

Convolution Neural Network (CNN)

- **Local details in CNN**
- **Properties of visual structure**
- **Layers of CNN**
- **Weights updating (Backward propagation)**
- **Examples**

CIFAR: from Wikipedia

In 2004, [Geoffrey Hinton](#) began leading CIFAR's Neural Computation & Adaptive Perception program. Its members included [Yoshua Bengio](#) and [Yann LeCun](#), among other neuroscientists, computer scientists, biologists, electrical engineers, physicists, and psychologists. Together, they confirmed Hinton's conviction about the power of neural networks when they created computing systems that mimicked human intelligence. Today, the three are widely acknowledged as the pioneers of [deep learning](#). In 2019, the [Association for Computing Machinery](#) (ACM) named Hinton, Bengio and LeCun as recipients of the 2018 [ACM A.M. Turing Award](#) for conceptual and engineering breakthroughs that have made deep [neural networks](#) a critical component of computing.^{[\[13\]](#)}

CIFAR-10:

A dataset contains 60,000 color images of 32x32 pixels in 3 channels divided into 10 classes. Each class contains 6,000 Images. The training set contains 50,000 images, while the test set provides 10,000 images.

Deep learning: In 2012, Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton published an article titled “ImageNet Classification with Deep Convolutional Neural Networks” in the Proceedings of Neural Information Processing Systems(NIPS). At the end of the paper, they wrote:

“It is notable that our network’s performance degrades if a single convolutional layer is removed. For example, removing any of the middle layers results in a loss of about 2% for the top-1 performance of the network. So the depth really is important for achieving our results.”

airplane



automobile



bird



cat



deer



dog



frog



horse



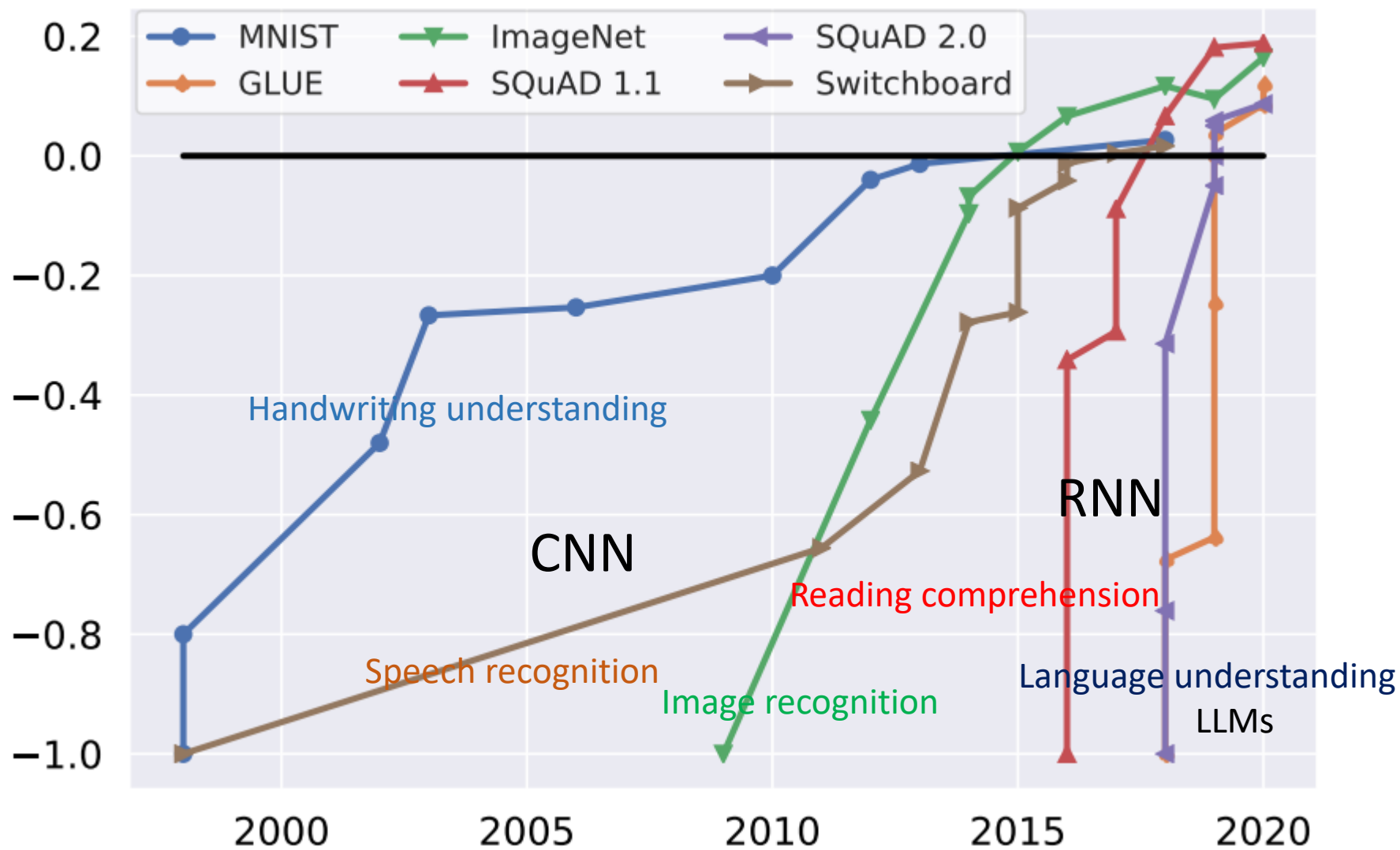
ship



truck



Advances in Deep Learning



Source: Dynabench: Rethinking Benchmarking in NLP, 2021

Advances of CNN on ImageNet

•

AlexNet-8 — Error rate: 16.4 %

VGGNet-19 — Error rate: 7.3 %

GoogleNet-22 — Error rate: 6.7 %

ResNet-152 — Error rate: 3.6 %

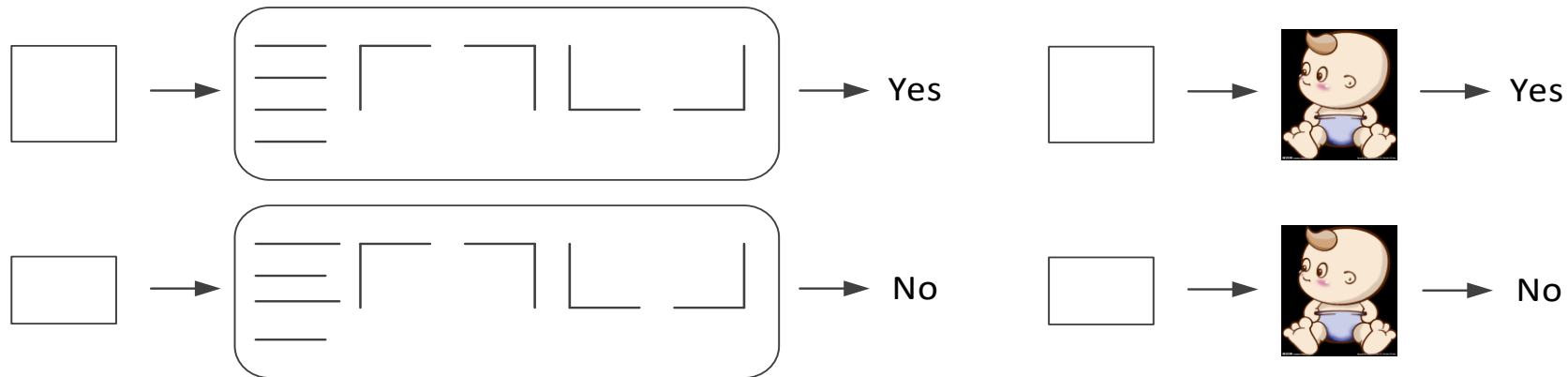
Error Rates



What is CNN?

If we want to judge whether a quadrangle is square or not, the rationale approach is to seek features of square, such as same length for four sides and four 90 degree corner angles

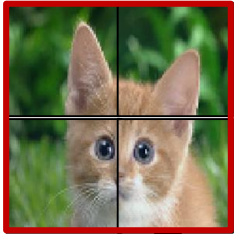
Use filters to extract features, and use features for classification.



Properties of visual structure

Local Processing: pixels close together go together

receptive fields capture local detail



Across Space: the same what, no matter where

recognize the same pattern in different places, translation invariant



Properties of visual structure- subsampling

Why CNN for Image

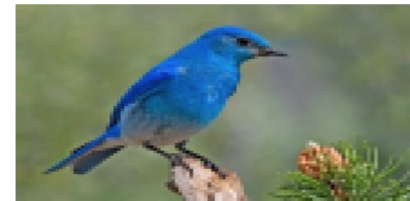
- Subsampling the pixels will not change the object

bird



subsampling

bird



We can subsample the pixels to make image smaller



Less parameters for the network to process the image

Basic Concept of Convolution and Polling in CNN Operations

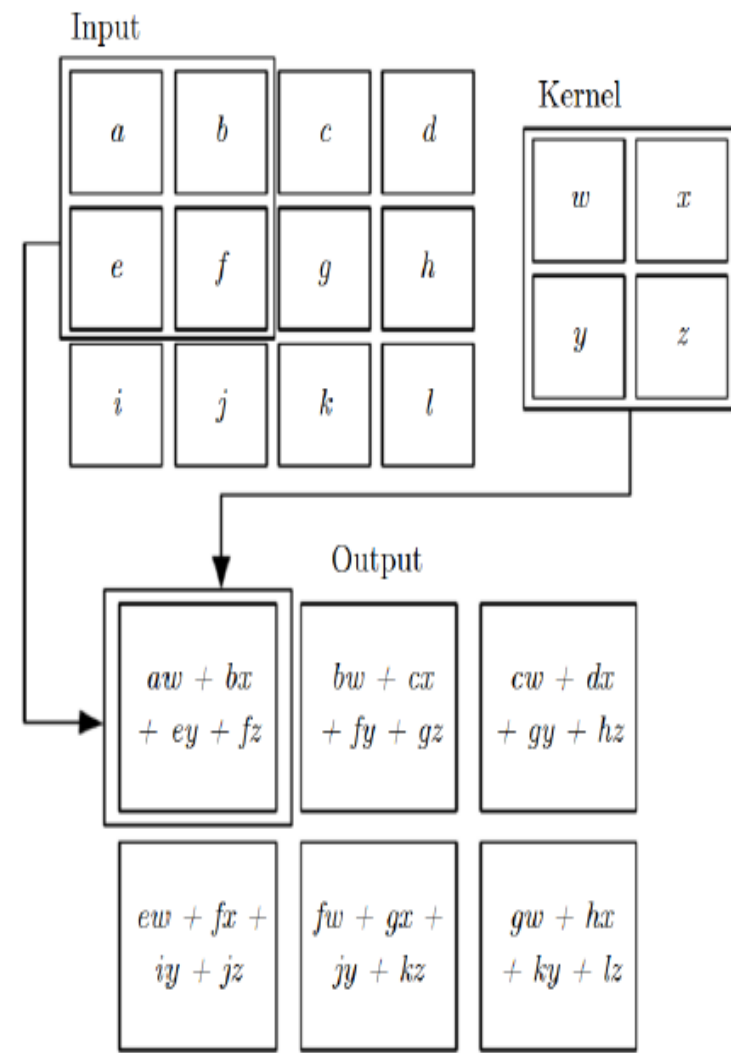
- Consider an image of 500x500 pixels. The no. of neurons in input layer is set to 500x500 with 10^8 neurons in hidden layer. Each connection between one input neuron and one neuron in hidden layer is set with a weight parameter
- The no.of weight parameters between input layer and hidden layer is $500 \times 500 \times 10^8 = 25 \times 10^{12}$. This may demand enormous amount of computations
- Design a 10x10 filter to extract the local features in an image. Then a hidden layer neuron is connected to a 10x10 area of the image through the filter

Properties of CNN

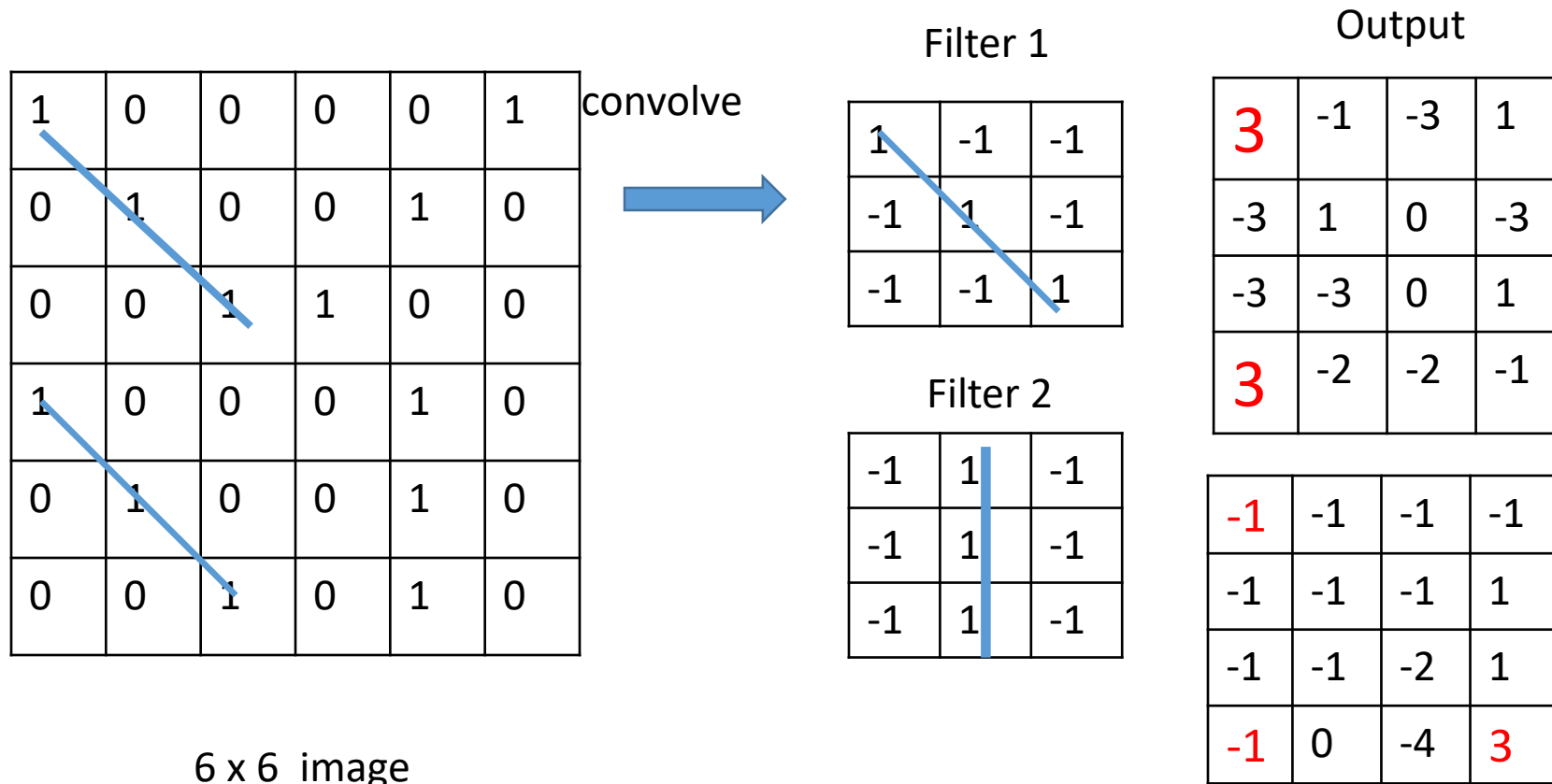
- Patterns are much smaller than the whole image; extract features for classification.
- The same pattern appears in different regions; convolution detects patterns in different regions of an image
- Weights sharing of neurons to reduce the number of neurons in CNN, (same filter to extract features at different locations)
- Use subsampling to reduce the amount of features (features in a bigger range) and the computational complexity of the network.

2-D Convolution (without kernel flipping)

Example of 'valid' 2-D convolution (without kernel flipping) where a 3x4 matrix convolved with a 2x2 kernel to output a 2x3 matrix



CNN Property 1



Same patterns produce the largest outputs

CNN – Convolution

-1	1	-1
-1	1	-1
-1	1	-1

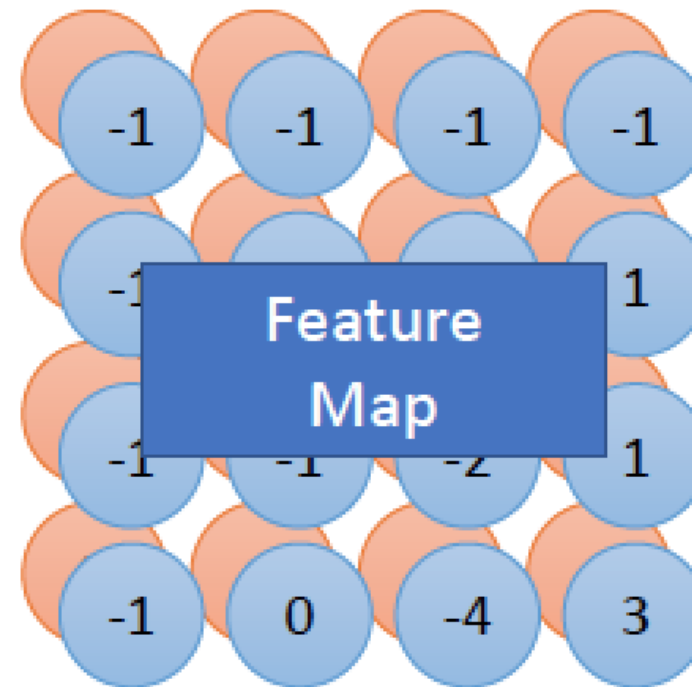
Filter 2

stride=1

1	0	0	0	0	1
0	1	0	0	1	0
0	0	1	1	0	0
1	0	0	0	1	0
0	1	0	0	1	0
0	0	1	0	1	0

6 x 6 image

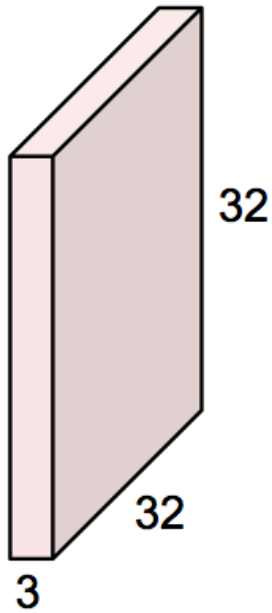
Do the same process for every filter



4 x 4 image

A Filter

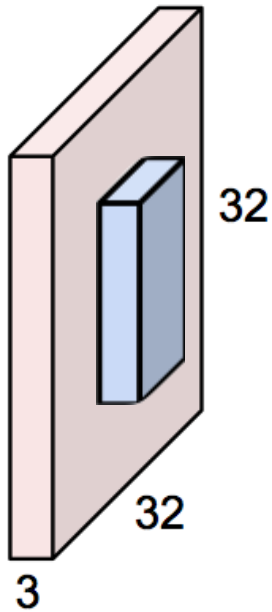
A filter is a **spatially local** and **cross-channel** template
Convnet filters are **learned**



input is **3x32x32** data
a color image (3 RGB channels) and square (32x32)

A Filter

A filter is a **spatially local** and **cross-channel** template
Convnet filters are **learned**



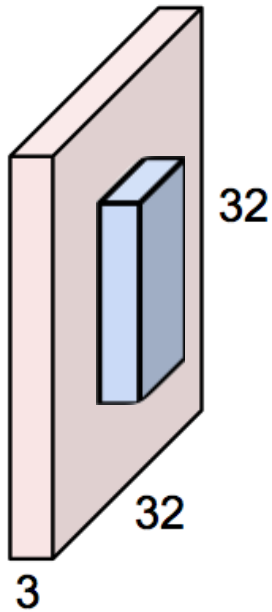
input is **3x32x32** data
a color image (3 RGB channels) and square (32x32)

filter is **3x5x5** weights

- spatially local: kernel size is 5x5
- cross-channel: connected across all input channels

A Filter

A filter is a **spatially local** and **cross-channel** template
Convnet filters are **learned**



input is **3x32x32** data
a color image (3 RGB channels) and square (32x32)

filter is **3x5x5** weights

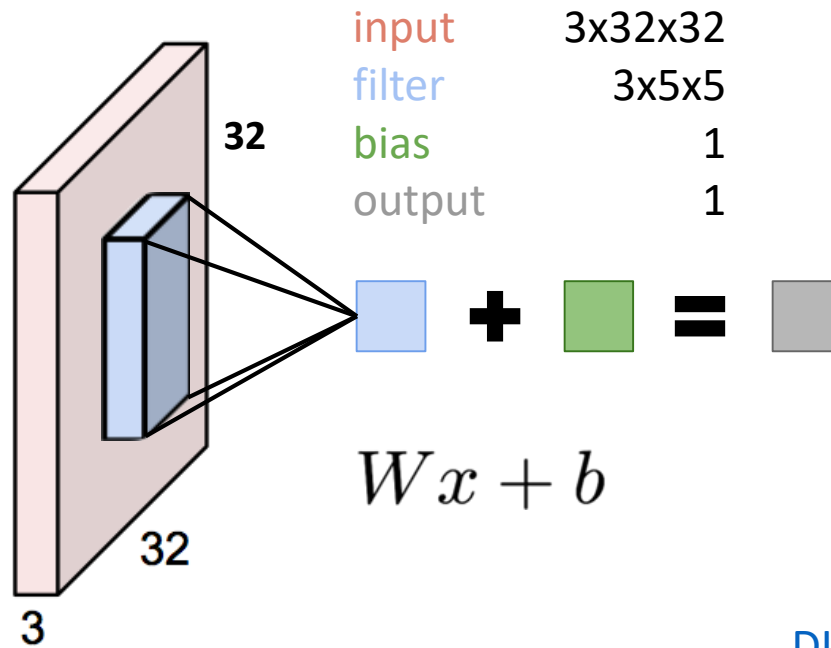
- spatially local: kernel size is 5x5
- cross-channel: connected across all input channels

total parameters:

$$3 \cdot 5^2 = 75 \text{ filter weights} + 1 \text{ bias}$$

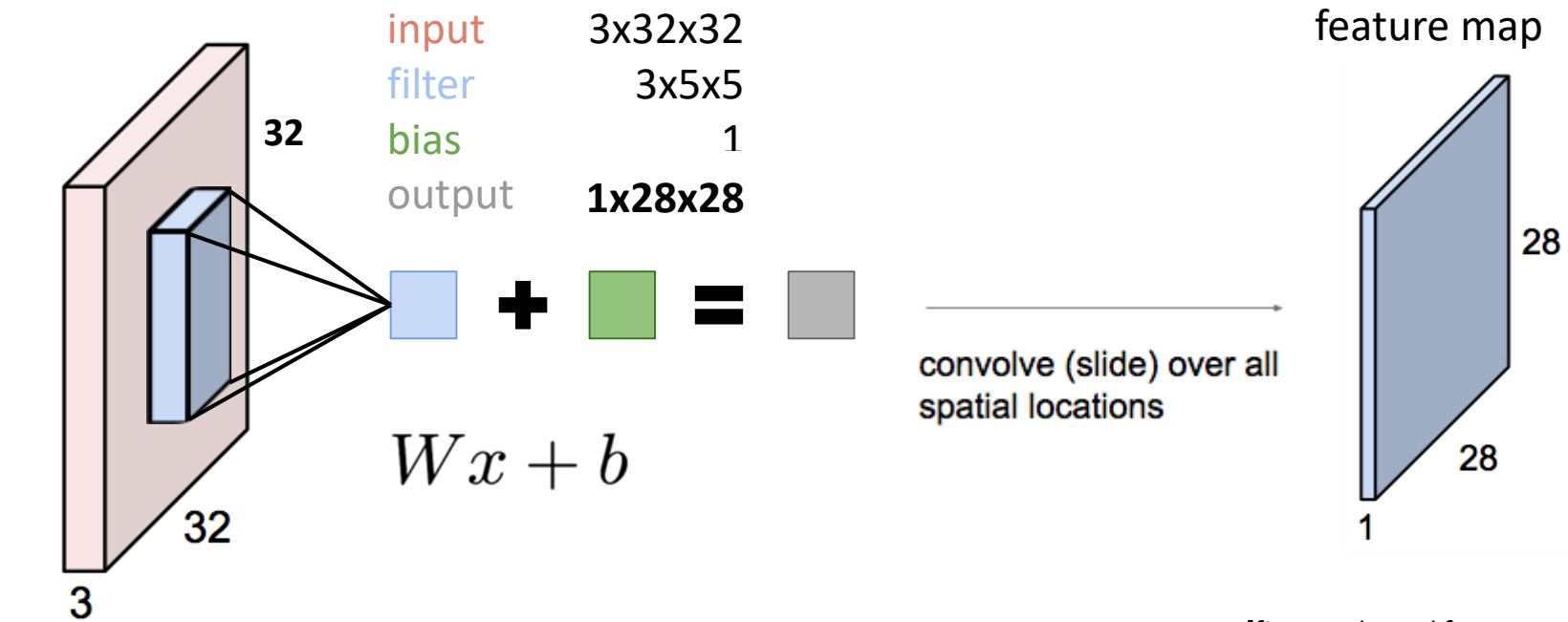
Convolution

One filter evaluation is a dot product between the input window and weights + bias



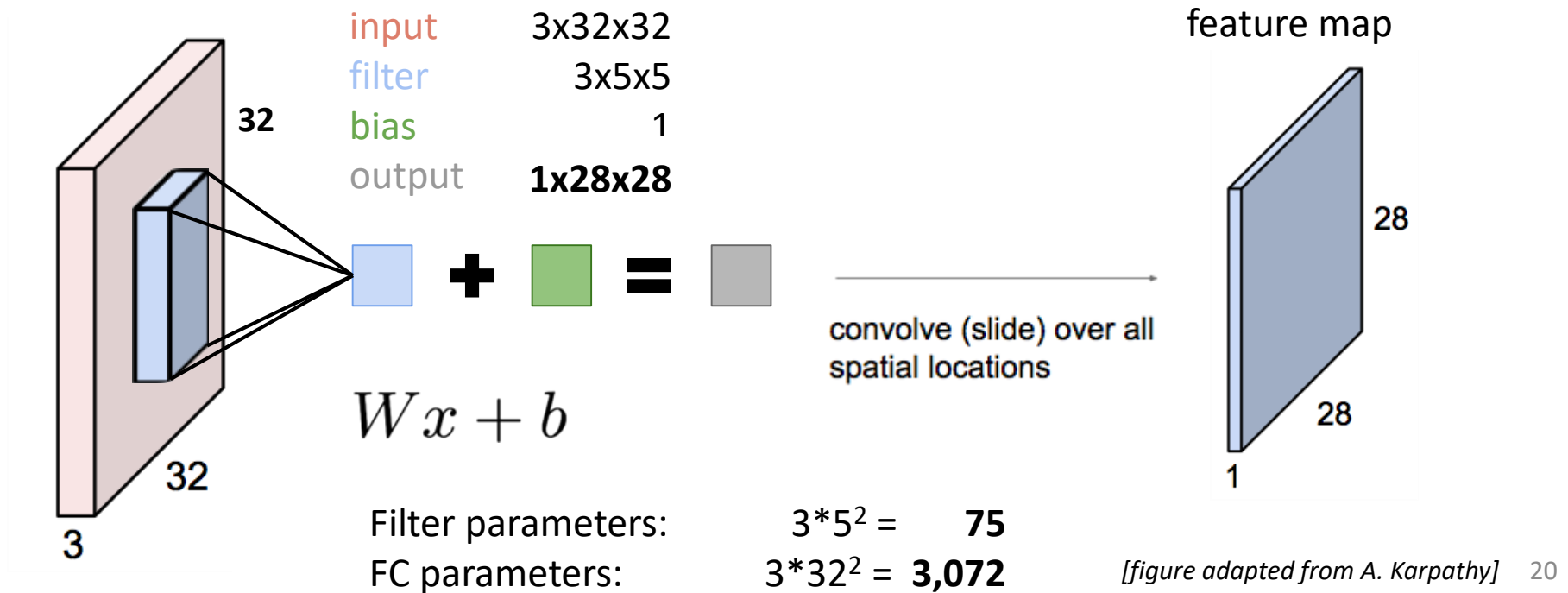
Convolution

Convolving the filter with the input gives a **feature map**.

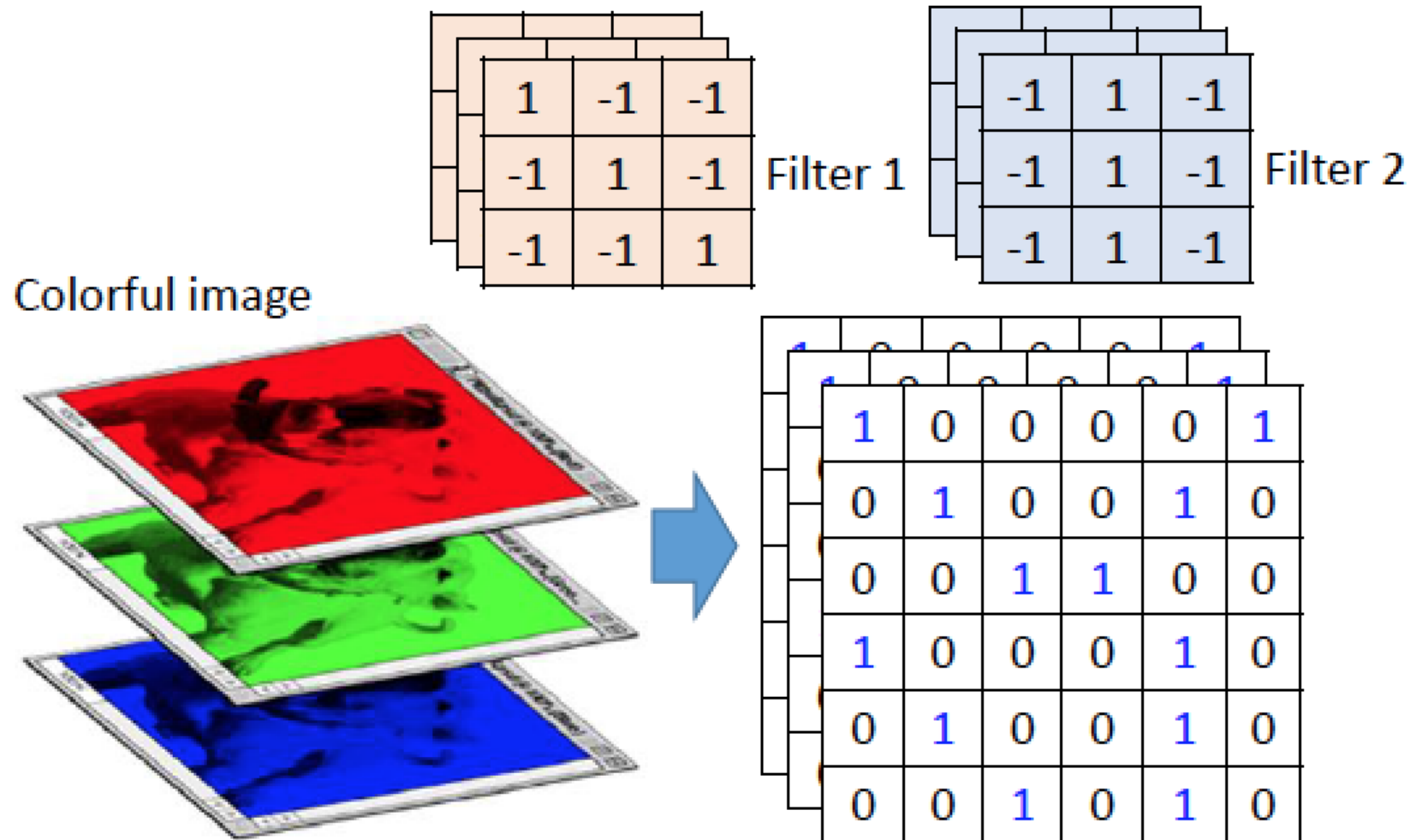


Convolution

Convolving the filter with the input gives a **feature map**.



CNN – Colorful image



Input Volume (+pad 1) (7x7x3)

x[:, :, 0]						
0	0	0	0	0	0	0
0	2	2	0	2	2	0
0	2	1	1	1	0	0
0	2	2	2	0	1	0
0	1	2	1	0	2	0
0	1	0	2	2	0	0
0	0	0	0	0	0	0
x[:, :, 1]						
0	0	0	0	0	0	0
0	0	0	2	0	1	0
0	0	0	0	1	0	0
0	2	1	0	2	1	0
0	2	0	2	0	0	0
0	1	0	1	2	1	0
0	0	0	0	0	0	0
x[:, :, 2]						
0	0	0	0	0	0	0
0	0	1	2	0	0	0
0	0	0	1	1	0	0
0	1	2	2	0	1	0
0	1	1	1	1	1	0
0	0	0	2	2	2	0
0	0	0	0	0	0	0

Filter W0 (3x3x3)

w0[:, :, 0]		
0	0	1
-1	0	-1
1	-1	0
w0[:, :, 1]		
-1	0	1
-1	-1	1
-1	1	-1
w0[:, :, 2]		
0	-1	1
0	0	0
1	0	-1
Bias b0 (1x1x1)		
b0[:, :, 0]		
1		

Filter W1 (3x3x3)

w1[:, :, 0]		
1	0	1
1	0	1
0	1	0
w1[:, :, 1]		
-1	0	1
1	0	0
-1	1	-1
w1[:, :, 2]		
0	-1	-1
1	0	1
0	0	1
Bias b1 (1x1x1)		
b1[:, :, 0]		
0		

Output Volume (3x3x2)

o[:, :, 0]		
-3	-7	-1
-1	5	-4
2	0	-5
o[:, :, 1]		
5	6	1
9	10	4
0	4	5

X

$InputShape = (7, 7, 3)$

$strid = 1$

$pad = 1$

W

$FilterShape = (3, 3, 3, 2)$

(inChannel, Shape, Shape, outChannel)

$bShape = (1, 1, 1)$

Y

$OutputShape = (3, 3, 2)$

FeatureMap

$$Channel1 = 1 - 2 + 1 = 0$$

$$Channel2 = 1 + 1 + 2 = 4$$

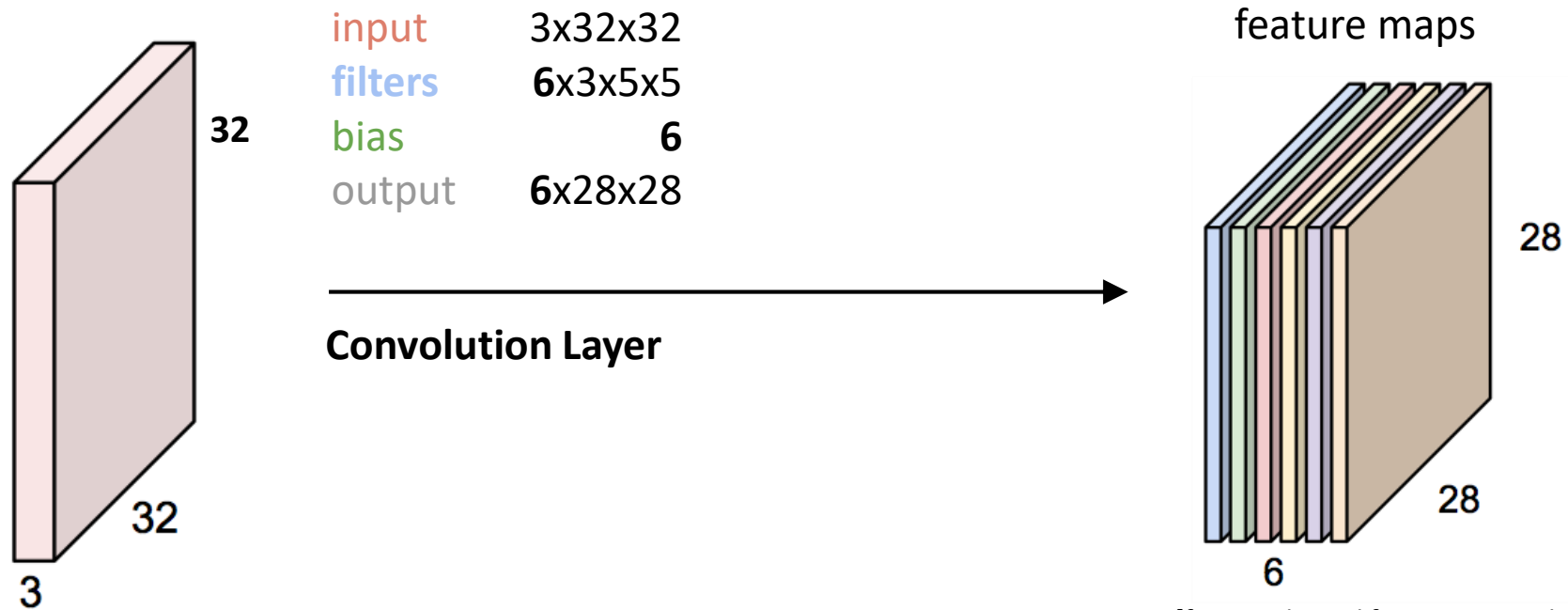
$$Channel3 = 1 + 0 + 0 = 1$$

$$b = 1$$

$$res = \sum(Channel_i) + b = 5$$

Convolution Layer (conv)

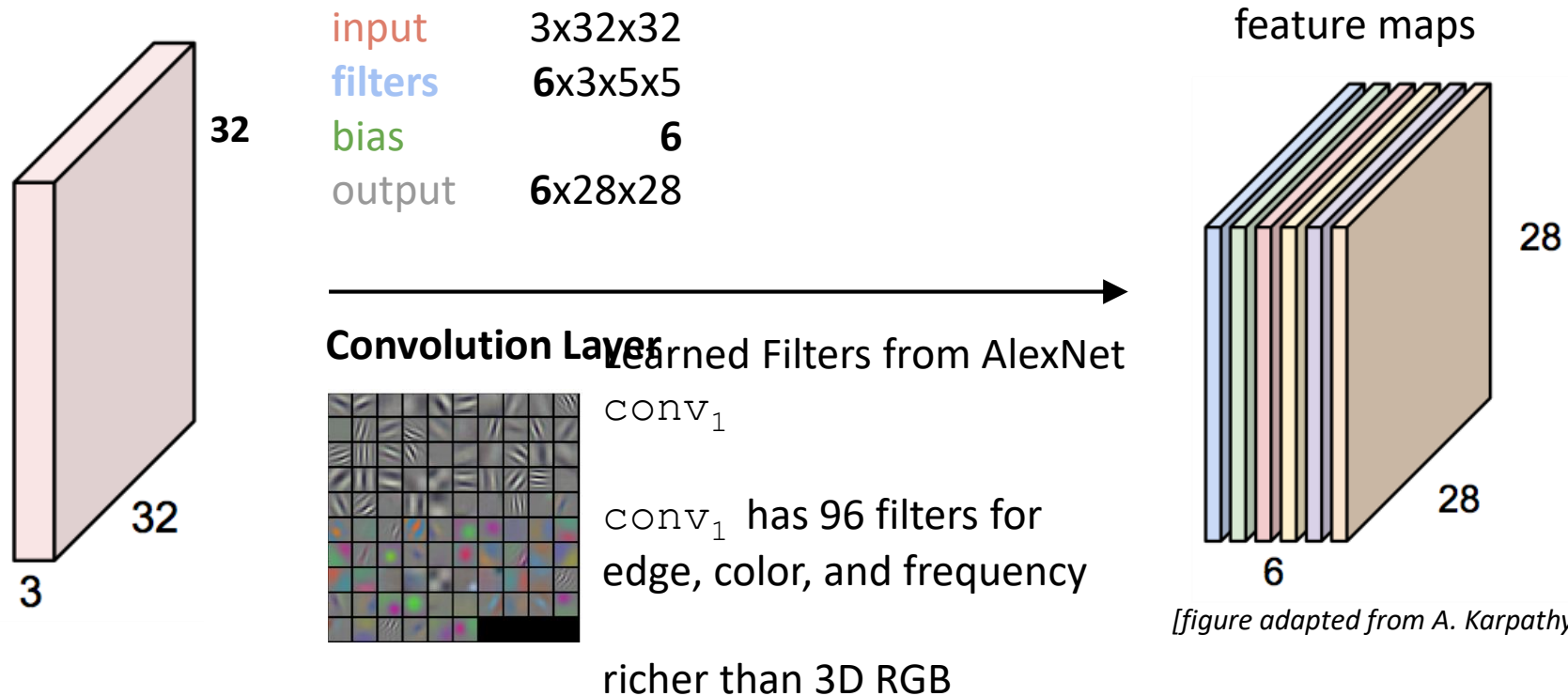
Convolution layers have multiple **filters** for more modeling capacity



[figure adapted from A. Karpathy] 23

Convolution Layer (conv)

Convolution layers have multiple **filters** for more modeling capacity



$m \times m$ image, $k \times k$ filter

Padding

Full

- Add zero-padding to the image enough for every pixel to be visited k times in each direction, with output size: $(m + k - 1) \times (m + k - 1)$

Valid

- With no zero-padding, kernel is restricted to traverse only within the image, with output size: $(m - k + 1) \times (m - k + 1)$

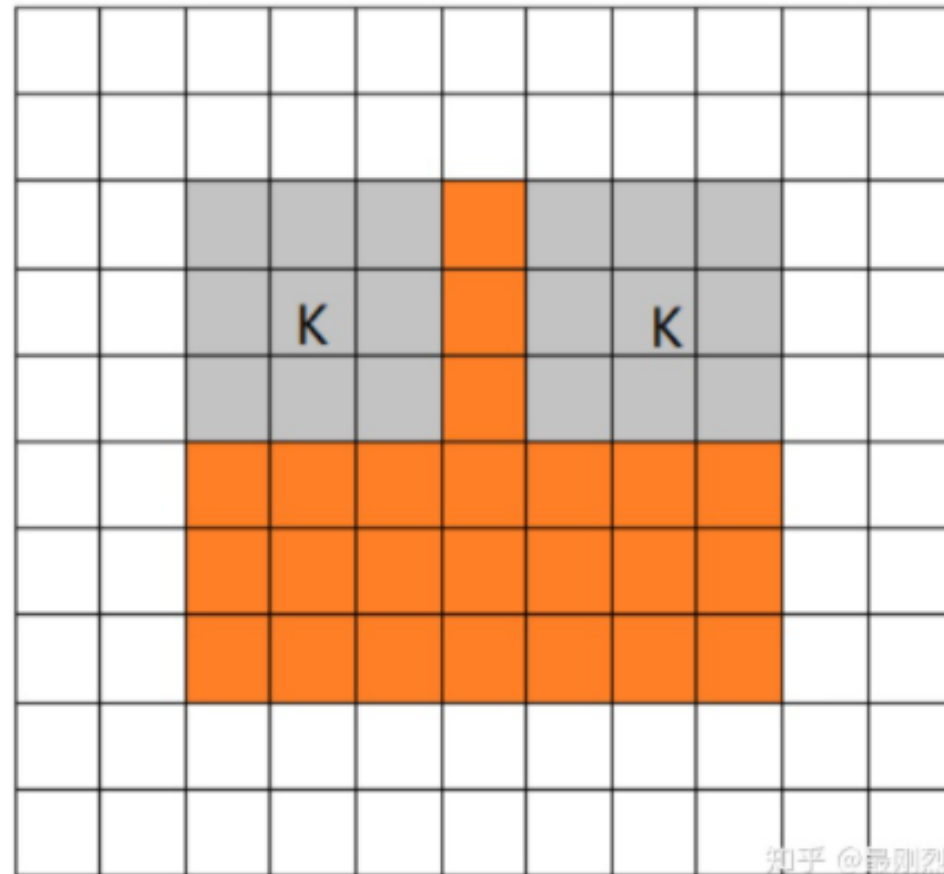
Same

- Add zero-padding to the image to have the output of the same size as the image, i.e., $m \times m$

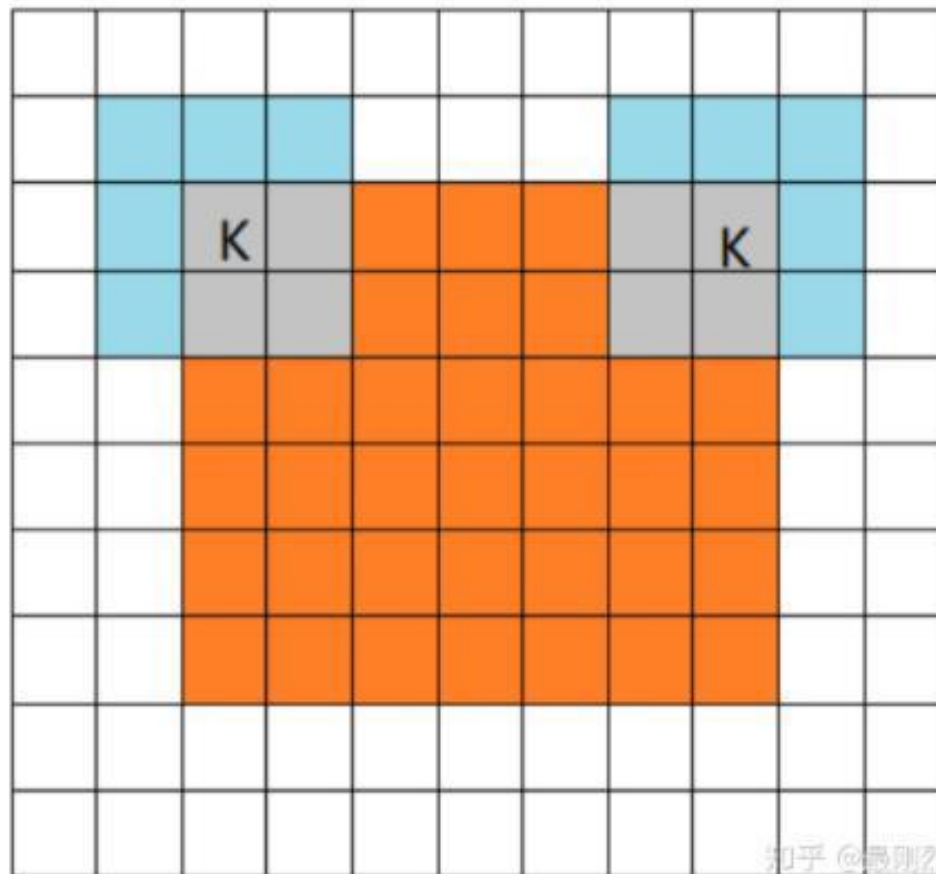
Stride s

- Down-sampling the output of convolution by sampling only every s pixels in each direction.
- For instance, the output of 'valid' convolution with stride s results in an output of size $\frac{m - k + s}{s} \times \frac{m - k + s}{s}$

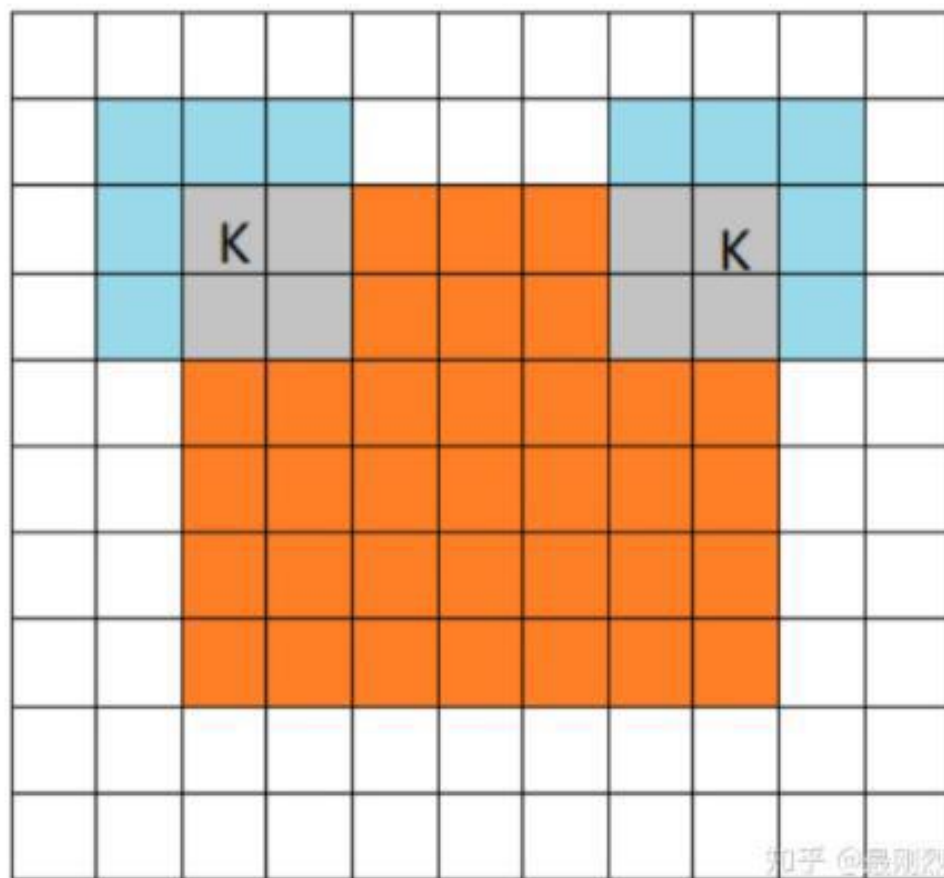
Valid, no zero-padding, $k \times k$ filter, $k=3$, $m=7$
Dim. of feature map $= (m-k+1) \times (m-k+1)$



Full, $k \times k$ filter, $k=3$, $m=7$, output
 $(m+k-1) \times (m+k-1)$



Same padding, $x-k+1=m$, $x=m+k-1$



Pooling (Pool)

Spatial summary by computing
operation over **window** with **stride**

- overlapping or non-overlapping
- separate across channels
- Current fashion:

**3x3 max pooling
with stride 2**

2x2 pooling, stride
2

23	7	7	8
10	1	9	0
4	4	11	6
2	5	12	7

Max pooling

23	9
5	12

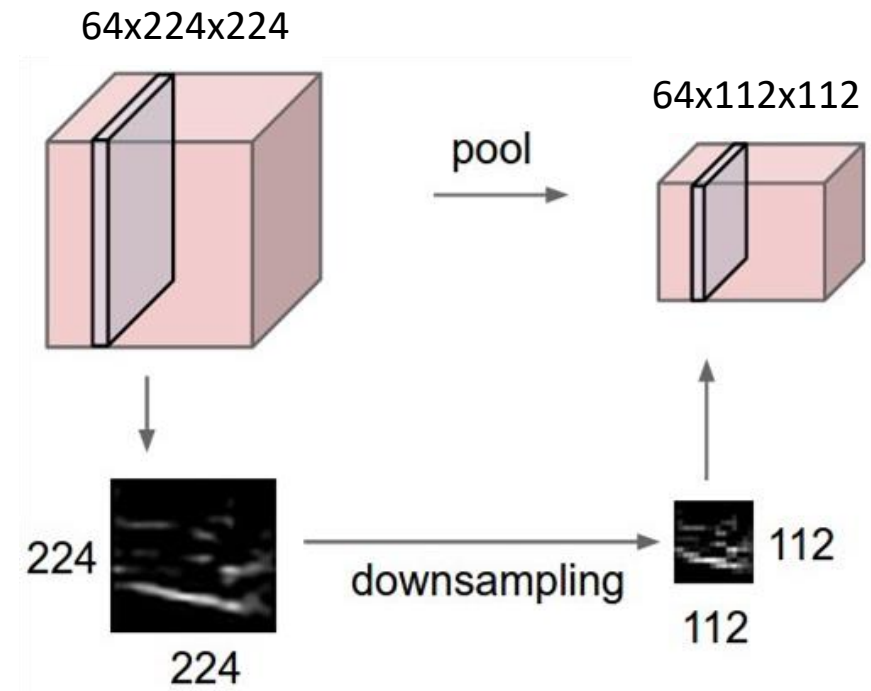
10	6
4	9

Average pooling

[figure credit BDTI]

Pooling

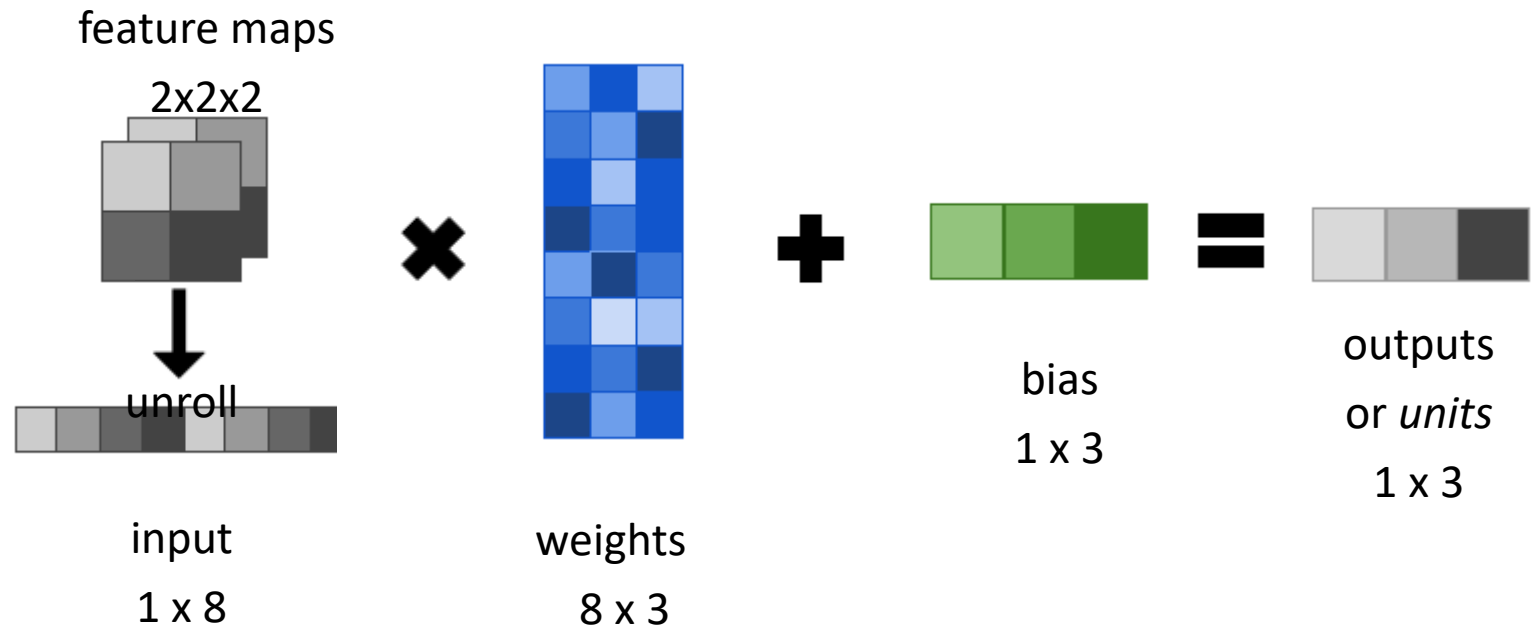
- reduce resolution
- increase receptive field size for later layers
- save computation
- add invariance to translation/noise within pooling window



Fully Connected Layers (FC)

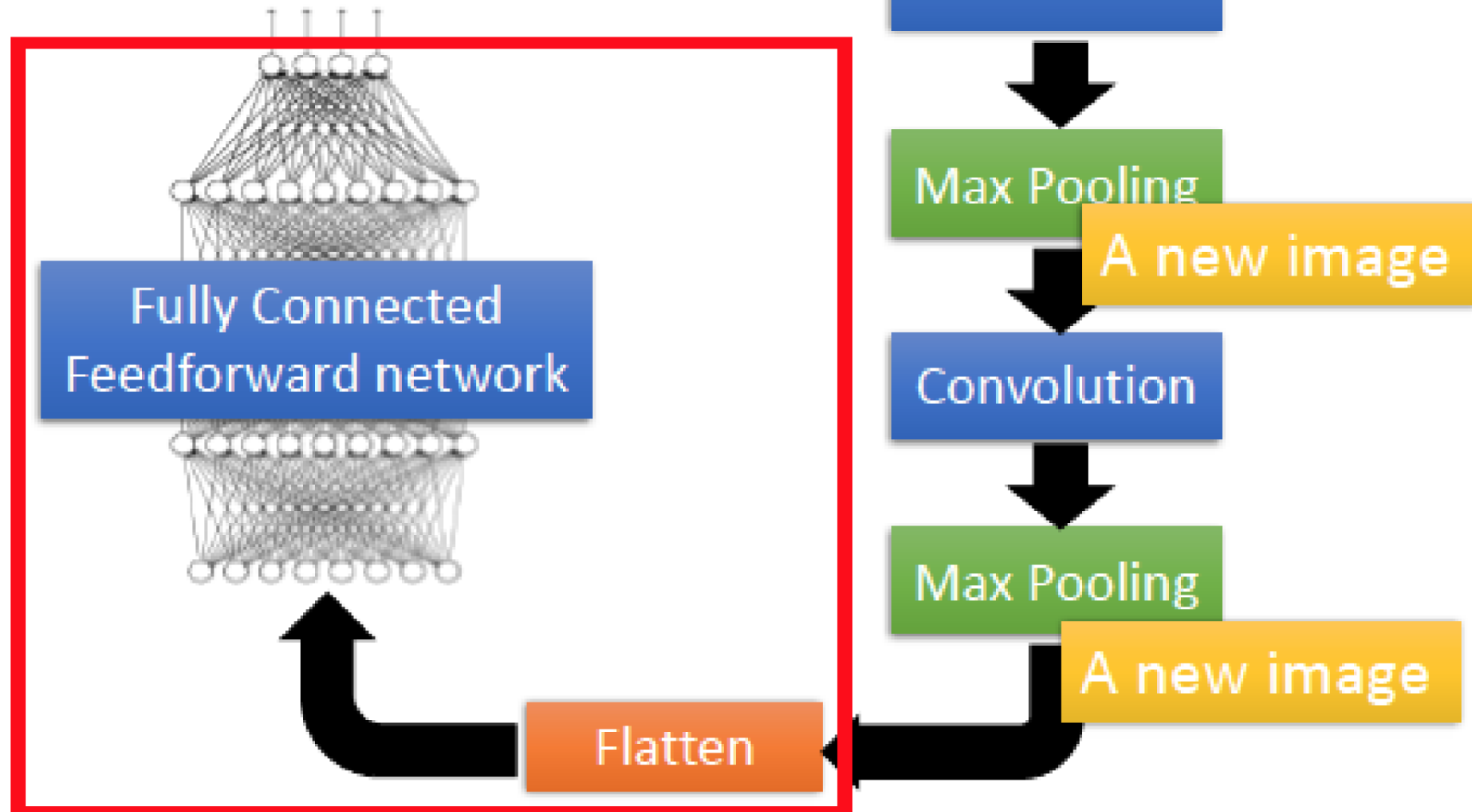
Learn a global feature from the full feature maps
Often found at the end of convnets

Note: unroll is also called **flatten**

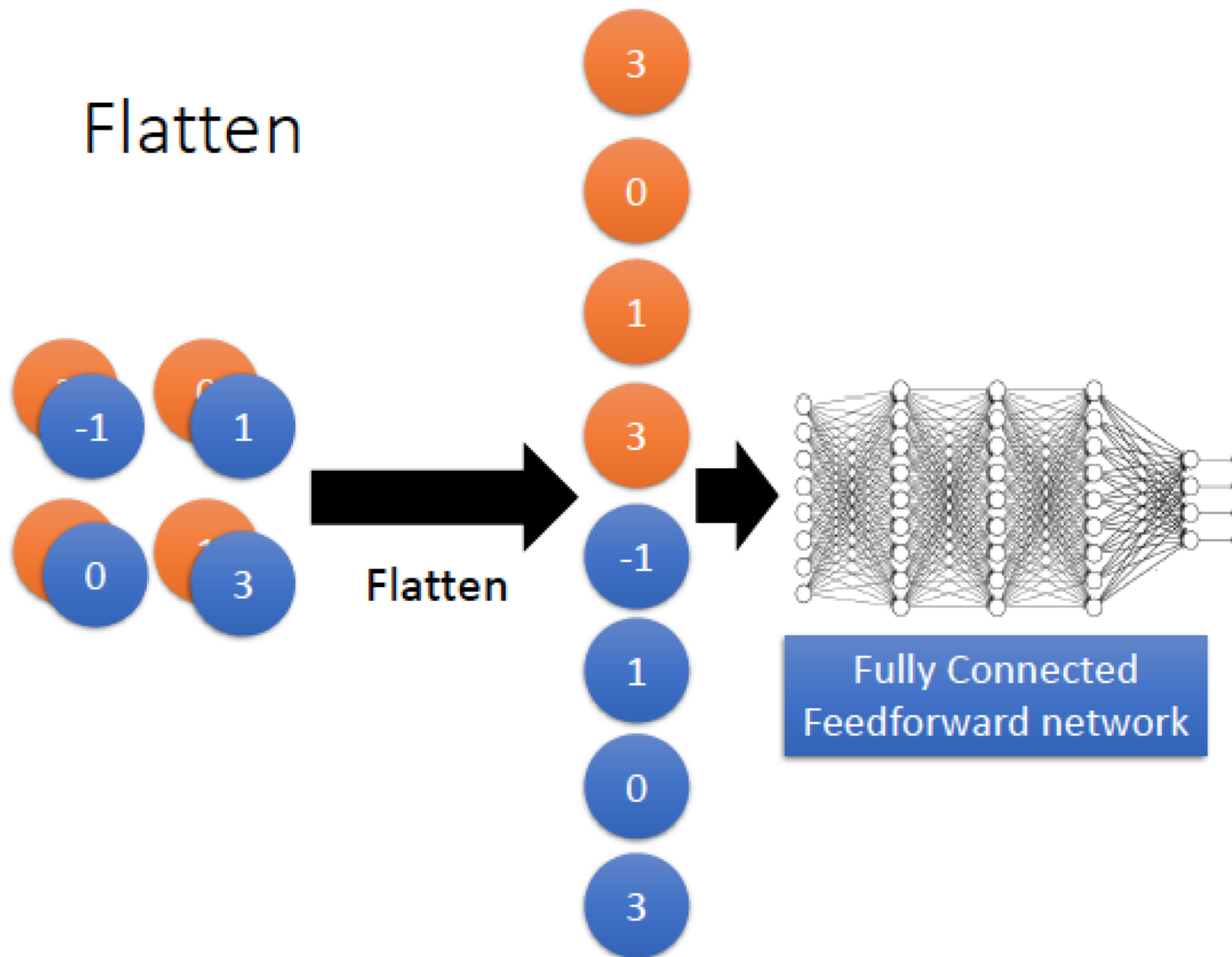


The whole CNN

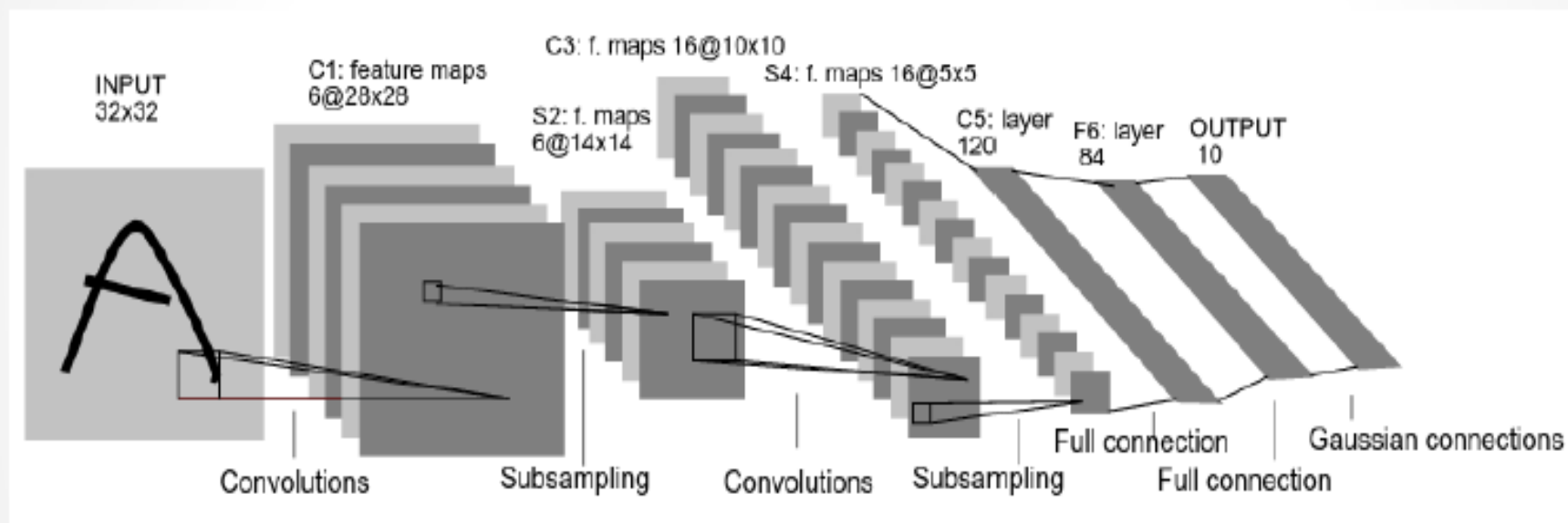
cat dog



Flatten

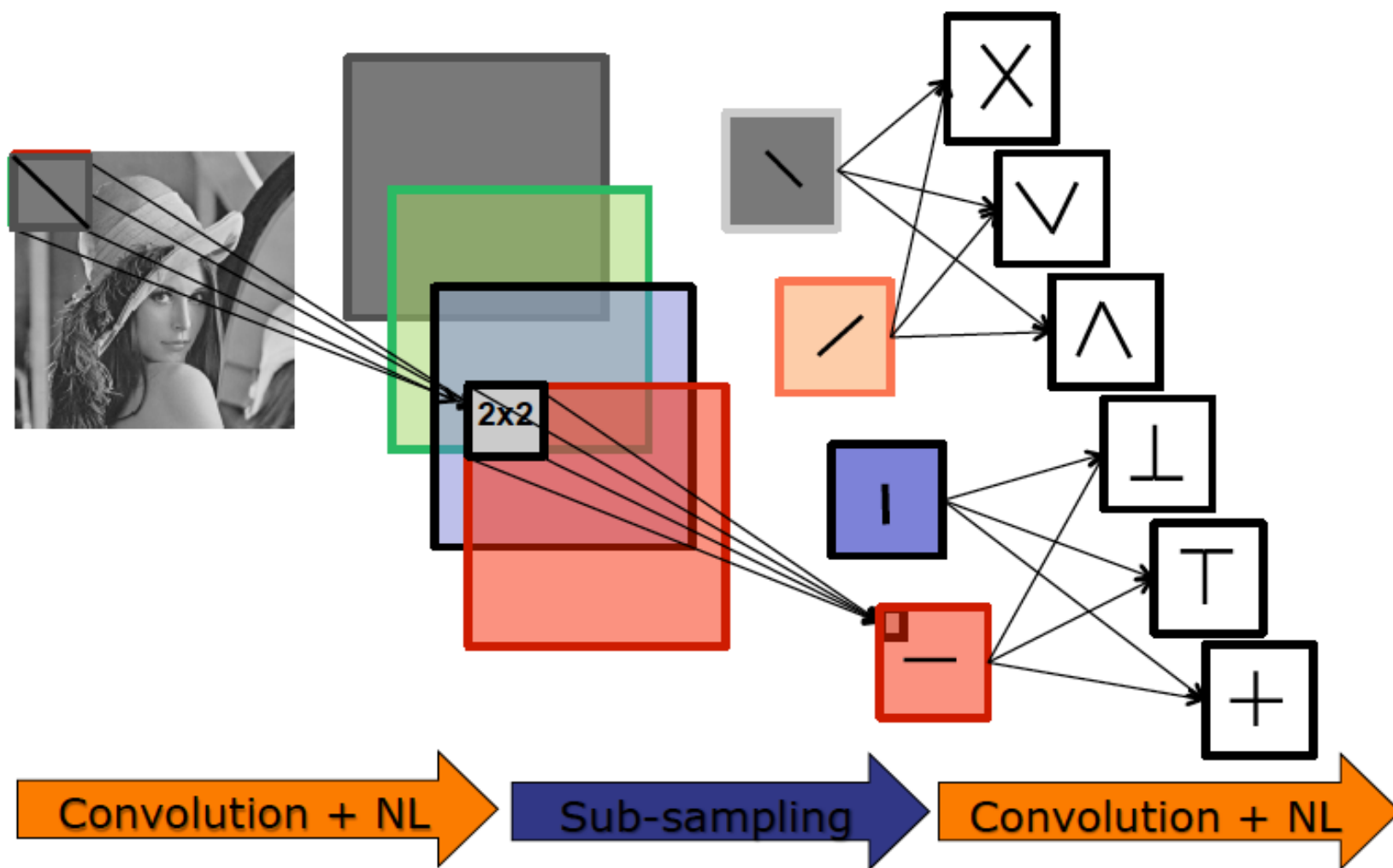


Putting It All Together!

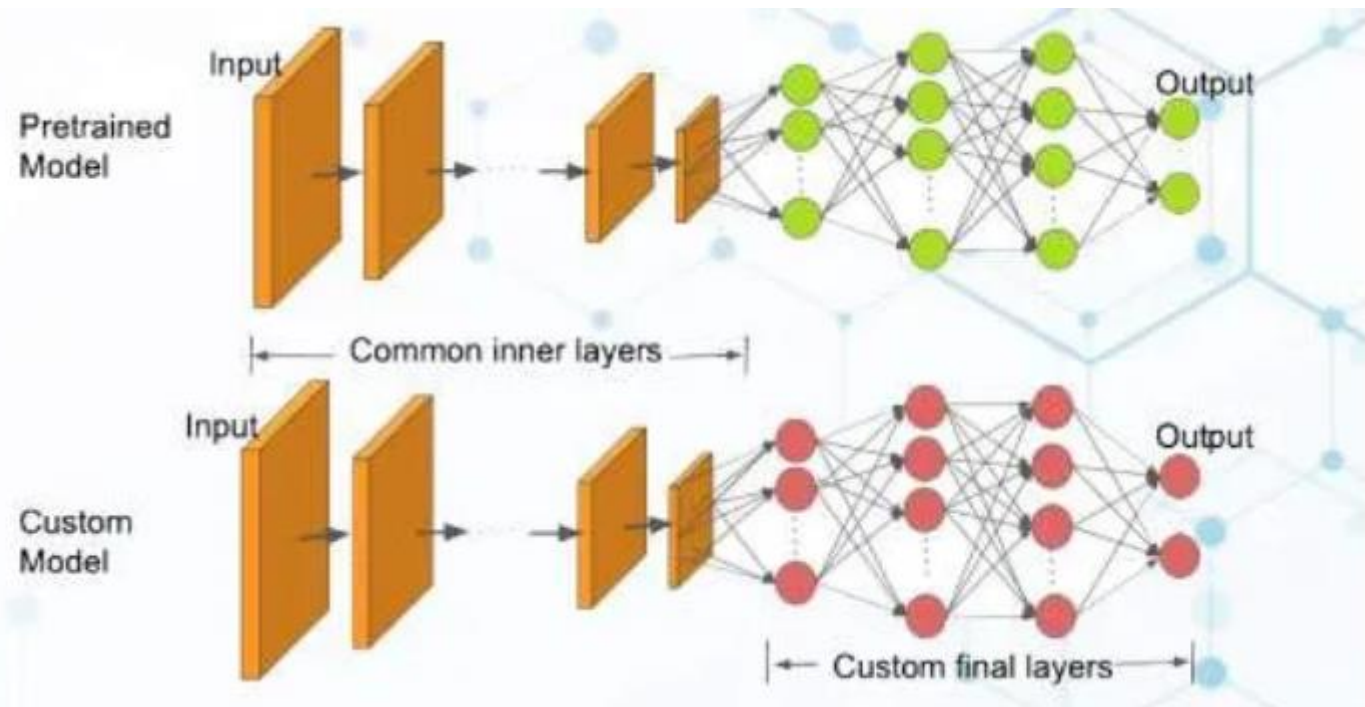


Lenet-5 (Lecun-98), Convolutional Neural Network for digits recognition

What is Convolutional NN ?

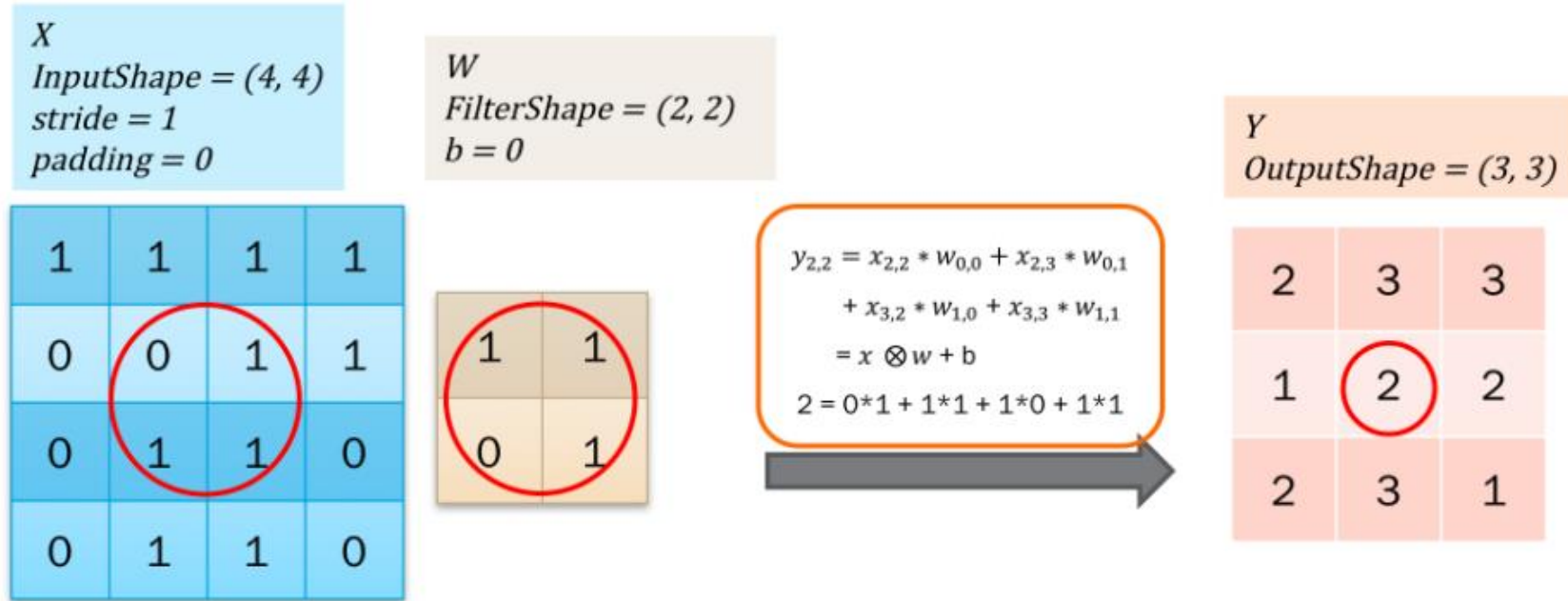


Transfer Learning



Learning of CNN

Forward propagation of CNN



$$OutputShape = \frac{InputShape - FilterShape + 2 * pad}{stride} + 1$$

Let's use the following notations:

$$\partial h_{ij} \text{ represents } \frac{\partial L}{\partial h_{ij}}$$

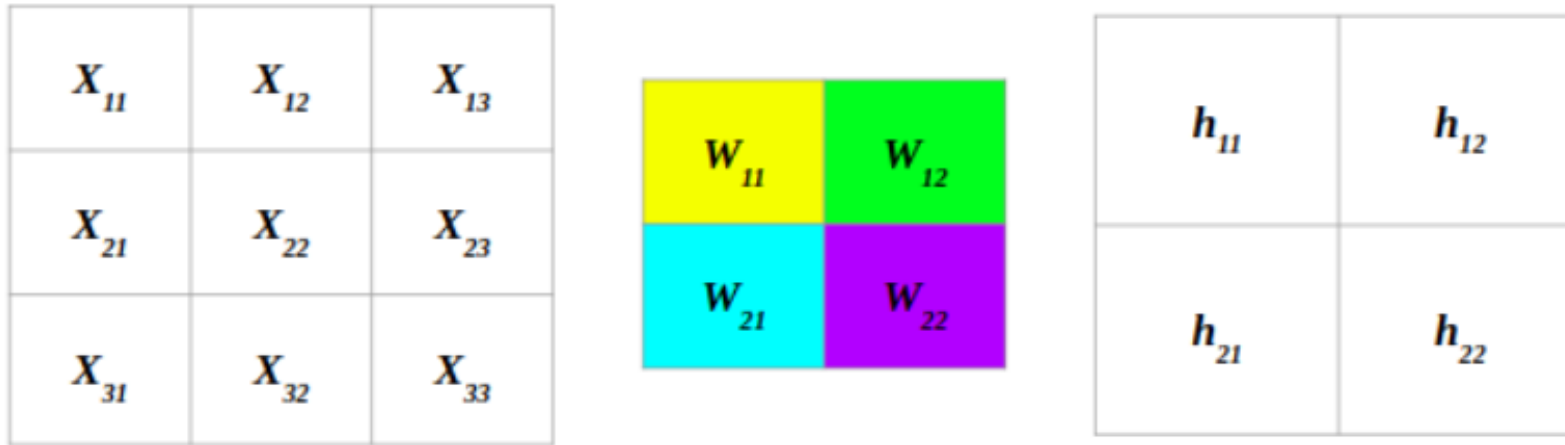
$$\partial w_{ij} \text{ represents } \frac{\partial L}{\partial w_{ij}}$$

Notations

Update on weight matrix of CNN

Note that there is only one way to reason the updates: Gradient descent, that is, moving along the opposite direction of the gradient vector.

Consider the following simple filter and the gradient h_{ij}



Input Size : 3x3, Filter Size : 2x2, Output Size : 2x2

Relationship of input, filter and output; note that Input and output could be just other hidden layers

Here, h_{11} is not the gradient, it is output!

$$h_{11} = W_{11}X_{11} + W_{12}X_{12} + W_{21}X_{21} + W_{22}X_{22}$$

$$h_{12} = W_{11}X_{12} + W_{12}X_{13} + W_{21}X_{22} + W_{22}X_{23}$$

$$h_{21} = W_{11}X_{21} + W_{12}X_{22} + W_{21}X_{31} + W_{22}X_{32}$$

$$h_{22} = W_{11}X_{22} + W_{12}X_{23} + W_{21}X_{32} + W_{22}X_{33}$$

Output Equations

Let's use the following notations:

$$\partial h_{ij} \text{ represents } \frac{\partial L}{\partial h_{ij}}$$

$$\partial w_{ij} \text{ represents } \frac{\partial L}{\partial w_{ij}}$$

Notations

To find Δw_{ij} we need the following partial derivative:

$$\frac{\partial L}{\partial w_{ij}} = \sum_{all\ related\ pk} \frac{\partial L}{\partial h_{pk}} \frac{\partial h_{pk}}{\partial w_{ij}}$$

So, we have:

$$\partial W_{11} = X_{11} \partial h_{11} + X_{12} \partial h_{12} + X_{21} \partial h_{21} + X_{22} \partial h_{22}$$

$$\partial W_{12} = X_{12} \partial h_{11} + X_{13} \partial h_{12} + X_{22} \partial h_{21} + X_{23} \partial h_{22}$$

$$\partial W_{21} = X_{21} \partial h_{11} + X_{22} \partial h_{12} + X_{31} \partial h_{21} + X_{32} \partial h_{22}$$

$$\partial W_{22} = X_{22} \partial h_{11} + X_{23} \partial h_{12} + X_{32} \partial h_{21} + X_{33} \partial h_{22}$$

X11	X12	X13
X21	X22	X23
X31	X32	X33

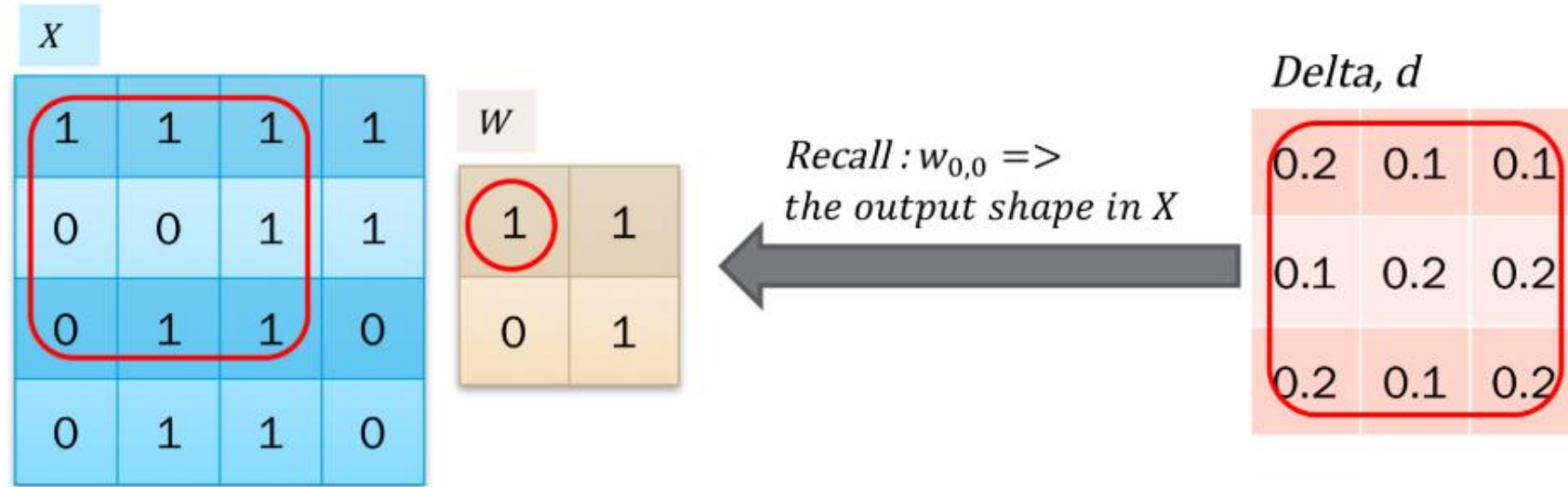
⊗

∂h11	∂h12
∂h21	∂h22

Final derivatives after performing back propagation

$$\Delta w_{ij} = -\eta Dw$$

Should also consider learning rate!



$$\frac{\partial d}{\partial w_{0,0}} = d_{0,0} * x_{0,0} + d_{0,1} * x_{1,0} + \dots + d_{3,3} * x_{3,3}$$

$$= x \otimes d = 0.4 + 0.2 + 0.3 = 0.9$$

$$\frac{\partial d}{\partial b} = \text{sum}(d)$$



Dw

0.9	1.1
0.8	0.8

Update on bias b of CNN

If bias is considered, we have:

$$h_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22} + b$$

$$h_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23} + b$$

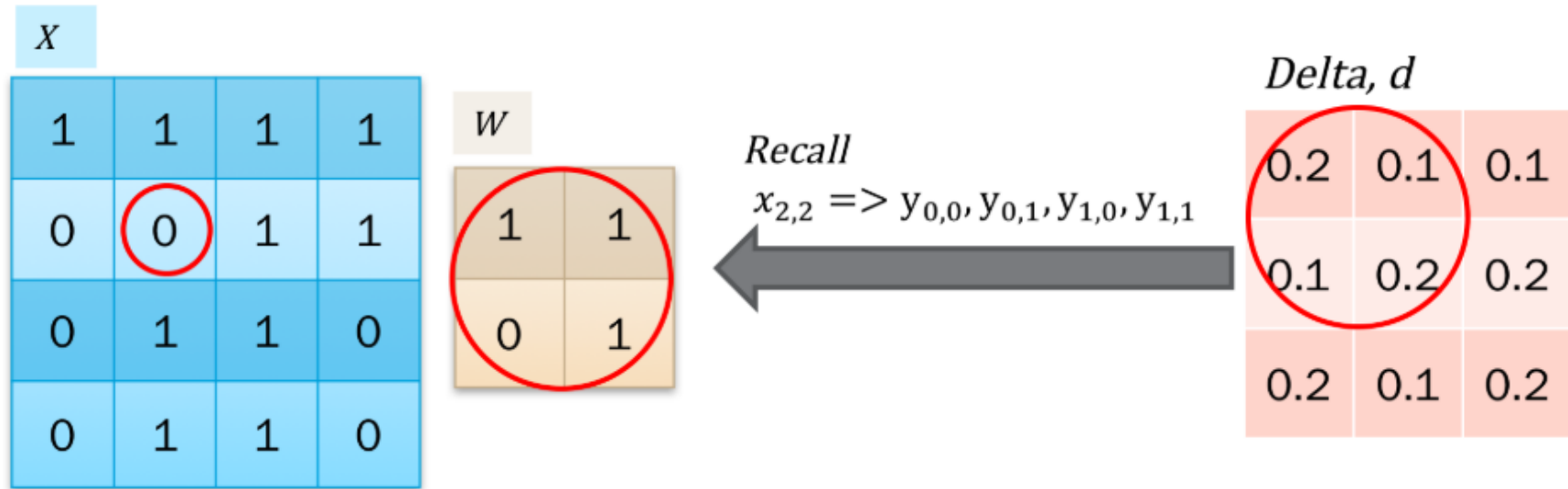
$$h_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32} + b$$

$$h_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33} + b$$

$$\frac{\partial L}{\partial b} = \sum_{all\ h_{ij}} \frac{\partial L}{\partial h_{ij}} \frac{\partial h_{ij}}{\partial b}$$

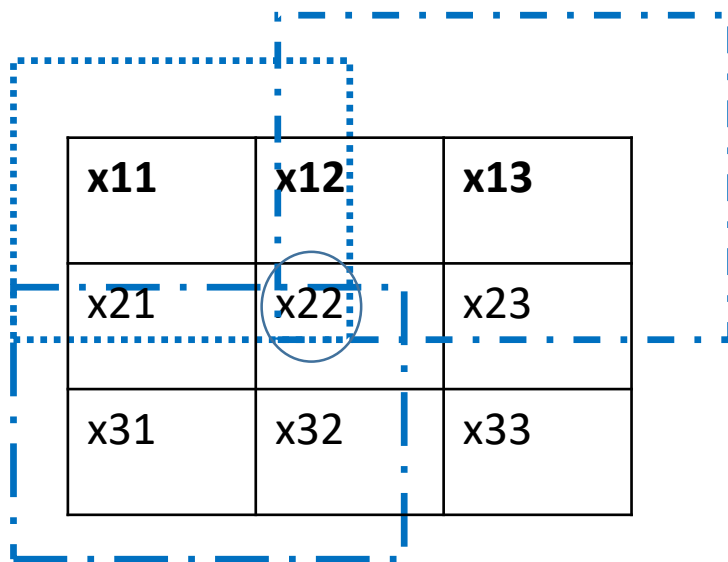
Gradients of the previous layer

Need to change the subscript of y, d, w by adding 1



$$\begin{aligned}\frac{\partial d}{\partial x_{2,2}} &= d_{0,0} * w_{1,1} + d_{0,1} * w_{1,0} + d_{1,0} * w_{0,1} + d_{1,1} * w_{0,0} \\ &= d \otimes w^{rotate} \\ &= 0.2 * 1 + 0.1 * 0 + 0.1 * 1 + 0.2 * 1 = 0.5\end{aligned}$$

Rotate (180 degree): switch rows and reverse order



w11	w12
w21	w22

∂h_{11}	∂h_{12}
∂h_{21}	∂h_{22}



h11	h12
h21	h22

w22	w21
w12	w11

w22

w21

$$\begin{aligned}
 x_{22} * w_{22} + \dots &= h_{11} \\
 x_{22} * w_{21} + \dots &= h_{12} \\
 x_{22} * w_{12} + \dots &= h_{21} \\
 x_{22} * w_{11} + \dots &= h_{22}
 \end{aligned}$$

$$\begin{aligned}
 \frac{\partial L}{\partial x_{22}} &= \frac{\partial L}{\partial h_{11}} * \frac{\partial h_{11}}{\partial x_{22}} + \frac{\partial L}{\partial h_{12}} * \frac{\partial h_{12}}{\partial x_{22}} + \\
 &\quad \frac{\partial L}{\partial h_{21}} * \frac{\partial h_{21}}{\partial x_{22}} + \frac{\partial L}{\partial h_{22}} * \frac{\partial h_{22}}{\partial x_{22}}
 \end{aligned}$$

dx

0.2	0.3	0.2	0.1
0.1	0.5	0.5	0.3
0.2	0.4	0.5	0.4
0.0	0.2	0.1	0.2

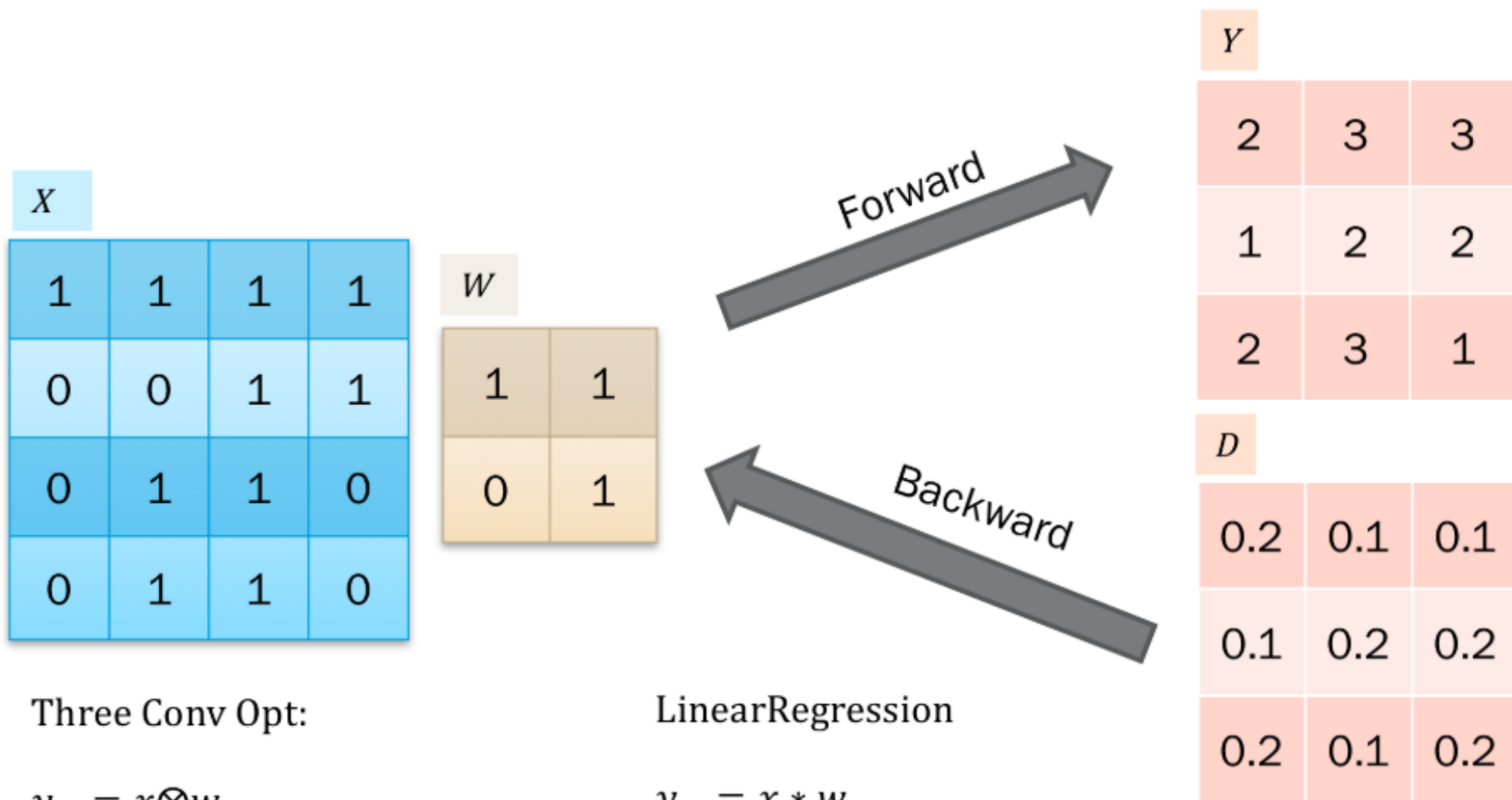
$NewFilter = w^{rotate}$

1	0
1	1



$Delta, d$

0	0	0	0	0
0	0.2	0.1	0.1	0
0	0.1	0.2	0.2	0
0	0.2	0.1	0.2	0
0	0	0	0	0



Three Conv Opt:

$$y = x \otimes w$$

$$dw = x \otimes d$$

$$dx = d \otimes w^{rotate}$$

LinearRegression

$$y = x * w$$

$$dw = x^T * d$$

$$dx = d * w^T$$

Pooling
 FilterShape = (2, 2)
 stride=1
 padding=0

X

2	3	3
1	2	2
2	3	1

MAX Pooling

AVG Pooling

L2 Pooling

Y

3	3
3	3

$$y_{0,0} = \text{Max}(x_{0,0}, x_{0,1}, x_{1,0}, x_{1,1})$$

$$= \text{Max}(2, 3, 1, 2)$$

$$= 3$$

2.0	2.5
2.0	2.0

$$y_{0,0} = \text{Mean}(x_{0,0}, x_{0,1}, x_{1,0}, x_{1,1})$$

$$= \text{sum}(2, 3, 1, 2) / 4.0$$

$$= 2.0$$

$\sqrt{18}$	$\sqrt{26}$
$\sqrt{18}$	$\sqrt{18}$

$$y_{0,0} = L_2(x_{0,0}, x_{0,1}, x_{1,0}, x_{1,1})$$

$$= \text{sqrt}(\sum x_i^2)$$

$$= \sqrt{(4 + 9 + 1 + 4)}$$

$$= \sqrt{18}$$

Forward

Backward

AVG Pooling = Cov

0.25	0.25
0.25	0.25

L2 Pooling =
square + Cov + sqrt

1	1
1	1