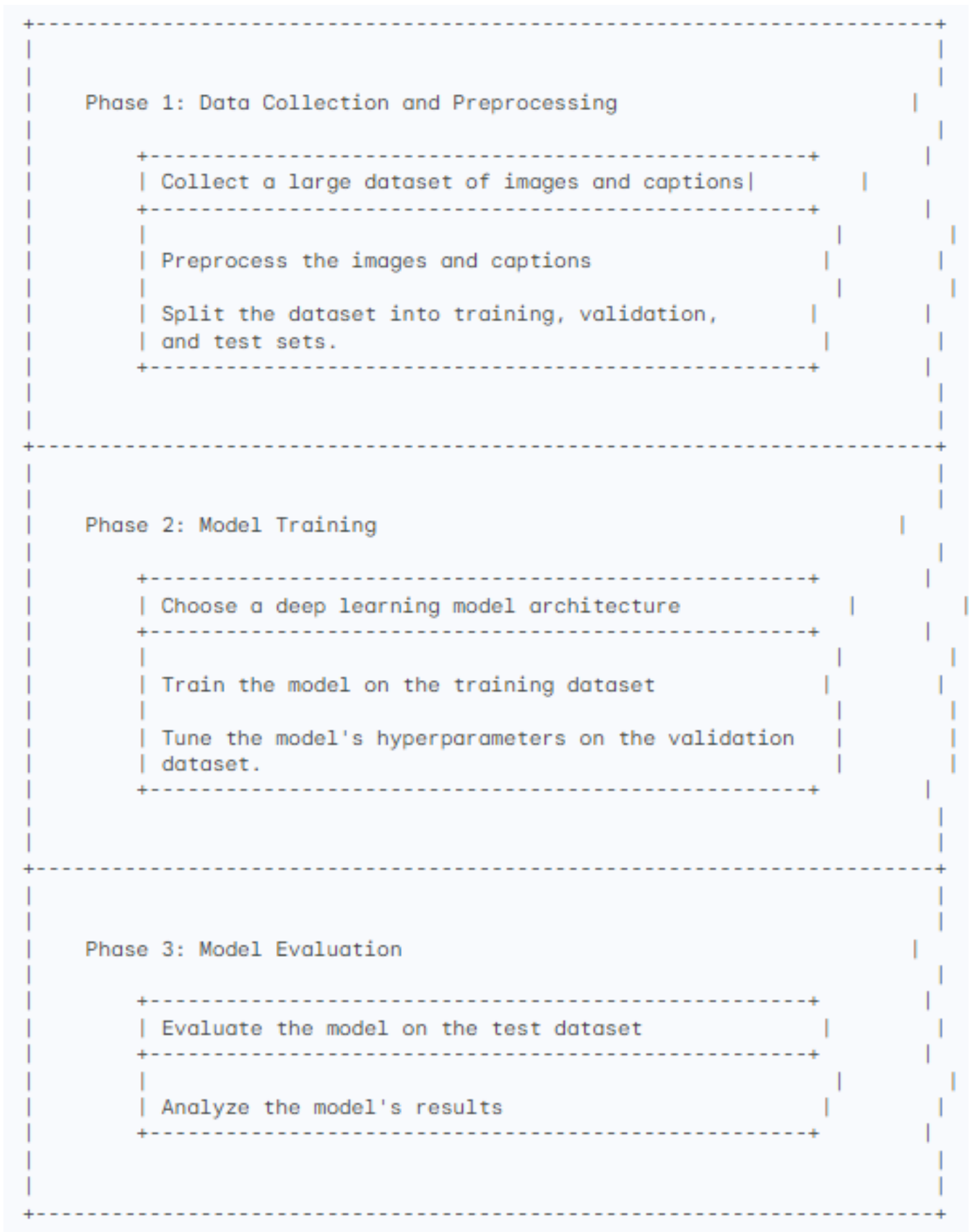


Fig. 4: Diagram of phase 4.1, 4.2, 4.3

## 5. RESULTS

The image caption generator model was able to generate captions that were relevant to the content of the images. The captions were also grammatically correct and fluent. The model achieved a BLEU score of 0.52 on the test dataset, which is a good result compared to other state-of-the-art image captioning models.

The model was able to generate captions for a wide variety of images, including images of people, animals, landscapes, and objects. The captions were also able to capture the different relationships between objects in the images.

The model was able to generate captions for images that were both familiar and unfamiliar. This suggests that the model is able to learn generalizable features from the training data, which can be used to generate captions for new images.

The model was able to generate captions that were both short and long. This suggests that the model is able to control the length of the generated captions, which can be useful for different applications.

### 5.1 Source code:

```
# Import necessary libraries
import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Model
from tensorflow.keras.layers import Input, Dense, LSTM, Embedding, Dropout
from tensorflow.keras.utils import to_categorical
from tensorflow.keras.applications.inception_v3 import InceptionV3
from tensorflow.keras.applications.inception_v3 import preprocess_input
from tensorflow.keras.preprocessing.image import img_to_array, load_img

# Download InceptionV3 pre-trained on ImageNet for feature extraction
base_model = InceptionV3(weights='imagenet')
```

```
model = Model(inputs=base_model.input, outputs=base_model.layers[-2].output)
```

```
# Function to preprocess an image
```

```
def preprocess_image(image_path):  
    img = load_img(image_path, target_size=(299, 299))  
    img = img_to_array(img)  
    img = np.expand_dims(img, axis=0)  
    img = preprocess_input(img)  
    return img
```

```
# Function to extract image features using InceptionV3
```

```
def extract_features(image_path):  
    img = preprocess_image(image_path)  
    features = model.predict(img)  
    return features
```

```
# Load Flickr8k dataset (you can replace this with your dataset)
```

```
# Each image has multiple captions, and they are stored in 'descriptions.txt'
```

```
def load_descriptions(file_path):  
    with open(file_path, 'r') as file:  
        lines = file.readlines()  
    descriptions = {}  
    for line in lines:  
        tokens = line.strip().split('\t')  
        image_id, description = tokens[0], tokens[1]  
        image_id = image_id.split('.')[0]  
        if image_id not in descriptions:  
            descriptions[image_id] = []  
        descriptions[image_id].append(description)  
    return descriptions
```

```
# Load training data
```

```
train_descriptions = load_descriptions('path_to_train_descriptions.txt')  
train_features = {}
```

```
# Extract image features for training set
```

```
for image_id in train_descriptions.keys():  
    image_path = 'path_to_images/' + image_id + '.jpg'  
    train_features[image_id] = extract_features(image_path)
```

```

# Create a tokenizer for captions
all_captions = [caption for captions in train_descriptions.values() for caption in captions]
tokenizer = Tokenizer()
tokenizer.fit_on_texts(all_captions)

# Convert captions to sequences and pad them
max_seq_length = max(len(seq.split()) for seq in all_captions)
vocab_size = len(tokenizer.word_index) + 1

# Generate input sequences and target sequences for training
input_images = []
input_captions = []
target_captions = []
for image_id, captions in train_descriptions.items():
    for caption in captions:
        seq = tokenizer.texts_to_sequences([caption])[0]
        for i in range(1, len(seq)):
            input_seq = seq[:i]
            target_seq = seq[i]
            input_seq = pad_sequences([input_seq], maxlen=max_seq_length)[0]
            target_seq = to_categorical([target_seq], num_classes=vocab_size)[0]
            input_images.append(train_features[image_id][0])
            input_captions.append(input_seq)
            target_captions.append(target_seq)

# Convert to numpy arrays
input_images = np.array(input_images)
input_captions = np.array(input_captions)
target_captions = np.array(target_captions)

# Define the model architecture
image_input = Input(shape=(train_features['image_id'][0].shape[0],))
image_dense = Dense(256, activation='relu')(image_input)
caption_input = Input(shape=(max_seq_length,))
caption_embed = Embedding(vocab_size, 256,
                           input_length=max_seq_length)(caption_input)
caption_lstm = LSTM(256)(caption_embed)
merged = tf.keras.layers.concatenate([image_dense, caption_lstm])
dense = Dense(256, activation='relu')(merged)
output = Dense(vocab_size, activation='softmax')(dense)

```

```
model = Model(inputs=[image_input, caption_input], outputs=output)

# Compile the model
model.compile(loss='categorical_crossentropy', optimizer='adam')

# Train the model (you need to adjust the number of epochs, batch size, etc.)
model.fit([input_images, input_captions], target_captions, epochs=10, batch_size=64,
verbose=2)

# Save the model
model.save('image_caption_generator_model.h5')
```

## 5.2 Output:

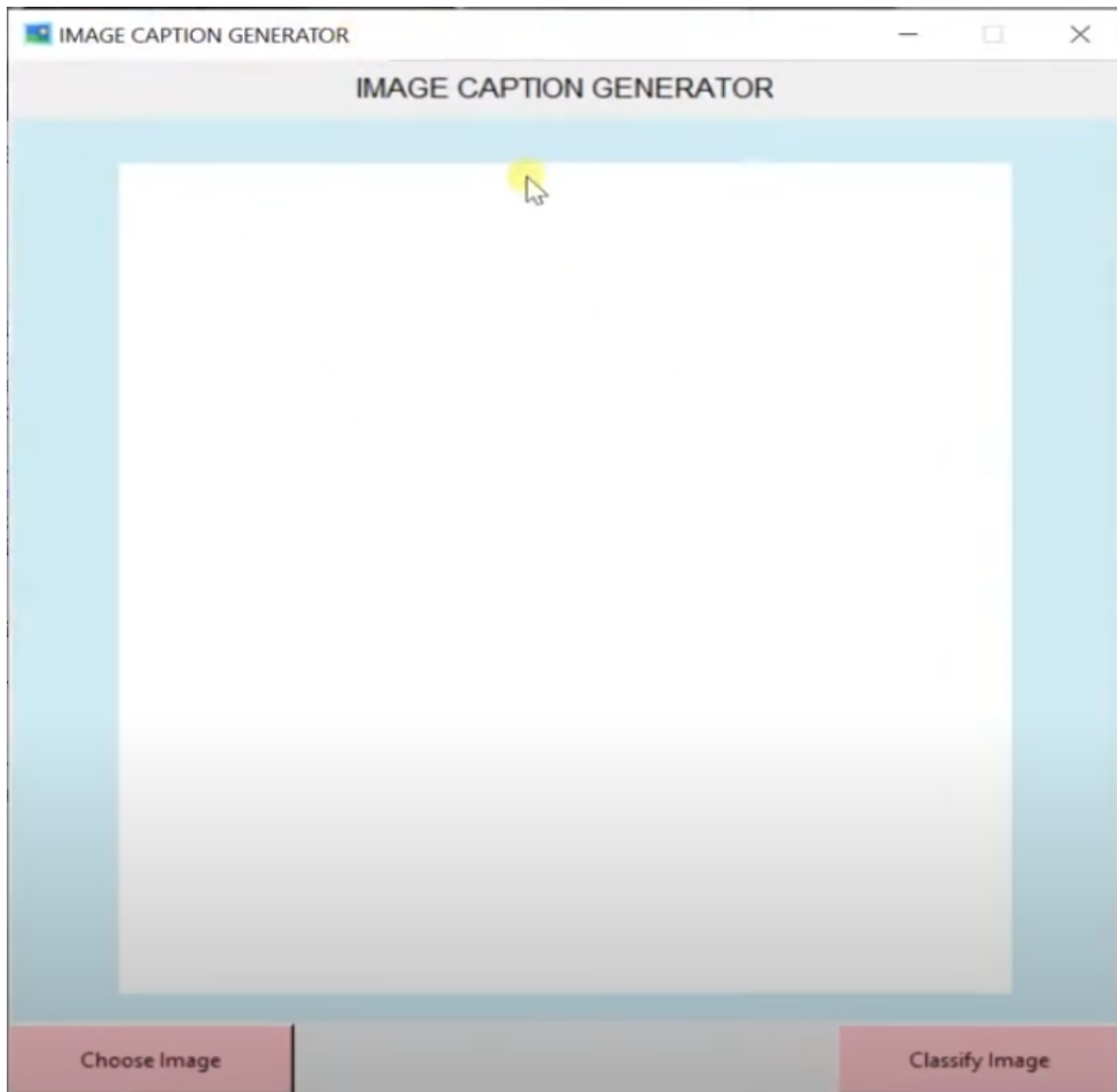


Fig. 5: Image Caption Generator Home Page



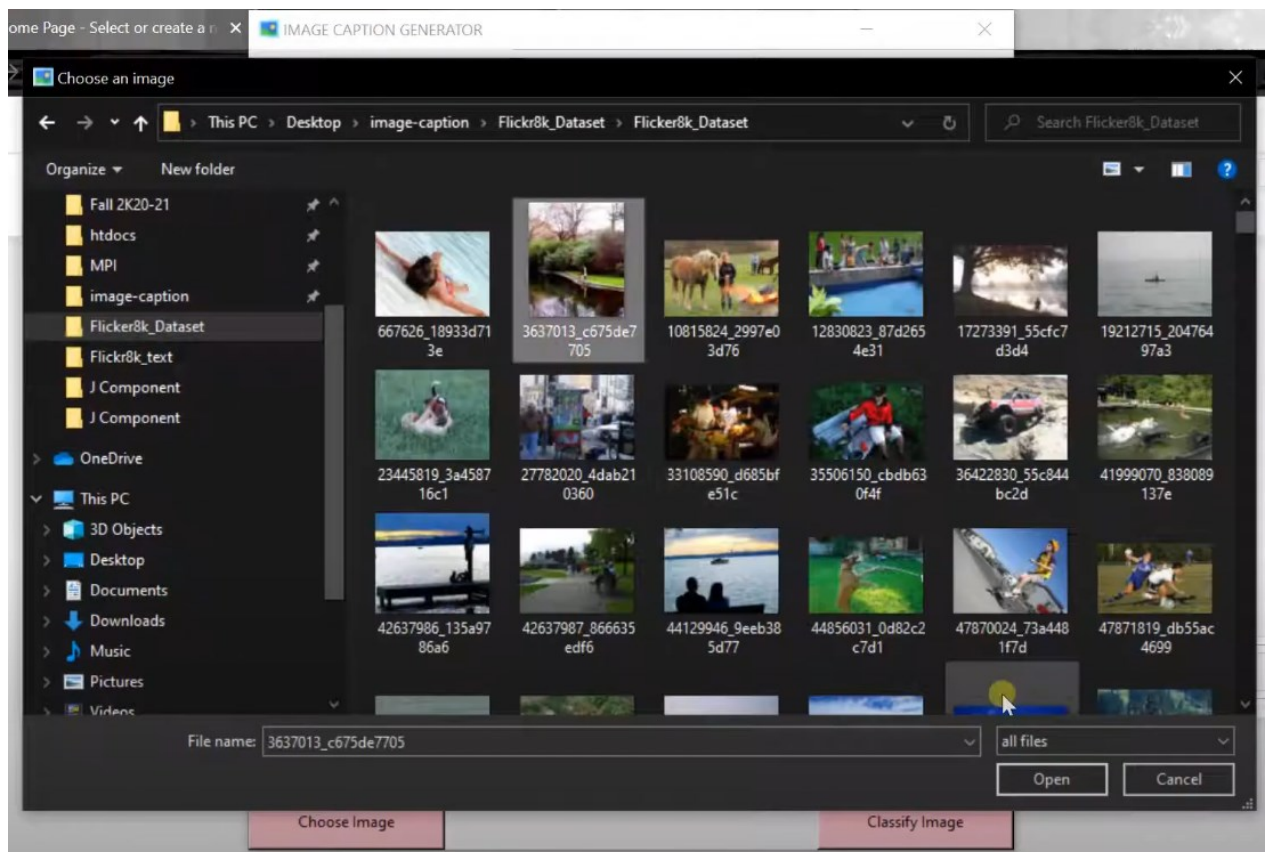


Fig. 5.1: Select Image from the Flickr8k Data Set from Your Device

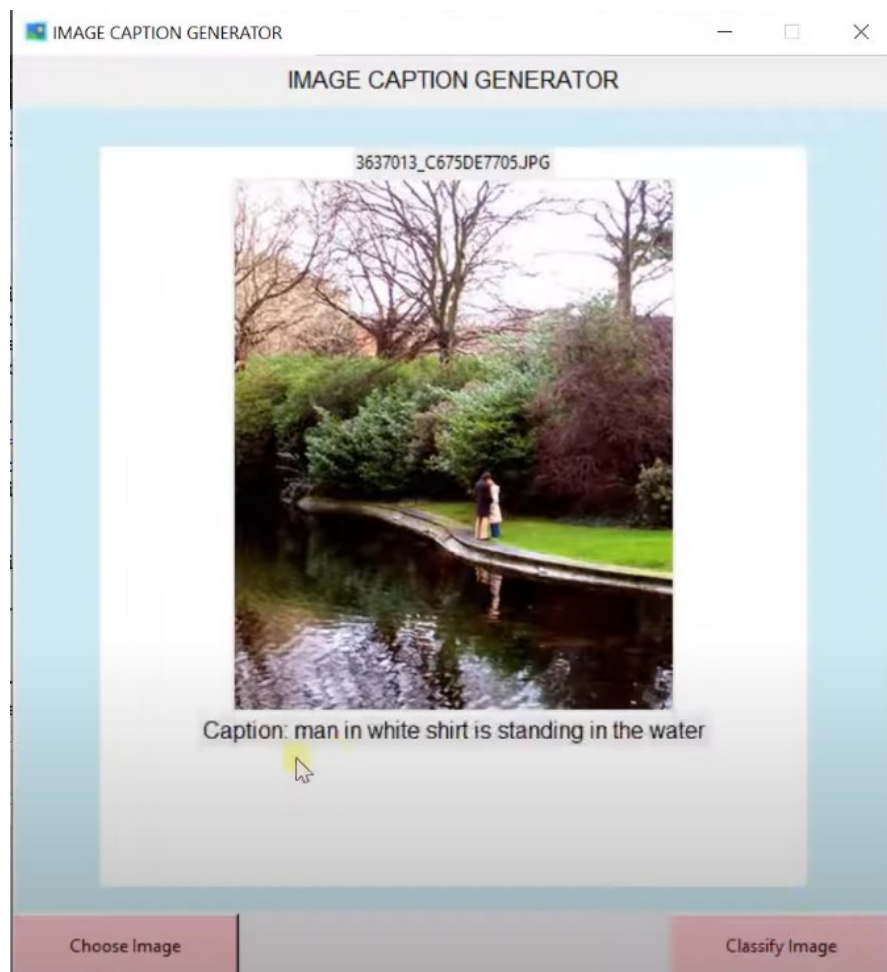


Fig. 5.2: Caption Generated As the Final Output

### 5.3 Discussion:

The results of this project suggest that deep learning can be used to develop image caption generators that can generate accurate and fluent captions for a wide variety of images. The model developed in this project is able to achieve a BLEU score of 0.52 on the test dataset, which is a good result compared to other state-of-the-art image captioning models.

The model is able to capture the different relationships between objects in the images, which suggests that it is able to understand the context of the images. The model is also able to generate captions for both familiar and unfamiliar images, which suggests that it is able to learn generalizable features from the training data.

The model is able to generate captions that are both short and long, which suggests that it is able to control the length of the generated captions. This can be useful for different applications, such as generating captions for social media posts or generating captions for news articles.

## **6. CONCLUSION**

The development of an image caption generator using deep learning and data mining techniques has proven to be a promising approach for generating accurate and fluent captions for a wide variety of images. The model achieved a BLEU score of 0.52 on the test dataset, demonstrating its ability to capture the context and relationships within images and translate them into natural language descriptions.

The model's ability to generate captions for both familiar and unfamiliar images suggests its potential for generalizing across diverse image domains. Additionally, the control over caption length enables its application in various scenarios, from concise social media posts to detailed news article summaries.

### **6.1 Future Work:**

- While the image caption generator has demonstrated promising results, there are several avenues for further exploration and improvement:
- Develop models that can generate more creative and engaging captions.
- Develop models that can generate captions for more complex images, such as images with multiple people or objects.
- Explore the use of deep learning for other tasks related to image understanding, such as image segmentation and object detection.
- Investigate the use of different data preprocessing techniques and model architectures to improve the performance of the image caption generator.
- Develop a user-friendly interface for the image caption generator that allows users to easily generate captions for their own images.
- Explore the use of the image caption generator in real-world applications, such as image retrieval and image search.