

# Deep Neural Networks Basics

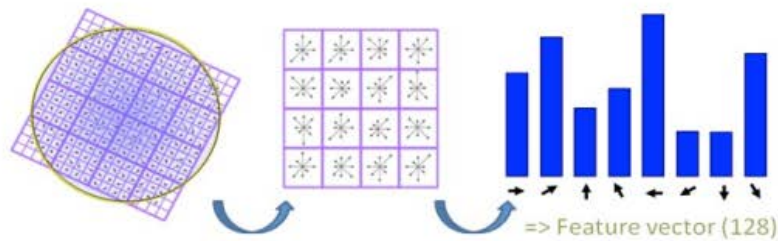
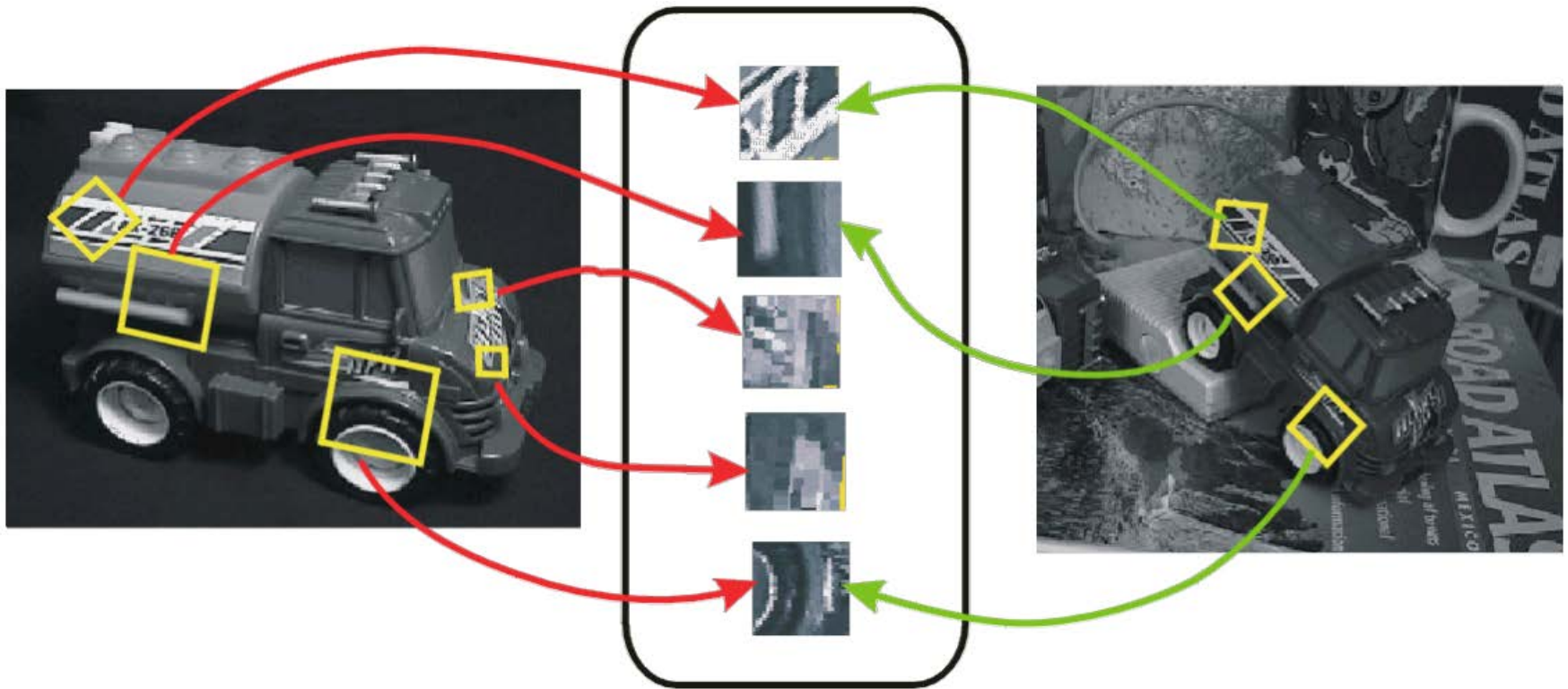
For ECS 289G  
Presented by Fanyi Xiao

# Computer Vision in the Pre-DNN Era



Face Detection, Viola & Jones, 2001

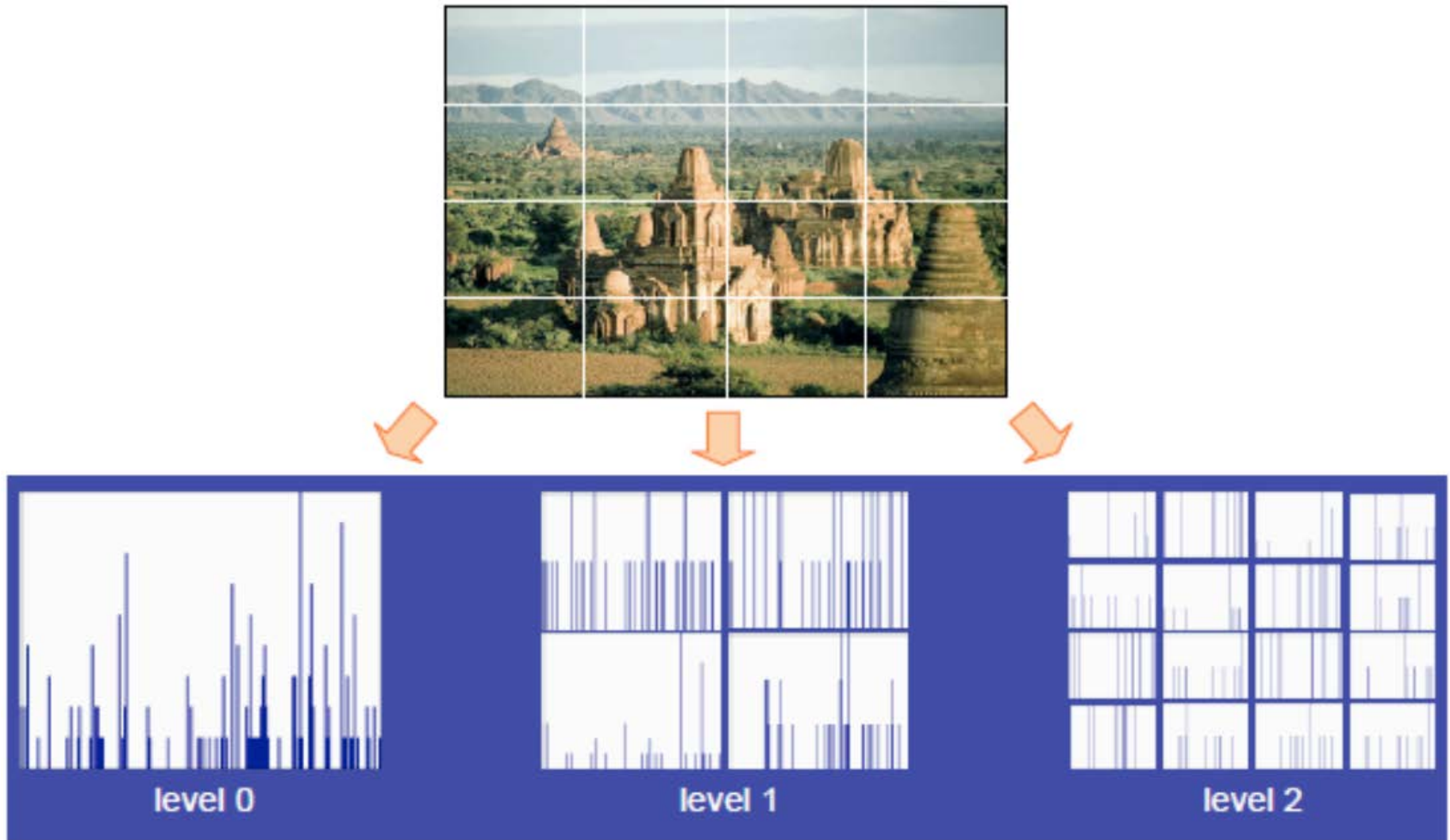
# Computer Vision in the Pre-DNN Era



"SIFT" & Object Recognition, David Lowe, 1999



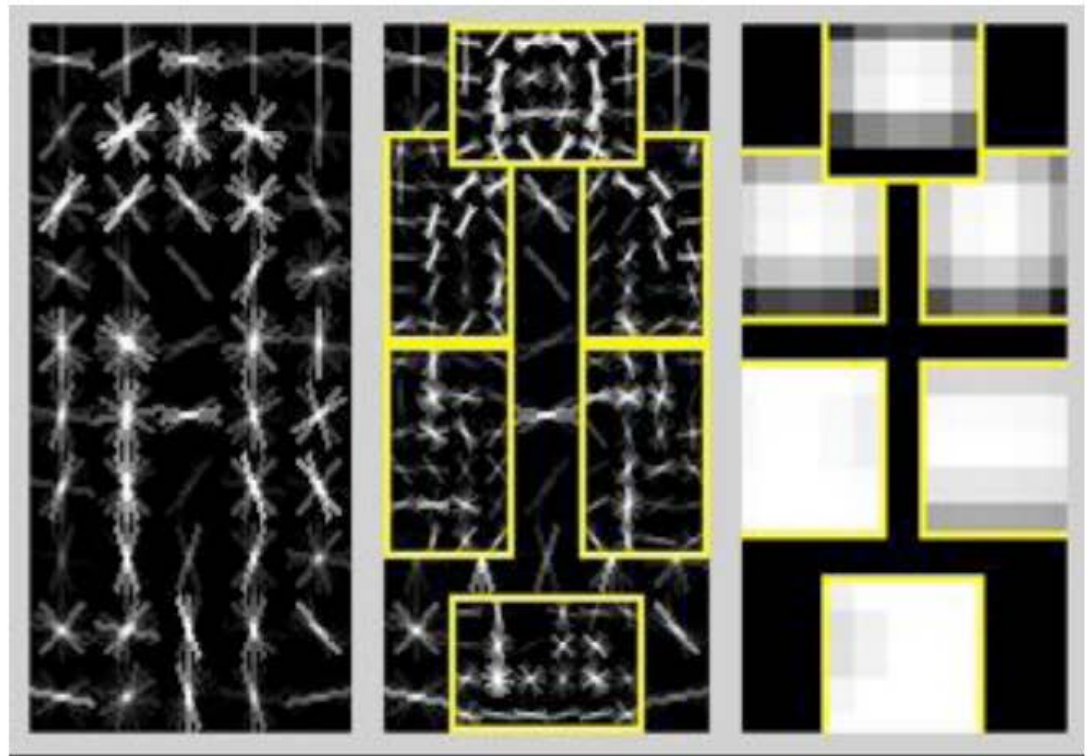
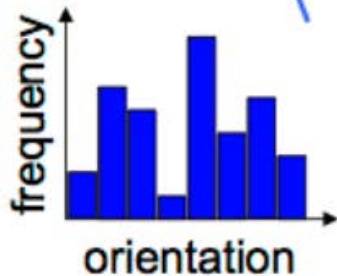
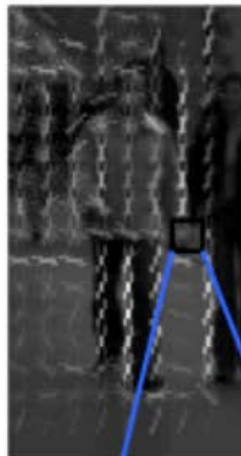
# Computer Vision in the Pre-DNN Era



Spatial Pyramid Matching, Lazebnik, Schmid & Ponce, 2006

Most materials taken from Andrej Karpathy/Richard Socher/Nando de Freitas/caffe CVPR tutorial

# Computer Vision in the Pre-DNN Era



Histogram of Gradients (HoG)  
Dalal & Triggs, 2005

Deformable Part Model  
Felzenswalb, McAllester, Ramanan, 2009

# Emergence of DNNs in Vision



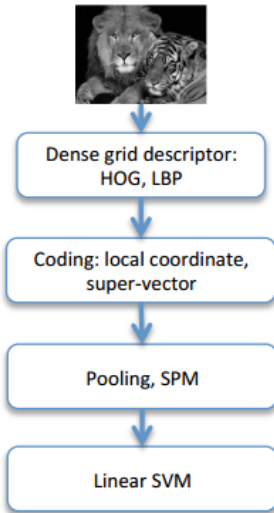
**IMAGENET**  
22K categories and

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plant
- Food
- Material

## IMAGENET Large Scale Visual Recognition Challenge

### Year 2010

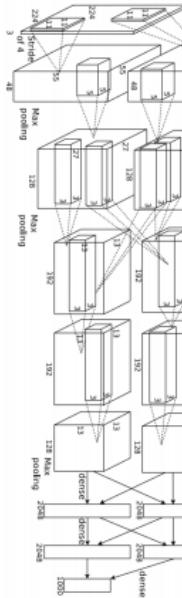
#### NEC-UIUC



[Lin CVPR 2011]

### Year 2012

#### SuperVision



[Krizhevsky NIPS 2012]

### Year 2014

#### GoogLeNet



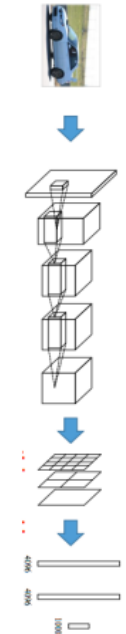
[Szegedy arxiv 2014]

#### VGG



[Simonyan arxiv 2014]

#### MSRA



[He arxiv 2014]



# Neural Networks

## Image Classification

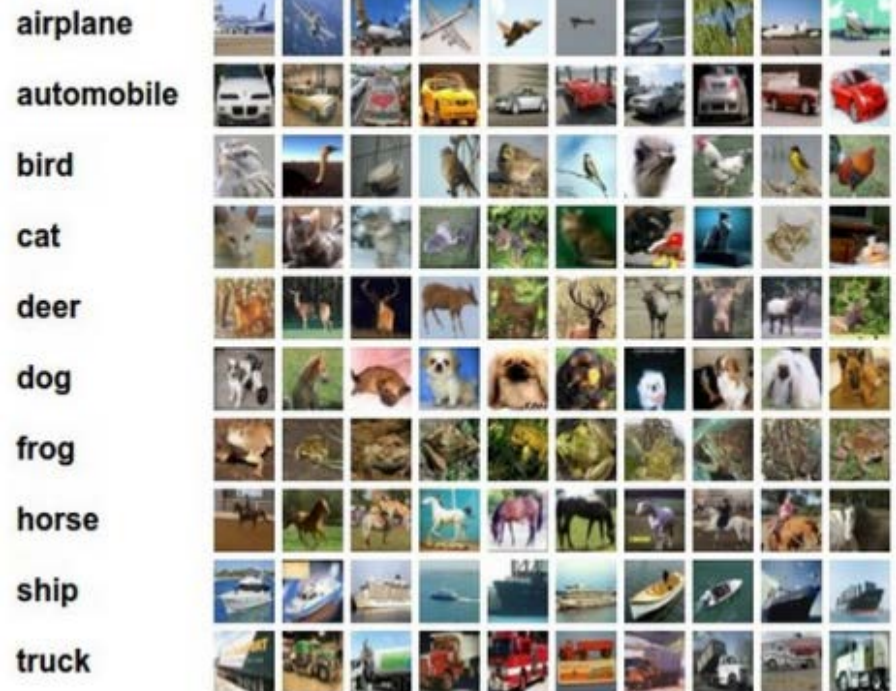


assume given set of discrete labels  
{dog, cat, truck, plane, ...}

→ cat

Learn visual features  
"end-to-end"

## Data-driven approach



# Neural Networks

## Compositional Models

Learned End-to-End

## Hierarchy of Representations

- vision: pixel, motif, part, object
- text: character, word, clause, sentence
- speech: audio, band, phone, word

concrete



abstract

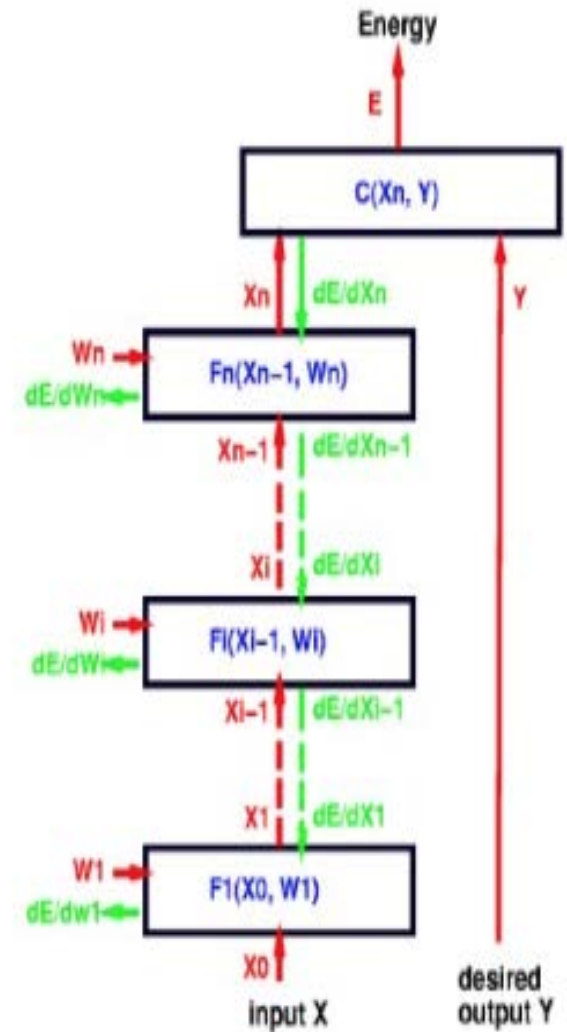


figure credit Yann LeCun, ICML '13 tutorial



# Neural Networks

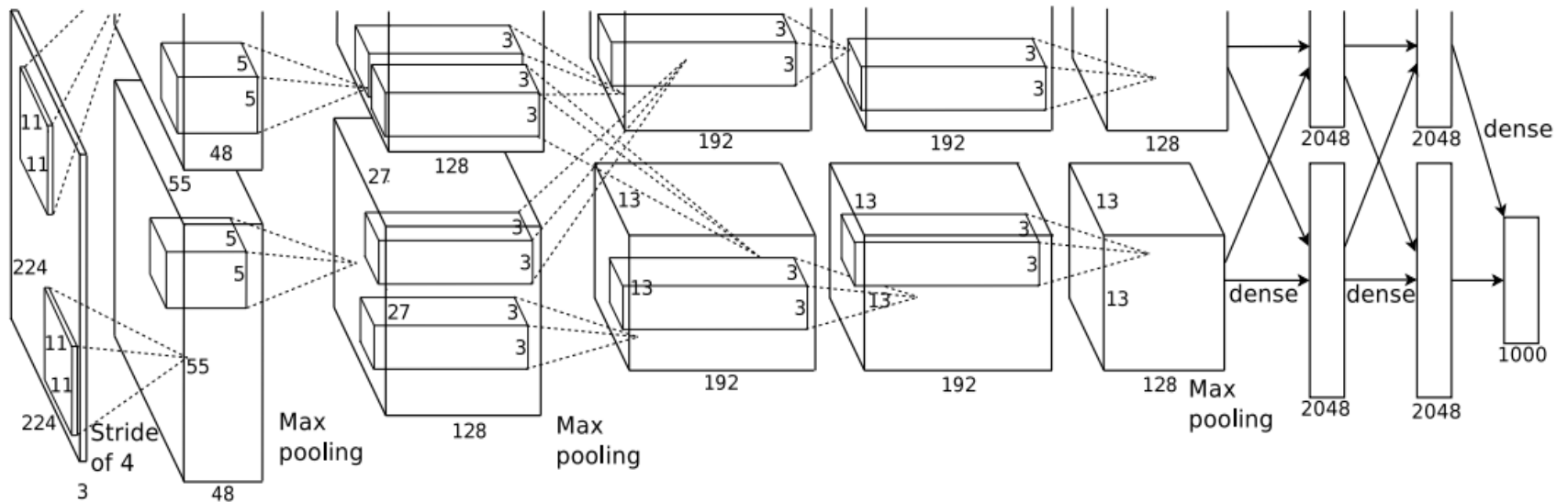
Three key ingredients for training an NN:

1. Score function
2. Loss function
3. Optimization

# Neural Networks

Three key ingredients for training an NN:

1. Score function:  $y=f(x,W)$



$x$  -- 224\*224\*3 image patch  
 $y$  -- 1000d vector

# Neural Networks

Three key ingredients for training an NN:

2. Loss function: for example max-margin loss and cross-entropy loss

$$L_i = \sum_{j \neq y_i} \max(0, f(x_i, W)_j - f(x_i, W)_{y_i} + \Delta)$$

$$L_i = -\log \left( \frac{e^{f_{y_i}}}{\sum_j e^{f_j}} \right)$$



# Neural Networks

Three key ingredients for training an NN:

3. Optimization: simple gradient descent!

$$\frac{df(x)}{dx}$$



# Neural Networks

Three key ingredients for training an NN:

3. Optimization: in practice, *stochastic (mini-batch) gradient descent!*

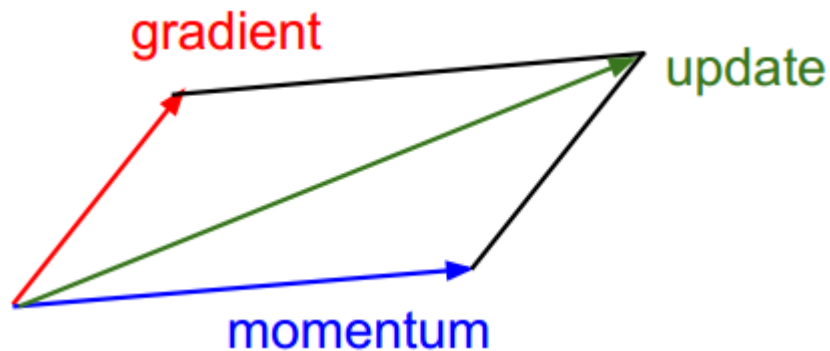
```
# Vanilla Minibatch Gradient Descent

while True:
    data_batch = sample_training_data(data, 256) # sample 256 examples
    weights_grad = evaluate_gradient(loss_fun, data_batch, weights)
    weights += - step_size * weights_grad # perform parameter update
```

# Neural Networks

Three key ingredients for training an NN:

3. Optimization: in practice, stochastic (mini-batch) gradient descent + *momentum*! (Many other optimization methods like adagrad/rmsprop)



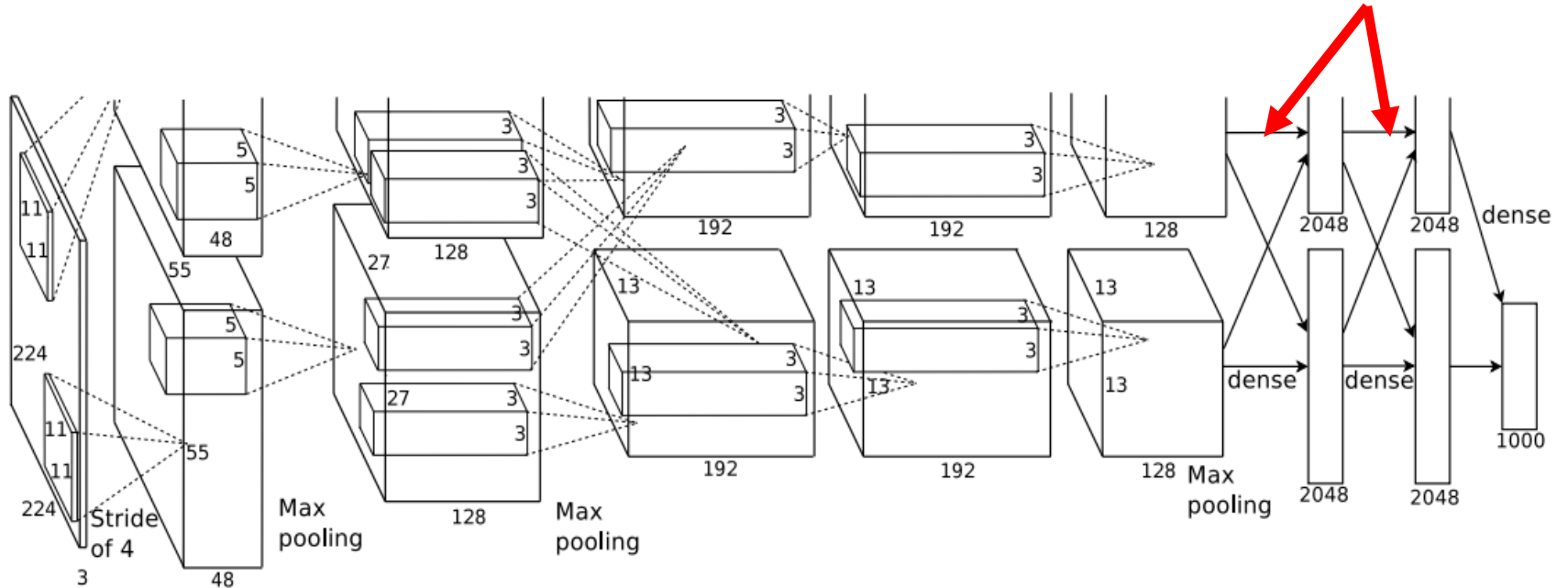
```
weights_grad = evaluate_gradient(loss_fun, data, weights)
vel = vel * 0.9 - step_size * weights_grad
weights += vel
```



# Convolution Neural Networks

Let's take a closer look at AlexNet

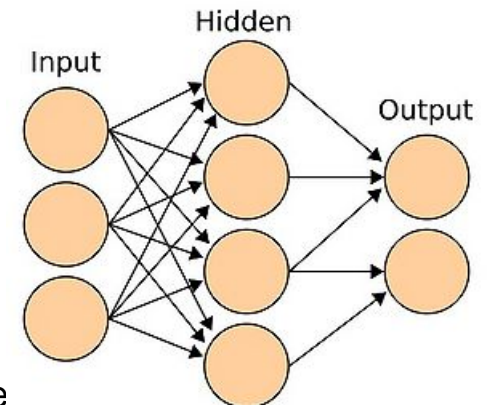
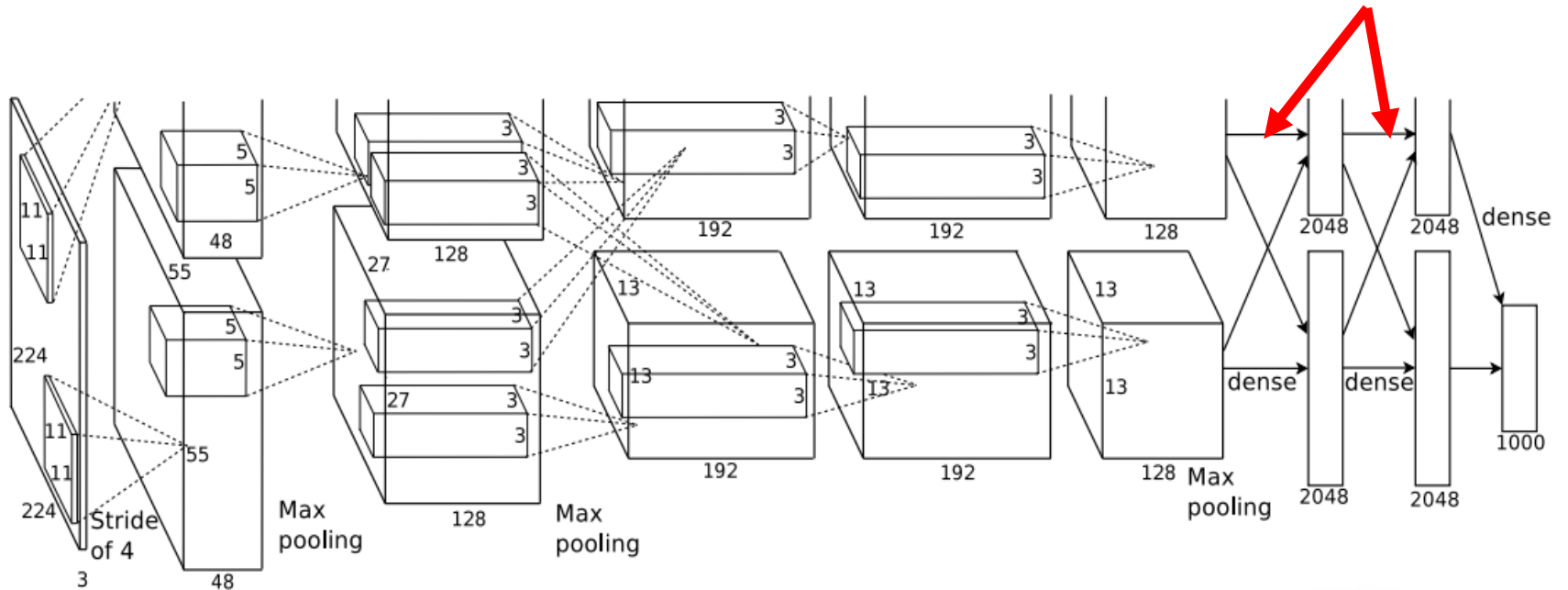
Linear transformation:  
 $y'=Wx+b$



# Convolution Neural Networks

Let's take a closer look at AlexNet

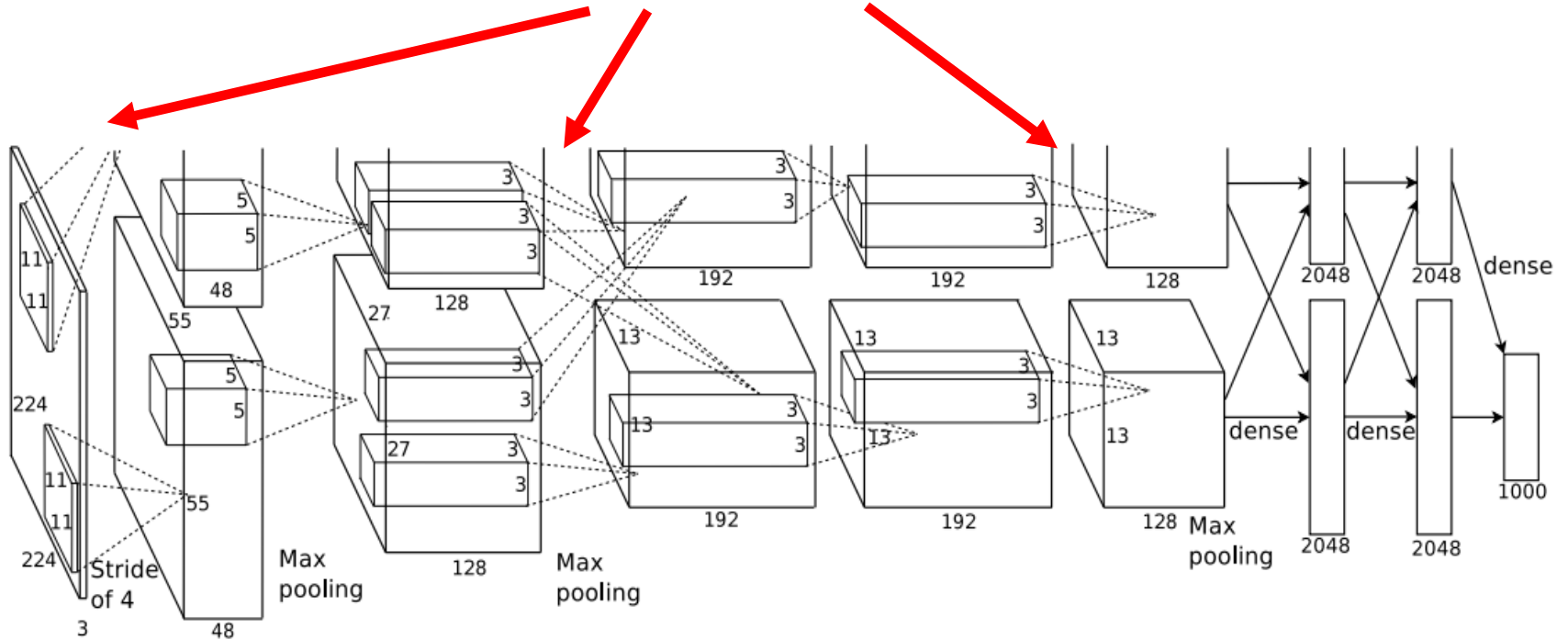
Linear transformation:  
 $y'=Wx+b$



# Convolution Neural Networks

Let's take a closer look at AlexNet

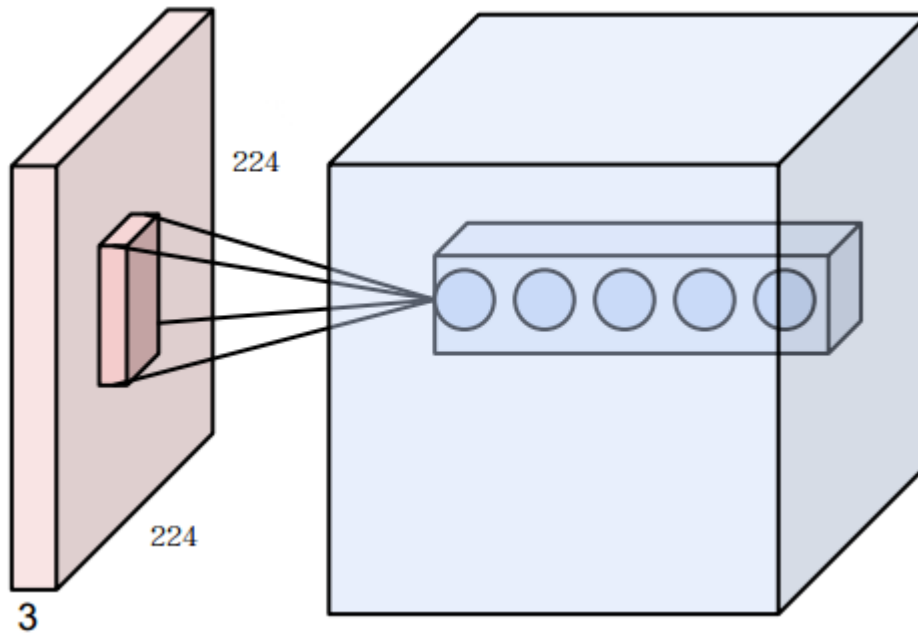
conv(h,w,stroke)





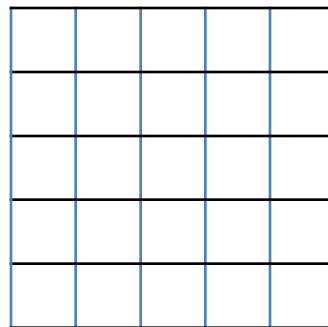
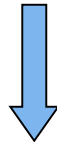
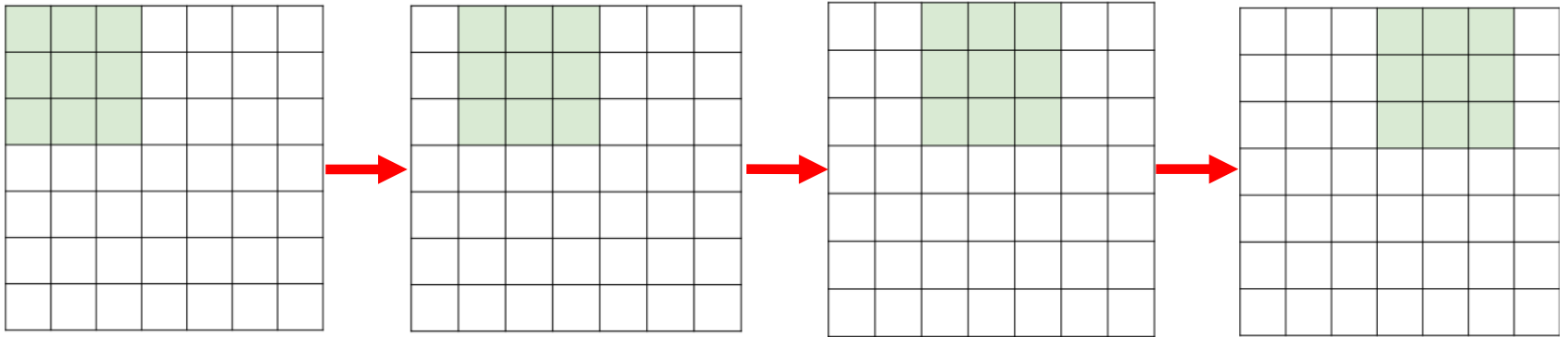
# Convolution Neural Networks

`conv(h,w,stroke)`



# Convolution Neural Networks

Example:  $\text{conv}(h=3, w=3, \text{stride}=1)$

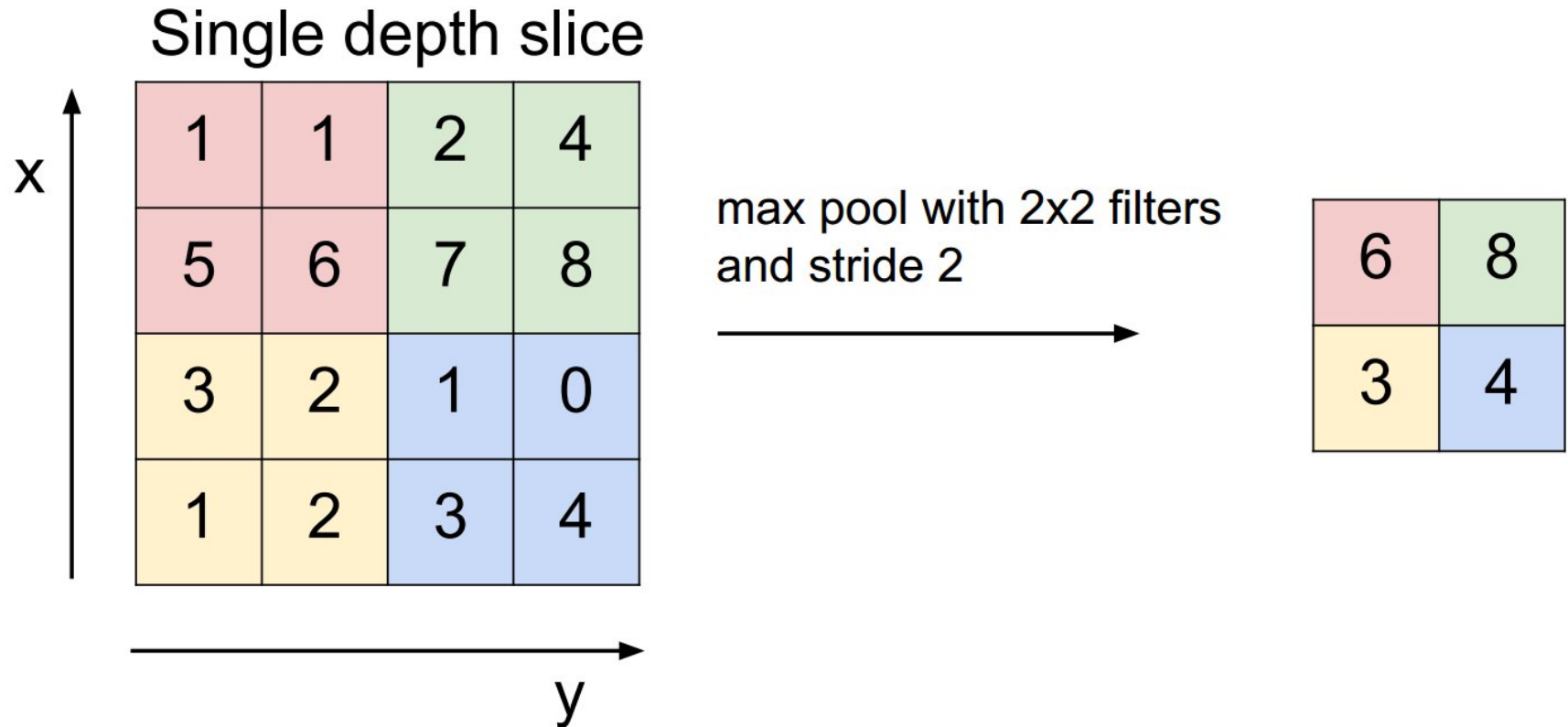


$(7-3)/1+1=5$   
End up as a 5\*5 feature map



# Convolution Neural Networks

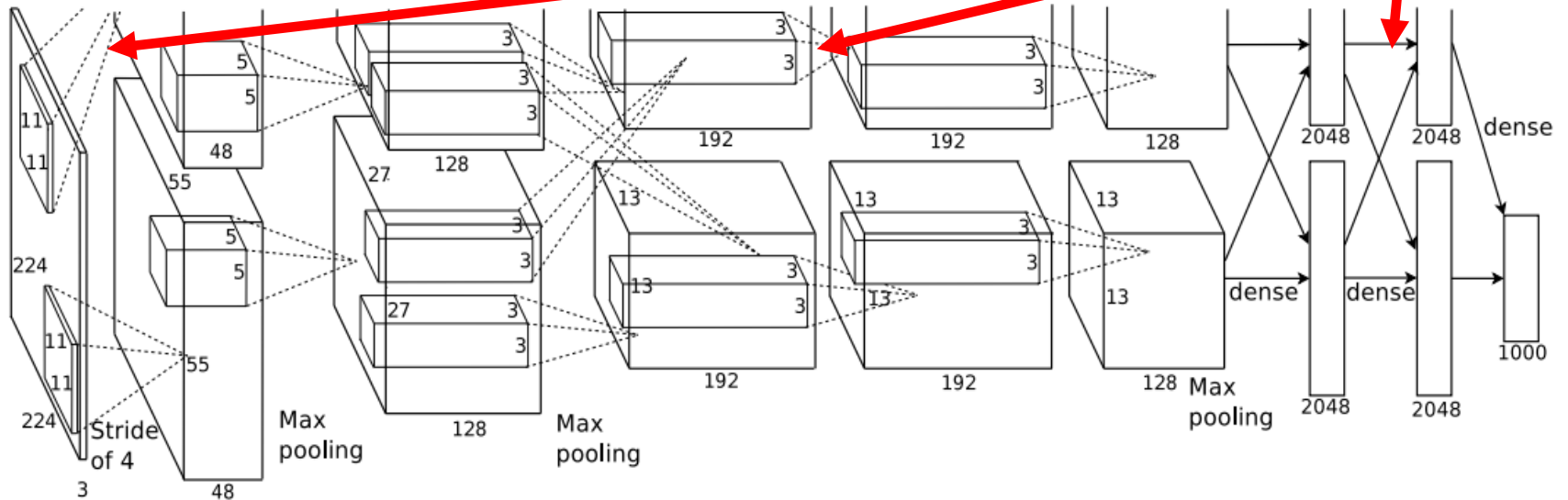
Example: maxpool(h=2,w=2, stride=2)



# Convolution Neural Networks

Let's take a closer look at AlexNet

Relu:  $y = \max(y', 0)$



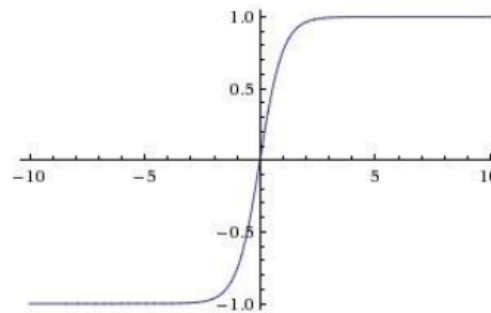


# Convolution Neural Networks

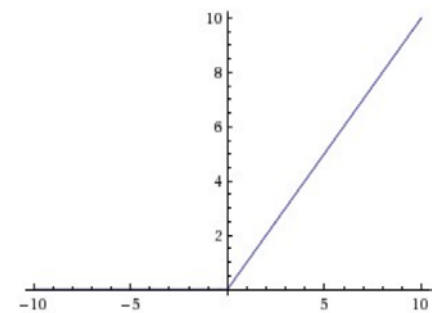
$$\tanh x = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = \max(0, x)$$

Problems with tanh:  
Saturated response



**tanh(x)**



**ReLU**

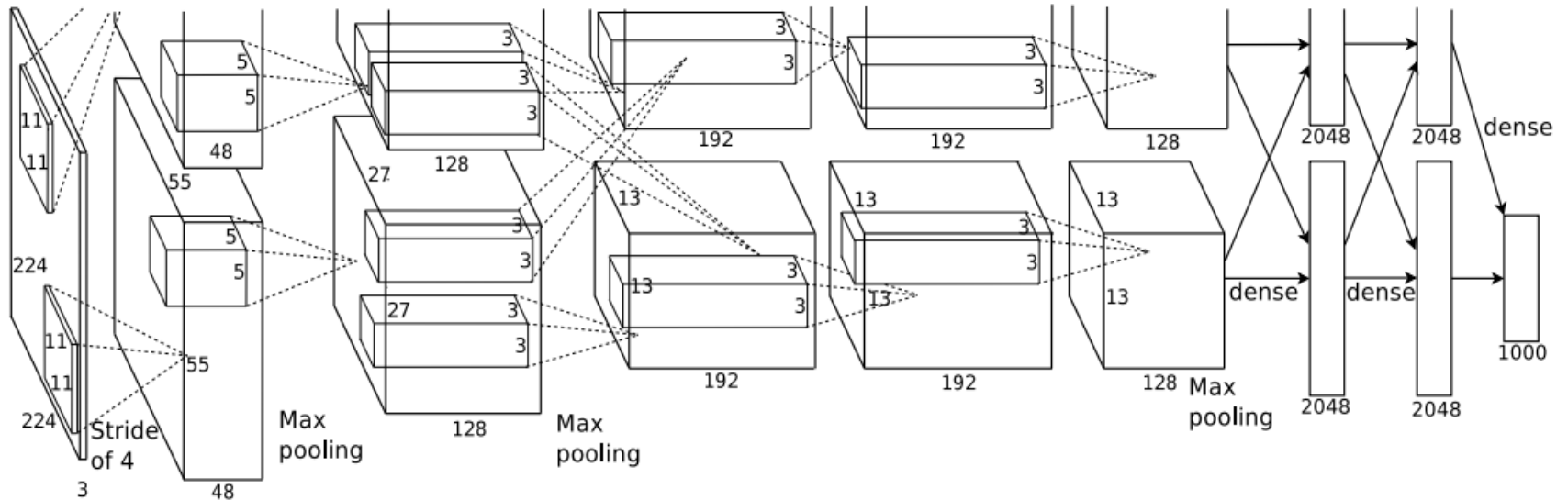
Relu:  $y = \max(y', 0)$

- Does not saturate
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice!

However, non-bounded response and dead when less than 0  
(improved version leaky ReLU)

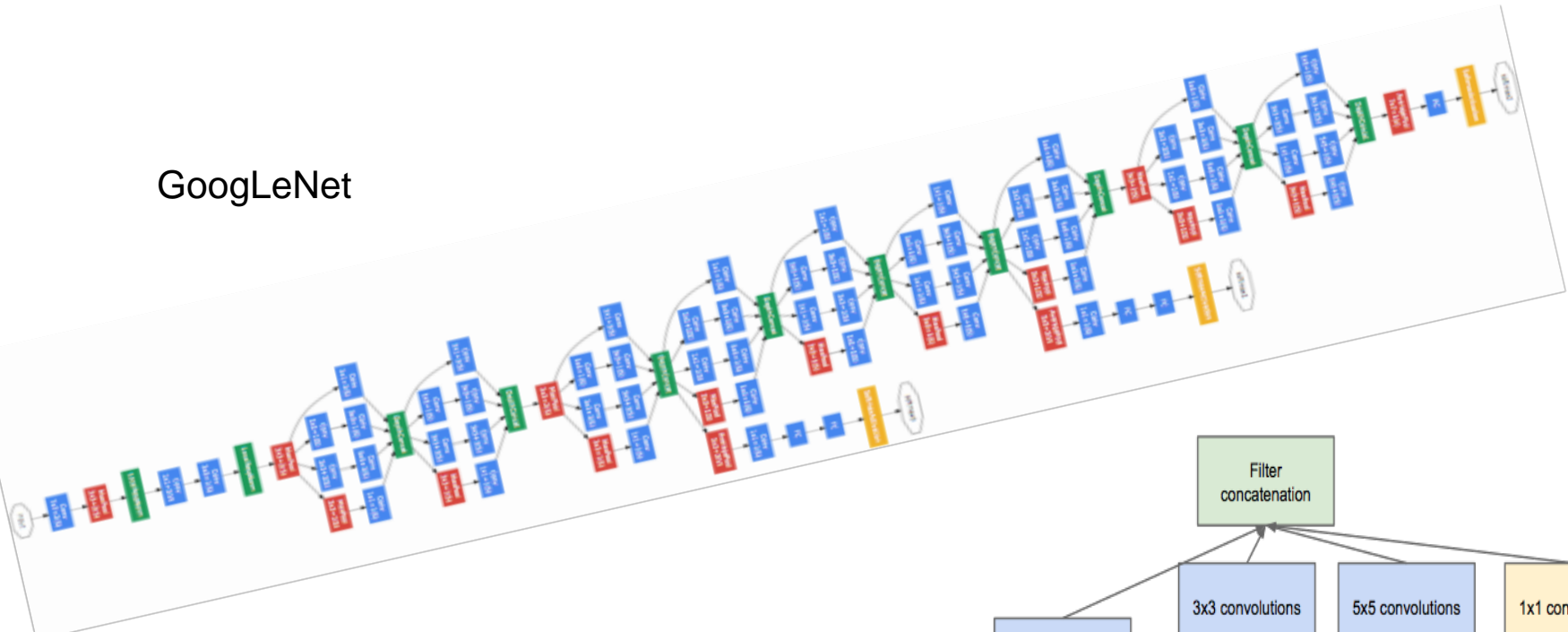
# Convolution Neural Networks

There are two key differences to Vanilla Neural Nets: neurons arranged in 3D volumes have local connectivity, share parameters



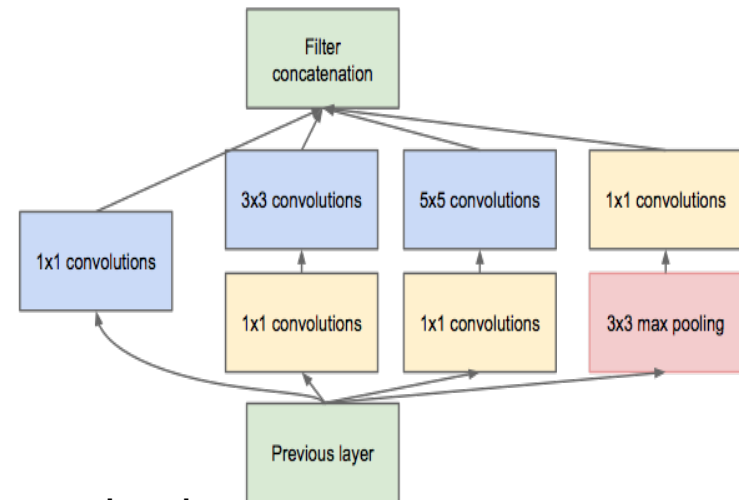
# Convolution Neural Networks

GoogLeNet



ILSVRC14 Winners: **~6.6% Top-5 error**

- GoogLeNet: composition of multi-scale dimension-reduced modules (pictured)
- VGG: 16 layers of 3x3 convolution interleaved with max pooling + 3 fully-connected layers



- + depth
- + data
- + dimensionality reduction

# Convolution Neural Networks

Object Detection

## R-CNN: Region-based Convolutional Networks

<http://nbviewer.ipython.org/github/BVLC/caffe/blob/master/examples/detection.ipynb>

Full R-CNN scripts available at

<https://github.com/rbgirshick/rcnn>

Ross Girshick et al.

*Rich feature hierarchies for accurate object detection and semantic segmentation. CVPR14.*

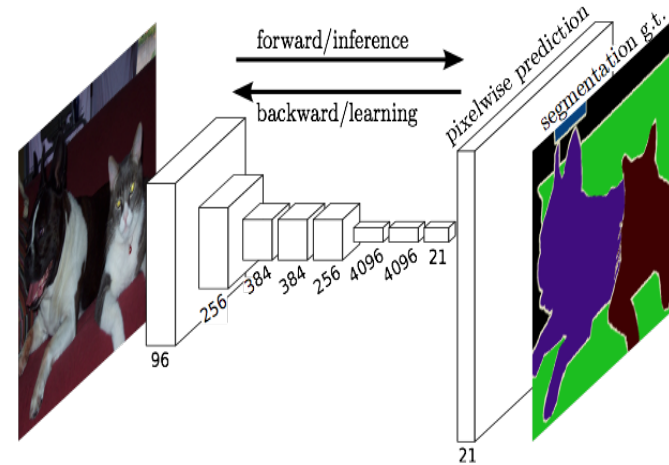
Fast R-CNN

[arXiv](#) and [code](#)



# Segmentation

Fully convolutional networks for pixel prediction  
applied to semantic segmentation  
end-to-end learning  
efficiency in inference and learning  
175 ms per-image prediction  
multi-modal, multi-task



Further applications

- depth estimation
- denoising

[arXiv](#) and [pre-release](#)



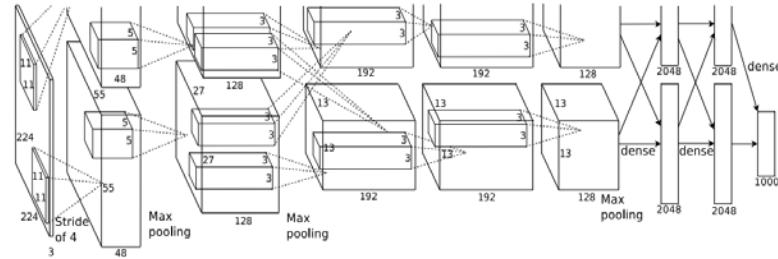
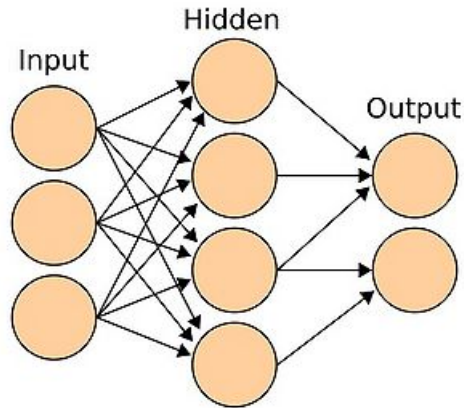
Jon Long\* & Evan

Shelhamer\*,





# Problem with Feed-forward Nets



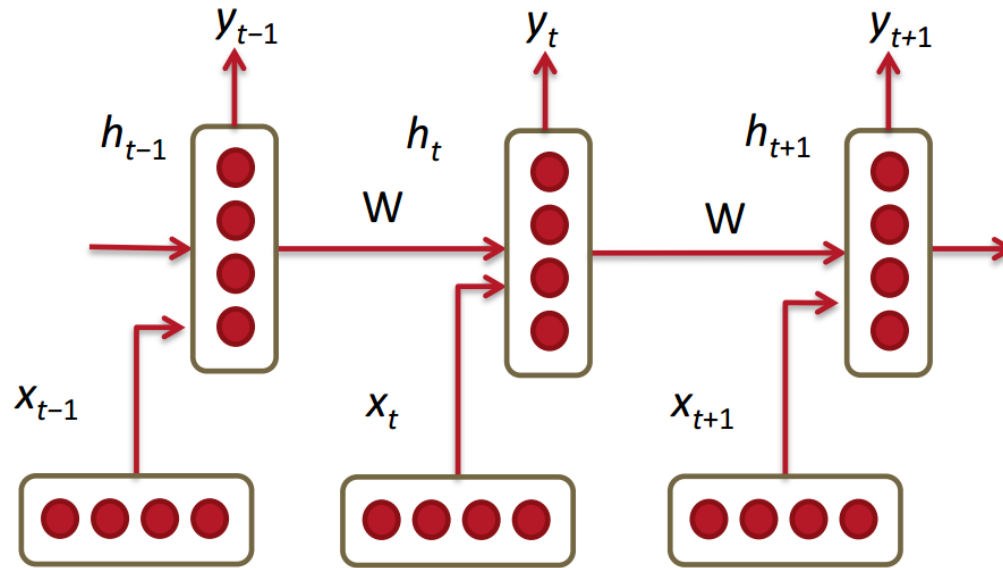
What if we want to be able to have a model telling us what's the probability of the following two sentences, respectively:

1. The cat sat on the mat
2. The mat is having dinner with the cat

Cannot handle variable length input

# Recurrent Neural Net

RNNs tie the weights at each time step

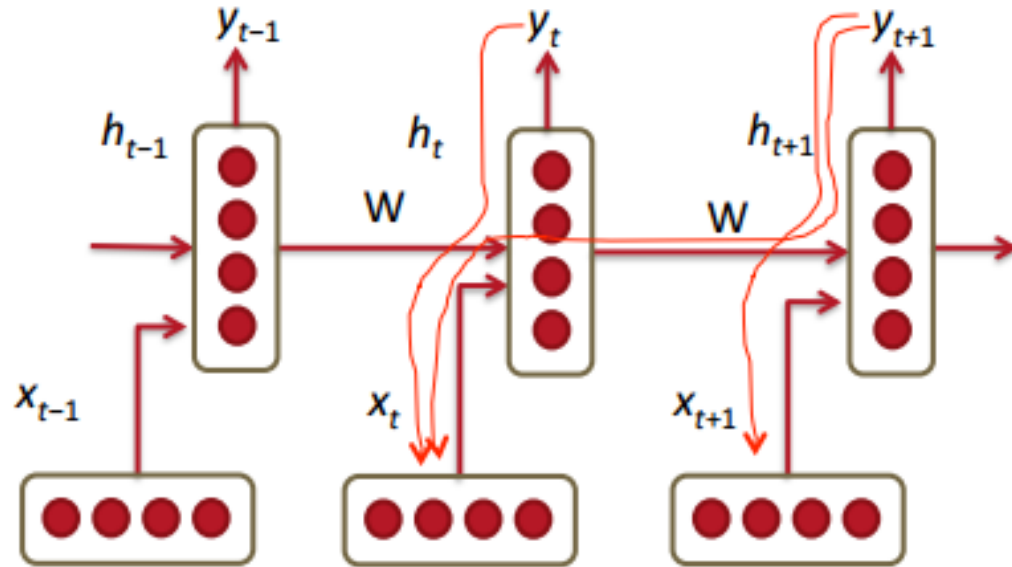


$$h_t = \sigma \left( W^{(hh)} h_{t-1} + W^{(hx)} x_{[t]} \right)$$

$$\hat{y}_t = \text{softmax} \left( W^{(S)} h_t \right)$$

# Recurrent Neural Net

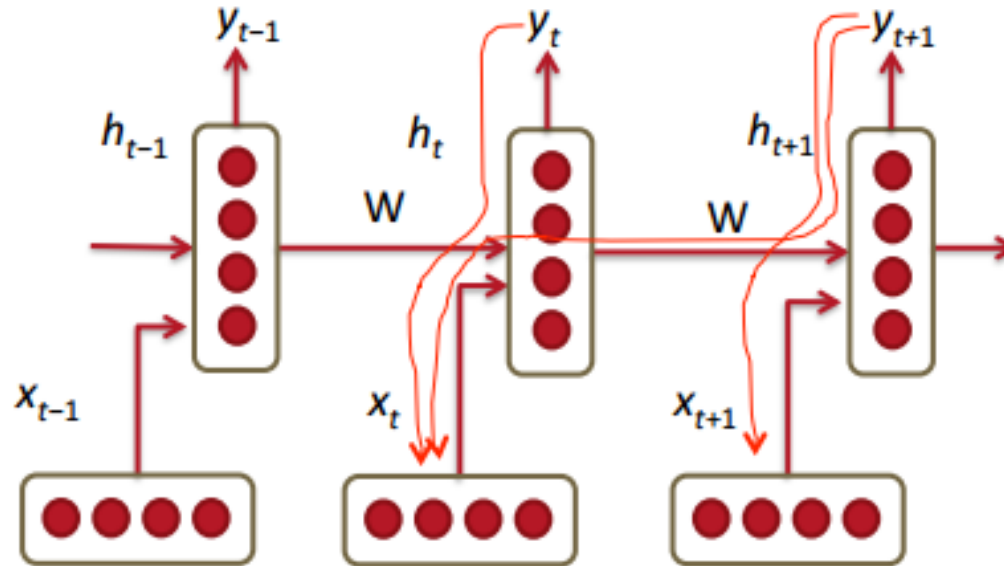
Training of RNNs is hard...



$$h_t = W f(h_{t-1}) + W^{(hx)} x[t]$$
$$\frac{\partial h_t}{\partial h_k} = \prod_{j=k+1}^t \frac{\partial h_j}{\partial h_{j-1}} = \prod_{j=k+1}^t W^T \text{diag}[f'(h_{j-1})]$$

# Recurrent Neural Net

Training of RNNs is hard...



**Solution 1: clip the gradient!**

---

**Algorithm 1** Pseudo-code for norm clipping the gradients whenever they explode

---

```
 $\hat{\mathbf{g}} \leftarrow \frac{\partial \mathcal{E}}{\partial \theta}$   
if  $\|\hat{\mathbf{g}}\| \geq \text{threshold}$  then  
   $\hat{\mathbf{g}} \leftarrow \frac{\text{threshold}}{\|\hat{\mathbf{g}}\|} \hat{\mathbf{g}}$   
end if
```

---

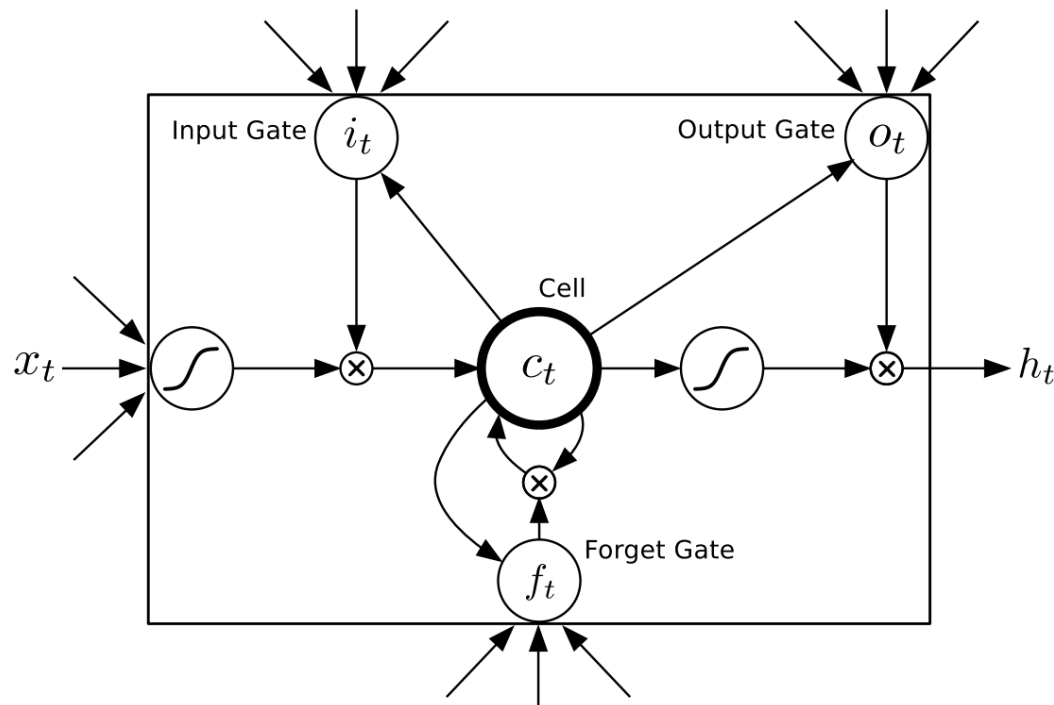
Some theory: On the difficulty of training recurrent neural networks, Pascanu et al. ICML2013



# Recurrent Neural Net

Training of RNNs is hard...

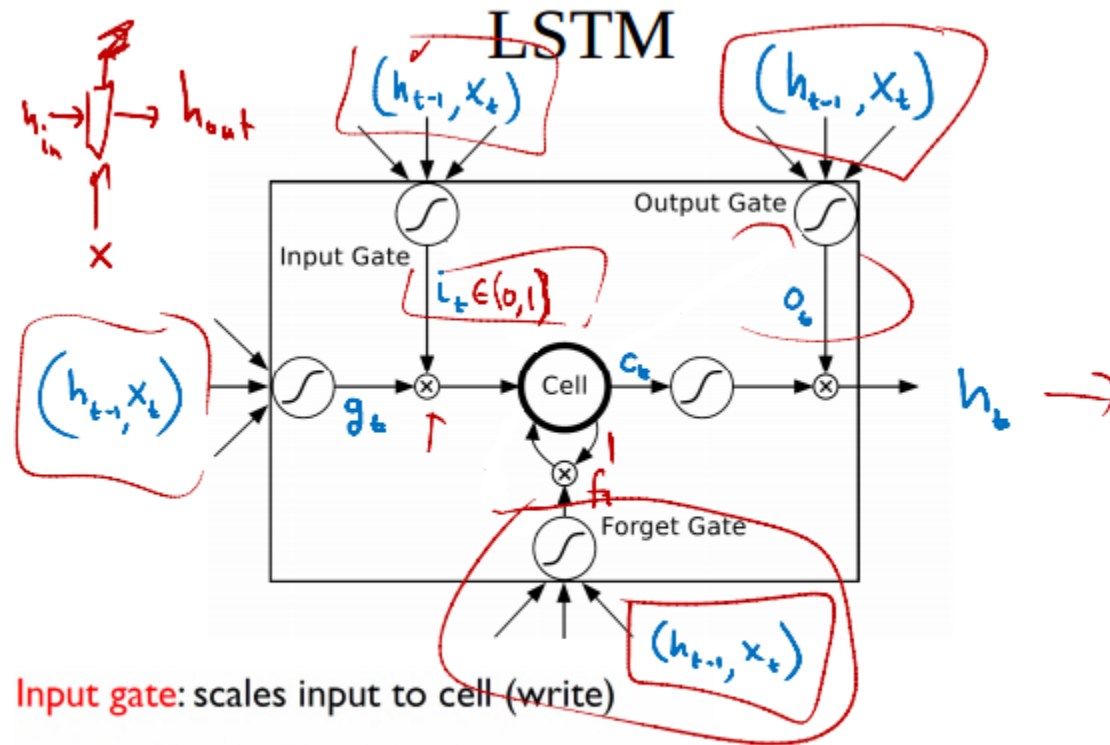
Solution 2: NNs with gating units (LSTM/GRU)



# Recurrent Neural Net

Training of RNNs is hard...

Solution 2: nets with gating units (LSTM/GRU)



**Input gate:** scales input to cell (write)

**Output gate:** scales output from cell (read)

**Forget gate:** scales old cell value (reset)

# Recurrent Neural Net

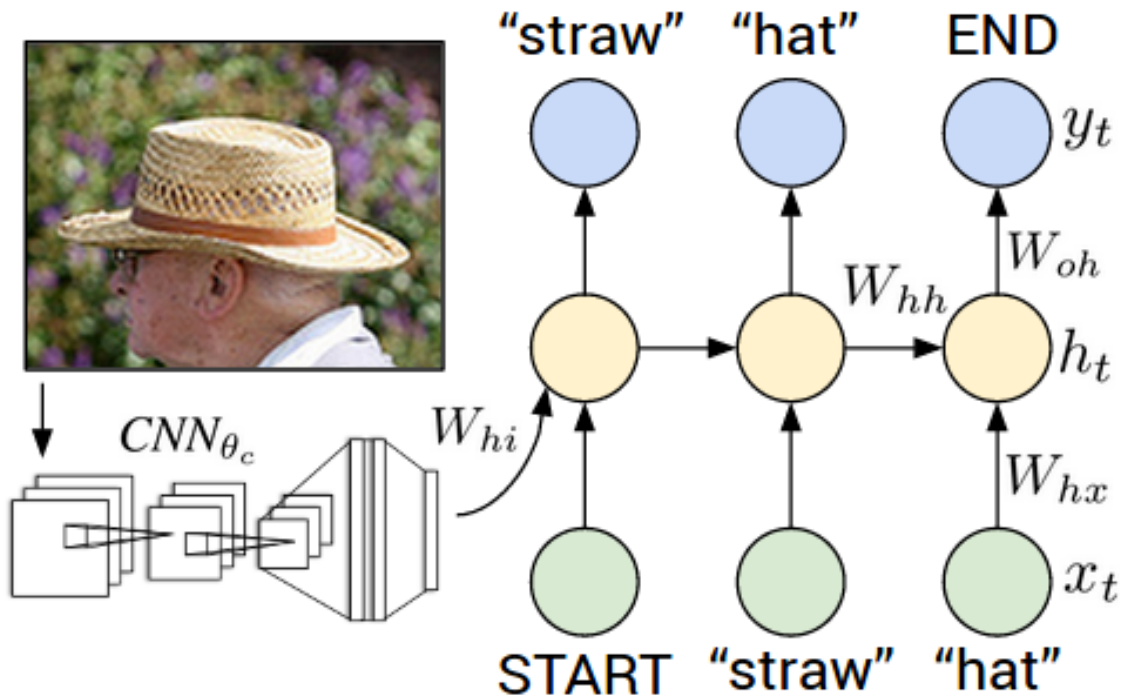
Training of RNNs is hard...

Solution 2: nets with gating units (LSTM/GRU)

$$\begin{aligned} \mathbf{i}_t &= \text{Sigm}(\boldsymbol{\theta}_{xi}\mathbf{x}_t + \boldsymbol{\theta}_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i) \\ \mathbf{f}_t &= \text{Sigm}(\boldsymbol{\theta}_{xf}\mathbf{x}_t + \boldsymbol{\theta}_{hf}\mathbf{h}_{t-1} + \mathbf{b}_f) \\ \mathbf{o}_t &= \text{Sigm}(\boldsymbol{\theta}_{xo}\mathbf{x}_t + \boldsymbol{\theta}_{ho}\mathbf{h}_{t-1} + \mathbf{b}_o) \\ \mathbf{g}_t &= \text{Tanh}(\boldsymbol{\theta}_{xg}\mathbf{x}_t + \boldsymbol{\theta}_{hg}\mathbf{h}_{t-1} + \mathbf{b}_g) \\ \mathbf{c}_t &= \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \mathbf{g}_t \\ \mathbf{h}_t &= \mathbf{o}_t \odot \text{Tanh}(\mathbf{c}_t) \end{aligned}$$
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \odot \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} x_1 y_1 \\ x_2 y_2 \end{bmatrix}$$

# RNN in vision

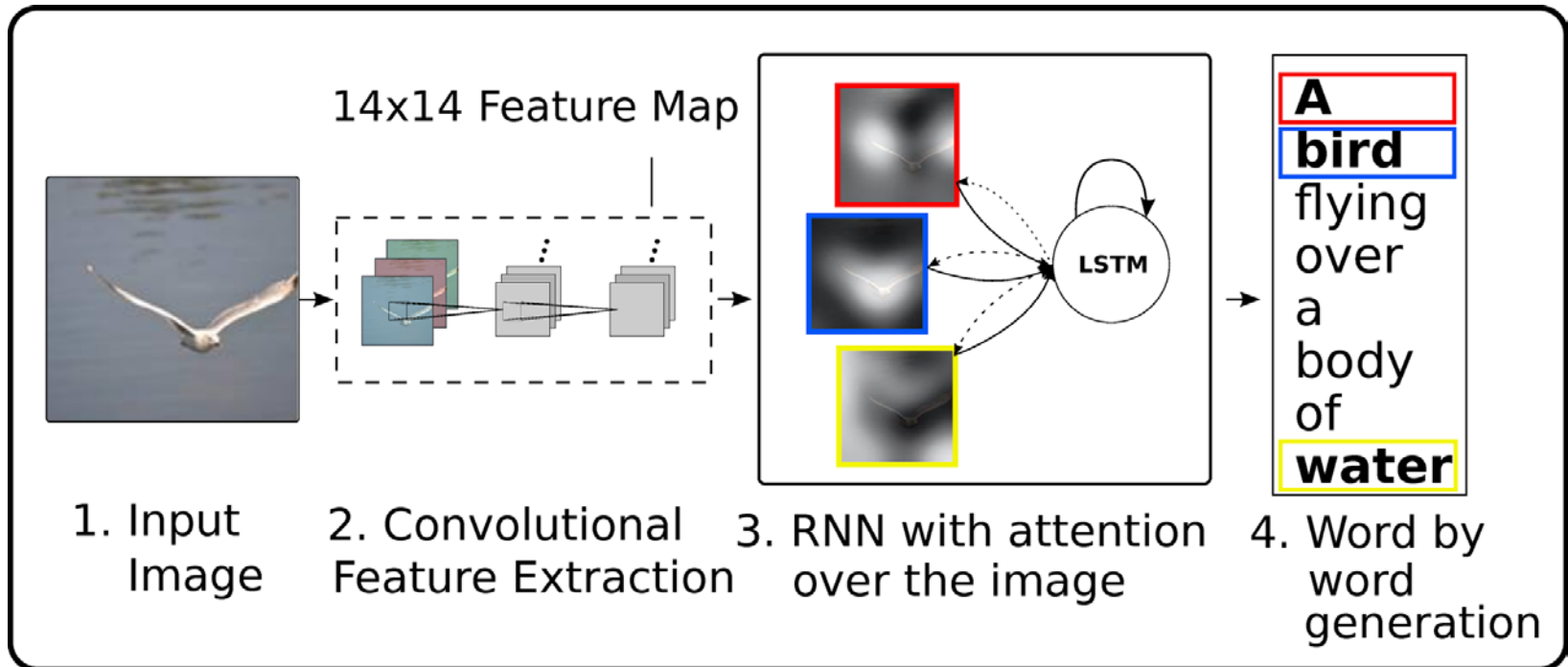
Image captioning



Deep Visual-Semantic Alignments for Generating Image Descriptions, Andrej Karpathy et al.

# RNN in vision

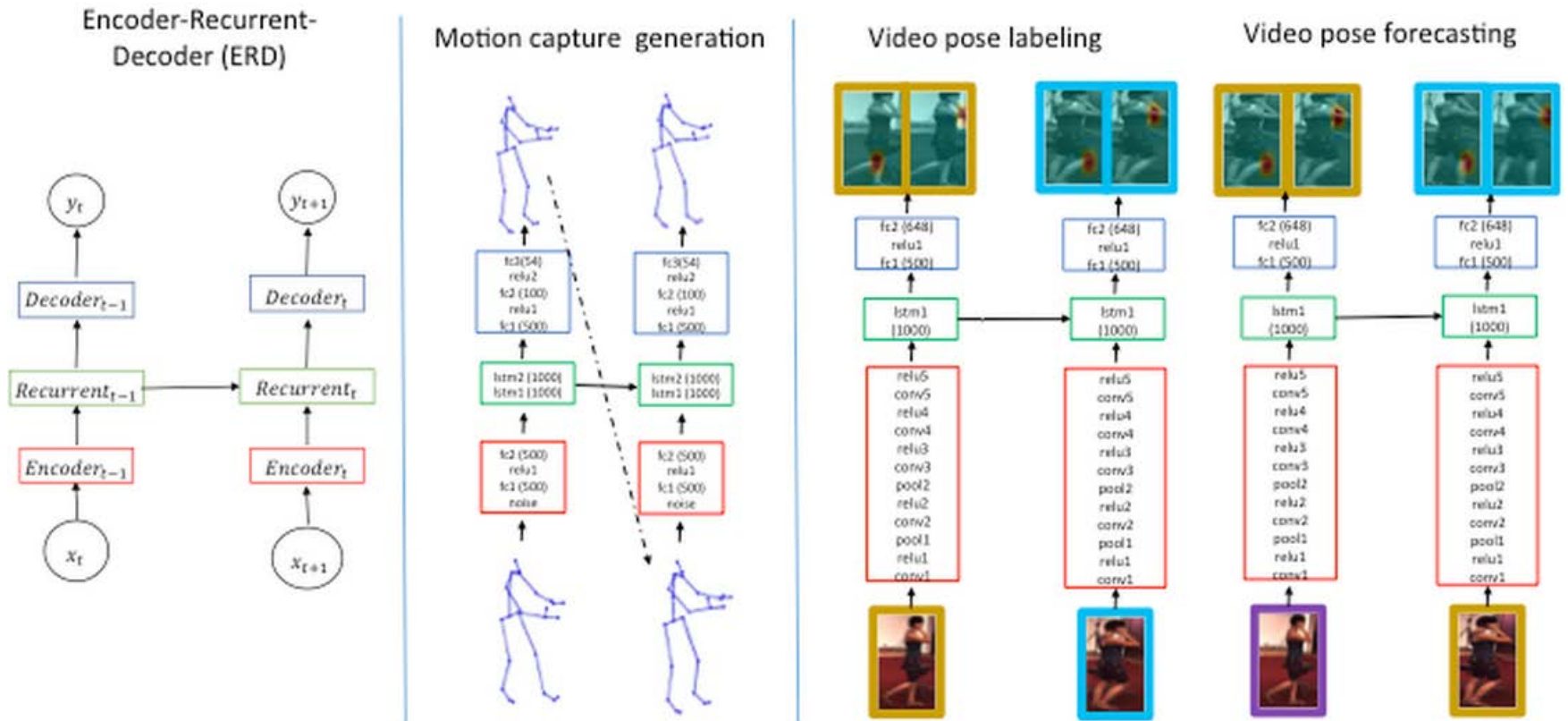
## Visual attention model



Show, Attend and Tell: Neural Image Caption Generation with Visual Attention, Kelvin Xu et al.

# RNN in vision

## RNNs for Human Dynamics



Recurrent Network Models for Human Dynamics, Katerina Fragkiadaki et al.

Most materials taken from Andrej Karpathy/Richard Socher/Nando de Freitas/caffe CVPR tutorial



# Tricks

## 1. Numerical gradient check

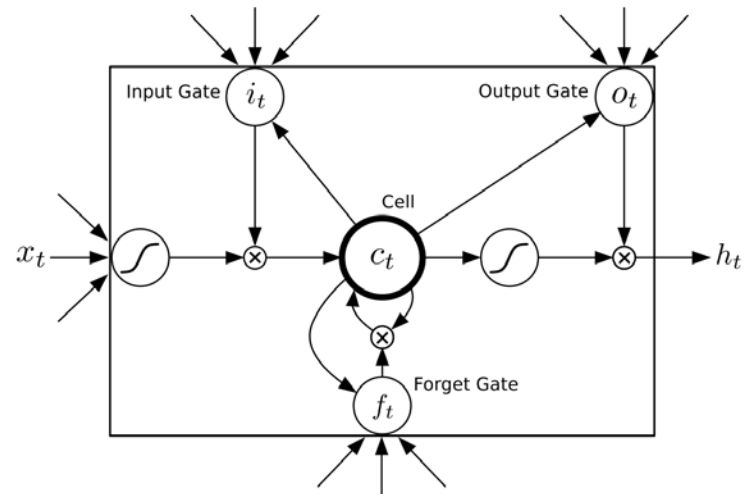
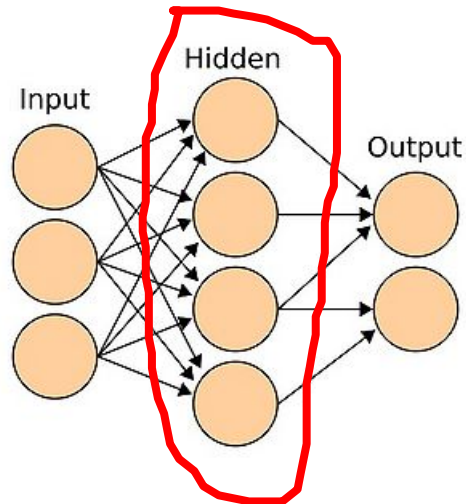
```
fx = f(x) # evaluate function value at original point
grad = np.zeros_like(x)
# iterate over all indexes in x
it = np.nditer(x, flags=['multi_index'], op_flags=['readwrite'])
while not it.finished:

    # evaluate function at x+h
    ix = it.multi_index
    oldval = x[ix]
    x[ix] = oldval + h # increment by h
    fxph = f(x) # evaluate f(x + h)
    x[ix] = oldval - h
    fxmh = f(x) # evaluate f(x - h)
    x[ix] = oldval # restore

    # compute the partial derivative with centered formula
    grad[ix] = (fxph - fxmh) / (2 * h) # the slope
    if verbose:
        print ix, grad[ix]
    it.iternext() # step to next dimension
```

# Tricks

1. Numerical gradient check
2. Modulize layers: only three functions needed
  - (1)  $\text{output} = \text{forward}(\text{input}, \text{model})$
  - (2)  $\text{dJ\_dW} = \text{computeParamGrad}(\text{input}, \text{outputGrad}, \text{model})$
  - (3)  $\text{dJ\_dInput} = \text{computeInputGrad}(\text{input}, \text{outputGrad}, \text{model})$Everything else is just putting together lego pieces



# Questions?

Thanks!