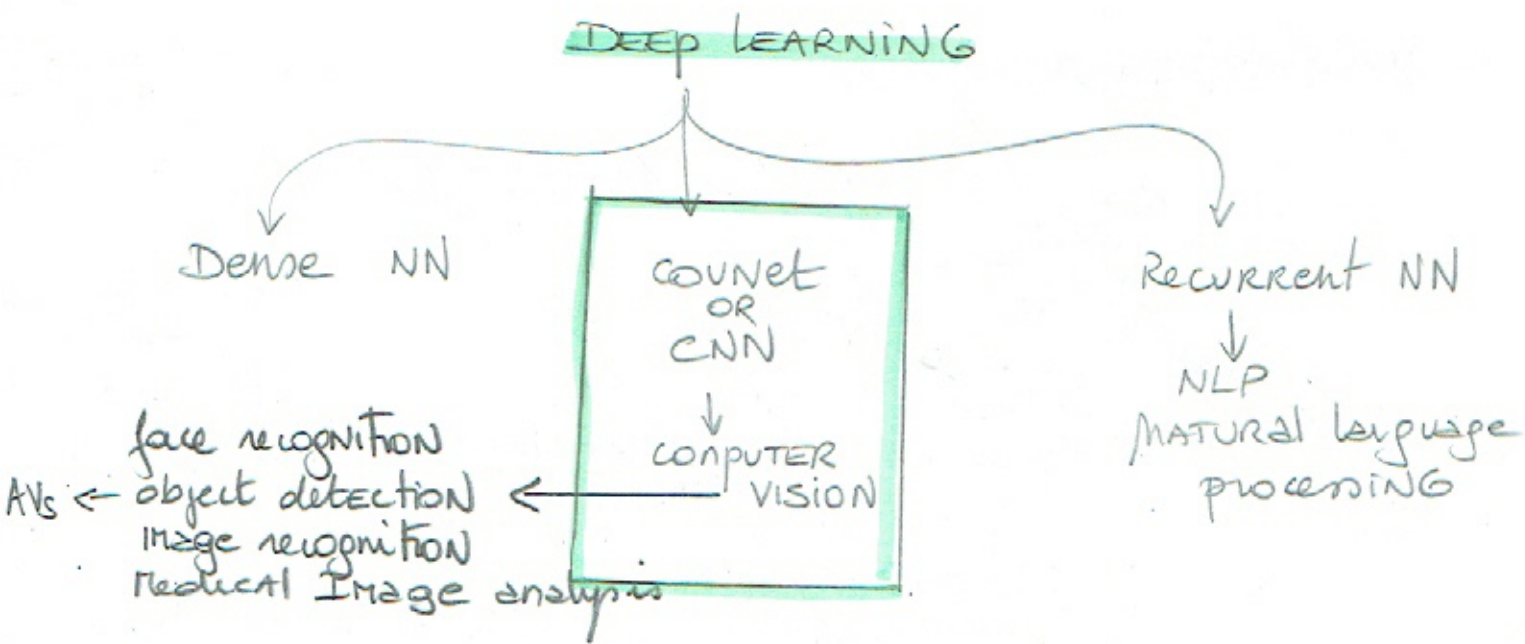
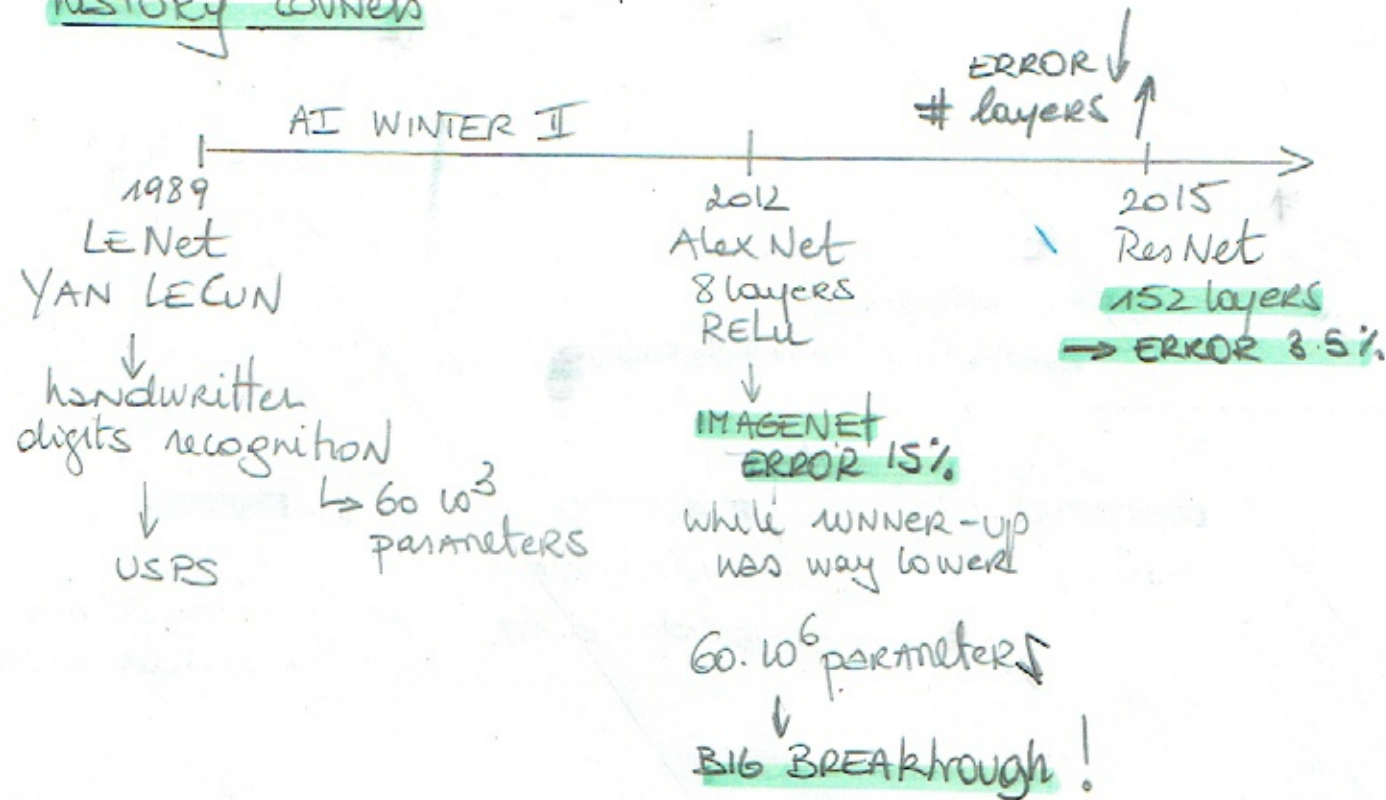


# Lesson 16: Convolutional Neural Networks: ConvNets ①



## history convnets



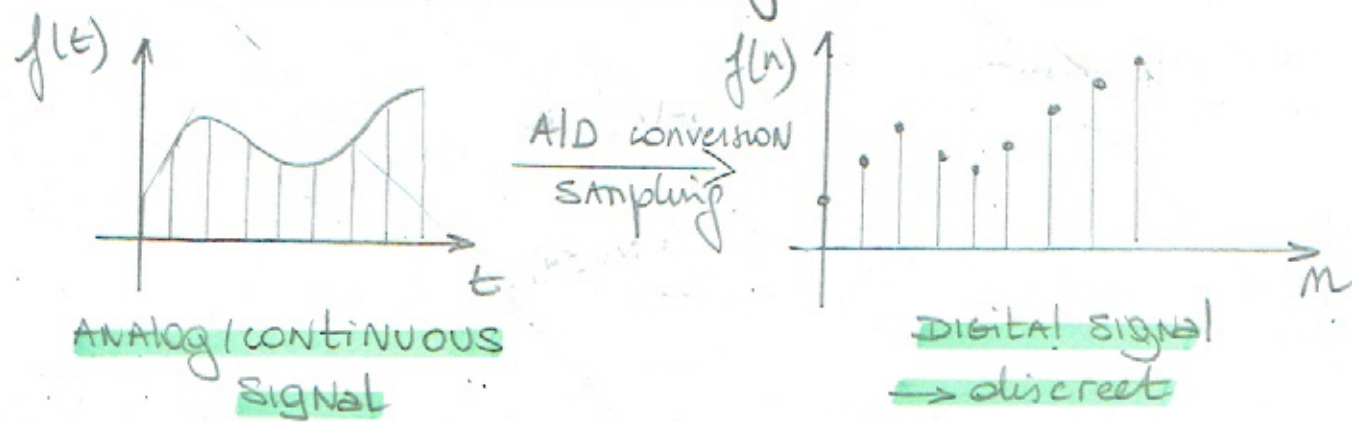
content ① CONVOLUTION → SIGNAL theory (optional)

② ConvNets → BASICS + CONCEPTS

- filters
- padding
- Strides
- convolution vs CROSS CORRELATION
- volume convolution / 2D convolution
- MAXpooling
- feature maps / feature extraction

- convnets for small datasets (2)
- transfer learning
  - fine tuning
  - data augmentation
- coding (keras) → CAT dog picture recognition

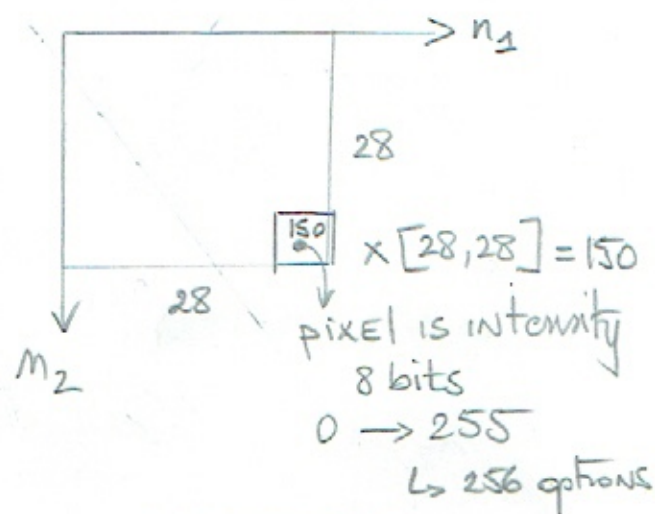
## ① Convolution → Signal Theory (optional)



## Image representation

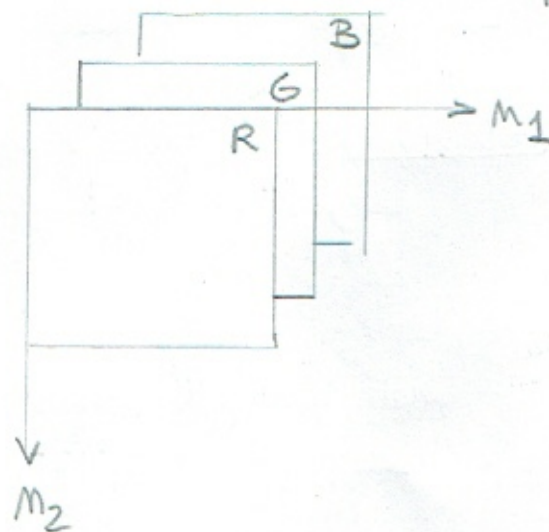
**grayscale image**

↓ SCALAR  
2D discrete signal.



**color image**

↓ VECTOR  
2D Discrete signal.



$$x(n_1, n_2) = [x_R(n_1, n_2), x_G(n_1, n_2), x_B(n_1, n_2)]$$

↓  
 $x_R, x_G, x_B$  represent pixel intensity along their axis

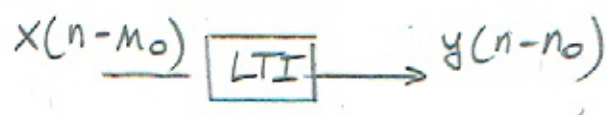
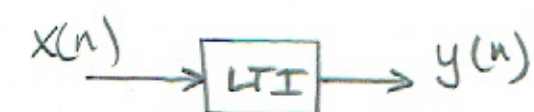
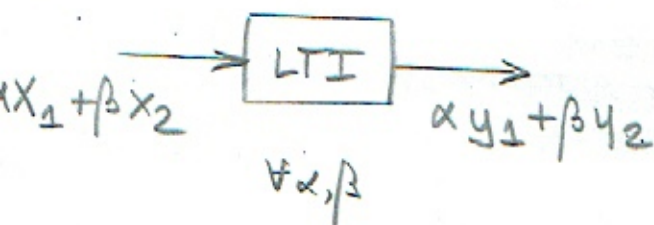
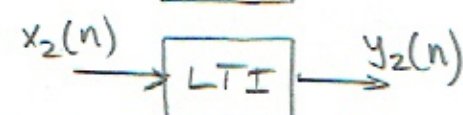
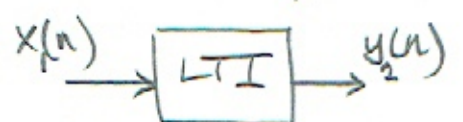


# convolution of 2 signals

(3)

↳ LINEAR, TIME-INVARIANT Systems (LTI)

superposition  
scaling



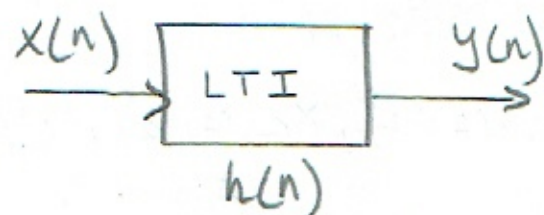
shift in input → same shift in output.

now what? when system is LTI

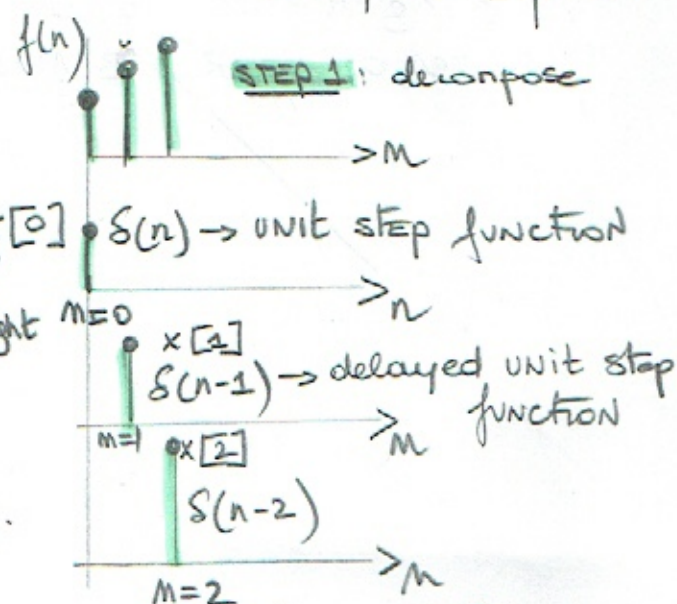
→ decompose an arbitrary signal into a set of basic signals

↳ delayed impulse responses  $h(n)$   
(convolution)

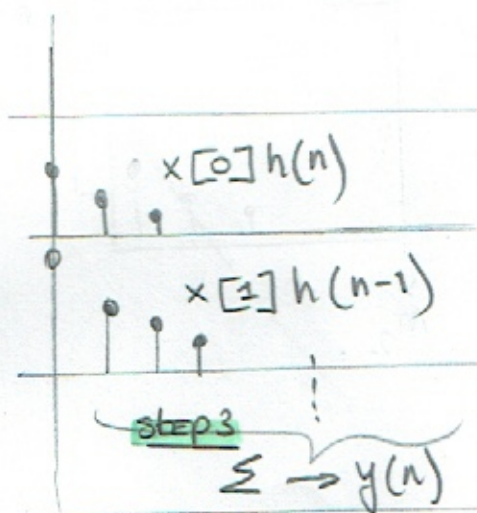
GOAL: determine  $y(n)$   
based on  $x(n)$  and  $h(n)$



↳ impulse response



unit step response  
STEP 2



$$\rightarrow x[n] = x[0] \delta(n) + x[1] \delta(n-1) + x[2] \delta(n-2) + \dots \rightarrow \sum_{k=-\infty}^{\infty} x[k] \delta(n-k) \quad (1)$$



(1) → we decomposed arbitrary sequence to a linear combination of delayed unit step functions

↓ LINEARITY

response to that linear combination is a linear combination of the responses

$$\rightarrow y(n) = \sum_{k=-\infty}^{+\infty} x(k) h_k(n)$$

↓ TIME-INVARIANT (response to impulse at time  $k$  = response to an impulse at time 0) shifted over to time  $k$

$$\rightarrow y(n) = \sum_{k=-\infty}^{+\infty} x_k h(n-k)$$

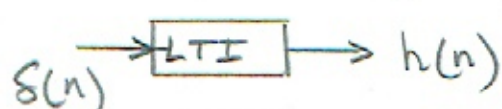
↳ CONVOLUTION SUM

$$y(n) = x * h(n)$$

$$\rightarrow h_k(n) = h_0(n-k)$$

SUMMARY:

→ impulse response

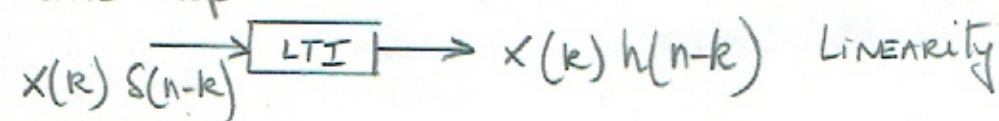
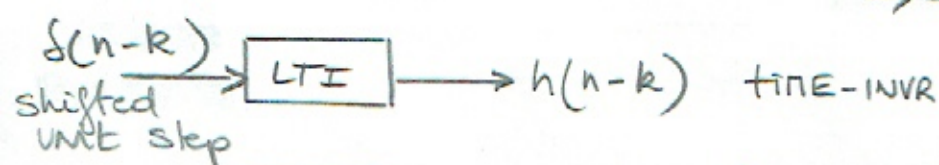


$\delta$ : unit step function

↳ DIRAC Delta function

$$\delta(n) = \begin{cases} +\infty & n=0 \\ 0 & n \neq 0 \end{cases}$$

(AREA = 1)



CONCLUSION for LTI

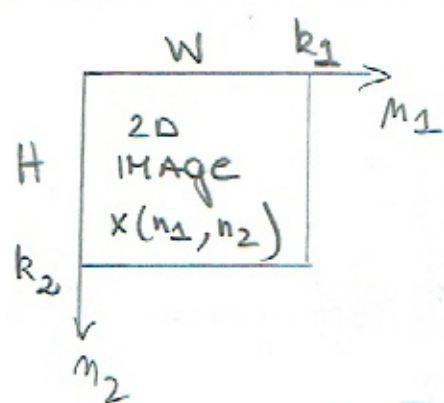
if you know its response to a unit step signal at  $n=0$  you can generate the response to an arbitrary sequence  $y(n)$  through the convolution sum

$h(n)$

$\delta(n)$

# IMAGE PROCESSING → 2D convolution

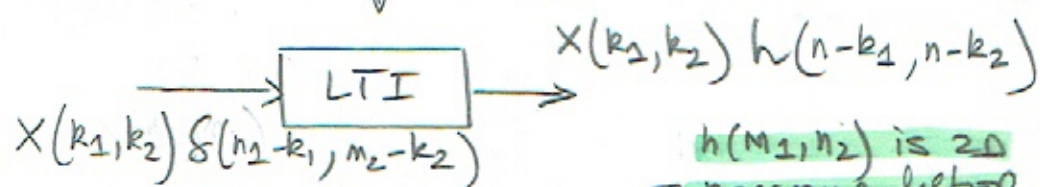
(5)



2D unit step function

$$S(n_1, n_2) = \begin{cases} +\infty & \text{if } n_1 = n_2 = 0 \\ 0 & \text{elsewhere} \end{cases}$$

1D ANALOGY

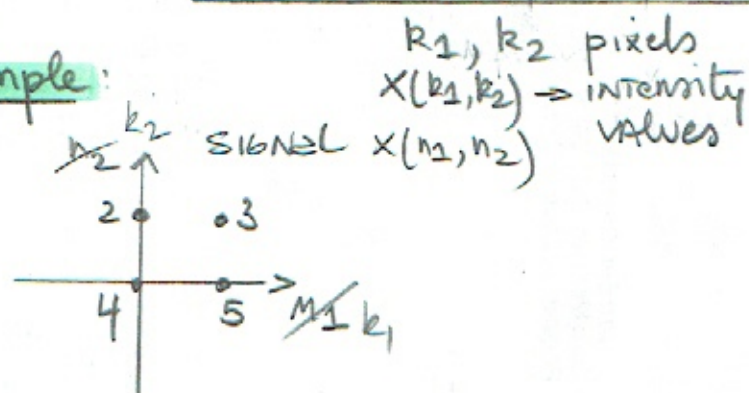


$h(n_1, n_2)$  is 2D processing filter

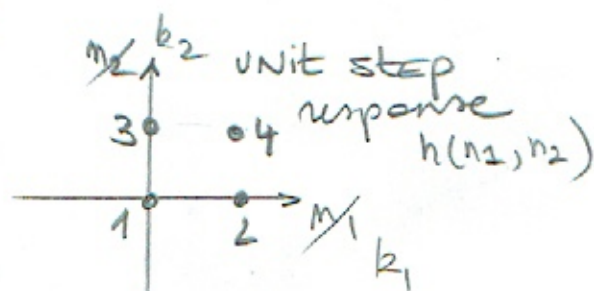
$$y(n_1, n_2) = \sum_{k_1=0}^{H-1} \sum_{k_2=0}^{W-1} x(k_1, k_2) h(n_1 - k_1, n_2 - k_2)$$

convolution sum

example:



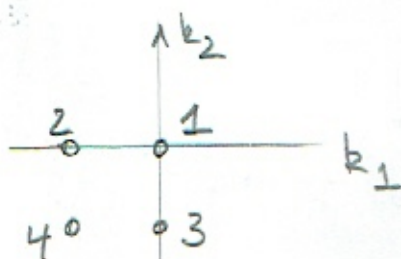
$$x(n_1, n_2) = \begin{cases} = 5 & n_1=1, n_2=0 \\ = 4 & n_1=n_2=0 \\ = 2 & n_1=0, n_2=1 \\ = 3 & n_1=n_2=1 \\ = 0 & \text{elsewhere} \end{cases}$$



$$h(n_1, n_2) = \begin{cases} = 2 & n_1=1, n_2=0 \\ = 1 & n_1=n_2=0 \\ = 3 & n_1=0, n_2=1 \\ = 4 & n_1=n_2=1 \\ = 0 & \text{elsewhere} \end{cases}$$

Step 1:  $n_1, n_2$  axes →  $k_1, k_2$  axes

Step 2: leave  $x(k_1, k_2)$  as is  
reverse impulse response  $h(k_1, k_2) \rightarrow h(-k_1, -k_2)$



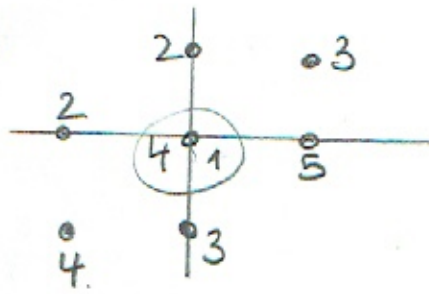
Step 3: → place  $h(-k_1, -k_2)$  at different offset values for  $n_1, n_2$   
→  $h(n_1 - k_1, n_2 - k_2)$



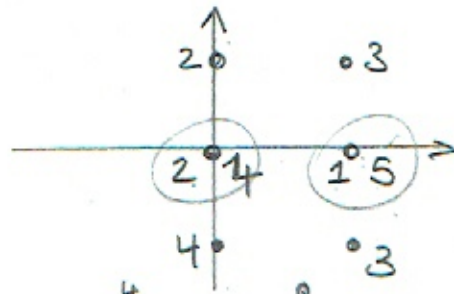
$$y(0,0) = x(0,0) \cdot h(0,0) = 4$$

we look for overlap  
otherwise conv. sum = 0! (6)

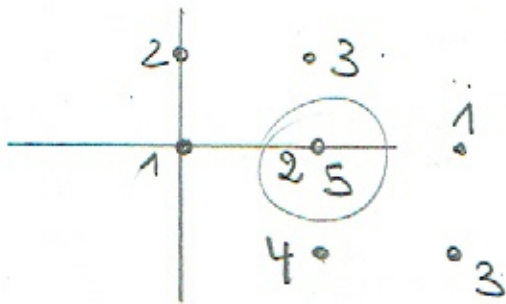
SITUATION 1



SITUATION 2



SITUATION 3



$$y(1,0) = x(0,0) \cdot h(0,0) \rightarrow 8$$

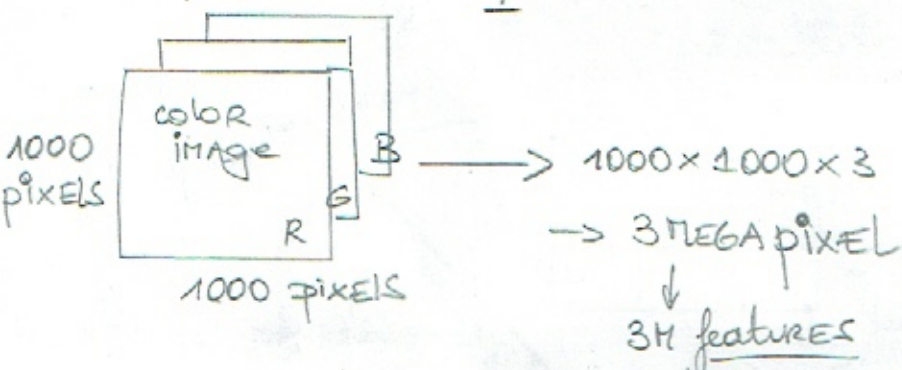
$$+ x(1,0) \cdot h(1,0) \rightarrow 5$$

$$= 13$$

$$y(2,0) = 10$$

ConvNets

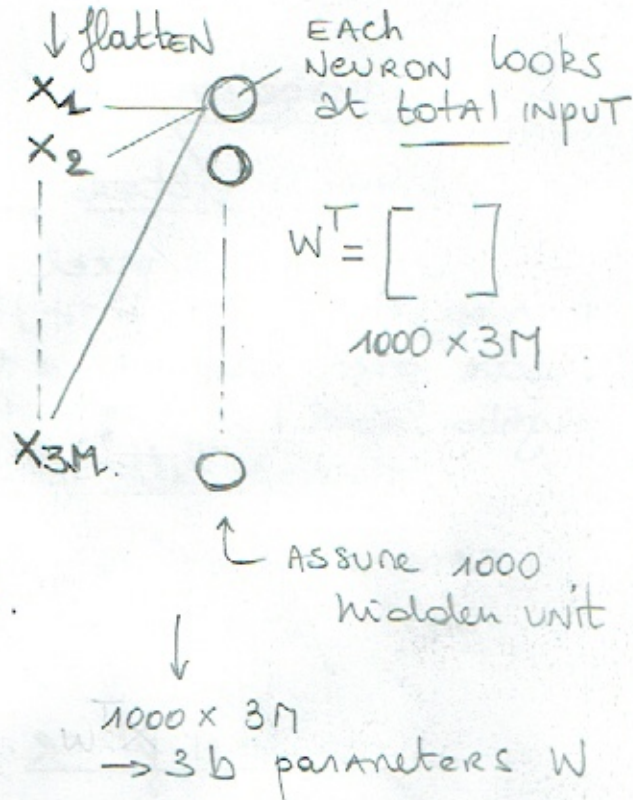
2 issues at hand with dense, fully connected NNs with respect to computer vision



(2) → flatten image  
↓  
spatial info is lost

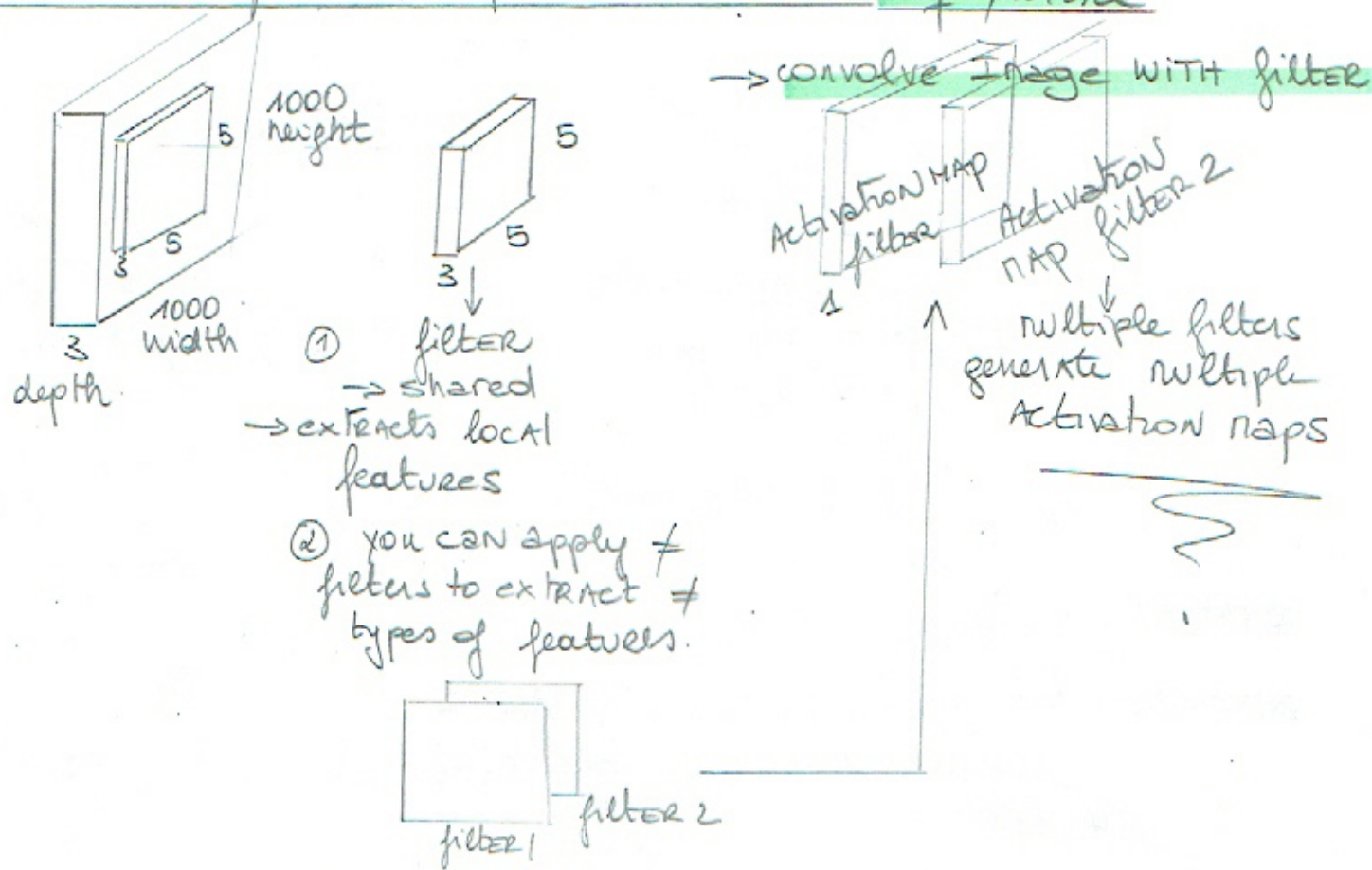
$$X = \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \quad 3M \times 1$$

3 Million Rows

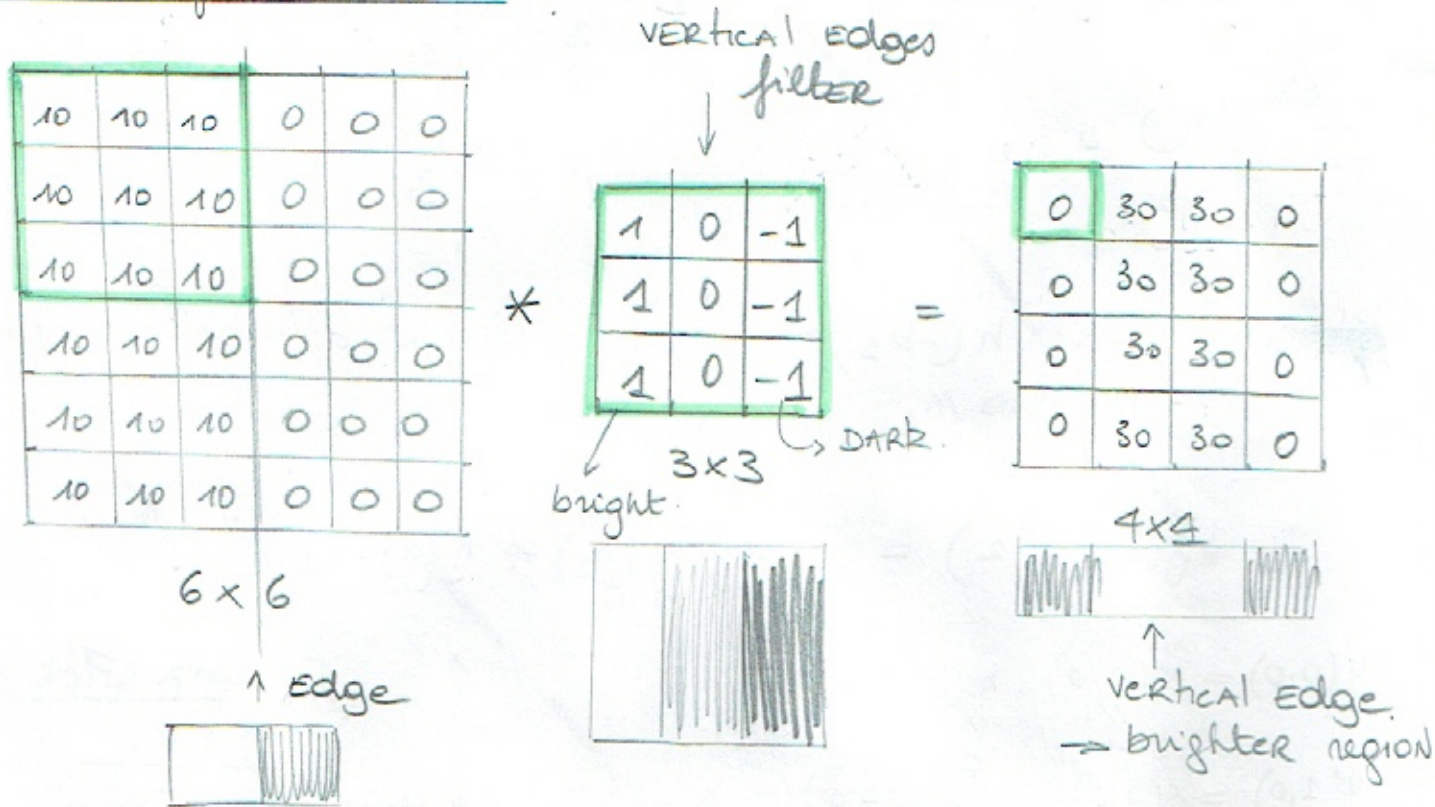


(1) → Dimensions quickly spin out of control

# We want to preserve spacial structure of picture (7)



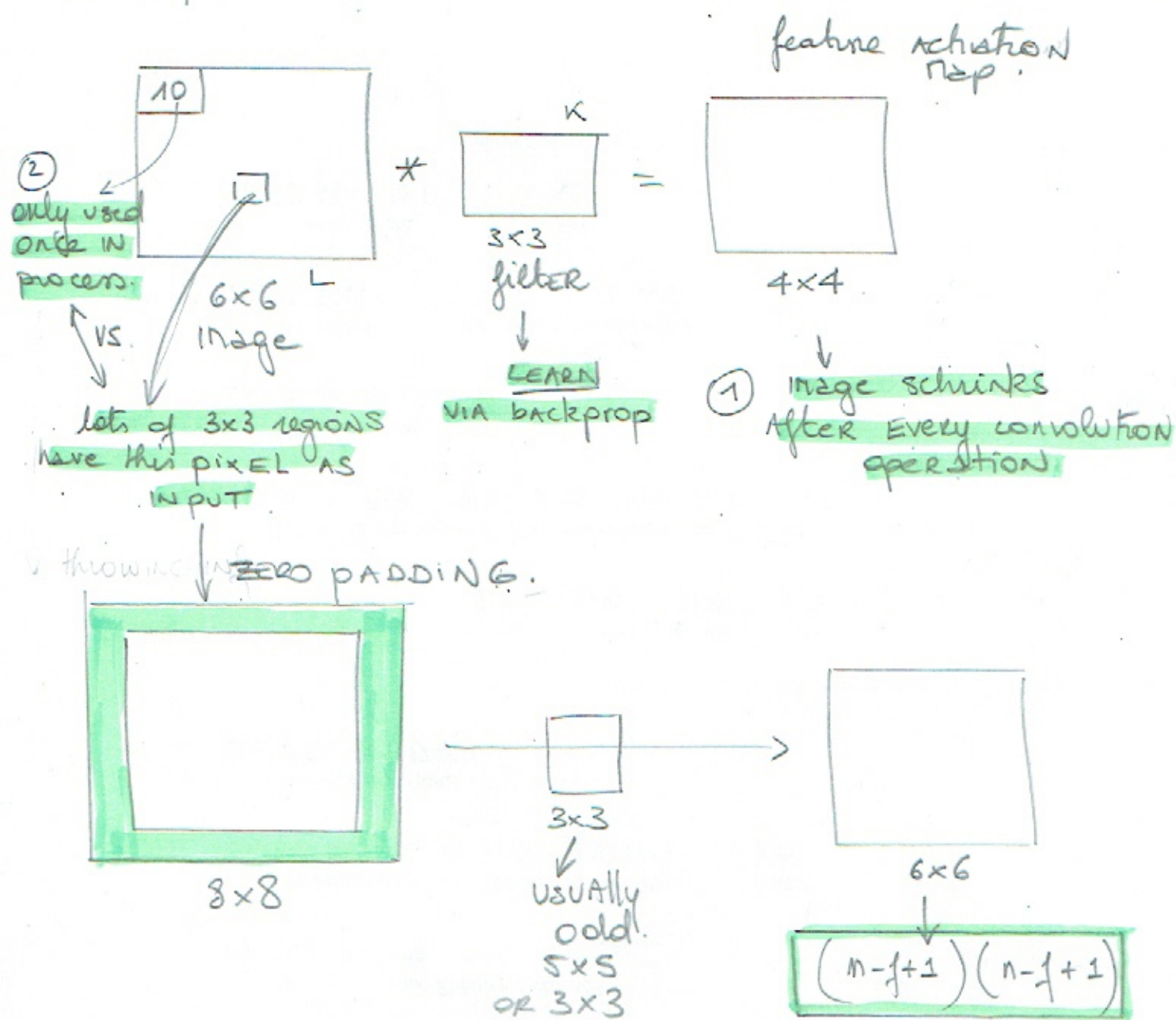
## How does filter work?





## ZERO PADDING

8



STRIDING  $\rightarrow$  default is 1

$\hookrightarrow$  number of pixels to move in each spatial direction

## OUTPUT IMAGE SIZE (convolved)

$L \rightarrow$  length of input image

$K \rightarrow$  length filter/kernel

$P \rightarrow$  # of zeros padded

$S \rightarrow$  Stride of the convolution

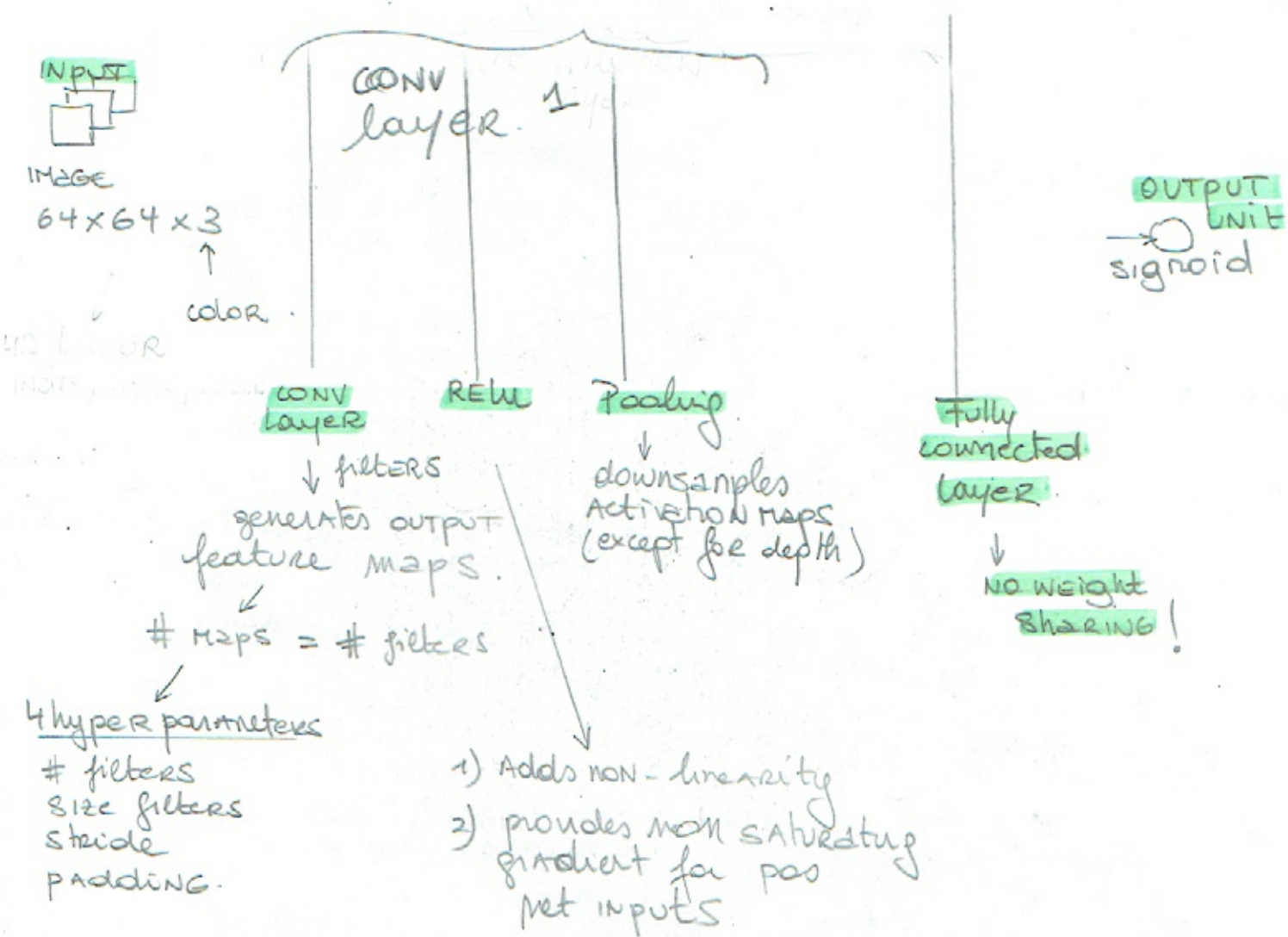
$$L' = \frac{L - K + 2P}{S} + 1$$

$\downarrow$   
length of convolved output image

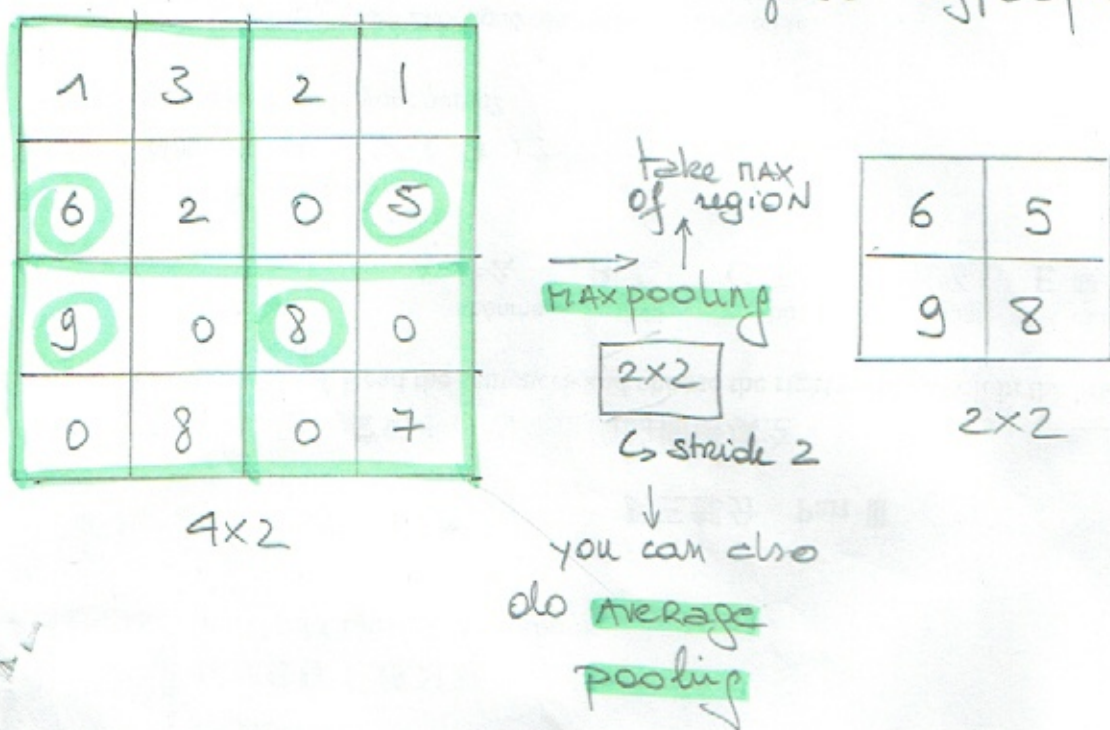


# Architecture of a convNet

9

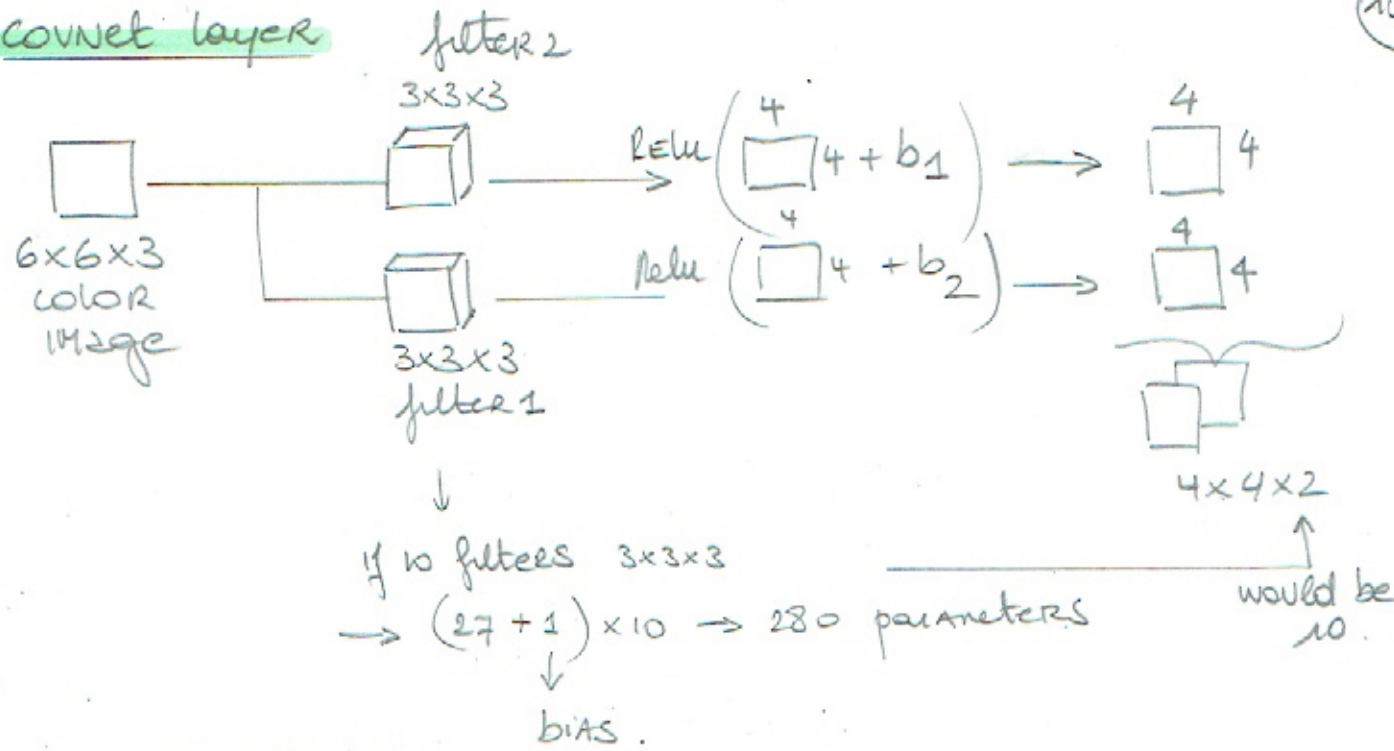


**Pooling** → downsample! → no parameters just hyperparameters

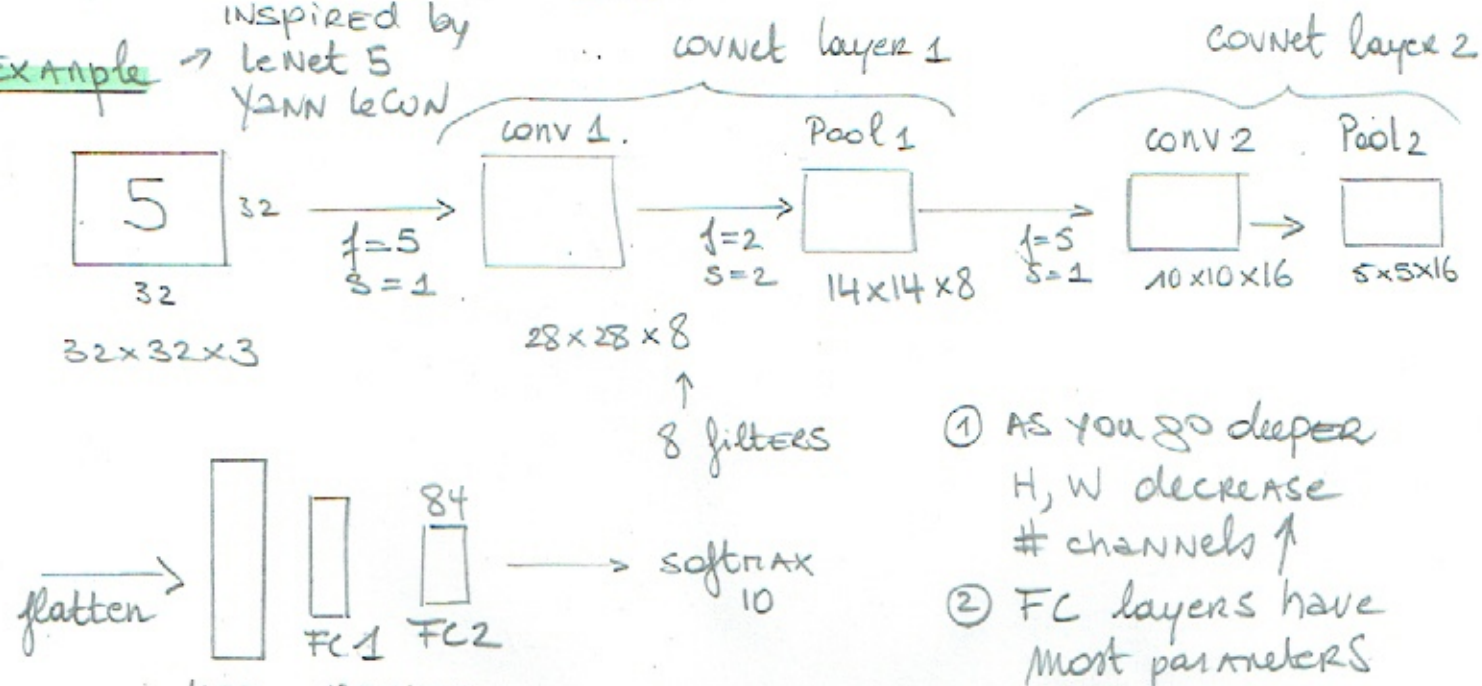


# convnet layer

10



Example → inspired by LeNet 5 YANN LECUN



hyperparameter.

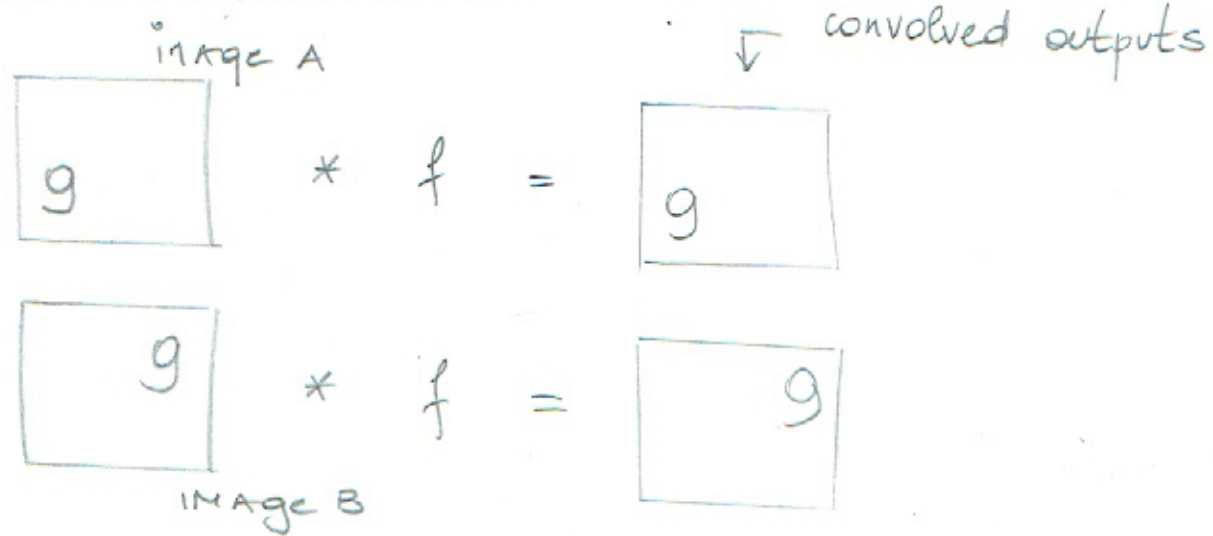
	Activation shape	Activation size	# parameters
INPUT	32x32x3	3072	0
conv 1	28x28x8	6272	$f=5 \rightarrow 5 \times 5 \times 3 = 75 \times 8 = 600$
Pool 1	14x14x8	1568	0
conv 2	10x10x16	1600	$f=5 \rightarrow 1216$
pool 2	5x5x16	400	0
FC1	120x1	120	48000 + 1
FC2	84x1	84	10080 + 1
Softmax	10x1	10	840 + 1

~ 60.10<sup>3</sup> parameters



# Translation Equivariance

11



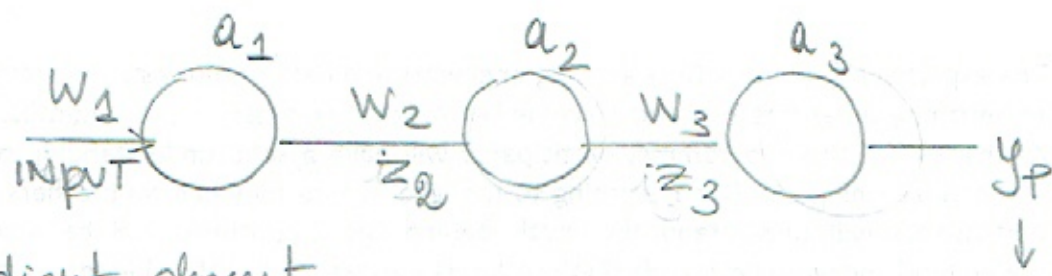
↳ same feature has been detected in both cases despite the fact that 9s are in different locations

Leaves a lot of information about edges

but not more complex features

# VANISHING GRADIENT

13

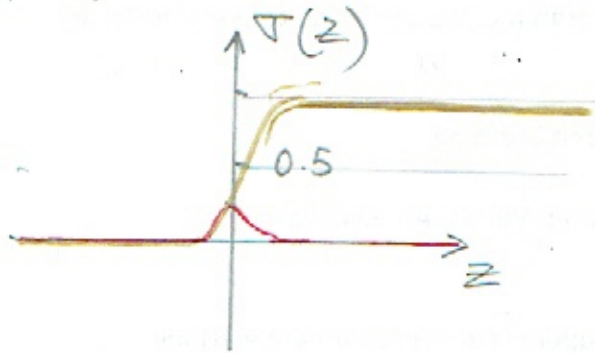


Gradient descent  
update rule

$$C \sim y_p - y$$

$$w_1^{\text{new}} = w_1^{\text{old}} - \alpha \frac{\partial C}{\partial w_1^{\text{old}}}$$

before 2000  $\rightarrow$  sigmoid activation:  $\sigma(z) = \frac{1}{1+e^{-z}}$



$$\sigma'(z) = \sigma(z)(1-\sigma(z))$$

$\rightarrow$  between  
0 - 0.25

$$\frac{\partial C}{\partial w_1} = \frac{\partial C}{\partial a_3} \cdot \frac{\partial a_3}{\partial a_2} \cdot \frac{\partial a_2}{\partial a_1} \cdot \frac{\partial a_1}{\partial w_1}$$

Chain rule.

multiplication of number between  
0 and 0.25

$\rightarrow$  Deeper  $\rightarrow$  smaller number

$\rightarrow$  VANISHING GRADIENT!

$$\text{because } w_1^{\text{new}} = w_1^{\text{old}} - \alpha$$

small  
number

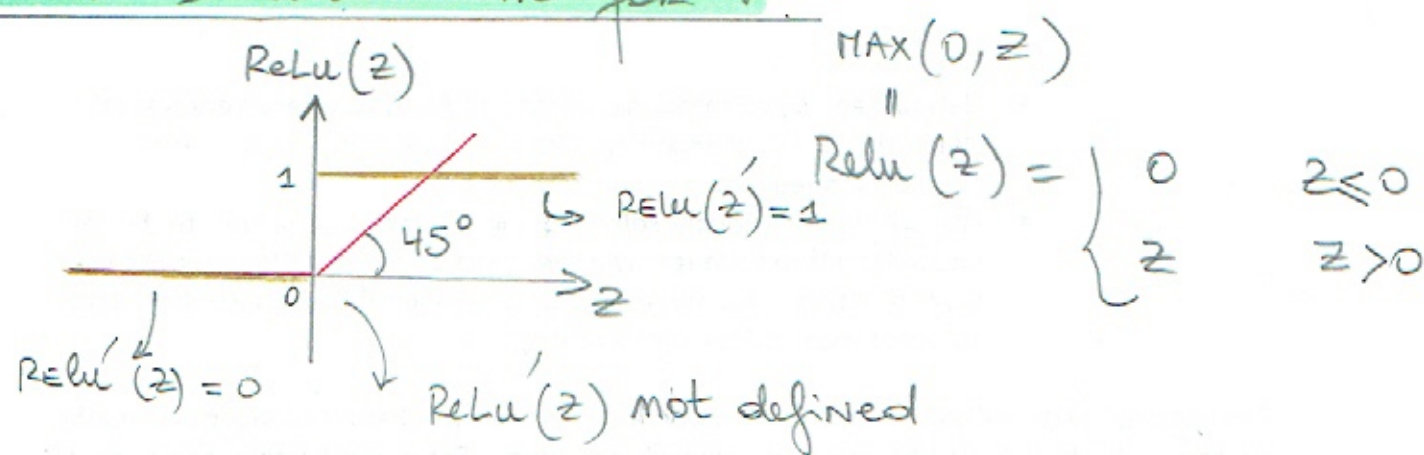
$$w_1^{\text{new}} \approx w_1^{\text{old}}$$

$\rightarrow$  algorithm stops



# Relu → alternative for $\sigma$

14



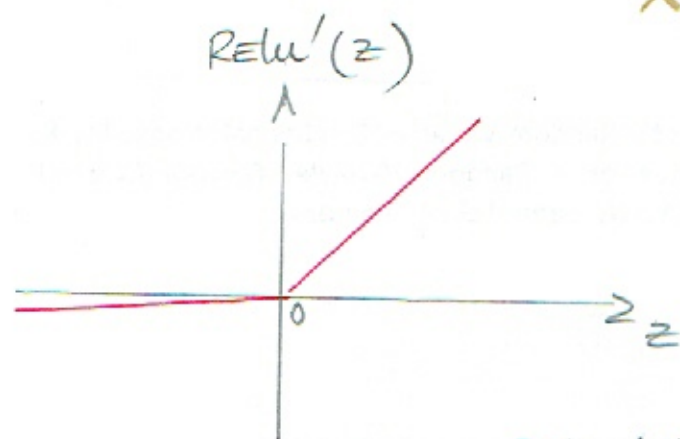
$$\frac{\partial C}{\partial w_1}_{old} = \frac{\partial C}{\partial a_3} \frac{\partial a_3}{\partial a_2} \frac{\partial a_2}{\partial a_1} \frac{\partial a_1}{\partial w_1}_{old}$$

1 x 1 x 1 for  $z > 0$

issue if  $ReLU'(z)$  is 0 or undefined

$\times 0 \rightarrow 0$

↓ DEAD NEURON!

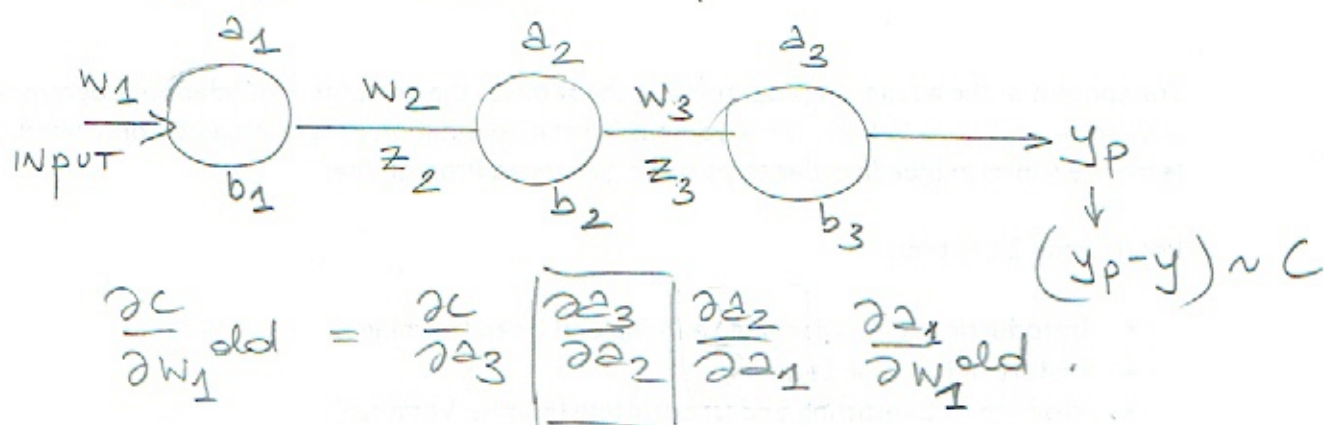


← solution  
leaky ReLU

$$ReLU(z) = \begin{cases} z & z > 0 \\ 0.01 \cdot z & z < 0 \end{cases}$$

# exploding gradients

15



Chain rule diagram:

$$\frac{\partial a_3}{\partial a_2} = \frac{\partial a_3}{\partial z_3} \left( \frac{\partial z_3}{\partial a_2} \right) \rightarrow w_3$$

Annotations:

- $\frac{\partial z_3}{\partial a_2}$  is output  $z_3$  (from the previous layer).
- Chain rule:  $0 \rightarrow 0.25$ .
- Initializer is important to avoid exploding gradient.
- ex:  $w_3 = 500$ .
- deeper  $\rightarrow$  will grow



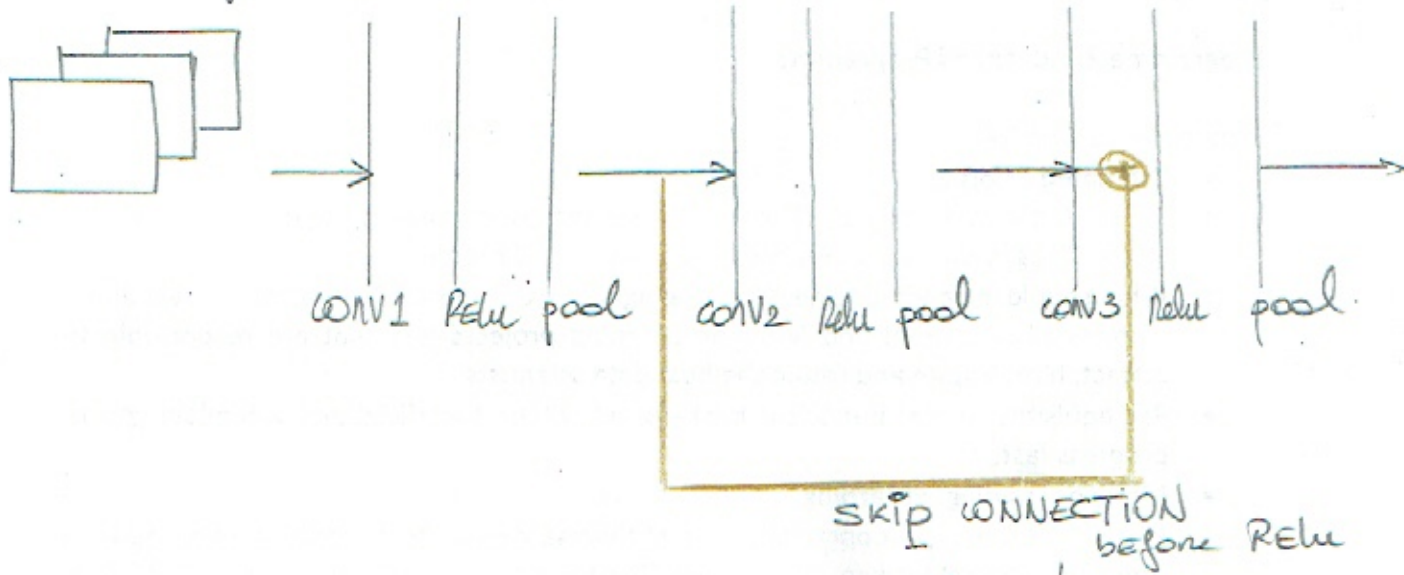
Resnet : 2015 (He et al)

(16)

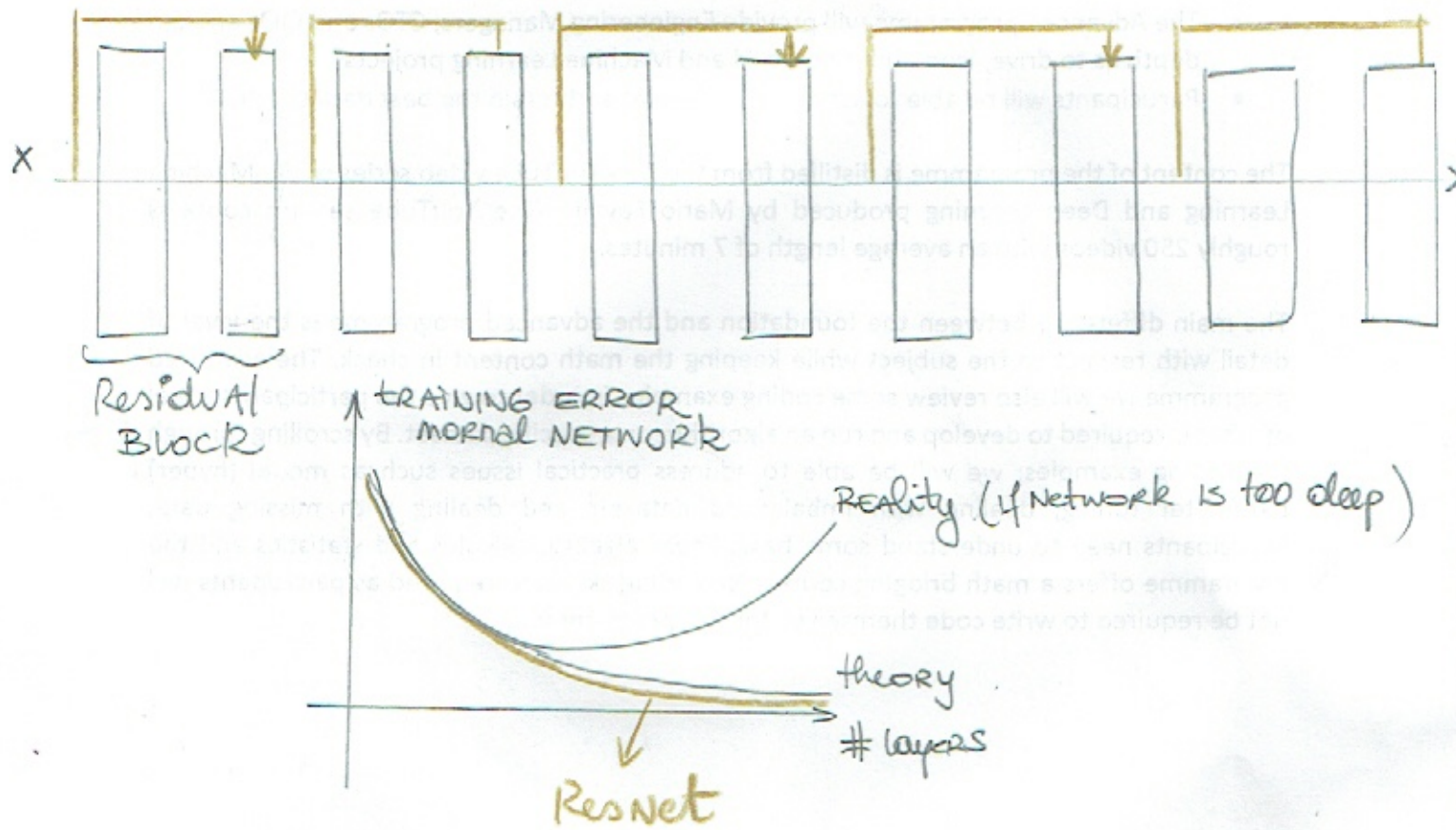
Microsoft  
152 layers

↓  
Residual Block

→ Resnet allows you to train deeper networks while avoiding vanishing or exploding gradients

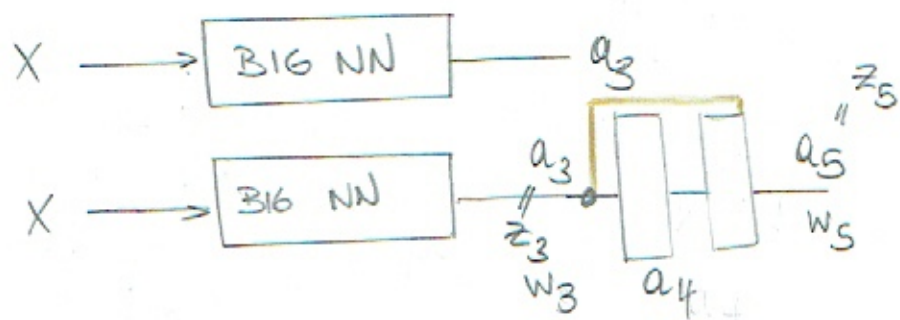


ResNet : → stacks residual blocks together !



Why do ResNets work well?

(14)



We assume that dimensions are same

$$a_5 = \text{relu}(z_5 + z_3)$$

$$= \text{relu}(w_5 a_4 + b_5 + z_3)$$

Make small  
using weight  
decay

$$\rightarrow \text{relu}(z_0 + z_3) = z_3$$

$$\rightarrow a_5 \approx a_3$$

identity function is easy  
to learn for residual block

add  $w_5$  MATRIX  
if dimensions are  
different  
 $\rightarrow$  zero padding