≡  **Navigation**

Want help with deep learning for text? Take the FREE Mini-Course

Search...

# How to Use Small Experiments to Develop a Caption Generation Model in Keras

by **Jason Brownlee** on November 24, 2017 in **Deep Learning for Natural Language Processing**

Caption generation is a challenging artificial intelligence problem where a textual description must be generated for a photograph.

It requires both methods from computer vision to understand the content of the image and a language model from the field of natural language processing to turn the understanding of the image into words in the right order. Recently, deep learning methods have achieved state of the art results on examples of this problem.

It can be hard to develop caption generating models on your own data, primarily because the datasets and the models are so large and take days to train. An alternative approach is to explore model configurations with a small sample of the fuller dataset.

In this tutorial, you will discover how you can use a small sample of a standard photo captioning dataset to explore different deep model designs.

After completing this tutorial, you will know:

- How to prepare data for photo captioning modeling.
- How to design a baseline and test harness to evaluate the skill of models and control for their stochastic nature.
- How to evaluate properties like model skill, feature extraction models, and word embeddings in order to lift model skill.

Let's get started.

How to Use Small Experiments to Develop a Caption Generation Model in Keras
Photo by Per, some rights reserved.

# Tutorial Overview

This tutorial is divided into 6 parts; they are:

1. Data Preparation
2. Baseline Caption Generation Model
3. Network Size Parameters
4. Configuring the Feature Extraction Model
5. Word Embedding Models
6. Analysis of Results

## Python Environment

This tutorial assumes you have a Python SciPy environment installed, ideally with Python 3.

You must have Keras (2.0 or higher) installed with either the TensorFlow or Theano backend.

The tutorial also assumes you have scikit-learn, Pandas, NumPy, and Matplotlib installed.

If you need help with your environment, see this tutorial:

- [How to Setup a Python Environment for Machine Learning and Deep Learning with Anaconda](#)

I recommend running the code on a system with a GPU.

You can access GPUs cheaply on Amazon Web Services. Learn how in this tutorial:

- [How To Develop and Evaluate Large Deep Learning Models with Keras on Amazon Web Services](#)

Let's dive in.

---

## Need help with Deep Learning for Text Data?

Take my free 7-day email crash course now (with code).

Click to sign-up and also get a free PDF Ebook version of the course.

**Start Your FREE Crash-Course Now**

---

# Data Preparation

First, we need to prepare the dataset for training the model.

We will use the Flickr8K dataset that is comprised of a little more than 8,000 photographs and their descriptions.

You can download the dataset from here:

- [Framing image description as a ranking task: data, models and evaluation metrics](#).

Unzip the photographs and descriptions into your current working directory into *Flicker8k_Dataset* and *Flickr8k_text* directories respectively.

There are two parts to the data preparation, they are:

1. Preparing the Text
2. Preparing the Photos

## Preparing the Text

The dataset contains multiple descriptions for each photograph and the text of the descriptions requires some minimal cleaning.

First, we will load the file containing all of the descriptions.

```
1  # load doc into memory
2  def load_doc(filename):
3      # open the file as read only
4      file = open(filename, 'r')
5      # read all text
6      text = file.read()
7      # close the file
8      file.close()
9      return text
10
11 filename = 'Flickr8k_text/Flickr8k.token.txt'
12 # load descriptions
13 doc = load_doc(filename)
```

Each photo has a unique identifier. This is used in the photo filename and in the text file of descriptions. Next, we will step through the list of photo descriptions and save the first description for each photo. Below defines a function named *load_descriptions()* that, given the loaded document text, will return a dictionary of photo identifiers to descriptions.

```
1  # extract descriptions for images
2  def load_descriptions(doc):
3      mapping = dict()
4      # process lines
5      for line in doc.split('\n'):
6          # split line by white space
7          tokens = line.split()
8          if len(line) < 2:
9              continue
10         # take the first token as the image id, the rest as the description
11         image_id, image_desc = tokens[0], tokens[1:]
12         # remove filename from image id
13         image_id = image_id.split('.')[0]
14         # convert description tokens back to string
15         image_desc = ' '.join(image_desc)
16         # store the first description for each image
17         if image_id not in mapping:
18             mapping[image_id] = image_desc
19     return mapping
20
21 # parse descriptions
22 descriptions = load_descriptions(doc)
23 print('Loaded: %d ' % len(descriptions))
```

Next, we need to clean the description text.

The descriptions are already tokenized and easy to work with. We will clean the text in the following ways in order to reduce the size of the vocabulary of words we will need to work with:

- Convert all words to lowercase.
- Remove all punctuation.
- Remove all words that are one character or less in length (e.g. 'a').

Below defines the *clean_descriptions()* function that, given the dictionary of image identifiers to descriptions, steps through each description and cleans the text.

```
1  import string
2
3  def clean_descriptions(descriptions):
4      # prepare translation table for removing punctuation
```

```
 5        table = str.maketrans('', '', string.punctuation)
 6        for key, desc in descriptions.items():
 7            # tokenize
 8            desc = desc.split()
 9            # convert to lower case
10            desc = [word.lower() for word in desc]
11            # remove punctuation from each token
12            desc = [w.translate(table) for w in desc]
13            # remove hanging 's' and 'a'
14            desc = [word for word in desc if len(word)>1]
15            # store as string
16            descriptions[key] =  ' '.join(desc)
17
18 # clean descriptions
19 clean_descriptions(descriptions)
20 # summarize vocabulary
21 all_tokens = ' '.join(descriptions.values()).split()
22 vocabulary = set(all_tokens)
23 print('Vocabulary Size: %d' % len(vocabulary))
```

Finally, we save the dictionary of image identifiers and descriptions to a new file named *descriptions.txt*, with one image identifier and description per line.

Below defines the *save_doc()* function that given a dictionary containing the mapping of identifiers to descriptions and a filename, saves the mapping to file.

```
 1  # save descriptions to file, one per line
 2  def save_doc(descriptions, filename):
 3      lines = list()
 4      for key, desc in descriptions.items():
 5          lines.append(key + ' ' + desc)
 6      data = '\n'.join(lines)
 7      file = open(filename, 'w')
 8      file.write(data)
 9      file.close()
10
11 # save descriptions
12 save_doc(descriptions, 'descriptions.txt')
```

Putting this all together, the complete listing is provided below.

```
 1  import string
 2
 3  # load doc into memory
 4  def load_doc(filename):
 5      # open the file as read only
 6      file = open(filename, 'r')
 7      # read all text
 8      text = file.read()
 9      # close the file
10      file.close()
11      return text
12
13 # extract descriptions for images
14 def load_descriptions(doc):
15     mapping = dict()
16     # process lines
17     for line in doc.split('\n'):
18         # split line by white space
19         tokens = line.split()
20         if len(line) < 2:
21             continue
```

```
22          # take the first token as the image id, the rest as the description
23          image_id, image_desc = tokens[0], tokens[1:]
24          # remove filename from image id
25          image_id = image_id.split('.')[0]
26          # convert description tokens back to string
27          image_desc = ' '.join(image_desc)
28          # store the first description for each image
29          if image_id not in mapping:
30              mapping[image_id] = image_desc
31      return mapping
32
33  def clean_descriptions(descriptions):
34      # prepare translation table for removing punctuation
35      table = str.maketrans('', '', string.punctuation)
36      for key, desc in descriptions.items():
37          # tokenize
38          desc = desc.split()
39          # convert to lower case
40          desc = [word.lower() for word in desc]
41          # remove punctuation from each token
42          desc = [w.translate(table) for w in desc]
43          # remove hanging 's' and 'a'
44          desc = [word for word in desc if len(word)>1]
45          # store as string
46          descriptions[key] =  ' '.join(desc)
47
48  # save descriptions to file, one per line
49  def save_doc(descriptions, filename):
50      lines = list()
51      for key, desc in descriptions.items():
52          lines.append(key + ' ' + desc)
53      data = '\n'.join(lines)
54      file = open(filename, 'w')
55      file.write(data)
56      file.close()
57
58  filename = 'Flickr8k_text/Flickr8k.token.txt'
59  # load descriptions
60  doc = load_doc(filename)
61  # parse descriptions
62  descriptions = load_descriptions(doc)
63  print('Loaded: %d ' % len(descriptions))
64  # clean descriptions
65  clean_descriptions(descriptions)
66  # summarize vocabulary
67  all_tokens = ' '.join(descriptions.values()).split()
68  vocabulary = set(all_tokens)
69  print('Vocabulary Size: %d' % len(vocabulary))
70  # save descriptions
71  save_doc(descriptions, 'descriptions.txt')
```

Running the example first prints the number of loaded photo descriptions (8,092) and the size of the clean vocabulary (4,484 words).

```
1  Loaded: 8092
2  Vocabulary Size: 4484
```

The clean descriptions are then written to '*descriptions.txt*'. Taking a look in the file, we can see that the descriptions are ready for modeling.

Taking a look in the file, we can see that the descriptions are ready for modeling.

```
1  3621647714_fc67ab2617 man is standing on snow with trees and mountains all around him
2  365128300_6966058139 group of people are rafting on river rapids
3  2751694538_fffa3d307d man and boy sit in the driver seat
4  537628742_146f2c24f8 little girl running in field
5  2320125735_27fe729948 black and brown dog with blue collar goes on alert by soccer ball in the g
6  ...
```

# Preparing the Photos

We will use a pre-trained model to interpret the content of the photos.

There are many models to choose from. In this case, we will use the Oxford Visual Geometry Group or VGG model that won the ImageNet competition in 2014. Learn more about the model here:

- Very Deep Convolutional Networks for Large-Scale Visual Recognition

Keras provides this pre-trained model directly. Note, the first time you use this model, Keras will download the model weights from the Internet, which are about 500 Megabytes. This may take a few minutes depending on your internet connection.

We could use this model as part of a broader image caption model. The problem is, it is a large model and running each photo through the network every time we want to test a new language model configuration (downstream) is redundant.

Instead, we can pre-compute the "photo features" using the pre-trained model and save them to file. We can then load these features later and feed them into our model as the interpretation of a given photo in the dataset. It is no different to running the photo through the full VGG model, it is just that we will have done it once in advance.

This is an optimization that will make training our models faster and consume less memory.

We can load the VGG model in Keras using the VGG class. We will load the model without the top; this means without the layers at the end of the network that are used to interpret the features extracted from the input and turn them into a class prediction. We are not interested in the image net classification of the photos and we will train our own interpretation of the image features.

Keras also provides tools for reshaping the loaded photo into the preferred size for the model (e.g. 3 channel 224 x 224 pixel image).

Below is a function named *extract_features()* that given a directory name will load each photo, prepare it for VGG and collect the predicted features from the VGG model. The image features are a 3-dimensional array with the shape (7, 7, 512).

The function returns a dictionary of image identifier to image features.

```
1  # extract features from each photo in the directory
2  def extract_features(directory):
3      # load the model
4      in_layer = Input(shape=(224, 224, 3))
5      model = VGG16(include_top=False, input_tensor=in_layer)
```

```
6      print(model.summary())
7      # extract features from each photo
8      features = dict()
9      for name in listdir(directory):
10         # load an image from file
11         filename = directory + '/' + name
12         image = load_img(filename, target_size=(224, 224))
13         # convert the image pixels to a numpy array
14         image = img_to_array(image)
15         # reshape data for the model
16         image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
17         # prepare the image for the VGG model
18         image = preprocess_input(image)
19         # get features
20         feature = model.predict(image, verbose=0)
21         # get image id
22         image_id = name.split('.')[0]
23         # store feature
24         features[image_id] = feature
25         print('>%s' % name)
26     return features
```

We can call this function to prepare the photo data for testing our models, then save the resulting dictionary to a file named '*features.pkl*'.

The complete example is listed below.

```
1  from os import listdir
2  from pickle import dump
3  from keras.applications.vgg16 import VGG16
4  from keras.preprocessing.image import load_img
5  from keras.preprocessing.image import img_to_array
6  from keras.applications.vgg16 import preprocess_input
7  from keras.layers import Input
8
9  # extract features from each photo in the directory
10 def extract_features(directory):
11     # load the model
12     in_layer = Input(shape=(224, 224, 3))
13     model = VGG16(include_top=False, input_tensor=in_layer)
14     print(model.summary())
15     # extract features from each photo
16     features = dict()
17     for name in listdir(directory):
18         # load an image from file
19         filename = directory + '/' + name
20         image = load_img(filename, target_size=(224, 224))
21         # convert the image pixels to a numpy array
22         image = img_to_array(image)
23         # reshape data for the model
24         image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
25         # prepare the image for the VGG model
26         image = preprocess_input(image)
27         # get features
28         feature = model.predict(image, verbose=0)
29         # get image id
30         image_id = name.split('.')[0]
31         # store feature
32         features[image_id] = feature
33         print('>%s' % name)
34     return features
35
36 # extract features from all images
```

```
37  directory = 'Flicker8k_Dataset'
38  features = extract_features(directory)
39  print('Extracted Features: %d' % len(features))
40  # save to file
41  dump(features, open('features.pkl', 'wb'))
```

Running this data preparation step may take a while depending on your hardware, perhaps one hour on the CPU with a modern workstation.

At the end of the run, you will have the extracted features stored in '*features.pkl*' for later use.

# Baseline Caption Generation Model

In this section, we will define a baseline model for generating captions for photos and how to evaluate it so that it can be compared to variations on this baseline.

This section is divided into 5 parts:

1. Load Data.
2. Fit Model.
3. Evaluate Model.
4. Complete Example
5. "A" versus "A" Test
6. Generate Photo Captions

## 1. Load Data

We are not going to fit the model on all of the caption data, or even on a large sample of the data.

In this tutorial, we are interested in quickly testing a suite of different configurations of a caption model to see what works on this data. That means we need the evaluation of one model configuration to happen quickly. Toward this end, we will train the models on 100 photographs and captions, then evaluate them on both the training dataset and on a new test set of 100 photographs and captions.

First, we need to load a pre-defined subset of photographs. The provided dataset has separate sets for train, test, and development, which are really just different groups of photo identifiers. We will load the development set and use the first 100 identifiers for train and the second 100 (e.g. from 100 to 200) as the test set.

The function *load_set()* below will load a pre-defined set of identifiers, and we will call it with the '*Flickr_8k.devImages.txt*' filename as an argument.

```
1  # load a pre-defined list of photo identifiers
2  def load_set(filename):
3      doc = load_doc(filename)
4      dataset = list()
5      # process line by line
6      for line in doc.split('\n'):
7          # skip empty lines
8          if len(line) < 1:
9              continue
```

```
10            # get the image identifier
11            identifier = line.split('.')[0]
12            dataset.append(identifier)
13     return set(dataset)
```

Next, we need to split the set into train and test sets.

We will start by ordering the identifiers by sorting them to ensure we always split them consistently across machines and runs, then take the first 100 for train and the next 100 for test.

The *train_test_split()* function below will create this split given the loaded set of identifiers as input.

```
1  # split a dataset into train/test elements
2  def train_test_split(dataset):
3      # order keys so the split is consistent
4      ordered = sorted(dataset)
5      # return split dataset as two new sets
6      return set(ordered[:100]), set(ordered[100:200])
```

Now, we can load the photo descriptions using the pre-defined set of train or test identifiers.

Below is the function *load_clean_descriptions()* that loads the cleaned text descriptions from '*descriptions.txt*' for a given set of identifiers and returns a dictionary of identifier to text.

The model we will develop will generate a caption given a photo, and the caption will be generated one word at a time. The sequence of previously generated words will be provided as input. Therefore, we will need a "*first word*" to kick-off the generation process and a '*last word*' to signal the end of the caption. We will use the strings '*startseq*' and '*endseq*' for this purpose.

```
1  # load clean descriptions into memory
2  def load_clean_descriptions(filename, dataset):
3      # load document
4      doc = load_doc(filename)
5      descriptions = dict()
6      for line in doc.split('\n'):
7          # split line by white space
8          tokens = line.split()
9          # split id from description
10         image_id, image_desc = tokens[0], tokens[1:]
11         # skip images not in the set
12         if image_id in dataset:
13             # store
14             descriptions[image_id] = 'startseq ' + ' '.join(image_desc) + ' endseq'
15     return descriptions
```

Next, we can load the photo features for a given dataset.

Below defines a function named *load_photo_features()* that loads the entire set of photo descriptions, then returns the subset of interest for a given set of photo identifiers. This is not very efficient as the loaded dictionary of all photo features is about 700 Megabytes. Nevertheless, this will get us up and running quickly.

Note, if you have a better approach, share it in the comments below.

```
1  # load photo features
```

```
2  def load_photo_features(filename, dataset):
3      # load all features
4      all_features = load(open(filename, 'rb'))
5      # filter features
6      features = {k: all_features[k] for k in dataset}
7      return features
```

We can pause here and test everything developed so far.

The complete code example is listed below.

```
1   from pickle import load
2
3   # load doc into memory
4   def load_doc(filename):
5       # open the file as read only
6       file = open(filename, 'r')
7       # read all text
8       text = file.read()
9       # close the file
10      file.close()
11      return text
12
13  # load a pre-defined list of photo identifiers
14  def load_set(filename):
15      doc = load_doc(filename)
16      dataset = list()
17      # process line by line
18      for line in doc.split('\n'):
19          # skip empty lines
20          if len(line) < 1:
21              continue
22          # get the image identifier
23          identifier = line.split('.')[0]
24          dataset.append(identifier)
25      return set(dataset)
26
27  # split a dataset into train/test elements
28  def train_test_split(dataset):
29      # order keys so the split is consistent
30      ordered = sorted(dataset)
31      # return split dataset as two new sets
32      return set(ordered[:100]), set(ordered[100:200])
33
34  # load clean descriptions into memory
35  def load_clean_descriptions(filename, dataset):
36      # load document
37      doc = load_doc(filename)
38      descriptions = dict()
39      for line in doc.split('\n'):
40          # split line by white space
41          tokens = line.split()
42          # split id from description
43          image_id, image_desc = tokens[0], tokens[1:]
44          # skip images not in the set
45          if image_id in dataset:
46              # store
47              descriptions[image_id] = 'startseq ' + ' '.join(image_desc) + ' endseq'
48      return descriptions
49
50  # load photo features
51  def load_photo_features(filename, dataset):
52      # load all features
53      all_features = load(open(filename, 'rb'))
```

```
54      # filter features
55      features = {k: all_features[k] for k in dataset}
56      return features
57
58  # load dev set
59  filename = 'Flickr8k_text/Flickr_8k.devImages.txt'
60  dataset = load_set(filename)
61  print('Dataset: %d' % len(dataset))
62  # train-test split
63  train, test = train_test_split(dataset)
64  print('Train=%d, Test=%d' % (len(train), len(test)))
65  # descriptions
66  train_descriptions = load_clean_descriptions('descriptions.txt', train)
67  test_descriptions = load_clean_descriptions('descriptions.txt', test)
68  print('Descriptions: train=%d, test=%d' % (len(train_descriptions), len(test_descriptions)))
69  # photo features
70  train_features = load_photo_features('features.pkl', train)
71  test_features = load_photo_features('features.pkl', test)
72  print('Photos: train=%d, test=%d' % (len(train_features), len(test_features)))
```

Running this example first loads the 1,000 photo identifiers in the development dataset. A train and test set is selected and used to filter the set of clean photo descriptions and prepared image features.

We are nearly there.

```
1  Dataset: 1,000
2  Train=100, Test=100
3  Descriptions: train=100, test=100
4  Photos: train=100, test=100
```

The description text will need to be encoded to numbers before it can be presented to the model as in input or compared to the model's predictions.

The first step in encoding the data is to create a consistent mapping from words to unique integer values. Keras provides the Tokenizer class that can learn this mapping from the loaded description data.

Below defines the *create_tokenizer()* that will fit a Tokenizer given the loaded photo description text.

```
1   # fit a tokenizer given caption descriptions
2   def create_tokenizer(descriptions):
3       lines = list(descriptions.values())
4       tokenizer = Tokenizer()
5       tokenizer.fit_on_texts(lines)
6       return tokenizer
7
8   # prepare tokenizer
9   tokenizer = create_tokenizer(descriptions)
10  vocab_size = len(tokenizer.word_index) + 1
11  print('Vocabulary Size: %d' % vocab_size)
```

We can now encode the text.

Each description will be split into words. The model will be provided one word and the photo and generate the next word. Then the first two words of the description will be provided to the model as input with the image to generate the next word. This is how the model will be trained.

For example, the input sequence "*little girl running in field*" would be split into 6 input-output pairs to train the model:

```
1  X1,       X2 (text sequence),                        y (word)
2  photo     startseq,                                  little
3  photo     startseq, little,                          girl
4  photo     startseq, little, girl,                    running
5  photo     startseq, little, girl, running,           in
6  photo     startseq, little, girl, running, in,       field
7  photo     startseq, little, girl, running, in, field, endseq
```

Later when the model is used to generate descriptions, the generated words will be concatenated and recursively provided as input to generate a caption for an image.

The function below named *create_sequences()* given the tokenizer, a single clean description, the features for a photo, and the maximum description length will prepare a set of input-output pairs for training a model. Calling this function will return *X1* and *X2* for the arrays of image data and input sequence data and the *y* value for the output word.

The input sequences are integer encoded and the output word is one-hot encoded to represent the probability distribution of the expected word across the whole vocabulary of possible words.

```python
1  # create sequences of images, input sequences and output words for an image
2  def create_sequences(tokenizer, desc, image, max_length):
3      Ximages, XSeq, y = list(), list(),list()
4      vocab_size = len(tokenizer.word_index) + 1
5      # integer encode the description
6      seq = tokenizer.texts_to_sequences([desc])[0]
7      # split one sequence into multiple X,y pairs
8      for i in range(1, len(seq)):
9          # select
10         in_seq, out_seq = seq[:i], seq[i]
11         # pad input sequence
12         in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
13         # encode output sequence
14         out_seq = to_categorical([out_seq], num_classes=vocab_size)[0]
15         # store
16         Ximages.append(image)
17         XSeq.append(in_seq)
18         y.append(out_seq)
19     # Ximages, XSeq, y = array(Ximages), array(XSeq), array(y)
20     return [Ximages, XSeq, y]
```

## 2. Fit Model

We are nearly ready to fit the model.

Parts of the model have already been discussed, but let's re-iterate.

The model is based on the example laid out in the paper "Show and Tell: A Neural Image Caption Generator", 2015.

The model involves three parts:

- **Photo Feature Extractor**. This is a 16-layer VGG model pre-trained on the ImageNet dataset. We have pre-processed the photos with a the VGG model (without the top) and will use the extracted features predicted by this model as input.
- **Sequence Processor**. This is a word embedding layer for handling the text input, followed by an LSTM layer. The LSTM output is interpreted by a Dense layer one output at a time.
- **Interpreter (for lack of a better name)**. Both the feature extractor and sequence processor output a fixed-length vector that is the length of a maximum sequence. These are concatenated together and processed by an LSTM and Dense layer before a final prediction is made.

A conservative number of neurons is used in the base model. Specifically, a 128 Dense layer after the feature extractor, a 50-dimensionality word embedding followed by a 256 unit LSTM and 128 neuron Dense after the sequence processor, and finally a 500 unit LSTM followed by a 500 neuron Dense at the end of the network.

The model predicts a probability distribution across the vocabulary, therefore a softmax activation function is used and a categorical cross entropy loss function is minimized while fitting the network.

The function *define_model()* defines the baseline model, given the size of the vocabulary and the maximum length of photo descriptions. The Keras functional API is used to define the model as it provides the flexibility needed to define a model that takes two input streams and combines them.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
5       fe1 = GlobalMaxPooling2D()(inputs1)
6       fe2 = Dense(128, activation='relu')(fe1)
7       fe3 = RepeatVector(max_length)(fe2)
8       # embedding
9       inputs2 = Input(shape=(max_length,))
10      emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11      emb3 = LSTM(256, return_sequences=True)(emb2)
12      emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13      # merge inputs
14      merged = concatenate([fe3, emb4])
15      # language model (decoder)
16      lm2 = LSTM(500)(merged)
17      lm3 = Dense(500, activation='relu')(lm2)
18      outputs = Dense(vocab_size, activation='softmax')(lm3)
19      # tie it together [image, seq] [word]
20      model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21      model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22      print(model.summary())
23      plot_model(model, show_shapes=True, to_file='plot.png')
24      return model
```

To get a sense for the structure of the model, specifically the shapes of the layers, see the summary listed below.

```
1   _____
2   Layer (type)                    Output Shape              Param #     Connected to
3   ====================================================================================
4   input_1 (InputLayer)            (None, 7, 7, 512)         0
5   _____
6   input_2 (InputLayer)            (None, 25)                0
```

```
 7   _____
 8   global_max_pooling2d_1 (GlobalMa (None, 512)              0           input_1[0][0]
 9   _____
10   embedding_1 (Embedding)          (None, 25, 50)           18300       input_2[0][0]
11   _____
12   dense_1 (Dense)                  (None, 128)              65664       global_max_pooling2d_1[0][0]
13   _____
14   lstm_1 (LSTM)                    (None, 25, 256)          314368      embedding_1[0][0]
15   _____
16   repeat_vector_1 (RepeatVector)   (None, 25, 128)          0           dense_1[0][0]
17   _____
18   time_distributed_1 (TimeDistribu (None, 25, 128)          32896       lstm_1[0][0]
19   _____
20   concatenate_1 (Concatenate)      (None, 25, 256)          0           repeat_vector_1[0][0]
21                                                                        time_distributed_1[0][0]
22   _____
23   lstm_2 (LSTM)                    (None, 500)              1514000     concatenate_1[0][0]
24   _____
25   dense_3 (Dense)                  (None, 500)              250500      lstm_2[0][0]
26   _____
27   dense_4 (Dense)                  (None, 366)              183366      dense_3[0][0]
28   ===================================================================
29   Total params: 2,379,094
30   Trainable params: 2,379,094
31   Non-trainable params: 0
32   _____
```

We also create a plot to visualize the structure of the network that better helps understand the two streams of input.

Plot of the Baseline Captioning Deep Learning Model

We will train the model using a data generator. This is strictly not required given that the captions and extracted photo features can probably fit into memory as a single dataset. Nevertheless, it is good practice for when you come to train the final model on the entire dataset.

A generator will yield a result when called. In Keras, it will yield a single batch of input-output samples that are used to estimate the error gradient and update the model weights.

The function *data_generator()* defines the data generator, given a dictionary of loaded photo descriptions, photo features, the tokenizer for integer encoding sequences, and the maximum sequence length in the dataset.

The generator loops forever and keeps yielding batches of input-output pairs when asked. We also have a *n_step* parameter that allows us to tune how many images worth of input-output pairs to generate for each batch. The average sequence has 10 words, that is 10 input-output pairs, and a good batch size might be 30 samples, which is about 2-to-3 images worth.

```
1   # data generator, intended to be used in a call to model.fit_generator()
2   def data_generator(descriptions, features, tokenizer, max_length, n_step):
3       # loop until we finish training
4       while 1:
5           # loop over photo identifiers in the dataset
6           keys = list(descriptions.keys())
7           for i in range(0, len(keys), n_step):
8               Ximages, XSeq, y = list(), list(),list()
9               for j in range(i, min(len(keys), i+n_step)):
10                  image_id = keys[j]
11                  # retrieve photo feature input
12                  image = features[image_id][0]
13                  # retrieve text input
14                  desc = descriptions[image_id]
15                  # generate input-output pairs
16                  in_img, in_seq, out_word = create_sequences(tokenizer, desc, image, max_length)
17                  for k in range(len(in_img)):
18                      Ximages.append(in_img[k])
19                      XSeq.append(in_seq[k])
20                      y.append(out_word[k])
21              # yield this batch of samples to the model
22              yield [[array(Ximages), array(XSeq)], array(y)]
```

The model can be fit by calling *fit_generator()* and passing it to the data generator, along with all of the parameters needed. When fitting the model, we can also specify the number of batches to run per epoch and the number of epochs.

```
1  model.fit_generator(data_generator(train_descriptions, train_features, tokenizer, max_length, n_
```

For these experiments, we will use 2 images per batch, 50 batches (or 100 images) per epoch, and 50 training epochs. You can experiment with different configurations in your own experiments.

## 3. Evaluate Model

Now that we know how to prepare the data and define a model, we must define a test harness to evaluate a given model.

We will evaluate a model by training it on the dataset, generating descriptions for all photos in the training dataset, evaluating those predictions with a cost function, and then repeating this evaluation process multiple times.

The outcome will be a distribution of skill scores for the model that we can summarize by calculating the mean and standard deviation. This is the preferred way to evaluate deep learning models. See this post:

- How to Evaluate the Skill of Deep Learning Models

First, we need to be able to generate a description for a photo using a trained model.

This involves passing in the start description token '*startseq*', generating one word, then calling the model recursively with generated words as input until the end of sequence token is reached '*endseq*' or the maximum description length is reached.

The function below named *generate_desc()* implements this behavior and generates a textual description given a trained model, and a given prepared photo as input. It calls the function *word_for_id()* in order to map an integer prediction back to a word.

```
1  # map an integer to a word
2  def word_for_id(integer, tokenizer):
3      for word, index in tokenizer.word_index.items():
4          if index == integer:
5              return word
6      return None
7
8  # generate a description for an image
9  def generate_desc(model, tokenizer, photo, max_length):
10     # seed the generation process
11     in_text = 'startseq'
12     # iterate over the whole length of the sequence
13     for i in range(max_length):
14         # integer encode input sequence
15         sequence = tokenizer.texts_to_sequences([in_text])[0]
16         # pad input
17         sequence = pad_sequences([sequence], maxlen=max_length)
18         # predict next word
19         yhat = model.predict([photo,sequence], verbose=0)
20         # convert probability to integer
21         yhat = argmax(yhat)
22         # map integer to word
23         word = word_for_id(yhat, tokenizer)
24         # stop if we cannot map the word
25         if word is None:
26             break
27         # append as input for generating the next word
28         in_text += ' ' + word
29         # stop if we predict the end of the sequence
30         if word == 'endseq':
31             break
32     return in_text
```

We will generate predictions for all photos in the training dataset and in the test dataset.

The function below named *evaluate_model()* will evaluate a trained model against a given dataset of photo descriptions and photo features. The actual and predicted descriptions are collected and evaluated

collectively using the corpus BLEU score that summarizes how close the generated text is to the expected text.

```
1  # evaluate the skill of the model
2  def evaluate_model(model, descriptions, photos, tokenizer, max_length):
3      actual, predicted = list(), list()
4      # step over the whole set
5      for key, desc in descriptions.items():
6          # generate description
7          yhat = generate_desc(model, tokenizer, photos[key], max_length)
8          # store actual and predicted
9          actual.append([desc.split()])
10         predicted.append(yhat.split())
11     # calculate BLEU score
12     bleu = corpus_bleu(actual, predicted)
13     return bleu
```

BLEU scores are used in text translation for evaluating translated text against one or more reference translations. We do in fact have access to multiple reference descriptions for each image that we could compare to, but for simplicity, we will use the first description for each photo in the dataset (e.g. the cleaned version).

You can learn more about the BLEU score here:

- BLEU (bilingual evaluation understudy) on Wikipedia

The NLTK Python library implements the BLEU score calculation in the *corpus_bleu()* function. A higher score close to 1.0 is better, a score closer to zero is worse.

Finally, all we need to do is define, fit, and evaluate the model multiple times in a loop then report the final average score.

Ideally, we would repeat the experiment 30 times or more, but this will take too long for our small test harness. Instead, will evaluate the model 3 times. It will be faster, but the mean score will have higher variance.

Below defines the model evaluation loop. At the end of the run, the distribution of BLEU scores for the train and test sets are saved to a file.

```
1  # run experiment
2  train_results, test_results = list(), list()
3  for i in range(n_repeats):
4      # define the model
5      model = define_model(vocab_size, max_length)
6      # fit model
7      model.fit_generator(data_generator(train_descriptions, train_features, tokenizer, max_lengt
8      # evaluate model on training data
9      train_score = evaluate_model(model, train_descriptions, train_features, tokenizer, max_leng
10     test_score = evaluate_model(model, test_descriptions, test_features, tokenizer, max_length)
11     # store
12     train_results.append(train_score)
13     test_results.append(test_score)
14     print('>%d: train=%f test=%f' % ((i+1), train_score, test_score))
15 # save results to file
16 df = DataFrame()
17 df['train'] = train_results
```

```
18  df['test'] = test_results
19  print(df.describe())
20  df.to_csv(model_name+'.csv', index=False)
```

We parameterize the run as follows, allowing us to name each run and save the result to separate files.

```
1  # define experiment
2  model_name = 'baseline1'
3  verbose = 2
4  n_epochs = 50
5  n_photos_per_update = 2
6  n_batches_per_epoch = int(len(train) / n_photos_per_update)
7  n_repeats = 3
```

## 4. Complete Example

The complete example is listed below.

```
 1  from os import listdir
 2  from numpy import array
 3  from numpy import argmax
 4  from pandas import DataFrame
 5  from nltk.translate.bleu_score import corpus_bleu
 6  from pickle import load
 7
 8  from keras.preprocessing.text import Tokenizer
 9  from keras.preprocessing.sequence import pad_sequences
10  from keras.utils import to_categorical
11  from keras.preprocessing.image import load_img
12  from keras.preprocessing.image import img_to_array
13  from keras.applications.vgg16 import preprocess_input
14  from keras.applications.vgg16 import VGG16
15  from keras.utils import plot_model
16  from keras.models import Model
17  from keras.layers import Input
18  from keras.layers import Dense
19  from keras.layers import Flatten
20  from keras.layers import LSTM
21  from keras.layers import RepeatVector
22  from keras.layers import TimeDistributed
23  from keras.layers import Embedding
24  from keras.layers.merge import concatenate
25  from keras.layers.pooling import GlobalMaxPooling2D
26
27  # load doc into memory
28  def load_doc(filename):
29      # open the file as read only
30      file = open(filename, 'r')
31      # read all text
32      text = file.read()
33      # close the file
34      file.close()
35      return text
36
37  # load a pre-defined list of photo identifiers
38  def load_set(filename):
39      doc = load_doc(filename)
40      dataset = list()
41      # process line by line
42      for line in doc.split('\n'):
43          # skip empty lines
44          if len(line) < 1:
45              continue
```

```
 46              # get the image identifier
 47              identifier = line.split('.')[0]
 48              dataset.append(identifier)
 49      return set(dataset)
 50
 51  # split a dataset into train/test elements
 52  def train_test_split(dataset):
 53      # order keys so the split is consistent
 54      ordered = sorted(dataset)
 55      # return split dataset as two new sets
 56      return set(ordered[:100]), set(ordered[100:200])
 57
 58  # load clean descriptions into memory
 59  def load_clean_descriptions(filename, dataset):
 60      # load document
 61      doc = load_doc(filename)
 62      descriptions = dict()
 63      for line in doc.split('\n'):
 64          # split line by white space
 65          tokens = line.split()
 66          # split id from description
 67          image_id, image_desc = tokens[0], tokens[1:]
 68          # skip images not in the set
 69          if image_id in dataset:
 70              # store
 71              descriptions[image_id] = 'startseq ' + ' '.join(image_desc) + ' endseq'
 72      return descriptions
 73
 74  # load photo features
 75  def load_photo_features(filename, dataset):
 76      # load all features
 77      all_features = load(open(filename, 'rb'))
 78      # filter features
 79      features = {k: all_features[k] for k in dataset}
 80      return features
 81
 82  # fit a tokenizer given caption descriptions
 83  def create_tokenizer(descriptions):
 84      lines = list(descriptions.values())
 85      tokenizer = Tokenizer()
 86      tokenizer.fit_on_texts(lines)
 87      return tokenizer
 88
 89  # create sequences of images, input sequences and output words for an image
 90  def create_sequences(tokenizer, desc, image, max_length):
 91      Ximages, XSeq, y = list(), list(),list()
 92      vocab_size = len(tokenizer.word_index) + 1
 93      # integer encode the description
 94      seq = tokenizer.texts_to_sequences([desc])[0]
 95      # split one sequence into multiple X,y pairs
 96      for i in range(1, len(seq)):
 97          # select
 98          in_seq, out_seq = seq[:i], seq[i]
 99          # pad input sequence
100          in_seq = pad_sequences([in_seq], maxlen=max_length)[0]
101          # encode output sequence
102          out_seq = to_categorical([out_seq], num_classes=vocab_size)[0]
103          # store
104          Ximages.append(image)
105          XSeq.append(in_seq)
106          y.append(out_seq)
107      # Ximages, XSeq, y = array(Ximages), array(XSeq), array(y)
108      return [Ximages, XSeq, y]
109
110  # define the captioning model
```

```python
111 def define_model(vocab_size, max_length):
112     # feature extractor (encoder)
113     inputs1 = Input(shape=(7, 7, 512))
114     fe1 = GlobalMaxPooling2D()(inputs1)
115     fe2 = Dense(128, activation='relu')(fe1)
116     fe3 = RepeatVector(max_length)(fe2)
117     # embedding
118     inputs2 = Input(shape=(max_length,))
119     emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
120     emb3 = LSTM(256, return_sequences=True)(emb2)
121     emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
122     # merge inputs
123     merged = concatenate([fe3, emb4])
124     # language model (decoder)
125     lm2 = LSTM(500)(merged)
126     lm3 = Dense(500, activation='relu')(lm2)
127     outputs = Dense(vocab_size, activation='softmax')(lm3)
128     # tie it together [image, seq] [word]
129     model = Model(inputs=[inputs1, inputs2], outputs=outputs)
130     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
131     print(model.summary())
132     plot_model(model, show_shapes=True, to_file='plot.png')
133     return model
134
135 # data generator, intended to be used in a call to model.fit_generator()
136 def data_generator(descriptions, features, tokenizer, max_length, n_step):
137     # loop until we finish training
138     while 1:
139         # loop over photo identifiers in the dataset
140         keys = list(descriptions.keys())
141         for i in range(0, len(keys), n_step):
142             Ximages, XSeq, y = list(), list(),list()
143             for j in range(i, min(len(keys), i+n_step)):
144                 image_id = keys[j]
145                 # retrieve photo feature input
146                 image = features[image_id][0]
147                 # retrieve text input
148                 desc = descriptions[image_id]
149                 # generate input-output pairs
150                 in_img, in_seq, out_word = create_sequences(tokenizer, desc, image, max_length
151                 for k in range(len(in_img)):
152                     Ximages.append(in_img[k])
153                     XSeq.append(in_seq[k])
154                     y.append(out_word[k])
155             # yield this batch of samples to the model
156             yield [[array(Ximages), array(XSeq)], array(y)]
157
158 # map an integer to a word
159 def word_for_id(integer, tokenizer):
160     for word, index in tokenizer.word_index.items():
161         if index == integer:
162             return word
163     return None
164
165 # generate a description for an image
166 def generate_desc(model, tokenizer, photo, max_length):
167     # seed the generation process
168     in_text = 'startseq'
169     # iterate over the whole length of the sequence
170     for i in range(max_length):
171         # integer encode input sequence
172         sequence = tokenizer.texts_to_sequences([in_text])[0]
173         # pad input
174         sequence = pad_sequences([sequence], maxlen=max_length)
175         # predict next word
```

```
176          yhat = model.predict([photo,sequence], verbose=0)
177          # convert probability to integer
178          yhat = argmax(yhat)
179          # map integer to word
180          word = word_for_id(yhat, tokenizer)
181          # stop if we cannot map the word
182          if word is None:
183              break
184          # append as input for generating the next word
185          in_text += ' ' + word
186          # stop if we predict the end of the sequence
187          if word == 'endseq':
188              break
189      return in_text
190
191  # evaluate the skill of the model
192  def evaluate_model(model, descriptions, photos, tokenizer, max_length):
193      actual, predicted = list(), list()
194      # step over the whole set
195      for key, desc in descriptions.items():
196          # generate description
197          yhat = generate_desc(model, tokenizer, photos[key], max_length)
198          # store actual and predicted
199          actual.append([desc.split()])
200          predicted.append(yhat.split())
201      # calculate BLEU score
202      bleu = corpus_bleu(actual, predicted)
203      return bleu
204
205  # load dev set
206  filename = 'Flickr8k_text/Flickr_8k.devImages.txt'
207  dataset = load_set(filename)
208  print('Dataset: %d' % len(dataset))
209  # train-test split
210  train, test = train_test_split(dataset)
211  # descriptions
212  train_descriptions = load_clean_descriptions('descriptions.txt', train)
213  test_descriptions = load_clean_descriptions('descriptions.txt', test)
214  print('Descriptions: train=%d, test=%d' % (len(train_descriptions), len(test_descriptions)))
215  # photo features
216  train_features = load_photo_features('features.pkl', train)
217  test_features = load_photo_features('features.pkl', test)
218  print('Photos: train=%d, test=%d' % (len(train_features), len(test_features)))
219  # prepare tokenizer
220  tokenizer = create_tokenizer(train_descriptions)
221  vocab_size = len(tokenizer.word_index) + 1
222  print('Vocabulary Size: %d' % vocab_size)
223  # determine the maximum sequence length
224  max_length = max(len(s.split()) for s in list(train_descriptions.values()))
225  print('Description Length: %d' % max_length)
226
227  # define experiment
228  model_name = 'baseline1'
229  verbose = 2
230  n_epochs = 50
231  n_photos_per_update = 2
232  n_batches_per_epoch = int(len(train) / n_photos_per_update)
233  n_repeats = 3
234
235  # run experiment
236  train_results, test_results = list(), list()
237  for i in range(n_repeats):
238      # define the model
239      model = define_model(vocab_size, max_length)
240      # fit model
```

```
241     model.fit_generator(data_generator(train_descriptions, train_features, tokenizer, max_leng
242     # evaluate model on training data
243     train_score = evaluate_model(model, train_descriptions, train_features, tokenizer, max_ler
244     test_score = evaluate_model(model, test_descriptions, test_features, tokenizer, max_length
245     # store
246     train_results.append(train_score)
247     test_results.append(test_score)
248     print('>%d: train=%f test=%f' % ((i+1), train_score, test_score))
249 # save results to file
250 df = DataFrame()
251 df['train'] = train_results
252 df['test'] = test_results
253 print(df.describe())
254 df.to_csv(model_name+'.csv', index=False)
```

Running the example first prints summary statistics for the loaded training data.

```
1 Dataset: 1,000
2 Descriptions: train=100, test=100
3 Photos: train=100, test=100
4 Vocabulary Size: 366
5 Description Length: 25
```

The example should take about 20 minutes on GPU hardware, a little longer on CPU hardware.

At the end of the run, a mean BLEU of 0.06 is reported on the training set and 0.04 on the test set. Results are stored in *baseline1.csv*.

```
1           train      test
2 count   3.000000   3.000000
3 mean    0.060617   0.040978
4 std     0.023498   0.025105
5 min     0.042882   0.012101
6 25%     0.047291   0.032658
7 50%     0.051701   0.053215
8 75%     0.069484   0.055416
9 max     0.087268   0.057617
```

This provides a baseline model for comparison to alternate configurations.

## "A" versus "A" Test

Before we start testing variations of the model, it is important to get an idea of whether or not the test harness is stable.

That is, whether the summarizing skill of the model over 5 runs is sufficient to control for the stochastic nature of the model.

We can get an idea of this by running the experiment again in what is called an A vs A test in A/B testing land. We would expect to get an equivalent result if we ran the same experiment again; if we don't, perhaps additional repeats would be required to control for the stochastic nature of the method and on the dataset.

Below are the results from a second run of the algorithm.

```
1           train      test
2 count   3.000000   3.000000
3 mean    0.036902   0.043003
```

```
4 std      0.020281  0.017295
5 min      0.018522  0.026055
6 25%      0.026023  0.034192
7 50%      0.033525  0.042329
8 75%      0.046093  0.051477
9 max      0.058660  0.060624
```

We can see that the run gets a very similar mean and standard deviation BLEU scores. Specifically, a mean BLEU of 0.03 vs 0.06 on train and 0.04 to 0.04 for test.

The harness is a little noisy, but stable enough for comparison.

Is the model any good?

## Generate Photo Captions

We expect the model is under-trained and maybe even under provisioned, but can it generate any kind of readable text at all?

It is important that the baseline model have some modicum of capability so that we can relate the BLEU scores of the baseline to an idea of what kind of quality of descriptions are being generated.

Let's train a single model and generate a few descriptions from the train and test sets as a sanity check.

Change the number of repeats to 1 and the name of the run to '*baseline_generate*'.

```
1  model_name = 'baseline_generate'
2  n_repeats = 1
```

Then update the *evaluate_model()* function to only evaluate the first 5 photos in the dataset and print the descriptions, as follows.

```
1   # evaluate the skill of the model
2   def evaluate_model(model, descriptions, photos, tokenizer, max_length):
3       actual, predicted = list(), list()
4       # step over the whole set
5       for key, desc in descriptions.items():
6           # generate description
7           yhat = generate_desc(model, tokenizer, photos[key], max_length)
8           # store actual and predicted
9           actual.append([desc.split()])
10          predicted.append(yhat.split())
11          print('Actual:    %s' % desc)
12          print('Predicted: %s' % yhat)
13          if len(actual) >= 5:
14              break
15      # calculate BLEU score
16      bleu = corpus_bleu(actual, predicted)
17      return bleu
```

Re-run the example.

You should see results for the train set like the following (the specific results will vary given the stochastic nature of the algorithm):

```
1  Actual:     startseq boy bites hard into treat while he sits outside endseq
2  Predicted: startseq boy boy while while he while outside endseq
3
4  Actual:     startseq man in field backed by american flags endseq
5  Predicted: startseq man in in standing city endseq
6
7  Actual:     startseq two girls are walking down dirt road in park endseq
8  Predicted: startseq man walking down down road in endseq
9
10 Actual:     startseq girl laying on the tree with boy kneeling before her endseq
11 Predicted: startseq boy while in up up up water endseq
12
13 Actual:     startseq boy in striped shirt is jumping in front of water fountain endseq
14 Predicted: startseq man is is shirt is on on on on bike endseq
```

You should see results on the test dataset as follows:

```
1  Actual:     startseq three people are looking into photographic equipment endseq
2  Predicted: startseq boy racer on on on on bike endseq
3
4  Actual:     startseq boy is leaning on chair whilst another boy pulls him around with rope endse
5  Predicted: startseq girl in playing on on on sword endseq
6
7  Actual:     startseq black and brown dog jumping in midair near field endseq
8  Predicted: startseq dog dog running running running and dog in grass endseq
9
10 Actual:     startseq dog places his head on man face endseq
11 Predicted: startseq brown dog dog to to to to to to to ball endseq
12
13 Actual:     startseq man in green hat is someplace up high endseq
14 Predicted: startseq man in up up his waves endseq
```

We can see that the descriptions are not perfect, some are a little rough, but generally the model is generating somewhat readable text. A good starting point for improvement.

Next, let's look at some experiments to vary the size or capacity of different sub-models.

# Network Size Parameters

In this section, we will see how gross variations to the network structure impact model skill.

We will look at the following aspects of the model size:

1. Size of the fixed-vector output from the 'encoders'.
2. Size of the sequence encoder model.
3. Size of the language model.

Let's dive in.

## Size of Fixed-Length Vector

In the baseline model, the photo feature extractor and the text sequence encoder both output a 128 element vector. These vectors are then concatenated to be processed by the language model.

The 128 element vector from each sub-model contains everything known about the input sequence and photo. We can vary the size of this vector to see if it impacts model skill

First, we can decrease the size by half from 128 elements to 64 elements.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
5       fe1 = GlobalMaxPooling2D()(inputs1)
6       fe2 = Dense(64, activation='relu')(fe1)
7       fe3 = RepeatVector(max_length)(fe2)
8       # embedding
9       inputs2 = Input(shape=(max_length,))
10      emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11      emb3 = LSTM(256, return_sequences=True)(emb2)
12      emb4 = TimeDistributed(Dense(64, activation='relu'))(emb3)
13      # merge inputs
14      merged = concatenate([fe3, emb4])
15      # language model (decoder)
16      lm2 = LSTM(500)(merged)
17      lm3 = Dense(500, activation='relu')(lm2)
18      outputs = Dense(vocab_size, activation='softmax')(lm3)
19      # tie it together [image, seq] [word]
20      model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21      model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22      return model
```

We will name this model '*size_sm_fixed_vec*'.

```
1  model_name = 'size_sm_fixed_vec'
```

Running this experiment produces the following BLEU scores, perhaps a small gain over baseline on the test set.

```
1              train       test
2  count    3.000000    3.000000
3  mean     0.204421    0.063148
4  std      0.026992    0.003264
5  min      0.174769    0.059391
6  25%      0.192849    0.062074
7  50%      0.210929    0.064757
8  75%      0.219246    0.065026
9  max      0.227564    0.065295
```

We can also double the size of the fixed-length vector from 128 to 256 units.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
5       fe1 = GlobalMaxPooling2D()(inputs1)
6       fe2 = Dense(256, activation='relu')(fe1)
7       fe3 = RepeatVector(max_length)(fe2)
8       # embedding
9       inputs2 = Input(shape=(max_length,))
10      emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11      emb3 = LSTM(256, return_sequences=True)(emb2)
12      emb4 = TimeDistributed(Dense(256, activation='relu'))(emb3)
13      # merge inputs
14      merged = concatenate([fe3, emb4])
```

```
15        # language model (decoder)
16        lm2 = LSTM(500)(merged)
17        lm3 = Dense(500, activation='relu')(lm2)
18        outputs = Dense(vocab_size, activation='softmax')(lm3)
19        # tie it together [image, seq] [word]
20        model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21        model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22        return model
```

We will name this configuration '*size_lg_fixed_vec*'.

```
1  model_name = 'size_lg_fixed_vec'
```

Running this experiment shows BLEU scores suggesting that the model is not better off.

It is possible that with more data and/or longer training, we may see a different story.

```
1               train       test
2  count   3.000000   3.000000
3  mean    0.023517   0.027813
4  std     0.009951   0.010525
5  min     0.012037   0.021737
6  25%     0.020435   0.021737
7  50%     0.028833   0.021737
8  75%     0.029257   0.030852
9  max     0.029682   0.039966
```

## Sequence Encoder Size

We can call the sub-model that interprets the input sequence of words generated so far as the sequence encoder.

First, we can try to see if decreasing the representational capacity of the sequence encoder impacts model skill. We can reduce the number of memory units in the LSTM layer from 256 to 128.

```
1  # define the captioning model
2  def define_model(vocab_size, max_length):
3        # feature extractor (encoder)
4        inputs1 = Input(shape=(7, 7, 512))
5        fe1 = GlobalMaxPooling2D()(inputs1)
6        fe2 = Dense(128, activation='relu')(fe1)
7        fe3 = RepeatVector(max_length)(fe2)
8        # embedding
9        inputs2 = Input(shape=(max_length,))
10       emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11       emb3 = LSTM(128, return_sequences=True)(emb2)
12       emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13       # merge inputs
14       merged = concatenate([fe3, emb4])
15       # language model (decoder)
16       lm2 = LSTM(500)(merged)
17       lm3 = Dense(500, activation='relu')(lm2)
18       outputs = Dense(vocab_size, activation='softmax')(lm3)
19       # tie it together [image, seq] [word]
20       model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21       model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22       return model
23
24  model_name = 'size_sm_seq_model'
```

Running this example, we can see perhaps a small bump on both train and test over baseline. This might be an artifact of the small training set size.

```
1            train       test
2  count   3.000000   3.000000
3  mean    0.074944   0.053917
4  std     0.014263   0.013264
5  min     0.066292   0.039142
6  25%     0.066713   0.048476
7  50%     0.067134   0.057810
8  75%     0.079270   0.061304
9  max     0.091406   0.064799
```

Going the other way, we can double the number of LSTM layers from one to two and see if that makes a dramatic difference.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
5       fe1 = GlobalMaxPooling2D()(inputs1)
6       fe2 = Dense(128, activation='relu')(fe1)
7       fe3 = RepeatVector(max_length)(fe2)
8       # embedding
9       inputs2 = Input(shape=(max_length,))
10      emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11      emb3 = LSTM(256, return_sequences=True)(emb2)
12      emb4 = LSTM(256, return_sequences=True)(emb3)
13      emb5 = TimeDistributed(Dense(128, activation='relu'))(emb4)
14      # merge inputs
15      merged = concatenate([fe3, emb5])
16      # language model (decoder)
17      lm2 = LSTM(500)(merged)
18      lm3 = Dense(500, activation='relu')(lm2)
19      outputs = Dense(vocab_size, activation='softmax')(lm3)
20      # tie it together [image, seq] [word]
21      model = Model(inputs=[inputs1, inputs2], outputs=outputs)
22      model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
23      return model
24
25  model_name = 'size_lg_seq_model'
```

Running this experiment shows a decent bump in BLEU on both train and test sets.

```
1            train       test
2  count   3.000000   3.000000
3  mean    0.094937   0.096970
4  std     0.022394   0.079270
5  min     0.069151   0.046722
6  25%     0.087656   0.051279
7  50%     0.106161   0.055836
8  75%     0.107830   0.122094
9  max     0.109499   0.188351
```

We can also try to increase the representational capacity of the word embedding by doubling it from 50-dimensions to 100-dimensions.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
```

```
 5        fe1 = GlobalMaxPooling2D()(inputs1)
 6        fe2 = Dense(128, activation='relu')(fe1)
 7        fe3 = RepeatVector(max_length)(fe2)
 8        # embedding
 9        inputs2 = Input(shape=(max_length,))
10        emb2 = Embedding(vocab_size, 100, mask_zero=True)(inputs2)
11        emb3 = LSTM(256, return_sequences=True)(emb2)
12        emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13        # merge inputs
14        merged = concatenate([fe3, emb4])
15        # language model (decoder)
16        lm2 = LSTM(500)(merged)
17        lm3 = Dense(500, activation='relu')(lm2)
18        outputs = Dense(vocab_size, activation='softmax')(lm3)
19        # tie it together [image, seq] [word]
20        model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21        model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22        return model
23
24 model_name = 'size_em_seq_model'
```

We see a large movement on the training dataset, but perhaps little movement on the test dataset.

```
1  count   3.000000   3.000000
2  mean    0.112743   0.050935
3  std     0.017136   0.006860
4  min     0.096121   0.043741
5  25%     0.103940   0.047701
6  50%     0.111759   0.051661
7  75%     0.121055   0.054533
8  max     0.130350   0.057404
```

## Size of Language Model

We can refer to the model that learns from the concatenated sequence and photo feature input as the language model. It is responsible for generating words.

First, we can look at the impact on model skill by cutting the LSTM and dense layers from 500 to 256 neurons.

```
 1  # define the captioning model
 2  def define_model(vocab_size, max_length):
 3        # feature extractor (encoder)
 4        inputs1 = Input(shape=(7, 7, 512))
 5        fe1 = GlobalMaxPooling2D()(inputs1)
 6        fe2 = Dense(128, activation='relu')(fe1)
 7        fe3 = RepeatVector(max_length)(fe2)
 8        # embedding
 9        inputs2 = Input(shape=(max_length,))
10        emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11        emb3 = LSTM(256, return_sequences=True)(emb2)
12        emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13        # merge inputs
14        merged = concatenate([fe3, emb4])
15        # language model (decoder)
16        lm2 = LSTM(256)(merged)
17        lm3 = Dense(256, activation='relu')(lm2)
18        outputs = Dense(vocab_size, activation='softmax')(lm3)
19        # tie it together [image, seq] [word]
20        model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21        model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
```

```
22      return model
23
24  model_name = 'size_sm_lang_model'
```

We can see that this has a small positive effect on BLEU for both training and test datasets, again, likely related to the small size of the datasets.

```
1            train       test
2  count  3.000000   3.000000
3  mean   0.063632   0.056059
4  std    0.018521   0.009064
5  min    0.045127   0.048916
6  25%    0.054363   0.050961
7  50%    0.063599   0.053005
8  75%    0.072884   0.059630
9  max    0.082169   0.066256
```

We can also look at the impact of doubling the capacity of the language model by adding a second LSTM layer of the same size.

```
1  # define the captioning model
2  def define_model(vocab_size, max_length):
3      # feature extractor (encoder)
4      inputs1 = Input(shape=(7, 7, 512))
5      fe1 = GlobalMaxPooling2D()(inputs1)
6      fe2 = Dense(128, activation='relu')(fe1)
7      fe3 = RepeatVector(max_length)(fe2)
8      # embedding
9      inputs2 = Input(shape=(max_length,))
10     emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11     emb3 = LSTM(256, return_sequences=True)(emb2)
12     emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13     # merge inputs
14     merged = concatenate([fe3, emb4])
15     # language model (decoder)
16     lm2 = LSTM(500, return_sequences=True)(merged)
17     lm3 = LSTM(500)(lm2)
18     lm4 = Dense(500, activation='relu')(lm3)
19     outputs = Dense(vocab_size, activation='softmax')(lm4)
20     # tie it together [image, seq] [word]
21     model = Model(inputs=[inputs1, inputs2], outputs=outputs)
22     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
23     return model
24
25  model_name = 'size_lg_lang_model'
```

Again, we see minor movements in BLEU, perhaps an artifact of noise and dataset size. The improvement on the test

The improvement on the test dataset may be a good sign. This might be a change worth exploring.

```
1            train       test
2  count  3.000000   3.000000
3  mean   0.043838   0.067658
4  std    0.037580   0.045813
5  min    0.017990   0.015757
6  25%    0.022284   0.050252
7  50%    0.026578   0.084748
8  75%    0.056763   0.093608
9  max    0.086948   0.102469
```

Tuning model size on a much smaller dataset is challenging.

# Configuring the Feature Extraction Model

The use of the pre-trained VGG16 model provides some additional points of configuration.

The baseline model removed the top from the VGG model, including a global max pooling layer, which then feeds into an encoding of the features to a 128 element vector.

In this section, we will look at the following modifications to the baseline model:

1. Using a global average pooling layer after the VGG model.
2. Not using any global pooling.

## Global Average Pooling

We can replace the GlobalMaxPooling2D layer with a GlobalAveragePooling2D to achieve average pooling.

Global average pooling was developed to reduce overfitting for image classification problems, but may offer some benefit in interpreting the features extracted from the image.

For more on global average pooling, see the paper:

- Network In Network, 2013.

The updated *define_model()* function and experiment name are listed below.

```
# define the captioning model
def define_model(vocab_size, max_length):
    # feature extractor (encoder)
    inputs1 = Input(shape=(7, 7, 512))
    fe1 = GlobalAveragePooling2D()(inputs1)
    fe2 = Dense(128, activation='relu')(fe1)
    fe3 = RepeatVector(max_length)(fe2)
    # embedding
    inputs2 = Input(shape=(max_length,))
    emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
    emb3 = LSTM(256, return_sequences=True)(emb2)
    emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
    # merge inputs
    merged = concatenate([fe3, emb4])
    # language model (decoder)
    lm2 = LSTM(500)(merged)
    lm3 = Dense(500, activation='relu')(lm2)
    outputs = Dense(vocab_size, activation='softmax')(lm3)
    # tie it together [image, seq] [word]
    model = Model(inputs=[inputs1, inputs2], outputs=outputs)
    model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
    return model

model_name = 'fe_avg_pool'
```

The results suggest a dramatic improvement on the training dataset, which may be a sign of overfitting. We also see a small lift on test skill. This might be a change worth exploring.

We also see a small lift on test skill. This might be a change worth exploring.

```
1            train      test
2  count  3.000000  3.000000
3  mean   0.834627  0.060847
4  std    0.083259  0.040463
5  min    0.745074  0.017705
6  25%    0.797096  0.042294
7  50%    0.849118  0.066884
8  75%    0.879404  0.082418
9  max    0.909690  0.097952
```

## No Pooling

We can remove the GlobalMaxPooling2D and flatten the 3D photo feature and feed it directly into a Dense layer.

I would not expect this to be a good model design, but it is worth testing this assumption.

```
1  # define the captioning model
2  def define_model(vocab_size, max_length):
3      # feature extractor (encoder)
4      inputs1 = Input(shape=(7, 7, 512))
5      fe1 = Flatten()(inputs1)
6      fe2 = Dense(128, activation='relu')(fe1)
7      fe3 = RepeatVector(max_length)(fe2)
8      # embedding
9      inputs2 = Input(shape=(max_length,))
10     emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
11     emb3 = LSTM(256, return_sequences=True)(emb2)
12     emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13     # merge inputs
14     merged = concatenate([fe3, emb4])
15     # language model (decoder)
16     lm2 = LSTM(500)(merged)
17     lm3 = Dense(500, activation='relu')(lm2)
18     outputs = Dense(vocab_size, activation='softmax')(lm3)
19     # tie it together [image, seq] [word]
20     model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22     return model
23
24 model_name = 'fe_flat'
```

Surprisingly, we see a small lift on training data and a large lift on test data. This is surprising (to me) and may be worth further investigation.

```
1            train      test
2  count  3.000000  3.000000
3  mean   0.055988  0.135231
4  std    0.017566  0.079714
5  min    0.038605  0.044177
6  25%    0.047116  0.106633
7  50%    0.055627  0.169089
8  75%    0.064679  0.180758
9  max    0.073731  0.192428
```

We can try repeating this experiment and provide more capacity for interpreting the extracted photo features. A new Dense layer with 500 neurons is added after the Flatten layer.

```
1   # define the captioning model
2   def define_model(vocab_size, max_length):
3       # feature extractor (encoder)
4       inputs1 = Input(shape=(7, 7, 512))
5       fe1 = Flatten()(inputs1)
6       fe2 = Dense(500, activation='relu')(fe1)
7       fe3 = Dense(128, activation='relu')(fe2)
8       fe4 = RepeatVector(max_length)(fe3)
9       # embedding
10      inputs2 = Input(shape=(max_length,))
11      emb2 = Embedding(vocab_size, 50, mask_zero=True)(inputs2)
12      emb3 = LSTM(256, return_sequences=True)(emb2)
13      emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
14      # merge inputs
15      merged = concatenate([fe4, emb4])
16      # language model (decoder)
17      lm2 = LSTM(500)(merged)
18      lm3 = Dense(500, activation='relu')(lm2)
19      outputs = Dense(vocab_size, activation='softmax')(lm3)
20      # tie it together [image, seq] [word]
21      model = Model(inputs=[inputs1, inputs2], outputs=outputs)
22      model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
23      return model
24
25  model_name = 'fe_flat2'
```

This results in a less impressive change and perhaps worse BLEU results on the test dataset.

```
1              train        test
2   count   3.000000    3.000000
3   mean    0.060126    0.029487
4   std     0.030300    0.013205
5   min     0.031235    0.020850
6   25%     0.044359    0.021887
7   50%     0.057483    0.022923
8   75%     0.074572    0.033805
9   max     0.091661    0.044688
```

# Word Embedding Models

A key part of the model is the sequence learning model that must interpret the sequence of words generated so far for a photo.

At the input to this sub-model is a word embedding and a good way to improve a word embedding over learning it from scratch as part of the model (as in the baseline model) is to use pre-trained word embeddings.

In this section, we will explore the impact of using a pre-trained word embedding on the model. Specifically:

1. Training a Word2Vec Model
2. Training a Word2Vec Model + Fine Tuning

## Trained word2vec Embedding

An efficient learning algorithm for pre-training a word embedding from a corpus of text is the word2vec algorithm.

You can learn more about the word2vec algorithm here:

- Word2Vec Google Code Project

We can use this algorithm to train a new standalone set of word vectors using the cleaned photo descriptions in the dataset.

The Gensim library provides access to an implementation of the algorithm that we can use to pre-train the embedding.

First, we must load the clean photo descriptions for the training dataset, as before.

Next, we can fit the word2vec model on all of the clean descriptions. We should note that this includes more descriptions than the 50 used in the training dataset. A fairer model for these experiments should only be trained on those descriptions in the training dataset.

Once fit, we can save the words and word vectors to an ASCII file, perhaps for later inspection or visualization.

```
1  # train word2vec model
2  lines = [s.split() for s in train_descriptions.values()]
3  model = Word2Vec(lines, size=100, window=5, workers=8, min_count=1)
4  # summarize vocabulary size in model
5  words = list(model.wv.vocab)
6  print('Vocabulary size: %d' % len(words))
7
8  # save model in ASCII (word2vec) format
9  filename = 'custom_embedding.txt'
10 model.wv.save_word2vec_format(filename, binary=False)
```

The word embedding is saved to the file '*custom_embedding.txt*'.

Now, we can load the embedding into memory, retrieve only the word vectors for the words in our vocabulary, then save them to a new file.

```
1  # load the whole embedding into memory
2  embedding = dict()
3  file = open('custom_embedding.txt')
4  for line in file:
5      values = line.split()
6      word = values[0]
7      coefs = asarray(values[1:], dtype='float32')
8      embedding[word] = coefs
9  file.close()
10 print('Embedding Size: %d' % len(embedding))
11
12 # summarize vocabulary
13 all_tokens = ' '.join(train_descriptions.values()).split()
14 vocabulary = set(all_tokens)
15 print('Vocabulary Size: %d' % len(vocabulary))
16
17 # get the vectors for words in our vocab
```

```
18  cust_embedding = dict()
19  for word in vocabulary:
20      # check if word in embedding
21      if word not in embedding:
22          continue
23      cust_embedding[word] = embedding[word]
24  print('Custom Embedding %d' % len(cust_embedding))
25
26  # save
27  dump(cust_embedding, open('word2vec_embedding.pkl', 'wb'))
28  print('Saved Embedding')
```

The complete example is listed below.

```
1   # prepare word vectors for captioning model
2
3   from numpy import asarray
4   from pickle import dump
5   from gensim.models import Word2Vec
6
7   # load doc into memory
8   def load_doc(filename):
9       # open the file as read only
10      file = open(filename, 'r')
11      # read all text
12      text = file.read()
13      # close the file
14      file.close()
15      return text
16
17  # load a pre-defined list of photo identifiers
18  def load_set(filename):
19      doc = load_doc(filename)
20      dataset = list()
21      # process line by line
22      for line in doc.split('\n'):
23          # skip empty lines
24          if len(line) < 1:
25              continue
26          # get the image identifier
27          identifier = line.split('.')[0]
28          dataset.append(identifier)
29      return set(dataset)
30
31  # split a dataset into train/test elements
32  def train_test_split(dataset):
33      # order keys so the split is consistent
34      ordered = sorted(dataset)
35      # return split dataset as two new sets
36      return set(ordered[:100]), set(ordered[100:200])
37
38  # load clean descriptions into memory
39  def load_clean_descriptions(filename, dataset):
40      # load document
41      doc = load_doc(filename)
42      descriptions = dict()
43      for line in doc.split('\n'):
44          # split line by white space
45          tokens = line.split()
46          # split id from description
47          image_id, image_desc = tokens[0], tokens[1:]
48          # skip images not in the set
49          if image_id in dataset:
50              # store
```

```
51            descriptions[image_id] = 'startseq ' + ' '.join(image_desc) + ' endseq'
52        return descriptions
53
54  # load dev set
55  filename = 'Flickr8k_text/Flickr_8k.devImages.txt'
56  dataset = load_set(filename)
57  print('Dataset: %d' % len(dataset))
58  # train-test split
59  train, test = train_test_split(dataset)
60  print('Train=%d, Test=%d' % (len(train), len(test)))
61  # descriptions
62  train_descriptions = load_clean_descriptions('descriptions.txt', train)
63  print('Descriptions: train=%d' % len(train_descriptions))
64
65  # train word2vec model
66  lines = [s.split() for s in train_descriptions.values()]
67  model = Word2Vec(lines, size=100, window=5, workers=8, min_count=1)
68  # summarize vocabulary size in model
69  words = list(model.wv.vocab)
70  print('Vocabulary size: %d' % len(words))
71
72  # save model in ASCII (word2vec) format
73  filename = 'custom_embedding.txt'
74  model.wv.save_word2vec_format(filename, binary=False)
75
76  # load the whole embedding into memory
77  embedding = dict()
78  file = open('custom_embedding.txt')
79  for line in file:
80      values = line.split()
81      word = values[0]
82      coefs = asarray(values[1:], dtype='float32')
83      embedding[word] = coefs
84  file.close()
85  print('Embedding Size: %d' % len(embedding))
86
87  # summarize vocabulary
88  all_tokens = ' '.join(train_descriptions.values()).split()
89  vocabulary = set(all_tokens)
90  print('Vocabulary Size: %d' % len(vocabulary))
91
92  # get the vectors for words in our vocab
93  cust_embedding = dict()
94  for word in vocabulary:
95      # check if word in embedding
96      if word not in embedding:
97          continue
98      cust_embedding[word] = embedding[word]
99  print('Custom Embedding %d' % len(cust_embedding))
100
101 # save
102 dump(cust_embedding, open('word2vec_embedding.pkl', 'wb'))
103 print('Saved Embedding')
```

Running this example creates a new dictionary mapping of word-to-word vectors stored in the file '*word2vec_embedding.pkl*'.

```
1  Dataset: 1000
2  Train=100, Test=100
3  Descriptions: train=100
4  Vocabulary size: 365
5  Embedding Size: 366
6  Vocabulary Size: 365
7  Custom Embedding 365
```

```
8  Saved Embedding
```

Next, we can load this embedding and use the word vectors as the fixed weights in an Embedding layer.

Below provides the *load_embedding()* function that loads the custom word2vec embedding and returns the new Embedding layer for use in the model.

```
1  # load a word embedding
2  def load_embedding(tokenizer, vocab_size, max_length):
3      # load the tokenizer
4      embedding = load(open('word2vec_embedding.pkl', 'rb'))
5      dimensions = 100
6      trainable = False
7      # create a weight matrix for words in training docs
8      weights = zeros((vocab_size, dimensions))
9      # walk words in order of tokenizer vocab to ensure vectors are in the right index
10     for word, i in tokenizer.word_index.items():
11         if word not in embedding:
12             continue
13         weights[i] = embedding[word]
14     layer = Embedding(vocab_size, dimensions, weights=[weights], input_length=max_length, train
15     return layer
```

We can use it in our model by calling the function directly from our *define_model()* function.

```
1  # define the captioning model
2  def define_model(tokenizer, vocab_size, max_length):
3      # feature extractor (encoder)
4      inputs1 = Input(shape=(7, 7, 512))
5      fe1 = GlobalMaxPooling2D()(inputs1)
6      fe2 = Dense(128, activation='relu')(fe1)
7      fe3 = RepeatVector(max_length)(fe2)
8      # embedding
9      inputs2 = Input(shape=(max_length,))
10     emb2 = load_embedding(tokenizer, vocab_size, max_length)(inputs2)
11     emb3 = LSTM(256, return_sequences=True)(emb2)
12     emb4 = TimeDistributed(Dense(128, activation='relu'))(emb3)
13     # merge inputs
14     merged = concatenate([fe3, emb4])
15     # language model (decoder)
16     lm2 = LSTM(500)(merged)
17     lm3 = Dense(500, activation='relu')(lm2)
18     outputs = Dense(vocab_size, activation='softmax')(lm3)
19     # tie it together [image, seq] [word]
20     model = Model(inputs=[inputs1, inputs2], outputs=outputs)
21     model.compile(loss='categorical_crossentropy', optimizer='adam', metrics=['accuracy'])
22     return model
23
24 model_name = 'seq_w2v_fixed'
```

We can see some lift on the training dataset, perhaps no real notable change on the test dataset.

```
1            train       test
2  count  3.000000   3.000000
3  mean   0.096780   0.047540
4  std    0.055073   0.008445
5  min    0.033511   0.038340
6  25%    0.078186   0.043840
7  50%    0.122861   0.049341
8  75%    0.128414   0.052140
9  max    0.133967   0.054939
```

## Trained word2vec Embedding with Fine Tuning

We can repeat the previous experiment and allow the model to tune the word vectors while fitting the model.

The updated *load_embedding()* function that permits the embedding layer to be fine-tuned is listed below.

```
1   # load a word embedding
2   def load_embedding(tokenizer, vocab_size, max_length):
3       # load the tokenizer
4       embedding = load(open('word2vec_embedding.pkl', 'rb'))
5       dimensions = 100
6       trainable = True
7       # create a weight matrix for words in training docs
8       weights = zeros((vocab_size, dimensions))
9       # walk words in order of tokenizer vocab to ensure vectors are in the right index
10      for word, i in tokenizer.word_index.items():
11          if word not in embedding:
12              continue
13          weights[i] = embedding[word]
14      layer = Embedding(vocab_size, dimensions, weights=[weights], input_length=max_length, trair
15      return layer
16
17  model_name = 'seq_w2v_tuned'
```

Again, we do not see much difference in using these pre-trained word embedding vectors over the baseline model.

```
1              train       test
2   count   3.000000   3.000000
3   mean    0.065297   0.042712
4   std     0.080194   0.007697
5   min     0.017675   0.034593
6   25%     0.019003   0.039117
7   50%     0.020332   0.043641
8   75%     0.089108   0.046772
9   max     0.157885   0.049904
```

# Analysis of Results

We have performed a few experiments on a very small sample (1.6%) from the Flickr8k training dataset of 8,000 photos.

It is possible that the sample is too small, that the models were not trained for long enough, and that 3 repeats of each model results in too much variance. These aspects can also be tested by evaluated by designing experiments such as:
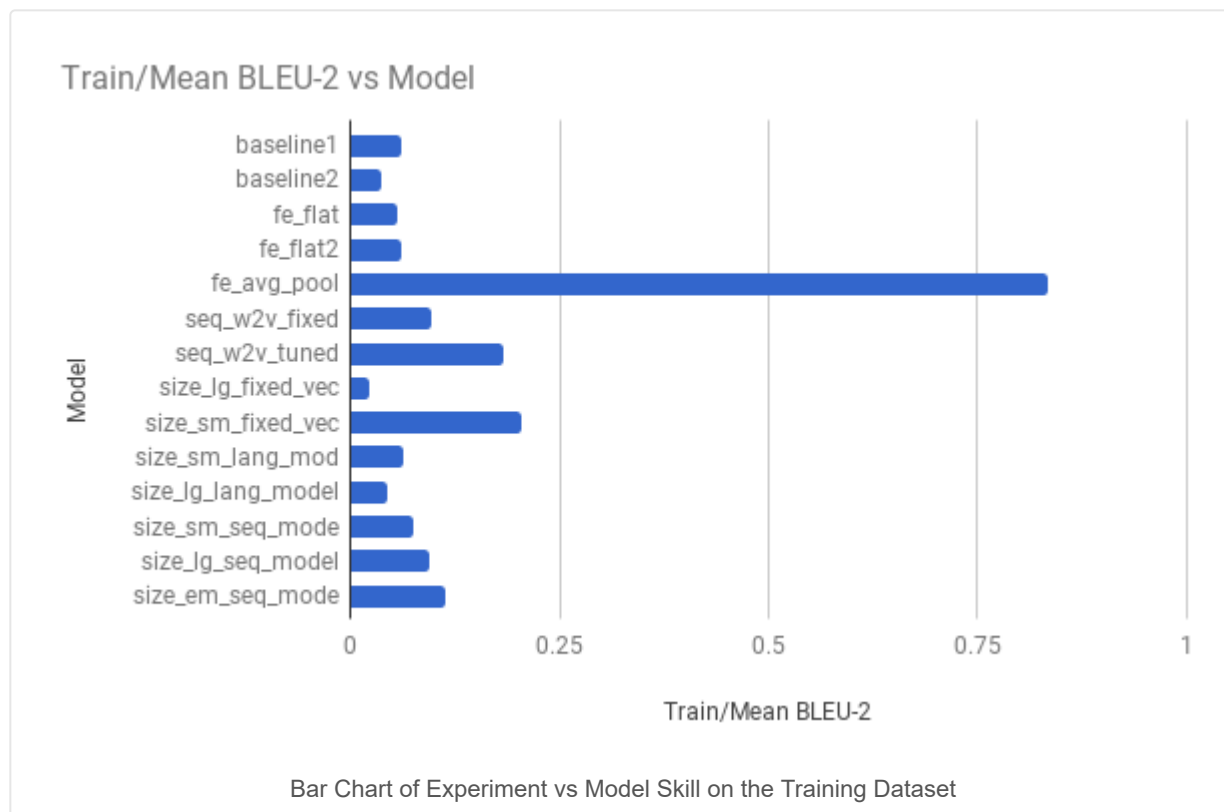
1. Does model skill scale with the size of the dataset?
2. Do more epochs result in better skill?
3. Do more repeats result in a skill with less variance?

Nevertheless, we have some ideas on how we might configure a model for the fuller dataset.
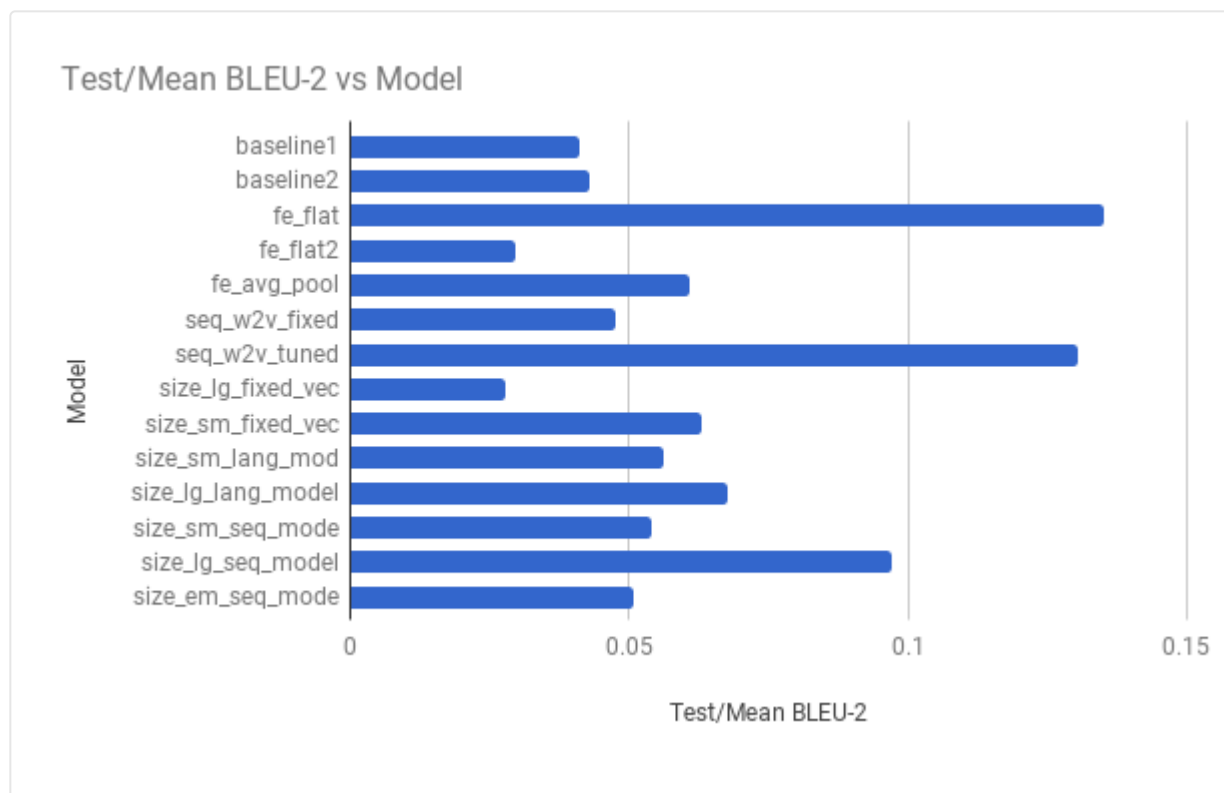
Below is a summary of the mean results from the experiments performed in this tutorial.

It is helpful to review a graph of the results. If we had more repeats, a box and whisker plot for each distribution of scores might be a good visualization. Here we use a simple bar graph. Remember, that larger BLEU scores are better.

Results on the training dataset:



Bar Chart of Experiment vs Model Skill on the Training Dataset

Results on the test dataset:

Bar Chart of Experiment vs Model Skill on the Test Dataset

From just looking at the mean results on the test dataset, we can suggest:

- Perhaps pooling is not required after the photo feature extractor (fe_flat at 0.135231).
- Perhaps average pooling offers an advantage over max pooling after the photo feature extractor (fe_avg_pool at 0.060847).
- Perhaps a smaller sized fixed-length vector after the sub-models is a good idea (size_sm_fixed_vec at 0.063148).
- Perhaps adding more layers to the language model offers some benefit (size_lg_lang_model at 0.067658).
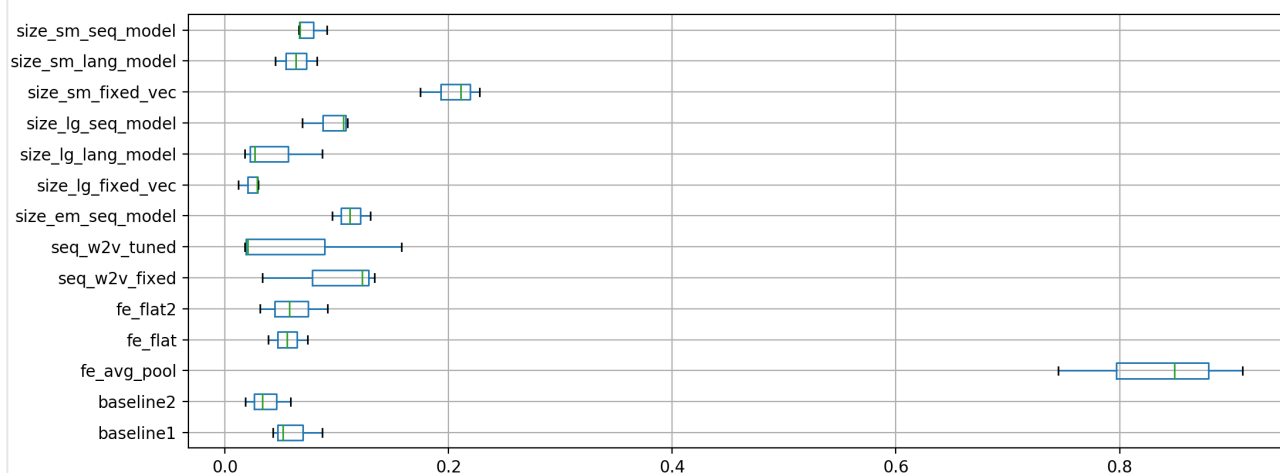- Perhaps adding more layers to the sequence model offers some benefit (size_lg_seq_model at 0.09697).

I would also recommend exploring combinations of these suggestions.

We can also review the distribution of results.

Below is some code to load the saved results from each experiment and create a box-and-whisker plot of results on the train and test sets for review.
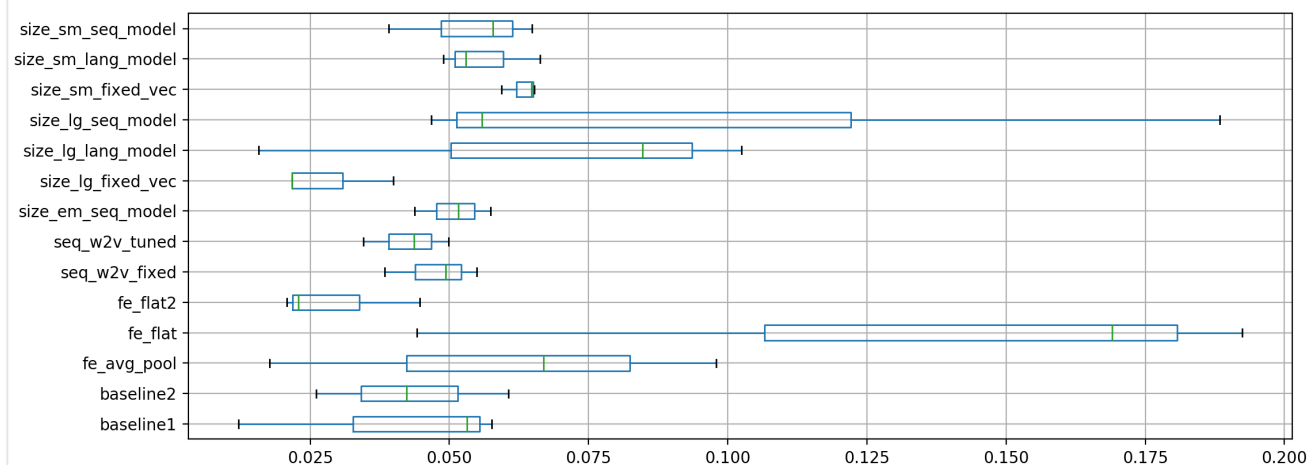
```
1   from os import listdir
2   from pandas import read_csv
3   from pandas import DataFrame
4   from matplotlib import pyplot
5
6   # load all .csv results into a dataframe
7   train, test = DataFrame(), DataFrame()
8   directory = 'results'
9   for name in listdir(directory):
10      if not name.endswith('csv'):
11          continue
12      filename = directory + '/' + name
13      data = read_csv(filename, header=0)
14      experiment = name.split('.')[0]
15      train[experiment] = data['train']
16      test[experiment] = data['test']
17
18  # plot results on train
19  train.boxplot(vert=False)
20  pyplot.show()
21  # plot results on test
22  test.boxplot(vert=False)
23  pyplot.show()
```

Distribution of results on the training dataset.

Box and Whisker Plot of Experiment vs Model Skill on the Training Dataset

Distribution of results on the test dataset.



Box and Whisker Plot of Experiment vs Model Skill on the Test Dataset

A review of these distributions suggests:

- The spread on the flat results is large; perhaps going with average pooling might be safer.
- The spread on the larger language model is large and skewed in the wrong/risky direction.
- The spread on the larger sequence model is large and skewed in the right direction.
- There may be some benefit in a smaller fixed-length vector size.

I would expect increasing repeats to 5, 10, or 30 would tighten up these distributions somewhat.

# Further Reading

This section provides more resources on the topic if you are looking go deeper.

## Papers

- Show and Tell: A Neural Image Caption Generator, 2015.
- Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, 2016.
- Network In Network, 2013.

## Related Captioning Projects

- caption_generator: An image captioning project
- Keras Image Caption
- Neural Image Captioning (NIC)
- Keras deep learning for image caption retrieval
- DataLab Cup 2: Image Captioning

## Other

- Framing image description as a ranking task: data, models and evaluation metrics.

## API

- Keras VGG16 API
- Gensim word2vec API

# Summary

In this tutorial, you discovered how you can use a small sample of the photo captioning dataset to explore different model designs.

Specifically, you learned:

- How to prepare data for photo captioning modeling.
- How to design a baseline and test harness to evaluate the skill of models and control for their stochastic nature.
- How to evaluate properties like model skill, feature extraction model, and word embeddings in order to lift model skill.
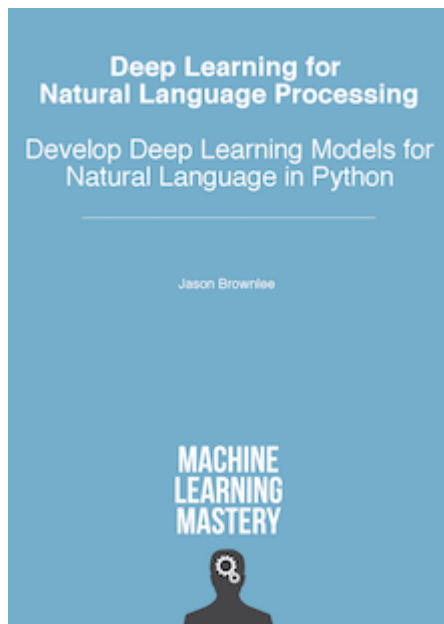
What experiments can you think up?
What else have you tried?
What are the best results you can get on the train and test dataset?

Let me know in the comments below.

---

# Develop Deep Learning models for Text Data Today!

## Develop Your Own Text models in Minutes

…with just a few lines of python code

Discover how in my new Ebook:
Deep Learning for Natural Language Processing

It provides **self-study tutorials** on topics like:
*Bag-of-Words, Word Embedding, Language Models, Caption Generation, Text Translation* and much more…

## Finally Bring Deep Learning to your Natural Language Processing Projects

Skip the Academics. Just Results.

Click to learn more.

---

### About Jason Brownlee

Jason Brownlee, PhD is a machine learning specialist who teaches developers how to get results with modern machine learning methods via hands-on tutorials.

View all posts by Jason Brownlee →

---

---

## 41 Responses to *How to Use Small Experiments to Develop a Caption Generation Model in Keras*

**Emil** November 24, 2017 at 6:39 am #       REPLY ↰

Hats off, another ace tutorial!

I'm curious how the TimeDistributed layer impacts the data before the concatenation. Is it possible to skip it? Also, is there a reason you are using VGG instead of the InceptionResNetV2 class other than memory/compute constraints.

Thanks!

**Jason Brownlee** November 24, 2017 at 9:52 am #

REPLY ↩

I chose VGG because it is smaller and simpler. You can use anything you wish.

You can skip the TimeDistributed as Dense can support time steps now I believe. I like it in there as it reminds me what is going on (e.g. outputting time steps).

**Alex** November 24, 2017 at 6:16 pm #

REPLY ↩

Ho jason, why don t you reset the LSTM states between the inputs related to different images? Ad they are not related to the same sequences.

**Jason Brownlee** November 25, 2017 at 10:14 am #

REPLY ↩

For speed of training.

It's a great suggestion though, try it and see if it lifts skill! Let me know how you go.

**Alex** November 28, 2017 at 1:39 am #

REPLY ↩

Thanks! In order to try this way, should I set stateful=True (avoiding the LSTM to reset itself automatically) and manually run model.reset_states() before training a single batch? (each batch is related to the sequence of a single image).

**Jason Brownlee** November 28, 2017 at 8:38 am #

REPLY ↩

Yes.

**Emil** December 13, 2017 at 12:10 am #

REPLY ↩

What's the logic of the +1 when you are creating the vocab len: "vocab_size = len(tokenizer.word_index) + 1"? Is it to leave remove for the 0?

Thanks

**Jason Brownlee** December 13, 2017 at 5:39 am #

REPLY ↩

Good question, to make space for 0 – words in the vocab start at 1.

**Yang Cheng** March 12, 2018 at 8:50 pm #

REPLY ↩

Does it mean we leave index 0 for 'endseq' token?

**Jason Brownlee** March 13, 2018 at 6:27 am #

REPLY ↩

No, the start and end tokens are legit parts of the problem. The model must specify when the sequence has ended.

**xiaolian** December 23, 2017 at 1:54 am #

REPLY ↩

i got a error : Error when checking input: expected input_11 to have 4 dimensions, but got array with shape (28, 4096)

**Jason Brownlee** December 23, 2017 at 5:21 am #

REPLY ↩

Are you able to confirm that your libraries are up to date and that you copied all of the code from the post?

**xiaolian** December 23, 2017 at 1:40 pm #

REPLY ↩

can i get the code in github?

**Jason Brownlee** December 24, 2017 at 4:51 am #

REPLY ↩

The code is on the post, why do you need it on github?

**Niels** February 24, 2018 at 12:41 am #

REPLY ↩

Hello Jason.

Thank you very much for your article.
I have a practical question regarding training an encoder decoder network.

So basically I have images of serial numbers and I want to predict the full serial number. (E.G 018F6176)

So if I train a network to predict the next character I would practically have to create a for loop predicting untill I reach max length or the stop word.

This I get. However how can i structure my data?

What I have is the following

data is a numpy float array of shape (nb_samples width, height, nb_channels=3,).

labels is a numpy integer array of shape (nb_samples, max_caption_len)

So if I were to construct a similar dataset with the structure you get from create_sequences() my number of samples would increase and how do I ensure that images are loaded so as the sequence for each images comes in the correct order(or does that matter?)

Hopefully this makes sense to you.

Best regards

Niels

---

**Jason Brownlee** February 24, 2018 at 9:17 am #                  REPLY ↩

Great question.

The caption model is the approach you want to use I believe.

See this post, specifically the section titled "ord-By-Word Model":
https://machinelearningmastery.com/prepare-photo-caption-dataset-training-deep-learning-model/

It will show you exactly how to prepare your data and how to think about it.

---

**Niels** February 26, 2018 at 8:01 pm #                  REPLY ↩

Thank you. This is exactly what I was looking for. One question though. Does the order of input matter, so let's say all samples for image 1 comes in order, all samples for image 2 next etc. Because then I would have to create a batch generator, however I would like to avoid that.

My initial thought is it doesn't really matter since you are just training the algorithm to recognize what comes next regarding previous input. Is that a correct assumption?

Best Regards

Niels

---

**Jason Brownlee** February 27, 2018 at 6:25 am #                  REPLY ↩

I think so, I think all samples for a give photo should be together by intuition (LSTMs have memory across the batch), but testing all assumptions is a good idea.

---

**Niels** February 27, 2018 at 9:27 pm #

Okay. I created a model both with the data_generator and randomly splitting in train,test(image order scrambled). The tran, test model did not seem to converge(max 17% validation accuracy). However the data_generator model(where order is preserved) reaches 99.5% validation accuracy , so that's pretty awesome.

Thank you for your time 🙂

**Jason Brownlee** February 28, 2018 at 6:03 am #

Very nice. Thanks for running the experiment!

REPLY ↰

**Binay** March 7, 2018 at 12:18 pm #

Why and How the dimension of input_1 is (7, 7, 512)

REPLY ↰

**Jason Brownlee** March 7, 2018 at 3:06 pm #

That is the shape that we saved the extracted features from the photos.

REPLY ↰

**Steven** March 23, 2018 at 9:03 pm #

Thanks for your tutorials, Jason. I can run this tutorial in contrary with the tutorial on the same dataset https://machinelearningmastery.com/develop-a-deep-learning-caption-generation-model-in-python/ where memory errors occur.

I compared both examples and found that the feature extraction of the images is different. This example the features are a 3-dimensional array with the shape (7, 7, 512) the other example the result is a 1-dimensional 4,096 element vector. The features.pkl file of this experiment is much larger compared to the other example. What is the reason of the different shape of the features?

REPLY ↰

**Jason Brownlee** March 24, 2018 at 6:29 am #

Perhaps the VGG model is cut at different points? e.g. keeping the dense or discarding it and working with the CNN outputs. I don't recall, but a comparison of the code would make it clearer.

REPLY ↰

**Ashish** March 28, 2018 at 4:56 pm #

Thanks for the tutorial, Jason. Here you are using features from a pretrained cnn model. But I want to backpropagate the error through the cnn as well and hence jointly train the cnn and lstm. i.e. load the cnn through a pretrained model and than train it further. Can you suggest a way to do that? Thanks.

**Jason Brownlee** March 29, 2018 at 6:31 am #                   REPLY ↩

Yes, you can load the CNN as part of your in memory model.

**Vector** April 7, 2018 at 3:09 pm #                   REPLY ↩

Hello Jason,

Thanks for this great posts and the one in here (https://machinelearningmastery.com/develop-a-deep-learning-caption-generation-model-in-python/)!!

I was trying to run the example you posted in here.I was able to train using different learning rate and more epoch on your "average pooling" model and saved them in a .h5 file.

However when I try to use the code in (https://machinelearningmastery.com/develop-a-deep-learning-caption-generation-model-in-python/) to Generate New Captions with the "Average Model" I trained, I got this error:
ValueError: Error when checking : expected input_1 to have 4 dimensions, but got array with shape (1, 4096)

Just wonder do you know how to use this .h5 file i just trained to make prediction on one other photo? Thanks!!

**Jason Brownlee** April 8, 2018 at 6:14 am #                   REPLY ↩

It looks like your photo features might have too many dimensions. Change your code to provide the photo pixels directly. e.g.: photo[0]

**Md. Zakir Hossain** May 25, 2018 at 6:46 pm #                   REPLY ↩

Hai Jason,
that's a great post. Really it is very helpful for me. However, I get the following error:

File "", line 1, in
runfile('C:/Users/33083707/Codes/Projects/Final.py', wdir='C:/Users/33083707/Codes/Projects')

File "C:\Users\33083707\Anaconda3\lib\site-packages\spyder\utils\site\sitecustomize.py", line 866, in runfile
execfile(filename, namespace)

File "C:\Users\33083707\Anaconda3\lib\site-packages\spyder\utils\site\sitecustomize.py", line 102, in execfile

exec(compile(f.read(), filename, 'exec'), namespace)

File "C:/Users/33083707/Codes/Projects/Final.py", line 244, in
test_score = evaluate_model(model, test_descriptions, test_features, tokenizer, max_length)

File "C:/Users/33083707/Codes/Projects/Final.py", line 202, in evaluate_model
bleu = corpus_bleu(actual, predicted)

File "C:\Users\33083707\Anaconda3\lib\site-packages\nltk\translate\bleu_score.py", line 146, in corpus_bleu
p_i = modified_precision(references, hypothesis, i)

File "C:\Users\33083707\Anaconda3\lib\site-packages\nltk\translate\bleu_score.py", line 287, in
modified_precision
return Fraction(numerator, denominator, _normalize=False)

File "C:\Users\33083707\Anaconda3\lib\fractions.py", line 186, in __new__
raise ZeroDivisionError('Fraction(%s, 0)' % numerator)

ZeroDivisionError: Fraction(0, 0)

---

**Jason Brownlee** May 26, 2018 at 5:51 am #                    REPLY ↩

Sorry to hear that, I have some suggestions to try here:

https://machinelearningmastery.com/faq/single-faq/why-does-the-code-in-the-tutorial-not-work-for-me

---

**Mhemmed Elly** June 29, 2018 at 9:29 pm #                    REPLY ↩

Hello Dr. Jason,

I want to try this on my own data, but I don't know how to prepare it. I have a folder of images and a CSV file. Each line in the csv file contains an image ID/name and a caption that describes the image. Any advice on how to proceed?

Thank you

---

**Jason Brownlee** June 30, 2018 at 6:07 am #                    REPLY ↩

Start by writing code to load the images into memory.

Perhaps you can use PIL or Pillo to load the images?

---

**Jeet Sen Sarma** July 6, 2018 at 6:53 pm #                    REPLY ↩

in the following part of the code :

```
# run experiment
train_results, test_results = list(), list()
for i in range(n_repeats):
# define the model
model = define_model(vocab_size, max_length)
# fit model
model.fit_generator(data_generator(train_descriptions, train_features, tokenizer, max_length,
n_photos_per_update), steps_per_epoch=n_batches_per_epoch, epochs=n_epochs, verbose=verbose)
# evaluate model on training data
train_score = evaluate_model(model, train_descriptions, train_features, tokenizer, max_length)
test_score = evaluate_model(model, test_descriptions, test_features, tokenizer, max_length)
# store
train_results.append(train_score)
test_results.append(test_score)
print('>%d: train=%f test=%f' % ((i+1), train_score, test_score))
```

why are you defining the model within the for loop ? shouldn't it be defined outside of it ?

---

**Jason Brownlee** July 7, 2018 at 6:13 am #                        REPLY ↰

Nope, I want a new model (random weights) for each repeat.

Why would I define outside the loop?

---

**Jeet** July 9, 2018 at 2:47 pm #                        REPLY ↰

Could you please explain why is that? I think that a model has to be defined, and the weights are to be trained by learning from examples. But defining a new model everytime would delete the previously learned weights. Where am I going wrong?

---

**Jason Brownlee** July 10, 2018 at 6:41 am #                        REPLY ↰

Yes, we want to throw it away. We are repeating the experiment with a model each time so we can see how well the model on our data performs on average. Not the performance of any one randomly fit model.

You can learn more here:
https://machinelearningmastery.com/evaluate-skill-deep-learning-models/

---

**Md. Zakir Hossain** July 9, 2018 at 6:13 pm #

Hi Jason, Thank you very much.

I tried to use DenseNet-121 model instead of VGG16 model to extract the features. But I got an error like:

ValueError: Error when checking input: expected input_1 to have 4 dimensions, but got array with shape (21, 1024).

Can you please suggest me anything about this.

**Jeet** July 10, 2018 at 4:26 am #

This is because both the CNN models have different output dimensions of their last layer from where you are pulling the features for your images. Best guess is to convert the output dimension of your DenseNet model into the same dimension of that of VGG-16.
Also, if you use Flatten() instead of GlobalMaxPooling2D() then I guess you will not face the same problem.

**Jason Brownlee** July 10, 2018 at 6:44 am #

Not really, I don't know about the things you are trying.

**Omnia** October 24, 2018 at 8:25 am #

Hi Jason

Thanks for the wonderful tutorial as always

I got this result on 70 epochs and 3 repeats
it seems that the accuracy increases when we have a large number of epochs

but I don't understand why I got the train and test = 0

this is my result

Epoch 70/70
– 14s – loss: 1.8873 – acc: 0.4100
Actual: startseq child and woman are at waters edge in big city endseq
Predicted: startseq child woman are edge edge in big endseq

Actual: startseq boy with stick kneeling in front of goalie net endseq
Predicted: startseq boy boy in in front and and and of of of of of in front and of of in front boy front and of of

Actual: startseq woman crouches near three dogs in field endseq
Predicted: startseq two dog dogs in near near near near near near near near near near near near near near near

near near near near near near near

Actual: startseq boy bites hard into treat while he sits outside endseq
Predicted: startseq boy bites sits while while outside endseq

Actual: startseq person eats takeout while watching small television endseq
Predicted: startseq person eats while takeout small small television endseq

Actual: startseq couple with young child wrapped in blanket sitting on concrete step endseq
Predicted: startseq boy girls in in in bike endseq

Actual: startseq adults and children stand and play in front of steps near wooded area endseq
Predicted: startseq boy and young boy and and and and and and and and and and and and and and
and and and and and and

Actual: startseq boy in grey pajamas is jumping on the couch endseq
Predicted: startseq boy man into if rough endseq

Actual: startseq boy holding kitchen utensils and making threatening face endseq
Predicted: startseq girl holds shorts shorts shorts with down her endseq

Actual: startseq man in green hat is someplace up high endseq
Predicted: startseq boy in with with waves endseq

>3: train=0.000000 test=0.000000
train test
count 3.000000e+00 3.000000e+00
mean 6.038250e-02 4.013559e-155
std 1.045856e-01 2.160777e-156
min 2.413005e-78 3.808820e-155
25% 2.799294e-78 3.900627e-155
50% 3.185584e-78 3.992435e-155
75% 9.057374e-02 4.115929e-155
max 1.811475e-01 4.239423e-155

Also, I got this warning

UserWarning:
The hypothesis contains 0 counts of 3-gram overlaps.
Therefore the BLEU score evaluates to 0, independently of
how many N-gram overlaps of lower order it contains.
Consider using lower n-gram order or use SmoothingFunction()
warnings.warn(_msg)

would you mind giving your opinion on my experiment and suggesting a better solution for the warning?

Thanks

---

**Jason Brownlee** October 24, 2018 at 2:41 pm #                    REPLY ↩

Nice work!

BLEU may give warnings when some predictions are shorter than expected. This might help understand the BLEU score calculation some more:

https://machinelearningmastery.com/calculate-bleu-score-for-text-python/

---

# Leave a Reply

Name (required)

Email (will not be published) (required)

Website
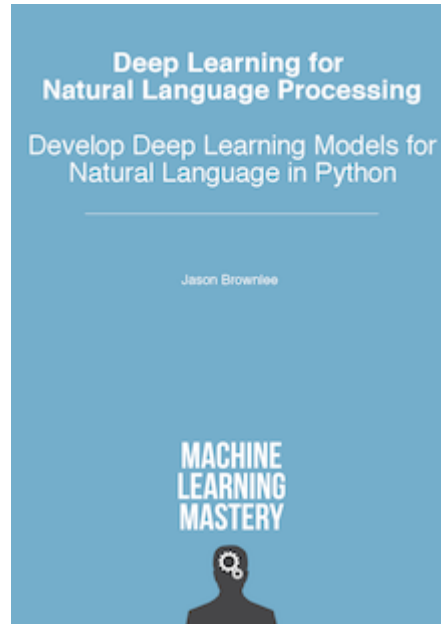
SUBMIT COMMENT

**Welcome to Machine Learning Mastery!**

Hi, I'm Jason Brownlee, PhD
I write tutorials to help developers (*like you*) get results with machine learning.

Read More

**Deep Learning for NLP**
Develop deep learning models for text data.

POPULAR

**How to Develop a Neural Machine Translation System from Scratch**
JANUARY 10, 2018

**So, You are Working on a Machine Learning Problem…**
APRIL 4, 2018

**How to Develop an N-gram Multichannel Convolutional Neural Network for Sentiment Analysis**
JANUARY 12, 2018

**How to Make Predictions with Keras**
APRIL 9, 2018

**11 Classical Time Series Forecasting Methods in Python (Cheat Sheet)**
AUGUST 6, 2018

**You're Doing it Wrong. Why Machine Learning Does Not Have to Be So Hard**
MARCH 28, 2018

**How to Develop LSTM Models for Multi-Step Time Series Forecasting of Household Power Consumption**
OCTOBER 10, 2018

## You might also like…

- How to Install Python for Machine Learning

- Your First Machine Learning Project in Python
- Your First Neural Network in Python
- Your First Classifier in Weka
- Your First Time Series Forecasting Project

- Your First Machine Learning Project in Python
- Your First Neural Network in Python