

Architecture Document

Image Captioning

Revision No: 1.0

Last Date of Revision: 23/11/2023.

Document Version Control-

<u>Date Issued</u>	<u>Version</u>	<u>Description</u>	<u>Author</u>
23/11/2023	1.0	LLD – V- 1.0	Mahesh. A

Contents -

Description	Page No.
Document Version Control	2
1. Introduction	4
1.1 What is Low-Level design document	4
1.2 Scope	4
2. Architecture	4
3. Architecture Description	5
3.1 Data Accessing	5
3.2 Data Pre-Processing	7
3.3 Splitting the Data	9
3.4 Model Building	9
3.5 Create Front End User Module using flask	9
3.6 Testing the Model	12
4.0 KPI	17

Abstract

By using the pre-trained of Image Captioning model user can get the auto caption for the uploaded image...

1 Introduction

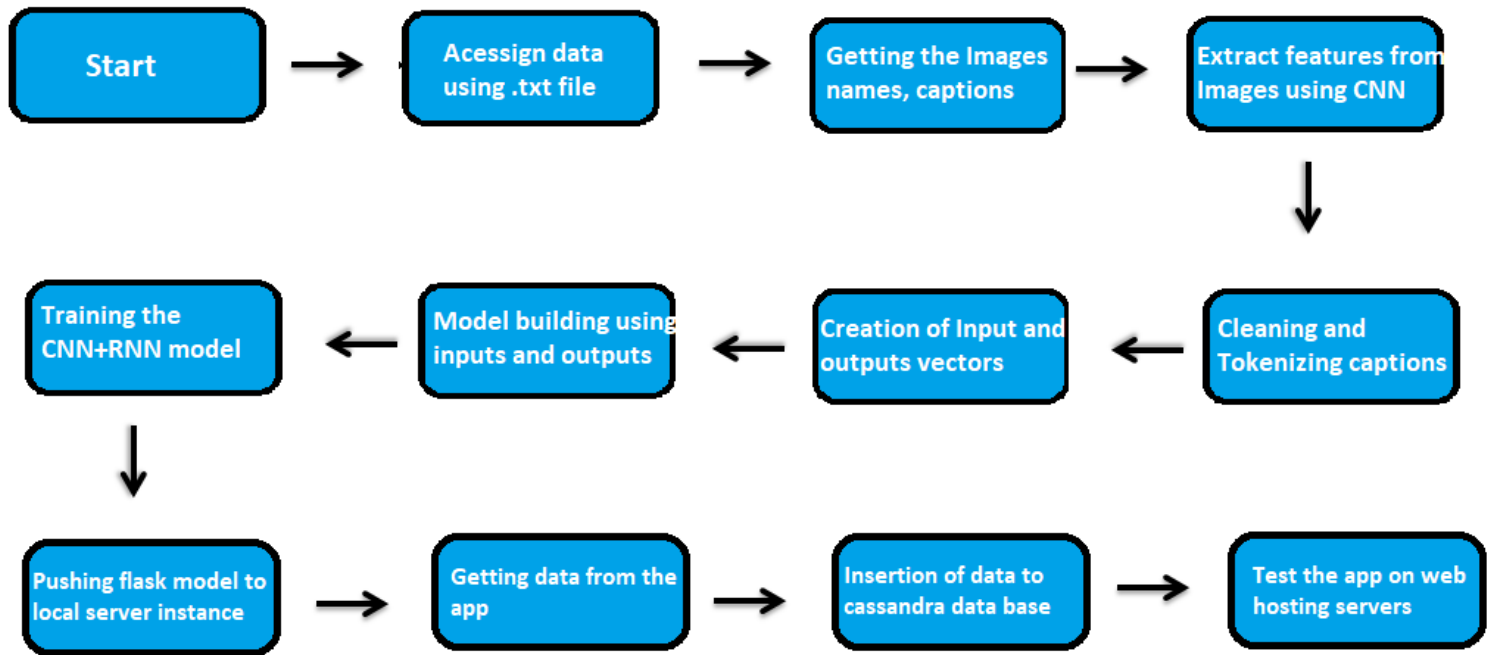
1.1. What is document required?

The goal of Architecture is to give the internal logical design of the actual programmed code for Image captioning model. Describes the captioned relations with images .It describes the modules so that the programmer can directly code the program from the document.

1.2. Scope

Architecture document is a component-level design process that follows a step-by-step process. This process can be used for designing data structures, required software architecture, source code and ultimately, performance algorithms. Overall, the data organization may be defined during requirement analysis and then refined during data design work.

2 Architecture



3 Architecture Description

3.1 Data Accessing

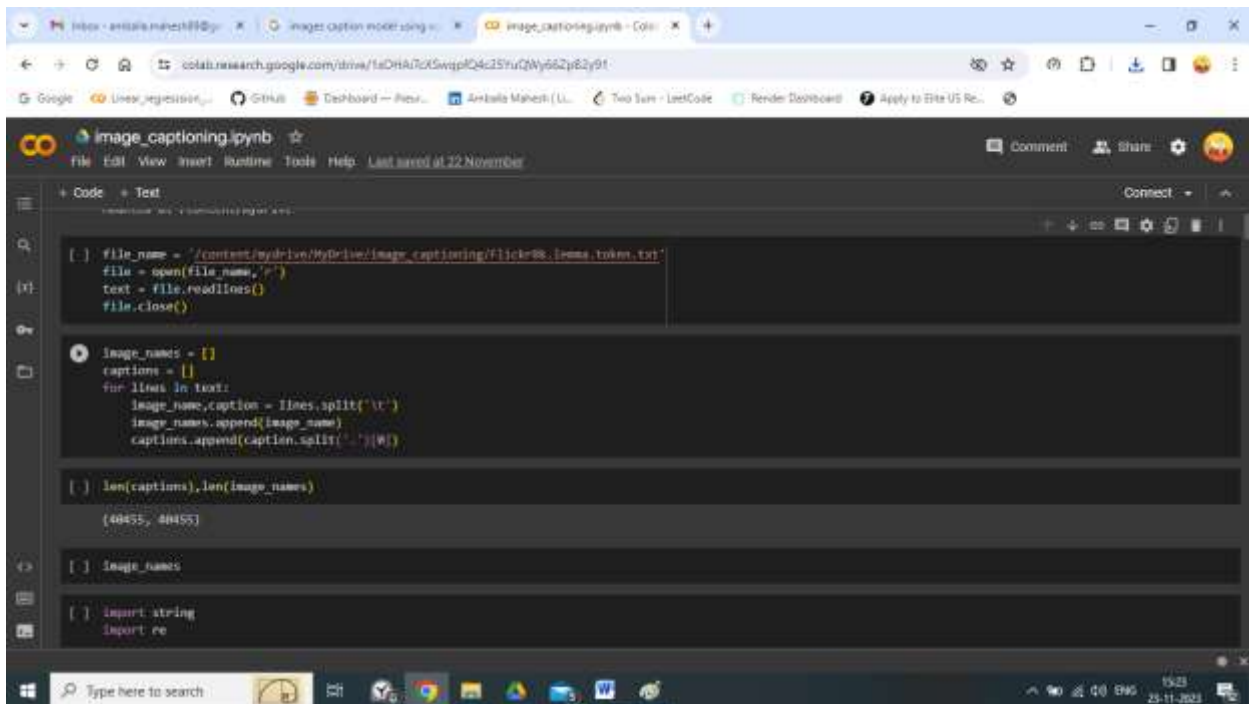
We can access the data in the from the download link as available in the project.

We load the data to the framework using the pandas read function.

file name	caption
1305564994_00513f9a5b.jpg#0	A man in street racer armor be examine the tire of another racer 's motorbike .
1305564994_00513f9a5b.jpg#1	Two racer drive a white bike down a road .
1305564994_00513f9a5b.jpg#2	Two motorist be ride along on their vehicle that be oddly design and color .
1305564994_00513f9a5b.jpg#3	Two person be in a small race car drive by a green hill .
1305564994_00513f9a5b.jpg#4	Two person in race uniform in a street car .
1351764581_4d4fb1b40f.jpg#0	A firefighter extinguish a fire under the hood of a car .
1351764581_4d4fb1b40f.jpg#1	a fireman spray water into the hood of small white car on a jack
1351764581_4d4fb1b40f.jpg#2	A fireman spray inside the open hood of small white car , on a jack .
1351764581_4d4fb1b40f.jpg#3	A fireman use a firehose on a car engine that be up on a carjack .
1351764581_4d4fb1b40f.jpg#4	Firefighter use water to extinguish a car that be on fire .
1358089136_976e3d2e30.jpg#0	A boy sand surf down a hill
1358089136_976e3d2e30.jpg#1	A man be attempt to surf down a hill make of sand on a sunny day .
1358089136_976e3d2e30.jpg#2	A man be slide down a huge sand dune on a sunny day .

3.2 Data Pre-Processing

By the usage of the different data manipulation techniques we will remove unwanted words by this reduce the dimensions and make all the features in numerical data type using Keras tokenizer to get the vector of caption and Keras Xception model to get features from images.



```
[ ] file_name = '/content/drive/MyDrive/image_captioning/flickr8k_image_tokens.txt'
file = open(file_name, 'r')
text = file.readlines()
file.close()

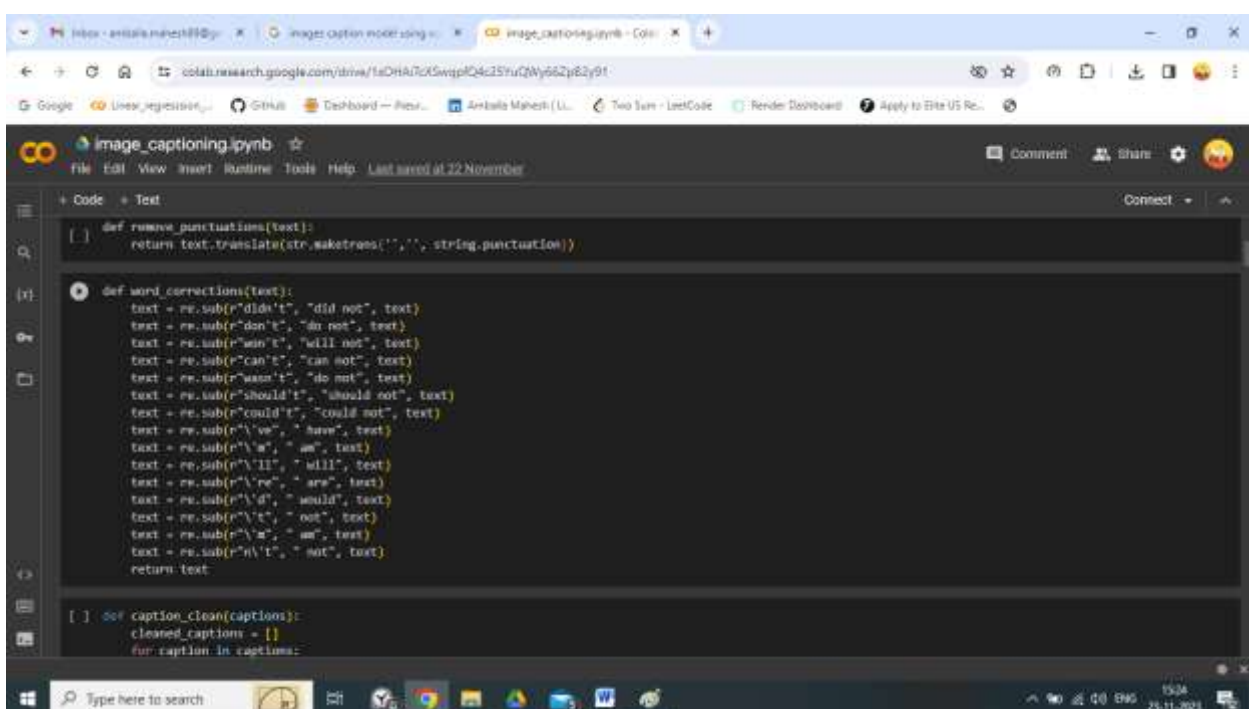
[ ] image_names = []
captions = []
for lines in text:
    image_name, caption = lines.split('\t')
    image_names.append(image_name)
    captions.append(caption.split('.')[0])

[ ] len(captions), len(image_names)

(40455, 40455)

[ ] image_names

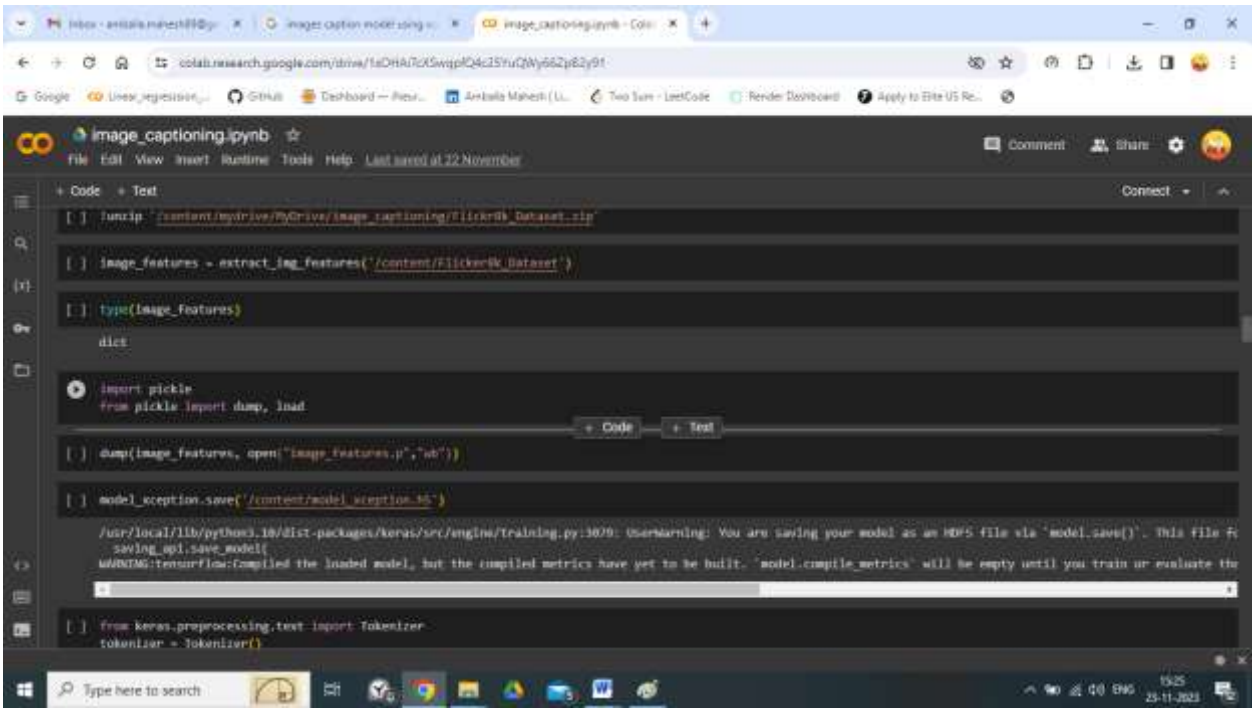
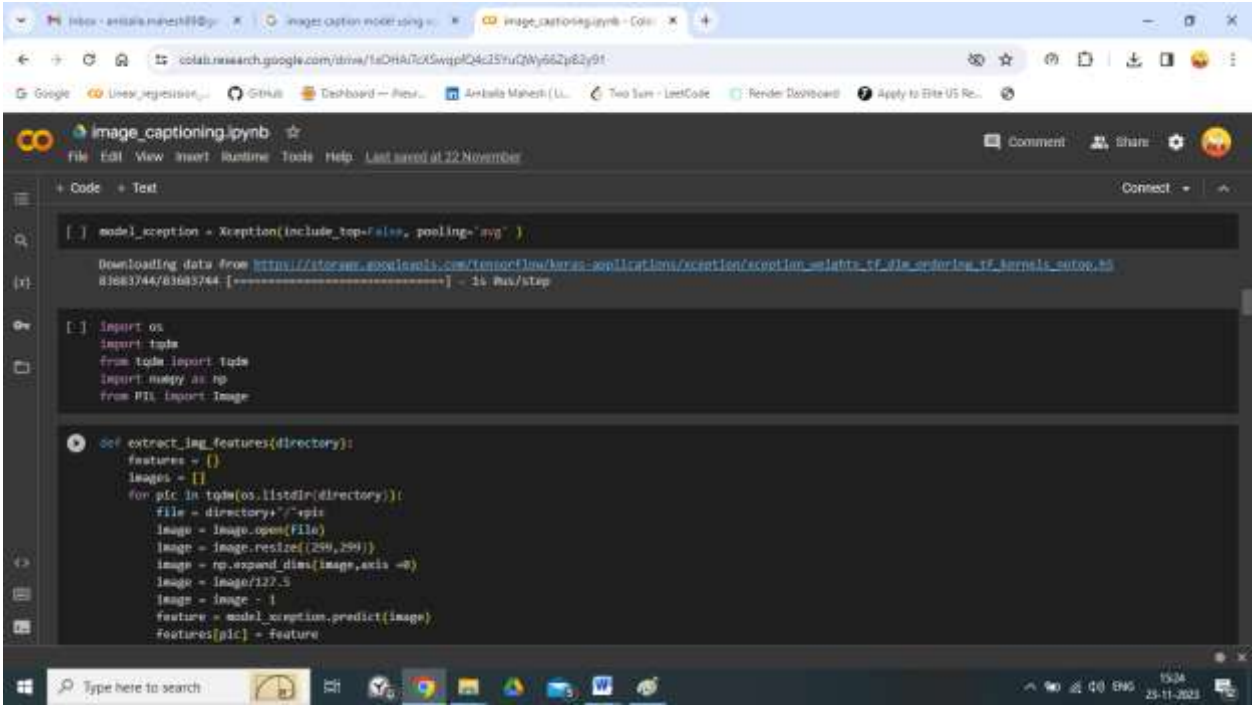
[ ] import string
import re
```

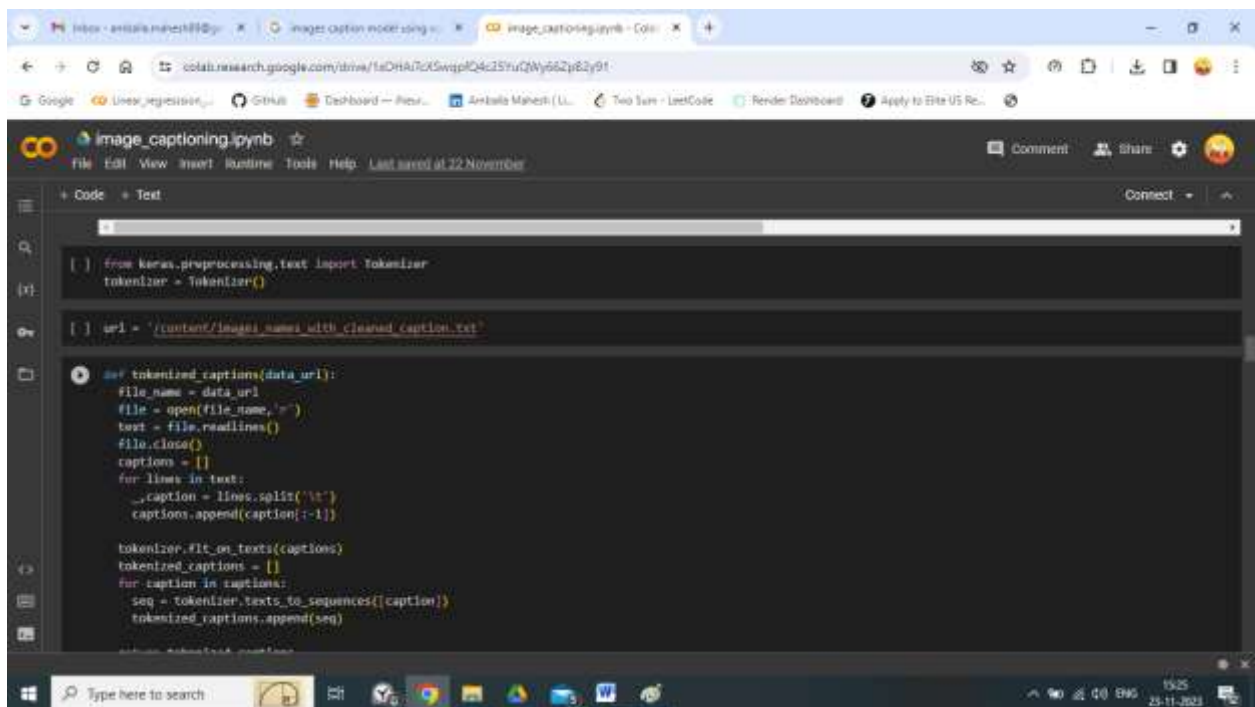


```
[ ] def remove_punctuations(text):
    return text.translate(str.maketrans('', '', string.punctuation))

[ ] def word_corrections(text):
    text = re.sub(r"didn't", "did not", text)
    text = re.sub(r"don't", "do not", text)
    text = re.sub(r"won't", "will not", text)
    text = re.sub(r"can't", "can not", text)
    text = re.sub(r"wasn't", "do not", text)
    text = re.sub(r"should't", "should not", text)
    text = re.sub(r"could't", "could not", text)
    text = re.sub(r"\'ve", " have", text)
    text = re.sub(r"\'m", " am", text)
    text = re.sub(r"\'ll", " will", text)
    text = re.sub(r"\'re", " are", text)
    text = re.sub(r"\'d", " would", text)
    text = re.sub(r"\'t", " not", text)
    text = re.sub(r"\'s", " am", text)
    text = re.sub(r"n't", " not", text)
    return text

[ ] def caption_clean(captions):
    cleaned_captions = []
    for caption in captions:
```





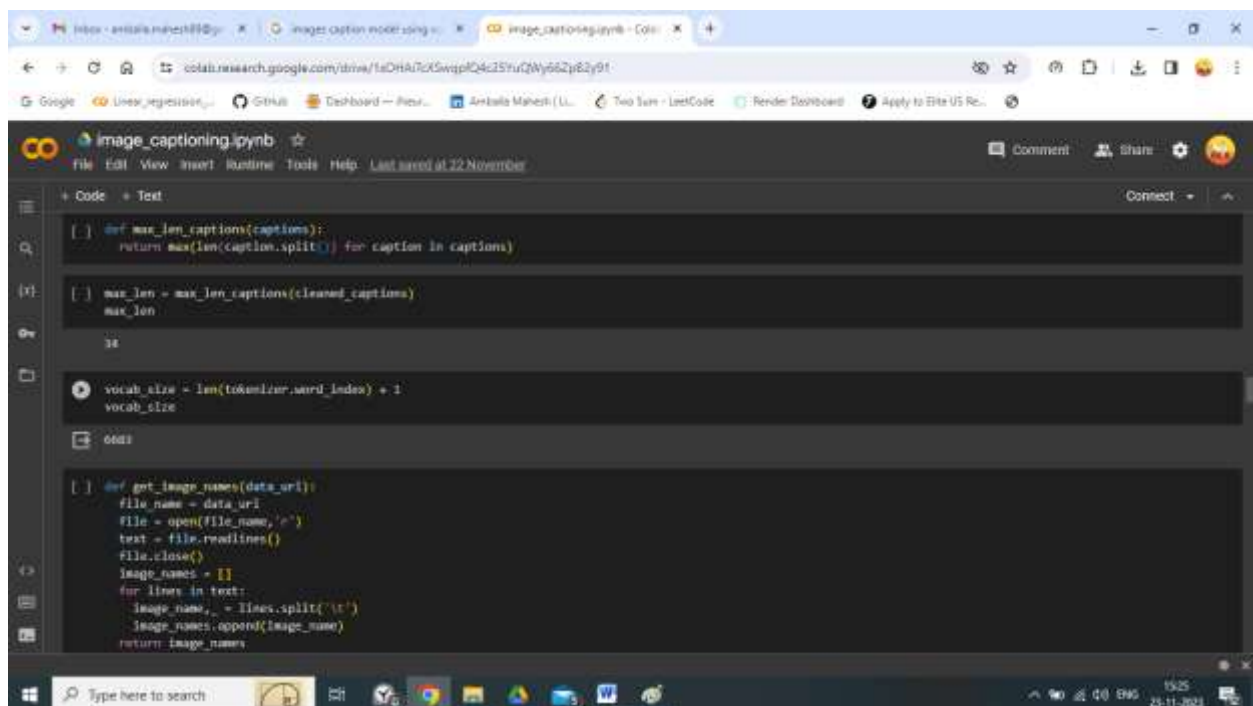
image_captioning.ipynb

```
[ ] from keras.preprocessing.text import Tokenizer
tokenizer = Tokenizer()

[ ] url = '/content/images_names_with_cleaned_caption.txt'

[ ] def tokenize_captions(data_url):
    file_name = data_url
    file = open(file_name, 'r')
    text = file.readlines()
    file.close()
    captions = []
    for lines in text:
        caption = lines.split('\t')
        captions.append(caption[:-1])

    tokenizer.fit_on_texts(captions)
    tokenized_captions = []
    for caption in captions:
        seq = tokenizer.texts_to_sequences([caption])
        tokenized_captions.append(seq)
```



image_captioning.ipynb

```
[ ] def max_len_captions(captions):
    return max(len(caption.split()) for caption in captions)

[ ] max_len = max_len_captions(cleaned_captions)
max_len
34

[ ] vocab_size = len(tokenizer.word_index) + 1
vocab_size
6043

[ ] def get_image_names(data_url):
    file_name = data_url
    file = open(file_name, 'r')
    text = file.readlines()
    file.close()
    image_names = []
    for lines in text:
        image_name, _ = lines.split('\t')
        image_names.append(image_name)
    return image_names
```


image_captioning.ipynb

```
(40455, 40455)

[ ] all_features = load(open("../content/aidrive/hydrive/image_captioning/image_features.p", "rb"))

[ ] features = []
for name in image_names:
    features.append(all_features[name][0])

[ ] len(features)

40455

[ ] len(caps)

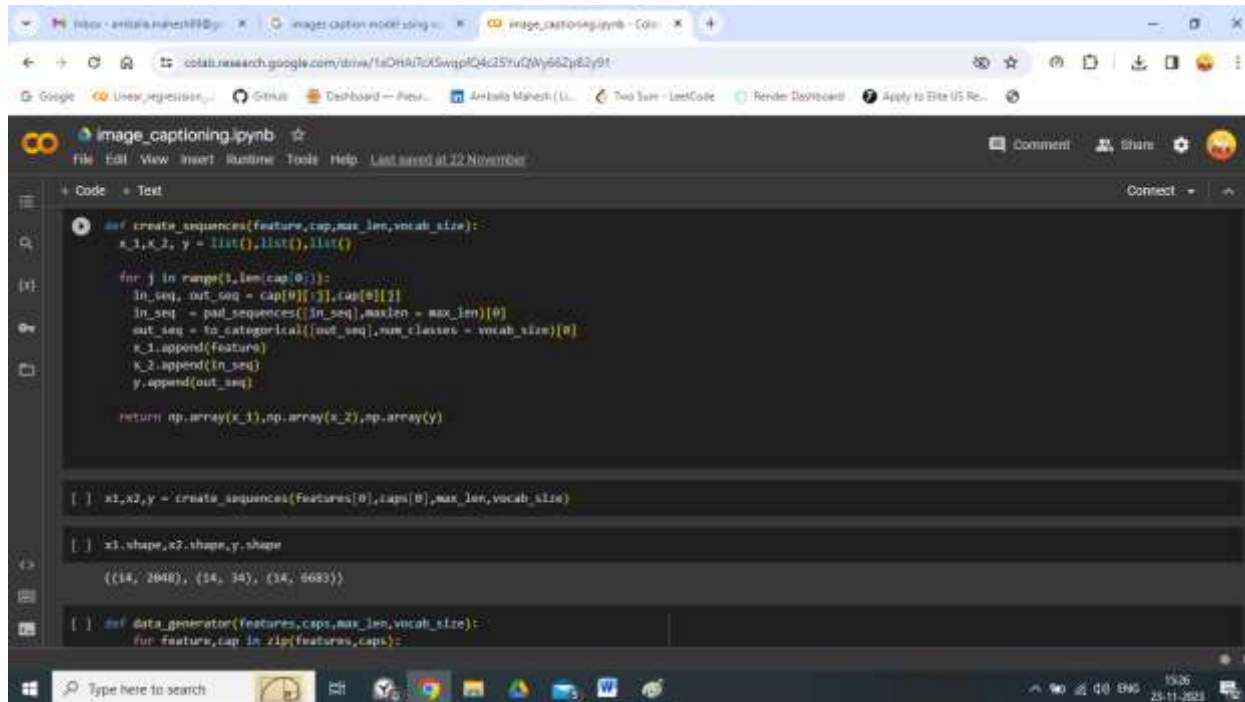
40455

[ ] caps[0]

[[2, 0, 3, 69, 531, 3879, 4, 989, 6, 359, 11, 70, 531, 663, 1]]

[ ] from keras.preprocessing.sequence import pad_sequences
from keras.utils import to_categorical
```

3.3 Creating Input and Output data



The screenshot shows a Jupyter Notebook interface with the following code and output:

```
def create_sequences(feature_cap,max_len,vocab_size):
    x_1,x_2,y = list(),list(),list()

    for j in range(1,len(cap[0])):
        in_seq, out_seq = cap[0][1:],cap[0][2:]
        in_seq = pad_sequences([in_seq],maxlen = max_len)[0]
        out_seq = to_categorical([out_seq],num_classes = vocab_size)[0]
        x_1.append(feature)
        x_2.append(in_seq)
        y.append(out_seq)

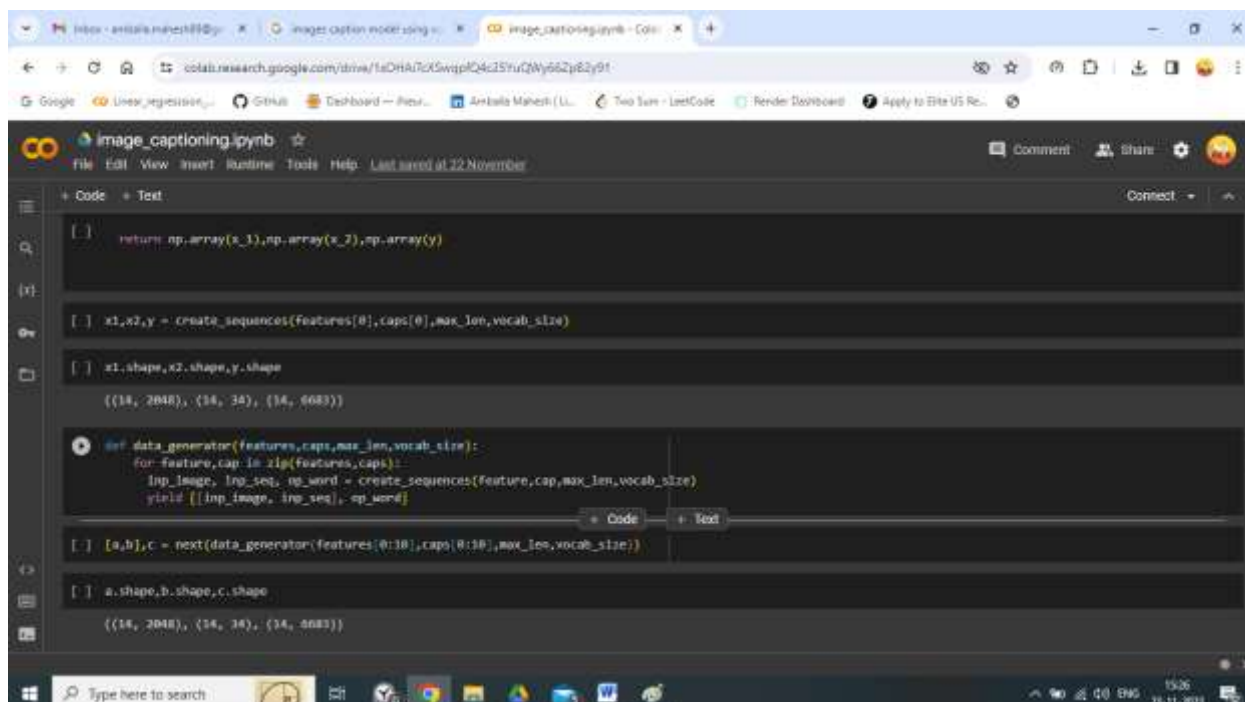
    return np.array(x_1),np.array(x_2),np.array(y)
```

```
[ ] x1,x2,y = create_sequences(features[0],caps[0],max_len,vocab_size)
```

```
[ ] x1.shape,x2.shape,y.shape
```

```
[(18, 2048), (14, 34), (14, 6683)]
```

```
[ ] def data_generator(features,caps,max_len,vocab_size):
    for feature,cap in zip(features,caps):
```



The screenshot shows the continuation of the Jupyter Notebook with the following code and output:

```
return np.array(x_1),np.array(x_2),np.array(y)
```

```
[ ] x1,x2,y = create_sequences(features[0],caps[0],max_len,vocab_size)
```

```
[ ] x1.shape,x2.shape,y.shape
```

```
[(18, 2048), (14, 34), (14, 6683)]
```

```
def data_generator(features,caps,max_len,vocab_size):
    for feature,cap in zip(features,caps):
        inp_image, inp_seq, np_word = create_sequences(feature,cap,max_len,vocab_size)
        yield [(inp_image, inp_seq), np_word]
```

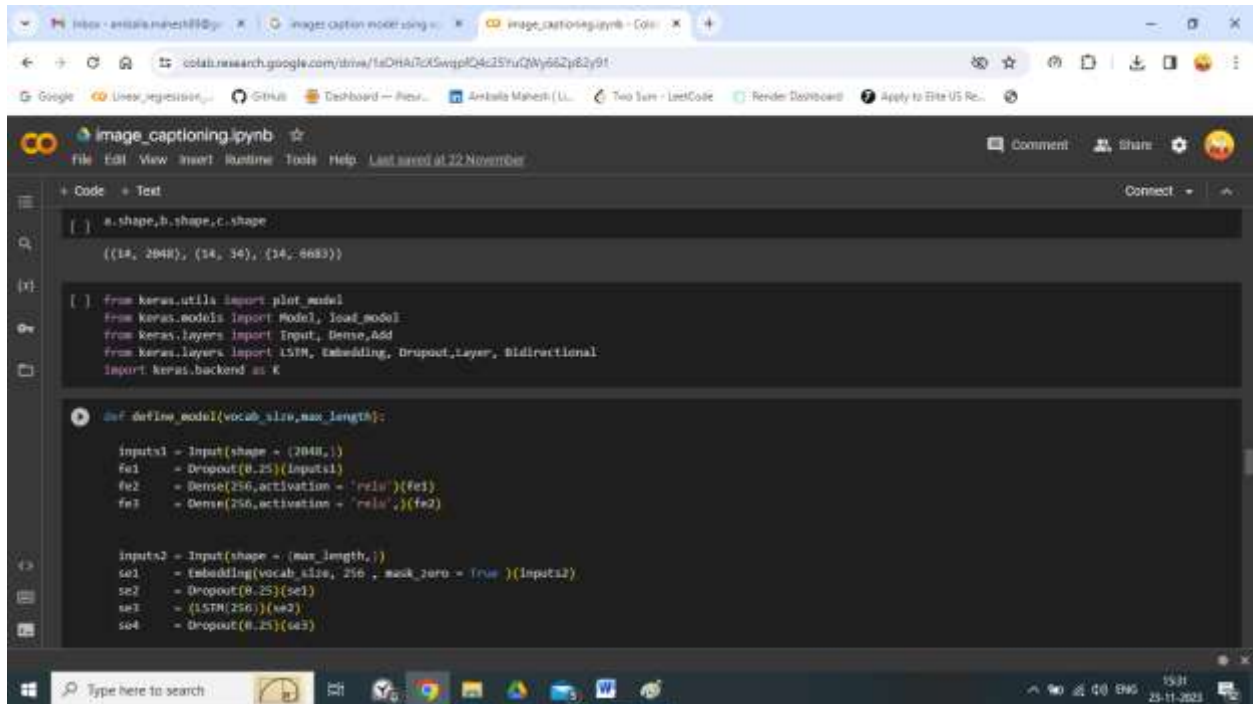
```
[ ] [a,b],c = next(data_generator(features[0:10],caps[0:10],max_len,vocab_size))
```

```
[ ] a.shape,b.shape,c.shape
```

```
[(18, 2048), (14, 34), (14, 6683)]
```

3.4 Model Building

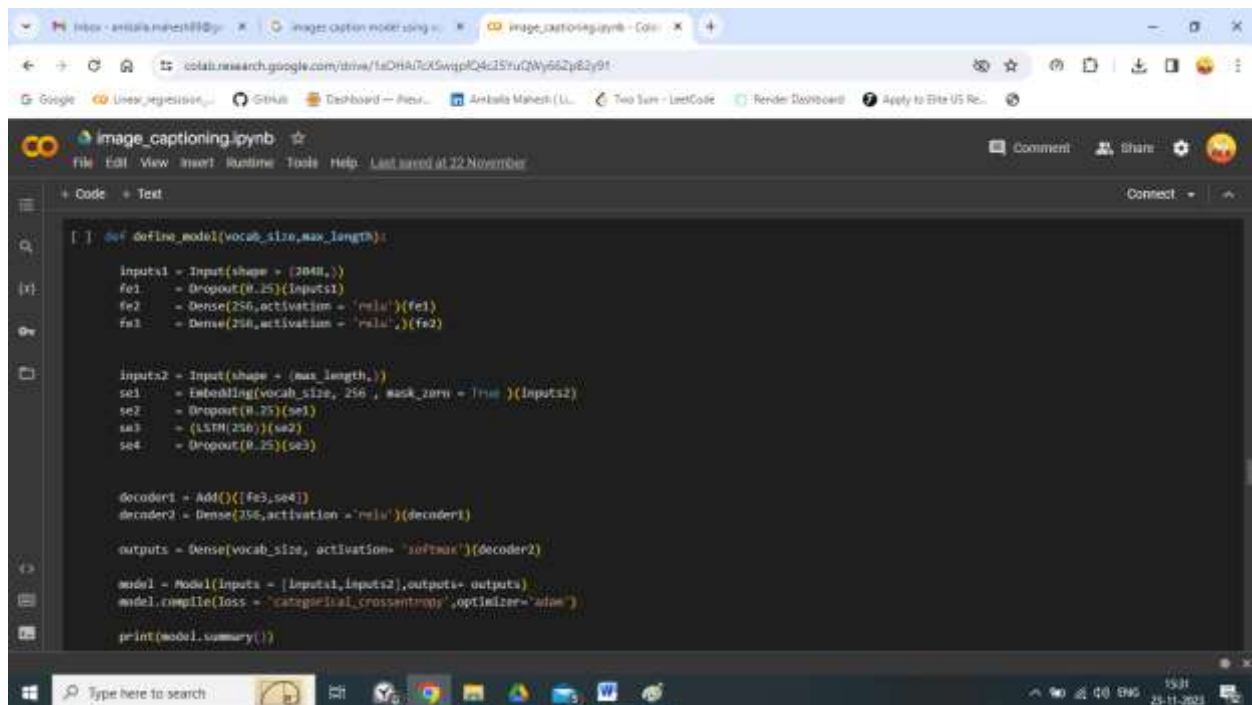
We create CNN + RNN base model using Keras as mentioned below.



```
from keras.utils import plot_model
from keras.models import Model, load_model
from keras.layers import Input, Dense, Add
from keras.layers import LSTM, Embedding, Dropout, Layer, Bidirectional
import keras.backend as K

def define_model(vocab_size, max_length):
    inputs1 = Input(shape = (2048,))
    fe1 = Dropout(0.25)(inputs1)
    fe2 = Dense(256, activation = 'relu')(fe1)
    fe3 = Dense(256, activation = 'relu')(fe2)

    inputs2 = Input(shape = (max_length,))
    se1 = Embedding(vocab_size, 256, mask_zero = True)(inputs2)
    se2 = Dropout(0.25)(se1)
    se3 = (LSTM(256))(se2)
    se4 = Dropout(0.25)(se3)
```

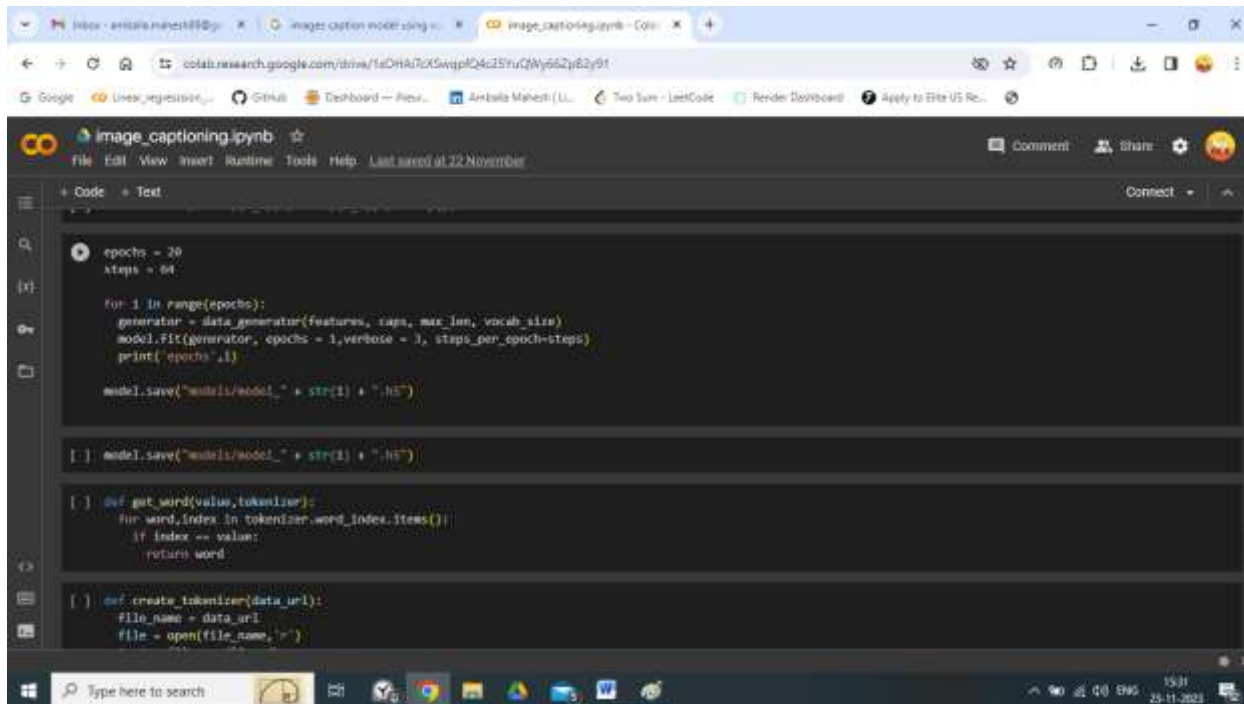


```
decoder1 = Add()([fe3, se4])
decoder2 = Dense(256, activation = 'relu')(decoder1)

outputs = Dense(vocab_size, activation = 'softmax')(decoder2)

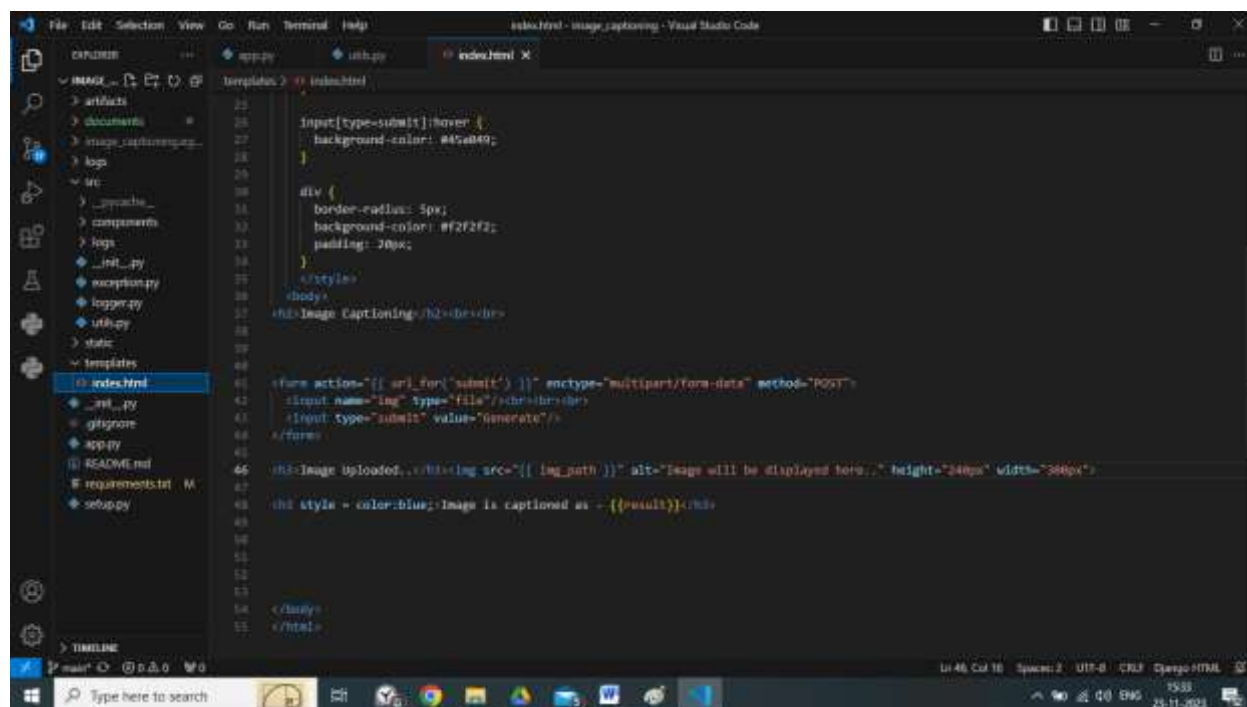
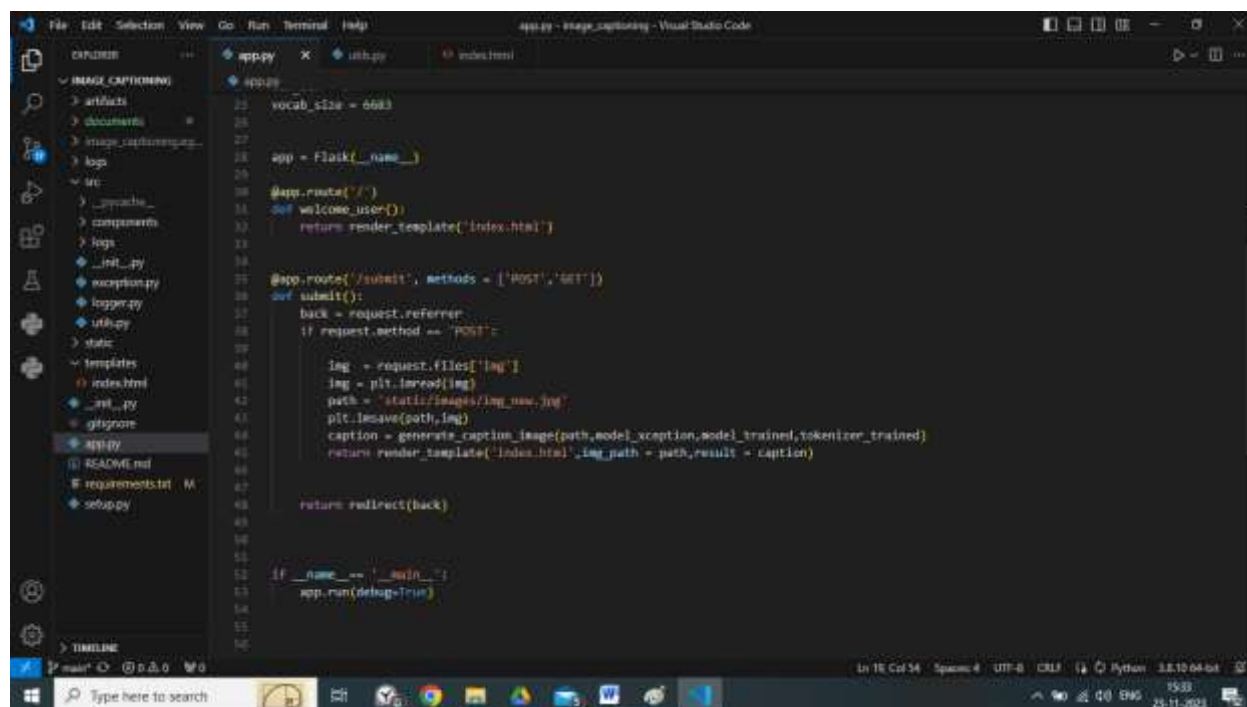
model = Model(inputs = [inputs1, inputs2], outputs = outputs)
model.compile(loss = 'categorical_crossentropy', optimizer = 'adam')

print(model.summary())
```



3.5 Create Front End User Module using flask

Once the model is created download and save the model and now we create GUI for front end user using the flask incorporated with HTML, CSS. Align and map the user data to the data base created. From user data create the data frame and load it to the model for the prediction the same prediction is send back to the user GUI and well saved in the data base (MySQL and Cassandra).



3.6 Testing the Model

- ✓ Verify whether the application is the loading on the local server instance.
- ✓ Verify whether the user can access the application.
- ✓ Verify the user can access the different fields for selection and can be visible
- ✓ Once the user selection the fields and made the submit
- ✓ Check the user can get the result or prediction.
- ✓ Once he gets the prediction.
- ✓ Check the data form the user and prediction from the model is loaded into the local MySQL and Cassandra Data base.
- ✓ Verify whether the application is the loading on the web service instance.
- ✓ Check the database and download the data...

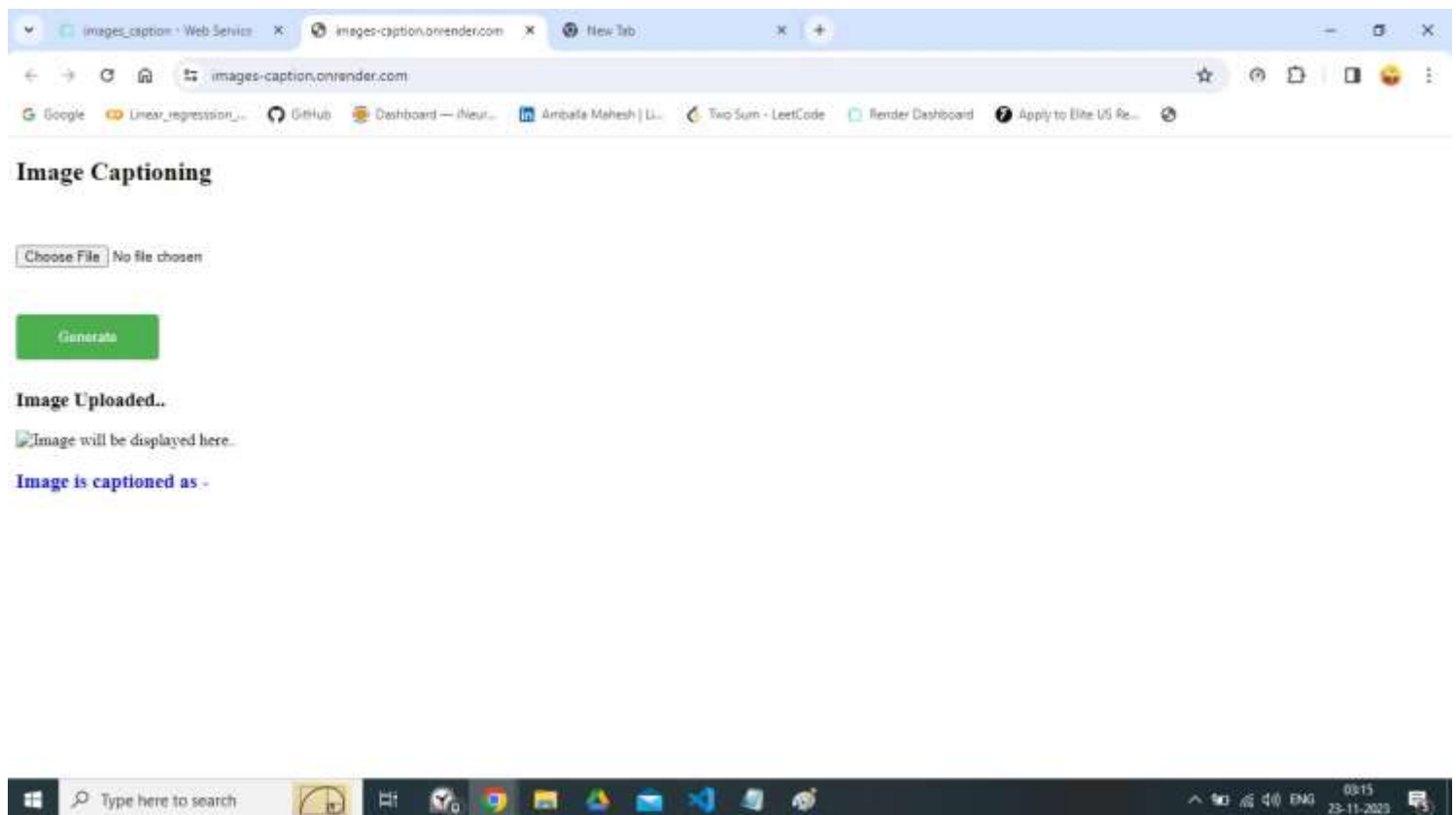


Image Captioning

Choose File No file chosen

Generate

Image Uploaded..



Image is captioned as - little league player point to another player

Image Captioning

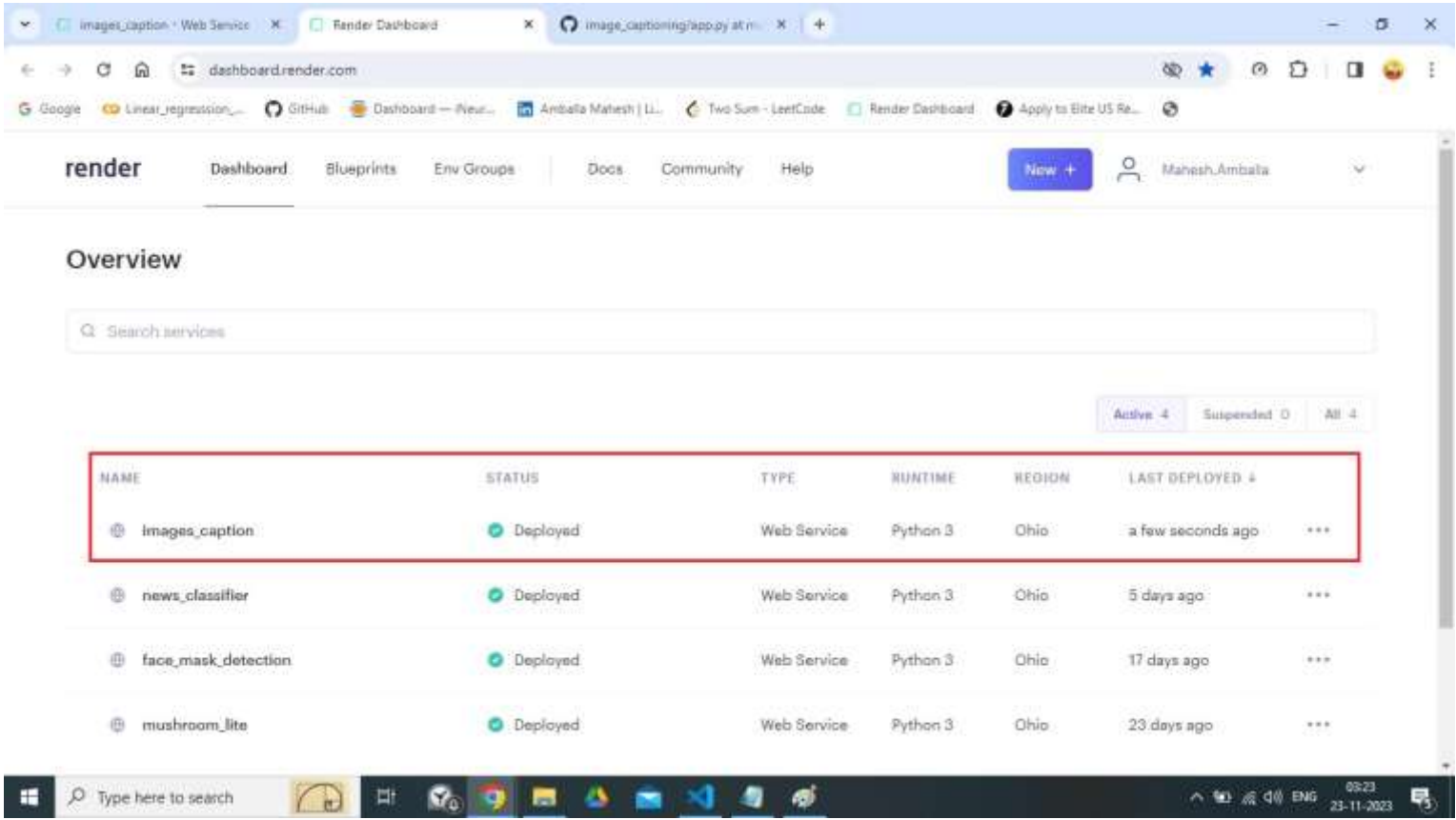
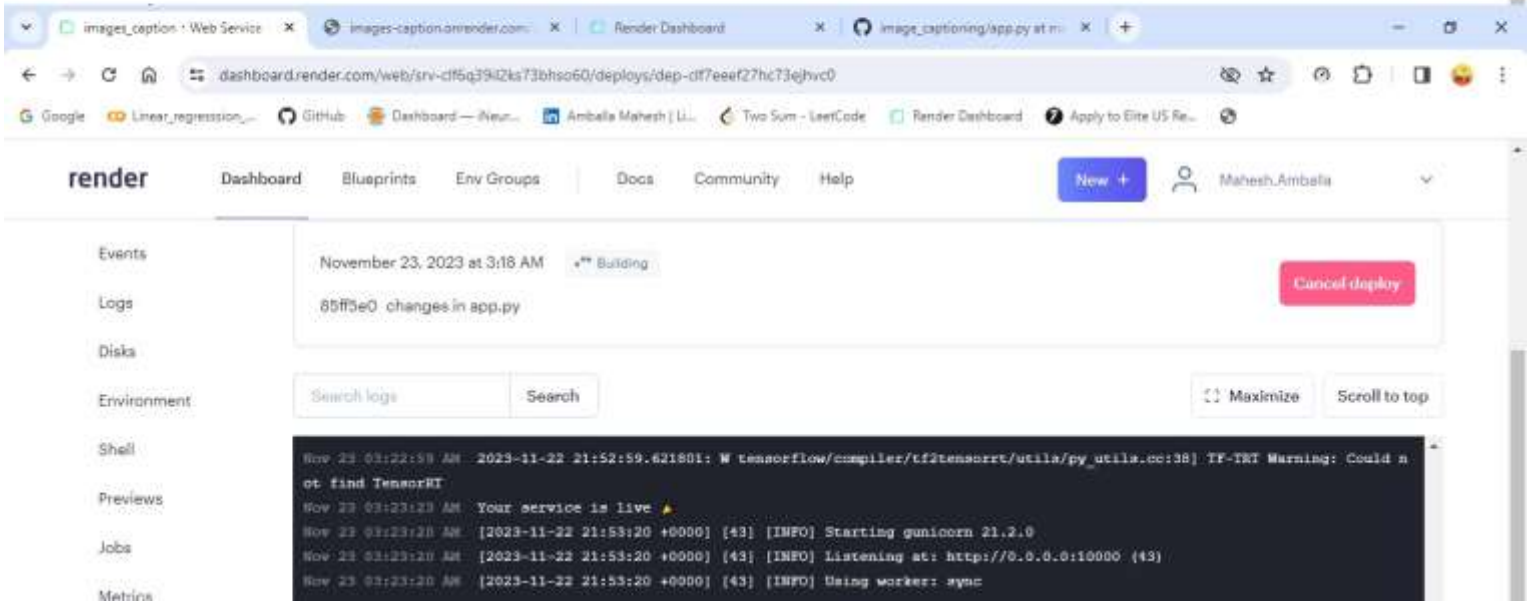
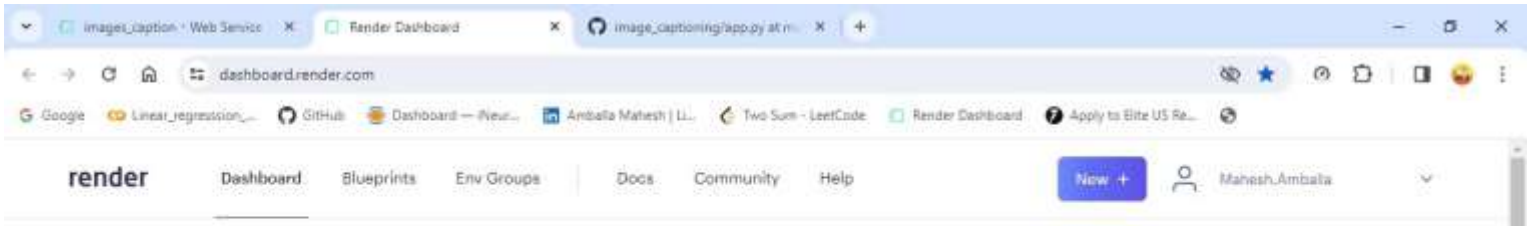
Choose File No file chosen

Generate

Image Uploaded..



Image is captioned as - little league player point to another player



HOST WEB ADDRESS: <https://images-caption.onrender.com>

4. Key performance indicators (KPI)

- Time and work load reduction by using the flask model.
- Compare the accuracy of model using prediction and actual results.
- Check for the wrong predictions
- If found any wrong predictions again train the model with the new data along with previous data
- Retest the model unless the productions attain the good results.