# Foundations of DL

Deep Learning

Alfredo Canziani, Ritchie Ng

@alfcnz, @RitchieNg

ALF

# Convolutional Neural Nets

Exploiting stationarity, locality, and compositionality of natural data

# Signals can be represented as vectors



$$\boldsymbol{x} = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_t & \dots \end{bmatrix}^\top$$

$x_t$ are waveform heights



$$\boldsymbol{x} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} & x_{21} & x_{22} & \dots \end{bmatrix}^\top$$

$x_{ij}$ are pixel values

"John picked up the apple"

$$\boldsymbol{x} = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \end{bmatrix}^\top$$

$x_t$ are one-hot vectors

# Signals can be represented as vectors

$$\boldsymbol{x} = [$$

$x_t$ are

$x_1$

e

$$[x_1 \;\; x_2 \;\; x_3 \;\; x_4 \;\; x_5]^\top$$

"John picked up the apple"

$x_t$ are one-hot vectors

# Signals can be represented as vectors



$$\boldsymbol{x} = [x_1 \quad x_2 \quad x_3 \quad \ldots \quad x_t \quad \ldots]^\top$$

$x_t$ are waveform heights

$$\boldsymbol{x} = [x_{11} \quad x_{12} \quad \ldots \quad x_{1n} \quad x_{21} \quad x_{22} \quad \ldots]^\top$$

$x_{ij}$ are pixel values

"John picked up the apple"

$$\boldsymbol{x} = [x_1 \quad x_2 \quad x_3 \quad x_4 \quad x_5]^\top$$

$x_t$ are one-hot vectors

# Fully connected (FC) layer

$$\hat{\boldsymbol{y}}$$

$$\boldsymbol{h} = f(\boldsymbol{W_h}\boldsymbol{x} + \boldsymbol{b_h})$$

$$\hat{\boldsymbol{y}} = g(\boldsymbol{W_y}\boldsymbol{h} + \boldsymbol{b_y})$$

*j*-th row of $W^{(1)}$

$$a_j^{(2)} = f(\boxed{\boldsymbol{w}^{(j)}}\boldsymbol{x} + b_j) = f\Big(\big(\sum_{i=1}^{n} w_i^{(j)} x_i\big) + b_j\Big)$$

$g$

$W_y$

$\boldsymbol{h}$

$f$

$W_h$

$\boldsymbol{x}$

$$f, g = (\cdot)^+, \sigma(\cdot),$$
$$\tanh(\cdot), \mathrm{softmax}(\cdot)$$



$\boldsymbol{x}$

$\boldsymbol{a}^{(1)}$

$W^{(1)}$

$\boldsymbol{h}^{(1)}$

$\boldsymbol{a}^{(2)}$

$W^{(2)}$

$\boldsymbol{h}^{(2)}$

$\boldsymbol{a}^{(3)}$

$W^{(3)}$

$\boldsymbol{h}^{(3)}$

$\boldsymbol{a}^{(\ell)}$

$W^{(\ell)}$

$\hat{\boldsymbol{y}}$

$\boldsymbol{a}^{(L)}$

# Locality ⇒ sparsity

global view

15

3

$\ell - 1$    $\ell$    $\ell + 1$

9

3

3 RF

3 RF
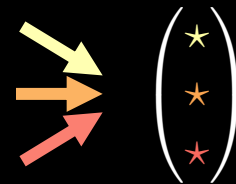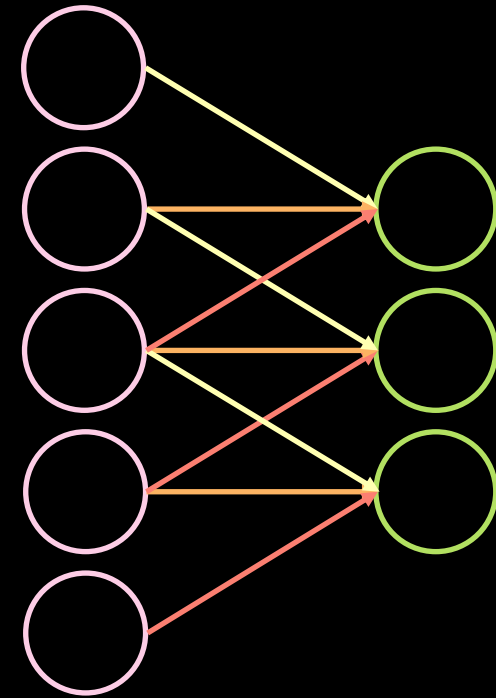
5 RF

# Stationarity ⇒ parameters sharing



**Parameters sharing**
- faster convergence
- better generalisation
- not constrained to input size
- kernel independence
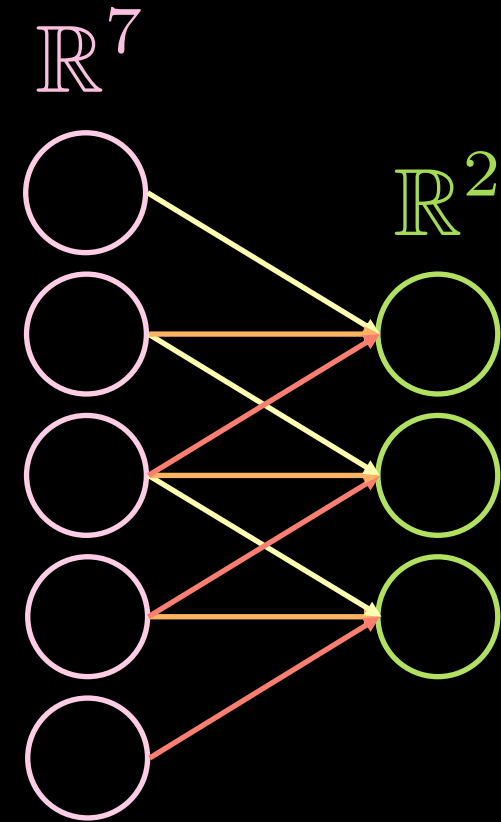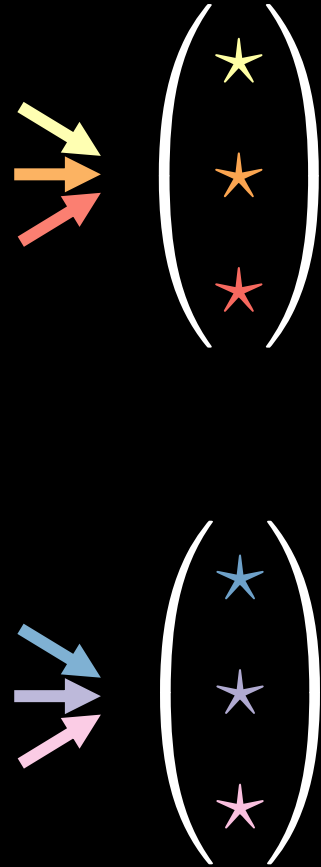  ⇒ high parallelisation
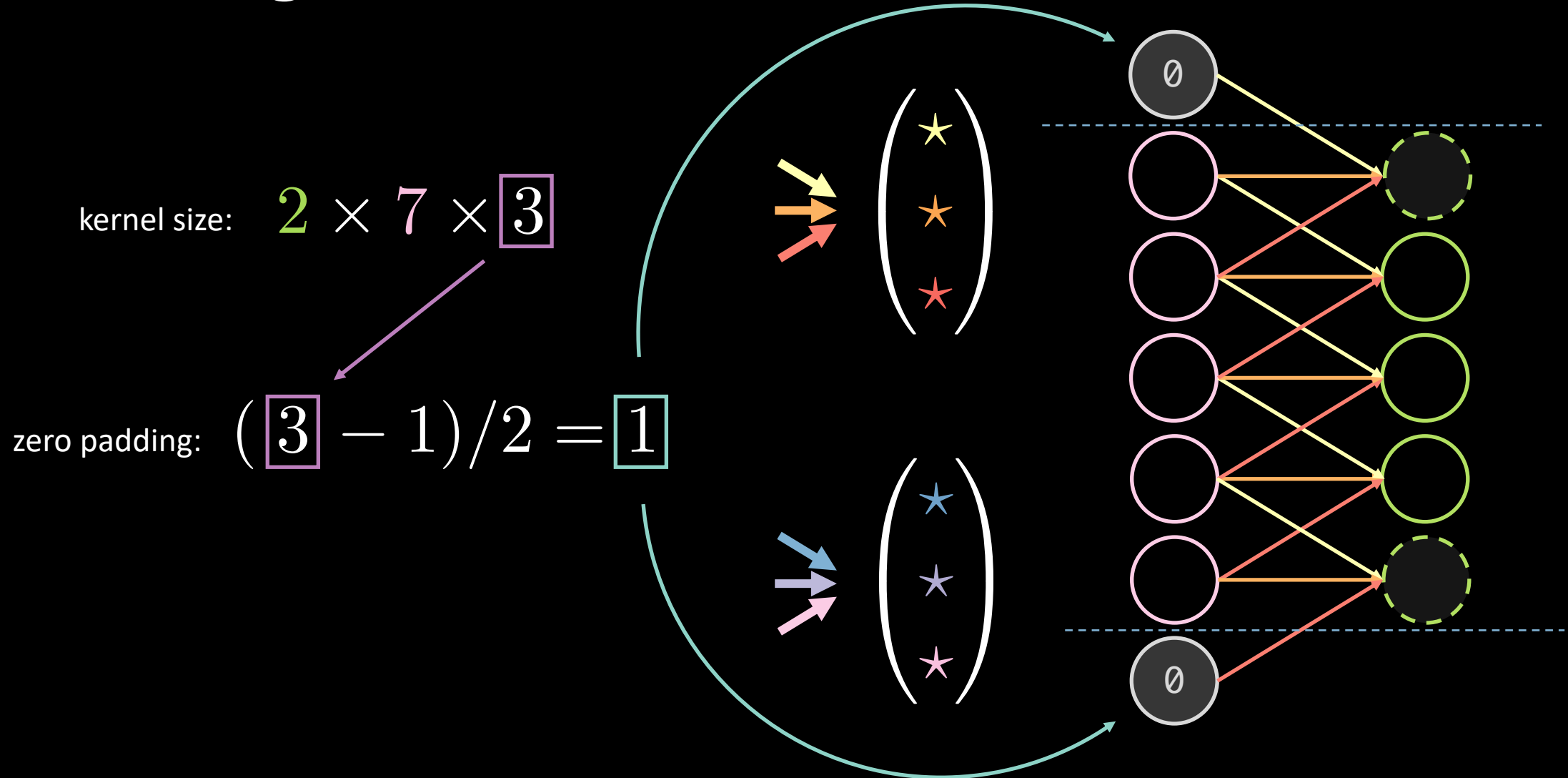
**Connection sparsity**
- reduced amount of computation

# Kernels − 1D data

kernel size: $2 \times 7 \times 3$

1D data uses 3D kernels!

$\mathbb{R}^7$

$\mathbb{R}^2$

# Padding – 1D data

kernel size: $2 \times 7 \times \boxed{3}$

zero padding: $(\boxed{3} - 1)/2 = \boxed{1}$
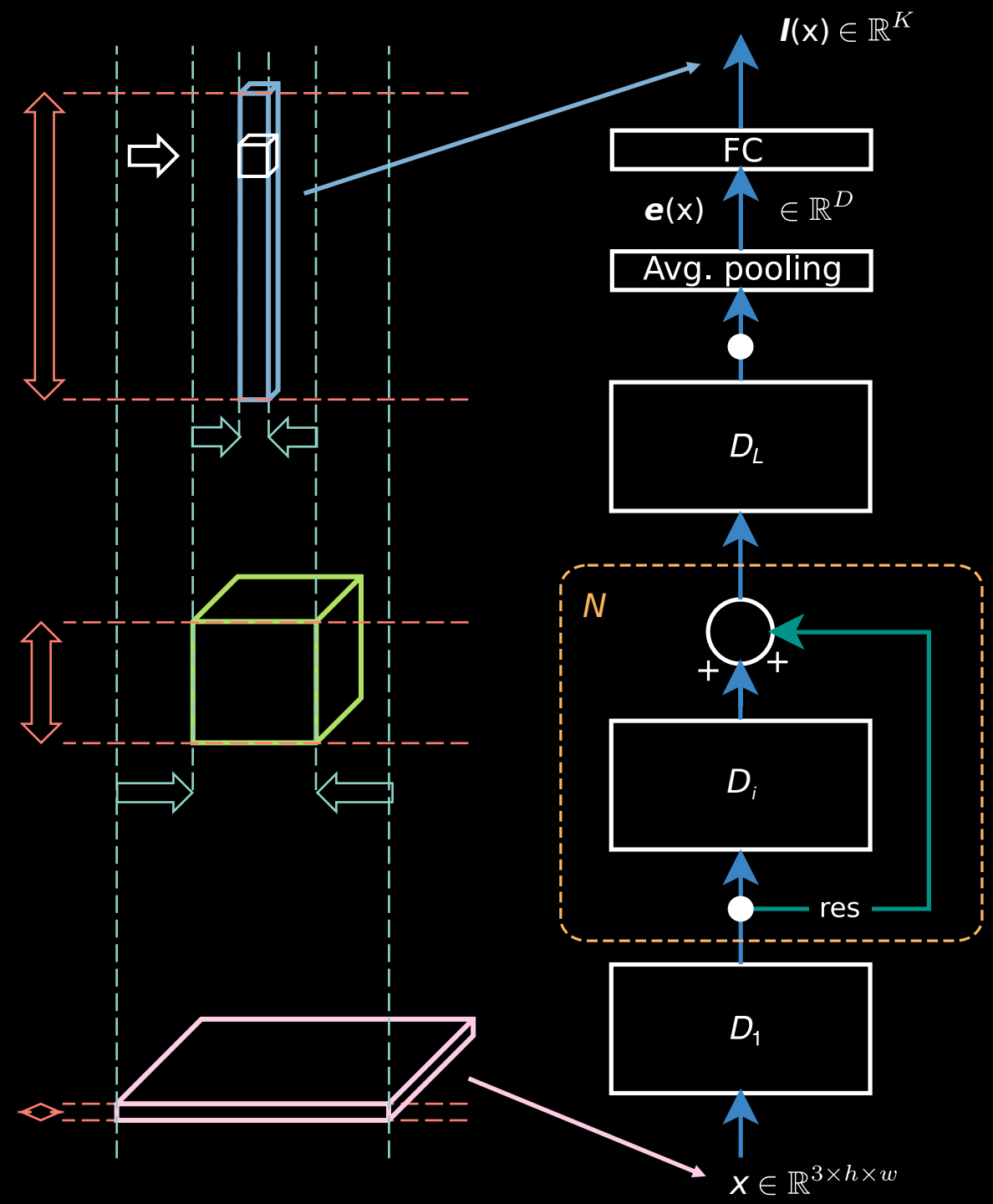
# Standard spatial CNN

- Multiple layers
  - Convolution
  - Non-linearity
  - Pooling
  - Batch normalisation

- Residual bypass connection

# Pooling

$$\|x\|_p := \left( \sum_i |x_i|^p \right)^{1/p}$$

$$\|x\|_p \to max(x), \; p \to +\infty$$

$L_p$-norm