

## Generation of Captions from Images – Report

Surya Menon – S-89A Deep Learning for Natural Language Processing, Summer 2019

### Abstract –

#### Problem Statement:

Use deep learning and natural language processing techniques to generate relevant captions from input images.

#### Approach:

Train a neural network with the Flickr 8k dataset, that takes an image and text input and produces an appropriate word sequence to describe the image.

#### Use / Benefits:

This process could be used to extract insights/summaries from images, provide image information to those with visual impairments, and to classify/identify patterns in social media photo albums (i.e. facial recognition), or perhaps even in medical imaging assessments.

#### Drawbacks / Challenges:

- There was necessary text and image preparation and heavy reliance on pretrained networks for extracting image features and representing words as vectors.
- Encountered overfitting (perhaps due to needing more data for training) despite applying regularization techniques and simplifying the network.
- The model's effectiveness is limited to the variety of images it is trained on.

#### Results:

We were able to get the validation loss of generating sequences similar to actual captions to around 3.2, but experienced overfitting. The corpus BLEU scores used to evaluate the model were not high, but conducting predictions on test images (using both greedy and beam searches) resulted in captions that generally effectively capture the context of the images. Training with larger samples or employing data augmentation, as well as including an attention mechanism, cross-validation, and further hyperparameter tuning may improve results.

#### Working Example Description:

The following notebooks were used to create the demonstration:

1. **project\_img-features.ipynb**: Processes the dataset images in a pretrained convolutional neural network (CNN) and extracts/saves image features
2. **model\_create.ipynb**: Processes and visualizes dataset captions; creates training, validation, test data; tokenizes words and establishes word embedding; trains a “merge” model with a data generator; visualizes losses; evaluates model with BLEU scores; demonstrates a few examples of captions using greedy and beam searches.
3. **model\_demo.ipynb**: Loads final model and generates captions for a few random test images with image plotting.

#### Data Source:

The Flickr 8k dataset consists of around 8000 images (of everyday people and situations from [Flickr](#)) with five associated captions per image, with predefined training, validation, and test sets.

- Officially request data here: <https://forms.illinois.edu/sec/1713398>
- Alternative to get immediately: <https://github.com/jbrownlee/Datasets/releases>

**Video:** <https://youtu.be/zPawquEoQFs>

## Introduction and Dataset –

To generate captions from images, we need to develop a model that includes image processing and sequence generation, employing both convolutional and recurrent neural networks<sup>1</sup>. We follow the logic employed by several resources for text and image preparation and techniques for training to create an effective model<sup>2</sup>.

The dataset used is called Flickr 8k, which consists of 6000 training images (JPEG), 1000 validation images, and 1000 test images, with each image having five associated captions; these images are predefined and indexed into these categories. The dataset comes in multiple zip files, which you should unzip and put in a file directory (i.e. project\_data). We must process both the image and text data before we use them in our model, which we see next.

**NOTE:** The data submitted in `CaptionGeneration_MenorSurya_Final.zip` (in the `project_data` folder) for this project only contains a **sample** of the dataset used for training and evaluating the model created. If you want to reproduce the results seen in the code below (including full training, BLEU evaluations, and specific test predictions with image plotting), you need to download the full dataset, which can be accessed [here](#), and place the downloaded folders into a `project_data` directory. When using Google Colab, you can upload the needed files (i.e. model, captions text, indexing files, image features dictionary, specific JPEG images) via the Files tab; the paths defined in the code in `model_create.ipynb` follow this process. Alternatively, we provide the paths to use with Jupyter Notebook (or configuring with Google Drive) in the comments and other code files. The artifacts produced in the following code (except the intermediate model files) are in the `artifacts` folder.

## Image Preparation –

Associated file: `project_img-features.ipynb`

```
1. # alternatively install keras separately
2. import tensorflow
3. from tensorflow import keras
4. from tensorflow.keras.applications.inception_v3 import InceptionV3
5. from tensorflow.keras.applications.inception_v3 import preprocess_input
6. from tensorflow.keras.models import Model
7. from tensorflow.keras.preprocessing import image
8. import numpy as np
9. import os
10. import pickle
```

We want to process our input images when training the model so they contain the most relevant features (i.e. capturing the most important context, content, patterns) but are in a compressed representation. For efficiency, we capture and save these image feature representations first, which we later pass into the model as inputs, rather than training the entire image with the text sequences all at once<sup>3</sup>.

We use the InceptionV3 model which is a convolutional neural network (CNN) that conducts image classification using pre-trained weights<sup>4</sup>. Specifically, this model uses ImageNet weights, which have been trained on millions of images and over thousands of categories<sup>5</sup>. With Keras we can load this model, remove the classification layer at the top of the network, and obtain outputs from this model that are compressed representations of the input images:

```
1. # process input for inception model
```

```

2. def img_size(img):
3.     # turn image into array
4.     img = image.img_to_array(img)
5.     # add dimension
6.     img = np.expand_dims(img, axis=0)
7.     # use InceptionV3 preprocess
8.     img = preprocess_input(img)
9.     return img
10.
11. # go through images in folder and load model
12. # using Jupyter Notebook here
13. img_folder = os.path.join(os.getcwd(), 'project_data\Flickr8k_Dataset\Flicker8k_Dataset
    \\')
14. model = InceptionV3(weights='imagenet')
15.
16. # not use to classify but get fixed length vector for images (compression)
17. # want output to be 2nd to last layer
18. feature_mod = Model(model.input, model.layers[-2].output)

```

After loading and initializing the model, we go through the dataset images folder and pass each image through the model; each image is first processed so it can be passed into the network (including resizing and formatting), and the resulting 2048-dimensional vector output is stored in a dictionary, which is indexed by the image identifier (i.e. image file name). We save this dictionary using the pickle package (under name `img_features.pkl`), and we later load these image representations as inputs for the model created in `model_create.ipynb`.

```

1. # feature dictionary (index by image name)
2. imgs = dict()
3.
4. # go through images in folder
5. for i in os.listdir(img_folder):
6.     # print image name - just for tracking
7.     print(i)
8.     # create path
9.     path = img_folder + i
10.    # load image with InceptionV3 default size
11.    img = image.load_img(path, target_size=(299,299))
12.    # prepare image for CNN model
13.    img = img_size(img)
14.    # get features
15.    feat = feature_mod.predict(img)
16.    feat = np.reshape(feat, feat.shape[1])
17.    # pass to dictionary
18.    imgs[i.split('.')[0]] = feat
19.
20. # save for later
21. pickle.dump(imgs, open('img_features.pkl', 'wb'))

```

## Text Preparation –

Associated file: `model_create.ipynb`

```

1. import os
2. import string
3. import pickle
4. import numpy as np

```

Each input image in the dataset has five associated captions, so we also create a dictionary out of these captions, where each image identifier key results in a list of associated captions. First we open the file `Flickr8k.token.txt`, which contains all of the captions data:

**Note:** When reading these lines in Jupyter Notebook we had to increase our data rate limit, but on Google Colab using GPU we did not encounter this issue<sup>6</sup>.

```

1. # load dataset with captions:
2. def get_captions(filename):
3.     # open captions file and return lines
4.     file = open(filename, 'r')
5.     text = file.read()
6.     file.close()
7.     return text
8.
9. # path to captions data
10. # Google Colab version (loaded into Files section)
11. filename = 'Flickr8k.token.txt'
12.
13. # Jupyter file path - put project_data folder in same directory as this file
14. # filename = 'project_data\Flickr8k_text\Flickr8k.token.txt'
15.
16. path = os.path.join(os.getcwd(), filename)
17. # get all captions
18. doc = get_captions(path)
19.
20. # preview original captions text
21. doc.split('\n')[:5]

Out[1]: ['1000268201_693b08cb0e.jpg#0\tA child in a pink dress is climbing up a set of stairs in an e
ntry way .',
 '1000268201_693b08cb0e.jpg#1\tA girl going into a wooden building .',
 '1000268201_693b08cb0e.jpg#2\tA little girl climbing into a wooden playhouse .',
 '1000268201_693b08cb0e.jpg#3\tA little girl climbing the stairs to her playhouse .',
 '1000268201_693b08cb0e.jpg#4\tA little girl in a pink dress going into a wooden cabin .']
```

For each line of text, we remove punctuation, casing, and words that have numbers or 1 letter; we do this to avoid duplicates and to keep the vocabulary to more common words. Additionally we add start and end tokens (< and >), which the model will need when a generating sequences. We pass each cleaned line to the list of the corresponding dictionary entry. We also provide the option to save this dictionary as `captions.pkl`, which can be reloaded later:

```

1. # make captions clean
2. def new_caps(s):
3.     # take out punctuation
4.     s = s.translate(str.maketrans('', '', string.punctuation))
5.     # make lower case (prevent duplicates)
6.     s = s.lower()
7.     # no words with numbers
8.     s = ' '.join(w for w in s.split() if not any(d.isdigit() for d in w))
9.     # no 1 letter words (i.e. s that follows removed apostrophe)
10.    s = ' '.join(w for w in s.split() if len(w)>1)
11.    return s
12.
13. # create dictionary of photo_id and captions
14. captions = dict()
15. # each line
```

```

16. for line in doc.split('\n'):
17.     # skip last empty line
18.     if len(line) < 1:
19.         continue
20.     # split on .jpg to get image name, 1st token is image id
21.     tokens = line.split('.')
22.     img_id = tokens[0]
23.     # check pics with no captions
24.     if '\t' not in tokens[1]:
25.         cap = tokens[1]
26.     # split 2nd token into descriptions, after \t, put in list
27.     else:
28.         cap = tokens[1].split('\t')[1]
29.         # make lower case, remove punctuation
30.         cap = new_caps(cap)
31.         # put start and end tokens
32.         cap = '<' + cap + '>'
33.     # append caption to list of associated image entry
34.     if img_id not in captions:
35.         captions[img_id] = list()
36.     captions[img_id].append(cap)
37.
38. # test captions dictionary
39. print(captions['997722733_0cb5439472'])
40.
41. # save captions dictionary in file
42. pickle.dump(captions, open('captions.pkl', 'wb'))

```

```

Out[2]: ['< man in pink shirt climbs rock face >',
        '< man is rock climbing high in the air >',
        '< person in red shirt climbing up rock face covered in assist handles >',
        '< rock climber in red shirt >',
        '< rock climber practices on rock climbing wall >']

```

## Visualizing Vocabulary –

Associated file: `model_create.ipynb`

```

1. import nltk
2. import pandas as pd
3. import matplotlib.pyplot as plt

```

We next visualize the vocabulary in our captions input text to assess whether it is worthwhile to limit the vocabulary size in our model, since noisier vocabulary with lots of uncommon words may result in longer model training<sup>3</sup>. First we go through all the text in the captions dictionary created above, and get a set of all of the words; we found that there were 8761 unique words in the captions “corpus”, and 413286 total words:

```

1. # create vocabulary out of all captions (8000 x # of captions per each)
2. vocab = set()
3. # go through each image caption list
4. for index in captions.keys():
5.     # add words to set from all captions
6.     [vocab.update(c.split()) for c in captions[index]]
7.
8. # remove start and end tokens for visualize
9. vocab.remove('<')
10. vocab.remove('>')
11.

```

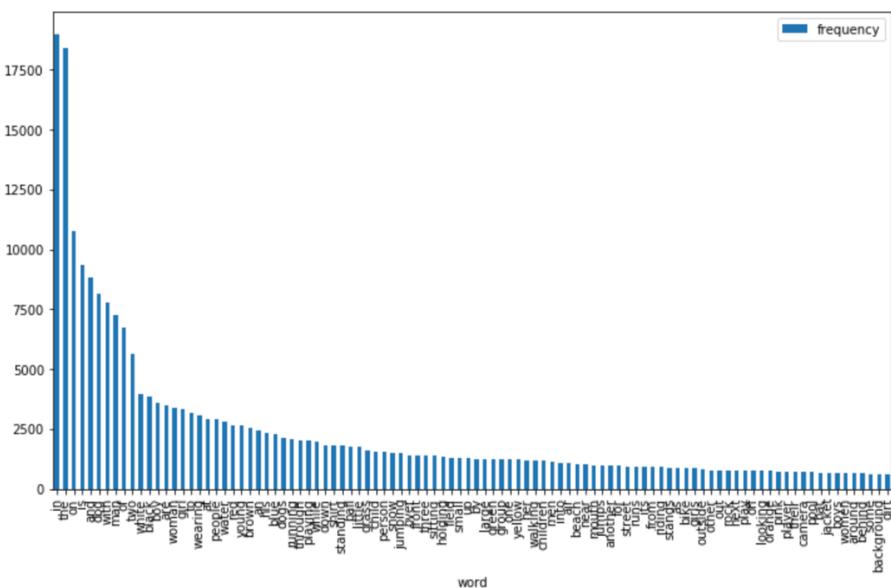
```
12. # get vocab size
13. print('Vocab size:', len(vocab))
14.
15. # put all text in string - to get frequencies
16. txt = ''
17. # go through each image caption list
18. for index in captions.keys():
19.     # add words to set from all captions
20.     t = ''.join(captions[index])
21.     # remove tokens from count
22.     t = t.translate(str.maketrans(' ', ' ', '<>'))
23.     txt += t
24.
25. txt = txt.split(' ')
26. print('Total words:', len(txt))
```

Vocab size: 8761  
Total words: 413286

Using the `nltk` library, we create frequency distributions of the top 100 words in the captions text data, and plot these results:

```
1. # get word frequencies
2. word_ct = nltk.FreqDist(txt)
3.
4. # look at top 5 words
5. print(word_ct.most_common(5))
6.
7. # select top 100 words (index 1 hold space so exclude)
8. wc = word_ct.most_common(101)[1:]
9.
10. wc = pd.DataFrame(wc, columns=['word', 'frequency'])
11. # plot frequencies
12. wc.plot.bar(x='word', y='frequency', figsize=(12,7))
```

```
[(' ', 40458), ('in', 18972), ('the', 18418), ('on', 10741), ('is', 9344)]  
Out[171]: <matplotlib.axes._subplots.AxesSubplot at 0x7fd412dc6198>
```



We see that a large bulk of the words are taken from the top 100 words, so we determine that we can safely reduce the vocabulary size, which can help speed up model training without detrimentally impacting the model quality.

## Loading Training Data and Tokenizing –

Associated file: `model_create.ipynb`

```
1. import keras
2. from keras.preprocessing.text import Tokenizer
3. from keras.preprocessing.sequence import pad_sequences
4. from keras.utils import to_categorical
```

Now we create functions to load the training (and later validation and test) data for our model. The file `Flickr_8k.trainImages.txt` in the text data folder indicates which 6000 images can be used for training the network; we create a function to get a list of these image identifiers. We then retrieve the associated training images and captions by going through the captions and image features dictionaries (created in the Image and Text Preparation sections above), and including entries that have keys with training image identifiers. The result is training image and caption dictionaries:

**Note:** Place `img_features.pkl` in the same directory as `model_create.ipynb` to run the code below:

```
1. # Google Colab version (loaded into Files section)
2. filename = 'Flickr_8k.trainImages.txt'
3.
4. # Jupyter file path - put project_data folder in same directory as this file
5. # filename = 'project_data\Flickr8k_text\Flickr_8k.trainImages.txt'
6.
7. f = os.path.join(os.getcwd(), filename)
8.
9. # get (train) photo ids - put in list
10. def get_ids(f):
11.     # open file - list of training image ids
12.     doc = get_captions(f)
13.     train_ids = list()
14.     # by line
15.     for line in doc.split('\n'):
16.         # skip empty lines
17.         if len(line) < 1:
18.             continue
19.         # take out jpg ending
20.         train_ids.append(line.split('.')[0])
21.     return train_ids
22.
23. # get captions and images with associated ids
24. def get_data(f, captions, img_file):
25.     idens = get_ids(f)
26.     # for each id, get associated captions
27.     caps = {i: list(captions[i]) for i in idens}
28.     # load image features file from pretrained CNN
29.     all_imgs = pickle.load(open(img_file, 'rb'))
30.     # associated id, get image feature
31.     imgs = {i: all_imgs[i] for i in idens}
32.     return idens, caps, imgs
33.
```

```

34. # get training data
35. train_ids, train_caps, train_imgs = get_data(f, captions, 'img_features.pkl')

```

We then have to encode the words (creating an index that maps words to integers) by utilizing the Tokenizer class in Keras. We first convert the training captions dictionary into a list (so it can be passed to the tokenizer as a “corpus” of sorts). Based on the results in the Visualizing Vocabulary section, we decide to reduce our vocabulary size to the top 3500 words. We create an index (which includes our start and end tokens) and can save this tokenizer as `tokenizer.pkl`, which can be loaded later (such as in the file, `model_demo.ipynb`):

```

1. # turn caption dictionary to caption list to encode
2. def dict_list(caps):
3.     cap_list = list()
4.     for index in caps.keys():
5.         [cap_list.append(c) for c in caps[index]]
6.     return cap_list
7.
8. # turn training captions into list
9. cap_list = dict_list(train_caps)
10.
11. # use reduced vocabulary size
12. n_words = 3500
13. # include <, > as words in tokenizer, create index
14. tokens = Tokenizer(num_words = n_words, filters='!"#$%&()*+,-./:;=?@[\\]^_`{|}~\\t\\n')
15.
16. # vectorize words
17. tokens.fit_on_texts(cap_list)
18.
19. # add 1 for 0 padding
20. vocab_size = n_words + 1
21.
22. # can save tokenizer for later
23. pickle.dump(tokens, open('tokenizer.pkl', 'wb'))

```

## Word Embeddings –

Associated file: `model_create.ipynb`

```
1. import io
```

In our model, the text input will be passed through an Embedding layer, which provides a dense, low-dimensional vector representation of the words based on semantics<sup>7</sup>. We decide to use pretrained word embeddings, since our model will generally be looking at generic everyday images (and thus common word/semantic structures), based on the dataset content<sup>8</sup>.

For our pre-trained word embedding space, we use the 300-dimensional GloVe word vector representations from the file `glove.6B.300d.txt`, which we can download [here](#) and put in the same directory as `model_create.ipynb`<sup>9</sup>. We create an embedding dimension of 300, go through our tokenizer word index, and add the vector representation of each word in our 3500-word vocabulary into an embedding matrix; this matrix will be passed into the Embedding layer of our model and frozen, so these pretrained representations are unchanged during training. We can save our matrix as `em3500_300d.npy`, so we can reload rather than recreate if needed.

```

1. # use pretrained word embeddings - follow lecture examples
2. embeddings_index = {}
3. # Google Colab file

```

```

4. with io.open('/content/drive/My Drive/nlp_project_data/glove.6B.300d.txt', encoding='utf8') as f:
5.     # Jupyter file path
6.     # with io.open('glove.6B.300d.txt', encoding='utf8') as f:
7.         for line in f:
8.             values = line.split()
9.             word = values[0]
10.            coefs = np.asarray(values[1:], dtype='float32')
11.            embeddings_index[word] = coefs
12.
13. # embedding dimension layer
14. embedding_dim = 300
15.
16. # create embedding matrix with weights, of appropriate size
17. embedding_matrix = np.zeros((vocab_size, embedding_dim))
18. # go through tokenizer
19. for w, i in tokens.word_index.items():
20.     embedding_vector = embeddings_index.get(w)
21.     # if within vocab size
22.     if i < vocab_size:
23.         if embedding_vector is not None:
24.             # add vector representation
25.             embedding_matrix[i] = embedding_vector
26.
27. # save embedding matrix
28. np.save('em3500_300d.npy', embedding_matrix)

```

## Model Creation and Training –

Associated file: model\_create.ipynb

```

1. from keras import optimizers
2. from keras.models import Model
3. from keras import backend as K
4. from keras.layers.merge import add
5. from keras.models import load_model
6. from keras.callbacks import ModelCheckpoint, EarlyStopping
7. from keras.layers import Input, Dense, Embedding, Dropout, RepeatVector, Bidirectional, GRU

```

Now that we have loaded our training data (both images and text), we begin building our caption generator model. To manage the amount of memory used during training, due to the size of data being processed, we employ a data generator function and use `fit_generator()` for a model in Keras; only the current batch of data being processed is stored in memory, rather than all of the data<sup>10</sup>.

We create a generator function which prepares all input (text and images) and output combinations for a specific batch size. During training, we go through each index and encode the current training captions text set (including padding for a fixed sized input), prepare the image features input, create classification labels for the actual word outputs, and generate all input-output pairs; when we reach the desired batch size, we yield this prepared data for training. We will create training and validation generators. We adapt a common generator structure design<sup>2</sup>:

```

1. # generator to run training - can alter batch size
2. def generator(captions, imgs, tokens, max_len, vocab_size, b_size=1):
3.     # store 2 inputs and 1 output
4.     inp_txt, inp_img, out = list(), list(), list()
5.     n = 0
6.     while 1:

```

```

7.      # go through captions and ids
8.      for index, caps in captions.items():
9.          n +=1
10.         # prep photo input
11.         img = np.array(imgs[index])
12.         # for each caption
13.         for cap in caps:
14.             # encode captions
15.             encode = tokens.texts_to_sequences([cap])[0]
16.             # create all possible pairs in a caption
17.             for i in range(1, len(encode)):
18.                 # get all possible input/output sequences
19.                 in_seq, out_seq = encode[:i], encode[i]
20.                 # pad so input is fixed length
21.                 in_seq = pad_sequences([in_seq], maxlen=max_len)[0]
22.                 # create classification labels for output
23.                 out_seq = to_categorical([out_seq], num_classes = vocab_size)[0]
24.                 inp_img.append(img)
25.                 inp_txt.append(in_seq)
26.                 out.append(out_seq)
27.             # return if reach batch size
28.             if n == b_size:
29.                 yield[[np.array(inp_img), np.array(inp_txt)], np.array(out)]
30.                 # reset
31.                 inp_txt, inp_img, out = list(), list(), list()
32.             n=0

```

We also define a function to calculate the maximum length, based on the largest caption in the training dataset; when the model is outputting a caption sequence, it is either looking for the end token or reaching a word limit.

```

1. # look at longest caption - use as max length in model
2. def max_length(caps):
3.     return max(len(c.split()) for c in dict_list(caps))
4.
5. # get max_len of sequence for training
6. max_len = max_length(train_caps)

```

Next we design our actual model structure, which takes an initial encoded word input (i.e. start token) and an image features input vector, and generates the next word in a sequence; the model will continue to pass the currently constructed sequence and image as inputs until a full caption is generated<sup>11</sup>.

We experimented with different model structures, complexities, and hyperparameters, but we only present the final implemented model here. For this model structure we chose to merge the text and image inputs first, before passing to a recurrent neural network (RNN) cell; however we also saw examples of the text input processed in the RNN first before being merged with the image features<sup>11</sup>.

We previously extracted image features in the file `project_img-features.ipynb` using a pretrained CNN, so we use training data in `img_features.pkl` as the image input into our model. Using functional API, we further process this input by reducing from a 2048 to 300 dimension representation with a Dense layer; we also use a Repeat Vector to ensure that we can merge this vector with the text vector later in the model.

```
1. # take photo features of (2048, ) input
```

```

2. img_inputs = Input(shape=(2048,))
3. # handle overfitting
4. feats1 = Dropout(0.5)(img_inputs)
5. # reduce size
6. feats2 = Dense(300, activation='relu')(feats1)
7. # so same size as text input when merge
8. feats3 = RepeatVector(max_len)(feats2)

```

The encoded text input is passed into an Embedding layer, where we freeze the pretrained GloVe weights during training. The output of this process is also represented in 300 dimensions:

```

1. # pass text sequences
2. # expect max_len inputs since padded
3. txt_inputs = Input(shape=(max_len,))
4. # put in embedding - size of vocab, 0 is padding
5. # add embedding matrix, freeze (don't update weights)
6. seq1 = Embedding(vocab_size, 300, weights=[embedding_matrix], trainable=False,
7.                  mask_zero=True)(txt_inputs)

```

Finally, we concatenate the image and text input vectors together, and pass them into a RNN layer, specifically a bidirectional GRU layer with 500 memory units (actually 1000 since we are creating hidden layers in opposite directions). We then feed this output into a final output Dense classification layer the size of the vocabulary, that predicts the next encoded word in the sequence, using softmax activation. Furthermore, we insert Dropout layers before the image processing and RNN layers, in attempts to alleviate some overfitting<sup>3</sup>. From different training runs, we also discovered that starting with an initial learning rate of 0.0005 when using an Adam optimizer resulted in improved performance. We initialize and compile this model and preview our structure:

```

1. # merge inputs
2. merge1 = add([feats3, seq1])
3. # handle overfitting
4. merge2 = Dropout(0.5)(merge1)
5. # gru model to read sequences, make bidirectional
6. merge3 = Bidirectional(GRU(500, return_sequences=False))(merge2)
7. # generate classification/output layer
8. outputs = Dense(vocab_size, activation='softmax')(merge3)
9.
10. # pass inputs and outputs into model
11. model = Model(inputs=[img_inputs, txt_inputs], outputs=outputs)
12.
13. # look at layers
14. print(model.summary())
15.
16. # compile model - set loss, optimizer, learning rate
17. model.compile(loss='categorical_crossentropy', optimizer=optimizers.Adam(lr=0.0005))

```

Layer (type)	Output Shape	Param #	Connected to
<hr/>			
input_1 (InputLayer)	(None, 2048)	0	
dropout_1 (Dropout)	(None, 2048)	0	input_1[0][0]
dense_1 (Dense)	(None, 300)	614700	dropout_1[0][0]
input_2 (InputLayer)	(None, 34)	0	
repeat_vector_1 (RepeatVector)	(None, 34, 300)	0	dense_1[0][0]
embedding_1 (Embedding)	(None, 34, 300)	1050300	input_2[0][0]
add_1 (Add)	(None, 34, 300)	0	repeat_vector_1[0][0] embedding_1[0][0]
dropout_2 (Dropout)	(None, 34, 300)	0	add_1[0][0]
bidirectional_1 (Bidirectional)	(None, 1000)	2403000	dropout_2[0][0]
dense_2 (Dense)	(None, 3501)	3504501	bidirectional_1[0][0]
<hr/>			
Total params:	7,572,501		
Trainable params:	6,522,201		
Non-trainable params:	1,050,300		

As mentioned above, we conducted training using different model variations and found this model structure produced the best results. Prior models included the following variations: Dropout layers in different locations, passing just the text input into a RNN structure and merging these results with the image features input, fewer and more units in the Dense and RNN layers, multiple RNN layers, lack of bidirectionality, LSTM versus GRU layer, use of default learning rate, etc. Additionally, we adjusted different parameters in training as well, such as epochs and batch size. Again, the procedure performed in this report presents the training conducted to create our final model.

First we load our validation dataset images and captions, following the same procedure used to set up our training data in the Loading Training Data and Tokenizing section above:

```

1. # Google Colab version (loaded into Files section)
2. val_file = 'Flickr_8k.devImages.txt'
3.
4. # Jupyter file path - put project_data folder in same directory as this file
5. # val_file = os.path.join(os.getcwd(), 'project_data\Flickr8k_text\Flickr_8k.devImages')
6.
7. # get validation data
8. val_ids, val_caps, val_imgs = get_data(val_file, captions, 'img_features.pkl')
```

We set up our training and validation generators and first train our model on 5 epochs, using a batch size of 5. Experimenting with batch size in particular was informative in improving initial training (i.e. larger batch sizes trained faster but often led to faster overfitting). We saved the current model state, along with the training and validation losses:

```

1. # define epochs and batches
2. epochs = 5
3. pics_batch = 5
4. t_steps = len(train_caps)//pics_batch
5. v_steps = len(val_caps)//pics_batch
6.
7. # create generators
8. train_g = generator(train_caps, train_imgs, tokens, max_len, vocab_size, pics_batch)
```

```

9. val_g = generator(val_caps, val_imgs, tokens, max_len, vocab_size, pics_batch)
10.
11. # run training
12. history = model.fit_generator(train_g, steps_per_epoch=t_steps, epochs=epochs,
13.                               validation_data=val_g, validation_steps=v_steps)
14.
15. # save first part of model, and history
16. model.save('mod-fin.h5')
17. with open('mod-fin_history', 'wb') as handle:
18.     pickle.dump(history.history, handle)

```

```

Epoch 1/5
1200/1200 [=====] - 183s 152ms/step - loss: 4.4144 - val_loss: 3.7254
Epoch 2/5
1200/1200 [=====] - 179s 149ms/step - loss: 3.6171 - val_loss: 3.4637
Epoch 3/5
1200/1200 [=====] - 179s 149ms/step - loss: 3.3415 - val_loss: 3.3691
Epoch 4/5
1200/1200 [=====] - 180s 150ms/step - loss: 3.1585 - val_loss: 3.3310
Epoch 5/5
1200/1200 [=====] - 179s 149ms/step - loss: 3.0128 - val_loss: 3.3170

```

Next we train on 5 more epochs (with the same batch size) with a further reduced learning rate, to allow for more fine-tuned weight adjustments<sup>12</sup>. Again, we save model weights and history:

```

1. # reduce learning rate
2. K.set_value(model.optimizer.lr, 0.0001)
3.
4. history = model.fit_generator(train_g, steps_per_epoch=t_steps, epochs=epochs,
5.                               validation_data=val_g, validation_steps=v_steps)
6.
7. model.save('mod-fin-1.h5')
8. with open('mod-fin-1_history', 'wb') as handle:
9.     pickle.dump(history.history, handle)

```

```

Epoch 1/5
1200/1200 [=====] - 182s 151ms/step - loss: 2.7795 - val_loss: 3.2341
Epoch 2/5
1200/1200 [=====] - 182s 152ms/step - loss: 2.7232 - val_loss: 3.2269
Epoch 3/5
1200/1200 [=====] - 180s 150ms/step - loss: 2.6858 - val_loss: 3.2221
Epoch 4/5
1200/1200 [=====] - 180s 150ms/step - loss: 2.6524 - val_loss: 3.2208
Epoch 5/5
1200/1200 [=====] - 178s 149ms/step - loss: 2.6249 - val_loss: 3.2189

```

We run our training on 5 more epochs, increasing batch size to 10, which has been discussed to sometimes improve performance in later training<sup>10</sup>. We employ model checkpoints and early stopping, only saving models that continued to reduce validation loss, and terminating training when improvements were not made after 2 epochs. We save the final loss history:

```

1. epochs = 5
2. pics_batch = 10
3. t_steps = len(train_caps)//pics_batch
4. v_steps = len(val_caps)//pics_batch
5.
6. train_g = generator(train_caps, train_imgs, tokens, max_len, vocab_size, pics_batch)
7. val_g = generator(val_caps, val_imgs, tokens, max_len, vocab_size, pics_batch)
8.

```

```

9. # save models with improved performance
10. filepath = 'mod-fin-ep{epoch:03d}-loss{loss:.3f}-val_loss{val_loss:.3f}.h5'
11.
12. # define callbacks for tracking loss
13. checkpoint = ModelCheckpoint(filepath, monitor='val_loss', verbose=1, save_best_only=True, mode='min')
14. early = EarlyStopping(monitor='val_loss', patience=2, verbose=1)
15.
16. history = model.fit_generator(train_g, steps_per_epoch=t_steps, epochs=epochs,
17.                               validation_data=val_g, validation_steps=v_steps,
18.                               callbacks=[early, checkpoint])
19.
20. with open('mod-fin-1-1_history', 'wb') as handle:
21.     pickle.dump(history.history, handle)

```

```

Epoch 1/5
600/600 [=====] - 123s 205ms/step - loss: 2.5863 - val_loss: 3.2138

Epoch 00001: val_loss improved from inf to 3.21377, saving model to mod-fin-ep001-loss2.588-val_loss3.214.h5
Epoch 2/5
600/600 [=====] - 123s 205ms/step - loss: 2.5686 - val_loss: 3.2147

Epoch 00002: val_loss did not improve from 3.21377
Epoch 3/5
600/600 [=====] - 122s 203ms/step - loss: 2.5537 - val_loss: 3.2158

Epoch 00003: val_loss did not improve from 3.21377
Epoch 00003: early stopping

```

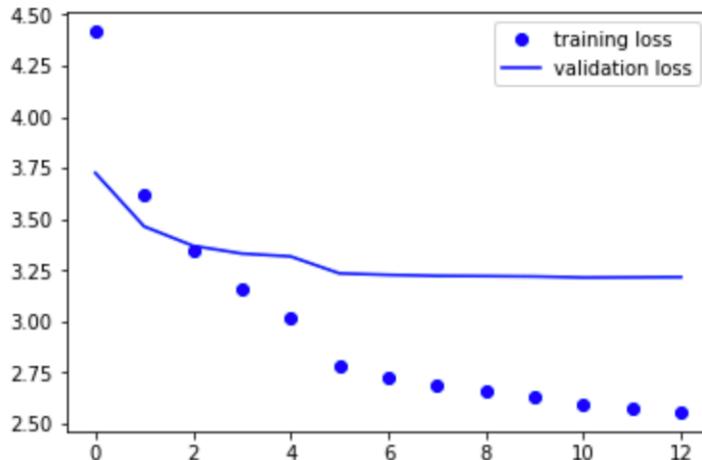
We see that our final model is saved as `mod-fin-ep001-loss2.588-val_loss3.214.h5`, and we stopped training after 3 more epochs; our full training was 13 epochs, with a final validation loss at around 3.21. Through our various training runs and using different structures, complexities, and hyperparameters, we could not really reach losses below this threshold.

To plot training and validations losses, we reload the saved loss histories:

```

1. # includes losses from all epochs
2. losses1 = pickle.load(open('mod-fin_history', 'rb'))
3. losses2 = pickle.load(open('mod-fin-1_history', 'rb'))
4. # alt - losses3 = pickle.load(open('mod-fin-1-1_history', 'rb'))
5. losses3 = history.history
6.
7. # combine losses
8. loss = losses1['loss'] + losses2['loss'] + losses3['loss']
9. val_loss = losses1['val_loss'] + losses2['val_loss'] + losses3['val_loss']
10.
11. def plot_loss(loss, val):
12.     epochs=range(len(loss))
13.     plt.figure()
14.     plt.plot(epochs, loss, 'bo', label='training loss')
15.     plt.plot(epochs, val, 'b', label='validation loss')
16.     plt.legend()
17.     plt.show()
18.
19. # plot training and validation losses
20. plot_loss(loss, val_loss)

```



We observe that the training loss declines as expected, and the validation loss also declines before starting to level off at around 3.2; after the second epoch, the validation loss remains larger than the training loss, which indicates overfitting in the model. Despite adjusting model complexity and parameters in previous training runs (such as reducing neurons, adjusting model structure, learning rates, batching, etc.) and playing around with regularization techniques (in particular Dropout layers), we ultimately found that our networks kept overfitting. This indicates that we may need to use a larger dataset to improve our training.

## Model Evaluation –

Associated file: `model_create.ipynb`

```
1. import matplotlib.pyplot as plt
2. from nltk.translate.bleu_score import corpus_bleu
```

We now evaluate our model on the test dataset, utilizing corpus BLEU scores, which compare a predicted sequence against an actual reference sentence; perfect scores are 1.0, and complete mismatches are 0.0<sup>13</sup>. It has been shown that caption generator models can be evaluated by looking at BLEU-(1,2,3,4) cumulative scores, by counting all matching n-grams in a sequence to determine overlap between generated and actual captions<sup>11</sup>.

First we load the test dataset captions and image features:

```
1. # Google Colab version (loaded into Files section)
2. test_file = 'Flickr_8k.testImages.txt'
3.
4. # Jupyter file path - put project_data folder in same directory as this file
5. # test_file = os.path.join(os.getcwd(), 'project_data\Flickr8k_text\Flickr_8k.testImage
   s')
6.
7. # get test data
8. test_ids, test_caps, test_imgs = get_data(test_file , captions, 'img_features.pkl')
```

We also define a function to decode output text inputs (from integers back to words), using the word index from our tokenizer. We reload our saved final model (placed in the same directory as `model_create.ipynb`), so we can run predictions:

```
1. # covert index to word
```

```

2. def idtw(num, tokens):
3.     for w, i in tokens.word_index.items():
4.         if i == num:
5.             return w
6.     return None
7.
8. # load best model
9. # Google Colab version (loaded into Files section)
10. model_name = 'mod-fin-ep001-loss2.588-val_loss3.214.h5'
11. mod = load_model(model_name)

```

To generate caption predictions, we create a function to which we can pass an image feature input and start token, and receive an output sequence. Specifically, this method performs a greedy search to generate a sequence by selecting the word with the highest probability from the softmax activation when creating a caption<sup>14</sup>:

```

1. # predict captions - greedy version
2. def make_cap(model, tokens, img, max_len):
3.     # start token
4.     in_txt='<'
5.     # go over entire possible sequence (of max length)
6.     for i in range(max_len):
7.         # encode input - similar set up to generator
8.         seq = tokens.texts_to_sequences([in_txt])[0]
9.         # pad
10.        seq = pad_sequences([seq], maxlen=max_len)
11.        # get next word in sequence
12.        pred = model.predict([np.array(img), seq])
13.        # get encoded of highest probability
14.        pred = np.argmax(pred)
15.        # decode output
16.        w = idtw(pred, tokens)
17.        # break if cannot map
18.        if w is None:
19.            break
20.        # append to sequence
21.        in_txt += ' ' + w
22.        # stop if reach end token
23.        if w == '>':
24.            break
25.    return in_txt

```

We can test this method on an image in the test dataset, by printing the caption and plotting the actual original JPEG image:

```

1. # Google Colab version (loaded into Files section)
2. i_path='2654514044_a70a6e2c21.jpg'
3.
4. # Jupyter file path - put project_data folder in same directory as this file
5. # i_path = os.path.join(os.getcwd(), 'project_data\Flickr8k_Dataset\Flicker8k_Dataset\\')
6. # i_path = os.path.join(img_folder, '2654514044_a70a6e2c21.jpg')
7.
8. # in case just loading model without running training
9. max_len = max_length(train_caps)
10.
11. # plot image
12. x=plt.imread(i_path)

```

```

13. plt.imshow(x)
14.
15. # get caption - take out start/end tokens
16. print(make_cap(mod, tokens, test_imgs['2654514044_a70a6e2c21'].reshape(1,2048), max_len
)[:-1])

```

brown dog is running through the grass



We follow a procedure to generate corpus BLEU scores similar to what was discussed in lecture, employing greedy search predictions with the test data:

```

1. # evaluate model on test images using greedy search
2. def bleu(mod, test_caps, test_imgs, tokens, max_len):
3.     act, pred = list(), list()
4.     for i, caps in test_caps.items():
5.         img = test_imgs[i].reshape(1,2048)
6.         p = make_cap(mod, tokens, img, max_len)
7.         act.append([c.split() for c in caps])
8.         pred.append(p.split())
9.     print('BLEU-1: %f' % corpus_bleu(act, pred, weights=(1.0, 0, 0, 0)))
10.    print('BLEU-2: %f' % corpus_bleu(act, pred, weights=(0.5, 0.5, 0, 0)))
11.    print('BLEU-3: %f' % corpus_bleu(act, pred, weights=(0.3, 0.3, 0.3, 0)))
12.    print('BLEU-4: %f' % corpus_bleu(act, pred, weights=(0.25, 0.25, 0.25, 0.25)))
13.
14. bleu(mod, test_caps, test_imgs, tokens, max_len)

```

```

BLEU-1: 0.583389
BLEU-2: 0.348187
BLEU-3: 0.250425
BLEU-4: 0.128517

```

The BLEU-1 score is the best (close to .6) with subsequent cumulative n-gram scores declining; the overall scores are slightly lower but generally comparative to results seen in previous models<sup>11</sup>. Perhaps a larger training set could result in improved performance. Additionally, since we are using a reduced vocabulary based only on training captions, it is unlikely we will be able to generate exact test caption matches.

An alternative approach, beam search, has been proposed to generate caption sequences<sup>3</sup>. Instead of picking the most probable next word, we instead consider and store multiple alternatives and track possible sequences generated at each time step<sup>15</sup>. Some research has demonstrated that beam search may sometimes generate better captions, so we incorporate a certain implementation of this process to see if there are improvements<sup>16</sup>.

We adapt a function to generate beam search captions below, which runs the model on start token and image features inputs, and goes through and stores the top possibilities (beam\_index) at each time step until we finally select the highest probable sequence at the end:

```

1. # create beam search
2. def beam_caps(model, tokens, img, max_len, beam_index=3):
3.     # start sequence
4.     in_txt = [[tokens.texts_to_sequences(['<'])[0], 0.0]]
5.
6.     # while less than max length
7.     while len(in_txt[0][0]) < max_len:
8.         temp = []
9.         # go through sequence
10.        for i in in_txt:
11.            # pad input sequences
12.            seq = pad_sequences([i[0]], maxlen=max_len, padding='post')
13.            # get predictions for next word
14.            pred = model.predict([np.array(img), seq])
15.            # sort and pick top beam_index possibilities
16.            beam_preds = np.argsort(pred[0])[-beam_index:]
17.
18.            # create new lists from each top pick - words and probabilities
19.            # build each sequence option
20.            for w in beam_preds:
21.                # get current seq, prob
22.                new_cap, prob = i[0][:], i[1]
23.                # add word and associated probability
24.                new_cap.append(w)
25.                prob += pred[0][w]
26.                temp.append([new_cap, prob])
27.
28.        # hold possibilities
29.        in_txt = temp
30.        # sort probabilities - lambda identifies value
31.        in_txt = sorted(in_txt, key=lambda p: p[1])
32.        # pick top in full list of options
33.        in_txt = in_txt[-beam_index:]
34.
35.        # pick top sequence in all
36.        in_txt = in_txt[-1][0]
37.        # decode
38.        inter_cap = [idtw(i, tokens) for i in in_txt]
39.
40.        # construct caption
41.        final_cap = []
42.        for i in inter_cap:
43.            # if not end token
44.            if i != '>':
45.                final_cap.append(i)
46.            else:
47.                break
48.        # take out start token, make string
49.        final_cap = ' '.join(final_cap[1:])
50.    return final_cap

```

We use another test image to run a function that generates caption predictions with greedy search, along with beam search results using indices of 3, 5, and 7:

```

1. # Google Colab version (loaded into Files section)
2. i_path='1322323208_c7ecb742c6.jpg'
3.
4. # Jupyter file path - put project_data folder in same directory as this file
5. # i_path = os.path.join(os.getcwd(), 'project_data\Flickr8k_Dataset\Flicker8k_Dataset\\')
6. # i_path = os.path.join(img_folder, '1322323208_c7ecb742c6.jpg')
7.
8. # acquire greedy and beam searches and display an image
9. def print_img_caps(path, mod, tokens, img, max_len):
10.    # plot image
11.    x=plt.imread(path)
12.    plt.imshow(x)
13.    # reshape image features
14.    img = img.reshape(1,2048)
15.    # get greedy caption
16.    txt = make_cap(mod, tokens, img, max_len)[1:-1]
17.    # get beam searches
18.    print('Greedy:', txt)
19.    print('Beam, k=3:', beam_caps(mod, tokens, img, max_len))
20.    print('Beam, k=5:', beam_caps(mod, tokens, img, max_len, beam_index=5))
21.    print('Beam, k=7:', beam_caps(mod, tokens, img, max_len, beam_index=7))
22.
23. print_img_caps(i_path, mod, tokens, test_imgs['1322323208_c7ecb742c6'], max_len)

```

Greedy: man in black shirt is standing on the beach with the sun behind him  
 Beam, k=3: two people are standing on the beach at dusk  
 Beam, k=5: two people are sitting on the beach next to the ocean  
 Beam, k=7: young boy running on the beach at sunset



In the above example, we see that beam search (in particular with index 3) results in a preferable caption based on the image, compared to the greedy search output. We decide to re-run our BLEU scores with beam search index 3 predictions to see if there is an improvement:

**Note:** The code below takes a while to run.

```

1. # evaluate model on test images using beam search, beam_index=3
2. def bleu_2(mod, test_caps, test_imgs, tokens, max_len):
3.     act, pred = list(), list()
4.     for i, caps in test_caps.items():
5.         img = test_imgs[i].reshape(1,2048)
6.         p = beam_caps(mod, tokens, img, max_len)
7.         # take out start/end tokens to compare
8.         act.append([c[1:-1].split() for c in caps])
9.         pred.append(p.split())
10.    print('BLEU-1: %f' % corpus_bleu(act, pred, weights=(1.0, 0, 0, 0)))
11.    print('BLEU-2: %f' % corpus_bleu(act, pred, weights=(0.5, 0.5, 0, 0)))
12.    print('BLEU-3: %f' % corpus_bleu(act, pred, weights=(0.3, 0.3, 0.3, 0)))

```

```

13.     print('BLEU-4: %f' % corpus_bleu(act, pred, weights=(0.25, 0.25, 0.25, 0.25)))
14.
15. bleu_2(mod, test_caps, test_imgs, tokens, max_len)

```

```

BLEU-1: 0.488283
BLEU-2: 0.307064
BLEU-3: 0.227320
BLEU-4: 0.116195

```

We see that these BLEU scores are lower than the ones we calculated above. Although we saw that beam search may sometimes result in more accurate captions (when looking at specific examples), they may not exactly match actual captions in the dataset, which could explain the lower scores.

Next we create a demo file, `model_demo.ipynb`, which will plot images and generate greedy and beam search captions for a random set of test images in the dataset. We save the image features dictionary for the test data so we can reuse them in the demo:

```

1. # save test imgs dictionary in file
2. pickle.dump(test_imgs, open('test_imgs.pkl', 'wb'))

```

## Demonstration –

Associated file: `model_demo.ipynb`

```

1. import os
2. import keras
3. import random
4. import string
5. import pickle
6. import numpy as np
7. import matplotlib.pyplot as plt
8. from keras.models import load_model
9. from keras.preprocessing.text import Tokenizer
10. from keras.preprocessing.sequence import pad_sequences

```

We want to see actual generated greedy and beam search captions for some images, so in the code below, we select 5 random test images to plot and produce captions. We load some of the saved dictionaries and files generated in the file `model_create.ipynb` (and placed in the same directory): `tokenizer.pkl`, `test_imgs.pkl`, and `mod-fin-ep001-loss2.588-val_loss3.214.h5`. We also set the maximum caption length based on training captions:

**Note:** Place these files in the same directory as `model_demo.ipynb` to run the code below:

```

1. # based on training of model - determined in model_create.ipynb
2. max_len = 34
3.
4. # load previous saved test image features
5. test_imgs = pickle.load(open('test_imgs.pkl', 'rb'))
6.
7. # load saved tokenizer
8. tokens = pickle.load(open('tokenizer.pkl', 'rb'))
9.
10. # load final model
11. model_name = 'mod-fin-ep001-loss2.588-val_loss3.214.h5'
12. mod = load_model(model_name)

```

Next we include functions created in `model_create.ipynb` to decode integer indices back to words and generate and print greedy and search captions: `idtw()`, `make_cap()`, `beam_caps()`, and `print_img_caps()`. See the Model Evaluation section for the code for these functions.

**Note:** We insert the following code in the `print_img_caps()` method before generating captions, to print the image identifier.

```
1. print('\033[1m', 'Image ID:', '\033[0m', path.split('\\')[-1])
```

We test the methods on an image, and then we select 5 random images from the `test_imgs` dictionary. We create a loop that goes through each image, and plots the image and its generated greedy and beam search captions:

```
1. # set up full images folder
2. img_folder = os.path.join(os.getcwd(), 'project_data\Flickr8k_Dataset\Flicker8k_Dataset
\\')
3.
4. # test image
5. i = '3462454965_a481809cea.jpg'
6. path = os.path.join(img_folder, i)
7. print_img_caps(path, mod, tokens, test_imgs[i.split('.')[0]], max_len)
```

```
Image ID: 3462454965_a481809cea.jpg
Greedy: black dog is running through the grass
Beam, k=3: black dog and brown dog are playing in the grass
Beam, k=5: the black dog is running through the grass
Beam, k=7: the black dog is running through the green grass
```

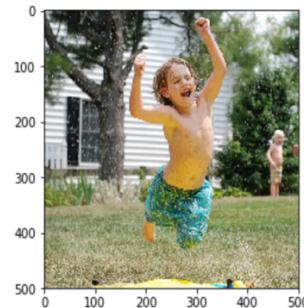


```
1. # to reproduce - can take out if want different images
2. random.seed(8)
3.
4. # pick 5 random images
5. keys = random.sample(list(test_imgs), 5)
6.
7. for k in keys:
8.     img = k + '.jpg'
9.     path = os.path.join(img_folder, img)
10.    print_img_caps(path, mod, tokens, test_imgs[img.split('.')[0]], max_len)
11.    plt.show()
12.    # add space between images
13.    print()
```

**Image ID:** 3708177171\_529bb4ff1d.jpg  
Greedy: skateboarder in the air  
Beam, k=3: skateboarder doing trick on ramp  
Beam, k=5: the skateboarder does trick on ramp  
Beam, k=7: the skateboarder does trick on ramp



**Image ID:** 2718024196\_3ff660416a.jpg  
Greedy: young boy in swim trunks is jumping into pool  
Beam, k=3: young boy wearing blue swim trunks is jumping into swimming pool  
Beam, k=5: young boy in swim trunks dives into pool  
Beam, k=7: young boy in swim trunks jumping into swimming pool



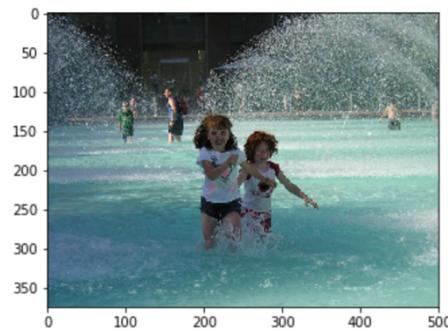
**Image ID:** 2167644298\_100ca79f54.jpg  
Greedy: man in black shirt and tie is smiling  
Beam, k=3: man and woman wearing black are standing in front of store  
Beam, k=5: the man in the black jacket is looking at the camera  
Beam, k=7: man and woman pose for picture in front of brick wall



**Image ID:** 3502993968\_4ee36afb0e.jpg  
Greedy: man in blue shirt and helmet is riding bike on the beach  
Beam, k=3: man on bike is jumping over hill  
Beam, k=5: the man is performing trick on bmx bike  
Beam, k=7: the man is performing trick on bmx bike



**Image ID:** 2453971388\_76616b6a82.jpg  
Greedy: boy in swimming trunks is jumping into pool  
Beam, k=3: young boy jumping into swimming pool  
Beam, k=5: young boy in swim trunks plays in the water  
Beam, k=7: young boy jumping into swimming pool



To generate captions on our own images, we would first have to generate image features from these images as inputs to our model. We would have to pass the images into the InceptionV3 model, running code similar to `project_img-features.ipynb`, to create an image features dictionary, and then we can pass these extracted features into our model to produce caption sequences, as seen above.

### Takeaways and Future Steps –

We were able to create a caption generator model that provides fairly decent captions for the provided input images. Although our final model did not have particularly high corpus BLEU scores, our evaluation metric, when looking at actual predicted caption results, our network generated sequences that generally captured the context of the original images.

Through the model creation and training process, we ran different iterations of model structures and complexities, as well as hyperparameters. We found that starting with lower learning rates, utilizing Dropout layers, and running the RNN on a bidirectional GRU layer (rather than multiple layers, or an LSTM layer) resulted in preferable results. However, more

experimentation with a different structure (i.e. passing only the text input into the RNN) and further hyperparameter tuning could be fruitful, especially with cross-validation.

We also discovered that our final model, as well as our previous training models, ended up overfitting, despite regularization techniques<sup>3</sup>. It appears that we need to conduct training on larger amounts of data, such as the COCO dataset, or perhaps implement data augmentation to improve results<sup>17</sup>. Additional model improvements could include: training on a larger variety of images (versus everyday images, to allow for more generalization, which could improve testing on our own images), further reducing vocabulary size to speed up training, learning word embeddings rather than using pretrained vectors, and implementing an attention mechanism.

**YouTube video link –**  
<https://youtu.be/zPawquEoQFs>

## References –

1. Karpathy, A., & Fei-Fei, L. (2015). Deep Visual-Semantic Alignments for Generating Image Descriptions. Retrieved 27 July 2019, from <https://cs.stanford.edu/people/karpathy/cvpr2015.pdf>.
2. Brownlee, J. (2017). How to Develop a Deep Learning Photo Caption Generator from Scratch. Retrieved 27 July 2019, from <https://machinelearningmastery.com/develop-a-deep-learning-caption-generation-model-in-python/>.
3. Vinyals, O., Toshev, A., Benigo, S., & Erhan, D. (2015). Show and Tell: A Neural Image Caption Generator. Retrieved 27 July 2019, from <https://arxiv.org/pdf/1411.4555.pdf>.
4. Applications - Keras Documentation. Retrieved 27 July 2019, from <https://keras.io/applications/#inceptionv3>.
5. ImageNet. Retrieved 27 July 2019, from <https://en.wikipedia.org/wiki/ImageNet>.
6. IOPub data rate exceeded. The notebook server will temporarily stop sending output to the client in order to avoid crashing it. (2019). Retrieved 27 July 2019, from <https://www.drjamesfroggatt.com/python-and-neural-networks/iopub-data-rate-exceeded-the-notebook-server-will-temporarily-stop-sending-output-to-the-client-in-order-to-avoid-crashing-it/>.
7. Koehrsen, W. (2018). Neural Network Embeddings Explained. Retrieved 27 July 2019, from <https://towardsdatascience.com/neural-network-embeddings-explained-4d028e6f0526>.
8. Using word embeddings. Retrieved 27 July 2019, from <https://jjallaire.github.io/deep-learning-with-r-notebooks/notebooks/6.1-using-word-embeddings.nb.html>.
9. Pennington, J., Socher, R., & Manning, C. (2014). GloVe: Global Vectors for Word Representation. Retrieved 27 July 2019, from <https://nlp.stanford.edu/projects/glove/>.
10. Lamba, H. (2018). Image Captioning with Keras. Retrieved 27 July 2019, from <https://towardsdatascience.com/image-captioning-with-keras-teaching-computers-to-describe-pictures-c88a46a311b8>.
11. Tanti, M., Gatt, A., & Camilleri, K. (2018). Where to put the Image in an Image Caption Generator. Retrieved 27 July 2019, from <https://arxiv.org/pdf/1703.09137.pdf>.
12. Surmenok, P. (2017). Estimating an Optimal Learning Rate For a Deep Neural Network. Retrieved 27 July 2019, from <https://towardsdatascience.com/estimating-optimal-learning-rate-for-a-deep-neural-network-ce32f2556ce0>.

13. Brownlee, J. (2017). A Gentle Introduction to Calculating the BLEU Score for Text in Python. Retrieved 27 July 2019, from <https://machinelearningmastery.com/calculate-bleu-score-for-text-python/>.
14. Why is beam search required in sequence to sequence transduction using recurrent neural networks? - Quora. (2017). Retrieved 27 July 2019, from <https://www.quora.com/Why-is-beam-search-required-in-sequence-to-sequence-transduction-using-recurrent-neural-networks>.
15. Coursera. *Beam Search* [Video]. Retrieved from <https://www.coursera.org/lecture/nlp-sequence-models/beam-search-4EtHZ>.
16. Mustafa, F. (2018). Keras implementation of Image Captioning Model. Retrieved 27 July 2019, from <https://medium.com/@faizanmustafa75/keras-implementation-of-image-captioning-model-3a7ab68e67d4>.
17. COCO - Common Objects in Context. Retrieved 27 July 2019, from <http://cocodataset.org/#overview>.