

29-08-2024

# Convolutional neural network (CNN)

- Convolution operation

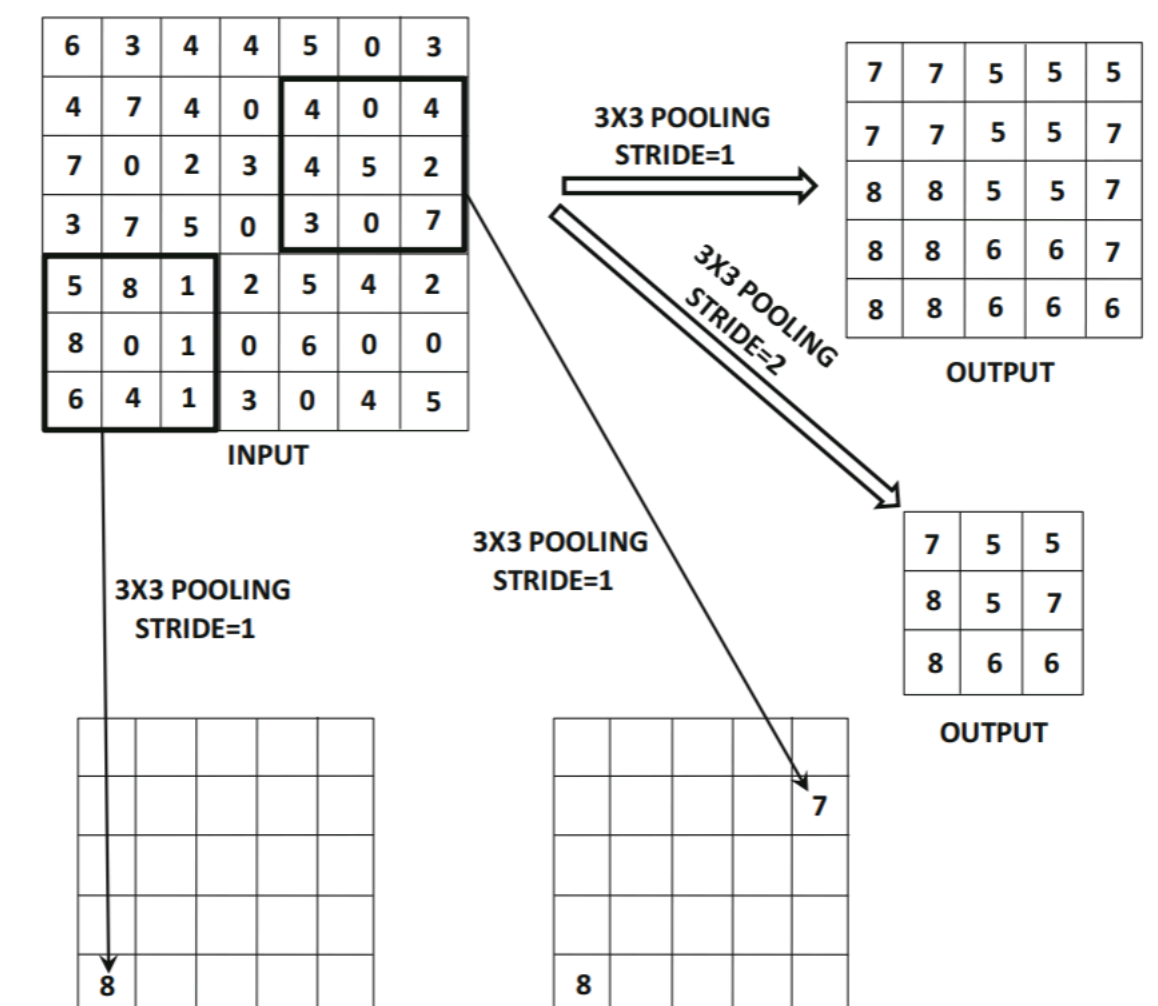
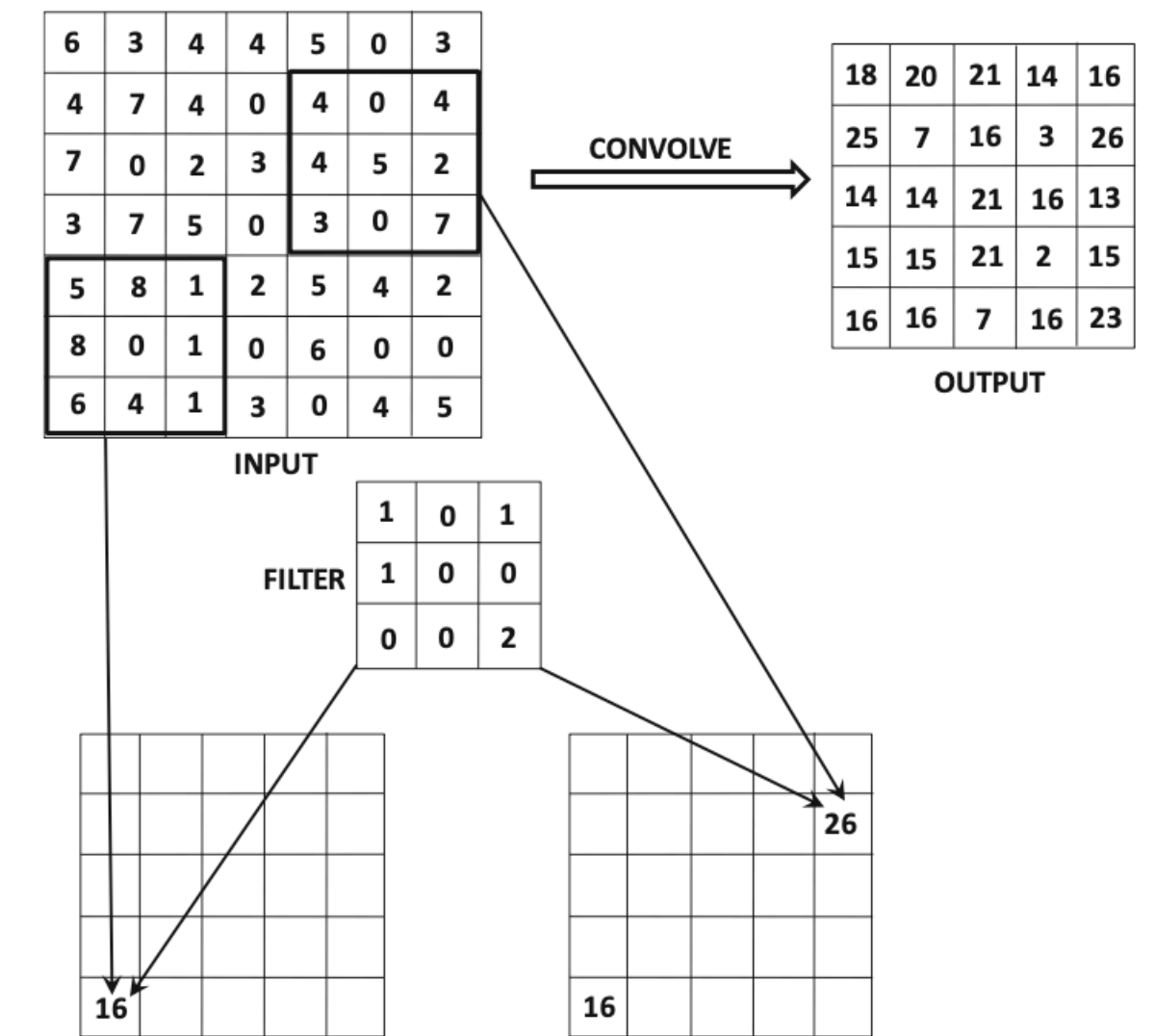
- Filters
- 2-D convolution?
- 3-D convolution ?
- Padding
  - Zero or duplicate border values
- Stride
  - Shift the convolution operations
  - 1, 2 not higher!
  - Depends on the filter size?

- Bias

- ReLU

- Pooling

- Max
- Average
- Simply down sample

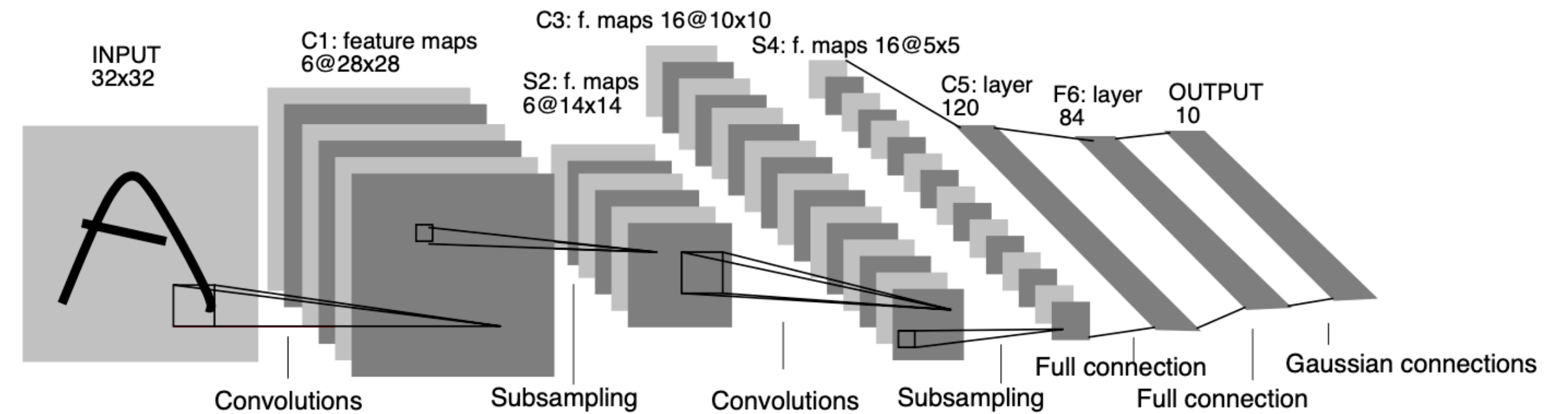


# CNN (Cont...)

- Convolution operation
- ReLU
- Pooling
- Fully connected layers
  - Multi-layer perceptron

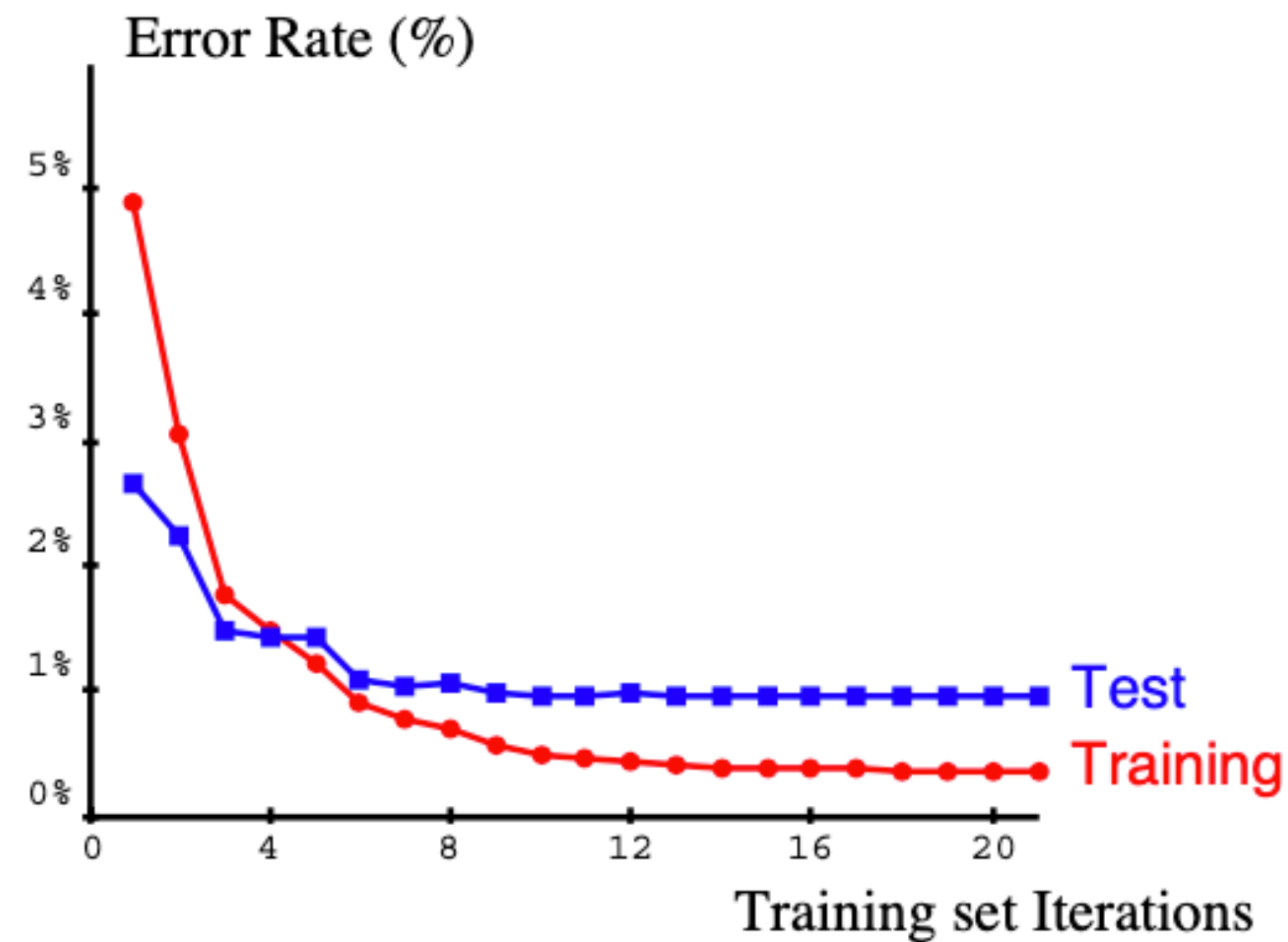
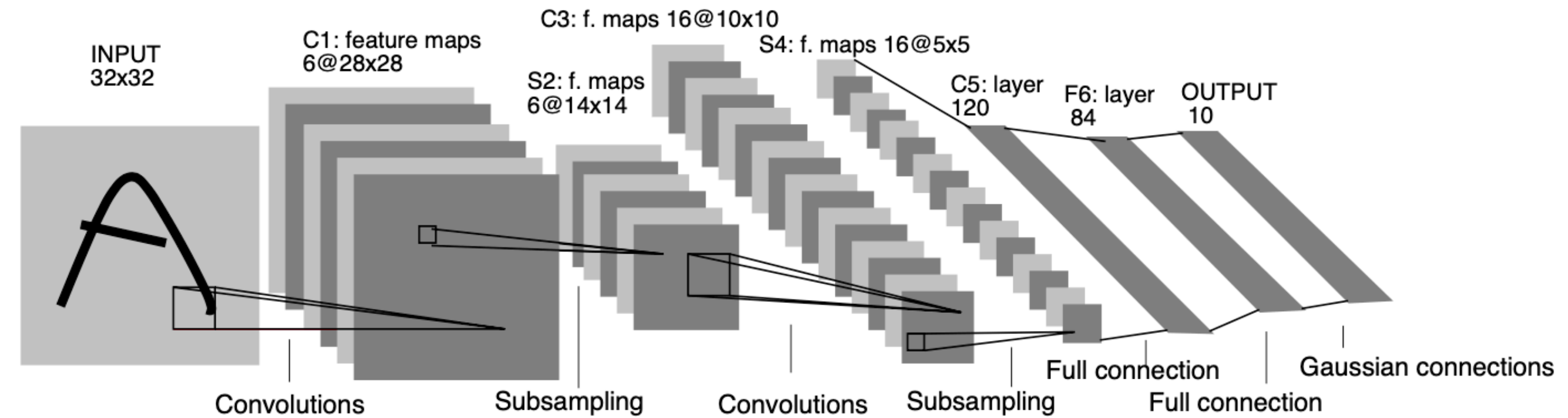
Feature  
detectors

Classifier



# LeNet-5

- LeNet-5: 1998
  - #parameters: 61,706



Network architecture (num channel=1, num class= 10) as follows:

```
LeNet5(  
  (features_detector): Sequential(  
    (0): Conv2d(1, 6, kernel_size=(5, 5), stride=(1, 1))  
    (1): AvgPool2d(kernel_size=2, stride=2, padding=0)  
    (2): Sigmoid()  
    (3): Conv2d(6, 16, kernel_size=(5, 5), stride=(1, 1))  
    (4): AvgPool2d(kernel_size=2, stride=2, padding=0)  
    (5): Sigmoid()  
  )  
  (classifier): Sequential(  
    (0): Linear(in_features=400, out_features=120, bias=True)  
    (1): Sigmoid()  
    (2): Linear(in_features=120, out_features=84, bias=True)  
    (3): Sigmoid()  
    (4): Linear(in_features=84, out_features=10, bias=True)  
  )  
)
```

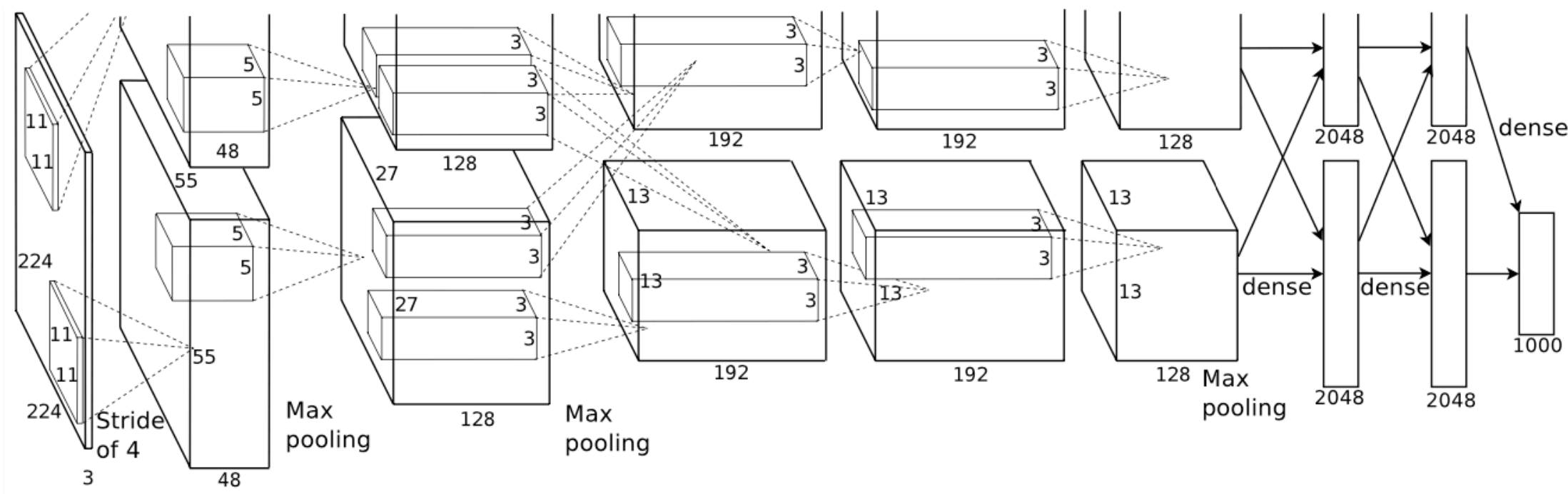
Network summary:

Layer (type)	Output Shape	Param #
Conv2d-1	[-1, 6, 28, 28]	156
AvgPool2d-2	[-1, 6, 14, 14]	0
Sigmoid-3	[-1, 6, 14, 14]	0
Conv2d-4	[-1, 16, 10, 10]	2,416
AvgPool2d-5	[-1, 16, 5, 5]	0
Sigmoid-6	[-1, 16, 5, 5]	0
Linear-7	[-1, 120]	48,120
Sigmoid-8	[-1, 120]	0
Linear-9	[-1, 84]	10,164
Sigmoid-10	[-1, 84]	0
Linear-11	[-1, 10]	850

Total params: 61,706

# AlexNet

- Breakthrough in ML, 2012
  - ▶ Error: 16.4
  - ▶ #parameters: 60,000,000



Network architecture (num channels=1, num classes= 10) as follows:

```
AlexNetR(
  (features): Sequential(
    (0): Conv2d(1, 64, kernel_size=(11, 11), stride=(4, 4), padding=(2, 2))
    (1): ReLU(inplace=True)
    (2): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1, ceil_mode=False)
    (3): Conv2d(64, 192, kernel_size=(5, 5), stride=(1, 1), padding=(2, 2))
    (4): ReLU(inplace=True)
    (5): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1, ceil_mode=False)
    (6): Conv2d(192, 384, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (7): ReLU(inplace=True)
    (8): Conv2d(384, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (9): ReLU(inplace=True)
    (10): Conv2d(256, 256, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1))
    (11): ReLU(inplace=True)
    (12): MaxPool2d(kernel_size=3, stride=2, padding=0, dilation=1, ceil_mode=False)
  )
  (avgpool): AdaptiveAvgPool2d(output_size=(6, 6))
  (classifier): Sequential(
    (0): Dropout(p=0.5, inplace=False)
    (1): Linear(in_features=9216, out_features=4096, bias=True)
    (2): ReLU(inplace=True)
    (3): Dropout(p=0.5, inplace=False)
    (4): Linear(in_features=4096, out_features=4096, bias=True)
    (5): ReLU(inplace=True)
    (6): Linear(in_features=4096, out_features=10, bias=True)
  )
)
```

Network summary:

Layer (type:depth-idx)	Output Shape	Param #
Sequential: 1-1	[-1, 256, 6, 6]	--
└─Conv2d: 2-1	[-1, 64, 56, 56]	7,808
└─ReLU: 2-2	[-1, 64, 56, 56]	--
└─MaxPool2d: 2-3	[-1, 64, 27, 27]	--
└─Conv2d: 2-4	[-1, 192, 27, 27]	307,392
└─ReLU: 2-5	[-1, 192, 27, 27]	--
└─MaxPool2d: 2-6	[-1, 192, 13, 13]	--
└─Conv2d: 2-7	[-1, 384, 13, 13]	663,936
└─ReLU: 2-8	[-1, 384, 13, 13]	--
└─Conv2d: 2-9	[-1, 256, 13, 13]	884,992
└─ReLU: 2-10	[-1, 256, 13, 13]	--
└─Conv2d: 2-11	[-1, 256, 13, 13]	590,080
└─ReLU: 2-12	[-1, 256, 13, 13]	--
└─MaxPool2d: 2-13	[-1, 256, 6, 6]	--
└─AdaptiveAvgPool2d: 1-2	[-1, 256, 6, 6]	--
Sequential: 1-3	[-1, 10]	--
└─Dropout: 2-14	[-1, 9216]	--
└─Linear: 2-15	[-1, 4096]	37,752,832
└─ReLU: 2-16	[-1, 4096]	--
└─Dropout: 2-17	[-1, 4096]	--
└─Linear: 2-18	[-1, 4096]	16,781,312
└─ReLU: 2-19	[-1, 4096]	--
└─Linear: 2-20	[-1, 10]	40,970

Total params: 57,029,322

# VggNet

- VggNet: 2015
  - Error: 7.3%
  - #parameters: 140,000,000

ConvNet Configuration					
A	A-LRN	B	C	D	E
11 weight layers	11 weight layers	13 weight layers	16 weight layers	16 weight layers	19 weight layers
input (224 × 224 RGB image)					
conv3-64	conv3-64 <b>LRN</b>	conv3-64 <b>conv3-64</b>	conv3-64 conv3-64	conv3-64 conv3-64	conv3-64 conv3-64
maxpool					
conv3-128	conv3-128	conv3-128 <b>conv3-128</b>	conv3-128 conv3-128	conv3-128 conv3-128	conv3-128 conv3-128
maxpool					
conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256	conv3-256 conv3-256 <b>conv1-256</b>	conv3-256 conv3-256 <b>conv3-256</b>	conv3-256 conv3-256 conv3-256 <b>conv3-256</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512	conv3-512 conv3-512 <b>conv1-512</b>	conv3-512 conv3-512 <b>conv3-512</b>	conv3-512 conv3-512 conv3-512 <b>conv3-512</b>
maxpool					
FC-4096					
FC-4096					
FC-1000					
soft-max					

Source: Karen Simonyan\* & Andrew Zisserman, ICLR, 2015



# ResNet

- ResNet: 2015
  - Error: 3.6%
  - #parameters: 58,161,162

layer name	output size	18-layer	34-layer	50-layer	101-layer	152-layer
conv1	112×112	7×7, 64, stride 2				
conv2_x	56×56	3×3 max pool, stride 2				
		$\begin{bmatrix} 3\times 3, 64 \\ 3\times 3, 64 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 64 \\ 3\times 3, 64 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 64 \\ 3\times 3, 64 \\ 1\times 1, 256 \end{bmatrix} \times 3$
conv3_x	28×28	$\begin{bmatrix} 3\times 3, 128 \\ 3\times 3, 128 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 128 \\ 3\times 3, 128 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 4$	$\begin{bmatrix} 1\times 1, 128 \\ 3\times 3, 128 \\ 1\times 1, 512 \end{bmatrix} \times 8$
conv4_x	14×14	$\begin{bmatrix} 3\times 3, 256 \\ 3\times 3, 256 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 256 \\ 3\times 3, 256 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 6$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 23$	$\begin{bmatrix} 1\times 1, 256 \\ 3\times 3, 256 \\ 1\times 1, 1024 \end{bmatrix} \times 36$
conv5_x	7×7	$\begin{bmatrix} 3\times 3, 512 \\ 3\times 3, 512 \end{bmatrix} \times 2$	$\begin{bmatrix} 3\times 3, 512 \\ 3\times 3, 512 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$	$\begin{bmatrix} 1\times 1, 512 \\ 3\times 3, 512 \\ 1\times 1, 2048 \end{bmatrix} \times 3$
	1×1	average pool, 1000-d fc, softmax				
FLOPs		$1.8 \times 10^9$	$3.6 \times 10^9$	$3.8 \times 10^9$	$7.6 \times 10^9$	$11.3 \times 10^9$

ResNet Architectures

Source: Hei et al., Deep Residual Learning for Image Recognition, CVPR 2016

# How to train a CNN ?

- What is different from our MLP training?
- Let's consider the operations in different layers?
  - ▶ Convolution
  - ▶ ReLU
  - ▶ Pooling
  - ▶ Fully connected layers