

Convolutional Neural Networks

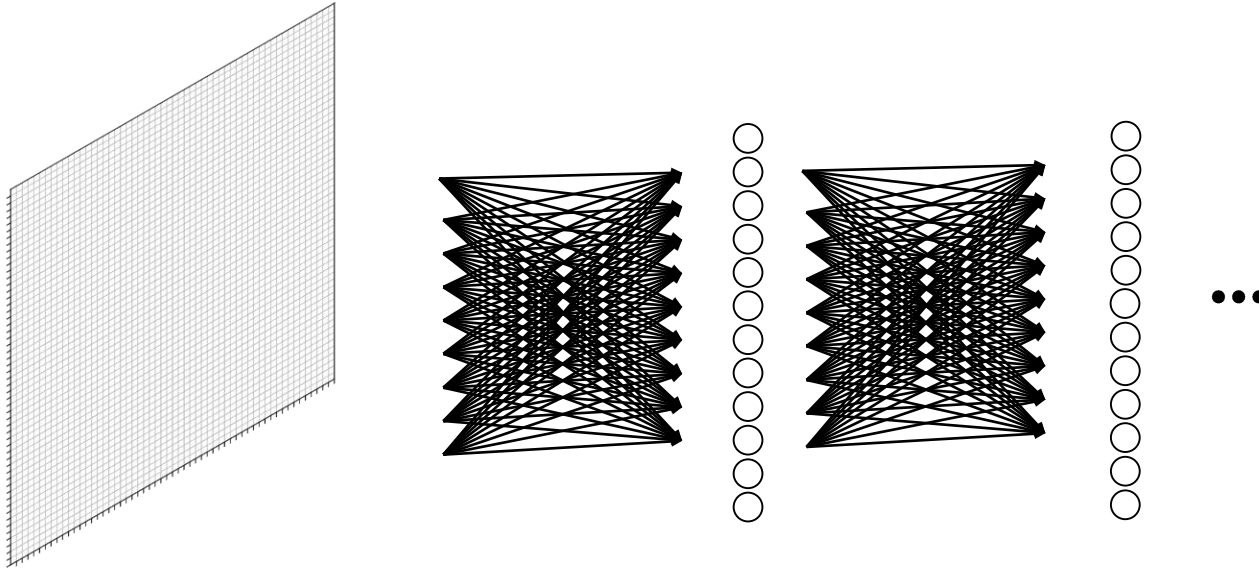
Neural networks inspired
by the **visual cortex**

Faculty of Mathematics and Computer Science, University of Bucharest
and
Sparktech Software

Academic Year 2018/2019, 1st Semester

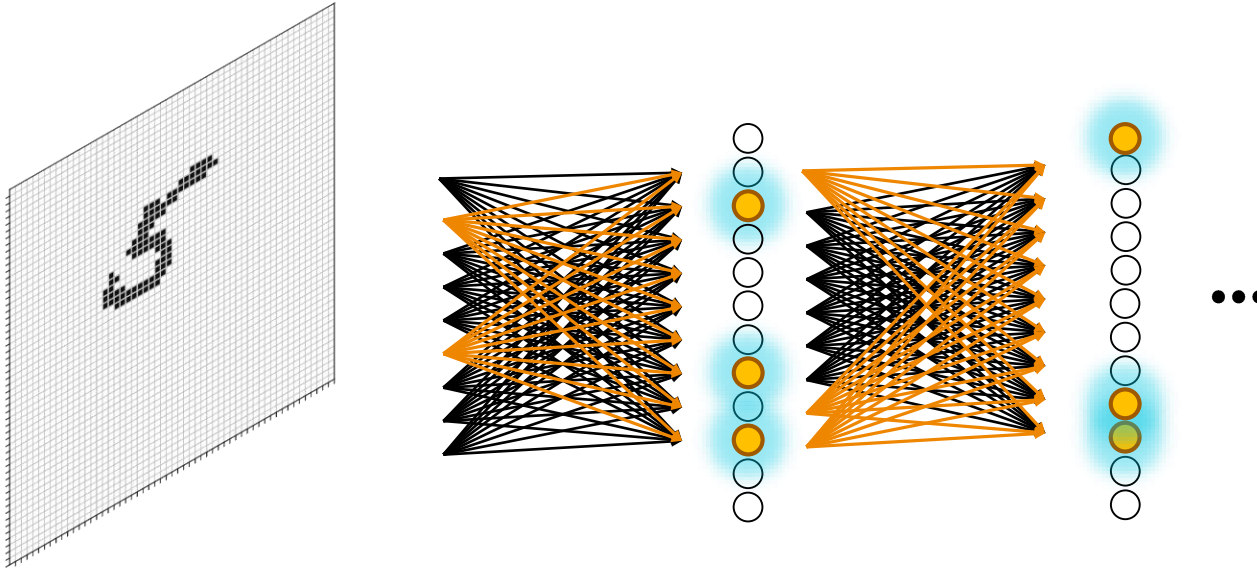
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).



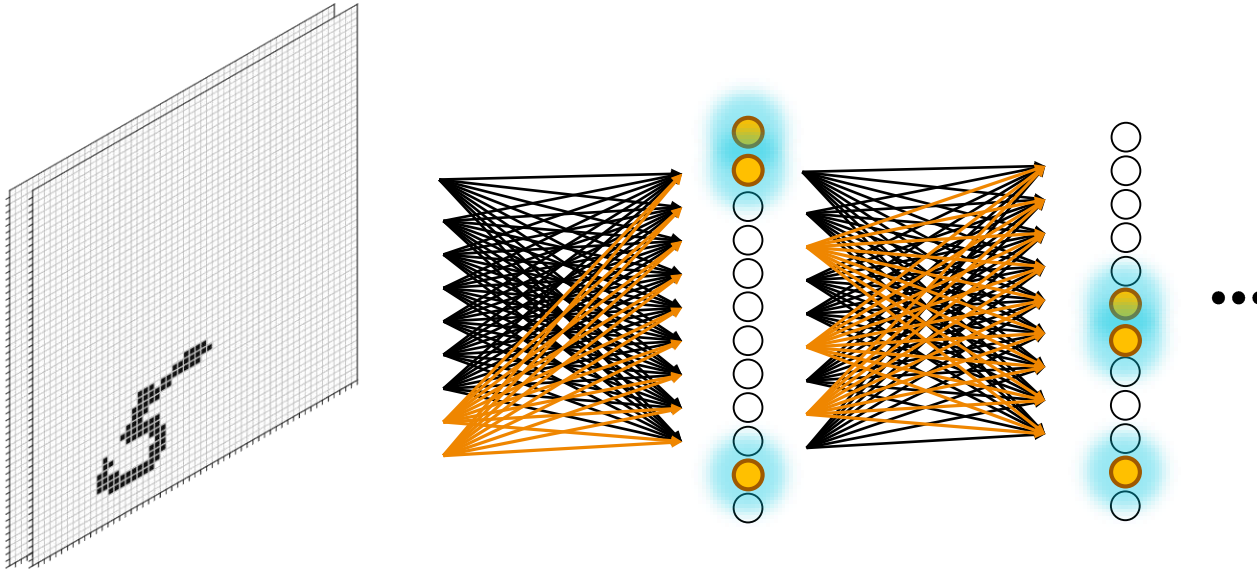
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).



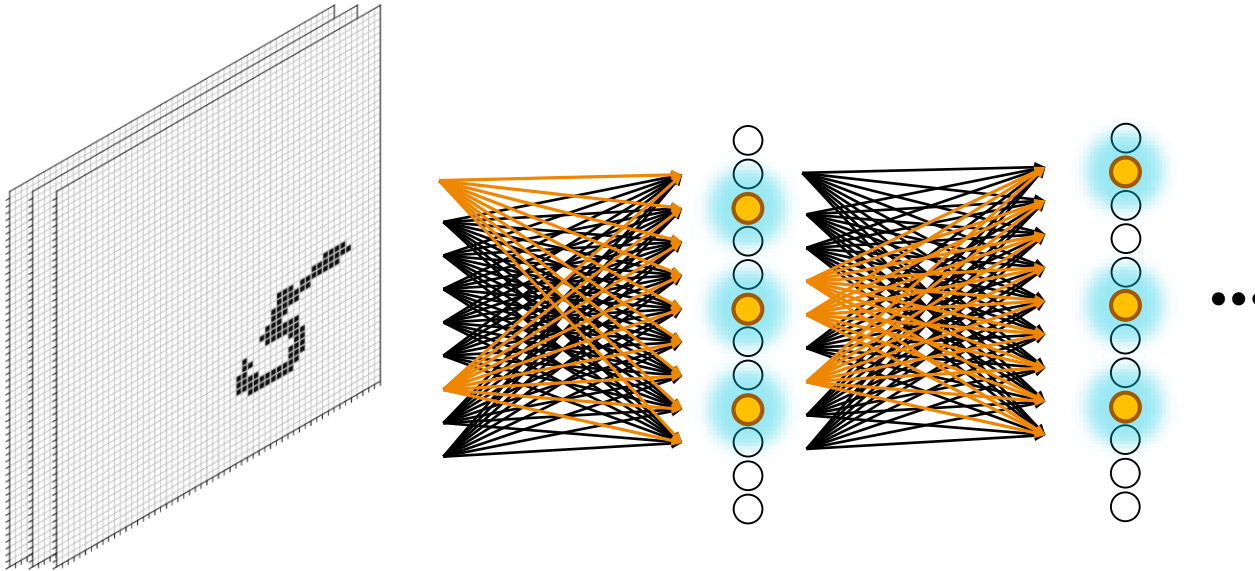
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).



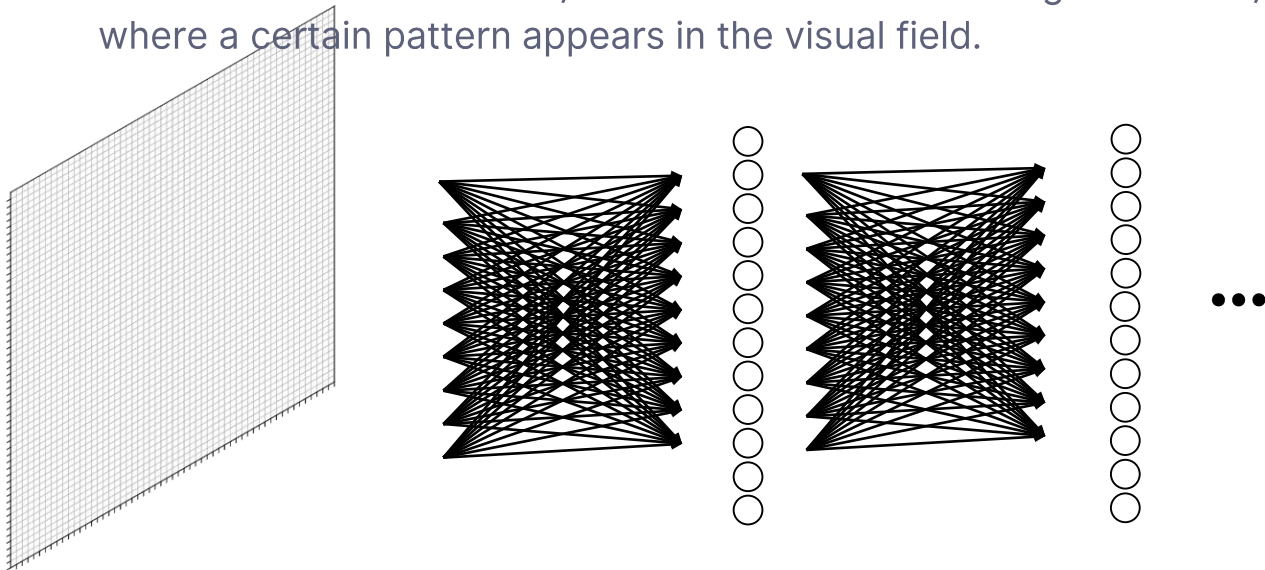
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).



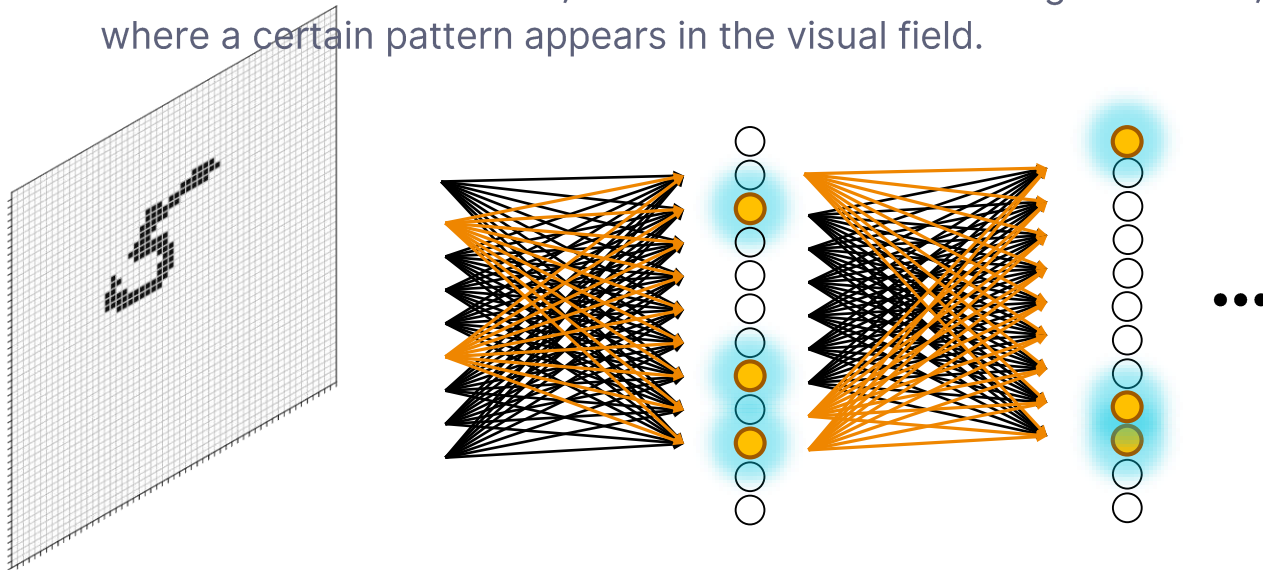
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).
- In the *human visual cortex*, the same neurons tend to get activate, regardless of where a certain pattern appears in the visual field.



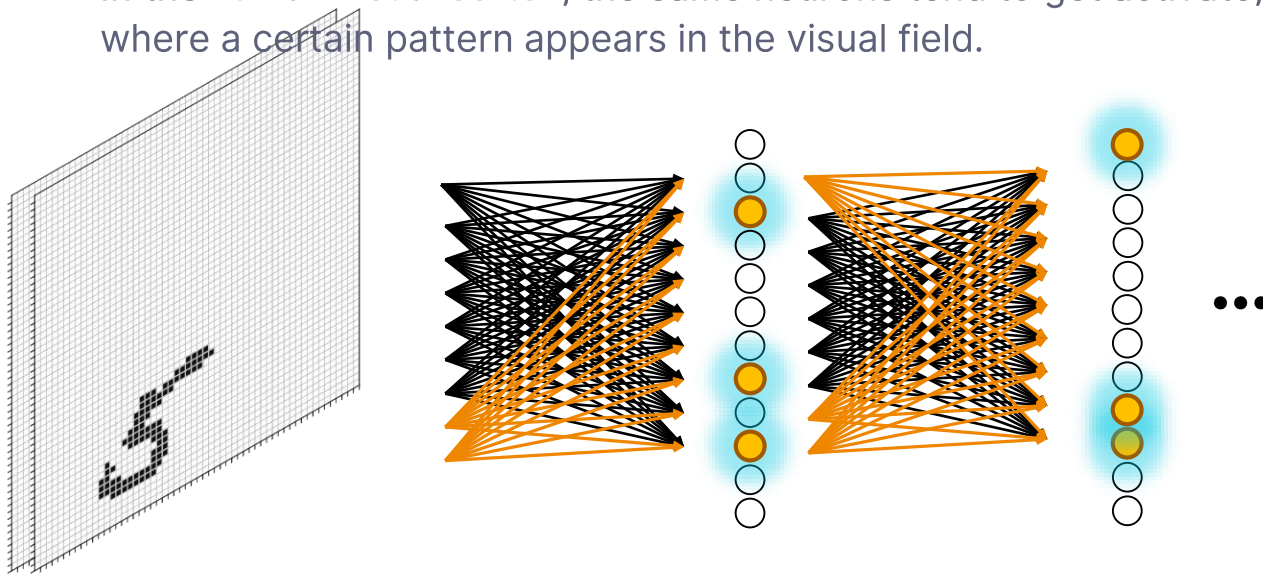
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).
- In the *human visual cortex*, the same neurons tend to get activate, regardless of where a certain pattern appears in the visual field.



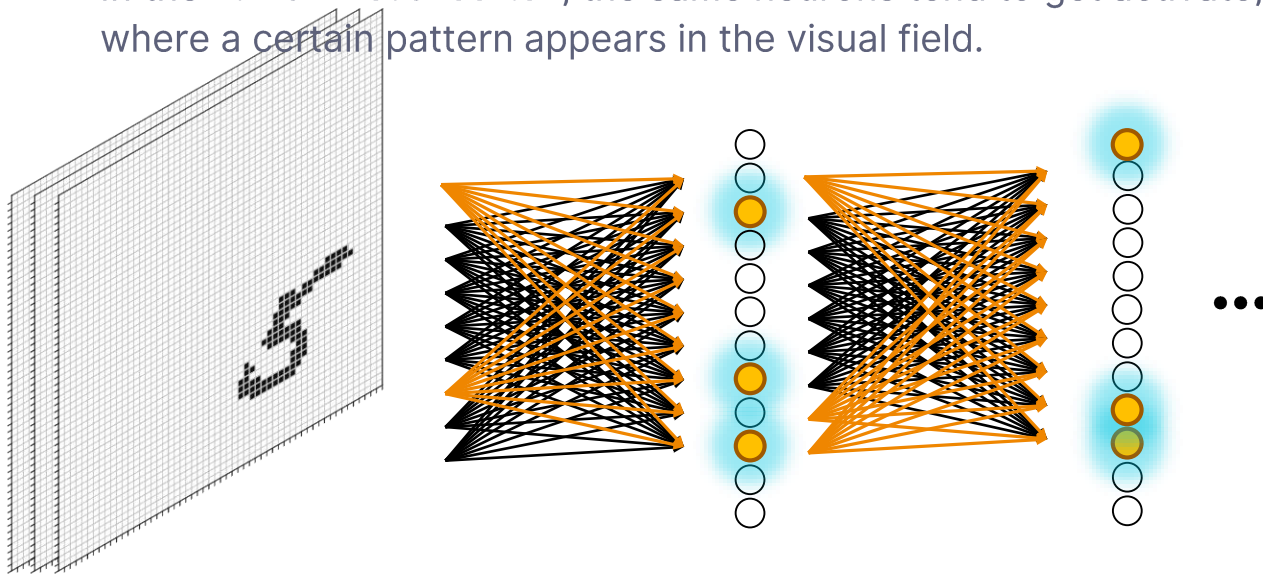
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).
- In the *human visual cortex*, the same neurons tend to get activate, regardless of where a certain pattern appears in the visual field.



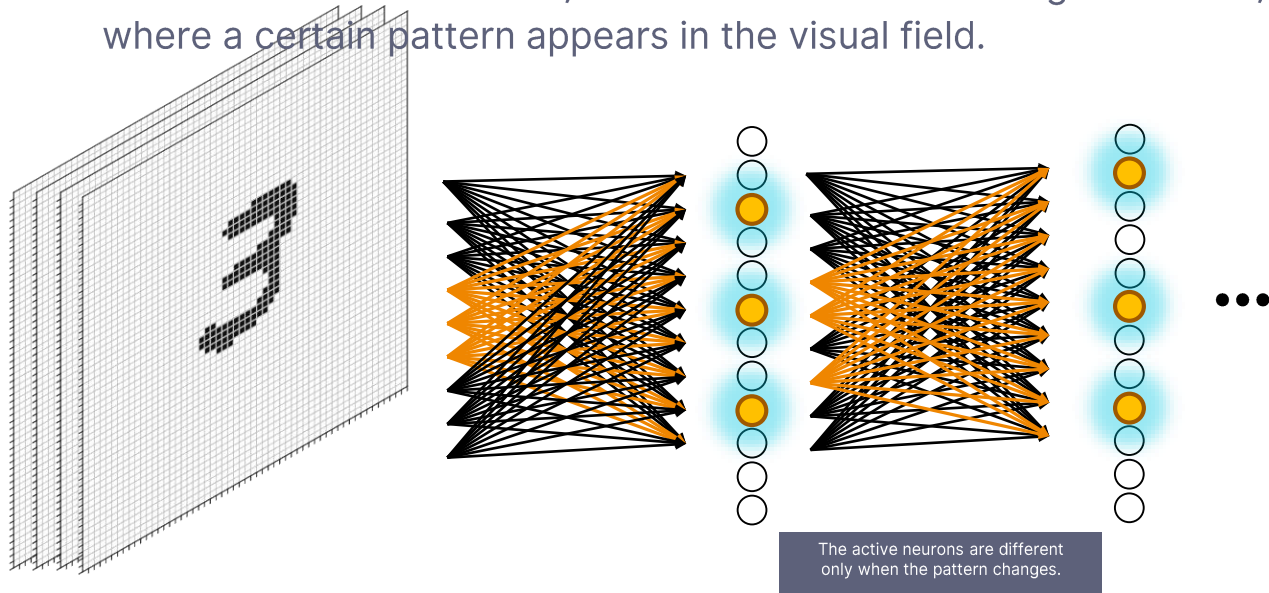
Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).
- In the *human visual cortex*, the same neurons tend to get activate, regardless of where a certain pattern appears in the visual field.



Motivation

- In a *Multilayer Perceptron*, completely different neurons get activated if the same pattern appears in different parts of the input image (even if the MLP correctly identifies the pattern).
- In the *human visual cortex*, the same neurons tend to get activate, regardless of where a certain pattern appears in the visual field.

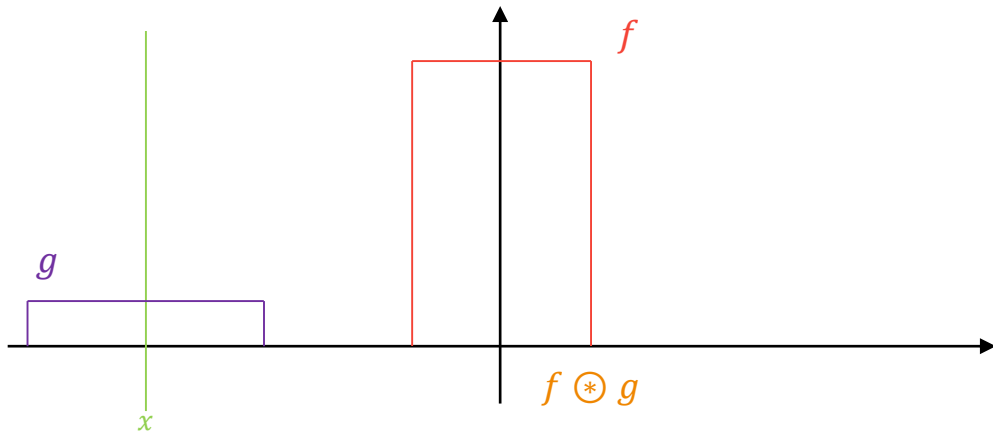


What is a convolution?

What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

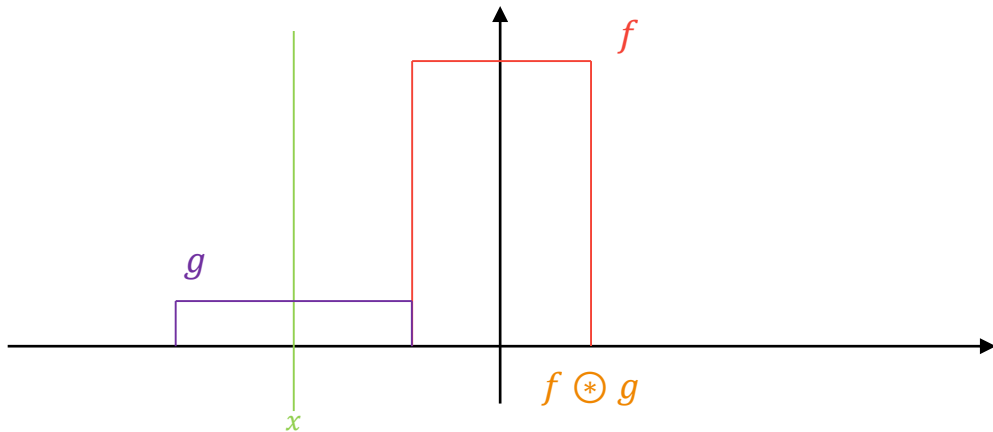
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

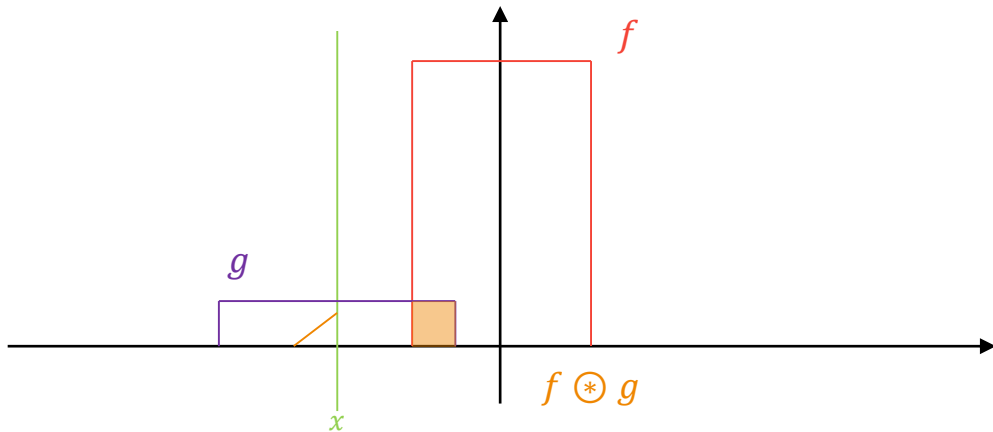
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

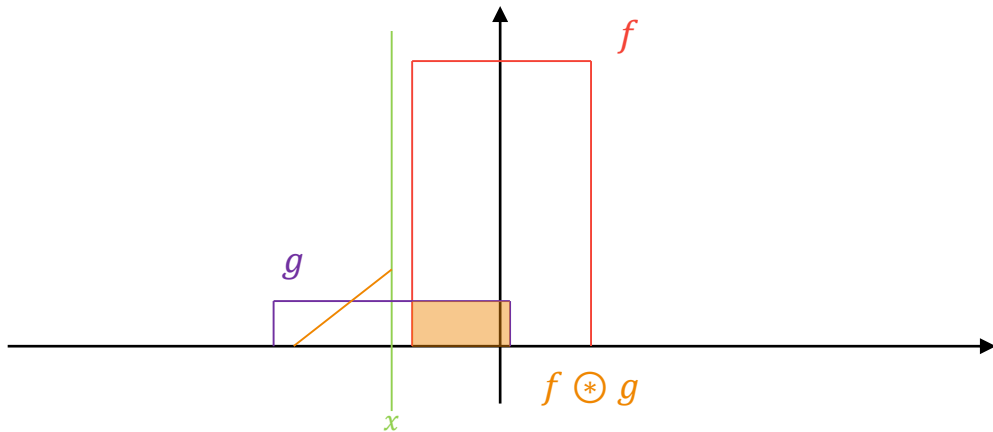
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

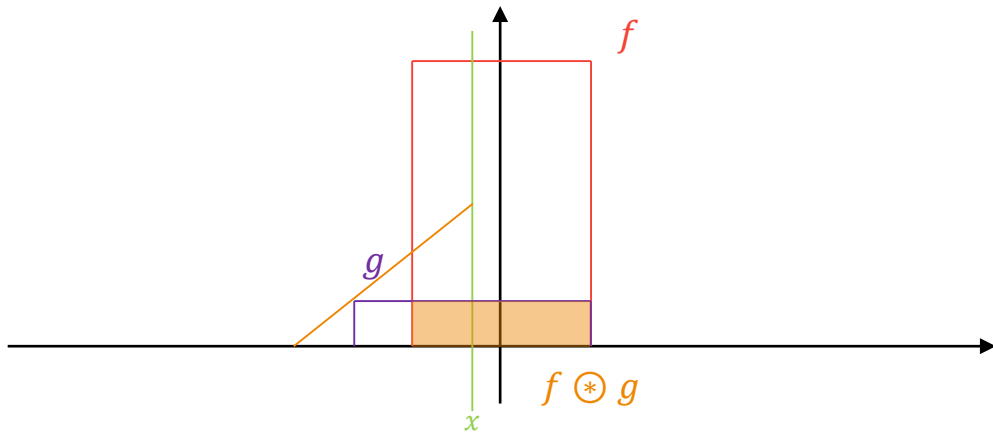
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

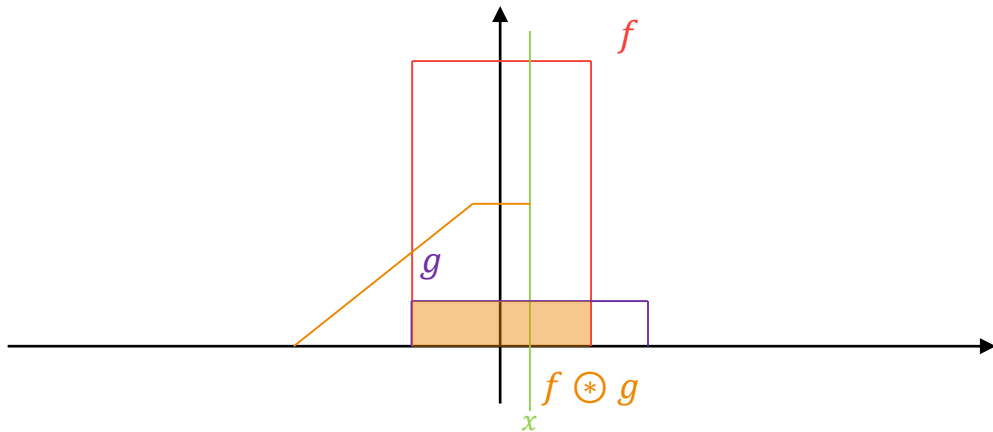
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

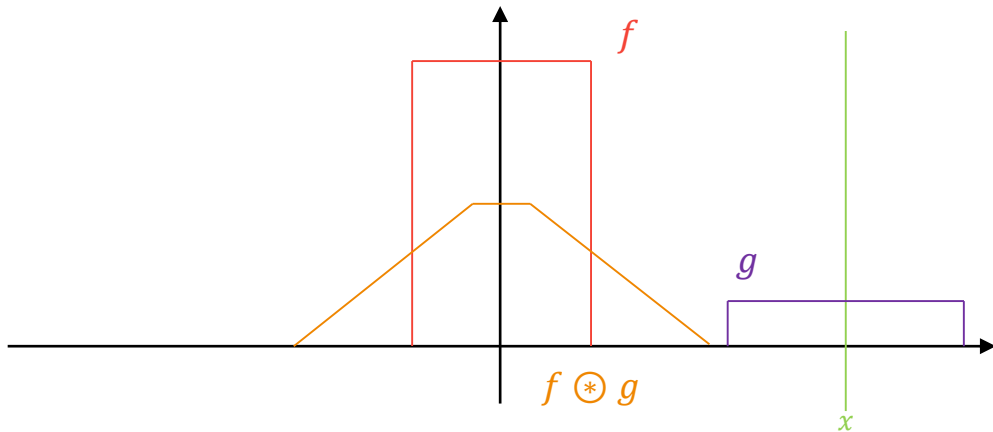
$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



What is a convolution?

- Mathematically speaking, a convolution on two functions f and g is an operation \circledast which produces a third function $f \circledast g$ that expresses how the shape of one is modified by the other.

$$(f \circledast g)(x) = \int_{-\infty}^{\infty} f(t)g(x-t)dt$$



Discrete Convolution

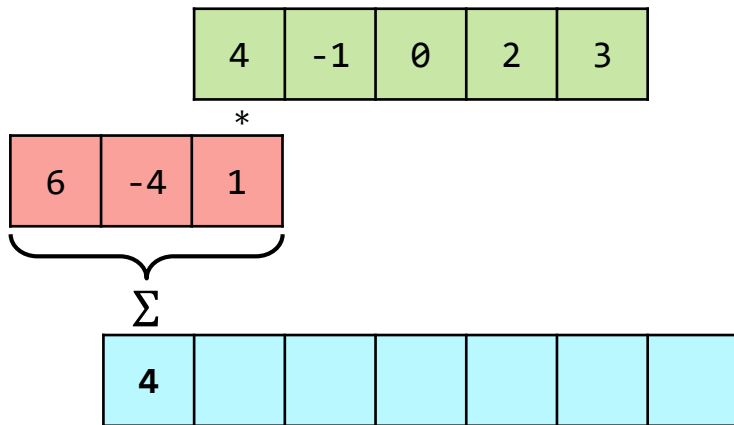
4	-1	0	2	3
---	----	---	---	---

6	-4	1
---	----	---

--	--	--	--	--	--	--

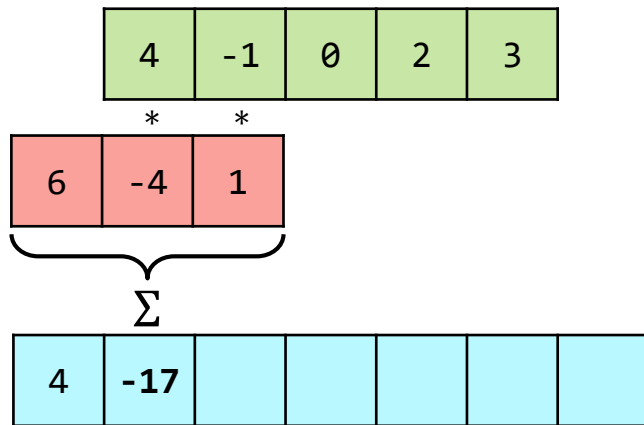
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Discrete Convolution



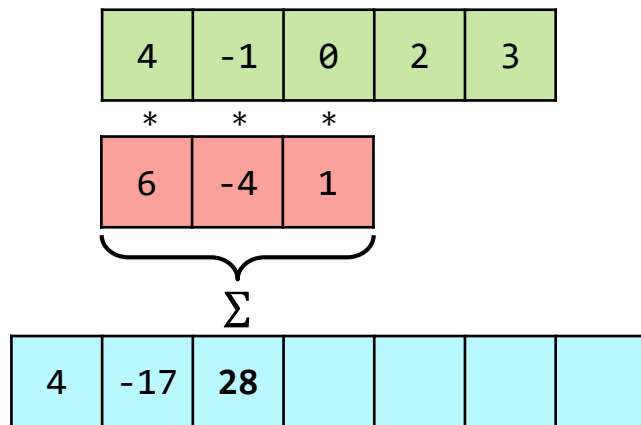
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Discrete Convolution



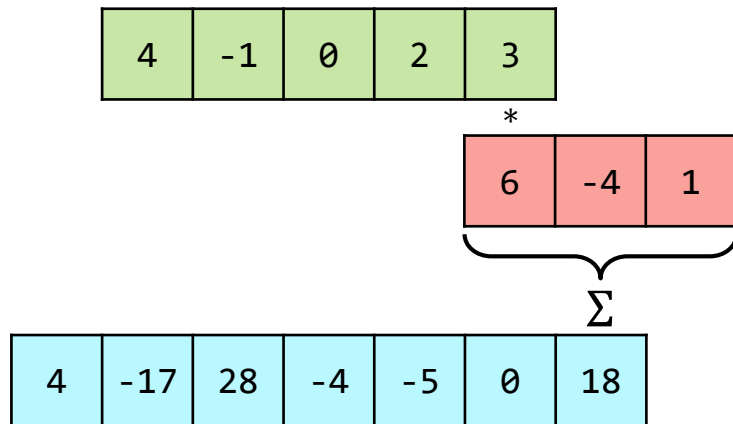
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Discrete Convolution



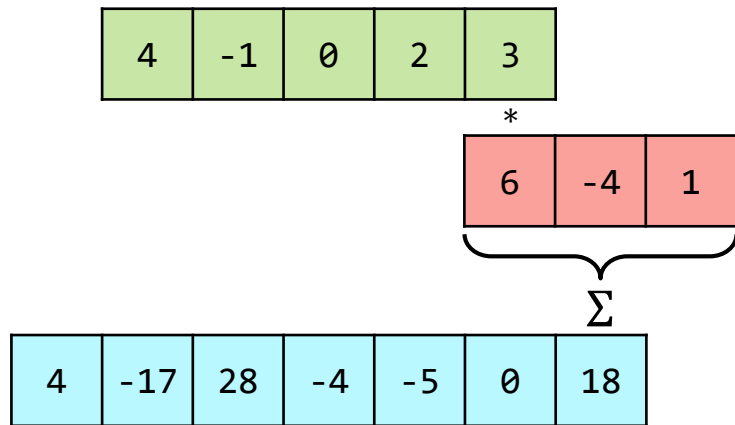
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Discrete Convolution



1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Discrete Convolution



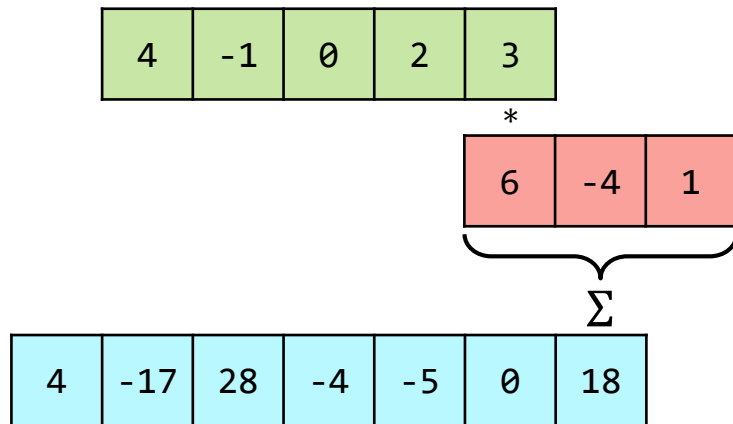
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

6	2	1
0	-2	3
2	-1	-4

4	1	0	2	5
1	-1	6	-2	4
-3	4	1	3	2
0	-3	0	5	-1

2D: $(f \circledast g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Discrete Convolution



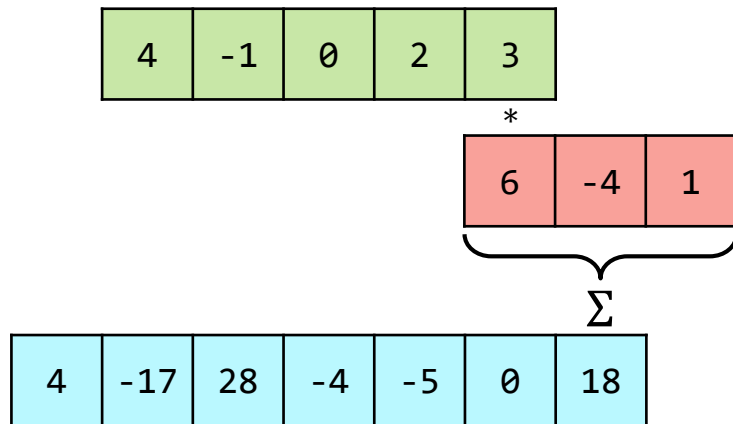
1D: $(f \otimes g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

6	2	1						
0	-2	3						
2	-1	4	-4	1	0	2	5	
				1	-1	6	-2	4
				-3	4	1	3	2
				0	-3	0	5	-1

-16						

2D: $(f \otimes g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Discrete Convolution



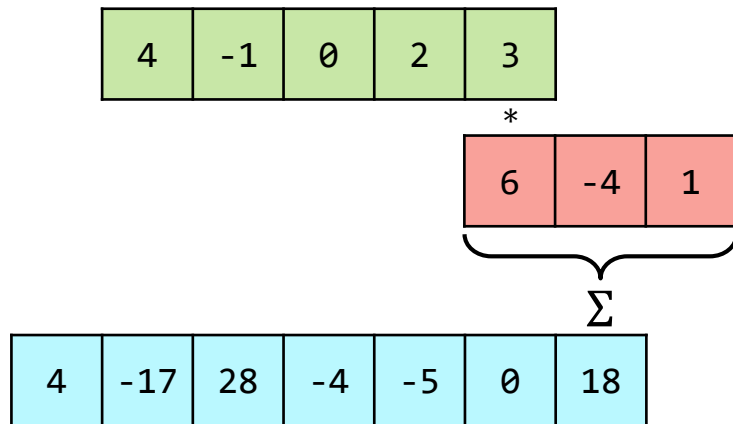
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

6	2	1				
0	-2	3				
2	4	-1	1	0	2	5
	1	-1	6	-2	4	
	-3	4	1	3	2	
	0	-3	0	5	-1	

-16	-8					

2D: $(f \circledast g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Discrete Convolution



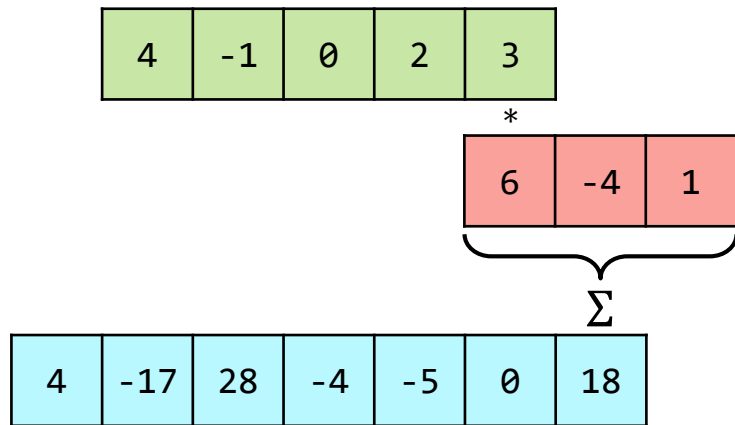
1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

6	2	1			
0	-2	3			
4	2	1	0	2	5
1	-1	6	-2	4	
-3	4	1	3	2	
0	-3	0	5	-1	

-16	-8	7				

2D: $(f \circledast g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Discrete Convolution



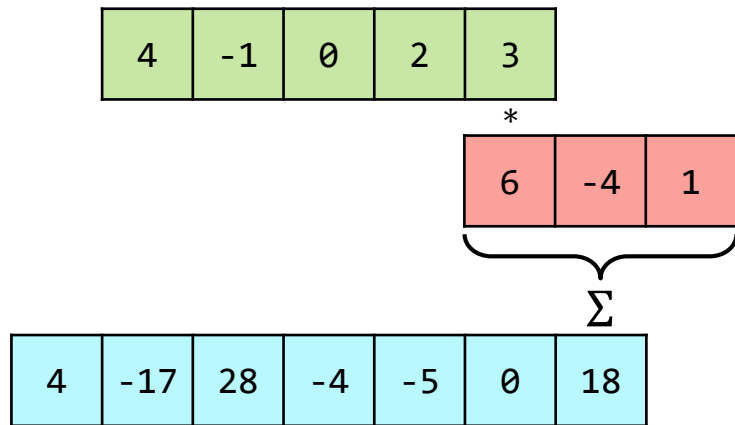
1D: $(f \otimes g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

4	1 ₆	0 ₂	2 ₁	5
1	-1 ₀	6 ₋₂	-2 ₃	4
-3	4 ₂	1 ₋₁	3 ₋₄	2
0	-3	0	5	-1

-16	-8	7	-6	-22	-1	10
8	-2	-23	6	9	-18	8
19	-9	32	-15			

2D: $(f \otimes g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1 - i, x_2 - j)$

Discrete Convolution



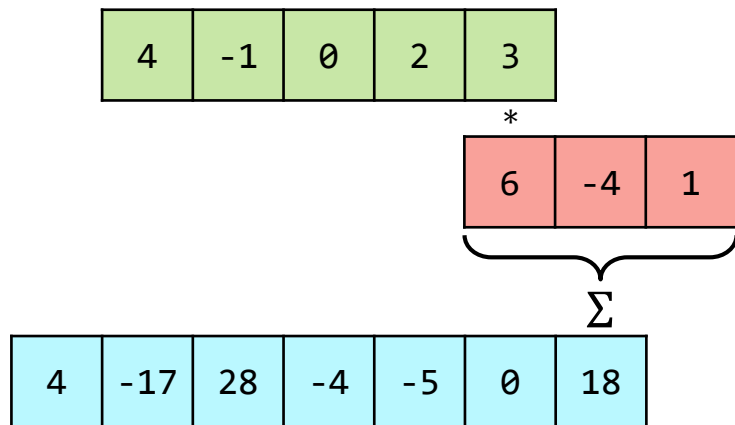
1D: $(f \otimes g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

4	1	0	2	5
1	-1	6	-2	4
-3	4	1	3	2
0	-3	0	5	-1

-16	-8	7	-6	-22	-1	10
8	-2	-23	6	9	-18	8
19	-9	32	-15	16		

2D: $(f \otimes g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1 - i, x_2 - j)$

Discrete Convolution



1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

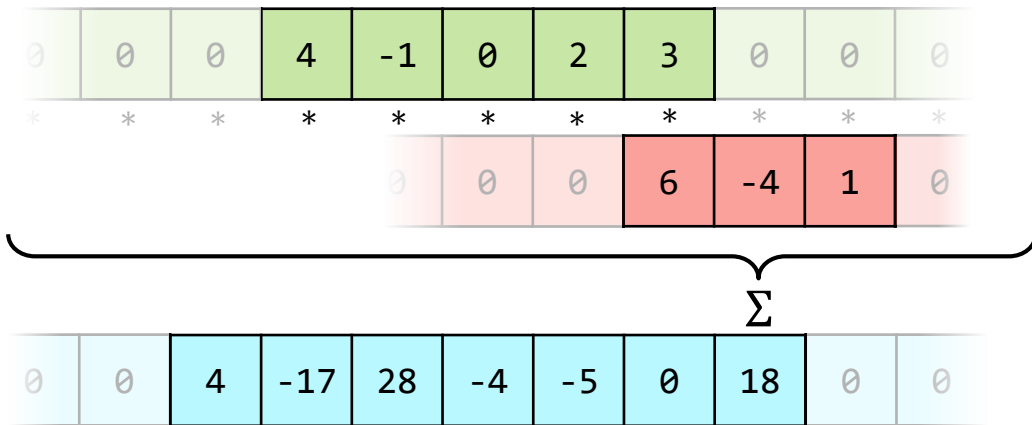
4	1	0	2	5		
1	-1	6	-2	4		
-3	4	1	3	2		
0	-3	0	5	-1	6	2
				0	-2	3
				2	-1	-4

-16	-8	7	-6	-22	-1	10
8	-2	-23	6	9	-18	8
19	-9	32	-15	16	18	34
-8	31	8	-15	35	3	22
-3	-11	-3	44	1	24	12
0	-3	-6	-13	9	28	-6

2D: $(f \circledast g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Discrete Convolution

The sums are *infinite*, but most of the elements are 0 and we are only interested in the *non-zero part of the result*.



1D: $(f \circledast g)(x) = \sum_{i=-\infty}^{\infty} f(i)g(x-i)$

Diagram illustrating 2D discrete convolution. The input is a 10x10 grid of values, and the output is a 7x7 grid of values. The kernel is a 7x7 grid of values. The resulting output grid is shown below.

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0

2D: $(f \circledast g)(x_1, x_2) = \sum_{i=-\infty}^{\infty} \sum_{j=-\infty}^{\infty} f(i, j)g(x_1-i, x_2-j)$

Convolutions in Image Processing

- Even though convolution is *commutative* ($f \circledast g = g \circledast f$), in practice we usually have a larger “input” matrix (typically an *image*) and a smaller matrix, called a **kernel**, which we “convolve” over the image.
- In the field of image processing, the **kernel** is also known as a **filter** or **mask** and it is used for *blurring*, *sharpening*, *edge detection* and other transformations of the input image.



$$\circledast \begin{bmatrix} 0 & -2 & 0 \\ -2 & 9 & -2 \\ 0 & -2 & 0 \end{bmatrix}$$



Sharpen

$$\circledast \frac{1}{256} \begin{bmatrix} 1 & 4 & 6 & 4 & 1 \\ 4 & 16 & 24 & 16 & 4 \\ 6 & 24 & 36 & 24 & 6 \\ 4 & 16 & 24 & 16 & 4 \\ 1 & 4 & 6 & 4 & 1 \end{bmatrix}$$



Blur

$$\circledast \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$$



Edge Detection

Convolution Parameters

Padding

Kernel matrix can be non-square, but it rarely happens in practice.

- Given an image $X \in \mathbb{R}^{n \times m}$ and a kernel $K \in \mathbb{R}^{k \times k}$, the mathematical definition of convolution will produce a result which is larger in size than X (it will be $(n + k - 1) \times (m + k - 1)$)

$$\underset{n \times m}{X} \circledast \underset{k \times k}{K}$$

=

$$\underset{(n+k-1) \times (m+k-1)}{R}$$

4	1	0	2	5		
1	-1	6	-2	4		
-3	4	1	3	2		
0	-3	0	5	-1	6	2
				0	-2	3
				2	-1	-4

-16	-8	7	-6	-22	-1	10
8	-2	-23	6	9	-18	8
19	-9	32	-15	16	18	34
-8	31	8	-15	35	3	22
-3	-11	-3	44	1	24	12
0	-3	-6	-13	9	28	-6

Padding

Kernel matrix can be non-square, but it rarely happens in practice.

- Given an image $X \in \mathbb{R}^{n \times m}$ and a kernel $K \in \mathbb{R}^{k \times k}$, the mathematical definition of convolution will produce a result which is larger in size than X (it will be $(n + k - 1) \times (m + k - 1)$)
- The *convolution* used in image processing is only applied to the area in which the two matrices **overlap completely**.
 - This will result in a matrix $R \in \mathbb{R}^{(n-k+1) \times (m-k+1)}$.

Padding

Kernel matrix can be non-square, but it rarely happens in practice.

- Given an image $X \in \mathbb{R}^{n \times m}$ and a kernel $K \in \mathbb{R}^{k \times k}$, the mathematical definition of convolution will produce a result which is larger in size than X (it will be $(n + k - 1) \times (m + k - 1)$)
- The *convolution* used in image processing is only applied to the area in which the two matrices **overlap completely**.
 - This will result in a matrix $R \in \mathbb{R}^{(n-k+1) \times (m-k+1)}$.

$$\begin{array}{c}
 \textcolor{green}{X} \circledast \textcolor{red}{K} \\
 n \times m \quad k \times k
 \end{array}
 =
 \begin{array}{c}
 \textcolor{blue}{R} \\
 (n - k + 1) \times (m - k + 1)
 \end{array}$$

4	1	0	2	5
1	-1	6 ₆	-2 ₂	4 ₁
-3	4	1 ₀	3 ₋₂	2 ₃
0	-3	0 ₂	5 ₋₁	-1 ₋₄

32	-15	16
8	-15	35

"Valid" Padding

- Basically, "no" padding.
- Sometimes used in practice.

Padding

Kernel matrix can be non-square, but it rarely happens in practice.

- Given an image $X \in \mathbb{R}^{n \times m}$ and a kernel $K \in \mathbb{R}^{k \times k}$, the mathematical definition of convolution will produce a result which is larger in size than X (it will be $(n + k - 1) \times (m + k - 1)$)
- The *convolution* used in image processing is only applied to the area in which the two matrices **overlap completely**.
 - This will result in a matrix $R \in \mathbb{R}^{(n-k+1) \times (m-k+1)}$.
- In practice, we sometimes want the output size to be the same as the input size ($R \in \mathbb{R}^{n \times m}$).
 - This means we need to **pad** the original image with a $k-1/2$ border of zeros.

$$\begin{matrix}
 X & \circledast & K \\
 n \times m & & k \times k
 \end{matrix}
 =
 \begin{matrix}
 R \\
 (n-k+1) \times (m-k+1)
 \end{matrix}$$

4	1	0	2	5
1	-1	6	-2	4
-3	4	1	3	2
0	-3	0	5	-1

32	-15	16
8	-15	35

"Valid" Padding

- Basically, "no" padding.
- Sometimes used in practice.

0	0	0	0	0	0	0
0	4	1	0	2	5	0
0	1	-1	6	-2	4	0
0	-3	4	1	3	2	0
0	0	-3	0	5	-1	0
0	0	0	0	0	-1	0

R
 $n \times m$

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24

"Same" Padding

- Pad such that the output is the "same" size as the input
- Very common in practice.

Padding

Kernel matrix can be non-square, but it rarely happens in practice.

- Given an image $X \in \mathbb{R}^{n \times m}$ and a kernel $K \in \mathbb{R}^{k \times k}$, the mathematical definition of convolution will produce a result which is larger in size than X (it will be $(n + k - 1) \times (m + k - 1)$)
- The *convolution* used in image processing is only applied to the area in which the two matrices **overlap completely**.
 - This will result in a matrix $R \in \mathbb{R}^{(n-k+1) \times (m-k+1)}$.
- In practice, we sometimes want the output size to be the same as the input size ($R \in \mathbb{R}^{n \times m}$).
 - This means we need to **pad** the original image with a $k^{-1}/2$ border of zeros.

$$\begin{matrix}
 \textcolor{green}{X} & \otimes & \textcolor{red}{K} & = & \textcolor{blue}{R} \\
 n \times m & & k \times k & & (n-k+1) \times (m-k+1)
 \end{matrix}$$

4	1	0	2	5
1	-1	6	-2	4
-3	4	1	3	2
0	-3	0	5	-1

32	-15	16
8	-15	35

"Valid" Padding

- Basically, "no" padding.
- Sometimes used in practice.

Kernel size k is usually odd. Otherwise we need unsymmetrical padding.

0	0	0	0	0	0	0
0	4	1	0	2	5	0
0	1	-1	6	-2	4	0
0	-3	4	1	3	2	0
0	0	-3	0	5	-1	0
0	0	0	0	0	-1	0

$\textcolor{blue}{R}$
 $n \times m$

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24

"Same" Padding

- Pad such that the output is the "same" size as the input
- Very common in practice.

Stride

- **Stride** is the amount by which the filter moves when it is convolved over the input image.

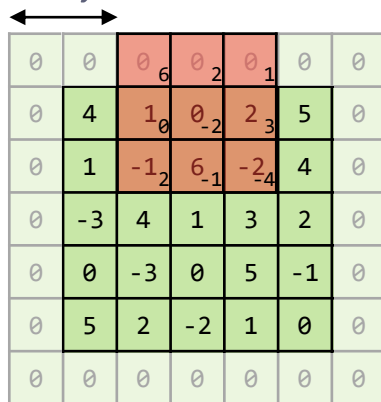
0 ₆	0 ₂	0 ₁	0	0	0	0
0 ₀	4 ₋₂	1 ₃	0	2	5	0
0 ₂	1 ₋₁	-1 ₋₄	6	-2	4	0
0	-3	4	1	3	2	0
0	0	-3	0	5	-1	0
0	5	2	-2	1	0	0
0	0	0	0	0	0	0

Stride $s = 2$

-2		

Stride

- **Stride** is the amount by which the filter moves when it is convolved over the input image.



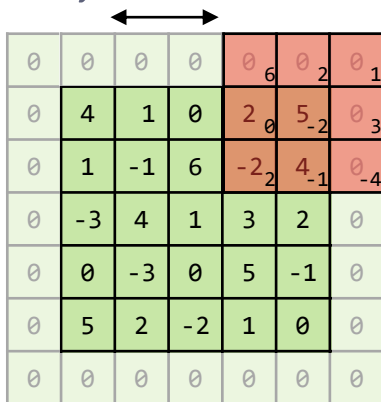
0	0	0 ₆	0 ₂	0 ₁	0	0
0	4	1 ₀	0 ₋₂	2 ₃	5	0
0	1	-1 ₂	6 ₋₁	-2 ₋₄	4	0
0	-3	4	1	3	2	0
0	0	-3	0	5	-1	0
0	5	2	-2	1	0	0
0	0	0	0	0	0	0

Stride $s = 2$

-2	6	

Stride

- **Stride** is the amount by which the filter moves when it is convolved over the input image.



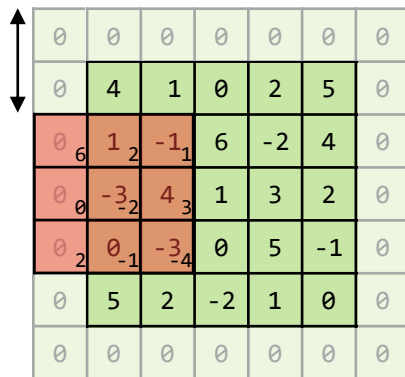
0	0	0	0	0	6	0	2	0	1
0	4	1	0	2	0	5	-2	0	3
0	1	-1	6	-2	2	4	-1	0	-4
0	-3	4	1	3	2	0			
0	0	-3	0	5	-1	0			
0	5	2	-2	1	0	0			
0	0	0	0	0	0	0			

Stride $s = 2$

-2	6	-18

Stride

- **Stride** is the amount by which the filter moves when it is convolved over the input image.



0	0	0	0	0	0	0
0	4	1	0	2	5	0
0	6	1	-1	6	-2	4
0	0	-3	4	1	3	2
0	2	0	-1	-3	0	5
0	5	2	-2	1	0	0
0	0	0	0	0	0	0

Stride $s = 2$

-2	6	-18
31		

Stride

- **Stride** is the amount by which the filter moves when it is convolved over the input image.

0	0	0	0	0	0	0
0	4	1	0	2	5	0
0	1	-1	6	-2	4	0
0	-3	4	1	3	2	0
0	0	-3	0	5	-1	0
0	5	2	-2	1	0	0
0	0	0	0	0	-2	0

-2	6	-18
31	-15	3
-7	-14	28

Stride $s = 2$

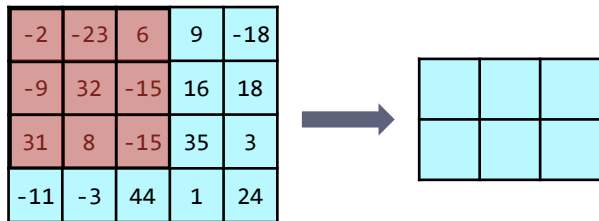
- If we have an image of size $n \times m$, a kernel of size $k \times k$ and we have padding p and stride s , the output will have size:

$$o = \left(\left\lfloor \frac{n - k + 2p}{s} \right\rfloor + 1 \right) \times \left(\left\lfloor \frac{m - k + 2p}{s} \right\rfloor + 1 \right)$$

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1




Average Pooling

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



1.4		


Average Pooling

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



1.4	5.8	


Average Pooling

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



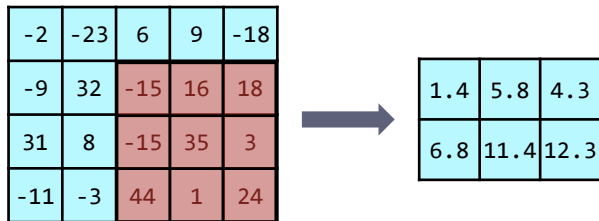
1.4	5.8	4.3
6.8	11.4	12.3

Average Pooling

Pooling

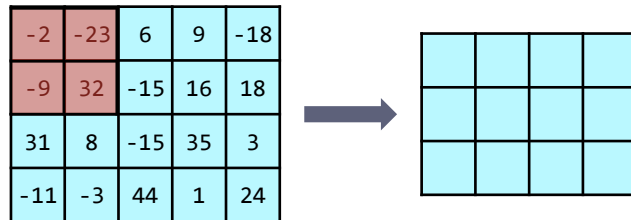
- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1



Average Pooling

size = 2x2, stride = 1



Max Pooling


Most common in practice

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24




1.4	5.8	4.3
6.8	11.4	12.3

Average Pooling

size = 2x2, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



32			

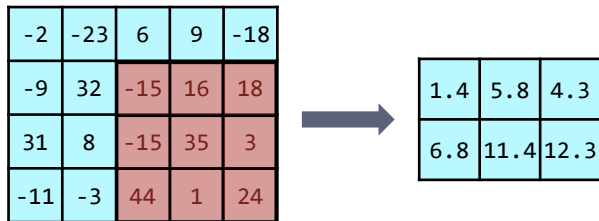
Max Pooling

Most common in practice

Pooling

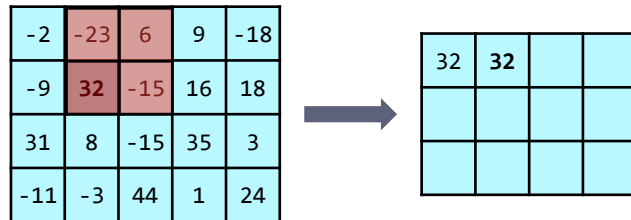
- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1



Average Pooling

size = 2x2, stride = 1



Max Pooling


Most common in practice

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24




1.4	5.8	4.3
6.8	11.4	12.3

Average Pooling

size = 2x2, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



32	32	16	

Max Pooling


Most common in practice

Pooling

- **Pooling** is an operation which *down-samples* an input image and it is typically applied after a filter convolution.
 - The reasoning is that a filter usually tries to *capture a specific aspect* of the input image (e.g. edge detection) and it does not take as many pixels (features) to “describe” this particular aspect as for the original image.
 - Similarly to a convolution operation, pooling has *size*, *padding* and *stride*.

size = 3x3, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24




1.4	5.8	4.3
6.8	11.4	12.3

Average Pooling

size = 2x2, stride = 1

-2	-23	6	9	-18
-9	32	-15	16	18
31	8	-15	35	3
-11	-3	44	1	24



32	32	16	18
32	32	35	35
31	44	44	35

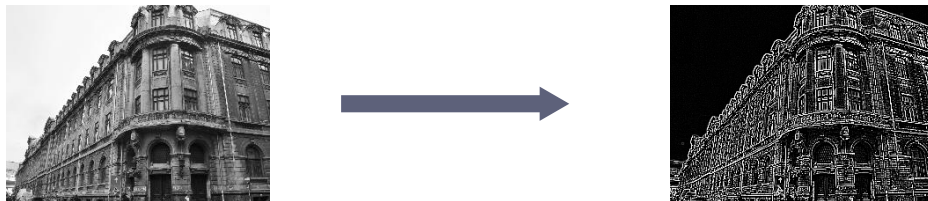
Max Pooling

Most common in practice

Neural Networks with Convolution

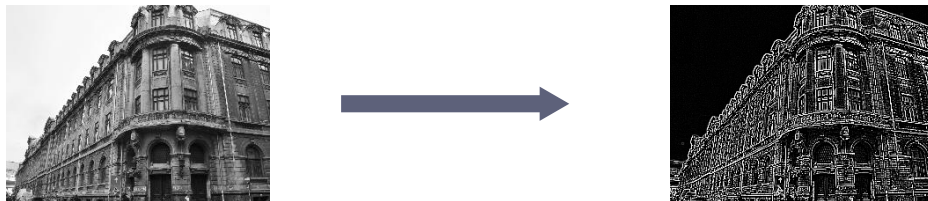
Convolutions as Features

- Convoluting a kernel over an input image transforms it into a *new image* with *new characteristics* which were not present in the original image (a different “**representation**” of the input).



Convolutions as Features

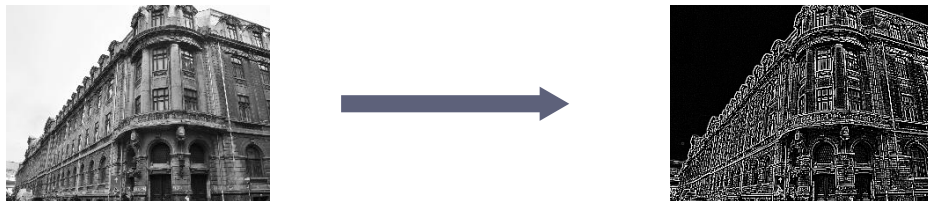
- Convoluting a kernel over an input image transforms it into a *new image with new characteristics* which were not present in the original image (a different “**representation**” of the input).



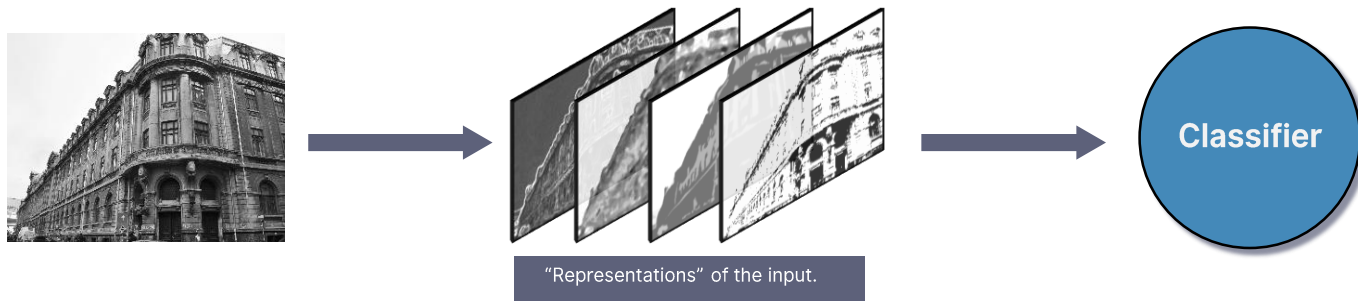
- These *newly obtained features* might contain *better discriminative information* about the labels we are trying to predict.
 - Passing them to a classifier instead of the original pixels could improve performance.

Convolutions as Features

- Convoluting a kernel over an input image transforms it into a *new image with new characteristics* which were not present in the original image (a different “**representation**” of the input).



- These *newly obtained features* might contain *better discriminative information* about the labels we are trying to predict.
 - Passing them to a classifier instead of the original pixels could improve performance.
- Better still, we could take a couple of different kernels and feed all their outputs to the classifier.



Convolutional Layer

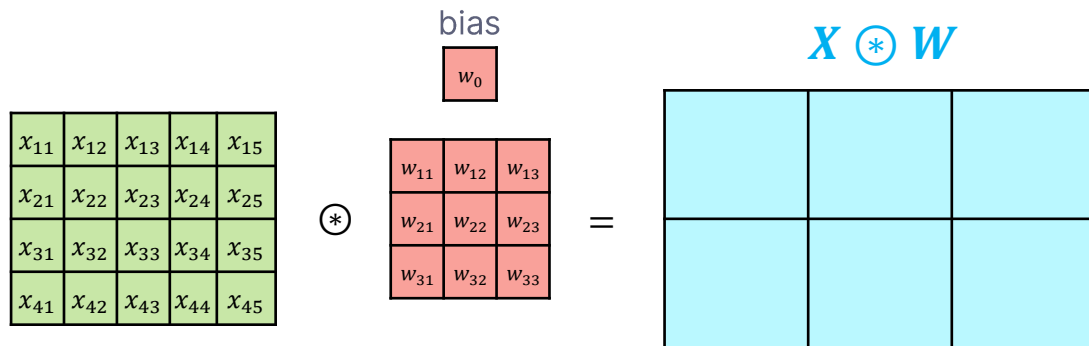
- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.

Convolutional Layer

- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, *"choosing the proper convolutions"* should be the task of the learning algorithm itself.

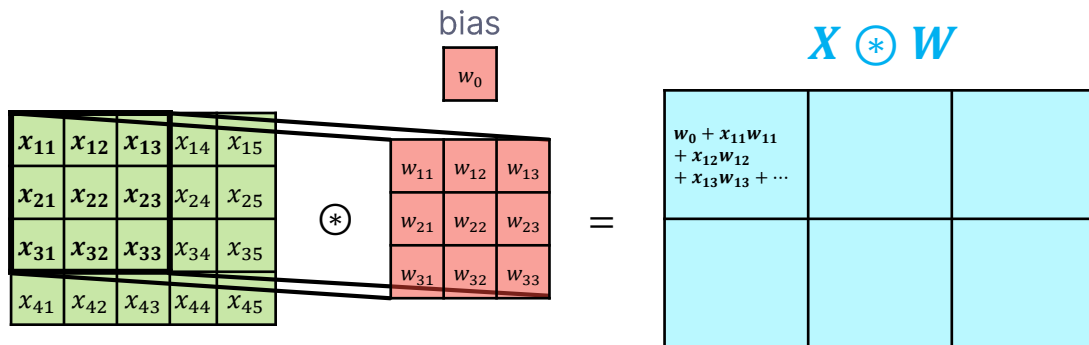
Convolutional Layer

- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, "*choosing the proper convolutions*" should be the task of the learning algorithm itself.
- A **convolutional layer** is a filter (or a set of filters) whose *weights are learned* through *gradient descent* in the *optimization process*.



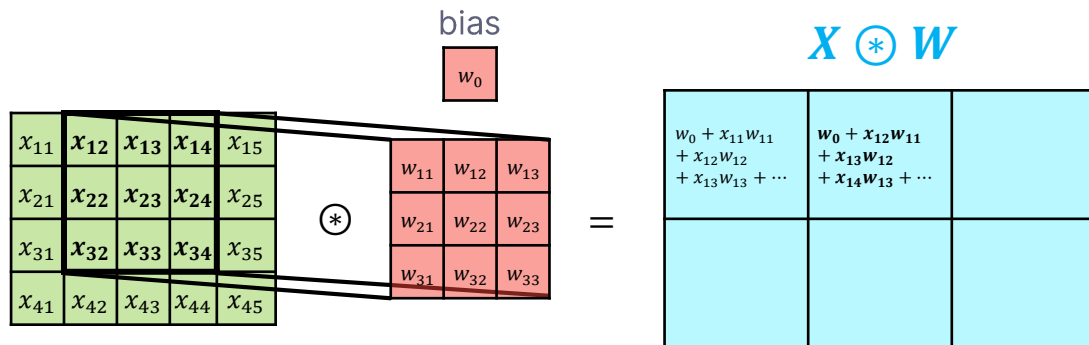
Convolutional Layer

- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, "*choosing the proper convolutions*" should be the task of the learning algorithm itself.
- A **convolutional layer** is a filter (or a set of filters) whose *weights are learned* through *gradient descent* in the *optimization process*.



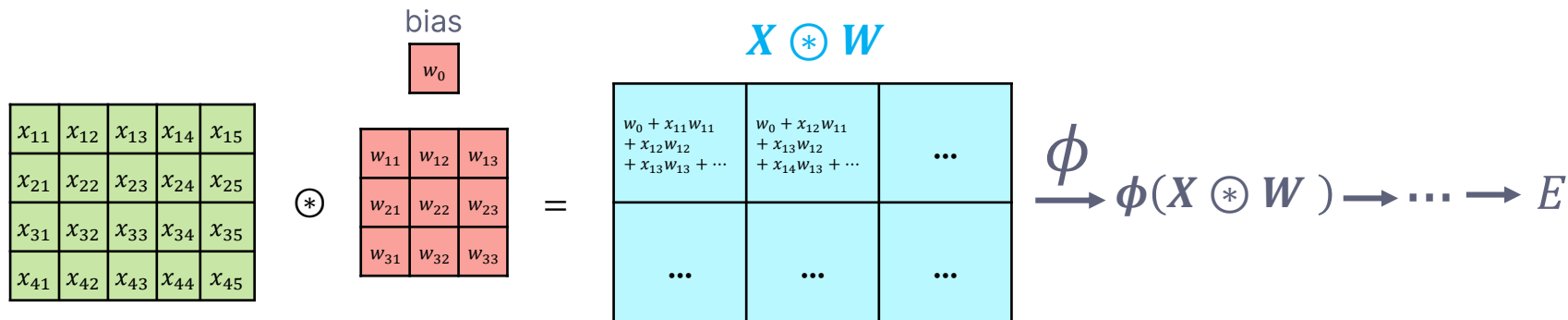
Convolutional Layer

- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, "choosing the proper convolutions" should be the task of the learning algorithm itself.
- A **convolutional layer** is a filter (or a set of filters) whose *weights are learned* through *gradient descent* in the *optimization process*.



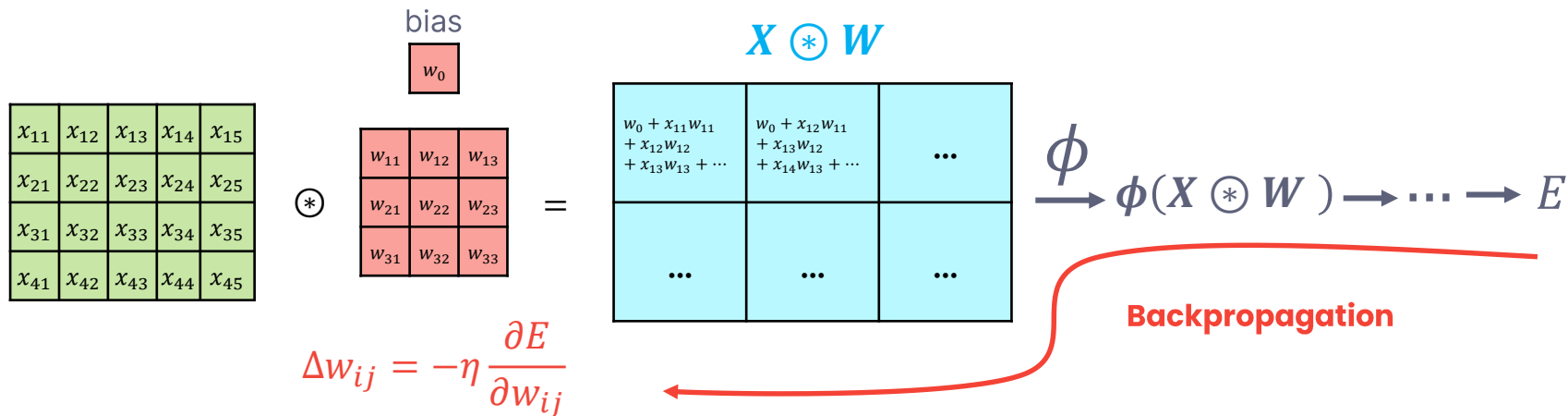
Convolutional Layer

- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, "choosing the proper convolutions" should be the task of the learning algorithm itself.
- A **convolutional layer** is a filter (or a set of filters) whose *weights are learned* through *gradient descent* in the *optimization process*.



Convolutional Layer

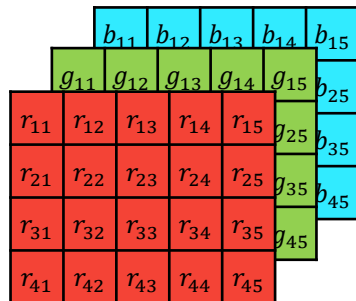
- By choosing proper convolution kernels, we could substantially improve the performance of the predictive algorithm.
 - But in **Deep Learning**, "choosing the proper convolutions" should be the task of the learning algorithm itself.
- A **convolutional layer** is a filter (or a set of filters) whose *weights are learned* through *gradient descent* in the *optimization process*.



Multiple Channels

- A color image has 3 channels (each pixel has 3 different values for *red*, *green* and *blue*).
- A filter on a “multi-channel” image needs weights for each channel.

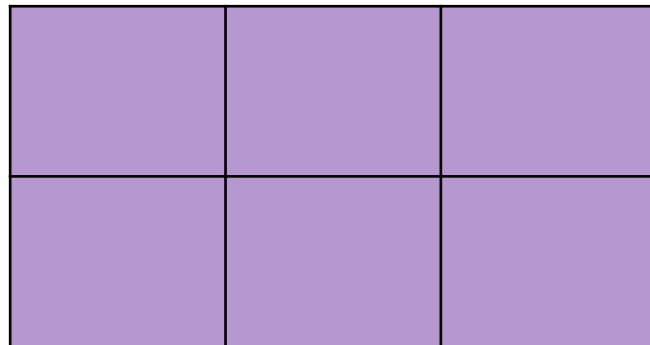
$$x_{ij} = (r_{ij}, g_{ij}, b_{ij})$$



(*)



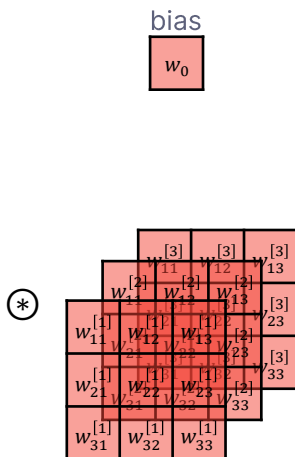
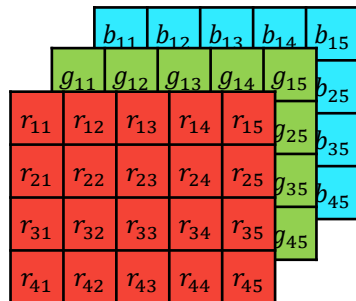
=



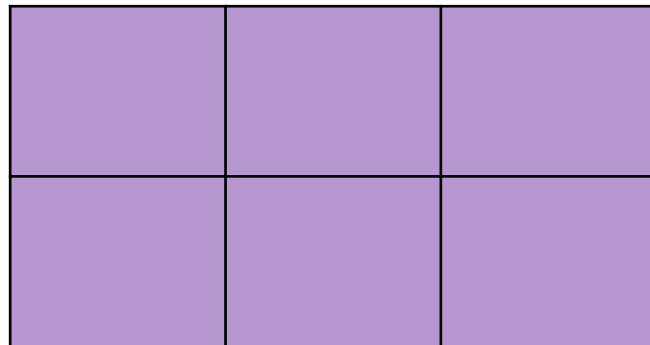
Multiple Channels

- A color image has 3 channels (each pixel has 3 different values for *red*, *green* and *blue*).
- A filter on a “multi-channel” image needs weights for each channel.

$$x_{ij} = (r_{ij}, g_{ij}, b_{ij})$$



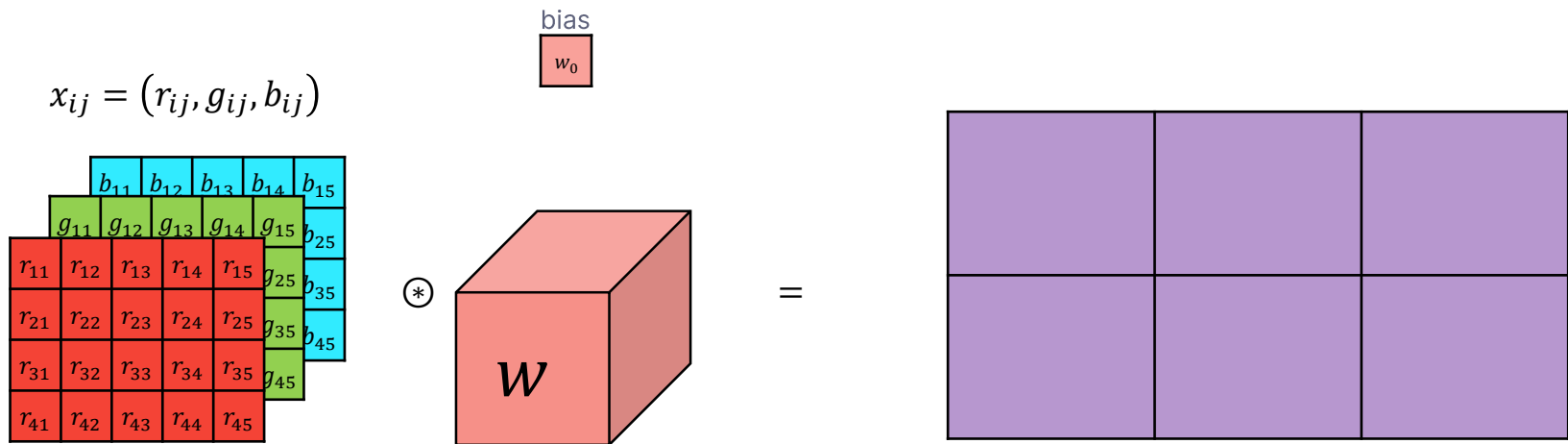
=



Multiple Channels

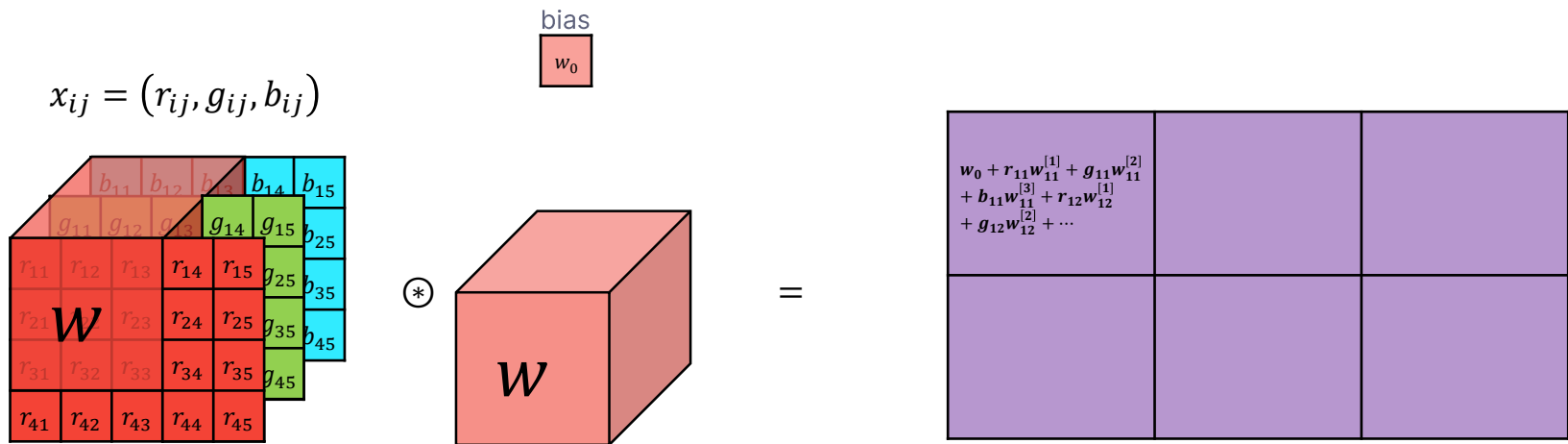
- A color image has 3 channels (each pixel has 3 different values for *red*, *green* and *blue*).
- A filter on a “multi-channel” image needs weights for each channel.

$$x_{ij} = (r_{ij}, g_{ij}, b_{ij})$$



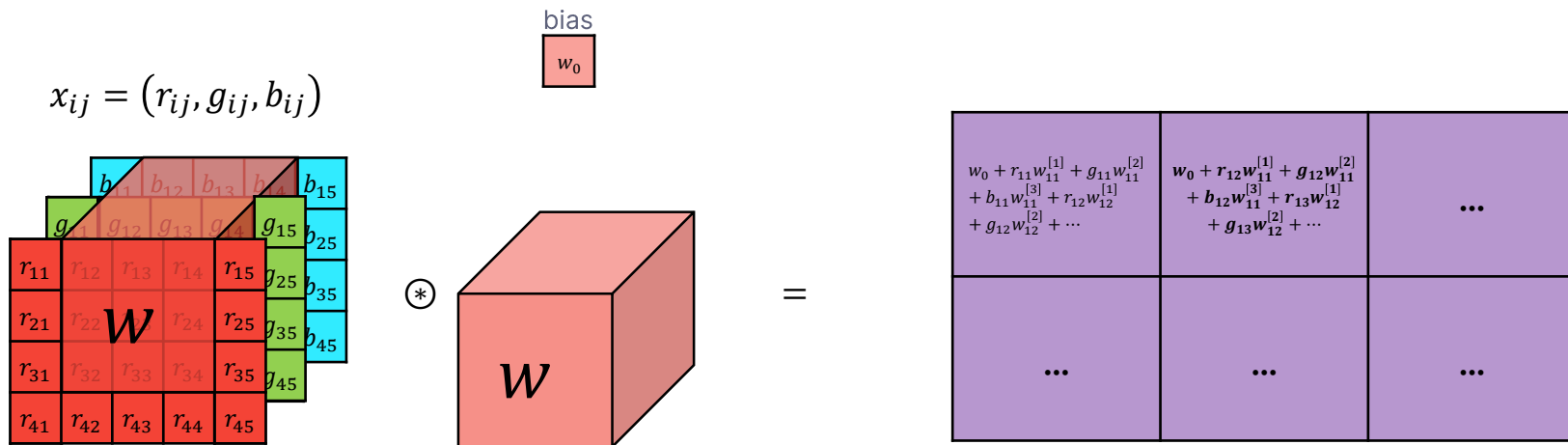
Multiple Channels

- A color image has 3 channels (each pixel has 3 different values for *red*, *green* and *blue*).
- A filter on a “multi-channel” image needs weights for each channel.



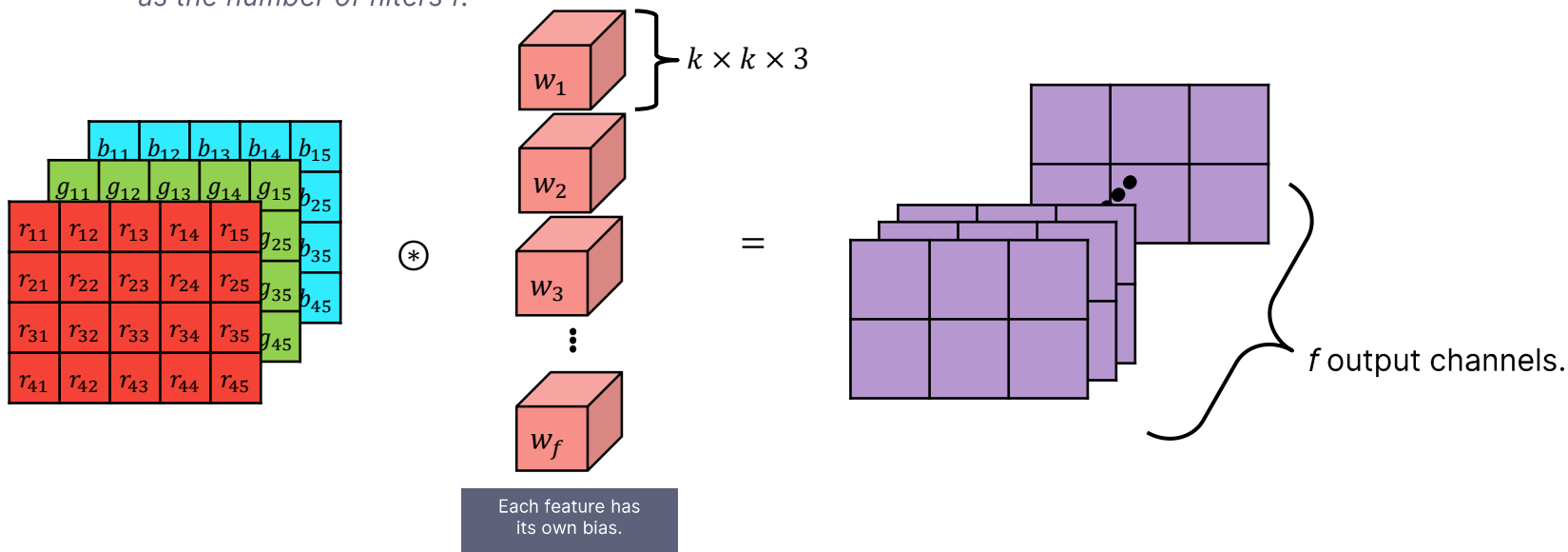
Multiple Channels

- A color image has 3 channels (each pixel has 3 different values for *red*, *green* and *blue*).
- A filter on a “multi-channel” image needs weights for each channel.

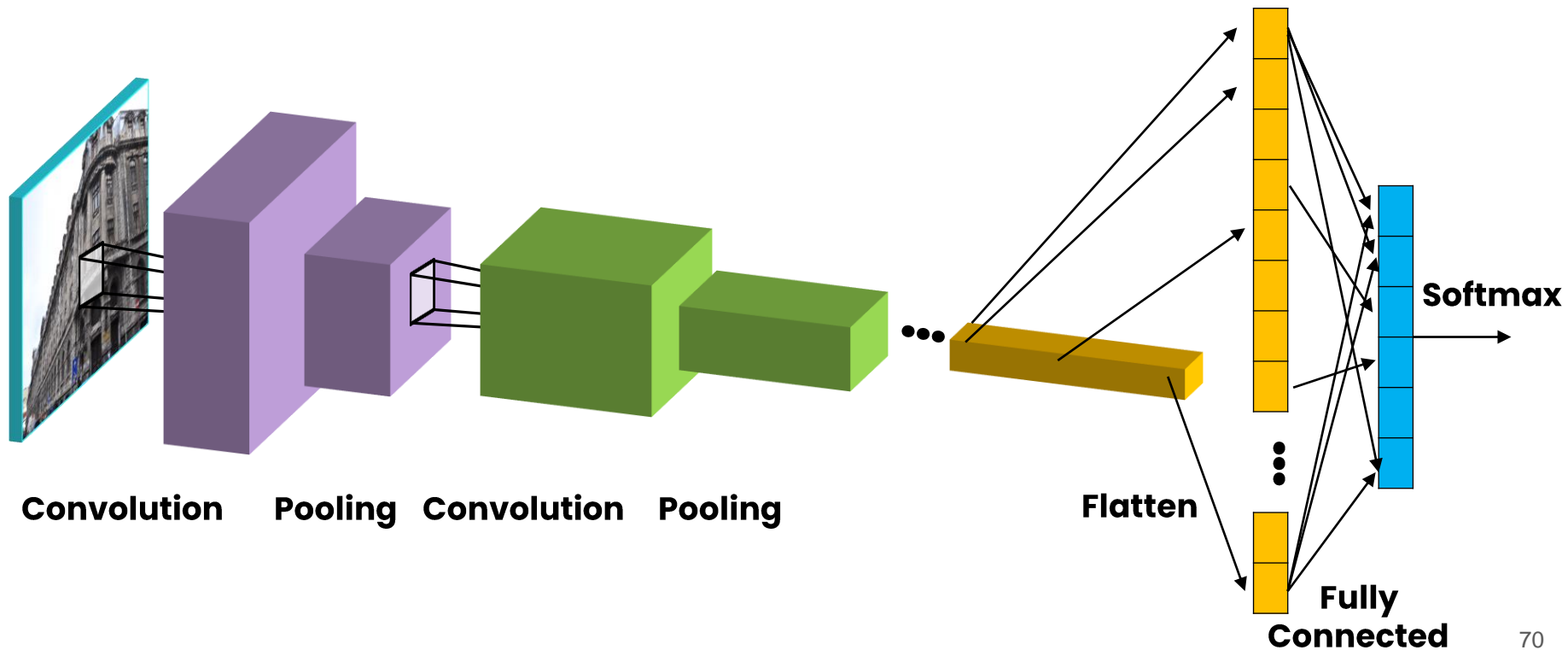


Multiple Channels

- One filter usually captures one feature of the input image, but we want to obtain a *representation with many features*.
 - A **convolutional layer** is usually made up of **multiple filters**.
 - Each filter has the *depth* (number of channels) of the input image and the output has *as many channels as the number of filters* f .



Typical CNN Architecture



Keywords

Convolution

Kernel

Filter

Mask

Padding

Valid Padding

Same Padding

Stride

Pooling

MaxPooling

Convolutional Neural Network (CNN)

Recurrent Neural Networks

Handling **sequences**
with neural networks

Faculty of Mathematics and Computer Science, University of Bucharest
and
Sparktech Software

Academic Year 2018/2019, 1st Semester

Motivation

- The problems we have seen so far are “one-to-one” (i.e. one input mapped to one output).
 - Better said, all inputs and all outputs always have the same fixed size.
 - e.g. In the digit recognition problem, all inputs had $28 \times 28 = 724$ pixels and all outputs had 10 values (one hot encoding of a digit).

Motivation

- The problems we have seen so far are “one-to-one” (i.e. one input mapped to one output).
 - Better said, all inputs and all outputs always have the same fixed size.
 - e.g. In the digit recognition problem, all inputs had $28 \times 28 = 784$ pixels and all outputs had 10 values (one hot encoding of a digit).
- Many real problems involve sequences of data:

Speech Recognition



“What is the weather going to be like tomorrow?”

Image Captioning



“Black and white dog jumps over bar.”

Machine Translation

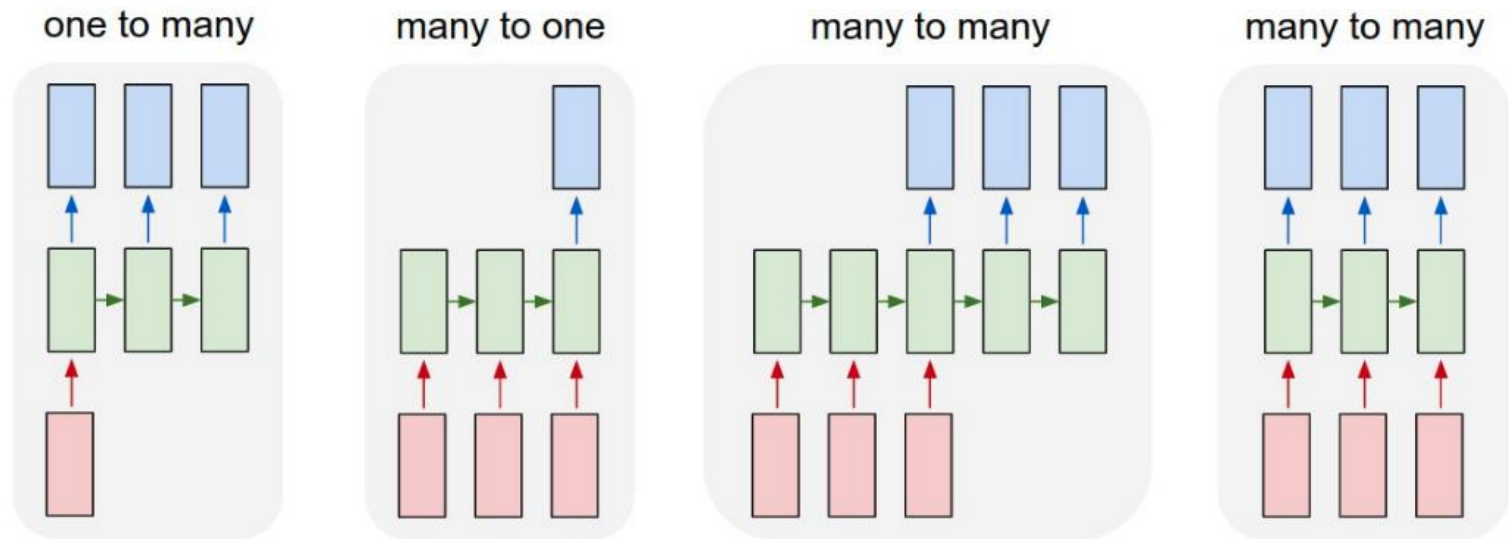
“I am going to the cinema.”



“Ich gehe ins Kino.”

Motivation

- Types of problems which involve sequences:



Andrej Karpathy

<http://karpathy.github.io/2015/05/21/rnn-effectiveness/>

Recurrent Layers

Notation

- One training example is now a **sequence of values**: $\vec{x} \in \mathbb{R}^{m \times n \times T}$
 - $\vec{x}_j^{(i)(t)}$ - component j of training sample i at timestep t

Notation

- One training example is now a **sequence of values**: $\vec{x} \in \mathbb{R}^{m \times n \times T}$
 - $\vec{x}_j^{(i)(t)}$ - component j of training sample i at timestep t

$$\vec{x}^{(t)} \in \mathbb{R}^n$$

input at "timestep" t

$$\hat{y}^{(t)} \in \mathbb{R}^o$$

output at timestep t



Notation

- One training example is now a **sequence of values**: $\vec{x} \in \mathbb{R}^{m \times n \times T}$
 - $\vec{x}_j^{(i)(t)}$ - component j of training sample i at timestep t

$$\vec{x}^{(t)} \in \mathbb{R}^n$$

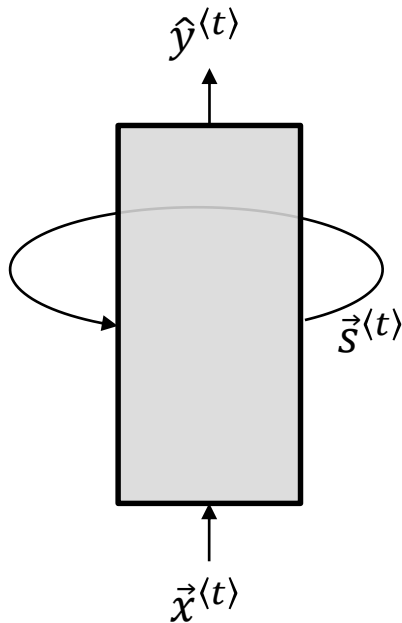
input at "timestep" t

$$\hat{y}^{(t)} \in \mathbb{R}^o$$

output at timestep t

$$\vec{s}^{(t)} \in \mathbb{R}^h$$

hidden state at "timestep" t



Notation

- One training example is now a **sequence of values**: $\vec{x} \in \mathbb{R}^{m \times n \times T}$
 - $\vec{x}_j^{(i)(t)}$ - component j of training sample i at timestep t

$$\vec{x}^{(t)} \in \mathbb{R}^n$$

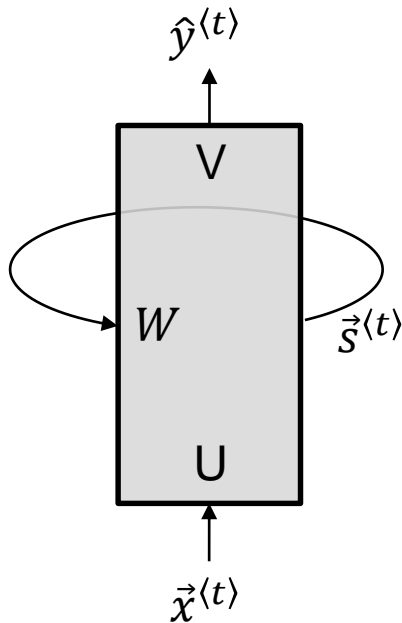
input at "timestep" t

$$\hat{y}^{(t)} \in \mathbb{R}^o$$

output at timestep t

$$\vec{s}^{(t)} \in \mathbb{R}^h$$

hidden state at "timestep" t



Three weight matrices:

- $U \in \mathbb{R}^{h \times n}$
- $W \in \mathbb{R}^{h \times h}$
- $V \in \mathbb{R}^{o \times h}$

Notation

- One training example is now a **sequence of values**: $\vec{x} \in \mathbb{R}^{m \times n \times T}$
 - $\vec{x}_j^{(i)(t)}$ - component j of training sample i at timestep t

$$\vec{x}^{(t)} \in \mathbb{R}^n$$

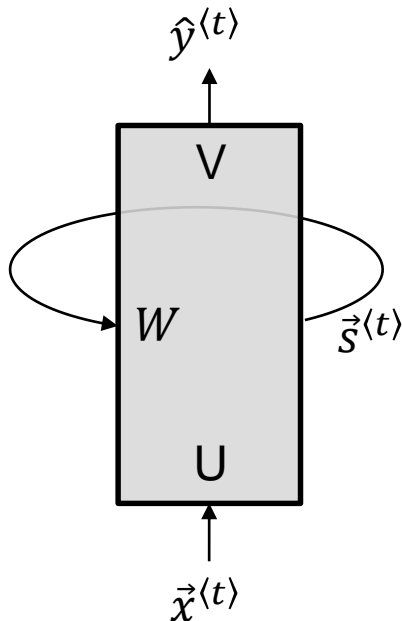
input at "timestep" t

$$\hat{y}^{(t)} \in \mathbb{R}^o$$

output at timestep t

$$\vec{s}^{(t)} \in \mathbb{R}^h$$

hidden state at "timestep" t



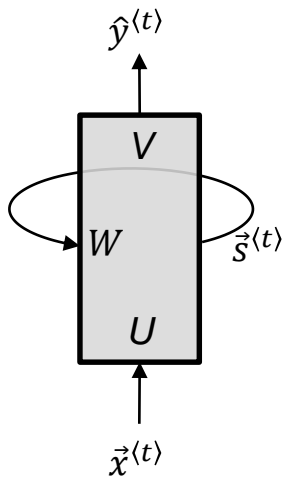
Three weight matrices:

- $U \in \mathbb{R}^{h \times n}$
- $W \in \mathbb{R}^{h \times h}$
- $V \in \mathbb{R}^{o \times h}$

$$\vec{s}^{(t)} = \phi_1(W\vec{s}^{(t-1)} + U\vec{x}^{(t)})$$
$$\hat{y}^{(t)} = \phi_2(V\vec{s}^{(t)})$$

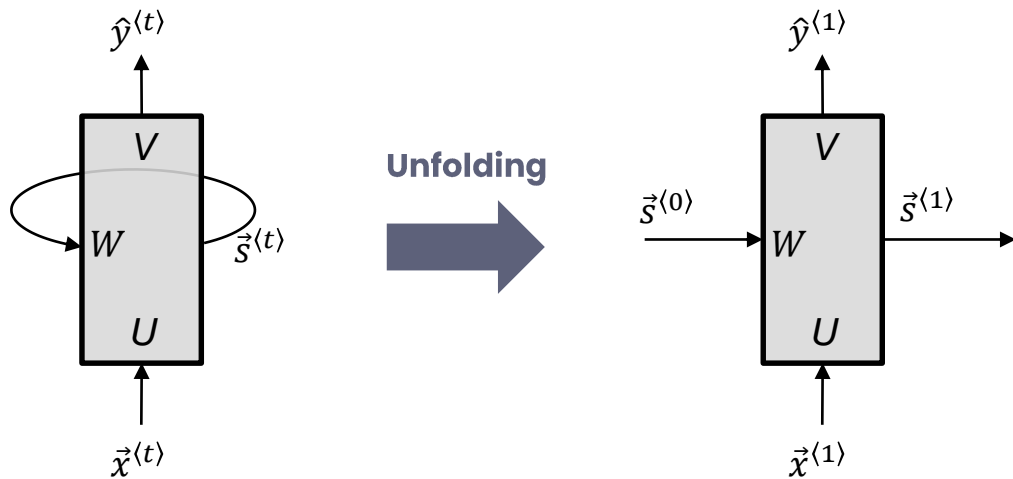
ϕ_1 and ϕ_2 are different activation functions.
Usually ϕ_1 is a tanh and ϕ_2 is a sigmoid (or softmax)

Forward propagation



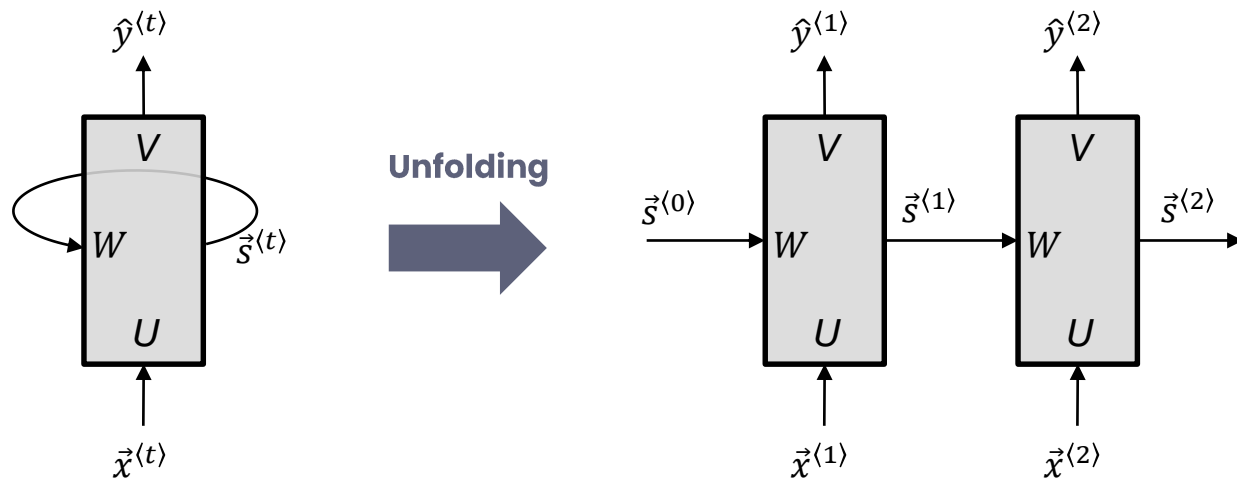
$$\begin{aligned}\vec{s}^{(0)} &= \vec{0} \\ s^{(t)} &= \phi_1(W\vec{s}^{(t-1)} + U\vec{x}^{(t)}) \\ \hat{y}^{(t)} &= \phi_2(V\vec{s}^{(t)})\end{aligned}$$

Forward propagation



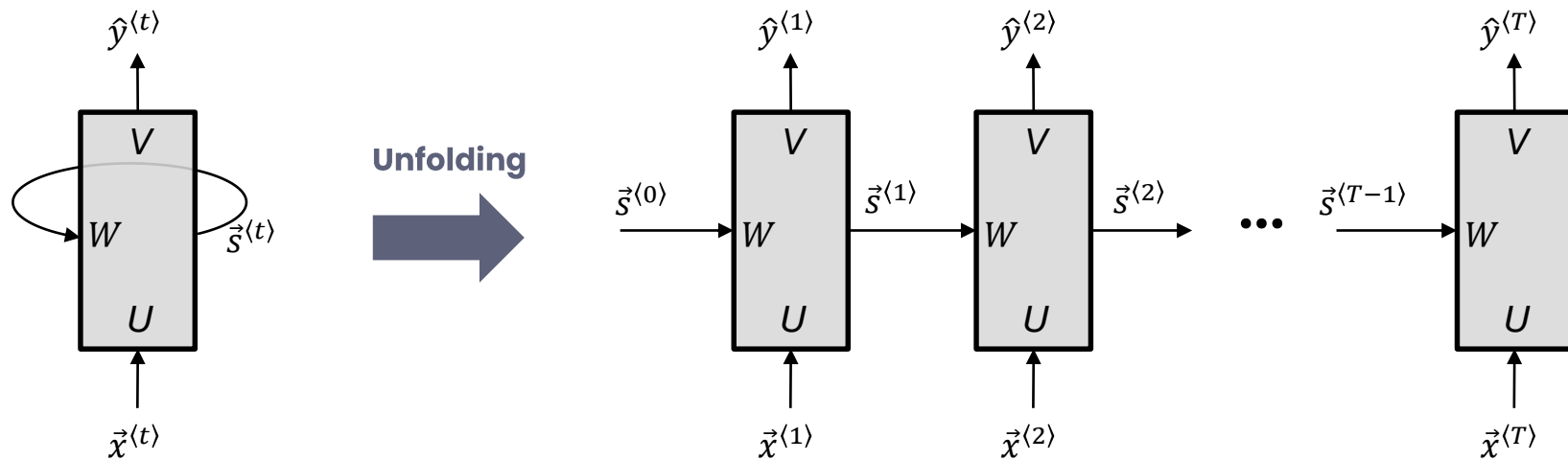
$$\begin{aligned}\vec{s}^{(0)} &= \vec{0} \\ s^{(t)} &= \phi_1(W\vec{s}^{(t-1)} + U\vec{x}^{(t)}) \\ \hat{y}^{(t)} &= \phi_2(V\vec{s}^{(t)})\end{aligned}$$

Forward propagation



$$\begin{aligned}\vec{s}^{(0)} &= \vec{0} \\ s^{(t)} &= \phi_1(W\vec{s}^{(t-1)} + U\vec{x}^{(t)}) \\ \hat{y}^{(t)} &= \phi_2(V\vec{s}^{(t)})\end{aligned}$$

Forward propagation



$$\begin{aligned}\vec{s}^{(0)} &= \vec{0} \\ s^{(t)} &= \phi_1(W\vec{s}^{(t-1)} + U\vec{x}^{(t)}) \\ \hat{y}^{(t)} &= \phi_2(V\vec{s}^{(t)})\end{aligned}$$

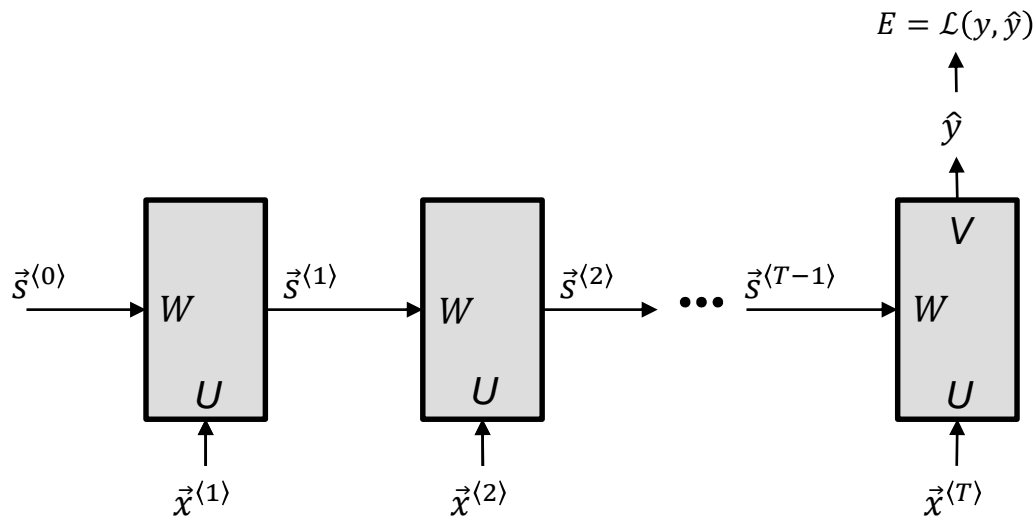
An RNN is a network, not just a neuron, even though it is sometimes called a "cell".

Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.

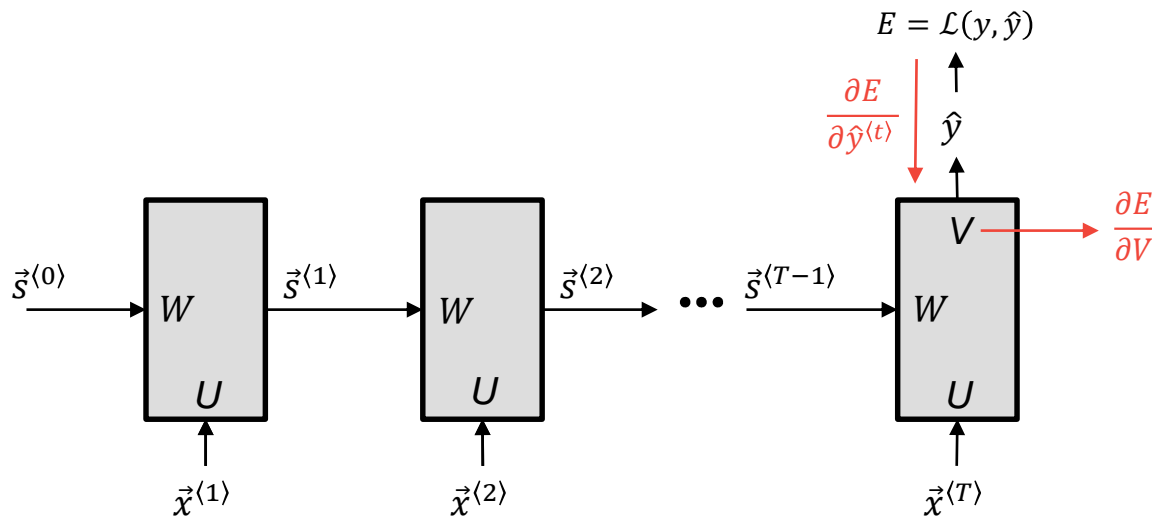


Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.



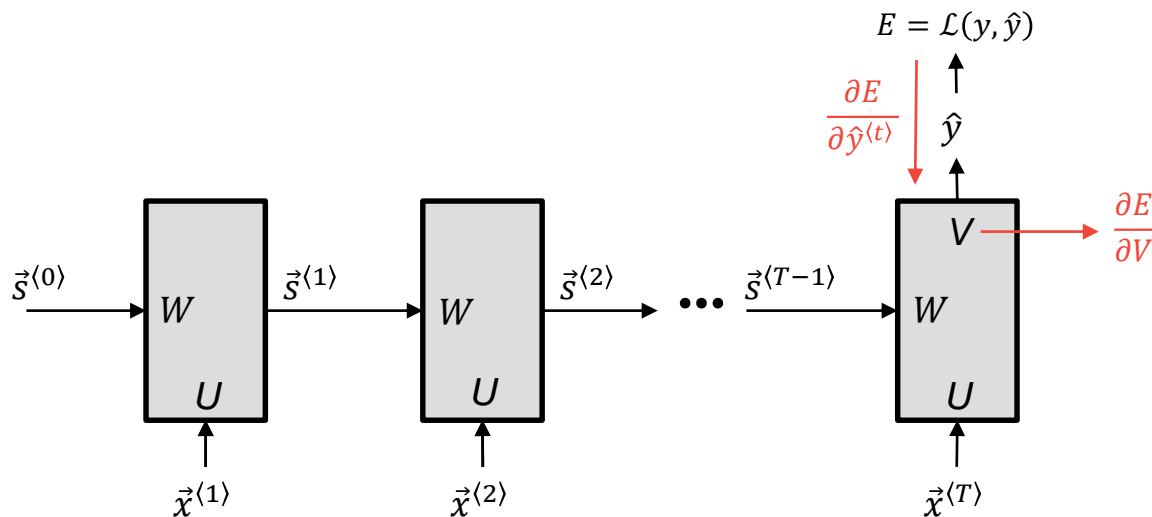
Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

- $\partial E / \partial W$ depends on $\partial \vec{s}^{(T)} / \partial W$, but $\vec{s}^{(T)}$ depends on $\vec{s}^{(T-1)}$, which itself depends on W , and so on.
 - Same for $\partial E / \partial U$.

This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.

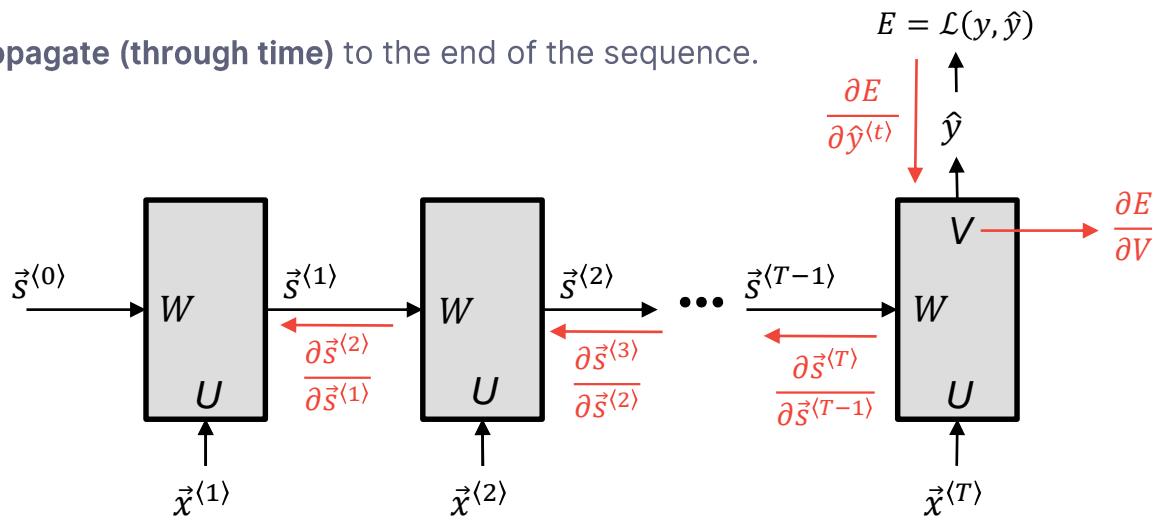


Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

- $\partial E / \partial W$ depends on $\partial \vec{s}^{(T)} / \partial W$, but $\vec{s}^{(T)}$ depends on $\vec{s}^{(T-1)}$, which itself depends on W , and so on.
 - Same for $\partial E / \partial U$.
 - We need to **backpropagate (through time)** to the end of the sequence.



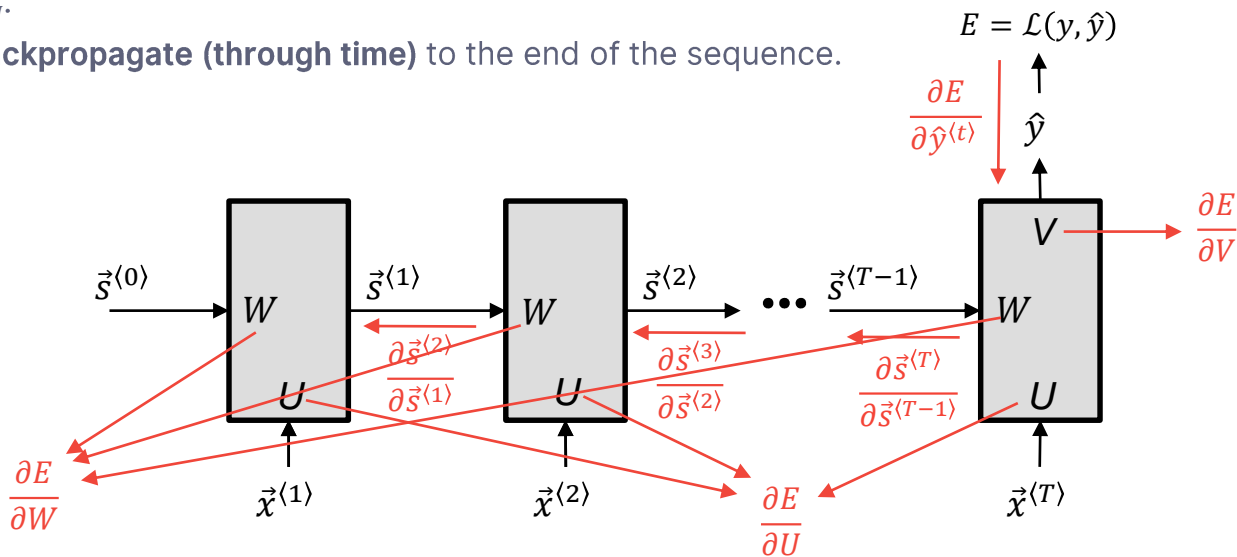
This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.

Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

- $\frac{\partial E}{\partial W}$ depends on $\frac{\partial \vec{s}^{(T)}}{\partial W}$, but $\vec{s}^{(T)}$ depends on $\vec{s}^{(T-1)}$, which itself depends on W , and so on.
 - Same for $\frac{\partial E}{\partial U}$.
 - We need to **backpropagate (through time)** to the end of the sequence.



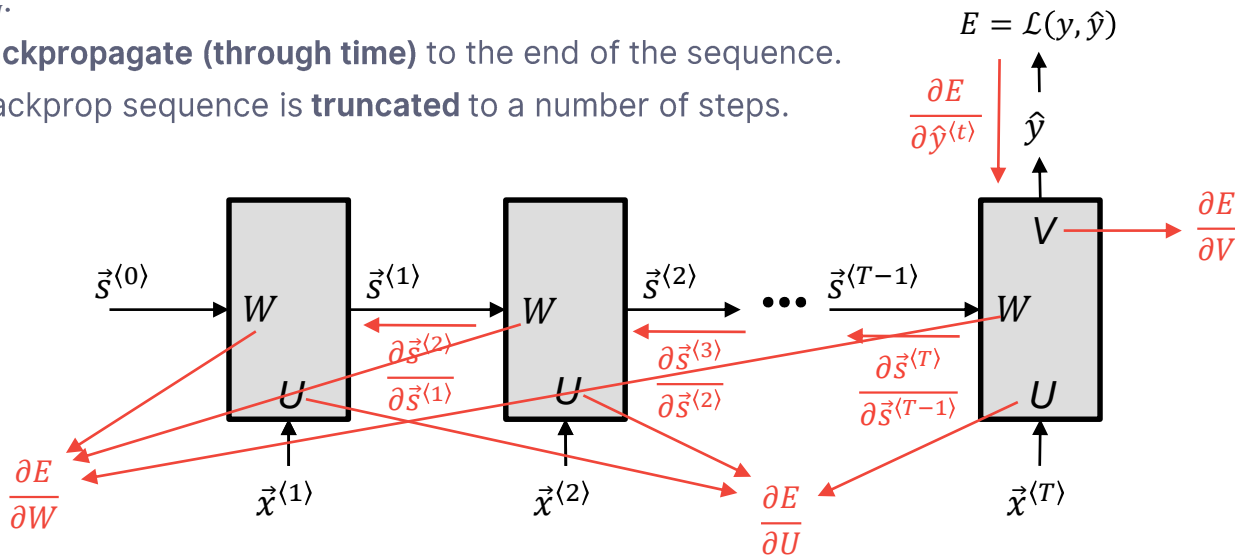
This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.

Backpropagation through time (BPTT)

- Like before, weights are updated through **backpropagation**:

$$\Delta V = -\eta \frac{\partial E}{\partial V} \quad \Delta W = -\eta \frac{\partial E}{\partial W} \quad \Delta U = -\eta \frac{\partial E}{\partial U}$$

- $\frac{\partial E}{\partial W}$ depends on $\frac{\partial \vec{s}^{(T)}}{\partial W}$, but $\vec{s}^{(T)}$ depends on $\vec{s}^{(T-1)}$, which itself depends on W , and so on.
 - Same for $\frac{\partial E}{\partial U}$.
 - We need to **backpropagate (through time)** to the end of the sequence.
 - Usually, the backprop sequence is **truncated** to a number of steps.



This is a simpler “many-to-one” case but it can be generalized to “many-to-many”.

Long Short-Term Memory

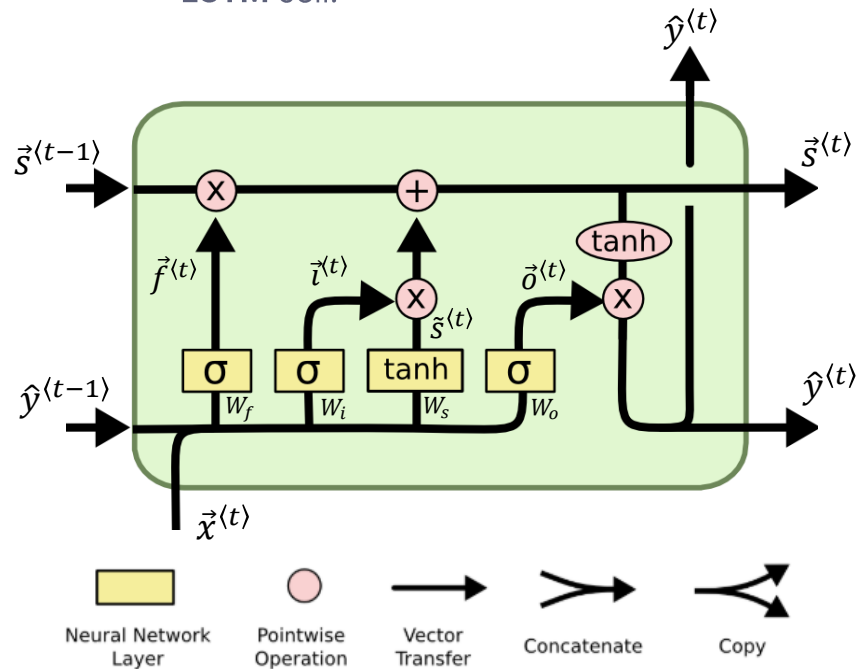
- Because of repeated application of the tanh activation, RNNs suffer from the **vanishing gradient** problem.

Long Short-Term Memory

- Because of repeated application of the tanh activation, RNNs suffer from the **vanishing gradient** problem.
- ReLUs can be used to ameliorate this, but a much better way is to use a different inner structure, called an **LSTM** cell.

Long Short-Term Memory

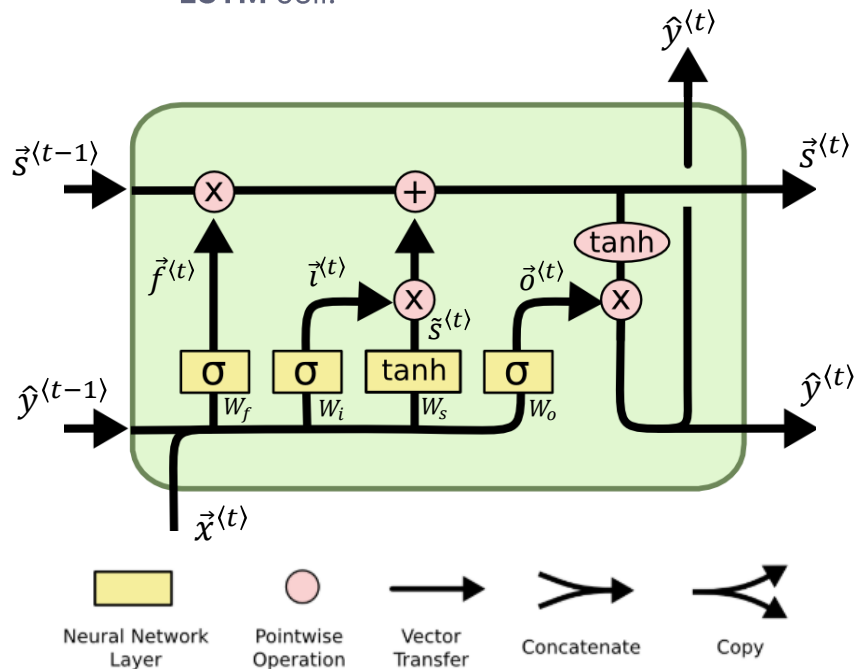
- Because of repeated application of the tanh activation, RNNs suffer from the **vanishing gradient** problem.
- ReLUs can be used to ameliorate this, but a much better way is to use a different inner structure, called an **LSTM** cell.



*Visuals are from <http://colah.github.io/posts/2015-08-Understanding-LSTMs>, but with different notation.

Long Short-Term Memory

- Because of repeated application of the tanh activation, RNNs suffer from the **vanishing gradient** problem.
- ReLUs can be used to ameliorate this, but a much better way is to use a different inner structure, called an **LSTM** cell.



$$\vec{f}^{(t)} = \sigma(W_f \cdot [\hat{y}^{(t-1)}, \vec{x}^{(t)}] + b_f) \quad \text{"Forget" gate}$$

$$\vec{i}^{(t)} = \sigma(W_i \cdot [\hat{y}^{(t-1)}, \vec{x}^{(t)}] + b_i) \quad \text{"Input" gate}$$

$$\tilde{s}^{(t)} = \tanh(W_s \cdot [\hat{y}^{(t-1)}, \vec{x}^{(t)}] + b_s) \quad \text{State candidates}$$

$$\vec{s}^{(t)} = \vec{f}^{(t)} * \vec{s}^{(t-1)} + \vec{i}^{(t)} * \tilde{s}^{(t)} \quad \text{Hidden state (Memory)}$$

$$\vec{o}^{(t)} = \sigma(W_o \cdot [\hat{y}^{(t-1)}, \vec{x}^{(t)}] + b_o) \quad \text{"Output" gate}$$

$$\hat{y}^{(t)} = \vec{o}^{(t)} * \tanh(\vec{s}^{(t)}) \quad \text{Output}$$

*Visuals are from <http://colah.github.io/posts/2015-08-Understanding-LSTMs>, but with different notation.

Keywords

Recurrent Neural Network

RNN

Backpropagation Through Time

BPTT

Truncated BPTT

Long Short-Term Memory

LSTM