

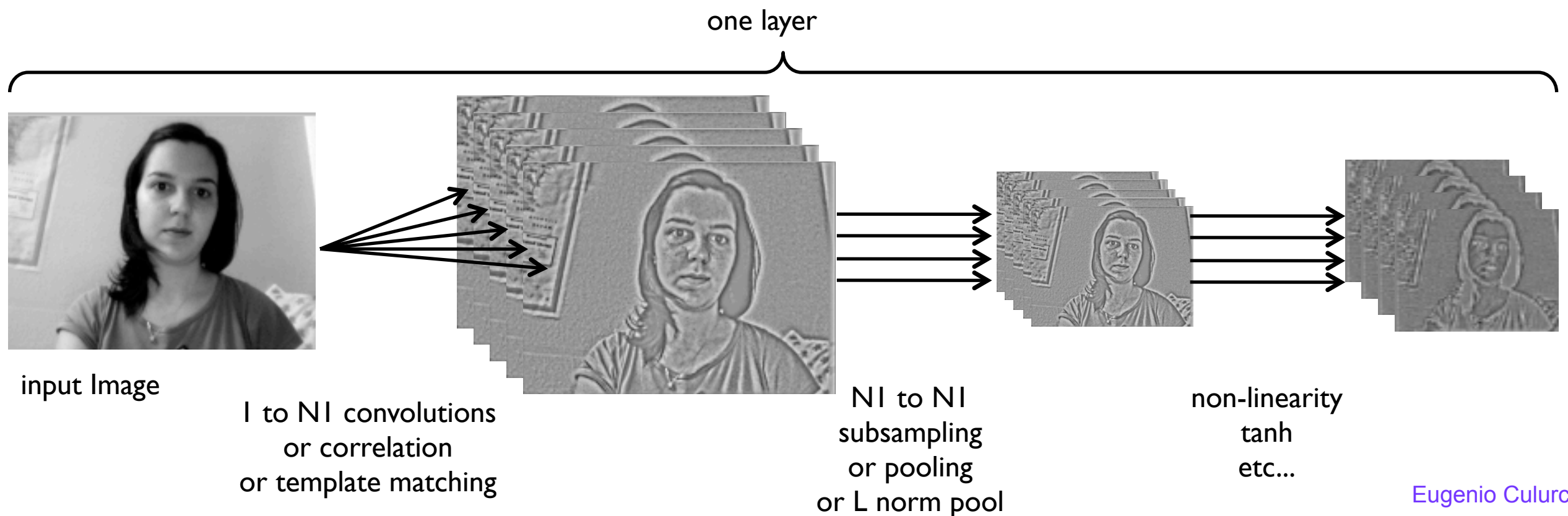
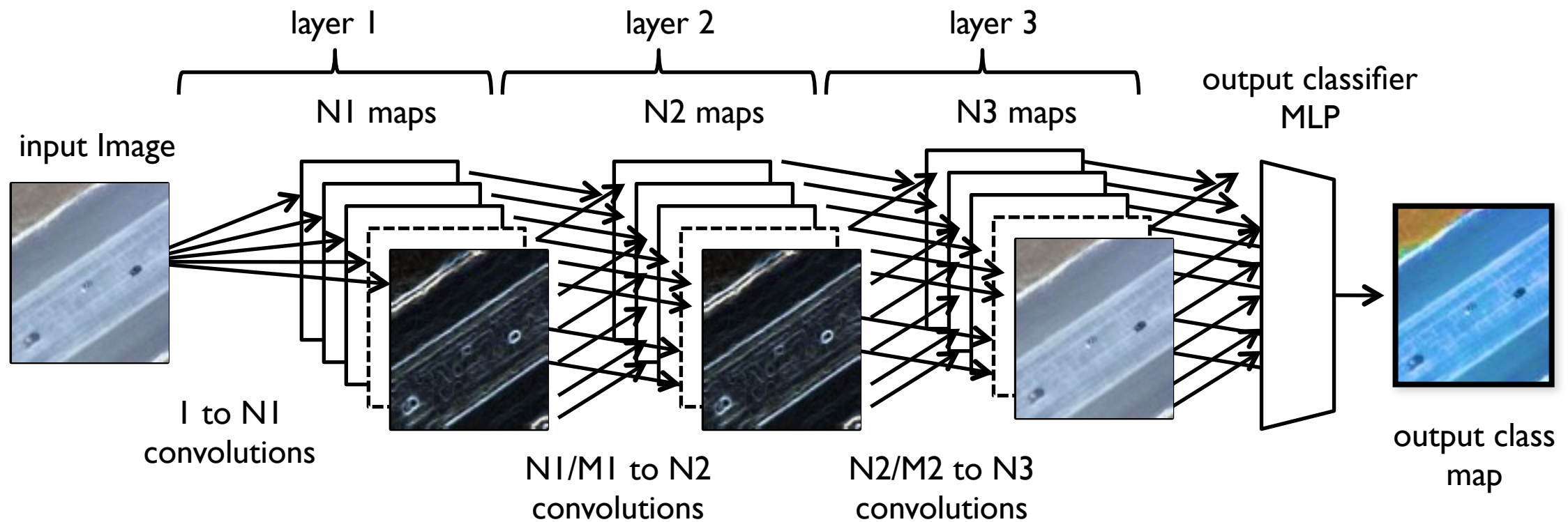
Artificial and robotic vision



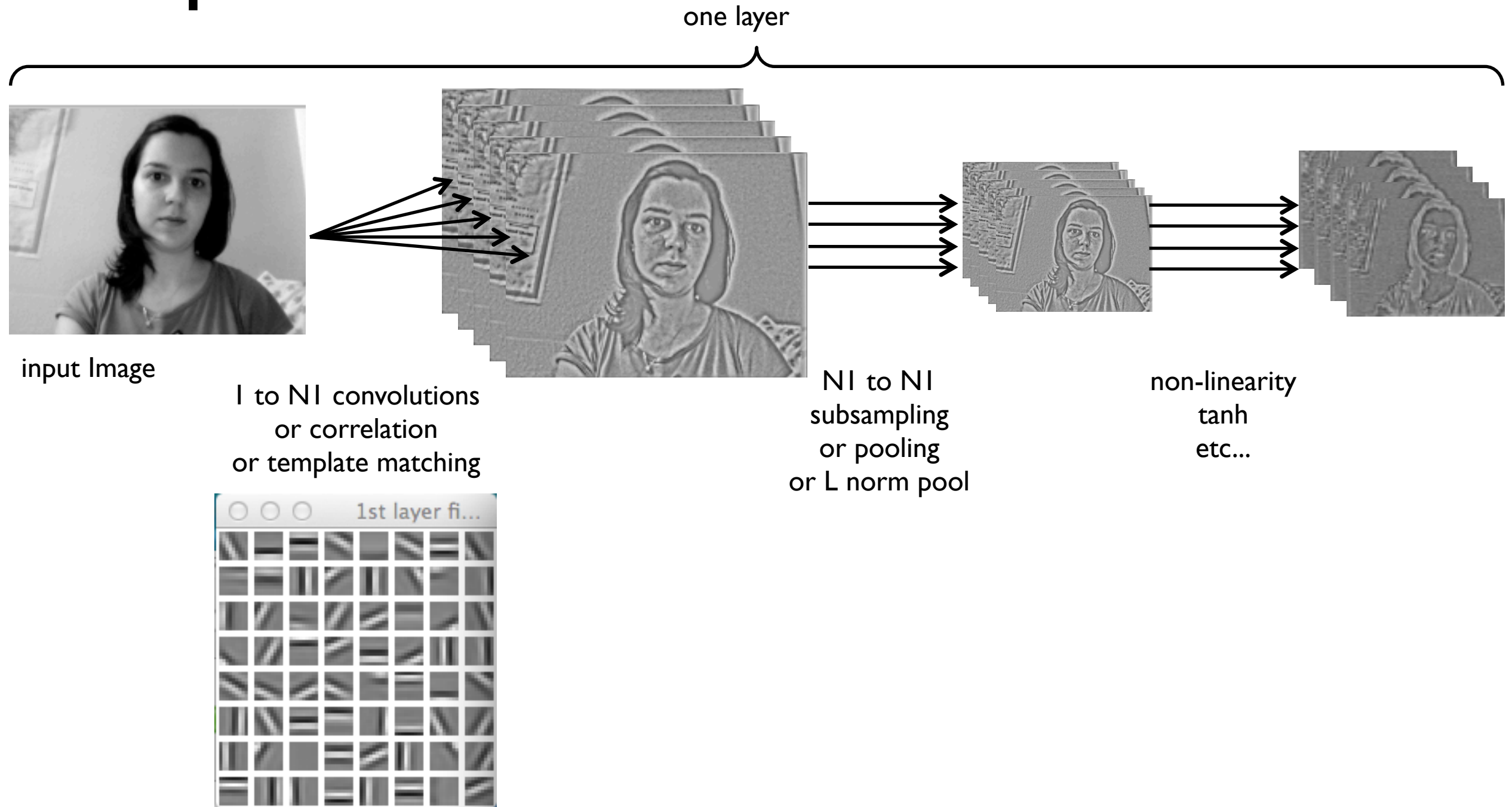
Spring 2013

Lecture 6: unsupervised learning

deep networks



deep networks

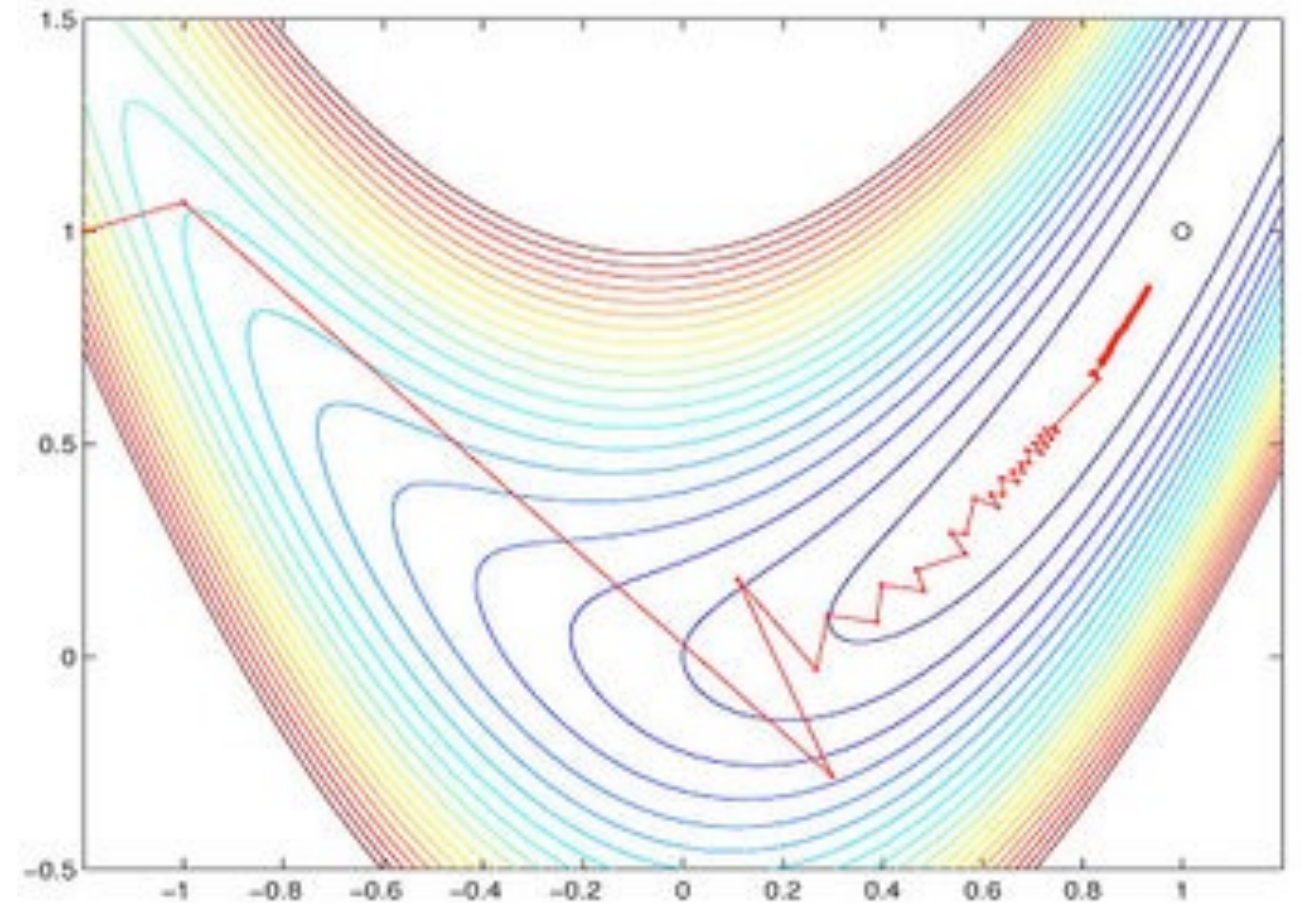


How do we compute these filters?

supervised training

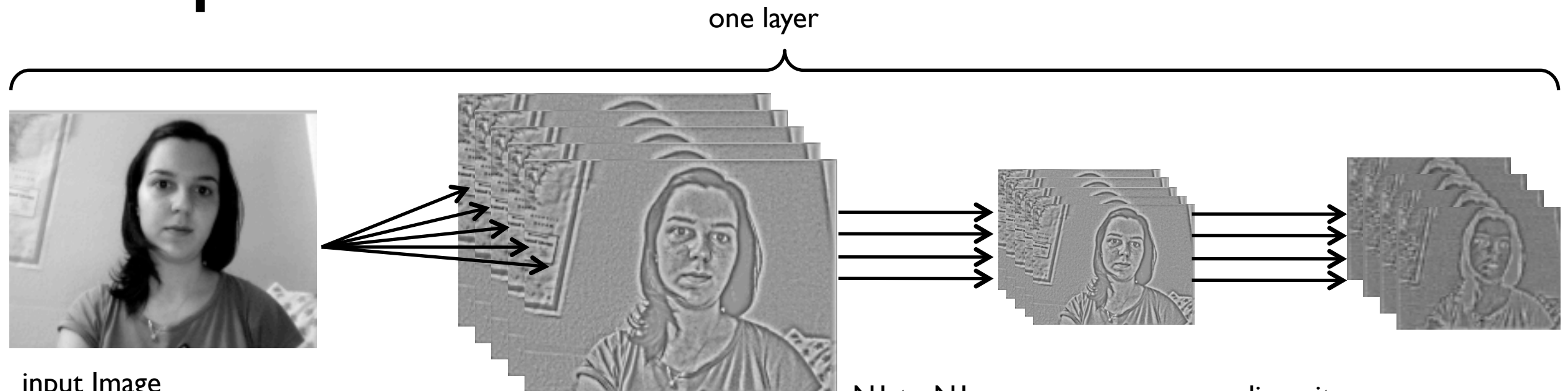


dataset



stochastic gradient descent

deep networks

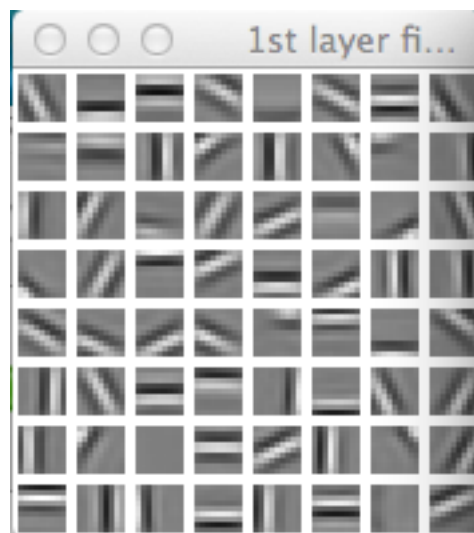


input Image

I to NI convolutions
or correlation
or template matching

NI to NI
subsampling
or pooling
or L norm pool

non-linearity
tanh
etc...



gradient descent, min cost function

$$\underset{W}{\text{minimize}} \quad \lambda \sum_{t=1}^{N-1} \|\mathbf{p}^{(t)} - \mathbf{p}^{(t+1)}\|_1 + \sum_{t=1}^N \|\mathbf{x}^{(t)} - W^T W \mathbf{x}^{(t)}\|_2^2 + \gamma \sum_{t=1}^{\bar{N}} \|\mathbf{p}^{(t)}\|_1$$

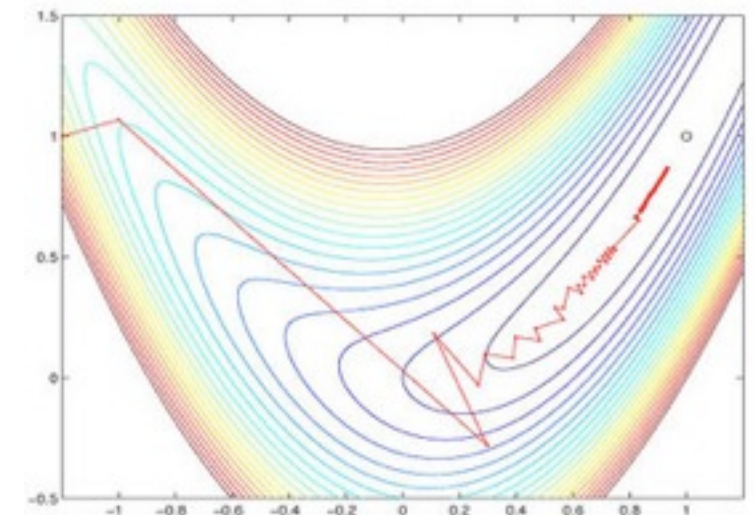
p = features, x = input

supervised training **issues:**



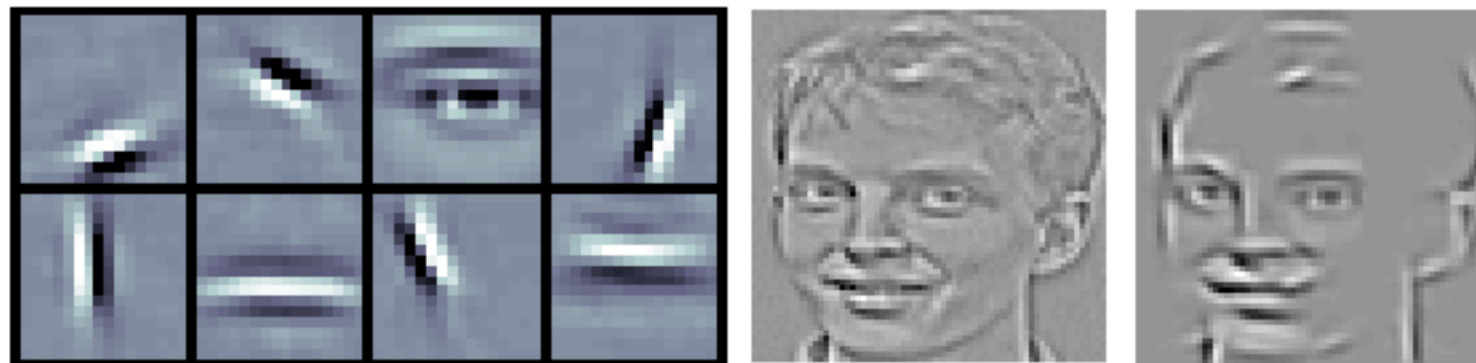
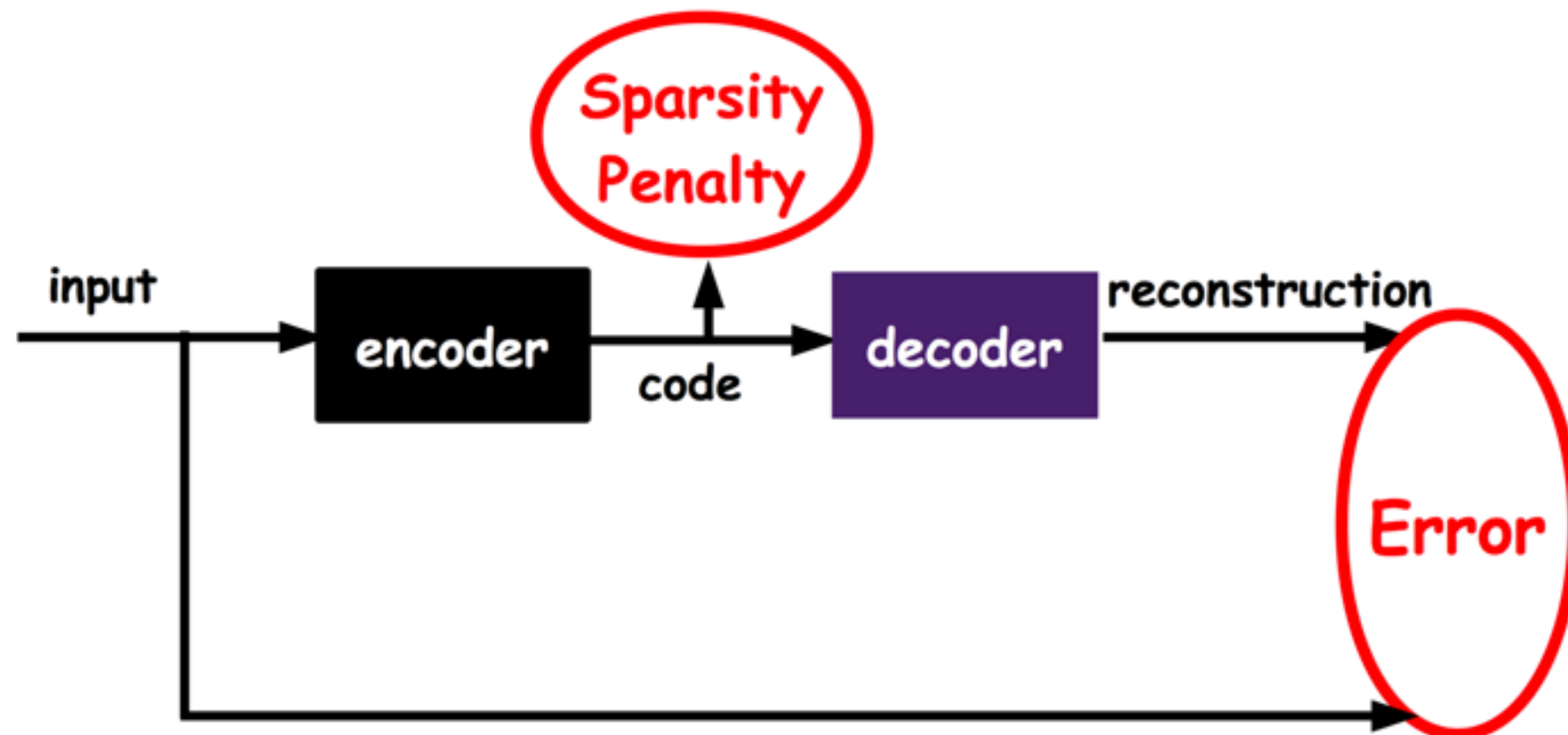
- need labeled data!
- takes a long time
- lots of effort
- in videos it is crazy
- we cannot scale up!

- ~all use gradient descent
- not related to learning in the brain
- use global error propagation
- not local err. prop.
- math-heavy techniques
- they take a long time to compute

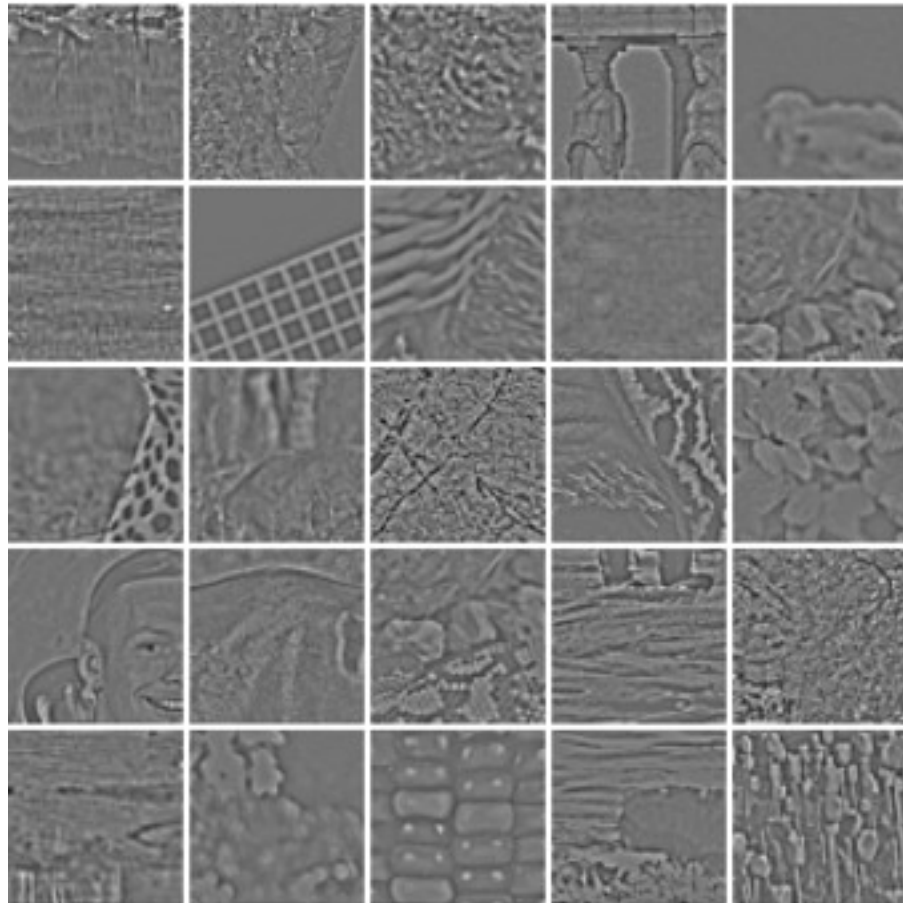


unsupervised training: autoencoders

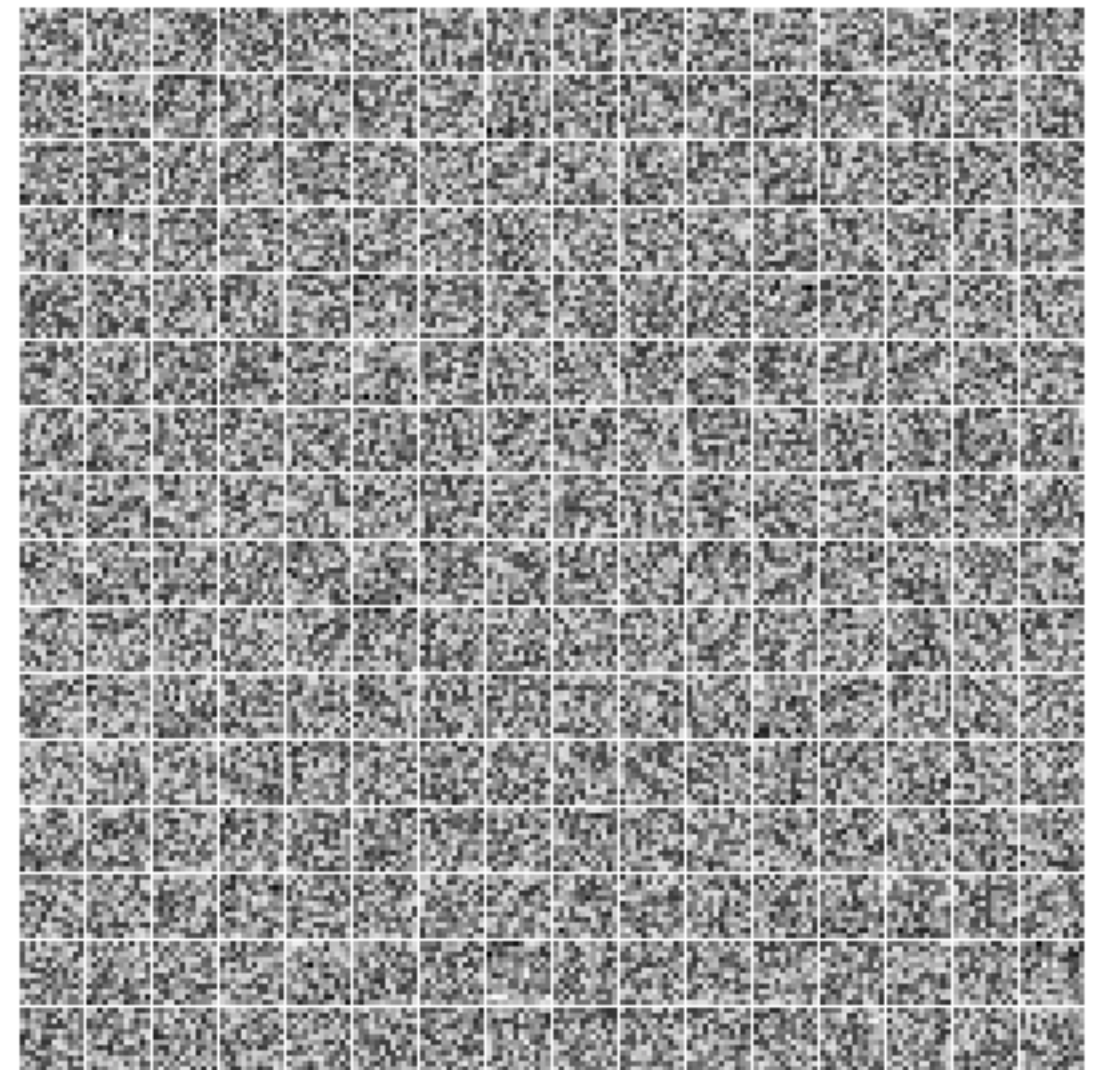
main idea: learn to reconstruct the input with some sparse base functions



unsupervised training: autoencoders



the dataset is just unlabeled data!
of any form: video, frames, etc
easy to get!
easy to consume



iteration no 0

train by finding base-functions
or the basic blocks of images

unsupervised training: autoencoders

related to:
vector quantization
clustering

used for:
compression
de-noising



sparse coding

Sparse coding (Olshausen & Field, 1996). Originally developed to explain early visual processing in the brain (edge detection).

Training: given a set of random patches x , learning a dictionary of bases $[\Phi_1, \Phi_2, \dots]$

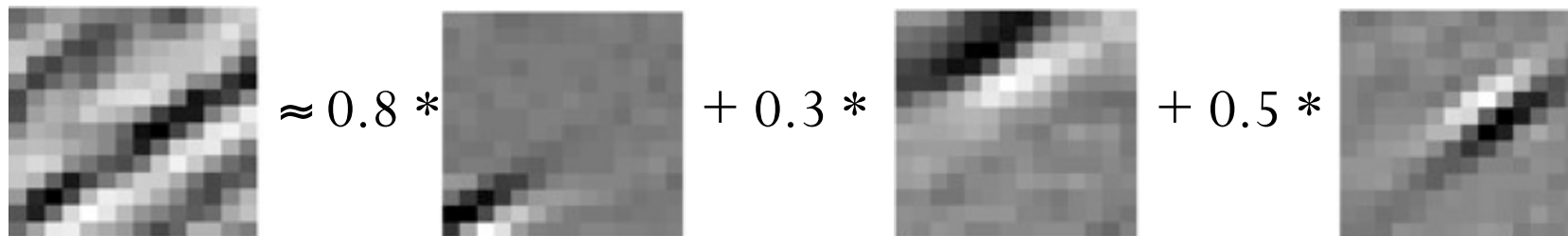
$$\min_{a, \phi} \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$

sparse coding

Input: Novel image patch x (in \mathbb{R}^d) and previously learned ϕ_i 's

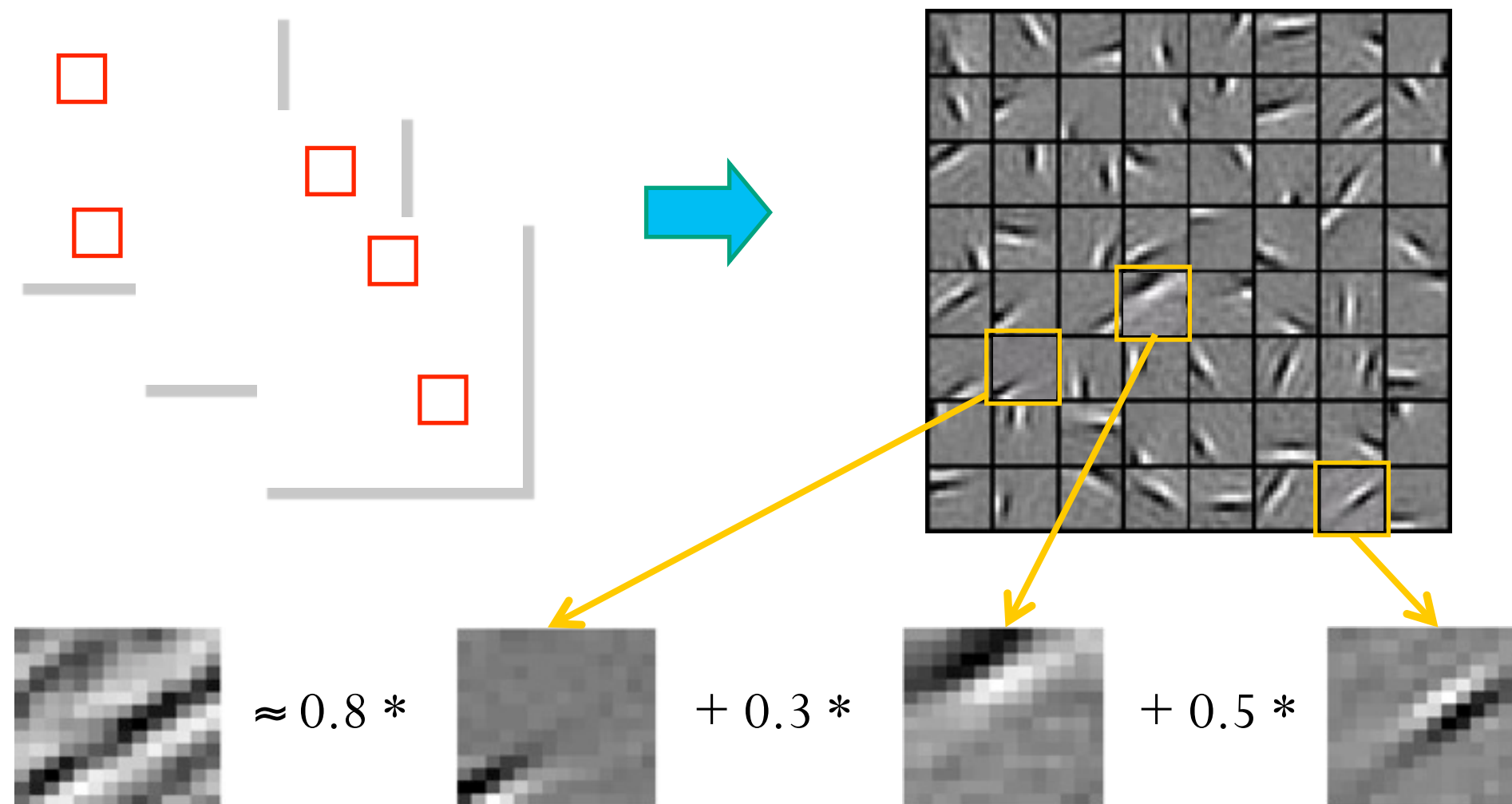
Output: Representation $[a_{i,1}, a_{i,2}, \dots, a_{i,K}]$ of image patch x_i .

$$\min_a \sum_{i=1}^m \left\| x_i - \sum_{j=1}^k a_{i,j} \phi_j \right\|^2 + \lambda \sum_{i=1}^m \sum_{j=1}^k |a_{i,j}|$$

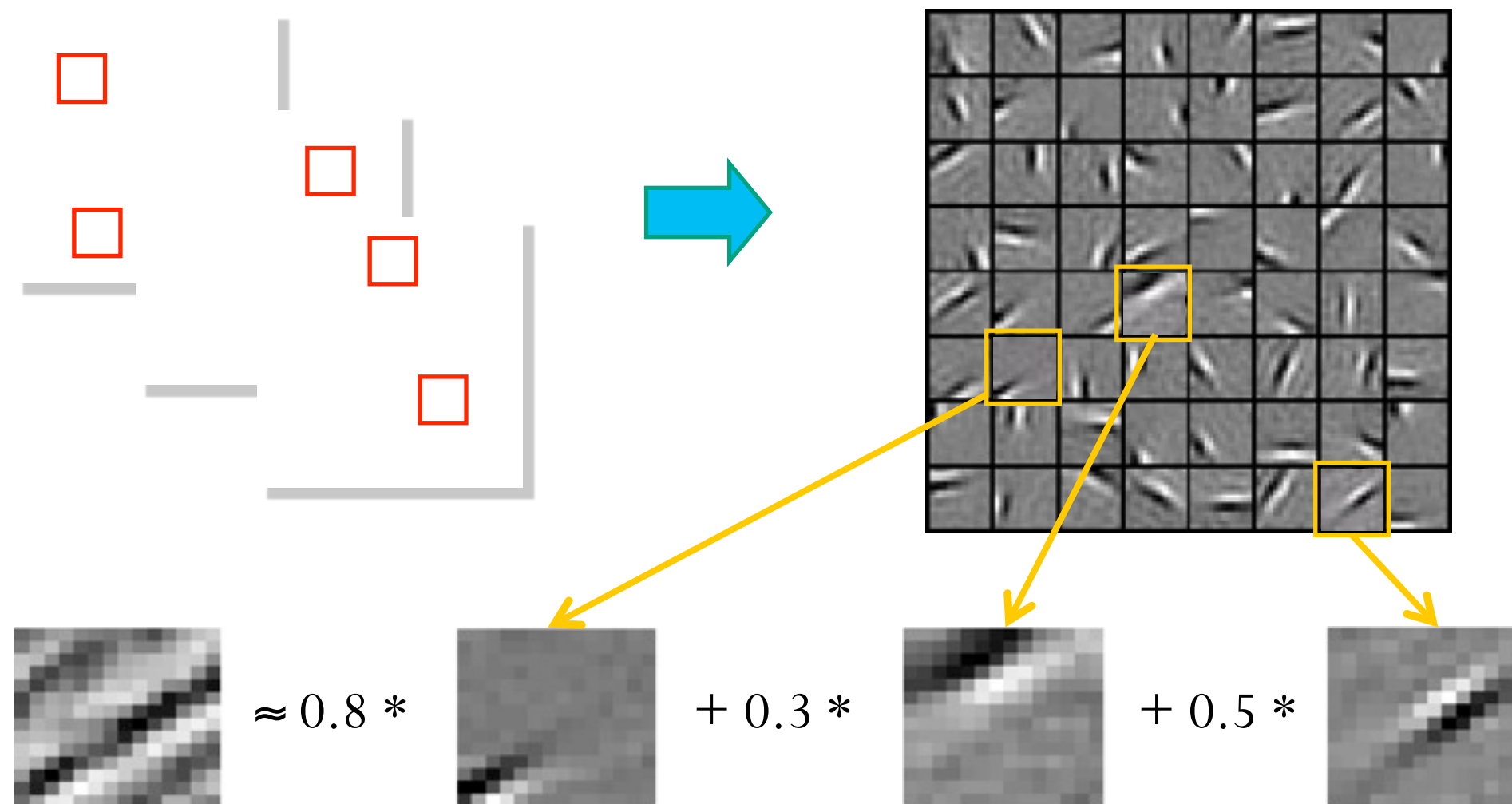


Represent x_i as: $a_i = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, \dots]$

sparse coding



sparse coding



$[a_1, \dots, a_{64}] = [0, 0, \dots, 0, \mathbf{0.8}, 0, \dots, 0, \mathbf{0.3}, 0, \dots, 0, \mathbf{0.5}, 0]$
(feature representation)

Compact & easily interpretable

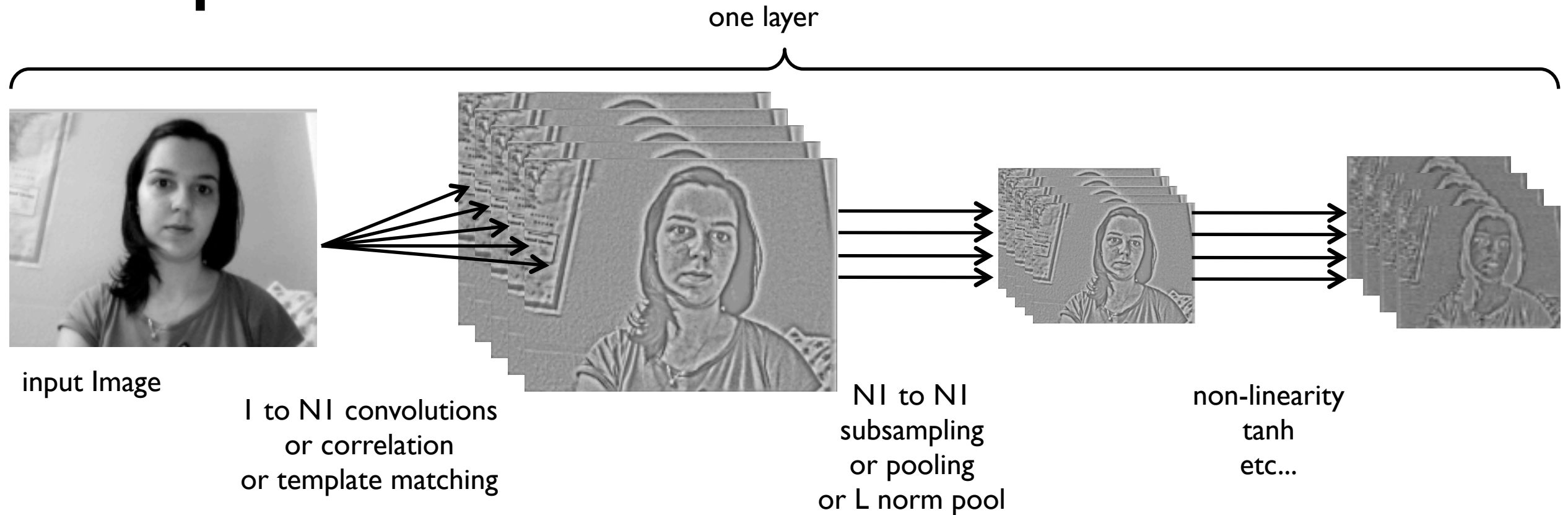
sparse coding: example



64 bases functions of 8x8 pixels

The bases seem to capture the intrinsic structure of the building elements, that are mainly composed of vertical, horizontal, slanting edges and corners.

deep networks



Multiple layers of deep network:

Repeat for each layer:

- 1- sample output of previous layer (new input)
- 2- learn dictionary of inputs = filters
- 3- use filters to generate outputs

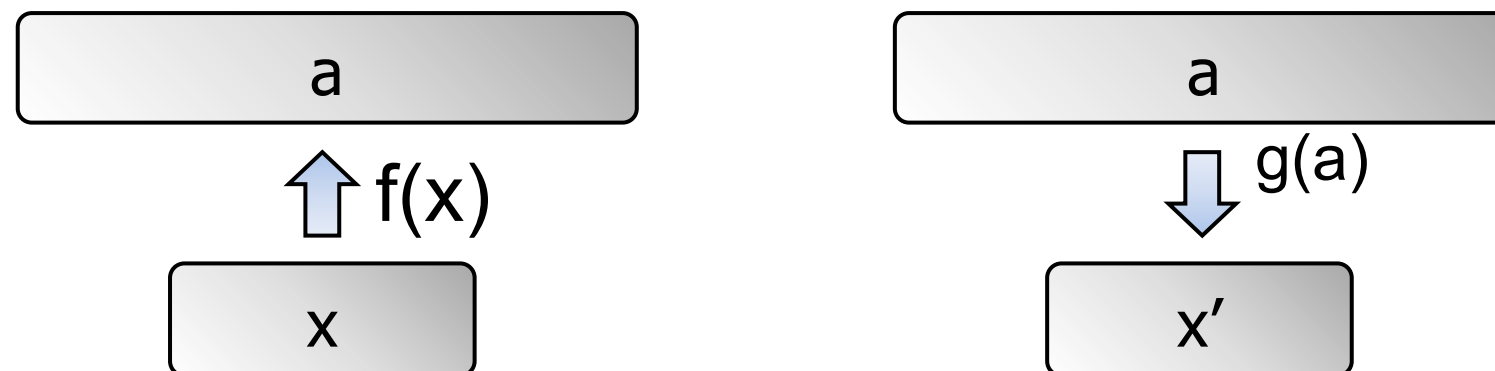
sparse coding: deep networks

Any feature mapping from x to a , i.e. $a = f(x)$, where

- a is sparse (and often higher dim. than x)

- $f(x)$ is nonlinear

- reconstruction $x' = g(a)$, such that $x' \approx x$



Therefore, sparse RBMs, sparse auto-encoder, even VQ can be viewed as a form of sparse coding.