

CNNs

January 31st, 2025

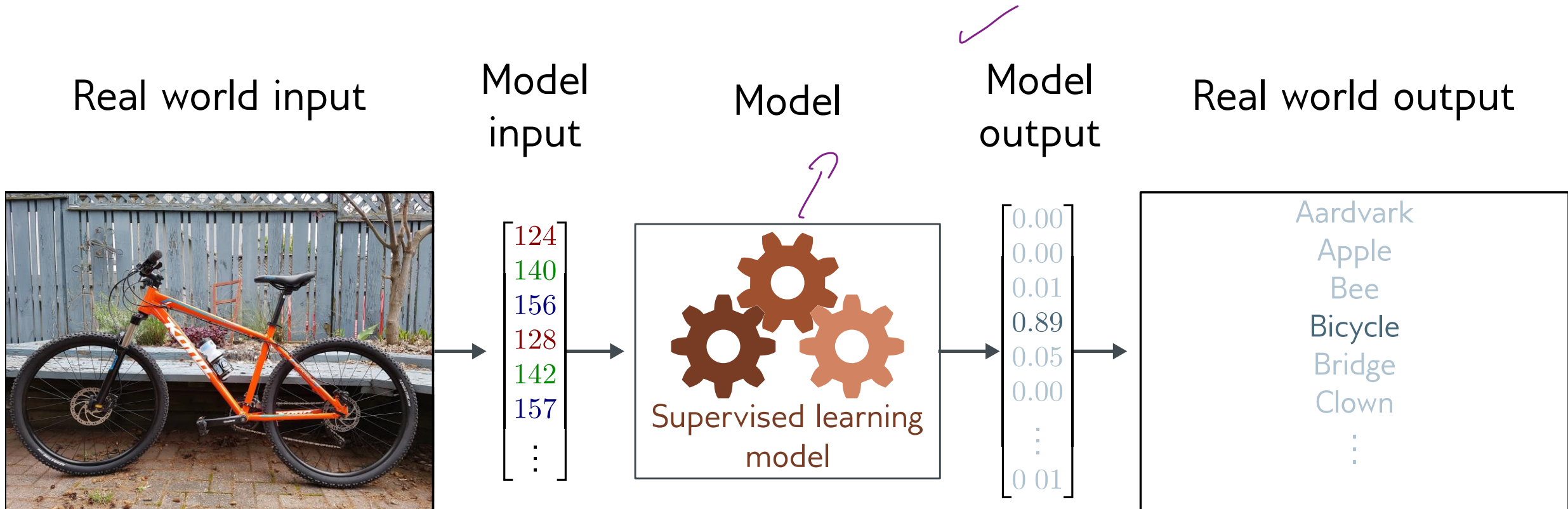
Deep Learning (CS60010)

**Slides adapted from <http://udlbook.com>*

Convolutional networks

- Networks for images
- Invariance and equivariance
- 1D convolution
- Convolutional layers
- Channels
- Receptive fields
- Convolutional network for MNIST 1D

Image classification



- Multiclass classification problem (discrete classes, >2 possible classes)
- Convolutional network

Object detection

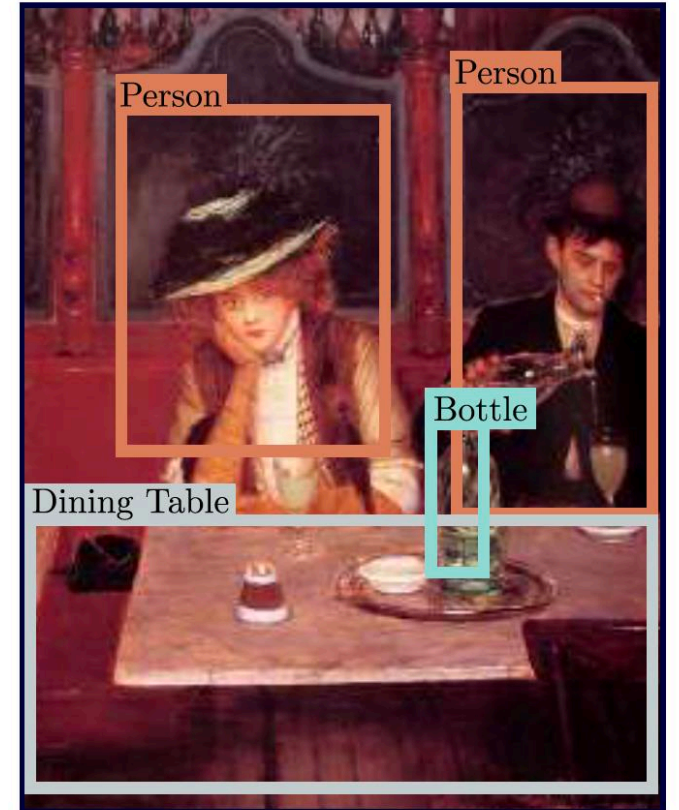
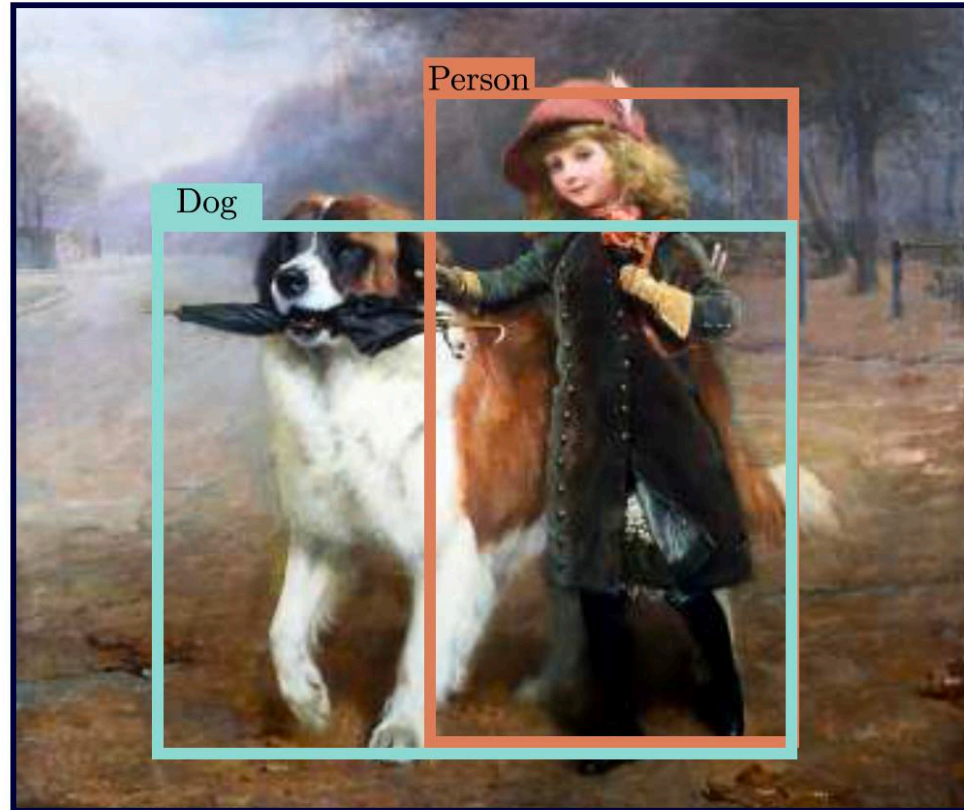
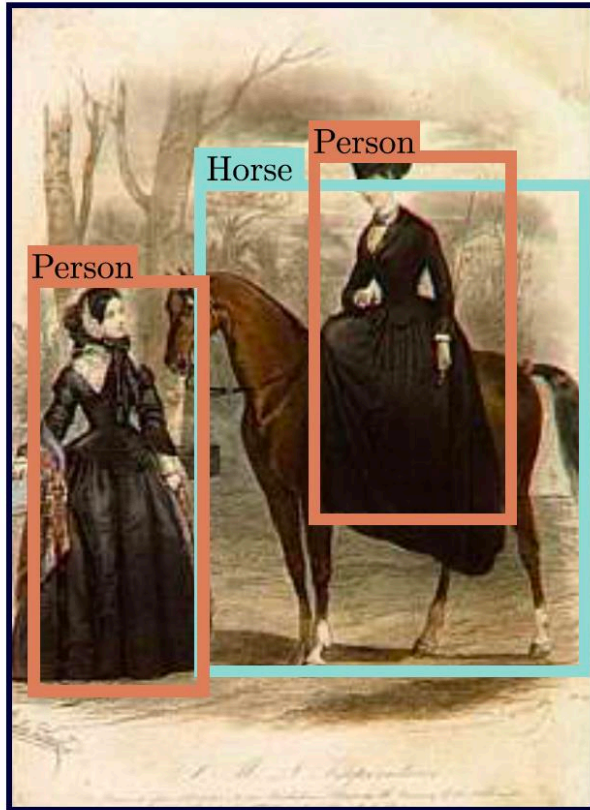
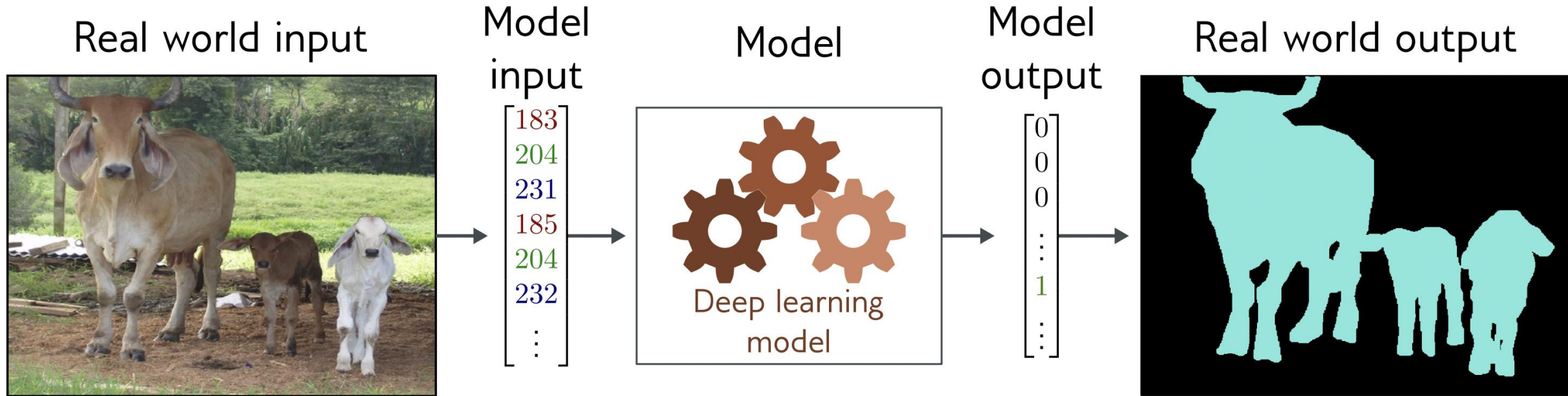


Image segmentation



- Multivariate binary classification problem (many outputs, two discrete classes)
- Convolutional encoder-decoder network

Networks for images

224×224 image
pixels $\times 3$ channels

each value
(0, 255)

- Problems with fully-connected networks

1. Size

- 224×224 RGB image = $150,528$ dimensions
- Hidden layers generally larger than inputs
- One hidden layer = $150,520 \times 150,528$ weights -- 22 billion

2. Nearby pixels statistically related

- But could permute pixels and relearn and get same results with FC

3. Should be stable under transformations

- Don't want to re-learn appearance at different parts of image

Convolutional networks

- Parameters only look at local image patches
- Share parameters across image

Convolutional networks

- Networks for images
- Invariance and equivariance
- 1D convolution
- Convolutional layers
- Channels
- Receptive fields
- Convolutional network for MNIST 1D

Invariance

- A function $f[x]$ is **invariant** to a transformation $t[]$ if:

$$f[t[x]] = f[x]$$

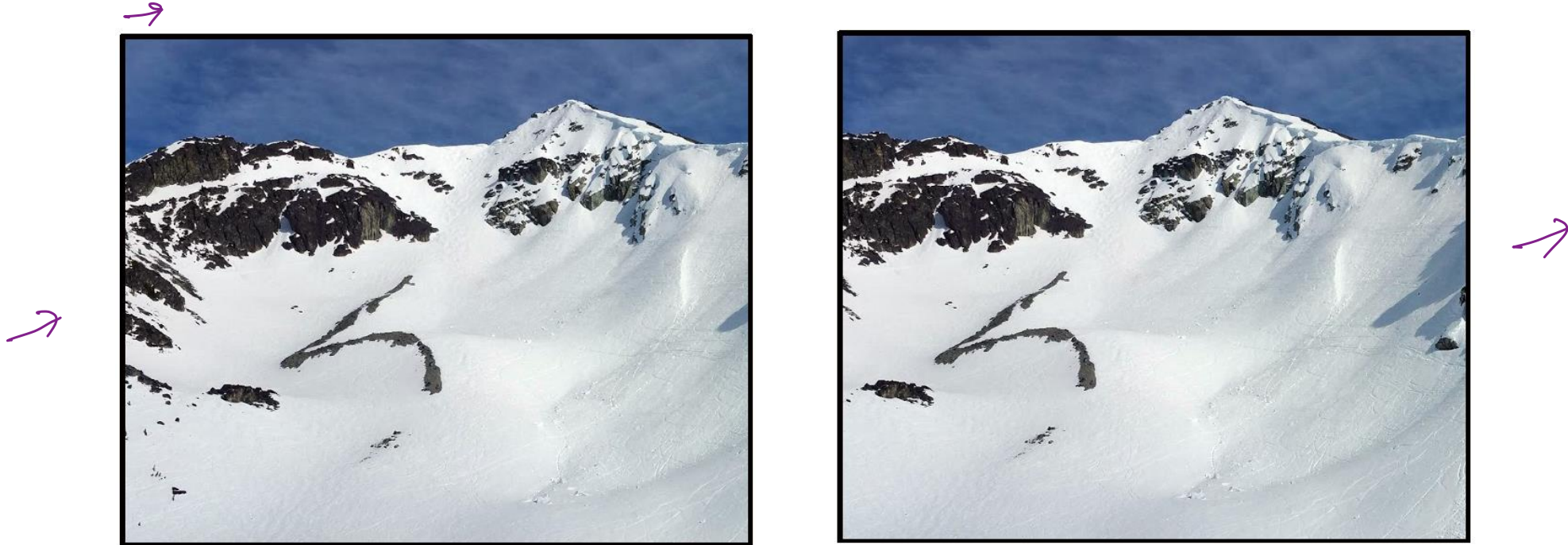
i.e., the function output is the same even after the transformation is applied.

shift

Invariance example

e.g., Image classification

- Image has been translated, but we want our classifier to give the same result



Equivariance

- A function $f[x]$ is **equivariant** to a transformation $t[]$ if:

$$\underset{\nearrow}{f}[\underset{\sim}{t[x]}] = \underset{\nearrow}{t}[\underset{\sim}{f[x]}]$$

i.e., the output is transformed in the same way as the input

Equivariance example


e.g., Image segmentation

- Image has been translated and we want segmentation to translate with it



Convolutional networks



- Networks for images
- Invariance and equivariance
- 1D convolution ✓ 
- Convolutional layers
- Channels
- Receptive fields
- Convolutional network for MNIST 1D

Convolution* in 1D

- Input vector \mathbf{x} :

$$\mathbf{x} = [x_1, x_2, \color{red}{x_i}, x_{I+1}, x_I]$$

$$z_i = \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}$$

- Output is weighted sum of neighbors:

$$z_i = \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}$$

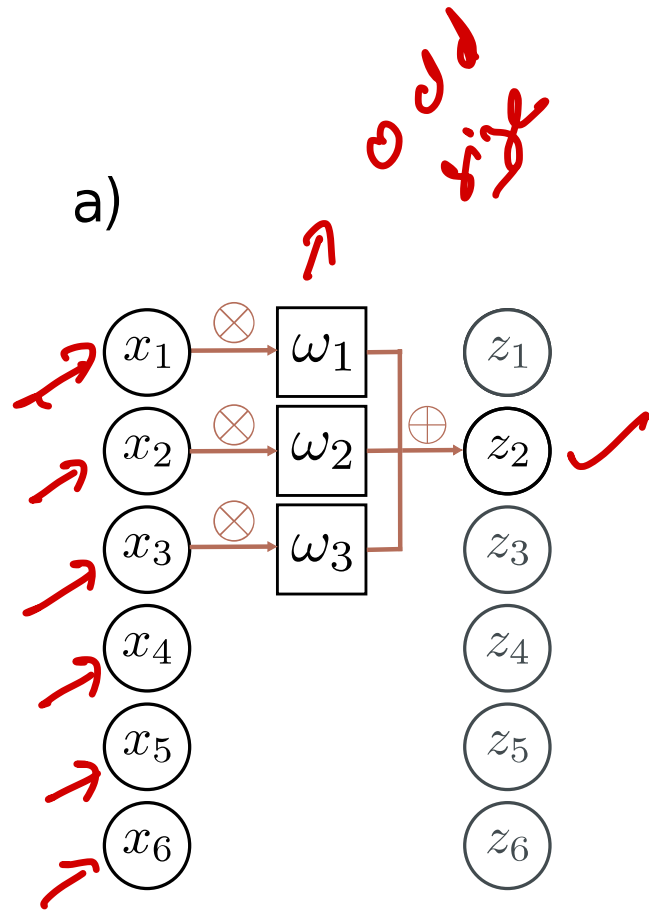
- Convolutional kernel or filter:

$$\boldsymbol{\omega} = [\omega_1, \omega_2, \omega_3]^T$$

Kernel size = 3

* Not really technically convolution

Convolution with kernel size 3



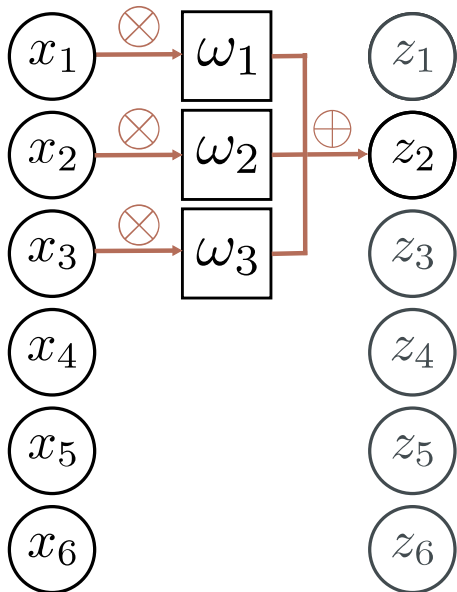
x_{-1}
 x_0
 x_1
 x_2
 x_3
 x_4

$$\begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix}$$

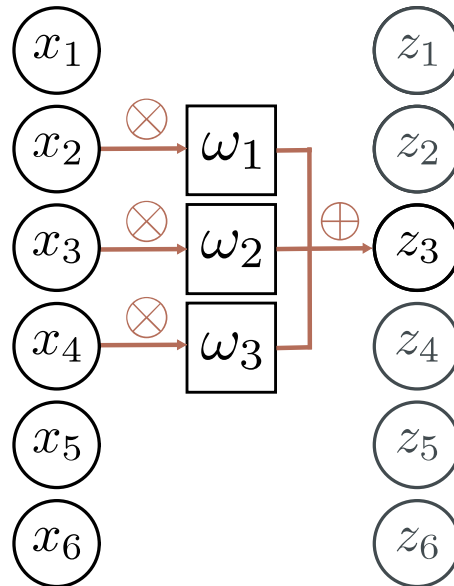
z_2 ✓

Convolution with kernel size 3

a)



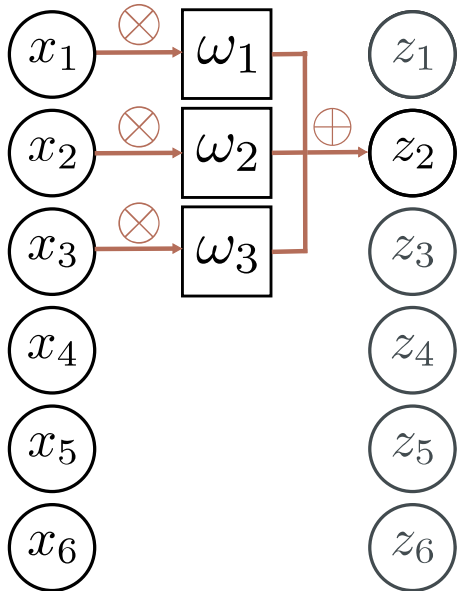
b)



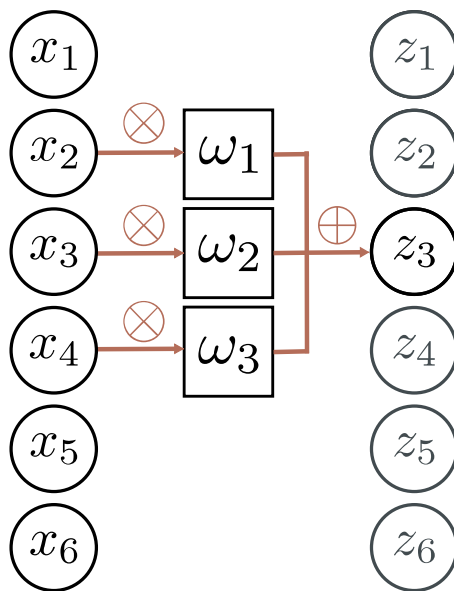
Equivariant to translation of input
 $\mathbf{f}[\mathbf{t}[\mathbf{x}]] = \mathbf{t}[\mathbf{f}[\mathbf{x}]]$

Zero padding

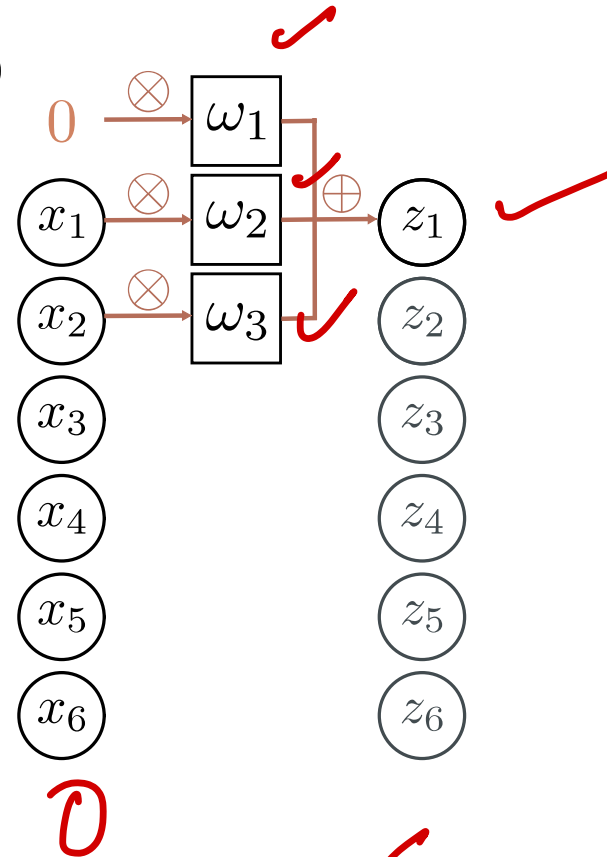
a)



b)

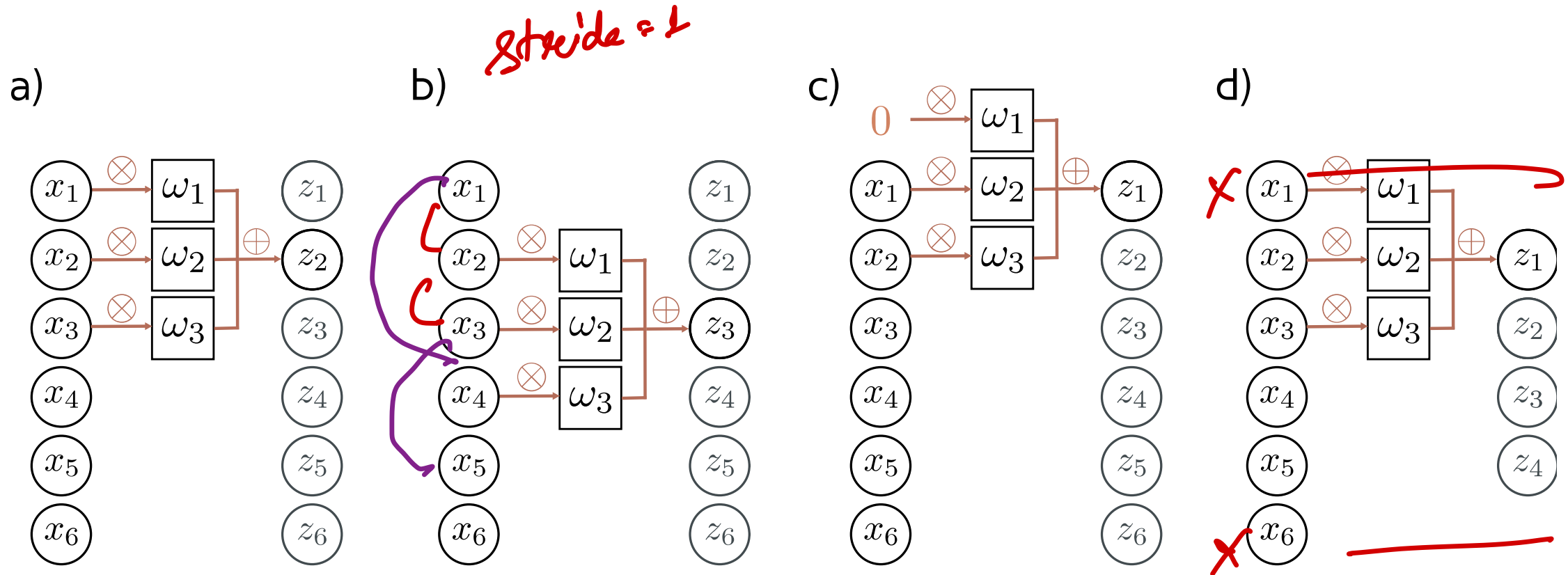


c)



Treat positions that are beyond end of the input as zero.

“Valid” convolutions

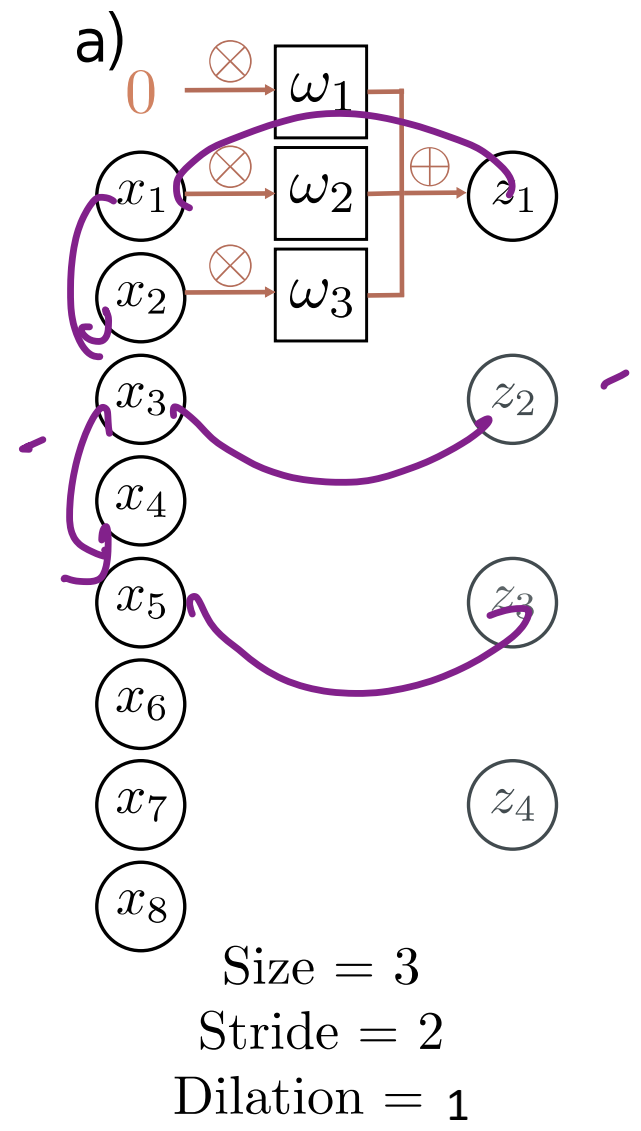


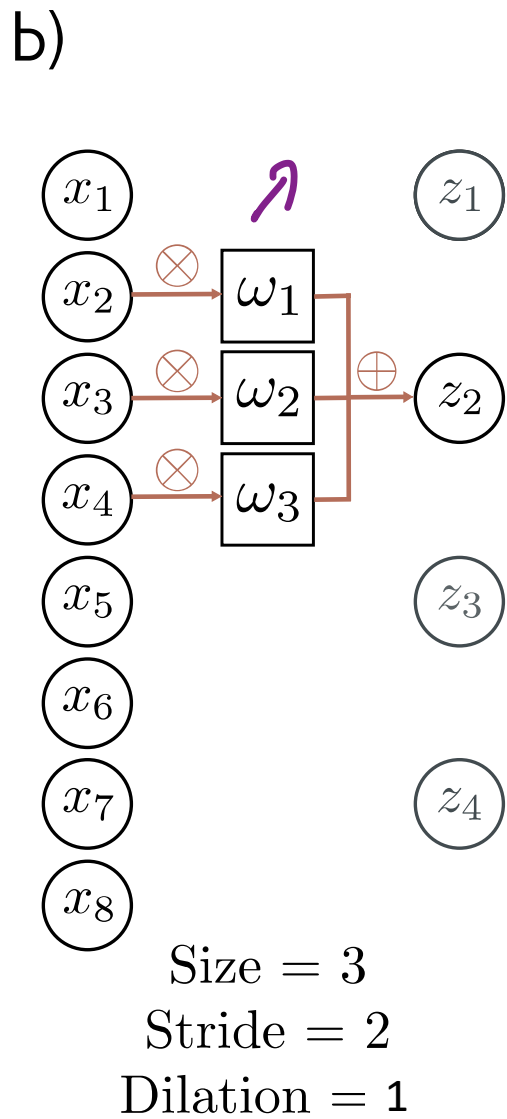
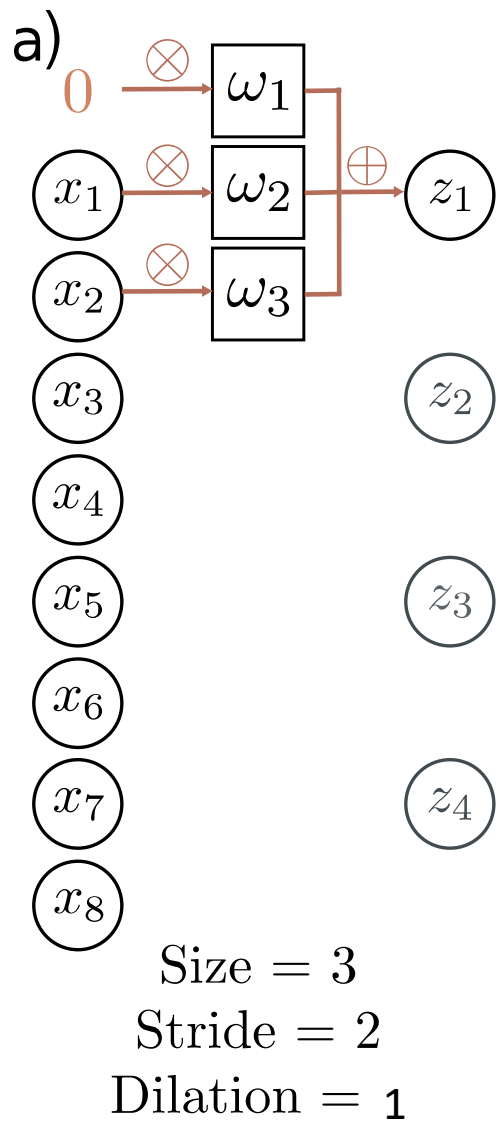
Stride, kernel size, and dilation

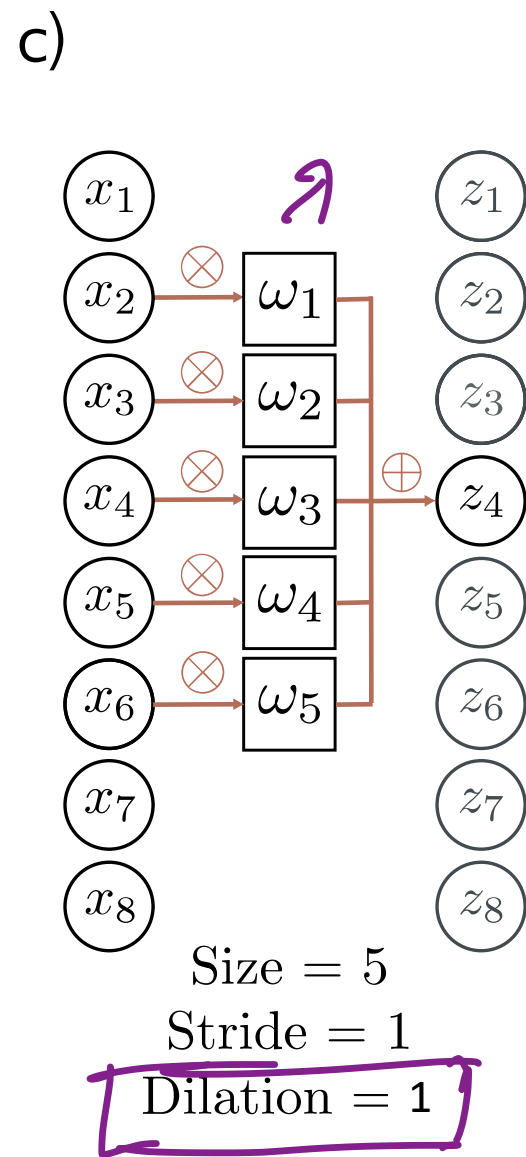
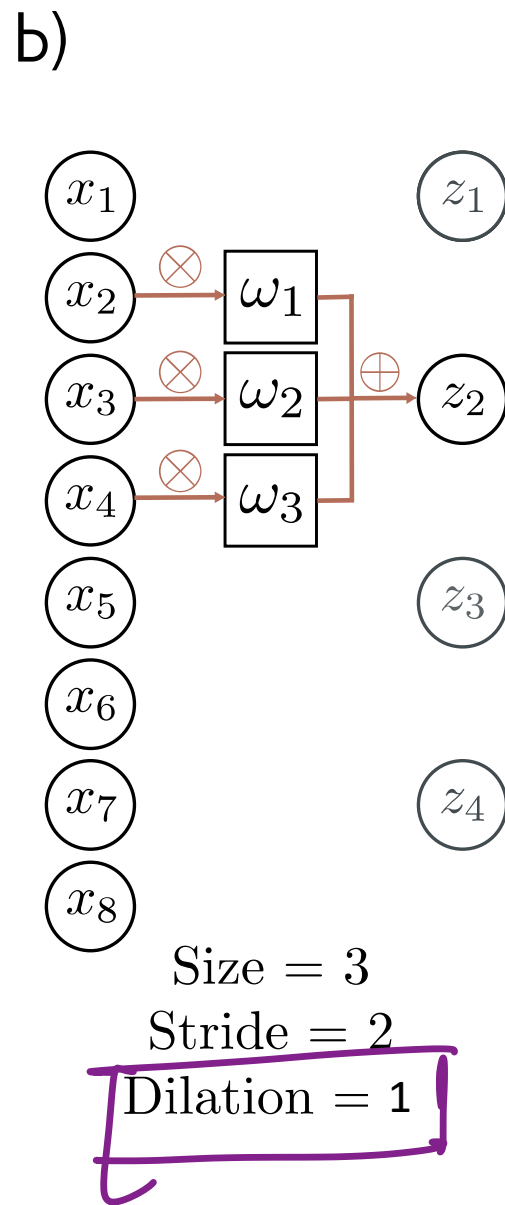
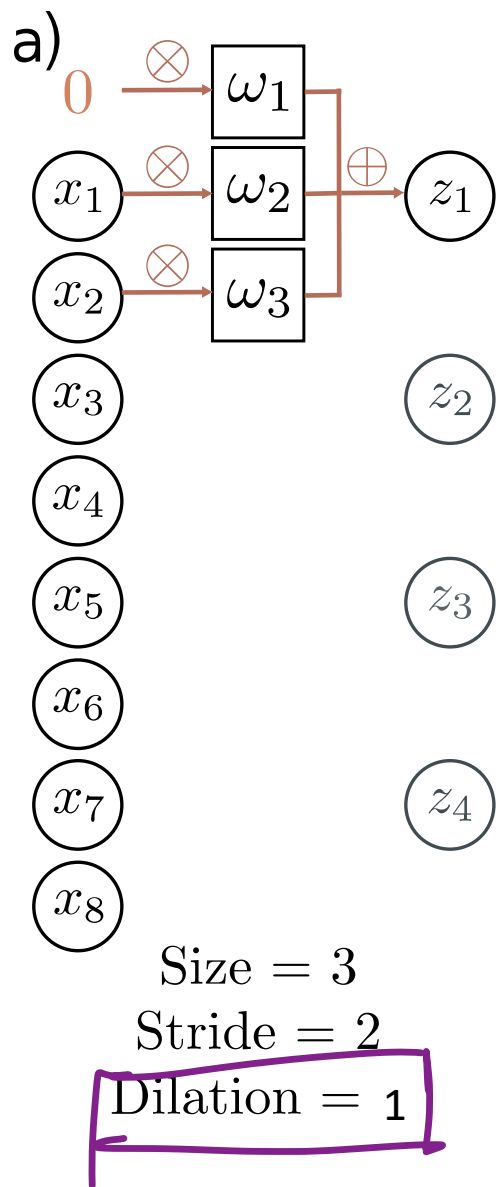
- **Stride** = shift by k positions for each output
 - Decreases size of output relative to input
- **Kernel size** = weight a different number of inputs for each output
 - Combine information from a larger area
 - But kernel size 5 uses 5 parameters ✓
- **Dilated** convolutions = intersperse kernel values with zeros
 - Combine information from a larger area
 - Fewer parameters

Kernel size = 3
 stride = 2
 dilation = 1

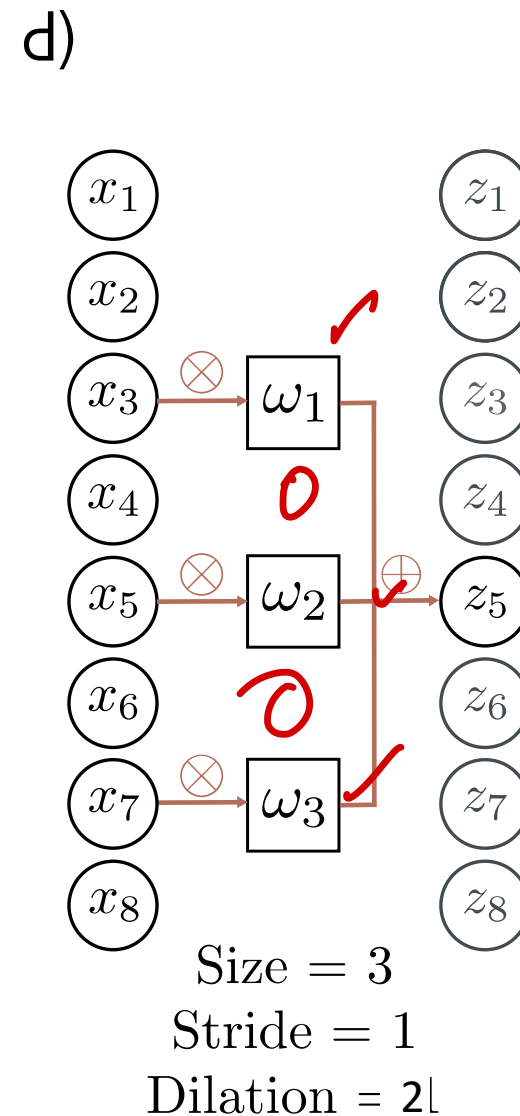
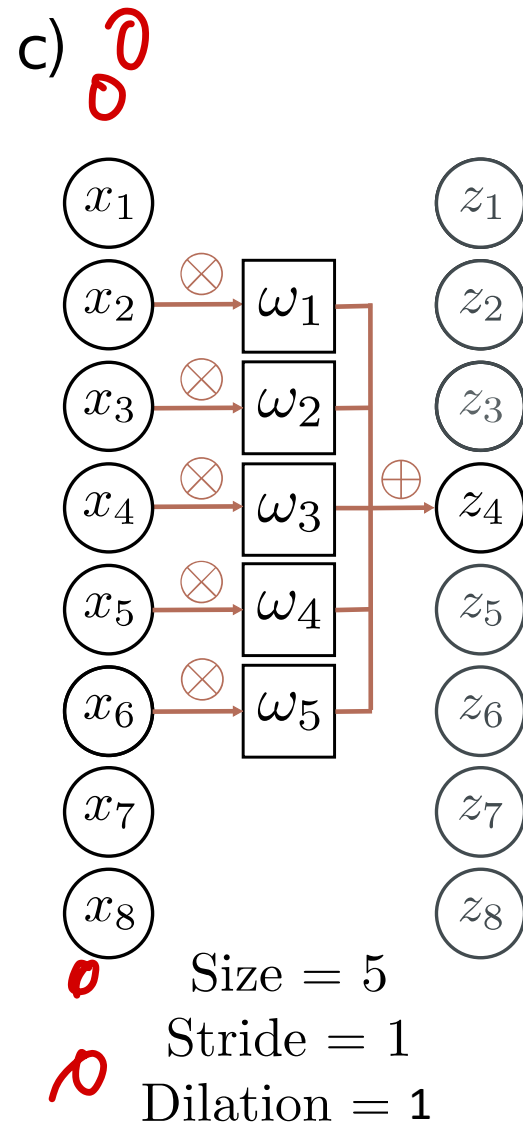
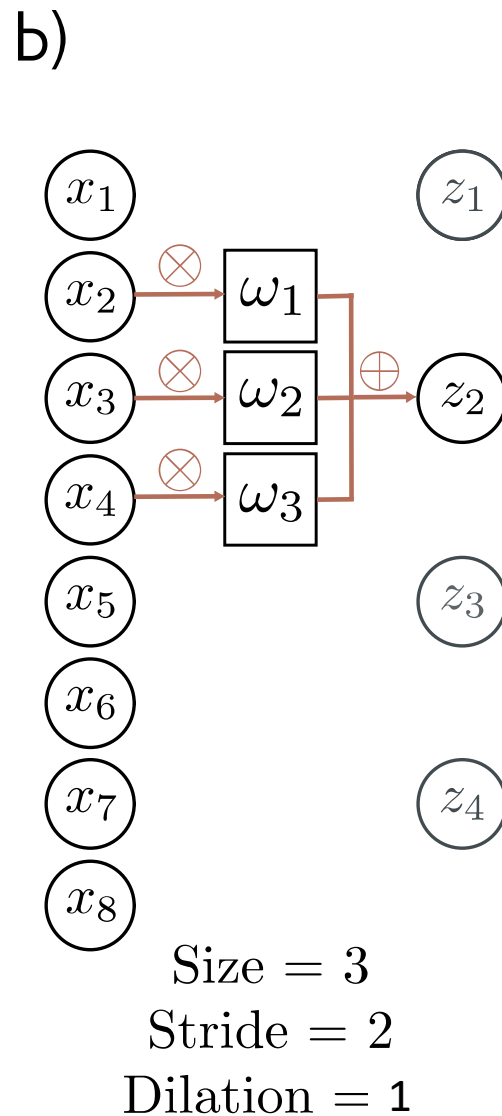
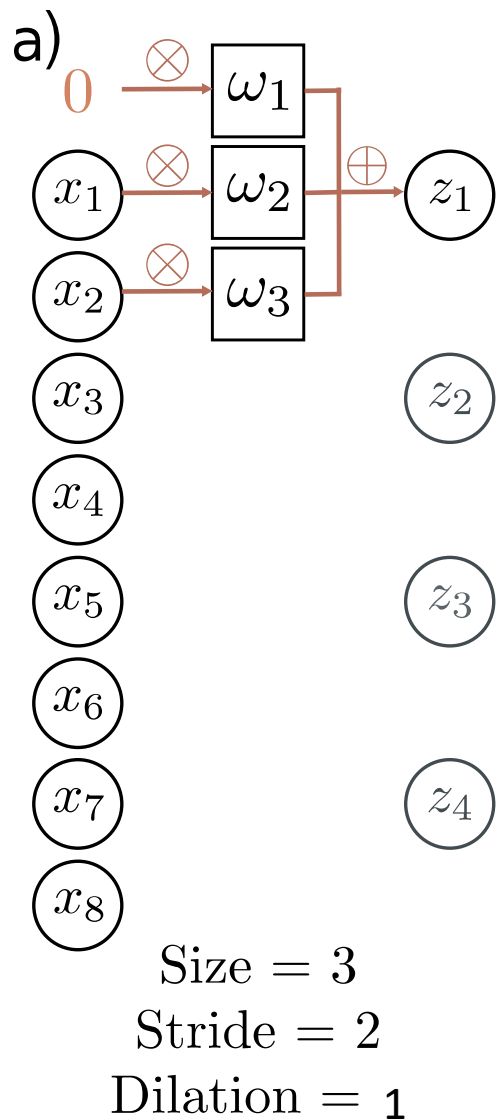
dilation 2	dilation 3
w_1	w_1
0	0
0	0
w_2	w_2
0	0
0	0
w_3	w_3
	0
	0
	w_4







Zero-padding



Convolutional networks

- Networks for images
- Invariance and equivariance
- 1D convolution
- Convolutional layers
- Channels
- Receptive fields
- Convolutional network for MNIST 1D

FC.

$\alpha_1 - -$

β

α_n

Convolutional layer

$$\begin{aligned} h_i &= a[\beta + \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}] \\ &= a \left[\beta + \sum_{j=1}^3 \omega_j x_{i+j-2} \right] \end{aligned}$$

Special case of fully-connected network

Convolutional network:

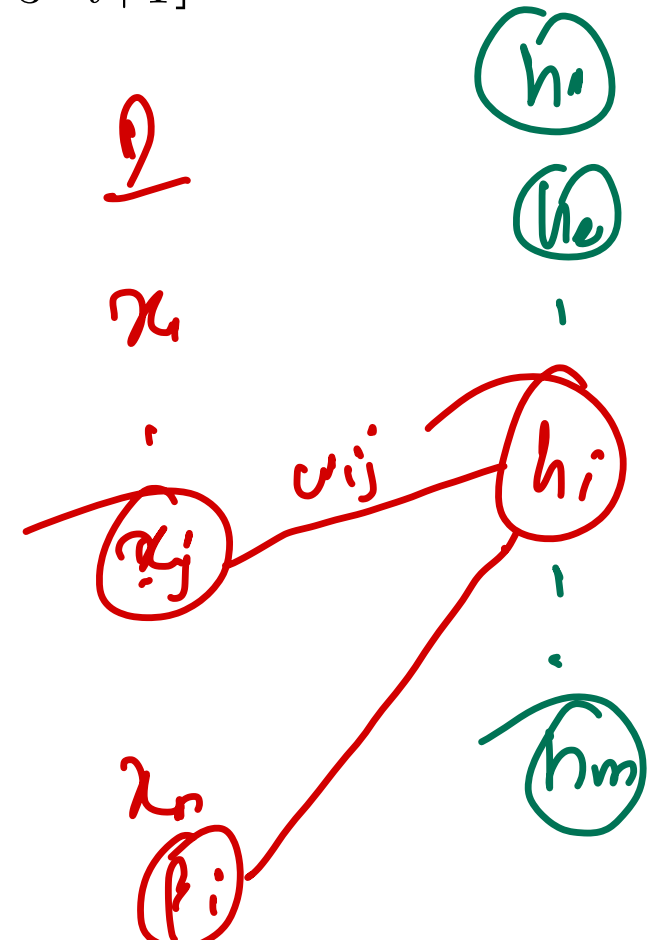
$$h_i = a[\beta + \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}]$$
$$= a\left[\beta + \sum_{j=1}^3 \omega_j x_{i+j-2}\right]$$

1D
1D
3 weights
1 bias

Fully connected network:

2 weights
1 bias

$$h_i = a\left[\beta_i + \sum_{j=1}^D \omega_{ij} x_j\right]$$



Special case of fully-connected network

Convolutional network:

$$h_i = a[\beta + \omega_1 x_{i-1} + \omega_2 x_i + \omega_3 x_{i+1}]$$

$$= a \left[\beta + \sum_{j=1}^3 \omega_j x_{i+j-2} \right]$$

inductive
bias

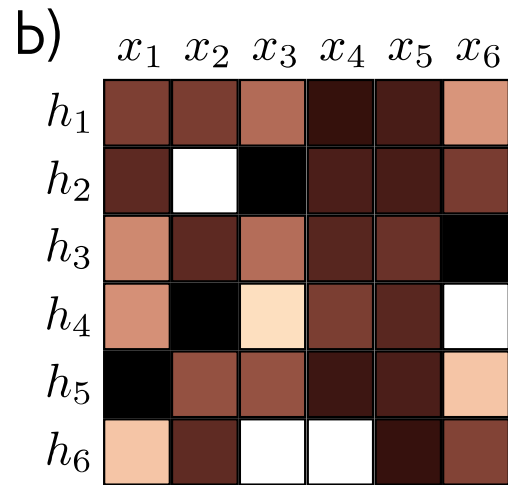
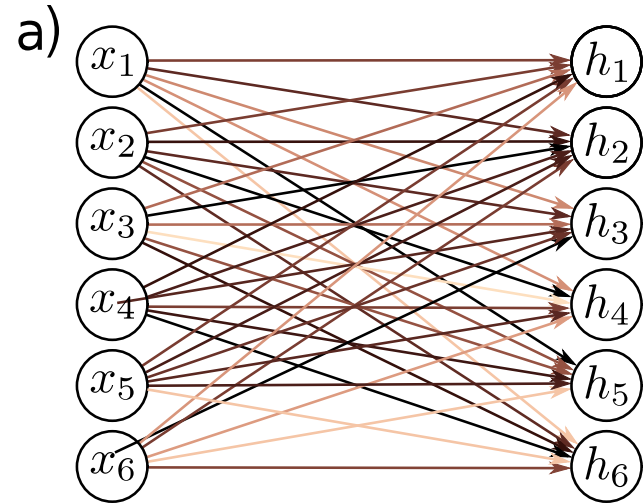
3 weights, 1 bias

Fully connected network:

$$h_i = a \left[\beta_i + \sum_{j=1}^D \omega_{ij} x_j \right]$$

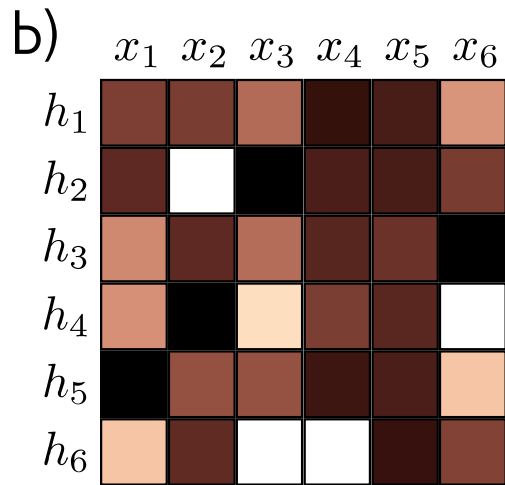
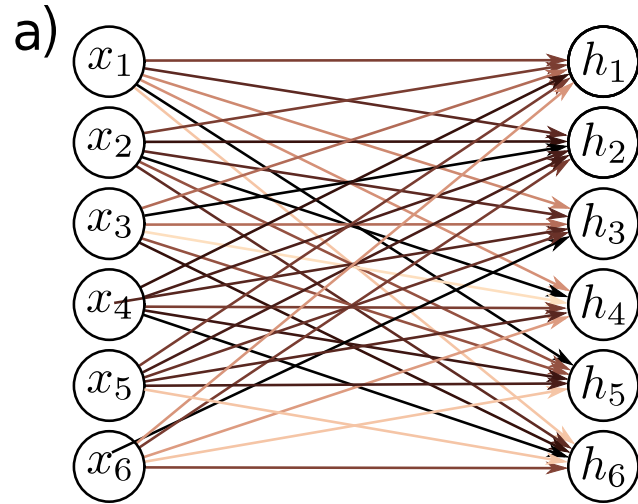
D^2 weights, D biases

Special case of fully-connected network

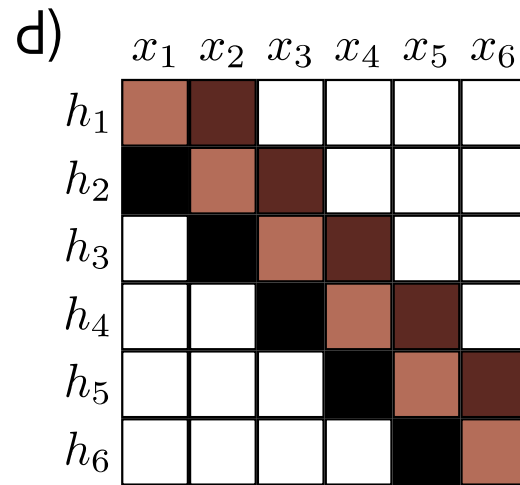
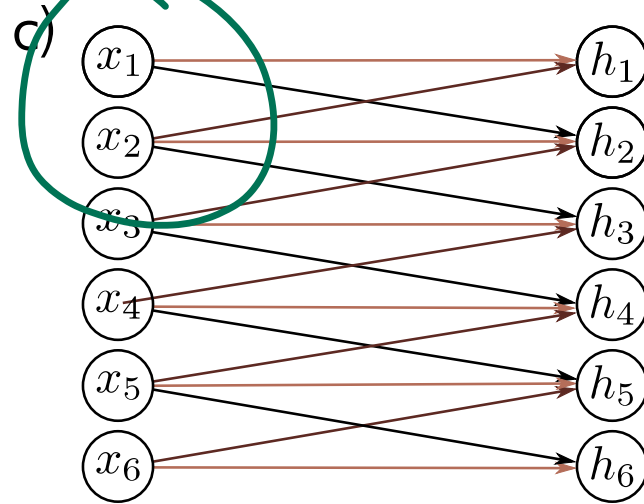


Fully connected network

Special case of fully-connected network

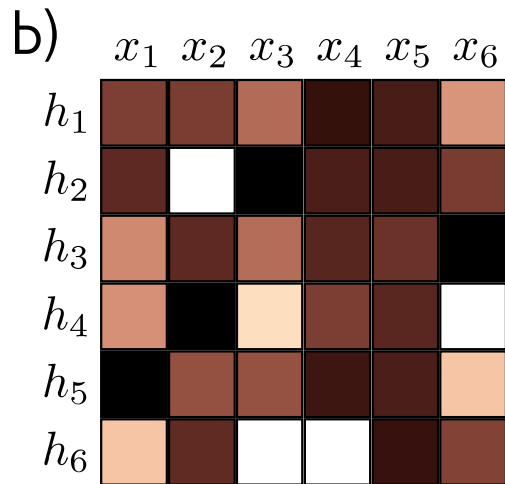
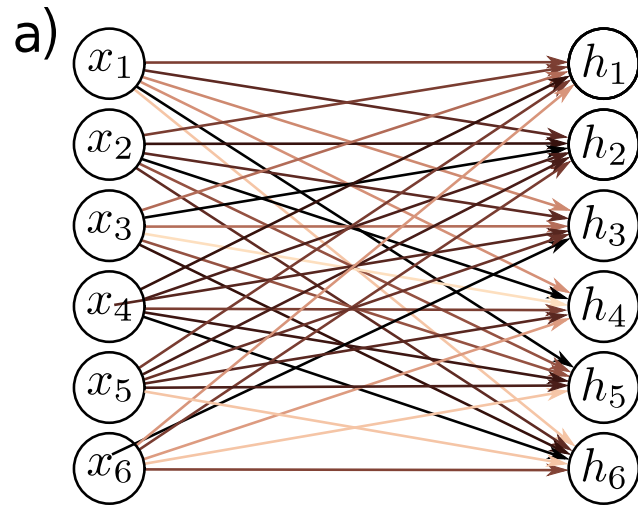


Fully connected network

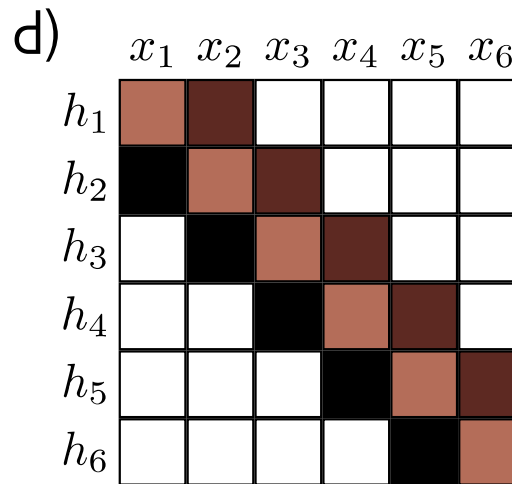
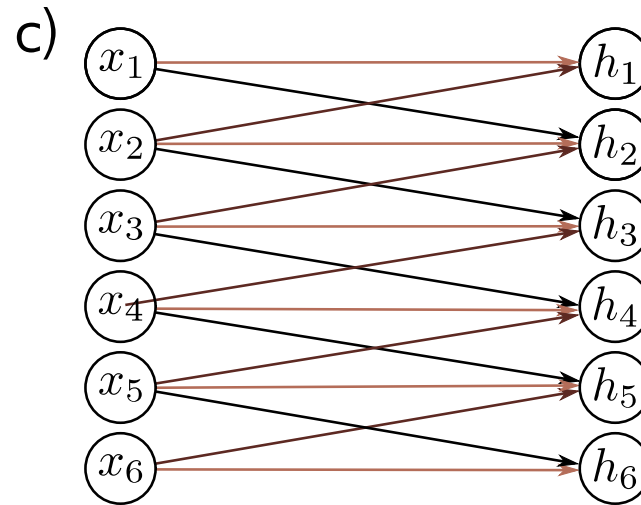


Convolution, kernel 3,
stride 1, dilation 1

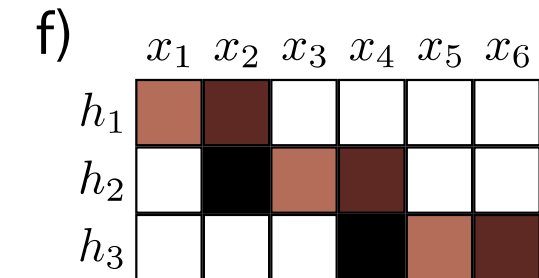
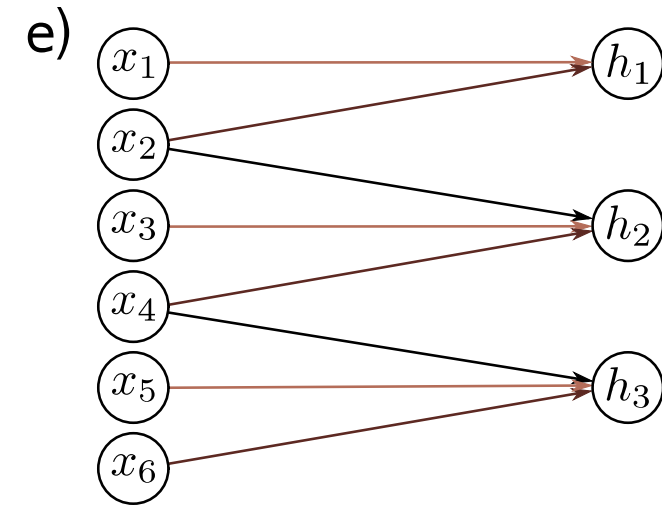
Special case of fully-connected network



Fully connected network



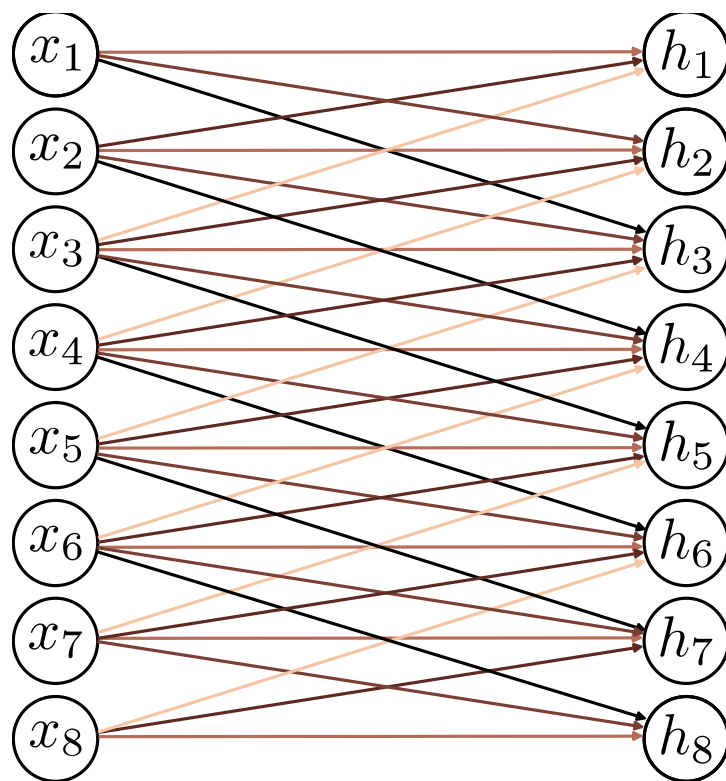
Convolution, size 3, stride 1,
dilation 1, zero padding



Convolution, size 3, stride 2,
dilation 1, zero padding

Question 1

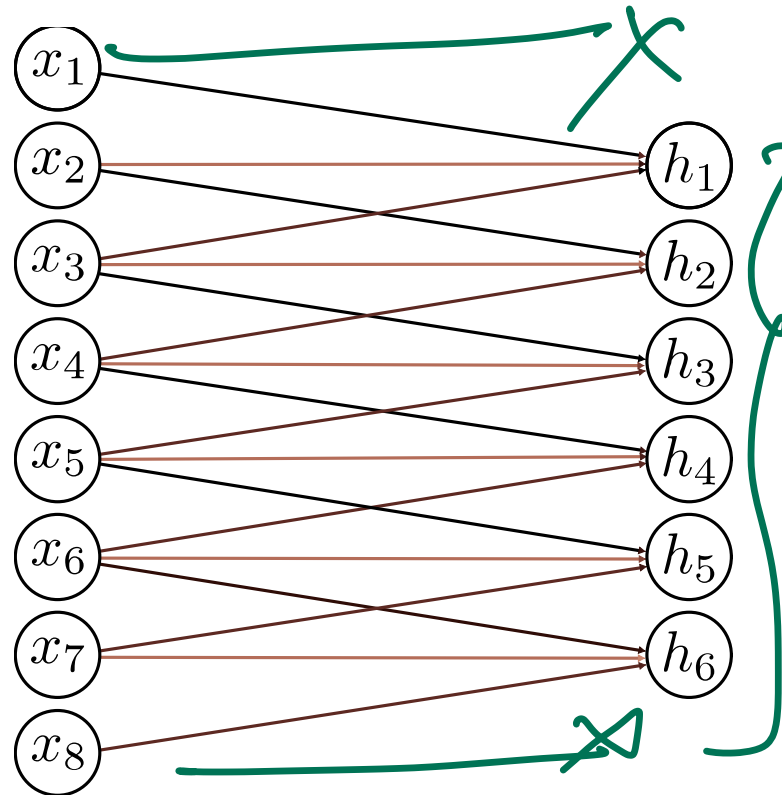
- Kernel size? 5
- Stride? 1
- Dilation? 1
- Zero padding / valid? \checkmark



	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
h_1								
h_2								
h_3								
h_4								
h_5								
h_6								
h_7								
h_8								

Question 2

- Kernel size? 3
- Stride? 1
- Dilation? 1
- Zero padding / valid? ☒ zero padding ☒ valid?

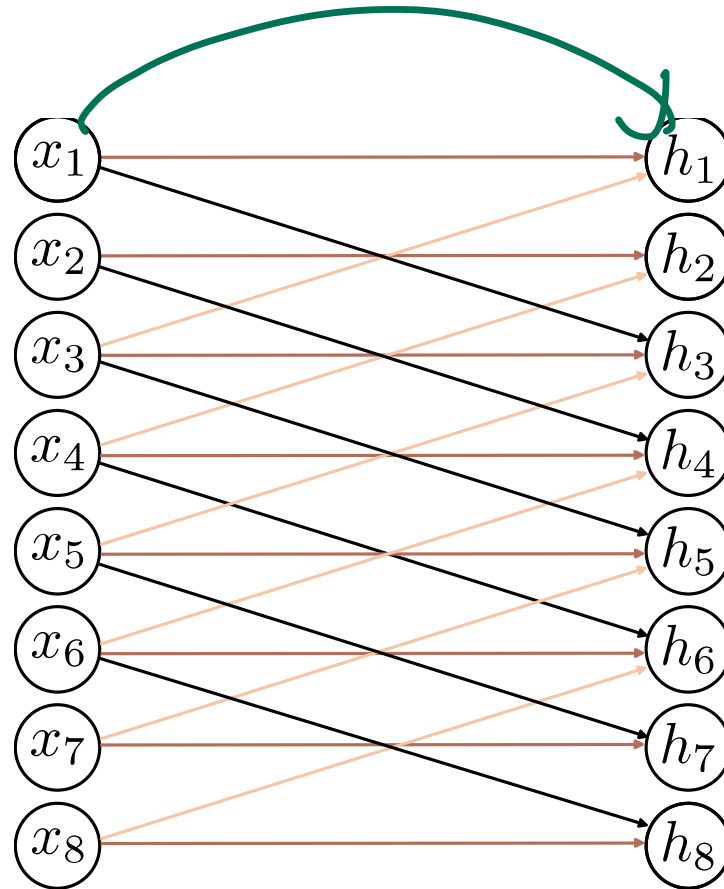


2

	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
h_1								
h_2								
h_3								
h_4								
h_5								
h_6								

Question 3

- Kernel size? 3
- Stride? 1
- Dilation? 2
- Zero padding / valid? ✓ ✗



	x_1	x_2	x_3	x_4	x_5	x_6	x_7	x_8
h_1	■		■					
h_2		■		■				
h_3	■		■		■			
h_4		■		■		■		
h_5			■		■		■	
h_6				■		■		■
h_7					■		■	
h_8						■		■

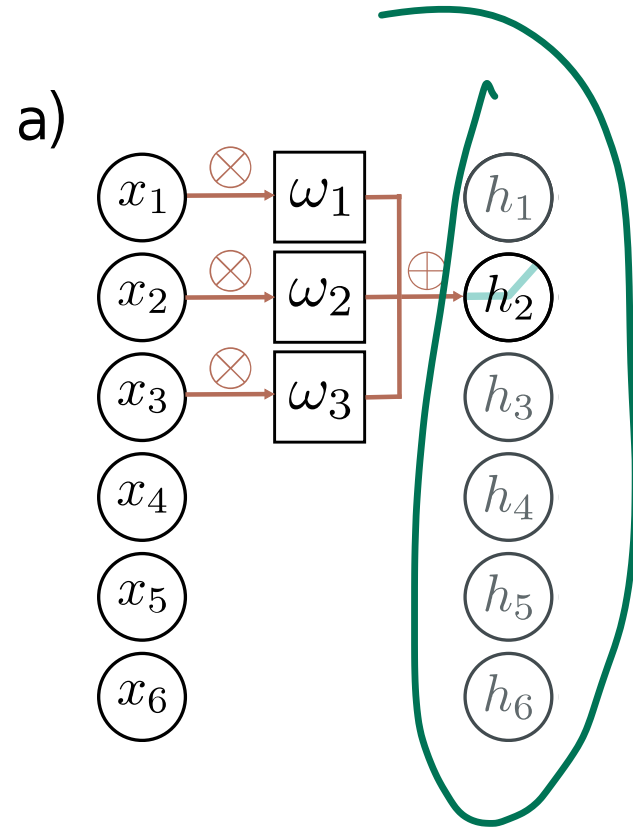
Convolutional networks

- Networks for images
- Invariance and equivariance
- 1D convolution
- Convolutional layers
- Channels
- Receptive fields
- Convolutional network for MNIST 1D

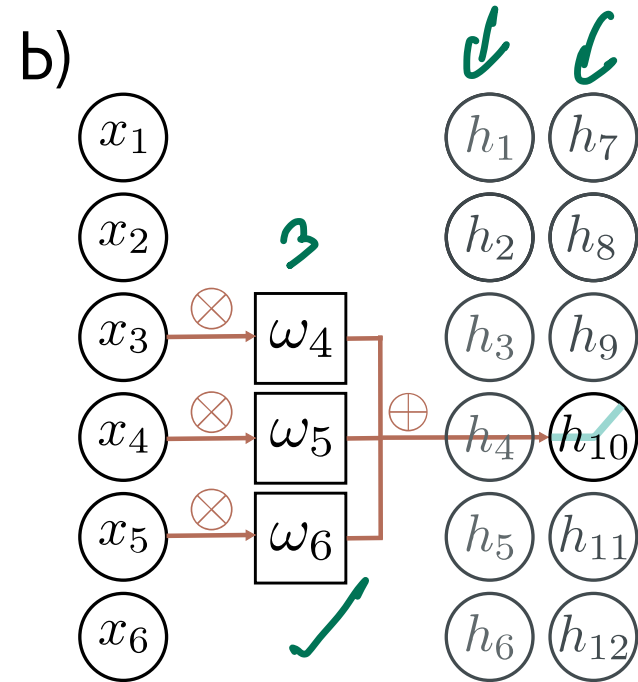
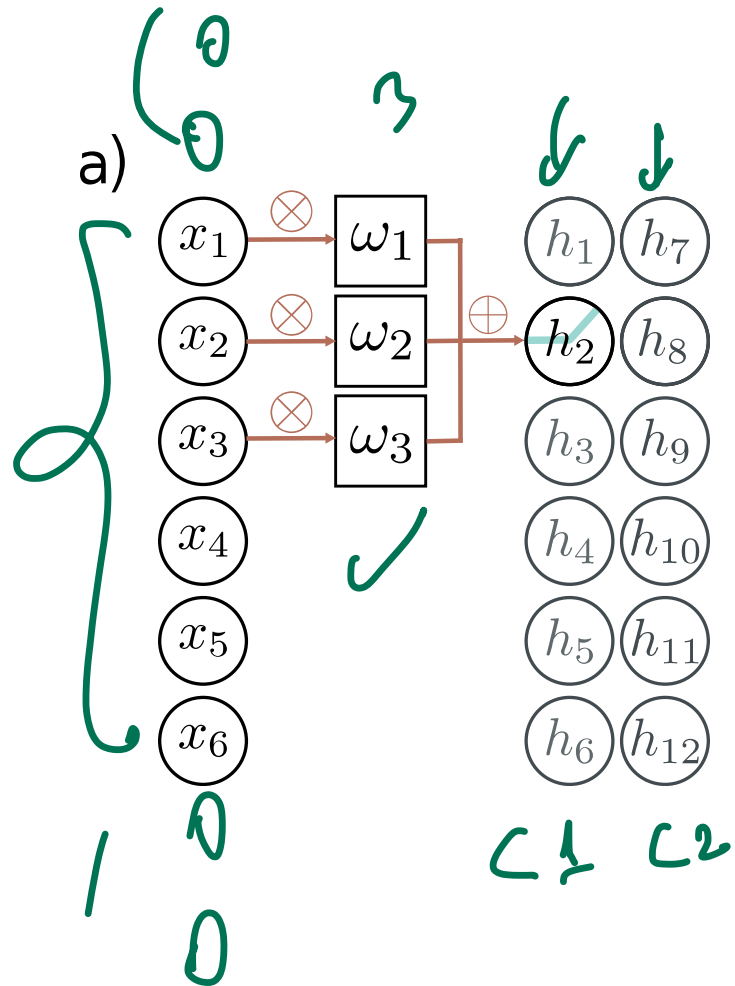
Channels

- The convolutional operation averages together the inputs
- Plus passes through ReLU function
- Has to lose information
- Solution:
 - apply several convolutions and stack them in channels
 - Sometimes also called feature maps

Two output channels, one input channel

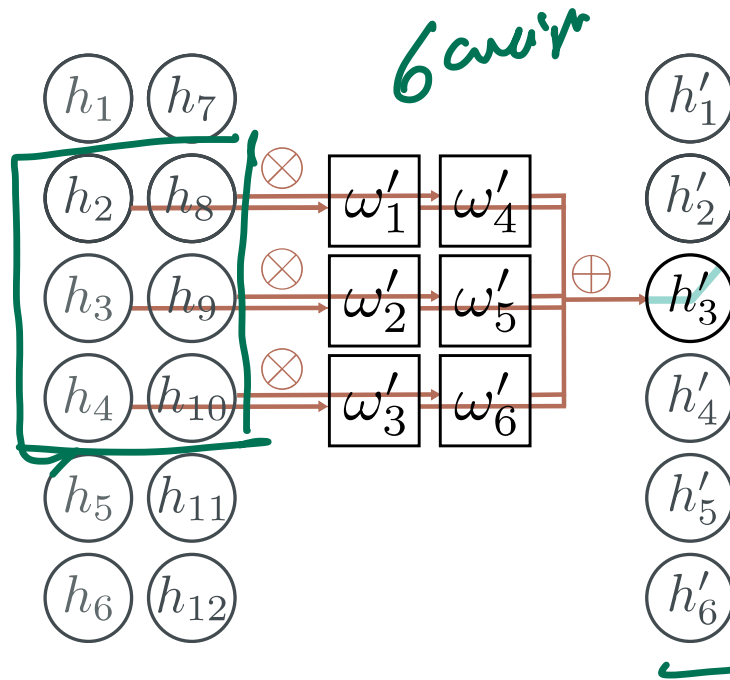


Two output channels, one input channel



6 - learnable parameters

Two input channels, one output channel



How many parameters?

$C_i \times K$

- If there are C_i input channels and kernel size K

$$\Omega \in \mathbb{R}^{C_i \times K}$$

$$\beta \in \mathbb{R}$$

- If there are C_i input channels and C_o output channels

$$\Omega \in \mathbb{R}^{C_i \times C_o \times K}$$

$$\beta \in \mathbb{R}^{C_o}$$