

Image Captioning with RNNs

In this exercise you will implement a vanilla recurrent neural networks and use them it to train a model that can generate novel captions for images.

Refresh Environment and Get Datasets from scripts:

```
# set workspace and get datasets
from google.colab import drive
drive.mount('/content/drive')
FOLDERNAME = "cs231n/assignments/assignment3/"
assert FOLDERNAME is not None, "[!] Enter the foldername."
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

Install h5py

The COCO dataset we will be using is stored in HDF5 format. To load HDF5 files, we will need to install the `h5py` Python package. From the command line, run:

```
pip install h5py
```

If you receive a permissions error, you may need to run the command as root:

```
sudo pip install h5py
```

You can also run commands directly from the Jupyter notebook by prefixing the command with the "!" character:

```
!pip install h5py
```

```
Requirement already satisfied: h5py in /usr/local/lib/python3.11/dist-packages
(3.13.0)
Requirement already satisfied: numpy>=1.19.3 in /usr/local/lib/python3.11/dist-
packages (from h5py) (2.0.2)
```

```
# As usual, a bit of setup
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient,
eval_numerical_gradient_array
from cs231n.rnn_layers import *
from cs231n.captioning_solver import CaptioningSolver
from cs231n.classifiers.rnn import CaptioningRNN
```

```

from cs231n.coco_utils import load_coco_data, sample_coco_minibatch,
decode_captions
from cs231n.image_utils import image_from_url

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-
ipython
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))

```

Microsoft COCO

For this exercise we will use the 2014 release of the [Microsoft COCO dataset](#) which has become the standard testbed for image captioning. The dataset consists of 80,000 training images and 40,000 validation images, each annotated with 5 captions written by workers on Amazon Mechanical Turk.

You should have already downloaded the data by changing to the `cs231n/datasets` directory and running the script `get_assignment3_data.sh`. If you haven't yet done so, run that script now. Warning: the COCO data download is ~1GB.

We have preprocessed the data and extracted features for you already. For all images we have extracted features from the fc7 layer of the VGG-16 network pretrained on ImageNet; these features are stored in the files `train2014_vgg16_fc7.h5` and `val2014_vgg16_fc7.h5` respectively. To cut down on processing time and memory requirements, we have reduced the dimensionality of the features from 4096 to 512; these features can be found in the files `train2014_vgg16_fc7_pca.h5` and `val2014_vgg16_fc7_pca.h5`.

The raw images take up a lot of space (nearly 20GB) so we have not included them in the download. However all images are taken from Flickr, and URLs of the training and validation images are stored in the files `train2014_urls.txt` and `val2014_urls.txt` respectively. This allows you to download images on the fly for visualization. Since images are downloaded on-the-fly, **you must be connected to the internet to view images**.

Dealing with strings is inefficient, so we will work with an encoded version of the captions. Each word is assigned an integer ID, allowing us to represent a caption by a sequence of integers. The mapping between integer IDs and words is in the file `coco2014_vocab.json`, and you can use the function `decode_captions` from the file `cs231n/coco_utils.py` to convert numpy arrays of integer IDs back into strings.

There are a couple special tokens that we add to the vocabulary. We prepend a special `<START>` token and append an `<END>` token to the beginning and end of each caption respectively. Rare words are replaced with a special `<UNK>` token (for "unknown"). In addition, since we want to train with minibatches containing captions of different lengths, we pad short captions with a special `<NULL>` token after the `<END>` token and don't compute loss or gradient for `<NULL>`

tokens. Since they are a bit of a pain, we have taken care of all implementation details around special tokens for you.

You can load all of the MS-COCO data (captions, features, URLs, and vocabulary) using the `load_coco_data` function from the file `cs231n/coco_utils.py`. Run the following cell to do so:

```
# Load COCO data from disk; this returns a dictionary
# we'll work with dimensionality-reduced features for this notebook, but feel
# free to experiment with the original features by changing the flag below.
data = load_coco_data(pca_features=True)

# Print out all the keys and values from the data dictionary
for k, v in data.items():
    if type(v) == np.ndarray:
        print(k, type(v), v.shape, v.dtype)
    else:
        print(k, type(v), len(v))
```

```
base dir
/content/drive/MyDrive/cs231n/assignments/assignment3/cs231n/datasets/coco_captioning
train_captions <class 'numpy.ndarray'> (400135, 17) int32
train_image_idxs <class 'numpy.ndarray'> (400135,) int32
val_captions <class 'numpy.ndarray'> (195954, 17) int32
val_image_idxs <class 'numpy.ndarray'> (195954,) int32
train_features <class 'numpy.ndarray'> (82783, 512) float32
val_features <class 'numpy.ndarray'> (40504, 512) float32
idx_to_word <class 'list'> 1004
word_to_idx <class 'dict'> 1004
train_urls <class 'numpy.ndarray'> (82783,) <U63
val_urls <class 'numpy.ndarray'> (40504,) <U63
```

Look at the data

It is always a good idea to look at examples from the dataset before working with it.

You can use the `sample_coco_minibatch` function from the file `cs231n/coco_utils.py` to sample minibatches of data from the data structure returned from `load_coco_data`. Run the following to sample a small minibatch of training data and show the images and their captions. Running it multiple times and looking at the results helps you to get a sense of the dataset.

Note that we decode the captions using the `decode_captions` function and that we download the images on-the-fly using their Flickr URL, so **you must be connected to the internet to view images**.

```
# Sample a minibatch and show the images and captions
# If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
batch_size = 4

captions, features, urls = sample_coco_minibatch(data, batch_size=batch_size)
for i, (caption, url) in enumerate(zip(captions, urls)):
    img = image_from_url(url)
    if img is None: continue
    plt.imshow(img)
```

```
plt.axis('off')
caption_str = decode_captions(caption, data['idx_to_word'])
plt.title(caption_str)
plt.show()
```

Output hidden; open in <https://colab.research.google.com> to view.

Recurrent Neural Networks

As discussed in lecture, we will use recurrent neural network (RNN) language models for image captioning. The file `cs231n/rnn_layers.py` contains implementations of different layer types that are needed for recurrent neural networks, and the file `cs231n/classifiers/rnn.py` uses these layers to implement an image captioning model.

We will first implement different types of RNN layers in `cs231n/rnn_layers.py`.

Vanilla RNN: step forward

Open the file `cs231n/rnn_layers.py`. This file implements the forward and backward passes for different types of layers that are commonly used in recurrent neural networks.

First implement the function `rnn_step_forward` which implements the forward pass for a single timestep of a vanilla recurrent neural network. After doing so run the following to check your implementation. You should see errors on the order of e-8 or less.

```
N, D, H = 3, 10, 4

x = np.linspace(-0.4, 0.7, num=N*D).reshape(N, D)
prev_h = np.linspace(-0.2, 0.5, num=N*H).reshape(N, H)
wx = np.linspace(-0.1, 0.9, num=D*H).reshape(D, H)
wh = np.linspace(-0.3, 0.7, num=H*H).reshape(H, H)
b = np.linspace(-0.2, 0.4, num=H)

next_h, _ = rnn_step_forward(x, prev_h, wx, wh, b)
expected_next_h = np.asarray([
    [-0.58172089, -0.50182032, -0.41232771, -0.31410098],
    [ 0.66854692,  0.79562378,  0.87755553,  0.92795967],
    [ 0.97934501,  0.99144213,  0.99646691,  0.99854353]])

print('next_h error: ', rel_error(expected_next_h, next_h))
```

next_h error: 6.292421426471037e-09

Vanilla RNN: step backward

In the file `cs231n/rnn_layers.py` implement the `rnn_step_backward` function. After doing so run the following to numerically gradient check your implementation. You should see errors on the order of e-8 or less.

```
from cs231n.rnn_layers import rnn_step_forward, rnn_step_backward
np.random.seed(231)
N, D, H = 4, 5, 6
```

```

x = np.random.randn(N, D)
h = np.random.randn(N, H)
wx = np.random.randn(D, H)
wh = np.random.randn(H, H)
b = np.random.randn(H)

out, cache = rnn_step_forward(x, h, wx, wh, b)

dnext_h = np.random.randn(*out.shape)

fx = lambda x: rnn_step_forward(x, h, wx, wh, b)[0]
fh = lambda prev_h: rnn_step_forward(x, h, wx, wh, b)[0]
fwx = lambda wx: rnn_step_forward(x, h, wx, wh, b)[0]
fwh = lambda wh: rnn_step_forward(x, h, wx, wh, b)[0]
fb = lambda b: rnn_step_forward(x, h, wx, wh, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dnext_h)
dprev_h_num = eval_numerical_gradient_array(fh, h, dnext_h)
dwx_num = eval_numerical_gradient_array(fwx, wx, dnext_h)
dwh_num = eval_numerical_gradient_array(fwh, wh, dnext_h)
db_num = eval_numerical_gradient_array(fb, b, dnext_h)

dx, dprev_h, dwx, dwh, db = rnn_step_backward(dnext_h, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dprev_h error: ', rel_error(dprev_h_num, dprev_h))
print('dwx error: ', rel_error(dwx_num, dwx))
print('dwh error: ', rel_error(dwh_num, dwh))
print('db error: ', rel_error(db_num, db))

```

```

dx error: 2.7795541640745535e-10
dprev_h error: 2.732467428030486e-10
dwx error: 9.709219069305414e-10
dwh error: 5.034262638717296e-10
db error: 1.708752322503098e-11

```

Vanilla RNN: forward

Now that you have implemented the forward and backward passes for a single timestep of a vanilla RNN, you will combine these pieces to implement a RNN that processes an entire sequence of data.

In the file `cs231n/rnn_layers.py`, implement the function `rnn_forward`. This should be implemented using the `rnn_step_forward` function that you defined above. After doing so run the following to check your implementation. You should see errors on the order of `e-7` or less.

```

N, T, D, H = 2, 3, 4, 5

x = np.linspace(-0.1, 0.3, num=N*T*D).reshape(N, T, D)
h0 = np.linspace(-0.3, 0.1, num=N*H).reshape(N, H)
wx = np.linspace(-0.2, 0.4, num=D*H).reshape(D, H)
wh = np.linspace(-0.4, 0.1, num=H*H).reshape(H, H)
b = np.linspace(-0.7, 0.1, num=H)

h, _ = rnn_forward(x, h0, wx, wh, b)

```

```

expected_h = np.asarray([
    [
        [-0.42070749, -0.27279261, -0.11074945, 0.05740409, 0.22236251],
        [-0.39525808, -0.22554661, -0.0409454, 0.14649412, 0.32397316],
        [-0.42305111, -0.24223728, -0.04287027, 0.15997045, 0.35014525],
    ],
    [
        [-0.55857474, -0.39065825, -0.19198182, 0.02378408, 0.23735671],
        [-0.27150199, -0.07088804, 0.13562939, 0.33099728, 0.50158768],
        [-0.51014825, -0.30524429, -0.06755202, 0.17806392, 0.40333043]]])
print('h error: ', rel_error(expected_h, h))

```

h error: 7.728466151011529e-08

Vanilla RNN: backward

In the file `cs231n/rnn_layers.py`, implement the backward pass for a vanilla RNN in the function `rnn_backward`. This should run back-propagation over the entire sequence, making calls to the `rnn_step_backward` function that you defined earlier. You should see errors on the order of e-6 or less.

```

np.random.seed(231)

N, D, T, H = 2, 3, 10, 5

x = np.random.randn(N, T, D)
h0 = np.random.randn(N, H)
wx = np.random.randn(D, H)
wh = np.random.randn(H, H)
b = np.random.randn(H)

out, cache = rnn_forward(x, h0, wx, wh, b)

dout = np.random.randn(*out.shape)

dx, dh0, dwx, dwh, db = rnn_backward(dout, cache)

fx = lambda x: rnn_forward(x, h0, wx, wh, b)[0]
fh0 = lambda h0: rnn_forward(x, h0, wx, wh, b)[0]
fwx = lambda wx: rnn_forward(x, h0, wx, wh, b)[0]
fwh = lambda wh: rnn_forward(x, h0, wx, wh, b)[0]
fb = lambda b: rnn_forward(x, h0, wx, wh, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
dh0_num = eval_numerical_gradient_array(fh0, h0, dout)
dwx_num = eval_numerical_gradient_array(fwx, wx, dout)
dwh_num = eval_numerical_gradient_array(fwh, wh, dout)
db_num = eval_numerical_gradient_array(fb, b, dout)

print('dx error: ', rel_error(dx_num, dx))
print('dh0 error: ', rel_error(dh0_num, dh0))
print('dwx error: ', rel_error(dwx_num, dwx))
print('dwh error: ', rel_error(dwh_num, dwh))
print('db error: ', rel_error(db_num, db))

```

```
dx error: 1.5354482248401769e-09
dh0 error: 3.3830821485562176e-09
dwx error: 7.23583883274483e-09
dwh error: 1.3049601378601992e-07
db error: 1.5197668388626435e-10
```

Word embedding: forward

In deep learning systems, we commonly represent words using vectors. Each word of the vocabulary will be associated with a vector, and these vectors will be learned jointly with the rest of the system.

In the file `cs231n/rnn_layers.py`, implement the function `word_embedding_forward` to convert words (represented by integers) into vectors. Run the following to check your implementation. You should see an error on the order of `e-8` or less.

```
N, T, V, D = 2, 4, 5, 3

x = np.asarray([[0, 3, 1, 2], [2, 1, 0, 3]])
w = np.linspace(0, 1, num=V*D).reshape(V, D)

out, _ = word_embedding_forward(x, w)
expected_out = np.asarray([
    [[ 0.,          0.07142857,  0.14285714],
     [ 0.64285714,  0.71428571,  0.78571429],
     [ 0.21428571,  0.28571429,  0.35714286],
     [ 0.42857143,  0.5,          0.57142857]],
    [[ 0.42857143,  0.5,          0.57142857],
     [ 0.21428571,  0.28571429,  0.35714286],
     [ 0.,          0.07142857,  0.14285714],
     [ 0.64285714,  0.71428571,  0.78571429]]])

print('out error: ', rel_error(expected_out, out))
```

```
out error: 1.000000094736443e-08
```

Word embedding: backward

Implement the backward pass for the word embedding function in the function `word_embedding_backward`. After doing so run the following to numerically gradient check your implementation. You should see an error on the order of `e-11` or less.

```
np.random.seed(231)

N, T, V, D = 50, 3, 5, 6
x = np.random.randint(V, size=(N, T))
w = np.random.randn(V, D)

out, cache = word_embedding_forward(x, w)
dout = np.random.randn(*out.shape)
dw = word_embedding_backward(dout, cache)

f = lambda w: word_embedding_forward(x, w)[0]
```

```
dw_num = eval_numerical_gradient_array(f, w, dout)

print('dw error: ', rel_error(dw, dw_num))
```

```
dw error: 3.2774595693100364e-12
```

Temporal Affine layer

At every timestep we use an affine function to transform the RNN hidden vector at that timestep into scores for each word in the vocabulary. Because this is very similar to the affine layer that you implemented in assignment 2, we have provided this function for you in the `temporal_affine_forward` and `temporal_affine_backward` functions in the file `cs231n/rnn_layers.py`. Run the following to perform numeric gradient checking on the implementation. You should see errors on the order of e-9 or less.

```
np.random.seed(231)

# Gradient check for temporal affine layer
N, T, D, M = 2, 3, 4, 5
x = np.random.randn(N, T, D)
w = np.random.randn(D, M)
b = np.random.randn(M)

out, cache = temporal_affine_forward(x, w, b)

dout = np.random.randn(*out.shape)

fx = lambda x: temporal_affine_forward(x, w, b)[0]
fw = lambda w: temporal_affine_forward(x, w, b)[0]
fb = lambda b: temporal_affine_forward(x, w, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
dw_num = eval_numerical_gradient_array(fw, w, dout)
db_num = eval_numerical_gradient_array(fb, b, dout)

dx, dw, db = temporal_affine_backward(dout, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dw error: ', rel_error(dw_num, dw))
print('db error: ', rel_error(db_num, db))
```

```
dx error: 2.9215945034030545e-10
dw error: 1.5772088618663602e-10
db error: 3.252200556967514e-11
```

Temporal Softmax loss

In an RNN language model, at every timestep we produce a score for each word in the vocabulary. We know the ground-truth word at each timestep, so we use a softmax loss function to compute loss and gradient at each timestep. We sum the losses over time and average them over the minibatch.

However there is one wrinkle: since we operate over minibatches and different captions may have different lengths, we append `<NULL>` tokens to the end of each caption so they all have the same length. We don't want these `<NULL>` tokens to count toward the loss or gradient, so in addition to scores and ground-truth labels our loss function also accepts a `mask` array that tells it which elements of the scores count towards the loss.

Since this is very similar to the softmax loss function you implemented in assignment 1, we have implemented this loss function for you; look at the `temporal_softmax_loss` function in the file `cs231n/rnn_layers.py`.

Run the following cell to sanity check the loss and perform numeric gradient checking on the function. You should see an error for `dx` on the order of e-7 or less.

```
# Sanity check for temporal softmax loss
from cs231n.rnn_layers import temporal_softmax_loss

N, T, V = 100, 1, 10

def check_loss(N, T, V, p):
    x = 0.001 * np.random.randn(N, T, V)
    y = np.random.randint(V, size=(N, T))
    mask = np.random.rand(N, T) <= p
    print(temporal_softmax_loss(x, y, mask)[0])

check_loss(100, 1, 10, 1.0)    # Should be about 2.3
check_loss(100, 10, 10, 1.0)   # Should be about 23
check_loss(5000, 10, 10, 0.1)  # Should be within 2.2-2.4

# Gradient check for temporal softmax loss
N, T, V = 7, 8, 9

x = np.random.randn(N, T, V)
y = np.random.randint(V, size=(N, T))
mask = (np.random.rand(N, T) > 0.5)

loss, dx = temporal_softmax_loss(x, y, mask, verbose=False)

dx_num = eval_numerical_gradient(lambda x: temporal_softmax_loss(x, y, mask)[0],
                                  x, verbose=False)

print('dx error: ', rel_error(dx, dx_num))
```

```
2.302778739129877
23.025691984820515
2.267110925052805
dx error: 2.6313495241754103e-08
```

RNN for image captioning

Now that you have implemented the necessary layers, you can combine them to build an image captioning model. Open the file `cs231n/classifiers/rnn.py` and look at the `CaptioningRNN` class.

Implement the forward and backward pass of the model in the `loss` function. For now you only need to implement the case where `cell_type='rnn'` for vanialla RNNs; you will implement the LSTM case later. After doing so, run the following to check your forward pass using a small test case; you should see error on the order of `e-10` or less.

```
N, D, W, H = 10, 20, 30, 40
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
V = len(word_to_idx)
T = 13

model = CaptioningRNN(word_to_idx,
                      input_dim=D,
                      wordvec_dim=W,
                      hidden_dim=H,
                      cell_type='rnn',
                      dtype=np.float64)

# Set all model parameters to fixed values
for k, v in model.params.items():
    model.params[k] = np.linspace(-1.4, 1.3, num=v.size).reshape(*v.shape)

features = np.linspace(-1.5, 0.3, num=(N * D)).reshape(N, D)
captions = (np.arange(N * T) % V).reshape(N, T)

loss, grads = model.loss(features, captions)
expected_loss = 9.83235591003

print('loss: ', loss)
print('expected loss: ', expected_loss)
print('difference: ', abs(loss - expected_loss))
```

```
loss: 9.832355910027387
expected loss: 9.83235591003
difference: 2.6130209107577684e-12
```

Run the following cell to perform numeric gradient checking on the `CaptioningRNN` class; you should see errors around the order of `e-6` or less.

```
np.random.seed(231)

batch_size = 2
timesteps = 3
input_dim = 4
wordvec_dim = 5
hidden_dim = 6
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
vocab_size = len(word_to_idx)

captions = np.random.randint(vocab_size, size=(batch_size, timesteps))
features = np.random.randn(batch_size, input_dim)

model = CaptioningRNN(word_to_idx,
                      input_dim=input_dim,
                      wordvec_dim=wordvec_dim,
                      hidden_dim=hidden_dim,
```

```

        cell_type='rnn',
        dtype=np.float64,
    )

loss, grads = model.loss(features, captions)

for param_name in sorted(grads):
    f = lambda _: model.loss(features, captions)[0]
    param_grad_num = eval_numerical_gradient(f, model.params[param_name],
    verbose=False, h=1e-6)
    e = rel_error(param_grad_num, grads[param_name])
    print('%s relative error: %e' % (param_name, e))

```

```

w_embed relative error: 2.331071e-09
w_proj relative error: 9.974425e-09
w_vocab relative error: 4.274378e-09
wh relative error: 1.313259e-08
wx relative error: 8.455229e-07
b relative error: 9.727212e-10
b_proj relative error: 1.934807e-08
b_vocab relative error: 7.087090e-11

```

Overfit small data

Similar to the `Solver` class that we used to train image classification models on the previous assignment, on this assignment we use a `CaptioningSolver` class to train image captioning models. Open the file `cs231n/captioning_solver.py` and read through the `CaptioningSolver` class; it should look very familiar.

Once you have familiarized yourself with the API, run the following to make sure your model overfits a small sample of 100 training examples. You should see a final loss of less than 0.1.

```

np.random.seed(231)

small_data = load_coco_data(max_train=50)

small_rnn_model = CaptioningRNN(
    cell_type='rnn',
    word_to_idx=data['word_to_idx'],
    input_dim=data['train_features'].shape[1],
    hidden_dim=512,
    wordvec_dim=256,
)

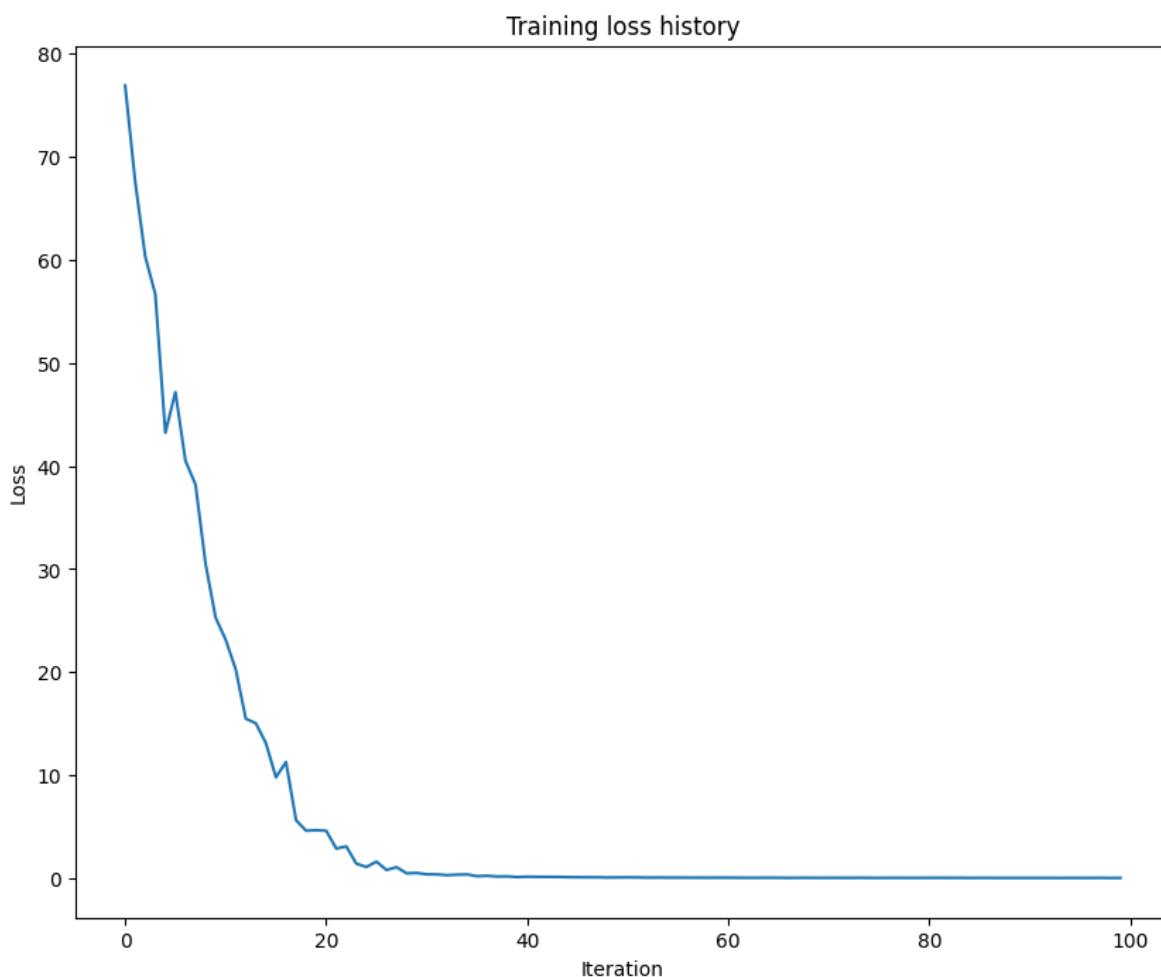
small_rnn_solver = CaptioningSolver(small_rnn_model, small_data,
    update_rule='adam',
    num_epochs=50,
    batch_size=25,
    optim_config={
        'learning_rate': 1e-2,
    },
    lr_decay=0.95,
    verbose=True, print_every=10,
)

```

```
small_rnn_solver.train()

# Plot the training losses
plt.plot(small_rnn_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```

```
base dir
/content/drive/MyDrive/cs231n/assignments/assignment3/cs231n/datasets/coco_capti
oning
(Iteration 1 / 100) loss: 76.913487
(Iteration 11 / 100) loss: 23.150531
(Iteration 21 / 100) loss: 4.616762
(Iteration 31 / 100) loss: 0.388363
(Iteration 41 / 100) loss: 0.158645
(Iteration 51 / 100) loss: 0.094666
(Iteration 61 / 100) loss: 0.063804
(Iteration 71 / 100) loss: 0.049477
(Iteration 81 / 100) loss: 0.047606
(Iteration 91 / 100) loss: 0.038166
```



Print final training loss. You should see a final loss of less than 0.1.

```
print('Final loss: ', small_rnn_solver.loss_history[-1])
```

```
Final loss: 0.04022374051571897
```

Test-time sampling

Unlike classification models, image captioning models behave very differently at training time and at test time. At training time, we have access to the ground-truth caption, so we feed ground-truth words as input to the RNN at each timestep. At test time, we sample from the distribution over the vocabulary at each timestep, and feed the sample as input to the RNN at the next timestep.

In the file `cs231n/classifiers/rnn.py`, implement the `sample` method for test-time sampling. After doing so, run the following to sample from your overfitted model on both training and validation data. The samples on training data should be very good; the samples on validation data probably won't make sense.

```
# If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
for split in ['train', 'val']:
    minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2)
    gt_captions, features, urls = minibatch
    gt_captions = decode_captions(gt_captions, data['idx_to_word'])

    sample_captions = small_rnn_model.sample(features)
    sample_captions = decode_captions(sample_captions, data['idx_to_word'])

    for gt_caption, sample_caption, url in zip(gt_captions, sample_captions,
                                                urls):
        img = image_from_url(url)
        # skip missing URLs.
        if img is None: continue
        plt.imshow(img)
        plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
        plt.axis('off')
        plt.show()
```

Output hidden; open in <https://colab.research.google.com> to view.

<START> a man who is wearing a <UNK> shirt and a tie <END>



<START> a <UNK> donuts in a <UNK> board box <END>





<START> this shows a street in a busy <UNK> area in <UNK> <END>



<START> two one way street signs <UNK> on a street sign <END>



INLINE QUESTION 1

In our current image captioning setup, our RNN language model produces a word at every timestep as its output. However, an alternate way to pose the problem is to train the network to operate over *characters* (e.g. 'a', 'b', etc.) as opposed to words, so that at every timestep, it receives the previous character as input and tries to predict the next character in the sequence. For example, the network might generate a caption like

'A', ' ', 'c', 'a', 't', ' ', 'o', 'n', ' ', 'a', ' ', 'b', 'e', 'd'

Can you describe one advantage of an image-captioning model that uses a character-level RNN? Can you also describe one disadvantage? HINT: there are several valid answers, but it might be useful to compare the parameter space of word-level and character-level models.

Your Answer:

Advantages

- **Smaller Vocabulary:** Character-level RNNs operate on a much smaller vocabulary (e.g., letters, punctuation) compared to word-level RNNs, which need to handle thousands or even tens of thousands of words. This reduces the model's parameter space and computational complexity during training and inference.
- **Flexibility in Generation:** Character-level RNNs can generate any possible word combination without being restricted by a predefined vocabulary. Word-level RNNs, on the other hand, are limited to the vocabulary defined during training and struggle with out-of-vocabulary (OOV) words.
- **Handling Spelling Variations:** Character-level RNNs can more easily handle spelling variations or errors, such as generating both "color" and "colour," since they operate at the character level rather than the word level.

Disadvantages

- **Longer Sequence Length:** Character-level RNNs process longer sequences because generating a single word requires multiple characters. For example, generating "cat" requires 3 characters, whereas a word-level model only needs 1 word. This can lead to increased computational overhead.
- **Weaker Semantic Understanding:** Character-level RNNs must learn semantic information from character sequences, which is more challenging compared to word-level RNNs that directly leverage word-level semantics.
- **Lower Generation Quality:** Character-level RNNs may generate more spelling errors or nonsensical character combinations due to the lack of word-level semantic constraints.

Image Captioning with LSTMs

In the previous exercise you implemented a vanilla RNN and applied it to image captioning. In this notebook you will implement the LSTM update rule and use it for image captioning.

```
# set workspace and get datasets
from google.colab import drive
drive.mount('/content/drive')
FOLDERNAME = "cs231n/assignments/assignment3/"
assert FOLDERNAME is not None, "[!] Enter the foldername."
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/cs231n/assignments/assignment3/cs231n/datasets
/content/drive/My Drive/cs231n/assignments/assignment3
```

```
# As usual, a bit of setup
import time, os, json
from email.mime import image

import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient,
eval_numerical_gradient_array
from cs231n.rnn_layers import *
from cs231n.captioning_solver import CaptioningSolver
from cs231n.classifiers.rnn import CaptioningRNN
from cs231n.coco_utils import load_coco_data, sample_coco_minibatch,
decode_captions
from cs231n.image_utils import image_from_url

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-
ipython
%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

Load MS-COCO data

As in the previous notebook, we will use the Microsoft COCO dataset for captioning.

```
# Load coco data from disk; this returns a dictionary
# we'll work with dimensionality-reduced features for this notebook, but feel
# free to experiment with the original features by changing the flag below.
data = load_coco_data(pca_features=True)

# Print out all the keys and values from the data dictionary
for k, v in data.items():
    if type(v) == np.ndarray:
        print(k, type(v), v.shape, v.dtype)
    else:
        print(k, type(v), len(v))
```

```
base dir /content/drive/My
Drive/cs231n/assignments/assignment3/cs231n/datasets/coco_captioning
trainCaptions <class 'numpy.ndarray'> (400135, 17) int32
trainImageIdxs <class 'numpy.ndarray'> (400135,) int32
valCaptions <class 'numpy.ndarray'> (195954, 17) int32
valImageIdxs <class 'numpy.ndarray'> (195954,) int32
trainFeatures <class 'numpy.ndarray'> (82783, 512) float32
valFeatures <class 'numpy.ndarray'> (40504, 512) float32
idxToWord <class 'list'> 1004
wordToIdx <class 'dict'> 1004
trainUrls <class 'numpy.ndarray'> (82783,) <U63
valUrls <class 'numpy.ndarray'> (40504,) <U63
```

LSTM

If you read recent papers, you'll see that many people use a variant on the vanilla RNN called Long-Short Term Memory (LSTM) RNNs. Vanilla RNNs can be tough to train on long sequences due to vanishing and exploding gradients caused by repeated matrix multiplication. LSTMs solve this problem by replacing the simple update rule of the vanilla RNN with a gating mechanism as follows.

Similar to the vanilla RNN, at each timestep we receive an input $x_t \in \mathbb{R}^D$ and the previous hidden state $h_{t-1} \in \mathbb{R}^H$; the LSTM also maintains an H -dimensional *cell state*, so we also receive the previous cell state $c_{t-1} \in \mathbb{R}^H$. The learnable parameters of the LSTM are an *input-to-hidden* matrix $W_x \in \mathbb{R}^{4H \times D}$, a *hidden-to-hidden* matrix $W_h \in \mathbb{R}^{4H \times H}$ and a *bias vector* $b \in \mathbb{R}^{4H}$.

At each timestep we first compute an *activation vector* $a \in \mathbb{R}^{4H}$ as $a = W_x x_t + W_h h_{t-1} + b$. We then divide this into four vectors $a_i, a_f, a_o, a_g \in \mathbb{R}^H$ where a_i consists of the first H elements of a , a_f is the next H elements of a , etc. We then compute the *input gate* $g \in \mathbb{R}^H$, *forget gate* $f \in \mathbb{R}^H$, *output gate* $o \in \mathbb{R}^H$ and *block input* $g \in \mathbb{R}^H$ as

$$i = \sigma(a_i) \quad f = \sigma(a_f) \quad o = \sigma(a_o) \quad g = \tanh(a_g)$$

where σ is the sigmoid function and \tanh is the hyperbolic tangent, both applied elementwise.

Finally we compute the next cell state c_t and next hidden state h_t as

$$c_t = f \odot c_{t-1} + i \odot g \quad h_t = o \odot \tanh(c_t)$$

where \odot is the elementwise product of vectors.

In the rest of the notebook we will implement the LSTM update rule and apply it to the image captioning task.

In the code, we assume that data is stored in batches so that $X_t \in \mathbb{R}^{N \times D}$, and will work with transposed versions of the parameters: $W_x \in \mathbb{R}^{D \times 4H}$, $W_h \in \mathbb{R}^{H \times 4H}$ so that activations $A \in \mathbb{R}^{N \times 4H}$ can be computed efficiently as $A = X_t W_x + H_{t-1} W_h$

LSTM: step forward

Implement the forward pass for a single timestep of an LSTM in the `lstm_step_forward` function in the file `cs231n/rnn_layers.py`. This should be similar to the `rnn_step_forward` function that you implemented above, but using the LSTM update rule instead.

Once you are done, run the following to perform a simple test of your implementation. You should see errors on the order of `e-8` or less.

```
N, D, H = 3, 4, 5
x = np.linspace(-0.4, 1.2, num=N*D).reshape(N, D)
prev_h = np.linspace(-0.3, 0.7, num=N*H).reshape(N, H)
prev_c = np.linspace(-0.4, 0.9, num=N*H).reshape(N, H)
wx = np.linspace(-2.1, 1.3, num=4*D*H).reshape(D, 4 * H)
wh = np.linspace(-0.7, 2.2, num=4*H*H).reshape(H, 4 * H)
b = np.linspace(0.3, 0.7, num=4*H)

next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, wx, wh, b)

expected_next_h = np.asarray([
    [ 0.24635157,  0.28610883,  0.32240467,  0.35525807,  0.38474904],
    [ 0.49223563,  0.55611431,  0.61507696,  0.66844003,  0.7159181 ],
    [ 0.56735664,  0.66310127,  0.74419266,  0.80889665,  0.858299  ]])
expected_next_c = np.asarray([
    [ 0.32986176,  0.39145139,  0.451556,    0.51014116,  0.56717407],
    [ 0.66382255,  0.76674007,  0.87195994,  0.97902709,  1.08751345],
    [ 0.74192008,  0.90592151,  1.07717006,  1.25120233,  1.42395676]])

print('next_h error: ', rel_error(expected_next_h, next_h))
print('next_c error: ', rel_error(expected_next_c, next_c))
```

```
next_h error:  5.7054131967097955e-09
next_c error:  5.8143123088804145e-09
```

LSTM: step backward

Implement the backward pass for a single LSTM timestep in the function `lstm_step_backward` in the file `cs231n/rnn_layers.py`. Once you are done, run the following to perform numeric gradient checking on your implementation. You should see errors on the order of `e-7` or less.

```
np.random.seed(231)

N, D, H = 4, 5, 6
x = np.random.randn(N, D)
prev_h = np.random.randn(N, H)
```

```

prev_c = np.random.randn(N, H)
wx = np.random.randn(D, 4 * H)
wh = np.random.randn(H, 4 * H)
b = np.random.randn(4 * H)

next_h, next_c, cache = lstm_step_forward(x, prev_h, prev_c, wx, wh, b)

dnext_h = np.random.randn(*next_h.shape)
dnext_c = np.random.randn(*next_c.shape)

fx_h = lambda x: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]
fh_h = lambda h: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]
fc_h = lambda c: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]
fwx_h = lambda wx: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]
fwh_h = lambda wh: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]
fb_h = lambda b: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[0]

fx_c = lambda x: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]
fh_c = lambda h: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]
fc_c = lambda c: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]
fwx_c = lambda wx: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]
fwh_c = lambda wh: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]
fb_c = lambda b: lstm_step_forward(x, prev_h, prev_c, wx, wh, b)[1]

num_grad = eval_numerical_gradient_array

dx_num = num_grad(fx_h, x, dnext_h) + num_grad(fx_c, x, dnext_c)
dh_num = num_grad(fh_h, prev_h, dnext_h) + num_grad(fh_c, prev_h, dnext_c)
dc_num = num_grad(fc_h, prev_c, dnext_h) + num_grad(fc_c, prev_c, dnext_c)
dwx_num = num_grad(fwx_h, wx, dnext_h) + num_grad(fwx_c, wx, dnext_c)
dwh_num = num_grad(fwh_h, wh, dnext_h) + num_grad(fwh_c, wh, dnext_c)
db_num = num_grad(fb_h, b, dnext_h) + num_grad(fb_c, b, dnext_c)

dx, dh, dc, dwx, dwh, db = lstm_step_backward(dnext_h, dnext_c, cache)

print('dx error: ', rel_error(dx_num, dx))
print('dh error: ', rel_error(dh_num, dh))
print('dc error: ', rel_error(dc_num, dc))
print('dwx error: ', rel_error(dwx_num, dwx))
print('dwh error: ', rel_error(dwh_num, dwh))
print('db error: ', rel_error(db_num, db))

```

```

dx error: 5.833093013472354e-10
dh error: 3.4168728051126624e-10
dc error: 1.5221723979041107e-10
dwx error: 1.6933643922734908e-09
dwh error: 2.7311400266248628e-08
db error: 1.7349356733443412e-10

```

LSTM: forward

In the function `lstm_forward` in the file `cs231n/rnn_layers.py`, implement the `lstm_forward` function to run an LSTM forward on an entire timeseries of data.

When you are done, run the following to check your implementation. You should see an error on the order of `e-7` or less.

```

N, D, H, T = 2, 5, 4, 3
x = np.linspace(-0.4, 0.6, num=N*T*D).reshape(N, T, D)
h0 = np.linspace(-0.4, 0.8, num=N*H).reshape(N, H)
wx = np.linspace(-0.2, 0.9, num=4*D*H).reshape(D, 4 * H)
wh = np.linspace(-0.3, 0.6, num=4*H*H).reshape(H, 4 * H)
b = np.linspace(0.2, 0.7, num=4*H)

h, cache = lstm_forward(x, h0, wx, wh, b)

expected_h = np.asarray([
[[ 0.01764008,  0.01823233,  0.01882671,  0.0194232 ],
 [ 0.11287491,  0.12146228,  0.13018446,  0.13902939],
 [ 0.31358768,  0.33338627,  0.35304453,  0.37250975]],
[[ 0.45767879,  0.4761092,   0.4936887,   0.51041945],
 [ 0.6704845,   0.69350089,  0.71486014,  0.7346449 ],
 [ 0.81733511,  0.83677871,  0.85403753,  0.86935314]]])

print('h error: ', rel_error(expected_h, h))

```

h error: 8.610537442272635e-08

LSTM: backward

Implement the backward pass for an LSTM over an entire timeseries of data in the function `lstm_backward` in the file `cs231n/rnn_layers.py`. When you are done, run the following to perform numeric gradient checking on your implementation. You should see errors on the order of `e-8` or less. (For `dwh`, it's fine if your error is on the order of `e-6` or less).

```

from cs231n.rnn_layers import lstm_forward, lstm_backward
np.random.seed(231)

N, D, T, H = 2, 3, 10, 6

x = np.random.randn(N, T, D)
h0 = np.random.randn(N, H)
wx = np.random.randn(D, 4 * H)
wh = np.random.randn(H, 4 * H)
b = np.random.randn(4 * H)

out, cache = lstm_forward(x, h0, wx, wh, b)

dout = np.random.randn(*out.shape)

dx, dh0, dwx, dwh, db = lstm_backward(dout, cache)

fx = lambda x: lstm_forward(x, h0, wx, wh, b)[0]
fh0 = lambda h0: lstm_forward(x, h0, wx, wh, b)[0]
fwx = lambda wx: lstm_forward(x, h0, wx, wh, b)[0]
fwh = lambda wh: lstm_forward(x, h0, wx, wh, b)[0]
fb = lambda b: lstm_forward(x, h0, wx, wh, b)[0]

dx_num = eval_numerical_gradient_array(fx, x, dout)
dh0_num = eval_numerical_gradient_array(fh0, h0, dout)
dwx_num = eval_numerical_gradient_array(fwx, wx, dout)
dwh_num = eval_numerical_gradient_array(fwh, wh, dout)

```

```

db_num = eval_numerical_gradient_array(fb, b, dout)

print('dx error: ', rel_error(dx_num, dx))
print('dh0 error: ', rel_error(dh0_num, dh0))
print('dwx error: ', rel_error(dwx_num, dwx))
print('dwh error: ', rel_error(dwh_num, dwh))
print('db error: ', rel_error(db_num, db))

```

```

dx error: 7.1480958161034e-10
dh0 error: 2.3791401549917165e-08
dwx error: 1.0057877119886127e-09
dwh error: 6.064216037522062e-07
db error: 6.252306209370124e-10

```

INLINE QUESTION

Recall that in an LSTM the input gate i , forget gate f , and output gate o are all outputs of a sigmoid function. Why don't we use the ReLU activation function instead of sigmoid to compute these values? Explain.

Your Answer:

In an LSTM, the input gate i , forget gate f , and output gate o are designed to control the flow of information into and out of the cell state. These gates need to produce values between 0 and 1 to act as "switches" that determine how much information should be retained, forgotten, or passed along. Here's why we use the sigmoid function instead of ReLU for these gates:

1. Output Range

Sigmoid: The sigmoid function outputs values in the range $[0,1]$, which is ideal for representing probabilities or proportions. This allows the gates to act as "filters" that control the proportion of information passed through.

ReLU: The ReLU function outputs values in the range $[0, \infty)$. This unbounded output makes it unsuitable for gates, as it cannot represent a proportion or probability.

2. Interpretation as Probabilities

The gates in an LSTM are often interpreted as probabilities (e.g., the probability that a piece of information should be retained or forgotten). The sigmoid function naturally maps inputs to probabilities, while ReLU does not.

3. Gradient Stability

Sigmoid: The derivative of the sigmoid function is well-behaved and avoids the problem of exploding gradients, which is critical for training deep recurrent networks.

ReLU: While ReLU is computationally efficient and mitigates vanishing gradients in feedforward networks, its unbounded output can lead to exploding gradients in the context of LSTM gates, where values are repeatedly multiplied over time steps.

4. Control Over Information Flow

The sigmoid function ensures that the gates produce values that can be directly multiplied with other values (e.g., cell state or hidden state) to control information flow. ReLU lacks this property because it does not constrain outputs to a fixed range.

LSTM captioning model

Now that you have implemented an LSTM, update the implementation of the `loss` method of the `CaptioningRNN` class in the file `cs231n/classifiers/rnn.py` to handle the case where `self.cell_type` is `lstm`. This should require adding less than 10 lines of code.

Once you have done so, run the following to check your implementation. You should see a difference on the order of `e-10` or less.

```
N, D, W, H = 10, 20, 30, 40
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
V = len(word_to_idx)
T = 13

model = CaptioningRNN(word_to_idx,
                      input_dim=D,
                      wordvec_dim=W,
                      hidden_dim=H,
                      cell_type='lstm',
                      dtype=np.float64)

# Set all model parameters to fixed values
for k, v in model.params.items():
    model.params[k] = np.linspace(-1.4, 1.3, num=v.size).reshape(*v.shape)

features = np.linspace(-0.5, 1.7, num=N*D).reshape(N, D)
captions = (np.arange(N * T) % V).reshape(N, T)

loss, grads = model.loss(features, captions)
expected_loss = 9.82445935443

print('loss: ', loss)
print('expected loss: ', expected_loss)
print('difference: ', abs(loss - expected_loss))
```

```
loss: 9.82445935443226
expected loss: 9.82445935443
difference: 2.261302256556519e-12
```

Overfit LSTM captioning model

Run the following to overfit an LSTM captioning model on the same small dataset as we used for the RNN previously. You should see a final loss less than 0.5.

```
np.random.seed(231)

small_data = load_coco_data(max_train=50)

small_lstm_model = CaptioningRNN(
    cell_type='lstm',
    word_to_idx=data['word_to_idx'],
    input_dim=data['train_features'].shape[1],
    hidden_dim=512,
    wordvec_dim=256,
    dtype=np.float32,
)

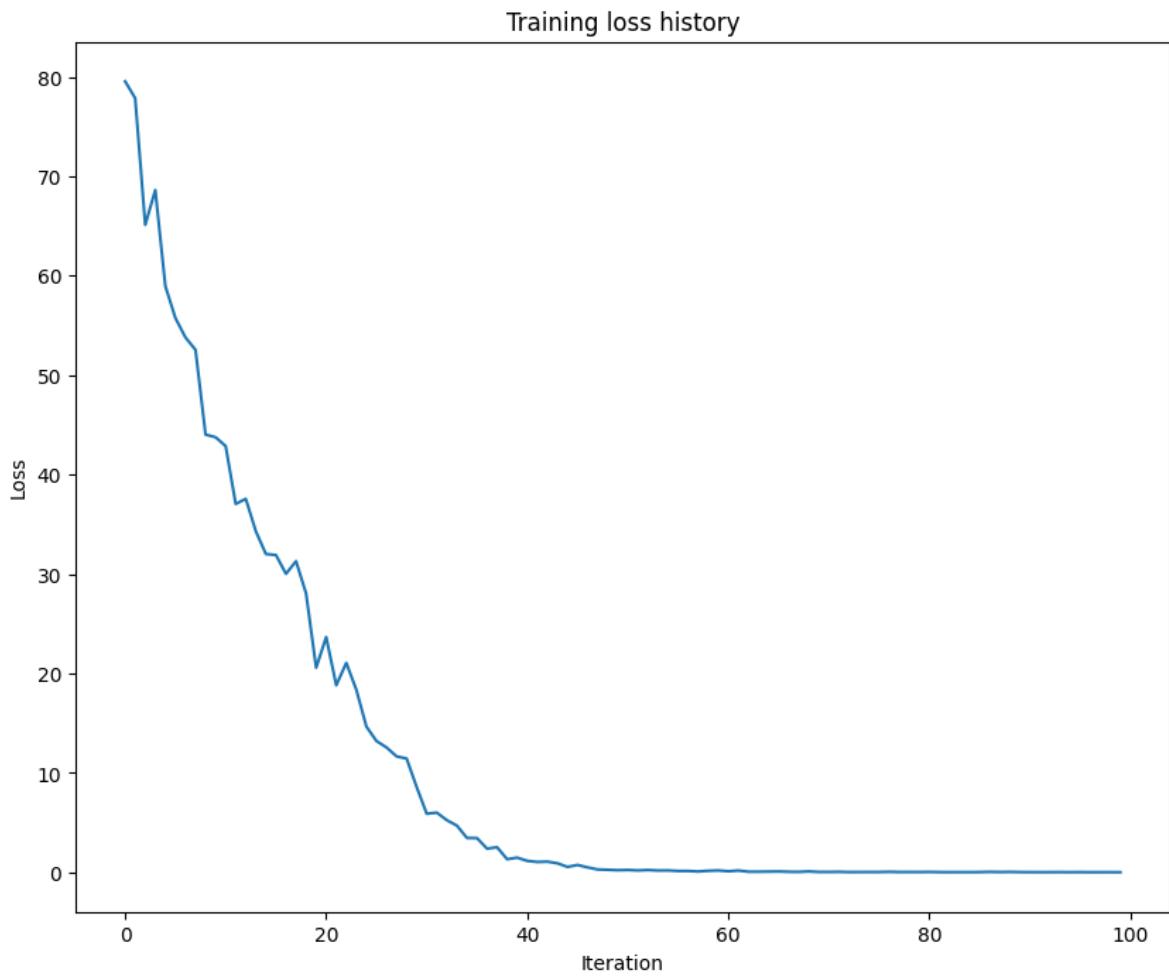
small_lstm_solver = CaptioningSolver(small_lstm_model, small_data,
                                    update_rule='adam',
                                    num_epochs=50,
                                    batch_size=25,
```

```
    optim_config={
        'learning_rate': 1e-2,
    },
    lr_decay=0.999,
    verbose=True, print_every=10,
)

small_lstm_solver.train()

# Plot the training losses
plt.plot(small_lstm_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```

```
base dir /content/drive/My
Drive/cs231n/assignments/assignment3/cs231n/datasets/coco_captioning
(Iteration 1 / 100) loss: 79.551150
(Iteration 11 / 100) loss: 42.869805
(Iteration 21 / 100) loss: 23.669663
(Iteration 31 / 100) loss: 5.899650
(Iteration 41 / 100) loss: 1.156738
(Iteration 51 / 100) loss: 0.246394
(Iteration 61 / 100) loss: 0.112114
(Iteration 71 / 100) loss: 0.045706
(Iteration 81 / 100) loss: 0.046131
(Iteration 91 / 100) loss: 0.026539
```



Print final training loss. You should see a final loss of less than 0.5.

```
print('Final loss: ', small_lstm_solver.loss_history[-1])
```

```
Final loss: 0.018814703881202818
```

LSTM test-time sampling

Modify the `sample` method of the `CaptioningRNN` class to handle the case where `self.cell_type` is `lstm`. This should take fewer than 10 lines of code.

When you are done run the following to sample from your overfit LSTM model on some training and validation set samples. As with the RNN, training results should be very good, and validation results probably won't make a lot of sense (because we're overfitting).

```
# If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
for split in ['train', 'val']:
    minibatch = sample_coco_minibatch(small_data, split=split, batch_size=2)
    gt_captions, features, urls = minibatch
    gt_captions = decode_captions(gt_captions, data['idx_to_word'])

    sample_captions = small_lstm_model.sample(features)
    sample_captions = decode_captions(sample_captions, data['idx_to_word'])
```

```
for gt_caption, sample_caption, url in zip(gt_captions, sample_captions,
urls):
    img = image_from_url(url)
    if img is None: continue
    plt.imshow(img)
    plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
    plt.axis('off')
    plt.show()
```

Output hidden; open in <https://colab.research.google.com> to view.

train

a man standing on the side of a road with bags of luggage <END>
GT:<START> a man standing on the side of a road with bags of luggage <END>



train

a man <UNK> with a bright colorful kite <END>
GT:<START> a man <UNK> with a bright colorful kite <END>





val

a plane that is about to take off flying <END>
GT:<START> a sign that is on the front of a train station <END>



val

a dog is <UNK> on a <UNK> <END>
GT:<START> a car is parked on a street at night <END>

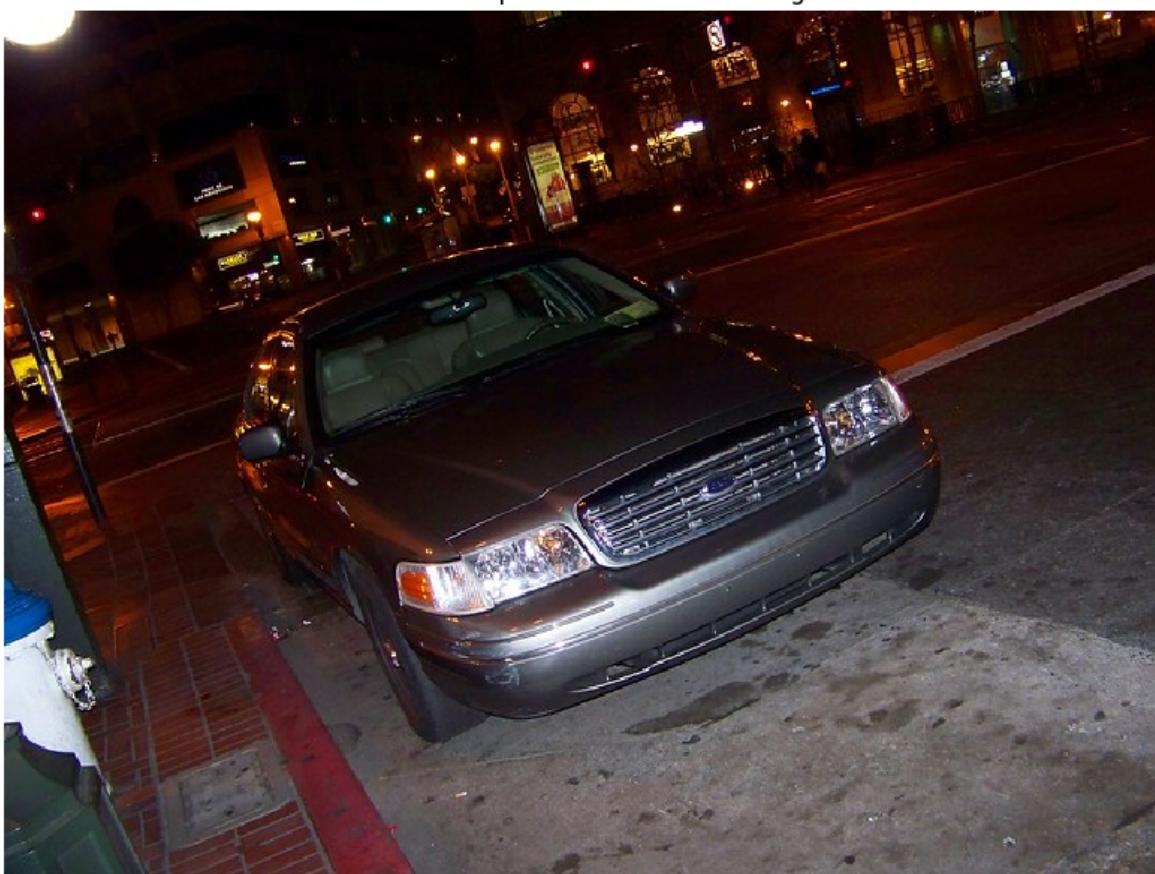


Image Captioning with Transformers

You have now implemented a vanilla RNN and for the task of image captioning. In this notebook you will implement key pieces of a transformer decoder to accomplish the same task.

NOTE: This notebook will be primarily written in PyTorch rather than NumPy, unlike the RNN notebook.

```
from google.colab import drive
drive.mount('/content/drive')
FOLDERNAME = "cs231n/assignments/assignment3/"
assert FOLDERNAME is not None, "[!] Enter the foldername."
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/cs231n/assignments/assignment3/cs231n/datasets
/content/drive/My Drive/cs231n/assignments/assignment3
```

```
# Setup cell.
import time, os, json
import numpy as np
import matplotlib.pyplot as plt

from cs231n.gradient_check import eval_numerical_gradient,
eval_numerical_gradient_array
from cs231n.transformer_layers import *
from cs231n.captioning_solver_transformer import CaptioningSolverTransformer
from cs231n.classifiers.transformer import CaptioningTransformer
from cs231n.coco_utils import load_coco_data, sample_coco_minibatch,
decode_captions
from cs231n.image_utils import image_from_url

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # Set default size of plots.
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

%load_ext autoreload
%autoreload 2

def rel_error(x, y):
    """ returns relative error """
    return np.max(np.abs(x - y) / (np.maximum(1e-8, np.abs(x) + np.abs(y))))
```

COCO Dataset

As in the previous notebooks, we will use the COCO dataset for captioning.

```
# Load COCO data from disk into a dictionary.  
data = load_coco_data(pca_features=True)  
  
# Print out all the keys and values from the data dictionary.  
for k, v in data.items():  
    if type(v) == np.ndarray:  
        print(k, type(v), v.shape, v.dtype)  
    else:  
        print(k, type(v), len(v))
```

```
base dir /content/drive/My  
Drive/cs231n/assignments/assignment3/cs231n/datasets/coco_captioning  
trainCaptions <class 'numpy.ndarray'> (400135, 17) int32  
trainImageIdxs <class 'numpy.ndarray'> (400135,) int32  
valCaptions <class 'numpy.ndarray'> (195954, 17) int32  
valImageIdxs <class 'numpy.ndarray'> (195954,) int32  
trainFeatures <class 'numpy.ndarray'> (82783, 512) float32  
valFeatures <class 'numpy.ndarray'> (40504, 512) float32  
idxToWord <class 'list'> 1004  
wordToIdx <class 'dict'> 1004  
trainUrls <class 'numpy.ndarray'> (82783,) <U63  
valUrls <class 'numpy.ndarray'> (40504,) <U63
```

Transformer

As you have seen, RNNs are incredibly powerful but often slow to train. Further, RNNs struggle to encode long-range dependencies (though LSTMs are one way of mitigating the issue). In 2017, Vaswani et al introduced the Transformer in their paper "[Attention Is All You Need](#)" to a) introduce parallelism and b) allow models to learn long-range dependencies. The paper not only led to famous models like BERT and GPT in the natural language processing community, but also an explosion of interest across fields, including vision. While here we introduce the model in the context of image captioning, the idea of attention itself is much more general.

Transformer: Multi-Headed Attention

Dot-Product Attention

Recall that attention can be viewed as an operation on a query $q \in \mathbb{R}^d$, a set of value vectors $\{v_1, \dots, v_n\}, v_i \in \mathbb{R}^d$, and a set of key vectors $\{k_1, \dots, k_n\}, k_i \in \mathbb{R}^d$, specified as

$$c = \sum_{i=1}^n v_i \alpha_i \alpha_i = \frac{\exp(k_i^\top q)}{\sum_{j=1}^n \exp(k_j^\top q)}$$

where α_i are frequently called the "attention weights", and the output $c \in \mathbb{R}^d$ is a correspondingly weighted average over the value vectors.

Self-Attention

In Transformers, we perform self-attention, which means that the values, keys and query are derived from the input $X \in \mathbb{R}^{\ell \times d}$, where ℓ is our sequence length. Specifically, we learn parameter matrices $V, K, Q \in \mathbb{R}^{d \times d}$ to map our input X as follows:

$$\begin{aligned} v_i &= Vx_i \quad i \in \{1, \dots, \ell\} \\ k_i &= Kx_i \quad i \in \{1, \dots, \ell\} \\ q_i &= Qx_i \quad i \in \{1, \dots, \ell\} \end{aligned}$$

Multi-Headed Scaled Dot-Product Attention

In the case of multi-headed attention, we learn a parameter matrix for each head, which gives the model more expressivity to attend to different parts of the input. Let h be number of heads, and Y_i be the attention output of head i . Thus we learn individual matrices Q_i, K_i and V_i . To keep our overall computation the same as the single-headed case, we choose $Q_i \in \mathbb{R}^{d \times d/h}$, $K_i \in \mathbb{R}^{d \times d/h}$ and $V_i \in \mathbb{R}^{d \times d/h}$. Adding in a scaling term $\frac{1}{\sqrt{d/h}}$ to our simple dot-product attention above, we have

$$\begin{aligned} \text{\begin{equation} \label{qkv_eqn} } \\ Y_i = \text{softmax} \left(\frac{XQ_i(XK_i)^T}{\sqrt{d/h}} \right) XVi \\ \text{\end{equation}} \end{aligned}$$

where $Y_i \in \mathbb{R}^{\ell \times d/h}$, where ℓ is our sequence length.

In our implementation, we apply dropout to the attention weights (though in practice it could be used at any step):

$$\begin{aligned} \text{\begin{equation} \label{qkvdropout_eqn} } \\ Y_i = \text{dropout} \left(\text{softmax} \left(\frac{XQ_i(XK_i)^T}{\sqrt{d/h}} \right) \right) XVi \\ \text{\end{equation}} \end{aligned}$$

Finally, then the output of the self-attention is a linear transformation of the concatenation of the heads:

$$Y = [Y_1; \dots; Y_h]A$$

where $A \in \mathbb{R}^{d \times d}$ and $[Y_1; \dots; Y_h] \in \mathbb{R}^{\ell \times d}$.

Implement multi-headed scaled dot-product attention in the `MultiHeadAttention` class in the file `cs231n/transformer_layers.py`. The code below will check your implementation. The relative error should be less than `e-3`.

```
torch.manual_seed(231)

# Choose dimensions such that they are all unique for easier debugging:
# Specifically, the following values correspond to N=1, H=2, T=3, E//H=4, and
# E=8.
batch_size = 1
sequence_length = 3
embed_dim = 8
attn = MultiHeadAttention(embed_dim, num_heads=2)

# Self-attention.
data = torch.randn(batch_size, sequence_length, embed_dim)
self_attn_output = attn(query=data, key=data, value=data)
```

```

# Masked self-attention.
mask = torch.randn(sequence_length, sequence_length) < 0.5
masked_self_attn_output = attn(query=data, key=data, value=data, attn_mask=mask)

# Attention using two inputs.
other_data = torch.randn(batch_size, sequence_length, embed_dim)
attn_output = attn(query=data, key=other_data, value=other_data)

expected_self_attn_output = np.asarray([[[-0.2494,  0.1396,  0.4323, -0.2411, -0.1547,  0.2329, -0.1936,
   -0.1444],
  [-0.1997,  0.1746,  0.7377, -0.3549, -0.2657,  0.2693, -0.2541,
   -0.2476],
  [-0.0625,  0.1503,  0.7572, -0.3974, -0.1681,  0.2168, -0.2478,
   -0.3038]]])

expected_masked_self_attn_output = np.asarray([[[-0.1347,  0.1934,  0.8628, -0.4903, -0.2614,  0.2798, -0.2586,
   -0.3019],
  [-0.1013,  0.3111,  0.5783, -0.3248, -0.3842,  0.1482, -0.3628,
   -0.1496],
  [-0.2071,  0.1669,  0.7097, -0.3152, -0.3136,  0.2520, -0.2774,
   -0.2208]]])

expected_attn_output = np.asarray([[[-0.1980,  0.4083,  0.1968, -0.3477,  0.0321,  0.4258, -0.8972,
   -0.2744],
  [-0.1603,  0.4155,  0.2295, -0.3485, -0.0341,  0.3929, -0.8248,
   -0.2767],
  [-0.0908,  0.4113,  0.3017, -0.3539, -0.1020,  0.3784, -0.7189,
   -0.2912]]])

print('self_attn_output error: ', rel_error(expected_self_attn_output,
self_attn_output.detach().numpy()))
print('masked_self_attn_output error: ',
rel_error(expected_masked_self_attn_output,
masked_self_attn_output.detach().numpy()))
print('attn_output error: ', rel_error(expected_attn_output,
attn_output.detach().numpy()))

```

```

self_attn_output error: 0.0003775124598178026
masked_self_attn_output error: 0.0001526367643724865
attn_output error: 0.0003527921483788199

```

Positional Encoding

While transformers are able to easily attend to any part of their input, the attention mechanism has no concept of token order. However, for many tasks (especially natural language processing), relative token order is very important. To recover this, the authors add a positional encoding to the embeddings of individual word tokens.

Let us define a matrix $P \in \mathbb{R}^{l \times d}$, where $P_{ij} =$

$$\begin{cases} \sin\left(i \cdot 10000^{-\frac{j}{d}}\right) & \text{if } j \text{ is even} \\ \cos\left(i \cdot 10000^{-\frac{(j-1)}{d}}\right) & \text{otherwise} \end{cases}$$

Rather than directly passing an input $X \in \mathbb{R}^{l \times d}$ to our network, we instead pass $X + P$.

Implement this layer in `PositionalEncoding` in `cs231n/transformer_layers.py`. Once you are done, run the following to perform a simple test of your implementation. You should see errors on the order of `e-3` or less.

```
torch.manual_seed(231)

batch_size = 1
sequence_length = 2
embed_dim = 6
data = torch.randn(batch_size, sequence_length, embed_dim)

pos_encoder = PositionalEncoding(embed_dim)
output = pos_encoder(data)

expected_pe_output = np.asarray([[[[-1.2340,  1.1127,  1.6978, -0.0865, -0.0000,
1.2728],
[ 0.9028, -0.4781,  0.5535,  0.8133,  1.2644,
1.7034]]])

print('pe_output error: ', rel_error(expected_pe_output,
output.detach().numpy()))
```

```
pe_output error:  0.00010421011374914356
```

Inline Question 1

Several key design decisions were made in designing the scaled dot product attention we introduced above. Explain why the following choices were beneficial:

1. Using multiple attention heads as opposed to one.
2. Dividing by $\sqrt{d/h}$ before applying the softmax function. Recall that d is the feature dimension and h is the number of heads.
3. Adding a linear transformation to the output of the attention operation.

Only one or two sentences per choice is necessary, but be sure to be specific in addressing what would have happened without each given implementation detail, why such a situation would be suboptimal, and how the proposed implementation improves the situation.

Your Answer:

1. **Using multiple attention heads** allows the model to jointly attend to information from different representation subspaces at different positions. A single attention head might not capture all relevant patterns or dependencies in the data, whereas multiple heads enable the model to learn diverse features and improve its ability to handle complex relationships.
2. **Dividing by ($\sqrt{d/h}$)** before applying the softmax function helps prevent the softmax from saturating and having extremely small gradients, especially when the feature dimension (d) is large. Without this scaling, the dot products of queries and keys can become large in magnitude, pushing the softmax function into regions where it becomes less sensitive to changes in input values. This scaling ensures that the attention weights are computed in a numerically stable manner.
3. **Adding a linear transformation to the output of the attention operation** allows the model to project the concatenated outputs of multiple attention heads into a space that matches the dimensionality of the input. This linear transformation helps in integrating

information from different heads and ensures that the output of the attention mechanism can be effectively used in subsequent layers of the model. Without this transformation, the model would not be able to combine the information from different heads in a meaningful way.

Transformer for Image Captioning

Now that you have implemented the previous layers, you can combine them to build a Transformer-based image captioning model. Open the file `cs231n/classifiers/transformer.py` and look at the `CaptioningTransformer` class.

Implement the `forward` function of the class. After doing so, run the following to check your forward pass using a small test case; you should see error on the order of `e-5` or less.

```
torch.manual_seed(231)
np.random.seed(231)

N, D, W = 4, 20, 30
word_to_idx = {'<NULL>': 0, 'cat': 2, 'dog': 3}
V = len(word_to_idx)
T = 3

transformer = CaptioningTransformer(
    word_to_idx,
    input_dim=D,
    wordvec_dim=W,
    num_heads=2,
    num_layers=2,
    max_length=30
)

# Set all model parameters to fixed values
for p in transformer.parameters():
    p.data = torch.tensor(np.linspace(-1.4, 1.3,
num=p.numel()).reshape(*p.shape))

features = torch.tensor(np.linspace(-1.5, 0.3, num=(N * D)).reshape(N, D))
captions = torch.tensor((np.arange(N * T) % V).reshape(N, T))

scores = transformer(features, captions)
expected_scores = np.asarray([[[-16.9532,     4.8261,   26.6054],
                               [-17.1033,     4.6906,   26.4844],
                               [-15.0708,     4.1108,   23.2924]],
                              [[-17.1767,     4.5897,   26.3562],
                               [-15.6017,     4.8693,   25.3403],
                               [-15.1028,     4.6905,   24.4839]],
                              [[-17.2172,     4.7701,   26.7574],
                               [-16.6755,     4.8500,   26.3754],
                               [-17.2172,     4.7701,   26.7574]],
                              [[-16.3669,     4.1602,   24.6872],
                               [-16.7897,     4.3467,   25.4831],
                               [-17.0103,     4.7775,   26.5652]]])
print('scores error: {:.2e}'.format(rel_error(expected_scores,
scores.detach().numpy())))
```

```
scores error: 5.06e-06
```

Overfit Transformer Captioning Model on Small Data

Run the following to overfit the Transformer-based captioning model on the same small dataset as we used for the RNN previously.

```
torch.manual_seed(231)
np.random.seed(231)

data = load_coco_data(max_train=50)

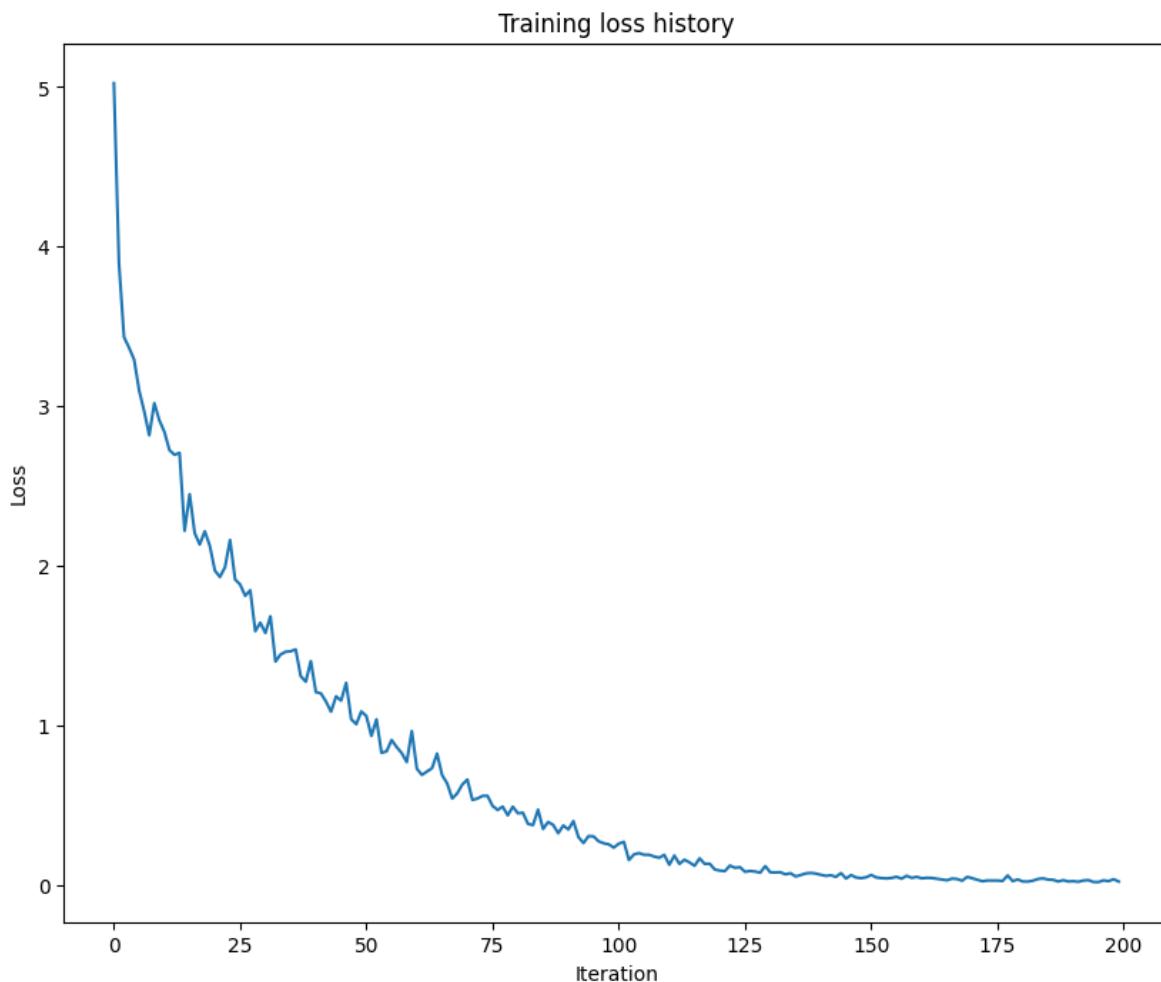
transformer = CaptioningTransformer(
    word_to_idx=data['word_to_idx'],
    input_dim=data['train_features'].shape[1],
    wordvec_dim=256,
    num_heads=2,
    num_layers=2,
    max_length=30
)

transformer_solver = CaptioningSolverTransformer(transformer, data,
idx_to_word=data['idx_to_word'],
    num_epochs=100,
    batch_size=25,
    learning_rate=0.001,
    verbose=True, print_every=10,
)
transformer_solver.train()

# Plot the training losses.
plt.plot(transformer_solver.loss_history)
plt.xlabel('Iteration')
plt.ylabel('Loss')
plt.title('Training loss history')
plt.show()
```

```
base dir /content/drive/My
Drive/cs231n/assignments/assignment3/cs231n/datasets/coco_captioning
(Iteration 1 / 200) loss: 5.023862
(Iteration 11 / 200) loss: 2.838942
(Iteration 21 / 200) loss: 1.969214
(Iteration 31 / 200) loss: 1.578360
(Iteration 41 / 200) loss: 1.207806
(Iteration 51 / 200) loss: 1.057401
(Iteration 61 / 200) loss: 0.727205
(Iteration 71 / 200) loss: 0.659644
(Iteration 81 / 200) loss: 0.448454
(Iteration 91 / 200) loss: 0.346852
(Iteration 101 / 200) loss: 0.257898
(Iteration 111 / 200) loss: 0.125795
```

```
(Iteration 121 / 200) loss: 0.088954
(Iteration 131 / 200) loss: 0.078130
(Iteration 141 / 200) loss: 0.062350
(Iteration 151 / 200) loss: 0.061456
(Iteration 161 / 200) loss: 0.040430
(Iteration 171 / 200) loss: 0.040600
(Iteration 181 / 200) loss: 0.020886
(Iteration 191 / 200) loss: 0.022638
```



Print final training loss. You should see a final loss of less than 0.03.

```
print('Final loss: ', transformer_solver.loss_history[-1])
```

```
Final loss: 0.020350398
```

Transformer Sampling at Test Time

The sampling code has been written for you. You can simply run the following to compare with the previous results with the RNN. As before the training results should be much better than the validation set results, given how little data we trained on.

```
# If you get an error, the URL just no longer exists, so don't worry!
# You can re-sample as many times as you want.
for split in ['train', 'val']:
```

```
minibatch = sample_coco_minibatch(data, split=split, batch_size=2)
gt_captions, features, urls = minibatch
gt_captions = decodeCaptions(gt_captions, data['idx_to_word'])

sampleCaptions = transformer.sample(features, max_length=30)
sampleCaptions = decodeCaptions(sampleCaptions, data['idx_to_word'])

for gt_caption, sample_caption, url in zip(gt_captions, sampleCaptions,
urls):
    img = imageFromUrl(url)
    # Skip missing URLs.
    if img is None: continue
    plt.imshow(img)
    plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_caption))
    plt.axis('off')
    plt.show()
```

Output hidden; open in <https://colab.research.google.com> to view.

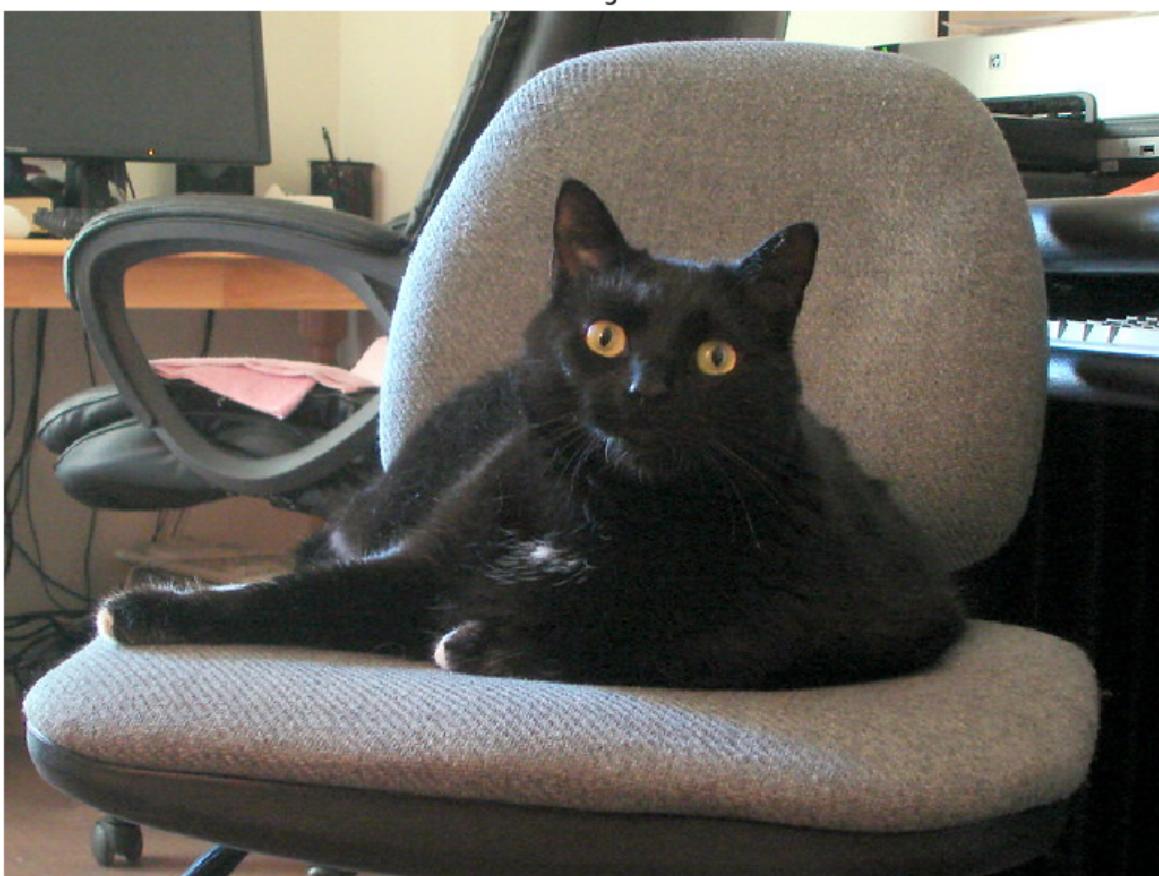
train

a large truck parked on the cement in front of a building <END>
GT:<START> a large truck parked on the cement in front of a building <END>



train

a black cat sitting in an office chair <END>
GT:<START> a black cat sitting in an office chair <END>



val

a plane that is <UNK> and a school coming bus <END>
GT:<START> a city train stopped at the train station <END>



val

a dog is <UNK> <UNK> during a bed of a <END>
GT:<START> two men playing tennis outside during the day <END>



Network Visualization (PyTorch)

In this notebook we will explore the use of *image gradients* for generating new images.

When training a model, we define a loss function which measures our current unhappiness with the model's performance; we then use backpropagation to compute the gradient of the loss with respect to the model parameters, and perform gradient descent on the model parameters to minimize the loss.

Here we will do something slightly different. We will start from a convolutional neural network model which has been pretrained to perform image classification on the ImageNet dataset. We will use this model to define a loss function which quantifies our current unhappiness with our image, then use backpropagation to compute the gradient of this loss with respect to the pixels of the image. We will then keep the model fixed, and perform gradient descent *on the image* to synthesize a new image which minimizes the loss.

In this notebook we will explore three techniques for image generation:

1. **Saliency Maps:** Saliency maps are a quick way to tell which part of the image influenced the classification decision made by the network.
2. **Fooling Images:** We can perturb an input image so that it appears the same to humans, but will be misclassified by the pretrained network.
3. **Class Visualization:** We can synthesize an image to maximize the classification score of a particular class; this can give us some sense of what the network is looking for when it classifies images of that class.

This notebook uses **PyTorch**; we have provided another notebook which explores the same concepts in TensorFlow. You only need to complete one of these two notebooks.

```
# set workspace and get datasets
from google.colab import drive
drive.mount('/content/drive')
FOLDERNAME = "cs231n/assignments/assignment3/"
assert FOLDERNAME is not None, "[!] Enter the foldername."
import sys
sys.path.append('/content/drive/My Drive/{}'.format(FOLDERNAME))
%cd /content/drive/My\ Drive/$FOLDERNAME/cs231n/datasets/
!bash get_datasets.sh
%cd /content/drive/My\ Drive/$FOLDERNAME
```

```
Mounted at /content/drive
/content/drive/My Drive/cs231n/assignments/assignment3/cs231n/datasets
/content/drive/My Drive/cs231n/assignments/assignment3
```

```
import torch
import torchvision
import numpy as np
import random
import matplotlib.pyplot as plt
from PIL import Image
```

```

from cs231n.image_utils import SQUEEZENET_MEAN, SQUEEZENET_STD

%matplotlib inline
plt.rcParams['figure.figsize'] = (10.0, 8.0) # set default size of plots
plt.rcParams['image.interpolation'] = 'nearest'
plt.rcParams['image.cmap'] = 'gray'

# for auto-reloading external modules
# see http://stackoverflow.com/questions/1907993/autoreload-of-modules-in-
ipython
%load_ext autoreload
%autoreload 2

```

Helper Functions

Our pretrained model was trained on images that had been preprocessed by subtracting the per-color mean and dividing by the per-color standard deviation. We define a few helper functions for performing and undoing this preprocessing in `cs231n/net_visualization_pytorch`. You don't need to do anything here.

```

from cs231n.net_visualization_pytorch import preprocess, deprocess, rescale,
blur_image

```

Pretrained Model

For all of our image generation experiments, we will start with a convolutional neural network which was pretrained to perform image classification on ImageNet. We can use any model here, but for the purposes of this assignment we will use SqueezeNet [1], which achieves accuracies comparable to AlexNet but with a significantly reduced parameter count and computational complexity.

Using SqueezeNet rather than AlexNet or VGG or ResNet means that we can easily perform all image generation experiments on CPU.

[1] Iandola et al, "SqueezeNet: AlexNet-level accuracy with 50x fewer parameters and < 0.5MB model size", arXiv 2016

```

# Download and load the pretrained SqueezeNet model.
model = torchvision.models.squeezenet1_1(pretrained=True)

# we don't want to train the model, so tell PyTorch not to compute gradients
# with respect to model parameters.
for param in model.parameters():
    param.requires_grad = False

# you may see warning regarding initialization deprecated, that's fine, please
continue to next steps

```

```
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:208:  
UserWarning: The parameter 'pretrained' is deprecated since 0.13 and may be  
removed in the future, please use 'weights' instead.  
    warnings.warn(  
/usr/local/lib/python3.11/dist-packages/torchvision/models/_utils.py:223:  
UserWarning: Arguments other than a weight enum or `None` for 'weights' are  
deprecated since 0.13 and may be removed in the future. The current behavior is  
equivalent to passing `weights=SqueezeNet1_1_Weights.IMGNET1K_V1`. You can also  
use `weights=SqueezeNet1_1_Weights.DEFAULT` to get the most up-to-date weights.  
    warnings.warn(msg)  
Downloading: "https://download.pytorch.org/models/squeezeimagenet1_1-b8a52dc0.pth" to  
/root/.cache/torch/hub/checkpoints/squeezeimagenet1_1-b8a52dc0.pth  
100%|██████████| 4.73M/4.73M [00:00<00:00, 91.1MB/s]
```

Load some ImageNet images

We have provided a few example images from the validation set of the ImageNet ILSVRC 2012 Classification dataset. To download these images, descend into `cs231n/datasets/` and run

```
get_imagenet_val.sh
```

Since they come from the validation set, our pretrained model did not see these images during training.

Run the following cell to visualize some of these images, along with their ground-truth labels.

```
from cs231n.data_utils import load_imagenet_val  
X, y, class_names = load_imagenet_val(num=5)  
  
plt.figure(figsize=(12, 6))  
for i in range(5):  
    plt.subplot(1, 5, i + 1)  
    plt.imshow(X[i])  
    plt.title(class_names[y[i]])  
    plt.axis('off')  
plt.gcf().tight_layout()
```



Saliency Maps

Using this pretrained model, we will compute class saliency maps as described in Section 3.1 of [2].

A **saliency map** tells us the degree to which each pixel in the image affects the classification score for that image. To compute it, we compute the gradient of the unnormalized score corresponding to the correct class (which is a scalar) with respect to the pixels of the image. If the image has shape `(3, H, W)` then this gradient will also have shape `(3, H, W)`; for each pixel in the image,

this gradient tells us the amount by which the classification score will change if the pixel changes by a small amount. To compute the saliency map, we take the absolute value of this gradient, then take the maximum value over the 3 input channels; the final saliency map thus has shape (H, W) and all entries are nonnegative.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

Hint: PyTorch gather method

Recall in Assignment 1 you needed to select one element from each row of a matrix; if s is an numpy array of shape (N, C) and y is a numpy array of shape $(N,)$ containing integers $0 \leq y[i] < C$, then $s[np.arange(N), y]$ is a numpy array of shape $(N,)$ which selects one element from each element in s using the indices in y .

In PyTorch you can perform the same operation using the `gather()` method. If s is a PyTorch Tensor of shape (N, C) and y is a PyTorch Tensor of shape $(N,)$ containing longs in the range $0 \leq y[i] < C$, then

```
s.gather(1, y.view(-1, 1)).squeeze()
```

will be a PyTorch Tensor of shape $(N,)$ containing one entry from each row of s , selected according to the indices in y .

run the following cell to see an example.

You can also read the documentation for [the gather method](#) and [the squeeze method](#).

```
# Example of using gather to select one entry from each row in PyTorch
def gather_example():
    N, C = 4, 5
    s = torch.randn(N, C)
    y = torch.LongTensor([1, 2, 1, 3])
    print(s)
    print(y)
    print(s.gather(1, y.view(-1, 1)).squeeze())
gather_example()
```

```
tensor([[-0.5223, -0.6737,  0.8914, -0.6262, -0.3851],
       [ 0.3857, -0.1525, -0.3810,  0.1732,  0.8880],
       [-0.1002,  0.8215,  1.4236, -0.5000,  0.8534],
       [-1.4488, -2.5667,  0.9185, -0.8976, -0.6915]])
tensor([1, 2, 1, 3])
tensor([-0.6737, -0.3810,  0.8215, -0.8976])
```

Implement `compute_saliency_maps` function inside `cs231n/net_visualization_pytorch.py`

```
# Load saliency maps computation function
from cs231n.net_visualization_pytorch import compute_saliency_maps
```

Once you have completed the implementation above, run the following to visualize some class saliency maps on our example images from the ImageNet validation set:

```

def show_saliency_maps(x, y):
    # Convert x and y from numpy arrays to Torch Tensors
    x_tensor = torch.cat([preprocess(Image.fromarray(x)) for x in X], dim=0)
    y_tensor = torch.LongTensor(y)

    # Compute saliency maps for images in x
    saliency = compute_saliency_maps(x_tensor, y_tensor, model)

    # Convert the saliency map from Torch Tensor to numpy array and show images
    # and saliency maps together.
    saliency = saliency.numpy()
    N = x.shape[0]
    for i in range(N):
        plt.subplot(2, N, i + 1)
        plt.imshow(x[i])
        plt.axis('off')
        plt.title(class_names[y[i]])
        plt.subplot(2, N, N + i + 1)
        plt.imshow(saliency[i], cmap=plt.cm.hot)
        plt.axis('off')
    plt.gcf().set_size_inches(12, 5)
    plt.show()

show_saliency_maps(x, y)

```



INLINE QUESTION

A friend of yours suggests that in order to find an image that maximizes the correct score, we can perform gradient ascent on the input image, but instead of the gradient we can actually use the saliency map in each step to update the image. Is this assertion true? Why or why not?

Your Answer:

This assertion is **not true**. The reasons are as follows:

1. Saliency Map vs. Gradient:

- o A saliency map is a visualization tool that shows which parts of the input image are most influential for the model's prediction. It typically represents the magnitude of the gradient but discards directional information (i.e., it's the absolute value or max across channels of the gradient).

- The gradient, on the other hand, contains both magnitude and direction information. It tells us not only how much the output changes with respect to the input but also **in which direction** (positive or negative) the input should be adjusted to increase the output.

2. Directional Information:

- When performing gradient ascent, the direction of the gradient is crucial. The gradient's sign indicates whether increasing or decreasing the input pixel values will increase the target score.
- The saliency map, by discarding the sign of the gradient, loses this directional information. Using it to update the image would not guide the update in the correct direction to maximize the correct score.

3. Update Mechanism:

- In gradient ascent, the image is updated by adding the gradient (scaled by a learning rate) to the image. This ensures that each update step moves the image in the direction that increases the target score.
- If we were to use the saliency map (which is the absolute gradient), the updates would not have a consistent direction. This could lead to updates that either increase or decrease the target score unpredictably.

Fooling Images

We can also use image gradients to generate "fooling images" as discussed in [3]. Given an image and a target class, we can perform gradient **ascent** over the image to maximize the target class, stopping when the network classifies the image as the target class. Implement the following function to generate fooling images.

[3] Szegedy et al, "Intriguing properties of neural networks", ICLR 2014

Implement `make_fooling_image` function inside `cs231n/net_visualization_pytorch.py`

Run the following cell to generate a fooling image. You should ideally see at first glance no major difference between the original and fooling images, and the network should now make an incorrect prediction on the fooling one. However you should see a bit of random noise if you look at the 10x magnified difference between the original and fooling images. Feel free to change the `idx` variable to explore other images.

```
from cs231n.net_visualization_pytorch import make_fooling_image
idx = 0
target_y = 3

X_tensor = torch.cat([preprocess(Image.fromarray(x)) for x in X], dim=0)
X_fooling = make_fooling_image(X_tensor[idx:idx+1], target_y, model)

scores = model(X_fooling)
assert target_y == scores.data.max(1)[0].item(), 'The model is not fooled!'
```

Fooling succeeded at iteration 14

After generating a fooling image, run the following cell to visualize the original image, the fooling image, as well as the difference between them.

```
X_fooling_np = deprocess(X_fooling.clone())
X_fooling_np = np.asarray(X_fooling_np).astype(np.uint8)
```

```

plt.subplot(1, 4, 1)
plt.imshow(x[idx])
plt.title(class_names[y[idx]])
plt.axis('off')

plt.subplot(1, 4, 2)
plt.imshow(x_fooling_np)
plt.title(class_names[target_y])
plt.axis('off')

plt.subplot(1, 4, 3)
x_pre = preprocess(Image.fromarray(x[idx]))
diff = np.asarray(deprocess(x_fooling - x_pre, should_rescale=False))
plt.imshow(diff)
plt.title('difference')
plt.axis('off')

plt.subplot(1, 4, 4)
diff = np.asarray(deprocess(10 * (x_fooling - x_pre), should_rescale=False))
plt.imshow(diff)
plt.title('Magnified difference (10x)')
plt.axis('off')

plt.gcf().set_size_inches(12, 5)
plt.show()

```



Class visualization

By starting with a random noise image and performing gradient ascent on a target class, we can generate an image that the network will recognize as the target class. This idea was first presented in [2]; [3] extended this idea by suggesting several regularization techniques that can improve the quality of the generated image.

Concretely, let I be an image and let y be a target class. Let $s_y(I)$ be the score that a convolutional network assigns to the image I for class y ; note that these are raw unnormalized scores, not class probabilities. We wish to generate an image I^* that achieves a high score for the class y by solving the problem

$$I^* = \arg \max_I (s_y(I) - R(I))$$

where R is a (possibly implicit) regularizer (note the sign of $R(I)$ in the argmax: we want to minimize this regularization term). We can solve this optimization problem using gradient ascent, computing gradients with respect to the generated image. We will use (explicit) L2 regularization of the form

$$R(I) = \lambda \|I\|_2^2$$

and implicit regularization as suggested by [3] by periodically blurring the generated image. We can solve this problem using gradient ascent on the generated image.

[2] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR Workshop 2014.

[3] Yosinski et al, "Understanding Neural Networks Through Deep Visualization", ICML 2015 Deep Learning Workshop

In `cs231n/net_visualization_pytorch.py` complete the implementation of the `image_visualization_update_step` used in the `create_class_visualization` function below. Once you have completed that implementation, run the following cells to generate an image of a Tarantula:

```
from cs231n.net_visualization_pytorch import class_visualization_update_step,
jitter, blur_image
def create_class_visualization(target_y, model, dtype, **kwargs):
    """
    Generate an image to maximize the score of target_y under a pretrained
    model.
    """

    Inputs:
    - target_y: Integer in the range [0, 1000) giving the index of the class
    - model: A pretrained CNN that will be used to generate the image
    - dtype: Torch datatype to use for computations

    Keyword arguments:
    - l2_reg: Strength of L2 regularization on the image
    - learning_rate: How big of a step to take
    - num_iterations: How many iterations to use
    - blur_every: How often to blur the image as an implicit regularizer
    - max_jitter: How much to jitter the image as an implicit regularizer
    - show_every: How often to show the intermediate result
    """
    model.type(dtype)
    l2_reg = kwargs.pop('l2_reg', 1e-3)
    learning_rate = kwargs.pop('learning_rate', 25)
    num_iterations = kwargs.pop('num_iterations', 200)
    blur_every = kwargs.pop('blur_every', 10)
    max_jitter = kwargs.pop('max_jitter', 16)
    show_every = kwargs.pop('show_every', 25)

    # Randomly initialize the image as a PyTorch Tensor, and make it requires
    # gradient.
    img = torch.randn(1, 3, 224, 224).mul_(1.0).type(dtype).requires_grad_()

    for t in range(num_iterations):
        # Randomly jitter the image a bit; this gives slightly nicer results
        ox, oy = random.randint(0, max_jitter), random.randint(0, max_jitter)
```

```

        img.data.copy_(jitter(img.data, ox, oy))
        class_visualization_update_step(img, model, target_y, l2_reg,
learning_rate)
        # Undo the random jitter
        img.data.copy_(jitter(img.data, -ox, -oy))

        # As regularizer, clamp and periodically blur the image
        for c in range(3):
            lo = float(-SQUEEZENET_MEAN[c] / SQUEEZENET_STD[c])
            hi = float((1.0 - SQUEEZENET_MEAN[c]) / SQUEEZENET_STD[c])
            img.data[:, c].clamp_(min=lo, max=hi)
        if t % blur_every == 0:
            blur_image(img.data, sigma=0.5)

        # Periodically show the image
        if t == 0 or (t + 1) % show_every == 0 or t == num_iterations - 1:
            plt.imshow(deprocess(img.data.clone().cpu()))
            class_name = class_names[target_y]
            plt.title('%s\nIteration %d / %d' % (class_name, t + 1,
num_iterations))
            plt.gcf().set_size_inches(4, 4)
            plt.axis('off')
            plt.show()

    return deprocess(img.data.cpu())

```

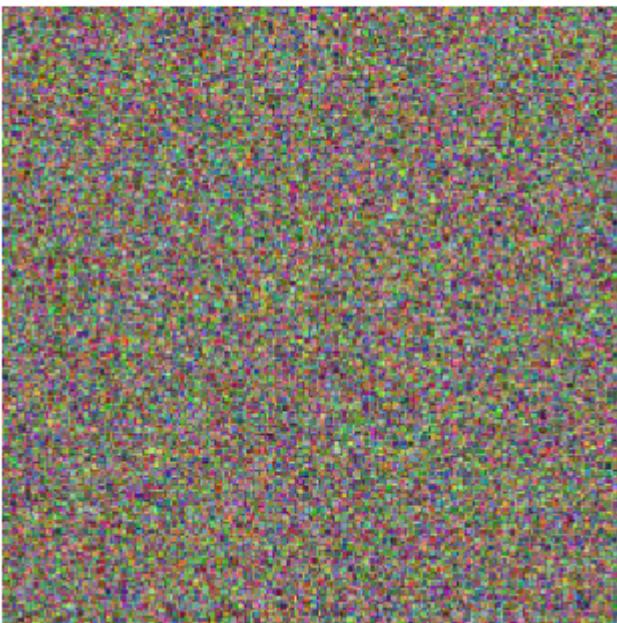
```

dtype = torch.FloatTensor
# dtype = torch.cuda.FloatTensor # Uncomment this to use GPU
model.type(dtype)

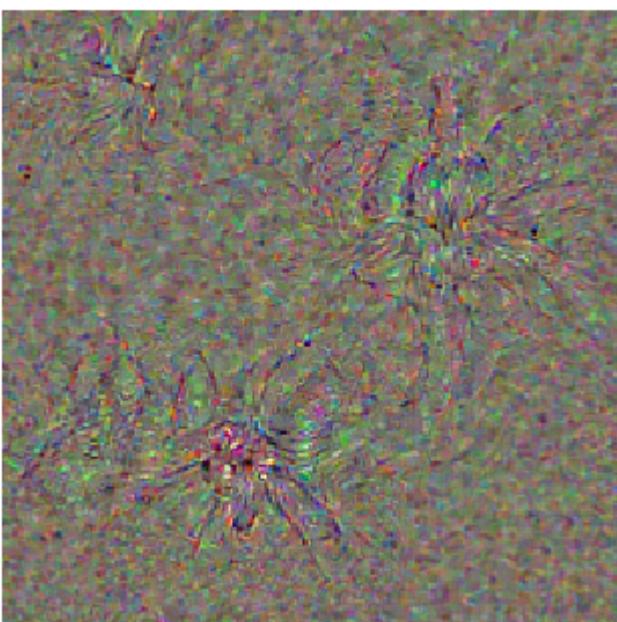
target_y = 76 # Tarantula
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
out = create_class_visualization(target_y, model, dtype)

```

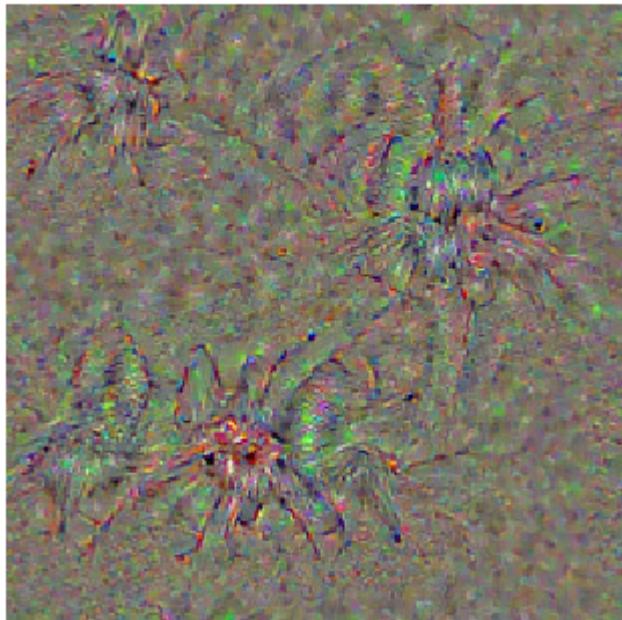
tarantula
Iteration 1 / 200



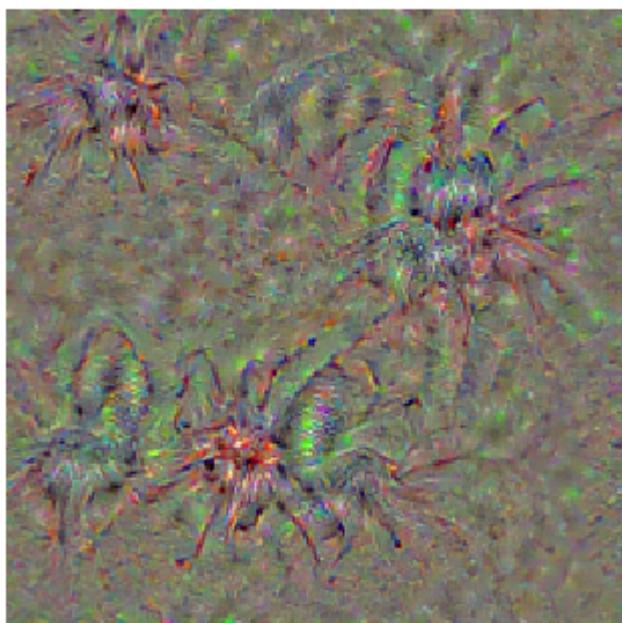
tarantula
Iteration 25 / 200



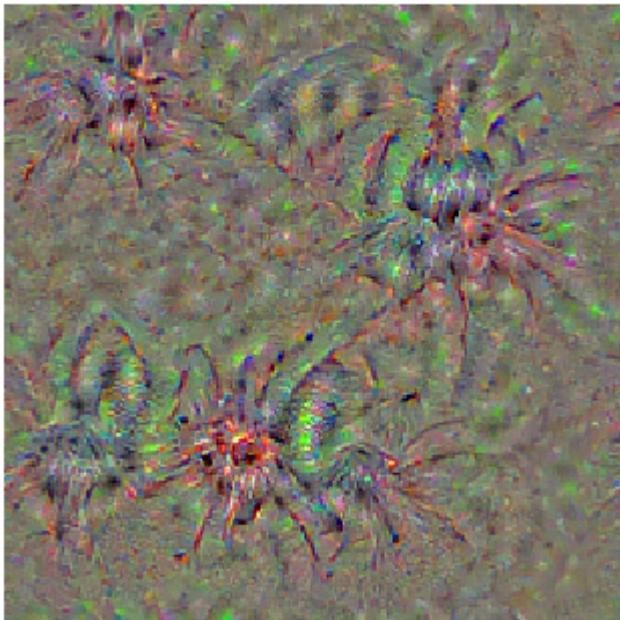
tarantula
Iteration 50 / 200



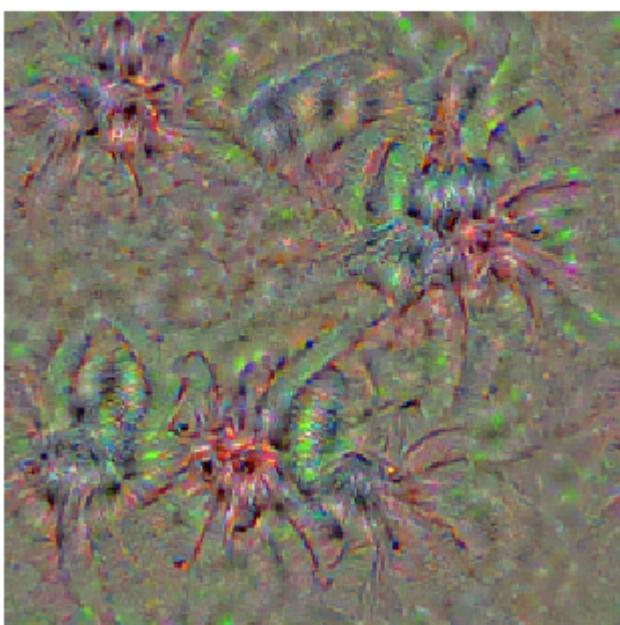
tarantula
Iteration 75 / 200



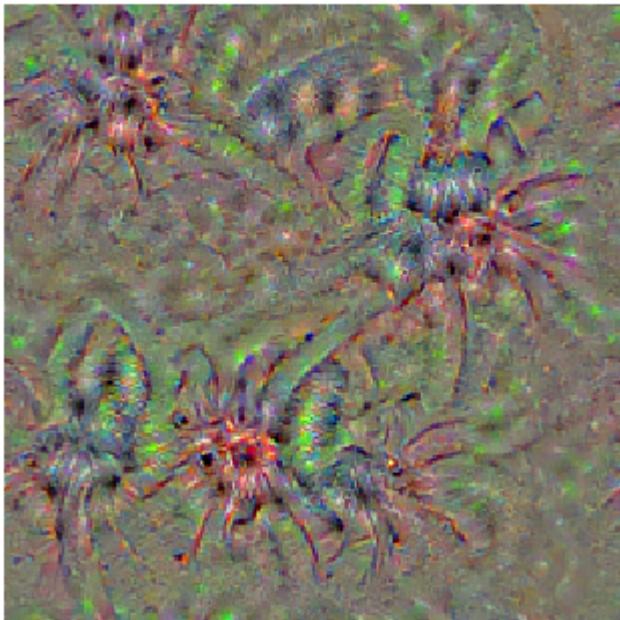
tarantula
Iteration 100 / 200



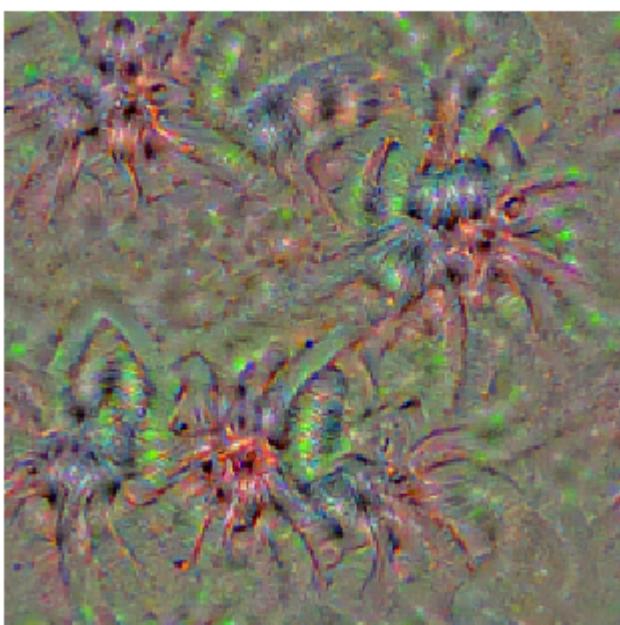
tarantula
Iteration 125 / 200



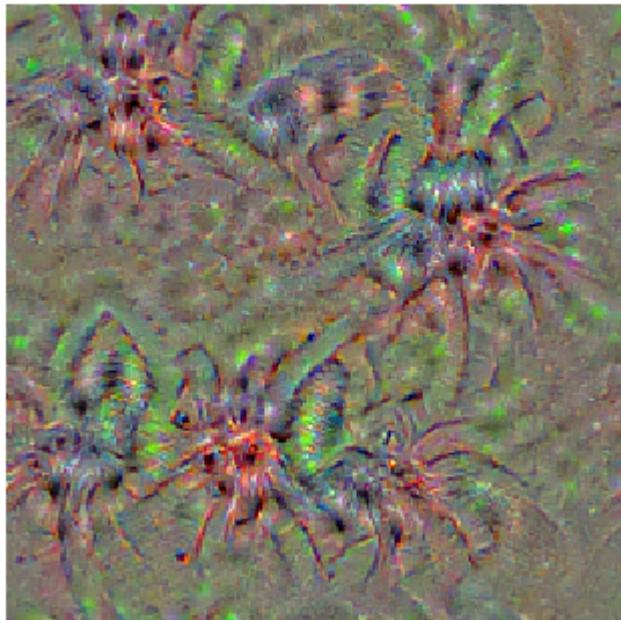
tarantula
Iteration 150 / 200



tarantula
Iteration 175 / 200



tarantula
Iteration 200 / 200

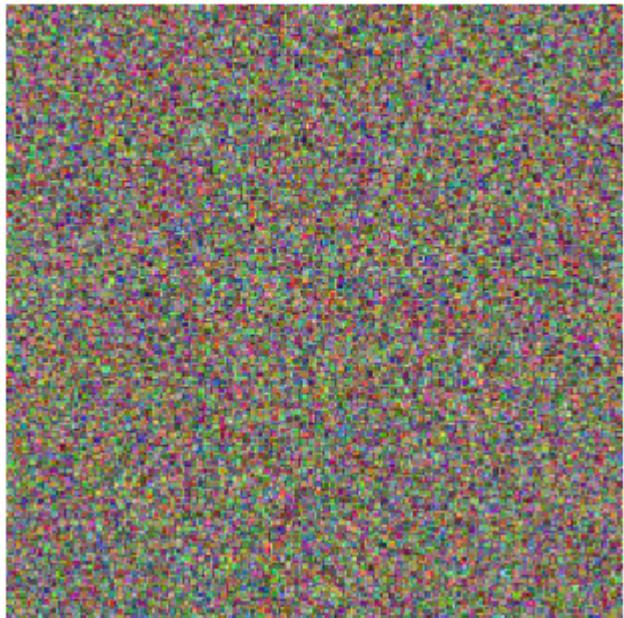


Try out your class visualization on other classes! You should also feel free to play with various hyperparameters to try and improve the quality of the generated image, but this is not required.

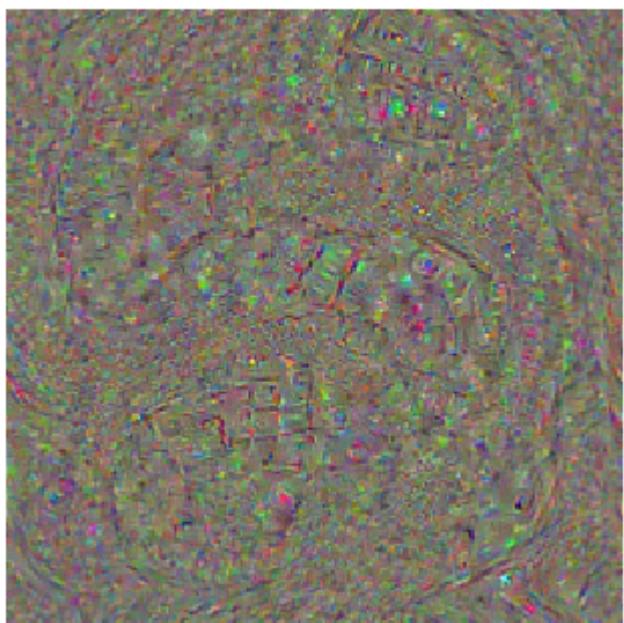
```
# target_y = 78 # Tick
# target_y = 187 # Yorkshire Terrier
# target_y = 683 # Oboe
# target_y = 366 # Gorilla
# target_y = 604 # Hourglass
target_y = np.random.randint(1000)
print(class_names[target_y])
X = create_class_visualization(target_y, model, dtype)
```

digital watch

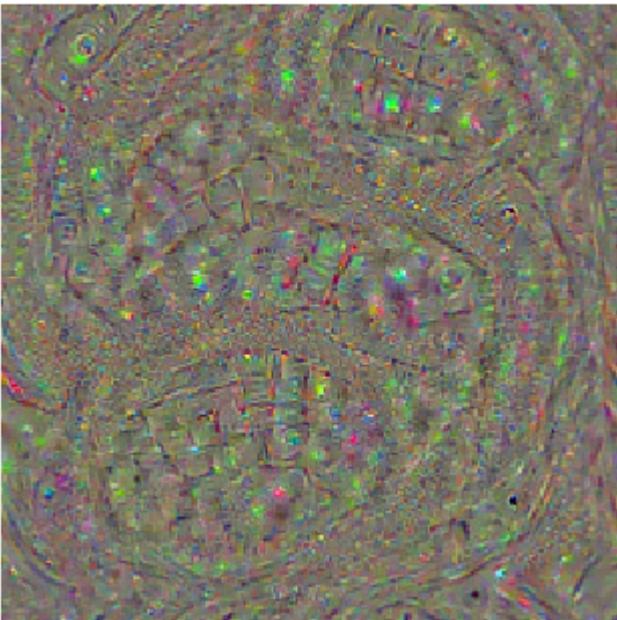
digital watch
Iteration 1 / 200



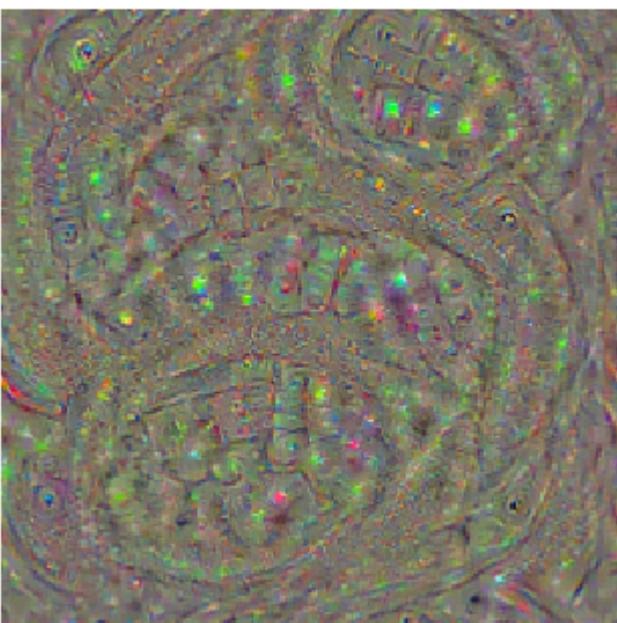
digital watch
Iteration 25 / 200



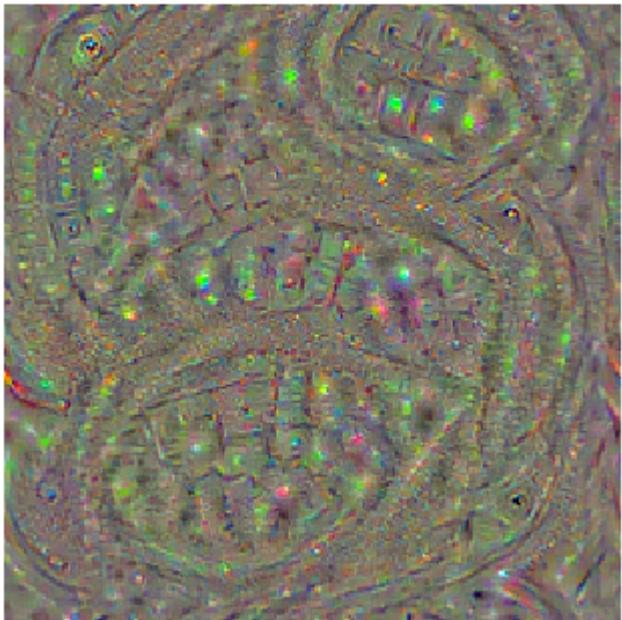
digital watch
Iteration 50 / 200



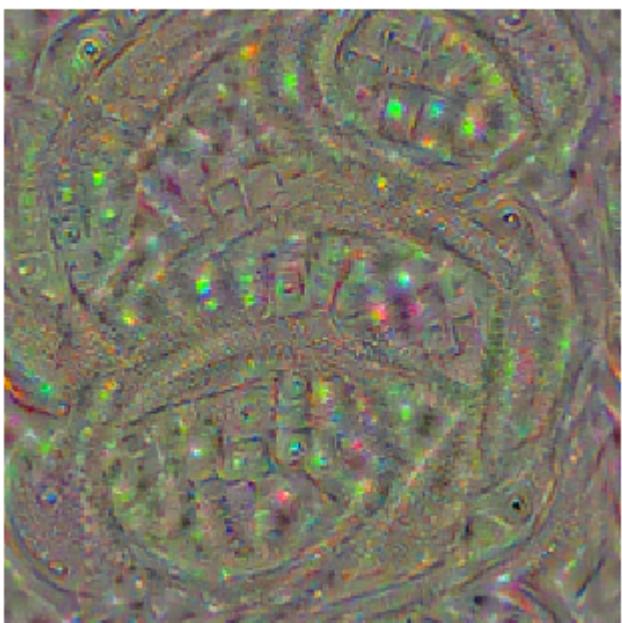
digital watch
Iteration 75 / 200



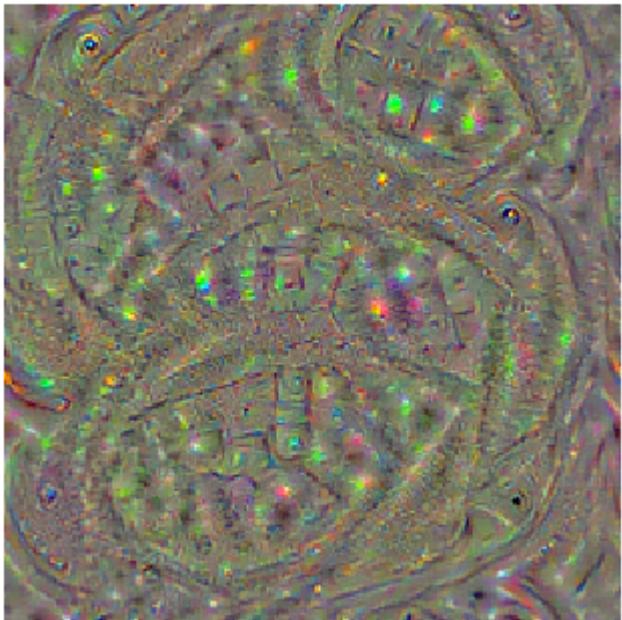
digital watch
Iteration 100 / 200



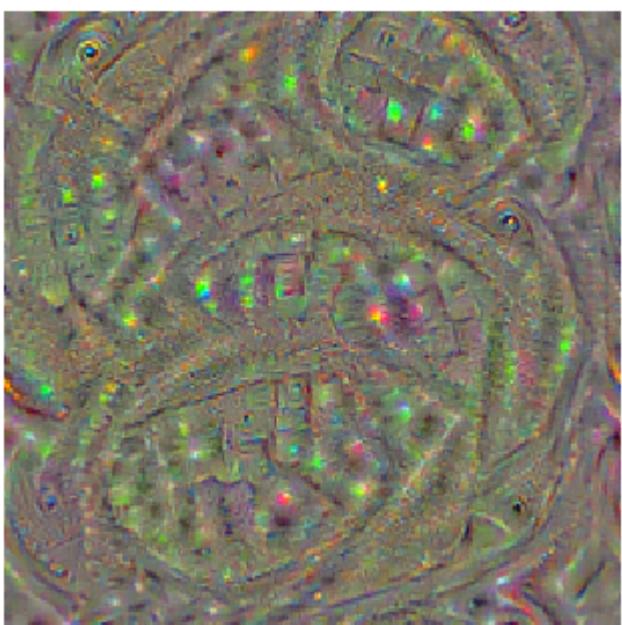
digital watch
Iteration 125 / 200



digital watch
Iteration 150 / 200



digital watch
Iteration 175 / 200



digital watch
Iteration 200 / 200

