

EEE543 Final Project Report

Doruk akmakçı 21502293

Furkan Özden 21502397

Mert Albaba 21501950

Tahir Ahmet Gölge 21501627

Elif Engin 21302193

Bilkent University — March 7, 2020

Abstract

Image captioning is a popular research area in recent years. In this project, an image captioning model is presented as the term project of EEE543 course at Bilkent University. Our architecture, which is inspired by state-of-art methods in the literature, combines deep, pre-trained convolutional neural networks and long short term memory networks in order to generate captions of given figures. As convolutional networks, InceptionV3 and VGG-16 are used to encode images. For both, an LSTM network is utilized to produce captions from given image embeddings. Proposed model is accurate in most of the images. In order to show the performance of the model, test results are presented for different images for our best working model with InceptionV3.

1 Introduction

Human may easily understand and describe the scene, or images, with a quick glance [1]. To make artificially intelligent systems have the same visual capability, there is a huge amount of research and interest in this area. Previously, image classification was the leading research area in this field. In order to find features and classify images correctly, several architectures have been proposed as a result of the ImageNet challenge. Stacked convolutional layers along with pooling layers are proposed as an architecture in [2] and became the baseline model in this area. In [3], rectified linear units are introduced to a convolutional network. Much deeper networks are proposed and utilized in [4]. In [5], a deep network that stacks blocks containing convolutional layers instead of stacking layers is proposed and achieved a great success.

In recent years, the main field of research in image understanding area is generation of descriptions of images. Primitive works, [6] and [7], in this area utilized hard coded mechanism to generate descriptions of images. On the other hand, most of the recently proposed architectures utilizes convolutional networks or blocks for image feature extraction, which have been developed for classification of images. In [8], an architecture that contains two sub-architectures for different purposes are proposed. The main goal of the first of these sub-architectures is detecting sets of words found in images and this architecture has utilized deep convolutional networks. Second sub-architecture utilizes a maximum entropy language model for generating sentences. Object recognizing deep convolutional networks are used as feature extractors in [9], which are explained before. For sentence generation and as a language model, a recurrent neural network is utilized in this work. In [10], previous two works are compared and advanced with the usage of BLEU scores on a large dataset. In [11], deep convolutional neural networks are combined with recurrent neural networks for generating image descriptions, i.e. captions. In this work, a deep convolutional neural network is used as an image encoder, which is pretrained for ImageNet challenge, and outputs of this deep network is used as inputs to caption generating architecture. To generate captions, a long short term memory networks are utilized. This architecture is simply presented in Figure 1 [11].

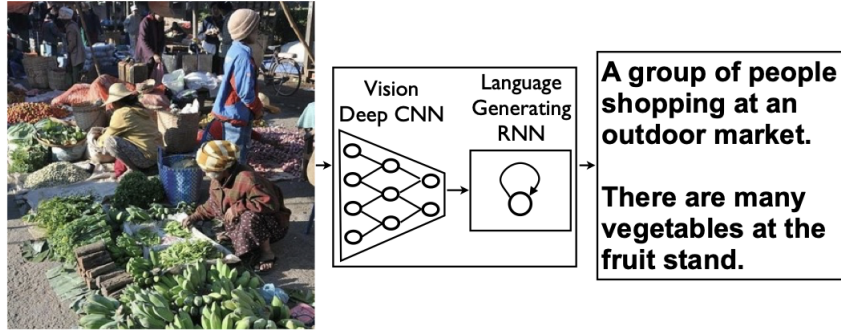


Figure 1: Simplified architecture used in [11].

In this project, an image captioning, i.e. description, program is developed. Dataset used in this work contains images and is provided by our professor and explained in Section 2. As explained before, [11] proposes an end-to-end system and achieved great BLEU scores, for instance 59 in Pascal dataset, thus in this project, we have used a very similar approach to [11]. For generating captions, an LSTM network is utilized. For image encoder convolutional neural network, we have tried two different pretrained networks: VGG16 [4] and Inception V3 [5]. Utilized architecture is detailly explained in Section 3 and for both of these, results are presented in Section 4. Codes for both architecture can be found in Section 5.

2 Methods

2.1 Dataset Description

Dataset used for this project contains real life images portraying interactions between different objects in an environment or interactions between an object and an environment. The objects of interest can vary from a human to a teddy bear. The environment can vary from a natural scene to indoors of a household. The wide array of object-environment combinations present in the dataset preserve an important situational information. Our goal in this project is to generate a model trained on the dataset which can describe the scene in an image

The dataset contains 82783 URLs from an FTP server with images and a total of 400135 captions for the images. Images may have more than one caption. Only 73,659 of the URLs given were valid. Hence, the number of images present in the dataset decreased to 73,659 and the number of captions decreased correspondingly. Captions are combination of words in a vocabulary of size 1004. Maximum length of a caption in the dataset is 17.

Glove word embedding is used to project vocabulary from word space to a vector space. Glove is an unsupervised learning algorithm for obtaining vector representation of words [12]. We used 200-dimensional vector representations of a vocabulary of size 400K, pretrained on Wikipedia. We used the subset of vector representations corresponding to our vocabulary. Nine words in our vocabulary were not present in the pretrained vectors. We set the word embeddings of missing words to 200-dimensional zero vectors. Word embedding is discussed in detail in the section InceptionV3.

After generating word embeddings of the vocabulary, we used t-SNE to visualize vocabulary space. t-SNE reduces the dimensionality of the word embedding space for visualization. In the leftmost figure below, we can see that words related to color are clustered. In the middle figure below, we can see that words related to cooking and food ingredients are clustered. In the rightmost figure below, we can see the missing words, which are set to zero vectors, are clustered. This visualization is a validation of the projection of semantic space to a space where distance between points represent semantic closeness.

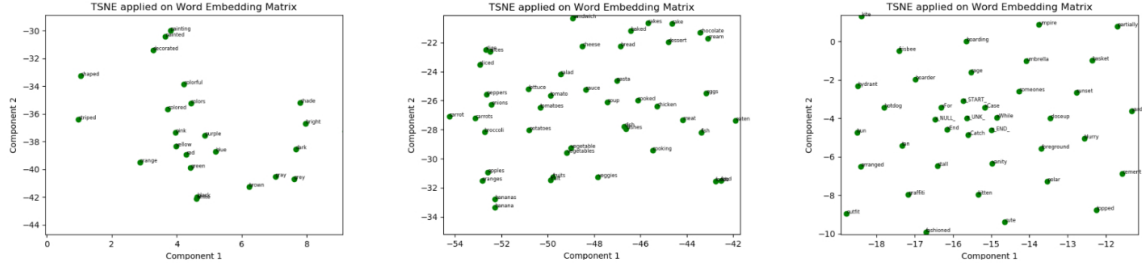


Figure 2: t-SNE visualizations of word embedding space for example clusters

2.2 Training Methodology

It is very well known that CNN's can encode images to a subspace that can efficiently encode the given image with its features. This is crucial since the images are embedded by the CNN's to a fixed length vector, that can later be used for visual tasks such as image captioning. Thus, it is very common to use a pre-trained CNN classifier with its final softmax layer popped out in order to encode the images efficiently for image captioning task. Our approach to the problem is as follows. Since the task includes both a computer vision and a natural language processing part, we combined two state-of-the-art models with their pre-trained weights for encoding purposes. After encoding processes are done, dropout layers are used to prevent overfitting. The encoded features are then decoded with dense layers with learnable parameters. In the final layer, our model chooses the next word in the caption sequence with a softmax activated dense layer with proper dimensions.

We chose Inception V3 [5] as an image feature extractor for our architecture. For encoding of the caption sequences, we used GLOVE pre-trained word embedding matrix. With these extracted features as inputs to our model architecture, our model design is summarized in the figures below.

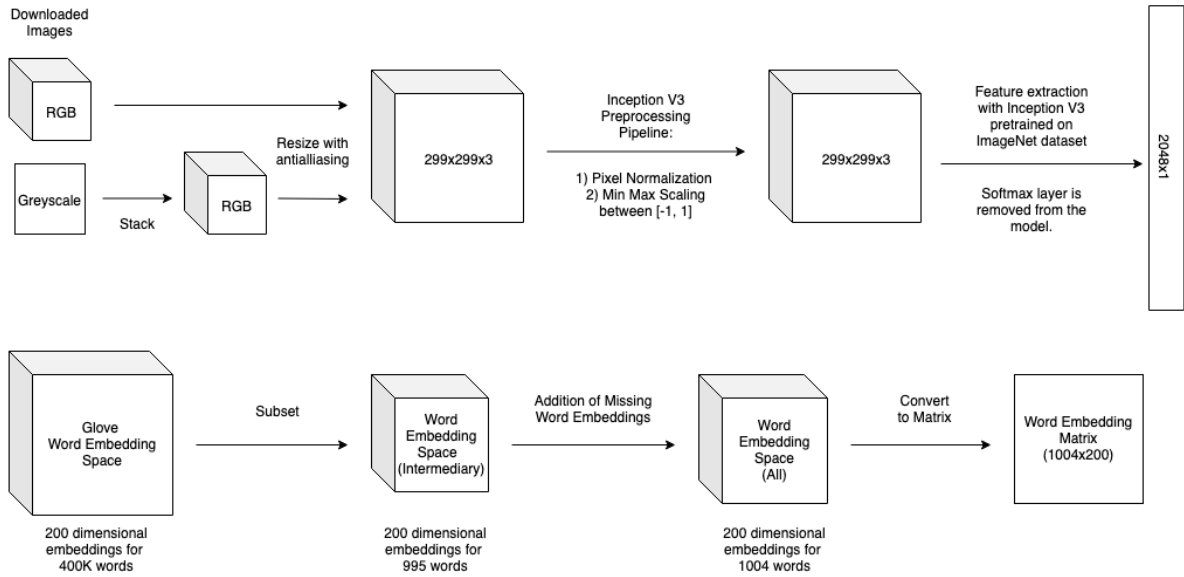


Figure 3: Data pre-processing with pre-trained models

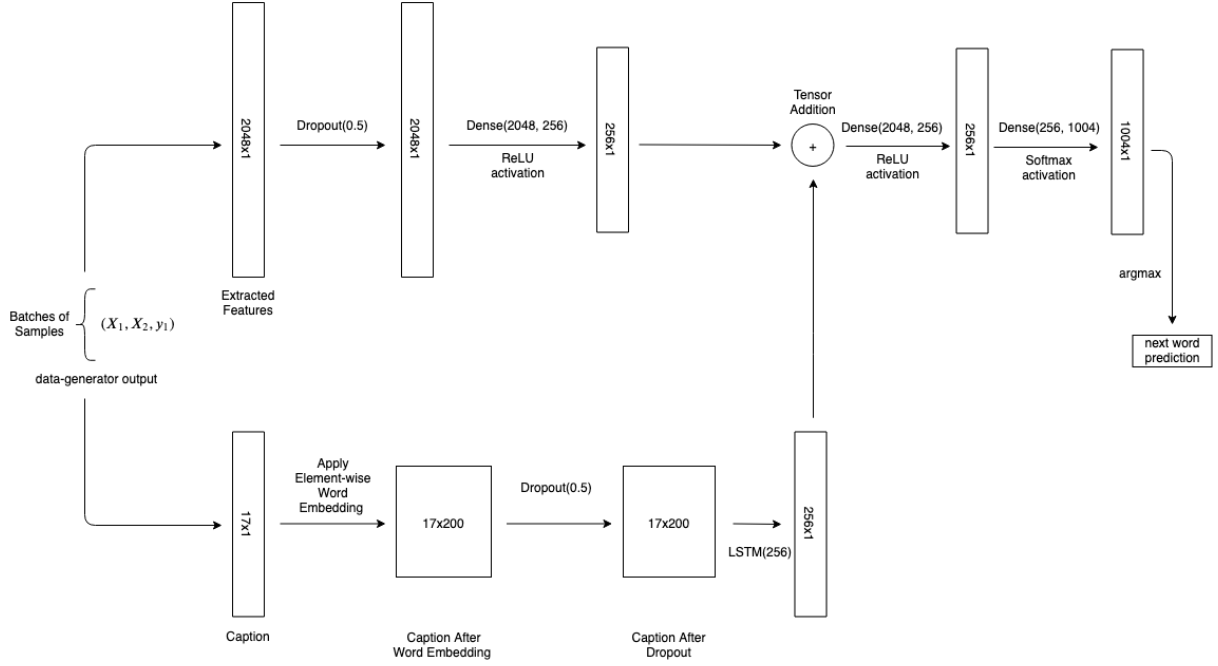


Figure 4: Model architecture

Although it can be seen from the figures above, we also discuss the model training pipeline in detail with the following items.

- Process all images in the dataset such that each of them is of shape $299 \times 299 \times 3$ (i.e, the input shape of InceptionV3 pre-trained model). If any image is gray-scale, we add another dummy axis to convert it to 3-channel rgb image. The images are pre-processed with reshaping using keras pre-process input method.
- After pre-processing step, we normalize each image between $[-1, 1]$ range for better feature extraction.
- All images are then encoded using InceptionV3 with final shape of 2048×1 . The encodings of images are stored separately.
- Our captions are in the form of integer sequences as it is explained in the Dataset Description section. First off, using GLOVE, we extracted pre-trained word embeddings for each word in the vocabulary of the given dataset. These word embeddings are also stored separately. We used GLOVE with embedding size 200. Some of the words in the vocabulary were not present in the GLOVE vocabulary. These words are:

- x_START_
- x_END_
- x_UNK_
- x_NULL_
- xFor
- xWhile
- xCase
- xCatch
- xEnd

The first four words that are not present in the GLOVE vocabulary are understandable since they are tokens for start of sequence, end of sequence, out of vocabulary word (i.e, o.o.v) and sequence padding encodings respectively. However all of the other words that are not present in the GLOVE dictionary are not of big importance in terms of captioning a given image. All of the o.o.v images are encoded with all zero vectors of size 200 (i.e, embedding size).

- With the above feature extractions are performed and saved, we feed them to our data-generator function.
- The data-generator function takes an image and a captioning sequence, it generates 3-tuples in the form $(X1, X2, Y1)$, where $X1$ is always the input image encoding. $X2$ is a fixed length integer encoding of a caption sequence. For an input image features, $X1$, assume that its caption sentence has n words. The generator function generates sequences with first i words from the caption sentence of $X1$ where $i \in [0, n]$ with $n < 17$ where 17 is the maximum length of captions in the dataset. The generated sequences are then padded with 0's (i.e, integer encoding of `x_NULL_`) for fixed length mapping. For each of the generated $X2$ sequence with first i words are taken from the original caption sequence, Y is assigned as the target variable with the next word after the last word in the $X2$ sequence, that is $(i + 1)$ 'st word. Hence, the data-generator function produces n training examples from an image with caption length n .
- The training samples generated by the data-generator function are then fed to our model. Hence, our model has two inputs and a single output word prediction.
- The encoded image features, $X1$'s, are directly fed to a Dropout layer with drop rate of 0.5 to prevent over-fitting. The drop rate choice of 0.5 is result of trial-error search. Different drop rates are tried.
- After the Dropout layer, we feed the outputs to a Dense layer with 256 neurons, each of which has ReLU activation, for further caption-specific feature extraction.
- The integer caption encodings, $X2$'s, are directly fed to the pre-trained and weight-frozen Embedding layer to generate corresponding embedding matrix of the sequence.
- Obtained embeddings are fed to a Dropout layer with drop rate of 0.5 to prevent over-fitting. The drop rate choice is again result of a trial-error search.
- After the Dropout layer, obtained outputs are fed to a LSTM layer with dimensionality of the output space 256, to match with the output dimensionality of the feature extractions of the image input.
- The 256 dimensional sequence features obtained from LSTM layer are then added element-wise to the 256 dimensional image features obtained from the dense layer. Let's call the result of this addition as total features.
- Total features are fed to a dense layer with 256 neurons each of which has ReLU activation again, to learn the mapping between prediction layer and total features of the training sample.
- Finally, the output is fed to a dense layer with 1004 neurons with softmax activation in order to make prediction of the next word in the captioning sequence, namely Y .
- Our architecture combines categorical cross-entropy loss with the softmax outputs since the problem is now expressed as a supervised classification problem. The initial learning rate of the model is 0.001. We used Adam optimizer [13] for optimization. We used a batch size of 6 images initially. Notice that, batch size of 6 images does not mean 6 training examples. Since each image has multiple captions with each length n caption having n training samples.
- Also we realized that model gets into a plateau after 4 epochs where the loss does not significantly decrease over 1 full epoch. Thus, we manually set learning rate of the model to 0.0001 after 4 epochs and increased the batch size to 12 for better loss function approximation. After 10 epochs, we set learning rate to 0.00001 and increased the batch size to 18 with a similar reasoning.
- For validation and test measures, we splitted the dataset with 70%, 20%, 10% ratios for training, validation and testing respectively. We used image count while performing the split.

2.3 Training Hardware, Software Specifications & Training Times

The described model is implemented in Python language, with Keras Deep Learning framework. We have used Pillow library in order to process, save, load images and for BLEU score metric calculations we have used NLTK library. The model is then trained on a above average laptop computer. The specifications of the computer is listed below:

- Model : Dell XPS9570
- Processor : Intel Core i7 8th Gen. 8750H 2.2GHz
- GPU : NVIDIA Geforce GTX1050Ti
- RAM : 8GB

We have used GPU support for training, with the GPU specified above. We trained the model with VGG-16 as feature extractor on 5 epochs with a total training time of ~ 1 hour. We chose the model with InceptionV3 as our best model and we trained it 32 epochs with a total training time of ~ 8 hours. The increase in batch size has a positive effect on training time, the batch size specifications are discussed in the above subsection.

3 Results

In this section we report following performance measures regarding our models:

- Learning curves (i.e, validation and loss graphs)
- Sample test caption predictions from our model
- BLEU score [14]

As we have already mentioned in the Introduction section, we have tried two different major architectures to the problem. The two approaches are listed below:

- CNN + LSTM architecture with CNN transfer learning using VGG-16
- CNN + LSTM architecture with CNN transfer learning using InceptionV3

Both of the above approaches had identical inner layers with different input shapes to the models. Since VGG-16 extracts features with shape 4096×1 and InceptionV3 extracts features with shape 2048×1 . In subsection Results VGG-16, we share the results regarding the model with VGG-16 as a feature extractor and in subsection Results InceptionV3, we share the results of the model with InceptionV3 as a feature extractor. Note that, in each of the models GLOVE word embedding is used.

3.1 Results InceptionV3

The training procedure of the model is discussed in detail in the Methods section. We chose this model as best model from our model inventory. Hence, we examined this models output in great detail.

3.1.1 Learning curve and BLEU score

Figure below shows the learning curves of the model. Since the dataset is quite big with around 72000 images, batch-wise loss decrease is a more smooth decrease. However we present here average of batch losses over 1 full epoch. The graph below illustrates the validation and training losses separately.

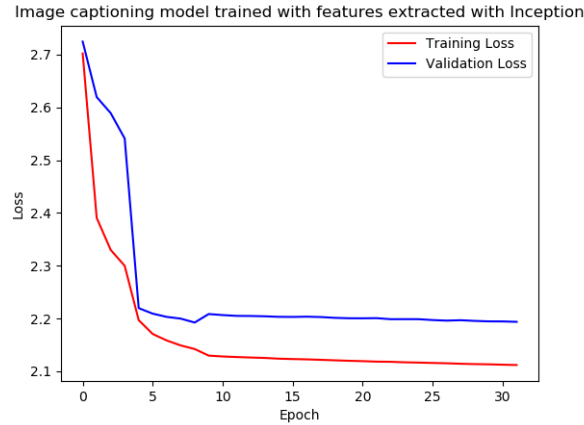


Figure 5: Learning curve of model with InceptionV3

As we have mentioned in the Methods section, we perform learning rate manipulations at epochs 4 and 10. The effects of these manipulations can be seen from the loss graphs. We also calculated BLEU score of our model using 10% of the given dataset. The BLEU score is calculated as 28.83.

3.1.2 Selected caption predictions from test set

We present some selected caption predictions from test set with their real captions in the below grid figures. We can see that model learned to generalize caption generation task very well in overall. It can recognize variety of actions along with environmental effects. It uses `x_UNK_` token very appropriately when generating a caption.

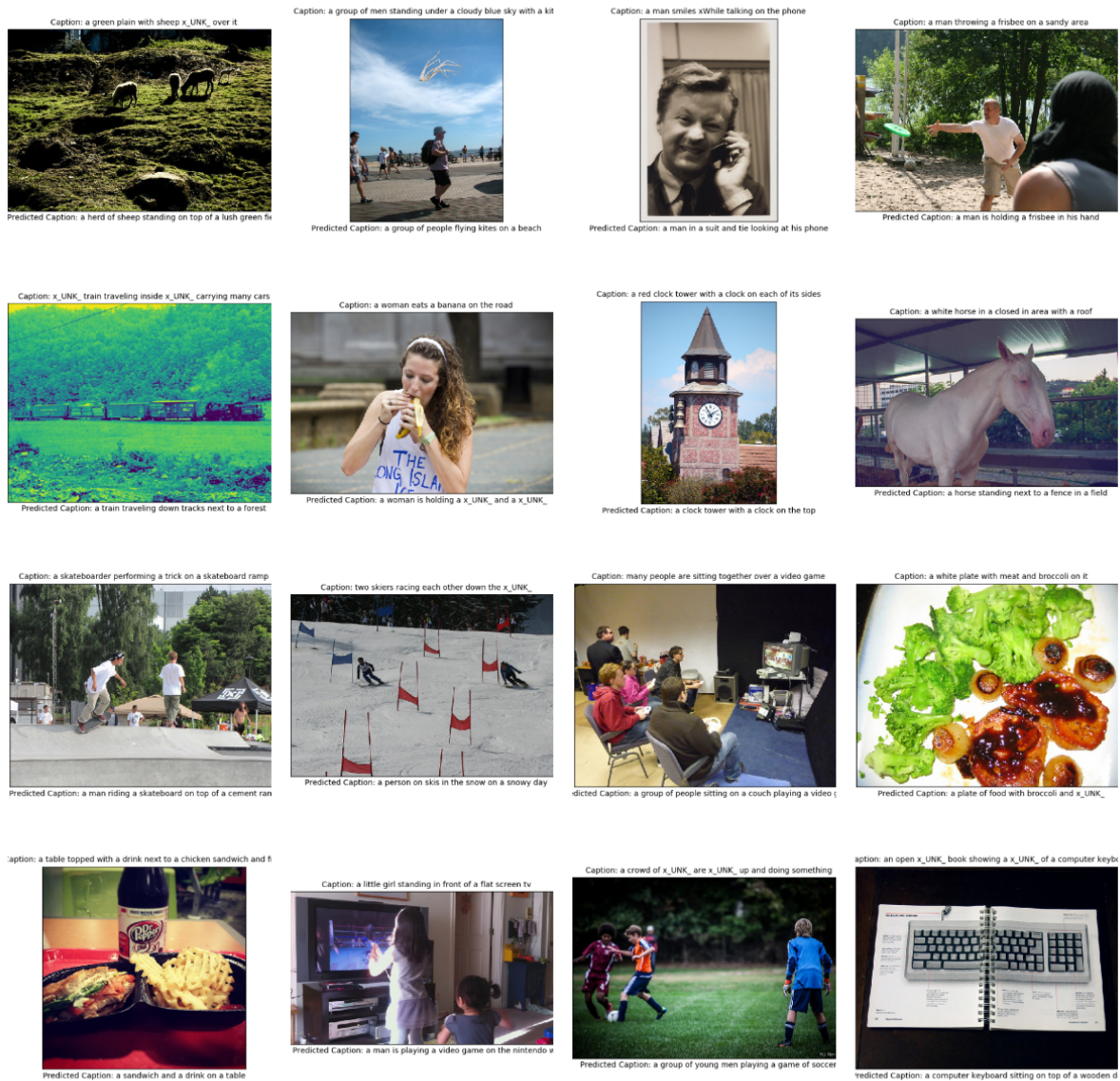
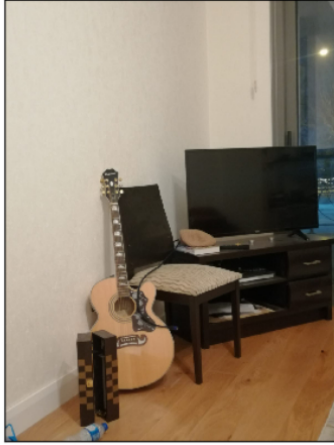


Figure 6: Some selected caption predictions of our model, for display purposes we took the first label caption of each image (i.e, images have multiple label captions we took first one)

3.1.3 Real life caption predictions (images are taken by our team members)

In this subsection, we wanted to test our model with pictures that are taken by our team members. Below pictures are taken using our cell-phone's camera, and tested with the final weights of our model. The predictions somewhat resembled to the real scenes in the photos.



Predicted Caption: a living room with a television and a television



Predicted Caption: a man sitting on a couch with a laptop



Predicted Caption: a car is parked in front of a x_UNK_



Predicted Caption: a man is sitting in front of a tv

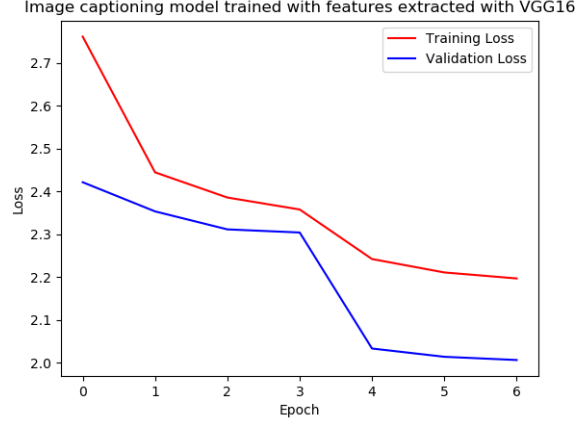
Figure 7: Caption predictions of photos taken by our team

3.2 Results VGG-16

In this sub-section we briefly summarize results obtained from the model uses VGG-16 as feature extractor.

3.2.1 Learning curve and BLEU score

The learning curve regarding the model with VGG-16 is given in the below figure. We calculated BLEU score in the test set also with this model. The BLEU score is calculated as 26. Because of faster training, fewer model parameters we chose to train the model with InceptionV3 as our main model. We trained VGG-16 only for 6 epochs. This model has a lot more learnable parameters when compared to the model with InceptionV3 as a feature extractor because VGG-16 gives features of images with 4096 dimensional vectors which drastically increases parameters required in the architecture.



Learning curve of model with VGG-16

4 Discussion

In this work, an image captioning algorithm, which is a combination of deep pre-trained convolutional neural network and long short term memory network, is proposed and implemented. As explained before, there is a huge amount of interest to this area, generation of image descriptions. Our model has been inspired by one of the state-of-art methods in the literature, [11], which also combines an LSTM network with a deep CNN, Inception V3. In [11], it is stated that the proposed model, which is very similar to our model, has achieved 63 and 66 as BLEU scores in Flickr 8k and Flickr 30k datasets, respectively. In this work, BLEU scores are also calculated as 26 and 28.83, for both of the proposed models, model with VGG16 and with Inception V3, respectively. In order to understand the reason of the difference between BLEU scores of [11] and our Inception V3 model, we have analyzed our dataset first and compared it with Flickr 8k dataset. According to our analysis, number and variety of figures in our dataset are enough, however, some captions are very short. As the result of the comparison with Flickr 8k dataset, it is deduced that although the dataset used in this work has more samples than the Flickr 8k dataset, captions are very limited. In other words, vocabulary of the dataset used is very limited when compared with the Flickr 8k dataset. To clarify, our dataset has a vocabulary formed by 1004 words and Flickr 8k dataset's vocabulary contains more than 4000 words. Furthermore, captions in our dataset are generally shorter than the ones in Flickr 8k dataset. To exemplify, the maximum caption length in Flickr 8k dataset was 34 and it is 17 in our dataset. Thus, we believe that the main reason of the difference in BLEU scores is captions in our dataset. In order to understand if Inception V3 fails in our dataset, VGG 16 is also utilized, however, produced worse results in terms of BLEU scores. Also along the training process of model with InceptionV3, we realized that the model was able to generate object names but without a proper natural language. As a solution, we added Dropout layers after the image feature extraction part, we increased the drop rate iteratively by trial-error search in order to achieve most robust results.

In order to reach a higher BLEU scores and produce better captions, image descriptions, i.e. captions, provided in our dataset may be enriched. Furthermore, this dataset may be combined with another dataset containing large vocabulary to generate better captions. Lastly, attention mechanism may also be introduced to our caption generating network in order to achieve better results.

All in all, the main purpose stated in the introduction in this work is generating meaningful captions for the images in the provided dataset. As presented results validates, we have reached the goals we have set and generated meaningful captions for the most of the images in the dataset. Thus, our proposed model generated the expected outputs, which are meaningful and accurate captions.

5 Appendix

Please note that, for easy running of all the codeblocks, we uploaded a .zip file containing a directory structure. We also included a README file which explains all of the code design structure of the project. One can run the main training file discriptor.py directly from the directory that was downloaded. All result

generating codes are included in the .zip file, explained in the README file and also given in the below listings.

5.1 Codes used for both models, result generating & dictionary preparing codes

```

1 import pickle
2 import os
3 from nltk.translate.bleu_score import corpus_bleu
4
5 # test dataset (10% of the data)
6 with open("./data/dataset/test_dataset_inception.pkl", "rb") as f:
7     test_features = pickle.load(f)
8     test_encoded_images = preprocess_dataset(test_features)
9 with open("./data/dataset/test_captions_inception.pkl", "rb") as f:
10     test_captions = pickle.load(f)
11     test_dataset = test_captions
12
13 # relevant dictionaries for word conversion
14 with open("./utils/word_to_idx.pkl", "rb") as f:
15     word_to_idx = pickle.load(f)
16 with open("./utils/idx_to_word.pkl", "rb") as f:
17     idx_to_word = pickle.load(f)
18
19 # other
20 inception = InceptionV3(weights='imagenet', include_top=True)
21 inception = Model(inception.input, inception.layers[-2].output)
22 model = load_model('./model_weights/model_inception_epoch1.h5')
23 max_length = 17
24 corpus_bleu_score = 0
25 references = []
26 candidates = []
27 for key in test_captions.keys():
28     candidate = ['x_START_']
29     reference = test_captions[key]
30     for i in range(max_length):
31         sequence = [word_to_idx[w] for w in candidate if w in word_to_idx.keys()]
32         sequence = pad_sequences([sequence], maxlen=max_length)
33         yhat = model.predict([img, sequence], verbose=0)
34         yhat = np.argmax(yhat)
35         word = idx_to_word[str(yhat)]
36         candidate.append(word)
37         if word == 'x_END_':
38             while len(candidate) != 17:
39                 candidate.append('x_NULL_')
40             break
41     candidates.append(candidate)
42     references.append(reference)
43
44 bleu = corpus_bleu(references, candidates)
45 print("Bleu Score: ", bleu)

```

```

1 import pickle
2 import h5py
3 import os
4 import pdb
5
6 import numpy as np
7
8 from gensim.models import Word2Vec
9
10 # input is a list of integers or words, and corresponding lookup table
11 def convert_caption(caption, dct):
12     result = []
13     for item in caption:
14         result.append(dct[str(item)])
15     return result
16
17 # import dataset
18 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
19     captions = f['train_cap'][:]

```

```

20 word_codes = f['word_code'][:]
21 image_ids = f['train_imid'][:]
22
23 vocabulary_size = len(word_codes[0].dtype)
24
25 # create dictionaries for conversion between words and corresponding indices
26 idx_to_word = {}
27 word_to_idx = {}
28 for i in range(vocabulary_size):
29     idx_to_word[str(int(word_codes[0][i]))] = word_codes[0].dtype.descr[i][0]
30     word_to_idx[word_codes[0].dtype.descr[i][0]] = str(int(word_codes[0][i]))
31
32 pdb.set_trace()
33
34 # save all dictionaries
35 with open('./utils/word_to_idx.pkl', 'wb') as f:
36     pickle.dump(word_to_idx, f)
37 with open('./utils/idx_to_word.pkl', 'wb') as f:
38     pickle.dump(idx_to_word, f)

```



```

1 import h5py
2 import os
3 import pdb
4 import urllib.request
5
6 import numpy as np
7
8 from multiprocessing.pool import ThreadPool
9
10
11
12 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
13     captions = f['train_cap'][:]
14     image_ids = f['train_imid'][:]
15     images = f['train_ims'][:]
16     urls = f['train_url'][:]
17     word_codes = f['word_code'][:]
18
19 print("Dataset Fields")
20 print("Captions: ", captions.shape)
21 print("Image IDs: ", image_ids.shape)
22 print("Images: ", images.shape)
23 print("URLs: ", urls.shape)
24 print("Word Codes: ", word_codes.shape)
25 print("Min image id:", np.min(image_ids))
26 print("Max image id:", np.max(image_ids))
27
28 dataset_path = "./data/dataset/"
29 images_path = os.path.join(dataset_path, "images/")
30
31 # first create a directory named "dataset/images/" to be used as download destination
32 # uncomment to download images
33 # error_count = 0
34 urls = list(enumerate(urls))
35 def download(idx, url):
36     try:
37         urllib.request.urlretrieve(url.decode('UTF-8'), os.path.join('./data/dataset/
38 images/', str(idx + 1) + ".jpg"))
39     except urllib.request.HTTPError:
40         pass
41 with ThreadPool(20) as p:
42     p.starmap(download, urls)
43
44 # for idx, url in enumerate(urls):
45 #     try:
46 #         urllib.request.urlretrieve(url.decode('UTF-8'), os.path.join(images_path, str(
47 idx) + ".jpg"))
48 #     except urllib.request.HTTPError as e:
49 #         error_count += 1
50 #         print("Error occured at idx: ", idx, " with source: ", url)
51 #         print("Error Status Code: ", e.code)

```

```

51 #         print("Error Message: ", e.read())
52
53 # print(error_count, "images were not downloaded due to errors.")

1 import pdb
2 import h5py
3 import pickle
4
5 import numpy as np
6
7 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
8     captions = f['train_cap'][:]
9     image_ids = f['train_imid'][:]
10    images = f['train_ims'][:]
11    urls = f['train_url'][:]
12    word_codes = f['word_code'][:]
13
14 glove_embedding_dict = {}
15 with open("./utils/glove.6B.200d.txt", 'r', encoding="utf-8") as f:
16     for line in f:
17         values = line.split()
18         word = values[0]
19         vector = np.asarray(values[1:], "float32")
20         glove_embedding_dict[word] = vector
21
22 # check if all vocabulary is present on glove embeddings
23 with open("./utils/word_to_idx.pkl", "rb") as f:
24     word_to_idx = pickle.load(f)
25     vocabulary = list(word_to_idx.keys())
26
27 embedding_dict = {}
28 error_cnt = 0
29 successful_cnt = 0
30 for word in vocabulary:
31     try:
32         # print(word, "\t", glove_embedding_dict[word].shape)
33         embedding_dict[word] = glove_embedding_dict[word]
34         successful_cnt += 1
35     except:
36         print("glove embedding not found for: ", word)
37         print("replaced with zeros")
38         embedding_dict[word] = np.zeros(200)
39         error_cnt += 1
40
41
42 print("Error count: ", error_cnt)
43 print("Successful count: ", successful_cnt)
44
45
46 embedding_matrix_size = len(list(embedding_dict.keys()))
47 embedding_matrix = np.zeros((embedding_matrix_size, 200))
48 for word in embedding_dict.keys():
49     embedding_matrix[int(word_to_idx[word]), :] = embedding_dict[word]
50
51 # set embedding of unknown tag as average
52 embedding_matrix[int(word_to_idx["x_UNK_"]), :] = np.mean(np.array(list(embedding_dict.values())), axis=0)
53
54 np.save("./utils/embedding_matrix.npy", embedding_matrix)

1
2 import pickle
3 from keras.preprocessing.sequence import pad_sequences
4 import numpy as np
5 import pdb
6 import random
7 from keras.applications.inception_v3 import InceptionV3, preprocess_input
8 import sys
9 from PIL import Image
10 from keras.models import load_model
11 from keras.models import Model
12 import os

```

```

13 import matplotlib.pyplot as plt
14
15 # import conversion dictionaries
16 with open("./utils/word_to_idx.pkl", "rb") as f:
17     word_to_idx = pickle.load(f)
18 with open("./utils/idx_to_word.pkl", "rb") as f:
19     idx_to_word = pickle.load(f)
20
21 # all dataset
22 with open("./data/dataset/processed_dataset_inception.pkl", "rb") as f:
23     features = pickle.load(f)
24 with open("./utils/captions.pkl", "rb") as f:
25     captions = pickle.load(f)
26
27 count = 0
28 for key in captions.keys():
29     for capt in captions[key]:
30         count+=1
31 pdb.set_trace()
32 inception = InceptionV3(weights='imagenet', include_top=True)
33 inception = Model(inception.input, inception.layers[-2].output)
34
35 model = load_model('./model_weights/model_inception_epoch1.h5')
36 max_length = 17
37
38 def get_image(image_id):
39     feature = features[str(image_id)]
40     caps = captions[str(image_id)]
41     cap = caps[0] # use first caption
42     return feature, cap
43
44 # generate caption for a given sample
45 def describe_image(img, model):
46     input_string = 'x_START_'
47     for i in range(max_length):
48         sequence = [word_to_idx[w] for w in input_string.split() if w in word_to_idx.keys
49                     ()]
50         sequence = pad_sequences([sequence], maxlen=max_length)
51         yhat = model.predict([img,sequence], verbose=0)
52         yhat = np.argmax(yhat)
53         word = idx_to_word[str(yhat)]
54         input_string += ' ' + word
55
56         if word == 'x_END_':
57             break
58
59     final = input_string.split()
60     final = final[1:-1]
61     final = ' '.join(final)
62     return final
63
64 def get_image_caption(image_id):
65     features, caption = get_image(image_id)
66     result = [idx_to_word[str(x)] for x in caption if idx_to_word[str(x)] != "x_NULL_"]
67     result = result[1:-1]
68     result = ' '.join(result)
69     return result
70
71 images = ["car_1", "doruk_1", "ev_1", "furkan_laptop"]
72
73 for id in images:
74     path= os.path.join('./data/real_life_images/'+id+".jpeg")
75     img_ = Image.open(path)
76     processed_img = img_.resize((299,299), resample=Image.ANTIALIAS)
77     temp = np.array(processed_img)
78     if temp.shape == (299,299):
79         temp_ = [temp, temp, temp]
80         temp = np.stack(temp_, axis=2)
81     img= np.array(temp, dtype=np.uint8)
82     img = preprocess_input(img)
83     res = inception.predict(img)

```

```

84     res = np.reshape(res, res.shape[1])
85     img_features = np.expand_dims(res, axis=0)
86
87     pred_cap = describe_image(img_features, model)
88     print(pred_cap)
89     fig, ax = plt.subplots()
90     ax.get_xaxis().set_ticks([])
91     ax.get_yaxis().set_ticks([])
92     ax.imshow(img_)
93     ax.axis('on')
94     ax.set_xlabel("Predicted Caption: " + pred_cap, fontsize=10)
95     plt.savefig(id+"_capt")
96     plt.close()

```

```

1
2 import pickle
3 from keras.preprocessing.sequence import pad_sequences
4 import numpy as np
5 import pdb
6 import random
7 from keras.applications.inception_v3 import InceptionV3, preprocess_input
8 import sys
9 from PIL import Image
10 from keras.models import load_model
11 from tensorflow.keras.models import Model
12
13 # import conversion dictionaries
14 with open("./utils/word_to_idx.pkl", "rb") as f:
15     word_to_idx = pickle.load(f)
16 with open("./utils/idx_to_word.pkl", "rb") as f:
17     idx_to_word = pickle.load(f)
18
19 # all dataset
20 with open("./data/dataset/processed_dataset_vgg.pkl", "rb") as f:
21     features = pickle.load(f)
22 with open("./utils/captions.pkl", "rb") as f:
23     captions = pickle.load(f)
24
25 model = load_model('./model_weights/model_inception_epoch1.h5')
26 max_length = 17
27
28 def get_image(image_id):
29     feature = features[str(image_id)]
30     caps = captions[str(image_id)]
31     cap = caps[0] # use first caption
32     return feature, cap
33
34 # generate caption for a given sample
35 def describe_image(img, model):
36     input_string = 'x_START_'
37     for i in range(max_length):
38         sequence = [word_to_idx[w] for w in input_string.split() if w in word_to_idx.keys()]
39         sequence = pad_sequences([sequence], maxlen=max_length)
40
41         yhat = model.predict([img, sequence], verbose=0)
42         yhat = np.argmax(yhat)
43         word = idx_to_word[str(yhat)]
44         input_string += ' ' + word
45
46         if word == 'x_END_':
47             break
48
49     final = input_string.split()
50     final = final[1:-1]
51     final = ' '.join(final)
52     return final
53
54 def get_image_caption(image_id):
55     features, caption = get_image(image_id)
56     result = [idx_to_word[str(x)] for x in caption]
57     result = result[1:-1]

```

```

58     result = ' '.join(result)
59     return result
60
61 images = [24182, 137, 30296, 10518, 25178, 16968, 21062, 30707, 27506, 18408, 39625,
62           860, 6086, 40825, 8690, 32552]
63
64 for id in images:
65     image = Image.open(os.path.join('./data/dataset/images/'+str(id)+".jpeg"))
66     img, cap = get_image(id)
67     pred_cap = describe_image(features[str(id)], model)
68     pdb.set_trace()

```

```

1  import os
2  import pdb
3  import h5py
4  import pickle
5
6  import numpy as np
7
8  from tqdm import tqdm
9
10 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
11     captions = f['train_cap'][:]
12     image_ids = f['train_imid'][:]
13
14 min_id = np.min(image_ids)
15 max_id = np.max(image_ids)
16
17 dct = {}
18 for i in tqdm(range(min_id, max_id + 1)):
19     idx = np.where(image_ids == i)
20     dct[str(i)] = [x.tolist() for x in captions[idx]]
21
22 pdb.set_trace()
23
24 with open("./utils/captions.pkl", "wb") as f:
25     pickle.dump(dct, f)

```

```

1  import pdb
2  import h5py
3  import pickle
4
5  import numpy as np
6  import matplotlib.pyplot as plt
7
8  from sklearn.manifold import TSNE
9  from mpl_toolkits.mplot3d import Axes3D
10
11 embedding_matrix = np.load("./utils/embedding_matrix.npy")
12 with open("./utils/idx_to_word.pkl", "rb") as f:
13     idx_to_word = pickle.load(f)
14
15 vocabulary_size = 1004
16 tsne = TSNE(n_components=2, n_iter=3000, perplexity=30, learning_rate=200)
17 embedding = tsne.fit_transform(embedding_matrix)
18
19
20 figure = plt.figure()
21 ax = figure.add_subplot(111)
22 ax.set_title("TSNE applied on Word Embedding Matrix")
23 ax.scatter(embedding[:,0], embedding[:,1], c='g')
24 ax.set_xlabel("Component 1")
25 ax.set_ylabel("Component 2")
26 for idx in range(vocabulary_size):
27     plt.annotate(idx_to_word[str(idx)], (embedding[idx,0], embedding[idx,1]), fontsize=6)
28
29 plt.savefig('./plots/tsne_word_embedding.png')

```

```

1  import os
2  import pdb
3  import h5py
4  import pickle

```



```

5
6 import numpy as np
7
8 from tqdm import tqdm
9
10 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
11     captions = f['train_cap'][:]
12     image_ids = f['train_imid'][:]
13     images = f['train_ims'][:]
14     urls = f['train_url'][:]
15     word_codes = f['word_code'][:]
16
17 with open('./utils/word_to_idx.pkl', 'rb') as f:
18     word_to_idx = pickle.load(f)
19 with open('./utils/idx_to_word.pkl', 'rb') as f:
20     idx_to_word = pickle.load(f)
21
22 word_frequency = {}
23 for word in list(word_to_idx.keys()):
24     word_frequency[word] = 0
25
26 for caption in tqdm(captions):
27     for word in caption:
28         word_frequency[idx_to_word[str(word)]] += 1
29
30 low_frequency_word_count = 0
31 for word in word_frequency.keys():
32     if word_frequency[word] <= 10:
33         print(word, "\t", word_frequency[word])
34         low_frequency_word_count += 1
35
36 print("Low Frequency Word Count: ", low_frequency_word_count)
37
38 words = list(word_frequency.keys())
39 frequencies = list(word_frequency.values())
40 frequencies, words = (list(t) for t in zip(*sorted(zip(frequencies, words))))
41 print(words[-10:])
42 print(frequencies[-10:])

```

5.2 Codes for Inception V3 + LSTM

```

1 # from tensorflow.python.framework import ops
2 # ops.reset_default_graph()
3
4 import numpy as np
5 import tensorflow as tf
6
7 from keras.layers import Conv1D
8 from keras.models import Sequential
9 from keras.applications.inception_v3 import InceptionV3
10 from keras.preprocessing.sequence import pad_sequences
11 from keras.utils import to_categorical
12 from keras.optimizers import Adam, RMSprop
13 from keras.layers.merge import add
14 from keras.layers import Input
15 from keras.layers import Dropout, Dense, LSTM, Embedding
16 from keras.models import Model, load_model
17 import keras.backend as K
18
19 import pickle
20 import pdb
21
22 def preprocess_dataset(dct):
23     result = {}
24     for key, value in zip(dct['keys'], dct['values']):
25         result[key] = value
26     return result
27
28 # # all dataset
29 # with open("./utils/processed_dataset_vgg.pkl", "rb") as f:
30 #     all_features = pickle.load(f)
31 # with open("./utils/captions.pkl", "rb") as f:

```

```

32 #     all_captions = pickle.load(f)
33
34 # train dataset (70% of the data)
35 with open("./data/dataset/train_dataset_inception.pkl", "rb") as f:
36     train_features = pickle.load(f)
37     train_encoded_images = preprocess_dataset(train_features)
38 with open("./data/dataset/train_captions_inception.pkl", "rb") as f:
39     train_captions = pickle.load(f)
40     train_dataset = train_captions
41     # train_dataset = preprocess_dataset(train_captions)
42
43 # validation dataset (20% of the data)
44 with open("./data/dataset/validation_dataset_inception.pkl", "rb") as f:
45     validation_features = pickle.load(f)
46     validation_encoded_images = preprocess_dataset(validation_features)
47 with open("./data/dataset/validation_captions_inception.pkl", "rb") as f:
48     validation_captions = pickle.load(f)
49     validation_dataset = validation_captions
50     # validation_dataset = preprocess_dataset(validation_captions)
51
52 # test dataset (10% of the data)
53 with open("./data/dataset/test_dataset_inception.pkl", "rb") as f:
54     test_features = pickle.load(f)
55     test_encoded_images = preprocess_dataset(test_features)
56 with open("./data/dataset/test_captions_inception.pkl", "rb") as f:
57     test_captions = pickle.load(f)
58     test_dataset = test_captions
59     # test_dataset = preprocess_dataset(test_captions)
60
61 # train_dataset = train_captions
62 # train_encoded_images = train_features
63
64 # test_dataset = test_captions
65 # test_encoded_images = test_features
66
67 # validation_dataset = validation_captions
68 # validation_encoded_images = validation_features
69
70
71 # model
72 max_length = 17
73 embedding_dim = 200
74 vocabulary_size = 1004
75
76 '''
77 Model takes 2 types of inputs, first being encoded photos by inception to 2048
78 dimensional vectors. Second input type to the model are sequences of words (i.e,
79 captions).
80 Captions are all zero padded to max_length. With first item being SOS token and last item
81 being EOS token. Dropout layers are used for prevention of overfitting. Photos are
82 processed
83 such that their dimensions are reduced to 256 by dense layer with relu activation.
84 Caption sequences are processed by the glove pretrained embedding matrix. Embedding
85 dimension
86 is chosen as 200. Vocabulary size in the problem is 1020.
87 '''
88
89 inputs1 = Input(shape=(2048,))
90 features1 = Dropout(0.5)(inputs1)
91 features2 = Dense(256, activation='relu')(features1)
92 inputs2 = Input(shape=(max_length,))
93 seqfeatures1 = Embedding(vocabulary_size, embedding_dim, mask_zero=True)(inputs2)
94 seqfeatures2 = Dropout(0.5)(seqfeatures1)
95 seqfeatures3 = LSTM(256)(seqfeatures2)
96
97 #add two different types of feature extractions into one tensor.
98 decoder1 = add([features2, seqfeatures3])
99
100 decoder2 = Dense(256, activation='relu')(decoder1)
101 #softmax to predict target word over vocabulary.
102 outputs = Dense(vocabulary_size, activation='softmax')(decoder2)
103 #construction of the model.

```

```

98 model = Model(inputs=[inputs1, inputs2], outputs=outputs)
99
100 '''
101 Data generator function streams the input data to the model in batches. This will be used
102 for fit_generator function with the model.
103 Note that, batch size is in unit of images. With batch_size = 3, there will be 3 images
104 with corresponding training sample generation.
105 '''
106 batch_size = 6
107 def data_generator(dataset, imgs, max_length, batch_size = batch_size):
108     X1, X2, y = list(), list(), list()
109     n = 0
110     #infinite loop over images
111     while True:
112         for key, captions_list in dataset.items():
113             n += 1
114             # retrieve the photo feature
115             try:
116                 image = imgs[key]
117             except:
118                 n -= 1
119                 continue
120             for capt in captions_list:
121                 # split one sequence into multiple X, y pairs
122                 for i in range(1, len(capt)):
123                     # split into input and output pair
124                     in_seq, out_seq = capt[:i], capt[i]
125                     if out_seq == 0:
126                         break
127                     # pad input sequence
128                     in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
129                     # encode output sequence
130                     out_seq = to_categorical([out_seq], num_classes=vocabulary_size)[0]
131                     # store
132                     X1.append(image)
133                     X2.append(in_seq)
134                     y.append(out_seq)
135                 # yield the batch data
136                 if n==batch_size:
137                     yield [[np.array(X1), np.array(X2)], np.array(y)]
138                     X1, X2, y = list(), list(), list()
139                     n=0
140
141 #get embedding matrix with shape 1024x200
142 print(model.summary())
143 embedding_matrix = np.load("./utils/embedding_matrix.npy")
144 model.layers[2].set_weights([embedding_matrix])
145 model.layers[2].trainable = False
146
147 #compile the model with categorical cross entropy, combined with softmax
148 model.compile(loss='categorical_crossentropy', optimizer='adam')
149
150 steps = len(train_dataset)//batch_size
151 # generator = data_generator(train_dataset, train_encoded_images, max_length, batch_size)
152 # model.optimizer.lr = 0.01
153 # model.fit_generator(generator, epochs=9, steps_per_epoch=steps, verbose=1)
154 # model.save('./model_' + '9epochs' + '.h5')
155 epochs = 50
156 # # load model
157 base = 0
158 # model = load_model('./model_weights/model_9.h5')
159 # model.optimizer.learning_rate = 0.00001
160
161 # print(K.eval(model.optimizer.lr))
162 # K.set_value(model.optimizer.lr, 0.000001)
163 # print(K.eval(model.optimizer.lr))
164
165
166 for i in range(1, epochs+1):
167     if i == 5:

```

```

168         batch_size = 12
169         K.set_value(model.optimizer.lr, 0.0001)
170     elif i == 10:
171         batch_size = 18
172         K.set_value(model.optimizer.lr, 0.00001)
173     train_generator = data_generator(train_dataset, train_encoded_images, max_length,
174                                     batch_size)
175     validation_generator = data_generator(validation_dataset, validation_encoded_images,
176                                         max_length, batch_size)
177     # pdb.set_trace()
178     model.fit_generator(train_generator, epochs=1, steps_per_epoch=steps, verbose=1,
179                       validation_data=validation_generator, validation_steps=steps, validation_freq=1)
180     model.save('./model_weights/model_inception_epoch' + str(base+i) + '.h5')

```

```

1 import os
2 import pdb
3 import pickle
4
5 import pandas as pd
6 import numpy as np
7
8 from keras.applications.inception_v3 import InceptionV3, preprocess_input
9 from keras.models import Model
10 from keras import backend as K
11
12 from tqdm import tqdm
13
14 inception = InceptionV3(weights='imagenet', include_top=True)
15 model = Model(inception.input, inception.layers[-2].output)
16
17 def extract_features(img):
18     #img are 299x299x3 images of rgb
19     img = preprocess_input(img)
20     img = np.expand_dims(img,axis=0)
21     res = model.predict(img)
22     res = np.reshape(res, res.shape[1])
23
24     return res
25
26 # import processed images
27 with open("./data/dataset/dataset_inception.pkl", "rb") as f:
28     dct = pickle.load(f)
29
30 images = list(dct.values())
31 ids = list(dct.keys())
32
33 dataset = {}
34 for (img, id) in tqdm(zip(images, ids)):
35     feature = extract_features(img)
36     dataset[id] = feature
37
38 # save image features
39 with open("./data/dataset/processed_dataset_inception.pkl", "wb") as f:
40     pickle.dump(dataset, f)

```

```

1 import os
2 import pdb
3 import h5py
4 import re
5 import pickle
6
7 import numpy as np
8
9 from PIL import Image
10 from tqdm import tqdm
11
12 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
13     captions = f['train_cap'][:]
14     image_ids = f['train_imid'][:]
15     images = f['train_ims'][:]
16     urls = f['train_url'][:]
17     word_codes = f['word_code'][:]

```

```

18
19 # preprocess and save images
20 url_count = 82783
21 filenames = os.listdir('./data/dataset/images/')
22 image_ids = list(range(1, url_count + 1))
23 dataset = {}
24 for filename in tqdm(filenames):
25     img = Image.open(os.path.join('./data/dataset/images/', filename))
26     img_id = int(re.split('\.', filename)[0])
27     try:
28         img = img.convert('RGB')
29         processed_img = img.resize((299,299), resample=Image.ANTIALIAS)
30         # print(filename, "--> ", img_id)
31         # image is greyscale
32         temp = np.array(processed_img)
33         if temp.shape == (299,299):
34             temp_ = [temp, temp, temp]
35             temp = np.stack(temp_, axis=2)
36             dataset[str(img_id)] = np.array(temp, dtype=np.uint8)
37     except:
38         print("Exception --> ", img_id)
39     pdb.set_trace()
40
41 # generate training and test datasets
42 with open('./data/dataset/dataset_inception.pkl', 'wb') as f:
43     pickle.dump(dataset, f)

```

```

1 import os
2 import pickle
3 import pdb
4
5 import numpy as np
6
7 from sklearn.model_selection import train_test_split
8
9 # dummy data
10 # dataset = dict(a=[1], b=[2], c=[3], d=[4], e=[5], f=[6], g=[7], h=[8], i=[9], j=[10])
11 with open("./data/dataset/processed_dataset_inception.pkl", "rb") as f:
12     dataset = pickle.load(f)
13
14 keys = list(dataset.keys())
15 values = list(dataset.values())
16
17 images_train_val, images_test, id_train_val, id_test = train_test_split(values, keys,
18     shuffle=True, train_size=9/10)
19 images_train, images_val, id_train, id_val = train_test_split(images_train_val,
20     id_train_val, shuffle=True, train_size = 7/9)
21
22 train_dataset = dict(keys=id_train, values=images_train)
23 validation_dataset = dict(keys=id_val, values=images_val)
24 test_dataset = dict(keys=id_test, values=images_test)
25
26 # save datasets
27 with open("./data/dataset/train_dataset_inception.pkl", "wb") as f:
28     pickle.dump(train_dataset, f)
29 with open("./data/dataset/validation_dataset_inception.pkl", "wb") as f:
30     pickle.dump(validation_dataset, f)
31 with open("./data/dataset/test_dataset_inception.pkl", "wb") as f:
32     pickle.dump(test_dataset, f)
33
34 # load captions and split to train test and validation in a corresponding manner
35 with open("./utils/captions.pkl", "rb") as f:
36     captions = pickle.load(f)
37
38 train_captions = {}
39 for key in id_train:
40     train_captions[key] = captions[key]
41
42 validation_captions = {}
43 for key in id_val:
44     validation_captions[key] = captions[key]

```

```

44 test_captions = {}
45 for key in id_test:
46     test_captions[key] = captions[key]
47
48 # save captions
49 with open("./data/dataset/train_captions_inception.pkl", "wb") as f:
50     pickle.dump(train_captions, f)
51 with open("./data/dataset/validation_captions_inception.pkl", "wb") as f:
52     pickle.dump(validation_captions, f)
53 with open("./data/dataset/test_captions_inception.pkl", "wb") as f:
54     pickle.dump(test_captions, f)

```

5.3 Codes for VGG16 + LSTM

```

1 # from tensorflow.python.framework import ops
2 # ops.reset_default_graph()
3
4 import numpy as np
5 import tensorflow as tf
6
7 from keras.layers import Conv1D
8 from keras.models import Sequential
9 from keras.applications.inception_v3 import InceptionV3
10 from keras.preprocessing.sequence import pad_sequences
11 from keras.utils import to_categorical
12 from keras.optimizers import Adam, RMSprop
13 from keras.layers.merge import add
14 from keras.layers import Input
15 from keras.layers import Dropout, Dense, LSTM, Embedding
16 from keras.models import Model, load_model
17 import keras.backend as K
18
19 import pickle
20 import pdb
21
22 def preprocess_dataset(dct):
23     result = {}
24     for key, value in zip(dct['keys'], dct['values']):
25         result[key] = value
26     return result
27
28 # all dataset
29 # with open("./utils/processed_dataset_vgg.pkl", "rb") as f:
30 #     all_features = pickle.load(f)
31 # with open("./utils/captions.pkl", "rb") as f:
32 #     all_captions = pickle.load(f)
33
34 # train dataset (70% of the data)
35 with open("./data/dataset/train_dataset_vgg.pkl", "rb") as f:
36     train_features = pickle.load(f)
37     train_encoded_images = preprocess_dataset(train_features)
38 with open("./data/dataset/train_captions_vgg.pkl", "rb") as f:
39     train_captions = pickle.load(f)
40     train_dataset = train_captions
41     # train_dataset = preprocess_dataset(train_captions)
42
43 # validation dataset (20% of the data)
44 with open("./data/dataset/validation_dataset_vgg.pkl", "rb") as f:
45     validation_features = pickle.load(f)
46     validation_encoded_images = preprocess_dataset(validation_features)
47 with open("./data/dataset/validation_captions_vgg.pkl", "rb") as f:
48     validation_captions = pickle.load(f)
49     validation_dataset = validation_captions
50     # validation_dataset = preprocess_dataset(validation_captions)
51
52 # test dataset (10% of the data)
53 with open("./data/dataset/test_dataset_vgg.pkl", "rb") as f:
54     test_features = pickle.load(f)
55     test_encoded_images = preprocess_dataset(test_features)
56 with open("./data/dataset/test_captions_vgg.pkl", "rb") as f:
57     test_captions = pickle.load(f)
58     test_dataset = test_captions

```

```

59     # test_dataset = preprocess_dataset(test_captions)
60
61 # train_dataset = train_captions
62 # train_encoded_images = train_features
63
64 # test_dataset = test_captions
65 # test_encoded_images = test_features
66
67 # validation_dataset = validation_captions
68 # validation_encoded_images = validation_features
69
70
71 # model
72 max_length = 17
73 embedding_dim = 200
74 vocabulary_size = 1004
75
76 '''
77 Model takes 2 types of inputs, first being encoded photos by inception to 2048
    dimensional vectors. Second input type to the model are sequences of words (i.e,
    captions).
78 Captions are all zero padded to max_length. With first item being SOS token and last item
    being EOS token. Dropout layers are used for prevention of overfitting. Photos are
    processed
79 such that their dimensions are reduced to 256 by dense layer with relu activation.
    Caption sequences are processed by the glove pretrained embedding matrix. Embedding
    dimension
80 is chosen as 200. Vocabulary size in the problem is 1020.
81 '''
82
83 inputs1 = Input(shape=(4096,))
84 features1 = Dropout(0.5)(inputs1)
85 features2 = Dense(256, activation='relu')(features1)
86 inputs2 = Input(shape=(max_length,))
87 seqfeatures1 = Embedding(vocabulary_size, embedding_dim, mask_zero=True)(inputs2)
88 seqfeatures2 = Dropout(0.5)(seqfeatures1)
89 seqfeatures3 = LSTM(256)(seqfeatures2)
90
91 #add two different types of feature extractions into one tensor.
92 decoder1 = add([features2, seqfeatures3])
93
94 decoder2 = Dense(256, activation='relu')(decoder1)
95 #softmax to predict target word over vocabulary.
96 outputs = Dense(vocabulary_size, activation='softmax')(decoder2)
97 #construction of the model.
98 model = Model(inputs=[inputs1, inputs2], outputs=outputs)
99
100
101 '''
102 Data generator function streams the input data to the model in batches. This will be used
    for fit_generator function with the model.
103 Note that, batch size is in unit of images. With batch_size = 3, there will be 3 images
    with corresponding training sample generation.
104 '''
105
106 batch_size = 6
107 def data_generator(dataset, imgs, max_length, batch_size = batch_size):
108     X1, X2, y = list(), list(), list()
109     n = 0
110     #infinite loop over images
111     while True:
112         for key, captions_list in dataset.items():
113             n += 1
114             # retrieve the photo feature
115             try:
116                 image = imgs[key]
117             except:
118                 n -= 1
119                 continue
120             for capt in captions_list:
121                 # split one sequence into multiple X, y pairs
122                 for i in range(1, len(capt)):

```

```

123         # split into input and output pair
124         in_seq, out_seq = capt[:i], capt[i]
125         if out_seq == 0:
126             break
127         # pad input sequence
128         in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
129         # encode output sequence
130         out_seq = to_categorical([out_seq], num_classes=vocabulary_size)[0]
131         # store
132         X1.append(image)
133         X2.append(in_seq)
134         y.append(out_seq)
135     # yield the batch data
136     if n==batch_size:
137         yield [[np.array(X1), np.array(X2)], np.array(y)]
138         X1, X2, y = list(), list(), list()
139         n=0
140
141 #get embedding matrix with shape 1024x200
142 print(model.summary())
143 embedding_matrix = np.load("./utils/embedding_matrix.npy")
144 model.layers[2].set_weights([embedding_matrix])
145 model.layers[2].trainable = False
146
147 #compile the model with categorical cross entropy, combined with softmax
148 model.compile(loss='categorical_crossentropy', optimizer='adam')
149
150 steps = len(train_dataset)//batch_size
151 # generator = data_generator(train_dataset, train_encoded_images, max_length, batch_size)
152 # model.optimizer.lr = 0.01
153 # model.fit_generator(generator, epochs=9, steps_per_epoch=steps, verbose=1)
154 # model.save('./model_' + '9epochs' + '.h5')
155 epochs = 50
156 # # load model
157 base = 0
158 # model = load_model('./model_weights/model_9.h5')
159 # model.optimizer.learning_rate = 0.00001
160
161 # print(K.eval(model.optimizer.lr))
162 # K.set_value(model.optimizer.lr, 0.000001)
163 # print(K.eval(model.optimizer.lr))
164
165
166 for i in range(1,epochs+1):
167     if i == 5:
168         batch_size = 12
169         K.set_value(model.optimizer.lr, 0.0001)
170     elif i == 10:
171         batch_size = 18
172         K.set_value(model.optimizer.lr, 0.00001)
173     train_generator = data_generator(train_dataset, train_encoded_images, max_length,
174                                     batch_size)
175     validation_generator = data_generator(validation_dataset, validation_encoded_images,
176                                         max_length, batch_size)
177     # pdb.set_trace()
178     model.fit_generator(train_generator, epochs=1, steps_per_epoch=steps, verbose=1,
179                       validation_data=validation_generator, validation_steps=steps, validation_freq=1)
180     model.save('./model_weights/model_vgg_epoch' + str(base+i) + '.h5')
181     base+=1

```

```

1 import os
2 import pdb
3 import pickle
4
5 import pandas as pd
6 import numpy as np
7
8 from keras.applications.vgg16 import VGG16, preprocess_input
9 from keras.models import Model
10 from tqdm import tqdm
11 from keras import backend as K
12
13

```



```

14 vgg = VGG16(weights='imagenet', include_top=True)
15 transfer_layer = vgg.get_layer('fc2') # one previous from last
16 model = Model(vgg.input, transfer_layer.output)
17
18 def extract_features(img):
19     #img are 299x299x3 images of rgb
20     img = preprocess_input(img)
21     img = np.expand_dims(img,axis=0)
22     res = model.predict(img)
23     res = np.reshape(res, res.shape[1])
24     return res
25
26 # import processed images
27 with open("./data/dataset/dataset_vgg.pkl", "rb") as f:
28     dct = pickle.load(f)
29
30 images = list(dct.values())
31 ids = list(dct.keys())
32
33 dataset = {}
34 for (img, id) in tqdm(zip(images, ids)):
35     feature = extract_features(img)
36     dataset[id] = feature
37
38 # save image features
39 with open("./data/dataset/processed_dataset_vgg.pkl", "wb") as f:
40     pickle.dump(dataset, f)

```

```

1 import os
2 import pdb
3 import h5py
4 import re
5 import pickle
6
7 import numpy as np
8
9 from PIL import Image
10 from tqdm import tqdm
11
12 with h5py.File('./data/eee443_project_dataset_train.h5', 'r') as f:
13     captions = f['train_cap'][:]
14     image_ids = f['train_imid'][:]
15     images = f['train_ims'][:]
16     urls = f['train_url'][:]
17     word_codes = f['word_code'][:]
18
19 # preprocess and save images
20 url_count = 82783
21 filenames = os.listdir('./data/dataset/images/')
22 image_ids = list(range(1, url_count + 1))
23 dataset = {}
24 for filename in tqdm(filenames):
25     img = Image.open(os.path.join('./data/dataset/images/', filename))
26     img_id = int(re.split('\.', filename)[0])
27     try:
28         img = img.convert('RGB')
29         processed_img = img.resize((224,224), resample=Image.ANTIALIAS)
30         # print(filename, "--> ", img_id)
31         # image is greyscale
32         temp = np.array(processed_img)
33         if temp.shape == (224,224):
34             temp_ = [temp, temp, temp]
35             temp = np.stack(temp_, axis=2)
36             dataset[str(img_id)] = np.array(temp, dtype=np.uint8)
37     except:
38         print("Exception --> ", img_id)
39 pdb.set_trace()
40
41 # generate training and test datasets
42 with open('./data/dataset/dataset_vgg.pkl', 'wb') as f:
43     pickle.dump(dataset, f)

```

```

1 import os
2 import pickle
3 import pdb
4
5 import numpy as np
6
7 from sklearn.model_selection import train_test_split
8
9 # dummy data
10 # dataset = dict(a=[1], b=[2], c=[3], d=[4], e=[5], f=[6], g=[7], h=[8], i=[9], j=[10])
11 with open("./data/dataset/processed_dataset_vgg.pkl", "rb") as f:
12     dataset = pickle.load(f)
13
14 keys = list(dataset.keys())
15 values = list(dataset.values())
16
17 images_train_val, images_test, id_train_val, id_test = train_test_split(values, keys,
18     shuffle=True, train_size=9/10)
19 images_train, images_val, id_train, id_val = train_test_split(images_train_val,
20     id_train_val, shuffle=True, train_size = 7/9)
21
22 train_dataset = dict(keys=id_train, values=images_train)
23 validation_dataset = dict(keys=id_val, values=images_val)
24 test_dataset = dict(keys=id_test, values=images_test)
25
26 # save datasets
27 with open("./data/dataset/train_dataset_vgg.pkl", "wb") as f:
28     pickle.dump(train_dataset, f)
29 with open("./data/dataset/validation_dataset_vgg.pkl", "wb") as f:
30     pickle.dump(validation_dataset, f)
31 with open("./data/dataset/test_dataset_vgg.pkl", "wb") as f:
32     pickle.dump(test_dataset, f)
33
34 # load captions and split to train test and validation in a corresponding manner
35 with open("./utils/captions.pkl", "rb") as f:
36     captions = pickle.load(f)
37
38 train_captions = {}
39 for key in id_train:
40     train_captions[key] = captions[key]
41
42 validation_captions = {}
43 for key in id_val:
44     validation_captions[key] = captions[key]
45
46 test_captions = {}
47 for key in id_test:
48     test_captions[key] = captions[key]
49
50 # save captions
51 with open("./data/dataset/train_captions_vgg.pkl", "wb") as f:
52     pickle.dump(train_captions, f)
53 with open("./data/dataset/validation_captions_vgg.pkl", "wb") as f:
54     pickle.dump(validation_captions, f)
55 with open("./data/dataset/test_captions_vgg.pkl", "wb") as f:
56     pickle.dump(test_captions, f)

```

References

- [1] L. Fei-Fei, A. Iyer, C. Koch, and P. Perona, "What do we perceive in a glance of a real-world scene?" *Journal of vision*, vol. 7, no. 1, pp. 10–10, 2007.
- [2] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [4] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

- [5] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2818–2826.
- [6] G. Kulkarni, V. Premraj, V. Ordonez, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg, “Babytalk: Understanding and generating simple image descriptions,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 12, pp. 2891–2903, 2013.
- [7] A. Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth, “Every picture tells a story: Generating sentences from images,” in *European conference on computer vision*. Springer, 2010, pp. 15–29.
- [8] H. Fang, S. Gupta, F. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt *et al.*, “From captions to visual concepts and back,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1473–1482.
- [9] A. Karpathy and L. Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3128–3137.
- [10] J. Devlin, H. Cheng, H. Fang, S. Gupta, L. Deng, X. He, G. Zweig, and M. Mitchell, “Language models for image captioning: The quirks and what works,” *arXiv preprint arXiv:1505.01809*, 2015.
- [11] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3156–3164.
- [12] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *In EMNLP*, 2014.
- [13] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [14] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 2002, pp. 311–318.