

COMS W4705 - Homework 4

Image Captioning with Conditioned LSTM Generators

Yassine Benajiba yb2235@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 submit a copy of this notebook only, including all outputs. Do not submit any of the data files.

Getting Started

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

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

from keras import Sequential, Model
from keras.layers import Embedding, LSTM, Dense, Input, Bidirectional, RepeatVector, TimeDistributed
from keras.activations import softmax
from keras.utils import to_categorical
from keras.preprocessing.sequence import pad_sequences

from keras.applications.inception_v3 import InceptionV3

from keras.optimizers import Adam

from google.colab import drive
```

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>
[\(http://www.jair.org/papers/paper3994.html\)](http://www.jair.org/papers/paper3994.html) when discussing our results

I have uploaded all the data and model files you'll need to my GDrive and you can access the folder here: [\(https://drive.google.com/drive/folders/1i9lun4h3EN1vSd1A1woez0mXJ9vRjFIT?usp=sharing\)](https://drive.google.com/drive/folders/1i9lun4h3EN1vSd1A1woez0mXJ9vRjFIT?usp=sharing)

Google Drive does not allow to copy a folder, so you'll need to download the whole folder and then upload it again to your own drive. Please assign the name you chose for this folder to the variable `my_data_dir` in the next cell.

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: [\(https://forms.illinois.edu/sec/1713398\)](https://forms.illinois.edu/sec/1713398).

In [2]: `#this is where you put the name of your data folder.
#Please make sure it's correct because it'll be used in many places later.
my_data_dir="hw5_data"`

Mounting your GDrive so you can access the files from Colab

In [14]: `#running this command will generate a message that will ask you to click on
#copy paste that code in the text box that will appear below
drive.flush_and_unmount()
drive.mount('/content/gdrive')`

Mounted at /content/gdrive

Please look at the 'Files' tab on the left side and make sure you can see the 'hw5_data' folder that you have in your GDrive.

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.

In [4]: `def load_image_list(filename):
 with open(filename, 'r') as image_list_f:
 return [line.strip() for line in image_list_f]`

```
In [5]: train_list = load_image_list('/content/gdrive/My Drive/' + my_data_dir + '/Flickr8k_Dataset/training')
dev_list = load_image_list('/content/gdrive/My Drive/' + my_data_dir + '/Flickr8k_Dataset/dev')
test_list = load_image_list('/content/gdrive/My Drive/' + my_data_dir + '/Flickr8k_Dataset/test')
```

Let's see how many images there are

```
In [9]: from google.colab import drive
drive.mount('/content/drive')
```

Mounted at /content/drive

```
In [10]: len(train_list), len(dev_list), len(test_list)
```

```
Out[10]: (6000, 1000, 1000)
```

Each entry is an image filename.

```
In [11]: dev_list[20]
```

```
Out[11]: '3693961165_9d6c333d5b.jpg'
```

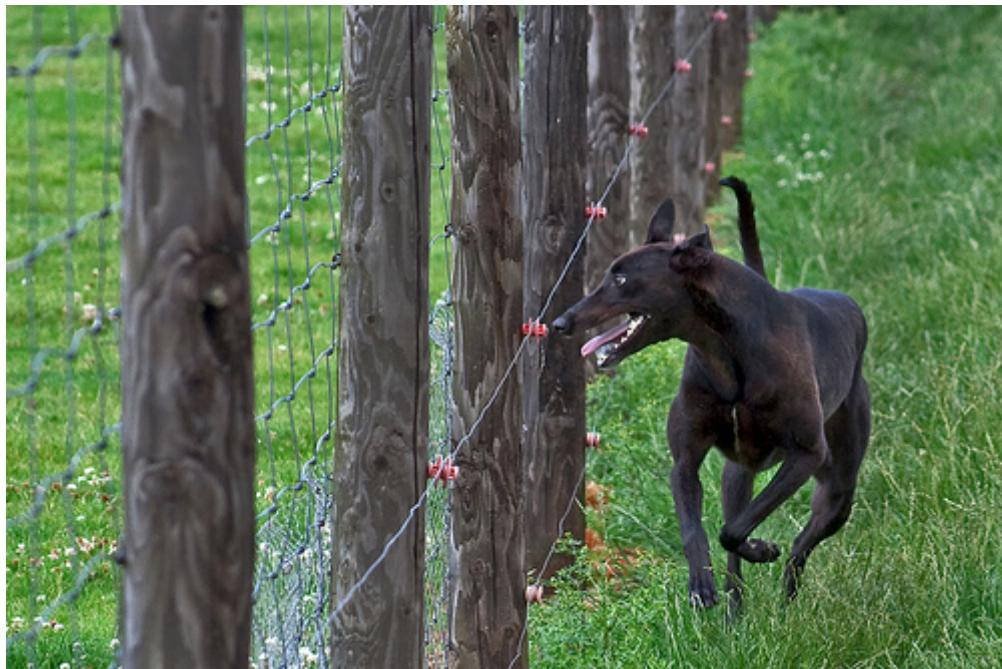
The images are located in a subdirectory.

```
In [12]: IMG_PATH = "Flickr8k_Dataset"
```

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

```
In [15]: image = PIL.Image.open(os.path.join('/content/gdrive/MyDrive/hw5_data/Flickr8k_Dataset', dev_list[0]))
```

```
Out[15]:
```



if you can't see the image, try

```
In [16]: plt.imshow(image)
```

```
Out[16]: <matplotlib.image.AxesImage at 0x7f664f11f208>
```



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_CV_\(https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Szegedy_Rethinking_the_Inception_CV.html](https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Szegedy_Rethinking_the_Inception_CV_(https://www.cv-foundation.org/openaccess/content_cvpr_2016/html/Szegedy_Rethinking_the_Inception_CV.html)

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.

```
In [17]: np.asarray(image).shape
```

```
Out[17]: (333, 500, 3)
```

The values range from 0 to 255.

```
In [18]: np.asarray(image)
```

```
Out[18]: 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]],

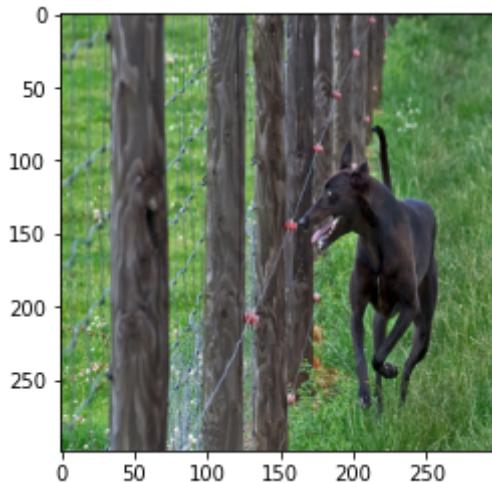
  [[100, 158, 97],
   [ 38, 100, 37],
   [ 46, 117, 51],
   ...,
   [ 25, 54,  8],
   [ 88, 112, 76],
   [ 65, 106, 48]],

  [[ 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.

```
In [19]: new_image = np.asarray(image.resize((299,299))) / 255.0  
plt.imshow(new_image)
```

```
Out[19]: <matplotlib.image.AxesImage at 0x7f664f0fdc18>
```



```
In [20]: new_image.shape
```

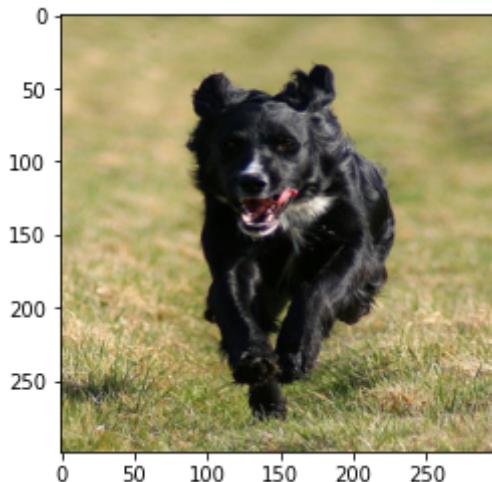
```
Out[20]: (299, 299, 3)
```

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

```
In [6]: def get_image(image_name):  
    image = PIL.Image.open(os.path.join('/content/gdrive/MyDrive/hw5_data/F  
    return np.asarray(image.resize((299,299))) / 255.0
```

```
In [8]: plt.imshow(get_image(dev_list[25]))
```

```
Out[8]: <matplotlib.image.AxesImage at 0x7f665162dbe0>
```



Next, we load the pre-trained Inception model.

```
In [21]: img_model = InceptionV3(weights='imagenet') # This will download the weight
          Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/inception_v3/inception_v3_weights_tf_dim_ordering_tf_kernels.h5
          (https://storage.googleapis.com/tensorflow/keras-applications/inception_v3/inception_v3_weights_tf_dim_ordering_tf_kernels.h5)
          96116736/96112376 [=====] - 1s 0us/step
```

```
In [22]: img_model.summary() # this is quite a complex model.
```

Model: "inception_v3"

Layer (type) d to	Output Shape	Param #	Connecte
input_1 (InputLayer)	[(None, 299, 299, 3) 0		
conv2d (Conv2D) [0][0]	(None, 149, 149, 32) 864		input_1
batch_normalization (BatchNorma [0][0]	(None, 149, 149, 32) 96		conv2d
activation (Activation)	(None, 149, 149, 32) 0		batch_no

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.

```
In [23]: 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
```

Let's try the encoder.

```
In [24]: encoded_image = img_encoder.predict(np.array([new_image]))
```

```
In [25]: encoded_image.shape
```

```
Out[25]: (1, 2048)
```

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: [\(https://wiki.python.org/moin/Generators\)](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).

```
In [26]: def img_generator(img_list):
    ...
    for img_name in img_list:
        image = get_image(img_name)
        yield np.array([image])
```

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

```
In [ ]: enc_train = img_encoder.predict_generator(img_generator(train_list), steps=6000/6000 [=====] - 5134s 856ms/step
```

```
In [ ]: enc_train[11]
```

```
Out[27]: array([0.2681858 , 1.032167 , 0.5851618 , ..., 1.2316743 , 0.17969358,
       0.2240533 ], dtype=float32)
```

```
In [ ]: enc_dev = img_encoder.predict_generator(img_generator(dev_list), steps=len(1000/1000 [=====] - 915s 915ms/step
                                                43/1000 [>.....] - ETA: 15s
```

```
In [ ]: enc_test = img_encoder.predict_generator(img_generator(test_list), steps=1000/1000 [=====] - 888s 888ms/step
```

WARNING:tensorflow:From <ipython-input-29-c4f3480df2d6>:1: Model.predict_generator (from tensorflow.python.keras.engine.training) is deprecated and will be removed in a future version.
 Instructions for updating:
 Please use Model.predict, which supports generators.

```
In [ ]: np.save("gdrive/My Drive/" + my_data_dir + "/outputs/encoded_images_train.npy",
              np.save("gdrive/My Drive/" + my_data_dir + "/outputs/encoded_images_dev.npy",
              np.save("gdrive/My Drive/" + my_data_dir + "/outputs/encoded_images_test.npy",
```

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

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

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.

```
In [27]: def read_image_descriptions(filename):
    image_descriptions = defaultdict(list)
    with open(filename, 'r') as lines:
        for line in lines:
            img_name, description = line.split('\t')
            tokens = ['<START>'] + description.lower().split() + ['<END>']
            image_descriptions[img_name.split('#')[0]].append(tokens)
    return image_descriptions
```

```
In [28]: descriptions = read_image_descriptions("gdrive/My Drive/" + my_data_dir + "/Fli
```

```
In [29]: print(descriptions[train_list[1]])
```

```
[[ '<START>', 'a', 'little', 'baby', 'plays', 'croquet', '.', '<END>'],
 ['<START>', 'a', 'little', 'girl', 'plays', 'croquet', 'next', 'to', 'a',
 'truck', '.', '<END>'],
 ['<START>', 'the', 'child', 'is', 'playing', 'croquette', 'by', 'the',
 'truck', '.', '<END>'],
 ['<START>', 'the', 'kid', 'is', 'in', 'front', 'of', 'a', 'car', 'with',
 'a', 'put', 'and', 'a', 'ball', '.', '<END>'],
 ['<START>', 'the', 'little', 'boy', 'is', 'playing', 'with', 'a', 'croquet',
 'hammer', 'and', 'ball', 'beside', 'the', 'car', '.', '<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>']]
```

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.

```
In [30]: tokens_set = set()
for img_name in train_list:
    for tokens in descriptions[img_name]:
        tokens_set.update(set(tokens))
words = sorted(list(tokens_set))
```

```
In [31]: print(words)

ork', 'as', 'ascend', 'ascending', 'ascends', 'ash', 'ashen', 'ashtray',
'ashy', 'asia', 'asian', 'aside', 'asking', 'asleep', 'asphalt', 'assembl
e', 'assist', 'assistance', 'assisting', 'associated', 'astonishment', 'a
stride', 'astro', 'at', 'at&t', 'ate', 'athelete', 'athletes', 'athlet
e', 'athletes', 'athletic', 'athletics', 'atm', 'atmosphere', 'atomic',
'atop', 'atrium', 'attached', 'attaches', 'attaching', 'attack', 'attacke
d', 'attacking', 'attacks', 'attampts', 'attemping', 'attempt', 'attempte
d', 'attempting', 'attempts', 'attention', 'attentive', 'attentively', 'a
ttire', 'attraction', 'attractive', 'attrative', 'atv', 'atvs', 'audienc
e', 'auditorium', 'australian', 'auto', 'automobile', 'automobiles', 'aut
umn', 'autumnal', 'avoid', 'avoiding', 'avoids', 'avrovulcan.com', 'awai
t', 'awaiting', 'awaits', 'award', 'awards', 'away', 'awe', 'awkward', 'a
wkwardly', 'awning', 'awnings', 'ax', 'axe', 'babies', 'baby', 'back', 'b
ack-to-back', 'backbends', 'backdrop', 'backed', 'backgound', 'backgrou
d', 'background', 'backgrounds', 'backhand', 'backing', 'backlegs', 'back
less', 'backlit', 'backpack', 'backpacker', 'backpackers', 'backpacking',
'backpacks', 'backround', 'backs', 'backseat', 'backside', 'backsides',
'backstage', 'backstand', 'backstroke', 'backstrokes', 'backward', 'backw
ards', 'backyard', 'badge', 'badges', 'badly', 'badminton', 'bag', 'bagga
ge', 'baggy', 'bagging', 'baggrins', 'baggrins', 'baggs', 'baked', 'bak
```

```
In [32]: id_to_word = {idx: word for idx, word in enumerate(words)}
```

```
In [33]: word_to_id = {word: idx for idx, word in enumerate(words)}
```

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

Out[34]: 1985

```
In [35]: id_to_word[59] # should print a token
```

Out[35]: '<START>'

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.

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.

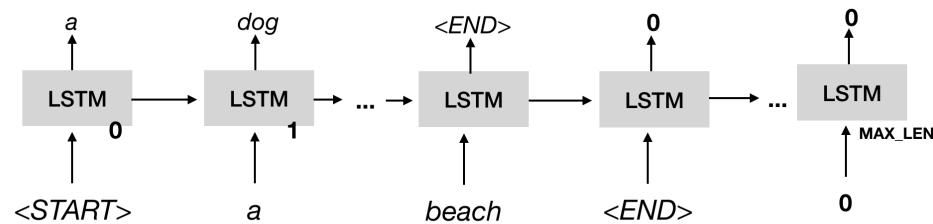
There are different ways to do this and our approach will be slightly different from the generator discussed in class.

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).

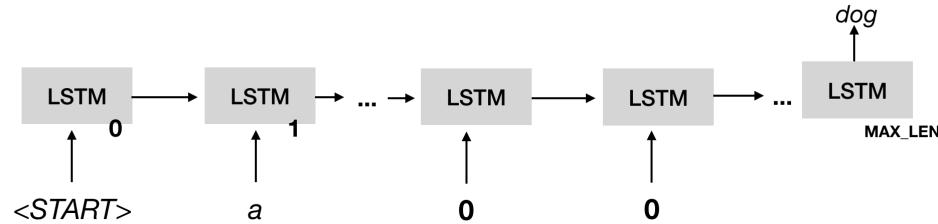
```
In [36]: max(len(description) for image_id in train_list for description in descriptions)
```

```
Out[36]: 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. 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.

```
In [37]: 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: "functional_3"

Layer (type)          Output Shape         Param #
=====
input_2 (InputLayer)   [(None, 40)]          0
embedding (Embedding) (None, 40, 300)        2312100
bidirectional (Bidirectional) (None, 1024)    3330048
dense (Dense)          (None, 7707)         7899675
=====
Total params: 13,541,823
Trainable params: 13,541,823
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.

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.

```
In [38]: def text_training_generator(batch_size=128):
    img_idx = 0
    des_idx = 0
    description = descriptions[train_list[img_idx]][des_idx]
    cur_input = 0
    cur_output = 1
    while True:
        input = np.zeros((batch_size, MAX_LEN))
        output = np.zeros((batch_size, vocab_size))
        for i in range(batch_size):
            input[i][0:cur_input+1] = [word_to_id[word] for word in description[0:cur_input+1]]
            output[i][word_to_id[description[cur_output]]] = 1
            if description[cur_output] == '<END>':
                cur_input = 0
                cur_output = 1
                des_idx += 1
                if des_idx == len(descriptions[train_list[img_idx]]):
                    des_idx = 0
                    img_idx += 1
                if img_idx == len(train_list):
                    img_idx = 0
                    description = descriptions[train_list[img_idx]][des_idx]
            else:
                cur_input += 1
                cur_output += 1
        yield (input, output)
    # ...
```

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`

```
In [39]: batch_size = 128
generator = text_training_generator(batch_size)
steps = len(train_list) * MAX_LEN // batch_size
```

```
In [40]: model.fit_generator(generator, steps_per_epoch=steps, verbose=True, epochs=10)

WARNING:tensorflow:From <ipython-input-40-e579278fa2eb>:1: Model.fit_generator (from tensorflow.python.keras.engine.training) is deprecated and will be removed in a future version.
Instructions for updating:
Please use Model.fit, which supports generators.

Epoch 1/10
1875/1875 [=====] - 103s 55ms/step - loss: 4.248
9 - accuracy: 0.2956
Epoch 2/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.687
7 - accuracy: 0.3587
Epoch 3/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.520
2 - accuracy: 0.3745
Epoch 4/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.412
8 - accuracy: 0.3838
Epoch 5/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.372
4 - accuracy: 0.3898
Epoch 6/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.295
2 - accuracy: 0.3958
Epoch 7/10
1875/1875 [=====] - 104s 55ms/step - loss: 3.276
3 - accuracy: 0.3977
Epoch 8/10
1875/1875 [=====] - 103s 55ms/step - loss: 3.249
2 - accuracy: 0.4017
Epoch 9/10
1875/1875 [=====] - 103s 55ms/step - loss: 3.258
9 - accuracy: 0.4023
Epoch 10/10
1875/1875 [=====] - 104s 56ms/step - loss: 3.244
7 - accuracy: 0.4041
```

```
Out[40]: <tensorflow.python.keras.callbacks.History at 0x7f65ea319f28>
```

```
In [41]: model.save_weights("gdrive/My Drive/"+my_data_dir+"/outputs/model_text.h5")
```

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

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.

```
In [42]: def decoder():
    # ...
    start = np.zeros((1, 40))
    start[0][0] = 59
    sentence = ['<START>']
    i = 1
    while True:
        output = model.predict(start)
        maxx = np.argmax(output)
        start[0][i] = maxx
        word = id_to_word[maxx]
        sentence.append(word)
        i += 1
        if word == '<END>':
            break
    return sentence
```

```
In [43]: print(decoder())
```

```
['<START>', 'a', 'man', 'is', 'standing', 'on', 'a', 'mountain', 'with',
 'his', 'arms', 'around', 'him', '.', '<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](#) (<https://docs.scipy.org/doc/numpy/reference/generated/numpy.random.multinomial.html>) function to do this.

```
In [44]: def sample_decoder():
    start = np.zeros((1,40))
    start[0][0] = 59
    sentence = ['<START>']
    i = 1
    while True:
        output = model.predict(start)
        output = output.reshape(vocab_size).astype('float64')
        output /= np.sum(output)
        choice = np.argmax(np.random.multinomial(1, output, size=1)[0])
        start[0][i] = choice
        word = id_to_word[choice]
        sentence.append(word)
        i += 1
        if word == '<END>':
            break
    return sentence
```

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.

```
In [46]: for i in range(10):
    print(sample_decoder())
['<START>', 'young', 'dogs', 'a', 'race', 'near', 'the', 'camera', 'with',
 'view', 'of', 'sunglasses', '.', '<END>']
['<START>', 'brown', 'dog', 'does', 'the', 'ball', '<END>']
['<START>', 'a', 'black', 'dog', 'runs', 'along', 'the', 'beach', '',
 'large', 'white', 'fence', 'in', 'the', 'background', '.', '<END>']
['<START>', 'a', 'smiling', 'little', 'boy', 'wearing', 'a', 'red', 'orange',
 'shirt', 'sits', 'on', 'the', 'playground', 'looks', 'past', 'a', 'yellow',
 'and', 'woman', 'in', 'blue', 'shirt', '.', '<END>']
['<START>', 'a', 'man', 'is', 'how', 'to', 'be', 'nearby', 'a', 'little',
 'boy', 'in', 'the', 'air', '.', '<END>']
['<START>', 'the', 'large', 'black', 'and', 'white', 'dog', 'is', 'walking',
 'in', 'the', 'sand', '.', '<END>']
['<START>', 'a', 'dog', 'a', 'handstand', 'to', 'catch', 'a', 'brown',
 'and', 'blue', 'toy', 'in', 'a', 'field', '.', '<END>']
['<START>', 'people', 'are', 'in', 'a', 'man', 'and', 'woman', 'in', 'she',
 '.', '<END>']
['<START>', 'a', 'group', 'of', 'people', 'sitting', 'in', 'the', 'snow',
 'walking', 'together', '.', '<END>']
['<START>', 'an', 'motorcycle', 'jumps', 'on', 'a', 'grass', '',
 'ground', 'is', 'field', '.', '<END>']
```

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:

In [25]:

```

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: "functional_3"

Layer (type)	Output Shape	Param #	Connected to
input_2 (InputLayer)	[(None, 2048)]	0	
dense_1 (Dense) [0][0]	(None, 300)	614700	input_2
input_3 (InputLayer)	[(None, 40)]	0	
repeat_vector (RepeatVector) [0][0]	(None, 40, 300)	0	dense_1
embedding_1 (Embedding) [0][0]	(None, 40, 300)	2312100	input_3
concatenate (Concatenate) concatenate[0][0]	(None, 40, 600)	0	repeat_vector embedding_1[0][0]
bidirectional_1 (Bidirectional) concatenate[0][0]	(None, 512)	1755136	concatenate
dense_2 (Dense) bidirectional_1[0][0]	(None, 7707)	3953691	bidirectional_1[0][0]

```
=====
Total params: 8,635,627
Trainable params: 8,635,627
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:

```
In [20]: enc_train = np.load("gdrive/My Drive/"+my_data_dir+"/outputs/encoded_images"
enc_dev = np.load("gdrive/My Drive/"+my_data_dir+"/outputs/encoded_images_d
```

```
In [21]: enc_train.shape
```

```
Out[21]: (6000, 2048)
```

```
In [22]: def training_generator(batch_size=128):
    img_idx = 0
    des_idx = 0
    description = descriptions[train_list[img_idx]][des_idx]
    cur_input = 0
    cur_output = 1
    while True:
        input = [np.zeros((batch_size, 2048)), np.zeros((batch_size, MAX_LEN))]
        output = np.zeros((batch_size, vocab_size))
        for i in range(batch_size):
            input[0][i] = enc_train[img_idx]
            input[1][i][0:cur_input+1] = [word_to_id[word] for word in description]
            output[i][word_to_id[description[cur_output]]] = 1
            if description[cur_output] == '<END>':
                cur_input = 0
                cur_output = 1
                des_idx += 1
                if des_idx == len(descriptions[train_list[img_idx]]):
                    des_idx = 0
                    img_idx += 1
                    if img_idx == len(train_list):
                        img_idx = 0
                    description = descriptions[train_list[img_idx]][des_idx]
                else:
                    cur_input += 1
                    cur_output += 1
        yield (input, output)
    # ...
```

```
In [56]: for i in training_generator():
    break
```

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

```
In [23]: batch_size = 128
generator = training_generator(batch_size)
steps = len(train_list) * MAX_LEN // batch_size
```

```
In [58]: model.fit_generator(generator, steps_per_epoch=steps, verbose=True, epochs=
```

```
Epoch 1/20
1875/1875 [=====] - 67s 36ms/step - loss: 4.2530
- accuracy: 0.2886
Epoch 2/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.6244
- accuracy: 0.3659
Epoch 3/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.4304
- accuracy: 0.3866
Epoch 4/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.3199
- accuracy: 0.3970
Epoch 5/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2591
- accuracy: 0.4038
Epoch 6/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2048
- accuracy: 0.4102
Epoch 7/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2219
- accuracy: 0.4088
Epoch 8/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2199
- accuracy: 0.4113
Epoch 9/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2058
- accuracy: 0.4138
Epoch 10/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2359
- accuracy: 0.4137
Epoch 11/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.2634
- accuracy: 0.4145
Epoch 12/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.2198
- accuracy: 0.4169
Epoch 13/20
1875/1875 [=====] - 67s 35ms/step - loss: 3.2235
- accuracy: 0.4196
Epoch 14/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.1989
- accuracy: 0.4225
Epoch 15/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.2066
- accuracy: 0.4227
Epoch 16/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.2210
- accuracy: 0.4223
Epoch 17/20
1875/1875 [=====] - 66s 35ms/step - loss: 3.2089
- accuracy: 0.4242
Epoch 18/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2390
- accuracy: 0.4230
Epoch 19/20
```

```
1875/1875 [=====] - 68s 36ms/step - loss: 3.2564
- accuracy: 0.4238
Epoch 20/20
1875/1875 [=====] - 67s 36ms/step - loss: 3.2496
- accuracy: 0.4253
```

Out[58]: <tensorflow.python.keras.callbacks.History at 0x7fe36a0a15f8>

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.

In [61]: `model.save_weights("gdrive/My Drive/" + my_data_dir + "/outputs/model.h5")`

to load the model:

In [26]: `model.load_weights("gdrive/My Drive/" + my_data_dir + "/outputs/model.h5")`

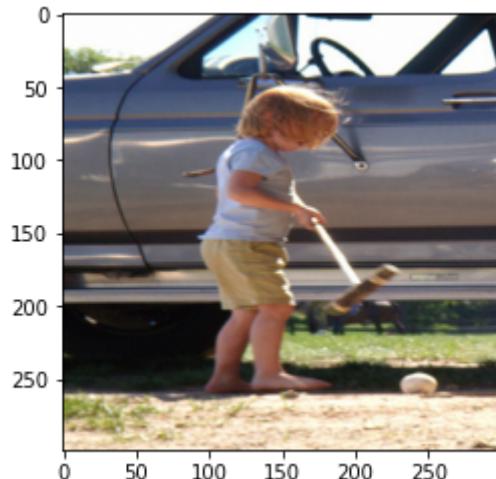
TODO: Now we are ready to actually generate image captions using the trained model. Modify the simple greedy decoder you wrote for the text-only generator, so that it takes an encoded image (a vector of length 2048) as input, and returns a sequence.

```
In [27]: def image_decoder(enc_image):
    start = [np.zeros((1, 2048)), np.zeros((1, 40))]
    start[0][0] = enc_image
    start[1][0][0] = 59
    sentence = ['<START>']
    i = 1
    while True:
        output = model.predict(start)
        output = output.reshape(vocab_size)
        maxx = np.argmax(output)
        start[1][0][i] = maxx
        word = id_to_word[maxx]
        sentence.append(word)
        i += 1
        if word == '<END>':
            break
    return sentence
```

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

```
In [94]: plt.imshow(get_image(train_list[1]))
image_decoder(enc_train[1])
```

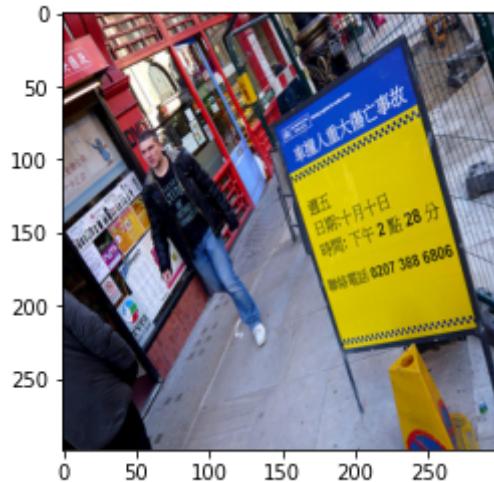
```
Out[94]: ['<START>',
 'a',
 'man',
 'in',
 'a',
 'blue',
 'shirt',
 'and',
 'a',
 'dog',
 'in',
 'a',
 'blue',
 'shirt',
 '.',
 '<END>']
```



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

```
In [29]: plt.imshow(get_image(dev_list[104]))
image_decoder(enc_dev[104])
```

```
Out[29]: ['<START>',
 'a',
 'man',
 'in',
 'a',
 'black',
 'shirt',
 'and',
 'a',
 'woman',
 'in',
 'a',
 'white',
 'shirt',
 'are',
 'in',
 'a',
 'white',
 'and',
 'blue',
 'and',
 'white',
 'and',
 'blue',
 'shirt',
 '.',
 '<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.

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.

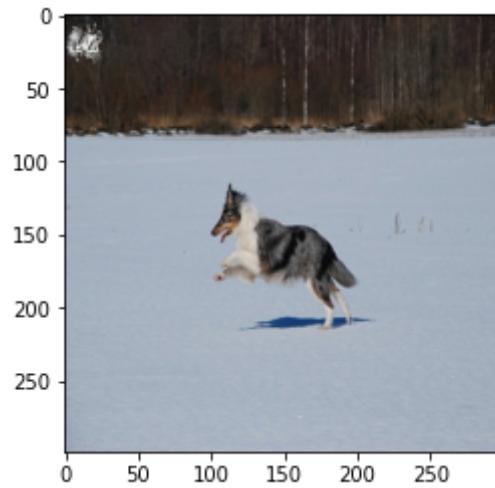
Note that you cannot use the occurrence of the "`<END>`" tag to terminate generation, because the tag may occur in different positions for different entries in the beam.

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

```
In [30]: import copy
```

```
In [31]: def img_beam_decoder(n, image_enc):
    start = [np.zeros((1, 2048)), np.zeros((1, 40))]
    start[0][0] = image_enc
    start[1][0][0] = 59
    sentence = ['<START>']
    i = 1
    beam = [(1.0, start)]
    while i < 40:
        cur_beam = []
        for probability, start in beam:
            output = model.predict(start)
            output = output.reshape(vocab_size)
            n_best = np.argsort(output)[-n:]
            for idx in n_best:
                cur_start = copy.deepcopy(start)
                cur_start[1][0][i] = idx
                cur_beam.append((probability * output[idx], cur_start))
        beam = sorted(cur_beam, key=lambda x: x[0])[-n:]
        i += 1
    sentence = [id_to_word[idx] for idx in beam[-1][1][1][0]]
    return sentence
```

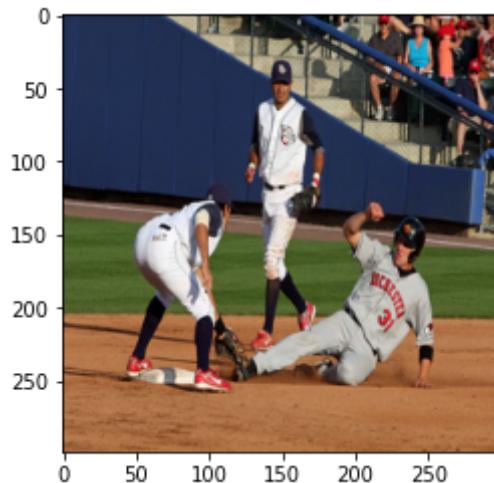
```
In [48]: plt.imshow(get_image(dev_list[700]))
img_beam_decoder(3, enc_dev[700])
```



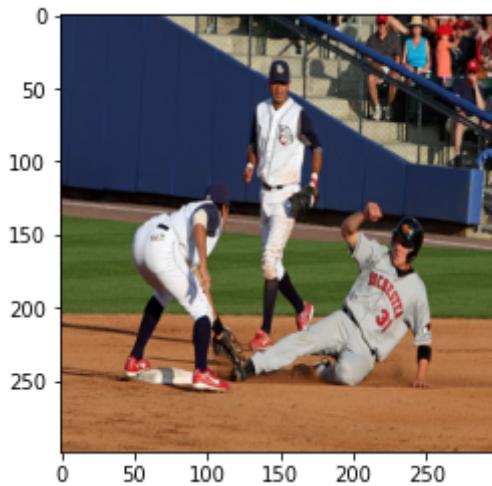
TODO Finally, before you submit this assignment, please show 5 development images, each with 1) their greedy output, 2) beam search at n=3 3) beam search at n=5.

```
In [ ]: plt.imshow(get_image(dev_list[104]))
image_decoder(enc_dev[104])
```

```
In [42]: plt.imshow(get_image(dev_list[100]))
img_beam_decoder(3, enc_dev[100])
```

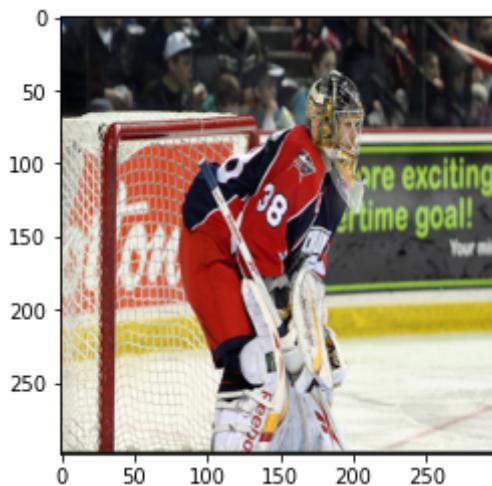


```
In [43]: plt.imshow(get_image(dev_list[100]))
img_beam_decoder(5, enc_dev[100])
```



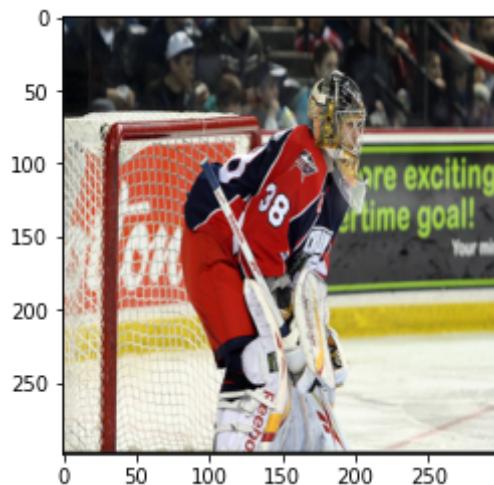
```
In [75]: plt.imshow(get_image(dev_list[200]))
image_decoder(enc_dev[200])
```

```
Out[75]: ['<START>',
 'a',
 'man',
 'in',
 'a',
 'red',
 'uniform',
 'is',
 'in',
 'a',
 'race',
 '.',
 '<END>']
```

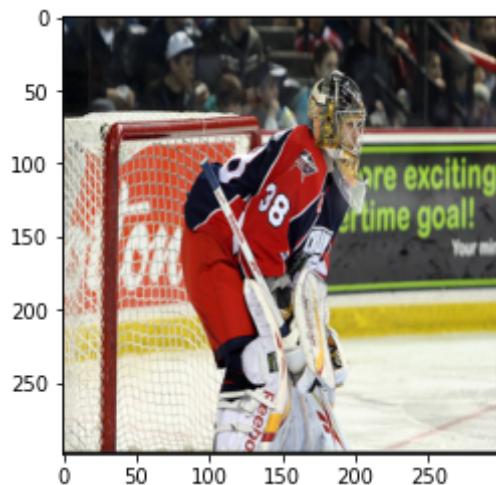


```
In [76]: plt.imshow(get_image(dev_list[200]))
img_beam_decoder(3, enc_dev[200])
```

```
Out[76]: [ '<START>' ,  
          'a' ,  
          'man' ,  
          'in' ,  
          'a' ,  
          'red' ,  
          'uniform' ,  
          'is' ,  
          'in' ,  
          'a' ,  
          'race' ,  
          '.' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '.' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '...' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '...' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '...' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ,  
          '<END>' ]
```

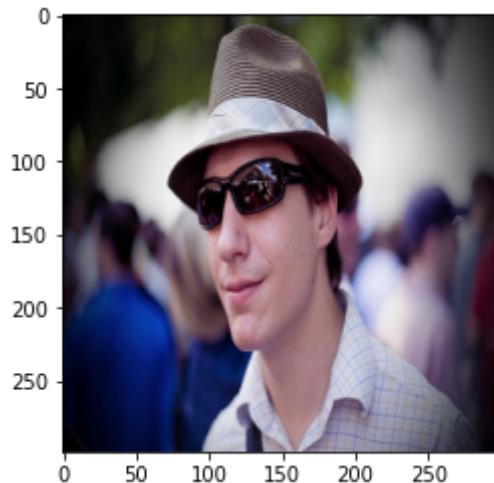


```
In [77]: plt.imshow(get_image(dev_list[200]))
img_beam_decoder(5, enc_dev[200])
```

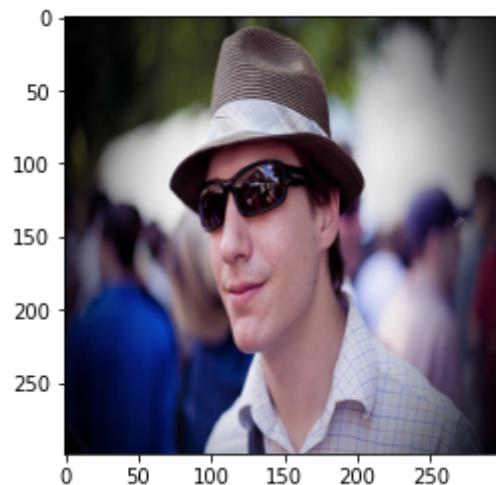


```
In [82]: plt.imshow(get_image(dev_list[300]))
image_decoder(enc_dev[300])
```

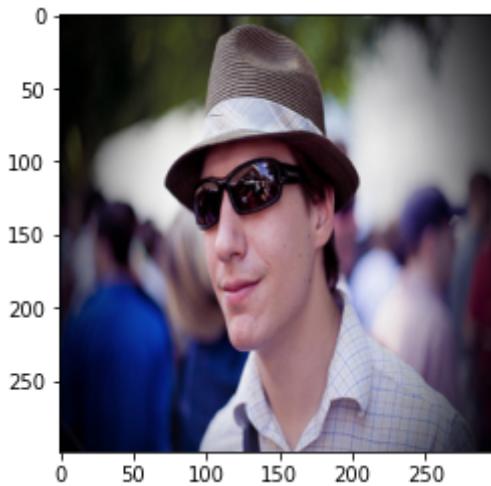
```
Out[82]: ['<START>',
 'a',
 'man',
 'in',
 'a',
 'black',
 'shirt',
 'and',
 'a',
 'woman',
 'in',
 'a',
 'white',
 'shirt',
 'are',
 'in',
 'a',
 'white',
 'and',
 'white',
 'and',
 'white',
 'shirt',
 '.',
 '<END>']
```



```
In [79]: plt.imshow(get_image(dev_list[300]))  
img_beam_decoder(3, enc_dev[300])
```

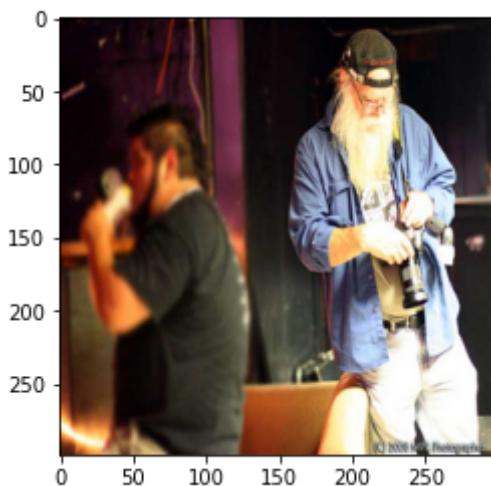


```
In [80]: plt.imshow(get_image(dev_list[300]))  
img_beam_decoder(5, enc_dev[300])
```

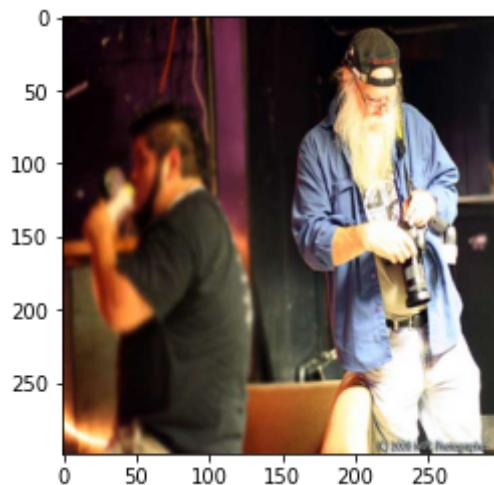


```
In [83]: plt.imshow(get_image(dev_list[400]))
image_decoder(enc_dev[400])
```

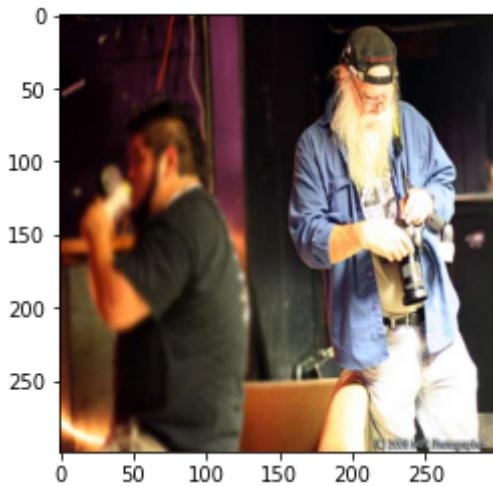
```
Out[83]: ['<START>',
 'a',
 'man',
 'in',
 'a',
 'black',
 'shirt',
 'and',
 'a',
 'man',
 'in',
 'a',
 'white',
 'shirt',
 '.',
 '<END>']
```



```
In [84]: plt.imshow(get_image(dev_list[400]))  
img_beam_decoder(3, enc_dev[400])
```

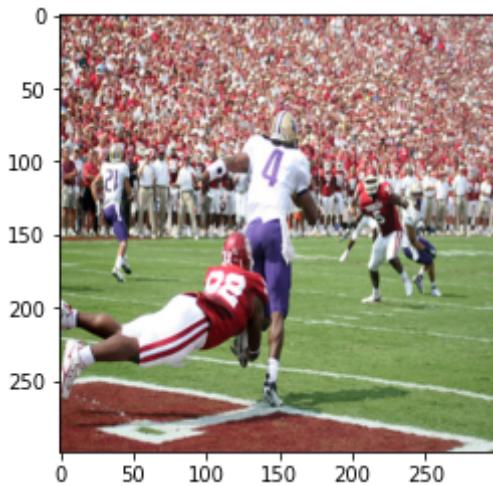


```
In [85]: plt.imshow(get_image(dev_list[400]))
img_beam_decoder(5, enc_dev[400])
```

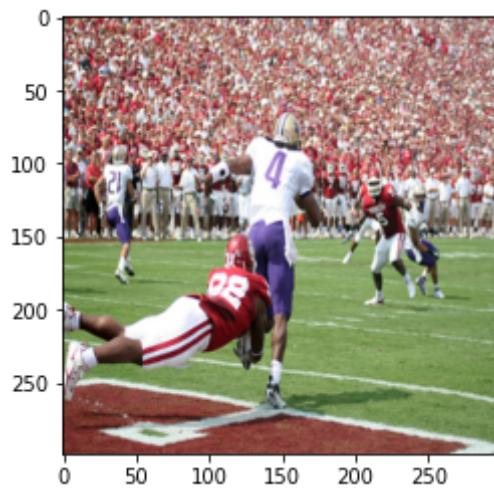


```
In [86]: plt.imshow(get_image(dev_list[500]))
image_decoder(enc_dev[500])
```

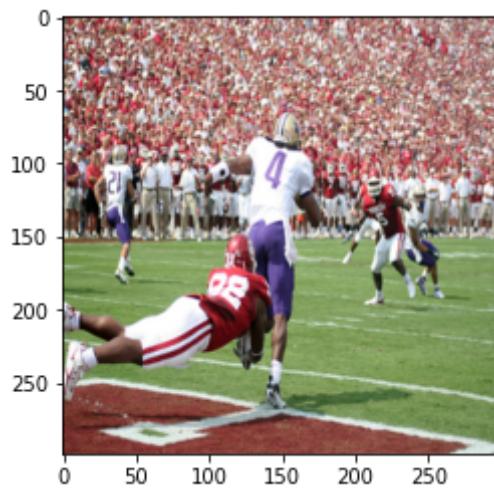
```
Out[86]: ['<START>',
 'a',
 'group',
 'of',
 'people',
 'are',
 'in',
 'a',
 'field',
 '.',
 '<END>']
```



```
In [87]: plt.imshow(get_image(dev_list[500]))
img_beam_decoder(3, enc_dev[500])
```



```
In [88]: plt.imshow(get_image(dev_list[500]))
img_beam_decoder(5, enc_dev[500])
```



In []: