

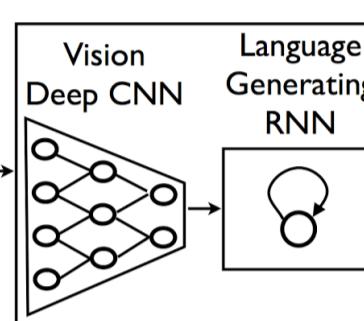
```
In [1]: 1 # set tf 1.x for colab  
2 %tensorflow_version 1.x
```

UsageError: Line magic function `%tensorflow_version` not found.

```
In [2]: 1 import warnings  
2 warnings.filterwarnings('ignore', category=DeprecationWarning)  
3 warnings.filterwarnings('ignore', category=FutureWarning)
```

Image Captioning Final Project

In this final project you will define and train an image-to-caption model, that can produce descriptions for real world images!



A group of people shopping at an outdoor market.
There are many vegetables at the fruit stand.

Model architecture: CNN encoder and RNN decoder. (<https://research.googleblog.com/2014/11/a-picture-is-worth-thousand-coherent.html>)
(<https://research.googleblog.com/2014/11/a-picture-is-worth-thousand-coherent.html>)

Import stuff

```
In [3]: 1 import sys  
2 sys.path.append("../")  
3 import grading  
4 import download_utils
```

```
In [4]: 1 download_utils.link_all_keras_resources()
```

```
In [5]: 1 import tensorflow as tf  
2 from tensorflow.contrib import keras  
3 import numpy as np  
4 %matplotlib inline  
5 import matplotlib.pyplot as plt  
6 L = keras.layers  
7 K = keras.backend  
8 import utils  
9 import time  
10 import zipfile  
11 import json  
12 from collections import defaultdict  
13 import re  
14 import random  
15 from random import choice  
16 import grading_utils  
17 import os  
18 from keras_utils import reset_tf_session  
19 import tqdm_utils
```

Using TensorFlow backend.

Prepare the storage for model checkpoints

```
In [6]: 1 # Leave USE_GOOGLE_DRIVE = False if you're running locally!
2 # We recommend to set USE_GOOGLE_DRIVE = True in Google Colab!
3 # If set to True, we will mount Google Drive, so that you can restore your checkpoint
4 # and continue training even if your previous Colab session dies.
5 # If set to True, follow on-screen instructions to access Google Drive (you must have a Google account).
6 USE_GOOGLE_DRIVE = False
7
8 def mount_google_drive():
9     from google.colab import drive
10    mount_directory = "/content/gdrive"
11    drive.mount(mount_directory)
12    drive_root = mount_directory + "/" + list(filter(lambda x: x[0] != '.', os.listdir(mount_directory)))[0] + "/col"
13    return drive_root
14
15 CHECKPOINT_ROOT = ""
16 if USE_GOOGLE_DRIVE:
17     CHECKPOINT_ROOT = mount_google_drive() + "/"
18
19 def get_checkpoint_path(epoch=None):
20     if epoch is None:
21         return os.path.abspath(CHECKPOINT_ROOT + "weights")
22     else:
23         return os.path.abspath(CHECKPOINT_ROOT + "weights_{}.format(epoch))")
24
25 # example of checkpoint dir
26 print(get_checkpoint_path(10))
```

D:\GoogleDrive\Study\Advanced Machine Learning - Coursera\1.intro-to-dl\week6\weights_10

Fill in your Coursera token and email

To successfully submit your answers to our grader, please fill in your Coursera submission token and email

```
In [7]: 1 grader = grading.Grader(assignment_key="NEDBg6CgEee8nQ6uE8a70A",
2                               all_parts=[ "19Wpv", "uJh73", "yiJkt", "rbpnH", "E2OIL", "YJR7z" ])
```

```
In [8]: 1 # token expires every 30 min
2 COURSERA_TOKEN = 'gu6es12ZRKHnD8sv'
3 COURSERA_EMAIL = 'lxwvictor@gmail.com'
```

Download data

Takes 10 hours and 20 GB. We've downloaded necessary files for you.

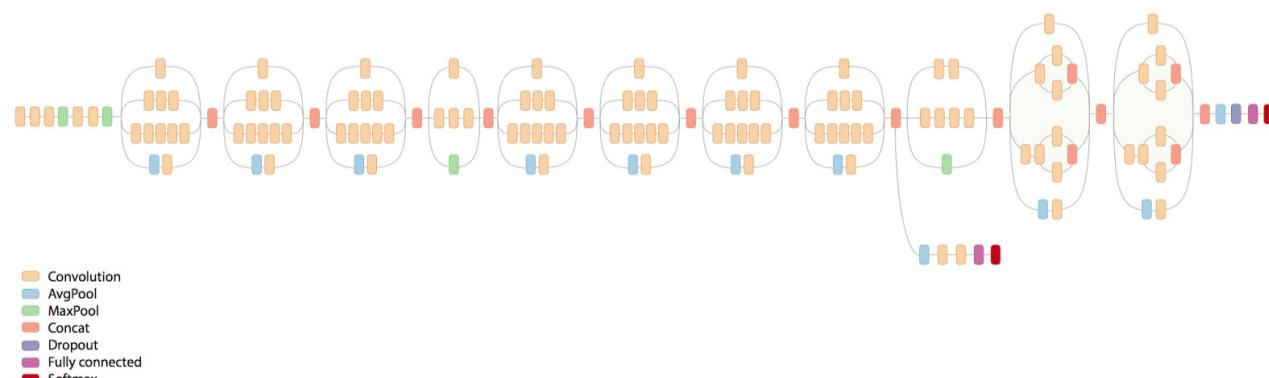
Relevant links (just in case):

- train images <http://msvocds.blob.core.windows.net/coco2014/train2014.zip> (<http://msvocds.blob.core.windows.net/coco2014/train2014.zip>)
- validation images <http://msvocds.blob.core.windows.net/coco2014/val2014.zip> (<http://msvocds.blob.core.windows.net/coco2014/val2014.zip>)
- captions for both train and validation http://msvocds.blob.core.windows.net/annotations-1-0-3/captions_train-val2014.zip (http://msvocds.blob.core.windows.net/annotations-1-0-3/captions_train-val2014.zip)

```
In [9]: 1 # we downloaded them for you, just link them here
2 download_utils.link_week_6_resources()
```

Extract image features

We will use pre-trained InceptionV3 model for CNN encoder (<https://research.googleblog.com/2016/03/train-your-own-image-classifier-with.html> (<https://research.googleblog.com/2016/03/train-your-own-image-classifier-with.html>)) and extract its last hidden layer as an embedding:



```
In [10]: 1 IMG_SIZE = 299
```

```
In [11]: 1 # we take the last hidden layer of InceptionV3 as an image embedding
2 def get_cnn_encoder():
3     K.set_learning_phase(False)
4     model = keras.applications.InceptionV3(include_top=False)
5     preprocess_for_model = keras.applications.inception_v3.preprocess_input
6
7     model = keras.models.Model(model.inputs, keras.layers.GlobalAveragePooling2D()(model.output))
8     return model, preprocess_for_model
```

Features extraction takes too much time on CPU:

- Takes 16 minutes on GPU.
- 25x slower (InceptionV3) on CPU and takes 7 hours.
- 10x slower (MobileNet) on CPU and takes 3 hours.

So we've done it for you with the following code:

```
# Load pre-trained model
reset_tf_session()
encoder, preprocess_for_model = get_cnn_encoder()

# extract train features
train_img_embeds, train_img_fns = utils.apply_model(
    "train2014.zip", encoder, preprocess_for_model, input_shape=(IMG_SIZE, IMG_SIZE))
utils.save_pickle(train_img_embeds, "train_img_embeds.pickle")
utils.save_pickle(train_img_fns, "train_img_fns.pickle")

# extract validation features
val_img_embeds, val_img_fns = utils.apply_model(
    "val2014.zip", encoder, preprocess_for_model, input_shape=(IMG_SIZE, IMG_SIZE))
utils.save_pickle(val_img_embeds, "val_img_embeds.pickle")
utils.save_pickle(val_img_fns, "val_img_fns.pickle")

# sample images for learners
def sample_zip(fn_in, fn_out, rate=0.01, seed=42):
    np.random.seed(seed)
    with zipfile.ZipFile(fn_in) as fin, zipfile.ZipFile(fn_out, "w") as fout:
        sampled = filter(lambda _: np.random.rand() < rate, fin.filelist)
        for zInfo in sampled:
            fout.writestr(zInfo, fin.read(zInfo))

sample_zip("train2014.zip", "train2014_sample.zip")
sample_zip("val2014.zip", "val2014_sample.zip")
```

```
In [12]: 1 # Load prepared embeddings
2 train_img_embeds = utils.read_pickle("train_img_embeds.pickle")
3 train_img_fns = utils.read_pickle("train_img_fns.pickle")
4 val_img_embeds = utils.read_pickle("val_img_embeds.pickle")
5 val_img_fns = utils.read_pickle("val_img_fns.pickle")
6 # check shapes
7 print(train_img_embeds.shape, len(train_img_fns))
8 print(val_img_embeds.shape, len(val_img_fns))
```

```
(82783, 2048) 82783
(40504, 2048) 40504
```

```
In [13]: 1 # check prepared samples of images
2 list(filter(lambda x: x.endswith("_sample.zip"), os.listdir(".")))
```

```
Out[13]: ['train2014_sample.zip', 'val2014_sample.zip']
```

Extract captions for images

In [14]:

```
1 # extract captions from zip
2 def get_captions_for_fns(fns, zip_fn, zip_json_path):
3     zf = zipfile.ZipFile(zip_fn)
4     j = json.loads(zf.read(zip_json_path).decode("utf8"))
5     id_to_fn = {img["id"]: img["file_name"] for img in j["images"]}
6     fn_to_caps = defaultdict(list)
7     for cap in j['annotations']:
8         fn_to_caps[id_to_fn[cap['image_id']]].append(cap['caption'])
9     fn_to_caps = dict(fn_to_caps)
10    return list(map(lambda x: fn_to_caps[x], fns))
11
12 trainCaptions = get_captions_for_fns(train_img_fns, "captions_train-val2014.zip",
13                                     "annotations/captions_train2014.json")
14
15 valCaptions = get_captions_for_fns(val_img_fns, "captions_train-val2014.zip",
16                                     "annotations/captions_val2014.json")
17
18 # check shape
19 print(len(train_img_fns), len(trainCaptions))
20 print(len(val_img_fns), len(valCaptions))
```

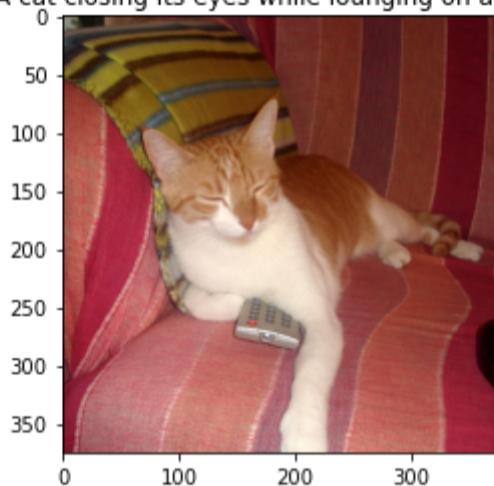
82783 82783
40504 40504

In [15]:

```
1 # Look at training example (each has 5 captions)
2 def show_trainig_example(train_img_fns, trainCaptions, example_idx=0):
3     """
4     You can change example_idx and see different images
5     """
6     zf = zipfile.ZipFile("train2014_sample.zip")
7     captions_by_file = dict(zip(train_img_fns, trainCaptions))
8     all_files = set(train_img_fns)
9     found_files = list(filter(lambda x: x.filename.rsplit("/")[-1] in all_files, zf.filelist))
10    example = found_files[example_idx]
11    img = utils.decode_image_from_buf(zf.read(example))
12    plt.imshow(utils.image_center_crop(img))
13    plt.title("\n".join(captions_by_file[example.filename.rsplit("/")[-1]]))
14    plt.show()
15
16 show_trainig_example(train_img_fns, trainCaptions, example_idx=142)
```

A cat sitting on a pink striped couch
An orange and white cat sitting in a striped chair.
An orange and white cat sleeping on a remote
Brown and white cat sleeping on couch while lying on remote.

A cat closing its eyes while lounging on a chair.



Prepare captions for training

In [16]:

```
1 # preview captions data
2 trainCaptions[:2]
```

Out[16]: [['A long dirt road going through a forest.',
'A SCENE OF WATER AND A PATH WAY',
'A sandy path surrounded by trees leads to a beach.',
'Ocean view through a dirt road surrounded by a forested area. ',
'dirt path leading beneath barren trees to open plains'],
['A group of zebra standing next to each other.',
'This is an image of zebras drinking',
'ZEBRAS AND BIRDS SHARING THE SAME WATERING HOLE',
'Zebras that are bent over and drinking water together.',
'a number of zebras drinking water near one another']]

```
In [17]: 1 # special tokens
2 PAD = "#PAD#"
3 UNK = "#UNK#"
4 START = "#START#"
5 END = "#END#"
6
7 # split sentence into tokens (split into lowercased words)
8 def split_sentence(sentence):
9     return list(filter(lambda x: len(x) > 0, re.split('\W+', sentence.lower())))
10
11 def generate_vocabulary(train_captions):
12 """
13     Return {token: index} for all train tokens (words) that occur 5 times or more,
14     `index` should be from 0 to N, where N is a number of unique tokens in the resulting dictionary.
15     Use `split_sentence` function to split sentence into tokens.
16     Also, add PAD (for batch padding), UNK (unknown, out of vocabulary),
17         START (start of sentence) and END (end of sentence) tokens into the vocabulary.
18 """
19 ### YOUR CODE HERE ###
20 from collections import Counter
21 sentence_list = [item for sublist in train_captions for item in sublist]
22 tokens_nested_list = list(map(lambda x: split_sentence(x), sentence_list))
23 tokens = [item for sublist in tokens_nested_list for item in sublist]
24 tc = Counter(tokens)
25
26 vocab = [item for item in tc if tc[item] >= 5] + [PAD, UNK, START, END]
27 return {token: index for index, token in enumerate(sorted(vocab))}
28
29 def caption_tokens_to_indices(captions, vocab):
30 """
31     `captions` argument is an array of arrays:
32     [
33         [
34             "image1 caption1",
35             "image1 caption2",
36             ...
37         ],
38         [
39             "image2 caption1",
40             "image2 caption2",
41             ...
42         ],
43         ...
44     ]
45     Use `split_sentence` function to split sentence into tokens.
46     Replace all tokens with vocabulary indices, use UNK for unknown words (out of vocabulary).
47     Add START and END tokens to start and end of each sentence respectively.
48     For the example above you should produce the following:
49     [
50         [
51             [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
52             [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
53             ...
54         ],
55         ...
56     ]
57 """
58 ### YOUR CODE HERE ###
59 res = []
60 for caption in captions:
61     tokens = list(map(lambda x: split_sentence(x), caption))
62     indices = [[vocab.get(item, vocab[UNK]) for item in inner] for inner in tokens]
63     indices = [[vocab[START]] + item + [vocab[END]] for item in indices]
64     res.append(indices)
65
66 return res
```

```
In [18]: 1 # prepare vocabulary
2 vocab = generate_vocabulary(train_captions)
3 vocab_inverse = {idx: w for w, idx in vocab.items()}
4 print(len(vocab))
```

8769

```
In [19]: 1 # replace tokens with indices
2 train_captions_indexed = caption_tokens_to_indices(train_captions, vocab)
3 val_captions_indexed = caption_tokens_to_indices(val_captions, vocab)
```

Captions have different length, but we need to batch them, that's why we will add PAD tokens so that all sentences have an equal length.

We will crunch LSTM through all the tokens, but we will ignore padding tokens during loss calculation.

```
In [20]: 1 # we will use this during training
2 def batch_captions_to_matrix(batch_captions, pad_idx, max_len=None):
3     """
4         `batch_captions` is an array of arrays:
5         [
6             [vocab[START], ..., vocab[END]],
7             [vocab[START], ..., vocab[END]],
8             ...
9         ]
10        Put vocabulary indexed captions into np.array of shape (len(batch_captions), columns),
11        where "columns" is max(map(len, batch_captions)) when max_len is None
12        and "columns" = min(max_len, max(map(len, batch_captions))) otherwise.
13        Add padding with pad_idx where necessary.
14        Input example: [[1, 2, 3], [4, 5]]
15        Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=None
16        Output example: np.array([[1, 2], [4, 5]]) if max_len=2
17        Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=100
18        Try to use numpy, we need this function to be fast!
19    """
20    #####YOUR CODE HERE#####
21    if not max_len:
22        max_len = max(map(len, batch_captions))
23    else:
24        max_len = min(max_len, max(map(len, batch_captions)))
25
26    matrix = np.array([x + [pad_idx]*(max_len-len(x))
27                      if max_len >= len(x)
28                      else x[:max_len]
29                      for x in batch_captions])
30
31    return matrix
```

```
In [21]: 1 batch_captions_to_matrix(train_captions_indexed[1], vocab[PAD], 15).shape
```

Out[21]: (5, 11)

```
In [22]: 1 ## GRADED PART, DO NOT CHANGE!
2 # Vocabulary creation
3 grader.set_answer("19Wpv", grading_utils.test_vocab(vocab, PAD, UNK, START, END))
4 # Captions indexing
5 grader.set_answer("uJh73", grading_utils.test_captions_indexing(train_captions_indexed, vocab, UNK))
6 # Captions batching
7 grader.set_answer("yiJkt", grading_utils.test_captions_batching(batch_captions_to_matrix))
```

```
In [23]: 1 # you can make submission with answers so far to check yourself at this stage
2 grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

You used an invalid email or your token may have expired. Please make sure you have entered all fields correctly. Try generating a new token if the issue still persists.

```
In [24]: 1 # make sure you use correct argument in caption_tokens_to_indices
2 assert len(caption_tokens_to_indices(train_captions[:10], vocab)) == 10
3 assert len(caption_tokens_to_indices(train_captions[:5], vocab)) == 5
```

Training

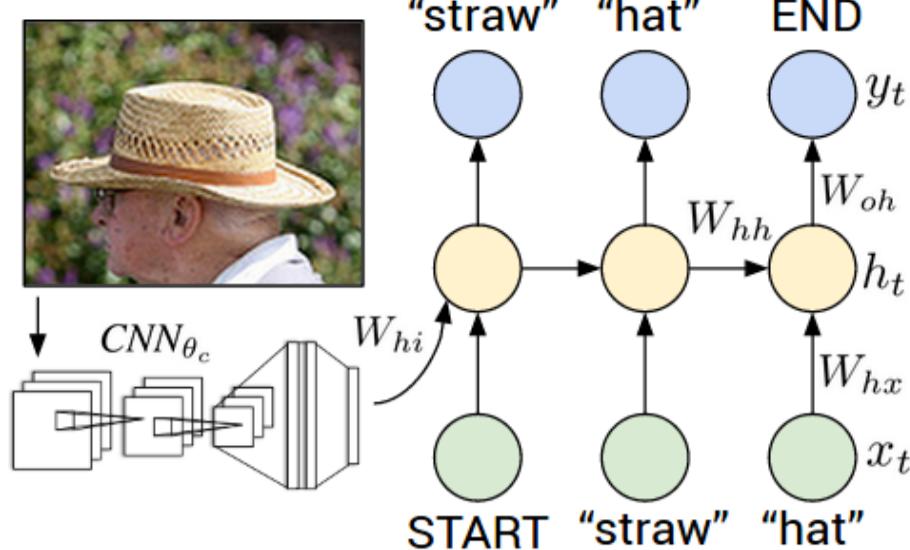
Define architecture

Since our problem is to generate image captions, RNN text generator should be conditioned on image. The idea is to use image features as an initial state for RNN instead of zeros.

Remember that you should transform image feature vector to RNN hidden state size by fully-connected layer and then pass it to RNN.

During training we will feed ground truth tokens into the lstm to get predictions of next tokens.

Notice that we don't need to feed last token (END) as input (<http://cs.stanford.edu/people/karpathy/>) (<http://cs.stanford.edu/people/karpathy/>):



```
In [25]: 1 IMG_EMBED_SIZE = train_img_embeds.shape[1]
2 IMG_EMBED_BOTTLENECK = 120
3 WORD_EMBED_SIZE = 100
4 LSTM_UNITS = 300
5 LOGIT_BOTTLENECK = 120
6 pad_idx = vocab[PAD]
```

```
In [26]: 1 train_img_embeds.shape, len(vocab)
```

Out[26]: ((82783, 2048), 8769)

```
In [27]: 1 # remember to reset your graph if you want to start building it from scratch!
2 s = reset_tf_session()
3 tf.set_random_seed(42)
```

WARNING:tensorflow:From ..\keras_utils.py:68: The name tf.get_default_session is deprecated. Please use tf.compat.v1.get_default_session instead.

WARNING:tensorflow:From ..\keras_utils.py:75: The name tf.ConfigProto is deprecated. Please use tf.compat.v1.ConfigProto instead.

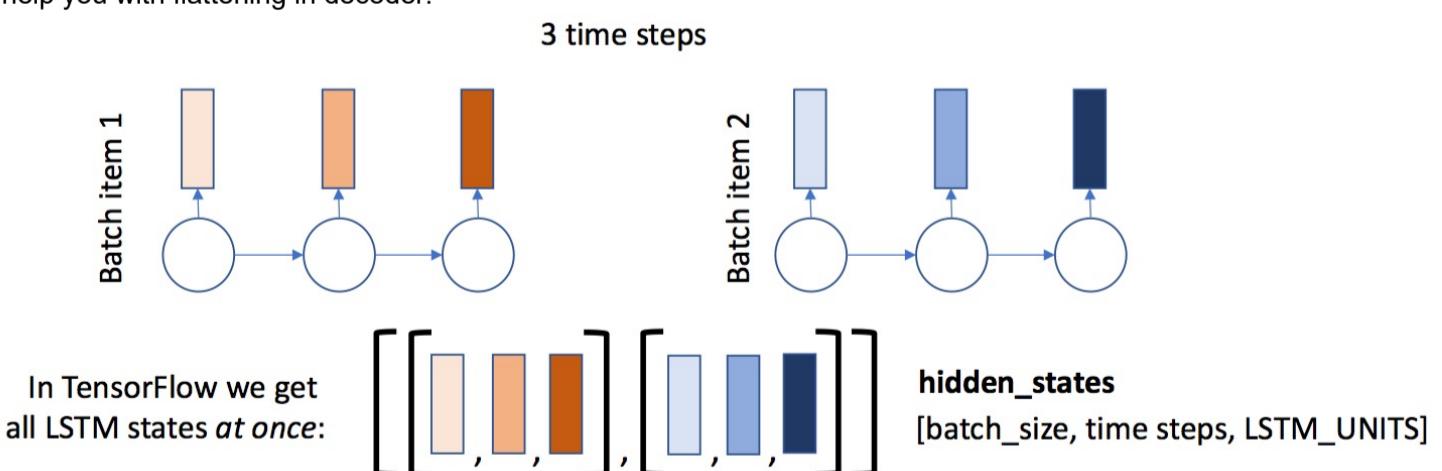
WARNING:tensorflow:From ..\keras_utils.py:77: The name tf.InteractiveSession is deprecated. Please use tf.compat.v1.InteractiveSession instead.

Here we define decoder graph.

We use Keras layers where possible because we can use them in functional style with weights reuse like this:

```
dense_layer = L.Dense(42, input_shape=(None, 100) activation='relu')
a = tf.placeholder('float32', [None, 100])
b = tf.placeholder('float32', [None, 100])
dense_layer(a) # that's how we applied dense layer!
dense_layer(b) # and again
```

Here's a figure to help you with flattening in decoder:



But we need to calculate token probability *for each time step of every example!*

That's why we want to *flatten* these states and apply dense layers to calculate *all token logits at once*:



In [28]:

```
1 # tf.reset_default_graph()
2 class decoder:
3     # [batch_size, IMG_EMBED_SIZE] of CNN image features
4     img_embeds = tf.placeholder('float32', [None, IMG_EMBED_SIZE])
5     # [batch_size, time steps] of word ids
6     sentences = tf.placeholder('int32', [None, None])
7
8     # we use bottleneck here to reduce the number of parameters
9     # image embedding -> bottleneck
10    img_embed_to_bottleneck = L.Dense(IMG_EMBED_BOTTLENECK,
11                                      input_shape=(None, IMG_EMBED_SIZE),
12                                      activation='elu')
13    # image embedding bottleneck -> lstm initial state
14    img_embed_bottleneck_to_h0 = L.Dense(LSTM_UNITS,
15                                         input_shape=(None, IMG_EMBED_BOTTLENECK),
16                                         activation='elu')
17    # word -> embedding
18    word_embed = L.Embedding(len(vocab), WORD_EMBED_SIZE)
19    # lstm cell (from tensorflow)
20    lstm = tf.nn.rnn_cell.LSTMCell(LSTM_UNITS)
21
22    # we use bottleneck here to reduce model complexity
23    # lstm output -> logits bottleneck
24    token_logits_bottleneck = L.Dense(LOGIT_BOTTLENECK,
25                                       input_shape=(None, LSTM_UNITS),
26                                       activation="elu")
27    # logits bottleneck -> Logits for next token prediction
28    token_logits = L.Dense(len(vocab),
29                           input_shape=(None, LOGIT_BOTTLENECK))
30
31    # initial lstm cell state of shape (None, LSTM_UNITS),
32    # we need to condition it on `img_embeds` placeholder.
33    ### YOUR CODE HERE ###
34    c0 = h0 = img_embed_bottleneck_to_h0(img_embed_to_bottleneck(img_embeds))
35
36    # embed all tokens but the last for lstm input,
37    # remember that L.Embedding is callable,
38    # use `sentences` placeholder as input.
39    ### YOUR CODE HERE ###
40    word_embeds = word_embed(sentences[:, :-1])
41
42    # during training we use ground truth tokens `word_embeds` as context for next token prediction.
43    # that means that we know all the inputs for our lstm and can get
44    # all the hidden states with one tensorflow operation (tf.nn.dynamic_rnn).
45    # `hidden_states` has a shape of [batch_size, time steps, LSTM_UNITS].
46    hidden_states, _ = tf.nn.dynamic_rnn(lstm, word_embeds,
47                                         initial_state=tf.nn.rnn_cell.LSTMStateTuple(c0, h0))
48
49    # now we need to calculate token logits for all the hidden states
50
51    # first, we reshape `hidden_states` to [-1, LSTM_UNITS]
52    ### YOUR CODE HERE ###
53    flat_hidden_states = tf.reshape(hidden_states, (-1, LSTM_UNITS))
54
55    # then, we calculate logits for next tokens using `token_logits_bottleneck` and `token_logits` layers
56    ### YOUR CODE HERE ###
57    flat_token_logits = token_logits(token_logits_bottleneck(flat_hidden_states))
58
59    # then, we flatten the ground truth token ids.
60    # remember, that we predict next tokens for each time step,
61    # use `sentences` placeholder.
62    ### YOUR CODE HERE ###
63    flat_ground_truth = tf.reshape(sentences[:, 1:], (-1,))
64
65    # we need to know where we have real tokens (not padding) in `flat_ground_truth`,
66    # we don't want to propagate the loss for padded output tokens,
67    # fill `flat_loss_mask` with 1.0 for real tokens (not pad_idx) and 0.0 otherwise.
68    ### YOUR CODE HERE ###
69    flat_loss_mask = tf.not_equal(flat_ground_truth, vocab[PAD])
70
71    # compute cross-entropy between `flat_ground_truth` and `flat_token_logits` predicted by lstm
72    xent = tf.nn.sparse_softmax_cross_entropy_with_logits(
73        labels=flat_ground_truth,
74        logits=flat_token_logits
75    )
76
77    # compute average `xent` over tokens with nonzero `flat_loss_mask`.
78    # we don't want to account misclassification of PAD tokens, because that doesn't make sense,
79    # we have PAD tokens for batching purposes only!
80    ### YOUR CODE HERE ###
81    loss = tf.reduce_mean(tf.boolean_mask(xent, flat_loss_mask))
```

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\init_ops.py:1251: calling VarianceScaling.__init__ (from tensorflow.python.ops.init_ops) with dtype is deprecated and will be removed in a future version.

Instructions for updating:

Call initializer instance with the dtype argument instead of passing it to the constructor

```
WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\keras\initialize.py:119: calling RandomUniform.__init__ (from tensorflow.python.ops.init_ops) with dtype is deprecated and will be removed in a future version.  
Instructions for updating:  
Call initializer instance with the dtype argument instead of passing it to the constructor  
WARNING:tensorflow:From <ipython-input-28-389ff8a7bcc0>:20: LSTMCell.__init__ (from tensorflow.python.ops.rnn_cell_impl) is deprecated and will be removed in a future version.  
Instructions for updating:  
This class is equivalent as tf.keras.layers.LSTMCell, and will be replaced by that in Tensorflow 2.0.  
WARNING:tensorflow:From <ipython-input-28-389ff8a7bcc0>:47: dynamic_rnn (from tensorflow.python.ops.rnn) is deprecated and will be removed in a future version.  
Instructions for updating:  
Please use `keras.layers.RNN(cell)`, which is equivalent to this API  
WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\rnn_cell_impl.py:961: calling Zeros.__init__ (from tensorflow.python.ops.init_ops) with dtype is deprecated and will be removed in a future version.  
Instructions for updating:  
Call initializer instance with the dtype argument instead of passing it to the constructor  
WARNING:tensorflow:Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x000020601442588>> could not be transformed and will be executed as-is. Please report this to the Autograph team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERTOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>>: AttributeError: module 'gast' has no attribute 'Num'  
WARNING: Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>> could not be transformed and will be executed as-is. Please report this to the Autograph team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERTOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>>: AttributeError: module 'gast' has no attribute 'Num'  
WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\ops\array_ops.py:1354: add_dispatch_support.<locals>.wrapper (from tensorflow.python.ops.array_ops) is deprecated and will be removed in a future version.  
Instructions for updating:  
Use tf.where in 2.0, which has the same broadcast rule as np.where
```

In [29]:

```
1 # define optimizer operation to minimize the loss  
2 optimizer = tf.train.AdamOptimizer(learning_rate=0.001)  
3 train_step = optimizer.minimize(decoder.loss)  
4  
5 # will be used to save/Load network weights.  
6 # you need to reset your default graph and define it in the same way to be able to load the saved weights!  
7 saver = tf.train.Saver()  
8  
9 # initialize all variables  
10 s.run(tf.global_variables_initializer())
```

In [30]:

```
1 ## GRADED PART, DO NOT CHANGE!  
2 # Decoder shapes test  
3 grader.set_answer("rbpnH", grading_utils.test_decoder_shapes(decoder, IMG_EMBED_SIZE, vocab, s))  
4 # Decoder random loss test  
5 grader.set_answer("E2OIL", grading_utils.test_random_decoder_loss(decoder, IMG_EMBED_SIZE, vocab, s))
```

In [31]:

```
1 # you can make submission with answers so far to check yourself at this stage  
2 grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

You used an invalid email or your token may have expired. Please make sure you have entered all fields correctly. Try generating a new token if the issue still persists.

Training loop

Evaluate train and validation metrics through training and log them. Ensure that loss decreases.

In [32]:

```
1 train_captions_indexed = np.array(train_captions_indexed)  
2 val_captions_indexed = np.array(val_captions_indexed)
```

In [33]:

```
1 # generate batch via random sampling of images and captions for them,
2 # we use `max_len` parameter to control the length of the captions (truncating long captions)
3 def generate_batch(images_embeddings, indexed_captions, batch_size, max_len=None):
4     """
5         `images_embeddings` is a np.array of shape [number of images, IMG_EMBED_SIZE].
6         `indexed_captions` holds 5 vocabulary indexed captions for each image:
7         [
8             [
9                 [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
10                [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
11                ...
12            ],
13            ...
14        ]
15     Generate a random batch of size `batch_size`.
16     Take random images and choose one random caption for each image.
17     Remember to use `batch_captions_to_matrix` for padding and respect `max_len` parameter.
18     Return feed dict {decoder.img_embeds: ..., decoder.sentences: ...}.
19     """
20     ### YOUR CODE HERE ###
21     total = images_embeddings.shape[0]
22     idx = np.random.randint(total, size = batch_size)
23     batch_image_embeddings = images_embeddings[idx, :]
24
25     ### YOUR CODE HERE ###
26     def sample_1_caption(image_captions):
27         rand_idx = np.random.randint(len(image_captions), size = 1)[0]
28         return image_captions[rand_idx]
29
30     batch_captions_matrix = batch_captions_to_matrix([sample_1_caption(item) for item in indexed_captions[idx]],
31                                                     vocab[PAD], max_len = max_len)
32
33     return {decoder.img_embeds: batch_image_embeddings,
34            decoder.sentences: batch_captions_matrix}
```

In [34]:

```
1 batch_size = 64
2 n_epochs = 12
3 n_batches_per_epoch = 1000
4 n_validation_batches = 100 # how many batches are used for validation after each epoch
```

In [35]:

```
1 # you can load trained weights here
2 # uncomment the next line if you need to load weights
3 # saver.restore(s, get_checkpoint_path(epoch=4))
```

Look at the training and validation loss, they should be decreasing!

In [36]:

```
1 # actual training loop
2 MAX_LEN = 20 # truncate long captions to speed up training
3
4 # to make training reproducible
5 np.random.seed(42)
6 random.seed(42)
7
8 for epoch in range(n_epochs):
9
10    train_loss = 0
11    pbar = tqdm_utils.tqdm_notebook_failsafe(range(n_batches_per_epoch))
12    counter = 0
13    for _ in pbar:
14        train_loss += s.run([decoder.loss, train_step],
15                            generate_batch(train_img_embeds,
16                                         train_captions_indexed,
17                                         batch_size,
18                                         MAX_LEN))[0]
19        counter += 1
20        pbar.set_description("Training loss: %f" % (train_loss / counter))
21
22    train_loss /= n_batches_per_epoch
23
24    val_loss = 0
25    for _ in range(n_validation_batches):
26        val_loss += s.run(decoder.loss, generate_batch(val_img_embeds,
27                           val_captions_indexed,
28                           batch_size,
29                           MAX_LEN))
30    val_loss /= n_validation_batches
31
32    print('Epoch: {}, train loss: {}, val loss: {}'.format(epoch, train_loss, val_loss))
33
34    # save weights after finishing epoch
35    saver.save(s, get_checkpoint_path(epoch))
36
37 print("Finished!")
```

Training loss: 4.259904: 100% 1000/1000 [05:34<00:00, 2.99it/s]

Epoch: 0, train loss: 4.2599036183357235, val loss: 3.668520784378052

Training loss: 3.352907: 100% 1000/1000 [00:53<00:00, 18.58it/s]

Epoch: 1, train loss: 3.3529067544937132, val loss: 3.1696912479400634

Training loss: 2.998766: 100% 1000/1000 [04:40<00:00, 3.56it/s]

Epoch: 2, train loss: 2.998766419649124, val loss: 2.9141571307182312

Training loss: 2.851692: 100% 1000/1000 [04:14<00:00, 3.94it/s]

Epoch: 3, train loss: 2.8516918346881868, val loss: 2.8400726103782654

Training loss: 2.763287: 100% 1000/1000 [01:18<00:00, 12.69it/s]

Epoch: 4, train loss: 2.763286616325378, val loss: 2.8050848507881163

Training loss: 2.681020: 100% 1000/1000 [00:52<00:00, 19.01it/s]

Epoch: 5, train loss: 2.681019562959671, val loss: 2.747049078941345

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\training\saver.py:9
60: remove_checkpoint (from tensorflow.python.training.checkpoint_management) is deprecated and will be removed in a future version.

Instructions for updating:

Use standard file APIs to delete files with this prefix.

Training loss: 2.636492: 100% 1000/1000 [02:54<00:00, 5.73it/s]

Epoch: 6, train loss: 2.636492332935333, val loss: 2.6975600266456605

Training loss: 2.593149: 100% 1000/1000 [02:28<00:00, 6.75it/s]

Epoch: 7, train loss: 2.5931488173007966, val loss: 2.644725923538208

Training loss: 2.564821: 100% 1000/1000 [01:21<00:00, 12.30it/s]

```
Epoch: 8, train loss: 2.5648211677074433, val loss: 2.627164342403412
```

```
Training loss: 2.528064: 100% 1000/1000 [00:54<00:00, 18.34it/s]
```

```
Epoch: 9, train loss: 2.5280644454956054, val loss: 2.59504102230072
```

```
Training loss: 2.502668: 100% 1000/1000 [01:06<00:00, 14.93it/s]
```

```
Epoch: 10, train loss: 2.502668368577957, val loss: 2.592053937911987
```

```
Training loss: 2.482272: 100% 1000/1000 [00:25<00:00, 39.68it/s]
```

```
Epoch: 11, train loss: 2.4822721979618074, val loss: 2.5885070610046386  
Finished!
```

```
In [37]: 1 ## GRADED PART, DO NOT CHANGE!  
2 # Validation Loss  
3 grader.set_answer("YJR7z", grading_utils.test_validation_loss(  
4     decoder, s, generate_batch, val_img_embeds, val_captions_indexed))
```

```
100% 1000/1000 [00:11<00:00, 84.16it/s]
```

```
In [38]: 1 # you can make submission with answers so far to check yourself at this stage  
2 grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

You used an invalid email or your token may have expired. Please make sure you have entered all fields correctly. Try generating a new token if the issue still persists.

```
In [39]: 1 # check that it's learnt something, outputs accuracy of next word prediction (should be around 0.5)  
2 from sklearn.metrics import accuracy_score, log_loss  
3  
4 def decode_sentence(sentence_indices):  
5     return " ".join(list(map(vocab_inverse.get, sentence_indices)))  
6  
7 def check_after_training(n_examples):  
8     fd = generate_batch(train_img_embeds, train_captions_indexed, batch_size)  
9     logits = decoder.flat_token_logits.eval(fd)  
10    truth = decoder.flat_ground_truth.eval(fd)  
11    mask = decoder.flat_loss_mask.eval(fd).astype(bool)  
12    print("Loss:", decoder.loss.eval(fd))  
13    print("Accuracy:", accuracy_score(logits.argmax(axis=1)[mask], truth[mask]))  
14    for example_idx in range(n_examples):  
15        print("Example", example_idx)  
16        print("Predicted:", decode_sentence(logits.argmax(axis=1).reshape((batch_size, -1))[example_idx]))  
17        print("Truth:", decode_sentence(truth.reshape((batch_size, -1))[example_idx]))  
18        print("")  
19  
20 check_after_training(3)
```

```
Loss: 2.4119074
```

```
Accuracy: 0.4722598105548038
```

```
Example 0
```

```
Predicted: a living room with a couch chair and a fireplace screen tv #END# #END# #END# #END# #END# #END# #END# #END#  
Truth: a living room with a sofa piano and large flat screen tv #END# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD#
```

```
Example 1
```

```
Predicted: a street flag station is is clearly up #END# night #END# #END# #END# #END# #END# #END# #END# #END#  
Truth: the canadian gas station sign is lit up at night #END# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD#
```

```
Example 2
```

```
Predicted: a teddy of stuffed bears sitting a floor next a pile #END# #END# #END# #END# #END# #END# #END# #END#  
Truth: a group of teddy bears on the floor near a chair #END# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD# #PAD#
```

```
In [40]: 1 # save last graph weights to file!  
2 saver.save(s, get_checkpoint_path())
```

```
Out[40]: 'D:\\\\GoogleDrive\\\\Study\\\\Advanced Machine Learning - Coursera\\\\1.intro-to-dl\\\\weights'
```

Applying model

Here we construct a graph for our final model.

It will work as follows:

- take an image as an input and embed it

- condition lstm on that embedding
- predict the next token given a START input token
- use predicted token as an input at next time step
- iterate until you predict an END token

```
In [41]: 1 class final_model:
2     # CNN encoder
3     encoder, preprocess_for_model = get_cnn_encoder()
4     saver.restore(s, get_checkpoint_path()) # keras applications corrupt our graph, so we restore trained weights
5
6     # containers for current lstm state
7     lstm_c = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="cell")
8     lstm_h = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="hidden")
9
10    # input images
11    input_images = tf.placeholder('float32', [1, IMG_SIZE, IMG_SIZE, 3], name='images')
12
13    # get image embeddings
14    img_embeds = encoder(input_images)
15
16    # initialize Lstm state conditioned on image
17    init_c = init_h = decoder.img_embed_bottleneck_to_h0(decoder.img_embed_to_bottleneck(img_embeds))
18    init_lstm = tf.assign(lstm_c, init_c), tf.assign(lstm_h, init_h)
19
20    # current word index
21    current_word = tf.placeholder('int32', [1], name='current_input')
22
23    # embedding for current word
24    word_embed = decoder.word_embed(current_word)
25
26    # apply Lstm cell, get new Lstm states
27    new_c, new_h = decoder.lstm(word_embed, tf.nn.rnn_cell.LSTMStateTuple(lstm_c, lstm_h))[1]
28
29    # compute logits for next token
30    new_logits = decoder.token_logits(decoder.token_logits_bottleneck(new_h))
31    # compute probabilities for next token
32    new_probs = tf.nn.softmax(new_logits)
33
34    # `one_step` outputs probabilities of next token and updates Lstm hidden state
35    one_step = new_probs, tf.assign(lstm_c, new_c), tf.assign(lstm_h, new_h)
```

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\training\saver.py:1276: checkpoint_exists (from tensorflow.python.training.checkpoint_management) is deprecated and will be removed in a future version.

Instructions for updating:

Use standard file APIs to check for files with this prefix.

INFO:tensorflow:Restoring parameters from D:\GoogleDrive\Study\Advanced Machine Learning - Coursera\1.intro-to-dl\week6\weights

WARNING:tensorflow:Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>> could not be transformed and will be executed as-is. Please report this to the Autograph team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERTOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>>: At attributeError: module 'gast' has no attribute 'Num'

WARNING: Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>> could not be transformed and will be executed as-is. Please report this to the Autograph team. When filing the bug, set the verbosity to 10 (on Linux, `export AUTOGRAPH_VERTOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x0000020601442588>>: AttributeError: module 'gast' has no attribute 'Num'

```
In [42]: 1 # Look at how temperature works for probability distributions
2 # for high temperature we have more uniform distribution
3 _ = np.array([0.5, 0.4, 0.1])
4 for t in [0.01, 0.1, 1, 10, 100]:
5     print(" ".join(map(str, _**(1/t) / np.sum(_**(1/t))))), "with temperature", t)
```

```
0.999999997962965 2.0370359759195462e-10 1.2676505999700117e-70 with temperature 0.01
0.9030370433250645 0.09696286420394223 9.247099323648666e-08 with temperature 0.1
0.5 0.4 0.1 with temperature 1
0.35344772639219624 0.34564811360592396 0.3009041600018798 with temperature 10
0.33536728048099185 0.33461976434857876 0.3300129551704294 with temperature 100
```

```
In [43]: 1 # this is an actual prediction loop
2 def generate_caption(image, t=1, sample=False, max_len=20):
3     """
4         Generate caption for given image.
5         if `sample` is True, we will sample next token from predicted probability distribution.
6         `t` is a temperature during that sampling,
7             higher `t` causes more uniform-like distribution = more chaos.
8     """
9     # condition Lstm on the image
10    s.run(final_model.init_lstm,
11          {final_model.input_images: [image]})
```

current caption
start with only START token
caption = [vocab[START]]

for _ in range(max_len):
 next_word_probs = s.run(final_model.one_step,
 {final_model.current_word: [caption[-1]]})[0]
 next_word_probs = next_word_probs.ravel()

apply temperature
next_word_probs = next_word_probs**(1/t) / np.sum(next_word_probs**(1/t))

if sample:
 next_word = np.random.choice(range(len(vocab)), p=next_word_probs)
else:
 next_word = np.argmax(next_word_probs)

caption.append(next_word)
if next_word == vocab[END]:
 break

return list(map(vocab_inverse.get, caption))

```
In [44]: 1 # Look at validation prediction example
2 def apply_model_to_image_raw_bytes(raw):
3     img = utils.decode_image_from_buf(raw)
4     fig = plt.figure(figsize=(7, 7))
5     plt.grid('off')
6     plt.axis('off')
7     plt.imshow(img)
8     img = utils.crop_and_preprocess(img, (IMG_SIZE, IMG_SIZE), final_model.preprocess_for_model)
9     print(' '.join(generate_caption(img)[1:-1]))
10    plt.show()
11
12 def show_valid_example(val_img_fns, example_idx=0):
13     zf = zipfile.ZipFile("val2014_sample.zip")
14     all_files = set(val_img_fns)
15     found_files = list(filter(lambda x: x.filename.rsplit("/")[-1] in all_files, zf.filelist))
16     example = found_files[example_idx]
17     apply_model_to_image_raw_bytes(zf.read(example))
18
19 show_valid_example(val_img_fns, example_idx=100)
```

a baseball player swinging a bat at a ball



In [45]:

```
1 # sample more images from validation
2 for idx in np.random.choice(range(len(zipfile.ZipFile("val2014_sample.zip").filelist) - 1), 10):
3     show_valid_example(val_img_fns, example_idx=idx)
4     time.sleep(1)
```

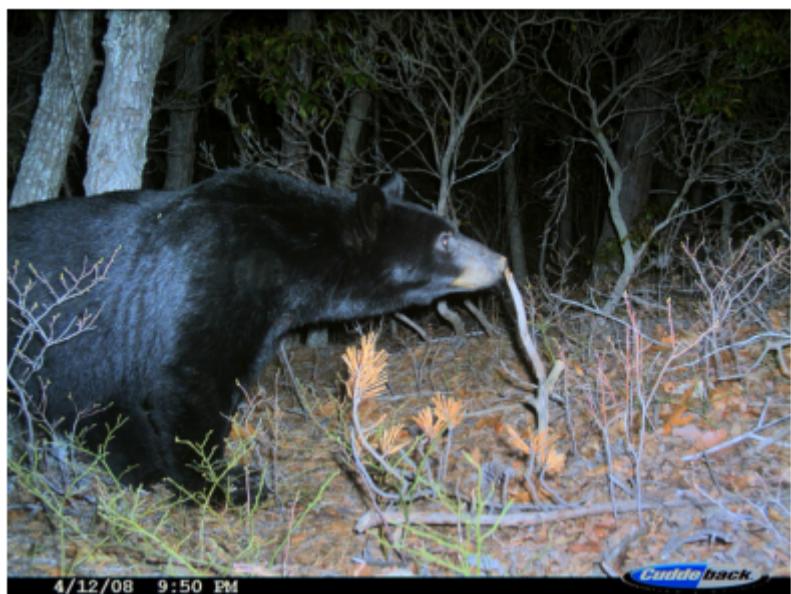
a plate with a sandwich and french fries on it



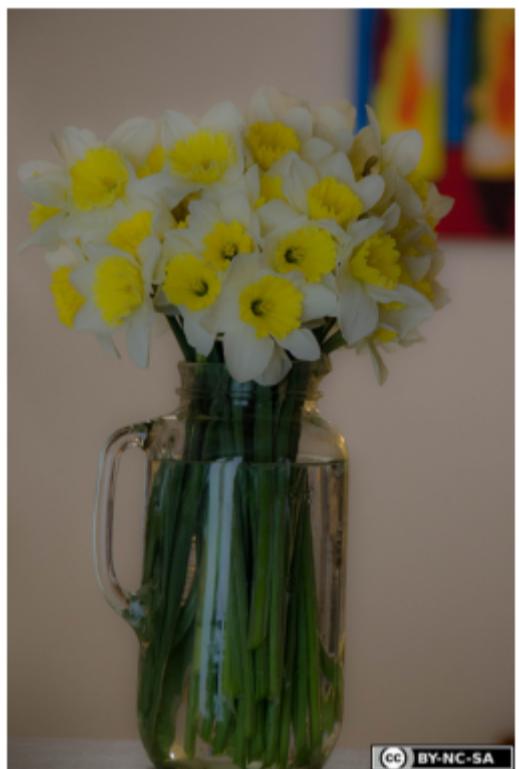
a man and a woman playing a video game



a black bear is standing in the grass



a vase with a flower arrangement in it



a chocolate cake with a fork and a fork



a cat laying on top of a white and black cat



a person on skis on a snowy slope



a stop sign with a sign on it



a group of people walking down a street with umbrellas



a table with a table with a table and a vase of food on it



You can download any image from the Internet and apply your model to it!

```
In [46]: 1 download_utils.download_file(  
2     "http://www.bijouxandbits.com/wp-content/uploads/2016/06/portal-cake-10.jpg",  
3     "portal-cake-10.jpg"  
4 )
```

portal-cake-10.jpg: 100% 273k/273k [00:01<00:00, 242kB/s]

```
In [47]: 1 apply_model_to_image_raw_bytes(open("portal-cake-10.jpg", "rb").read())
```

a cake with a fork and fork on a plate



Now it's time to find 10 examples where your model works good and 10 examples where it fails!

You can use images from validation set as follows:

```
show_valid_example(val_img_fns, example_idx=...)
```

You can use images from the Internet as follows:

```
! wget ...  
apply_model_to_image_raw_bytes(open("...", "rb").read())
```

If you use these functions, the output will be embedded into your notebook and will be visible during peer review!

When you're done, download your noteboook using "File" -> "Download as" -> "Notebook" and prepare that file for peer review!

In [48]:

```
1 # The 10 good
2 good_idx = [355, 203, 87, 24, 373, 340, 29, 214, 363, 140]
3 for idx in good_idx:
4     print("Index:", idx)
5     show_valid_example(val_img_fns, example_idx=idx)
6     time.sleep(1)
```

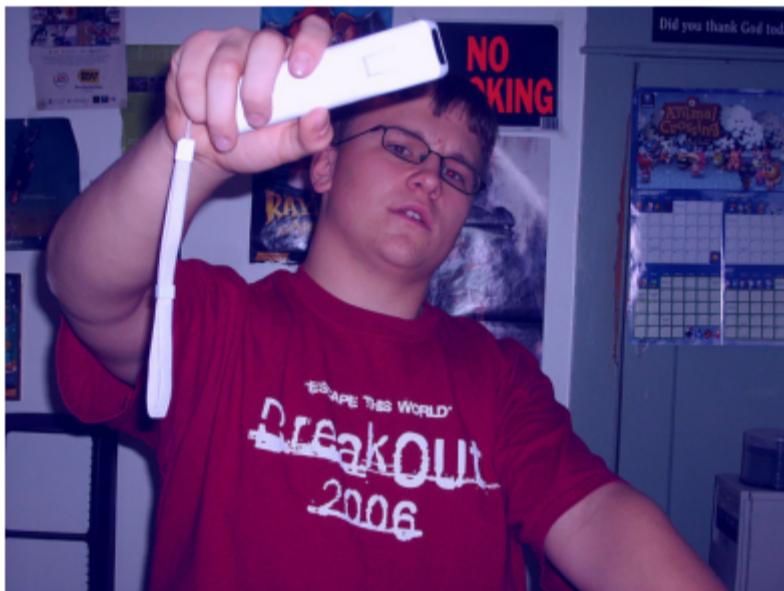
Index: 355

a stove with a pizza cutter on top of it



Index: 203

a man is holding a cell phone to his ear



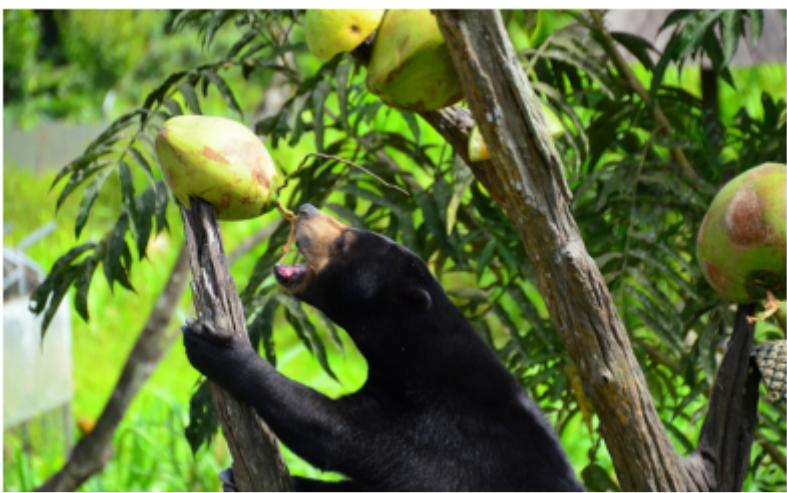
Index: 87

a woman sitting on a skateboard in the middle of a room



Index: 24

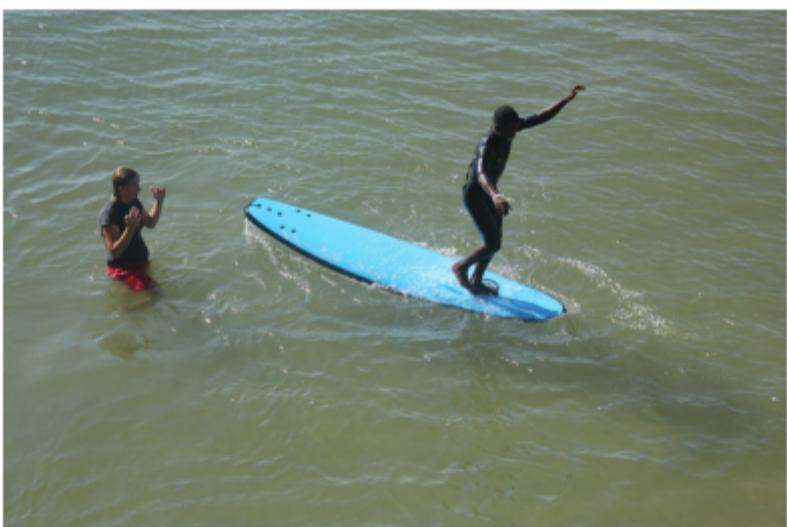
a black bear is sitting on a tree branch



Index: 373
a cat sitting on a toilet bowl in a bathroom



Index: 340
a man is surfing on a wave in the ocean



Index: 29
a bunch of vegetables are on a cutting board



Index: 214
a man is kicking a soccer ball in a field



Index: 363
a bus that is parked in the street



Index: 140
a man swinging a tennis racquet on a tennis court



In [49]:

```
1 # The 10 bad
2 bad_idx = [177, 133, 190 ,312, 247, 310, 293, 39, 107, 78]
3 for idx in bad_idx:
4     print("Index:", idx)
5     show_valid_example(val_img_fns, example_idx=idx)
6     time.sleep(1)
```

Index: 177

a dog is standing in a field with a dog



Index: 133

a room with a tv and a tv



Index: 190

a table with a bunch of fruits and vegetables



Index: 312

a group of people standing around a motorcycle



Index: 247
a stop sign with a sign on it



Index: 310
a woman in a dress playing a game with a wii controller



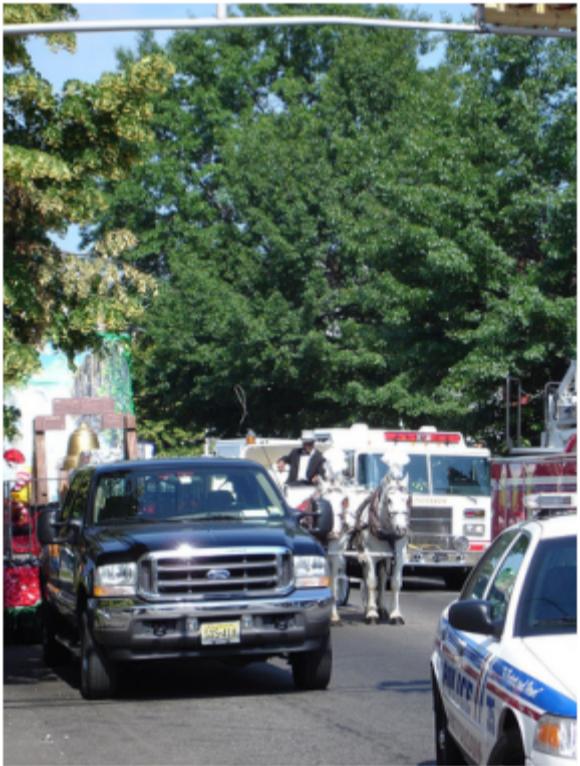
Index: 293
a woman in a wheelchair holding a tennis racket



Index: 39
a dog sitting on a chair next to a laptop



Index: 107
a truck driving down a street with a sign on it



Index: 78
a person holding a cell phone to a window



In [50]: 1 *### YOUR EXAMPLES HERE ###*

That's it!

Congratulations, you've trained your image captioning model and now can produce captions for any picture from the Internet!

In []:

1

