

Convolutional Neural Nets

Exploiting stationarity, locality, and compositionality of natural data

Input layer / samples

$$\mathcal{X} = \{\mathbf{x}^{(p)} \in \mathbb{R}^n \mid \mathbf{x}^{(p)} \text{ is a data sample}\}_{p=1}^P \quad \text{input samples}$$

$$\mathcal{X} = \{ \mathbf{x}^{(p)} : \overset{\text{domain}}{\Omega} \rightarrow \overset{\text{channels}}{\mathbb{R}^c}, \omega \mapsto \mathbf{x}^{(p)}(\omega) \}_{p=1}^P$$

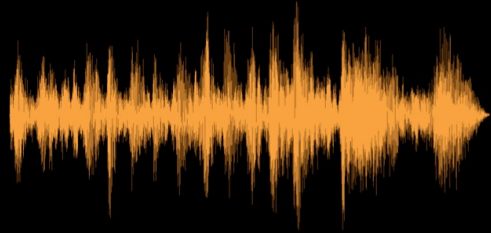
$$\Omega = \{1, 2, \dots, \overset{\text{total time}}{T/\Delta t}\} \subset \mathbb{N}, \quad c \in \{1, 2, \overset{\text{stereo}}{5+1}, \dots\}$$

sampling interval mono Dolby 5.1

$$\Omega = \{1, \dots, \overset{\text{height}}{h}\} \times \{1, \dots, \overset{\text{width}}{w}\} \subset \mathbb{N}^2, \quad c \in \{1, \overset{\text{grey scale}}{3}, \overset{\text{colour}}{20}, \dots\} \quad \text{hyperspectral}$$

$$\Omega = \underset{\text{space-time}}{\mathbb{R}^4} \times \overset{\text{four-momentum}}{\mathbb{R}^4}, \quad \underset{\text{Hamiltonian}}{c} = 1 \quad \quad \quad \boldsymbol{x}(\omega_1, \omega_2) = \begin{pmatrix} r(\omega_1, \omega_2) \\ g(\omega_1, \omega_2) \\ b(\omega_1, \omega_2) \end{pmatrix}$$

Signals can be represented as vectors



$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_t \ \dots]^\top$$

x_t are waveform heights



$$\mathbf{x} = [x_{11} \ x_{12} \ \dots \ x_{1n} \ x_{21} \ x_{22} \ \dots]^\top$$

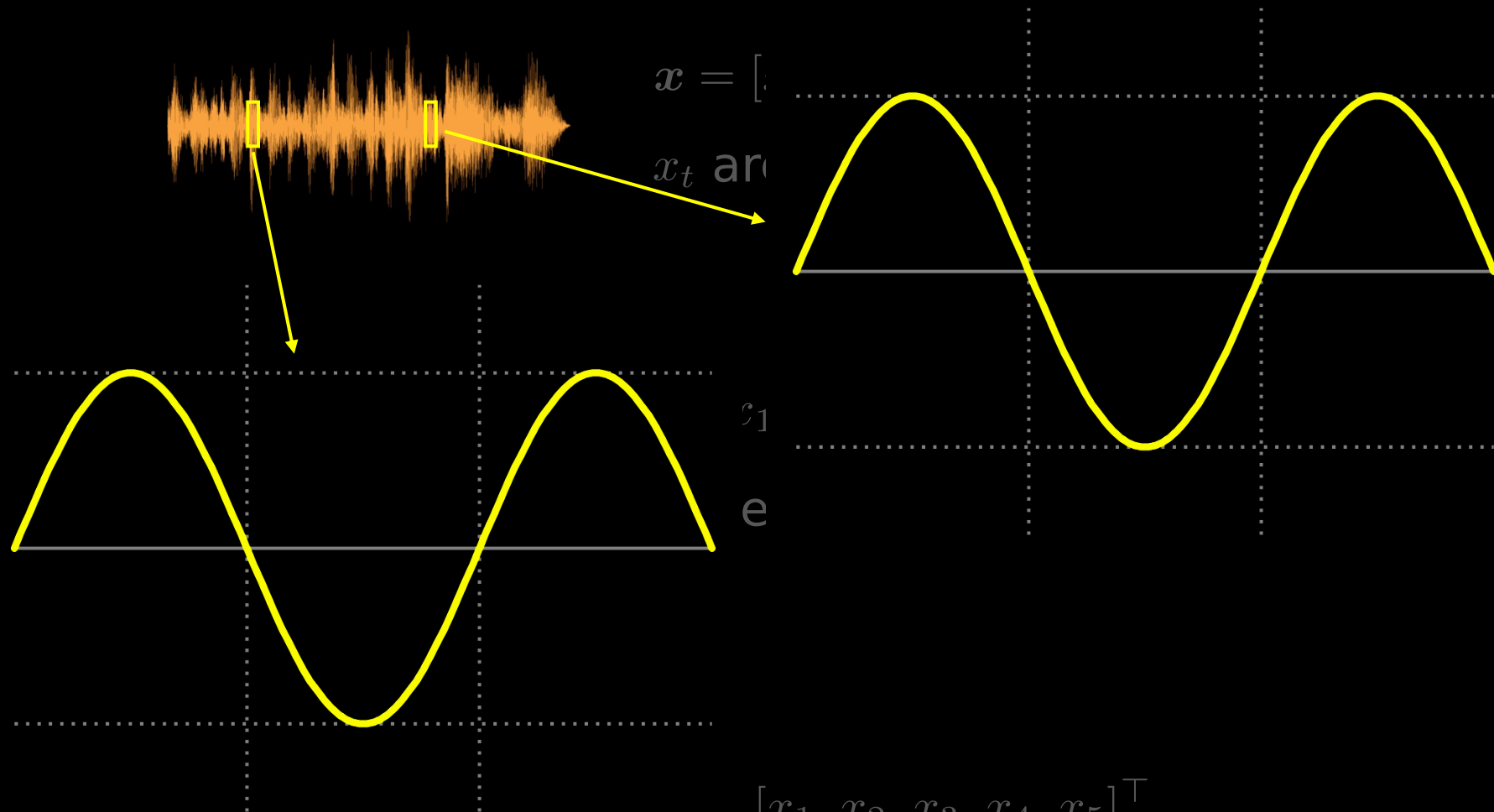
x_{ij} are pixel values

“John picked up the apple”

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

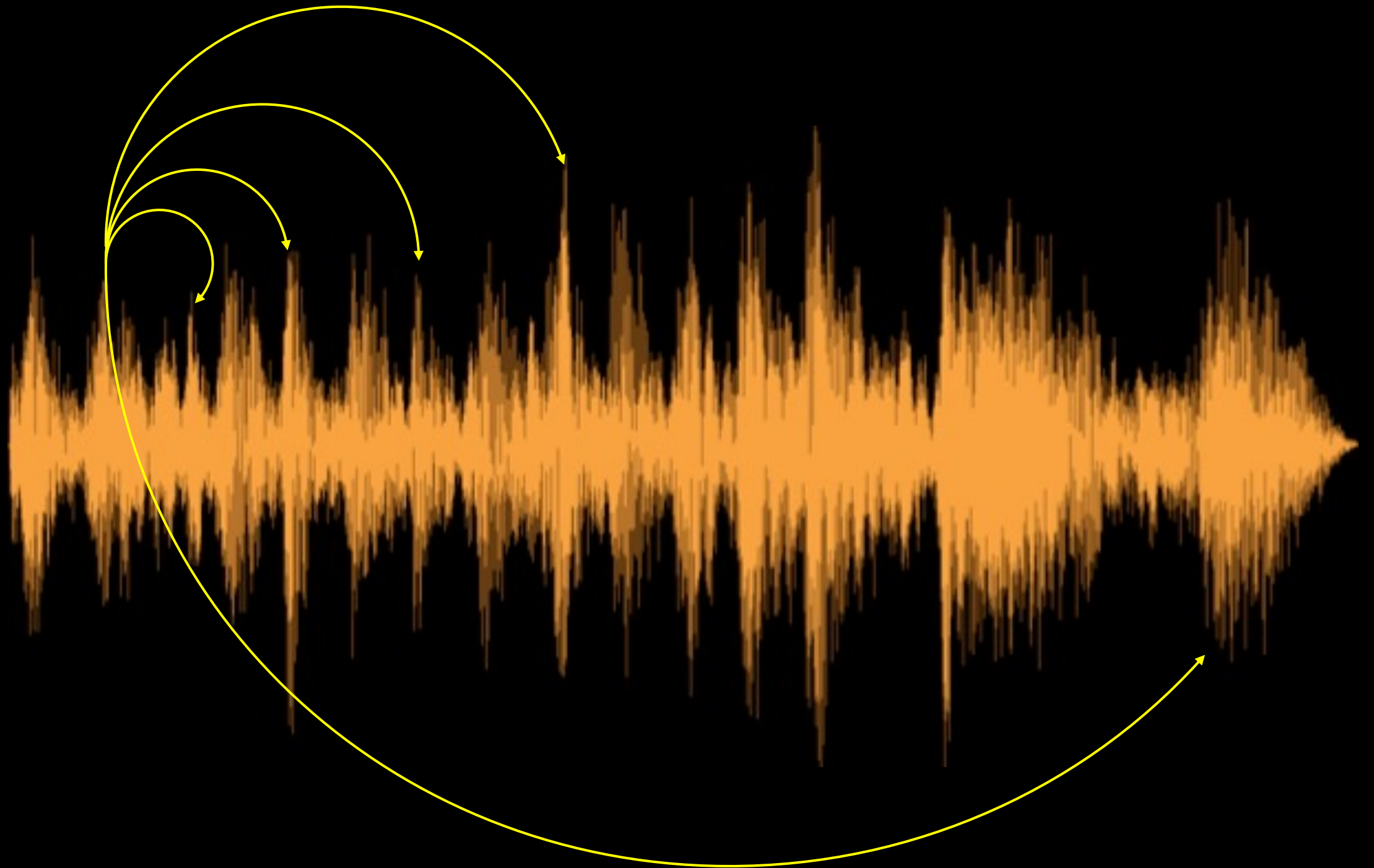
x_t are one-hot vectors

Signals can be represented as vectors

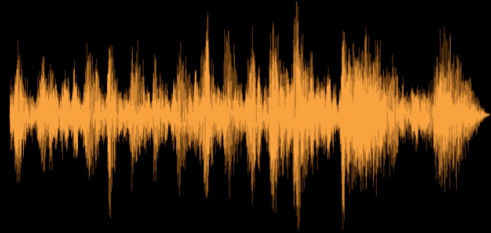


“John picked up the apple”

x_t are one-hot vectors



Signals can be represented as vectors



$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_t \ \dots]^\top$$

x_t are waveform heights



$$\mathbf{x} = [x_{11} \ x_{12} \ \dots \ x_{1n} \ x_{21} \ x_{22} \ \dots]^\top$$

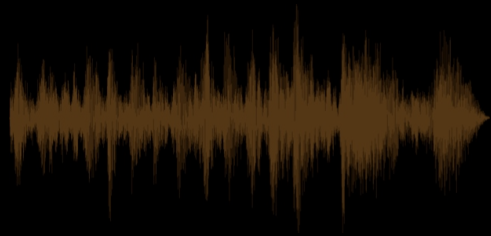
x_{ij} are pixel values

“John picked up the apple”

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

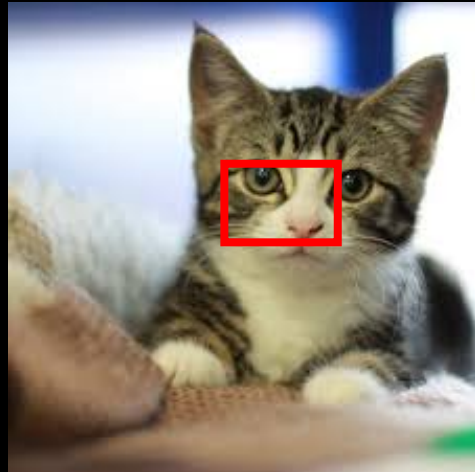
x_t are one-hot vectors

Signals can be represented as vectors



$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_t \ \dots]^\top$$

x_t are waveform heights



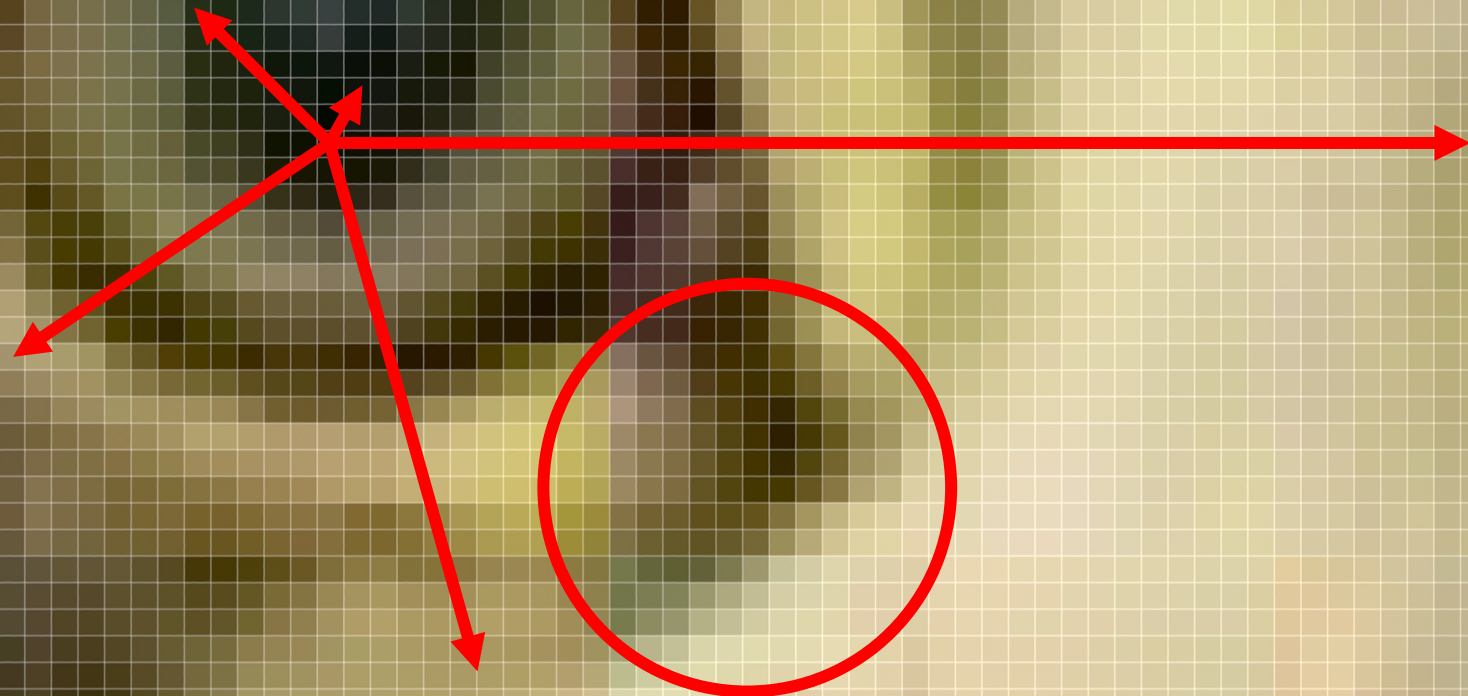
$$\mathbf{x} = [x_{11} \ x_{12} \ \dots \ x_{1n} \ x_{21} \ x_{22} \ \dots]^\top$$

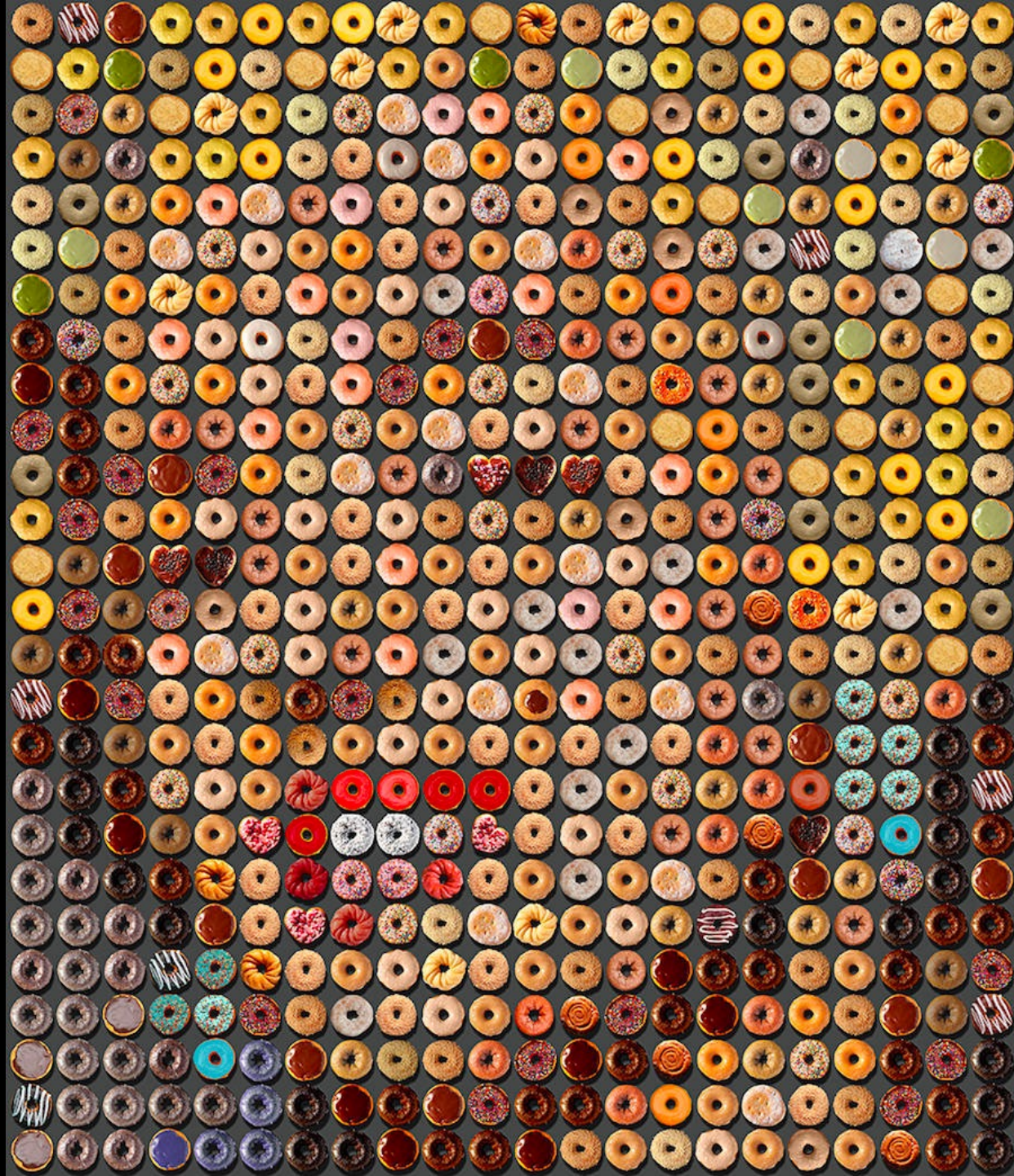
x_{ij} are pixel values

“John picked up the apple”

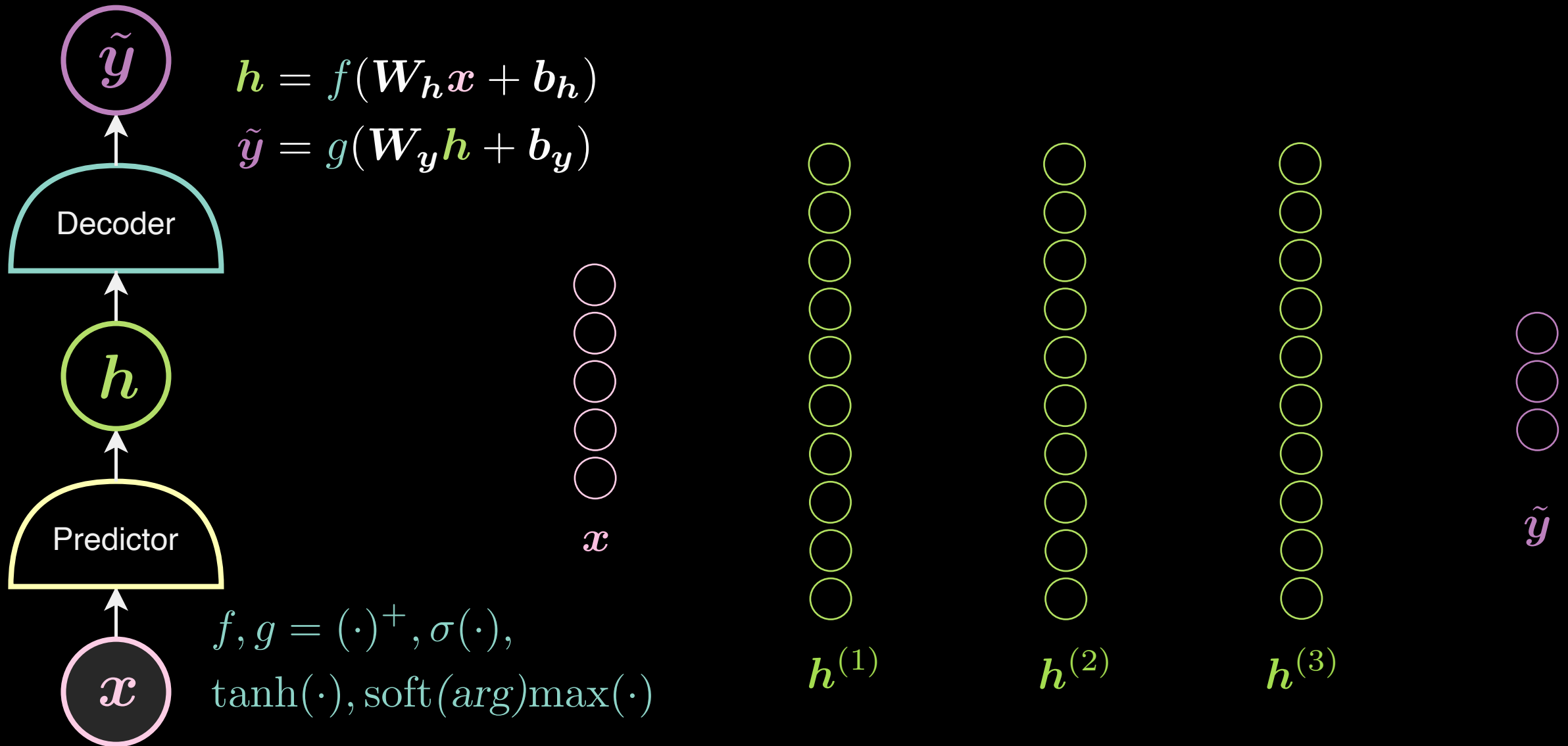
$$\mathbf{x} = [x_1 \ x_2 \ x_3 \ x_4 \ x_5]^\top$$

x_t are one-hot vectors





Fully connected (FC) layer

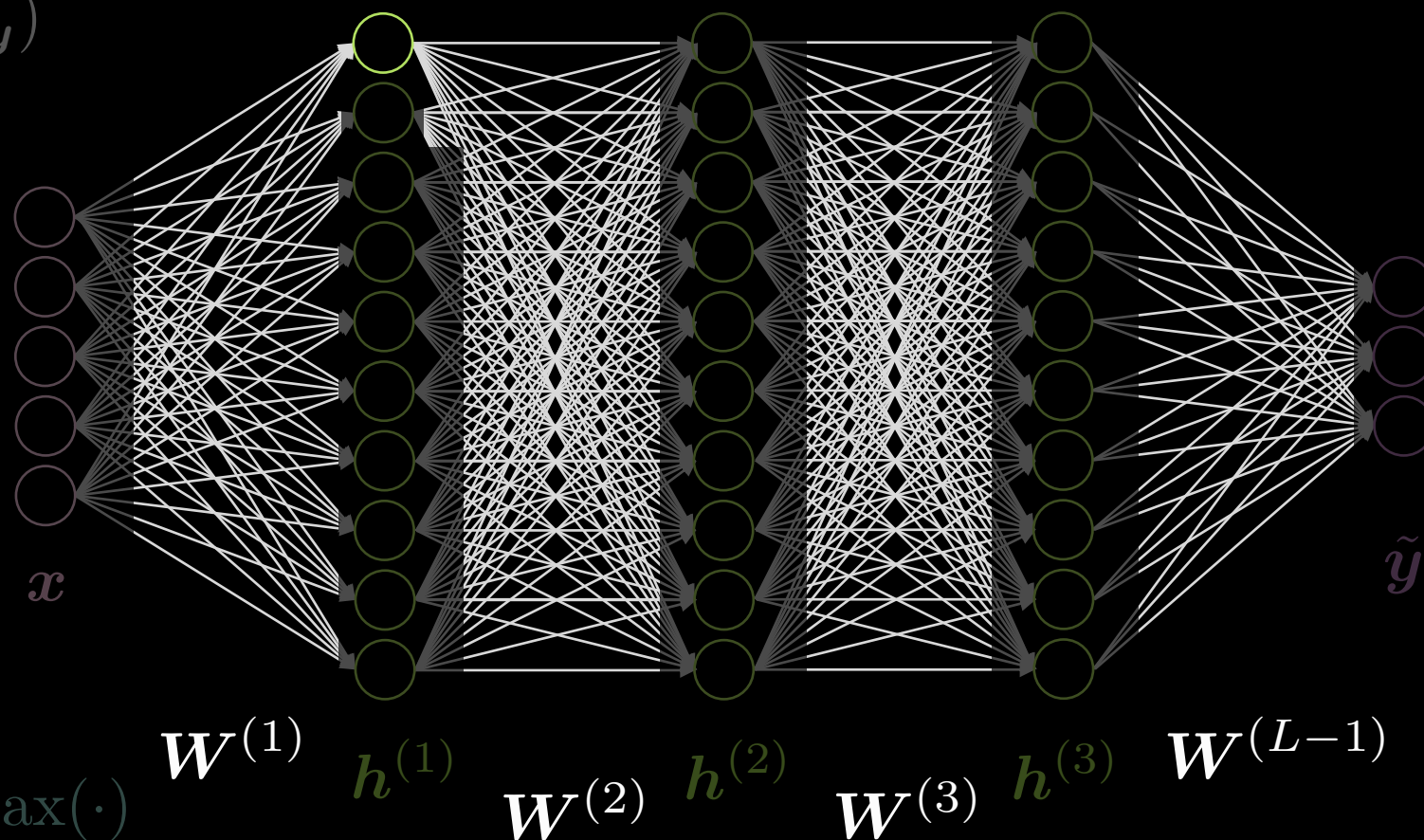
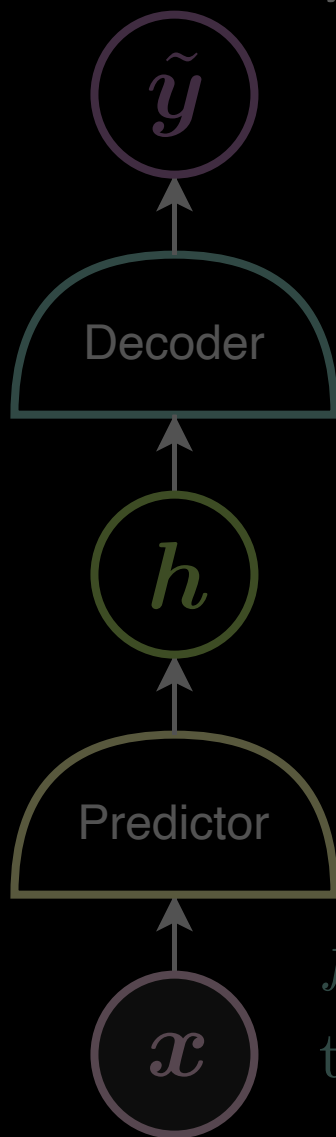


Fully connected (FC) layer

$$\mathbf{h} = f(\mathbf{W}_h \mathbf{x} + \mathbf{b}_h)$$

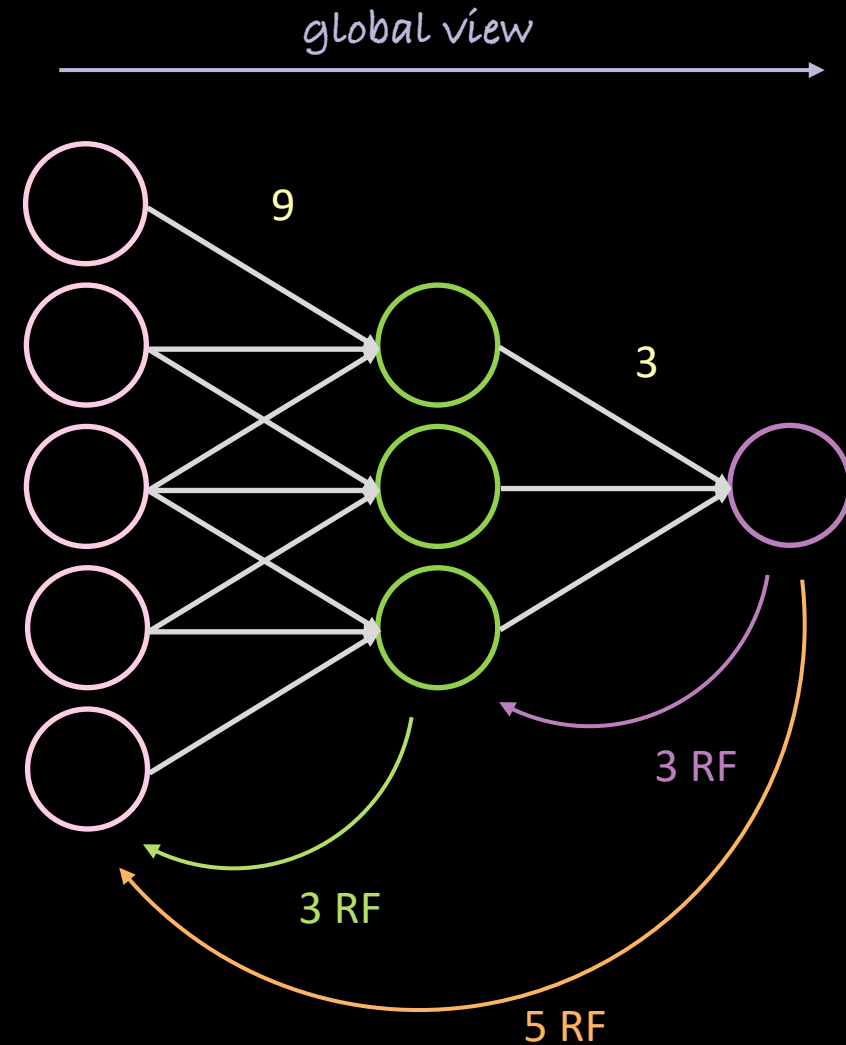
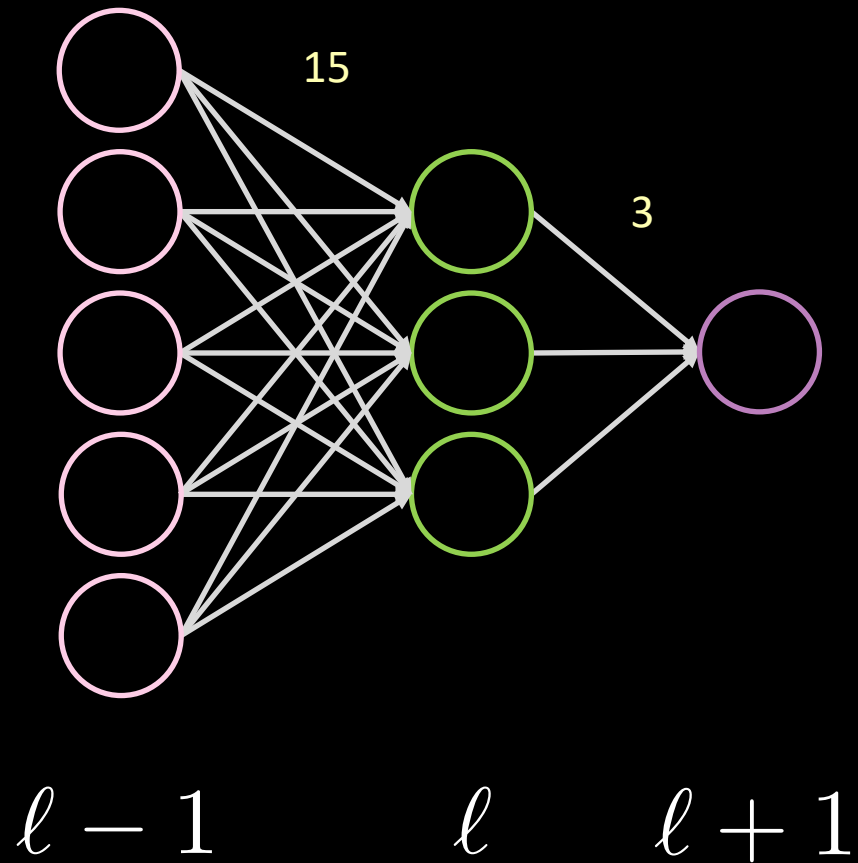
$$\tilde{\mathbf{y}} = g(\mathbf{W}_y \mathbf{h} + \mathbf{b}_y)$$

$$h_j^{(1)} = f(\boxed{w^{(j)}} \mathbf{x} + b_j) = f\left[\left(\sum_{i=1}^n w_i^{(j)} x_i\right) + b_j\right]$$

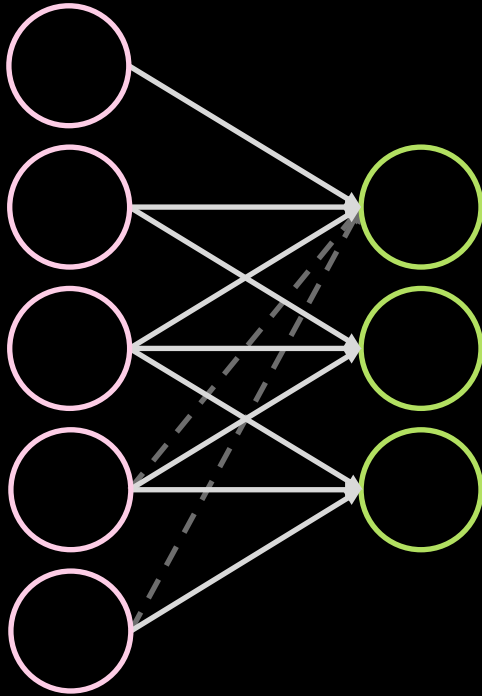


$f, g = (\cdot)^+, \sigma(\cdot), \tanh(\cdot), \text{soft}(\text{arg})\text{max}(\cdot)$

Locality \Rightarrow sparsity



Stationarity \Rightarrow parameters sharing

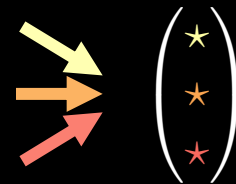
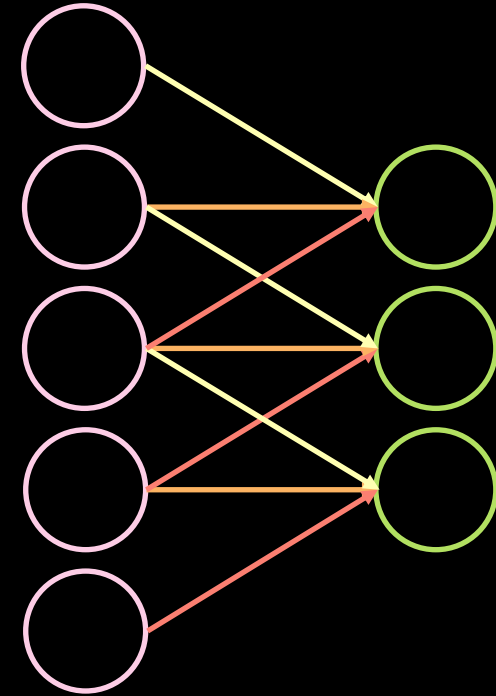


Parameters sharing

- faster convergence
- better generalisation
- not constrained to input size
- kernel independence
 \Rightarrow high parallelisation

Connection sparsity

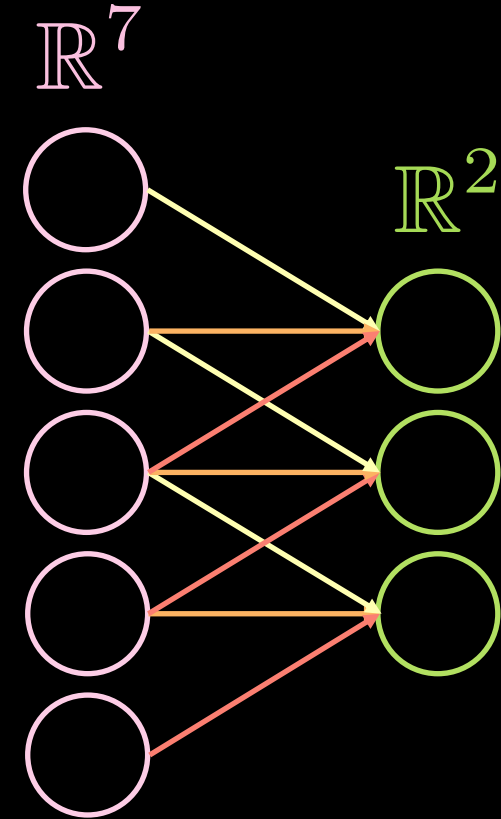
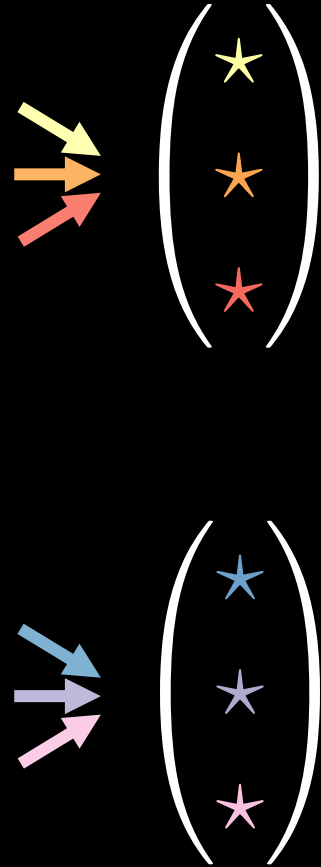
- reduced amount of computation



Kernels – 1D data

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

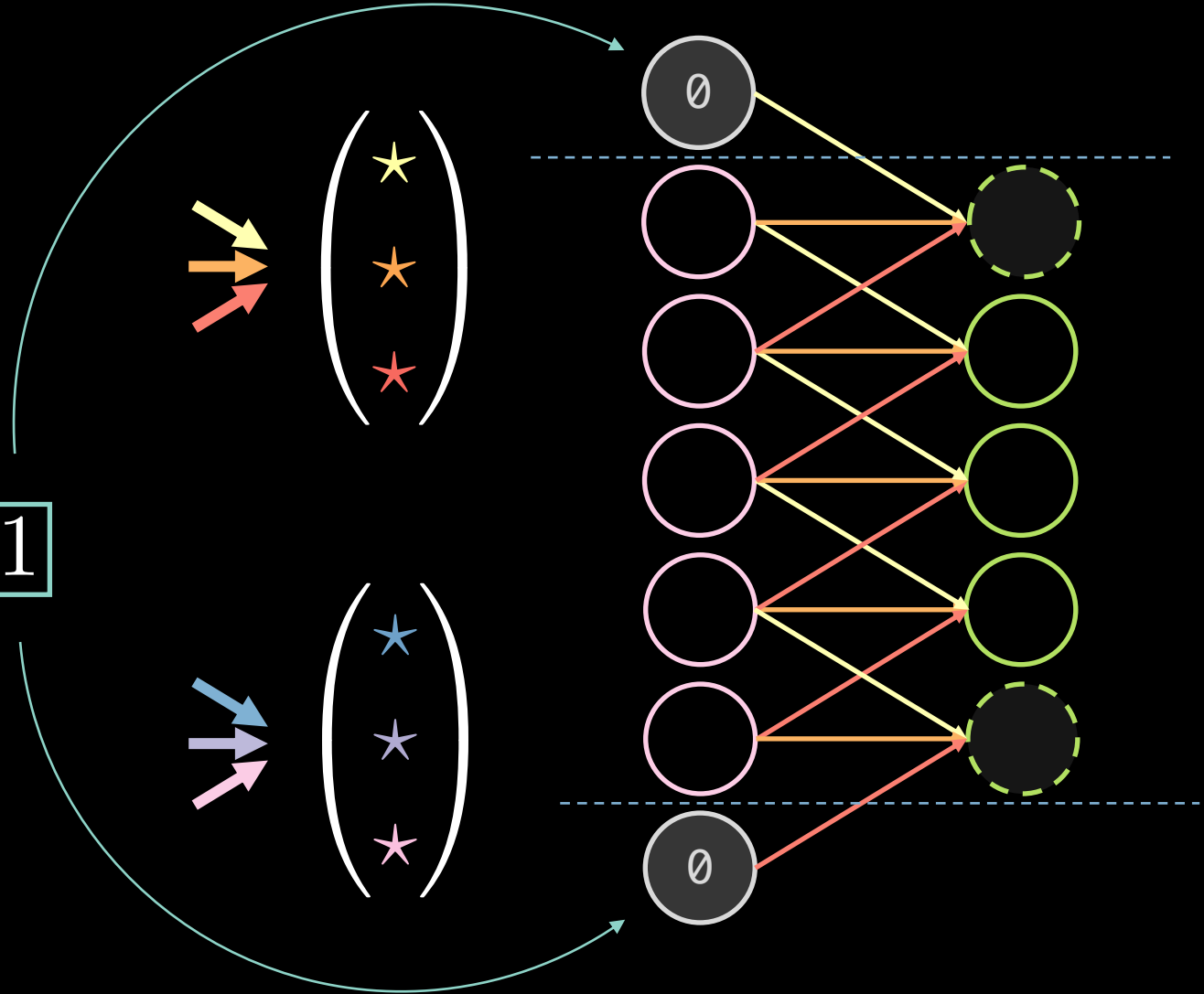
1D data uses 3D kernels-collection!



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

