# Configuration Manual

MSc Research Project
Data Analytics

# Namita Mohan

Student ID: 19212500

School of Computing
National College of Ireland

Supervisor:     Dr Paul Stynes and Dr Pramod Pathak

**National College of Ireland**
**Project Submission Sheet**
**School of Computing**

| | |
|---|---|
| **Student Name:** | Namita Mohan |
| **Student ID:** | 19212500 |
| **Programme:** | Data Analytics |
| **Year:** | 2020/2021 |
| **Module:** | MSc Research Project |
| **Supervisor:** | Dr Paul Stynes and Dr Pramod Pathak |
| **Submission Due Date:** | 16/08/2021 |
| **Project Title:** | Configuration Manual |
| **Word Count:** | 1109 |
| **Page Count:** | 11 |

I hereby certify that the information contained in this (my submission) is information pertaining to research I conducted for this project. All information other than my own contribution will be fully referenced and listed in the relevant bibliography section at the rear of the project.

__ALL__ internet material must be referenced in the bibliography section. Students are required to use the Referencing Standard specified in the report template. To use other author's written or electronic work is illegal (plagiarism) and may result in disciplinary action.

| **Signature:** | Namita Mohan |
|---|---|
| **Date:** | 18th September 2021 |

**PLEASE READ THE FOLLOWING INSTRUCTIONS AND CHECKLIST:**

| | |
|---|---|
| Attach a completed copy of this sheet to each project (including multiple copies). | ☐ |
| **Attach a Moodle submission receipt of the online project submission**, to each project (including multiple copies). | ☐ |
| **You must ensure that you retain a HARD COPY of the project**, both for your own reference and in case a project is lost or mislaid. It is not sufficient to keep a copy on computer. | ☐ |

Assignments that are submitted to the Programme Coordinator office must be placed into the assignment box located outside the office.

| **Office Use Only** | |
|---|---|
| Signature: | |
| Date: | |
| Penalty Applied (if applicable): | |

# Configuration Manual

Namita Mohan
19212500

## 1 Introduction

This configuration manual will discuss steps followed in the implementation of this research in detail. The document covers steps from data collection to the evaluation of the implemented remote sensing image captioning framework. It includes system configuration, code snippets and directions to execute to ensure that this work can be reproduced in future easily.

## 2 Hardware and Software configurations

All the execution for this research is done on the local machine. To track progress and save time, processed data and models were saved locally as pickle files. Table 1 and Table 2 presents hardware and software specifications for this research. It is also suggested to use the latest version of Tensorflow and Keras for this implementation.

Table 1: Hardware specifications

| Host machine/Operating System | MacBook Pro/MacOS Mojave |
|---|---|
| RAM | 8GB, 2.7 GHz Intel Core i5 |
| Hard Disk | 256GB |
| GPU | Intel Iris Graphics 6100 1536 MB |

Table 2: Software specifications

| Programming language | Python (Anaconda distribution) |
|---|---|
| IDE | Jupyter notebook |
| Browser | Google chrome |

## 3 Data Preparation

### 3.1 Collecting Data

Download the RSICD dataset on your local machine from GitHub using the URL https://github.com/201528014227051/RSICD_optimal. The dataset includes images folder

and a JSON file with key as image name and value as an array of text descriptions.

## 3.2   Reading Data

1. Create a new Jupyter notebook environment in the same folder where the dataset is stored.

2. Import all required libraries for data processing.

**Importing all libraries**

```
In [3]:    1  # Load libraries
           2  import os
           3  import cv2
           4  import json
           5  from PIL import Image
           6  import numpy as np
           7  from numpy import array, prod
           8  import random
           9  from collections import defaultdict
          10  from collections import OrderedDict
          11  import pickle
          12  from pickle import load, dump
          13  import tensorflow as tf
          14  from tensorflow.keras import backend as K
          15  from tensorflow.keras.applications.inception_v3 import InceptionV3
          16  from tensorflow.keras.models import Sequential, Model
          17  from tensorflow.keras.layers import Dense, Flatten,Input, Convolution2D, Dropout, LSTM, \
          18  TimeDistributed, Embedding, Bidirectional, Activation, RepeatVector, Concatenate
          19  from tensorflow.keras.layers import add
          20  from tensorflow.keras.optimizers import Adam
          21  from tensorflow.keras.callbacks import ModelCheckpoint, EarlyStopping
          22  from tensorflow.keras.preprocessing import image, sequence
          23  from keras.preprocessing.image import img_to_array, load_img
          24  from tensorflow.keras.utils import to_categorical
          25  from keras.utils.vis_utils import plot_model
          26  from tensorflow.python.keras.utils import np_utils
          27  from keras.preprocessing.sequence import pad_sequences
```

Figure 1: Import required libraries.

3. Define directory parth for images and json file to read data.

**Reading image and text data**

```
In [4]:    1  images_folder = 'RSICD/RSICD_images/'
           2  annotations_name = 'RSICD/dataset_rsicd.json'
           3
           4  images = [] # save images
           5  annotations = [] # save five sentences for every image.
           6  images_list = []
           7  train_list = [] # image index for train
           8  val_list = [] # image index for validation
           9  test_list = [] # image index for test
```

Figure 2: Defing directory paths for the data.

4. Load images into notebook using OS package of python. The get_image function defined in Figure 3 creates a global list containng names of the images in the given directory. Also, a list of images and empty list of same length is also created to store reference descriptions.

5. To load text data into notebook use JSON package of python. Split the loaded data into three parts using the 'split' key of the json for training, testing, and validation of the model. The fuction get_annotations is used to parse json data and create three lists containing index of the image in the global image_list for training, testing and validation.

2

```
In [5]:   1  def get_images():
          2      """
          3      Get names of all images to a list.
          4      """
          5      global images_list
          6      images_list = os.listdir(images_folder)
          7      image_number = len(images_list)
          8      # Create annotations_list for all images.
          9      print('Load images...')
         10      for image_name in images_list:
         11          image_name = os.path.join(images_folder, image_name)
         12          images.append(cv2.imread(image_name))
         13          annotations.append([])
         14      print(
         15          'Images path:', images_folder,
         16          '\nGot', image_number, 'images.',
         17          '\nAnnotations length:', len(annotations))
```

Figure 3: Parsing image data.

```
In [6]:   1  def get_annotations():
          2      """
          3      Get five sentences of every image in the '.json' file.
          4      """
          5      global sentence_length
          6      global train_list
          7      global val_list
          8      global test_list
          9
         10      file = open(annotations_name, encoding='utf-8')
         11      load_dict = json.load(file)
         12
         13      # The file has two element: dataset, and images, we take images.
         14      for element in load_dict:
         15          print(annotations_name, 'has element', '[', element, ']')
         16
         17      print('Getting annotations for [ images ]')
         18      for image in load_dict['images']:
         19          """
         20          # 'images' element has many key.
         21          # key filename: save the name of image's.
         22          # key sentences: save five captions for the image(filename)
         23          # key sentences['raw']: five str.
         24          # key sentences['tokens']: five list.
         25          """
         26          filename = image['filename']
         27          if filename not in images_list:
         28              continue
         29          index = images_list.index(filename)
         30          if image['split'] == 'train':
         31              train_list.append(index)
         32          elif image['split'] == 'val':
         33              val_list.append(index)
         34          else:
         35              test_list.append(index)
         36          for sentence in image['sentences']:
         37              """
         38              # In sentences['raw'], a word may be connected by a ','.
         39              # For example, 'I am fine, and you.', 'fine,' is a word.
         40              # But, we want divide 'fine,' to two words: ['fine', ',']
         41              # So, we add a ' ' into 'fine' and ','.
         42              """
         43              sentence_temp = ''
         44              for word in sentence['raw']:
         45                  if ',' in word:
         46                      word = ' ' + ','
         47                  sentence_temp += word
         48              annotations[index].append(sentence_temp)
         49
         50      print('The number of train data: ', len(train_list))
         51      print('The number of validation: ', len(val_list))
         52      print('The number of test data: ', len(test_list))
```

Figure 4: Parsing text data.

6. Figure 5 represents a sample image from the loaded RSICD dataset along with references text descriptions.

```
In [8]:   1  #view an image from the training set
          2  image = Image.open(os.path.join(images_folder, 'sparseresidential_150.jpg'))
          3  image
```

Out[8]:



```
In [9]:   1  index = images_list.index('sparseresidential_150.jpg')
          2  annotations[index]
```

Out[9]: ['a building with a swimming pool is near many green trees .',
         'a building with a swimming pool is near many green trees .',
         'a building with a swimming pool is near many green trees .',
         'a building with a swimming pool is near many green trees .',
         'a building with a swimming pool is near many green trees .']

Figure 5: Model Architecture.

# 4   Image Processing

The next step is image processing and transformation to generate image feature vectors. In image processing, loaded image are resized as required by the encoder model. Figure 6 shows the code for feature extraction including resizing and getting feature vector from the model.

```
In [26]:   1  images_features = {}
           2
           3  count = 0
           4  for i in train_list:
           5      img_name = images_list[i]
           6      img = cv2.imread(images_folder + img_name)
           7      img = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
           8      img = cv2.resize(img, (299,299))
           9
          10      img = img.reshape(1,299,299,3)
          11      pred = img_encoder.predict(img).reshape(2048,)
          12      images_features[img_name] = pred
          13
          14      count += 1
```

Figure 6: Image feature extraction process.

This research used two type of encoder model CNN and Capsule Network. For the CNN, a pre-trained InceptionV3 model is used which is avaible in keras applications [1]. For Capsule Network, a encoder model is created including 4 casule units which are

---
[1]https://keras.io/api/applications/inceptionv3/

flatten together in a layer. The architecture of the capsule network is inspired from the model proposed in the study Deka (2020), which is available on GitHub [2].

# 5 Text Processing

1. In this step, text descriptions are cleaned and transformed for modelling. For cleaning text descriptions first each word is converted into the lower case to avoid replication of words in the vocabulary and <START> and <END> tokens are added to mark the beginning and ending of a sentence. Figure 7 presents the code for this part.

**Text processing**

```
In [10]:   1  # Create a dictonary with key as name of the image and value as array of a
           2  def list_all_text():
           3      image_descriptions = defaultdict(list)
           4
           5      for index in range(len(images_list)):
           6          filename = images_list[index]
           7          sentences = annotations[index]
           8          sentence_list = []
           9          for sentence in sentences:
          10              val = sentence.split()
          11              new_val = []
          12              for word in val:
          13                  new_val.append(word.lower())
          14              new_val.insert(0, "<START>")
          15              new_val.append("<END>")
          16              sentence_list.append(new_val)
          17          image_descriptions[filename] = sentence_list
          18      return image_descriptions
```

Figure 7: Coverting to lower case and adding start and end tokens.

2. The next step is to remove punctuations from the description. Code used for removing punctuations is shown in Figure 8.

```
In [12]:   1  # Clean descriptions
           2
           3  import string
           4  des = {}
           5  for image, caps in descriptions.items():
           6      cap_list = []
           7      for cap in caps:
           8          clean_cap =[word for word in cap if word not in string.punctuation]
           9          for word in clean_cap:
          10              try:
          11                  int(word)
          12                  clean_cap.remove(word)
          13              except:
          14                  continue
          15          cap_list.append(clean_cap)
          16      des[image] = cap_list
```

Figure 8: Removing punctuations.

3. After cleaning the text descriptions, the next step is to calculate vocabulary size for modelling. The code in Figure 9 is followed for this step.

4. The last step in text processing is the tokenization of the text descriptions. A dictionary is created to convert words to numbers as showing in Figure 10.

---

[2]https://github.com/jaydeepdeka/Flickr8k_image_captioning_using_CapsNet

**Create Vocabulary**

```
In [15]:   1  count_words = {}
           2  for image, desc_list in des.items():
           3      for desc in desc_list:
           4          for word in desc:
           5              if word not in count_words:
           6                  count_words[word] = 0
           7              else:
           8                  count_words[word] += 1
```

Figure 9: Counting total words in vocabulary.

```
In [17]:   1  THRESH = -1
           2  count = 1
           3  new_dict = {}
           4  for k,v in count_words.items():
           5      if count_words[k] > THRESH:
           6          new_dict[k] = count
           7          count += 1
```

Figure 10: Creating word to number dictionary.

# 6   Modelling

1. Importing remaining libraries required for modelling.

```
In [41]:    1  from tensorflow.keras import layers, initializers, regularizers, constraints, models
            2  from tensorflow.keras.models import Sequential, Model
            3  from tensorflow.python.keras.utils import conv_utils
            4  from tensorflow.keras.layers import InputSpec
            5  from tensorflow.python.keras.utils.conv_utils import conv_output_length, normalize_data_format
            6  from tensorflow.keras.layers import Input, Conv2D, Activation, Dense, Reshape
            7  from tensorflow.keras.layers import BatchNormalization, Conv2DTranspose
            8
            9  cf = K.image_data_format() == '..'
           10  useGPU = True
```

Figure 11: Importing libraries for modelling.

2. Creating a model as shown in Figure 12, which takes image feature vector and word vector to generate a description. The architecture of this model is inspired by the study of (Agughalam et al.; 2021).

3. For the second experiment, an encoder-decoder framework uses a feature vector generated by CapsNet encoder to generate text sequences. The architecture of the model is presented in Figure 13.

4. For using CNN and CapsNet features together for generating text descriptions in experiment three, the architecture in Figure 14 is implemented. It takes three inputs, first for CNN feature vector, second for CapsNet feature vector and third for word vector.

5. In the training of the models 10 epochs are used. The model with the lowest validation loss is saved and used for generating text sequences. The code in Figure 15 is used for the training of models.
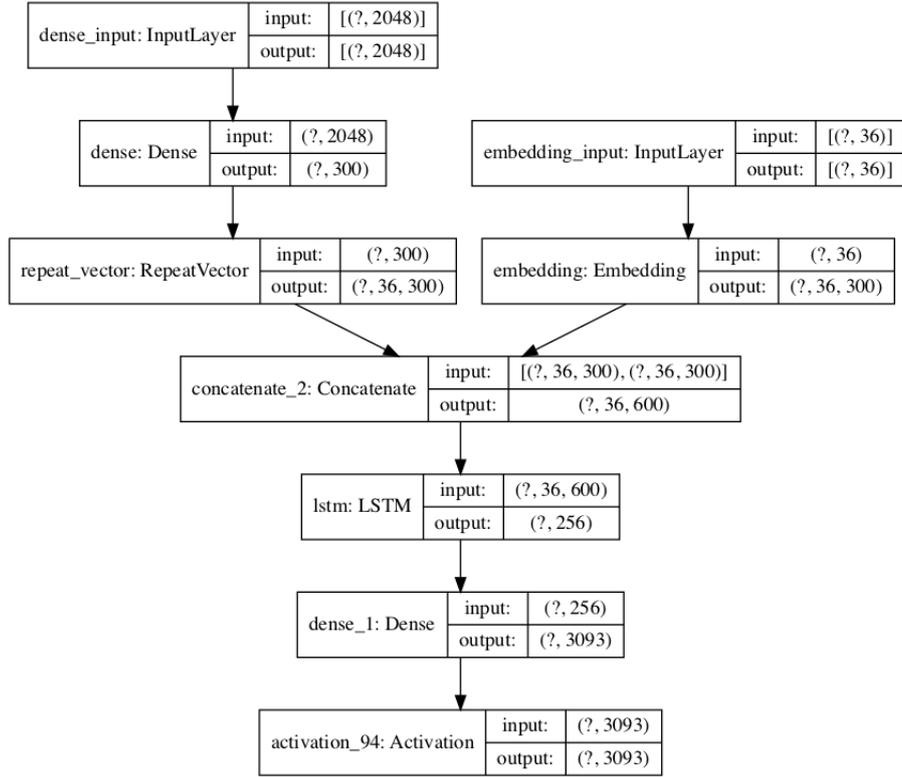
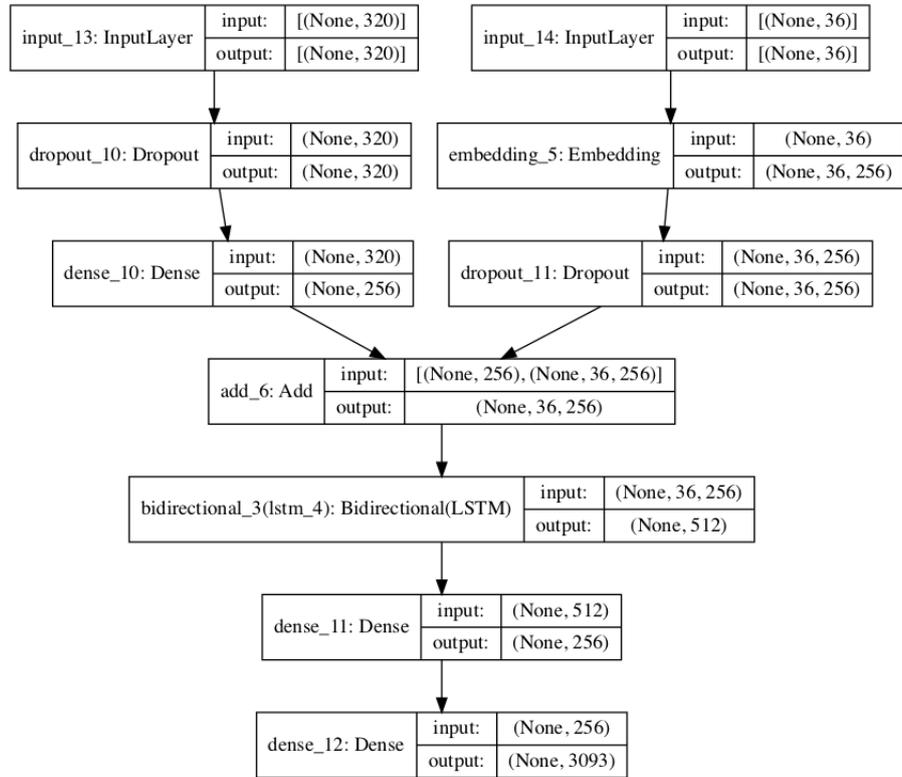Figure 12: CNN-LSTM encoder-decoder architecture.



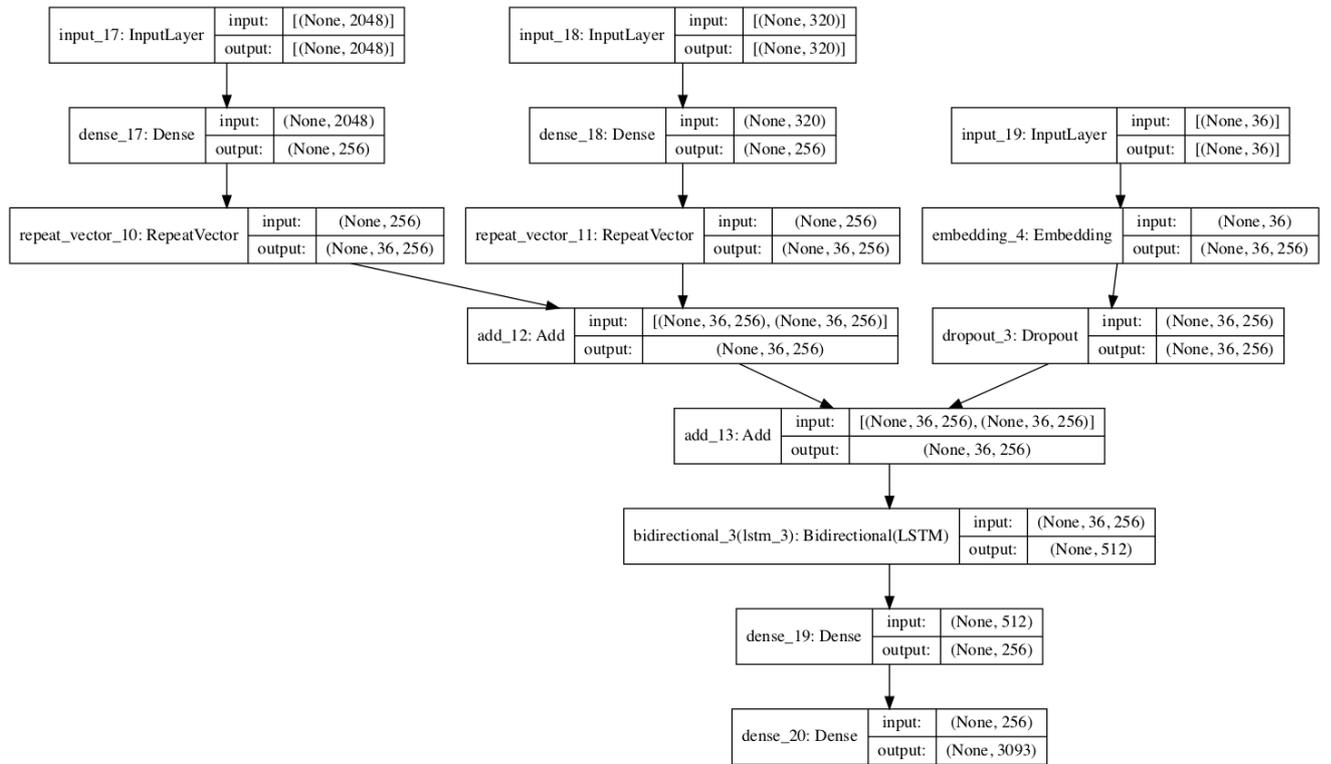Figure 13: CapsNet-BiLSTM encoder-decoder architecture.

Figure 14: CNN-CapsNet-BiLSTM encoder-decoder architecture.



```
# Train model for CapsNet
epochs = 10
train_steps = len(train_desc)
val_steps = len(val_desc)
filepath = 'model-loss{loss:.3f}-val_loss{val_loss:.3f}_Per_CapsNet_par.h5'
checkpoint = ModelCheckpoint(filepath, monitor='val_loss', verbose=1, save_best_only=True, mode='min')

# fit for one epoch
history = capsule_bilstm_model.fit(train_generator, epochs=epochs, steps_per_epoch=train_steps, \
          verbose=1, validation_data=valid_generator, validation_steps=val_steps, callbacks=[checkpoint])
```

Figure 15: Training model.

# 7 Inference and Evaluation

## 7.1 Inference

1. The function evaluation_model takes a list of test images and generate a text description for each image. First, the image is resized as required for the encoder model and the feature vector is generated. Then text sequence is predicted using that feature vector.

```python
def evaluate_model(test_list):

    actual, predicted = list(), list()

    for index in range(len(test_list)):

        test_img_path = images_folder + images_list[test_list[index]]
        test_img = cv2.imread(test_img_path)
        test_img = cv2.cvtColor(test_img, cv2.COLOR_BGR2RGB)
        test_image = cv2.resize(test_img, (299,299))
        test_image = np.reshape(test_image, (1,299,299,3))

        test_features = img_encoder.predict(test_image).reshape(1,2048)

        text_inp = ['<START>']

        count = 0
        caption = ''
        while count < MAX_LEN:
            count += 1

            encoded = []
            for i in text_inp:
                encoded.append(new_dict[i])

            encoded = [encoded]
            encoded = pad_sequences(encoded, padding='post', truncating='post', maxlen=MAX_LEN)
            prediction = np.argmax(model_1_10.predict([test_features, encoded]))
            sampled_word = inv_dict[prediction]
            caption = caption + ' ' + sampled_word
            text_inp.append(sampled_word)

            if sampled_word == '<END>':
                break

        image = images_list[test_list[index]]
        references = [d for d in reference_list[image]]
        actual.append(references)
        predicted.append(caption)

    return actual, predicted
```

Figure 16: .

2. The word with the highest probability is selected as the next word in a sentence. Like this, all words are appended together to form a sentence.

3. For all three experiments predicted text description is generated for the test images and stored in pickle files.

## 7.2 Evaluation

1. Remove <START> and <END> tokens from begining and ending of predicted descriptions for evaluation.

2. Calculate BLEU and ROUGE-L scores for each model. Figure 17 and Figure 18 presents the code for BLEU and ROUGE metrices.

3. For qualitative analysis of the research results, predicted sentences by all models are compared with respect to the content of test image. Figure 19 shown one of test image with actual and predicted sentences by each model.

```
In [44]:    1  print('BLEU-1: %f' % corpus_bleu(act_1_10, pred_1_10, weights=(1.0, 0, 0, 0)))
            2  print('BLEU-2: %f' % corpus_bleu(act_1_10, pred_1_10, weights=(0.5, 0.5, 0, 0)))
            3  print('BLEU-3: %f' % corpus_bleu(act_1_10, pred_1_10, weights=(0.3, 0.3, 0.3, 0)))
            4  print('BLEU-4: %f' % corpus_bleu(act_1_10, pred_1_10, weights=(0.25, 0.25, 0.25, 0.25)))

        BLEU-1: 0.874560
        BLEU-2: 0.766410
        BLEU-3: 0.691050
        BLEU-4: 0.586885
```

Figure 17: Model Architecture.

```
In [115]:   1  from rouge_score import rouge_scorer
            2  scorer = rouge_scorer.RougeScorer(['rouge1', 'rougeL'], use_stemmer=True)
            3
            4  def calc_rouge(reference,pred):
            5      vals = []
            6      for i in range(len(pred)):
            7          text = ' '.join(pred[i])
            8          refs = reference[i]
            9          rouges = []
           10          for ref in refs:
           11              scores = scorer.score(' '.join(ref),text)
           12              rouges.append(scores['rougeL'].fmeasure)
           13          vals.append(np.mean(rouges))
           14
           15      return np.mean(vals)
```

Figure 18: Model Architecture.

**Reference Descriptions**

1. a rectangular church with a big and four small round roofs is located on a rounded triangle red square enclosed by three roads with trees planted along
2. a rectangle church with a big and four small round roofs is located on a rounded triangle red square enclosed by three roads with trees planted along
3. the long black church with a umbelliform roof is in a delta region between two rows of trees
4. it is an i shaped church in the middle of two roads
5. three lines of green trees are near a church

**CNN (InceptionV3)-LSTM**

many buildings and green trees are in a medium residential area

**CapsNet-BiLSTM**

many green trees are around a building

Figure 19: Model Architecture.

# References

Agughalam, D., Pathak, P. and Stynes, P. (2021). Bidirectional LSTM approach to image captioning with scene features, *in* X. Jiang and H. Fujita (eds), *Thirteenth International Conference on Digital Image Processing (ICDIP 2021)*, Vol. 11878, International Society for Optics and Photonics, SPIE, pp. 81 – 88.
   **URL:** *https://doi.org/10.1117/12.2600465*

Deka, J. (2020). *Image captioning: Capsule network vs cnn approach*, Master's thesis, Dublin, National College of Ireland.
   **URL:** *http://norma.ncirl.ie/4298/*