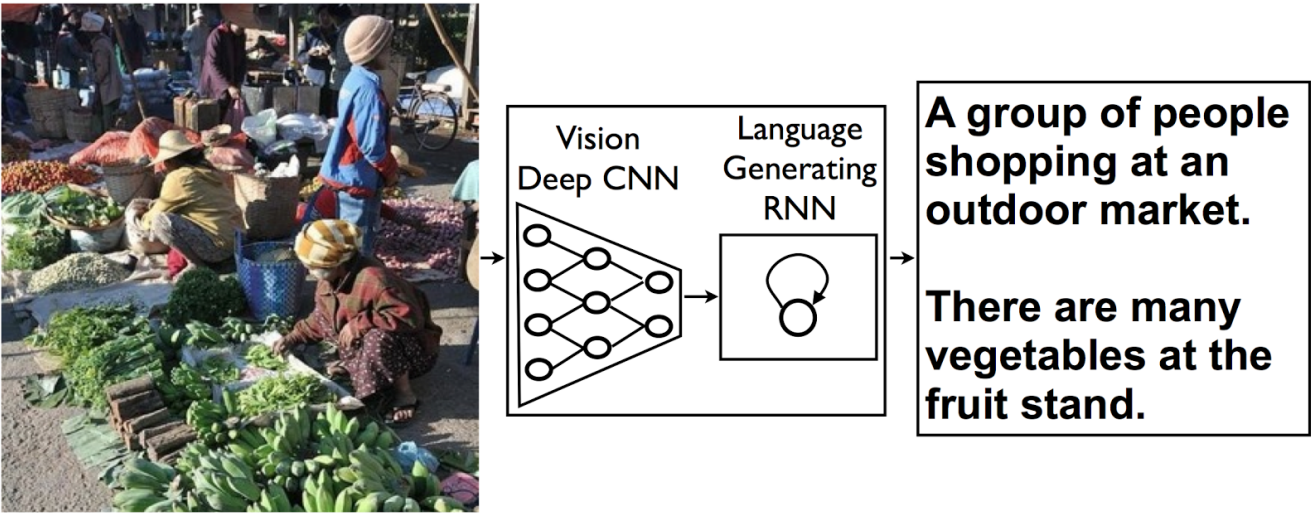```
In [1]:    1  # set tf 1.x for colab
           2  %tensorflow_version 1.x
```

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

# 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!



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

```
In [3]:    1  download_utils.link_all_keras_resources()
```

```
In [4]:    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
```

```
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:516: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_qint8 = np.dtype([("qint8", np.int8, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:517: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_quint8 = np.dtype([("quint8", np.uint8, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:518: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_qint16 = np.dtype([("qint16", np.int16, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:519: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_quint16 = np.dtype([("quint16", np.uint16, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:520: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  _np_qint32 = np.dtype([("qint32", np.int32, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorflow\python\framework\dtypes.py:525: FutureWarning: Passing (type,
1) or '1type' as a synonym of type is deprecated; in a future version of numpy, it will be understood as (type, (1,)) / '(1,)type'.
  np_resource = np.dtype([("resource", np.ubyte, 1)])
C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\tensorboard\compat\tensorflow_stub\dtypes.py:541: FutureWarning: Passing
```

# Prepare the storage for model checkpoints

```
In [5]:    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 trainig 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] + "/colab"
          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:\Google Drive\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 [6]:    1  grader = grading.Grader(assignment_key="NEDBg6CgEee8nQ6uE8a7OA",
           2                          all_parts=["19Wpv", "uJh73", "yiJkt", "rbpnH", "E2OIL", "YJR7z"])
```

```
In [7]:    1  # token expires every 30 min
           2  COURSERA_TOKEN = 'gu6esl2ZRKHnD8sv'
           3  COURSERA_EMAIL = 'lxwvictor@gmail.com'
```

## Download data

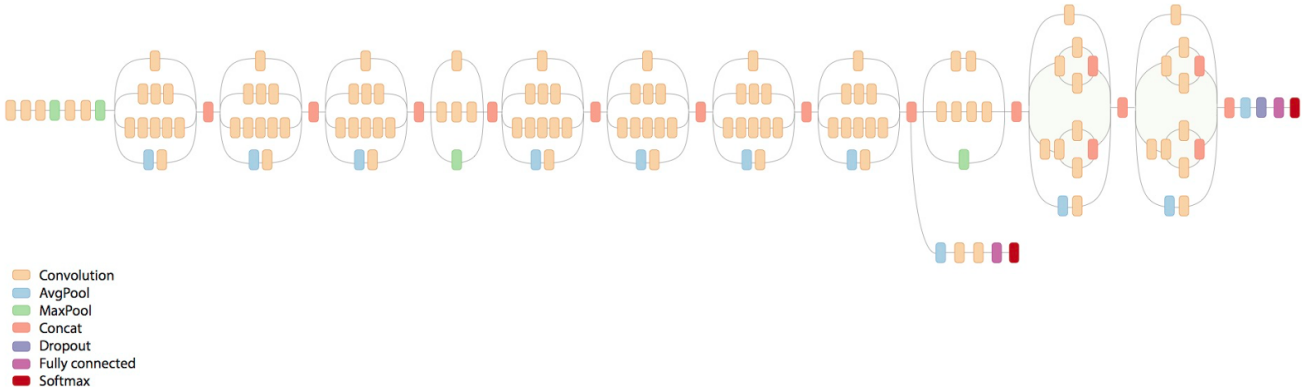Takes 10 hours and 20 GB. We've downloaded necessary files for you.

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 [8]:    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 [9]:    1  IMG_SIZE = 299
```

```
In [10]:   1  # we take the last hidden layer of IncetionV3 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:

```python
# 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")
```

```python
In [11]:
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
```

```python
In [12]:
1  # check prepared samples of images
2  list(filter(lambda x: x.endswith("_sample.zip"), os.listdir(".")))
```

```
Out[12]: ['train2014_sample.zip', 'val2014_sample.zip']
```
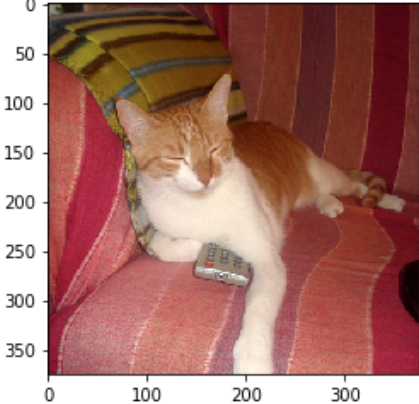
## Extract captions for images

```python
In [13]:
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  train_captions = get_captions_for_fns(train_img_fns, "captions_train-val2014.zip",
13                                         "annotations/captions_train2014.json")
14
15  val_captions = 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(train_captions))
20  print(len(val_img_fns), len(val_captions))
```

```
82783 82783
40504 40504
```

```
In [14]:    1  # look at training example (each has 5 captions)
            2  def show_trainig_example(train_img_fns, train_captions, 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, train_captions))
            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, train_captions, example_idx=142)
```



## Prepare captions for training

```
In [15]:    1  # preview captions data
            2  train_captions[:2]
```

```
Out[15]: [['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 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']]
```

```python
In [16]:    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
```

```python
In [17]:    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
```

```python
In [18]:    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.

```python
In [19]:    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      return matrix
```

```python
In [20]:    1  batch_captions_to_matrix(train_captions_indexed[1], vocab[PAD], 15).shape
```

```
Out[20]:  (5, 11)
```

```python
In [21]:    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))
```

```python
In [22]:    1  # you can make submission with answers so far to check yourself at this stage
            2  grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

```
Submitted to Coursera platform. See results on assignment page!
```

```python
In [23]:    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
```
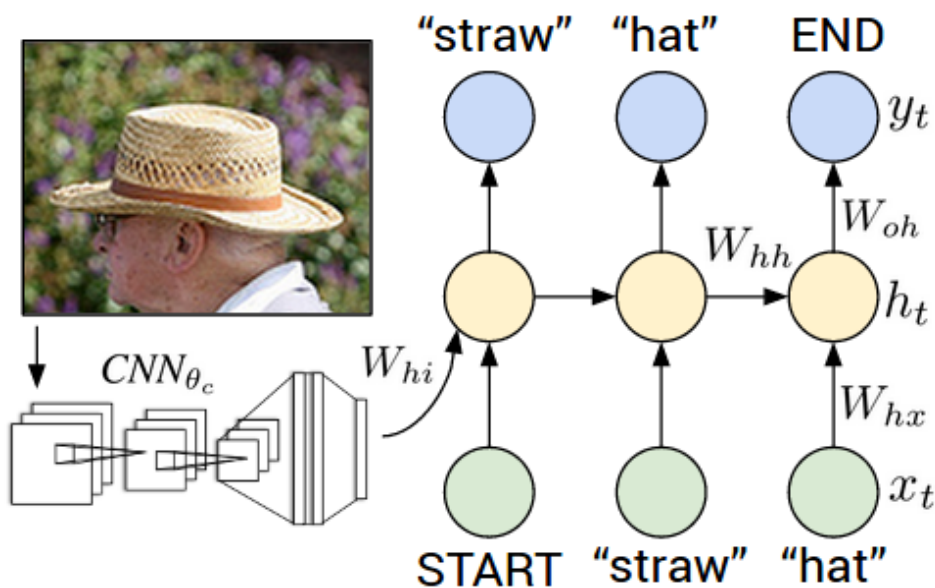
# Training

## 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/)](http://cs.stanford.edu/people/karpathy/)):



```python
In [24]:    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]
```

```python
In [25]:    1  train_img_embeds.shape, len(vocab)
```

```
Out[25]:  ((82783, 2048), 8769)
```

```
In [26]:    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 C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\keras\backend\tensorflow_backend.py:95: The name tf.reset_default_graph is deprecated. Please use tf.compat.v1.reset_default_graph instead.

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\keras\backend\tensorflow_backend.py:98: The name tf.placeholder_with_default is deprecated. Please use tf.compat.v1.placeholder_with_default instead.

WARNING:tensorflow:From C:\Users\Xiaowei\Anaconda3\envs\tfspark\lib\site-packages\keras\backend\tensorflow_backend.py:102: The name tf.get_default_graph is deprecated. Please use tf.compat.v1.get_default_graph instead.
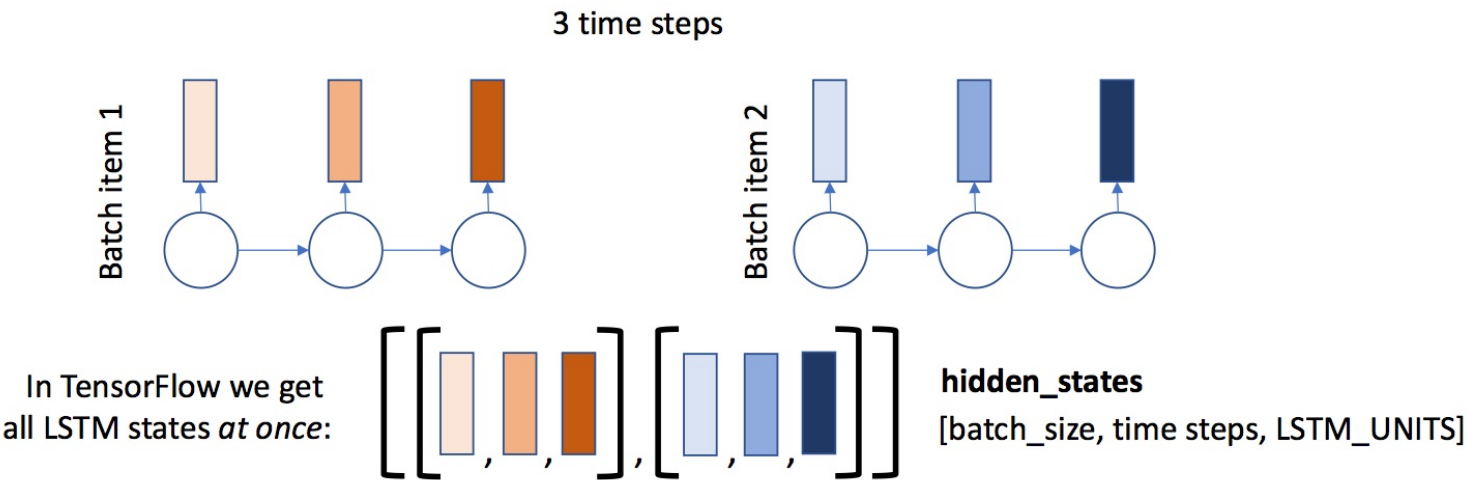
WARNING:tensorflow:From ..\keras_utils.py:75: The name tf.ConfigProto is deprecated. Please use tf.compat.v1.ConfigProto instead.
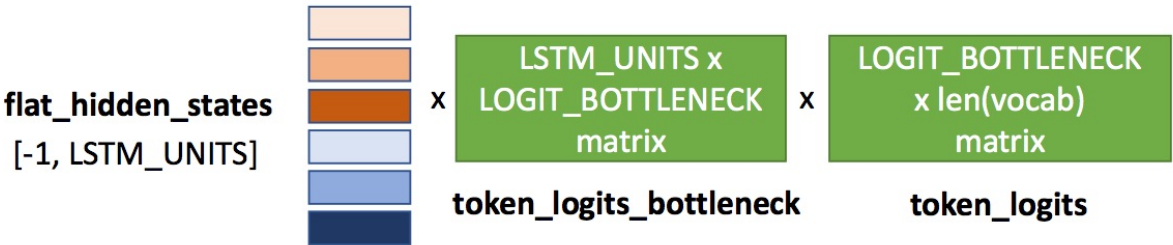
Here we define decoder graph.

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:

In [27]:

```python
# tf.reset_default_graph()
class decoder:
    # [batch_size, IMG_EMBED_SIZE] of CNN image features
    img_embeds = tf.placeholder('float32', [None, IMG_EMBED_SIZE])
    # [batch_size, time steps] of word ids
    sentences = tf.placeholder('int32', [None, None])

    # we use bottleneck here to reduce the number of parameters
    # image embedding -> bottleneck
    img_embed_to_bottleneck = L.Dense(IMG_EMBED_BOTTLENECK,
                                      input_shape=(None, IMG_EMBED_SIZE),
                                      activation='elu')
    # image embedding bottleneck -> lstm initial state
    img_embed_bottleneck_to_h0 = L.Dense(LSTM_UNITS,
                                         input_shape=(None, IMG_EMBED_BOTTLENECK),
                                         activation='elu')
    # word -> embedding
    word_embed = L.Embedding(len(vocab), WORD_EMBED_SIZE)
    # lstm cell (from tensorflow)
    lstm = tf.nn.rnn_cell.LSTMCell(LSTM_UNITS)

    # we use bottleneck here to reduce model complexity
    # lstm output -> logits bottleneck
    token_logits_bottleneck = L.Dense(LOGIT_BOTTLENECK,
                                      input_shape=(None, LSTM_UNITS),
                                      activation="elu")
    # logits bottleneck -> logits for next token prediction
    token_logits = L.Dense(len(vocab),
                           input_shape=(None, LOGIT_BOTTLENECK))

    # initial lstm cell state of shape (None, LSTM_UNITS),
    # we need to condition it on `img_embeds` placeholder.
    ### YOUR CODE HERE ###
    c0 = h0 = img_embed_bottleneck_to_h0(img_embed_to_bottleneck(img_embeds))

    # embed all tokens but the last for lstm input,
    # remember that L.Embedding is callable,
    # use `sentences` placeholder as input.
    ### YOUR CODE HERE ###
    word_embeds = word_embed(sentences[:, :-1])

    # during training we use ground truth tokens `word_embeds` as context for next token prediction.
    # that means that we know all the inputs for our lstm and can get
    # all the hidden states with one tensorflow operation (tf.nn.dynamic_rnn).
    # `hidden_states` has a shape of [batch_size, time steps, LSTM_UNITS].
    hidden_states, _ = tf.nn.dynamic_rnn(lstm, word_embeds,
                                         initial_state=tf.nn.rnn_cell.LSTMStateTuple(c0, h0))

    # now we need to calculate token logits for all the hidden states

    # first, we reshape `hidden_states` to [-1, LSTM_UNITS]
    ### YOUR CODE HERE ###
    flat_hidden_states = tf.reshape(hidden_states, (-1, LSTM_UNITS))

    # then, we calculate logits for next tokens using `token_logits_bottleneck` and `token_logits` layers
    ### YOUR CODE HERE ###
    flat_token_logits = token_logits(token_logits_bottleneck(flat_hidden_states))

    # then, we flatten the ground truth token ids.
    # remember, that we predict next tokens for each time step,
    # use `sentences` placeholder.
    ### YOUR CODE HERE ###
    flat_ground_truth = tf.reshape(sentences[:, 1:], (-1,))

    # we need to know where we have real tokens (not padding) in `flat_ground_truth`,
    # we don't want to propagate the loss for padded output tokens,
    # fill `flat_loss_mask` with 1.0 for real tokens (not pad_idx) and 0.0 otherwise.
    ### YOUR CODE HERE ###
    flat_loss_mask = tf.not_equal(flat_ground_truth, vocab[PAD])

    # compute cross-entropy between `flat_ground_truth` and `flat_token_logits` predicted by lstm
    xent = tf.nn.sparse_softmax_cross_entropy_with_logits(
        labels=flat_ground_truth,
        logits=flat_token_logits
    )

    # compute average `xent` over tokens with nonzero `flat_loss_mask`.
    # we don't want to account misclassification of PAD tokens, because that doesn't make sense,
    # we have PAD tokens for batching purposes only!
    ### YOUR CODE HERE ###
    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 V
arianceScaling.__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\initializers.py:119: call
ing 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-27-389ff8a7bcc0>:20: LSTMCell.__init__ (from tensorflow.python.ops.rnn_cell_impl) is depreca
ted 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-27-389ff8a7bcc0>: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: calli
ng 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 0x000001FB326D4548
>> could not be transformed and will be executed as-is. Please report this to the AutgoGraph team. When filing the bug, set the ver

```
bosity to 10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call
of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>>: AttributeError: module 'gast' has no attribute 'Nu
m'
WARNING: Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>> could n
ot be transformed and will be executed as-is. Please report this to the AutgoGraph team. When filing the bug, set the verbosity to
10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensor
flow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>>: 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_disp
atch_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
```

```python
In [28]:  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  # intialize all variables
         10  s.run(tf.global_variables_initializer())
```

```python
In [29]:  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))
```

```python
In [30]:  1  # you can make submission with answers so far to check yourself at this stage
          2  grader.submit(COURSERA_EMAIL, COURSERA_TOKEN)
```

Submitted to Coursera platform. See results on assignment page!

## Training loop

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

```python
In [31]:  1  train_captions_indexed = np.array(train_captions_indexed)
          2  val_captions_indexed = np.array(val_captions_indexed)
```

```python
In [32]:  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}
```

```python
In [33]:  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
```

```python
In [34]:  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!

```python
# actual training loop
MAX_LEN = 20  # truncate long captions to speed up training

# to make training reproducible
np.random.seed(42)
random.seed(42)

for epoch in range(n_epochs):

    train_loss = 0
    pbar = tqdm_utils.tqdm_notebook_failsafe(range(n_batches_per_epoch))
    counter = 0
    for _ in pbar:
        train_loss += s.run([decoder.loss, train_step],
                            generate_batch(train_img_embeds,
                                           train_captions_indexed,
                                           batch_size,
                                           MAX_LEN))[0]
        counter += 1
        pbar.set_description("Training loss: %f" % (train_loss / counter))

    train_loss /= n_batches_per_epoch

    val_loss = 0
    for _ in range(n_validation_batches):
        val_loss += s.run(decoder.loss, generate_batch(val_img_embeds,
                                                        val_captions_indexed,
                                                        batch_size,
                                                        MAX_LEN))
    val_loss /= n_validation_batches

    print('Epoch: {}, train loss: {}, val loss: {}'.format(epoch, train_loss, val_loss))

    # save weights after finishing epoch
    saver.save(s, get_checkpoint_path(epoch))

print("Finished!")
```

Error rendering Jupyter widget: missing widget manager


Epoch: 0, train loss: 4.259566727161407, val loss: 3.6712243938446045

Error rendering Jupyter widget: missing widget manager


Epoch: 1, train loss: 3.361799514055252, val loss: 3.16270094871521

Error rendering Jupyter widget: missing widget manager


Epoch: 2, train loss: 2.9963913245201113, val loss: 2.907550594806671

Error rendering Jupyter widget: missing widget manager


Epoch: 3, train loss: 2.844430368900299, val loss: 2.8292782354354857

Error rendering Jupyter widget: missing widget manager


```python
## GRADED PART, DO NOT CHANGE!
# Validation loss
grader.set_answer("YJR7z", grading_utils.test_validation_loss(
    decoder, s, generate_batch, val_img_embeds, val_captions_indexed))
```

Error rendering Jupyter widget: missing widget manager


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

Submitted to Coursera platform. See results on assignment page!

```python
# check that it's learnt something, outputs accuracy of next word prediction (should be around 0.5)
from sklearn.metrics import accuracy_score, log_loss

def decode_sentence(sentence_indices):
    return " ".join(list(map(vocab_inverse.get, sentence_indices)))

def check_after_training(n_examples):
    fd = generate_batch(train_img_embeds, train_captions_indexed, batch_size)
    logits = decoder.flat_token_logits.eval(fd)
    truth = decoder.flat_ground_truth.eval(fd)
    mask = decoder.flat_loss_mask.eval(fd).astype(bool)
    print("Loss:", decoder.loss.eval(fd))
    print("Accuracy:", accuracy_score(logits.argmax(axis=1)[mask], truth[mask]))
    for example_idx in range(n_examples):
        print("Example", example_idx)
        print("Predicted:", decode_sentence(logits.argmax(axis=1).reshape((batch_size, -1))[example_idx]))
        print("Truth:", decode_sentence(truth.reshape((batch_size, -1))[example_idx]))
        print("")

check_after_training(3)
```

```
Loss: 2.3909705
Accuracy: 0.4884979702300406
Example 0
Predicted: a living room with a tv chair and tv tv screen tv #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#

Example 1
Predicted: a street sign station sign is clearly up #END# night #END# #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 #END# a table #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#
```

```python
# save last graph weights to file!
saver.save(s, get_checkpoint_path())
```

Out[39]: 'D:\\Google Drive\\Study\\Advanced Machine Learning - Coursera\\1.intro-to-dl\\week6\\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

```python
In [40]:   class final_model:
               # CNN encoder
               encoder, preprocess_for_model = get_cnn_encoder()
               saver.restore(s, get_checkpoint_path())  # keras applications corrupt our graph, so we restore trained weights

               # containers for current lstm state
               lstm_c = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="cell")
               lstm_h = tf.Variable(tf.zeros([1, LSTM_UNITS]), name="hidden")

               # input images
               input_images = tf.placeholder('float32', [1, IMG_SIZE, IMG_SIZE, 3], name='images')

               # get image embeddings
               img_embeds = encoder(input_images)

               # initialize lstm state conditioned on image
               init_c = init_h = decoder.img_embed_bottleneck_to_h0(decoder.img_embed_to_bottleneck(img_embeds))
               init_lstm = tf.assign(lstm_c, init_c), tf.assign(lstm_h, init_h)

               # current word index
               current_word = tf.placeholder('int32', [1], name='current_input')

               # embedding for current word
               word_embed = decoder.word_embed(current_word)

               # apply lstm cell, get new lstm states
               new_c, new_h = decoder.lstm(word_embed, tf.nn.rnn_cell.LSTMStateTuple(lstm_c, lstm_h))[1]

               # compute logits for next token
               new_logits = decoder.token_logits(decoder.token_logits_bottleneck(new_h))
               # compute probabilities for next token
               new_probs = tf.nn.softmax(new_logits)

               # `one_step` outputs probabilities of next token and updates lstm hidden state
               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: checkpo
int_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:\Google Drive\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 0x000001FB326D4548
>> could not be transformed and will be executed as-is. Please report this to the AutgoGraph team. When filing the bug, set the ver
bosity to 10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call
of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>>: AttributeError: module 'gast' has no attribute 'Nu
m'
WARNING: Entity <bound method LSTMCell.call of <tensorflow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>> could n
ot be transformed and will be executed as-is. Please report this to the AutgoGraph team. When filing the bug, set the verbosity to
10 (on Linux, `export AUTOGRAPH_VERBOSITY=10`) and attach the full output. Cause: converting <bound method LSTMCell.call of <tensor
flow.python.ops.rnn_cell_impl.LSTMCell object at 0x000001FB326D4548>>: AttributeError: module 'gast' has no attribute 'Num'
```

```python
In [41]:   # look at how temperature works for probability distributions
           # for high temperature we have more uniform distribution
           _ = np.array([0.5, 0.4, 0.1])
           for t in [0.01, 0.1, 1, 10, 100]:
               print(" ".join(map(str, _**(1/t) / np.sum(_**(1/t)))), "with temperature", t)
```

```
0.9999999997962965 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
```

```python
In [42]:   # this is an actual prediction loop
           def generate_caption(image, t=1, sample=False, max_len=20):
               """
               Generate caption for given image.
               if `sample` is True, we will sample next token from predicted probability distribution.
               `t` is a temperature during that sampling,
                   higher `t` causes more uniform-like distribution = more chaos.
               """
               # condition lstm on the image
               s.run(final_model.init_lstm,
                     {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))
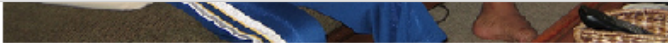```

```
In [43]:    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)
```
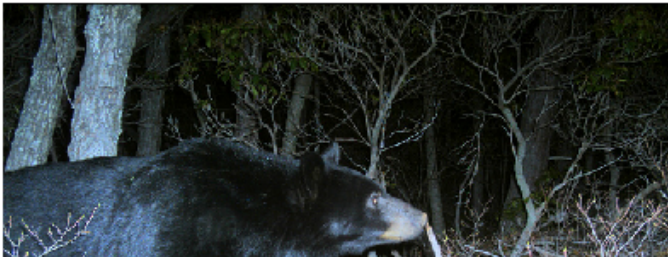
a baseball player swinging a bat at a ball



```
In [44]:    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 black bear is standing in a field



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

```
In [ ]:    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   )
```

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

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



```
Index: 140
a man is playing tennis on a tennis court
```



In [49]:
```python
# The 10 bad
bad_idx = [177, 133, 190 ,312, 247, 310, 293, 39, 107, 78]
for idx in bad_idx:
    print("Index:", idx)
    show_valid_example(val_img_fns, example_idx=idx)
    time.sleep(1)
```



```
Index: 107
a bus that is driving down a street
```



In [ ]:
```python
### 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 [ ]: