

# Deep Learning for Vision

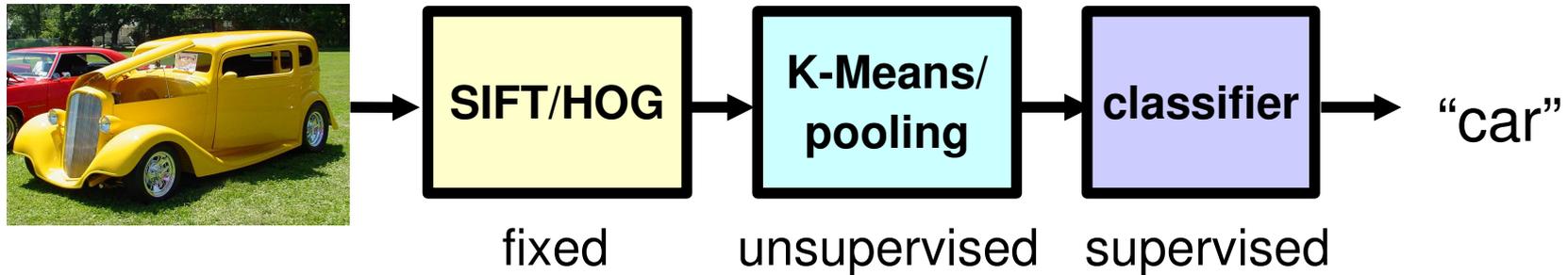
Marc'Aurelio Ranzato



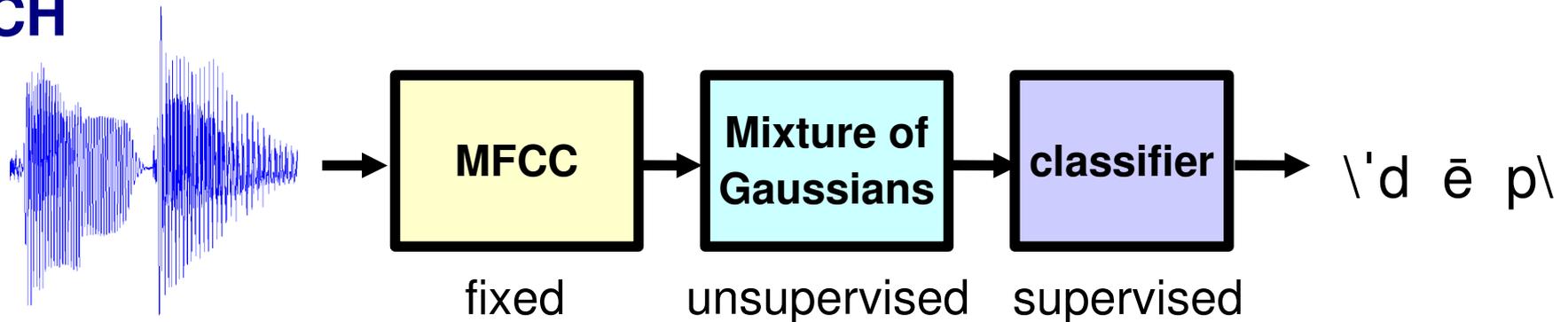
Facebook, AI Group

# Traditional Pattern Recognition

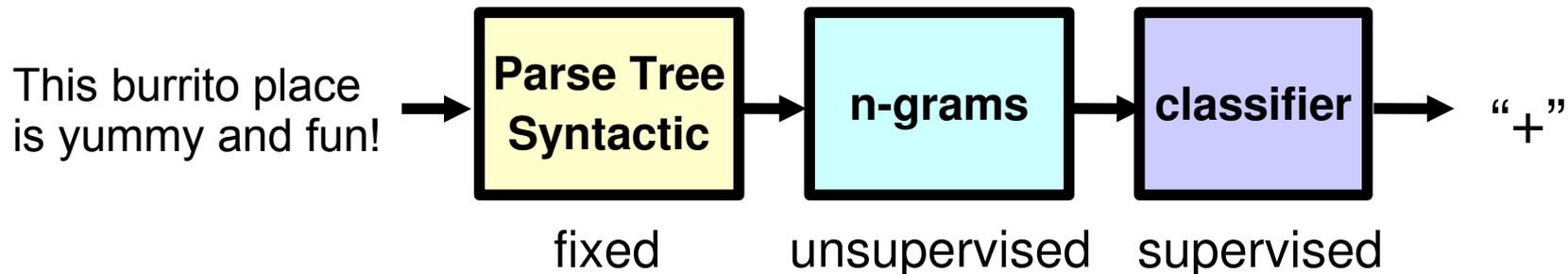
## VISION



## SPEECH



## NLP



# Hierarchical Compositionality (DEEP)

## VISION

pixels → edge → texton → motif → part → object

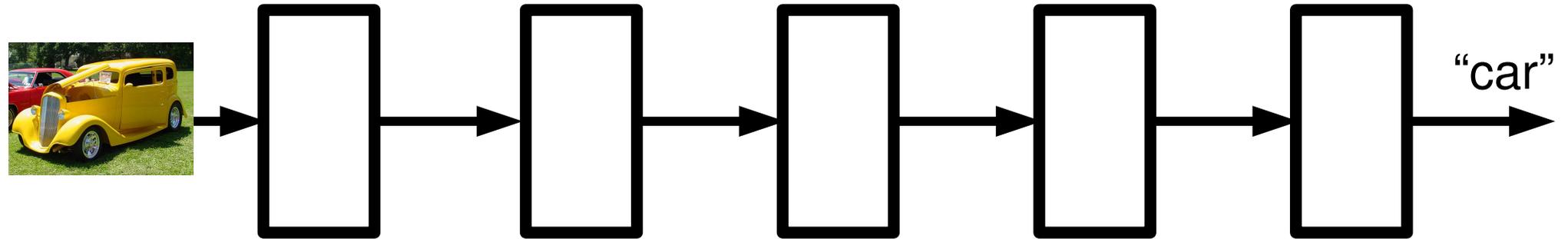
## SPEECH

sample → spectral  
band → formant → motif → phone → word

## NLP

character → word → NP/VP/.. → clause → sentence → story

# Deep Learning



## What is Deep Learning

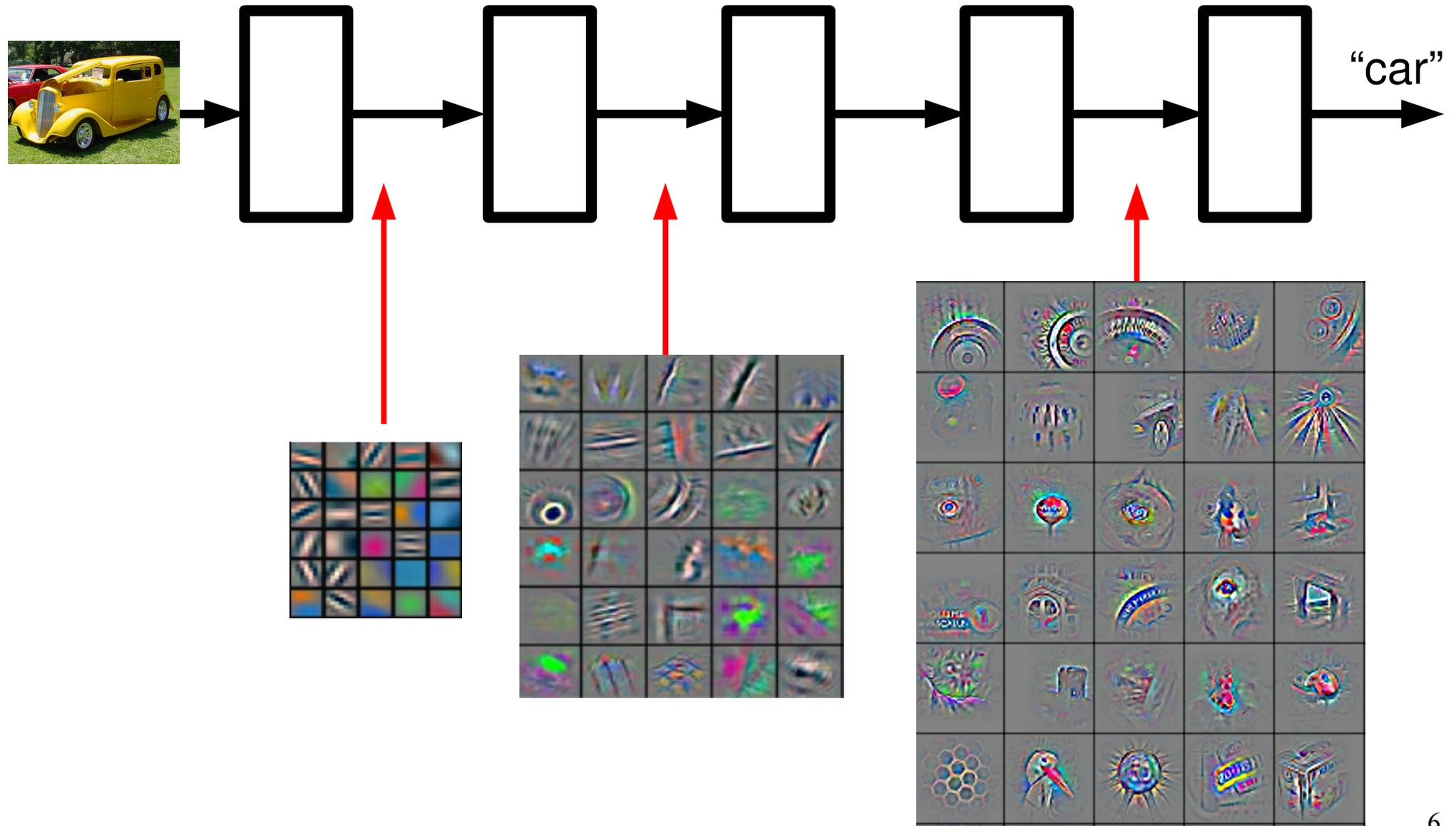
- Cascade of non-linear transformations
- End to end learning
- General framework (any hierarchical model is deep)

# Deep Learning VS Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels → edge → texton → motif → part → object

# Deep Learning



# Deep Learning VS Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels → edge → texture → motif → part → object

- It is more efficient.

E.g.: Checking N-bit parity requires N-1 gates laid out on a tree of depth  $\log(N-1)$ . The same would require  $O(\exp(N))$  with a two layer architecture.

$$p = \sum_i \alpha_i f_i(x) \quad \text{VS} \quad p = \alpha_n f_n(\alpha_{n-1} f_{n-1}(\dots \alpha_1 f_1(x) \dots))$$

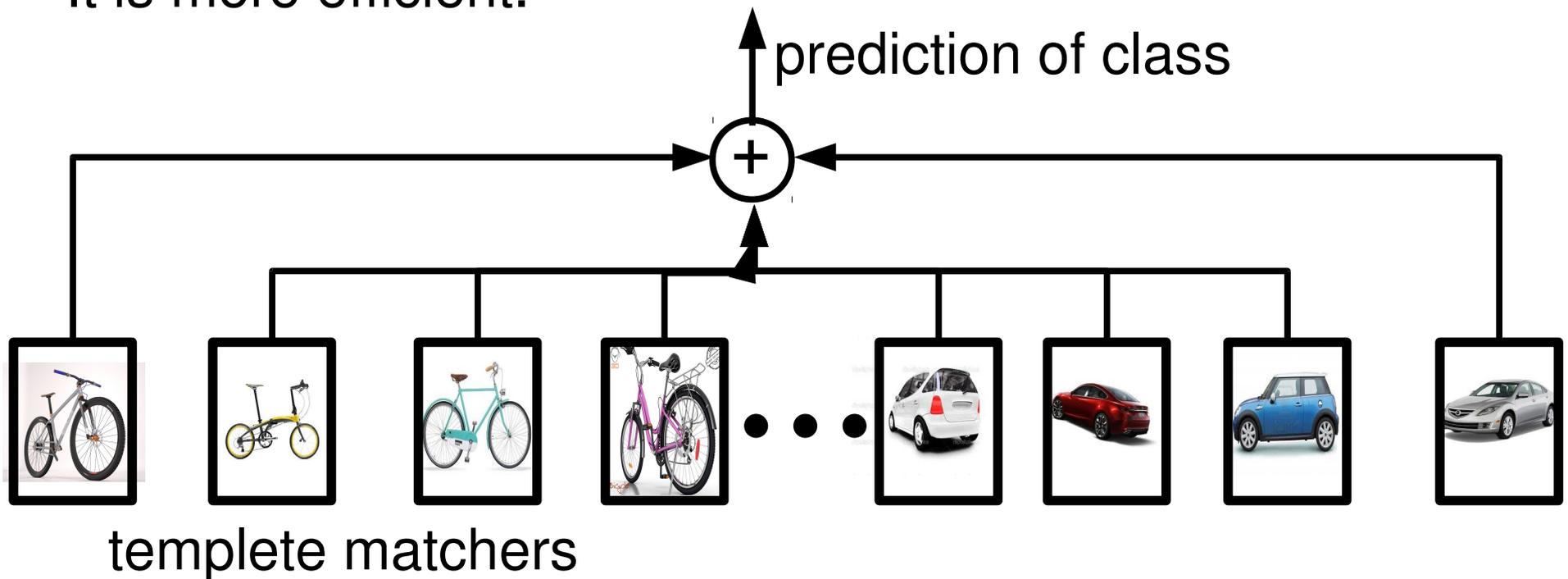
Shallow learner is often inefficient: it requires exponential number of templates (basis functions).

# Deep Learning VS Shallow Learning

- Structure of the system naturally matches the problem which is inherently hierarchical.

pixels  $\rightarrow$  edge  $\rightarrow$  texture  $\rightarrow$  motif  $\rightarrow$  part  $\rightarrow$  object

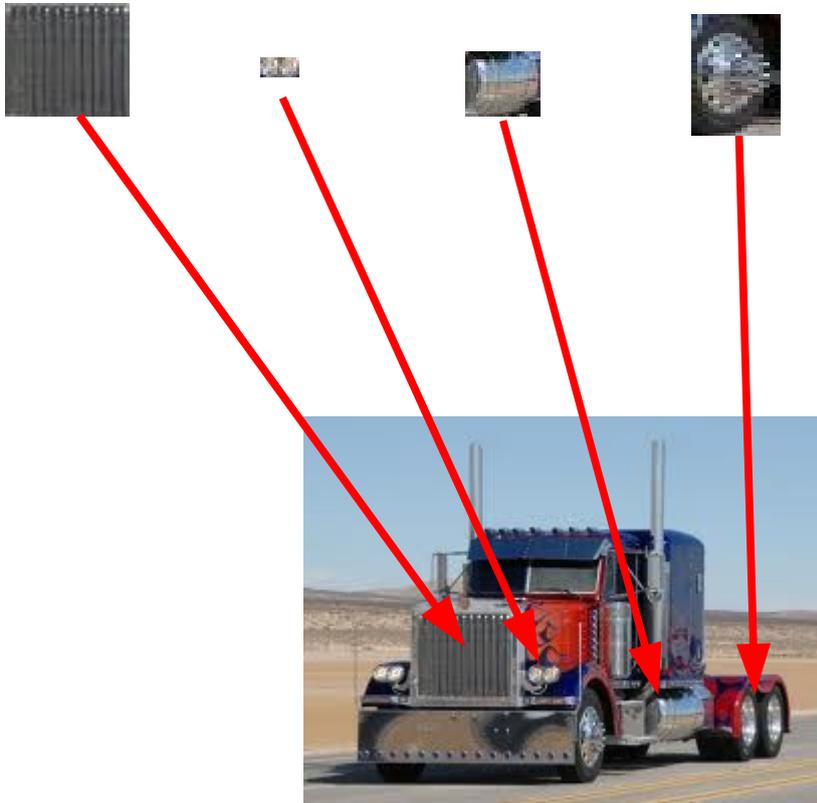
- It is more efficient.



Shallow learner is inefficient.

# Composition: distributed representations

[0 0 **1** 0 0 0 0 **1** 0 0 **1 1** 0 0 **1** 0 ... ] truck feature

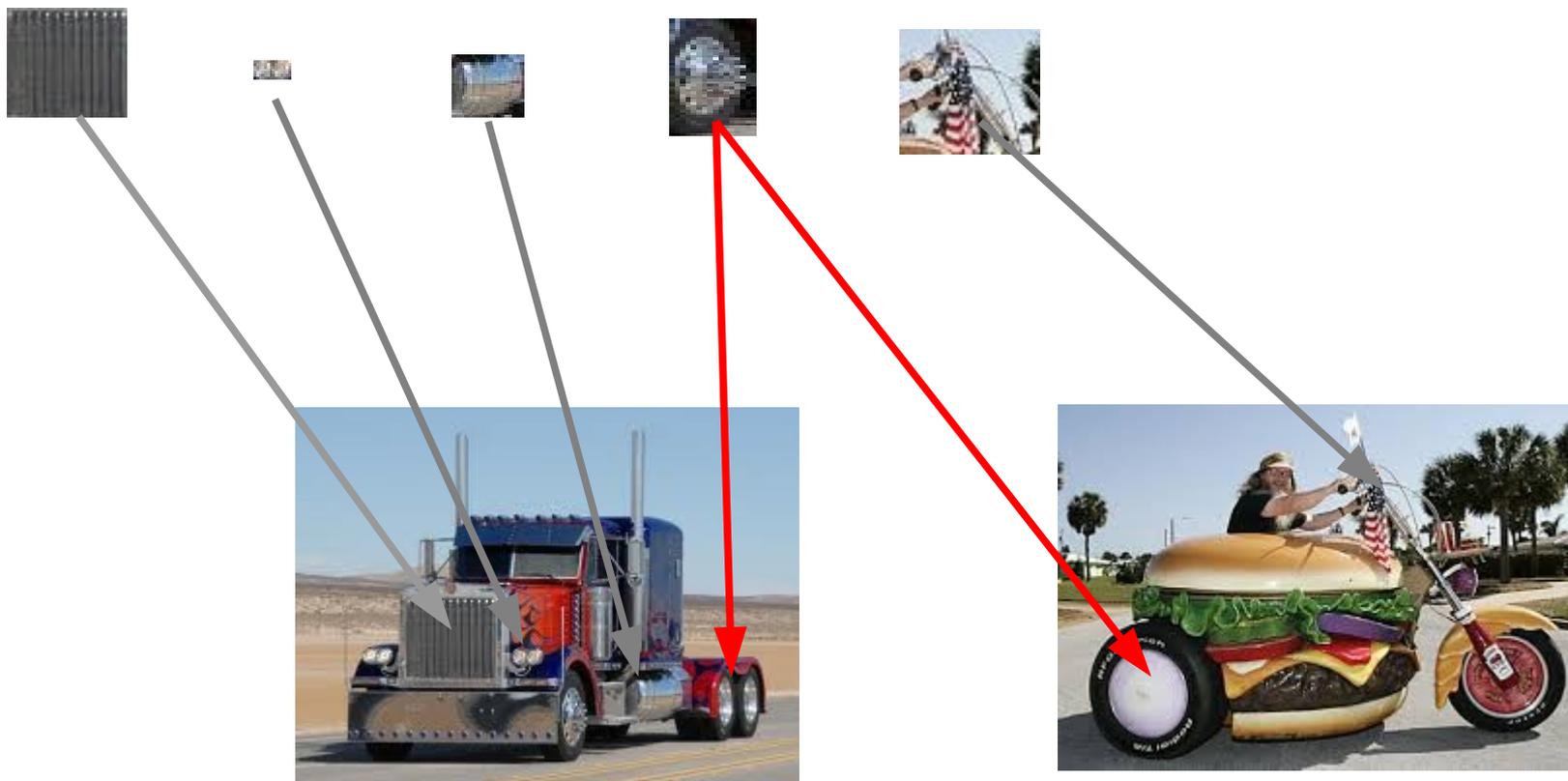


Exponentially more efficient than a 1-of-N representation (a la k-means)

# Composition: sharing

[1 1 0 0 0 1 0 **1** 0 0 0 0 1 1 0 1... ] motorbike

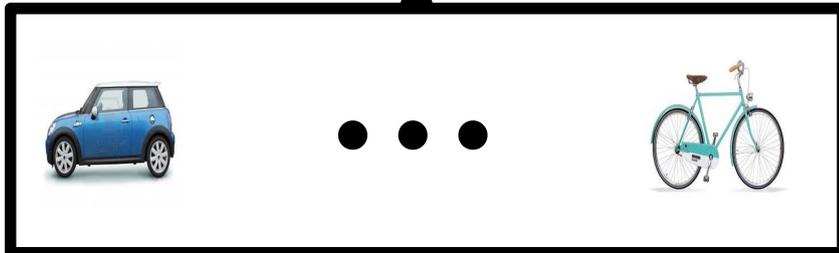
[0 0 1 0 0 0 0 **1** 0 0 1 1 0 0 1 0 ... ] truck



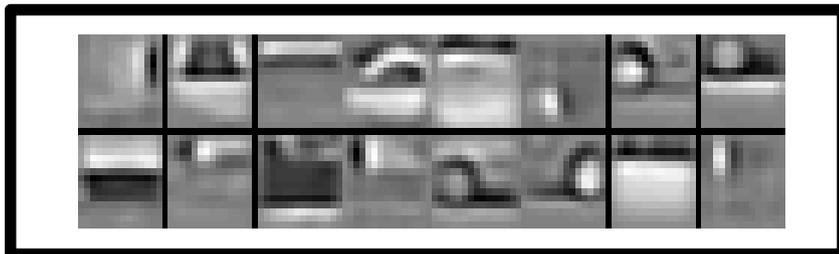
# Composition

prediction of class

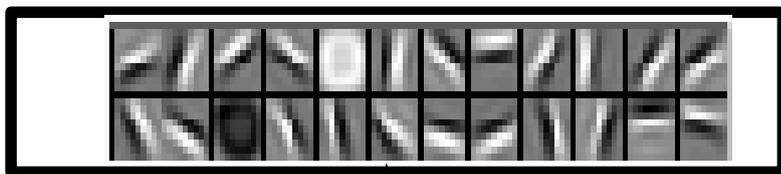
high-level parts



mid-level parts



low level parts



- distributed representations
- feature sharing
- compositionality

Input image



**GOOD: (exponentially)  
more efficient**

**Deep Learning**

**=**

**Representation Learning**

# Ideal Features



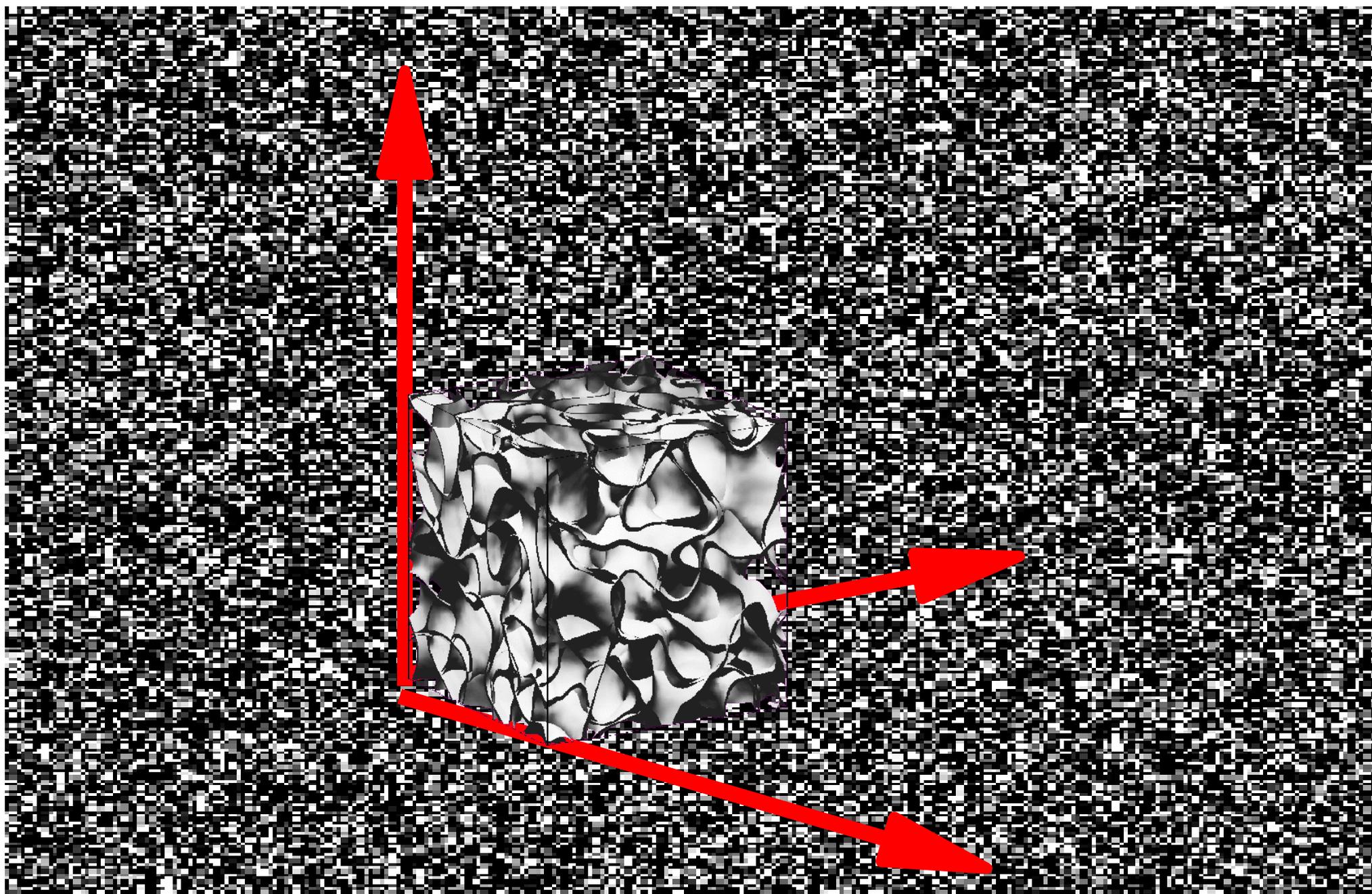
- window, right
- chair, left
- monitor, top of shelf
- carpet, bottom
- drums, corner
- ...

**Ideal  
Feature  
Extractor**

- pillows on couch

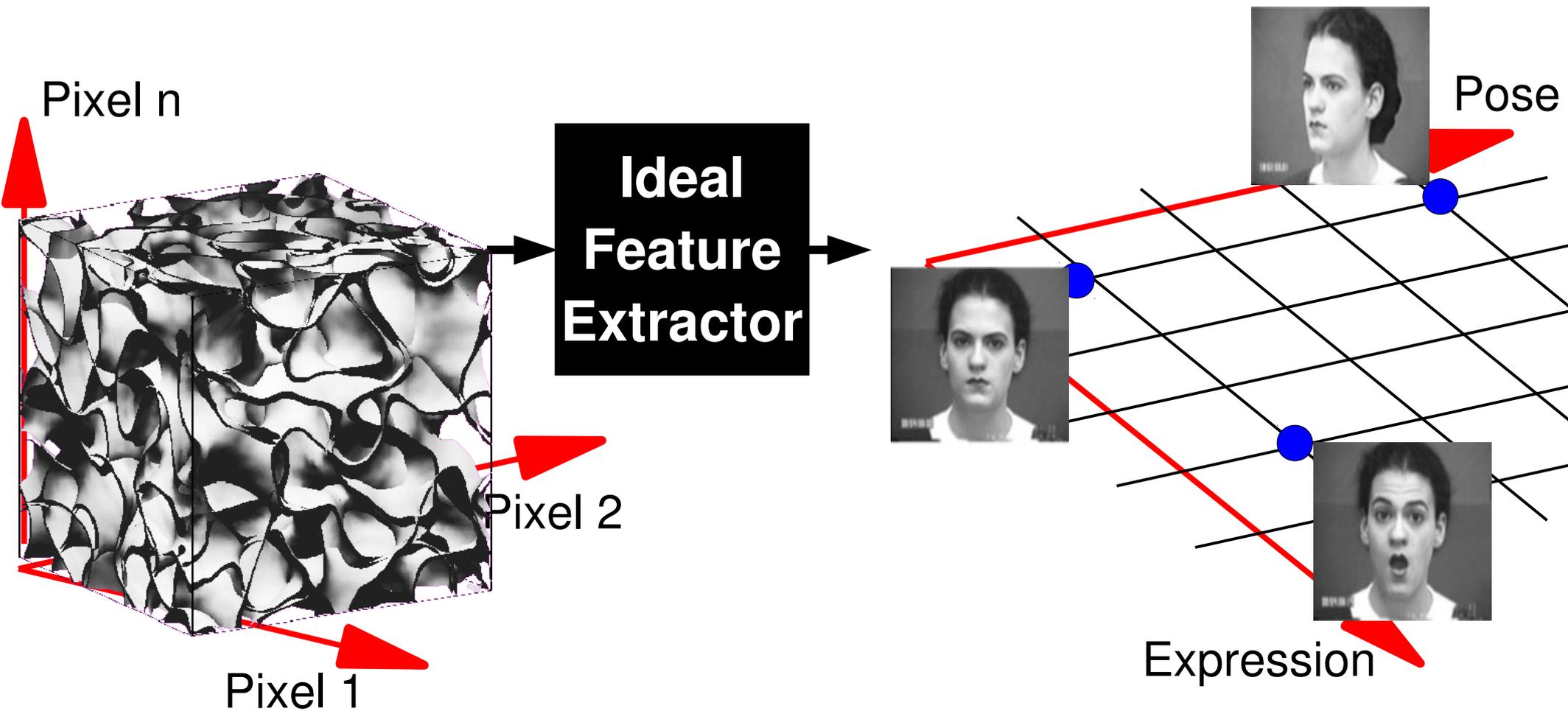
**Q.:** What objects are in the image? Where is the lamp?  
What is on the couch? ...

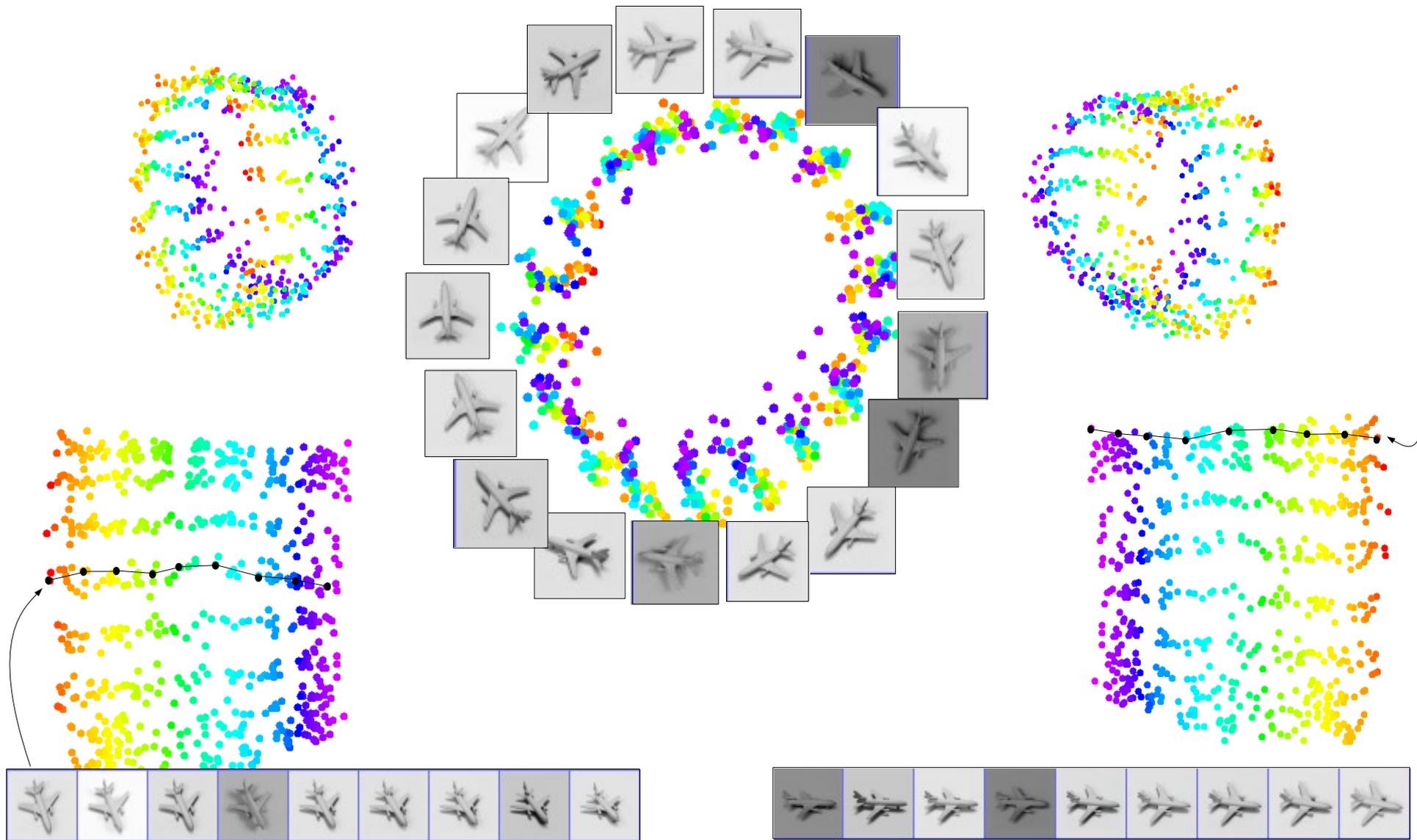
# The Manifold of Natural Images



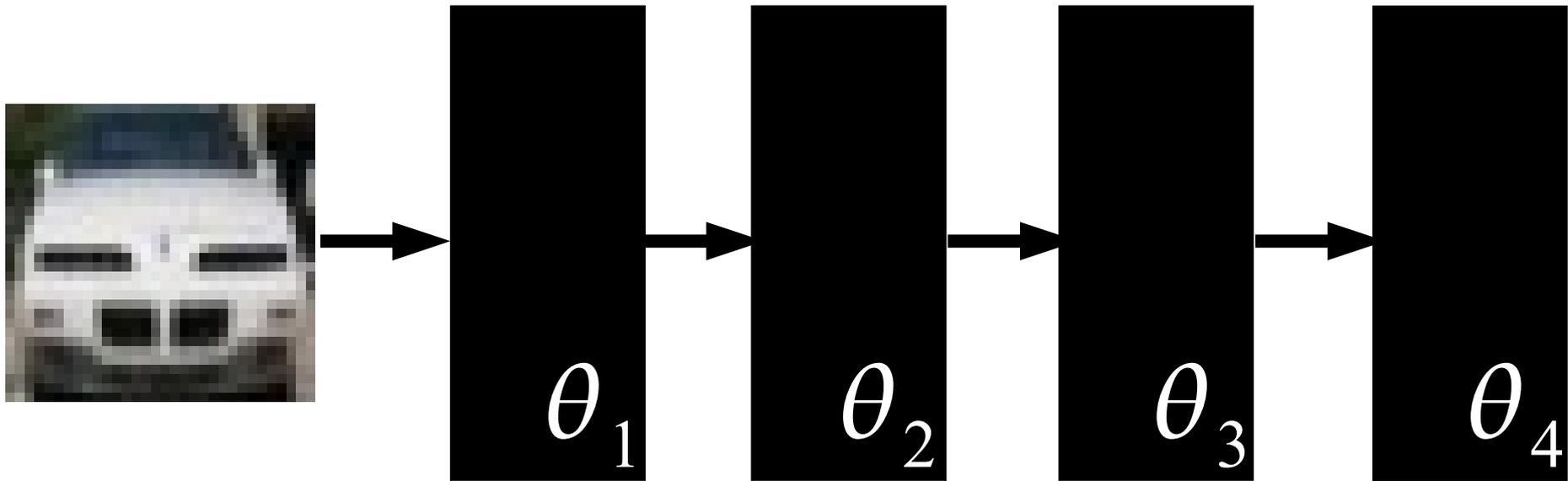
# Ideal Feature Extraction

E.g.: face images live in about 60-D manifold (x,y,z, pitch, yaw, roll, 53 muscles).





# Deep Learning



*Given lots of data, engineer less and learn more!!  
Let the data find the structure (intrinsic dimensions).*

# Deep Learning in Practice

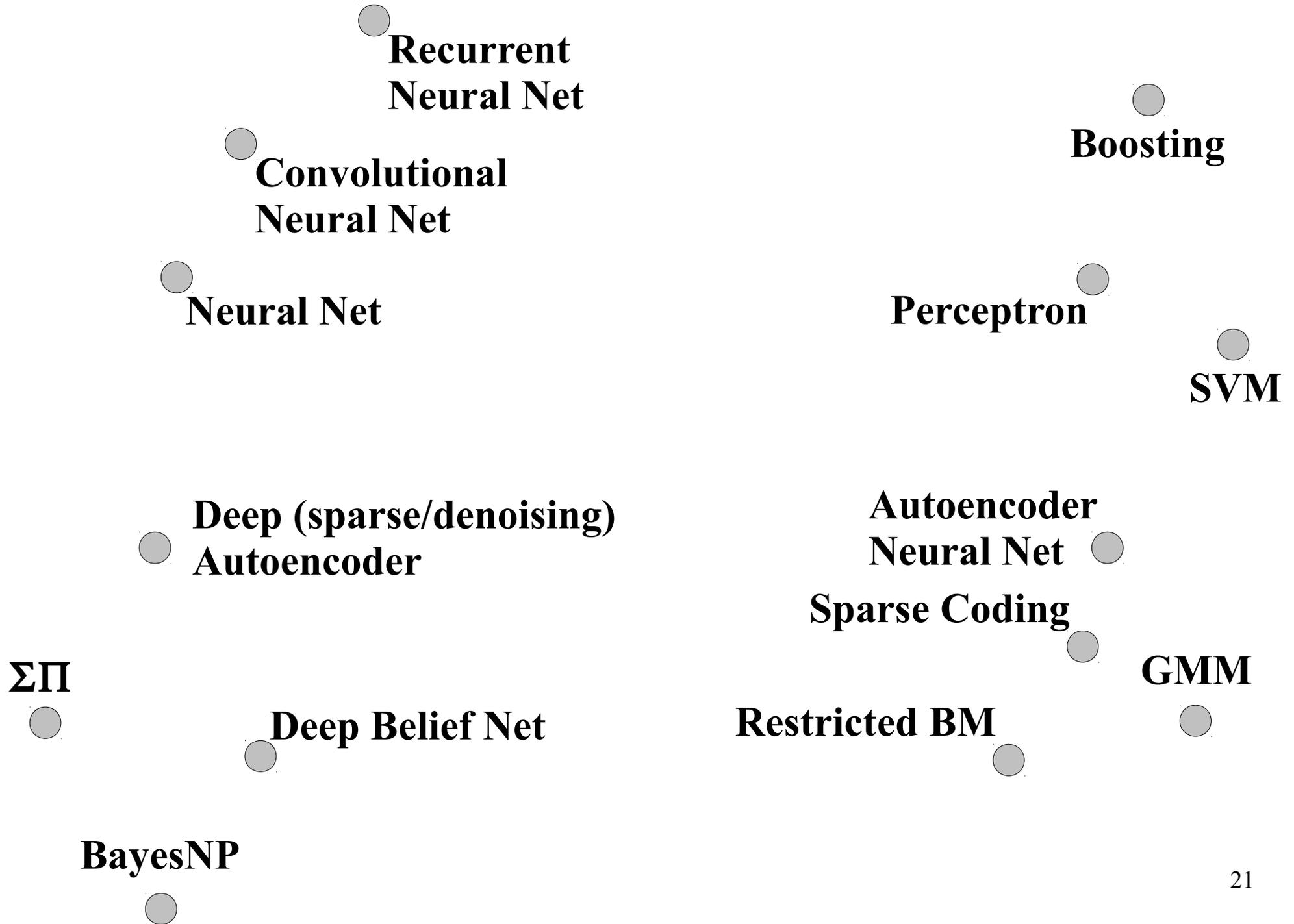
It works very well in practice:

The image is a collage illustrating the practical application of deep learning in image search. At the top center is a Google search interface with the query "my photos sunset". Below the search bar, there are navigation tabs for "Web", "Images", "Maps", "Shopping", and "More", along with "Search tools". The search results show "70 personal results. 137,000,000 other results." and a grid of sunset photos. A large, semi-transparent "facebook" logo is overlaid on the center of the image. To the left, there is a Baidu search interface with a grid of fashion-related images. To the right, there is a dark blue area featuring the Android logo (a green robot) and the word "ANDROID" in white. At the bottom, there are logos for Microsoft, IBM, and Facebook. The word "amazon.com" is written in a white box with a black border, tilted at an angle in the upper left. The word "NEC" is written in a white box with a black border, tilted at an angle in the upper right.

# KEY IDEAS OF DEEP LEARNING

- Hierarchical non-linear system
  - Distributed representations
  - Sharing
- End-to-end learning
  - Joint optimization of features and classifier
  - Good features are learned as a side product of the learning process

# THE SPACE OF MACHINE LEARNING METHODS



**SHALLOW**

**Recurrent  
Neural Net**

**Convolutional  
Neural Net**

**Neural Net**

**Boosting**

**Perceptron**

**SVM**

**Deep (sparse/denoising)  
Autoencoder**

**Autoencoder  
Neural Net**

**Sparse Coding**

**GMM**

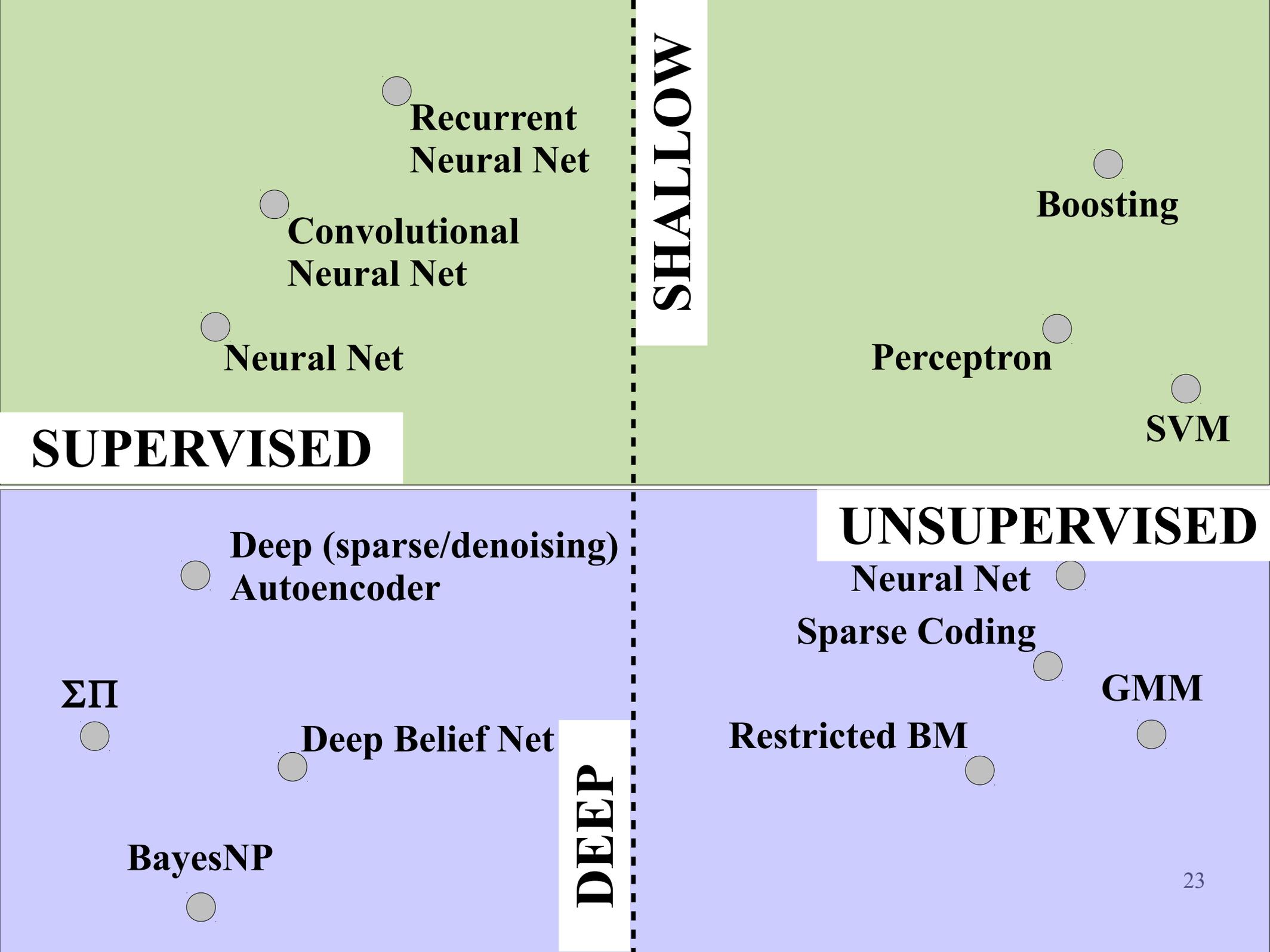
$\Sigma\Pi$

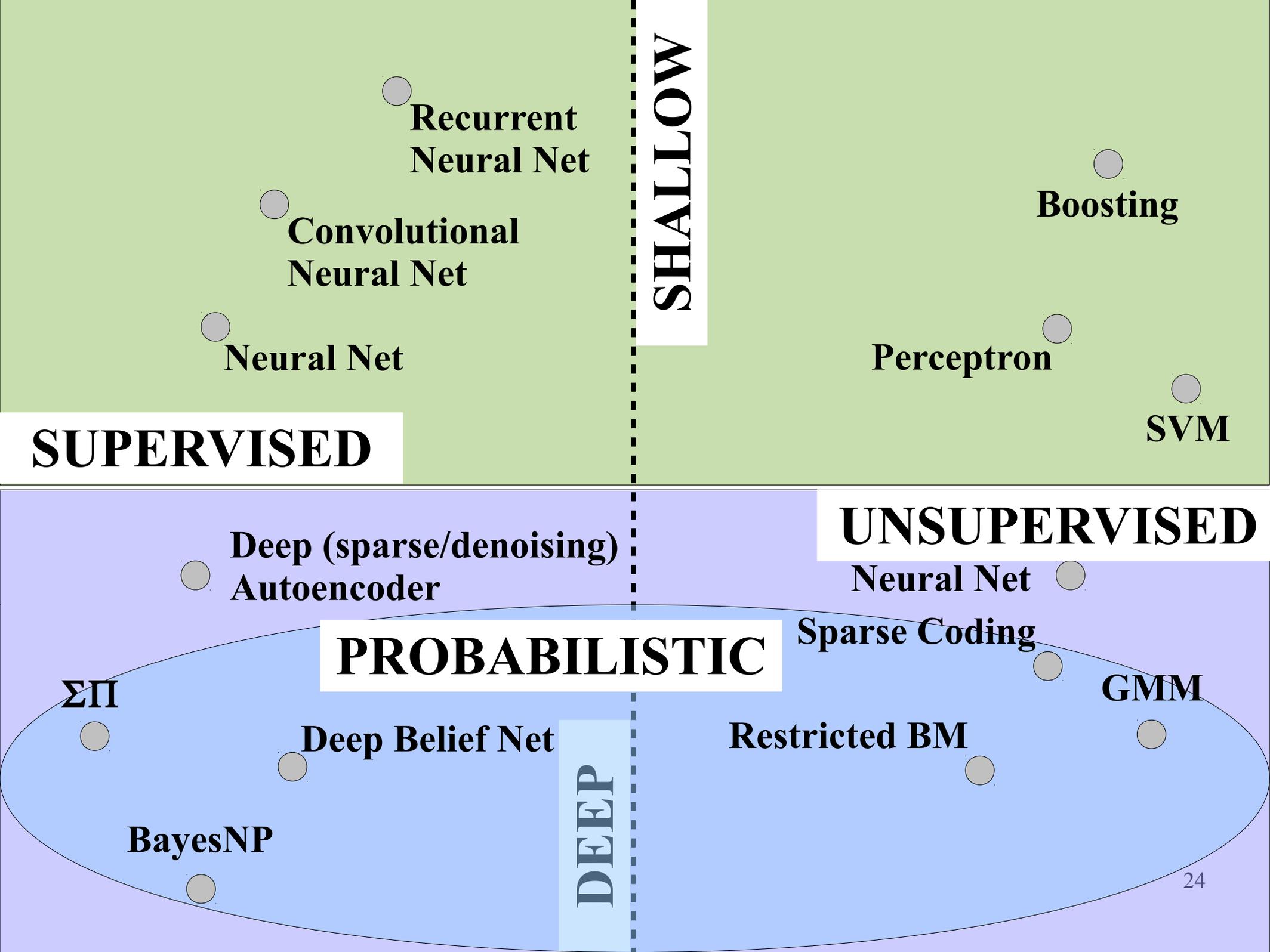
**Deep Belief Net**

**Restricted BM**

**BayesNP**

**DEEP**





**SHALLOW**

Recurrent  
Neural Net

Convolutional  
Neural Net

Neural Net

Boosting

Perceptron

SVM

**SUPERVISED**

**UNSUPERVISED**

Deep (sparse/denoising)  
Autoencoder

Neural Net

**PROBABILISTIC**

Sparse Coding

$\Sigma\Pi$

Deep Belief Net

GMM

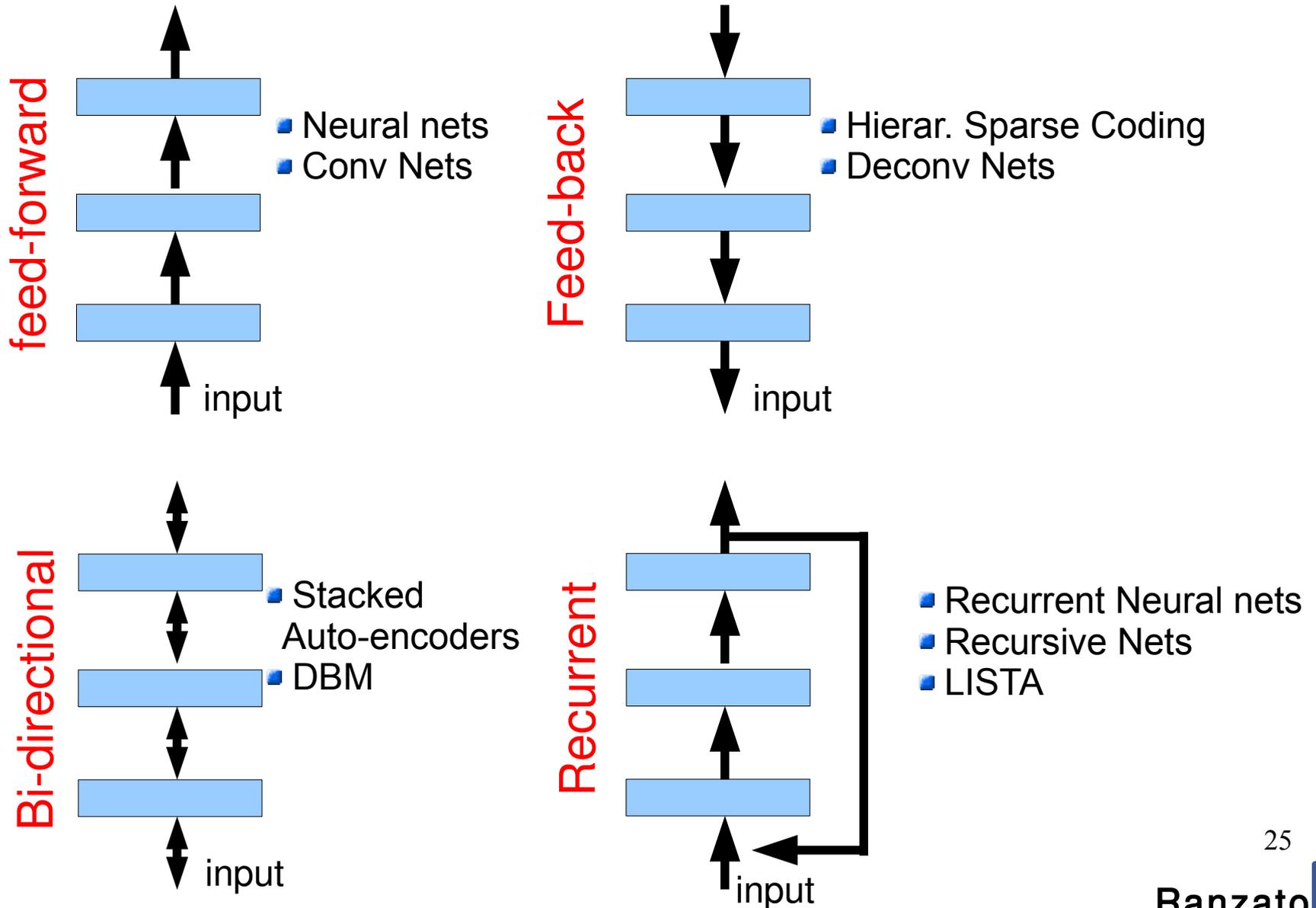
BayesNP

Restricted BM

**DEEP**

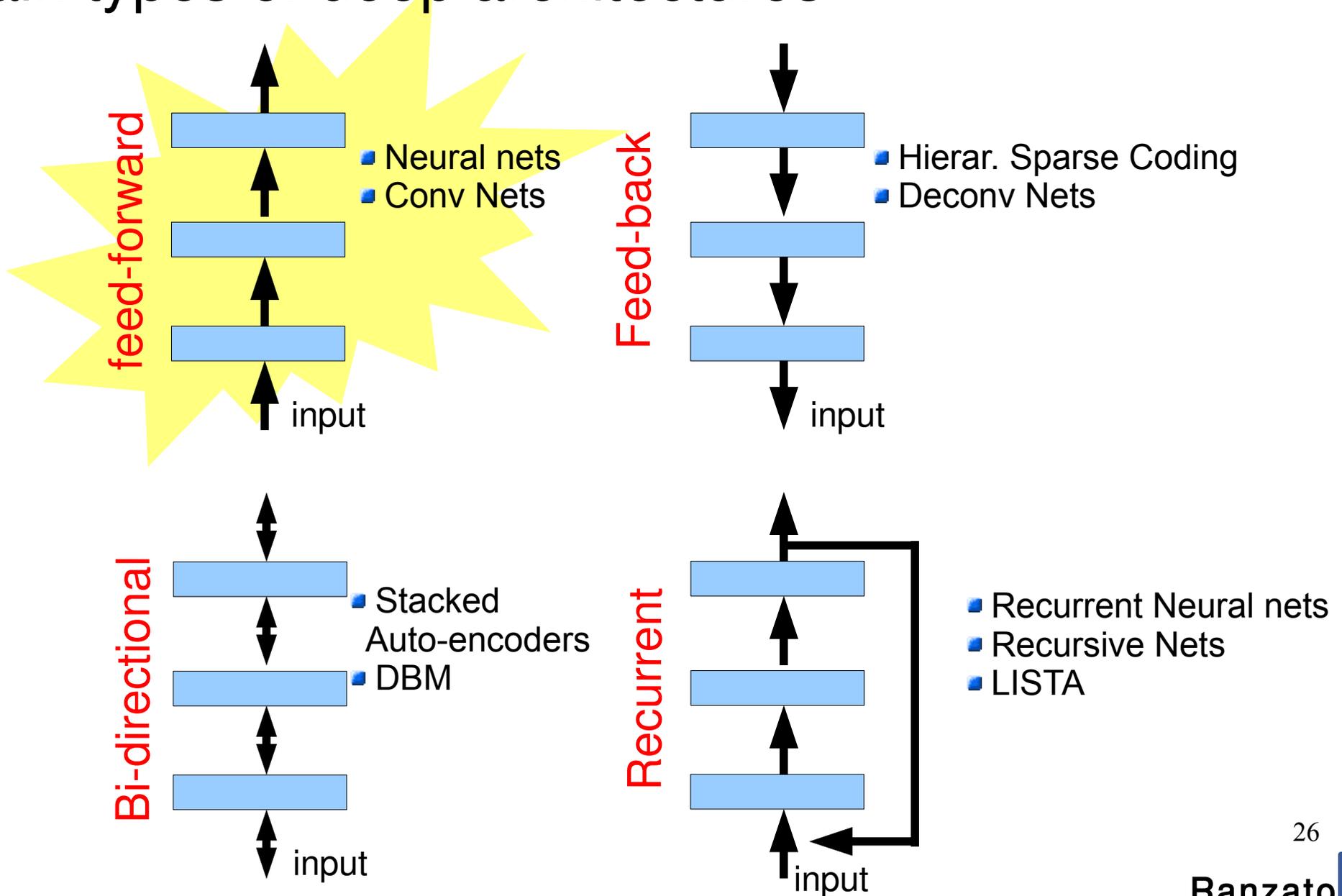
# Deep Learning is B I G

- Main types of deep architectures



# Deep Learning is B I G

- Main types of deep architectures



# Deep Learning is BIG

- Main types of learning protocols
  - Purely supervised
    - Backprop + SGD
    - Good when there is lots of labeled data.
  - Layer-wise unsupervised + superv. linear classifier
    - Train each layer in sequence using regularized auto-encoders or RBMs
    - Hold fix the feature extractor, train linear classifier on features
    - Good when labeled data is scarce but there is lots of unlabeled data.
  - Layer-wise unsupervised + supervised backprop
    - Train each layer in sequence
    - Backprop through the whole system
    - Good when learning problem is very difficult.

# Deep Learning is B I G

- Main types of learning protocols
  - Purely supervised
    - Backprop + SGD
    - Good when there is lots of labeled data.
  - Layer-wise unsupervised + superv. linear classifier
    - Train each layer in sequence using regularized auto-encoders or RBMs
    - Hold fix the feature extractor, train linear classifier on features
    - Good when labeled data is scarce but there is lots of unlabeled data.
  - Layer-wise unsupervised + supervised backprop
    - Train each layer in sequence
    - Backprop through the whole system
    - Good when learning problem is very difficult.

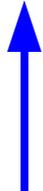
# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- Examples:
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - Unsupervised learning: sparse coding
  - Unsupervised learning: gated MRF
- Other examples
- Practical tricks

# Energy-Based Models: Energy Function

- Energy:

$$E(y; \theta) \quad \text{or} \quad E(y; x, \theta)$$

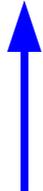
 or 

unsupervised supervised

# Energy-Based Models: Energy Function

- Energy:

$$E(y; \theta) \quad \text{or} \quad E(y; x, \theta)$$

 unsupervised                       supervised

$y$  can be  $\left\{ \begin{array}{l} \text{discrete} \\ \text{continuous} \end{array} \right.$

# Energy-Based Models: Energy Function

- Energy:

$$E(y; \theta) \quad \text{or} \quad E(y; x, \theta)$$

↑  
unsupervised

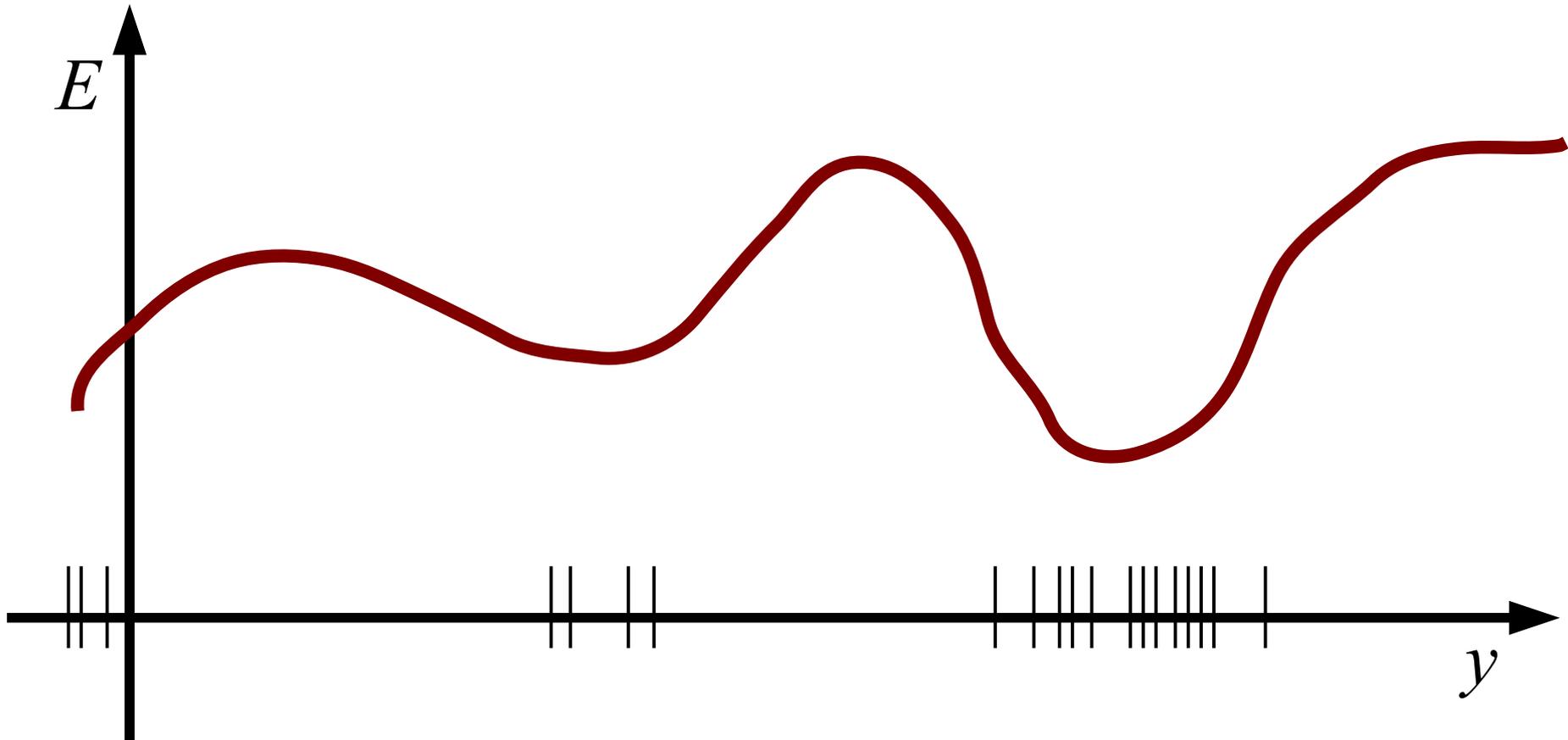
↑  
supervised

$y$  can be { discrete  
continuous

**We will refer to the unsupervised/continuous case, but much of the following applies to the other cases as well.**

# Energy-Based Models: Energy Function

- Energy should be lower for desired output



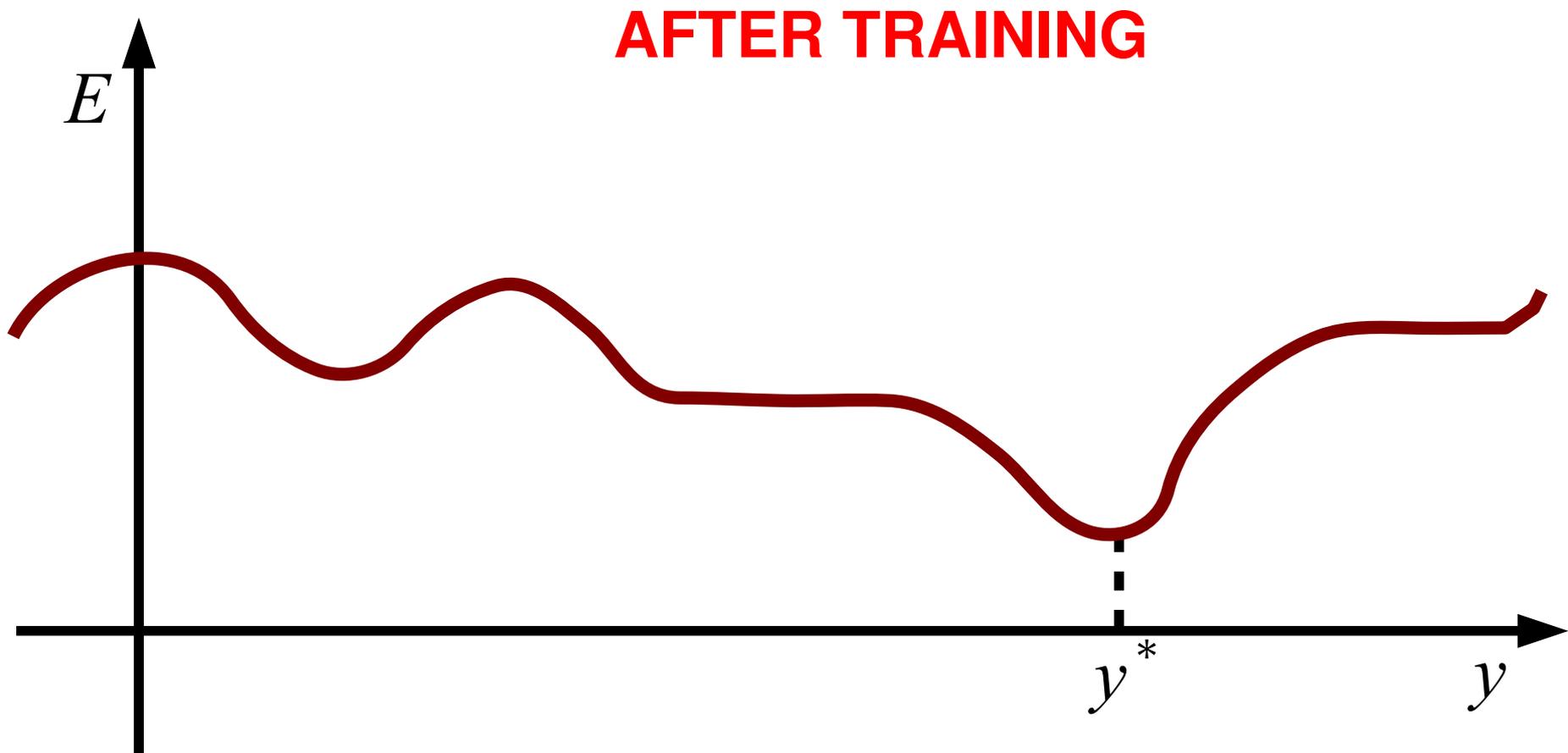
# Energy-Based Models: Energy Function

- Make energy lower at the desired output



# Energy-Based Models: Energy Function

- Make energy lower at the desired output



# Energy-Based Models: Energy Function

- Examples of energy function:

- PCA

$$E(y) = \|y - W W^T y\|_2^2$$

- 
- Linear binary classifier  $y \in \{-1, +1\}$

$$E(y; x) = -y(W^T x)$$

- 
- Neural net binary classifier

$$E(y; x) = -y(W_2^T f(x; W_1))$$

# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

- Examples of loss function:

- PCA

$$E(y) = \|y - W W^T y\|_2^2 \quad L(E(y)) = E(y)$$

- Logistic regression classifier

$$E(y; x) = -y(W^T x)$$

$$L(E(y; x)) = \log(1 + \exp(E(y; x)))$$

# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

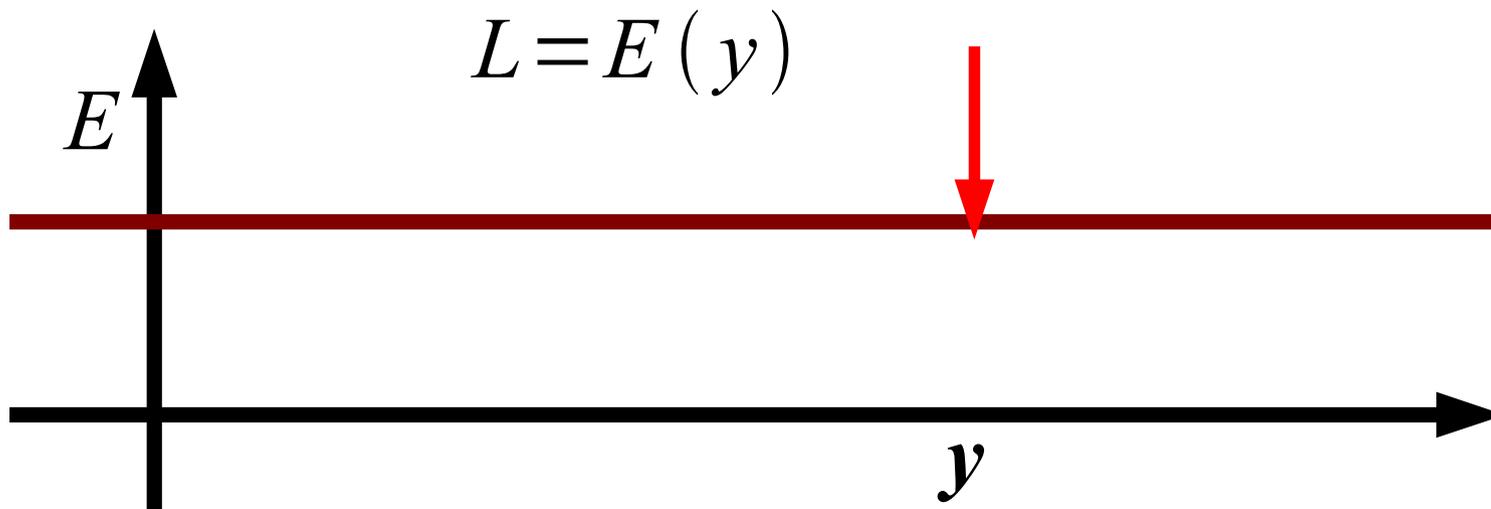
**How to design loss good functions?**

# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

**How to design loss good functions?**

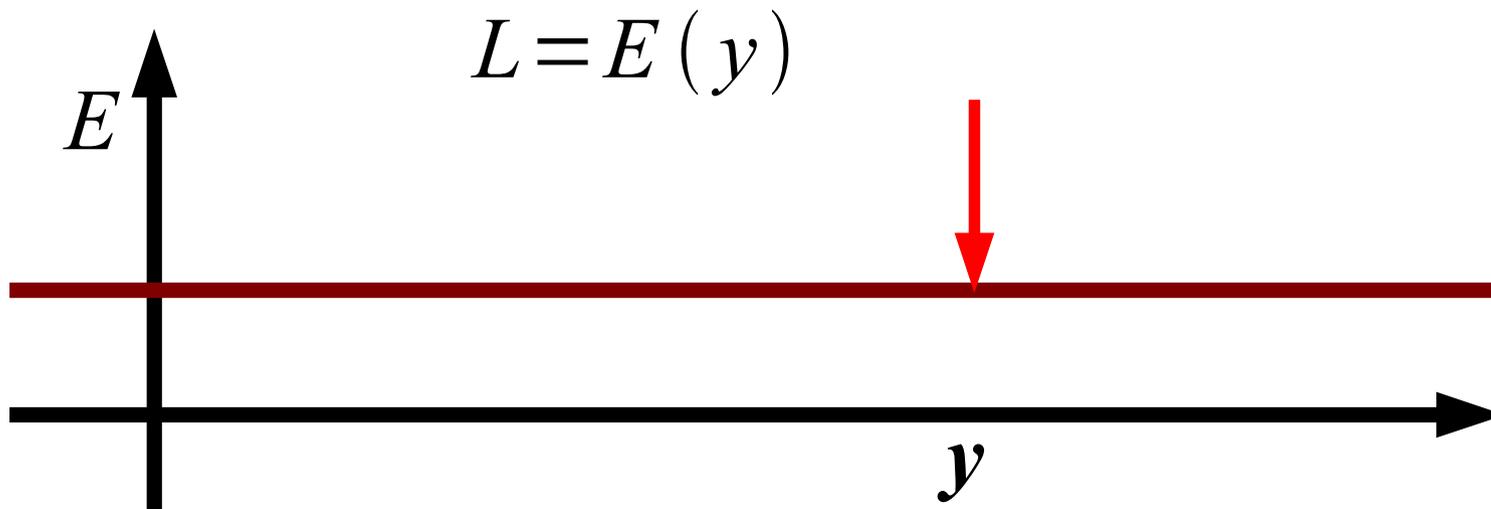


# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

**How to design loss good functions?**

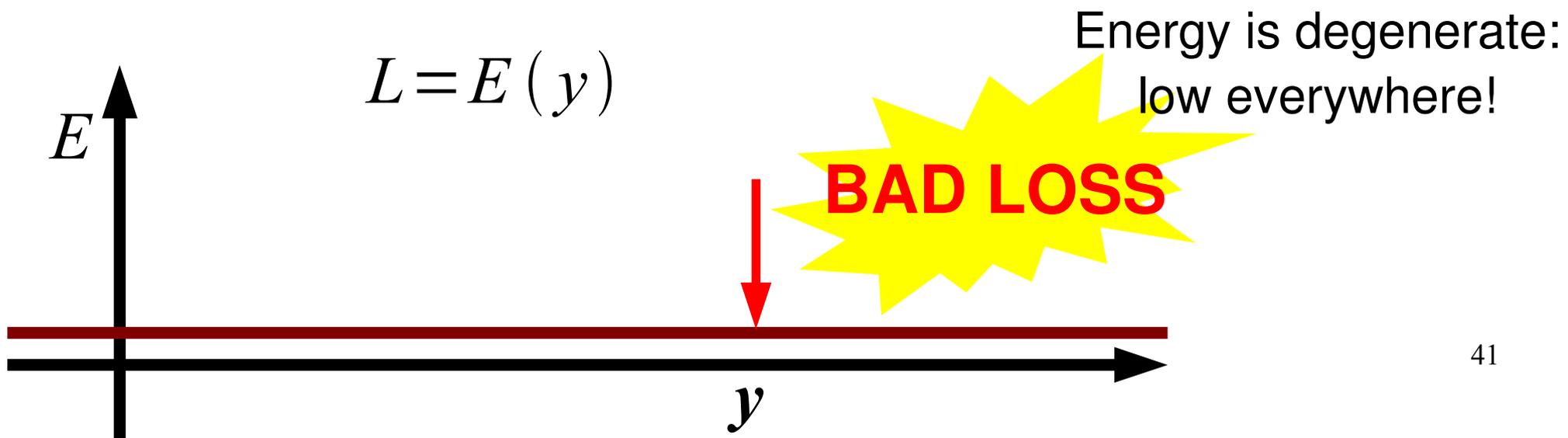


# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

**How to design loss good functions?**



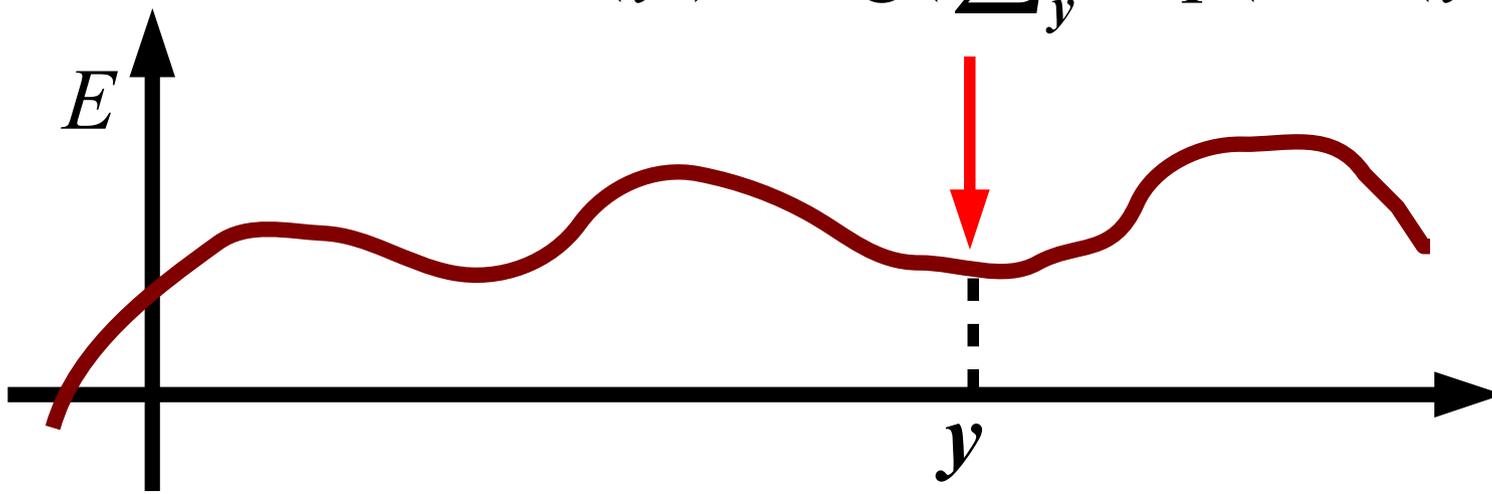
# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

**How to design loss good functions?**

$$L = E(y) + \log\left(\sum_{\bar{y}} \exp(-E(\bar{y}))\right)$$



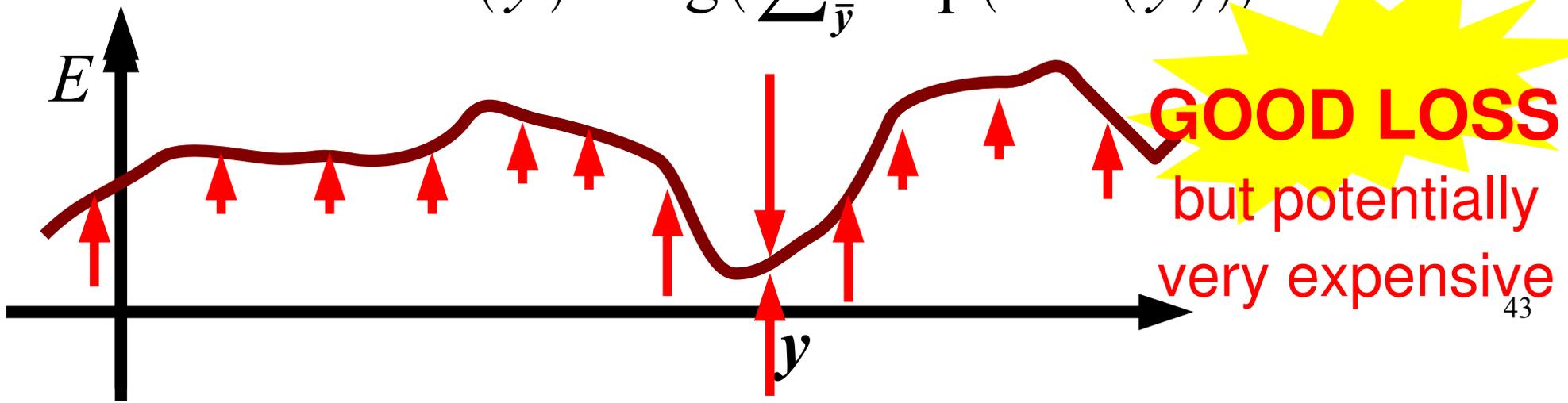
# Energy-Based Models: Loss Function

- Loss is a function of the energy
- Minimizing the loss over the training set yields the desired energy landscape.

$$\theta^* = \min_{\theta} \sum_p L(E(y^p; \theta))$$

**How to design loss good functions?**

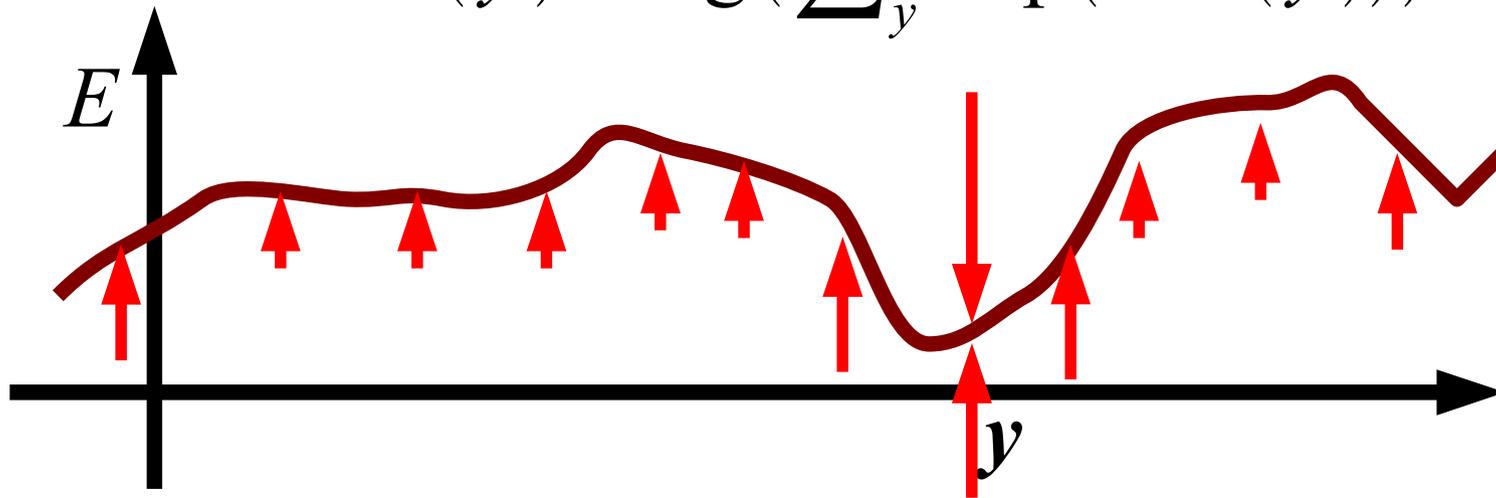
$$L = E(y) + \log\left(\sum_{\bar{y}} \exp(-E(\bar{y}))\right)$$



# Strategies to Shape E: #1

- Pull down the correct answer and pull up everywhere else.

$$L = E(y) + \log\left(\sum_{\bar{y}} \exp(-E(\bar{y}))\right)$$



Negative  
Log-Likelihood

$$p(y) = \frac{e^{-E(y)}}{\sum_u e^{-E(u)}}$$

## PROS

It produces calibrated probabilities.

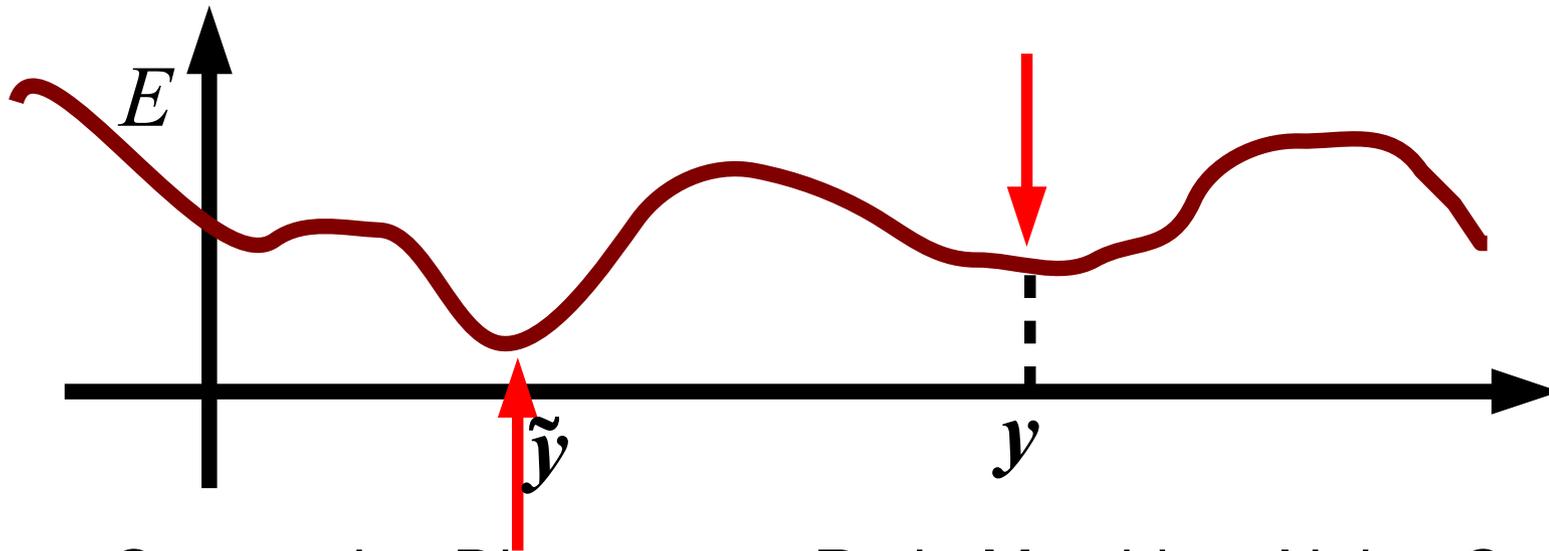
## CONS

Expensive to compute when  $y$  is discrete and high dimensional.  
Generally intractable when  $y$  is continuous.

# Strategies to Shape E: #2

- Pull down the correct answer and pull up carefully chosen points.

$$L = \max(0, m - (E(\bar{y}) - E(y))), \text{ e.g. } \bar{y} = \min_{\tilde{y} \neq y} E(\tilde{y})$$



E.g.: Contrastive Divergence, Ratio Matching, Noise Contrastive Estimation, Minimum Probability Flow...

Hinton et al. "A fast learning algorithm for DBNs" Neural Comp. 2008

Hyvarinen "Some extensions of score matching" Comp Stats 2007

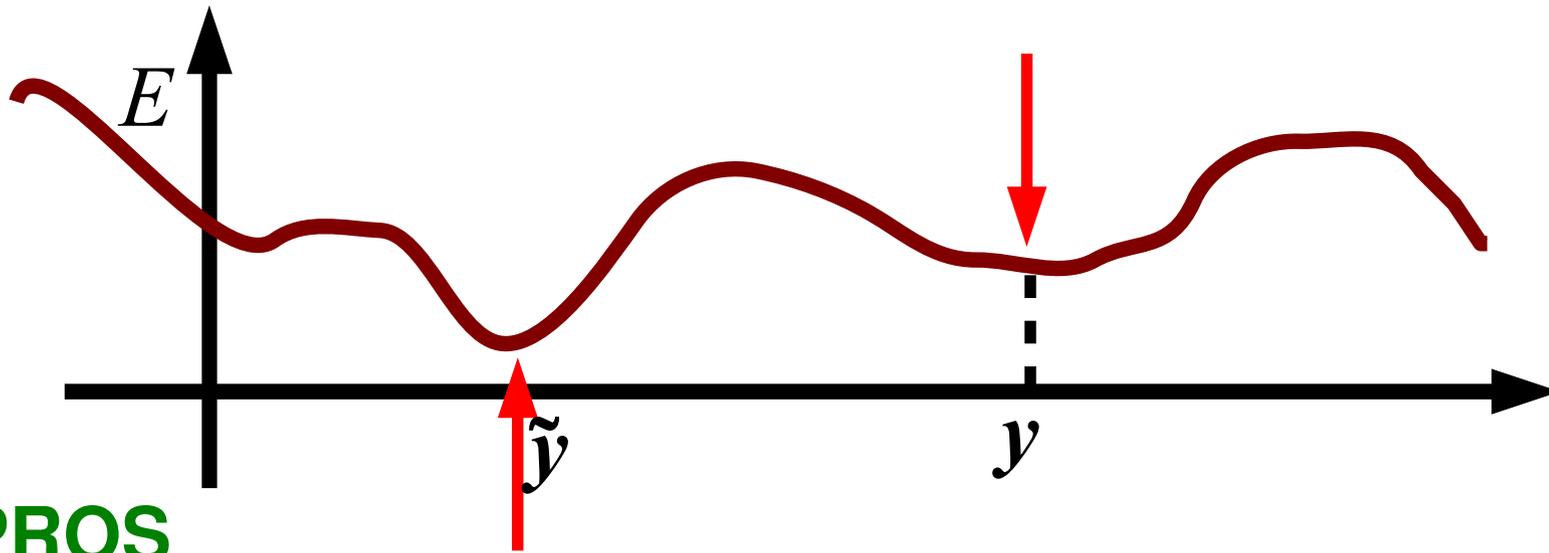
Gutmann et al. "Noise contrastive estimation of unnormalized..." JMLR 2012

Sohl-Dickstein et al. "Minimum probability flow learning" ICML 2011

# Strategies to Shape E: #2

- Pull down the correct answer and pull up carefully chosen points.

$$L = \max(0, m - (E(\bar{y}) - E(y))), \text{ e.g. } \bar{y} = \min_{\tilde{y} \neq y} E(\tilde{y})$$



## PROS

Efficient.

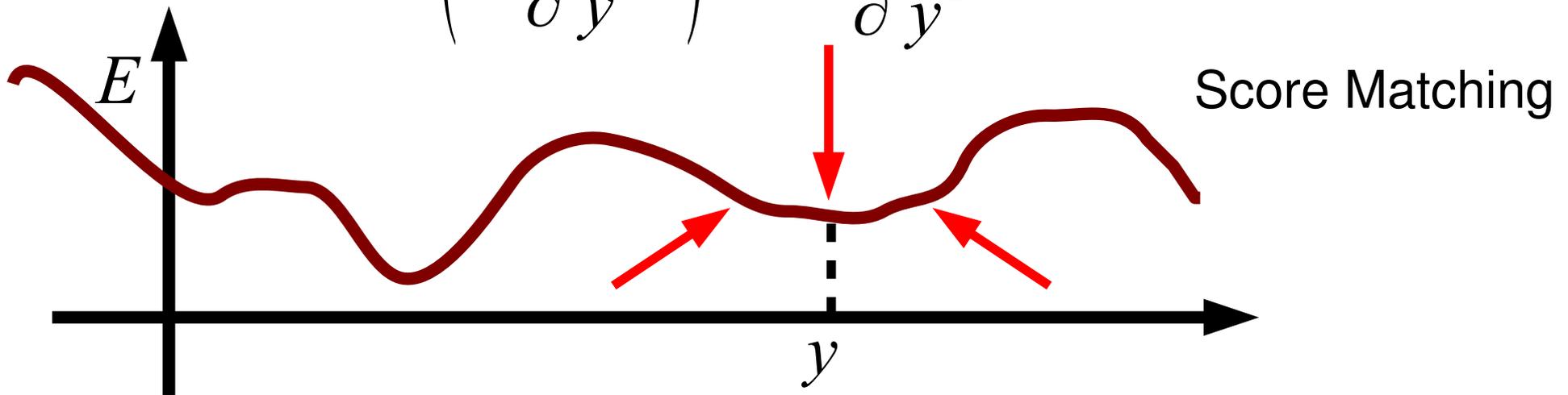
## CONS

The criterion to pick where to pull up is tricky (overall in high dimensional spaces): trades-off computational and statistical efficiency.

# Strategies to Shape E: #3

- Pull down the correct answer and increase local curvature.

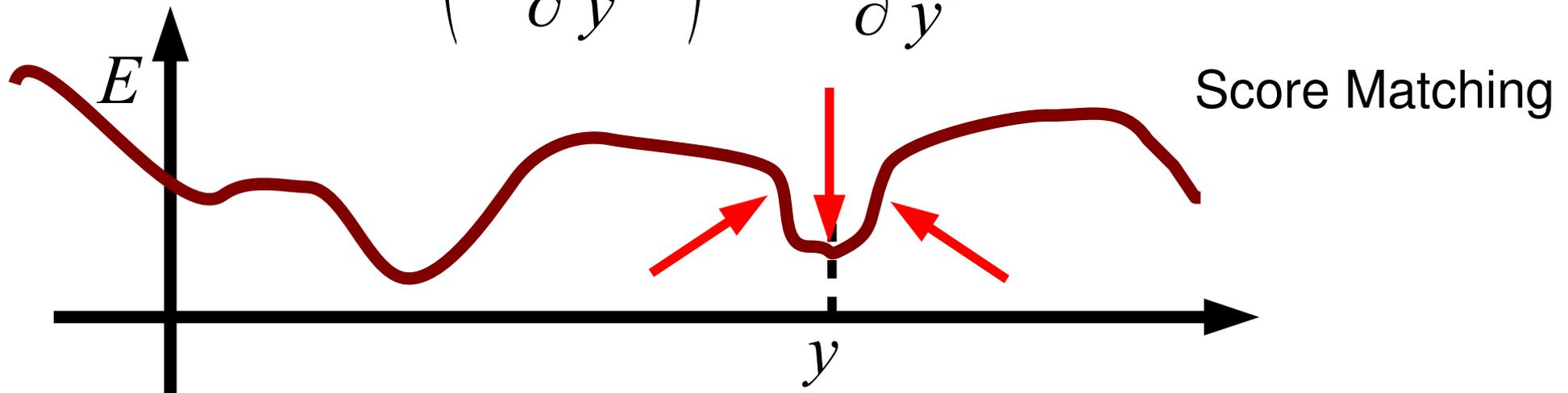
$$L = \left( \frac{\partial E(y)}{\partial y} \right)^2 + \frac{\partial^2 E(y)}{\partial y^2}$$



# Strategies to Shape E: #3

- Pull down the correct answer and increase local curvature.

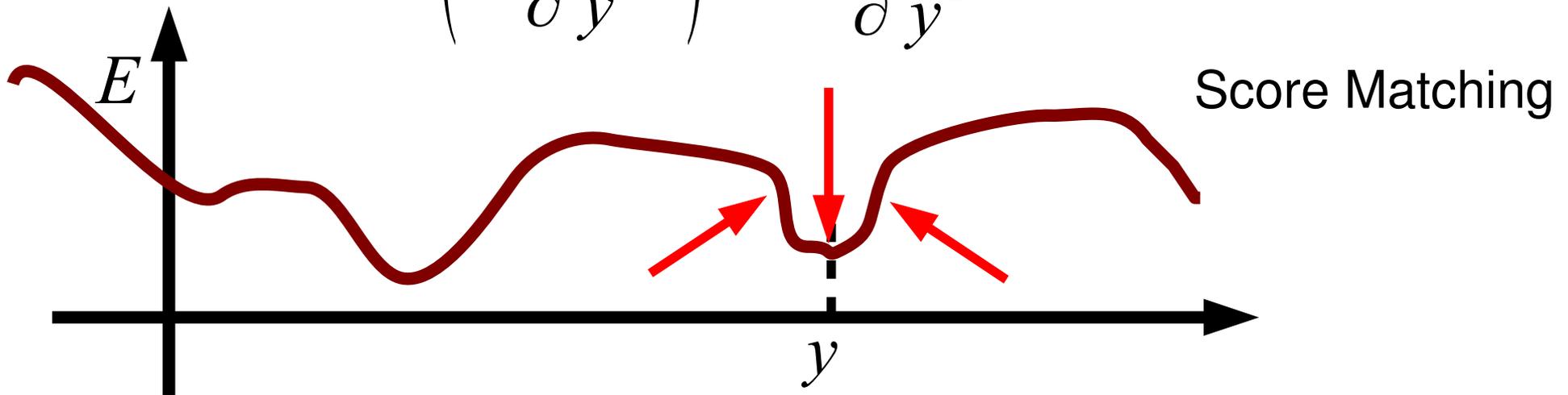
$$L = \left( \frac{\partial E(y)}{\partial y} \right)^2 + \frac{\partial^2 E(y)}{\partial y^2}$$



# Strategies to Shape E: #3

- Pull down the correct answer and increase local curvature.

$$L = \left( \frac{\partial E(y)}{\partial y} \right)^2 + \frac{\partial^2 E(y)}{\partial y^2}$$



## PROS

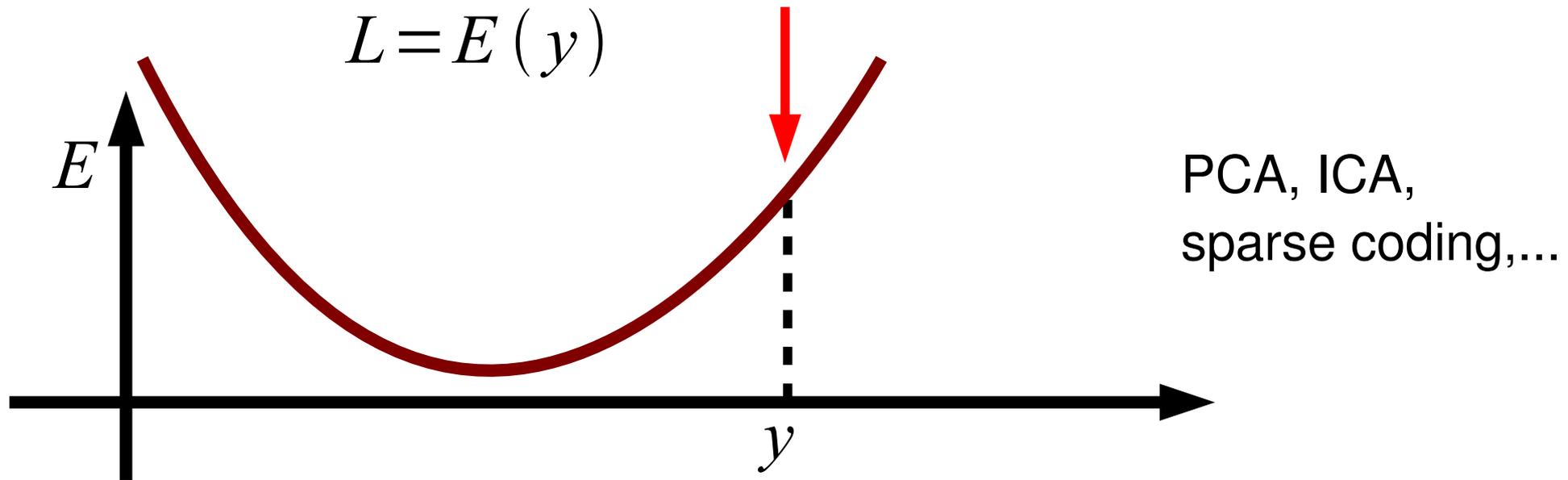
Efficient in continuous but not too high dimensional spaces.

## CONS

Very complicated to compute and not practical in very high dimensional spaces. Not applicable in discrete spaces.

# Strategies to Shape E: #4

- Pull down correct answer and have global constrain on the energy: only few minima exist



## PROS

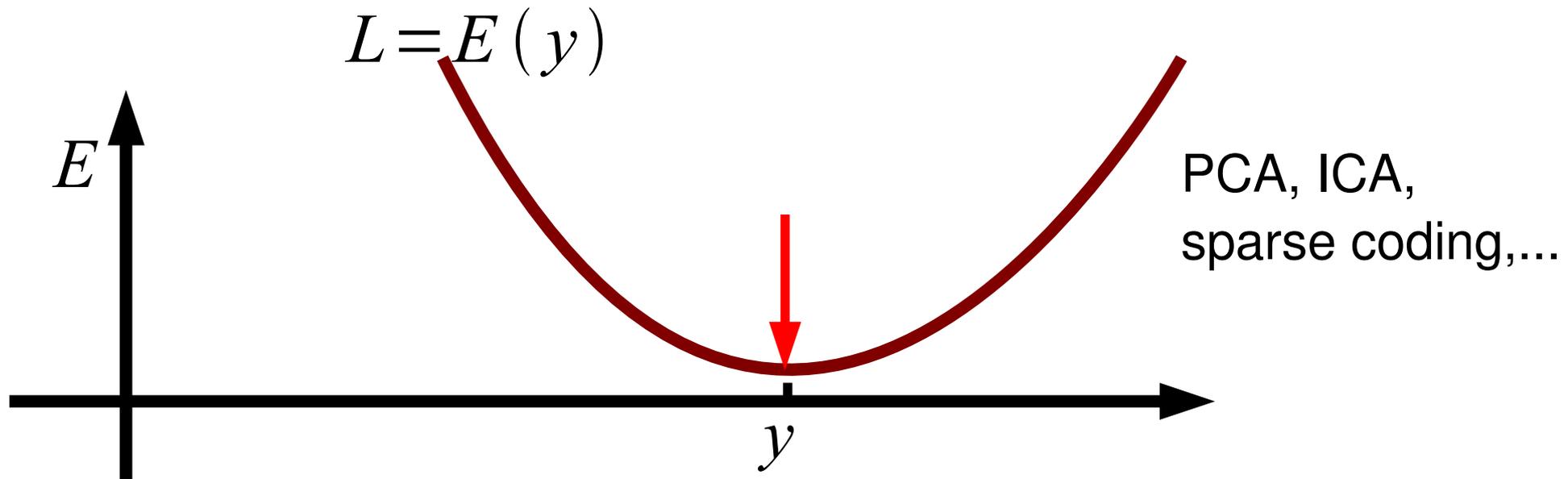
Efficient in continuous, high dimensional spaces.

## CONS

Need to design good global constraints. Used in unsup. learning only.

# Strategies to Shape E: #4

- Pull down correct answer and have global constrain on the energy: only few minima exist



## PROS

Efficient in continuous, high dimensional spaces.

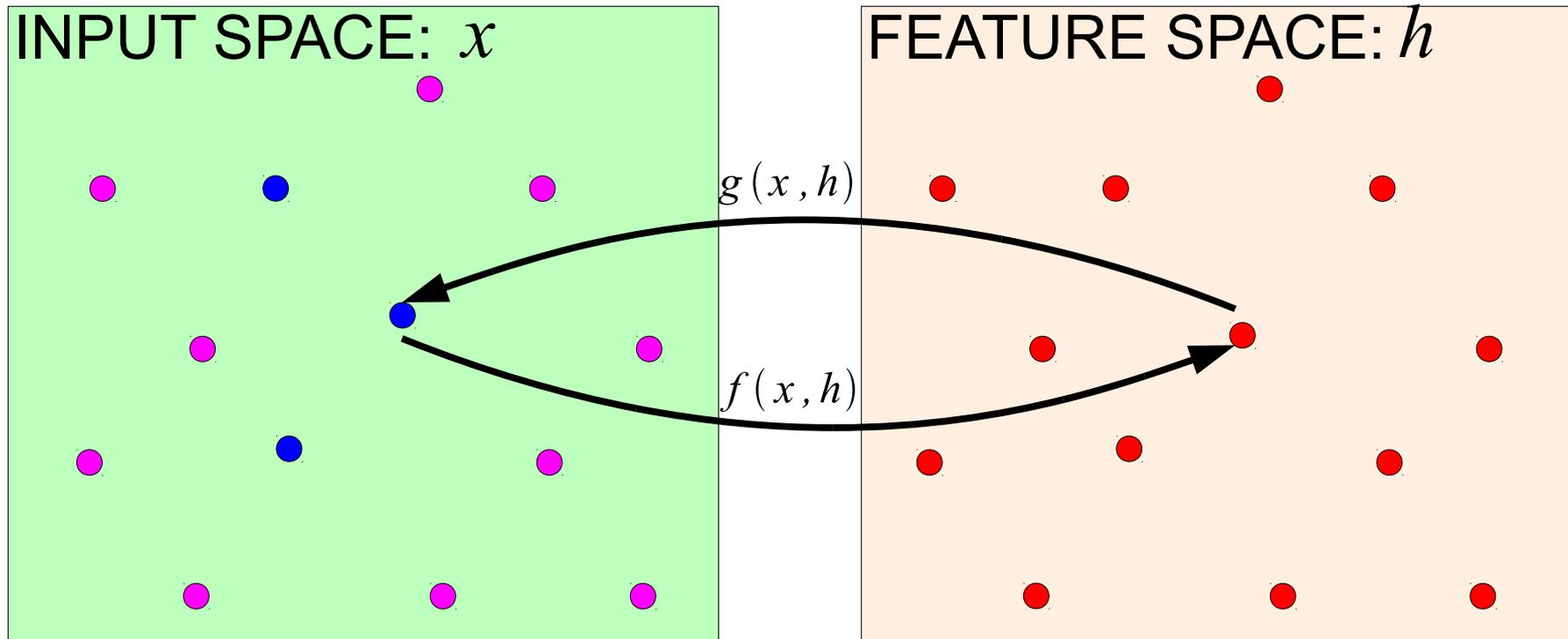
## CONS

Need to design good global constraints. Used in un-sup. learning only.

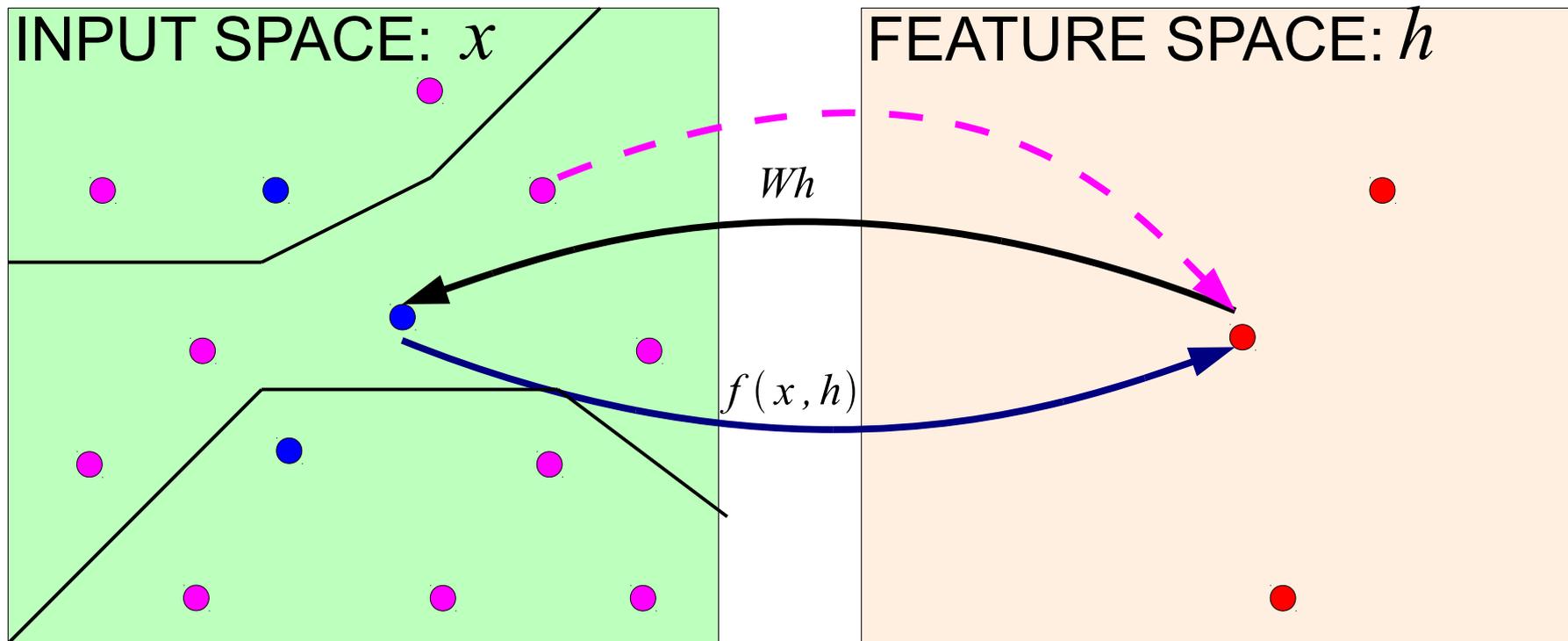
# Strategies to Shape E: #4

- Pull down correct answer and have global constrain on the energy: only few minima exist
- Typical methods (unsup. learning):
  - Use compact internal representation (PCA)
  - Have finite number of internal states (K-Means)
  - Use sparse codes (ICA, sparse coding)

- training sample
- input data point which is not a training sample
- feature (code)



E.g. K-means:  $h$  is 1-of-N.

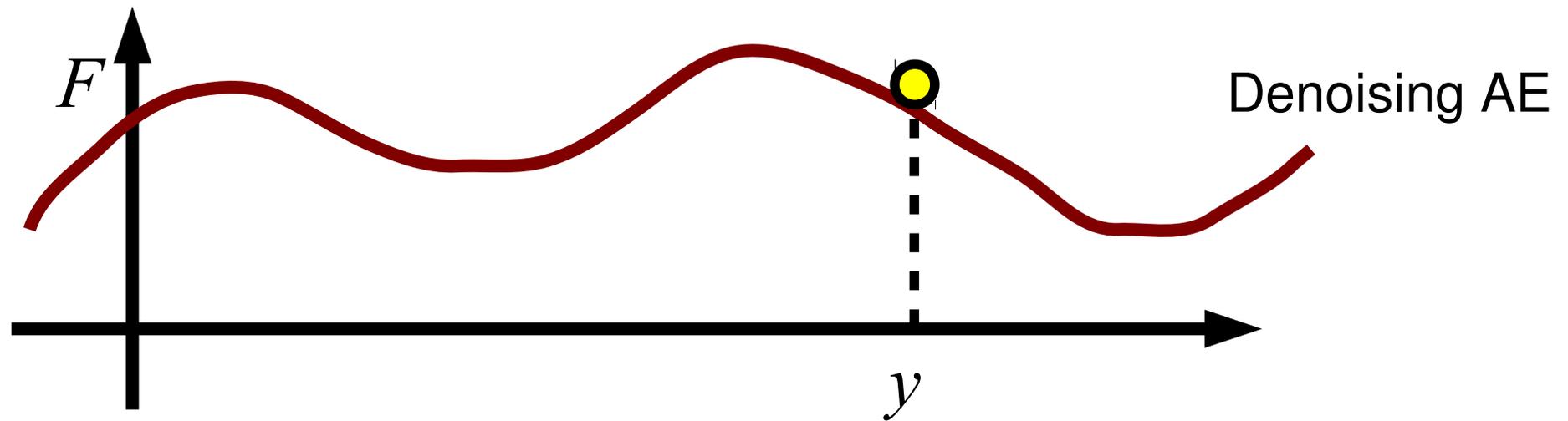


Since there are very few “codes” available and the energy (MSE) is minimized on the training set, the energy must be higher elsewhere.

# Strategies to Shape E: #5

- Make the observed  $y$  an attractor state of some energy function. Need only to define a dynamics.

$$L = E(y)$$



E.g.:

$$L = \frac{1}{2} \|y - \hat{y}\|^2$$

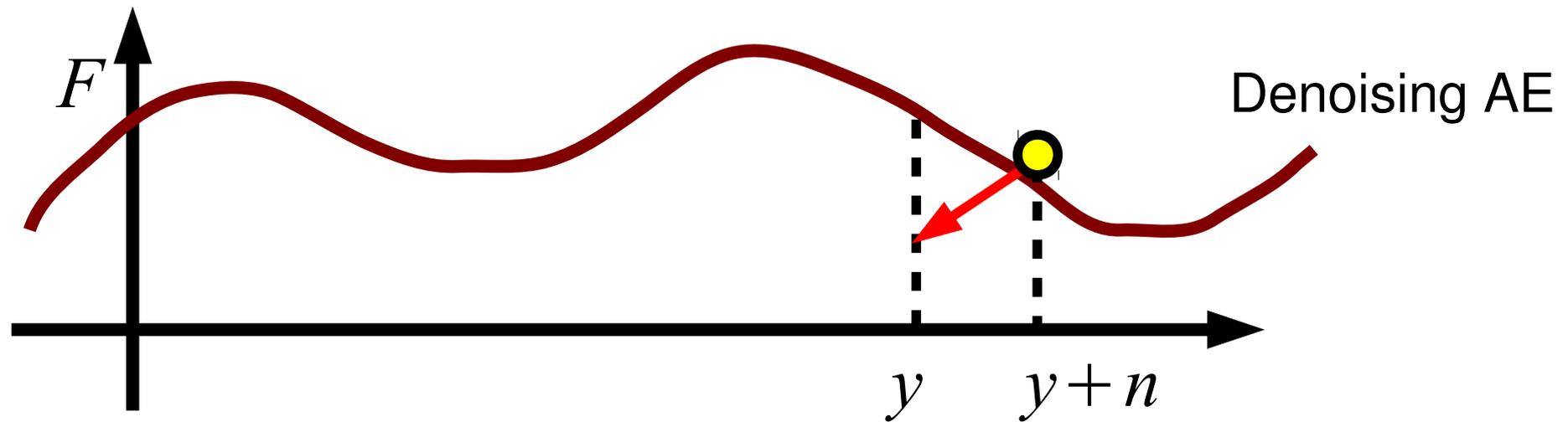
$$\hat{y} = W_2 \sigma(W_1 y + n), n \in N(0, I)$$

Training samples must be stable points (**local minima**) and **attractors**.

# Strategies to Shape E: #5

- Make the observed  $y$  an attractor state of some energy function. Need only to define a dynamics.

$$L = E(y)$$



E.g.:

$$L = \frac{1}{2} \|y - \hat{y}\|^2$$

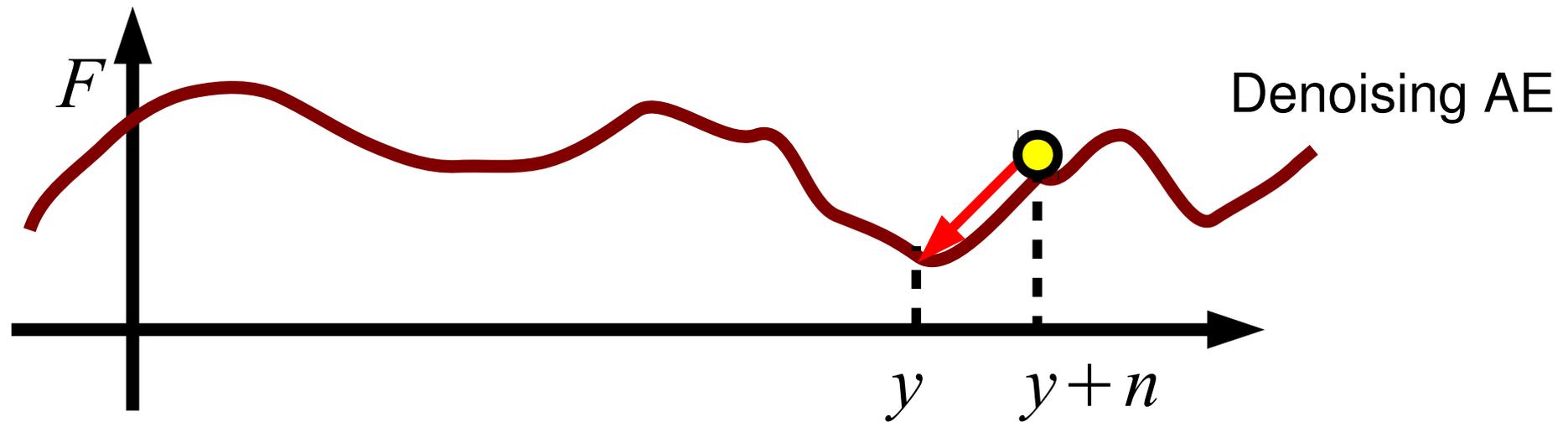
$$\hat{y} = W_2 \sigma(W_1 y + n), n \in N(0, I)$$

Training samples must be stable points (**local minima**) and **attractors**.

# Strategies to Shape E: #5

- Make the observed  $y$  an attractor state of some energy function. Need only to define a dynamics.

$$L = E(y)$$



E.g.:

$$L = \frac{1}{2} \|y - \hat{y}\|^2$$

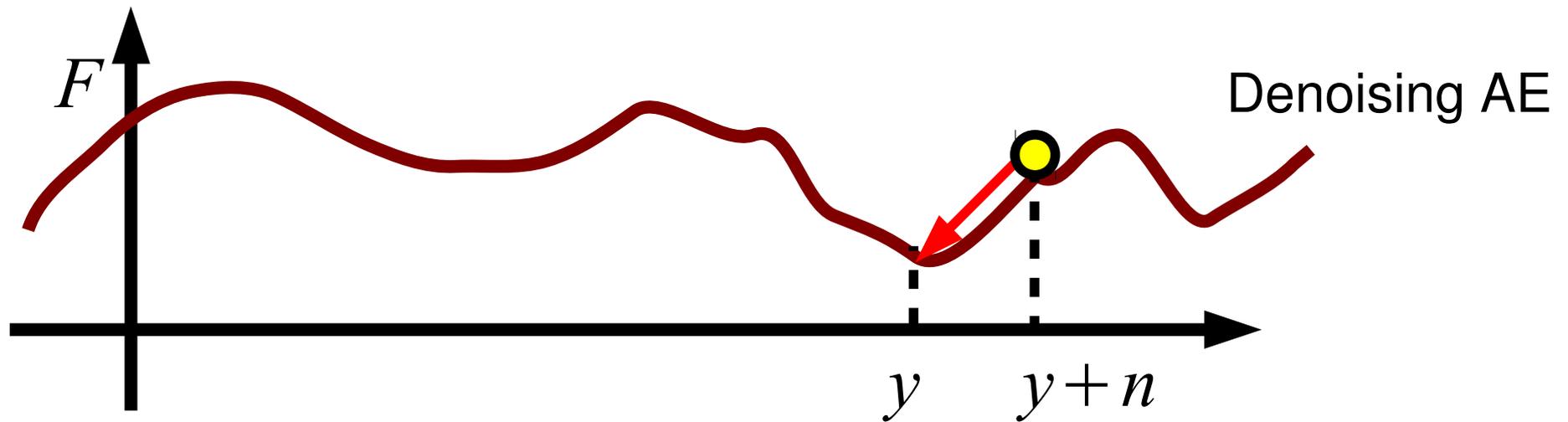
$$\hat{y} = W_2 \sigma(W_1 y + n), n \in N(0, I)$$

Training samples must be stable points (**local minima**) and **attractors**.

# Strategies to Shape E: #5

- Make the observed  $y$  an attractor state of some energy function. Need only to define a dynamics.

$$L = E(y)$$



## PROS

Efficient in high dimensional spaces.

## CONS

Need to pick noise distribution. May need to pick many noisy points.

# Loss: summary

- Goal of loss: make energy lower for correct answer.
- Different losses choose differently how to “pull-up”
  - Pull-up all points
  - Pull up one or a few points only
  - Make observations minima & increase curvature
  - Add global constraints/penalties on internal states
  - Define a dynamics with observations at the minima
- Choice of loss depends on desired landscape, computational budget, domain of input (discrete/continuous), task, etc.

# Final Notes on EBMs

- EBMs apply to any predictor (shallow & deep).
- EBMs subsume graphical models (e.g., use strategy #1 or #2 to “pull-up”).
- EBM is general framework to design good loss functions.

# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- **Examples:**
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - Unsupervised learning: sparse coding
  - Unsupervised learning: gated MRF
- Other examples
- Practical tricks

Recurrent  
Neural Net

CNN

SHALLOW

Boosting

Perceptron

SVM

Neural Net

SUPERVISED

UNSUPERVISED

Deep (sparse/denoising)  
Autoencoder

Sparse Coding

Autoencoder  
Neural Net

PROBABILISTIC

$\Sigma\Pi$

DBN

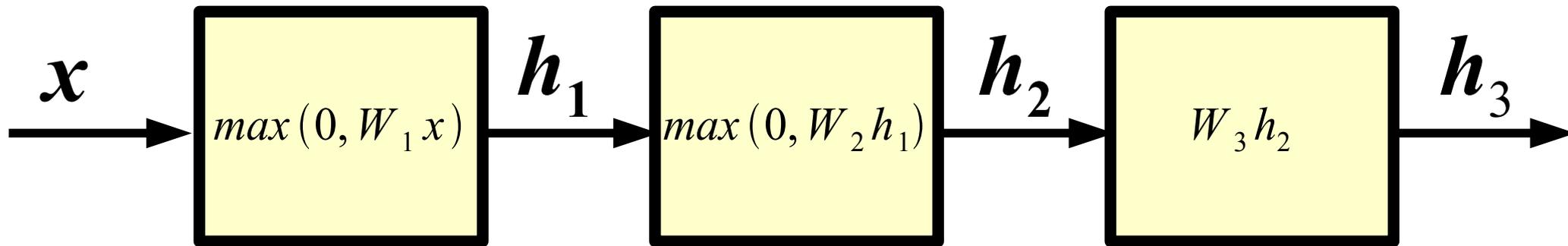
Restricted BM

GMM

BayesNP

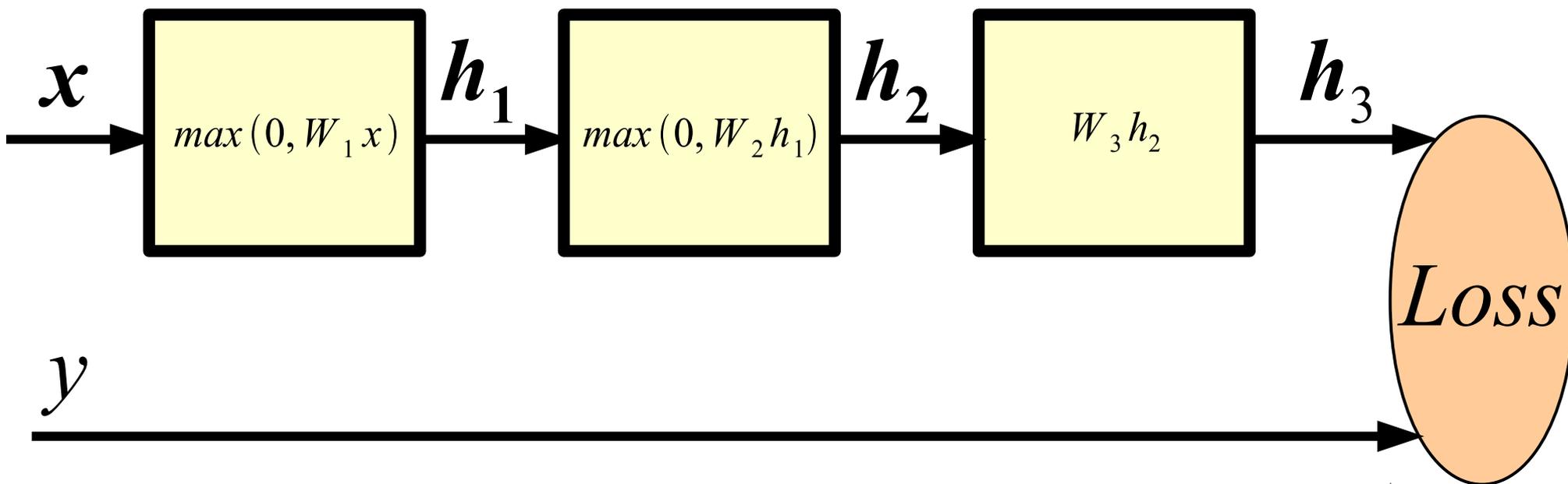
DEEP

# Neural Nets



**NOTE:** In practice, any (a.e. differentiable) non-linear transformation can be used.

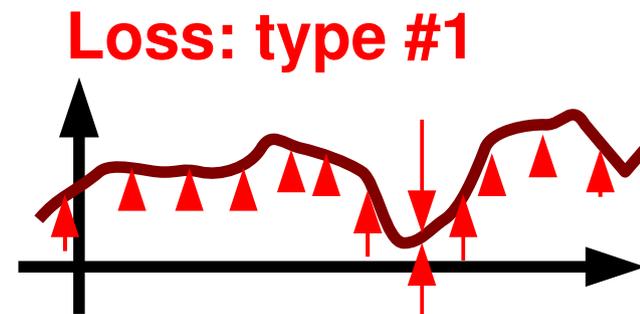
# Loss



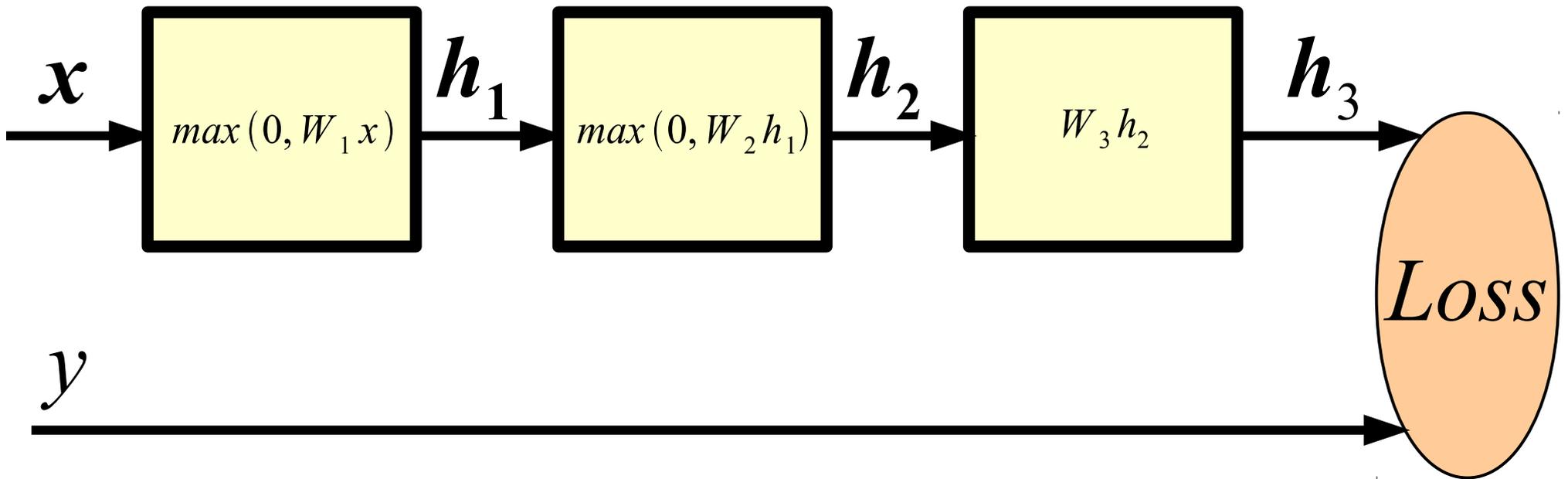
$$E(y; x) = -y h_3$$

$$L(E(y; x)) = \log(1 + \exp(E(y; x)))$$

$$\{W_1^*, W_2^*, W_3^*\} = \arg \min_{W_1, W_2, W_3} L(E(y; x))$$



# Loss

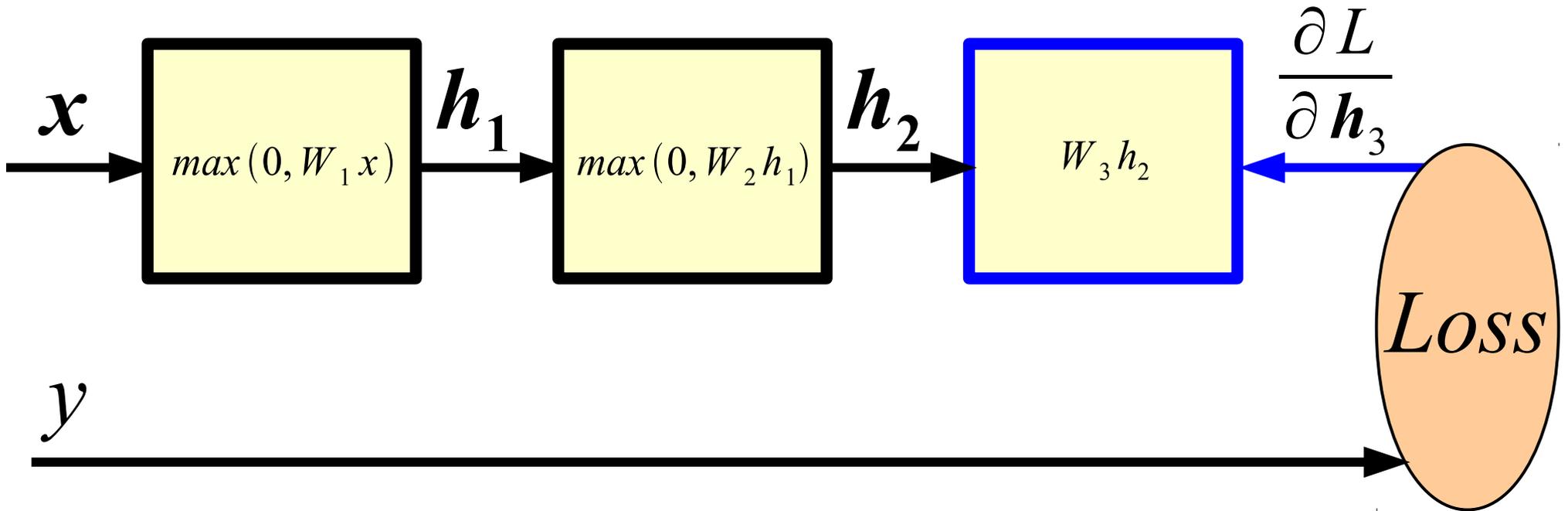


**Q.:** how to tune the parameters to decrease the loss?

If loss is (a.e.) differentiable we can compute gradients.

We can use chain-rule, a.k.a. **back-propagation**, to compute the gradients w.r.t. parameters at the lower layers.

# Backward Propagation

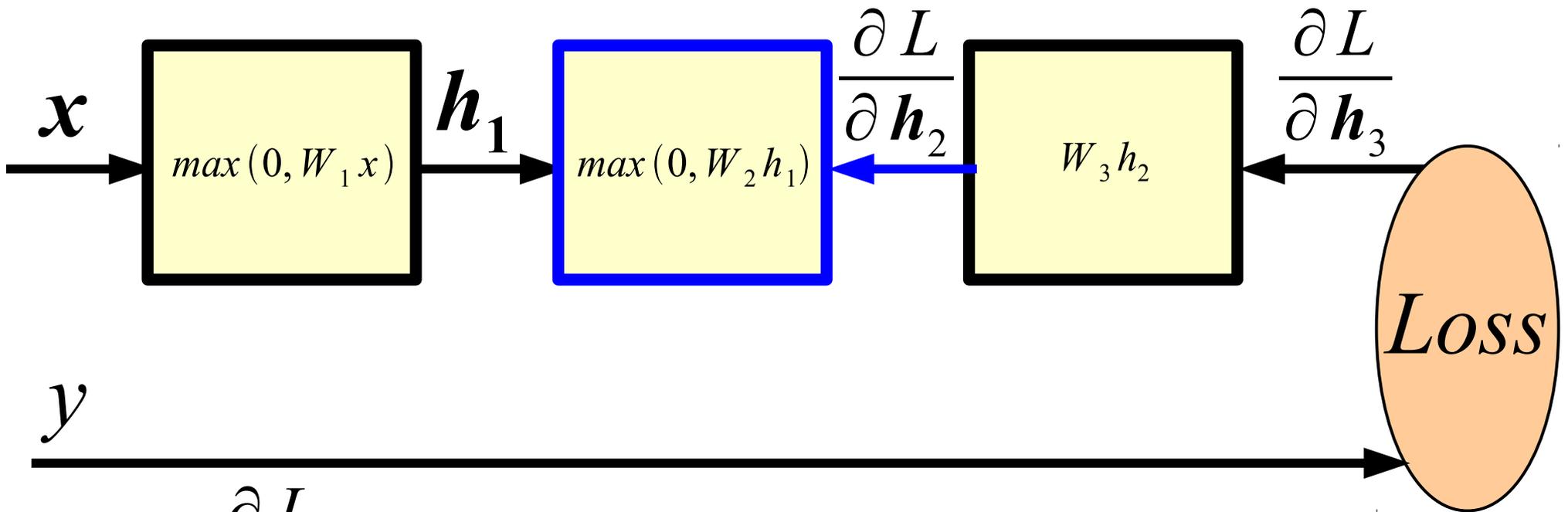


Given  $\frac{\partial L}{\partial h_3}$  and assuming we can easily compute the Jacobian of each module, we have:

$$\frac{\partial L}{\partial W_3} = \frac{\partial L}{\partial h_3} \frac{\partial h_3}{\partial W_3} \qquad \frac{\partial L}{\partial h_2} = \frac{\partial L}{\partial h_3} \frac{\partial h_3}{\partial h_2}$$

$$\frac{\partial L}{\partial W_3} = (\sigma(h_3) - y) h_2^T \qquad \frac{\partial L}{\partial h_2} = W_3^T (\sigma(h_3) - y)$$

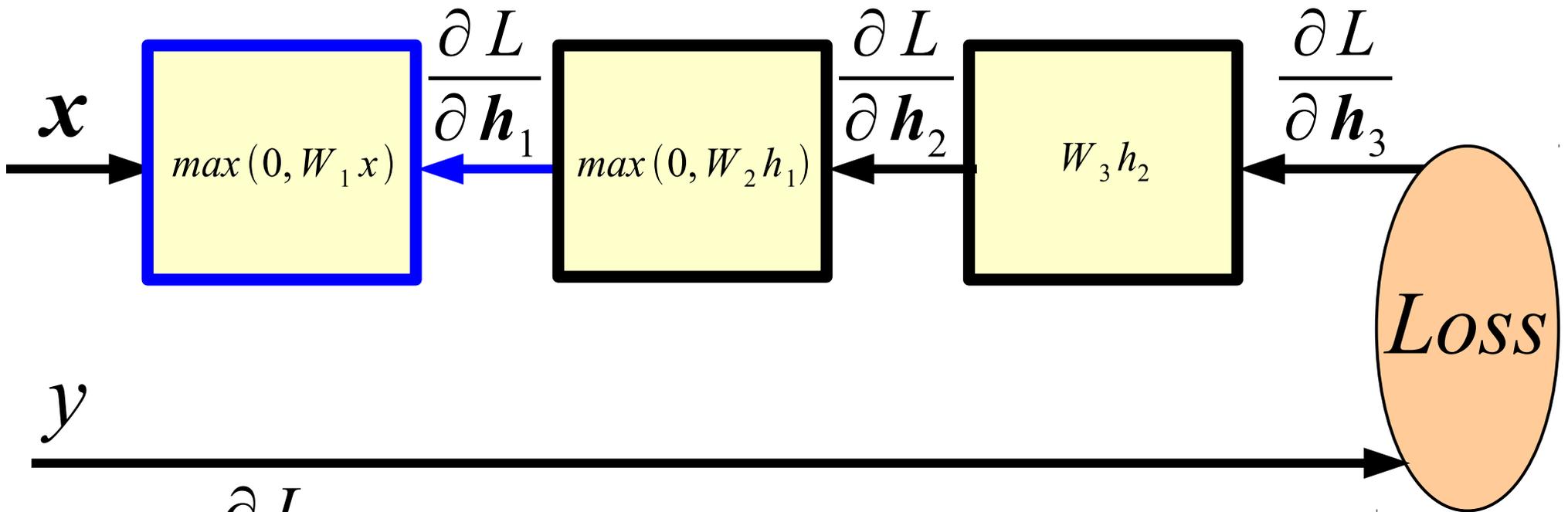
# Backward Propagation



Given  $\frac{\partial L}{\partial h_2}$  we can compute now:

$$\frac{\partial L}{\partial W_2} = \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial W_2} \qquad \frac{\partial L}{\partial h_1} = \frac{\partial L}{\partial h_2} \frac{\partial h_2}{\partial h_1}$$

# Backward Propagation



Given  $\frac{\partial L}{\partial h_1}$  we can compute now:

$$\frac{\partial L}{\partial W_1} = \frac{\partial L}{\partial h_1} \frac{\partial h_1}{\partial W_1}$$

# Optimization

**Stochastic Gradient Descent (on mini-batches):**

$$\theta \leftarrow \theta - \eta \frac{\partial L}{\partial \theta}, \eta \in R$$

**Stochastic Gradient Descent with Momentum:**

$$\theta \leftarrow \theta - \eta \Delta$$

$$\Delta \leftarrow 0.9 \Delta + \frac{\partial L}{\partial \theta}$$

# Toy Code (Matlab): Neural Net Trainer

```
% F-PROP
```

```
for i = 1 : nr_layers - 1
    [h{i} jac{i}] = nonlinearity(W{i} * h{i-1} + b{i});
end
h{nr_layers-1} = W{nr_layers-1} * h{nr_layers-2} + b{nr_layers-1};
prediction = softmax(h{l-1});
```

```
% CROSS ENTROPY LOSS
```

```
loss = - sum(sum(log(prediction) .* target)) / batch_size;
```

```
% B-PROP
```

```
dh{l-1} = prediction - target;
for i = nr_layers - 1 : -1 : 1
    Wgrad{i} = dh{i} * h{i-1}';
    bgrad{i} = sum(dh{i}, 2);
    dh{i-1} = (W{i}' * dh{i}) .* jac{i-1};
end
```

```
% UPDATE
```

```
for i = 1 : nr_layers - 1
    W{i} = W{i} - (lr / batch_size) * Wgrad{i};
    b{i} = b{i} - (lr / batch_size) * bgrad{i};
end
```

Recurrent  
Neural Net

Loss: type #1



**CNN**

Neural Net

**SHALLOW**

Boosting

Perceptron

SVM

**SUPERVISED**

**UNSUPERVISED**

Deep (sparse/denoising)  
Autoencoder

Sparse Coding

Autoencoder  
Neural Net

**PROBABILISTIC**

$\Sigma\Pi$

DBN

Restricted BM

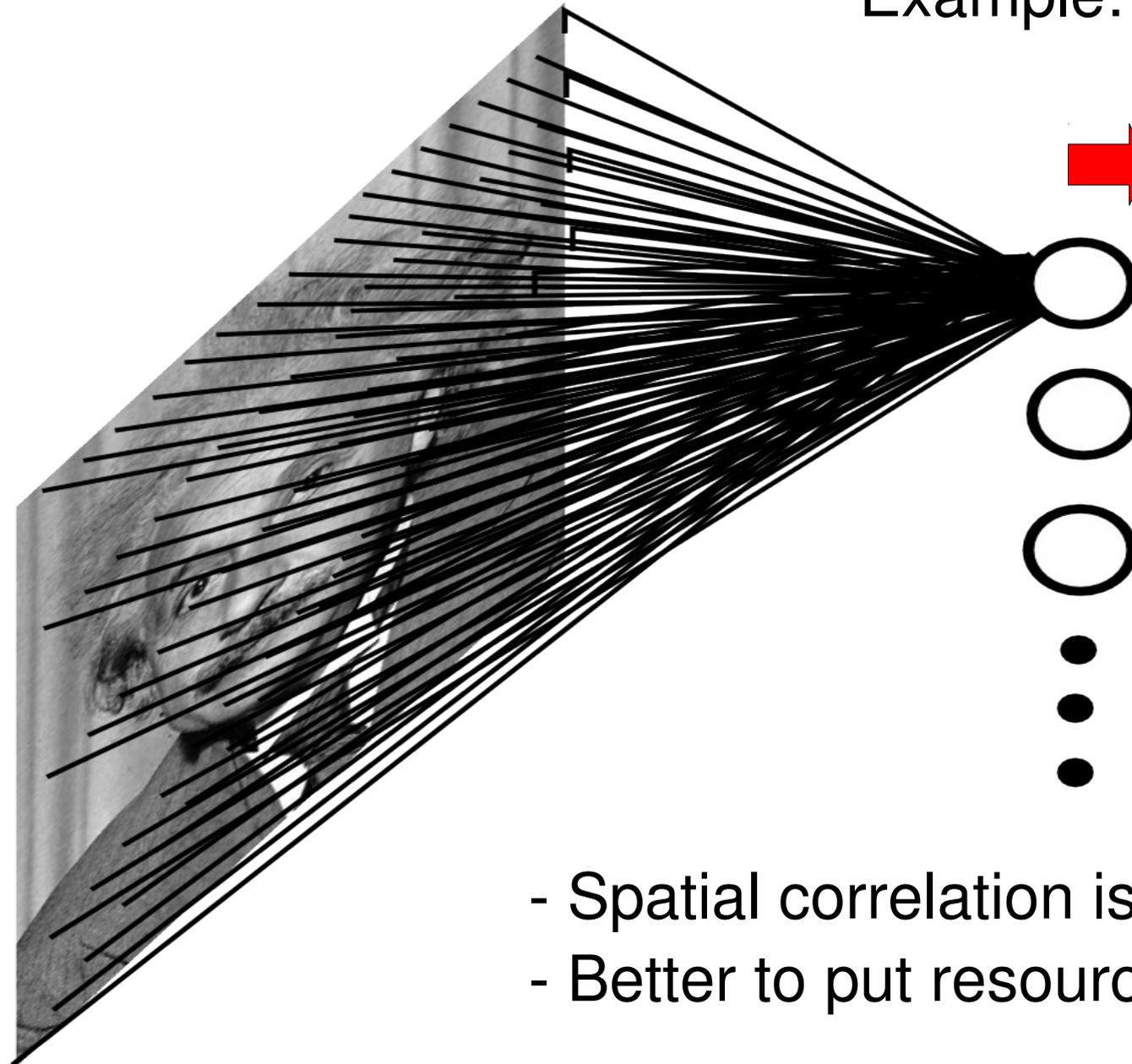
GMM

BayesNP

**DEEP**

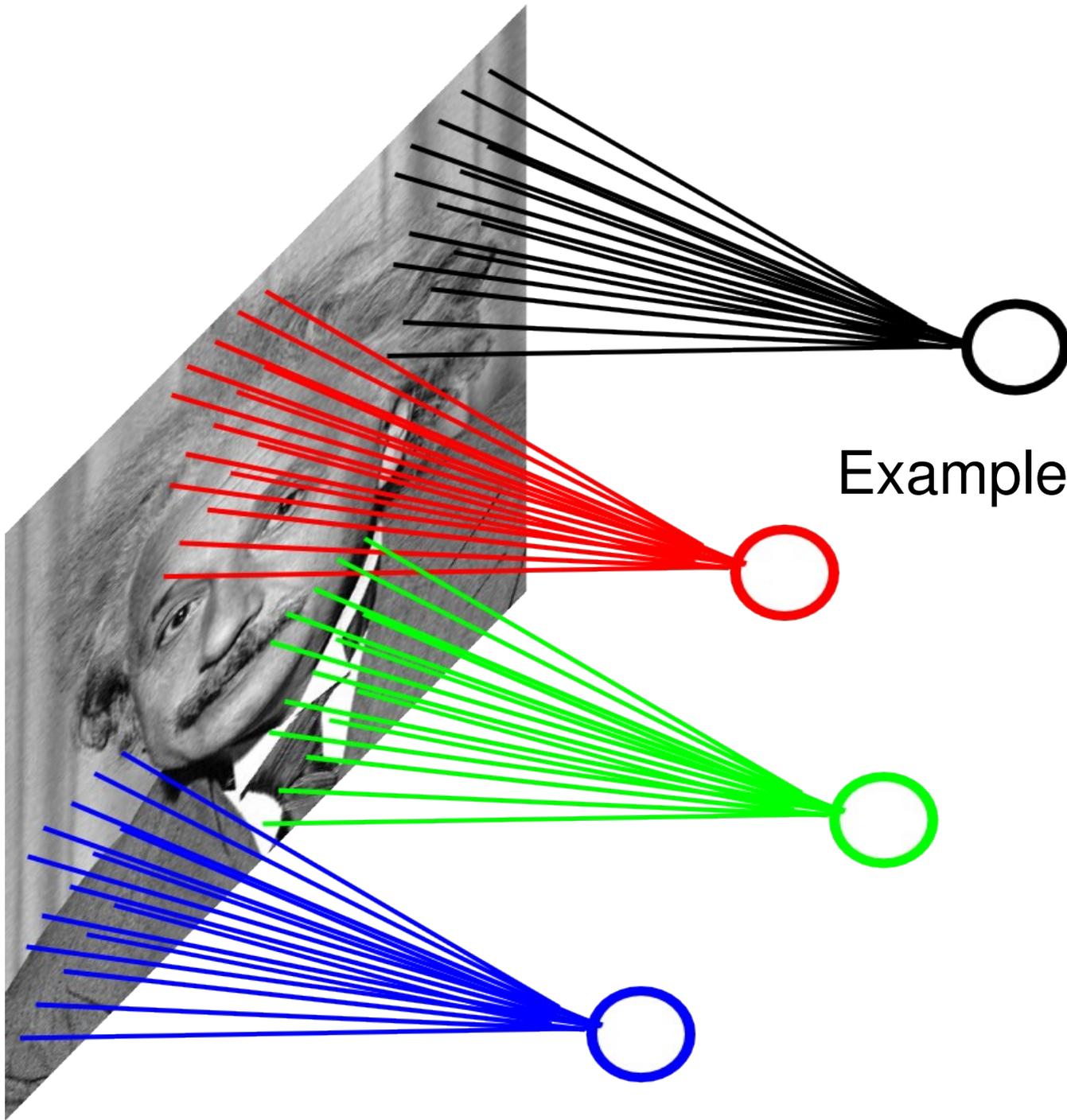
# FULLY CONNECTED NEURAL NET

Example: 200x200 image  
40K hidden units  
→ **~2B parameters!!!**



- Spatial correlation is local
- Better to put resources elsewhere!

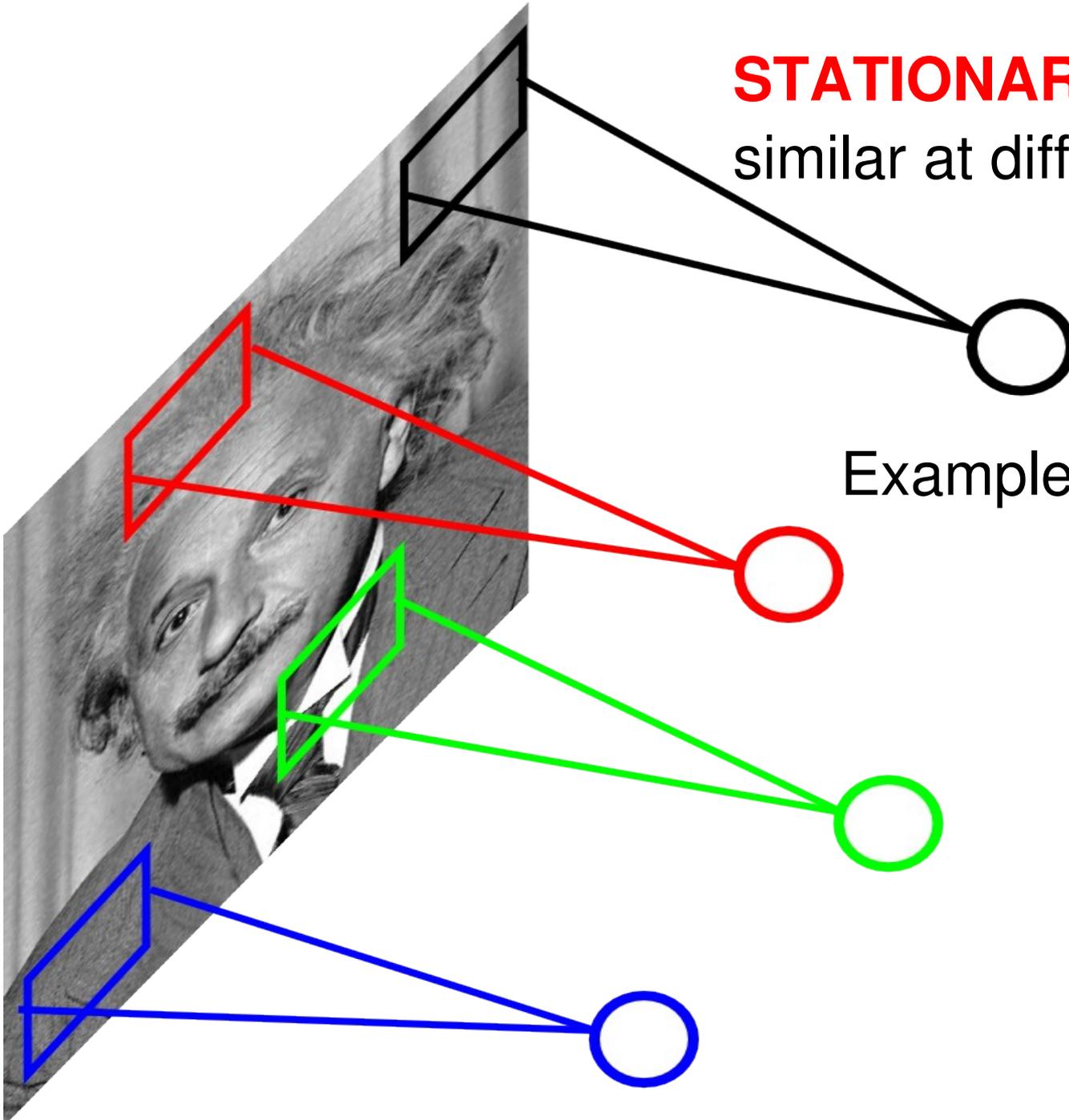
# LOCALLY CONNECTED NEURAL NET



Example: 200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

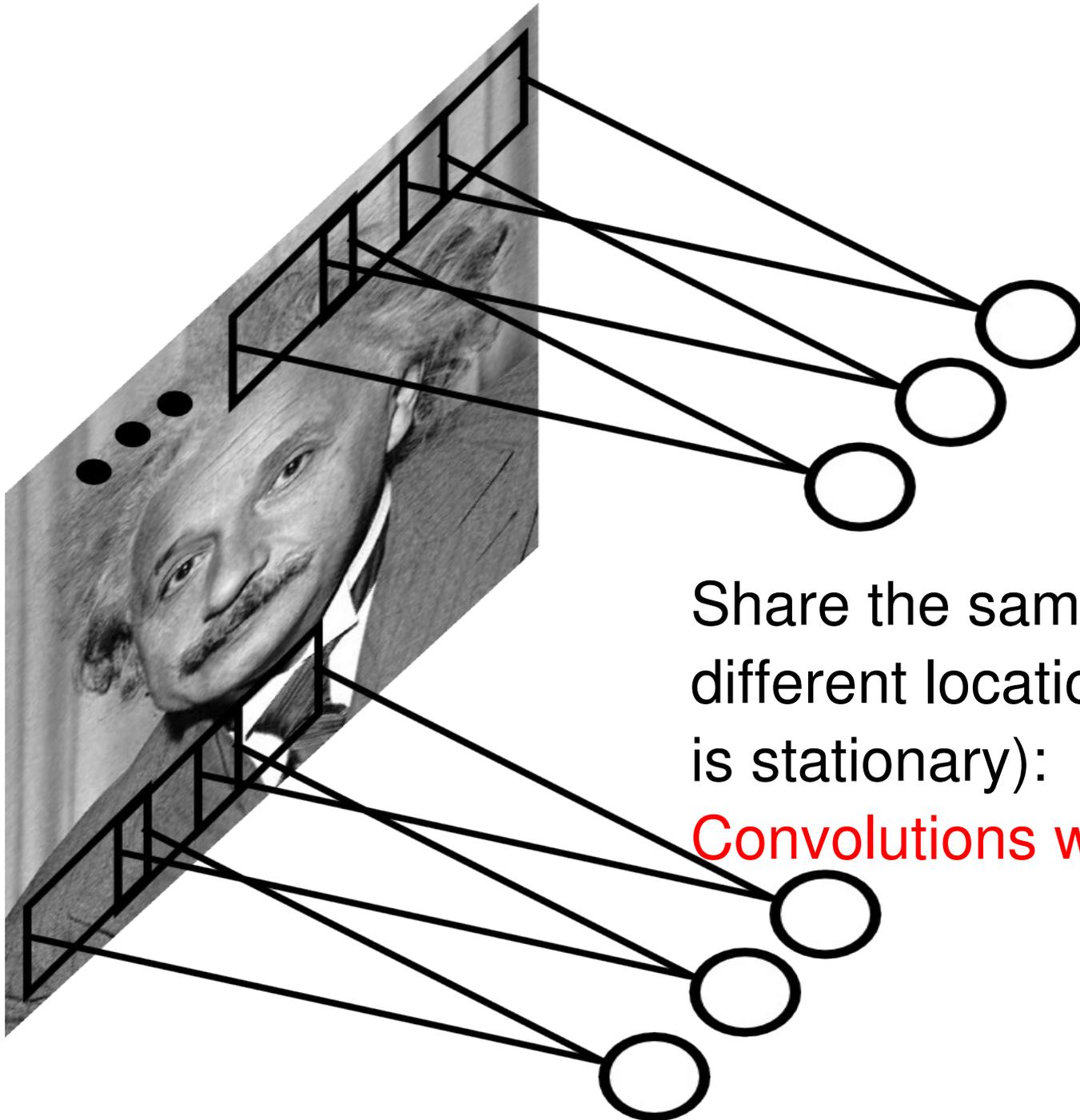
# LOCALLY CONNECTED NEURAL NET

**STATIONARITY?** Statistics are similar at different locations



Example: 200x200 image  
40K hidden units  
Filter size: 10x10  
4M parameters

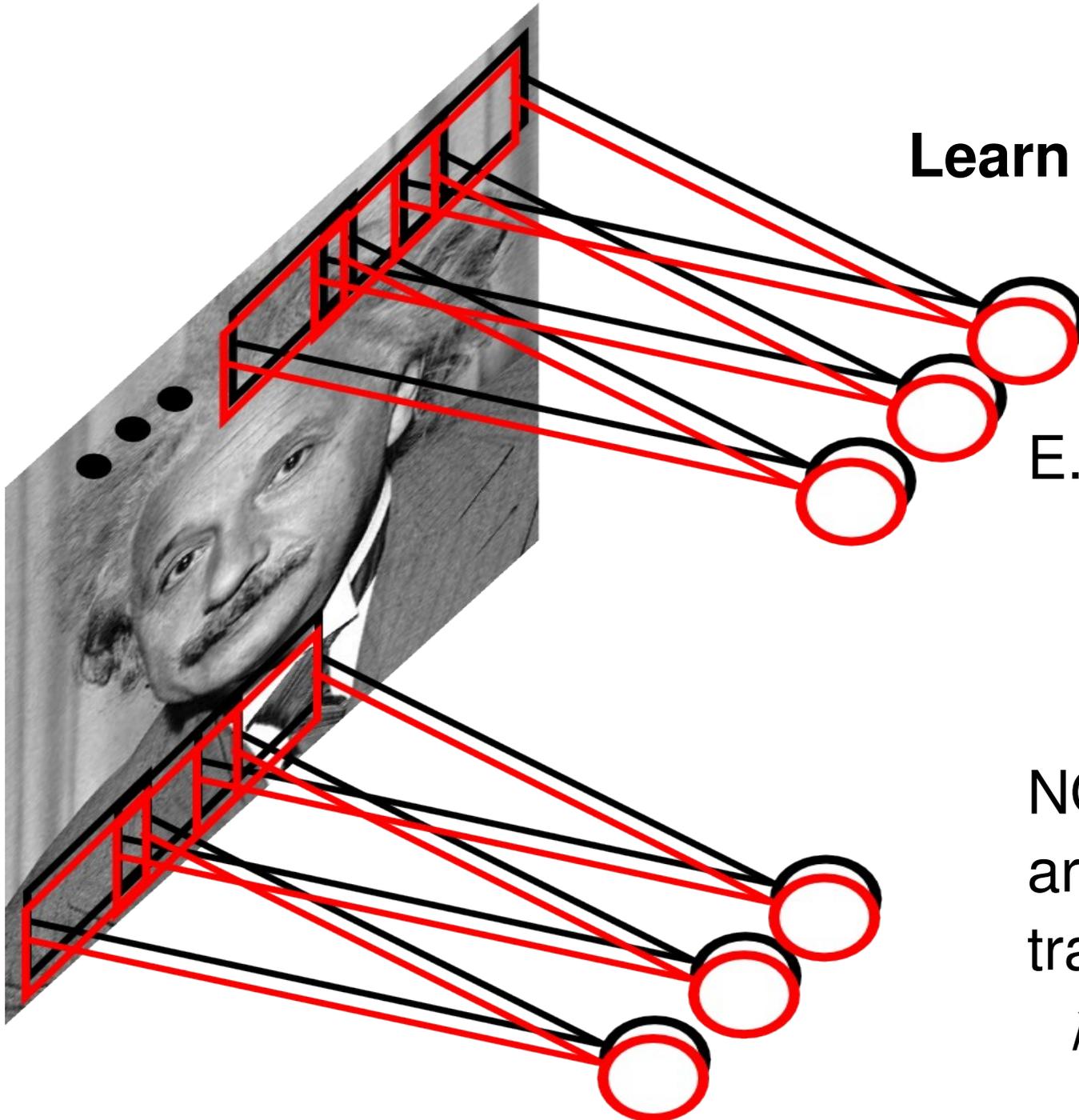
# CONVOLUTIONAL NET



Share the same parameters across different locations (assuming input is stationary):

**Convolutions with learned kernels**

# CONVOLUTIONAL NET



Learn **multiple filters**.

E.g.: 200x200 image  
100 Filters  
Filter size: 10x10  
10K parameters

NOTE: filter responses  
are non-linearly  
transformed:

$$h = \max(0, x * w)$$

76

# KEY IDEAS

A standard neural net applied to images:

- scales quadratically with the size of the input
- does not leverage stationarity

Solution:

- connect each hidden unit to a small patch of the input
- share the weight across hidden units

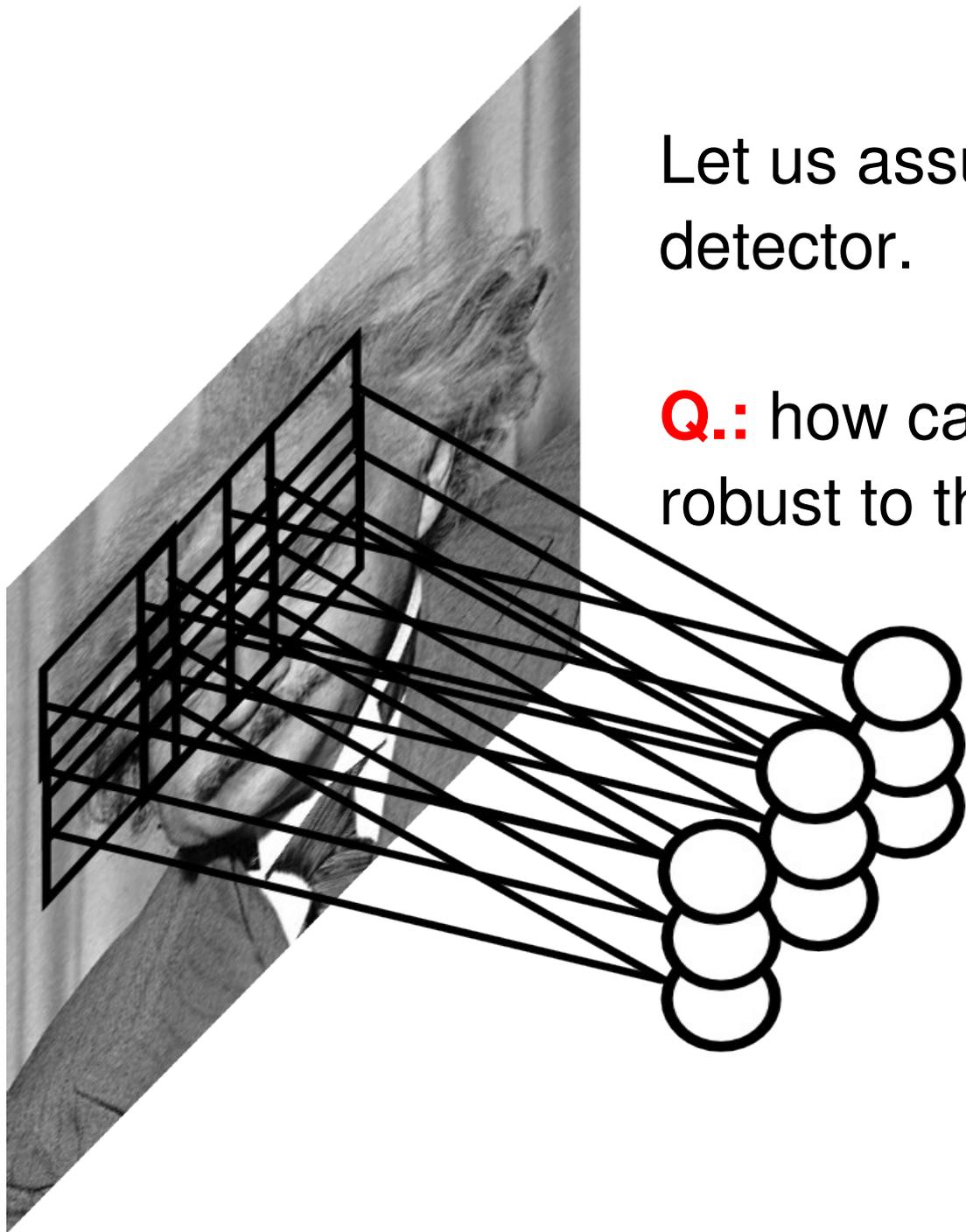
This is called: **convolutional layer.**

A network with convolutional layers is called **convolutional network.**

# POOLING

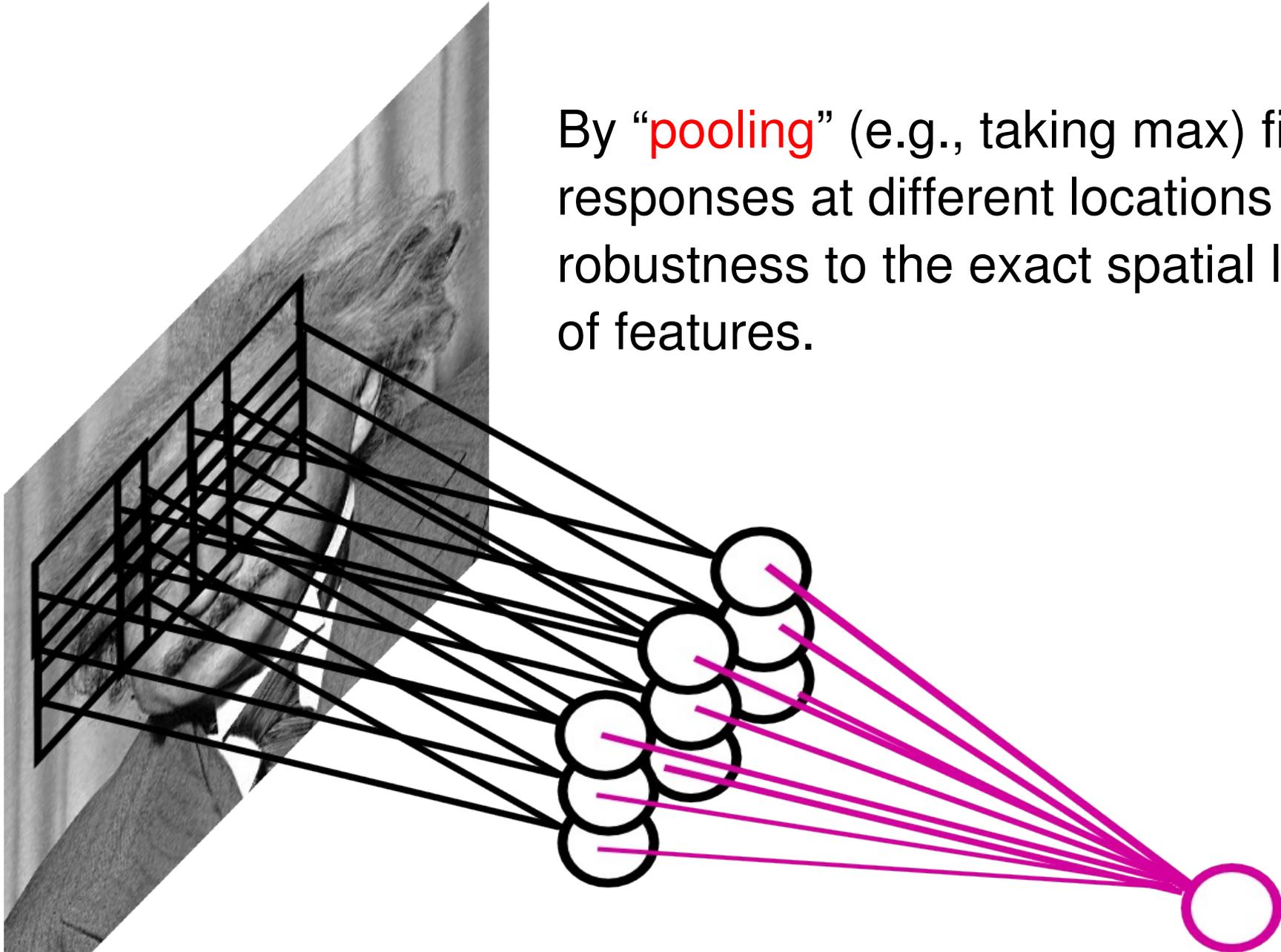
Let us assume filter is an “eye” detector.

**Q.:** how can we make the detection robust to the exact location of the eye?



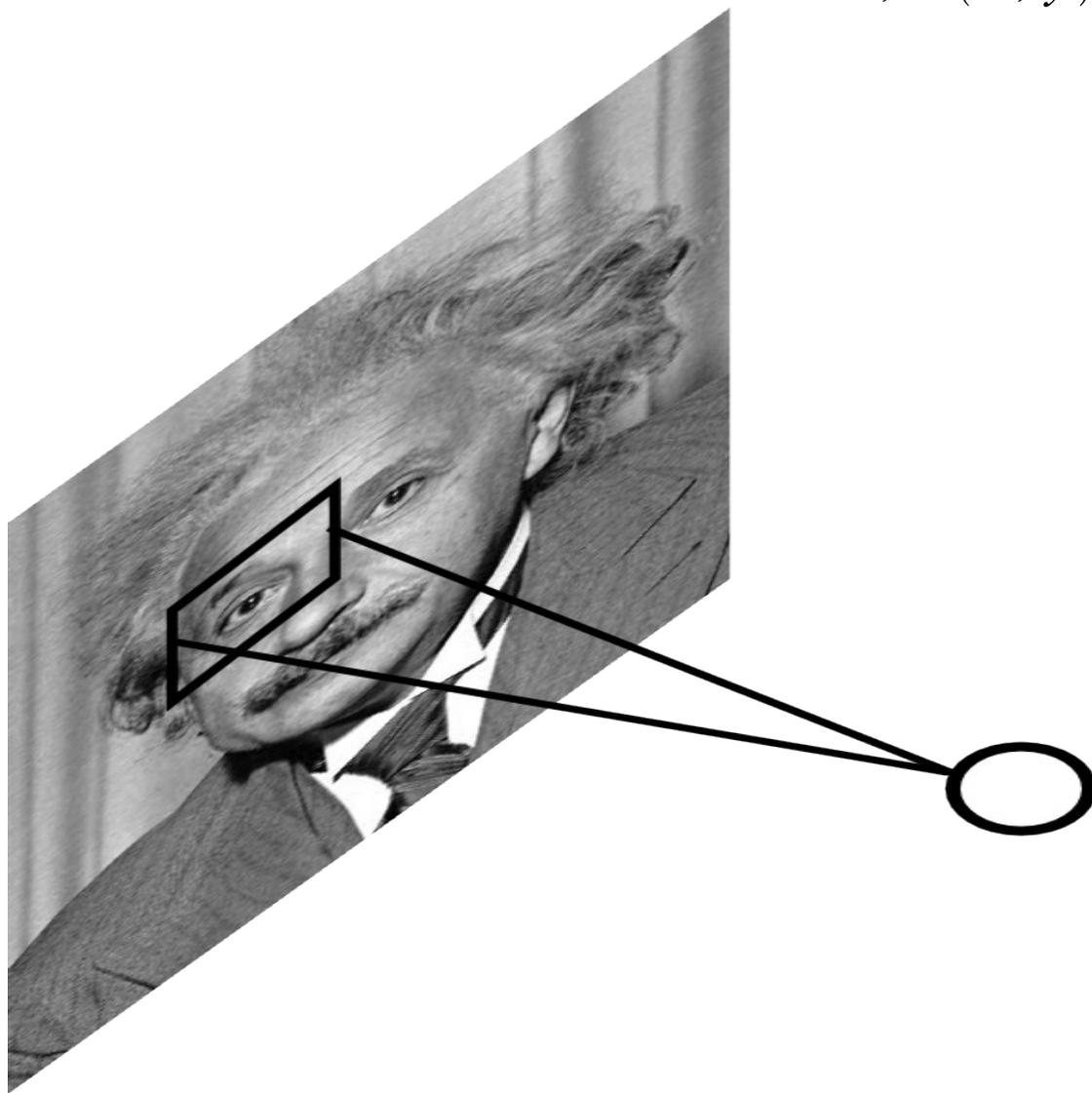
# POOLING

By “pooling” (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.



# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$



# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$



We want the same response.

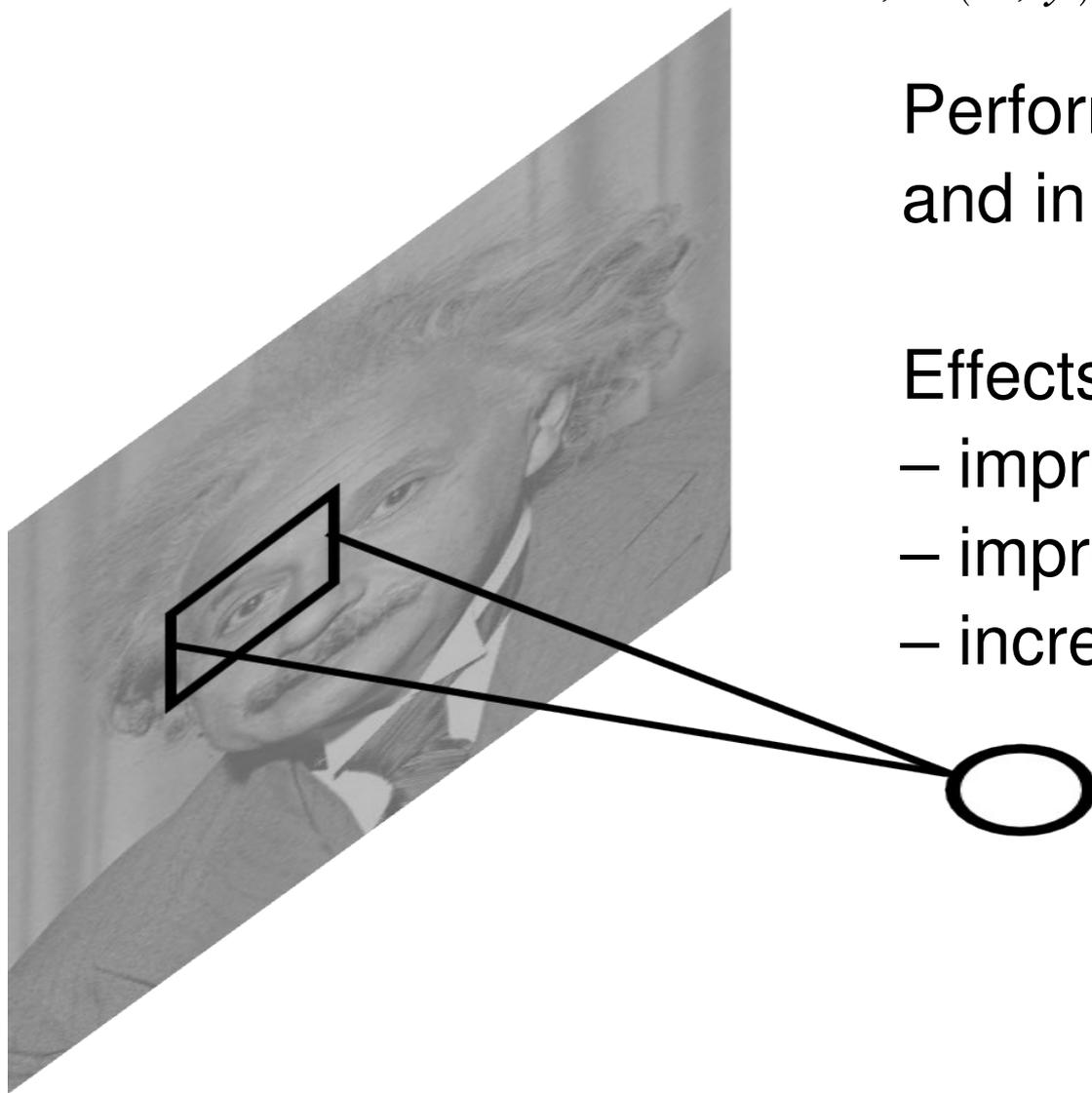
# LOCAL CONTRAST NORMALIZATION

$$h_{i+1,x,y} = \frac{h_{i,x,y} - m_{i,N(x,y)}}{\sigma_{i,N(x,y)}}$$

Performed also across features and in the higher layers.

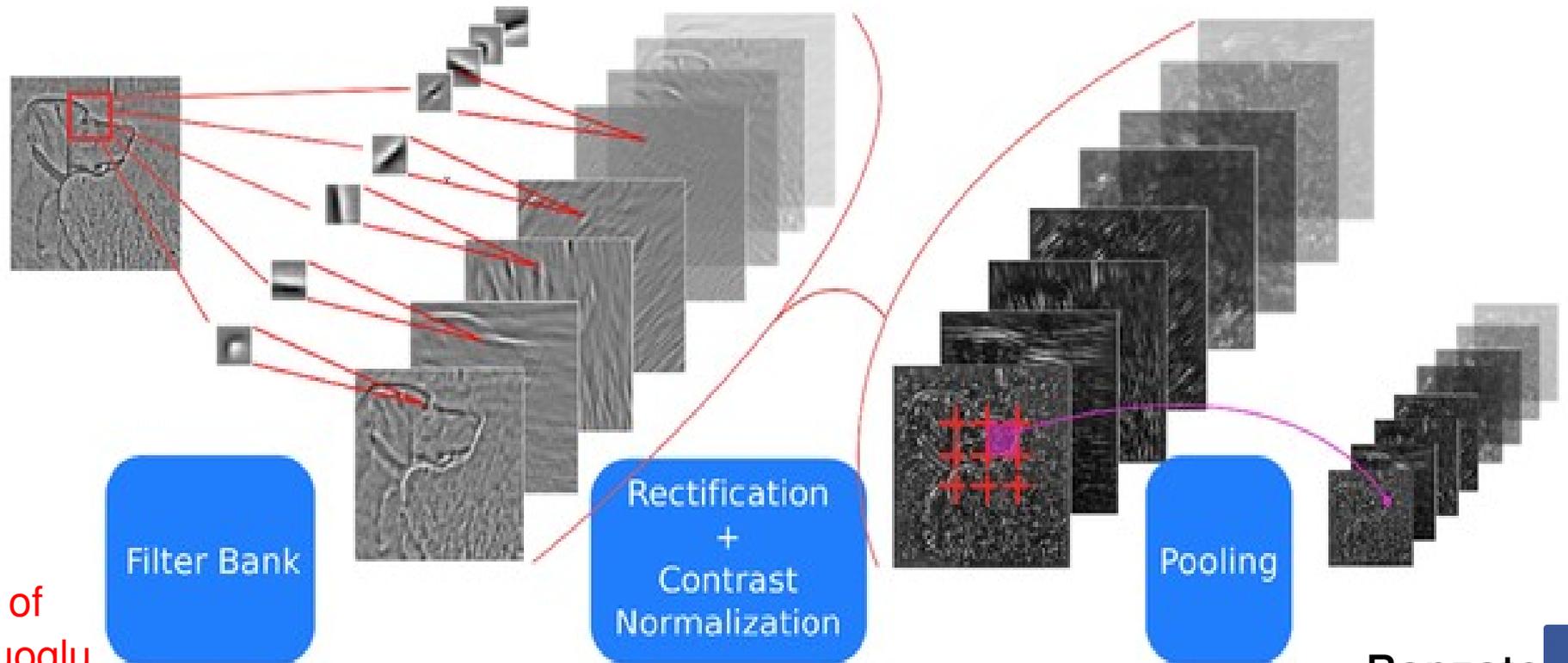
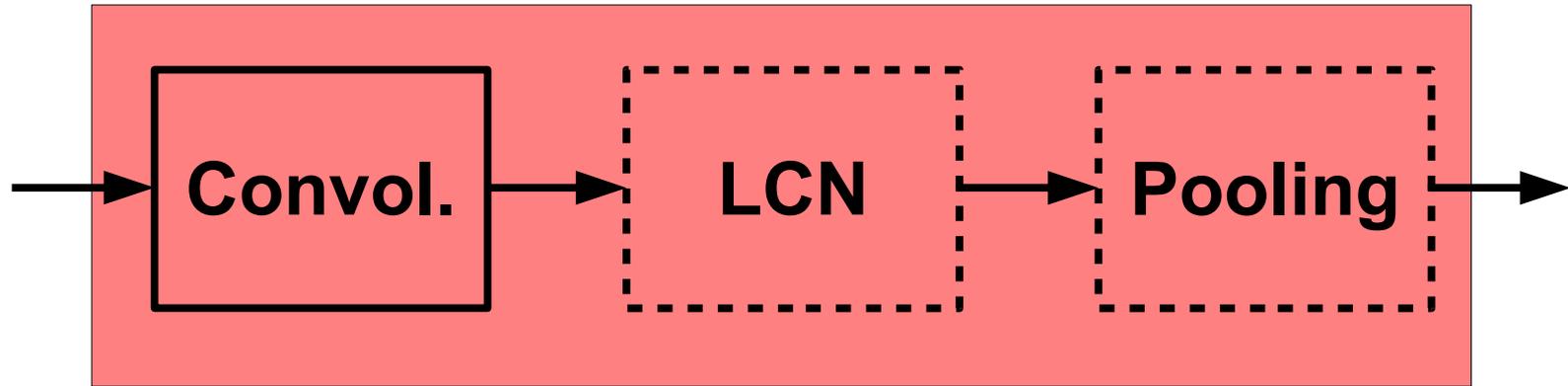
Effects:

- improves invariance
- improves optimization
- increases sparsity



# CONV NETS: TYPICAL ARCHITECTURE

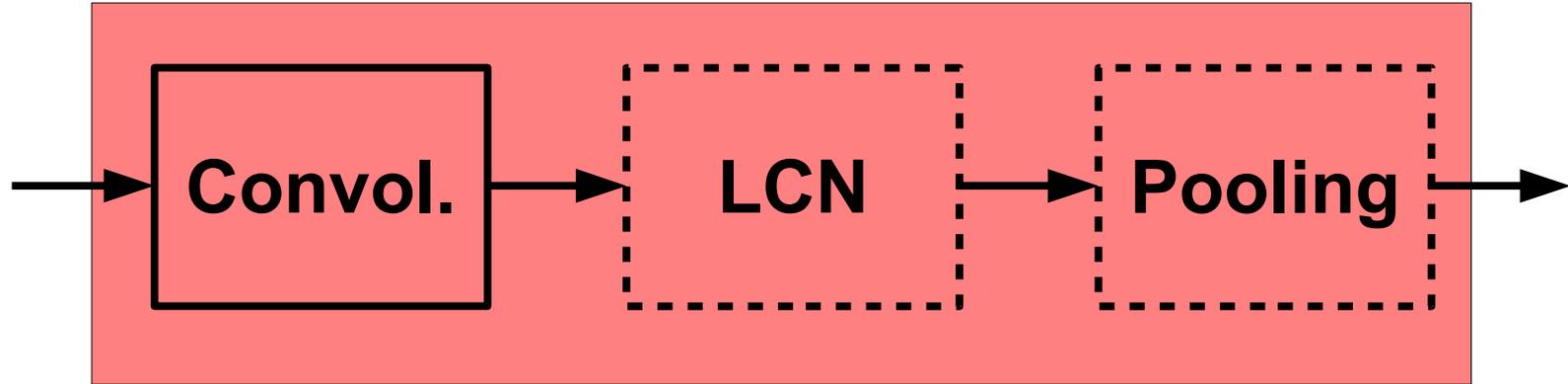
One stage (zoom)



courtesy of  
K. Kavukcuoglu

# CONV NETS: TYPICAL ARCHITECTURE

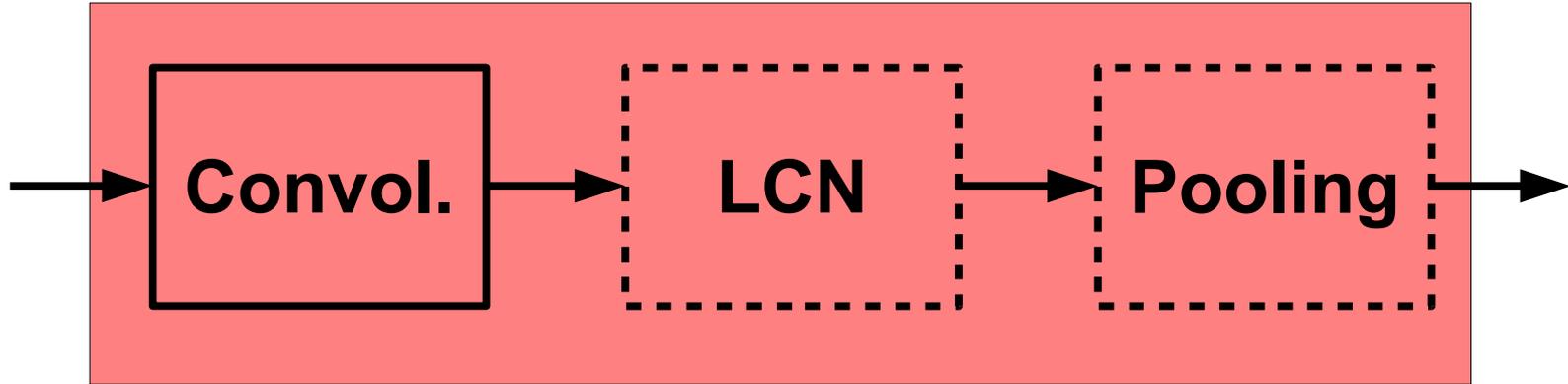
One stage (zoom)



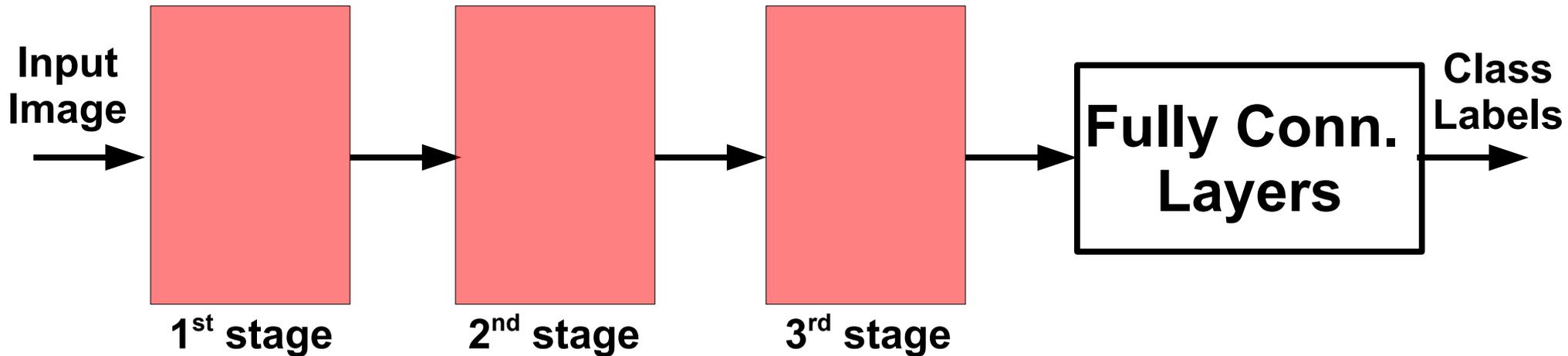
Conceptually similar to: SIFT, HoG, etc.

# CONV NETS: TYPICAL ARCHITECTURE

## One stage (zoom)

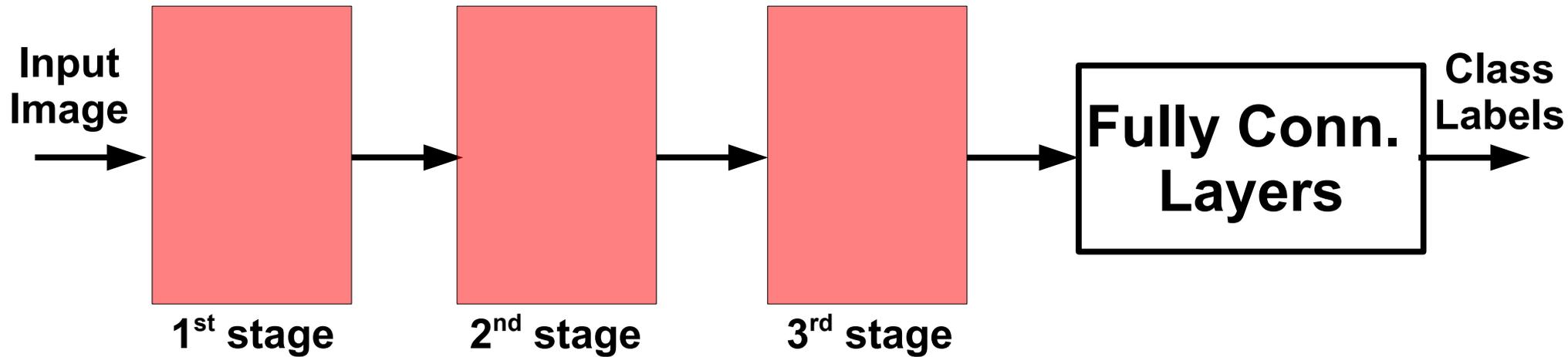


## Whole system



# CONV NETS: TYPICAL ARCHITECTURE

## Whole system



Conceptually similar to:

SIFT → K-Means → Pyramid Pooling → SVM

Lazebnik et al. "...Spatial Pyramid Matching..." CVPR 2006

SIFT → Fisher Vect. → Pooling → SVM

Sanchez et al. "Image classification with F.V.: Theory and practice" IJCV 2012

# CONV NETS: TRAINING

All layers are differentiable (a.e.).

We can use standard back-propagation.

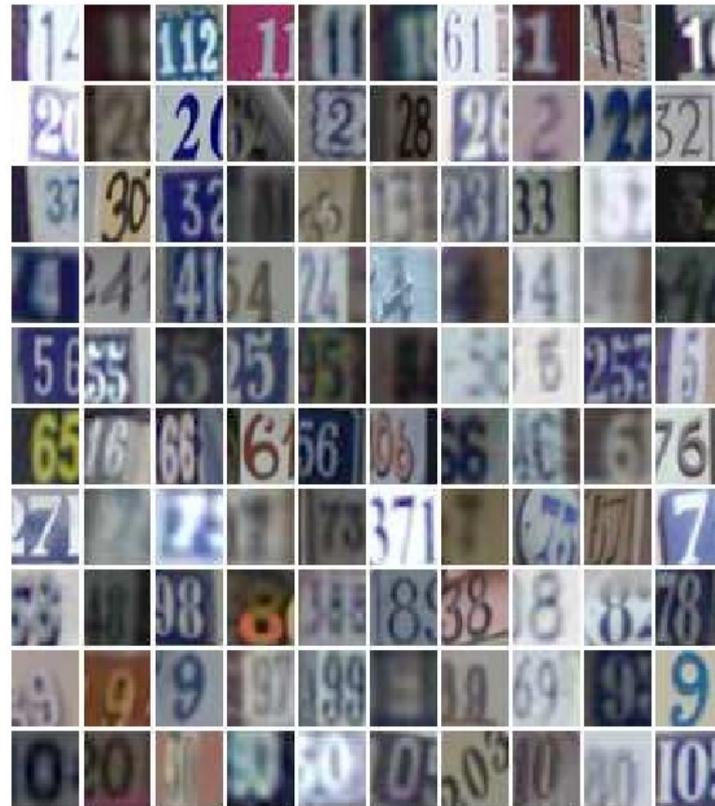
**Algorithm:**

**Given a small mini-batch**

- F-PROP**
- B-PROP**
- PARAMETER UPDATE**

# CONV NETS: EXAMPLES

- OCR / House number & Traffic sign classification



Ciresan et al. "MCDNN for image classification" CVPR 2012

Wan et al. "Regularization of neural networks using dropconnect" ICML 2013

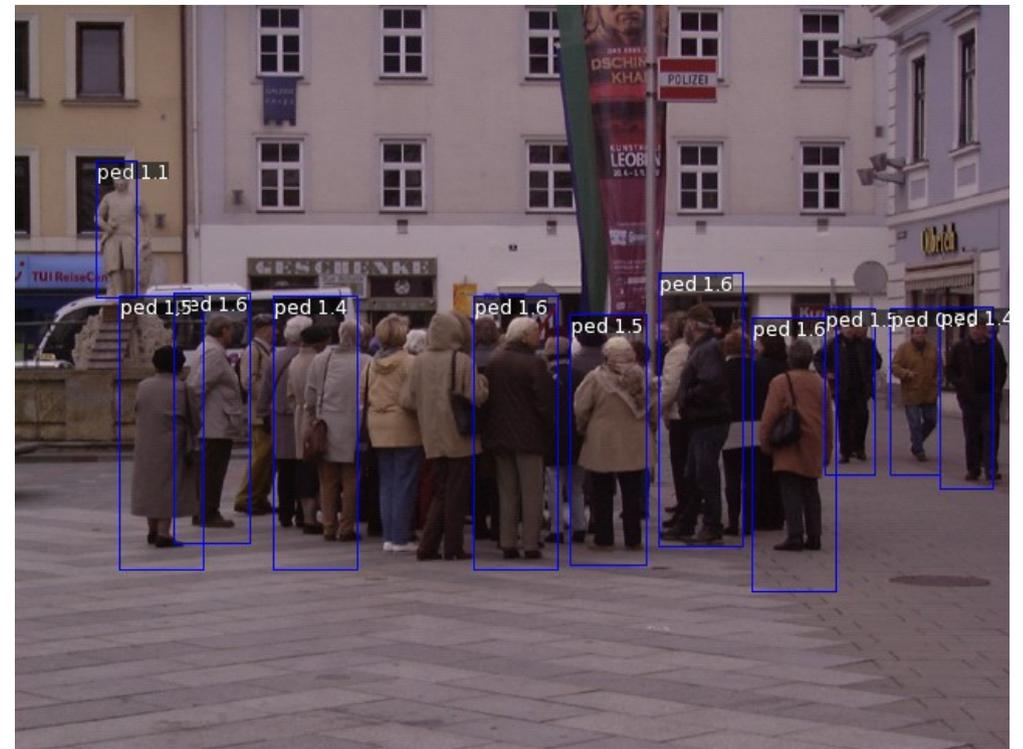
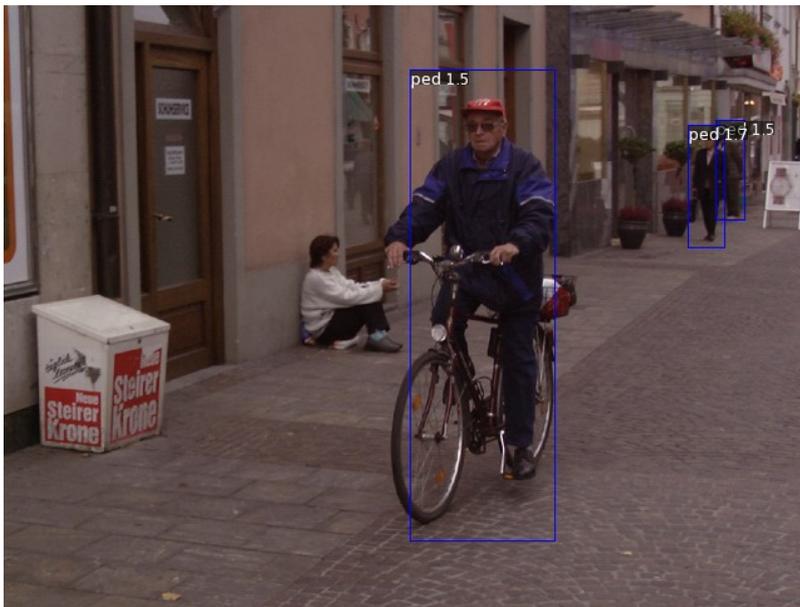
# CONV NETS: EXAMPLES

## - Texture classification



# CONV NETS: EXAMPLES

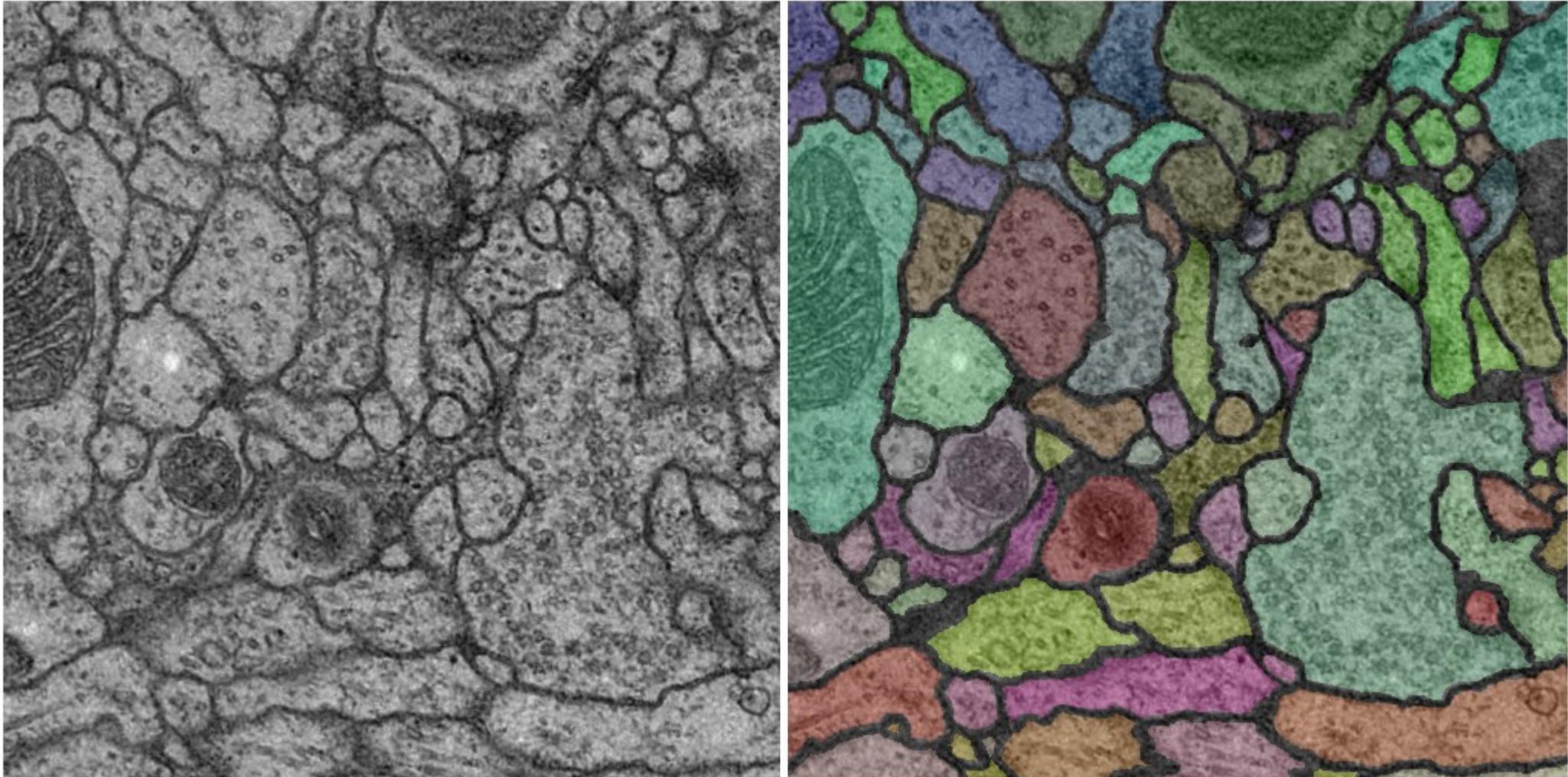
## - Pedestrian detection





# CONV NETS: EXAMPLES

## - Segmentation 3D volumetric images



Ciresan et al. "DNN segment neuronal membranes..." NIPS 2012

Turaga et al. "Maximin learning of image segmentation" NIPS 2009

# CONV NETS: EXAMPLES

- Action recognition from videos



# CONV NETS: EXAMPLES

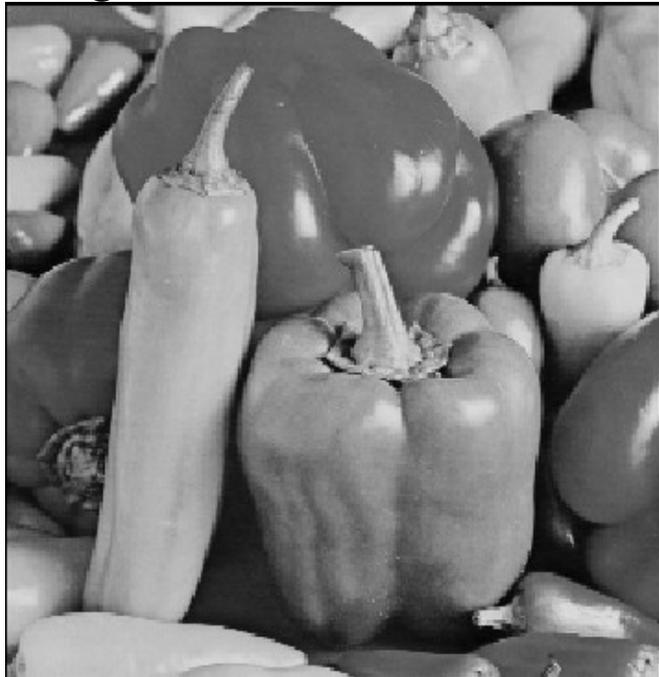
## - Robotics



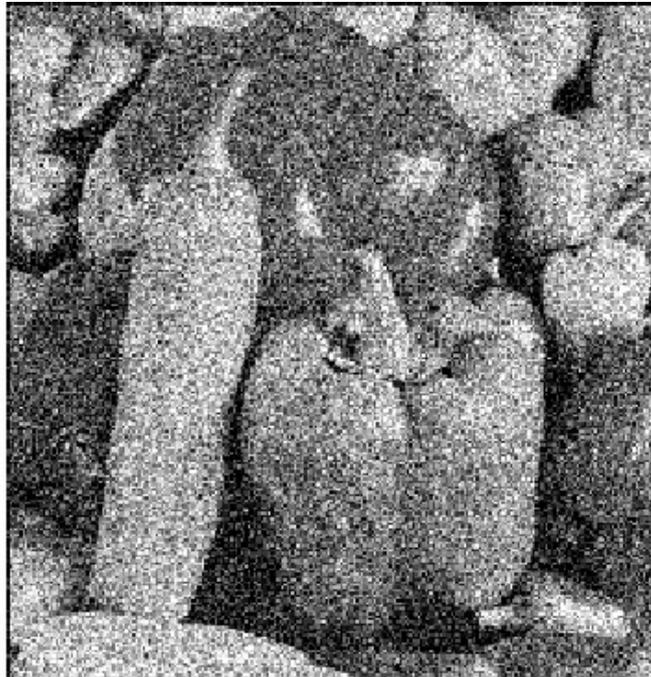
# CONV NETS: EXAMPLES

## - Denoising

original



noised

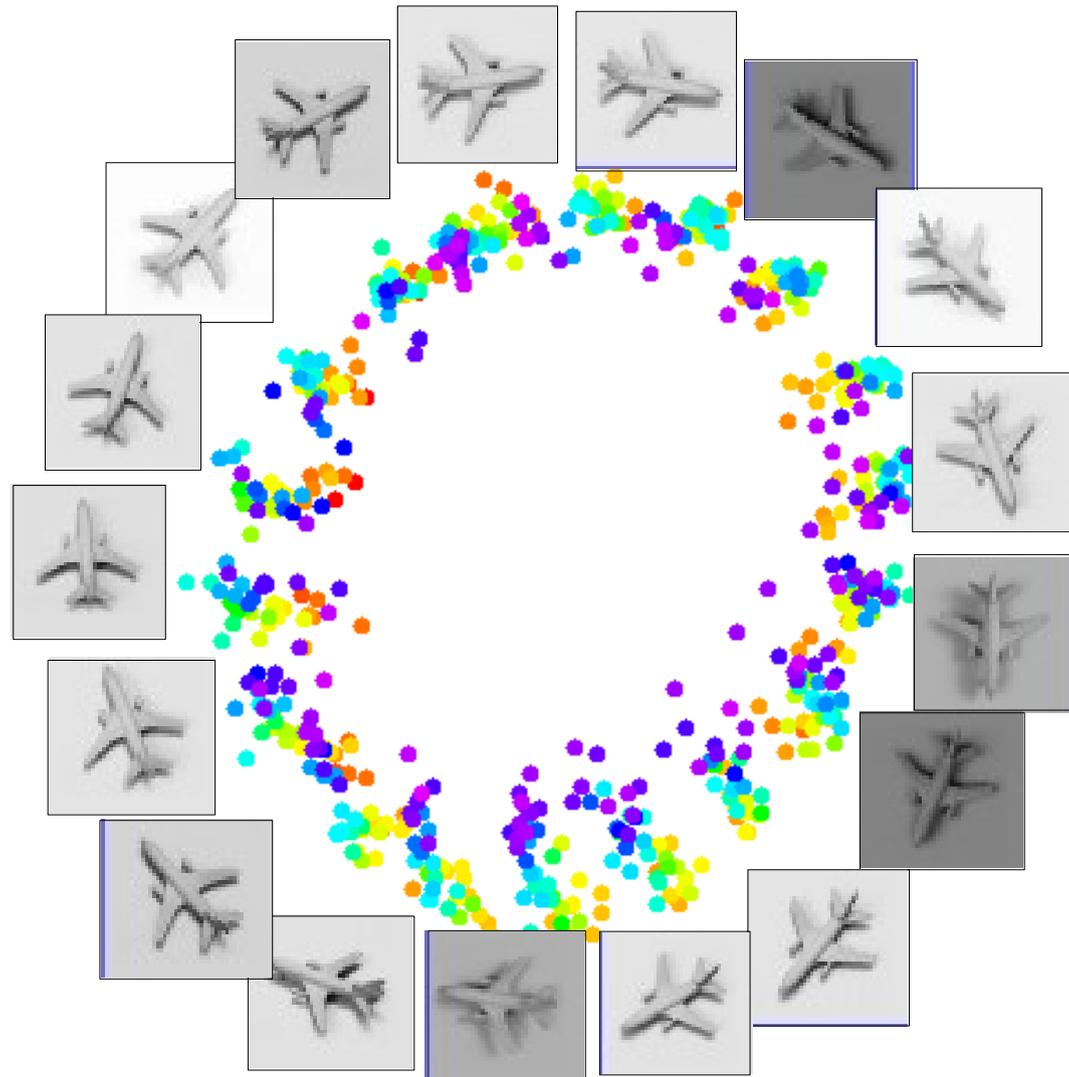


denoised



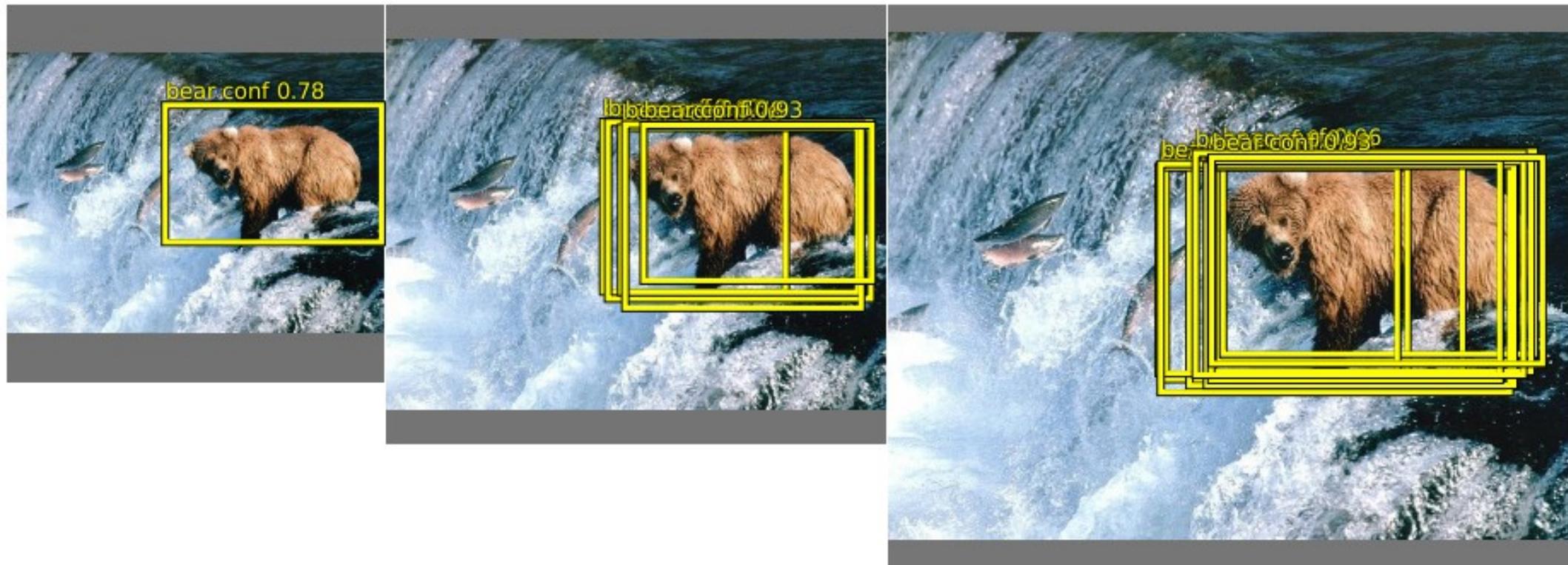
# CONV NETS: EXAMPLES

- Dimensionality reduction / learning embeddings



# CONV NETS: EXAMPLES

## - Object detection

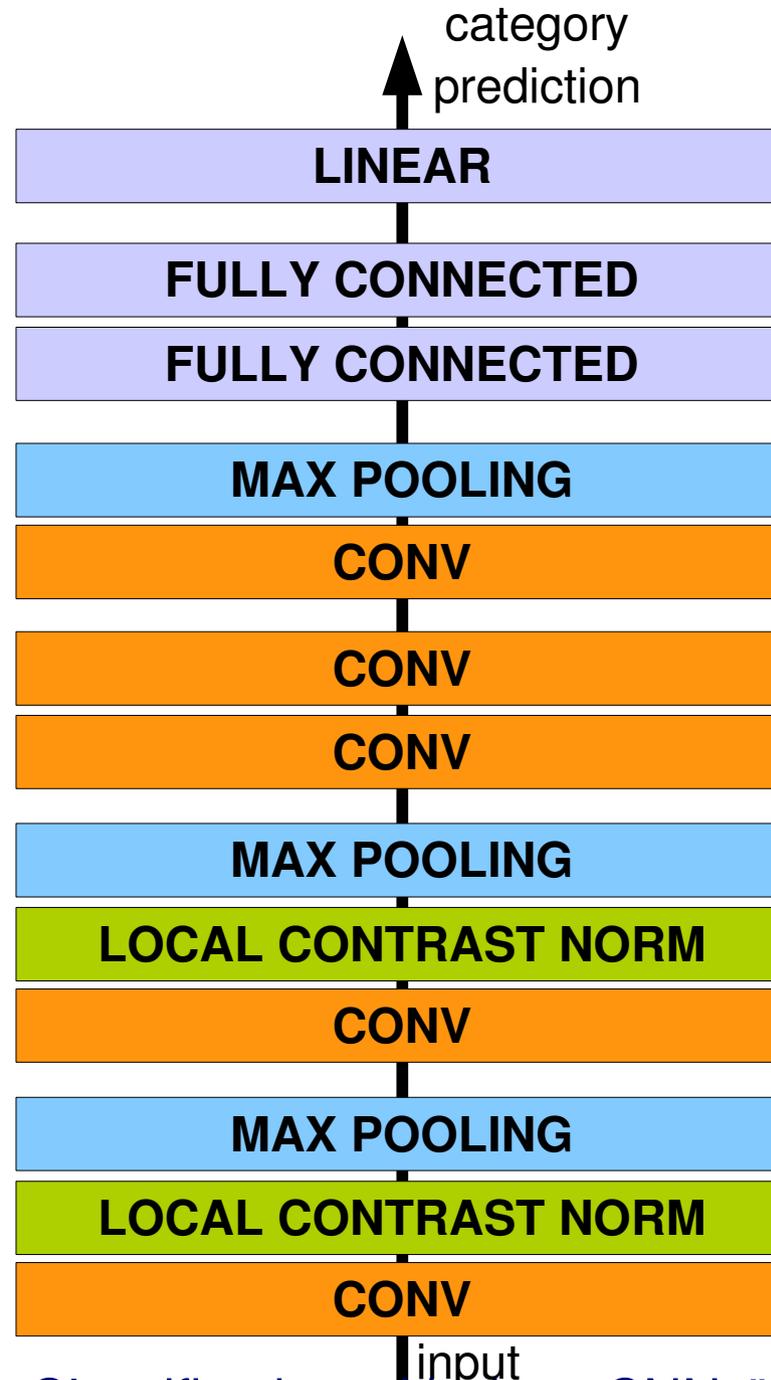


Sermanet et al. "OverFeat: Integrated recognition, localization, ..." arxiv 2013

Girshick et al. "Rich feature hierarchies for accurate object detection..." arxiv 2013 <sup>97</sup>

Szegedy et al. "DNN for object detection" NIPS 2013

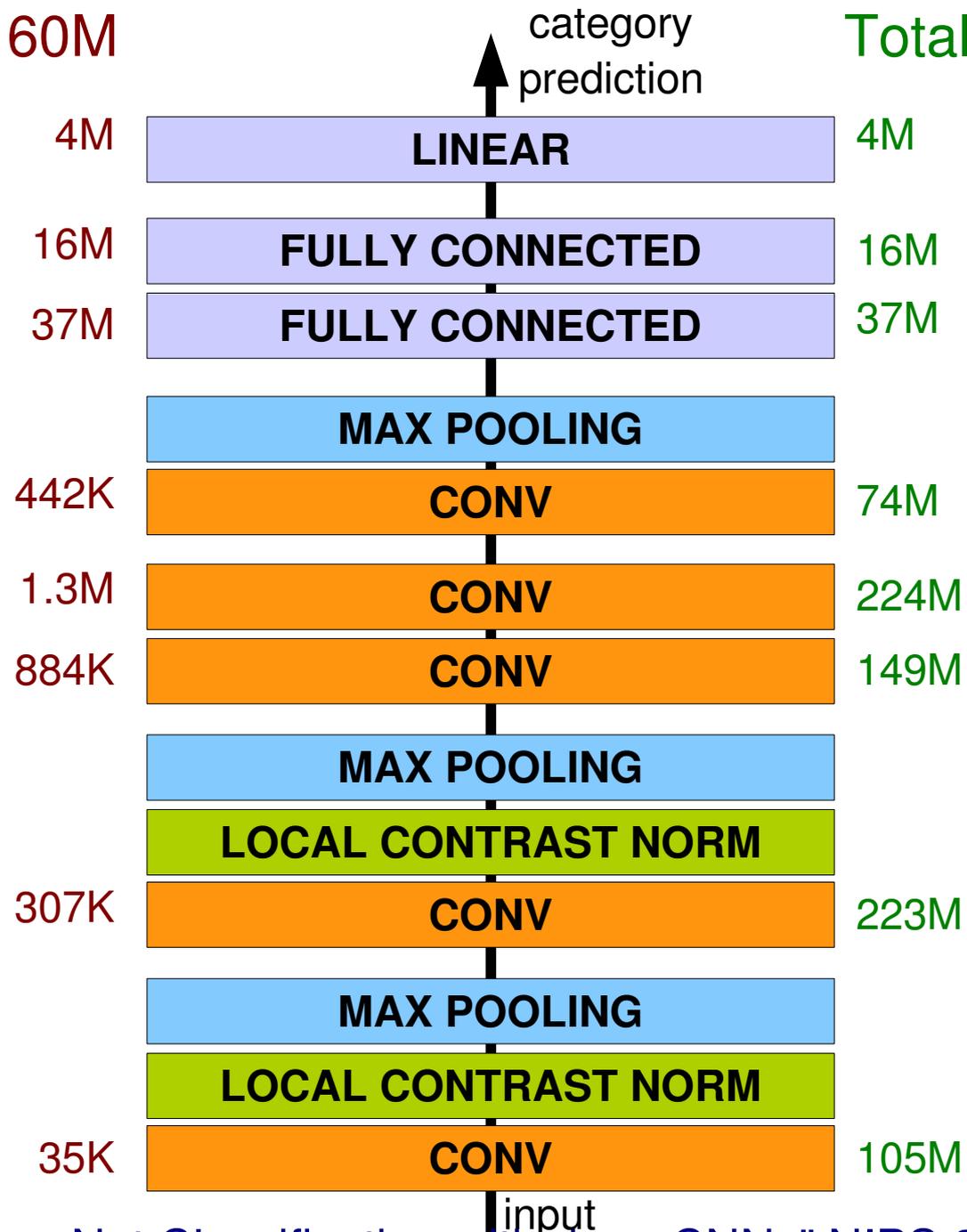
# Architecture for Classification



# Architecture for Classification

Total nr. params: 60M

Total nr. flops: 832M



# Optimization

## SGD with momentum:

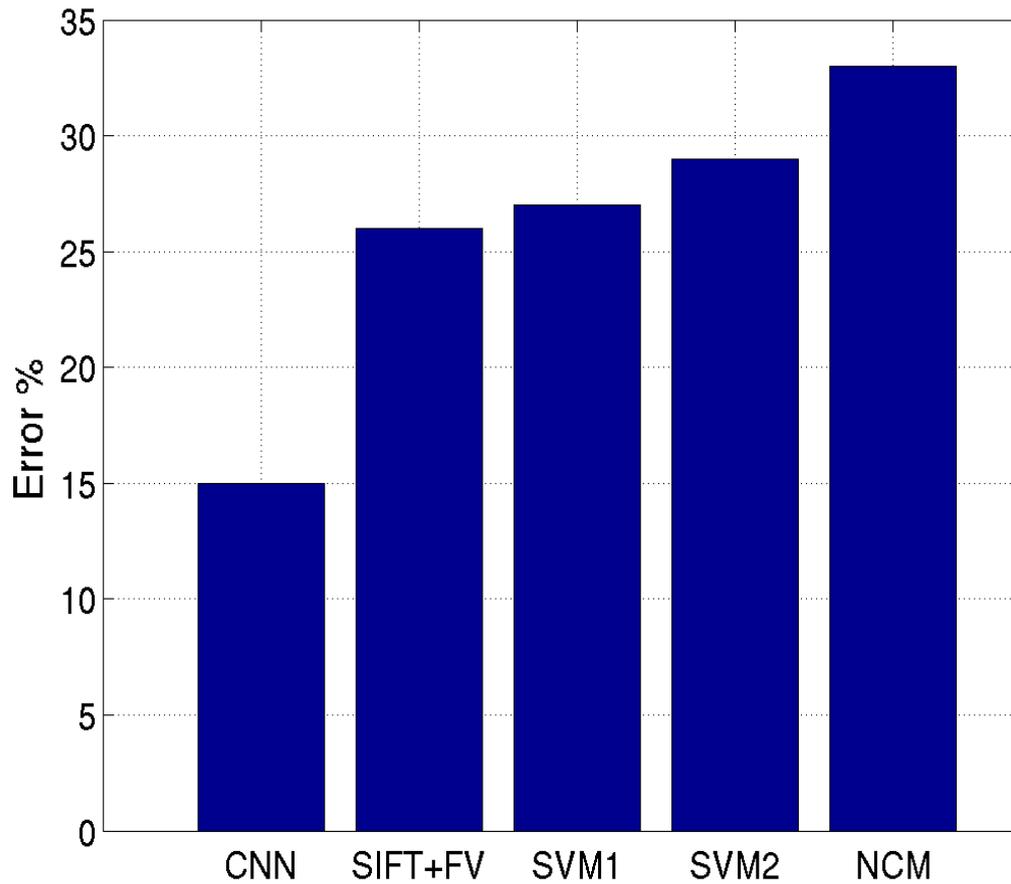
- Learning rate = 0.01
- Momentum = 0.9

## Improving generalization by:

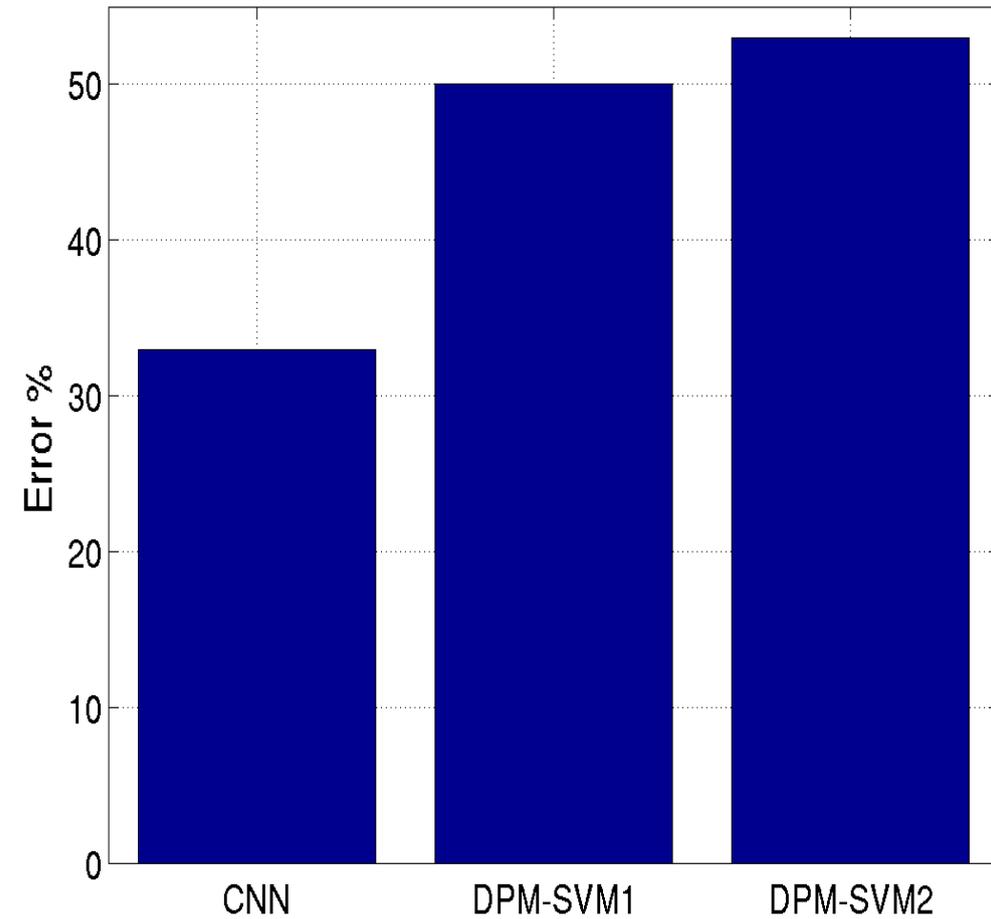
- Weight sharing (convolution)
- Input distortions
- Dropout = 0.5
- Weight decay = 0.0005

# Results: ILSVRC 2012

## TASK 1 - CLASSIFICATION



## TASK 2 - DETECTION



# Results



First layer learned filters (processing raw pixel values).



**mite**

**container ship**

**motor scooter**

**leopard**

	mite
	black widow
	cockroach
	tick
	starfish

	container ship
	lifeboat
	amphibian
	fireboat
	drilling platform

	motor scooter
	go-kart
	moped
	bumper car
	golfcart

	leopard
	jaguar
	cheetah
	snow leopard
	Egyptian cat



**grille**

**mushroom**

**cherry**

**Madagascar cat**

	convertible
	grille
	pickup
	beach wagon
	fire engine

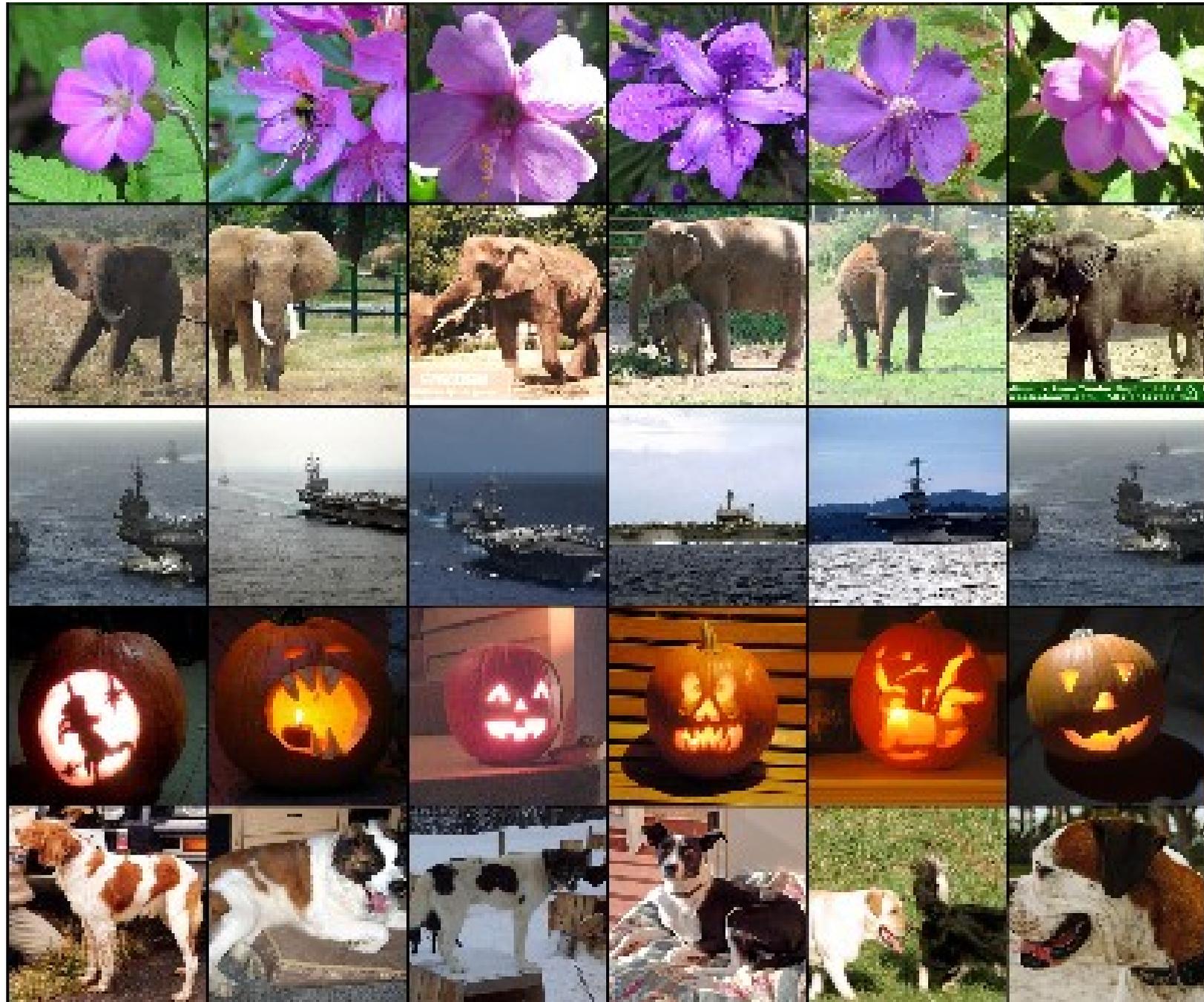
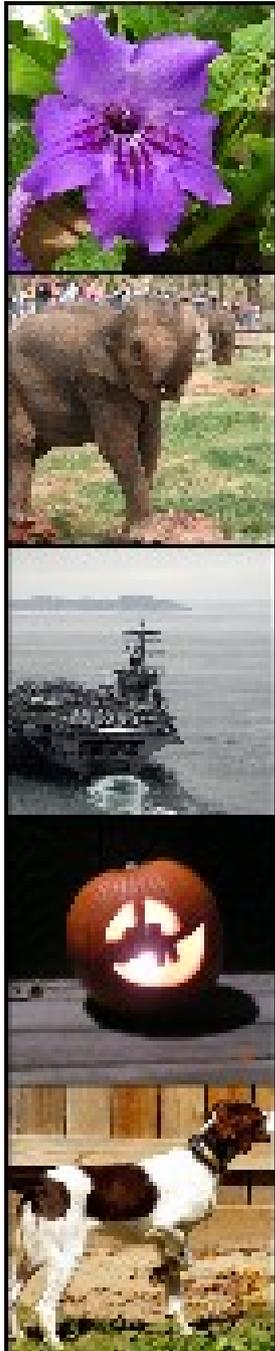
	agaric
	mushroom
	jelly fungus
	gill fungus
	dead-man's-fingers

	dalmatian
	grape
	elderberry
	ffordshire bullterrier
	currant

	squirrel monkey
	spider monkey
	titi
	indri
	howler monkey

# TEST IMAGE

# RETRIEVED IMAGES



# Demo of classifier by Matt Zeiler & Rob Fergus:

<http://horatio.cs.nyu.edu/>

The screenshot shows a web browser window with the URL `horatio.cs.nyu.edu/index.html`. The browser's address bar includes search engines like Google and Bing. The page has a navigation menu with "Image Classifier Demo", "Demo", "About", and "Terms". The main heading is "Image Classifier Demo" with the NYU logo in the top right. A text block explains the upload process: "Upload your images to have them classified by a machine! Upload multiple images using the button below or dropping them on this page. The predicted objects will be refreshed automatically. Images are resized such that the smallest dimension becomes 256, then the center 256x256 crop is used. More about the demo can be found [here](#) ."

Below the text are three buttons: a blue "+ Upload Images" button, a red "Remove All" button, and a checkbox for "Show help tips". A checked checkbox indicates agreement to the "Terms of Use".

The main content area features a small image of a lion on the left. To its right, under the heading "Predicted objects:", is a list of five items, each with a green checkmark, a red 'x', and an orange question mark icon, followed by the object name and its confidence score:

- 1. Lion, King Of Beasts, Panthera Leo (0.34)
- 2. Hartebeest (0.19)
- 3. Hyena, Hyaena (0.16)
- 4. Arabian Camel, Dromedary, Camelus Dromedarius (0.06)
- 5. Dingo, Warrigal, Warragal, Canis Dingo (0.04)

Below this list is a section for "Other objects:" with an empty text input field. A red "Remove" button is located to the right of the list.

At the bottom of the page, the text "Demo Notes" is partially visible.

# Demo of classifier by Yangqing Jia & Trevor Darrell:

<http://decafberkeleyvision.org/>

The screenshot shows a web browser window with the following elements:

- Browser tabs: Yangqing Jia, Decaf Demo, DeCAF: A Deep Convolution...
- Address bar: `decaf.berkeleyvision.org/classify_url?imageurl=http%3A%2F%2Fwww.careerealism.com%2Fhome%2Fjtdonnell%2Fcareerealism.i`
- Navigation: Most Visited, Getting Started, Latest Headlines, Click Here!
- Page Title: Decaf Image Classifier
- Text: New: get the [software](#) and [tech report](#) that we have released!  
[\[About this demo\]](#) [\[Sign up for Updates!\]](#)
- Text: Provide an image and have it classified by decaf. [Click for a Quick Example](#)
- Image: A photograph of a lion's head in a savanna setting.
- Buttons: Flat Prediction, Maximize Infogain
- Classification Results Table:

lion	0.87769
dhole	0.01987
red fox	0.01606
coyote	0.01503
red wolf	0.01321

CNN took 0.462 seconds.

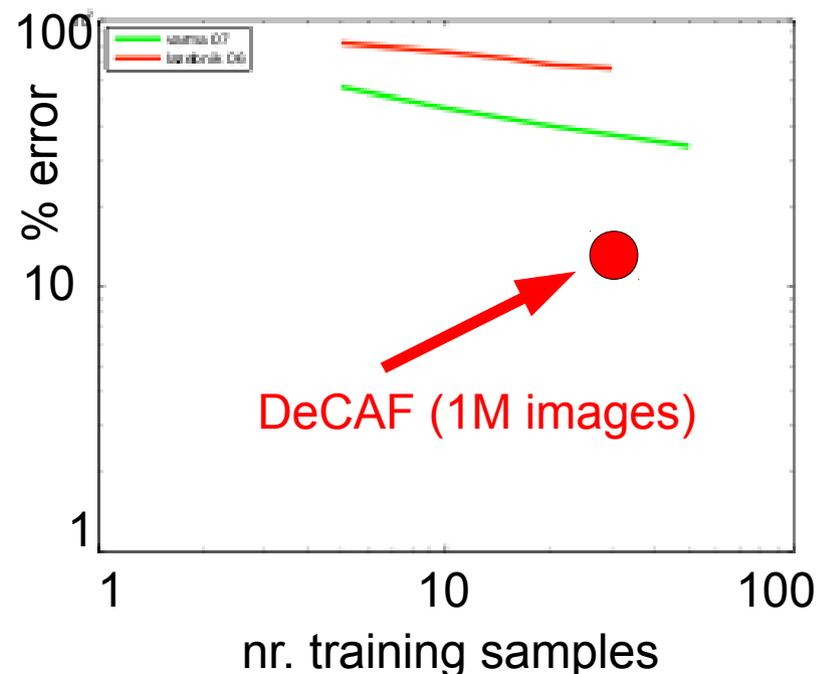
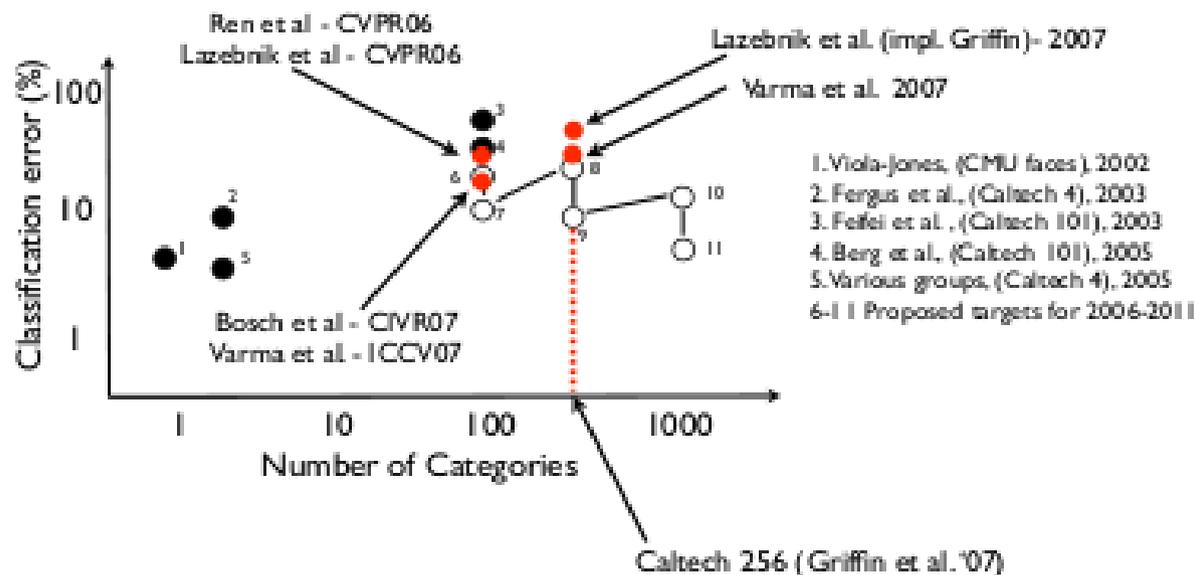
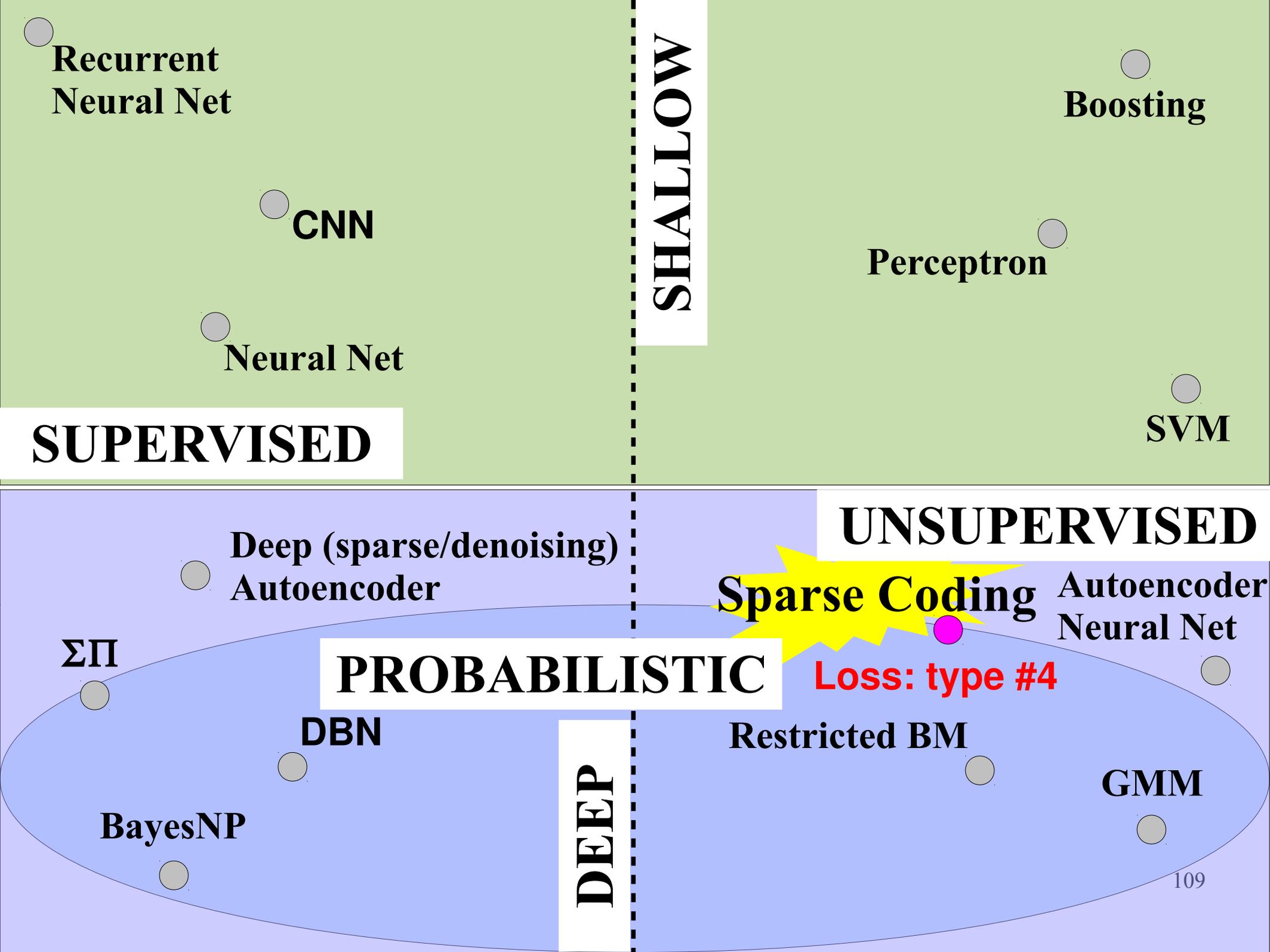


Figure 3: How well are we doing? (Left) Classification performance has seen steady improvement in the last few years, both in the number of categories on which algorithms are tested and in classification error rates. (Right) Performance of the best 2006 [Lazebnik et al., 2006] and the best 2007 algorithm [Varma, 2007] are compared here (classification error rates vs number of training examples). One may notice the significant year-on-year progress (see also left panel). Extrapolation enthusiasts may calculate that  $10^8$  training examples would be sufficient to achieve 1% error rates with current algorithms. Furthermore, if the pace of year-on-year progress is constant on this log scale chart, 1% error rates with 30 training examples will be achieved in 8-10 years.

# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- **Examples:**
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - **Unsupervised learning: sparse coding**
  - Unsupervised learning: gated MRF
- Other examples
- Practical tricks



Recurrent  
Neural Net

CNN

Neural Net

SHALLOW

Boosting

Perceptron

SVM

SUPERVISED

UNSUPERVISED

Deep (sparse/denoising)  
Autoencoder

Sparse Coding

Autoencoder  
Neural Net

$\Sigma\Pi$

PROBABILISTIC

Loss: type #4

DBN

Restricted BM

BayesNP

GMM

DEEP

# Energy & latent variables

- Energy may have latent variables.
- Two major approaches:
  - Marginalization (intractable if space is large)

$$\tilde{E}(\mathbf{y}) = \log \sum_{\mathbf{h}} \exp(-E(\mathbf{y}, \mathbf{h}))$$

- Minimization

$$\tilde{E}(\mathbf{y}) = \min_{\mathbf{h}} E(\mathbf{y}, \mathbf{h})$$

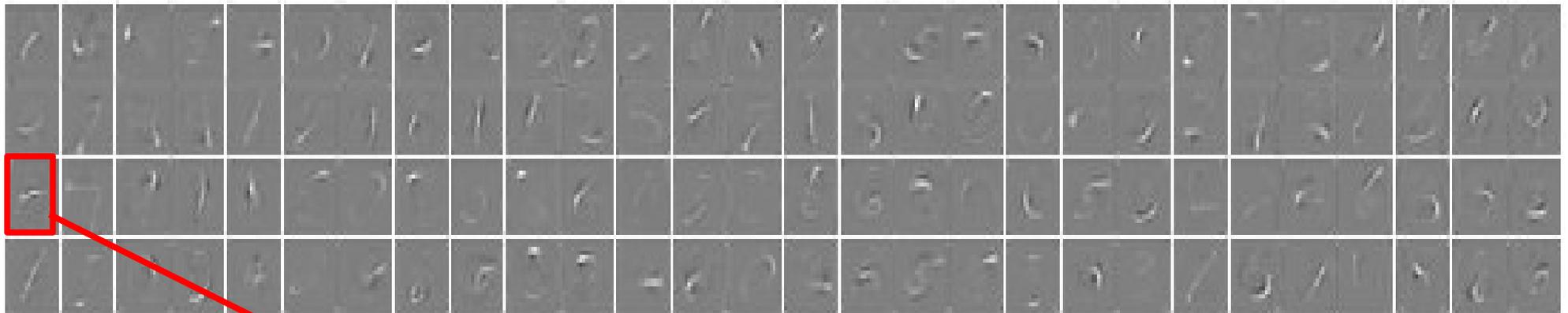
# Sparse Coding

$$E(\mathbf{x}, \mathbf{h}; W) = \frac{1}{2} \|\mathbf{x} - W\mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1$$

$$\tilde{E}(\mathbf{x}; W) = \min_{\mathbf{h}} E(\mathbf{x}, \mathbf{h}; W)$$

**Loss type #4:** energy loss with (sparsity) constraints.

$$L = \tilde{E}(\mathbf{x}; W)$$



$$\boxed{7} = 1 \boxed{9} + 1 \boxed{7} + 1 \boxed{2} + 1 \boxed{9} + 1 \boxed{1} + 1 \boxed{1} + 1 \boxed{7} + 0.8 \boxed{1} + 0.8 \boxed{7} + \dots$$

# Inference

- Prediction of latent variables

$$E(\mathbf{x}, \mathbf{h}; W) = \frac{1}{2} \|\mathbf{x} - W \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1$$

$$\mathbf{h}^* = \arg \min_{\mathbf{h}} E(\mathbf{x}, \mathbf{h}; W)$$

Inference is an iterative process.

# Inference

- Prediction of latent variables

$$E(\mathbf{x}, \mathbf{h}; W) = \frac{1}{2} \|\mathbf{x} - W \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1$$

$$\mathbf{h}^* = \arg \min_{\mathbf{h}} E(\mathbf{x}, \mathbf{h}; W)$$

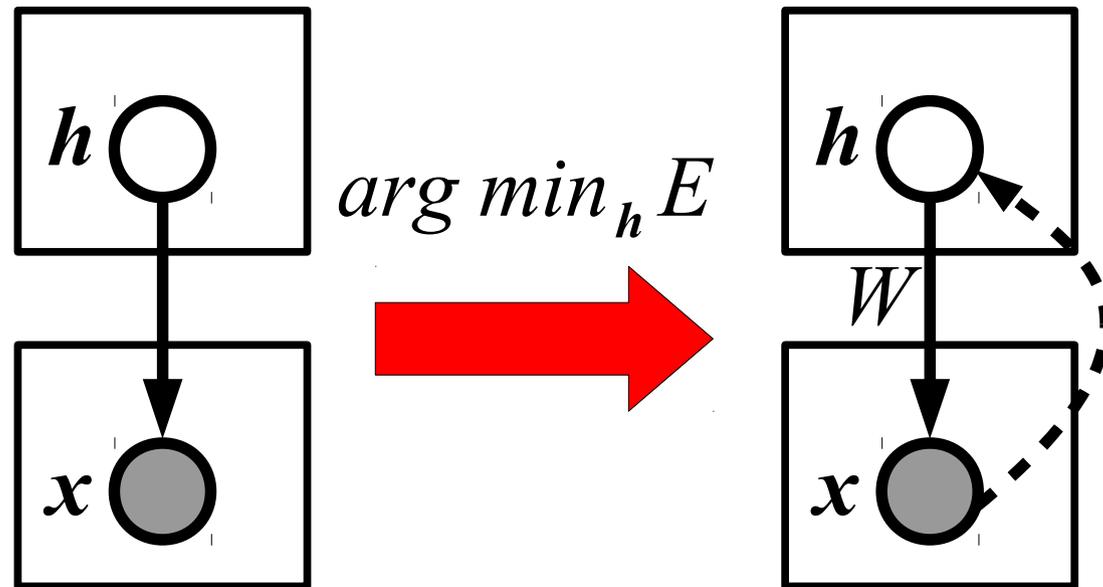
Inference is an iterative process.

**Q.:** Is it possible to make inference more efficient?

**A.:** Yes, by training another module to directly predict

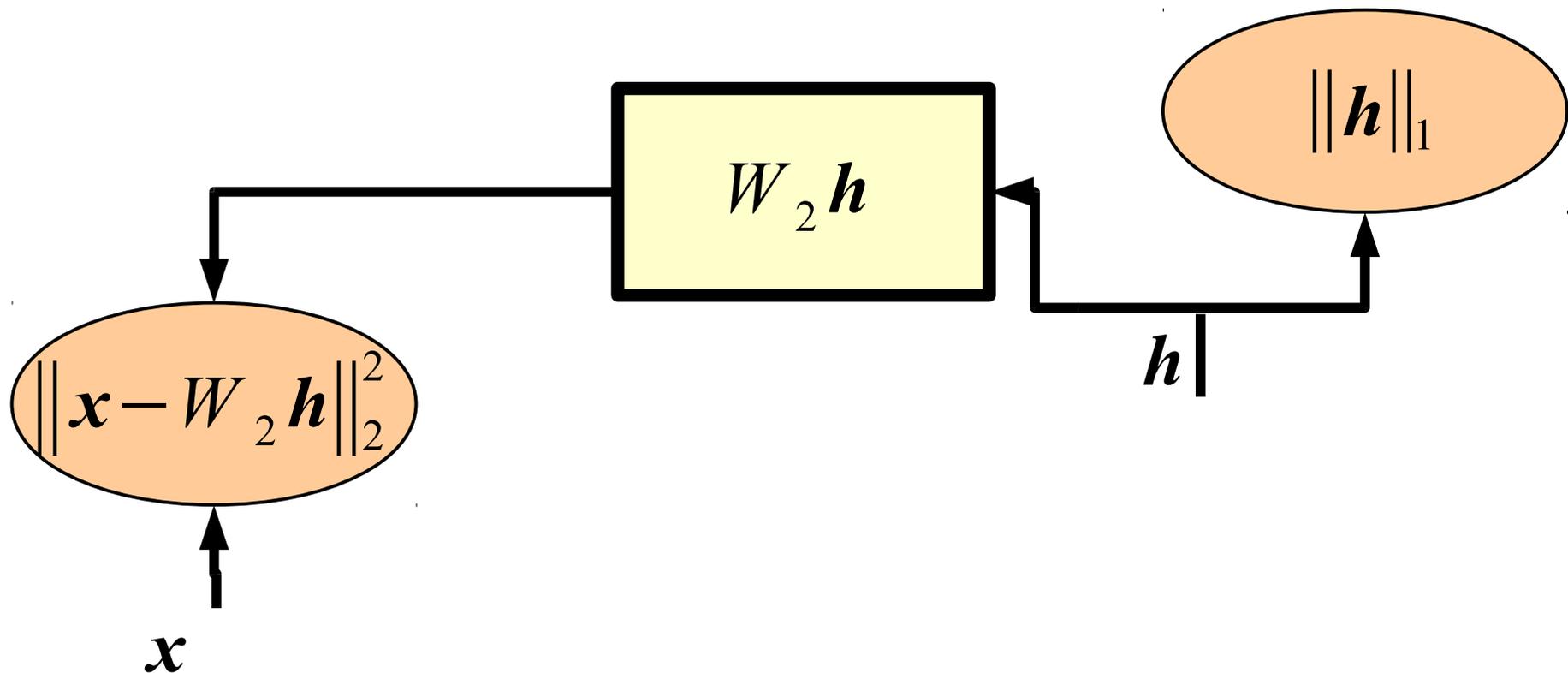
# Inference in Sparse Coding

$$E(\mathbf{x}, \mathbf{h}) = \frac{1}{2} \|\mathbf{x} - W_2 \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1$$



# Inference in Sparse Coding

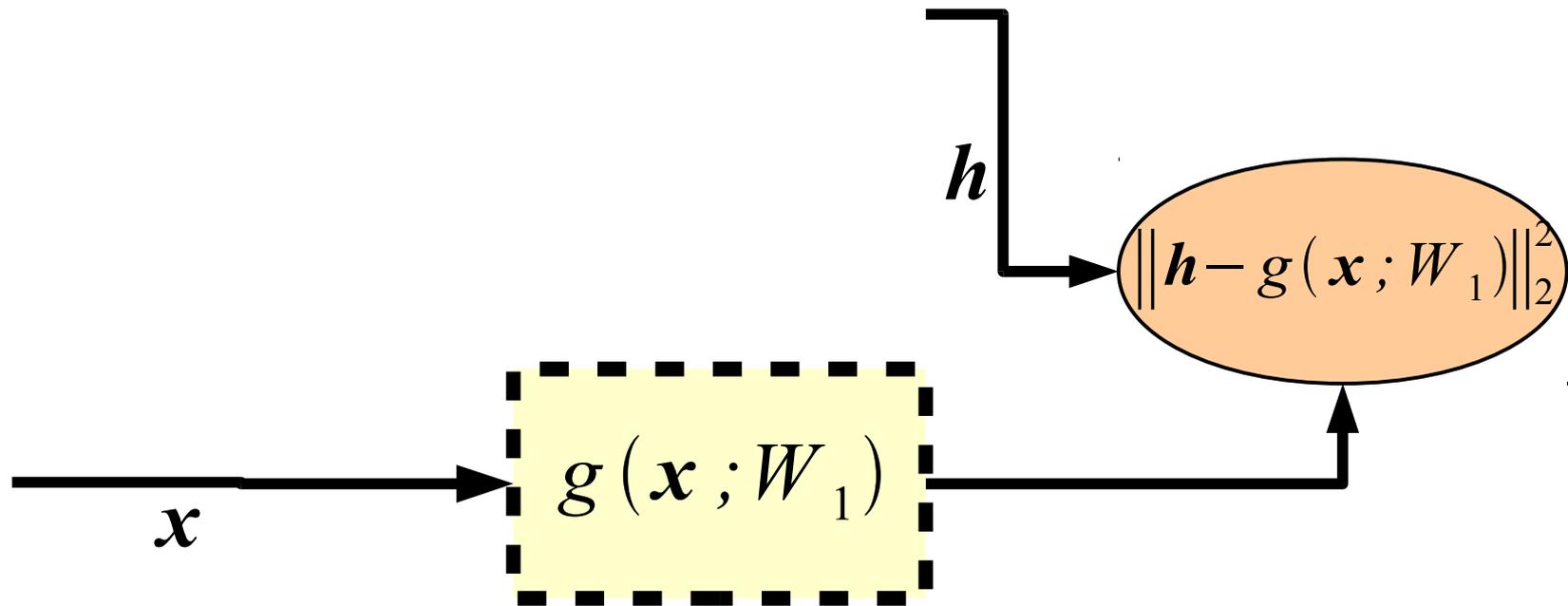
$$E(\mathbf{x}, \mathbf{h}) = \frac{1}{2} \|\mathbf{x} - W_2 \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1$$



# Learning To Perform Fast Inference

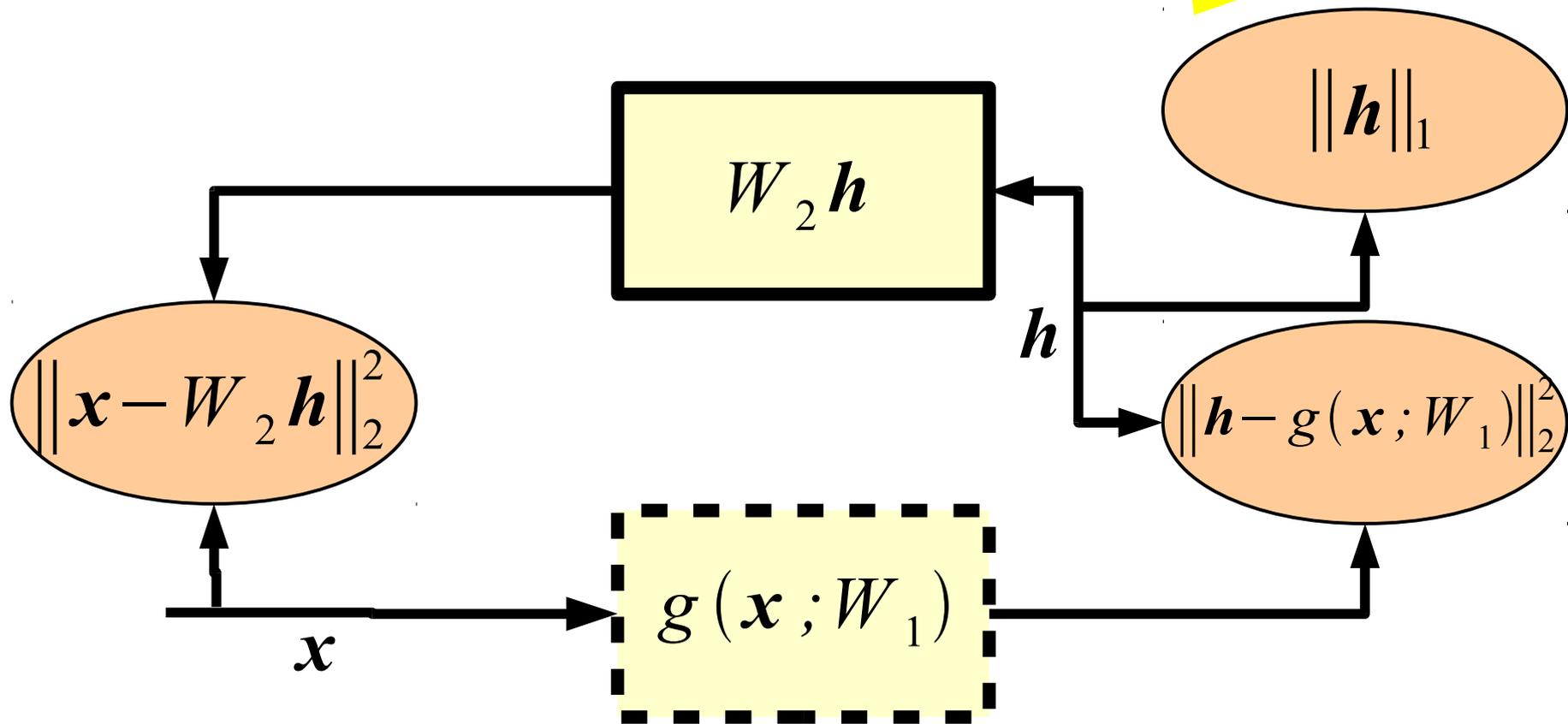
$$E(\mathbf{x}, \mathbf{h}) =$$

$$\frac{1}{2} \|\mathbf{h} - g(\mathbf{x}; W_1)\|_2^2$$



# Predictive Sparse Decomposition

$$E(\mathbf{x}, \mathbf{h}) = \frac{1}{2} \|\mathbf{x} - W_2 \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1 + \frac{1}{2} \|\mathbf{h} - g(\mathbf{x}; W_1)\|_2^2$$



# Sparse Auto-Encoders

- Example: Predictive Sparse Decomposition

$$E(\mathbf{x}, \mathbf{h}) = \frac{1}{2} \|\mathbf{x} - W_2 \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1 + \frac{1}{2} \|\mathbf{h} - g(\mathbf{x}; W_1)\|_2^2$$

## TRAINING:

For every sample:

Initialize  $\mathbf{h} = g(\mathbf{x}; W_1)$

(E) Infer optimal latent variables:  $\mathbf{h}^* = \min_{\mathbf{h}} E(\mathbf{x}, \mathbf{h}; W)$

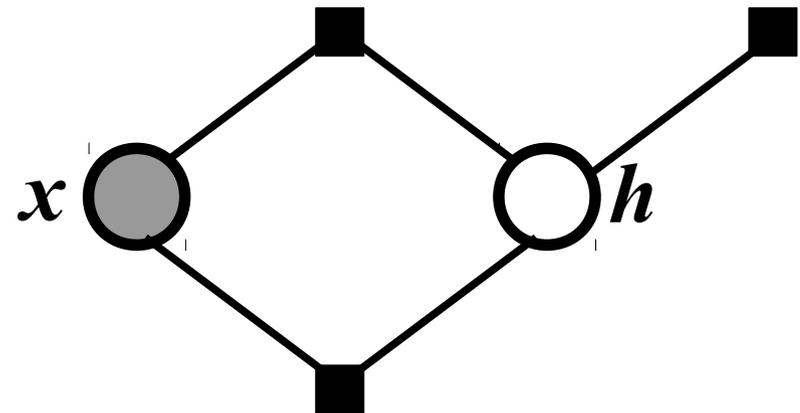
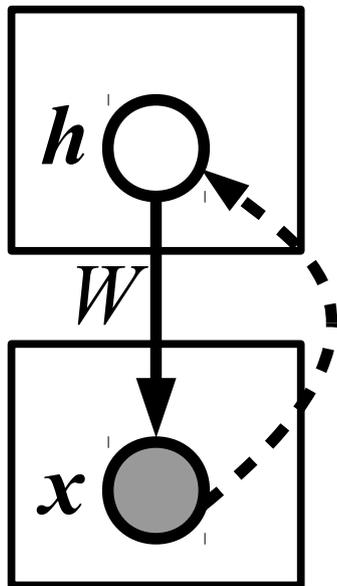
(M) Update parameters  $W_1, W_2$

## Inference at test time (fast):

$$\mathbf{h}^* \approx g(\mathbf{x}; W_1)$$

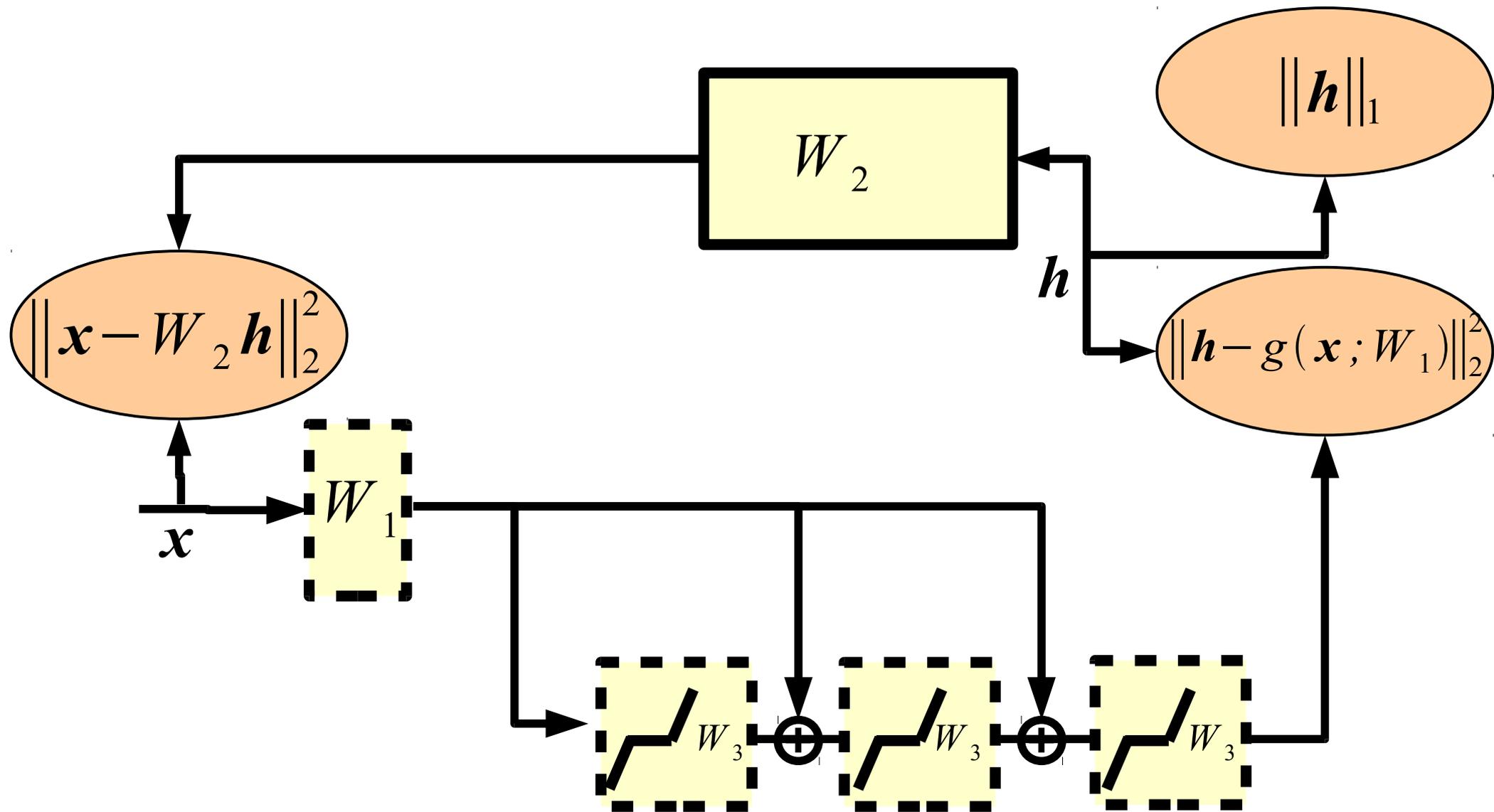
# Predictive Sparse Decomposition

$$E(\mathbf{x}, \mathbf{h}) = \frac{1}{2} \|\mathbf{x} - W_2 \mathbf{h}\|_2^2 + \lambda \|\mathbf{h}\|_1 + \frac{1}{2} \|\mathbf{h} - g(\mathbf{x}; W_1)\|_2^2$$



alternative graphical representations

# LISTA



# KEY IDEAS

- Inference can require expensive optimization
- We may approximate exact inference well by using a non-linear function (learn optimal approximation to perform fast inference)
- The original model and the fast predictor can be trained jointly

Kavukcuoglu et al. “Predictive Sparse Decomposition” ArXiv 2008

Kavukcuoglu et al. “Learning convolutional feature hierarchies..” NIPS 2010

Gregor et al. “Structured sparse coding via lateral inhibition” NIPS 2011

Szlam et al. “Fast approximations to structured sparse coding...” ECCV 2012

Rolfe et al. “Discriminative Recurrent Sparse Autoencoders” ICLR 2013

# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- **Examples:**
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - Unsupervised learning: sparse coding
  - **Unsupervised learning: gated MRF**
- Other examples
- Practical tricks

Recurrent  
Neural Net

CNN

Neural Net

SHALLOW

Boosting

Perceptron

SVM

SUPERVISED

UNSUPERVISED

Deep (sparse/denoising)  
Autoencoder

Sparse Coding

Autoencoder  
Neural Net

$\Sigma\Pi$

PROBABILISTIC

DBN

Loss: type #2

BayesNP

Restricted BM

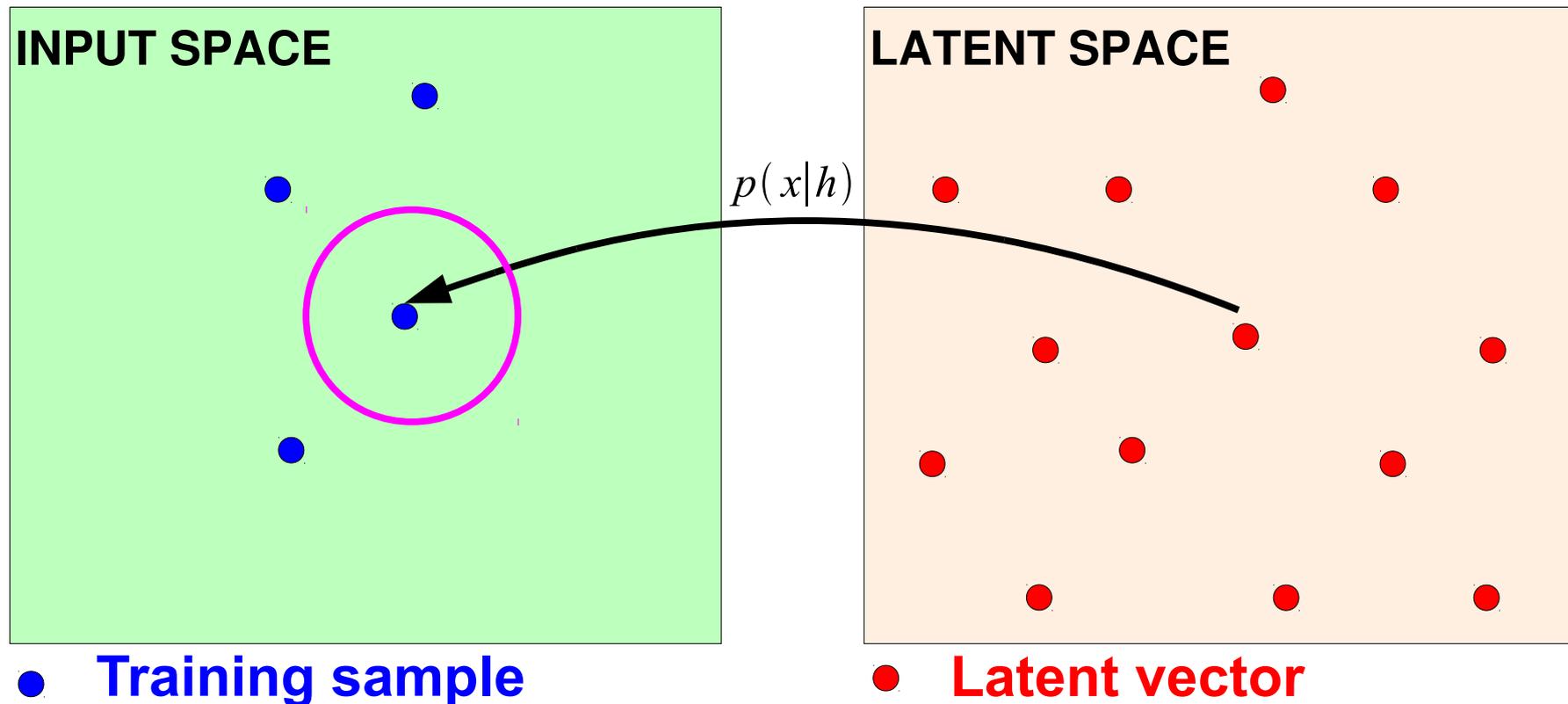
GMM

DEEP

# Probabilistic Models of Natural Images

$$p(x|h) = N(\text{mean}(h), D)$$

- examples: PPCA, Factor Analysis, ICA, Gaussian RBM



# Probabilistic Models of Natural Images

$$p(x|h) = N(\text{mean}(h), D)$$

- examples: PPCA, Factor Analysis, ICA, Gaussian RBM



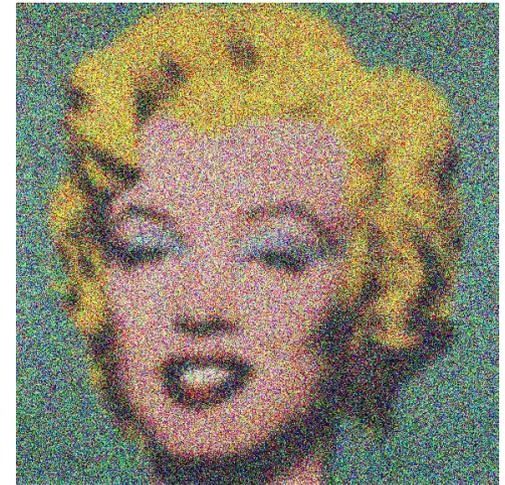
input image

$$p(h|x)$$



latent variables

$$p(x|h)$$



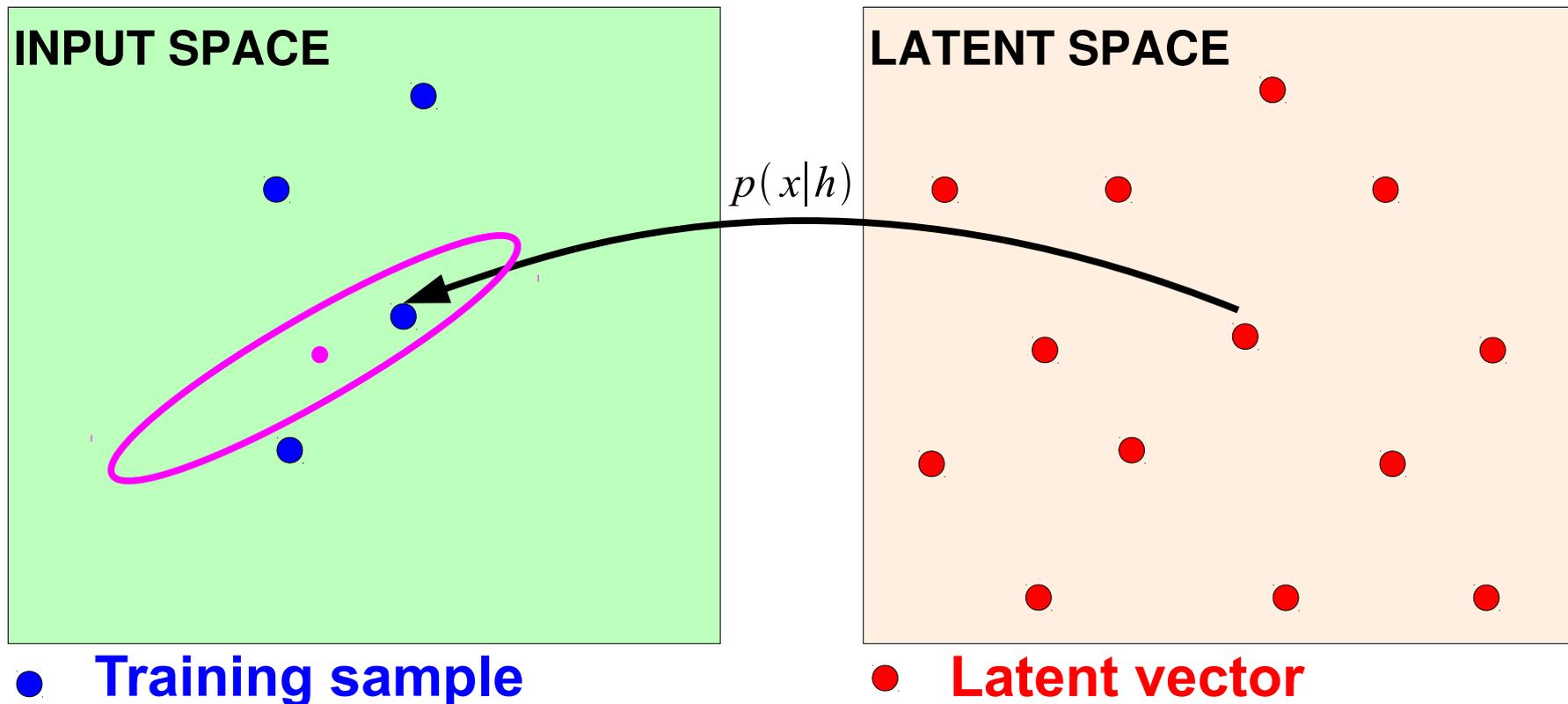
generated image

model does not represent well dependencies, only mean intensity

# Probabilistic Models of Natural Images

$$p(x|h) = N(0, \text{Covariance}(h))$$

- examples: PoT, cRBM



# Probabilistic Models of Natural Images

$$p(x|h) = N(0, \text{Covariance}(h))$$

- examples: PoT, cRBM



input image

$$p(h|x)$$



latent variables

$$p(x|h)$$



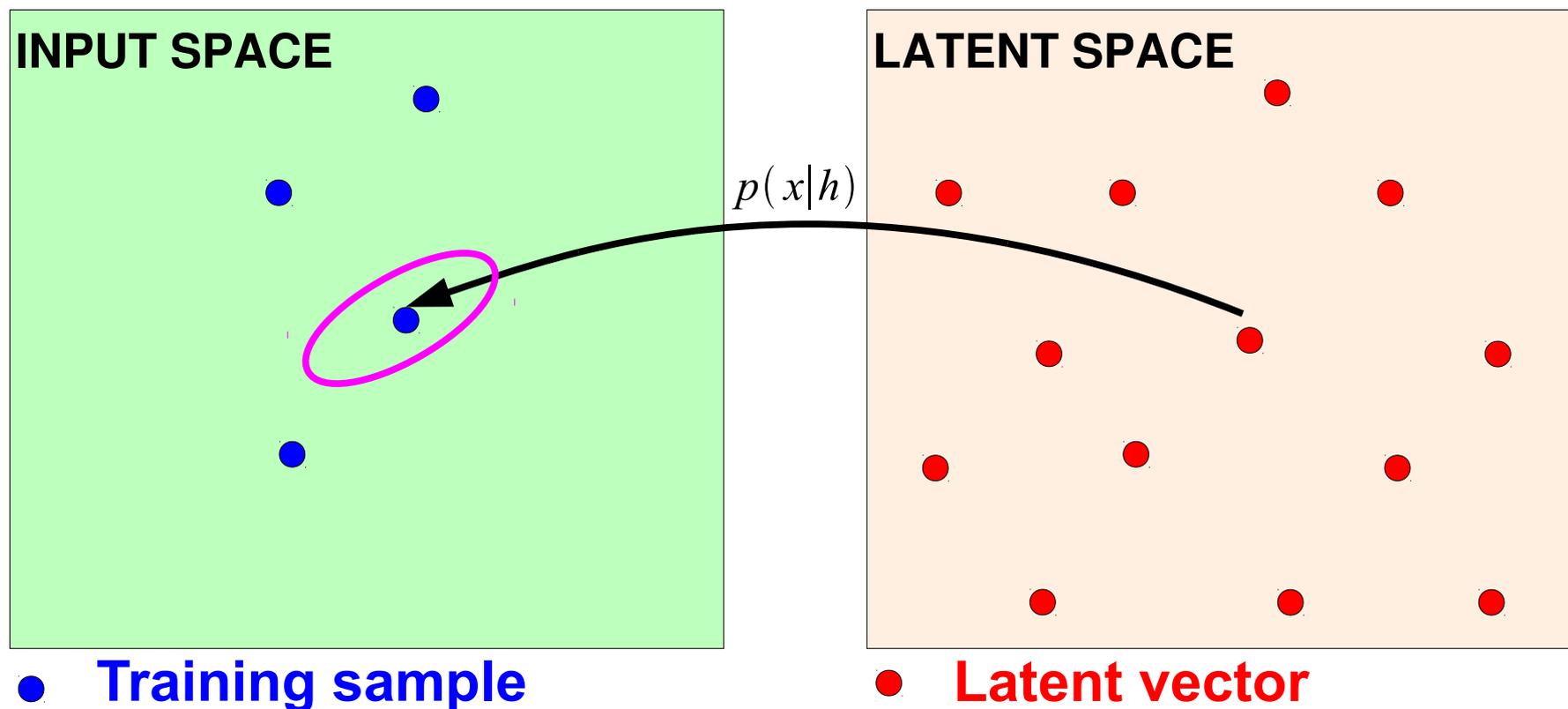
generated image

model does not represent well mean intensity, only dependencies

# Probabilistic Models of Natural Images

$$p(x|h) = N(\text{mean}(h), \text{Covariance}(h))$$

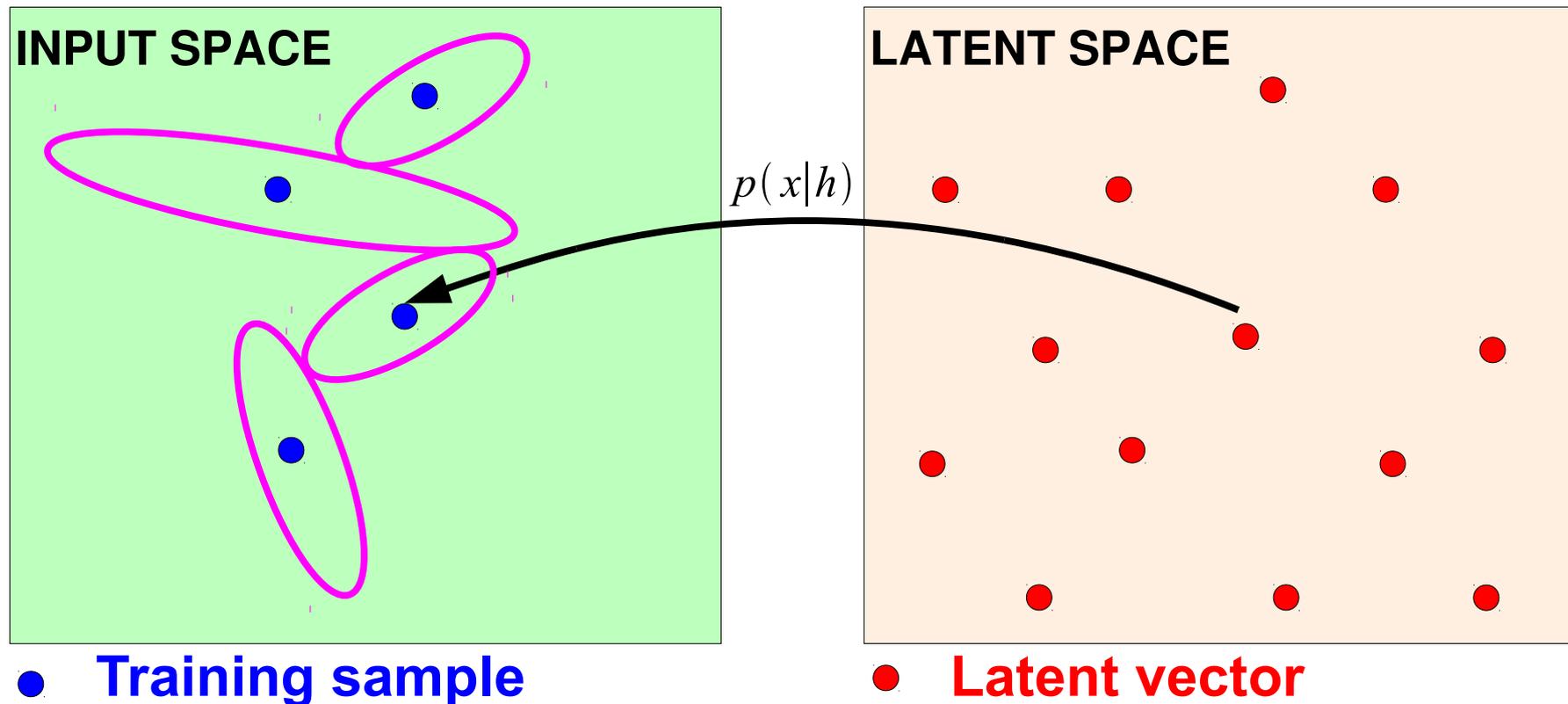
- this is what we propose: mcRBM, mPoT



# Probabilistic Models of Natural Images

$$p(x|h) = N(\text{mean}(h), \text{Covariance}(h))$$

- this is what we propose: mcRBM, mPoT



PoT



$N(0, \Sigma)$

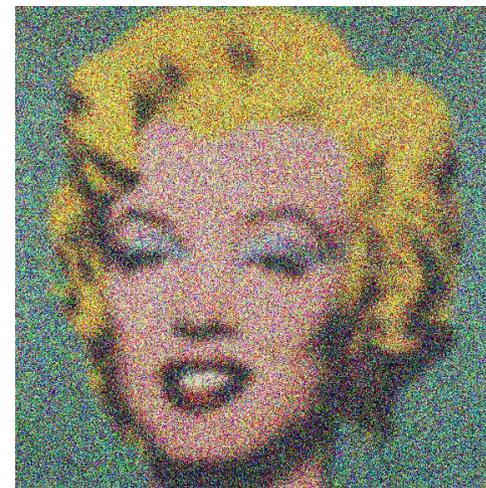


Our model

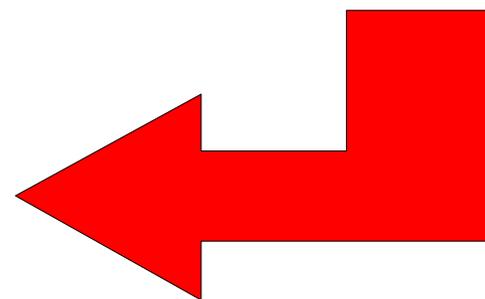
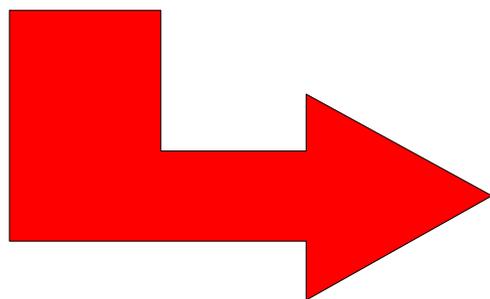


$N(m, \Sigma)$

PPCA



$N(m, I)$



# Deep Gated MRF

Layer 1:

$$E(x, h^c, h^m) = \frac{1}{2} x' \Sigma^{-1} x \quad p(x, h^c, h^m) \propto e^{-E(x, h^c, h^m)}$$

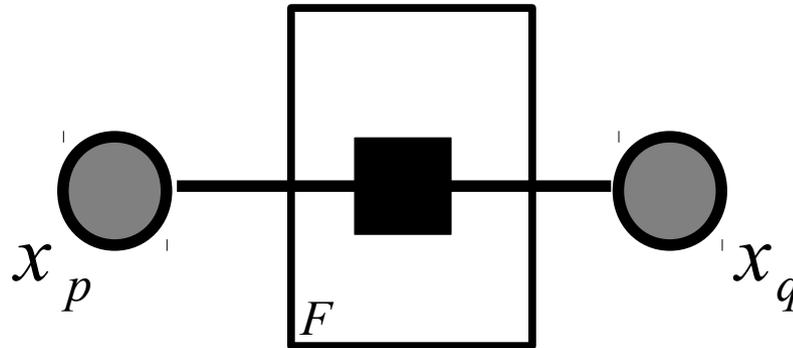


pair-wise MRF

# Deep Gated MRF

Layer 1:

$$E(x, h^c, h^m) = \frac{1}{2} x' C C' x$$



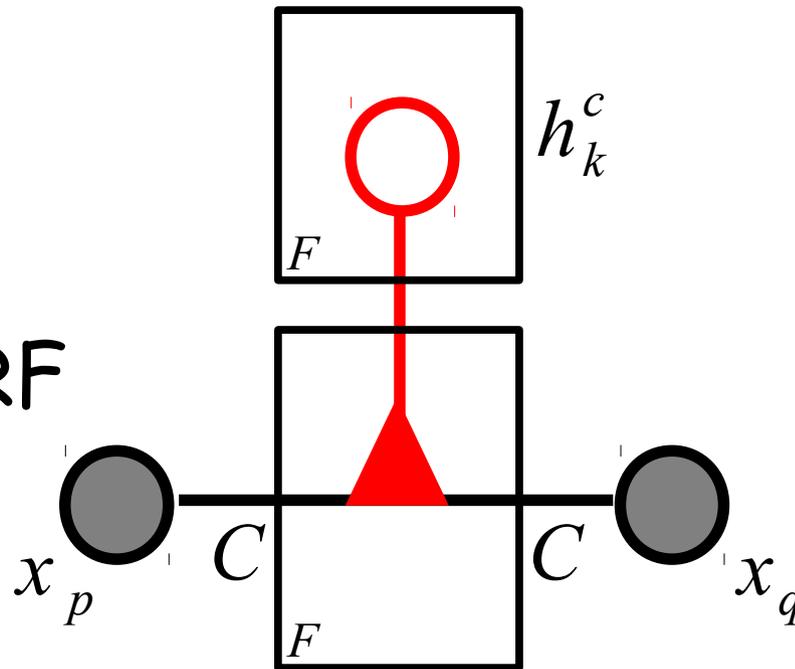
pair-wise MRF

# Deep Gated MRF

Layer 1:

$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x$$

gated MRF

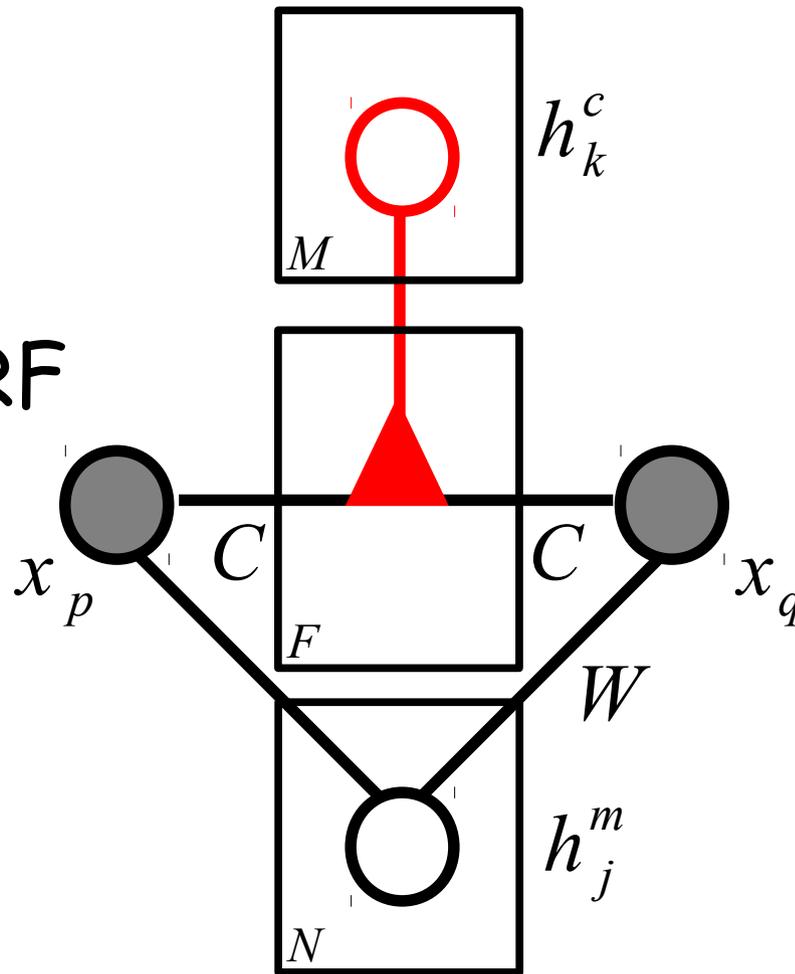


# Deep Gated MRF

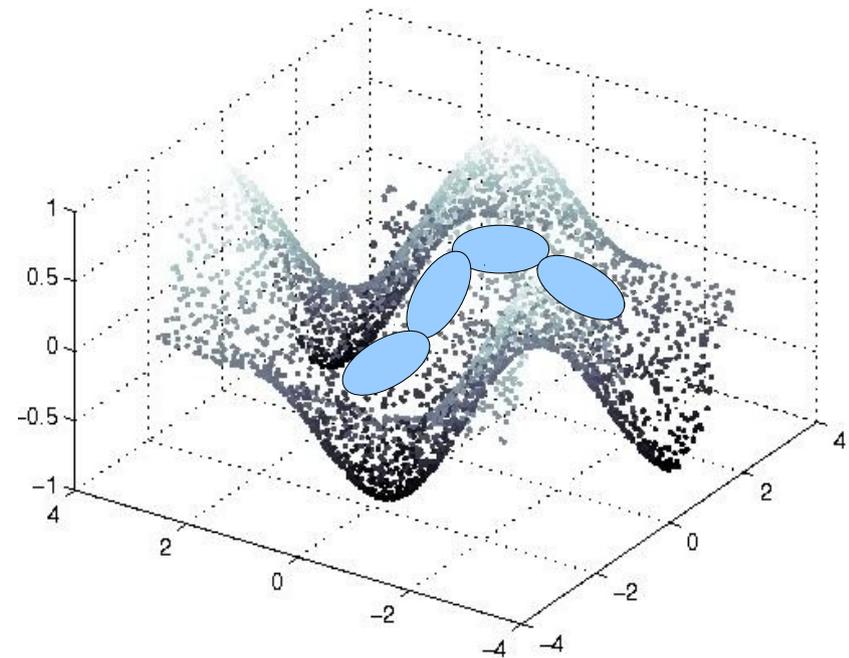
Layer 1:

$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x + \frac{1}{2} x' x - x' W h^m$$

gated MRF



$$p(x) \propto \int_{h^c} \int_{h^m} e^{-E(x, h^c, h^m)}$$



# Deep Gated MRF

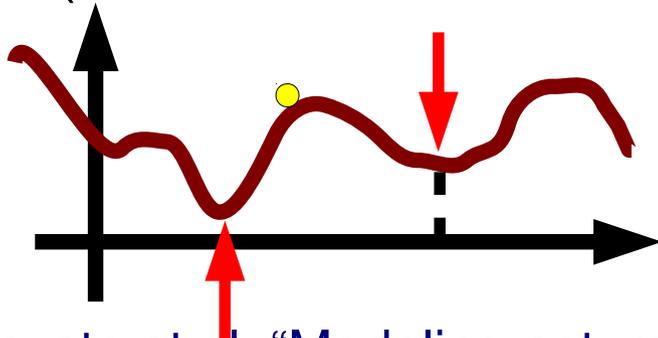
**Layer 1:**

$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x + \frac{1}{2} x' x - x' W h^m$$

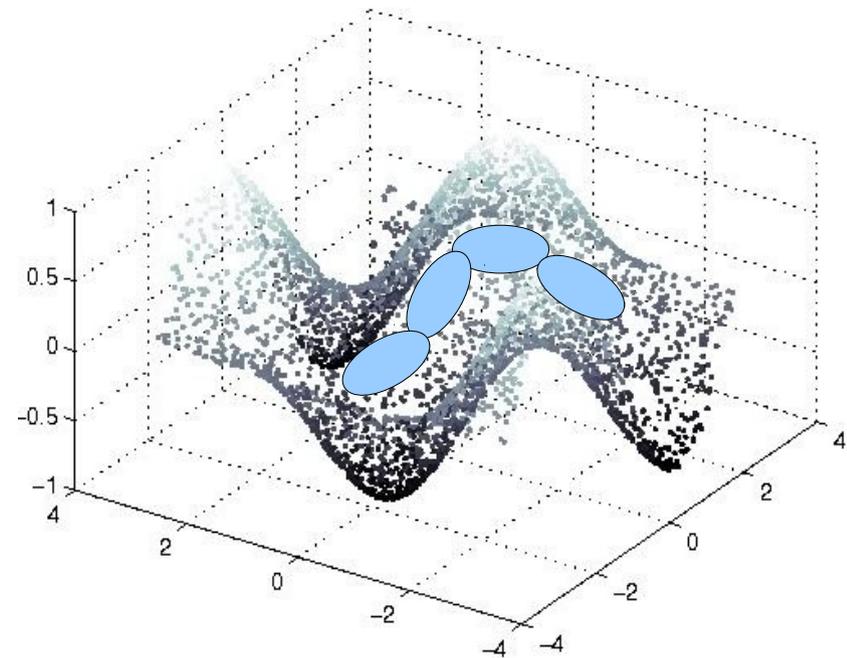
**Inference of latent variables:**  
just a forward pass

$$p(x) \propto \int_{h^c} \int_{h^m} e^{-E(x, h^c, h^m)}$$

**Training:**  
requires approximations  
(here we used HMC with PCD)



**Loss: type #2**



# Deep Gated MRF

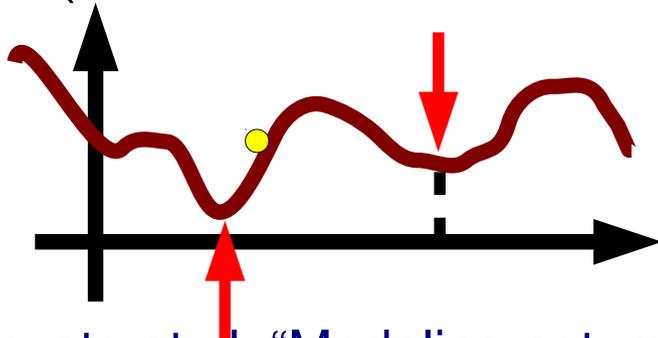
**Layer 1:**

$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x + \frac{1}{2} x' x - x' W h^m$$

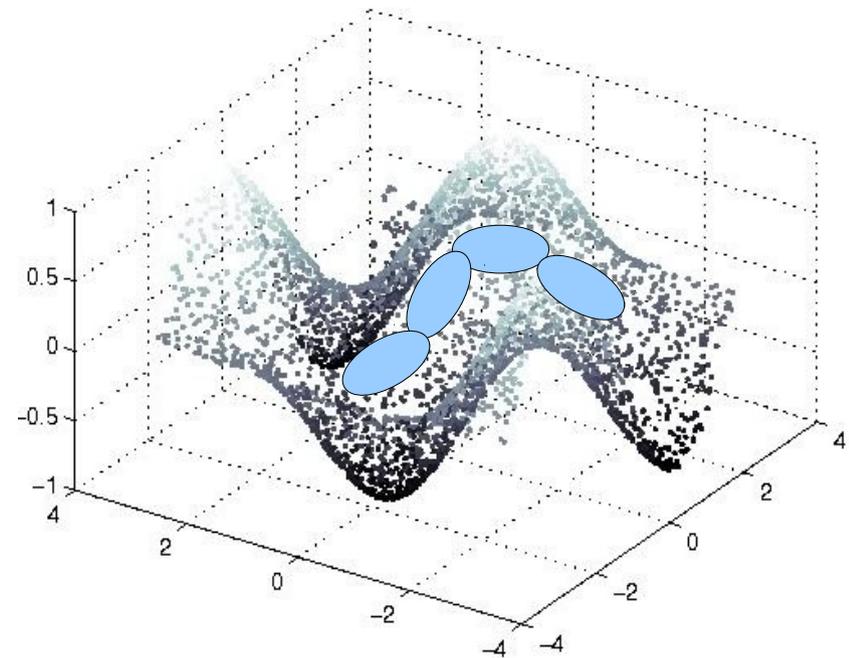
**Inference of latent variables:**  
just a forward pass

$$p(x) \propto \int_{h^c} \int_{h^m} e^{-E(x, h^c, h^m)}$$

**Training:**  
requires approximations  
(here we used HMC with PCD)



**Loss: type #2**



# Deep Gated MRF

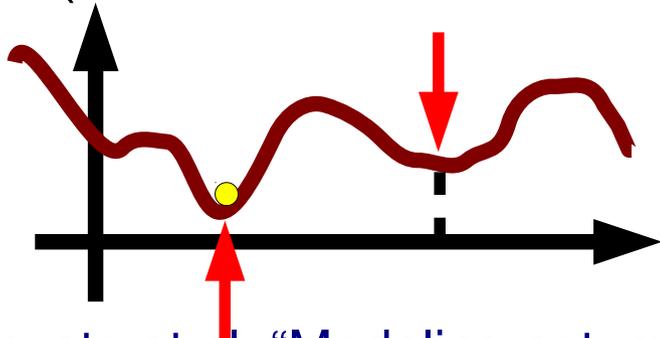
**Layer 1:**

$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x + \frac{1}{2} x' x - x' W h^m$$

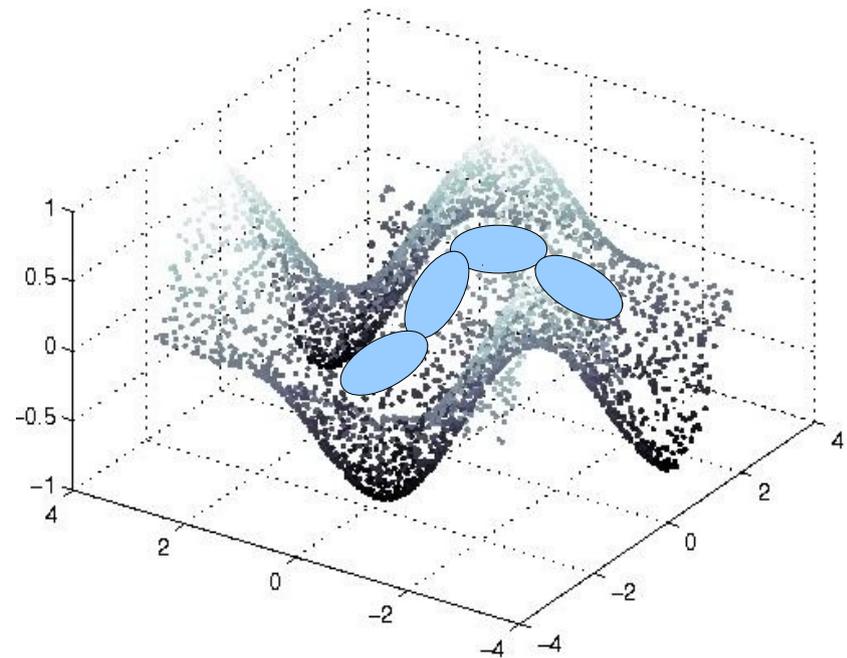
**Inference of latent variables:**  
just a forward pass

$$p(x) \propto \int_{h^c} \int_{h^m} e^{-E(x, h^c, h^m)}$$

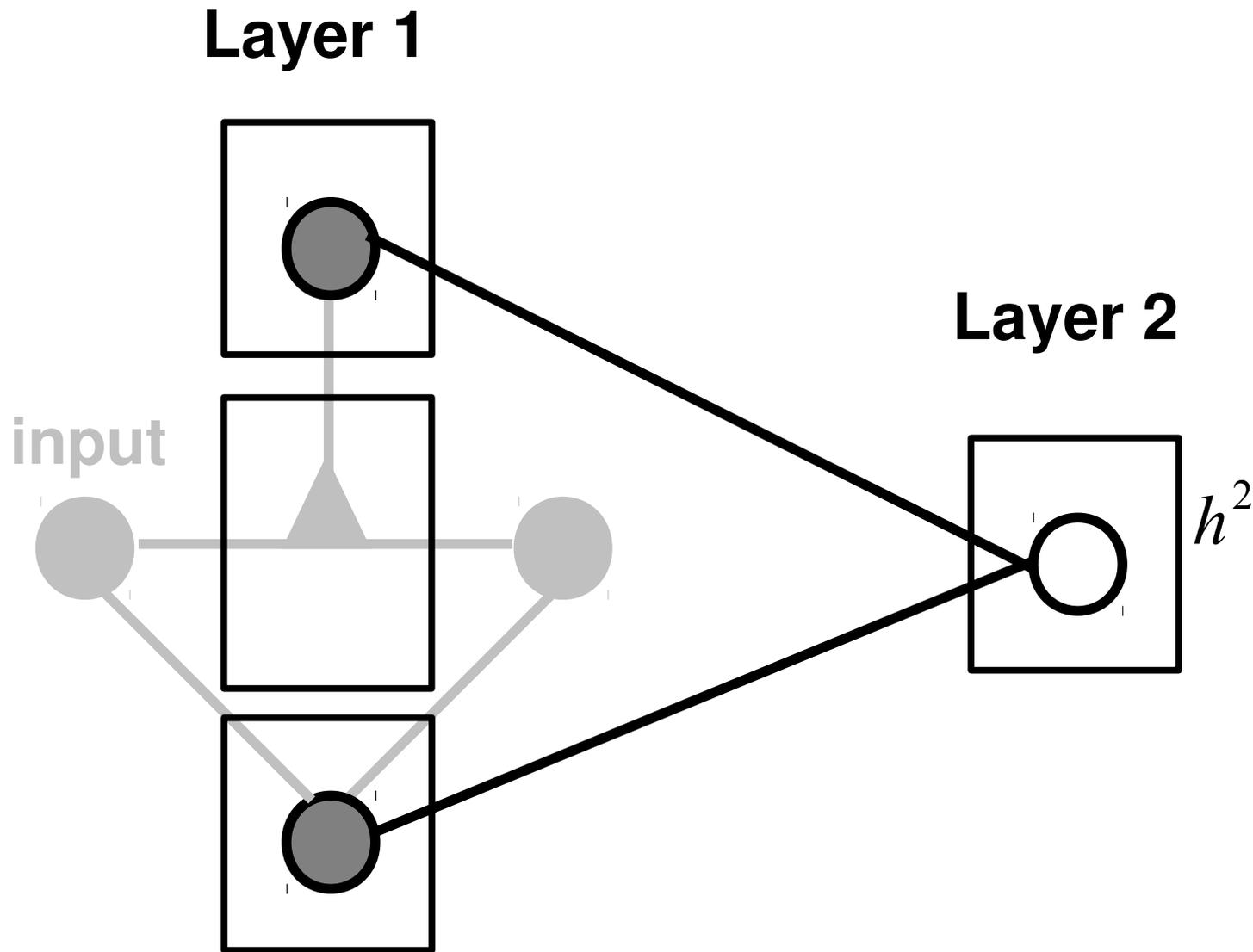
**Training:**  
requires approximations  
(here we used HMC with PCD)



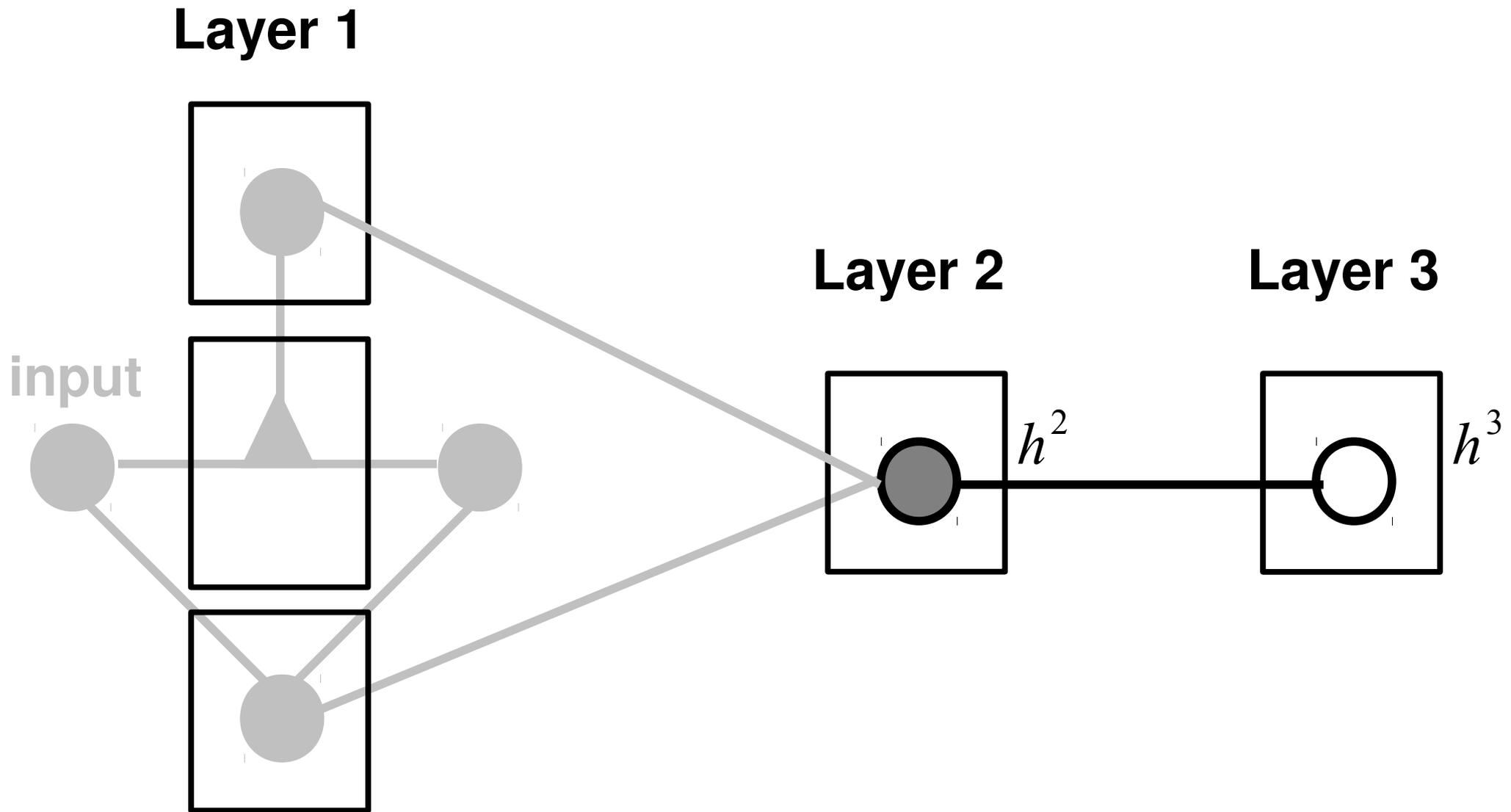
**Loss: type #2**



# Deep Gated MRF

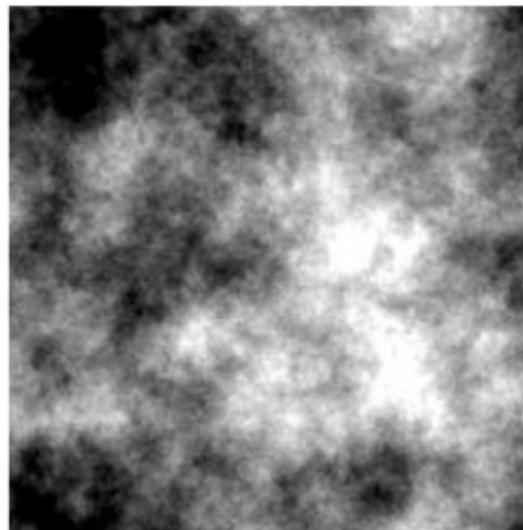


# Deep Gated MRF

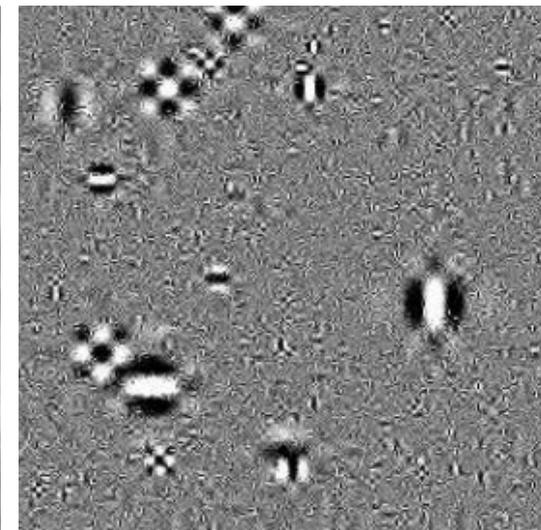


# Sampling High-Resolution Images

Gaussian model

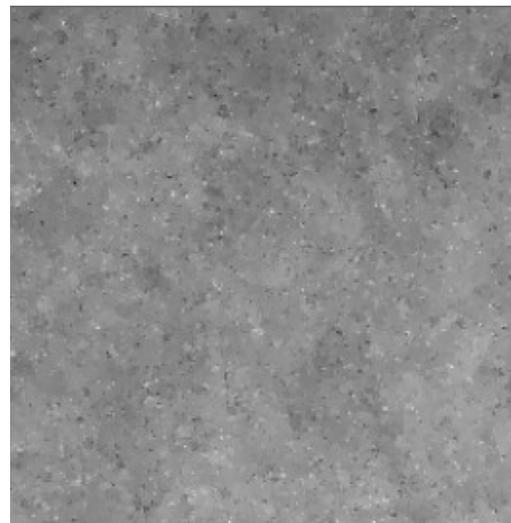


marginal wavelet

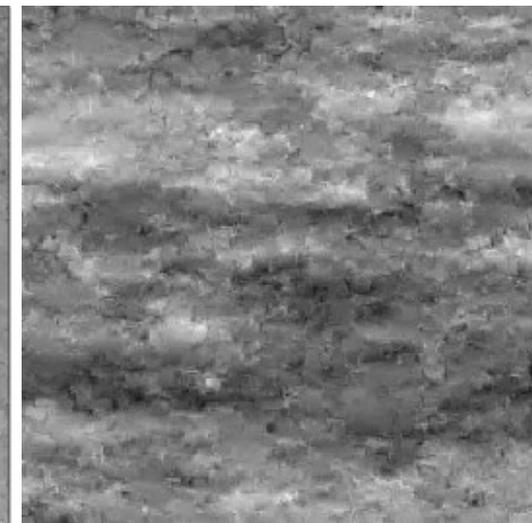


*from Simoncelli 2005*

Pair-wise MRF



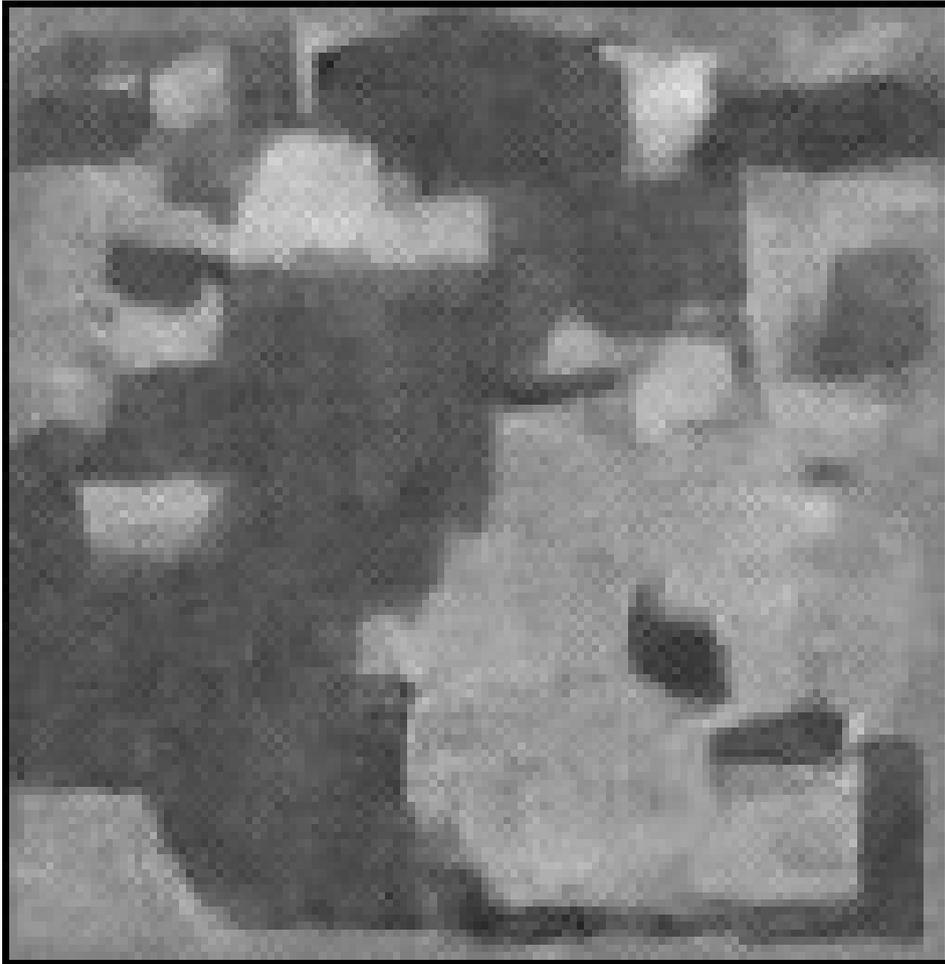
FoE



*from Schmidt, Gao, Roth CVPR 2010*

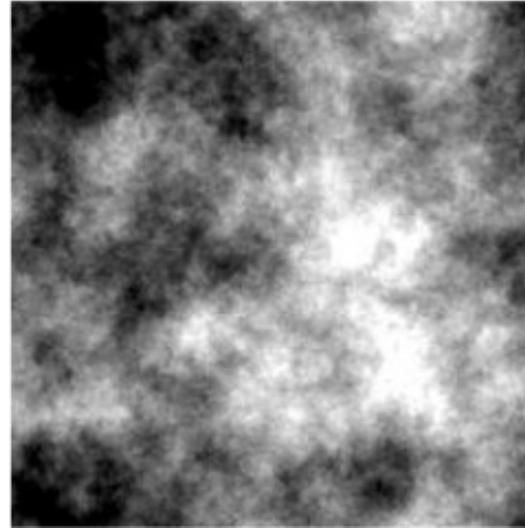
# Sampling High-Resolution Images

gMRF: 1 layer



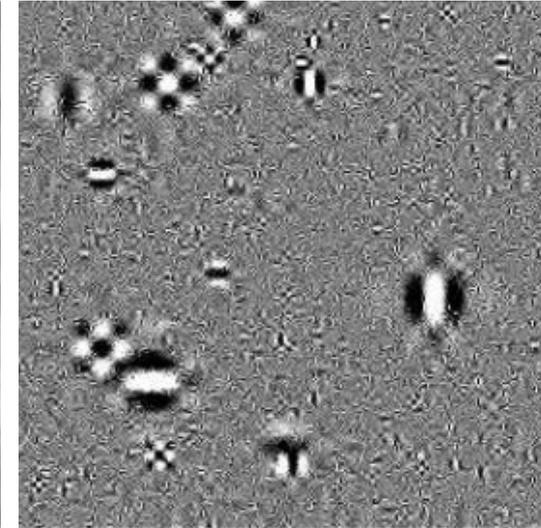
Ranzato et al. PAMI 2013

Gaussian model

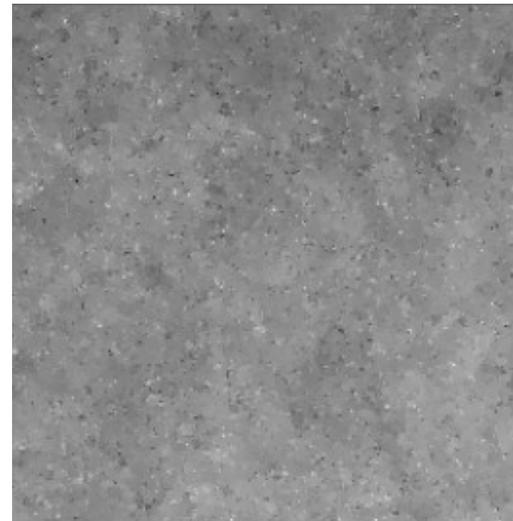


from Simoncelli 2005

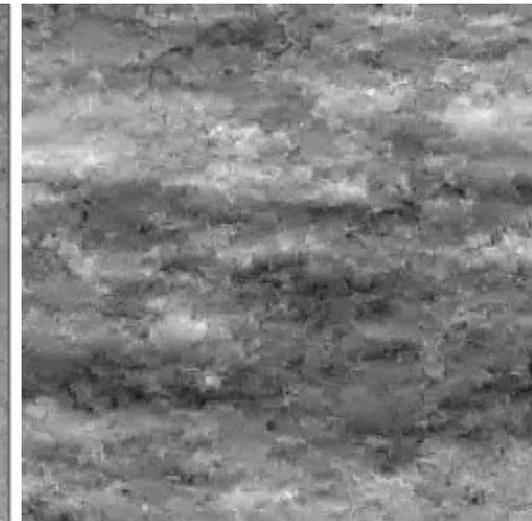
marginal wavelet



Pair-wise MRF



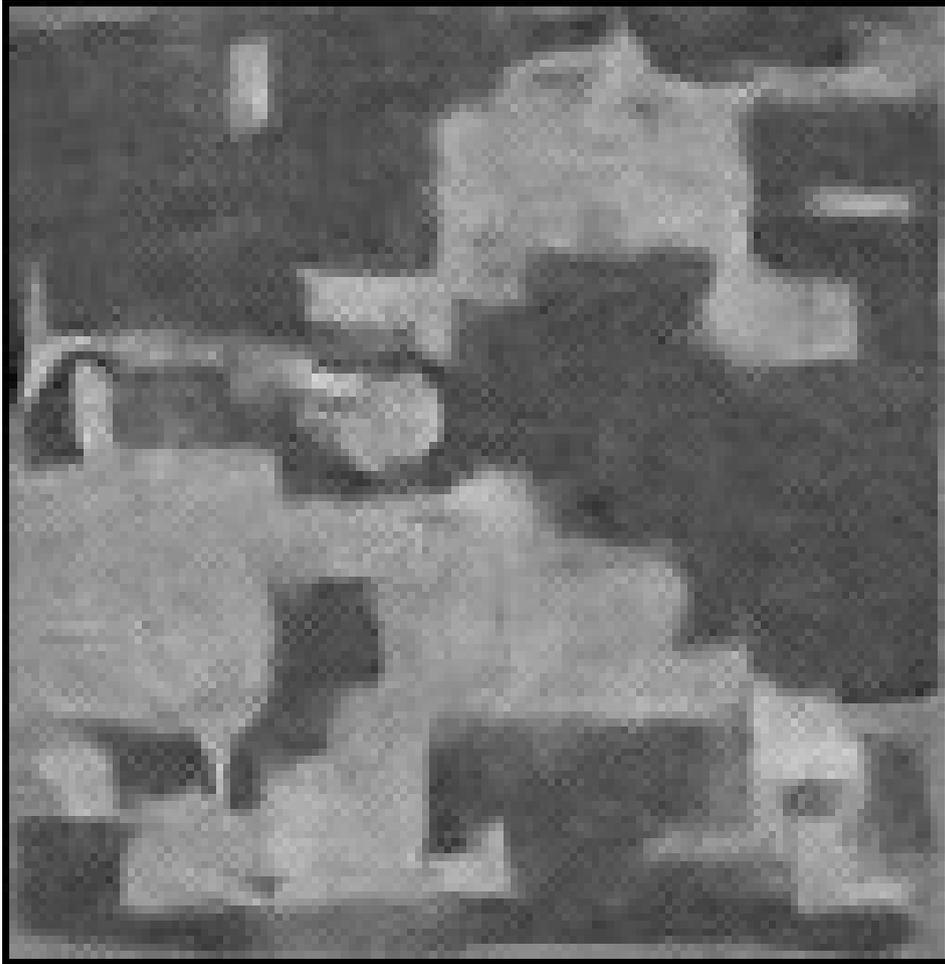
FoE



from Schmidt, Gao, Roth CVPR 2010

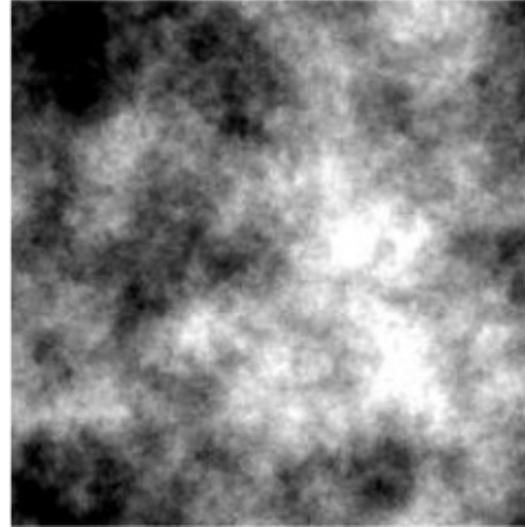
# Sampling High-Resolution Images

gMRF: 1 layer



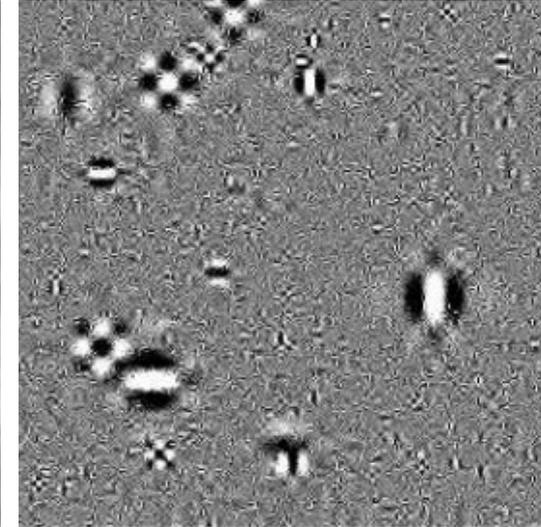
Ranzato et al. PAMI 2013

Gaussian model

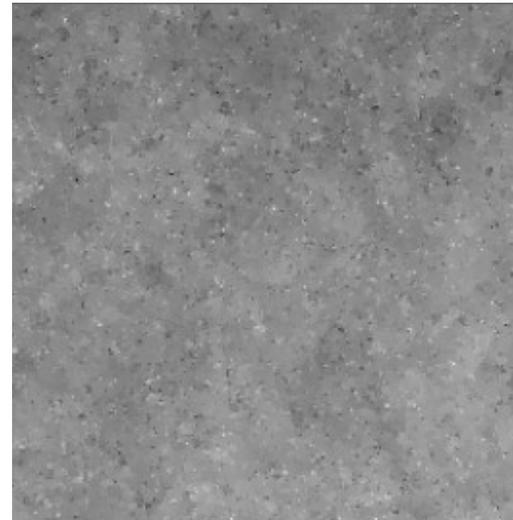


*from Simoncelli 2005*

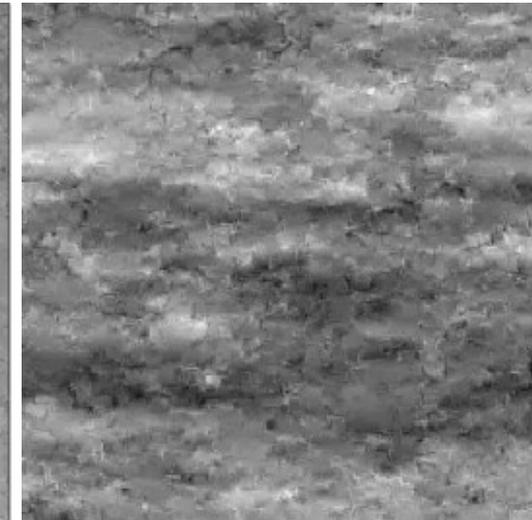
marginal wavelet



Pair-wise MRF



FoE



*from Schmidt, Gao, Roth CVPR 2010*

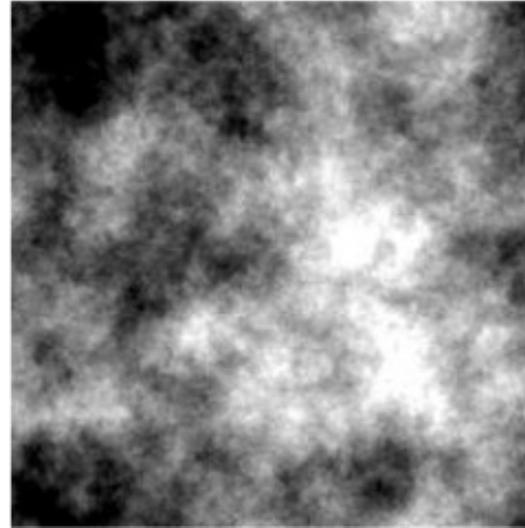
# Sampling High-Resolution Images

**gMRF: 1 layer**



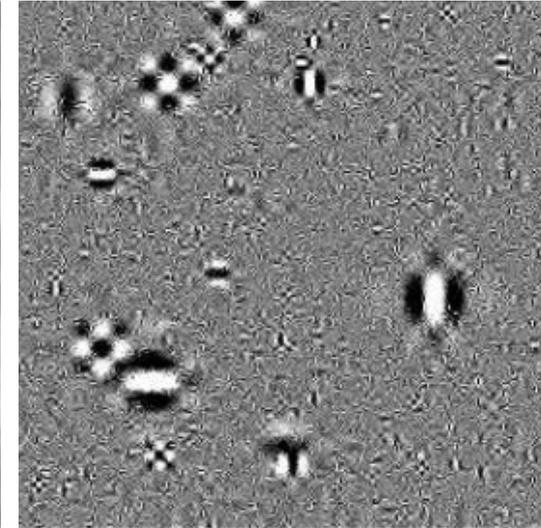
Ranzato et al. PAMI 2013

Gaussian model

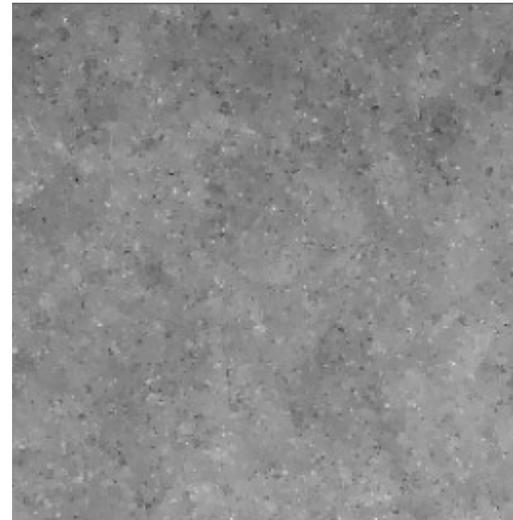


*from Simoncelli 2005*

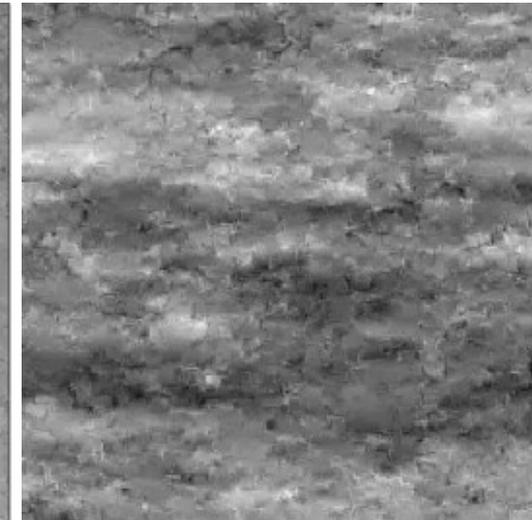
marginal wavelet



Pair-wise MRF



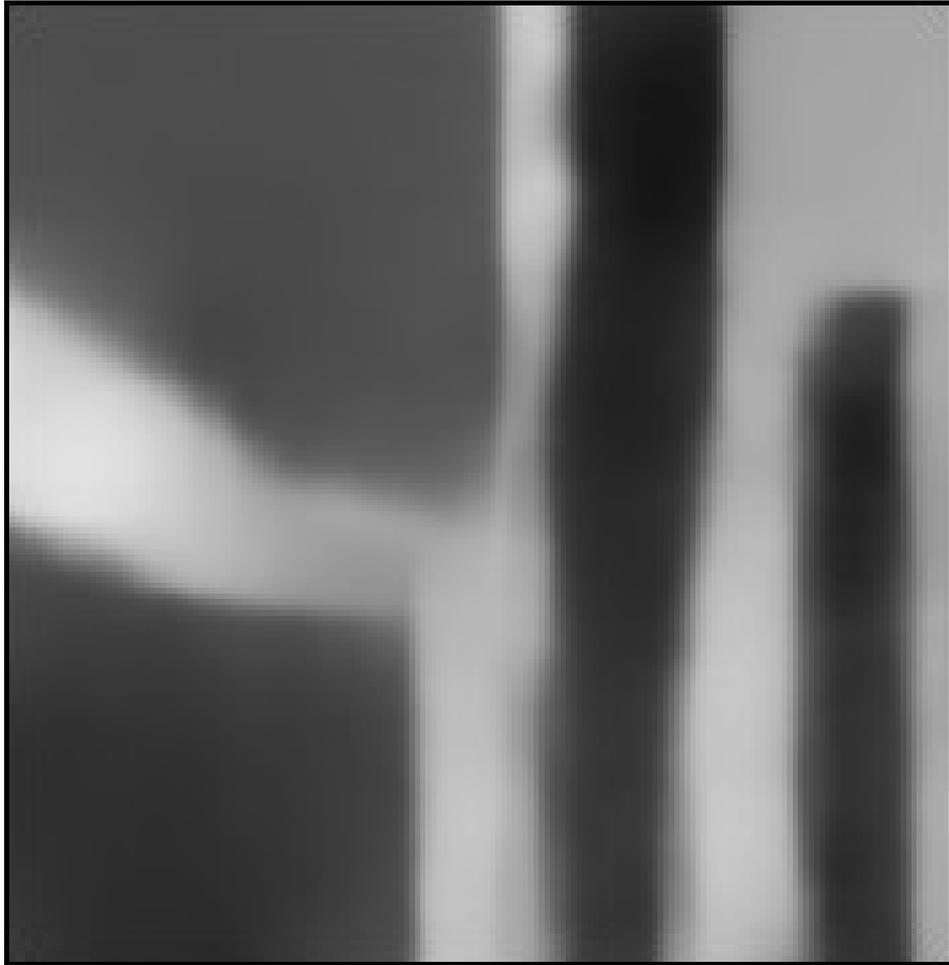
FoE



*from Schmidt, Gao, Roth CVPR 2010*

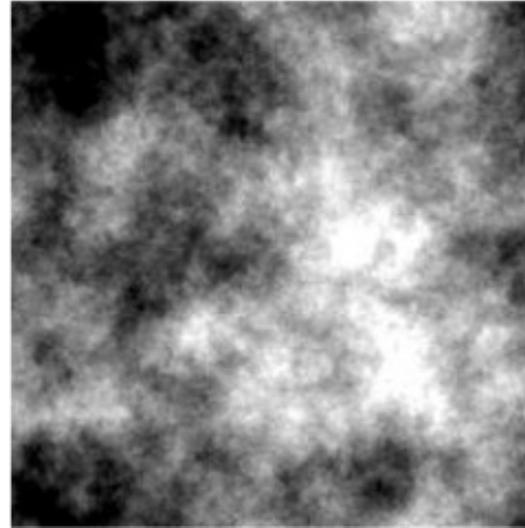
# Sampling High-Resolution Images

gMRF: 3 layer



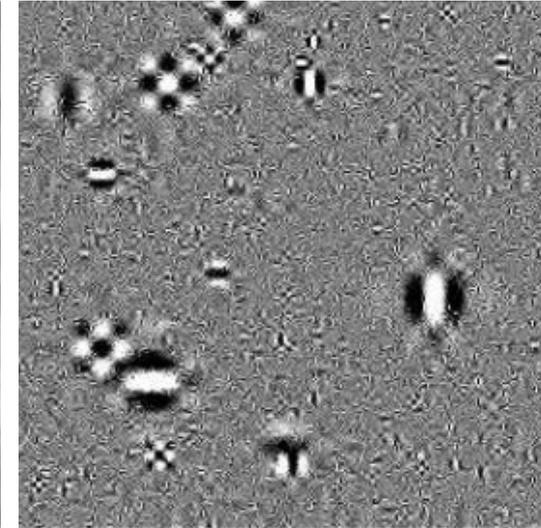
Ranzato et al. PAMI 2013

Gaussian model

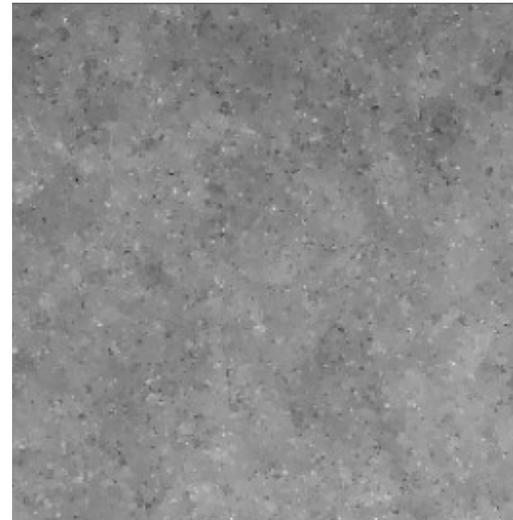


*from Simoncelli 2005*

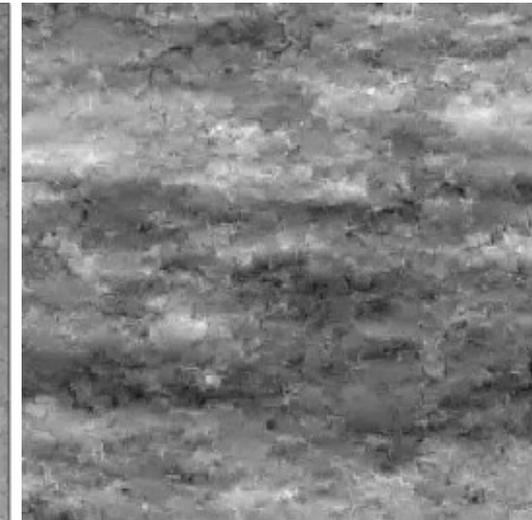
marginal wavelet



Pair-wise MRF



FoE



*from Schmidt, Gao, Roth CVPR 2010*

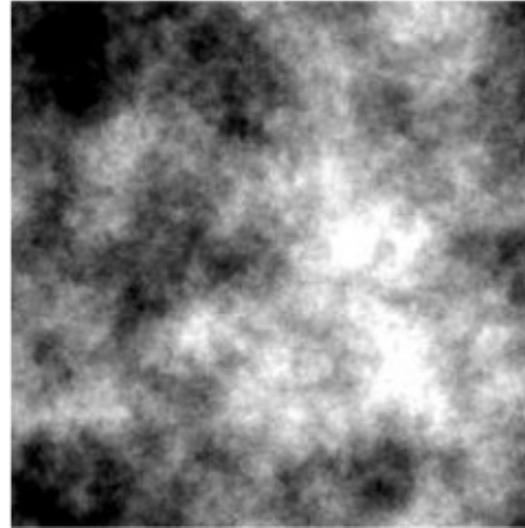
# Sampling High-Resolution Images

gMRF: 3 layer



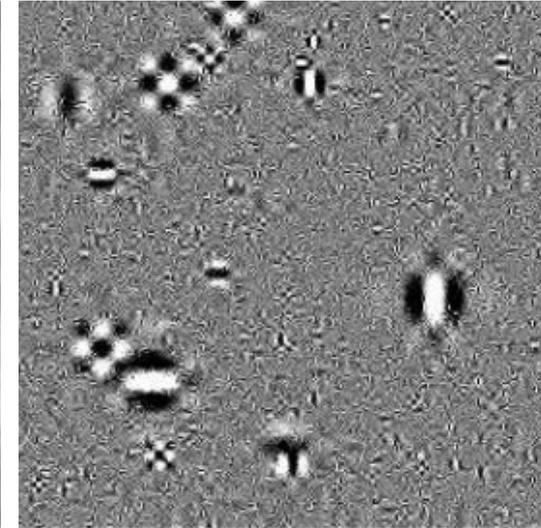
Ranzato et al. PAMI 2013

Gaussian model

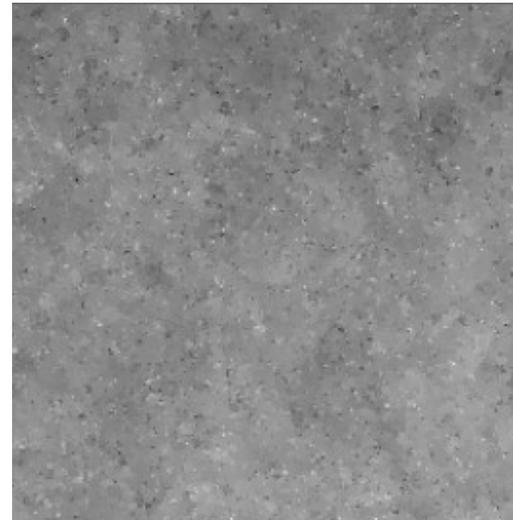


from Simoncelli 2005

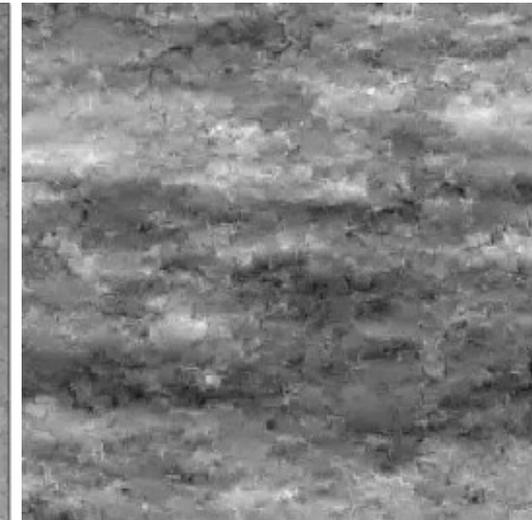
marginal wavelet



Pair-wise MRF



FoE



from Schmidt, Gao, Roth CVPR 2010

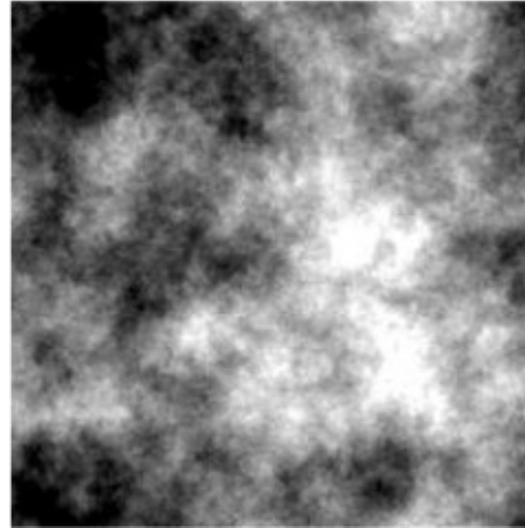
# Sampling High-Resolution Images

gMRF: 3 layer

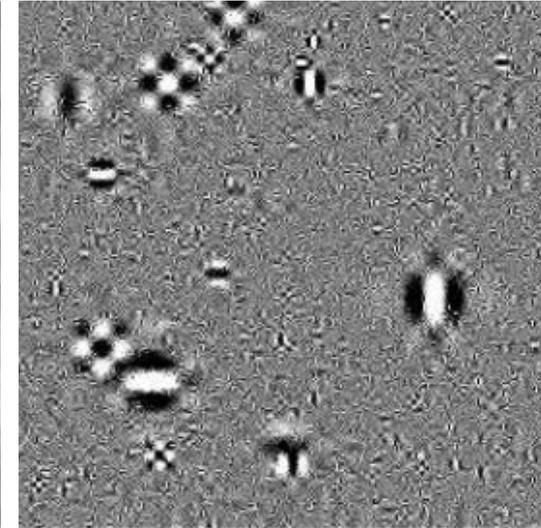


Ranzato et al. PAMI 2013

Gaussian model

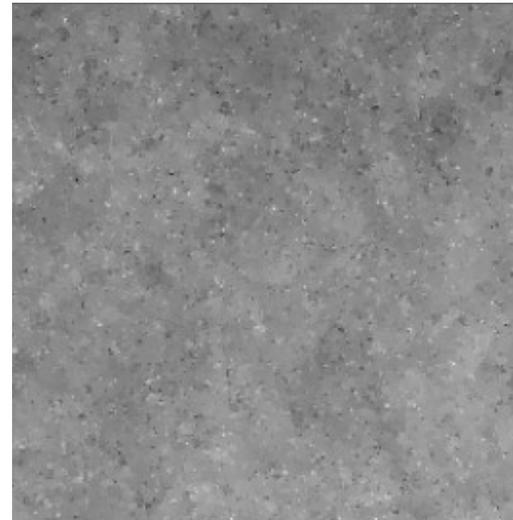


marginal wavelet

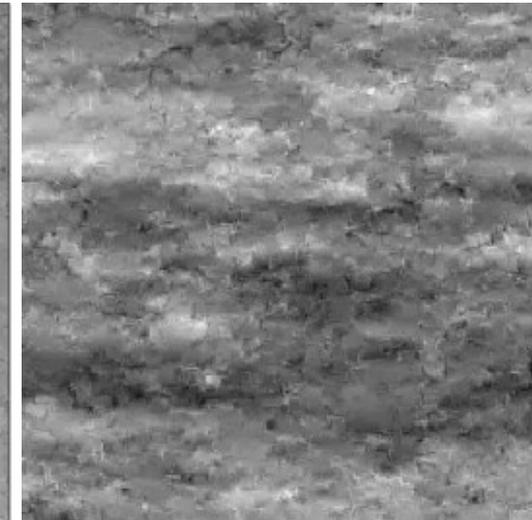


from Simoncelli 2005

Pair-wise MRF



FoE



from Schmidt, Gao, Roth CVPR 2010

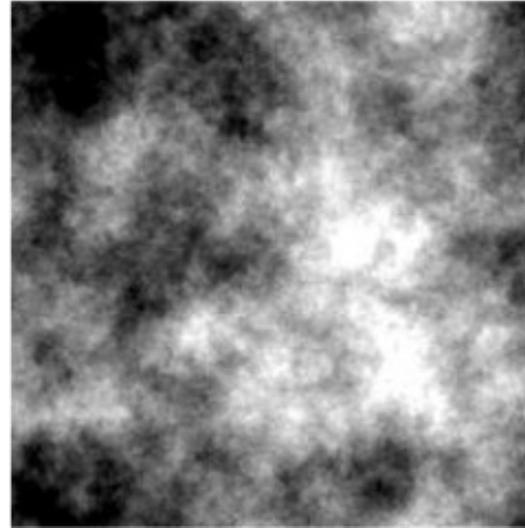
# Sampling High-Resolution Images

gMRF: 3 layer



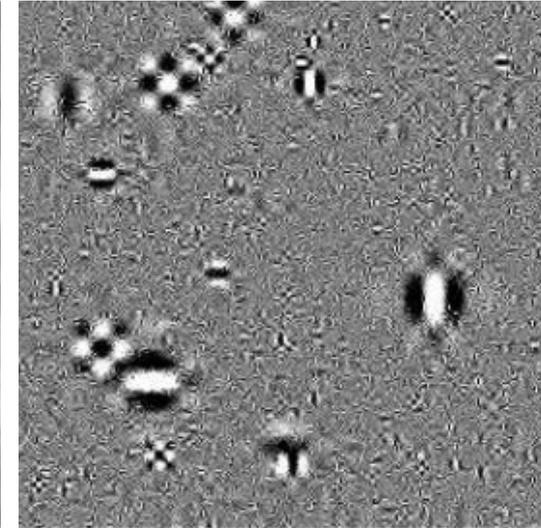
Ranzato et al. PAMI 2013

Gaussian model

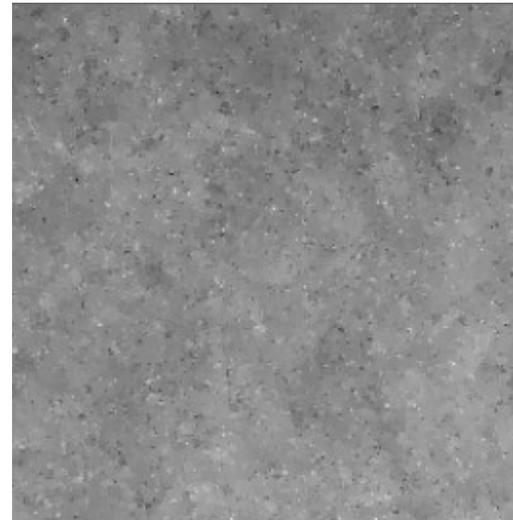


from Simoncelli 2005

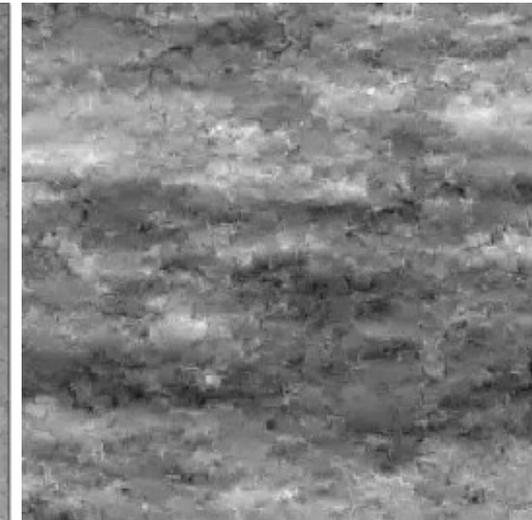
marginal wavelet



Pair-wise MRF

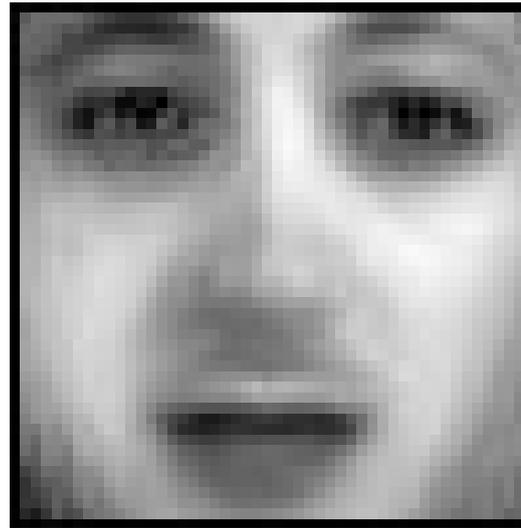


FoE

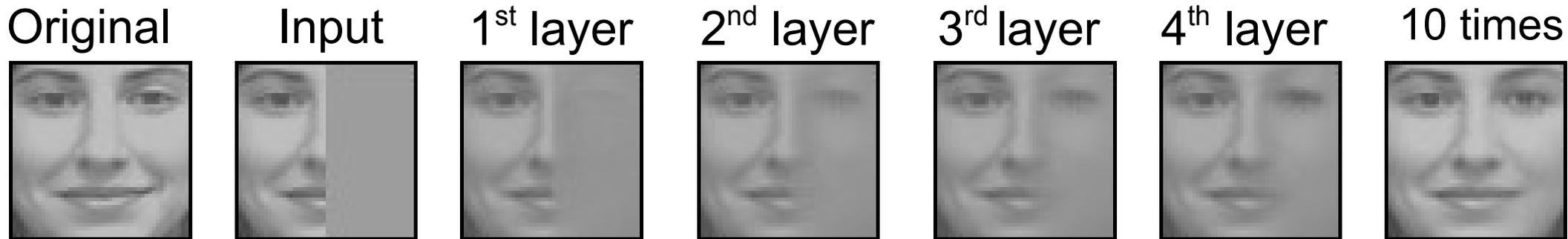


from Schmidt, Gao, Roth CVPR 2010

# Sampling After Training on Face Images

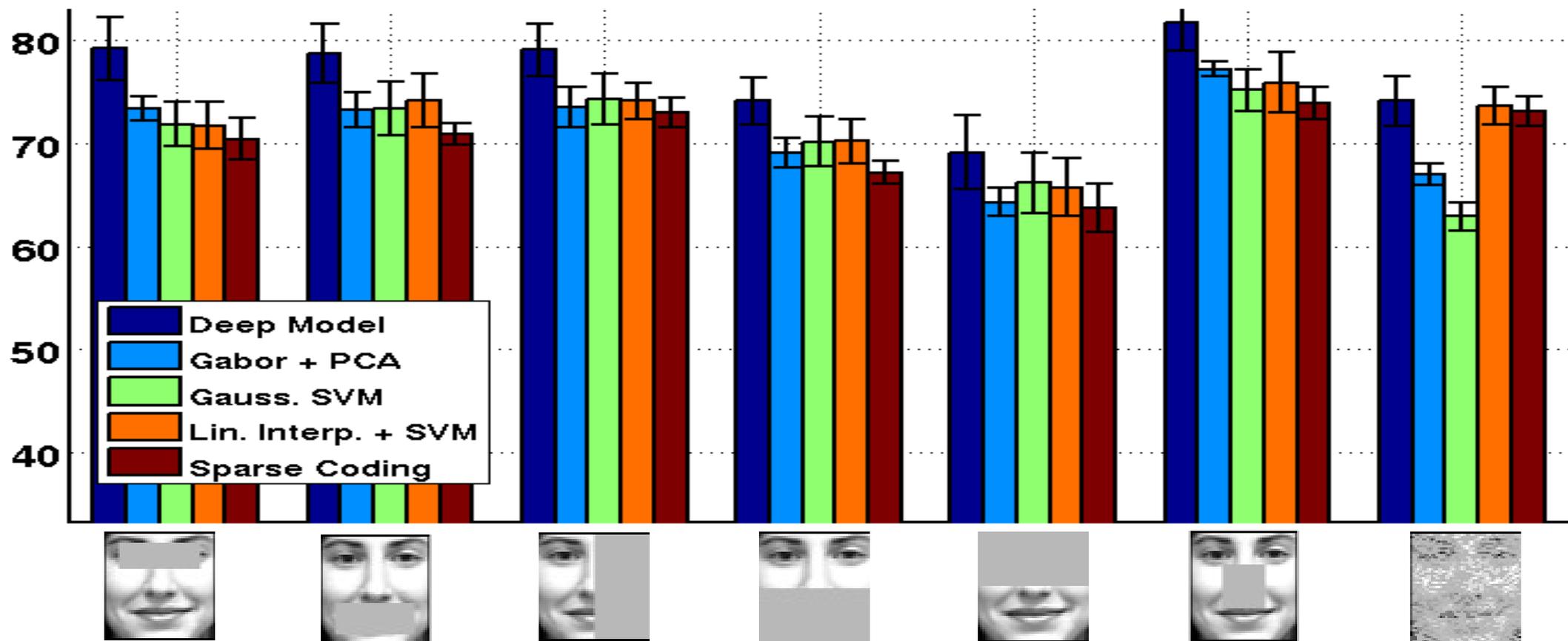


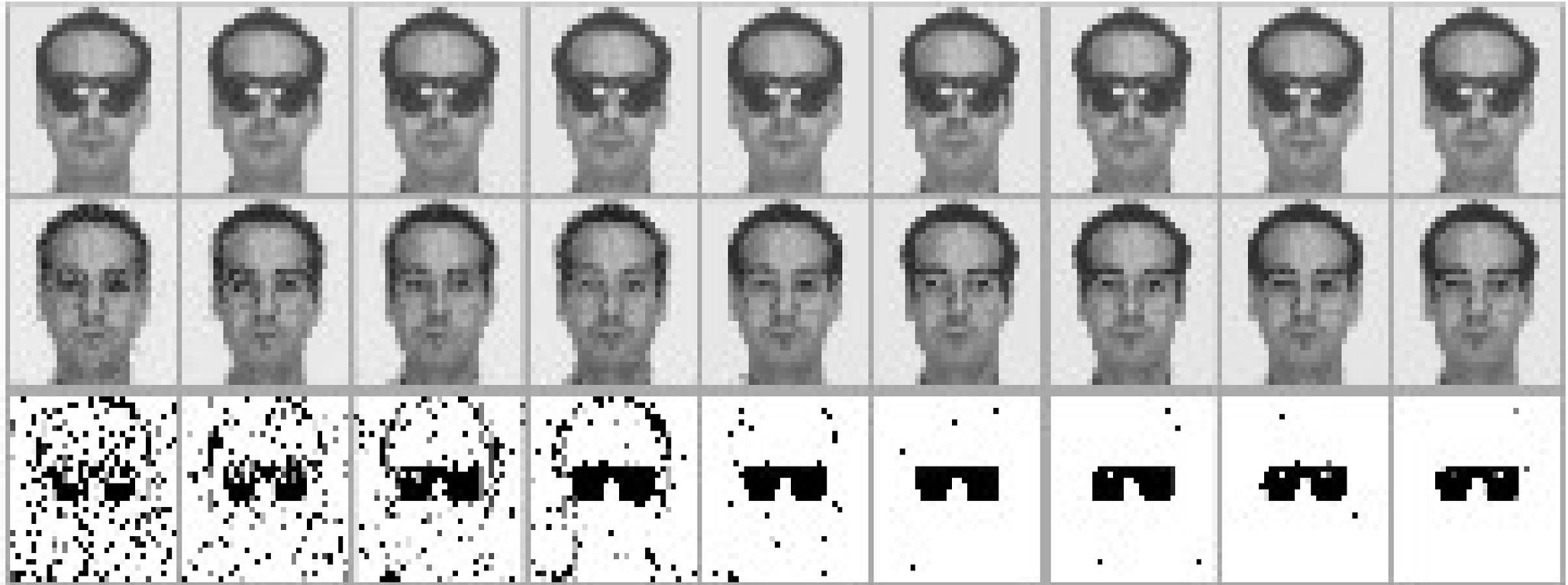
unconstrained samples



conditional (on the left part of the face) samples

# Expression Recognition Under Occlusion





1      3      5      7      10      20      30      40      50

## Pros

- Feature extraction is fast
- Unprecedented generation quality
- Advances models of natural images
- Trains without labeled data

## Cons

- Training is inefficient
  - Slow
  - Tricky
- Sampling scales badly with dimensionality
- What's the use case of generative models?

## Conclusion

- If generation is not required, other feature learning methods are more efficient (e.g., sparse auto-encoders).
- What's the use case of generative models?
- Given enough labeled data, unsup. learning methods have not produced more useful features.

# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- Examples:
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - Unsupervised learning: sparse coding
  - Unsupervised learning: gated MRF
- Other examples
- Practical tricks

# RNNs

recurrent neural network handwriting generation demo

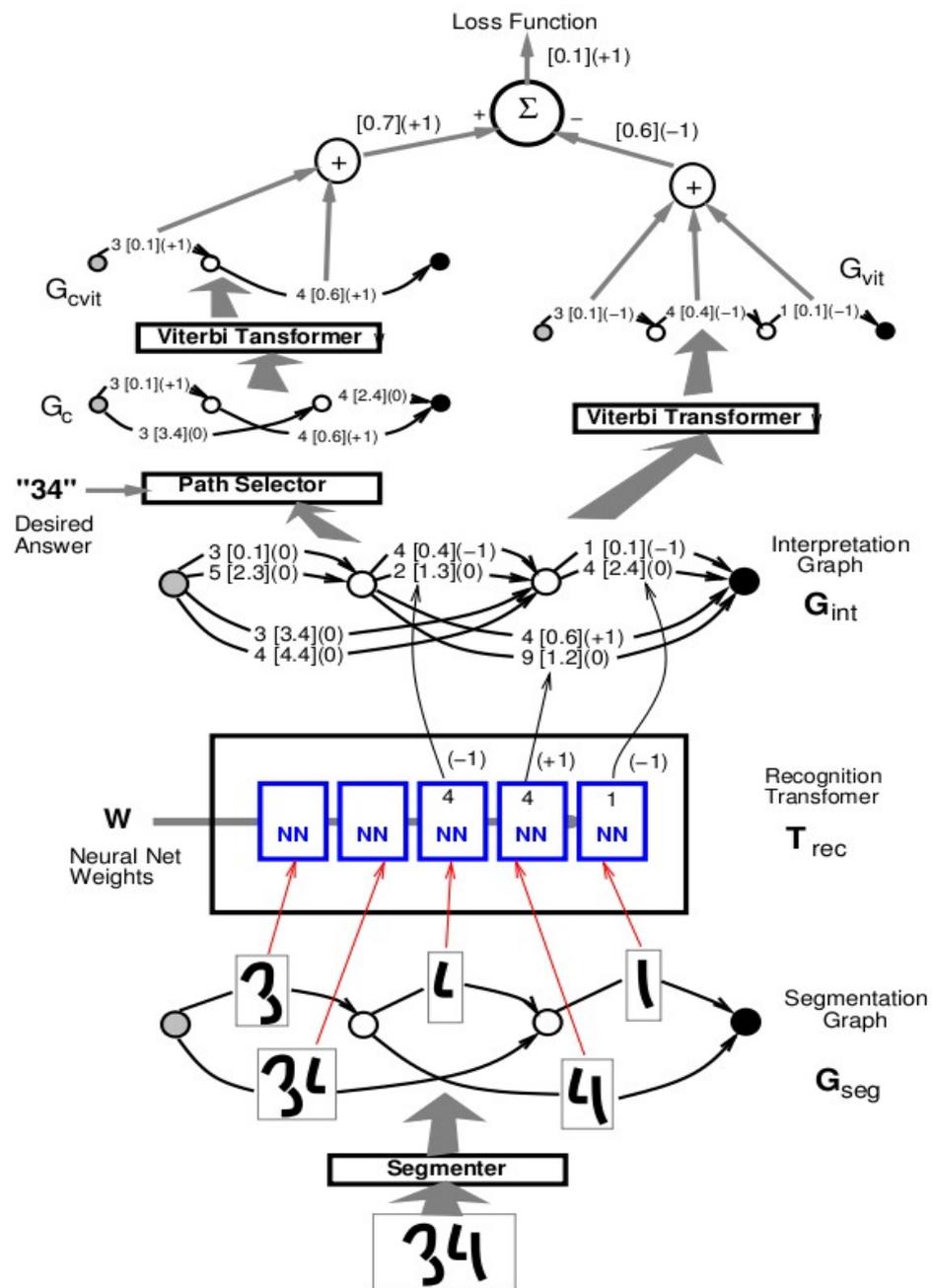
Type a message into the text box, and the network will try to write it out longhand ([this paper](#) explains how it works). Be patient, it can take a while!

**Text** --- up to 100 characters, lower case letters work best

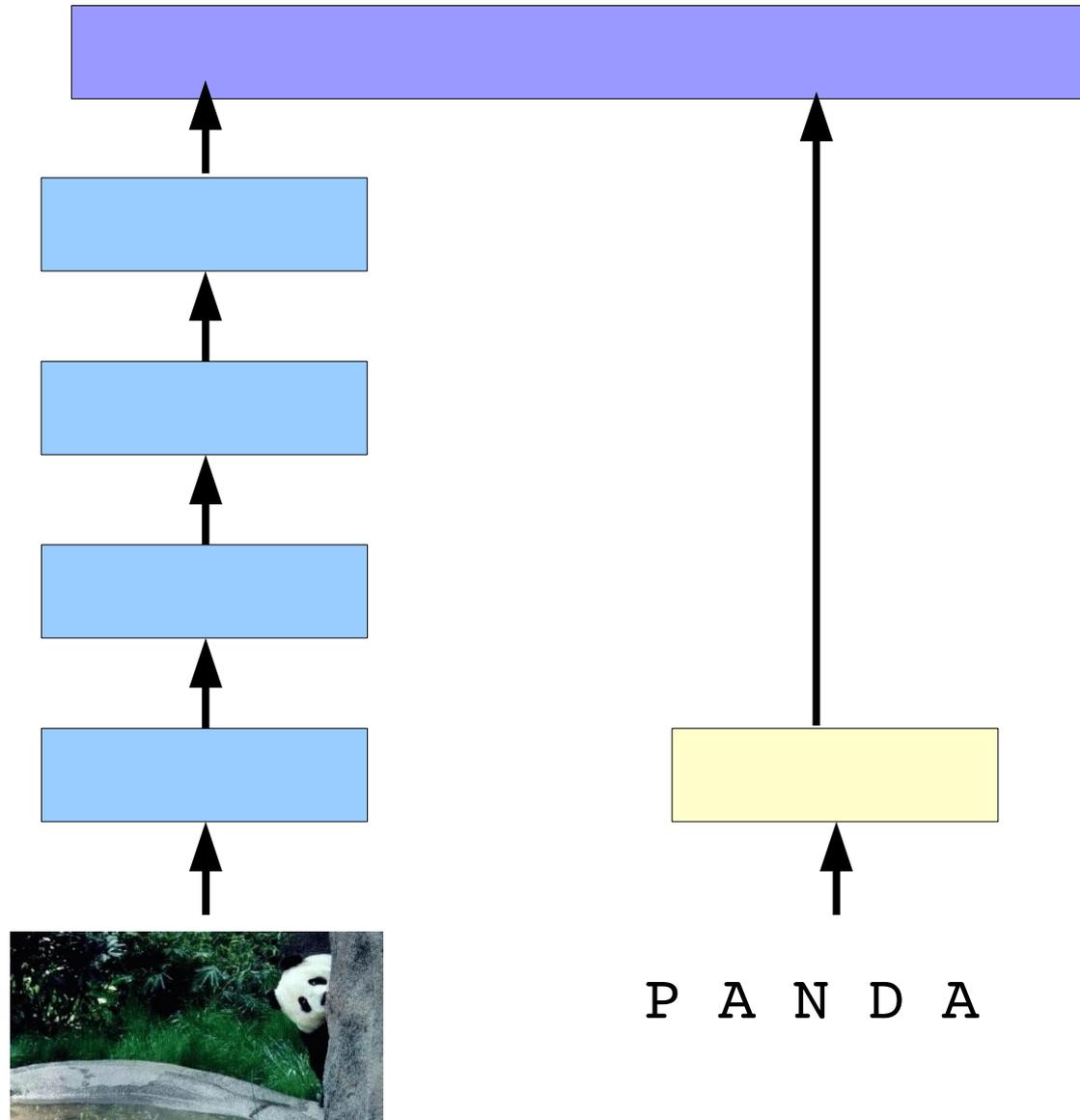
**Style** --- either let the network choose a writing style at random or prime it with a real sequence to make it mimic that writer's style.

- Take the breath away when they are
- He dismissed the idea
- prison welfare Officer complement
- She looked closely as she
- at Humbercombe is being adapted for
- random style

# Structured Prediction



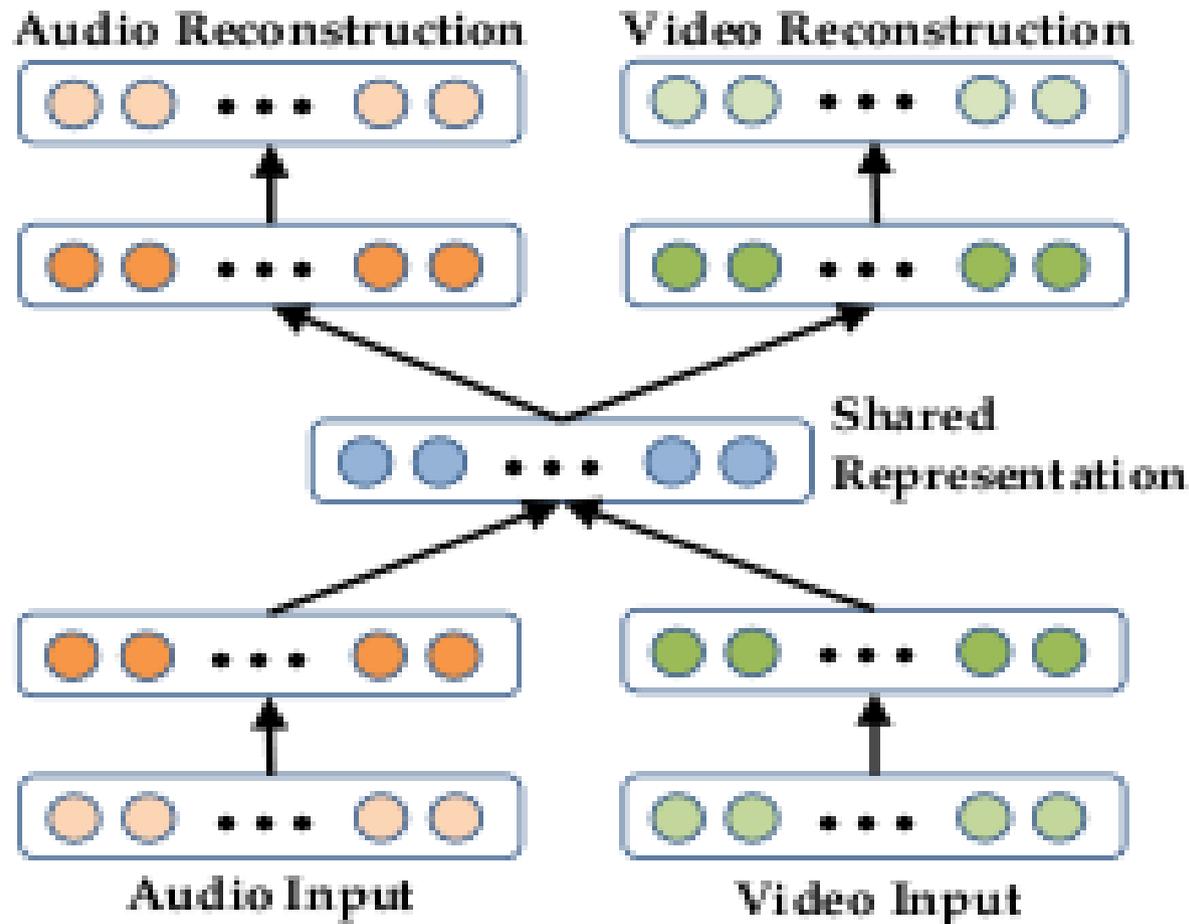
# Multi-Modal Learning



Frome et al. "DeVISE: A deep visual semantic embedding model" NIPS 2013

Socher et al. "Zero-shot learning through cross modal transfer" NIPS 2013

# Multi-Modal Learning



(b) Bimodal Deep Autoencoder

Ngiam et al. "Multimodal deep learning" ICML 2011

Srivastava et al. "Multi-modal learning with DBM" ICML 2012

# Outline

- Theory: Energy-Based Models
  - Energy function
  - Loss function
- Examples:
  - Supervised learning: neural nets
  - Supervised learning: convnets
  - Unsupervised learning: sparse coding
  - Unsupervised learning: gated MRF
- Other examples
- **Practical tricks for CNNs**

# CHOOSING THE ARCHITECTURE

- Task dependent
- Cross-validation
- [Convolution → LCN → pooling]\* + fully connected layer
- The more data: the more layers and the more kernels
  - Look at the number of parameters at each layer
  - Look at the number of flops at each layer
- Computational cost
- Be creative :)

# HOW TO OPTIMIZE

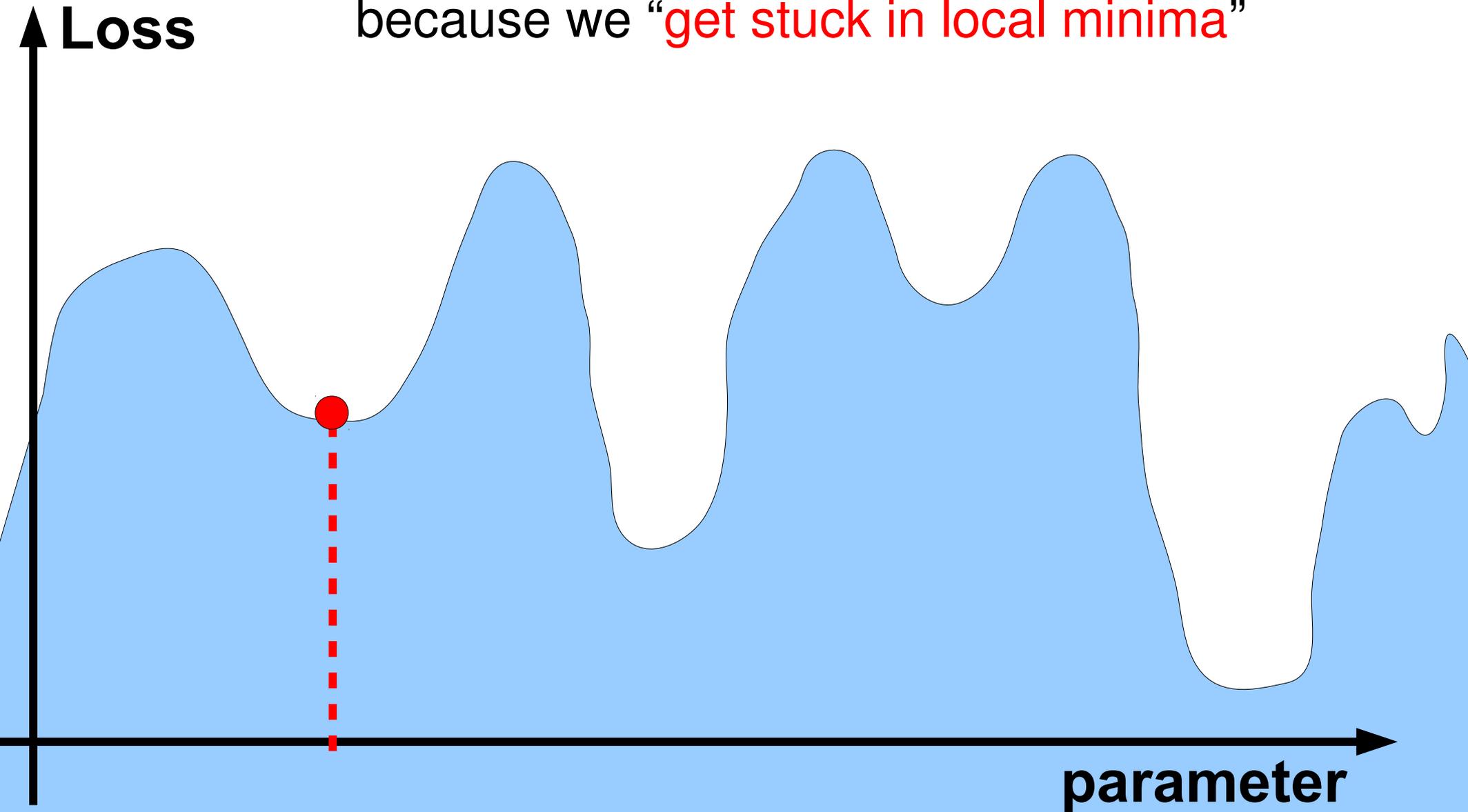
- SGD (with momentum) usually works very well
- Pick learning rate by running on a subset of the data
  - Bottou “Stochastic Gradient Tricks” Neural Networks 2012
    - Start with large learning rate and divide by 2 until loss does not diverge
    - Decay learning rate by a factor of  $\sim 1000$  or more by the end of training
- Use  non-linearity
- Initialize parameters so that each feature across layers has similar variance. Avoid units in saturation.

# HOW TO IMPROVE GENERALIZATION

- Weight sharing (greatly reduce the number of parameters)
- Data augmentation (e.g., jittering, noise injection, etc.)
- Dropout
  - Hinton et al. “Improving Nns by preventing co-adaptation of feature detectors”  
arxiv 2012
- Weight decay (L2, L1)
- Sparsity in the hidden units
- Multi-task (unsupervised learning)

# ConvNets: till 2012

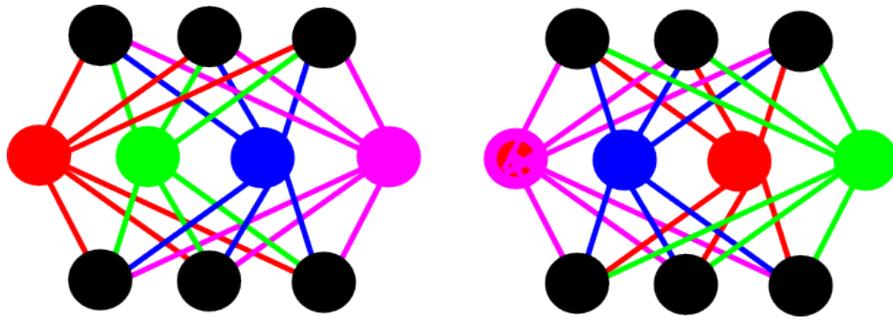
Common wisdom: training does not work because we “get stuck in local minima”



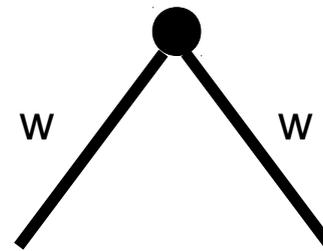
# ConvNets: today

Local minima are all similar, there are long plateaus, it can take long time to break symmetries.

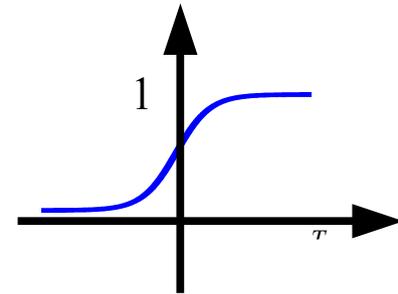
Loss



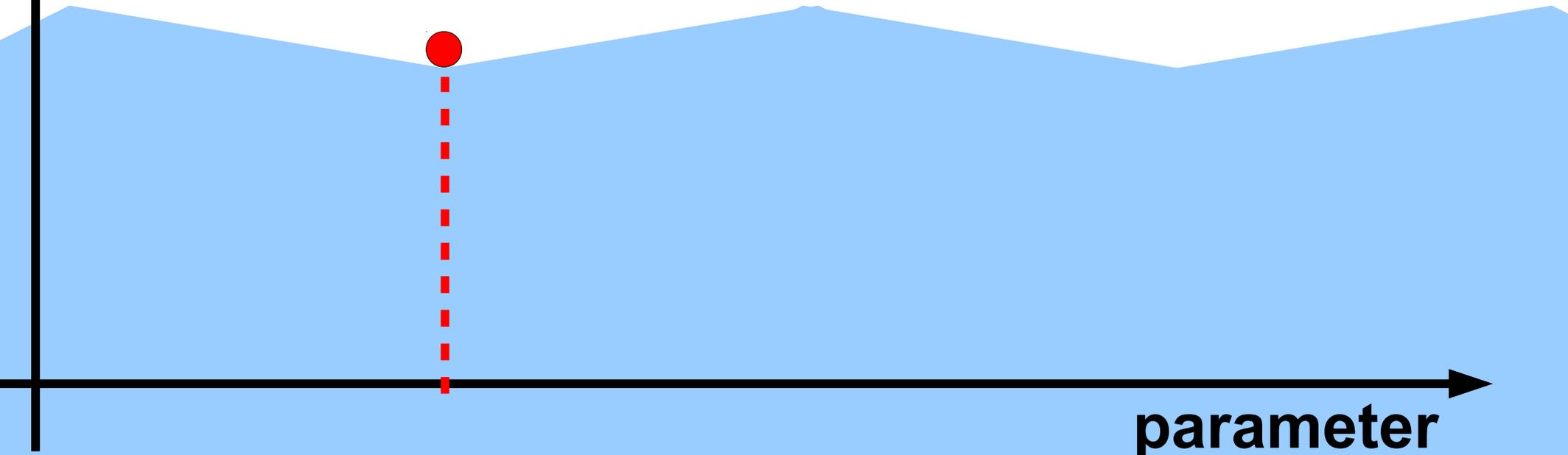
input/output invariant to permutations



breaking ties between parameters



Saturating units



parameter

# Neural Net Optimization is...

Like walking on a ridge between valleys

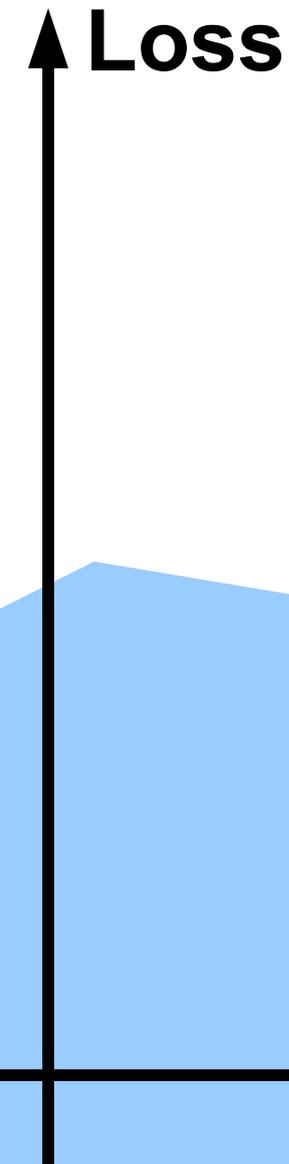


# ConvNets: today

Local minima are all similar, there are long plateaus, it can take long to break symmetries.

Optimization is not the real problem when:

- dataset is large
- unit do not saturate too much
- normalization layer



parameter

# ConvNets: today

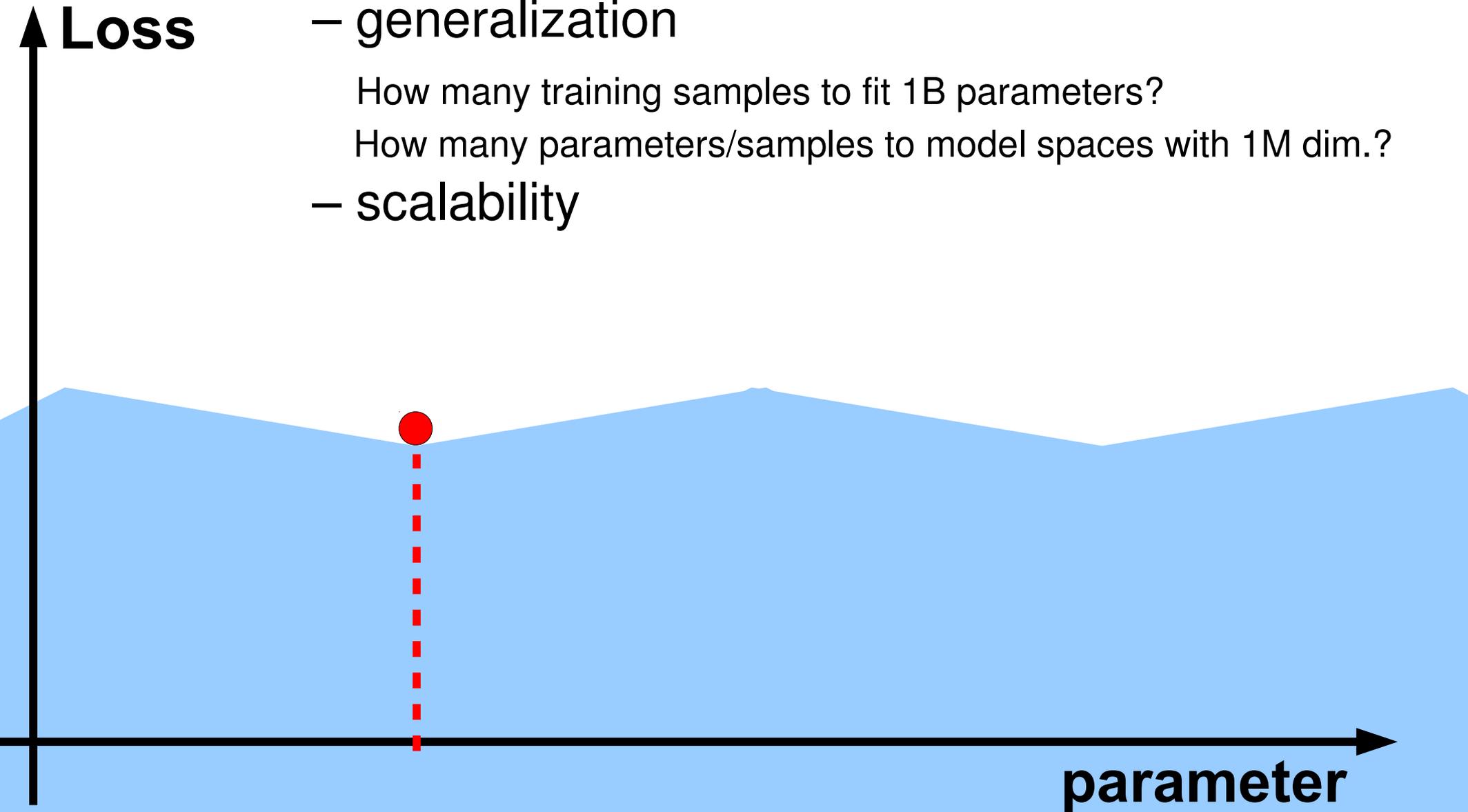
Today's belief is that the challenge is about:

- generalization

How many training samples to fit 1B parameters?

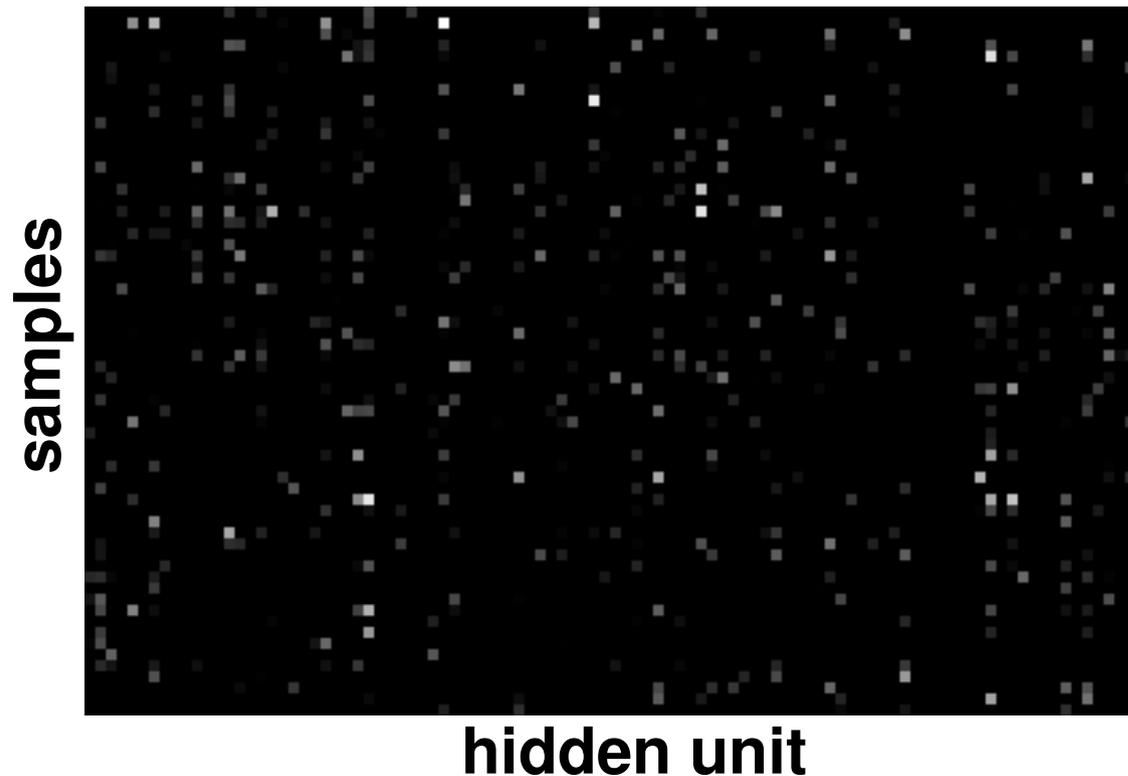
How many parameters/samples to model spaces with 1M dim.?

- scalability



# OTHER THINGS GOOD TO KNOW

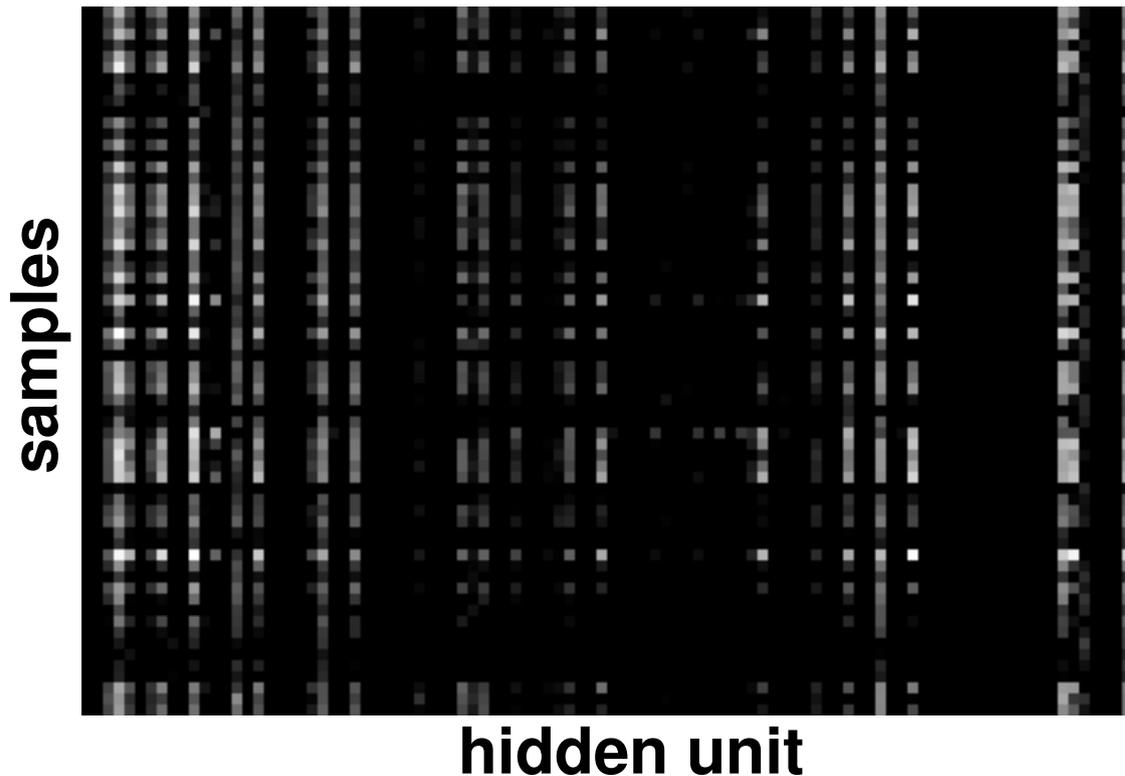
- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.



**Good training:** hidden units are sparse across samples and across features.

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.

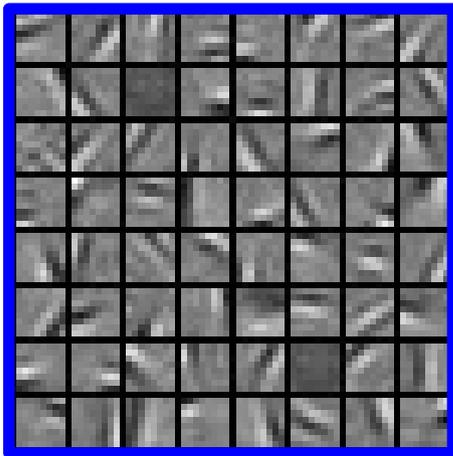


**Bad training:** many hidden units ignore the input and/or exhibit strong correlations.

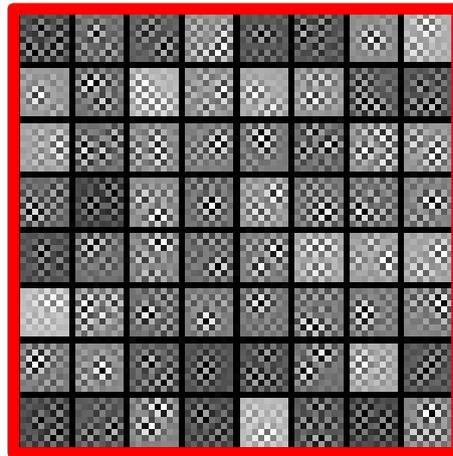
# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.
- Visualize parameters

GOOD

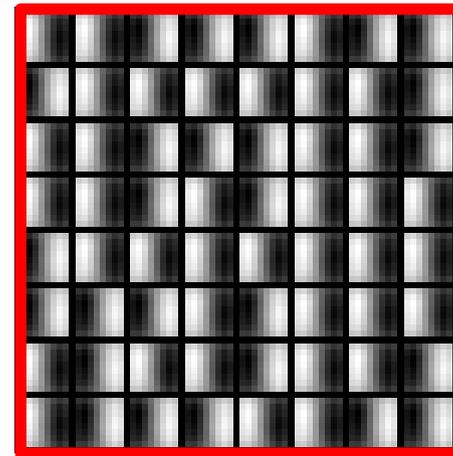


BAD



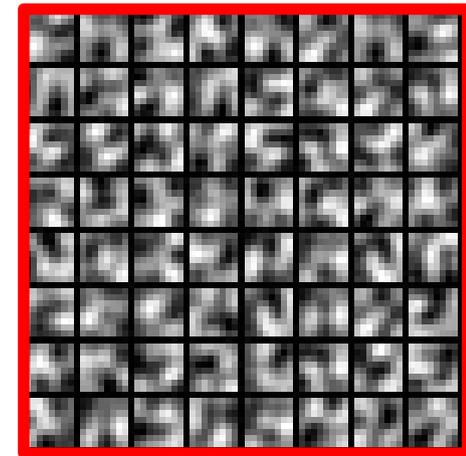
too noisy

BAD



too correlated

BAD



lack structure

**Good training:** learned filters exhibit structure and are uncorrelated.

# OTHER THINGS GOOD TO KNOW

- Check gradients numerically by finite differences
- Visualize features (feature maps need to be uncorrelated) and have high variance.
- Visualize parameters
- Measure error on both training and validation set.
- Test on a small subset of the data and check the error  $\rightarrow 0$ .

# WHAT IF IT DOES NOT WORK?

- Training diverges:
  - Learning rate may be too large → decrease learning rate
  - BPROP is buggy → numerical gradient checking
- Parameters collapse / loss is minimized but accuracy is low
  - Check loss function:
    - Is it appropriate for the task you want to solve?
    - Does it have degenerate solutions? Check “pull-up” term.
- Network is underperforming
  - Compute flops and nr. params. → if too small, make net larger
  - Visualize hidden units/params → fix optimization
- Network is too slow
  - Compute flops and nr. params. → GPU, distrib. framework, make net smaller

# SUMMARY

- Deep Learning = Learning Hierarchical representations. Leverage compositionality to gain efficiency.
- Unsupervised learning: active research topic.
- Supervised learning: most successful set up today.
- Optimization
  - Don't we get stuck in local minima? No, they are all the same!
  - In large scale applications, local minima are even less of an issue.
- Scaling
  - GPUs
  - Distributed framework (Google)
  - Better optimization techniques
- Generalization on small datasets (curse of dimensionality):
  - Input distortions
  - weight decay
  - dropout

# THANK YOU!

**NOTE: IJCV Special Issue on Deep Learning.  
Deadline: 9 Feb. 2014.**

# SOFTWARE

## **Torch7: learning library that supports neural net training**

---

<http://www.torch.ch>

<http://code.cogbits.com/wiki/doku.php> (tutorial with demos by C. Farabet)

## **Python-based learning library (U. Montreal)**

---

- <http://deeplearning.net/software/theano/> (does automatic differentiation)

## **Efficient CUDA kernels for ConvNets (Krizhevsky)**

---

– [code.google.com/p/cuda-convnet](http://code.google.com/p/cuda-convnet)

## **Caffe (Yangqing Jia)**

---

– <http://caffe.berkeleyvision.org>

# REFERENCES

## Convolutional Nets

- LeCun, Bottou, Bengio and Haffner: Gradient-Based Learning Applied to Document Recognition, Proceedings of the IEEE, 86(11):2278-2324, November 1998
- Krizhevsky, Sutskever, Hinton “ImageNet Classification with deep convolutional neural networks” NIPS 2012
- Jarrett, Kavukcuoglu, Ranzato, LeCun: What is the Best Multi-Stage Architecture for Object Recognition?, Proc. International Conference on Computer Vision (ICCV'09), IEEE, 2009
- Kavukcuoglu, Sermanet, Boureau, Gregor, Mathieu, LeCun: Learning Convolutional Feature Hierarchies for Visual Recognition, Advances in Neural Information Processing Systems (NIPS 2010), 23, 2010
- see [yann.lecun.com/exdb/publis](http://yann.lecun.com/exdb/publis) for references on many different kinds of convnets.
- see <http://www.cmap.polytechnique.fr/scattering/> for scattering networks (similar to convnets but with less learning and stronger mathematical foundations)
- see <http://www.idsia.ch/~juergen/> for other references to ConvNets and LSTMs.

# REFERENCES

## Applications of Convolutional Nets

- Farabet, Couprie, Najman, LeCun. Scene Parsing with Multiscale Feature Learning, Purity Trees, and Optimal Covers”, ICML 2012
- Pierre Sermanet, Koray Kavukcuoglu, Soumith Chintala and Yann LeCun: Pedestrian Detection with Unsupervised Multi-Stage Feature Learning, CVPR 2013
- D. Ciresan, A. Giusti, L. Gambardella, J. Schmidhuber. Deep Neural Networks Segment Neuronal Membranes in Electron Microscopy Images. NIPS 2012
- Raia Hadsell, Pierre Sermanet, Marco Scoffier, Ayse Erkan, Koray Kavackuoglu, Urs Muller and Yann LeCun. Learning Long-Range Vision for Autonomous Off-Road Driving, Journal of Field Robotics, 26(2):120-144, 2009
- Burger, Schuler, Harmeling. Image Denoisng: Can Plain Neural Networks Compete with BM3D?, CVPR 2012
- Hadsell, Chopra, LeCun. Dimensionality reduction by learning an invariant mapping, CVPR 2006
- Bergstra et al. Making a science of model search: hyperparameter optimization in hundred of dimensions for vision architectures, ICML 2013

# REFERENCES

## Deep Learning in general

---

- deep learning tutorial slides at ICML 2013
- Yoshua Bengio, Learning Deep Architectures for AI, Foundations and Trends in Machine Learning, 2(1), pp.1-127, 2009.
- LeCun, Chopra, Hadsell, Ranzato, Huang: A Tutorial on Energy-Based Learning, in Bakir, G. and Hofman, T. and Schölkopf, B. and Smola, A. and Taskar, B. (Eds), Predicting Structured Data, MIT Press, 2006

# THANK YOU

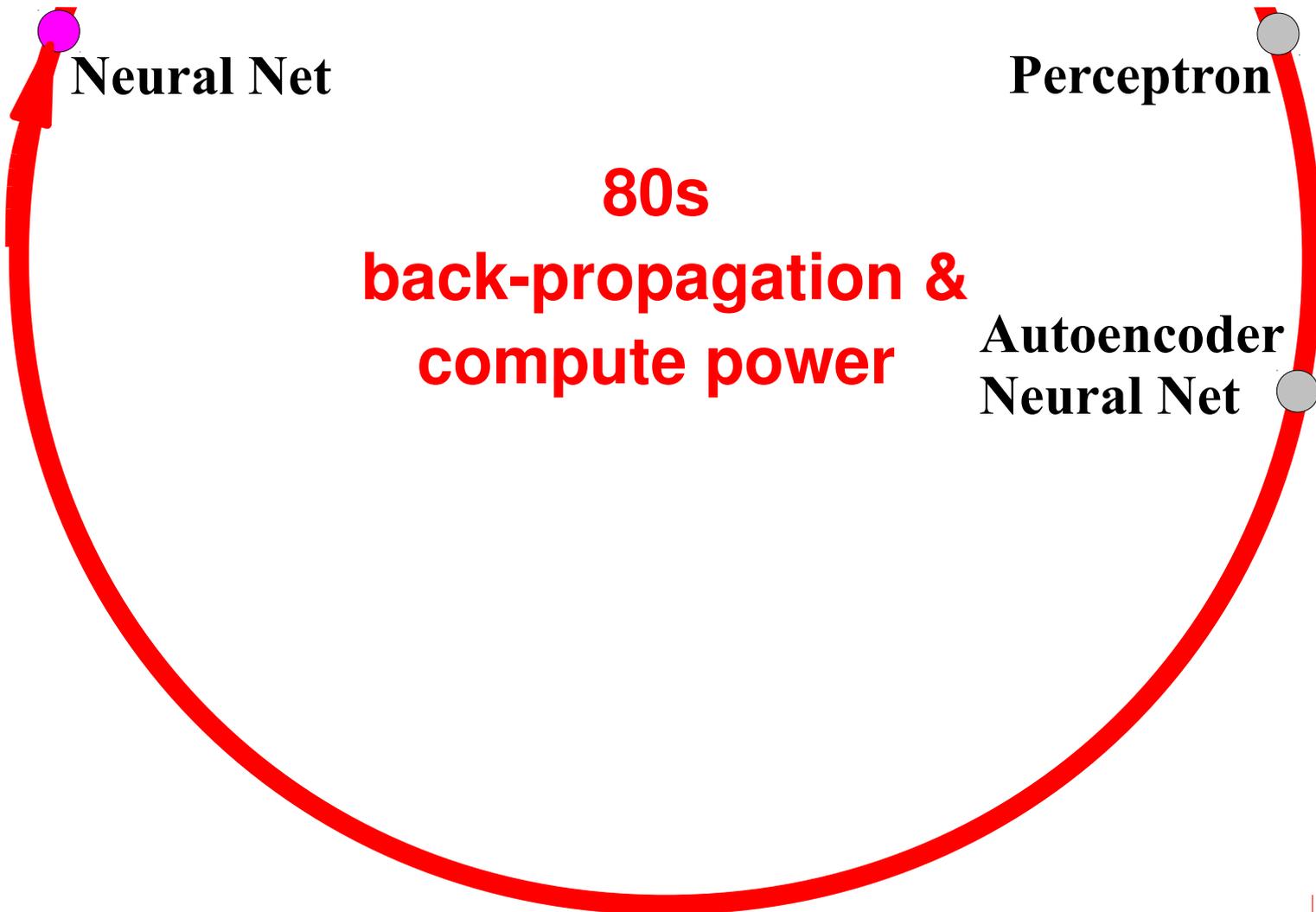
*Acknowledgements to Yann LeCun. Many slides from ICML 2013 and CVPR 2013 tutorial on deep learning.*

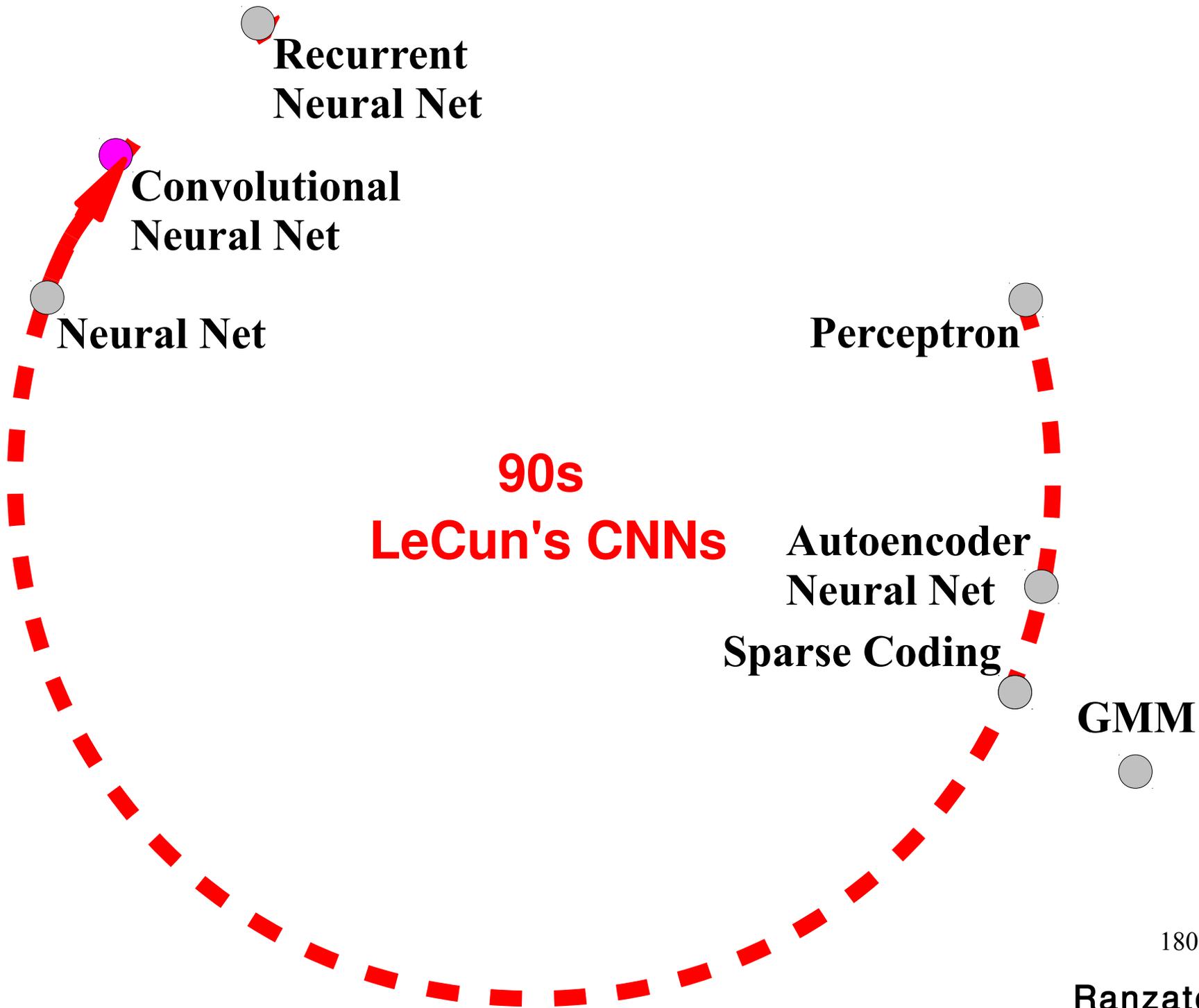
*IJCV SPECIAL ISSUE ON DEEP LEARNING: 9 February 2014.*

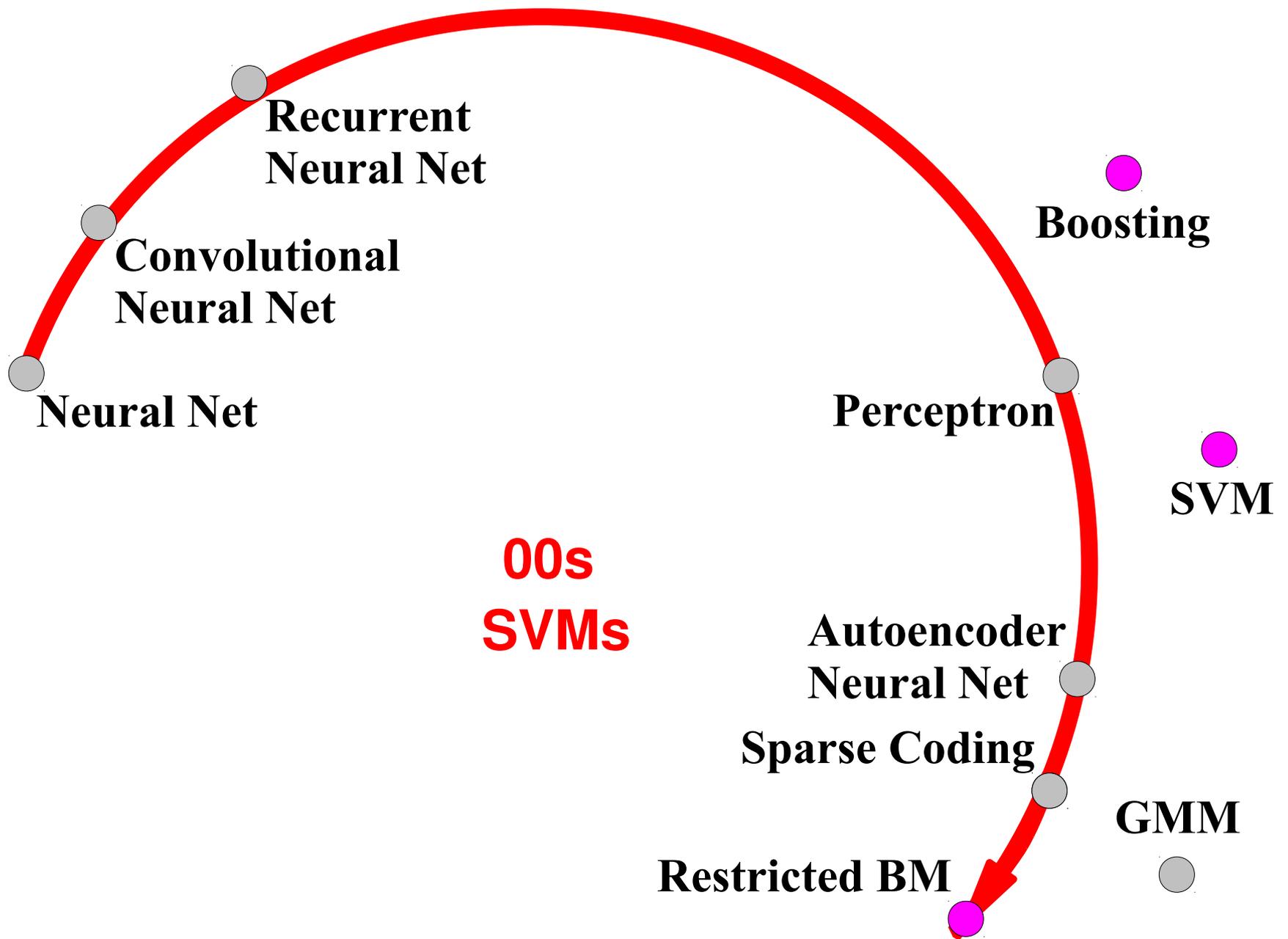
**1957**  
**Rosenblatt**

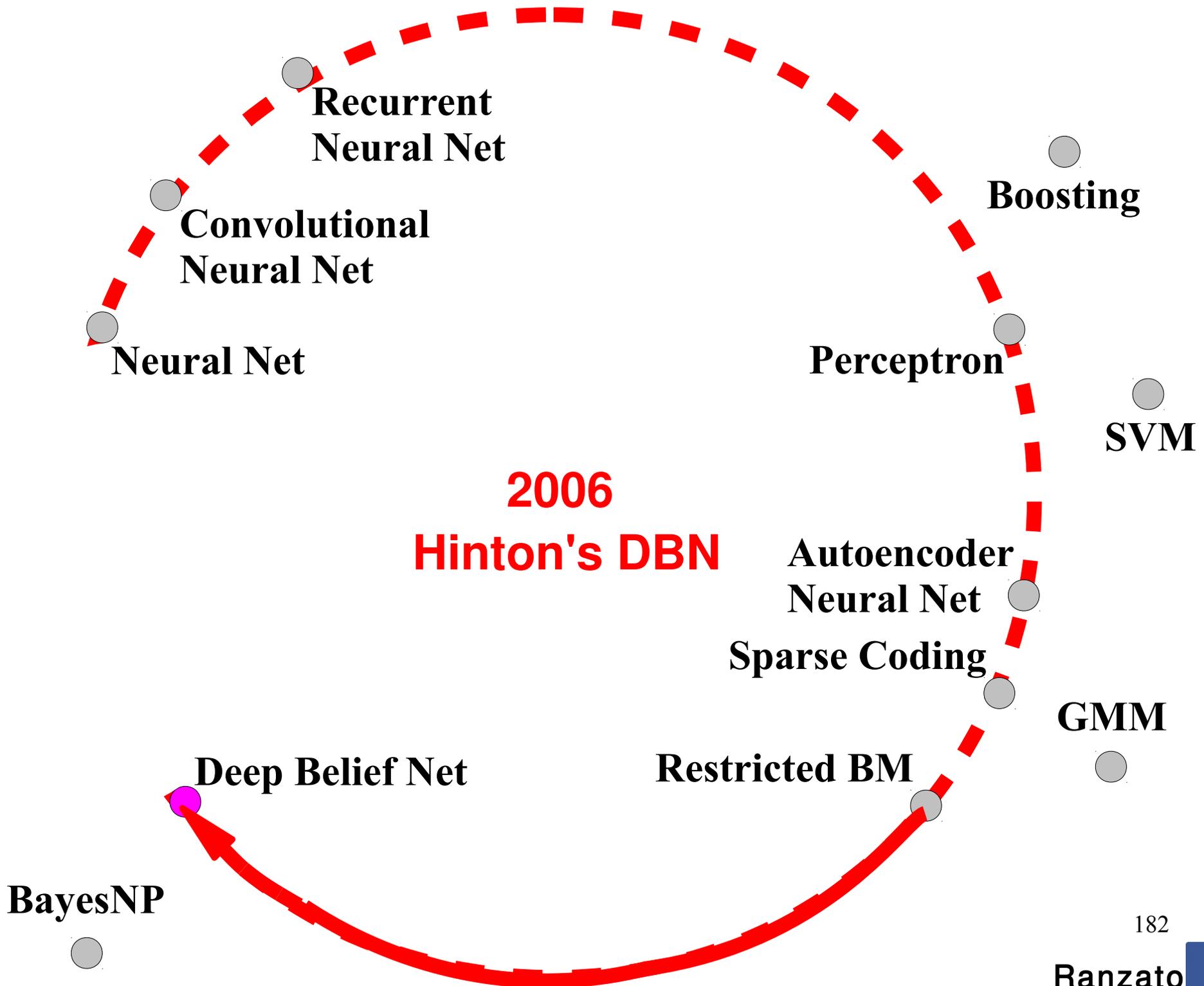
Perceptron 

# **THE SPACE OF MACHINE LEARNING METHODS**









$\Sigma\Pi$

BayesNP

Deep Belief Net

Deep (sparse/denoising)  
Autoencoder

Neural Net

Convolutional  
Neural Net

Recurrent  
Neural Net

2009

ASR (data + GPU)

Restricted BM

Sparse Coding

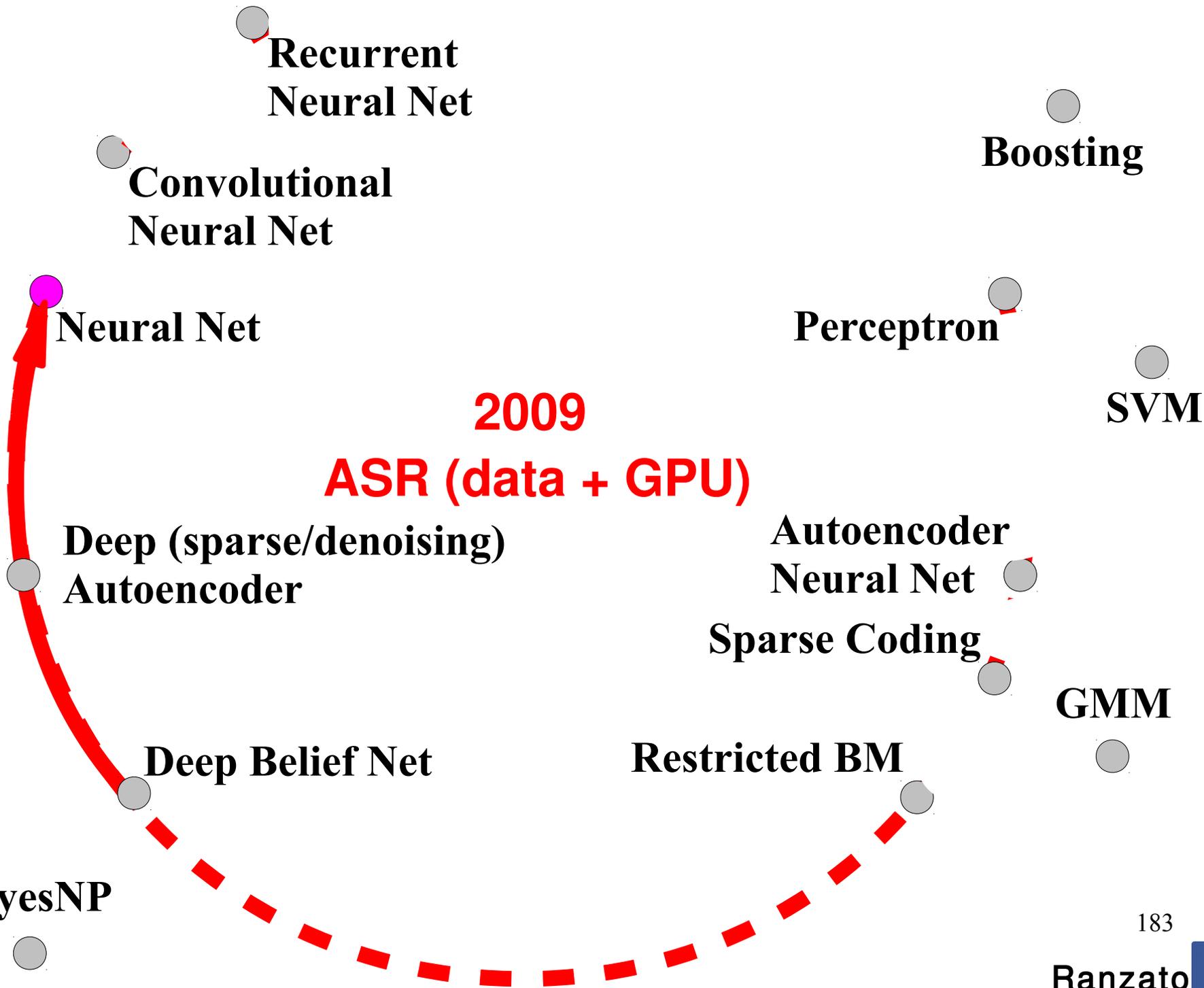
Autoencoder  
Neural Net

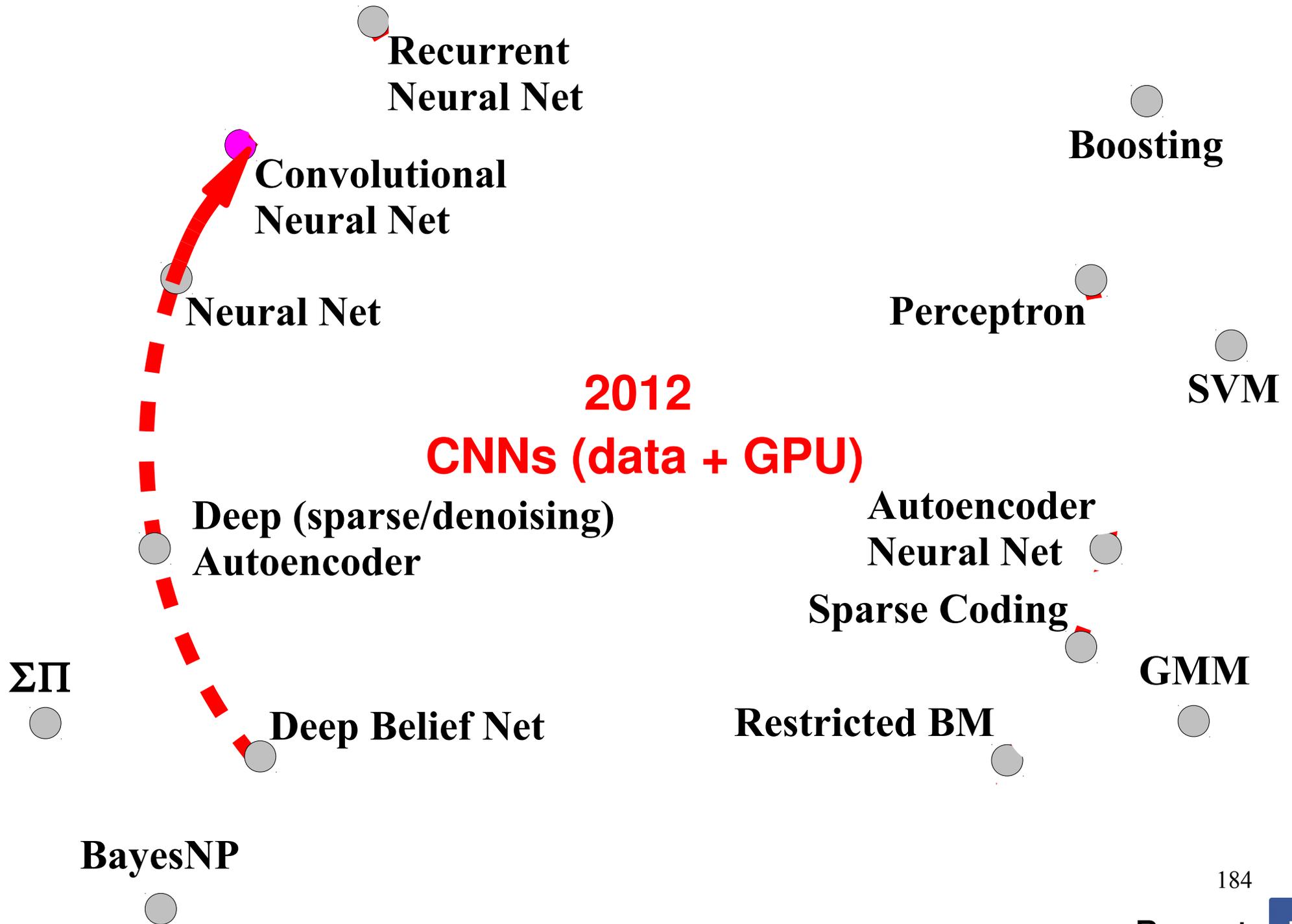
Perceptron

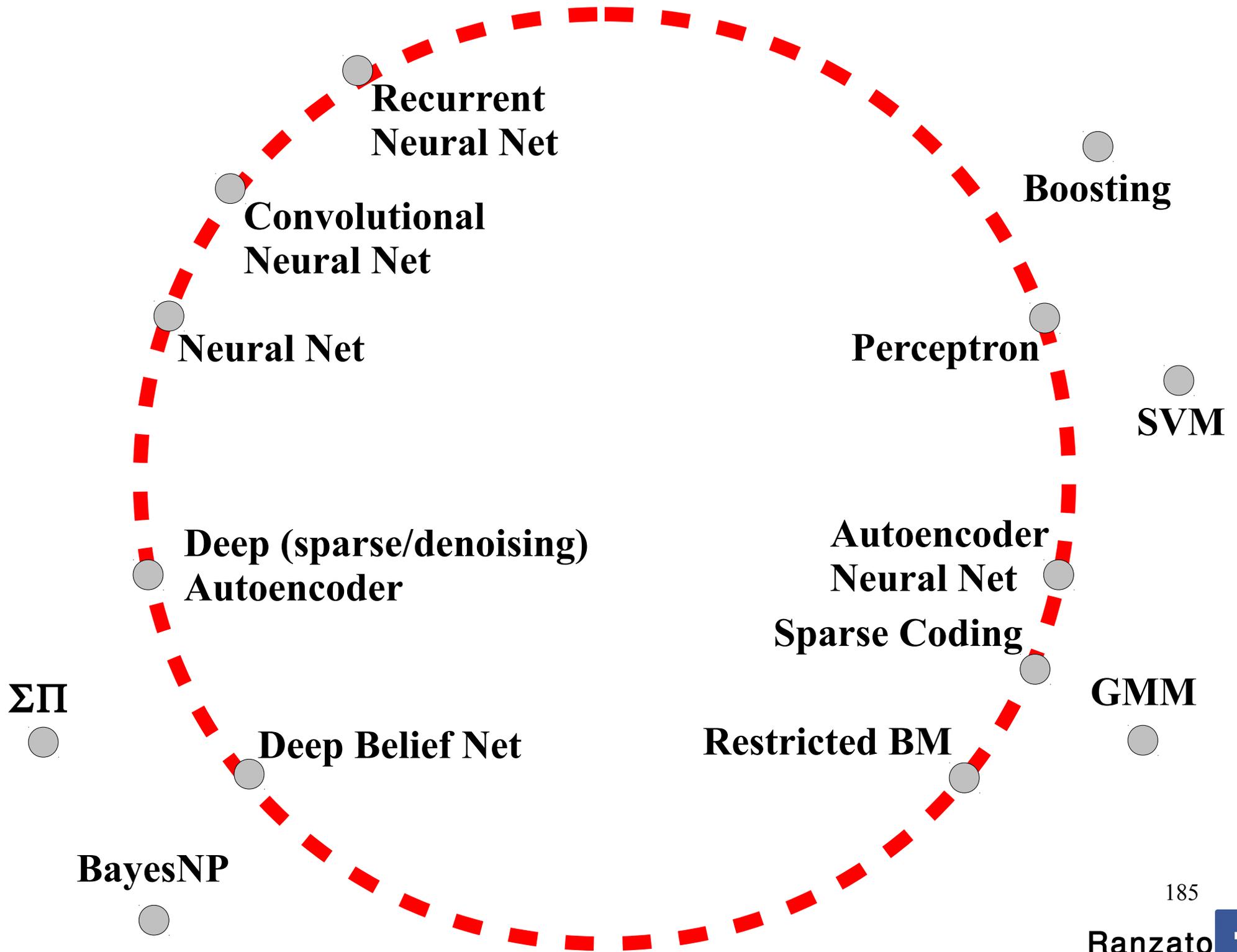
SVM

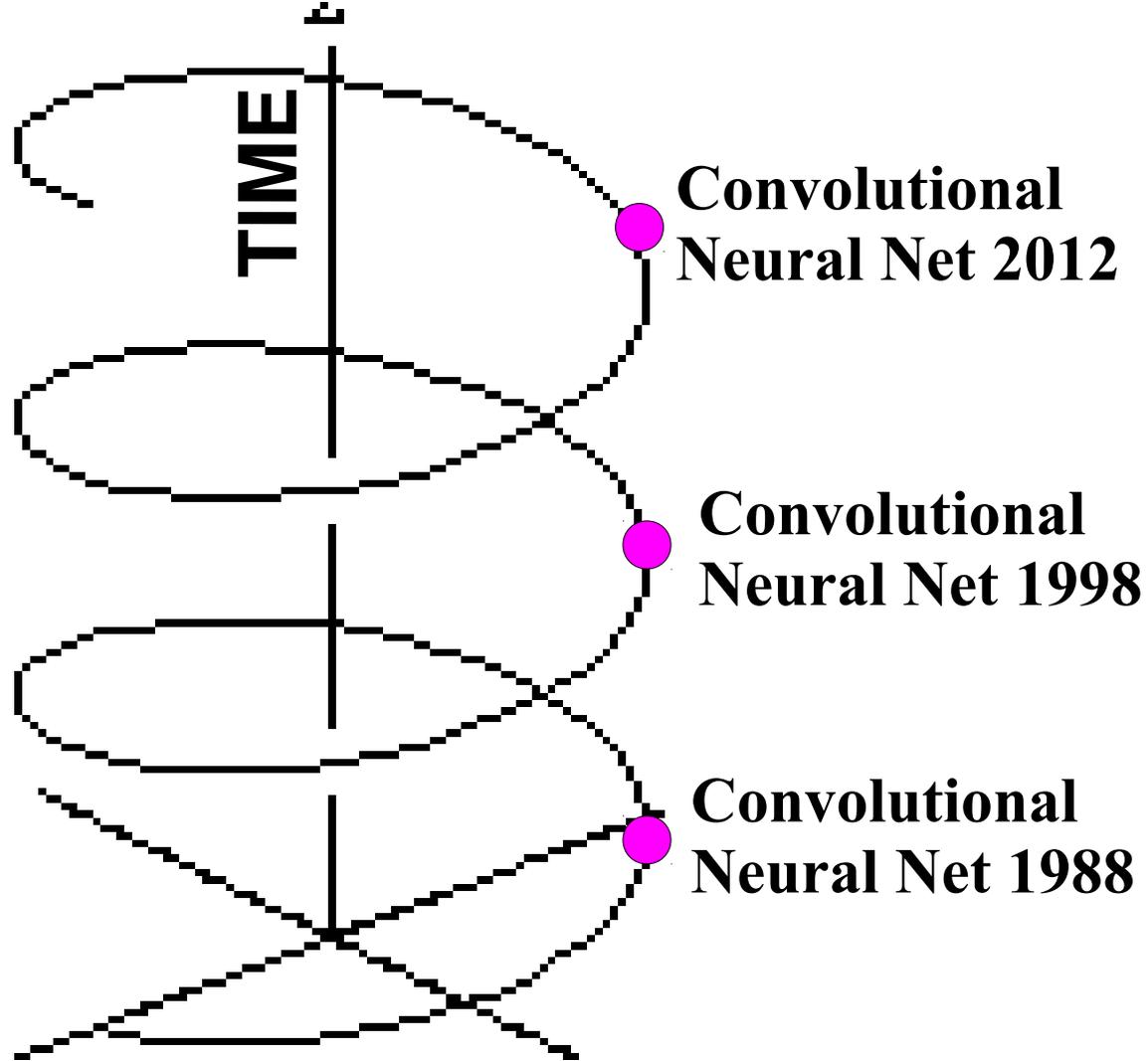
GMM

Boosting









**Q.:** Did we make any progress since then?

**A.:** The main reason for the breakthrough is: **data** and **GPU**, but we have also made networks **deeper** and more **non-linear**.

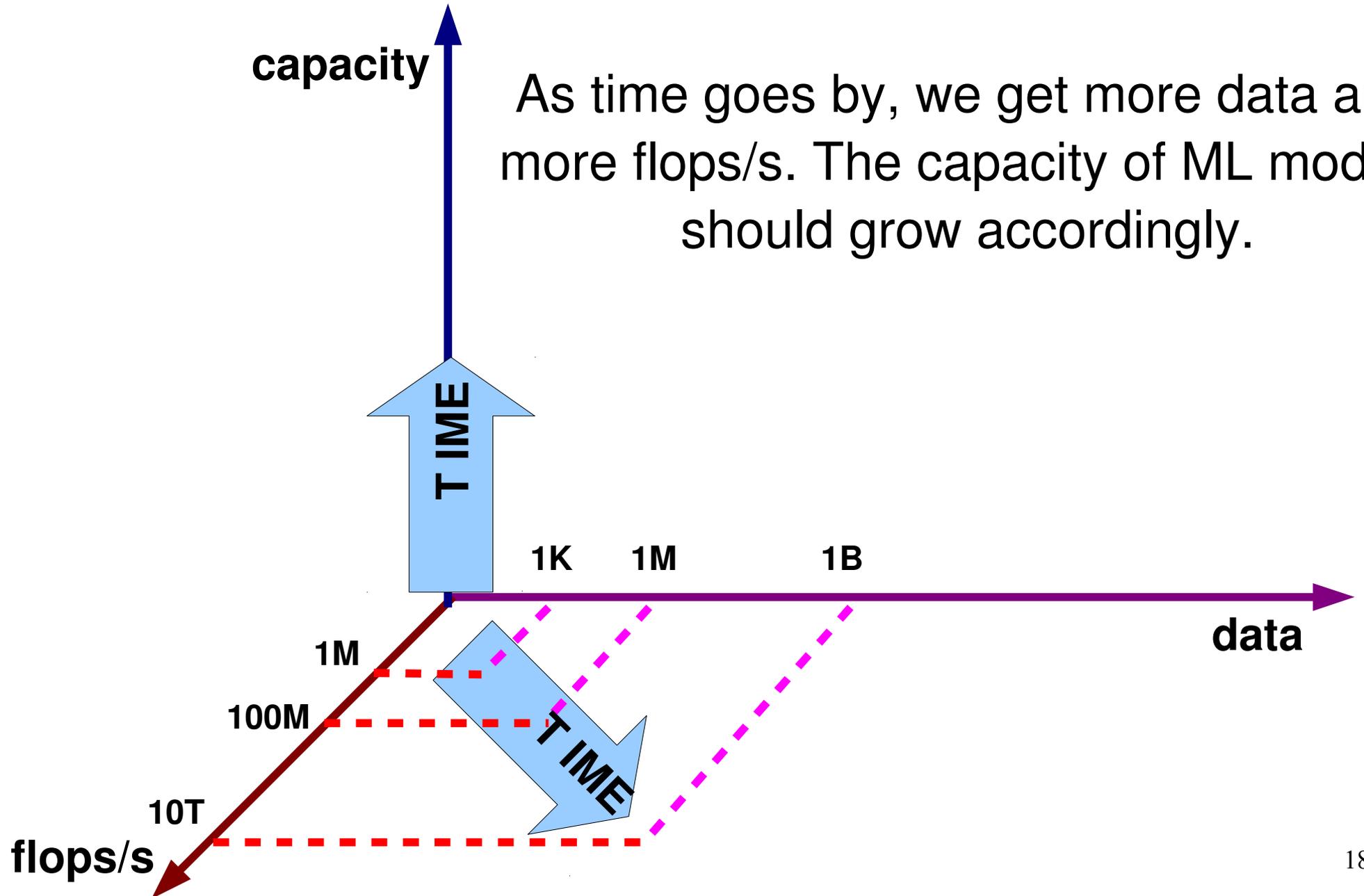
# ConvNets: History

- **1980 Fukushima**: designed network with same basic structure but did not train by backpropagation.
- **late 80s LeCun**: figured out backpropagation for CNN, popularized and deployed CNN for OCR applications and others.
- **1999 Poggio**: same basic structure but learning is restricted to top layer (k-means at second stage)
- **2006 LeCun**: unsupervised feature learning
- **2008 DiCarlo**: large scale experiments, normalization layer
- **2009 LeCun**: harsher non-linearities, normalization layer, learning unsupervised and supervised.
- **2011 Mallat**: provides a theory behind the architecture
- **2012 Hinton**: use bigger & deeper nets, GPUs, more data

# ConvNets: Why so successful now?

capacity

As time goes by, we get more data and more flops/s. The capacity of ML models should grow accordingly.



# ConvNets: Why so successful now?

