



UNIVERSITÀ  
DEGLI STUDI  
DI PADOVA

III  
DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE

**Static code  
893289**



Simone Milani  
Room 216  
DEI A

Phone: 049 827 7641  
E-mail:  
[simone.milani@dei.unipd.it](mailto:simone.milani@dei.unipd.it)

**3D AUGMENTED REALITY  
A.A. 2020/2021**

**LECTURE 17 – HAPTIC/DL**

# Haptic interfaces



So far we have seen,

- Visors (standards screens, stereo VR, stereo see-through, projectors, ...)
- 3D scanning camera (to model the space, to acquire movements, ...)
- Audio (voice-based) interfaces
- controllers, joysticks (used by hands)

Other interfaces can be possible ...



Haptic technology, also known as kinaesthetic communication or 3D touch, refers to any technology that can create an experience of touch by applying forces, vibrations, or motions to the user

**Actuators:** create the sense of touch

**Sensors:** measure motion and poses with inverse process



Usually they combine force, vibration, motions, thermal stimuli: from voltage to stimulus

Manipulate the movements of the user (go beyond a mere alert)

Not all the sensors require touch.

- Tactile
  - Vibrotactile
  - Thermal
  
  - Eccentric Rotating Mass Vibration (ERMV)  
*Offcentered motion of a rotating element*
  - Linear Resonant Actuators (LRA)  
*Similar to ERMV but motion is created by an EM field*
  - Piezoelectric Actuators  
*Material reacts to current: expand and contract*
  - TENS
  - Hydraulic
  - Ultrasound (no touch)
- Vibration-based

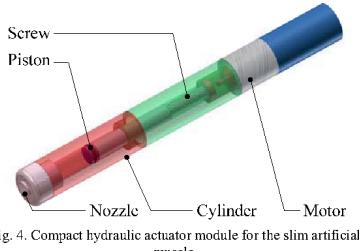
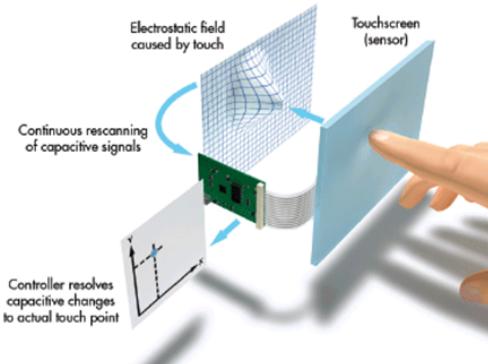


Fig. 4. Compact hydraulic actuator module for the slim artificial muscle.

Invert the previous processes: from motion/pressure to voltage=data

### The most widely-used haptic sensor: **touch screens**

- Resistive
- Capacitive
- Infrared
- Surface acoustic wave
- Acoustic-based



- **Electromyography (EMG)-based**  
Myo armband controller
- **Phantom arm based**  
Phantom Omni (touch feedback)
- **Infrared-based**  
Etee VR controlled, similar to leap-motion



## Omni-directional threadmill (illusion of walking)

Omnipad, Infinadeck, Cyberith

- Used also in CAVE
- Crucial if you do not want to limit the motion possibilities of the user
- Natural navigation of the scene
- Used in rehab as well

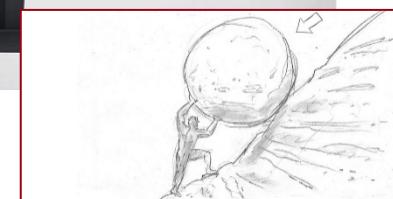


## Haptic rigs (VR env pushes back against you)

- Full-body
- suspension



Axon VR



## Haptic suits

### Tesla

- Movements and commands can be captured by sensors
- Suit provide a feedback to the user: sense of touching
- Temperature control



### *Transcutaneous electrical nerve stimulation (TENS or TNS/EMS)*

- Introduced to relieve pain
- Electrical stimulation by a couple of electrodes
- Stimulating the motor nerves, muscular contractions can be triggered that can be matched to a haptic event
- Avoid sensitive parts like brain, spinal column

Other products ...



Tactile Gaming Vest  
(2010)



SoundShirt (2016)



NeoSensory exoskin  
(2018)

# Other interfaces



## Smelling

Smell-integrating sensors exploit previous knowledge on smell synthesis.

Scentography is the technique of creating and storing odor by artificially recreating a smell using chemical and electronic means.

Previous experience: Smell-O-Vision, AromaRama

(*Scent of Mystery*, Todd Jr., '60; *Behind the Great Wall*, Lizzani, '59)

Cartridges of different smells

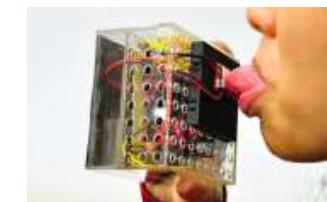


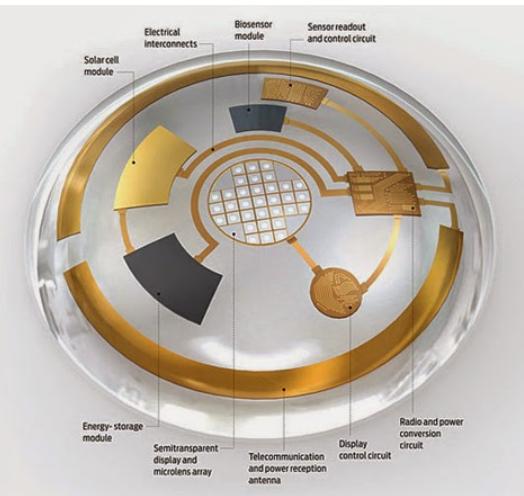
## Taste

Experimental actuators that emulate all the fundamental tastes that our tongue can perceive: sour, salty, bitter, sweet + sensations like mint and spicy.

**electrical stimuli** (the different taste depends on amplitude and frequency of the electrical wave going through the tongue), while **sweet, mint and spicy can be emulated using thermal stimuli** (for instance, a cold metal can simulate mint).

[FeelReal](#)





AR contact lens



Smart Hardhats



more to come ....



Windows (e.g., [WayRay](#), [hyperloop train](#), buildings)

## EXOTIC VISORS

# MACHINE LEARNING IN AR



**Supervised classification or learning:** the classifier is trained assuming that external teaching data are available (i.e., there is a set of training data where at each sample corresponds the correct classification label)

- Very powerful
- Requires extensive dataset and manual labelling of the samples to generate the training data



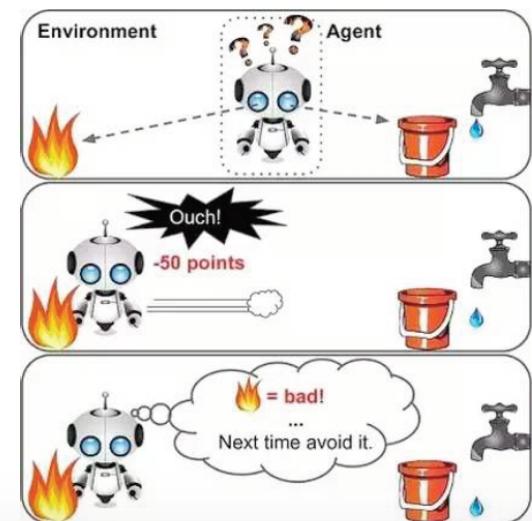
**Unsupervised classification or learning:** the classifier learns the underlying regularities and structure of the data without any prior information about their classes.

- Model building process
- Less accurate than supervised



**Reinforcement learning:** the classifier learns the underlying regularities and structure of the data by actively interacting with the environment.

- Extremely useful
- Computationally demanding



- Basic 3D CV tasks
  - (Feature computation, orientation, localization, 3D mapping, SLAM, depth estimation, image stitching)
- Object classification and identification (3D/2D)
- Gesture/speech recognition (interfaces)
  - (voice commands recognition, hand pose identification)
- Data augmentation and enhancement (Image conversion)
- Optimizing data (coding/compression)
- Enhancing user experience



## EXAMPLE: OBJECT DETECTION

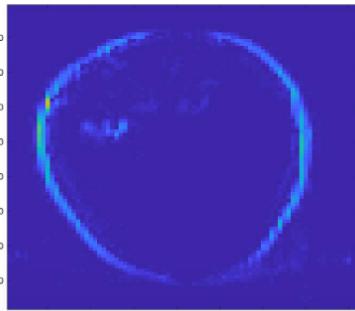
Possible features: vertical and horizontal gradients  
(gradients  $\approx$  object edges)

$$G_x(x, y) = I(x+1, y) - I(x, y) \quad \text{horizontal}$$

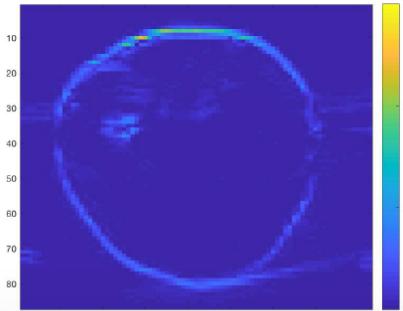
$$G_y(x, y) = I(x, y+1) - I(x, y) \quad \text{vertical}$$



$G_x(x, y)$



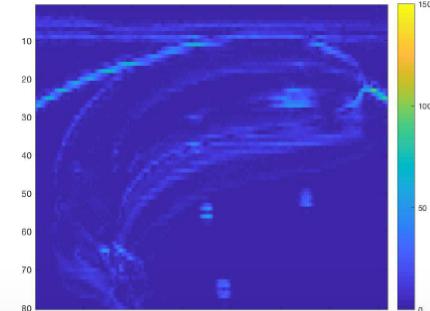
$G_y(x, y)$



$G_x(x, y)$



$G_y(x, y)$



Let's try to classify object (belonging to a finite set of classes) from their images!



Apple

Edge are equally oriented along all the directions

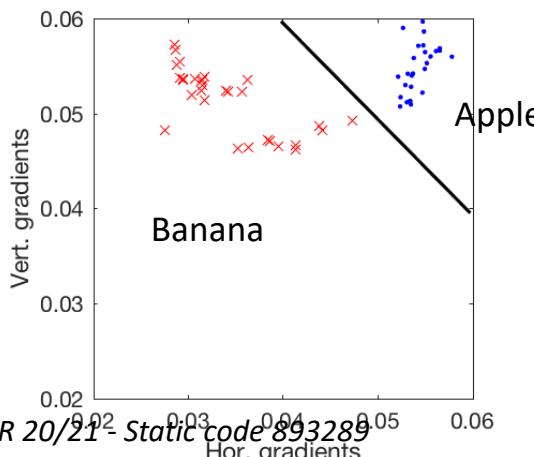
Averages  $x_1$  and  $x_2$  are similar.



Banana

Edge are oriented along a main direction

Averages  $x_1$  and  $x_2$  are not similar.



Assuming that the line equation is  $w_1x_1 + w_2x_2 = b$ , the problem to be solved is the

$$\min \|\mathbf{w}\|^2$$

$$\text{s.t. } y_i (\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1 \quad \forall i$$

Possible features: vertical and horizontal gradients

(gradients  $\approx$  object edges)



$$G_x(x, y) = I(x+1, y) - I(x, y) \quad \text{horizontal}$$

$$G_y(x, y) = I(x, y+1) - I(x, y) \quad \text{vertical}$$

Adopted features

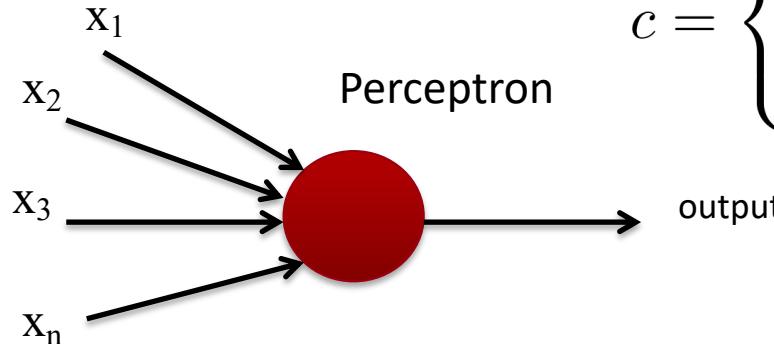
$$x_1 = E[|G_x(x, y)|]_{x, y}$$

$$x_2 = E[|G_y(x, y)|]_{x, y}$$

An  $R^n$  to binary maps

N input features  binary output decision/classification

### Perceptron

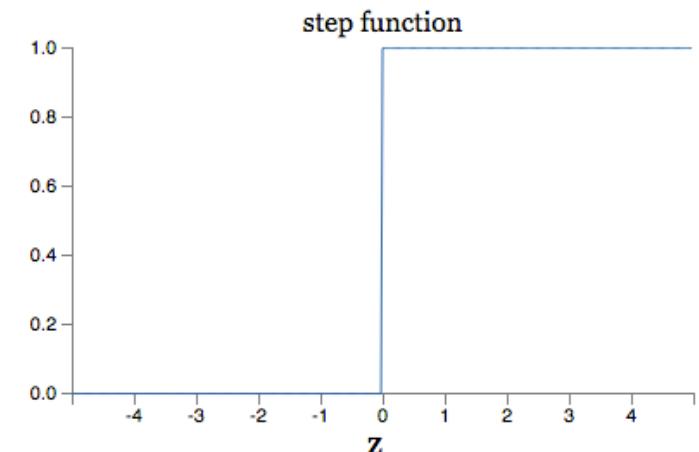


$$c = \begin{cases} 0 & \text{if } \sum_i w_i x_i - b = \mathbf{w}^T \mathbf{x} - b \leq 0 \\ 1 & \text{if } \sum_i w_i x_i - b = \mathbf{w}^T \mathbf{x} - b > 0 \end{cases}$$

It can be decomposed into a linear operator  $z +$   
non-linear decision function  $f()$ : activation function

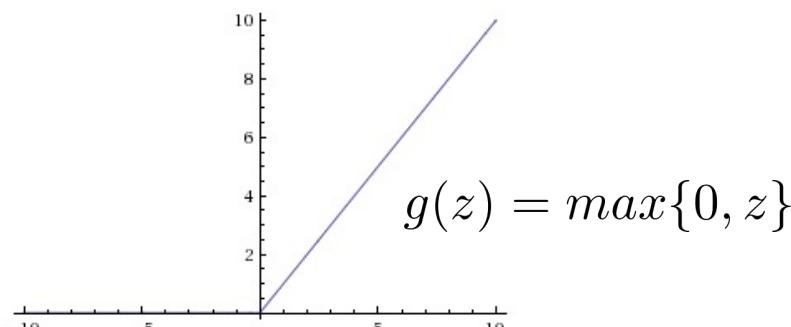
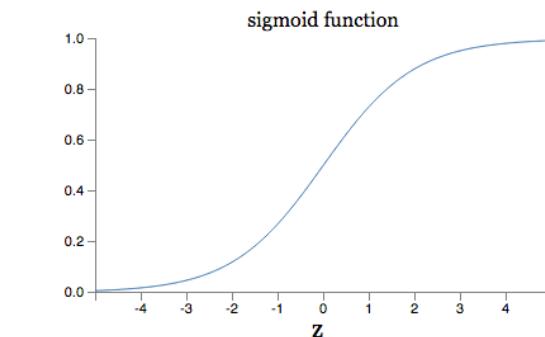
$$z = \sum_i w_i x_i$$

When  $z$  is higher than a threshold  $b$ , the perceptron/neuron fires

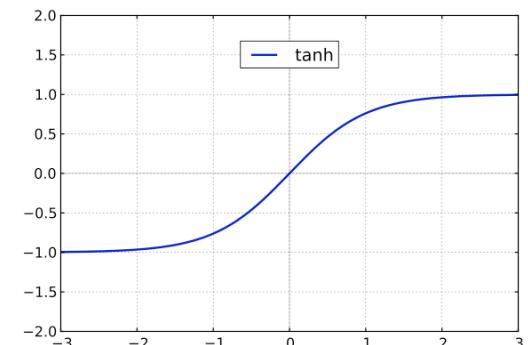


Sometimes a “real” output is required (or simply we would like to make the neuron fire with a different sensitivity).

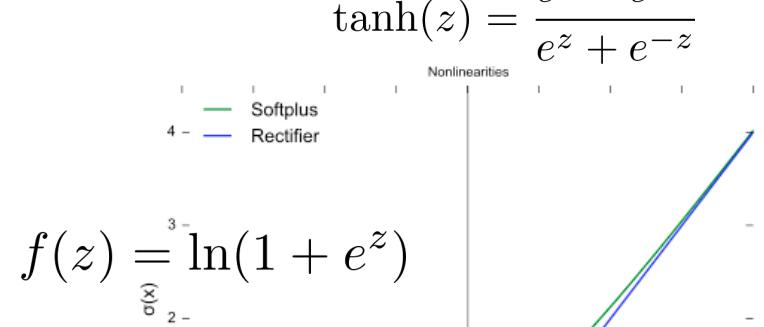
Different activation functions have been introduced.



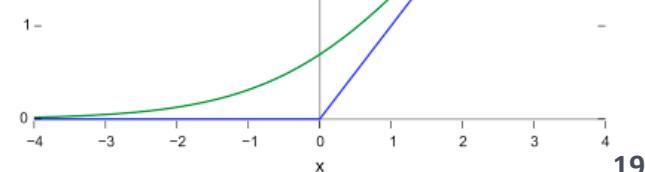
Rectifies Linear Units (ReLU)



$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Softplus



I want an multi-dimensional output.  $\mathbf{z}$  is a vector  $\mathbf{z}=[z_i]$

e.g. given a set of possible output options I want as output the probability of each option.

$$\text{softmax}(\mathbf{z})_j = \frac{\exp(z_i)}{\sum_j \exp z_j}$$



SoftMax	Sigmoid
multi-classification	binary classification
The sum of values will be 1	Values can be real number in [0, 1].
Used in the different layers of neural networks.	Used as activation function while building neural networks.

Logistic regression is a probabilistic, linear classifier.

The probability that  $\mathbf{x}$  belongs to class  $i$  can be expressed as

$$P(y = i|\mathbf{x}, W, b) = \text{softmax}(W\mathbf{x} + b) = \frac{e^{W_i\mathbf{x}+b}}{\sum_j e^{W_j\mathbf{x}+b}}$$

s.t. the prediction becomes

$$\hat{y} = \arg \max_i P(y = i|\mathbf{x}, W, b)$$

Ok, but how to find good values for  $W$  and  $b$ ?

We need a loss function!

likelihood

$$\mathcal{L}(W, b, \mathcal{D}) = \sum_{i=0}^{N-1} \log P(\hat{y} = y_i | \mathbf{x}_i, W, b)$$

Loss or cost function

$$\ell(W, b, \mathcal{D}) = -\mathcal{L}(W, b, \mathcal{D})$$

N: size of training set D

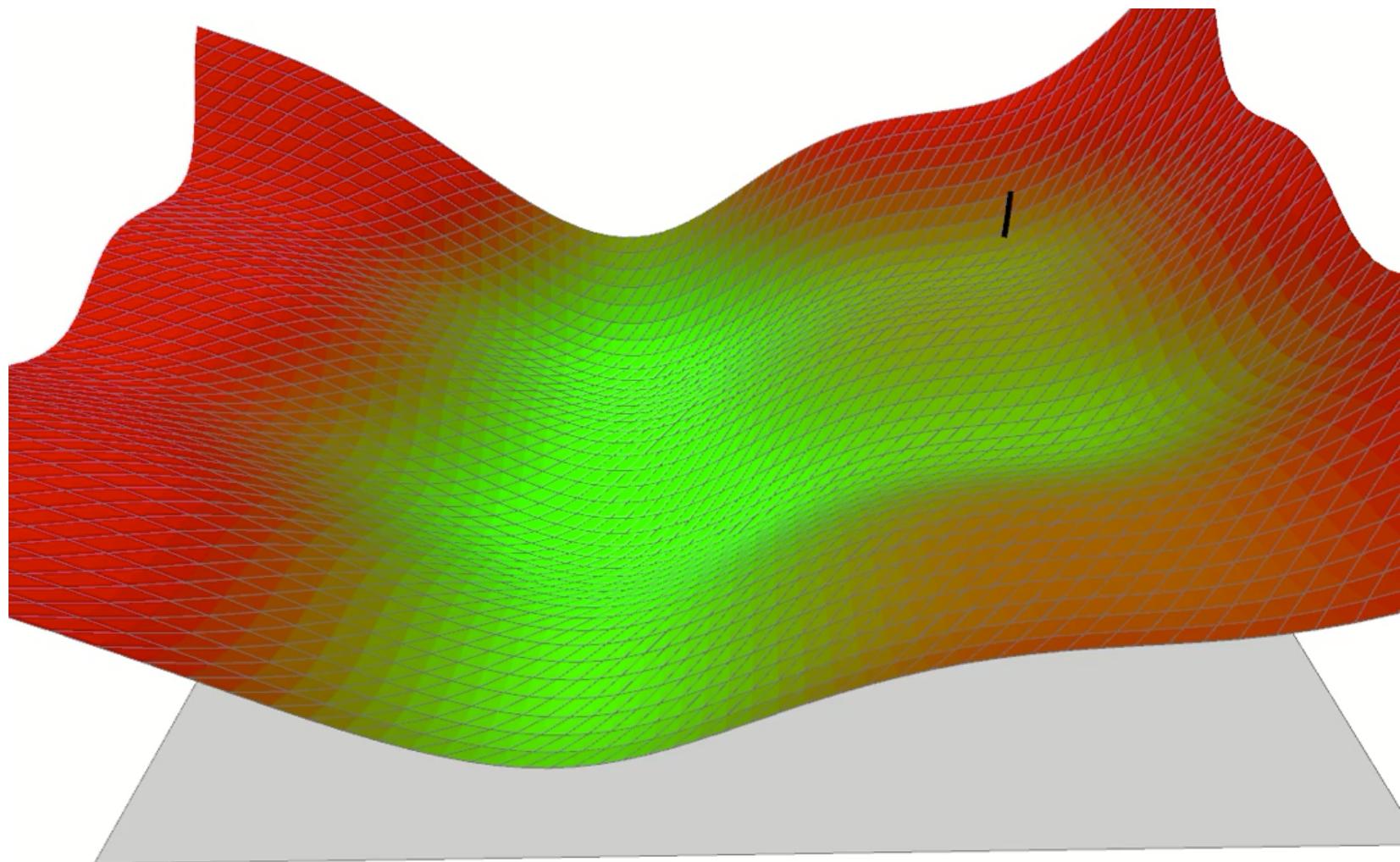
$$f(x_1, \dots, x_n) = f(W, \mathbf{b}) = \max_{W, \mathbf{b}} \ell(W, \mathbf{b}, \mathcal{D}) = - \sum_{i=0}^{N-1} \log P(\hat{y} = y_i | \mathbf{x}_i, W, \mathbf{b})$$

Minimize the function using gradient descend: move in the direction defined by the derivative, ex. 2D case

$$f(x_1, x_2) \quad \longrightarrow \quad \nabla f(x_1, x_2) = \left( \frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2} \right)^T$$

- ① Start at given location  $\mathbf{x}_0 = \begin{bmatrix} x_{0,1} \\ x_{0,2} \end{bmatrix}$
- ② Compute the gradient  $\nabla f(\mathbf{x}_0)$
- ③ Move in the direction opposite to gradient, i.e.,  $\mathbf{x}_1 \leftarrow \mathbf{x}_0 + \alpha \nabla f(\mathbf{x}_0)$
- ④ Iterate steps 1-3 until it does not move anymore!

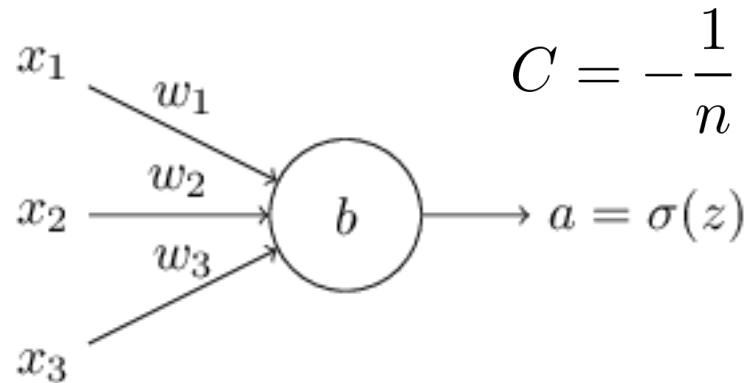
# COORDINATES DESCENT (2/2)



It may happen that learning speed goes down ... the learning is too slow.

How can we address the learning slowdown?

**quadratic cost** -> **cross-entropy**.



$$C = -\frac{1}{n} \sum_i y_i \ln(a(\mathbf{x}_i)) + (1 - y_i) \ln(1 - a(\mathbf{x}_i))$$

$$a(\mathbf{x}_i) = \sigma(\mathbf{w}^T \mathbf{x}_i + b)$$

Partial derivatives then becomes

$$\begin{aligned} \frac{\partial C}{\partial w_j} &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left( \frac{y}{\sigma(z)} - \frac{(1-y)}{1-\sigma(z)} \right) \sigma'(z) x_j \end{aligned}$$

$$\begin{aligned}\frac{\partial C}{\partial w_j} &= \frac{1}{n} \sum_x \frac{\sigma'(z)x_j}{\sigma(z)(1 - \sigma(z))} (\sigma(z) - y) \\ &= \boxed{\frac{1}{n} \sum_x x_j(\sigma(z) - y)}.\end{aligned}$$

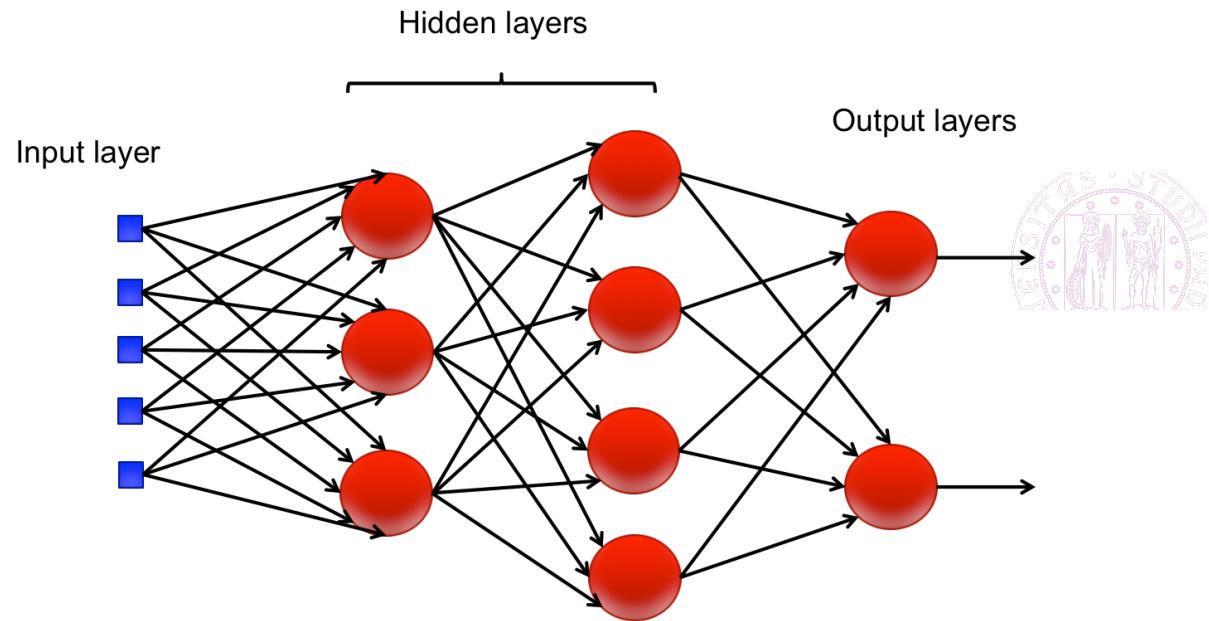
Error defines the learning speed!



It can be generalized to all the neurons in a network!

$$C = -\frac{1}{n} \sum_x \sum_j [y_j \ln a_j^L + (1 - y_j) \ln(1 - a_j^L)].$$

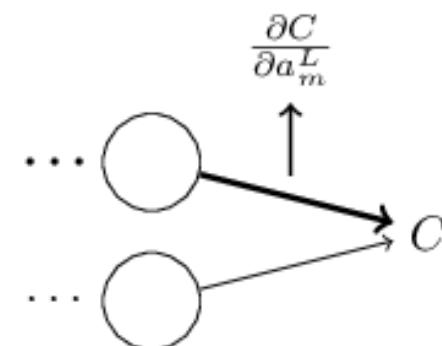
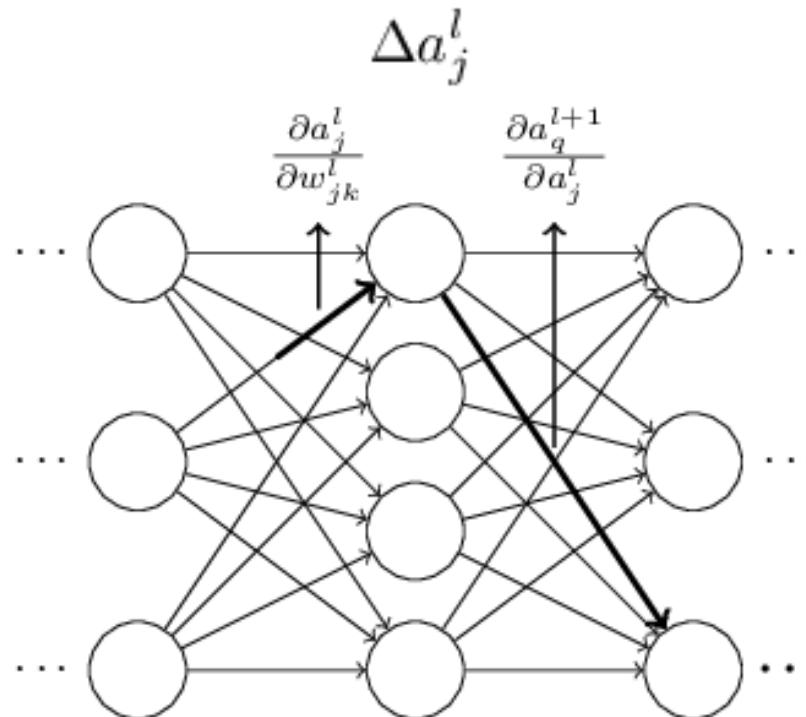
Combining multiple perceptron models together, it is possible to get a neural network.



Error propagation depends on derivatives!

$$a_m^l = f(z_m^l)$$

activation



# HOW TO OPTIMIZE WEIGHTS PARAMETERS?

$$C(\bar{\mathbf{W}}, \bar{\mathbf{b}}) \triangleq \frac{1}{M} \sum_k \|\mathbf{y}(\mathbf{x}_k) - \mathbf{c}_k\|^2$$

Cost function to be minimized:  
difference between real and  
estimated classes.  $\mathbf{y}(\mathbf{x}_k), \mathbf{c}_k$

Minimize C using gradient descend!

We need to compute derivatives!

$$a_m^l = f(z_m^l)$$

activation

$$\mathbf{z}^{l+1} = (W^l)^T \mathbf{a}^l(\mathbf{z}^l) + \mathbf{b}^l,$$

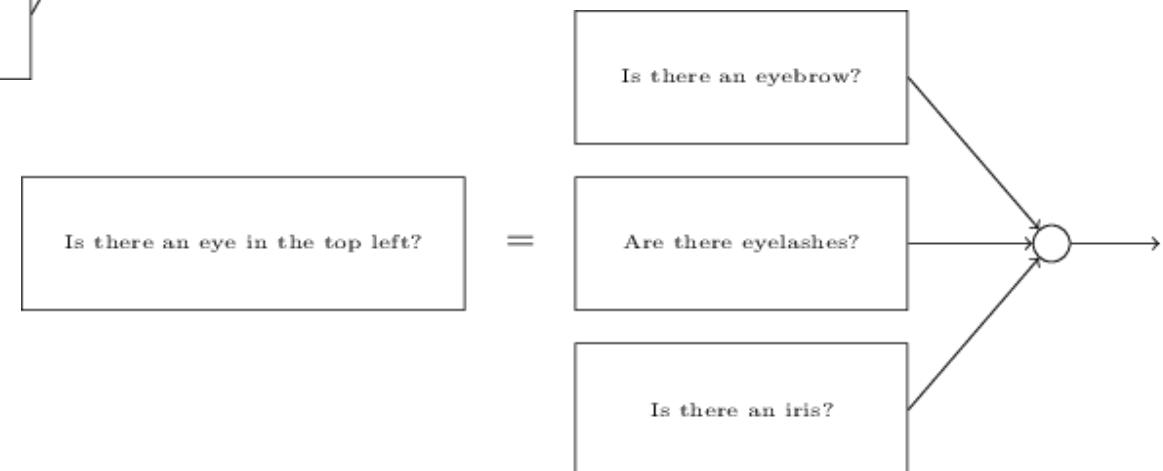
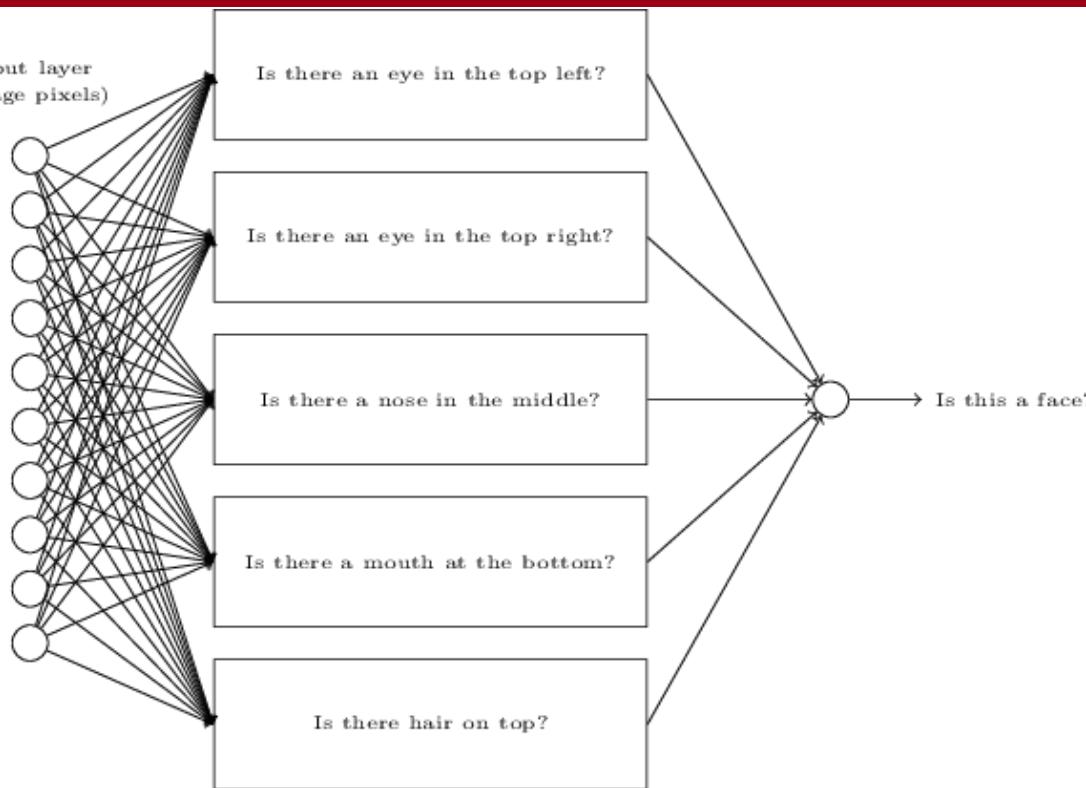
Setting  $\delta_m^l = \frac{\partial C}{\partial z_m^l}$  the last layer allows to write (for quadratic distance)

$$\delta_m^L = \frac{\partial C}{\partial a_m^L} a'(z_m^L) = \sum_k (\mathbf{c}_k - a(\mathbf{z}_m^L)) \odot \mathbf{a}'(\mathbf{z}^L)$$

Then it is possible to back-propagate the derivative to

$$\frac{\partial C}{\partial w_{m,i}^l} = \delta_m^l a_m^{l-1}$$

$$\frac{\partial C}{\partial b_m^l} = \delta_m^l$$



## CONVOLUTIONAL LAYERS (1/2)

It is possible to extend the linear functional for  $z$  using a linear kernel: filtering operations!

Input image: R, G, B, components (padded with 0 on borders)

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

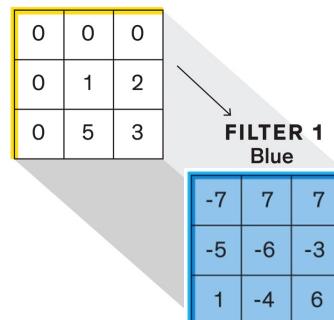
0	0	0	0	0	0	0
0	5	5	5	3	0	
0	7	5	2	3	0	
0	0	6	3	1	0	
0	0	0	0	1	0	
0	0	0	0	0	0	0

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



Filters with a finite support  
(every yellow square outputs one value)



e.g. using the blue filter we have

0	0	0
0	-6	-6
0	-20	18

Total:  
**-14**

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

0	0	0	0	0	0	0
0	5	5	5	3	0	
0	7	5	2	3	0	
0	0	6	3	1	0	
0	0	0	0	1	0	
0	0	0	0	0	0	0

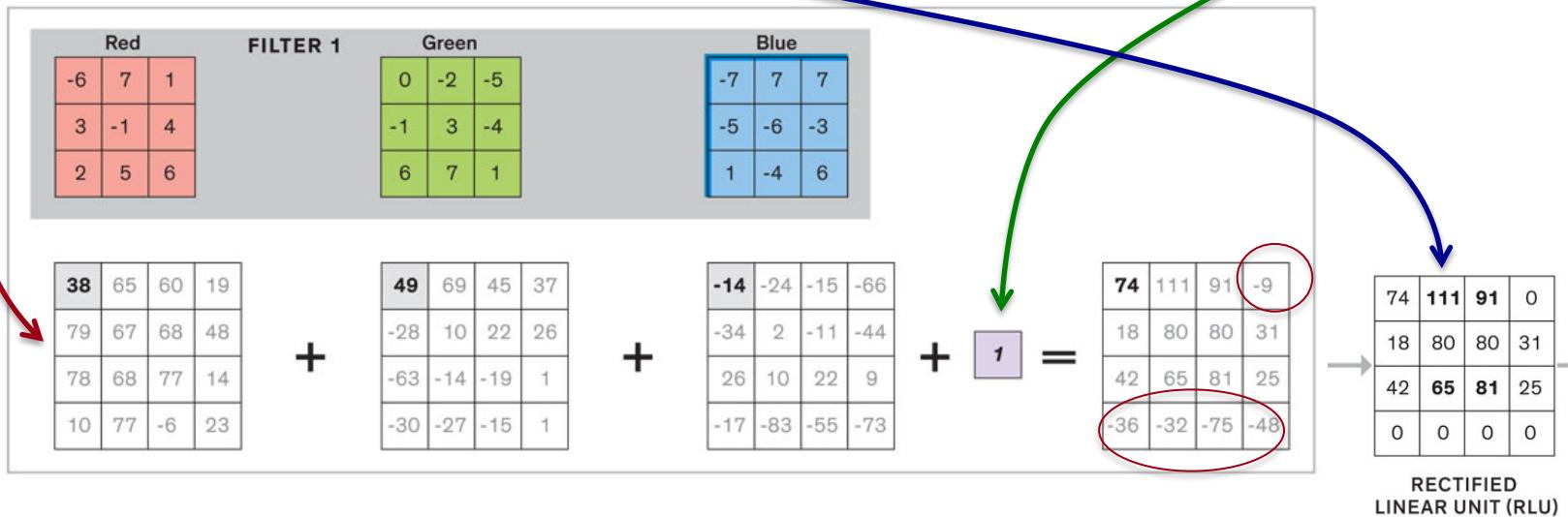
  

0	0	0	0	0	0	0
0	1	2	4	4	0	
0	5	3	2	6	0	
0	4	0	2	0	0	
0	6	3	7	4	0	
0	0	0	0	0	0	0

0	0	0	0	0	0	0
0	5	5	5	3	0	
0	7	5	2	3	0	
0	0	6	3	1	0	
0	0	0	0	1	0	
0	0	0	0	0	0	0

The operations are repeating  
sliding the support window  
(filtering)

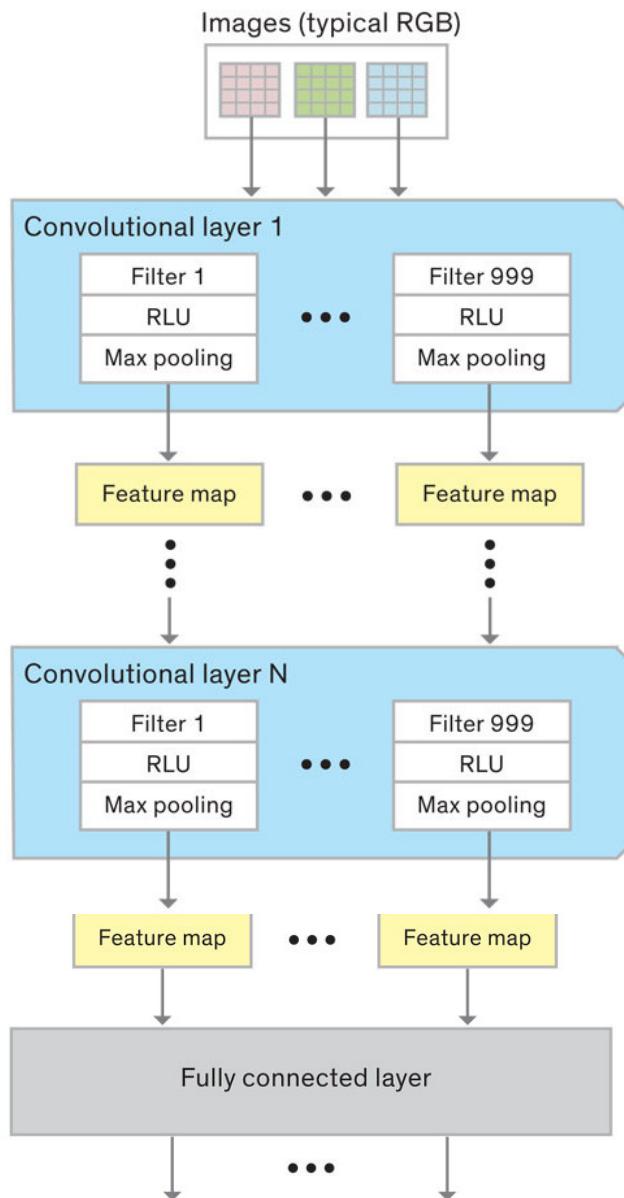
We obtain **filtered components**, which are then shifted by a **constant** and modified by the activation function (ReLU)



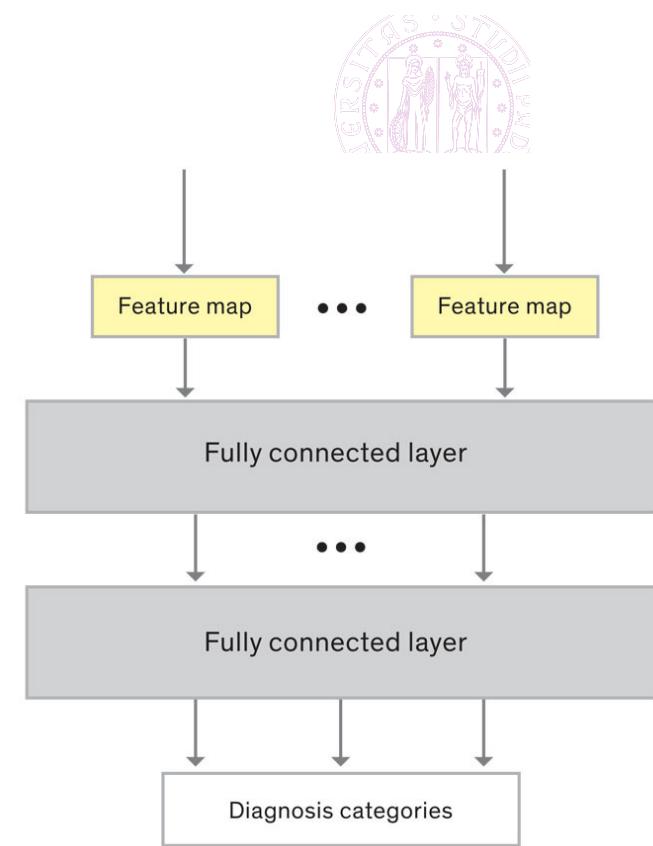
Then, max-pooling operation sample the input image into a smaller one, which is passed to the next layer.

Several filters are used for each color component.

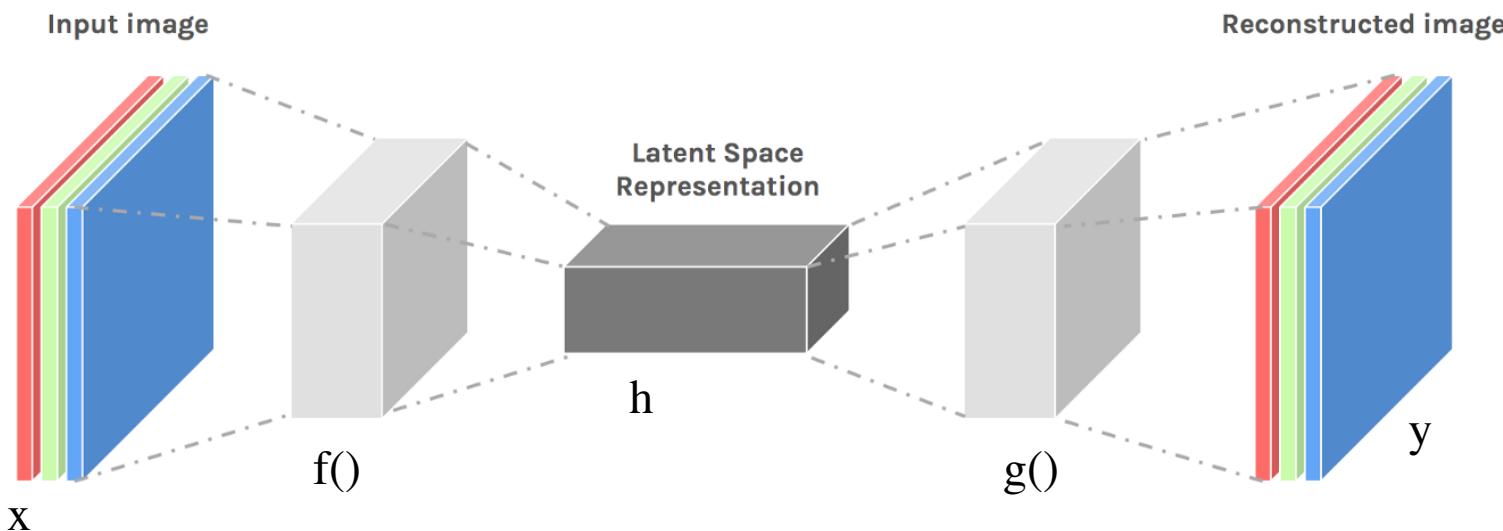
# CONVOLUTIONAL NEURAL NETWORK



Several layers can be combined (with different input sizes), terminated by a fully connected layer (1 or more)



- NN trained to copy its input in the output  $g(f(x))=x$
- Created in order to represent the input data by a set of features  $h$  (with a lower dimension)

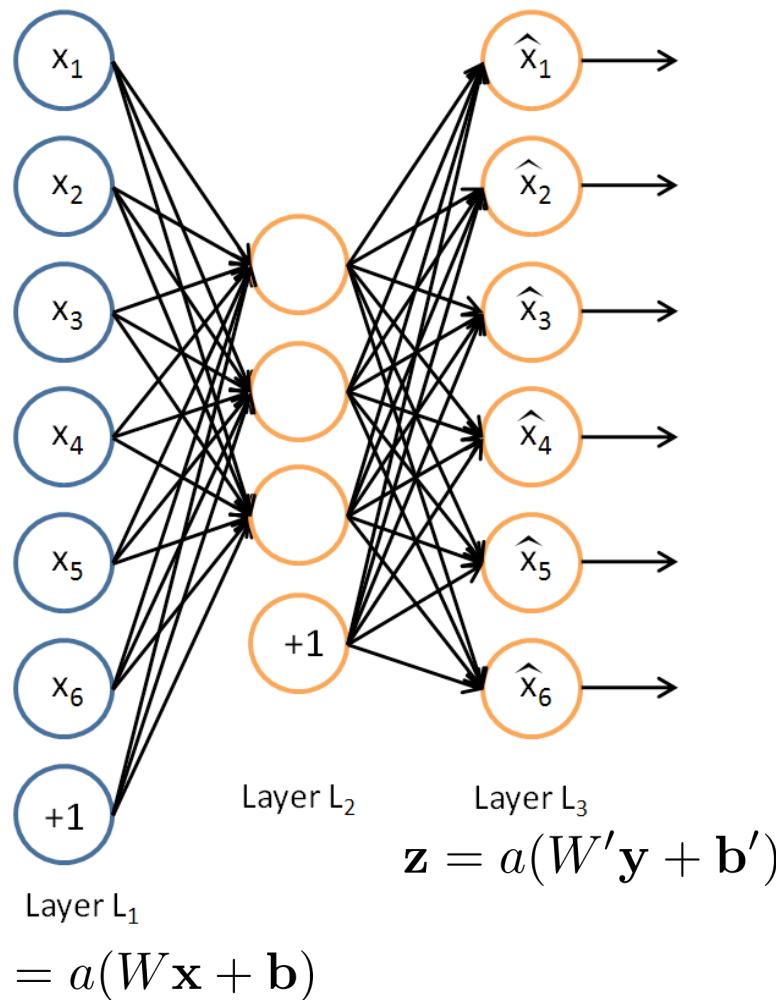


Minimize the cost function

$$\mathcal{L}(\mathbf{x}, g(f(\mathbf{x})))$$

In case a sparsity constraint is added

$$\mathcal{L}(\mathbf{x}, g(f(\mathbf{x}))) + \Omega(\mathbf{h})$$



Minimize the distortion or cost function

$$L(\mathbf{x}, \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|^2$$

or

$$L(\mathbf{x}, \mathbf{z}) = - \sum_{k=1}^d x_k \log z_k + (1 - x_k) \log(1 - z_k)$$



Hidden layer define a set of features for the input data

Can be used to denoise input-data:  
exploits the regularity of statistics.

## THE SIMPLEST POSSIBLE KERAS AUTOENCODER

```
#import declarations
import keras
from keras import layers

# This is the size of our encoded representations (in floats)
encoding_dim = 32    # 32 floats

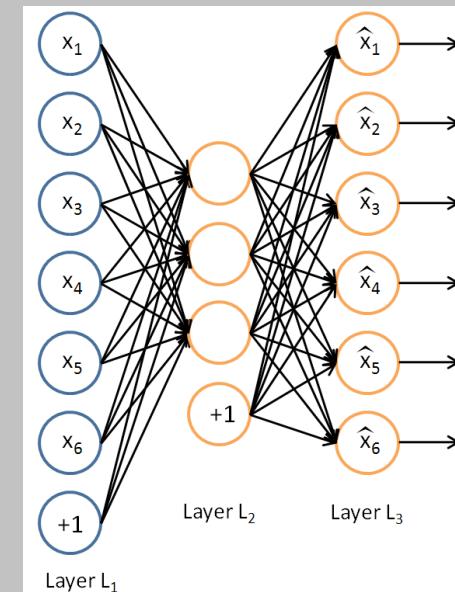
# This is our input vector
input_vet = keras.Input(shape=(64,))

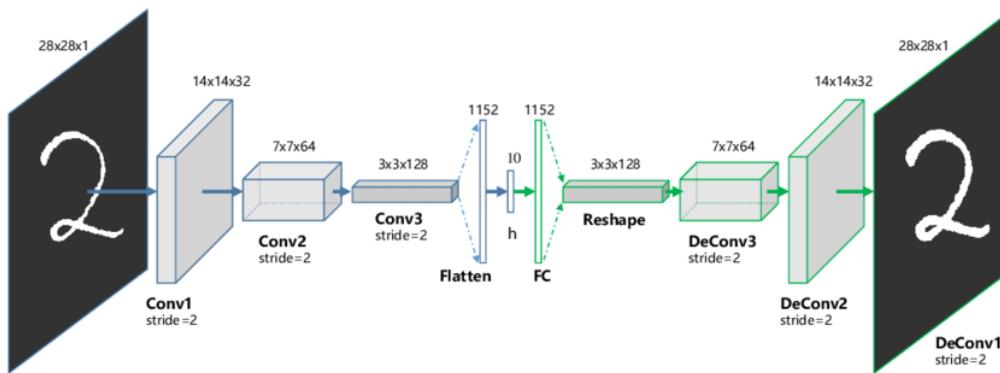
# "encoded" is the coded representation: hidden features
encoded = layers.Dense(encoding_dim, activation='relu')(input_img)
# "decoded" is the lossy reconstruction
decoded = layers.Dense(64, activation='sigmoid')(encoded)

# This model maps an input to its reconstruction
autoencoder = keras.Model(input_vet, decoded)
#Create the autoencoder
autoencoder.compile(optimizer='adam', loss='binary_crossentropy')

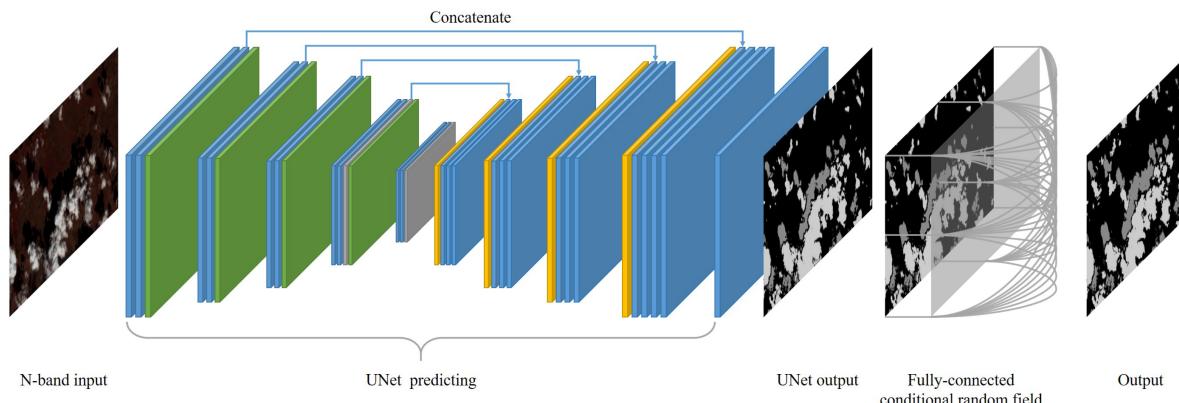
#Train the autoencoder
autoencoder.fit(x_train, x_train, epochs=50, batch_size=256, shuffle=True, validation_data=(x_test, x_test))

#Reconstruct the vector
decoded_vets = autoencoder.predict(x_test)
```

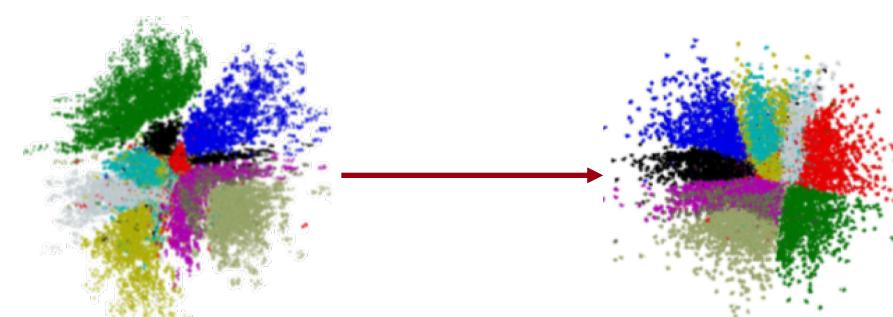




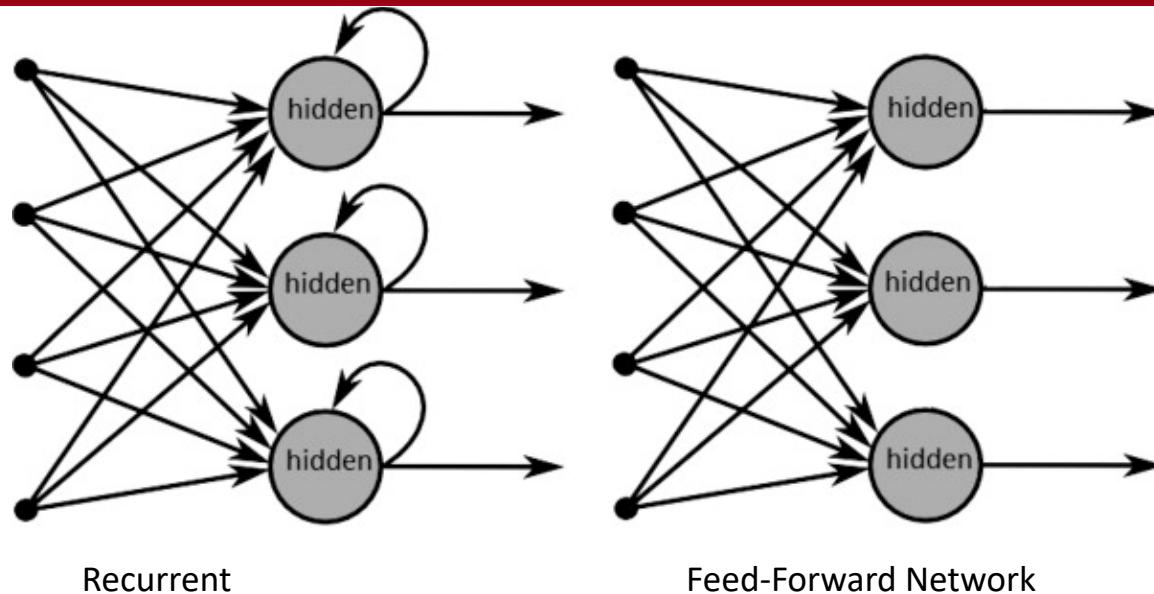
Convolutional



U-Net

Adversarial  
Autoencoder

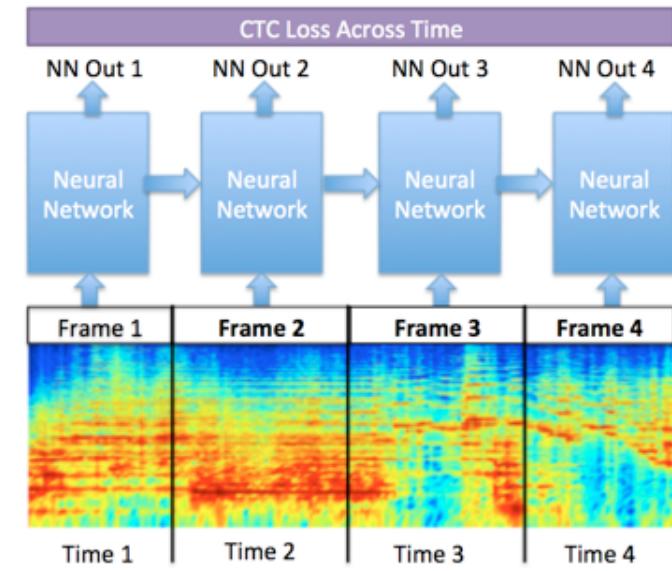
## RECURRENT NEURAL NETWORKS (RNN) -1/2



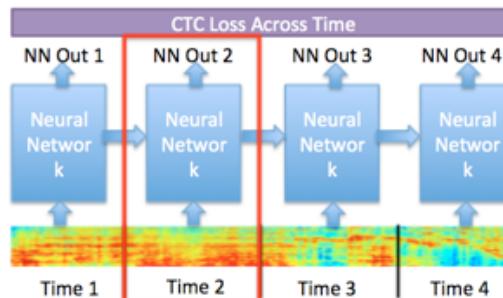
The input of a neuron depends on its output (at different time instants): cyclical connections

Example: speech recognition

Input: arrays of variables (samples) distributed over time.  
Output: detected digit

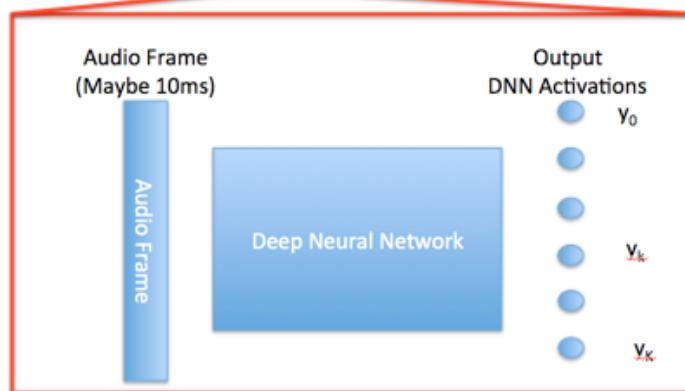


# RECURRENT NEURAL NETWORKS (RNN) -2/2



Feedback from previous instants conditions the output for the next instants.

There is a memory process in the detection.



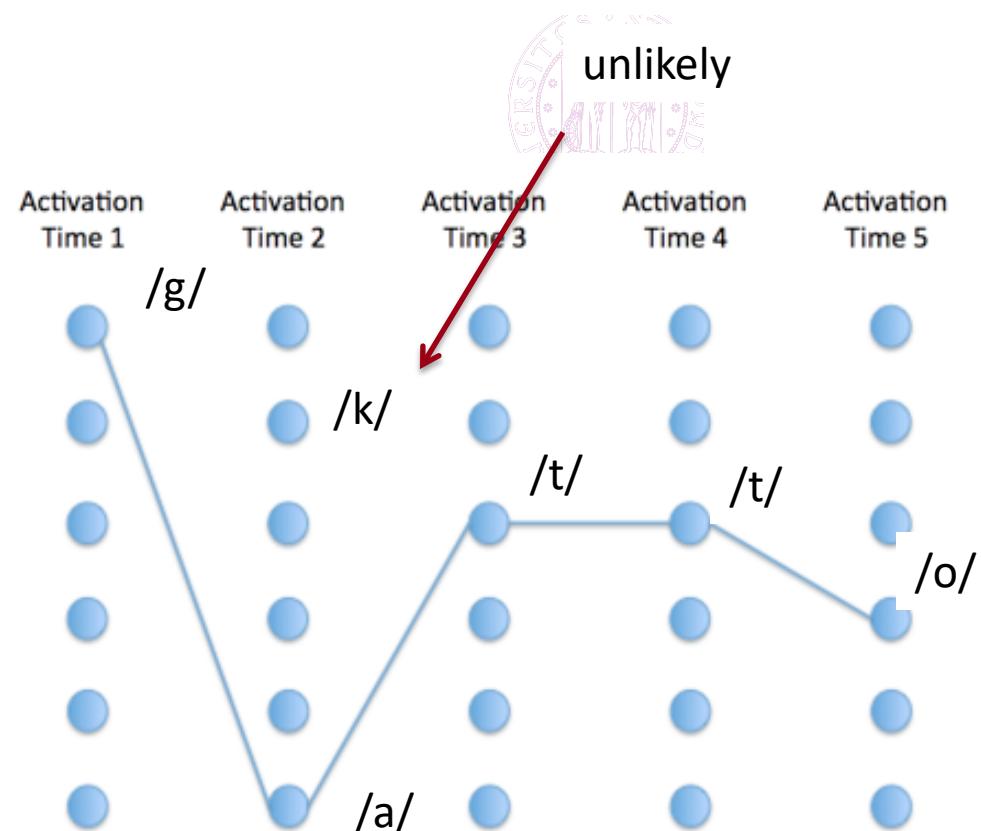
$$\mathbf{y}(t) = \sigma (\mathbf{W}^L \mathbf{z}^L(t) + \mathbf{b}^L)$$

Elman's Network

$$\mathbf{z}_k^l(t+1) = \sigma (\mathbf{W}_k^l \mathbf{z}_k^{l-1}(t+1) + \mathbf{U}_k^l \mathbf{z}_k^l(t) + \mathbf{b}_k^l)$$

Jordan's Network

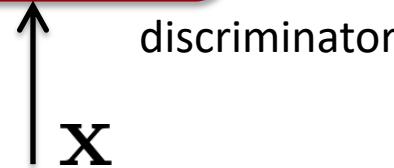
$$\mathbf{z}_k^l(t+1) = \sigma (\mathbf{W}_k^l \mathbf{z}_k^{l-1}(t+1) + \mathbf{U}_k^l \mathbf{y}(t) + \mathbf{b}_k^l)$$



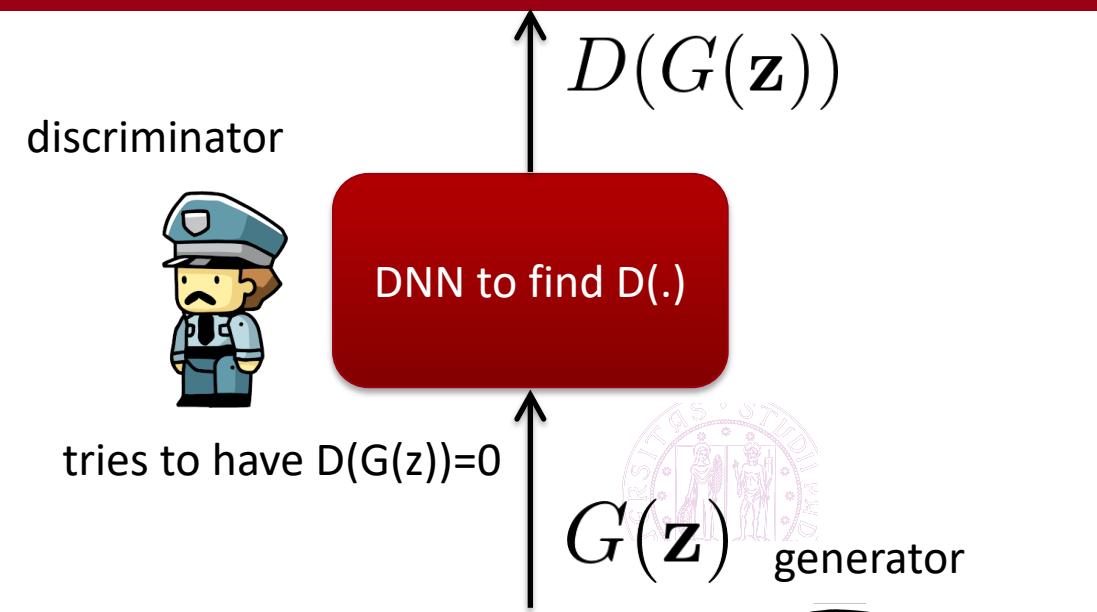
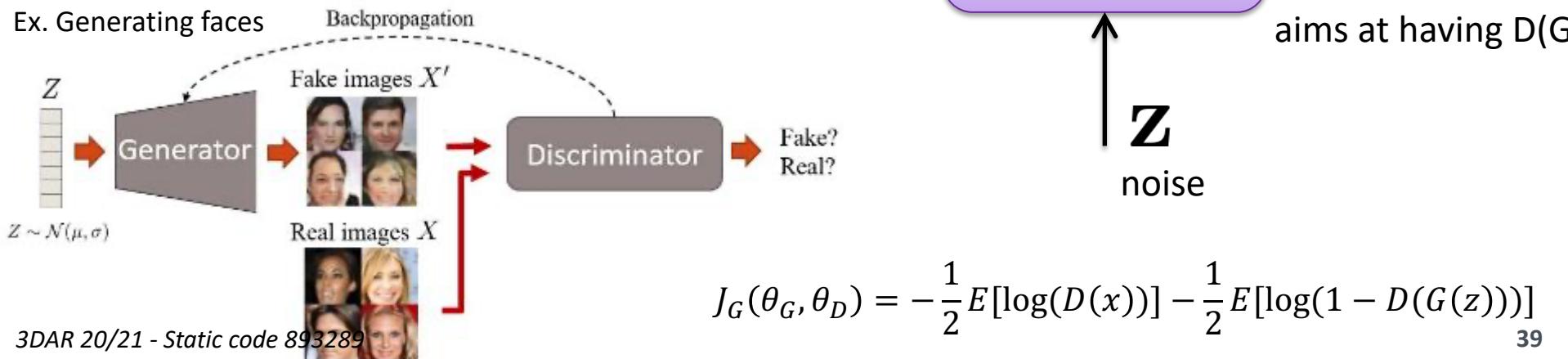
# GENERATIVE ADVERSARIAL NETWORKS (GANs)

$$J_D(\theta_G, \theta_D) = -J_D(\theta_G, \theta_D)$$

Real  $D(x) = 1$   
 or  
 Fake  $D(x) = 0$



Ex. Generating faces



$$D(G(z))$$

$$G(z)$$

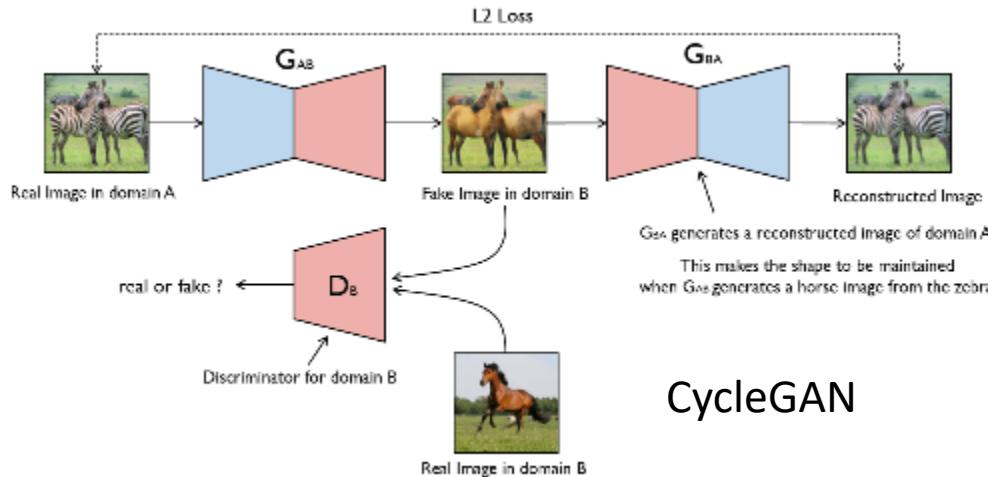


$$J_G(\theta_G, \theta_D) = -\frac{1}{2} E[\log(D(x))] - \frac{1}{2} E[\log(1 - D(G(z)))]$$

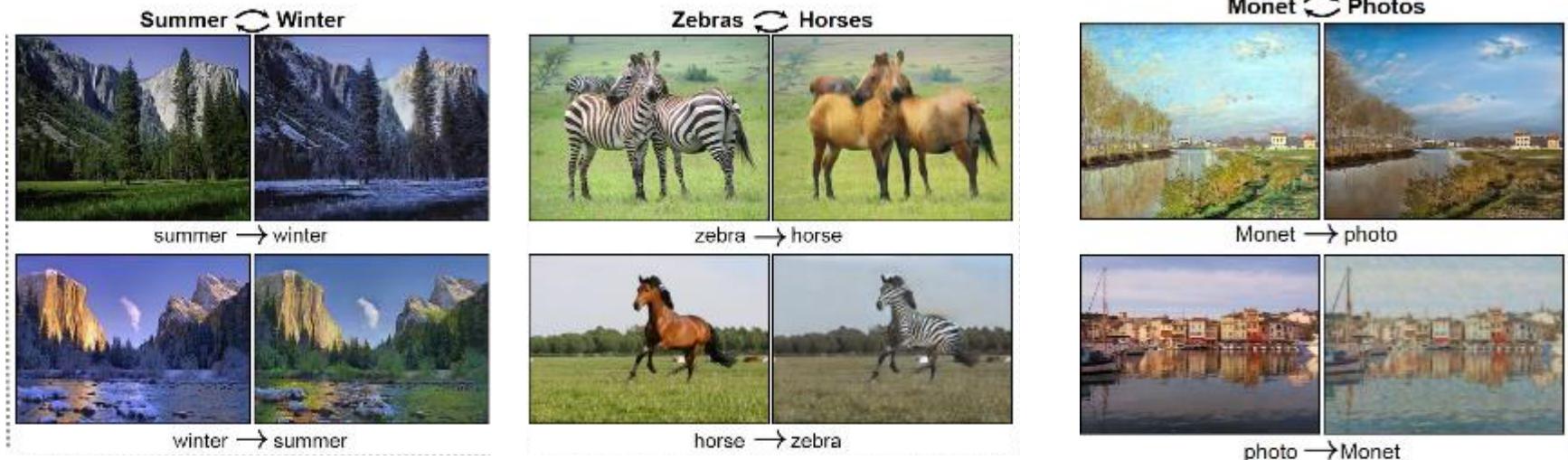
## GENERATIVE ADVERSARIAL NETWORKS (GANs)



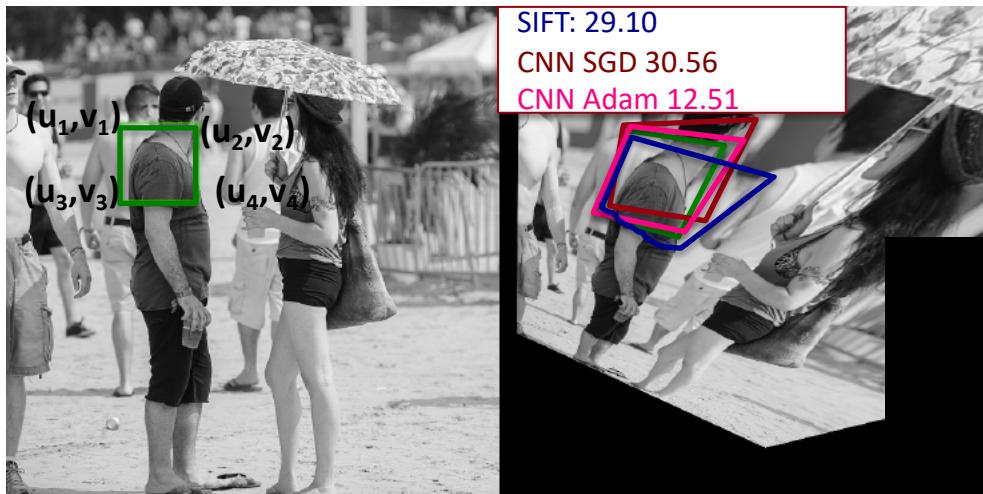
GANs can be used to map an image or face over another



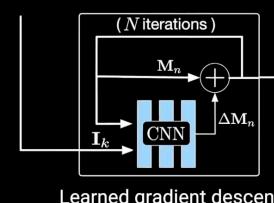
CycleGAN



- Panoramic video and image stitching
- Feature computation, orientation, localization, 3D mapping, SLAM (see SfM lecture)
- Depth estimation
- 3D reconstruction and rendering
- Pose/Gaze estimation
- Scene composition



A sparse set of input images from different viewpoints



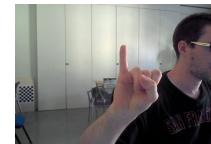
Learned gradient descent

- Hand pose estimation
- Gaze estimation
- Voice processing
- BCI

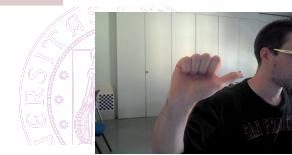
	G1	G10	G11
G1	91.04%	0.83%	8.13%
G10	0%	88.96%	11.04%
G11	0.21%	5.83%	93.86%



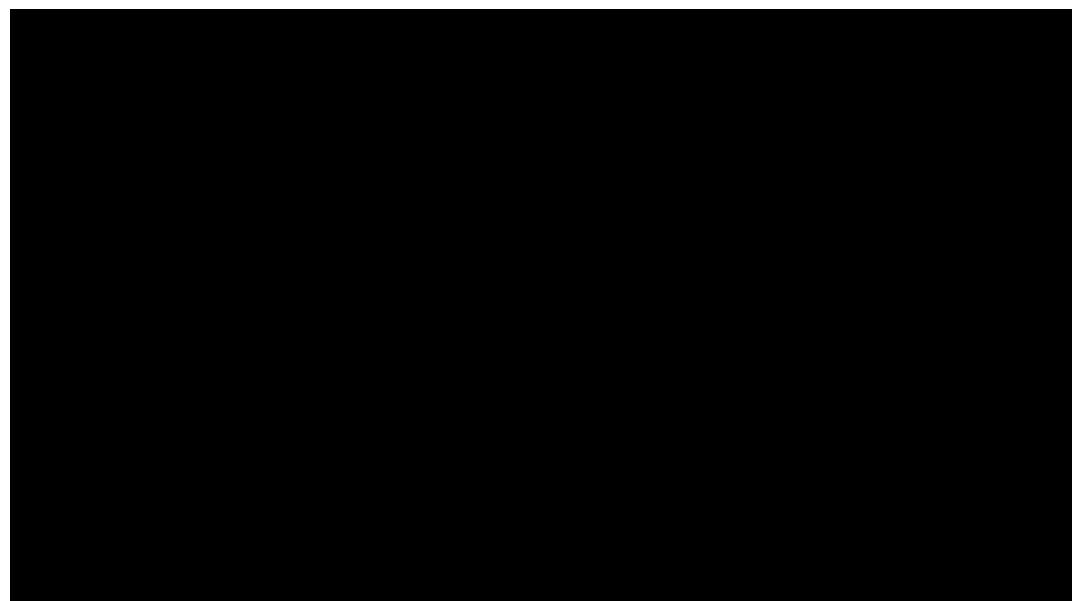
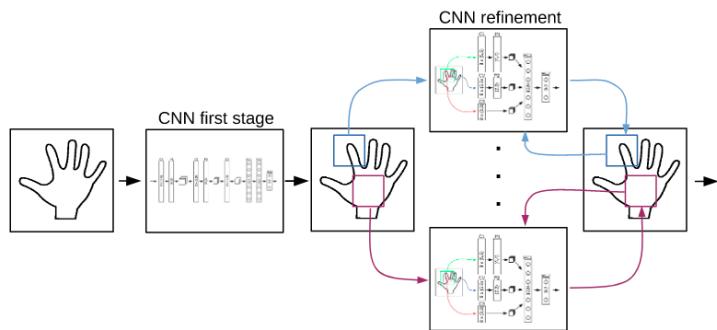
G1



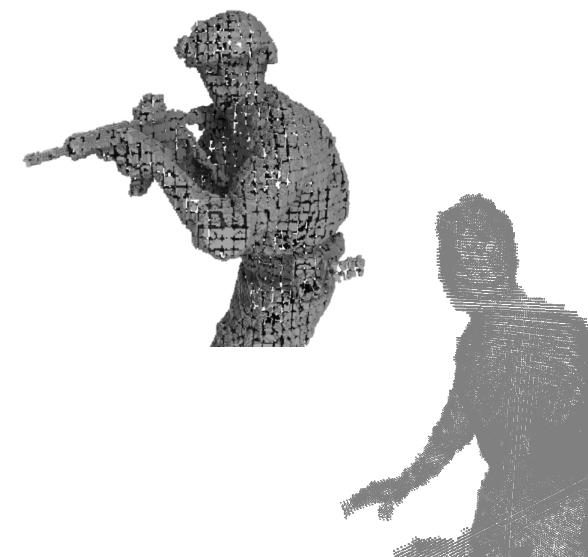
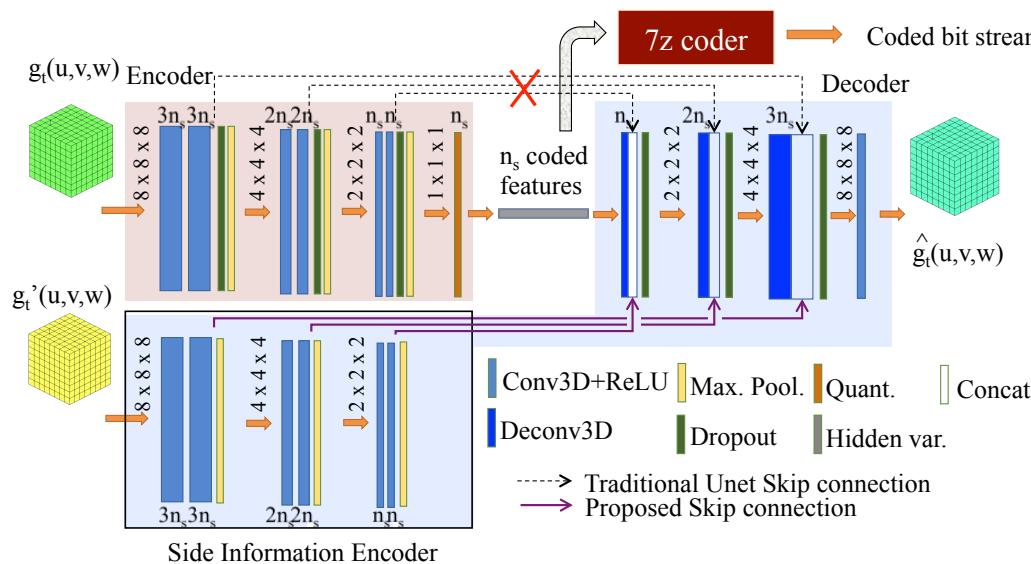
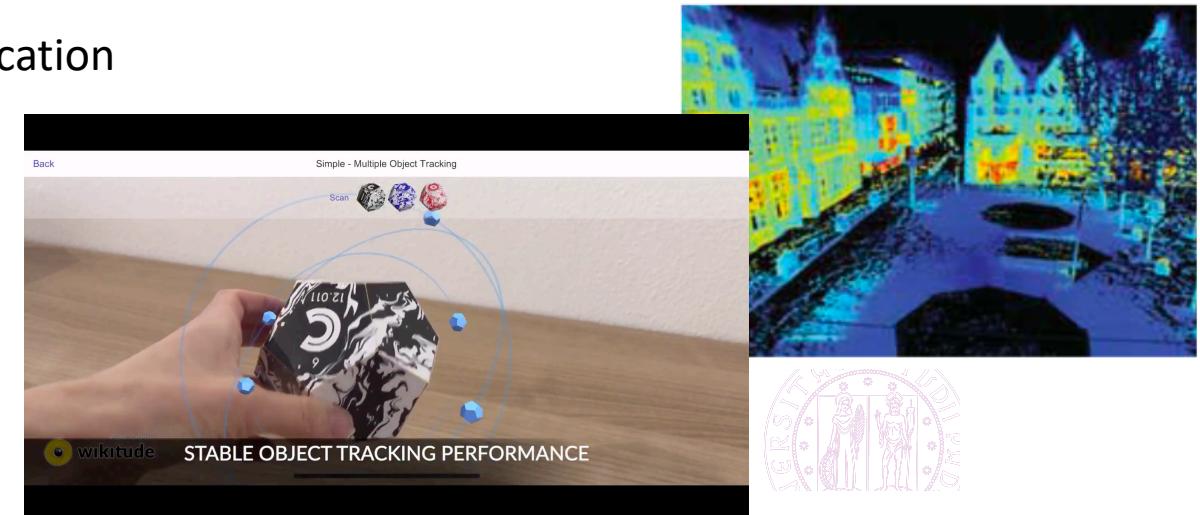
G11



G10



- Object detection and classification
- Semantic segmentation
- Point cloud processing



- Model rendering/coding and optimization
- QoE, Cybersickness management
- Image enhancement

e.g. GAN-based real-time food-to-food translation

