

# Foundations of Machine Learning (ECE 5984)

## - Convolutional Neural Networks -

Eunbyung Park

Assistant Professor

School of Electronic and Electrical Engineering

[Eunbyung Park \(silverbottlep.github.io\)](https://silverbottlep.github.io)

# Convolution

# 1D Convolution

- Convolution is a mathematical operation on two functions  $(f, g)$  that produces a third function  $f * g$

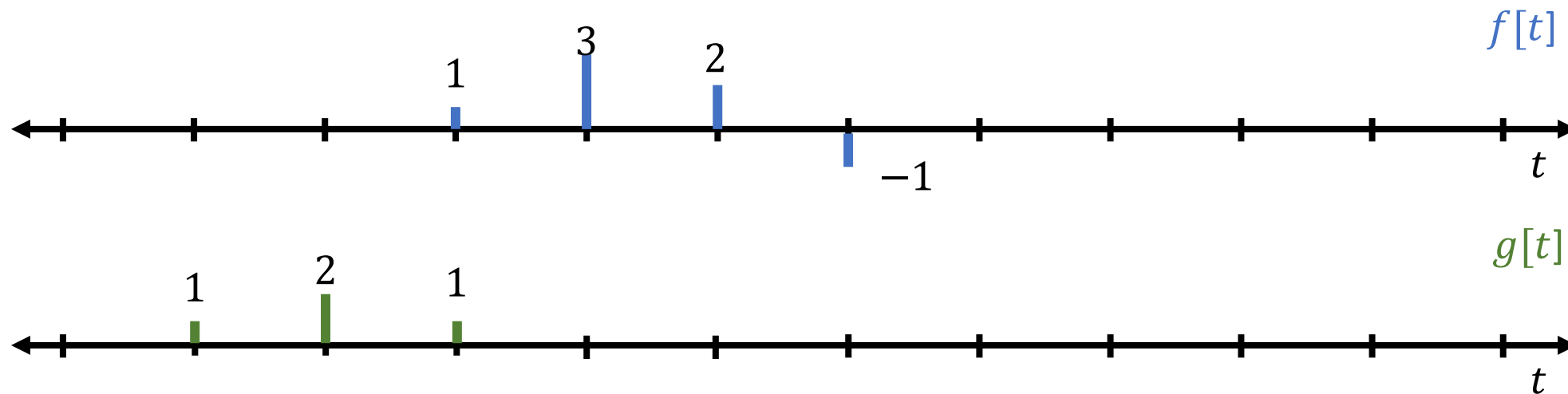
$$(f * g)(t) := \int_{-\infty}^{\infty} f(t - \tau)g(\tau)d\tau$$

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding



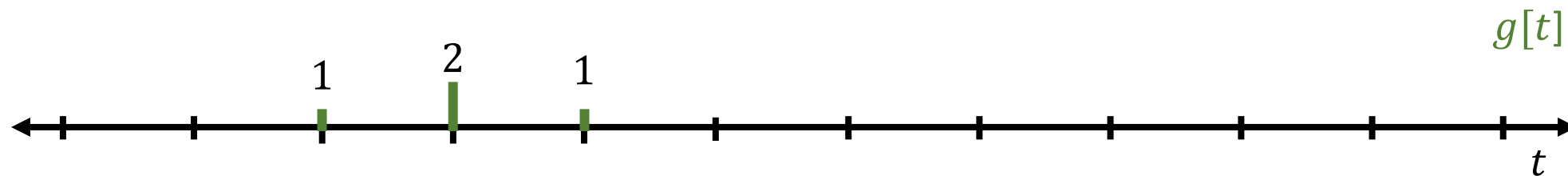
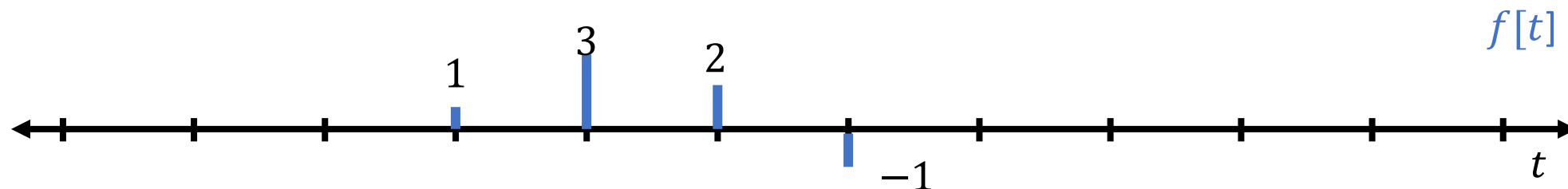
$$0 \cdot 1 + 0 \cdot 2 + 1 \cdot 1 = 1$$



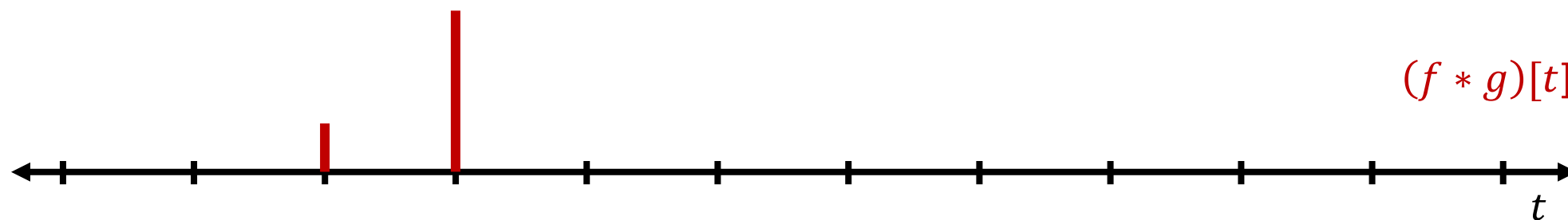
# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

- Flip the filter and sliding



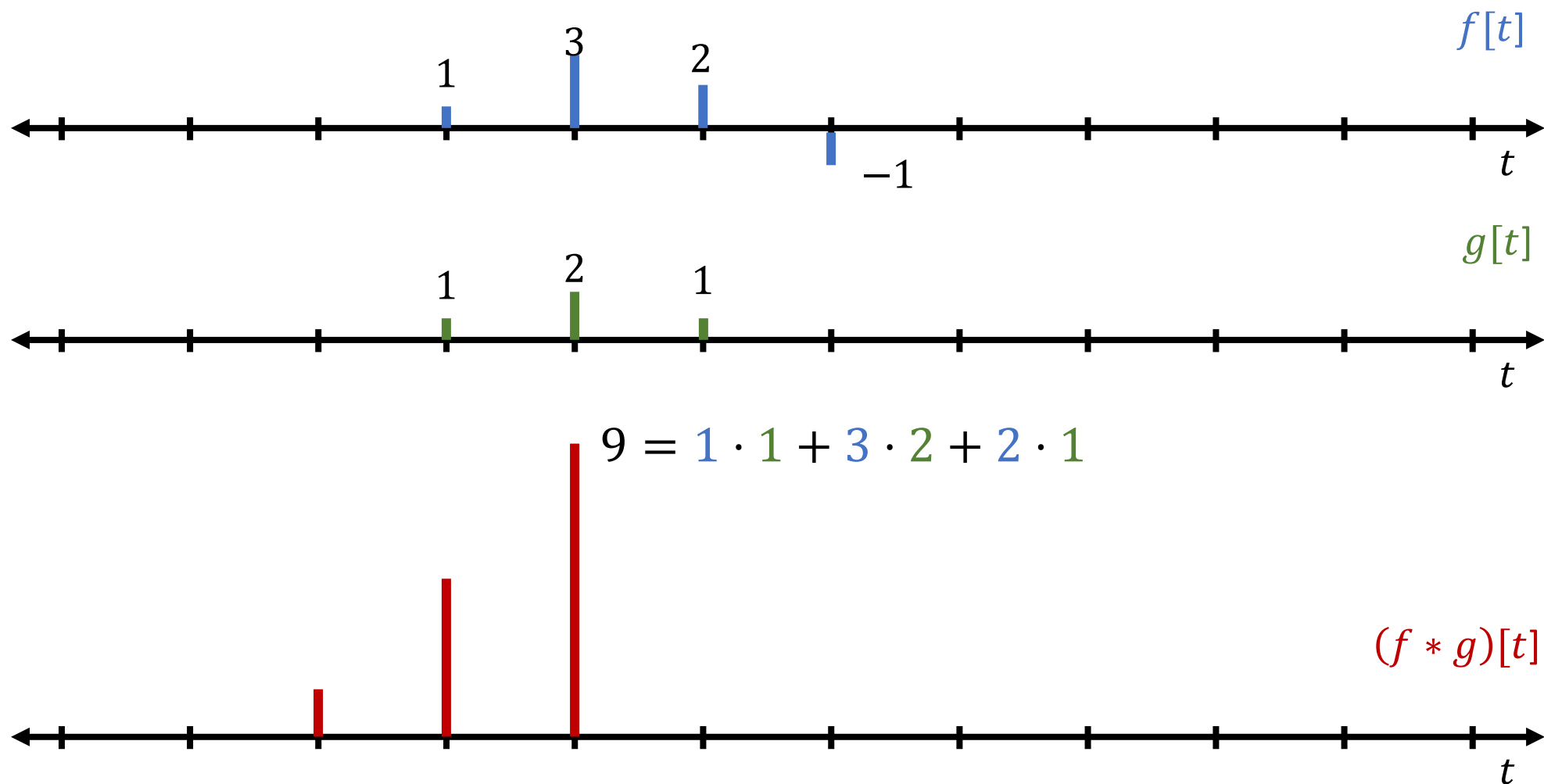
$$0 \cdot 1 + 1 \cdot 2 + 3 \cdot 1 = 5$$



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

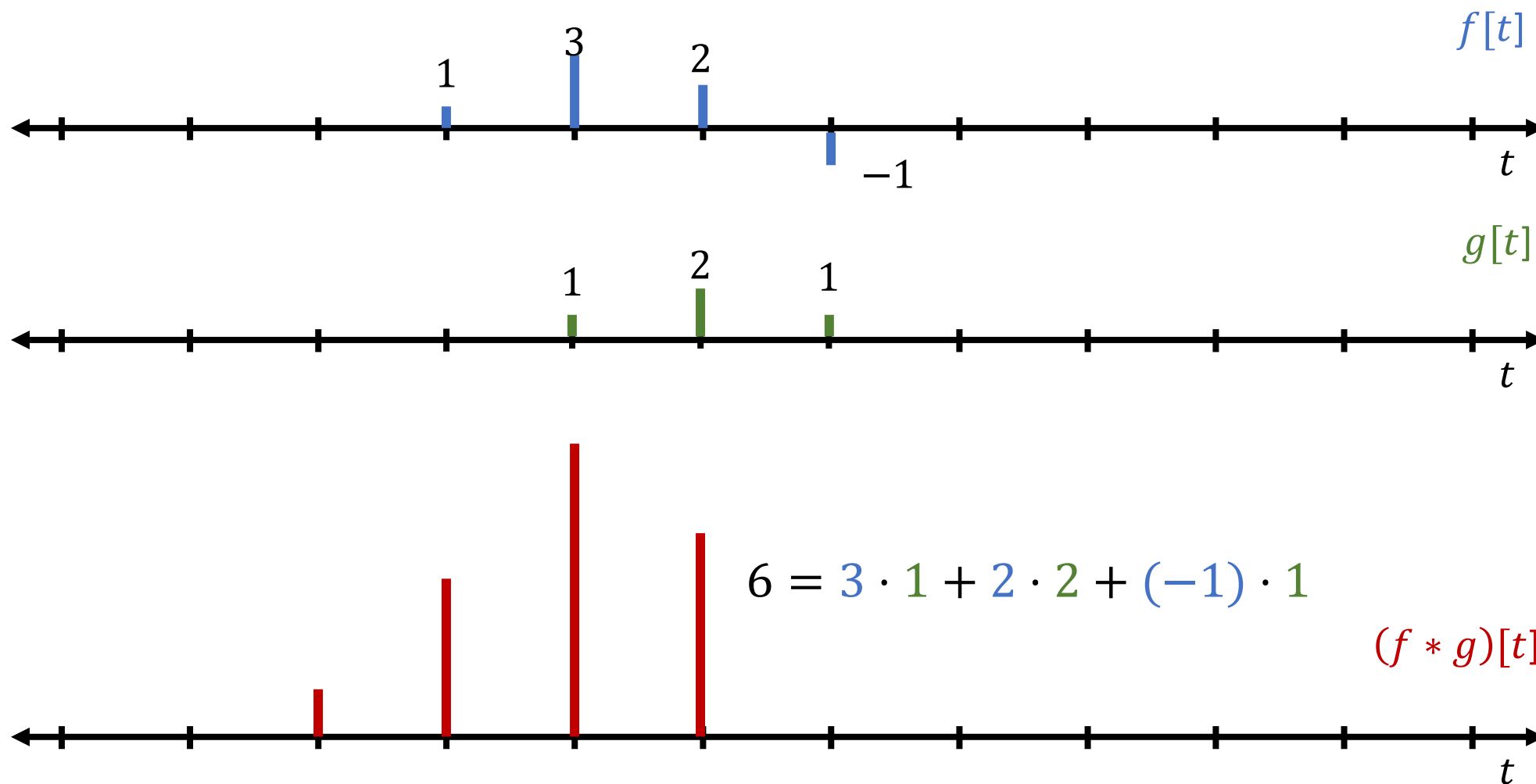
- Flip the filter and sliding



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

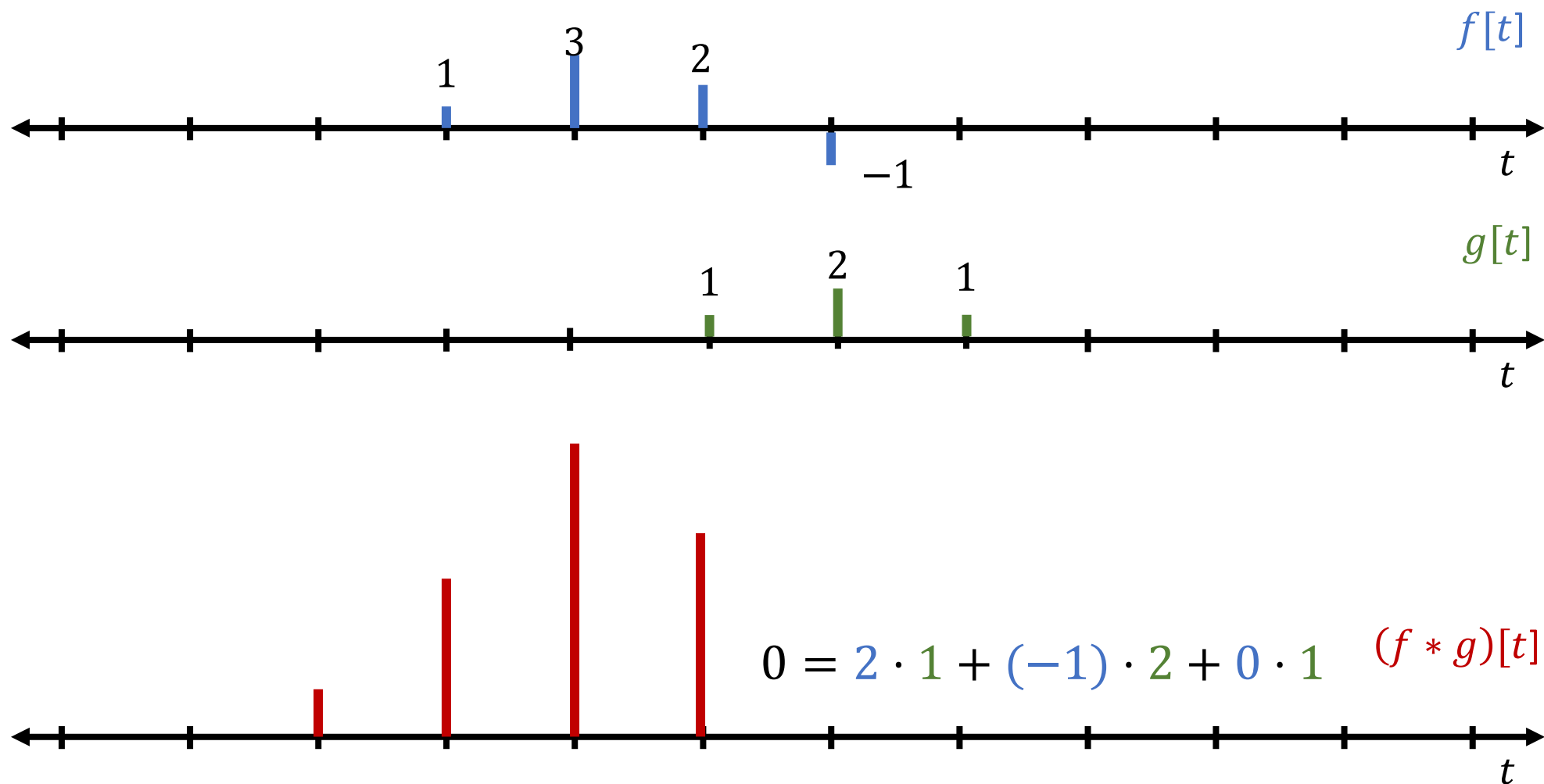
- Flip the filter and sliding



# 1D Convolution

$$(f * g)[t] := \sum_{\tau} f[t - \tau]g[\tau]$$

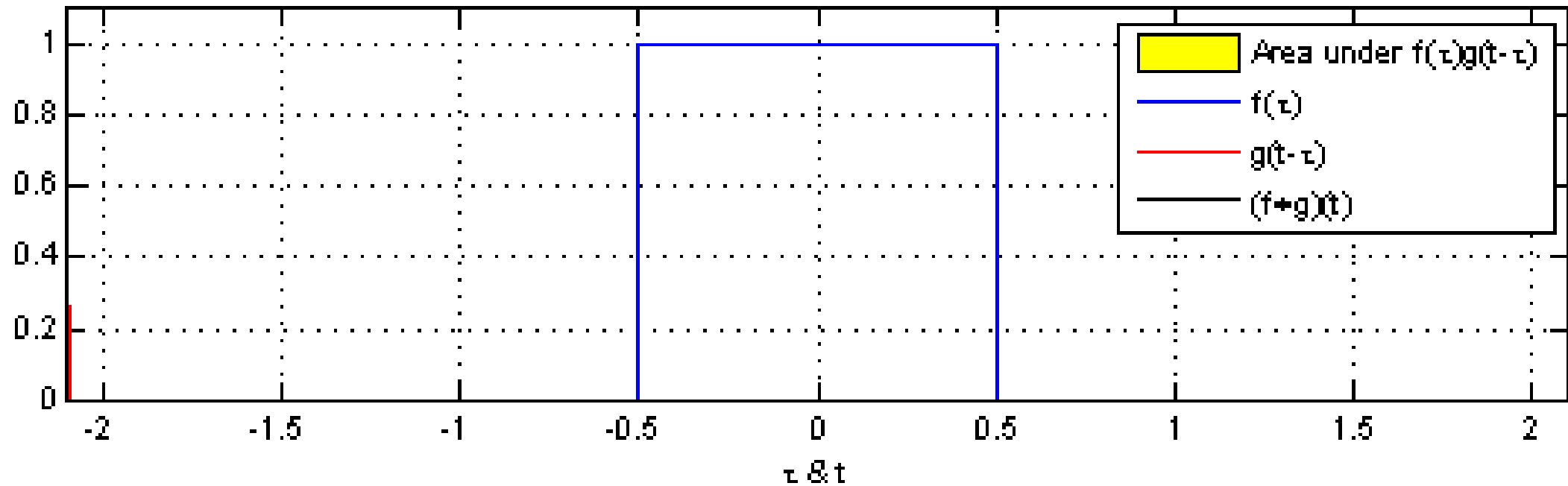
- Flip the filter and sliding





# 1D Convolution

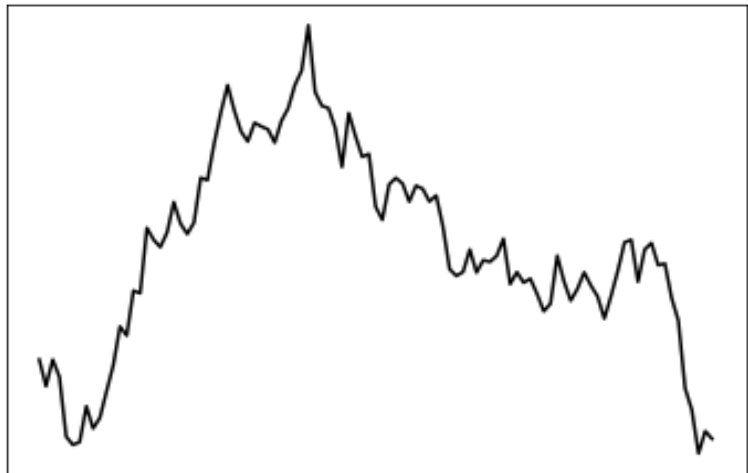
- Example



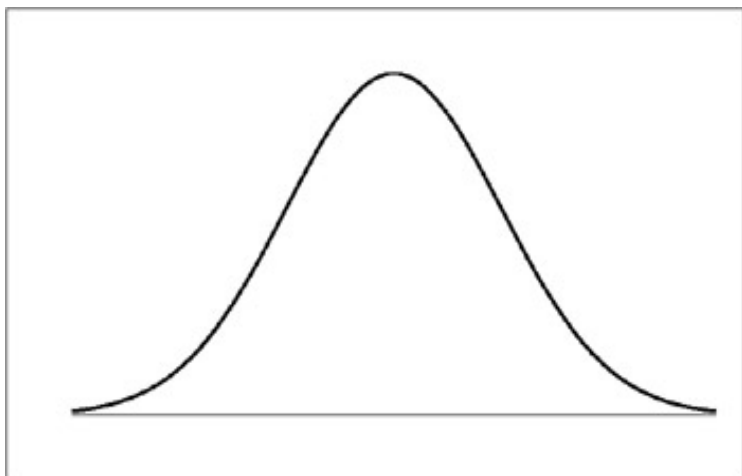
# 1D Convolution

- Gaussian filter

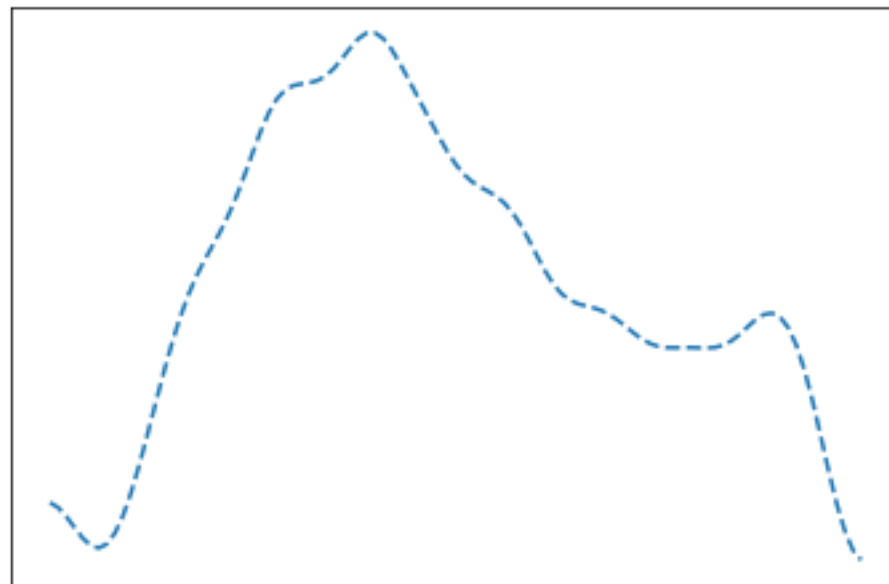
$f(t)$



$g(t)$



$f(t) * g(t)$



## 2D Convolution

$$(f * g)(s, t) := \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(s - \tau_1, t - \tau_2) g(\tau_1, \tau_2) d\tau_1 d\tau_2$$

$$(f * g)[s, t] := \sum_{\tau_1} \sum_{\tau_2} f[s - \tau_1, t - \tau_2] g[\tau_1, \tau_2]$$

# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

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

16				

# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

0	1	3	2	0	0	0
0	1	3	3	3	3	0
0	3	1	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28			

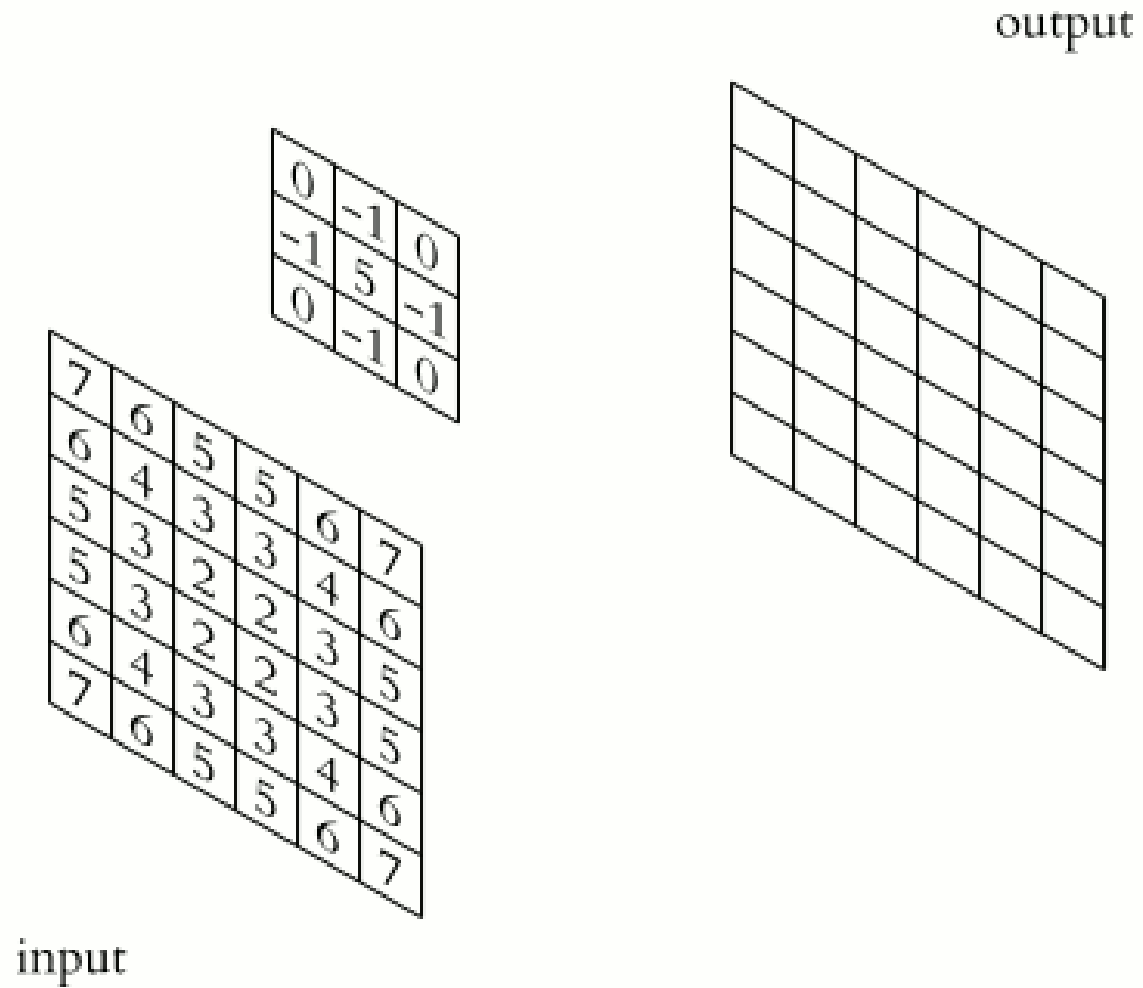
# 2D Convolution

- One input channel, e.g. gray color image
  - Padding=1, stride=1

0	0	1	3	2	0	0
0	1	1	3	3	3	0
0	3	3	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

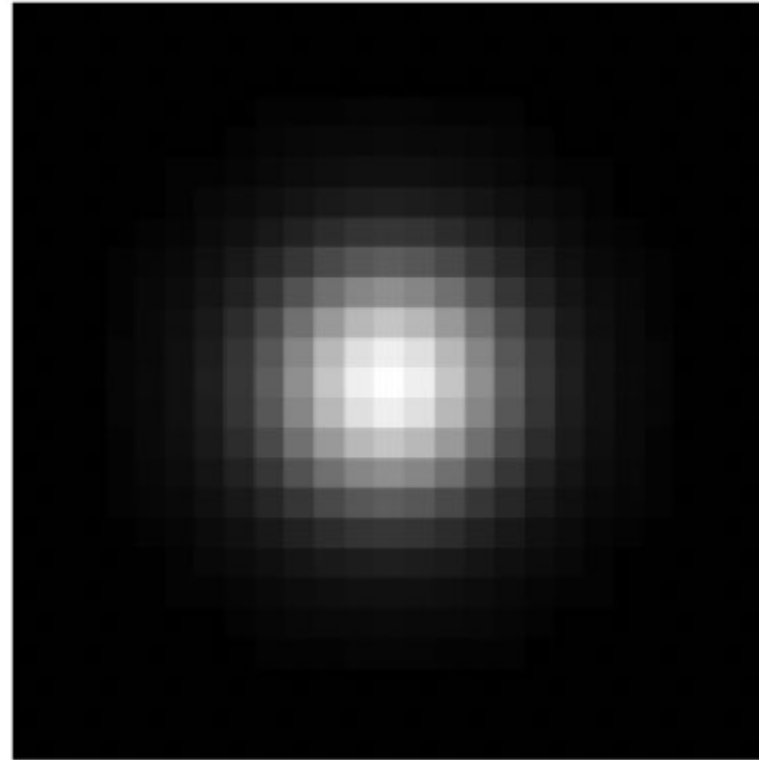
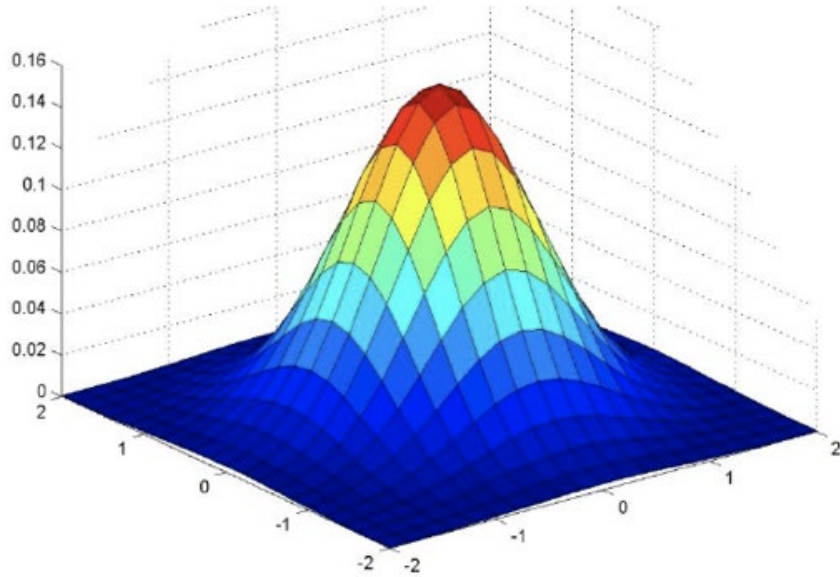
16	28	24		

# 2D Convolution



# 2D Convolution

- 2D gaussian filter





# 2D Convolution

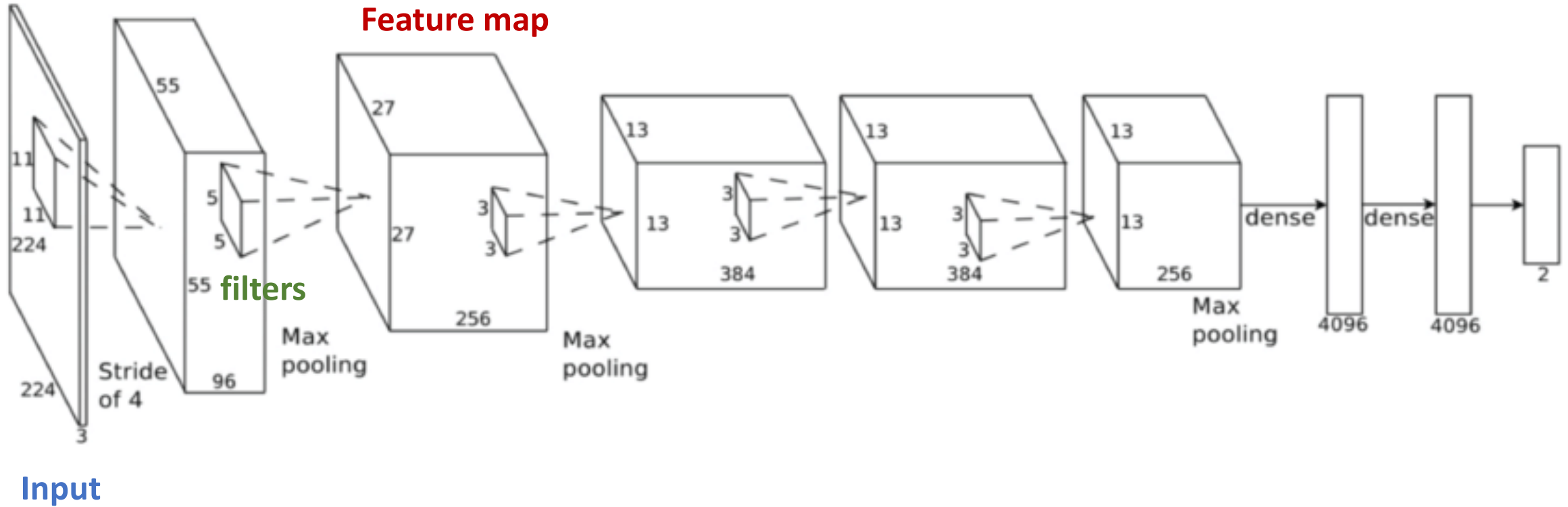
- 2D gaussian filter

Original Image (Left) Vs. Gaussian Filtered Image (Right)



# Convolutional Neural Networks

# Convolutional Neural Network



# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

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

16				

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	1	3	2	0	0	0
0	1	3	3	3	3	0
0	3	1	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28			

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	0	1	3	2	0	0
0	1	1	3	3	3	0
0	3	3	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	28	24		

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

Input

\*

1	3	2
1	3	3
3	1	1

filter

=

16	28	24	28	16
27	41	38	37	21
33	40	33	25	18
32	40	37	29	17
25	27	21	20	12

Output

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

1	3	2	0	0	0	0
1	3	3	2	3	3	0
3	1	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16		



# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	1	3	2	0	0
0	1	1	3	3	3	0
0	3	3	1	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, **stride=2**

0	0	0	0	1	3	2
0	1	3	2	1	3	3
0	3	1	2	3	1	1
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	16

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, **stride=2**

0	0	0	0	0	0	0
0	1	3	2	3	3	0
1	3	2	2	1	1	0
1	3	3	3	1	2	0
3	1	1	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

16	24	16
33		

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=1, stride=2

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

Input

\*

1	3	2
1	3	3
3	1	1

filter

=

16	24	16
33	33	18
25	21	12

Output

# 2D Convolution

- Input\_channel=1, output\_channel=1, padding=2, stride=1

0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0
0	0	1	3	2	3	3	0	0
0	0	3	1	2	1	1	0	0
0	0	3	3	3	1	2	0	0
0	0	2	2	1	2	1	0	0
0	0	2	3	2	1	2	0	0
0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0

\*

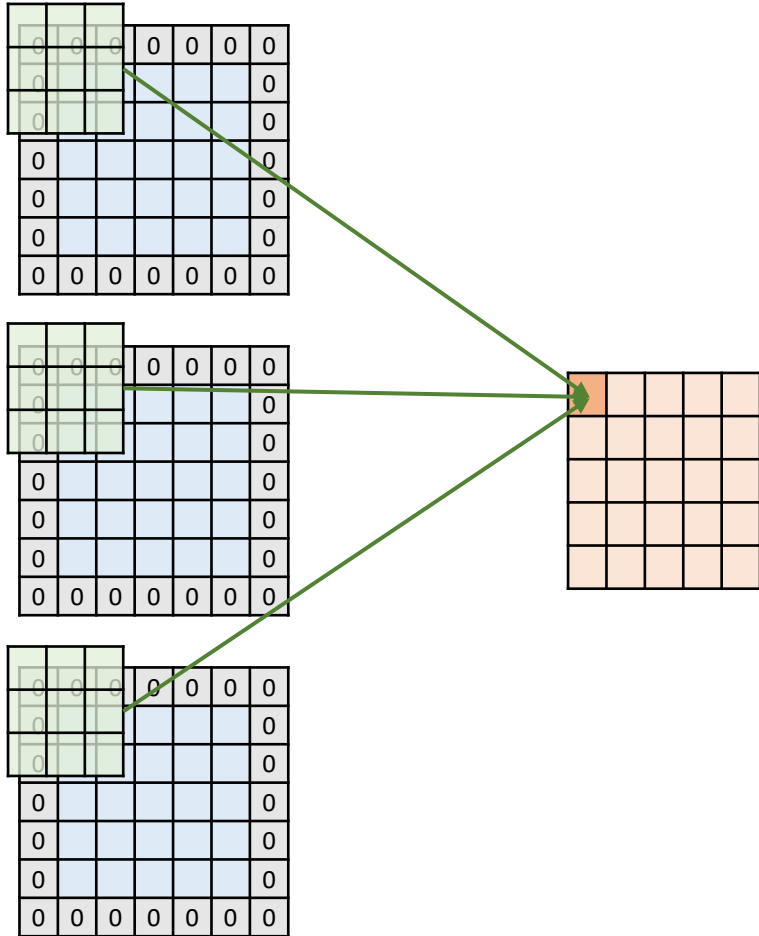
1	3	2
1	3	3
3	1	1

=

1	4	8	14	12	12	9
6	16	28	24	28	16	6
14	27	41	38	37	21	10
17	33	40	33	25	18	6
14	32	40	37	29	17	9
10	25	27	21	20	12	3
4	12	15	11	9	7	2

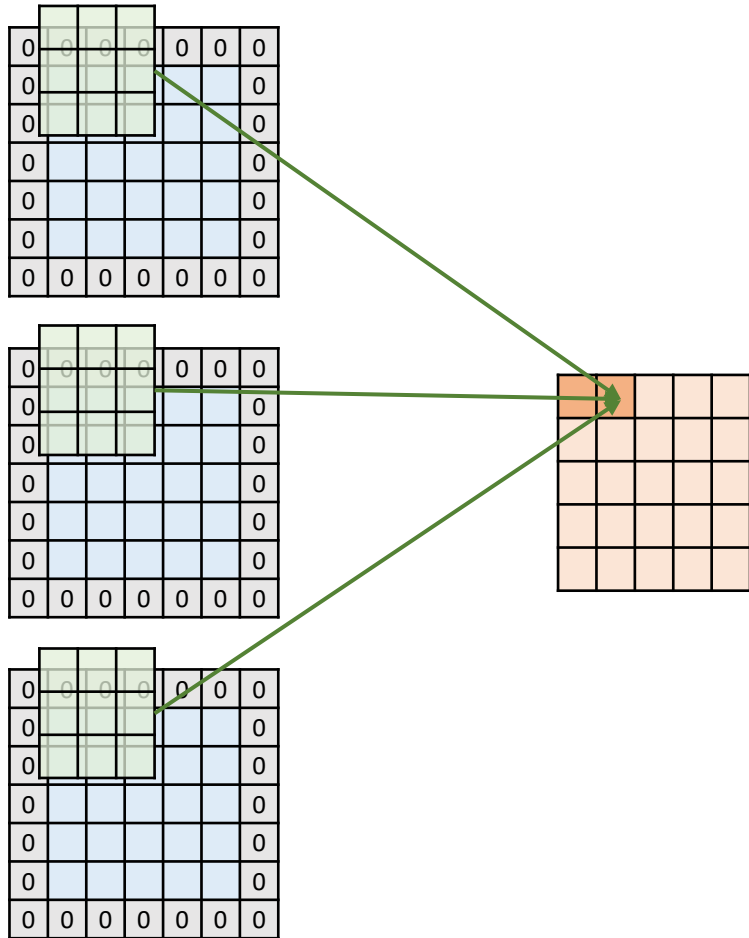
# 2D Convolution

- **Input\_channel=3**, output\_channel=1, padding=1, stride=1



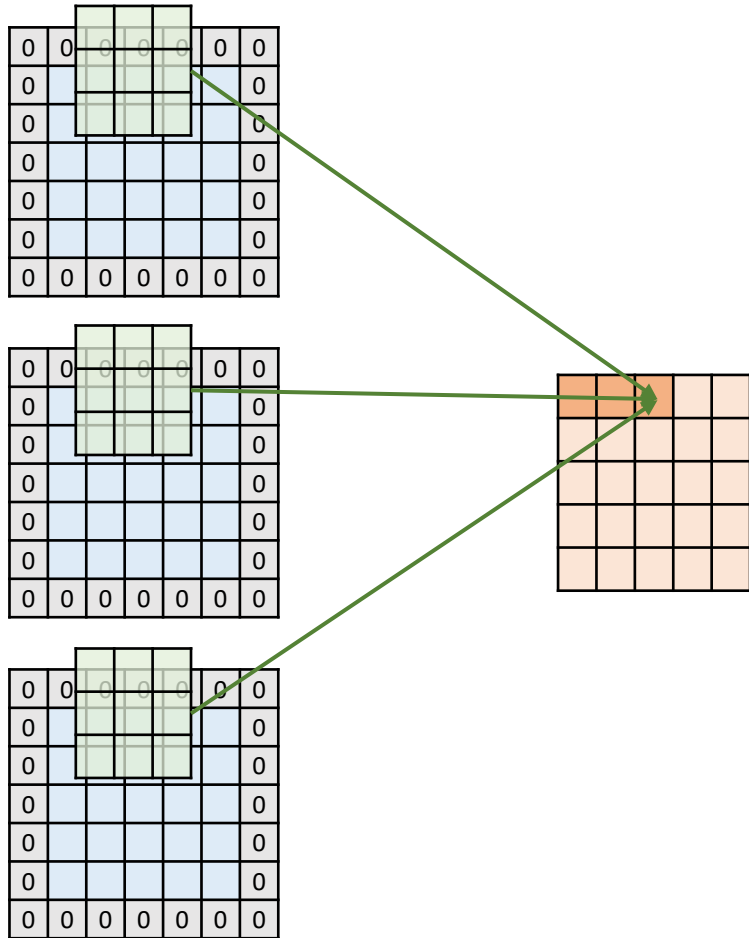
# 2D Convolution

- **Input\_channel=3**, output\_channel=1, padding=1, stride=1



# 2D Convolution

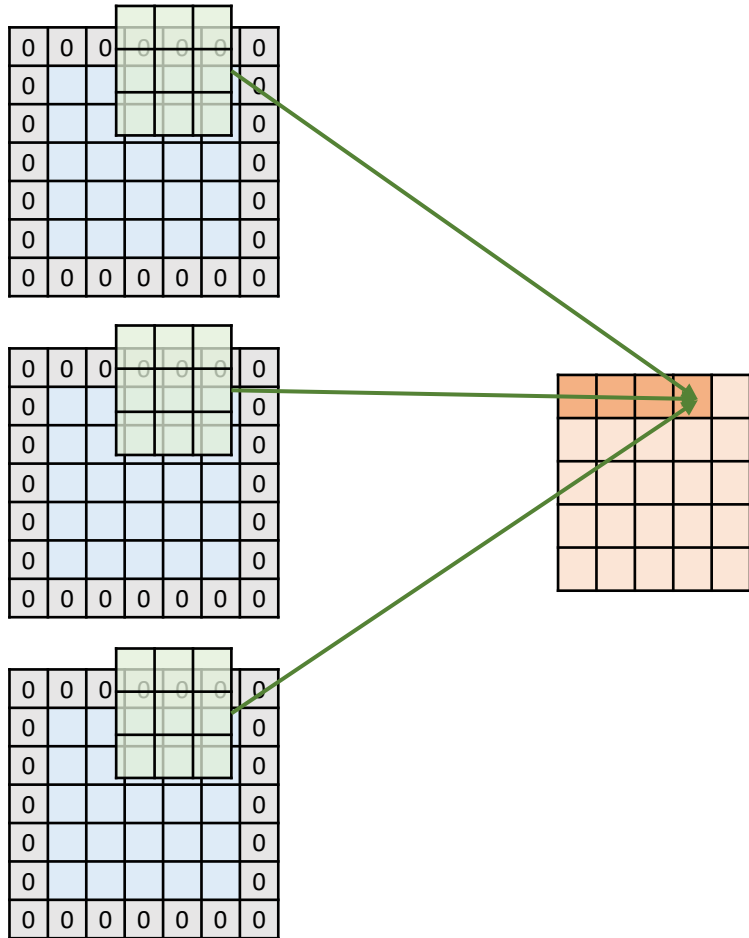
- **Input\_channel=3**, output\_channel=1, padding=1, stride=1





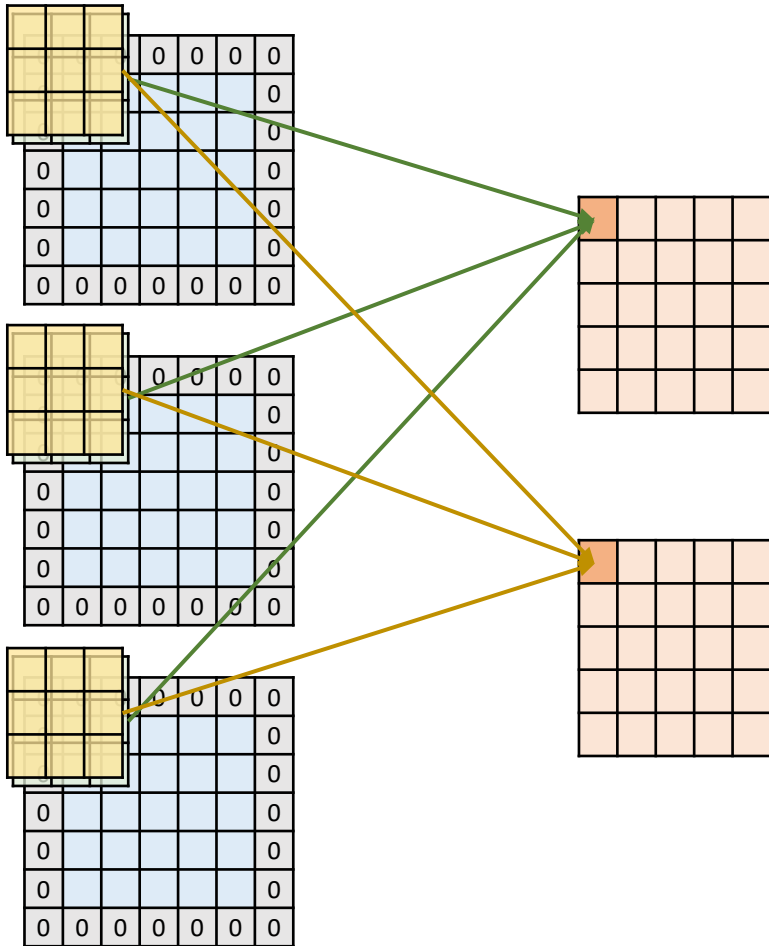
# 2D Convolution

- **Input\_channel=3**, output\_channel=1, padding=1, stride=1



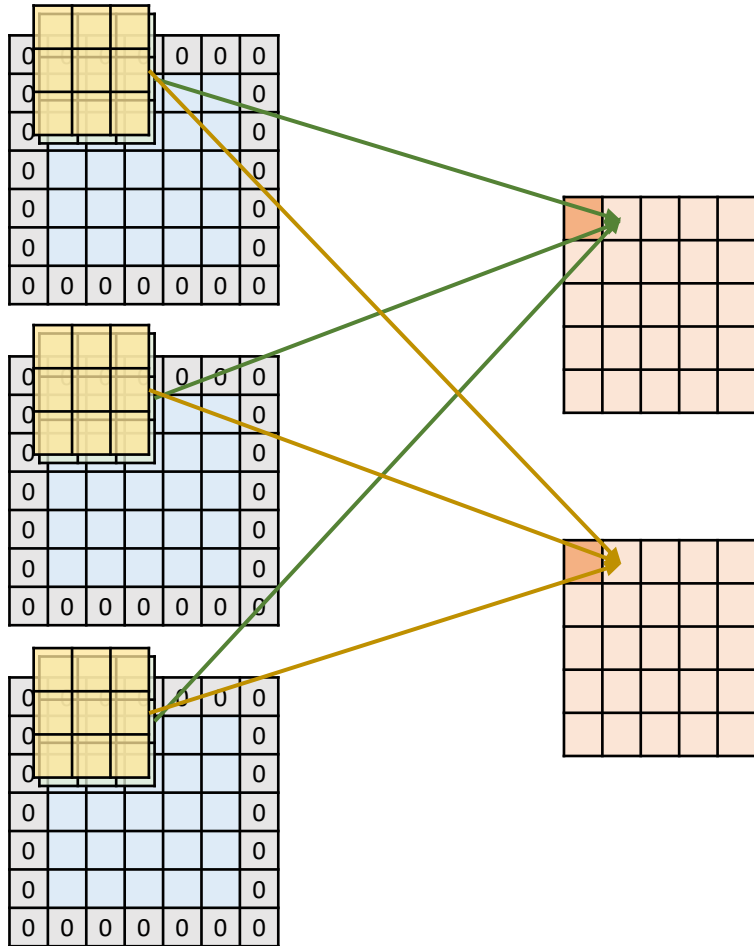
# 2D Convolution

- Input\_channel=3, output\_channel=2, padding=1, stride=1



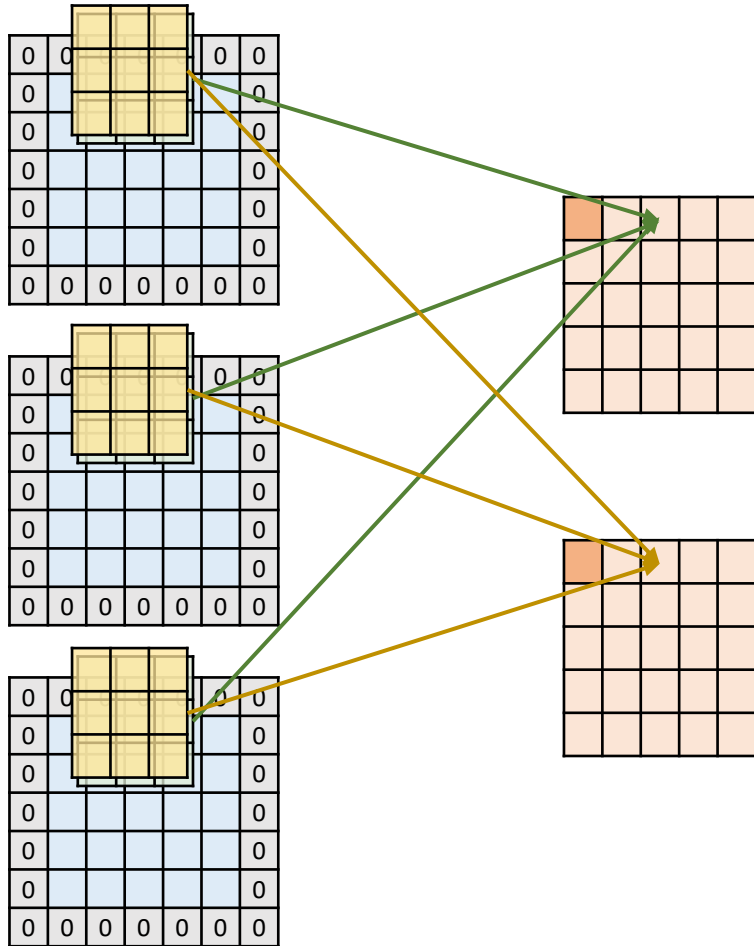
# 2D Convolution

- Input\_channel=3, output\_channel=2, padding=1, stride=1



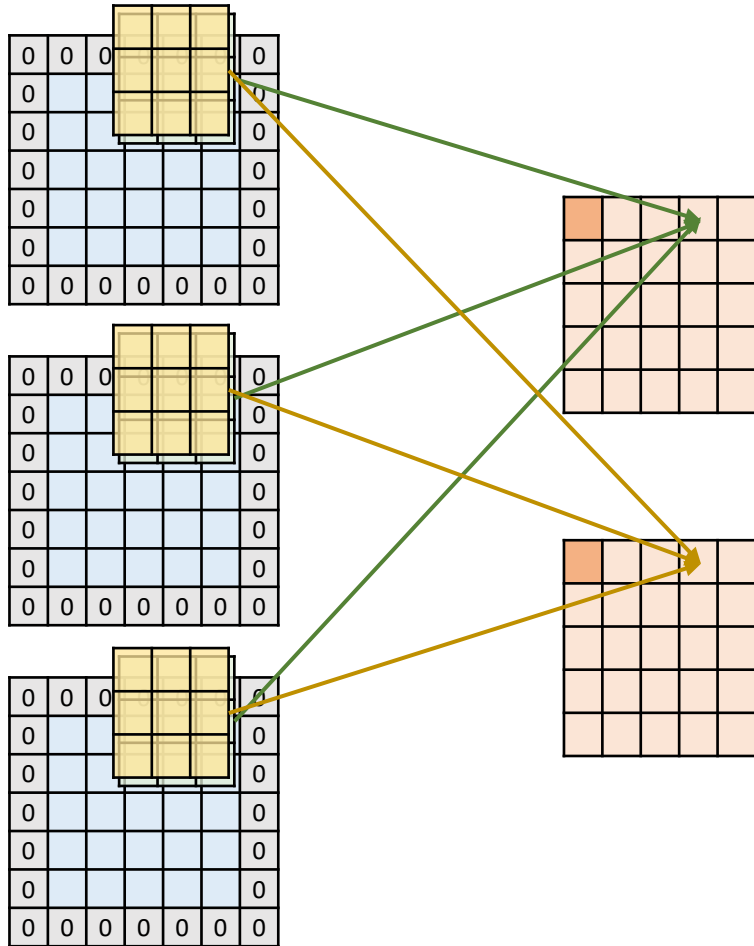
# 2D Convolution

- Input\_channel=3, output\_channel=2, padding=1, stride=1



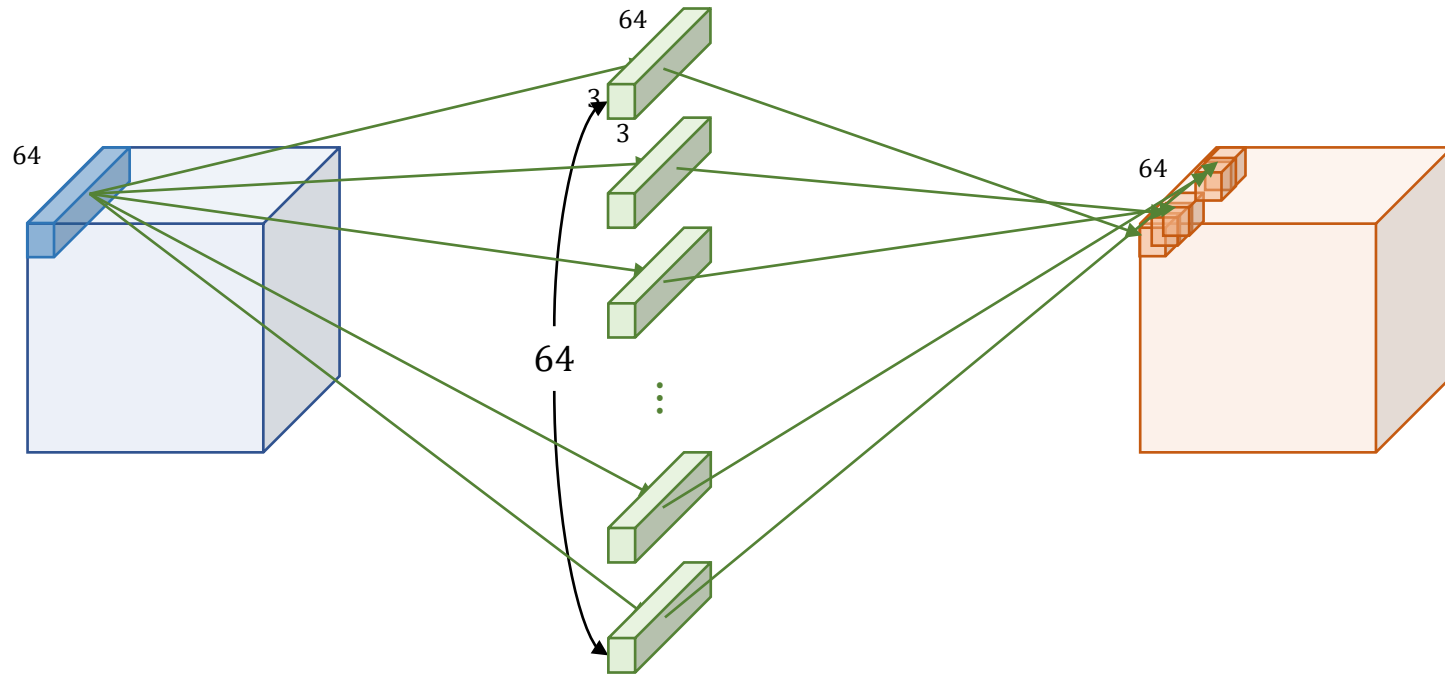
# 2D Convolution

- Input\_channel=3, output\_channel=2, padding=1, stride=1



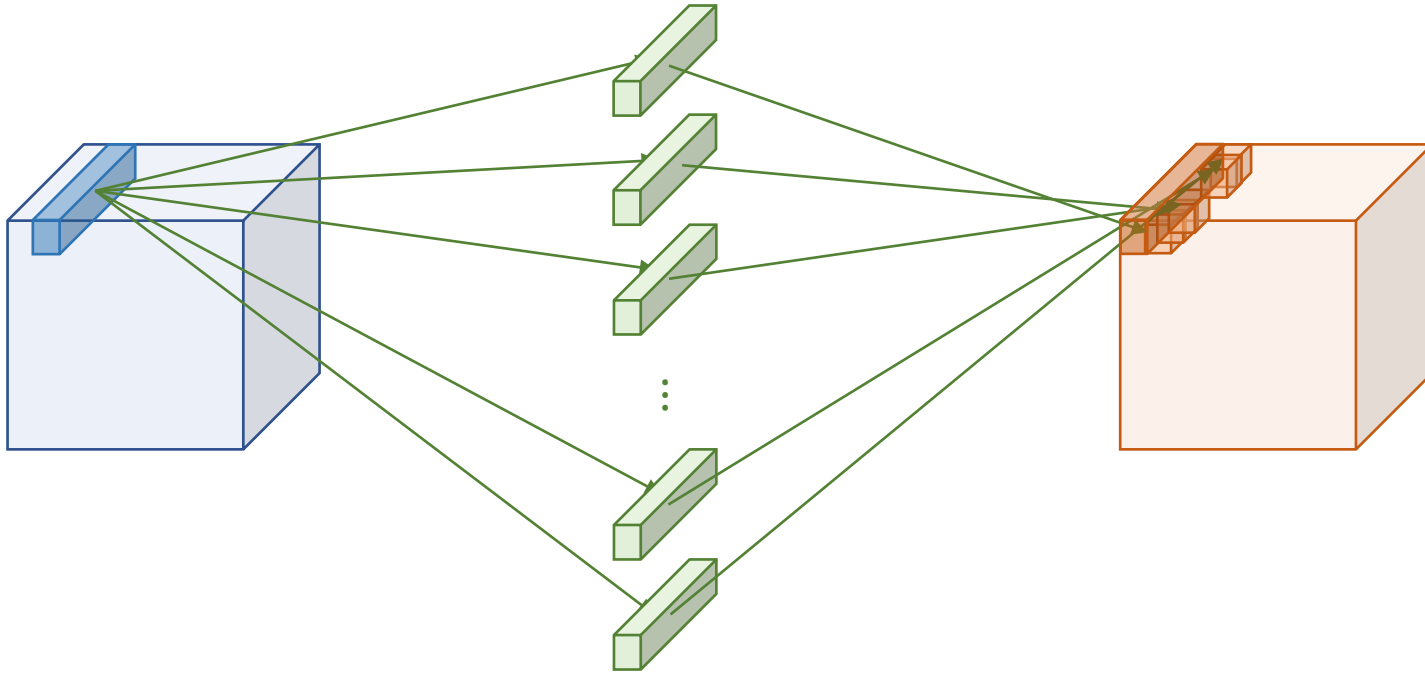
# 2D Convolution

- Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1



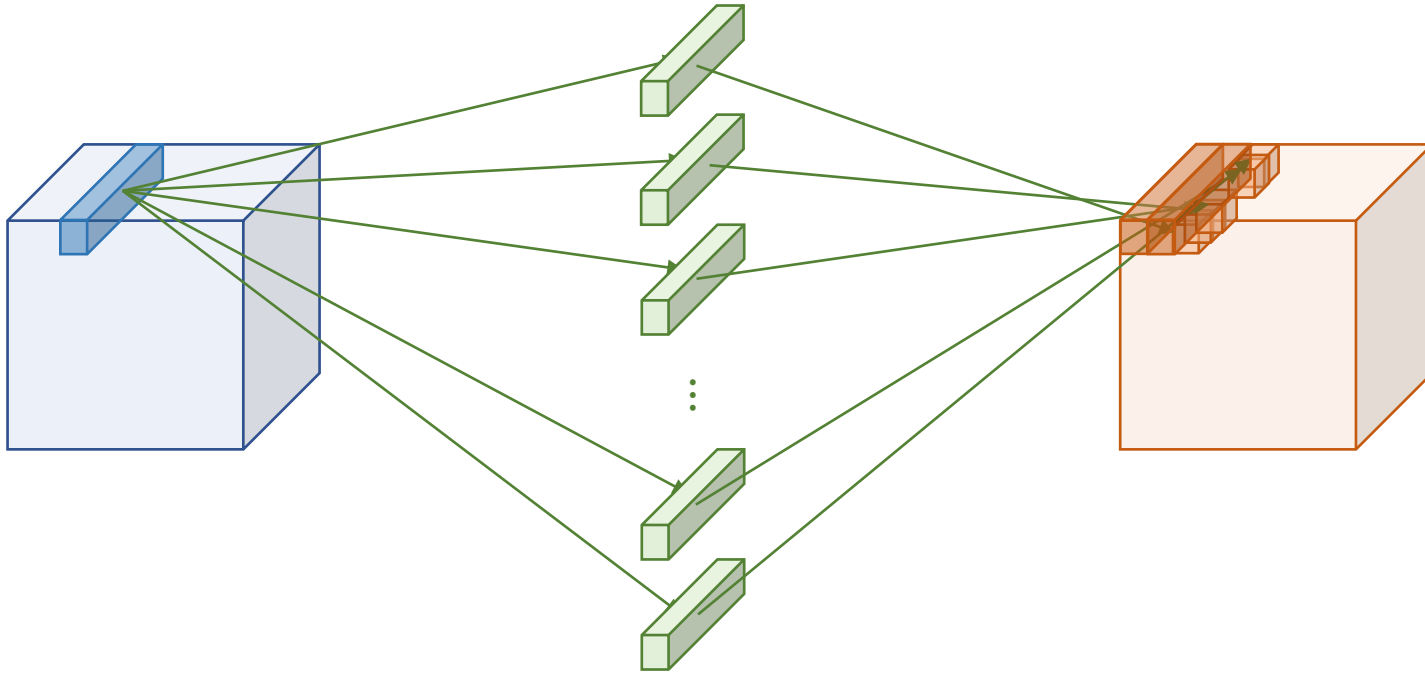
# 2D Convolution

- Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1



# 2D Convolution

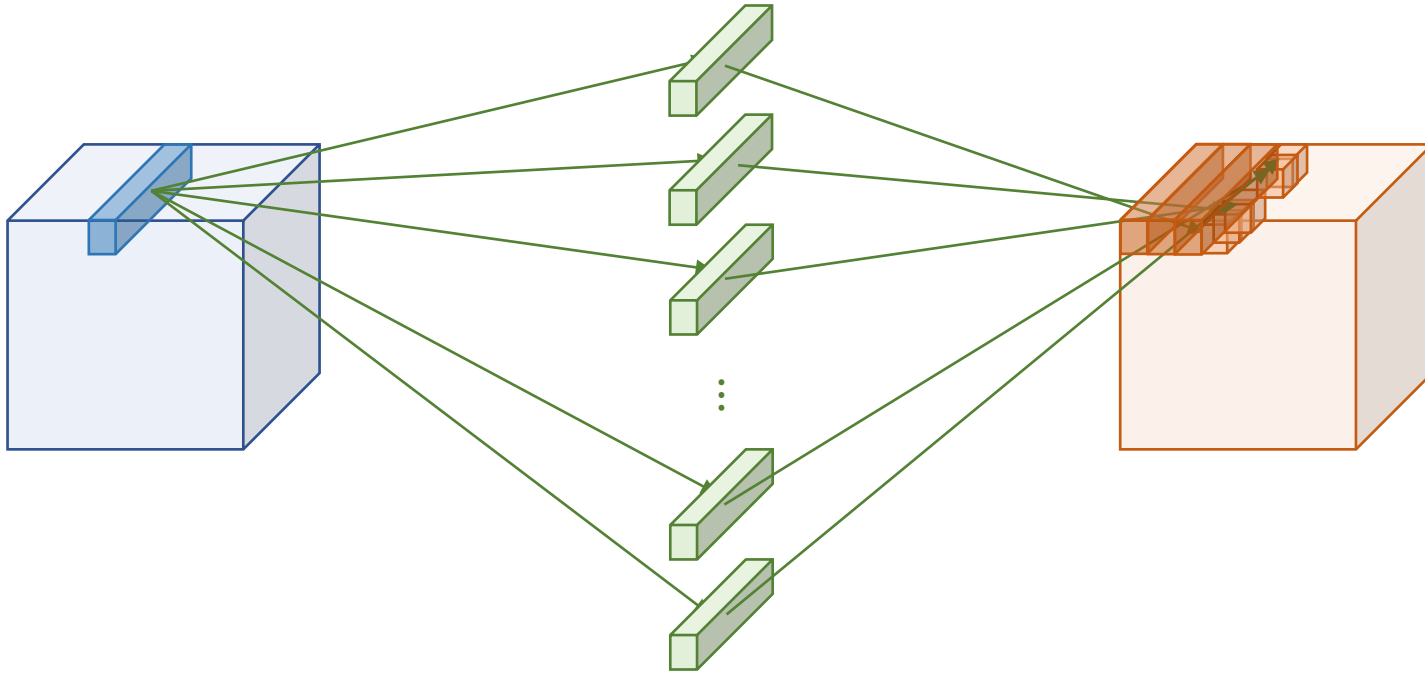
- Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1





# 2D Convolution

- Input\_channel=64, output\_channel=64, kernel\_size=3, padding=1, stride=1



# Convolutions in PyTorch

---

**CLASS** `torch.nn.Conv1d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

---

**CLASS** `torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

---

**CLASS** `torch.nn.Conv3d(in_channels, out_channels, kernel_size, stride=1, padding=0, dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)` [\[SOURCE\]](#)

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0	0
0	1	3	2	3	3	0	
0	3	1	2	1	1	0	
	0	3	3	3	1	2	0
	0	2	2	1	2	1	0
	0	2	3	2	1	2	0
	0	0	0	0	0	0	0

3		

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

3	3	

# Max Pooling

- Pooling a maximum value given the window
- Used to reduce the size of feature maps
- Example) stride=2, padding=1

0	0	0	0	0	0	0
0	1	3	2	3	3	0
0	3	1	2	1	1	0
0	3	3	3	1	2	0
0	2	2	1	2	1	0
0	2	3	2	1	2	0
0	0	0	0	0	0	0

3	3	3

# Max Pooling in PyTorch

```
CLASS torch.nn.MaxPool1d(kernel_size, stride=None, padding=0, dilation=1,  
    return_indices=False, ceil_mode=False)
```

[\[SOURCE\]](#)

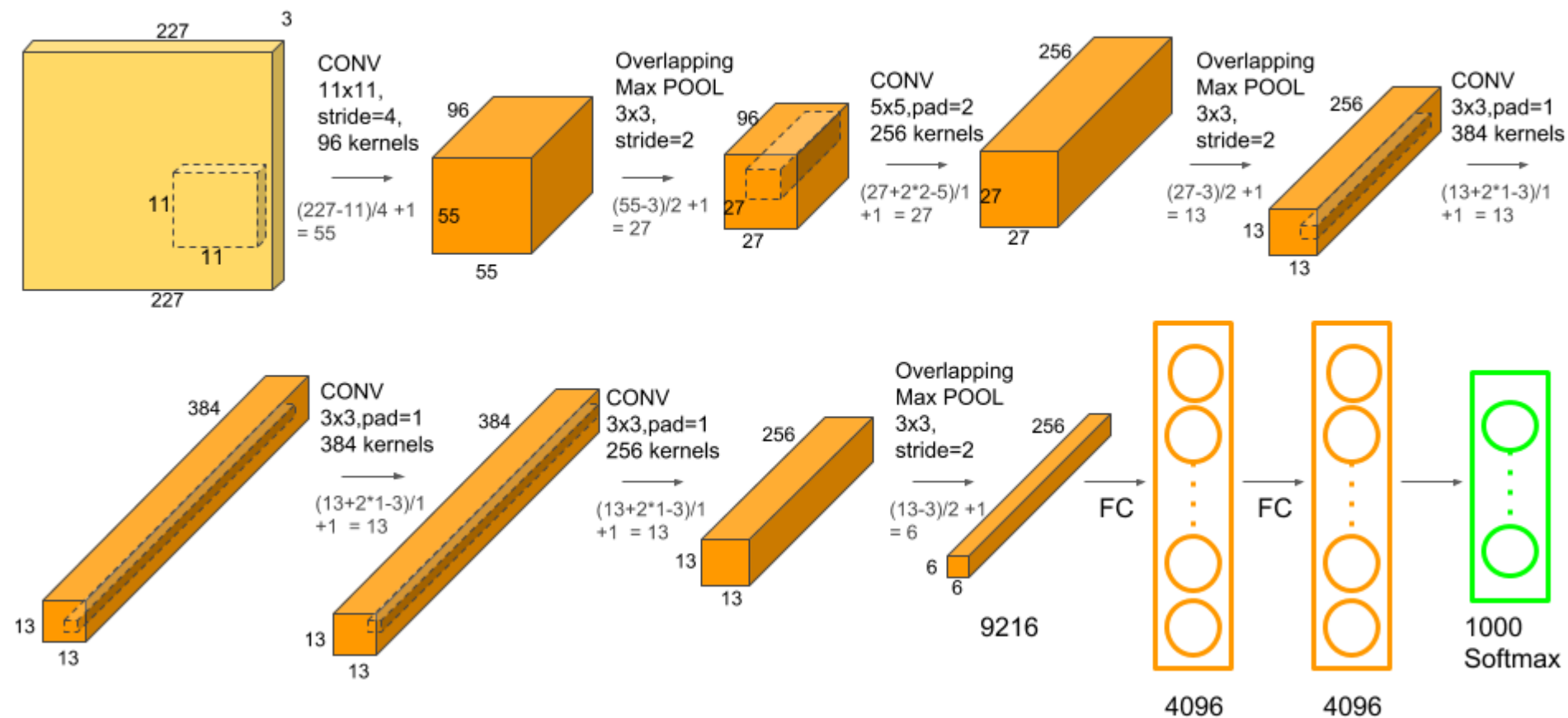
```
CLASS torch.nn.MaxPool2d(kernel_size, stride=None, padding=0, dilation=1,  
    return_indices=False, ceil_mode=False)
```

[\[SOURCE\]](#)

```
CLASS torch.nn.MaxPool3d(kernel_size, stride=None, padding=0, dilation=1,  
    return_indices=False, ceil_mode=False)
```

[\[SOURCE\]](#)

# AlexNet

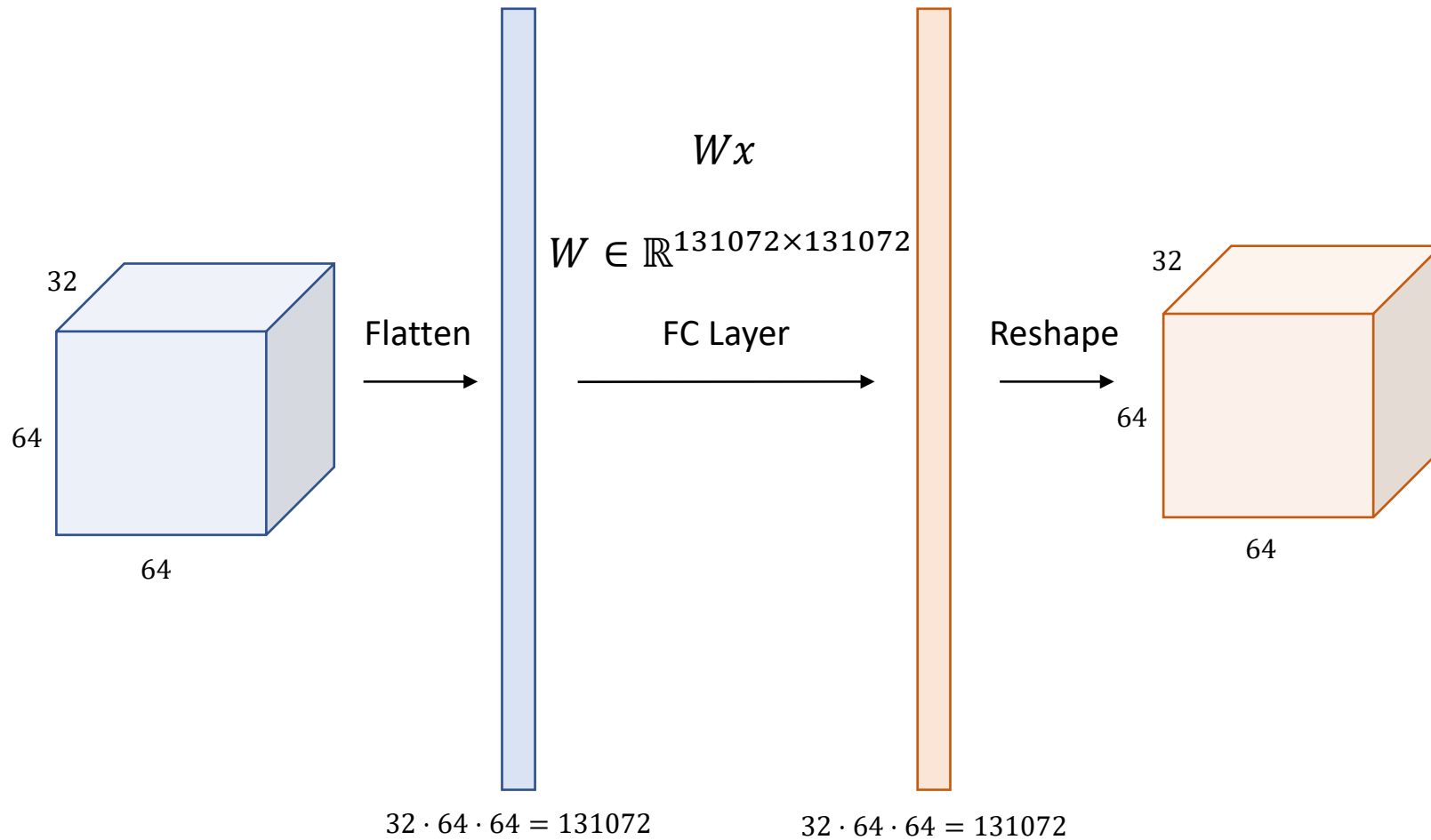


Dropout

Dropout

# Fully Connected Layer vs Convolutional Layer

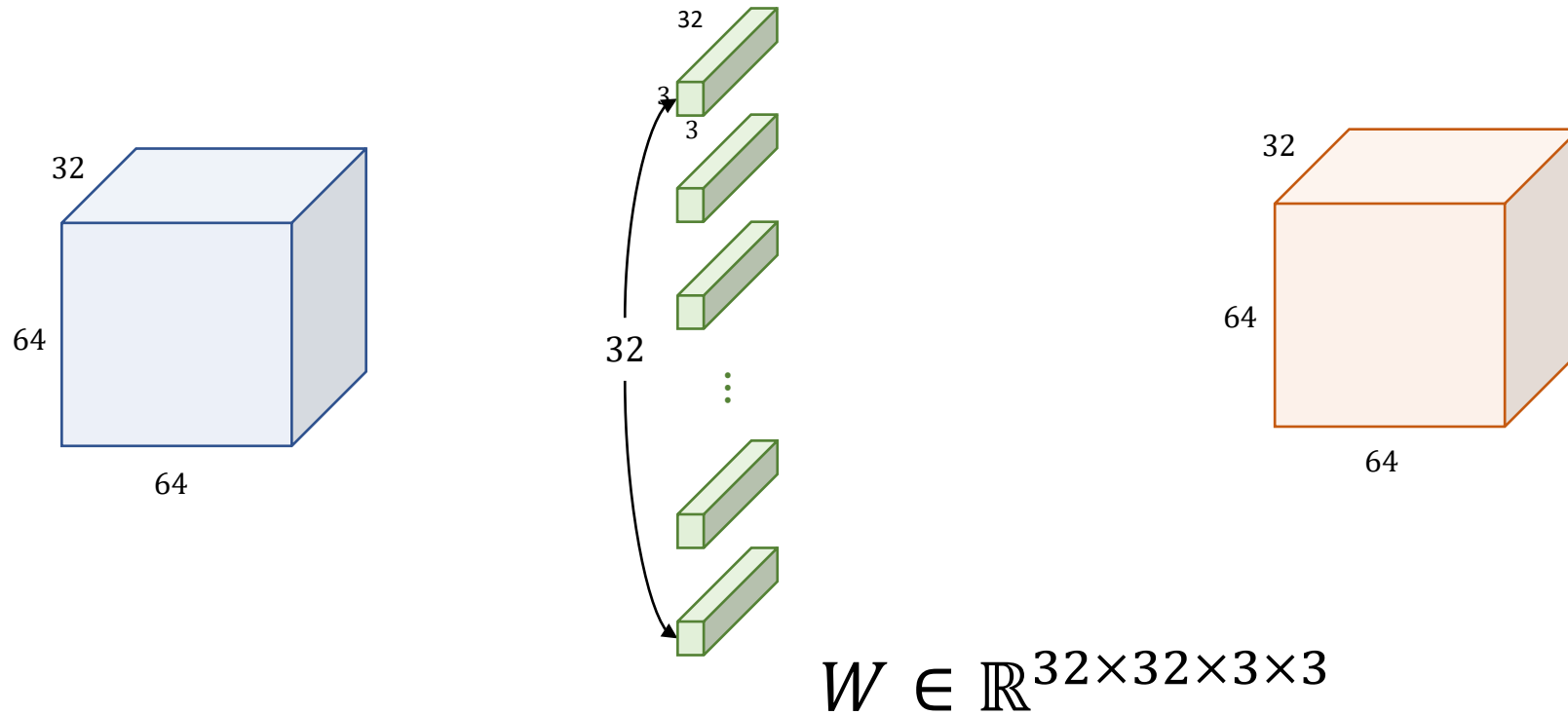
- Translation equivariance and parameter sharing





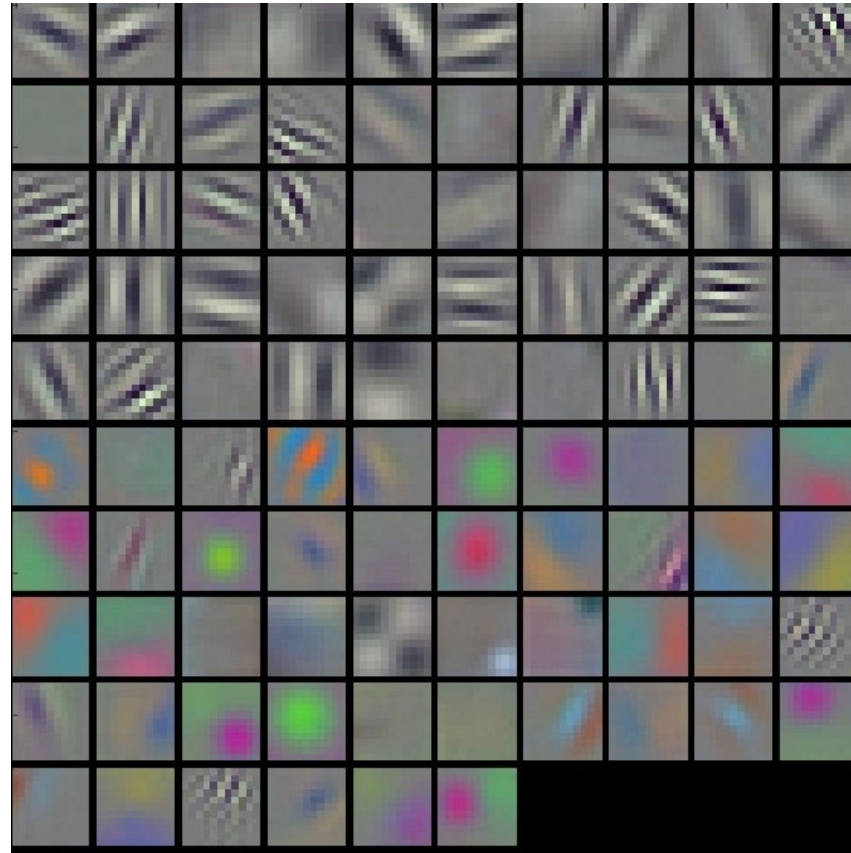
# Fully Connected Layer vs Convolutional Layer

- Translation equivariance and parameter sharing

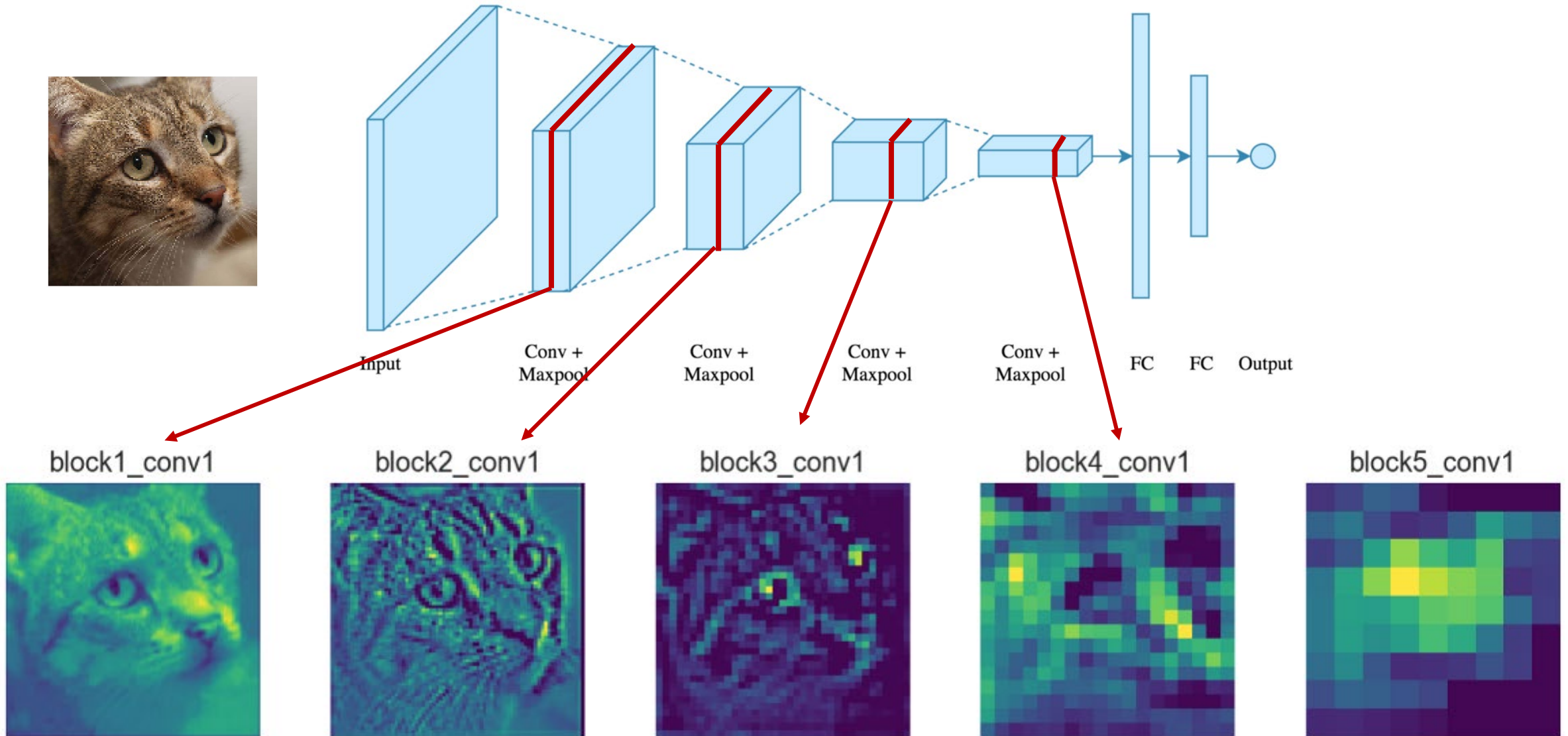


# Visualization of Learned Filter

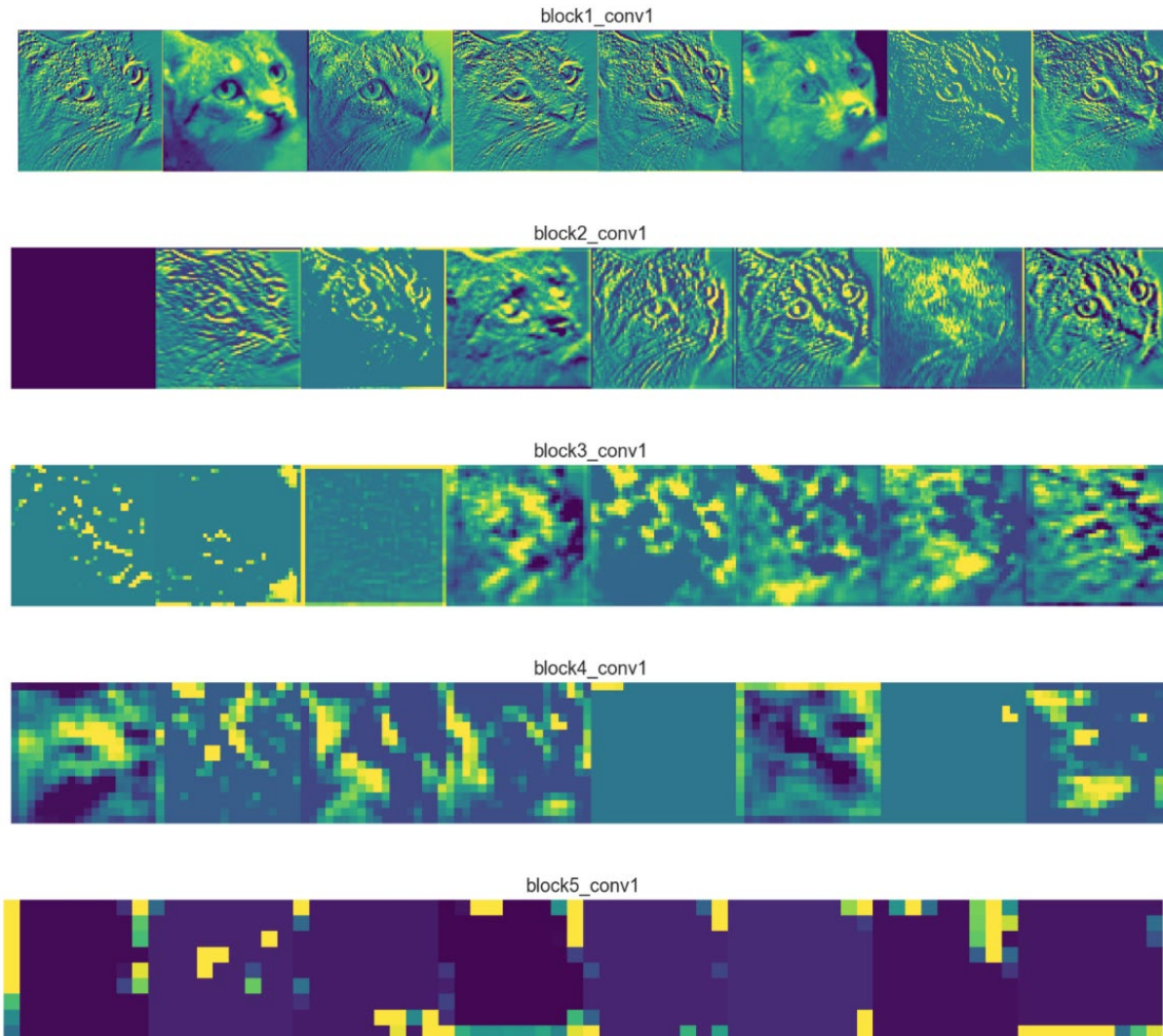
- First layer conv filters



# Visualization of Learned Feature Maps



# Visualization of Learned Feature Maps



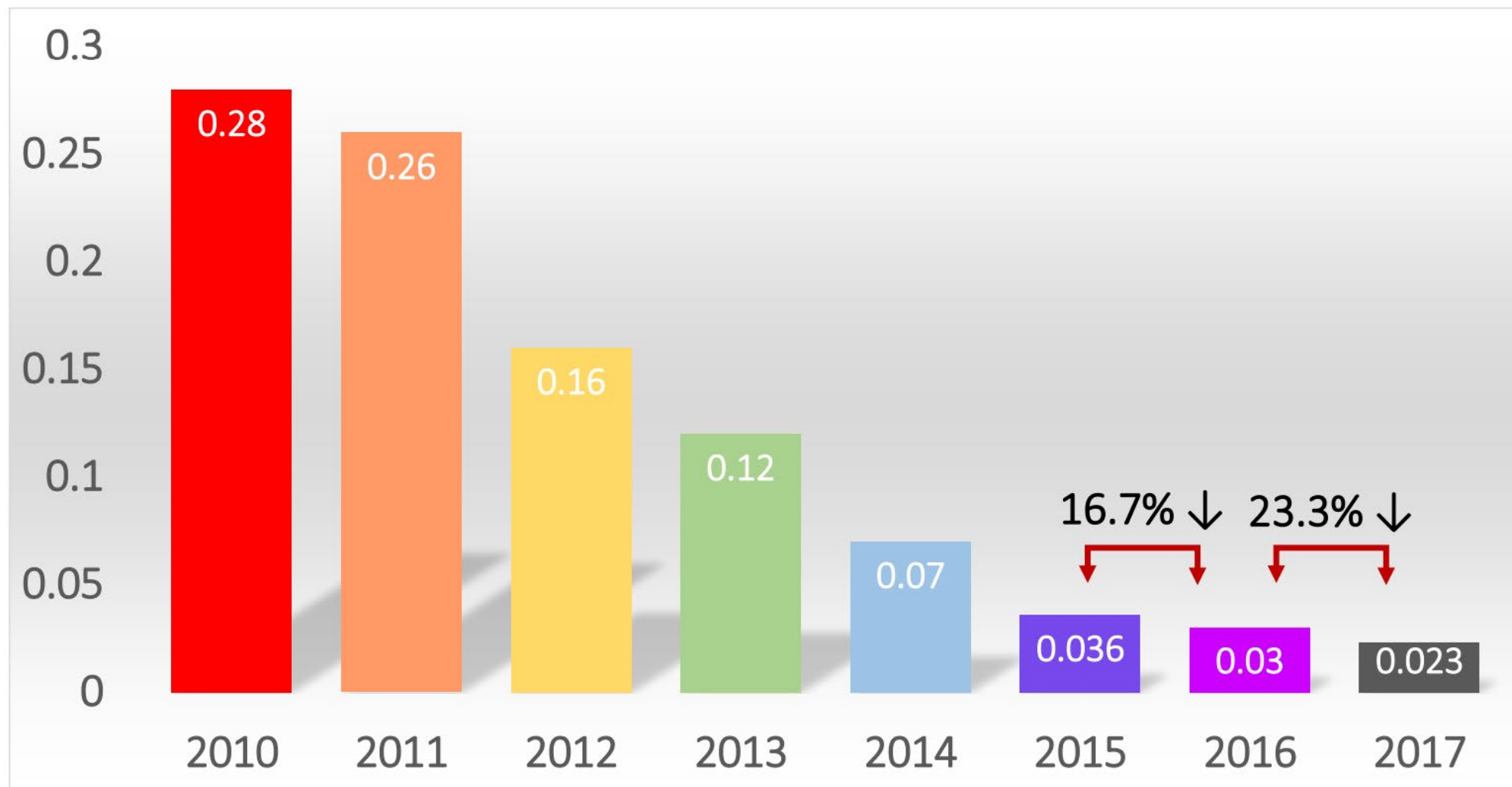
# ImageNet Large Scale Visual Recognition Challenge (ILSVRC)

# ILSVRC

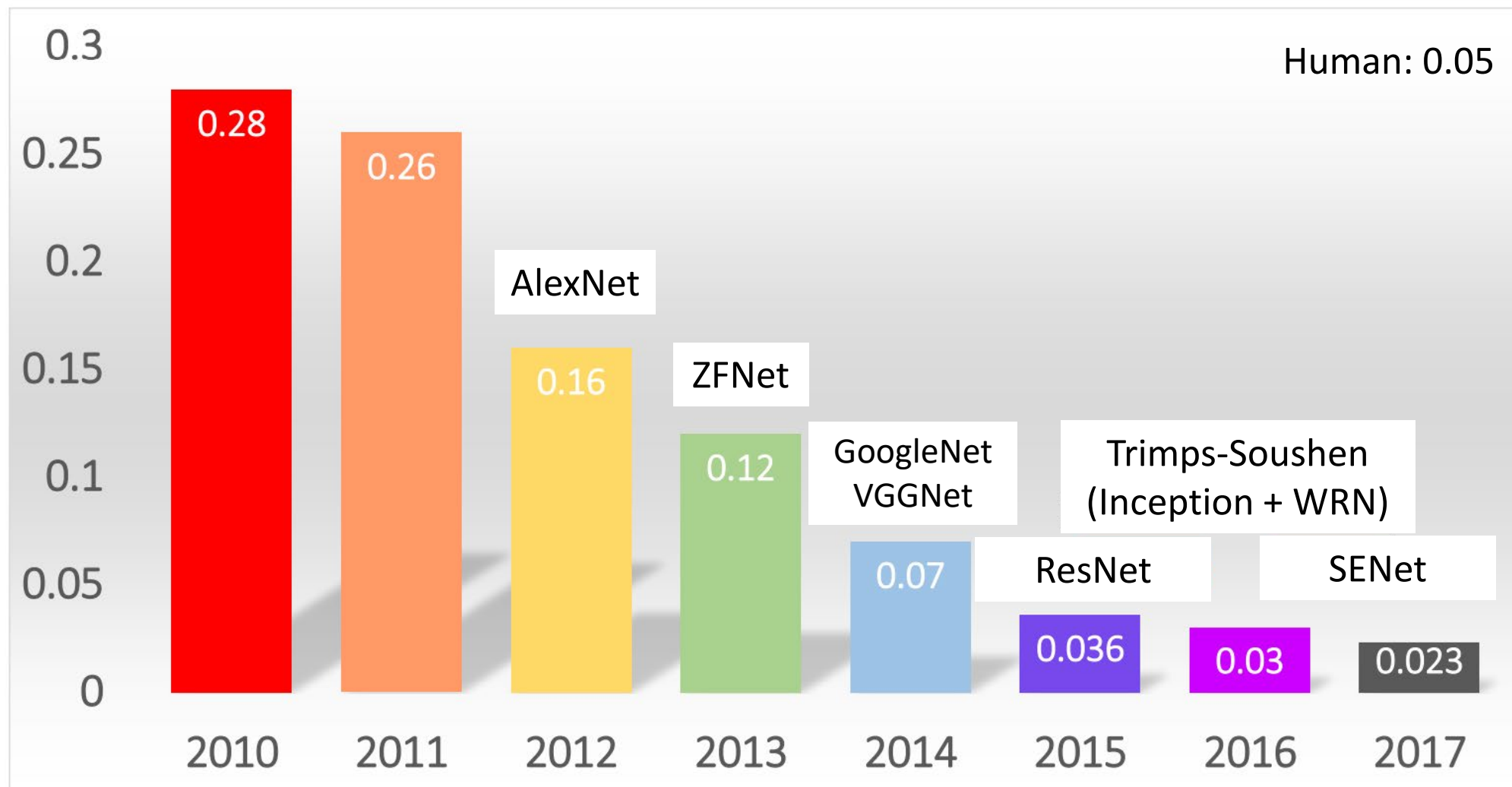
- ImageNet is an image database organized according to the WordNet hierarchy (nouns)
  - 1000 object classes
  - About 1.2M training images, 50K validation images, 100K test images
- The ImageNet Large Scale Visual Recognition Challenge (ILSVRC)
  - 8 years history (2010 – 2017)
  - It was the most powerful driving force to facilitate deep learning research



# Classification Results



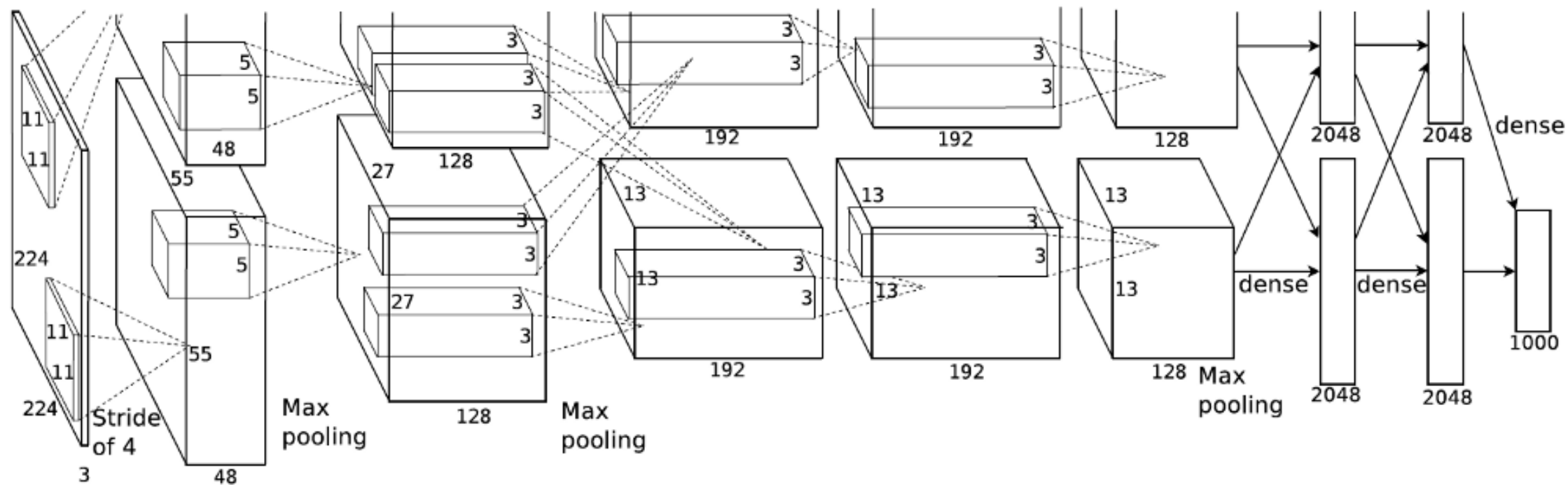
# Classification Results





# AlexNet

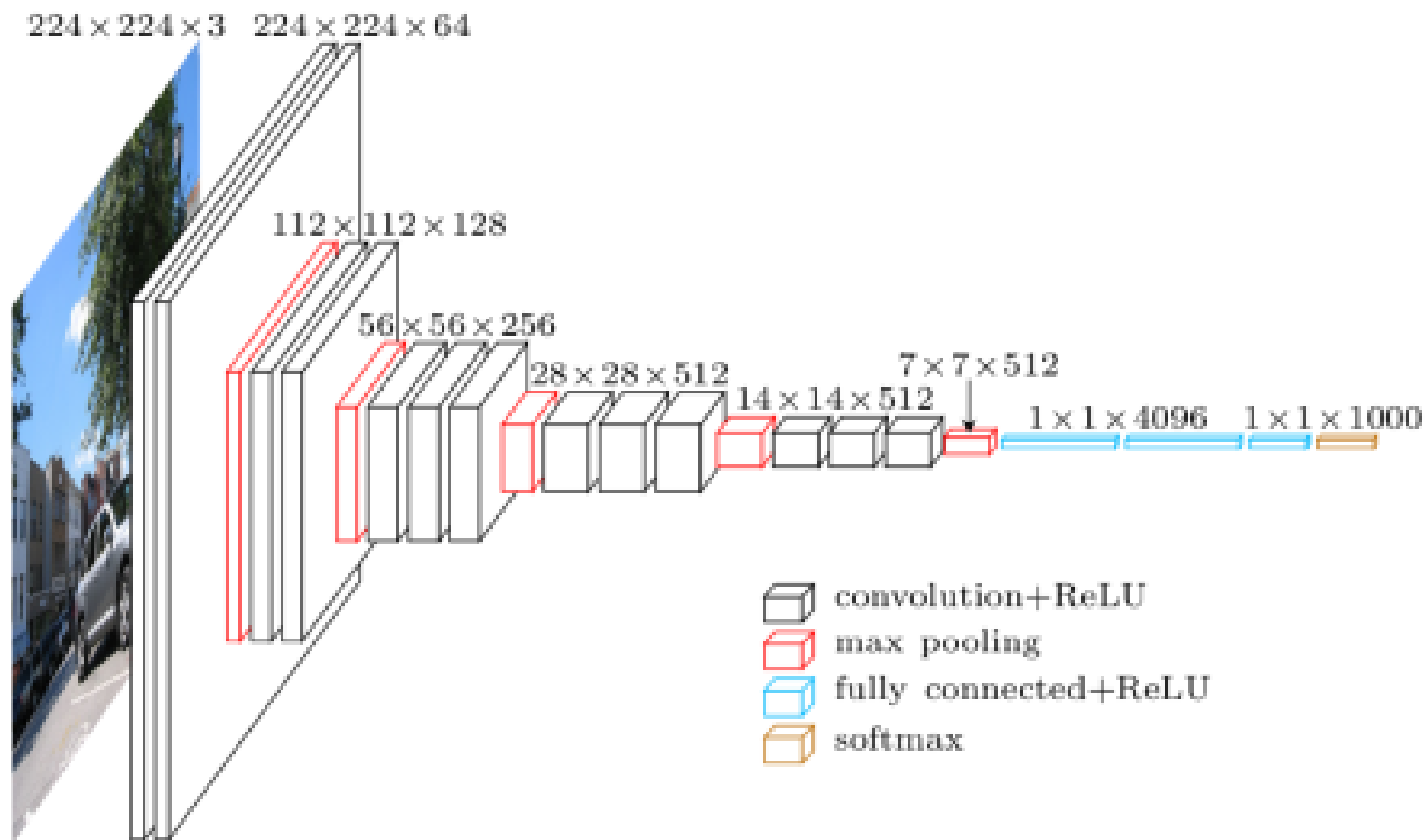
- The winner of ILSVRC 2012
- It changed the entire computer vision research



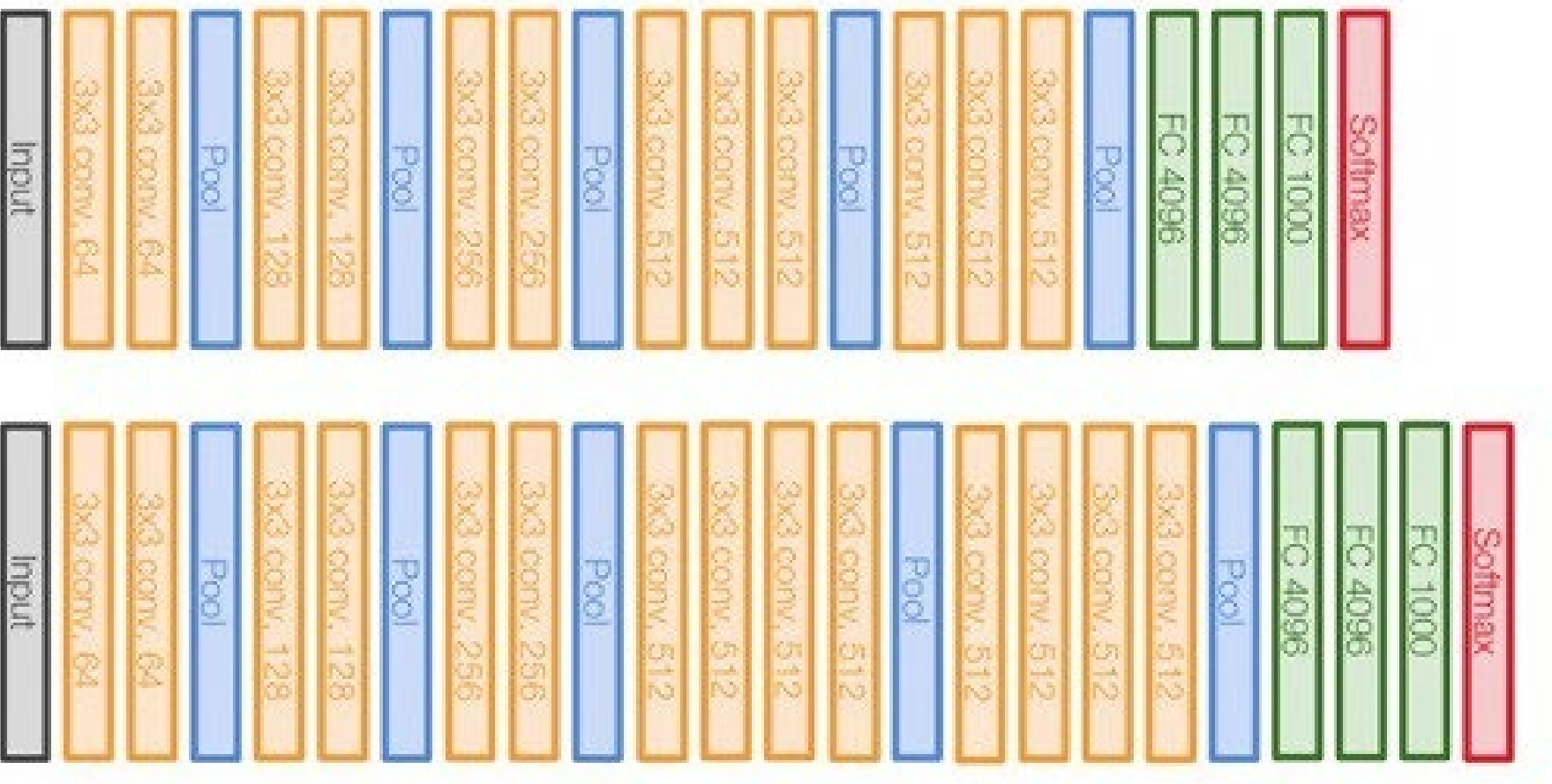
# ZFNet

- The winner of ILSVRC 2013
- The network architectures were developed by using the visualization techniques
  - Visualizing and Understanding Convolutional Networks, Zeiler et al, ECCV 2014
- Reduced the 1<sup>st</sup> layer filter size from 11x11 to 7x7
- 1<sup>st</sup> layer stride from 4 -> 2

# VGGNet



# VGGNet



# VGGNet

INPUT: [224x224x3] memory:  $224*224*3=150K$  params: 0 (not counting biases)

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*3)*64 = 1,728$

CONV3-64: [224x224x64] memory:  $224*224*64=3.2M$  params:  $(3*3*64)*64 = 36,864$

POOL2: [112x112x64] memory:  $112*112*64=800K$  params: 0

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*64)*128 = 73,728$

CONV3-128: [112x112x128] memory:  $112*112*128=1.6M$  params:  $(3*3*128)*128 = 147,456$

POOL2: [56x56x128] memory:  $56*56*128=400K$  params: 0

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*128)*256 = 294,912$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

CONV3-256: [56x56x256] memory:  $56*56*256=800K$  params:  $(3*3*256)*256 = 589,824$

POOL2: [28x28x256] memory:  $28*28*256=200K$  params: 0

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*256)*512 = 1,179,648$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [28x28x512] memory:  $28*28*512=400K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [14x14x512] memory:  $14*14*512=100K$  params: 0

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

CONV3-512: [14x14x512] memory:  $14*14*512=100K$  params:  $(3*3*512)*512 = 2,359,296$

POOL2: [7x7x512] memory:  $7*7*512=25K$  params: 0

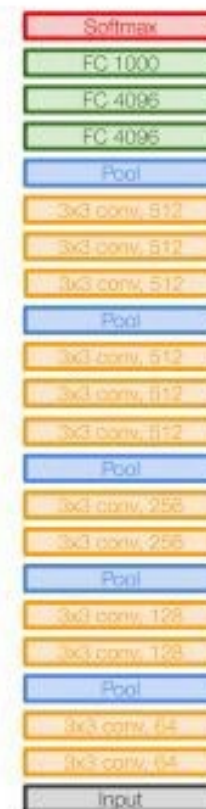
FC: [1x1x4096] memory: 4096 params:  $7*7*512*4096 = 102,760,448$

FC: [1x1x4096] memory: 4096 params:  $4096*4096 = 16,777,216$

FC: [1x1x1000] memory: 1000 params:  $4096*1000 = 4,096,000$

**TOTAL memory:**  $24M * 4 \text{ bytes} \approx 96MB$  / image (for a forward pass)

**TOTAL params:** 138M parameters



VGG16



# GoogLeNet

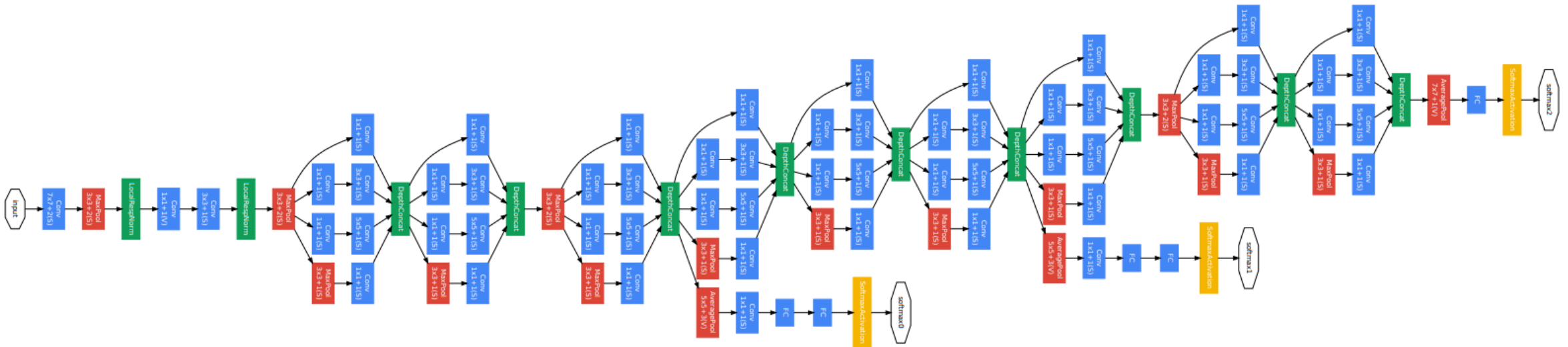
- Winner of ISLVR 2014
- Also called 'Inception'



Max pooling

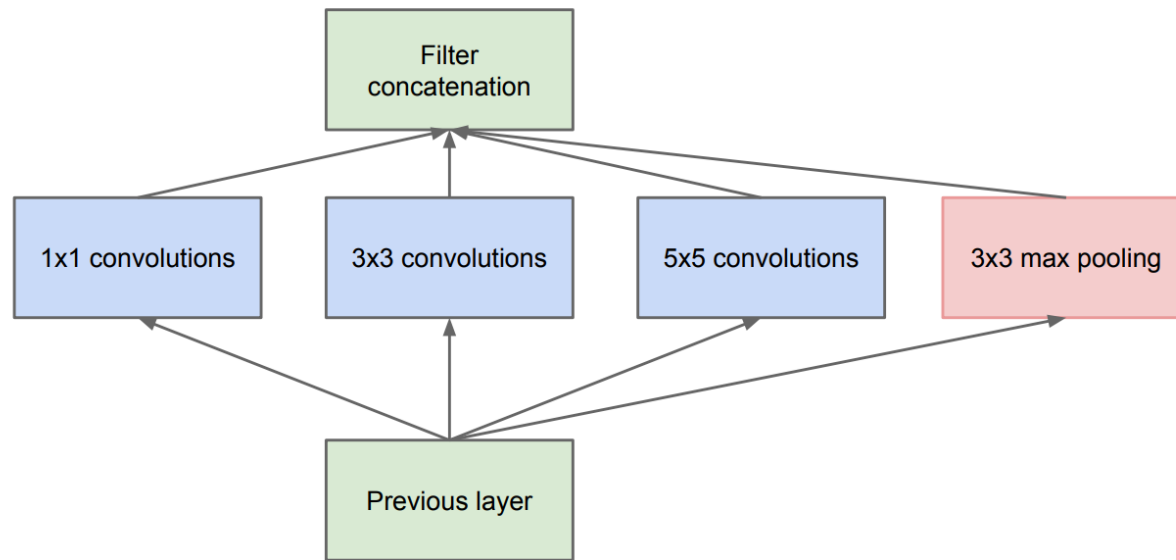
Concatenation

Convolution

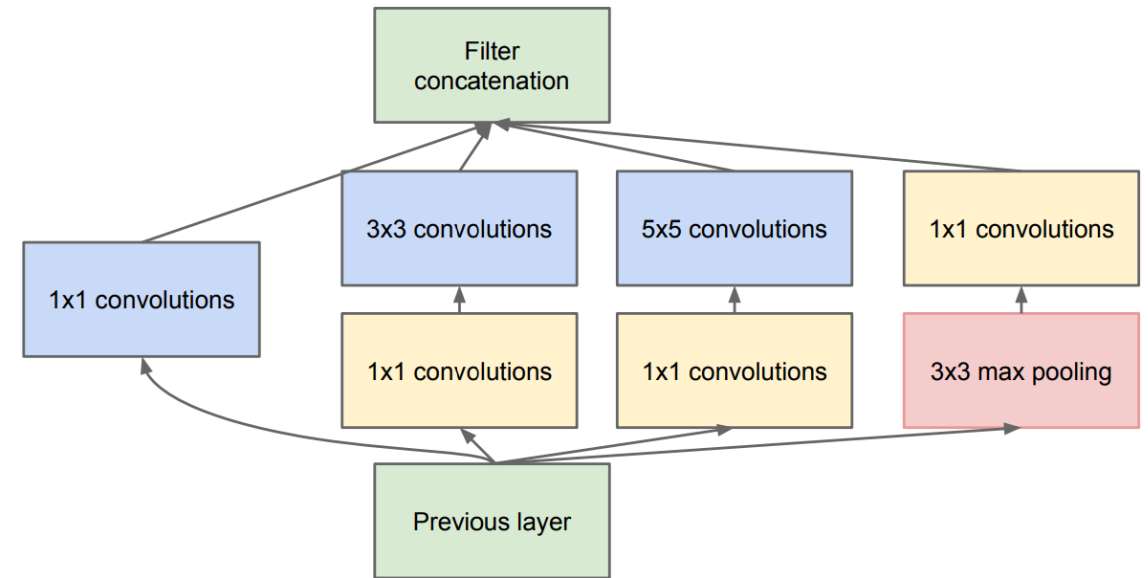


# GoogLeNet

- Inception module



(a) Inception module, naïve version



(b) Inception module with dimensionality reduction

# GoogLeNet

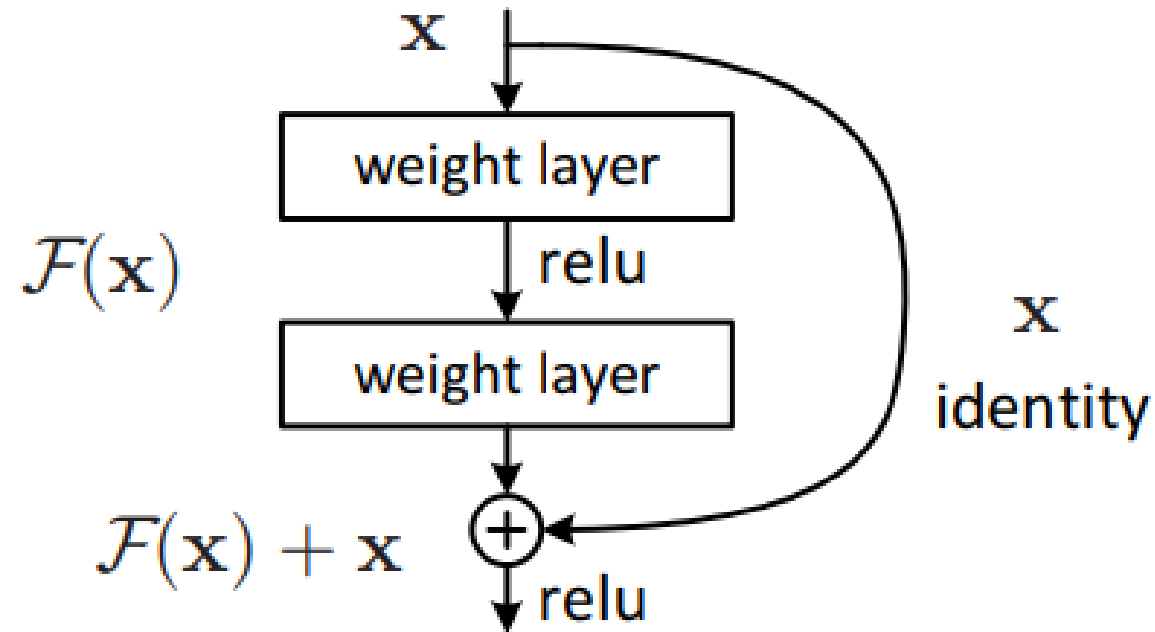
- ILSVRC 2014 classification results

Team	Year	Place	Error (top-5)	Uses external data
SuperVision	2012	1st	16.4%	no
SuperVision	2012	1st	15.3%	Imagenet 22k
Clarifai	2013	1st	11.7%	no
Clarifai	2013	1st	11.2%	Imagenet 22k
MSRA	2014	3rd	7.35%	no
VGG	2014	2nd	7.32%	no
GoogLeNet	2014	1st	6.67%	no

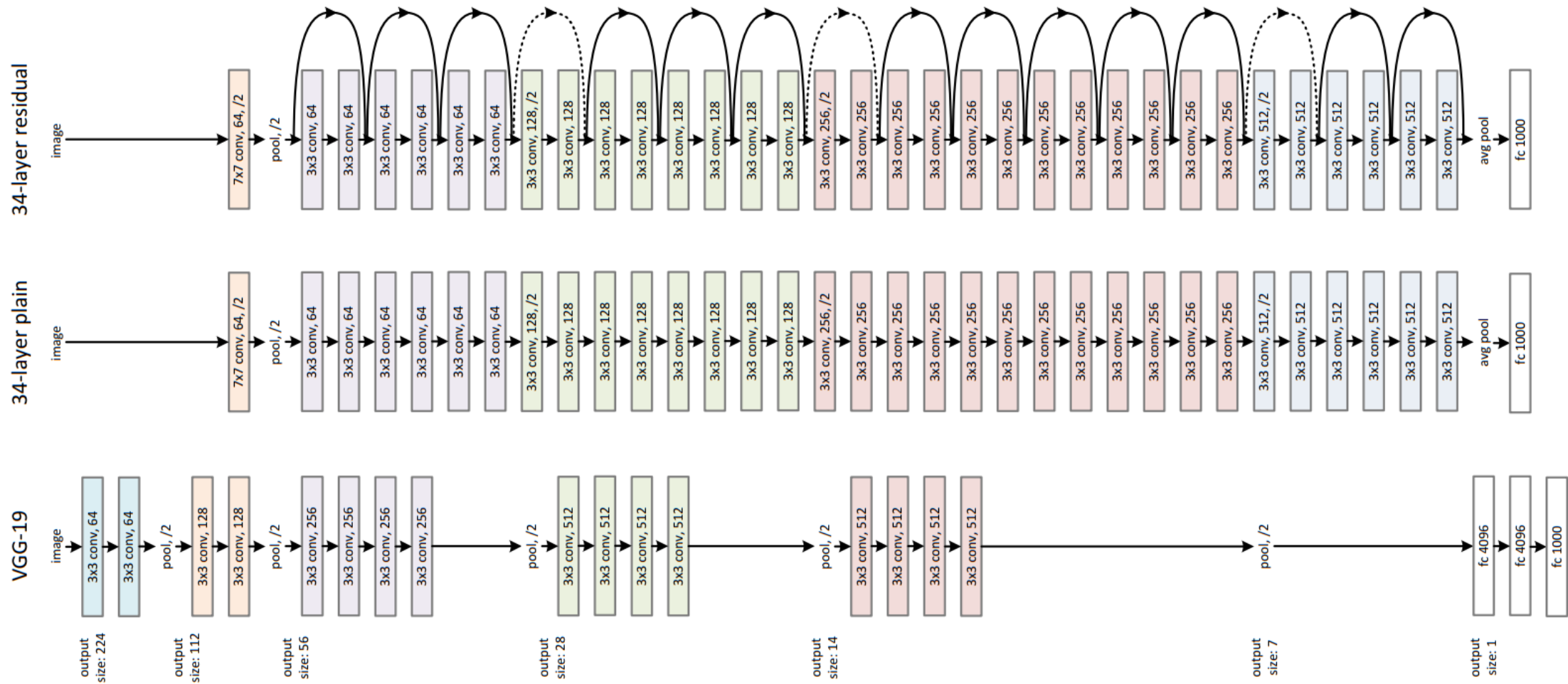


# ResNet

- The winner of ILSVRC 2015
- Residual building block



# ResNet



# ResNet

- Training on ImageNet

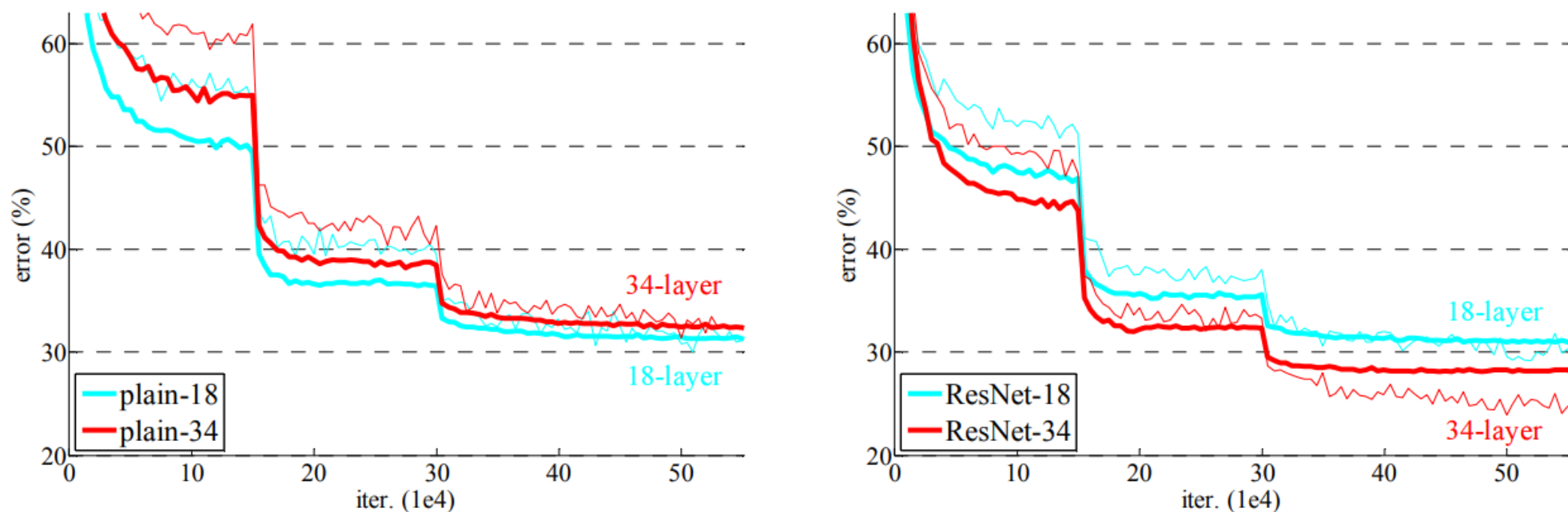


Figure 4. Training on **ImageNet**. Thin curves denote training error, and bold curves denote validation error of the center crops. Left: plain networks of 18 and 34 layers. Right: ResNets of 18 and 34 layers. In this plot, the residual networks have no extra parameter compared to their plain counterparts.

# ResNet

- ILSVRC 2015 classification results

method	top-5 err. ( <b>test</b> )
VGG [40] (ILSVRC'14)	7.32
GoogLeNet [43] (ILSVRC'14)	6.66
VGG [40] (v5)	6.8
PReLU-net [12]	4.94
BN-inception [16]	4.82
<b>ResNet (ILSVRC'15)</b>	<b>3.57</b>

Table 5. Error rates (%) of **ensembles**. The top-5 error is on the test set of ImageNet and reported by the test server.