



# A FASTER INTRODUCTION TO IMAGE PROCESSING

... than early « »

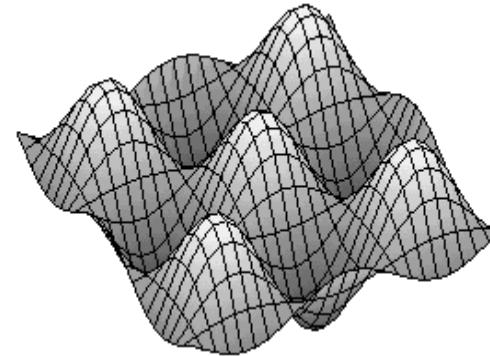
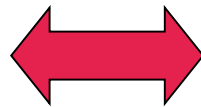
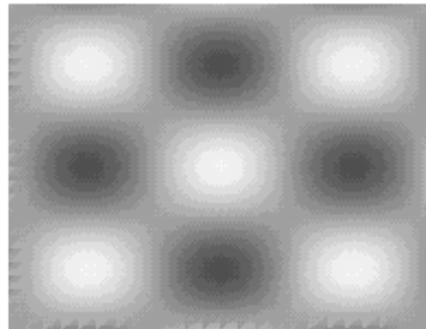
thomas.grenier@insa-lyon.fr

# IMAGE PROCESSING



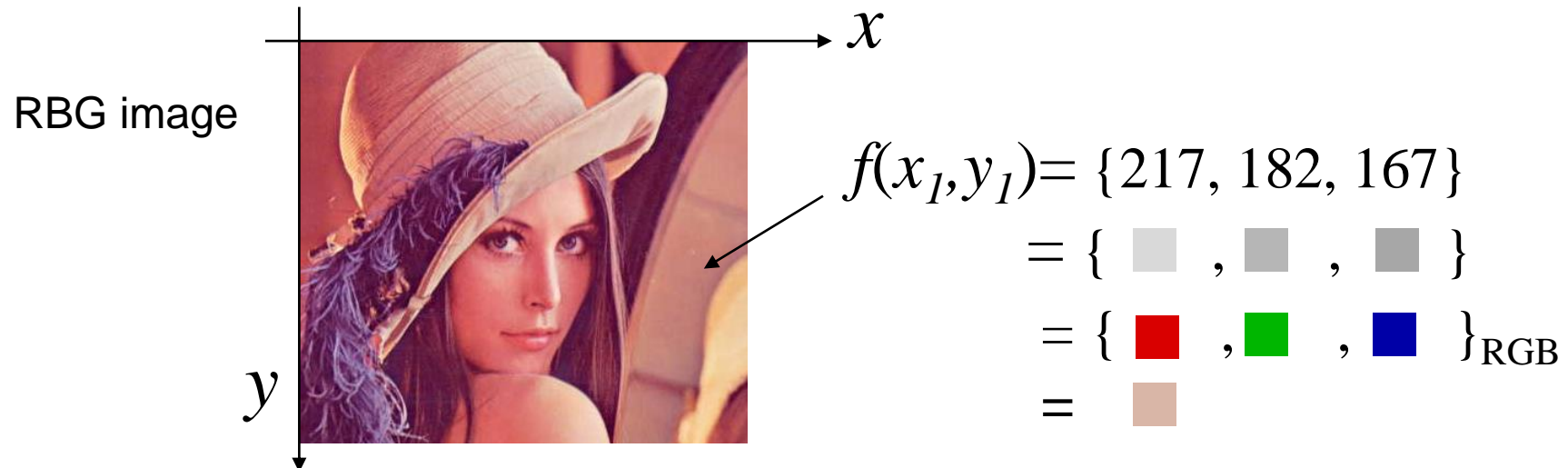
# WHAT IS AN *IMAGE*?

- Image definition
  - An image may be defined as a two-dimensional function,  $f(x,y)$ 
    - $x$  and  $y$  are spatial (plane) coordinates
    - the amplitude of  $f$  at any pair of coordinates  $(x,y)$  is called **intensity** or **gray level** of the image at that point

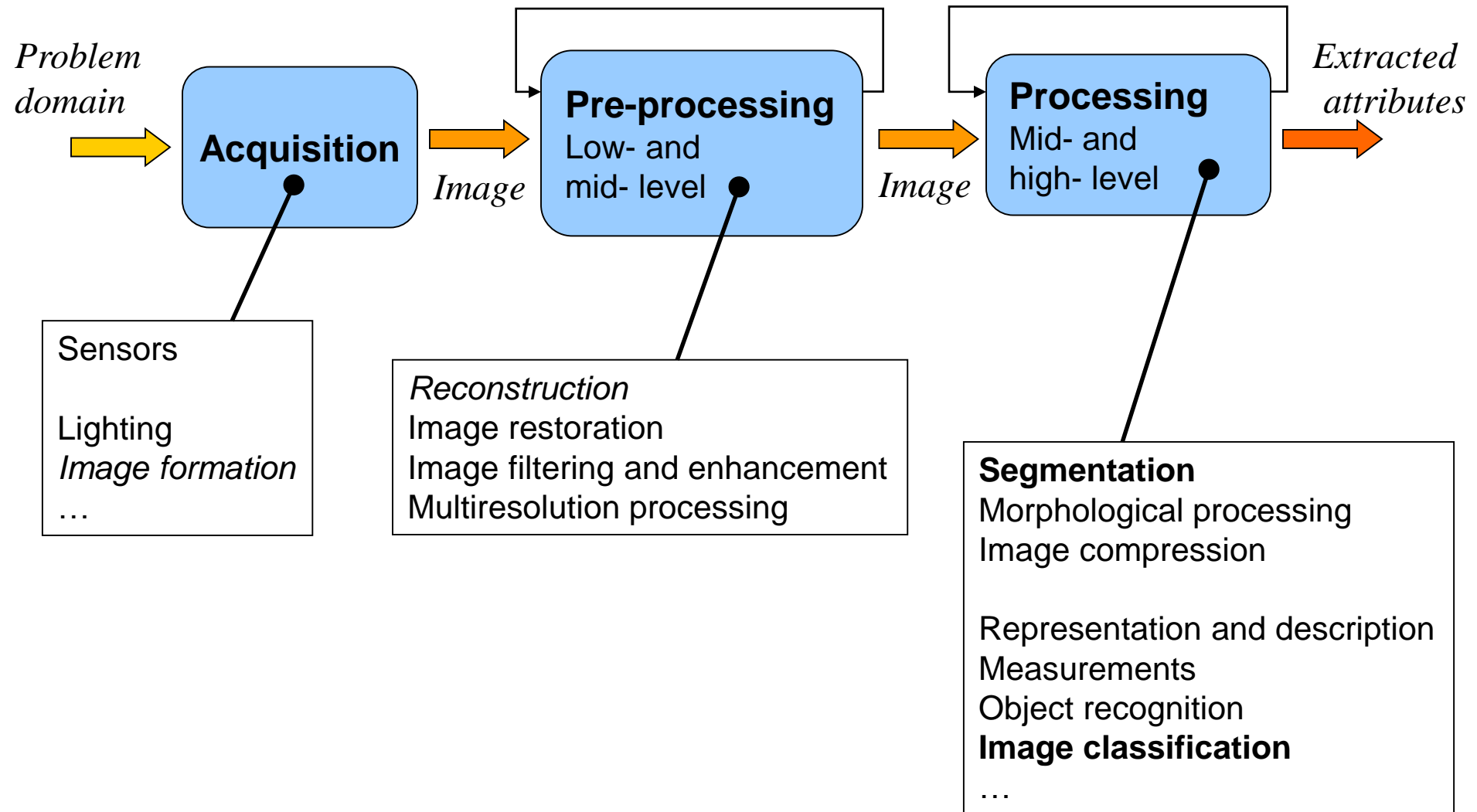


# WHAT IS AN *IMAGE*?

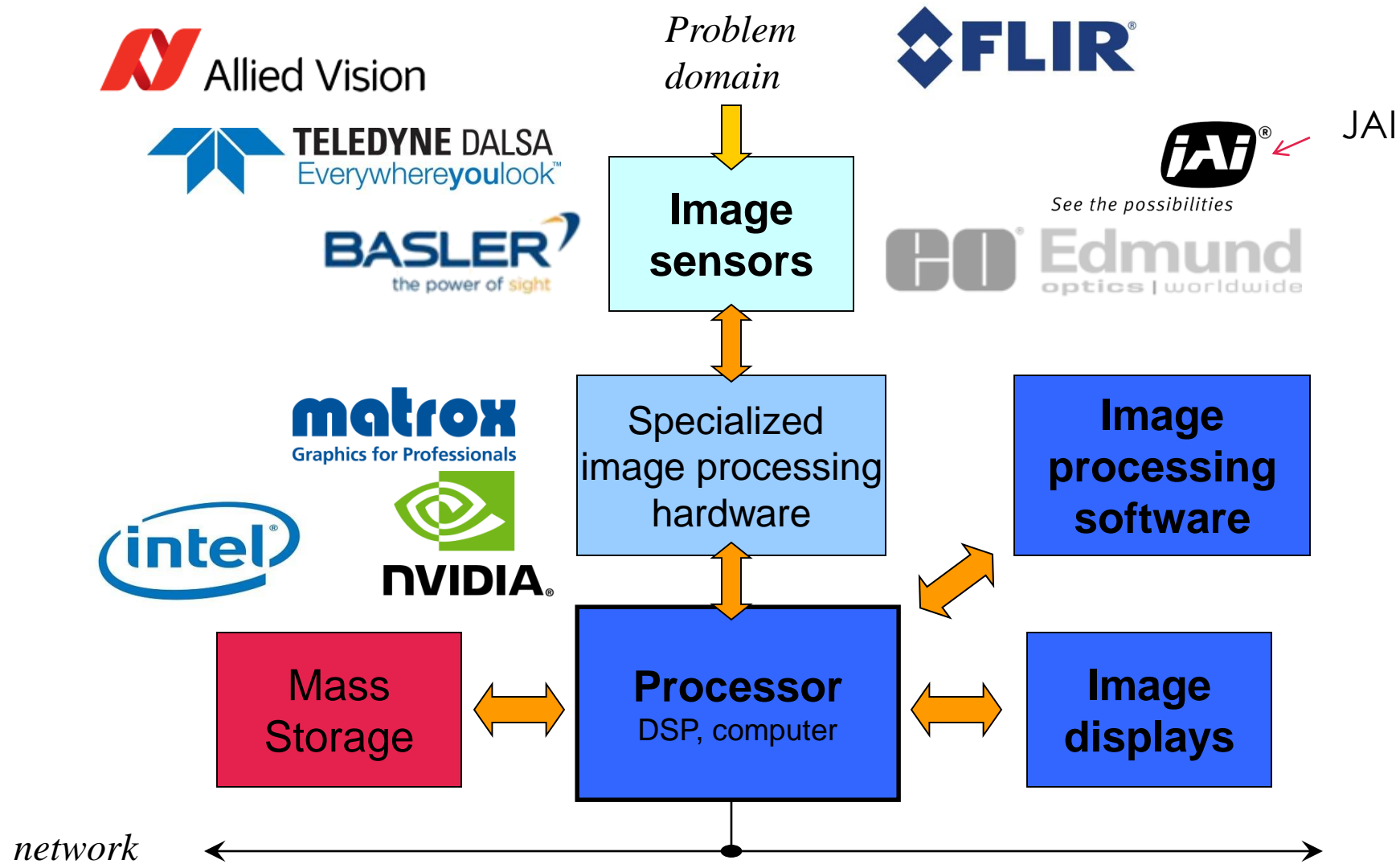
- Digital Image definition
  - The definition of  $f$  may be extended:
    - as a  $n$ -dimensional function,
      - i.e. 3D:  $f(x,y,z)$  or image sequence  $f(x,y,t)$
    - with amplitudes composed as a vector of data,
      - i.e. Color image: 3 components at each point, Complex number



# FUNDAMENTAL STEPS IN DIP



# COMPONENTS OF AN IMAGE PROCESSING SYSTEM



# LINEAR AND NONLINEAR OPERATIONS

- Important classification of an image-processing method
  - Is it a linear or a nonlinear method ?
- Let  $H$  be a general operator

$$H[f(x, y)] = g(x, y)$$

- $H$  is said to be a **linear operator** if

$$\begin{aligned} H[a_i f_i(x, y) + a_j f_j(x, y)] &= a_i H[f_i(x, y)] + a_j H[f_j(x, y)] \\ &= a_i g_i(x, y) + a_j g_j(x, y) \end{aligned}$$

➔ Tools: vector-matrix, PDE, set, and related operators (+, \*, convolution)...

# OPERATIONS ON IMAGES

- Arithmetic Operations  $+, -, \times, \div$ 
  - Application: Corrupted image  $g$  obtained by adding the noise  $\eta$  to a noiseless image  $f$

$$g(x, y) = f(x, y) + \eta(x, y)$$

- assumptions :
  - at every pair of coordinates  $(x, y)$  the noise is uncorrelated
  - the noise has zero average value
- ➔ Noise reduction by summing (averaging) a set of noisy images

$$\sigma_{\bar{g}(x, y)} = \frac{1}{\sqrt{K}} \cdot \sigma_{\eta(x, y)}$$



# Application : Noise Reduction

Original



$\sigma = 3,55$

Mean



1  $\sigma = 25,45$



2  $\sigma = 18,3$



4  $\sigma = 13,5$



6  $\sigma = 11,6$

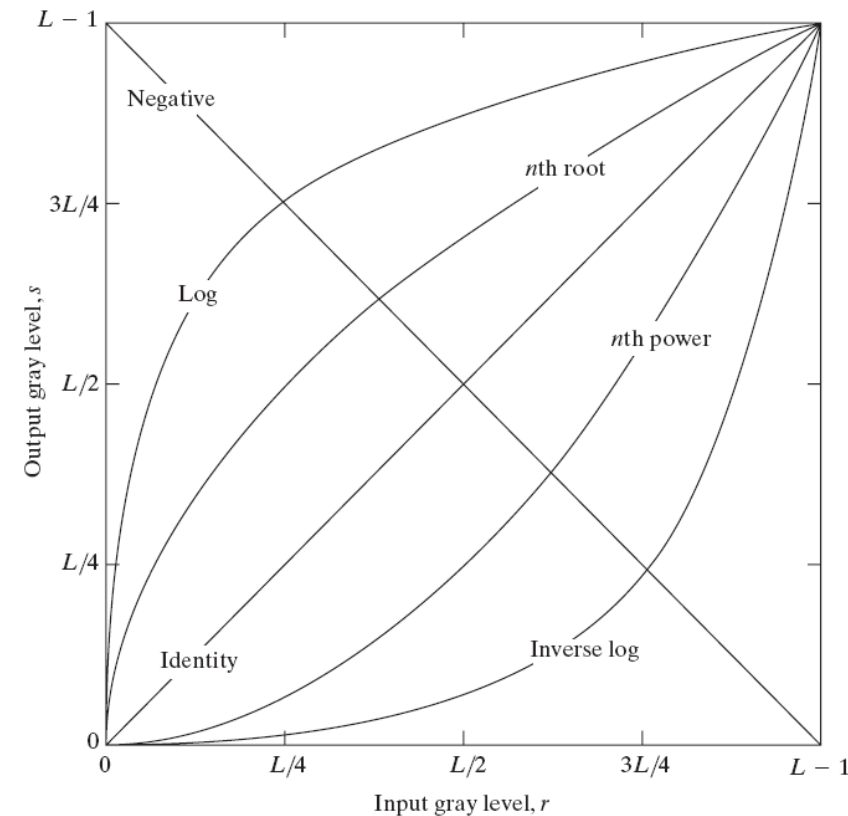
# OPERATIONS ON IMAGES

- Spatial Operations - Single pixel operations
  - $r$ : original intensity,
  - $s$ : new intensity,
  - $T$ : a transformation function

$$s = T(r)$$

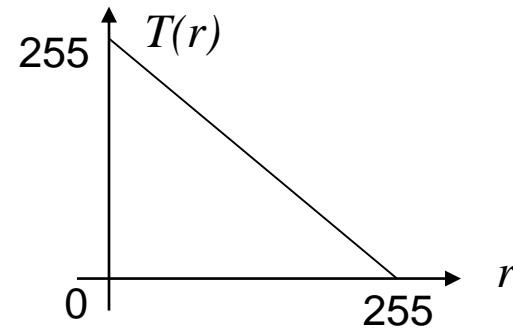
*Some basic grey-levels transformation functions used for image enhancement*

Range of intensities:  
 $[0, L-1]$



# OPERATIONS ON IMAGES

- Spatial Operations - Single pixel operations  
Example



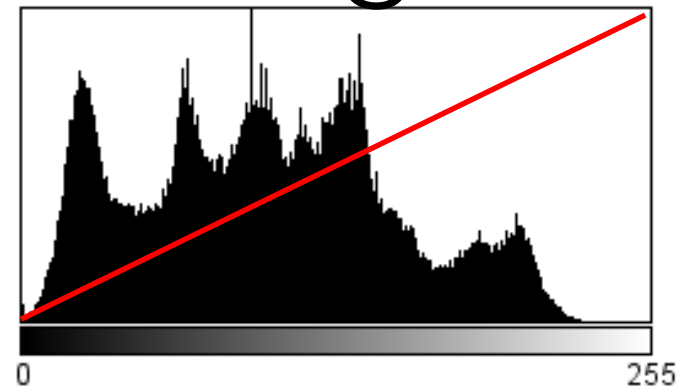
Original



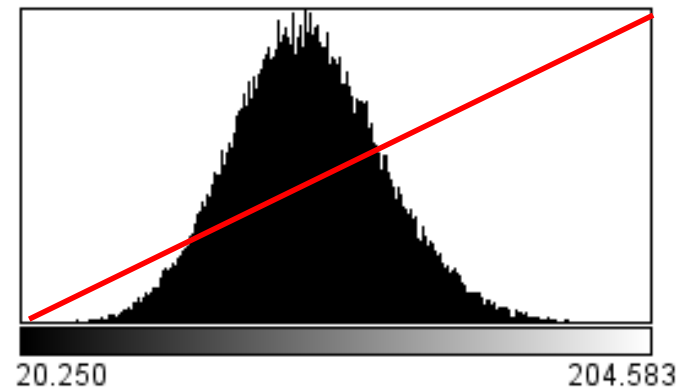
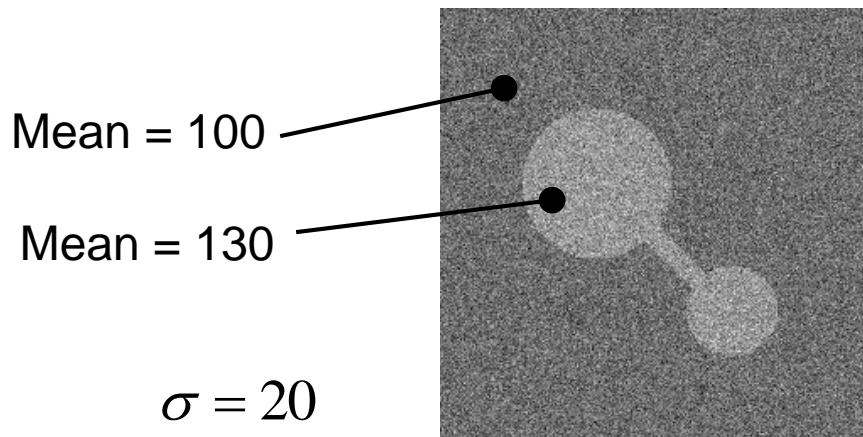
Negative

# Operations on images : Histogram

- Examples

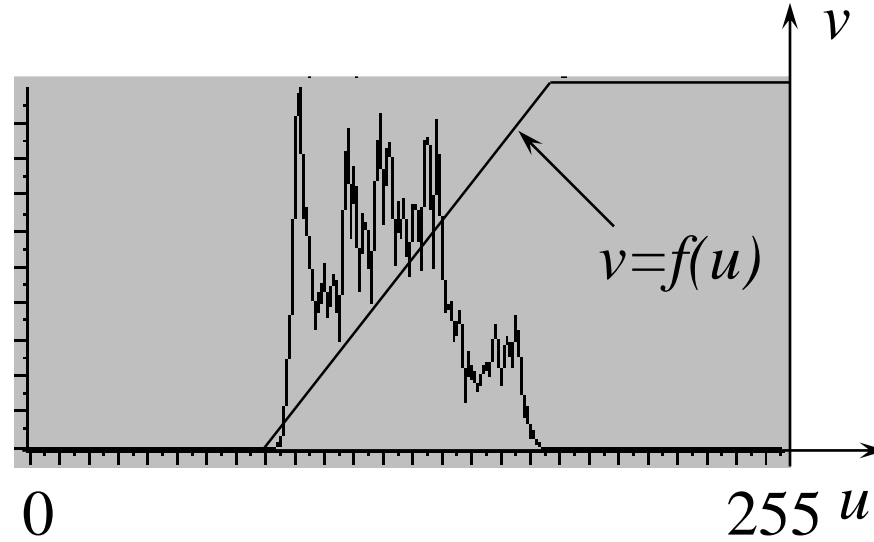


Count: 262144      Min: 0  
Mean: 99.434      Max: 243  
StdDev: 52.585      Mode: 93 (2760)

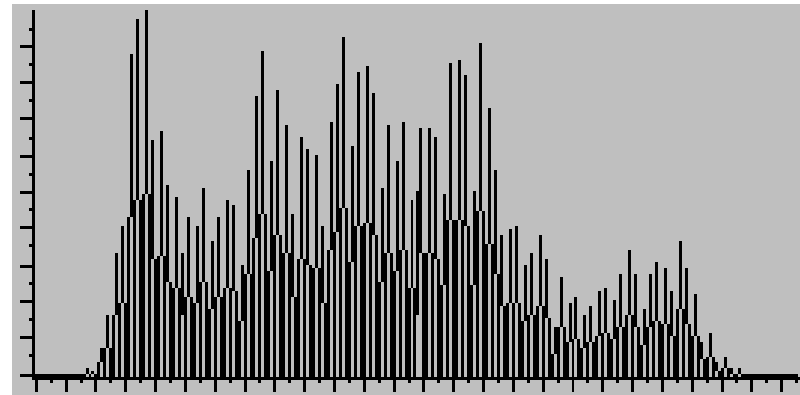
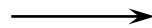


Count: 65536      Min: 20.250  
Mean: 104.219      Max: 204.583  
StdDev: 22.501      Mode: 103.416 (908)  
Bins: 256      Bin Width: 0.720

# HISTOGRAM MANIPULATION



$$v=f(u)$$

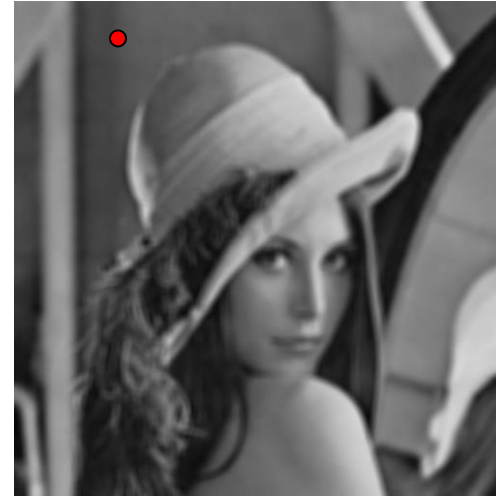
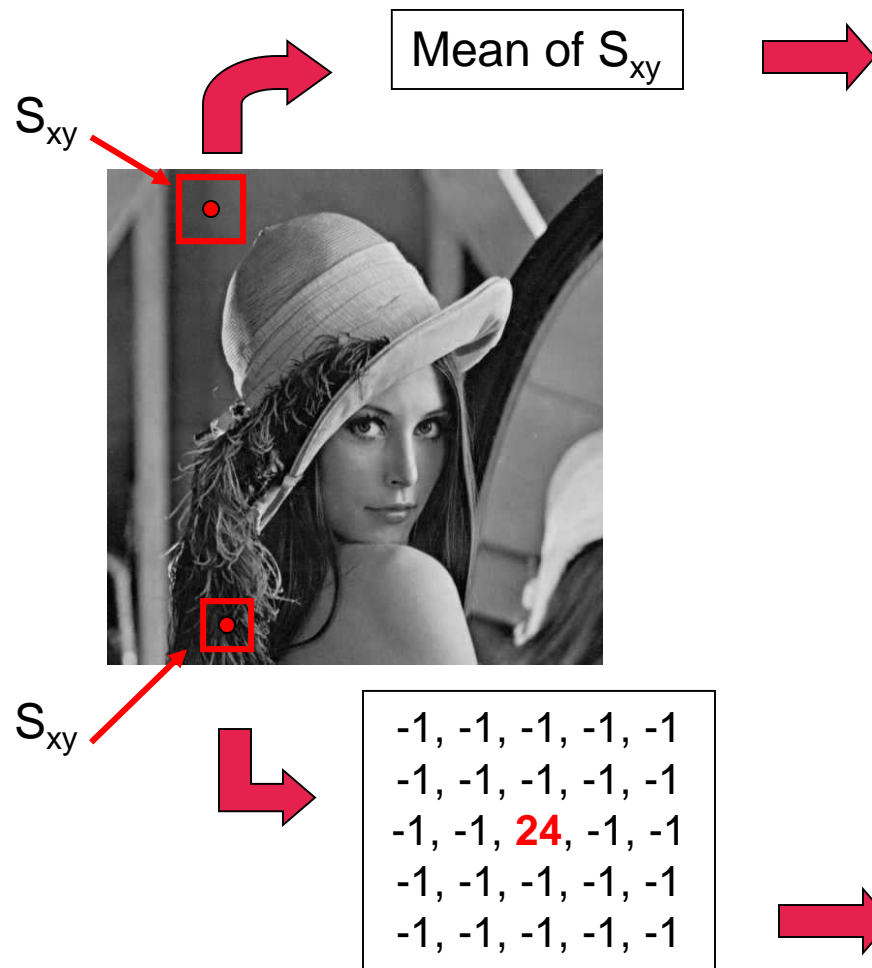


# OPERATIONS ON IMAGES

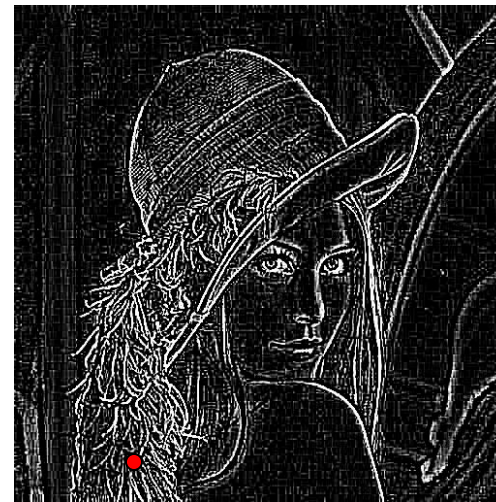
- Spatial Operations - Neighborhood operations
  - $S_{xy}$ : set of coordinates of a neighborhood centered on a point  $(x,y)$  in an image  $f$ .
  - Neighborhood processing generates one corresponding pixel in the output image  $g$  at the same  $(x,y)$  coordinates.
  - The value of that pixel in  $g$  is determined by a operation involving the pixels in  $S_{xy}$ .

# OPERATIONS ON IMAGES

- Spatial Operations - Neighborhood operations



$$g(x, y) = \frac{1}{m.n} \sum_{(r,c) \in S_{xy}} f(r, c)$$



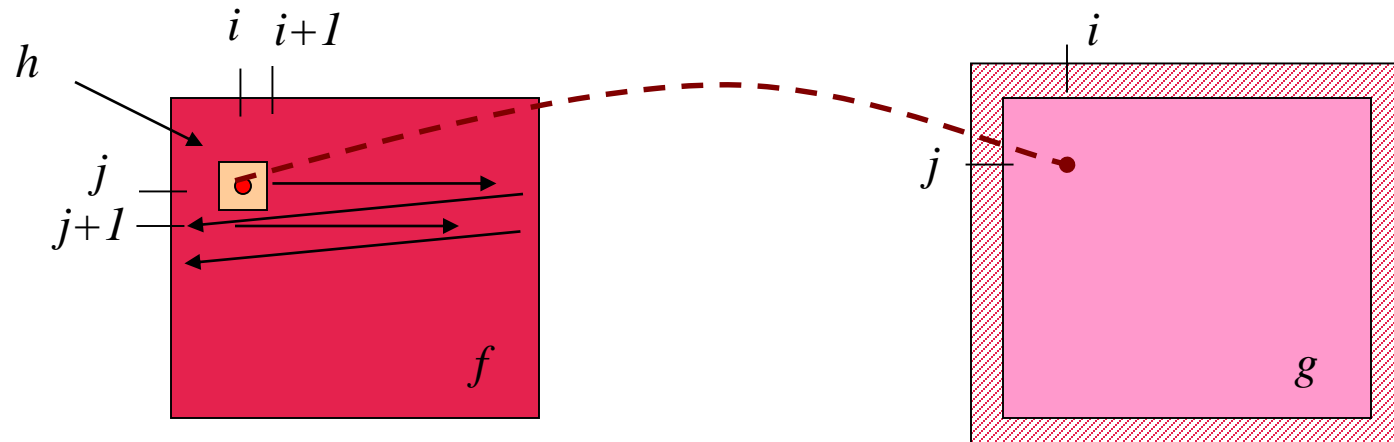
FIR filters...

# OPERATIONS ON IMAGES

- Convolution  $g(x,y) = h(x,y)*f(x,y)$  (two-dimensional convolution)

$$g(i, j) = \sum_{(k,l) \in H} \sum h(k, l) f(i-k, j-l)$$

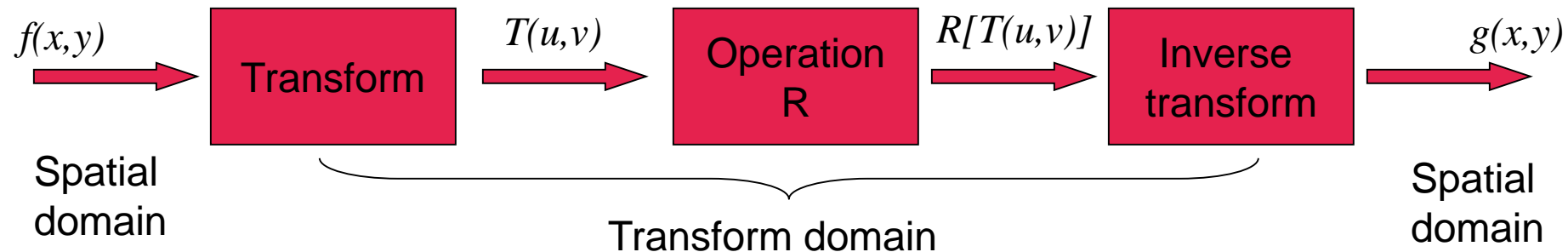
Output image  $\nearrow$   $g(i, j)$   $\nwarrow$  Convolution mask / Impulse response  $h(k, l)$   $\swarrow$  Input image  $f(i-k, j-l)$





# IMAGES TRANSFORMATION

- Previous methods work in spatial domain
  - In some cases, image processing tasks are best formulated in a transform domain.
    - i.e. frequency  $\rightarrow$  Fourier
- $\rightarrow$  Many transformations exist



# IMAGES TRANSFORMATION

- A particularly important class of 2D linear transforms can be expressed in the general form

$$T(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) \cdot r(x, y; u, v)$$

forward transform      Input image      Forward transformation kernel

$$f(x, y) = \sum_{u=0}^{M-1} \sum_{v=0}^{N-1} T(u, v) \cdot s(x, y; u, v)$$

Recovered image      Forward transform      inverse transformation kernel

Separable kernel:  $r(x, y, u, v) = r_1(x, u) \cdot r_2(y, v)$

Symmetric kernel:  $r(x, y, u, v) = r_1(x, u) \cdot r_1(y, v)$

# IMAGES TRANSFORMATION, DFT

- 2-D Discrete Fourier Transform (**DFT**)

Forward kernel  $r(x, y, u, v) = e^{-j2\pi(ux/M+vy/N)}$

Inverse kernel  $s(x, y, u, v) = \frac{1}{MN} e^{j2\pi(ux/M+vy/N)}$

$$T(u, v) = \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} f(x, y) \cdot e^{-2j\pi(ux/M+vy/N)}$$

$$f(x, y) = \frac{1}{MN} \sum_{x=0}^{M-1} \sum_{y=0}^{N-1} T(u, v) \cdot e^{2j\pi(ux/M+vy/N)}$$

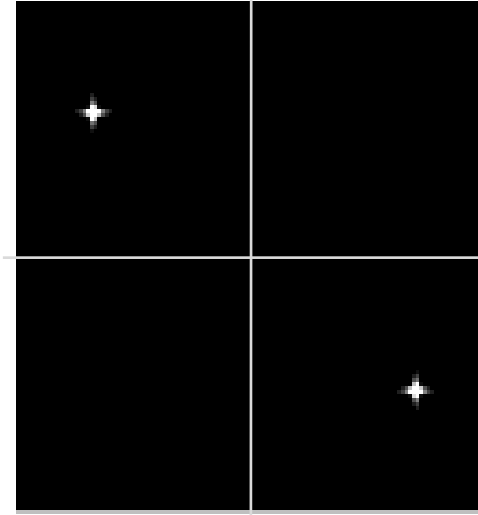
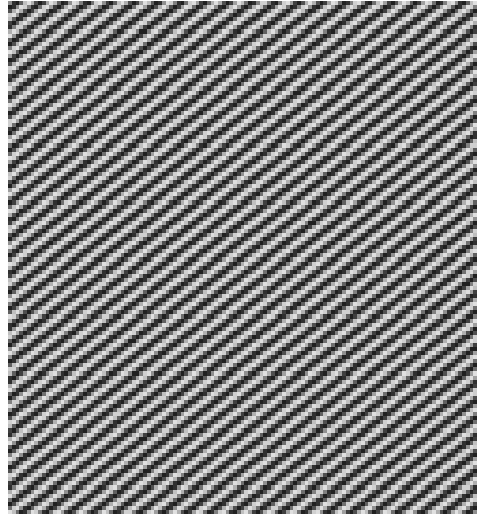
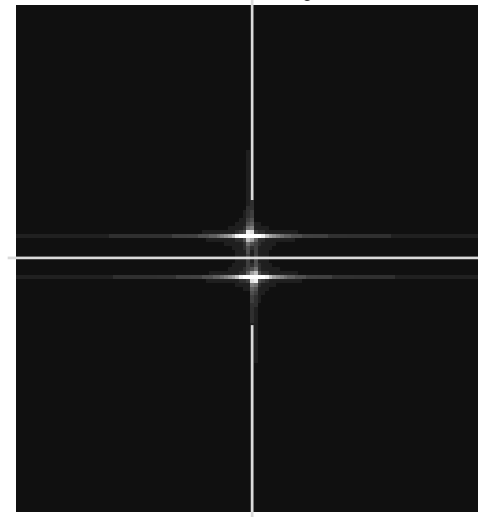
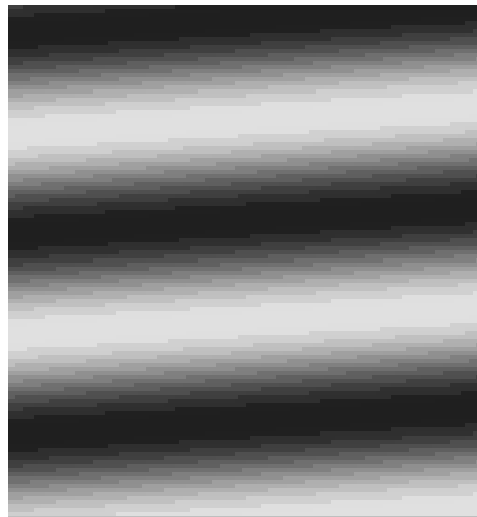
→ Complex numbers...

→ Modulus and phase (angle)

Sine images



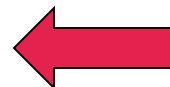
2 Dirac delta functions

 $f_i$ High  
frequency $f_j$  $f_i$ Low  
frequency $f_j$

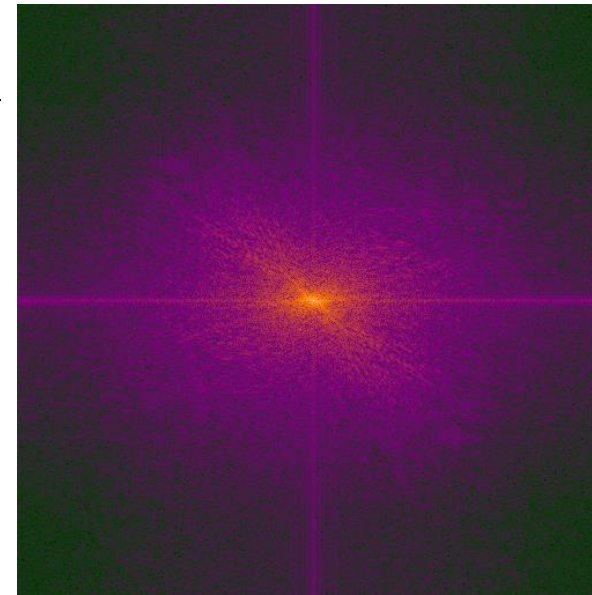
# Image example: 2-D Fourier transform



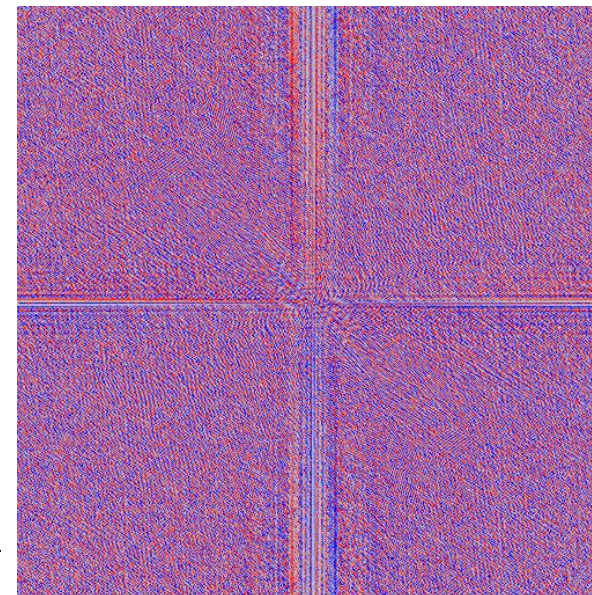
DFT



DFT<sup>-1</sup>

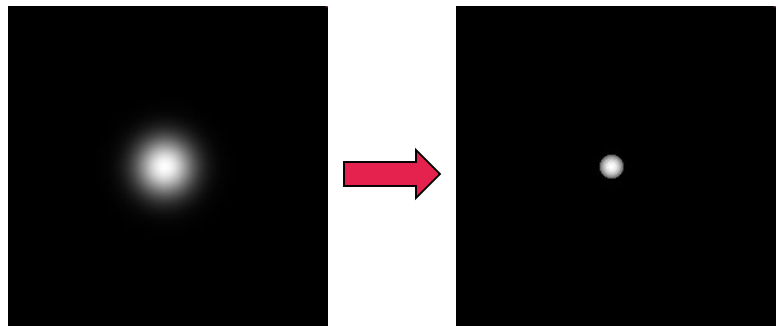
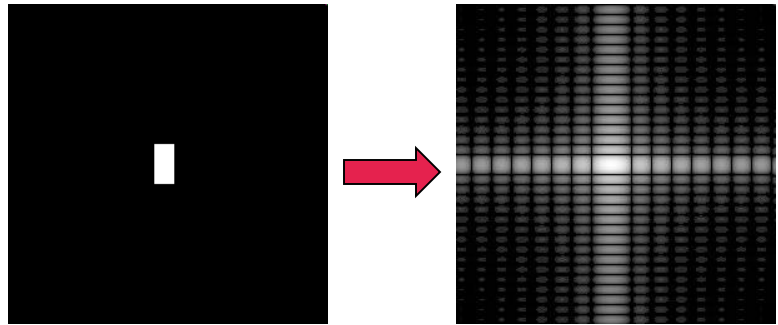
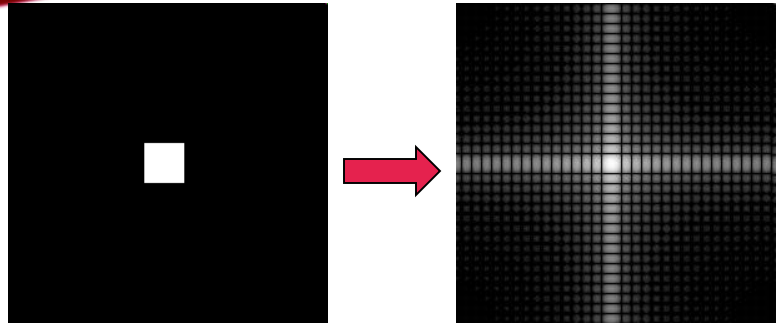


modulus

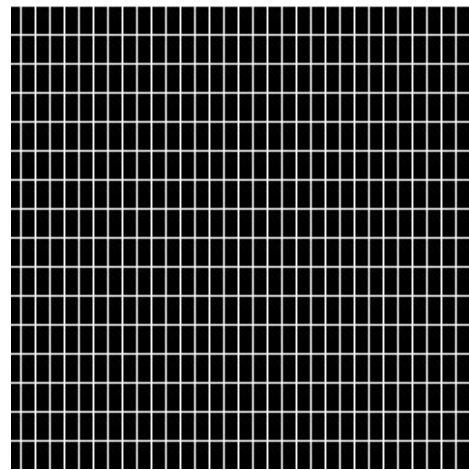
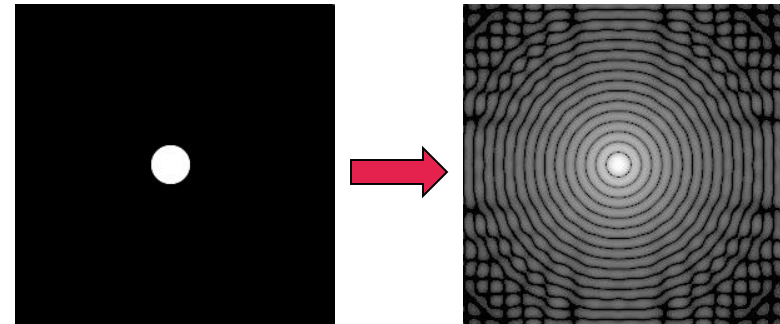
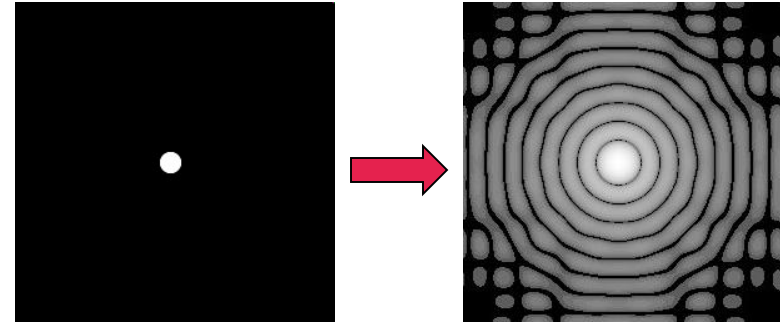


phase

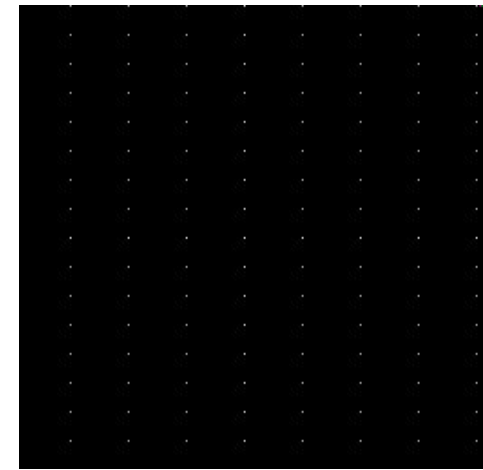
# Some basic DFTs



gauss

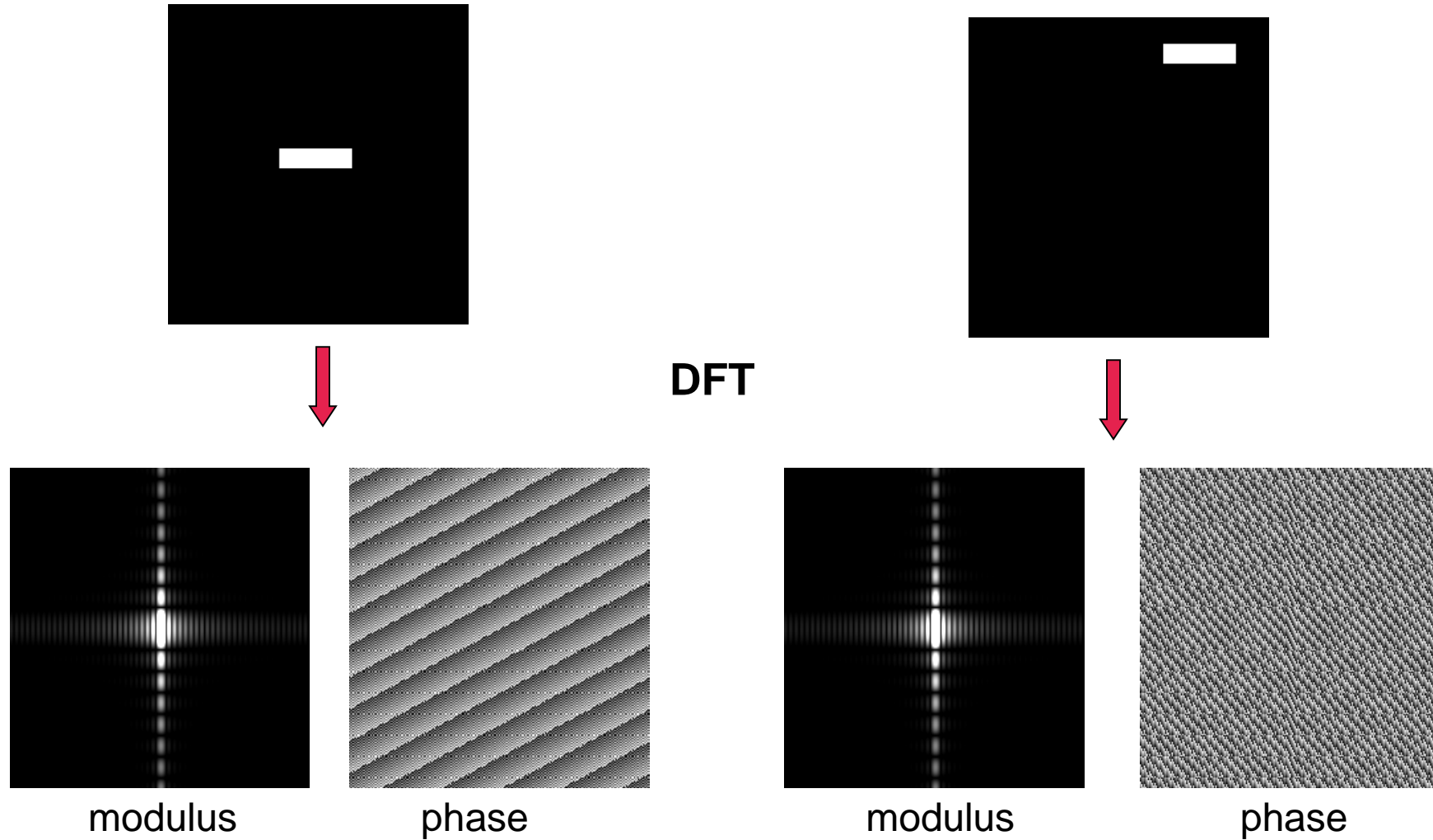


grid



Weighted 2D Dirac comb

# Notes on Phase influence



# GEOMETRIC AND SPATIAL TRANSFORMATIONS

- Spatial transformations
  - Example
    - Shrink image to half its size

$$(x', y') = T\{(x, y)\} = (x/2, y/2)$$

- **Affine transform:**

$$\mathbf{x}' = \mathbf{A} \cdot \mathbf{x} + \mathbf{t}$$

$$[x', y', 1] = [x, y, 1] \cdot \begin{bmatrix} t_{11} & t_{12} & 0 \\ t_{21} & t_{22} & 0 \\ t_{31} & t_{32} & 1 \end{bmatrix}$$

- Homographic transform
  - Two cameras looking points on a plane :  $x, y, \dots, z$  !
- Higher order

$$[x', y', 1] = [x, y, x^2, y^2, xy, \dots, 1] \cdot \mathbf{T}$$



## Affine transform

identity  $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x' = x \\ y' = y \end{cases}$

scaling  $\begin{bmatrix} c_x & 0 & 0 \\ 0 & c_y & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x' = c_x \cdot x \\ y' = c_y \cdot y \end{cases}$

translation  $\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ t_x & t_y & 1 \end{bmatrix} \Rightarrow \begin{cases} x' = x + t_x \\ y' = y + t_y \end{cases}$

Shear  
(vertical)  $\begin{bmatrix} 1 & 0 & 0 \\ s_x & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x' = x + s_x \cdot y \\ y' = y \end{cases}$

Shear  
(horizontal)  $\begin{bmatrix} 1 & s_y & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \Rightarrow \begin{cases} x' = x \\ y' = y + s_y \cdot x \end{cases}$

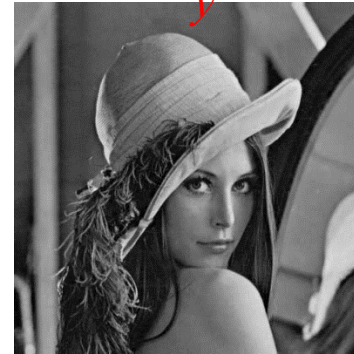
identity



$x'$

$y'$

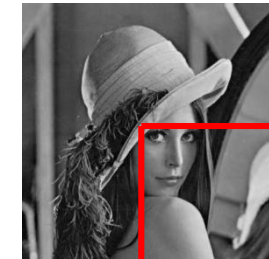
translation



$x'$

$y'$

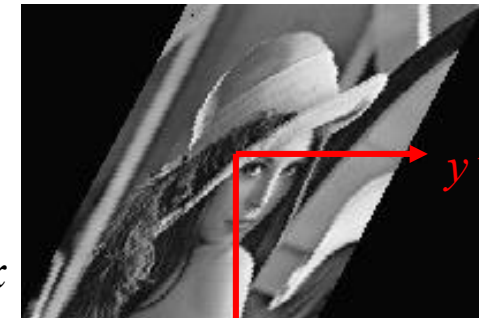
scaling



$x'$

$y'$

horizontal shearing



$x'$

$y'$

$$\begin{cases} x' = x \\ y' = y + 0.5x \end{cases}$$

- Affine transform: rotation

Rotation

$$T = \begin{bmatrix} \cos \theta & \sin \theta & 0 \\ -\sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

35° degrees rotation (from image center)



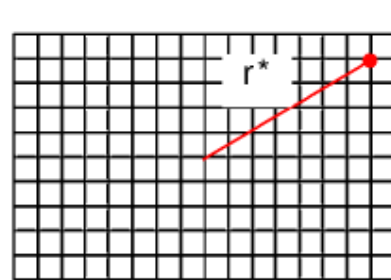
Remark : interpolation is necessary to avoid aliasing

- Higher order transforms

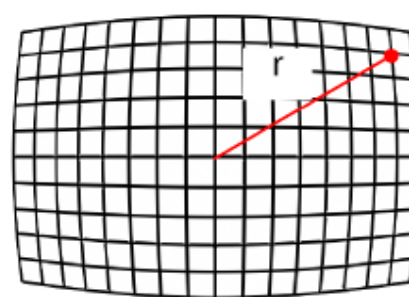


Applications :

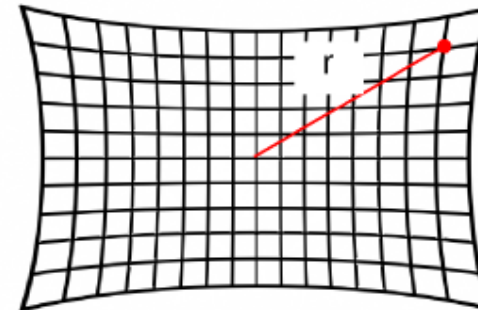
Lens distortion correction, perspective



Orthoscopic  
projection



Barrel  
distortion



Pincushion  
distorsion

# APPLICATION : IMAGE FILTERING



# NOISE CLEANING

- Many types of noise...



Gaussian Noise  
(sd 25)



Salt and pepper noise

# LINEAR FILTERING

- Mean filter

$$g(i, j) = \sum_{(k,l) \in W} h(k, l) f(i - k, j - l)$$

W: 25 neighbors

$H = (1/25) \cdot I$

Result on Gaussian noise



Result on Salt & Pepper noise



# NON LINEAR FILTERING

- Median filter

Result on Gaussian noise

W=25



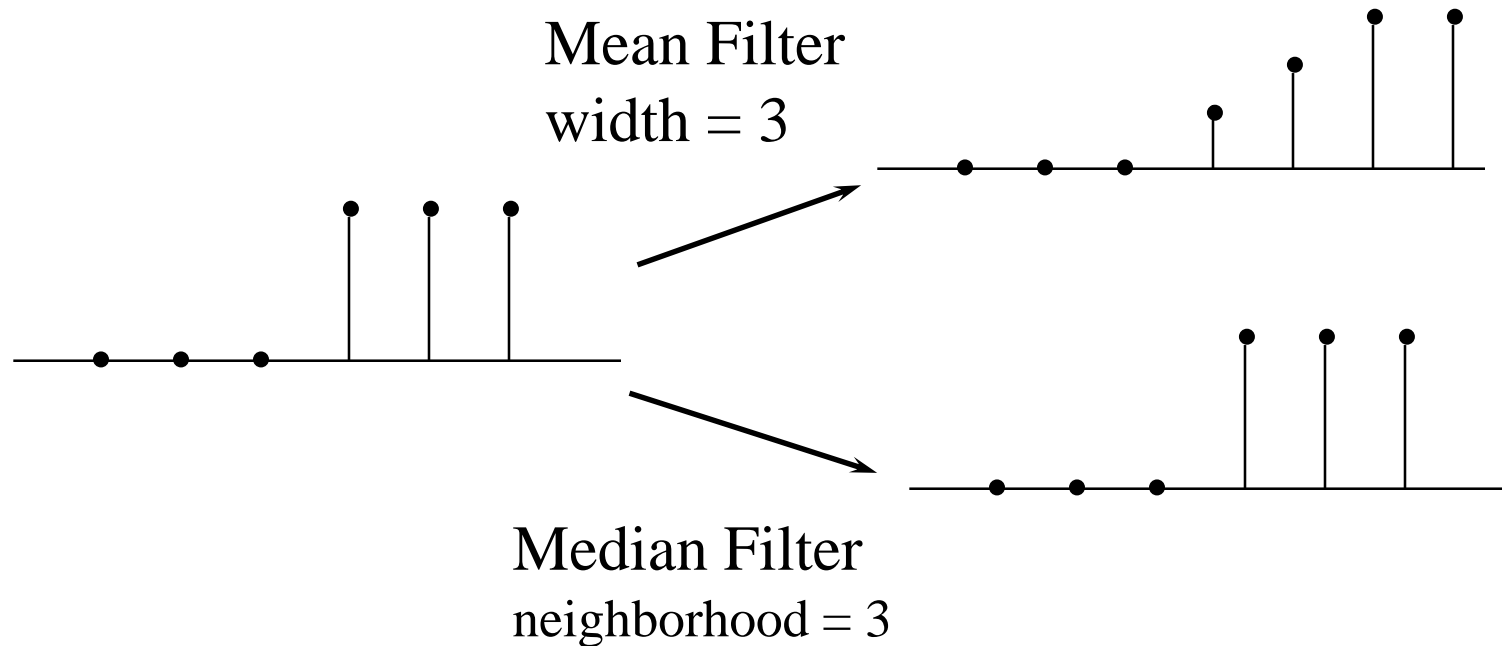
W=9

Result on salt & pepper noise



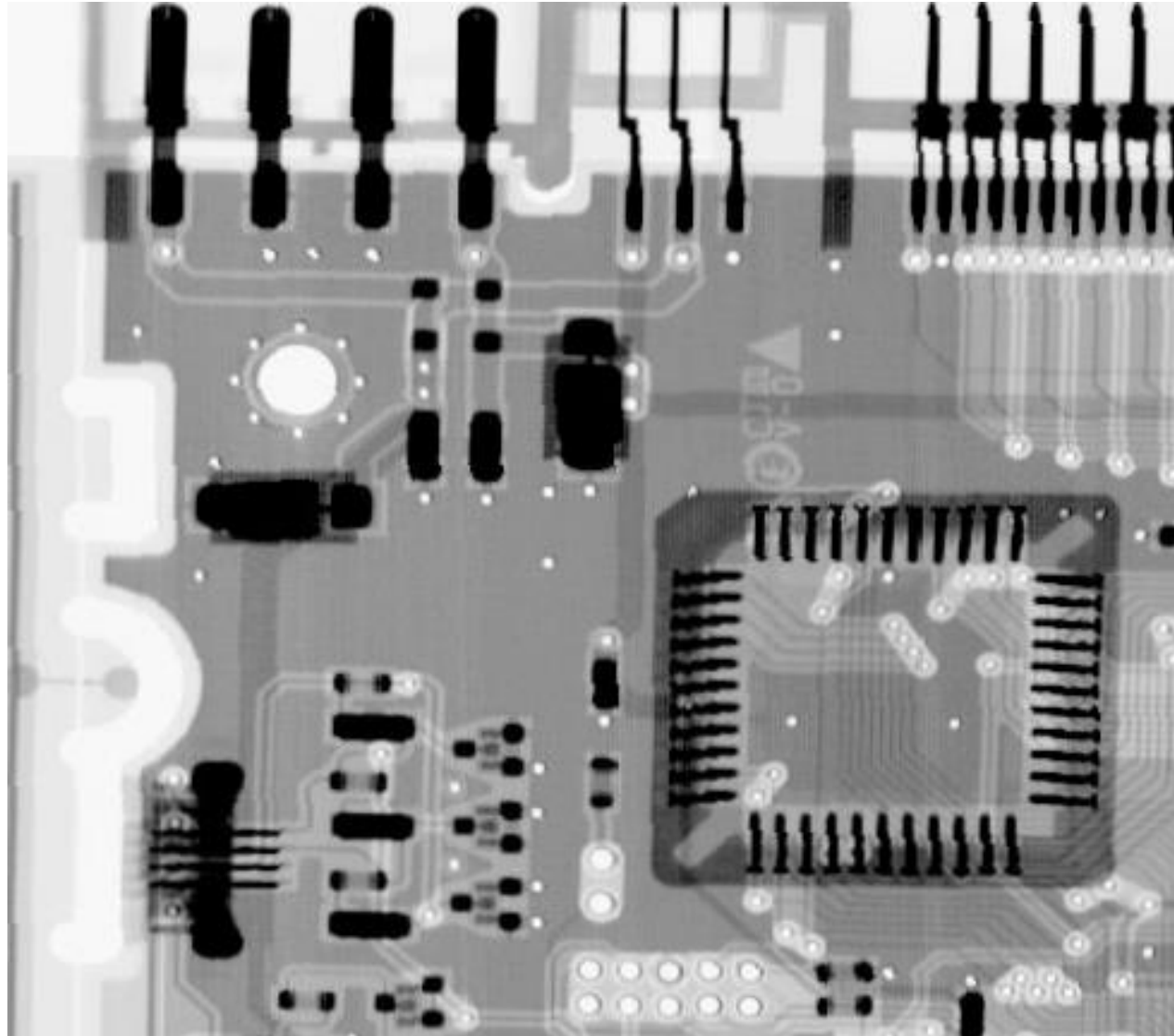
# NON LINEAR FILTERING

- Median filter
  - Advantage of median filtering over linear filtering: edges are preserved



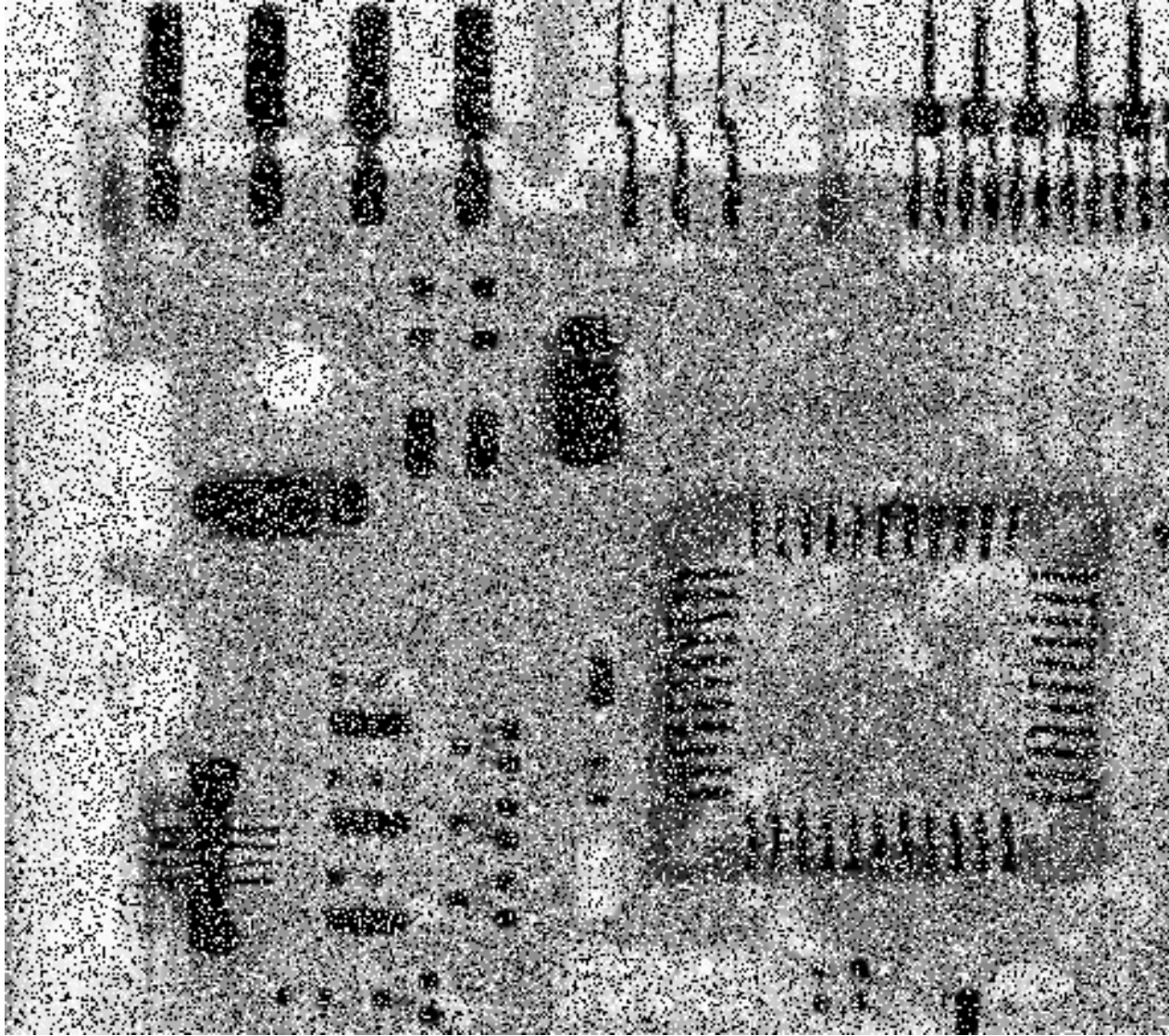


- Adaptive Median filter : Reference

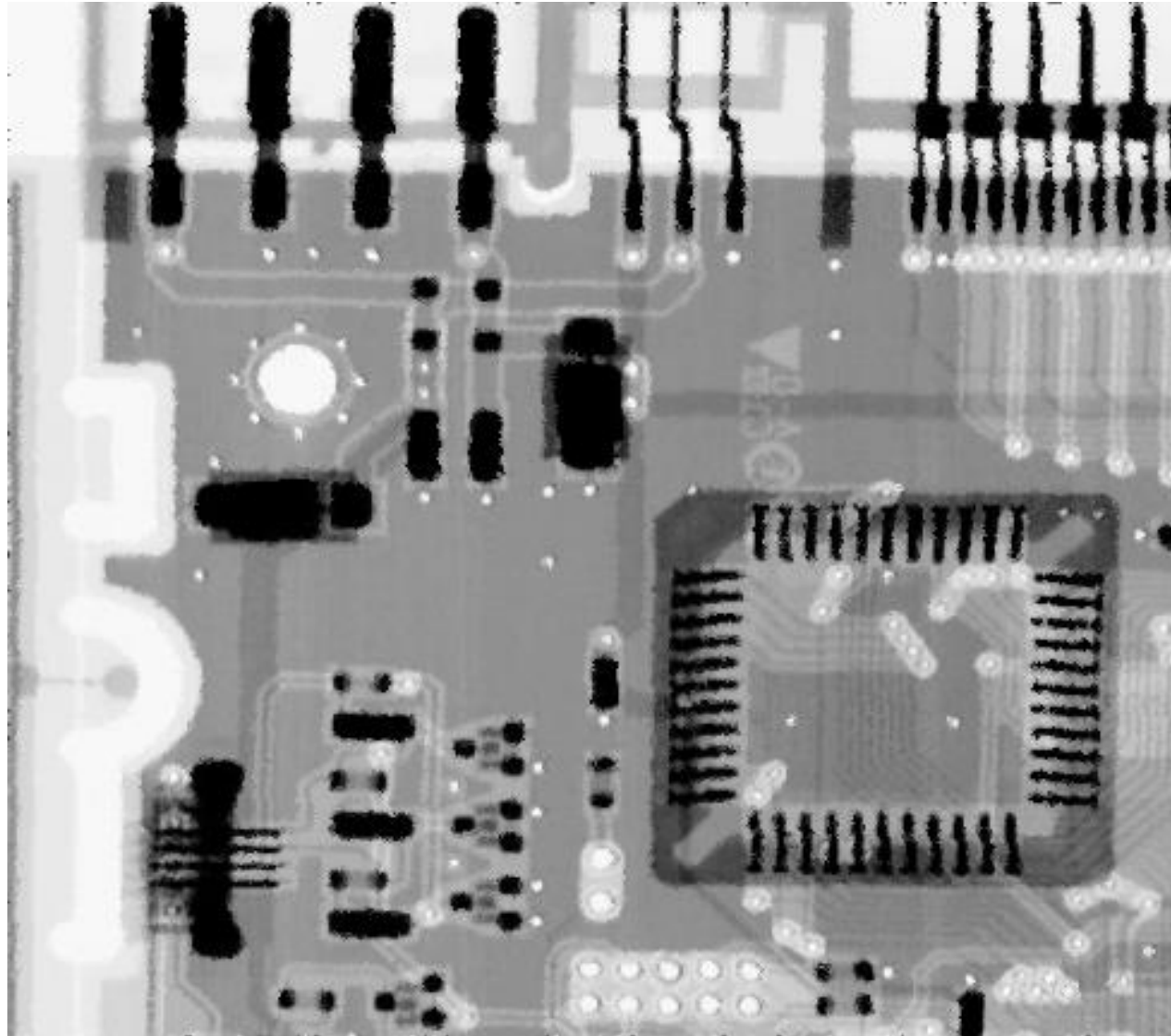


- Adaptive Median filter : Input

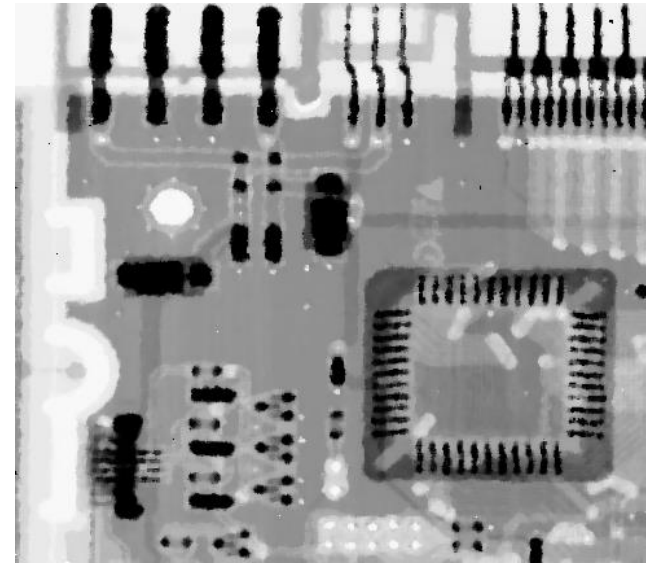
Noisy image (S&P, 0.35)



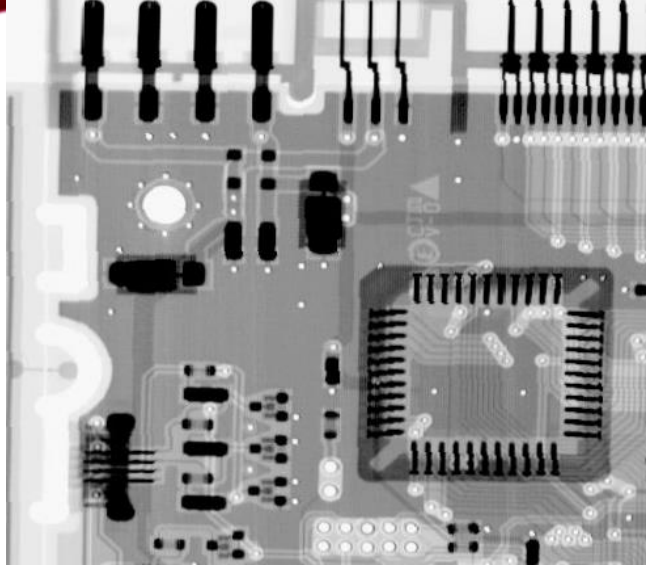
- **Adaptive Median filter**



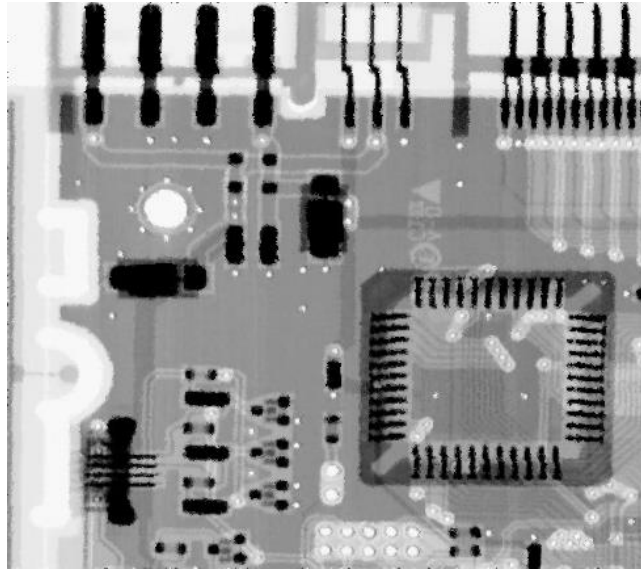
Median Filter (5x5)



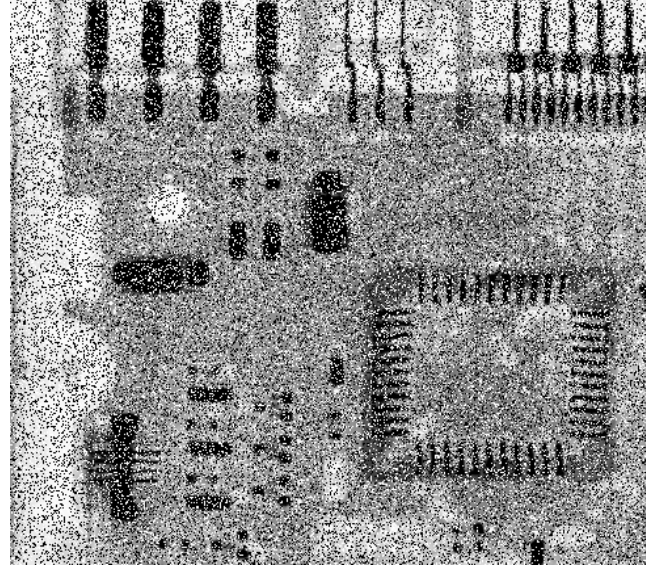
Reference



- Result of Adaptive Median filter

Adaptive Median ( $S_{\max}$  9x9)

Noisy image (S&amp;P, 0.35)



# RESTORATION: PERIODIC NOISE REDUCTION

$$g(x, y) = f(x, y) + \eta(x, y)$$

- By frequency domain filtering
  - Bandreject filter
  - Notch filter  $\rightarrow$  optimum notch

Build  $H_{NP}$  (Notch Pass) by placing a notch pass filter at the location of each spike.  
Interference noise pattern is:

$$N(u, v) = H_{NP}(u, v) \cdot G(u, v)$$

then  $\eta(x, y) = FT^{-1}[N(u, v)]$

thus  $\hat{f}(x, y) = g(x, y) - w(x, y) \cdot \eta(x, y)$

Estimate of  $f(x, y)$

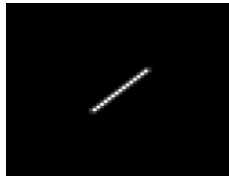
Weighted function (minimizes the effect of  
components not present in the estimate of  $\eta$ )

**$\rightarrow$  How to select  $w(x, y)$  ?**

# RESTORATION



Reference Image  $L$



Shift filter

$f$



Noise

$n$



Corrupted image  $I$

$$I = f \otimes L + n$$

$$(I, f) \xrightarrow{?} L$$

# HOW TO 'DECONVOLUTION' ?

$$\hat{F}(u, v) = \frac{G(u, v)}{H(u, v)}$$

$$\hat{F}(u, v) = F(u, v) + \frac{N(u, v)}{H(u, v)}$$

- Inverse filtering
  - Without noise
  - **With noise**
    - We have to know N!
    - What happens for small values of  $H(u, v)$  ?

## → Solutions

- Minimum Mean Square Error (Wiener) Filtering
- Constrained Least Squares Filtering
- **When H is unknown** : blind deconvolution

Many other approaches exist, ie the 'blind' ones

# WIENER (MINIMUM MSE) AND CLS

- Wiener (Fourier), Minimum MSE

$$\operatorname{argmin}_{\hat{L}} = E\{(L - \hat{L})^2\} \quad \longrightarrow \quad L(u, v) = \frac{1}{f(u, v)} \frac{|f(u, v)|^2}{|f(u, v)|^2 + \frac{S_n(u, v)}{S_L(u, v)}} I(u, v)$$

- Constrained Least Square (Fourier)

$$\operatorname{argmin}_{\hat{L}} = \|I - f \otimes \hat{L}\|_2^2 + \gamma \|\Delta \hat{L}\|_2^2 \quad \longrightarrow \quad L(u, v) = \frac{1}{f(u, v)} \frac{|f(u, v)|^2}{|f(u, v)|^2 + \gamma |P(u, v)|^2} I(u, v)$$

$f(u, v)$  = degradation function in the Fourier space

$|f(u, v)|^2 = f^*(u, v) \cdot f(u, v)$ , with  $f^*(u, v)$  the complex conjugate of  $f$

$S_n(u, v) = |N(u, v)|^2$  = noise power spectrum

$S_L(u, v) = |L(u, v)|^2$  = power spectrum of undegraded image

$\gamma$  = parameter to tune

$$p(x, y) = \begin{vmatrix} 0 & -1 & 0 \\ -1 & 4 & -1 \\ 0 & -1 & 0 \end{vmatrix}$$



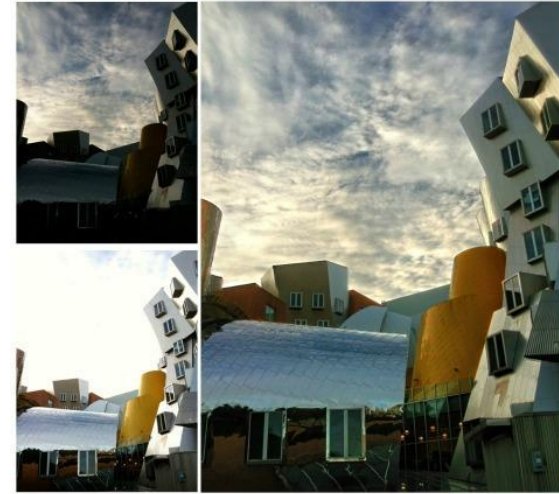
# APPLICATION : IMAGE SEGMENTATION



# SEGMENTATION

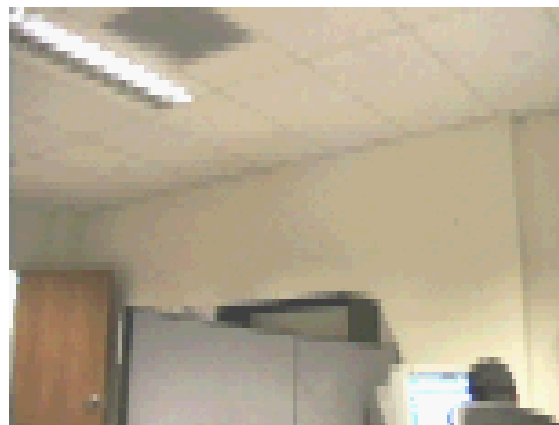
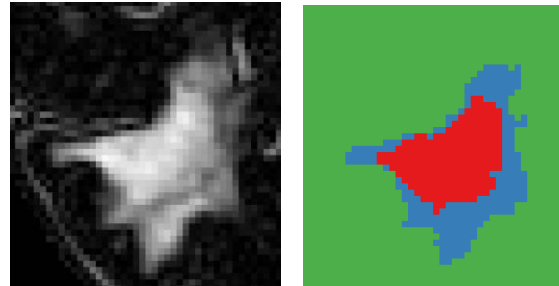
- Visualization
- Counting
- Identification
- Measurements (shape, volumes ...)
- Tracking

*HDR mode*



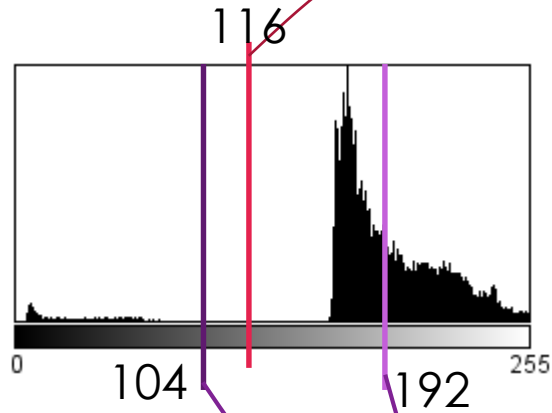
Top-left: iPhone photo #1. Bottom-left: iPhone photo #2. Right: TrueHDR result.

*IRM SEP*



*[Comaniciu 2003]*

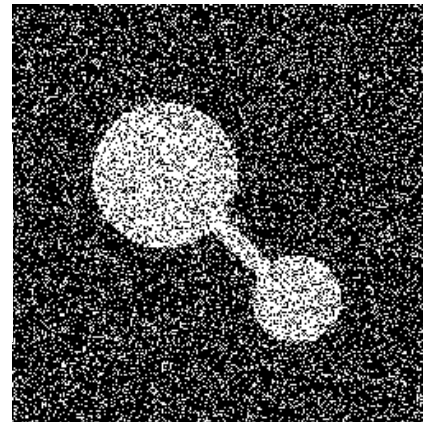
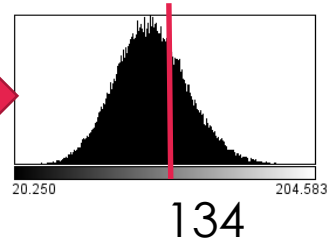
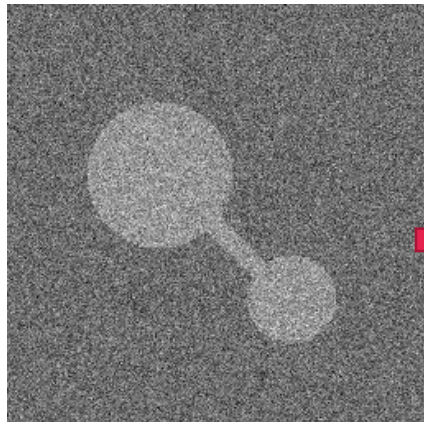
• Ostu Thresholding: 1 et 2 seuils



1 threshold (116)



2 thresholds (104 ; 192)



# REGION GROWING

- Algorithm

- Determine **C neighbors** of region R
  - Add neighbors **similar** to region R
- Iterate while R is growing

*Spatial proximity*

« *Ressemblance* »



## Exemple

Mean of intensities of R

$$\boldsymbol{\mu} = E[\mathbf{I}(R)]$$

Ressemblance ...

$$d = \|\boldsymbol{\mu} - \mathbf{I}(\mathbf{x})\|$$

Critère :  $\mathbf{x}$  is **similar** to R if :  $d < \text{seuil}$

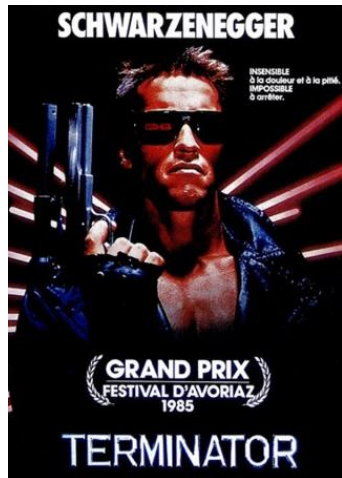
# MACHINE LEARNING

As fast as image introduction



# ARTIFICIAL INTELLIGENCE AND SOCIETY

Until few years ago only science fiction...



...



Today : leaders on personal data use it extensively



... *BUT ALSO*



- Why is it so exciting for value creation?
  - Proposes solutions to (all) problems that have been partially solved or not
  - Obtain better performances than conventional approaches (quality and/or use of resources) → 2012
  - Intrinsically allows to have a usable solution after the 'research' and development phase

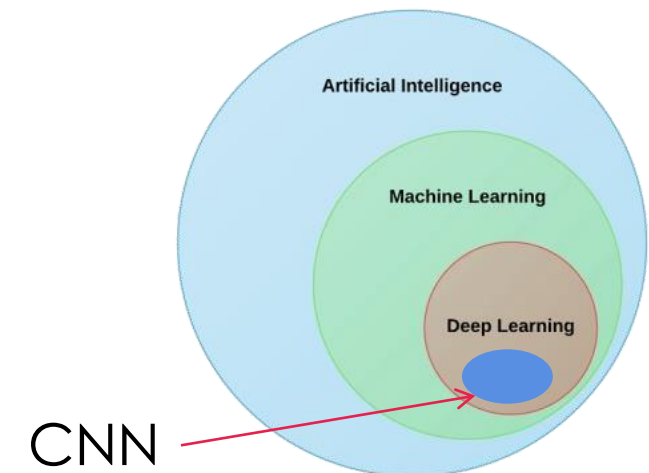




# MACHINE LEARNING APPROACHES

- Supervised approaches (*non ANN*)
  - Bayes, kNN, Support Vector Machine, Random Forest, LDA, Regression Logistique, ...
- Un-supervised approaches (*non ANN*)
  - *K-means, GMM, Hierarchical approach, DBScan, Laio-Rodriguez, ...*

➔ ANN (Artificial Neural Network)  
➔ « DNN » (**Deep Neural Network**)  
**MLP / CNN / RNN / BM**



# NOTATIONS : DATA

- Instance of all features values :  $\mathbf{x}$   $\mathbf{x} \in \mathcal{X}$
- The corresponding predicted label :  $\hat{u}$   $u \in \mathcal{U}$
- Examples  $\mathbf{x}$  :  $\hat{u}$ 
  - Bank statement : solvency
  - Flower picture : identification (rose, dandelion, iris...)
  - Address, surface, year , ... : apartment price

→ a Dataset consists of  $N$  instances  $(\mathbf{x}, u) \rightarrow \mathbf{X}, \mathbf{u}$

$$\mathcal{D} = \{\mathbf{X}, \mathbf{u}\} \quad \mathcal{D} = \{(\mathbf{x}_1, u_1), (\mathbf{x}_2, u_2), \dots, (\mathbf{x}_N, u_N)\}$$

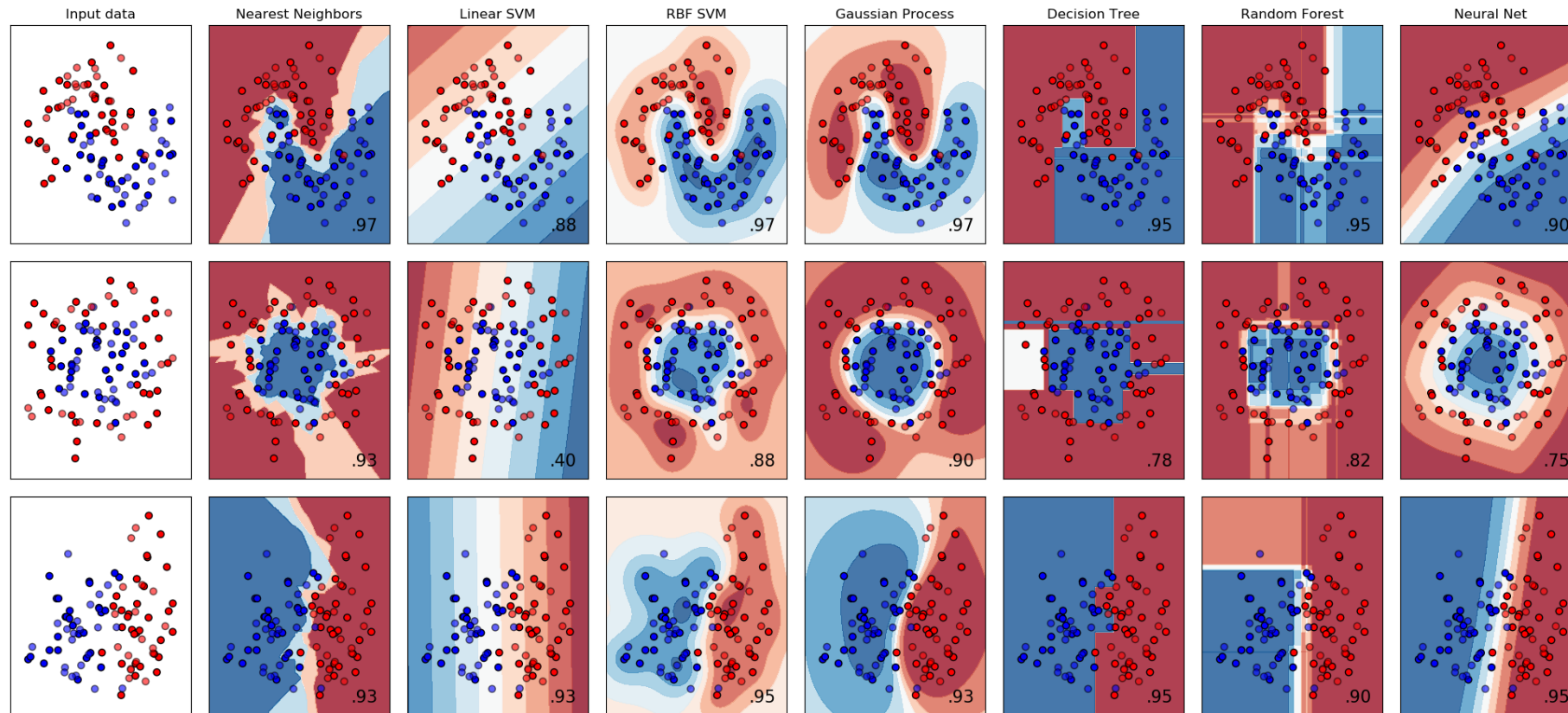
→ **Training set, validation set**

→ **Test set (assessment)**

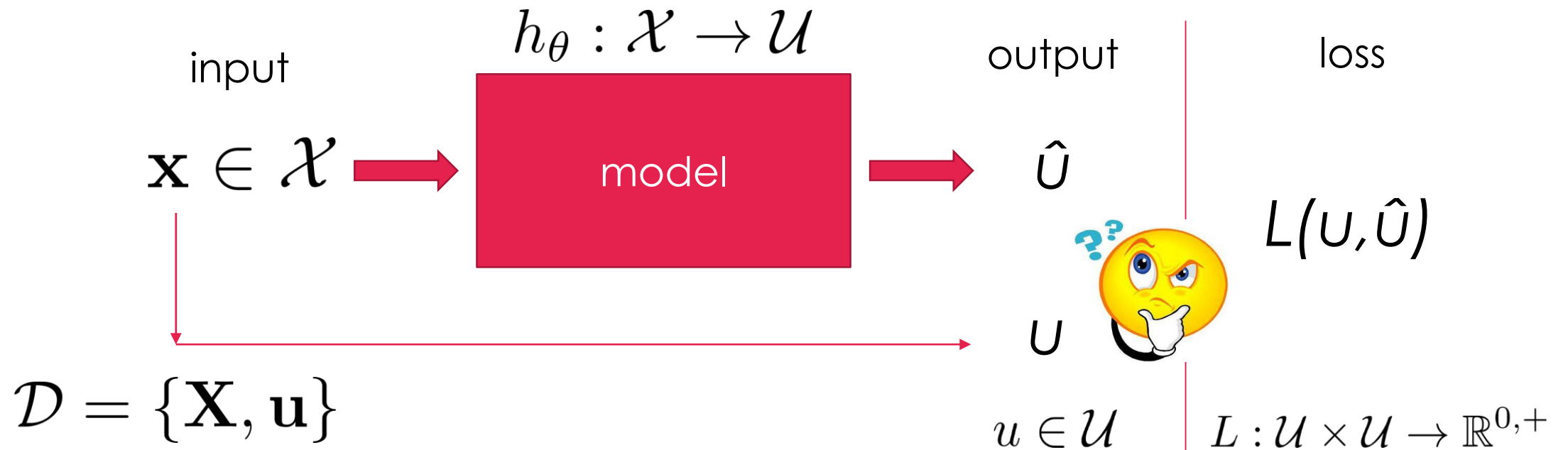
# CLASSIFICATION PROBLEM

- Objective : split the space to identify new incoming instance

scikit-learn.org

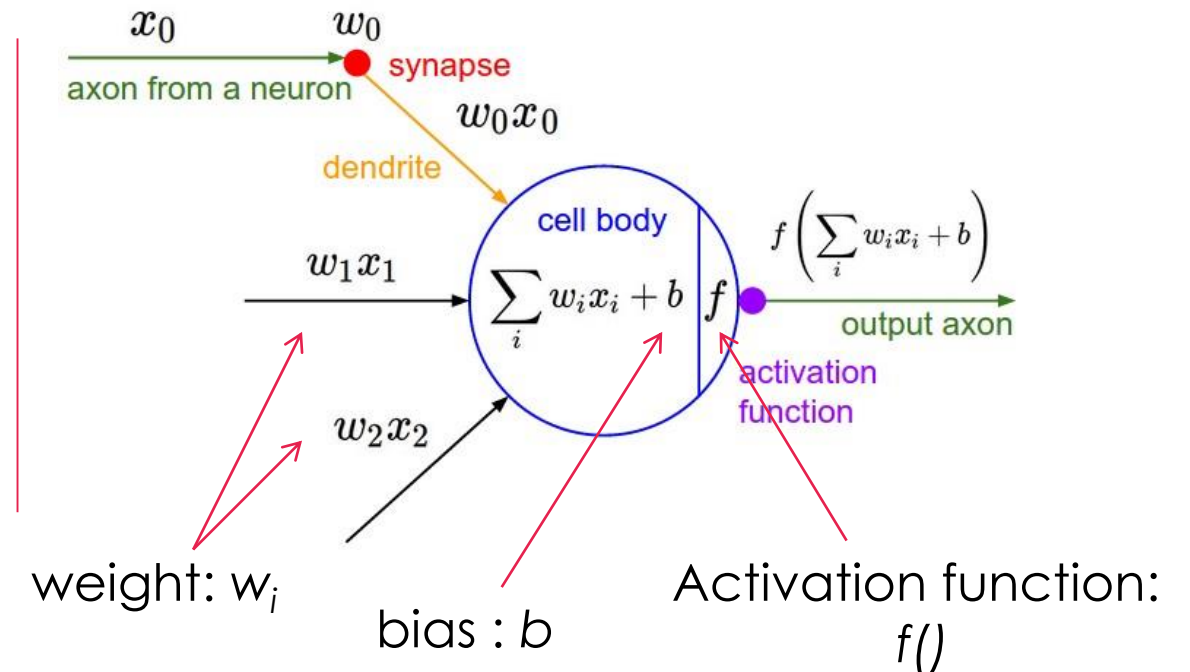
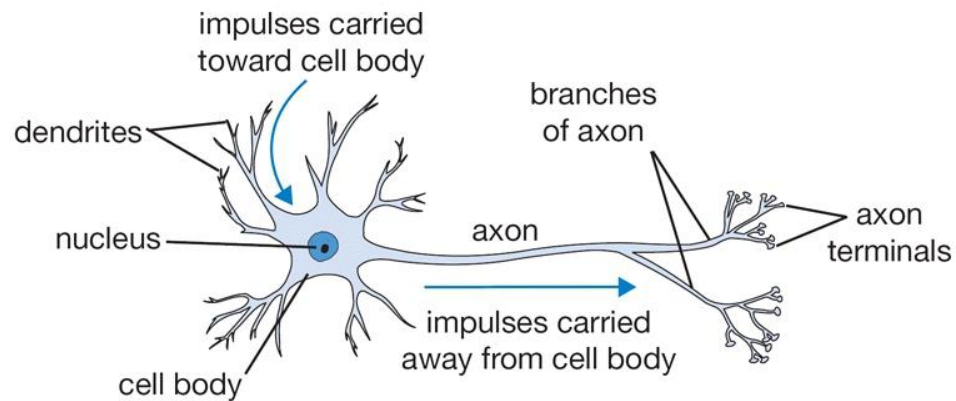


## SCHEMA



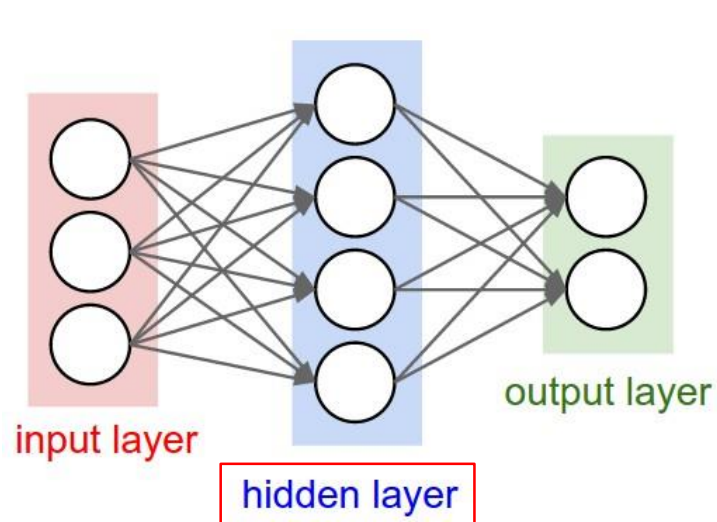
# BRIEFLY, ANN ARE

- Collection of connected **neurons**, structured by layers



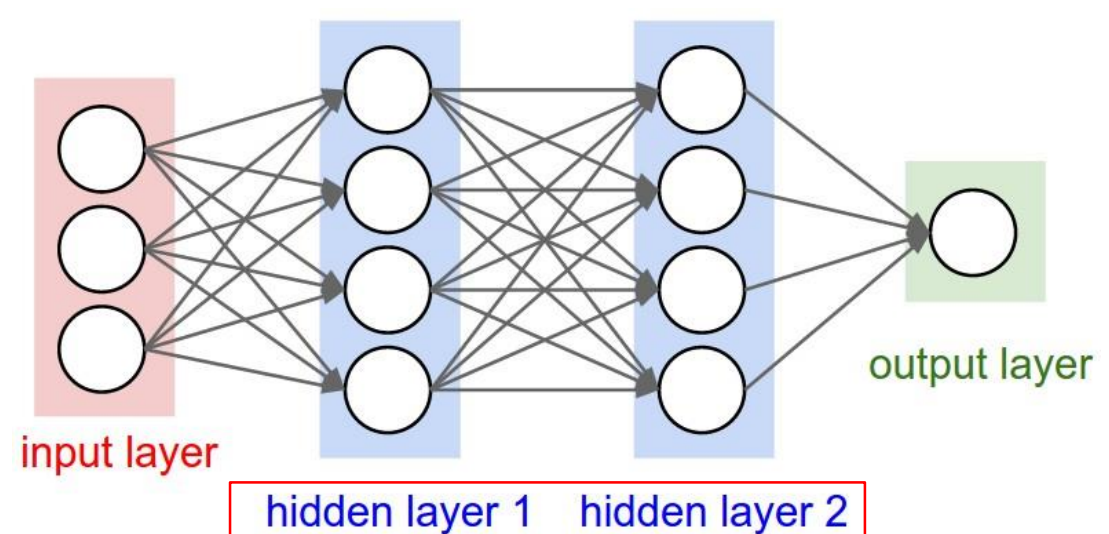
# BRIEFLY, ANN ARE

- Collection of connected **neurons**, structured by **layers**



16 + 10 = 26 variables

$$\hat{u} = h_{\theta}(\mathbf{x})$$

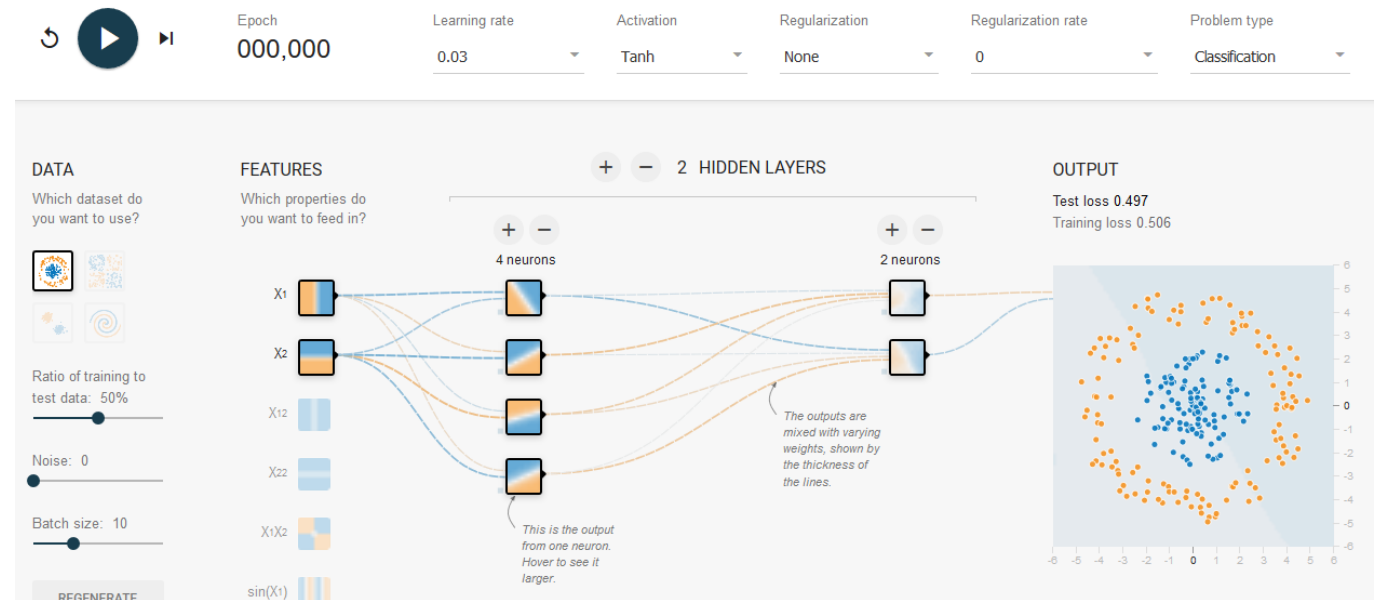


16 + 20 + 4 = 40 variables

# PLAYGROUND TENSORFLOW

<http://playground.tensorflow.org>

- Demo
  - One : simple example
  - Two : XOR case



# MODELING ANN TOWARD MATRIX/VECTOR MANIPULATION

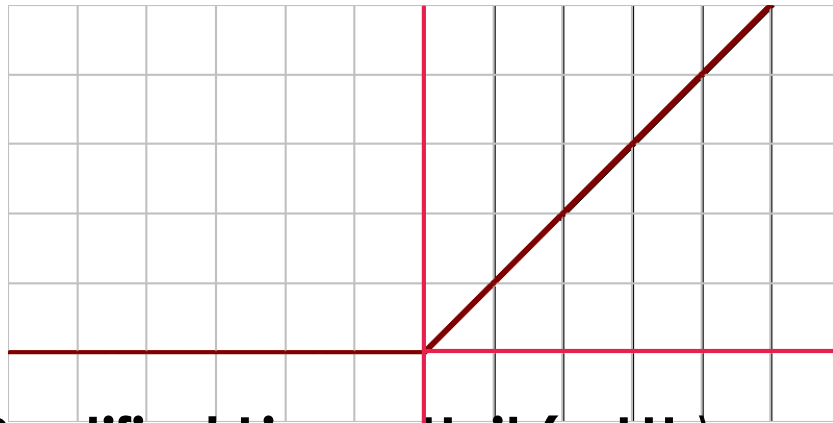
- Neuron with bias :  $\mathbf{w}_k, b_k$  input :  $\mathbf{x}_k$  output :  $z_k$ 
  - $z_k = \mathbf{w}^T \mathbf{x}_k + b$   $\mathbf{x}$  of size  $K_{l-1}$
- A layer  $l$  (without activation) :  $\mathbf{W}_l, \mathbf{b}_l$ 
  - $\mathbf{z}_l = \mathbf{W}_l \mathbf{x} + \mathbf{b}_l$   $\mathbf{W}_l$  of  $K_{l-1}$  row by  $K_l$  lines
  - With  $\mathbf{x}^+$  and  $\mathbf{W}^+$  are  $[1, \mathbf{x}]$  and  $[\mathbf{b}, \mathbf{W}]$
  - $\mathbf{z}_l = \mathbf{W}^+ \mathbf{x}^+$
- Activation layer
  - $\mathbf{y}_l = f_l(\mathbf{z}_l)$

$$z_k^l = \mathbf{w}_k^{lT} \mathbf{x}_k^l$$

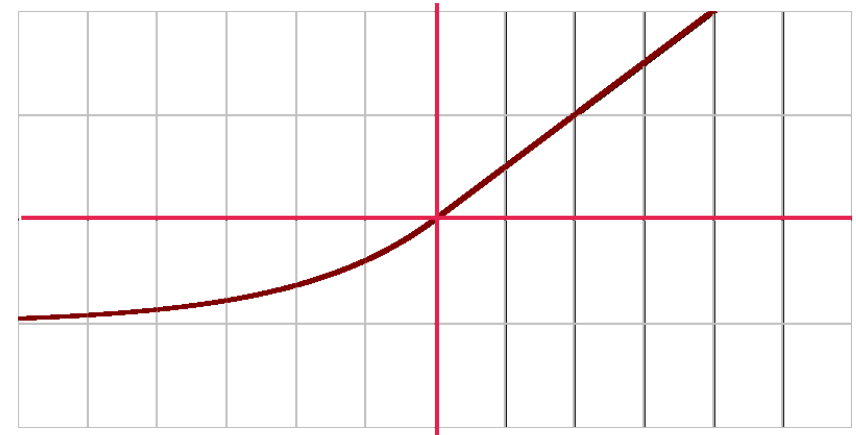
$$y_k^l = f_l(z_k^l)$$



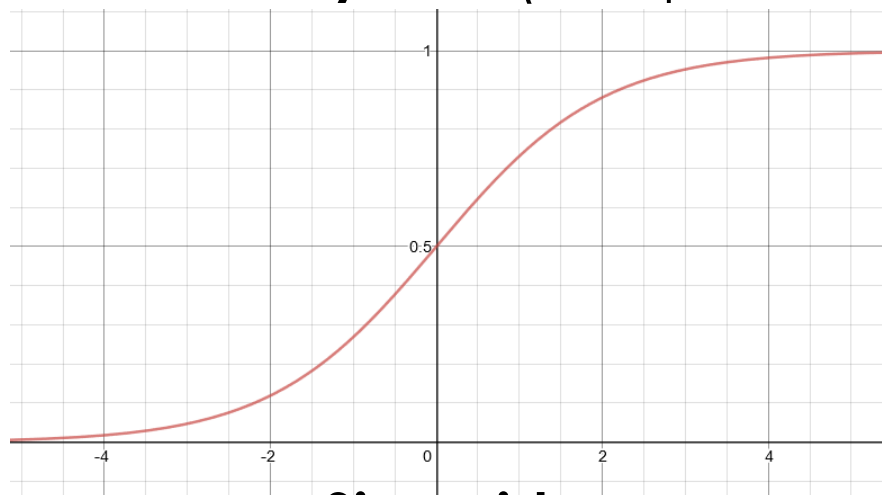
# ACTIVATION FUNCTIONS (LOCAL)



**Rectified Linear Unit (ReLU )**  
and **Leaky reLU** (one parameter)

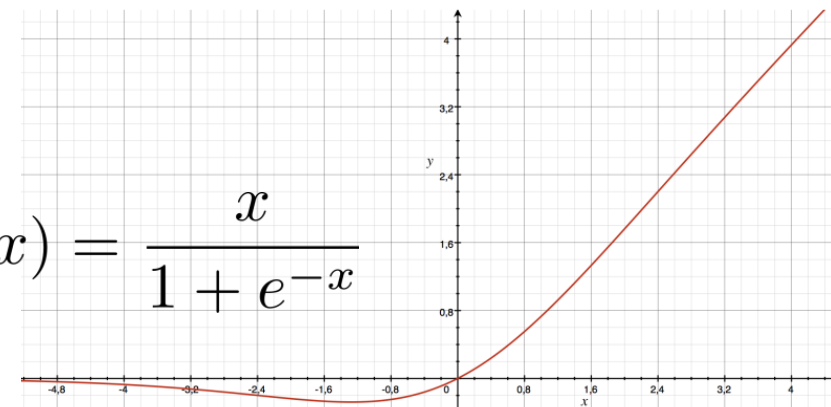


Exponential Linear Unit (eLU )



**Sigmoid**

$$f(x) = \frac{x}{1 + e^{-x}}$$

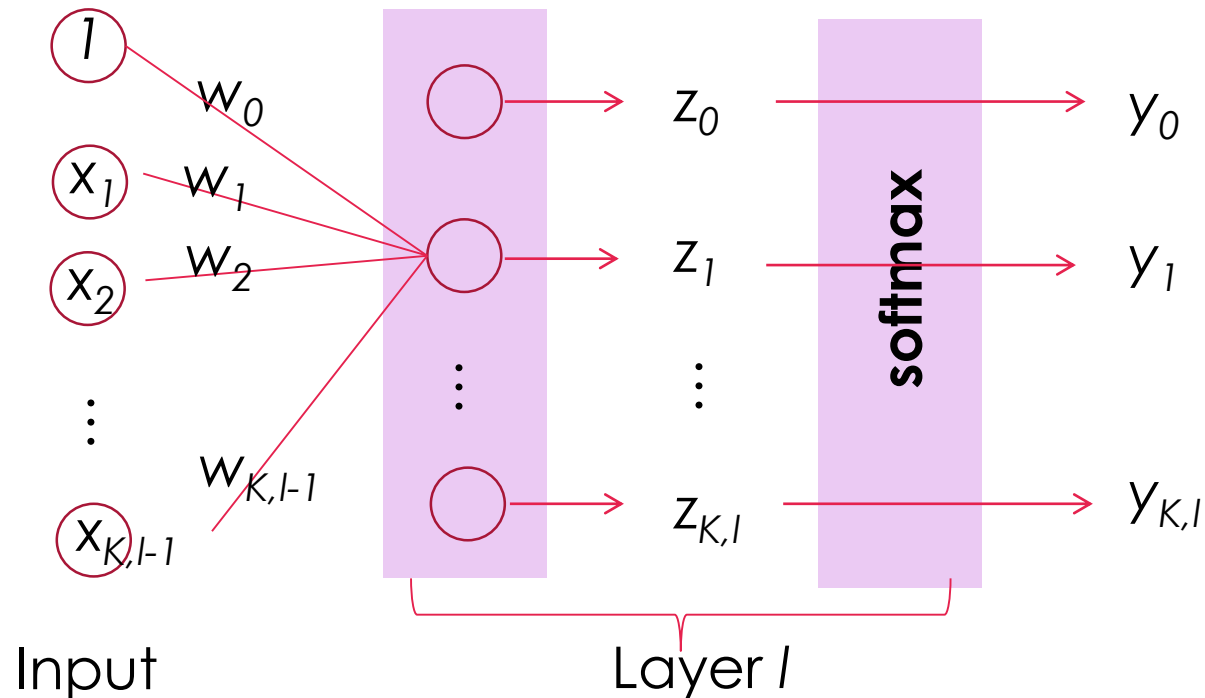


Activation Google (swish)

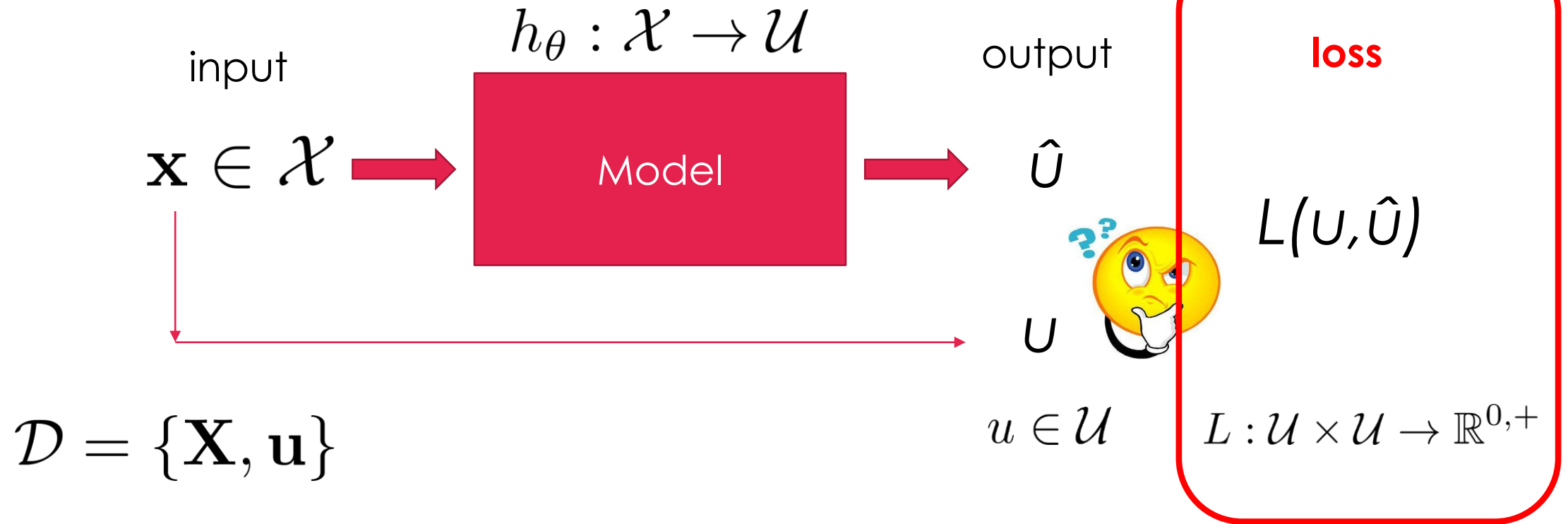
# ACTIVATION FUNCTIONS (OUTPUT LAYER)

**softmax**

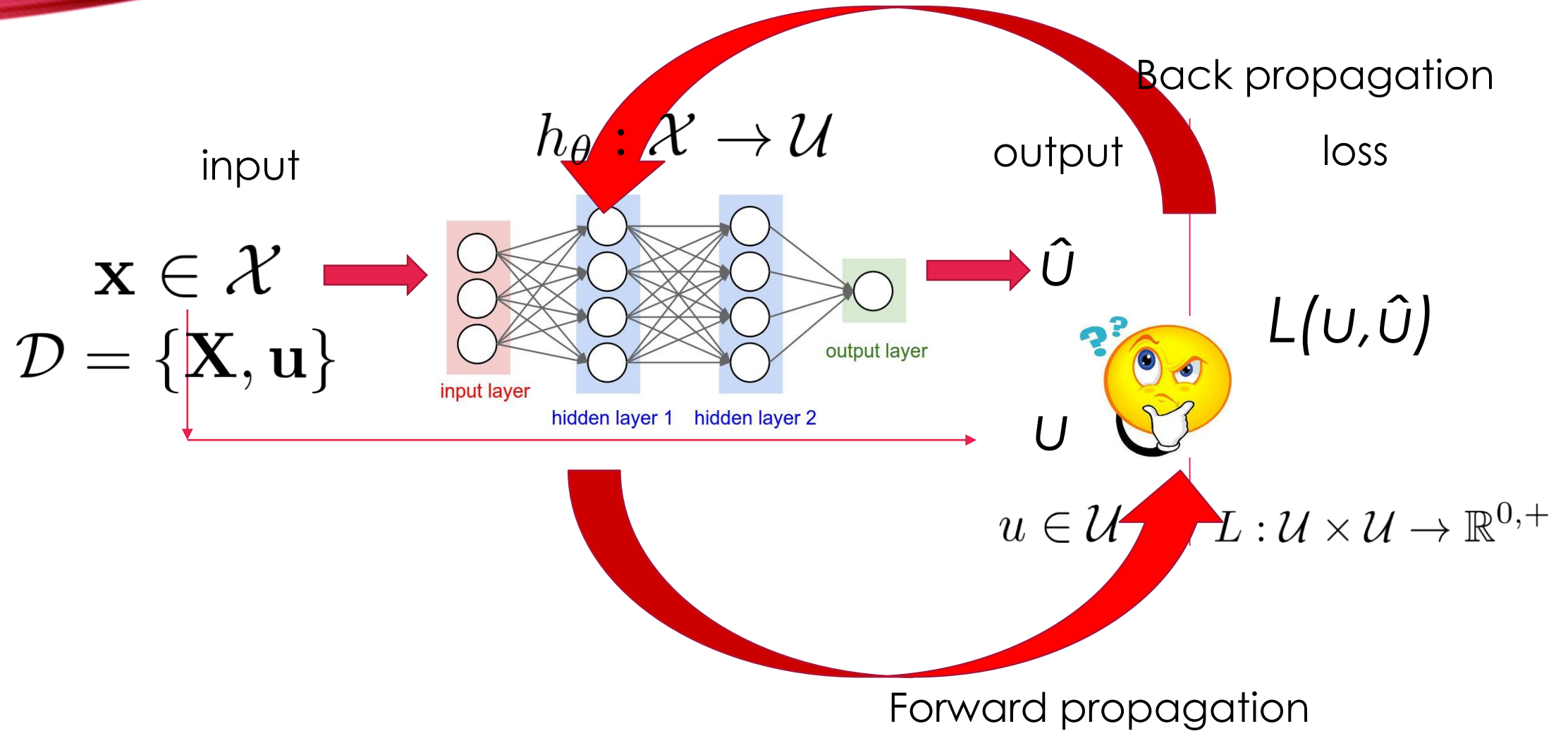
$$f(z_k)|_l = \frac{\exp(z_k)}{\sum_{i=1}^{K_l} \exp(z_i)}$$



## SCHEMA



# SCHEMA, TRAINING



# OPTIMIZATION OF $\mathbf{W}$ : TRAINING

- Loss function must be minimized and will drive the evolution of  $\mathbf{W}$ : gradient descent

$$\mathbf{W}^{[t+1]} = \mathbf{W}^{[t]} - \eta \nabla_{\mathbf{W}} L(\mathbf{W})$$

Learning rate

- In the perceptron case, Novikoff (1962) prove the convergence in  $(D/g)^2$  iterations

- $\|\mathbf{x}^i\|_2 < D$
- $y^i \mathbf{u}^T \mathbf{x}^i \geq g$  with  $\|\mathbf{u}\|_2 = 1$
- $D$  and  $g$  in  $\mathbb{R}^+_{*}$

$$\mathbf{w}_k^{[t+1]} = \mathbf{w}_k^{[t]} - \eta \frac{\partial L(f(\mathbf{w}_k \mathbf{x}^i + b_k), u^i)}{\partial \mathbf{w}_k}$$

# OPTIMIZATION W : DERIVATION?

$$\mathbf{W}^{[t+1]} = \mathbf{W}^{[t]} - \eta \nabla_{\mathbf{W}} L(\mathbf{W})$$

- To update the  $w_{q,k}^l$  the weight of neuron  $q$  of layer  $l-1$  to neuron  $k$  of layer  $l$  we need to compute

$$\frac{\partial L(h_{\mathbf{W}}(\mathbf{x}^i), u^i)}{\partial w_{q,k}^l}$$

$$z_k^l = \mathbf{w}_k^{lT} \mathbf{x}_k^l$$

$$y_k^l = f_l(z_k^l)$$

- Usage of chain rule

$$\begin{aligned} \frac{\partial L(h_{\mathbf{W}}(\mathbf{x}^i), u^i)}{\partial w_{q,k}^l} &= \frac{\partial L(h_{\mathbf{W}}(\mathbf{x}^i), u^i)}{\partial z_k^l} \frac{\partial z_k^l}{\partial w_{q,k}^l} \\ &= \frac{\partial L(h_{\mathbf{W}}(\mathbf{x}^i), u^i)}{\partial y_k^l} \frac{\partial y_k^l}{\partial z_k^l} \frac{\partial z_k^l}{\partial w_{q,k}^l} = \frac{\partial L(h_{\mathbf{W}}(\mathbf{x}^i), u^i)}{\partial y_k^l} \frac{\partial y_k^l}{\partial z_k^l} \frac{\partial z_k^l}{\partial w_{q,k}^l} \end{aligned}$$

# OPTIMIZATION W : DERIVATION?

- Matrice view chain rule

$$\mathbf{z}^l = \mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l$$

$$\mathbf{y}^l = f_l(\mathbf{z}^l)$$

$$\mathbf{x}^{l+1} = \mathbf{y}^l$$

$$\frac{\partial L}{\partial \mathbf{W}^l} = \frac{\partial L}{\partial \mathbf{y}^l} \frac{\partial \mathbf{y}^l}{\partial \mathbf{W}^l} = \frac{\partial L}{\partial \mathbf{y}^l} \frac{\partial \mathbf{y}^l}{\partial \mathbf{z}^l} \frac{\partial \mathbf{z}^l}{\partial \mathbf{W}^l}$$

$$= \frac{\partial L}{\partial \mathbf{x}^{l+1}} \frac{\partial \mathbf{y}^l}{\partial \mathbf{z}^l} \frac{\partial \mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l}{\partial \mathbf{W}^l}$$

$$= d'_{l+1} \odot f'_l \frac{\partial \mathbf{W}^l \mathbf{x}^l + \mathbf{b}^l}{\partial \mathbf{W}^l} = (\mathbf{x}^l)^T (d'_{l+1} \odot f'_l)$$

Matrix transpose

Derivative computed for layer l+1  
(wrt to its input  $\mathbf{x}^{l+1}$ )

$$\frac{\partial L}{\partial \mathbf{x}^l} = (\mathbf{W}^l)^T (d'_{l+1} \odot f'_l)$$

Hadamard product  
(element wise multiplication)

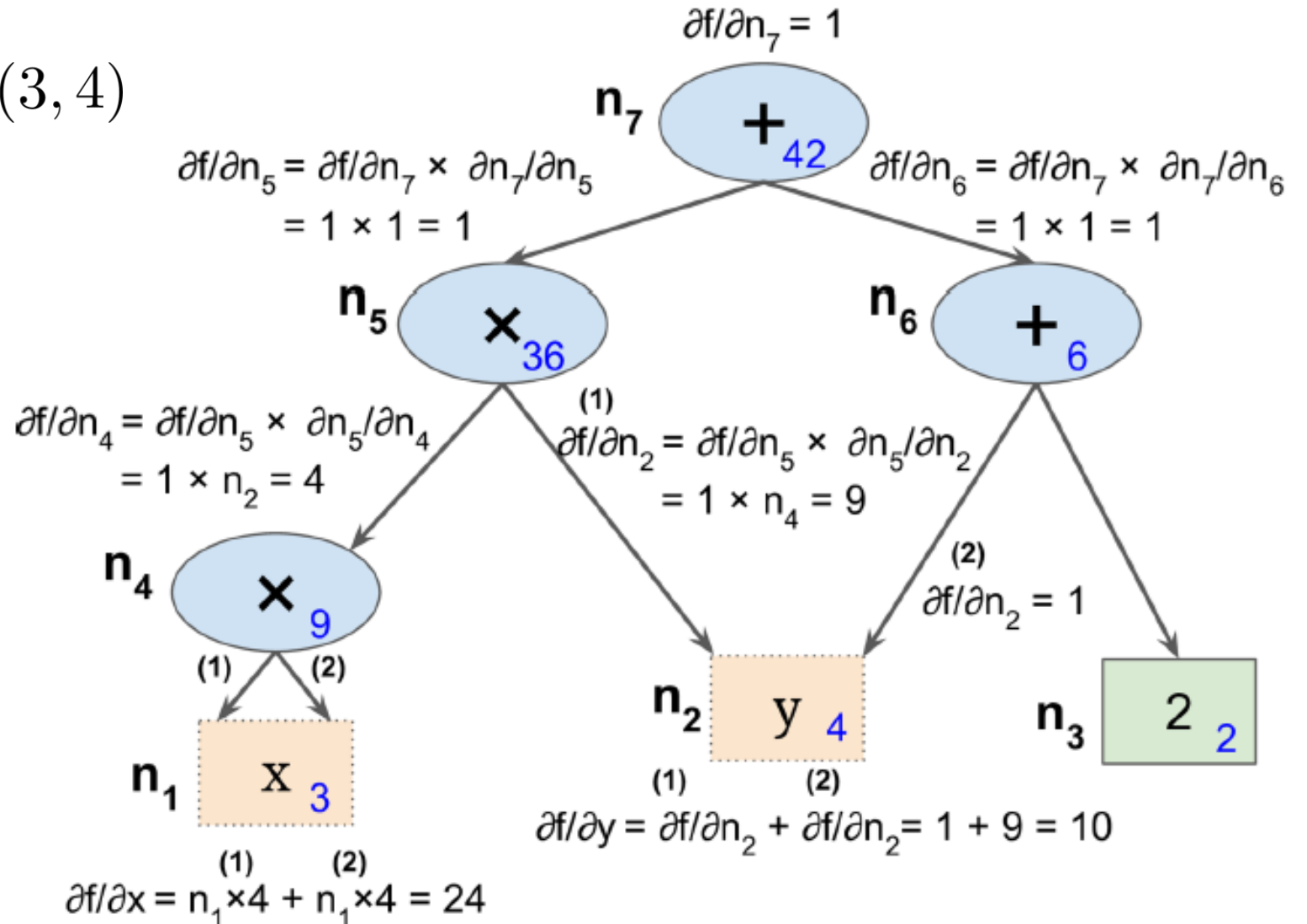
# BACKPROPAGATION EXAMPLE

- $f(3,4)$  and

$$f(x, y) = x^2y + y + 2$$

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial n_i} \frac{\partial n_i}{\partial x}$$

$$\nabla f(3, 4)$$



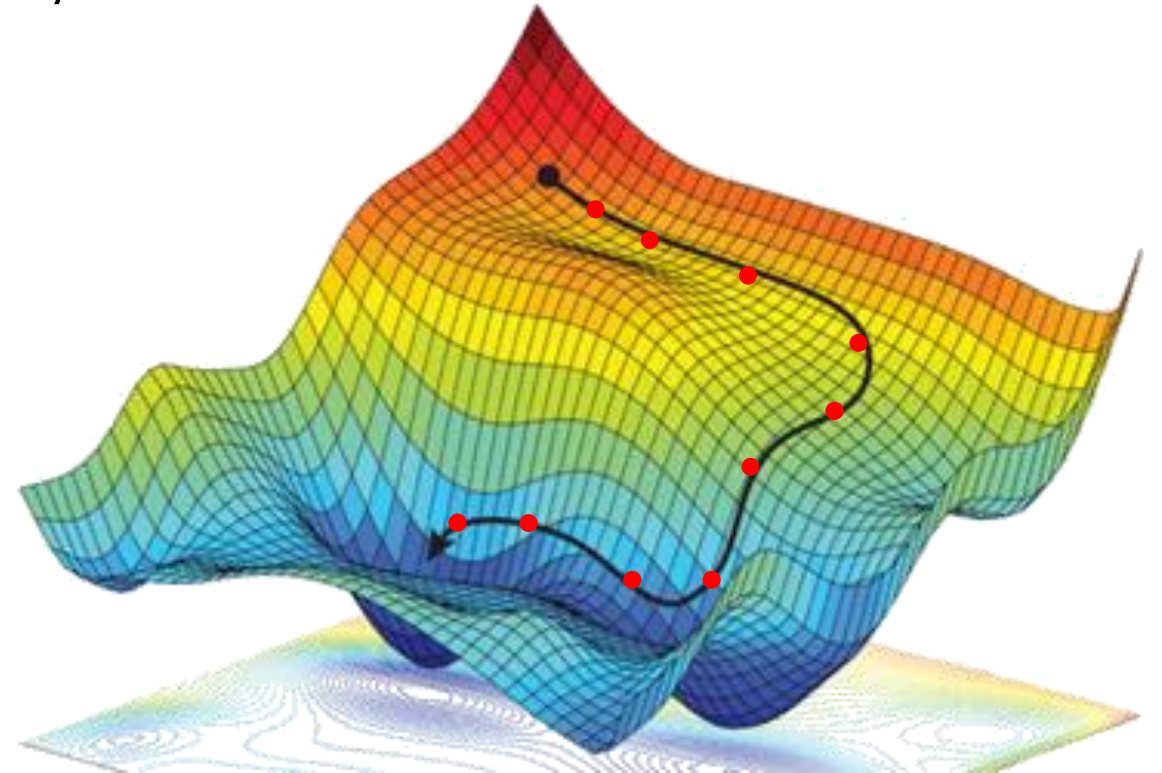


# GRADIENT DESCENT

$$\mathbf{W}^{[t+1]} = \mathbf{W}^{[t]} - \eta \nabla_{\mathbf{W}} L(\mathbf{W})$$

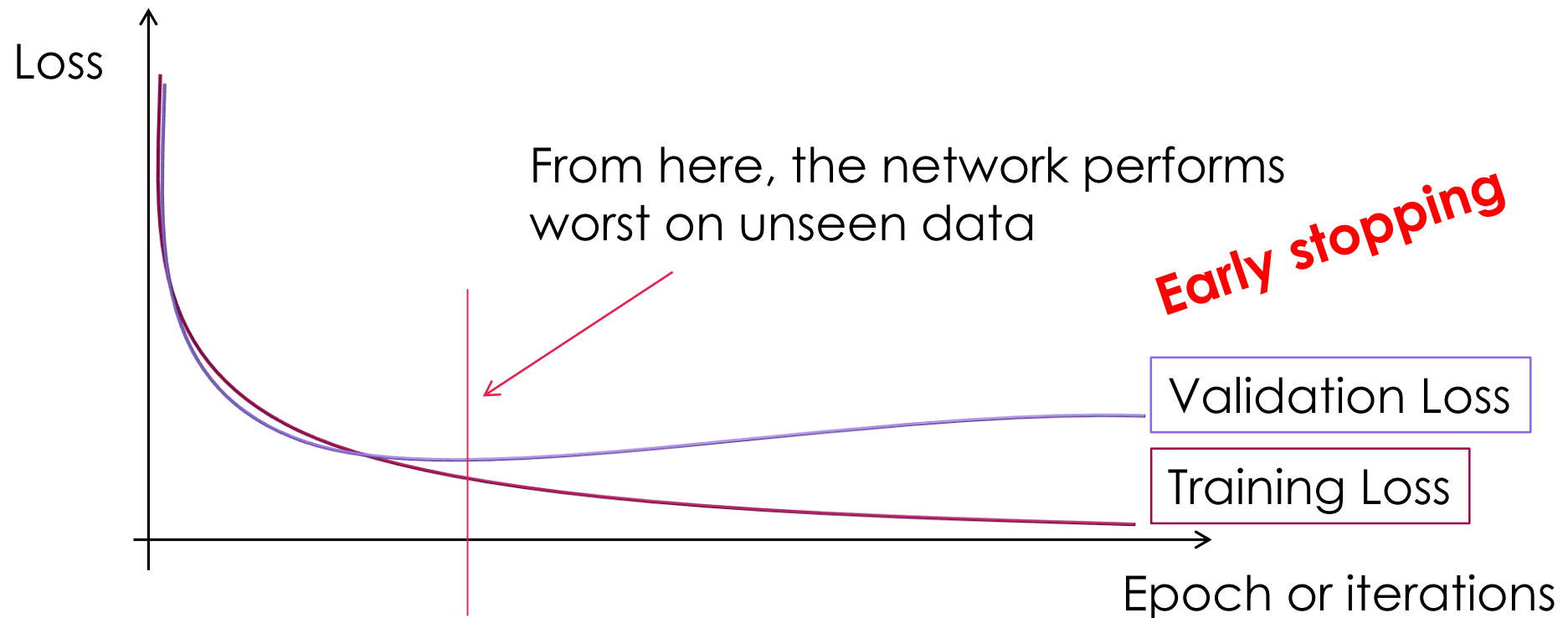
$= 0 \ ? \ \rightarrow \text{minimum}$

- Iterative process
- Can be regularized (**momentum**, **Adam**)
- 3 ways of using data
  - Use the whole dataset
    - Long, small update
  - Use one instance (Stochastic)
    - Fast but noisy
  - Use a couple of instances
    - $\rightarrow$  mini batch
    - Quite fast, less noisy



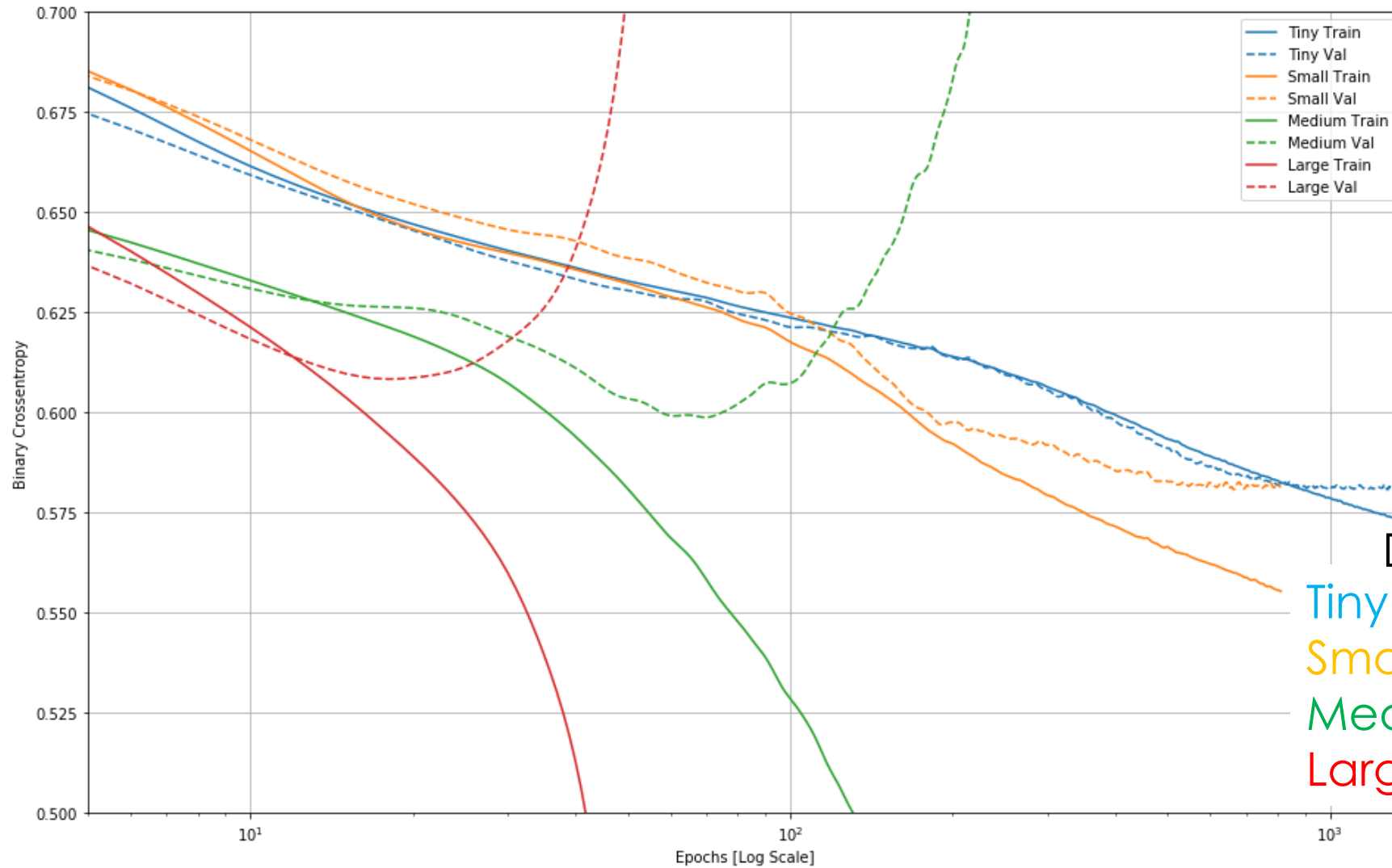
# LOSSES EVOLUTION DURING TRAINING

- Loss on training set → learning the problem
- Loss on validation set → generalization ability



# REAL EVOLUTION CURVES...

- [https://www.tensorflow.org/tutorials/keras/overfit\\_and\\_underfit](https://www.tensorflow.org/tutorials/keras/overfit_and_underfit)



Dense networks

Tiny : 16

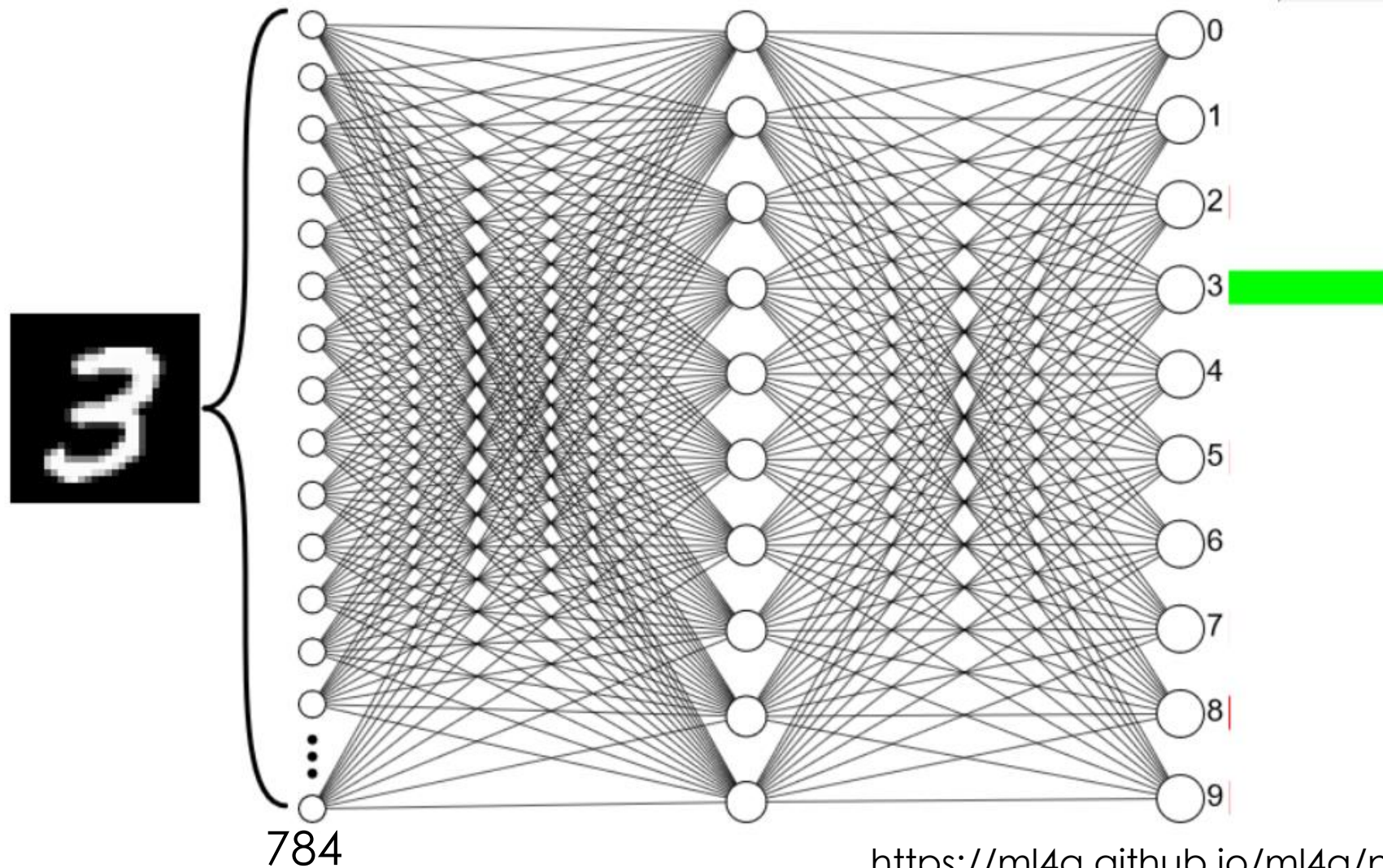
Small : 16, 16

Medium: 64, 64, 64

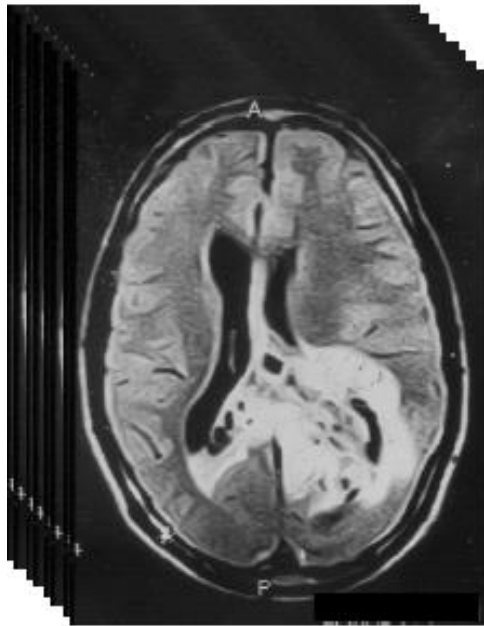
Large: 512, 512, 512, 512



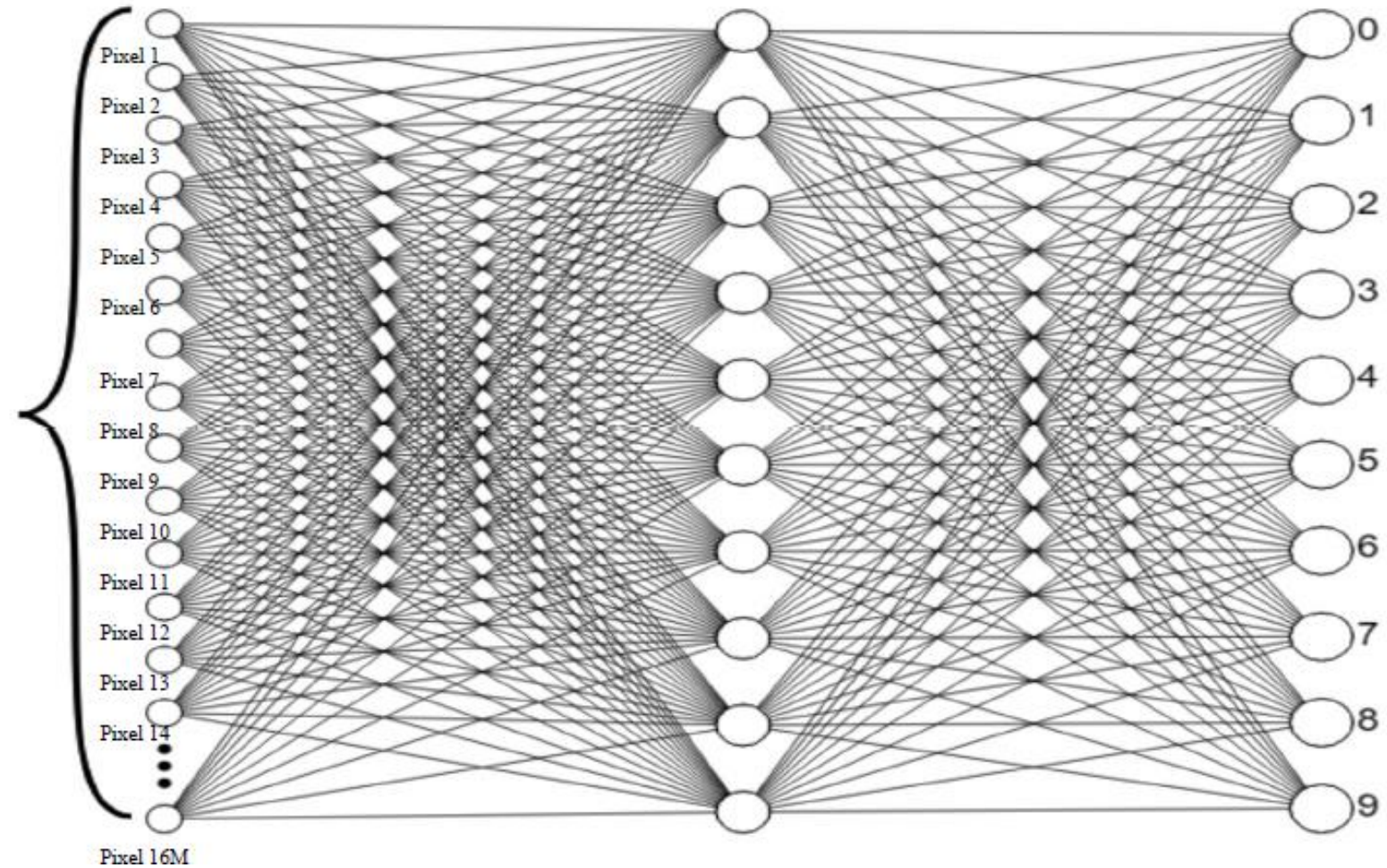
# MLP : CLASSIFICATION EXAMPLE



# MLP: LIMITATIONS



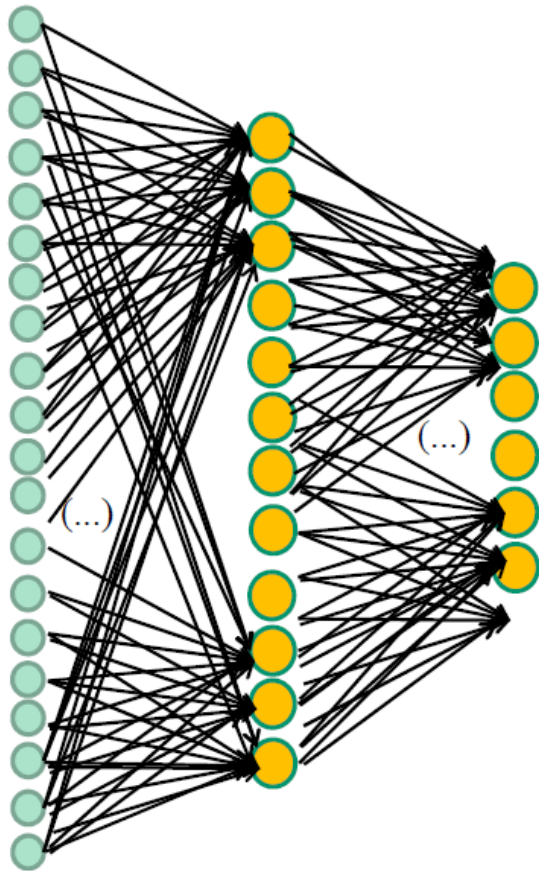
256x256x256



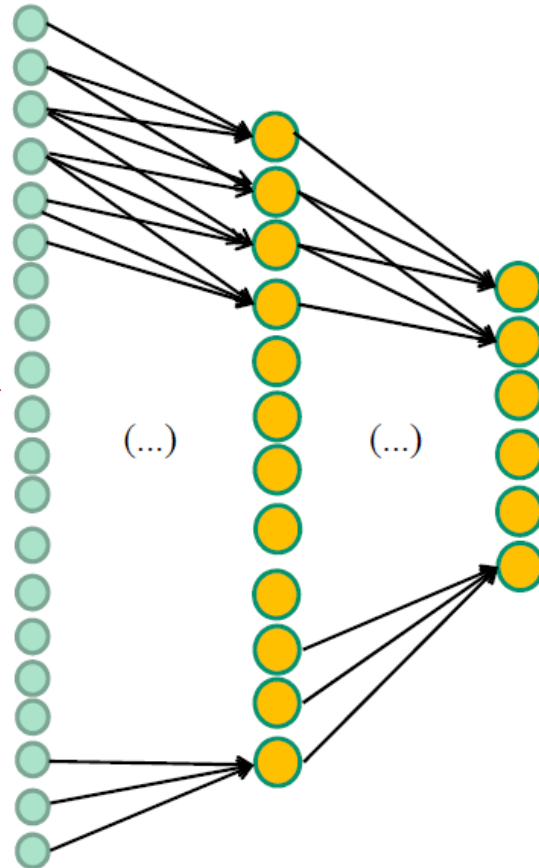
**160M** parameters for the first layer!

# MLP LIMITATIONS AND SOLUTION!

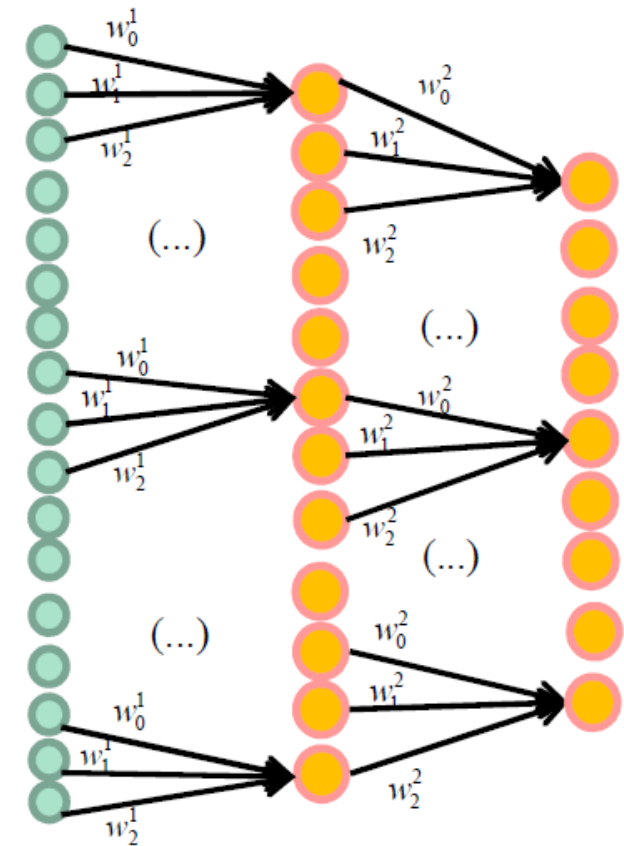
Fully connected



Not fully connected



**Sharing weights**



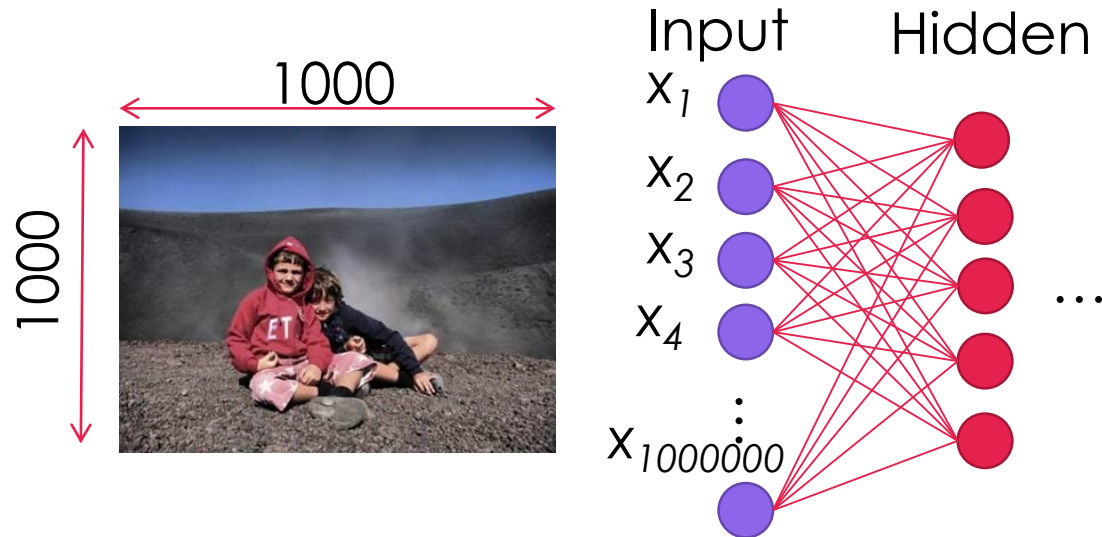
# II – CONVOLUTIONAL NEURAL NETWORK : ELEMENTS

Convolution and related dedicated layers





# CNN, YANN LECUN

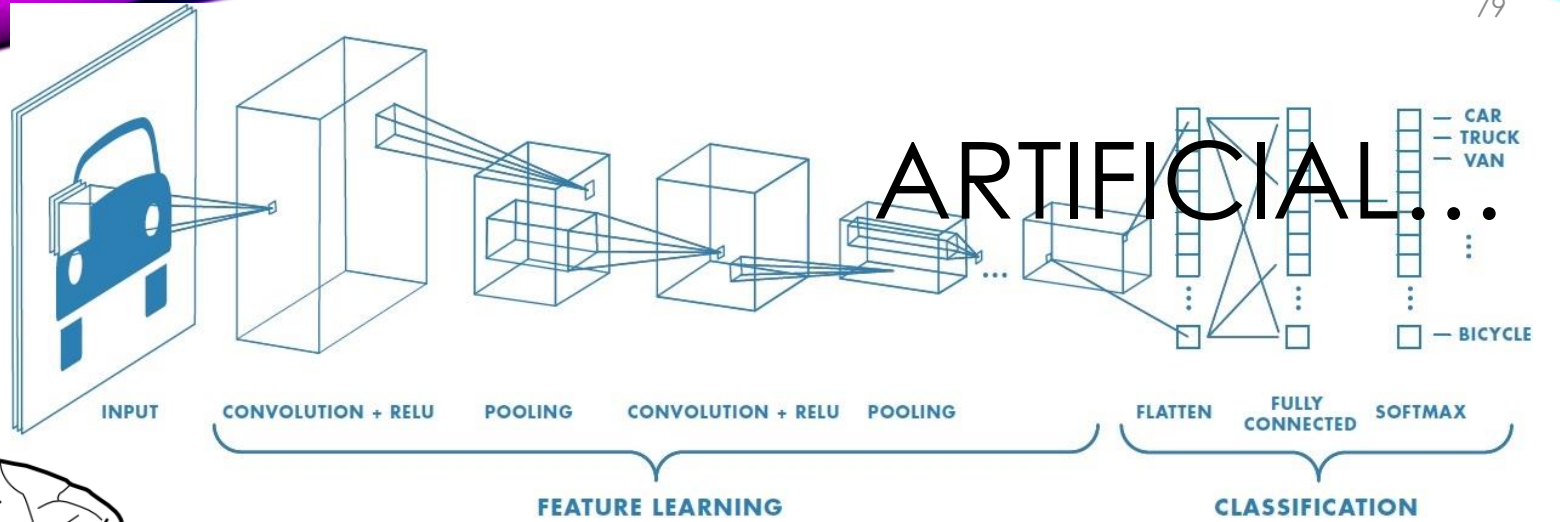


- Too much dimensions!
- Use carefully the *fully connected*

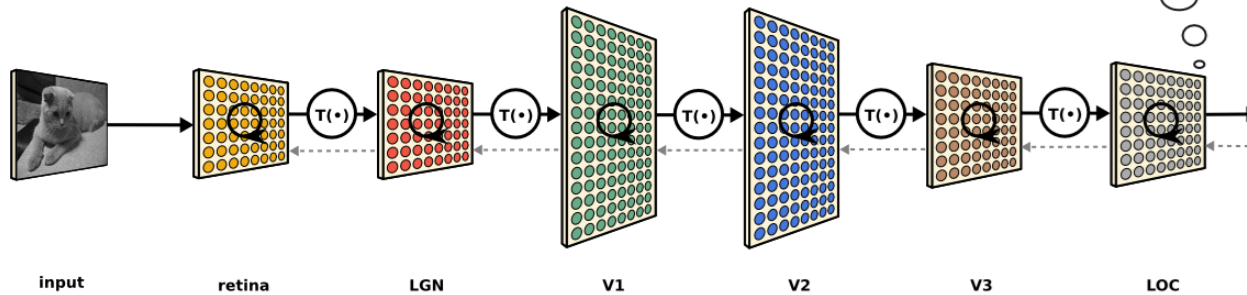
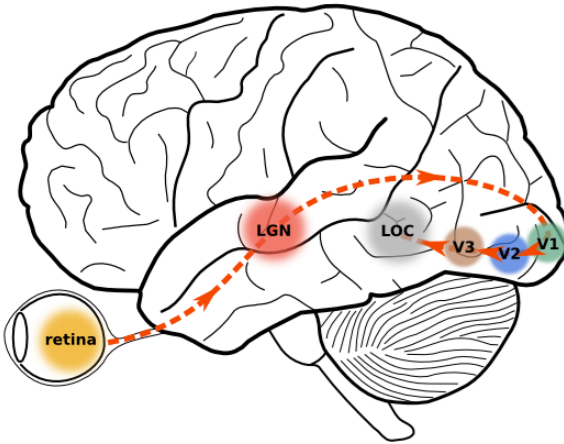
« malédiction de la grande dimension »

Le nombre d'exemples d'apprentissage doit augmenter exponentiellement avec le nombre de dimensions

- We have to reduce the number of dimensions if we don't want some influence on result
- 2 properties of natural and numeric data (link to human)
  - Multiscale hierarchic organization
  - **Symetries**



ARTIFICIAL...

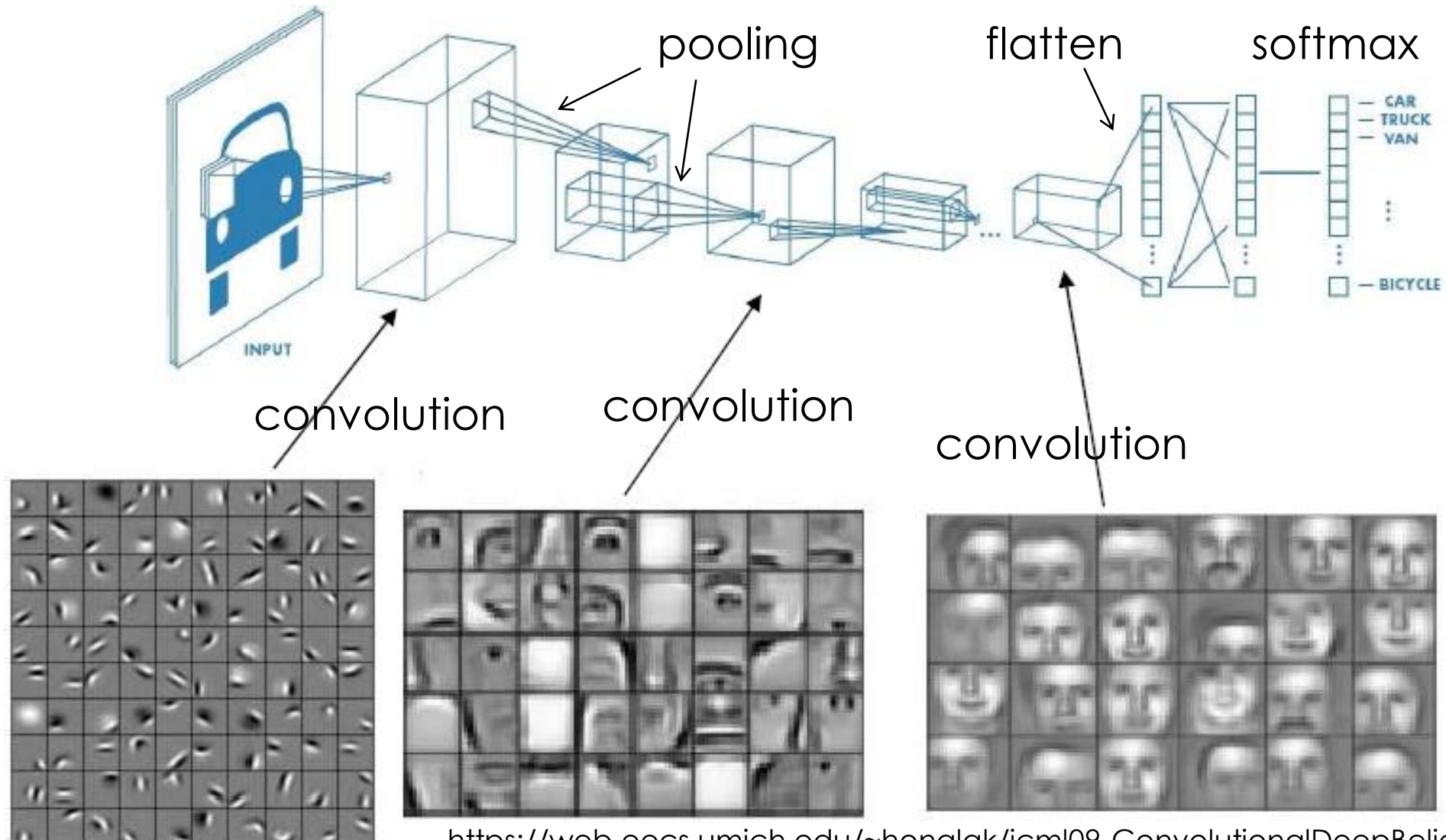


... vs natural

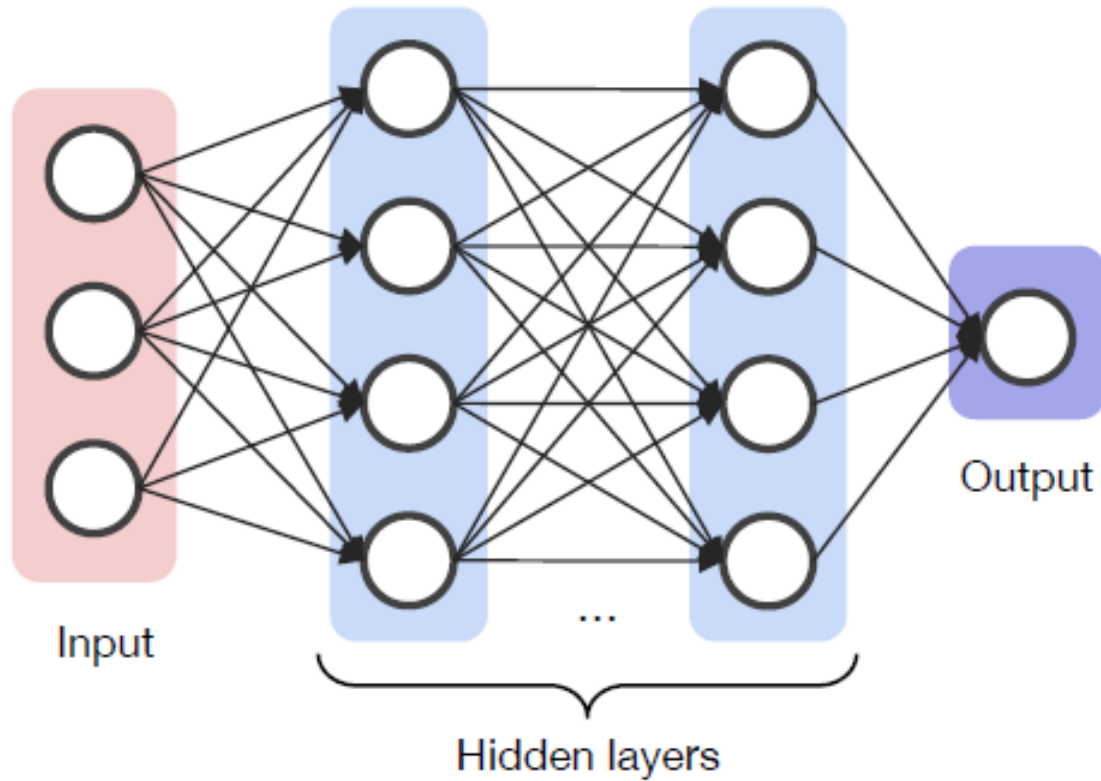
[https://neuwritesd.files.wordpress.com/2015/10/visual\\_stream\\_small.png](https://neuwritesd.files.wordpress.com/2015/10/visual_stream_small.png)

<https://fr.mathworks.com/solutions/deep-learning/convolutional-neural-network.htm>

# TYPICAL CNN (CLASSIFICATION)

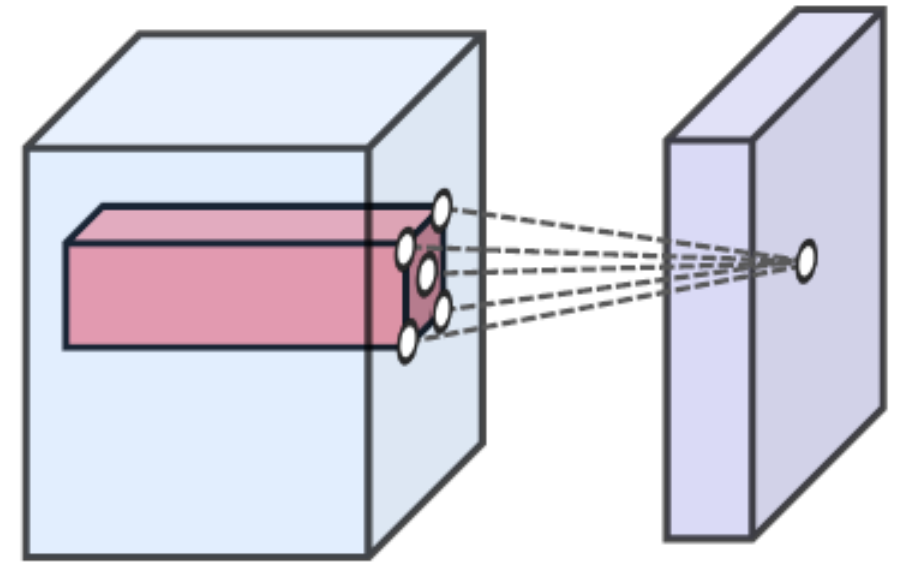


## Multi-layer perception (MLP)



Many parameters  
Limited depth

## CNN

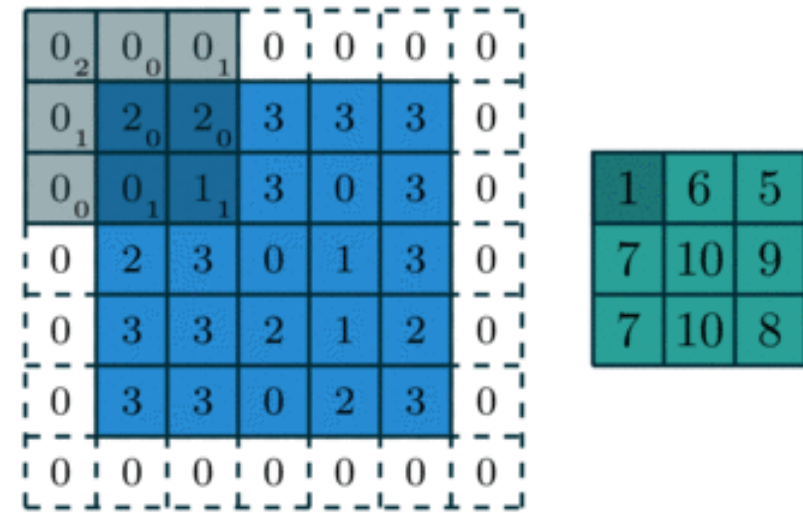


- Shared parameters
- Local connectivity

Less parameters  
Deep architectures possible

# CONVOLUTION FOR NEURAL NETWORK

- **Convolution**
- Pad
  - border condition
- Stride
  - jump
- Atrous
  - empty pixels in mask



[http://deeplearning.net/software/theano/tutorial/conv\\_arithmetic.html#zero-padding-unit-strides](http://deeplearning.net/software/theano/tutorial/conv_arithmetic.html#zero-padding-unit-strides)

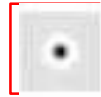
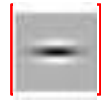
# CONVOLUTION, FILTERS AND CHANNELS



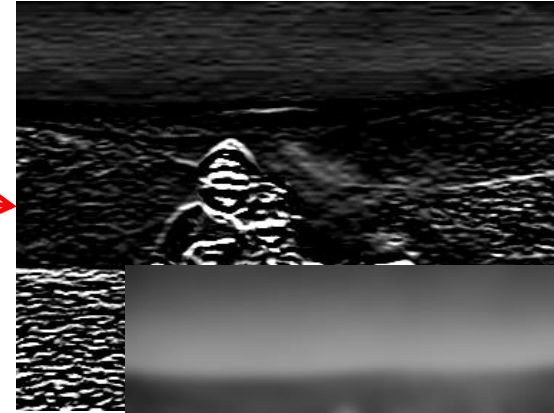
One input channel



convolution



Filters have the same number of channels than the inputs



Three filters → Three output channels

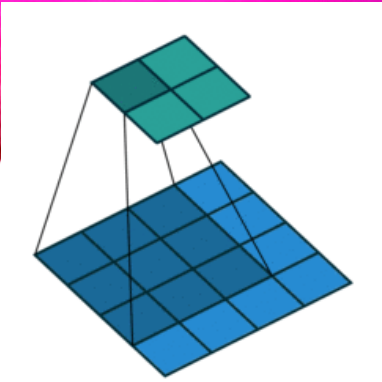
# PRE-PROCESSING 'LAYERS'

Intensity values of input data are critical: how to warrant that the dynamic and statistic are similar on all inputs ? → **pre-process all data!**

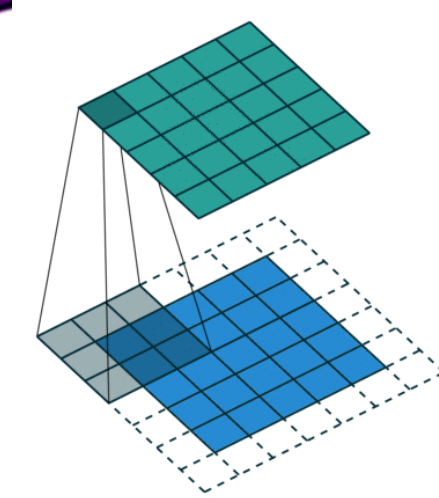
- Mean-substraction
- Normalization
- Local contrast normalisation
- (size adaptation and interpolation)

$$\mathbf{x}_0 = \mathbf{x} - \hat{\mathbf{x}}$$

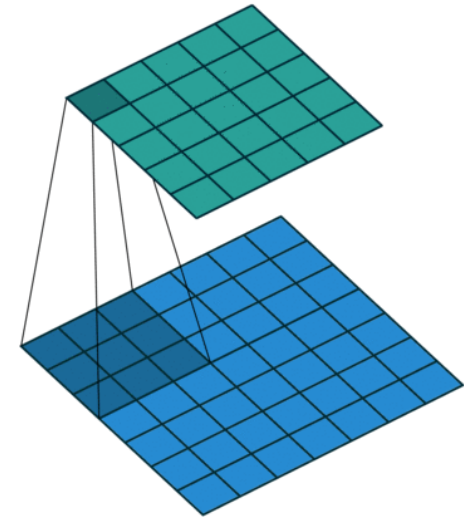
$$\mathbf{x}_1 = \frac{\mathbf{x}_0}{\sigma}$$



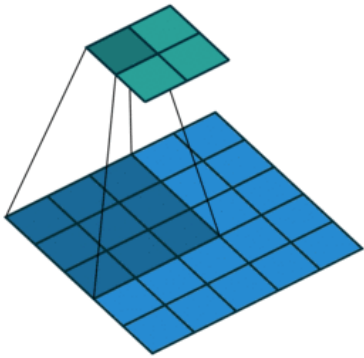
No zero padding,  
unit stride



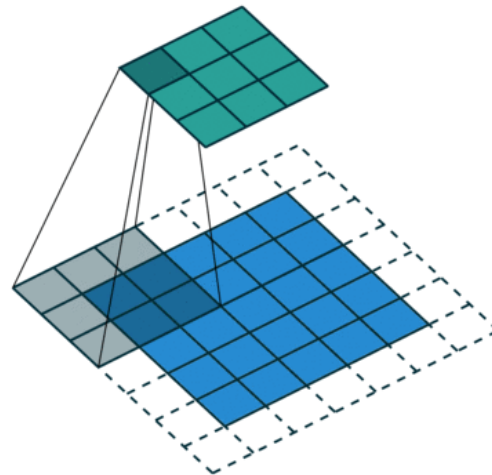
Half size padding ('same' size)  
Unit stride



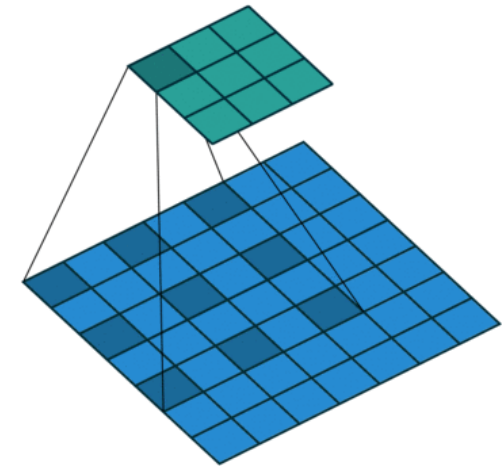
Full padding, Unit stride



No zero padding,  
Non-unit stride



Half size padding ('same' size)  
Non-unit stride



Atrous or dilated filter

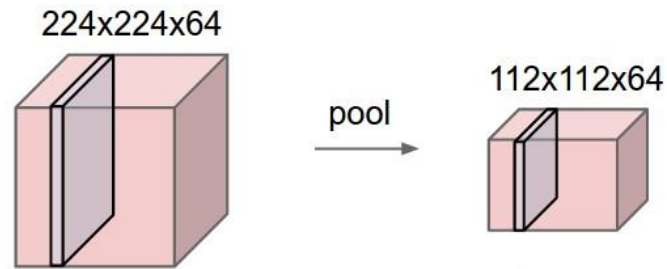


# POOLING

<https://cs231n.github.io/convolutional-networks/>

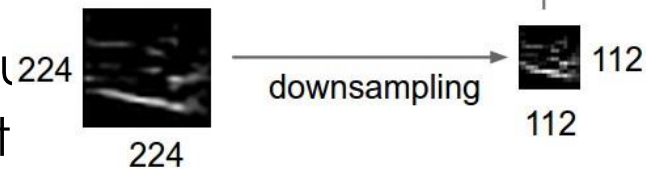
- Down sampling

- **max**,
- mean
- ...



- Up sampling

- Nearest Neighbor
- deconvolution – †



## Max pooling

Single depth slice

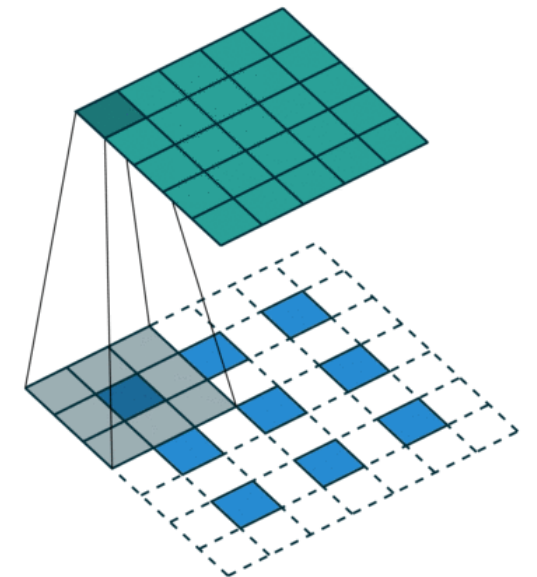
1	1	2	4
5	6	7	8
3	2	1	0
1	2	3	4

x

y

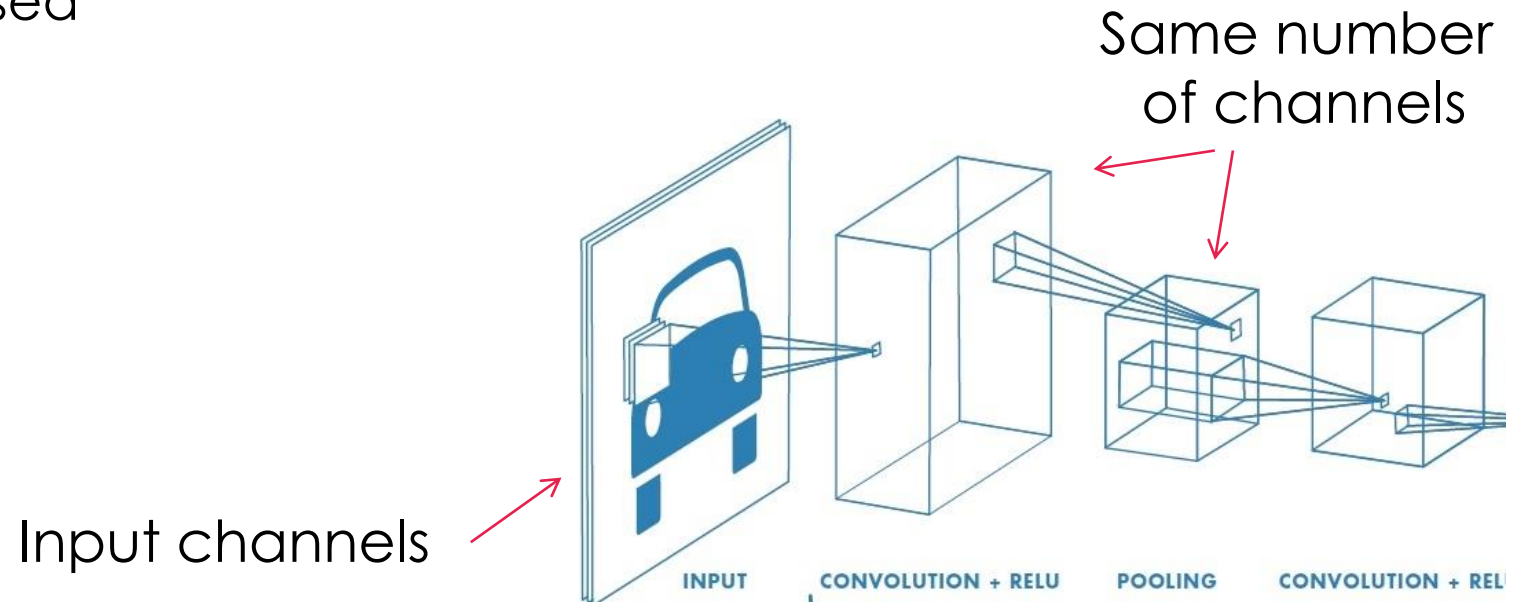
max pool with 2x2 filters  
and stride 2

6	8
3	4



# DON'T MISS IT

- Convolution works on all channels
  - Generally 3 by 3, 5 by 5, 7 by 7 or 1 by 1
  - At a given layer, many channels are obtained, one per convolution
- Pooling works on channel per channel
  - Generally size is decreased by 4 (2 by 2)
  - (it can be increased)
- Receptive field

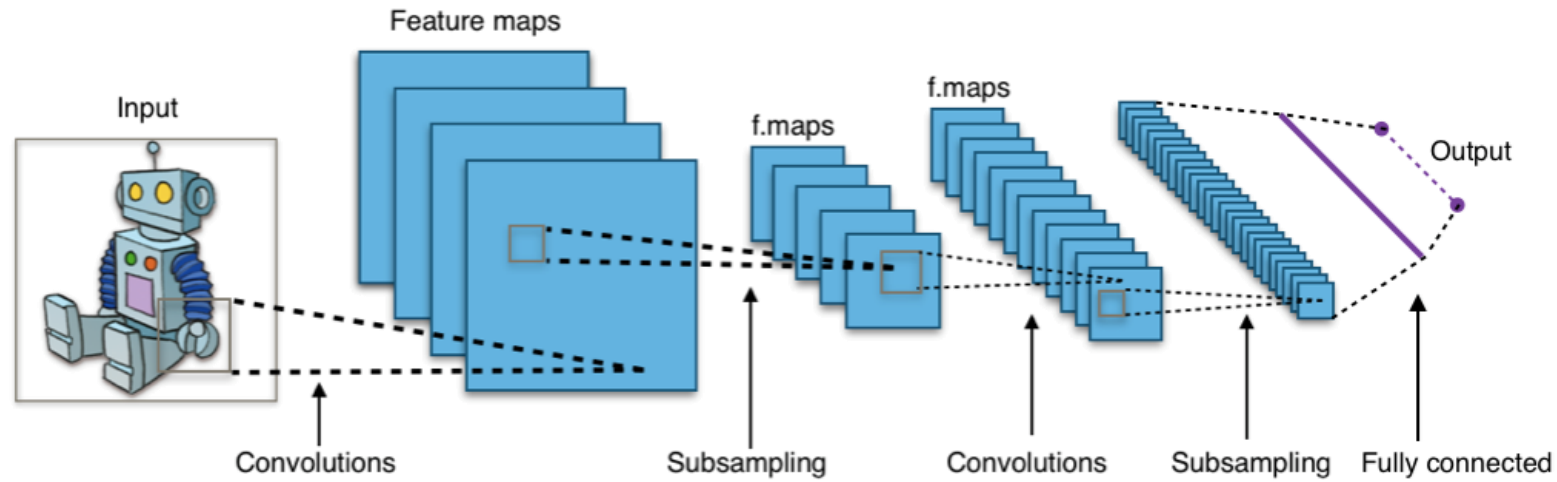


# APPLICATIONS: CLASSIFICATION AND SEGMENTATION



# 1. CNN

- Classification, détection, régression

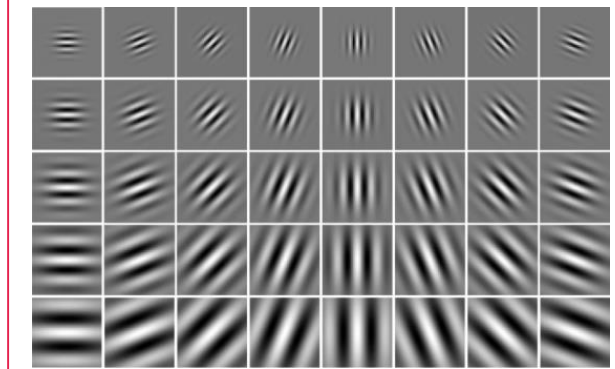


96 filtres  
(ImageNET)



[Krizhevsky 12]

*Filtres Gabor*



## 2. U-NET (2015)

- Segmentation
- Filtering

➔ Output

Same support as input

