

```
! shred -u setup_google_colab.py
! wget https://raw.githubusercontent.com/hse-aml/intro-to-dl/master/setup_google_colab.py
import setup_google_colab #connect to google_colab
setup_google_colab.setup_week6() # connect to google files

!pip uninstall keras-nightly
!pip uninstall -y tensorflow
!pip install tensorflow==1.15.0
!pip install keras==2.1.6
!pip install h5py==2.10.0

! pip install -U numpy==1.18.5 # for the purpose of compatibility

# set tf 1.x for colab
%tensorflow_version 1.x

*****
inception_v3_weights_tf_dim_ordering_tf_kernels_notop.h5
*****
cifar-10-batches-py.tar.gz
*****
mnist.npz
WARNING: Skipping keras-nightly as it is not installed.
Found existing installation: tensorflow 1.15.0
Uninstalling tensorflow-1.15.0:
  Successfully uninstalled tensorflow-1.15.0
Collecting tensorflow==1.15.0
  Using cached tensorflow-1.15.0-cp37-cp37m-manylinux2010_x86_64.whl (412.3 kB)
Requirement already satisfied: gast==0.2.2 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: google-pasta>=0.1.6 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: opt-einsum>=2.3.2 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: termcolor>=1.1.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: protobuf>=3.6.1 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: astor>=0.6.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: tensorflow-estimator==1.15.1 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: absl-py>=0.7.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: keras-applications>=1.0.8 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: grpcio>=1.8.6 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: wrapt>=1.11.1 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: wheel>=0.26 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: keras-preprocessing>=1.0.5 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: tensorboard<1.16.0,>=1.15.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: numpy<2.0,>=1.16.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: six>=1.10.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: h5py in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: werkzeug>=0.11.15 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: setuptools>=41.0.0 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: markdown>=2.6.8 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: importlib-metadata>=4.4 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: zipp>=0.5 in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: typing-extensions>=3.6.4 in /usr/local/lib/python3.7/dist-packages
Installing collected packages: tensorflow
ERROR: pip's dependency resolver does not currently take into account all the packages in the requirements file, but you have tensorflow 1.15.0 which is required by -r requirements.txt (line 1).
Successfully installed tensorflow-1.15.0
Collecting keras==2.1.6
  Using cached Keras-2.1.6-py2.py3-none-any.whl (339 kB)
Requirement already satisfied: six>=1.9.0 in /usr/local/lib/python3.7/dist-packages
```

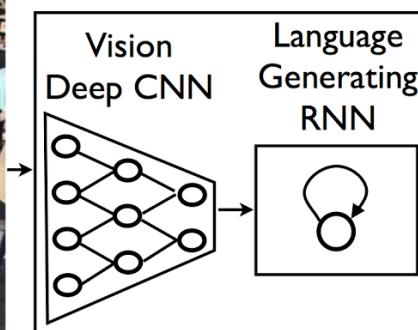
```

Requirement already satisfied: h5py in /usr/local/lib/python3.7/dist-packages
Requirement already satisfied: numpy>=1.9.1 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: pyyaml in /usr/local/lib/python3.7/dist-package
Requirement already satisfied: scipy>=0.14 in /usr/local/lib/python3.7/dist-pa
Installing collected packages: keras
  Attempting uninstall: keras
    Found existing installation: Keras 2.0.6
      Uninstalling Keras-2.0.6:
        Successfully uninstalled Keras-2.0.6
Successfully installed keras-2.1.6
Requirement already satisfied: install in /usr/local/lib/python3.7/dist-packag
Requirement already satisfied: h5py==2.10.0 in /usr/local/lib/python3.7/dist-p
Requirement already satisfied: six in /usr/local/lib/python3.7/dist-packages (
Requirement already satisfied: numpy>=1.7 in /usr/local/lib/python3.7/dist-pac
Requirement already satisfied: numpy==1.18.5 in /usr/local/lib/python3.7/dist-
TensorFlow 1.x selected.

```

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>)

▼ Import stuff

```

import sys
sys.path.append("..")
import grading
import download_utils

```

```
download_utils.link_all_keras_resources()
```

```

import tensorflow as tf
from tensorflow.contrib import keras
import numpy as np
%matplotlib inline
import matplotlib.pyplot as plt
L = keras.layers
K = keras.backend
import utils
import time
import zipfile
import json
from collections import defaultdict
import re
import random
from random import choice
import grading_utils
import os
from keras_utils import reset_tf_session
import tqdm_utils

```

Using TensorFlow backend.

▼ Prepare the storage for model checkpoints

```

# Leave USE_GOOGLE_DRIVE = False if you're running locally!
# We recommend to set USE_GOOGLE_DRIVE = True in Google Colab!
# If set to True, we will mount Google Drive, so that you can restore your checkpo
# and continue training even if your previous Colab session dies.
# If set to True, follow on-screen instructions to access Google Drive (you must ha
USE_GOOGLE_DRIVE = False

def mount_google_drive():
    from google.colab import drive
    mount_directory = "/content/gdrive"
    drive.mount(mount_directory)
    drive_root = mount_directory + "/" + list(filter(lambda x: x[0] != '.', os.listdir(
        return drive_root

CHECKPOINT_ROOT = ""
if USE_GOOGLE_DRIVE:
    CHECKPOINT_ROOT = mount_google_drive() + "/"

def get_checkpoint_path(epoch=None):
    if epoch is None:
        return os.path.abspath(CHECKPOINT_ROOT + "weights")
    else:
        return os.path.abspath(CHECKPOINT_ROOT + "weights_{}".format(epoch))

# example of checkpoint dir
print(get_checkpoint_path(10))

```

```
/content/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

```
grader = grading.Grader(assignment_key="NEDBg6CgEee8nQ6uE8a7OA",
                        all_parts=[ "19Wpv", "uJh73", "yiJkt", "rbpnH", "E2OIL", "YJ
```

```
# token expires every 30 min
COURSERA_TOKEN = '''### YOUR TOKEN HERE
COURSERA_EMAIL = '''### YOUR EMAIL HERE
```

▼ Download data

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

Relevant links (just in case):

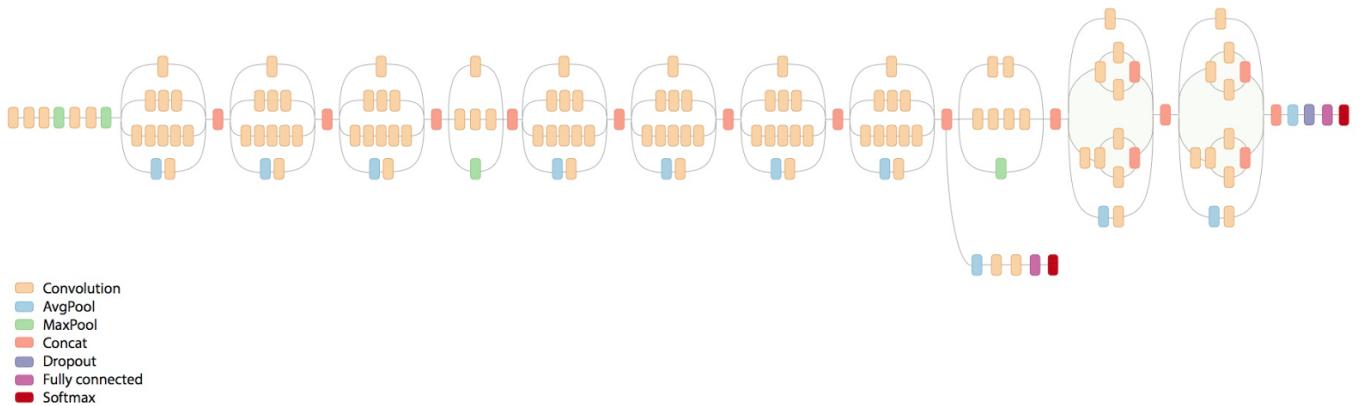
- train images <http://msvocds.blob.core.windows.net/coco2014/train2014.zip>
- validation images <http://msvocds.blob.core.windows.net/coco2014/val2014.zip>
- captions for both train and validation http://msvocds.blob.core.windows.net/annotations-1-0-3/captions_train-val2014.zip

```
# we downloaded them for you, just link them here
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>) and extract its last hidden layer as an embedding:



```
IMG_SIZE = 299
```

```
# we take the last hidden layer of InceptionV3 as an image embedding
def get_cnn_encoder():
    K.set_learning_phase(False)
    model = keras.applications.InceptionV3(include_top=False)
    preprocess_for_model = keras.applications.inception_v3.preprocess_input

    model = keras.models.Model(model.inputs, keras.layers.GlobalAveragePooling2D())
    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")

```

```

# load prepared embeddings
train_img_embeds = utils.read_pickle("train_img_embeds.pickle")
train_img_fns = utils.read_pickle("train_img_fns.pickle")
val_img_embeds = utils.read_pickle("val_img_embeds.pickle")
val_img_fns = utils.read_pickle("val_img_fns.pickle")
# check shapes
print(train_img_embeds.shape, len(train_img_fns))
print(val_img_embeds.shape, len(val_img_fns))

```

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

```

# check prepared samples of images
list(filter(lambda x: x.endswith("_sample.zip"), os.listdir(".")))

['val2014_sample.zip', 'train2014_sample.zip']

```

▼ Extract captions for images

```

# extract captions from zip
def get_captions_for_fns(fns, zip_fn, zip_json_path):
    zf = zipfile.ZipFile(zip_fn)
    j = json.loads(zf.read(zip_json_path).decode("utf8"))
    id_to_fn = {img["id"]: img["file_name"] for img in j["images"]}
    fn_to_caps = defaultdict(list)
    for cap in j['annotations']:
        fn_to_caps[id_to_fn[cap['image_id']]].append(cap['caption'])
    fn_to_caps = dict(fn_to_caps)
    return list(map(lambda x: fn_to_caps[x], fns))

train_captions = get_captions_for_fns(train_img_fns, "captions_train-val2014.zip",
                                       "annotations/captions_train2014.json")

val_captions = get_captions_for_fns(val_img_fns, "captions_train-val2014.zip",
                                       "annotations/captions_val2014.json")

```

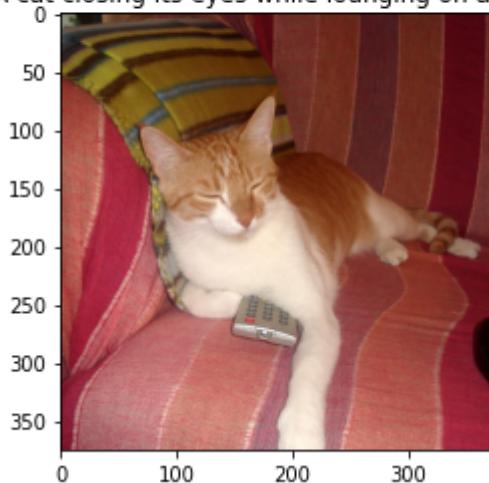
```
# check shape
print(len(train_img_fns), len(train_captions))
print(len(val_img_fns), len(val_captions))

82783 82783
40504 40504
```

```
# look at training example (each has 5 captions)
def show_trainig_example(train_img_fns, train_captions, example_idx=0):
    """
    You can change example_idx and see different images
    """
    zf = zipfile.ZipFile("train2014_sample.zip")
    captions_by_file = dict(zip(train_img_fns, train_captions))
    all_files = set(train_img_fns)
    found_files = list(filter(lambda x: x.filename.rsplit("/")[-1] in all_files, zf))
    example = found_files[example_idx]
    img = utils.decode_image_from_buf(zf.read(example))
    plt.imshow(utils.image_center_crop(img))
    plt.title("\n".join(captions_by_file[example.filename.rsplit("/")[-1]]))
    plt.show()

show_trainig_example(train_img_fns, train_captions, 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

```
# preview captions data
train_captions[:2]

[['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']]
```

```
from collections import Counter  
# special tokens  
PAD = "#PAD#"  
UNK = "#UNK#"  
START = "#START#"  
END = "#END#"  
  
# split sentence into tokens (split into lowercased words)  
def split_sentence(sentence):  
    return list(filter(lambda x: len(x) > 0, re.split('\w+', sentence.lower())))  
  
def generate_vocabulary(train_captions):  
    """  
        Return {token: index} for all train tokens (words) that occur 5 times or more,  
        `index` should be from 0 to N, where N is a number of unique tokens in the  
        Use `split_sentence` function to split sentence into tokens.  
        Also, add PAD (for batch padding), UNK (unknown, out of vocabulary),  
        START (start of sentence) and END (end of sentence) tokens into the vocabul  
    """  
  
    # vocab = ### YOUR CODE HERE ###  
  
    words = [sentence for caption in train_captions for sentence in caption]  
    words = split_sentence(' '.join(words))  
  
    counter = Counter(words)  
  
    vocab = [token for token, count in counter.items() if count>=5]  
    vocab.extend([PAD, UNK, START, END])  
  
    return {token: index for index, token in enumerate(sorted(vocab))}  
  
def caption_tokens_to_indices(captions, vocab):  
    """  
        `captions` argument is an array of arrays:  
        [  
            [  
                "image1 caption1",  
                "image1 caption2",  
                ...  
            ],  
            [  
                "image2 caption1",  
                "image2 caption2",  
                ...  
            ],  
            ...  
        ]  
    """
```

```

    ...
]

Use `split_sentence` function to split sentence into tokens.
Replace all tokens with vocabulary indices, use UNK for unknown words (out of v
Add START and END tokens to start and end of each sentence respectively.
For the example above you should produce the following:
[
    [
        [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
        [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
        ...
    ],
    ...
]
"""

#YOUR OWN CODE
vocabing_words = lambda words: [vocab[START]] + [vocab[w] if w in vocab else vo
ing_words = lambda captions: [vocabing_words(split_sentence(caption.lower())) f
res = [ing_words(caption_list) for caption_list in captions]
return res

```

```
# prepare vocabulary
vocab = generate_vocabulary(train_captions)
vocab_inverse = {idx: w for w, idx in vocab.items()}
print(len(vocab))
```

8769

```
# replace tokens with indices
train_captions_indexed = caption_tokens_to_indices(train_captions, vocab)
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.

```
# we will use this during training
def batch_captions_to_matrix(batch_captions, pad_idx, max_len=None):
    """
    `batch_captions` is an array of arrays:
    [
        [vocab[START], ..., vocab[END]],
        [vocab[START], ..., vocab[END]],
        ...
    ]
    Put vocabulary indexed captions into np.array of shape (len(batch_captions), co
        where "columns" is max(map(len, batch_captions)) when max_len is None
        and "columns" = min(max_len, max(map(len, batch_captions))) otherwise.
    
```

```

Add padding with pad_idx where necessary.
Input example: [[1, 2, 3], [4, 5]]
Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=None
Output example: np.array([[1, 2], [4, 5]]) if max_len=2
Output example: np.array([[1, 2, 3], [4, 5, pad_idx]]) if max_len=100
Try to use numpy, we need this function to be fast!
"""
#      matrix = ### YOUR CODE HERE ####
columns = max(map(len, batch_captions)) if not max_len else min(max_len, max(ma
matrix = pad_idx * np.ones((len(batch_captions), columns))

for i, caption in enumerate(batch_captions):
    for j, value in enumerate(caption):
        if (j < columns):
            matrix[i][j] = value
        else:
            break
return matrix

```

```

## GRADED PART, DO NOT CHANGE!
# Vocabulary creation
grader.set_answer("19Wpv", grading_utils.test_vocab(vocab, PAD, UNK, START, END))
# Captions indexing
grader.set_answer("uJh73", grading_utils.test_captions_indexing(train_captions_inde
# Captions batching
grader.set_answer("yiJkt", grading_utils.test_captions_batching(batch_captions_to_m

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

```

You used an invalid email or your token may have expired. Please make sure you

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

▼ Training

▼ Define architecture

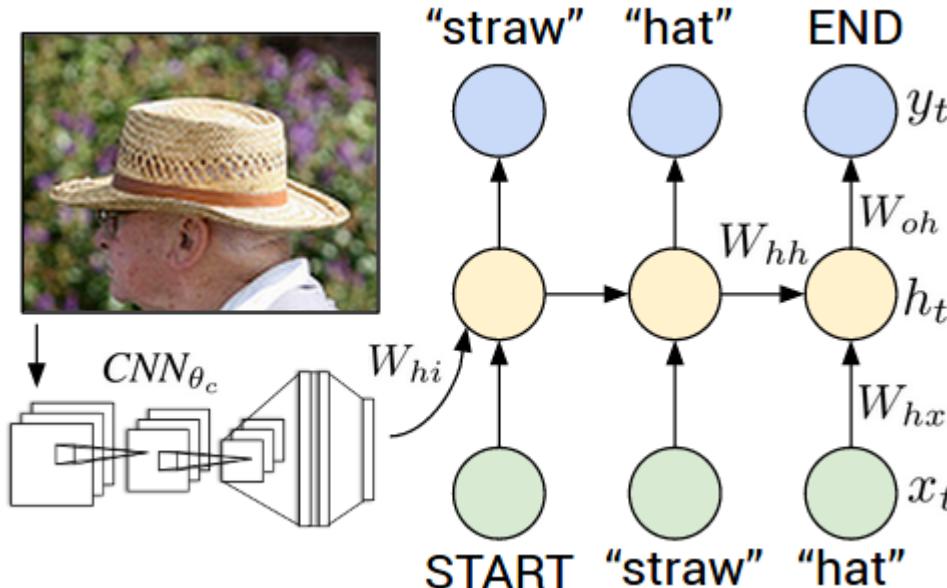
Since our problem is to generate image captions, RNN text generator should be conditioned on image. The idea is to use image features as an initial state for RNN instead of zeros.

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/>):



```
IMG_EMBED_SIZE = train_img_embeds.shape[1]
IMG_EMBED_BOTTLENECK = 120
WORD_EMBED_SIZE = 100
LSTM_UNITS = 300
LOGIT_BOTTLENECK = 120
pad_idx = vocab[PAD]
```

```
# remember to reset your graph if you want to start building it from scratch!
s = reset_tf_session()
tf.set_random_seed(42)

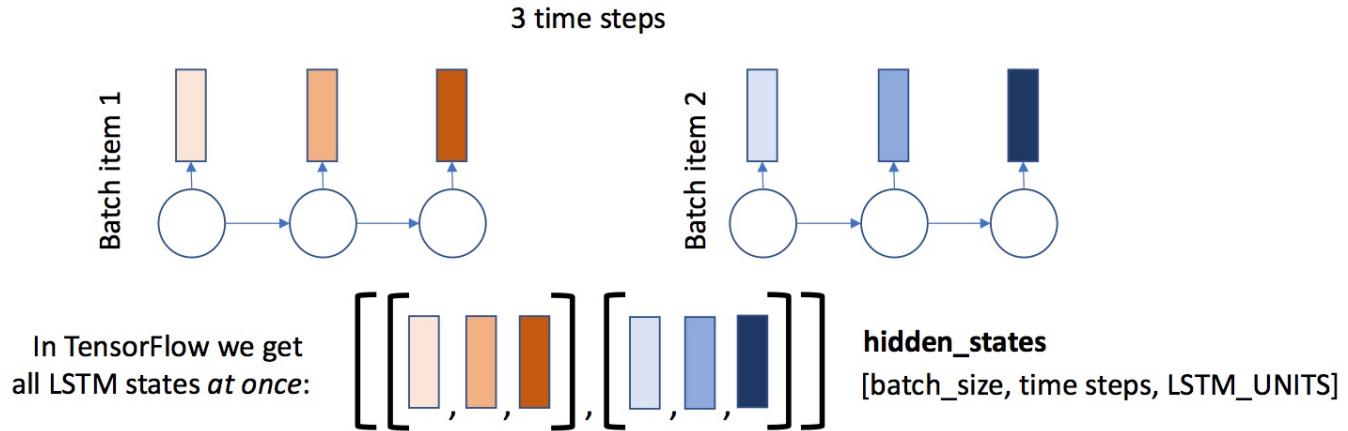
WARNING:tensorflow:From /content/keras_utils.py:68: The name tf.get_default_se
WARNING:tensorflow:From /content/keras_utils.py:75: The name tf.ConfigProto is
WARNING:tensorflow:From /content/keras_utils.py:77: The name tf.InteractiveSes
```

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*:



```
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.
c0 = h0 = img_embed_bottleneck_to_h0(img_embed_to_bottleneck(img_embeds)) ### YOUR CODE HERE ###

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

# during training we use ground truth tokens `word_embeds` as context for next
# 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)

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

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

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

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

# 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
flat_loss_mask = tf.where (tf.equal( flat_ground_truth, tf.constant( pad_idx ) )
                           tf.zeros_like(flat_ground_truth),
                           tf.ones_like(flat_ground_truth))

### YOUR CODE HERE ###

# compute cross-entropy between `flat_ground_truth` and `flat_token_logits` present
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
# we have PAD tokens for batching purposes only!

loss = tf.reduce_mean(tf.boolean_mask(xent,flat_loss_mask)) ### YOUR CODE HERE ###

```

WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/keras/layers/recurrent.py:1055: LSTMCell.__init__(...): The class LSTMCell 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
 WARNING:tensorflow:From <ipython-input-25-1fa841146597>:19: LSTMCell.__init__
 Instructions for updating:

This class is equivalent as `tf.keras.layers.LSTMCell`, and will be replaced by `WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/or`
 Instructions for updating:
 If using Keras pass `*_constraint` arguments to layers.
`WARNING:tensorflow:From <ipython-input-25-1fa841146597>:44: dynamic_rnn` (from
 Instructions for updating:
 Please use ``keras.layers.RNN(cell)``, which is equivalent to this API
`WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/or`
 Instructions for updating:
 Please use ``layer.add_weight`` method instead.
`WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/or`
 Instructions for updating:
 Call initializer instance with the `dtype` argument instead of passing it to the
`WARNING:tensorflow:From <ipython-input-25-1fa841146597>:64: where` (from tensor
 Instructions for updating:
 Use `tf.where` in 2.0, which has the same broadcast rule as `np.where`

```
# define optimizer operation to minimize the loss
optimizer = tf.train.AdamOptimizer(learning_rate=0.001)
train_step = optimizer.minimize(decoder.loss)

# will be used to save/load network weights.
# you need to reset your default graph and define it in the same way to be able to
saver = tf.train.Saver()

# initialize all variables
s.run(tf.global_variables_initializer())
```

```
## GRADED PART, DO NOT CHANGE!
# Decoder shapes test
grader.set_answer("rbpnH", grading_utils.test_decoder_shapes(decoder, IMG_EMBED_SIZE))
# Decoder random loss test
grader.set_answer("E2OIL", grading_utils.test_random_decoder_loss(decoder, IMG_EMBED_SIZE))
```

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

You used an invalid email or your token may have expired. Please make sure you

▼ Training loop

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

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

```
# generate batch via random sampling of images and captions for them,
# we use `max_len` parameter to control the length of the captions (truncating long
def generate_batch(images_embeddings, indexed_captions, batch_size, max_len=None):
    """
    `images_embeddings` is a np.array of shape [number of images, IMG_EMBED_SIZE].
```

```

`indexed_captions` holds 5 vocabulary indexed captions for each image:
[
    [
        [vocab[START], vocab["image1"], vocab["caption1"], vocab[END]],
        [vocab[START], vocab["image1"], vocab["caption2"], vocab[END]],
        ...
    ],
    ...
]
Generate a random batch of size `batch_size`.
Take random images and choose one random caption for each image.
Remember to use `batch_captions_to_matrix` for padding and respect `max_len` pa
Return feed dict {decoder.img_embeds: ..., decoder.sentences: ...}.
"""

batch = np.random.choice(len(images_embeddings), batch_size, replace=False)
batch_captions = [caption[np.random.randint(5)] for caption in indexed_captions

batch_image_embeddings = images_embeddings[batch] #YOUR OWN CODE

batch_captions_matrix = batch_captions_to_matrix(batch_captions, pad_idx, max_l

return {decoder.img_embeds: batch_image_embeddings,
        decoder.sentences: batch_captions_matrix}

batch_size = 64
n_epochs = 11
n_batches_per_epoch = 1000
n_validation_batches = 100 # how many batches are used for validation after each e

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

```

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

```

# 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_l

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

print("Finished!")

```

```

*****
Training loss: 4.260807
Epoch: 0, train loss: 4.2608068609237675, val loss: 3.6993525576591493
*****
Training loss: 3.393728
Epoch: 1, train loss: 3.3937281548976896, val loss: 3.1605656480789186
*****
Training loss: 3.031979
Epoch: 2, train loss: 3.031978526353836, val loss: 2.962694671154022
*****
Training loss: 2.858394
Epoch: 3, train loss: 2.8583938920497896, val loss: 2.857237856388092
*****
Training loss: 2.765346
Epoch: 4, train loss: 2.765346033334732, val loss: 2.7924067163467408
*****
Training loss: 2.701538
Epoch: 5, train loss: 2.701538345813751, val loss: 2.7356078386306764
WARNING:tensorflow:From /tensorflow-1.15.2/python3.7/tensorflow_core/python/tr
Instructions for updating:
Use standard file APIs to delete files with this prefix.
*****
Training loss: 2.642332
Epoch: 6, train loss: 2.6423319716453553, val loss: 2.6959809923171996
*****
Training loss: 2.598755
Epoch: 7, train loss: 2.5987552185058593, val loss: 2.6712168765068056
*****
Training loss: 2.568427
Epoch: 8, train loss: 2.5684266645908354, val loss: 2.6184411668777465
*****
Training loss: 2.532141
Epoch: 9, train loss: 2.5321407120227812, val loss: 2.621663691997528
*****
Training loss: 2.504851

```

```
Epoch: 10, train loss: 2.504851251363754, val loss: 2.596101016998291
Finished!
```

```
## 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))
```

```
*****
```

```
# 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!

```
# check that it's learnt something, outputs accuracy of next word prediction (should be ~0.4)
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: 2.5658107
Accuracy: 0.489501312335958
Example 0
Predicted: a man and a from a a a #END# #END# #END# #END# #END# #END# #END# #END#
Truth: a man pouring wine from #UNK# for patrons #END# #PAD# #PAD# #PAD# #PAD# #PAD#
Example 1
Predicted: a man and to to bike to on #END# #END# #END# #END# #END# #END# #END#
Truth: a man tries out a bicycle powered blender #END# #PAD# #PAD# #PAD# #PAD# #PAD#
Example 2
Predicted: a white up of a toilet toilet toilet #END# #END# #END# #END# #END# #END#
Truth: a close up of a clean white toilet #END# #PAD# #PAD# #PAD# #PAD# #PAD#
```

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

```
'/content/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

```
class final_model:
    # CNN encoder
    encoder, preprocess_for_model = get_cnn_encoder()
    saver.restore(s, get_checkpoint_path()) # keras applications corrupt our graph

    # 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='image')

    # 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)
    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))

    # 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)
```

INFO:tensorflow:Restoring parameters from /content/weights

```
# 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

```
# 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))
```

```
# look at validation prediction example
def apply_model_to_image_raw_bytes(raw):
    img = utils.decode_image_from_buf(raw)
    fig = plt.figure(figsize=(7, 7))
```

```

plt.grid('off')
plt.axis('off')
plt.imshow(img)
img = utils.crop_and_preprocess(img, (IMG_SIZE, IMG_SIZE), final_model.preproce
print(' '.join(generate_caption(img)[1:-1]))
plt.show()

def show_valid_example(val_img_fns, example_idx=0):
    zf = zipfile.ZipFile("val2014_sample.zip")
    all_files = set(val_img_fns)
    found_files = list(filter(lambda x: x.filename.rsplit('/')[-1] in all_files, zf
example = found_files[example_idx]
apply_model_to_image_raw_bytes(zf.read(example))

show_valid_example(val_img_fns, example_idx=100)

```

a baseball player swinging a bat at a baseball game

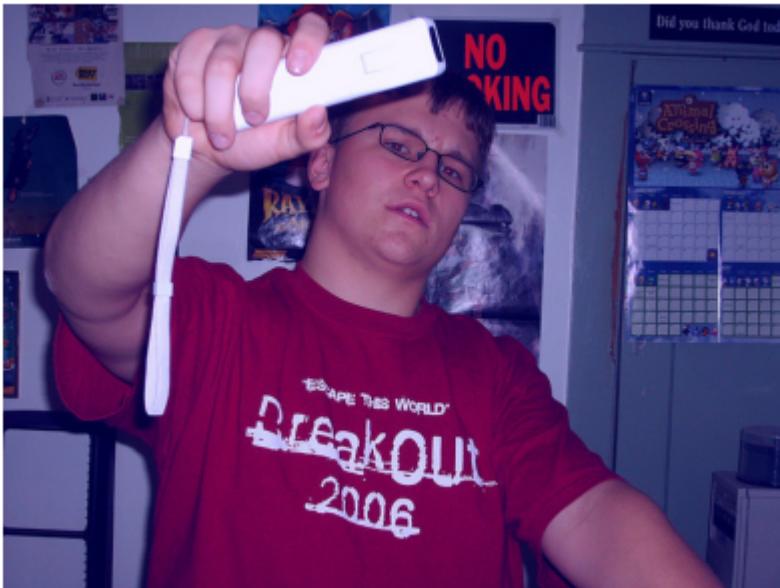


```

# sample more images from validation
for idx in np.random.choice(range(len(zipfile.ZipFile("val2014_sample.zip").fil
    show_valid_example(val_img_fns, example_idx=idx)
    time.sleep(1)

```

a man holding a wii remote in his hand



a sheep standing next to a fence in a field



a street sign on a pole with a street sign



a bowl of fruit with a bunch of fruit on it



a man and woman are playing a game of frisbee



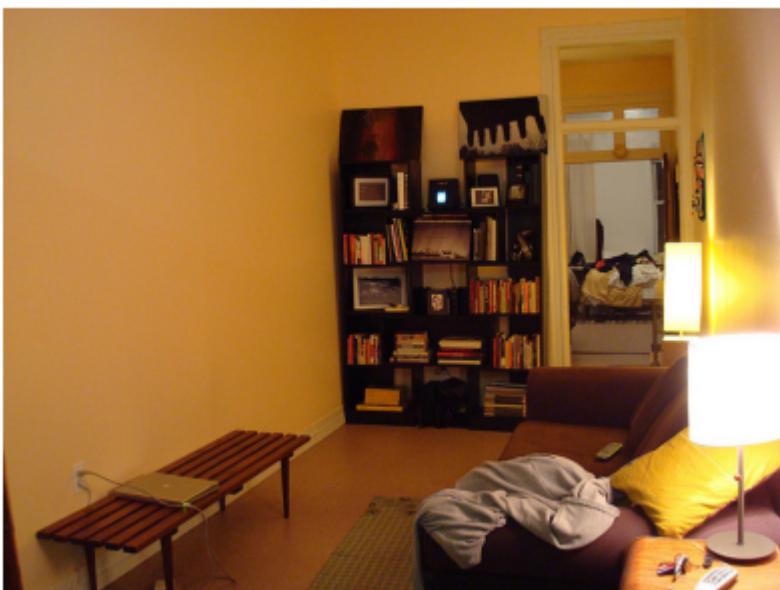
a kitchen with a stove top oven and a stove top oven



a double decker bus is driving down a street



a bedroom with a bed and a bed



a bicycle is parked next to a building



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

```
download_utils.download_file(
    "http://www.bijouxandbits.com/wp-content/uploads/2016/06/portal-cake-10.jpg",
    "portal-cake-10.jpg"
)
```

[portal-cake-10.jpg](#)

```
apply_model_to_image_raw_bytes(open("portal-cake-10.jpg", "rb").read())
```

a cake decorated with a bunch of different types of candies



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!

```
### YOUR EXAMPLES HERE ###
```

```
!wget -O 4.jpg http://www.creatingcomforts.co.uk/blog/wp-content/uploads/2013/12/Ma

--2022-04-25 07:44:08-- http://www.creatingcomforts.co.uk/blog/wp-content/uploads/2013/12/M
Resolving www.creatingcomforts.co.uk (www.creatingcomforts.co.uk)... 89.34.17.
Connecting to www.creatingcomforts.co.uk (www.creatingcomforts.co.uk).|89.34.17|
HTTP request sent, awaiting response... 301 Moved Permanently
Location: https://www.creatingcomforts.co.uk/blog/wp-content/uploads/2013/12/M
--2022-04-25 07:44:10-- https://www.creatingcomforts.co.uk/blog/wp-content/uploads/2013/12/M
Connecting to www.creatingcomforts.co.uk (www.creatingcomforts.co.uk).|89.34.17|
HTTP request sent, awaiting response... 200 OK
Length: 16906 (17K) [image/jpeg]
Saving to: '4.jpg'

4.jpg          100%[=====] 16.51K  --.-KB/s    in 0s

2022-04-25 07:44:10 (213 MB/s) - '4.jpg' saved [16906/16906]
```

```
apply_model_to_image_raw_bytes(open("{}{}.jpg".format(4), "rb").read())
```

a man holding a plate of food with a dog



Image by stockimages www.freigitalphotos.net

That's it!

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