# Lecture 9:
# Visualizing CNNs
# and
# Recurrent Neural Networks

*Tuesday February 28, 2017*

comp150dl

# Announcements!

- HW #3 is out

- Final Project proposals due this **Thursday March 2**

- Papers to read: Students should read all papers on the **Schedule** tab, and are encouraged to read as many papers as possible from the **Papers** tab.

- Next paper: **March 7** *You Only Look Once: Unified, Real-Time Object Detection.* If this paper seems too deep or confusing, look at *Fast R-CNN, Faster R-CNN*
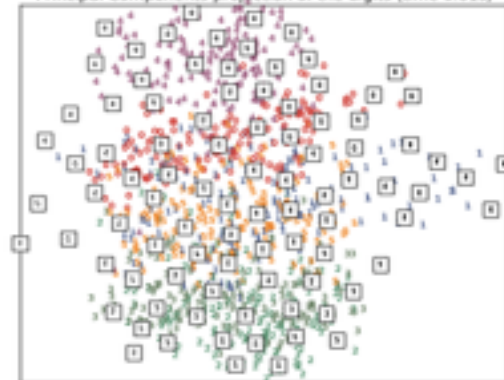
# Python/Numpy of the Day

- t-SNE (t-Distributed Stochastic Nearest Neighbor Embedding)

  - Scikit-Learn t-SNE

  - Examples of 2D Embedding Visualizations of MNIST dataset
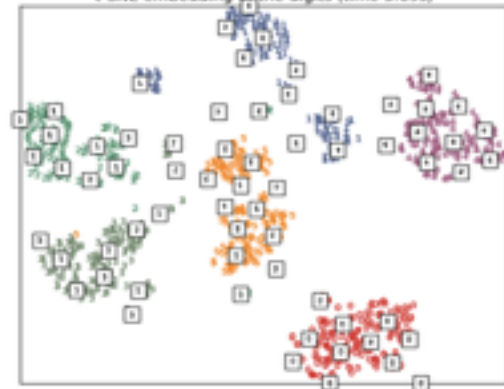
  - Other Embedding functions in Scikit-Learn

A selection from the 64-dimensional digits dataset
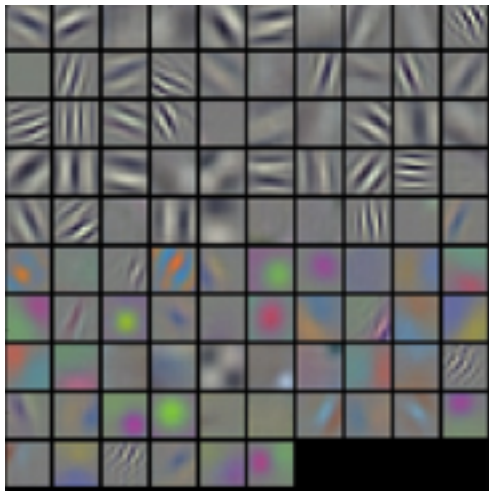
Principal Components projection of the digits (time 0.01s)

t-SNE embedding of the digits (time 5.69s)

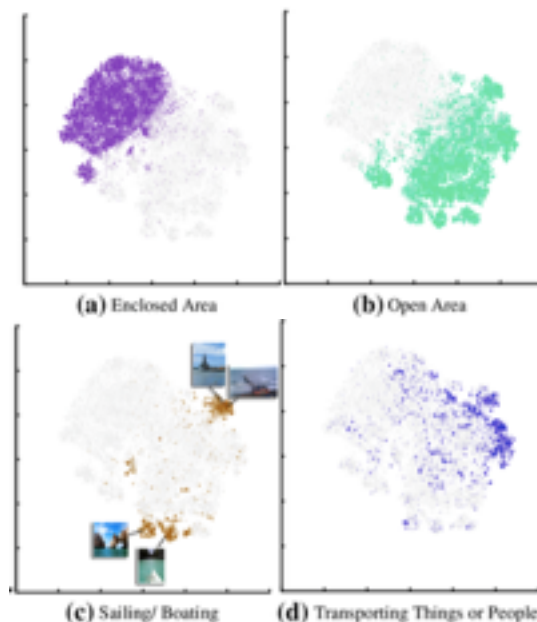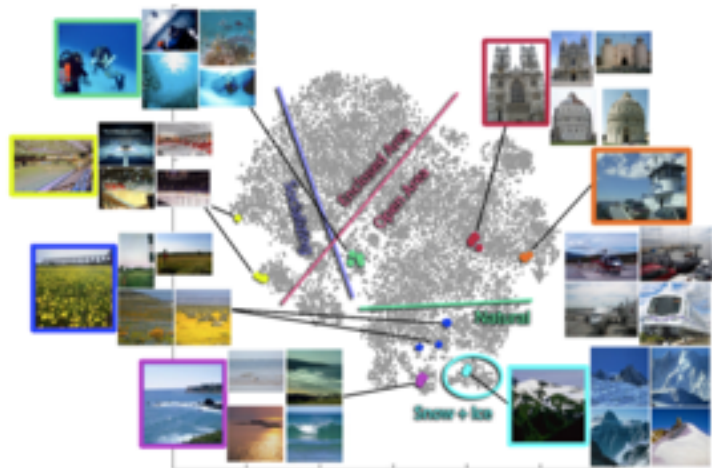# Visualizing CNN Behavior

- How can we see what's going on in a CNN?
  - *Stuff we've already done:*
    - Visualize the weights
    - Occlusion experiments — ex. Jason and Lisa's AlexNet Occlusion Tests
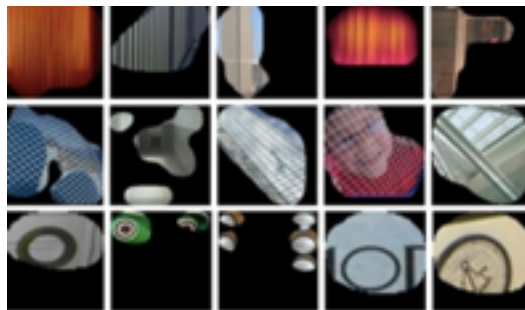
# Visualizing CNN Behavior

- How can we see what's going on in a CNN?
  - *Straightforward stuff to try in the future:*
    - Visualize the representation space (e.g. with t-SNE)
    - Human experiment comparisons





(a) Enclosed Area

(b) Open Area

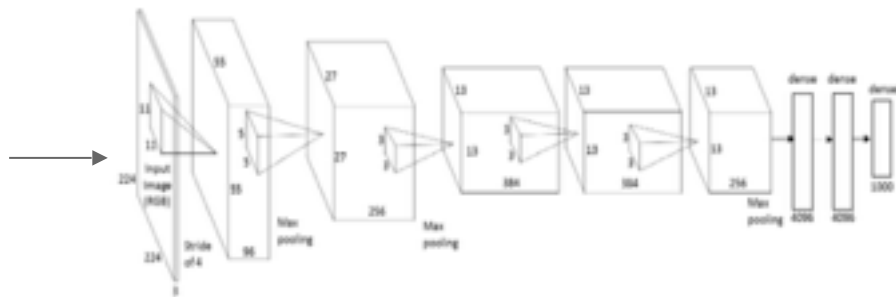(c) Sailing/ Boating

(d) Transporting Things or People

# Visualizing CNN Behavior

- How can we see what's going on in a CNN?
  - *More sophisticated approaches (HW #4)*



- Visualize patches that maximally activate neurons
- Optimization over image approaches (optimization)
- Deconv approaches (single backward pass)

# Deconv approaches -

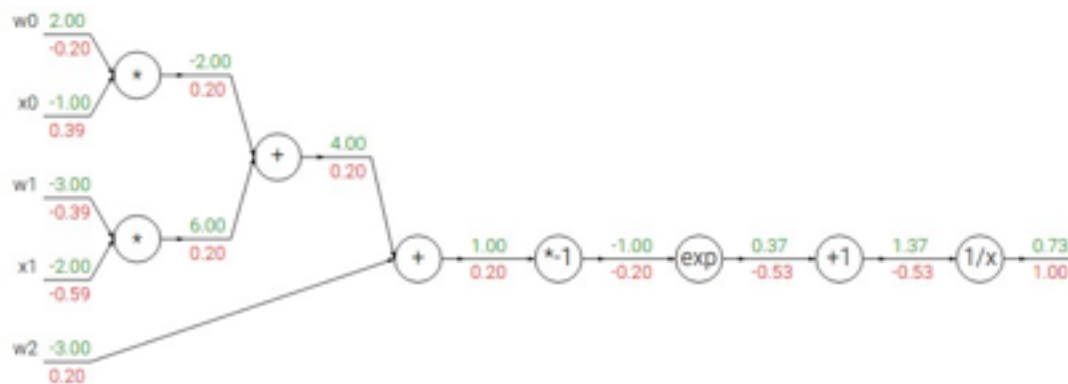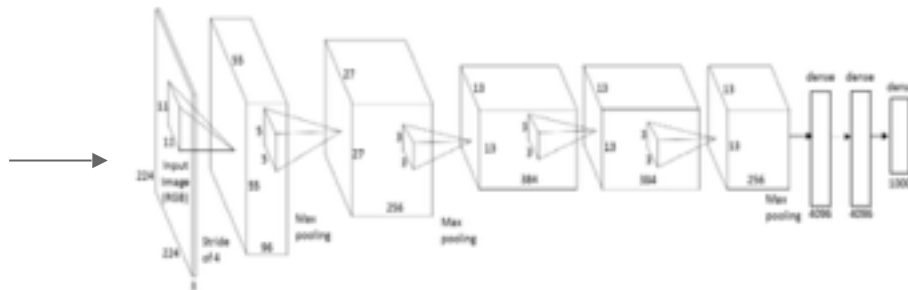projecting backward from one neuron to see what is activating it

1. *Feed image into net*



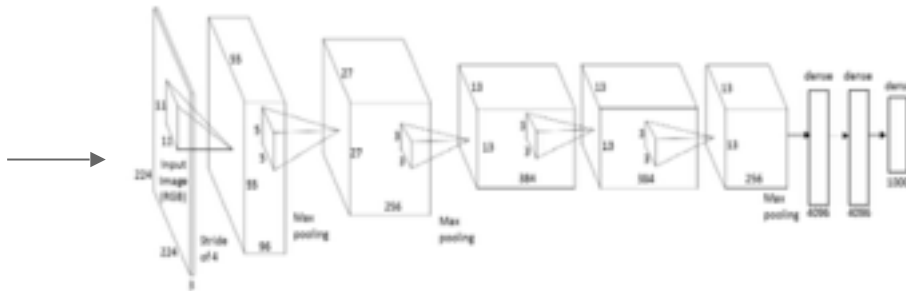*Q: how can we compute the gradient of any arbitrary neuron in the network w.r.t. the image?*

comp150dl

Tufts
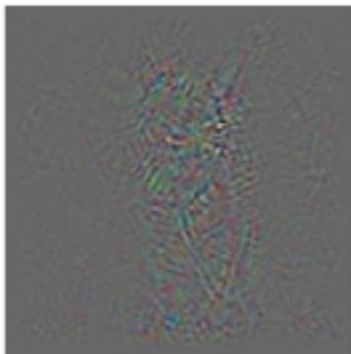
# Deconv approaches

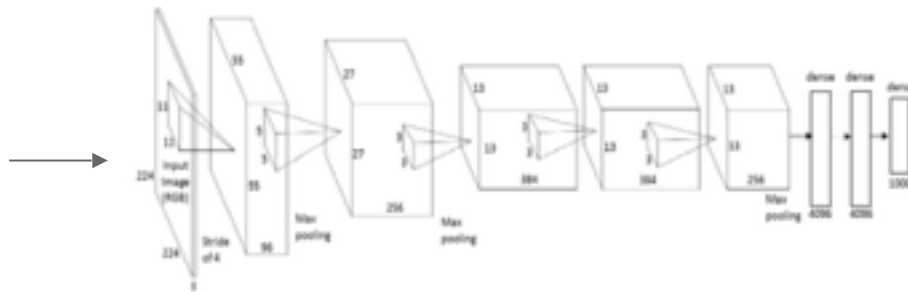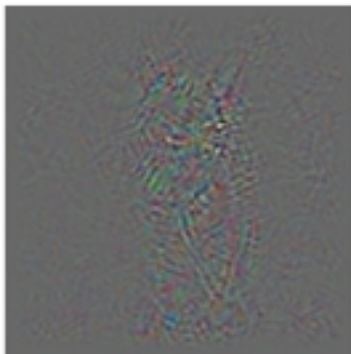## 1. Feed image into net

comp150dl

# Deconv approaches

*1. Feed image into net*



*2. Pick a layer, set the gradient there to be all zero except for one 1 for some neuron of interest*

*3. Backprop to image:*

comp150dl

# Deconv approaches

*1. Feed image into net*



*2. Pick a layer, set the gradient there to be all zero except for one 1 for some neuron of interest*
*3. Backprop to image:*
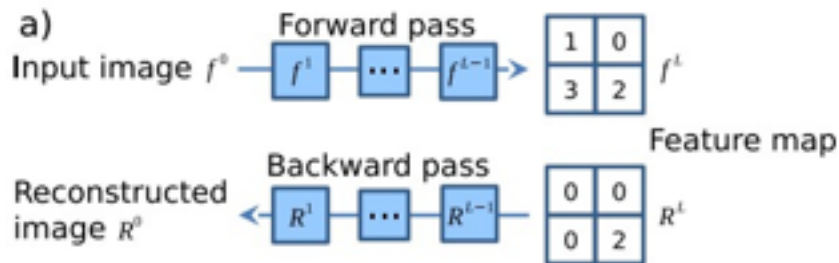


***"Guided backpropagation:"*** *only propagate positive gradients*
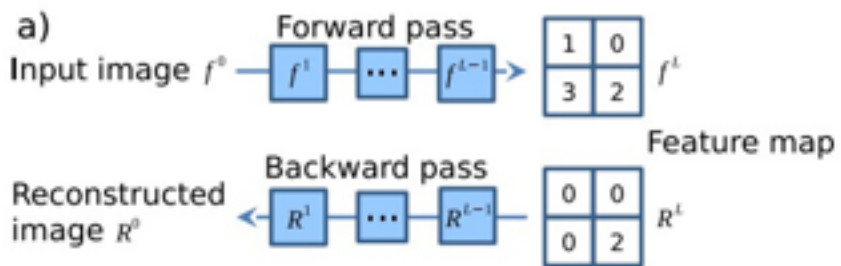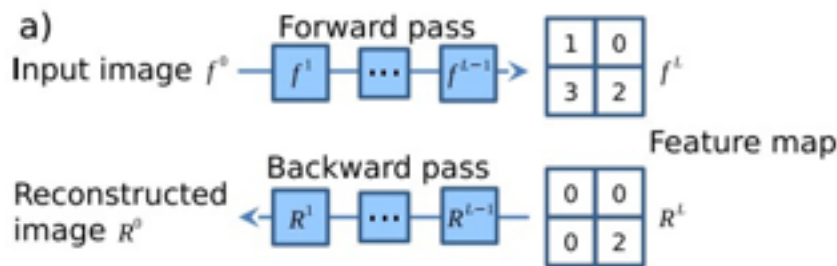
comp150dl

# Deconv approaches

*[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]*
*[Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al., 2014]*
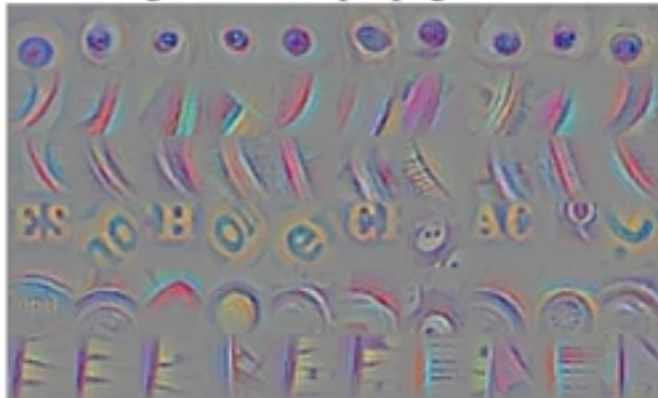*[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]*

# Deconv approaches

*[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]*
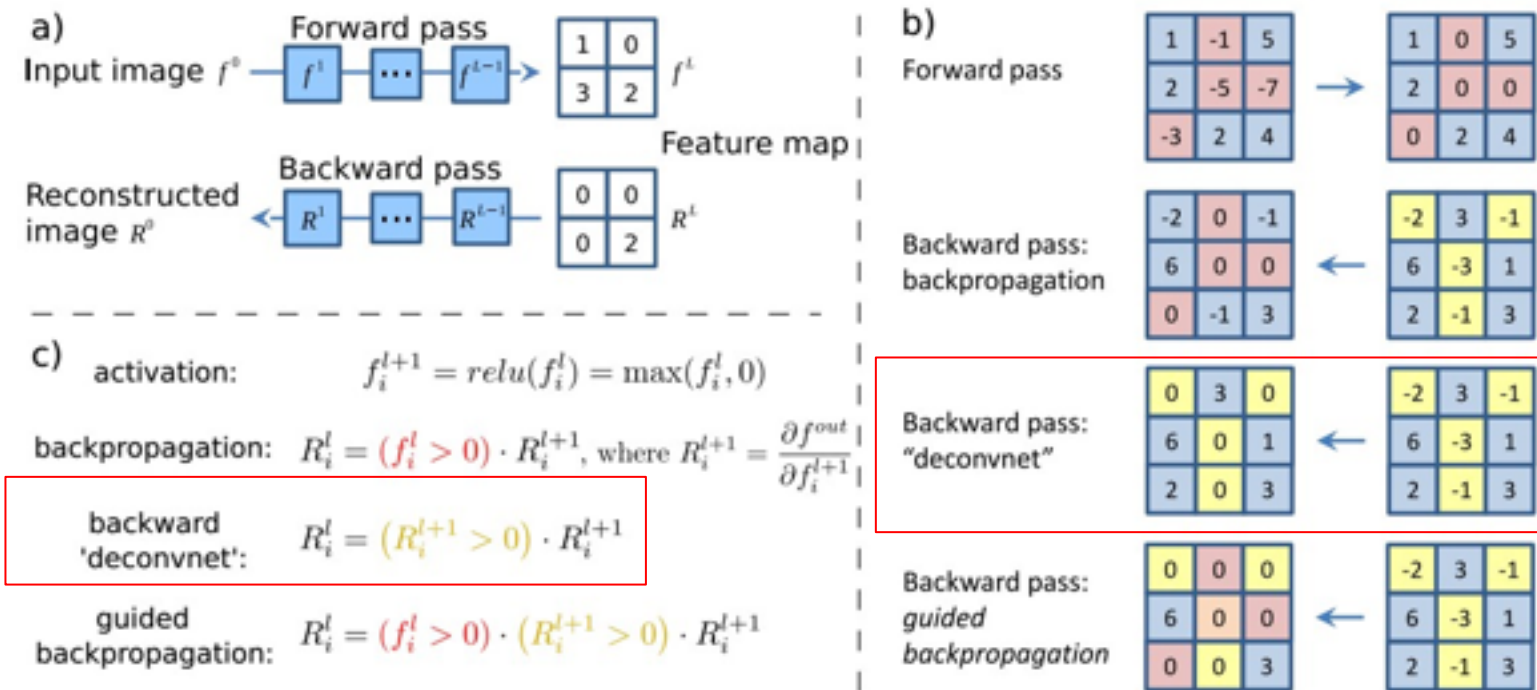*[Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al., 2014]*
*[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]*

a)
Input image $f^0$ — Forward pass — $f^1$ — $\cdots$ — $f^{L-1}$ → $f^L$ Feature map

Reconstructed image $R^0$ ← Backward pass — $R^1$ — $\cdots$ — $R^{L-1}$ ← $R^L$

b)
Forward pass

Backward pass: backpropagation

c) activation: $f_i^{l+1} = relu(f_i^l) = \max(f_i^l, 0)$

backpropagation: $R_i^l = (f_i^l > 0) \cdot R_i^{l+1}$, where $R_i^{l+1} = \dfrac{\partial f^{out}}{\partial f_i^{l+1}}$

*Backward pass for a ReLU (will be changed in Guided Backprop)*

# Deconv approaches

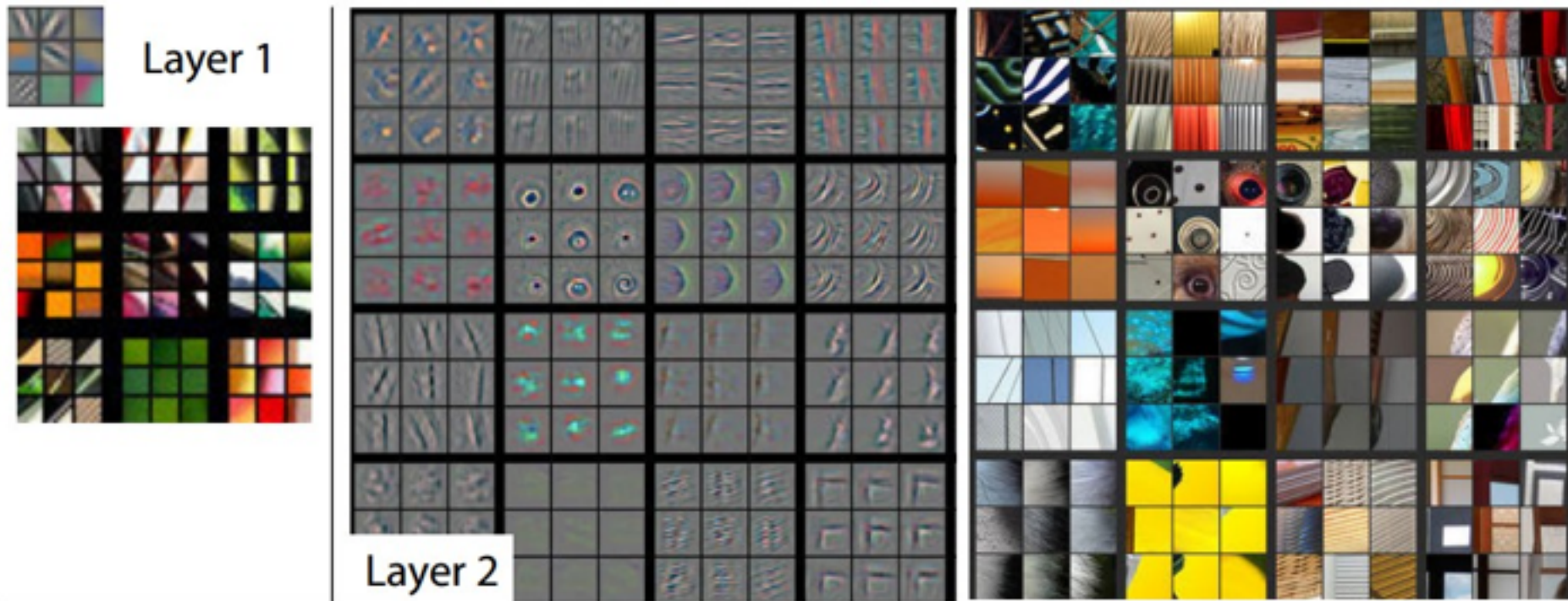*[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]*
*[Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al., 2014]*
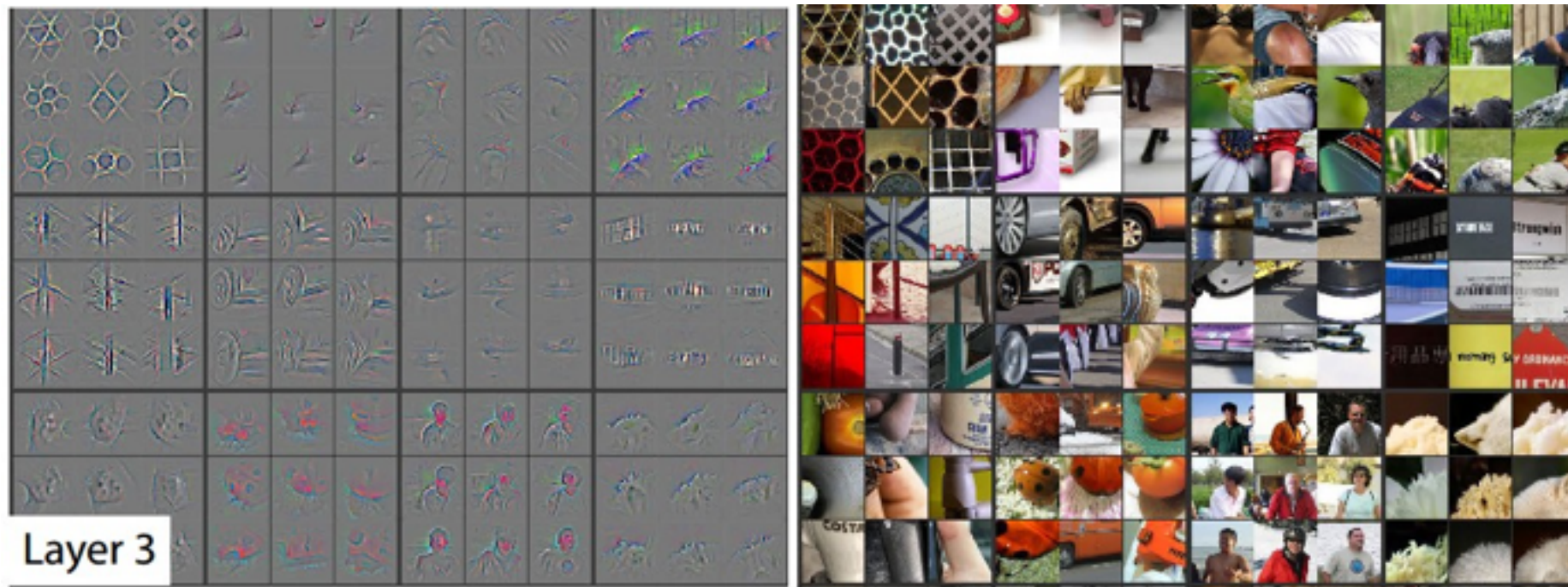*[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]*



a)
Input image $f^0$ — Forward pass

Reconstructed image $R^0$ — Backward pass — Feature map

c)
activation: $f_i^{l+1} = relu(f_i^l) = \max(f_i^l, 0)$

backpropagation: $R_i^l = (f_i^l > 0) \cdot R_i^{l+1}$, where $R_i^{l+1} = \frac{\partial f^{out}}{\partial f_i^{l+1}}$

guided backpropagation: $R_i^l = (f_i^l > 0) \cdot (R_i^{l+1} > 0) \cdot R_i^{l+1}$

b)
Forward pass

Backward pass: backpropagation

Backward pass: *guided backpropagation*

comp150dl  **Tufts** UNIVERSITY

*Visualization of patterns learned by the layer **conv6** (top) and layer **conv9** (bottom) of the network trained on ImageNet.*

*Each row corresponds to one filter.*

*The visualization using "guided backpropagation" is based on the top 10 image patches activating this filter taken from the ImageNet dataset.*



*[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]*

# Deconv approaches

*[Visualizing and Understanding Convolutional Networks, Zeiler and Fergus 2013]*
*[Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps, Simonyan et al., 2014]*
*[Striving for Simplicity: The all convolutional net, Springenberg, Dosovitskiy, et al., 2015]*



*backprops to weights that were zero-d out by ReLu*

comp150dl  Tufts

## *Visualizing arbitrary neurons along the way to the top...*



Layer 1

Layer 2

# Visualizing arbitrary neurons along the way to the top...



Layer 3

comp150dl

**Tufts**
UNIVERSITY

*Visualizing arbitrary neurons along the way to the top...*



Layer 4

Layer 5

comp150dl

Tufts

# Optimization to Image



*Q: can we find an image that maximizes some class score?*

# Optimization to Image

$$\arg \max_I \boxed{S_c(I)} - \lambda \|I\|_2^2$$

*score for class c (before Softmax)*



*Q: can we find an image that maximizes some class score?*

comp150dl

**Tufts**

# Optimization to Image

zero image

*2. set the gradient of the scores vector to be [0,0,....1,....,0], then backprop to image*

comp150dl

# Optimization to Image

*1. feed in zeros.*



zero image

*2. set the gradient of the scores vector to be [0,0,....1,....,0], then backprop to image*
*3. do a small "image update"*
*4. forward the image through the network.*
*5. go back to 2.*

$$\arg\max_{I} \boxed{S_c(I)} - \lambda \|I\|_2^2$$

*score for class c (before Softmax)*

comp150dl

# 1. Find images that maximize some class score:



dumbbell          cup          dalmatian

bell pepper       lemon        husky

comp150dl

# 1. Find images that maximize some class score:



washing machine     computer keyboard     kit fox

goose     ostrich     limousine

comp150dl

# 2. Visualize the Data gradient:

*(note that the gradient on data has three channels. Here they visualize M, s.t.:*

$$M_{ij} = \max_c \left| w_{h(i,j,c)} \right|$$

*(at each pixel take abs val, and max over channels)*

M = ?

comp150dl

# 2. Visualize the Data gradient:

*(note that the gradient on data has three channels. Here they visualize M, s.t.:*

$$M_{ij} = \max_c \left| w_{h(i,j,c)} \right|$$

*(at each pixel take abs val, and max over channels)*

comp150dl

- Use **grabcut** for segmentation
- This optimization can be done for arbitrary neurons in the CNN

# Question: Given a CNN code, is it possible to reconstruct the original image?

*Understanding Deep Image Representations by Inverting Them
[Mahendran and Vedaldi, 2014]*

# Find an image such that:
- Its code is similar to a given code
- It "looks natural" (image prior regularization)

*original image*



*reconstructions
from the 1000
log probabilities
for ImageNet
(ILSVRC)
classes*

comp150dl

DeepDream  https://github.com/google/deepdream

comp150dl  Tufts

```python
def objective_L2(dst):
    dst.diff[:] = dst.data

def make_step(net, step_size=1.5, end='inception_4c/output',
              jitter=32, clip=True, objective=objective_L2):
    '''Basic gradient ascent step.'''

    src = net.blobs['data'] # input image is stored in Net's 'data' blob
    dst = net.blobs[end]

    ox, oy = np.random.randint(-jitter, jitter+1, 2)
    src.data[0] = np.roll(np.roll(src.data[0], ox, -1), oy, -2) # apply jitter shift

    net.forward(end=end)
    objective(dst)  # specify the optimization objective
    net.backward(start=end)
    g = src.diff[0]
    # apply normalized ascent step to the input image
    src.data[:] += step_size/np.abs(g).mean() * g

    src.data[0] = np.roll(np.roll(src.data[0], -ox, -1), -oy, -2) # unshift image

    if clip:
        bias = net.transformer.mean['data']
        src.data[:] = np.clip(src.data, -bias, 255-bias)
```

```
def objective_L2(dst):
    dst.diff[:] = dst.data          DeepDream:  set dx = x :)

def make_step(net, step_size=1.5, end='inception_4c/output',
              jitter=32, clip=True, objective=objective_L2):
    '''Basic gradient ascent step.'''

    src = net.blobs['data'] # input image is stored in Net's 'data' blob
    dst = net.blobs[end]

    ox, oy = np.random.randint(-jitter, jitter+1, 2)
    src.data[0] = np.roll(np.roll(src.data[0], ox, -1), oy, -2) # apply jitter shift

    net.forward(end=end)
    objective(dst)  # specify the optimization objective
    net.backward(start=end)
    g = src.diff[0]
    # apply normalized ascent step to the input image
    src.data[:] += step_size/np.abs(g).mean() * g

    src.data[0] = np.roll(np.roll(src.data[0], -ox, -1), -oy, -2) # unshift image

    if clip:
        bias = net.transformer.mean['data']
        src.data[:] = np.clip(src.data, -bias, 255-bias)
```

*jitter regularizer*

*"image update"*

*inception_4c/output*

DeepDream modifies the image in a way that "boosts" all activations, at any layer

this creates a <u>feedback loop</u>: e.g. any slightly detected dog face will be made more and more dog like over time

comp150dl

*inception_4c/output*

"Admiral Dog!"  "The Pig-Snail"  "The Camel-Bird"  "The Dog-Fish"

*DeepDream modifies the image in a way that "boosts" all activations, at any layer*

comp150dl

*inception_3b/5x5_reduce*

*DeepDream modifies the image in a way that "boosts" all activations, at any layer*

comp150dl  **Tufts** UNIVERSITY

# NeuralStyle

[ A Neural Algorithm of Artistic Style by Leon A. Gatys,
Alexander S. Ecker, and Matthias Bethge, 2015]
**good implementation by Justin in Torch:**
**https://github.com/jcjohnson/neural-style**

We can pose an optimization over the input image to maximize any class score.
That seems useful.

Question: Can we use this to "fool" ConvNets?

comp150dl

*[Intriguing properties of neural networks, Szegedy et al., 2013]*



*correct*          *+distort*          *ostrich*          *correct*          *+distort*          *ostrich*

comp150dl

*[Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images Nguyen, Yosinski, Clune, 2014]*

*>99.6% confidences*



| robin | cheetah | armadillo | lesser panda |
| centipede | peacock | jackfruit | bubble |

*These kinds of results were around even before ConvNets…*
*[Exploring the Representation Capabilities of the HOG Descriptor, Tatu et al., 2011]*



*Identical HOG represention*

# EXPLAINING AND HARNESSING ADVERSARIAL EXAMPLES
[Goodfellow, Shlens & Szegedy, 2014]

*"primary cause of neural networks' vulnerability to adversarial perturbation is their **linear nature**"*

# Lets fool a binary linear classifier:

| x | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 | ← input example |
|---|---|----|---|----|---|---|---|----|---|---|------|
| w | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 | ← weights |

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

comp150dl  **Tufts**

# Lets fool a binary linear classifier:

| X | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 | ← input example |
|---|---|----|---|----|---|---|---|----|---|---|---|
| W | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 | ← weights |

*class 1 score = dot product:*
*= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3*
*=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474*
*i.e. the classifier is **95%** certain that this is class 0 example.*

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| **x** | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |
| **w** | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |
| *adversarial x* | ? | ? | ? | ? | ? | ? | ? | ? | ? | ? |

← *input example*

← *weights*

*class 1 score = dot product:*
*= -2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3*
*=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474*
*i.e. the classifier is **95%** certain that this is class 0 example.*

$$P(y=1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| x | 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |
|---|---|----|---|----|---|---|---|----|---|---|

← *input example*

| w | -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |
|---|----|----|---|----|---|----|---|---|----|---|

← *weights*

*adversarial x*

| 1.5 | -1.5 | 3.5 | -2.5 | 2.5 | 1.5 | 1.5 | -3.5 | 4.5 | 1.5 |
|-----|------|-----|------|-----|-----|-----|------|-----|-----|

*class 1 score before:*
*-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3*
*=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474*
*-1.5+1.5+3.5+2.5+2.5-1.5+1.5-3.5-4.5+1.5 = 2*
*=> probability of class 1 is now 1/(1+e^(-(2))) = 0.88*
***i.e. we improved the class 1 probability from 5% to 88%***

$$P(y = 1 \mid x; w, b) = \frac{1}{1 + e^{-(w^T x + b)}} = \sigma(w^T x + b)$$

# Lets fool a binary linear classifier:

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 2 | -1 | 3 | -2 | 2 | 2 | 1 | -4 | 5 | 1 |

$x$ ← *input example*

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| -1 | -1 | 1 | -1 | 1 | -1 | 1 | 1 | -1 | 1 |

$w$ ← *weights*

*adversarial x*

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 1.5 | -1.5 | 3.5 | -2.5 | 2.5 | 1.5 | 1.5 | -3.5 | 4.5 | 1.5 |

*class 1 score before:*
*-2 + 1 + 3 + 2 + 2 - 2 + 1 - 4 - 5 + 1 = -3*
*=> probability of class 1 is 1/(1+e^(-(-3))) = 0.0474*
*-1.5+1.5+3.5+2.5+2.5-1.5+1.5-3.5-4.5+1.5 = 2*
*=> probability of class 1 is now 1/(1+e^(-(2))) = 0.88*
**i.e. we improved the class 1 probability from 5% to 88%**

*This was only with 10 input dimensions. A 224x224 input image has 150,528.*

*(It's significantly easier with more numbers, need smaller nudge for each)*

# *Andrej Karpathy Blog post: Breaking Linear Classifiers on ImageNet*

*Recall CIFAR-10 linear classifiers:*



*ImageNet classifiers:*

comp150dl

*mix in a tiny bit of Goldfish classifier weights*

0.9% bobsled + 100.0% goldfish = **100% Goldfish**

comp150dl **Tufts** UNIVERSITY

1.0% kit fox

8.0% goldfish

comp150dl

# Recurrent Neural Networks

# Recurrent Networks offer a lot of flexibility:

one to one    one to many    many to one    many to many    many to many

**Vanilla Neural Networks**

comp150dl    Tufts

# Recurrent Networks offer a lot of flexibility:

one to one    one to many    many to one    many to many    many to many

e.g. **Image Captioning**
image -> sequence of words

comp150dl

# Recurrent Networks offer a lot of flexibility:



one to one | one to many | many to one | many to many | many to many

e.g. **Sentiment Classification**
sequence of words -> sentiment

# Recurrent Networks offer a lot of flexibility:



one to one | one to many | many to one | many to many | many to many

e.g. **Machine Translation**
seq of words -> seq of words

comp150dl

# Recurrent Networks offer a lot of flexibility:

one to one    one to many    many to one    many to many    many to many

e.g. **Video classification on frame level**

# Recurrent Networks



**Recurrent Neural Networks have loops.**

# RNN - at each time step



An unrolled recurrent neural network.

$$h_t = f_W(h_{t-1}, x_t)$$

new state

some function
with parameters W

old state

Notice: the same
function and the same
set of parameters are
used at every time step.

comp150dl **Tufts**

* figure courtesy Chris Olah

# (Vanilla) Recurrent Neural Network

The state consists of a single *"hidden"* vector **h**:

$$h_t = f_W(h_{t-1}, x_t)$$

$$\downarrow$$

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

$$y_t = W_{hy}h_t$$

comp150dl **Tufts**

**Character-level language model example**

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

comp150dl **Tufts** UNIVERSITY

# Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

# Character-level language model example

Vocabulary:
[h,e,l,o]

Example training sequence:
**"hello"**

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$

# Character-level language model example

Vocabulary: [h,e,l,o]

Example training sequence: **"hello"**

# min-char-rnn.py gist: 112 lines of Python



(https://gist.github.com/karpathy/
d4dee566867f8291f086)

# Data I/O

```
"""
Minimal character-level Vanilla RNN model. Written by Andrej Karpathy (@karpathy)
BSD License
"""
import numpy as np

# data I/O
data = open('input.txt', 'r').read() # should be simple plain text file
chars = list(set(data))
data_size, vocab_size = len(data), len(chars)
print 'data has %d characters, %d unique.' % (data_size, vocab_size)
char_to_ix = { ch:i for i,ch in enumerate(chars) }
ix_to_char = { i:ch for i,ch in enumerate(chars) }
```

n Andrej Karpathy
1n

comp150dl

# Initializations

```
15    # hyperparameters
16    hidden_size = 100 # size of hidden layer of neurons
17    seq_length = 25 # number of steps to unroll the RNN for
18    learning_rate = 1e-1
19
20    # model parameters
21    Wxh = np.random.randn(hidden_size, vocab_size)*0.01 # input to hidden
22    Whh = np.random.randn(hidden_size, hidden_size)*0.01 # hidden to hidden
23    Why = np.random.randn(vocab_size, hidden_size)*0.01 # hidden to output
24    bh = np.zeros((hidden_size, 1)) # hidden bias
25    by = np.zeros((vocab_size, 1)) # output bias
```

recall:



n Andrej Karpathy
1n

comp150dl

# Main loop

```
81   n, p = 0, 0
82   mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
83   mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84   smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85   while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88       hprev = np.zeros((hidden_size,1)) # reset RNN memory
89       p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95       sample_ix = sample(hprev, inputs[0], 200)
96       txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97       print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
106                                  [dWxh, dWhh, dWhy, dbh, dby],
107                                  [mWxh, mWhh, mWhy, mbh, mby]):
108      mem += dparam * dparam
109      param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

n Andrej Karpathy
1n

# Main loop



```
81   n, p = 0, 0
82   mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
83   mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84   smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85   while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88       hprev = np.zeros((hidden_size,1)) # reset RNN memory
89       p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95       sample_ix = sample(hprev, inputs[0], 200)
96       txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97       print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
106                                  [dWxh, dWhh, dWhy, dbh, dby],
107                                  [mWxh, mWhh, mWhy, mbh, mby]):
108      mem += dparam * dparam
109      param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

n Andrej Karpathy
1n

# Main loop

```
81   n, p = 0, 0
82   mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
83   mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84   smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85   while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88       hprev = np.zeros((hidden_size,1)) # reset RNN memory
89       p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95       sample_ix = sample(hprev, inputs[0], 200)
96       txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97       print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
106                                  [dWxh, dWhh, dWhy, dbh, dby],
107                                  [mWxh, mWhh, mWhy, mbh, mby]):
108      mem += dparam * dparam
109      param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

n Andrej Karpathy
1n

# Main loop

```python
81   n, p = 0, 0
82   mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
83   mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84   smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85   while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88       hprev = np.zeros((hidden_size,1)) # reset RNN memory
89       p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95       sample_ix = sample(hprev, inputs[0], 200)
96       txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97       print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
106                                  [dWxh, dWhh, dWhy, dbh, dby],
107                                  [mWxh, mWhh, mWhy, mbh, mby]):
108      mem += dparam * dparam
109      param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

n Andrej Karpathy
1n

# Main loop



```python
81   n, p = 0, 0
82   mWxh, mWhh, mWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
83   mbh, mby = np.zeros_like(bh), np.zeros_like(by) # memory variables for Adagrad
84   smooth_loss = -np.log(1.0/vocab_size)*seq_length # loss at iteration 0
85   while True:
86     # prepare inputs (we're sweeping from left to right in steps seq_length long)
87     if p+seq_length+1 >= len(data) or n == 0:
88       hprev = np.zeros((hidden_size,1)) # reset RNN memory
89       p = 0 # go from start of data
90     inputs = [char_to_ix[ch] for ch in data[p:p+seq_length]]
91     targets = [char_to_ix[ch] for ch in data[p+1:p+seq_length+1]]
92
93     # sample from the model now and then
94     if n % 100 == 0:
95       sample_ix = sample(hprev, inputs[0], 200)
96       txt = ''.join(ix_to_char[ix] for ix in sample_ix)
97       print '----\n %s \n----' % (txt, )
98
99     # forward seq_length characters through the net and fetch gradient
100    loss, dWxh, dWhh, dWhy, dbh, dby, hprev = lossFun(inputs, targets, hprev)
101    smooth_loss = smooth_loss * 0.999 + loss * 0.001
102    if n % 100 == 0: print 'iter %d, loss: %f' % (n, smooth_loss) # print progress
103
104    # perform parameter update with Adagrad
105    for param, dparam, mem in zip([Wxh, Whh, Why, bh, by],
106                                  [dWxh, dWhh, dWhy, dbh, dby],
107                                  [mWxh, mWhh, mWhy, mbh, mby]):
108      mem += dparam * dparam
109      param += -learning_rate * dparam / np.sqrt(mem + 1e-8) # adagrad update
110
111    p += seq_length # move data pointer
112    n += 1 # iteration counter
```

n Andrej Karpathy
1n

## Loss function

- forward pass (compute loss)
- backward pass (compute param gradient)

```python
27    def lossFun(inputs, targets, hprev):
28      """
29      inputs,targets are both list of integers.
30      hprev is Hx1 array of initial hidden state
31      returns the loss, gradients on model parameters, and last hidden state
32      """
33      xs, hs, ys, ps = {}, {}, {}, {}
34      hs[-1] = np.copy(hprev)
35      loss = 0
36      # forward pass
37      for t in xrange(len(inputs)):
38        xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39        xs[t][inputs[t]] = 1
40        hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh) # hidden state
41        ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
42        ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43        loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
44      # backward pass: compute gradients going backwards
45      dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
46      dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47      dhnext = np.zeros_like(hs[0])
48      for t in reversed(xrange(len(inputs))):
49        dy = np.copy(ps[t])
50        dy[targets[t]] -= 1 # backprop into y
51        dWhy += np.dot(dy, hs[t].T)
52        dby += dy
53        dh = np.dot(Why.T, dy) + dhnext # backprop into h
54        dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55        dbh += dhraw
56        dWxh += np.dot(dhraw, xs[t].T)
57        dWhh += np.dot(dhraw, hs[t-1].T)
58        dhnext = np.dot(Whh.T, dhraw)
59      for dparam in [dWxh, dWhh, dWhy, dbh, dby]:
60        np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61      return loss, dWxh, dWhh, dWhy, dbh, dby, hs[len(inputs)-1]
```

n Andrej Karpathy
1n

cor

```
27   def lossFun(inputs, targets, hprev):
28     """
29     inputs,targets are both list of integers.
30     hprev is Hx1 array of initial hidden state
31     returns the loss, gradients on model parameters, and last hidden state
32     """
33     xs, hs, ys, ps = {}, {}, {}, {}
34     hs[-1] = np.copy(hprev)
35     loss = 0
36     # forward pass
37     for t in xrange(len(inputs)):
38       xs[t] = np.zeros((vocab_size,1)) # encode in 1-of-k representation
39       xs[t][inputs[t]] = 1
40       hs[t] = np.tanh(np.dot(Wxh, xs[t]) + np.dot(Whh, hs[t-1]) + bh) # hidden state
41       ys[t] = np.dot(Why, hs[t]) + by # unnormalized log probabilities for next chars
42       ps[t] = np.exp(ys[t]) / np.sum(np.exp(ys[t])) # probabilities for next chars
43       loss += -np.log(ps[t][targets[t],0]) # softmax (cross-entropy loss)
```

$$h_t = \tanh(W_{hh}h_{t-1} + W_{xh}x_t)$$
$$y_t = W_{hy}h_t$$

Softmax classifier

Andrej Karpathy
1n

comp150dl  Tufts

```
44    # backward pass: compute gradients going backwards
45    dWxh, dWhh, dWhy = np.zeros_like(Wxh), np.zeros_like(Whh), np.zeros_like(Why)
46    dbh, dby = np.zeros_like(bh), np.zeros_like(by)
47    dhnext = np.zeros_like(hs[0])
48    for t in reversed(xrange(len(inputs))):
49      dy = np.copy(ps[t])
50      dy[targets[t]] -= 1 # backprop into y
51      dWhy += np.dot(dy, hs[t].T)
52      dby += dy
53      dh = np.dot(Why.T, dy) + dhnext # backprop into h
54      dhraw = (1 - hs[t] * hs[t]) * dh # backprop through tanh nonlinearity
55      dbh += dhraw
56      dWxh += np.dot(dhraw, xs[t].T)
57      dWhh += np.dot(dhraw, hs[t-1].T)
58      dhnext = np.dot(Whh.T, dhraw)
59    for dparam in [dWxh, dWhh, dWhy, dbh, dby]:
60      np.clip(dparam, -5, 5, out=dparam) # clip to mitigate exploding gradients
61    return loss, dWxh, dWhh, dWhy, dbh, dby, hs[len(inputs)-1]
```

recall:



n Andrej Karpathy
1n

comp150dl    Tufts

```python
def sample(h, seed_ix, n):
    """
    sample a sequence of integers from the model
    h is memory state, seed_ix is seed letter for first time step
    """
    x = np.zeros((vocab_size, 1))
    x[seed_ix] = 1
    ixes = []
    for t in xrange(n):
        h = np.tanh(np.dot(Wxh, x) + np.dot(Whh, h) + bh)
        y = np.dot(Why, h) + by
        p = np.exp(y) / np.sum(np.exp(y))
        ix = np.random.choice(range(vocab_size), p=p.ravel())
        x = np.zeros((vocab_size, 1))
        x[ix] = 1
        ixes.append(ix)
    return ixes
```

n Andrej Karpathy
1n

comp150dl   Tufts

## Sonnet 116 – Let me not …

*by William Shakespeare*

Let me not to the marriage of true minds
    Admit impediments. Love is not love
Which alters when it alteration finds,
    Or bends with the remover to remove:
O no! it is an ever-fixed mark
    That looks on tempests and is never shaken;
It is the star to every wandering bark,
    Whose worth's unknown, although his height be taken.
Love's not Time's fool, though rosy lips and cheeks
    Within his bending sickle's compass come:
Love alters not with his brief hours and weeks,
    But bears it out even to the edge of doom.
If this be error and upon me proved,
    I never writ, nor no man ever loved.

comp150dl  Tufts

at first:

```
tyntd-iafhatawiaoihrdemot  lytdws  e ,tfti, astai f ogoh eoase rrranbyne 'nhthnee e
plia tklrgd t o idoe ns,smtt   h ne etie h,hregtrs nigtike,aoaenns lng
```

↓ train more

```
"Tmont thithey" fomesscerliund
Keushey. Thom here
sheulke, anmerenith ol sivh I lalterthend Bleipile shuwy fil on aseterlome
coaniogennc Phe lism thond hon at. MeiDimorotion in ther thize."
```

↓ train more

```
Aftair fall unsuch that the hall for Prince Velzonski's that me of
her hearly, and behs to so arwage fiving were to it beloge, pavu say falling misfort
how, and Gogition is so overelical and ofter.
```

↓ train more

```
"Why do what that day," replied Natasha, and wishing to himself the fact the
princess, Princess Mary was easier, fed in had oftened him.
Pierre aking his soul came to the packs and drove up his father-in-law women.
```

```
PANDARUS:
Alas, I think he shall be come approached and the day
When little srain would be attain'd into being never fed,
And who is but a chain and subjects of his death,
I should not sleep.

Second Senator:
They are away this miseries, produced upon my soul,
Breaking and strongly should be buried, when I perish
The earth and thoughts of many states.

DUKE VINCENTIO:
Well, your wit is in the care of side and that.

Second Lord:
They would be ruled after this chamber, and
my fair nues begun out of the fact, to be conveyed,
Whose noble souls I'll have the heart of the wars.

Clown:
Come, sir, I will make did behold your worship.

VIOLA:
I'll drink it.
```

```
VIOLA:
Why, Salisbury must find his flesh and thought
That which I am not aps, not a man and in fire,
To show the reining of the raven and the wars
To grace my hand reproach within, and not a fair are hand,
That Caesar and my goodly father's world;
When I was heaven of presence and our fleets,
We spare with hours, but cut thy council I am great,
Murdered and by thy master's ready there
My power to give thee but so much as hell:
Some service in the noble bondman here,
Would show him to her wine.

KING LEAR:
O, if you were a feeble sight, the courtesy of your law,
Your sight and several breath, will wear the gods
With his heads, and my hands are wonder'd at the deeds,
So drop upon your lordship's head, and your opinion
Shall be against your honour.
```

# open source textbook on algebraic geometry



## The Stacks Project

home    about    tags explained    tag lookup    browse    search    bibliography    recent comments    blog    add slogans

### Browse chapters

| Part | Chapter | online | TeX source | view pdf |
|------|---------|--------|-----------|----------|
| Preliminaries | | | | |
| | 1. Introduction | online | tex | pdf |
| | 2. Conventions | online | tex | pdf |
| | 3. Set Theory | online | tex | pdf |
| | 4. Categories | online | tex | pdf |
| | 5. Topology | online | tex | pdf |
| | 6. Sheaves on Spaces | online | tex | pdf |
| | 7. Sites and Sheaves | online | tex | pdf |
| | 8. Stacks | online | tex | pdf |
| | 9. Fields | online | tex | pdf |
| | 10. Commutative Algebra | online | tex | pdf |

### Parts

1. Preliminaries
2. Schemes
3. Topics in Scheme Theory
4. Algebraic Spaces
5. Topics in Geometry
6. Deformation Theory
7. Algebraic Stacks
8. Miscellany

### Statistics

The Stacks project now consists of

- 455910 lines of code
- 14221 tags (56 inactive tags)
- 2366 sections

## Latex source

comp150dl

**Tufts**

For $\bigoplus_{n=1,\ldots,m}$ where $\mathcal{L}_{m_*} = 0$, hence we can find a closed subset $\mathcal{H}$ in $\mathcal{H}$ and any sets $\mathcal{F}$ on $X$, $U$ is a closed immersion of $S$, then $U \to T$ is a separated algebraic space.

*Proof.* Proof of (1). It also start we get

$$S = \mathrm{Spec}(R) = U \times_X U \times_X U$$

and the comparicoly in the fibre product covering we have to prove the lemma generated by $\coprod Z \times_U U \to V$. Consider the maps $M$ along the set of points $Sch_{fppf}$ and $U \to U$ is the fibre category of $S$ in $U$ in Section, ?? and the fact that any $U$ affine, see Morphisms, Lemma ??. Hence we obtain a scheme $S$ and any open subset $W \subset U$ in $Sh(G)$ such that $\mathrm{Spec}(R') \to S$ is smooth or an

$$U = \bigcup U_i \times_{S_i} U_i$$

which has a nonzero morphism we may assume that $f_i$ is of finite presentation over $S$. We claim that $\mathcal{O}_{X,x}$ is a scheme where $x, x', s'' \in S'$ such that $\mathcal{O}_{X,x'} \to \mathcal{O}'_{X',x'}$ is separated. By Algebra, Lemma ?? we can define a map of complexes $\mathrm{GL}_{S'}(x'/S'')$ and we win.

To prove study we see that $\mathcal{F}|_U$ is a covering of $\mathcal{X}^o$, and $\mathcal{T}_i$ is an object of $\mathcal{F}_{X/S}$ for $i > 0$ and $\mathcal{F}_p$ exists and let $\mathcal{F}_i$ be a presheaf of $\mathcal{O}_X$-modules on $\mathcal{C}$ as a $\mathcal{F}$-module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\widetilde{M}^\bullet = \mathcal{I}^\bullet \otimes_{\mathrm{Spec}(k)} \mathcal{O}_{S,s} - i_X^{-1}\mathcal{F})$$

is a unique morphism of algebraic stacks. Note that

$$\mathrm{Arrows} = (Sch/S)^{opp}_{fppf}(Sch/S)_{fppf}$$

and

$$V = \Gamma(S, \mathcal{O}) \longmapsto (U, \mathrm{Spec}(A))$$

is an open subset of $X$. Thus $U$ is affine. This is a continuous map of $X$ is the inverse, the groupoid scheme $S$.

*Proof.* See discussion of sheaves of sets. $\square$

The result for prove any open covering follows from the less of Example ??. It may replace $S$ by $X_{spaces,\text{é}tale}$ which gives an open subspace of $X$ and $T$ equal to $S_{Zar}$, see Descent, Lemma ??. Namely, by Lemma ?? we see that $R$ is geometrically regular over $S$.

**Lemma 0.1.** *Assume (3) and (3) by the construction in the description.*

*Suppose $X = \lim |X|$ (by the formal open covering $X$ and a single map $\underline{\mathrm{Proj}}_X(A) = \mathrm{Spec}(B)$ over $U$ compatible with the complex*

$$\mathrm{Set}(A) = \Gamma(X, \mathcal{O}_{X,\mathcal{O}_X}).$$

*When in this case of to show that $Q \to \mathcal{C}_{Z/X}$ is stable under the following result in the second conditions of (1), and (3). This finishes the proof. By Definition ?? (without element is when the closed subschemes are catenary. If $T$ is surjective we may assume that $T$ is connected with residue fields of $S$. Moreover there exists a closed subspace $Z \subset X$ of $X$ where $U$ in $X'$ is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem*

   (1) *$f$ is locally of finite type. Since $S = \mathrm{Spec}(R)$ and $Y = \mathrm{Spec}(R)$.*

*Proof.* This is form all sheaves of sheaves on $X$. But given a scheme $U$ and a surjective étale morphism $U \to X$. Let $U \cap U = \coprod_{i=1,\ldots,n} U_i$ be the scheme $X$ over $S$ at the schemes $X_i \to X$ and $U = \lim_i X_i$. $\square$

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{x_0} = \mathcal{F}_{x_0} = \mathcal{F}_{X,\ldots,0}$.

**Lemma 0.2.** *Let $X$ be a locally Noetherian scheme over $S$, $E = \mathcal{F}_{X/S}$. Set $\mathcal{I} = \mathcal{J}_1 \subset \mathcal{I}'_n$. Since $\mathcal{I}^n \subset \mathcal{I}^n$ are nonzero over $i_0 \leq \mathfrak{p}$ is a subset of $\mathcal{J}_{n,0} \circ \overline{A}_2$ works.*

**Lemma 0.3.** *In Situation ??. Hence we may assume $\mathfrak{q}' = 0$.*

*Proof.* We will use the property we see that $\mathfrak{p}$ is the mext functor (??). On the other hand, by Lemma ?? we see that

$$D(\mathcal{O}_{X'}) = \mathcal{O}_X(D)$$

where $K$ is an $F$-algebra where $\delta_{n+1}$ is a scheme over $S$. $\square$

comp150dl    Tufts

*Proof.* Omitted. □

**Lemma 0.1.** *Let $\mathcal{C}$ be a set of the construction.*

*Let $\mathcal{C}$ be a gerber covering. Let $\mathcal{F}$ be a quasi-coherent sheaves of $\mathcal{O}$-modules. We have to show that*

$$\mathcal{O}_{\mathcal{O}_X} = \mathcal{O}_X(\mathcal{L})$$

.

*Proof.* This is an algebraic space with the composition of sheaves $\mathcal{F}$ on $X_{\text{étale}}$ we have

$$\mathcal{O}_X(\mathcal{F}) = \{morph_1 \times_{\mathcal{O}_X} (\mathcal{G}, \mathcal{F})\}$$

where $\mathcal{G}$ defines an isomorphism $\mathcal{F} \to \mathcal{F}$ of $\mathcal{O}$-modules. □

**Lemma 0.2.** *This is an integer $\mathcal{Z}$ is injective.*

*Proof.* See Spaces, Lemma ??. □

**Lemma 0.3.** *Let $S$ be a scheme. Let $X$ be a scheme and $X$ is an affine open covering. Let $\mathcal{U} \subset X$ be a canonical and locally of finite type. Let $X$ be a scheme. Let $X$ be a scheme which is equal to the formal complex.*

*The following to the construction of the lemma follows.*

*Let $X$ be a scheme. Let $X$ be a scheme covering. Let*

$$b : X \to Y' \to Y \to Y \to Y' \times_X Y \to X.$$

*be a morphism of algebraic spaces over $S$ and $Y$.*

*Proof.* Let $X$ be a nonzero scheme of $X$. Let $X$ be an algebraic space. Let $\mathcal{F}$ be a quasi-coherent sheaf of $\mathcal{O}_X$-modules. The following are equivalent

(1) $\mathcal{F}$ is an algebraic space over $S$.
(2) If $X$ is an affine open covering.

Consider a common structure on $X$ and $X$ the functor $\mathcal{O}_X(U)$ which is locally of finite type. □

This since $\mathcal{F} \in \mathcal{F}$ and $x \in \mathcal{G}$ the diagram



is a limit. Then $\mathcal{G}$ is a finite type and assume $S$ is a flat and $\mathcal{F}$ and $\mathcal{G}$ is a finite type $f_*$. This is of finite type diagrams, and

- the composition of $\mathcal{G}$ is a regular sequence,
- $\mathcal{O}_{X'}$ is a sheaf of rings.

□

*Proof.* We have see that $X = \operatorname{Spec}(R)$ and $\mathcal{F}$ is a finite type representable by algebraic space. The property $\mathcal{F}$ is a finite morphism of algebraic stacks. Then the cohomology of $X$ is an open neighbourhood of $U$. □

*Proof.* This is clear that $\mathcal{G}$ is a finite presentation, see Lemmas ??. A reduced above we conclude that $U$ is an open covering of $\mathcal{C}$. The functor $\mathcal{F}$ is a "field

$$\mathcal{O}_{X,x} \longrightarrow \mathcal{F}_{\overline{x}} \quad -1(\mathcal{O}_{X_{\text{étale}}}) \longrightarrow \mathcal{O}_{X,x}^{-1}\mathcal{O}_{X_\lambda}(\mathcal{O}_{X_\lambda}^\vee)$$

is an isomorphism of covering of $\mathcal{O}_{X_i}$. If $\mathcal{F}$ is the unique element of $\mathcal{F}$ such that $X$ is an isomorphism.

The property $\mathcal{F}$ is a disjoint union of Proposition ?? and we can filtered set of presentations of a scheme $\mathcal{O}_X$-algebra with $\mathcal{F}$ are opens of finite type over $S$. If $\mathcal{F}$ is a scheme theoretic image points. □

If $\mathcal{F}$ is a finite direct sum $\mathcal{O}_{X_\lambda}$ is a closed immersion, see Lemma ??. This is a sequence of $\mathcal{F}$ is a similar morphism.

comp150dl **Tufts** UNIVERSITY

# Generated C code

```c
static void do_command(struct seq_file *m, void *v)
{
  int column = 32 << (cmd[2] & 0x80);
  if (state)
    cmd = (int)(int_state ^ (in_8(&ch->ch_flags) & Cmd) ? 2 : 1);
  else
    seq = 1;
  for (i = 0; i < 16; i++) {
    if (k & (1 << 1))
      pipe = (in_use & UMXTHREAD_UNCCA) +
        ((count & 0x00000000ffffff8) & 0x000000f) << 8;
    if (count == 0)
      sub(pid, ppc_md.kexec_handle, 0x20000000);
    pipe_set_bytes(i, 0);
  }
  /* Free our user pages pointer to place camera if all dash */
  subsystem_info = &of_changes[PAGE_SIZE];
  rek_controls(offset, idx, &soffset);
  /* Now we want to deliberately put it to device */
  control_check_polarity(&context, val, 0);
  for (i = 0; i < COUNTER; i++)
    seq_puts(s, "policy ");
}
```

```c
/*
 *  Copyright (c) 2006-2010, Intel Mobile Communications.  All rights reserved.
 *
 *   This program is free software; you can redistribute it and/or modify it
 * under the terms of the GNU General Public License version 2 as published by
 * the Free Software Foundation.
 *
 *       This program is distributed in the hope that it will be useful,
 * but WITHOUT ANY WARRANTY; without even the implied warranty of
 *   MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE.  See the
 *
 *  GNU General Public License for more details.
 *
 *   You should have received a copy of the GNU General Public License
 *    along with this program; if not, write to the Free Software Foundation,
 *  Inc., 675 Mass Ave, Cambridge, MA 02139, USA.
 */

#include <linux/kexec.h>
#include <linux/errno.h>
#include <linux/io.h>
#include <linux/platform_device.h>
#include <linux/multi.h>
#include <linux/ckevent.h>

#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>
```

comp150dl       **Tufts**       *84*

```c
#include <asm/io.h>
#include <asm/prom.h>
#include <asm/e820.h>
#include <asm/system_info.h>
#include <asm/setew.h>
#include <asm/pgproto.h>

#define REG_PG     vesa_slot_addr_pack
#define PFM_NOCOMP  AFSR(0, load)
#define STACK_DDR(type)      (func)

#define SWAP_ALLOCATE(nr)       (e)
#define emulate_sigs()  arch_get_unaligned_child()
#define access_rw(TST)  asm volatile("movd %%esp, %0, %3" : : "r" (0));   \
  if (__type & DO_READ)

static void stat_PC_SEC __read_mostly offsetof(struct seq_argsqueue, \
        pC>[1]);

static void
os_prefix(unsigned long sys)
{
#ifdef CONFIG_PREEMPT
  PUT_PARAM_RAID(2, sel) = get_state_state();
  set_pid_sum((unsigned long)state, current_state_str(),
        (unsigned long)-1->lr_full; low;
}
```

# Recommended Reading:
## *Visualizing and Understanding Recurrent Networks*



*[Visualizing and Understanding Recurrent Networks, Andrej Karpathy\*, Justin Johnson\*, Li Fei-Fei]*

comp150dl

# Image Captioning



*Explain Images with Multimodal Recurrent Neural Networks*, Mao et al.
*Deep Visual-Semantic Alignments for Generating Image Descriptions*, Karpathy and Fei-Fei
*Show and Tell: A Neural Image Caption Generator*, Vinyals et al.
*Long-term Recurrent Convolutional Networks for Visual Recognition and Description*, Donahue et al.
*Learning a Recurrent Visual Representation for Image Caption Generation*, Chen and Zitnick

comp150dl

**Tufts**
UNIVERSITY

**Recurrent Neural Network**

**Convolutional Neural Network**

test image

comp150dl

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096
FC-1000
softmax

test image

comp150dl  **Tufts**
UNIVERSITY

| image |

| conv-64 |
| conv-64 |
| maxpool |

| conv-128 |
| conv-128 |
| maxpool |

| conv-256 |
| conv-256 |
| maxpool |

| conv-512 |
| conv-512 |
| maxpool |

| conv-512 |
| conv-512 |
| maxpool |

| FC-4096 |
| FC-4096 |
| FC-1000 |
| softmax |

X

test image

comp150dl   Tufts UNIVERSITY

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

x0
<STA
RT>

<START>

comp150dl

Tufts
UNIVERSITY

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

v

test image

**Wih**

y0

h0

x0
<STA
RT>

<START>

**before:**

h = tanh(Wxh * x + Whh * h)

**now:**

h = tanh(Wxh * x + Whh * h **+ Wih * v**)

comp150dl

Tufts
UNIVERSITY

test image

sample!

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

y0

h0

x0
<START>

straw

<START>

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

y0    y1

h0 → h1

x0
<START>    straw

<START>    comp150dl

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

sample!

y0    y1

h0 → h1

x0
<START>

straw

hat

<START>    comp150dl

test image

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

conv-512
conv-512
maxpool

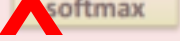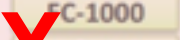conv-512
conv-512
maxpool

FC-4096
FC-4096

y0    y1    y2

h0 → h1 → h2

x0
<START>    straw    hat

<START>    comp150dl

image

conv-64
conv-64
maxpool

conv-128
conv-128
maxpool

conv-256
conv-256
maxpool

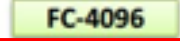conv-512
conv-512
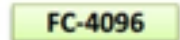maxpool

conv-512
conv-512
maxpool

FC-4096
FC-4096

test image

y0    y1    y2

sample
<END> token
=> finish.

h0 → h1 → h2

x0
<START>    straw    hat

<START>

comp150dl    Tufts

# Image Sentence Datasets

a man riding a bike on a dirt path through a forest.
bicyclist raises his fist as he rides on desert dirt trail.
this dirt bike rider is smiling and raising his fist in triumph.
a man riding a bicycle while pumping his fist in the air.
a mountain biker pumps his fist in celebration.



Microsoft COCO
*[Tsung-Yi Lin et al. 2014]*
mscoco.org

currently:
~120K images
~5 sentences each

"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

comp150dl

"man in black shirt is playing guitar."

"construction worker in orange safety vest is working on road."

"two young girls are playing with lego toy."

"boy is doing backflip on wakeboard."

"a young boy is holding a baseball bat."

"a cat is sitting on a couch with a remote control."

"a woman holding a teddy bear in front of a mirror."

"a horse is standing in the middle of a road."

# Preview of fancier architectures

RNN attends spatially to different parts of images while generating each word of the sentence:



*Show Attend and Tell, Xu et al., 2015*

comp150dl

# ION: INSIDE-OUTSIDE NET



conv1 conv2  conv3  conv4  conv5  conv7  4-dir RNN  4-dir RNN  context features

ROI Pooling

L2 normalize

concat

scale

1x1 conv

fc  fc  fc  softmax

fc  bbox

For each ROI

* slide courtesy Sean Bell

Base ConvNet: VGG16 [Simonyan 2014]

# Limitations of RNNs



"I grew up in France… I speak fluent **French**."

# Long Short Term Memory Networks



The repeating module in a standard RNN contains a single layer.

The repeating module in an LSTM contains four interacting layers.

* figures courtesy Chris Olah

# RNN:

$$h_t^l = \tanh W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$h \in \mathbb{R}^n. \qquad W^l \ [n \times 2n]$$

---

# LSTM:

$$W^l \ [4n \times 2n]$$

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$



depth

time

comp150dl **Tufts**

# LSTM: Cell State
## long running memory of the network

# LSTM: Forget Gate *f*

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \textbf{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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

$$f_t = \sigma\left(W_f \cdot [h_{t-1}, x_t] + b_f\right)$$



* figures courtesy Chris Olah

Tufts UNIVERSITY

# LSTM: Ignore Gate *i*



$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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

$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] \ + \ b_i\right)$$
$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \ + \ b_C)$$

* figures courtesy Chris Olah

Tufts
UNIVERSITY

# LSTM: Block Gate *g*
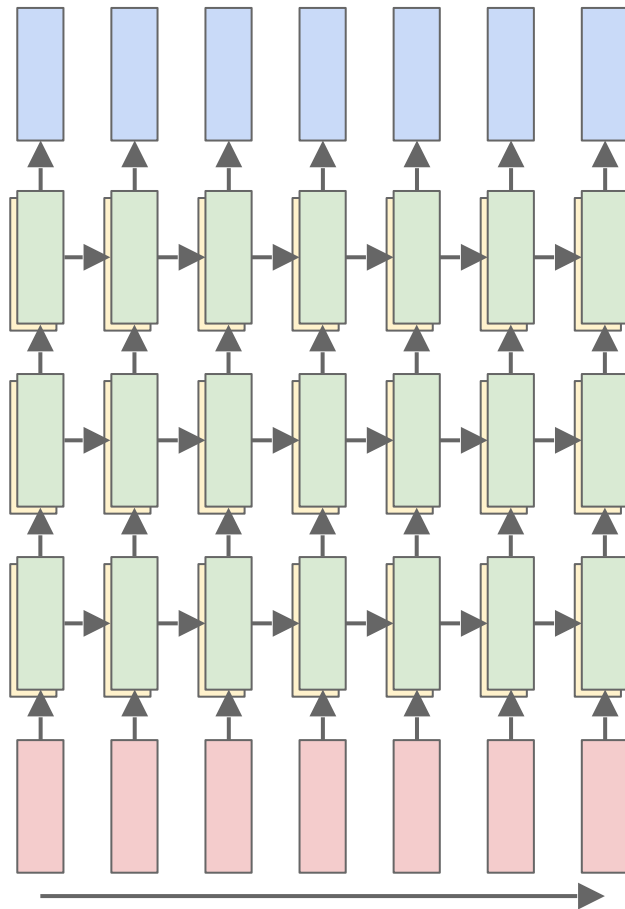
$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \text{sigm} \\ \tanh \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

$$c_t^l = f \odot c_{t-1}^l + i \odot g$$

$$h_t^l = o \odot \tanh(c_t^l)$$
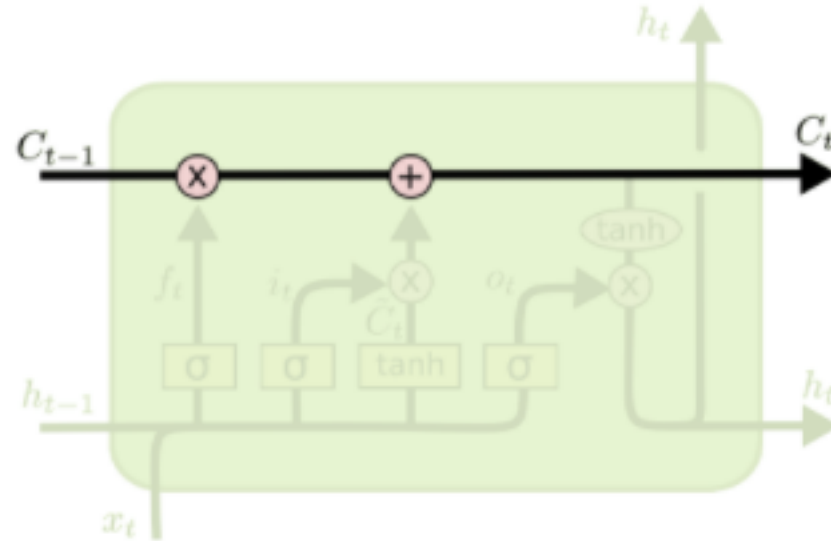


$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$
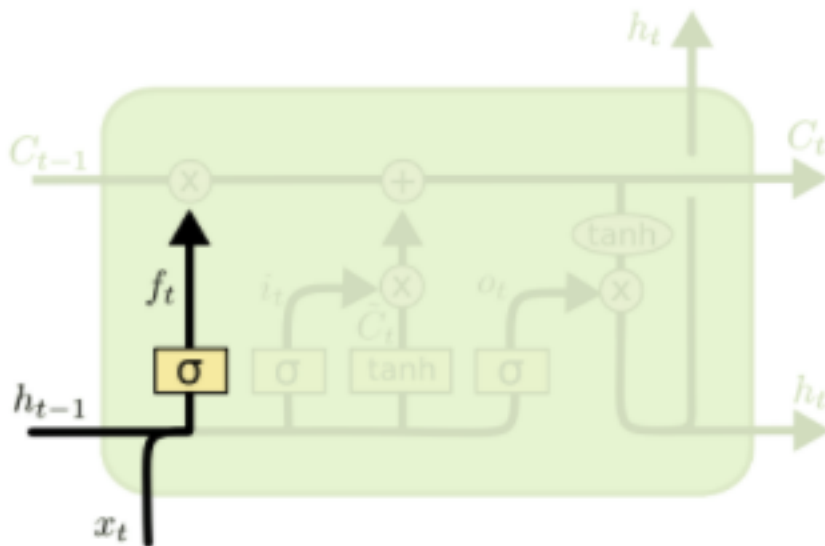
comp150dl  **Tufts** UNIVERSITY

# LSTM: Output Gate *o*

$$\begin{pmatrix} i \\ f \\ o \\ g \end{pmatrix} = \begin{pmatrix} \text{sigm} \\ \text{sigm} \\ \textbf{sigm} \\ \text{tanh} \end{pmatrix} W^l \begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$$

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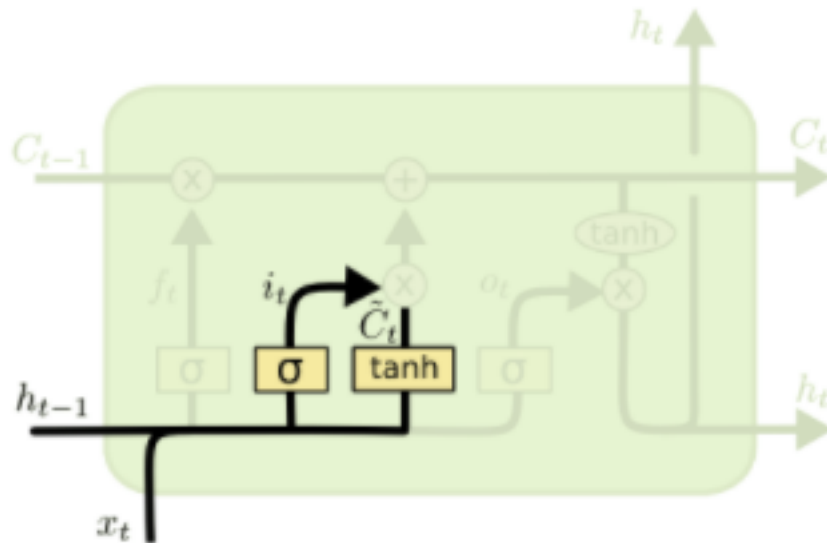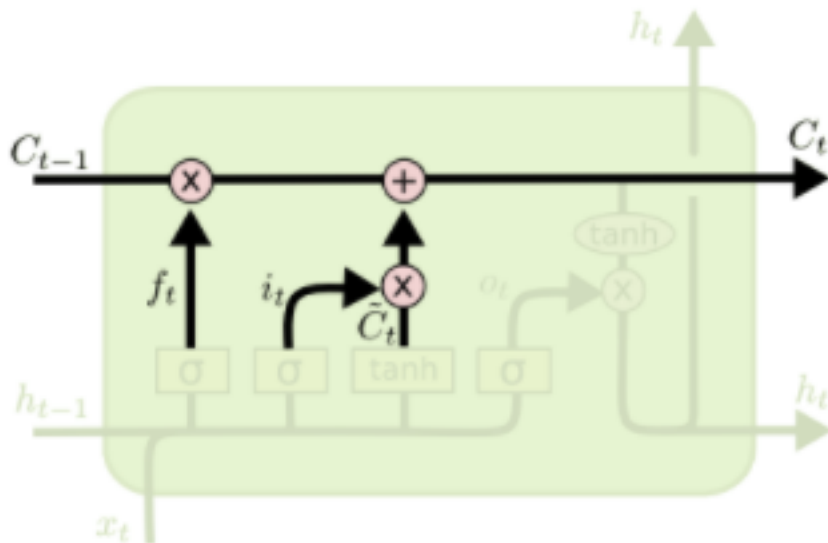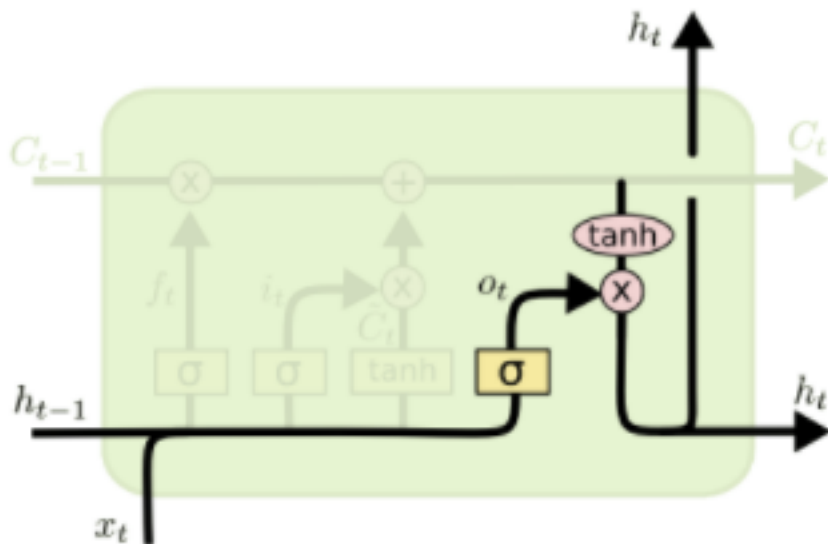$$h_t^l = o \odot \tanh(c_t^l)$$



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

* figures courtesy Chris Olah

comp150dl

# Summary

-   RNNs allow a lot of flexibility in architecture design
-   Vanilla RNNs are simple but don't work very well
-   Common to use LSTM: their additive interactions improve gradient flow
-   Backward flow of gradients in RNN can explode or vanish. Exploding is controlled with gradient clipping. Vanishing is controlled with additive interactions (LSTM)
-   Additional resource for RNNs and LSTMs for Deep NLP: [cs224d.stanford.edu](cs224d.stanford.edu)

comp150dl