

The background of the slide features a large, semi-transparent watermark of the University of Pisa crest. The crest is circular and contains a central figure, likely a personification of Justice or Liberty, surrounded by Latin text. The watermark is rendered in a dark blue color that blends with the dark blue background.

Deep Learning – Autoencoder Models

INTELLIGENT SYSTEMS FOR PATTERN RECOGNITION (ISPR)

DAVIDE BACCIU – DIPARTIMENTO DI INFORMATICA - UNIVERSITA' DI PISA

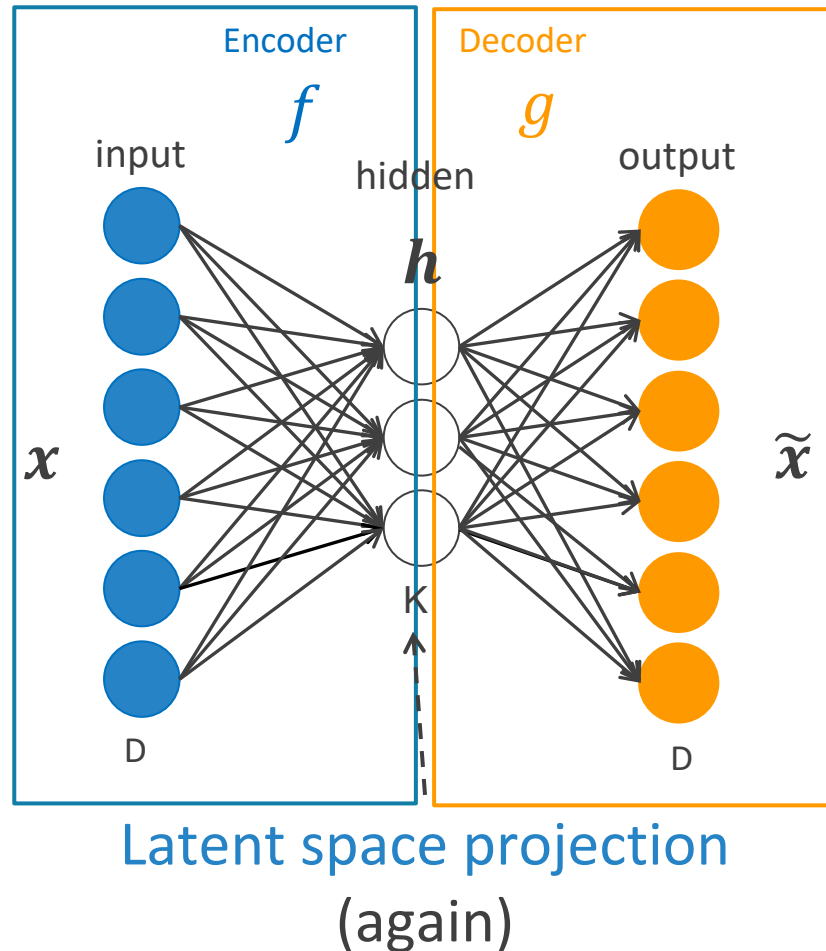
DAVIDE.BACCIU@UNIFI.IT

Lecture Outline

Autoencoders a.k.a. The first and the latest deep learning model

- Autoencoders and dimensionality reduction
- Deep **neural** autoencoders
 - Sparse
 - Denoising
 - Contractive
- Deep **generative-based** autoencoders
 - Deep Belief Networks
 - Deep Boltzmann Machines
- Application Examples

Basic Autoencoder (AE)



- Train a model to **reconstruct the input**
- Passing through some form of **information bottleneck**
 - $K \ll D$, or?
 - h sparsely active
- Train by loss minimization
$$L(\mathbf{x}, \tilde{\mathbf{x}}) = L(\mathbf{x}, g(f(\mathbf{x})))$$

Neural Autoencoders

Generally, we would like to train nonlinear AEs, with possibly $K > D$, that do not learn trivial identity

- Regularized autoencoders
 - Sparse AE
 - Denoising AE
 - Contractive AE
- Autoencoders with **dropout** layers

Sparse Autoencoder

Add a term to the cost function to penalize \mathbf{h} (want the number of active units to be small)

$$J_{SAE}(\theta) = \sum_{\mathbf{x} \in S} (L(\mathbf{x}, \tilde{\mathbf{x}}) + \lambda \Omega(\mathbf{h}))$$

Typically

$$\Omega(\mathbf{h}) = \Omega(f(\mathbf{x})) = \sum_j |h_j(\mathbf{x})|$$

Probabilistic Interpretation (Oh No, Again!)

Training with regularization is MAP inference

$$\max \log P(\mathbf{h}, \mathbf{x}) = \max (\log P(\mathbf{x}|\mathbf{h}) + \log P(\mathbf{h}))$$

Likelihood

$$P(\mathbf{h}) = \frac{\lambda}{2} \exp\left(-\frac{\lambda}{2} \|\mathbf{h}\|_1\right) \xrightarrow{\text{Prior}} \Omega(\mathbf{h}) = \lambda \|\mathbf{h}\|_1$$

Laplace



UNIVERSITÀ DI PISA

Denoising Autoencoder (DAE)

Train the AE to minimize the function

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

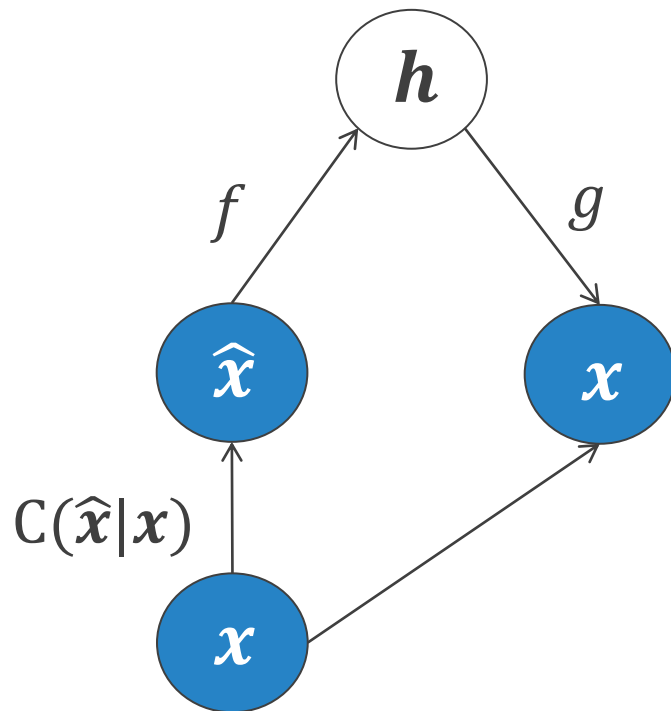
where $\hat{\mathbf{x}}$ is a version of original input \mathbf{x} corrupted by some noise process $C(\hat{\mathbf{x}}|\mathbf{x})$

Key Intuition - Learned representations should be robust to partial destruction of the input



Another Interpretation...

...yes, exactly the one you are thinking of



Learns the **denoising distribution**

$$P(x|\tilde{x})$$

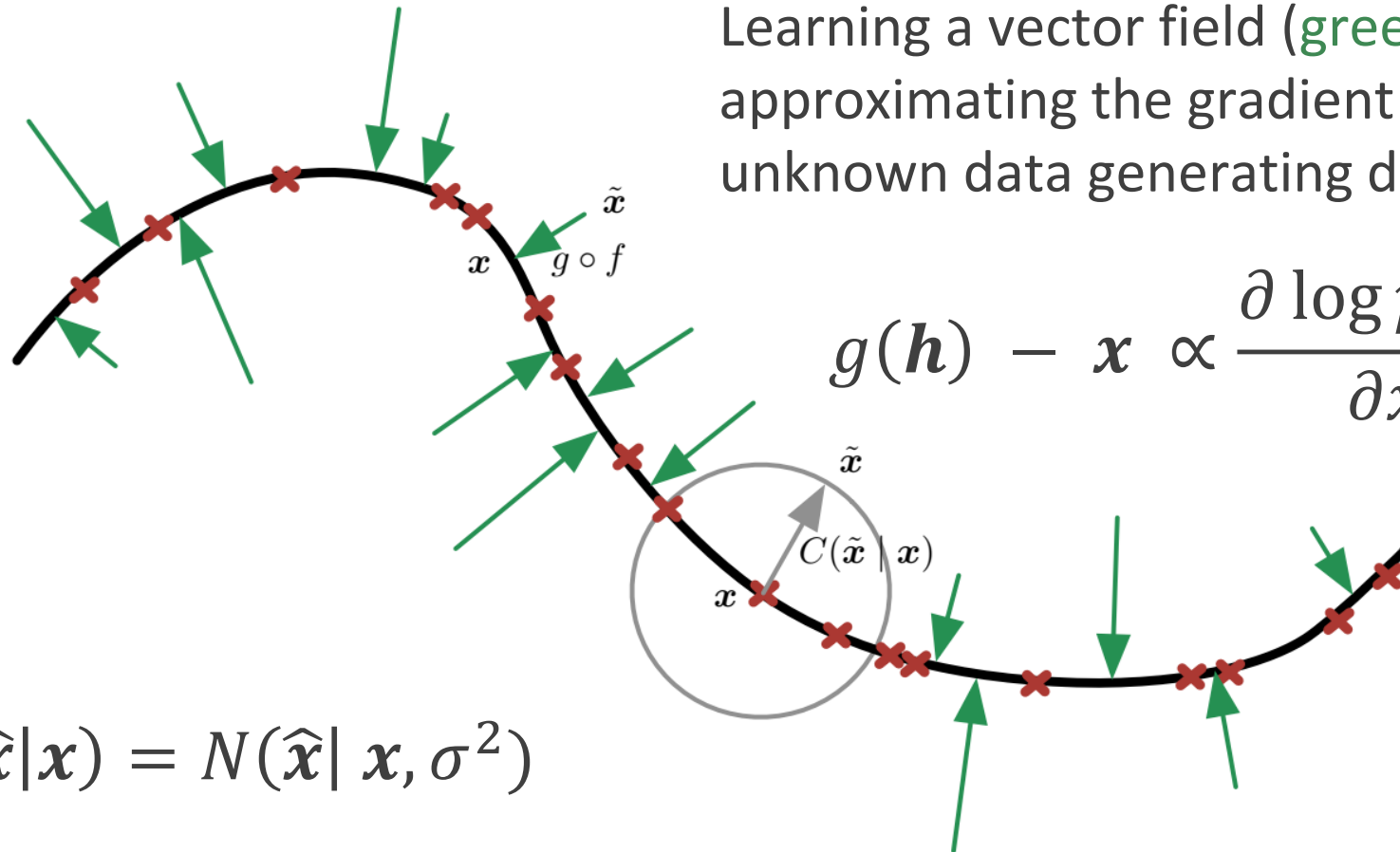
By minimizing

$$-\log P_d(x|h = f(\tilde{x}))$$



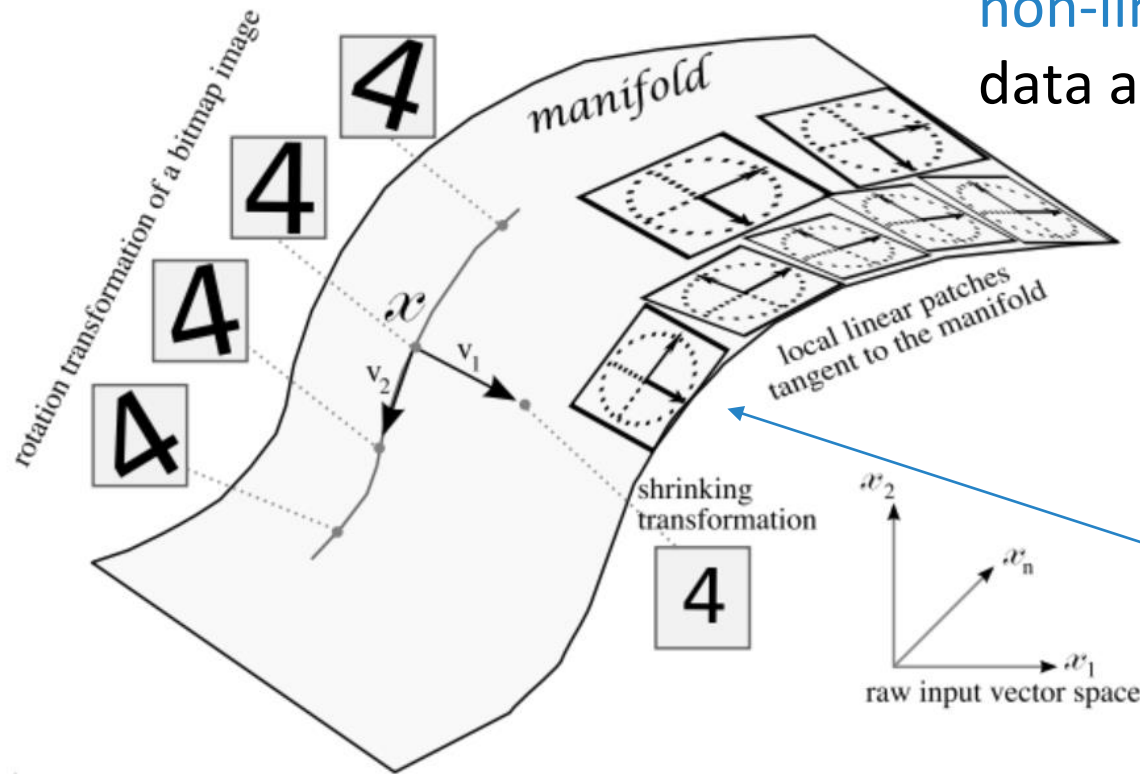
DAE as Manifold Learning

Learning a vector field (green arrows) approximating the gradient of the unknown data generating distribution



$$C(\hat{\mathbf{x}} | \mathbf{x}) = N(\hat{\mathbf{x}} | \mathbf{x}, \sigma^2)$$

The Manifold Assumption



Assume data lies on a lower dimensional **non-linear manifold** since variables in data are typically dependent

Regularized AE can afford to represent **only variations that are needed to reconstruct training examples**

AE mapping is sensitive only to **changes in manifold direction**

Yoshua Bengio, Learning Deep Architectures for AI, Foundations and Trends in Machine Learning, 2009.



UNIVERSITÀ DI PISA

Contractive Autoencoder

Penalize **encoding function** for input sensitivity

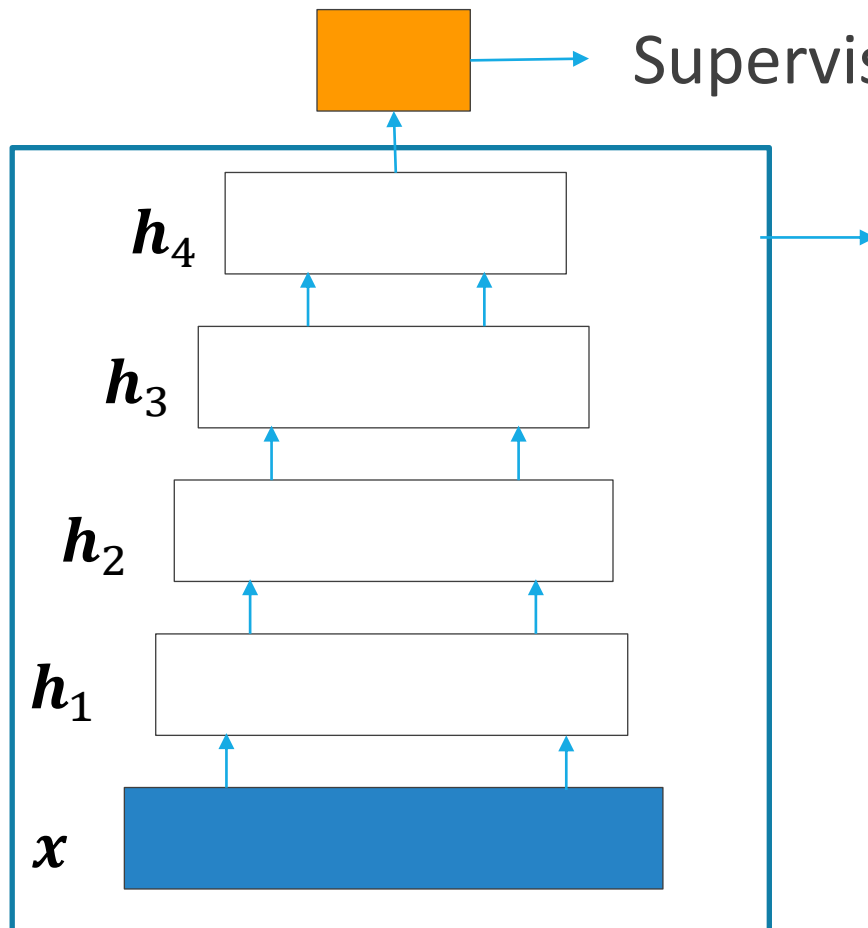
$$J_{CAE}(\theta) = \sum_{\mathbf{x} \in \mathcal{S}} (L(\mathbf{x}, \tilde{\mathbf{x}}) + \lambda \Omega(\mathbf{h}))$$

$$\Omega(\mathbf{h}) = \Omega(f(\mathbf{x})) = \left\| \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} \right\|_F$$

You can as well **penalize on higher order derivatives**



Deep Autoencoder (AE)



Supervised learning

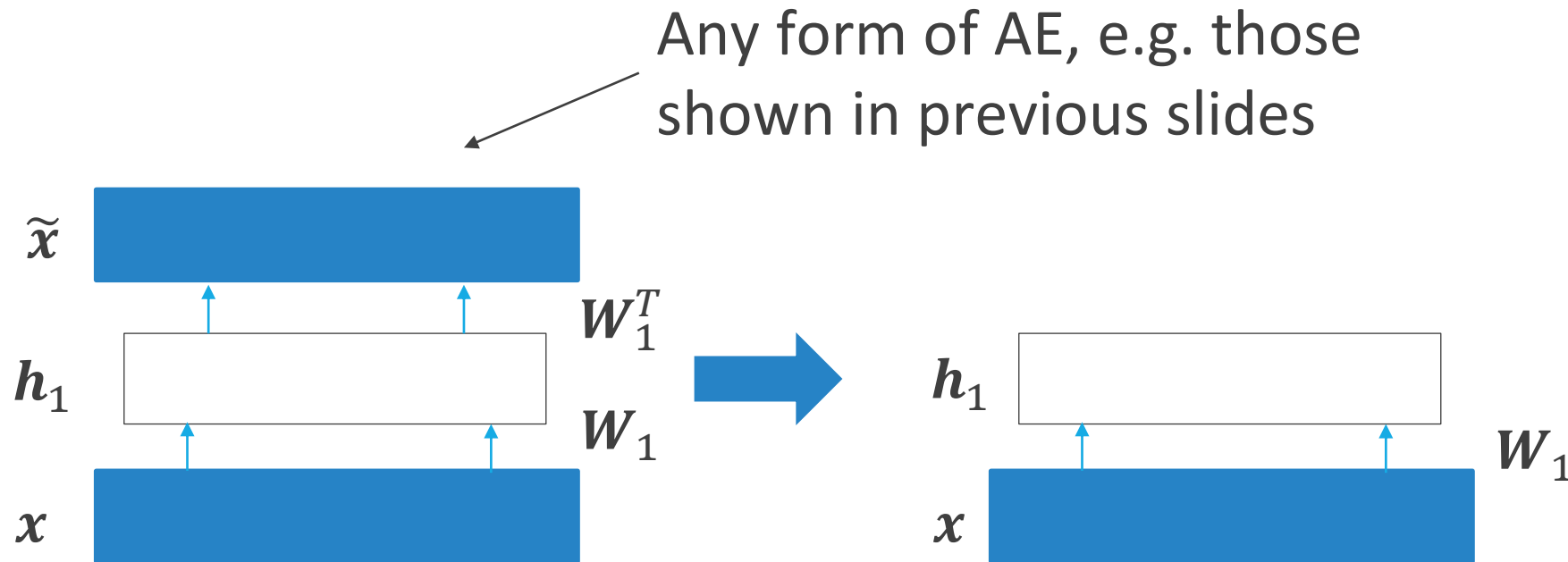
- Unsupervised training
- Hierarchical autoencoder
- Extracts a **representation of inputs** that facilitates
 - Data **visualization**, exploration, indexing,...
 - Realization of a **supervised task**



UNIVERSITÀ DI PISA

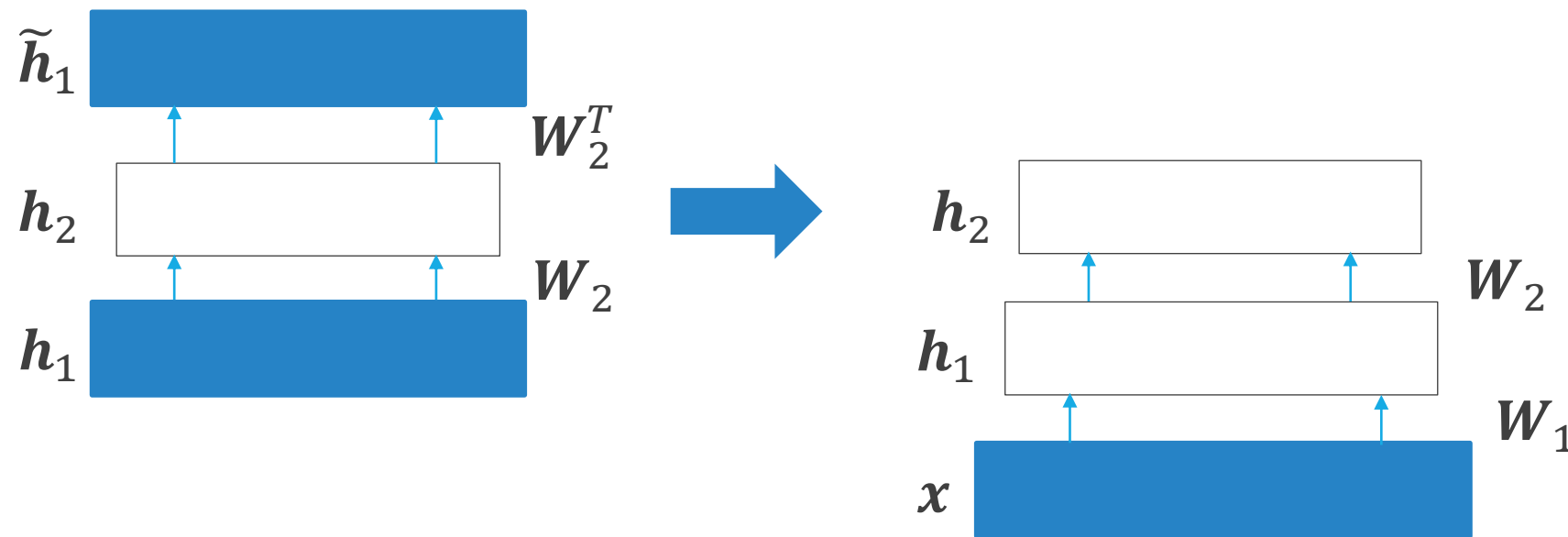
Unsupervised Layerwise Pretraining

Incremental unsupervised construction of the Deep AE



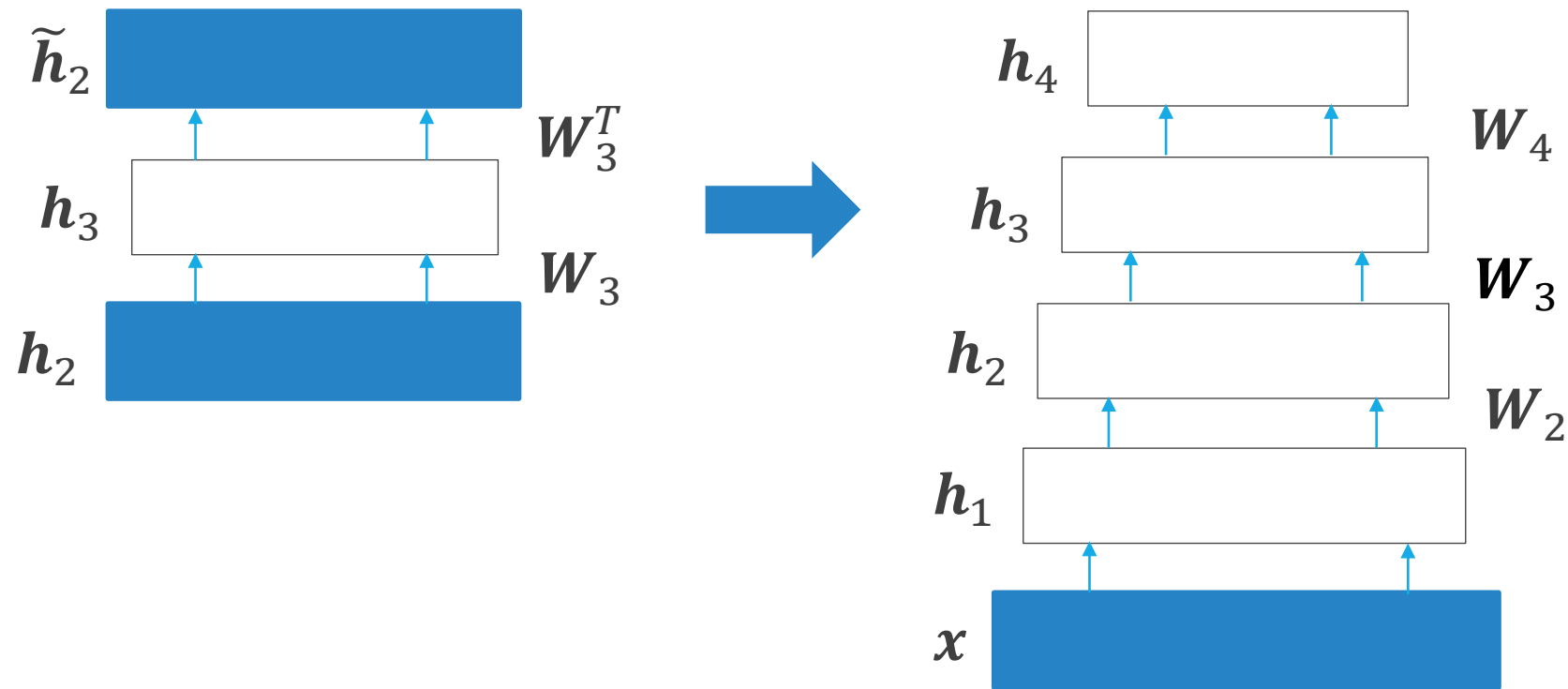
Unsupervised Layerwise Pretraining

Incremental unsupervised construction of the Deep AE



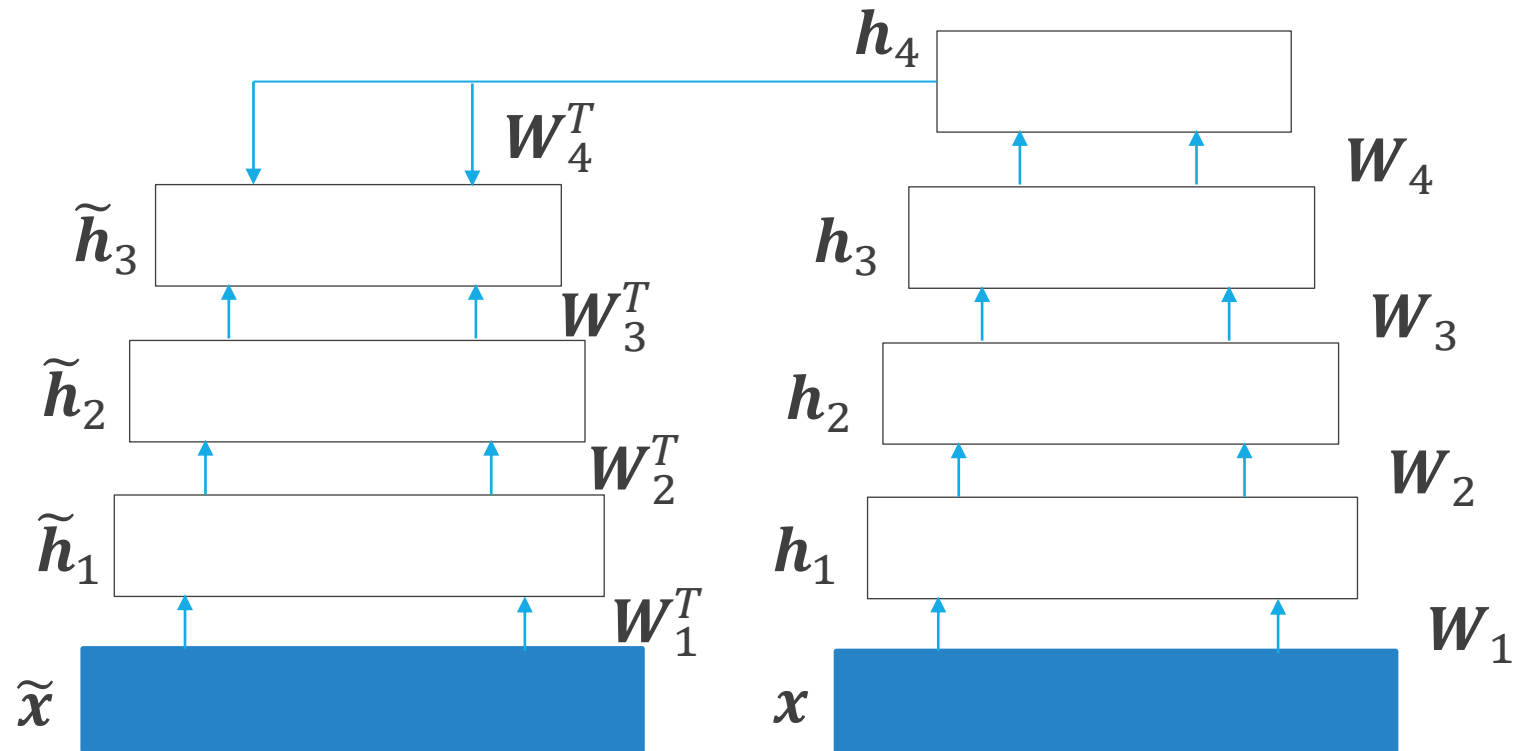
Unsupervised Layerwise Pretraining

Incremental unsupervised construction of the Deep AE



Optional Fine Tuning

Fine tune the whole autoencoder to optimize input reconstruction
You can use backpropagation, but it remains an unsupervised task

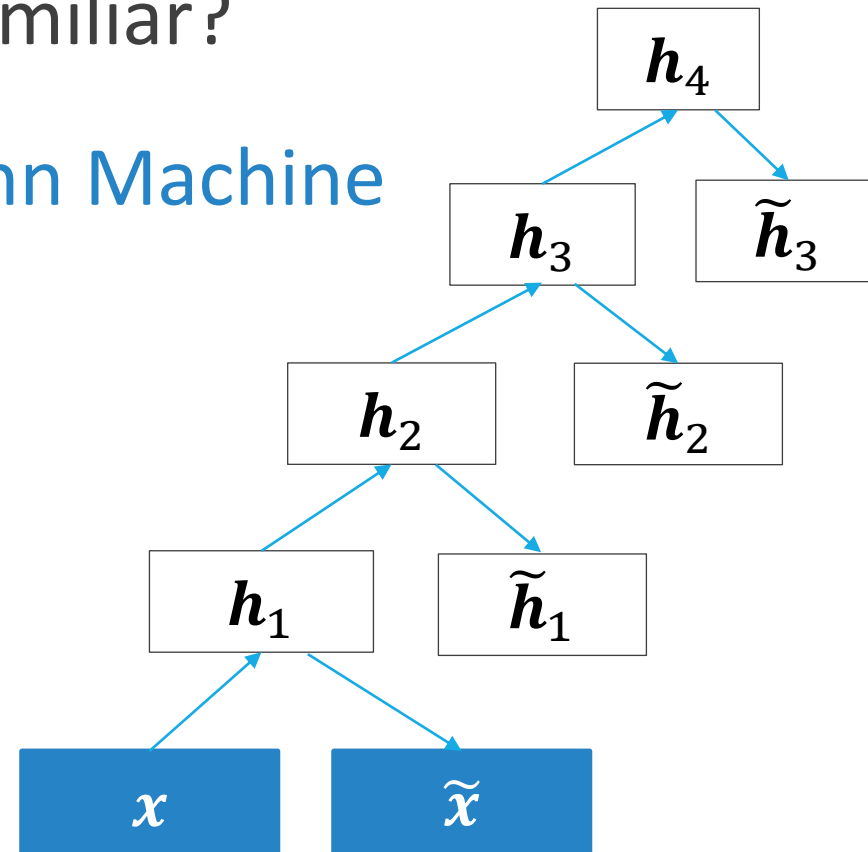


Rearranging the Graphics

Does it look like something familiar?

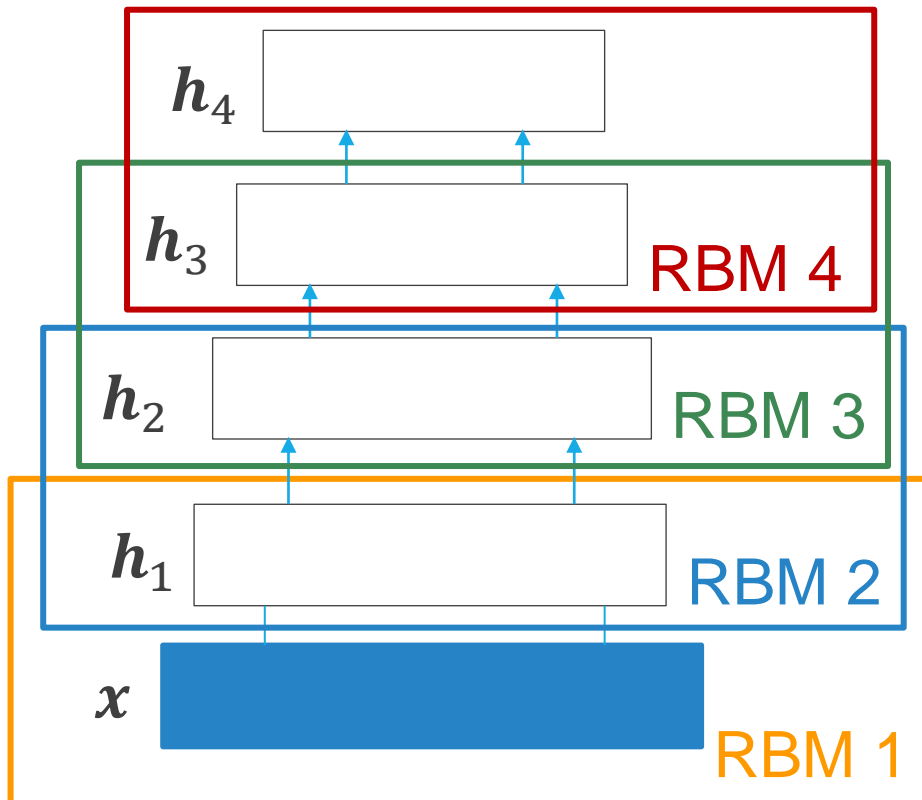
A layered **Restricted Boltzmann Machine**

Can use RBM to perform **layerwise pretraining** and learn the matrices W_i



Deep Belief Network (DBN)

A stack of pairwise RBM



IMPORTANT NOTE

A DBM is a deep autoencoder
but it is **NOT** a deep RBM

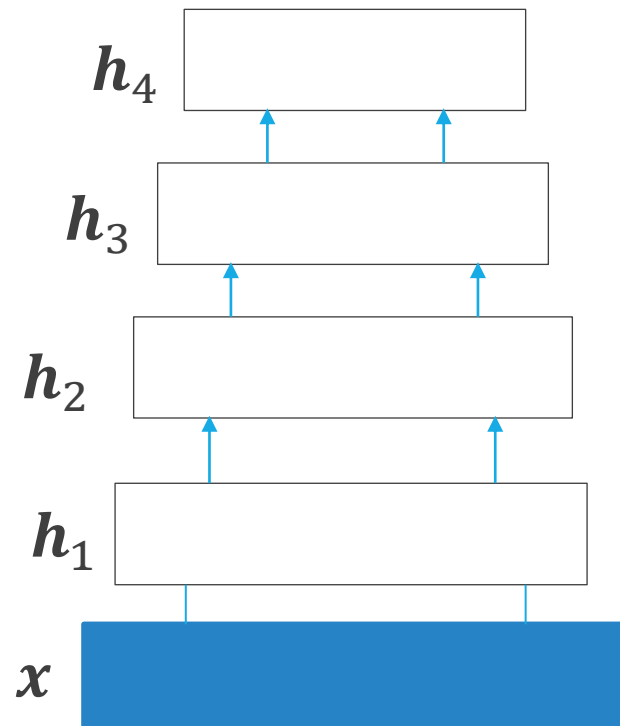
It is (mostly) **directed!**



UNIVERSITÀ DI PISA

Deep Boltzmann Machine (DBM)

How do we get this?



Training requires some attention because of the **recurrent interactions from higher layers** to the bottom

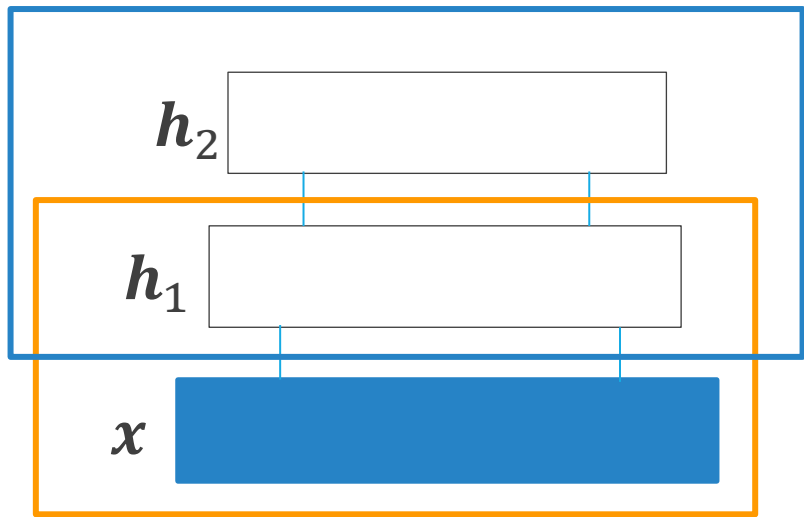
$$P(h_j^1 | \mathbf{x}, \mathbf{h}^2) = \sigma \left(\sum_i W_{ij}^1 x_i + \sum_m W_{jm}^2 h_m^2 \right)$$

$$P(x_i | \mathbf{h}^1) = \sigma \left(\sum_j W_{ij}^1 h_j^1 \right)$$



Pretraining DBM

How do we get this?



1) (Pre)training the first layer entails fitting this model

2) (Pre)training the second layer **changes** h^1 prior by

$$P(h^1 | W^2) = \sum_{h^2} P(h^1, h^2 | W^2)$$

When putting things together, we need to **average** between the two

$$P(h^1 | W^1) = \sum_x P(h^1, x | W^1)$$

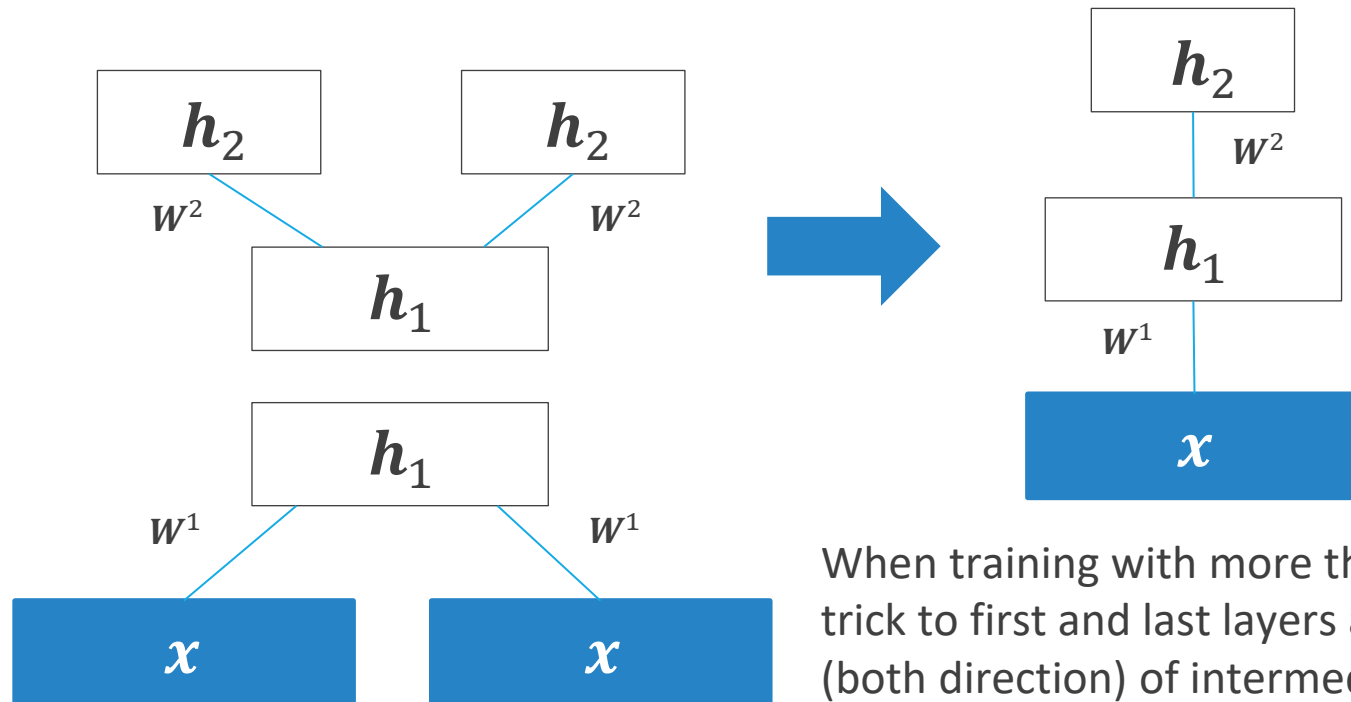
$$\rightarrow P(x | \theta) = \sum_{h^1} P(h^1 | W^1) P(x | h^1, W^1)$$



Pretraining DBM - Trick

Averaging the two models of h^1 can be approximated by taking half contribution from W^1 and half from W^2

- Using full W^1 and W^2 would double count x contribution as h^2 depends on x



When training with more than two RBMs apply trick to first and last layers and halve weights (both direction) of intermediate RBM

Software - Deep Neural Autoencoders

- All deep learning frameworks offer facilities to build (deep) AEs
- Check out classic Theano-based tutorials for [denoising autoencoders](#) and their [stacked version](#)
- A variety of deep AE in [Keras](#) and their counterpart in [Torch](#) (plus a selection in [Pytorch](#))
- Stacked autoencoders built with [official Matlab](#) toolbox functions

Matlab - Deep Generative Models

- [Matlab code](#) for the DBN paper with a demo on MNIST data
- [Matlab code](#) for Deep Boltzmann Machines with a demo on MNIST data
- [Deepmat](#) – Matlab library for deep generative models
- [DeeBNet](#) – Matlab/Octave toolbox for deep generative models with GPU support



Python - Deep Generative Models

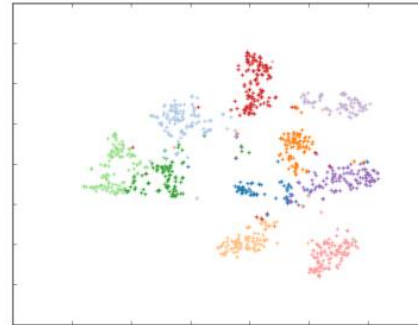
- DBN and DBM implementations exist for [all major deep learning libraries](#)
- [Deep Boltzmann machine implementation](#) (Tensorflow-based) with image processing application, pre-trained networks and notebooks
- [Deepnet](#) – A Toronto based implementation of deep autoencoders (neural and generative)
- Check out classic Theano-based tutorials for [deep belief networks](#) and [RBM](#)



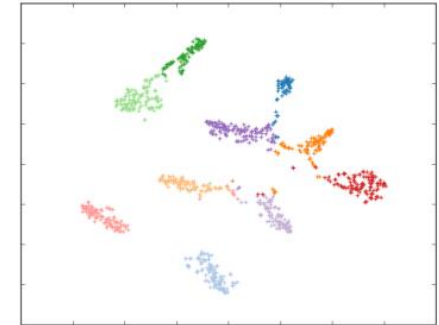
AE Applications - Visualization



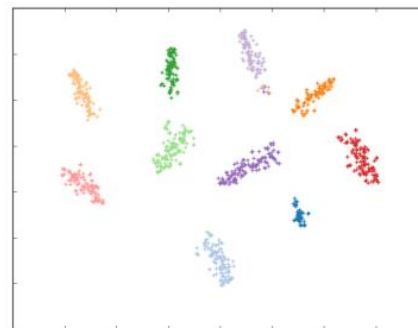
Visualizing complex data in learned latent space



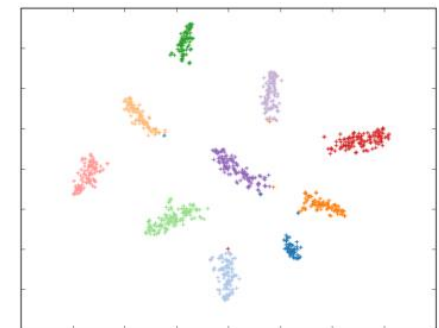
(a) Epoch 0



(b) Epoch 3



(d) Epoch 9



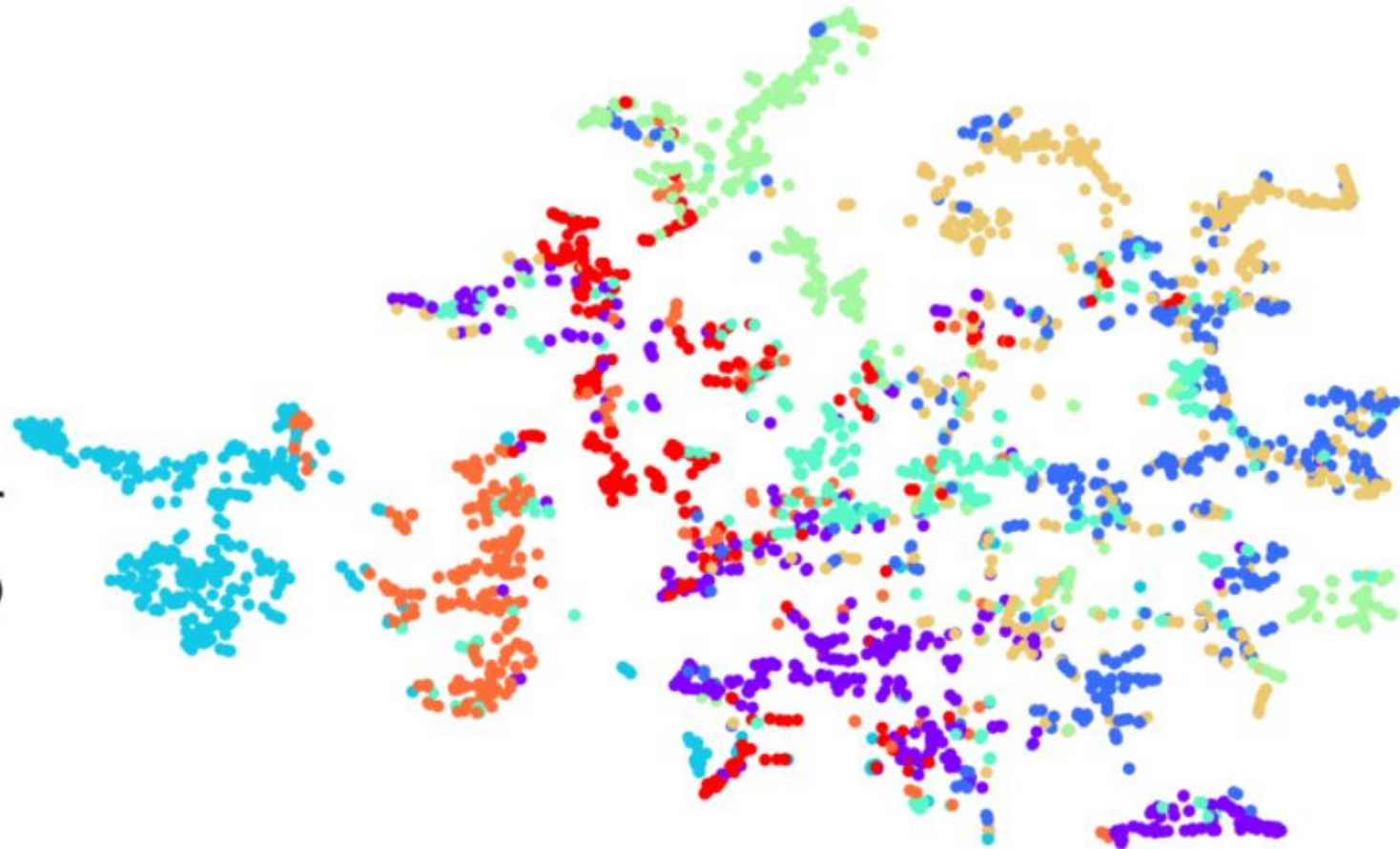
(e) Epoch 12



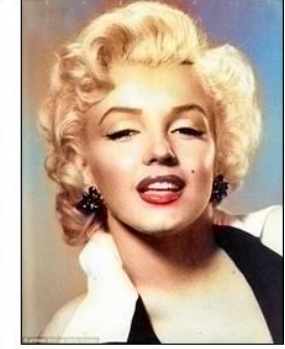
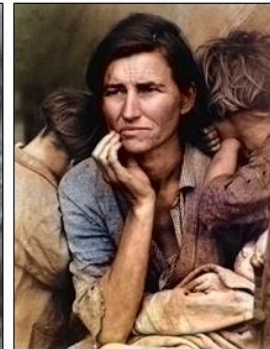
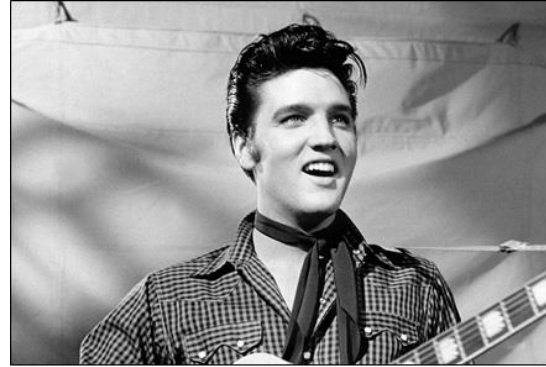
UNIVERSITÀ DI PISA

Visualizing Sound

- laughter
- rustle
- guitar
- cat
- helicopter
- water_tap
- child
- speech



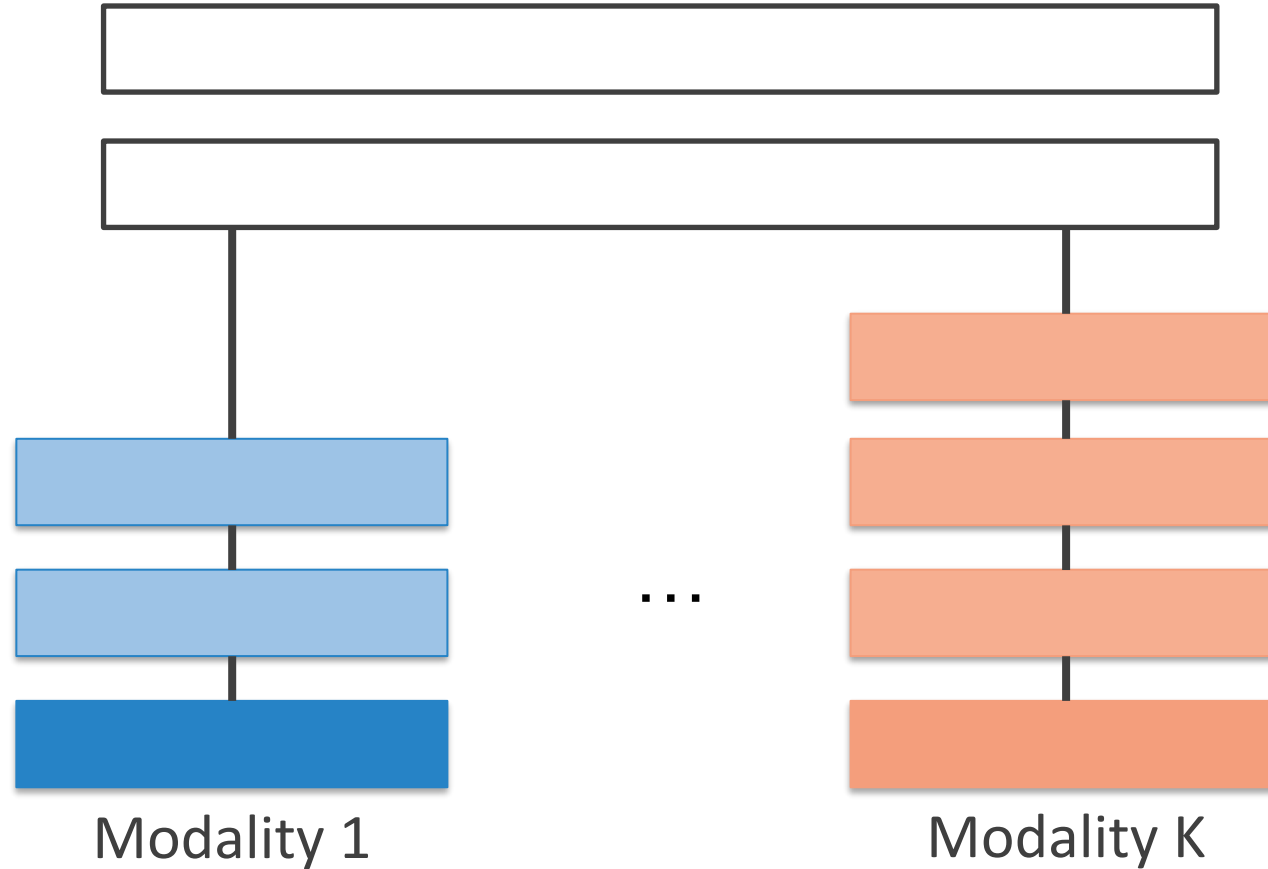
AE Applications – Image Restoration/Colorization



Apply autoencoder construction with advanced building blocks (e.g. CNN layers)

Multimodal DBM

Modality
fusion
layers















N. Srivastava, R. Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines, JMLR 2014



UNIVERSITÀ DI PISA











Multimodal DBM – Image and Text

	Image	Given Tags	Generated Tags	Input Tags	Nearest neighbors to generated image features	
$P(\text{txt} \text{img})$		pentax, k10d, kangarooisland, southaustralia, sa, 300mm, australia, australiansealion	beach, sea, surf, strand, shore, wave, seascape, sand, ocean, waves	nature, hill, scenery, green, clouds	 	$P(\text{img} \text{txt})$
		< no text >	night, lights, christmas, nightshot, nacht, nuit, notte, longexposure, noche, nocturna	flower, nature, green, flowers, petal, petals, bud	 	
		aheram, 0505, sarahc, moo	portrait, bw, balckandwhite, people, faces, girl, blackwhite, person, man	blue, red, art, artwork, painted, paint, artistic, surreal, gallery, bleu	 	
		unseulpixel, naturey crap	fall, autumn, trees, leaves, foliage, forest, woods, branches, path	bw, blackandwhite, noiretblanc, bianconero, blancoynegro	 	

N. Srivastava, R. Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines, JMLR 2014

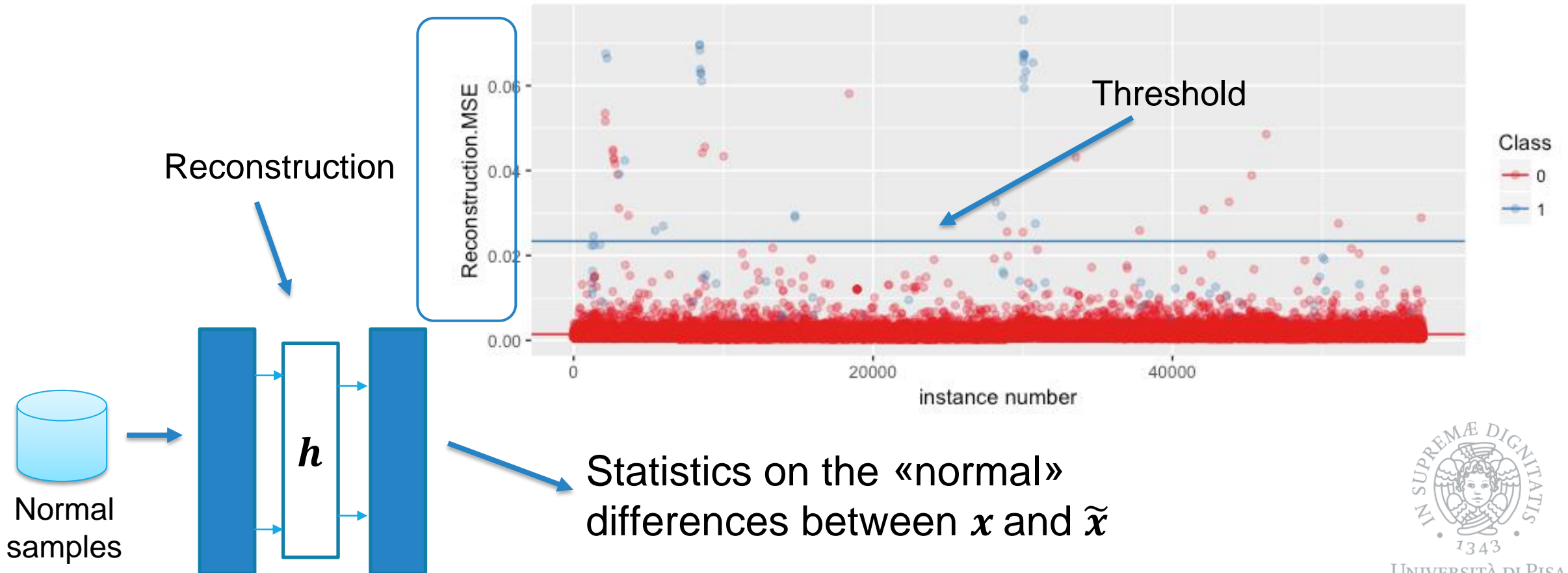


Multimodal DBM – Multimodal Quering

Multimodal Query	Top 4 retrieved results				
 <p data-bbox="397 733 659 891">hongkong, causewaybay, shoppingcentre, building, mall</p>	 <p data-bbox="749 733 1090 891">howell, bridge, genesee, river, rochester, downtown, building</p>	 <p data-bbox="1141 733 1480 891">london, uk, night, skyline, river, thames, lights, bridge</p>	 <p data-bbox="1538 758 1786 868">edinburgh, scotland, dusk, bank</p>	 <p data-bbox="1849 758 2181 868">arcoiris, fincadehierro, lluvia, sannicolos, valencia</p>	
 <p data-bbox="387 1143 670 1219">me, myself, eyes, blue, hair</p>	 <p data-bbox="754 1143 1085 1219">urban, me, abigfave, fiveflickrfav,</p>	 <p data-bbox="1144 1125 1475 1235">trisha, mynewcamera, lake, field, girl</p>	 <p data-bbox="1544 1143 1781 1219">me, ofme, self, selfportrait</p>	 <p data-bbox="1862 1143 2168 1219">pink, prettyinpink, explored</p>	

N. Srivastava, R. Salakhutdinov, Multimodal Learning with Deep Boltzmann Machines, JMLR 2014

Anomaly Detection



Take Home Messages

- Regularized autoencoder
 - Optimize reconstruction quality
 - Constrain stored information
- Autoencoder training is **manifold learning**
 - Learn a latent space manifold where input data resides
 - Store only **variations that are useful** to represent training data
- Autoencoders **learn a (conditional) distribution** of input data $P(\hat{\mathbf{x}} | \dots)$
- Deep AE: pretraining, fine tuning, supervised optimization
- Use AE for finding new/useful **data representations**
 - Or to learn its distribution

Next Lecture

Gated Recurrent Networks

- Learning with sequential data
- Gradient issues
- Gated RNN
 - Long-Short Term Memories (LSTM)
 - Gated Recurrent Units (GRU)
- Advanced topics
 - Understanding and exploiting memory encoding
 - Applications

PART I (tomorrow)

PART II (friday)

