

Convolutional Neural Nets

Exploiting stationarity, locality, and compositionality of natural data

Input layer / samples

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

We don't want pixels in images defined as points in an NxN dimensional space

(eg: 1 Mega Pixel = 1 million dimensional space)

because all data points are going to be really close in this huge dimensional space, and to separate them, we'll want to make this space bigger, which is computationally very expensive.

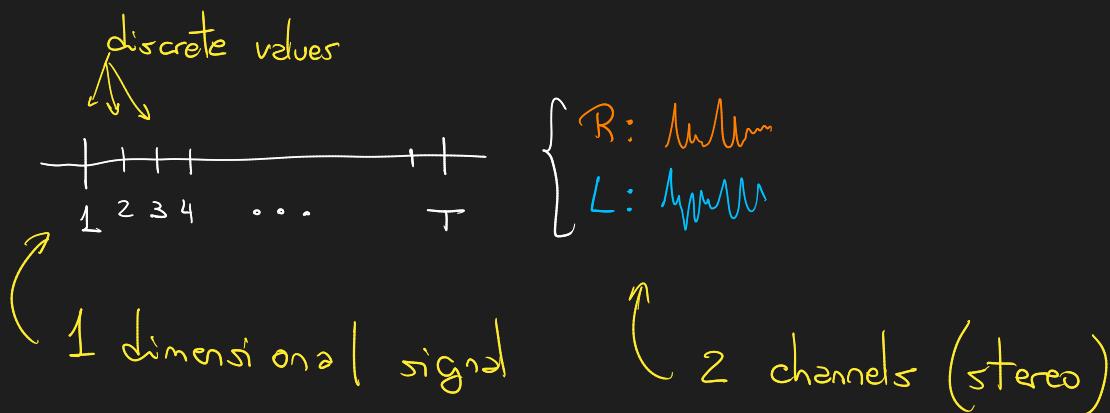
Instead, we can define the data points as functions

$$\mathcal{X} = \{\mathbf{x}^{(i)} : \Omega \xrightarrow{\text{domain}} \mathbb{R}^c, \omega \mapsto \mathbf{x}^{(i)}(\omega)\}_{i=1}^m$$

Some examples

- Audio

$$\Omega = \{1, 2, \dots, T/\Delta t\} \subset \mathbb{N}, \quad c \in \begin{cases} \text{mono} & 1 \\ \text{stereo} & 2, 5+1, \dots \\ \text{Dolby 5.1} & \end{cases}$$



- Images

$$\Omega = \{1, \dots, h\} \times \{1, \dots, w\} \subset \mathbb{N}^2, \quad c \in \begin{cases} \text{grey scale} & 1 \\ \text{colour} & 3 \\ \text{hyperspectral} & 20, \dots \end{cases}$$

ordered!

$$\mathbf{x}(\omega_1, \omega_2) = \begin{pmatrix} r(\omega_1, \omega_2) \\ g(\omega_1, \omega_2) \\ b(\omega_1, \omega_2) \end{pmatrix}$$

- Physics

$$\Omega = \underset{\text{space-time}}{\mathbb{R}^4} \times \underset{\text{four-momentum}}{\mathbb{R}^4}, \quad c = \underset{\text{Hamiltonian}}{1}$$

1D signal

$$x[k] = (x[1] \quad x[2] \quad \dots \quad)$$

Signals can be represented as vectors



$$x = [x_1 \ x_2 \ x_3 \ \dots \ x_t \ \dots]^T$$

x_t are waveform heights



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

x_{ij} are pixel values

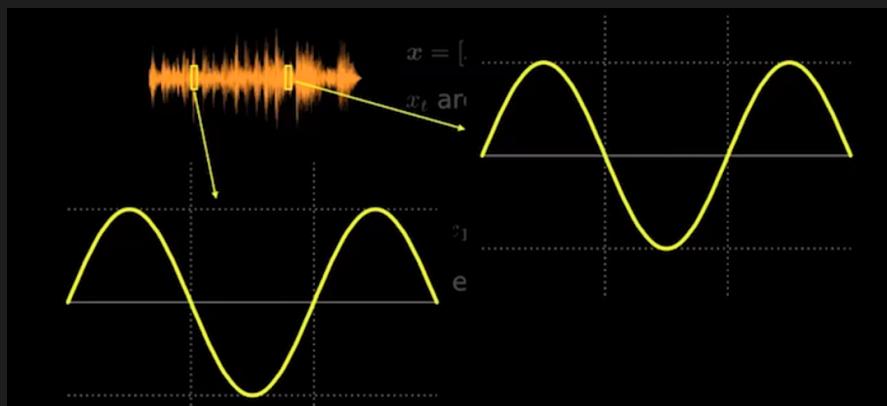
"John picked up the apple"

$$x = [x_1 \ x_2 \ x_3 \ x_4 \ x_5]^T$$

x_t are one-hot vectors

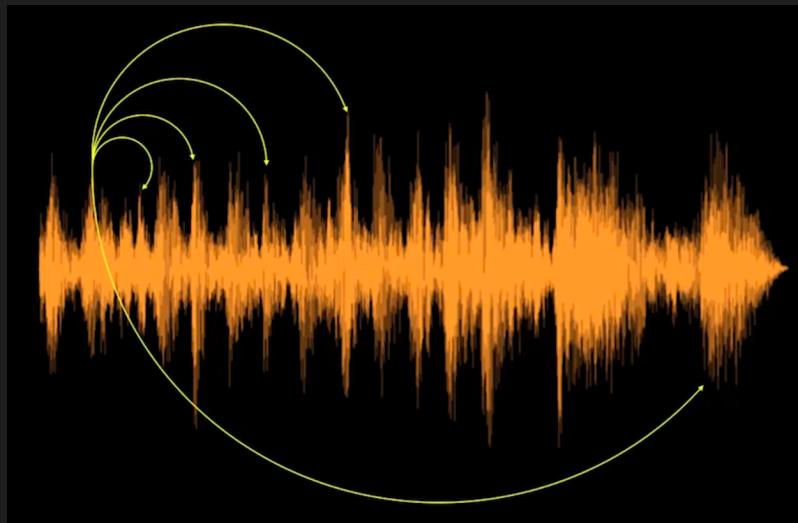
Stationarity :

- Patterns repeat



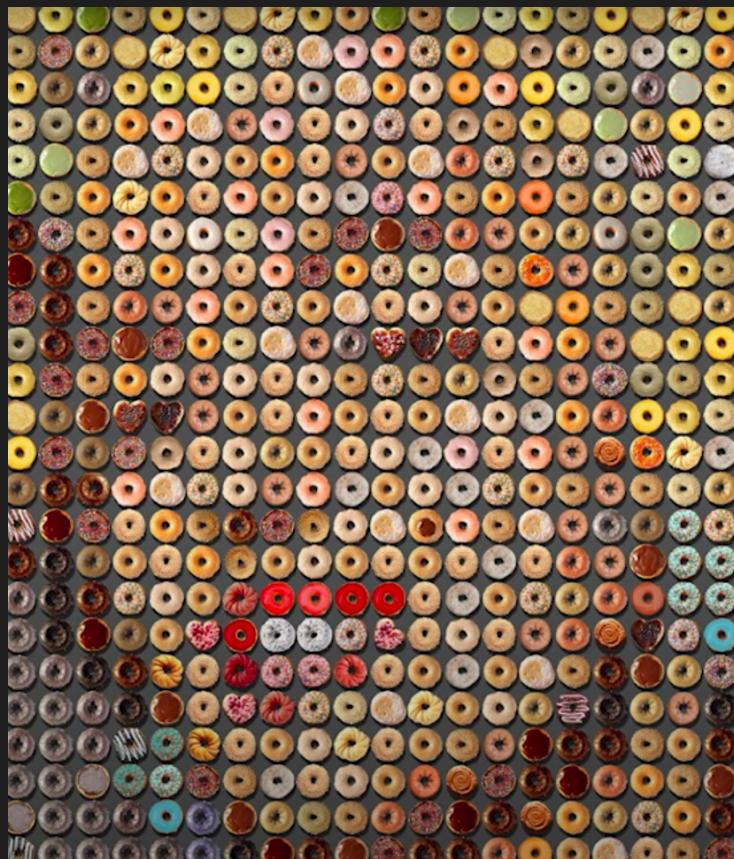
Locality :

- Strong correlation when closer.



Compositionality

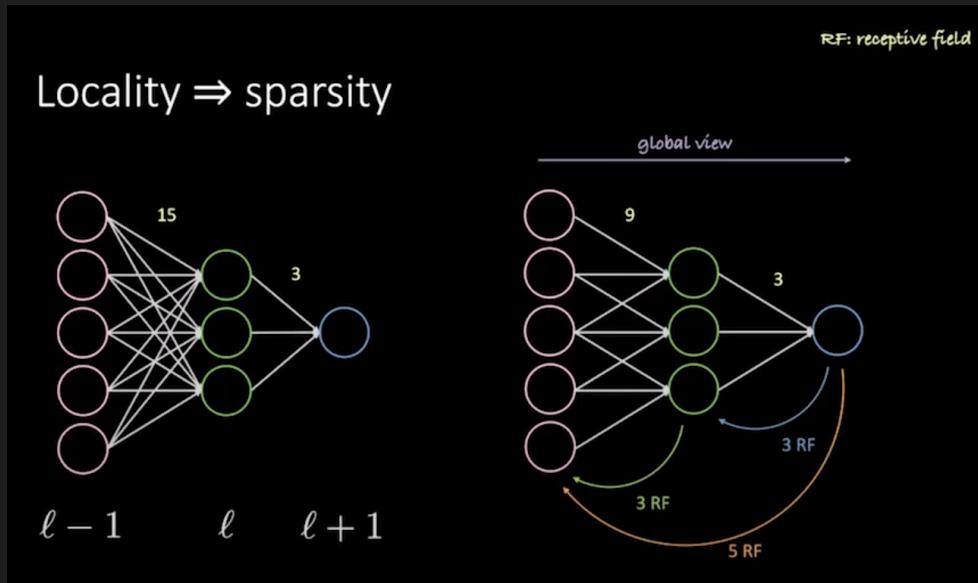
- Meaning comes from hierarchical structure



Review:

- CNNs work best when input has:
 - ↳ Locality (Close \equiv Similar)
 - ↳ Stationarity (Patterns repeat)
 - ↳ Compositionality (Big picture = Small pieces together)

Locality \Rightarrow sparsity

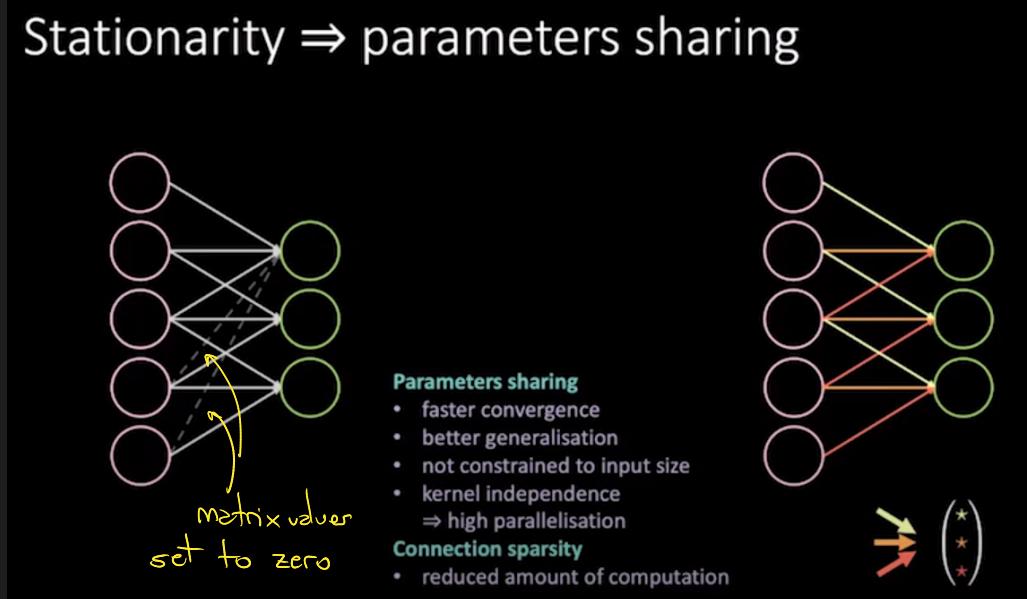


Recall

NNs : { Rotations
Squeezing }

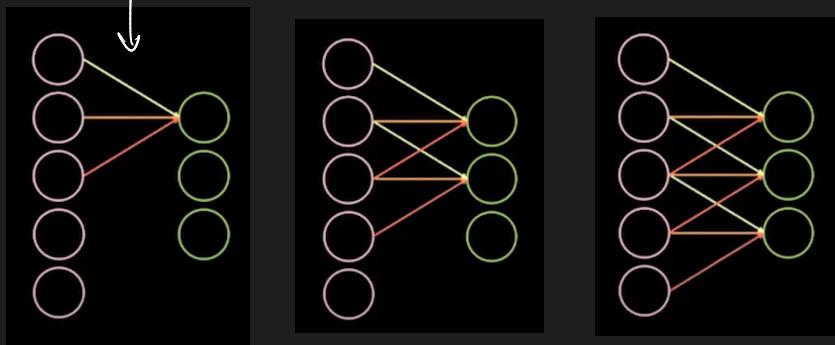
So CNN do these two) only over the RF)

Stationarity \Rightarrow parameters sharing



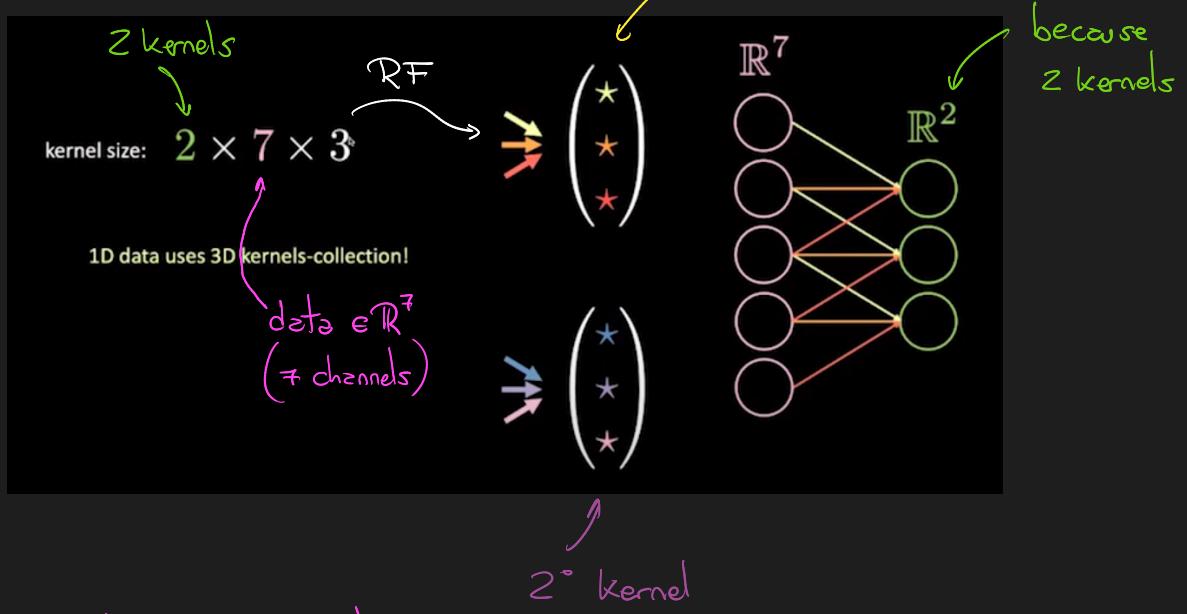
Kernels - 1D data

Kernel of dim. 3

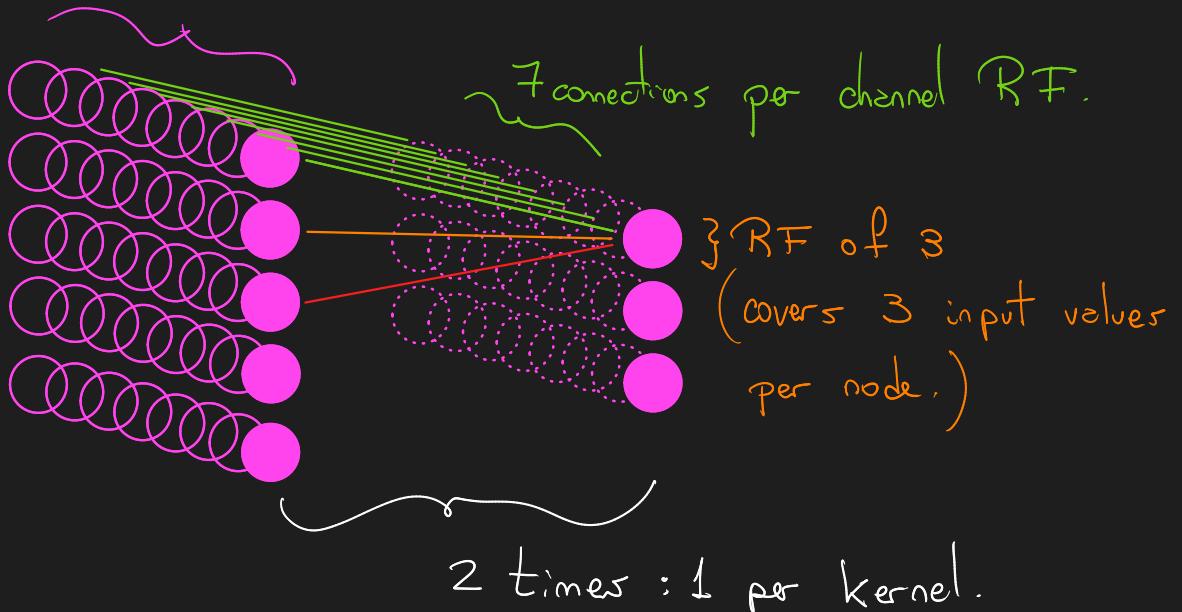


Lets add a second kernel

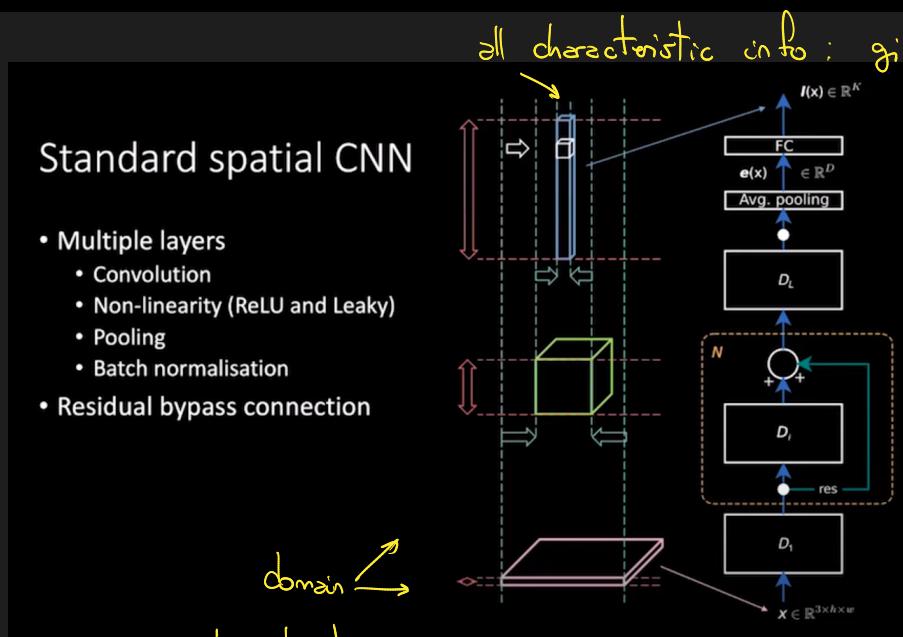
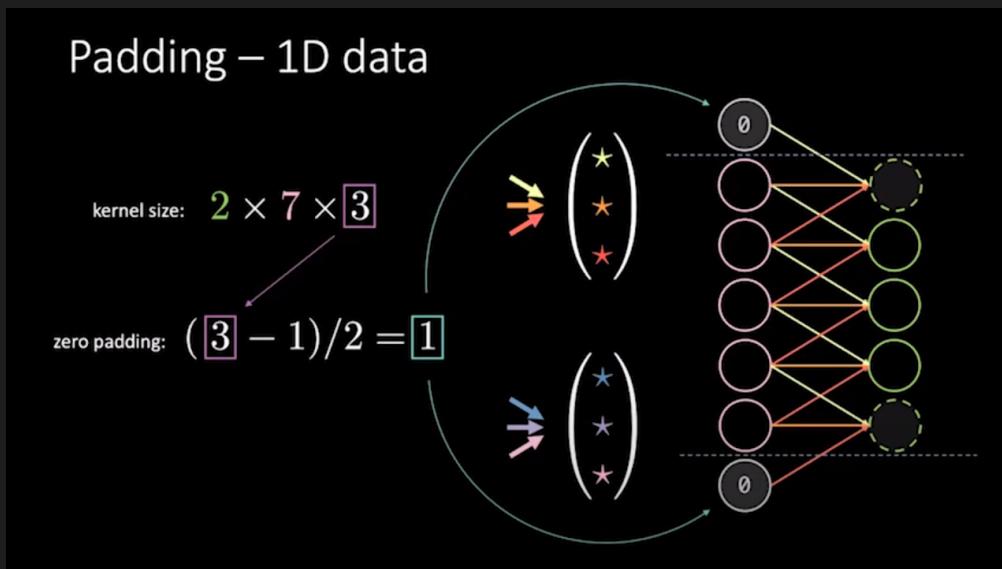
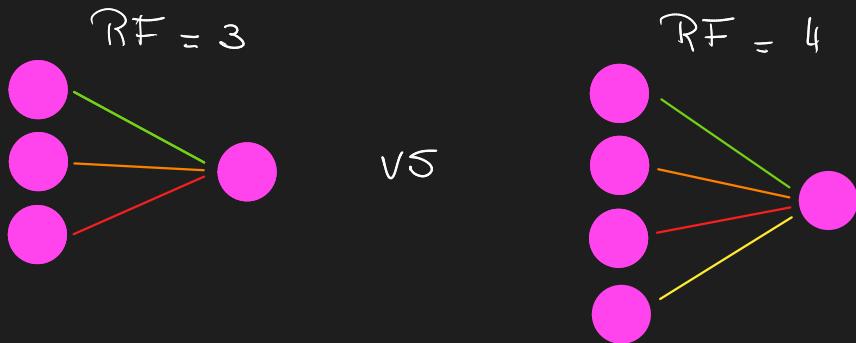
- 1° Kernel



Input has 7 channels

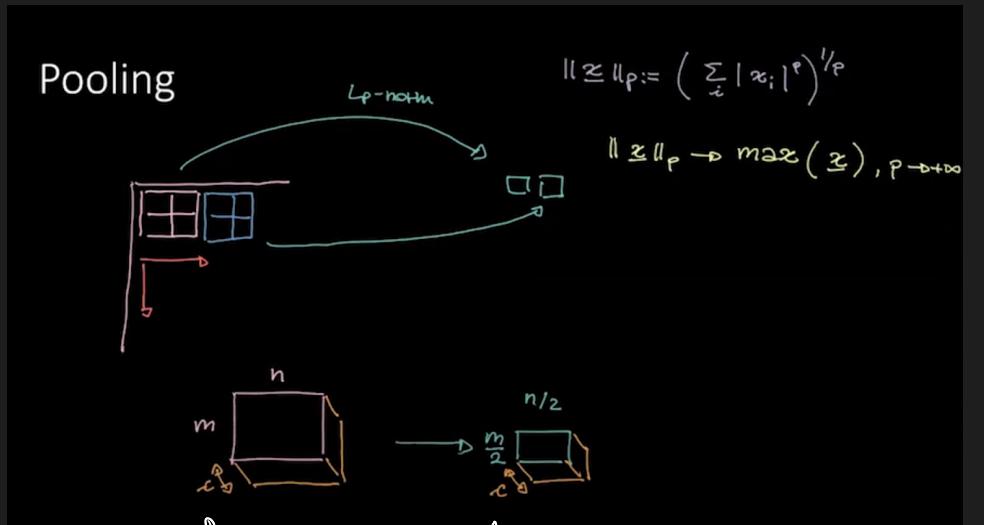


RF usually is odd, so the center correspond to a single value, and not 2



Characteristic information
(in this case, the color of the pixel)

Network: Reduction in domain, increasing characteristic information,

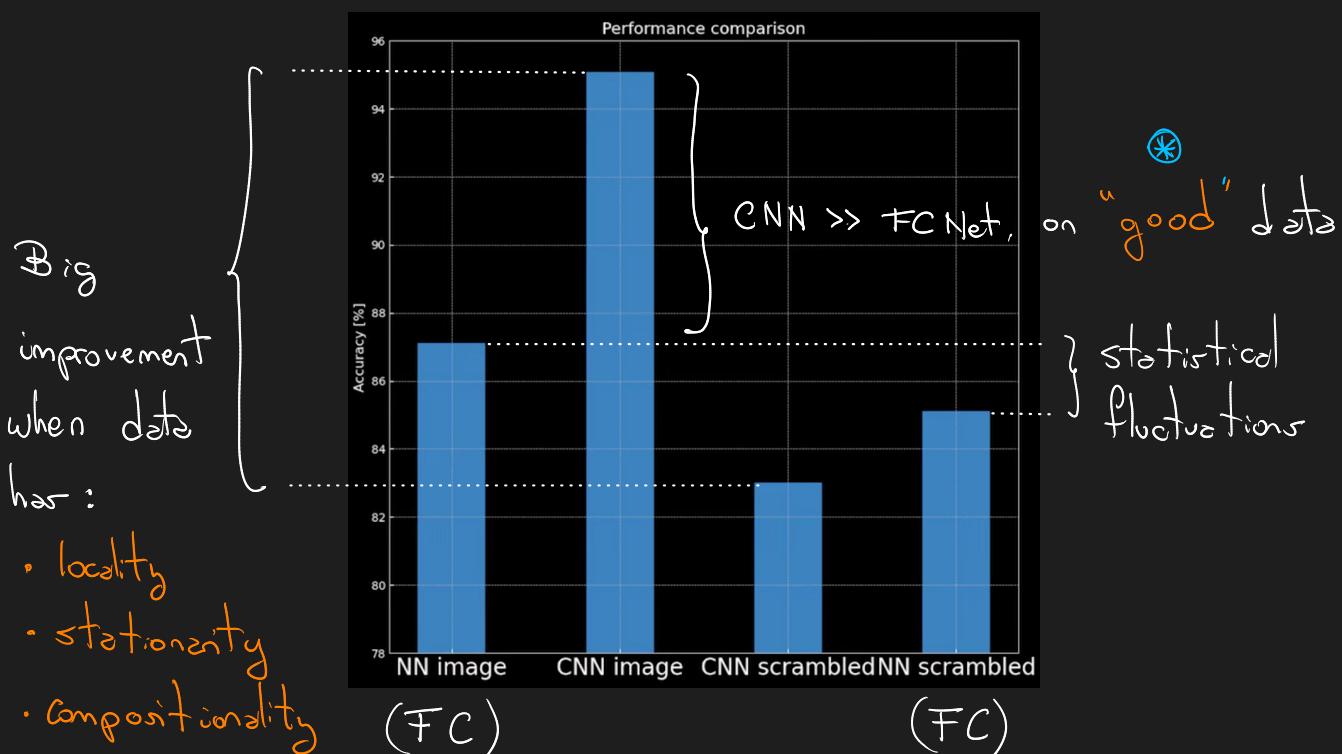


$$\begin{matrix} n \\ m \end{matrix} \xrightarrow{\quad} \begin{matrix} n/2 \\ m/2 \end{matrix}$$

channels after pooling stays the same.

Note book demo:

Alf shows how when permuting the pixels of the image (in a deterministic way), we loss all advantages of the CNN priors, resulting in similar performance to a Fully Connected Network, which perform exactly the same for both input cases.



⊗ where "good" means

it has the 3 properties