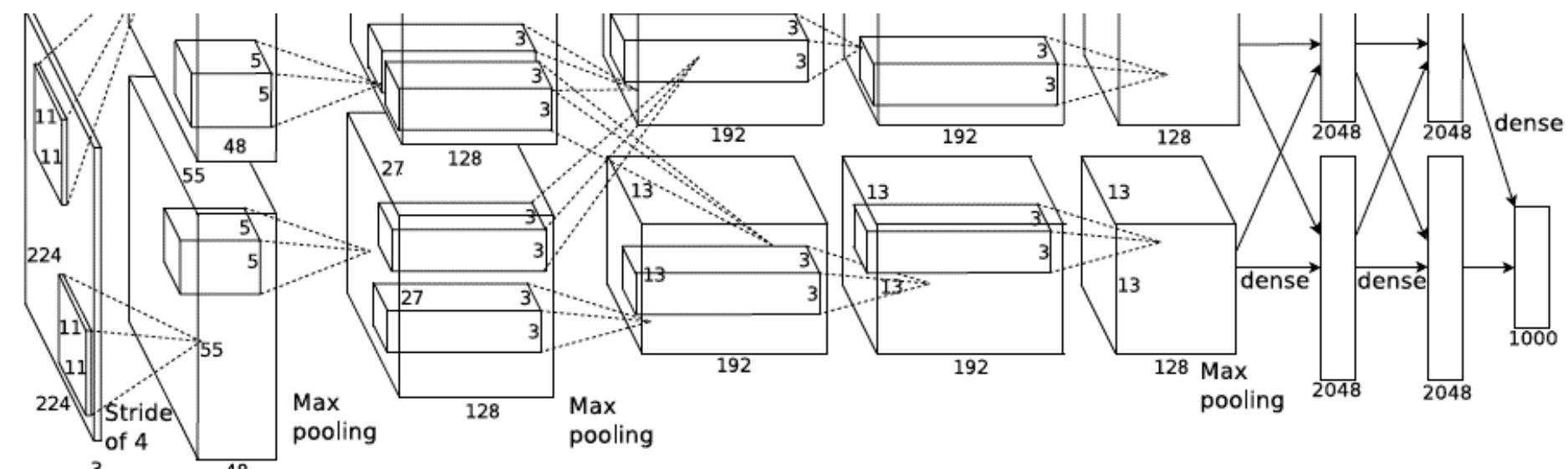


# Convolutional Neural Networks



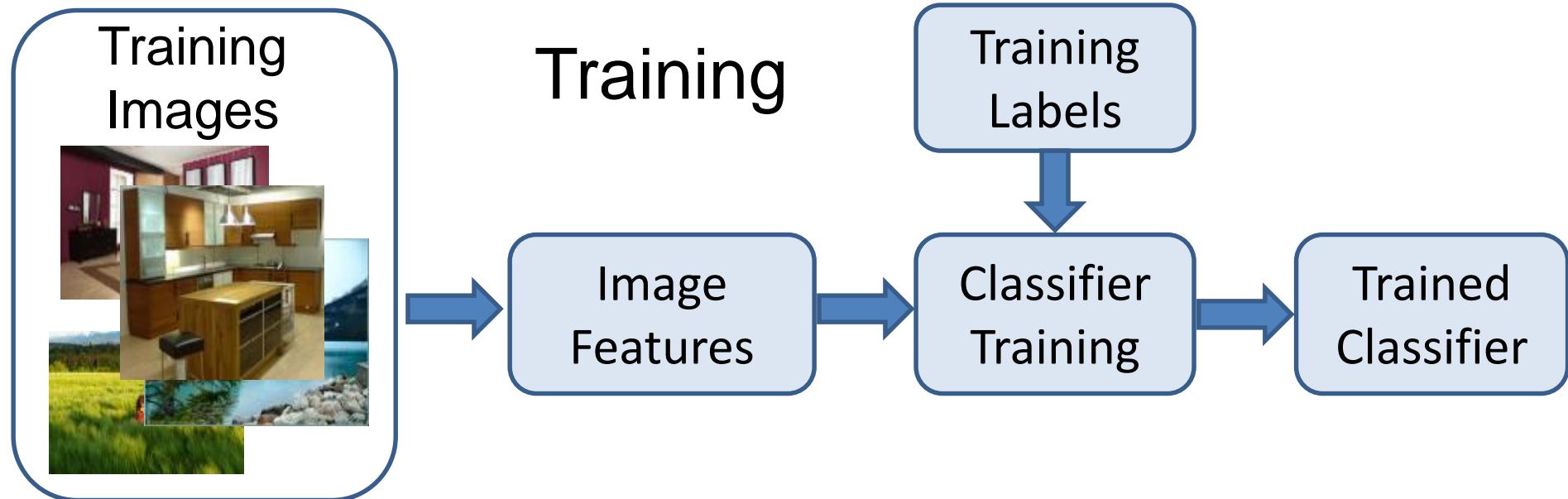
Computer Vision

Jia-Bin Huang, Virginia Tech

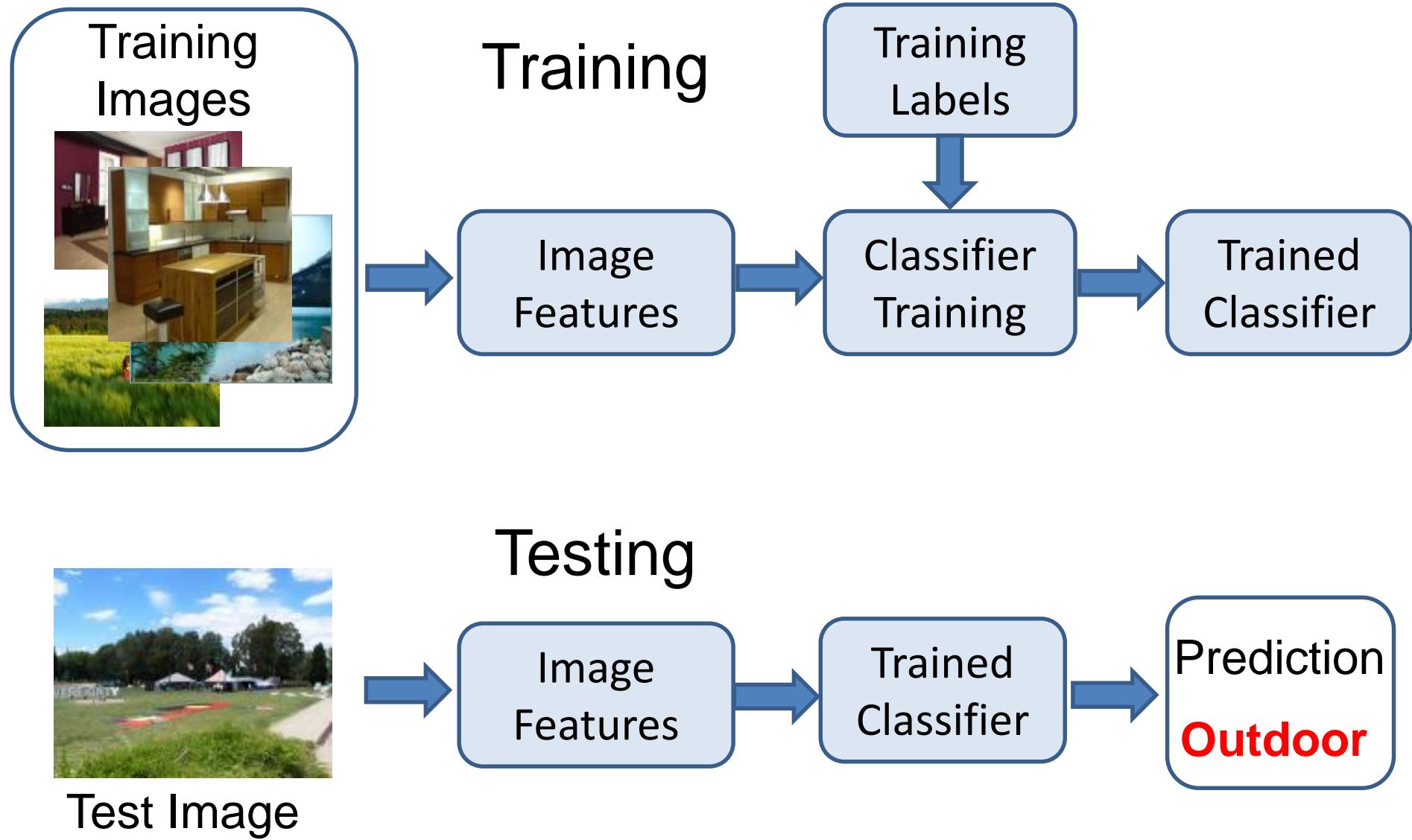
# Today's class

- Overview
- Convolutional Neural Network (CNN)
- Training CNN
- Understanding and Visualizing CNN

# Image Categorization: Training phase



# Image Categorization: Testing phase



# Features are the Keys

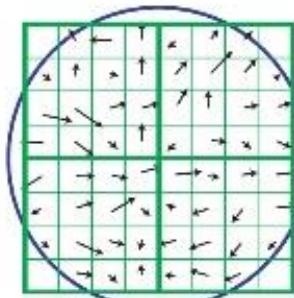
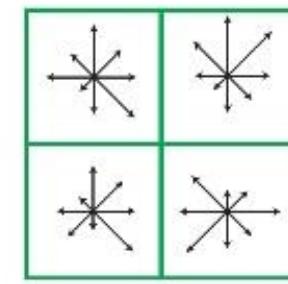
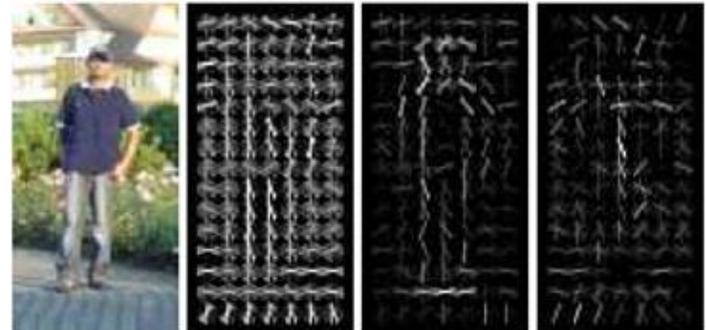


Image gradients

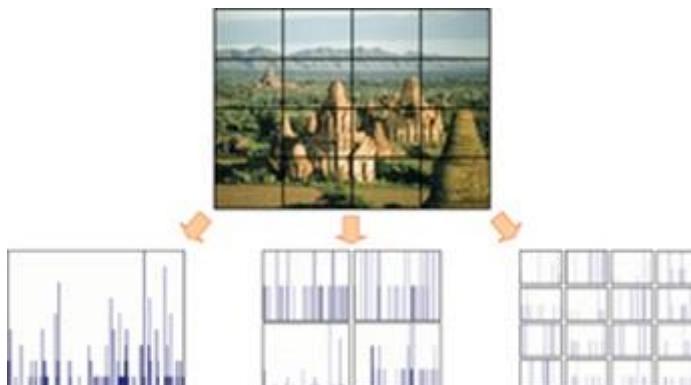


Keypoint descriptor

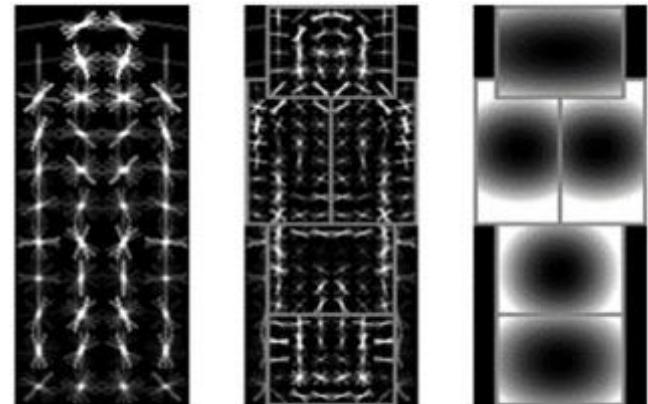
SIFT [Loeve IJCV 04]



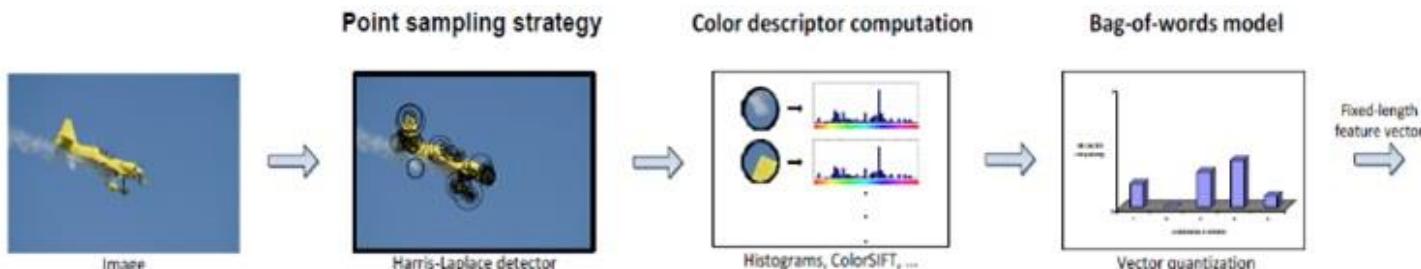
HOG [Dalal and Triggs CVPR 05]



SPM [Lazebnik et al. CVPR 06]



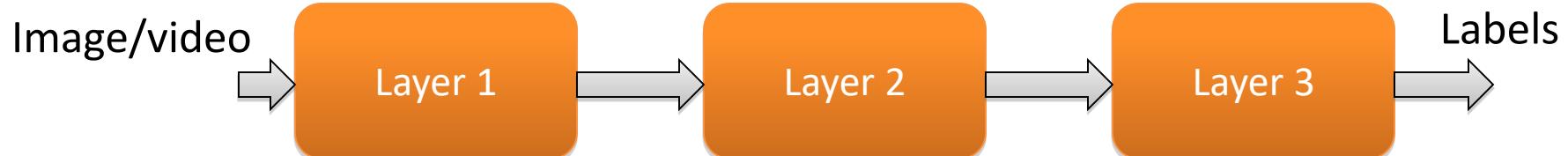
DPM [Felzenszwalb et al. PAMI 10]



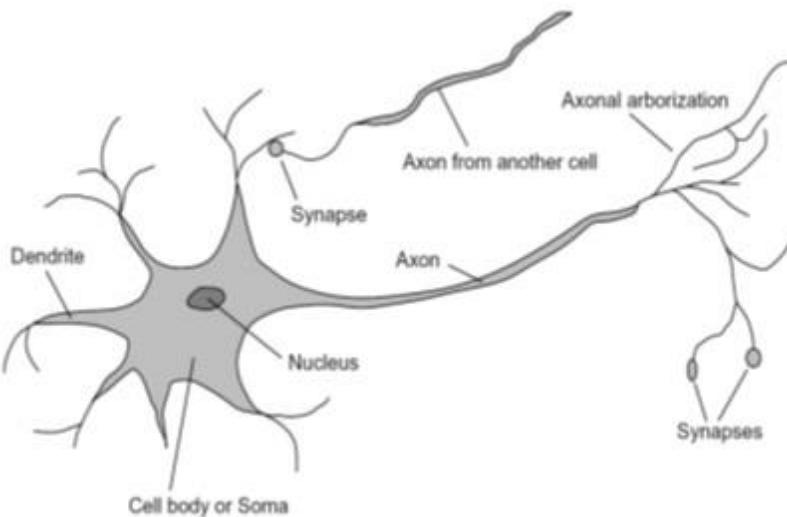
Color Descriptor [Van De Sande et al. PAMI 10]

# Learning a Hierarchy of Feature Extractors

- Each layer of hierarchy extracts features from output of previous layer
- All the way from pixels → classifier
- Layers have the (nearly) same structure



# Biological neuron and Perceptrons



A biological neuron

Input

Weights

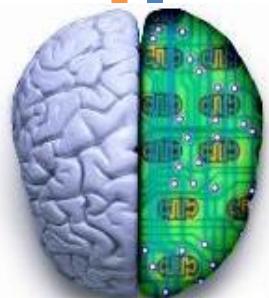
$$\begin{aligned}x_1 &\rightarrow w_1 \\x_2 &\rightarrow w_2 \\x_3 &\rightarrow w_3 \\\vdots &\\\vdots &\\\boldsymbol{x}_d &\rightarrow w_d\end{aligned}$$

Output:  $\sigma(\mathbf{w} \cdot \mathbf{x} + b)$

Sigmoid function:

$$\sigma(t) = \frac{1}{1 + e^{-t}}$$

An artificial neuron (Perceptron)  
- a linear classifier



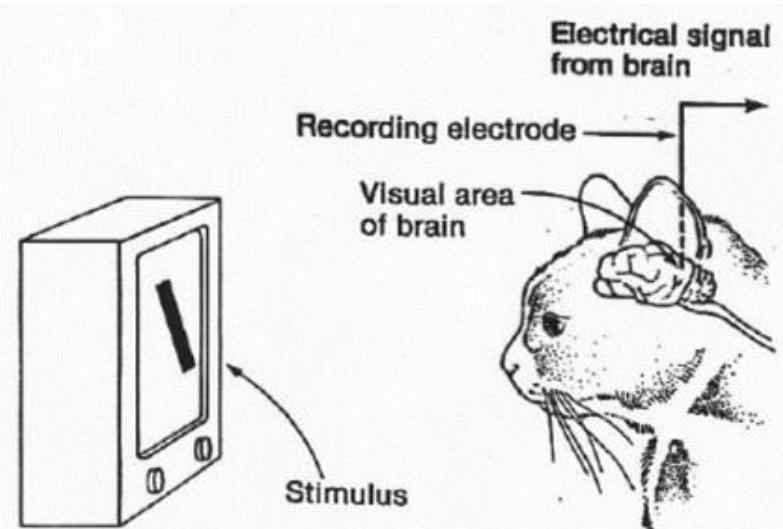
# Simple, Complex and Hypercomplex cells



David H. Hubel and Torsten Wiesel

Suggested a **hierarchy of feature detectors** in the visual cortex, with higher level features responding to patterns of activation in lower level cells, and propagating activation upwards to still higher level cells.

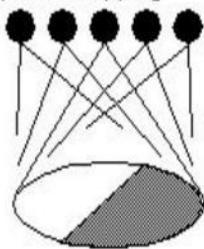
David Hubel's [Eye, Brain, and Vision](#)



# Hubel/Wiesel Architecture and Multi-layer Neural Network

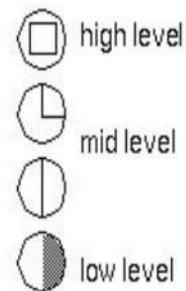
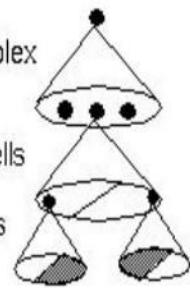
Hubel & Weisel

topographical mapping

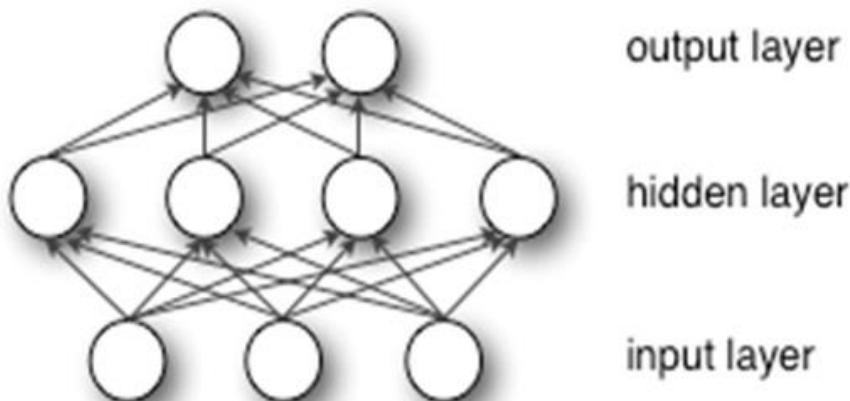


featural hierarchy

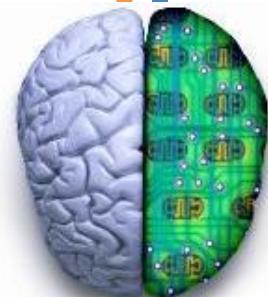
hyper-complex  
cells  
complex cells  
simple cells



Hubel and Weisel's architecture



Multi-layer Neural Network  
- A *non-linear* classifier

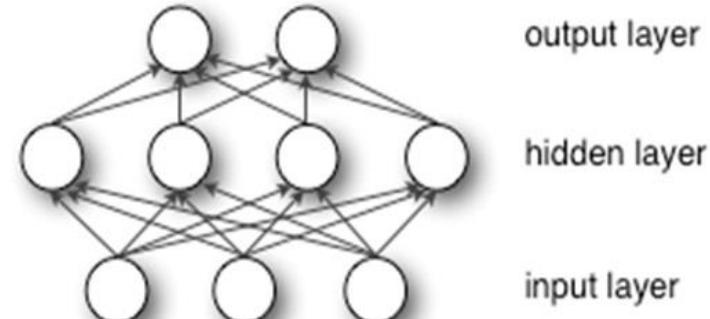


# Multi-layer Neural Network

- A non-linear classifier
- **Training:** find network weights  $\mathbf{w}$  to minimize the error between true training labels  $y_i$  and estimated labels  $f_{\mathbf{w}}(\mathbf{x}_i)$

$$E(\mathbf{w}) = \sum_{i=1}^N (y_i - f_{\mathbf{w}}(\mathbf{x}_i))^2$$

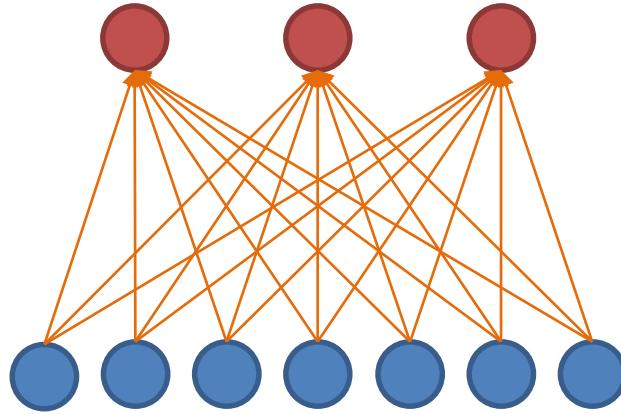
- Minimization can be done by gradient descent provided  $f$  is differentiable
- This training method is called **back-propagation**



# Convolutional Neural Networks

- Also known as CNN, ConvNet, DCN
- CNN = a multi-layer neural network with
  1. Local connectivity
  2. Weight sharing

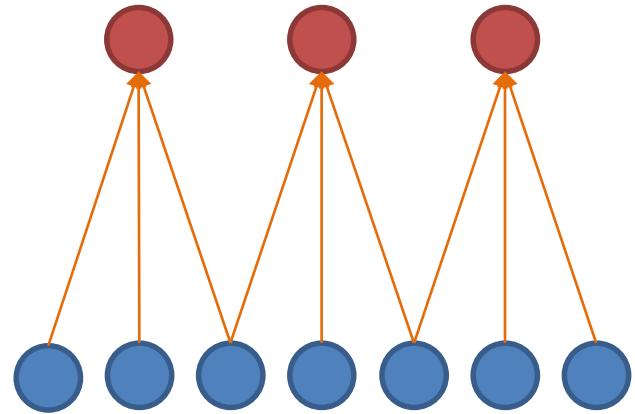
# CNN: Local Connectivity



**Global connectivity**

Hidden layer

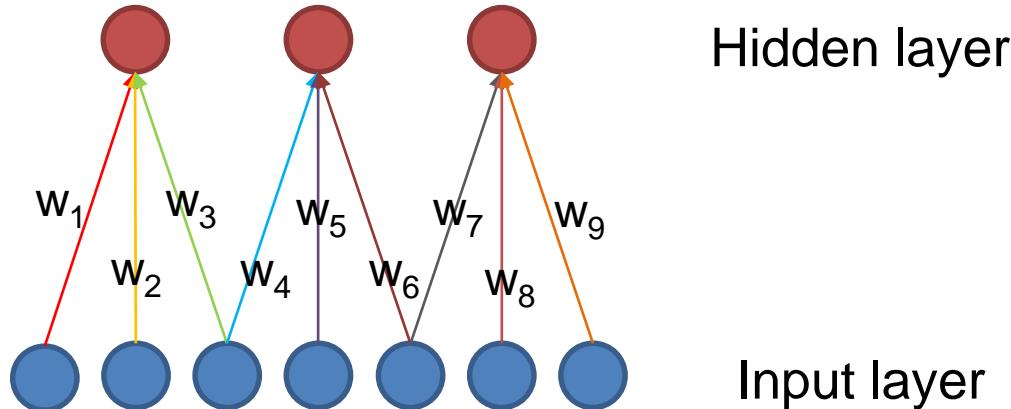
Input layer



**Local connectivity**

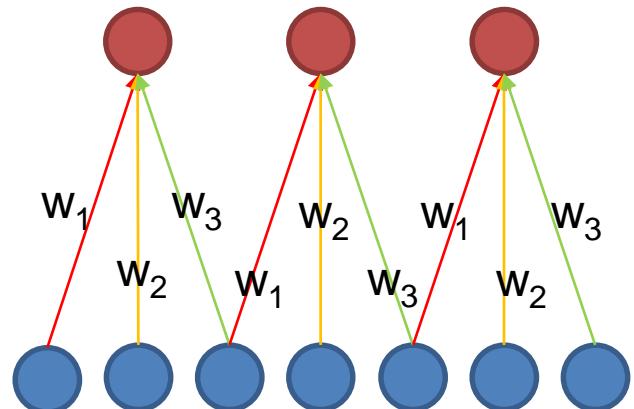
- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
  - Global connectivity:  $3 \times 7 = 21$
  - Local connectivity:  $3 \times 3 = 9$

# CNN: Weight Sharing



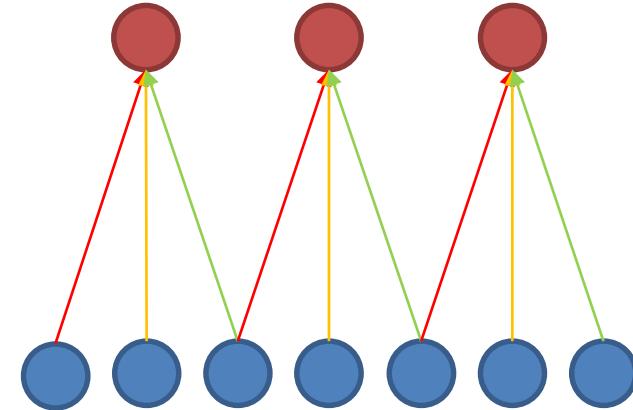
**Without weight sharing**

- # input units (neurons): 7
- # hidden units: 3
- Number of parameters
  - Without weight sharing:  $3 \times 3 = 9$
  - With weight sharing :  $3 \times 1 = 3$



**With weight sharing**

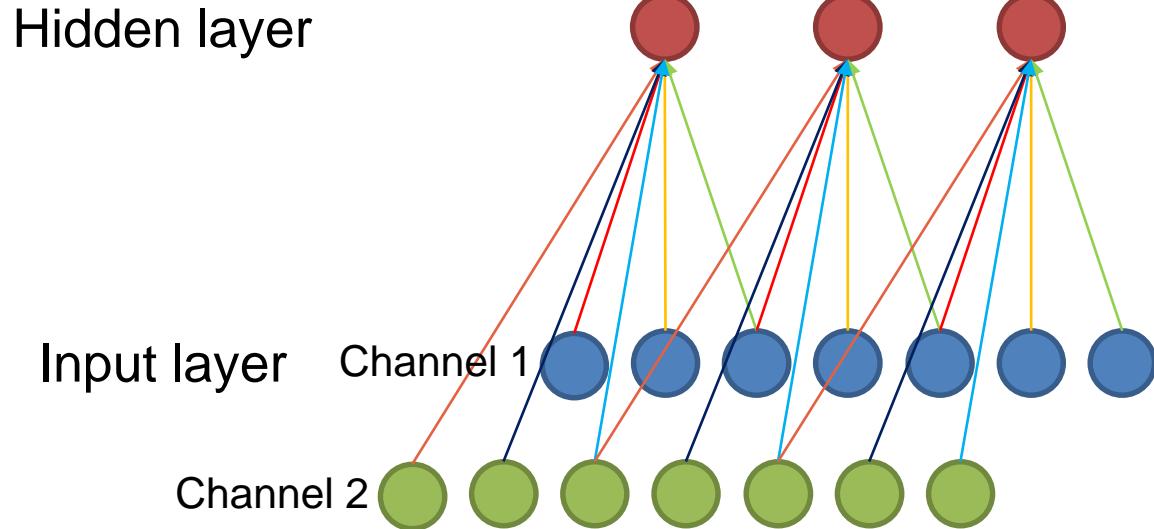
# CNN with multiple input channels



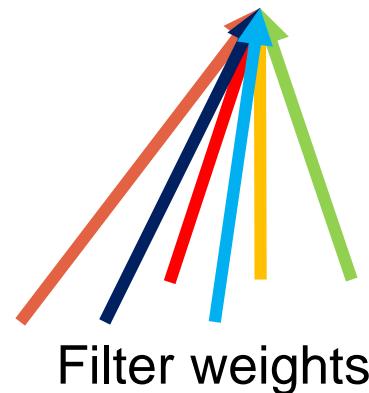
**Single input channel**



**Filter weights**

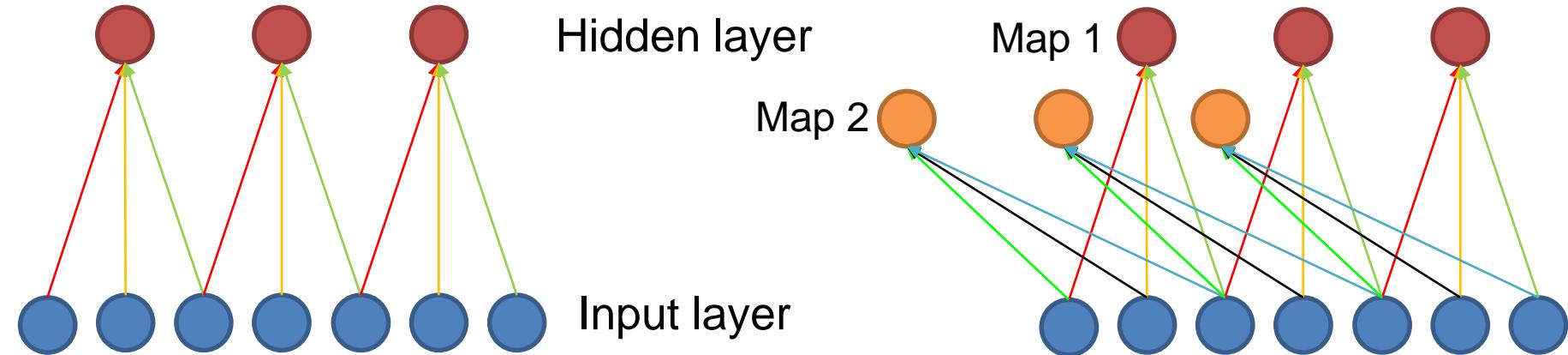


**Multiple input channels**



**Filter weights**

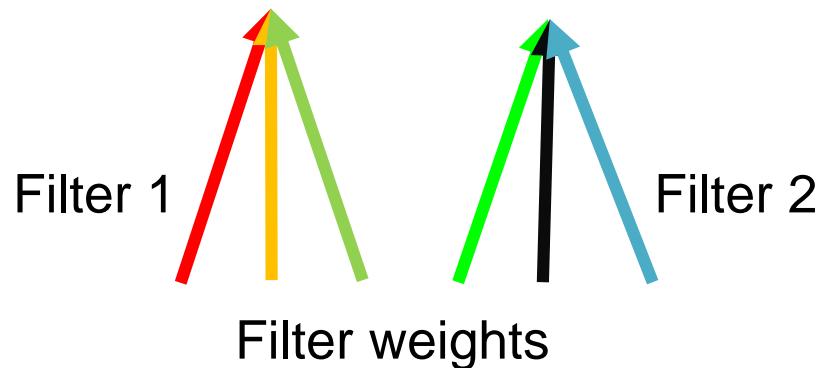
# CNN with multiple output maps



**Single** output map



**Multiple** output maps



# Putting them together

- Local connectivity
- Weight sharing
- Handling multiple input channels
- Handling multiple output maps

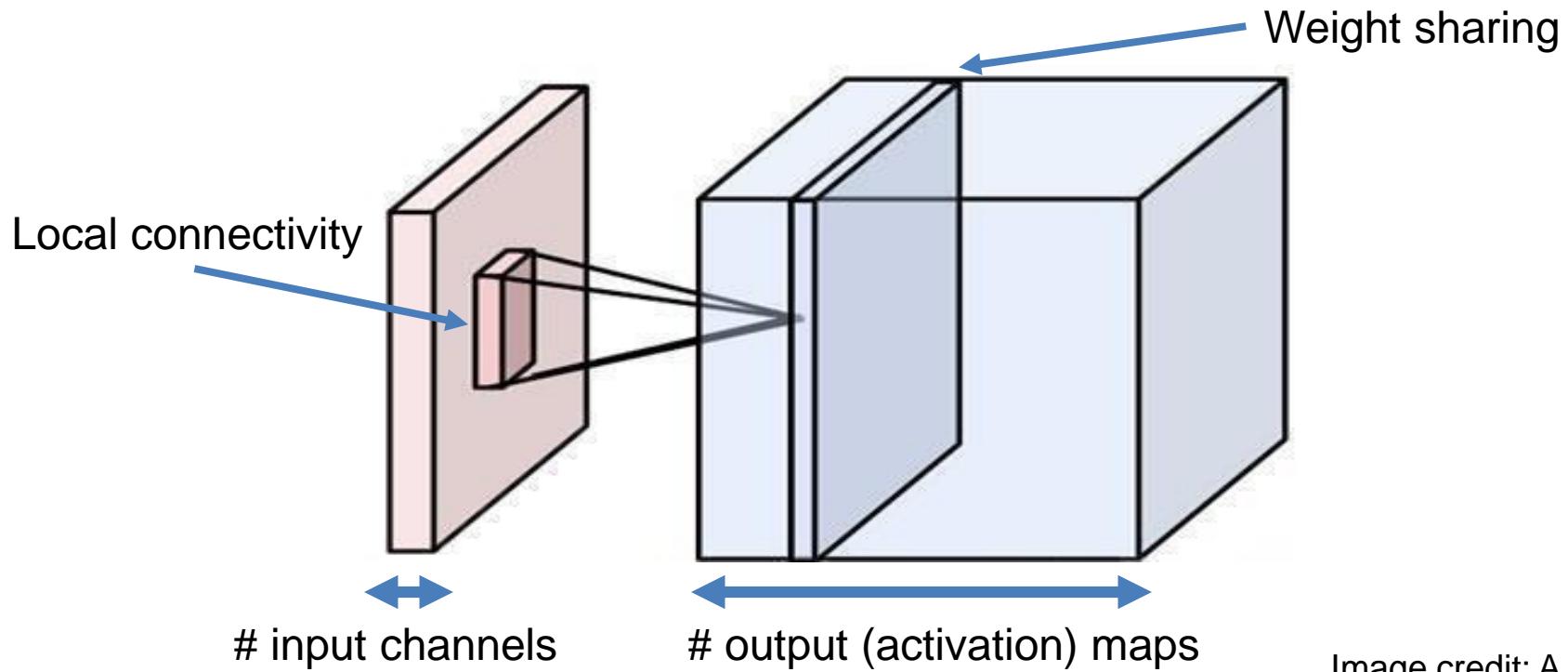


Image credit: A. Karpathy

# Neocognitron [Fukushima, Biological Cybernetics 1980]

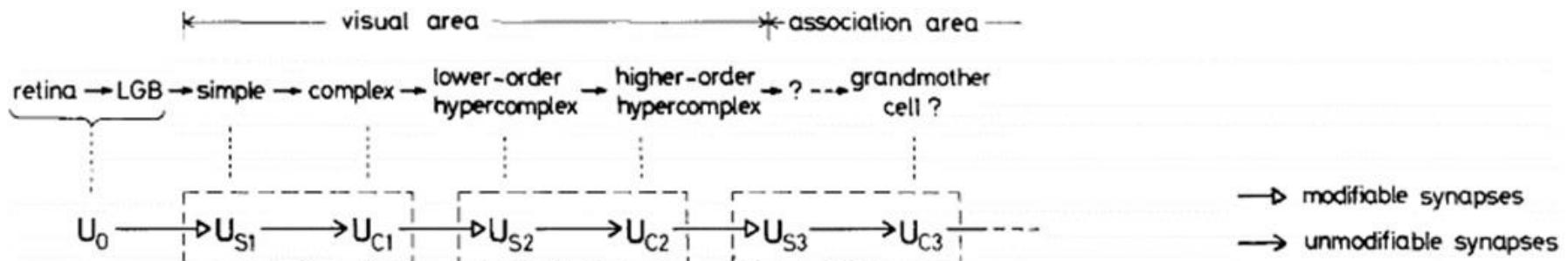
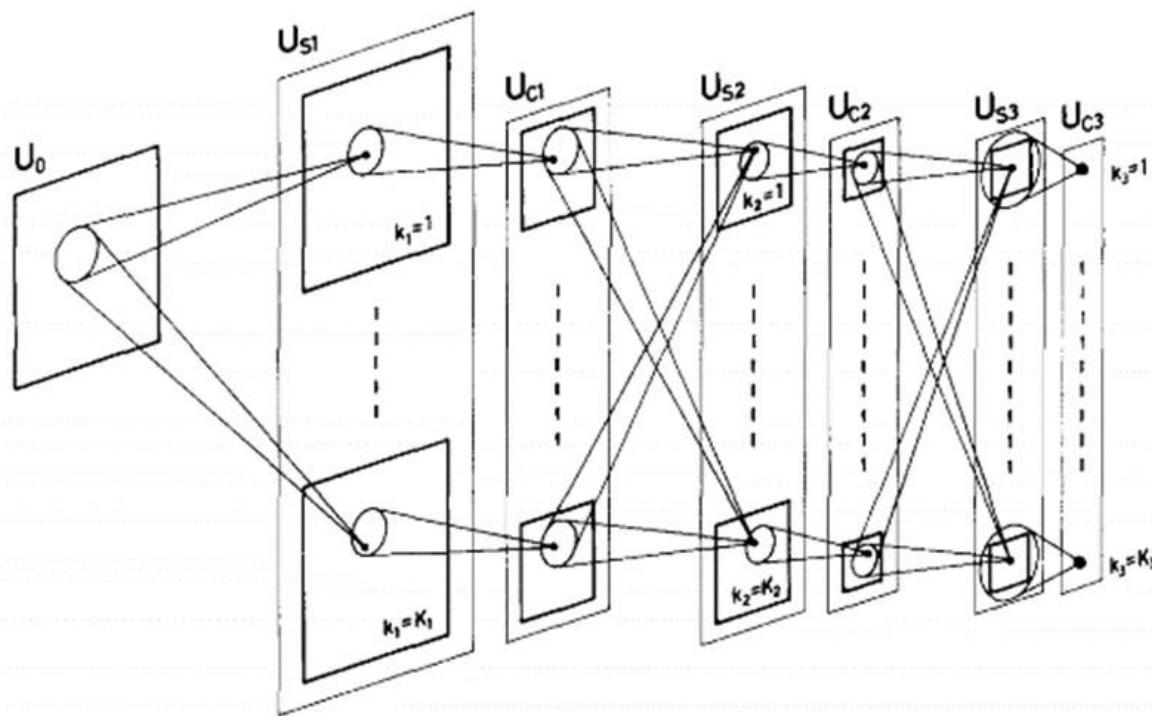


Fig. 1. Correspondence between the hierarchy model by Hubel and Wiesel, and the neural network of the neocognitron

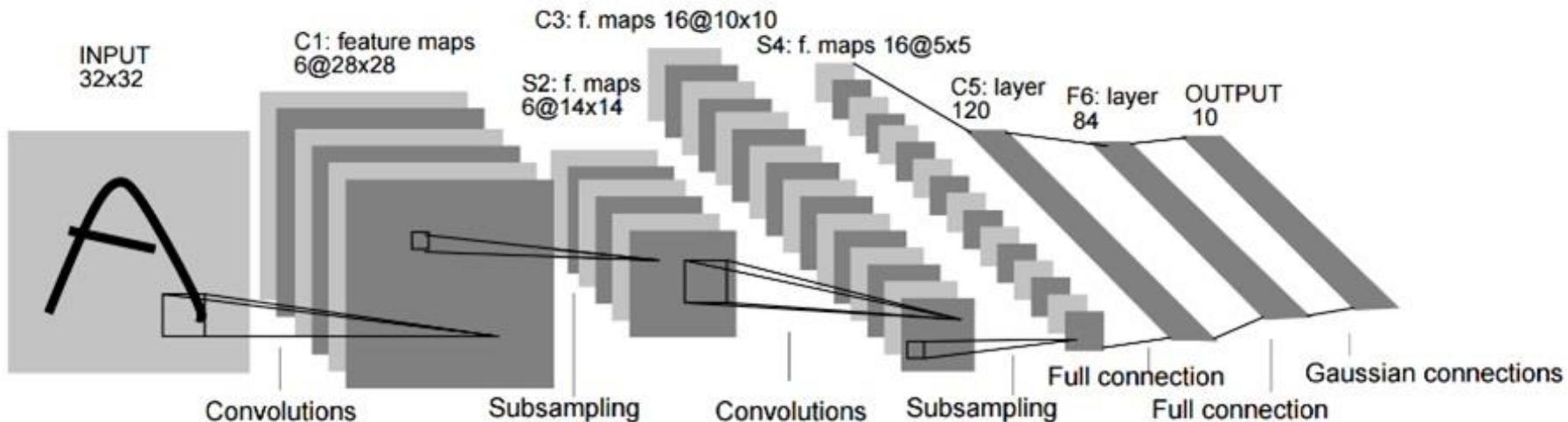


Deformation-Resistant  
Recognition

S-cells: (simple)  
- extract local features

C-cells: (complex)  
- allow for positional errors

# LeNet [LeCun et al. 1998]



Gradient-based learning applied to document  
recognition [[LeCun, Bottou, Bengio, Haffner 1998](#)]

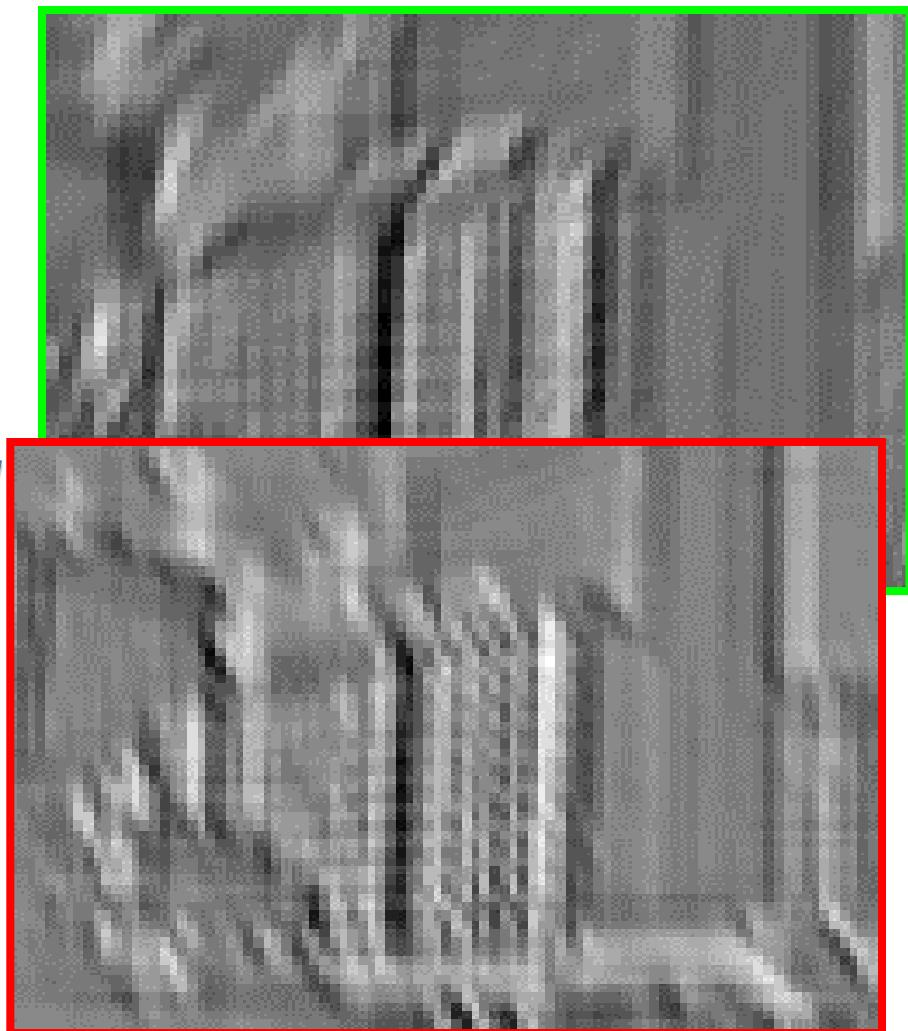
LeNet-1 from 1993

# What is a Convolution?

- Weighted moving sum



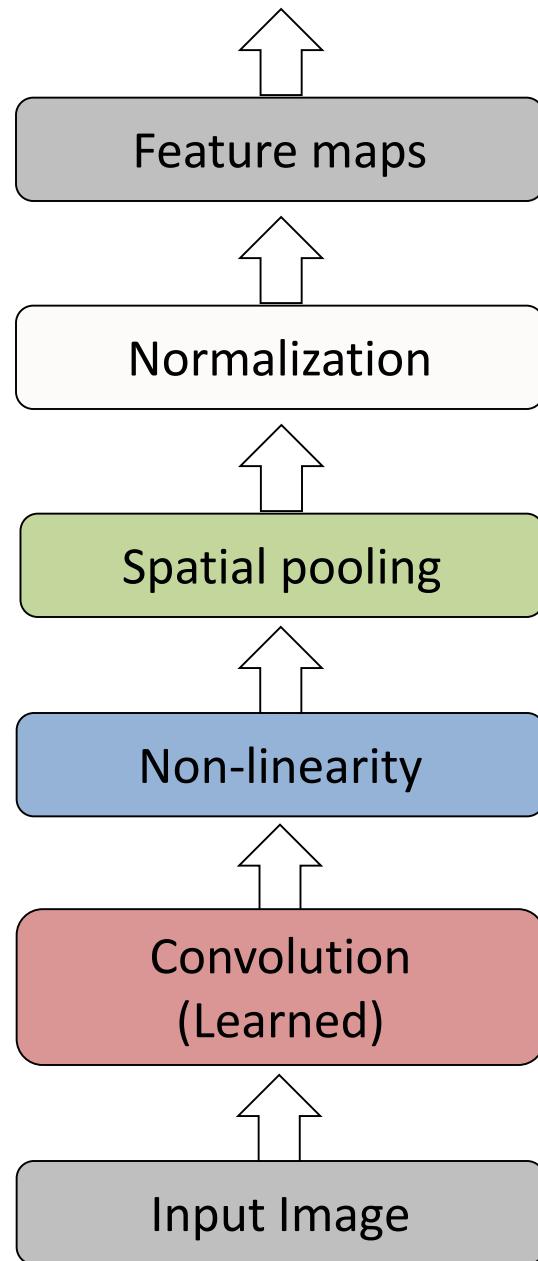
Input



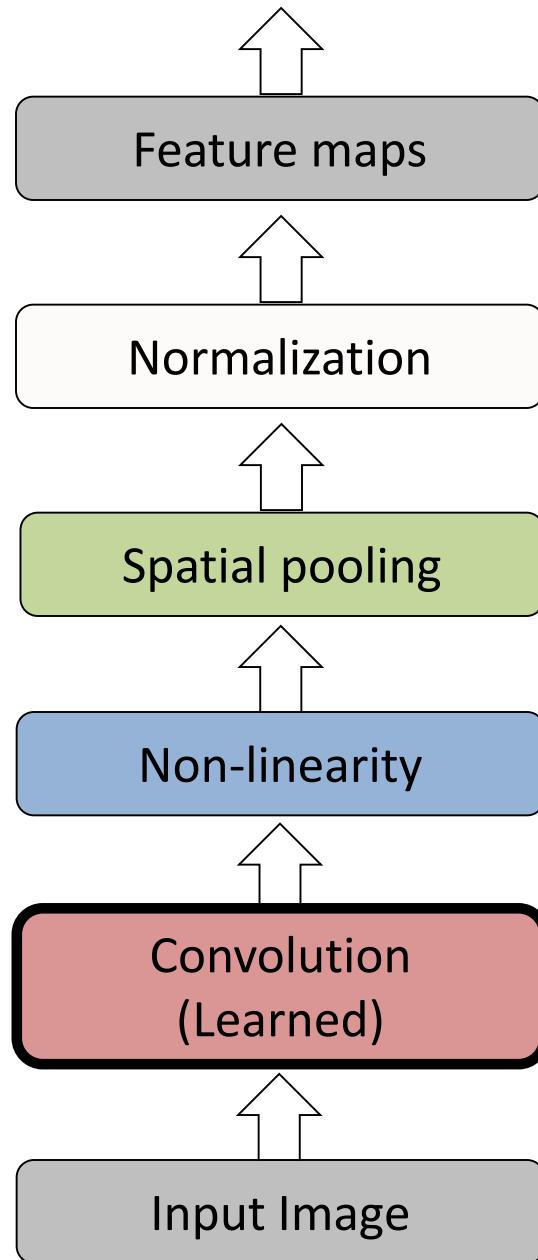
Feature Activation Map

slide credit: S. Lazebnik

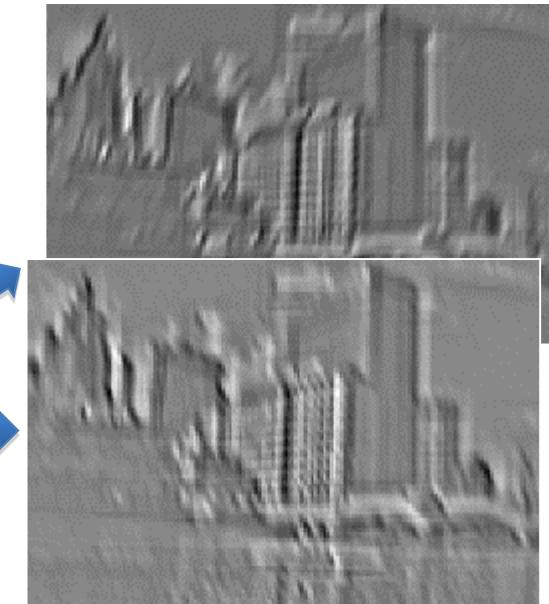
# Convolutional Neural Networks



# Convolutional Neural Networks



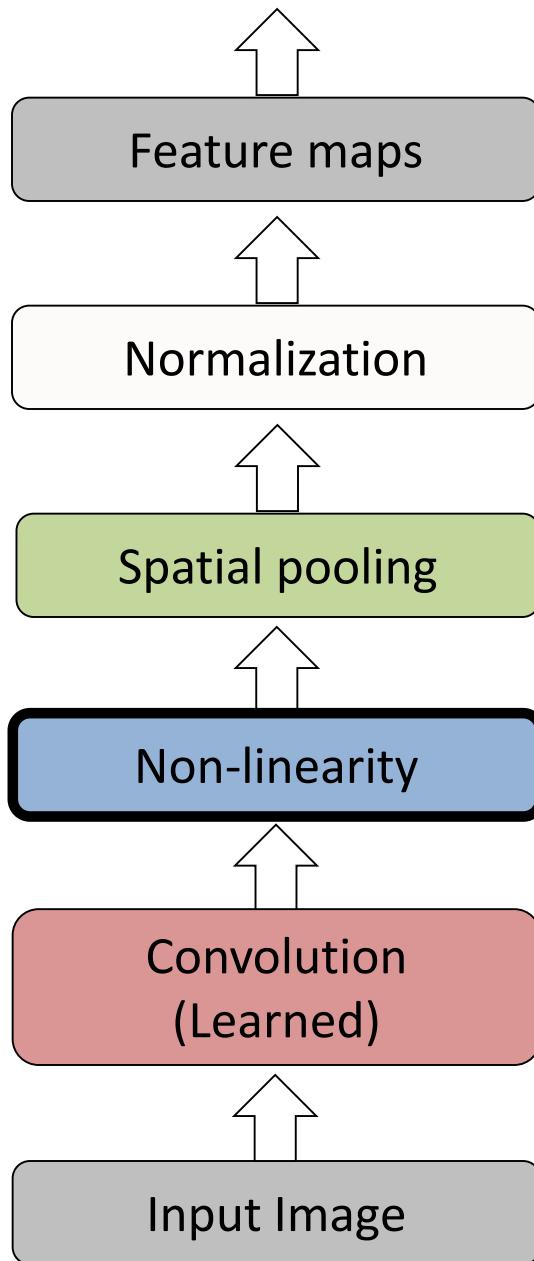
Input



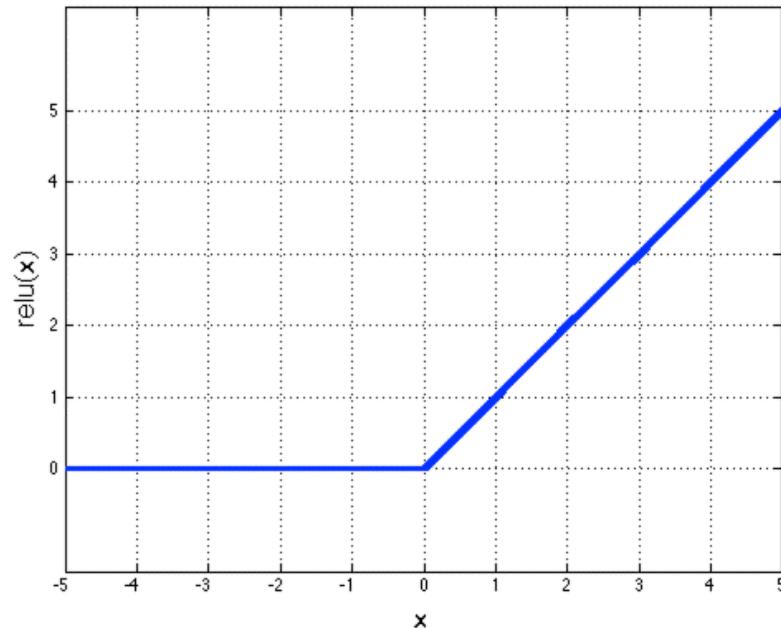
Feature Map

slide credit: S. Lazebnik

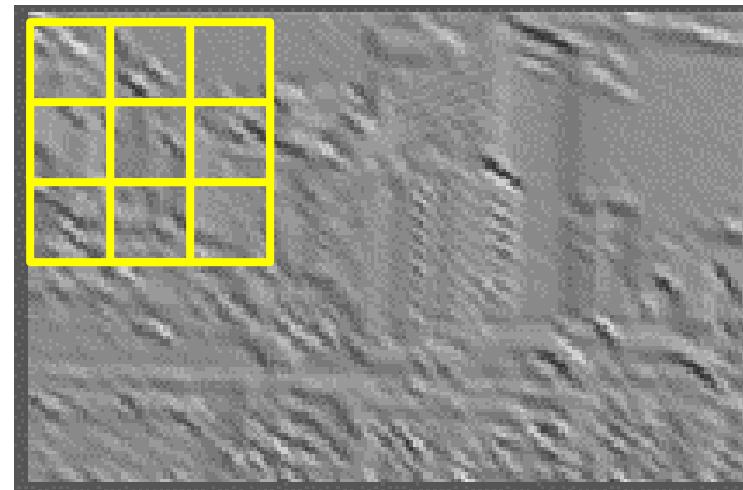
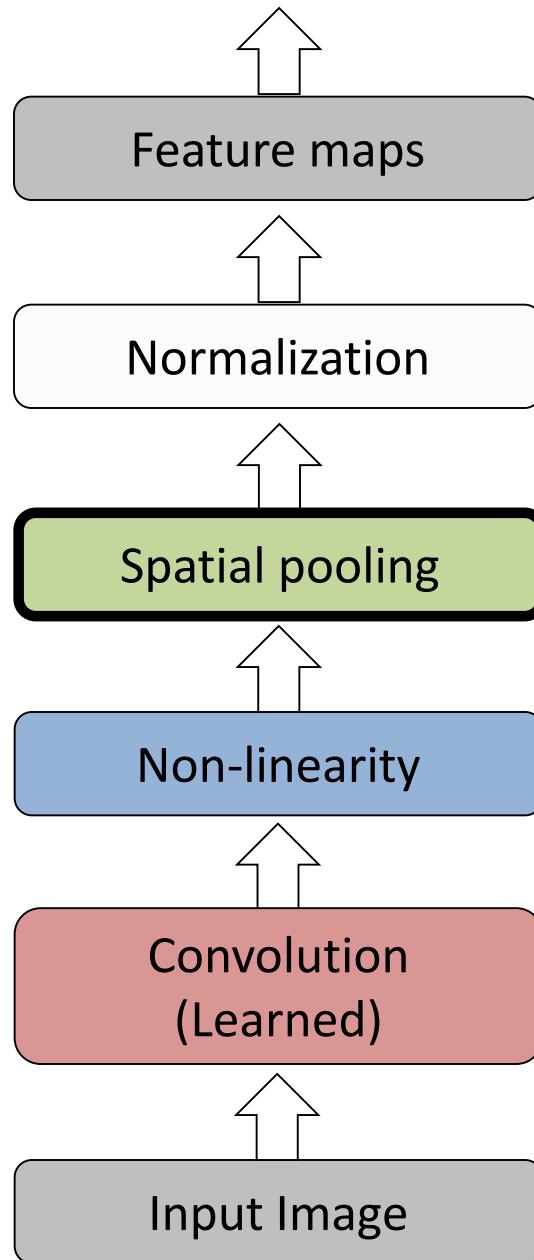
# Convolutional Neural Networks



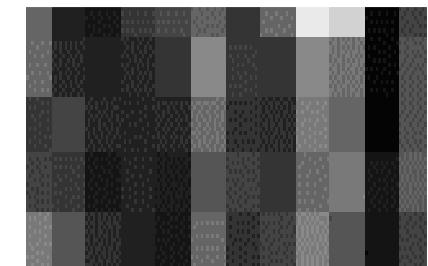
Rectified Linear Unit (ReLU)



# Convolutional Neural Networks



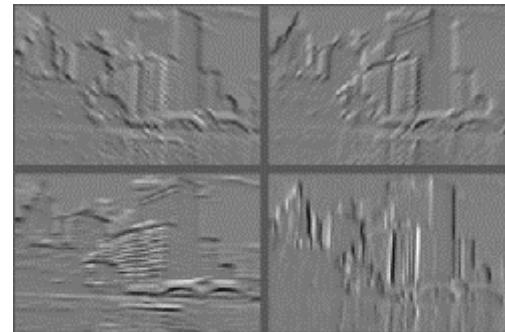
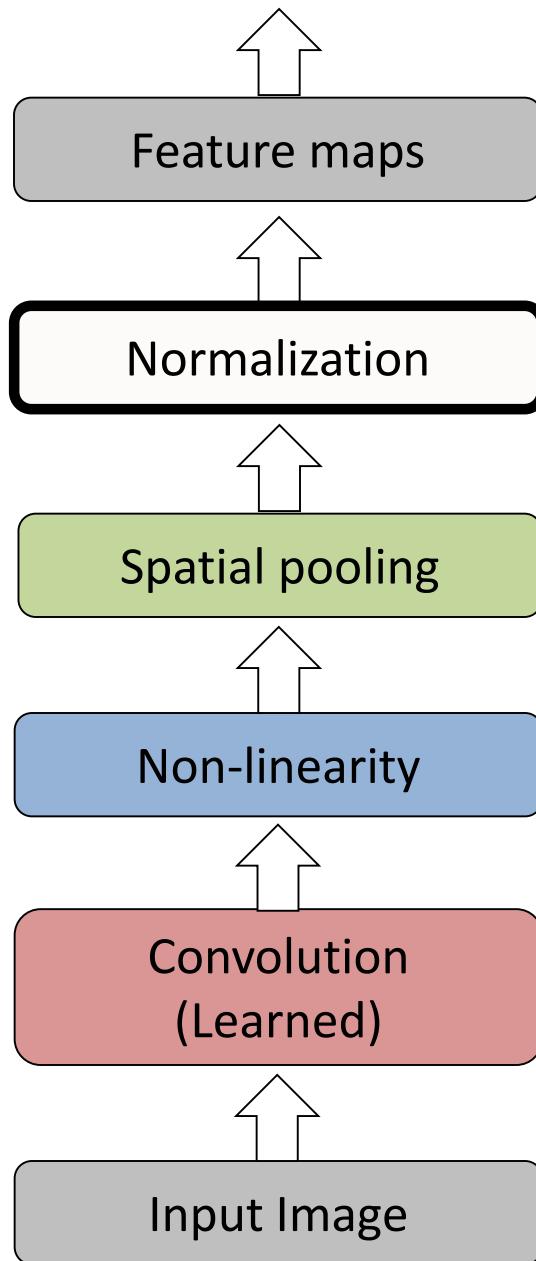
Max pooling



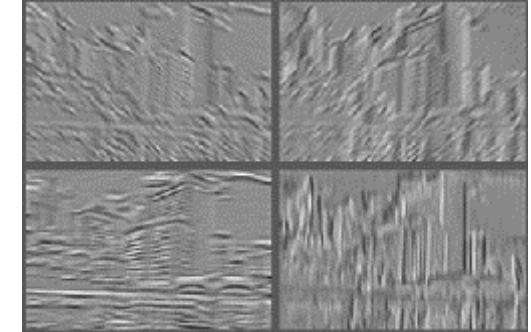
Max-pooling: a non-linear down-sampling

Provide *translation invariance*

# Convolutional Neural Networks

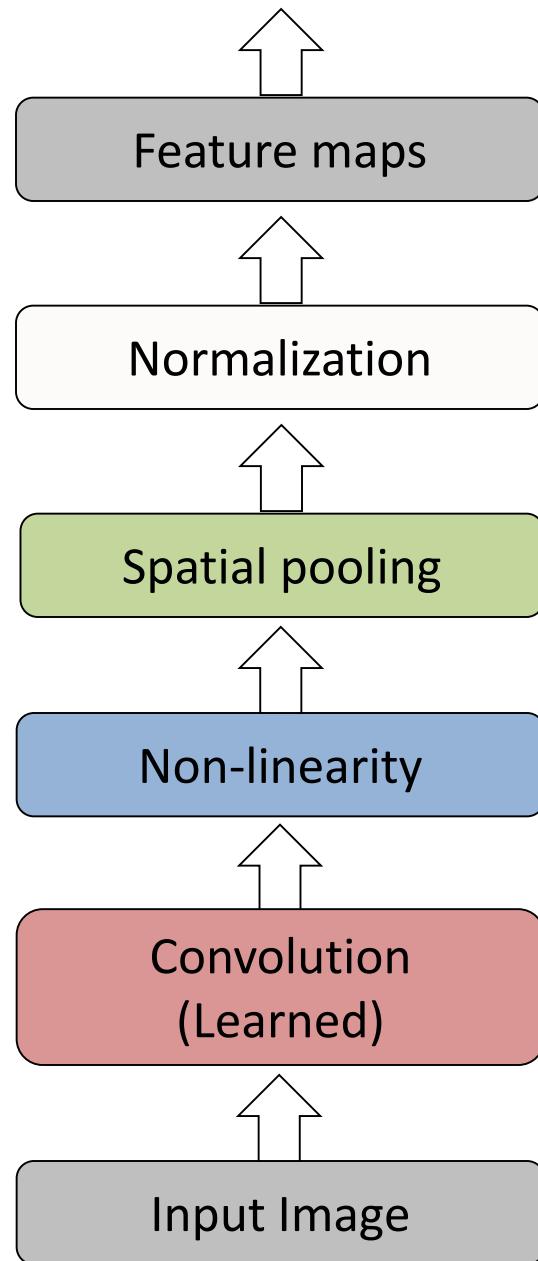


Feature Maps



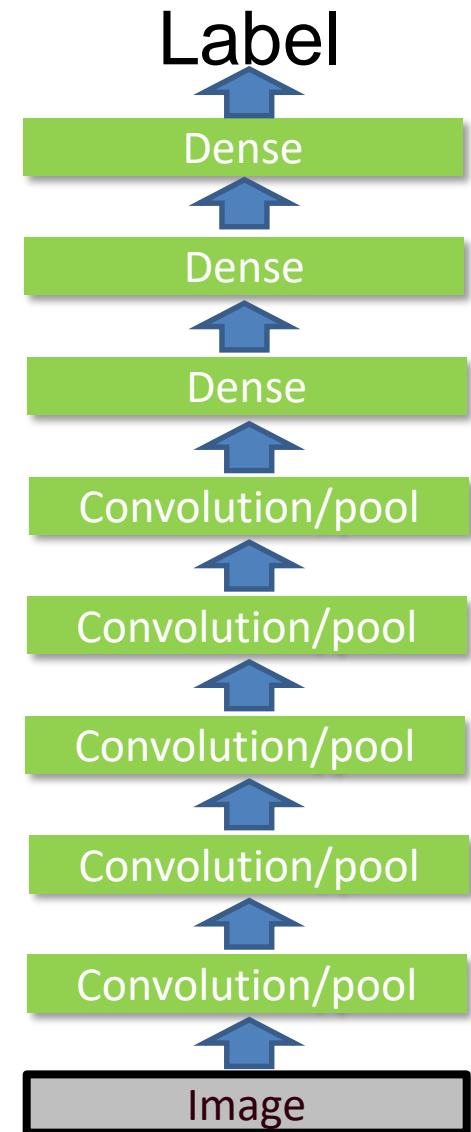
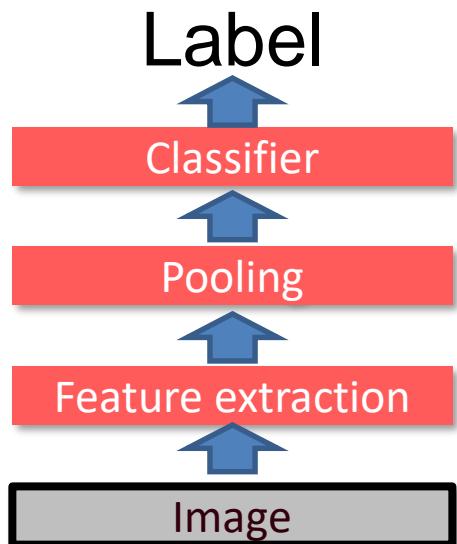
Feature Maps  
After Contrast  
Normalization

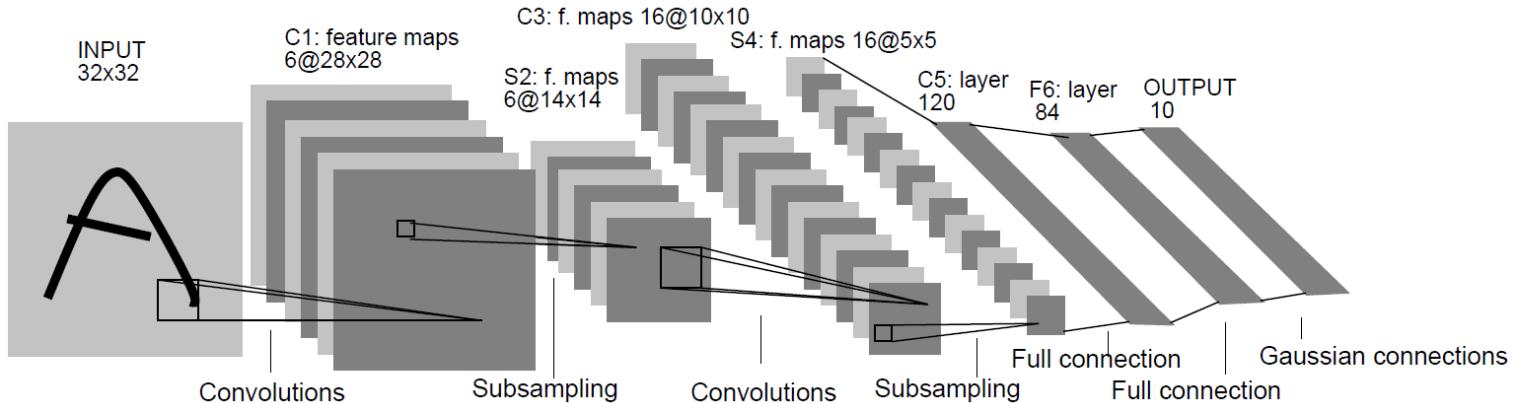
# Convolutional Neural Networks



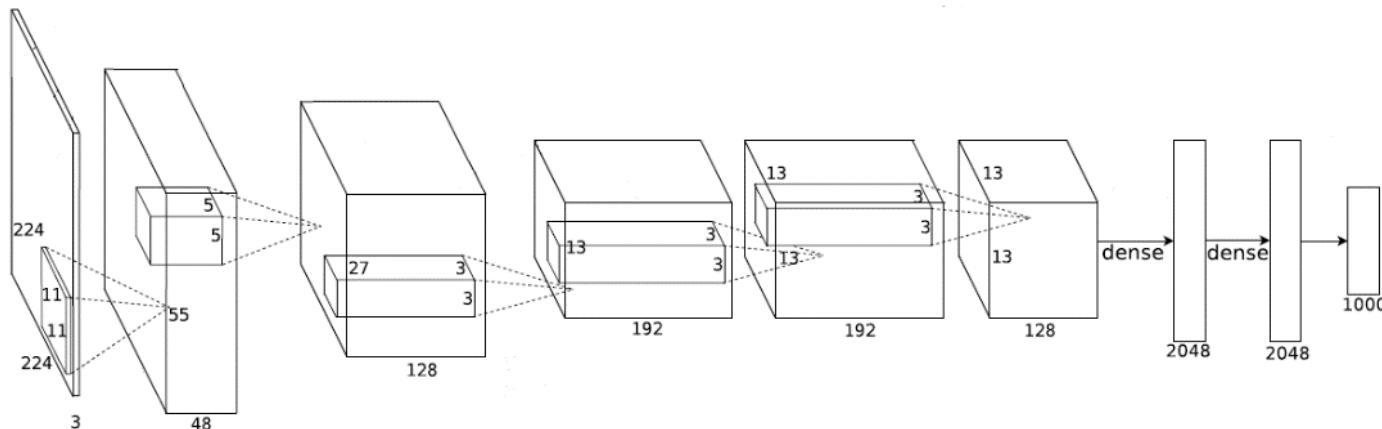
# Engineered vs. learned features

Convolutional filters are trained in a supervised manner by back-propagating classification error



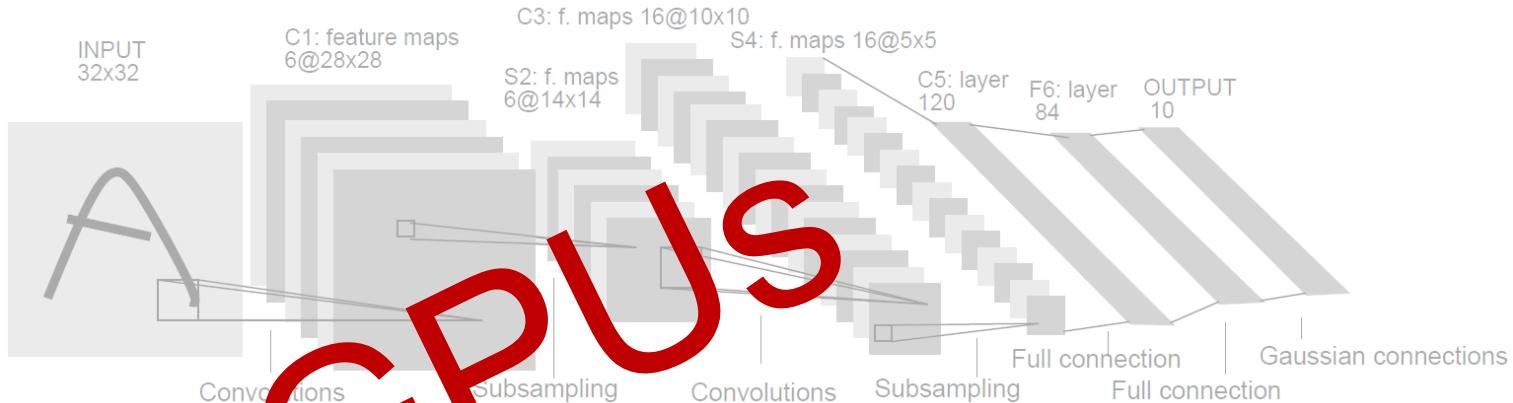


**Gradient-Based Learning Applied to Document Recognition**, LeCun, Bottou, Bengio and Haffner, Proc. of the IEEE, **1998**



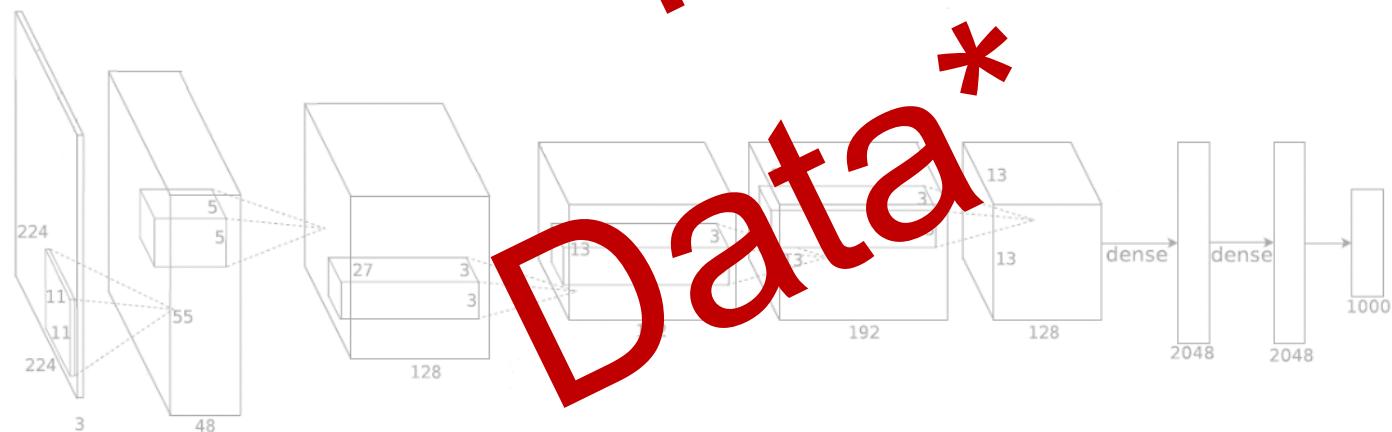
**Imagenet Classification with Deep Convolutional Neural Networks**, Krizhevsky, Sutskever, and Hinton, NIPS **2012**

Slide Credit: L. Zitnick



Gradient-Based Learning Applied to Document  
Recognition, LeCun, Bottou, Bengio and Haffner, Proc. of  
the IEEE, 1998

+



Imagenet Class  
Networks, Kriz

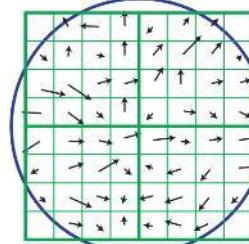
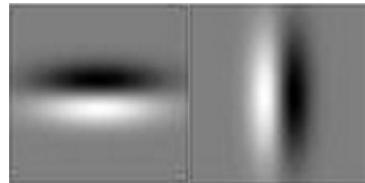
\* Rectified activations and dropout

Slide Credit: L. Zitnick

# SIFT Descriptor

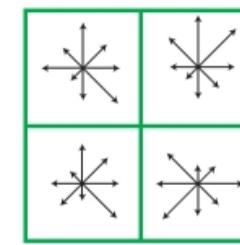
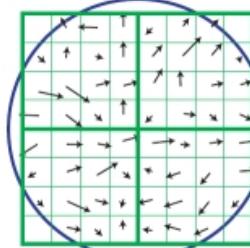
Image  
Pixels

Apply gradient  
filters

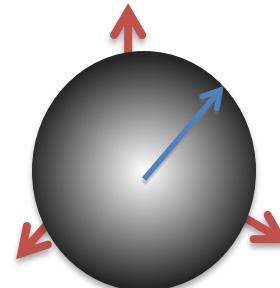


Lowe [IJCV 2004]

Spatial pool  
(Sum)



Normalize to unit  
length

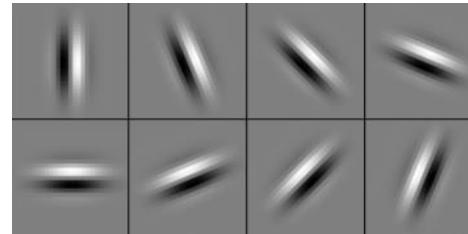


Feature  
Vector

# SIFT Descriptor

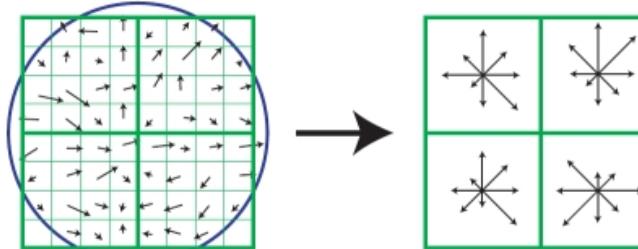
Image  
Pixels

Apply  
oriented filters

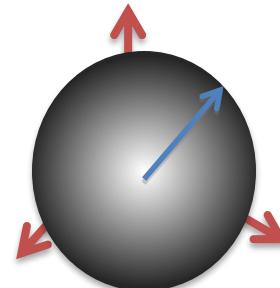


Lowe [IJCV 2004]

Spatial pool  
(Sum)



Normalize to unit  
length



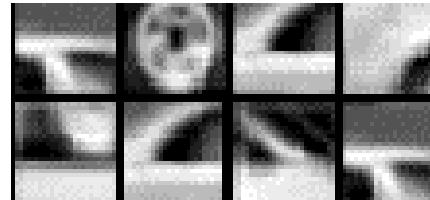
Feature  
Vector

slide credit: R. Fergus

# Spatial Pyramid Matching

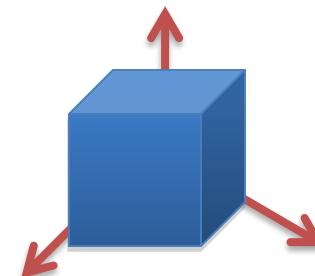
SIFT  
Features

Filter with  
Visual Words



Lazebnik,  
Schmid,  
Ponce  
[CVPR 2006]

Max



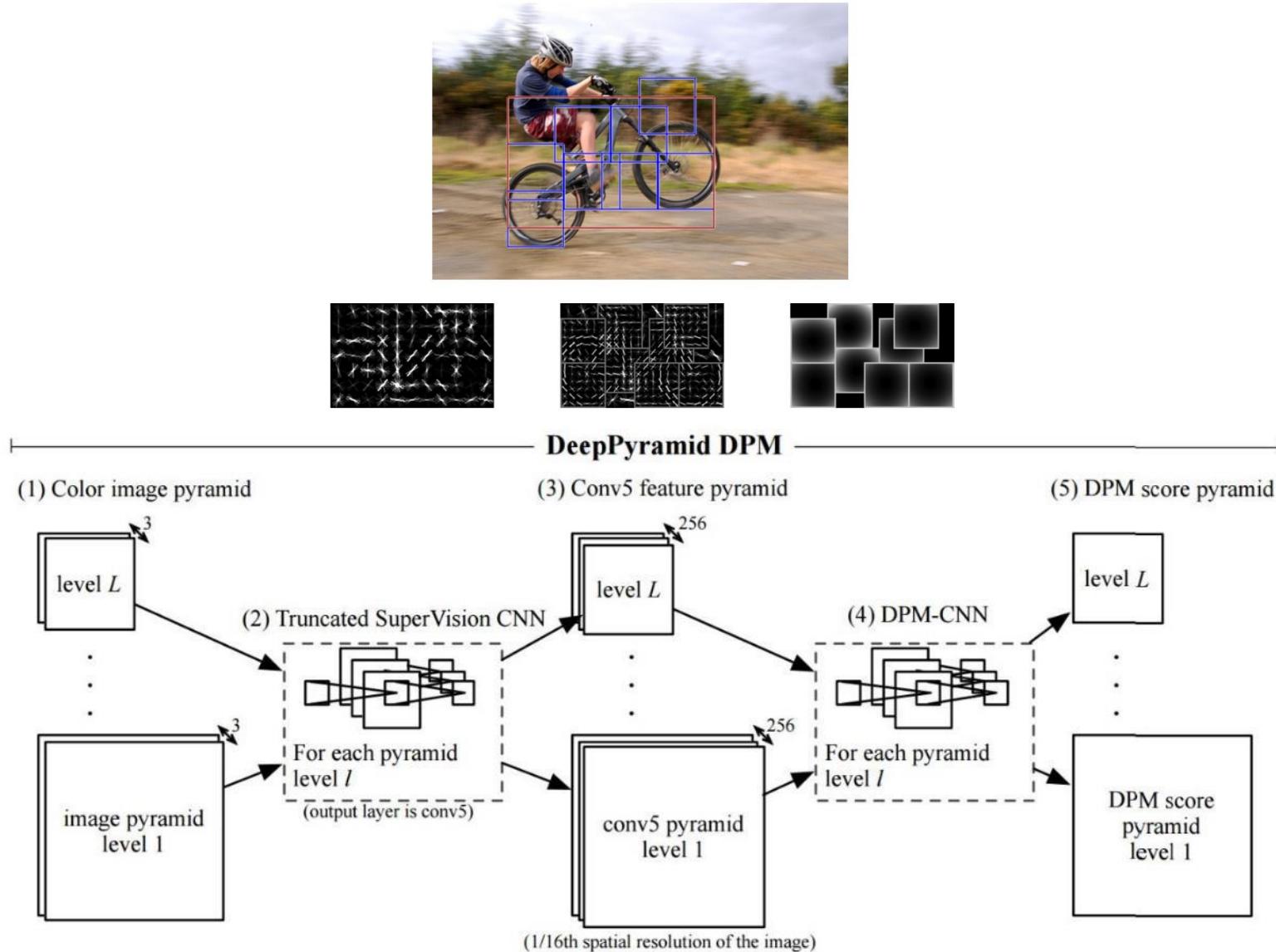
Multi-scale  
spatial pool  
(Sum)



Classifier

slide credit: R. Fergus

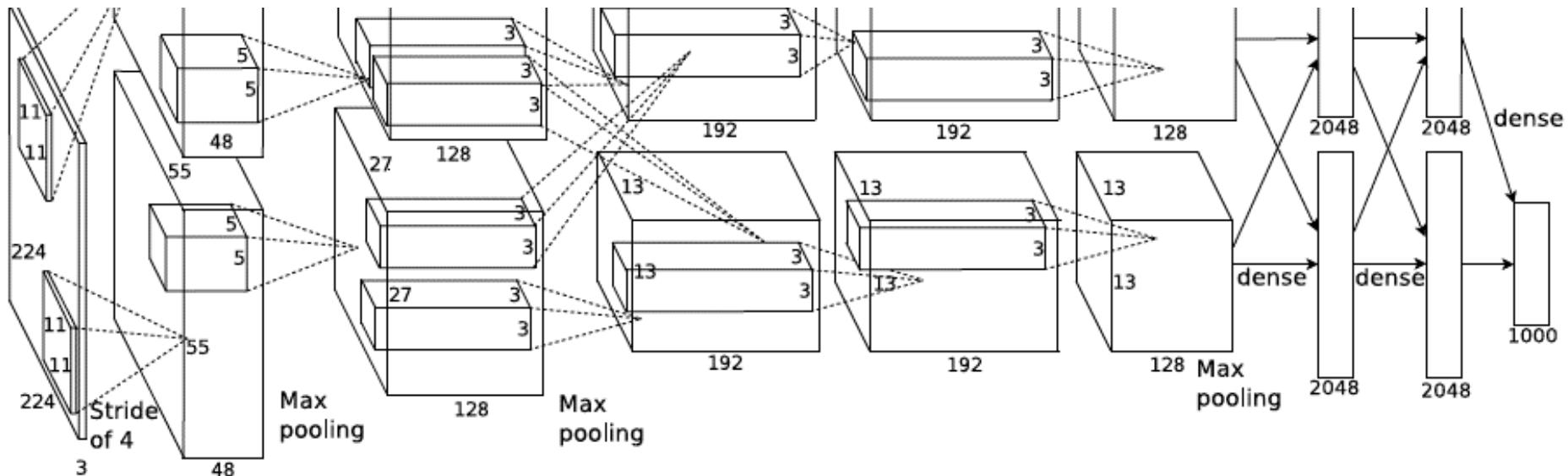
# Deformable Part Model



Deformable Part Models are Convolutional Neural Networks [Girshick et al. CVPR 15]

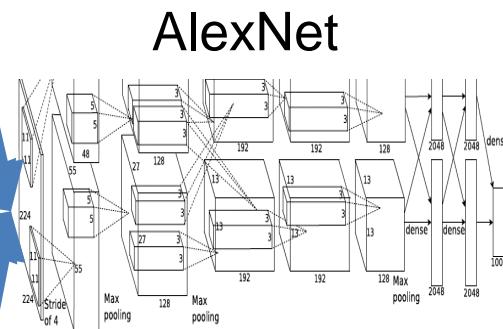
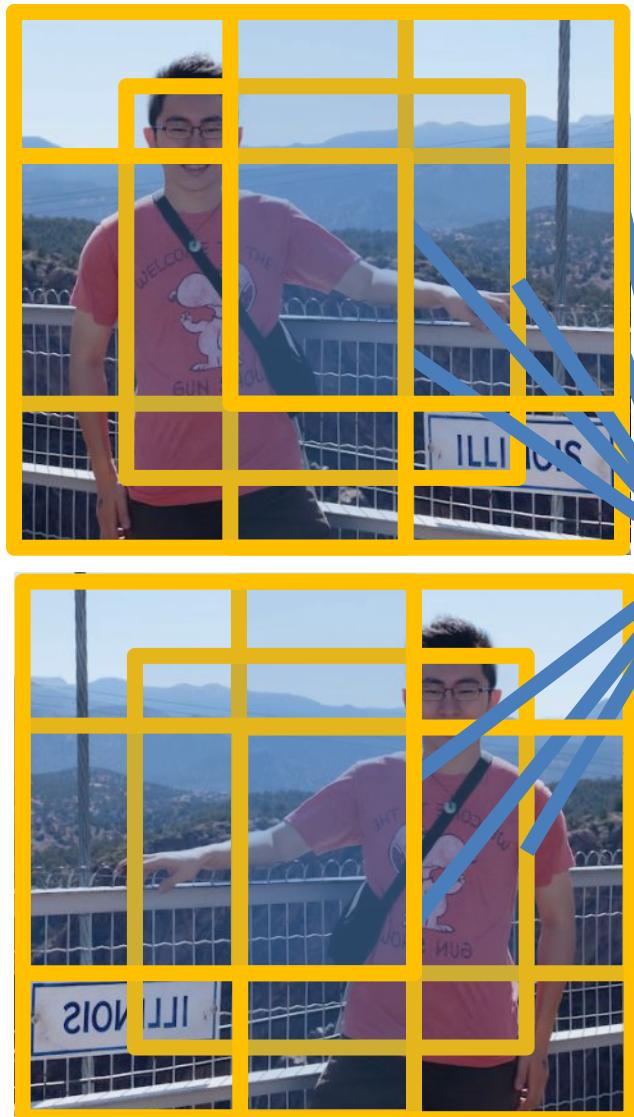
# AlexNet

- Similar framework to LeCun'98 but:
  - Bigger model (7 hidden layers, 650,000 units, 60,000,000 params)
  - More data ( $10^6$  vs.  $10^3$  images)
  - GPU implementation (50x speedup over CPU)
    - Trained on two GPUs for a week



A. Krizhevsky, I. Sutskever, and G. Hinton,  
ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012

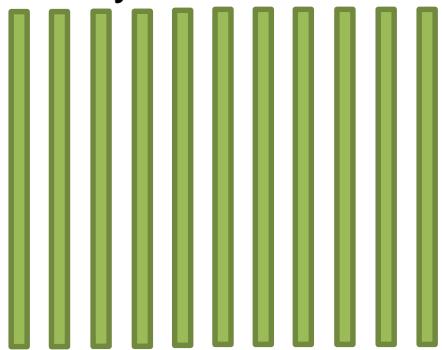
# Using CNN for Image Classification



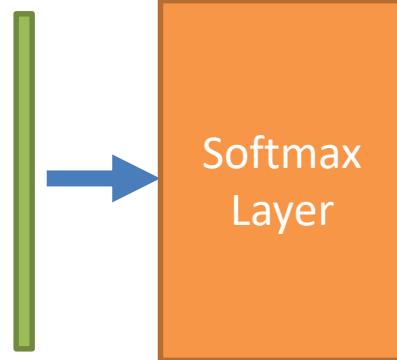
Fixed input size:  
224x224x3

$d = 4096$

Fully connected layer Fc7  
 $d = 4096$



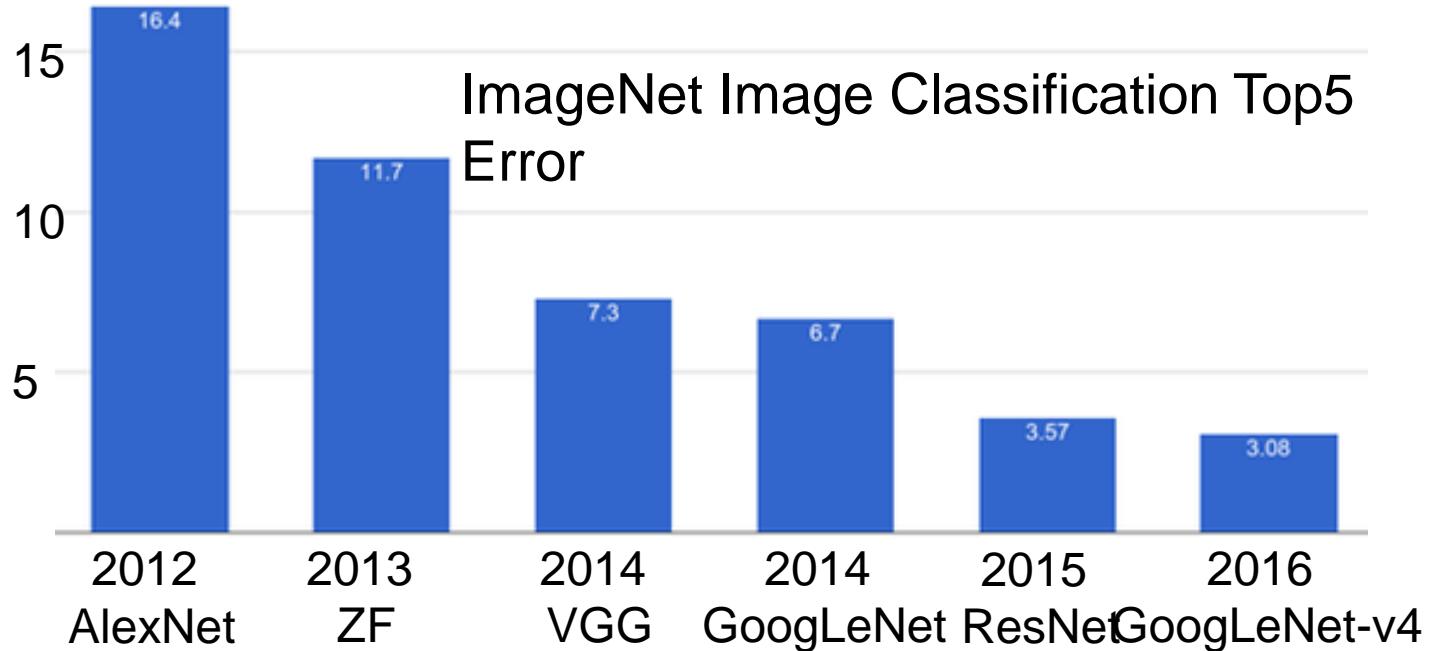
Averaging



Softmax  
Layer

"Jia-Bin"

# Progress on ImageNet



# VGG-Net

- The deeper, the better
- Key design choices:
  - 3x3 conv. Kernels
    - very small
  - conv. stride 1
    - no loss of information
- Other details:
  - Rectification (ReLU) non-linearity
  - 5 max-pool layers (x2 reduction)
  - no normalization
  - 3 fully-connected (FC) layers

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

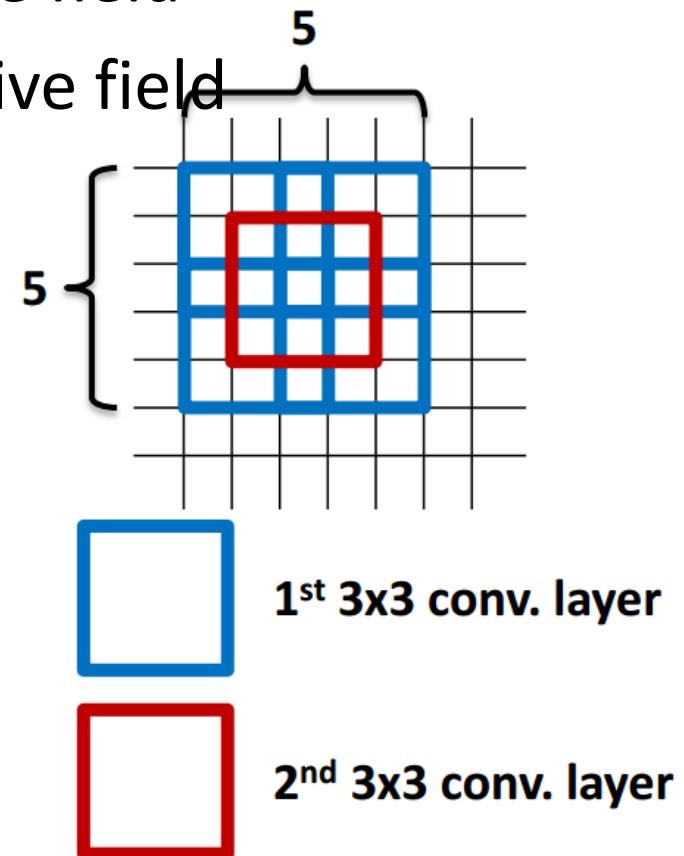
FC-4096

FC-1000

softmax

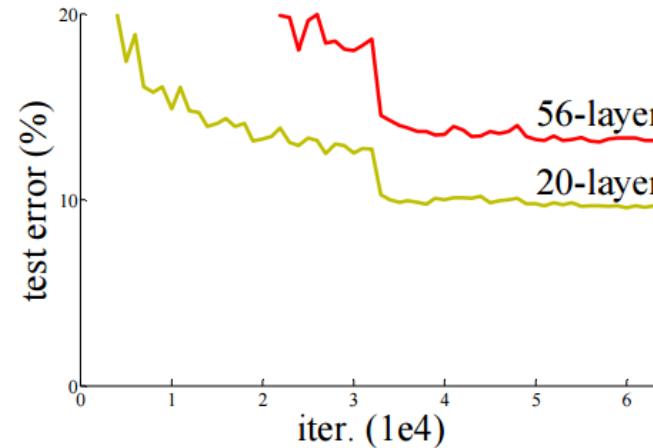
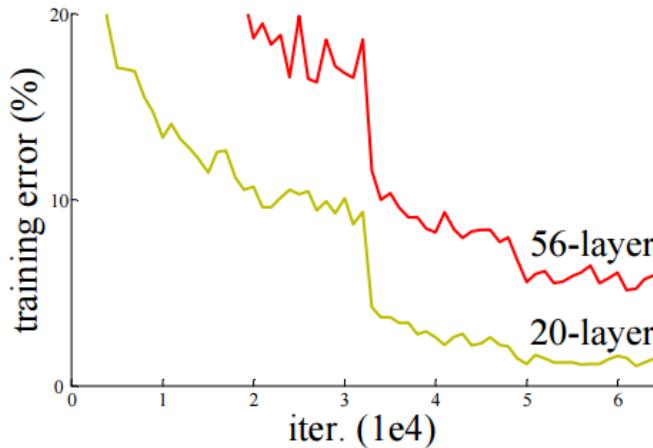
# VGG-Net

- Why 3x3 layers?
  - Stacked conv. layers have a large receptive field
  - two 3x3 layers – 5x5 receptive field
  - three 3x3 layers – 7x7 receptive field
- More non-linearity
  - Less parameters to learn
  - ~140M per net

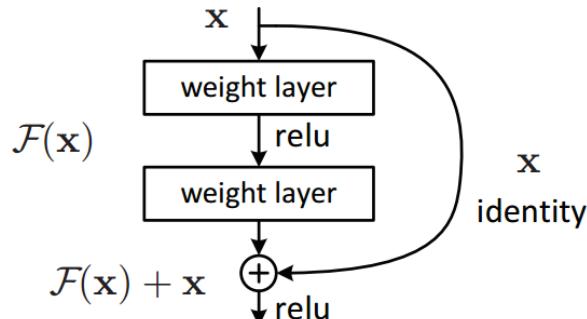


# ResNet

- Can we just increase the #layer?



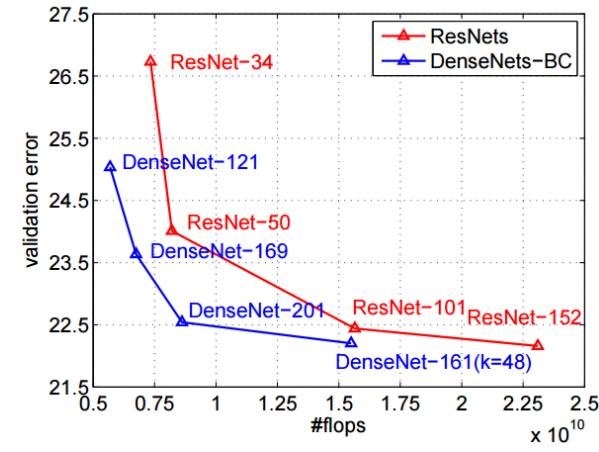
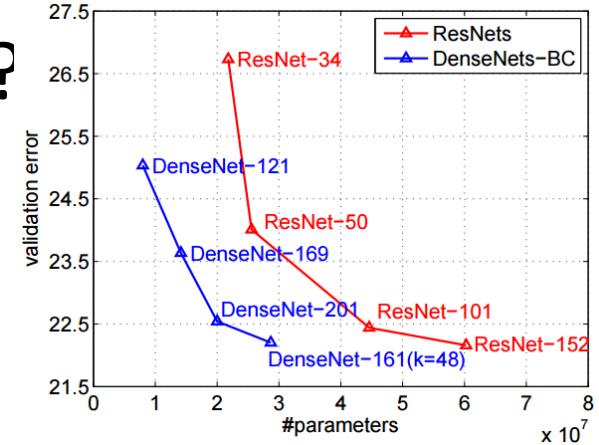
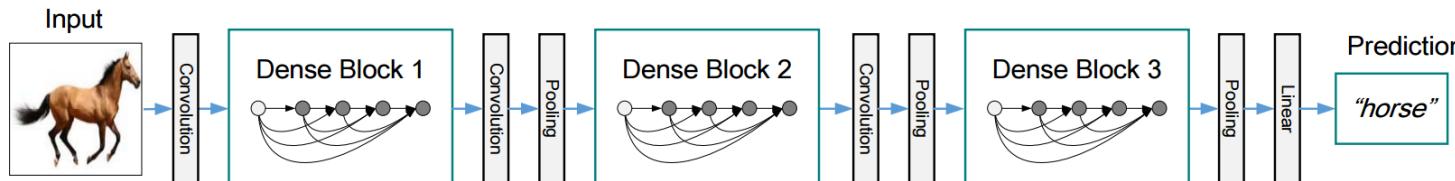
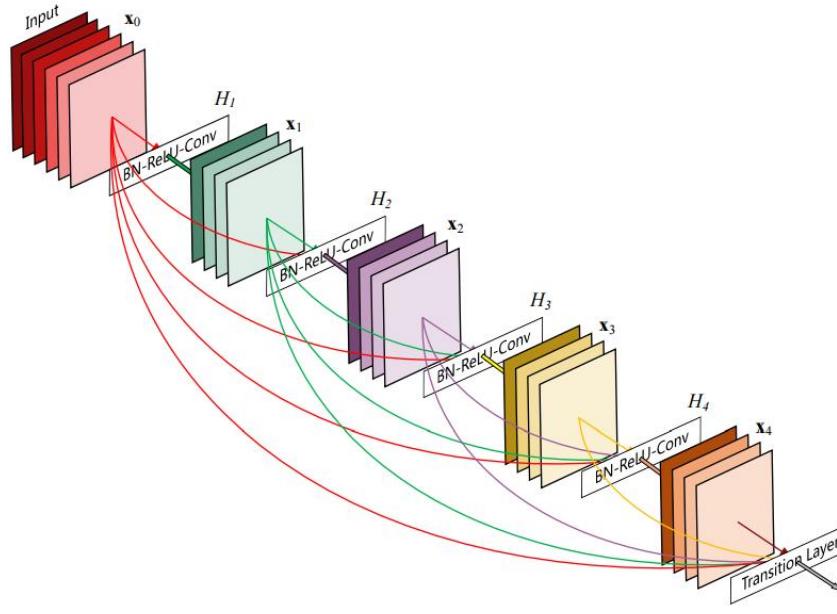
- How can we train very deep network?
  - Residual learning



method	top-5 err. (test)
VGG [41] (ILSVRC'14)	7.32
GoogLeNet [44] (ILSVRC'14)	6.66
VGG [41] (v5)	6.8
PReLU-net [13]	4.94
BN-inception [16]	4.82
<b>ResNet (ILSVRC'15)</b>	<b>3.57</b>

# DenseNet

- Shorter connections (like ResNet) help
- Why not just connect them all?



# Training Convolutional Neural Networks

- Backpropagation + stochastic gradient descent with momentum
  - [Neural Networks: Tricks of the Trade](#)
- Dropout
- Data augmentation
- Batch normalization
- Initialization
  - Transfer learning

# Training CNN with gradient descent

- A CNN as composition of functions

$$f_{\mathbf{w}}(\mathbf{x}) = f_L(\dots (f_2(f_1(\mathbf{x}; \mathbf{w}_1); \mathbf{w}_2) \dots; \mathbf{w}_L)$$

- Parameters

$$\mathbf{w} = (\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_L)$$

- Empirical loss function

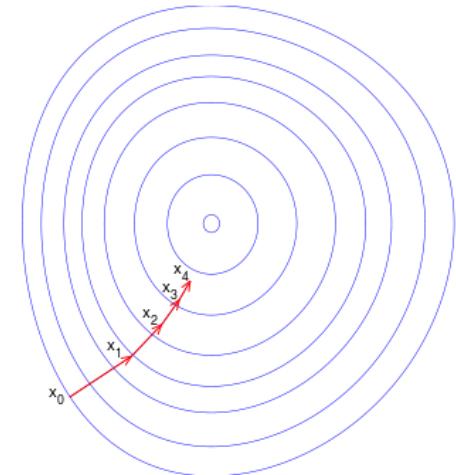
$$L(\mathbf{w}) = \frac{1}{n} \sum_i l(z_i, f_{\mathbf{w}}(\mathbf{x}_i))$$

- Gradient descent

$$\text{New weight} \rightarrow \mathbf{w}^{t+1} = \mathbf{w}^t - \eta_t \frac{\partial f}{\partial \mathbf{w}}(\mathbf{w}^t)$$

Diagram illustrating the gradient descent process:

- New weight**: Yellow box containing the next weight value.
- Old weight**: Red box containing the current weight value.
- Learning rate**: Green box containing the step size  $\eta_t$ .
- Gradient**: Blue box pointing to the derivative of the loss function with respect to the weight.



# An Illustrative example

$$f(x, y) = xy, \quad \frac{\partial f}{\partial x} = y, \frac{\partial f}{\partial y} = x$$

Example:  $x = 4, y = -3 \Rightarrow f(x, y) = -12$

Partial derivatives

$$\frac{\partial f}{\partial x} = -3, \quad \frac{\partial f}{\partial y} = 4$$

Gradient

$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y} \right]$$

$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$

$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

Goal: compute the gradient

$$\nabla f = \left[ \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z} \right]$$

$$f(x, y, z) = (x + y)z = qz$$

$$q = x + y$$

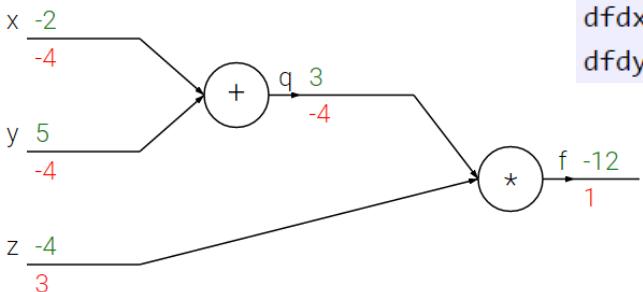
$$\frac{\partial q}{\partial x} = 1, \quad \frac{\partial q}{\partial y} = 1$$

$$f = qz$$

$$\frac{\partial f}{\partial q} = z, \quad \frac{\partial f}{\partial z} = q$$

## Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$



```

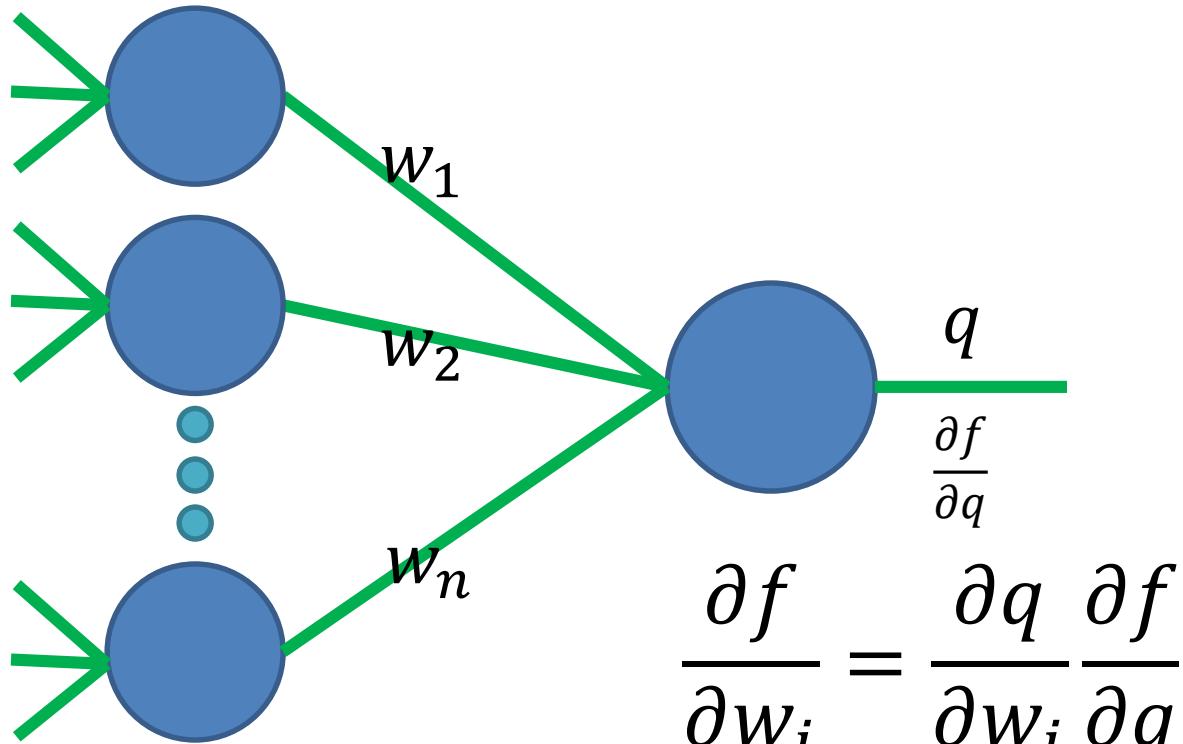
# set some inputs
x = -2; y = 5; z = -4

# perform the forward pass
q = x + y # q becomes 3
f = q * z # f becomes -12

# perform the backward pass (backpropagation) in reverse order:
# first backprop through f = q * z
dfdz = q # df/dz = q, so gradient on z becomes 3
dfdq = z # df/dq = z, so gradient on q becomes -4
# now backprop through q = x + y
dfdx = 1.0 * dfdq # dq/dx = 1. And the multiplication here is the chain rule!
dfdy = 1.0 * dfdq # dq/dy = 1
  
```

Example credit: Andrej Karpathy

# Backpropagation (recursive chain rule)



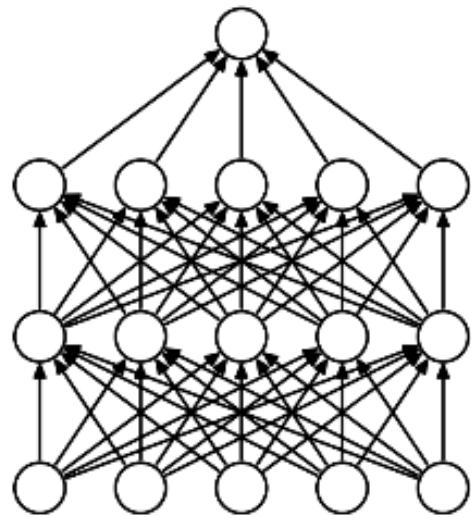
Local gradient

Gate gradient

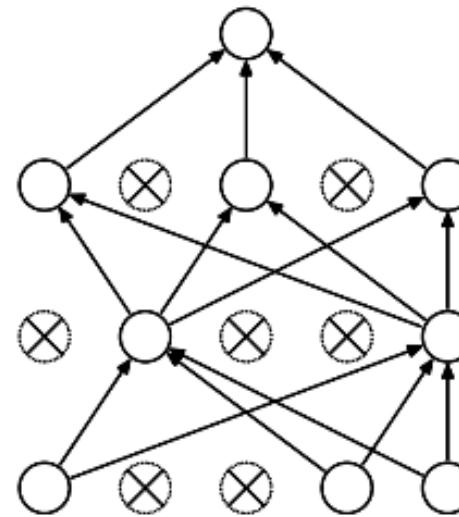
Can be computed during forward pass

The gate receives this during backprop

# Dropout



(a) Standard Neural Net

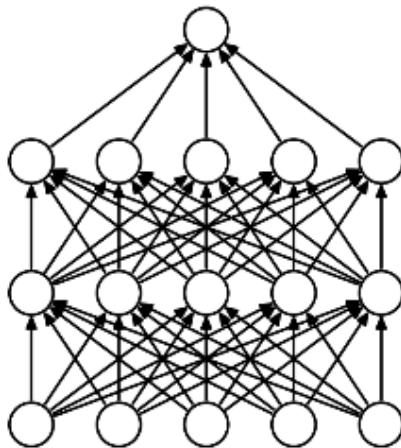


(b) After applying dropout.

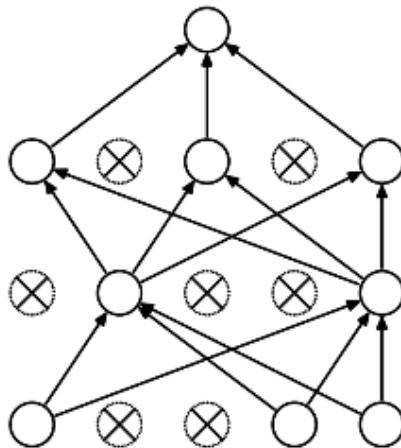
Intuition: successful conspiracies

- 50 people planning a conspiracy
- Strategy A: plan a big conspiracy involving 50 people
  - Likely to fail. 50 people need to play their parts correctly.
- Strategy B: plan 10 conspiracies each involving 5 people
  - Likely to succeed!

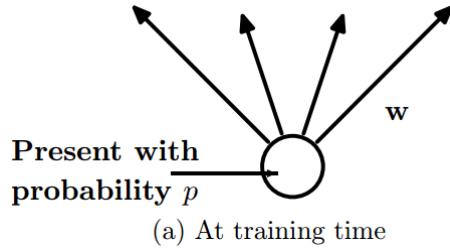
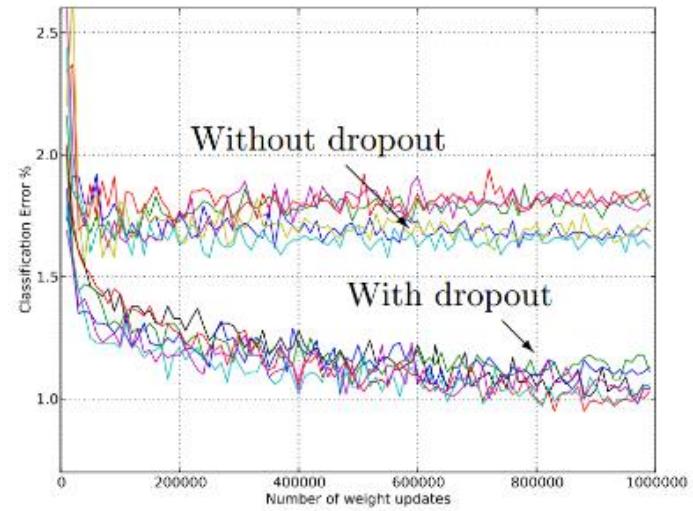
# Dropout



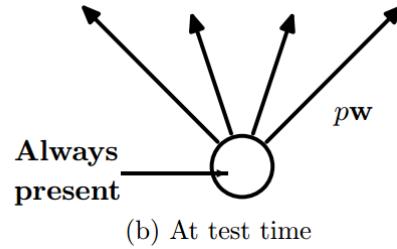
(a) Standard Neural Net



(b) After applying dropout.



(a) At training time



(b) At test time

**Main Idea:** approximately combining exponentially many different neural network architectures efficiently

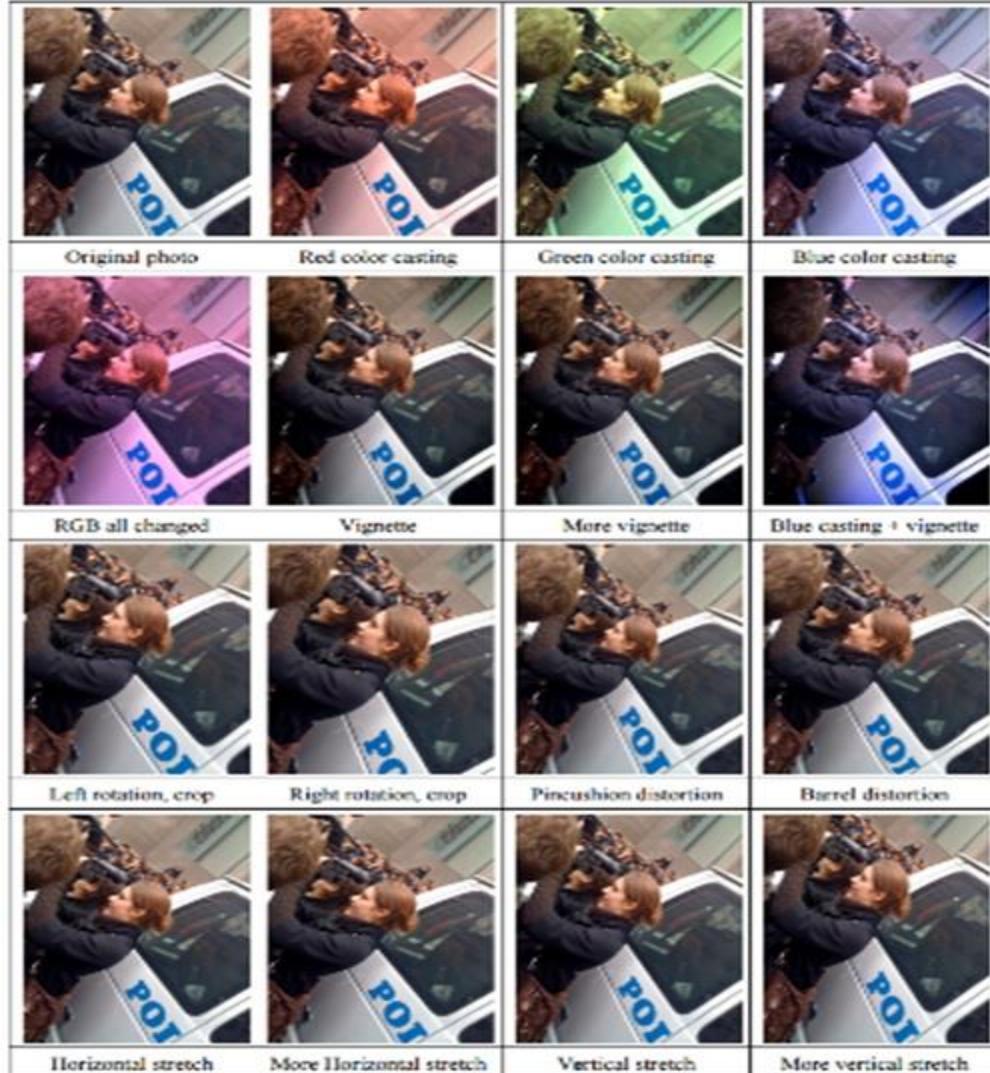
Model	Top-1 (val)	Top-5 (val)	Top-5 (test)
SVM on Fisher Vectors of Dense SIFT and Color Statistics	-	-	27.3
Avg of classifiers over FVs of SIFT, LBP, GIST and CSIFT	-	-	26.2
Conv Net + dropout (Krizhevsky et al., 2012)	40.7	18.2	-
Avg of 5 Conv Nets + dropout (Krizhevsky et al., 2012)	38.1	16.4	16.4

Table 6: Results on the ILSVRC-2012 validation/test set.

Dropout: A simple way to prevent neural networks from overfitting [Srivastava JMLR 2014]

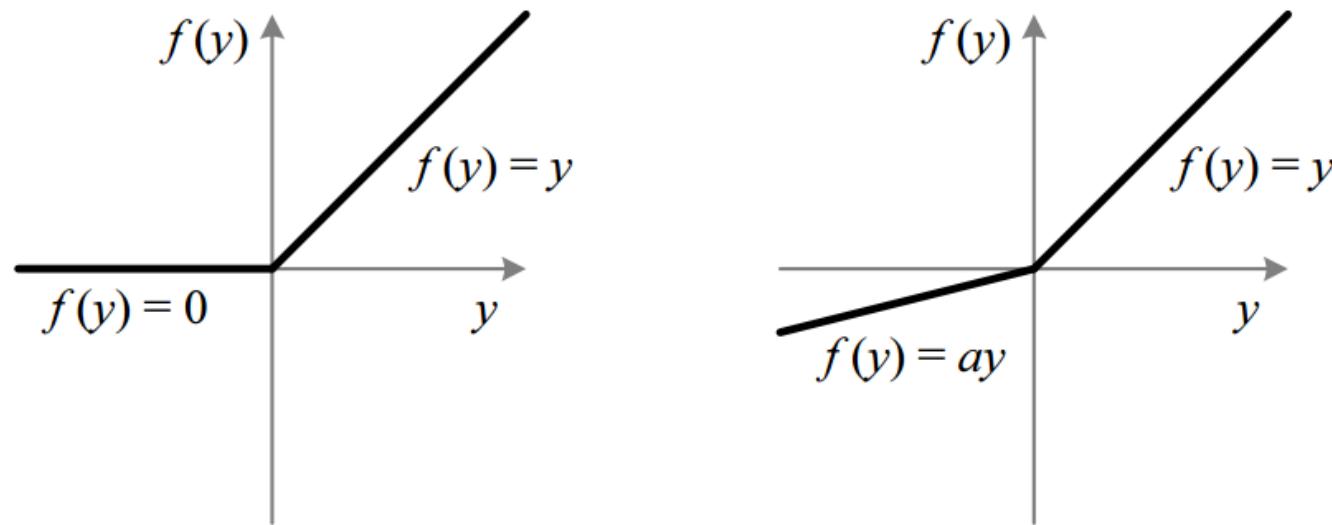
# Data Augmentation (Jittering)

- Create *virtual* training samples
  - Horizontal flip
  - Random crop
  - Color casting
  - Geometric distortion



Deep Image [Wu et al. 2015]

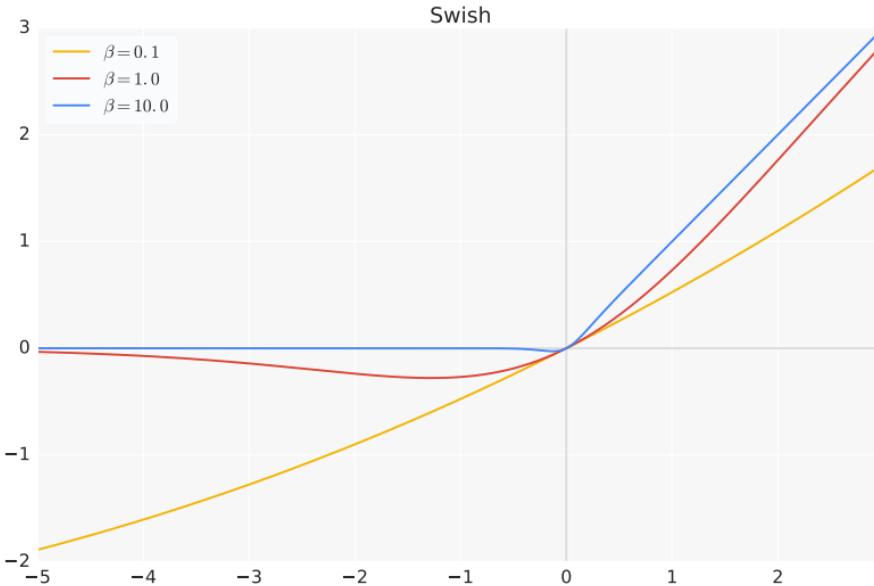
# Parametric Rectified Linear Unit



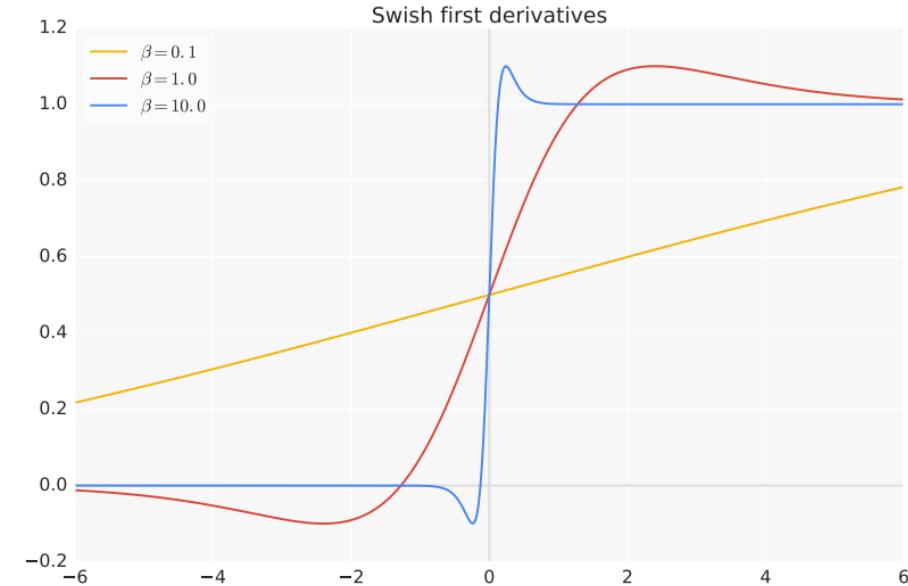
	team	top-5 ( <b>test</b> )
in competition ILSVRC 14	MSRA, SPP-nets [11]	8.06
	VGG [25]	7.32
	GoogLeNet [29]	6.66
post-competition	VGG [25] (arXiv v5)	6.8
	Baidu [32]	5.98
	<b>MSRA, PReLU-nets</b>	<b>4.94</b>

Delving Deep into Rectifiers: Surpassing Human-Level Performance on  
ImageNet Classification [[He et al. 2015](#)]

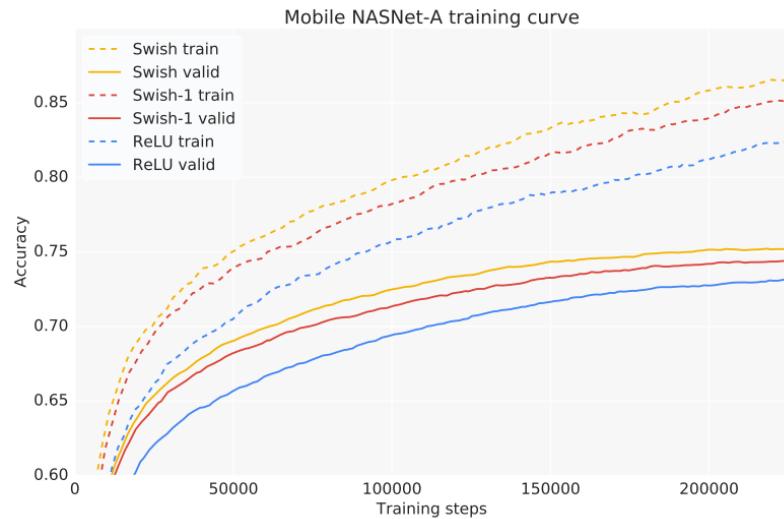
# Swish



The Swish activation function



First derivatives of Swish



Model	Top-1 Acc. (%)			Top-5 Acc. (%)		
LReLU	73.8	73.9	74.2	91.6	91.9	91.9
PReLU	74.6	74.7	74.7	92.4	92.3	92.3
Softplus	74.0	74.2	74.2	91.6	91.8	91.9
ELU	74.1	74.2	74.2	91.8	91.8	91.8
SELU	73.6	73.7	73.7	91.6	91.7	91.7
GELU	74.6	-	-	92.0	-	-
ReLU	73.5	73.6	73.8	91.4	91.5	91.6
Swish-1	74.6	74.7	74.7	92.1	92.0	92.0
Swish	74.9	74.9	75.2	92.3	92.4	92.4

# Batch Normalization

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;

Parameters to be learned:  $\gamma, \beta$

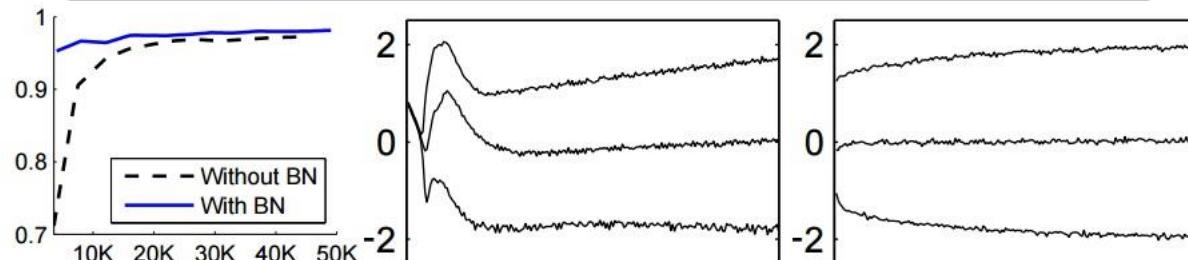
**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$



(a)

(b) Without BN

(c) With BN

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift [[Ioffe and Szegedy 2015](#)]

# Understanding and Visualizing CNN

- Find images that maximize some class scores
- Individual neuron activation
- Breaking CNNs

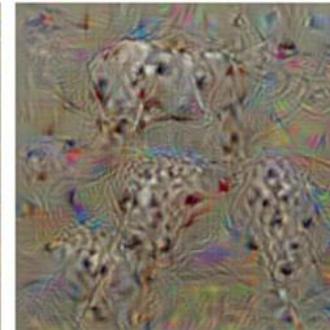
# Find images that maximize some class scores



dumbbell



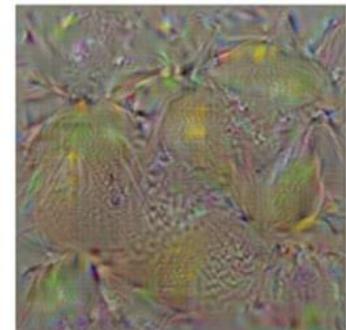
cup



dalmatian



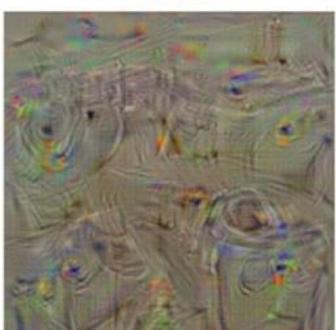
bell pepper



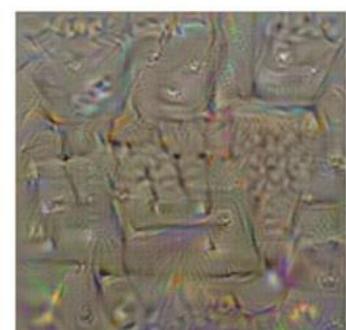
lemon



husky



washing machine



computer keyboard

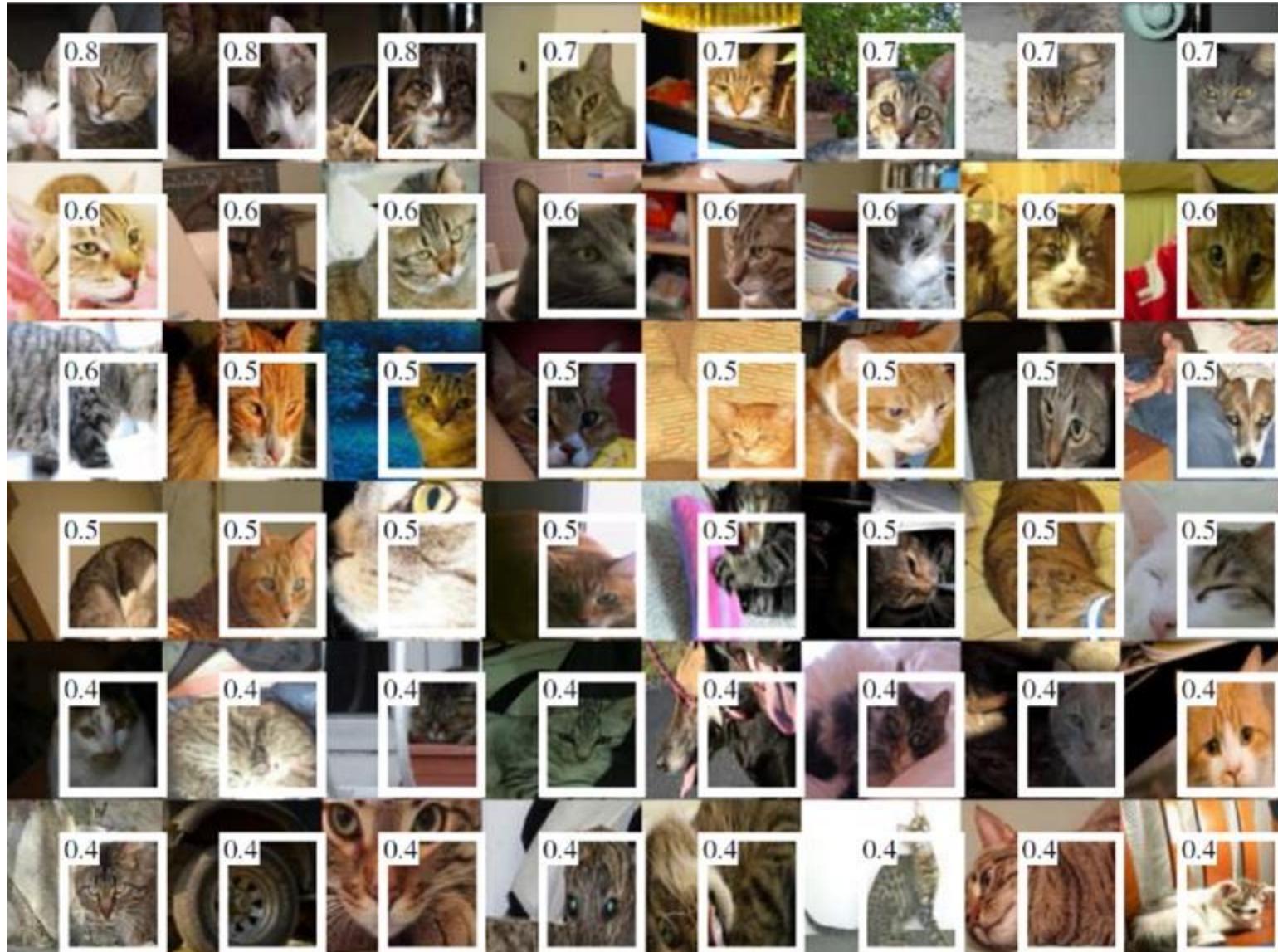


kit fox



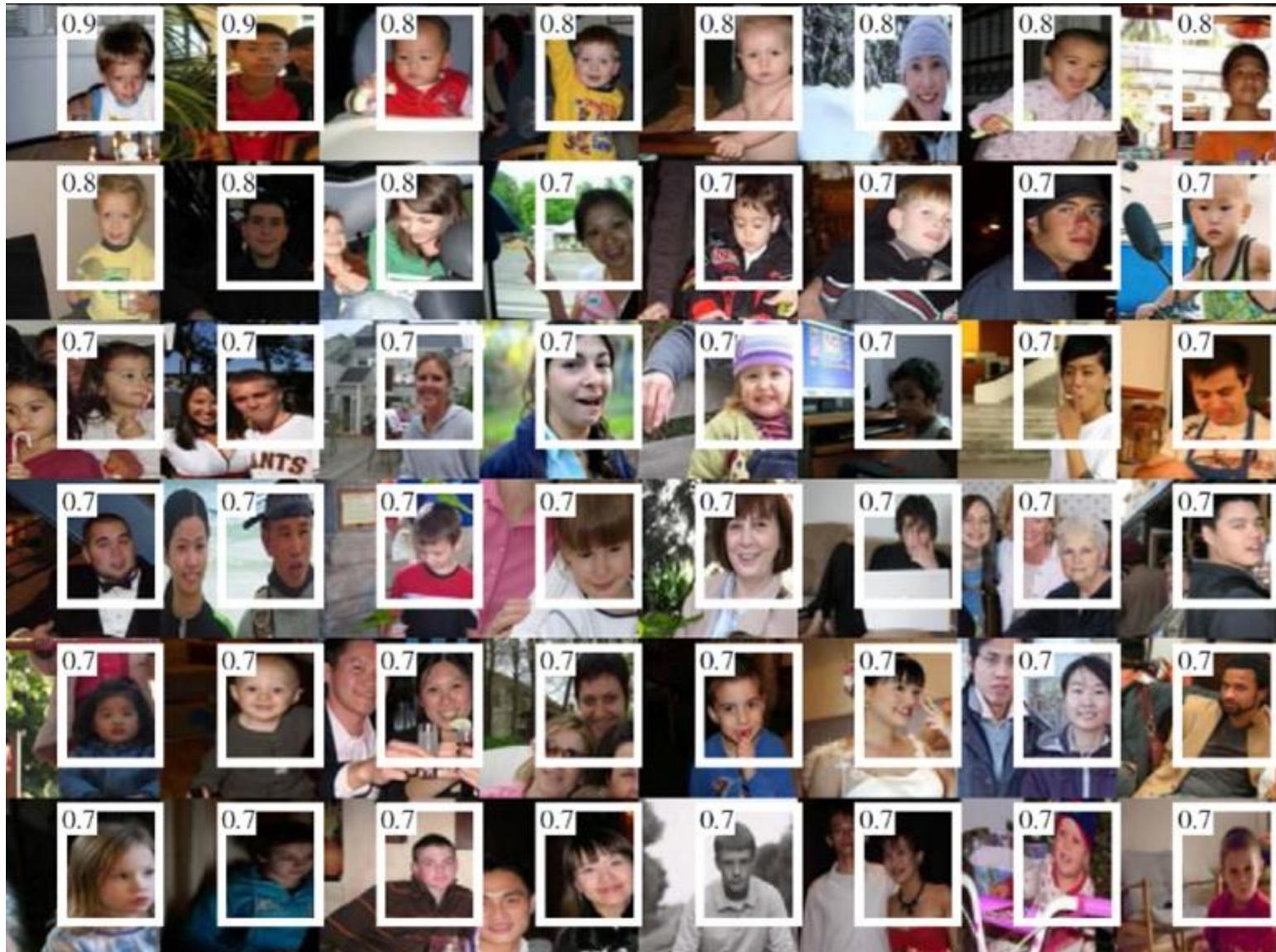
person: HOG template

# Individual Neuron Activation



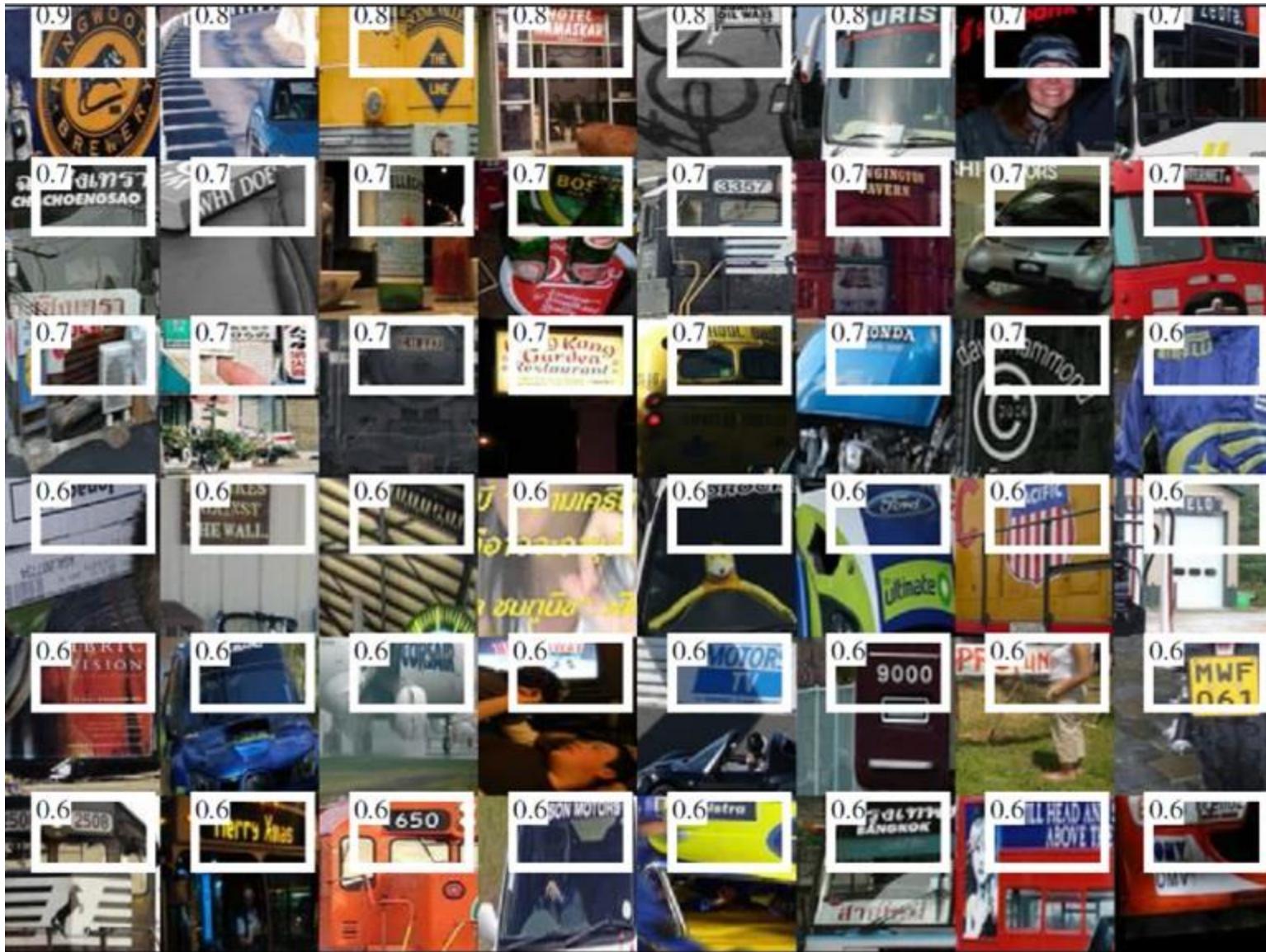
RCNN [Girshick et al. CVPR 2014]

# Individual Neuron Activation



RCNN [Girshick et al. CVPR 2014]

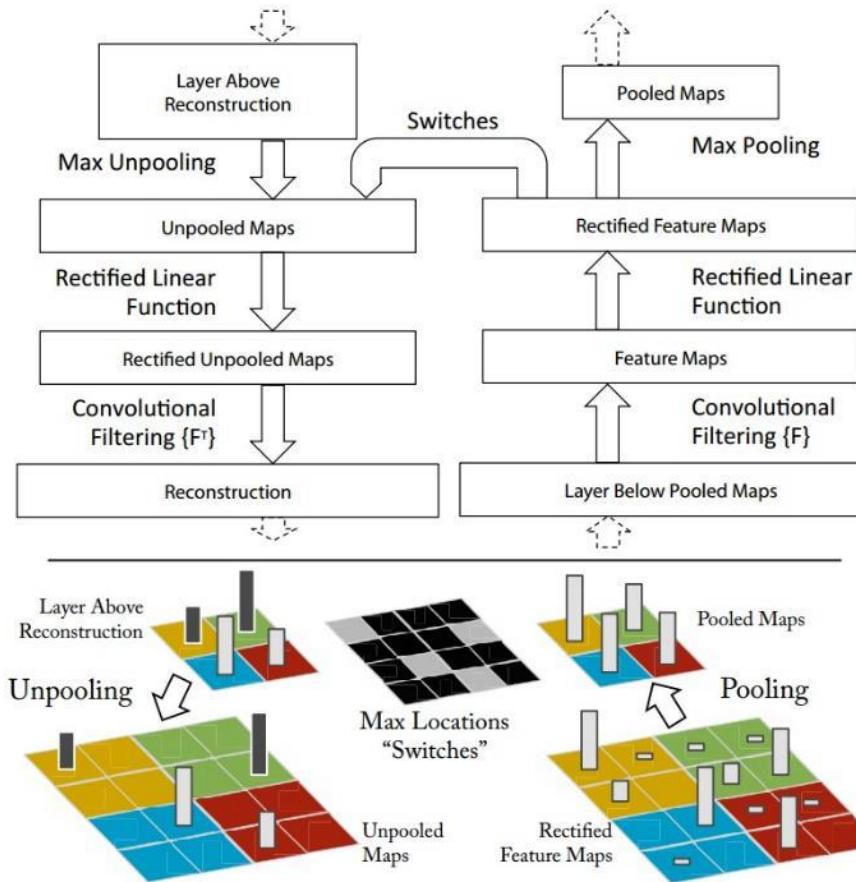
# Individual Neuron Activation



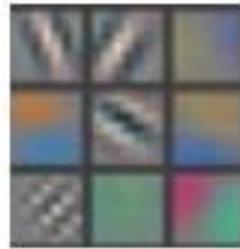
RCNN [Girshick et al. CVPR 2014]

# Map activation back to the input pixel space

- What input pattern originally caused a given activation in the feature maps?



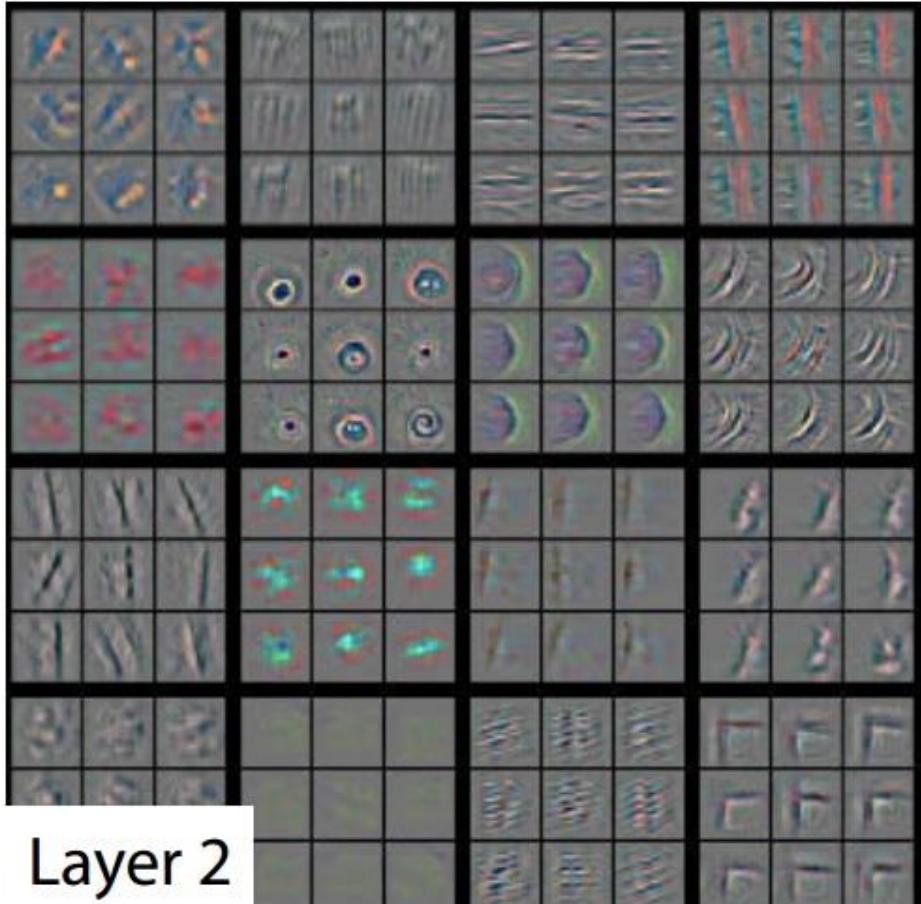
# Layer 1



Layer 1



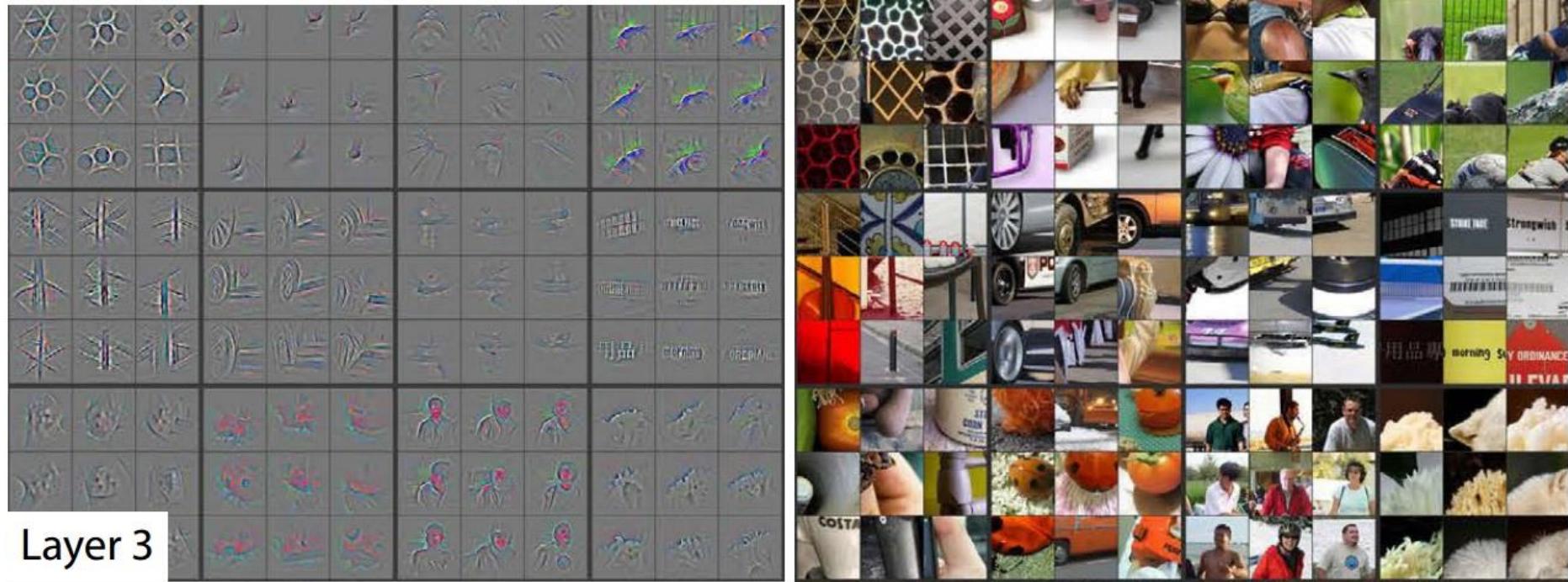
# Layer 2



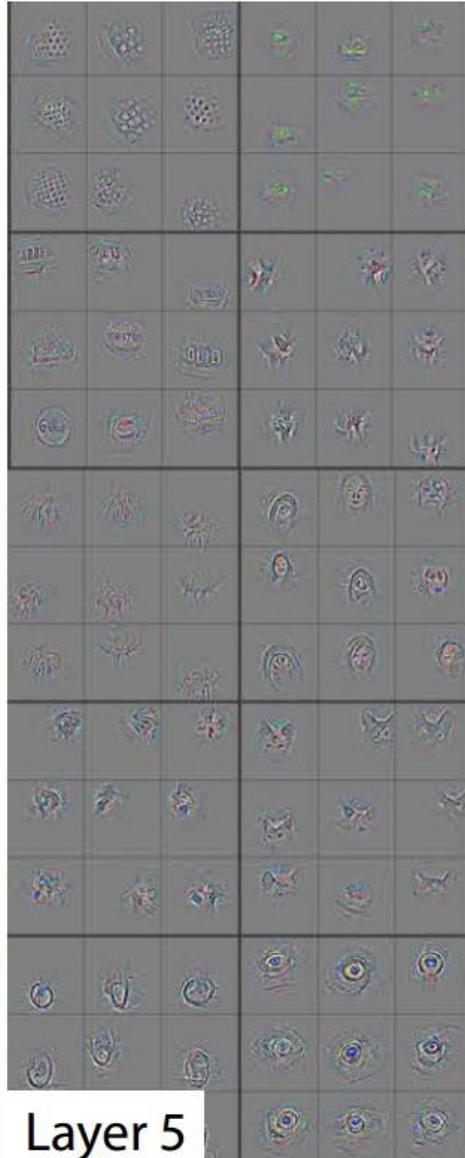
Layer 2



# Layer 3



# Layer 4 and 5



Visualizing and Understanding Convolutional Networks [Zeiler and Fergus, ECCV 2014]

# Network Dissection

House

res5c unit 1410



Dog

IoU=0.142 res5c unit 1573



Train

IoU=0.216 res5c unit 924



Plant

IoU=0.293 res5c unit 264



Airplane

IoU=0.126 res5c unit 1243



IoU=0.172

ResNet-152

res5c unit 301



IoU=0.087 res5c unit 1718



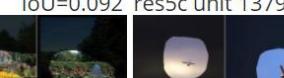
IoU=0.193 res5c unit 2001



IoU=0.255 res5c unit 766



IoU=0.092 res5c unit 1379



IoU=0.156

GoogLeNet

inception\_4e unit 789



IoU=0.137 inception\_4e unit 750



IoU=0.203 inception\_5b unit 626



IoU=0.145 inception\_4e unit 56



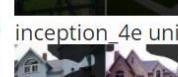
IoU=0.139 inception\_4e unit 92



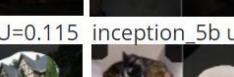
IoU=0.164

VGG-16

inception\_4e unit 175



IoU=0.115 inception\_5b unit 437



IoU=0.108 inception\_5b unit 415



IoU=0.143 inception\_4e unit 714

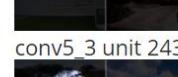


IoU=0.105 inception\_4e unit 759

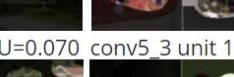


IoU=0.144

conv5\_3 unit 243



IoU=0.070 conv5\_3 unit 142



IoU=0.205 conv5\_3 unit 463



IoU=0.126 conv5\_3 unit 85



IoU=0.086 conv5\_3 unit 151



IoU=0.150

conv5\_3 unit 102



IoU=0.070 conv5\_3 unit 491



IoU=0.112 conv5\_3 unit 402



IoU=0.058 conv4\_3 unit 336



IoU=0.068 conv5\_3 unit 204

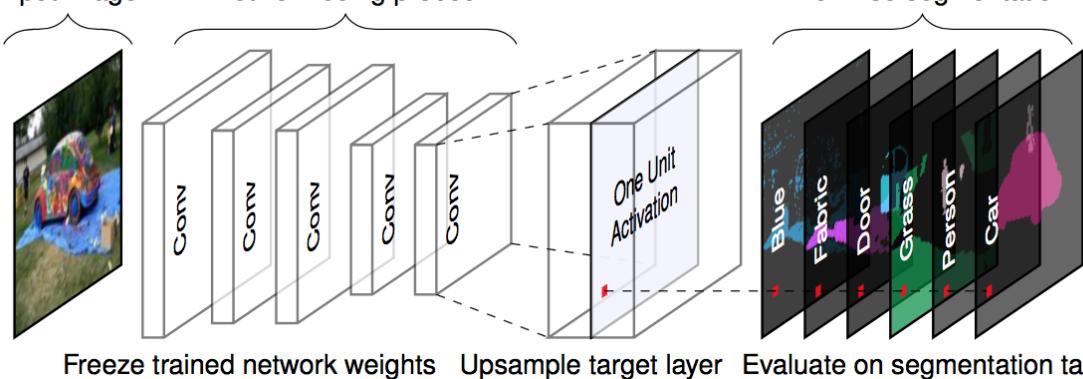


IoU=0.077

Input image

Network being probed

Pixel-wise segmentation



# Deep learning library

- TensorFlow
  - Research + Production
- PyTorch
  - Research
- Caffe2
  - Production



# Things to remember

- Convolutional neural networks
  - A cascade of conv + ReLU + pool
  - Representation learning
  - Advanced architectures
  - Tricks for training CNN
- Visualizing CNN
  - Activation
  - Dissection

# Resources

- <http://deeplearning.net/>
  - Hub to many other deep learning resources
- <https://github.com/ChristosChristofidis/awesome-deep-learning>
  - A resource collection deep learning
- <https://github.com/kjw0612/awesome-deep-vision>
  - A resource collection deep learning for computer vision
- <http://cs231n.stanford.edu/syllabus.html>
  - Nice course on CNN for visual recognition

# Things to remember

- Overview
  - Neuroscience, Perceptron, multi-layer neural networks
- Convolutional neural network (CNN)
  - Convolution, nonlinearity, max pooling
  - CNN for classification and beyond
- Understanding and visualizing CNN
  - Find images that maximize some class scores; visualize individual neuron activation, input pattern and images; breaking CNNs
- Training CNN
  - Dropout; data augmentation; batch normalization; transfer learning

