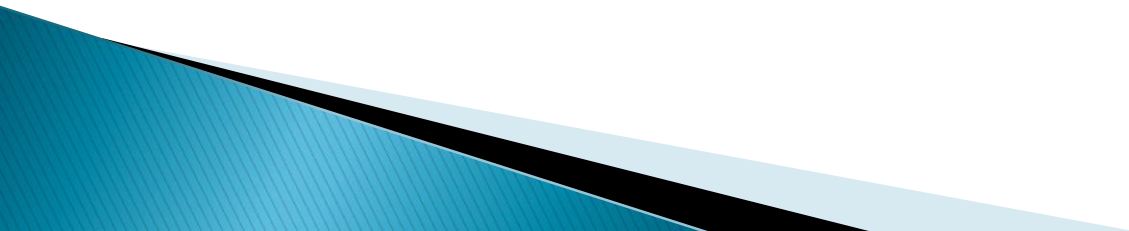


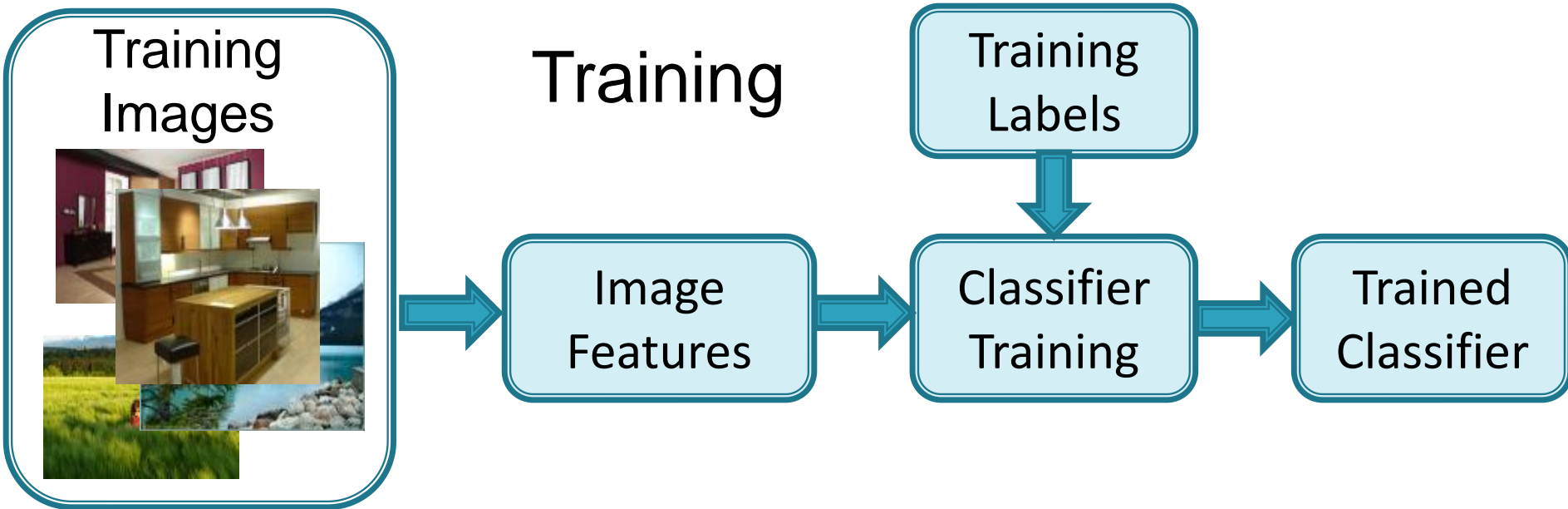
# Deep Learning for Visual Recognition



# Introduction



# Traditional Image Categorization: Training phase



# Traditional Image Categorization: Testing phase

Training  
Images

Training

Training  
Labels

Image  
Features

Classifier  
Training

Trained  
Classifier

Testing

Image  
Features

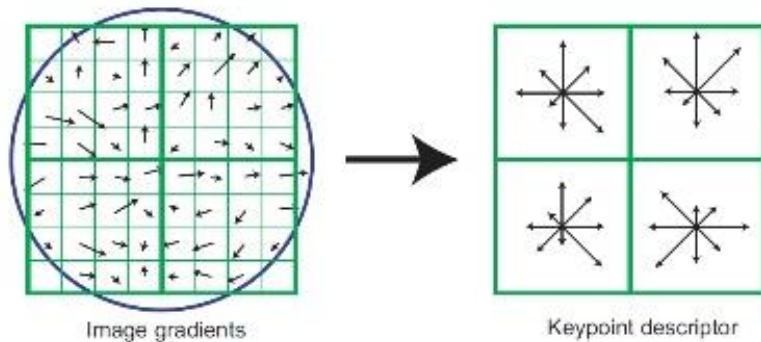
Trained  
Classifier

Prediction  
**Outdoor**

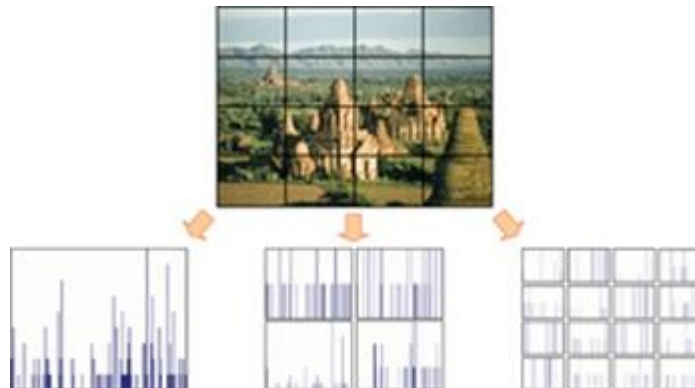
Test Image

Slide credit: Jia-Bin Huang

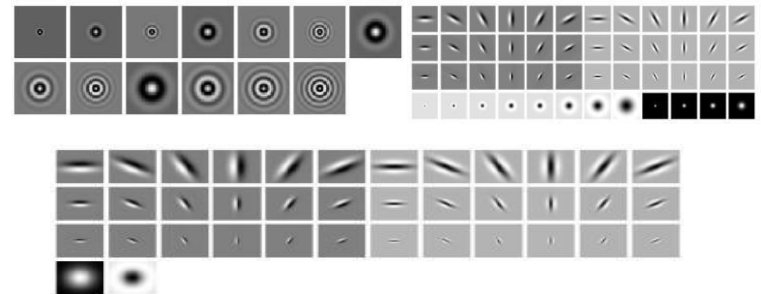
# Feature Design



SIFT [Lowe IJCV 04]



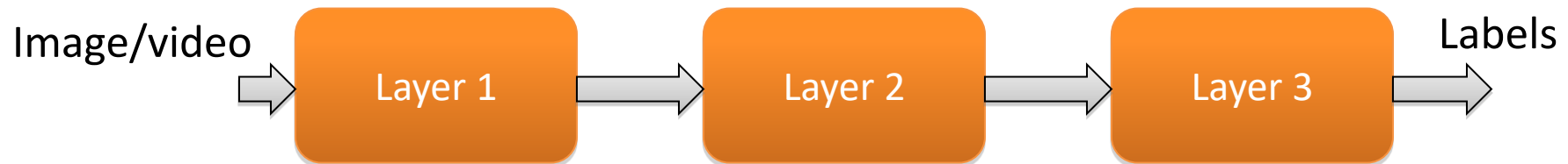
Harris Corner



Textons

# Hierarchy of Feature Extractors

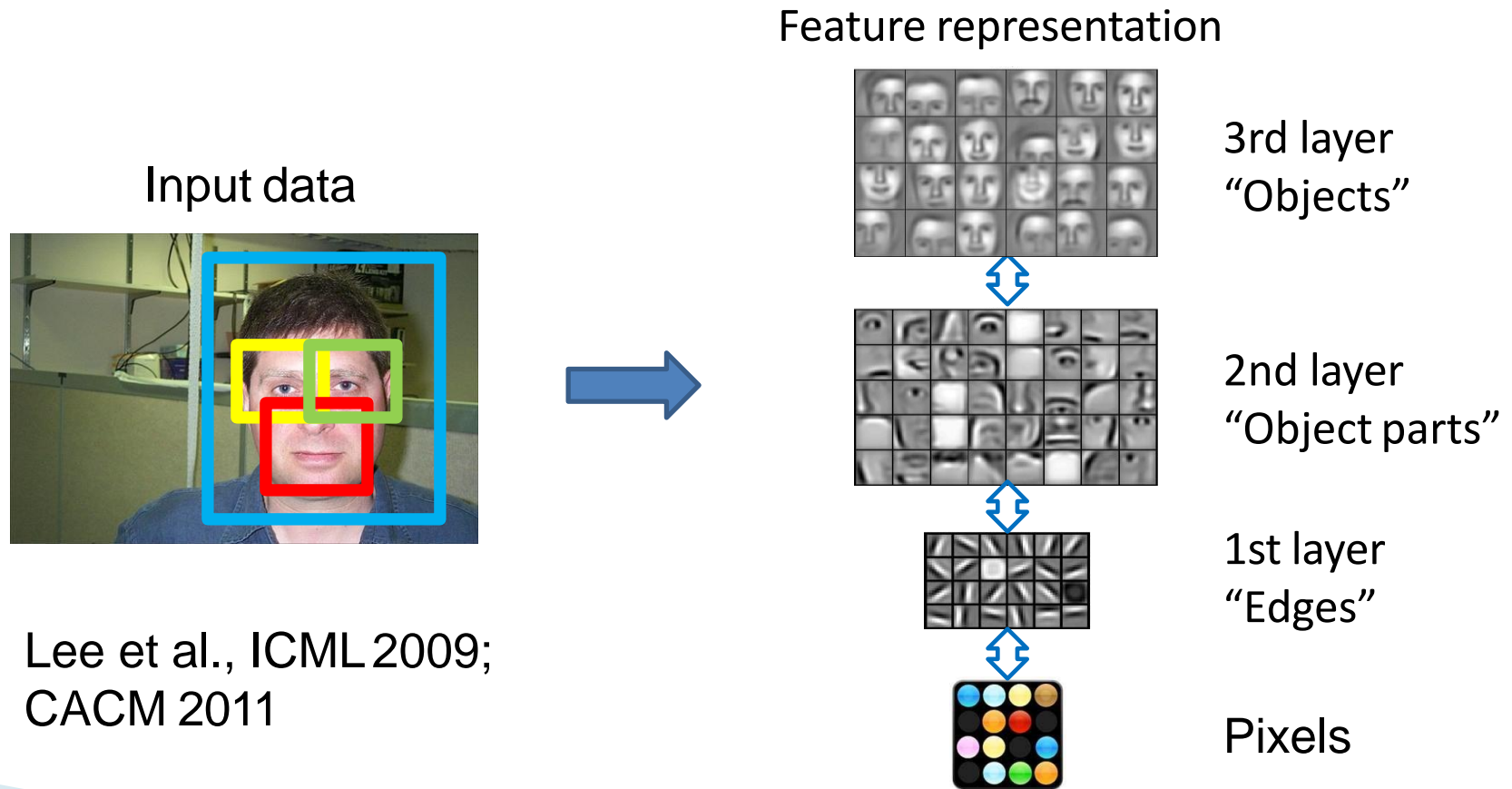
- Each layer of hierarchy extracts features from output of previous layer
- All the way from pixels → classifier
- Layers have the (nearly) same structure



- Train all layers jointly

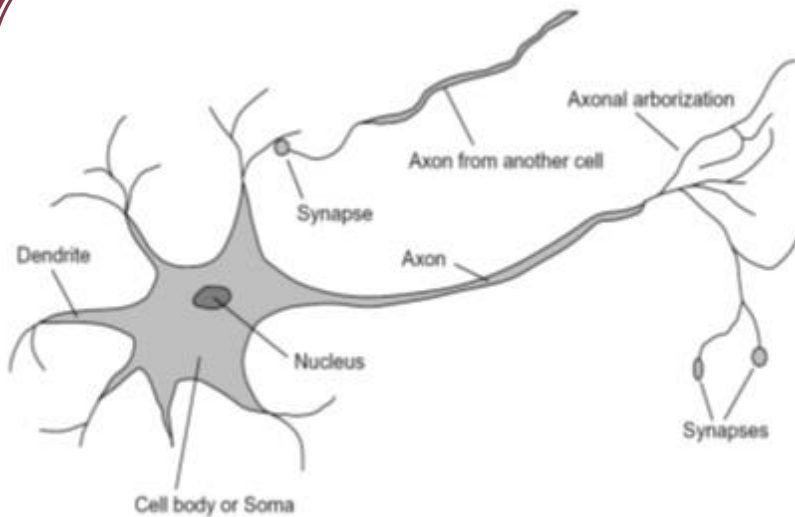
# Learning Feature Hierarchy

Goal: **Learn** useful higher-level features from images

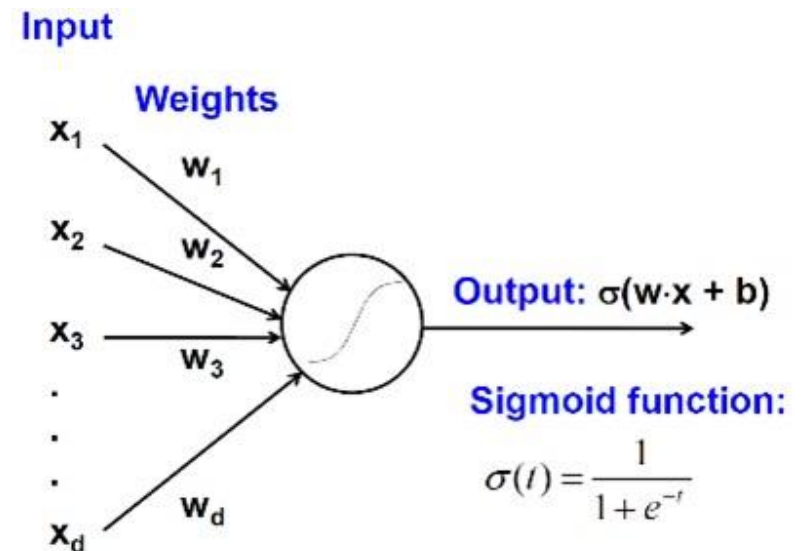




# Biological neuron and Perceptrons



A biological neuron



An artificial neuron (Perceptron)  
- a linear classifier

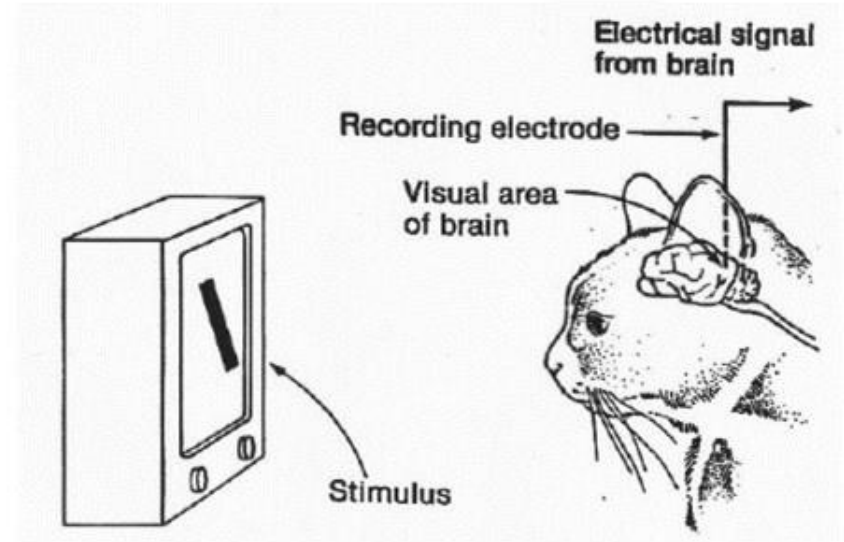




# Simple, Complex and Hypercomplex cells

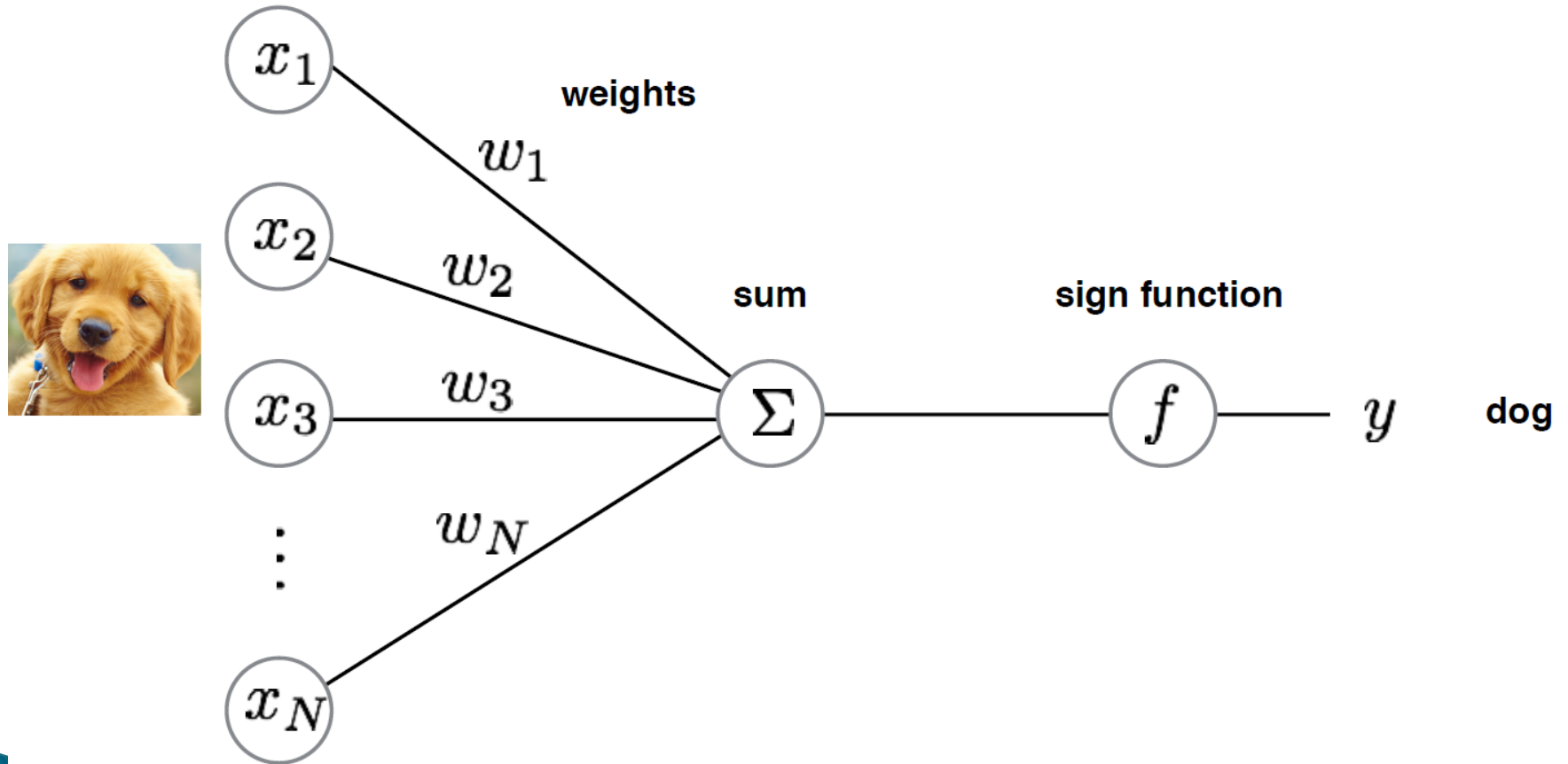


David H. Hubel and Torsten Wiesel



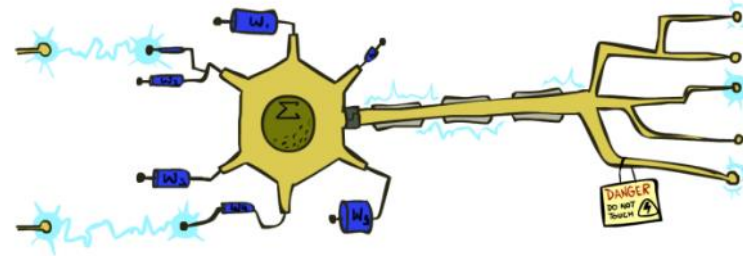
Suggested a **hierarchy** of **feature detectors** in the visual cortex, with higher level features responding to patterns of activation in lower level cells, and propagating activation upwards to still higher level cells.

# Perceptron: for image classification



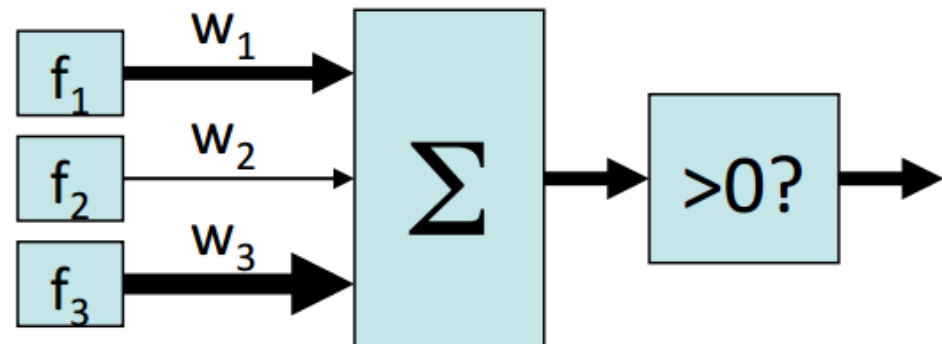
# Neuron: Linear Perceptron

- ▶ Inputs are **feature values**
- ▶ Each feature has a **weight**
- ▶ Sum is the **activation**



$$\text{activation}_w(x) = \sum_i w_i \cdot f_i(x) = w \cdot f(x)$$

- ▶ If the activation is:
  - Positive, output +1
  - Negative, output -1



# Perceptron training algorithm

- Initialize weights  $\mathbf{w}$  randomly
- Cycle through training examples in multiple passes (*epochs*)
- For each training example  $\mathbf{x}$  with label  $y$ :
  - Classify with current weights:

$$y' = \text{sgn}(\mathbf{w} \times \mathbf{x})$$

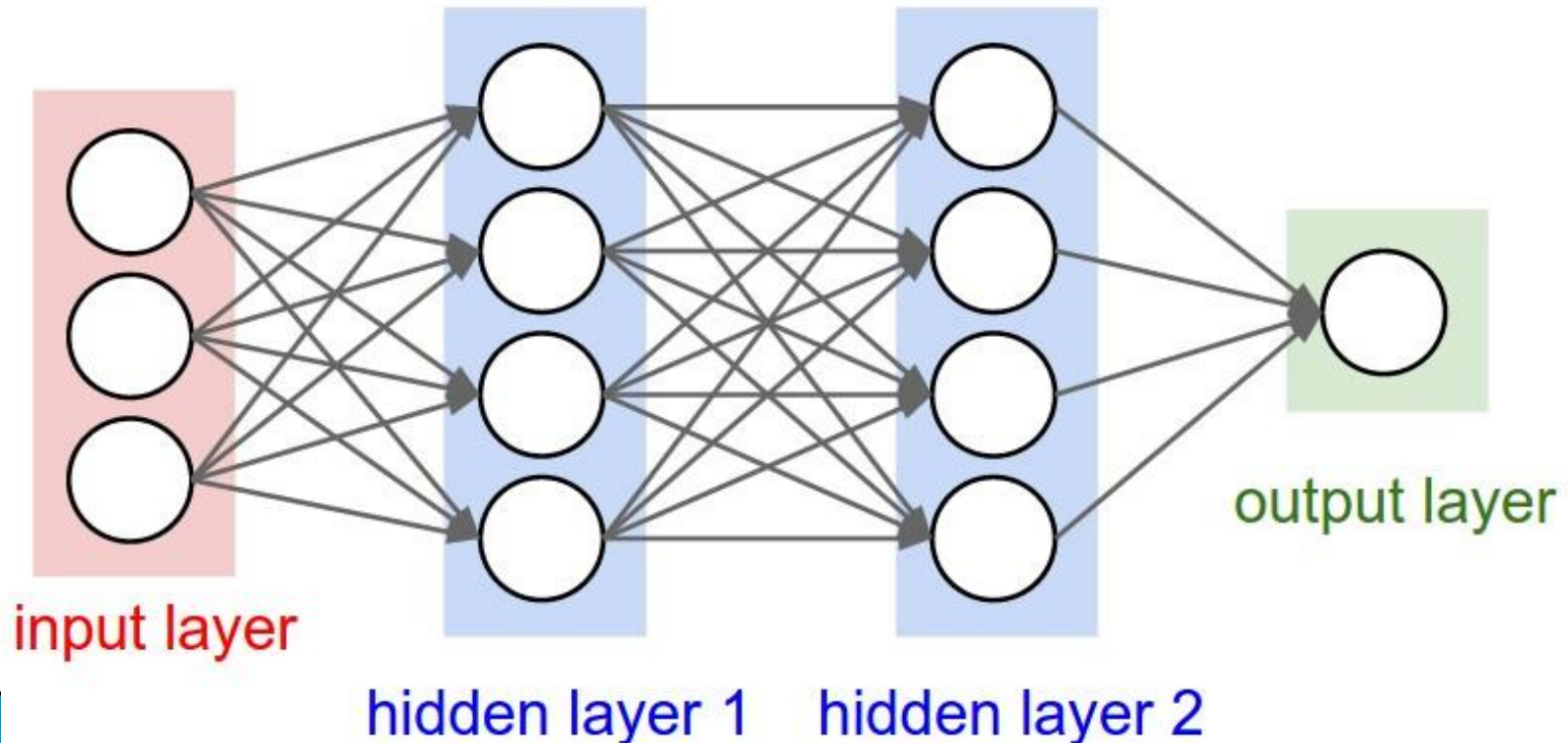
- If classified incorrectly, update weights:

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha (y - y') \mathbf{x}$$

( $\alpha$  is a positive *learning rate* that decays over time)

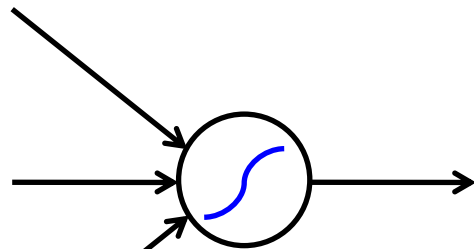
# Multi-layer perceptrons

- ▶ To make nonlinear classifiers out of perceptrons, build a multi-layer neural network!
  - This requires each perceptron to have a nonlinearity

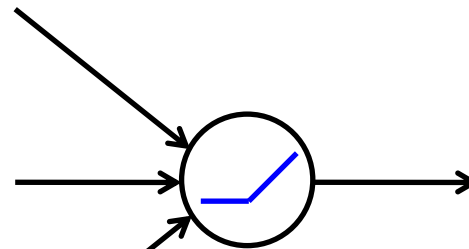


# Multi-layer perceptrons

- ▶ To make nonlinear classifiers out of perceptrons, build a multi-layer neural network!
  - This requires each perceptron to have a nonlinearity
  - To be trainable, the nonlinearity should be differentiable



**Sigmoid:**  $g(t) = \frac{1}{1 + e^{-t}}$



**Rectified linear unit (ReLU):**  $g(t) = \max(0, t)$

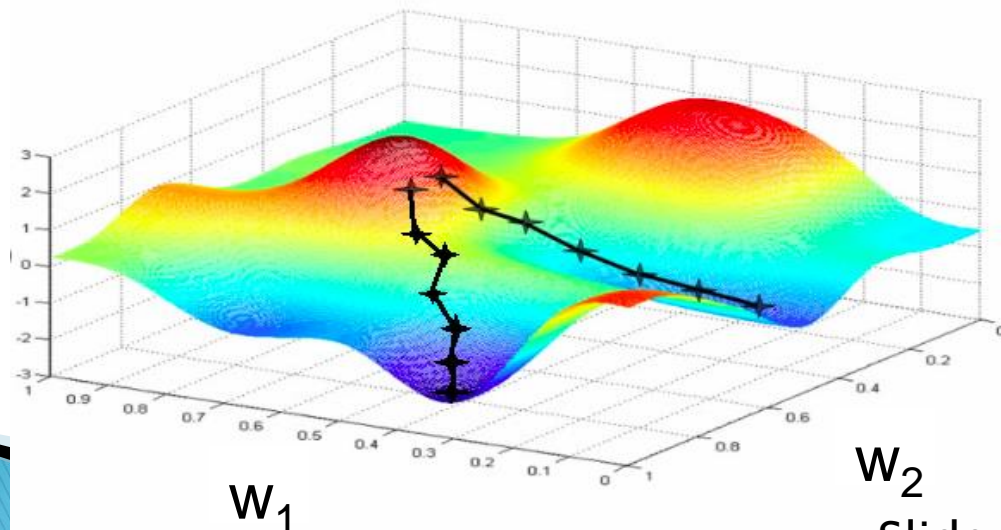


# Training of multi-layer networks

- Find network weights to minimize the error between true and estimated labels of training examples:

$$E(\mathbf{w}) = \sum_{j=1}^N \left( y_j - f_{\mathbf{w}}(\mathbf{x}_j) \right)^2$$

- Update weights by gradient descent:  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$



# Training of multi-layer networks

- ▶ Find network weights to minimize the error between true and estimated labels of training examples:

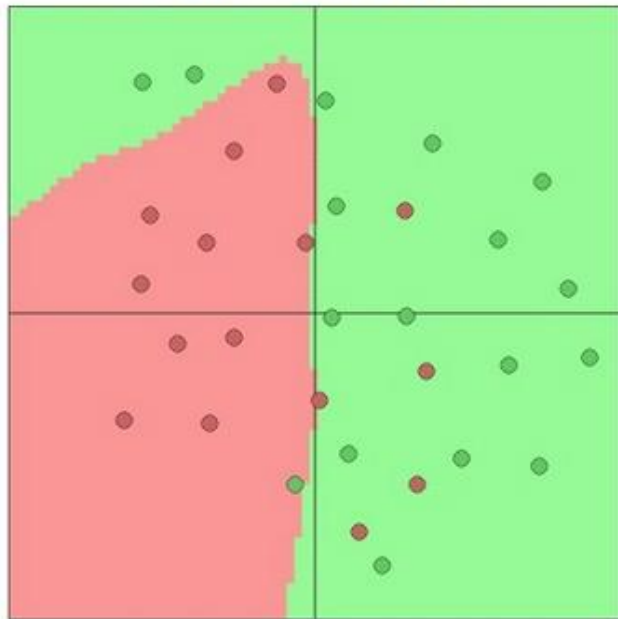
$$E(\mathbf{w}) = \sum_{j=1}^N \left( y_j - f_{\mathbf{w}}(\mathbf{x}_j) \right)^2$$

- ▶ Update weights by gradient descent:  $\mathbf{w} \leftarrow \mathbf{w} - \alpha \frac{\partial E}{\partial \mathbf{w}}$
- ▶ Back-propagation: gradients are computed in the direction from output to input layers and combined using chain rule
- ▶ Stochastic gradient descent: compute the weight update w.r.t. one training example (or a small batch of examples) at a time, cycle through training examples in random order in multiple epochs

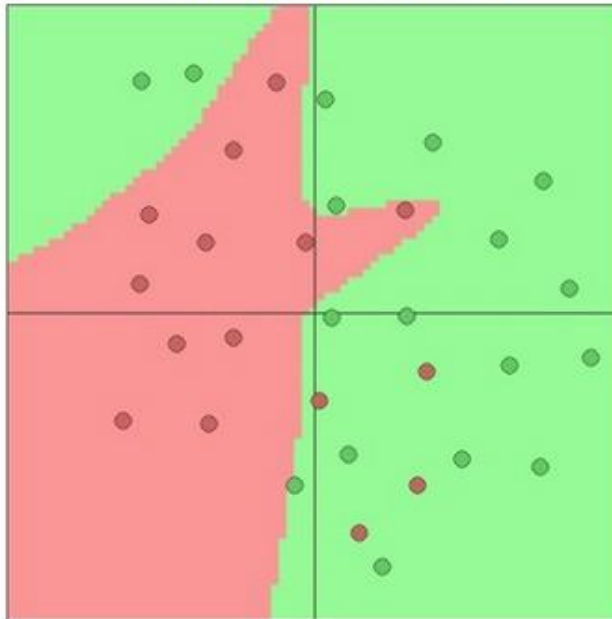
# Network with a single hidden layer

- ▶ Hidden layer size and network capacity:

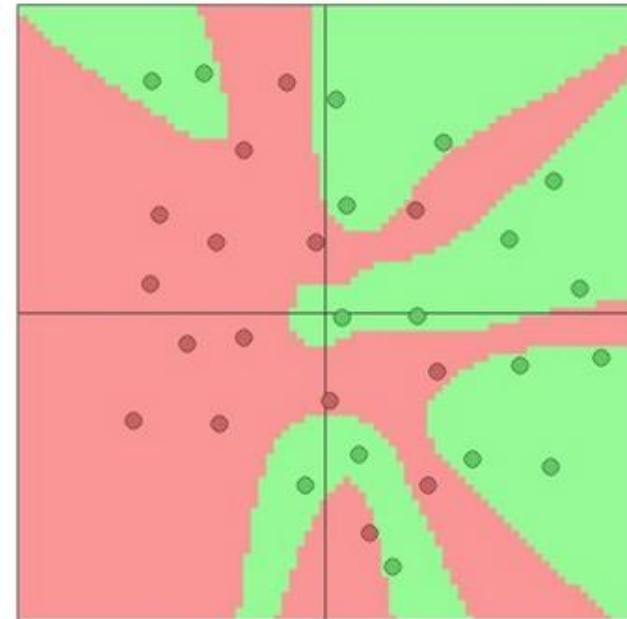
3 hidden neurons



6 hidden neurons



20 hidden neurons



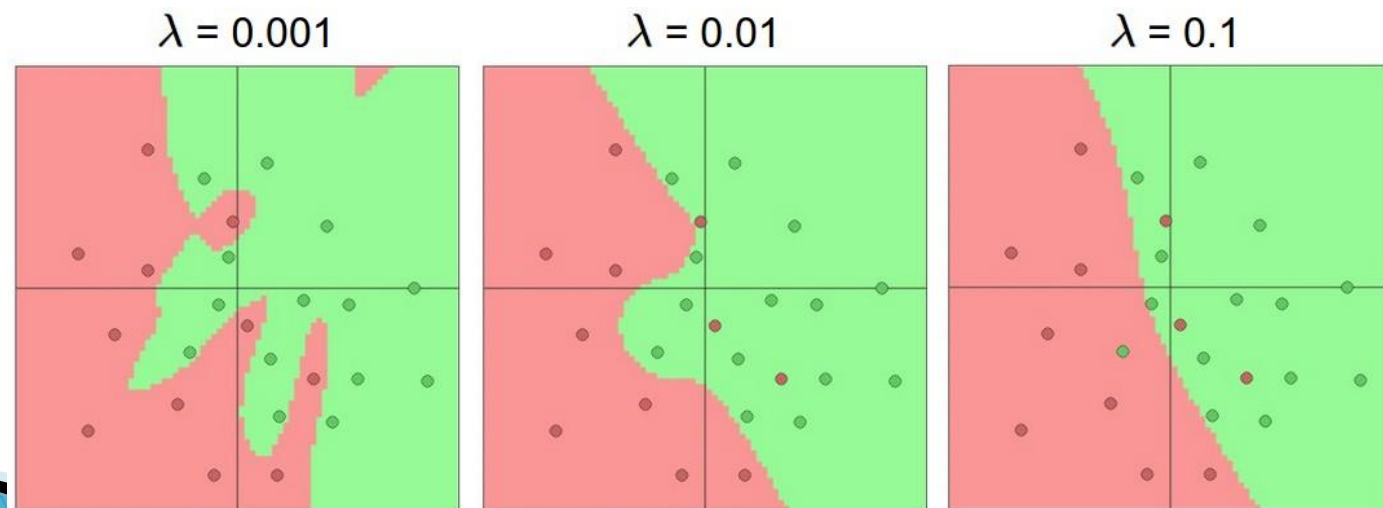
Source: <http://cs231n.github.io/neural-networks-1/>

# Regularization

- ▶ It is common to add a penalty on weight magnitudes to the objective function:

$$E(f) = \sum_{i=1}^N (y_i - f(\mathbf{x}_i))^2 + \frac{\lambda}{2} \sum_j w_j^2$$

- This encourages network to use all of its inputs “a little” rather than a few inputs “a lot”

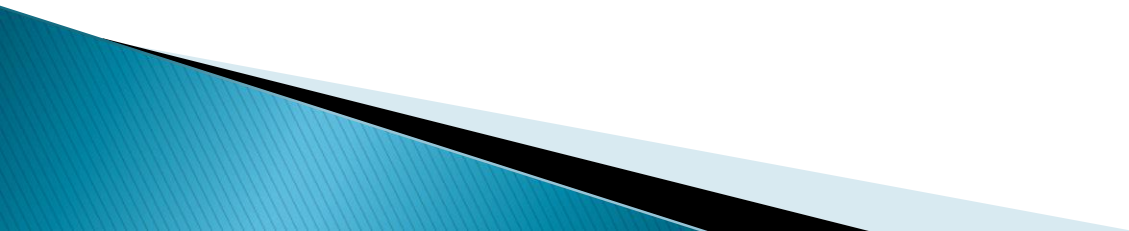


Source: <http://cs231n.github.io/neural-networks-1/>

# Neural networks: Pros and cons

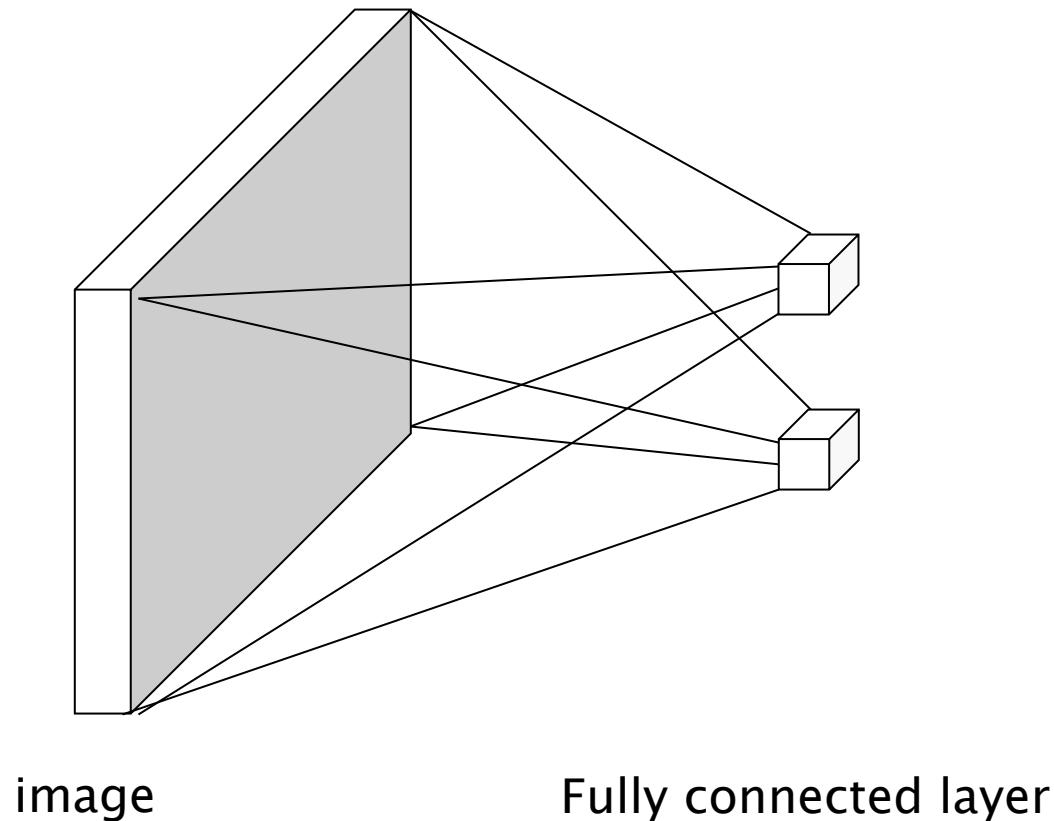
- **Pros**
  - Flexible and general function approximation framework
  - Can build extremely powerful models by adding more layers
- **Cons**
  - Hard to analyze theoretically (e.g., training is prone to local optima)
  - Huge amount of training data, computing power may be required to get good performance
  - The space of implementation choices is huge (network architectures, parameters)

# Convolutional Neural Network

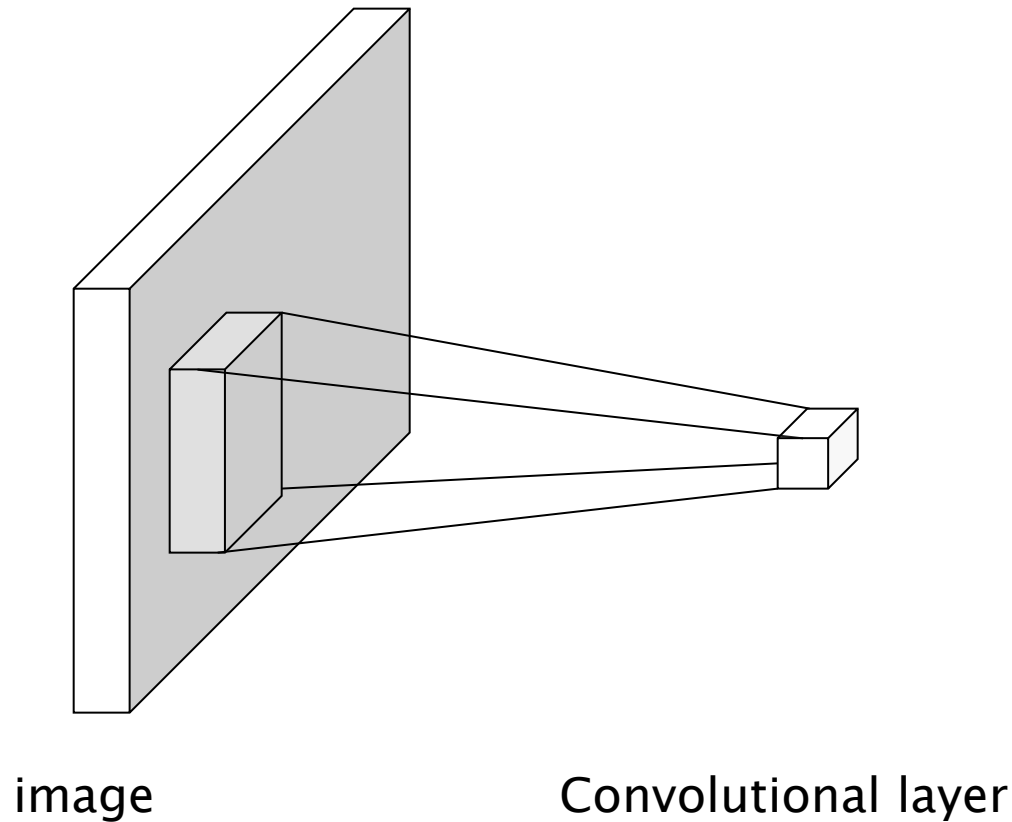




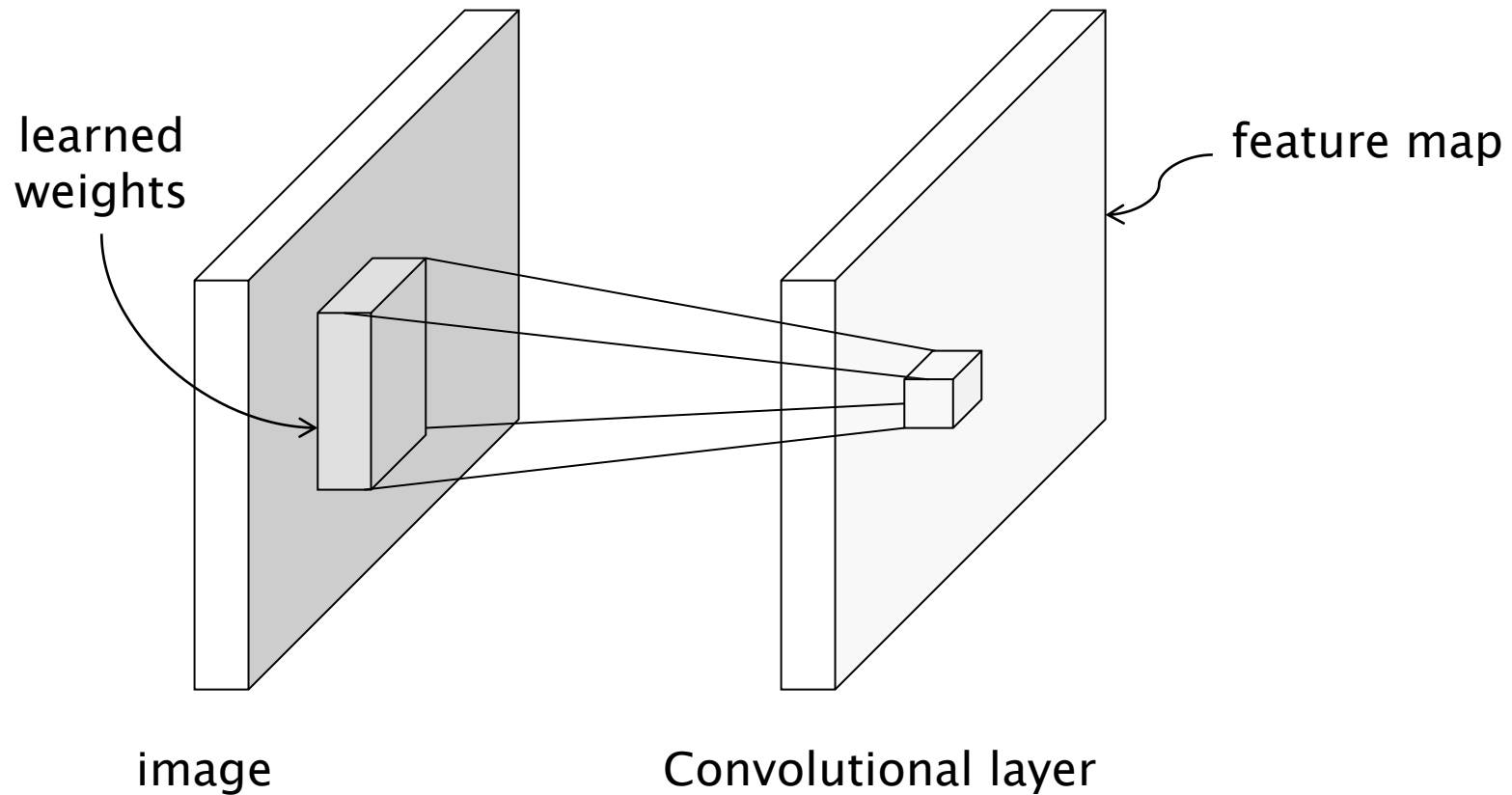
# Neural networks for images



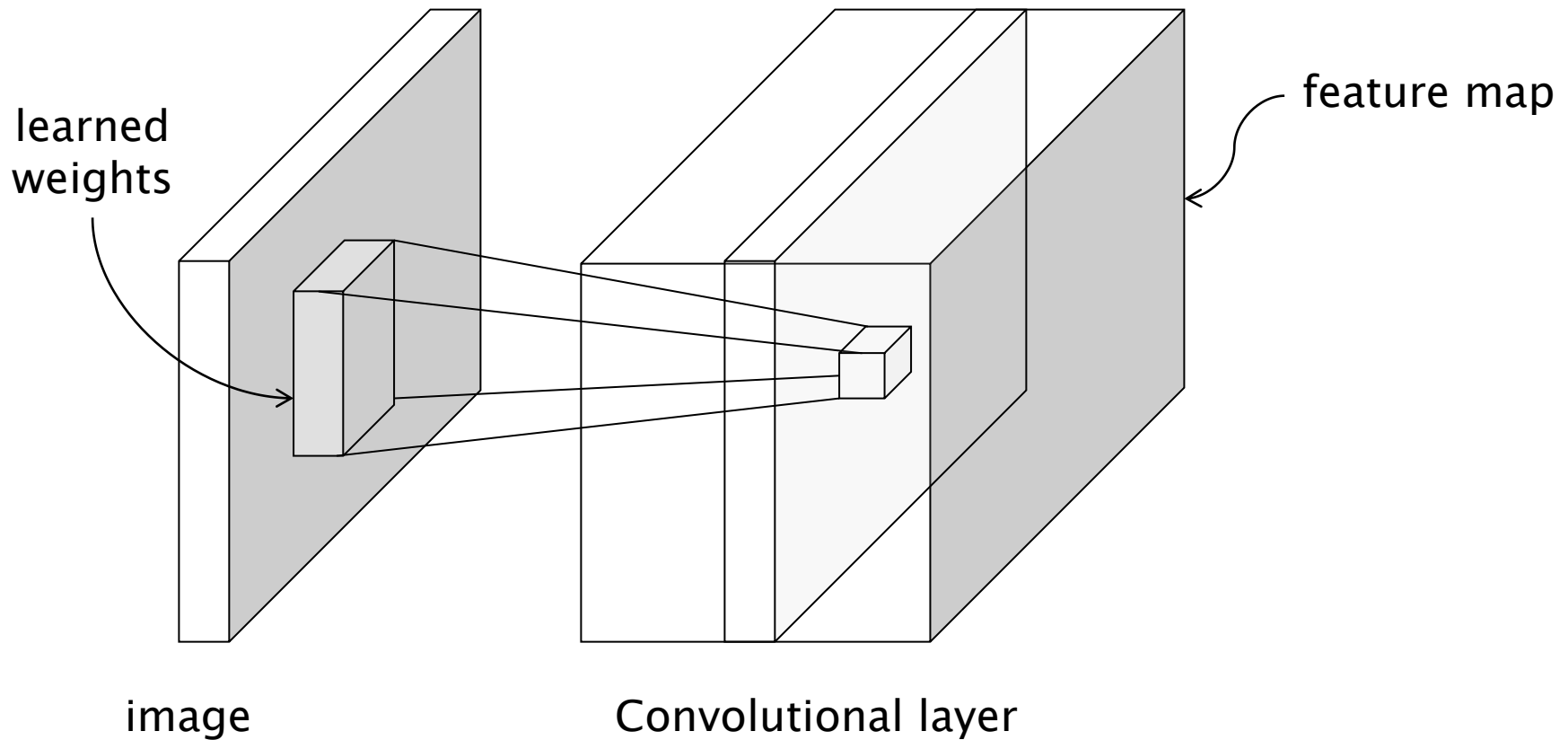
# Neural networks for images



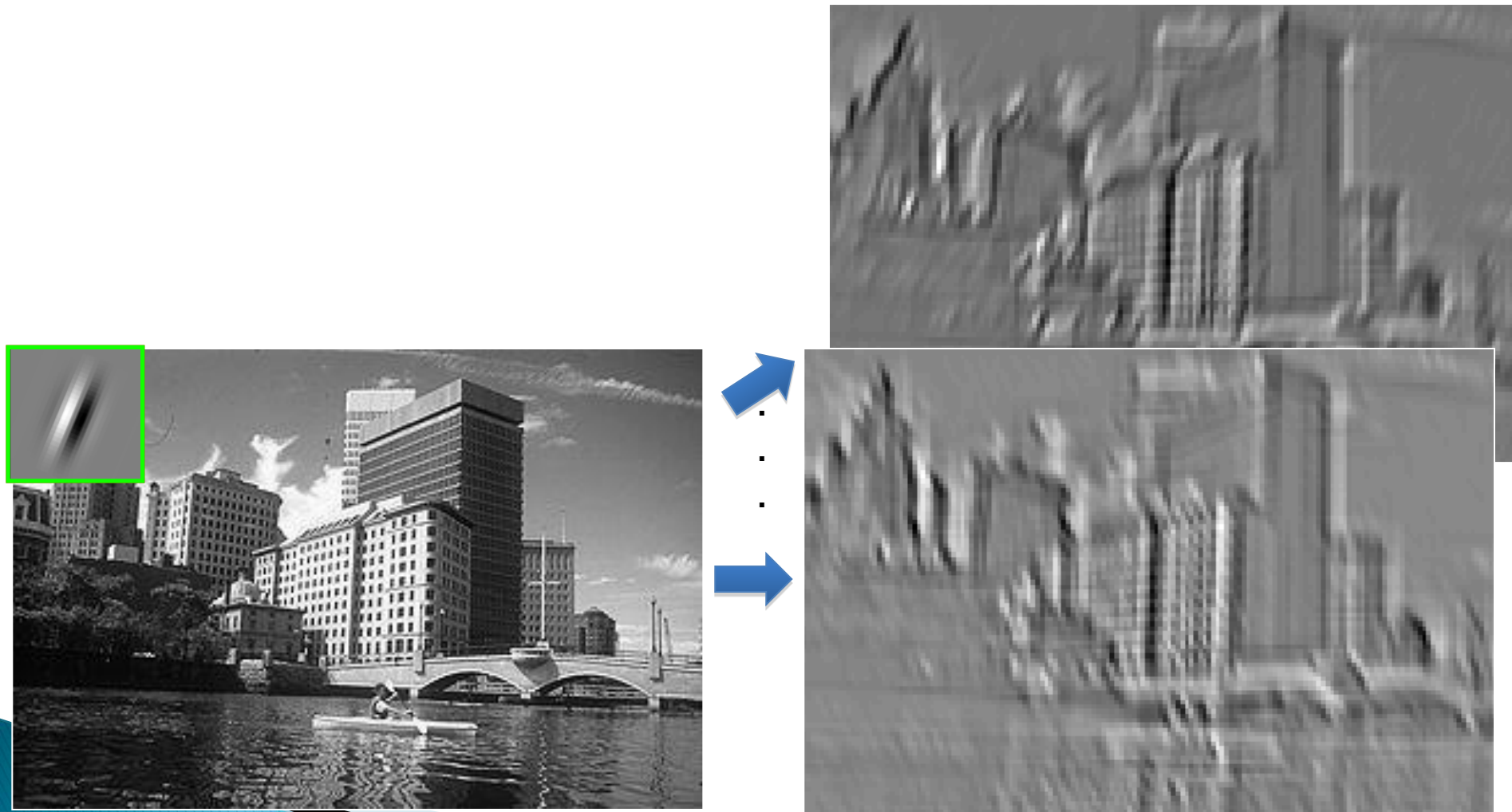
# Neural networks for images



# Neural networks for images



# Convolution as feature extraction

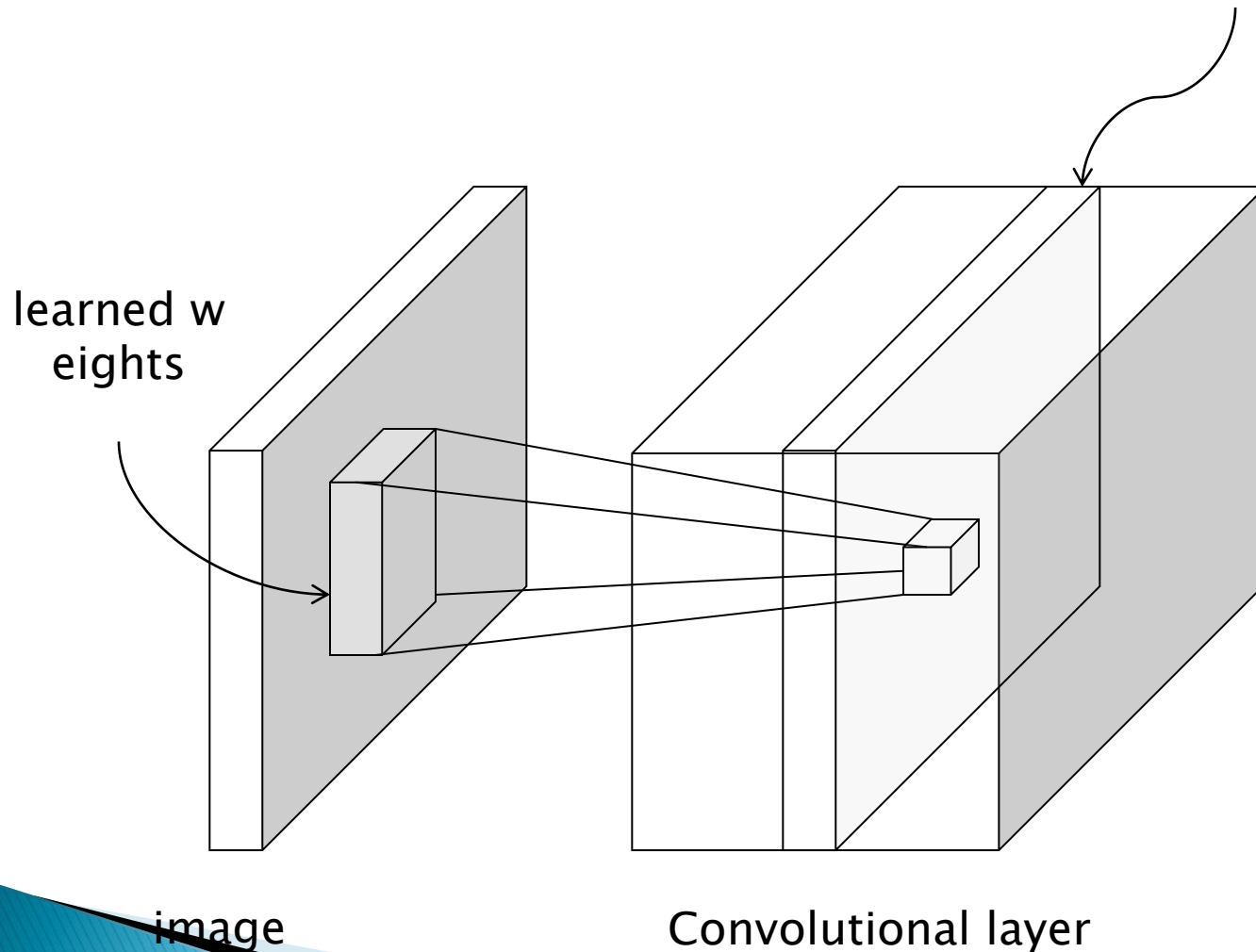


Input

Feature Map

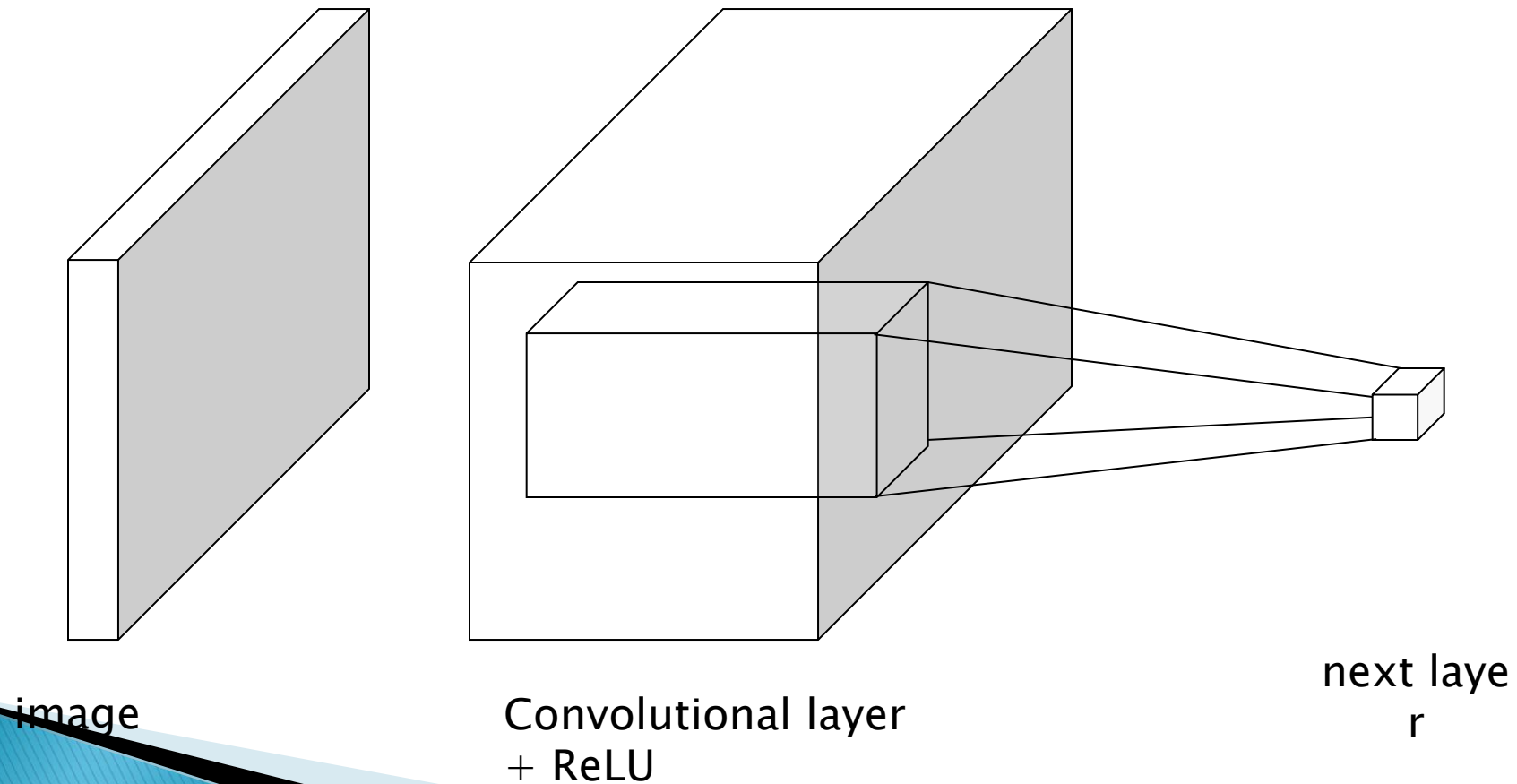
Slide credit: Svetlana Lazebnik

# Neural networks for images



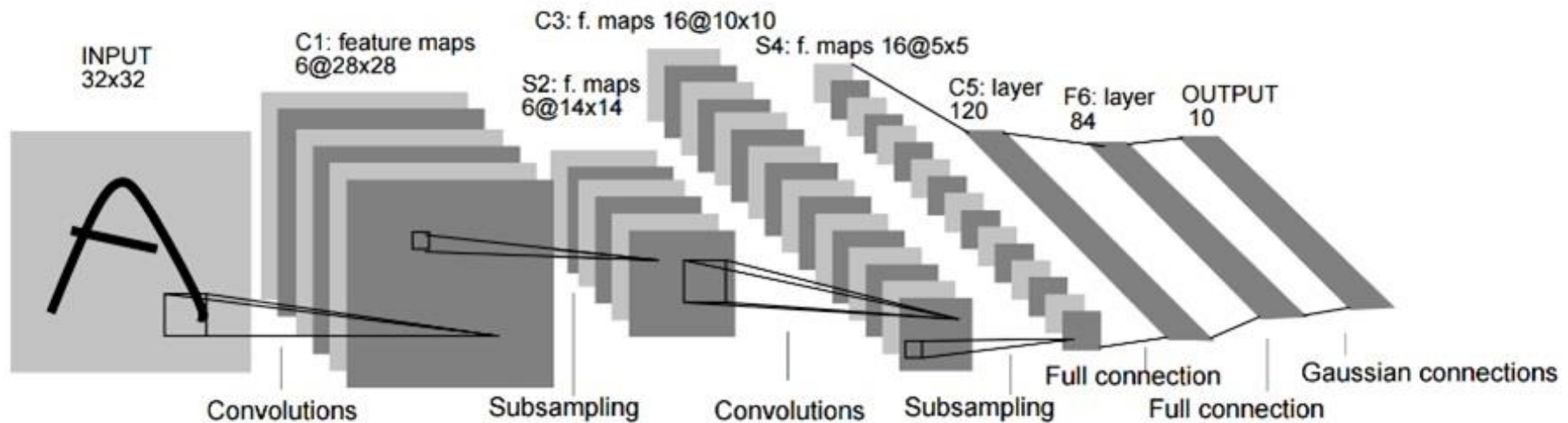


# Neural networks for images



Slide credit: Svetlana Lazebnik

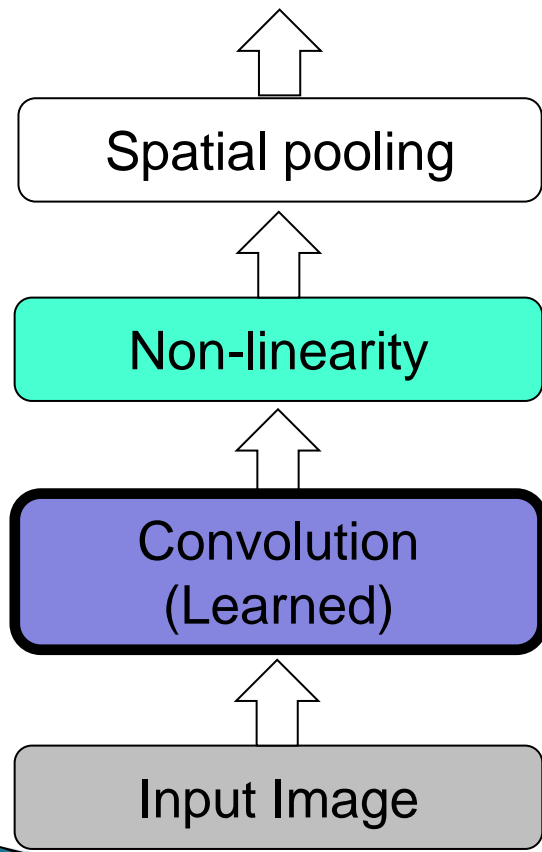
# LeNet [LeCun et al. 1998]



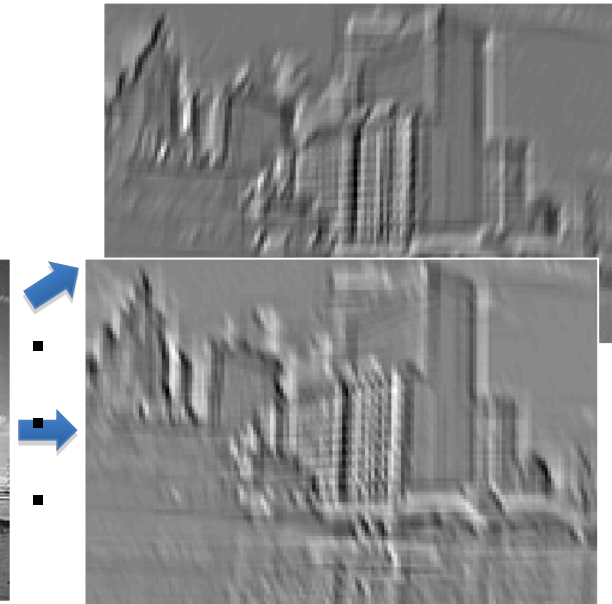
LeNet-1 from 1993

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

# Key operations in a CNN

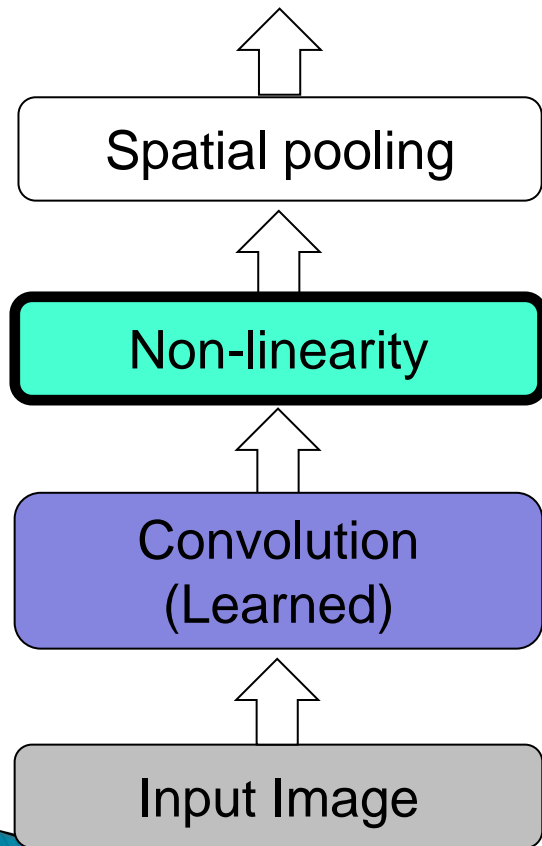


Input

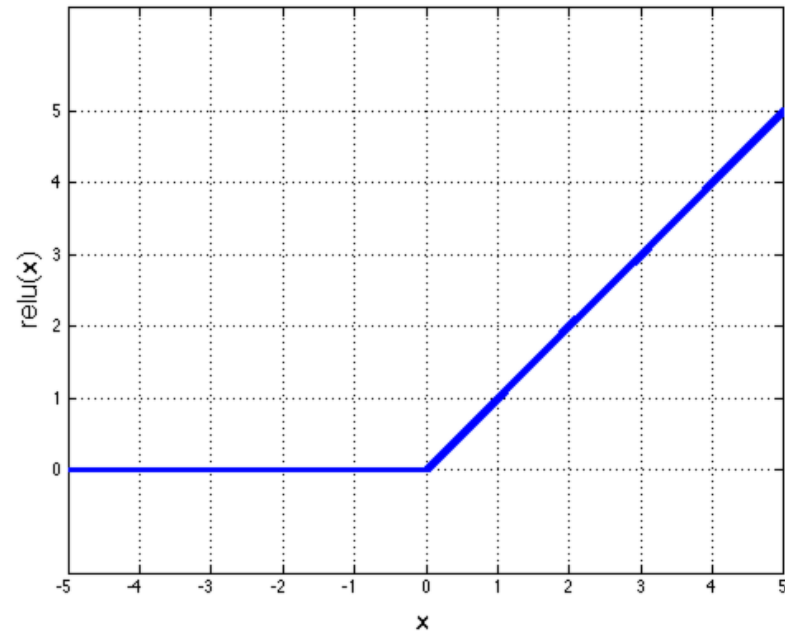


Feature Map

# Key operations in a CNN

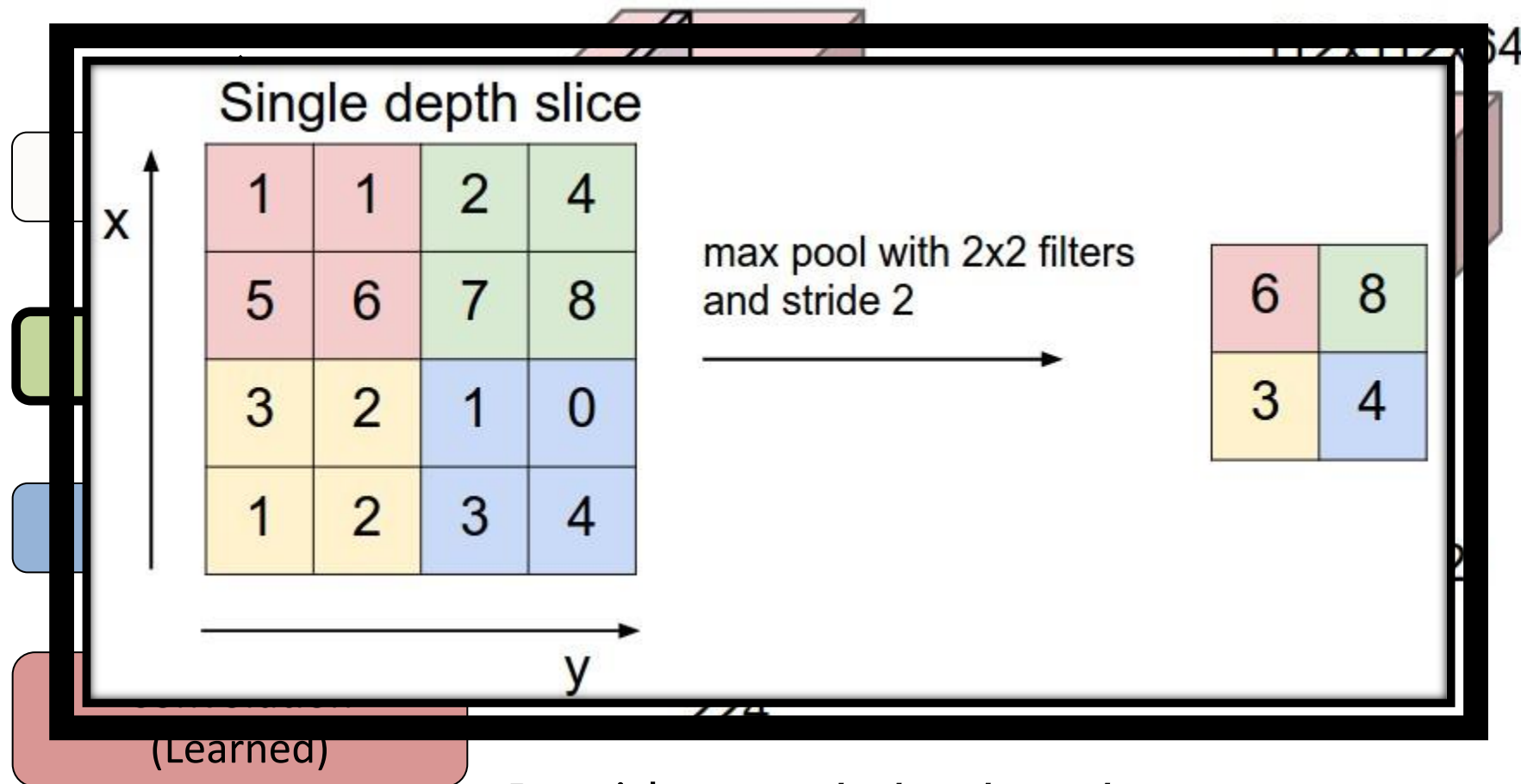


Rectified Linear Unit (ReLU)



# Key operations in a CNN

224x224x64

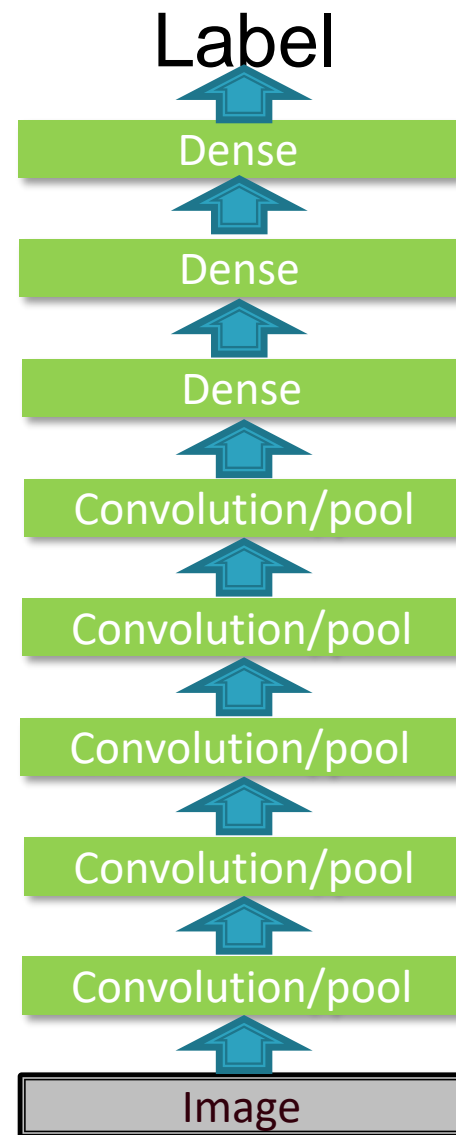
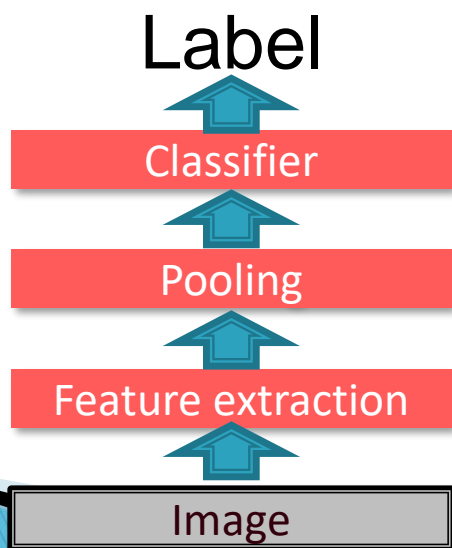


Provide *translation invariance*

Input Image

# Engineered vs. learned features

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



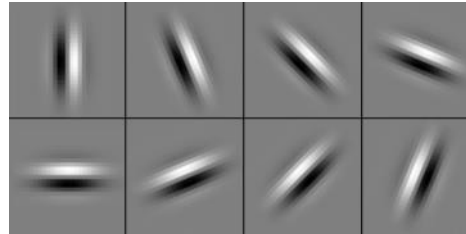


# SIFT Descriptor

Image  
Pixels



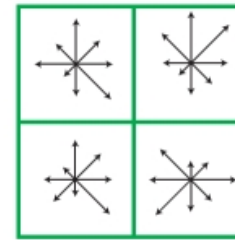
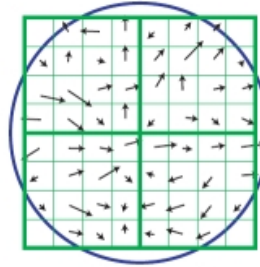
Apply  
oriented filters



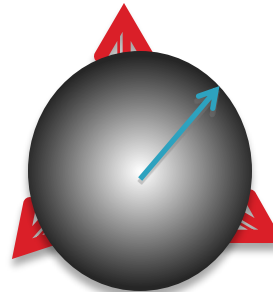
Lowe [IJCV 2004]



Spatial pool  
(Sum)



Normalize to unit  
length



Feature  
Vector

slide credit: R. Fergus

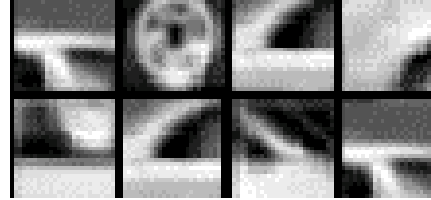
# Spatial Pyramid Matching

Lazebnik,  
Schmid,  
Ponce  
[CVPR 2006]

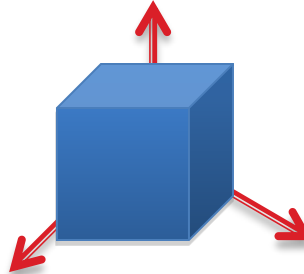
SIFT  
Features



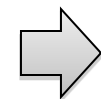
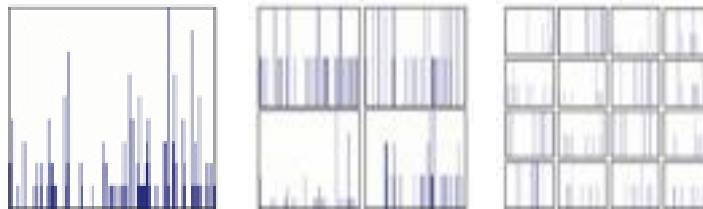
Filter with  
Visual Words



Take max VW  
response (L-inf  
normalization)



Multi-scale  
spatial pool  
(Sum)

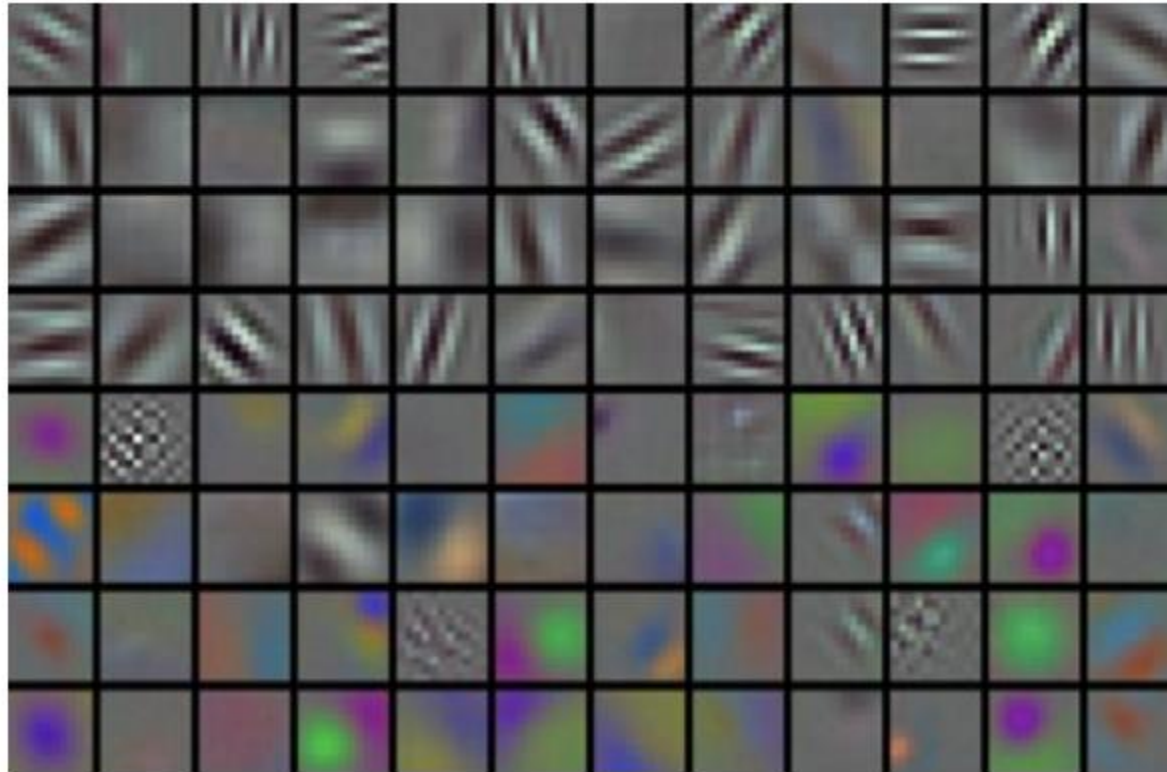


Classifier

slide credit: R. Fergus

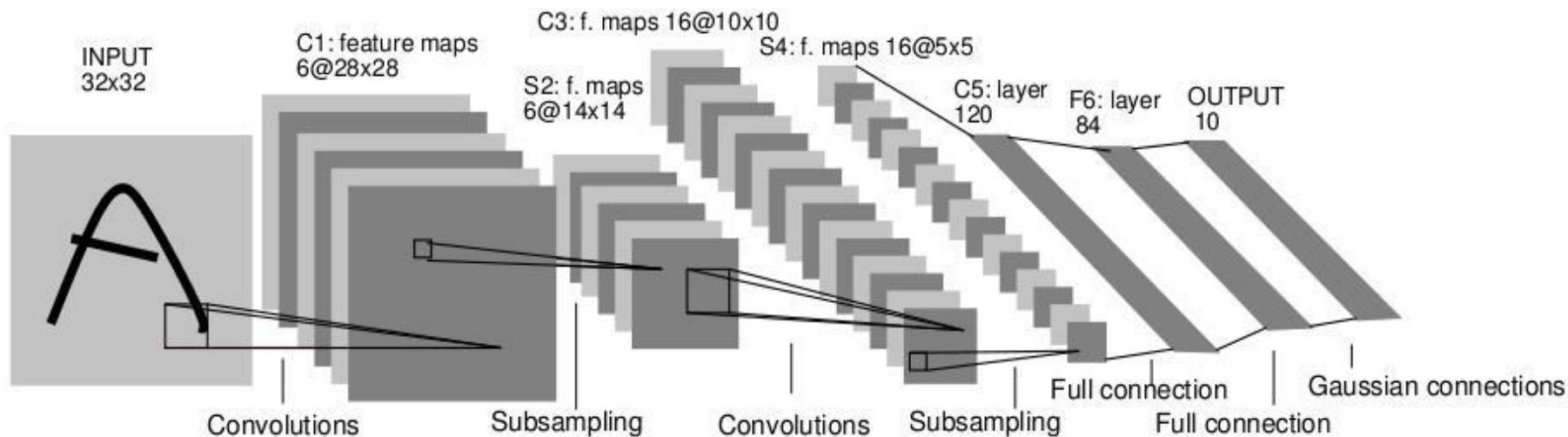
# Visualizing what was learned

- ▶ What do the learned filters look like?



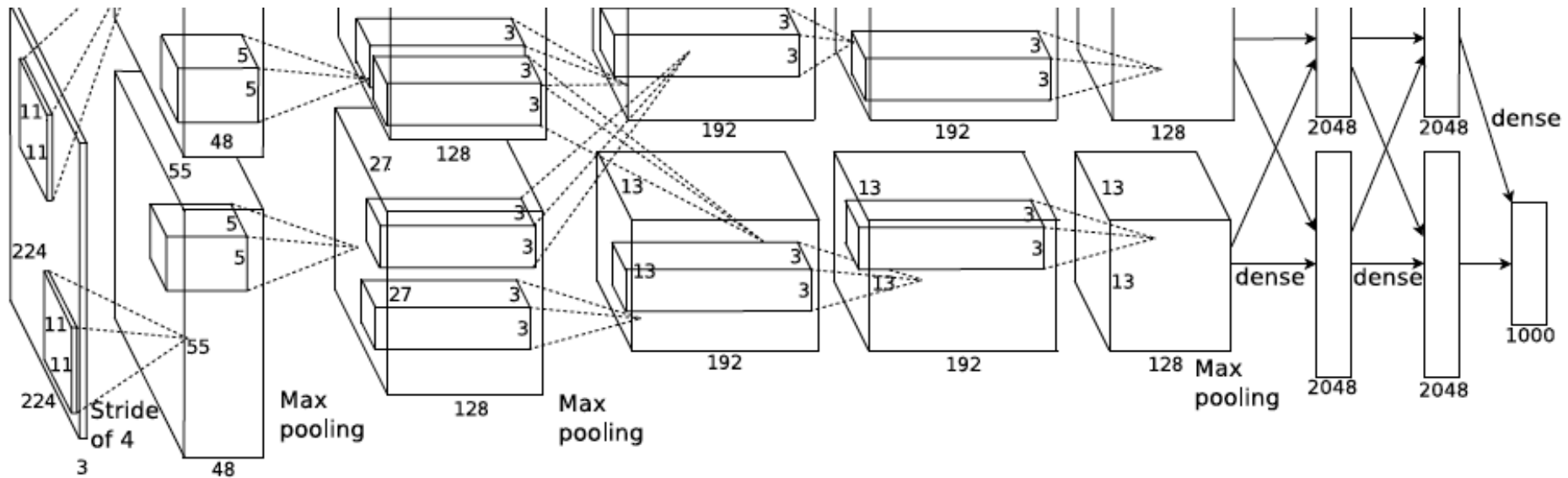
Typical first layer filters

# History of CNNs: LeNet-5



- Average pooling
- Sigmoid or tanh nonlinearity
- Fully connected layers at the end
- Trained on MNIST digit dataset with 60K training examples

# AlexNet: ILSVRC 2012 winner



- Similar framework to LeNet but:
    - Max pooling, ReLU nonlinearity
    - More data and bigger model (7 hidden layers, 650K units, 60M params)
    - GPU implementation (50x speedup over CPU)
      - Trained on two GPUs for a week
    - Dropout regularization
- A. Krizhevsky, I. Sutskever, and G. Hinton, [ImageNet Classification with Deep Convolutional Neural Networks](#), NIPS 2012