Double-click (or enter) to edit

# ▾ Image Captioning - DL Assignment2

## Assignment Details

1. **Group:90**
2. **Members**

   - Afrah Khan - 2019AD04096
   - Inderdeep Singh - 2019AD04061
   - Partha Pratim Saha - 2019AD04100

**Import Libraries/Dataset**

```
 1 import pandas as pd
 2 import numpy as np
 3 import re
 4 import os.path
 5 import tensorflow as tf
 6 from tensorflow.keras.preprocessing.image import load_img
 7 import matplotlib.pyplot as plt
 8 import glob
 9 from PIL import Image
10 import string
11 import time
12 pd.set_option('display.max_colwidth', None)
```

## ▾ Checking GPU availability

```
1 print("Version: ", tf.__version__)
2 print("Eager mode: ", tf.executing_eagerly())
3 print("GPU is", "available" if tf.config.list_physical_devices("GPU") else "NOT AVAILAB
4 print("GPU device name:", tf.test.gpu_device_name())
```

```
    Version:  2.5.0
    Eager mode:  True
    GPU is available
    GPU device name: /device:GPU:0
```

```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

```
    Mounted at /content/drive
```

# ▾ Data Visualization and augmentation:

## ▾ Read the pickle file

```
1 pkl_df = pd.read_pickle("/content/set_4.pkl")
```

```
1 pkl_list=[]
2 for d in pkl_df:
3   pkl_list.append(re.split('#|\t',d,maxsplit=2))
```

```
1 mdf = pd.DataFrame(pkl_list,columns =['ImageName', 'id','Caption'])
```

```
1 mdf = mdf.groupby(['ImageName']).sum()
2 mdf.drop(columns=['id'],inplace=True)
3 mdf.reset_index(inplace=True)
```

```
1 mdf.head()
```

| | ImageName | Caption |
|---|---|---|
| **0** | 1000268201_693b08cb0e.jpg | A girl go into a wooden building .A child in a pink dress be climb up a set of stair in an entry way .A little girl climb into a wooden playhouse .A little girl climb the stair to her playhouse .A little girl in a pink dress go into a wooden cabin . |
| **1** | 1001773457_577c3a7d70.jpg | A black dog and a tri-colored dog play with each other on a road .A black dog and a white dog with brown spot be stare at each other in a street .Two dog on pavement move toward each other . |
| **2** | 1002674143_1b742ab4b8.jpg | There be a girl with pigtail sit in front of a rainbow paint .A little girl be sit in front of a large painted rainbow . |
| | | A man lay on a bench while his dog sit by him .A man lay on a |

## ▾ Image dataset

```
1 import os
2 os.getcwd()
```

```
'/content/drive/My Drive/Flicker8k_Dataset'
```

**Convert the data into the correct format which could be used for ML model.**

```
1 from tqdm import tqdm
2 image_dict = {}
3 for filename in tqdm(glob.glob('*.jpg')):
4     img=load_img(filename,target_size=(224,224))
5     image_dict[filename]=img
```

```
9     image_dict[filename]=img
```

```
100%|████████| 8091/8091 [46:45<00:00,  2.88it/s]
```

```
1 len(image_dict)
```

```
8091
```

```
1 dic={}
2 for i in mdf.ImageName:
3   for j in image_dict:
4     if i in j:
5       dic[i]=image_dict[j]
```

```
1 df_test = pd.DataFrame.from_dict(dic,orient='index',columns=['Image'])
2 df_test.reset_index(inplace=True)
3 df_test.rename(columns={'index':'ImageName'},inplace=True)
```

```
1 main_df = pd.merge(mdf,df_test,on=['ImageName'])
```

```
1 main_df.head()
```

| | ImageName | Caption | Image |
|---|---|---|---|
| 0 | 1000268201_693b08cb0e.jpg | A girl go into a wooden building .A child in a pink dress be climb up a set of stair in an entry way .A little girl climb into a wooden playhouse .A little girl climb the stair to her playhouse .A little girl in a pink dress go into a wooden cabin . | <PIL.Image.Image image mode=RGB size=224x224 at 0x7F069F0719D0> |
| 1 | 1001773457_577c3a7d70.jpg | A black dog and a tri-colored dog play with each other on a road .A black dog and a white dog with brown spot be stare at each other in a street .Two dog on pavement move toward each other . | <PIL.Image.Image image mode=RGB size=224x224 at 0x7F069EEF1550> |
| 2 | 1002674143_1b742ab4b8.jpg | There be a girl with pigtail sit in front of a rainbow paint .A little girl be sit in front of a | <PIL.Image.Image image mode=RGB size=224x224 at |

```
1 main_df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 8032 entries, 0 to 8031
Data columns (total 3 columns):
 #   Column     Non-Null Count  Dtype
---  ------     --------------  -----
 0   ImageName  8032 non-null   object
 1   Caption    8032 non-null   object
 2   Image      8032 non-null   object
dtypes: object(3)
memory usage: 251.0+ KB
```

## Plot at least two samples and their captions (use matplotlib/seaborn/any other library)

```
1 for i in range(len(main_df.head())):
2   print("Caption: "+ str(main_df['Caption'][i].split(".")[:-1]))
3   plt.imshow(main_df['Image'][i])
4   plt.axis(False)
5   plt.show()
```

Caption: ['A child in a pink dress be climb up a set of stair in an entry way ', 'A ]



Caption: ['Two dog of different breed look at each other on a road ']



## Bring the train and test data in the required format

```
1 image = main_df['Image']
2 for i in range(len(image)):
3   image[i]=np.asarray(image[i])
```

```
1 from sklearn.model_selection import train_test_split
2 # keeping 50% size, as colab session crashes if we increase the train size
3 train, test = train_test_split(main_df, test_size=0.5)
```

```
1 caption_dict={}
2 for i in range(len(train)):
3   caption_dict[train.ImageName.values[i]]=["START "+sent.strip()+" END" for sent in tra
```

```
1 caption_dict.get('1000268201_693b08cb0e.jpg')
```

```
['START a girl go into a wooden building END',
 'START a child in a pink dress be climb up a set of stair in an entry way END',
 'START a little girl climb into a wooden playhouse END',
 'START a little girl climb the stair to her playhouse END',
 'START a little girl in a pink dress go into a wooden cabin END']
```

```
1 count_words = {}
2 count=1
3 for k,vv in caption_dict.items():
4     for v in vv:
```

```
5        for word in v.split():
6            if word not in count_words:
7                count_words[word] = count
8                count +=1
```

```
1 len(count_words)
```

```
4453
```

```
1 for k,vv in caption_dict.items():
2    for v in vv:
3        encode=[]
4        for word in v.split():
5          encode.append(count_words[word])
6        caption_dict[k][vv.index(v)]=encode
```

```
1 caption_dict.get('1000268201_693b08cb0e.jpg')
```

```
[[1, 2, 12, 473, 85, 2, 72, 382, 11],
 [1,
  2,
  37,
  4,
  2,
  249,
  119,
  38,
  168,
  158,
  2,
  182,
  17,
  181,
  4,
  223,
  1356,
  457,
  11],
 [1, 2, 331, 12, 168, 85, 2, 72, 1357, 11],
 [1, 2, 331, 12, 168, 9, 181, 78, 210, 1357, 11],
 [1, 2, 331, 12, 4, 2, 249, 119, 473, 85, 2, 72, 1358, 11]]
```

## Model Building

▼ **Use Pretrained Resnet-50 model trained on ImageNet dataset**

```
1 from tensorflow.keras.applications import ResNet50
2 resnet_model = ResNet50()
```

```
Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/re
102973440/102967424 [==============================] - 3s 0us/step
```

```python
1 from tensorflow.keras.models import Model
```

```python
1 image_model = Model(inputs=resnet_model.input, outputs=resnet_model.layers[-2].output)
```

## ▾ Bring the train in the required format.

```python
1 image_feat = {}
2 for im in range(len(train)):
3   img = train['Image'].values[im].reshape(1,224,224,3)
4   pred = image_model.predict(img).reshape(2048,)
5   image_feat[train.ImageName.values[im]] = pred
```

```python
1 from tensorflow.keras.utils import to_categorical
2 from tensorflow.keras.preprocessing.sequence import pad_sequences
```

```python
1 MAX_LEN = 0
2 for k, vv in caption_dict.items():
3     for v in vv:
4         if len(v) > MAX_LEN:
5             MAX_LEN = len(v)
```

```python
1 MAX_LEN
```

```
46
```

```python
1 VOCAB_SIZE = len(count_words)
2 def generator(photo, caption):
3     n_samples = 0
4     X = []
5     y_in = []
6     y_out = []
7     for k, vv in caption.items():
8         for v in vv:
9             for i in range(1, len(v)):
10                try:
11                    X.append(photo[k])
12                    in_seq= [v[:i]]
13                    out_seq = v[i]
14                    in_seq = pad_sequences(in_seq, maxlen=MAX_LEN, padding='post', truncati
15                    out_seq = to_categorical([out_seq], num_classes=VOCAB_SIZE+1)[0]
16                    y_in.append(in_seq)
17                    y_out.append(out_seq)
18                except:
19                    pass
20
21     return X, y_in, y_out
```

```python
1 X, y_in, y_out = generator(image_feat, caption_dict)
```

```
1 X = np.asarray(X)
2 y_in = np.asarray(y_in)
3 y_out = np.asarray(y_out)
```

```
1 X.shape,y_in.shape,y_out.shape
```

```
  ((143214, 2048), (143214, 46), (143214, 4454))
```

```
1 from tensorflow.keras.models import Model, Sequential
2 from tensorflow.keras import layers
3 from tensorflow.keras.layers import Dense
4 from tensorflow.keras.layers import LSTM
5 from tensorflow.keras.layers import Embedding
6 from tensorflow.keras.layers import Dropout
7 from tensorflow.keras.layers import RepeatVector
8 from tensorflow.keras.layers import TimeDistributed
9 from tensorflow.keras.layers import Concatenate
```

## ▼ Four(4) layered RNN - We are using LSTM as it is an RNN architecture

**Add one layer of dropout at the appropriate position and give reasons.**

Adding **L2 regularization** to all the **RNN** layers and one layer of 20% dropout - adding before last layer of **LSTM** abecause they are the one with the greater number of parameters and thus they're likely to excessively co-adapting themselves causing overfitting. However, since it's a stochastic regularization technique, we can really place it everywhere.

**Choose the appropriate activation function for all the layers.**

Adding **softmax** as appropriate activation function - **Adam** optimizer as it converges really faster and adding **categorical crossentropy** as loss function with **accuracy** as performance metric

## ▼ Give reasons for the choice of learning rate and its value: Adding learning rate as 0.0001 as it was hypertuned and gave better results

**Model Compilation**

```
1 embedding_size = 128
2 max_len = MAX_LEN
3 vocab_size = len(count_words) + 1
4 image_model1 = Sequential()
5 image_model1.add(Dense(embedding_size, input_shape=(2048,), activation='relu'))
6 image_model1.add(Dropout(0.2))
```

```
 7 image_model1.add(RepeatVector(max_len))
 8 language_model = Sequential()
 9 language_model.add(Embedding(input_dim=vocab_size, output_dim=embedding_size, input_len
10 language_model.add(LSTM(128, kernel_regularizer=tf.keras.regularizers.l2(l2=0.0001), re
11 language_model.add(Dropout(0.2))
12 language_model.add(TimeDistributed(Dense(embedding_size)))
13 conca = Concatenate()([image_model1.output, language_model.output])2
14 #Add L2 regularization to all the RNN layers.
15 x = LSTM(128, kernel_regularizer=tf.keras.regularizers.l2(l2=0.0001), return_sequences=
16 x = LSTM(128, kernel_regularizer=tf.keras.regularizers.l2(l2=0.0001), return_sequences=
17 out = Dense(vocab_size,activation='softmax')(x)
18 model = Model(inputs=[image_model1.input, language_model.input], outputs = out)
19
20 #Compiling the model with above parameters
21 model.compile(optimizer=tf.keras.optimizers.Adam(1e-3), loss=tf.losses.categorical_cros
22 #Print the model summary
23 model.summary()
```
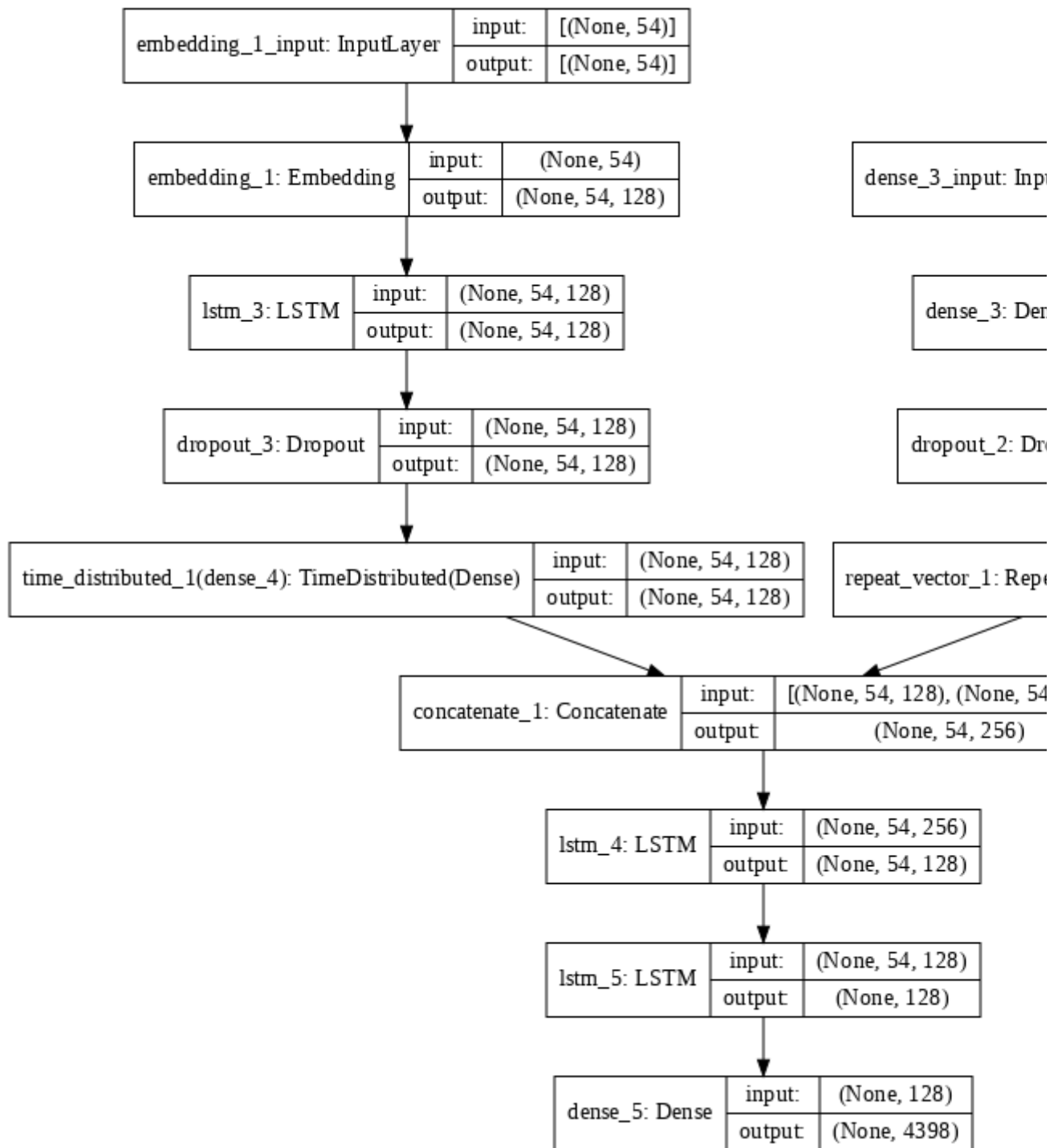
Model: "model_2"

| Layer (type) | Output Shape | Param # | Connected to |
|---|---|---|---|
| embedding_1_input (InputLayer) | [(None, 46)] | 0 | |
| dense_3_input (InputLayer) | [(None, 2048)] | 0 | |
| embedding_1 (Embedding) | (None, 46, 128) | 570112 | embedding_1_input[0] |
| dense_3 (Dense) | (None, 128) | 262272 | dense_3_input[0][0] |
| lstm_3 (LSTM) | (None, 46, 128) | 131584 | embedding_1[0][0] |
| dropout_2 (Dropout) | (None, 128) | 0 | dense_3[0][0] |
| dropout_3 (Dropout) | (None, 46, 128) | 0 | lstm_3[0][0] |
| repeat_vector_1 (RepeatVector) | (None, 46, 128) | 0 | dropout_2[0][0] |
| time_distributed_1 (TimeDistrib | (None, 46, 128) | 16512 | dropout_3[0][0] |
| concatenate_1 (Concatenate) | (None, 46, 256) | 0 | repeat_vector_1[0][0 time_distributed_1[0 |
| lstm_4 (LSTM) | (None, 46, 128) | 197120 | concatenate_1[0][0] |
| lstm_5 (LSTM) | (None, 128) | 131584 | lstm_4[0][0] |
| dense_5 (Dense) | (None, 4454) | 574566 | lstm_5[0][0] |

```
Total params: 1,883,750
Trainable params: 1,883,750
Non-trainable params: 0
```

```
1 from tensorflow.keras.utils import plot_model
2 plot_model(model,show_shapes=True,dpi=72)
```

| embedding_1_input: InputLayer | input: | [(None, 54)] |
|---|---|---|
| | output: | [(None, 54)] |

| embedding_1: Embedding | input: | (None, 54) |
|---|---|---|
| | output: | (None, 54, 128) |

| dense_3_input: Inp |
|---|

| lstm_3: LSTM | input: | (None, 54, 128) |
|---|---|---|
| | output: | (None, 54, 128) |

| dense_3: Der |
|---|

| dropout_3: Dropout | input: | (None, 54, 128) |
|---|---|---|
| | output: | (None, 54, 128) |

| dropout_2: Dr |
|---|

| time_distributed_1(dense_4): TimeDistributed(Dense) | input: | (None, 54, 128) |
|---|---|---|
| | output: | (None, 54, 128) |

| repeat_vector_1: Rep |
|---|

| concatenate_1: Concatenate | input: | [(None, 54, 128), (None, 54 |
|---|---|---|
| | output: | (None, 54, 256) |

| lstm_4: LSTM | input: | (None, 54, 256) |
|---|---|---|
| | output: | (None, 54, 128) |

| lstm_5: LSTM | input: | (None, 54, 128) |
|---|---|---|
| | output: | (None, 128) |

| dense_5: Dense | input: | (None, 128) |
|---|---|---|
| | output: | (None, 4398) |

**Model Training**

**Print the train and validation loss for each epoch. Use the appropriate batch size**

```
1 start = time.time()
2 history = model.fit([X, y_in], y_out, batch_size=720, epochs=35, validation_split=0.1,
3 end = time.time()
```

```
Epoch 1/35
180/180 [==============================] - 33s 146ms/step - loss: 5.5007 - accuracy
Epoch 2/35
180/180 [==============================] - 25s 136ms/step - loss: 5.2176 - accuracy
Epoch 3/35
180/180 [==============================] - 25s 137ms/step - loss: 5.1917 - accuracy
```

```
Epoch 4/35
180/180 [==============================] - 25s 138ms/step - loss: 5.1286 - accuracy
Epoch 5/35
180/180 [==============================] - 25s 138ms/step - loss: 5.0974 - accuracy
Epoch 6/35
180/180 [==============================] - 25s 136ms/step - loss: 5.0780 - accuracy
Epoch 7/35
180/180 [==============================] - 25s 136ms/step - loss: 5.0485 - accuracy
Epoch 8/35
180/180 [==============================] - 25s 139ms/step - loss: 5.0293 - accuracy
Epoch 9/35
180/180 [==============================] - 25s 139ms/step - loss: 5.0088 - accuracy
Epoch 10/35
180/180 [==============================] - 25s 138ms/step - loss: 4.9890 - accuracy
Epoch 11/35
180/180 [==============================] - 25s 138ms/step - loss: 4.9679 - accuracy
Epoch 12/35
180/180 [==============================] - 25s 139ms/step - loss: 4.9458 - accuracy
Epoch 13/35
180/180 [==============================] - 25s 138ms/step - loss: 4.9252 - accuracy
Epoch 14/35
180/180 [==============================] - 25s 137ms/step - loss: 4.8974 - accuracy
Epoch 15/35
180/180 [==============================] - 25s 137ms/step - loss: 4.8752 - accuracy
Epoch 16/35
180/180 [==============================] - 25s 138ms/step - loss: 4.8574 - accuracy
Epoch 17/35
180/180 [==============================] - 25s 140ms/step - loss: 4.8331 - accuracy
Epoch 18/35
180/180 [==============================] - 25s 138ms/step - loss: 4.8143 - accuracy
Epoch 19/35
180/180 [==============================] - 25s 138ms/step - loss: 4.7929 - accuracy
Epoch 20/35
180/180 [==============================] - 25s 139ms/step - loss: 4.7757 - accuracy
Epoch 21/35
180/180 [==============================] - 25s 139ms/step - loss: 4.7597 - accuracy
Epoch 22/35
180/180 [==============================] - 25s 138ms/step - loss: 4.7465 - accuracy
Epoch 23/35
180/180 [==============================] - 25s 137ms/step - loss: 4.7344 - accuracy
Epoch 24/35
180/180 [==============================] - 25s 138ms/step - loss: 4.7160 - accuracy
Epoch 25/35
180/180 [==============================] - 25s 139ms/step - loss: 4.7021 - accuracy
Epoch 26/35
180/180 [==============================] - 25s 138ms/step - loss: 4.6927 - accuracy
Epoch 27/35
180/180 [==============================] - 25s 138ms/step - loss: 4.6763 - accuracy
Epoch 28/35
180/180 [==============================] - 25s 139ms/step - loss: 4.6613 - accuracy
Epoch 29/35
180/180 [==============================] - 25s 139ms/step - loss: 4.6546 - accuracy
```

**total time** taken for training: ~15 mins

```
1 # Print the total time taken for training
2 print("Total time taken for training: ", (end - start)/60, "mins")
```
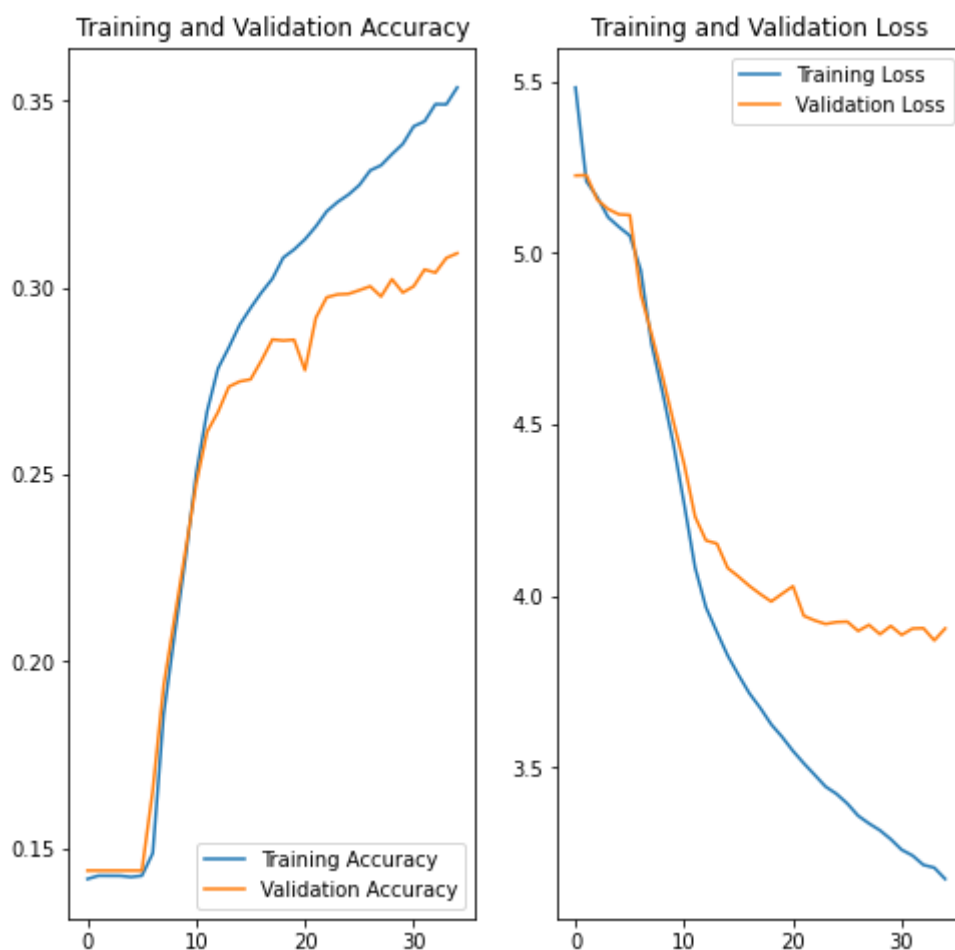
    Total time taken for training:  15.453375518321991 mins

## Plot the loss and accuracy history graphs for both train and validation set

```
1 acc = history.history['accuracy']
2 val_acc = history.history['val_accuracy']
3 loss = history.history['loss']
4 val_loss = history.history['val_loss']
5 epochs_range = range(len(history.epoch))
6
7 plt.figure(figsize=(8, 8))
8 plt.subplot(1, 2, 1)
9 plt.plot(epochs_range, acc, label='Training Accuracy')
10 plt.plot(epochs_range, val_acc, label='Validation Accuracy')
11 plt.legend(loc='lower right')
12 plt.title('Training and Validation Accuracy')
13
14 plt.subplot(1, 2, 2)
15 plt.plot(epochs_range, loss, label='Training Loss')
16 plt.plot(epochs_range, val_loss, label='Validation Loss')
17 plt.legend(loc='upper right')
18 plt.title('Training and Validation Loss')
19 plt.show()
```



```
1 inv_dict = {v:k for k, v in count_words.items()}
```

```
1 def getImage(x):
```

```
1 def getImage(x):
2     test_img = np.reshape(x, (1,224,224,3))
3     return test_img
```

## Model Evaluation: Take a random image from google and generate caption for that image
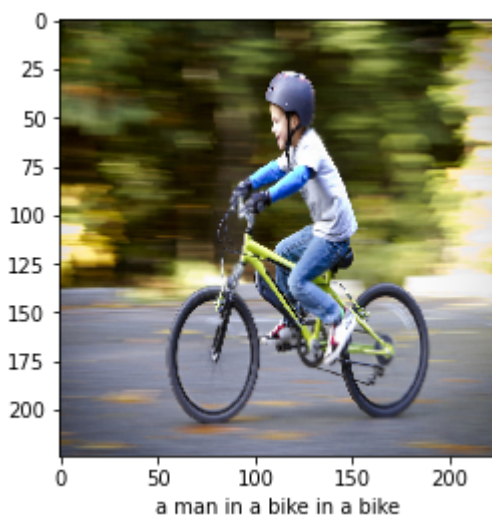
```
1 testImg = load_img('/content/test_image3.jpg',target_size=(224,224,3))
```

```
 1 test_pred = image_model.predict(getImage(testImg)).reshape(1,2048)
 2 text_inp = ['START']
 3 count = 0
 4 caption = ''
 5 while count < 25:
 6             count += 1
 7             encoded = []
 8             for i in text_inp:
 9                 encoded.append(count_words[i])
10             encoded = [encoded]
11             encoded = pad_sequences(encoded, padding='post', truncating='post', maxlen=
12             prediction = np.argmax(model.predict([test_pred, encoded]))
13             sampled_word = inv_dict[prediction]
14             caption = caption + ' ' + sampled_word
15             if sampled_word == 'END':
16                 break
17             text_inp.append(sampled_word)
18 plt.figure()
19 plt.imshow(testImg)
20 plt.xlabel(caption.replace("END","").strip())
```

```
Text(0.5, 0, 'a man in a bike in a bike')
```



a man in a bike in a bike

✓ 1s    completed at 23:02    ● ✕