

Practical 3.0

Convolutional Neural Networks – Basics

Overview (I)

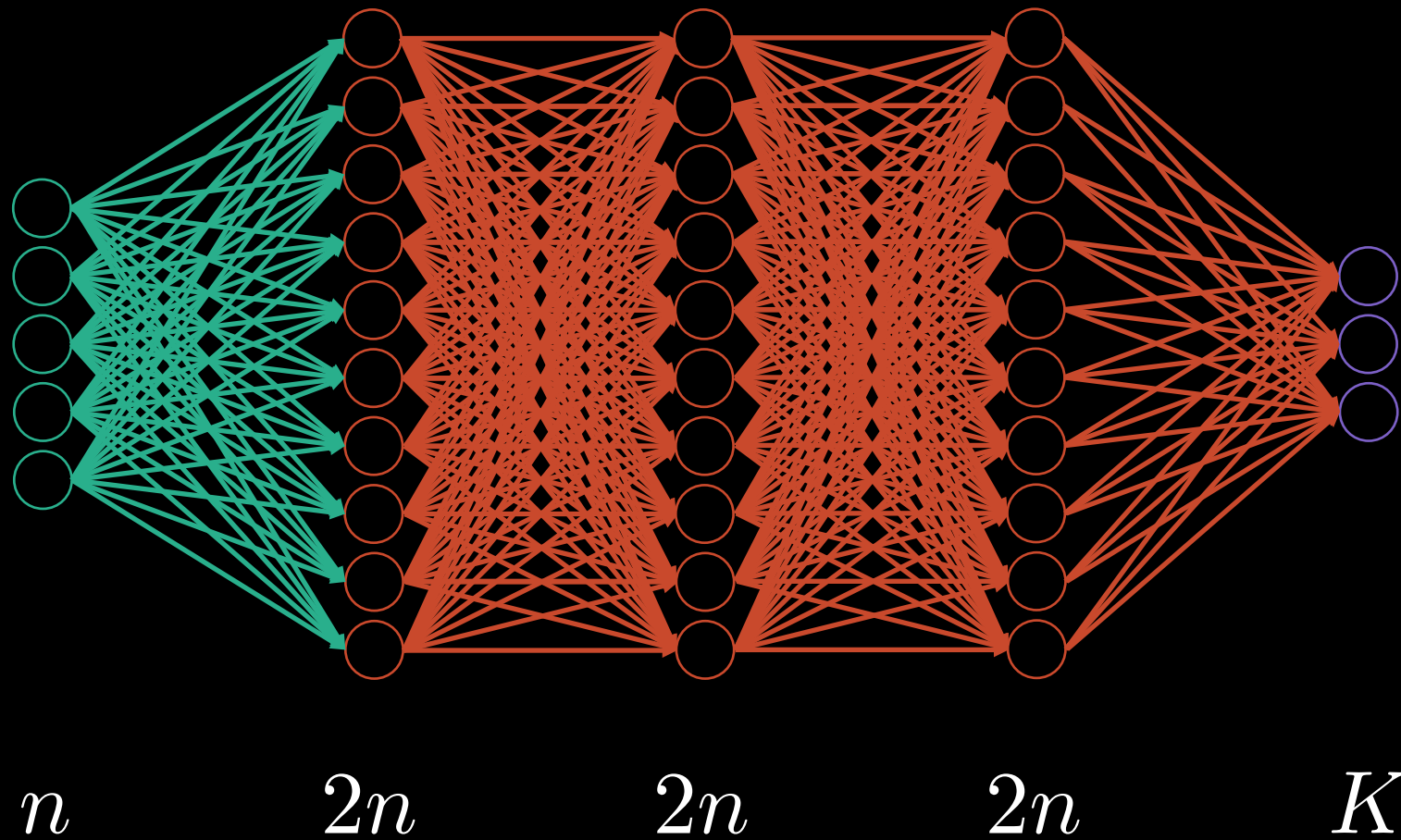
- Rationale / motivation
- Sparsity (locality)
 - Receptive field
 - Hierarchical view
- Parameters sharing (homogeneity)
- Convolutional layer (3D conv.)
- Non linearity layer (decaying learning speed)
 - Logistic sigmoid
 - Rectifying linear unit

Overview (II)

- L-p pooling layer
 - Average pooling
 - Max pooling
- Rationale / conclusion
 - Convolutional benefit
 - Pooling benefit
- Network colour coding with **pretty-nn**
- **e-lab/Torch7-profiling** tool
 - Profile CNN in Torch

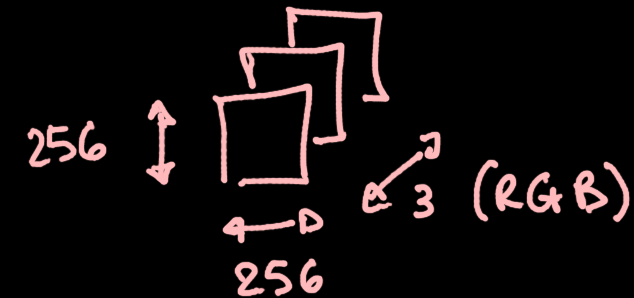
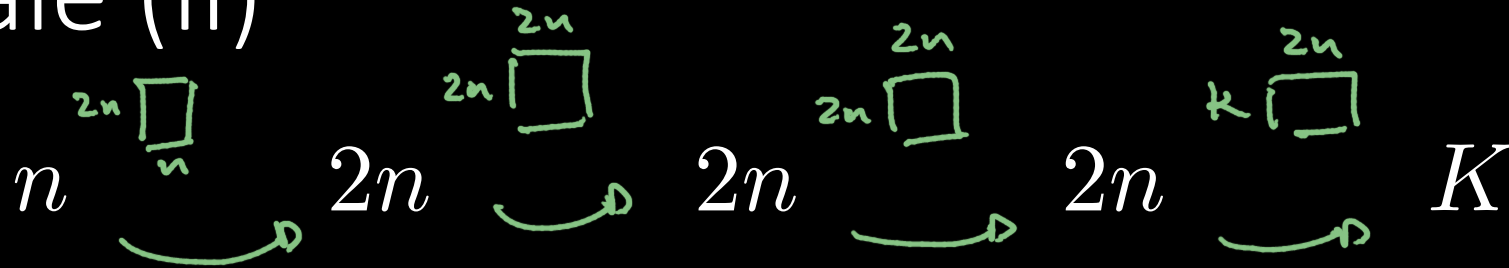


Rationale (I)



Rationale (II)

IMAGES \rightarrow Hi spatial correlation
 \rightarrow Objects are made of parts



$$n = 3 \cdot 256 \cdot 256 = 3 \cdot 2^8 \cdot 2^8 = 3 \cdot 2^{16}$$

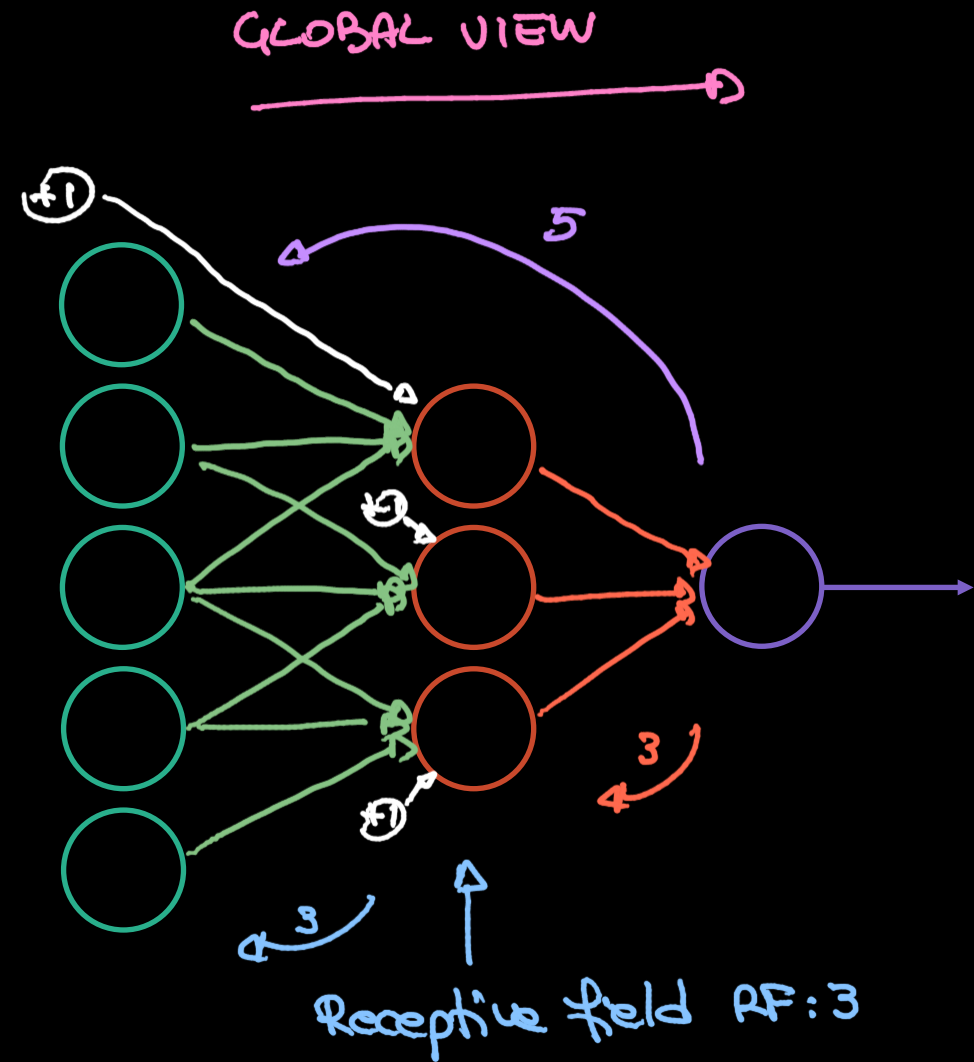
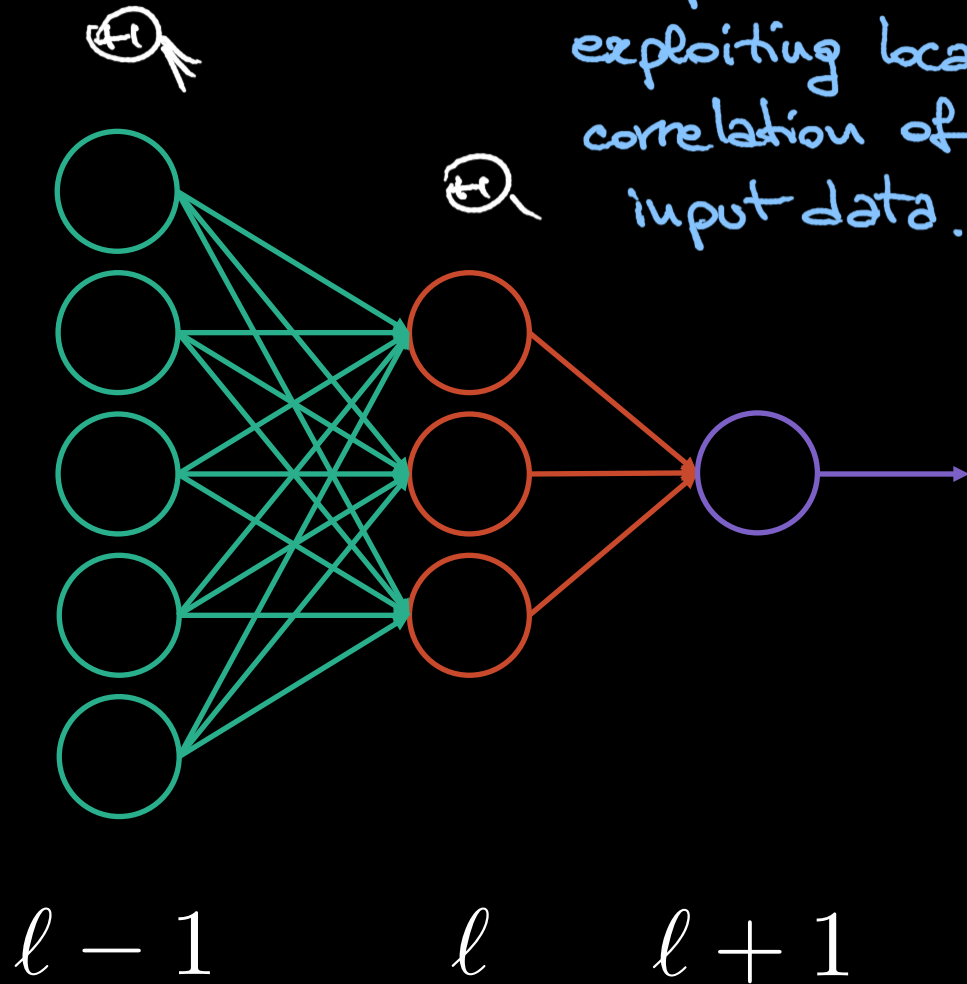
$$7n + K \approx 7 \cdot 3 \cdot 2^{16} = 1.4M$$

$$2n^2 + 2 \cdot 4n^2 + 2 \cdot 32n \approx 10n^2 = 2.5 \cdot 2^{32} \cdot 3 \approx 390G \text{ MAC}$$

Multiplication and accumulation

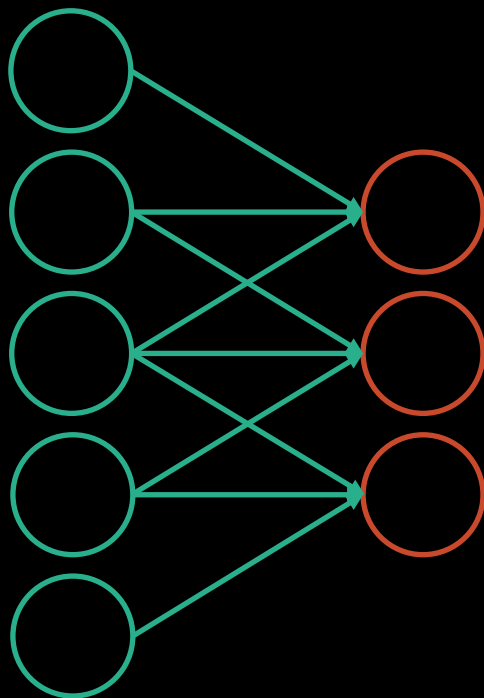
MacBook PRO : Intel Core i7 2.2 GHz $\Rightarrow 390G / 2.2GHz \approx 180s \approx \underline{\underline{3min}}$

Sparsity \Rightarrow REDUCTION in computation by exploiting local correlation of the input data.

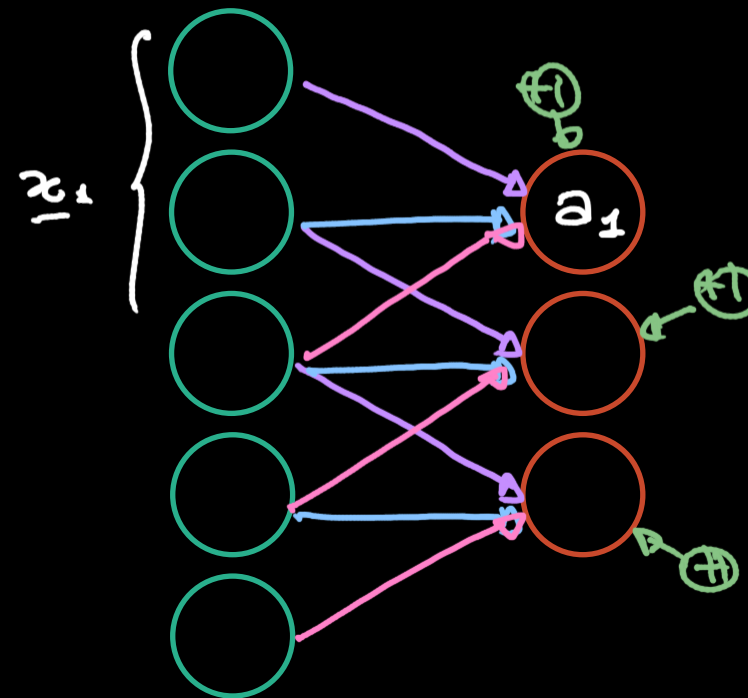


Parameters sharing \Rightarrow

Reduce $\# \Theta \Rightarrow$ helps
convergence \rightarrow time \downarrow
 \rightarrow error \downarrow



$$a_1 = \sigma(\underline{\theta}^T \underline{x}_1)$$



Convolutional layer (I)

\underline{x} 3D: n_1 2D $n_2 \times n_3$ (image input $\Rightarrow n_1 = 3$, RGB) $\uparrow 256 \times 256$

\underline{y} 3D: m_1 2D $m_2 \times m_3$ y_i feature maps ($m_2 \times m_3$)

\underline{k} 4D: m_1 3D $n_1 \times p_1 \times p_2$ kernels k_i , $i = 1, \dots, m_1$

\underline{b} 1D: m_1 bias term

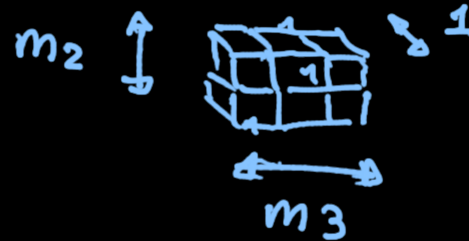
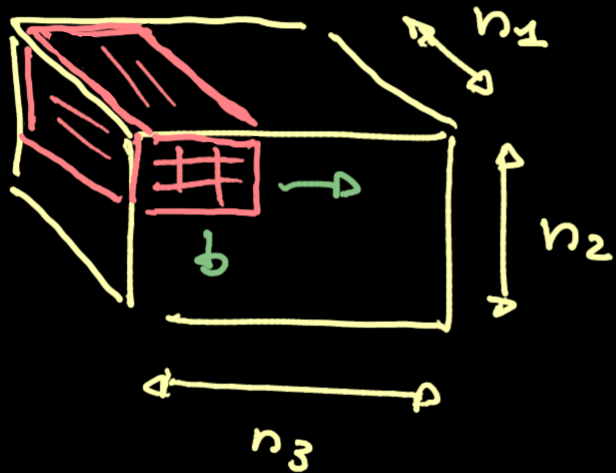
$$y_i = \underline{x} \overset{\text{CONVOLUTION}}{\star} \underline{k_i} + \underline{b_i}, \quad i = 1, \dots, m_1$$

$$(\underline{z} = \ominus \cdot \underline{a} \text{ or } \ominus \cdot \underline{x})$$

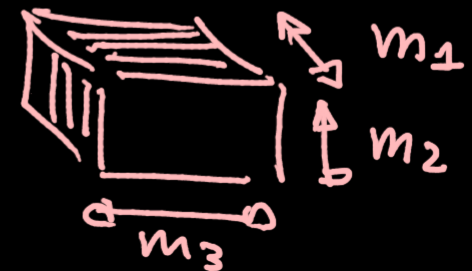
Convolutional layer (II)

$$y_i = x \star k_i + b_i, \quad i = 1, 2, \dots, m_1$$

$$f[l, m, n] \star g[l, m, n] = \sum_u \sum_v \sum_w f[u, v, w] \cdot g[l-u, m-v, n-w]$$



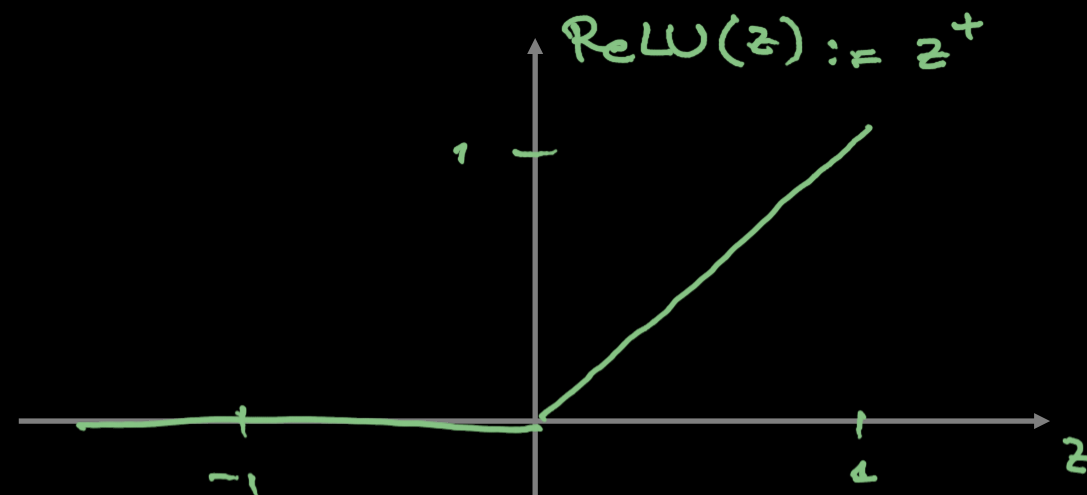
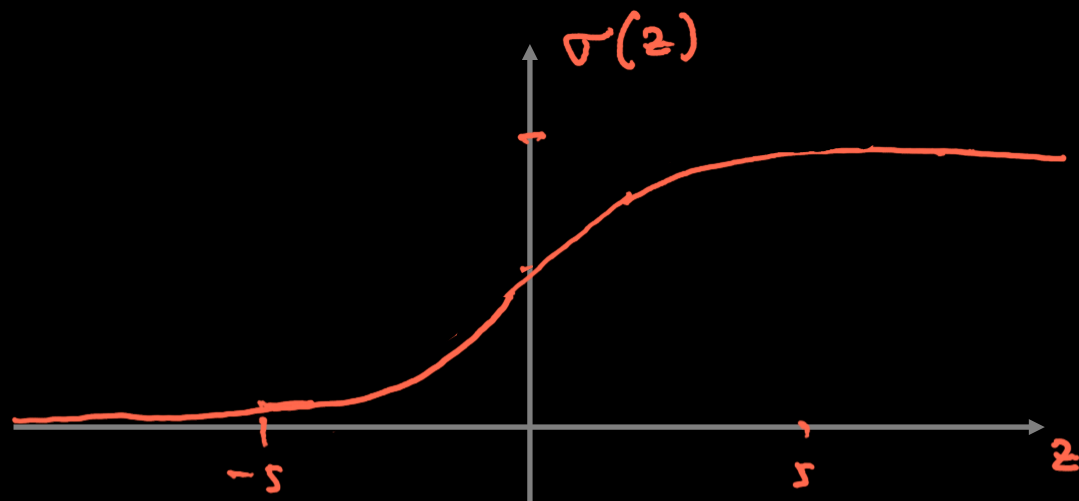
$\times m_1$ times



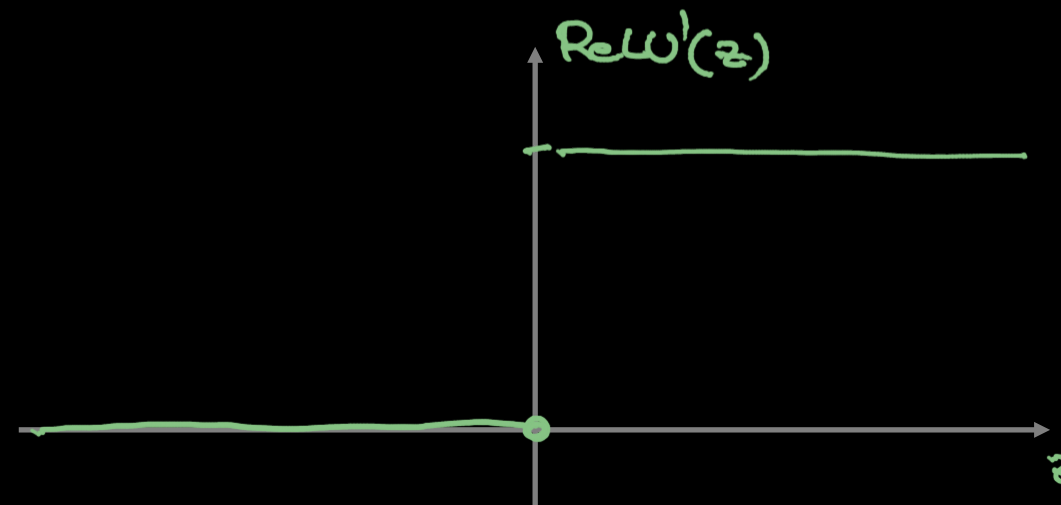
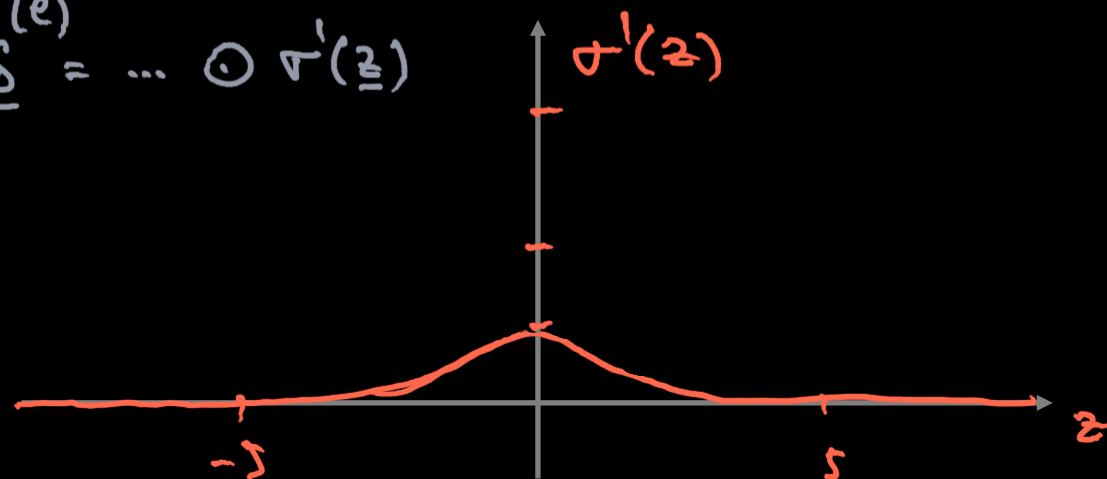
Non linearity

$$\underline{z} = \ominus \underline{\hat{z}}^{(l-1)} \rightsquigarrow \underline{a} = \sigma(\underline{z}) \text{ feature map.}$$

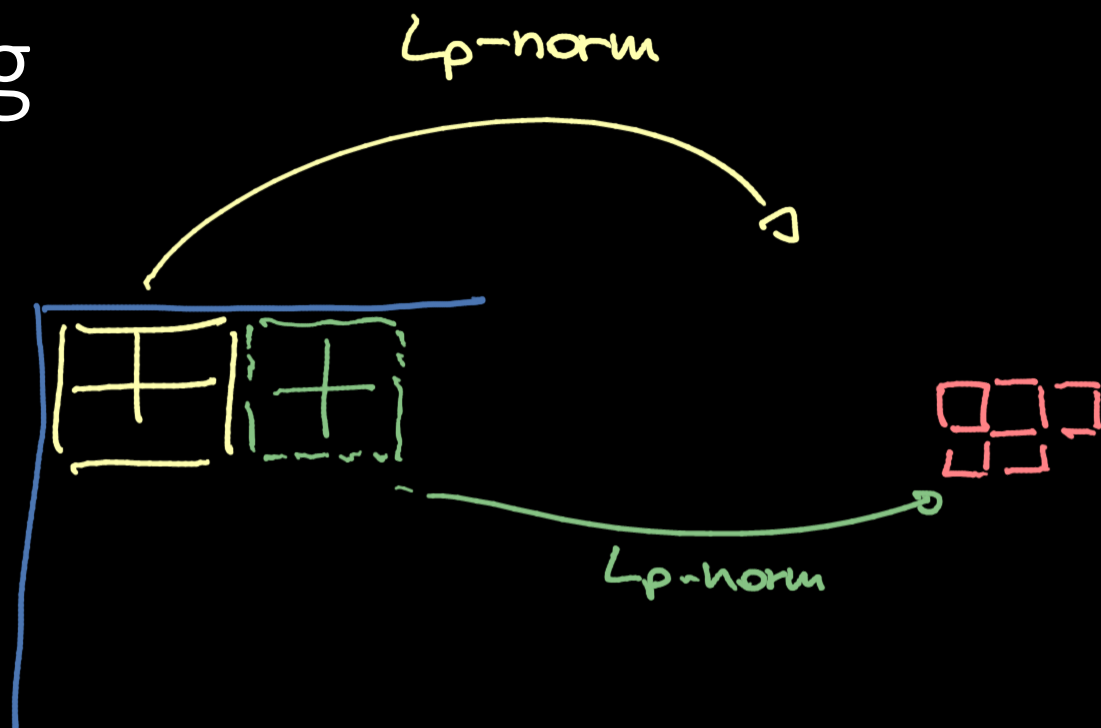
\uparrow projection by conv. \curvearrowright



$$\underline{\delta}^{(l)} = \dots \odot \sigma'(\underline{z})$$

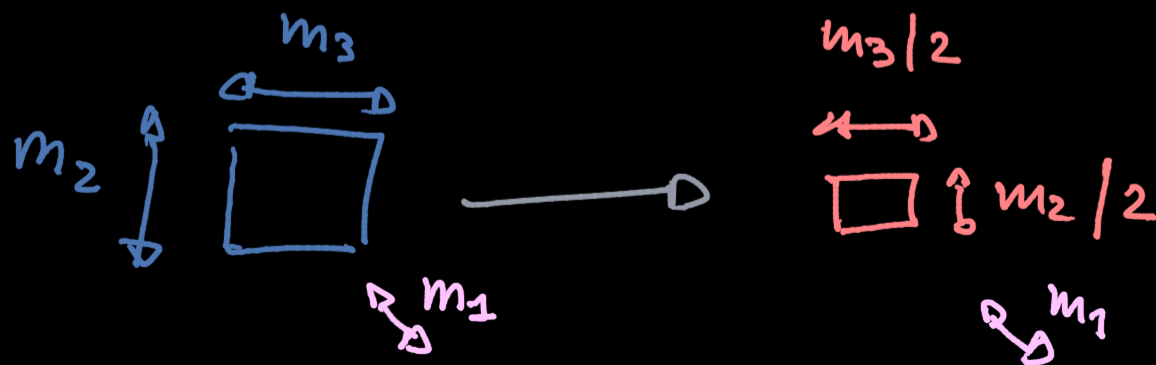


Pooling



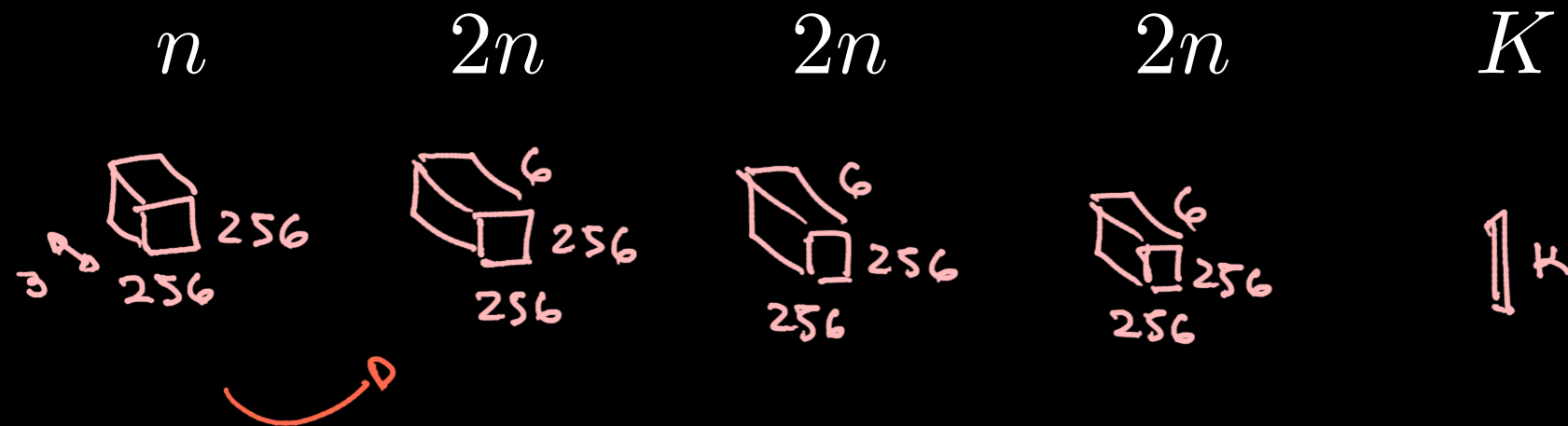
$$\| \underline{x} \|_p = \left(\sum_i |x_i|^p \right)^{1/p}$$

$$\| \underline{x} \|_p \rightarrow \max(\underline{x}), p \rightarrow +\infty$$



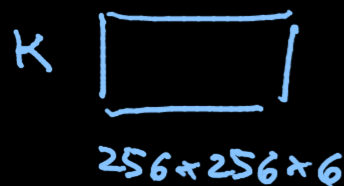
$k_i : 3/6$
 $k_i : 5 \times 5$

Rationale (III)



$$\underbrace{(256 \times 256 \times 6)}_{\# \text{ activations}} \times 3 \times 5 \times 5 \times \underbrace{(1+2+2)}_5 = 147 \text{ M MAC}$$

540 M MAC



$$\rightarrow K \cdot 6 \cdot 256 \cdot 256 = K \cdot 393k = 393M$$

ImageNet : $K=1k$

$\swarrow 2.2 \text{ GHz} \rightarrow$
250ms

Rationale (IV)

$$31,7 \text{ M Mac} / 2.2 \text{ GHz} = 14 \text{ ms}$$



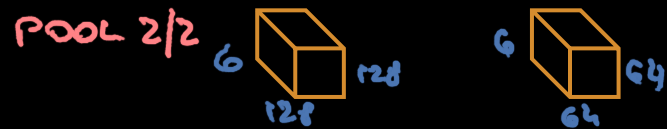
$$128 \times 128 \times 6 \times 5 \times 5 \times 3 \approx 7,3 \text{ M Mac} \quad \}$$

98k NL Mac

25,7 M CONV

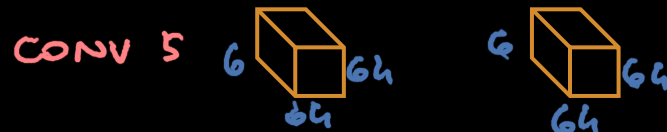


$$128 \times 128 \times 6 \times 5 \times 5 \times 6 \approx 14,7 \text{ M Mac} \quad \}$$



$$64 \times 64 \times 6 \times 2 \times 2 \approx 98 \text{ k Mac}$$

LINEAR



$$64 \times 64 \times 6 \times 5 \times 5 \times 6 \approx 3,7 \text{ M Mac} \quad \}$$

$$32 \times 32 \times 6 = 6 \text{ k}$$



$$32 \times 32 \times 6 \times 2 \times 2 \approx 25 \text{ k Mac}$$

