



INSTITUTE FOR ADVANCED
COMPUTING AND
SOFTWARE
DEVELOPMENT
AKURDI, PUNE

Documentation On
“Image Caption Generation using Deep Learning”
PG-DBDA SEP 2022

Submitted By:

Group No: 03

Ajit Yadav

Roll no. 229302

Mr. Rohit Puranik
Centre Coordinator

Dr. Shantanu Pathak
Project Guide

Contents:

1. Introduction.....	1
1.1 PROBLEM STATEMENT	1
1.2 Abstract	1
1.3 Use case	2
2. Overall Description.....	3
2.1 Workflow of Project:.....	3
2.2 Data Preprocessing and Cleaning.....	3
2.3 VGG16.....	5
2.3.1 Import Modules.....	7
2.3.2 Extract Image Features.....	9
2.3.3 Load the Captions Data.....	11
2.3.4 Preprocess Text Data	12
2.3.5 Train Test Split.....	14
2.3.6 Model Creation	15
2.3.7 Visualize the Results.....	19
2.3.8 Test with Real Image.....	21
3. Flow of Project.....	24
4. Applications and Future scope	25
5. Conclusion	26
6. References	27

1. INTRODUCTION

1.1 Problem Statement:

To develop a system for users, that can automatically generate a textual description of an image. This involves teaching a machine learning model to understand the visual content of an image and use that understanding to produce a coherent sentence that describes the objects, people, and actions depicted in the image.

1.2 Abstract:

Image captioning is an interesting and challenging task with applications in diverse domains such as image retrieval, organizing and locating images of users' interest etc. It has huge potential for replacing manual caption generation for images and is especially suitable for large scale image data. Recently, deep neural network based methods have achieved great success in the field of computer vision, machine translation and language generation. In this project, we propose an encoder-decoder based model that is capable of generating grammatically correct captions for images. This model makes use of VGG16 (Pre-Trained Model) as encoder and LSTM as decoder. To ensure the complete ground truth accuracy, the model is trained on the labelled Flickr8k dataset.

1.3 Use Case:

An image caption generator has several use cases across various fields.

Social Media: Social media platforms such as Facebook, Instagram, and Twitter can use image caption generators to help users with visual impairments access the content shared on the platform. The image caption generator can automatically generate a description of the images that are shared on the platform.

E-commerce: E-commerce websites can use image caption generators to improve the search ability and accessibility of their product catalogs. The generator can automatically generate descriptions of the products and help users find what they are looking for.

Healthcare: Medical imaging such as X-rays, CT scans, and MRIs can be described using an image caption generator. This can help doctors and healthcare professionals better understand and diagnose medical conditions.

Entertainment: Image caption generators can be used to enhance the user experience of online media platforms such as YouTube and Netflix. The generator can automatically generate captions for the videos and make them more accessible to a wider audience.

Education: Image caption generators can be used in educational settings to help students with visual impairments access visual content such as charts, graphs, and diagrams. The generator can automatically generate descriptions of the visual content and help students understand the material.

Autonomous Vehicles: Image caption generators can be used in autonomous vehicles to help them better understand and navigate their surroundings. The generator can automatically generate descriptions of the objects and scenes that the vehicle encounters and help it make better decisions.

2. Overall Description

2.1 Workflow of Project:

The diagram below shows the workflow of this project.

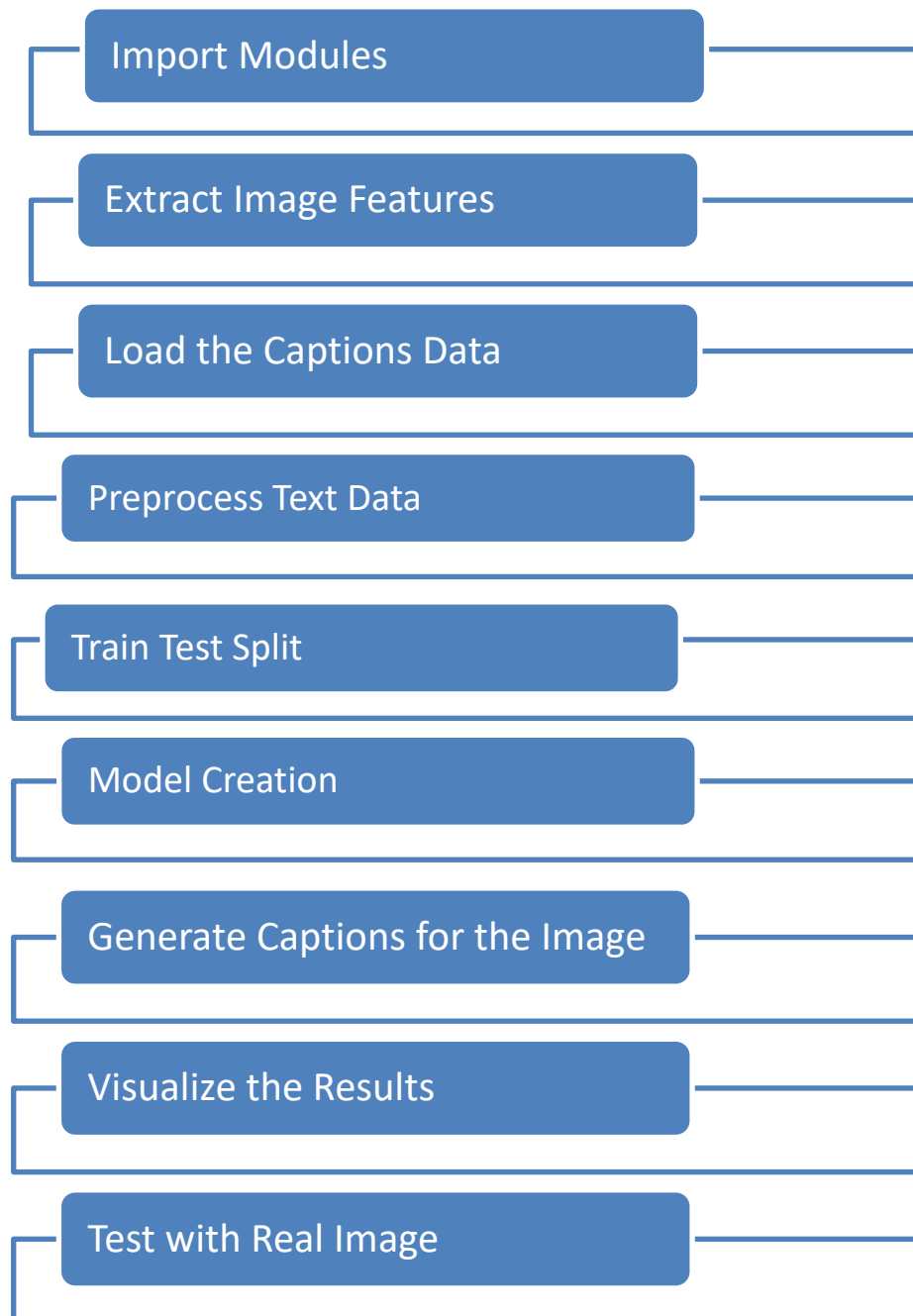


Figure 2.1 Workflow Diagram

2.2 About Datasets:

The Flickr8K dataset is a publicly available dataset consisting of 8,091 images that have been annotated with textual descriptions. The dataset was created by gathering images from the photo-sharing website Flickr and then annotating each image with five different captions. The captions were written by different people to capture a range of perspectives and interpretations of the same image. The dataset contains a wide variety of images, including landscapes, animals, people, and objects. The images have been resized to a resolution of 500 pixels on the longest side, and each image is accompanied by five captions in plain text format.

The Flickr8K dataset has been widely used for research in computer vision and natural language processing, particularly for tasks such as image captioning and multimodal machine learning. Text file contains 5 different caption for each image.

```

1 image,caption
2 1000268201_693b08cb0e.jpg,A child in a pink dress is climbing up a set of stairs in an entry way .
3 1000268201_693b08cb0e.jpg,A girl going into a wooden building .
4 1000268201_693b08cb0e.jpg,A little girl climbing into a wooden playhouse .
5 1000268201_693b08cb0e.jpg,A little girl climbing the stairs to her playhouse .
6 1000268201_693b08cb0e.jpg,A little girl in a pink dress going into a wooden cabin .
7 1001773457_577c3a7d70.jpg,A black dog and a spotted dog are fighting
8 1001773457_577c3a7d70.jpg,A black dog and a tri-colored dog playing with each other on the road .
9 1001773457_577c3a7d70.jpg,A black dog and a white dog with brown spots are staring at each other in the street .
10 1001773457_577c3a7d70.jpg,Two dogs of different breeds looking at each other on the road .
11 1001773457_577c3a7d70.jpg,Two dogs on pavement moving toward each other .
12 1002674143_1b742ab4b8.jpg,A little girl covered in paint sits in front of a painted rainbow with her hands in a bowl .
13 1002674143_1b742ab4b8.jpg,A little girl is sitting in front of a large painted rainbow .
14 1002674143_1b742ab4b8.jpg,A small girl in the grass plays with fingerpaints in front of a white canvas with a rainbow on it .
15 1002674143_1b742ab4b8.jpg,There is a girl with pigtails sitting in front of a rainbow painting .
16 1002674143_1b742ab4b8.jpg,Young girl with pigtails painting outside in the grass .

```

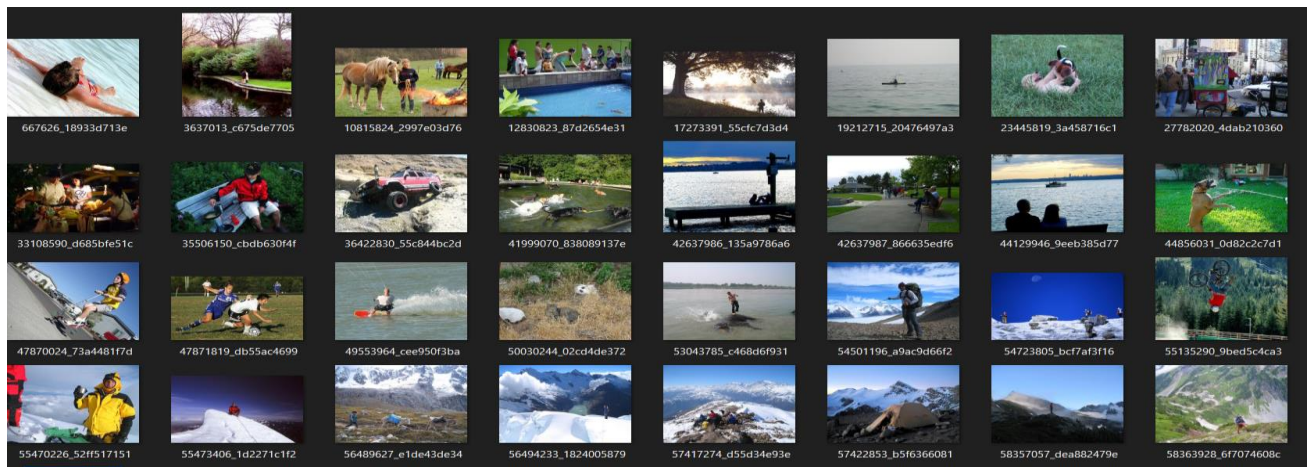


Figure 2.2 About Dataset

2.2.1 VGG16 Model:

The VGG16 model is a convolutional neural network (CNN) architecture that was proposed by the Visual Geometry Group (VGG) at the University of Oxford in 2014. It achieved state-of-the-art performance on the ImageNet Large Scale Visual Recognition Challenge (ILSVRC) in 2014, and has since been widely used as a pre-trained model for various computer vision tasks.

The VGG16 architecture consists of 16 layers, including 13 convolutional layers and 3 fully connected layers. It has a fixed input size of 224x224 RGB images and produces a 1000-dimensional output vector that represents the probabilities of the input image belonging to each of the 1000 classes in the ImageNet dataset.

The key innovation of the VGG16 architecture is the use of small 3x3 convolutional filters throughout the network, which allows the model to learn more complex and non-linear features while keeping the number of parameters manageable. VGG16 also uses max pooling and dropout layers to prevent overfitting.

Because of its strong performance on the ILSVRC, the pre-trained VGG16 model has been

used as a feature extractor for various computer vision tasks, such as image classification, object detection, and image segmentation. The model is available in the TensorFlow and Keras libraries, and can be easily fine-tuned on a new dataset for a specific task.

VGG16 Architecture:

The 16 in VGG16 refers to 16 layers that have weights. In VGG16 there are thirteen convolutional layers, five Max Pooling layers, and three Dense layers which sum up to 21 layers but it has only sixteen weight layers i.e., learnable parameters layer. VGG16 takes input tensor size as 224, 244 with 3 RGB channel.

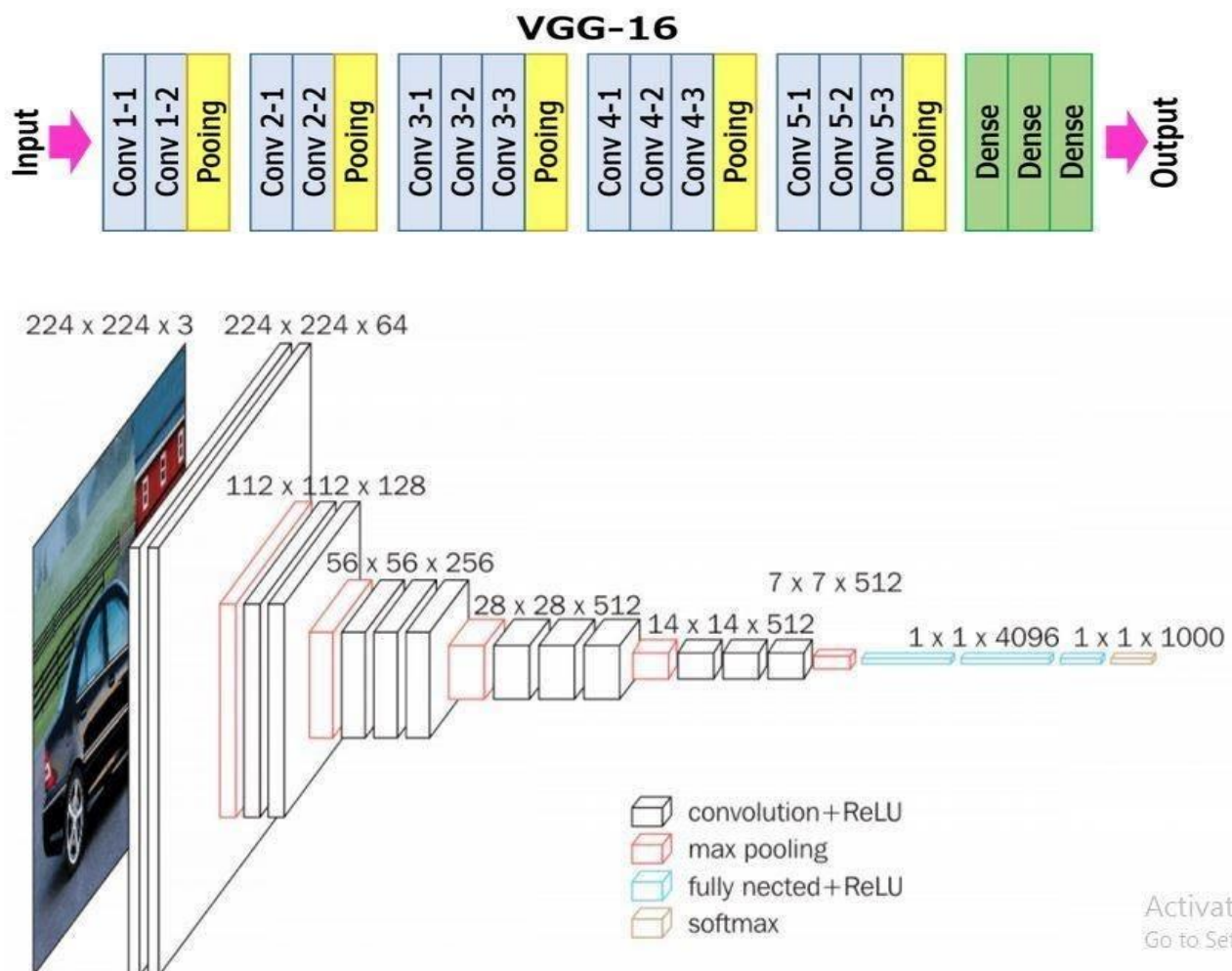


Figure 2.3 Architecture of VGG16 Model

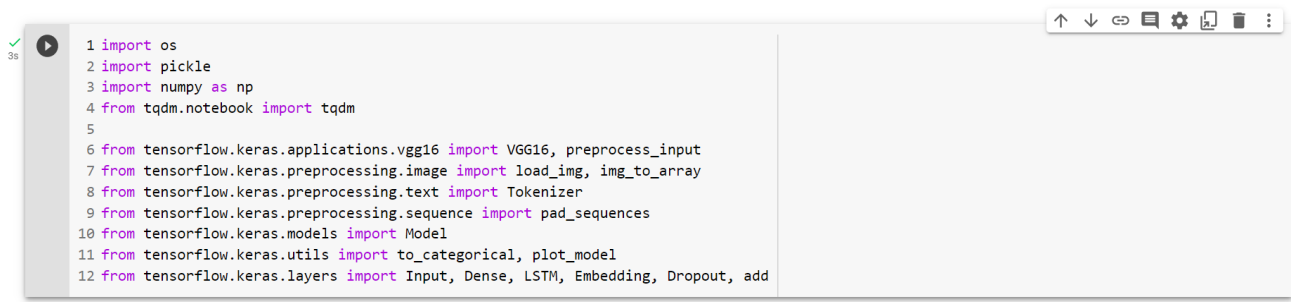
Most unique thing about VGG16 is that instead of having a large number of hyper-parameters they focused on having convolution layers of 3x3 filter with stride 1 and always used the same padding and maxpool layer of 2x2 filter of stride 2.

The convolution and max pool layers are consistently arranged throughout the whole architecture. Conv-1 Layer has 64 number of filters, Conv-2 has 128 filters, Conv-3 has 256 filters, Conv 4 and Conv 5 has 512 filters.

Three Fully-Connected (FC) layers follow a stack of convolutional layers: the first two have 4096 channels each, the third performs 1000-way ILSVRC classification and thus contains 1000 channels (one for each class). The final layer is the soft-max layer.

2.3.1 Import Modules:

▼ Import Modules



```

1 import os
2 import pickle
3 import numpy as np
4 from tqdm.notebook import tqdm
5
6 from tensorflow.keras.applications.vgg16 import VGG16, preprocess_input
7 from tensorflow.keras.preprocessing.image import load_img, img_to_array
8 from tensorflow.keras.preprocessing.text import Tokenizer
9 from tensorflow.keras.preprocessing.sequence import pad_sequences
10 from tensorflow.keras.models import Model
11 from tensorflow.keras.utils import to_categorical, plot_model
12 from tensorflow.keras.layers import Input, Dense, LSTM, Embedding, Dropout, add

```

OS Module:

OS provides a way of using operating system dependent functionality like reading or writing to the file system, working with environment variables, and executing system commands.

`os.path.join(path1, path2)`: Joins two paths together to form a complete path

Pickle Module:

The pickle module is used for object serialization and deserialization. It allows you to convert a Python object hierarchy into a byte stream that can be stored or transmitted, and then reconstruct the object hierarchy from the byte stream. This process is known as "pickling" and "unpickling".

The pickle module is commonly used for tasks such as:

Saving and loading machine learning models, so that they can be reused later
 Caching expensive computations, so that they can be reused instead of recalculated every time
 Saving and loading program state, so that a program can be resumed from where it left off after a crash or restart. The pickle module provides two main functions: **dump** and **load**.

```
with open('/content/drive/MyDrive/Colab Notebooks/Project/working/features.pkl', 'rb') as f:
    features = pickle.load(f)
```

```
pickle.dump(features, open(os.path.join(WORKING_DIR, '/content/drive/MyDrive/Colab Notebooks/Project/working/features_ver3.pkl'), 'wb'))
```

Numpy module:

NumPy is a powerful library for numerical computing in Python and is widely used for tasks such as scientific computing, data analysis, and machine learning.

Tensorflow.keras module:

TensorFlow is a popular open-source library for numerical computation and machine learning, which provides a flexible platform for building and deploying machine learning models. TensorFlow includes a wide range of tools and libraries for building and training deep learning models, including the Keras API. Keras is a high-level deep learning API written in Python, which provides a user-friendly interface for building, training, and evaluating deep learning models. Once VGG16 and the preprocess_input function are imported, you can use them to classify images or fine-tune the pre-trained model on a new

dataset.

img_to_array, load_img :

img_to_array is used to import two functions from the Keras preprocessing module in TensorFlow: load_img and img_to_array.

load_img is a function that loads an image file from disk and returns a PIL (Python Imaging Library) image object. It takes two arguments: path is the path to the image file, and target_size is a tuple specifying the size to which the image should be resized. If target_size is not specified, the function returns the original size of the image.

img_to_array is a function that converts a PIL image object into a NumPy array. It takes a single argument: img is the PIL image object to be converted.

These two functions are often used together to load and pre-process image data for deep learning models.

2.3.2 Extract Image Features:

In image caption generation, one common approach for extracting image features is to use a pre-trained convolutional neural network (CNN) such as VGG16. These CNNs are trained on large datasets (such as ImageNet) to classify images into different categories.

To extract image features, the pre-trained CNN is typically used as a feature extractor. The input image is passed through the CNN, and the output of one of the intermediate layers (before the fully connected layers) is extracted as the image feature representation. This feature representation can then be used as input to a separate natural language processing (NLP) model, such as a recurrent neural network (RNN), to generate captions

▼ Extract Image Features

34s

+ Code

+ Text

↑

↓

↶

💬

⚙️

📄

🗑️

⋮

```

1 # load vgg16 model
2 model = VGG16()
3 # restructure the model
4 model = Model(inputs=model.inputs, outputs=model.layers[-2].output)
5 # summarize
6 print(model.summary())

```

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_weights_tf_dim_ordering_tf_kernels.h5
553467096/553467096 [=====] - 28s 0us/step
Model: "model"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(None, 224, 224, 3)]	0
block1_conv1 (Conv2D)	(None, 224, 224, 64)	1792
block1_conv2 (Conv2D)	(None, 224, 224, 64)	36928
block1_pool (MaxPooling2D)	(None, 112, 112, 64)	0
block2_conv1 (Conv2D)	(None, 112, 112, 128)	73856
block2_conv2 (Conv2D)	(None, 112, 112, 128)	147584
block2_pool (MaxPooling2D)	(None, 56, 56, 128)	0
block3_conv1 (Conv2D)	(None, 56, 56, 256)	295168
block3_conv2 (Conv2D)	(None, 56, 56, 256)	590080
block3_conv3 (Conv2D)	(None, 56, 56, 256)	590080
block3_pool (MaxPooling2D)	(None, 28, 28, 256)	0
block4_conv1 (Conv2D)	(None, 28, 28, 512)	1180160
block4_conv2 (Conv2D)	(None, 28, 28, 512)	2359808
block4_conv3 (Conv2D)	(None, 28, 28, 512)	2359808
block4_pool (MaxPooling2D)	(None, 14, 14, 512)	0
block5_conv1 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv2 (Conv2D)	(None, 14, 14, 512)	2359808
block5_conv3 (Conv2D)	(None, 14, 14, 512)	2359808
block5_pool (MaxPooling2D)	(None, 7, 7, 512)	0
flatten (Flatten)	(None, 25088)	0
fc1 (Dense)	(None, 4096)	102764544
fc2 (Dense)	(None, 4096)	16781312

=====
Total params: 134,260,544
Trainable params: 134,260,544
Non-trainable params: 0

10

IACSD-PG-DBDA-SEP-22

```

1 # extract features from image
2 features = {}
3 directory = os.path.join(BASE_DIR, 'Images')
4
5 for img_name in tqdm(os.listdir(directory)):
6     # load the image from file
7     img_path = directory + '/' + img_name
8     image = load_img(img_path, target_size=(224, 224))
9     # convert image pixels to numpy array
10    image = img_to_array(image)
11    # reshape data for model
12    image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2])) # rgb image
13    # preprocess image for vgg
14    image = preprocess_input(image)
15    # extract features
16    feature = model.predict(image, verbose=0)
17    # get image ID
18    image_id = img_name.split('.')[0]
19    # store feature
20    features[image_id] = feature

```

2.3.3 Load the Captions Data

We read the contents of a file named 'captions.txt' located in the directory specified by the variable `BASE_DIR`. It then skips the first line of the file using the `next(f)` function call, and reads the rest of the file into a string variable named `captions_doc`. And we are Mapping each caption to its image ID.

▼ Load the Captions Data

```

[10] 1 with open(os.path.join(BASE_DIR, 'captions.txt'), 'r') as f:
2     next(f)
3     captions_doc = f.read()

```

```

1 captions_doc

```

```

'1000268201_693b08cb0e.jpg,A child in a pink dress is climbing up a set of stairs in an entry way .\n1000268201_693b08cb0e.jpg,A girl going into a wooden building .\n1000268201_693b08cb0e.jpg,A little girl climbing into a wooden playhouse .\n1000268201_693b08cb0e.jpg,A little girl climbing the stairs to her playhouse .\n1000268201_693b08cb0e.jpg,A little girl in a pink dress going into a wooden cabin .\n1001773457_577c3a7d70.jpg,A black dog and a spotted dog are fighting\n1001773457_577c3a7d70.jpg,A black dog and a tri-colored dog playing with each other on the road .\n1001773457_577c3a7d70.jpg,A black dog and a white dog with brown spots are staring at each other in the street .\n1001773457_577c3a7d70.jpg,Two dogs of different breeds looking at each other on the road .\n1001773457_577c3a7d70.jpg,Two dogs on pavement moving toward each other .\n1002674143_1b742ab4b8.jpg,A little girl covered in paint sits in front of a painted rainbow with her hands in a bowl .\n1002674143_1b742ab4b8.jp...

```

```

1 # create mapping of image to captions
2 mapping = {}
3 # process lines
4 for line in tqdm(captions_doc.split('\n')):
5     # split the line by comma(,)
6     tokens = line.split(',')
7     if len(line) < 2:
8         continue
9     image_id, caption = tokens[0], tokens[1:]
10    # remove extension from image ID
11    image_id = image_id.split('.')[0]
12    # convert caption list to string
13    caption = " ".join(caption)
14    # create list if needed
15    if image_id not in mapping:
16        mapping[image_id] = []
17    # store the caption
18    mapping[image_id].append(caption)

```

100% 40456/40456 [00:00<00:00, 259385.00it/s]

2.3.4 Preprocess Text Data:

We defined a function named `clean` that takes a single argument `mapping`, which is assumed to be a dictionary where each key corresponds to an image and the associated value is a list of captions describing that image. The function cleans and preprocesses each caption in the mapping by performing the following steps:

1. Converts each caption to lowercase using the `.lower()` method.
2. Removes any non-alphabetic characters (e.g. digits, special characters) from the caption using the `.replace()` method with a regular expression pattern `[\^w]`.
3. Removes any extra whitespace characters (e.g. multiple spaces in a row) from the caption using the `.replace()` method with the pattern `\s+`.
4. Adds special "start" and "end" tags to the caption to indicate the beginning and end of the caption using the `str.join()` method with a list comprehension.

The cleaned captions are stored back in the original mapping dictionary, overwriting the original values for each key.

▼ Preprocess Text Data

```

[13] 1 def clean(mapping):
      2     for key, captions in mapping.items():
      3         for i in range(len(captions)):
      4             # take one caption at a time
      5             caption = captions[i]
      6             # preprocessing steps
      7             # convert to lowercase
      8             caption = caption.lower()
      9             # delete digits, special chars, etc.,
      10            caption = caption.replace('[^A-Za-z]', '')
      11            # delete additional spaces
      12            caption = caption.replace('\s+', ' ')
      13            # add start and end tags to the caption
      14            caption = 'startseq ' + " ".join([word for word in caption.split() if len(word)>1]) + ' endseq'
      15            captions[i] = caption

```

We have displayed before and after Preprocessing of text data:

```

1 # before preprocess of text
2 mapping['1000268201_693b08cb0e']

['A child in a pink dress is climbing up a set of stairs in an entry way .',
'A girl going into a wooden building .',
'A little girl climbing into a wooden playhouse .',
'A little girl climbing the stairs to her playhouse .',
'A little girl in a pink dress going into a wooden cabin .']

[15] 1 # preprocess the text
      2 clean(mapping)

1 # after preprocess of text
2 mapping['1000268201_693b08cb0e']

['startseq child in pink dress is climbing up set of stairs in an entry way endseq',
'startseq girl going into wooden building endseq',
'startseq little girl climbing into wooden playhouse endseq',
'startseq little girl climbing the stairs to her playhouse endseq',
'startseq little girl in pink dress going into wooden cabin endseq']

```

Tokenize the text data:

We are creating Keras Tokenizer class by calling the constructor with no arguments. This creates a new Tokenizer object that can be used to convert sequences of text into sequences of integers suitable for use as input to a neural network.

The `fit_on_texts` method is then called on the Tokenizer object, passing in a list of all the captions in the dataset as its argument. This method updates the internal state of the Tokenizer object to create a vocabulary of all the unique words in the text, and assigns a

unique integer index to each word.

Finally, the `vocab_size` variable is set to the size of the vocabulary created by the `Tokenizer`, which is the total number of unique words plus one (to account for the special "out of vocabulary" token).

To prepare the vocabulary for the captions in the dataset, which will be used to represent each caption as a sequence of integers for input into a neural network. The `Tokenizer` object provides a convenient way to perform this transformation, and the resulting vocabulary size is used to set the size of the embedding layer in the neural network.

```
[ ] 1 # tokenize the text
2 tokenizer = Tokenizer()
3 tokenizer.fit_on_texts(all_captions) #fit_on_texts is used to produces the one-hot encoding for the original set of texts.
4 vocab_size = len(tokenizer.word_index) + 1
5
6 # Because Tokenizer.word_index is a python dictionary that contains token keys (string) and token ID values (integer),
7 # and where the first token ID is 1 (not zero) and where the token IDs are assigned incrementally. Therefore,
8 # the greatest token ID in word_index is len(word_index).
9 # Therefore, we need vocabulary of size len(word_index) + 1 to be able to index up to the greatest token ID.
```

```
[ ] 1 tokenizer.word_index
'him': 148,
'it': 149,
's': 150,
'road': 151,
'area': 152,
'that': 153,
'basketball': 154,
'tan': 155,
'back': 156,
'trick': 157,
'race': 158,
'swing': 159,
'head': 160,
'shorts': 161,
```

2.3.5 Train Test Split:

We are splitting the image ids into training and testing sets for use in a machine learning model. The split is based on a percentage of the total number of image ids, with 85% of the image ids assigned to the training set and 15% assigned to the testing set.

▼ Train Test Split

```

1 image_ids = list(mapping.keys())
2 split = int(len(image_ids) * 0.85)
3 train = image_ids[:split]
4 test = image_ids[split:]

```

data_generator():

We have created data generator function to get data in batch (avoids session crash).

```

1 # create data generator to get data in batch (avoids session crash)
2 def data_generator(data_keys, mapping, features, tokenizer, max_length, vocab_size, batch_size):
3     # loop over images
4     X1, X2, y = list(), list(), list()
5     n = 0
6     while 1:
7         for key in data_keys:
8             n += 1
9             captions = mapping[key]
10            # process each caption
11            for caption in captions:
12                # encode the sequence
13                seq = tokenizer.texts_to_sequences([caption])[0]
14
15            # split the sequence into X, y pairs
16            for i in range(1, len(seq)):
17                # split into input and output pairs
18                in_seq, out_seq = seq[:i], seq[i]
19                # pad input sequence
20                in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
21                # encode output sequence
22                out_seq = to_categorical([out_seq], num_classes=vocab_size)[0]
23
24            # store the sequences
25            X1.append(features[key][0])
26            X2.append(in_seq)
27            y.append(out_seq)
28
29            if n == batch_size:
30                X1, X2, y = np.array(X1), np.array(X2), np.array(y)
31                yield [X1, X2], y
32                X1, X2, y = list(), list(), list()
33                n = 0

```

2.3.6 Model Creation:

We have defined an image captioning model with an encoder-decoder architecture. The encoder takes in the image features as input and passes them through a series of layers, including dropout and a fully connected layer, to extract a feature vector. The decoder takes in the captions as input and passes them through an embedding layer and an LSTM layer

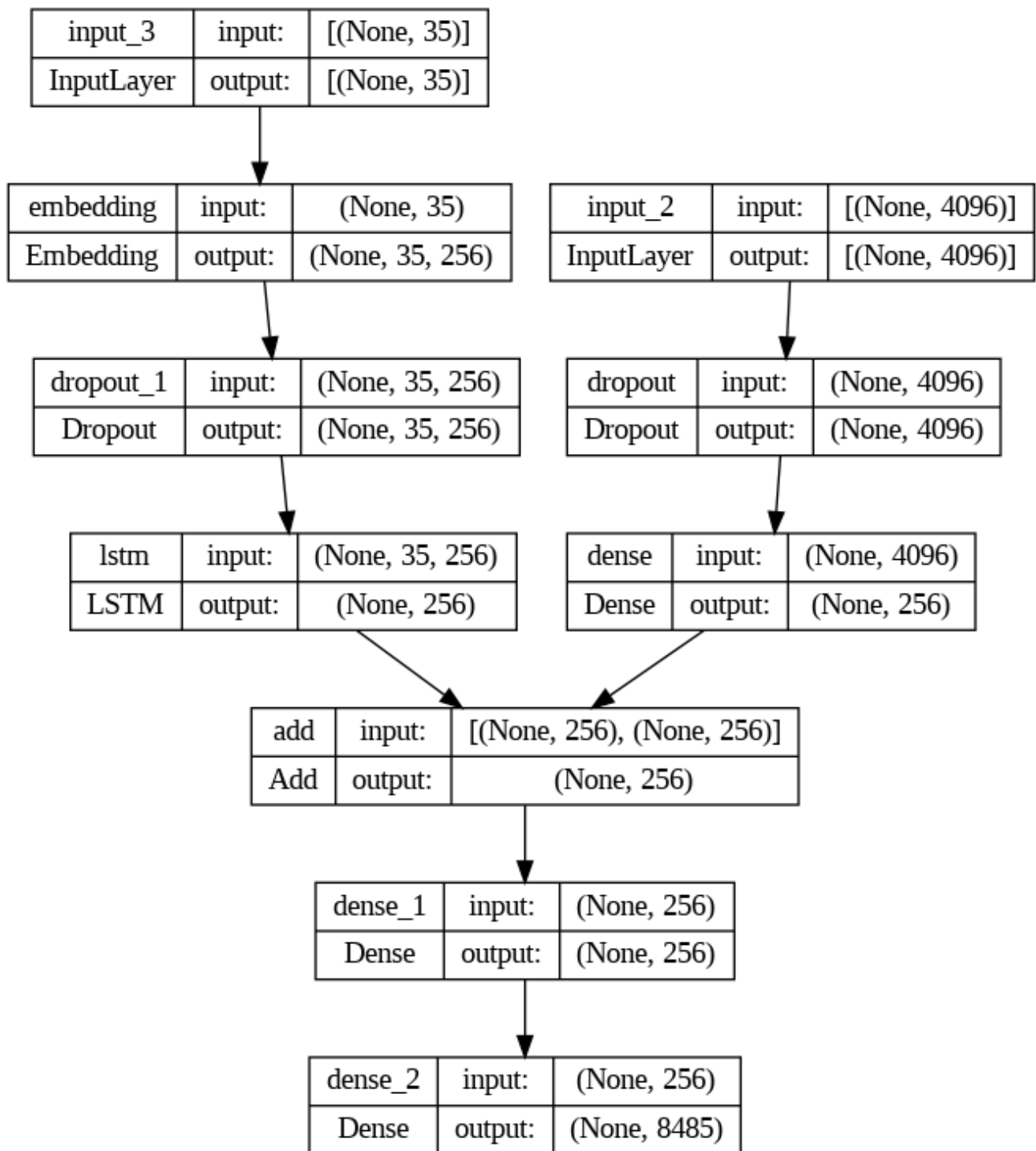
to extract a sequence of feature vectors. The decoder then combines the image feature vector with the sequence of feature vectors from the LSTM layer using an element-wise addition and passes the resulting vector through a series of fully connected layers to generate the final output. The model is compiled using the Adam optimizer and categorical cross-entropy loss function.

▼ Model Creation

```

1 # encoder model
2 # image feature layers
3 inputs1 = Input(shape=(4096,))
4 fe1 = Dropout(0.30)(inputs1)
5 fe2 = Dense(256, activation='relu')(fe1)
6 # sequence feature layers
7 inputs2 = Input(shape=(max_length,))
8 se1 = Embedding(vocab_size, 256, mask_zero=True)(inputs2)
9 se2 = Dropout(0.30)(se1)
10 se3 = LSTM(256)(se2)
11
12 # decoder model
13 decoder1 = add([fe2, se3])
14 decoder2 = Dense(256, activation='relu')(decoder1)
15 outputs = Dense(vocab_size, activation='softmax')(decoder2)
16
17 model = Model(inputs=[inputs1, inputs2], outputs=outputs)
18
19 # Unlike maintaining a single learning rate through training in SGD, Adam optimizer updates the learning rate for each network weight indivi
20 # Categorical cross entropy loss function is used to compute the quantity that the the model should seek to minimize during training
21 model.compile(loss='categorical_crossentropy', optimizer='adam')
22
23 # plot the model
24 plot_model(model, show_shapes=True)

```



```

1 # train the model
2 epochs = 25
3 batch_size = 30
4 steps = len(train) // batch_size
5
6 for i in range(epochs):
7     # create data generator
8     generator = data_generator(train, mapping, features, tokenizer, max_length, vocab_size, batch_size)
9     # fit for one epoch
10    model.fit(generator, epochs=1, steps_per_epoch=steps, verbose=1)

```

```

229/229 [=====] - 76s 331ms/step - loss: 3.7133
229/229 [=====] - 67s 292ms/step - loss: 3.4110
229/229 [=====] - 69s 299ms/step - loss: 3.1835
229/229 [=====] - 67s 293ms/step - loss: 3.0095
229/229 [=====] - 68s 298ms/step - loss: 2.8785
229/229 [=====] - 69s 301ms/step - loss: 2.7816
229/229 [=====] - 69s 301ms/step - loss: 2.6947
229/229 [=====] - 68s 298ms/step - loss: 2.6118
229/229 [=====] - 67s 294ms/step - loss: 2.5419
229/229 [=====] - 67s 293ms/step - loss: 2.4759
229/229 [=====] - 68s 296ms/step - loss: 2.4119
229/229 [=====] - 71s 309ms/step - loss: 2.3539
229/229 [=====] - 69s 302ms/step - loss: 2.3063
229/229 [=====] - 69s 302ms/step - loss: 2.2603
229/229 [=====] - 67s 291ms/step - loss: 2.2233
229/229 [=====] - 68s 297ms/step - loss: 2.1877
229/229 [=====] - 67s 294ms/step - loss: 2.1497
229/229 [=====] - 71s 308ms/step - loss: 2.1164
229/229 [=====] - 70s 306ms/step - loss: 2.0884

```

We defined the maximum number of epochs to 25 and the batch size to 30. We calculated the number of steps needed in each epoch for training and validation data.

And we use a data generator to generate arrays of these sequence of the size of the batch progressively. We used this approach to avoid reaching RAM and GPU limits. Otherwise, we weren't able to train the model with the RAM available in Google Collaboratory. The data generator receives the training data shuffled and works with image ids to take the images that correspond to the captions.

Rnn Model1: It has two inputs. In the text submodel it has an Embedding, Dropout and LSTM layers. In the image submodel, it has a Dropout and a Dense layer. Then it adds these two submodels and finally there are two Dense layers, the last one having vocabulary size with softmax to train the model with the RAM available in Google Colaboratoty. The data generator receives the training data shuffled and works with image ids to take the images that correspond to the captions.

2.3.7 Generate Captions for the Image:

The **idx_to_word** function takes an integer representing a word index and a tokenizer object as inputs, and returns the corresponding word. It iterates over the `word_index` dictionary of the tokenizer, which maps words to integer indices, and returns the word that has the given integer index. If the integer index is not found in the `word_index`, it returns `None`. This function can be used to convert the output of the model, which is a sequence of integer indices representing the predicted words, back into a sequence of words.

The function **predict_caption()** takes as input the trained model, an image, a tokenizer, and the maximum length of the sequence, and returns the predicted caption for the image.

The function starts by adding a special start tag `startseq` to the input sequence. It then iterates over the maximum length of the sequence, encoding the input sequence, padding it, and using the model to predict the next word. The predicted word is then converted to its corresponding string form using the `idx_to_word()` function. If the predicted word is not found, the loop is stopped. Otherwise, the predicted word is appended to the input sequence for generating the next word. If the predicted word is the end tag `endseq`, the loop is stopped.

▼ Generate Captions for the Image

```

1 def idx_to_word(integer, tokenizer):
2     for word, index in tokenizer.word_index.items():
3         if index == integer:
4             return word
5     return None

[ ] 1 # generate caption for an image
2 def predict_caption(model, image, tokenizer, max_length):
3     # add start tag for generation process
4     in_text = 'startseq'
5     # iterate over the max length of sequence
6     for i in range(max_length):
7         # encode input sequence
8         sequence = tokenizer.texts_to_sequences([in_text])[0]
9         # pad the sequence
10        sequence = pad_sequences([sequence], max_length)
11        # predict next word
12        yhat = model.predict([image, sequence], verbose=0)
13        # get index with high probability
14        yhat = np.argmax(yhat)
15        # convert index to word
16        word = idx_to_word(yhat, tokenizer)
17        # stop if word not found
18        if word is None:
19            break
20        # append word as input for generating next word
21        in_text += " " + word
22        # stop if we reach end tag
23        if word == 'endseq':
24            break
25
26    return in_text

```

The BLEU score is a metric used to evaluate the quality of generated text, such as machine translation or image captioning. It measures how similar the generated text is to a set of reference texts. The score ranges between 0 and 1, where a higher score indicates a better match with the reference texts. The BLEU score is computed using n-grams, which are sequences of n words. The score is a weighted average of the n-gram precision, which is the percentage of n-grams in the generated text that appear in the reference texts. The weights give more importance to longer n-grams.

In this code, we are computing the BLEU-1 and BLEU-2 scores using the `corpus_bleu` function from the `nltk.translate.bleu_score` module. The actual list contains the reference texts, which are the captions associated with the test images. The predicted list contains

the generated texts, which are the captions predicted by the model for the test images.

The BLEU-1 score is computed using only unigrams ($n=1$), while the BLEU-2 score is computed using both unigrams and bigrams ($n=2$). The weights used in the computation are (1,0,0,0) for BLEU-1 and (0.5,0.5,0,0) for BLEU-2, which give equal importance to unigrams and bigrams.

```
1 from nltk.translate.bleu_score import corpus_bleu
2 # validate with test data
3 actual, predicted = list(), list()
4
5 for key in tqdm(test):
6     # get actual caption
7     captions = mapping[key]
8     # predict the caption for image
9     y_pred = predict_caption(model, features[key], tokenizer, max_length)
10    # split into words
11    actual_captions = [caption.split() for caption in captions]
12    y_pred = y_pred.split()
13    # append to the list
14    actual.append(actual_captions)
15    predicted.append(y_pred)
16
17 # calculate BLEU score
18 print("BLEU-1: %f" % corpus_bleu(actual, predicted, weights=(1.0, 0, 0, 0)))
19 print("BLEU-2: %f" % corpus_bleu(actual, predicted, weights=(0.5, 0.5, 0, 0)))
```

100% 1214/1214 [13:40<00:00, 1.62it/s]
BLEU-1: 0.534889
BLEU-2: 0.313522

2.3.8 Visualize the Results:

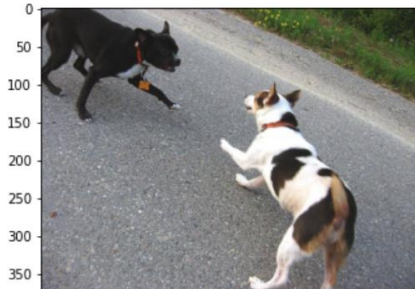
We are using generate_caption() function for caption generation , and we are displaying actual vs predicted caption .

▸ Visualize the Results

```
1 from PIL import Image
2 import matplotlib.pyplot as plt
3 def generate_caption(image_name):
4     # load the image
5     # image_name = "1001773457_577c3a7d70.jpg"
6     image_id = image_name.split('.')[0]
7     img_path = os.path.join(BASE_DIR, "Images", image_name)
8     image = Image.open(img_path)
9     captions = mapping[image_id]
10    print('-----Actual-----')
11    for caption in captions:
12        print(caption)
13    # predict the caption
14    y_pred = predict_caption(model, features[image_id], tokenizer, max_length)
15    print('-----Predicted-----')
16    print(y_pred)
17    plt.imshow(image)
```

```
1 generate_caption("1001773457_577c3a7d70.jpg")
```

```
-----Actual-----
startseq black dog and spotted dog are fighting endseq
startseq black dog and tri-colored dog playing with each other on the road endseq
startseq black dog and white dog with brown spots are staring at each other in the street endseq
startseq two dogs of different breeds looking at each other on the road endseq
startseq two dogs on pavement moving toward each other endseq
-----Predicted-----
startseq black dog and tri colored dog are playing with each other on the grass endseq
```



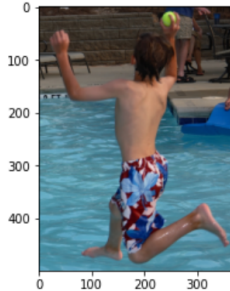
```
1 generate_caption("161669933_3e7d8c7e2c.jpg")
```

```
-----Actual-----
startseq competitive motorcycle racer prepares to make left turn along paved road endseq
startseq man rides motorcycle with the number on it endseq
startseq motorcycle driver swerves to the left endseq
startseq motorcycle with number five on front is being ridden by rider wearing red helmet endseq
startseq professional motorcycle racer turning corner endseq
-----Predicted-----
startseq motorcycle rider rides his bike on track endseq
```




```
1 generate_caption("524036004_6747cf909b.jpg")
```

```
-----Actual-----
startseq "a boy in blue and red trunks with green ball in his hand is jumping into pool ." endseq
startseq boy is jumping into the water with ball in his hand endseq
startseq boy jumps into pool while holding ball endseq
startseq boy in floral swim trunks jumps into the pool endseq
startseq boy in red and blue flowered trunks jumps into pool holding green ball endseq
-----Predicted-----
startseq boy in bathing suit is jumping into pool while another boy watches endseq
```



2.3.9 Test with Real Image:

Test with Real Image

```
[ ] 1 vgg_model = VGG16()
    2 # restructure the model
    3 vgg_model = Model(inputs=vgg_model.inputs, outputs=vgg_model.layers[-2].output)
```

```
[ ] 1 image_path = '/content/images2.jpeg'
    2 # load image
    3 image = load_img(image_path, target_size=(224, 224))
    4 # convert image pixels to numpy array
    5 image = img_to_array(image)
    6 # reshape data for model
    7 image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
    8 # preprocess image for vgg
    9 image = preprocess_input(image)
   10 # extract features
   11 feature = vgg_model.predict(image, verbose=0)
   12 # predict from the trained model
   13 predict_caption(model, feature, tokenizer, max_length)
```

```
'startseq black and white dog is running through the grass endseq'
```



161669933_3e7d8c7e2c



download



download2



images2



images3

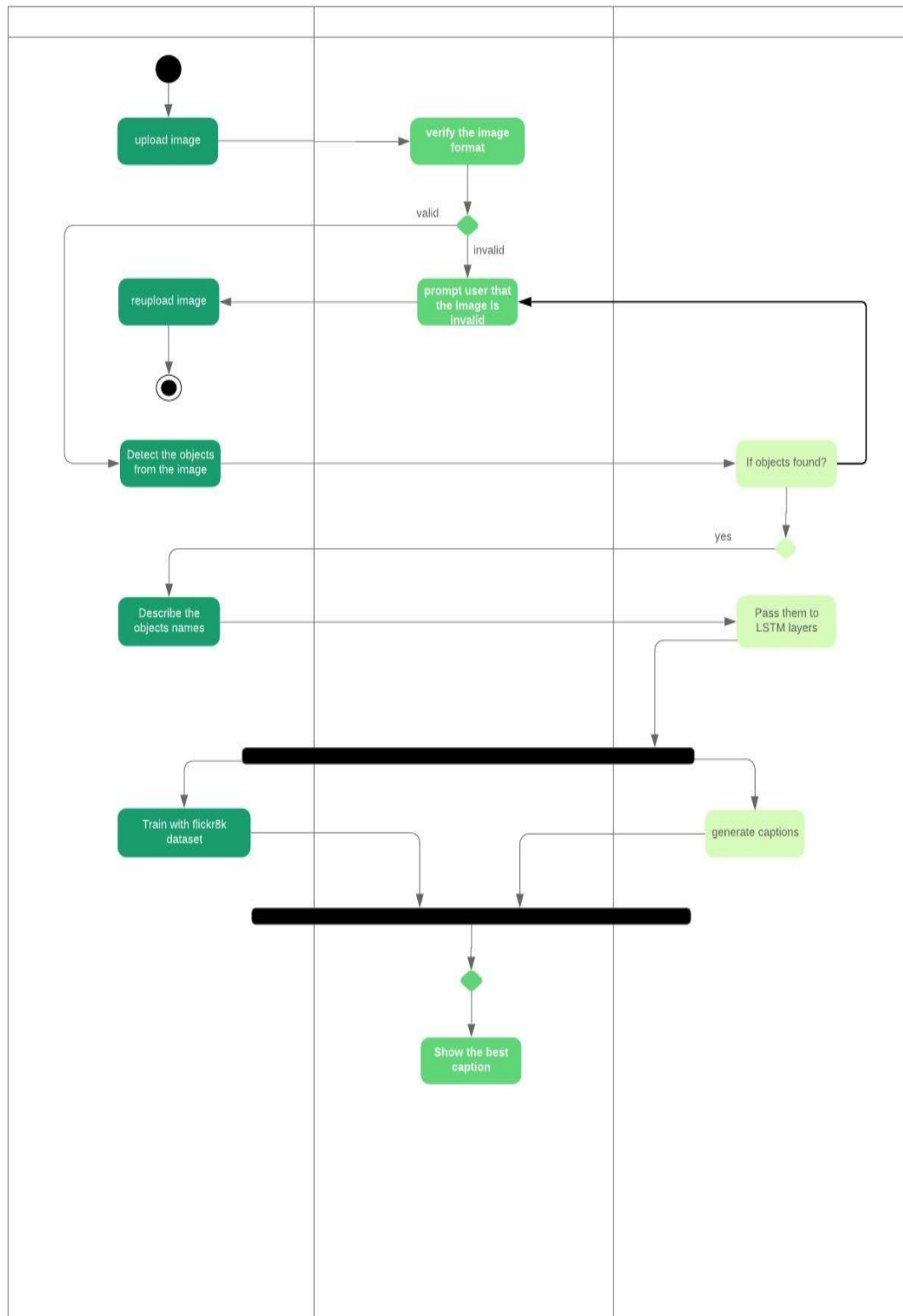


Fig : Flow of the project

Applications and Future scope of Image Caption Generator:

- It can be used to assist the blind using text-to-speech.
- It can be used to convert captions for images in social feed as well as messages to speech which will enhance social medial experience of users.
- It can also be used for educational purposes as young children can be assisted about recognition of objects and learning the English language.
- It is also helpful in field of robotics as environmental insights can be provided through natural language representation.
- It can also be used for image searches and indexing purposes on internet if images present on the internet have captions

Conclusion

In conclusion, an image caption generator is a powerful tool that uses deep learning algorithms to analyze an image and generate a descriptive and accurate caption for it. This technology has various applications, including improving accessibility for the visually impaired, enhancing social media engagement, and supporting image retrieval in large databases.

Image caption generators have been developed using various deep learning architectures, including Convolutional Neural Networks (CNNs) and Long Short Term Memory (LSTM), and have been trained on large datasets of images with corresponding captions.

While image caption generators have shown impressive results, there is still room for improvement in terms of the accuracy and diversity of generated captions. Researchers are continuing to work on developing more advanced models that can capture more nuanced and contextual information from images, as well as incorporating more knowledge about language and grammar.

Overall, image caption generators have the potential to revolutionize the way we interact with and understand visual content, and will likely continue to evolve and improve in the years to come.

References

<https://www.kaggle.com/code/aswintechguy/image-caption-generator-tutorial-flickr-dataset>

<https://www.youtube.com/watch?v=fUSTbGrL1tc>

<https://www.kaggle.com/datasets/shadabhussain/flickr8k>

<https://machinelearningmastery.com/develop-a-deep-learning-caption-generation-model-in-python/>

<https://www.analyticsvidhya.com/blog/2021/12/step-by-step-guide-to-build-image-caption-generator-using-deep-learning/>