

# Computer vision

---

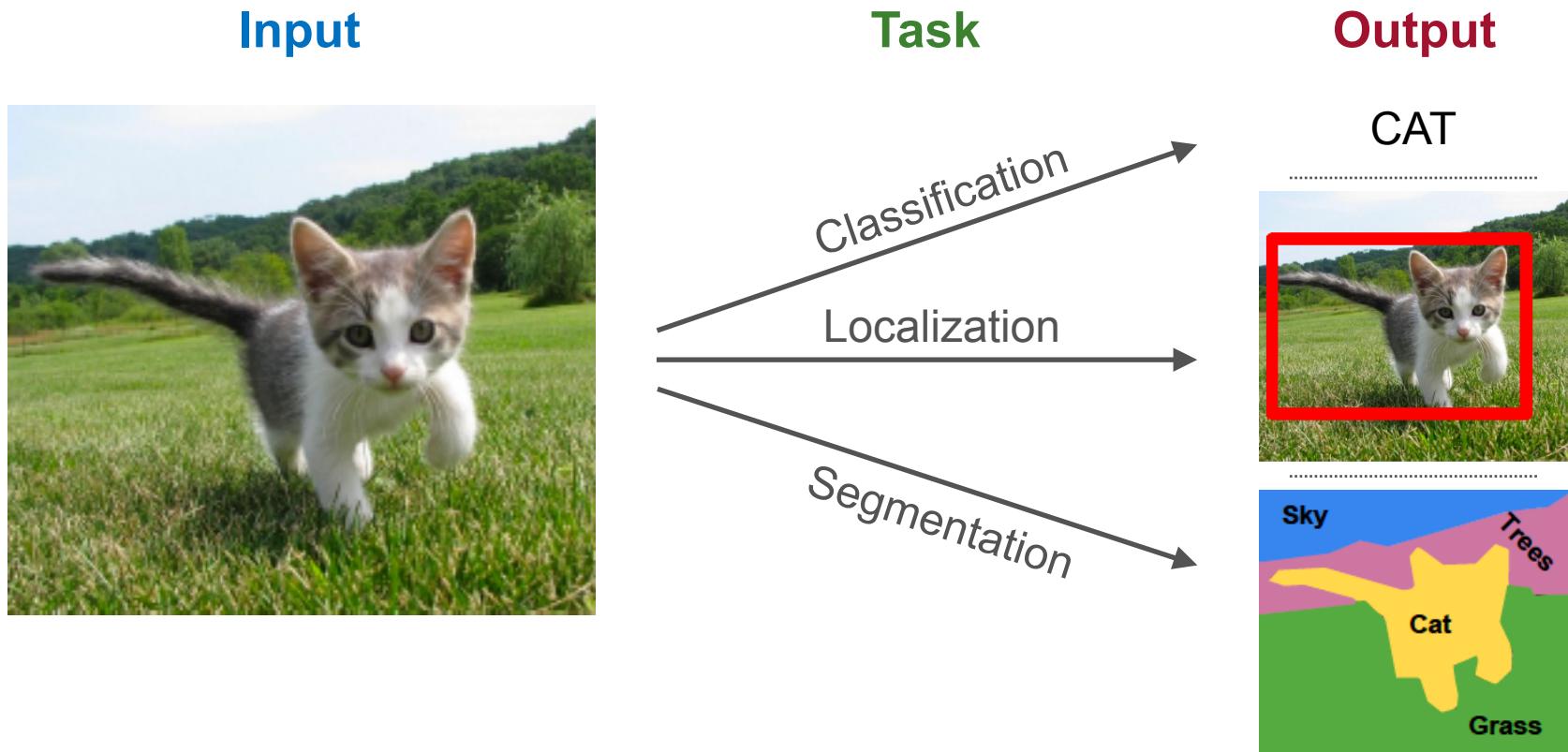
Image-related tasks

Fully-connected layer

Bottleneck of standard networks

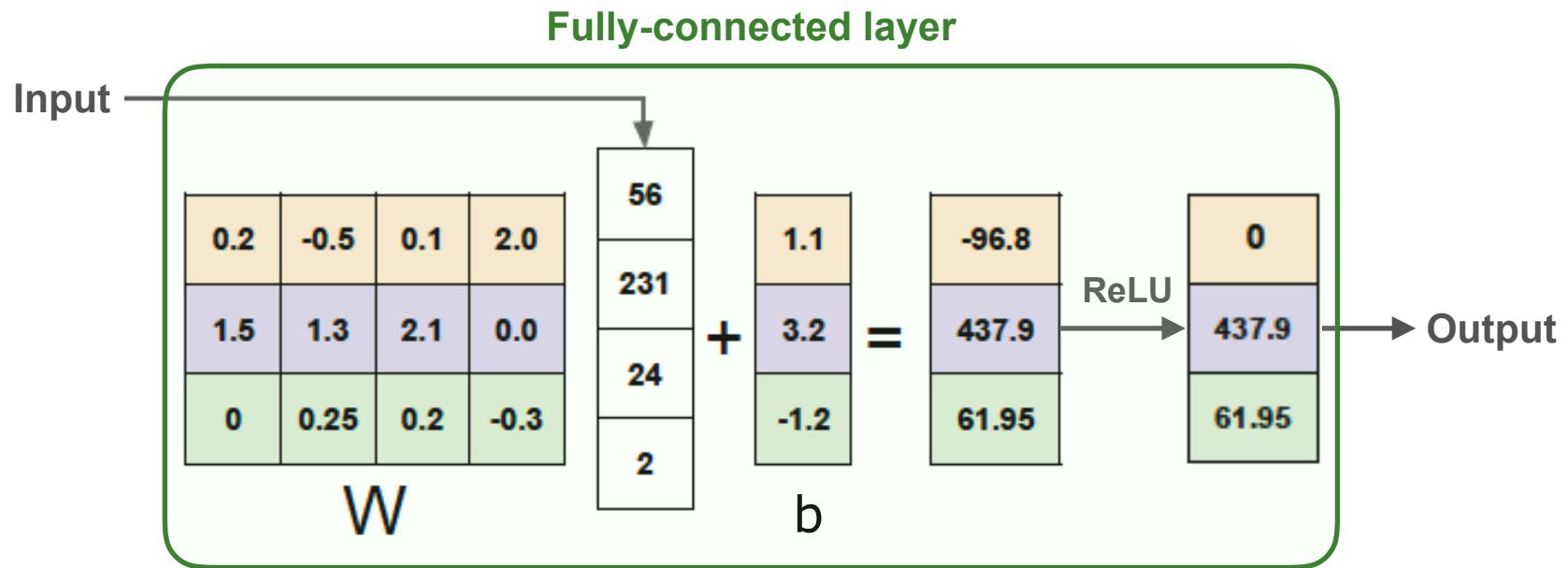
# Image-related tasks

- Neural networks can be used for many things...



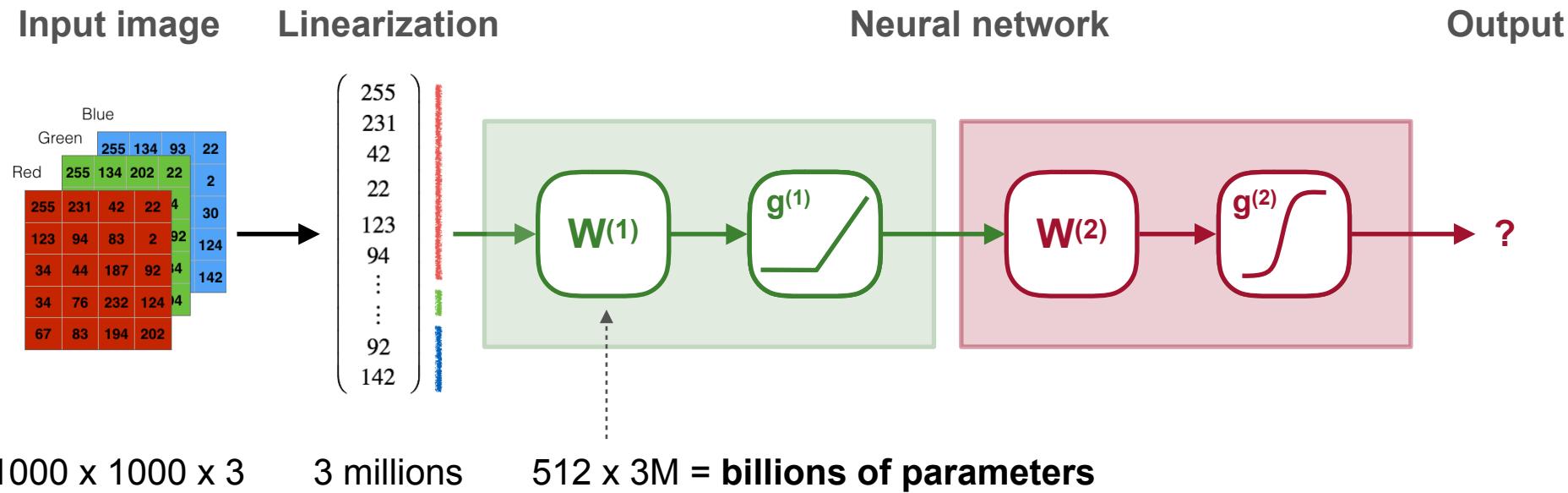
# Fully-connected layer

- Standard networks are made of **fully-connected layers**
  - *Each layer is a matrix multiplication*
  - *Each layer contains “(input dim. + 1) x hidden dim.” parameters*



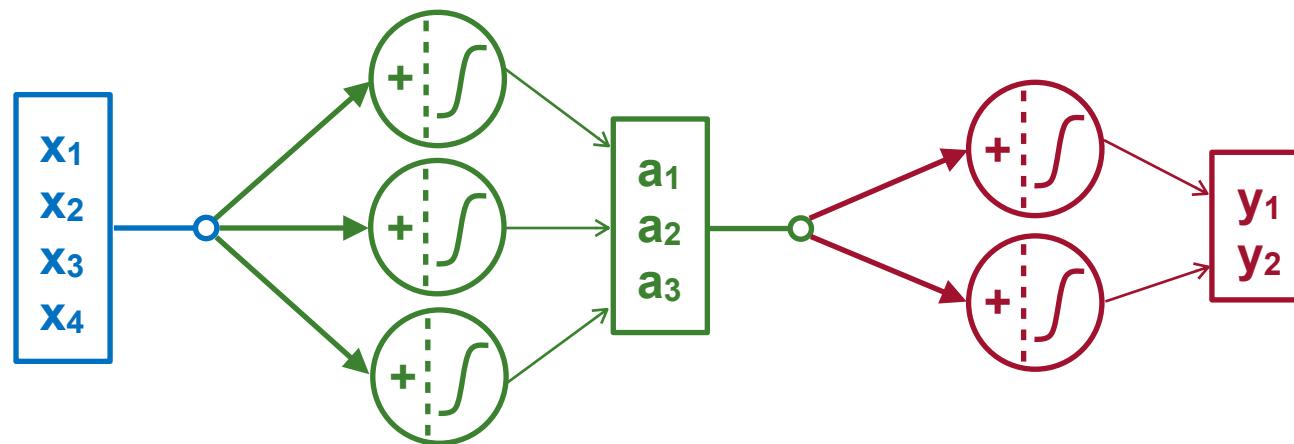
# Bottleneck of standard networks

- Fully-connected layers are **unsuitable for images**
  - Too many parameters to train*
  - Difficult to prevent overfitting*



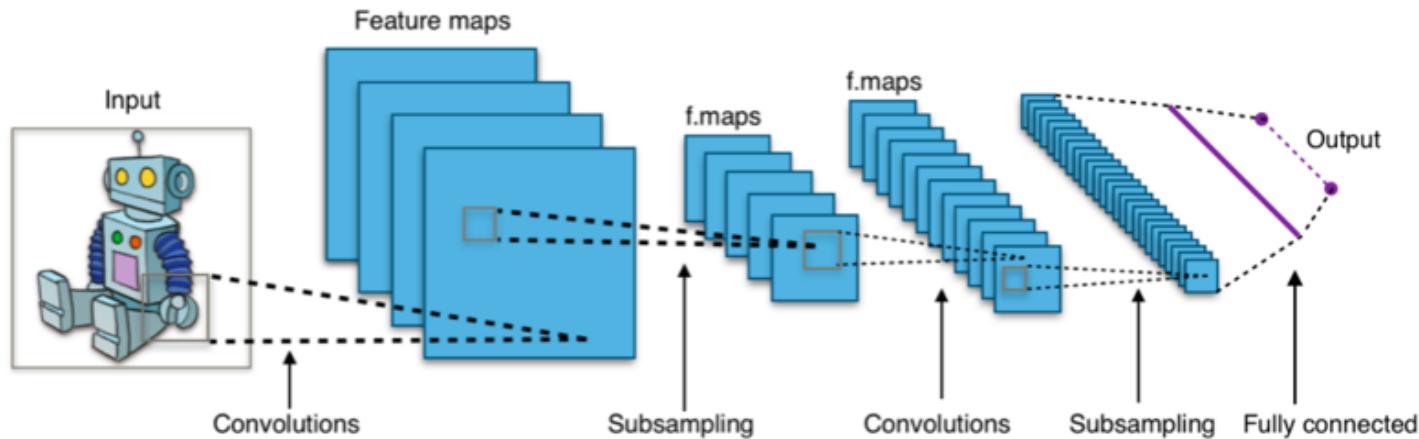
# Quiz

- Consider the two-layer network shown below.
  - How many parameters have the hidden layer?*
  - How many parameters have the output layer?*



# What we have seen so far...

- Standard networks are limited to **low-dimensional data**
  - *Too many parameters when the input is high-dimensional*
  - *How to deal with high-resolution images?*
- Solution for image inputs → **Convolutional networks**



# Convolution

---

Convolution in 1D/2D/3D  
Edge detection

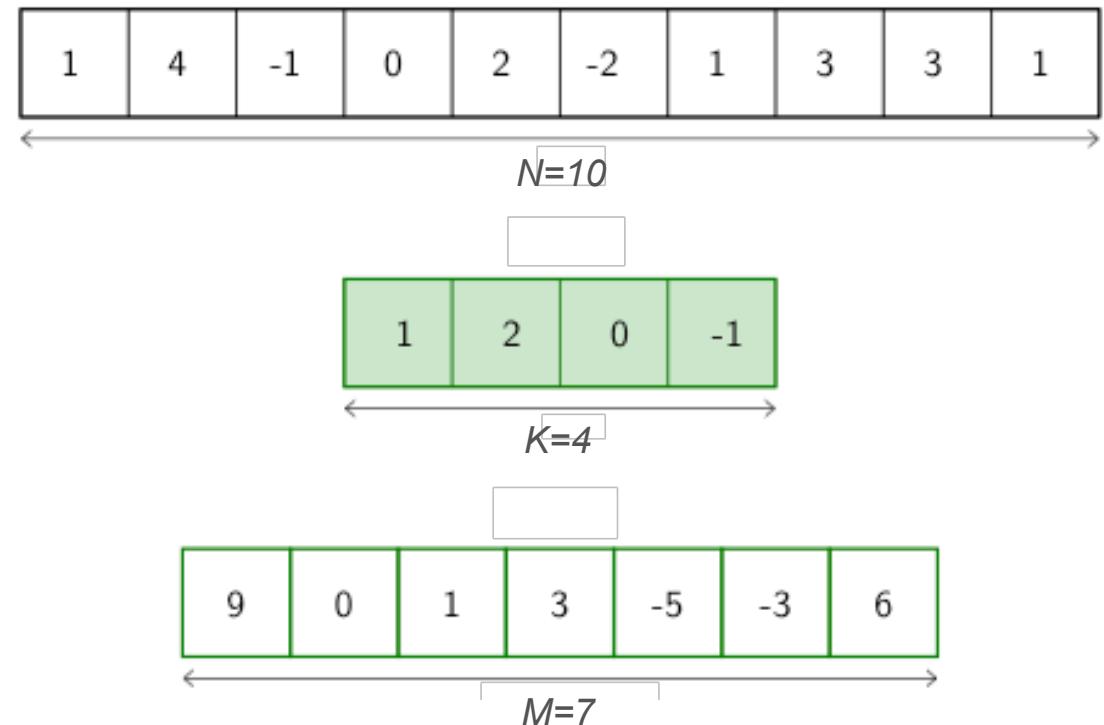
# Convolution 1D (1/2)

- Convolution with **vectors**

**Input**  
(vector of size  $N$ )

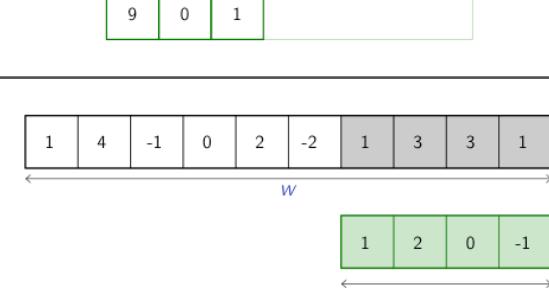
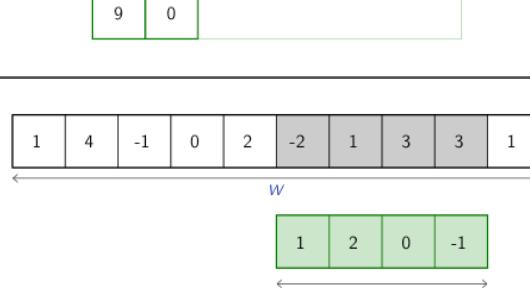
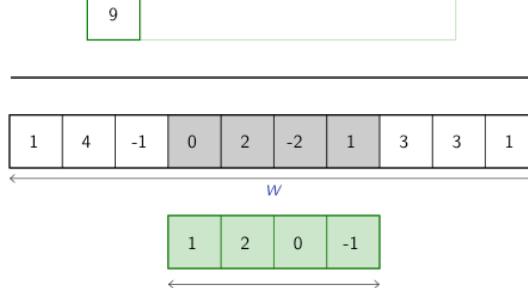
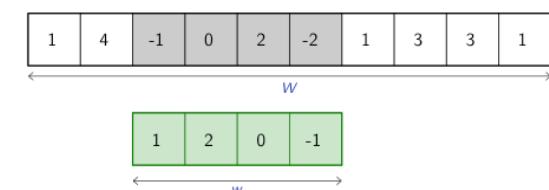
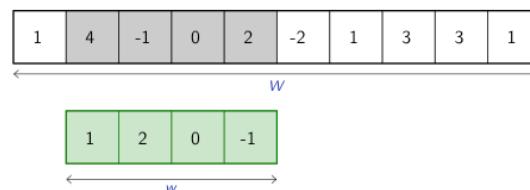
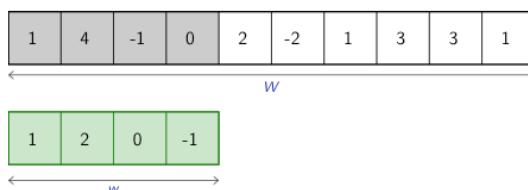
**Kernel**  
(vector of size  $K < N$ )

**Output**  
(vector of size  $N-K+1$ )



# Convolution 1D (2/2)

- Visual explanation of convolution in 1D
  - At each step, the kernel is multiplied to a chunk of the input...
  - ... and the resulting coefficients are summed



# Convolution 2D (1/2)

- Convolution with **images**

<i>Input</i> (size $N_1 \times N_2$ )					
3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

\*

*Kernel*  
(size  $K_1 \times K_2$ )

1	0	-1
1	0	-1
1	0	-1

=

-5	-4	0	8
-10	-2	2	3
0	-2	-4	-7
-3	-2	-3	-16

*Output*

(size  $N_1-K_1+1 \times N_2-K_2+1$ )

# Convolution 2D (2/2)

- Visual explanation of convolution in 2D

3 <sub>0</sub>	3 <sub>1</sub>	2 <sub>2</sub>	1	0
0 <sub>2</sub>	0 <sub>2</sub>	1 <sub>0</sub>	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3 <sub>0</sub>	2 <sub>1</sub>	1 <sub>2</sub>	0
0	0 <sub>2</sub>	1 <sub>2</sub>	3 <sub>0</sub>	1
3	1 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2 <sub>0</sub>	1 <sub>1</sub>	0 <sub>2</sub>
0	0	1 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>
3	1	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>
2	0	0	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>2</sub>	3	1
3 <sub>2</sub>	1 <sub>2</sub>	2 <sub>0</sub>	2	3
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	2	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0 <sub>0</sub>	1 <sub>1</sub>	3 <sub>2</sub>	1
3	1 <sub>2</sub>	2 <sub>2</sub>	2 <sub>0</sub>	3
2	0 <sub>0</sub>	0 <sub>1</sub>	2 <sub>2</sub>	2
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1 <sub>0</sub>	3 <sub>1</sub>	1 <sub>2</sub>
3	1	2 <sub>2</sub>	2 <sub>1</sub>	3 <sub>0</sub>
2	0	0 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>
2	0	0	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3 <sub>0</sub>	1 <sub>1</sub>	2 <sub>2</sub>	2	3
2 <sub>2</sub>	0 <sub>2</sub>	0 <sub>0</sub>	2	2
2 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	1

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3	1 <sub>0</sub>	2 <sub>1</sub>	2 <sub>2</sub>	3
2	0 <sub>2</sub>	0 <sub>2</sub>	2 <sub>0</sub>	2
2	0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1

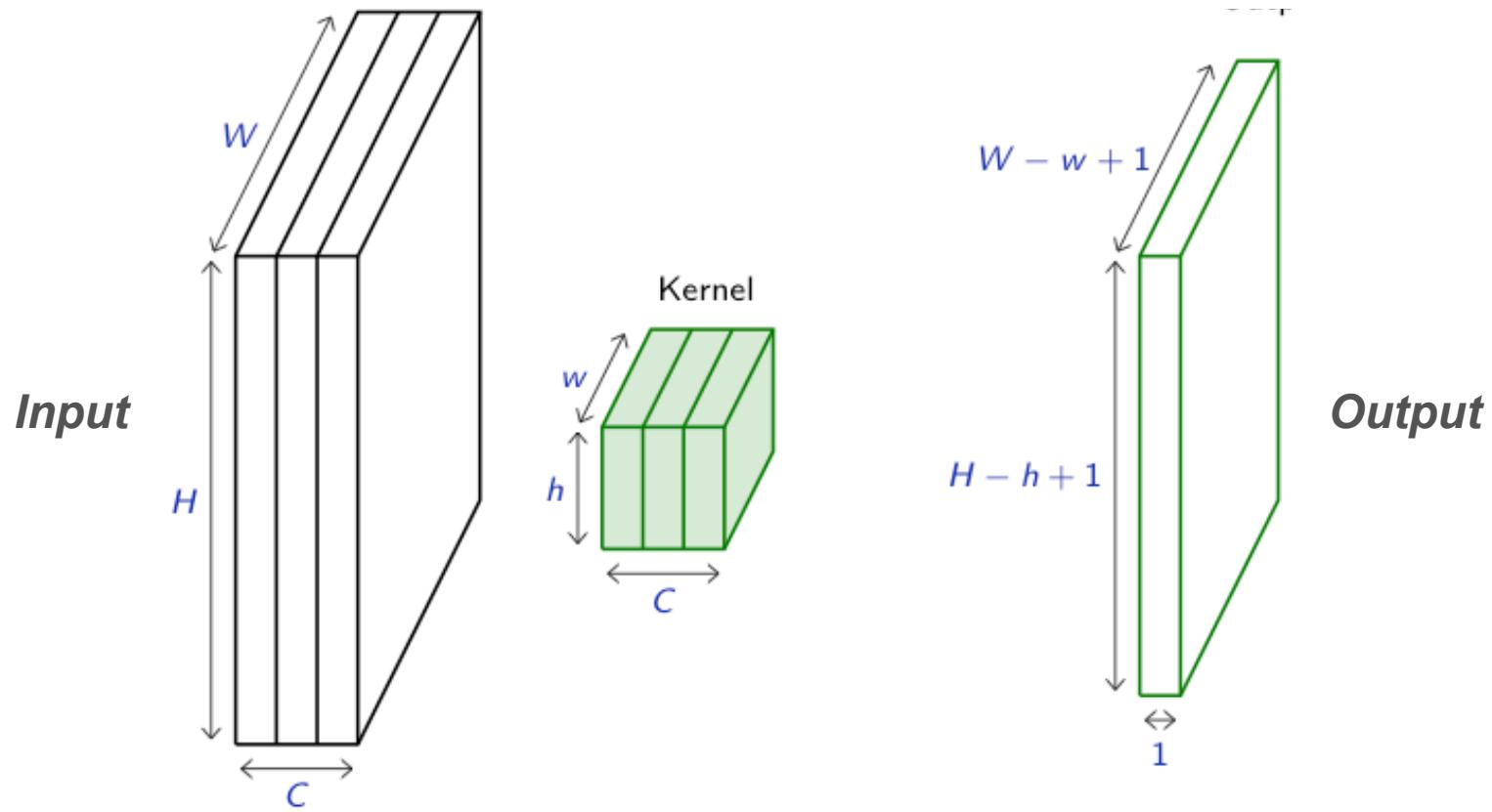
12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

3	3	2	1	0
0	0	1	3	1
3	1	2 <sub>0</sub>	2 <sub>1</sub>	3 <sub>2</sub>
2	0	0 <sub>2</sub>	2 <sub>2</sub>	2 <sub>0</sub>
2	0	0 <sub>0</sub>	0 <sub>1</sub>	1 <sub>2</sub>

12.0	12.0	17.0
10.0	17.0	19.0
9.0	6.0	14.0

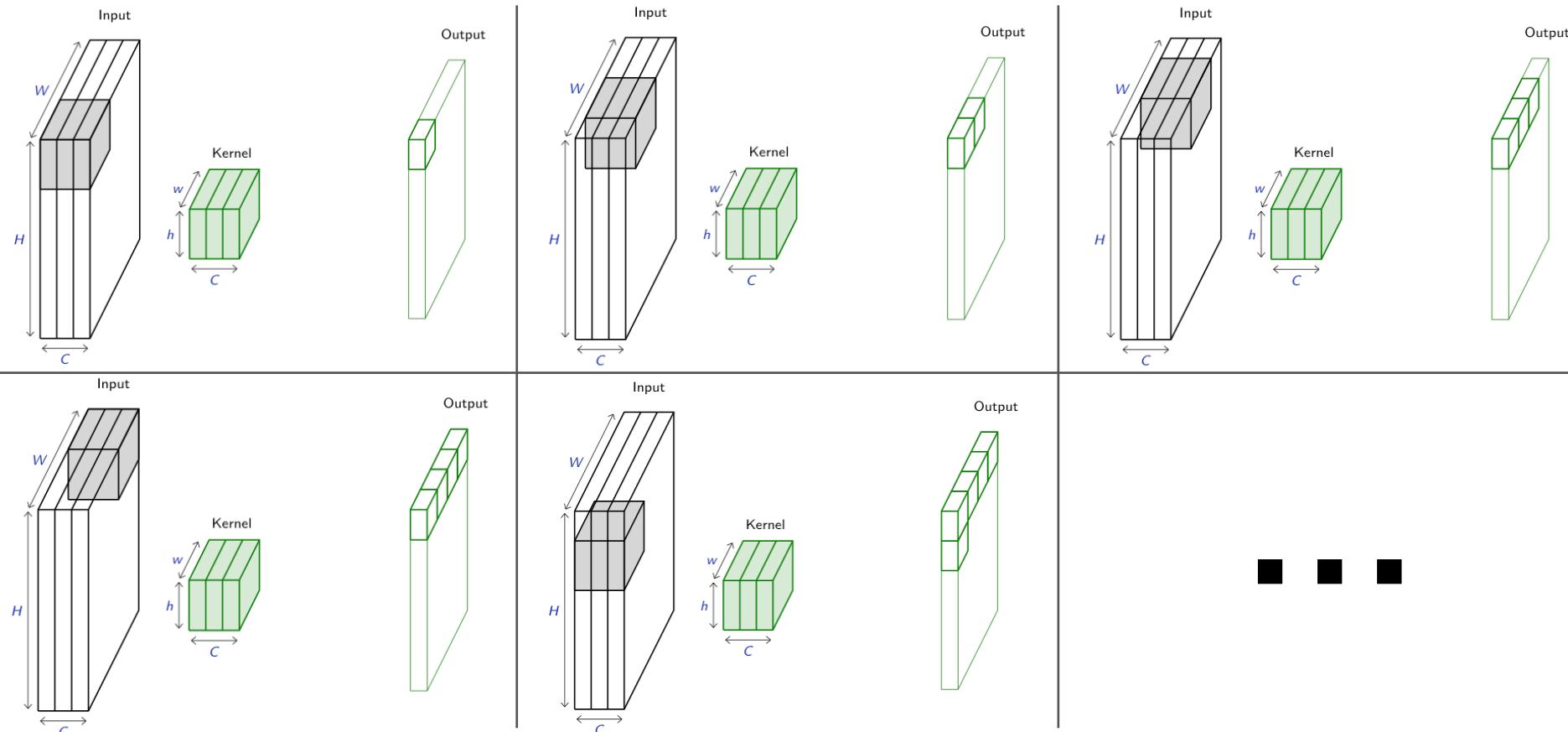
# Convolution 3D (1/2)

- Convolution with **multi-channel images**



# Convolution 3D (2/2)

- Visual explanation of convolution in 3D



# Edge detection (1/3)

- Vertical edge detection

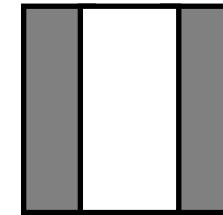
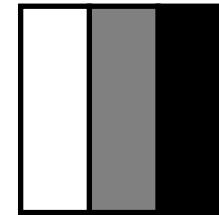
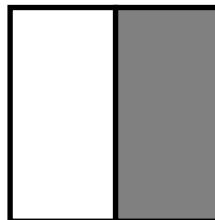
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0

\*

1	0	-1
1	0	-1
1	0	-1

=

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0



# Edge detection (2/3)

- Horizontal edge detection

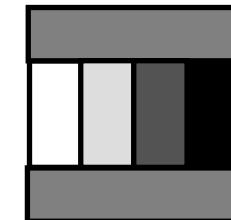
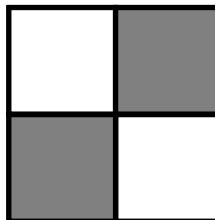
10	10	10	0	0	0
10	10	10	0	0	0
10	10	10	0	0	0
0	0	0	10	10	10
0	0	0	10	10	10
0	0	0	10	10	10

\*

1	1	1
0	0	0
-1	-1	-1

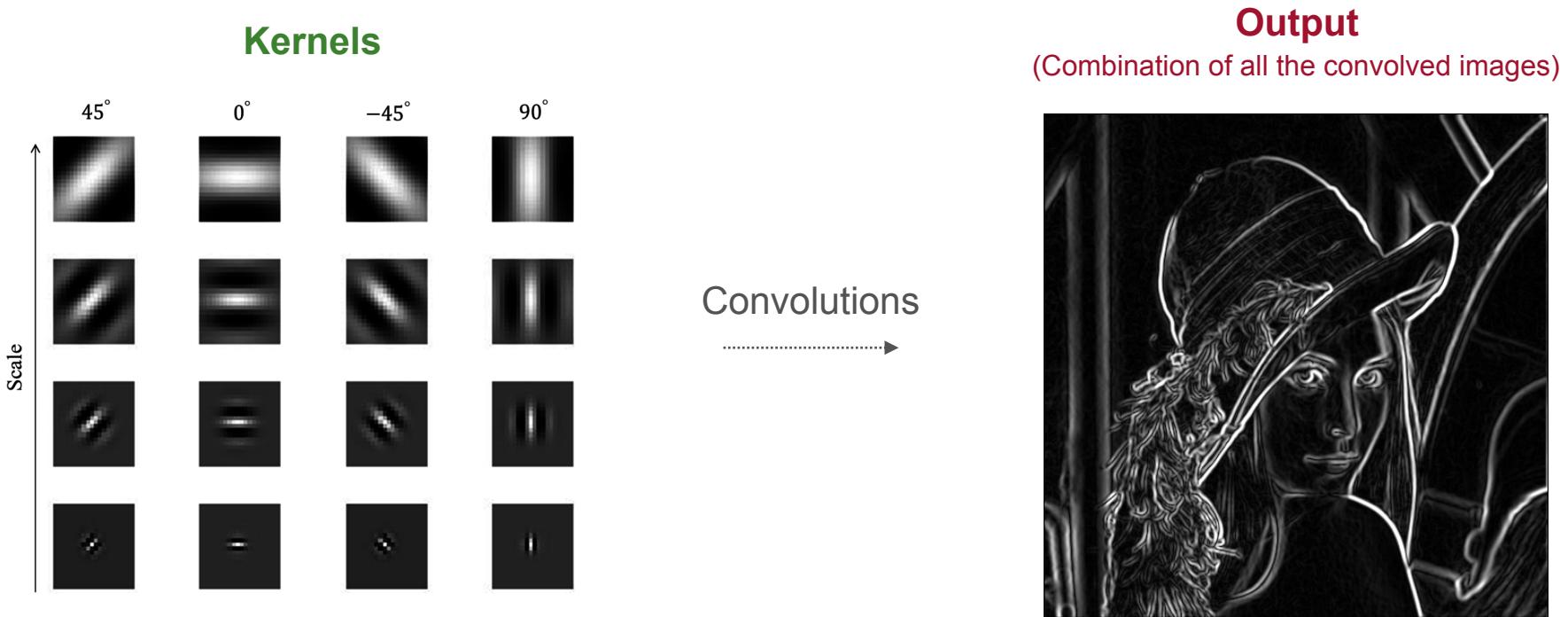
=

0	0	0	0
30	10	-10	-30
30	10	-10	-30
0	0	0	0



# Edge detection (3/3)

- **Putting all together**
  - *A different kernel for each orientation and scale*



# Quiz

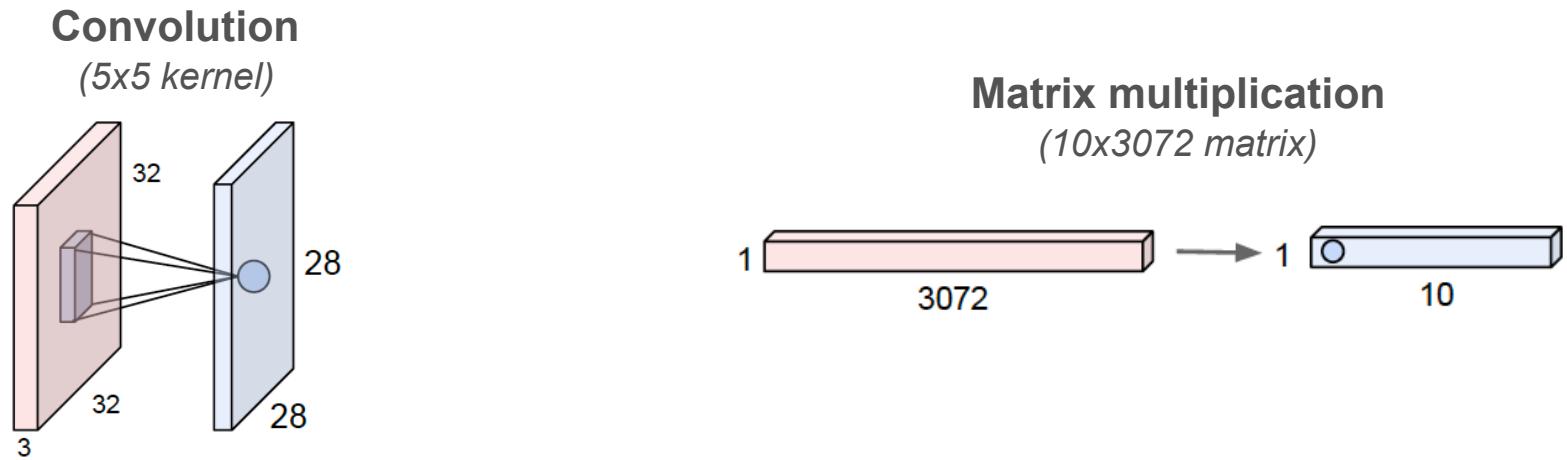
- 1) What is the output size of the following convolutions?
  - ❑ *10x12 matrix convolved by 5x5 kernel*
  - ❑ *15x15x10 volume convolved by 3x3x3 kernel*
- 2) Compute the following convolution.

3	0	1	2	7	4
1	5	8	9	3	1
2	7	2	5	1	3
0	1	3	1	7	8
4	2	1	6	2	8
2	4	5	2	3	9

$$\begin{matrix} & * & \begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{matrix} & = & \begin{matrix} & & & \\ & & & \\ & & & \\ & & & \\ & & & \end{matrix} \end{matrix}$$

# What we have seen so far...

- **Convolution** → spatially-invariant operation
  - *The same weights are used for computing the output coefficients*
- **Matrix multiplication** → general operation
  - *Different weights are used for computing the output coefficients*



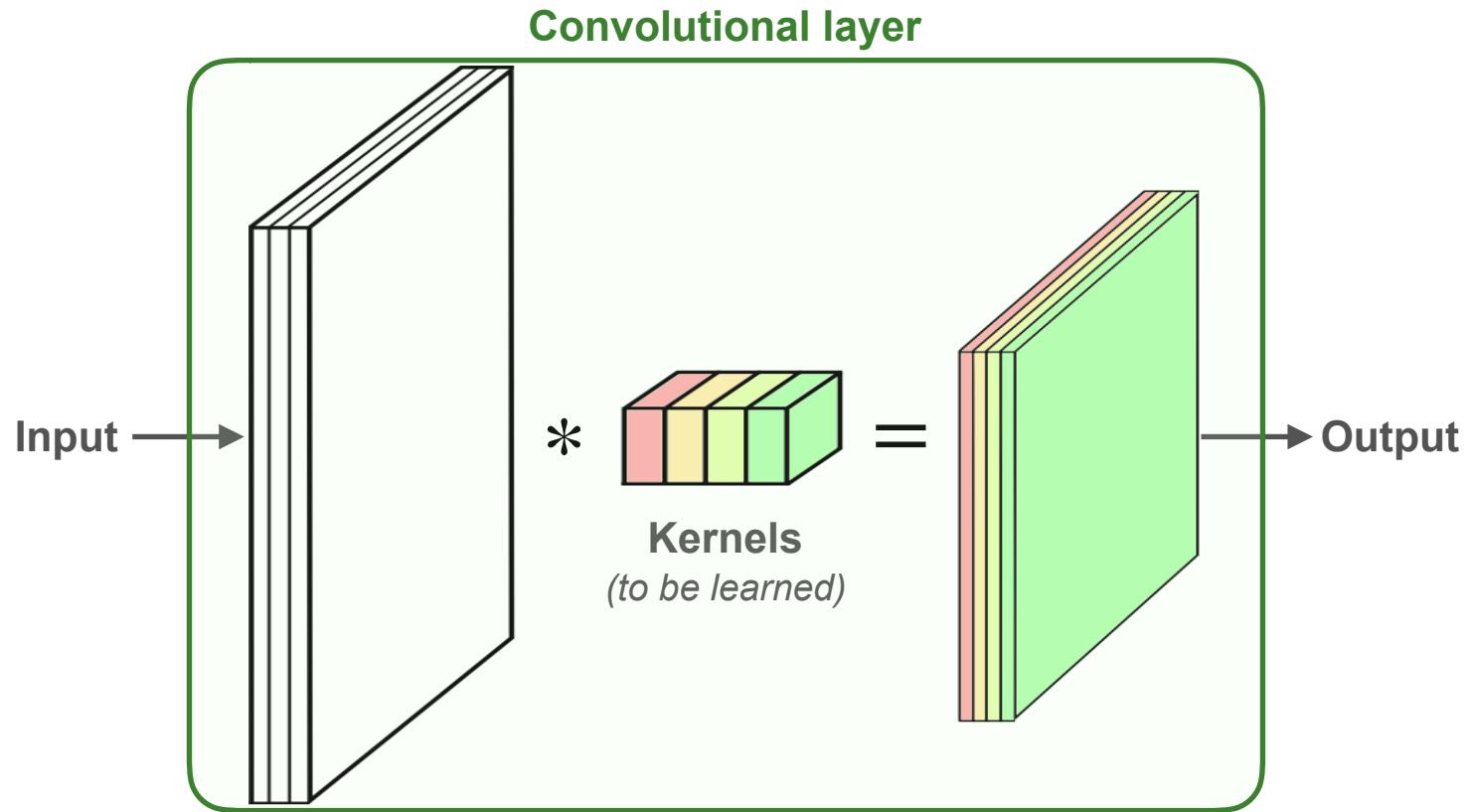
# Convolutional layer

---

- Multiple kernels
- Hyper-parameters
- Padding & Stride

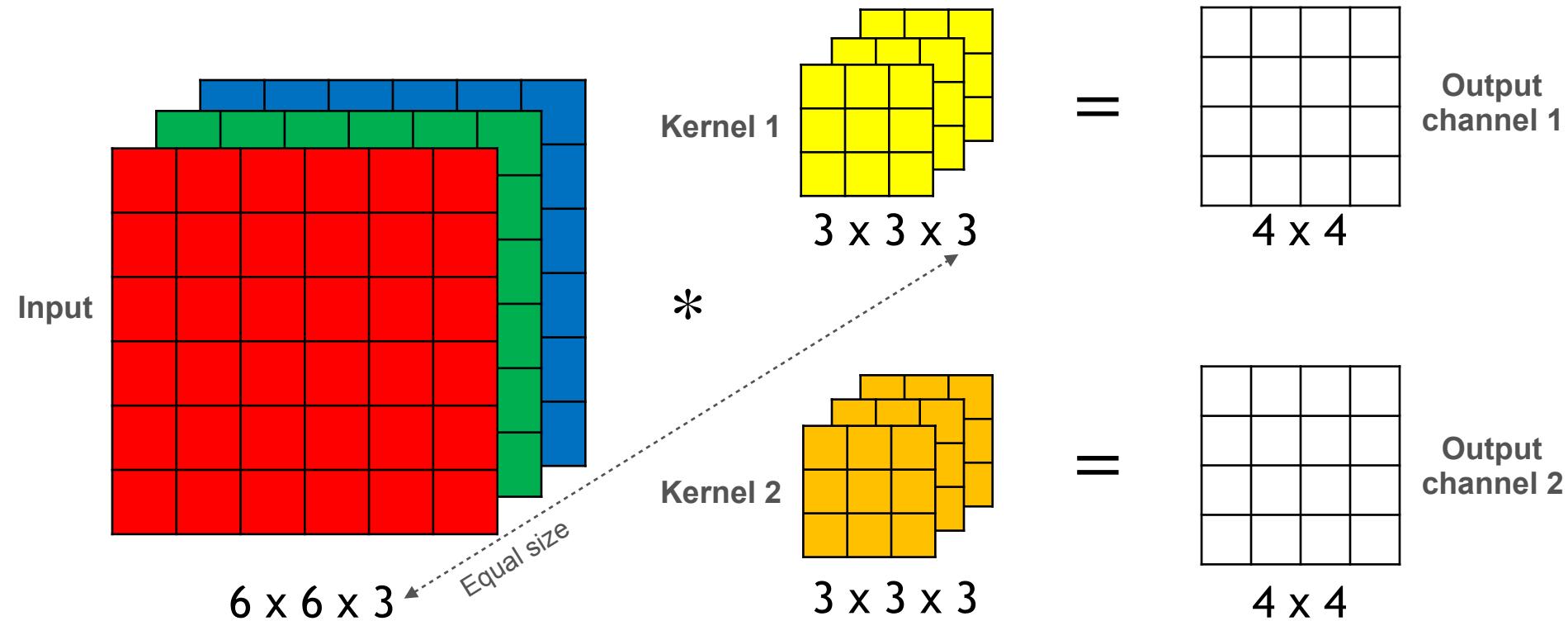
# Multiple kernels (1/2)

- The input is convolved with **multiple kernels**



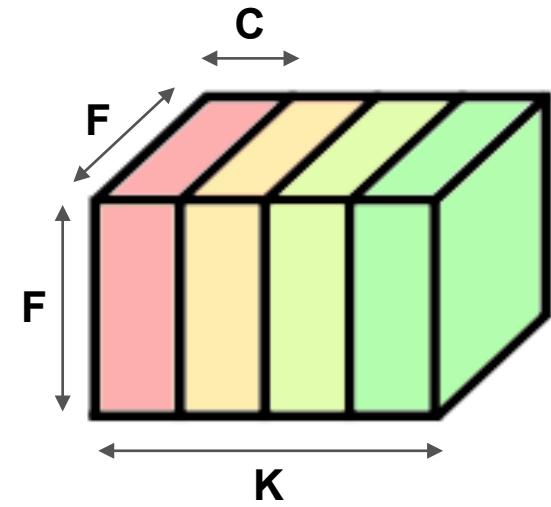
# Multiple kernels (2/2)

- Each convolution produces a channel for the output
  - The **kernel depth** must be equal to the number of **input channels***



# Hyper-parameters

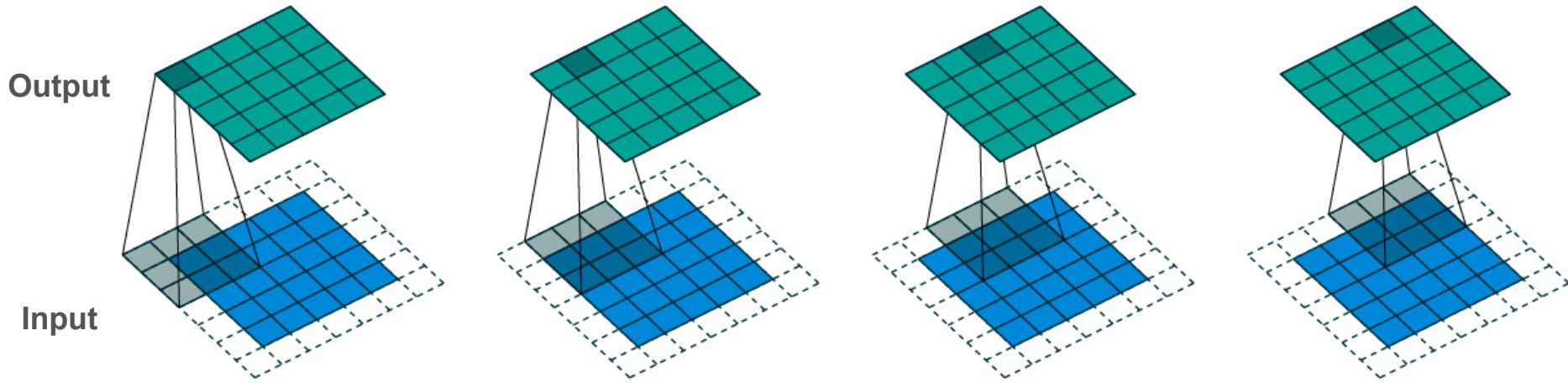
- Main hyper-parameters
  - $F \rightarrow$  *Spatial size of kernels*
  - $K \rightarrow$  *Number of kernels*
- Other hyper-parameters
  - $P \rightarrow$  *Padding*
  - $S \rightarrow$  *Stride*
- Fixed value
  - $C \rightarrow$  *Kernel depth* (*equal to the number of input channels*)



# Padding & stride (1/3)

- The input can be **padded** with zero rows/columns
  - The output size is extended by  $P > 0$*

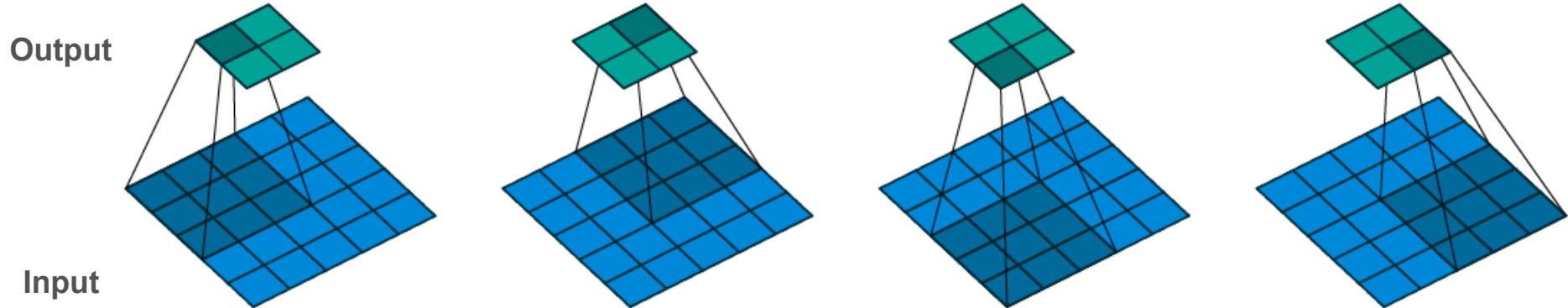
Example with  $P=1$



# Padding & stride (2/3)

- The output can be **downsampled** along the rows/columns
  - The output size is reduced by a factor  $S > 1$*

Example with  $S=2$



# Padding & stride (3/3)

- Example with **padding (P=1)** and **stride (S=2)**

0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	0	0	0	0
0 <sub>2</sub>	3 <sub>2</sub>	3 <sub>0</sub>	2	1	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3	1	0	0
0	3	1	2	2	3	0	0
0	2	0	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1	0	0
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>	2	2	3	0	0
0	2 <sub>1</sub>	0 <sub>2</sub>	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1 <sub>1</sub>	0 <sub>2</sub>	0 <sub>0</sub>
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>1</sub>	2	2	3 <sub>0</sub>	0 <sub>0</sub>	0 <sub>2</sub>
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1	0	0
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>	2	2	3	0	0
0	2 <sub>1</sub>	0 <sub>2</sub>	0	2	2	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1 <sub>1</sub>	0 <sub>2</sub>	0 <sub>0</sub>
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>1</sub>	2	2	3 <sub>0</sub>	0 <sub>0</sub>	0 <sub>2</sub>
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1 <sub>1</sub>	0 <sub>2</sub>	0 <sub>0</sub>
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>1</sub>	2	2	3 <sub>0</sub>	0 <sub>0</sub>	0 <sub>2</sub>
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1	0	0
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>0</sub>	2	2	3	0	0
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1 <sub>1</sub>	0 <sub>2</sub>	0 <sub>0</sub>
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>1</sub>	2	2	3 <sub>0</sub>	0 <sub>0</sub>	0 <sub>2</sub>
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

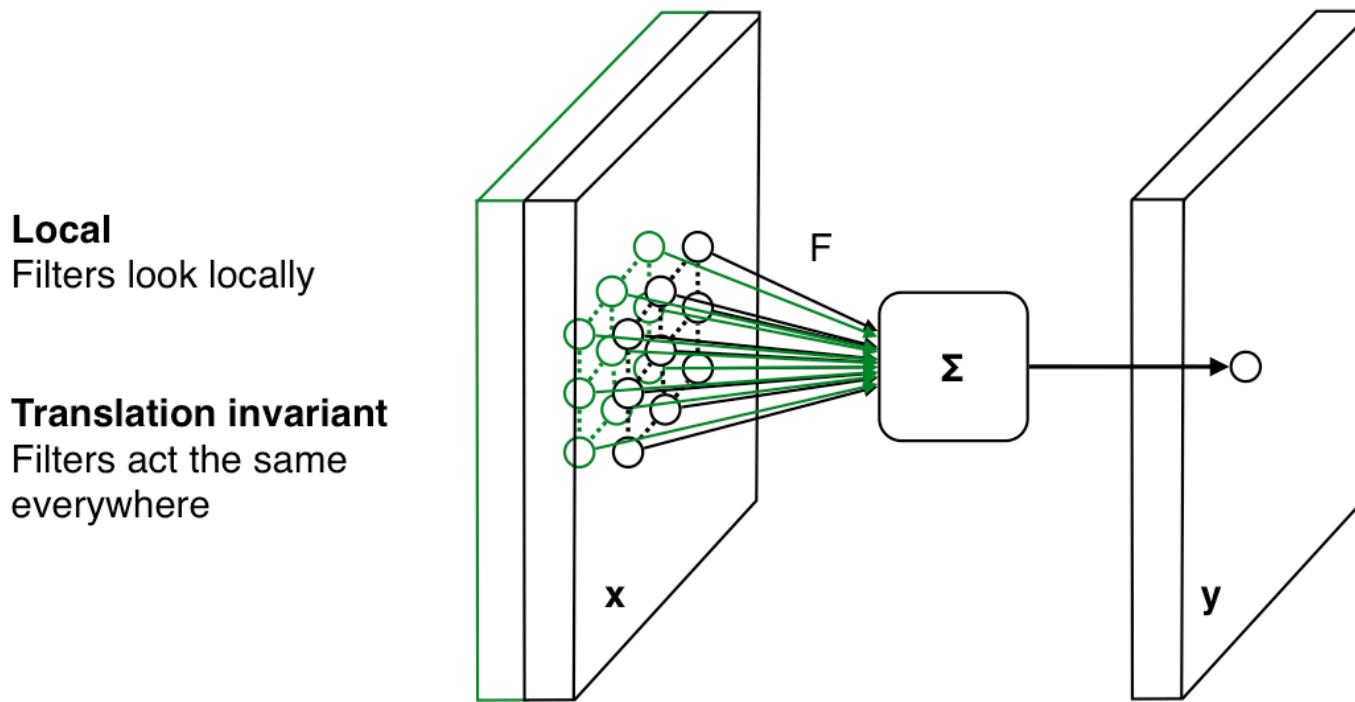
6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

0	0	0	0	0	0	0	0
0	3	3 <sub>2</sub>	2 <sub>1</sub>	1 <sub>0</sub>	0	0	0
0 <sub>0</sub>	0 <sub>1</sub>	0 <sub>2</sub>	1	3 <sub>1</sub>	1 <sub>1</sub>	0 <sub>2</sub>	0 <sub>0</sub>
0 <sub>2</sub>	3 <sub>2</sub>	1 <sub>1</sub>	2	2	3 <sub>0</sub>	0 <sub>0</sub>	0 <sub>2</sub>
0	3 <sub>1</sub>	1 <sub>2</sub>	2 <sub>2</sub>	2	3	0	0
0	2 <sub>0</sub>	0 <sub>1</sub>	2 <sub>1</sub>	2	0	0	0
0	2	0	0	0	1	0	0
0	0	0	0	0	0	0	0

6.0	17.0	3.0
8.0	17.0	13.0
6.0	4.0	4.0

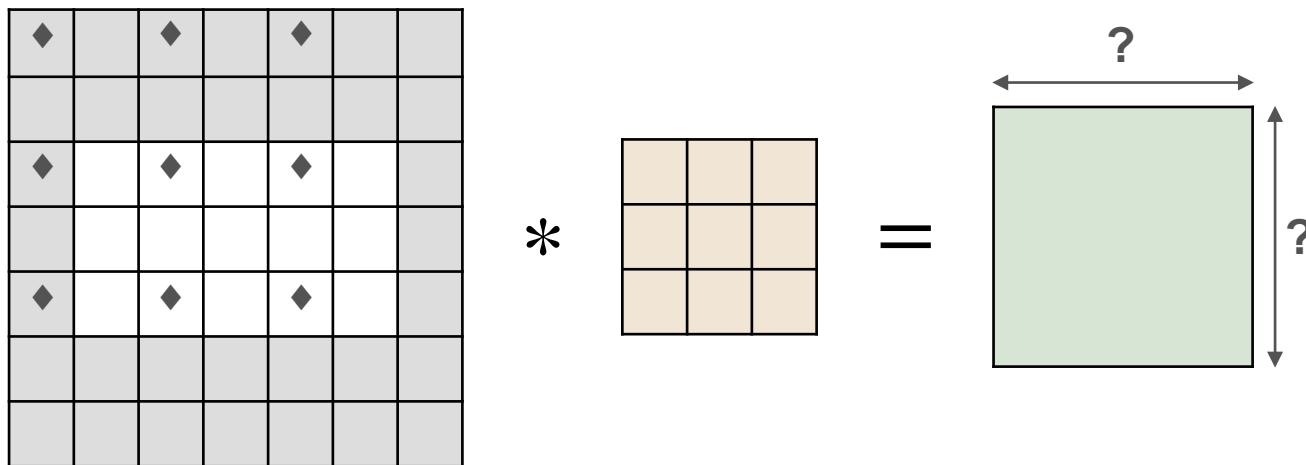
# Neural interpretation

- A convolution is a neuron with limited reception field
  - *The field limit is defined by the kernel size*
  - *The same neuron is "fired" over multiple areas from the input.*



# Quiz

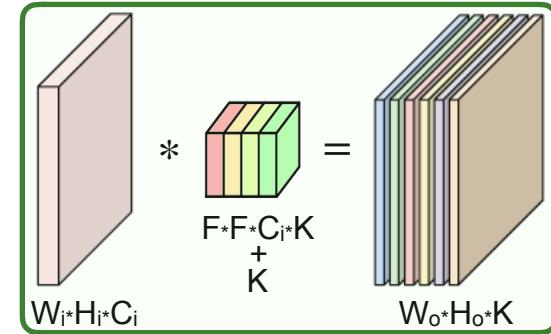
- Compute the output shape of a convolutional layer with
  - Input** →  $3 \times 5 \times C$
  - Kernel** →  $3 \times 3 \times C \times 2$
  - Padding** →  $2 \times 1$
  - Stride** →  $2 \times 2$



# What we have seen so far...

## ▪ Convolutional layer

- The input is a volume of size  $W_i \times H_i \times C_i$
- Four hyper-parameters are required
  - ★  $K \rightarrow$  Number of kernels
  - ★  $F \rightarrow$  Spatial size
  - ★  $S \rightarrow$  Stride
  - ★  $P \rightarrow$  Padding
- The output is a volume of size  $W_o \times H_o \times C_o$ 
  - ★  $W_o = (W_i - F + 2P) / S + 1$
  - ★  $H_o = (H_i - F + 2P) / S + 1$
  - ★  $C_o = K$
- The number of parameters to be learned is  $(F \times F \times C_i + 1) \times K$



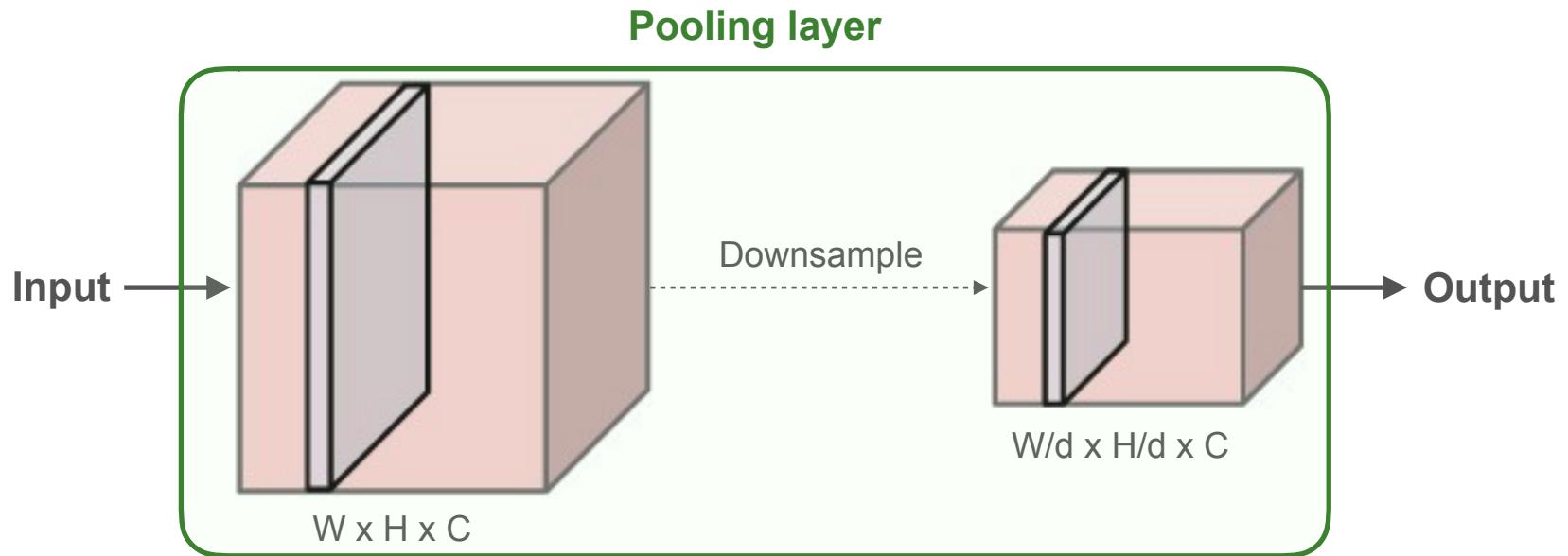
# Pooling layer

---

Downsampling  
Max-pooling  
Translation invariance

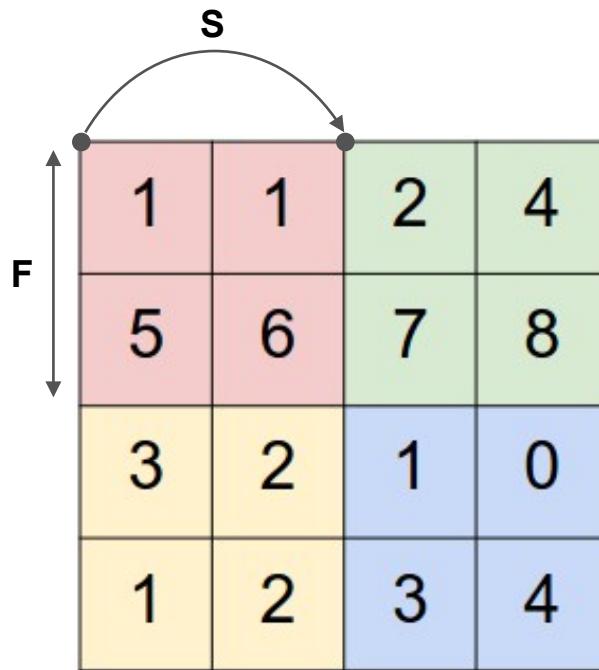
# Downsampling

- The **pooling layer** reduces the spatial size of the input
  - It operates over each channel map independently*
  - It has no parameter to be learned*



# Max-pooling (1 / 2)

- The most-common variant is **max-pooling**
  - It computes the max over a sliding window*
  - Hyper-parameters** → window size ( $F$ ) & stride ( $S$ )



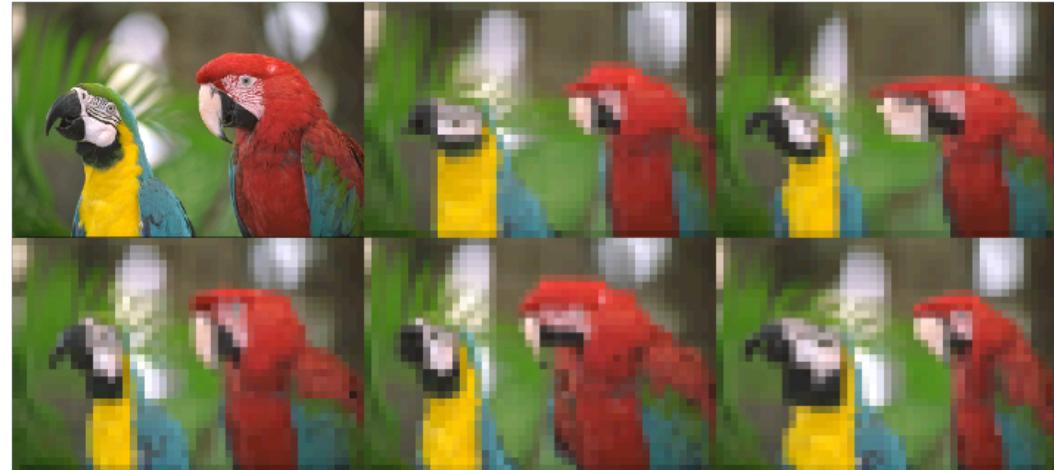
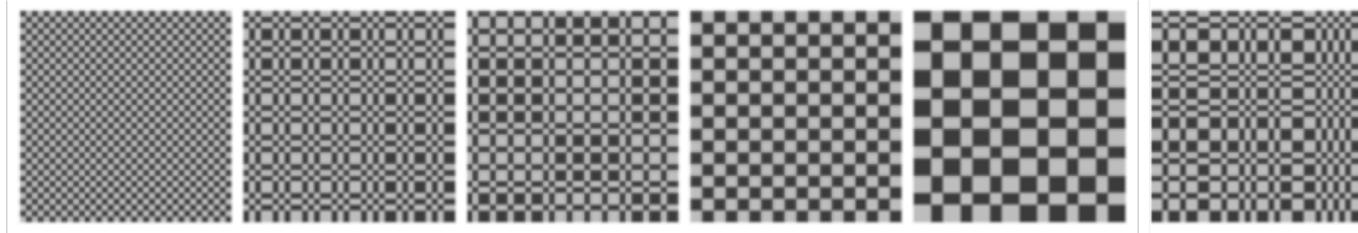
- $F = 2 \rightarrow$  window size  
•  $S = 2 \rightarrow$  stride

The resulting 2x2 output matrix is:

6	8
3	4

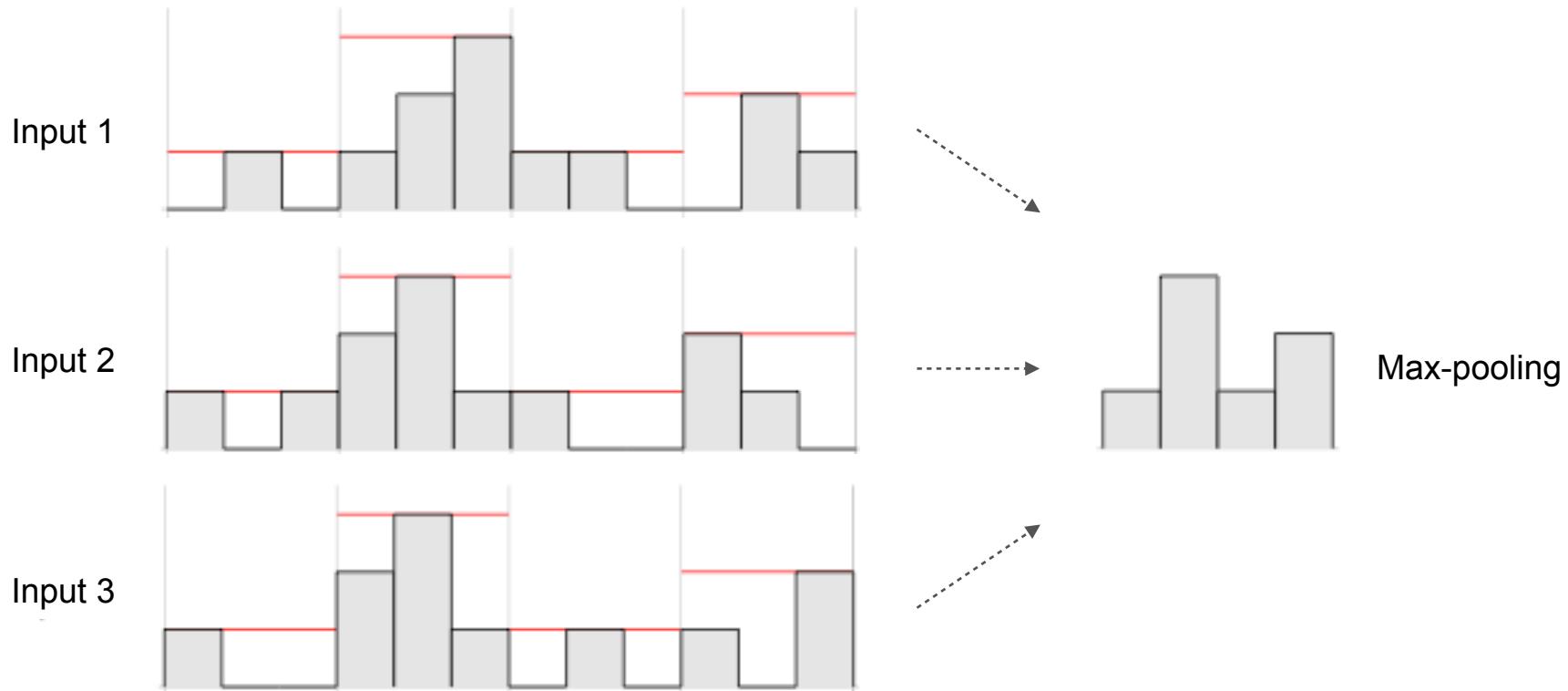
# Max-pooling (2/2)

- **Stochastic pooling** → Random pooling mask at each pass



# Translation invariance

- Pooling and convolution are **translation invariant**
  - *A spatial shift in the input doesn't change the output*

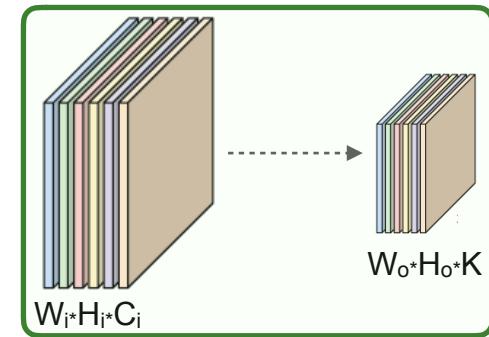


# Quiz

- 1) Because pooling layers do not have parameters, they do not affect the gradient calculation.
  - A. *True*
  - B. *False*
- 2) You have an input volume that is  **$32 \times 32 \times 16$** , and apply max-pooling with a **stride of 2** and a **window size of 2**. What is the output volume?
  - A.  $16 \times 16 \times 8$
  - B.  $32 \times 32 \times 8$
  - C.  $16 \times 16 \times 16$
  - D.  $15 \times 15 \times 16$

# What we have seen so far...

- **Pooling layer**
  - The input is a volume of size  $W_i \times H_i \times C_i$
  - Two hyper-parameters are required
    - ★  $F \rightarrow$  Window size
    - ★  $S \rightarrow$  Stride
  - The output is a volume of size  $W_o \times H_o \times C_o$ 
    - ★  $W_o = (W_i - F) / S + 1$
    - ★  $H_o = (H_i - F) / S + 1$
    - ★  $C_o = C_i$
  - No parameters to be learned
- **Note →** Pooling can be replaced with a stride in convolutional layers!



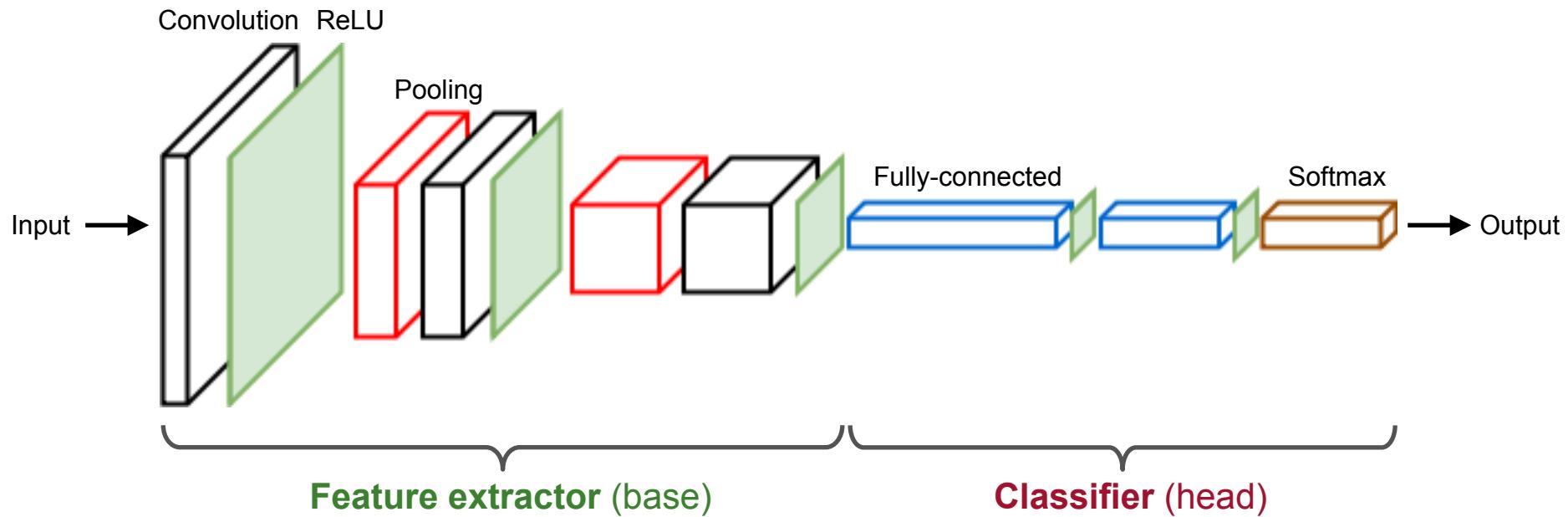
# Convolutional network

---

- Building blocks
- Classic architectures
- Modern architectures

# Building blocks (1 / 2)

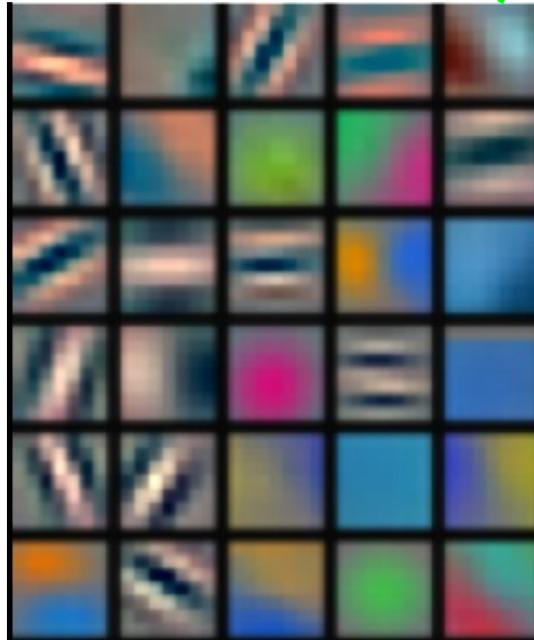
- Architecture of a **convolutional neural network**
  - Base** → Convolution + pooling
  - Head** → Fully-connected + output layer



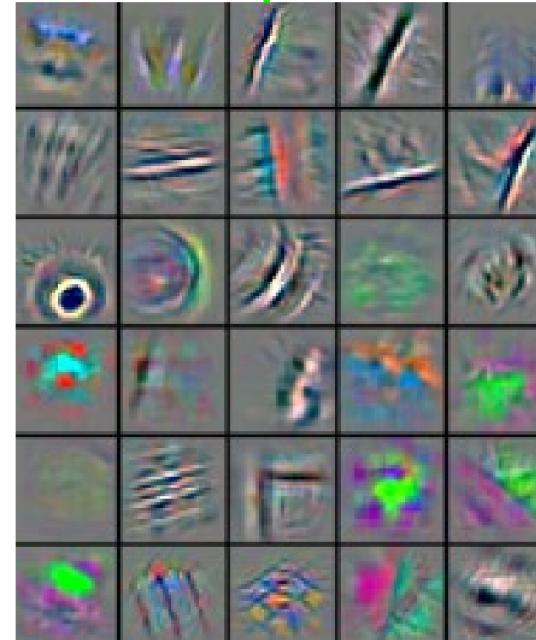
# Building blocks (2/2)

- ConvNets learn spatial **hierarchies of patterns**

**Low-level patterns**  
(conv. layers in the front)



**Mid-level patterns**  
(conv. layers in the middle)

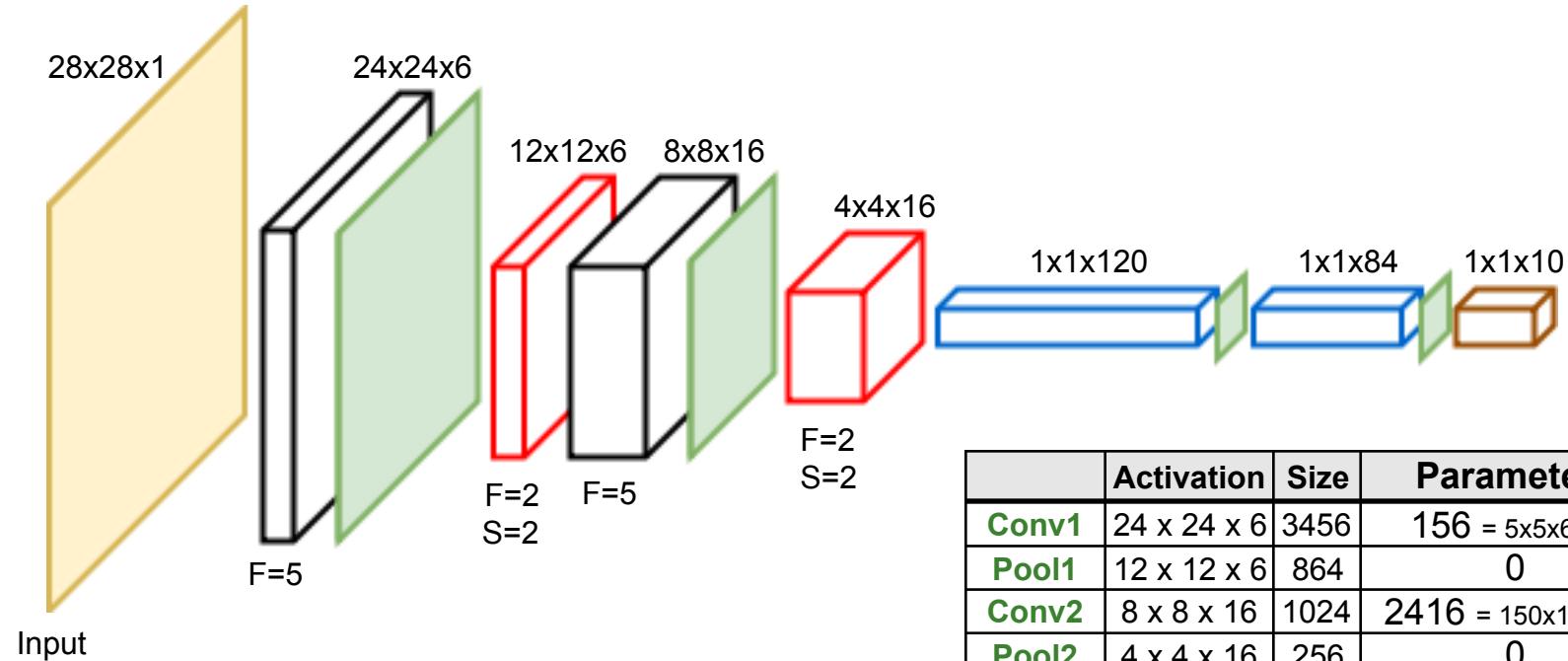


**High-level patterns**  
(conv. layers in the end)



# Classic architectures (1 / 3)

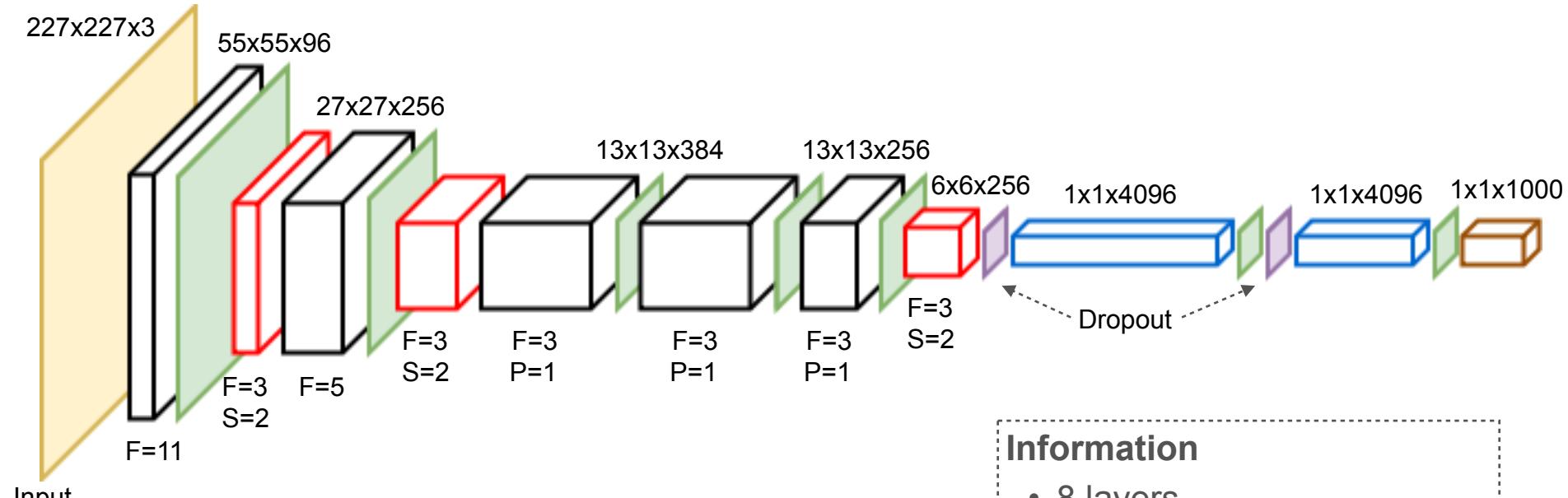
- **LeNet-5 (1998) → 0.05 millions of parameters**



	Activation	Size	Parameters
<b>Conv1</b>	$24 \times 24 \times 6$	3456	$156 = 5 \times 5 \times 6 + 6$
<b>Pool1</b>	$12 \times 12 \times 6$	864	0
<b>Conv2</b>	$8 \times 8 \times 16$	1024	$2416 = 150 \times 16 + 16$
<b>Pool2</b>	$4 \times 4 \times 16$	256	0
<b>Flatten</b>	$1 \times 1 \times 256$	256	0
<b>FC1</b>	$1 \times 1 \times 120$	120	$30840 = (256+1) \times 120$
<b>FC2</b>	$1 \times 1 \times 84$	84	$10164 = (120+1) \times 84$
<b>Softmax</b>	$1 \times 1 \times 10$	10	$850 = (84+1) \times 10$

# Classic architectures (2/3)

- **AlexNet (2012)**

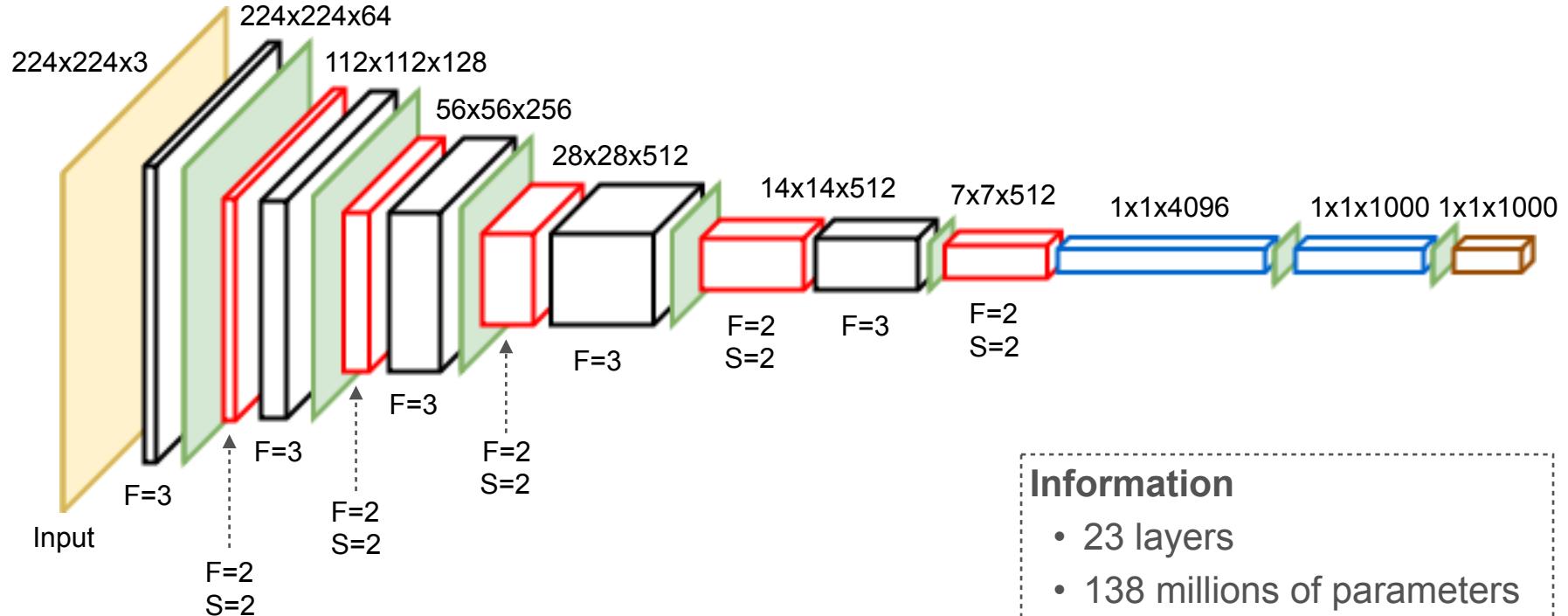


**Information**

- 8 layers
- 60 millions of parameters
- 84% accuracy on ImageNet

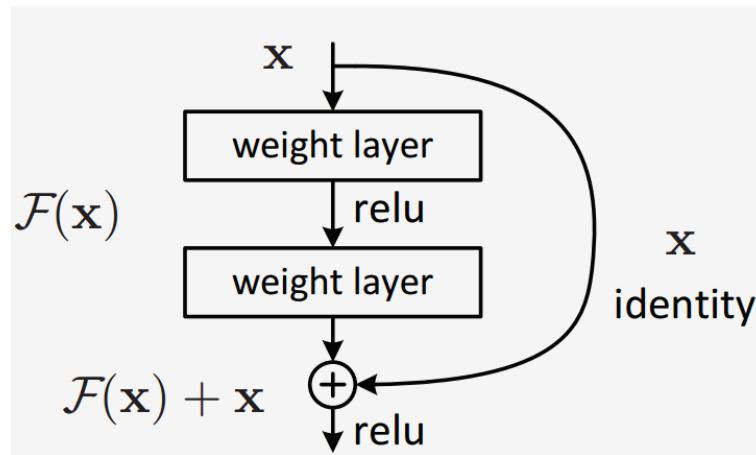
# Classic architectures (3/3)

- **VGG-16 (2015)**



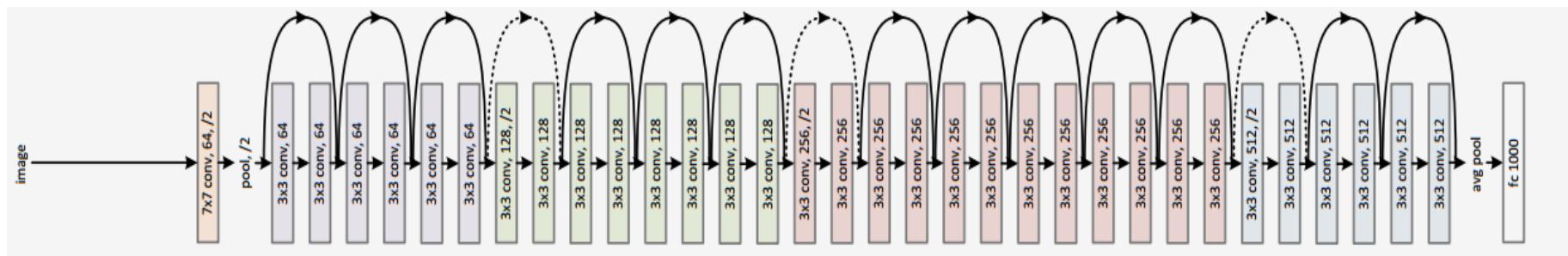
# Modern architectures (1/2)

- **Residual network (2016)**



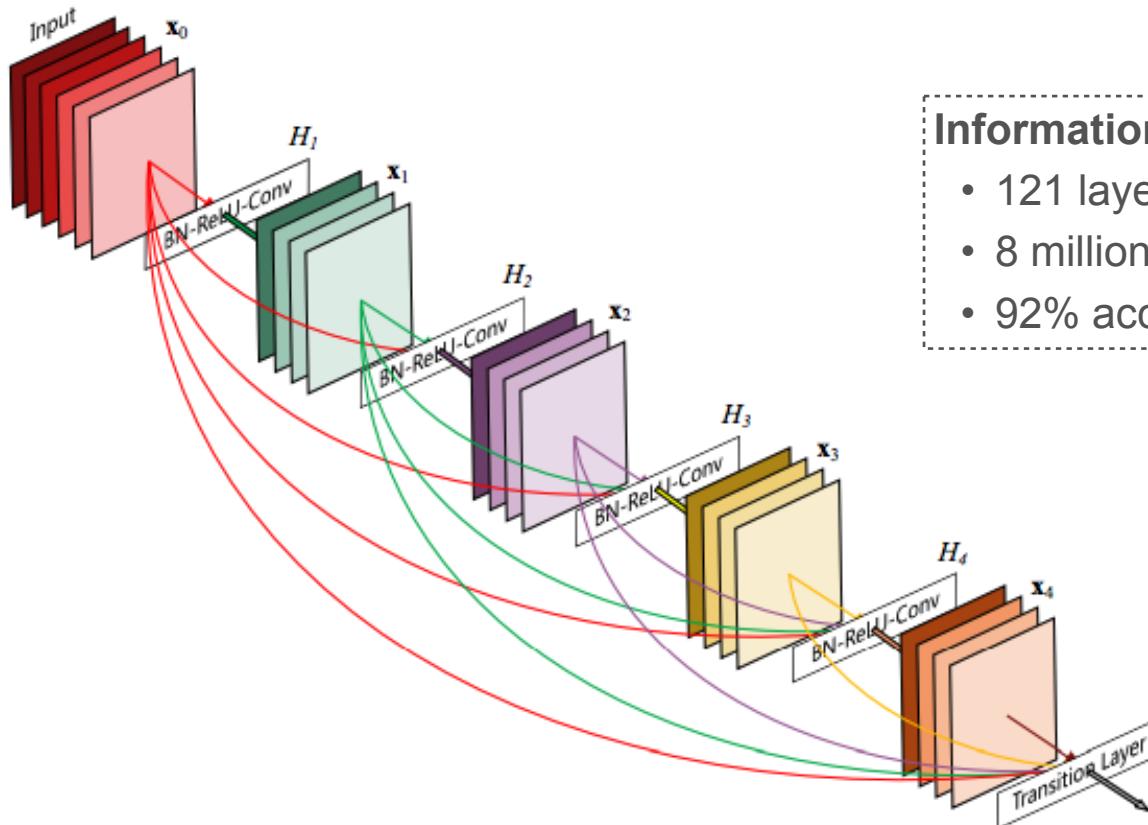
## Information

- 168 layers
- 25 millions of parameters
- 93% accuracy on ImageNet



# Modern architectures (2/2)

- Dense network (2017)



**Information**

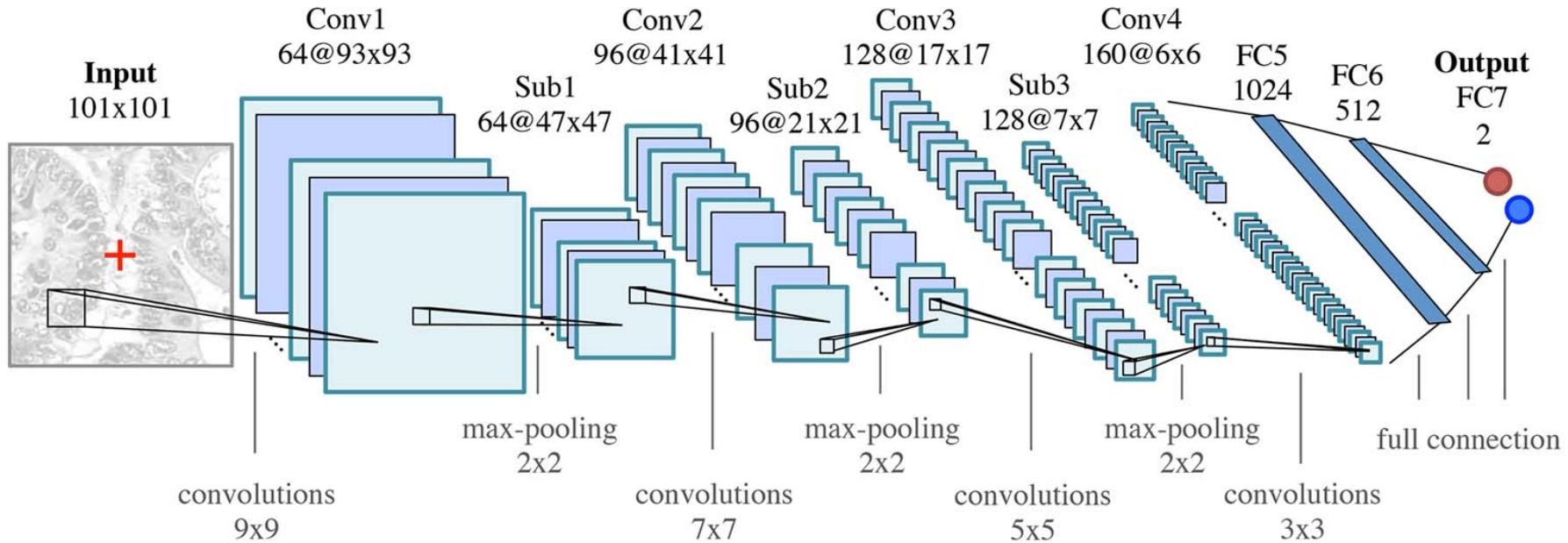
- 121 layers
- 8 millions of parameters
- 92% accuracy on ImageNet

# Quiz

- 1) Which of the following do you typically see as you move to deeper layers in a ConvNet? ( $W \times H \times C$  is the layers' output size)
  - A.  $W$  and  $H$  decreases, while  $C$  also decreases
  - B.  $W$  and  $H$  increases, while  $C$  decreases
  - C.  $W$  and  $H$  decreases, while  $C$  increases
  - D.  $W$  and  $H$  increases, while  $C$  also increases
- 2) Which of the following do you typically see in a ConvNet?
  - A. Multiple CONV layers follows by a POOL layer
  - B. Multiple POOL layers follows by a CONV layer
  - C. FC layers in the **last** few layers
  - D. FC layers in the **first** few layers

# What we have seen so far...

- ConvNets make the assumption that the **inputs are images**
  - **New layers → Convolution & Pooling**
  - **Architecture → Feature extractor + Classifier**



# Practical advice

- **Use whatever works best on ImageNet**
  - If you're feeling a bit of a fatigue in thinking about the architectural decisions, you'll be pleased to know that in 90% or more of applications you should not have to worry about these. Instead of rolling your own architecture for a problem, you should look at whatever architecture currently works best on ImageNet, download a pretrained model and fine-tune it on your data. You should rarely ever have to train a CNN from scratch or design one from scratch.
- **Prefer a stack of small CONV layers to one large CONV layer**
  - Suppose that you stack three 3x3 CONV layers on top of each other (with non-linearities in between, of course). In this arrangement, each neuron on the first CONV layer has a 3x3 view of the input volume. A neuron on the second CONV layer has a 3x3 view of the first CONV layer, and hence by extension a 5x5 view of the input volume. Similarly, a neuron on the third CONV layer has a 3x3 view of the 2nd CONV layer, and hence a 7x7 view of the input volume. Suppose that instead of these three layers of 3x3 CONV, we only wanted to use a single CONV layer with 7x7 receptive fields. These neurons would have a receptive field size of the input volume that is identical in spatial extent (7x7). However, the neurons would be computing a linear function over the input, while the three stacks of CONV layers contain non-linearities that make their features more expressive.