

Convolutions and CNNs

Badri N Subudhi

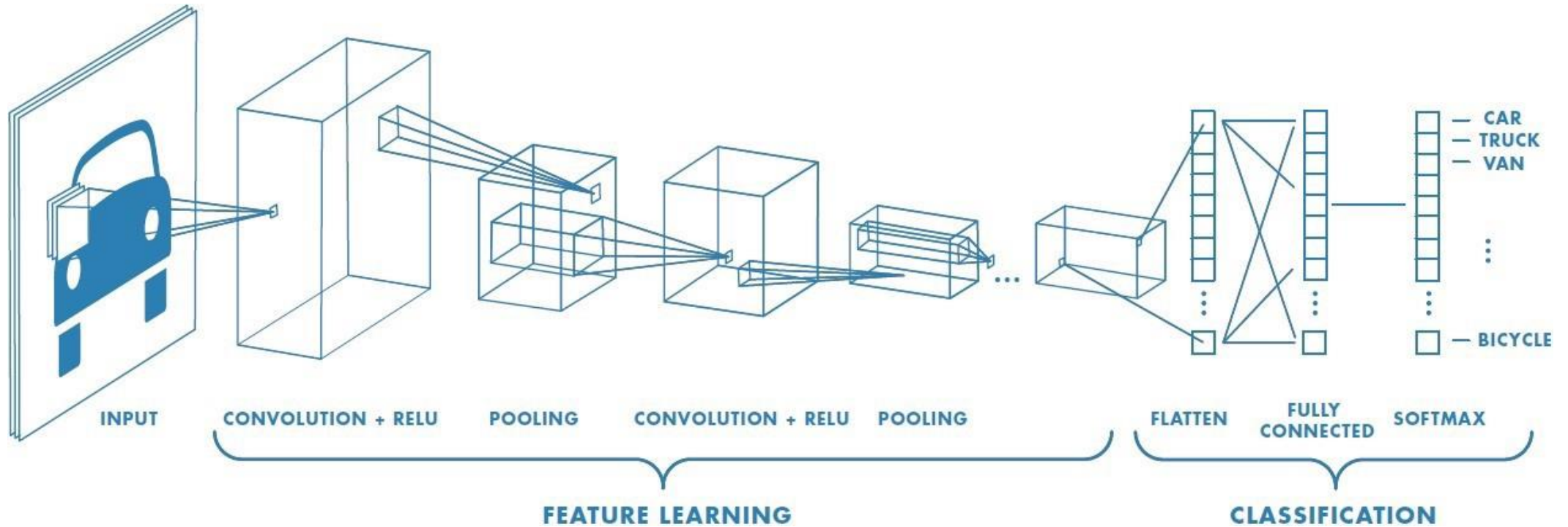
Associate Professor

IIT Jammu

CNNs ?

CNN : Neural Network with a convolution operation instead of matrix multiplication in at least one of the layers

A typical CNN architecture



Convolution and Correlation

Convolution and correlation operations are basically used to extract information from images.

Convolution and correlation are basically linear and shift-invariant operations

The term *linear* indicates that a pixel is replaced by the linear combination of its neighbors.

The term *shift invariant* indicates that the same operation is performed at every point in the image.

Convolution

Convolution is basically a mathematical operation where each value in the output is expressed as the sum of values in the input multiplied by a set of weighting coefficients.

Depending upon the weighting coefficients, convolution operation is used to perform spatial domain low-pass and high-pass filtering of the image.

An image can be either smoothened or sharpened by convolving the image with respect to low-pass and high-pass filter mask respectively.

Convolution has a multitude of applications including image filtering, image enhancement, image restoration, feature extraction and template matching.

$$y[n_1, n_2] = \sum_{k_1=-\infty}^{\infty} \sum_{k_2=-\infty}^{\infty} x(k_1, k_2) h(n_1 - k_1, n_2 - k_2)$$

The input matrix $x(m, n)$ and $h(m, n)$. Perform the linear convolution between these two matrices.

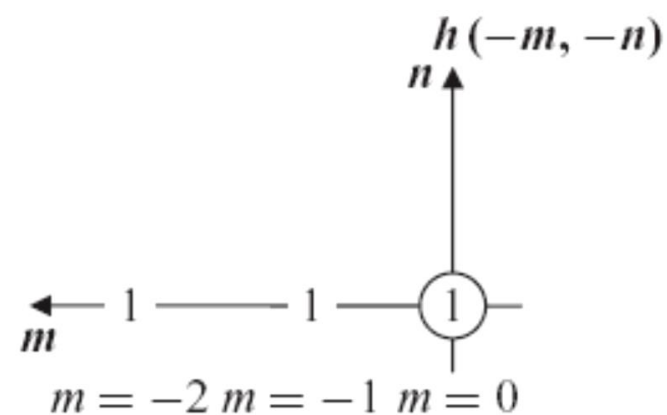
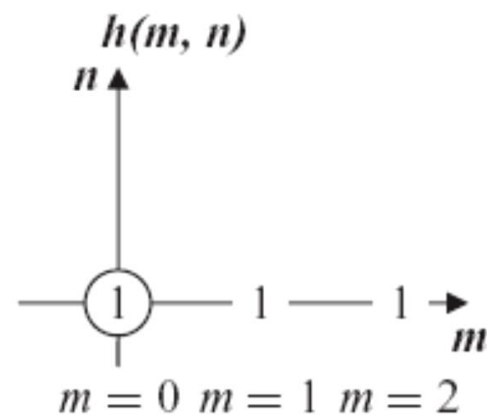
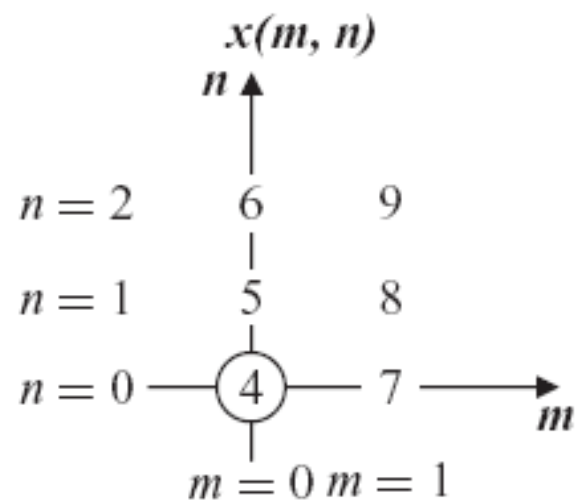
$$x(m, n) = \begin{pmatrix} 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \quad h(m, n) = \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

$$x(m, n) = \begin{pmatrix} (0,0) & (0,1) & (0,2) \\ 4 & 5 & 6 \\ (1,0) & (1,1) & (1,2) \\ 7 & 8 & 9 \end{pmatrix} \quad h(m, n) = \begin{pmatrix} (0,0) \\ 1 \\ (1,0) \\ 1 \\ (2,0) \\ 1 \end{pmatrix}$$

$$\text{Dimension of resultant matrix} = \begin{cases} (\text{No. of rows of } x(m, n) + \text{No. of rows of } h(m, n) - 1) \\ \times \\ (\text{No. of columns of } x(m, n) + \text{No. of columns of } h(m, n) - 1) \end{cases}$$

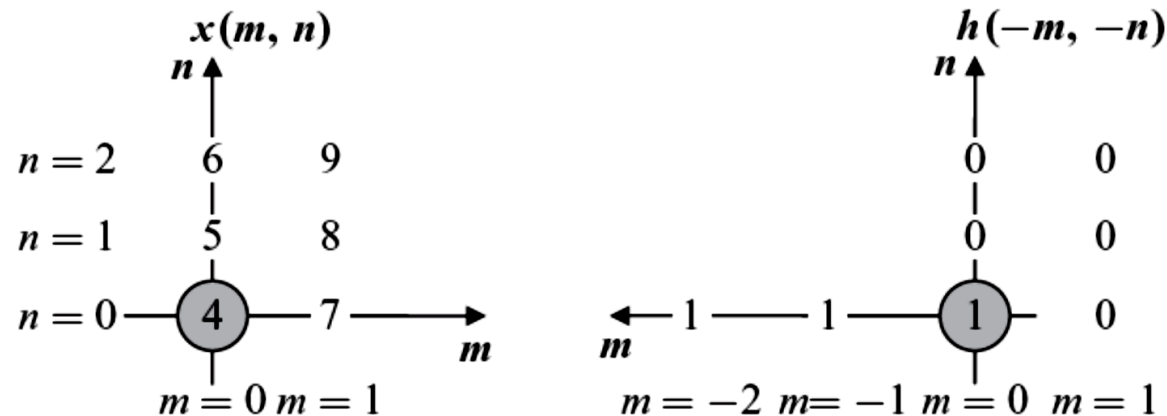
$$\text{Dimension of resultant matrix} = (2 + 3 - 1) \times (3 + 1 - 1) = 4 \times 3$$

$$y(m, n) = \begin{pmatrix} y(0, 0) & y(0, 1) & y(0, 2) \\ y(1, 0) & y(1, 1) & y(1, 2) \\ y(2, 0) & y(2, 1) & y(2, 2) \\ y(3, 0) & y(3, 1) & y(3, 2) \end{pmatrix}$$



To determine the value of $y(0, 0)$

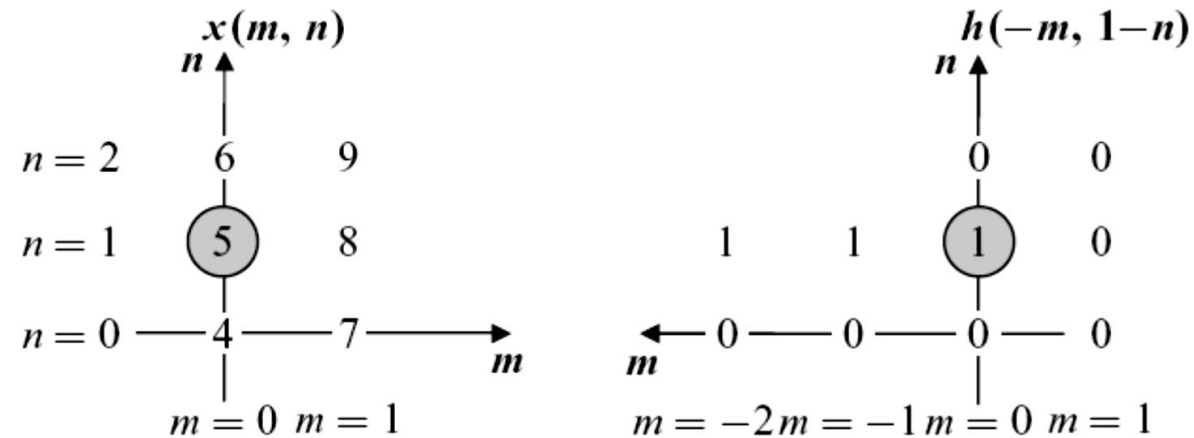
The common values of $x(m, n)$ and $h(-m, -n)$ are multiplied and then added to get the value of $y(0, 0)$. The shaded circle indicates the common area between the two signals.



The value $y(0, 0)$ is obtained as $y(0, 0) = 4 \times 1 = 4$.

To determine the value of $y(0,1)$

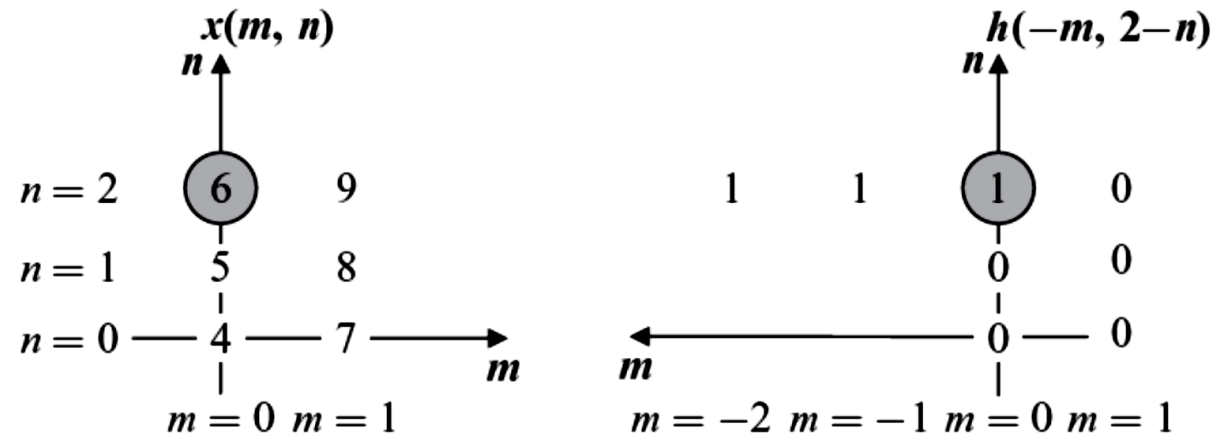
Now the signal $h(-m, -n)$ is shifted along the ' n ' axis by one unit to get $h(-m, 1-n)$ and $x(m, n)$ is unaltered. The common value between the two signals is multiplied and then added to get $y(0, 1)$.



The value of $y(0, 1)$ is given as $y(0, 1) = 5 \times 1 = 5$.

To determine the value of $y(0, 2)$

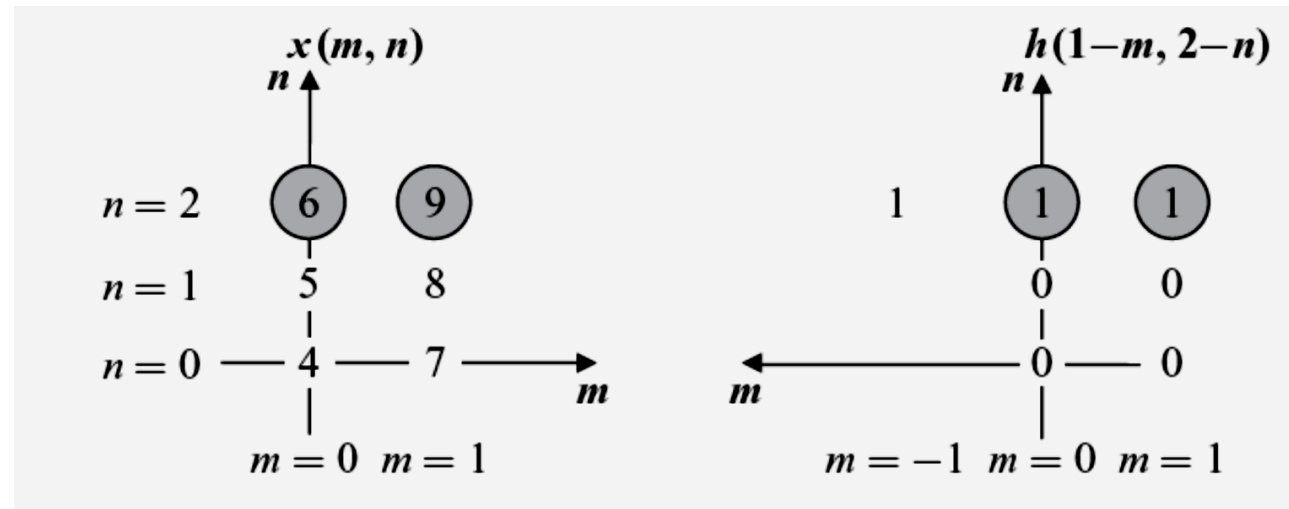
The value of $h(-m, -n)$ is shifted by two units along the ' n ' axis to get $h(-m, 2-n)$. Then, the common values between $x(m, n)$ and $h(-m, 2-n)$ are multiplied and then added to get $y(0, 2)$.



The resultant value of $y(0, 2)$ is $y(0, 2) = 6 \times 1 = 6$.

To determine the value of $y(1, 2)$

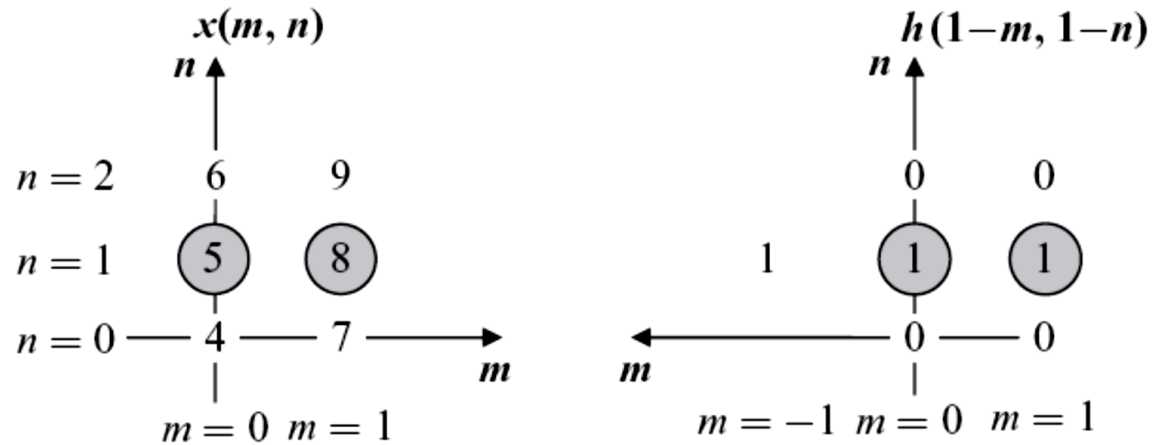
Here, $h(-m, 2-n)$ is shifted along the ' m ' axis to one unit towards right to get $h(1-m, 2-n)$. Then the common values between $x(m, n)$ and $h(1-m, 2-n)$ are multiplied and added to get $y(1, 2)$.



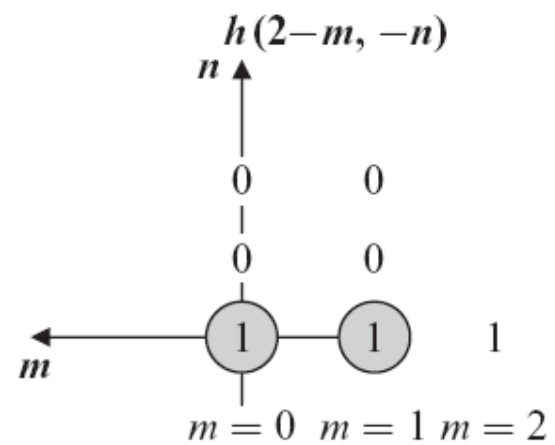
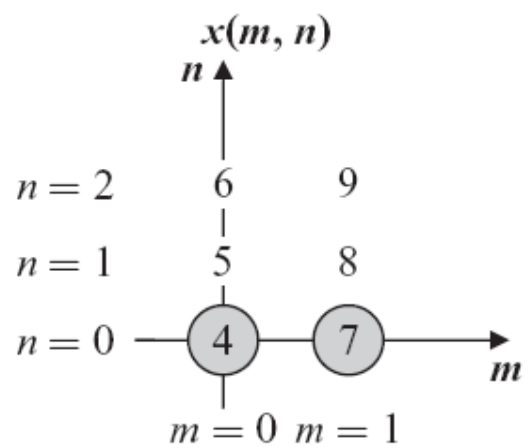
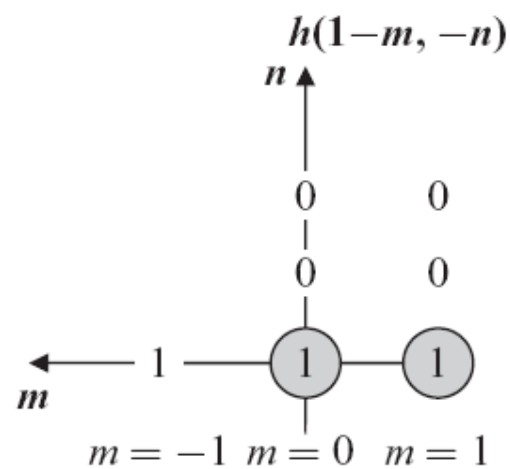
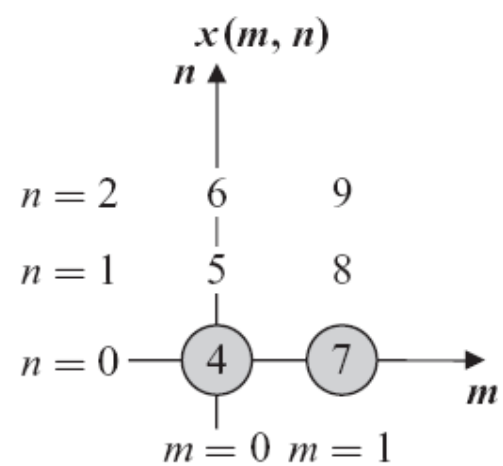
The final value of $y(1, 2) = 6 \times 1 + 9 \times 1 = 15$.

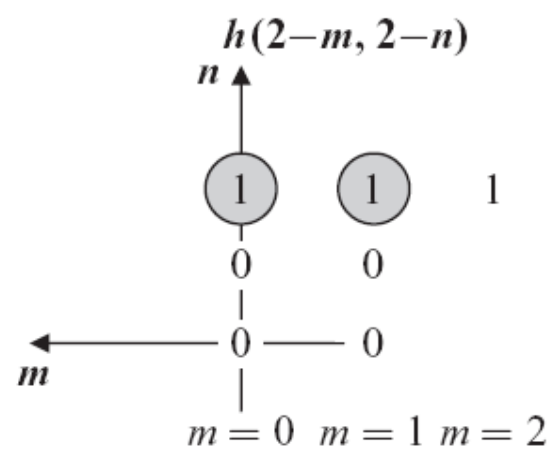
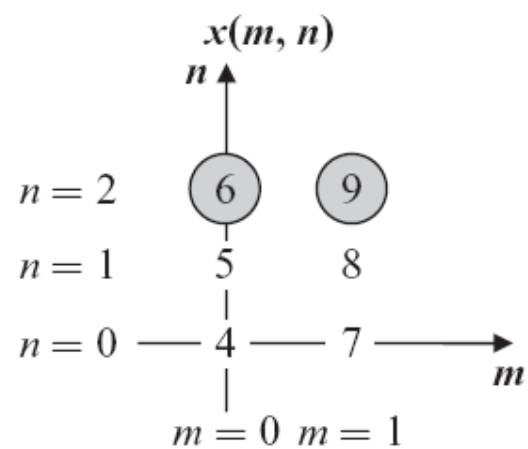
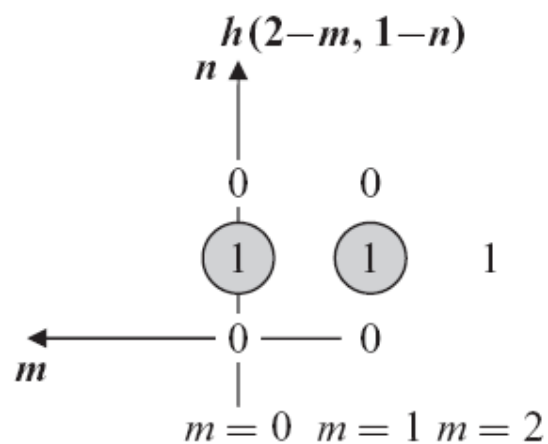
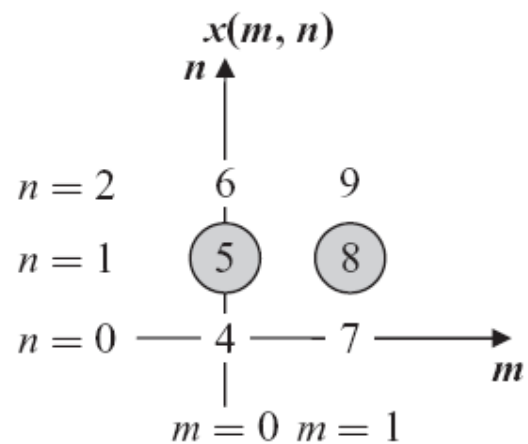
To determine the value of $y(1, 1)$

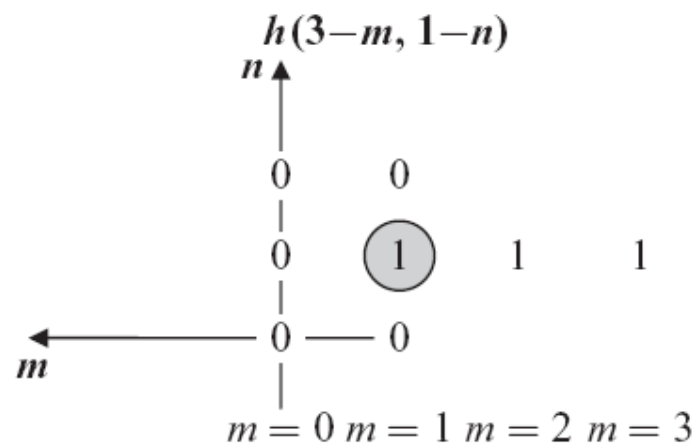
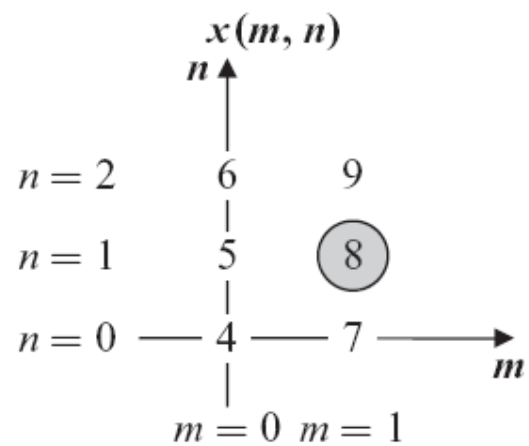
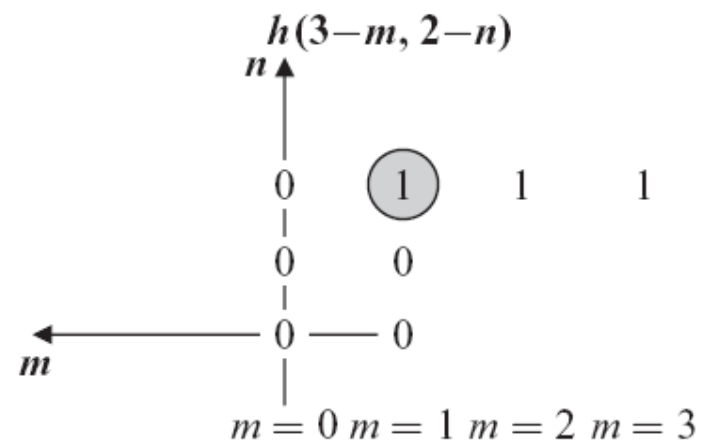
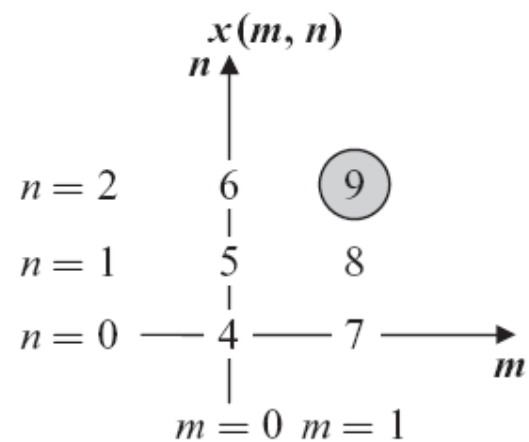
Now the value of $h(1-m, 2-n)$ is shifted down along the 'n' axis to get $h(1-m, 1-n)$. The common values between $x(m, n)$ and $h(1-m, 1-n)$ are multiplied and added to get $y(1, 1)$.

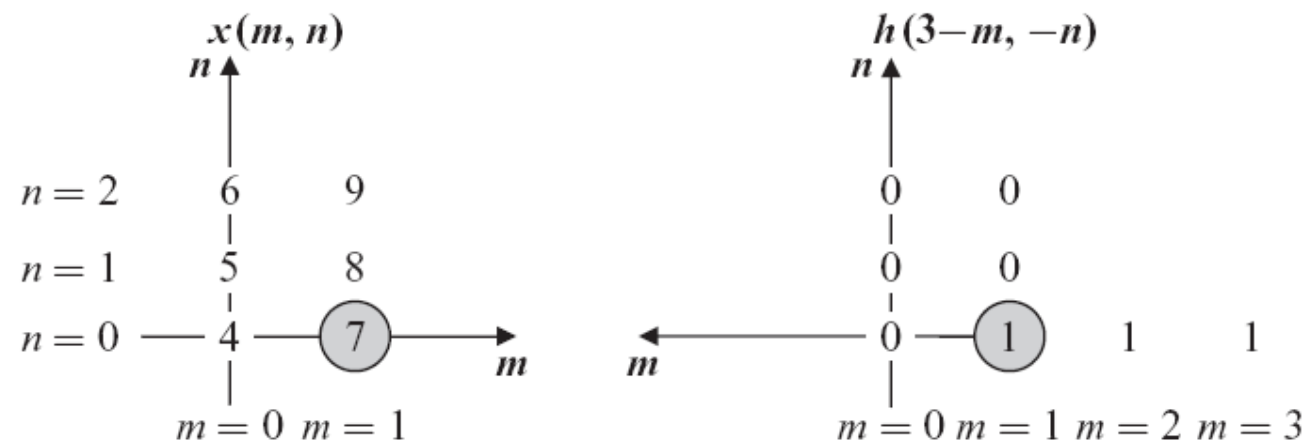


The value of $y(1, 1)$ is obtained as $y(1, 1) = 5 \times 1 + 8 \times 1 = 13$.









$$y(m, n) = \begin{pmatrix} 4 & 5 & 6 \\ 11 & 13 & 15 \\ 11 & 13 & 15 \\ 7 & 8 & 9 \end{pmatrix}$$

Matrix representation

Perform the linear convolution between the two matrices $x(m, n)$ and $h(m, n)$ given below

$$x(m, n) = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \quad h(m, n) = \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$y(m, n) = \begin{pmatrix} 1 & 3 & 5 & 3 \\ 5 & 12 & 16 & 9 \\ 12 & 27 & 33 & 18 \\ 1 & 24 & 28 & 15 \\ 7 & 15 & 17 & 9 \end{pmatrix}$$

Matrix representation

$$x(m, n) = \begin{pmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 7 & 8 & 9 \end{pmatrix} \quad h(m, n) = (3 \quad 4 \quad 5)$$

$$y(m, n) = \begin{pmatrix} 3 & 10 & 22 & 22 & 15 \\ 12 & 31 & 58 & 49 & 30 \\ 21 & 52 & 94 & 76 & 45 \end{pmatrix}$$

Matrix representation

Correlation

Correlation is a mathematical operation that is similar to convolution. Correlation is basically used to obtain similarity between two signals.

Autocorrelation is basically a signal correlated with itself.

The amplitude of each sample in the cross-correlation is a measure of how much one signal resembles the other.

$$x[n] * x^*[-n] = \sum_{k=-\infty}^{\infty} x[k]x^*[-(n-k)]$$

$$x[m, n] * x^*[-m, -n] = \sum_a \sum_b x(a, b)x^*[-(m-a), -(n-b)]$$

$$x[m, n] * x^*[-m, -n] = \sum_a \sum_b x(a, b)x^*[(a-m), (b-n)] = r_{xx}(m, n)$$

Determine the correlation between two matrices

$$x_1[m, n] = \begin{bmatrix} 3 & 1 \\ 2 & 4 \end{bmatrix}$$

$$x_2[m, n] = \begin{bmatrix} 1 & 5 \\ 2 & 3 \end{bmatrix}.$$

Step 1

From the given $x_2[m, n]$ determine $x_2[-m, -n]$.

Step 1a To get the folded version of $x_2[m, n]$, first fold $x_2[m, n]$ column wise to get

$$x_2[m, -n] = \begin{bmatrix} 5 & 1 \\ 3 & 2 \end{bmatrix}$$

Step 1b Then fold $x_2[m, -n]$ along row-wise to get $x_2[-m, -n]$ which is given by

$$x_2[-m, -n] = \begin{bmatrix} 3 & 2 \\ 5 & 1 \end{bmatrix}$$

Step 2

Now we have to perform linear convolution between $x_1[m, n]$ and $x_2[-m, -n]$.

Step 2a Formation of block matrix The number of block matrices depends on the number of columns of $x_1[m, n]$. In this case, the number of columns of $x_2[m, n]$ is two. Hence, two block matrices have to be formed. Let the two block matrices be H_0 and H_1 respectively.

Step 2b Formation of block matrix H_0 The block matrix H_0 formed from the first row of $x_1[m, n]$ is given below.

$$H_0 = \begin{bmatrix} 3 & 0 \\ 1 & 3 \\ 0 & 1 \end{bmatrix}$$

Step 2b Formation of block matrix H_1 The block matrix H_1 formed from the second row of $x_1[m, n]$ is given below.

$$H_1 = \begin{bmatrix} 2 & 0 \\ 4 & 2 \\ 0 & 4 \end{bmatrix}$$

Step 3 Formation of block Toeplitz matrix

Let the block Toeplitz matrix be denoted by A . The number of zeros to be appended in the block matrix A depends on the number of rows of $x_2[m, n]$. In this case, $x_2[m, n]$ has two rows; hence one zero has to be appended. The matrix A is given by

$$A = \begin{bmatrix} H_0 & 0 \\ H_1 & H_0 \\ 0 & H_1 \end{bmatrix}.$$

Then, substituting the values of H_0 and H_1 in the matrix A , we get

$$A = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 2 & 0 & 3 & 0 \\ 4 & 2 & 1 & 3 \\ 0 & 4 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 4 & 2 \\ 0 & 0 & 0 & 4 \end{bmatrix}$$

The resultant matrix $y[m, n]$ is obtained by multiplying the matrix A with a matrix whose elements are from $x_2[-m, -n]$.

$$y[m, n] = \begin{bmatrix} 3 & 0 & 0 & 0 \\ 1 & 3 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 2 & 0 & 3 & 0 \\ 4 & 2 & 1 & 3 \\ 0 & 4 & 0 & 1 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 4 & 2 \\ 0 & 0 & 0 & 4 \end{bmatrix} \times \begin{bmatrix} 3 \\ 2 \\ 5 \\ 1 \end{bmatrix} = \begin{bmatrix} 9 \\ 9 \\ 2 \\ 21 \\ 24 \\ 9 \\ 10 \\ 22 \\ 4 \end{bmatrix}$$

The result $y[m, n]$ is a 3×3 matrix which is given by

$$y[m, n] = \begin{bmatrix} 9 & 9 & 2 \\ 21 & 24 & 9 \\ 10 & 22 & 4 \end{bmatrix}$$

Determine the correlation between two matrices

$$x_1[m, n] = \begin{bmatrix} 1 & 5 \\ 2 & 4 \end{bmatrix}$$

$$x_2[m, n] = \begin{bmatrix} 3 & 2 \\ 4 & 1 \end{bmatrix}$$

Determine the correlation between two matrices

$$x_1[m, n] = \begin{bmatrix} 3 & 2 \\ 1 & 5 \end{bmatrix} \quad x_2[m, n] = \begin{bmatrix} 3 & 2 \\ 1 & 5 \end{bmatrix}$$

The Convolution operation

0	0	0	0	0	0	...
0	156	155	156	158	158	...
0	153	154	157	159	159	...
0	149	151	155	158	159	...
0	146	146	149	153	158	...
0	145	143	143	148	158	...
...

Input Channel #1 (Red)

0	0	0	0	0	0	...
0	167	166	167	169	169	...
0	164	165	168	170	170	...
0	160	162	166	169	170	...
0	156	156	159	163	168	...
0	155	153	153	158	168	...
...

Input Channel #2 (Green)

0	0	0	0	0	0	...
0	163	162	163	165	165	...
0	160	161	164	166	166	...
0	156	158	162	165	166	...
0	155	155	158	162	167	...
0	154	152	152	157	167	...
...

Input Channel #3 (Blue)

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

Kernel Channel #1



308

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

Kernel Channel #2



-498

0	1	1
0	1	0
1	-1	1

Kernel Channel #3



164

+

+

+ 1 = -25

Bias = 1

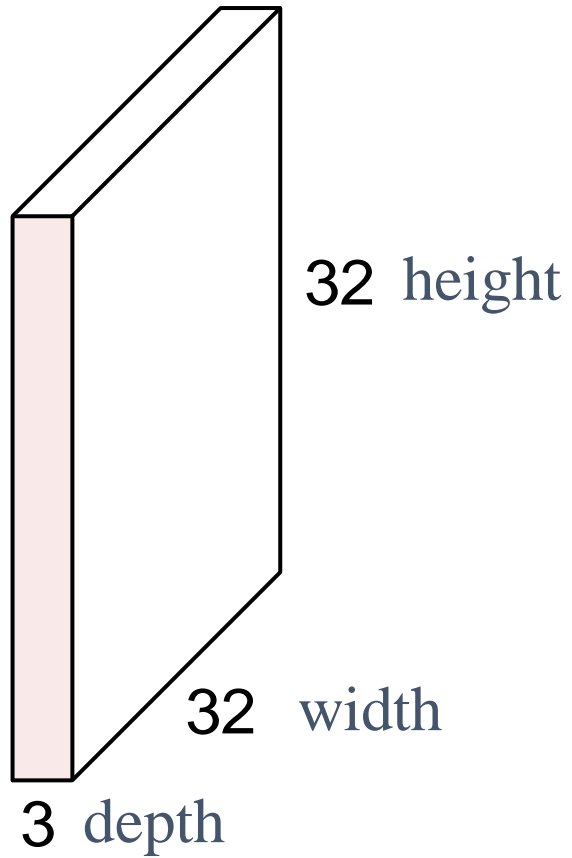
Output

-25				...
				...
				...
				...
...

Convolution Layers

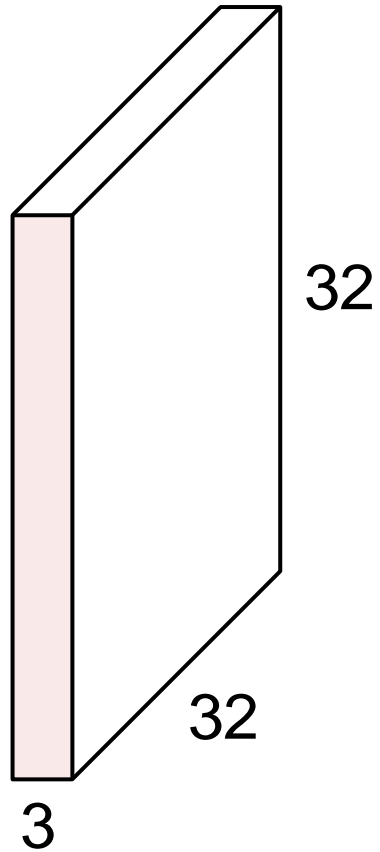
Convolution Layer

32x32x3 image



Convolution Layer

32x32x3 image



5x5x3 filter

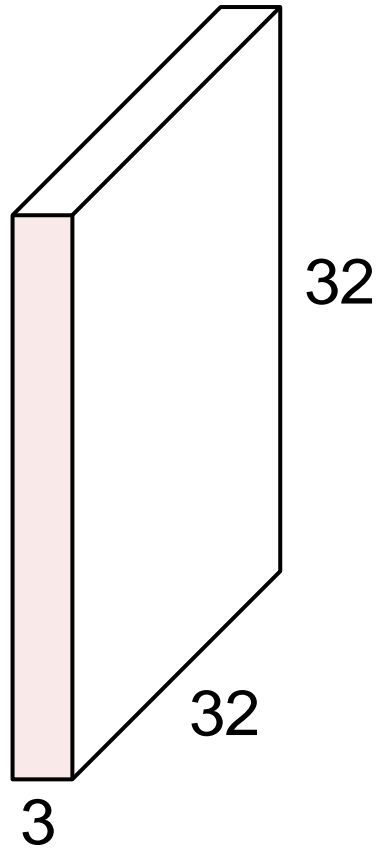


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

Convolution Layer

32x32x3 image

Filters always extend the full depth of the input volume

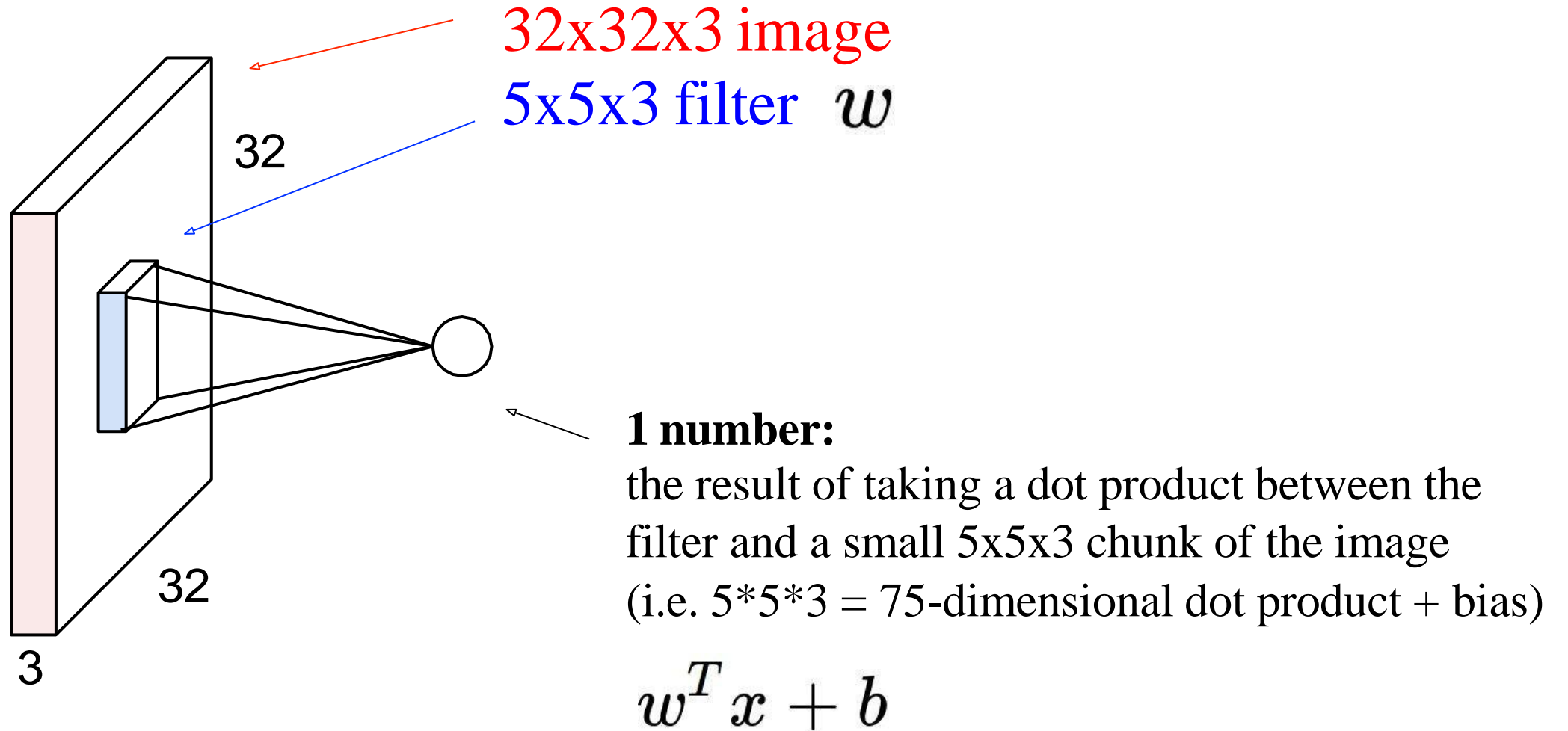


5x5x3 filter

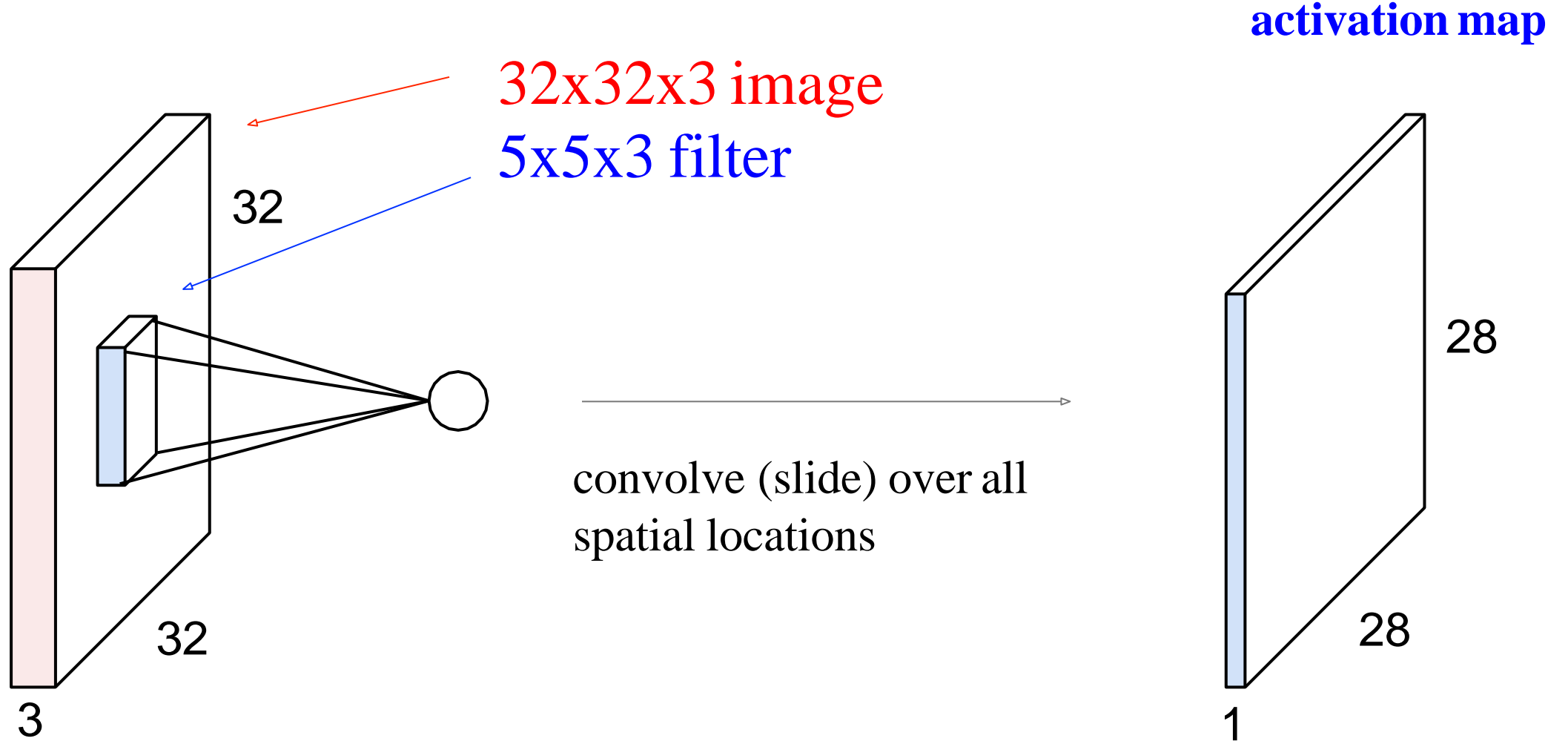


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

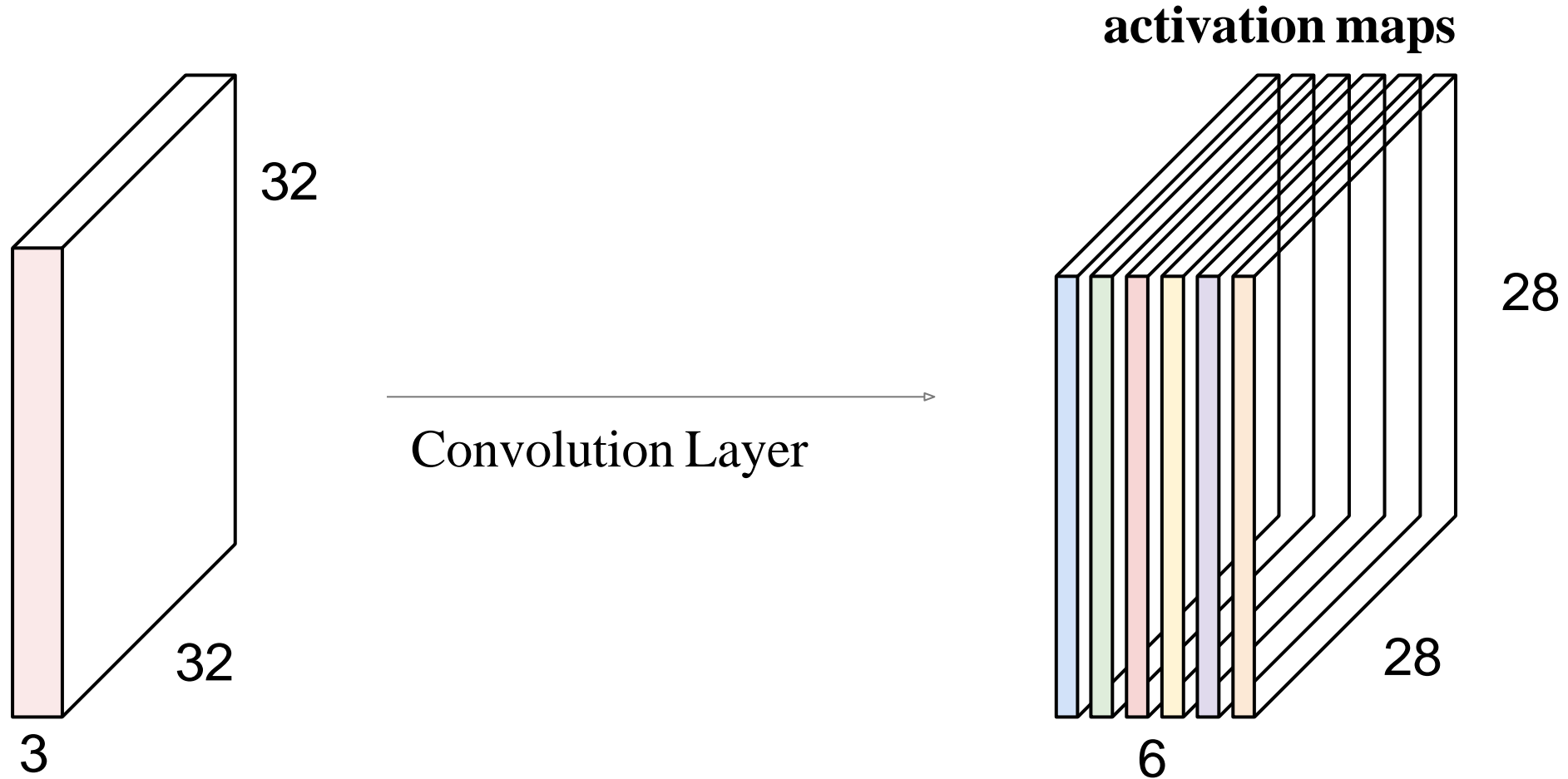
Convolution Layer



Convolution Layer

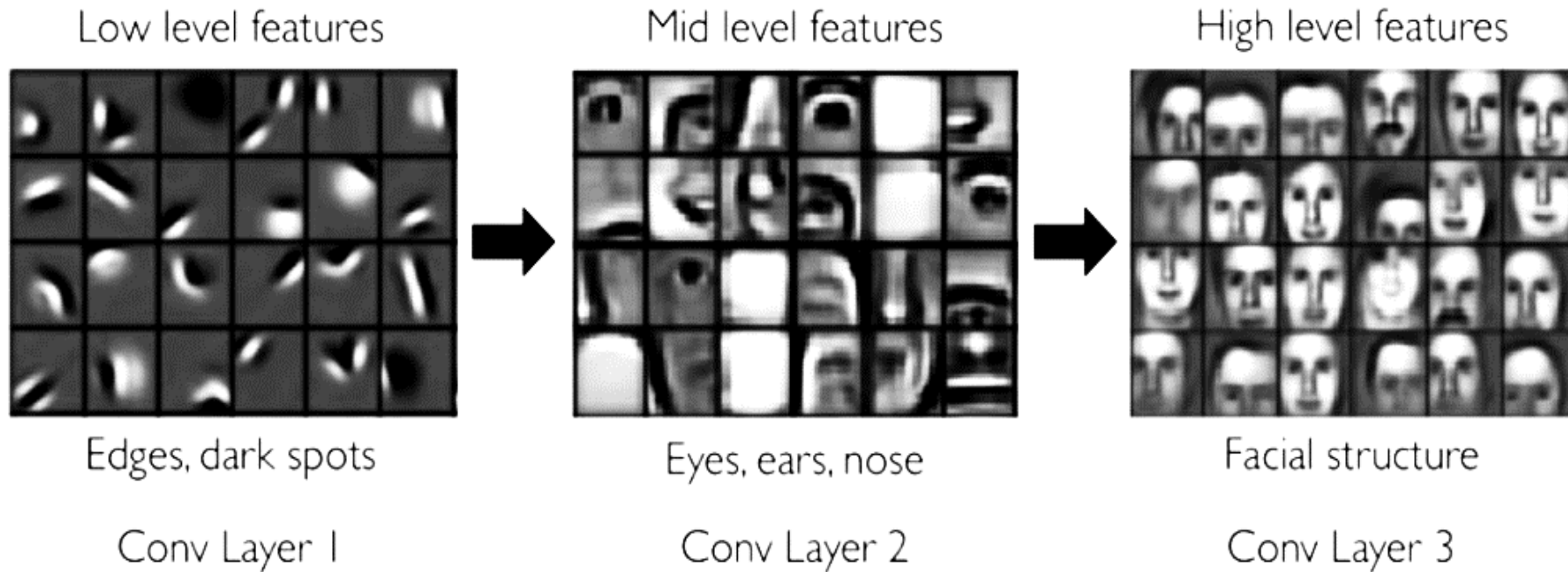


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



We stack these up to get a “new image” of size 28x28x6!

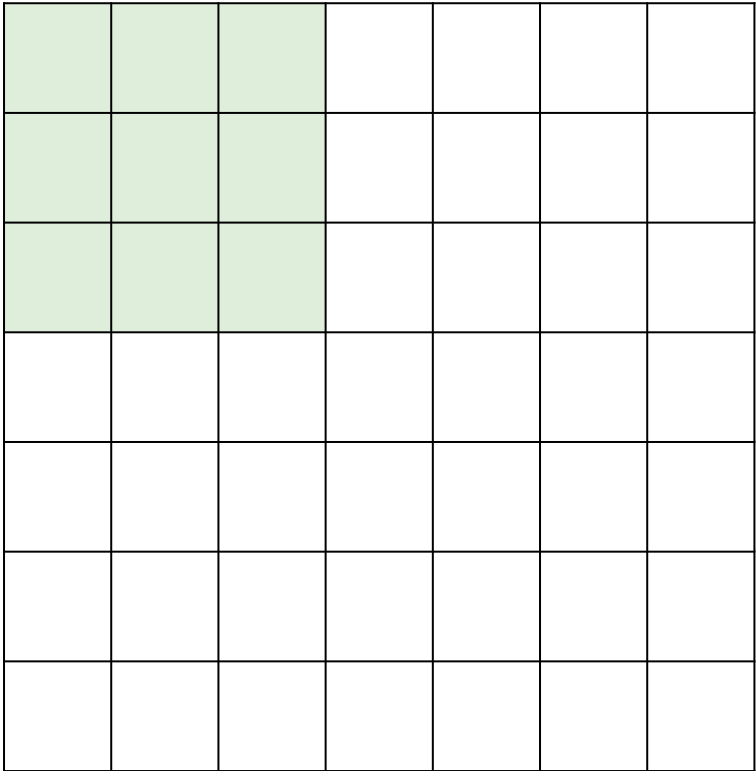
Representation Learning in Deep CNNs



Strides

A closer look at spatial dimensions:

7

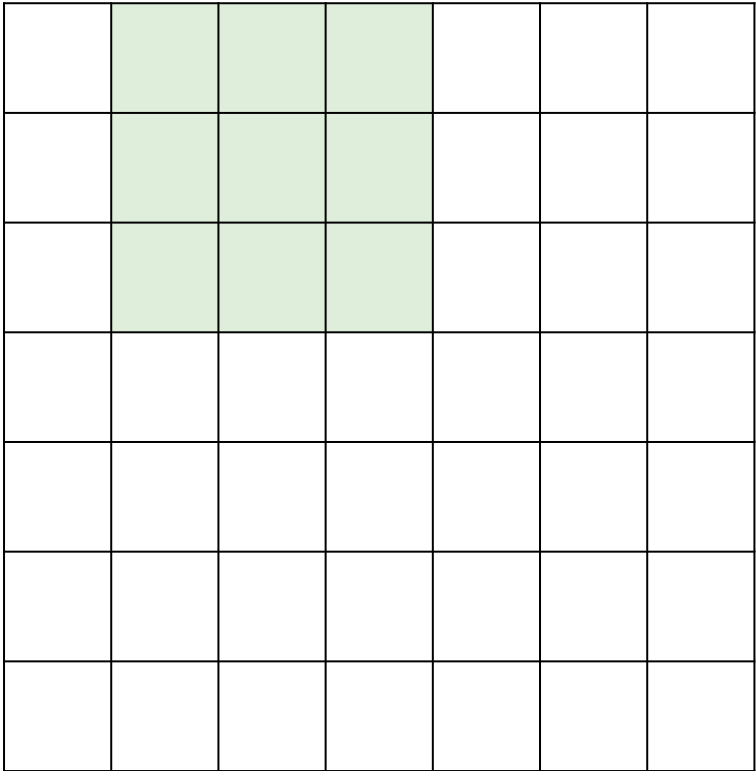


7x7 input (spatially) assume 3x3 filter

7

A closer look at spatial dimensions:

7

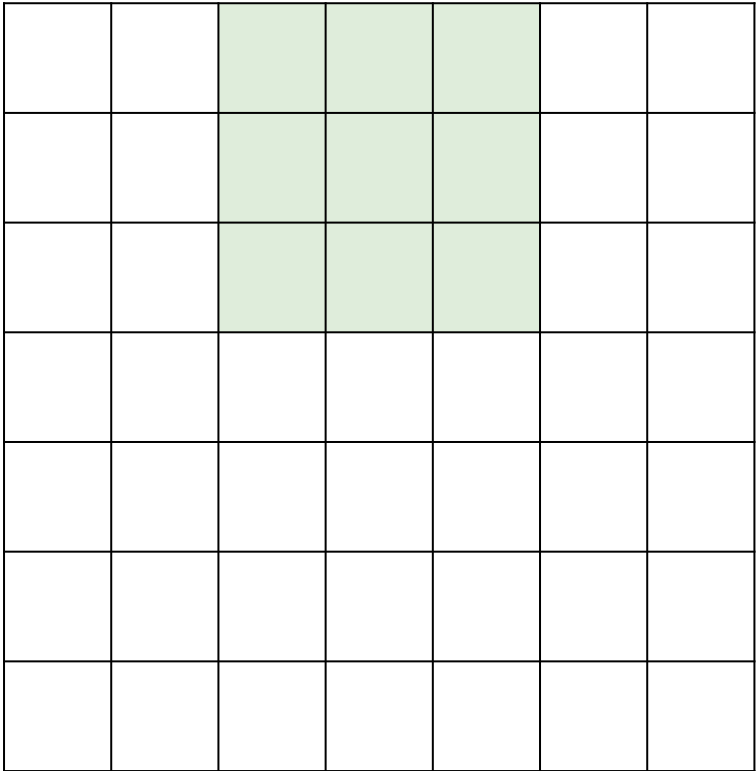


7

7x7 input (spatially) assume
3x3 filter

A closer look at spatial dimensions:

7



7

7x7 input (spatially) assume
3x3 filter

A closer look at spatial dimensions:

7

7

7x7 input (spatially) assume
3x3 filter

A closer look at spatial dimensions:

7

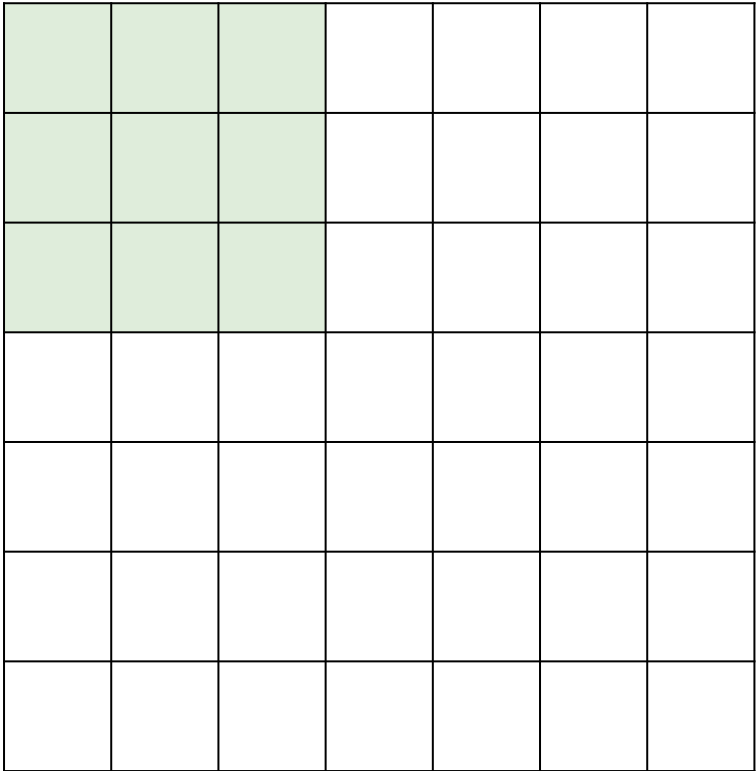
7

7x7 input (spatially) assume
3x3 filter

=> **5x5 output**

A closer look at spatial dimensions:

7



7

7x7 input (spatially) assume 3x3
filter applied **with stride 2**

A closer look at spatial dimensions:

7

7

7x7 input (spatially) assume
3x3 filter applied **with stride 2**

A closer look at spatial dimensions:

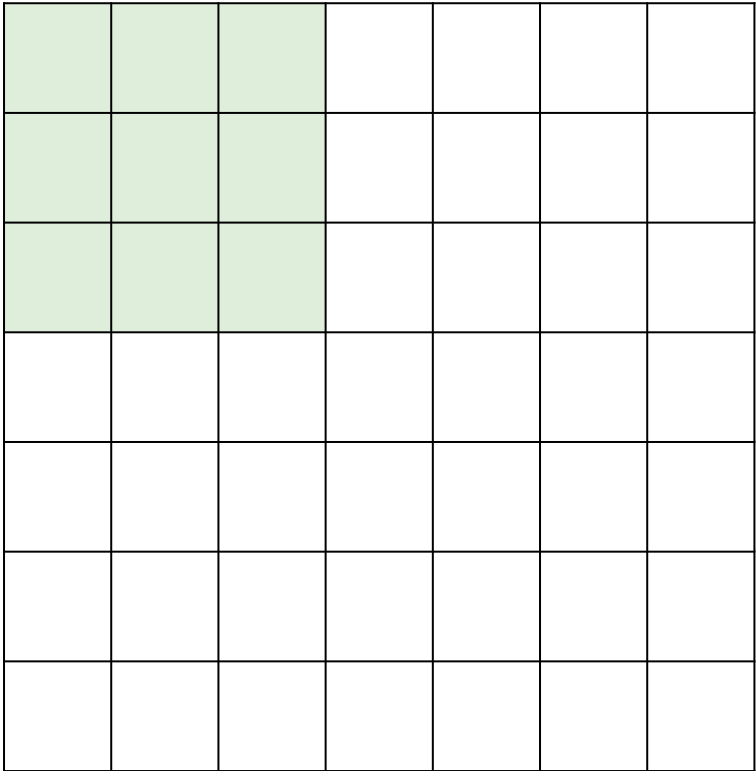
7

7

7x7 input (spatially) assume
3x3 filter applied **with stride 2**
=> **3x3 output!**

A closer look at spatial dimensions:

7



7

7x7 input (spatially) assume 3x3
filter applied **with stride 3?**

A closer look at spatial dimensions:

7

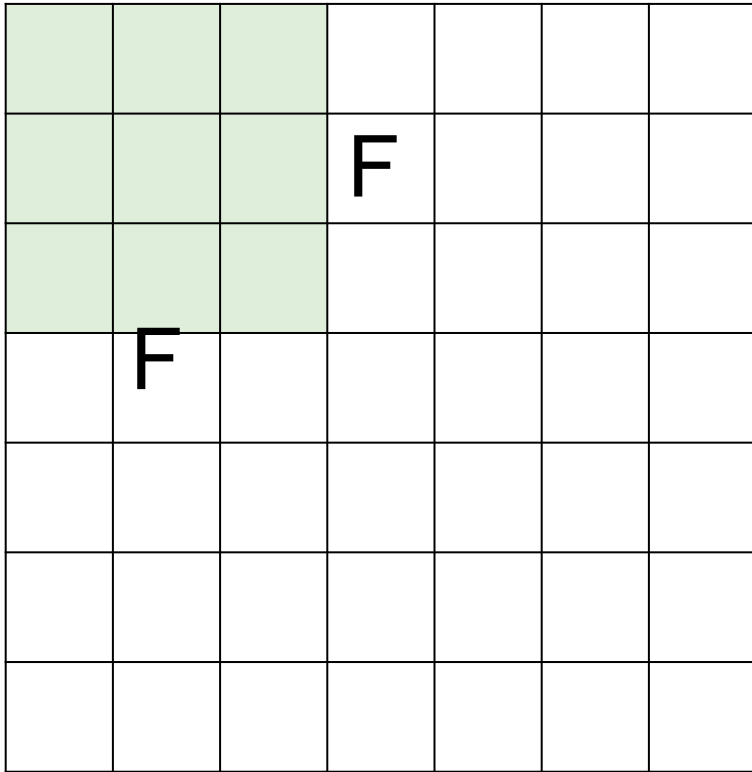
7

7x7 input (spatially) assume 3x3
filter applied **with stride 3**?

doesn't fit!

cannot apply 3x3 filter on 7x7 input
with stride 3.

N



Output size:

$$(\mathbf{N} - \mathbf{F}) / \mathbf{stride} + 1$$

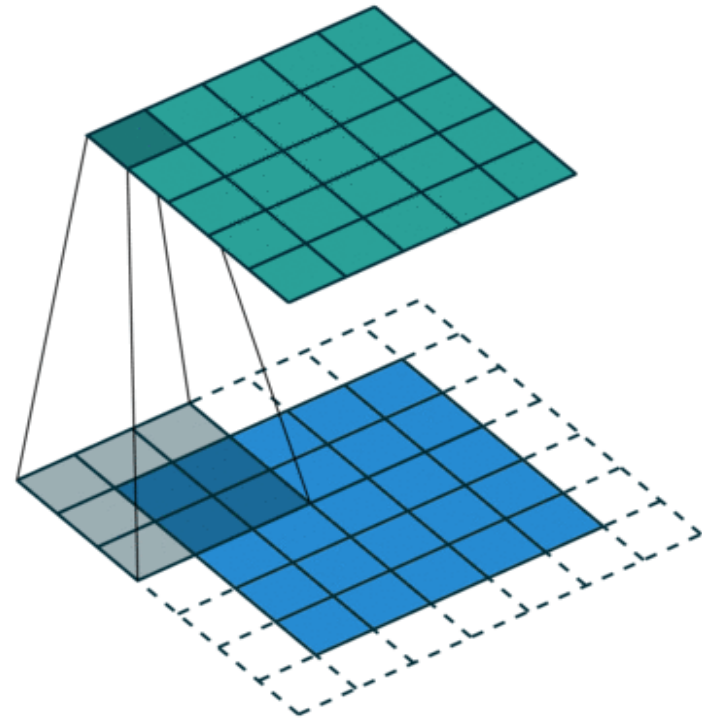
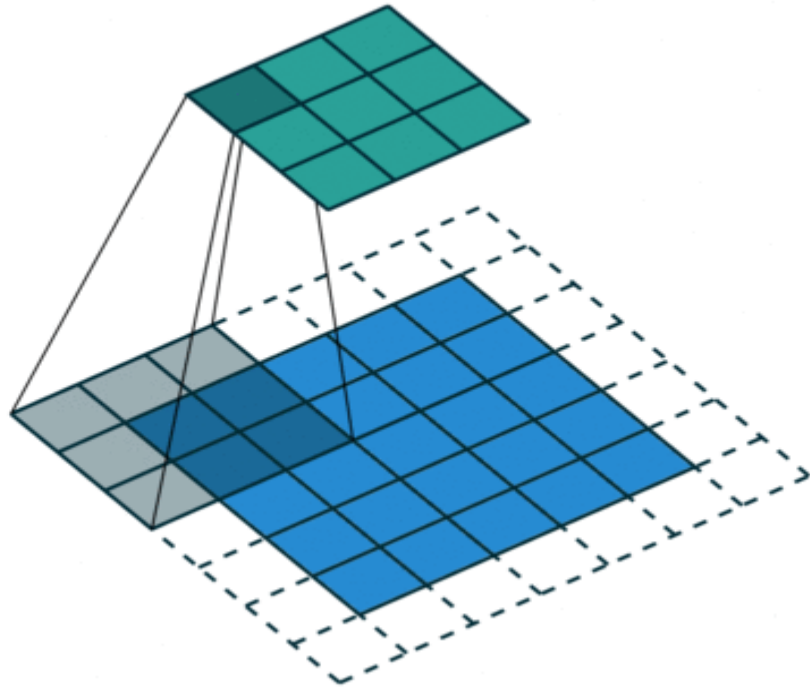
e.g. $N = 7, F = 3$:

$$\text{stride } 1 \Rightarrow (7 - 3) / 1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3) / 2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3) / 3 + 1 = 2.33$$

Zero-Padding



Zero-Padding: common to the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

(recall:)

$$(N - F) / \text{stride} + 1$$

Zero-Padding: common to the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

Zero-Padding: common to the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size $F \times F$, and zero-padding with $(F-1)/2$. (will preserve size spatially)

e.g. $F = 3 \Rightarrow$ zero pad with 1 $F = 5$

\Rightarrow zero pad with 2 $F = 7 \Rightarrow$

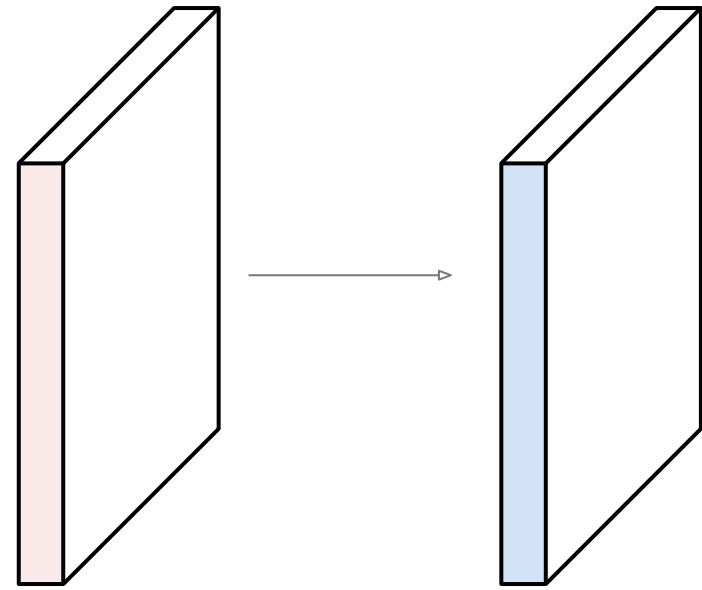
zero pad with 3

Examples:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Output volume size: ?



Examples time:

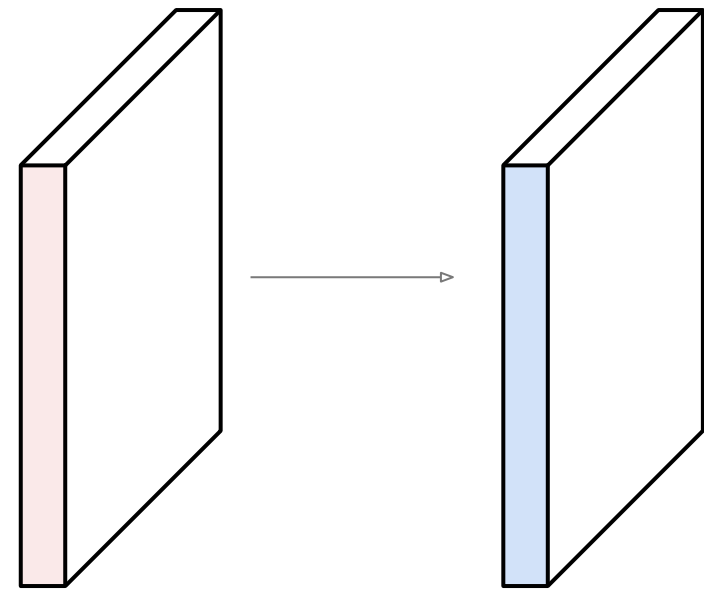
Input volume: **32x32x3**

10 **5x5** filters with stride **1**, pad **2**

Output volume size:

$(32+2*2-5)/1+1 = 32$ spatially, so

32x32x10

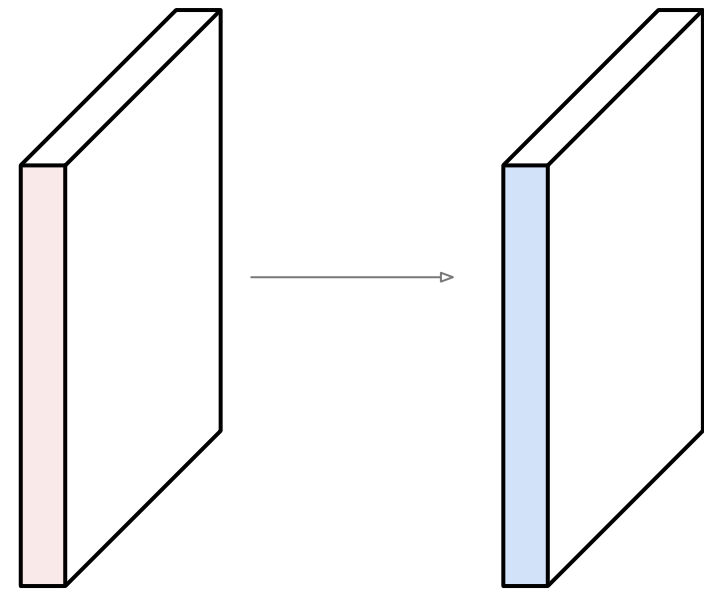


Examples :

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2

Number of parameters in this layer?



Examples :

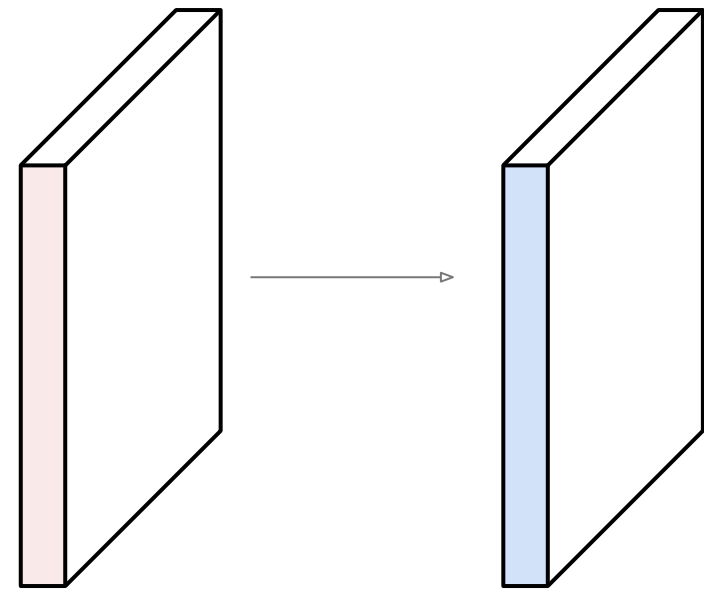
Input volume: **32x32x3**

10 **5x5** filters with stride 1, pad 2

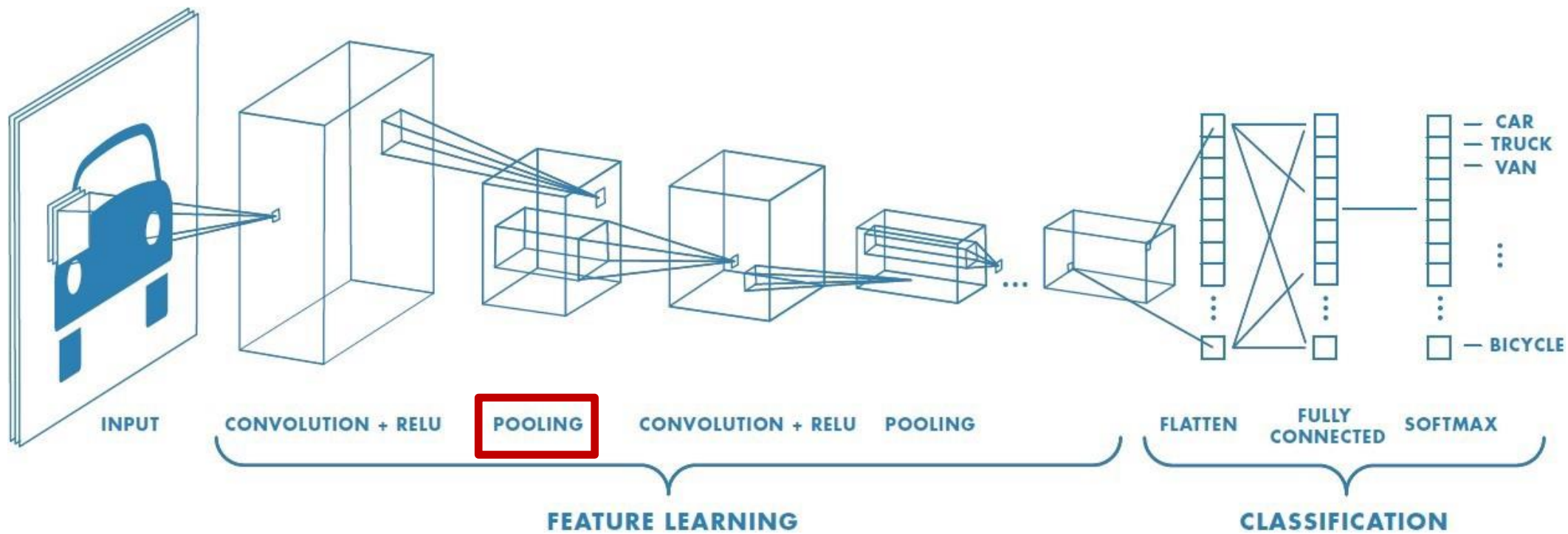
Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params

=> $76*10 = 760$

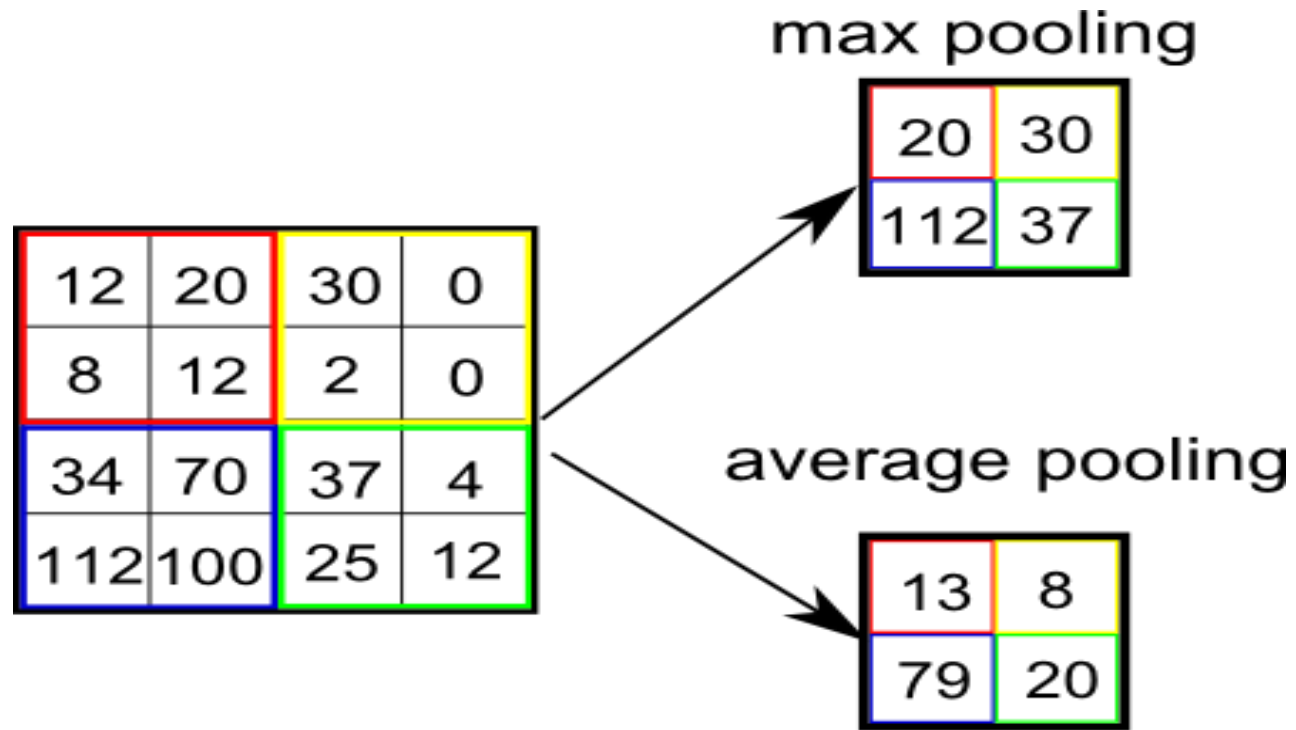


(+1 for bias)



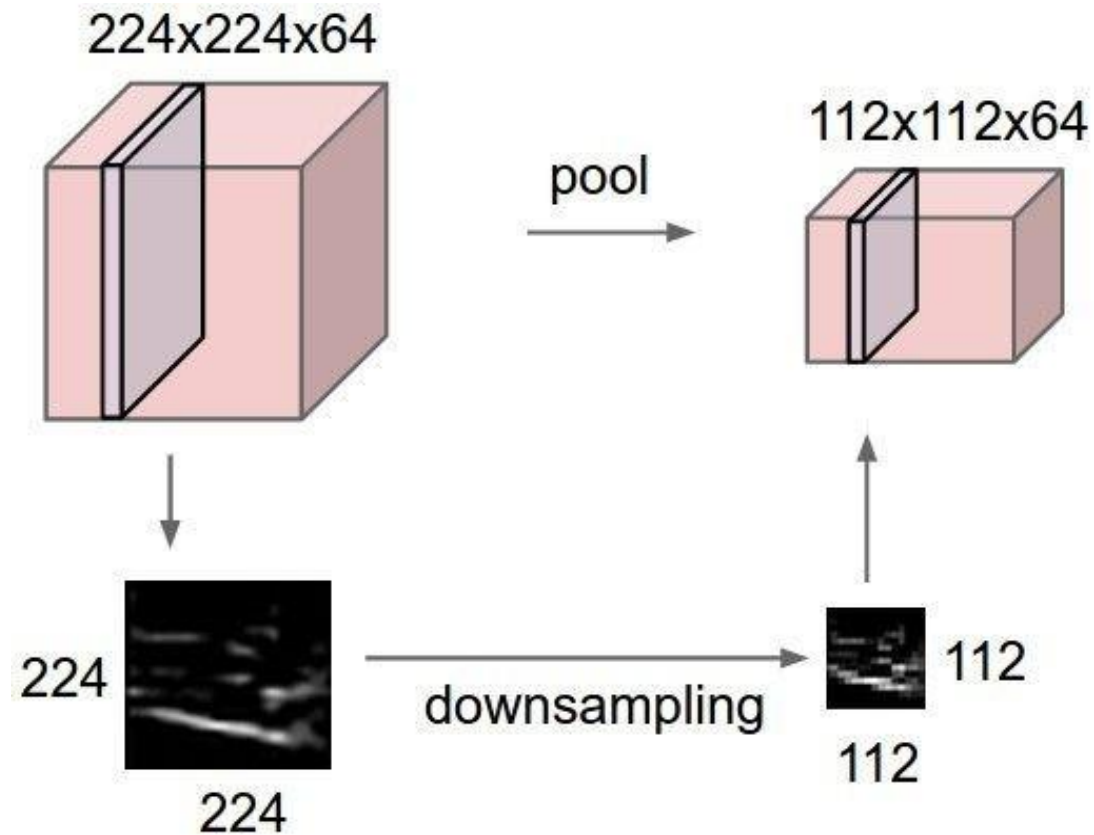
Pooling

Pooling



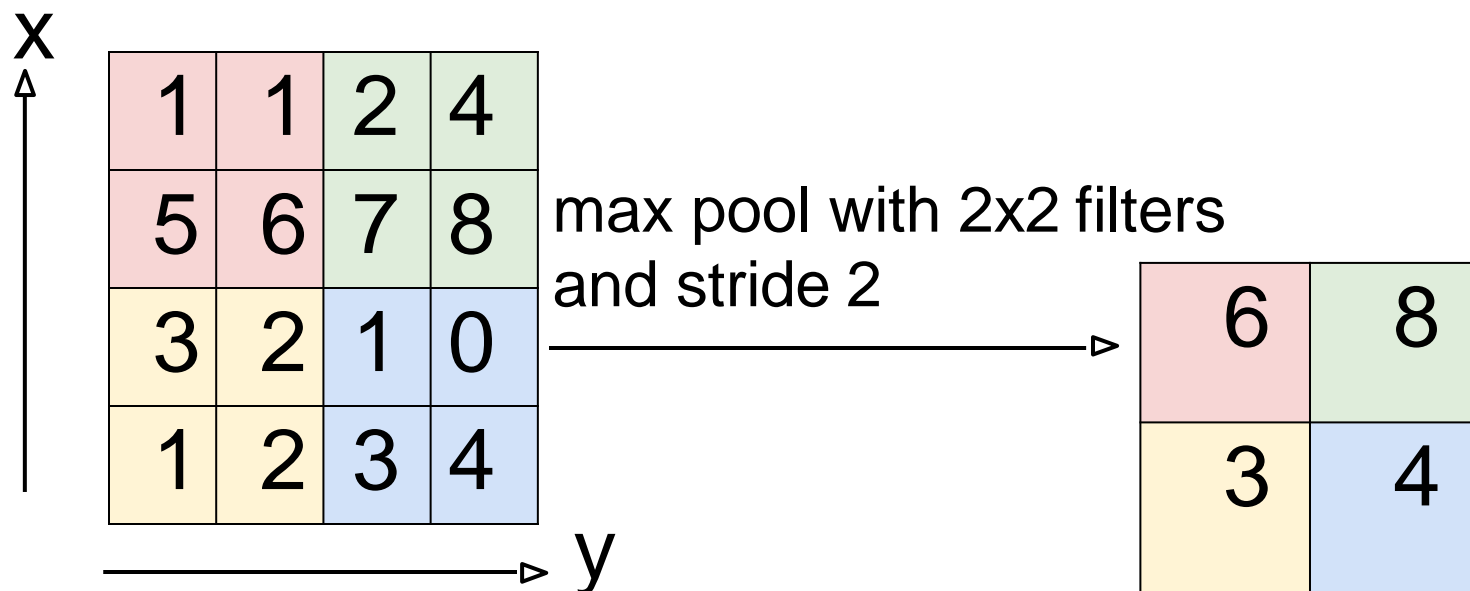
Pooling

- makes the representations smaller and more manageable
- operates over each activation map independently



Max Pooling

Single depth slice

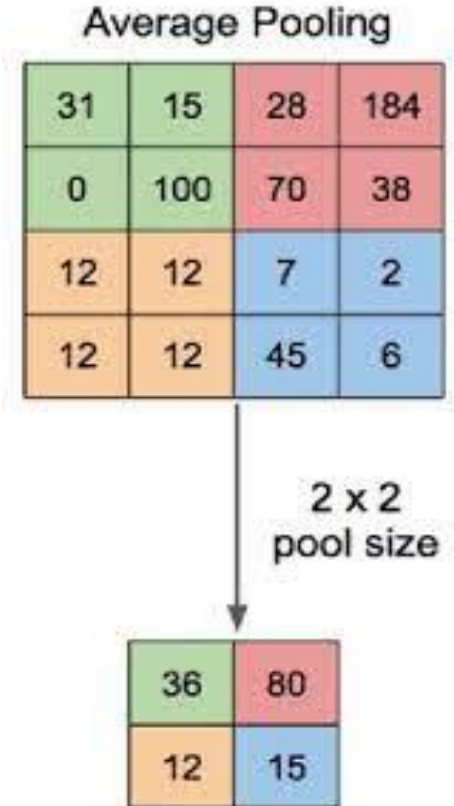


3.0	3.0	3.0
3.0	3.0	3.0
3.0	2.0	3.0

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

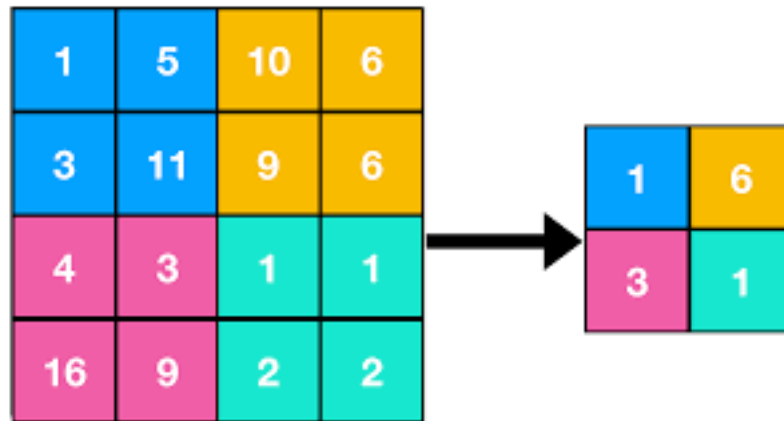
Average Pooling

- Image patches equivalent to kernel size are made
- **Average value** in each patch is the pooled value

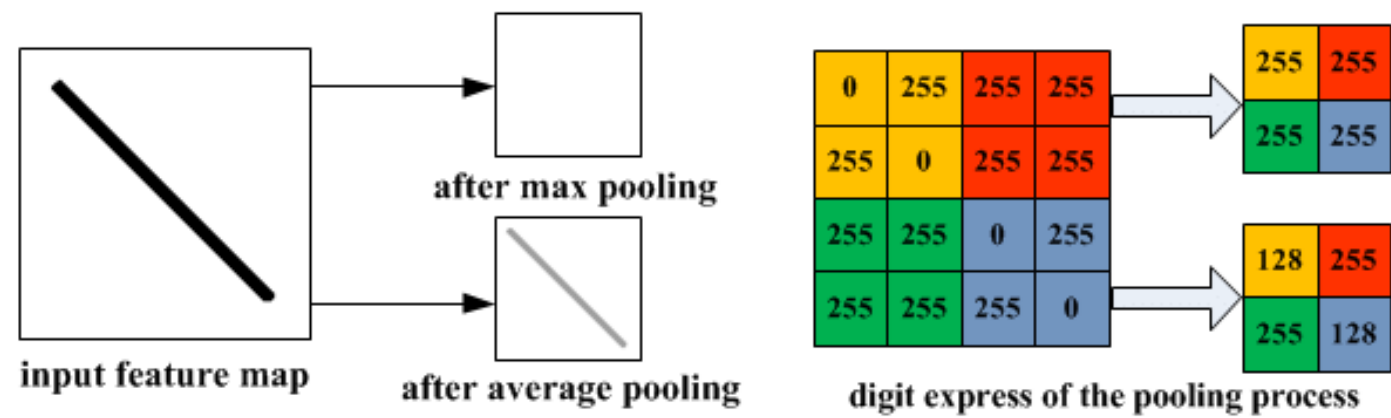


Min Pooling

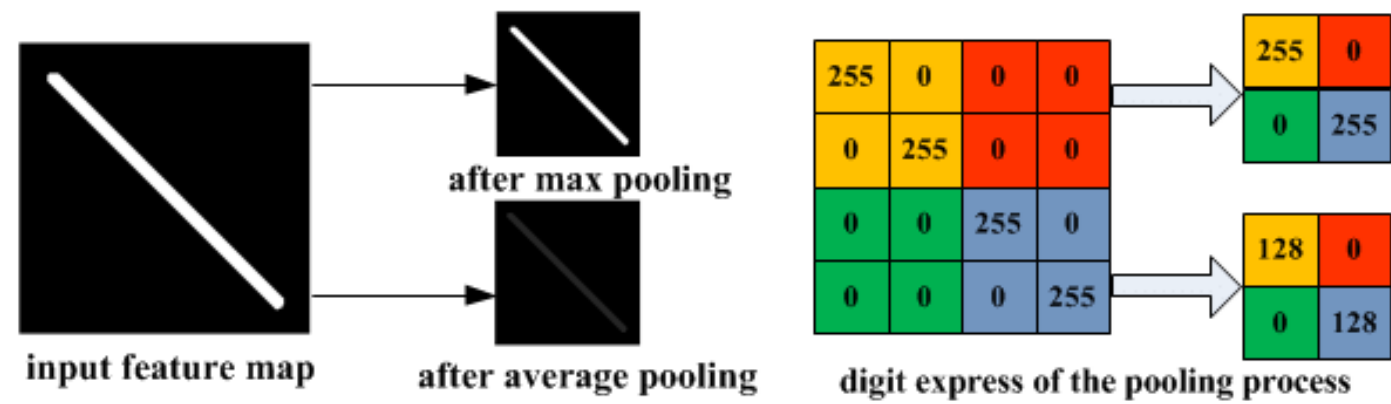
- Image patches equivalent to kernel size are made
- **Min value** in each patch is the pooled value



Pooling comparison

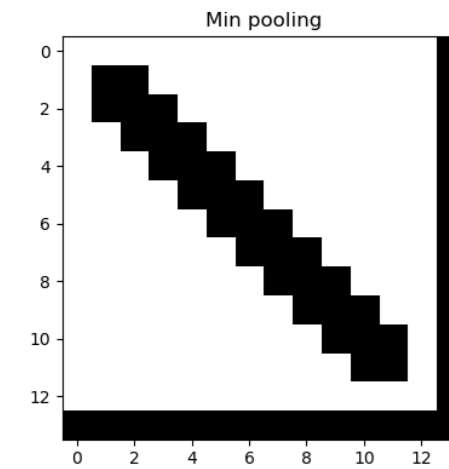
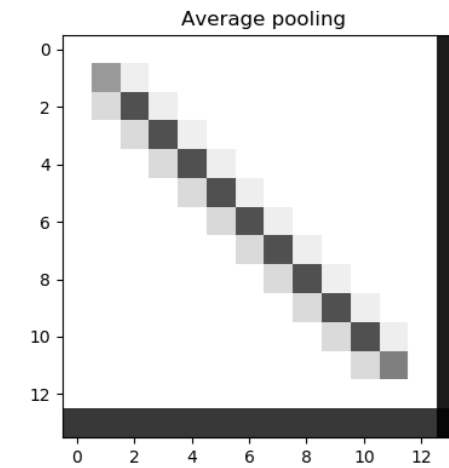
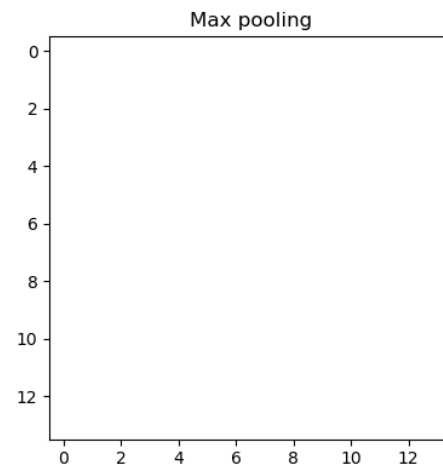
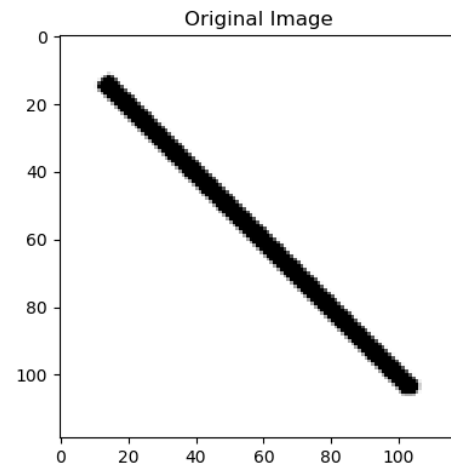


(a) Illustration of max pooling drawback



(b) Illustration of average pooling drawback

Pooling comparison

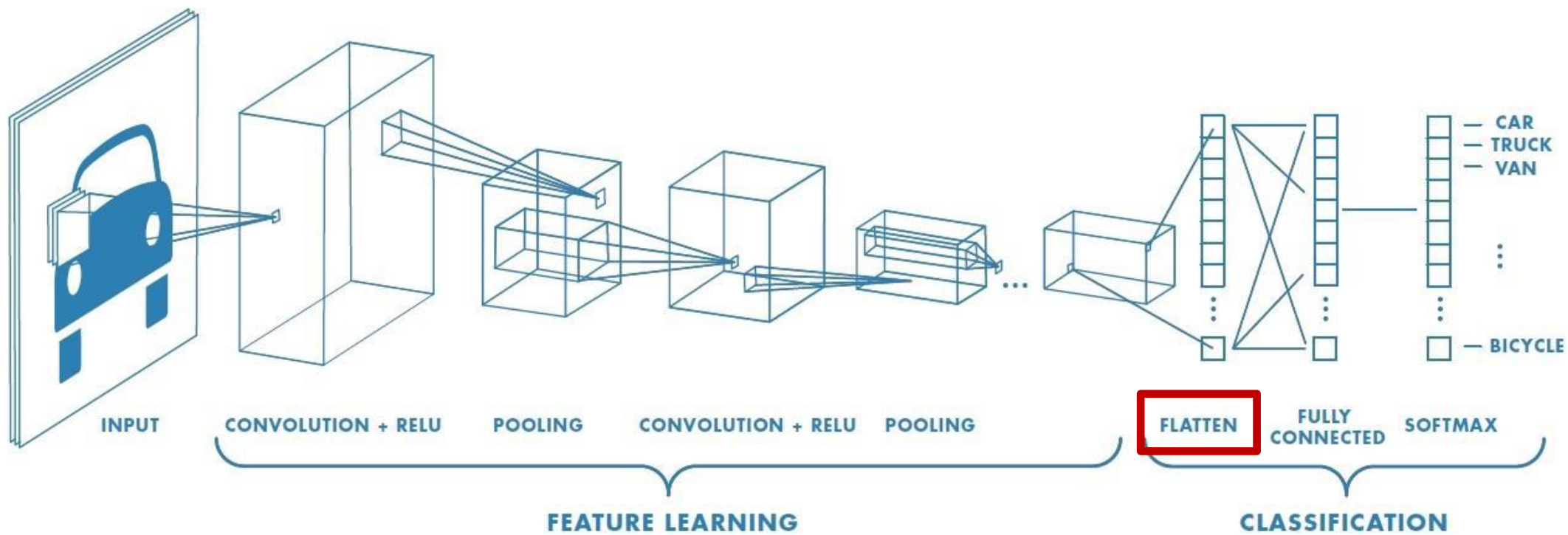


Pros

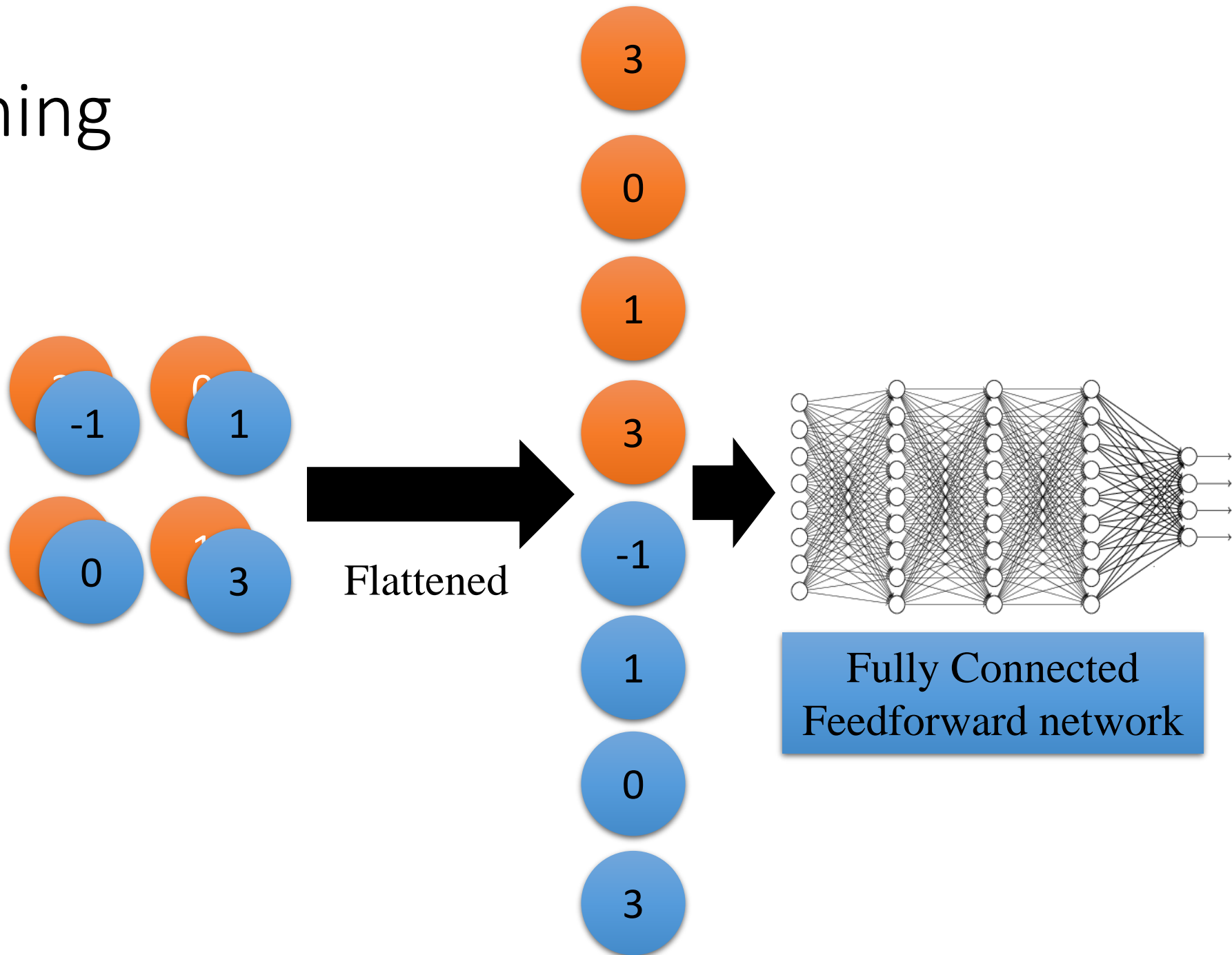
- Computational Efficiency
- Dimensionality Reduction
- Noise Reduction
- Translation invariance

Cons

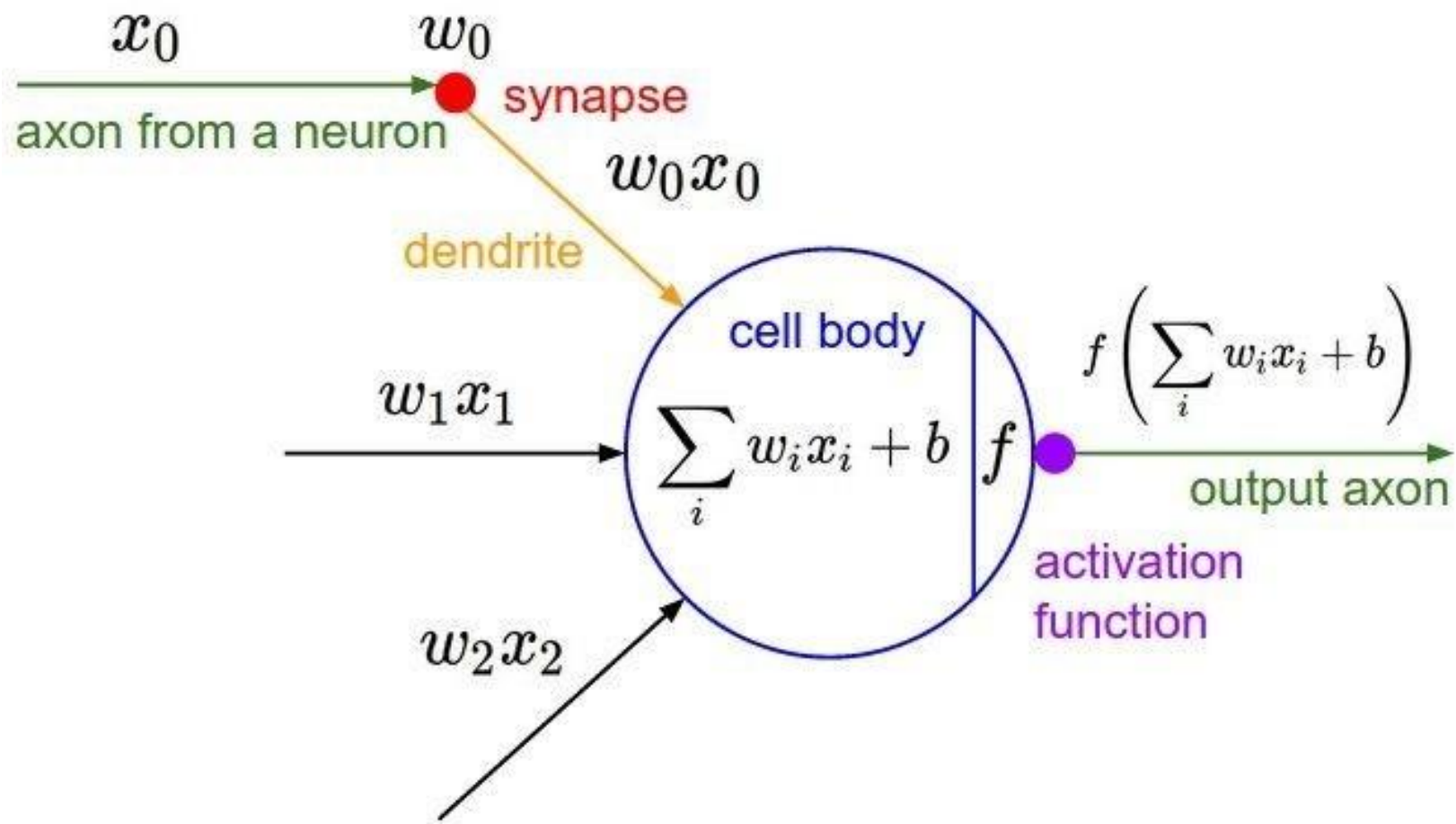
- Information loss
- Sensitivity to Hyperparameters
- Pooling artifacts



Flattening



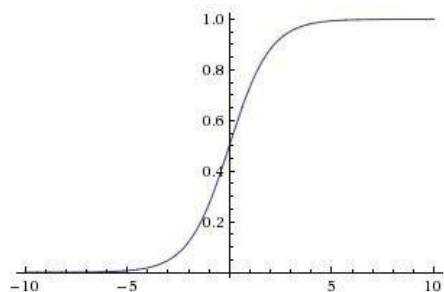
Activation Functions



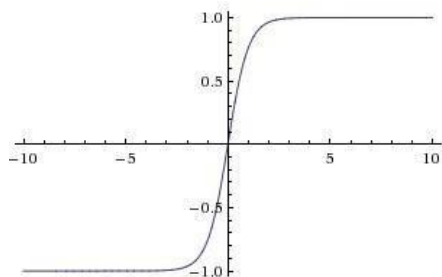
Activation Functions

Sigmoid

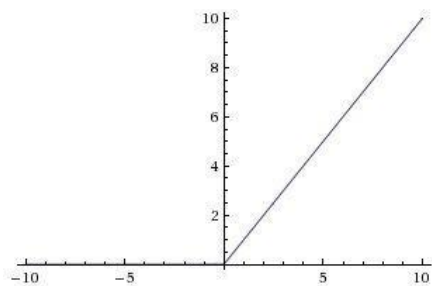
$$\sigma(x) = 1/(1 + e^{-x})$$



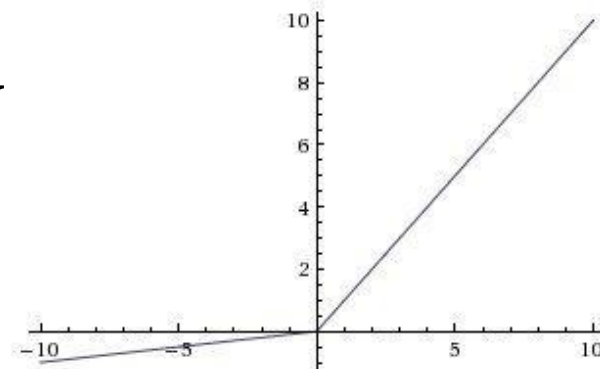
tanh $\tanh(x)$



ReLU $\max(0, x)$

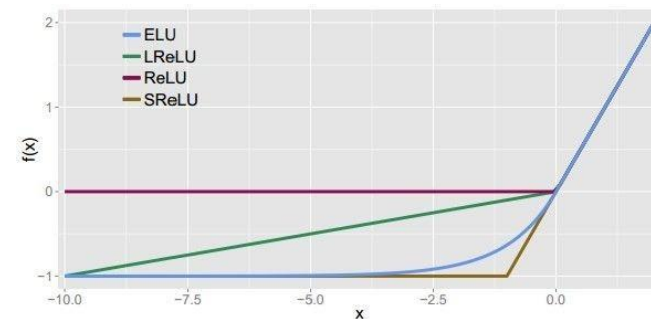


Leaky ReLU

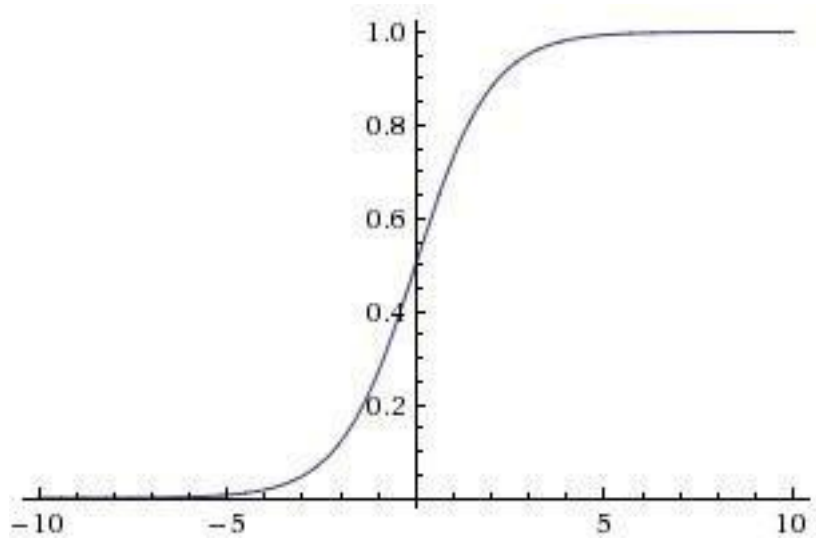


Maxout ELU $\max(w_1^T x + b_1, w_2^T x + b_2)$

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



Activation Functions

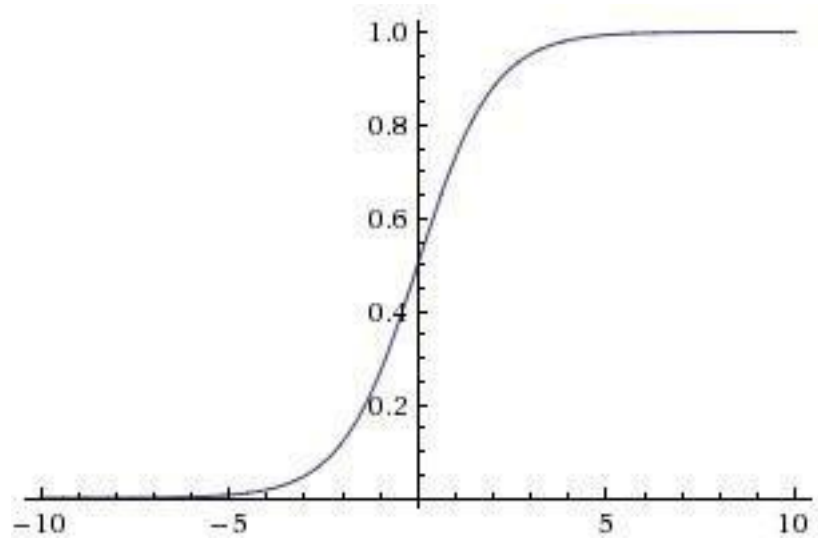


- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

Activation Functions

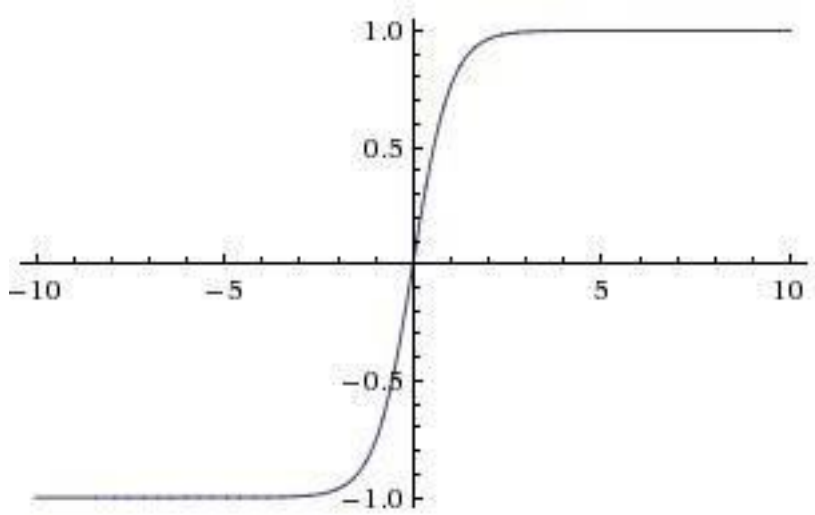


Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
 - Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron
1. Sigmoid outputs are not zero-centered
 2. $\exp()$ is a bit compute expensive

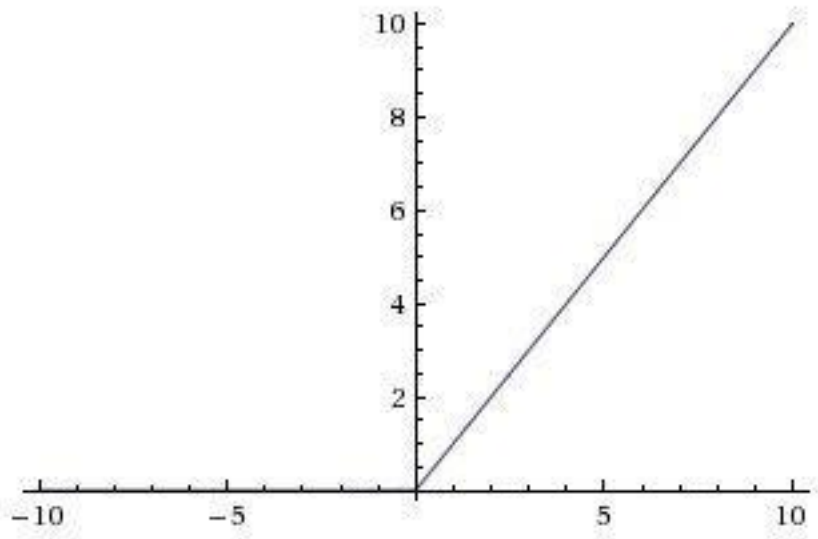
Activation Functions



$\tanh(x)$

- Squashes numbers to range $[-1,1]$
- zero centered (nice)
- still kills gradients when saturated :(

Activation Functions



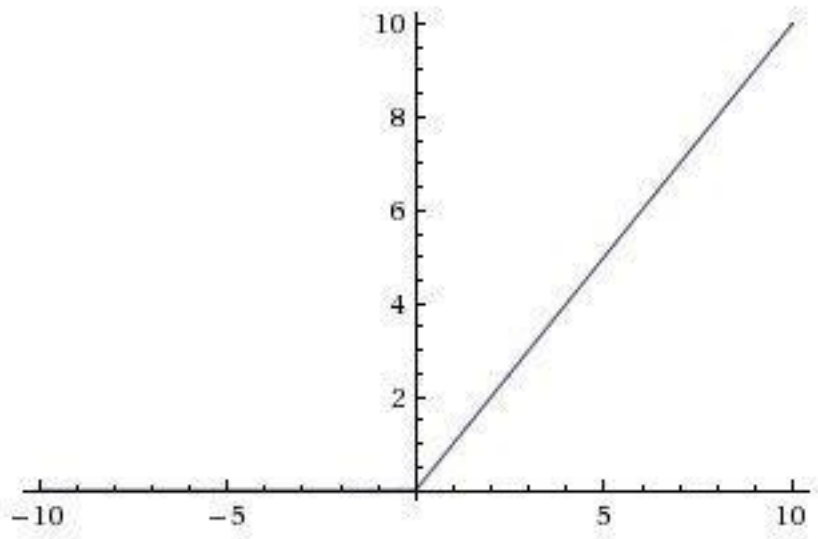
ReLU

(Rectified Linear Unit)

Computes $f(x) = \max(0, x)$

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

Activation Functions



ReLU

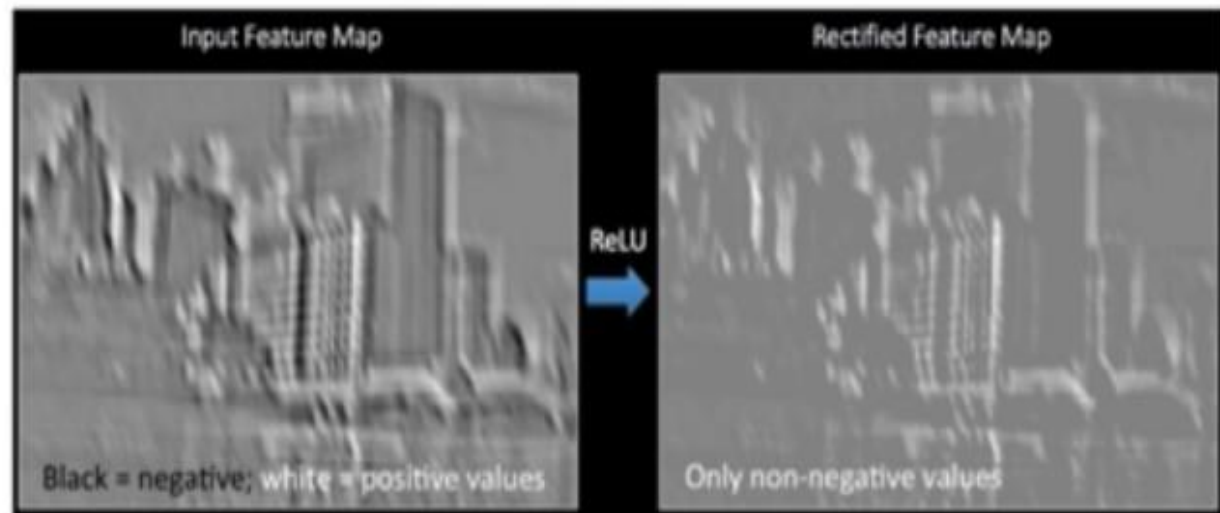
(Rectified Linear Unit)

Computes $f(x) = \max(0, x)$

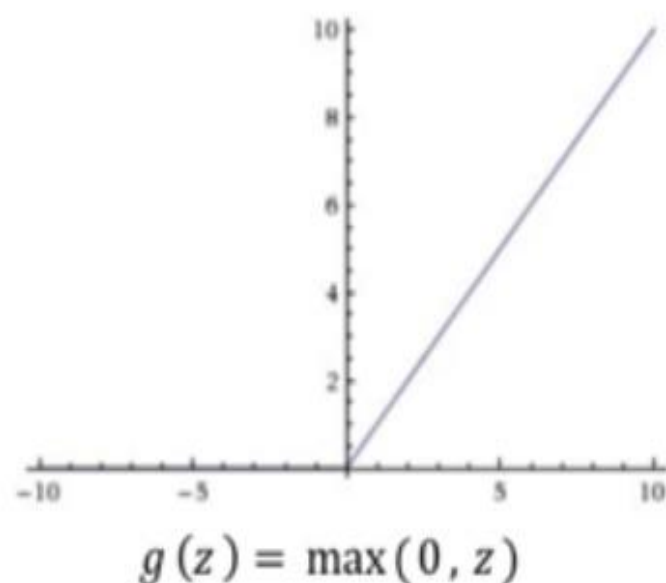
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Not zero-centered output

Introducing Non-Linearity

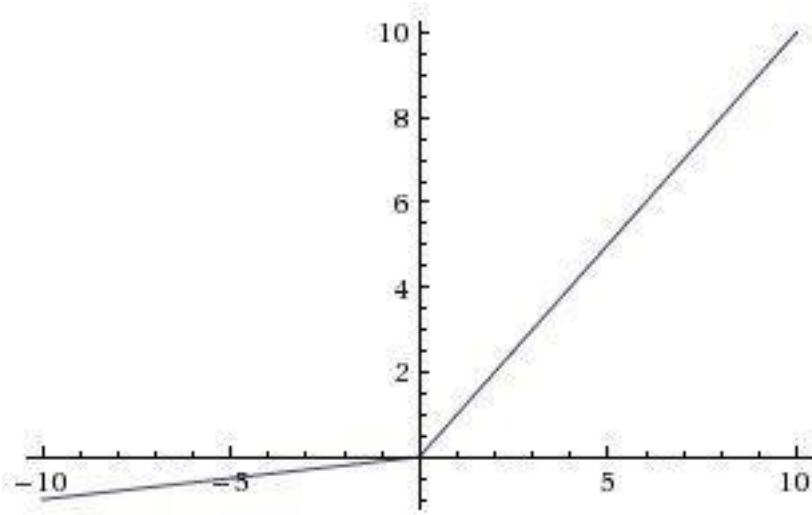
- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**



Rectified Linear Unit (ReLU)



Activation Functions



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Leaky ReLU

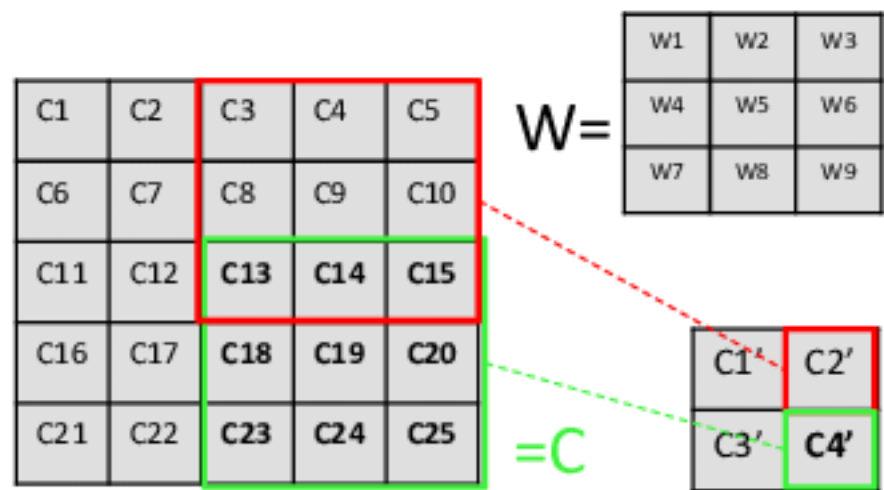
$$f(x) = \max(0.01x, x)$$

Deconvolution: An Introduction

- Convolution reduces the dimension of data and extracts the features
- The image is translated to feature space
- How to represent features in image domain?

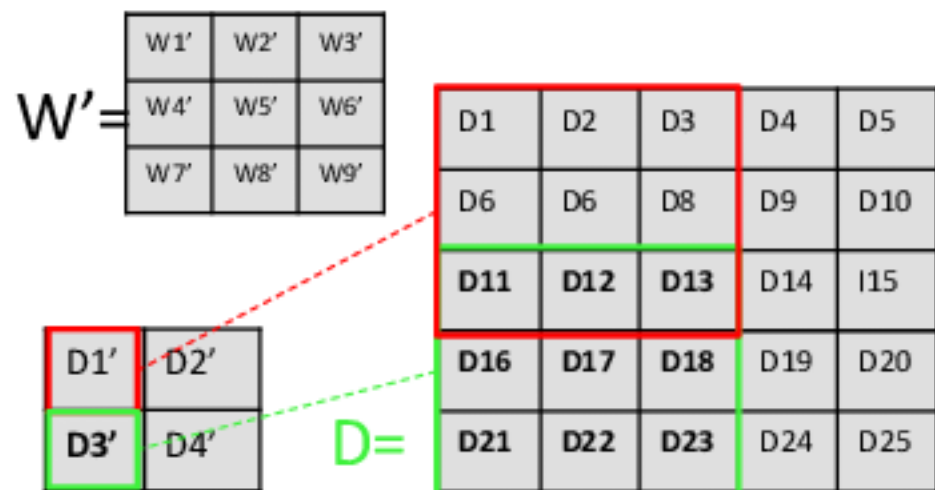
Deconvolution: An Introduction

- Deconvolution helps in achieving the inverse operation



$$C4' = W * C$$

Convolution (stride=2)

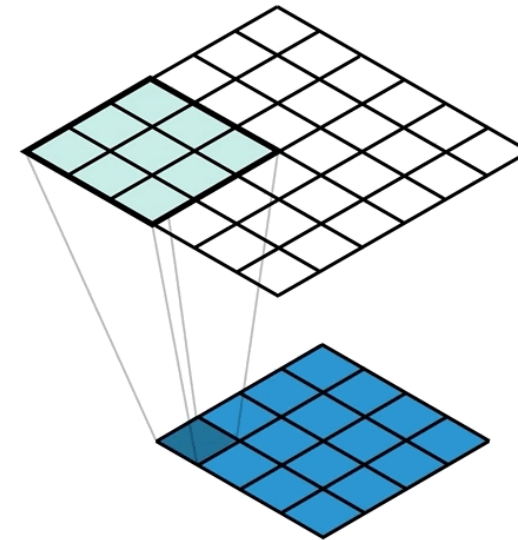


$$D = D3' \times W'$$

Deconvolution (stride=2)

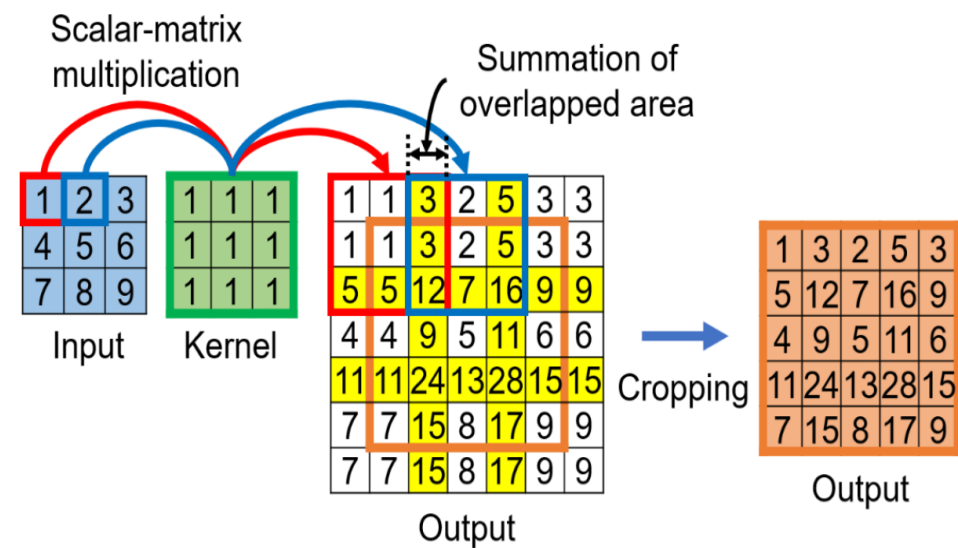
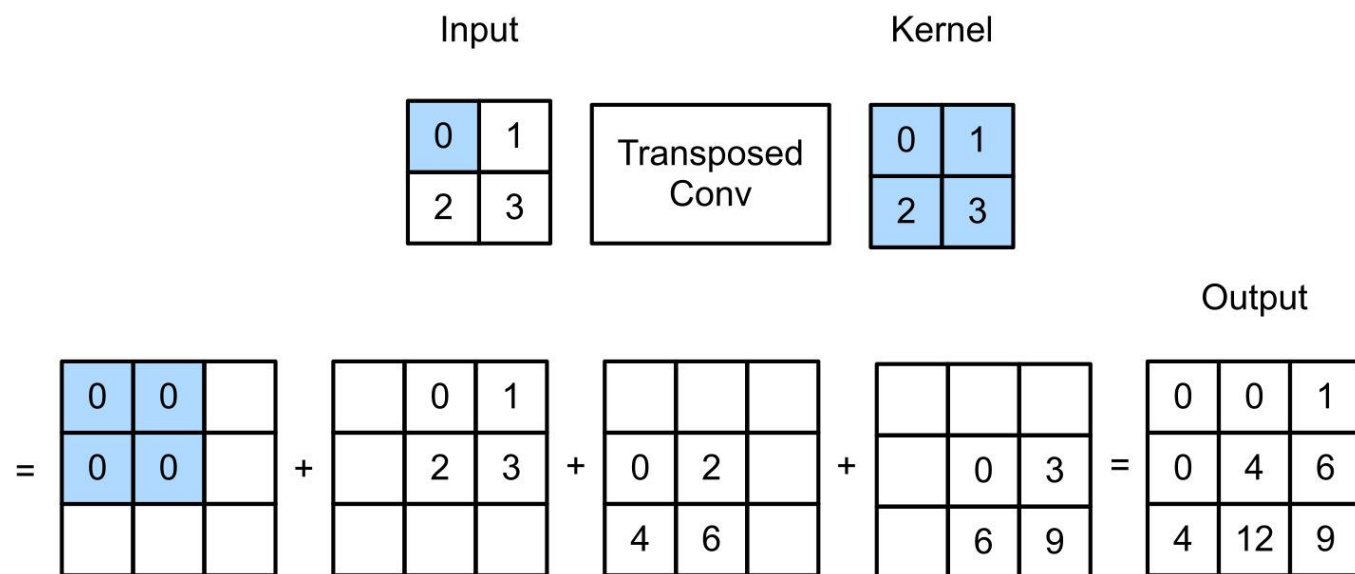
Deconvolution: Brief Concept

- Unlike convolution layers, which reduces spatial dimension
- Deconvolution Upsample the spatial resolution of input data

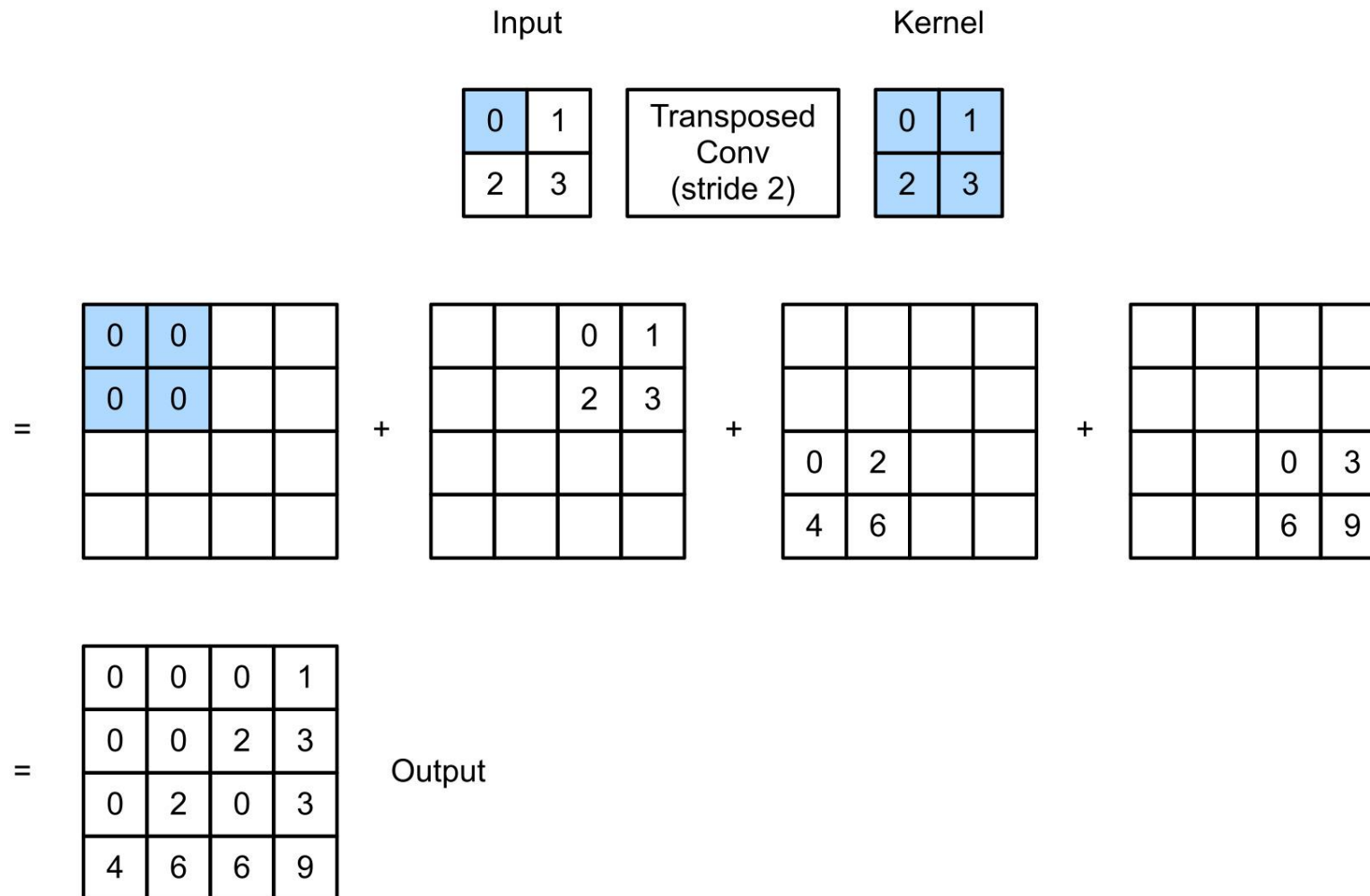


Types of Deconvolution: Padding-free

- No padding is required
- The output is sum of overlapped region
- Output Size= stride(inputsize-1)+filtersize-(2*padding)

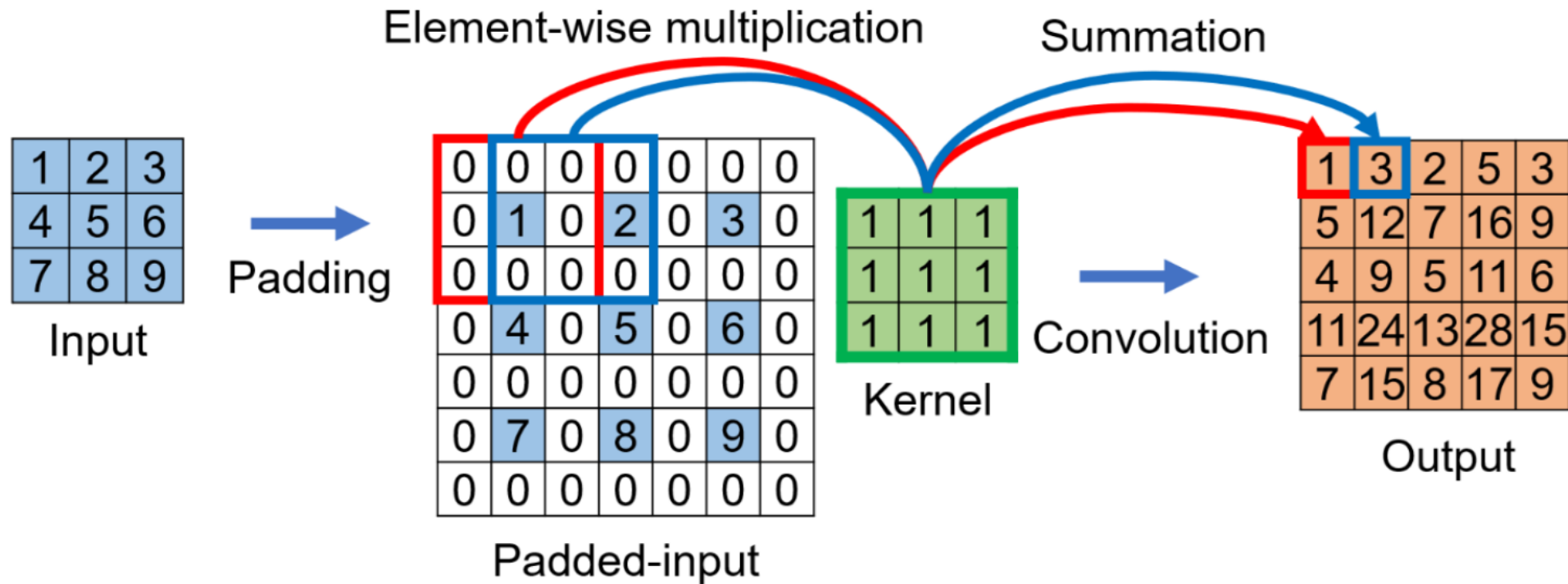


Types of Deconvolution: Padding-free with stride-2



Types of Deconvolution: Zero Padding

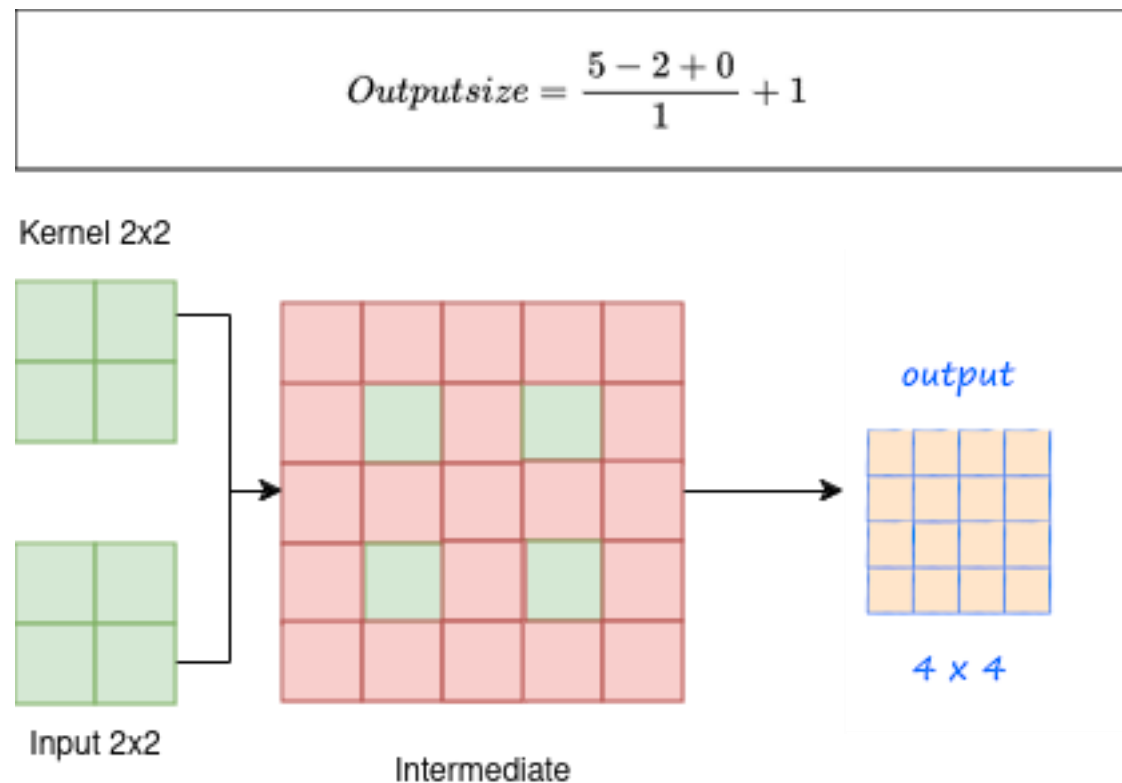
- Input image is padded across each element and convolution is applied



Output size calculation

$$n_{out} = \left\lfloor \frac{n_{in} + 2p - k}{s} \right\rfloor + 1$$

n_{in} : number of input features
 n_{out} : number of output features
 k : convolution kernel size
 p : convolution padding size
 s : convolution stride size



Thank You