

```
# Set the path to your zip file in Google Drive
zip_path = '/content/drive/MyDrive/image caption data /Images.zip'

# Set the path to the directory where you want to unzip the files
output_path = '/content/image data'

# Unzip the files using the zipfile module
import zipfile
with zipfile.ZipFile(zip_path, 'r') as zip_ref:
    zip_ref.extractall(output_path)

import os
import pickle
import numpy as np
from tqdm.notebook import tqdm

from tensorflow.keras.applications.vgg16 import VGG16, preprocess_input
from tensorflow.keras.preprocessing.image import load_img, img_to_array
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Model
from tensorflow.keras.utils import to_categorical, plot_model
from tensorflow.keras.layers import Input, Dense, LSTM, Embedding, Dropout, add

BASE_DIR = '/content/image data'
```

## ▼ impoting the libraries

```
import os
import pickle
import numpy as np
from tqdm.notebook import tqdm

from tensorflow.keras.applications.vgg16 import VGG16, preprocess_input
from tensorflow.keras.preprocessing.image import load_img, img_to_array
from tensorflow.keras.preprocessing.text import Tokenizer
from tensorflow.keras.preprocessing.sequence import pad_sequences
from tensorflow.keras.models import Model
from tensorflow.keras.utils import to_categorical, plot_model
from tensorflow.keras.layers import Input, Dense, LSTM, Embedding, Dropout, add
```

## ▼ extracting the image features

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

Downloading data from [https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16\\_weights\\_tf\\_dim\\_ordering\\_tf\\_kernels/553467096/553467096](https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_weights_tf_dim_ordering_tf_kernels/553467096/553467096) [=====] - 4s 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

## ▼ EXTRACTING AND PREPROCESSING THE IMAGE FEATURES

```
# extract features from image
features = {}
directory = os.path.join(BASE_DIR, 'Images')

for img_name in tqdm(os.listdir(directory)):
    # load the image from file
    img_path = directory + '/' + img_name
    image = load_img(img_path, target_size=(224, 224))
    # convert image pixels to numpy array
    image = img_to_array(image)
    # reshape data for model
    image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
    # preprocess image for vgg
    image = preprocess_input(image)
    # extract features
    feature = model.predict(image, verbose=0)
    # get image ID
    image_id = img_name.split('.')[0]
    # store feature
    features[image_id] = feature
```

100%

8091/8091 [10:50&lt;00:00, 13.58it/s]

```
WORKING_DIR = '/content/current working dir'
```

```
# store features in pickle
pickle.dump(features, open(os.path.join(WORKING_DIR, 'features.pkl'), 'wb'))
```

```
# load features from pickle
with open(os.path.join(WORKING_DIR, 'features.pkl'), 'rb') as f:
    features = pickle.load(f)
```

```
with open(os.path.join(BASE_DIR, 'captions.txt'), 'r') as f:
    next(f)
    captions_doc = f.read()
```

```
# create mapping of image to captions
mapping = {}
# process lines
for line in tqdm(captions_doc.split('\n')):
    # split the line by comma(,)
    tokens = line.split(',')
    # process tokens
```

```

if len(line) < 2:
    continue
image_id, caption = tokens[0], tokens[1:]
# remove extension from image ID
image_id = image_id.split('.')[0]
# convert caption list to string
caption = " ".join(caption)
# create list if needed
if image_id not in mapping:
    mapping[image_id] = []
# store the caption
mapping[image_id].append(caption)

```

100%

40456/40456 [00:00&lt;00:00, 263205.97it/s]

len(mapping)

8091

```

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

```

```

# before preprocess of text
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 .']

```

```

# preprocess the text
clean(mapping)

```

```

# after preprocess of text
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']

```

```

all_captions = []
for key in mapping:
    for caption in mapping[key]:
        all_captions.append(caption)

```

len(all\_captions)

40455

all\_captions[:20]

```

['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',
'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',
'startseq little girl covered in paint sits in front of painted rainbow with her hands in bowl endseq',

```

```
'startseq little girl is sitting in front of large painted rainbow endseq',
'startseq small girl in the grass plays with fingerpaints in front of white canvas with rainbow on it endseq',
'startseq there is girl with pigtails sitting in front of rainbow painting endseq',
'startseq young girl with pigtails painting outside in the grass endseq',
'startseq man lays on bench while his dog sits by him endseq',
'startseq man lays on the bench to which white dog is also tied endseq',
'startseq man sleeping on bench outside with white and black dog sitting next to him endseq',
'startseq shirtless man lies on park bench with his dog endseq',
'startseq man laying on bench holding leash of dog sitting on ground endseq']
```

```
# tokenize the text
tokenizer = Tokenizer()
tokenizer.fit_on_texts(all_captions)
vocab_size = len(tokenizer.word_index) + 1
```

```
vocab_size
```

```
8485
```

```
# get maximum length of the caption available
max_length = max(len(caption.split()) for caption in all_captions)
max_length
```

```
35
```

## ▼ Train Test Split

```
image_ids = list(mapping.keys())
split = int(len(image_ids) * 0.90)
train = image_ids[:split]
test = image_ids[split:]
```

```
# create data generator to get data in batch (avoids session crash)
def data_generator(data_keys, mapping, features, tokenizer, max_length, vocab_size, batch_size):
    # loop over images
    X1, X2, y = list(), list(), list()
    n = 0
    while 1:
        for key in data_keys:
            n += 1
            captions = mapping[key]
            # process each caption
            for caption in captions:
                # encode the sequence
                seq = tokenizer.texts_to_sequences([caption])[0]
                # split the sequence into X, y pairs
                for i in range(1, len(seq)):
                    # split into input and output pairs
                    in_seq, out_seq = seq[:i], seq[i]
                    # pad input sequence
                    in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
                    # encode output sequence
                    out_seq = to_categorical([out_seq], num_classes=vocab_size)[0]

                # store the sequences
                X1.append(features[key][0])
                X2.append(in_seq)
                y.append(out_seq)
            if n == batch_size:
                X1, X2, y = np.array(X1), np.array(X2), np.array(y)
                yield [X1, X2], y
                X1, X2, y = list(), list(), list()
                n = 0
```

## ▼ Model Creation

```
# encoder model
# image feature layers
inputs1 = Input(shape=(4096,))
fe1 = Dropout(0.4)(inputs1)
fe2 = Dense(256, activation='relu')(fe1)
# sequence feature layers
inputs2 = Input(shape=(max_length,))
se1 = Embedding(vocab_size, 256, mask_zero=True)(inputs2)
```

```

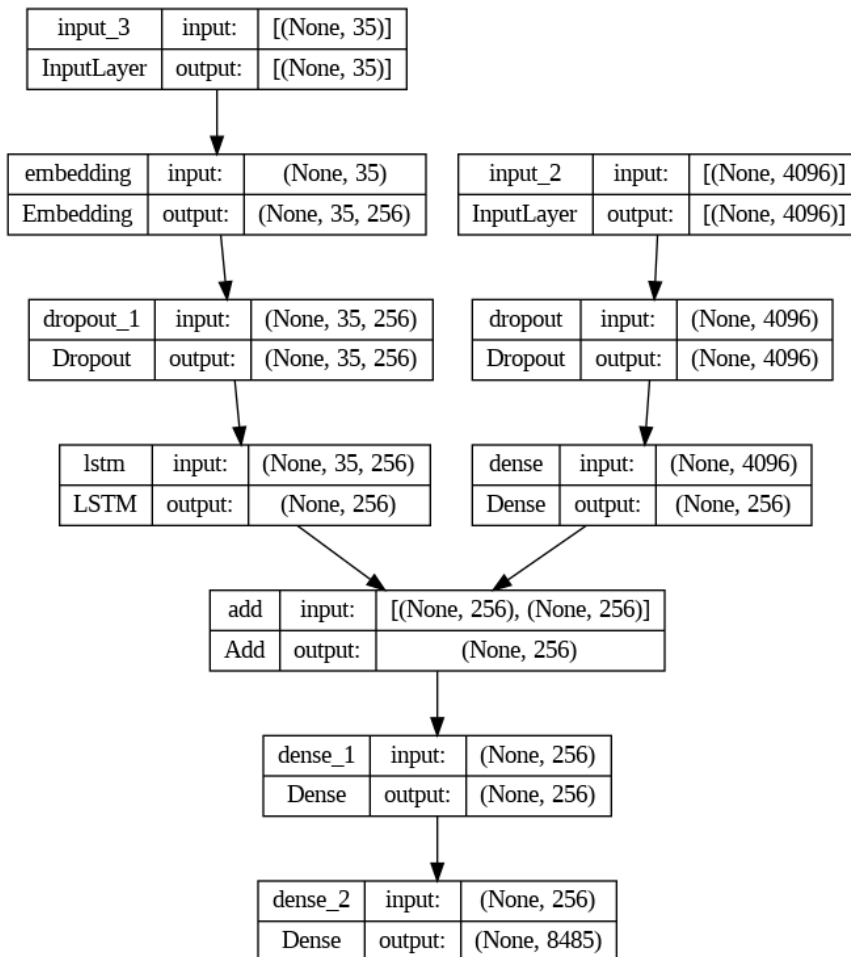
se2 = Dropout(0.4)(se1)
se3 = LSTM(256)(se2)

# decoder model
decoder1 = add([fe2, se3])
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')

# plot the model
plot_model(model, show_shapes=True)

```



```

# train the model
epochs = 20
batch_size = 32
steps = len(train) // batch_size

for i in range(epochs):
    # create data generator
    generator = data_generator(train, mapping, features, tokenizer, max_length, vocab_size, batch_size)
    # fit for one epoch
    model.fit(generator, epochs=1, steps_per_epoch=steps, verbose=1)

```

```

227/227 [=====] - 87s 357ms/step - loss: 5.2145
227/227 [=====] - 64s 282ms/step - loss: 3.9901
227/227 [=====] - 66s 289ms/step - loss: 3.5726
227/227 [=====] - 67s 294ms/step - loss: 3.3034
227/227 [=====] - 66s 293ms/step - loss: 3.1052
227/227 [=====] - 67s 295ms/step - loss: 2.9590
227/227 [=====] - 66s 290ms/step - loss: 2.8406
227/227 [=====] - 65s 285ms/step - loss: 2.7486
227/227 [=====] - 69s 302ms/step - loss: 2.6655
227/227 [=====] - 67s 297ms/step - loss: 2.5928
227/227 [=====] - 67s 296ms/step - loss: 2.5321
227/227 [=====] - 68s 297ms/step - loss: 2.4758
227/227 [=====] - 67s 295ms/step - loss: 2.4256
227/227 [=====] - 67s 296ms/step - loss: 2.3802
227/227 [=====] - 68s 299ms/step - loss: 2.3377
227/227 [=====] - 67s 296ms/step - loss: 2.2981
227/227 [=====] - 67s 295ms/step - loss: 2.2633

```

```

227/227 [=====] - 67s 293ms/step - loss: 2.2349
227/227 [=====] - 66s 289ms/step - loss: 2.2065
227/227 [=====] - 67s 297ms/step - loss: 2.1795

```

```

# save the model
model.save(WORKING_DIR+'/best_model.h5')

```

## ▼ Generate Captions for the Image

```

def idx_to_word(integer, tokenizer):
    for word, index in tokenizer.word_index.items():
        if index == integer:
            return word
    return None

```

```

# generate caption for an image
def predict_caption(model, image, tokenizer, max_length):
    # add start tag for generation process
    in_text = 'startseq'
    # iterate over the max length of sequence
    for i in range(max_length):
        # encode input sequence
        sequence = tokenizer.texts_to_sequences([in_text])[0]
        # pad the sequence
        sequence = pad_sequences([sequence], max_length)
        # predict next word
        yhat = model.predict([image, sequence], verbose=0)
        # get index with high probability
        yhat = np.argmax(yhat)
        # convert index to word
        word = idx_to_word(yhat, tokenizer)
        # stop if word not found
        if word is None:
            break
        # append word as input for generating next word
        in_text += " " + word
        # stop if we reach end tag
        if word == 'endseq':
            break

    return in_text

```

```

from nltk.translate.bleu_score import corpus_bleu
# validate with test data
actual, predicted = list(), list()

for key in tqdm(test):
    # get actual caption
    captions = mapping[key]
    # predict the caption for image
    y_pred = predict_caption(model, features[key], tokenizer, max_length)
    # split into words
    actual_captions = [caption.split() for caption in captions]
    y_pred = y_pred.split()
    # append to the list
    actual.append(actual_captions)
    predicted.append(y_pred)

# calculate BLEU score
print("BLEU-1: %f" % corpus_bleu(actual, predicted, weights=(1.0, 0, 0, 0)))
print("BLEU-2: %f" % corpus_bleu(actual, predicted, weights=(0.5, 0.5, 0, 0)))

```

```

100%
810/810 [08:42<00:00, 2.00it/s]
BLEU-1: 0.542755
BLEU-2: 0.316532

```

## ▼ Visualize the Results

```

from PIL import Image
import matplotlib.pyplot as plt
def generate_caption(image_name):
    # load the image
    # image_name = "1001773457_577c3a7d70.jpg"
    image_id = image_name.split('.')[0]

```

```

img_path = os.path.join(BASE_DIR, "Images", image_name)
image = Image.open(img_path)
captions = mapping[image_id]
print('-----Actual-----')
for caption in captions:
    print(caption)
# predict the caption
y_pred = predict_caption(model, features[image_id], tokenizer, max_length)
print('-----Predicted-----')
print(y_pred)
plt.imshow(image)

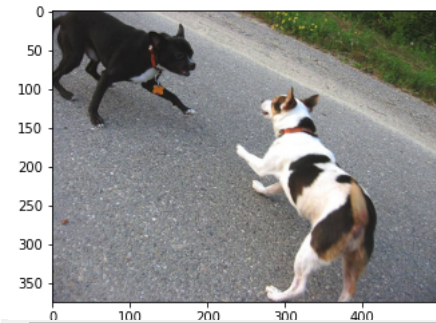
```

```
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 s
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 two dogs are playing tug of war with red toy endseq

```



```
generate_caption("1002674143_1b742ab4b8.jpg")
```

```

-----Actual-----
startseq little girl covered in paint sits in front of painted rainbow with her hands
startseq little girl is sitting in front of large painted rainbow endseq
startseq small girl in the grass plays with fingerpaints in front of white canvas wit
startseq there is girl with pigtails sitting in front of rainbow painting endseq
startseq young girl with pigtails painting outside in the grass endseq
-----Predicted-----
startseq two children and little girl in pink outfits are playing in the snow endseq

```

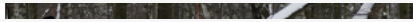


```
generate_caption("101669240_b2d3e7f17b.jpg")
```

```
-----Actual-----  
startseq man in hat is displaying pictures next to skier in blue hat endseq
```

## ▼ Test with Real Image

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



```
image_path = '/kaggle/input/flickr8k/Images/1000268201_693b08cb0e.jpg'  
# load image  
image = load_img(image_path, target_size=(224, 224))  
# convert image pixels to numpy array  
image = img_to_array(image)  
# reshape data for model  
image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))  
# preprocess image for vgg  
image = preprocess_input(image)  
# extract features  
feature = vgg_model.predict(image, verbose=0)  
# predict from the trained model  
predict_caption(model, feature, tokenizer, max_length)
```

9s completed at 9:56 AM

