# Image Captioning with LSTMs

```
In [12]:  # As usual, a bit of setup
          from __future__ import print_function
          import time, os, json
          import numpy as np
          import matplotlib.pyplot as plt

          from deeplearning.gradient_check import eval_numerical_gradient, eval
          _numerical_gradient_array
          from deeplearning.rnn_layers import *
          from deeplearning.captioning_solver import CaptioningSolver
          from deeplearning.classifiers.rnn import CaptioningRNN
          from deeplearning.coco_utils import load_coco_data, sample_coco_minib
          atch, decode_captions
          from deeplearning.image_utils import image_from_url

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

          # for auto-reloading external modules
          # see http://stackoverflow.com/questions/1907993/autoreload-of-module
          s-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.ab
          s(y))))
```

```
The autoreload extension is already loaded. To reload it, use:
  %reload_ext autoreload
```

# Load MS-COCO data

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

In [13]:
```python
# Load COCO data from disk; this returns a dictionary
# We'll work with dimensionality-reduced features for this notebook
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))
```

```
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
```

In [ ]:

In [ ]:

In [ ]:

In [ ]:

In [ ]:

# Overfit LSTM captioning model

In [19]:
```python
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=100,
          batch_size=50,
          optim_config={
            'learning_rate': 5e-3,
          },
          lr_decay=0.995,
          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()
```
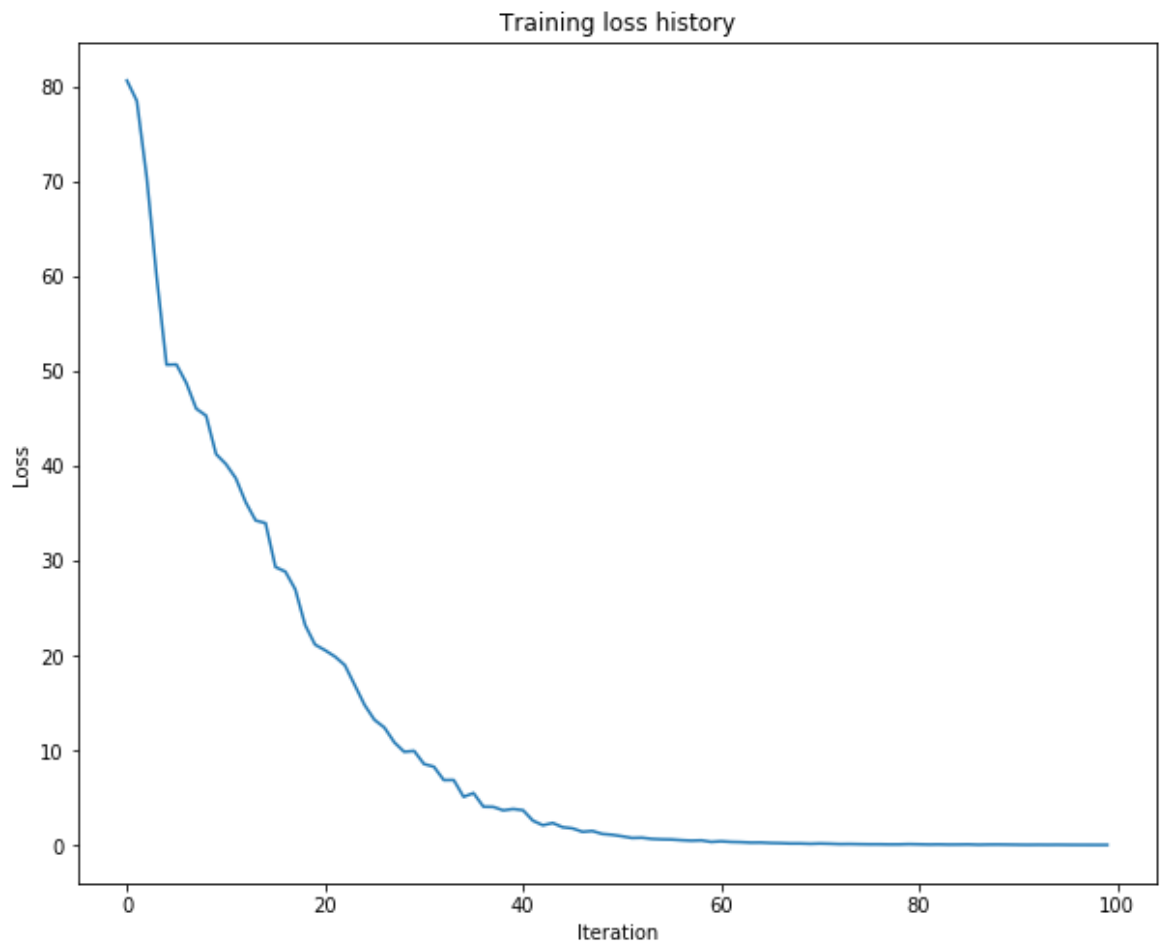
```
(Iteration 1 / 100) loss: 80.600535
(Iteration 11 / 100) loss: 40.194296
(Iteration 21 / 100) loss: 20.591769
(Iteration 31 / 100) loss: 8.602131
(Iteration 41 / 100) loss: 3.736282
(Iteration 51 / 100) loss: 0.991312
(Iteration 61 / 100) loss: 0.463356
(Iteration 71 / 100) loss: 0.227835
(Iteration 81 / 100) loss: 0.142447
(Iteration 91 / 100) loss: 0.098880
```



Training loss history

# LSTM test-time sampling

```
In [20]:  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_w
          ord'])

              for gt_caption, sample_caption, url in zip(gt_captions, sample_ca
          ptions, urls):
                  plt.imshow(image_from_url(url))
                  plt.title('%s\n%s\nGT:%s' % (split, sample_caption, gt_captio
          n))
                  plt.axis('off')
                  plt.show()
```

train
a bedroom with a striped &lt;UNK&gt; and red walls &lt;END&gt;
GT:&lt;START&gt; a bedroom with a striped &lt;UNK&gt; and red walls &lt;END&gt;

train
a truck that is &lt;UNK&gt; &lt;UNK&gt; and a school bus &lt;END&gt;
GT:&lt;START&gt; a truck that is &lt;UNK&gt; &lt;UNK&gt; and a school bus &lt;END&gt;

val
a boy is sitting on the cement in his hand <END>
GT:<START> a man is standing near a <UNK> bed <END>

val
half a &lt;UNK&gt; &lt;UNK&gt; on the &lt;UNK&gt; &lt;END&gt;
GT:&lt;START&gt; a &lt;UNK&gt; in a small boat leaves &lt;UNK&gt; in the water &lt;END&gt;

Train a good captioning model!

In [21]:
```python
import nltk

def BLEU_score(gt_caption, sample_caption):
    """
    gt_caption: string, ground-truth caption
    sample_caption: string, your model's predicted caption
    Returns unigram BLEU score.
    """
    reference = [x for x in gt_caption.split(' ')
                 if ('<END>' not in x and '<START>' not in x and '<UN
K>' not in x)]
    hypothesis = [x for x in sample_caption.split(' ')
                  if ('<END>' not in x and '<START>' not in x and '<U
NK>' not in x)]
    BLEUscore = nltk.translate.bleu_score.sentence_bleu([reference],
hypothesis, weights = [1])
    return BLEUscore

def evaluate_model(model):
    """
    model: CaptioningRNN model
    Prints unigram BLEU score averaged over 1000 training and val exa
mples.
    """
    BLEUscores = {}
    for split in ['train', 'val']:
        minibatch = sample_coco_minibatch(data, split=split, batch_si
ze=1000)
        gt_captions, features, urls = minibatch
        gt_captions = decode_captions(gt_captions, data['idx_to_word'
])

        sample_captions = model.sample(features)
        sample_captions = decode_captions(sample_captions, data['idx_
to_word'])

        total_score = 0.0
        for gt_caption, sample_caption, url in zip(gt_captions, sampl
e_captions, urls):
            total_score += BLEU_score(gt_caption, sample_caption)

        BLEUscores[split] = total_score / len(sample_captions)

    for split in BLEUscores:
        print('Average BLEU score for %s: %f' % (split, BLEUscores[sp
lit]))
# smaller_data=load_coco_data(max_train=10000)
# 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,
#       )
```

```python
# lstm_solver = CaptioningSolver(lstm_model, smaller_data,
#           update_rule='adam',
#           num_epochs=15,
#           batch_size=50,
#           optim_config={
#               'learning_rate': 10e-3,
#           },
#           lr_decay=0.8,
#           verbose=True, print_every=100,
#         )

# lstm_solver.train()

# evaluate_model(lstm_model)
```

In [ ]:

In [ ]:
```python
smaller_data=load_coco_data()
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,
        )

lstm_solver = CaptioningSolver(lstm_model, smaller_data,
          update_rule='adam',
          num_epochs=15,
          batch_size=100,
          optim_config={
              'learning_rate': 8e-3,
          },
          lr_decay=0.4,
          verbose=True, print_every=100,
        )

lstm_solver.train()

evaluate_model(lstm_model)
```

(Iteration 1 / 60015) loss: 74.972106

In [ ]: