

flickr_caption_generator

August 19, 2022

1 COMS W4705 - Homework 5

1.1 Image Captioning with Conditioned LSTM Generators

Daniel Bauer bauer@cs.columbia.edu

Follow the instructions in this notebook step-by step. Much of the code is provided, but some sections are marked with **todo**.

Specifically, you will build the following components:

- Create matrices of image representations using an off-the-shelf image encoder.
- Read and preprocess the image captions.
- Write a generator function that returns one training instance (input/output sequence pair) at a time.
- Train an LSTM language generator on the caption data.
- Write a decoder function for the language generator.
- Add the image input to write an LSTM caption generator.
- Implement beam search for the image caption generator.

Please see the special submission instructions at the bottom.

This notebook assumes that you are running it on a machine with the possibility of using a GPU (like the VM on GCP as instructed in tutorial).

1.1.1 Getting Started

First, run the following commands to make sure you have all required packages.

```
[1]: import os
from collections import defaultdict
import numpy as np
import PIL
from matplotlib import pyplot as plt
%matplotlib inline

from tensorflow.keras import Sequential, Model
from tensorflow.keras.layers import Embedding, LSTM, Dense, Input,
    Bidirectional, RepeatVector, Concatenate, Activation
from tensorflow.keras.activations import softmax
from tensorflow.keras.utils import to_categorical
```

```

from tensorflow.keras.preprocessing.sequence import pad_sequences

from tensorflow.keras.applications.inception_v3 import InceptionV3

from tensorflow.keras.optimizers import Adam

```

1.1.2 Access to the flickr8k data

We will use the flickr8k data set, described here in more detail:

M. Hodosh, P. Young and J. Hockenmaier (2013) “Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics”, Journal of Artificial Intelligence Research, Volume 47, pages 853-899 <http://www.jair.org/papers/paper3994.html> when discussing our results

The data is available here: <https://storage.googleapis.com/4705-hw5-data/hw5data-20220809T182644Z-001.zip>

N.B.: Usage of this data is limited to this homework assignment. If you would like to experiment with the data set beyond this course, I suggest that you submit your own download request here (it's free): <https://forms.illinois.edu/sec/1713398>

[2]: *#You can download the data by running the following command (it's about 1GB)*

```

! wget https://storage.googleapis.com/4705-hw5-data/
↪hw5data-20220809T182644Z-001.zip

```

```

--2022-08-18 15:11:37--
https://storage.googleapis.com/4705-hw5-data/hw5data-20220809T182644Z-001.zip
Resolving storage.googleapis.com (storage.googleapis.com)... 74.125.134.128,
74.125.139.128, 74.125.141.128, ...
Connecting to storage.googleapis.com
(storage.googleapis.com)|74.125.134.128|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 1120506096 (1.0G) [application/zip]
Saving to: 'hw5data-20220809T182644Z-001.zip'

hw5data-20220809T18 100%[=====>]    1.04G   280MB/s   in 3.8s

2022-08-18 15:11:40 (281 MB/s) - 'hw5data-20220809T182644Z-001.zip' saved
[1120506096/1120506096]

```

If you are using a separate persistent disk to store the data, ssh into your VM instance and copy the data directory there. Either way, you want the following variable to point to the location of the data directory.

[2]: `FLICKR_PATH="hw5data"`

1.2 Part I: Image Encodings (14 pts)

The files Flickr_8k.trainImages.txt Flickr_8k.devImages.txt Flickr_8k.testImages.txt, contain a list of training, development, and test images, respectively. Let's load these lists.

```
[3]: def load_image_list(filename):  
      with open(filename, 'r') as image_list_f:  
          return [line.strip() for line in image_list_f]  
  
[4]: train_list = load_image_list(os.path.join(FLICKR_PATH, 'Flickr_8k.trainImages.  
      ↪txt'))  
      dev_list = load_image_list(os.path.join(FLICKR_PATH, 'Flickr_8k.devImages.txt'))  
      test_list = load_image_list(os.path.join(FLICKR_PATH, 'Flickr_8k.testImages.  
      ↪txt'))
```

Let's see how many images there are

```
[5]: len(train_list), len(dev_list), len(test_list)
```

```
[5]: (6000, 1000, 1000)
```

Each entry is an image filename.

```
[6]: dev_list[20]
```

```
[6]: '3693961165_9d6c333d5b.jpg'
```

The images are located in a subdirectory.

```
[7]: IMG_PATH = os.path.join(FLICKR_PATH, "Flickr8k_Dataset")
```

We can use PIL to open the image and matplotlib to display it.

```
[8]: image = PIL.Image.open(os.path.join(IMG_PATH, dev_list[20]))  
      image
```

```
[8]:
```



if you can't see the image, try

```
[9]: plt.imshow(image)
```

```
[9]: <matplotlib.image.AxesImage at 0x7ff0ee09a250>
```



We are going to use an off-the-shelf pre-trained image encoder, the Inception V3 network. The model is a version of a convolution neural network for object detection. Here is more detail about this model (not required for this project):

Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., & Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 2818-2826). https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Szegedy_Rethinking_the_Inception_CVPR_2016_

The model requires that input images are presented as 299x299 pixels, with 3 color channels (RGB). The individual RGB values need to range between 0 and 1.0. The flickr images don't fit.

```
[10]: np.asarray(image).shape
```

```
[10]: (333, 500, 3)
```

The values range from 0 to 255.

```
[11]: np.asarray(image)
```

```
[11]: array([[118, 161,  89],
          [120, 164,  89],
          [111, 157,  82],
          ...,
          [ 68, 106,  65],
          [ 64, 102,  61],
          [ 65, 104,  60]],
```

```

[[125, 168, 96],
 [121, 164, 92],
 [119, 165, 90],
 ...,
 [ 72, 115, 72],
 [ 65, 108, 65],
 [ 72, 115, 70]],

[[129, 175, 102],
 [123, 169, 96],
 [115, 161, 88],
 ...,
 [ 88, 129, 87],
 [ 75, 116, 72],
 [ 75, 116, 72]],

...,

[[ 41, 118, 46],
 [ 36, 113, 41],
 [ 45, 111, 49],
 ...,
 [ 23, 77, 15],
 [ 60, 114, 62],
 [ 19, 59, 0]],

[[100, 158, 97],
 [ 38, 100, 37],
 [ 46, 117, 51],
 ...,
 [ 25, 54, 8],
 [ 88, 112, 76],
 [ 65, 106, 48]],

[[ 89, 148, 84],
 [ 44, 112, 35],
 [ 71, 130, 72],
 ...,
 [152, 188, 142],
 [113, 151, 110],
 [ 94, 138, 75]]], dtype=uint8)

```

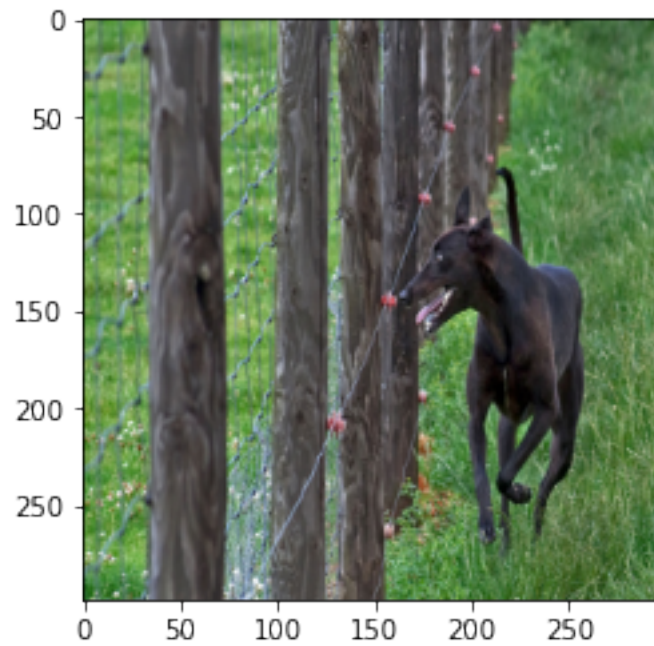
We can use PIL to resize the image and then divide every value by 255.

```

[12]: new_image = np.asarray(image.resize((299,299))) / 255.0
      plt.imshow(new_image)

```

```
[12]: <matplotlib.image.AxesImage at 0x7ff0edf6d250>
```



```
[13]: new_image.shape
```

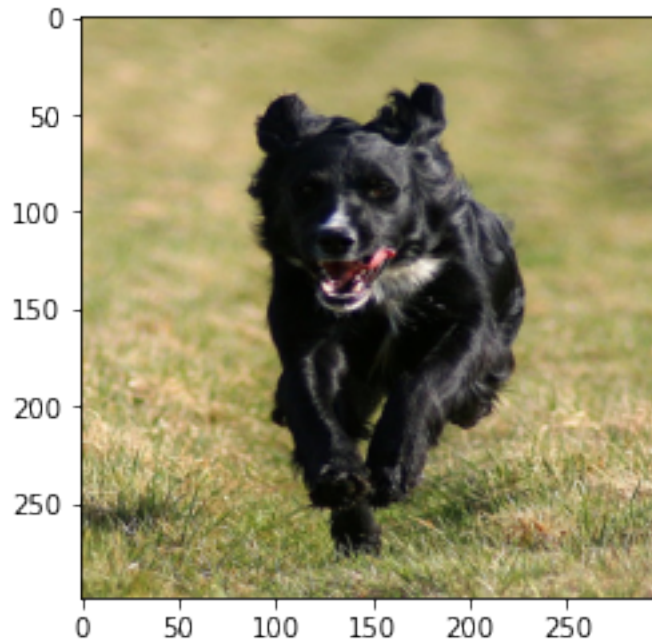
```
[13]: (299, 299, 3)
```

Let's put this all in a function for convenience.

```
[14]: def get_image(image_name):  
      image = PIL.Image.open(os.path.join(IMG_PATH, image_name))  
      return np.asarray(image.resize((299,299))) / 255.0
```

```
[15]: plt.imshow(get_image(dev_list[25]))
```

```
[15]: <matplotlib.image.AxesImage at 0x7ff0edf04290>
```

Next, we load the pre-trained Inception model.

```
[18]: img_model = InceptionV3(weights='imagenet') # This will download the weights
      ↪ files for you and might take a while.
```

```
[17]: img_model.summary() # this is quite a complex model.
```

Model: "inception_v3"

Layer (type)	Output Shape	Param #	Connected to
input_1 (InputLayer)	[(None, 299, 299, 3)]	0	[]
conv2d (Conv2D)	(None, 149, 149, 32)	864	['input_1[0][0]']
batch_normalization (Batch Normalization)	(None, 149, 149, 32)	96	['conv2d[0][0]']
activation (Activation)	(None, 149, 149, 32)	0	['batch_normalization[0][0]']


```

)

conv2d_1 (Conv2D)          (None, 147, 147, 32) 9216
['activation_0[0][0]']
)

batch_normalization_1 (BatchNormaliz (None, 147, 147, 32) 96
['conv2d_1[0][0]']
alization)
)

activation_1 (Activation)    (None, 147, 147, 32) 0
['batch_normalization_1[0][0]']
)

conv2d_2 (Conv2D)          (None, 147, 147, 64) 18432
['activation_1[0][0]']
)

batch_normalization_2 (BatchNormaliz (None, 147, 147, 64) 192
['conv2d_2[0][0]']
alization)
)

activation_2 (Activation)    (None, 147, 147, 64) 0
['batch_normalization_2[0][0]']
)

max_pooling2d (MaxPooling2D) (None, 73, 73, 64) 0
['activation_2[0][0]']

conv2d_3 (Conv2D)          (None, 73, 73, 80) 5120
['max_pooling2d[0][0]']

batch_normalization_3 (BatchNormaliz (None, 73, 73, 80) 240
['conv2d_3[0][0]']
alization)

activation_3 (Activation)    (None, 73, 73, 80) 0
['batch_normalization_3[0][0]']

conv2d_4 (Conv2D)          (None, 71, 71, 192) 138240
['activation_3[0][0]']

batch_normalization_4 (BatchNormaliz (None, 71, 71, 192) 576
['conv2d_4[0][0]']
alization)

activation_4 (Activation)    (None, 71, 71, 192) 0
['batch_normalization_4[0][0]']

```

```

max_pooling2d_1 (MaxPooling2D) (None, 35, 35, 192) 0
['activation_4[0][0]']

conv2d_8 (Conv2D) (None, 35, 35, 64) 12288
['max_pooling2d_1[0][0]']

batch_normalization_8 (BatchNormaliz (None, 35, 35, 64) 192
['conv2d_8[0][0]']
alization)

activation_8 (Activation) (None, 35, 35, 64) 0
['batch_normalization_8[0][0]']

conv2d_6 (Conv2D) (None, 35, 35, 48) 9216
['max_pooling2d_1[0][0]']

conv2d_9 (Conv2D) (None, 35, 35, 96) 55296
['activation_8[0][0]']

batch_normalization_6 (BatchNormaliz (None, 35, 35, 48) 144
['conv2d_6[0][0]']
alization)

batch_normalization_9 (BatchNormaliz (None, 35, 35, 96) 288
['conv2d_9[0][0]']
alization)

activation_6 (Activation) (None, 35, 35, 48) 0
['batch_normalization_6[0][0]']

activation_9 (Activation) (None, 35, 35, 96) 0
['batch_normalization_9[0][0]']

average_pooling2d (AveragePooling2D) (None, 35, 35, 192) 0
['max_pooling2d_1[0][0]']
ing2D)

conv2d_5 (Conv2D) (None, 35, 35, 64) 12288
['max_pooling2d_1[0][0]']

conv2d_7 (Conv2D) (None, 35, 35, 64) 76800
['activation_6[0][0]']

conv2d_10 (Conv2D) (None, 35, 35, 96) 82944
['activation_9[0][0]']

conv2d_11 (Conv2D) (None, 35, 35, 32) 6144

```

```

['average_pooling2d[0][0]']

batch_normalization_5 (BatchNo (None, 35, 35, 64) 192
['conv2d_5[0][0]']
rmalization)

batch_normalization_7 (BatchNo (None, 35, 35, 64) 192
['conv2d_7[0][0]']
rmalization)

batch_normalization_10 (BatchN (None, 35, 35, 96) 288
['conv2d_10[0][0]']
ormalization)

batch_normalization_11 (BatchN (None, 35, 35, 32) 96
['conv2d_11[0][0]']
ormalization)

activation_5 (Activation) (None, 35, 35, 64) 0
['batch_normalization_5[0][0]']

activation_7 (Activation) (None, 35, 35, 64) 0
['batch_normalization_7[0][0]']

activation_10 (Activation) (None, 35, 35, 96) 0
['batch_normalization_10[0][0]']

activation_11 (Activation) (None, 35, 35, 32) 0
['batch_normalization_11[0][0]']

mixed0 (Concatenate) (None, 35, 35, 256) 0
['activation_5[0][0]',
'activation_7[0][0]',
'activation_10[0][0]',
'activation_11[0][0]']

conv2d_15 (Conv2D) (None, 35, 35, 64) 16384
['mixed0[0][0]']

batch_normalization_15 (BatchN (None, 35, 35, 64) 192
['conv2d_15[0][0]']
ormalization)

activation_15 (Activation) (None, 35, 35, 64) 0
['batch_normalization_15[0][0]']

conv2d_13 (Conv2D) (None, 35, 35, 48) 12288
['mixed0[0][0]']

```

conv2d_16 (Conv2D)	(None, 35, 35, 96)	55296
['activation_15[0][0]']		
batch_normalization_13 (Batch Normalization)	(None, 35, 35, 48)	144
['conv2d_13[0][0]']		
batch_normalization_16 (Batch Normalization)	(None, 35, 35, 96)	288
['conv2d_16[0][0]']		
activation_13 (Activation)	(None, 35, 35, 48)	0
['batch_normalization_13[0][0]']		
activation_16 (Activation)	(None, 35, 35, 96)	0
['batch_normalization_16[0][0]']		
average_pooling2d_1 (Average Pooling2D)	(None, 35, 35, 256)	0
['mixed0[0][0]']		
conv2d_12 (Conv2D)	(None, 35, 35, 64)	16384
['mixed0[0][0]']		
conv2d_14 (Conv2D)	(None, 35, 35, 64)	76800
['activation_13[0][0]']		
conv2d_17 (Conv2D)	(None, 35, 35, 96)	82944
['activation_16[0][0]']		
conv2d_18 (Conv2D)	(None, 35, 35, 64)	16384
['average_pooling2d_1[0][0]']		
batch_normalization_12 (Batch Normalization)	(None, 35, 35, 64)	192
['conv2d_12[0][0]']		
batch_normalization_14 (Batch Normalization)	(None, 35, 35, 64)	192
['conv2d_14[0][0]']		
batch_normalization_17 (Batch Normalization)	(None, 35, 35, 96)	288
['conv2d_17[0][0]']		
batch_normalization_18 (Batch Normalization)	(None, 35, 35, 64)	192
['conv2d_18[0][0]']		

```

ormalization)

activation_12 (Activation)      (None, 35, 35, 64)    0
['batch_normalization_12[0][0]']

activation_14 (Activation)      (None, 35, 35, 64)    0
['batch_normalization_14[0][0]']

activation_17 (Activation)      (None, 35, 35, 96)    0
['batch_normalization_17[0][0]']

activation_18 (Activation)      (None, 35, 35, 64)    0
['batch_normalization_18[0][0]']

mixed1 (Concatenate)           (None, 35, 35, 288)   0
['activation_12[0][0]',
'activation_14[0][0]',
'activation_17[0][0]',
'activation_18[0][0]']

conv2d_22 (Conv2D)              (None, 35, 35, 64)   18432
['mixed1[0][0]']

batch_normalization_22 (BatchN  (None, 35, 35, 64)   192
['conv2d_22[0][0]']
ormalization)

activation_22 (Activation)      (None, 35, 35, 64)    0
['batch_normalization_22[0][0]']

conv2d_20 (Conv2D)              (None, 35, 35, 48)   13824
['mixed1[0][0]']

conv2d_23 (Conv2D)              (None, 35, 35, 96)   55296
['activation_22[0][0]']

batch_normalization_20 (BatchN  (None, 35, 35, 48)   144
['conv2d_20[0][0]']
ormalization)

batch_normalization_23 (BatchN  (None, 35, 35, 96)   288
['conv2d_23[0][0]']
ormalization)

activation_20 (Activation)      (None, 35, 35, 48)    0
['batch_normalization_20[0][0]']

activation_23 (Activation)      (None, 35, 35, 96)    0

```

```

['batch_normalization_23[0][0]']

average_pooling2d_2 (AveragePo (None, 35, 35, 288) 0
['mixed1[0][0]']
oling2D)

conv2d_19 (Conv2D) (None, 35, 35, 64) 18432
['mixed1[0][0]']

conv2d_21 (Conv2D) (None, 35, 35, 64) 76800
['activation_20[0][0]']

conv2d_24 (Conv2D) (None, 35, 35, 96) 82944
['activation_23[0][0]']

conv2d_25 (Conv2D) (None, 35, 35, 64) 18432
['average_pooling2d_2[0][0]']

batch_normalization_19 (BatchN (None, 35, 35, 64) 192
['conv2d_19[0][0]']
ormalization)

batch_normalization_21 (BatchN (None, 35, 35, 64) 192
['conv2d_21[0][0]']
ormalization)

batch_normalization_24 (BatchN (None, 35, 35, 96) 288
['conv2d_24[0][0]']
ormalization)

batch_normalization_25 (BatchN (None, 35, 35, 64) 192
['conv2d_25[0][0]']
ormalization)

activation_19 (Activation) (None, 35, 35, 64) 0
['batch_normalization_19[0][0]']

activation_21 (Activation) (None, 35, 35, 64) 0
['batch_normalization_21[0][0]']

activation_24 (Activation) (None, 35, 35, 96) 0
['batch_normalization_24[0][0]']

activation_25 (Activation) (None, 35, 35, 64) 0
['batch_normalization_25[0][0]']

mixed2 (Concatenate) (None, 35, 35, 288) 0
['activation_19[0][0]',

```

```

'activation_21[0][0]',
'activation_24[0][0]',
'activation_25[0][0]'

conv2d_27 (Conv2D)          (None, 35, 35, 64)    18432
['mixed2[0][0]']

batch_normalization_27 (BatchN (None, 35, 35, 64)    192
['conv2d_27[0][0]']
ormalization)

activation_27 (Activation)    (None, 35, 35, 64)    0
['batch_normalization_27[0][0]']

conv2d_28 (Conv2D)          (None, 35, 35, 96)    55296
['activation_27[0][0]']

batch_normalization_28 (BatchN (None, 35, 35, 96)    288
['conv2d_28[0][0]']
ormalization)

activation_28 (Activation)    (None, 35, 35, 96)    0
['batch_normalization_28[0][0]']

conv2d_26 (Conv2D)          (None, 17, 17, 384)   995328
['mixed2[0][0]']

conv2d_29 (Conv2D)          (None, 17, 17, 96)    82944
['activation_28[0][0]']

batch_normalization_26 (BatchN (None, 17, 17, 384)   1152
['conv2d_26[0][0]']
ormalization)

batch_normalization_29 (BatchN (None, 17, 17, 96)    288
['conv2d_29[0][0]']
ormalization)

activation_26 (Activation)    (None, 17, 17, 384)   0
['batch_normalization_26[0][0]']

activation_29 (Activation)    (None, 17, 17, 96)    0
['batch_normalization_29[0][0]']

max_pooling2d_2 (MaxPooling2D) (None, 17, 17, 288)   0
['mixed2[0][0]']

mixed3 (Concatenate)         (None, 17, 17, 768)   0

```



```

['activation_26[0][0]',
'activation_29[0][0]',
'max_pooling2d_2[0][0]']

conv2d_34 (Conv2D)          (None, 17, 17, 128) 98304
['mixed3[0][0]']

batch_normalization_34 (BatchN (None, 17, 17, 128) 384
['conv2d_34[0][0]']
ormalization)

activation_34 (Activation)    (None, 17, 17, 128) 0
['batch_normalization_34[0][0]']

conv2d_35 (Conv2D)          (None, 17, 17, 128) 114688
['activation_34[0][0]']

batch_normalization_35 (BatchN (None, 17, 17, 128) 384
['conv2d_35[0][0]']
ormalization)

activation_35 (Activation)    (None, 17, 17, 128) 0
['batch_normalization_35[0][0]']

conv2d_31 (Conv2D)          (None, 17, 17, 128) 98304
['mixed3[0][0]']

conv2d_36 (Conv2D)          (None, 17, 17, 128) 114688
['activation_35[0][0]']

batch_normalization_31 (BatchN (None, 17, 17, 128) 384
['conv2d_31[0][0]']
ormalization)

batch_normalization_36 (BatchN (None, 17, 17, 128) 384
['conv2d_36[0][0]']
ormalization)

activation_31 (Activation)    (None, 17, 17, 128) 0
['batch_normalization_31[0][0]']

activation_36 (Activation)    (None, 17, 17, 128) 0
['batch_normalization_36[0][0]']

conv2d_32 (Conv2D)          (None, 17, 17, 128) 114688
['activation_31[0][0]']

conv2d_37 (Conv2D)          (None, 17, 17, 128) 114688

```

```

['activation_36[0][0]']

batch_normalization_32 (BatchN (None, 17, 17, 128) 384
['conv2d_32[0][0]']
ormalization)

batch_normalization_37 (BatchN (None, 17, 17, 128) 384
['conv2d_37[0][0]']
ormalization)

activation_32 (Activation) (None, 17, 17, 128) 0
['batch_normalization_32[0][0]']

activation_37 (Activation) (None, 17, 17, 128) 0
['batch_normalization_37[0][0]']

average_pooling2d_3 (AveragePo (None, 17, 17, 768) 0
['mixed3[0][0]']
oling2D)

conv2d_30 (Conv2D) (None, 17, 17, 192) 147456
['mixed3[0][0]']

conv2d_33 (Conv2D) (None, 17, 17, 192) 172032
['activation_32[0][0]']

conv2d_38 (Conv2D) (None, 17, 17, 192) 172032
['activation_37[0][0]']

conv2d_39 (Conv2D) (None, 17, 17, 192) 147456
['average_pooling2d_3[0][0]']

batch_normalization_30 (BatchN (None, 17, 17, 192) 576
['conv2d_30[0][0]']
ormalization)

batch_normalization_33 (BatchN (None, 17, 17, 192) 576
['conv2d_33[0][0]']
ormalization)

batch_normalization_38 (BatchN (None, 17, 17, 192) 576
['conv2d_38[0][0]']
ormalization)

batch_normalization_39 (BatchN (None, 17, 17, 192) 576
['conv2d_39[0][0]']
ormalization)

```

```

activation_30 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_30[0][0]']

activation_33 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_33[0][0]']

activation_38 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_38[0][0]']

activation_39 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_39[0][0]']

mixed4 (Concatenate)           (None, 17, 17, 768)  0
['activation_30[0][0]',
'activation_33[0][0]',
'activation_38[0][0]',
'activation_39[0][0]']

conv2d_44 (Conv2D)              (None, 17, 17, 160) 122880
['mixed4[0][0]']

batch_normalization_44 (BatchN (None, 17, 17, 160) 480
['conv2d_44[0][0]']
ormalization)

activation_44 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_44[0][0]']

conv2d_45 (Conv2D)              (None, 17, 17, 160) 179200
['activation_44[0][0]']

batch_normalization_45 (BatchN (None, 17, 17, 160) 480
['conv2d_45[0][0]']
ormalization)

activation_45 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_45[0][0]']

conv2d_41 (Conv2D)              (None, 17, 17, 160) 122880
['mixed4[0][0]']

conv2d_46 (Conv2D)              (None, 17, 17, 160) 179200
['activation_45[0][0]']

batch_normalization_41 (BatchN (None, 17, 17, 160) 480
['conv2d_41[0][0]']
ormalization)

```

```

batch_normalization_46 (BatchN (None, 17, 17, 160) 480
['conv2d_46[0][0]']
ormalization)

activation_41 (Activation) (None, 17, 17, 160) 0
['batch_normalization_41[0][0]']

activation_46 (Activation) (None, 17, 17, 160) 0
['batch_normalization_46[0][0]']

conv2d_42 (Conv2D) (None, 17, 17, 160) 179200
['activation_41[0][0]']

conv2d_47 (Conv2D) (None, 17, 17, 160) 179200
['activation_46[0][0]']

batch_normalization_42 (BatchN (None, 17, 17, 160) 480
['conv2d_42[0][0]']
ormalization)

batch_normalization_47 (BatchN (None, 17, 17, 160) 480
['conv2d_47[0][0]']
ormalization)

activation_42 (Activation) (None, 17, 17, 160) 0
['batch_normalization_42[0][0]']

activation_47 (Activation) (None, 17, 17, 160) 0
['batch_normalization_47[0][0]']

average_pooling2d_4 (AveragePo (None, 17, 17, 768) 0
['mixed4[0][0]']
oling2D)

conv2d_40 (Conv2D) (None, 17, 17, 192) 147456
['mixed4[0][0]']

conv2d_43 (Conv2D) (None, 17, 17, 192) 215040
['activation_42[0][0]']

conv2d_48 (Conv2D) (None, 17, 17, 192) 215040
['activation_47[0][0]']

conv2d_49 (Conv2D) (None, 17, 17, 192) 147456
['average_pooling2d_4[0][0]']

batch_normalization_40 (BatchN (None, 17, 17, 192) 576
['conv2d_40[0][0]']

```

```

ormalization)

batch_normalization_43 (BatchN (None, 17, 17, 192) 576
['conv2d_43[0][0]'])
ormalization)

batch_normalization_48 (BatchN (None, 17, 17, 192) 576
['conv2d_48[0][0]'])
ormalization)

batch_normalization_49 (BatchN (None, 17, 17, 192) 576
['conv2d_49[0][0]'])
ormalization)

activation_40 (Activation) (None, 17, 17, 192) 0
['batch_normalization_40[0][0]'])

activation_43 (Activation) (None, 17, 17, 192) 0
['batch_normalization_43[0][0]'])

activation_48 (Activation) (None, 17, 17, 192) 0
['batch_normalization_48[0][0]'])

activation_49 (Activation) (None, 17, 17, 192) 0
['batch_normalization_49[0][0]'])

mixed5 (Concatenate) (None, 17, 17, 768) 0
['activation_40[0][0]',
'activation_43[0][0]',
'activation_48[0][0]',
'activation_49[0][0]']

conv2d_54 (Conv2D) (None, 17, 17, 160) 122880
['mixed5[0][0]']

batch_normalization_54 (BatchN (None, 17, 17, 160) 480
['conv2d_54[0][0]'])
ormalization)

activation_54 (Activation) (None, 17, 17, 160) 0
['batch_normalization_54[0][0]'])

conv2d_55 (Conv2D) (None, 17, 17, 160) 179200
['activation_54[0][0]']

batch_normalization_55 (BatchN (None, 17, 17, 160) 480
['conv2d_55[0][0]'])
ormalization)

```

```

activation_55 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_55[0][0]']

conv2d_51 (Conv2D)              (None, 17, 17, 160) 122880
['mixed5[0][0]']

conv2d_56 (Conv2D)              (None, 17, 17, 160) 179200
['activation_55[0][0]']

batch_normalization_51 (BatchN (None, 17, 17, 160) 480
['conv2d_51[0][0]']
ormalization)

batch_normalization_56 (BatchN (None, 17, 17, 160) 480
['conv2d_56[0][0]']
ormalization)

activation_51 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_51[0][0]']

activation_56 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_56[0][0]']

conv2d_52 (Conv2D)              (None, 17, 17, 160) 179200
['activation_51[0][0]']

conv2d_57 (Conv2D)              (None, 17, 17, 160) 179200
['activation_56[0][0]']

batch_normalization_52 (BatchN (None, 17, 17, 160) 480
['conv2d_52[0][0]']
ormalization)

batch_normalization_57 (BatchN (None, 17, 17, 160) 480
['conv2d_57[0][0]']
ormalization)

activation_52 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_52[0][0]']

activation_57 (Activation)      (None, 17, 17, 160)  0
['batch_normalization_57[0][0]']

average_pooling2d_5 (AveragePo (None, 17, 17, 768)  0
['mixed5[0][0]']
oling2D)

```

conv2d_50 (Conv2D)	(None, 17, 17, 192)	147456
['mixed5[0][0]']		
conv2d_53 (Conv2D)	(None, 17, 17, 192)	215040
['activation_52[0][0]']		
conv2d_58 (Conv2D)	(None, 17, 17, 192)	215040
['activation_57[0][0]']		
conv2d_59 (Conv2D)	(None, 17, 17, 192)	147456
['average_pooling2d_5[0][0]']		
batch_normalization_50 (Batch Normalization)	(None, 17, 17, 192)	576
['conv2d_50[0][0]']		
batch_normalization_53 (Batch Normalization)	(None, 17, 17, 192)	576
['conv2d_53[0][0]']		
batch_normalization_58 (Batch Normalization)	(None, 17, 17, 192)	576
['conv2d_58[0][0]']		
batch_normalization_59 (Batch Normalization)	(None, 17, 17, 192)	576
['conv2d_59[0][0]']		
activation_50 (Activation)	(None, 17, 17, 192)	0
['batch_normalization_50[0][0]']		
activation_53 (Activation)	(None, 17, 17, 192)	0
['batch_normalization_53[0][0]']		
activation_58 (Activation)	(None, 17, 17, 192)	0
['batch_normalization_58[0][0]']		
activation_59 (Activation)	(None, 17, 17, 192)	0
['batch_normalization_59[0][0]']		
mixed6 (Concatenate)	(None, 17, 17, 768)	0
['activation_50[0][0]', 'activation_53[0][0]', 'activation_58[0][0]', 'activation_59[0][0]']		
conv2d_64 (Conv2D)	(None, 17, 17, 192)	147456
['mixed6[0][0]']		


```

batch_normalization_64 (BatchN (None, 17, 17, 192) 576
['conv2d_64[0][0]']
ormalization)

activation_64 (Activation) (None, 17, 17, 192) 0
['batch_normalization_64[0][0]']

conv2d_65 (Conv2D) (None, 17, 17, 192) 258048
['activation_64[0][0]']

batch_normalization_65 (BatchN (None, 17, 17, 192) 576
['conv2d_65[0][0]']
ormalization)

activation_65 (Activation) (None, 17, 17, 192) 0
['batch_normalization_65[0][0]']

conv2d_61 (Conv2D) (None, 17, 17, 192) 147456
['mixed6[0][0]']

conv2d_66 (Conv2D) (None, 17, 17, 192) 258048
['activation_65[0][0]']

batch_normalization_61 (BatchN (None, 17, 17, 192) 576
['conv2d_61[0][0]']
ormalization)

batch_normalization_66 (BatchN (None, 17, 17, 192) 576
['conv2d_66[0][0]']
ormalization)

activation_61 (Activation) (None, 17, 17, 192) 0
['batch_normalization_61[0][0]']

activation_66 (Activation) (None, 17, 17, 192) 0
['batch_normalization_66[0][0]']

conv2d_62 (Conv2D) (None, 17, 17, 192) 258048
['activation_61[0][0]']

conv2d_67 (Conv2D) (None, 17, 17, 192) 258048
['activation_66[0][0]']

batch_normalization_62 (BatchN (None, 17, 17, 192) 576
['conv2d_62[0][0]']
ormalization)

```

```

batch_normalization_67 (BatchN (None, 17, 17, 192) 576
['conv2d_67[0][0]']
ormalization)

activation_62 (Activation) (None, 17, 17, 192) 0
['batch_normalization_62[0][0]']

activation_67 (Activation) (None, 17, 17, 192) 0
['batch_normalization_67[0][0]']

average_pooling2d_6 (AveragePo (None, 17, 17, 768) 0
['mixed6[0][0]']
oling2D)

conv2d_60 (Conv2D) (None, 17, 17, 192) 147456
['mixed6[0][0]']

conv2d_63 (Conv2D) (None, 17, 17, 192) 258048
['activation_62[0][0]']

conv2d_68 (Conv2D) (None, 17, 17, 192) 258048
['activation_67[0][0]']

conv2d_69 (Conv2D) (None, 17, 17, 192) 147456
['average_pooling2d_6[0][0]']

batch_normalization_60 (BatchN (None, 17, 17, 192) 576
['conv2d_60[0][0]']
ormalization)

batch_normalization_63 (BatchN (None, 17, 17, 192) 576
['conv2d_63[0][0]']
ormalization)

batch_normalization_68 (BatchN (None, 17, 17, 192) 576
['conv2d_68[0][0]']
ormalization)

batch_normalization_69 (BatchN (None, 17, 17, 192) 576
['conv2d_69[0][0]']
ormalization)

activation_60 (Activation) (None, 17, 17, 192) 0
['batch_normalization_60[0][0]']

activation_63 (Activation) (None, 17, 17, 192) 0
['batch_normalization_63[0][0]']

```

```

activation_68 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_68[0][0]']

activation_69 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_69[0][0]']

mixed7 (Concatenate)           (None, 17, 17, 768)  0
['activation_60[0][0]',
'activation_63[0][0]',
'activation_68[0][0]',
'activation_69[0][0]']

conv2d_72 (Conv2D)              (None, 17, 17, 192)  147456
['mixed7[0][0]']

batch_normalization_72 (BatchN (None, 17, 17, 192)  576
['conv2d_72[0][0]']
ormalization)

activation_72 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_72[0][0]']

conv2d_73 (Conv2D)              (None, 17, 17, 192)  258048
['activation_72[0][0]']

batch_normalization_73 (BatchN (None, 17, 17, 192)  576
['conv2d_73[0][0]']
ormalization)

activation_73 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_73[0][0]']

conv2d_70 (Conv2D)              (None, 17, 17, 192)  147456
['mixed7[0][0]']

conv2d_74 (Conv2D)              (None, 17, 17, 192)  258048
['activation_73[0][0]']

batch_normalization_70 (BatchN (None, 17, 17, 192)  576
['conv2d_70[0][0]']
ormalization)

batch_normalization_74 (BatchN (None, 17, 17, 192)  576
['conv2d_74[0][0]']
ormalization)

activation_70 (Activation)      (None, 17, 17, 192)  0
['batch_normalization_70[0][0]']

```

activation_74 (Activation)	(None, 17, 17, 192)	0
['batch_normalization_74[0][0]']		
conv2d_71 (Conv2D)	(None, 8, 8, 320)	552960
['activation_70[0][0]']		
conv2d_75 (Conv2D)	(None, 8, 8, 192)	331776
['activation_74[0][0]']		
batch_normalization_71 (BatchN	(None, 8, 8, 320)	960
ormalization)	['conv2d_71[0][0]']	
batch_normalization_75 (BatchN	(None, 8, 8, 192)	576
ormalization)	['conv2d_75[0][0]']	
activation_71 (Activation)	(None, 8, 8, 320)	0
['batch_normalization_71[0][0]']		
activation_75 (Activation)	(None, 8, 8, 192)	0
['batch_normalization_75[0][0]']		
max_pooling2d_3 (MaxPooling2D)	(None, 8, 8, 768)	0
['mixed7[0][0]']		
mixed8 (Concatenate)	(None, 8, 8, 1280)	0
['activation_71[0][0]',		
'activation_75[0][0]',		
'max_pooling2d_3[0][0]']		
conv2d_80 (Conv2D)	(None, 8, 8, 448)	573440
['mixed8[0][0]']		
batch_normalization_80 (BatchN	(None, 8, 8, 448)	1344
ormalization)	['conv2d_80[0][0]']	
activation_80 (Activation)	(None, 8, 8, 448)	0
['batch_normalization_80[0][0]']		
conv2d_77 (Conv2D)	(None, 8, 8, 384)	491520
['mixed8[0][0]']		
conv2d_81 (Conv2D)	(None, 8, 8, 384)	1548288
['activation_80[0][0]']		

batch_normalization_77 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_77[0][0]']		
batch_normalization_81 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_81[0][0]']		
activation_77 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_77[0][0]']		
activation_81 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_81[0][0]']		
conv2d_78 (Conv2D)	(None, 8, 8, 384)	442368
['activation_77[0][0]']		
conv2d_79 (Conv2D)	(None, 8, 8, 384)	442368
['activation_77[0][0]']		
conv2d_82 (Conv2D)	(None, 8, 8, 384)	442368
['activation_81[0][0]']		
conv2d_83 (Conv2D)	(None, 8, 8, 384)	442368
['activation_81[0][0]']		
average_pooling2d_7 (Average Pooling2D)	(None, 8, 8, 1280)	0
['mixed8[0][0]']		
conv2d_76 (Conv2D)	(None, 8, 8, 320)	409600
['mixed8[0][0]']		
batch_normalization_78 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_78[0][0]']		
batch_normalization_79 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_79[0][0]']		
batch_normalization_82 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_82[0][0]']		
batch_normalization_83 (Batch Normalization)	(None, 8, 8, 384)	1152
['conv2d_83[0][0]']		

conv2d_84 (Conv2D)	(None, 8, 8, 192)	245760
['average_pooling2d_7[0][0]']		
batch_normalization_76 (Batch Normalization)	(None, 8, 8, 320)	960
['conv2d_76[0][0]']		
activation_78 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_76[0][0]']		
activation_79 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_76[0][0]']		
activation_82 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_76[0][0]']		
activation_83 (Activation)	(None, 8, 8, 384)	0
['batch_normalization_76[0][0]']		
batch_normalization_84 (Batch Normalization)	(None, 8, 8, 192)	576
['conv2d_84[0][0]']		
activation_76 (Activation)	(None, 8, 8, 320)	0
['batch_normalization_84[0][0]']		
mixed9_0 (Concatenate)	(None, 8, 8, 768)	0
['activation_78[0][0]', 'activation_79[0][0]']		
concatenate (Concatenate)	(None, 8, 8, 768)	0
['activation_82[0][0]', 'activation_83[0][0]']		
activation_84 (Activation)	(None, 8, 8, 192)	0
['batch_normalization_84[0][0]']		
mixed9 (Concatenate)	(None, 8, 8, 2048)	0
['activation_76[0][0]', 'mixed9_0[0][0]', 'concatenate[0][0]', 'activation_84[0][0]']		
conv2d_89 (Conv2D)	(None, 8, 8, 448)	917504
['mixed9[0][0]']		
batch_normalization_89 (Batch Normalization)	(None, 8, 8, 448)	1344

```

['conv2d_89[0][0]']
ormalization)

activation_89 (Activation)      (None, 8, 8, 448)      0
['batch_normalization_89[0][0]']

conv2d_86 (Conv2D)              (None, 8, 8, 384)      786432
['mixed9[0][0]']

conv2d_90 (Conv2D)              (None, 8, 8, 384)      1548288
['activation_89[0][0]']

batch_normalization_86 (BatchN (None, 8, 8, 384)      1152
['conv2d_86[0][0]']
ormalization)

batch_normalization_90 (BatchN (None, 8, 8, 384)      1152
['conv2d_90[0][0]']
ormalization)

activation_86 (Activation)      (None, 8, 8, 384)      0
['batch_normalization_86[0][0]']

activation_90 (Activation)      (None, 8, 8, 384)      0
['batch_normalization_90[0][0]']

conv2d_87 (Conv2D)              (None, 8, 8, 384)      442368
['activation_86[0][0]']

conv2d_88 (Conv2D)              (None, 8, 8, 384)      442368
['activation_86[0][0]']

conv2d_91 (Conv2D)              (None, 8, 8, 384)      442368
['activation_90[0][0]']

conv2d_92 (Conv2D)              (None, 8, 8, 384)      442368
['activation_90[0][0]']

average_pooling2d_8 (AveragePo (None, 8, 8, 2048)      0
['mixed9[0][0]']
oling2D)

conv2d_85 (Conv2D)              (None, 8, 8, 320)      655360
['mixed9[0][0]']

batch_normalization_87 (BatchN (None, 8, 8, 384)      1152
['conv2d_87[0][0]']
ormalization)

```


batch_normalization_88 (Batch Normalization)	(None, 8, 8, 384)	1152
batch_normalization_91 (Batch Normalization)	(None, 8, 8, 384)	1152
batch_normalization_92 (Batch Normalization)	(None, 8, 8, 384)	1152
conv2d_93 (Conv2D)	(None, 8, 8, 192)	393216
batch_normalization_85 (Batch Normalization)	(None, 8, 8, 320)	960
activation_87 (Activation)	(None, 8, 8, 384)	0
activation_88 (Activation)	(None, 8, 8, 384)	0
activation_91 (Activation)	(None, 8, 8, 384)	0
activation_92 (Activation)	(None, 8, 8, 384)	0
batch_normalization_93 (Batch Normalization)	(None, 8, 8, 192)	576
activation_85 (Activation)	(None, 8, 8, 320)	0
mixed9_1 (Concatenate)	(None, 8, 8, 768)	0
concatenate_1 (Concatenate)	(None, 8, 8, 768)	0
activation_93 (Activation)	(None, 8, 8, 192)	0

```
['batch_normalization_93[0][0]']

mixed10 (Concatenate)          (None, 8, 8, 2048)    0
['activation_85[0][0]',
'mixed9_1[0][0]',
'concatenate_1[0][0]',
'activation_93[0][0]']

avg_pool (GlobalAveragePooling (None, 2048)    0
['mixed10[0][0]']
2D)

predictions (Dense)            (None, 1000)          2049000
['avg_pool[0][0]']
```

```
=====
=====
Total params: 23,851,784
Trainable params: 23,817,352
Non-trainable params: 34,432
-----
-----
```

This is a prediction model, so the output is typically a softmax-activated vector representing 1000 possible object types. Because we are interested in an encoded representation of the image we are just going to use the second-to-last layer as a source of image encodings. Each image will be encoded as a vector of size 2048.

We will use the following hack: hook up the input into a new Keras model and use the penultimate layer of the existing model as output.

```
[19]: new_input = img_model.input
      new_output = img_model.layers[-2].output
      img_encoder = Model(new_input, new_output) # This is the final Keras image_
      ↪ encoder model we will use.
```

Let's try the encoder. At this point, you may want to add a GPU to the VM you are using (if not using already).

```
[22]: encoded_image = img_encoder.predict(np.array([new_image]))
```

```
[21]: encoded_image
```

```
[21]: array([[0.63806653, 0.48873004, 0.05526242, ..., 0.64255667, 0.29595238,
              0.49004298]], dtype=float32)
```

TODO: We will need to create encodings for all images and store them in one big matrix (one for each dataset, train, dev, test). We can then save the matrices so that we never have to touch the bulky image data again.

To save memory (but slow the process down a little bit) we will read in the images lazily using a generator. We will encounter generators again later when we train the LSTM. If you are unfamiliar with generators, take a look at this page: <https://wiki.python.org/moin/Generators>

Write the following generator function, which should return one image at a time. `img_list` is a list of image file names (i.e. the train, dev, or test set). The return value should be a numpy array of shape (1,299,299,3).

```
[23]: def img_generator(img_list):
      for i in img_list:
          image = get_image(i)
          image_sh = image.reshape(1,299,299,3)
          yield image_sh
```

Now we can encode all images (this takes a few minutes).

```
[24]: enc_train = img_encoder.predict_generator(img_generator(train_list),
      ↪steps=len(train_list), verbose=1)
```

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:1: UserWarning:
`Model.predict_generator` is deprecated and will be removed in a future version.
Please use `Model.predict`, which supports generators.
```

```
"""Entry point for launching an IPython kernel.
```

```
6000/6000 [=====] - 176s 29ms/step
```

```
[25]: enc_train[11]
```

```
[25]: array([0.26818588, 1.0321662 , 0.5851617 , ..., 1.231674 , 0.1796931 ,
      0.22405314], dtype=float32)
```

```
[26]: enc_dev = img_encoder.predict_generator(img_generator(dev_list),
      ↪steps=len(dev_list), verbose=1)
```

```
5/1000 [...] - ETA: 27s
```

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:1: UserWarning:
`Model.predict_generator` is deprecated and will be removed in a future version.
Please use `Model.predict`, which supports generators.
```

```
"""Entry point for launching an IPython kernel.
```

```
1000/1000 [=====] - 29s 29ms/step
```

```
[27]: enc_test = img_encoder.predict_generator(img_generator(test_list),
      ↪steps=len(test_list), verbose=1)
```

```
5/1000 [...] - ETA: 28s
```

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:1: UserWarning:
`Model.predict_generator` is deprecated and will be removed in a future version.
```

Please use `Model.predict``, which supports generators.

```
"""Entry point for launching an IPython kernel.
```

```
1000/1000 [=====] - 29s 29ms/step
```

It's a good idea to save the resulting matrices, so we do not have to run the encoder again.

```
[28]: # Choose a suitable location here, please do NOT attempt to write your output_
      ↪ files to the shared data directory.
      OUTPUT_PATH = "hw5output"
      if not os.path.exists(OUTPUT_PATH):
          os.mkdir(OUTPUT_PATH)
```

```
[29]: np.save(os.path.join(OUTPUT_PATH, "encoded_images_train.npy"), enc_train)
      np.save(os.path.join(OUTPUT_PATH, "encoded_images_dev.npy"), enc_dev)
      np.save(os.path.join(OUTPUT_PATH, "encoded_images_test.npy"), enc_test)
```

1.3 Part II Text (Caption) Data Preparation (14 pts)

Next, we need to load the image captions and generate training data for the generator model.

1.3.1 Reading image descriptions

TODO: Write the following function that reads the image descriptions from the file `filename` and returns a dictionary in the following format. Take a look at the file `Flickr8k.token.txt` for the format of the input file. The keys of the dictionary should be image filenames. Each value should be a list of 5 captions. Each caption should be a list of tokens.

The captions in the file are already tokenized, so you can just split them at white spaces. You should convert each token to lower case. You should then pad each caption with a `START` token on the left and an `END` token on the right.

```
[199]: def read_image_descriptions(filename):
      image_descriptions = defaultdict(list)
      file = open(filename, "r")
      for line in file:
          token = line.strip().split("\t")
          name_img = token[0].split('#')[0] # key value of the dict
          des_img = token[1] #description of the image
          des_split = ['<START>'] + list(map(lambda p: p.lower(), des_img.split(' '
          ↪ ')))) + ['<END>']
          if name_img not in image_descriptions:
              image_descriptions[name_img] = list()
          image_descriptions[name_img].append(des_split)
      return image_descriptions
```

```
[200]: descriptions = read_image_descriptions(f"{FLICKR_PATH}/Flickr8k.token.txt")
```

```
[201]: print(descriptions[dev_list[0]])
```

```
[['<START>', 'the', 'boy', 'laying', 'face', 'down', 'on', 'a', 'skateboard',
'is', 'being', 'pushed', 'along', 'the', 'ground', 'by', 'another', 'boy', '.',
'<END>'], ['<START>', 'two', 'girls', 'play', 'on', 'a', 'skateboard', 'in',
'a', 'courtyard', '.', '<END>'], ['<START>', 'two', 'people', 'play', 'on', 'a',
'long', 'skateboard', '.', '<END>'], ['<START>', 'two', 'small', 'children',
'in', 'red', 'shirts', 'playing', 'on', 'a', 'skateboard', '.', '<END>'],
['<START>', 'two', 'young', 'children', 'on', 'a', 'skateboard', 'going',
'across', 'a', 'sidewalk', '<END>']]
```

Running the previous cell should print:

```
[['<START>', 'the', 'boy', 'laying', 'face', 'down', 'on', 'a', 'skateboard',
'is', 'being', 'pushed', 'along', 'the', 'ground', 'by', 'another', 'boy',
'.', '<END>'], ['<START>', 'two', 'girls', 'play', 'on', 'a', 'skateboard',
'in', 'a', 'courtyard', '.', '<END>'], ['<START>', 'two', 'people', 'play',
'on', 'a', 'long', 'skateboard', '.', '<END>'], ['<START>', 'two', 'small',
'children', 'in', 'red', 'shirts', 'playing', 'on', 'a', 'skateboard', '.',
'<END>'], ['<START>', 'two', 'young', 'children', 'on', 'a', 'skateboard',
'going', 'across', 'a', 'sidewalk', '<END>']]
```

1.3.2 Creating Word Indices

Next, we need to create a lookup table from the **training** data mapping words to integer indices, so we can encode input and output sequences using numeric representations. **TODO** create the dictionaries `id_to_word` and `word_to_id`, which should map tokens to numeric ids and numeric ids to tokens.

Hint: Create a set of tokens in the training data first, then convert the set into a list and sort it. This way if you run the code multiple times, you will always get the same dictionaries.

```
[33]: set_tokens = set()
id_to_word = defaultdict(int)
word_to_id = defaultdict(int)
for name_img in descriptions:
    for description in descriptions[name_img]:
        for token in description:
            set_tokens.add(token)
list_tokens = list(set_tokens) #create set of tokens in training data
list_tokens.sort() #sort it
l=len(list_tokens)
for i in range(l):
    id_to_word[i] = list_tokens[i]
m = len(id_to_word)
for i in range(m):
    word_to_id[id_to_word[i]] = i
```

```
[34]: word_to_id['dog'] # should print an integer
```

```
[34]: 2309
```

```
[35]: id_to_word[2309]
```

```
[35]: 'dog'
```

```
[36]: id_to_word[1985] # should print a token
```

```
[36]: 'crucified'
```

Note that we do not need an UNK word token because we are generating. The generated text will only contain tokens seen at training time.

1.4 Part III Basic Decoder Model (24 pts)

For now, we will just train a model for text generation without conditioning the generator on the image input.

There are different ways to do this and our approach will be slightly different from the generator discussed in class.

The core idea here is that the Keras recurrent layers (including LSTM) create an “unrolled” RNN. Each time-step is represented as a different unit, but the weights for these units are shared. We are going to use the constant MAX_LEN to refer to the maximum length of a sequence, which turns out to be 40 words in this data set (including START and END).

```
[37]: max(len(description) for image_id in train_list for description in_
      ↪ descriptions[image_id])
```

```
[37]: 40
```

In class, we discussed LSTM generators as transducers that map each word in the input sequence to the next word.

Instead, we will use the model to predict one word at a time, given a partial sequence. For example, given the sequence [“START”, “a”], the model might predict “dog” as the most likely word. We are basically using the LSTM to encode the input sequence up to this point.

To train the model, we will convert each description into a set of input output pairs as follows. For example, consider the sequence

```
['<START>', 'a', 'black', 'dog', '.', '<END>']
```

We would train the model using the following input/output pairs

i	input	output
0	[START]	a
1	[START,a]	black
2	[START,a, black]	dog
3	[START,a, black, dog]	END

Here is the model in Keras Keras. Note that we are using a Bidirectional LSTM, which encodes the

sequence from both directions and then predicts the output. Also note the `return_sequence=False` parameter, which causes the LSTM to return a single output instead of one output per state.

Note also that we use an embedding layer for the input words. The weights are shared between all units of the unrolled LSTM. We will train these embeddings with the model.

```
[181]: MAX_LEN = 40
EMBEDDING_DIM=300
vocab_size = len(word_to_id)

# Text input
text_input = Input(shape=(MAX_LEN,))
embedding = Embedding(vocab_size, EMBEDDING_DIM,
    ↪input_length=MAX_LEN)(text_input)
x = Bidirectional(LSTM(512, return_sequences=False))(embedding)
pred = Dense(vocab_size, activation='softmax')(x)
model = Model(inputs=[text_input], outputs=pred)
model.compile(loss='categorical_crossentropy', optimizer='RMSprop',
    ↪metrics=['accuracy'])

model.summary()
```

Model: "model_9"

Layer (type)	Output Shape	Param #
input_14 (InputLayer)	[(None, 40)]	0
embedding_8 (Embedding)	(None, 40, 300)	2676300
bidirectional_8 (Bidirectional)	(None, 1024)	3330048
dense_11 (Dense)	(None, 8921)	9144025

=====
 Total params: 15,150,373
 Trainable params: 15,150,373
 Non-trainable params: 0
 =====

The model input is a numpy ndarray (a tensor) of size `(batch_size, MAX_LEN)`. Each row is a vector of size `MAX_LEN` in which each entry is an integer representing a word (according to the `word_to_id` dictionary). If the input sequence is shorter than `MAX_LEN`, the remaining entries should be padded with 0.

For each input example, the model returns a softmax activated vector (a probability distribution) over possible output words. The model output is a numpy ndarray of size `(batch_size, vocab_size)`. `vocab_size` is the number of vocabulary words.

1.4.1 Creating a Generator for the Training Data

TODO:

We could simply create one large numpy ndarray for all the training data. Because we have a lot of training instances (each training sentence will produce up to MAX_LEN input/output pairs, one for each word), it is better to produce the training examples *lazily*, i.e. in batches using a generator (recall the image generator in part I).

Write the function `text_training_generator` below, that takes as a parameter the `batch_size` and returns an (input, output) pair. `input` is a (`batch_size`, MAX_LEN) ndarray of partial input sequences, `output` contains the next words predicted for each partial input sequence, encoded as a (`batch_size`, `vocab_size`) ndarray.

Each time the `next()` function is called on the generator instance, it should return a new batch of the *training* data. You can use `train_list` as a list of training images. A batch may contain input/output examples extracted from different descriptions or even from different images.

You can just refer back to the variables you have defined above, including `descriptions`, `train_list`, `vocab_size`, etc.

Hint: To prevent issues with having to reset the generator for each epoch and to make sure the generator can always return exactly `batch_size` input/output pairs in each step, wrap your code into a `while True:` loop. This way, when you reach the end of the training data, you will just continue adding training data from the beginning into the batch.

```
[205]: def text_training_generator(batch_size=128):
    input = []
    output = []
    num = 0
    while True:
        for i in train_list:
            for j in descriptions[i]:
                sentence = np.full(MAX_LEN, word_to_id[' '])
                sentence[0] = word_to_id[j[0]]
                for k in range(1, len(j)):
                    count = word_to_id[j[k]]
                    input.append(sentence.copy())
                    output.append(to_categorical(count, vocab_size))
                num += 1
            if num >= batch_size:
                input = np.array(input)
                output = np.array(output)
                yield(input, output)
                input = []
                output = []
                num = 0
            sentence[k] = word_to_id[j[k]]
```

1.4.2 Training the Model

We will use the `fit_generator` method of the model to train the model. `fit_generator` needs to know how many iterator steps there are per epoch.

Because there are `len(train_list)` training samples with up to `MAX_LEN` words, an upper bound for the number of total training instances is `len(train_list)*MAX_LEN`. Because the generator returns these in batches, the number of steps is `len(train_list) * MAX_LEN // batch_size`

```
[206]: batch_size = 128
generator = text_training_generator(batch_size)
steps = len(train_list) * MAX_LEN // batch_size
```

```
[207]: model.fit_generator(generator, steps_per_epoch=steps, verbose=True, epochs=10)
```

Epoch 1/10

```
/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:1: UserWarning:
`Model.fit_generator` is deprecated and will be removed in a future version.
Please use `Model.fit`, which supports generators.
```

```
"""Entry point for launching an IPython kernel.
```

```
1875/1875 [=====] - 288s 142ms/step - loss: 4.2645 -
accuracy: 0.2953
```

Epoch 2/10

```
1875/1875 [=====] - 266s 142ms/step - loss: 3.6967 -
accuracy: 0.3582
```

Epoch 3/10

```
1875/1875 [=====] - 266s 142ms/step - loss: 3.5298 -
accuracy: 0.3734
```

Epoch 4/10

```
1875/1875 [=====] - 265s 141ms/step - loss: 3.4276 -
accuracy: 0.3831
```

Epoch 5/10

```
1875/1875 [=====] - 265s 141ms/step - loss: 3.3752 -
accuracy: 0.3895
```

Epoch 6/10

```
1875/1875 [=====] - 265s 141ms/step - loss: 3.3106 -
accuracy: 0.3955
```

Epoch 7/10

```
1875/1875 [=====] - 265s 142ms/step - loss: 3.2983 -
accuracy: 0.3978
```

Epoch 8/10

```
1875/1875 [=====] - 265s 141ms/step - loss: 3.2928 -
accuracy: 0.3999
```

Epoch 9/10

```
1875/1875 [=====] - 265s 142ms/step - loss: 3.2770 -
accuracy: 0.4027
```

Epoch 10/10

```
1875/1875 [=====] - 265s 141ms/step - loss: 3.4336 -
```

accuracy: 0.4002

[207]: <keras.callbacks.History at 0x7ff0edc855d0>

Continue to train the model until you reach an accuracy of at least 40%.

1.4.3 Greedy Decoder

TODO Next, you will write a decoder. The decoder should start with the sequence ["<START>"], use the model to predict the most likely word, append the word to the sequence and then continue until "<END>" is predicted or the sequence reaches MAX_LEN words.

```
[208]: def decoder():
        input = np.zeros((1, MAX_LEN))
        input[0,0]= word_to_id["<START>"]
        count = 1
        output = ["<START>"]
        while count < MAX_LEN:
            predict = np.argmax(model.predict(input))
            output.append(id_to_word[predict])
            input[0, len(output)-1] = word_to_id[output[-1]]
            count+=1
            if predict == word_to_id["<END>"]:
                break
        return output
```

```
[209]: print(decoder())
```

```
['<START>', 'a', 'man', 'and', 'woman', 'are', 'standing', 'in', 'front', 'of',
'a', 'building', '.', '<END>']
```

This simple decoder will of course always predict the same sequence (and it's not necessarily a good one).

Modify the decoder as follows. Instead of choosing the most likely word in each step, sample the next word from the distribution (i.e. the softmax activated output) returned by the model. Take a look at the `np.random.multinomial` function to do this.

```
[212]: def sample_decoder():
        input = np.zeros((1, MAX_LEN))
        input[0,0]= word_to_id["<START>"]
        count = 1
        output = ["<START>"]
        while count < MAX_LEN:
            predict = model.predict(input)[0].astype('float64')
            predict /= sum(predict)
            predict_final = np.argmax(np.random.multinomial(1, predict))
            output.append(id_to_word[predict_final])
            input[0, len(output)-1] = word_to_id[output[-1]]
```

```

        count +=1
        if predict_final == word_to_id("<END>"):
            break
    return output

```

You should now be able to see some interesting output that looks a lot like flickr8k image captions – only that the captions are generated randomly without any image input.

```

[214]: for i in range(10):
        print(sample_decoder())

```

```

['<START>', 'a', 'person', 'is', 'standing', 'in', 'a', 'small', 'blue', ',', 'yellow', 'and', 'blue', 'crowd', 'in', 'matching', '.', '<END>']
['<START>', 'people', 'stand', 'on', 'police', 'men', '.', '<END>']
['<START>', 'a', 'child', 'bared', 'with', 'a', 'black', 'shirt', '<END>']
['<START>', 'a', 'boy', 'sitting', 'in', 'a', 'very', 'field', 'surrounded', 'by', 'people', 'standing', 'in', 'front', 'skateboarding', '.', '<END>']
['<START>', 'several', 'people', 'in', 'helmet', 'are', 'inline', 'doing', 'a', 'handstand', 'on', 'a', 'rail', '.', '<END>']
['<START>', 'a', 'person', 'is', 'attempting', 'to', 'wooden', 'caught', '.', '<END>']
['<START>', 'a', 'young', 'child', 'is', 'sliding', 'down', 'a', 'pole', 'in', 'a', 'green', 'outfit', '.', '<END>']
['<START>', 'a', 'young', 'child', 'and', 'a', 'brown', 'dog', 'sitting', 'on', 'a', 'rocky', 'road', '.', '<END>']
['<START>', 'two', 'little', 'girls', 'are', 'vampires', 'in', 'a', 'red', 'rough', 'field', '.', '<END>']
['<START>', 'skier', 'playing', 'person', 'black', 'dogs', '.', '<END>']

```

1.5 Part III - Conditioning on the Image (24 pts)

We will now extend the model to condition the next word not only on the partial sequence, but also on the encoded image.

We will project the 2048-dimensional image encoding to a 300-dimensional hidden layer. We then concatenate this vector with each embedded input word, before applying the LSTM.

Here is what the Keras model looks like:

```

[216]: MAX_LEN = 40
        EMBEDDING_DIM=300
        IMAGE_ENC_DIM=300

        # Image input
        img_input = Input(shape=(2048,))
        img_enc = Dense(300, activation="relu")(img_input)
        images = RepeatVector(MAX_LEN)(img_enc)

        # Text input

```

```

text_input = Input(shape=(MAX_LEN,))
embedding = Embedding(vocab_size, EMBEDDING_DIM,
    ↪input_length=MAX_LEN)(text_input)
x = Concatenate()([images, embedding])
y = Bidirectional(LSTM(256, return_sequences=False))(x)
pred = Dense(vocab_size, activation='softmax')(y)
model = Model(inputs=[img_input, text_input], outputs=pred)
model.compile(loss='categorical_crossentropy', optimizer="RMSProp",
    ↪metrics=['accuracy'])

model.summary()

```

Model: "model_11"

```

-----
Layer (type)                 Output Shape              Param #   Connected to
=====
input_17 (InputLayer)        [(None, 2048)]            0         []
dense_14 (Dense)              (None, 300)               614700    ['input_17[0][0]']
input_18 (InputLayer)        [(None, 40)]              0         []
repeat_vector_4 (RepeatVector) (None, 40, 300)           0         ['dense_14[0][0]']
embedding_10 (Embedding)     (None, 40, 300)          2676300   ['input_18[0][0]']
concatenate_8 (Concatenate)  (None, 40, 600)           0         ['repeat_vector_4[0][0]',
'embedding_10[0][0]']
bidirectional_10 (Bidirection (None, 512)              1755136   ['concatenate_8[0][0]']
1)
dense_15 (Dense)              (None, 8921)              4576473   ['bidirectional_10[0][0]']
=====
Total params: 9,622,609
Trainable params: 9,622,609
Non-trainable params: 0

```


The model now takes two inputs:

1. a (batch_size, 2048) ndarray of image encodings.
2. a (batch_size, MAX_LEN) ndarray of partial input sequences.

And one output as before: a (batch_size, vocab_size) ndarray of predicted word distributions.

TODO: Modify the training data generator to include the image with each input/output pair. Your generator needs to return an object of the following format: ([image_inputs, text_inputs], next_words). Where each element is an ndarray of the type described above.

You need to find the image encoding that belongs to each image. You can use the fact that the index of the image in train_list is the same as the index in enc_train and enc_dev.

If you have previously saved the image encodings, you can load them from disk:

```
[217]: enc_train = np.load(f"{OUTPUT_PATH}/encoded_images_train.npy")
enc_dev = np.load(f"{OUTPUT_PATH}/encoded_images_dev.npy")
```

```
[221]: def training_generator(batch_size=128):
    input_text = []
    output = []
    input_img = []
    num = 0
    l = len(train_list)
    while True:
        for i in range(l):
            for j in descriptions[train_list[i]]:
                sentence = np.full(MAX_LEN, word_to_id[' '])
                sentence[0] = word_to_id[j[0]]
                for k in range(1, len(j)):
                    count = word_to_id[j[k]]
                    input_text.append(sentence.copy())
                    output.append(to_categorical(count, vocab_size))
                    input_img.append(enc_train[i])
                num += 1
            if num >= batch_size:
                input_text_final = np.array(input_text)
                output_final = np.array(output)
                input_img_final = np.array(input_img)
                yield([input_img_final, input_text_final], output_final)
                input_text = []
                output = []
                input_img = []
                num = 0
            sentence[k] = word_to_id[j[k]]
```

You should now be able to train the model as before:

```
[222]: batch_size = 128
generator = training_generator(batch_size)
steps = len(train_list) * MAX_LEN // batch_size
```

```
[223]: model.fit_generator(generator, steps_per_epoch=steps, verbose=True, epochs=20)
```

Epoch 1/20

/opt/conda/lib/python3.7/site-packages/ipykernel_launcher.py:1: UserWarning:
`Model.fit_generator` is deprecated and will be removed in a future version.
Please use `Model.fit`, which supports generators.

"""Entry point for launching an IPython kernel.

1875/1875 [=====] - 173s 90ms/step - loss: 4.4433 -
accuracy: 0.2645

Epoch 2/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.6341 -
accuracy: 0.3665

Epoch 3/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.4374 -
accuracy: 0.3872

Epoch 4/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.3211 -
accuracy: 0.3971

Epoch 5/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.2760 -
accuracy: 0.4033

Epoch 6/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.2160 -
accuracy: 0.4088

Epoch 7/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.2080 -
accuracy: 0.4102

Epoch 8/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.1932 -
accuracy: 0.4147

Epoch 9/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.1957 -
accuracy: 0.4160

Epoch 10/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.2064 -
accuracy: 0.4170

Epoch 11/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.1999 -
accuracy: 0.4190

Epoch 12/20

1875/1875 [=====] - 169s 90ms/step - loss: 3.2214 -
accuracy: 0.4185

```

Epoch 13/20
1875/1875 [=====] - 170s 91ms/step - loss: 3.2258 -
accuracy: 0.4195
Epoch 14/20
1875/1875 [=====] - 170s 90ms/step - loss: 3.1899 -
accuracy: 0.4217
Epoch 15/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.1983 -
accuracy: 0.4210
Epoch 16/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.2114 -
accuracy: 0.4224
Epoch 17/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.1952 -
accuracy: 0.4249
Epoch 18/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.2110 -
accuracy: 0.4251
Epoch 19/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.2004 -
accuracy: 0.4249
Epoch 20/20
1875/1875 [=====] - 169s 90ms/step - loss: 3.2393 -
accuracy: 0.4253

```

[223]: <keras.callbacks.History at 0x7ff044dfef50>

Again, continue to train the model until you hit an accuracy of about 40%. This may take a while. I strongly encourage you to experiment with cloud GPUs using the GCP voucher for the class.

You can save your model weights to disk and continue at a later time.

[224]: `model.save_weights(f"{OUTPUT_PATH}/model.h5")`

to load the model:

[225]: `model.load_weights(f"{OUTPUT_PATH}/model.h5")`

TODO: Now we are ready to actually generate image captions using the trained model. Modify the simple greedy decoder you wrote for the text-only generator, so that it takes an encoded image (a vector of length 2048) as input, and returns a sequence.

[226]:

```
def image_decoder(enc_image):
    num = 1
    input_text = np.zeros((1, MAX_LEN))
    input_text[0, 0] = word_to_id['<START>']
    input_img = np.asarray(enc_image).reshape((-1, 2048))
    output_id = 0
    output = ['<START>']
```



```

while num < MAX_LEN:
    predict = np.argmax(model.predict([input_img, input_text]))
    output.append(id_to_word[predict])
    num+=1
    input_text[0, len(output)-1] = word_to_id[output[-1]]
    if predict == word_to_id["<END>"]:
        break
return output

```

As a sanity check, you should now be able to reproduce (approximately) captions for the training images.

```

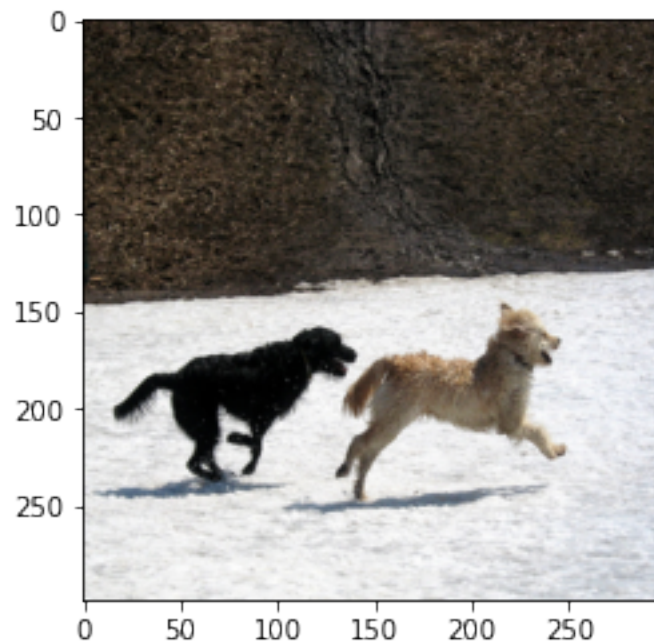
[227]: plt.imshow(get_image(train_list[0]))
       image_decoder(enc_train[0])

```

```

[227]: ['<START>',
        'a',
        'dog',
        'jumps',
        'over',
        'a',
        'and',
        'red',
        'ball',
        '.',
        '<END>']

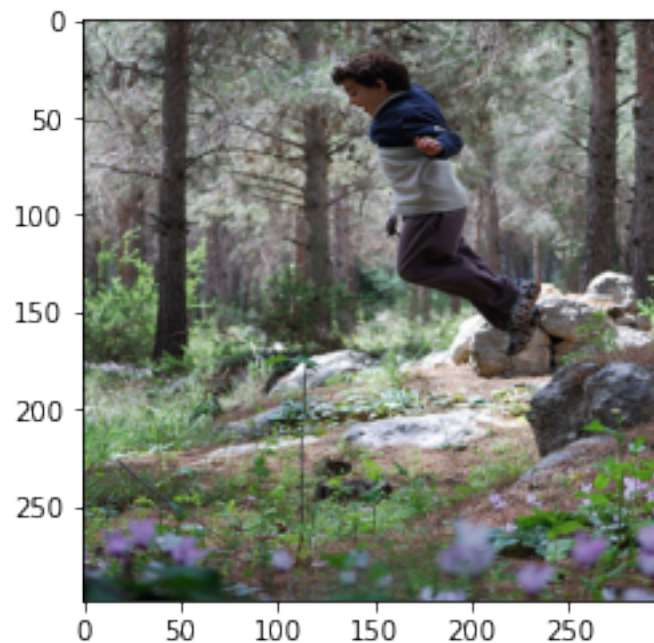
```



You should also be able to apply the model to dev images and get reasonable captions:

```
[228]: plt.imshow(get_image(dev_list[1]))  
       image_decoder(enc_dev[0])
```

```
[228]: ['<START>',  
       'a',  
       'man',  
       'in',  
       'a',  
       'grey',  
       'shirt',  
       'and',  
       'a',  
       'grey',  
       'shirt',  
       'is',  
       'climbing',  
       'a',  
       'rock',  
       '.',  
       '<END>']
```



For this assignment we will not perform a formal evaluation.

Feel free to experiment with the parameters of the model or continue training the model. At some point, the model will overfit and will no longer produce good descriptions for the dev images.

1.6 Part IV - Beam Search Decoder (24 pts)

TODO Modify the simple greedy decoder for the caption generator to use beam search. Instead of always selecting the most probable word, use a *beam*, which contains the n highest-scoring sequences so far and their total probability (i.e. the product of all word probabilities). I recommend that you use a list of (**probability**, **sequence**) tuples. After each time-step, prune the list to include only the n most probable sequences.

Then, for each sequence, compute the n most likely successor words. Append the word to produce n new sequences and compute their score. This way, you create a new list of $n*n$ candidates.

Prune this list to the best n as before and continue until **MAX_LEN** words have been generated.

Note that you cannot use the occurrence of the "<END>" tag to terminate generation, because the tag may occur in different positions for different entries in the beam.

Once **MAX_LEN** has been reached, return the most likely sequence out of the current n .

```
[ ]: def img_beam_decoder(n, image_enc):
```

```
    #...
```

```
[ ]: beam_decoder(3, dev_list[1])
```

TODO Finally, before you submit this assignment, please show 5 development images, each with 1) their greedy output, 2) beam search at $n=3$ 3) beam search at $n=5$.