
Deep Learning IC Design

Outline

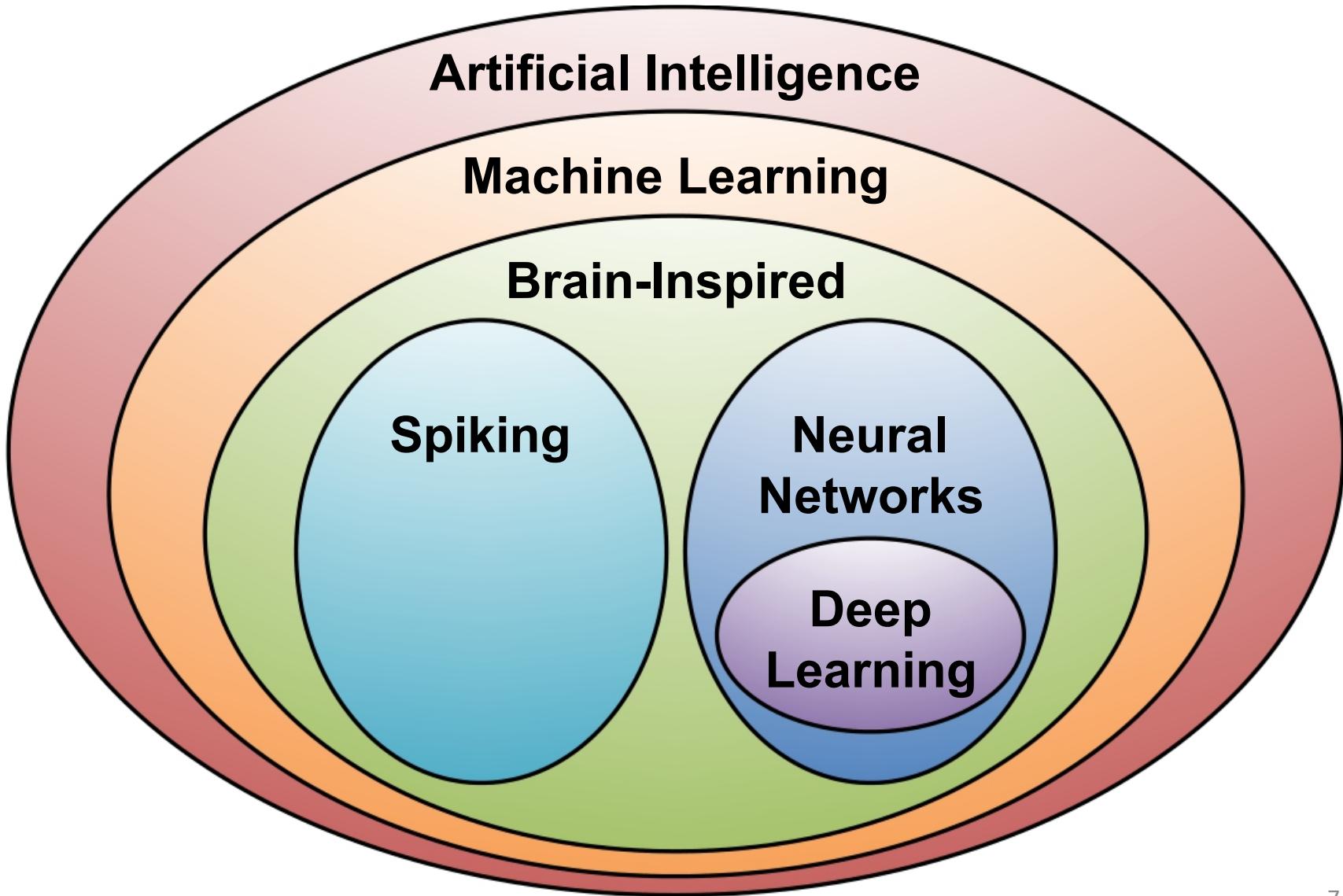
- **Overview of Deep Neural Networks**
- **DNN Development Resources**
- **Survey of DNN Hardware**
- **DNN Accelerators**
- **DNN Model and Hardware Co-Design**

Slides are modified from ISCA 2017 tutorial
by Vivienne Sze

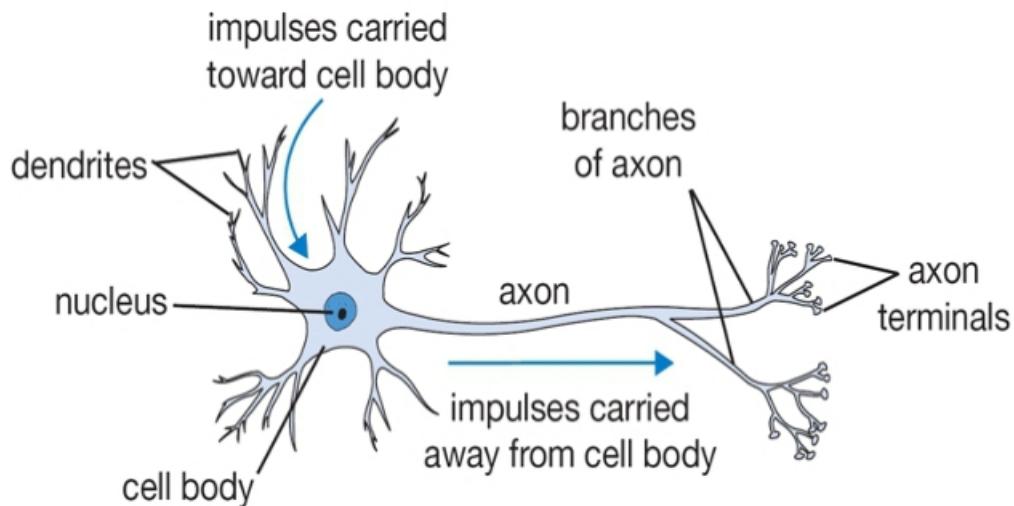
Website: <http://eyeriss.mit.edu/tutorial.html>

Background of Deep Neural Networks

Deep Learning

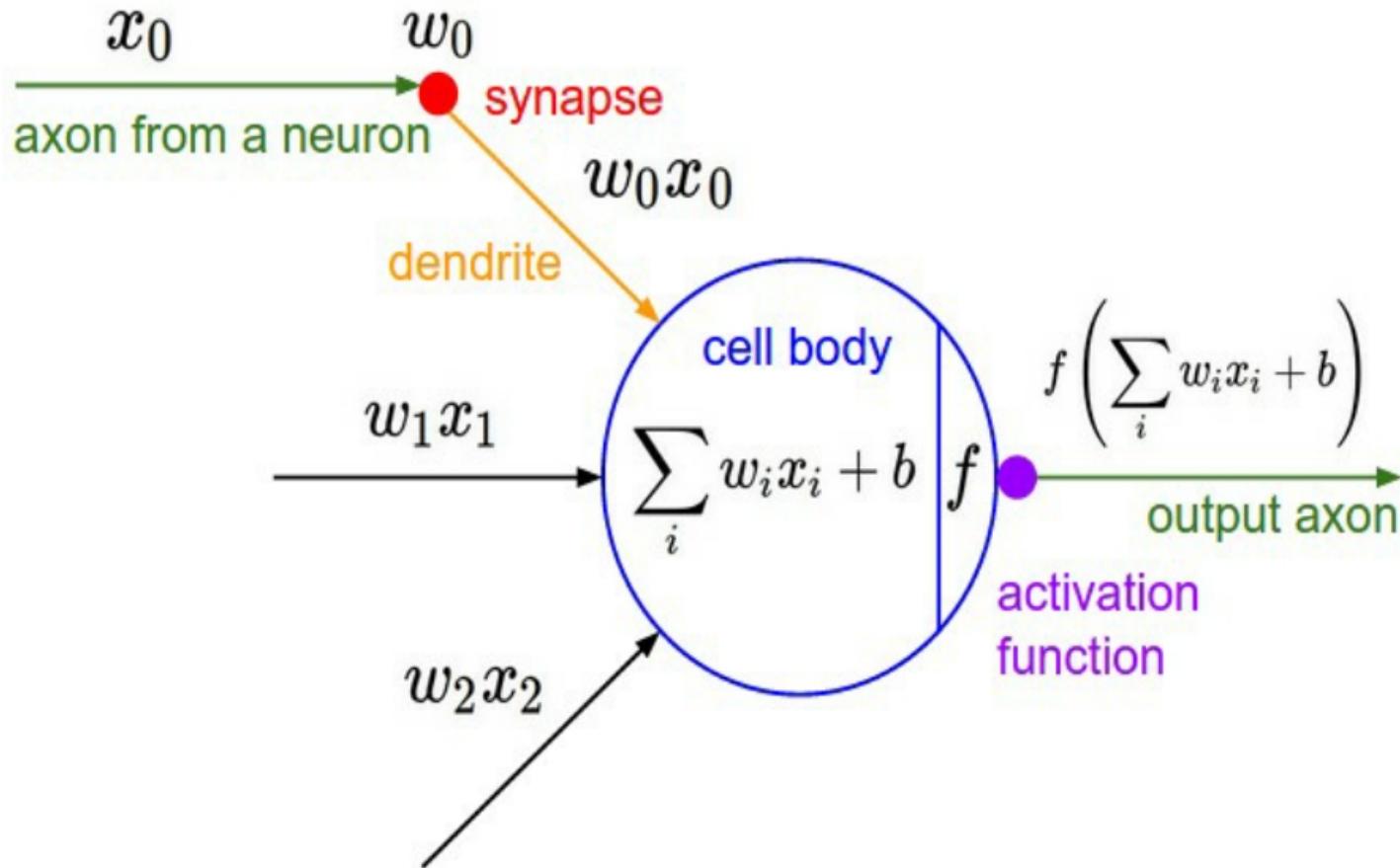


How Does the Brain Work?

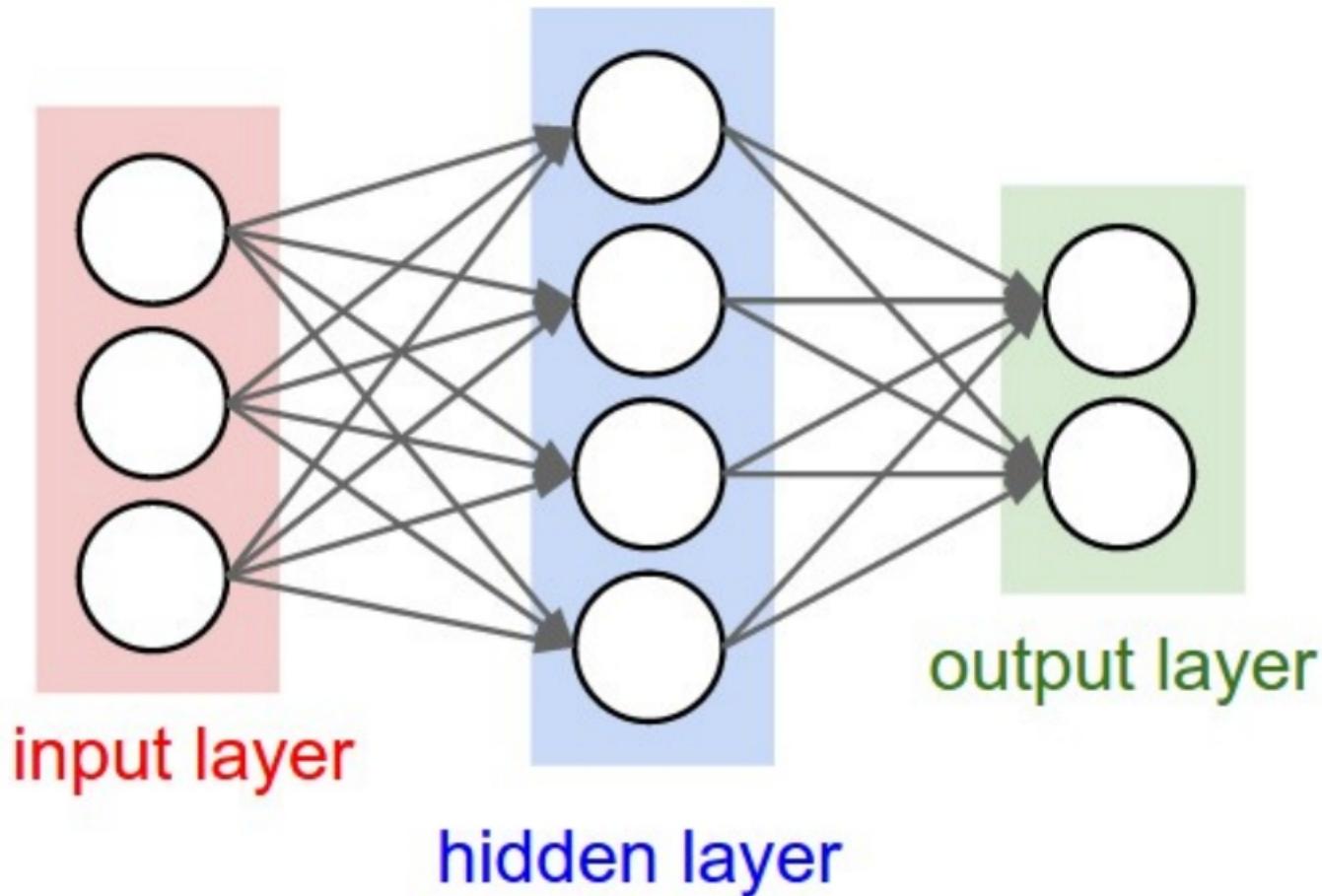


- The basic computational unit of the brain is a **neuron**
→ 86B neurons in the brain
- Neurons are connected with nearly **$10^{14} – 10^{15}$ synapses**
- Neurons receive input signal from **dendrites** and produce output signal along **axon**, which interact with the dendrites of other neurons via **synaptic weights**
- Synaptic weights – learnable & control influence strength

Neural Networks: Weighted Sum



Many Weighted Sums



What is Deep Learning?

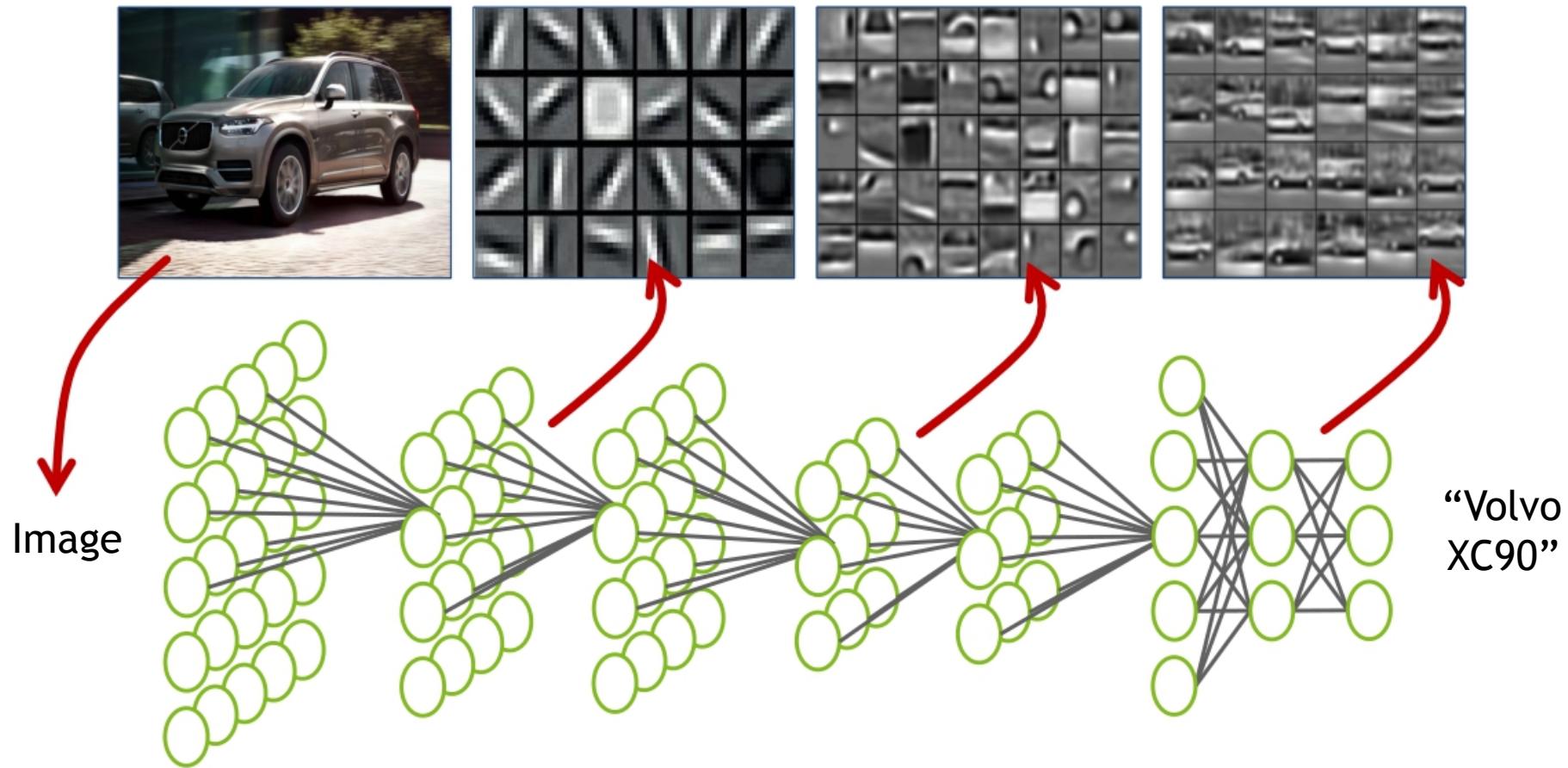


Image Source: [Lee et al., Comm. ACM 2011]

Why is Deep Learning Hot Now?

Big Data Availability

GPU Acceleration

New ML Techniques



350M images uploaded per day



2.5 Petabytes of customer data hourly



300 hours of video uploaded every minute



ImageNet Challenge

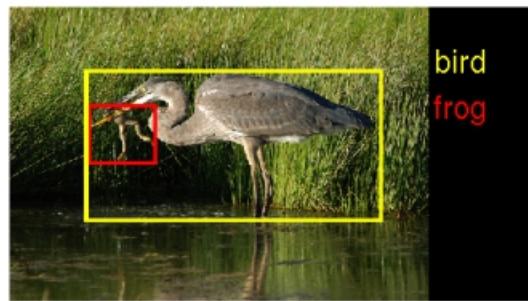
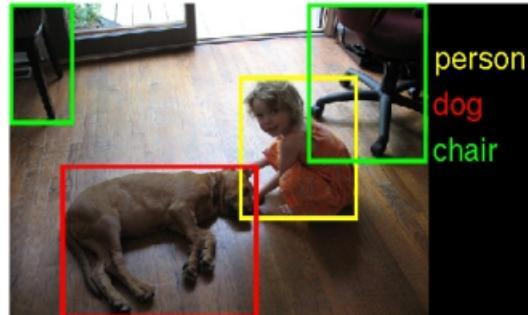


Image Classification Task:

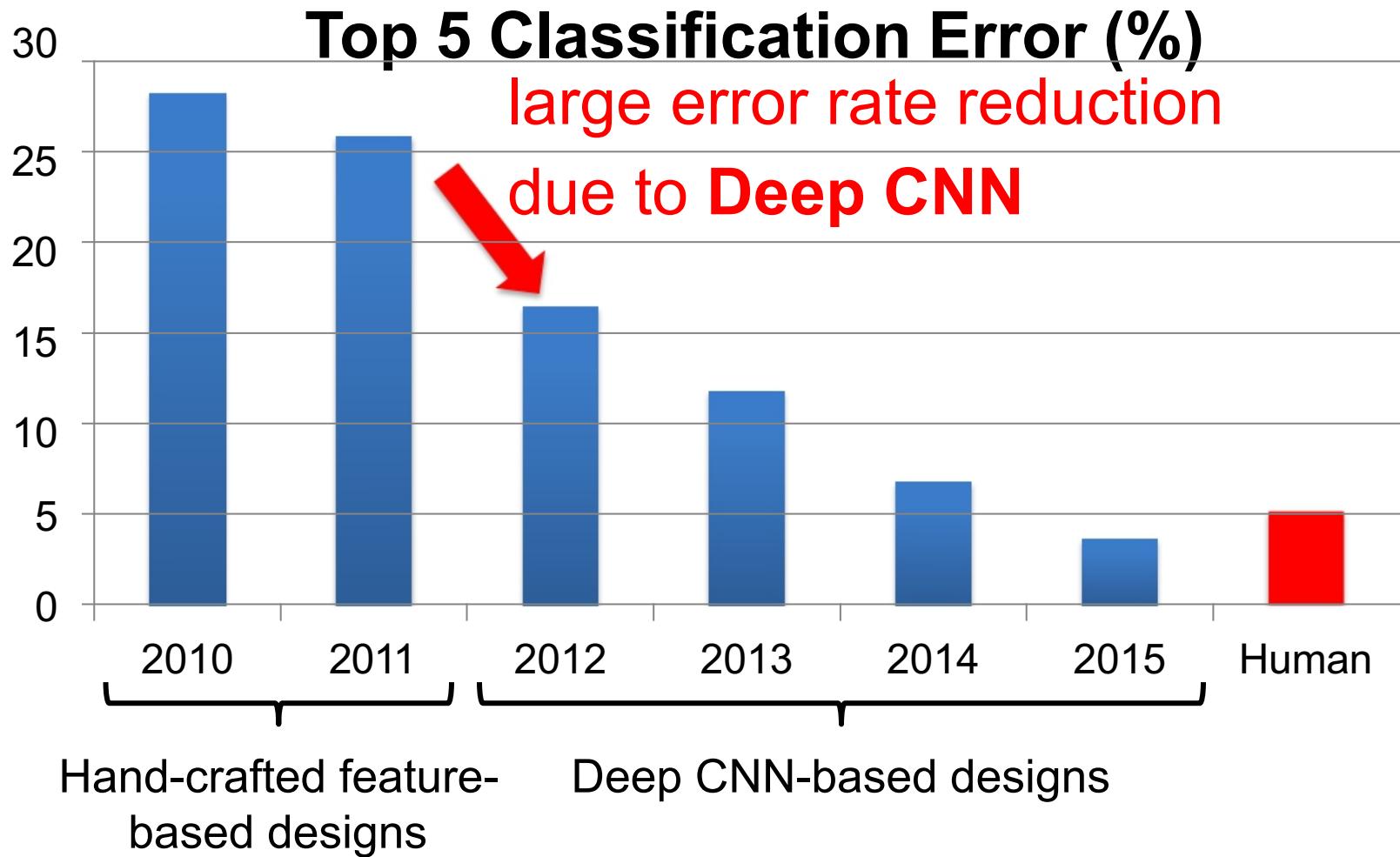
1.2M training images • 1000 object categories

Object Detection Task:

456k training images • 200 object categories



ImageNet: Image Classification Task



Established Applications

- **Image**
 - Classification: image to object class
 - Recognition: same as classification (except for faces)
 - Detection: assigning bounding boxes to objects
 - Segmentation: assigning object class to every pixel
- **Speech & Language**
 - Speech Recognition: audio to text
 - Translation
 - Natural Language Processing: text to meaning
 - Audio Generation: text to audio
- **Games**



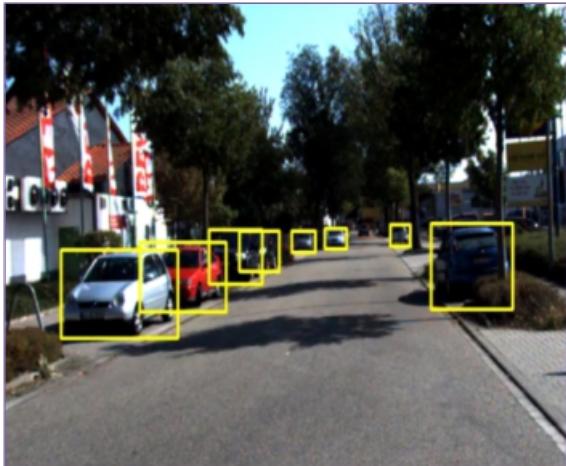
Deep Learning on Games

Google DeepMind AlphaGo



At last – a computer program that
can beat a champion Go player PAGE 404
ALL SYSTEMS GO

Deep Learning for Self-driving Cars



Overview of Deep Neural Networks

DNN Timeline

- **1940s: Neural networks were proposed**
- **1960s: Deep neural networks were proposed**
- **1989: Neural network for recognizing digits (LeNet)**
- **1990s: Hardware for shallow neural nets**
 - Example: Intel ETANN (1992)
- **2011: Breakthrough DNN-based speech recognition**
 - Microsoft real-time speech translation
- **2012: DNNs for vision supplanting traditional ML**
 - AlexNet for image classification
- **2014+: Rise of DNN accelerator research**
 - Examples: Neuflow, DianNao, etc.

So Many Neural Networks!

A mostly complete chart of

Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

○ Backfed Input Cell

○ Input Cell

△ Noisy Input Cell

● Hidden Cell

○ Probabilistic Hidden Cell

△ Spiking Hidden Cell

● Output Cell

○ Match Input Output Cell

● Recurrent Cell

○ Memory Cell

△ Different Memory Cell

● Kernel

○ Convolution or Pool

Perceptron (P)



Feed Forward (FF)



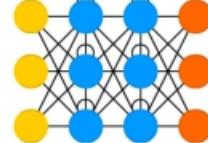
Radial Basis Network (RBF)



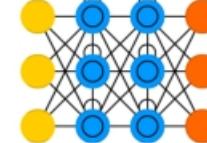
Deep Feed Forward (DFF)



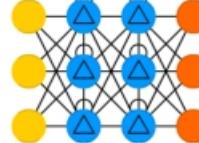
Recurrent Neural Network (RNN)



Long / Short Term Memory (LSTM)



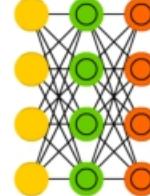
Gated Recurrent Unit (GRU)



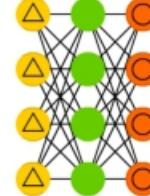
Auto Encoder (AE)



Variational AE (VAE)



Denoising AE (DAE)



Sparse AE (SAE)



Markov Chain (MC)



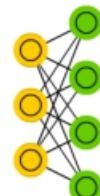
Hopfield Network (HN)



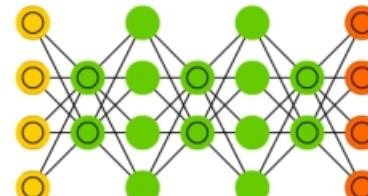
Boltzmann Machine (BM)



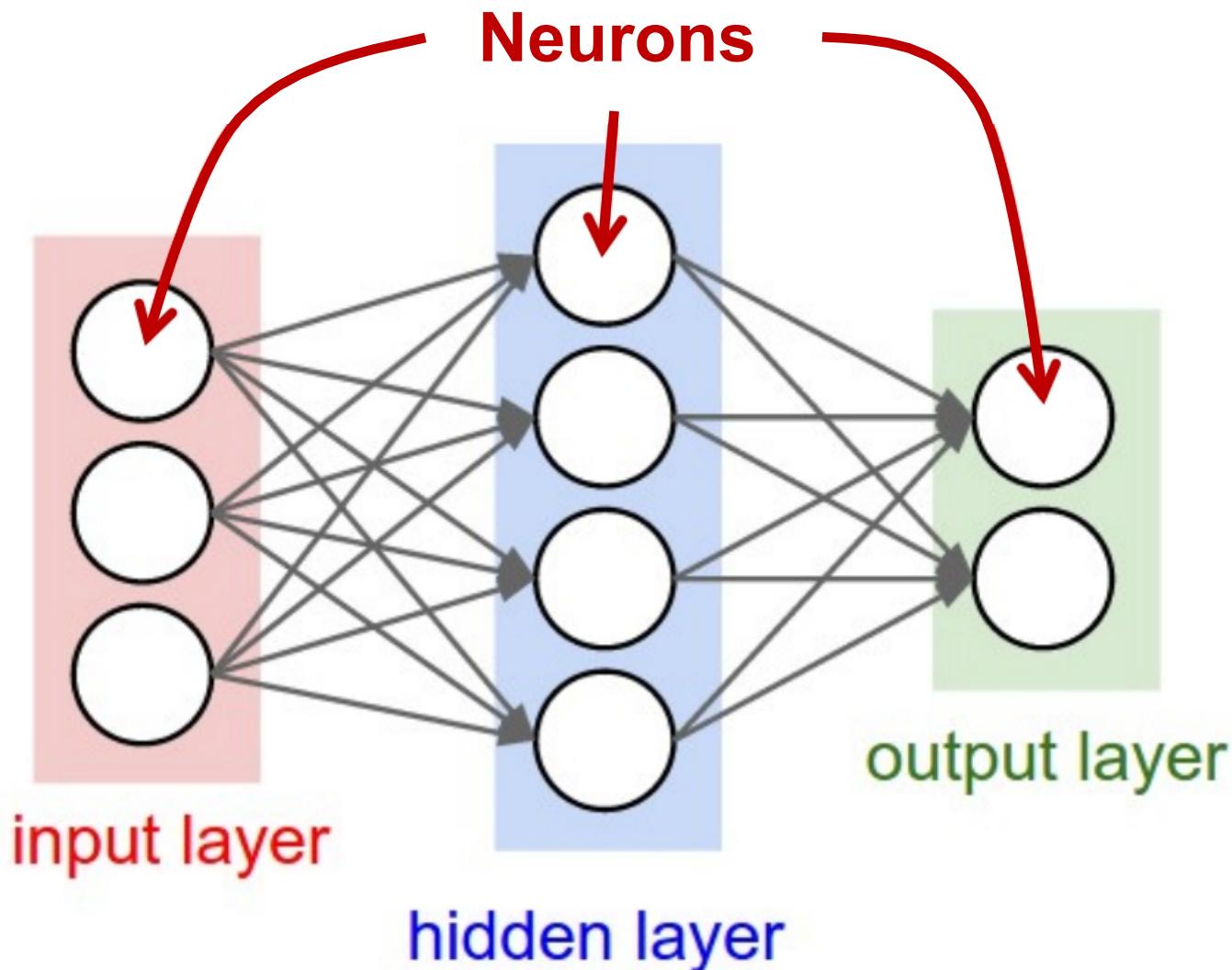
Restricted BM (RBM)



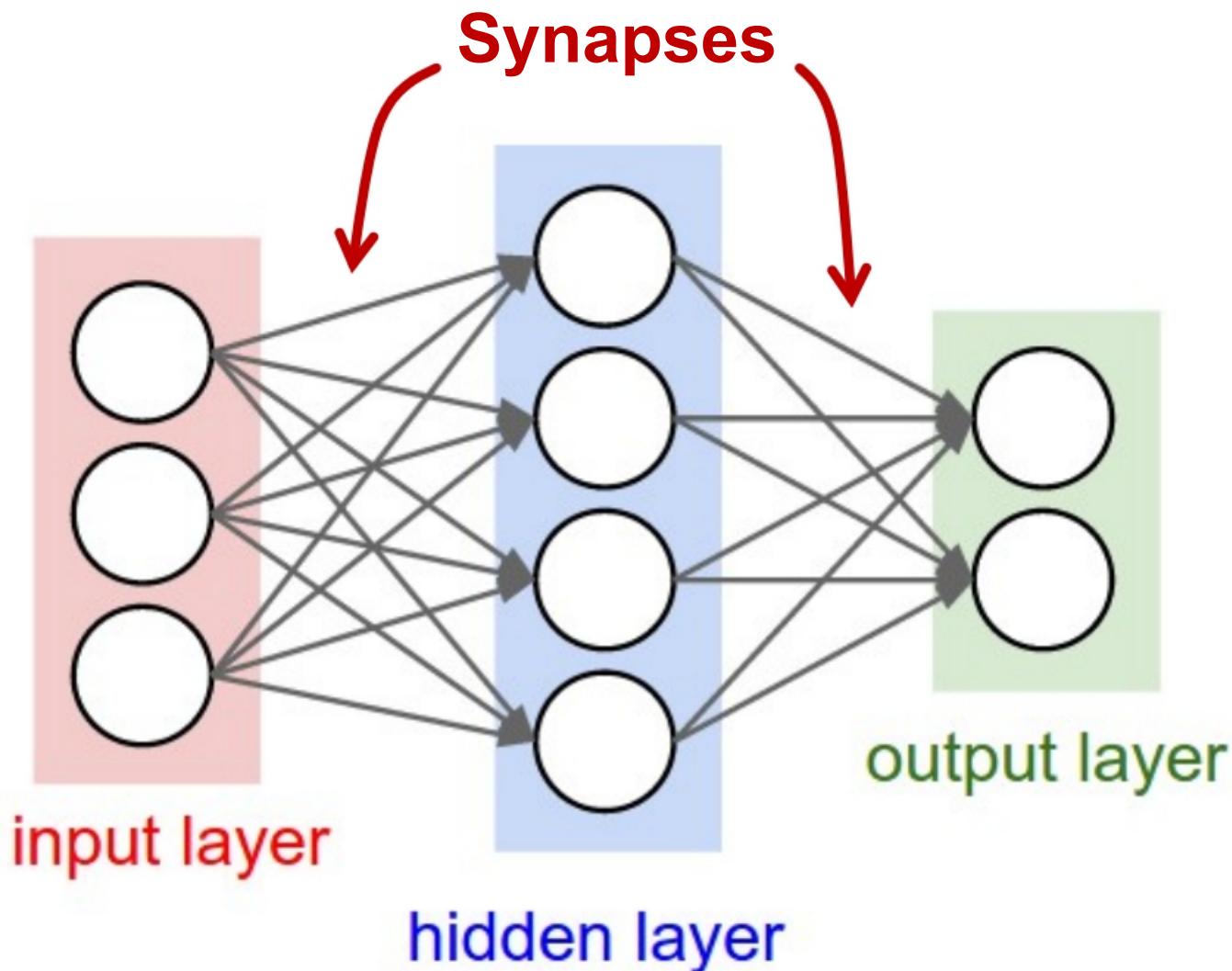
Deep Belief Network (DBN)



DNN Terminology 101

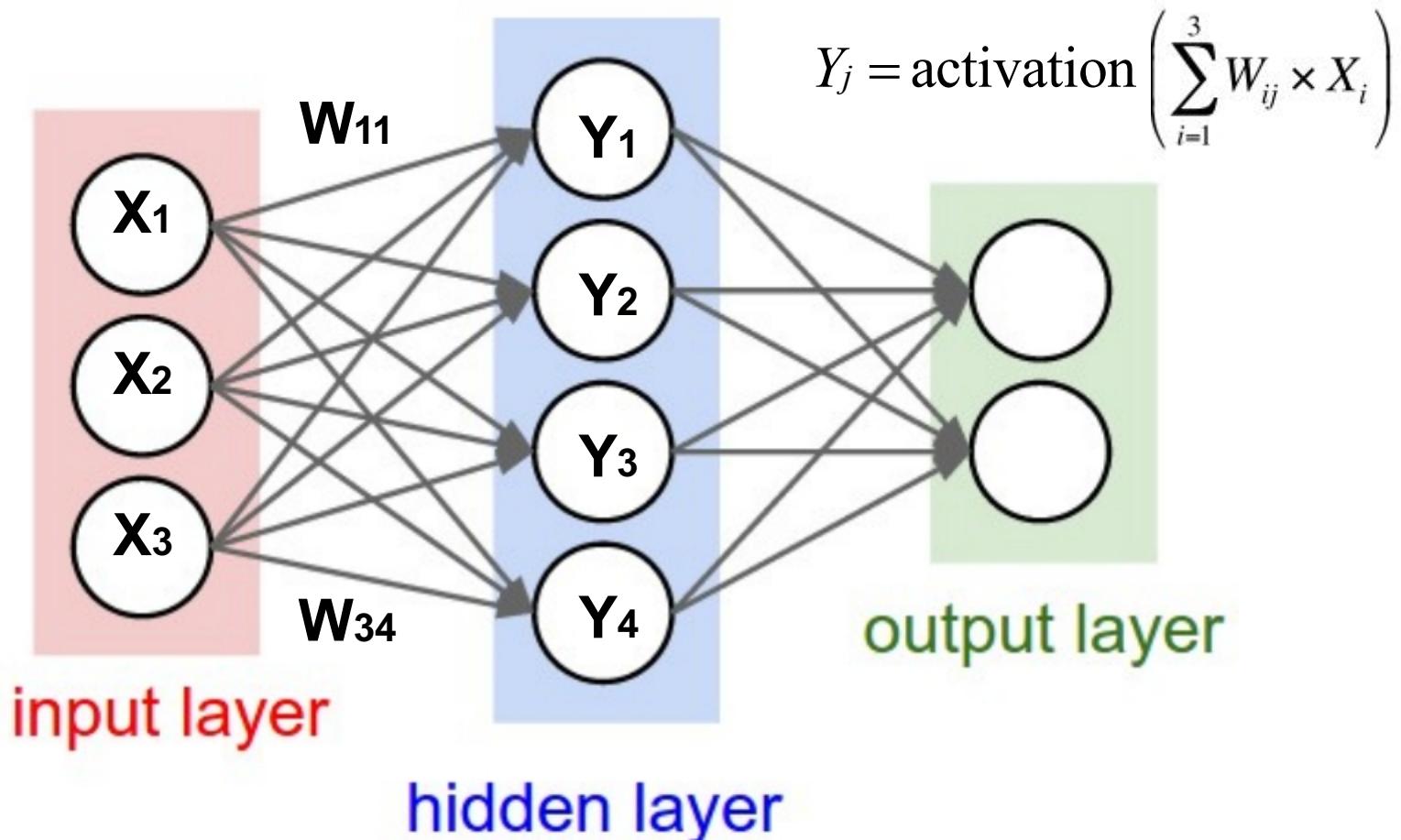


DNN Terminology 101



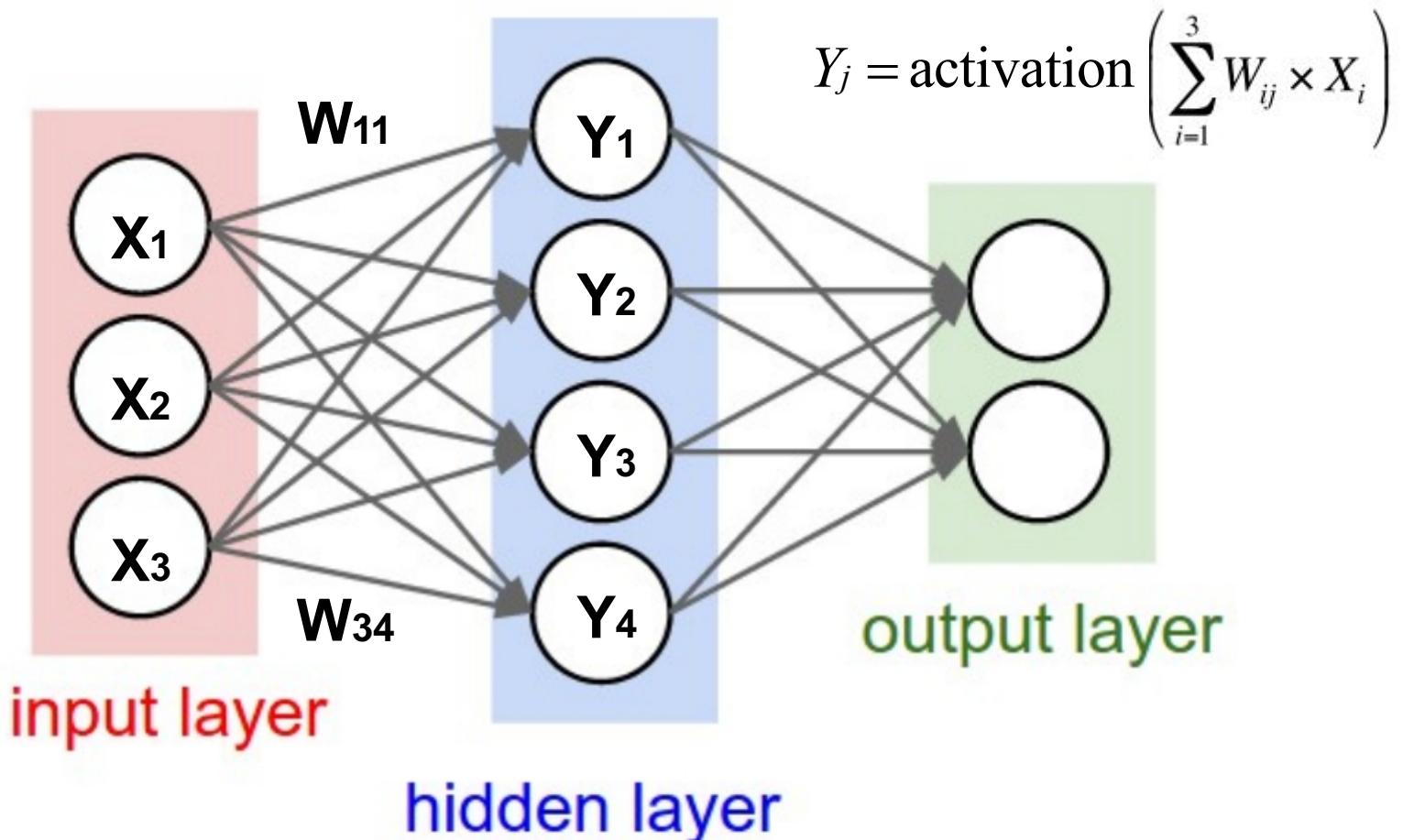
DNN Terminology 101

Each **synapse** has a **weight** for neuron **activation**

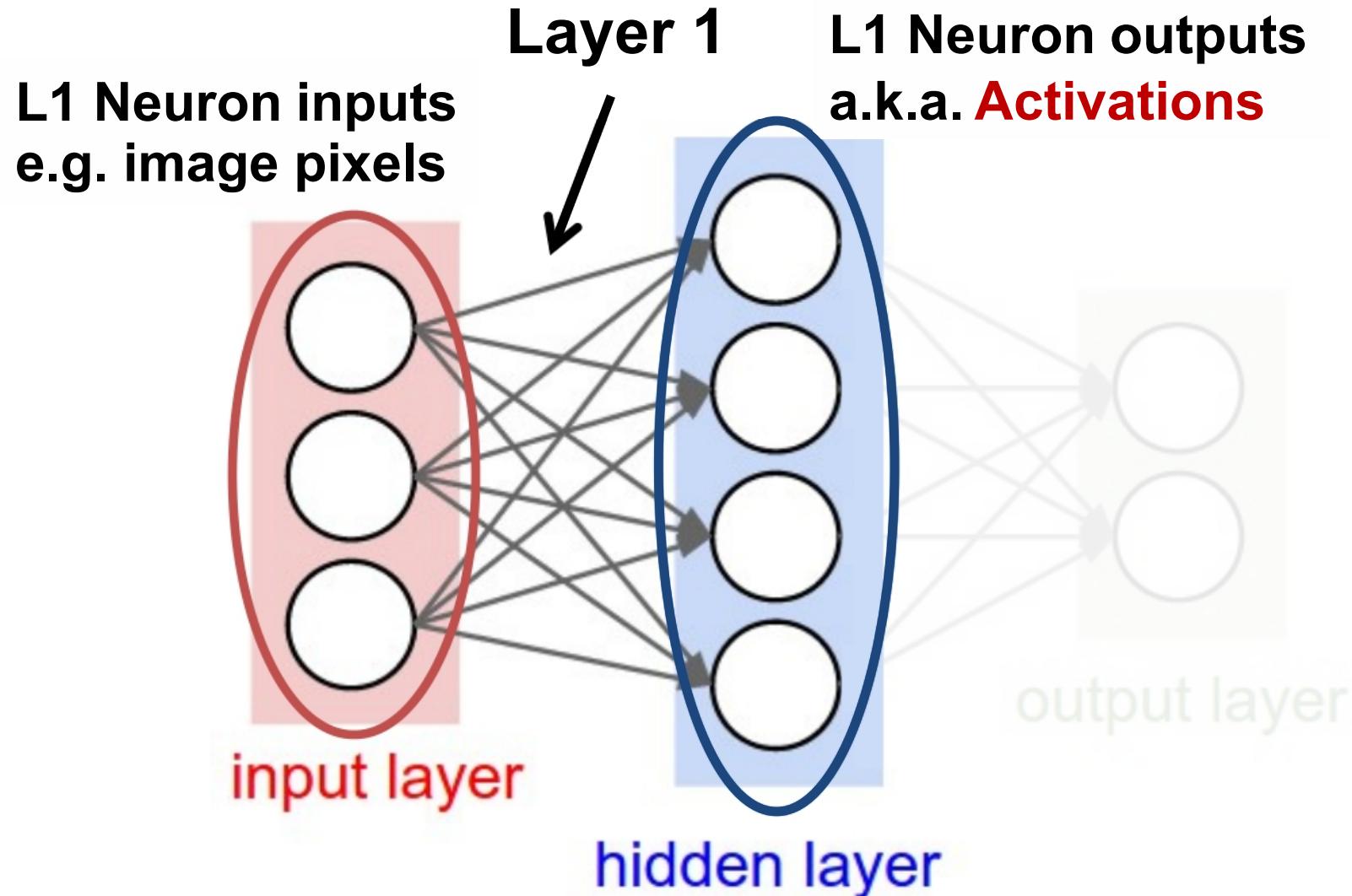


DNN Terminology 101

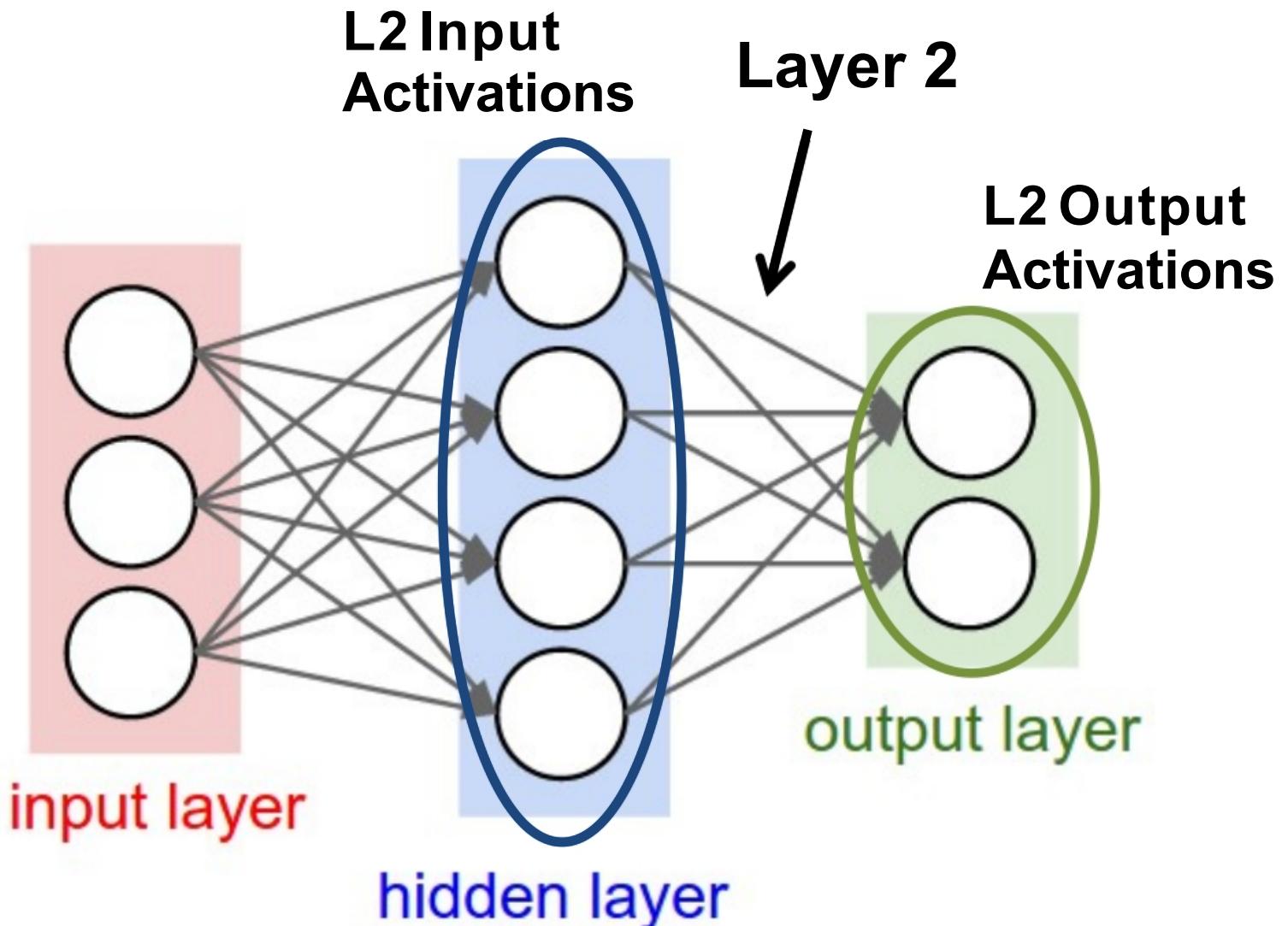
Weight Sharing: multiple synapses use the **same weight value**



DNN Terminology 101

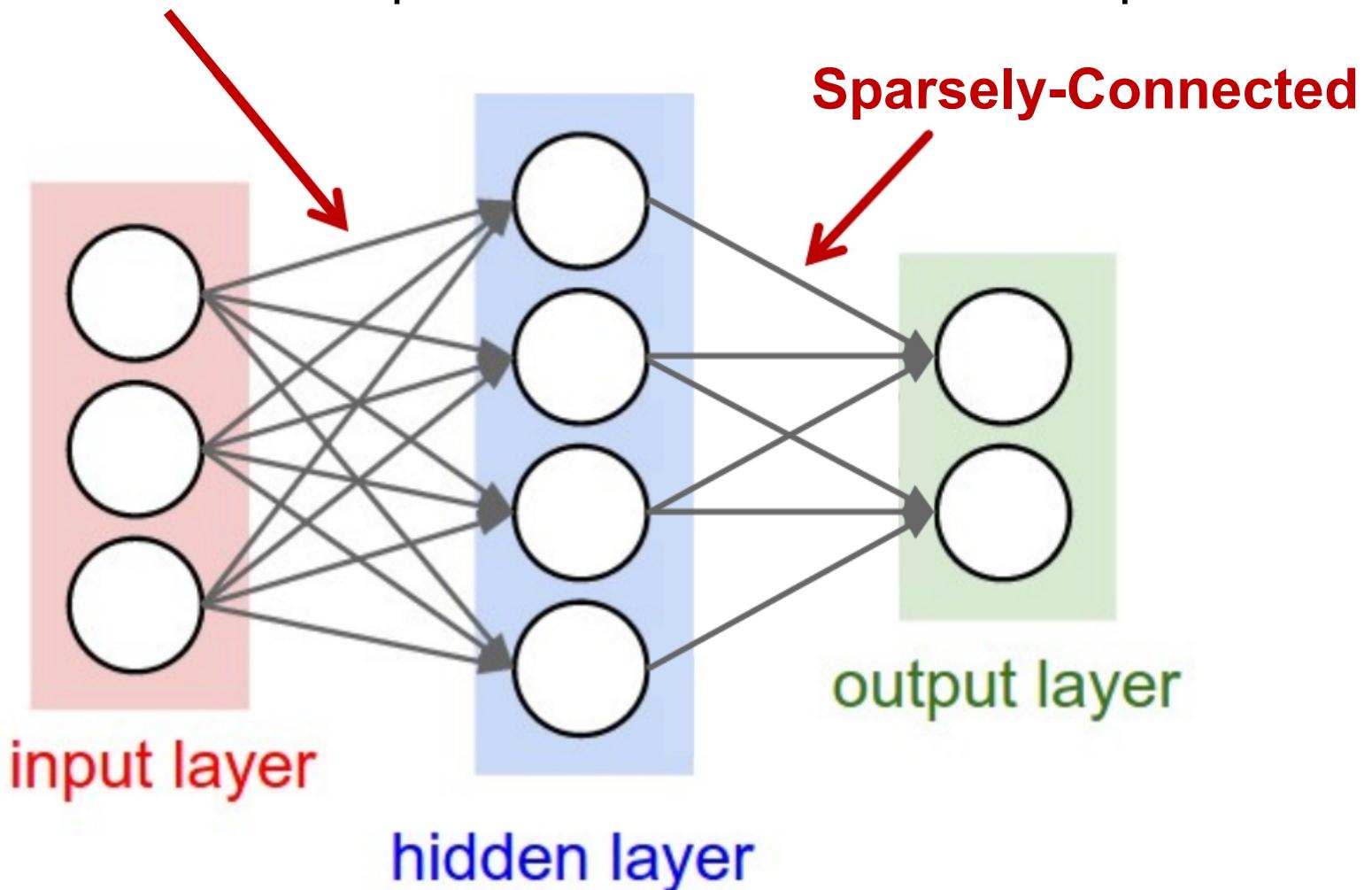


DNN Terminology 101

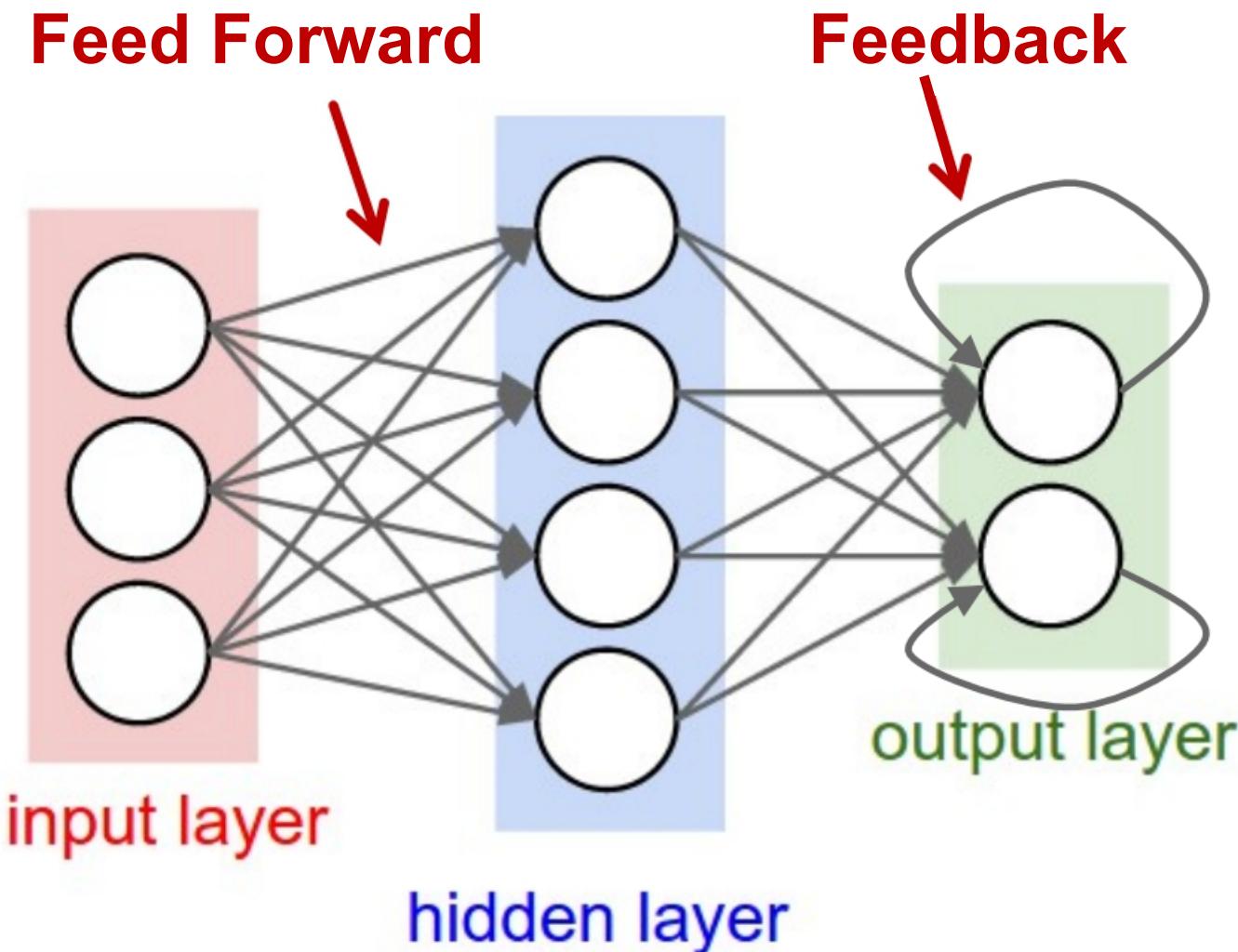


DNN Terminology 101

Fully-Connected: all i/p neurons connected to all o/p neurons



DNN Terminology 101



Popular Types of DNNs

- **Fully-Connected NN**
 - feed forward, a.k.a. multilayer perceptron (MLP)
- **Convolutional NN (CNN)**
 - feed forward, sparsely-connected w/ weight sharing
- **Recurrent NN (RNN)**
 - feedback
- **Long Short-Term Memory (LSTM)**
 - feedback + storage

Inference vs. Training

- **Training:** Determine weights
 - **Supervised:**
 - Training set has inputs and outputs, i.e., labeled
 - **Unsupervised:**
 - Training set is unlabeled
 - **Semi-supervised:**
 - Training set is partially labeled
 - **Reinforcement:**
 - Output assessed via rewards and punishments
- **Inference:** Apply weights to determine output

Convolutional Neural Network

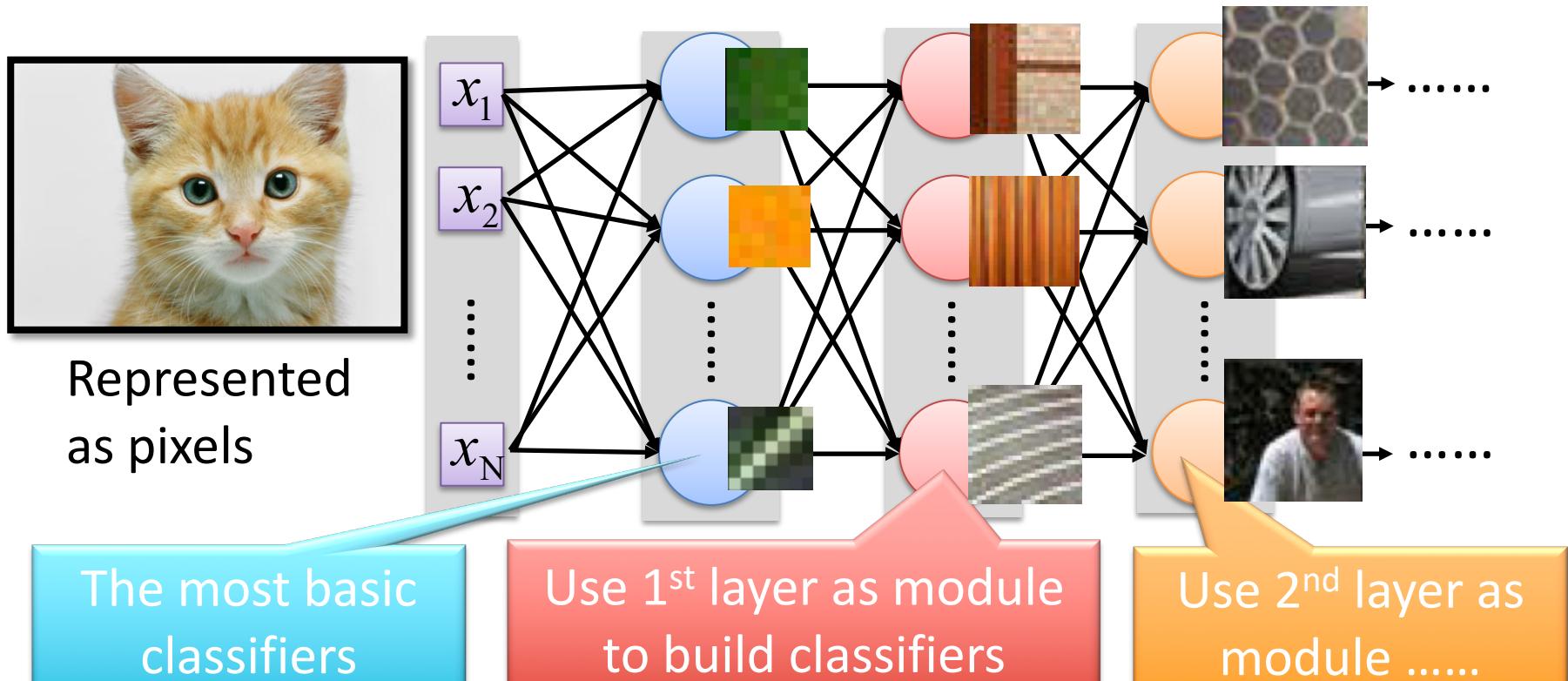
Slides are modified from Prof. Hung-yi
Lee's CNN tutorial

Available at

<http://speech.ee.ntu.edu.tw/~tlkagk/>

Why CNN for Image?

[Zeiler, M. D., ECCV 2014]



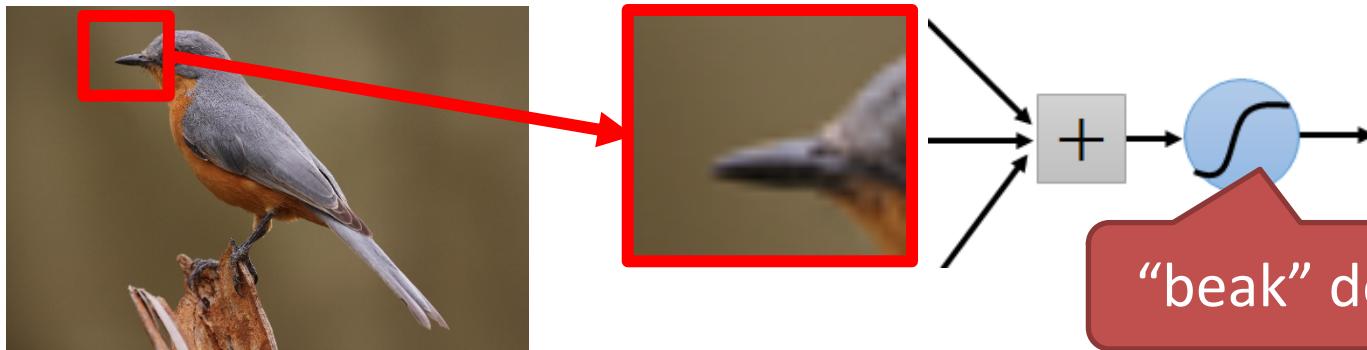
Can the network be simplified by considering the properties of images?

Why CNN for Image

- Some patterns are much smaller than the whole image

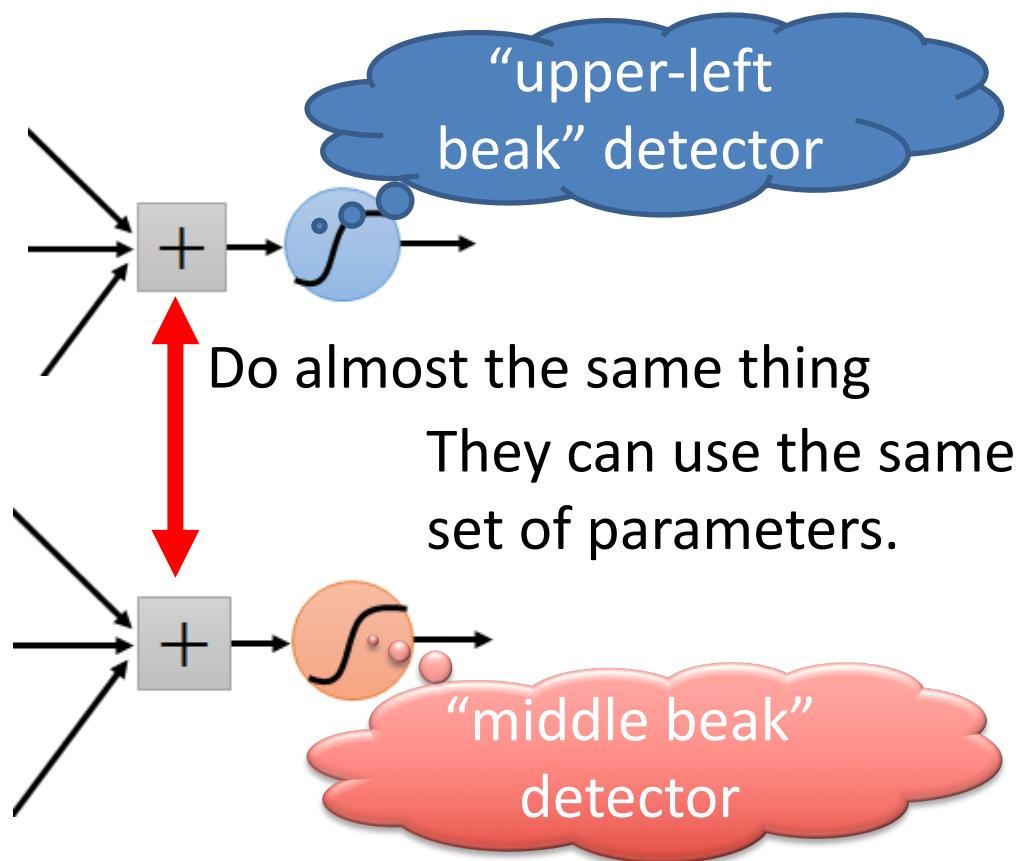
A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



Why CNN for Image

- The same patterns appear in different regions.



Why CNN for Image

- Subsampling the pixels will not change the object

bird



subsampling

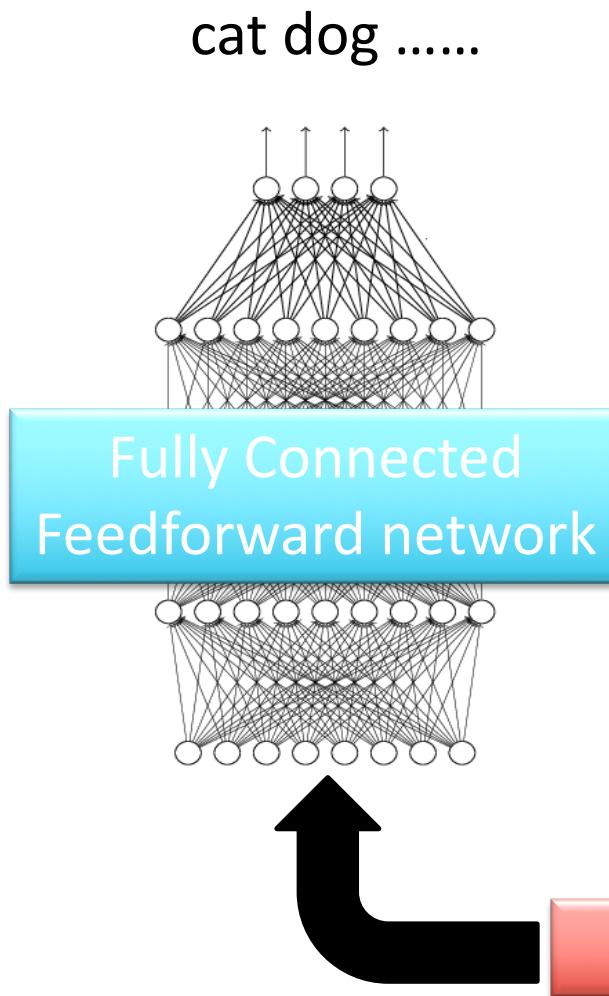
bird



We can subsample the pixels to make image smaller

→ Less parameters for the network to process the image

The whole CNN



Convolution

Max Pooling

Convolution

Max Pooling

Can repeat
many times

Flatten

The whole CNN



Property 1

- Some patterns are much smaller than the whole image

Property 2

- The same patterns appear in different regions.

Property 3

- Subsampling the pixels will not change the object

Convolution

Max Pooling

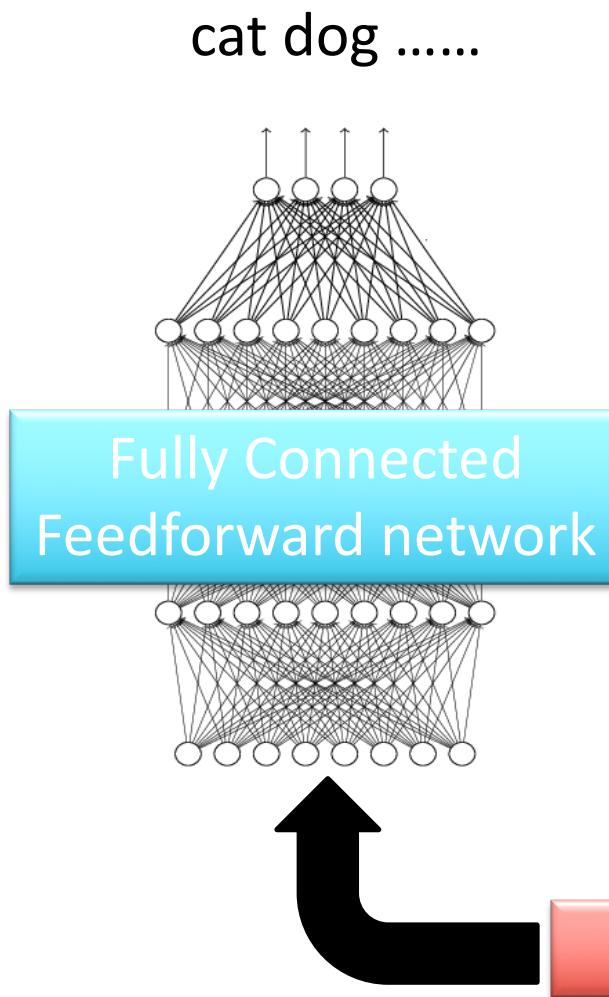
Convolution

Max Pooling

Flatten

Can repeat
many times

The whole CNN



Convolution

Max Pooling

Convolution

Max Pooling

Flatten

Can repeat
many times

CNN – Convolution

Those are the **network parameters** to be learned.

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1
Matrix

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2
Matrix

⋮

Property 1

Each filter detects a small pattern (3 x 3).

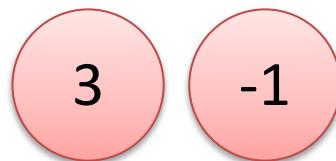
CNN – Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0



6 x 6 image

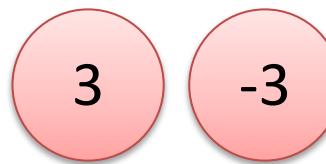
CNN – Convolution

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

If stride=2

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

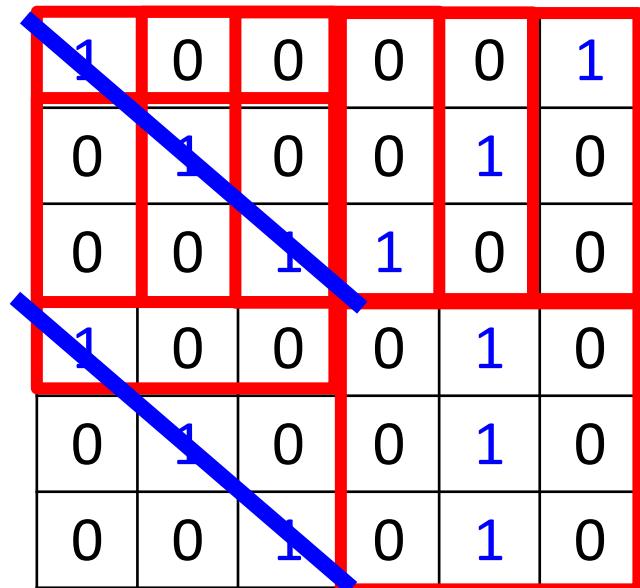


We set stride=1 below

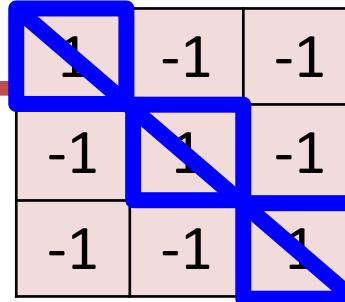
6 x 6 image

CNN – Convolution

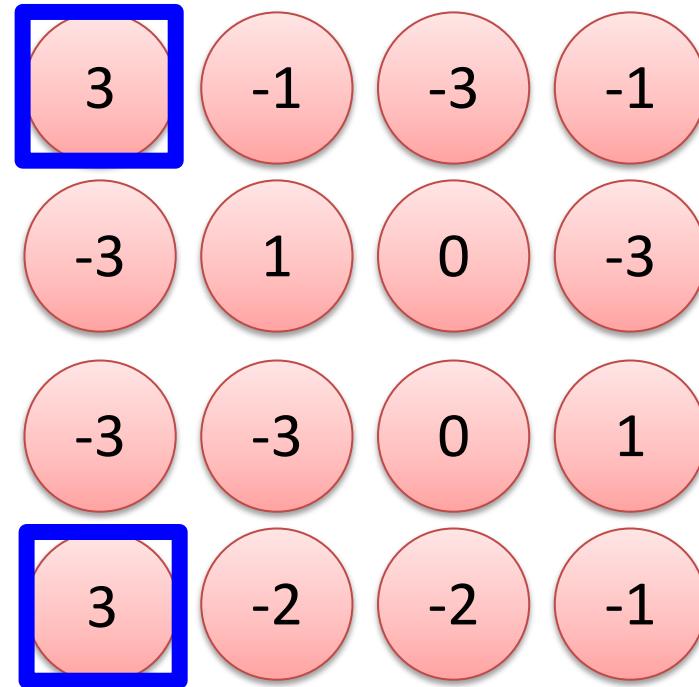
stride=1



6 x 6 image



Filter 1



Property 2

CNN – Convolution

stride=1

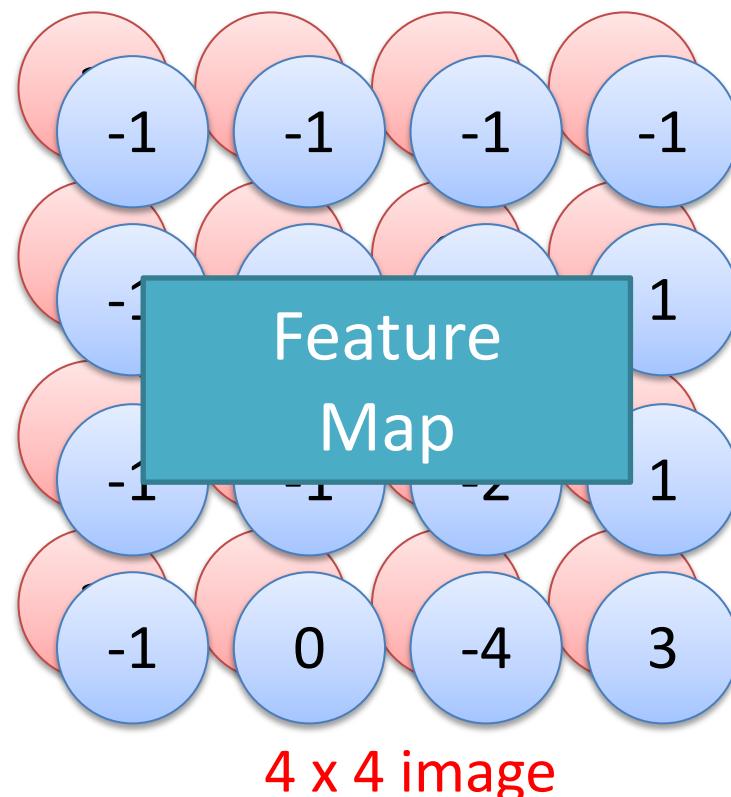
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

-1	1	-1
-1	1	-1
-1	1	-1

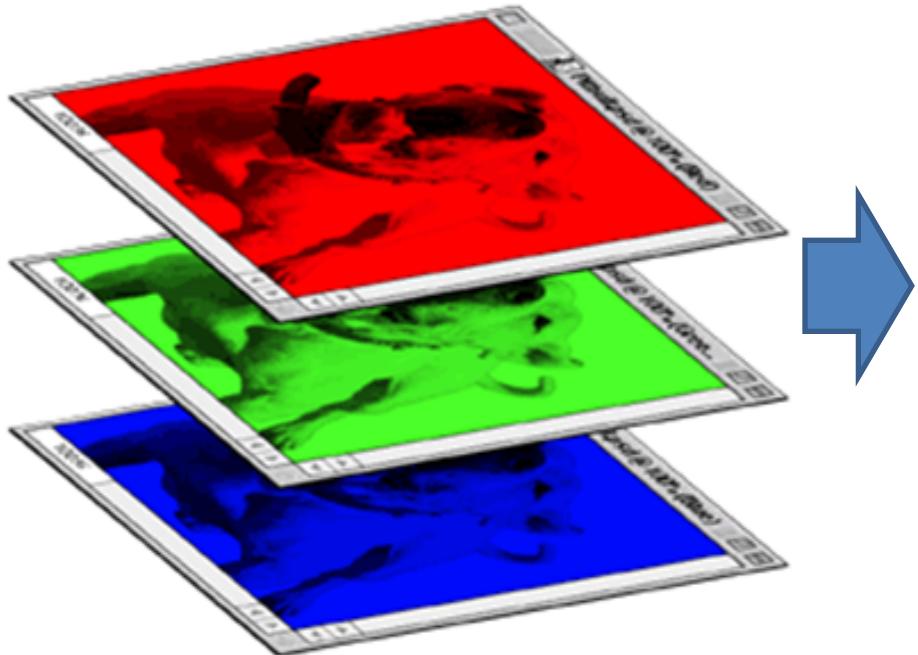
Filter 2

Do the same process for every filter



CNN – Colorful image

Colorful image



A diagram of a convolutional filter kernel. It is a 3x3 grid of weights with values: top row [1, -1, -1], middle row [-1, 1, -1], bottom row [-1, -1, 1]. The kernel is surrounded by a pink border.

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

A diagram of a convolutional filter kernel. It is a 3x3 grid of weights with values: top row [-1, 1, -1], middle row [-1, 1, -1], bottom row [-1, 1, -1]. The kernel is surrounded by a light blue border.

-1	1	-1
-1	1	-1
-1	1	-1

Filter 2

A diagram of the output feature map. It is a 7x6 grid of binary values (0 or 1). Blue numbers highlight specific values: the first column has values 1, 0, 0, 1, 0, 1, 0; the second column has values 0, 1, 0, 0, 1, 0, 1; the third column has values 0, 0, 1, 1, 0, 0, 1; the fourth column has values 1, 0, 0, 0, 1, 0, 1; the fifth column has values 0, 1, 0, 0, 1, 0, 1; the sixth column has values 0, 0, 1, 0, 1, 0, 1. The feature map is surrounded by a black border.

1	0	0	0	0	1
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

Convolution v.s. Fully Connected

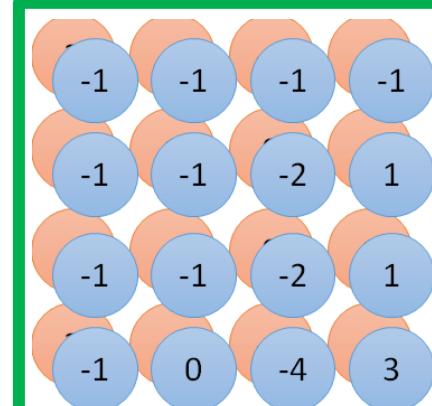
1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

image

1	-1	-1
-1	1	-1
-1	-1	1

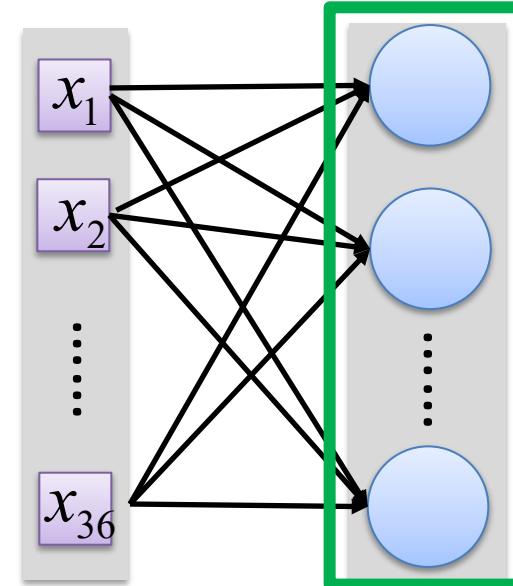
-1	1	-1
-1	1	-1
-1	1	-1

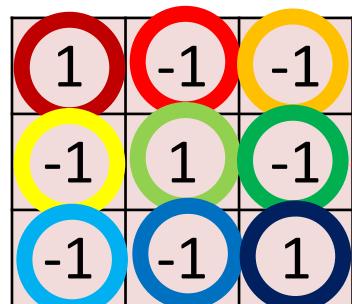
convolution



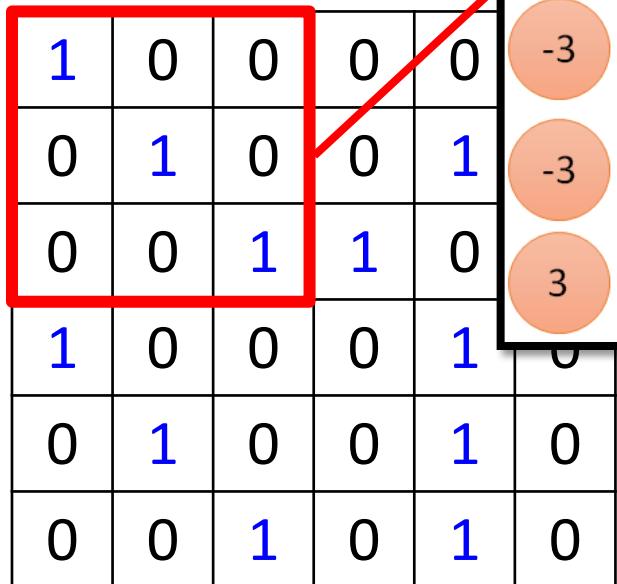
Fully-
connected

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0



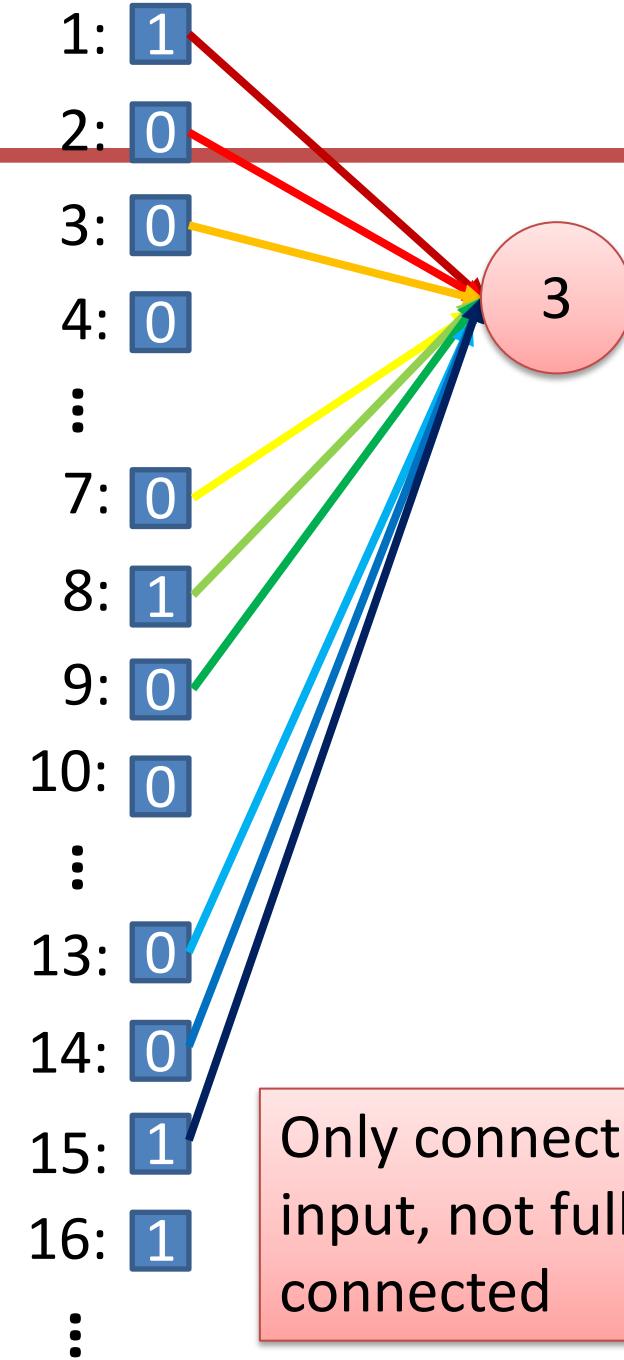
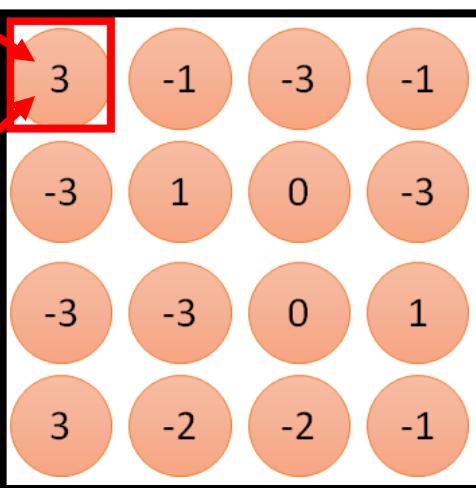


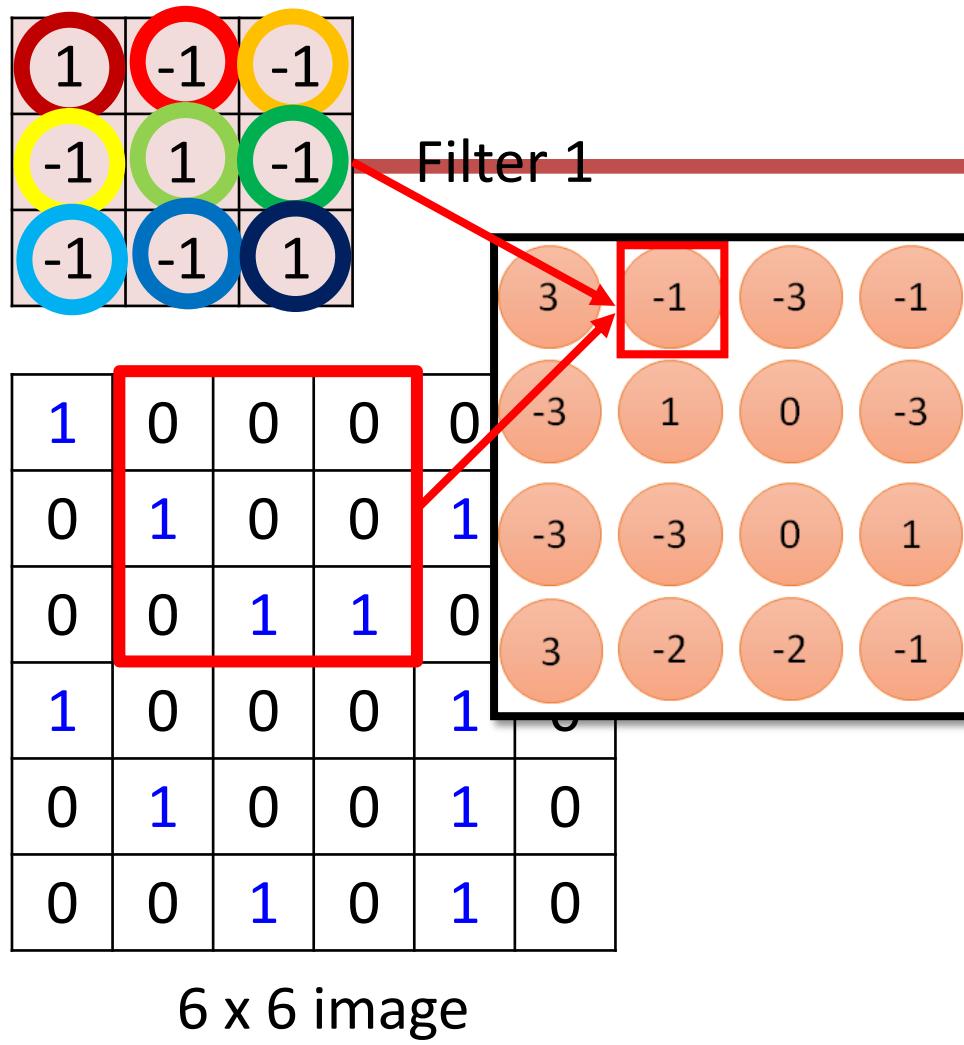
Filter 1



6×6 image

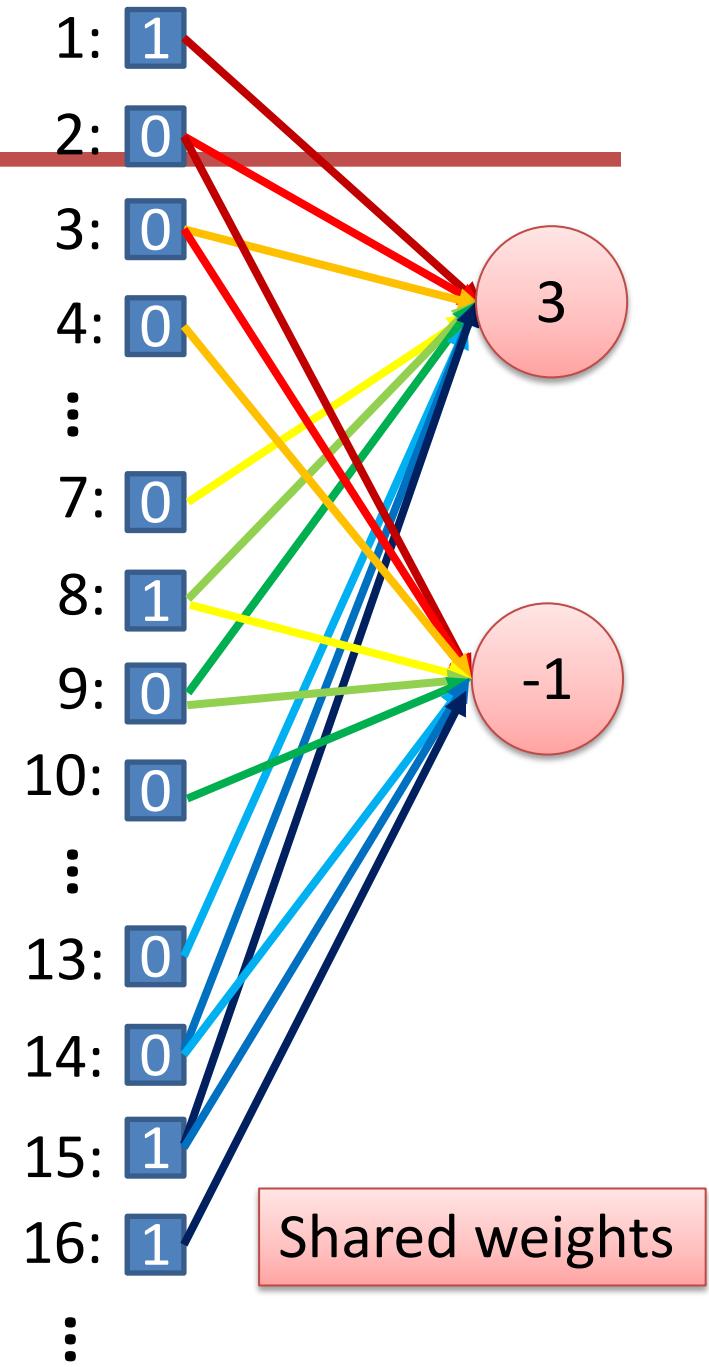
Less parameters!



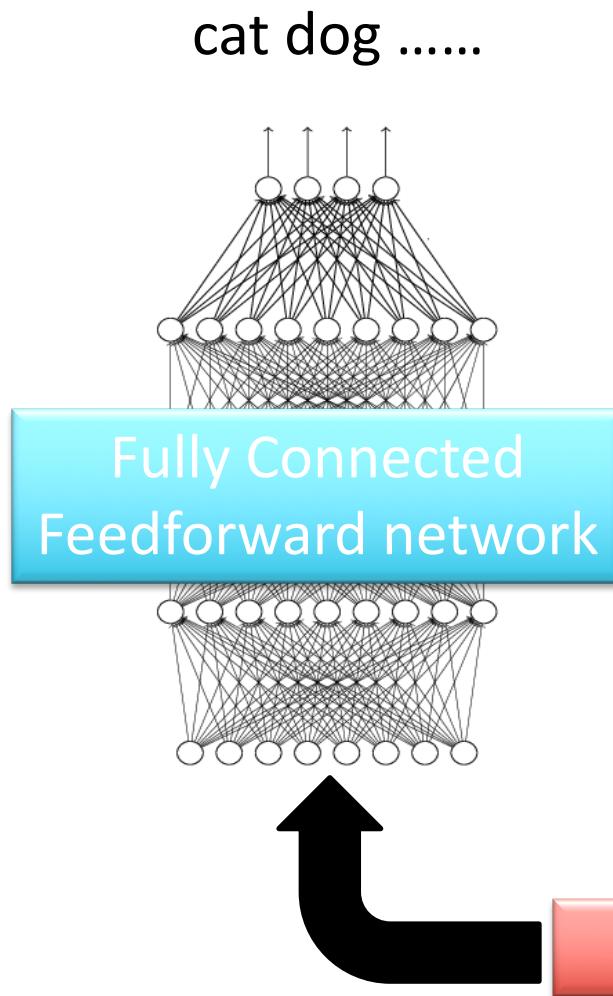


Less parameters!

Even less parameters!



The whole CNN



Convolution

Max Pooling

Convolution

Max Pooling

Can repeat
many times

Flatten

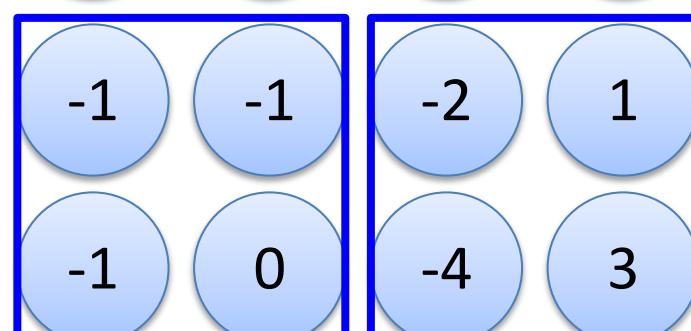
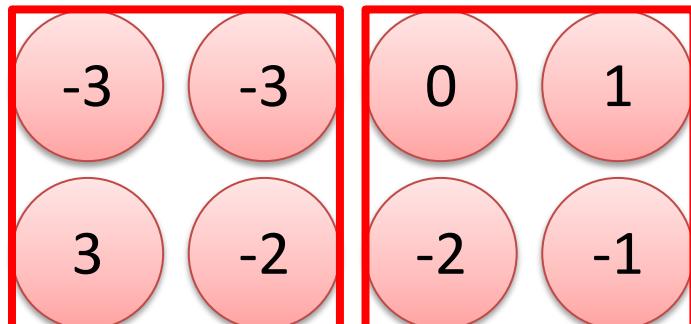
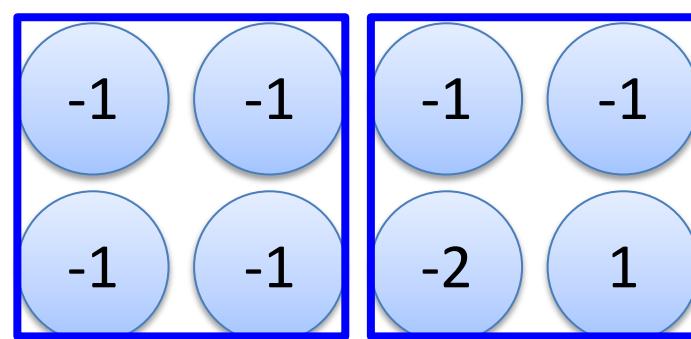
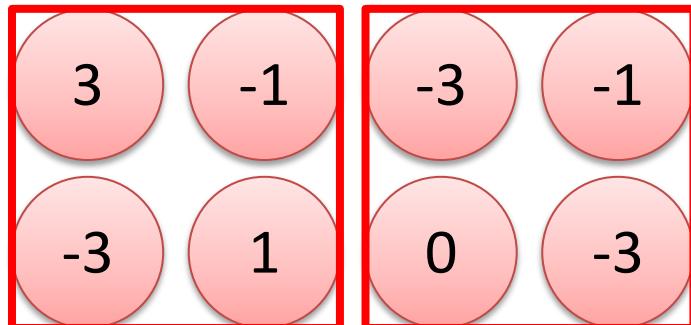
CNN – Max Pooling

1	-1	-1
-1	1	-1
-1	-1	1

Filter 1

-1	1	-1
-1	1	-1
-1	1	-1

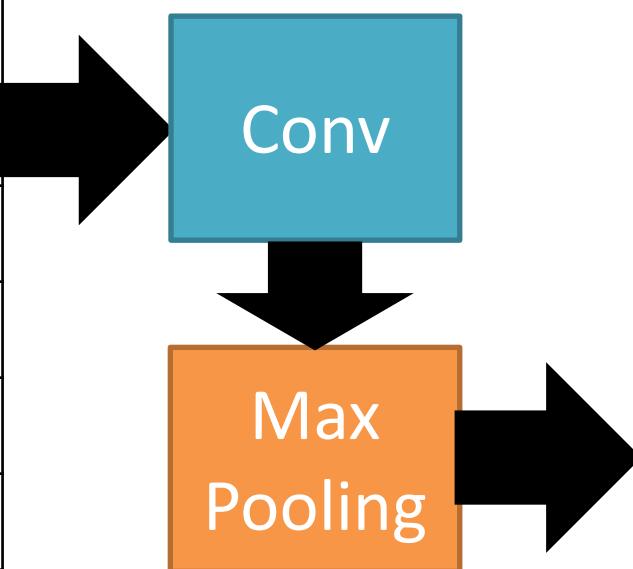
Filter 2



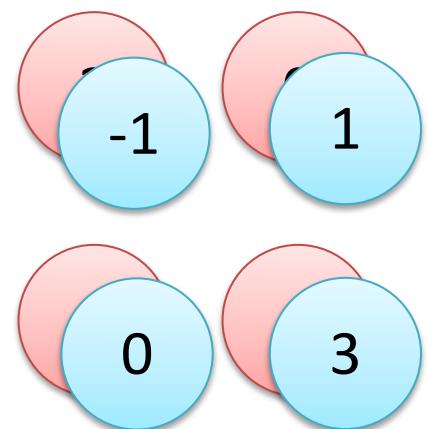
CNN – Max Pooling

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image



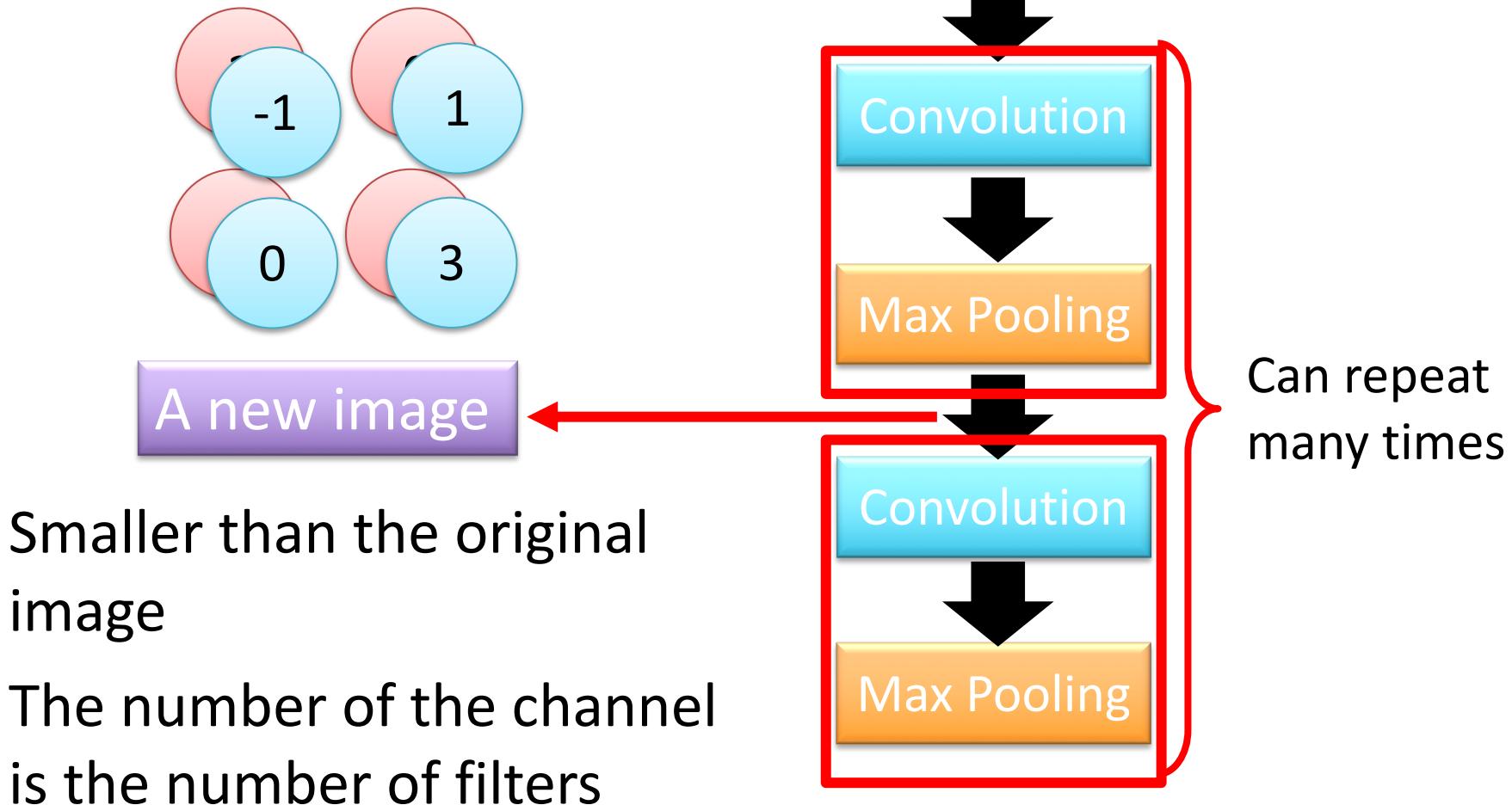
New image
but smaller



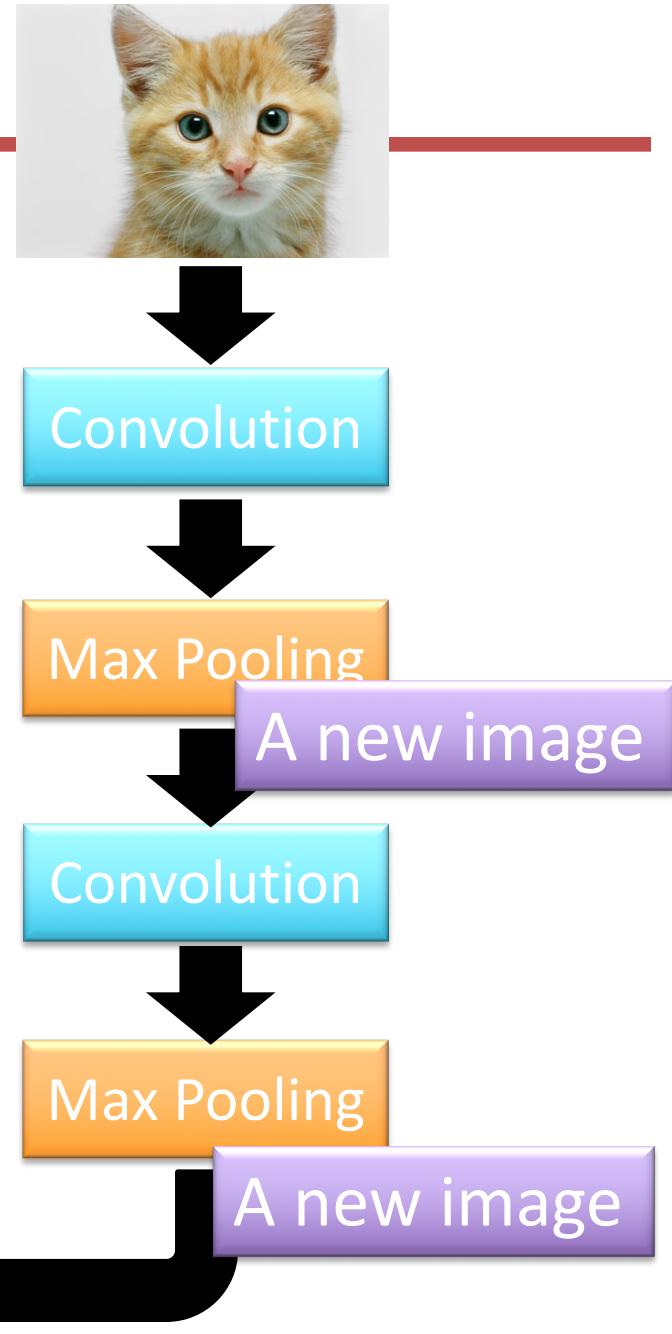
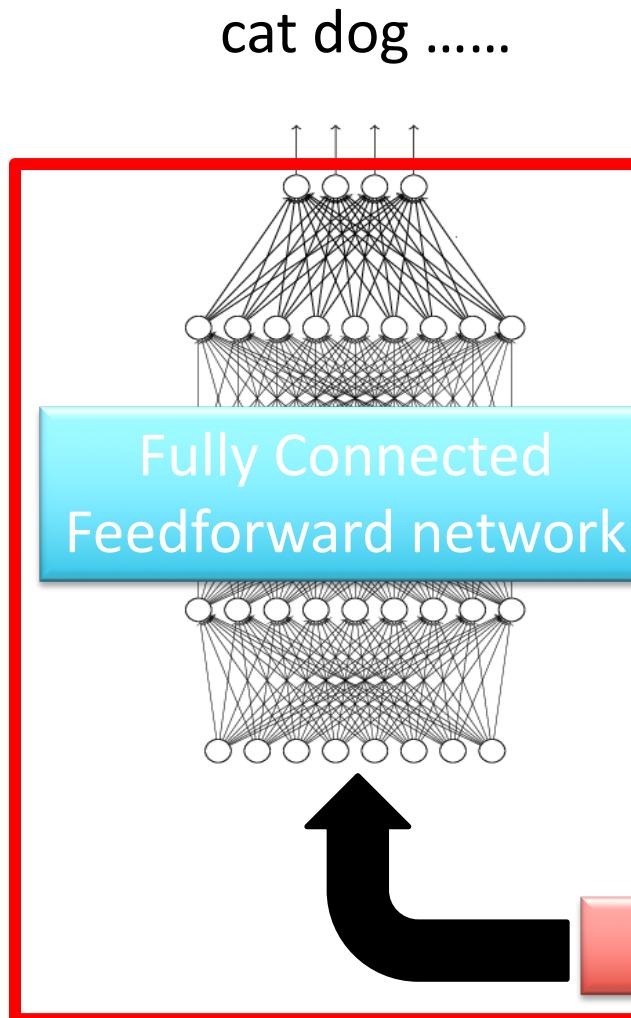
2 x 2 image

Each filter
is a channel

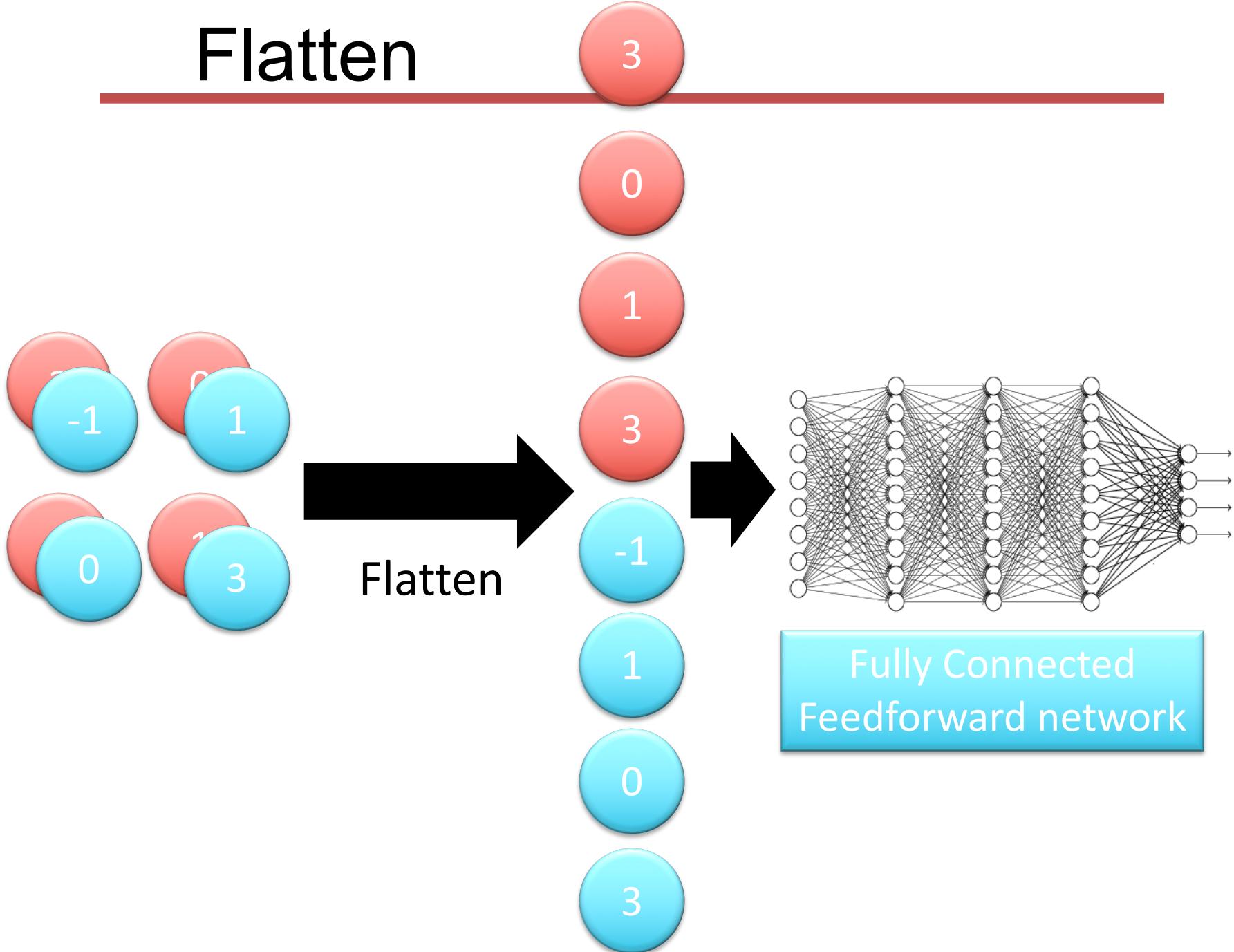
The whole CNN



The whole CNN

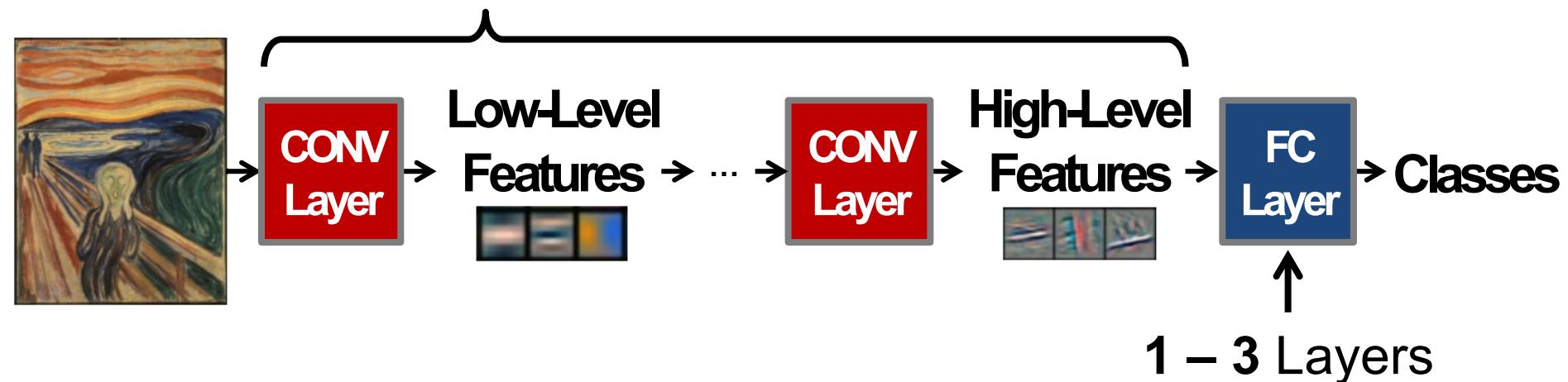


Flatten

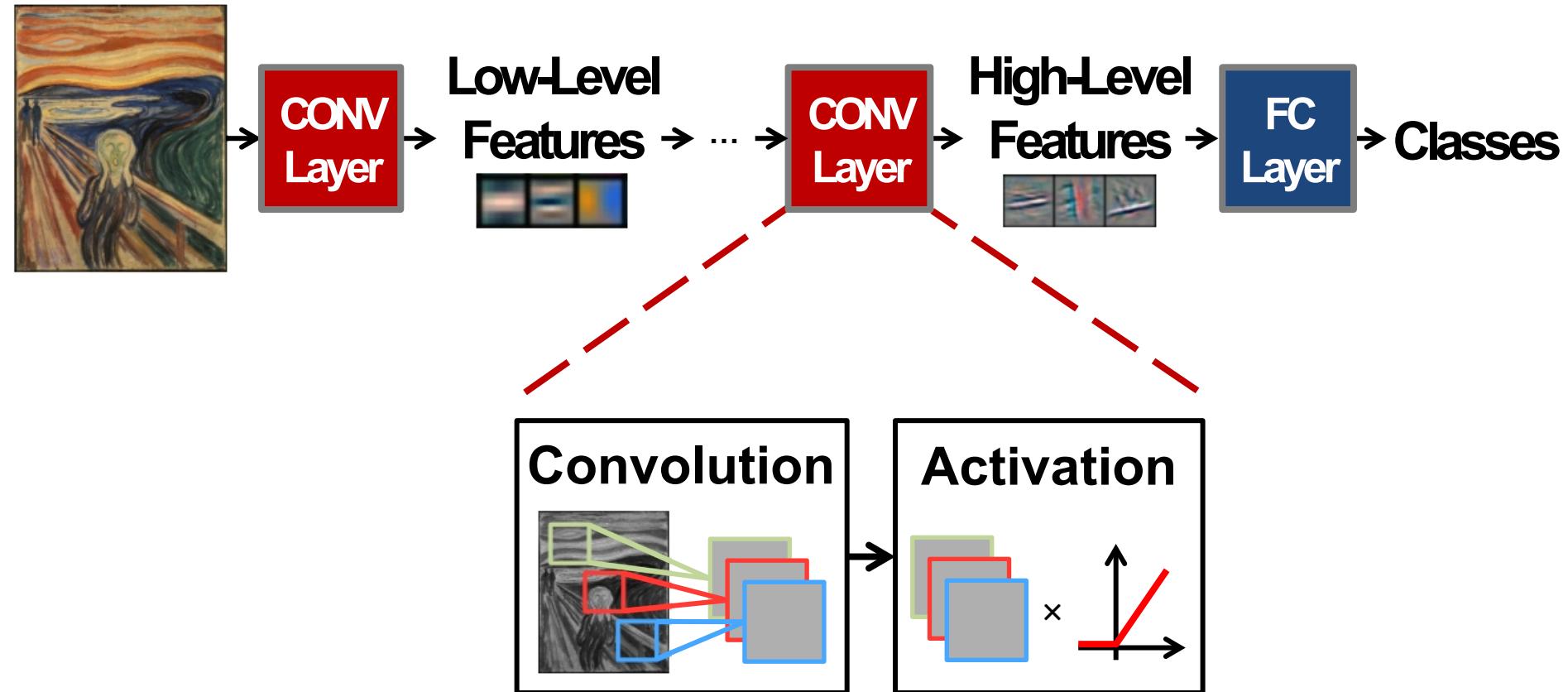


Deep Convolutional Neural Networks

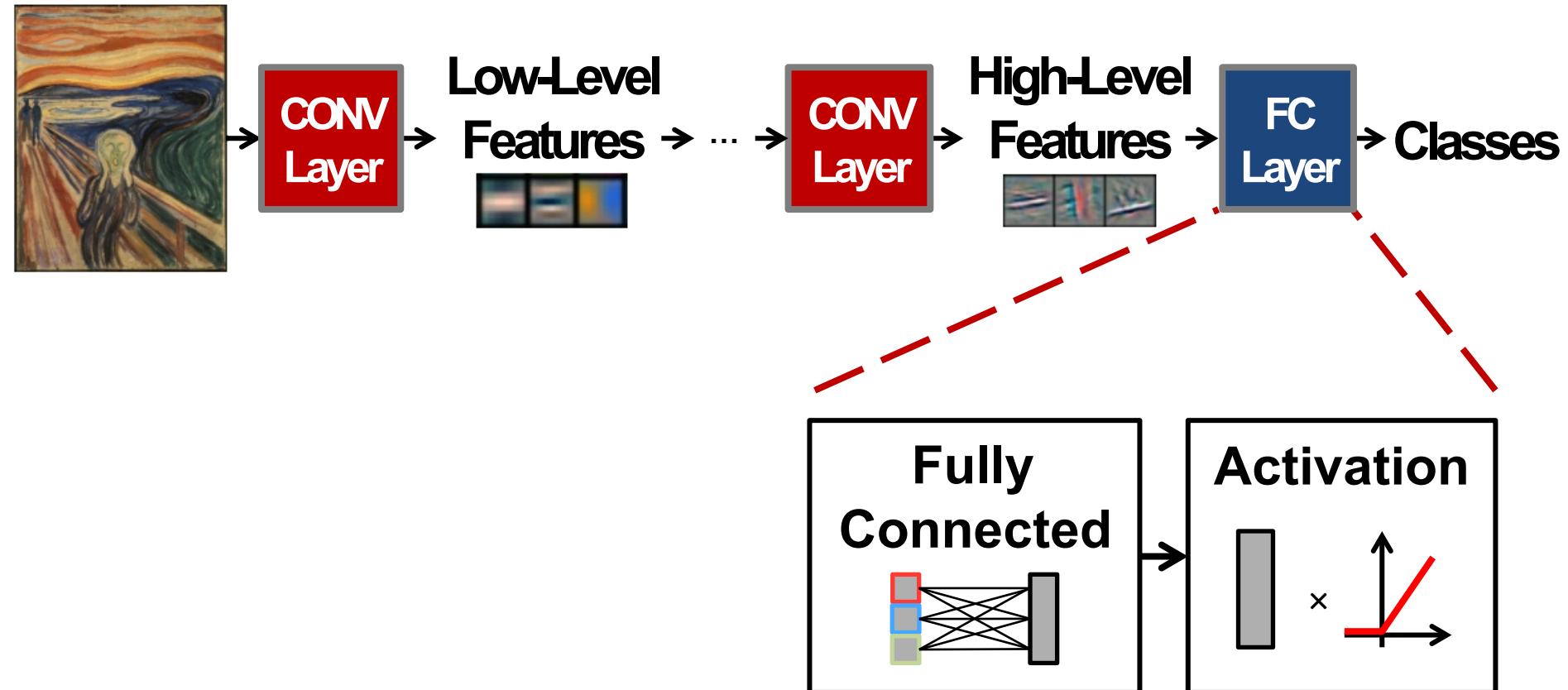
Modern Deep CNN: 5 – 1000 Layers



Deep Convolutional Neural Networks

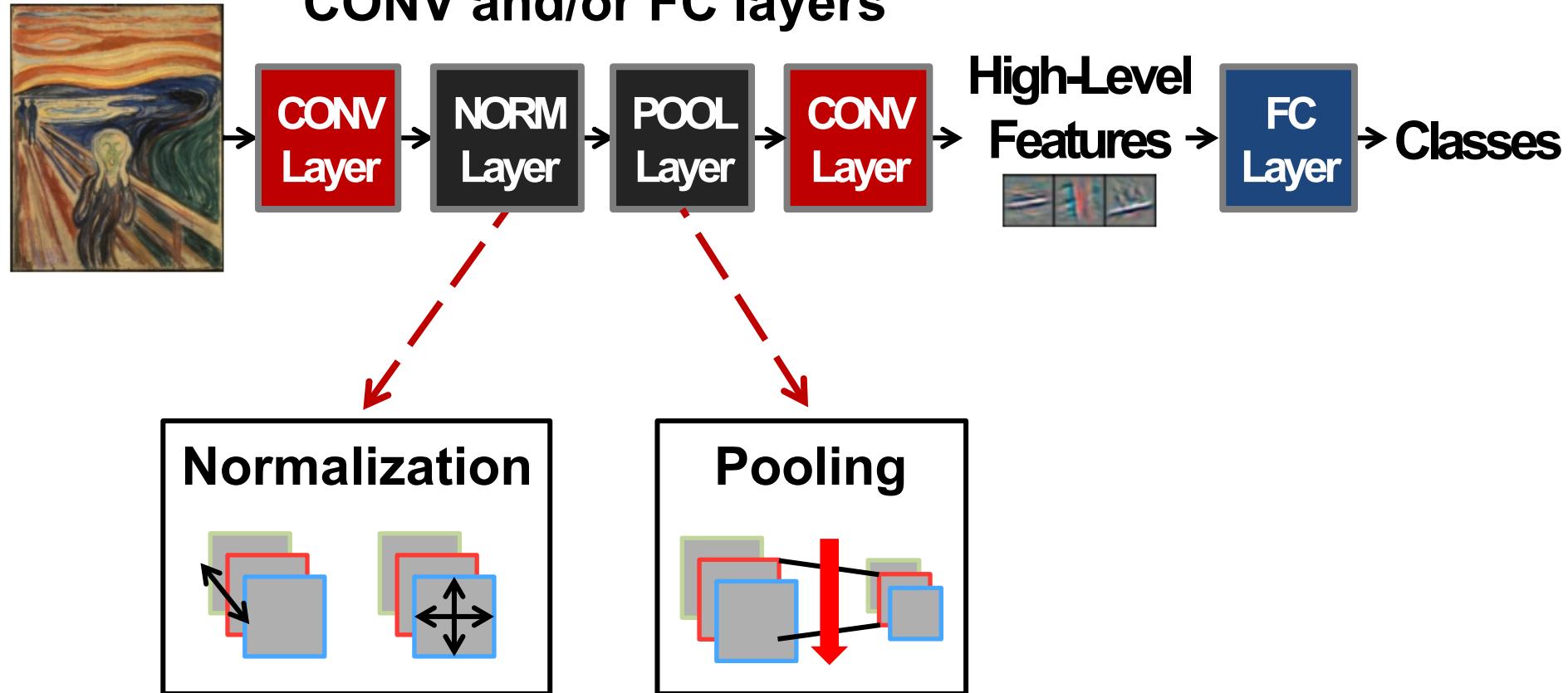


Deep Convolutional Neural Networks

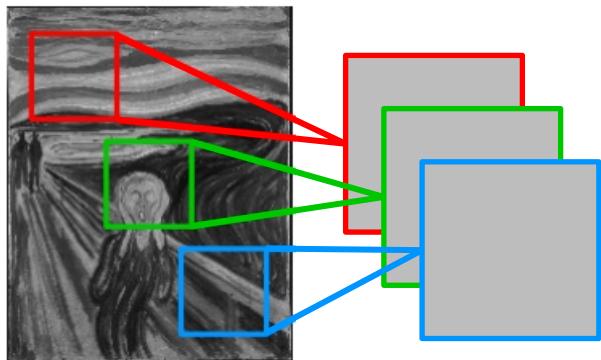
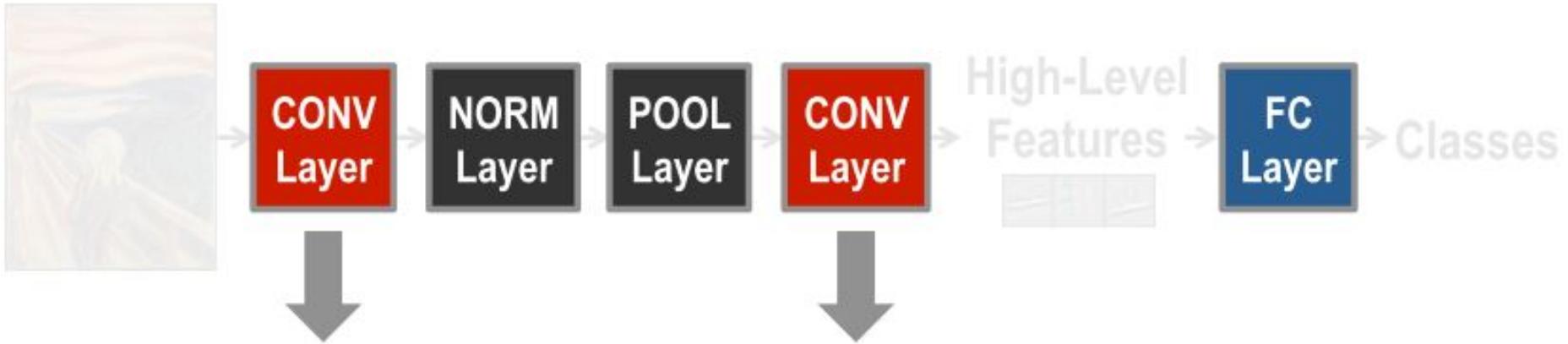


Deep Convolutional Neural Networks

Optional layers in between
CONV and/or FC layers



Deep Convolutional Neural Networks



Convolutions account for more than 90% of overall computation, dominating **runtime** and **energy consumption**

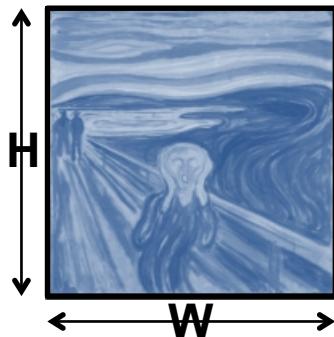
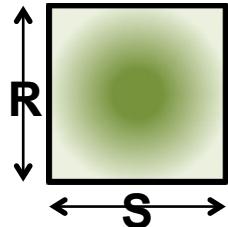
CNN Decoder Ring

- N – Number of **input fmmaps/output fmmaps** (batch size)
- C – Number of 2-D **input fmmaps /filters** (channels)
- H – Height of **input fmap** (activations)
- W – Width of **input fmap** (activations)
- R – Height of 2-D **filter** (weights)
- S – Width of 2-D **filter** (weights)
- M – Number of 2-D **output fmmaps** (channels)
- E – Height of **output fmap** (activations)
- F – Width of **output fmap** (activations)

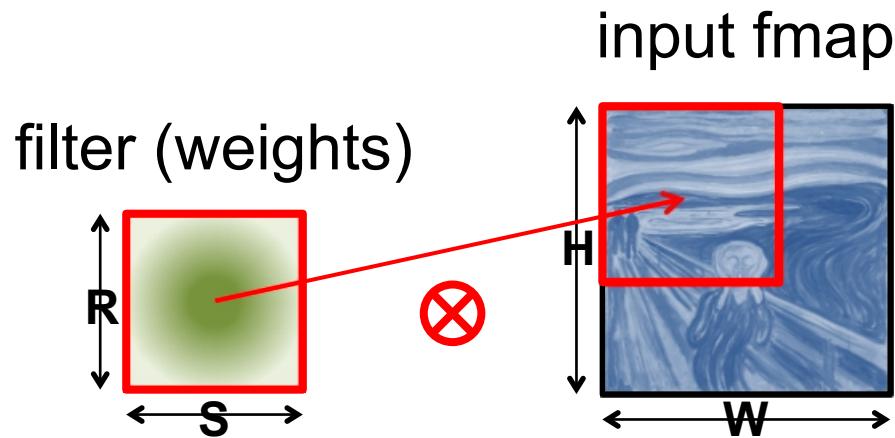
Convolution (CONV) Layer

a plane of input activations
a.k.a. **input feature map (fmap)**

filter (weights)

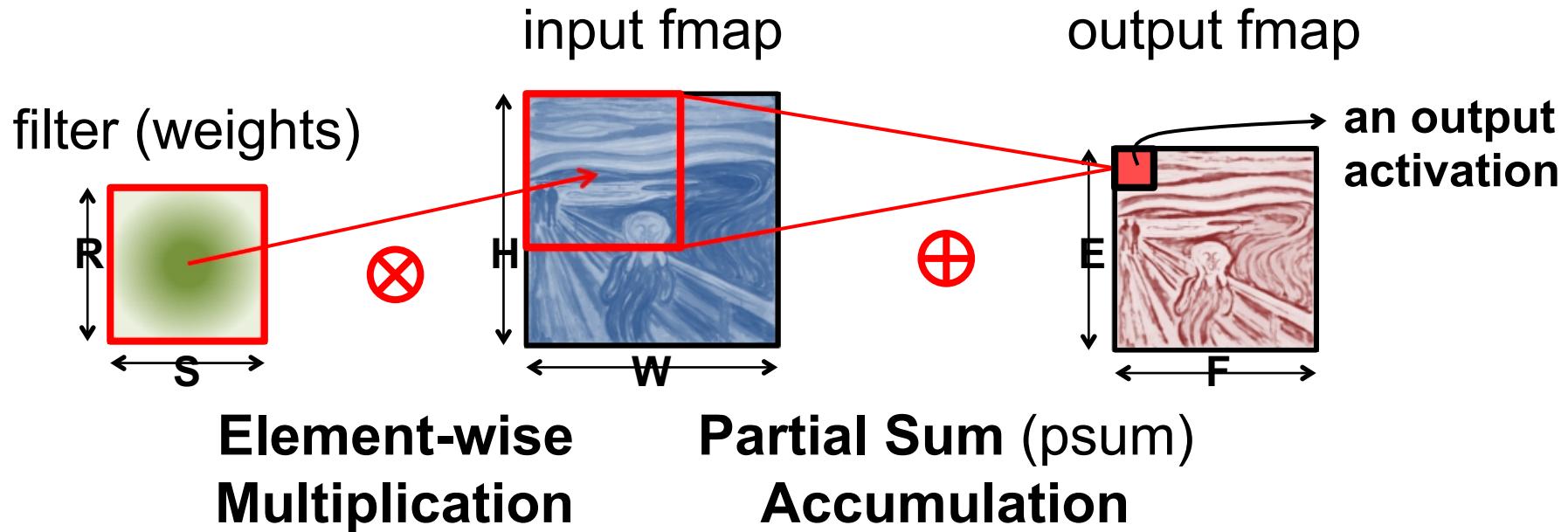


Convolution (CONV) Layer

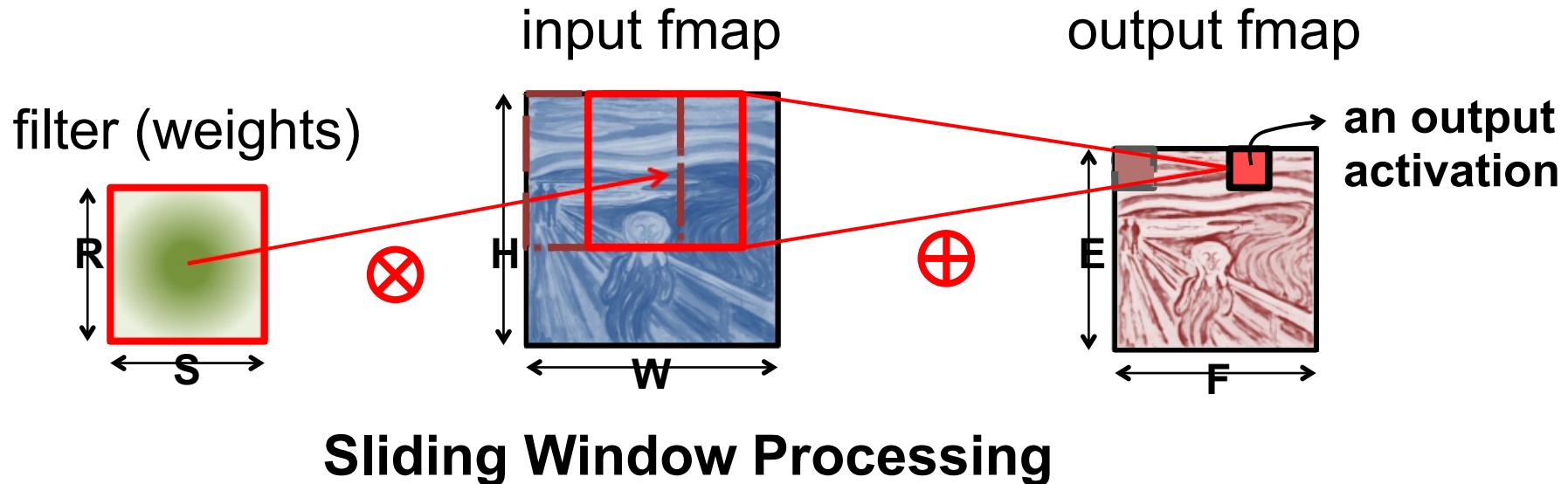


**Element-wise
Multiplication**

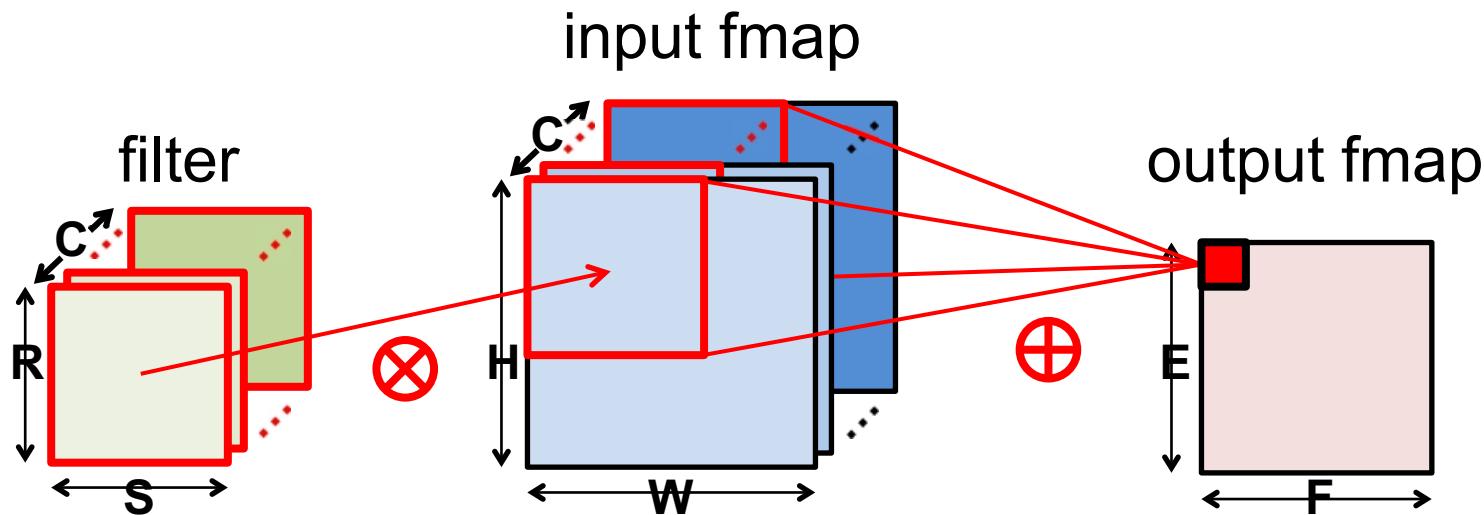
Convolution (CONV) Layer



Convolution (CONV) Layer

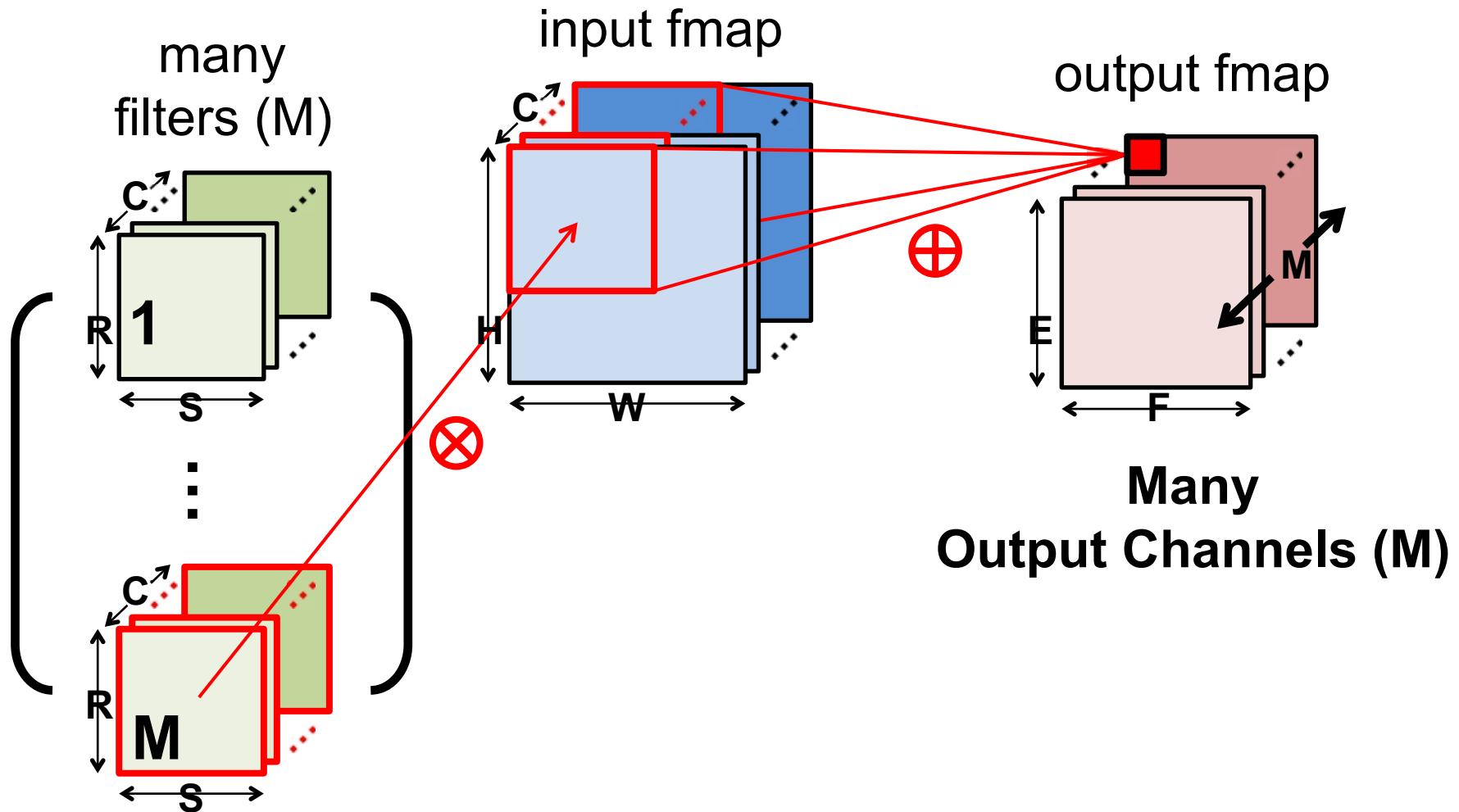


Convolution (CONV) Layer

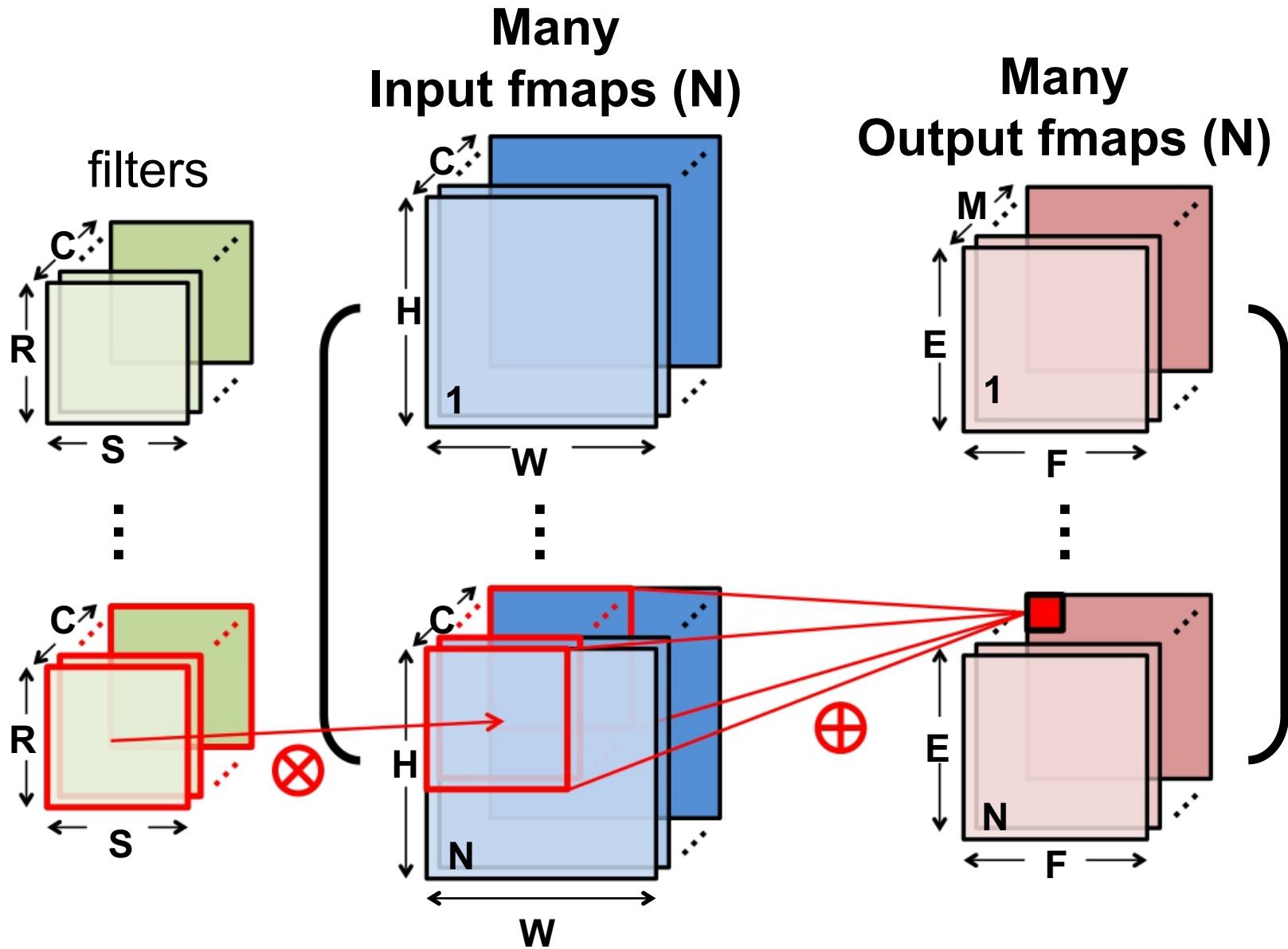


Many Input Channels (C)

Convolution (CONV) Layer



Convolution (CONV) Layer



CNN Decoder Ring

- N – Number of **input fmmaps/output fmmaps** (batch size)
- C – Number of 2-D **input fmmaps /filters** (channels)
- H – Height of **input fmap** (activations)
- W – Width of **input fmap** (activations)
- R – Height of 2-D **filter** (weights)
- S – Width of 2-D **filter** (weights)
- M – Number of 2-D **output fmmaps** (channels)
- E – Height of **output fmap** (activations)
- F – Width of **output fmap** (activations)

CONV Layer Tensor Computation

Output fmmaps (O)

Input fmmaps (I)

Biases (B)

Filter weights (W)

$$\underline{\mathbf{O}[n][m][x][y]} = \text{Activation}(\underline{\mathbf{B}[m]} + \sum_{i=0}^{R-1} \sum_{j=0}^{S-1} \sum_{k=0}^{C-1} \underline{\mathbf{I}[n][k][Ux+i][Uy+j]} \times \underline{\mathbf{W}[m][k][i][j]}),$$

$$0 \leq n < N, 0 \leq m < M, 0 \leq y < E, 0 \leq x < F,$$

$$E = (H - R + U)/U, F = (W - S + U)/U.$$

Shape Parameter	Description
N	fmap batch size
M	# of filters / # of output fmap channels
C	# of input fmap/filter channels
H/W	input fmap height/width
R/S	filter height/width
E/F	output fmap height/width
U	convolution stride

CONV Layer Implementation

Naïve 7-layer for-loop implementation:

```

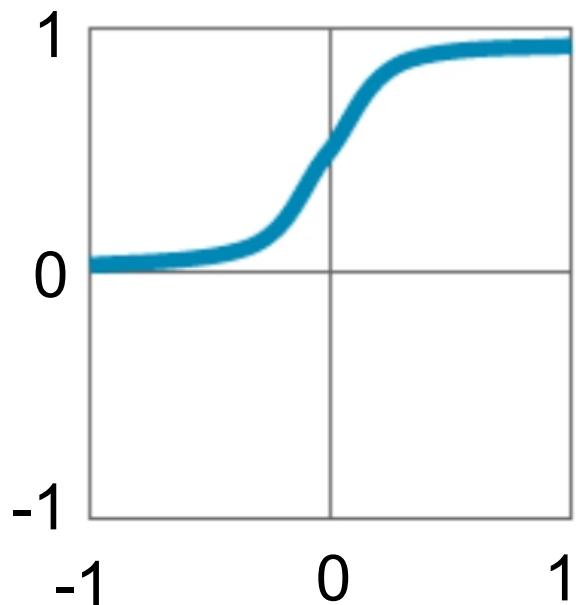
for (n=0; n<N; n++) {
    for (m=0; m<M; m++) {
        for (x=0; x<F; x++) {
            for (y=0; y<E; y++) {
                o[n][m][x][y] = B[m];
                for (i=0; i<R; i++) {
                    for (j=0; j<S; j++) {
                        for (k=0; k<C; k++) {
                            o[n][m][x][y] += I[n][k][Ux+i][Uy+j] * W[m][k][i][j];
                        }
                    }
                }
                o[n][m][x][y] = Activation(o[n][m][x][y]);
            }
        }
    }
}

```

} for each output fmap value

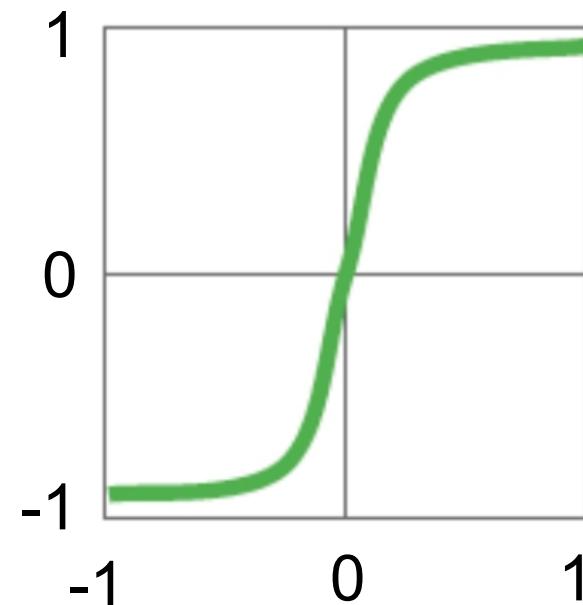
Traditional Activation Functions

Sigmoid



$$y = 1 / (1 + e^{-x})$$

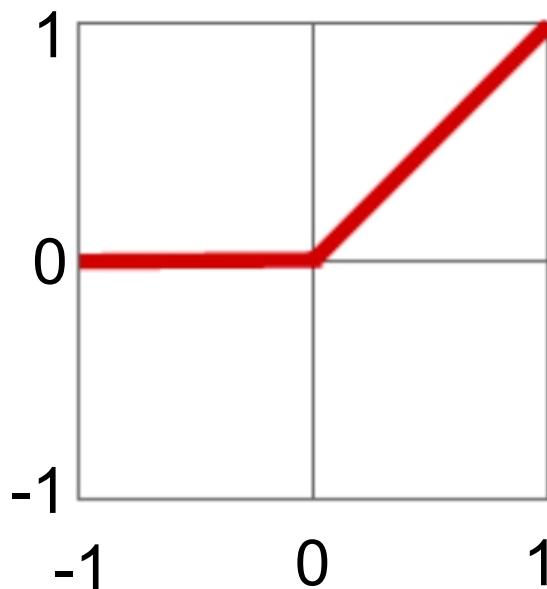
Hyperbolic Tangent



$$y = (e^x - e^{-x}) / (e^x + e^{-x})$$

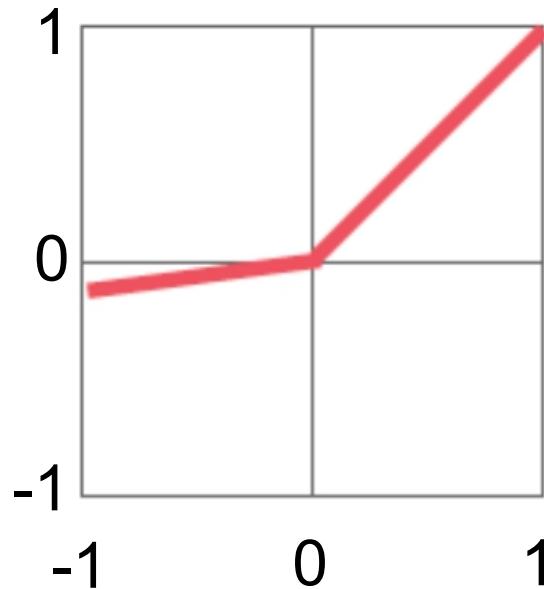
Modern Activation Functions

Rectified Linear Unit
(ReLU)



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

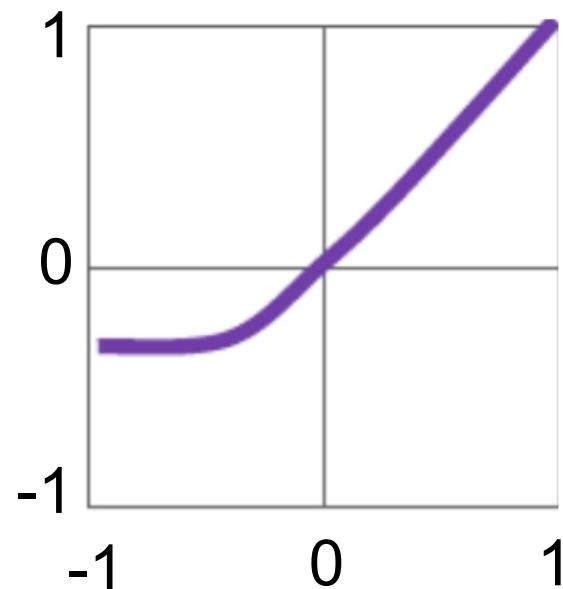
Leaky ReLU



$$y = \max(\alpha x, x)$$

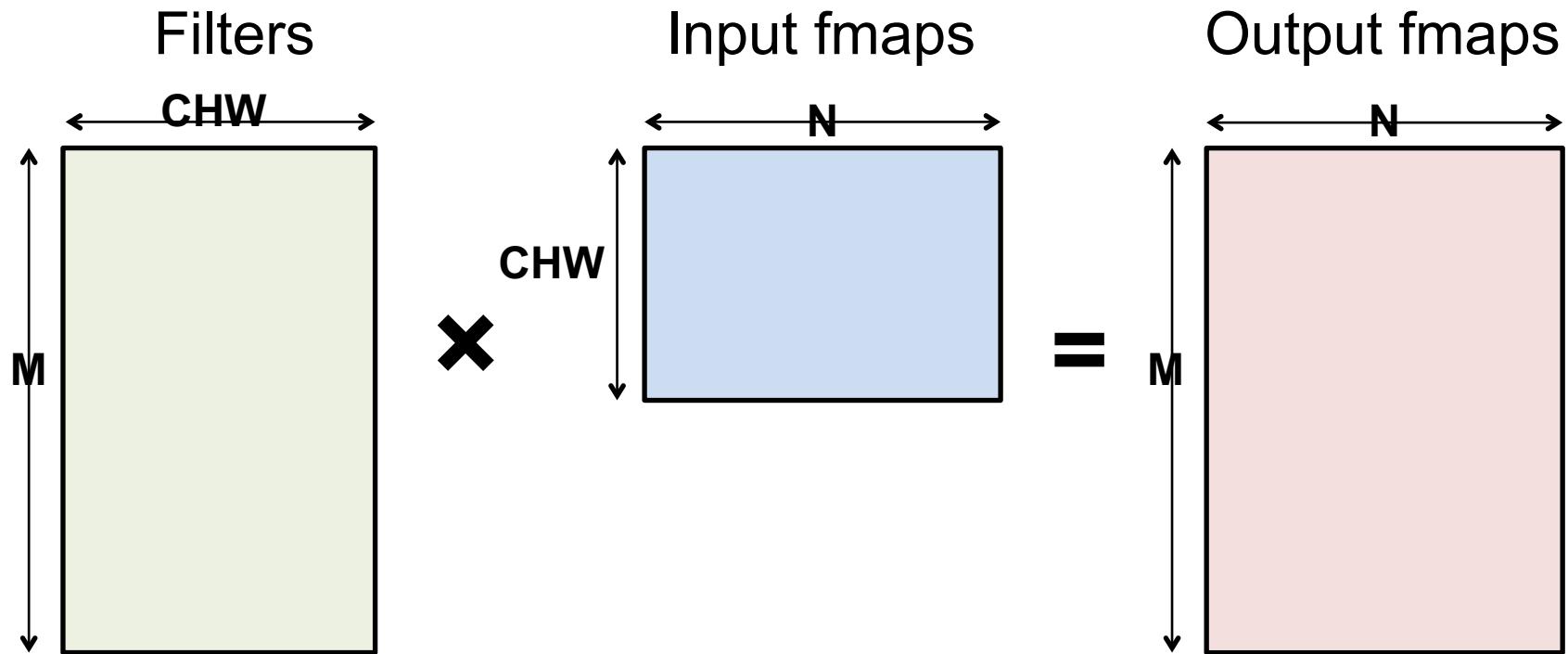
$\alpha = \text{small const. (e.g. 0.1)}$

Exponential LU

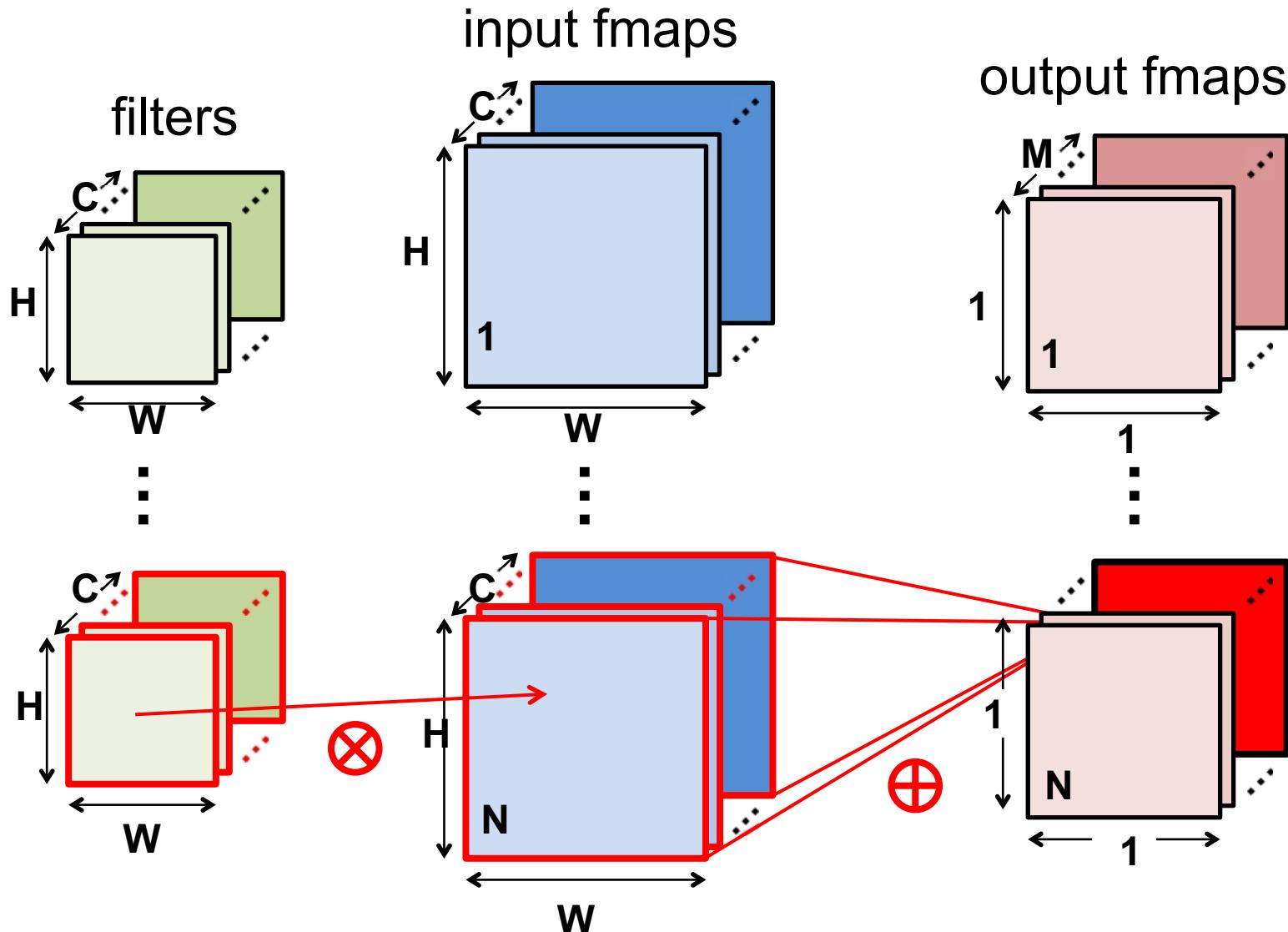


Fully-Connected (FC) Layer

- Height and width of output fmmaps are 1 ($E = F = 1$)
- Filters as large as input fmmaps ($R = H, S = W$)
- Implementation: **Matrix Multiplication**

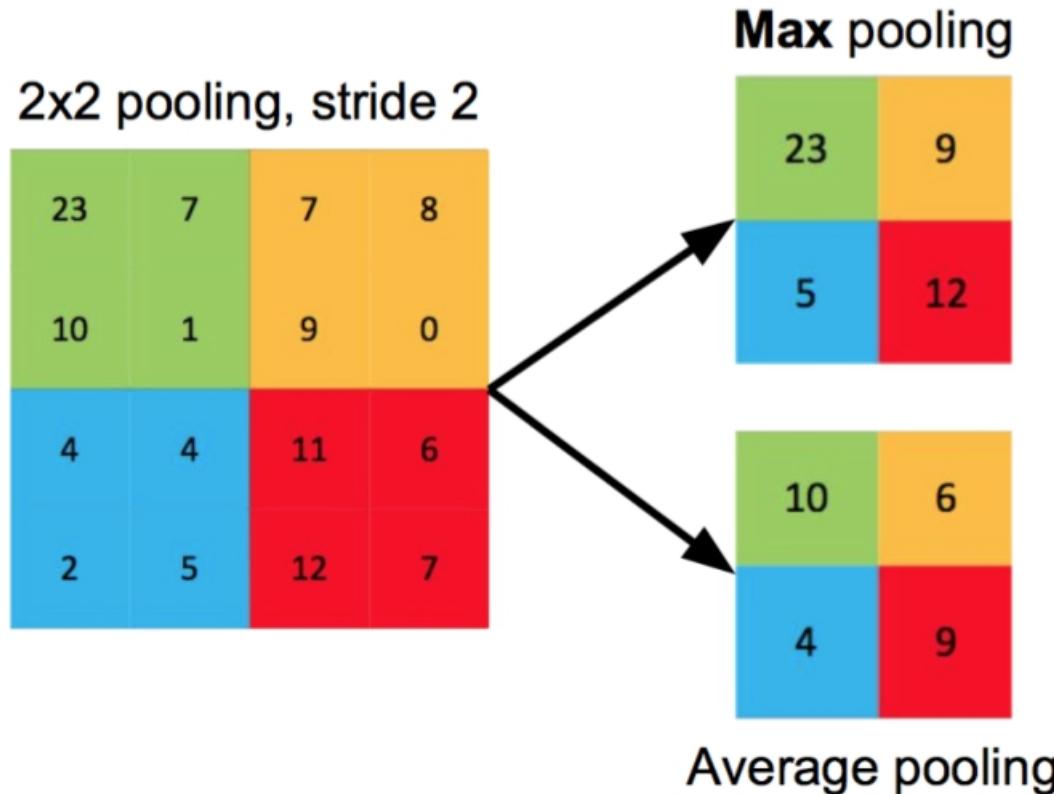


FC Layer – from CONV Layer POV



Pooling (POOL) Layer

- Reduce resolution of each channel independently
- Overlapping or non-overlapping → depending on stride



Increases translation-invariance and noise-resilience

POOL Layer Implementation

Naïve 6-layer for-loop max-pooling implementation:

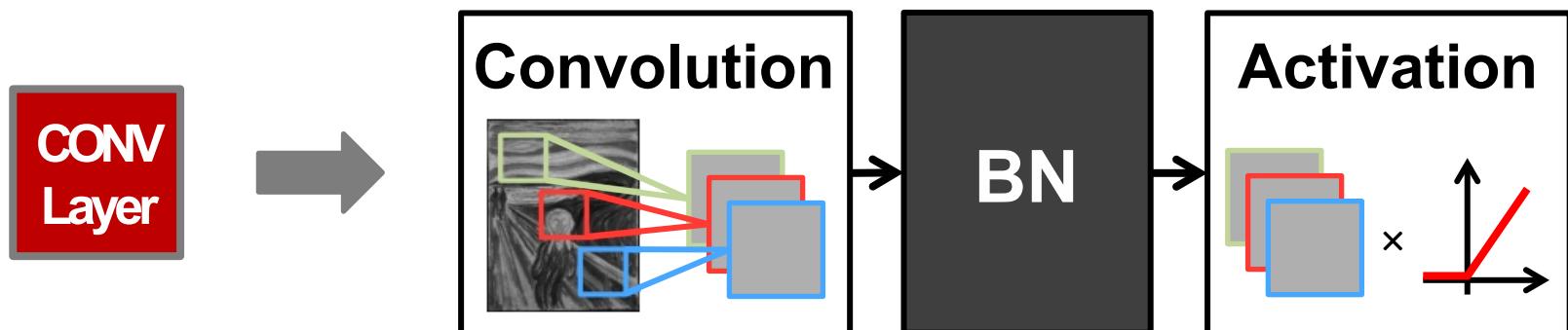
```
for (n=0; n<N; n++) {  
    for (m=0; m<M; m++) {  
        for (x=0; x<F; x++) {  
            for (y=0; y<E; y++) {  
  
                max = -Inf;  
                for (i=0; i<R; i++) {  
                    for (j=0; j<S; j++) {  
                        if (I[n][m][Ux+i][Uy+j] > max)  
                        {  
                            max = I[n][m][Ux+i][Uy+j];  
                        }  
                    }  
                }  
                o[n][m][x][y] = max;  
            }  
        }  
    }  
}
```

} for each pooled value

} find the max
with in a window

Normalization (NORM) Layer

- **Batch Normalization (BN)**
 - Normalize activations towards mean=0 and std. dev.=1 based on the statistics of the training dataset – put **in between** CONV/FC and Activation function



Believed to be key to getting high accuracy and faster training on very deep neural networks.

BN Layer Implementation

- The normalized value is further scaled and shifted, the parameters of which are learned from training

$$y = \frac{x - \mu}{\sqrt{\sigma^2 + \epsilon}} \gamma + \beta$$

Annotations for the BN formula:

- data mean**: μ (red arrow)
- learned scale factor**: γ (blue arrow)
- learned shift factor**: β (blue arrow)
- data std. dev.**: σ (red arrow)
- small const. to avoid numerical problems**: ϵ (grey arrow)

Relevant Components for Tutorial

- Typical operations that we will discuss:
 - Convolution (CONV)
 - Fully-Connected (FC)
 - Max Pooling
 - ReLU

Backup Slides