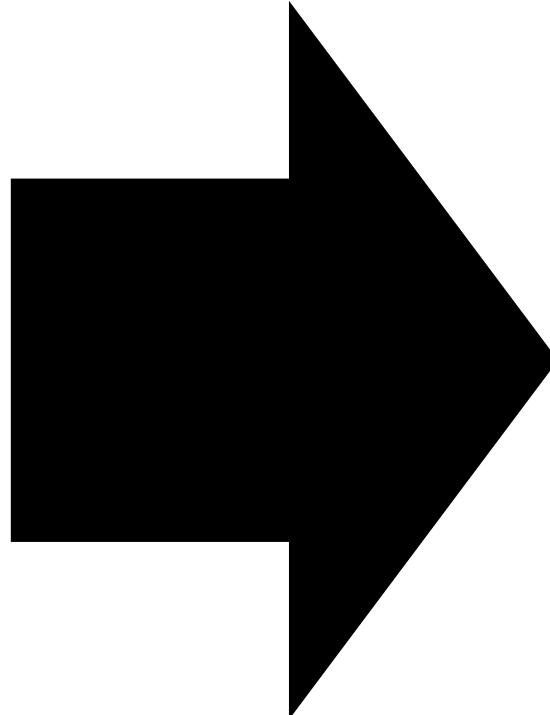


Neuroengineering (I)

4. Deep Learning with ANN

- **Scuola di Ingegneria Industriale e dell'Informazione**
 - Politecnico di Milano
- Prof. Pietro Cerveri

Cognition is based on abstract models

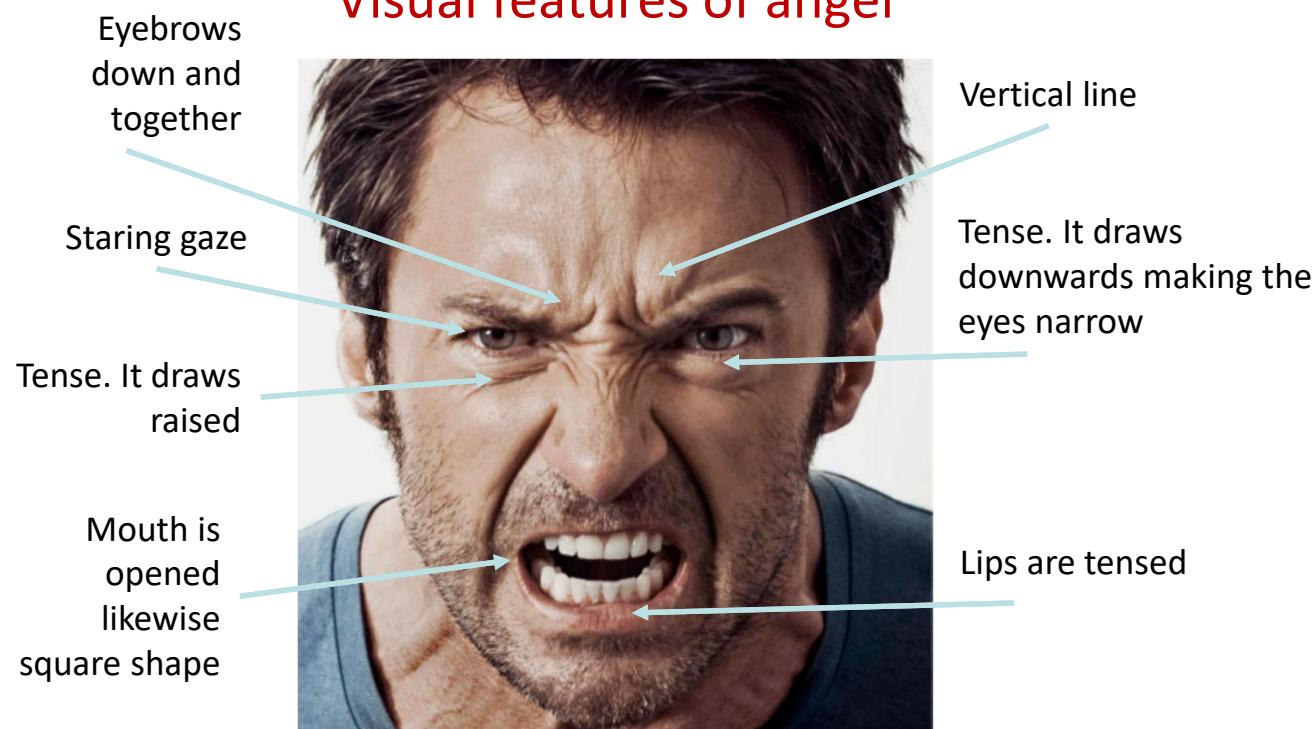


From sensation to semantics
Synthesizing features



Visual models

Visual features of anger



Dense scalar signal

Deconstruction

Features

Deconstruction – Learning in depth

Dense scalar signal



Features



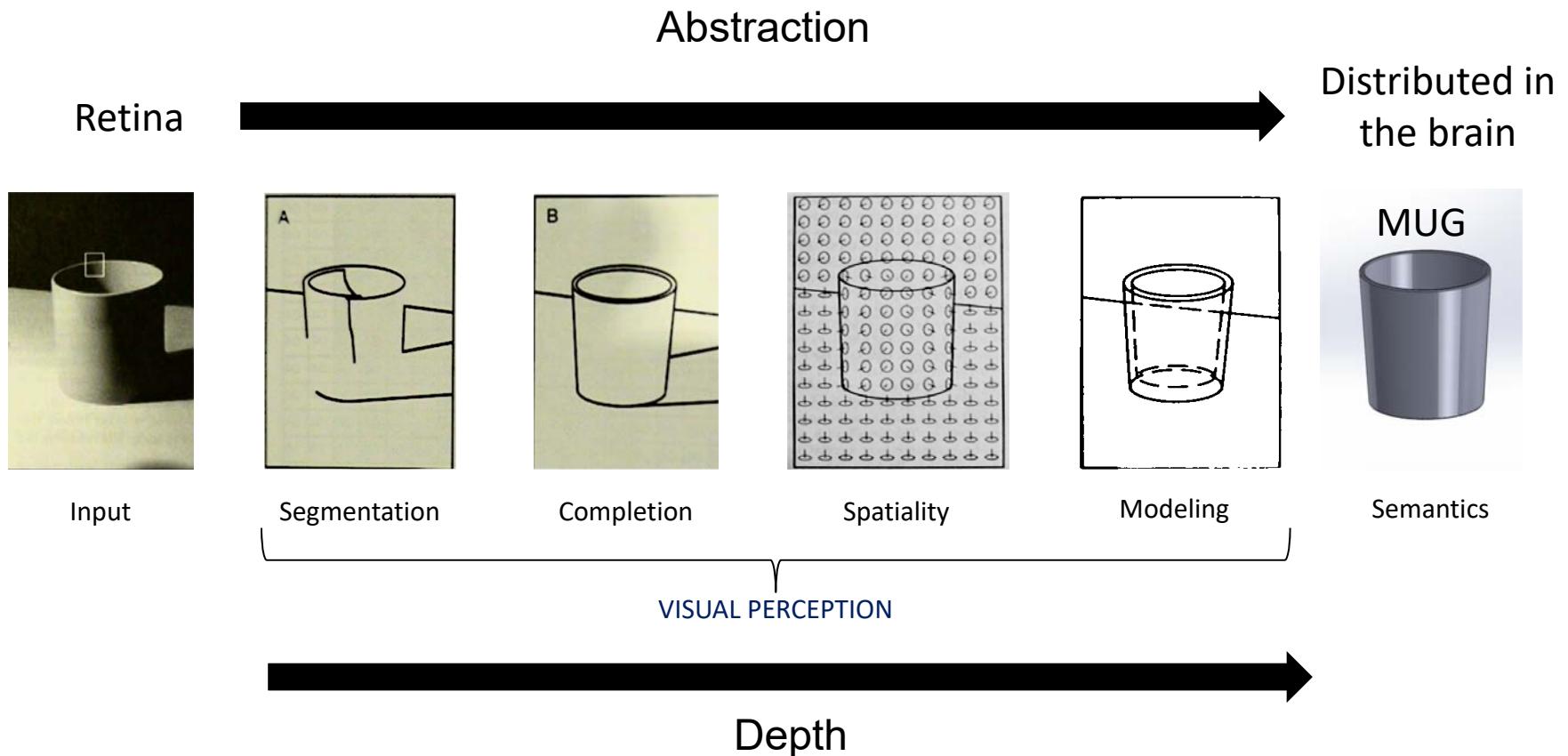
Image processing



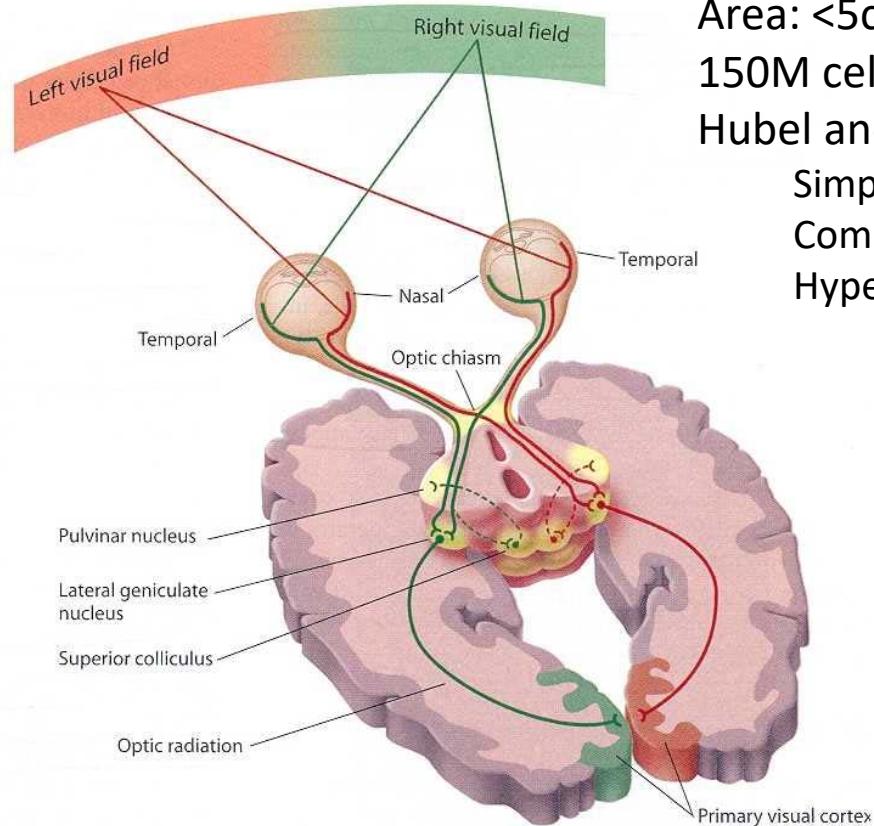
Deep learning

- From sensation to semantics
 - Increasing abstraction
 - Heterogeneous components (feature maps)
 - Physical (pixel intensity, signal frequencies, ...)
 - Scalar maps (lines, contours, curves, colors, ..)
 - Vector descriptors (geometries)
 - Perceptual (completion, depth, motion,...)
 - Emotional (aesthetics, affection, personal experiences,..)
 - Semantic (language/concepts)
-
1. Sensorial experience is deconstructed into the brain
 2. Representation made of a cascade of layers
 3. Mechanisms of knowledge reconstruction

Deconstruction (progressive encoding)



Primary visual cortex



Thickness: 2mm

Area: $<5\text{cm}^2$

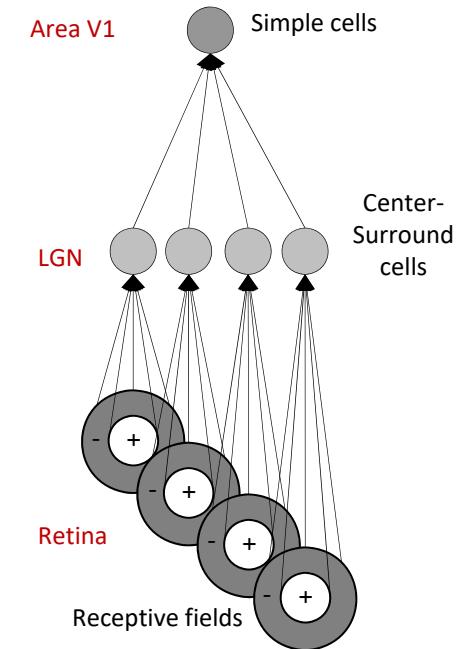
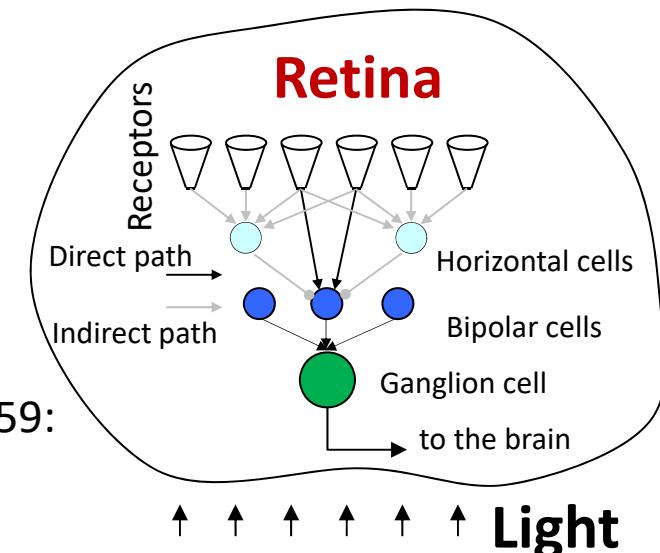
150M cells

Hubel and Wiesel, 1959:

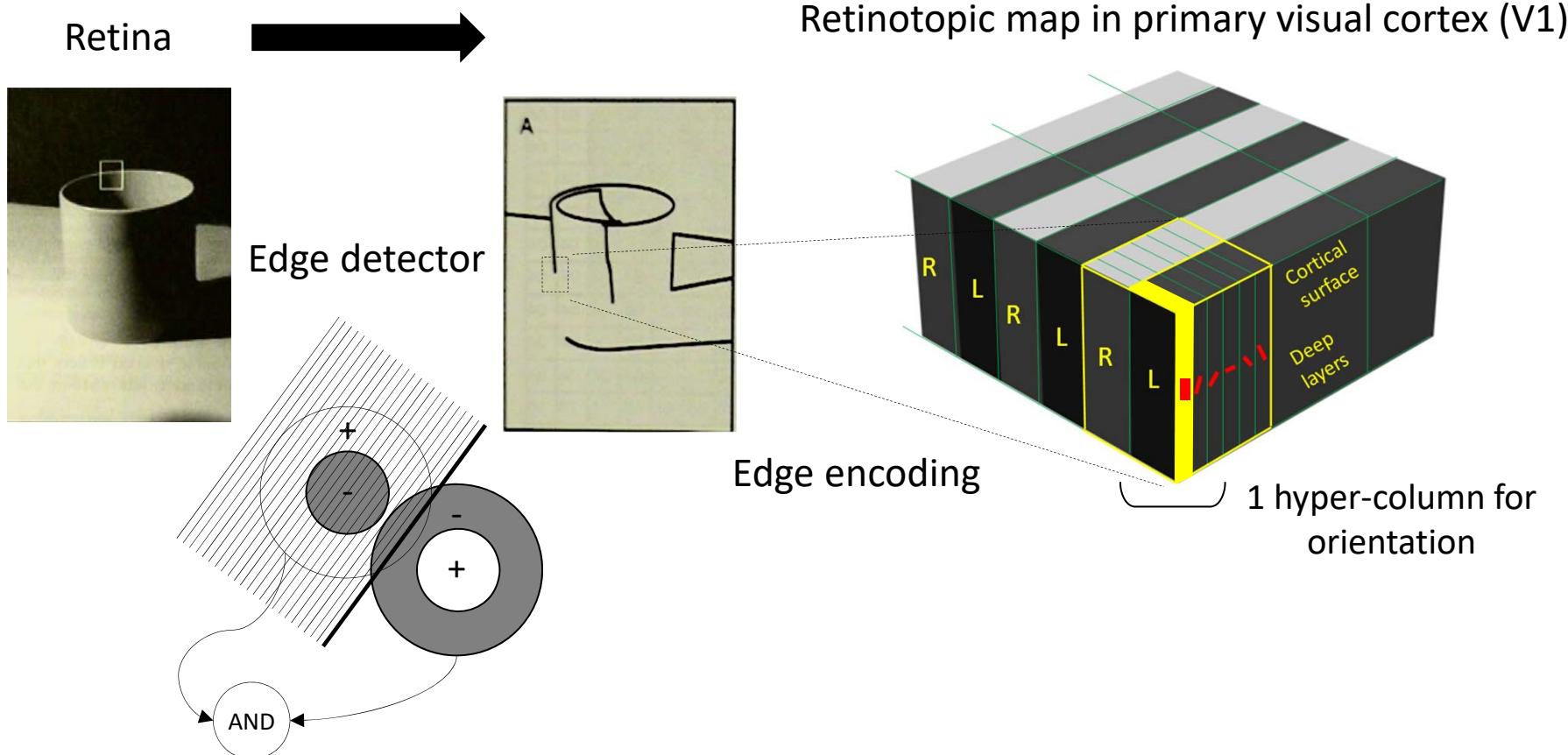
Simple cells

Complex cells

Hyper-complex cells



Visual encoding

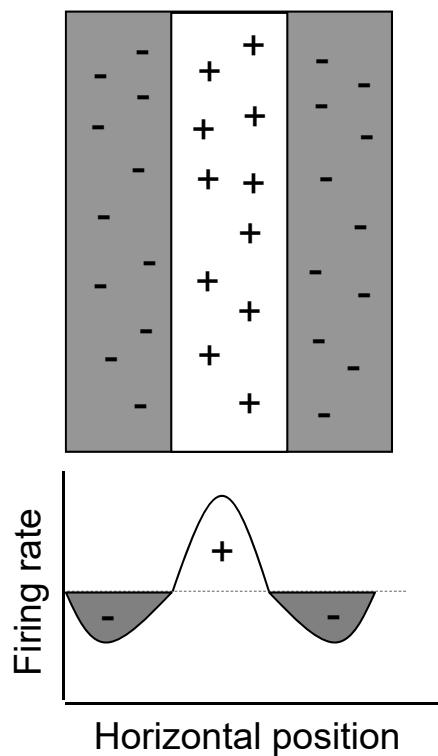


Neural implementation in LGN
(signals from ganglion cells)

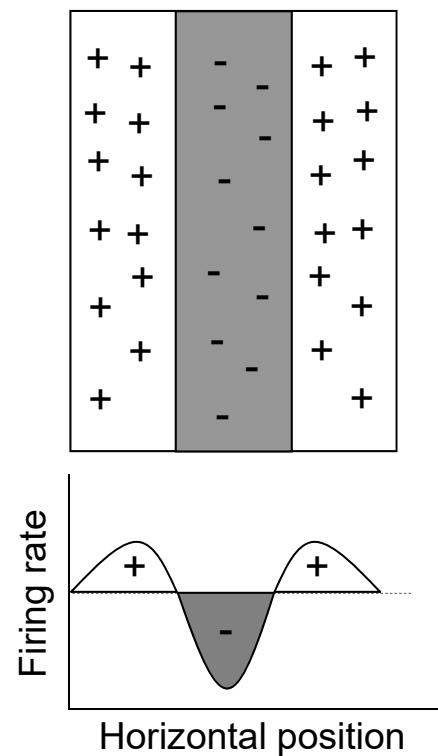
- On-center \cap Off-center

Receptive field of simple cells

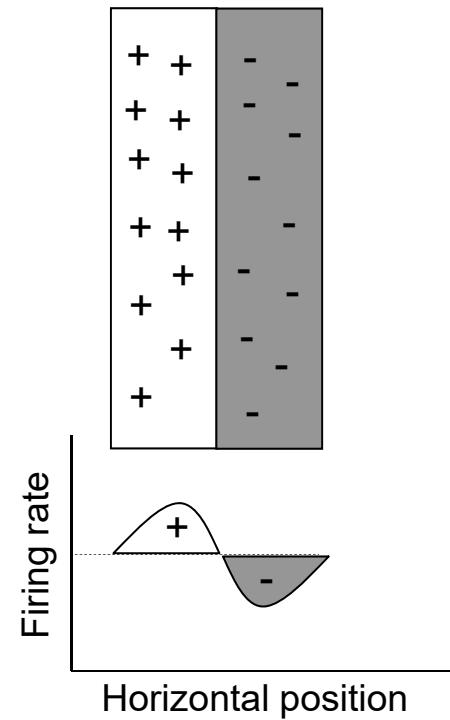
Detector for bright vertical line



Detector for dark vertical line



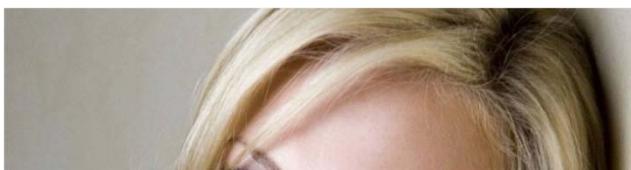
Detector for dark-bright border



Convolutional neural networks (CNN)

- Artificial implementation of physiological deep learning
- Based on mathematical convolution
- CNN use three basic ideas:
 - *local receptive fields*
 - *shared weights*
 - *rectified linear units and pooling*

Input image



Output of the convolution operation



Recall that

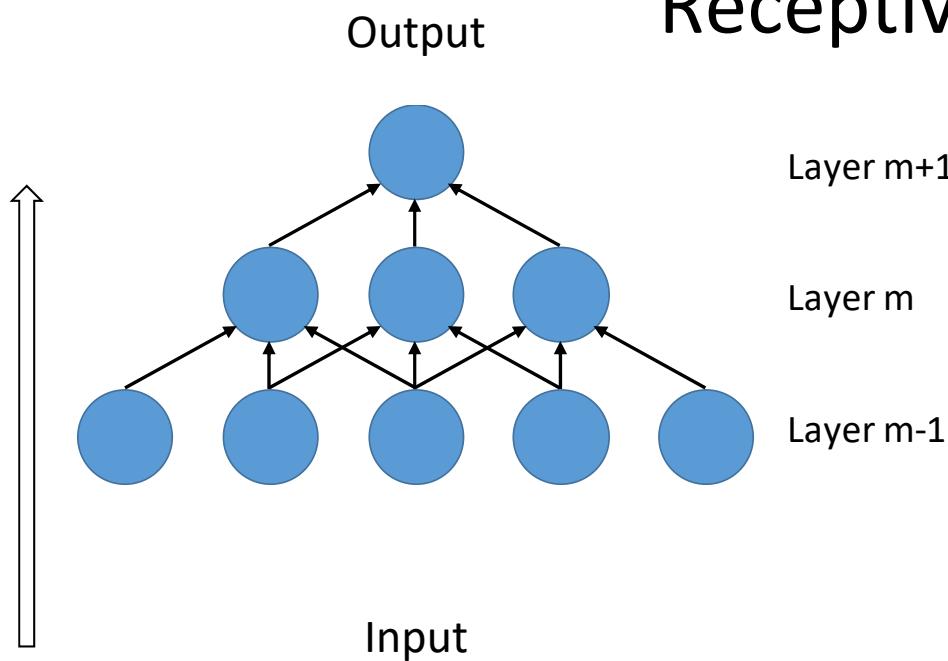
$$1D: o[n] = f[n]*g[n] = \sum_{u=-\infty}^{\infty} f[u]g[n-u] = \sum_{u=-\infty}^{\infty} f[n-u]g[u]$$

$$2D: o[m, n] = f[m, n]*g[m, n] = \sum_{u=-\infty}^{\infty} \sum_{v=-\infty}^{\infty} f[u, v]g[m-u, n-v]$$



Filtering and down-sampling

Receptive field



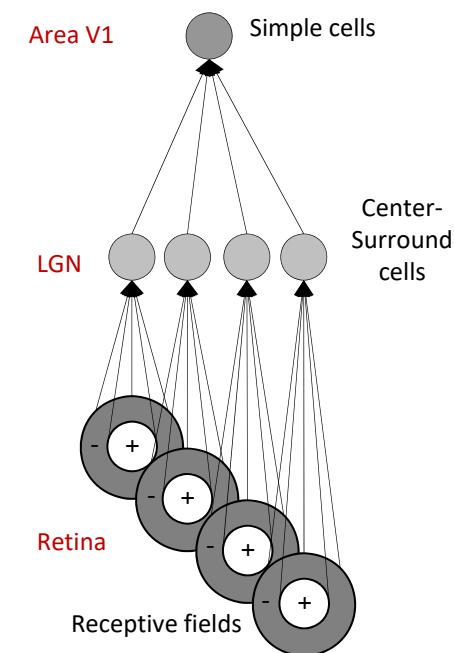
Layer m+1

Layer m

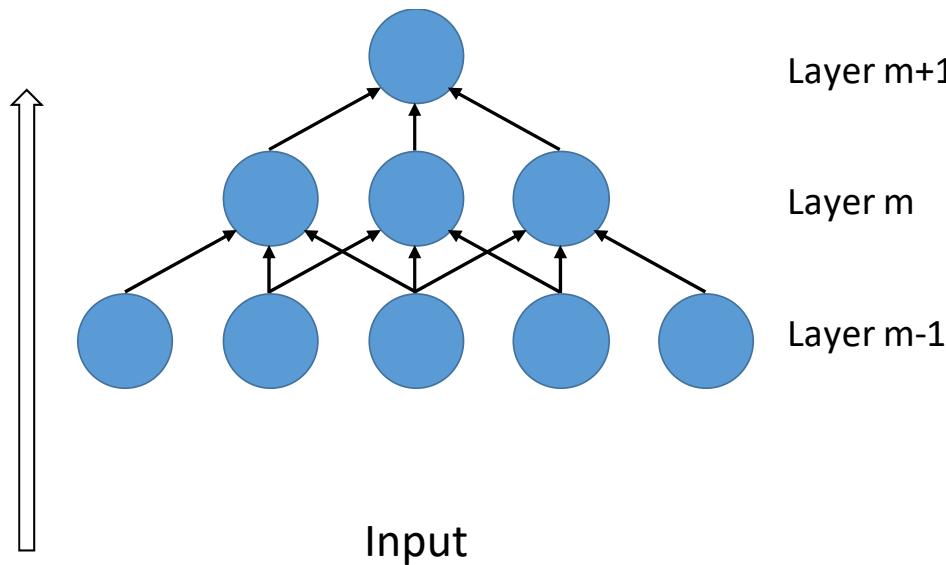
Layer m-1

n-unit local receptive field

Each unit is unresponsive to variations outside of its receptive field



Receptive field



A convolution operation forms a many-to-one relationship

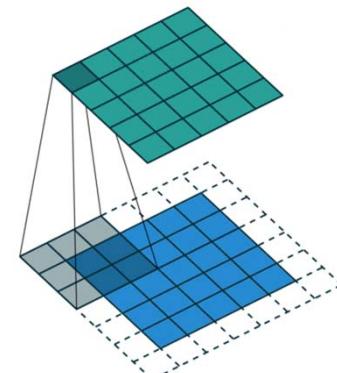
Implemented by means of a spatial filter

Size and stride

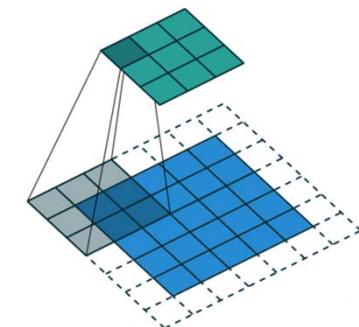
n-unit local receptive field

Each unit is unresponsive to variations outside of its receptive field

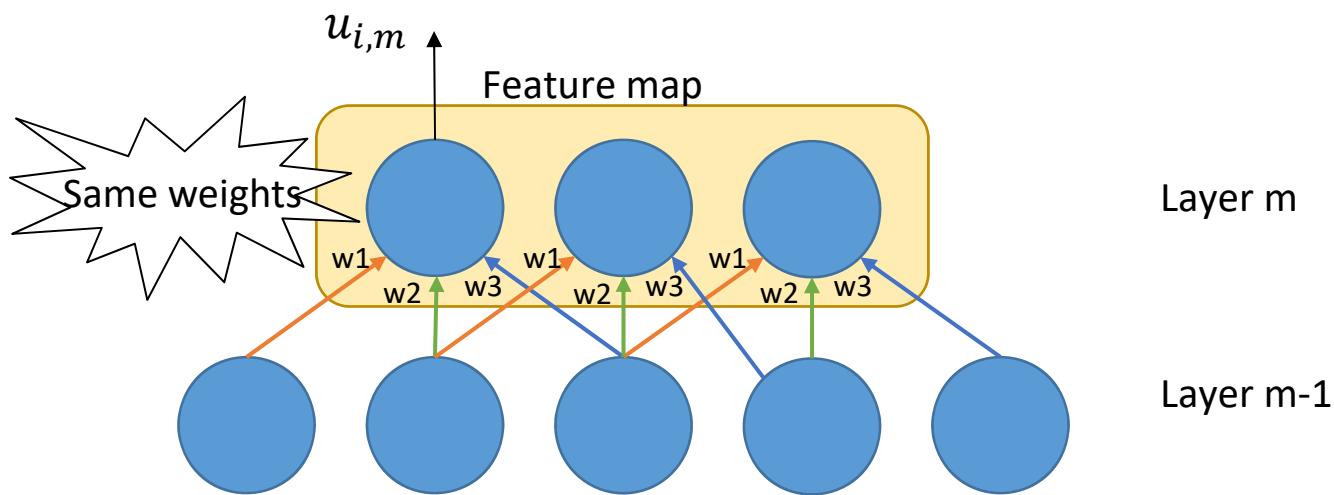
Filtering



Filtering and down-sampling



Shared weights and feature map

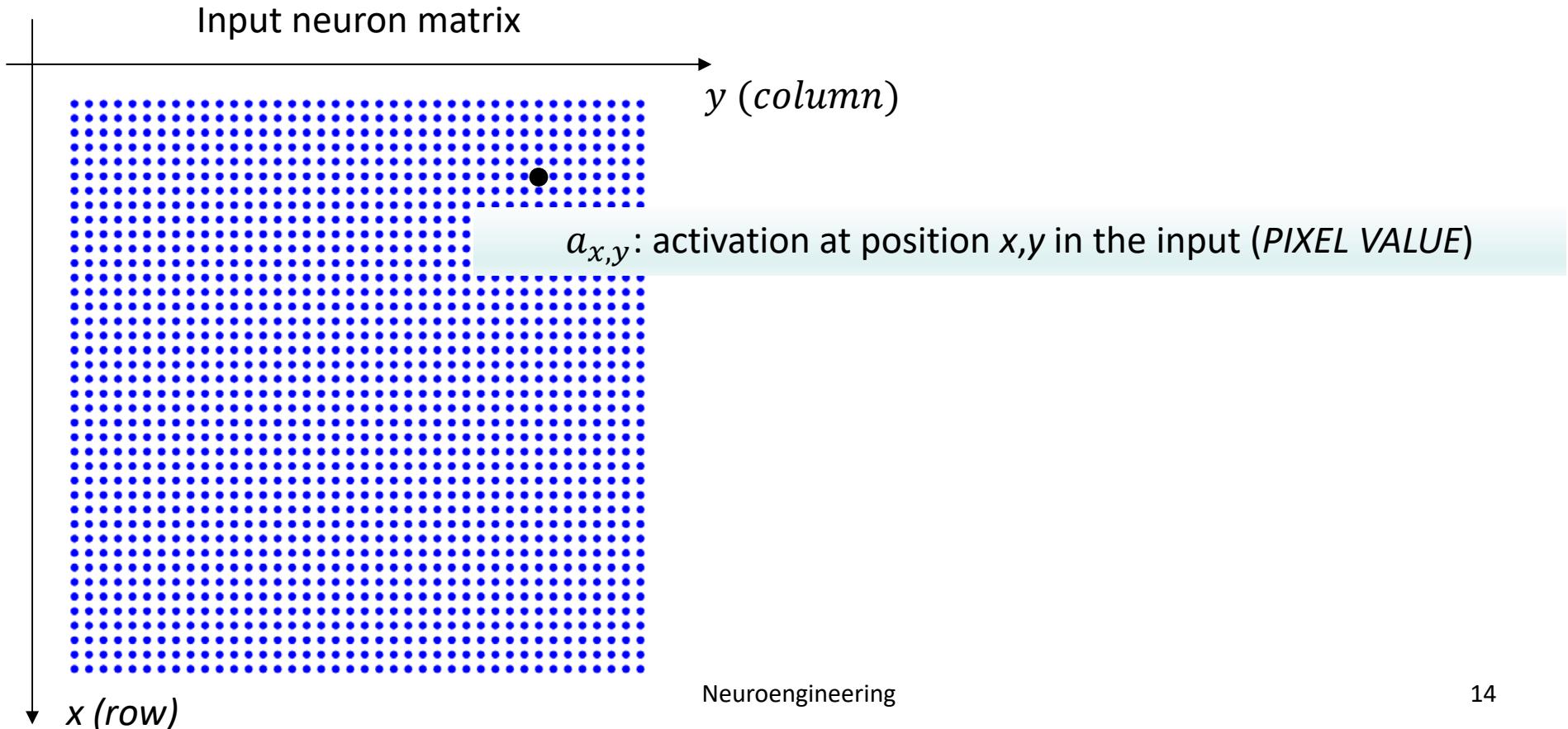


Each neuron shares the same weights in the layer (feature map)

$$u_{i,m} = \tanh \left(b + \sum_{n=1}^N w_n a_{n,m-1} \right)$$

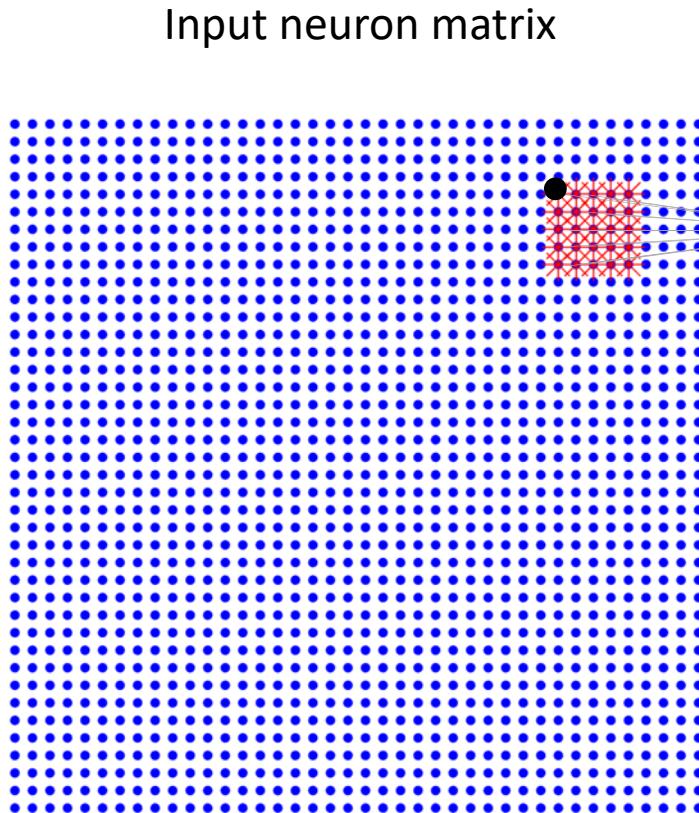
Convolutional layer: image processing

- $K \times K$ pixel input image



Convolutional layer: image processing

- $K \times K$ pixel input image and $n \times n$ LCF (local receptive field)

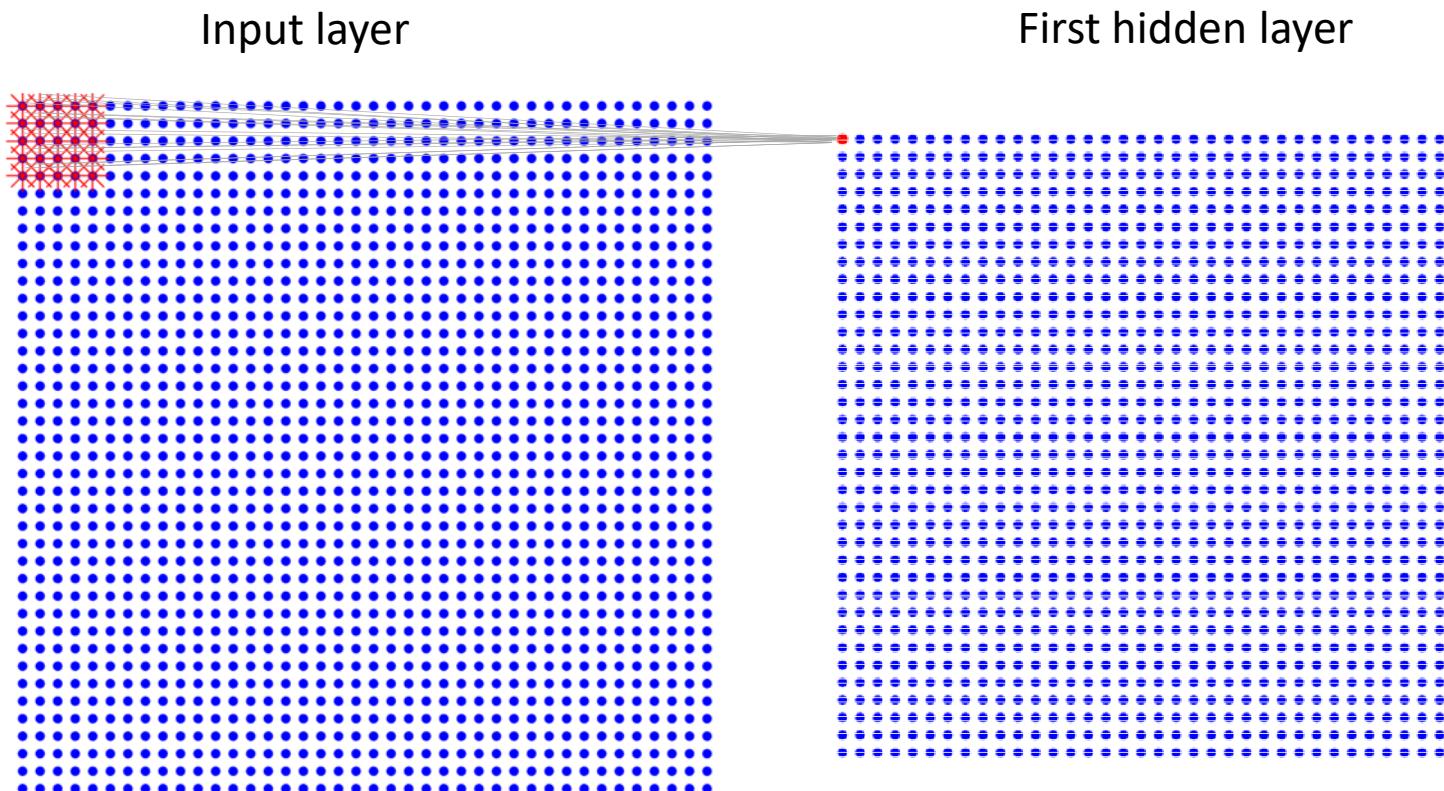


$$u_{j,k} = \tanh \left(b + \sum_{l=1}^L \sum_{m=1}^M w_{l,m} a_{j-1+l, k-1+m} \right)$$

Neuron j, k in the first hidden layer (matrix representation)

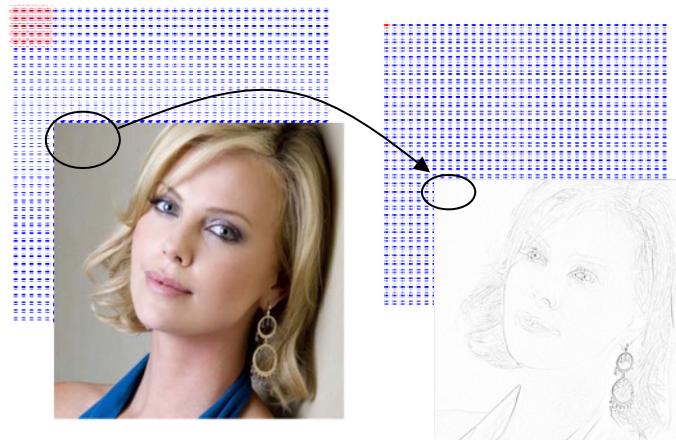
Scanning the image using the LCF

- $K \times K$ pixel input image and $n \times n$ LCF



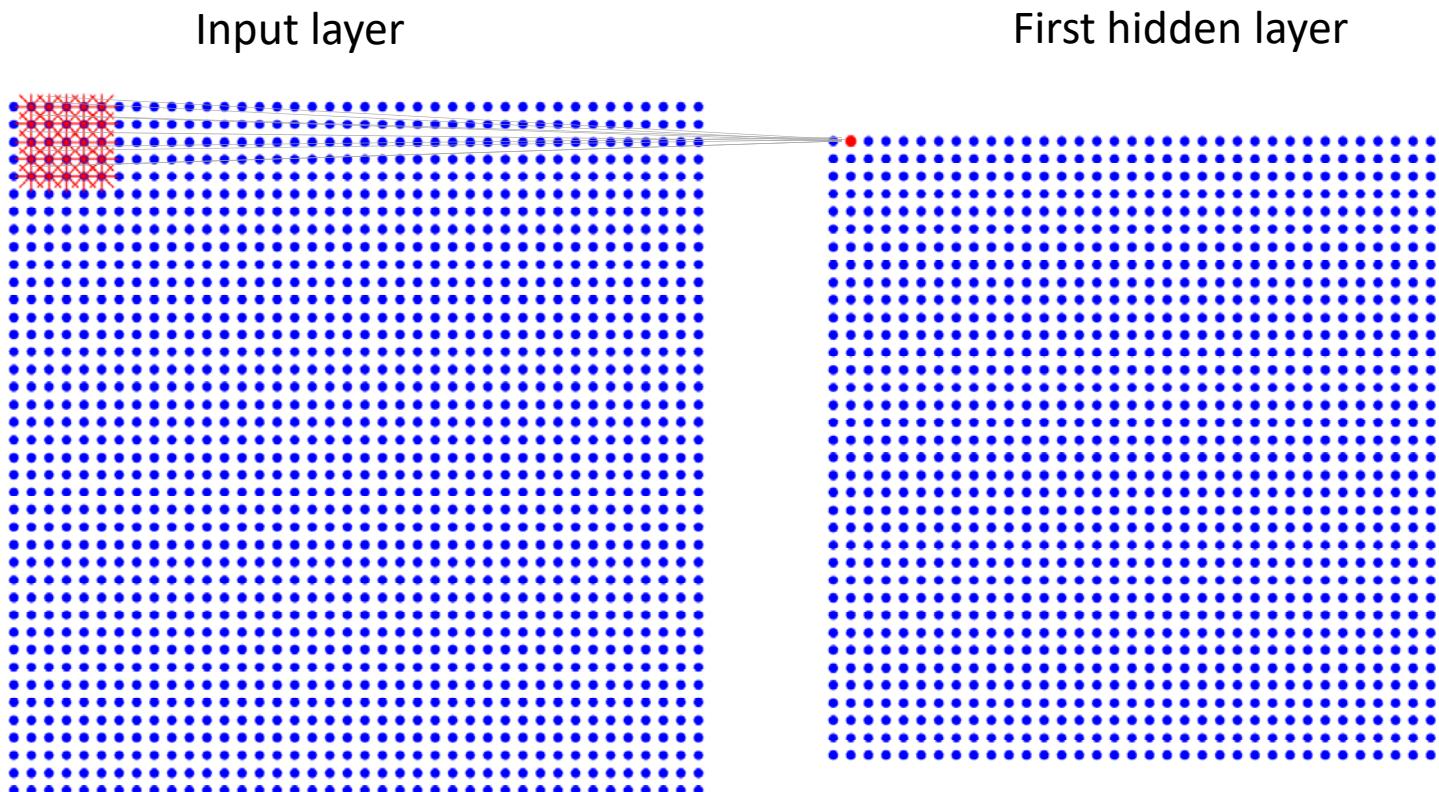
Scanning the image using the LCF

- One important point of such convolution operation is that the positional connectivity exists between the input values and the output values
- For example, the top left values in the input matrix affect the top left values of the output matrix



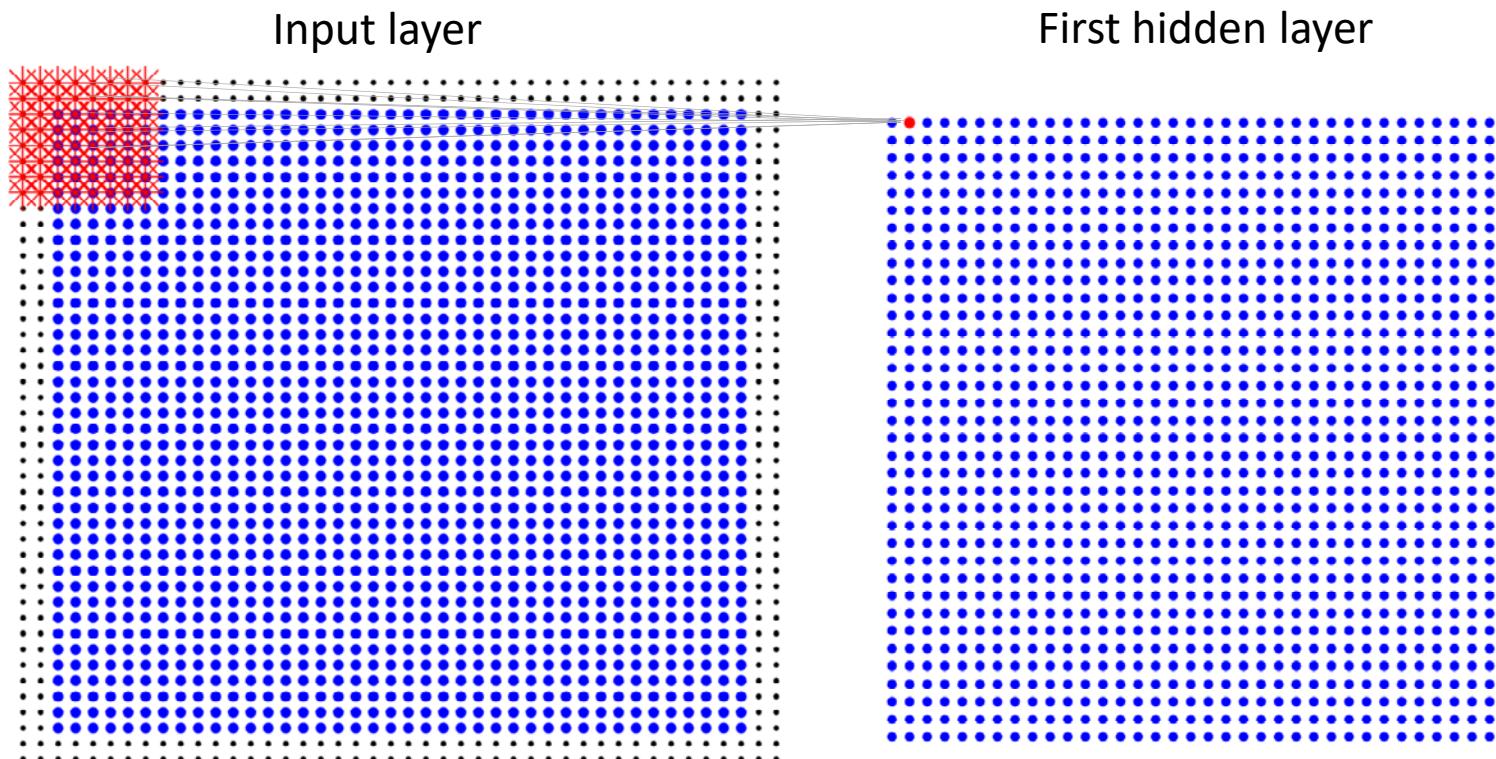
Scanning the image using the LCF

- $K \times K$ pixel input image with $n \times n$ LCF
- Stride length: 1 (pixel by pixel)



Scanning the image using the LCF

- $K \times K$ pixel input image with $n \times n$ LCF
- Zero-padding



The output size (s_x, s_y) of each convolutional layer can be computed as:

$$s_x = \frac{(I_x + 2z_x - f_x)}{L_x} + 1$$
$$s_y = \frac{(I_y + 2z_y - f_y)}{L_y} + 1$$

I_x and I_y : input image sizes in x and y directions

z_x and z_y : zero-padding sizes

f_x and f_y : filtering sizes

L_x and L_y : stride lengths.

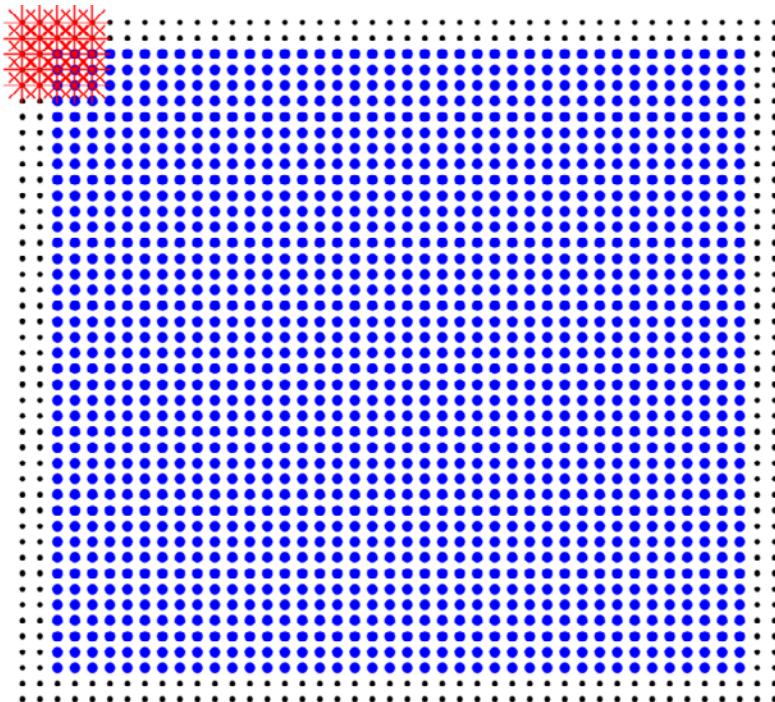
Let us assume a 240×240 pixel images, a convolutional filter of 7×7 pixels with a stride length of 3 and a 2-pixel zero padding both in x and y directions.

The output will be then 80×80 pixels.

Scanning the image using the LCF

- 40×40 pixel original image with 5×5 LCF
- 2 pixel zero-padding
- Stride length: 1 (pixel by pixel)

$$s_x = \frac{(I_x + 2z_x - f_x)}{L_x} + 1 = \frac{(40 + 4 - 5)}{1} + 1 = 40$$

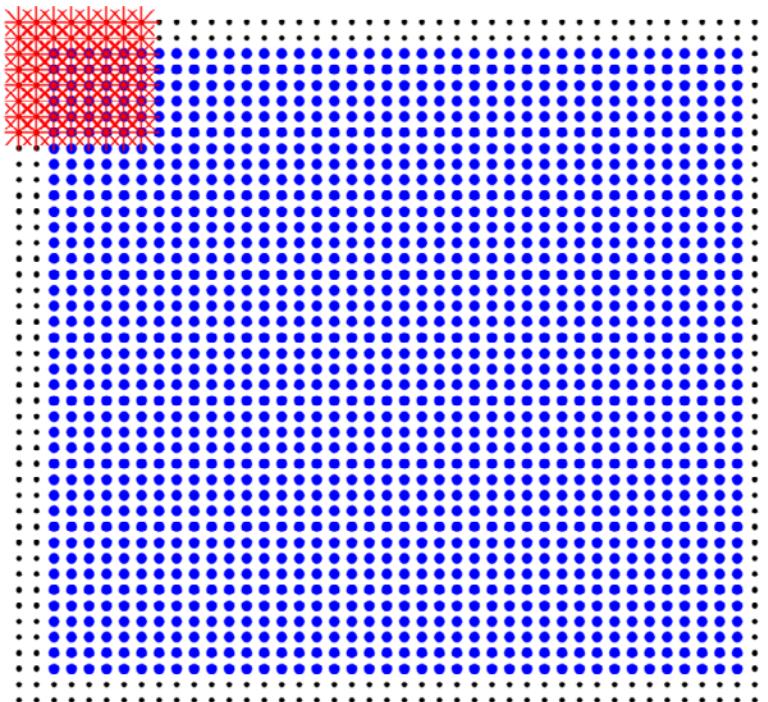


First hidden layer (40×40 neurons)
NO size compression

Scanning the image using the LCF

- 40×40 pixel original image with 8×8 LCF
- 2 pixel zero-padding
- Stride length: 2

$$s_x = \frac{(I_x + 2z_x - f_x)}{L_x} + 1 = \frac{(40 + 4 - 8)}{2} + 1 = 19$$



First hidden layer (19×19 neurons)

To understand and remember: each neuron in the convolutional layer will be forced to have the same weights and the same threshold

Shared weights and biases

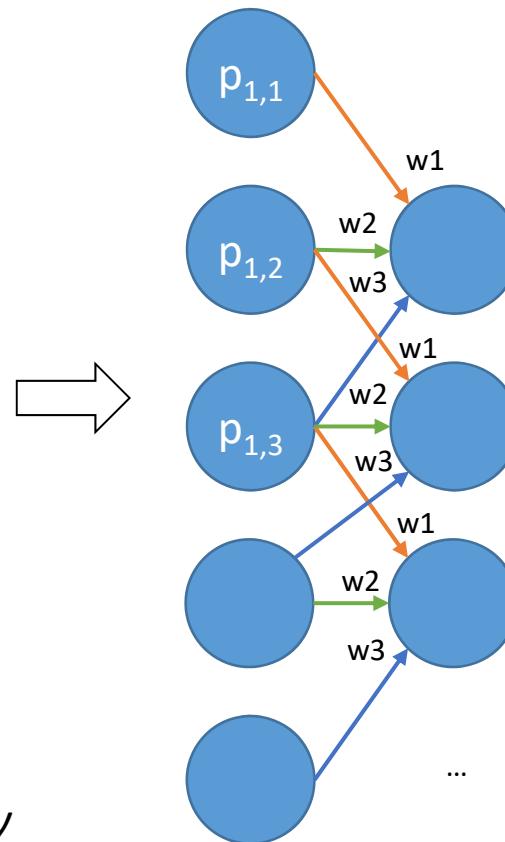
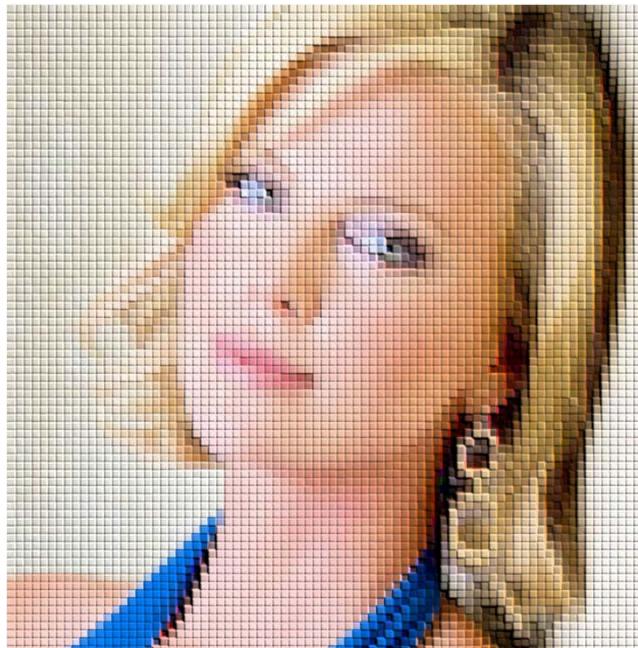
- 40×40 pixel image and 5×5 LCF
- **Feature map:** same weights and biases for all (36×36) the hidden neurons
- For neuron (j, k)

$$u_{j,k} = \tanh \left(b + \sum_{l=1}^5 \sum_{m=1}^5 w_{l,m} a_{j-1+l, k-1+m} \right)$$

$a_{x,y}$ Input activation at position x,y

- 5×5 LCF requires only 26 parameters (25 weights + 1 bias/threshold)

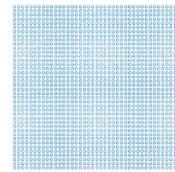
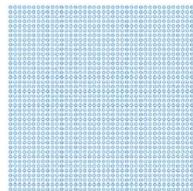
Feeding an image into a network



$p_{x,y}$ = pixel value at row x and column y

Image flattening (transform into a vector) with pixel connectivity constraint

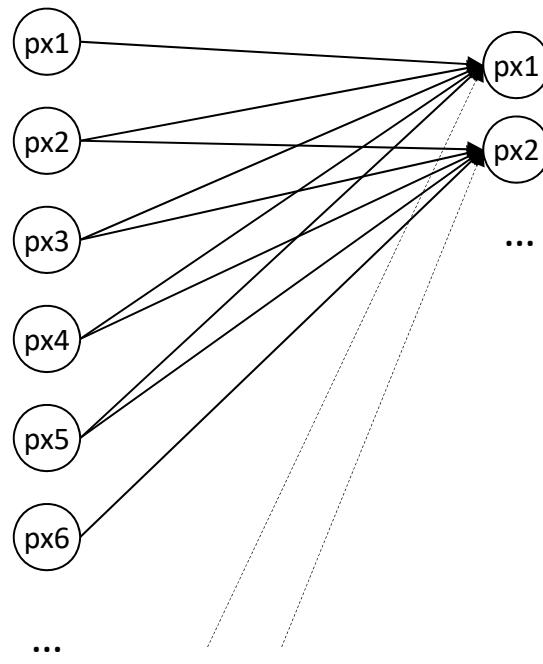
$40 \times 40 = 1600$ inputs $36 \times 36 = 1296$ units



40x40 pixel input image with 5x5 LCF
stride 1 pixel and no padding

Input layer (row order)

Hidden layer



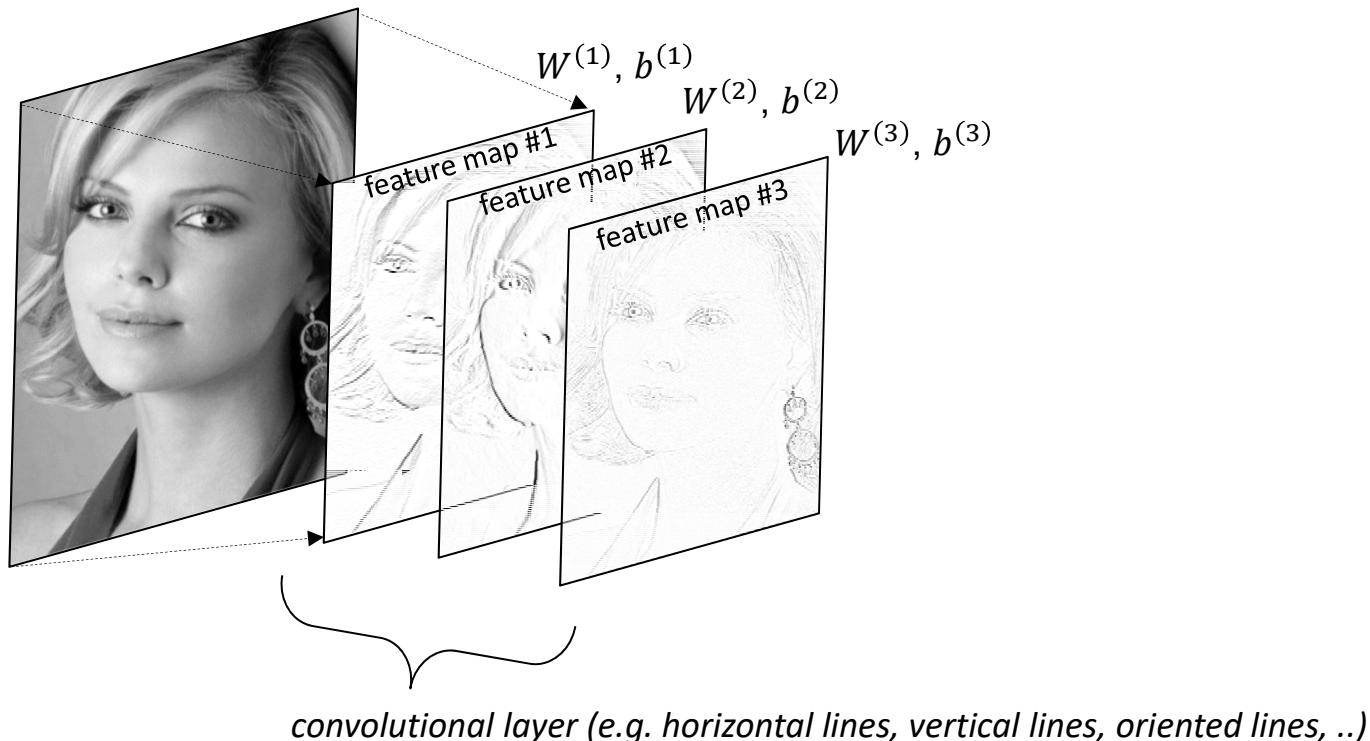
From neuron in the next rows

Wiring the net

Not fully
connected

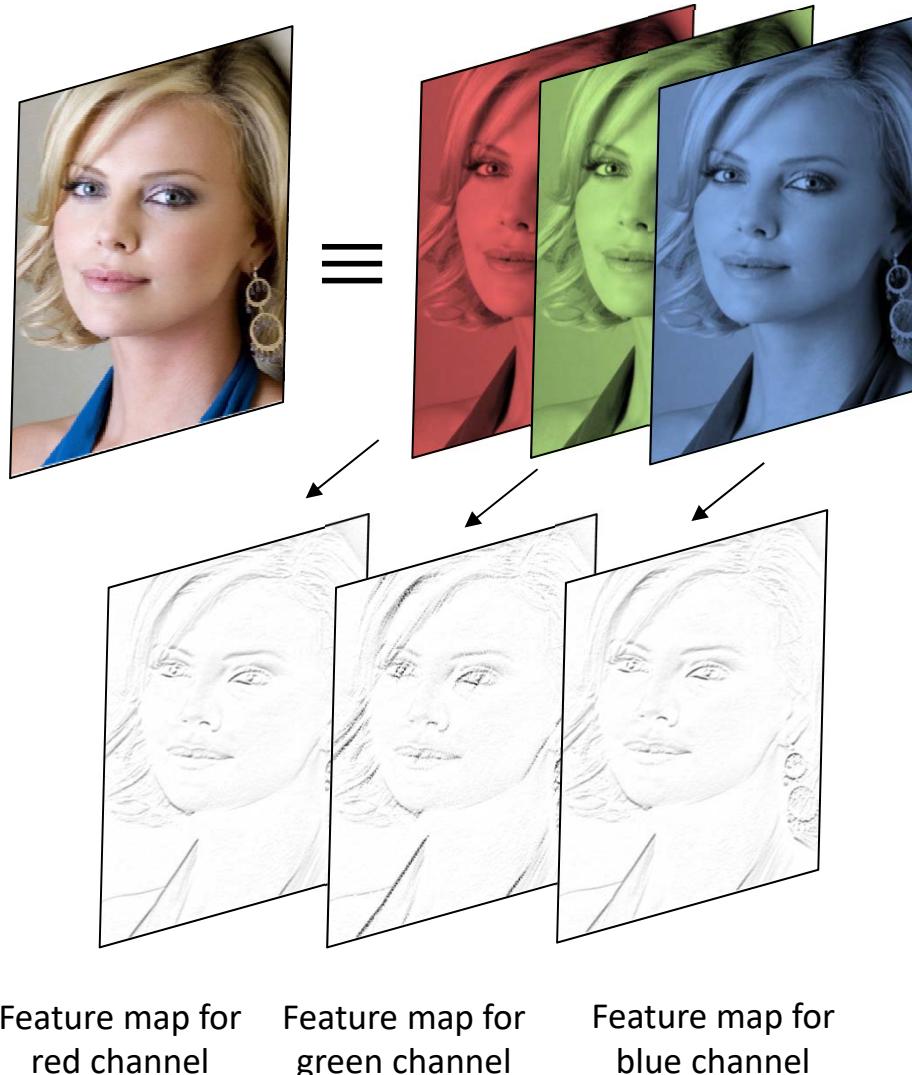
Many feature maps

- A *convolutional layer* can include several different feature maps, each corresponding to a filter
- E.g. $n \times n$ weights and a single shared bias identify one feature map
- For instance: oriented line detectors but they can be used to learn more complex patterns

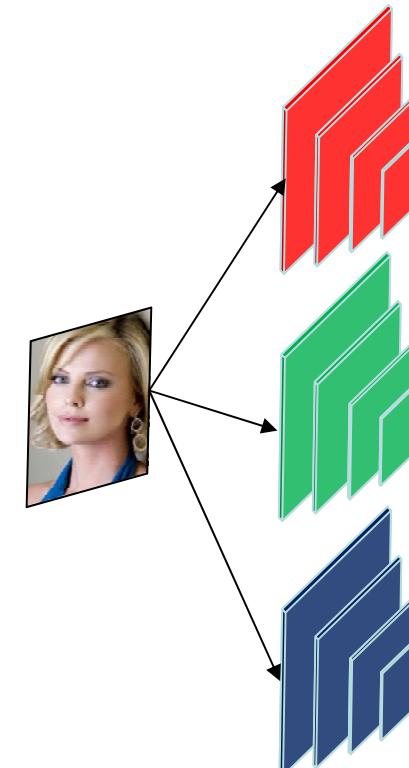


2D vs 3D convolution

RGB images



Three parallel lines of processing



DepthwiseConv2D

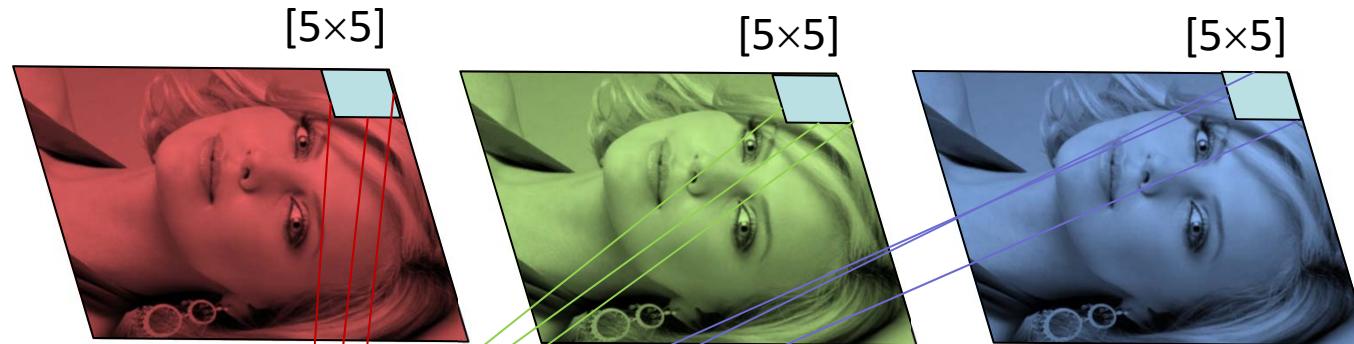
RGB images



40×40 pixel × 3 channels.

2D vs 3D convolution

??



1st neuron



Number of weights

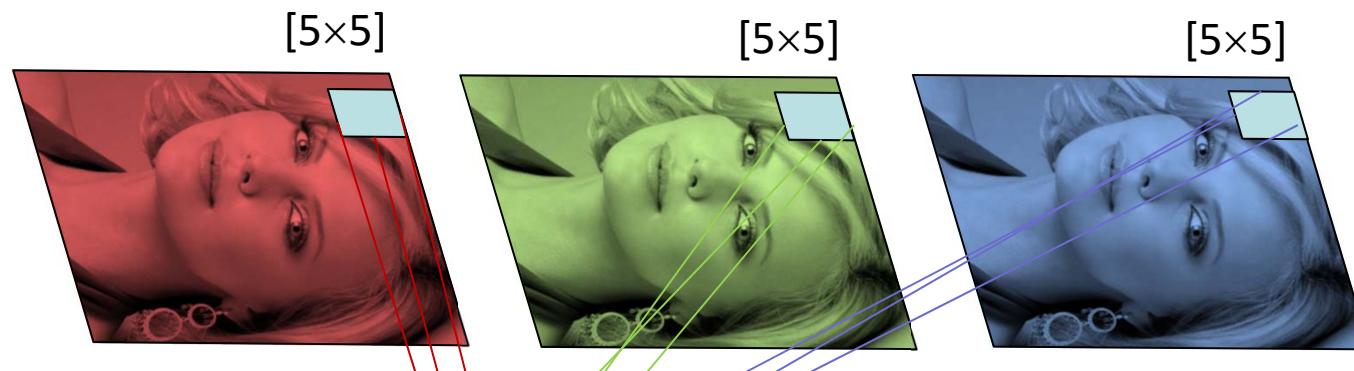
??

Convolutional layer
Feature map #1

Output of the 1st neuron in the 1st feature map in the 1st convolutional layer

2D vs 3D convolution

RGB images



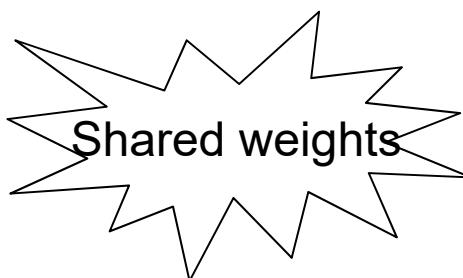
2nd neuron



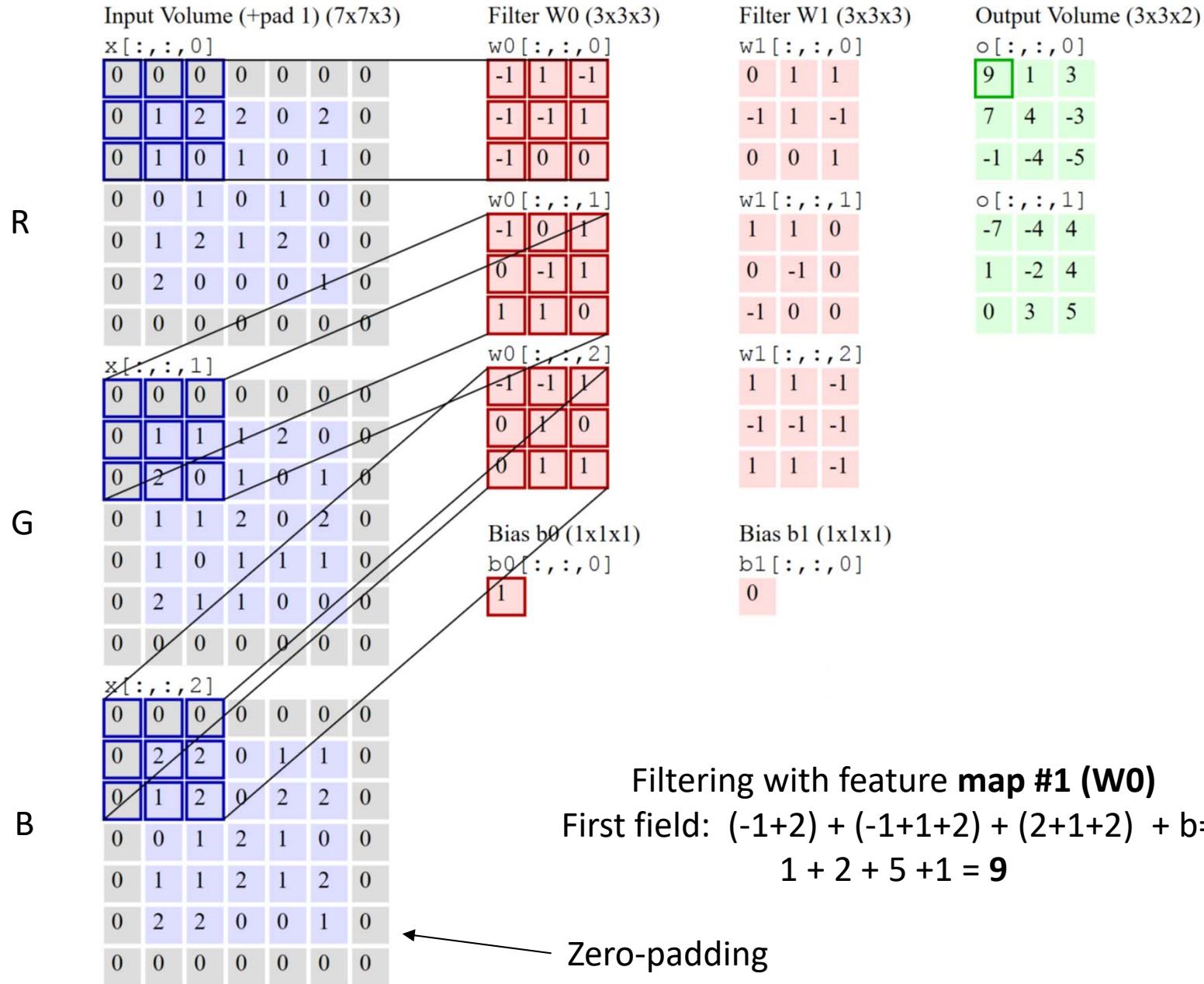
Number of weights
 $[5 \times 5 \times 3] + 1 = 76$

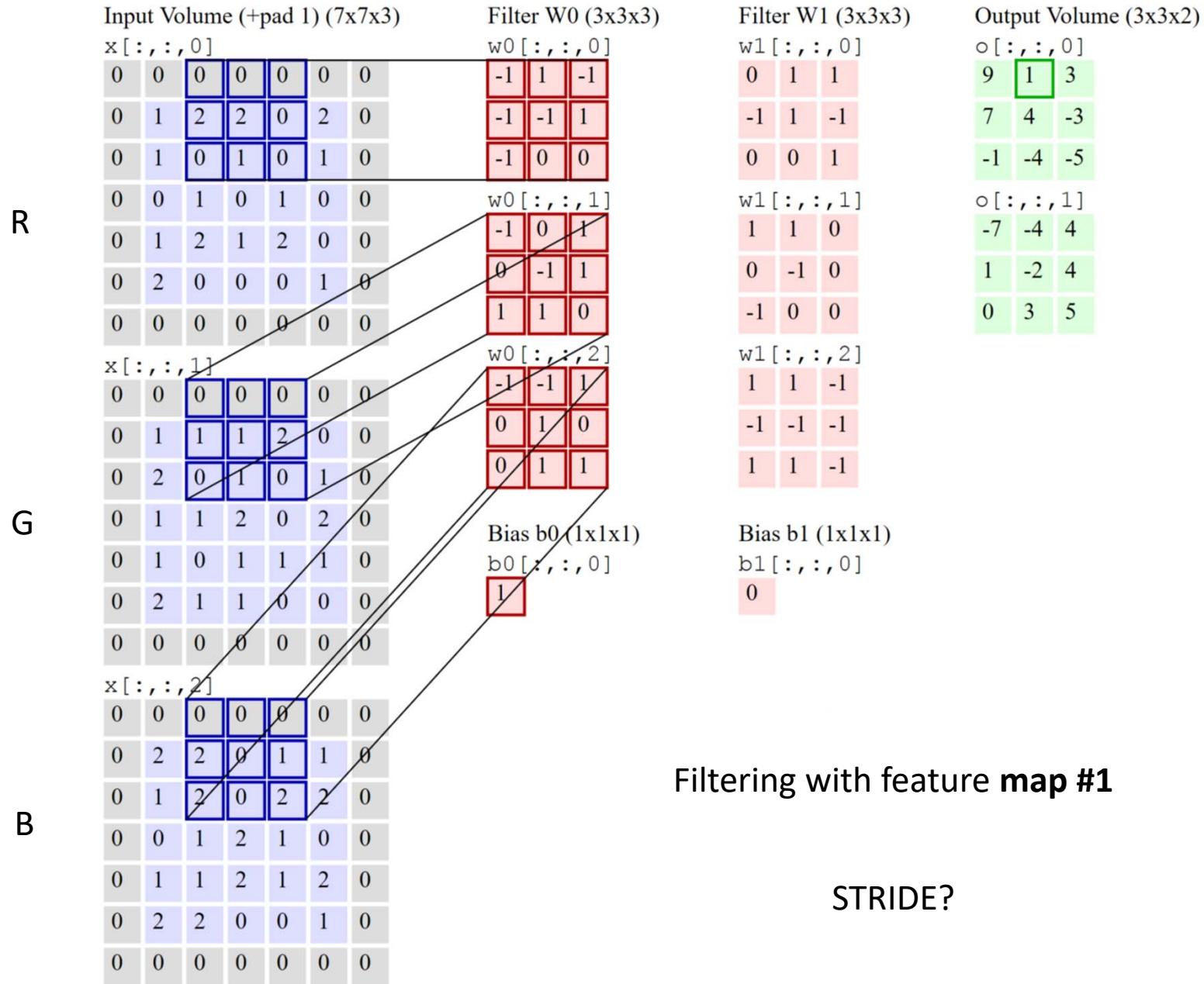
Convolutional layer
Feature map #1

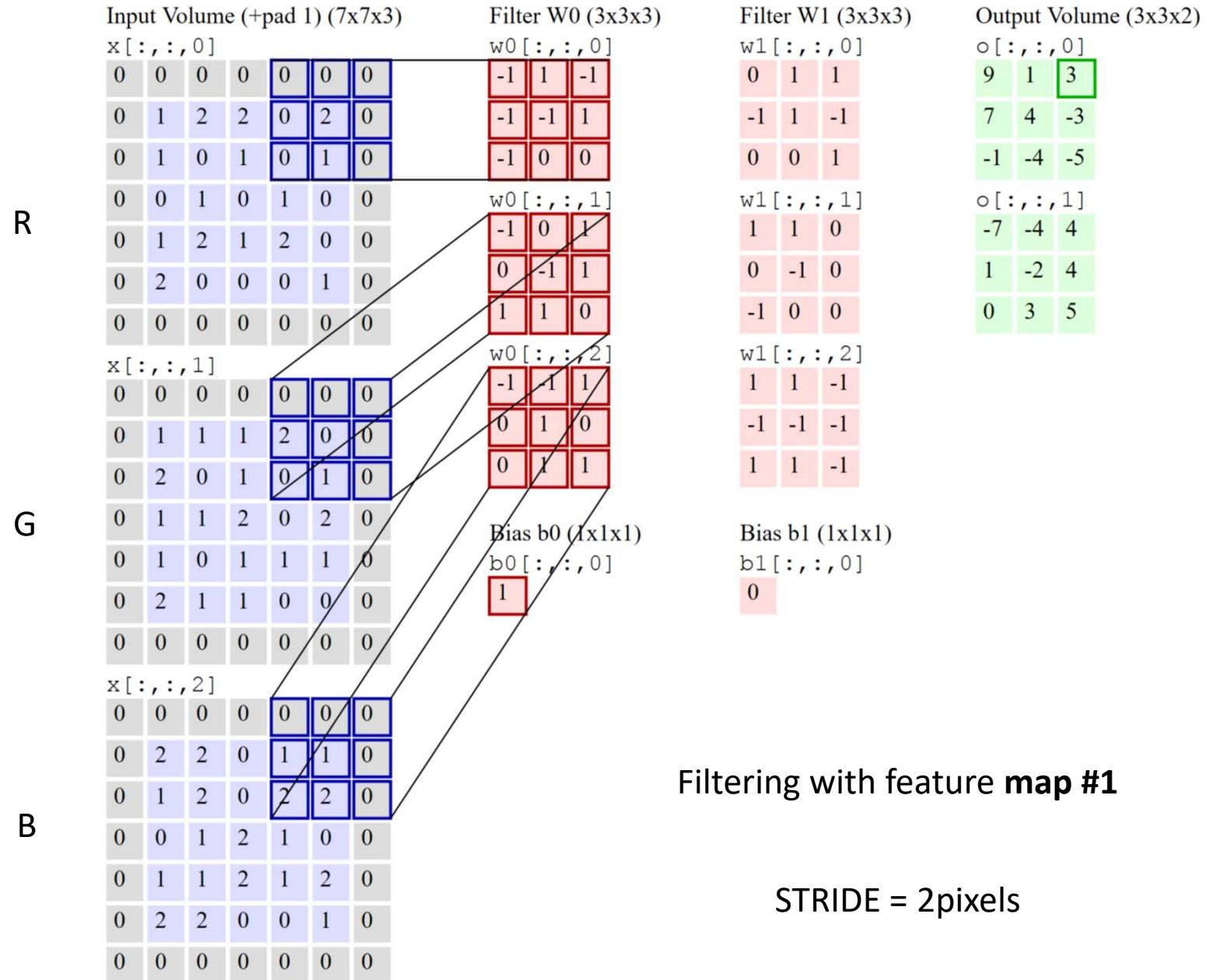
Output of the 2nd neuron in the 1st feature map in the 1st convolutional layer

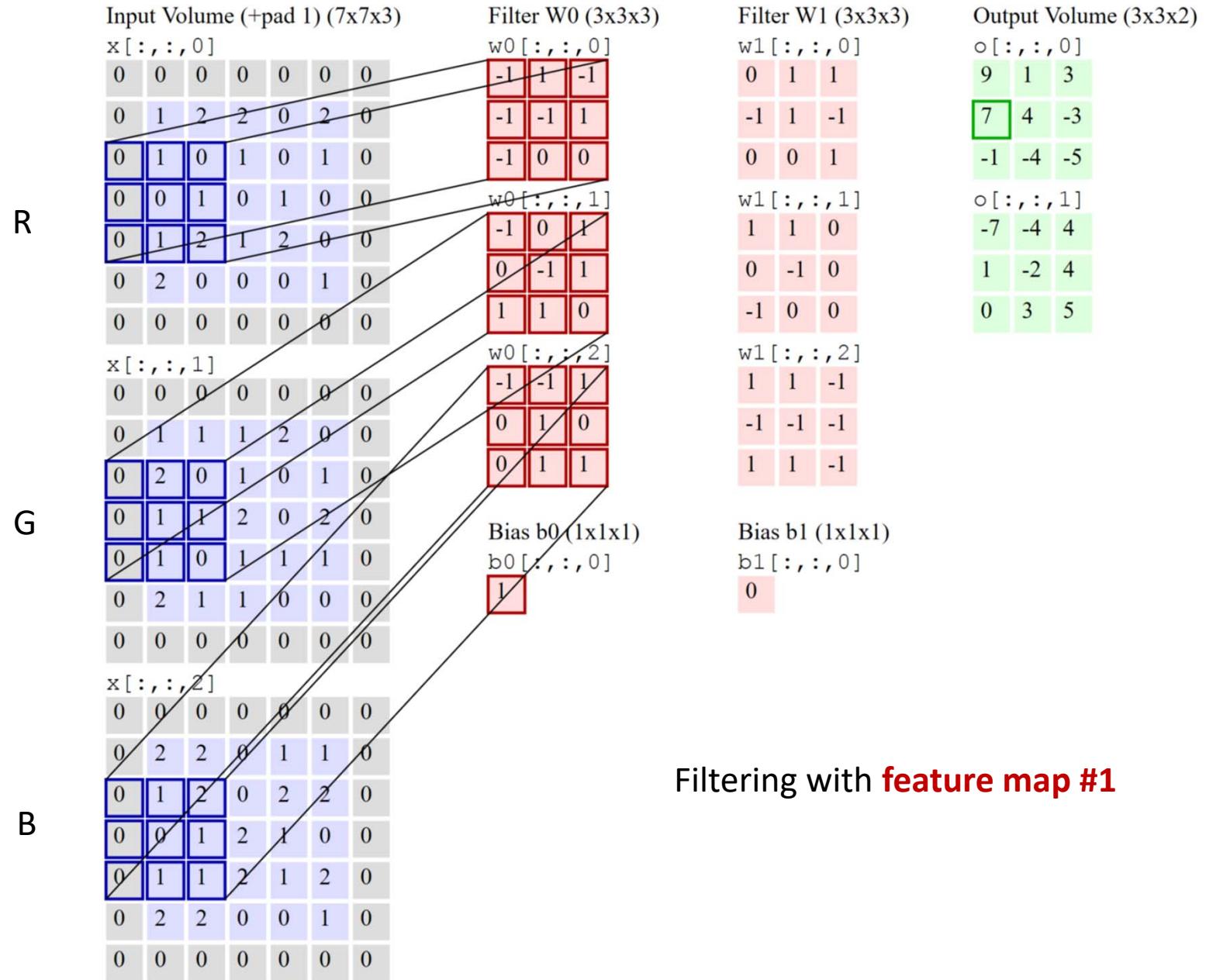


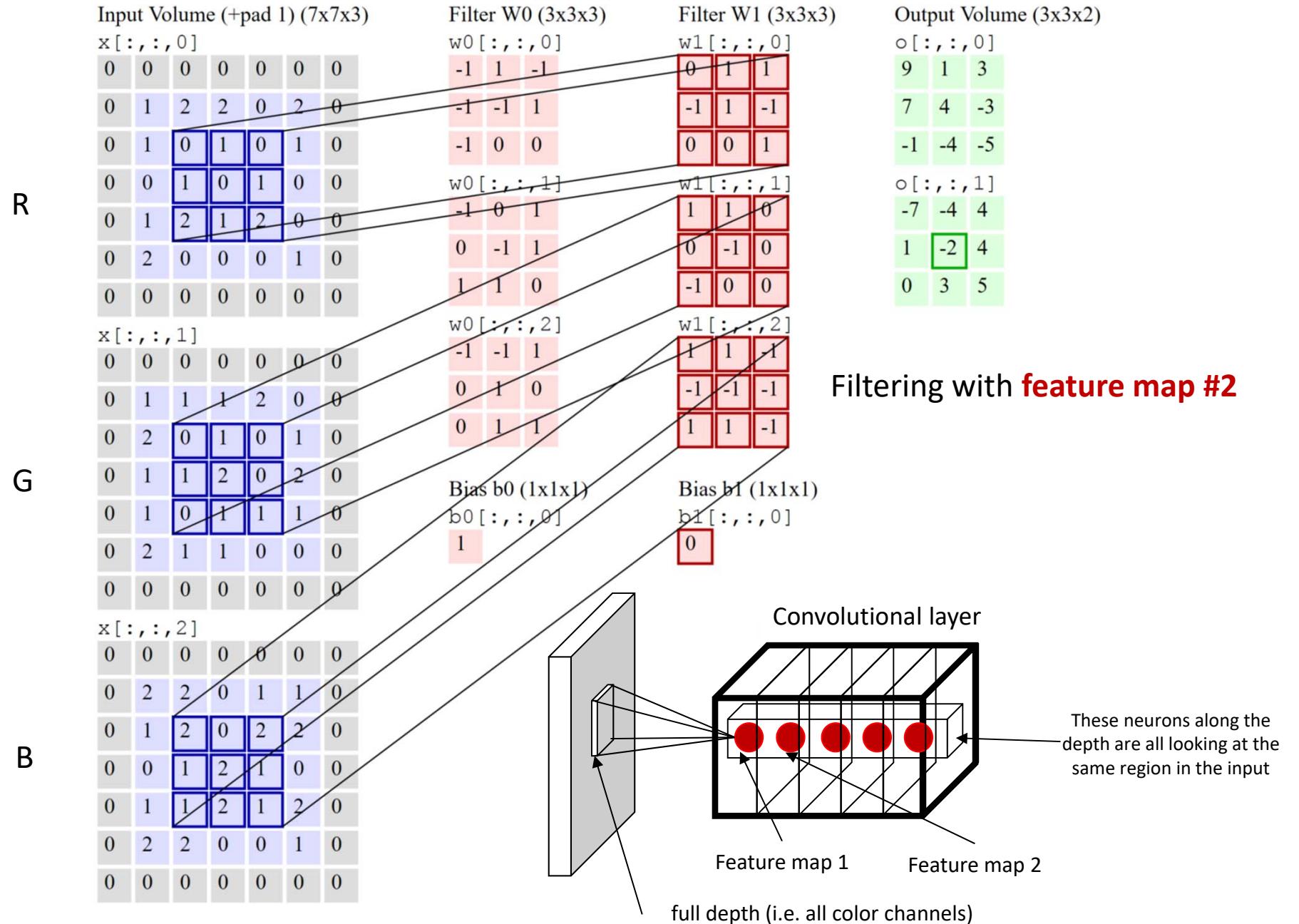
75 weights + 1 bias parameter







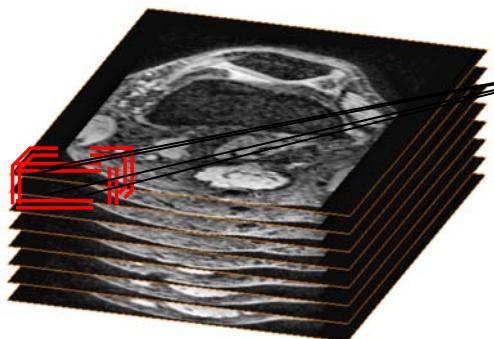




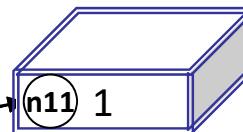
Volumetric biomedical images

3D scan: **33800 voxels**

Size of the input: $65 \times 65 \times 8$



Neuron #1 of the
feature map #1



50 weights +
1 threshold

How many neurons in the map?

Filter size: $5 \times 5 \times 2$

5 pixel in x direction

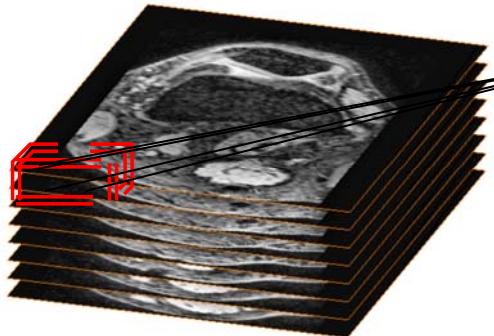
5 pixel in y direction

2 slices

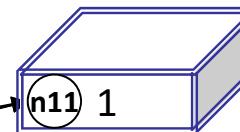
Volumetric biomedical images

3D scan: **33800 voxels**

Size of the input: $65 \times 65 \times 8$



Neuron #1 of the
feature map #1



50 weights +
1 threshold

How many neurons in the map?

$$s_x = \frac{(I_x + 2z_x - f_x)}{L_x} + 1 = 31$$

$$s_y = \frac{(I_y + 2z_y - f_y)}{L_y} + 1 = 31$$

$$s_z = \frac{(n_s + 2z_z - f_z)}{L_z} + 1 = 4$$

3844

Filter size: $5 \times 5 \times 2$

5 pixel in x direction

5 pixel in y direction

2 slices

Padding: none

Stride: 2 voxel

I_x and I_y : input image sizes in x and y directions

z_x , z_y and z_z : zero-padding sizes

f_x , f_y and f_z : filtering sizes

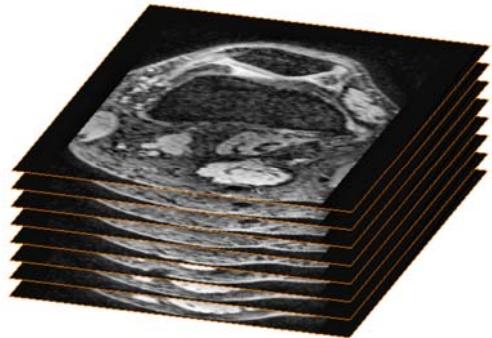
L_x , L_y , and L_z : stride lengths

n_s : number of slices

Volumetric biomedical images

3D scan: **33800 voxels**

Size of the input: $65 \times 65 \times 8$



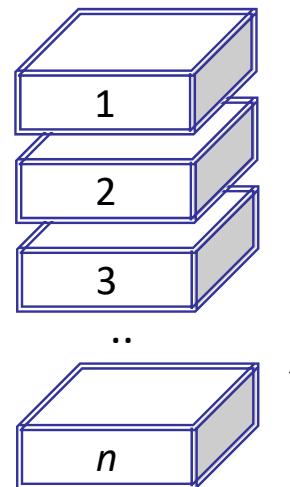
Filter size: $5 \times 5 \times 2$

5 pixel in x direction

5 pixel in y direction

2 slices

First convolutional
layer



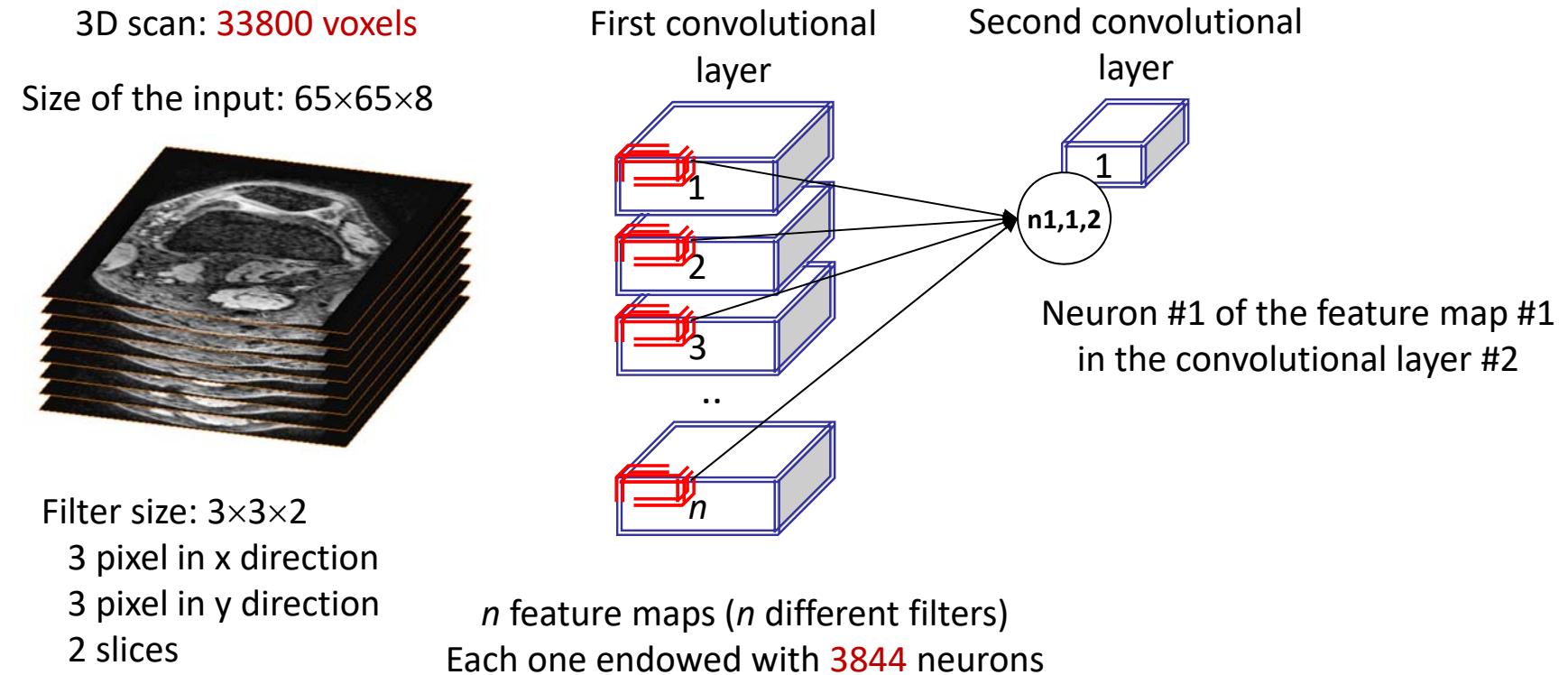
How many free parameter in
the layer?

$51 * n$

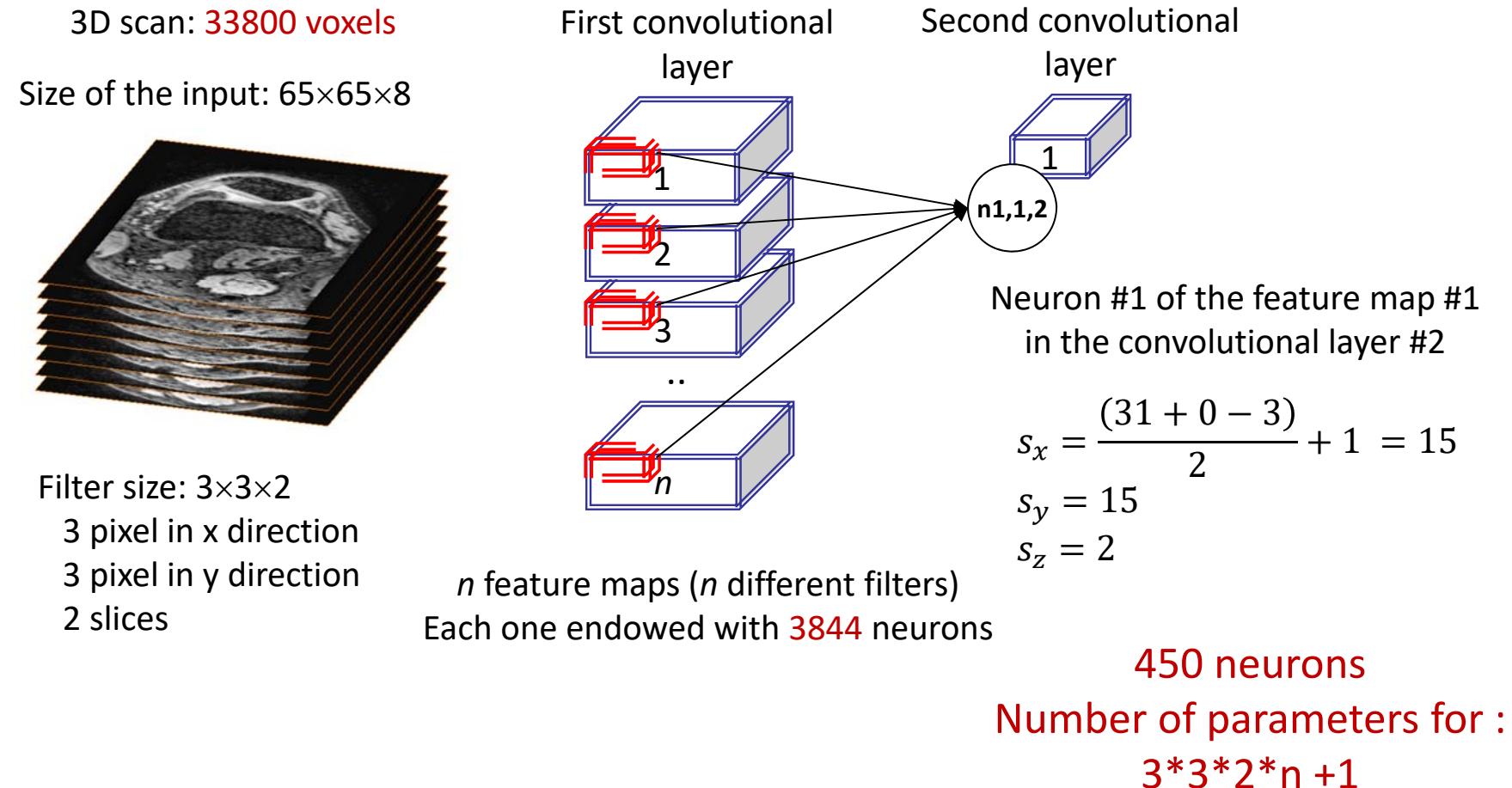
and the next layer?

n feature maps (n different filters)
Each one endowed with **3844** neurons

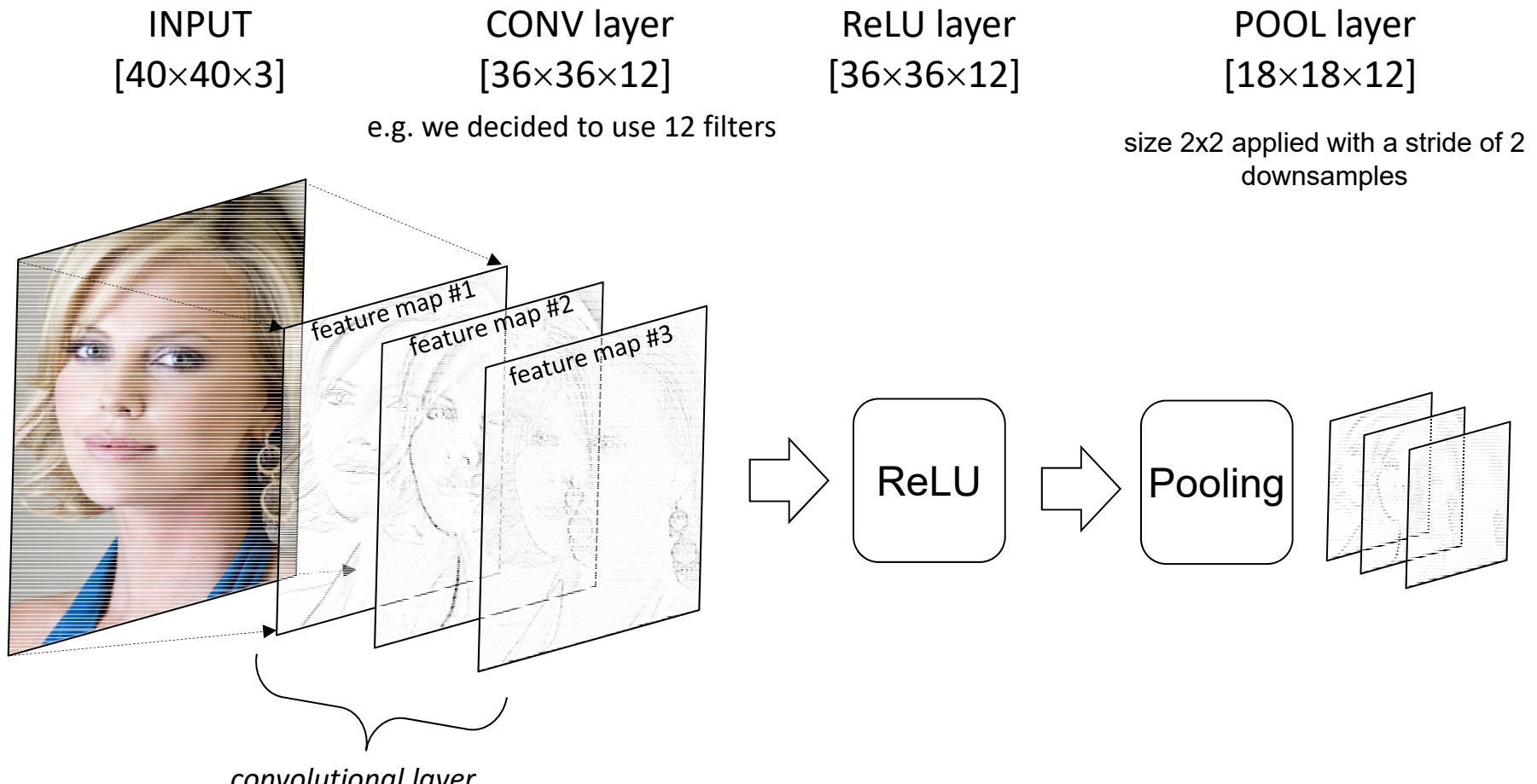
Volumetric biomedical images



Volumetric biomedical images

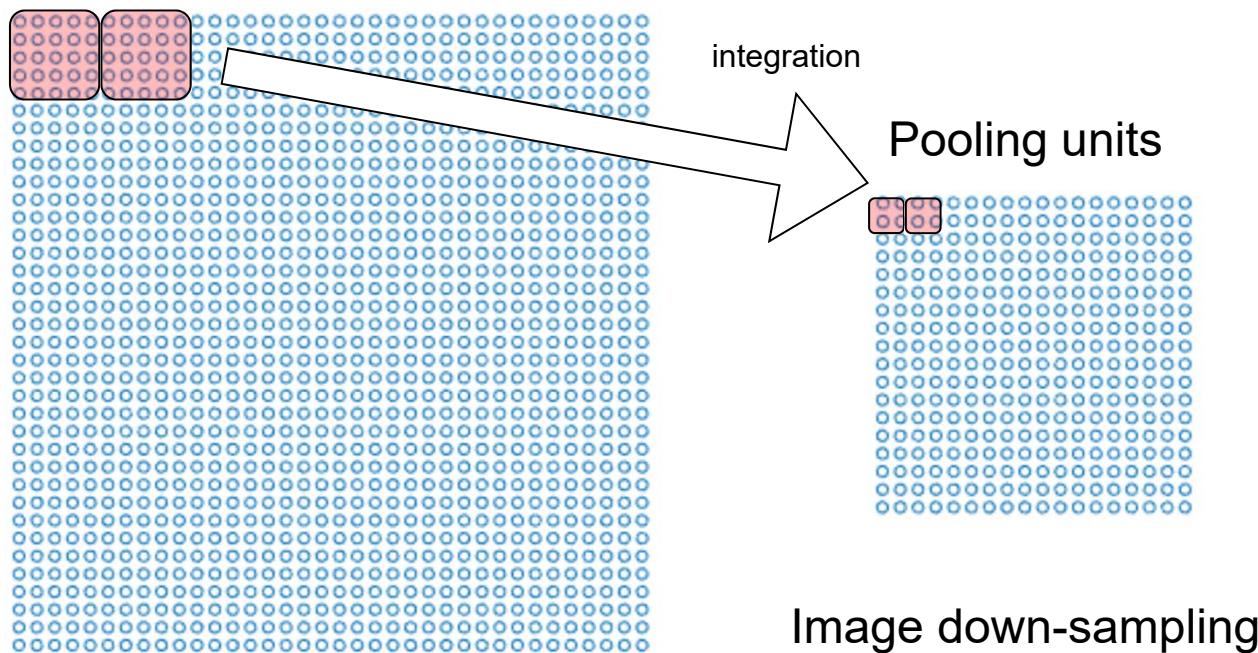


Conv+ReLU+Pooling

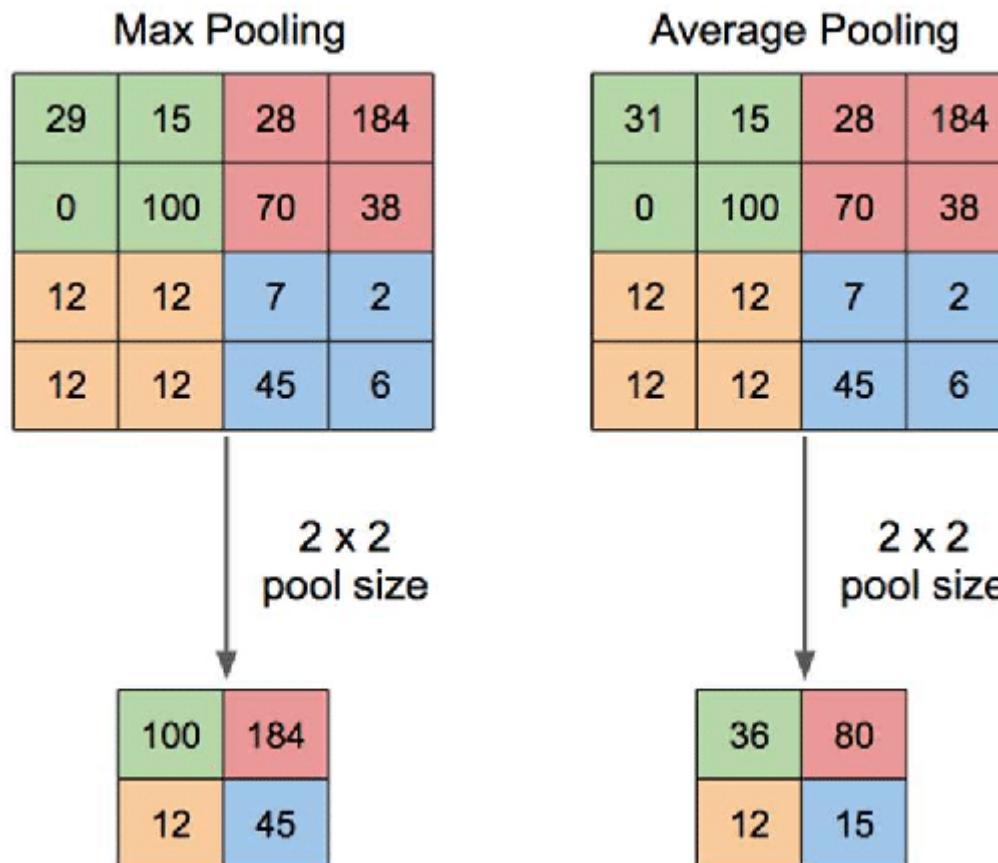


Pooling layer (integration/smoothing)

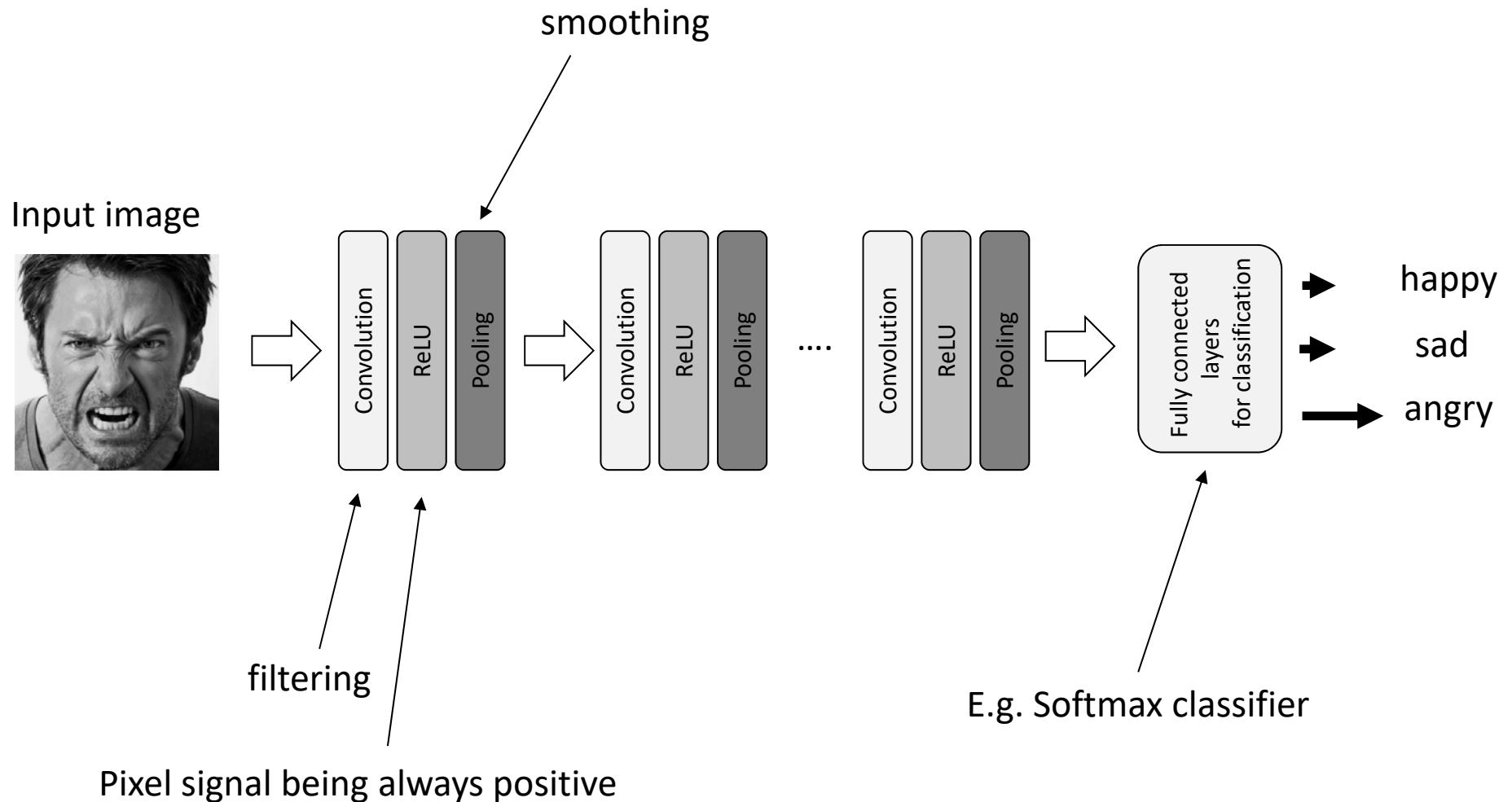
Hidden layer (output after feature map and ReLU)



Pooling layer (integration/smoothing)



ANN able to recognize emotions



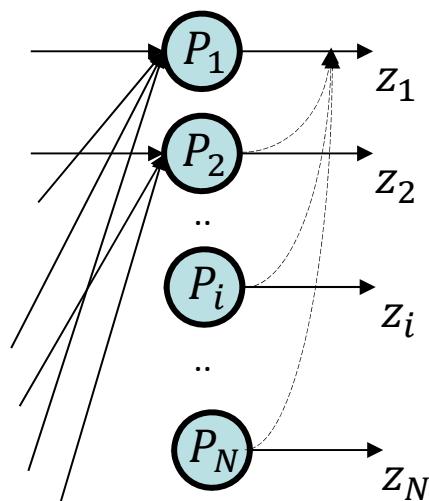
Softmax classifier

Only two classes: traditional logsig classifier

- One output neuron with sigmoidal activation
- Training: Delta/Backprop rule applied to error function

More than two classes: Softmax classifier

- Normalized exponential activation
- Number of output neurons as much as the number of classes
- Training: cost function is the negative log probability of the correct label and cross-entropy regime (very similar to KL divergence)
- Expected output signals:
 - 1 0 0 0 .. 0: class #1
 - 0 1 0 0 .. 0: class #2
 - 0 0 1 0 .. 0: class #3
 - ..
 - ..
 - 0 0 0 0.. 1: class #N



“maps” a N -dimensional vector of arbitrary real values to a N -dimensional vector of real values in the range $(0, 1)$ that add up to 1.

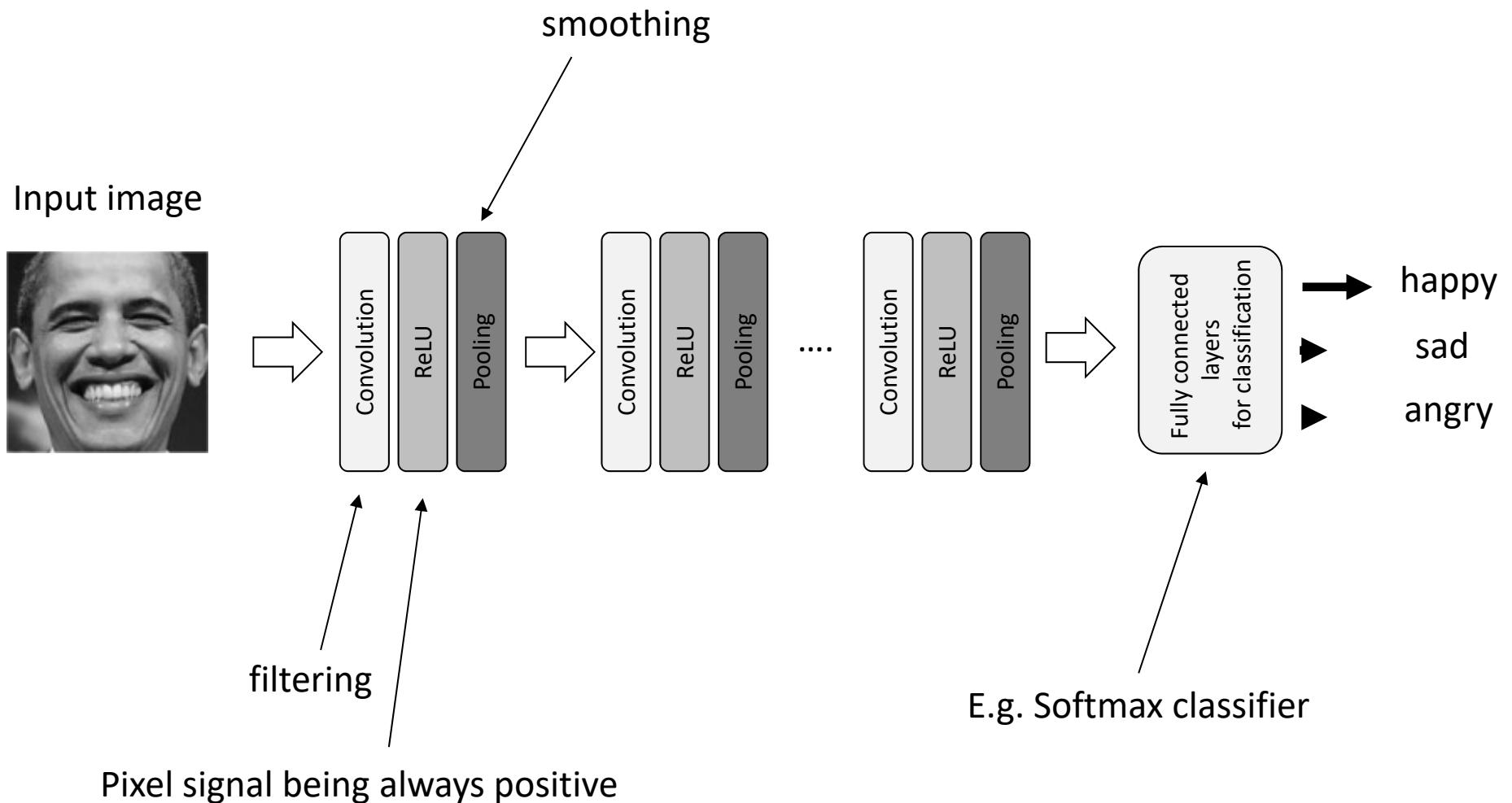
$$z_i = \frac{e^{P_i}}{\sum_j^N e^{P_j}}$$

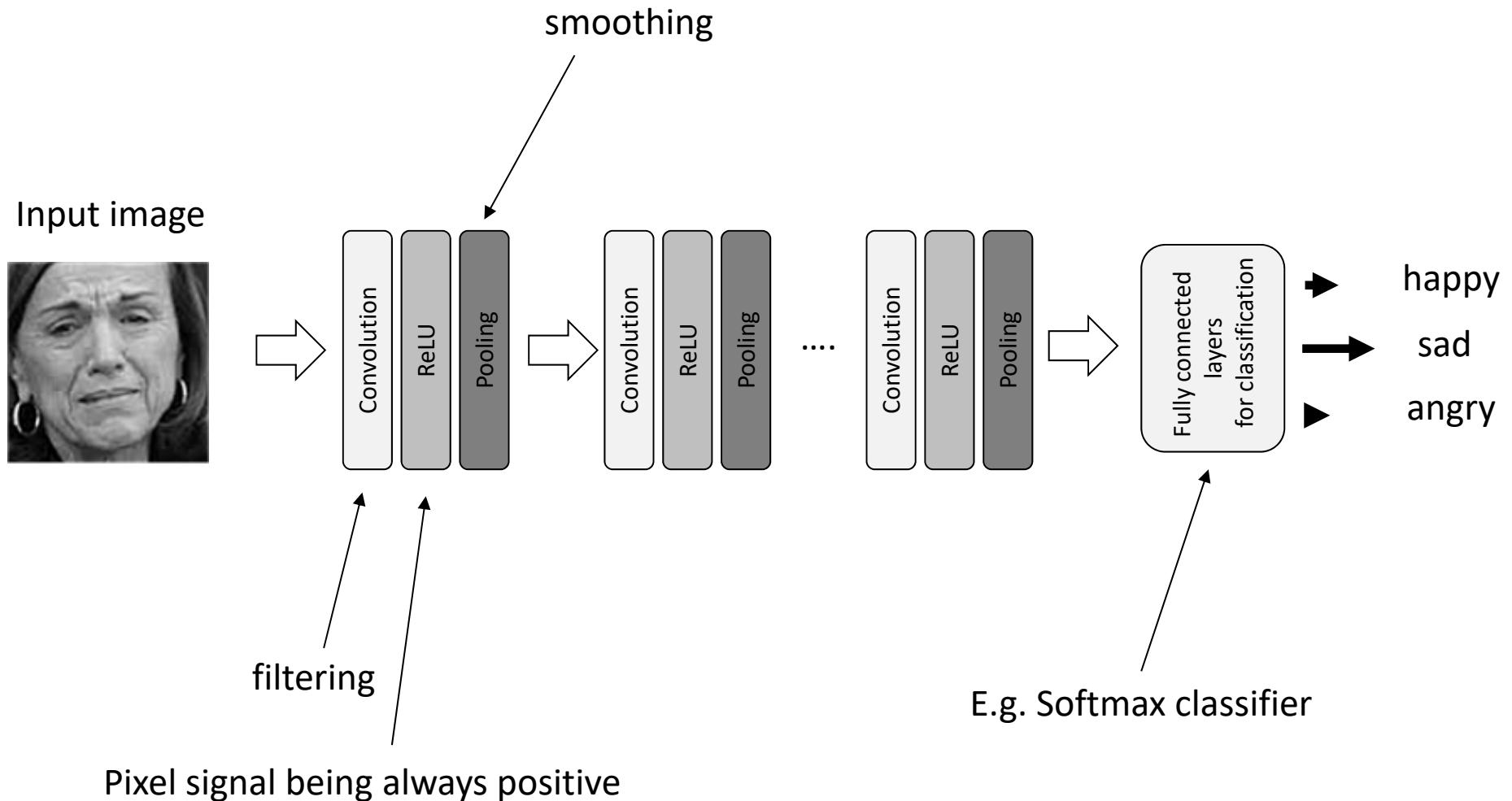
$$\frac{\partial z_i}{\partial P_i} = z_i (1 - z_i)$$

$$E = - \sum_j^N t_j \ln z_j$$

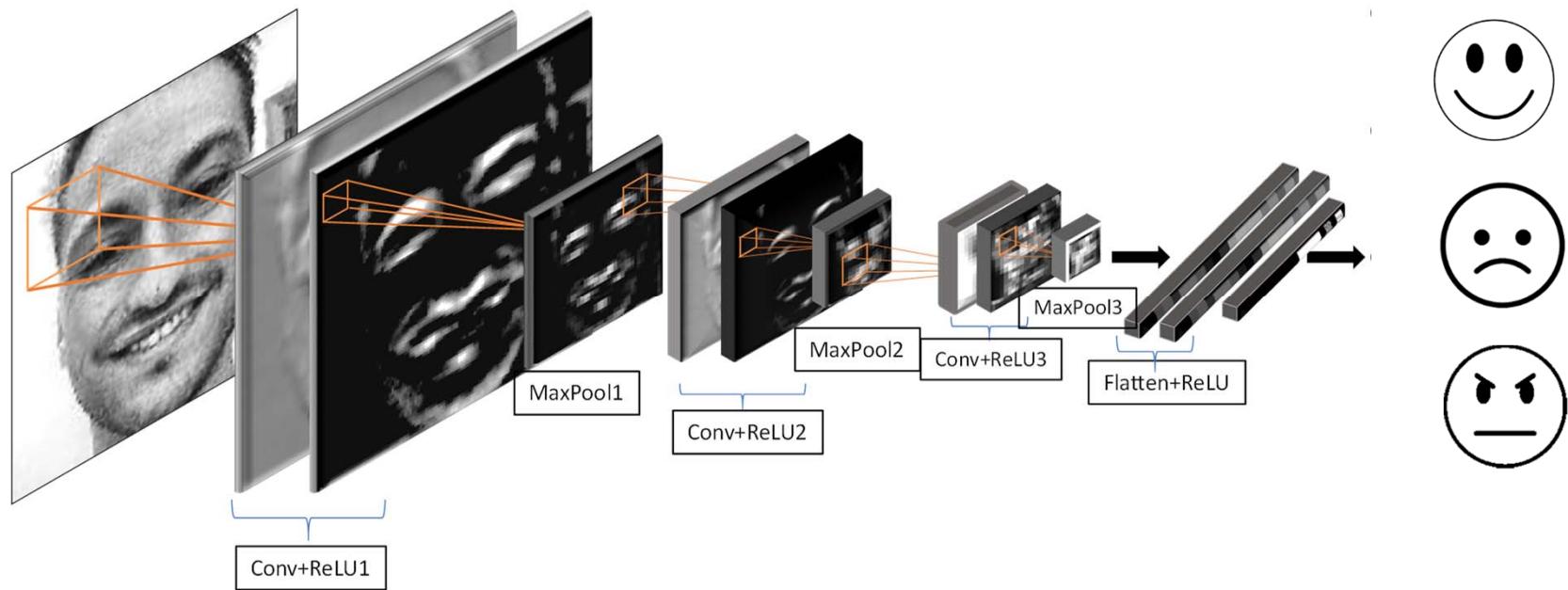
$$\frac{\partial E}{\partial P_i} = \sum_j^N \frac{\partial E}{\partial z_j} \frac{\partial z_j}{\partial P_i} = z_i - t_i$$

Apply deep analysis to classification problem





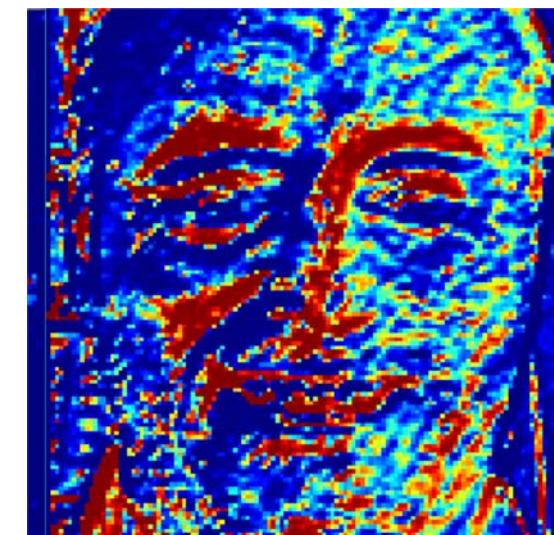
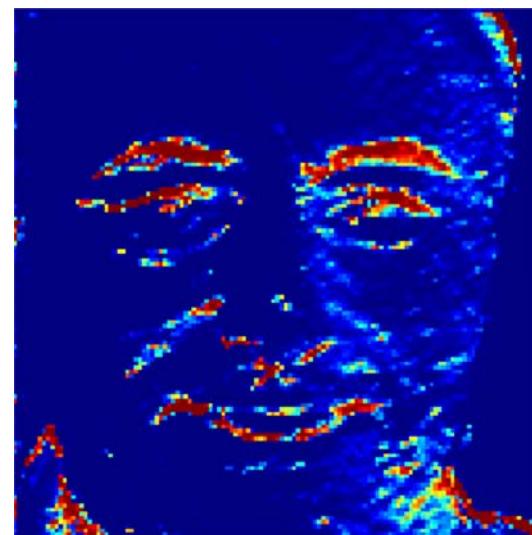
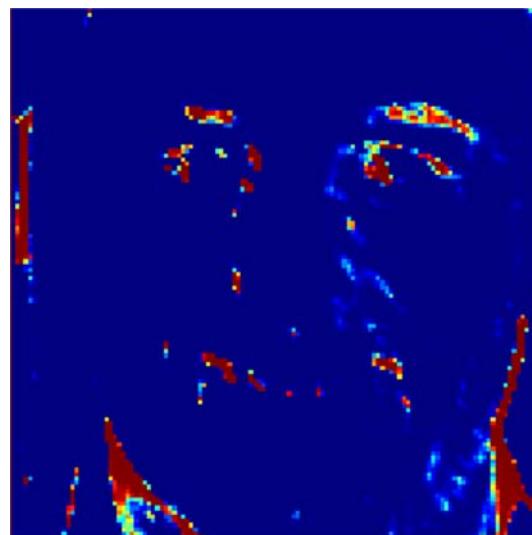
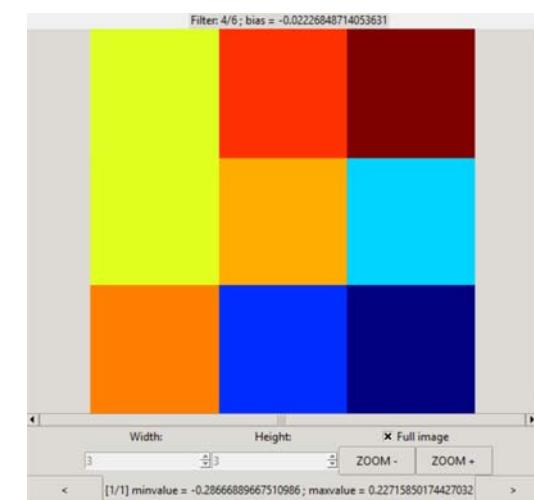
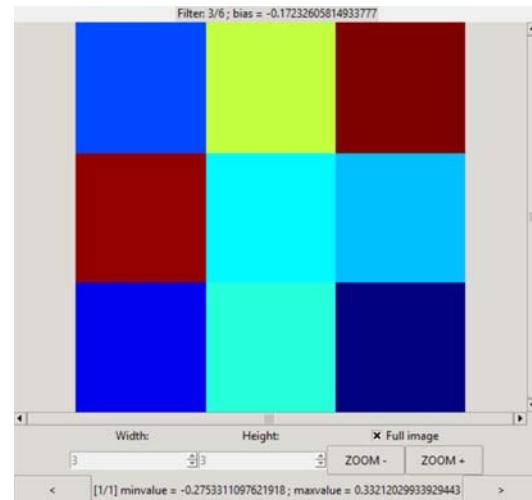
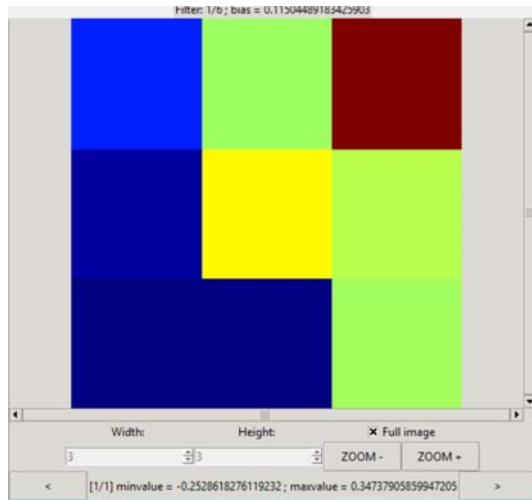
ANN able to recognize emotions



Training data: image with the corresponding label

Loss function for multi-class classification: Categorical Cross-Entropy

Three feature maps of the 1st layer



Three feature maps of the 2nd layer

