

23/10/2024

# Convolutional Neural Networks

Giacomo Boracchi,  
DEIB, Politecnico di Milano

[giacomo.boracchi@polimi.it](mailto:giacomo.boracchi@polimi.it)

<https://boracchi.faculty.polimi.it/>

# The Feature Extraction Perspective

# The Feature Extraction Perspective

Images can not be directly fed to a classifier

We need some intermediate step to:

- Extract meaningful information (to our understanding)
- Reduce data-dimension

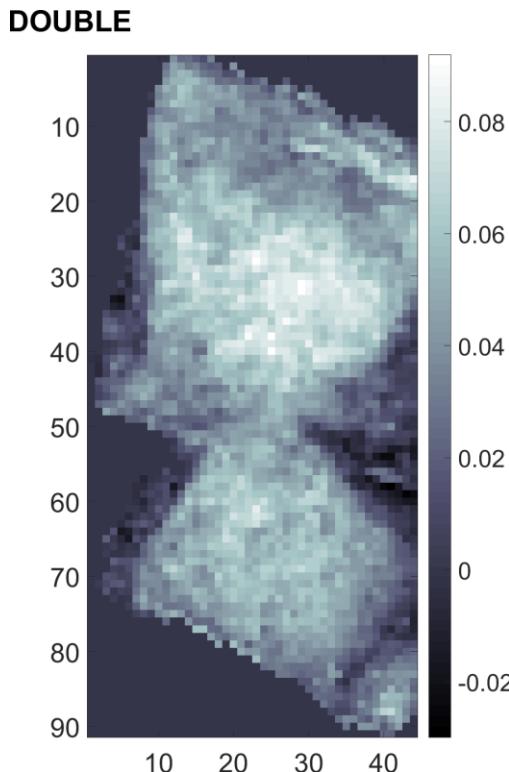
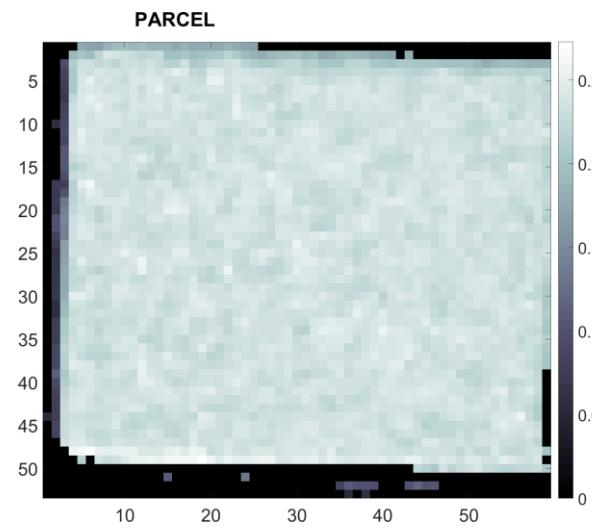
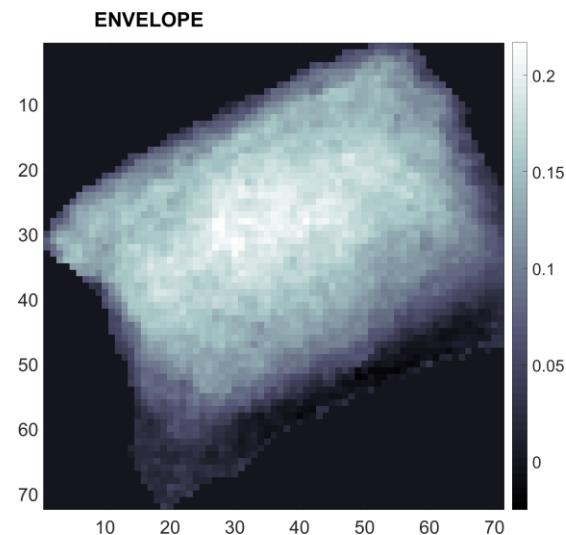
We need to extract features:

- The better our features, the better the classifier

# Example of Hand-Crafted Features

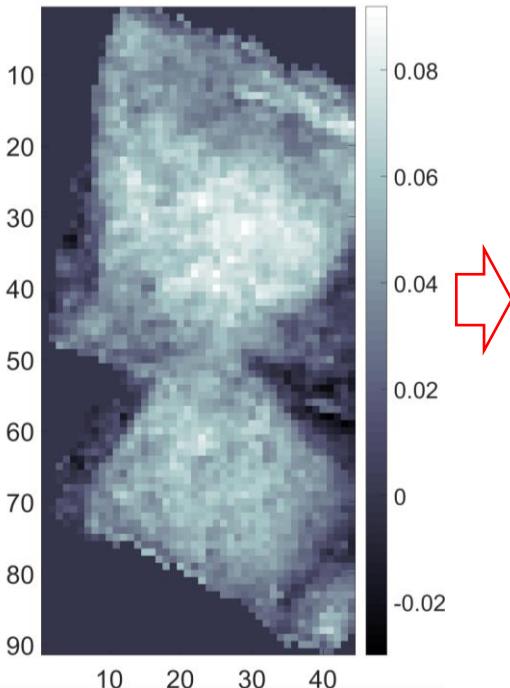
Example of features:

- Average height
- Area (coverage with nonzero measurements)
- Distribution of heights
- Perimeter
- Diagonals



# Neural Networks

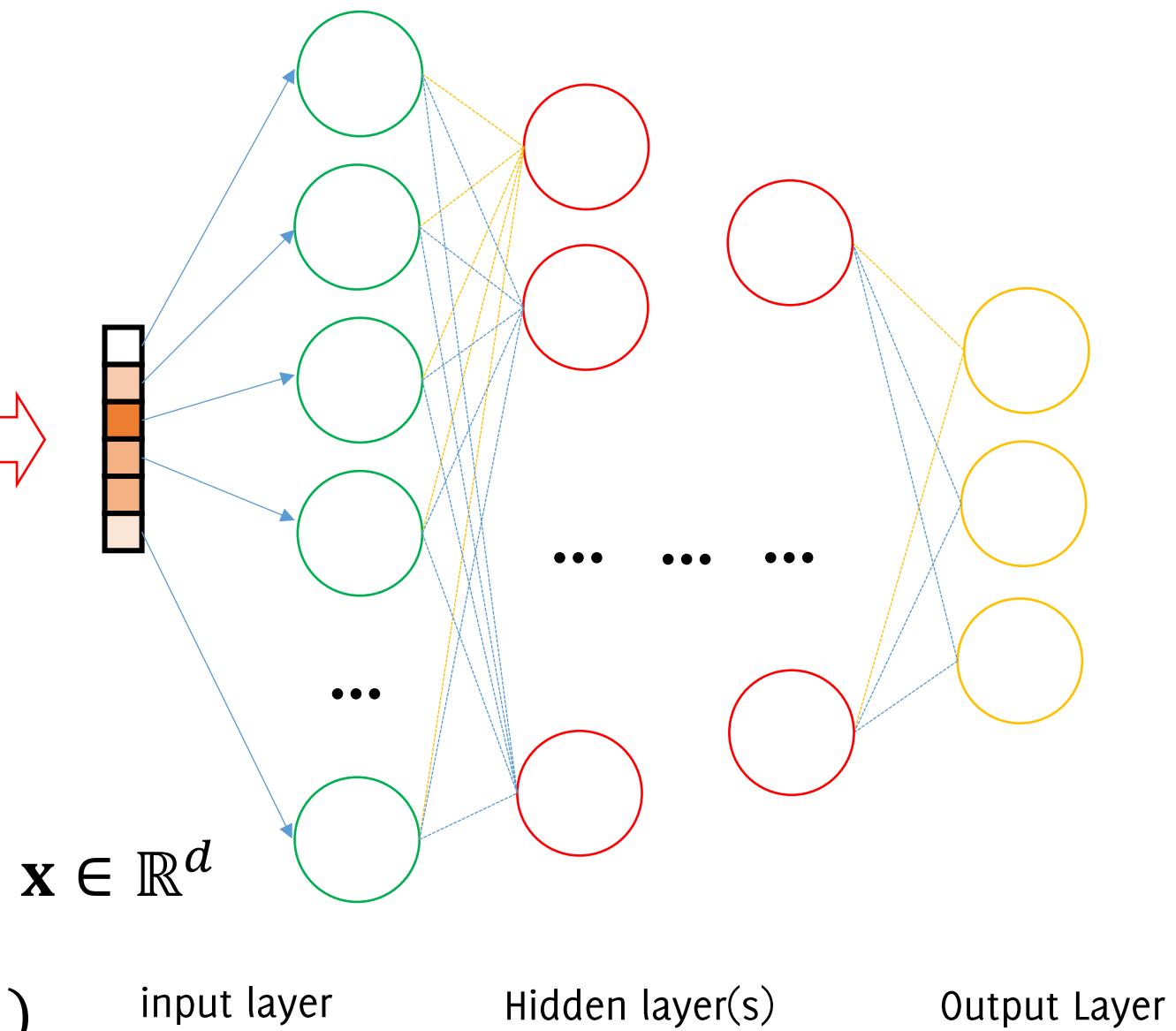
Input image



$$I_1 \in \mathbb{R}^{r_1 \times c_1}$$

Feature Extraction Algorithm

$$(d \ll r_1 \times c_1)$$



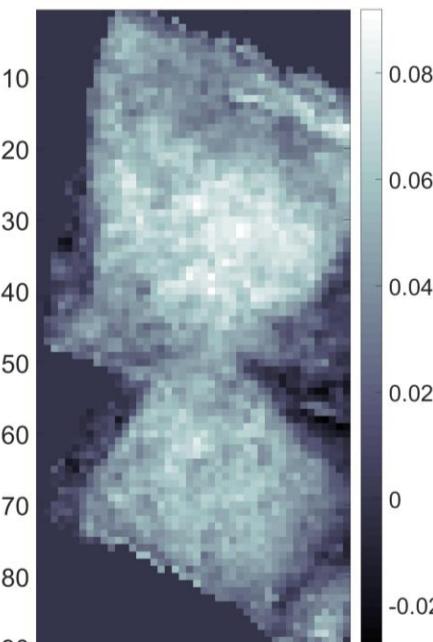
input layer

Hidden layer(s)

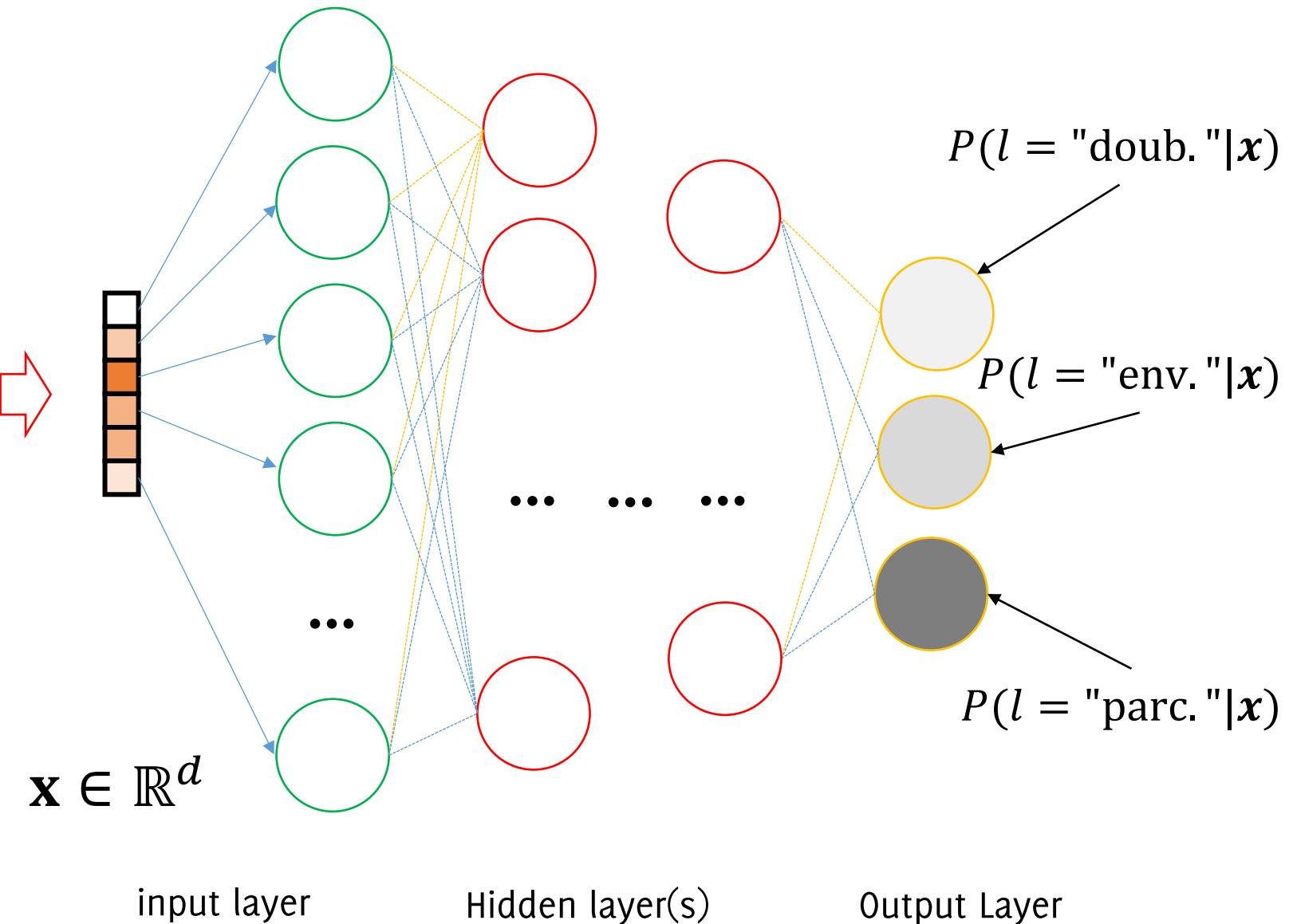
Output Layer

# Neural Network: After Training

Input image



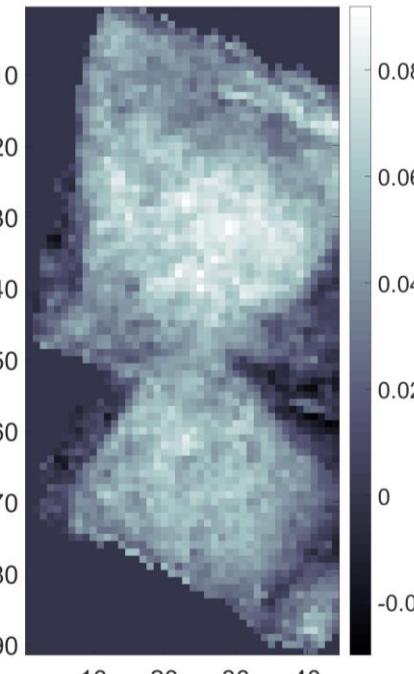
Feature Extraction Algorithm



# Image Classification by Hand Crafted Features



Input image

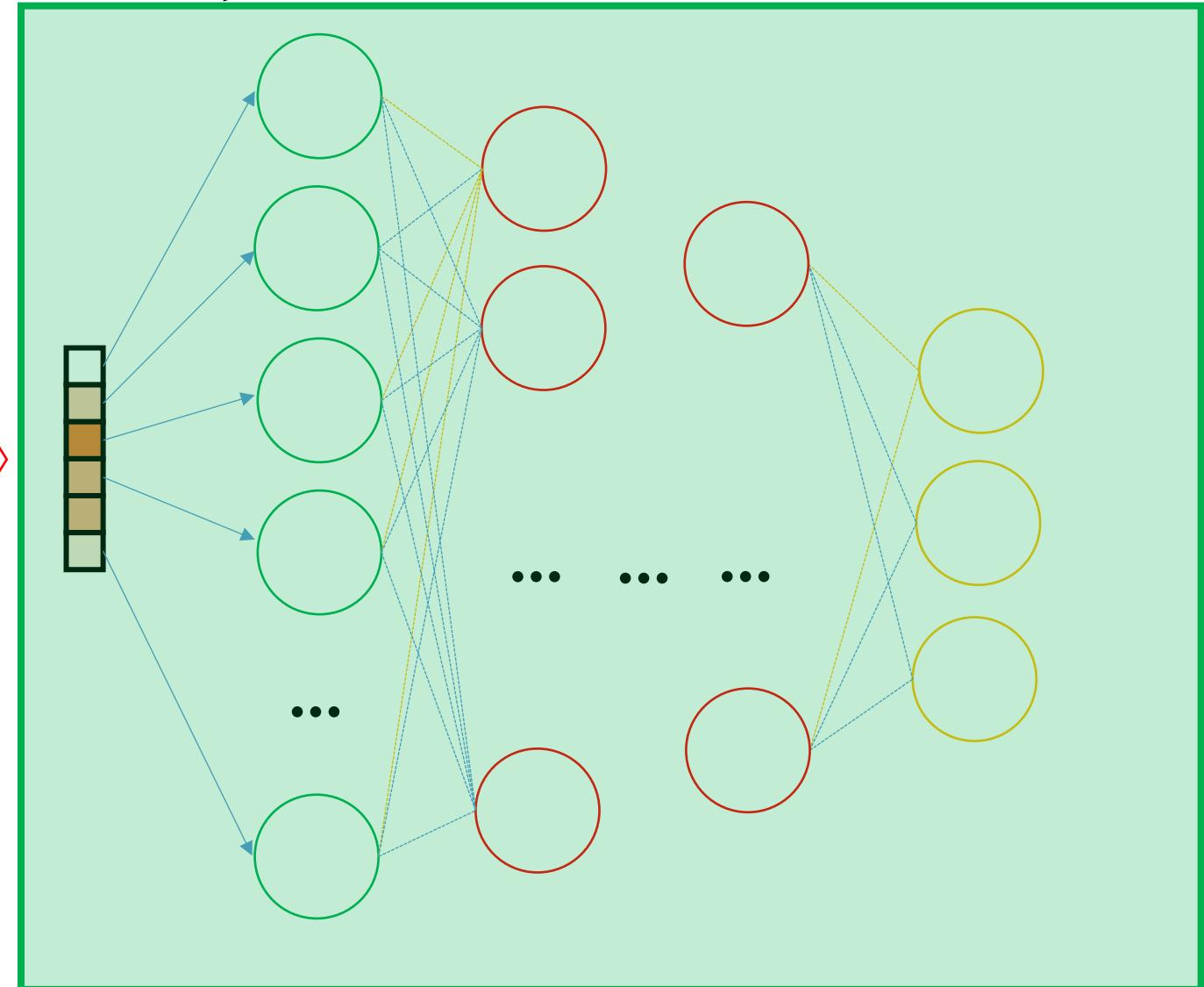


$$I_1 \in \mathbb{R}^{r_1 \times c_1}$$

Feature Extraction Algorithm



Hand Crafted



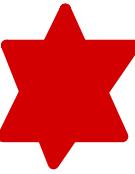
Data Driven

# Hand Crafted Features, pros:



- Exploit *a priori* / expert information
- Features are **interpretable** (you might understand why they are not working)
- You can **adjust features** to improve your performance
- **Limited amount of training data** needed
- You can give more relevance to some features

# Hand Crafted Features, cons:



- Requires a lot of **design/programming** efforts
- **Not viable** in many **visual recognition** tasks (e.g. on natural images) which are easily performed by humans
- **Risk of overfitting** the training set used in the design
- **Not very general** and "portable"

# Data-Driven Features

... the advent of deep learning

# Data-Driven Features



Input image

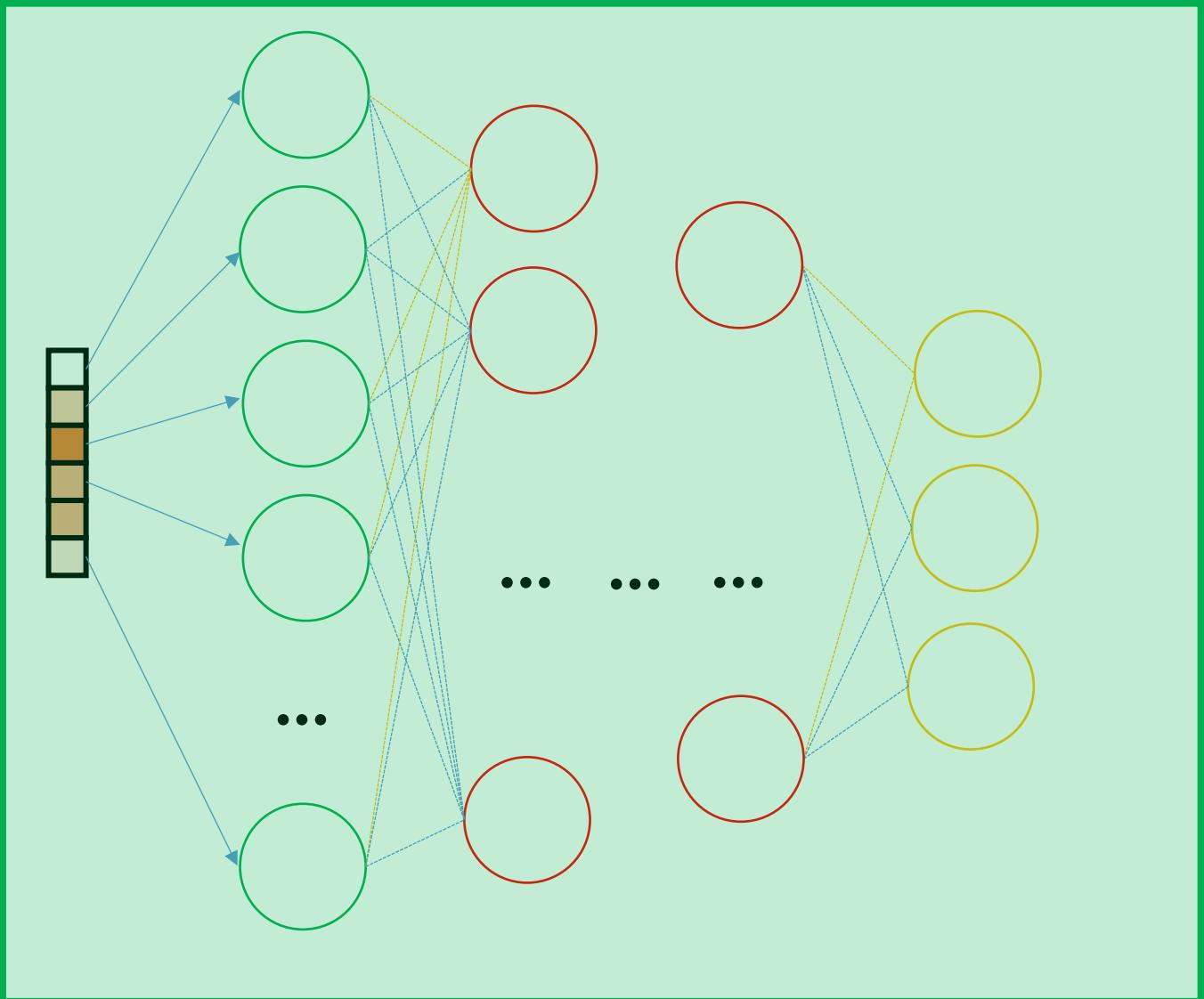


$$I_1 \in \mathbb{R}^{r_1 \times c_1}$$



Data Driven

Feature Extraction



Data Driven

# Convolutional Neural Networks

Setting up the stage

# Local Linear Filters RECAP

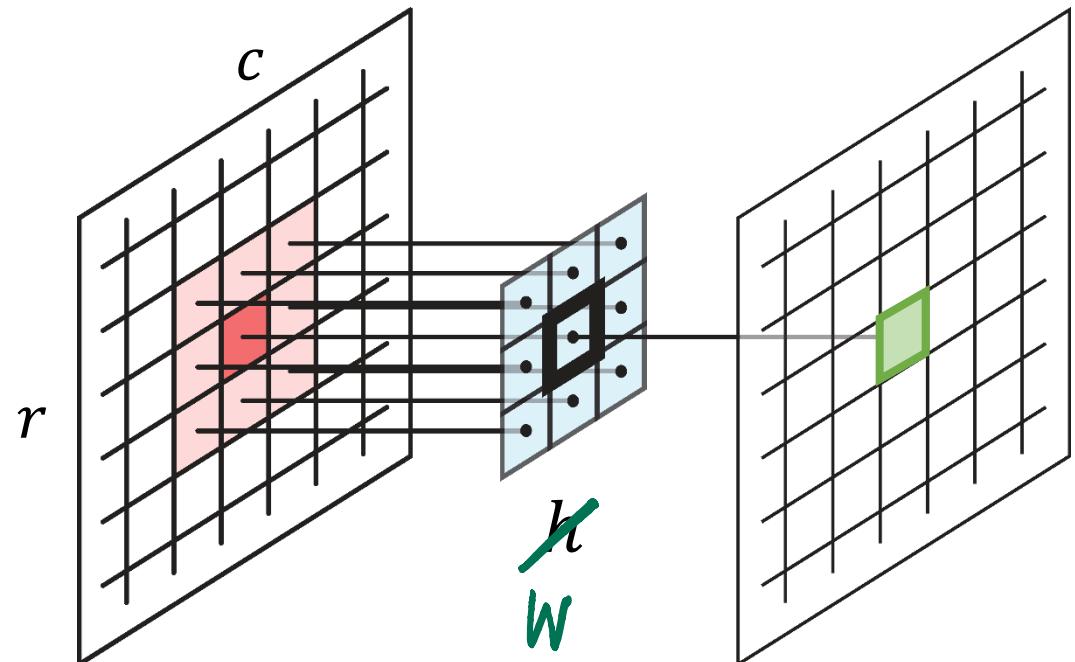
**Linear Transformation:** Linearity implies that the **output**  $T[I](r, c)$  is a linear combination of the pixels in  $U$ :

$$T[I](r, c) = \sum_{(u,v) \in U} w_i(u, v) * I(r + u, c + v)$$

Considering some weights  $\{w_i\}$

We can consider weights as an image, or a filter  $h$

The filter  $h$  entirely defines this operation



# Local Linear Filters RECAP

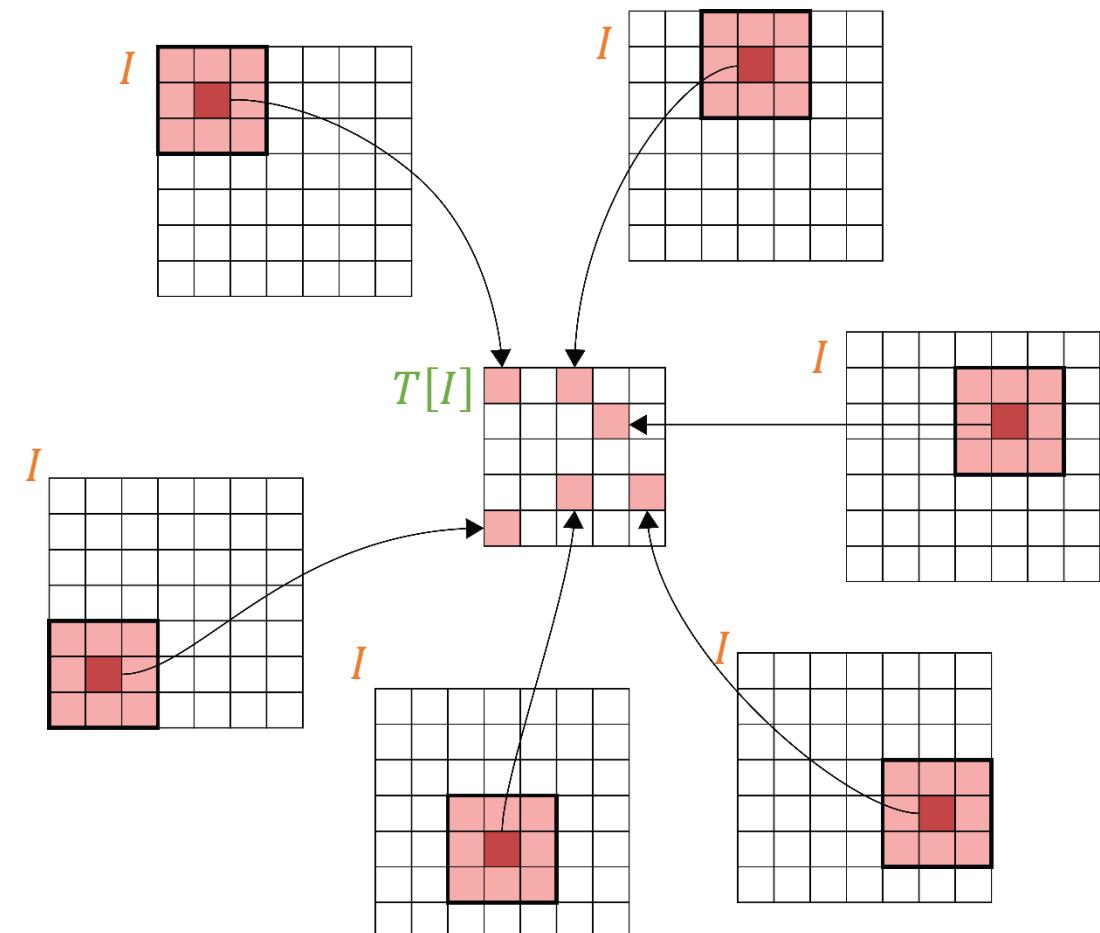
Linear Transformation: the filter weights can be associated to a matrix  $w$

$$T[I](r, c) = \sum_{(u,v) \in U} w_i(u, v) * I(r + u, c + v)$$

$w$

$w(-1, -1)$	$w(-1, 0)$	$w(-1, 1)$
$w(0, -1)$	$w(0, 0)$	$w(0, 1)$
$w(1, -1)$	$w(1, 0)$	$w(1, 1)$

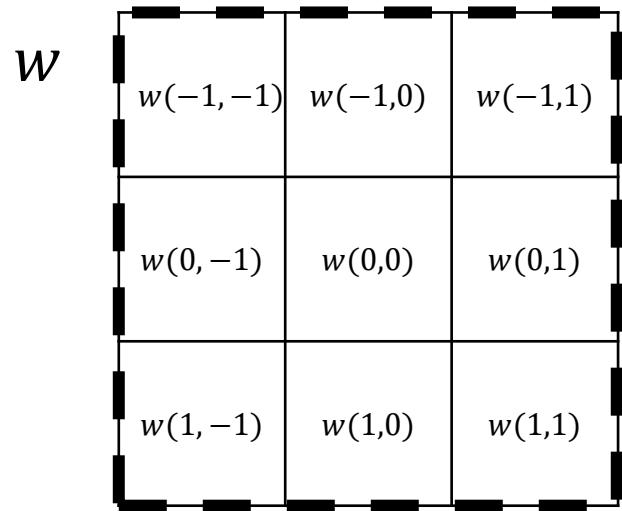
This operation is repeated for each pixel in the input image



# 2D Correlation

Convolution is a linear transformation. Linearity implies that

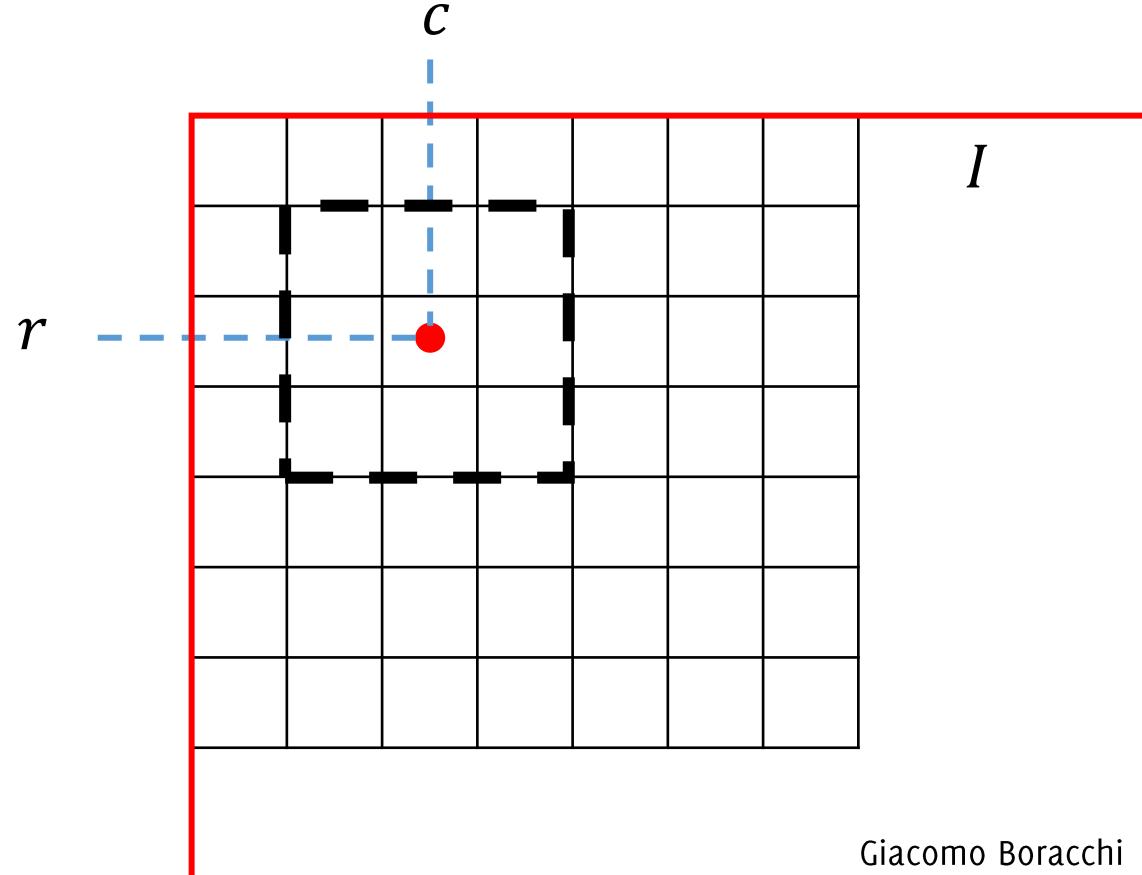
$$(I \otimes w)(r, c) = \sum_{(u,v) \in U} w(u, v) I(r + u, c + v)$$



We can consider weights as a filter  $h$

The filter  $h$  entirely defines convolution

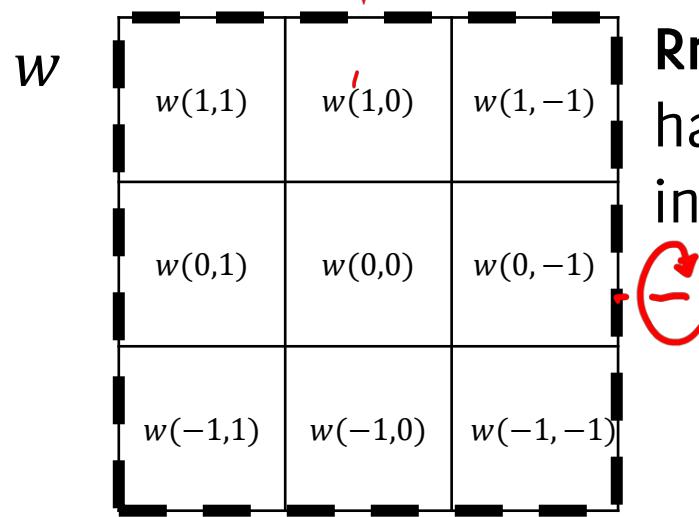
Convolution operates the same in each pixel



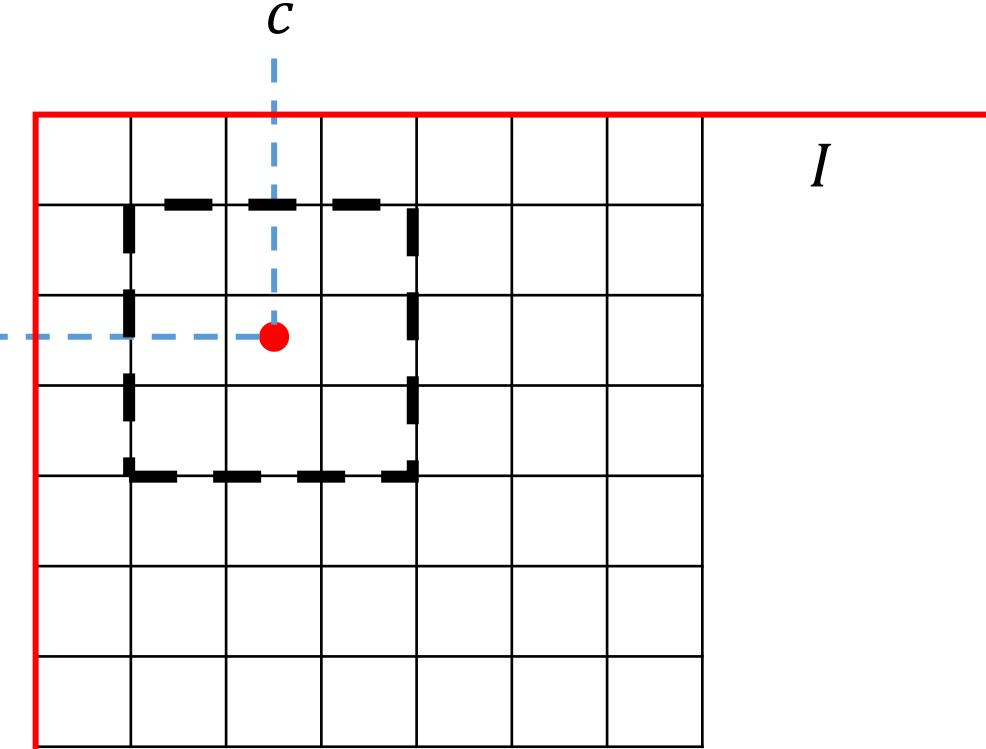
# 2D Convolution

Convolution is a linear transformation. Linearity implies that

$$(I \circledast w)(r, c) = \sum_{(u,v) \in U} w(u, v) I(r - u, c - v)$$



Rmk: indexes  
have been shifted  
in the filter  $w$



We can consider weights as a filter  $h$

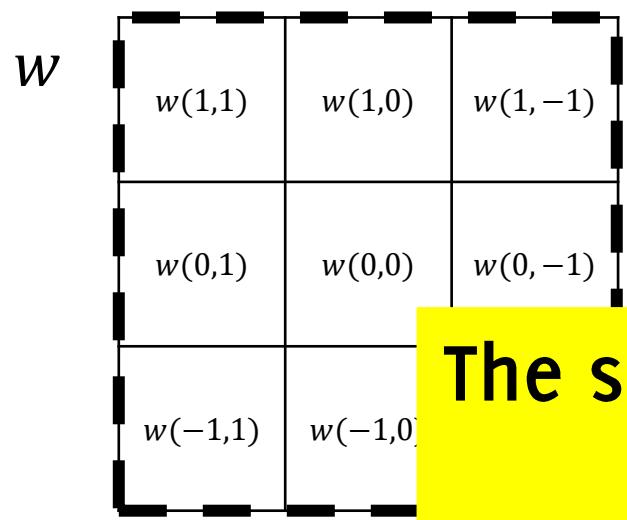
The filter  $h$  entirely defines convolution

Convolution operates the same in each pixel

# 2D Convolution

Convolution is a linear transformation. Linearity implies that

$$(I \circledast w)(r, c) = \sum_{(u,v) \in U} w(u, v) I(r - u, c - v)$$

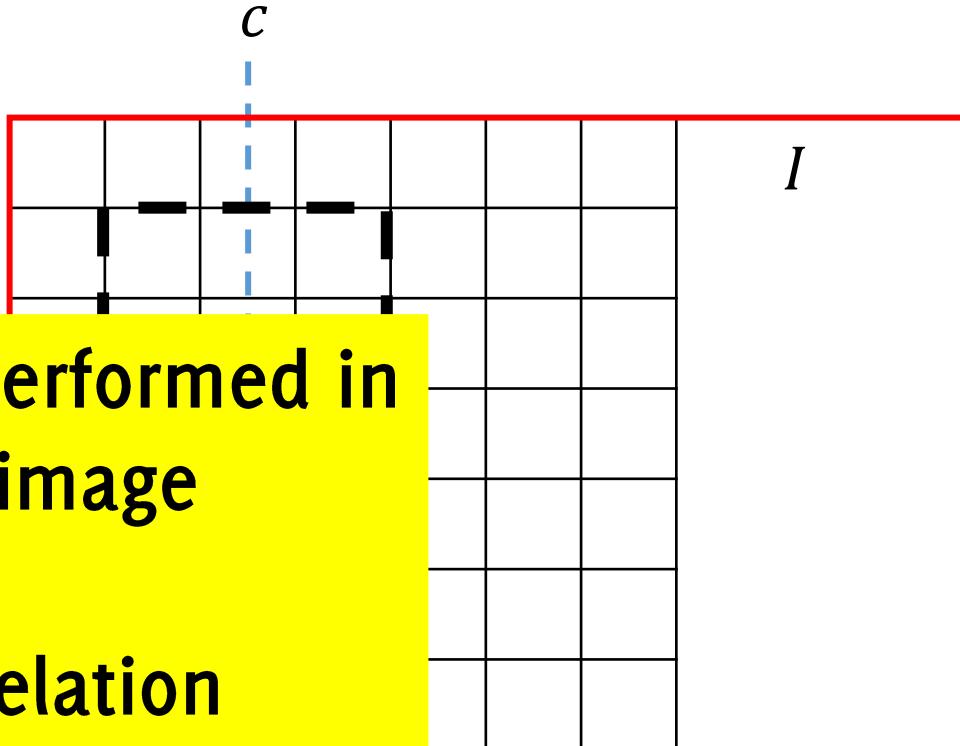


Rmk: indexes  
have been shifted  
in the filter  $w$

We can con  
The filter  $h$  e  
Convolution op

The same operation is being performed in each pixel of the input image

It is equivalent to 2D Correlation  
up to a «flip» in the filter  $w$



# 2D Convolution

Convolution is a linear transformation. Linearity implies that

$$(I \circledast w)(r, c) = \sum_{(u,v) \in U} w(u, v) * I(r - u, c - v)$$

Convolution is defined up to the “filter flip” for the Fourier Theorem to apply. Filter flip must be considered when computing convolution in Fourier domain and when designing filters.

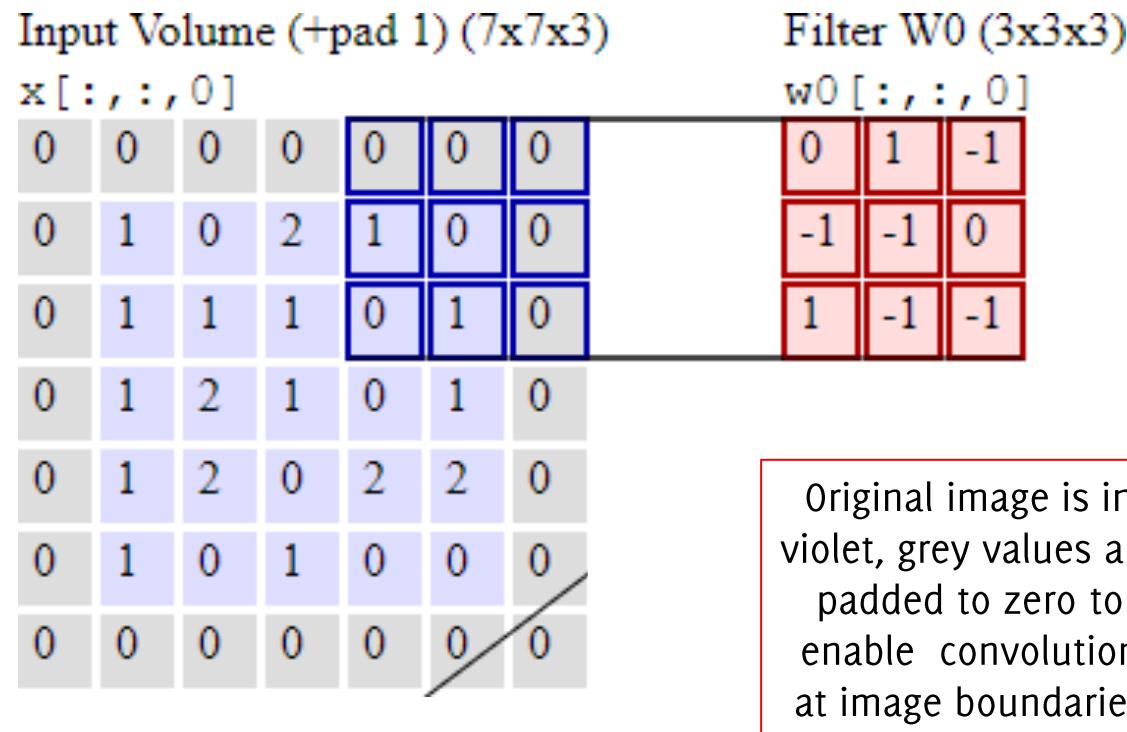
However, in CNN, convolutional filters are being learned from data, thus it is only important to use these in a consistent way.

In practice, in CNN arithmetic there is no flip!

# Convolution: Padding

How to define convolution output close to image boundaries?

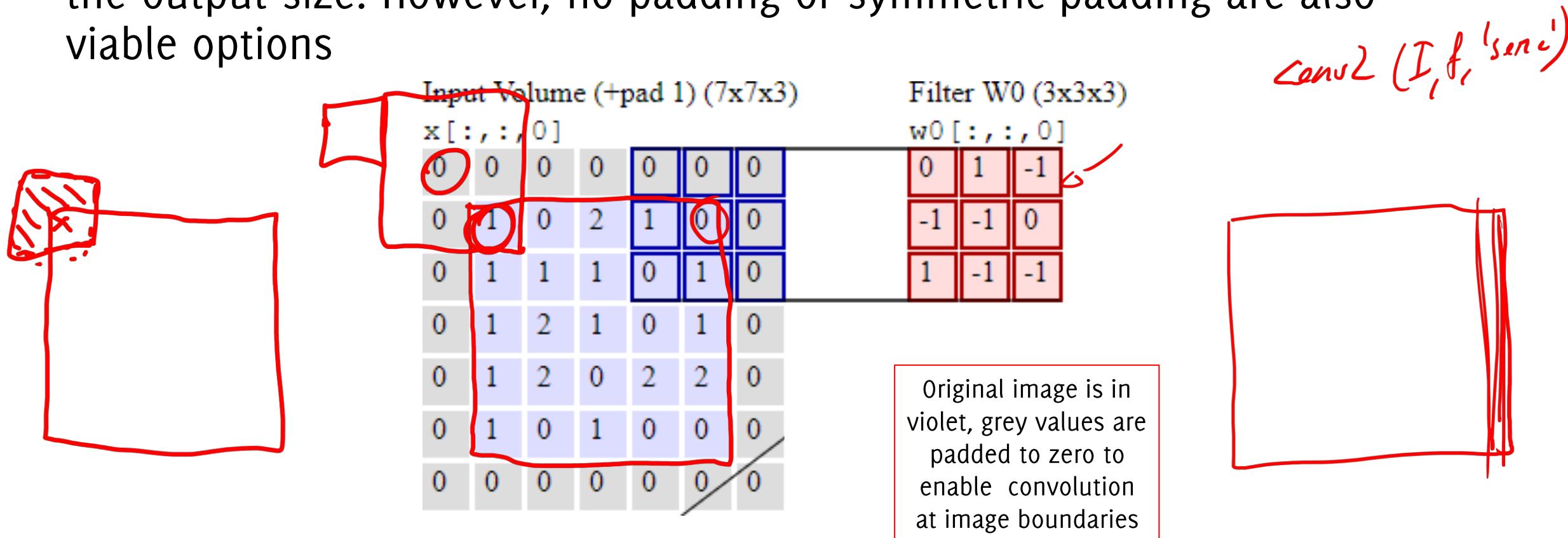
Padding with zero is the most frequent option, as this does not change the output size. However, no padding or symmetric padding are also viable options



# Convolution: Padding

How to define convolution output close to image boundaries?

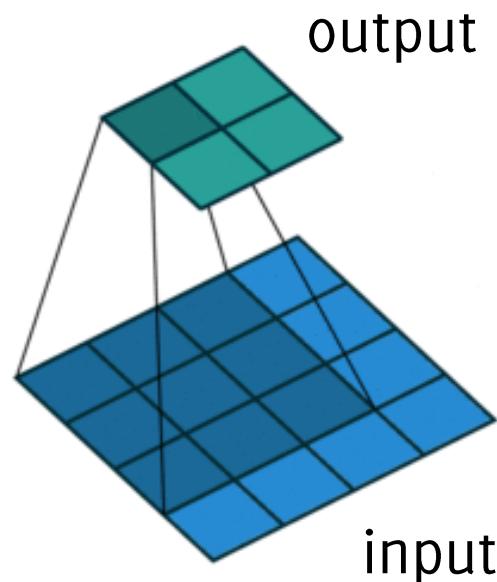
Padding with zero is the most frequent option, as this does not change the output size. However, no padding or symmetric padding are also viable options



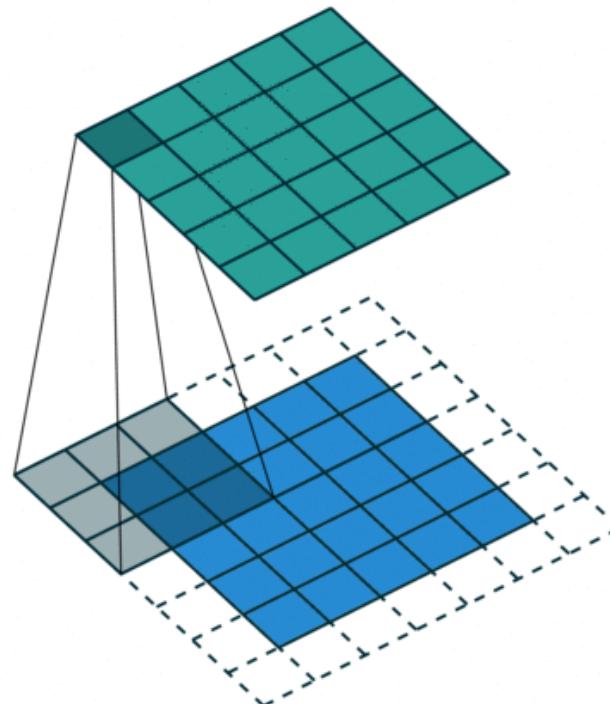
# Padding Options in Convolution Animation

Rmk: Blue maps are inputs, and cyan maps the outputs.

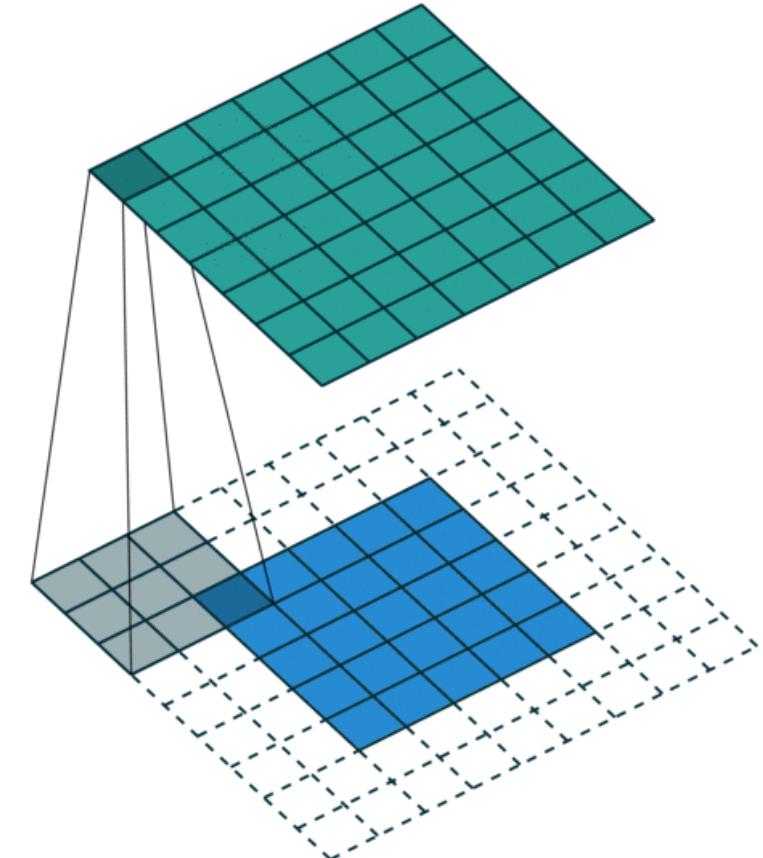
Rmk: the filter here is  $3 \times 3$



No padding  
«valid»



Half padding  
«same»

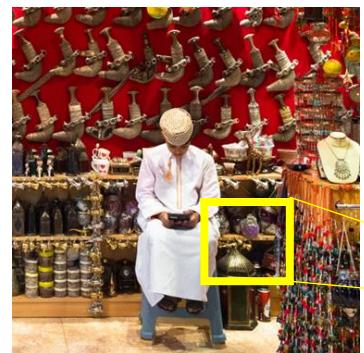


full padding  
«full»

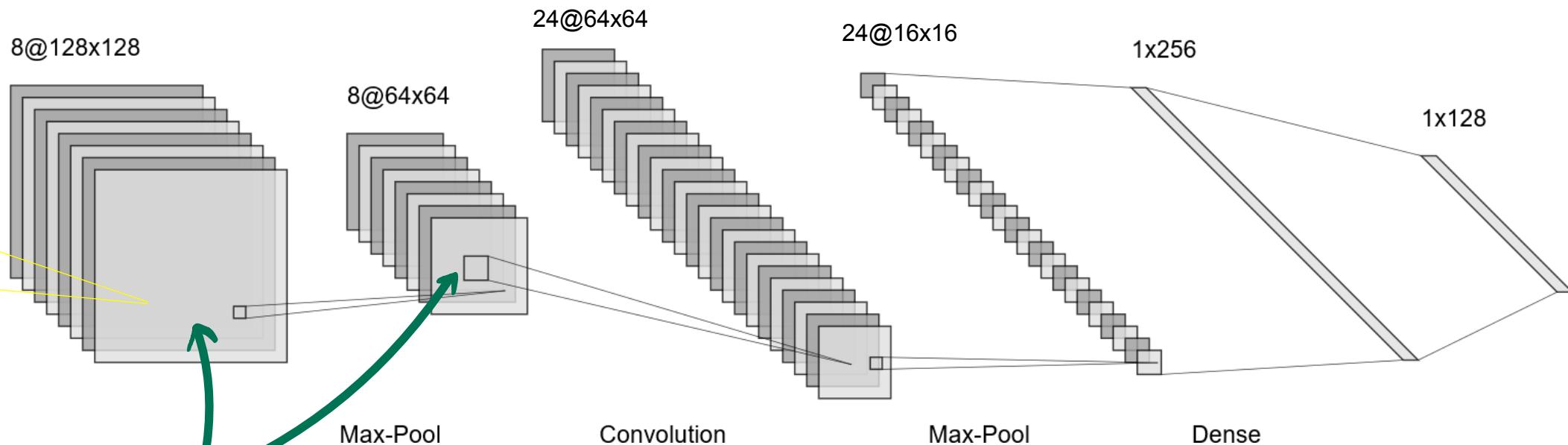
# Convolutional Neural Networks

CNNs

# The typical architecture of a CNN

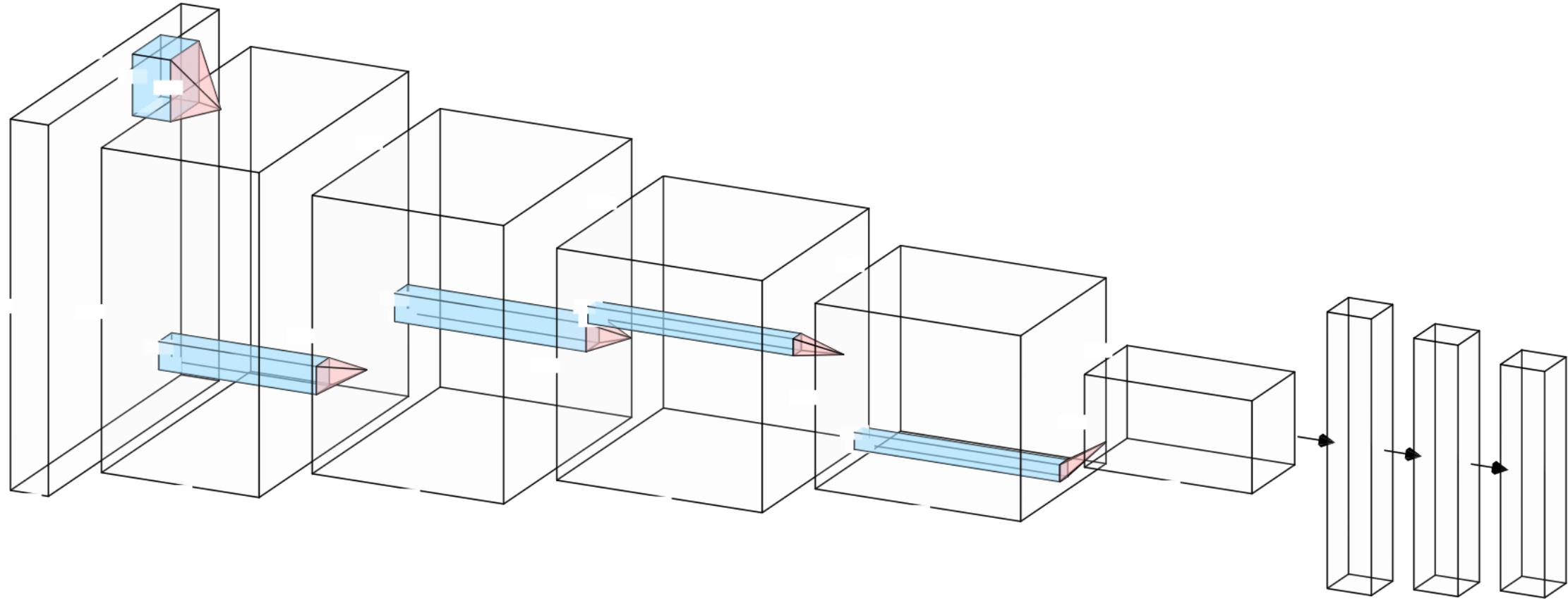


RGB

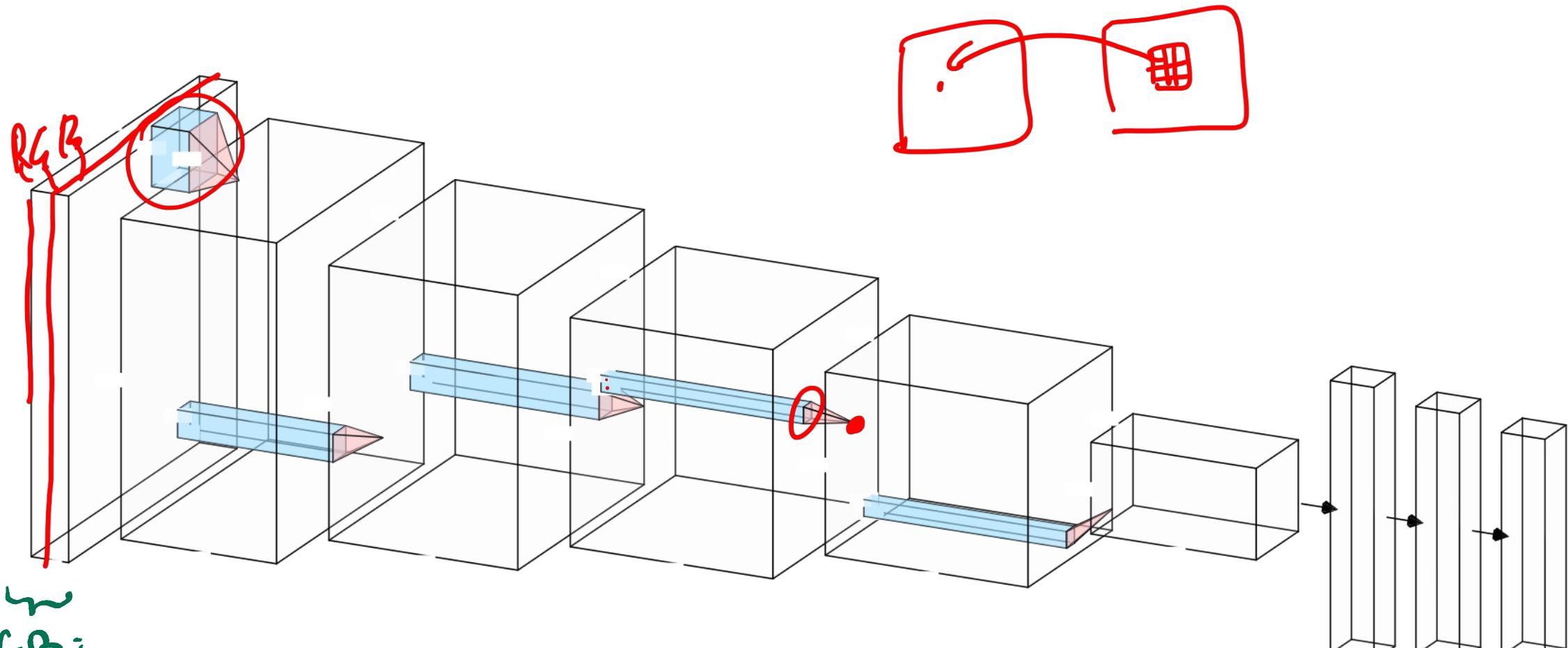


the "spatial size" (2D) ↓  
but the nb of layers ↑ (3<sup>rd</sup> dim).

# The typical architecture of a CNN



# The typical architecture of a CNN

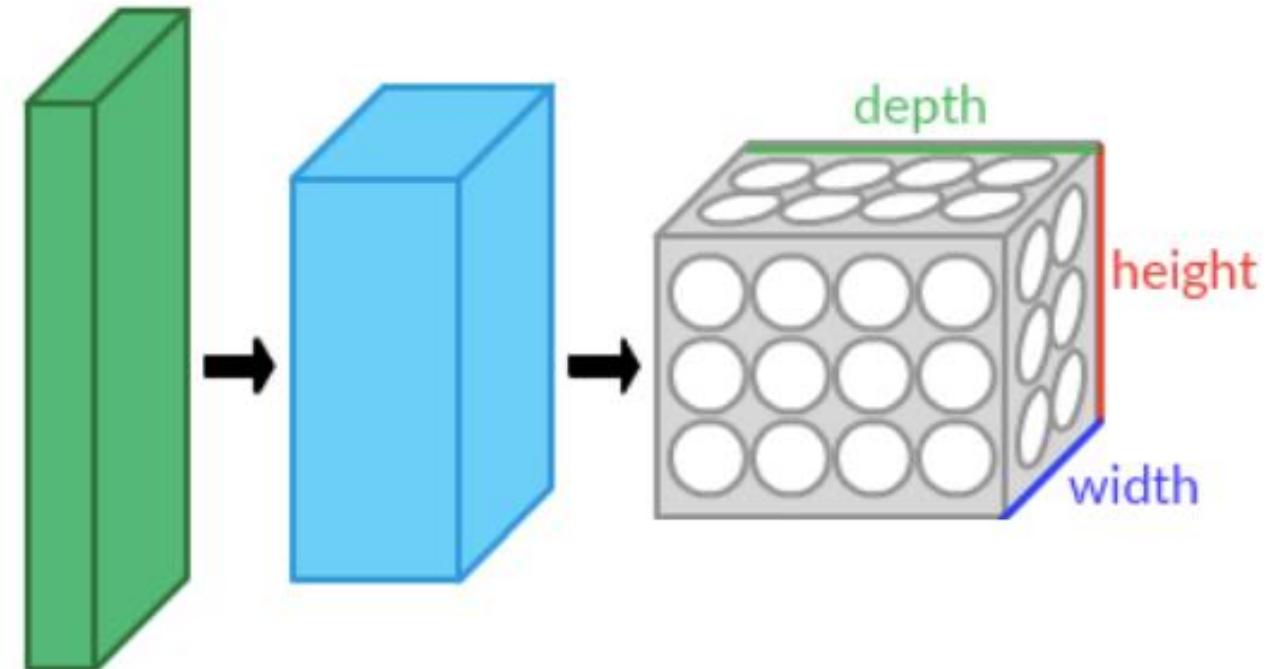


RGB:  
3 layers

# Convolutional Neural Networks (CNN)

CNN are typically made of blocks that include:

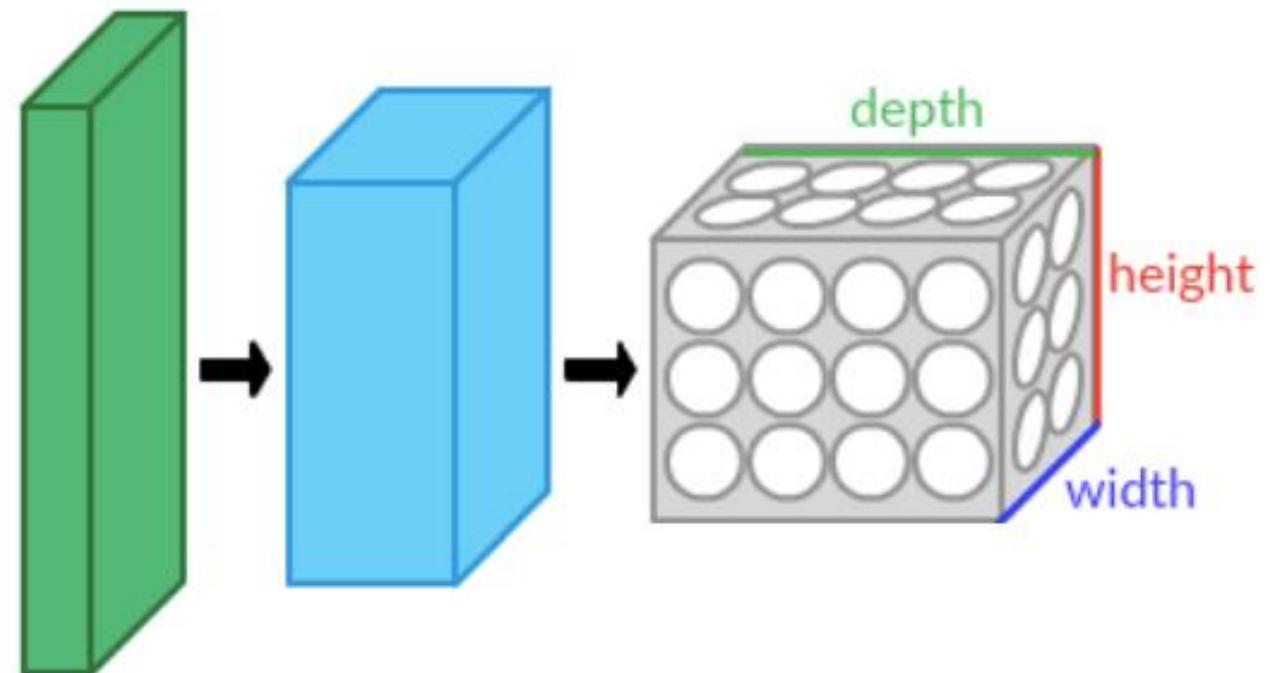
- Convolutional layers
- Nonlinearities (activation functions)
- Pooling Layers (Subsampling / maxpooling)



By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45661858>

# Convolutional Neural Networks (CNN)

- An image passing through a CNN is transformed in a sequence of volumes.
- As the depth increases, the height and width of the volume decreases
- Each layer takes as input and returns a volume



# Convolutional Layers

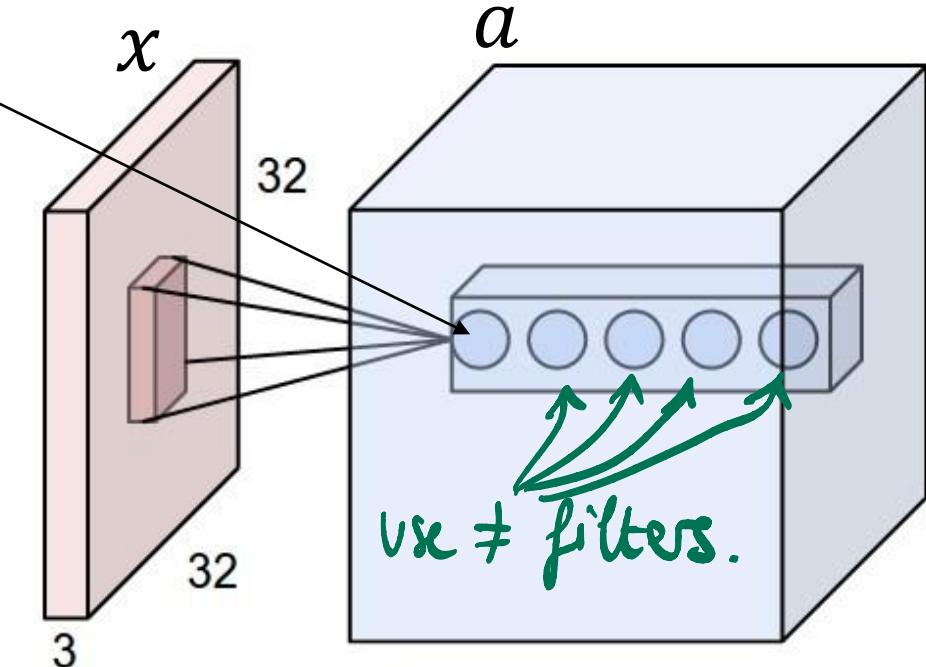


# Convolutional Layers

Convolutional layers "mix" all the input components

The output is a linear combination of all the values in a region of the input, considering all the channels

$$a(r, c, 1) = \sum_{(u,v) \in U, k} w^1(u, v, k) x(r + u, c + v, k) + b^1$$



Filters need to have the same number of channels as the input, to process all the values from the input layer



# Convolutional Layers

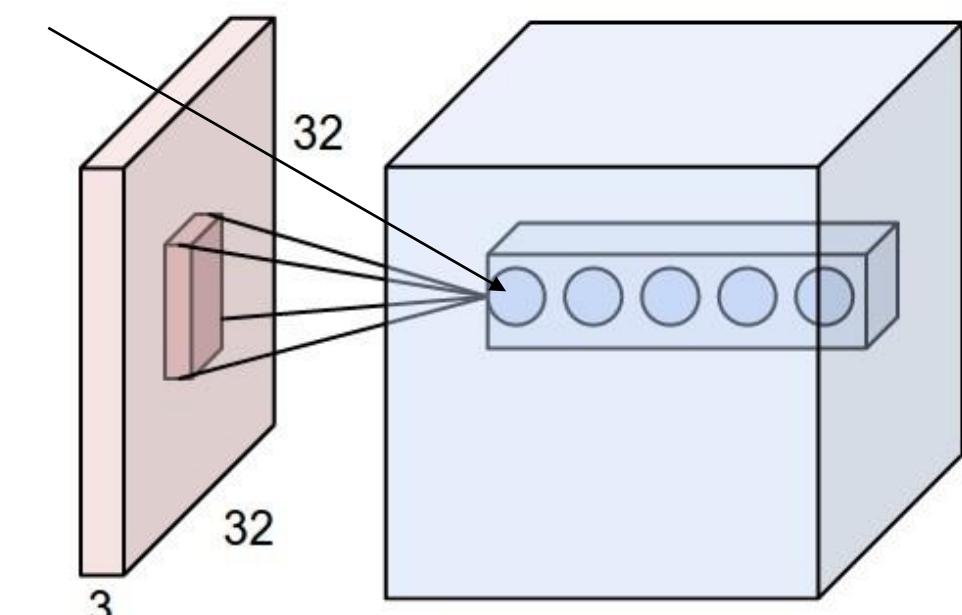
Convolutional layers "mix" all the input components

The output is a linear combination of all the values in a region of the input, considering all the channels *learned parameters*.

$$a(r, c, 1) = \sum_{(u,v) \in U, k} w^1(u, v, k) x(r + u, c + v, k) + b^1$$

The parameters of this layer are called filters.

The same filter is used through the whole spatial extent of the input





# Convolutional Layers

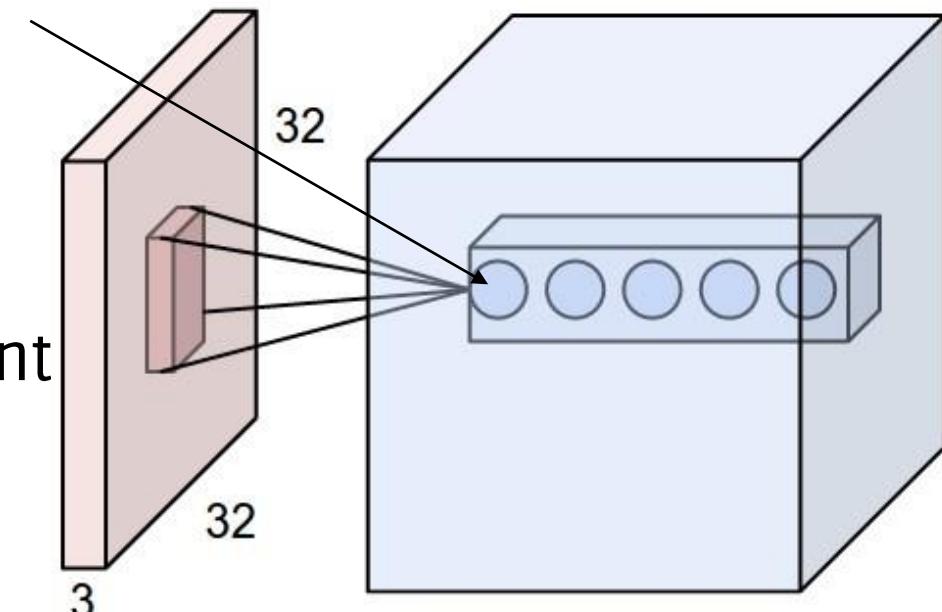
Convolutional layers "mix" all the input components

The output is a linear combination of all the values in a region of the input, considering all the channels

$$a(r, c, 1) = \sum_{(u,v) \in U, k} w^1(u, v, k) x(r + u, c + v, k) + b^1$$

The **spatial dimension**:

- spans a small neighborhood  $U$  (local processing, it's a convolution)
- $U$  needs to be specified, it is a very important attribute of the filter





# Convolutional Layers

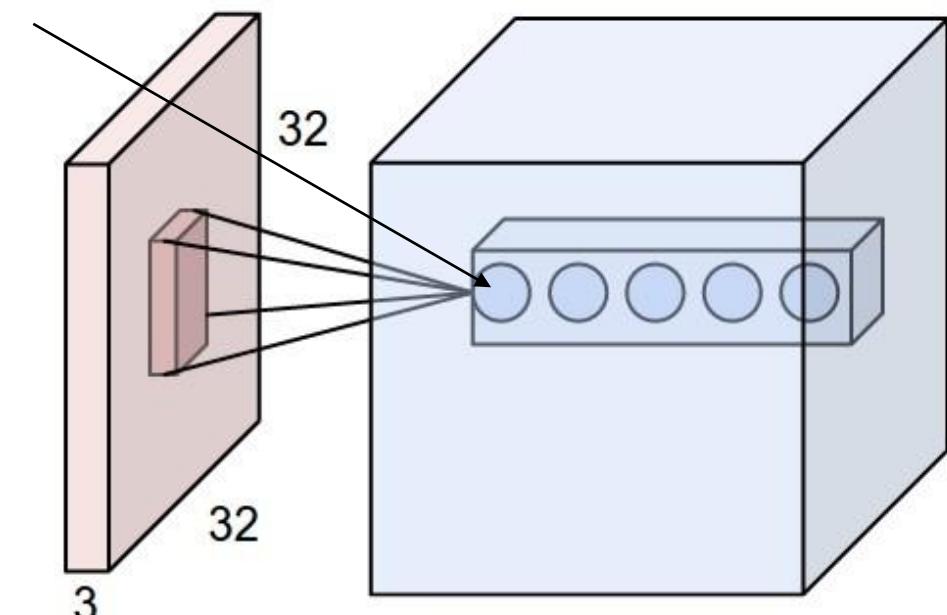
Convolutional layers "mix" all the input components

The output is a linear combination of all the values in a region of the input, considering all the channels

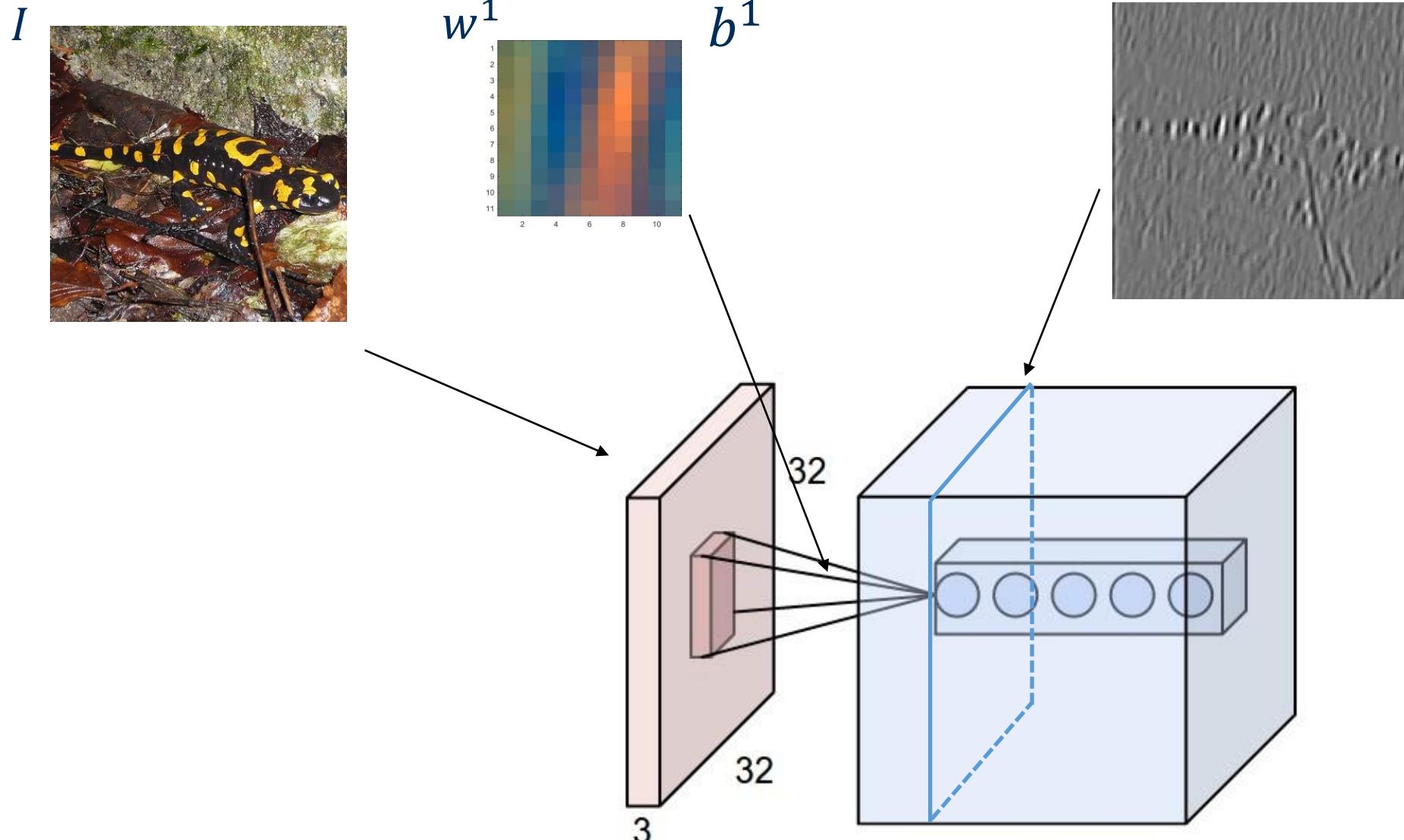
$$a(r, c, 1) = \sum_{(u,v) \in U, k} w^1(u, v, k) x(r + u, c + v, k) + b^1$$

The channel dimension:

- spans the entire input depth (no local processing, like spatial dimension)
- there is no need to specify that in the filter attributes



# Convolutional Layers



By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45659236>

# Convolutional Layers



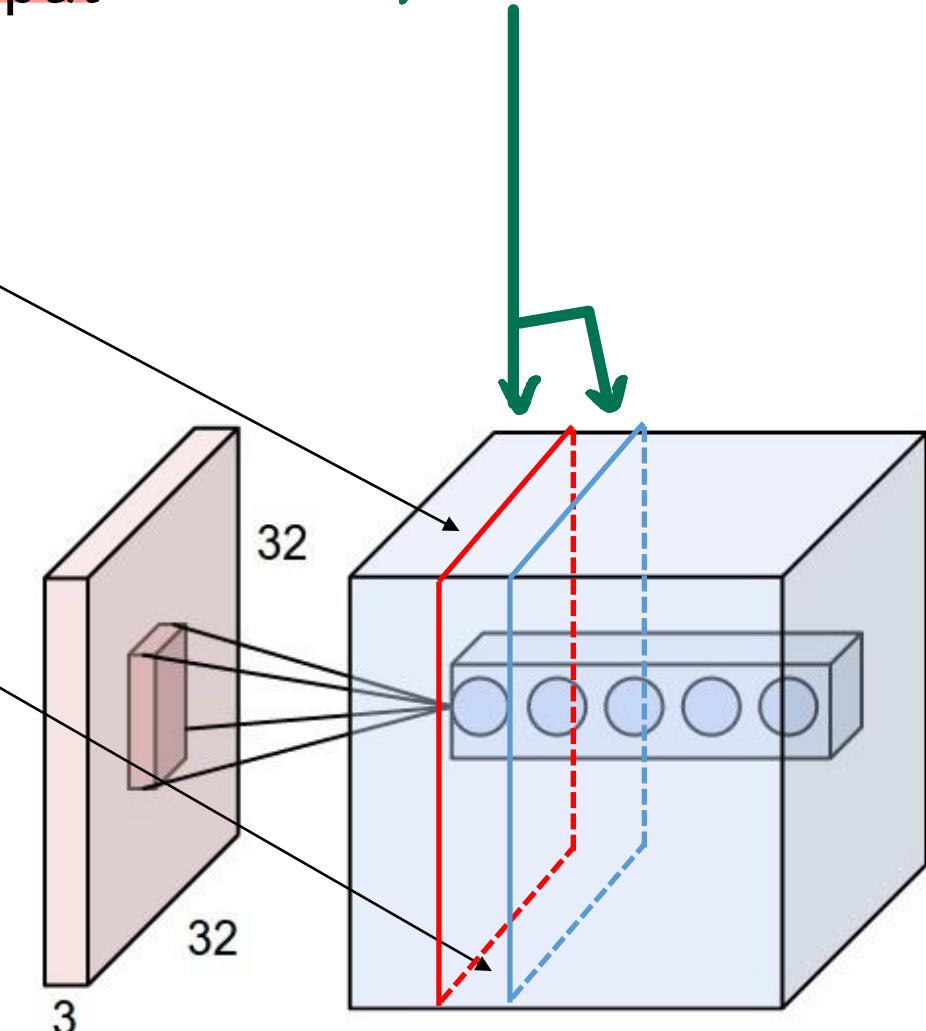
Different filters yield different layers in the output

$$a(r, c, 1) = \sum_{(u,v) \in U,k} w^1(u, v, k) x(r + u, c + v, k) + b^1$$

$$a(r, c, 2) = \sum_{(u,v) \in U,k} w^2(u, v, k) x(r + u, c + v, k) + b^2$$

Different filters of the same layer have the same spatial extent

use + filters to increase the depth!



# Convolutional Layers

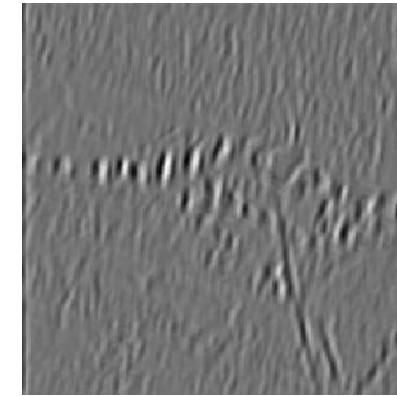
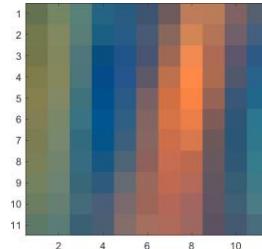
$a(:,:,1)$

$a(:,:,n)$

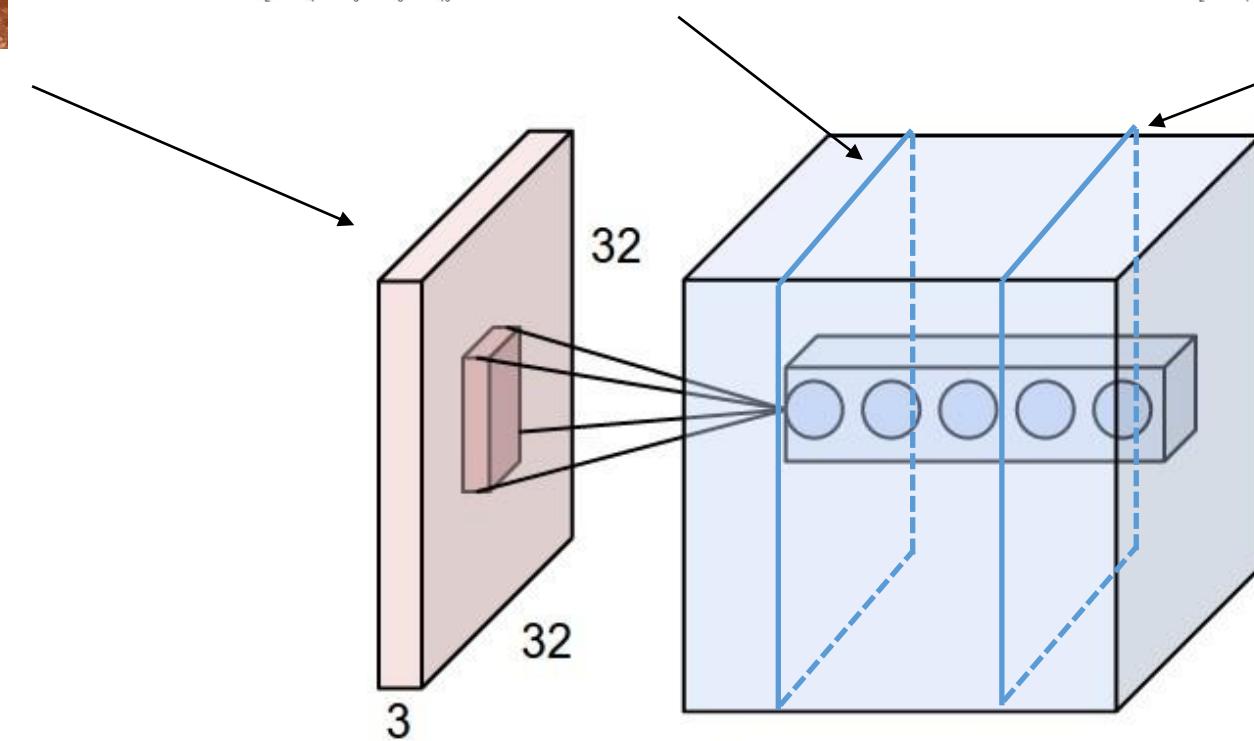
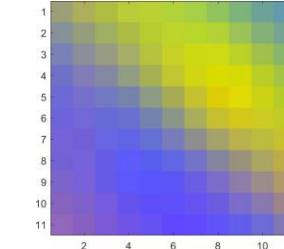
$I$



$w^1, b^1$



$w^n, b^n$



By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45659236>

Giacomo Boracchi



# Convolutional Layers, general formula:

$$a(r, c, l) = \sum_{\substack{(u,v) \in U, \\ k=1, \dots, C}} w^l(u, v, k) x(r + u, c + v, k) + b^l, \forall (r, c), l = 1, \dots, N_F$$

Where  $r, c$  denote a spatial location in the output activation  $a$ , while  $l$  is the channel.



The filter  $w$  is identified by

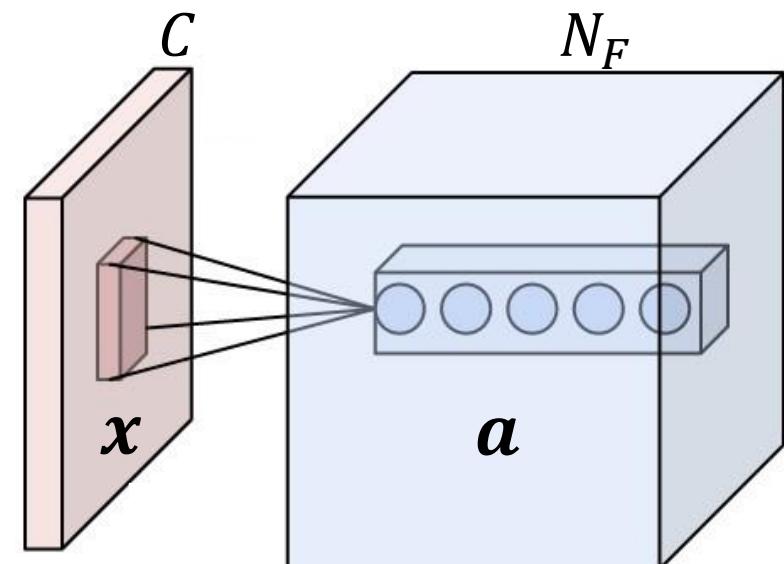
- a spatial neighborhood  $U$  having size  $h_r \cdot h_c$
- the same number of channels  $C$  as the input activations

The parameters are the weights + one bias per filter

The overall number of parameters is

$$(h_r \cdot h_c \cdot C) \cdot N_F + N_F$$

$$\forall l : h_r \times h_c \times C + 1 \quad \text{bias} .$$



Layers with the same hyperparameters can have different number of parameters depending on where these are located in the network

# Convolutional Layers, remarks:

Given:

```
conv2 = tfkl.Conv2D(  
    filters = n_f,  
    kernel_size = (h_r,h_c),  
    activation = 'relu',  
    strides = (1,1), → see 3 slides later.  
    padding = 'same',  
    name = 'conv2'  
)
```

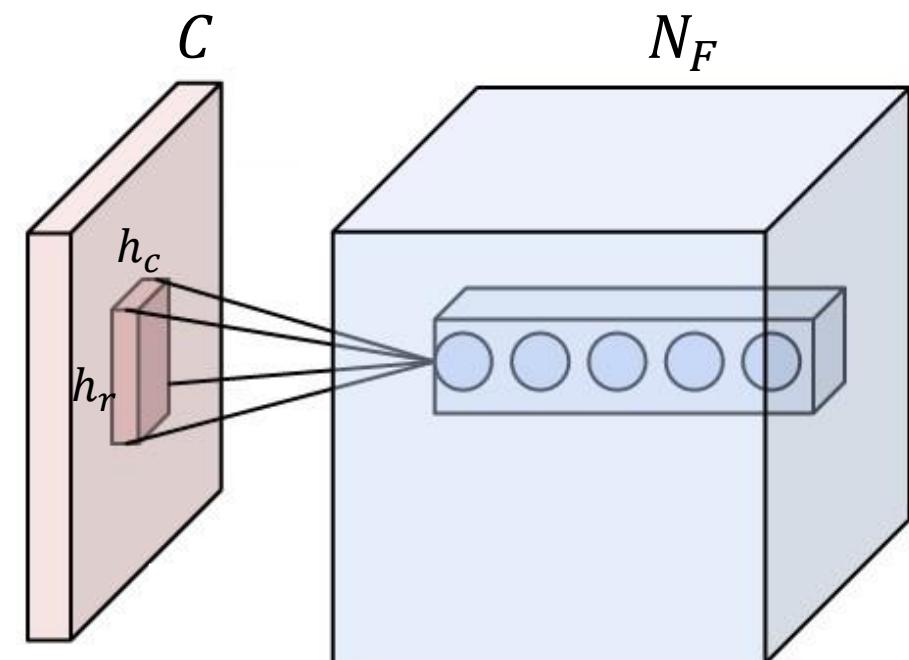
The parameters are the weights + one bias per filter

The overall number of parameters is

$$(h_r \cdot h_c \cdot C) \cdot N_F + N_F$$

Where  $d$  is the **depth of the input activation**

Layers with the same attribute can have different number of parameters depending on where these are located



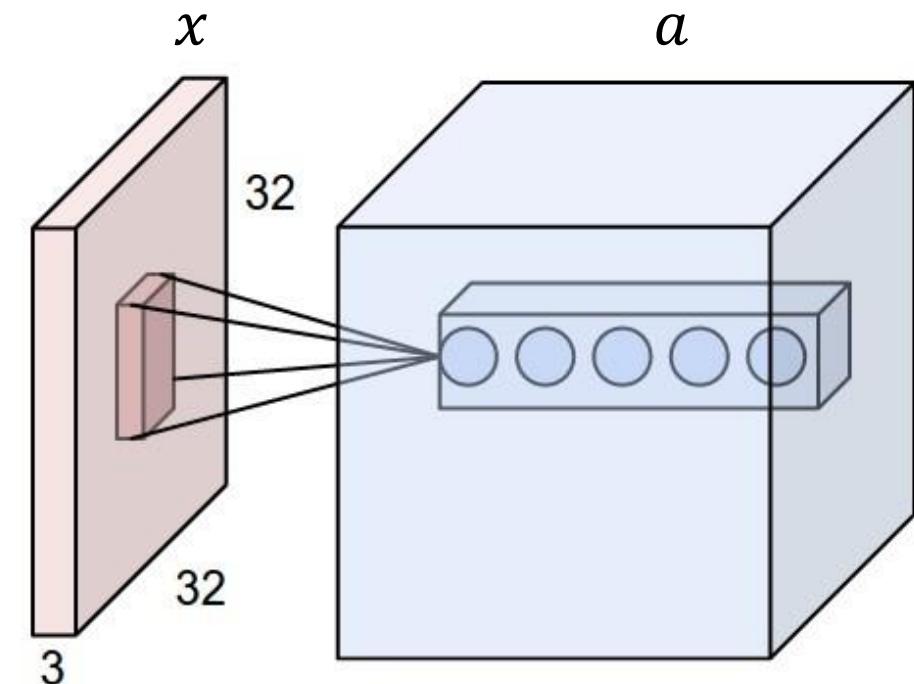
By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45659236>



# Recap: Convolutional Layers

Convolutional layers "mix" all the input components

- The output is also called **volume** or **activation maps**
- Each filter yields a different slice of the output volume
- Each filter has depth equal to the depth of the input volume



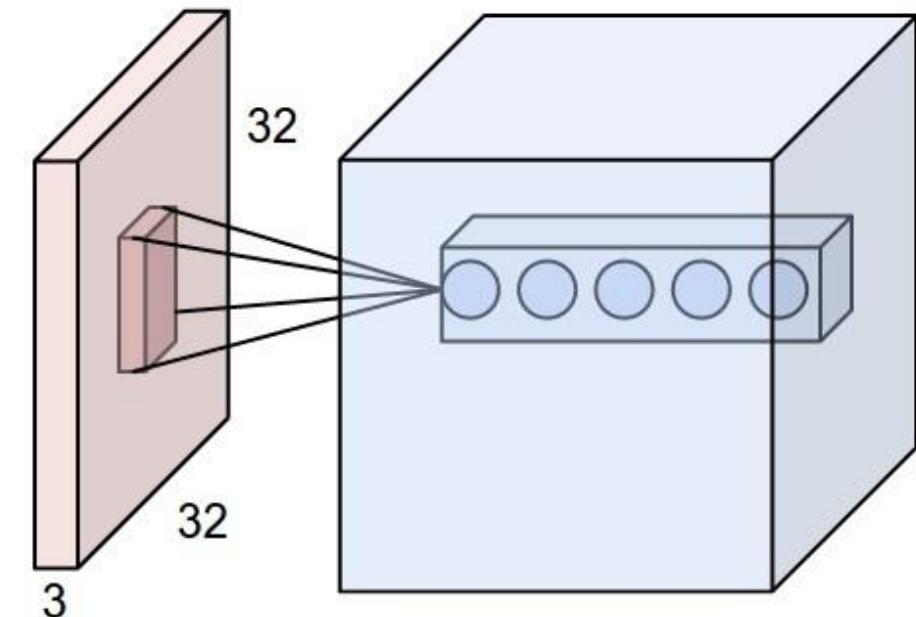
By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45659236>

# Convolutional Layers, remarks:



- Convolutional Layers are described by a set of filters
- Filters represent the weights of these linear combination.
- Filters have very small spatial extent and large depth extent,
- The filter depth is typically not specified, as it corresponds to the number of layers of the input volume

The output of the convolution against a filter becomes a slice in the volume feed to the next layer



By Aphex34 - Own work, CC BY-SA 4.0,  
<https://commons.wikimedia.org/w/index.php?curid=45659236>



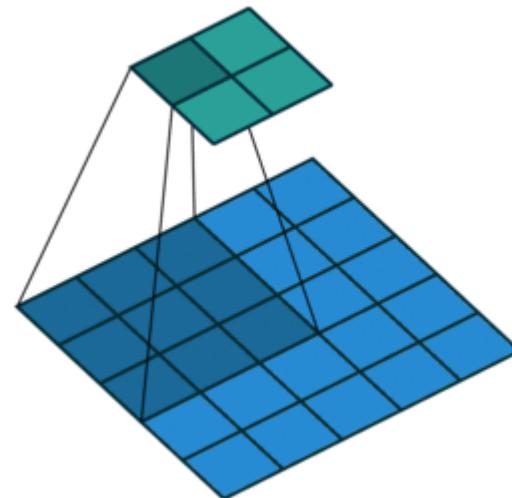
# Stride Options in Convolutions

This is not very popular in image processing, but **used in CNNs**.

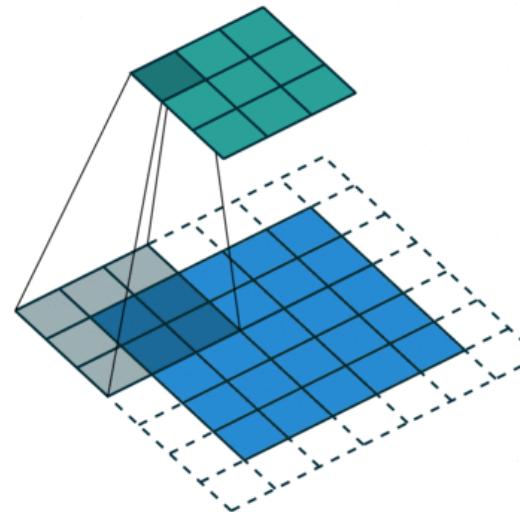
It is possible to **reduce the spatial extent** by «sliding» the filter over each channel of the input volume with a ***stride***. The larger the stride, the more the spatial dimensions get shrunk.

Stride = (1,1) corresponds to standard convolution, Stride = (2,2) means convolutional filter «jumps» every second row and every second column

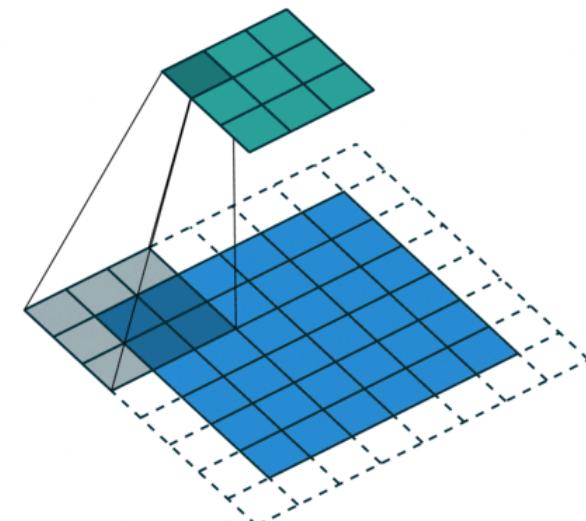
In principle this corresponds to convolution followed by downsampling (less efficient).



No Padding, stride 2

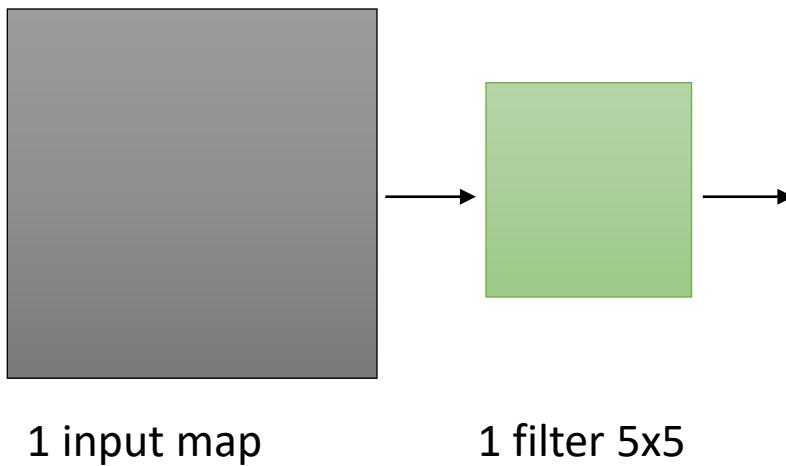


Padding, stride 2

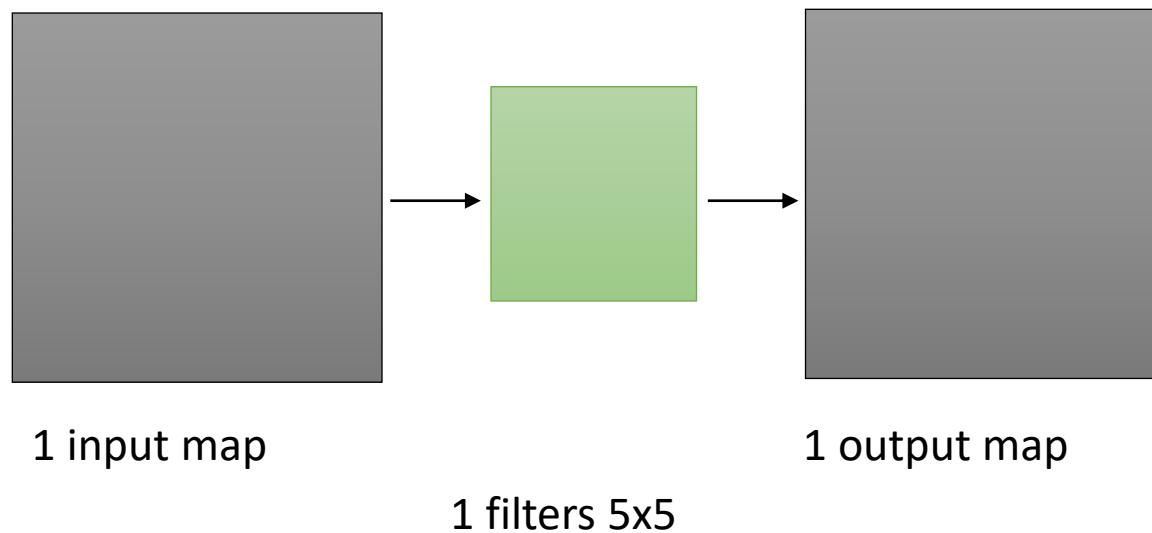


Padding, stride 2 (odd)

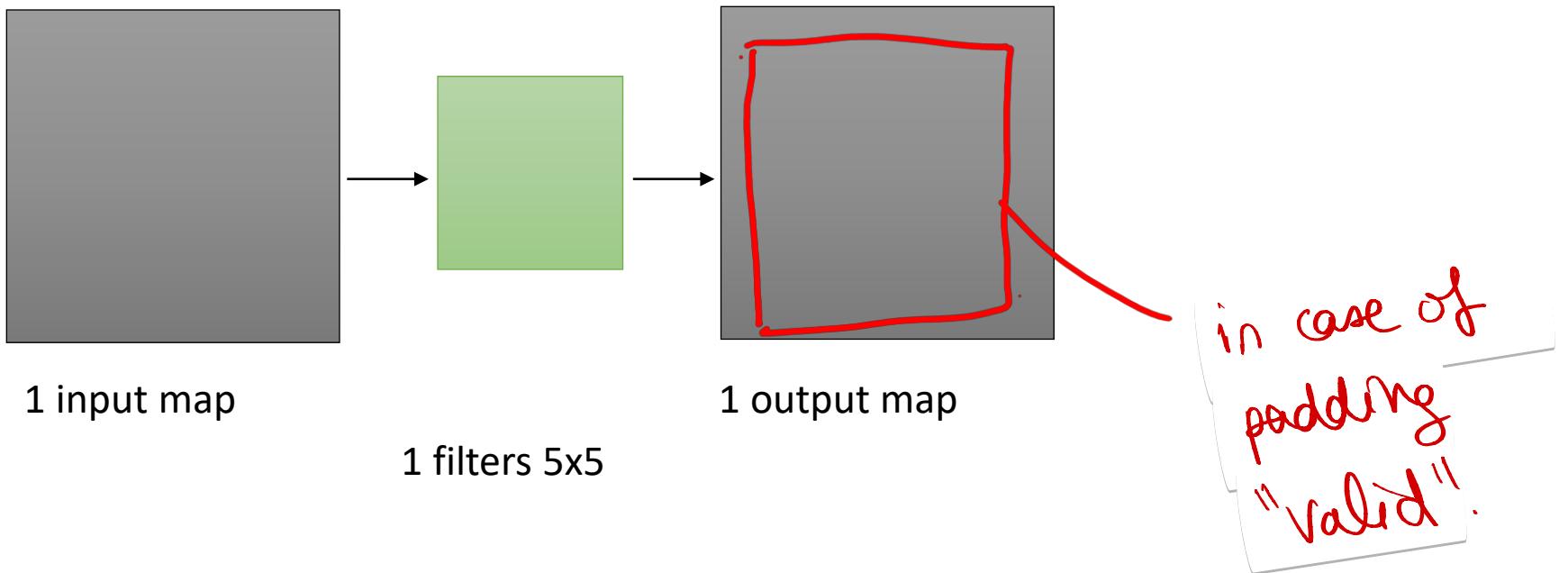
# CNN Arithmetic



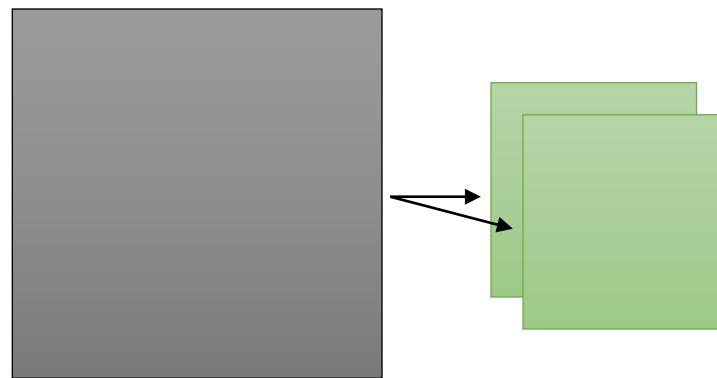
# CNN Arithmetic



# CNN Arithmetic



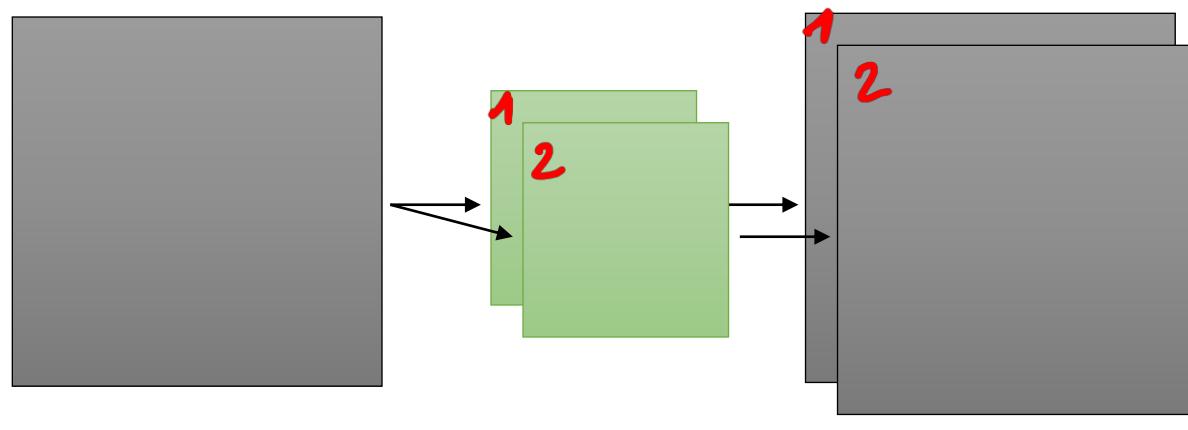
# CNN Arithmetic



1 input map

2 filters 5x5

# CNN Arithmetic

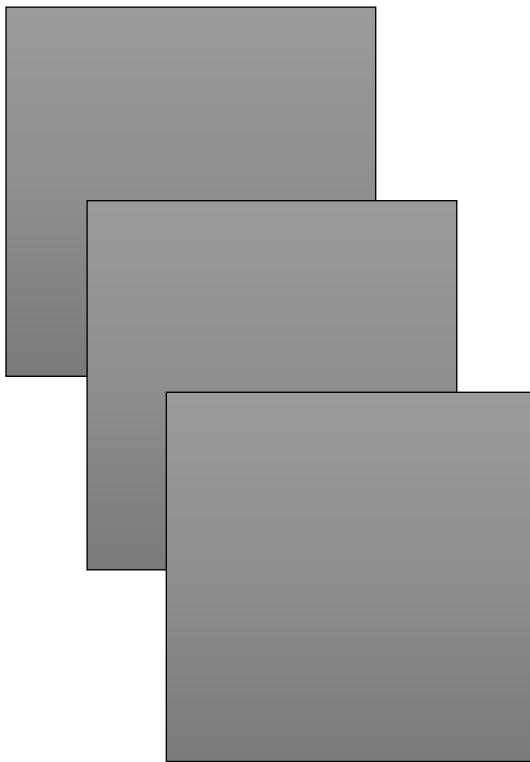


1 input map

2 filters 5x5

2 output maps

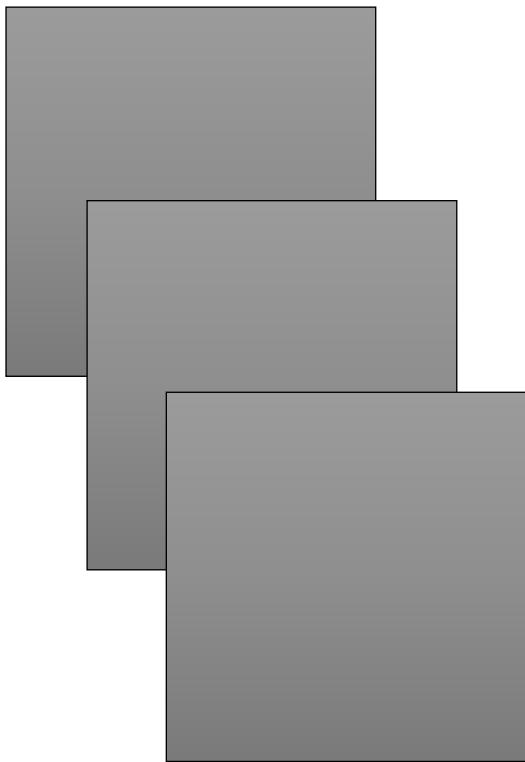
# CNN Arithmetic



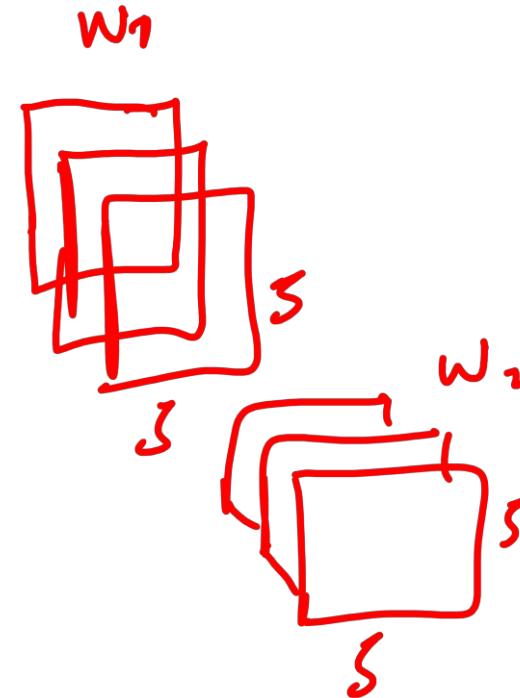
3 input maps

2 filters 5x5

# CNN Arithmetic



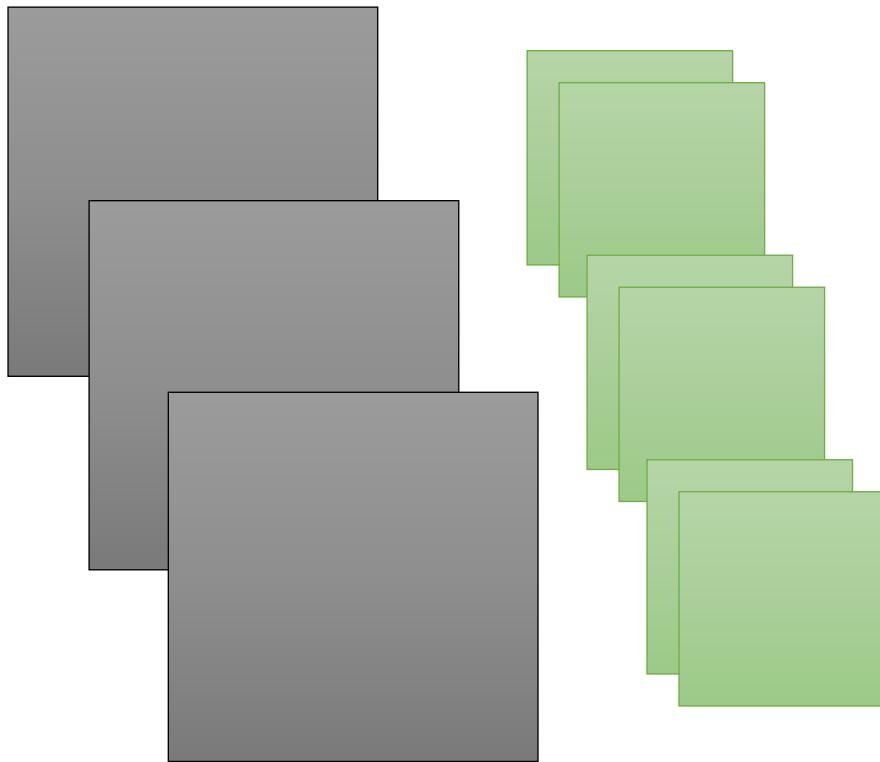
3 input maps



2 filters 5x5

↓  
3 channels per filter.

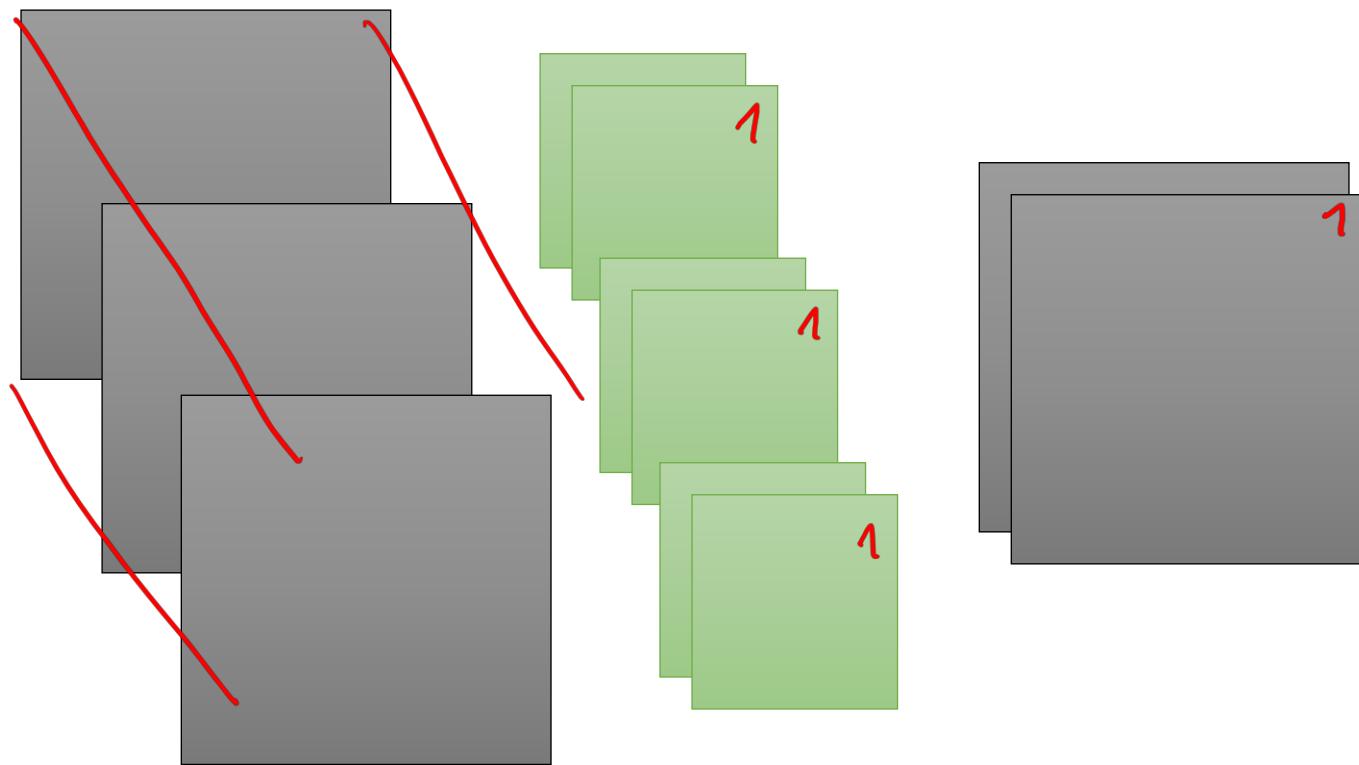
# CNN Arithmetic



3 input maps

2 filters 5x5

# CNN Arithmetic

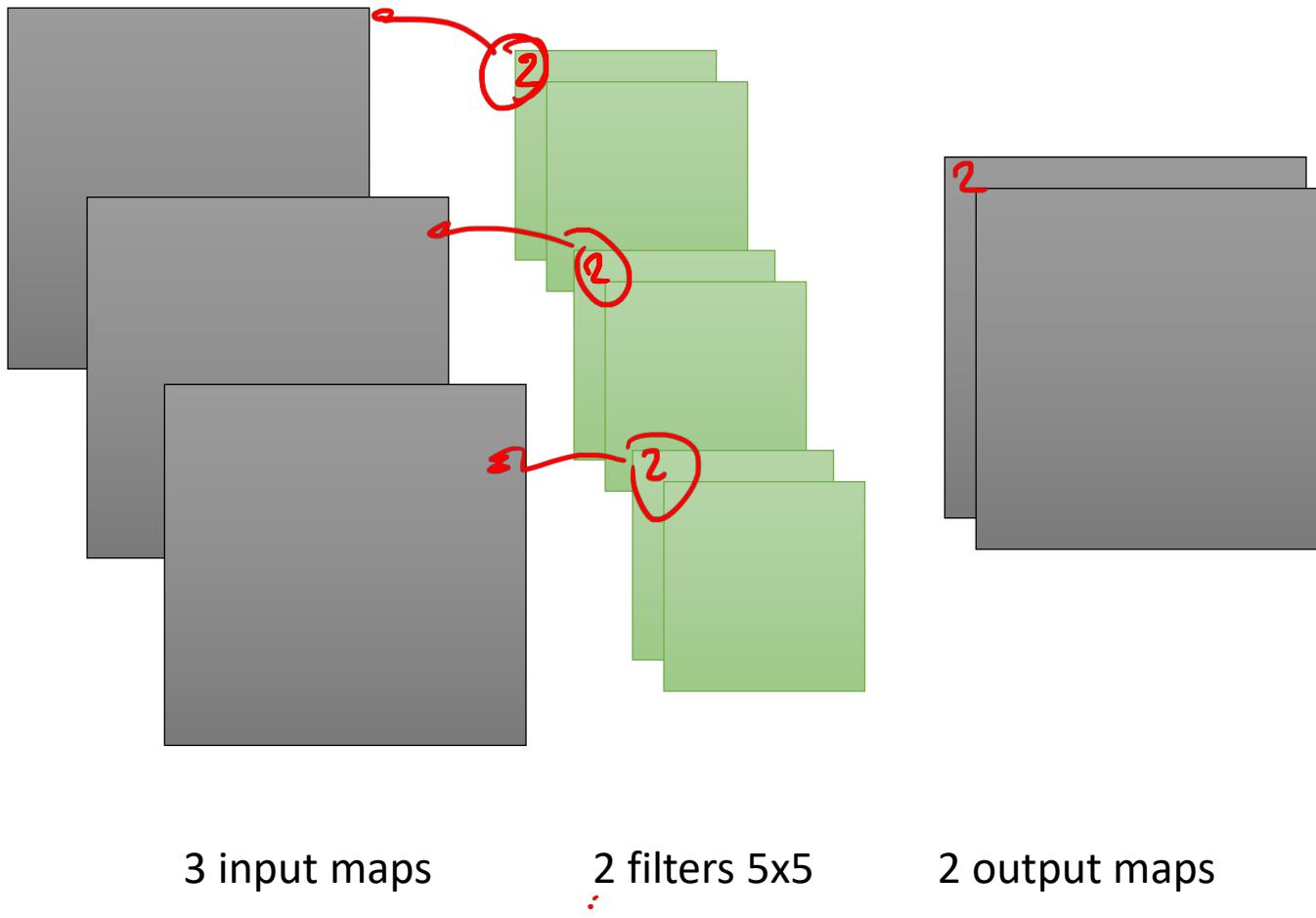


3 input maps

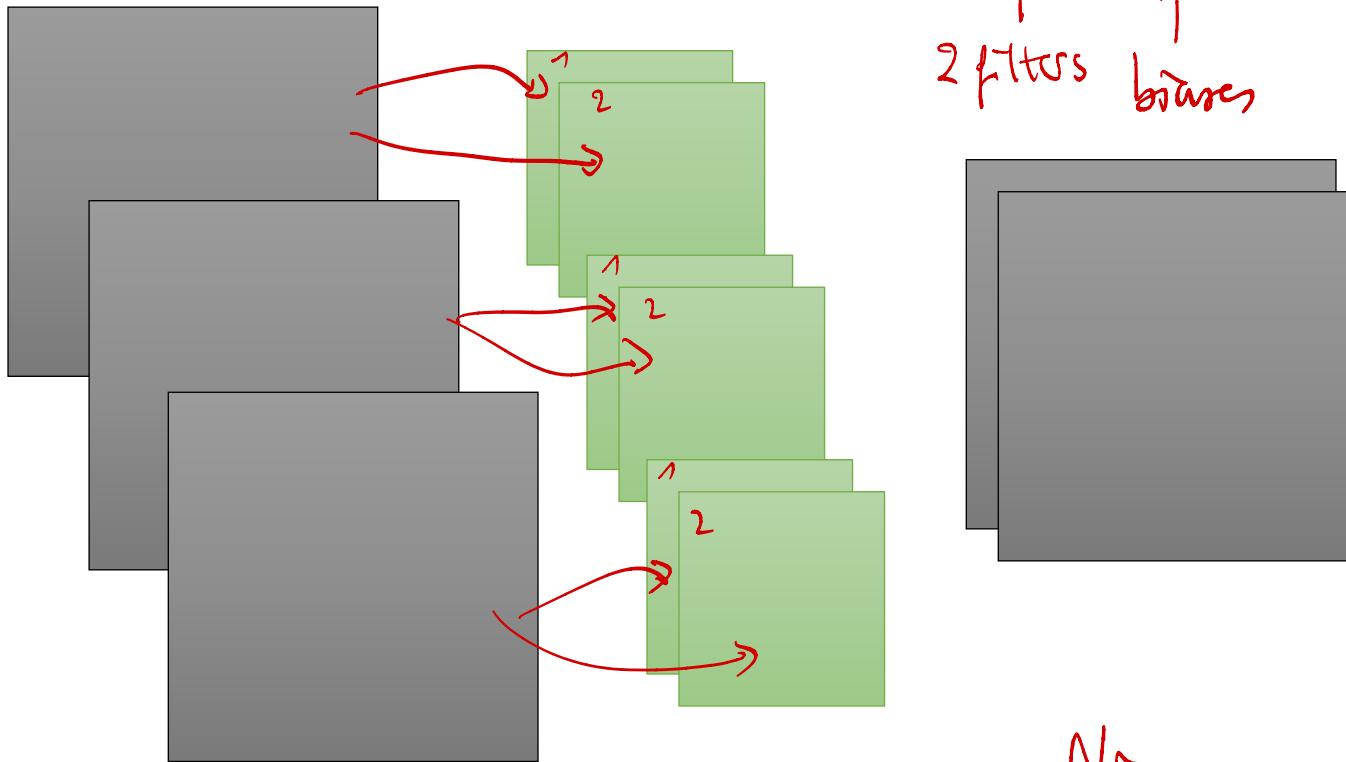
2 filters 5x5

2 output maps

# CNN Arithmetic



# CNN Arithmetic



3 channels  $\Rightarrow$  3 input images .

$$(5 \times 5 \times 3) \times 2 + 2 = 152$$

↑  
2 filters      ↑      biases

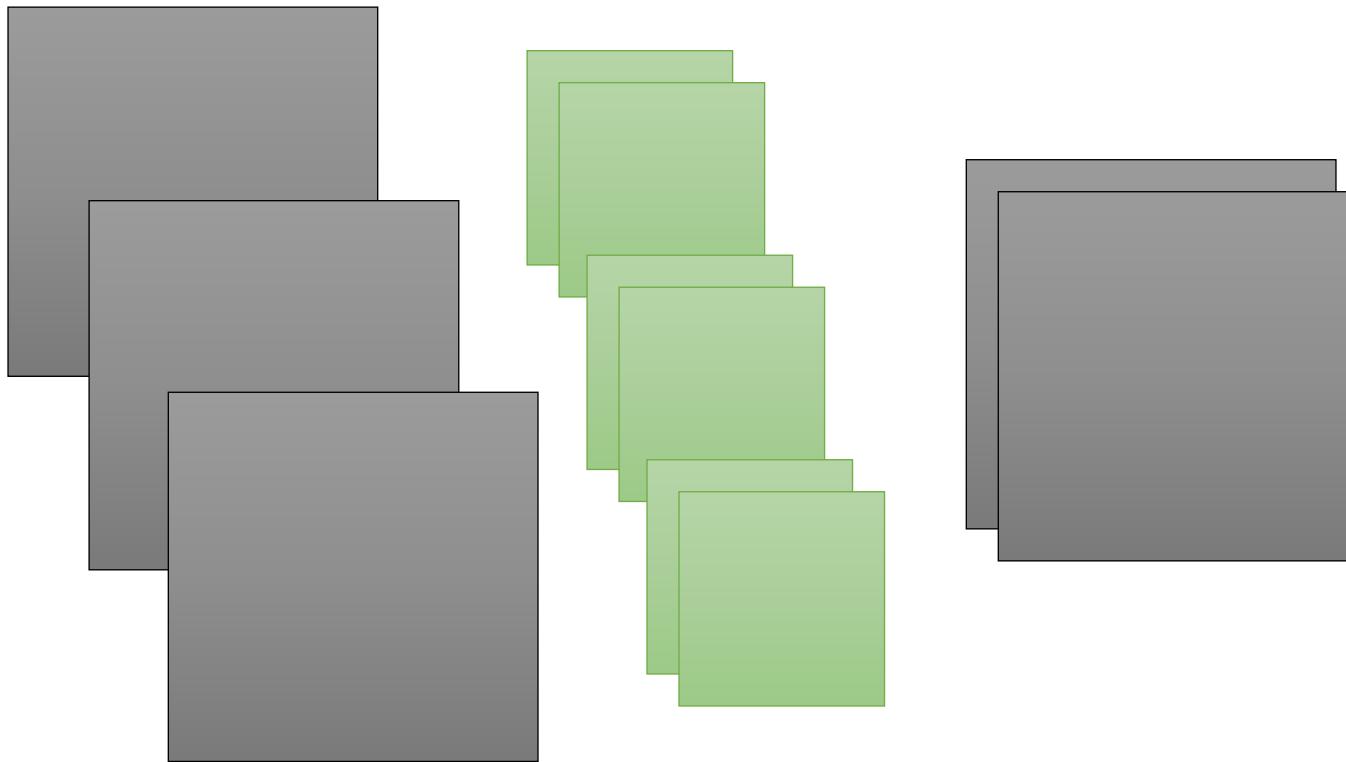
3 input maps

2 filters 5x5

2 output maps  
 $N_f$

Quiz: how many parameters  
does this layer have?

# CNN Arithmetic



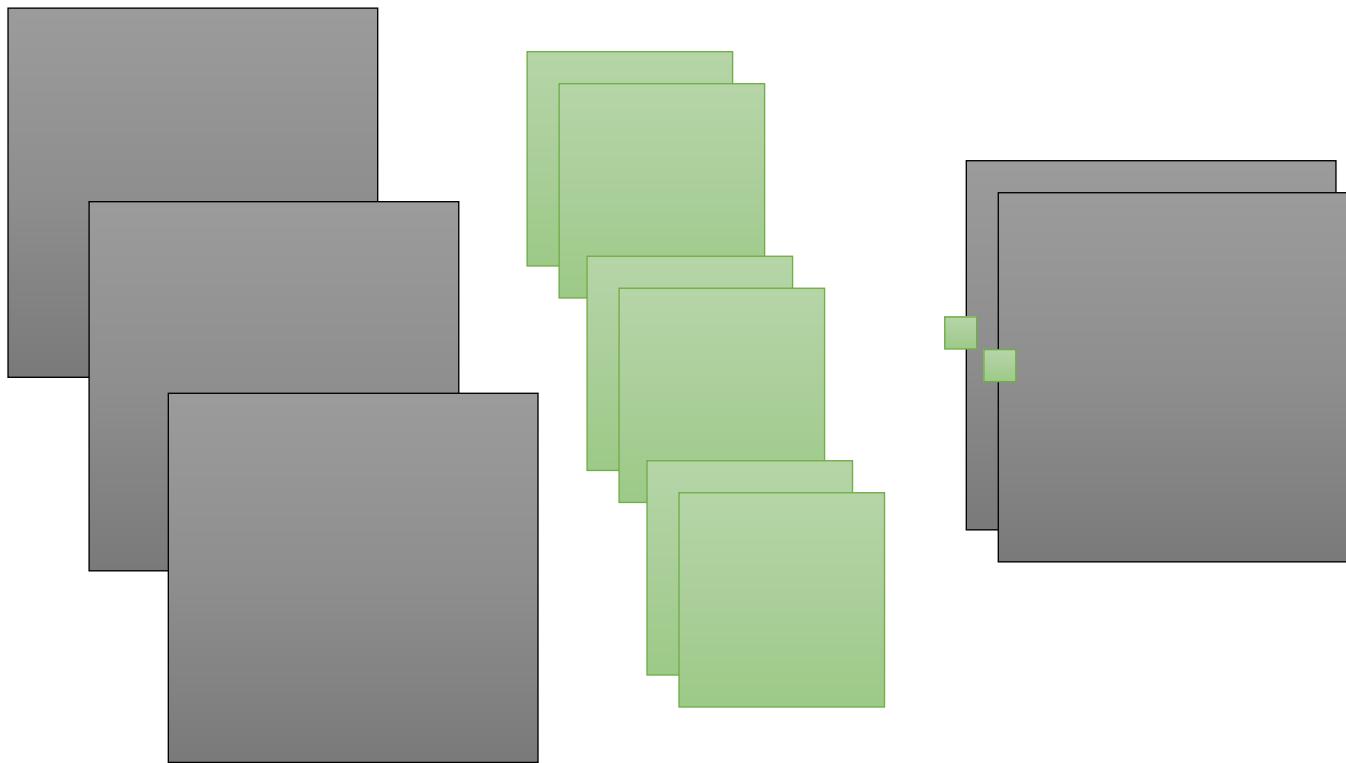
3 input maps

2 filters 5x5

= 150 parameters in the filters

2 output maps

# CNN Arithmetic

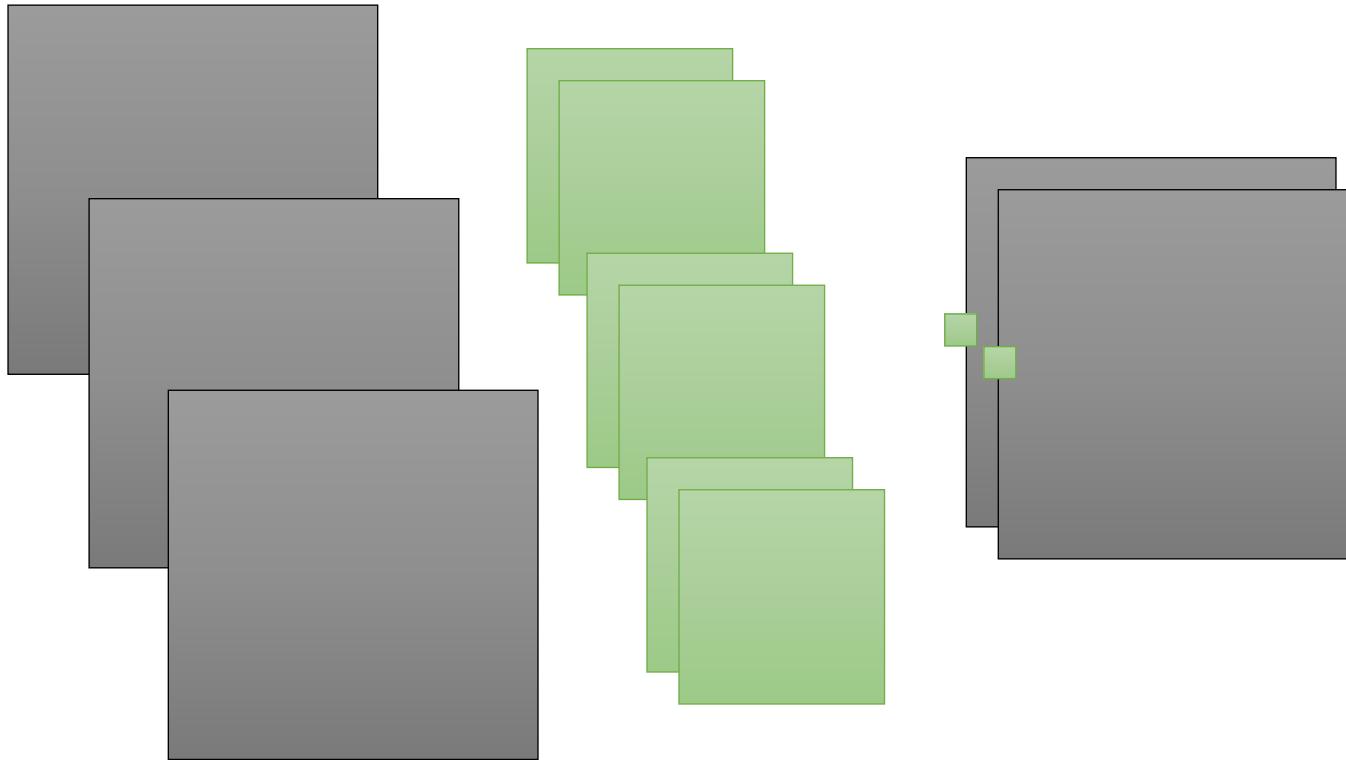


3 input maps

2 filters 5x5  
= 150 ...

2 output maps  
+ 2 biases

# CNN Arithmetic



3 input maps

2 filters 5x5

= 150 ...

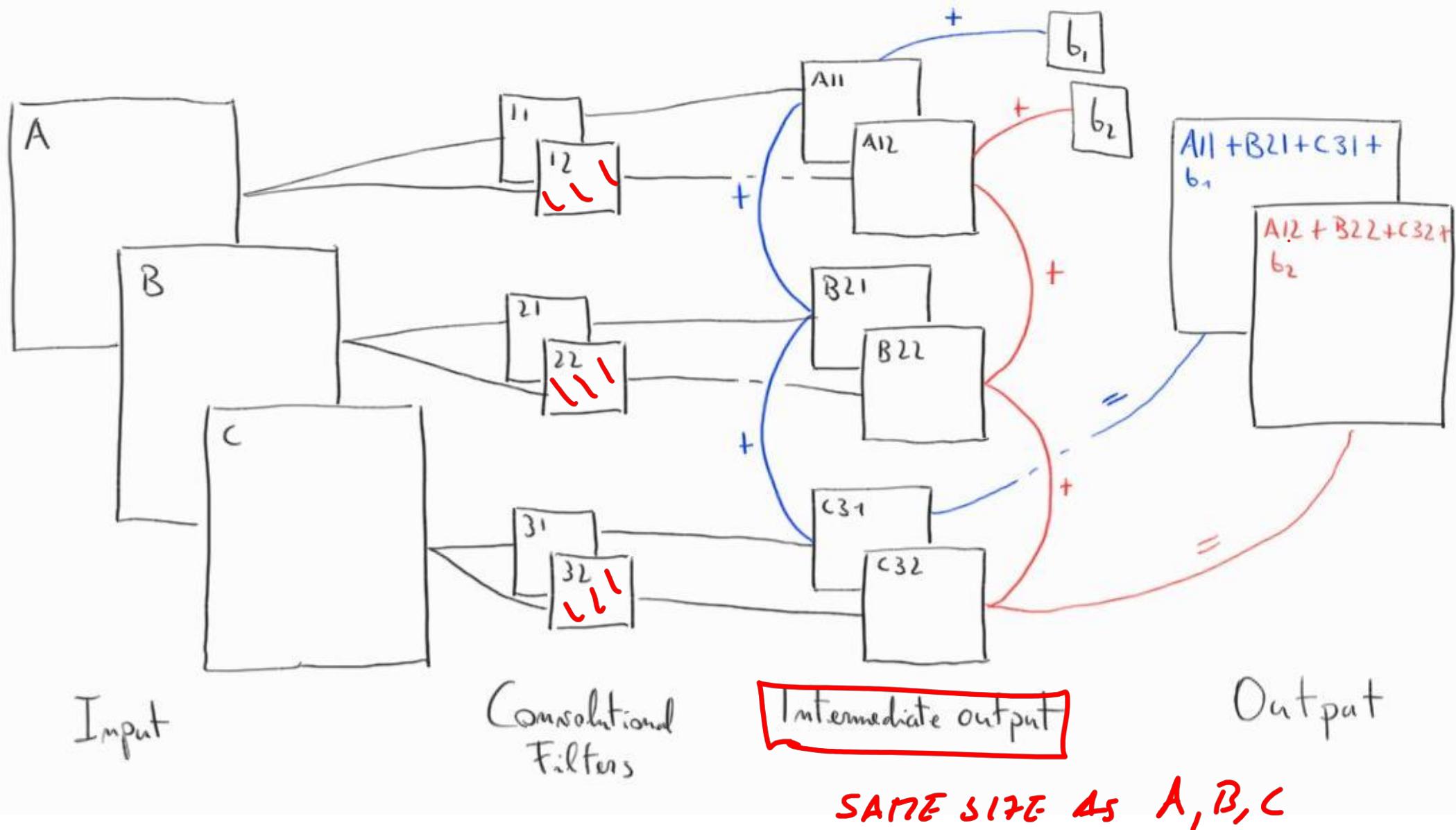
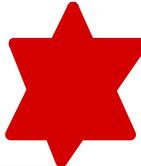
= 152 trainable parameters (weights)

2 output maps

+ 2 biases

# output maps  
= filters  
# of in Re layer  
=  $N_F$

# To Recap...



(\*) They are given by activation functions.



## Other Layers

Activation and Pooling

Not only convolution, otherwise  
everything will be linear...  
→ we need non-linearities (\*)



# Activation Layers



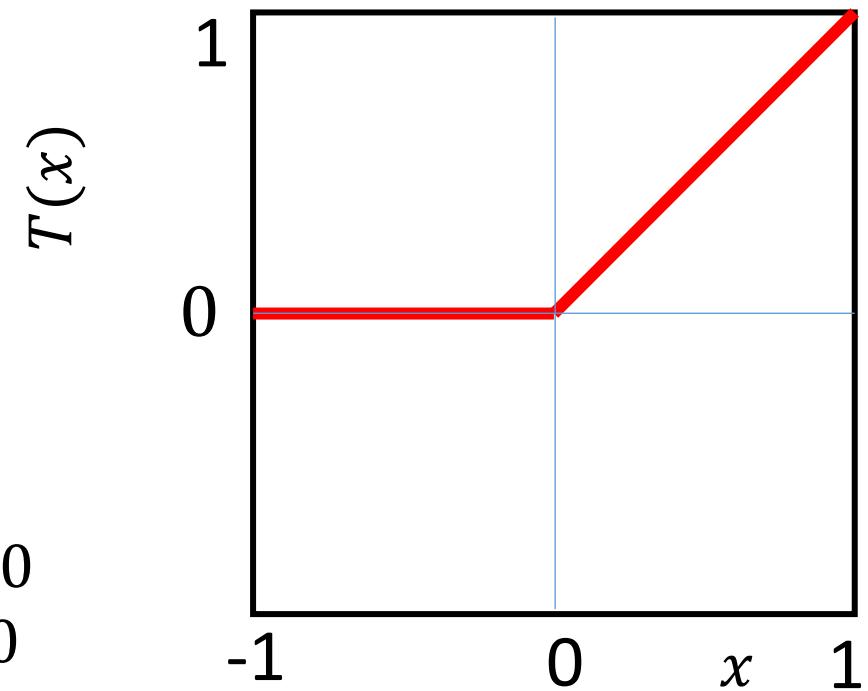
Introduce nonlinearities in the network, otherwise the CNN might be equivalent to a linear classifier...

Activation functions are scalar functions, namely they operate on each single value of the volume. Activations don't change volume size

RELU (Rectifier Linear Units): it's a thresholding on the feature maps, i.e., a  $\max(0, \cdot)$  operator.

- By far the most popular activation function in deep NN (since when it has been used in AlexNet)
- Dying neuron problem: a few neurons become insensitive to the input (vanishing gradient problem)

$$T(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0, & \text{if } x < 0 \end{cases}$$





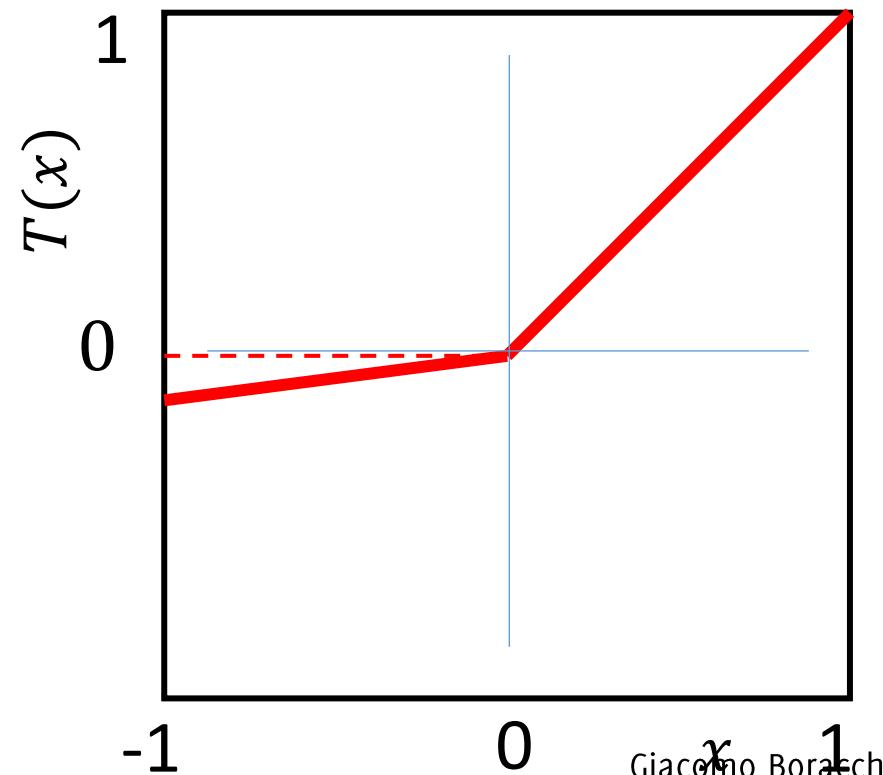
# Activation Layers



Introduce nonlinearities in the network, otherwise the CNN might be equivalent to a linear classifier...

**LEAKY RELU:** like the relu but include a small slope for negative values

$$T(x) = \begin{cases} x, & \text{if } x \geq 0 \\ 0.01 * x & \text{if } x < 0 \end{cases}$$

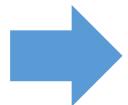
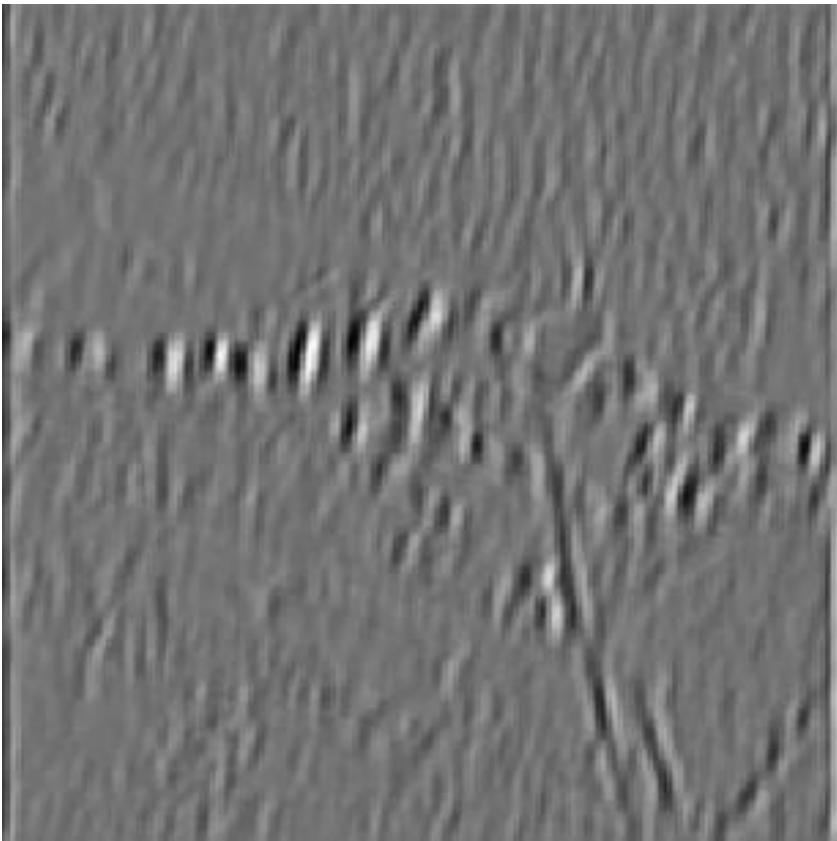


# ReLU

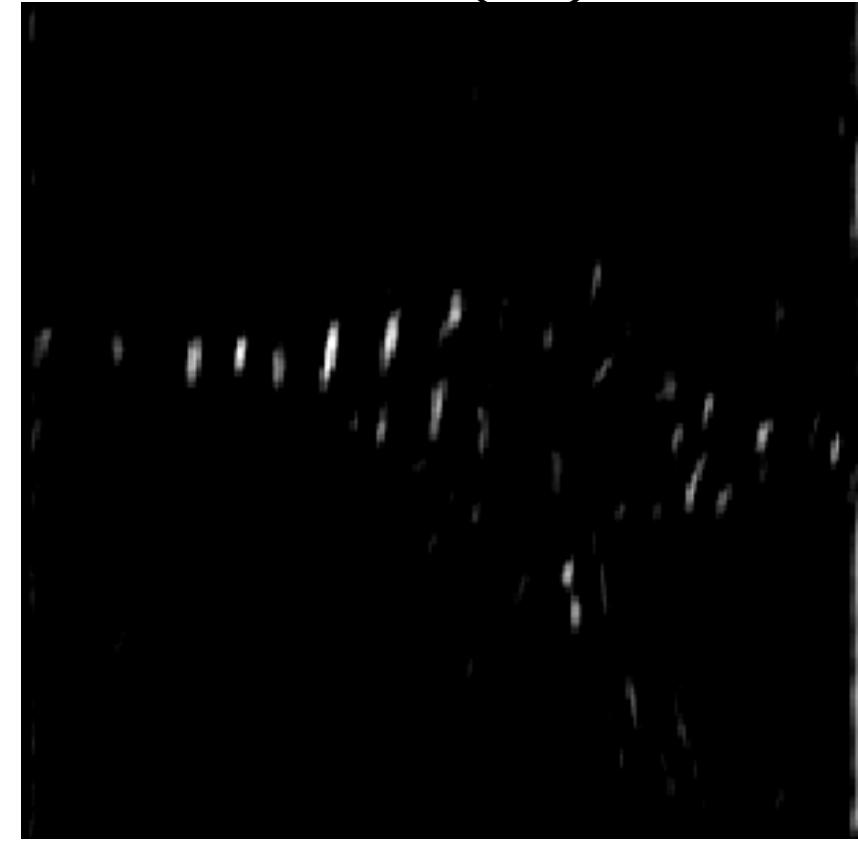


Acts separately on each layer

$$a^1$$



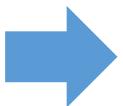
$$\text{ReLU}(a^1)$$



# ReLU

Acts separately on each layer

$$a^2$$



$$\text{ReLU}(a^2)$$





# Activation Layers

Introduce nonlinearities in the network, otherwise the CNN might be equivalent to a linear classifier...

**TANH** (**hyperbolic Tangent**): has a range (-1,1), continuous and differentiable

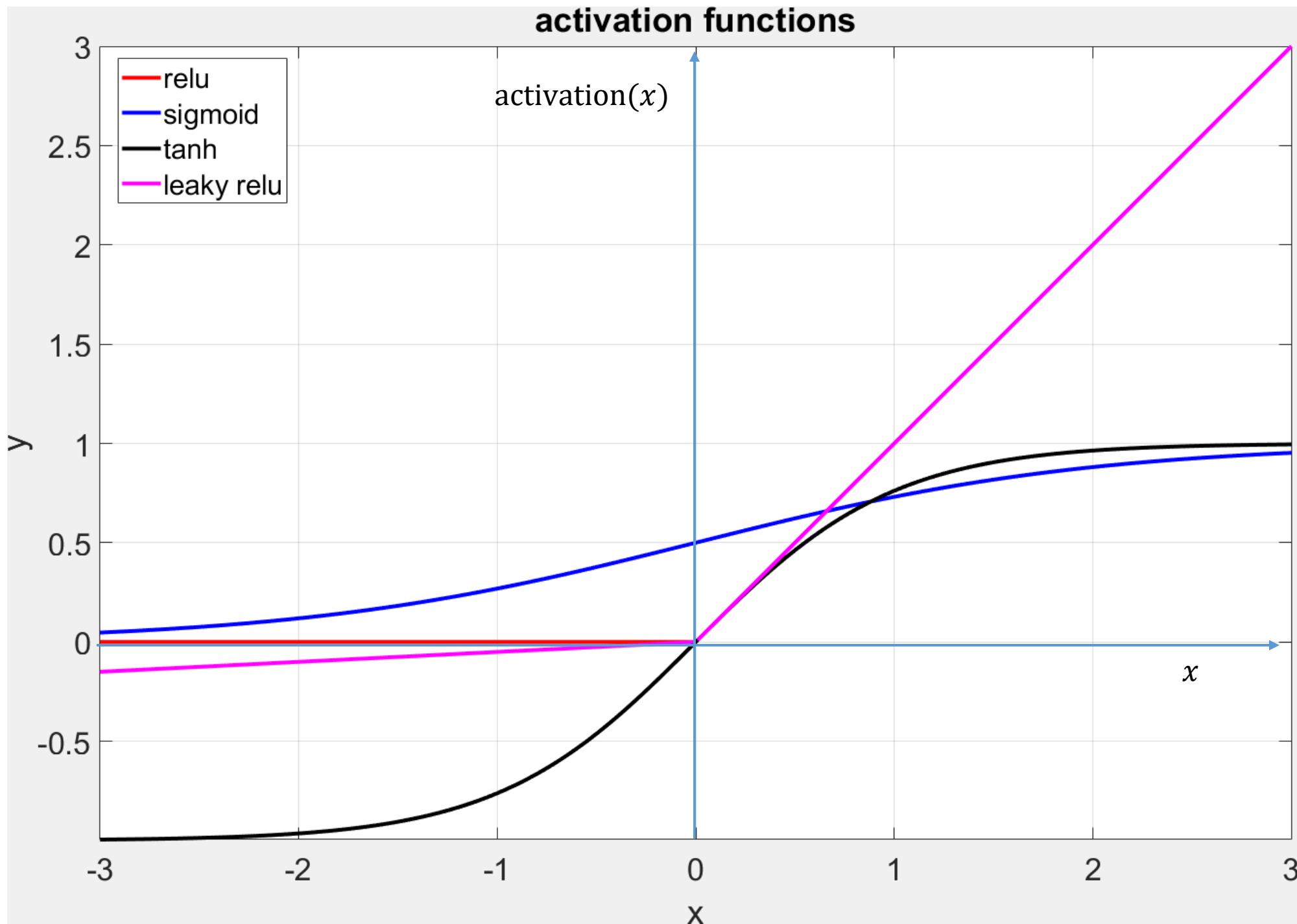
$$T(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} - 1$$

**SIGMOID**: has a range (0,1), continuous and differentiable

$$S(x) = \frac{1}{1 + e^{-x}}$$

These activation functions are mostly popular in **MLP architectures**

# activation functions

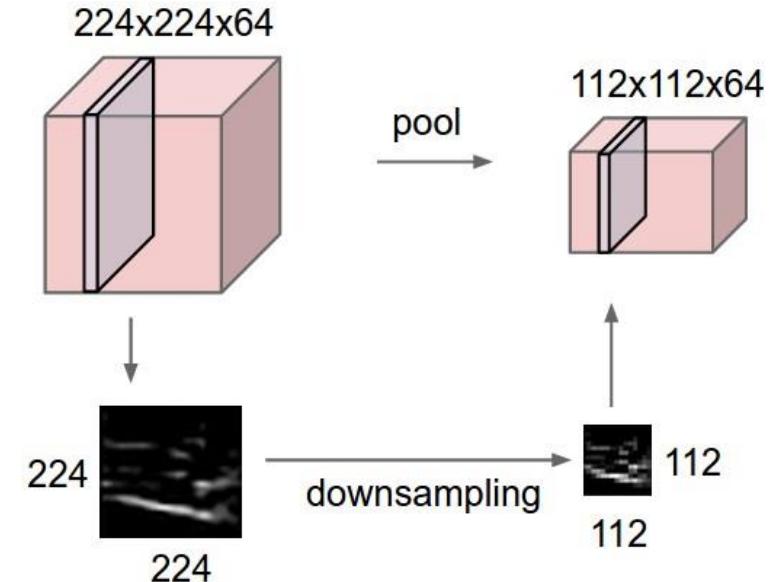
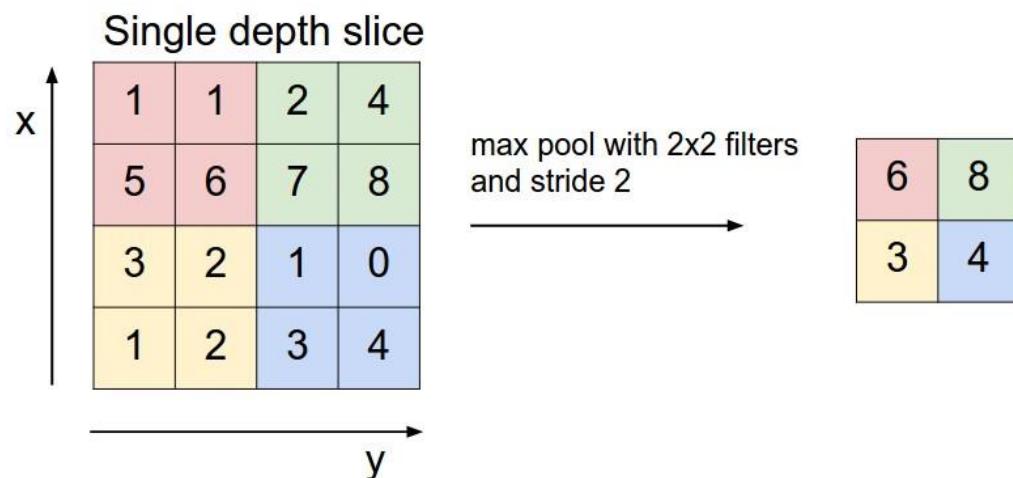




# Pooling Layers

Pooling Layers reduce the spatial size of the volume.

The Pooling Layer operates independently on every depth slice of the input and resizes it spatially, often using the MAX operation.



In a 2x2 support it discards 75% of samples in a volume

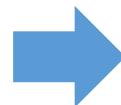
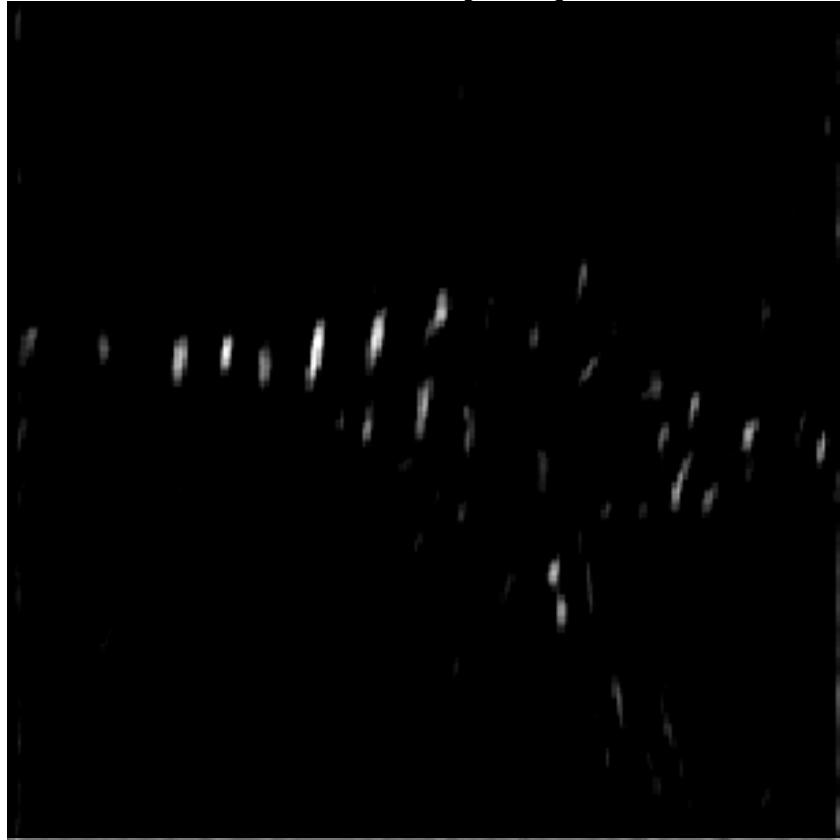


# Max-Pooling (MP)

Acts separately on each layer

It's an additional non-linearity.

$$ReLU(a^1)$$



$$MP(ReLU(a^1))$$





# Strides in Pooling Layers

so it reduces  
the size ...

Typically, the stride is assumed equal to the pooling size

- Where note specified, maxpooling has stride  $2 \times 2$  and reduces image size to 25%

It is also possible to use a different stride. In particular, it is possible to adopt stride = 1, which does not reduce the spatial size, but just perform pooling on each pixel

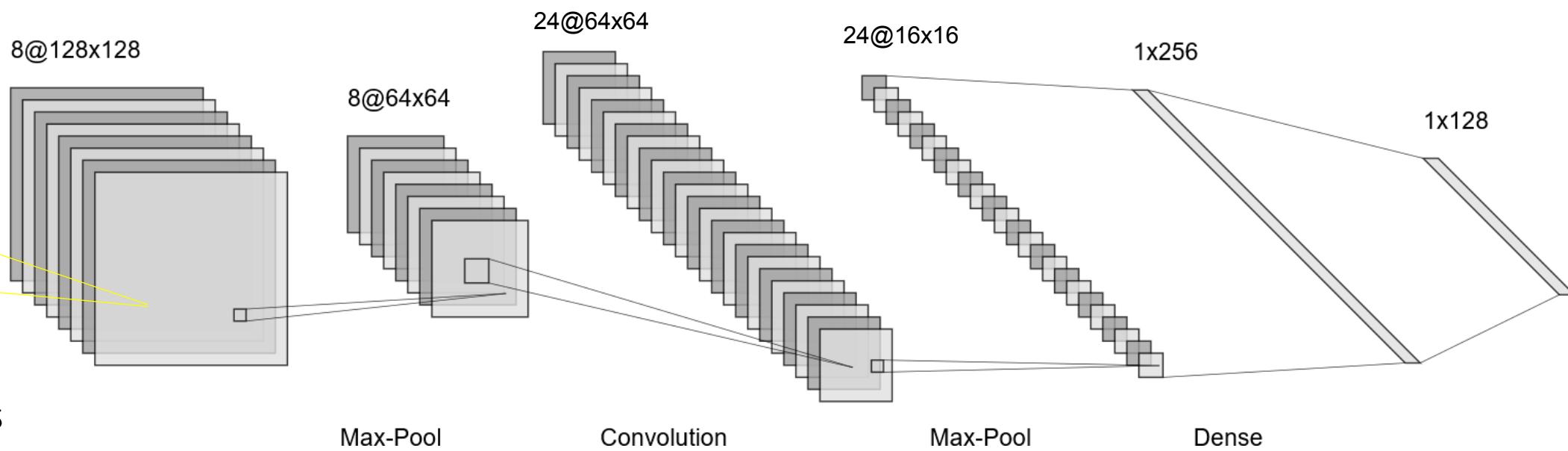
- this operation makes sense with nonlinear pooling (max-pooling)

# Dense Layers

As in feed-forward NN

# The typical architecture of a CNN

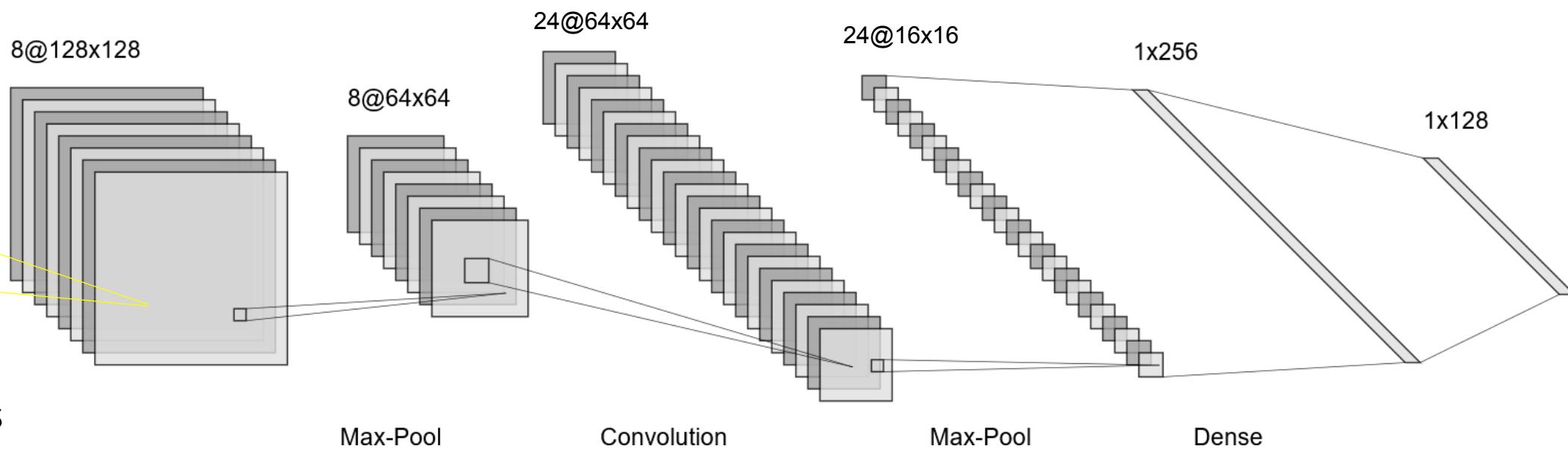
The input of a CNN  
is an entire image



The image gets  
convolved against  
many filters

# The typical architecture of a CNN

The input of a CNN  
is an entire image



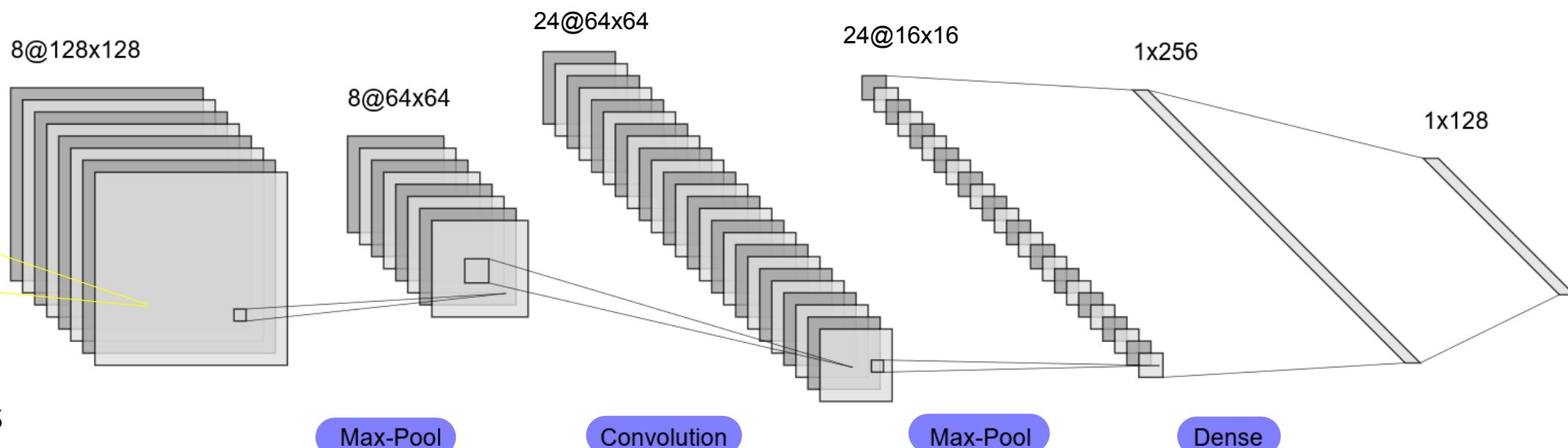
The image gets  
convolved against  
many filters

When progressing along the network, the «number of images» or the «number of channels in the images» increases, while the image size decreases



# The typical architecture of a CNN

The input of a CNN  
is an entire image

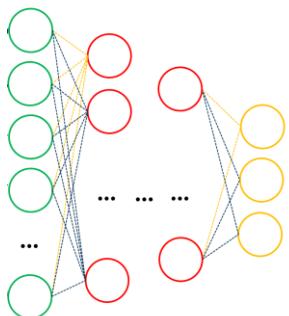


The image gets  
convolved against  
many filters

When progressing along the network, the «number of images» or the «number of channels in the images» increases, while the image size decreases



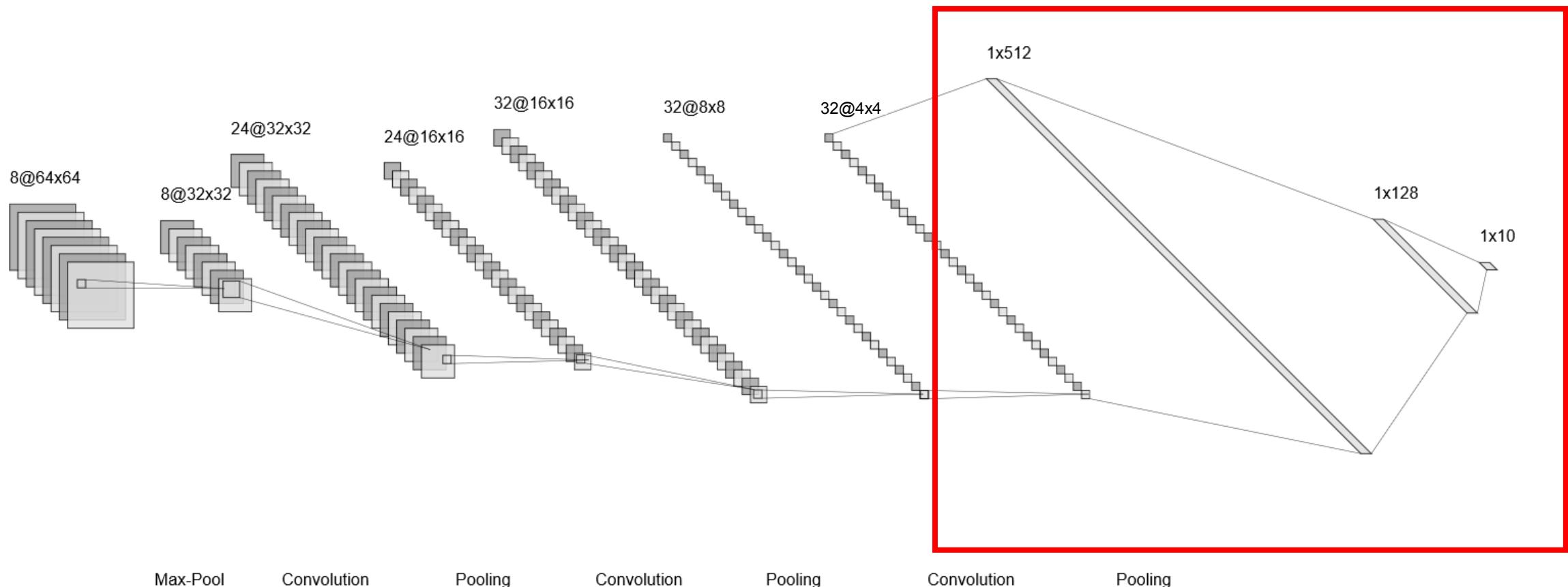
Once the image  
gets to a vector,  
this is fed to a  
traditional  
neural network



# The Dense Layers



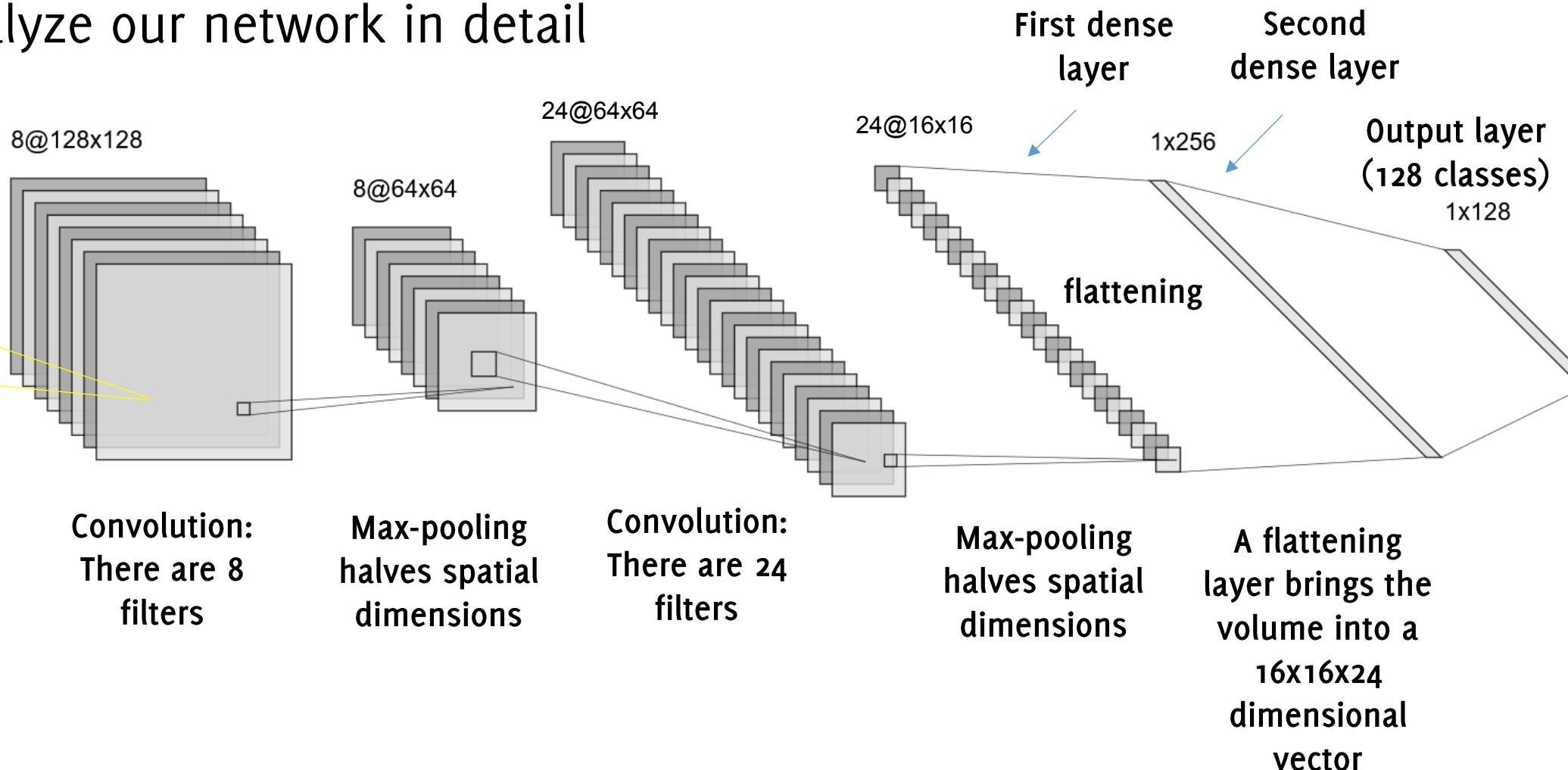
Here the spatial dimension is lost, the CNN stacks hidden layers from a MLP NN.  
It is called Dense as each output neuron is connected to each input neuron



# Convolutional Neural Networks (CNN)



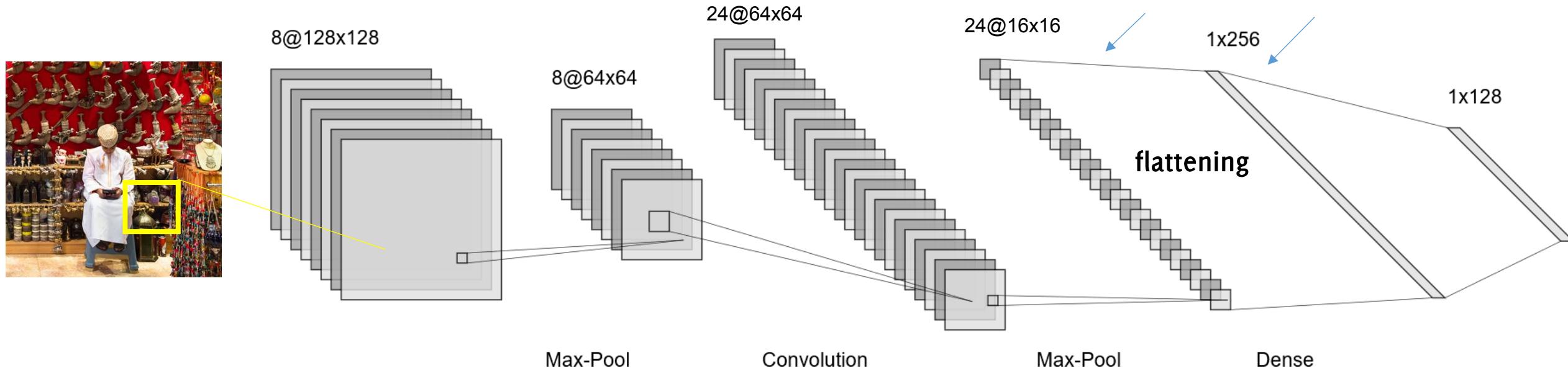
Let's analyze our network in detail



# Convolutional Neural Networks (CNN)



Let's analyze our network in detail



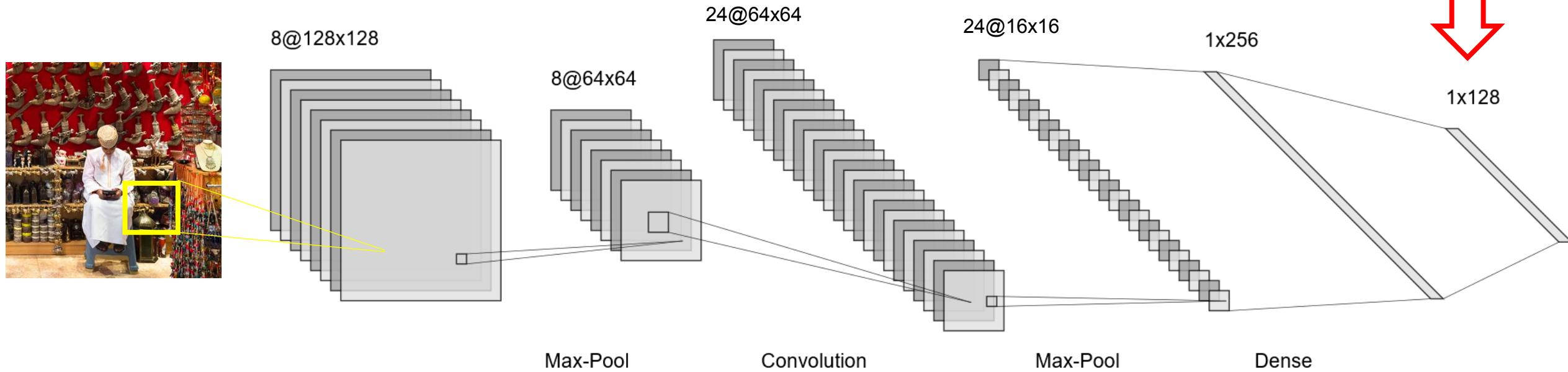
Flattening is needed to obtain a vector to be fed to the Dense layers in charge of classification.

Flattening is a good alternative to perform pooling till the spatial dimension gets reduced to 1

# Convolutional Neural Networks (CNN)

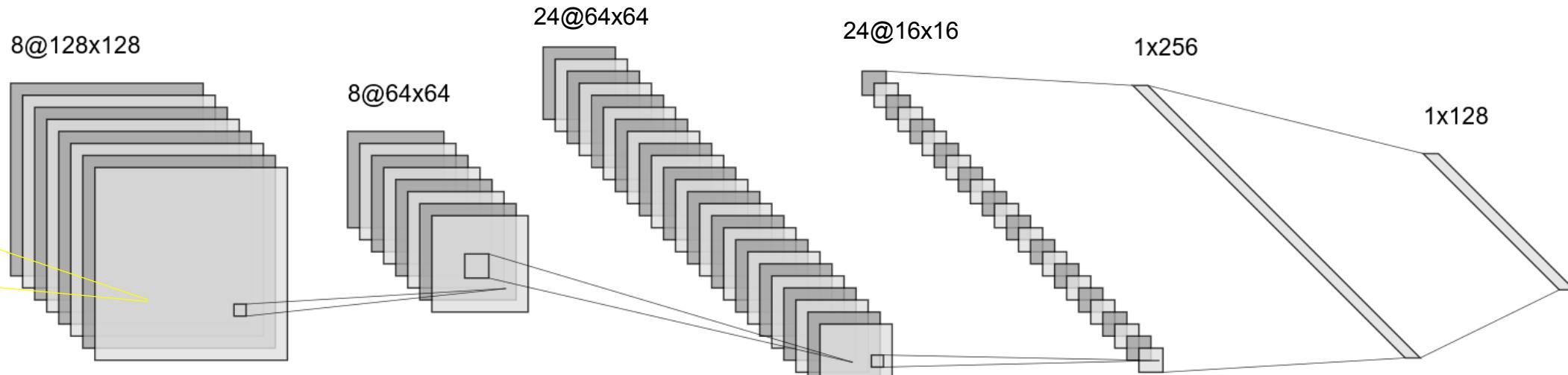


Let's analyze our network in detail



The output of the **fully connected (FC) layer** has the same size as the **number of classes**, and provides a **score** for the input image to belong to each class

# Convolutional Neural Networks (CNN)



Conv

Subsampling

Conv.

Subsampling

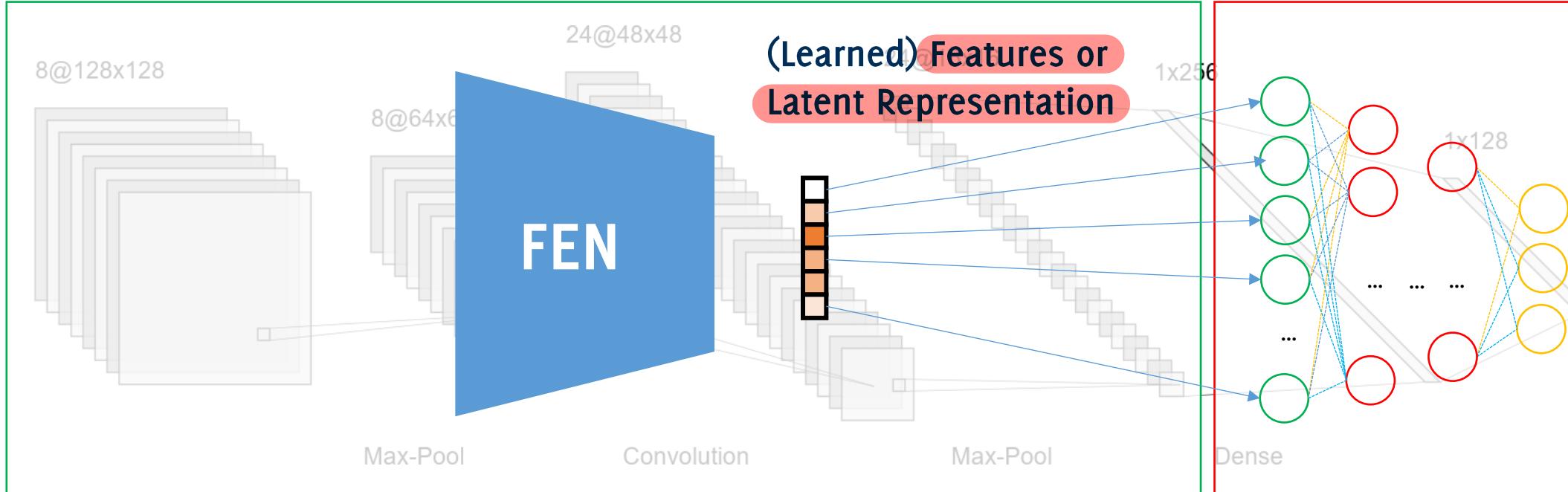
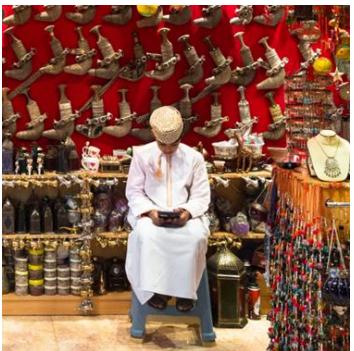
Dense

Convolution filters are learned for the classification task at hand

Thresholding +  
Downsampling  
(ReLU + Maxpooling)

Fully Connected  
Neural Network  
providing as output  
class probabilities

# The typical architecture of a CNN

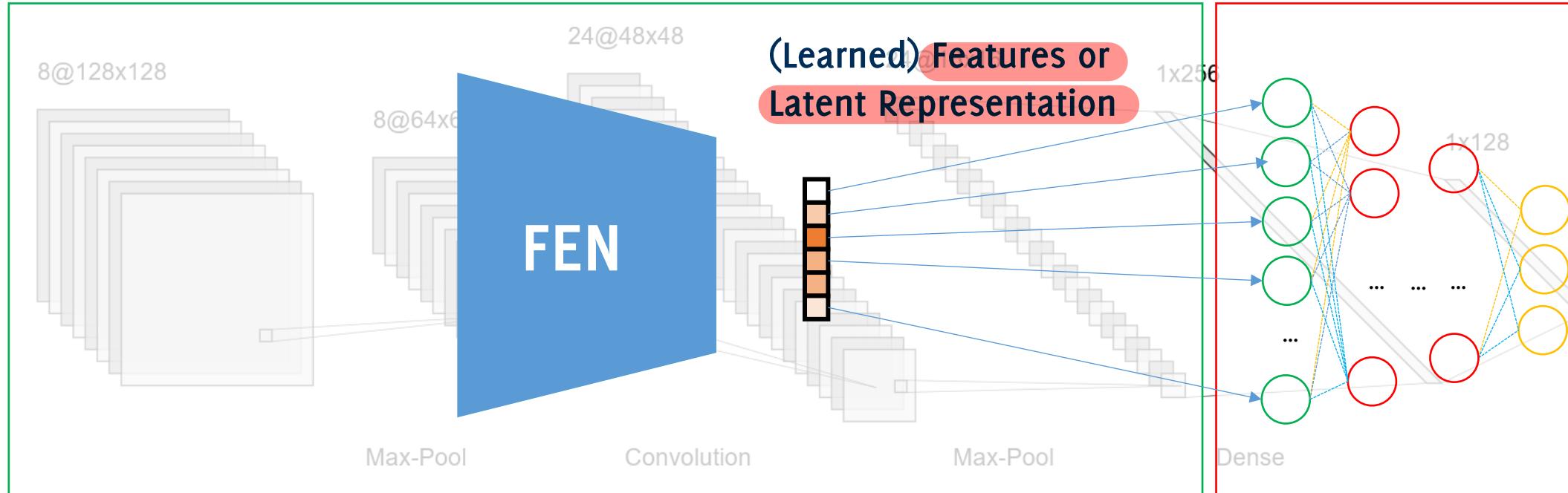
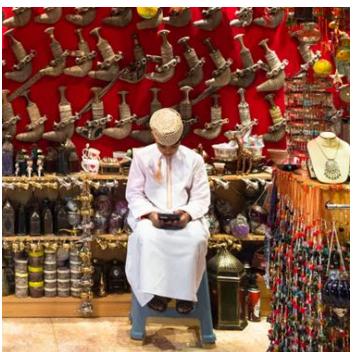
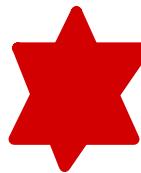


Data-Driven Feature  
extraction

Feature  
Classification

**FEN: FEATURE EXTRACTION NETWORK**, the convolutional block of CNN

# The typical architecture of a CNN



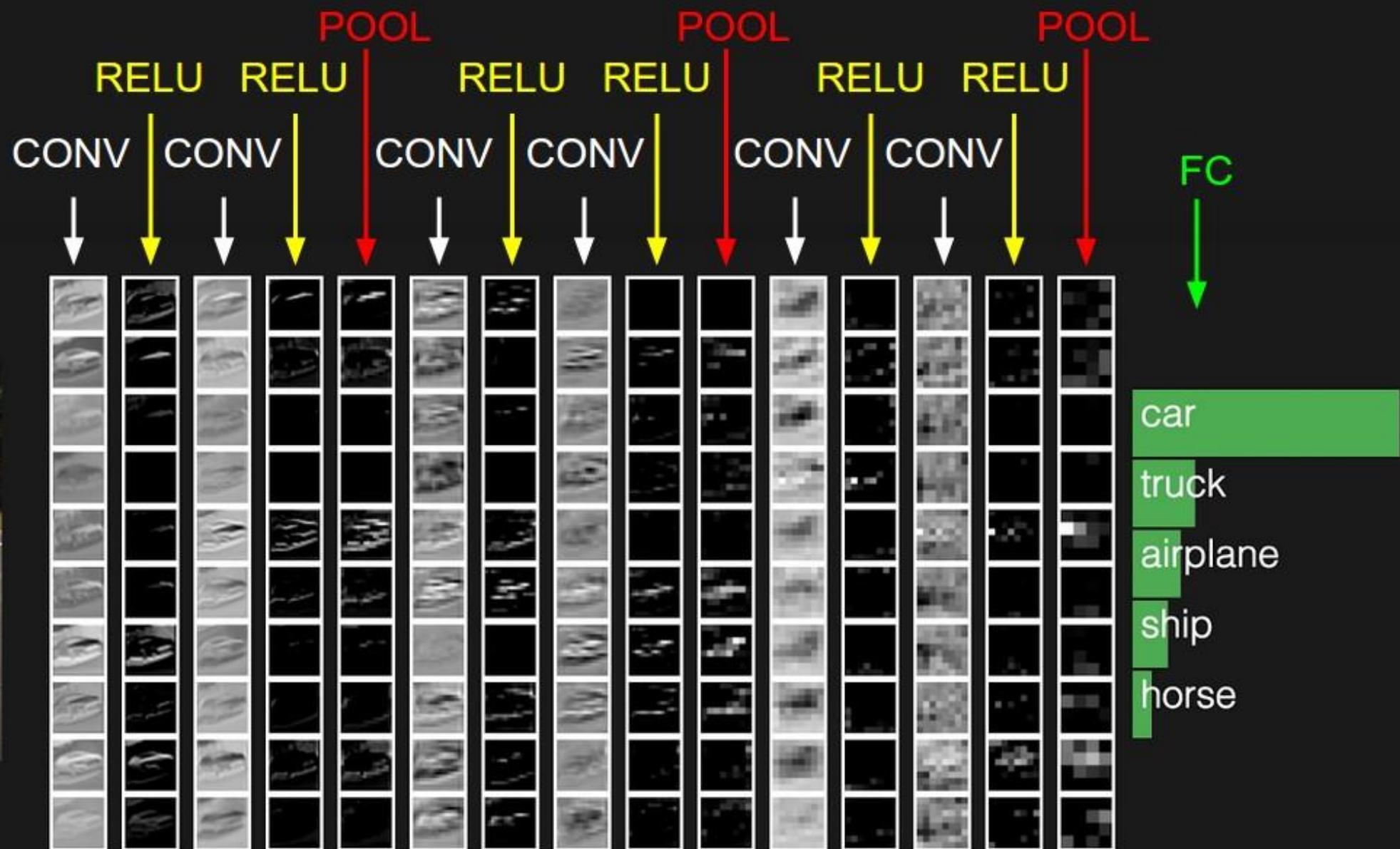
Data-Driven Feature  
extraction

Feature  
Classification

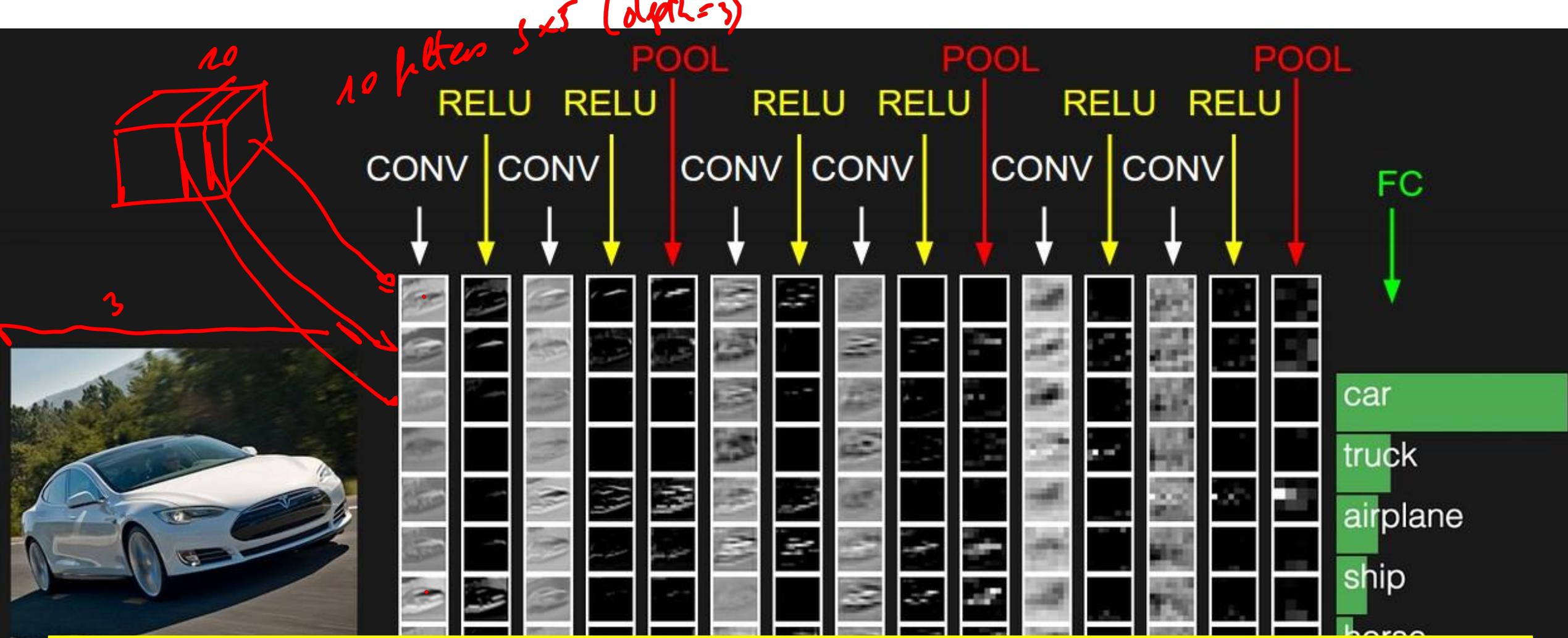
Typically, to learn meaningful representations, many layers are required  
The network becomes deep

CNN «in action»

# Activations in a convolutional network

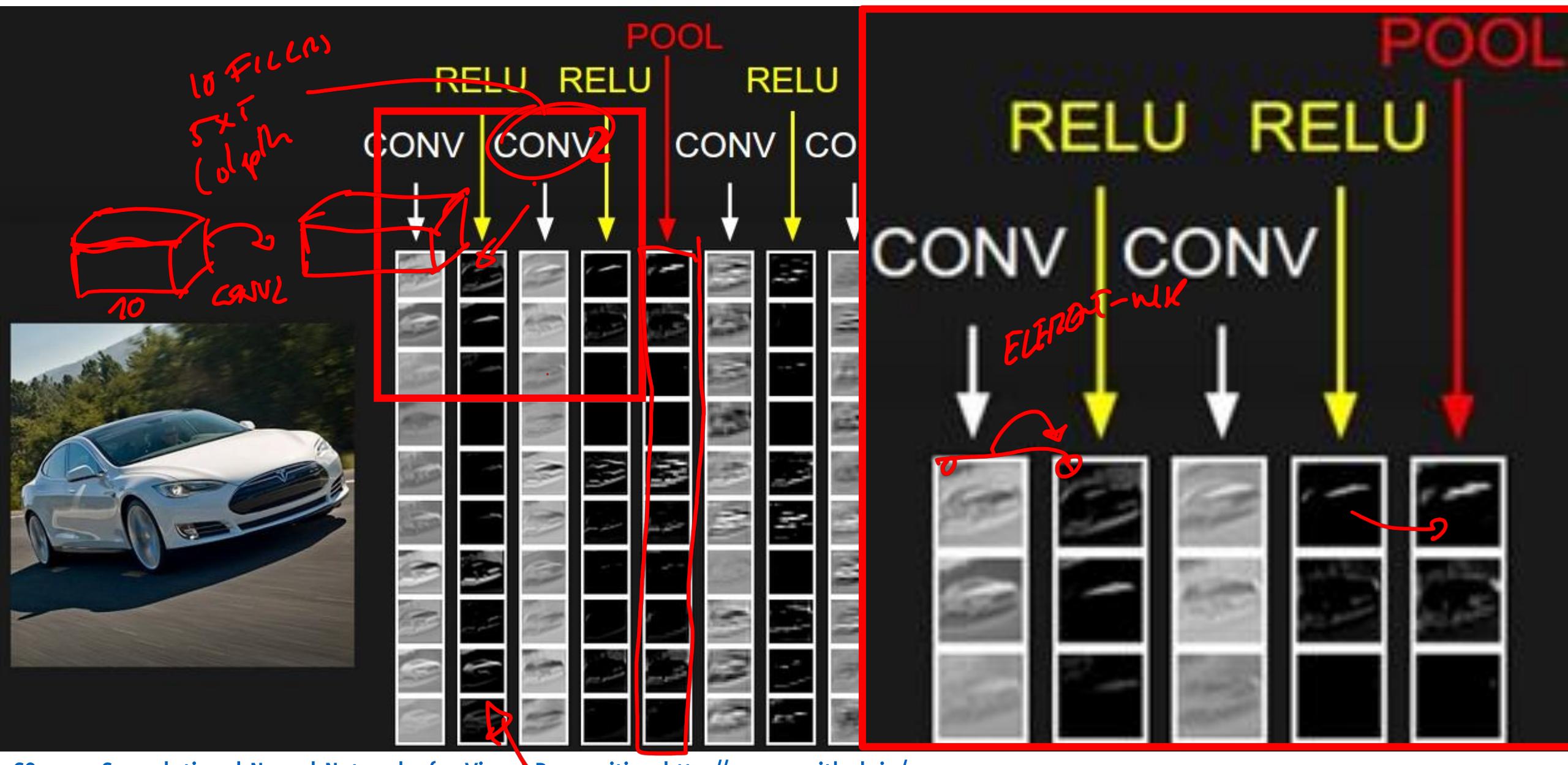


# Activations in a convolutional network

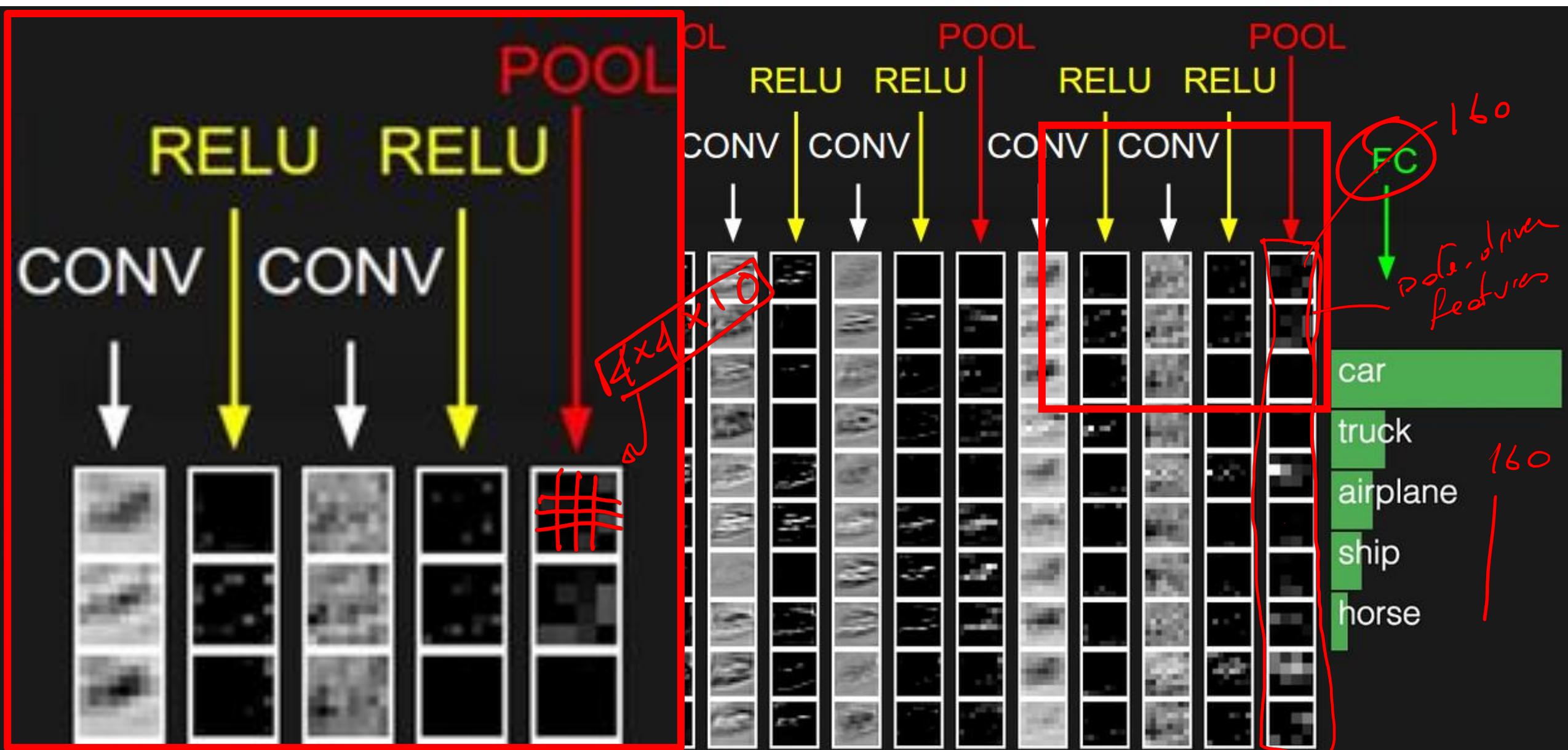


Each layer in the volume is represented as an image here  
(using the same size but different resolution for visualization sake)

# Activations in a convolutional network

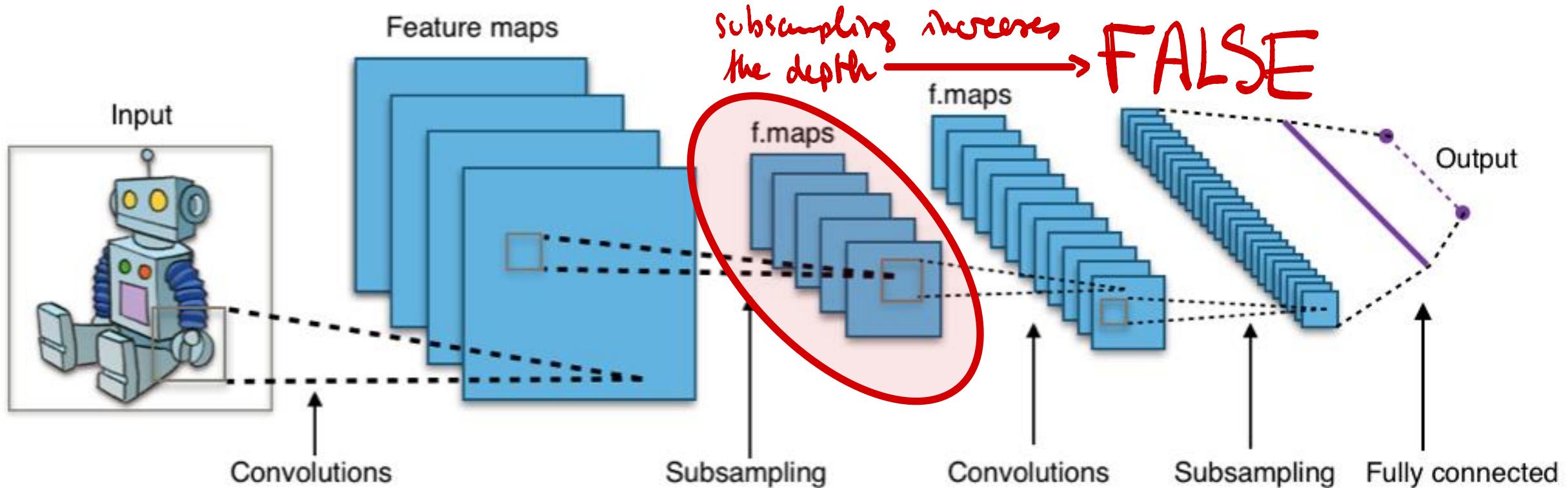


# Activations in a convolutional network



# Convolutional Neural Networks (CNN)

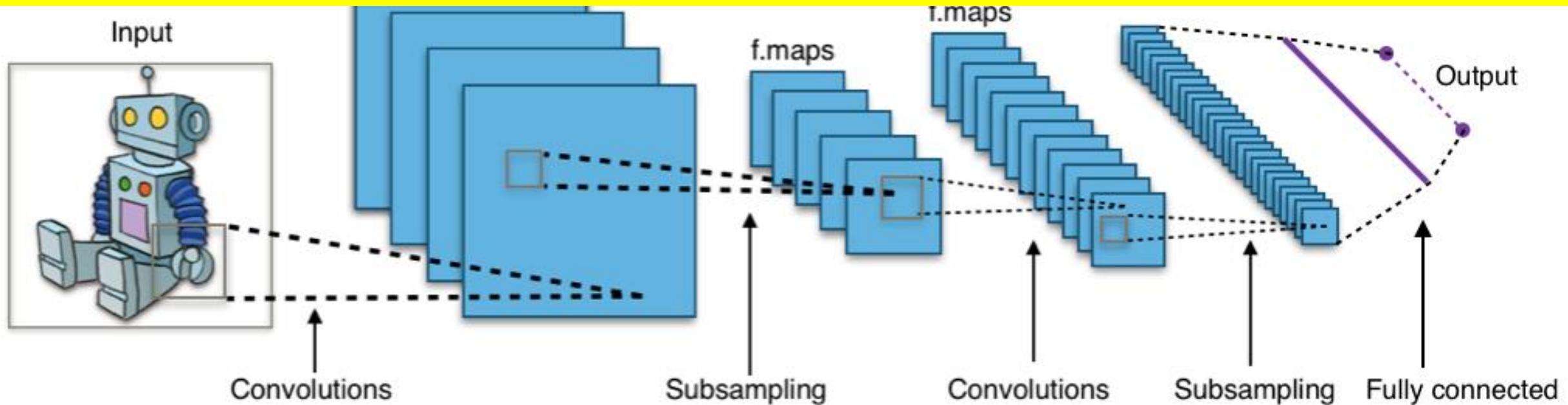
Btw, this figure contains an error.  
If you are CNN-Pro, you should spot it!



# Convolutional Neural Networks (CNN)

Btw, this figure contains an error.  
If you are CNN-Pro, you should spot it!

However, that could be in principle addressable with strided convolution,...  
but I don't think the designer was aware of that



# The First CNN

LeNet



# Gradient-Based Learning Applied to Document Recognition

Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner

## *Abstract—*

Multilayer Neural Networks trained with the backpropagation algorithm constitute the best example of a successful Gradient-Based Learning technique. Given an appropriate network architecture, Gradient-Based Learning algorithms can be used to synthesize a complex decision surface that can classify high-dimensional patterns such as handwritten characters, with minimal preprocessing. This paper reviews various methods applied to handwritten character recognition and compares them on a standard handwritten digit recognition task. Convolutional Neural Networks, that are specifically designed to deal with the variability of 2D shapes, are shown to outperform all other techniques.

## I. INTRODUCTION

Over the last several years, machine learning techniques, particularly when applied to neural networks, have played an increasingly important role in the design of pattern recognition systems. In fact, it could be argued that the availability of learning techniques has been a crucial factor in the recent success of pattern recognition applications such as continuous speech recognition and handwriting recognition.

# Fathers of the Deep Learning Revolution Receive ACM A.M. Turing Award

Bengio, Hinton and LeCun Ushered in Major Breakthroughs in Artificial Intelligence

<https://awards.acm.org/about/2018-turing>

# Handwritten digit recognition

- Very small grayscale images.
- The task is to classify each image in the corresponding digit
- This can be used to classify numbers as sequences of digits

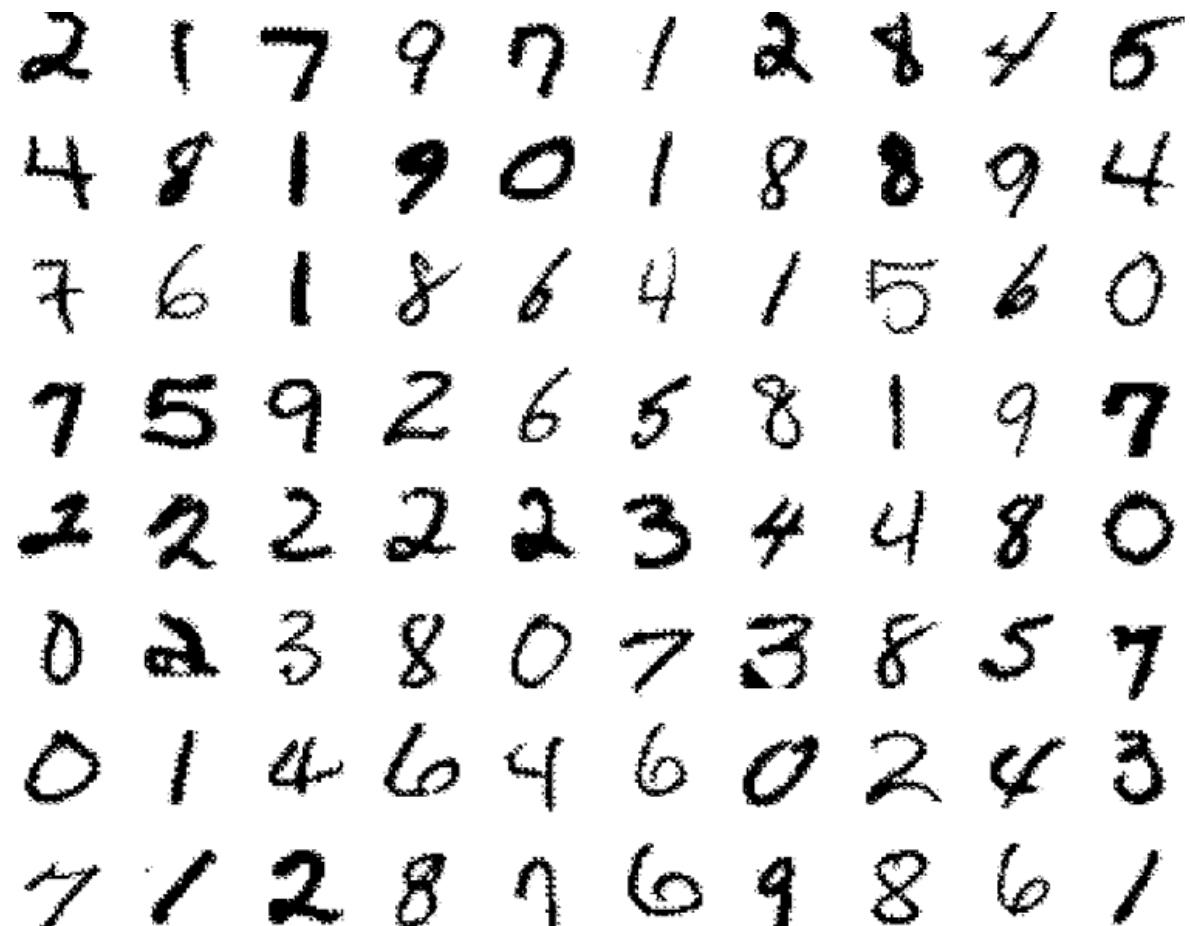
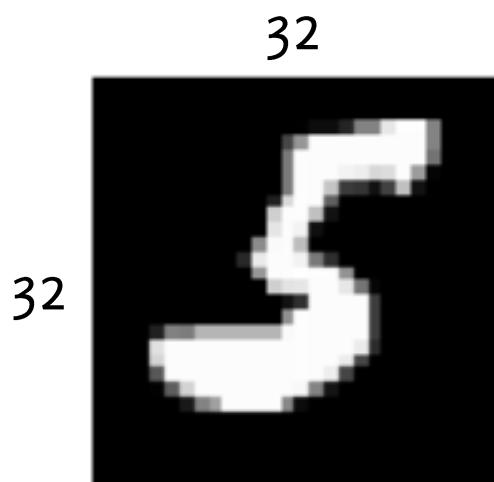
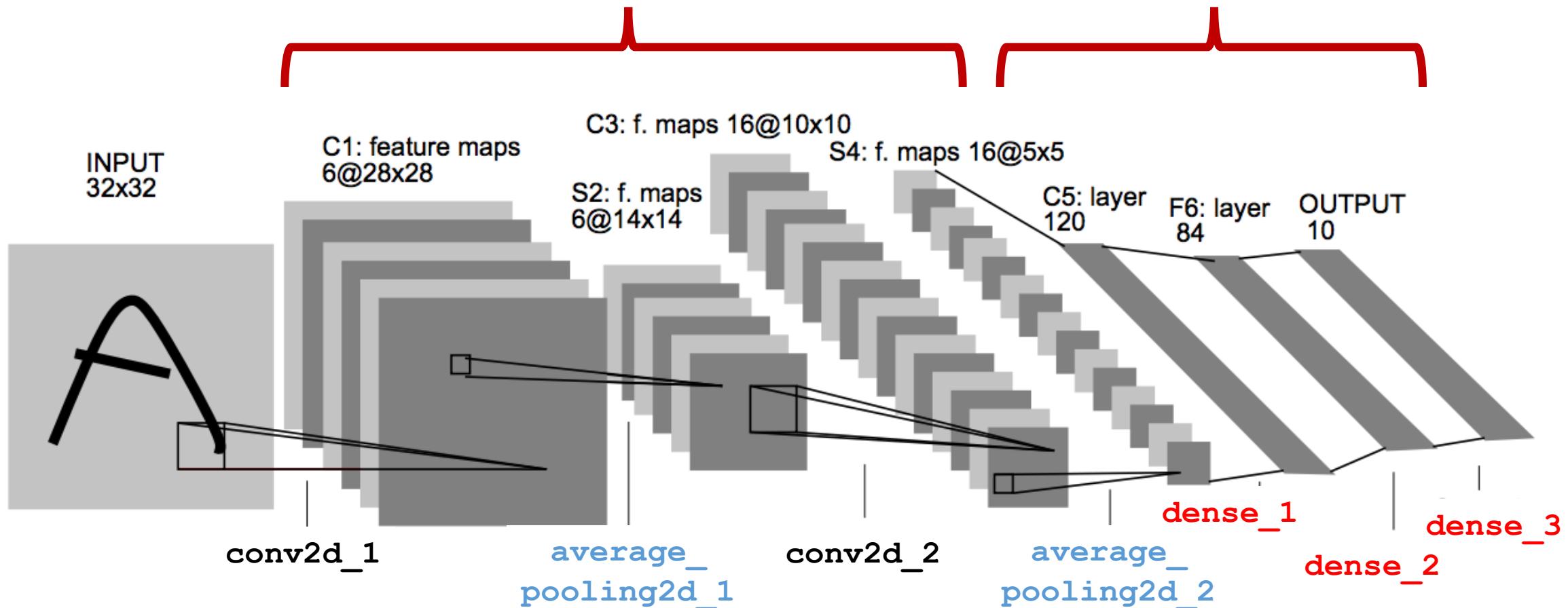


Fig. 4. Size-normalized examples from the MNIST database.

# LeNet-5 (1998)

Stack of Conv2D + RELU + AVG-POOLING      A TRADITIONAL MLP



# The First CNN

Do not use each pixel as a separate input of a large MLP, because:

- images are highly spatially correlated,
- using individual pixel of the image as separate input features would not take advantage of these correlations.

The first convolutional layer: 6 filters  $5 \times 5$

The second convolutional layer: 16 filters  $5 \times 5$

# LeNet-5 in Keras

```
from keras.models import Sequential
from keras.layers import Dense, Flatten, Conv2D, AveragePooling2D

num_classes = 10;
input_shape=(32, 32, 1);

model = Sequential()
model.add(Conv2D(filters = 6, kernel_size = (5, 5), activation='tanh', input_shape=input_shape, padding = 'valid'))
model.add(AveragePooling2D(pool_size=(2, 2)))
model.add(Conv2D(filters = 16, kernel_size = (5, 5), activation='tanh', padding = 'valid'))
model.add(AveragePooling2D(pool_size=(2, 2)))
model.add(Flatten())
model.add(Dense(120, activation='relu'))
model.add(Dense(84, activation='relu'))
model.add(Dense(num_classes, activation='softmax'))
```

# model.summary()

Layer (type)	Output Shape	Param #
=====		
conv2d_1 (Conv2D)	(None, 28, 28, 6)	...
average_pooling2d_1 (Average)	(None, 14, 14, 6)	...
conv2d_2 (Conv2D)	(None, 10, 10, 16)	...
average_pooling2d_2 (Average)	(None, 5, 5, 16)	...
flatten_1 (Flatten)	(None, 400)	...
dense_1 (Dense)	(None, 120)	...
dense_2 (Dense)	(None, 84)	...
dense_3 (Dense)	(None, 10)	...
=====		

Total params: 61,706  
Trainable params: 61,706  
Non-trainable params: 0

# model.summary()

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 28, 28, 6)	$6 \times (5 \times 5 \times 1) + 6 = 156$
average_pooling2d_1 (Average)	(None, 14, 14, 6)	0
conv2d_2 (Conv2D)	(None, 10, 10, 16)	$16 \times (3 \times 3 \times 6) + 16$
average_pooling2d_2 (Average)	(None, 5, 5, 16)	0
flatten_1 (Flatten)	(None, 400)	0
dense_1 (Dense)	(None, 120)	$400 \times 120 + 120$
dense_2 (Dense)	(None, 84)	$120 \times 84 + 84$
dense_3 (Dense)	(None, 10)	$84 \times 10 + 10$

Total params: 61,706

Trainable params: 61,706

Non-trainable params: 0

# model.summary()

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 28, 28, 6)	156 (6 x 5 x 5 + 6)
average_pooling2d_1 (Average)	(None, 14, 14, 6)	0
conv2d_2 (Conv2D)	(None, 10, 10, 16)	2416 (16 x 5 x 5 x 6 + 16)
average_pooling2d_2 (Average)	(None, 5, 5, 16)	0
flatten_1 (Flatten)	(None, 400)	0
dense_1 (Dense)	(None, 120)	$(400+1) \times 120$
dense_2 (Dense)	(None, 84)	48120
dense_3 (Dense)	(None, 10)	10164
Total params:	61,706	850
Trainable params:	61,706	
Non-trainable params:	0	

Input is a grayscale image

The input is a volume having depth = 6

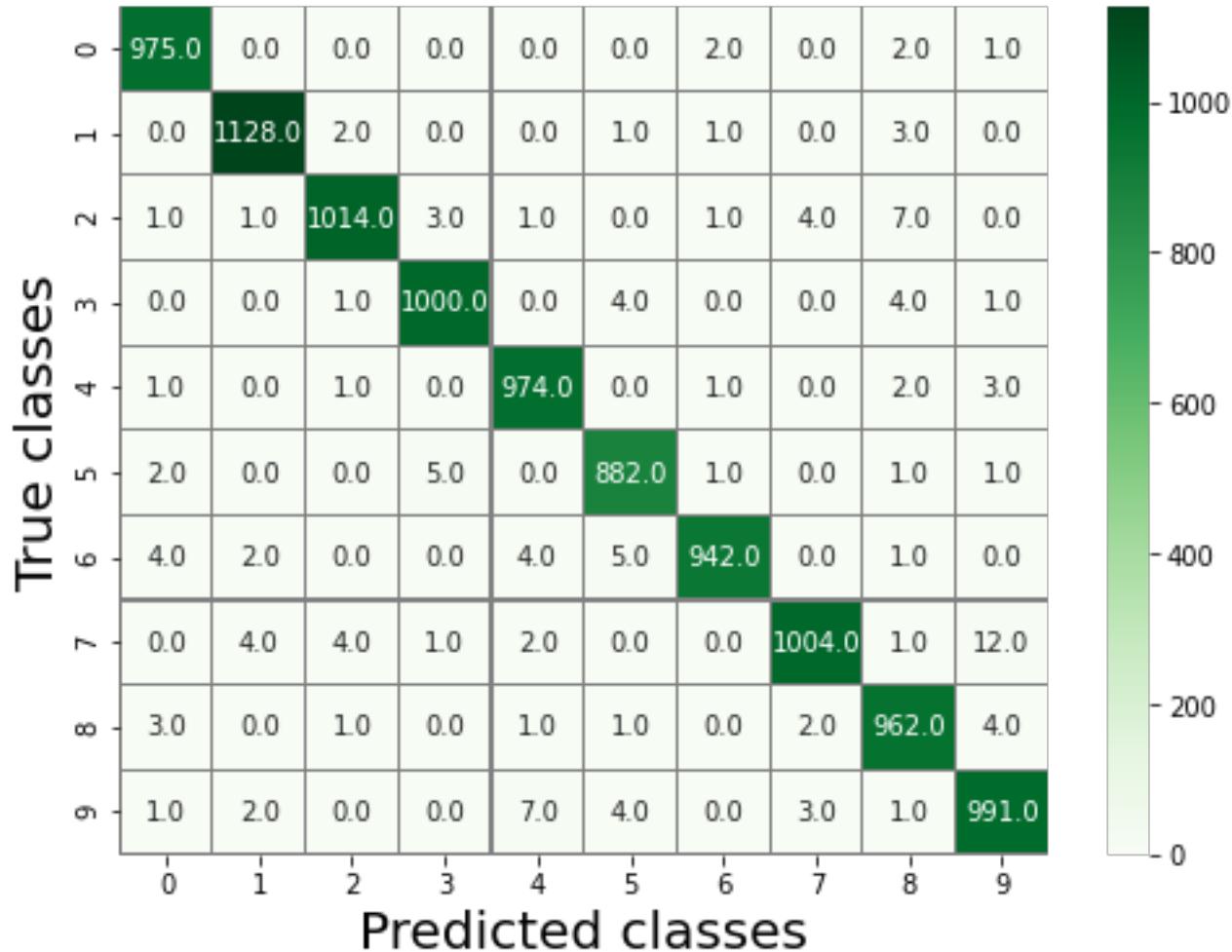
Most parameters are still in the MLP

# model.summary()

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 28, 28, 6)	156 (6 x 5 x 5 + 6)
average_pooling2d_1 (Average)	(None, 14, 14, 6)	0
conv2d_2 (Conv2D)	(None, 10, 10, 16)	2416 (16 x 5 x 5 x 6 + 16)
average_pooling2d_2 (Average)	(None, 5, 5, 16)	0
flatten_1 (Flatten)	(None, 400)	0
dense_1 (Dense)	(None, 120)	48120
dense_2 (Dense)	(None, 84)	10164
dense_3 (Dense)	(None, 10)	850
Total params:	61,706	
Trainable params:	61,706	

Here, no-padding at the first layer is necessary to reduce the size of the latent representation... and has no loss of information since images are black there!

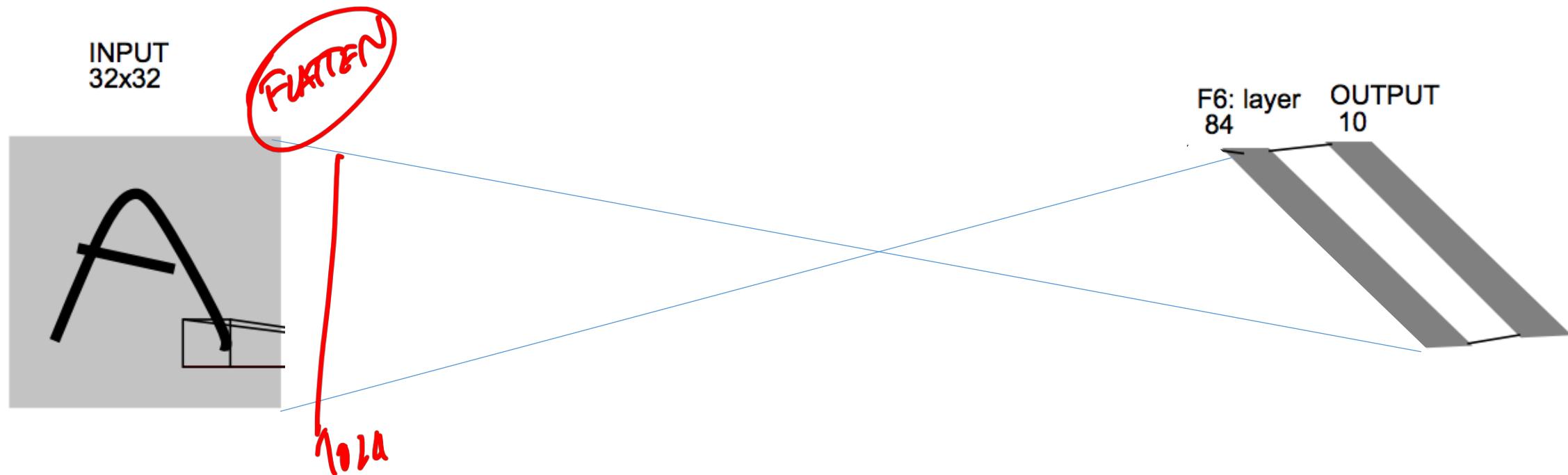
# LeNet-5 performance



# Most of parameters are in MLP

What about a MLP taking as input the whole image?

Input  $32 \times 32 = 1024$  pixels, fed to a 84 neurons (the last FC layers of the network) -> 86950 parameters:  $1024 * 84 + 84 + 84 * 10 + 10$





# Most of parameters are in MLP

What about a MLP taking as input the whole image?

Input  $32 \times 32 = 1024$  pixels, fed to a 84 neurons (the last FC layers of the network)  $\rightarrow$  86950 parameters

But.. If you take an RGB input:  $32 \times 32 \times 3$ ,

CNN: only the nr. of parameters in the filters at the first layer increases

$$\begin{aligned} 156 + 61550 &\rightarrow 456 + 61550 \\ (6 \times 5 \times 5) &\rightarrow (6 \times 5 \times 5 \times 3) \end{aligned}$$

MLP: only the first layer increases the # of parameters by a factor 3

$$1024 \times 84 \rightarrow 1024 \times 84 \times 3$$

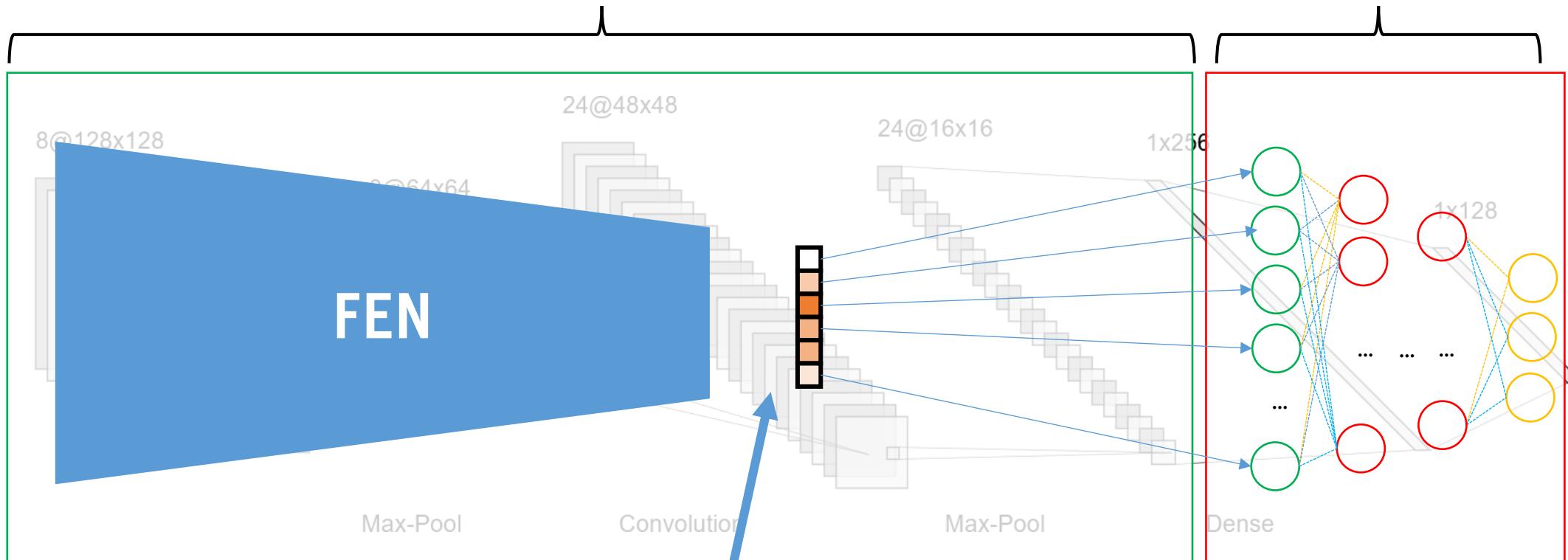
# Latent representation in CNNs

Repeat the «t-SNE experiment» on the CIFAR dataset,  
using the last layer of the CNN as vectors

# The typical architecture of a CNN

**Convolutional Layers**  
Extract high-level features from pixels

**Classify**



**Latent Representation:**  
Data-Driven Feature Vector

**MLP for feature  
classification**

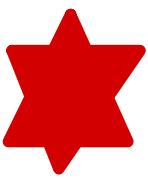


t-SNE of the 4096 dimensional latent representations of AlexNet, right before the classifier (test on ImageNet images)  
 $d(I_1, I_2) = \|FEN(I_1) - FEN(I_2)\|_2$



*It seems  
really  
nice.*





Distances in the latent  
representation of a CNN are  
much more meaningful than  
data itself

# t-SNE representation using $\ell_2$ distance



Bad!



# CNNs in Keras

# What is Keras?

An open-source library providing **high-level building blocks** for developing deep-learning models in Python

Designed to enable **fast experimentation with deep neural networks**, it focuses on being **user-friendly, modular, and extensible**

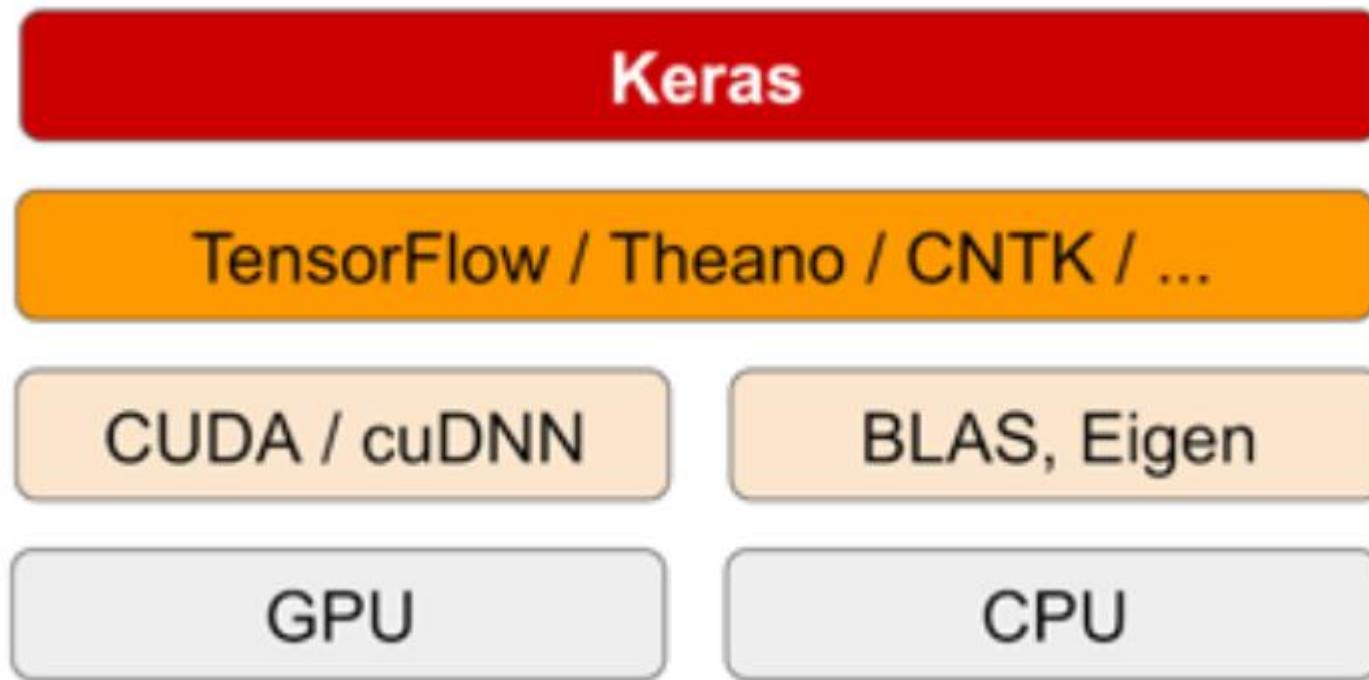
**Doesn't handle low-level operations** such as tensor manipulation and differentiation.

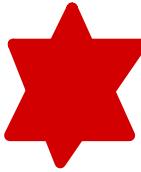
Relies on **backends** (TensorFlow, Microsoft Cognitive Toolkit, Theano, or PlaidML)

Enables full access to the backend



# The software stack





# Why Keras?

## Pros:

Higher level → fewer lines of code

Modular backend → not tied to tensorflow

Way to go if you focus on applications

## Cons:

Not as flexible

Need more flexibility? Access the backend directly!

# We will manipulate 4D tensors

Images are represented in 4D tensors:

Tensorflow convention: (samples, height, width, channels)

A detailed presentation on  
CNN in Keras will be made in  
the lab session

