# Foundations of Deep Learning

Alfredo Canziani

@alfcnz

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

$$\boldsymbol{x} = [$$

$x_t$ ar

$x_1$

e

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

$x_t$ are one-hot vectors

"John picked up the apple"

# 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



"John picked up the apple"

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

$x_{ij}$ are pixel values

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

$x_t$ are one-hot vectors

# Fully connected (FC) layer

$$\hat{y}$$

$$h = f(W_h x + b_h)$$

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

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

$$W_y$$

$$h$$

$$f$$

$$W_h$$

$$x$$

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



$$x$$
$$a^{(1)}$$

$$W^{(1)}$$

$$h^{(1)}$$
$$a^{(2)}$$

$$W^{(2)}$$

$$h^{(2)}$$
$$a^{(3)}$$

$$W^{(3)}$$

$$h^{(3)}$$
$$a^{(\ell)}$$

$$W^{(\ell)}$$
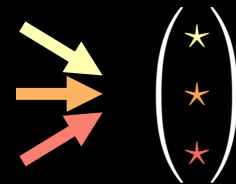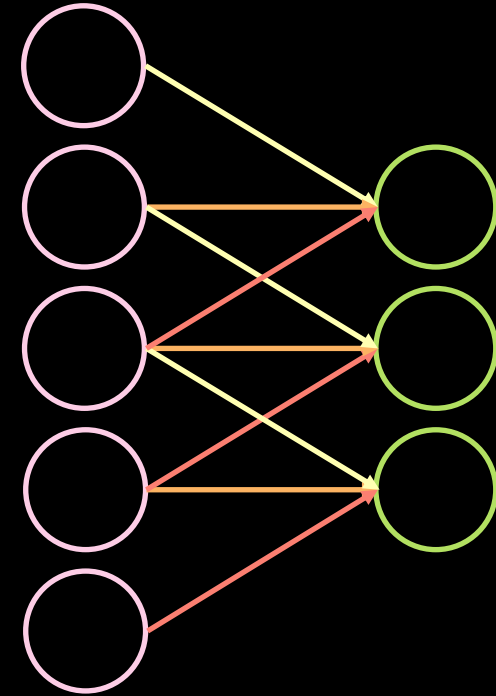
$$\hat{y}$$
$$a^{(L)}$$
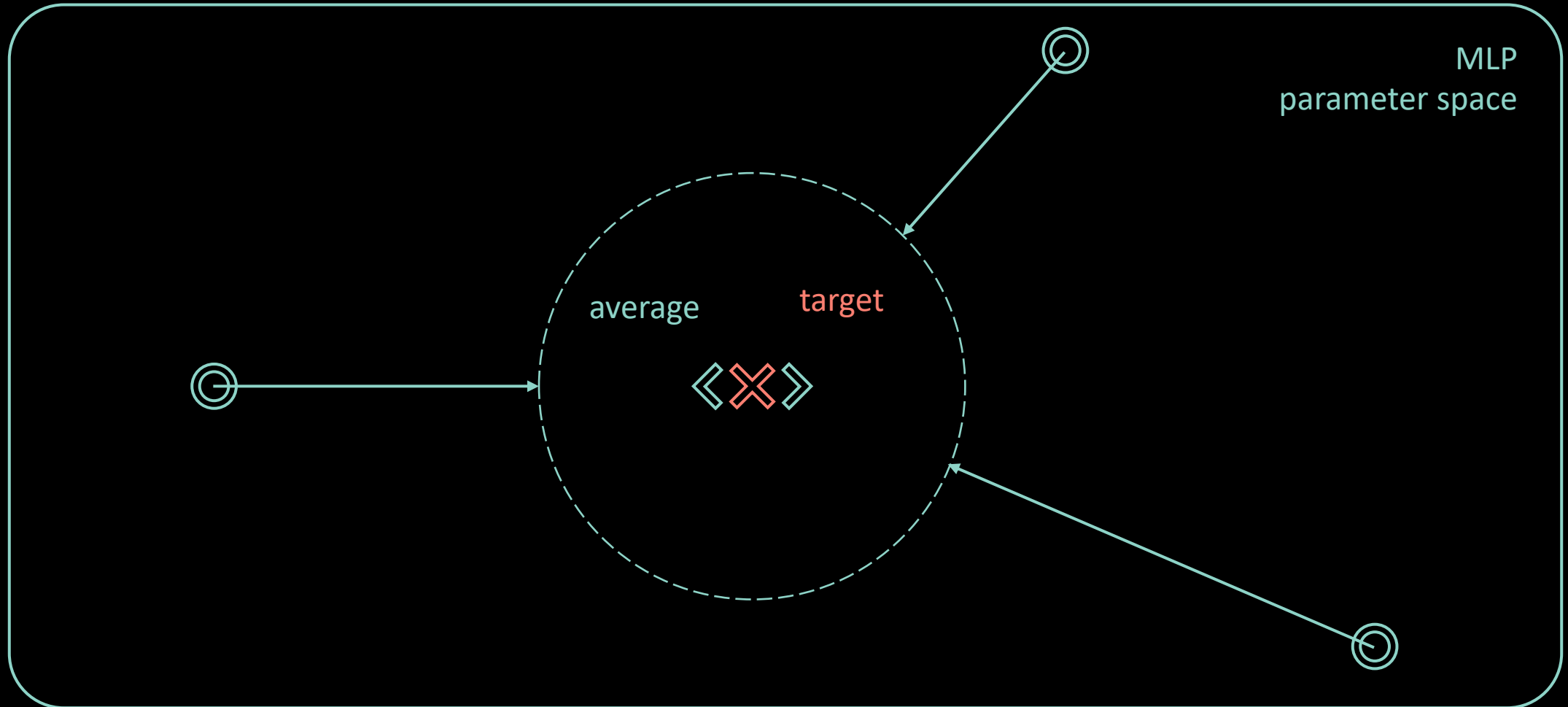
# Stationarity ⇒ parameters sharing

**Parameters sharing**
- faster convergence
- better generalisation
- not constrained to input size
- kernel independence
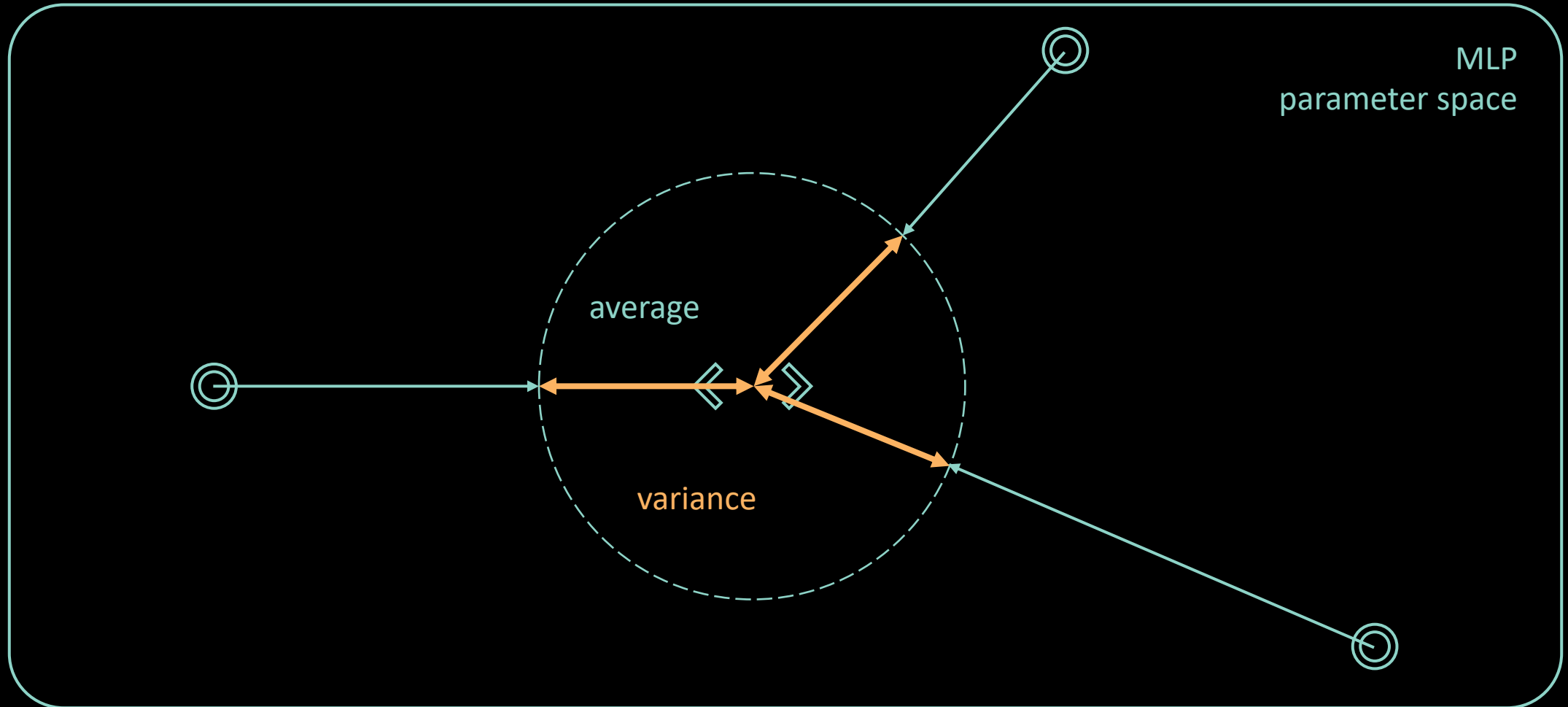  ⇒ high parallelisation

**Connection sparsity**
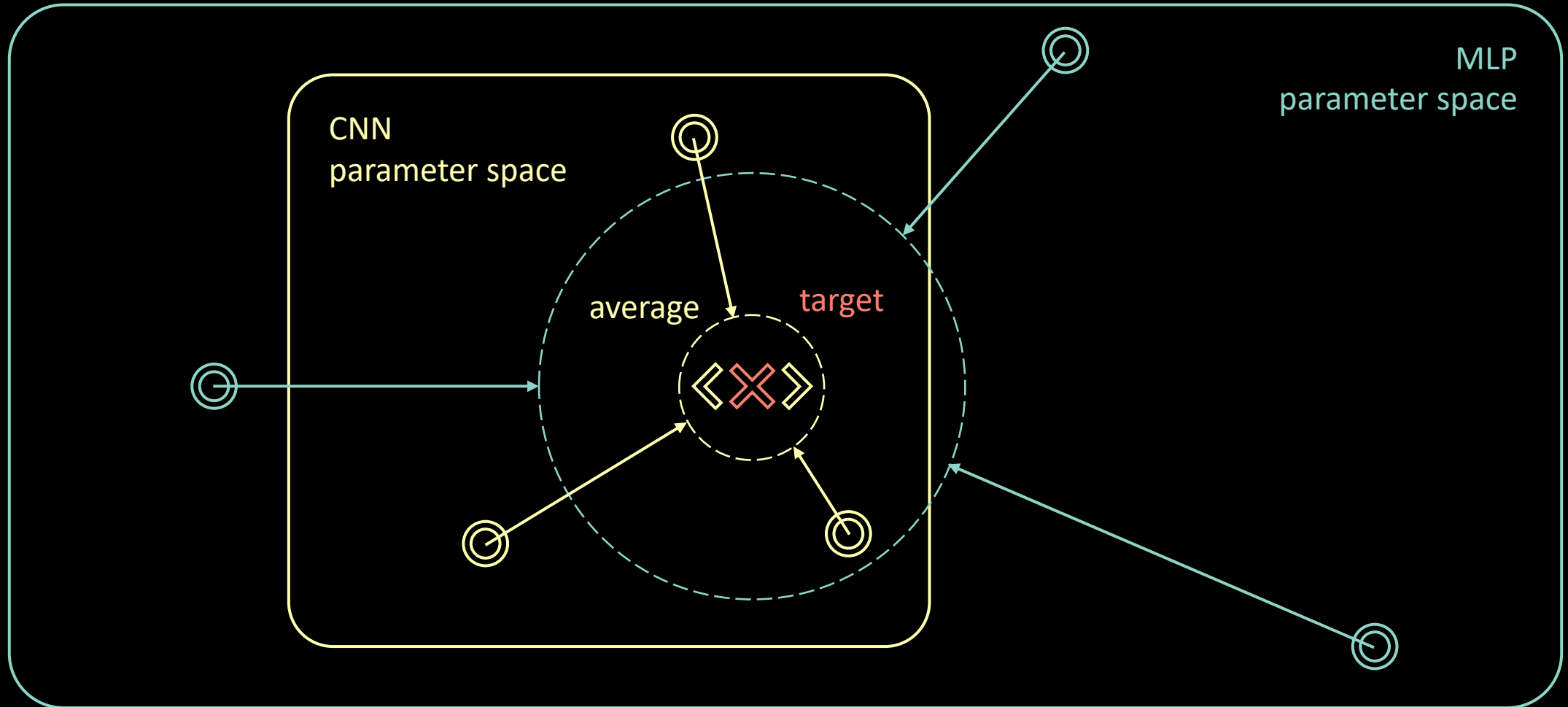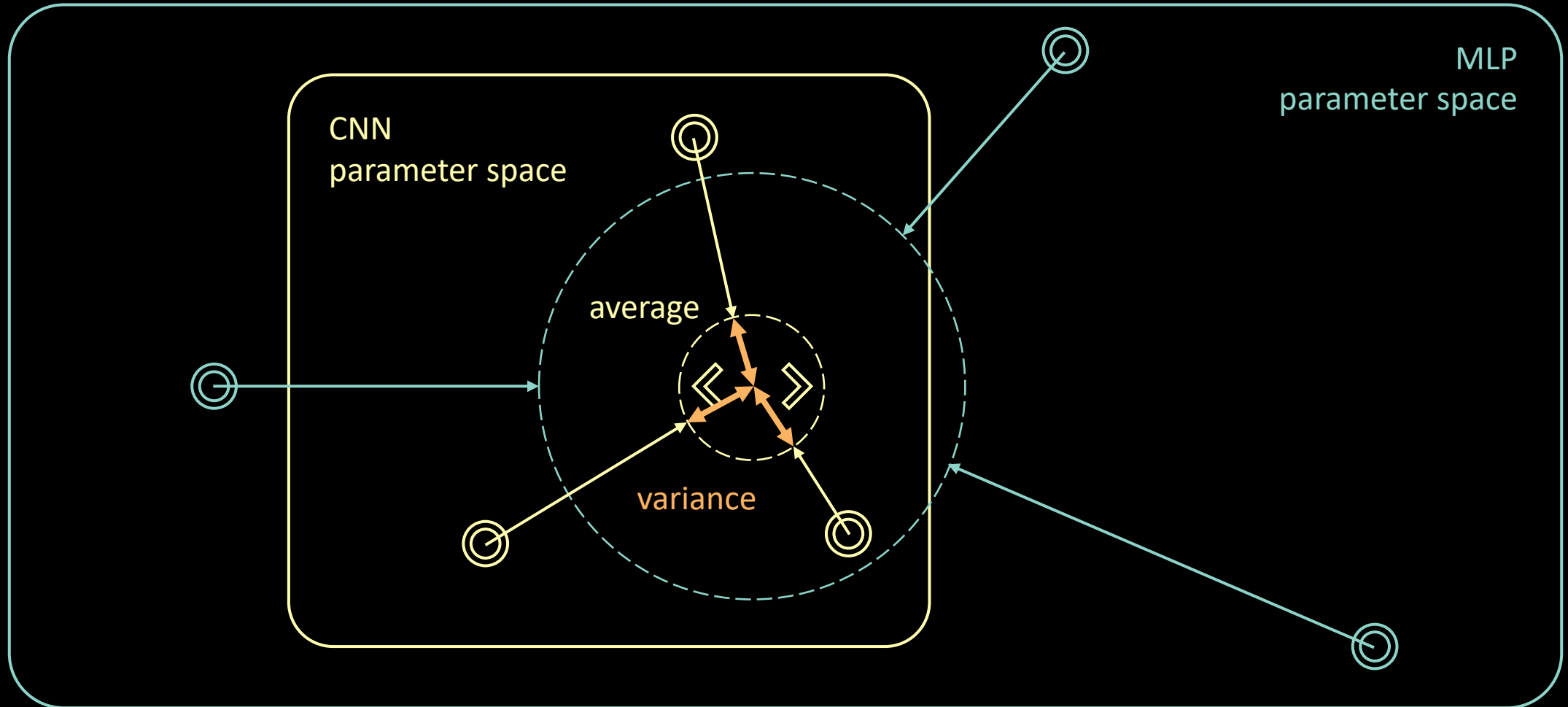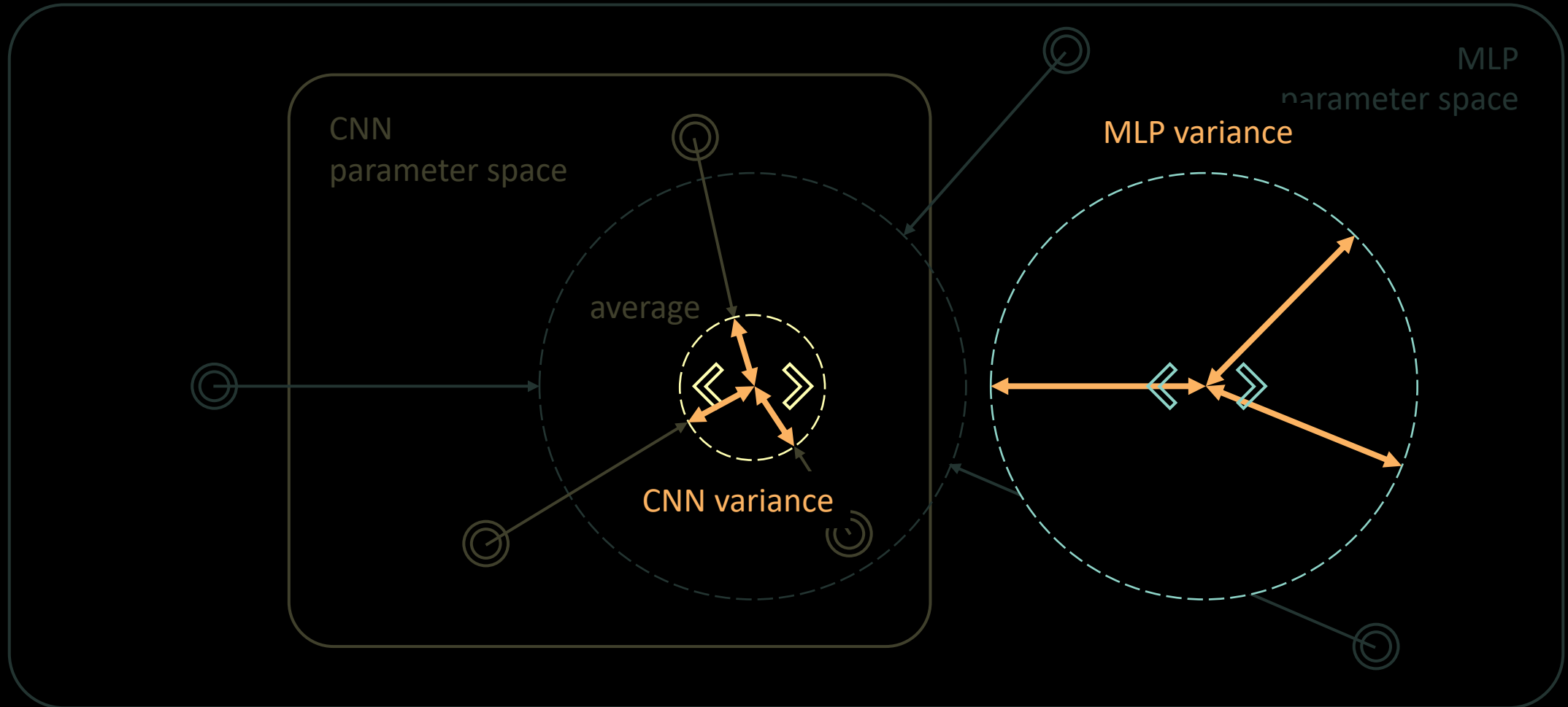- reduced amount of computation

# Generalisation error reduction

# Generalisation error reduction

# Generalisation error reduction



CNN parameter space

MLP parameter space

average  target

# Generalisation error reduction



MLP
parameter space

CNN
parameter space

average

variance

# Generalisation error reduction

CNN
parameter space

MLP
parameter space

MLP variance

average

CNN variance
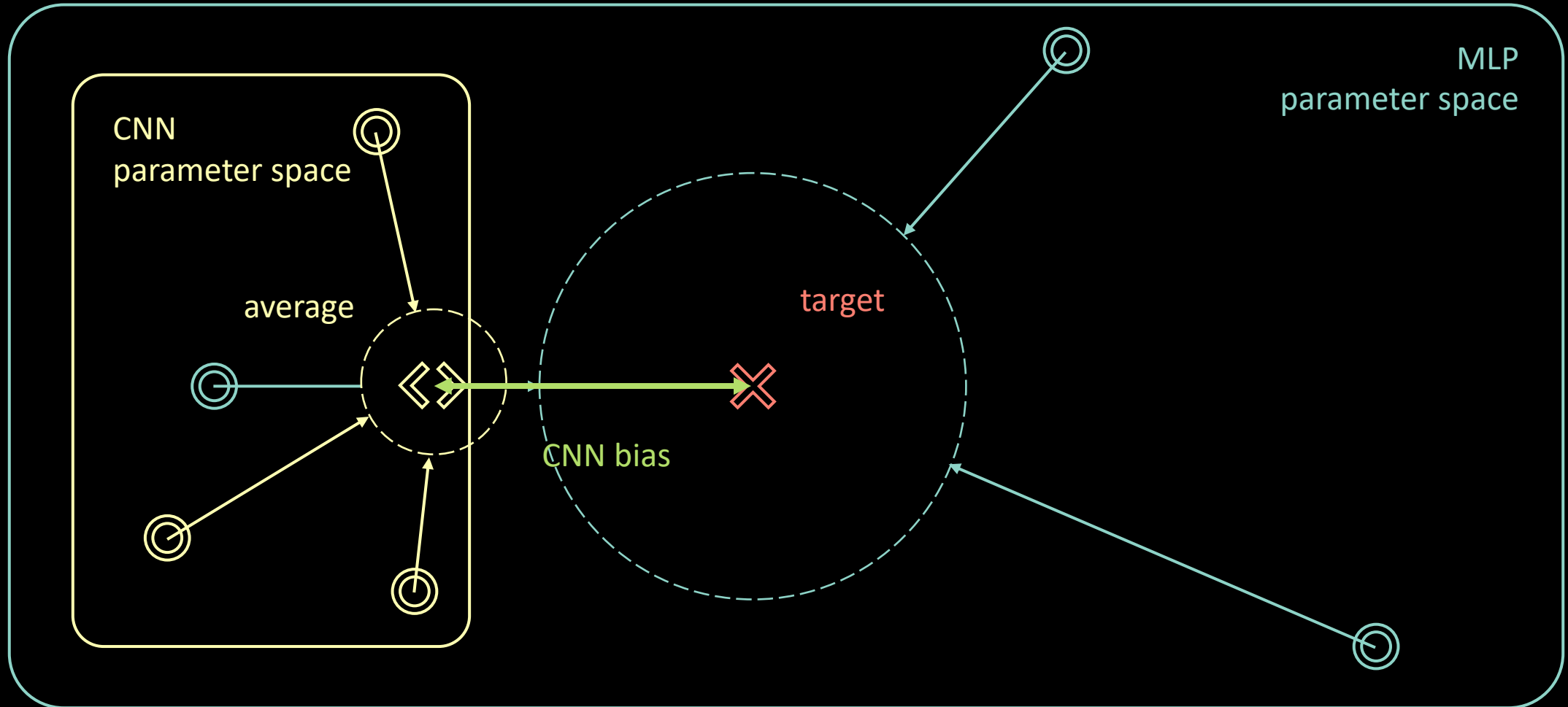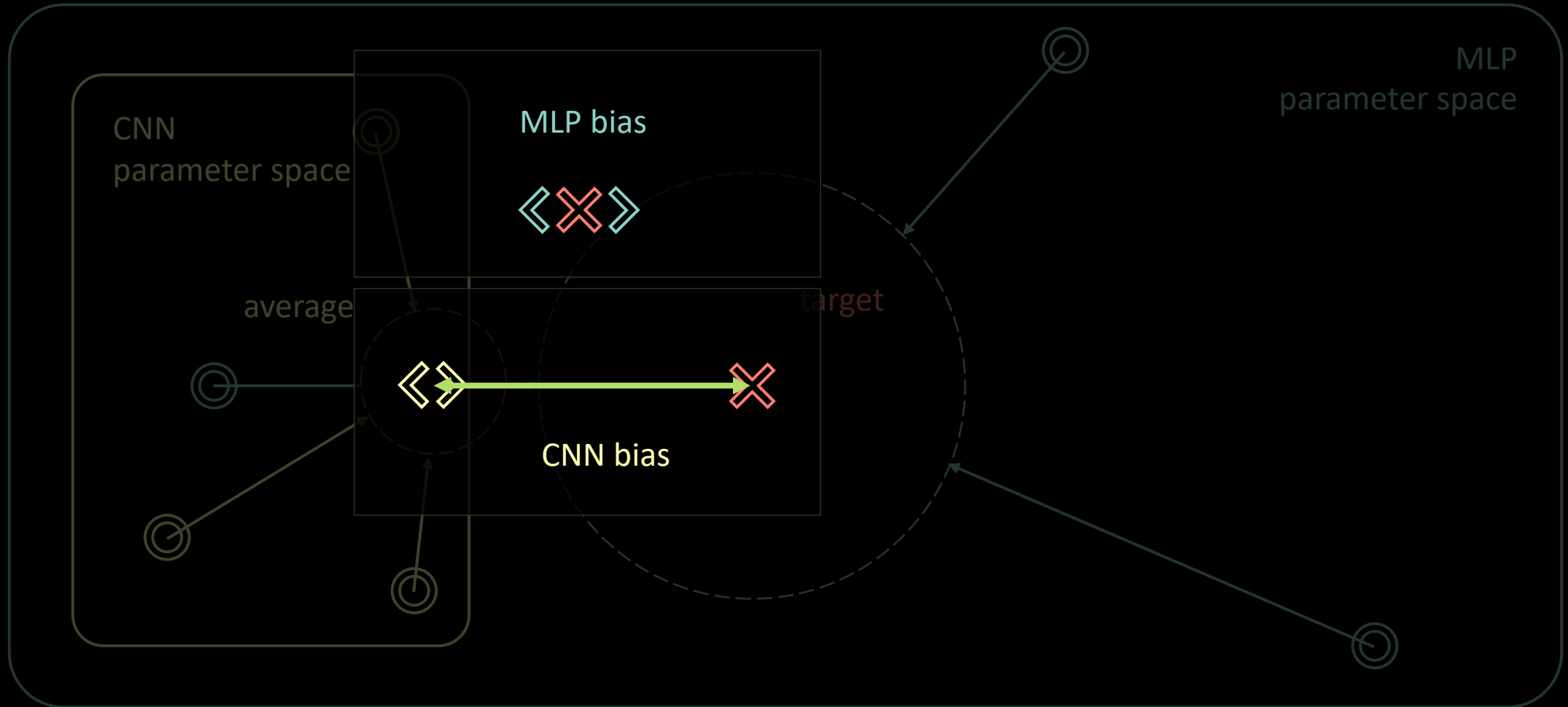
# Misspecification of model constraints

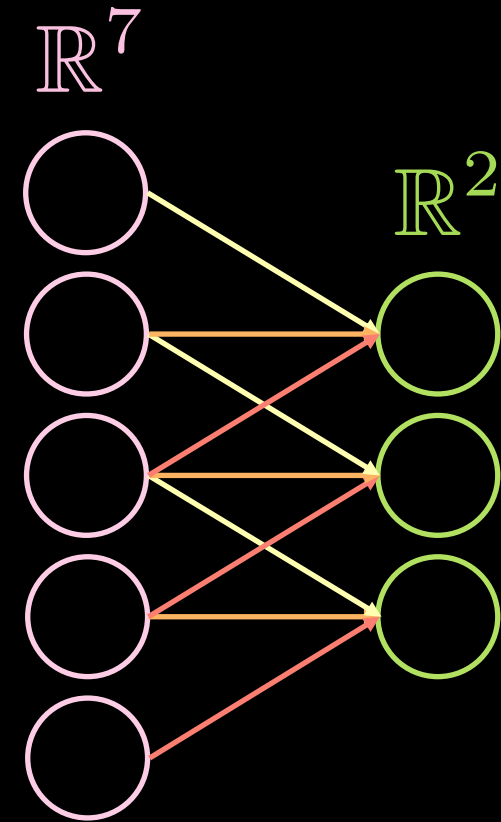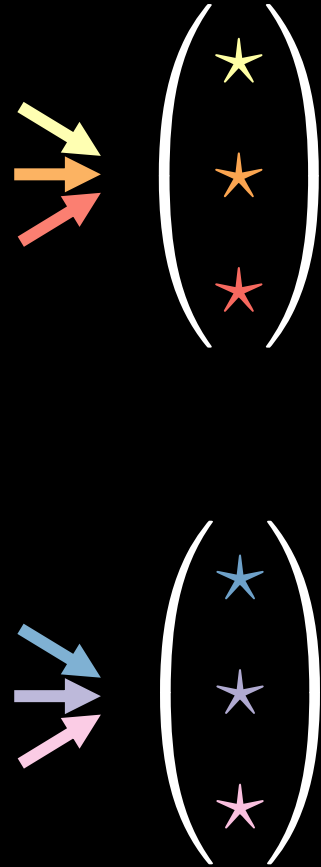# Misspecification of model constraints

# Misspecification of model constraints

# Kernels – 1D data

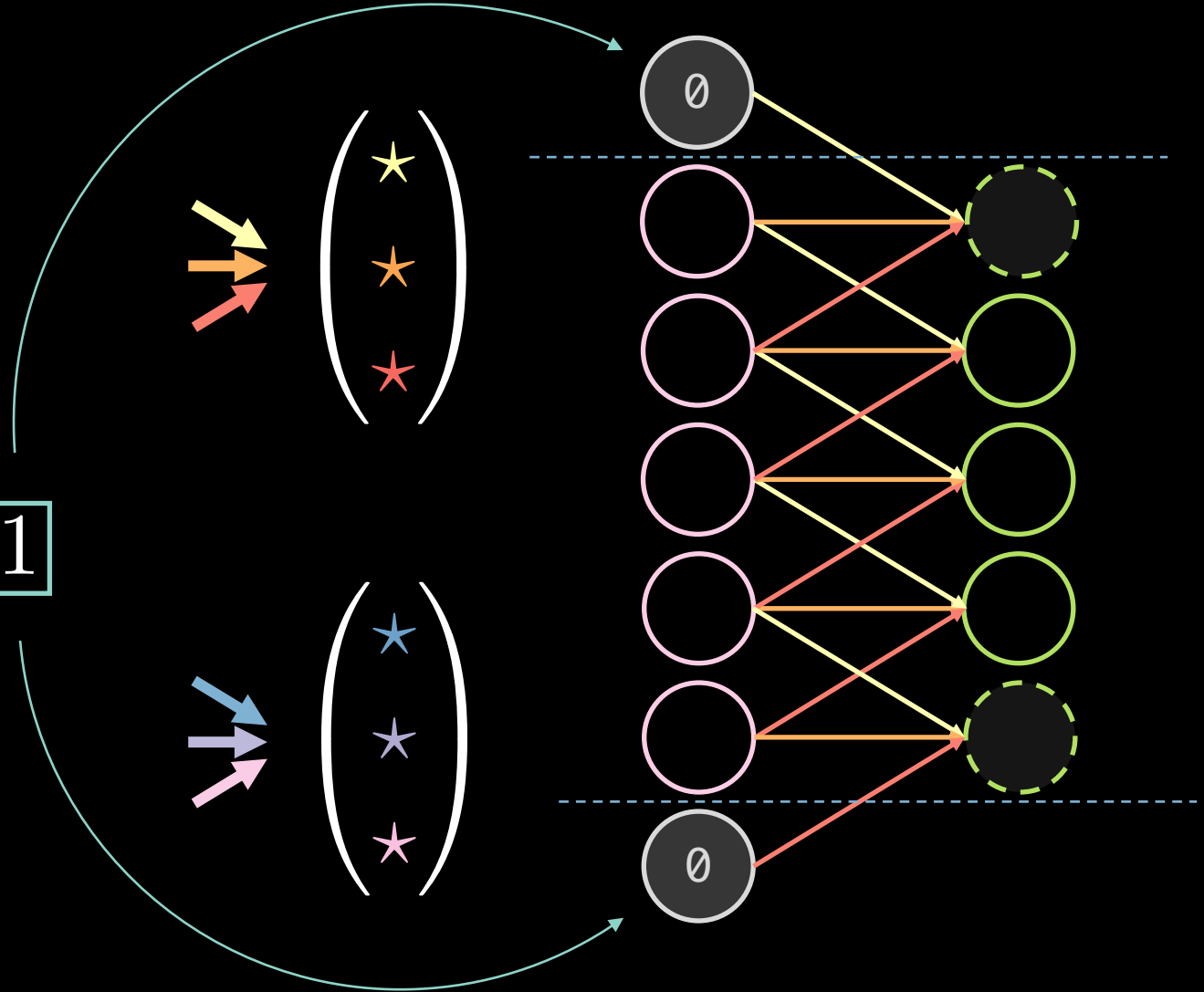kernel size: $2 \times 7 \times 3$

1D data uses 3D kernels-collection!
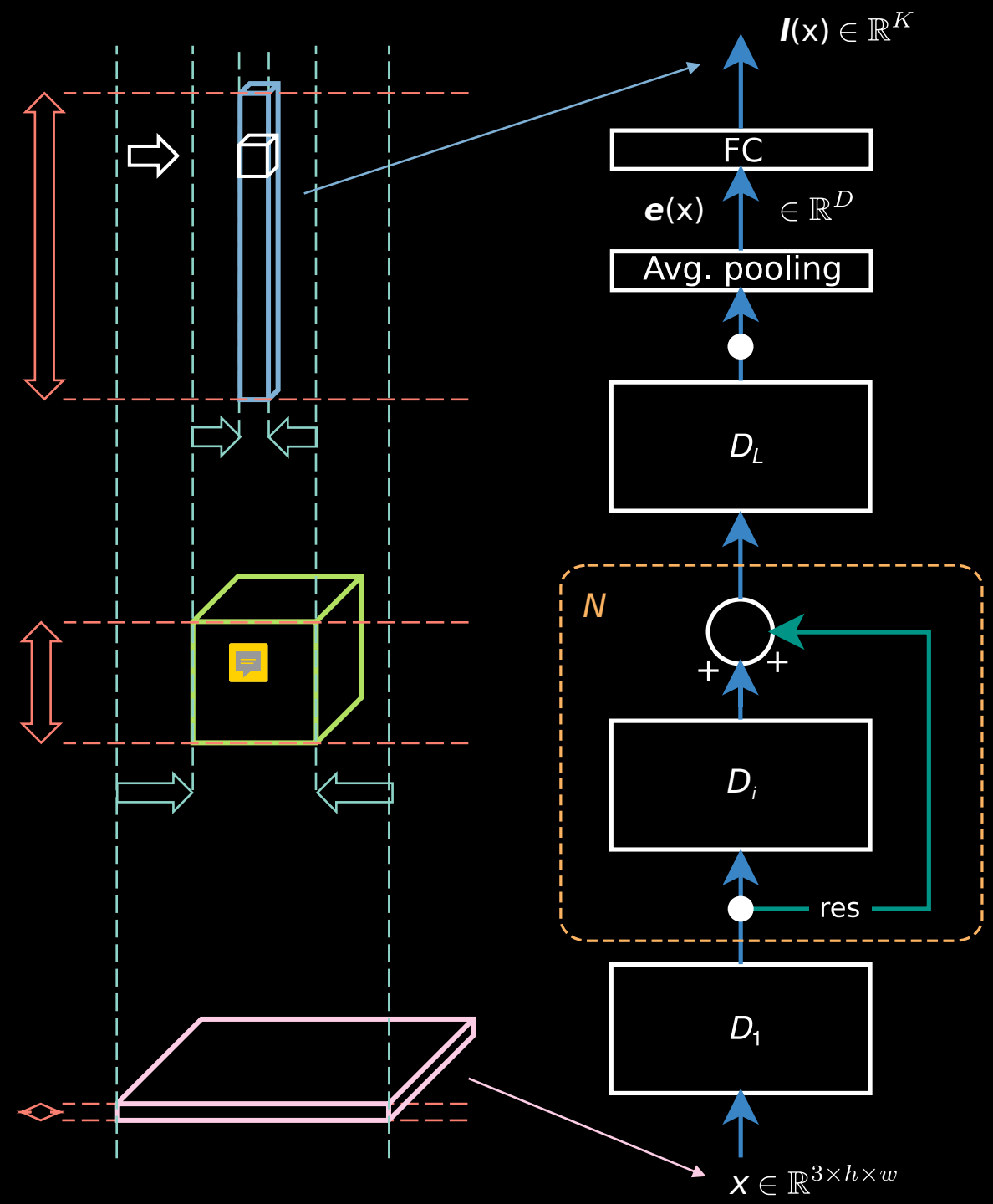
# Padding − 1D data

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

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

# Standard spatial CNN

- Multiple layers
  - Convolution
  - Non-linearity (ReLU and Leaky)
  - Pooling
  - Batch normalisation

- Residual bypass connection

# Pooling

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

$$\|x\|_p \longrightarrow max(x), \ p \rightarrow +\infty$$

$L_p$-norm

$n$

$m$

$c$

$n/2$

$\frac{m}{2}$

$c$