

Unsupervised Learning of Hierarchical Models

Marc'Aurelio Ranzato - Geoff Hinton

in collaboration with Josh Susskind and Vlad Mnih

Example: facial expression recognition

Unlabeled face images



Few labeled face images



happiness



neutral



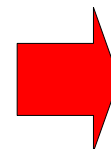
anger



fear



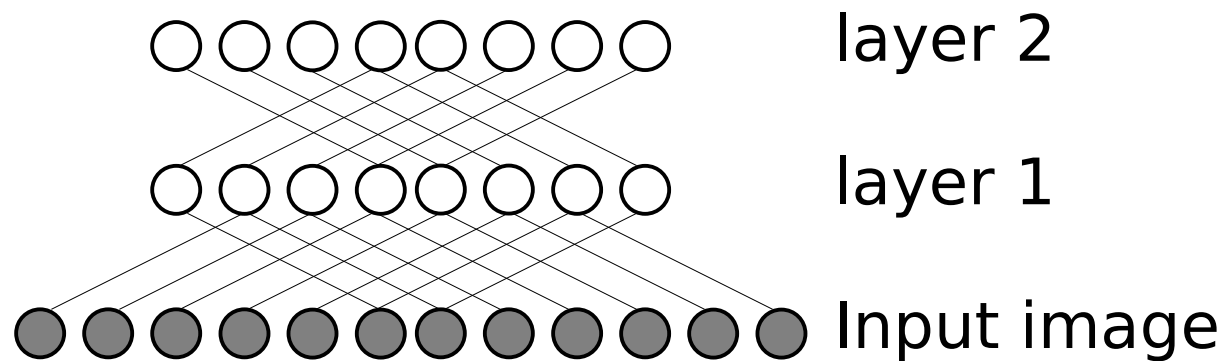
sadness



What is the expression?
(unseen occluded images)

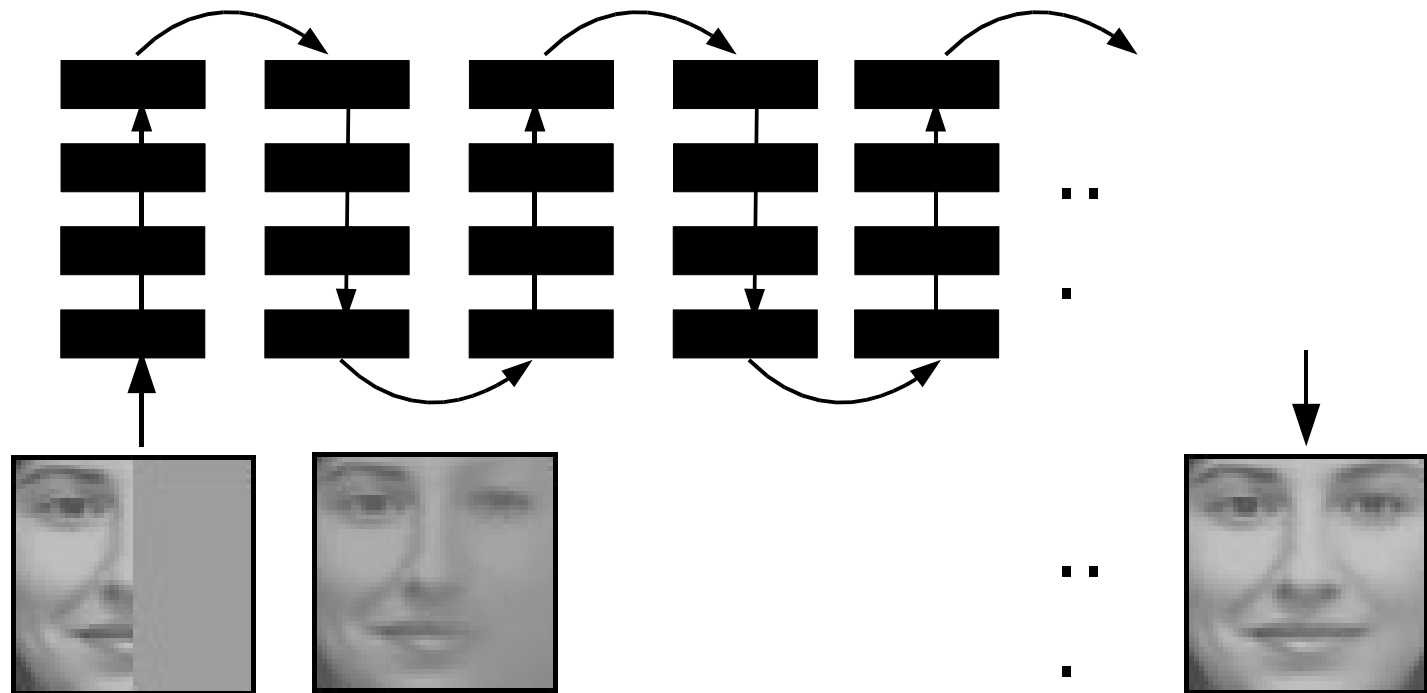
Example: facial expression recognition

1) training on unlabeled data: maximum likelihood



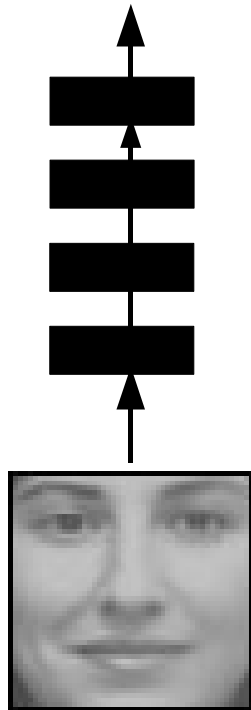
Example: facial expression recognition

- 1) training on unlabeled data: maximum likelihood
- 2) use the generative model to impute missing values

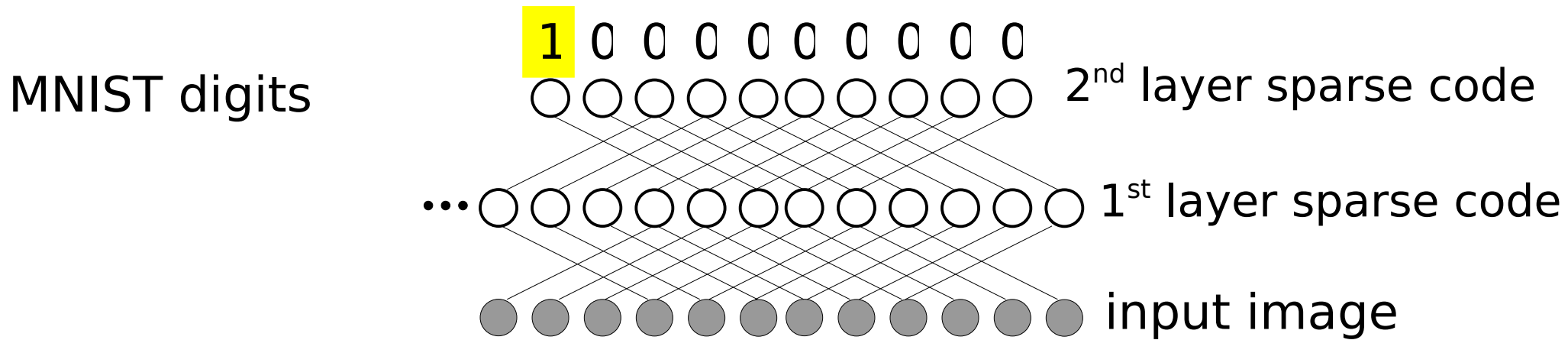


Example: facial expression recognition

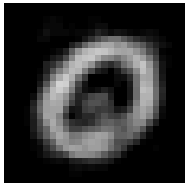
- 1) training on unlabeled data: maximum likelihood
- 2) use the generative model to impute missing values
- 3) use the latent representation as features
- 4) train classifier on labeled features



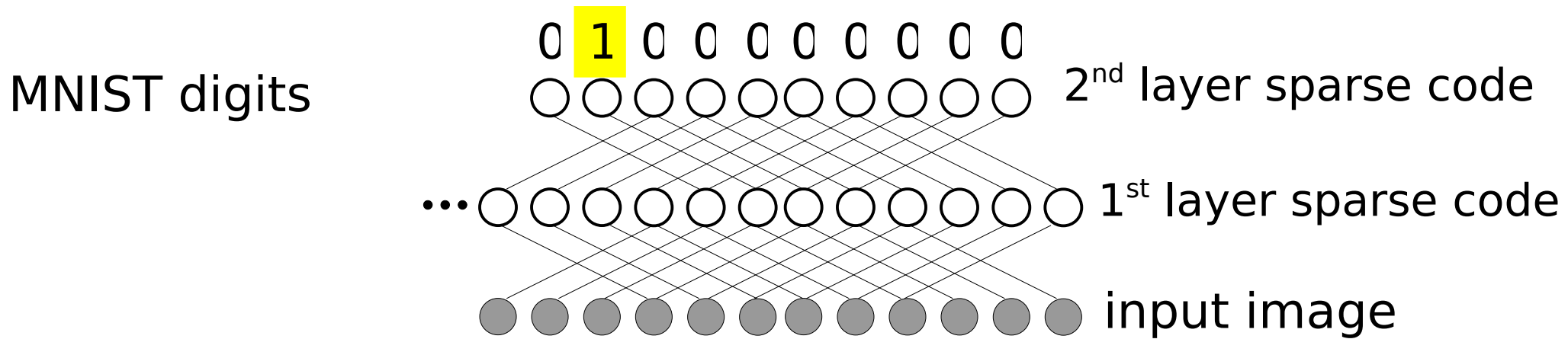
The power of hierarchical representations



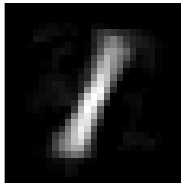
Receptive field of 2nd layer units



The power of hierarchical representations

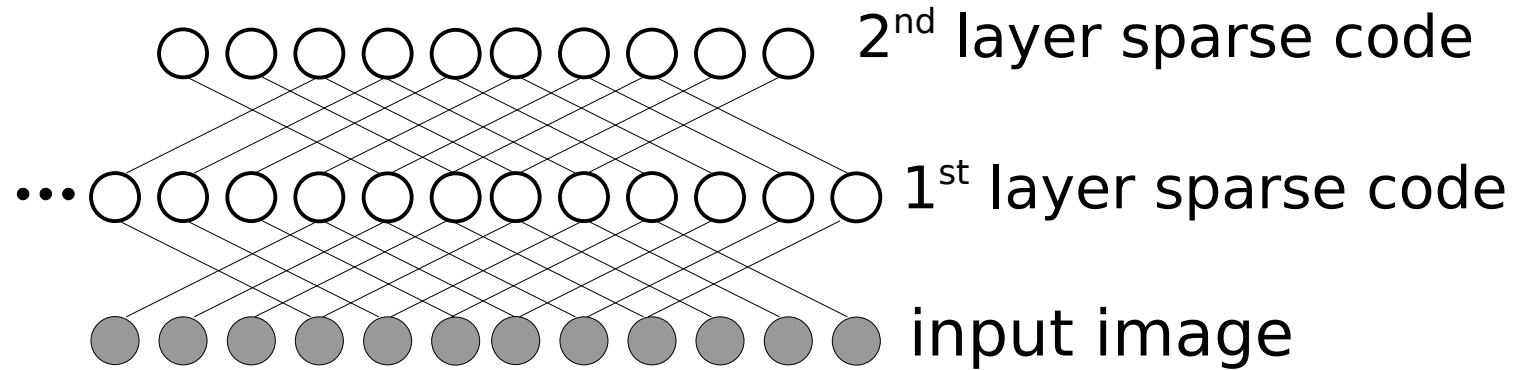


Receptive field of 2nd layer units

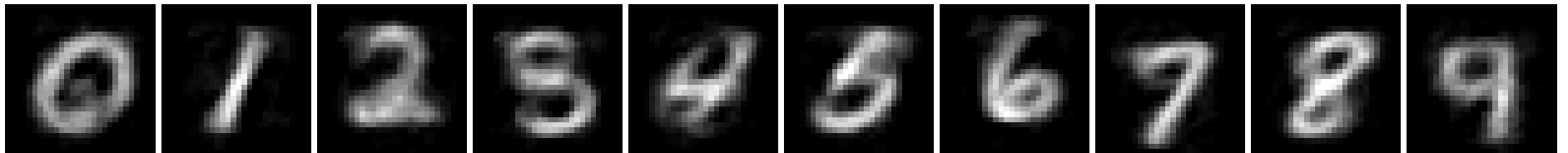


The power of hierarchical representations

MNIST digits



Receptive fields of 2nd layer units



Hierarchical representations can discover complex dependencies!

Unsupervised Learning

- we will learn hierarchical representations in a greedy fashion using unlabeled data (little if any assumption on the input)
- we can focus on single layer unsupervised algorithms
 - how can we interpret unsupervised algorithms?
 - how can they learn representations?
 - how can we compare them?

Two Approaches to Unsupervised Learning

1st strategy: constrain latent representation & optimize score only at training samples

- K-Means
- sparse coding

Ranzato et al. NIPS 06, Ranzato et al. CVPR 07, Ranzato et al. NIPS 07, Kavukcuoglu et al. CVPR 09

- use lower dimensional representations

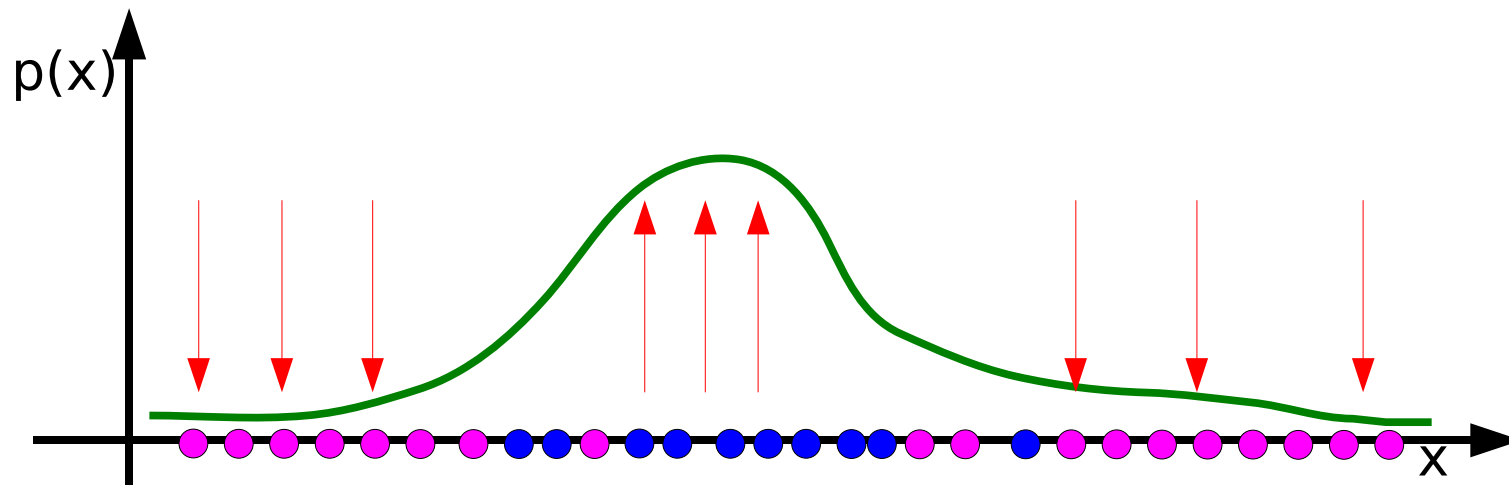
Ranzato et al. ICML 08

Two Approaches to Unsupervised Learning

- 1st strategy:** constrain latent representation & optimize score only at training samples
- K-Means
 - sparse coding
 - use lower dimensional representations

- 2nd strategy:** optimize score for training samples while normalizing the score over the whole space (maximum likelihood)

Ranzato et al. AISTATS 10, Ranzato et al. CVPR 10, Ranzato et al. NIPS 10, Ranzato et al. CVPR 11



Two Approaches to Unsupervised Learning

	Pros	Cons
1st strategy (constrain code)	<ul style="list-style-type: none">- efficient learning- defined through objective function (monitor training)	<ul style="list-style-type: none">- choice of functions- not robust to missing inputs- not good for generation
2nd strategy (maximum likelihood)	<ul style="list-style-type: none">- intuitive design- it can generate- compositionality	<ul style="list-style-type: none">- hard to train (variational inference, MCMC, etc.)

Outline

- mathematical formulation of the model
- training
- learning acoustic features for speech recognition
- generation of natural images
- recognition of facial expression under occlusion
- conclusion

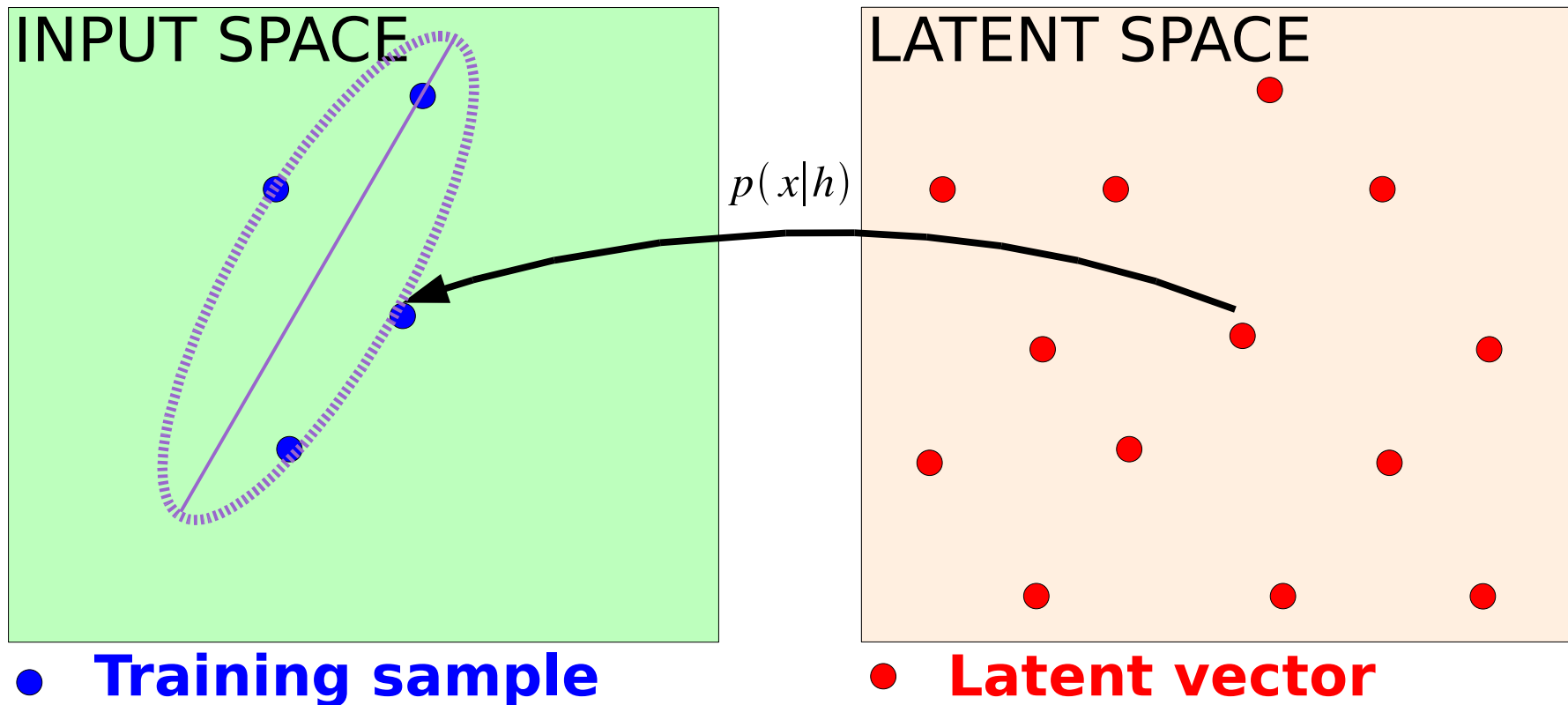
Outline

- mathematical formulation of the model
- training
- learning acoustic features for speech recognition
- generation of natural images
- recognition of facial expression under occlusion
- conclusion

Means and Covariances: PPCA

$$p(x|h) \neq p(x)$$

$p(x)$ specifies covariance across all data points

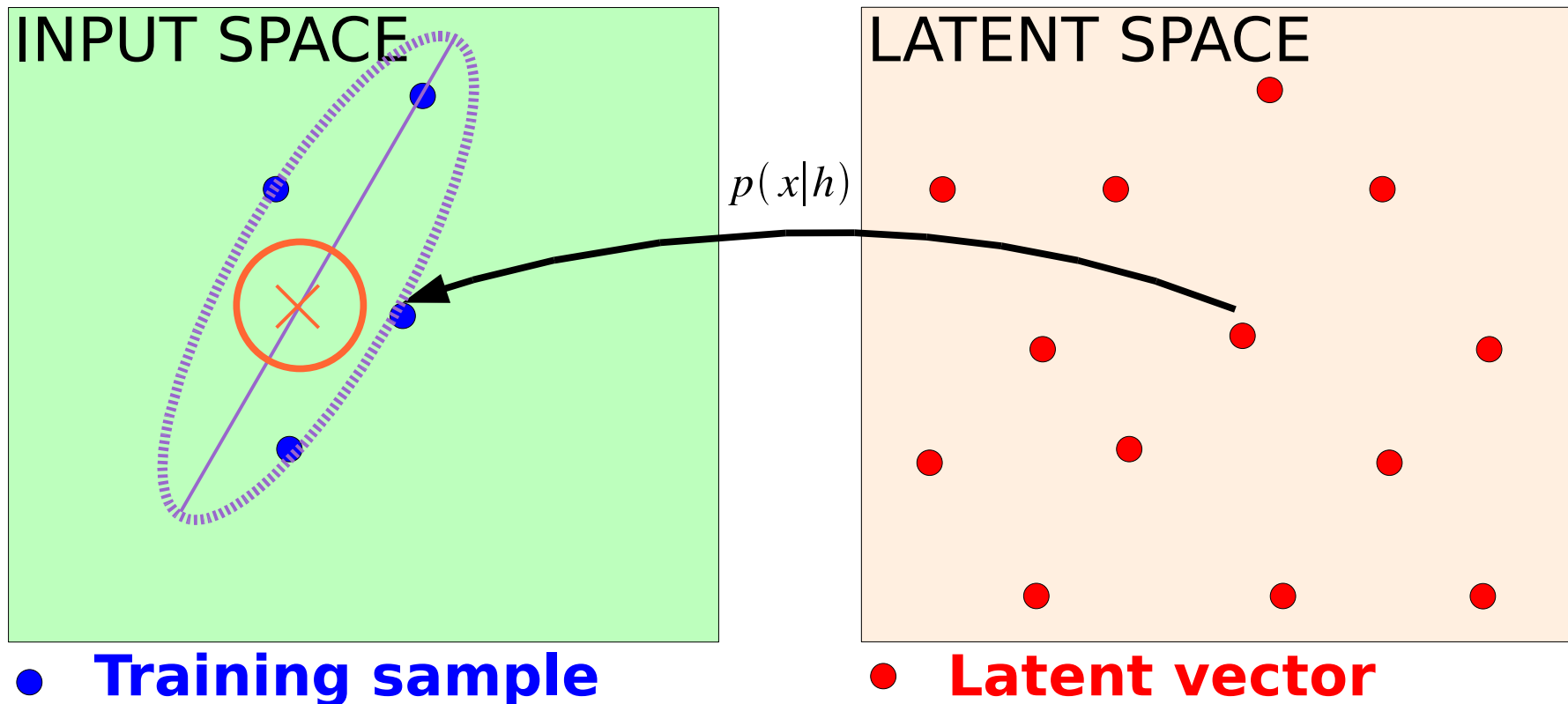


Means and Covariances: PPCA

$$p(x|h) \neq p(x)$$

$$h_t \sim p(h|x_t)$$

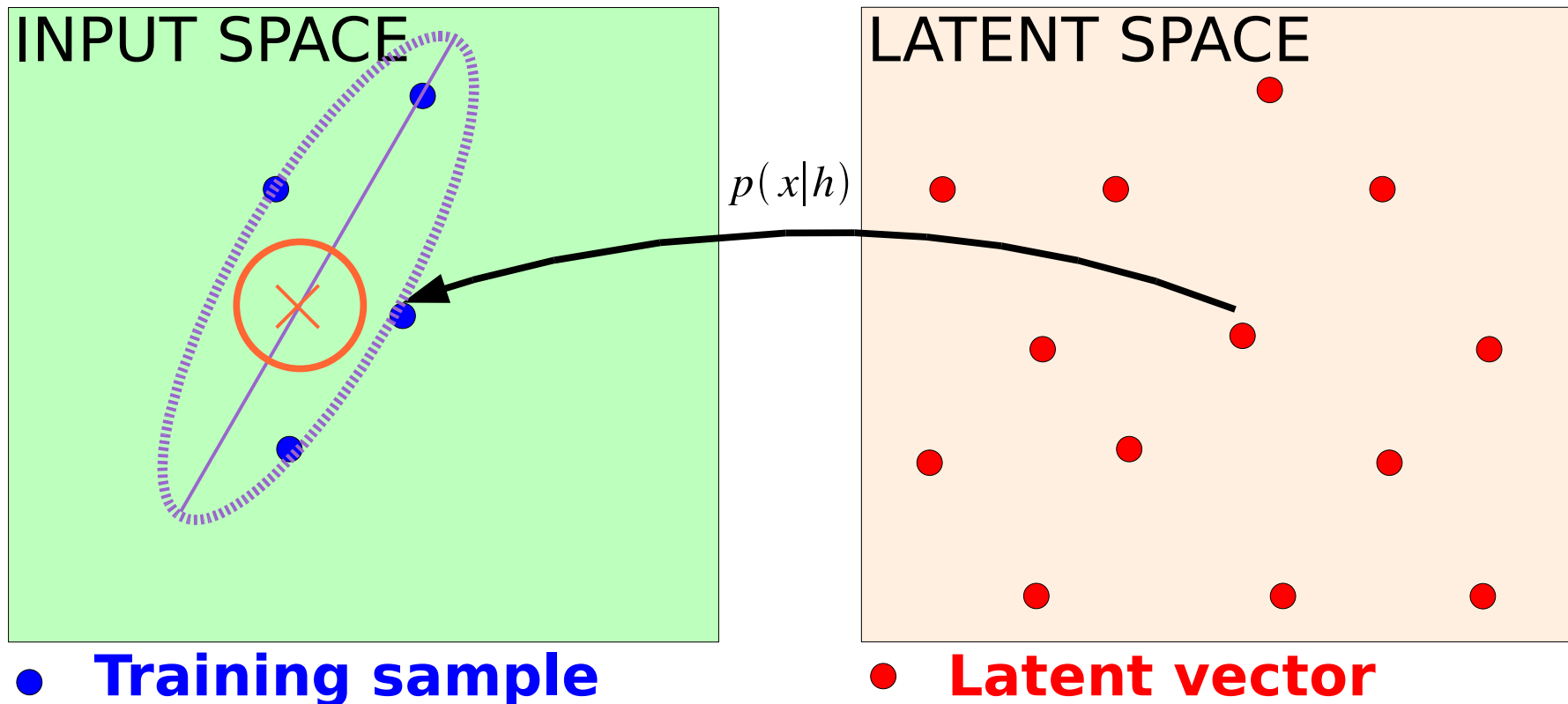
$p(x|h_t)$ specifies distribution for each "data point"



Means and Covariances: PPCA

$$p(x) = \int_h p(x|h) p(h)$$

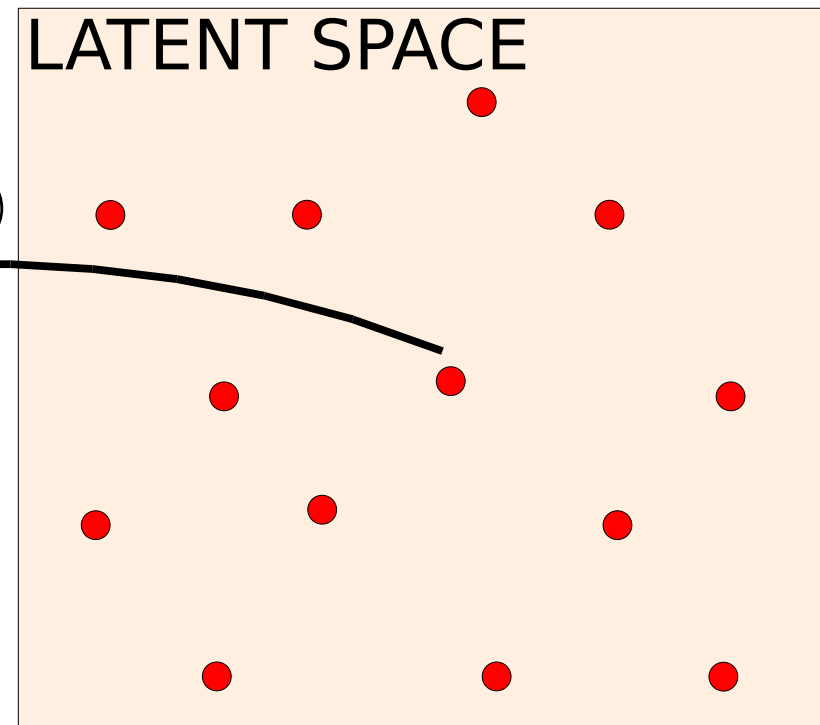
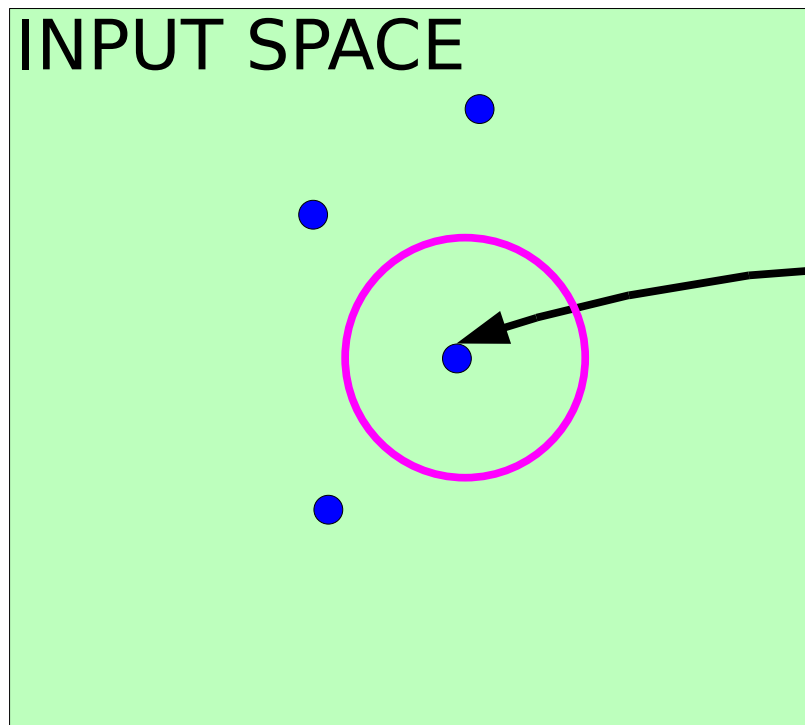
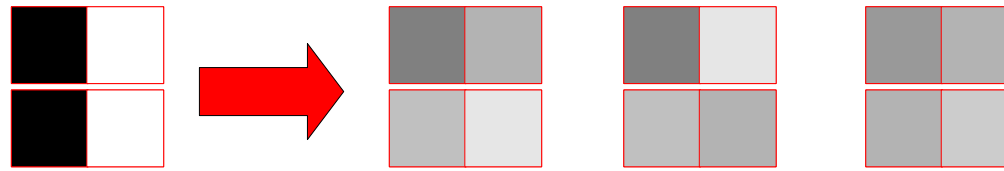
$p(x|h)$ specifies distribution for each “data point”



Conditional Distribution Over Input

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

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



$p(x|h)$

● Training sample

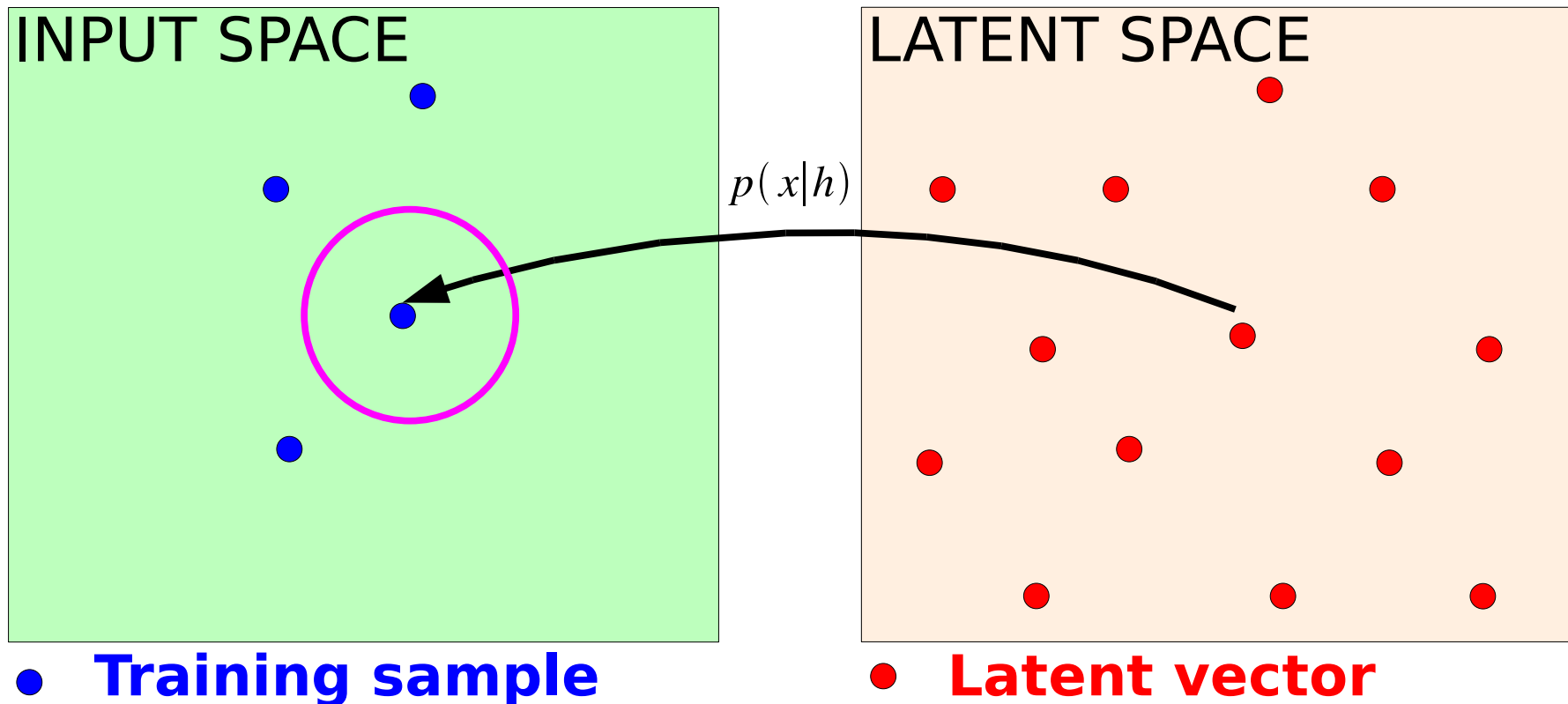
● Latent vector

Conditional Distribution Over Input

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

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

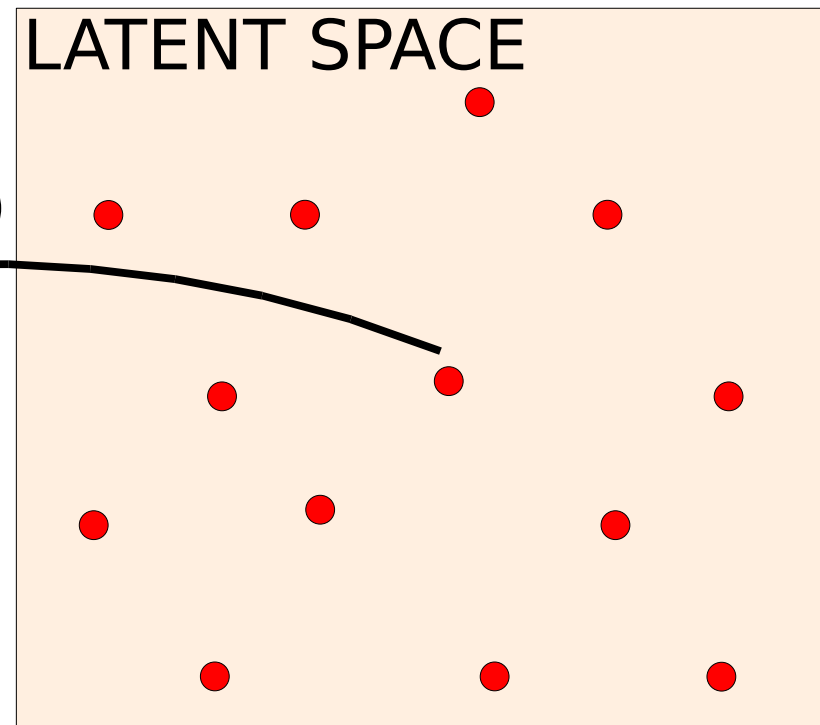
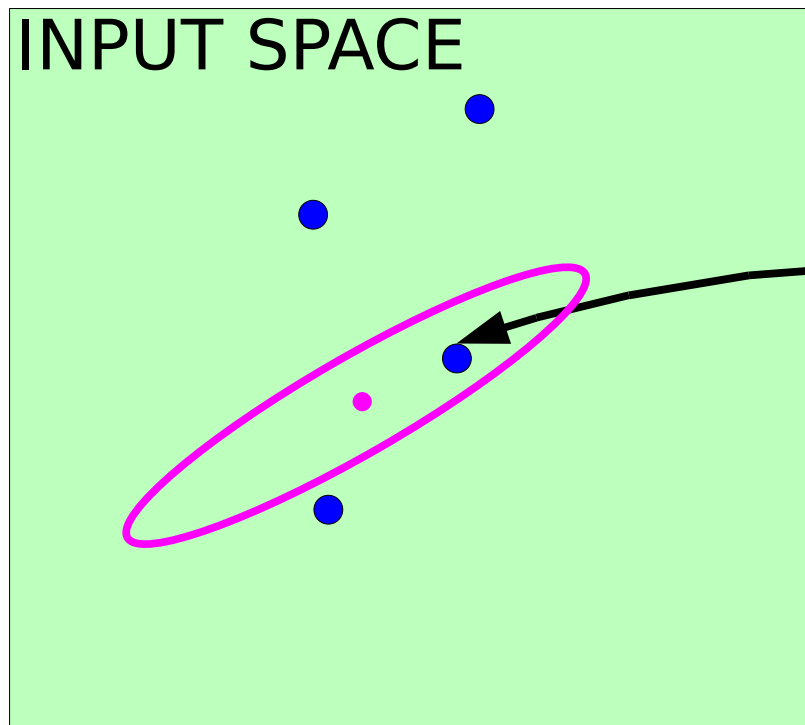
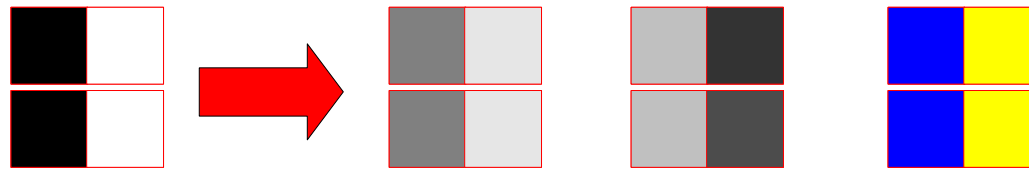
Poor model of dependency between input variables



Conditional Distribution Over Input

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

- examples: PoT, cRBM



● **Training sample**

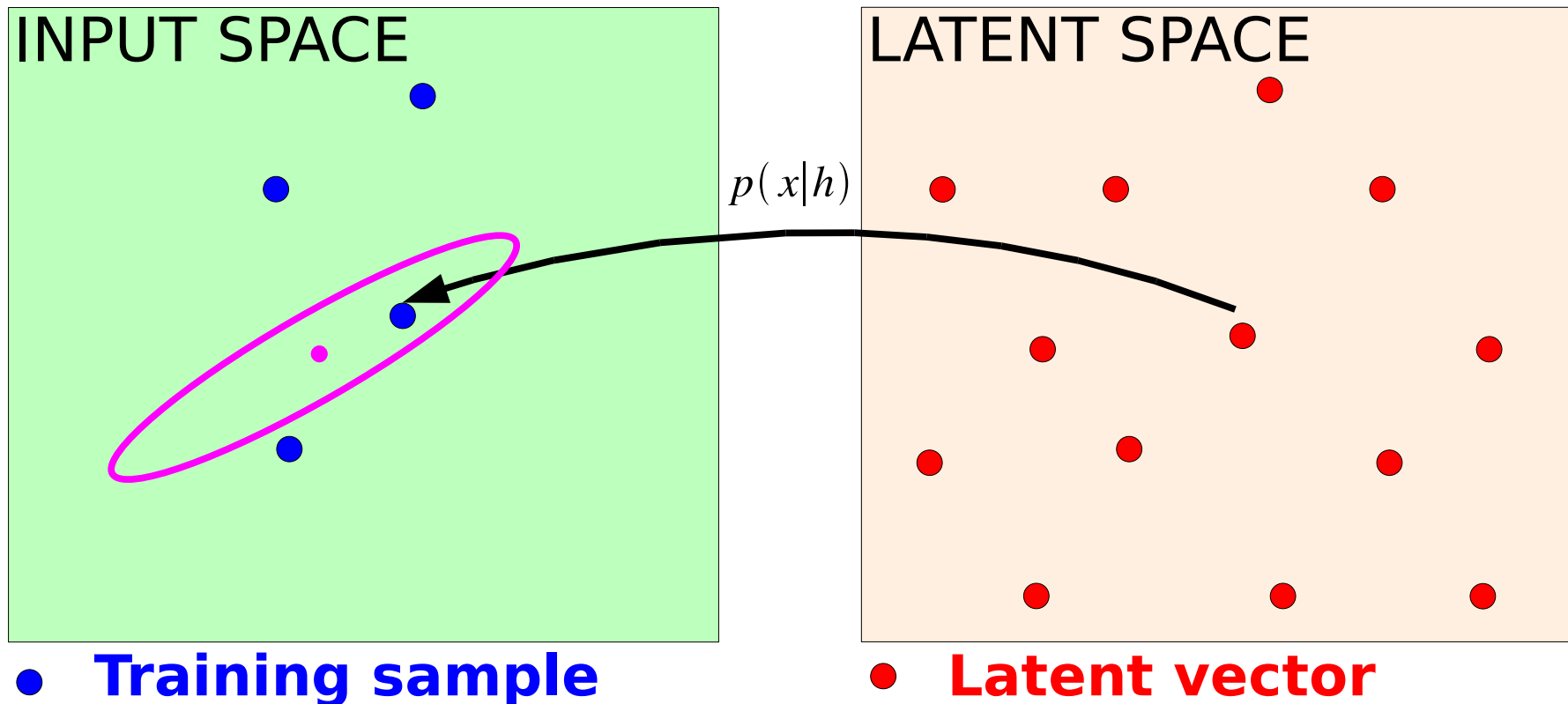
● **Latent vector**

Conditional Distribution Over Input

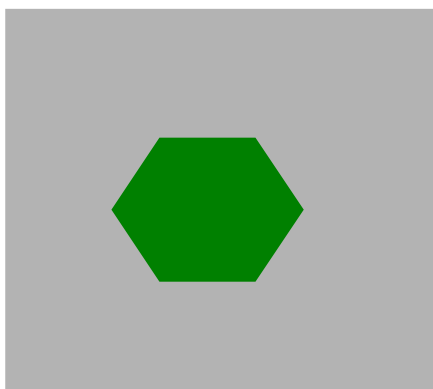
$$p(x|h) = N(0, \text{Covariance}(h))$$

- examples: PoT, cRBM

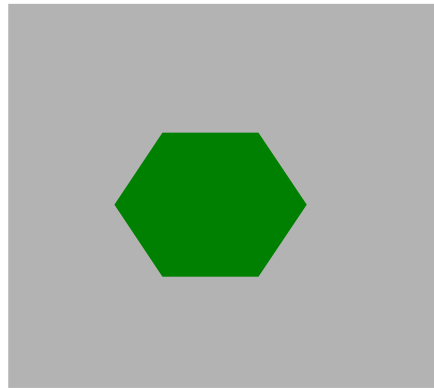
Poor model of mean intensity



input image

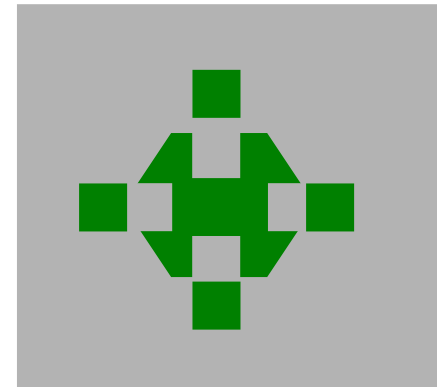
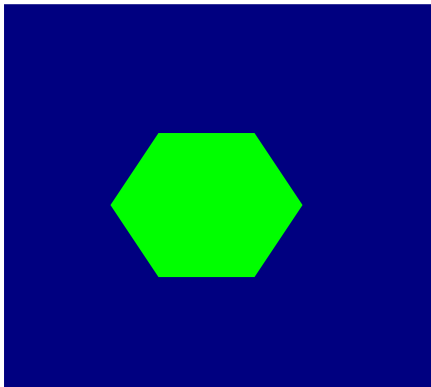


input image

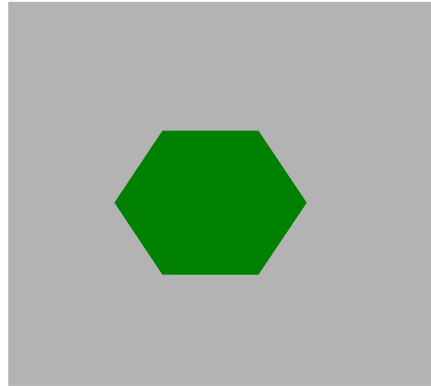


Which image has similar mean (but different covariance)?

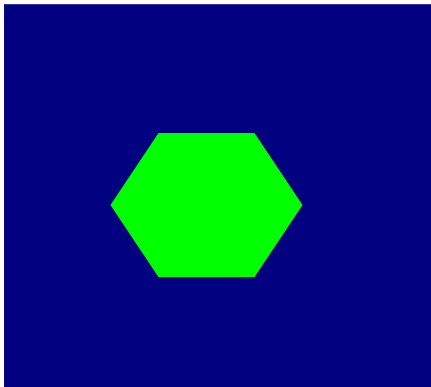
Which image has same covariance (but different mean)?



input image

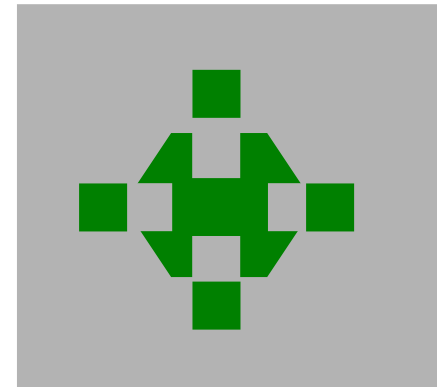


same covariance
but different mean



Equivalent image according to PoT

similar mean
but different covariance

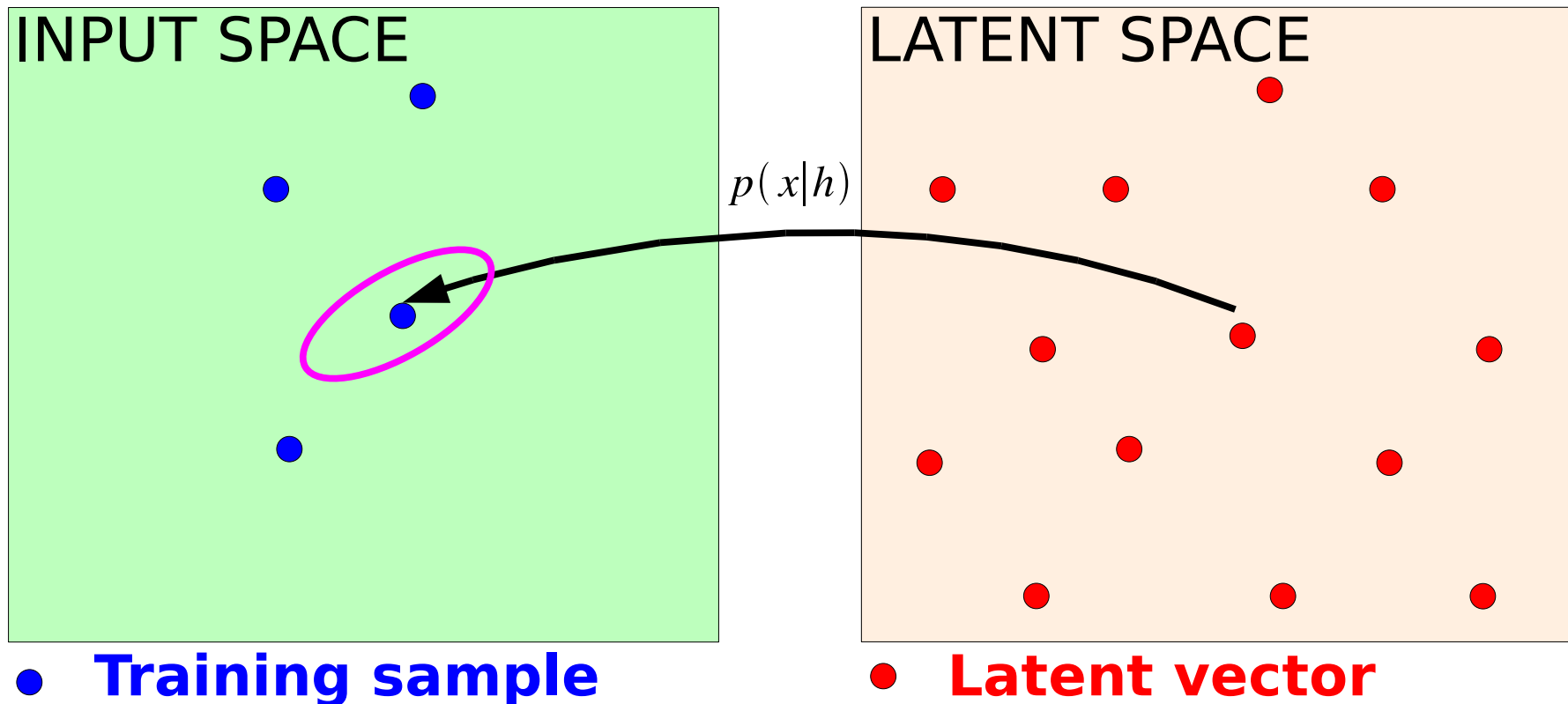
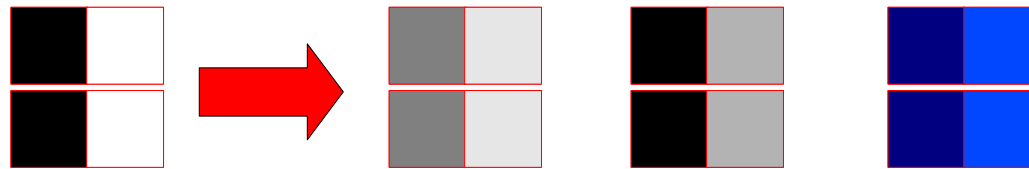


Equivalent image according to FA

Conditional Distribution Over Input

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

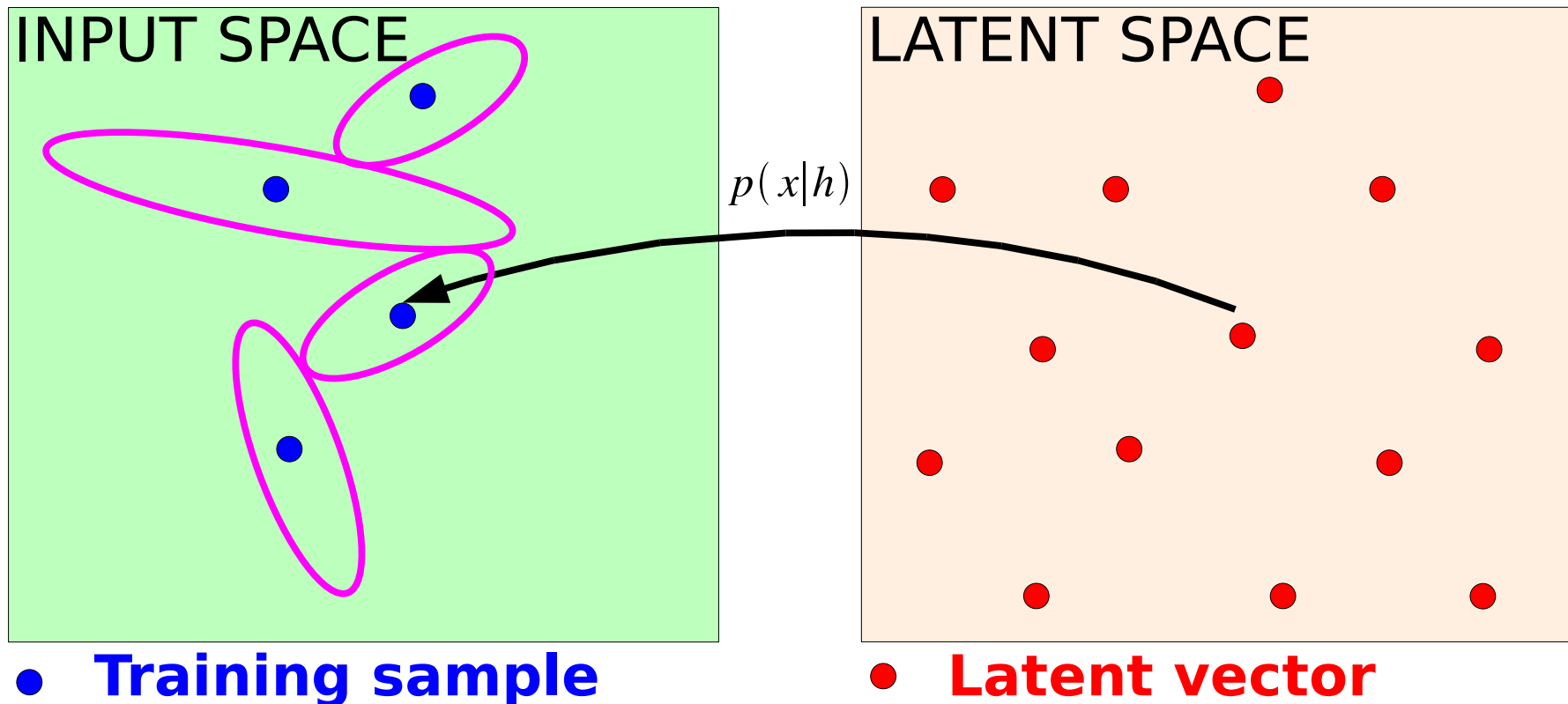
- this is what we propose: mcRBM, mPoT

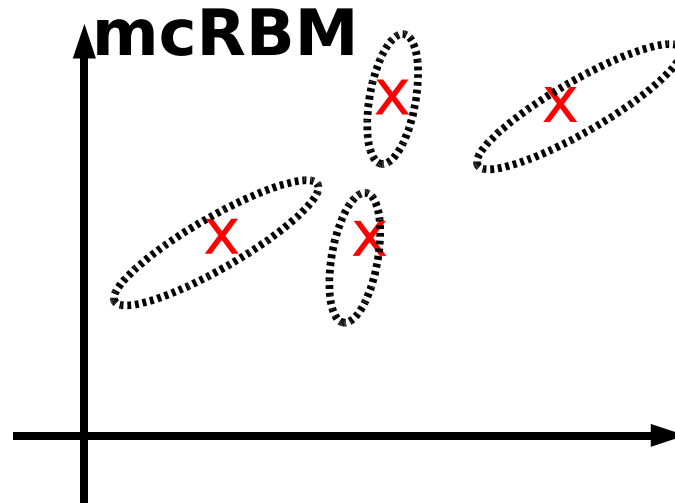


Conditional Distribution Over Input

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

- this is what we propose: mcRBM, mPoT





- two sets of latent variables to modulate mean and covariance of the conditional distribution over the input
- energy-based model

$$p(x, h^m, h^c) \propto \exp(-E(x, h^m, h^c))$$

$$x \in \mathbb{R}^D$$

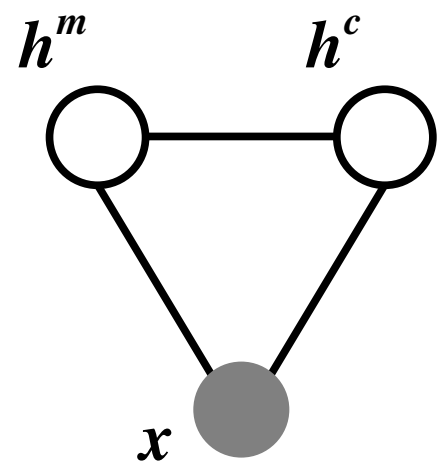
$$h^c \in \{0, 1\}^M$$

$$h^m \in \{0, 1\}^N$$

- easy generation
- slow inference

$$p(x|h^m, h^c) = N(m(h^m), \Sigma(h^c))$$

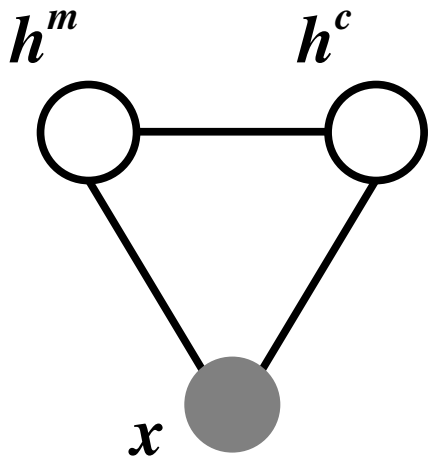
$$E = \frac{1}{2}(x - m)' \Sigma^{-1}(x - m)$$



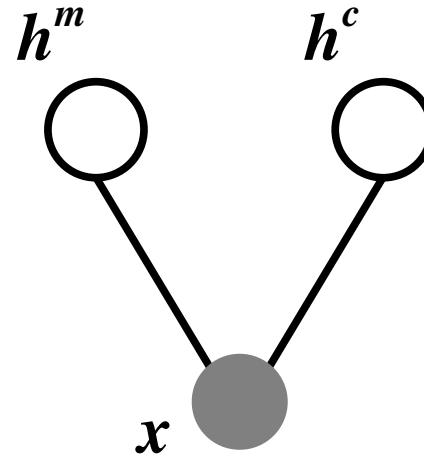
- easy generation
- slow inference

$$p(x|h^m, h^c) = N(m(h^m), \Sigma(h^c))$$

$$E = \frac{1}{2}(x - m)' \Sigma^{-1}(x - m)$$



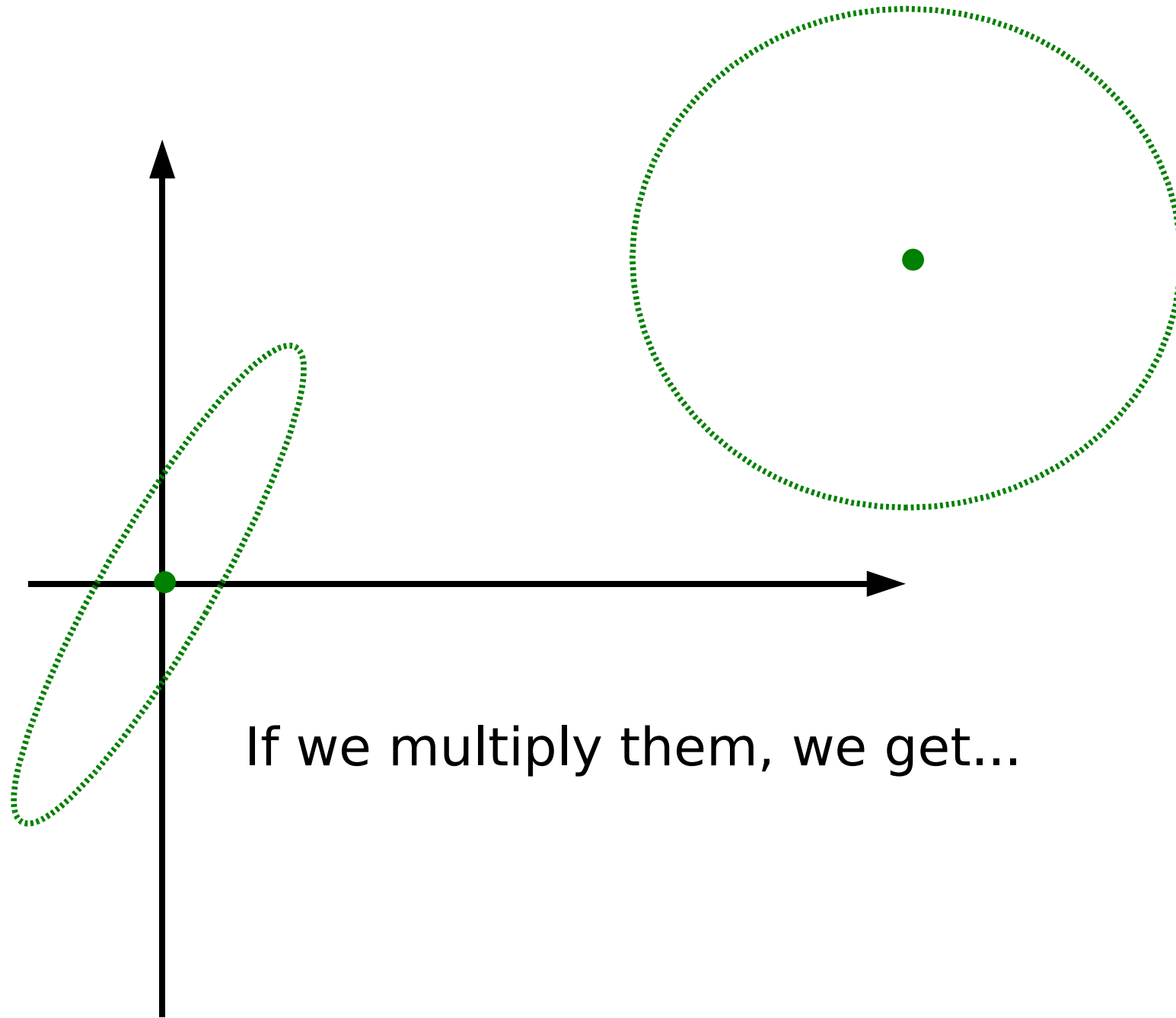
- less easy generation
- fast inference



$$E = \frac{1}{2} x' \Sigma^{-1} x - m x$$

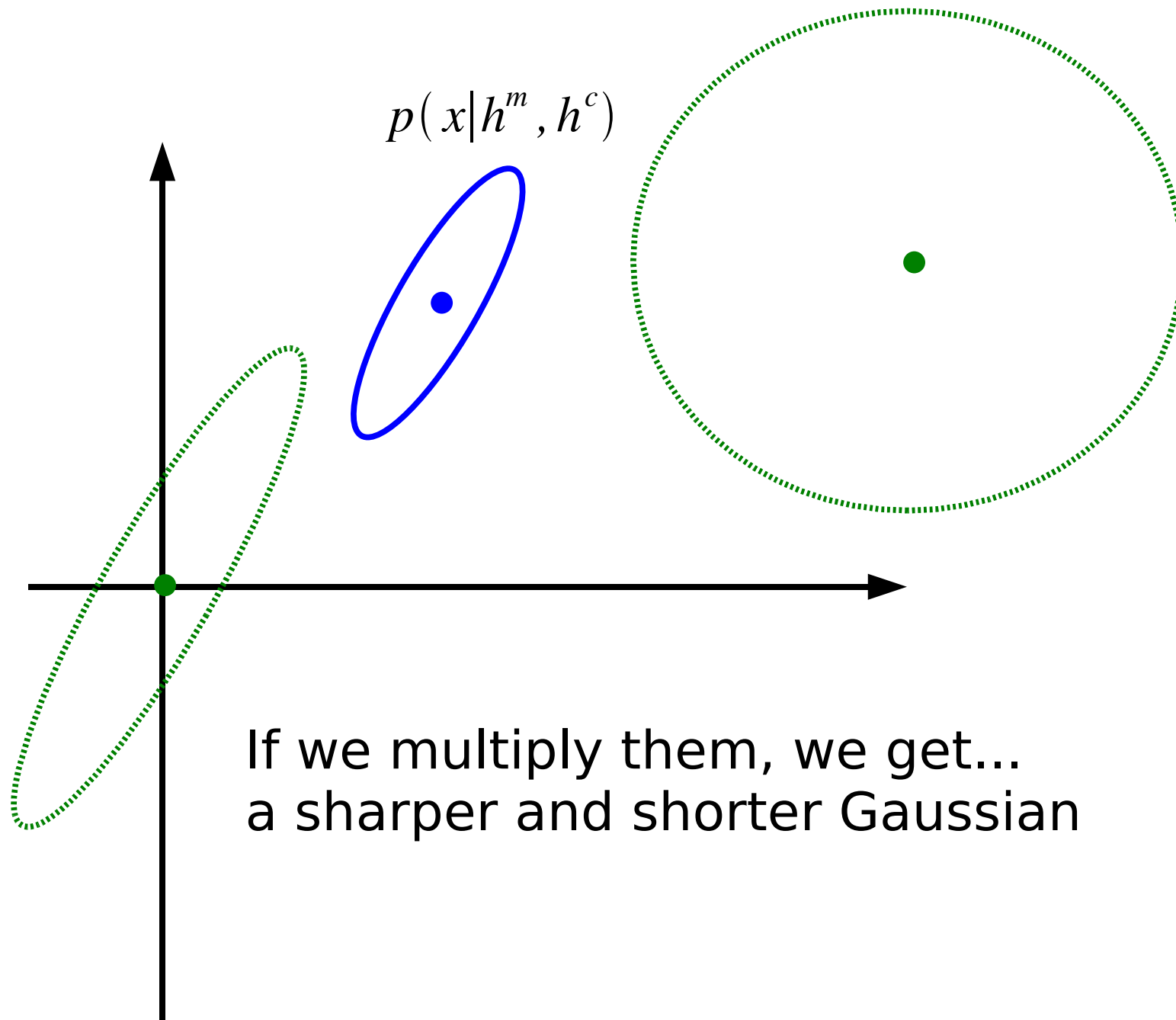
$$p(x|h^m, h^c) = N(\Sigma(h^c)^{-1} m(h^m), \Sigma(h^c))$$

Geometric interpretation



If we multiply them, we get...

Geometric interpretation

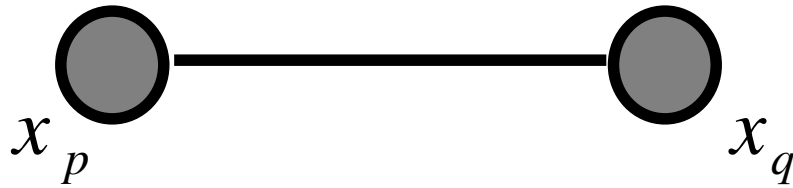


Covariance part of the energy function:

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

$$x \in \mathbb{R}^D$$

$$\Sigma^{-1} \in \mathbb{R}^{D \times D}$$



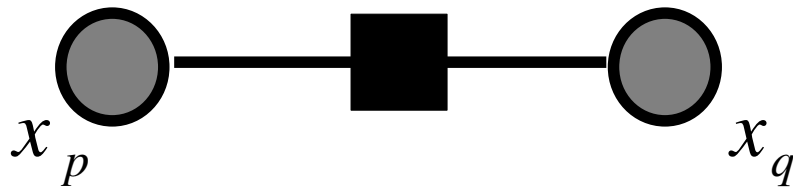
pair-wise MRF

Covariance part of the energy function:

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

$x \in \mathbb{R}^D$ factorization

$C \in \mathbb{R}^{D \times F}$



pair-wise MRF

Covariance part of the energy function:

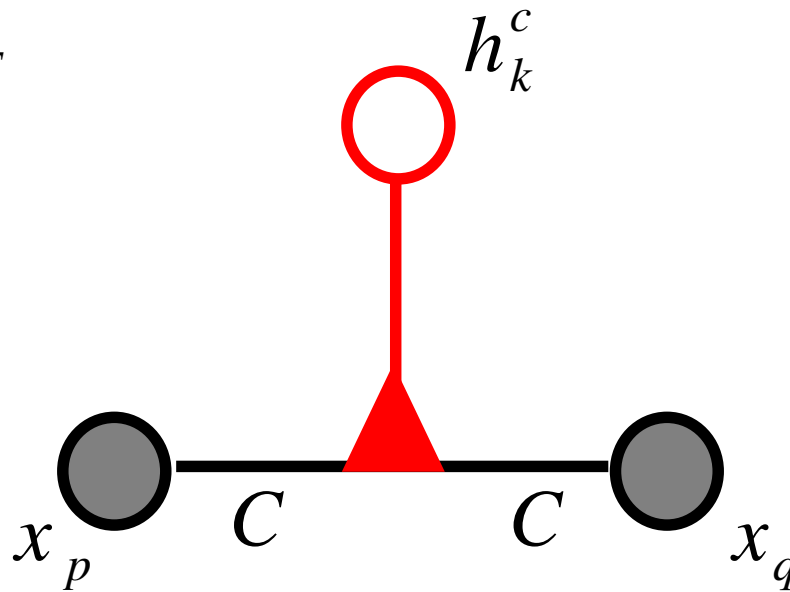
$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(h^c)] C' x$$

$$x \in \mathbb{R}^D$$

factorization + hidden

$$C \in \mathbb{R}^{D \times F}$$

$$h^c \in \{0, 1\}^F$$



gated MRF

Covariance part of the energy function:

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

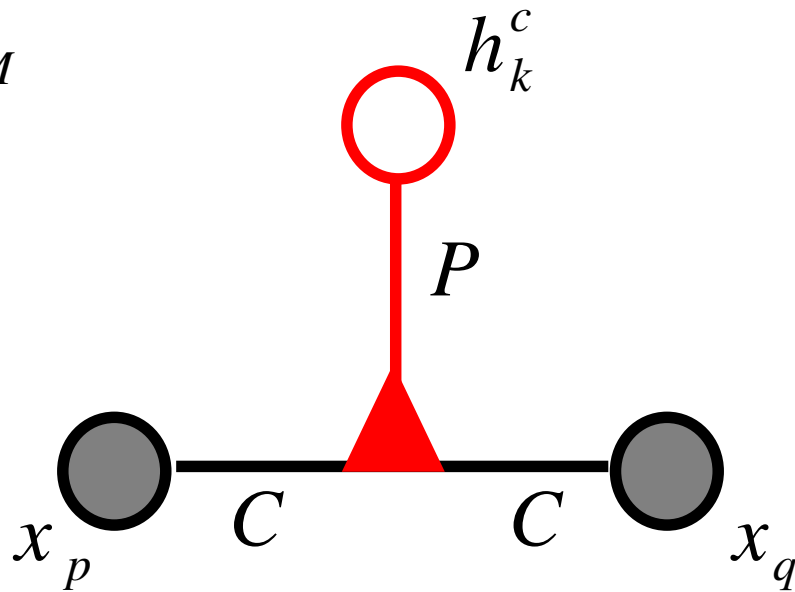
$$x \in \mathbb{R}^D$$

factorization + hidden

$$C \in \mathbb{R}^{D \times F}$$

$$h^c \in \{0, 1\}^M$$

$$P \in \mathbb{R}^{F \times M}$$



gated MRF

Overall energy function:

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

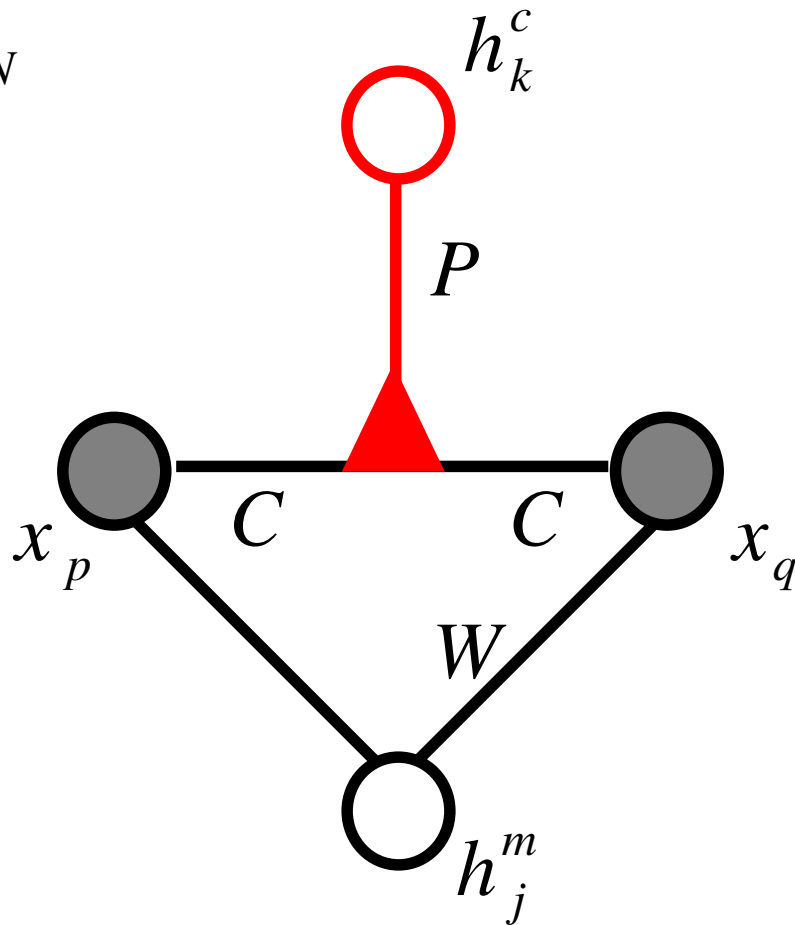
$$x \in \mathbb{R}^D$$

$$W \in \mathbb{R}^{D \times N}$$

$$h^m \in \{0, 1\}^N$$

covariance
part

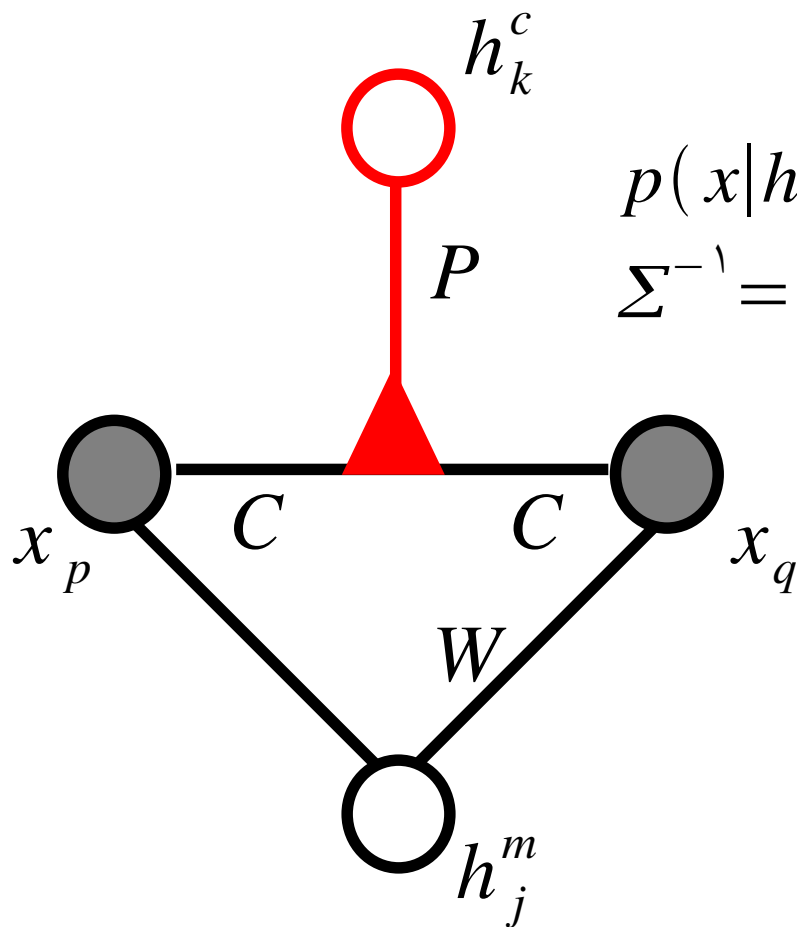
mean part



gated MRF

Overall energy function:

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



$$p(x|h^c, h^m) = N(\Sigma(W h^m), \Sigma)$$

$$\Sigma^{-1} = C \text{diag}[P h^c] C' + I$$

gated MRF

Overall energy function:

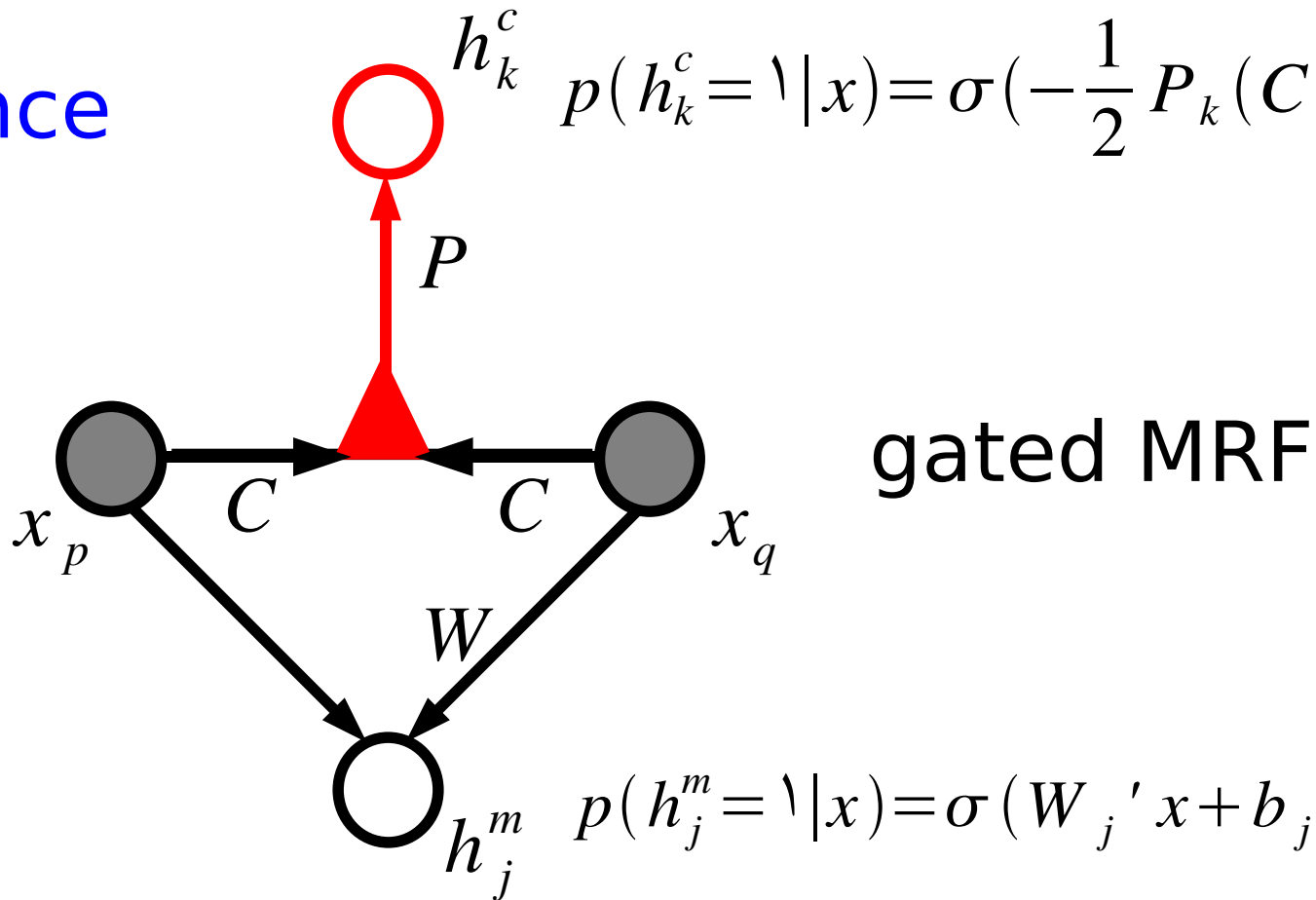
$$E(x, h^c, h^m) = \frac{1}{2} x' C [\text{diag}(P h^c)] C' x + \frac{1}{2} x' x - x' W h^m$$

covariance
part

mean part

inference

$$p(h_k^c = \lambda | x) = \sigma\left(-\frac{1}{2} P_k (C' x)^2 + b_k\right)$$



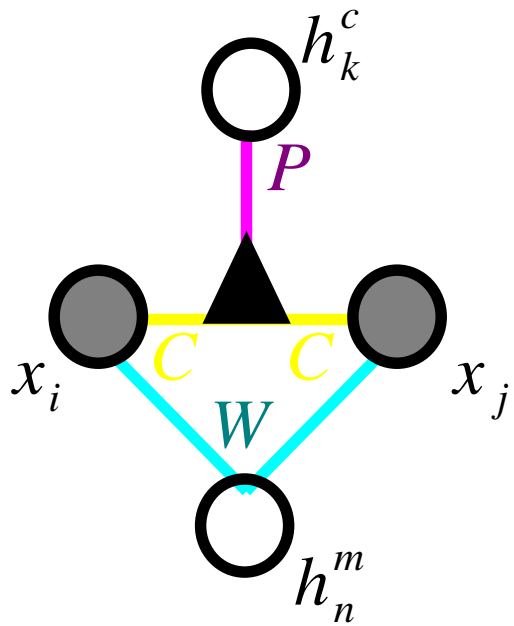
$$p(h_j^m = \lambda | x) = \sigma(W_j' x + b_j)$$

Outline

- mathematical formulation of the model
- **training**
- learning acoustic features for speech recognition
- generation of natural images
- recognition of facial expression under occlusion
- conclusion

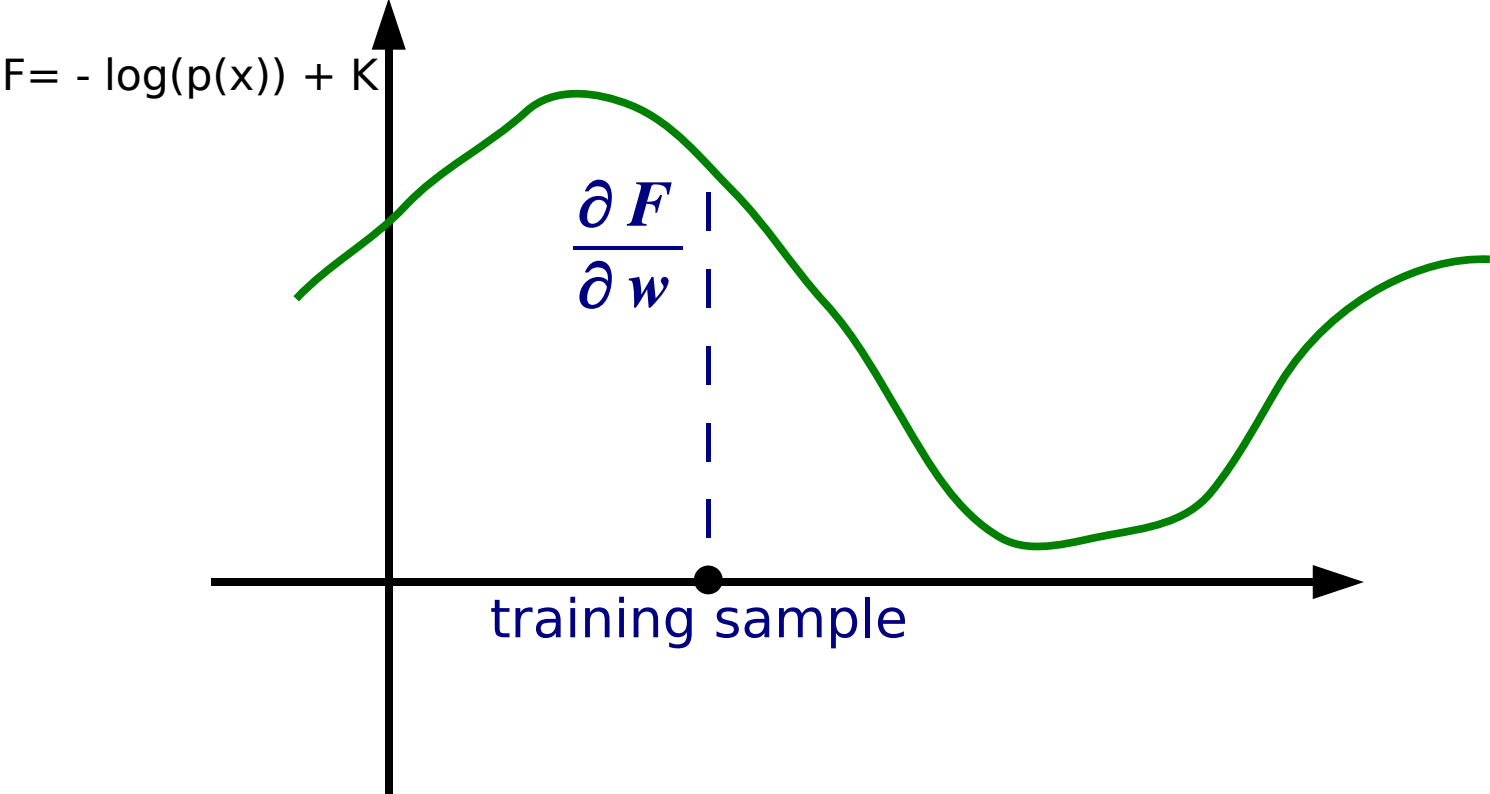
Learning

- maximum likelihood $p(x) = \frac{\int_{h^m, h^c} e^{-E(x, h^m, h^c)}}{\int_{x, h^m, h^c} e^{-E(x, h^m, h^c)}}$
- Fast Persistent Contrastive Divergence
- Hybrid Monte Carlo to draw samples

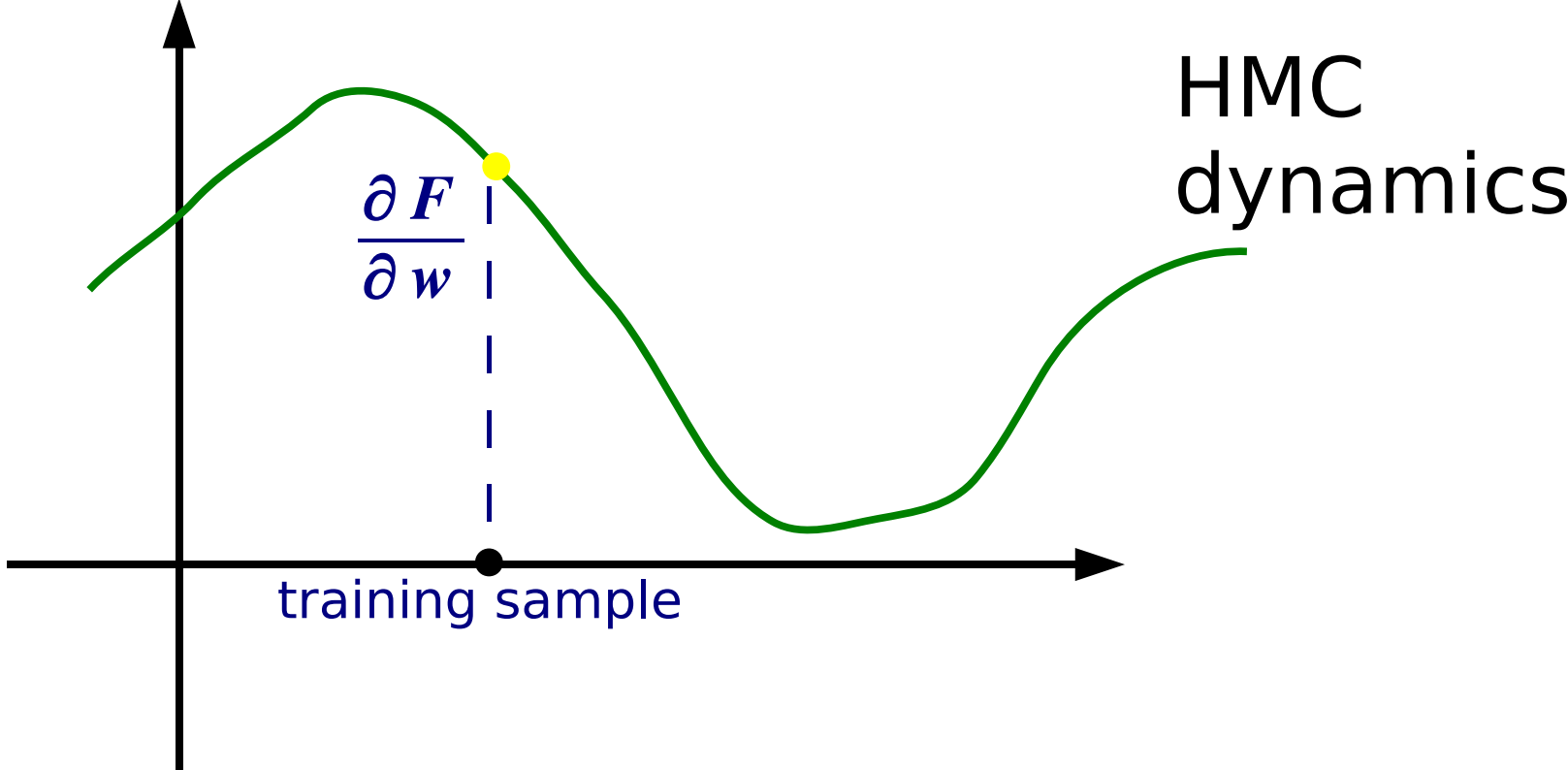


$$E = \frac{1}{\tau} x' C [\text{diag} (P h^c)] C' x - x' W h^m + \dots$$

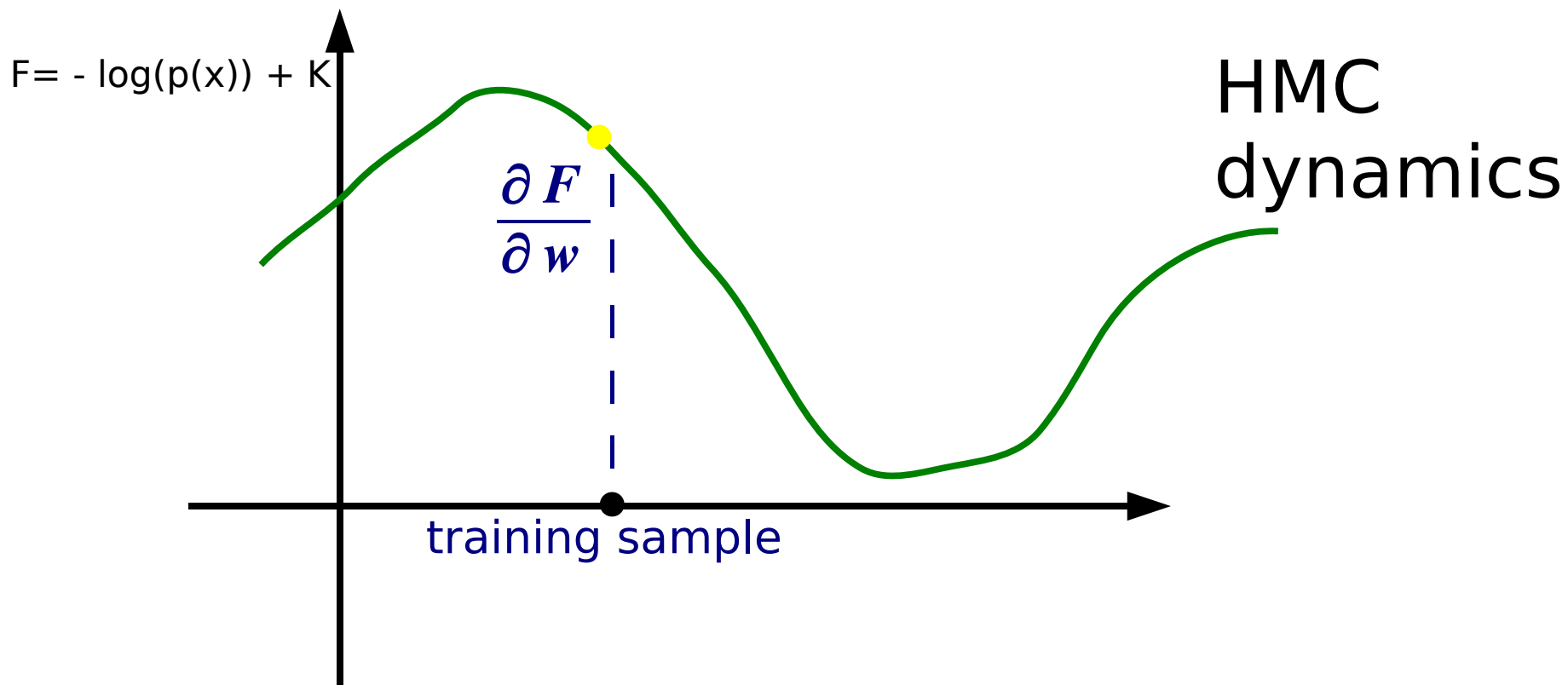
Learning



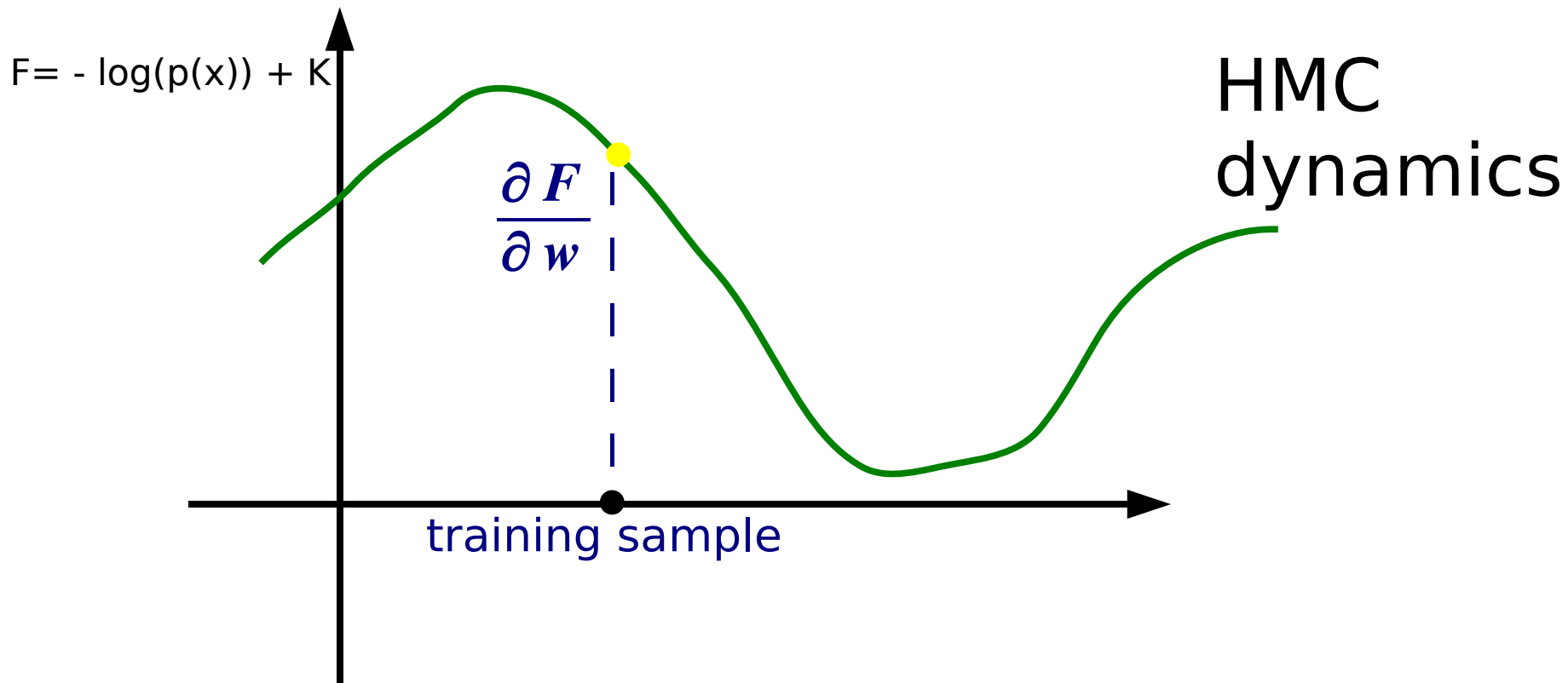
Learning



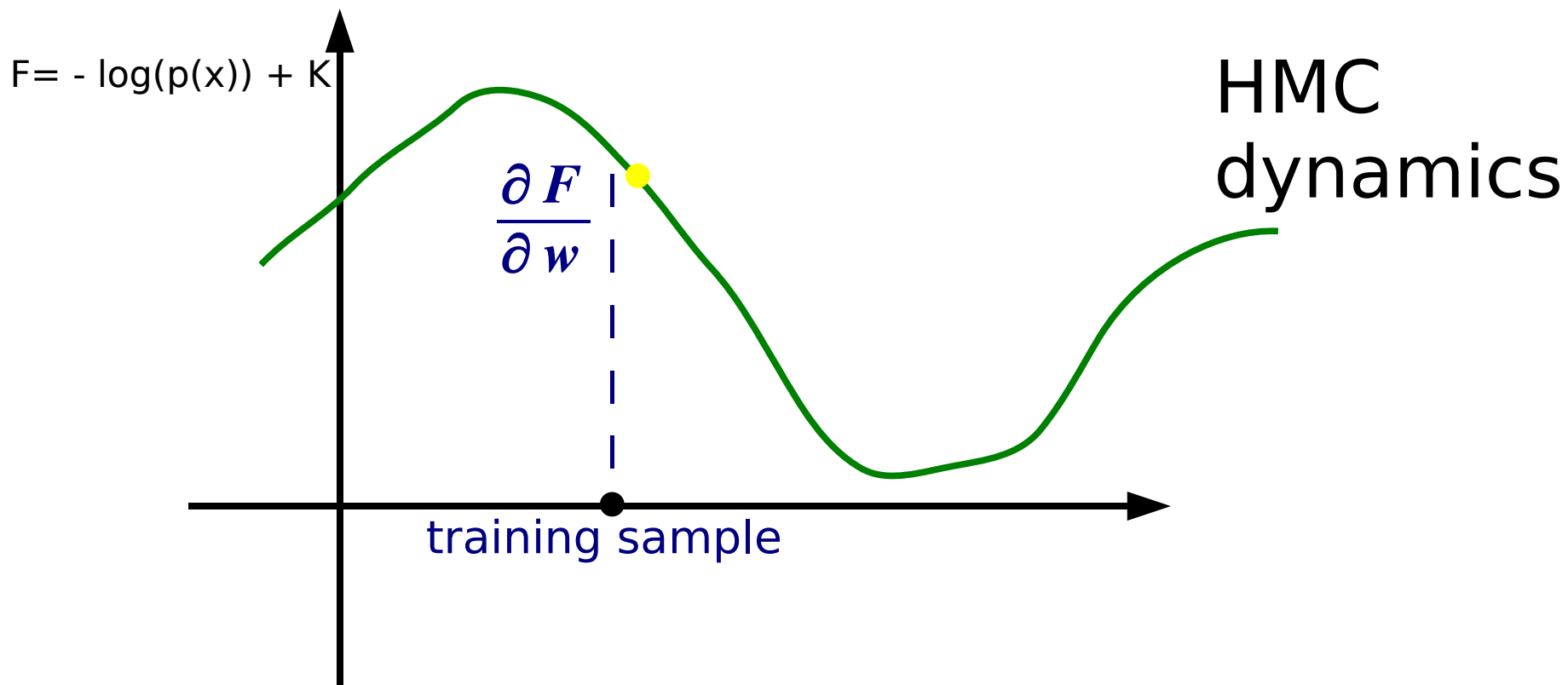
Learning



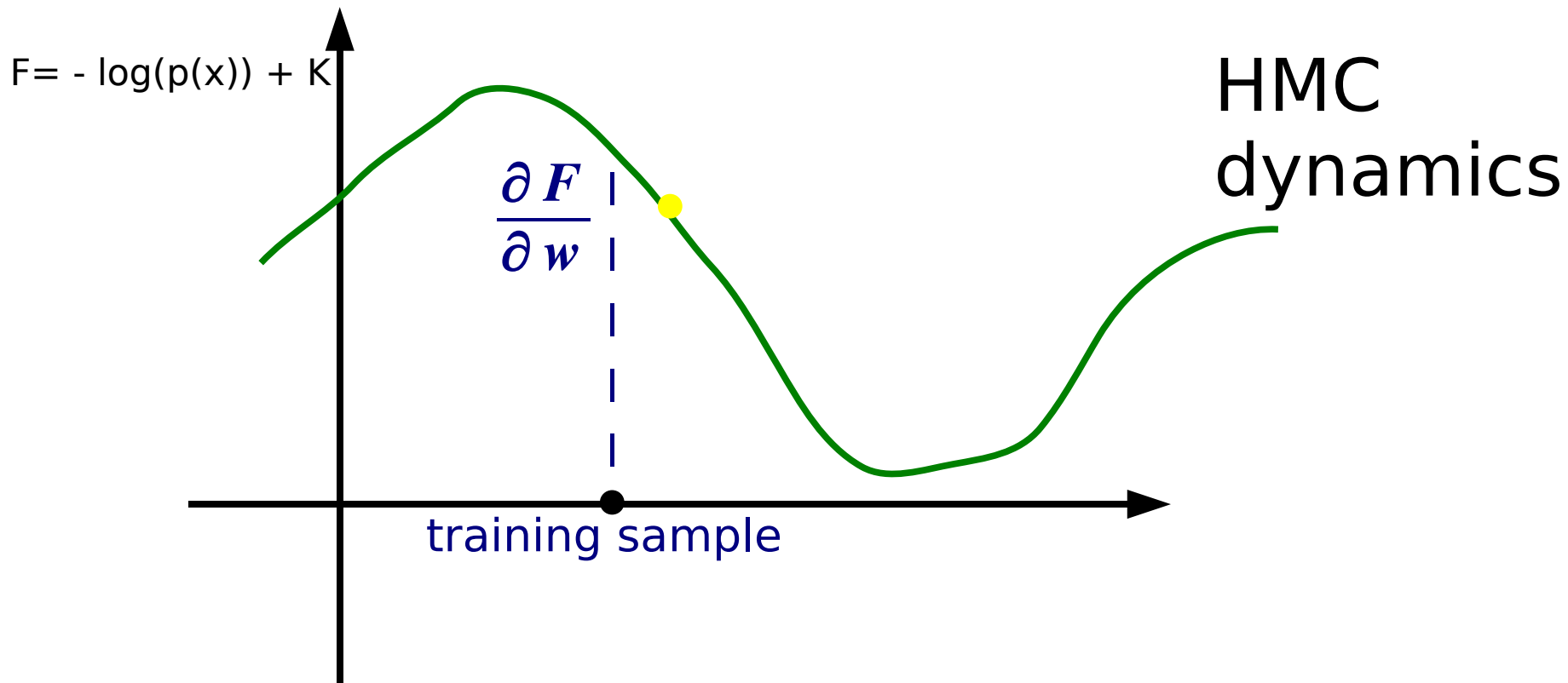
Learning



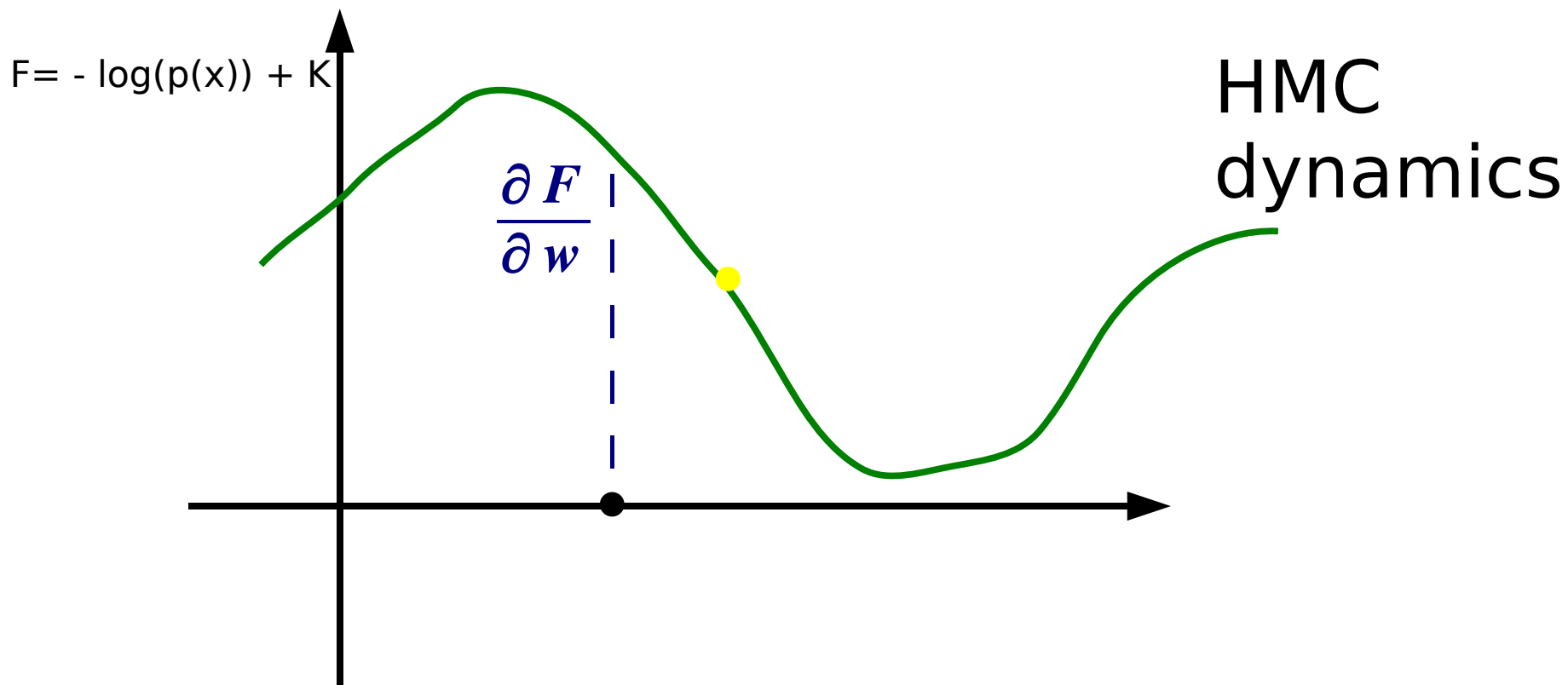
Learning



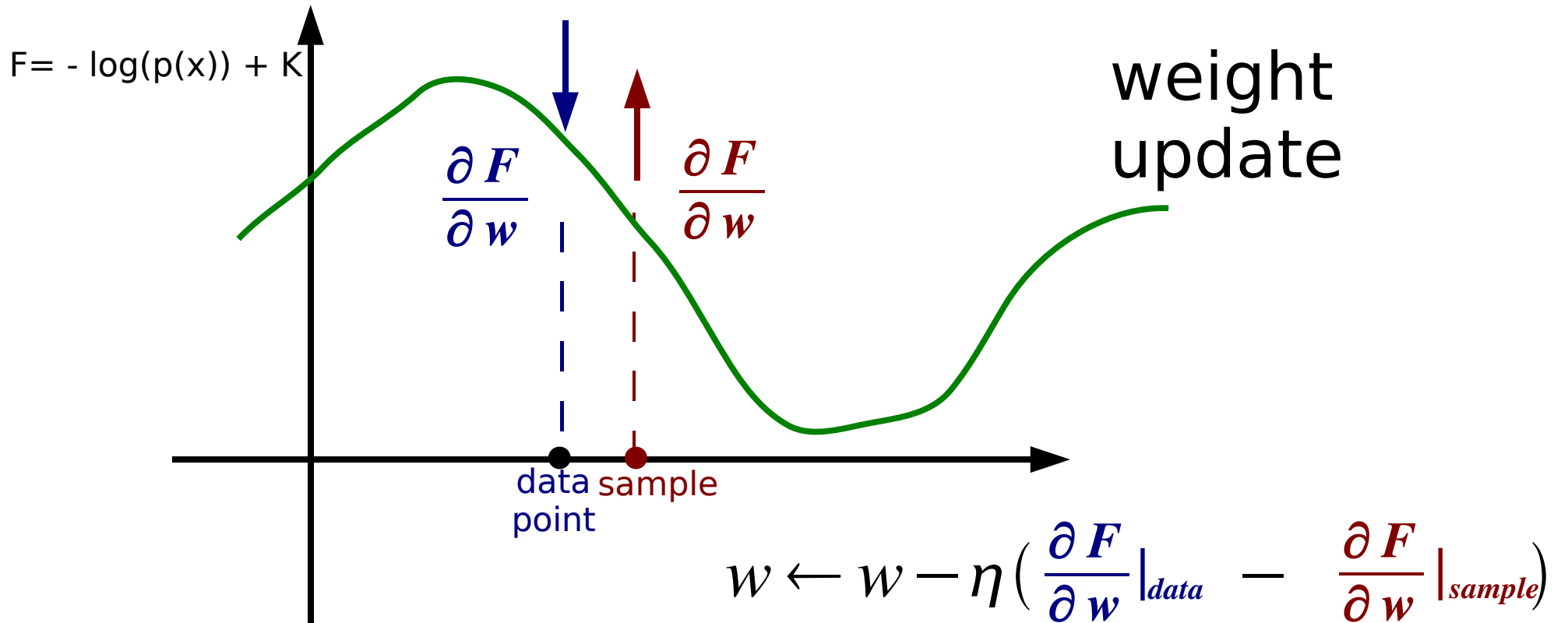
Learning



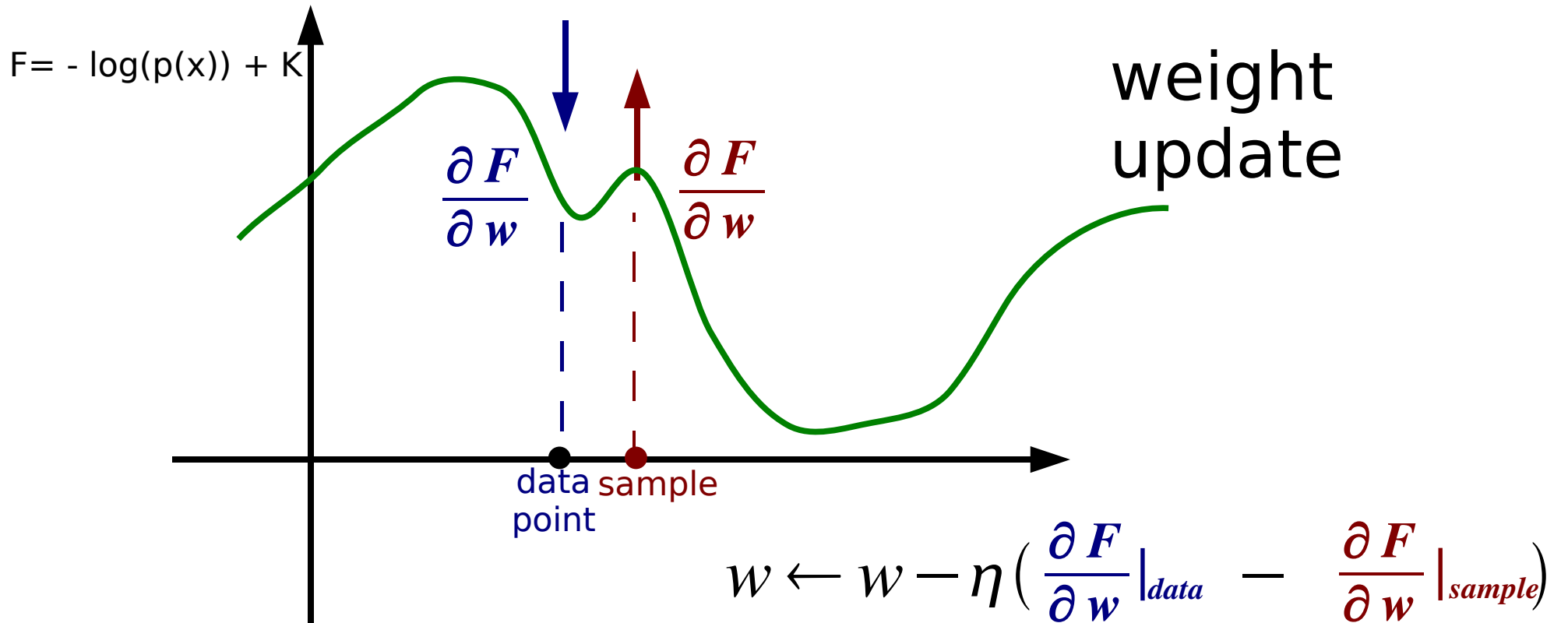
Learning



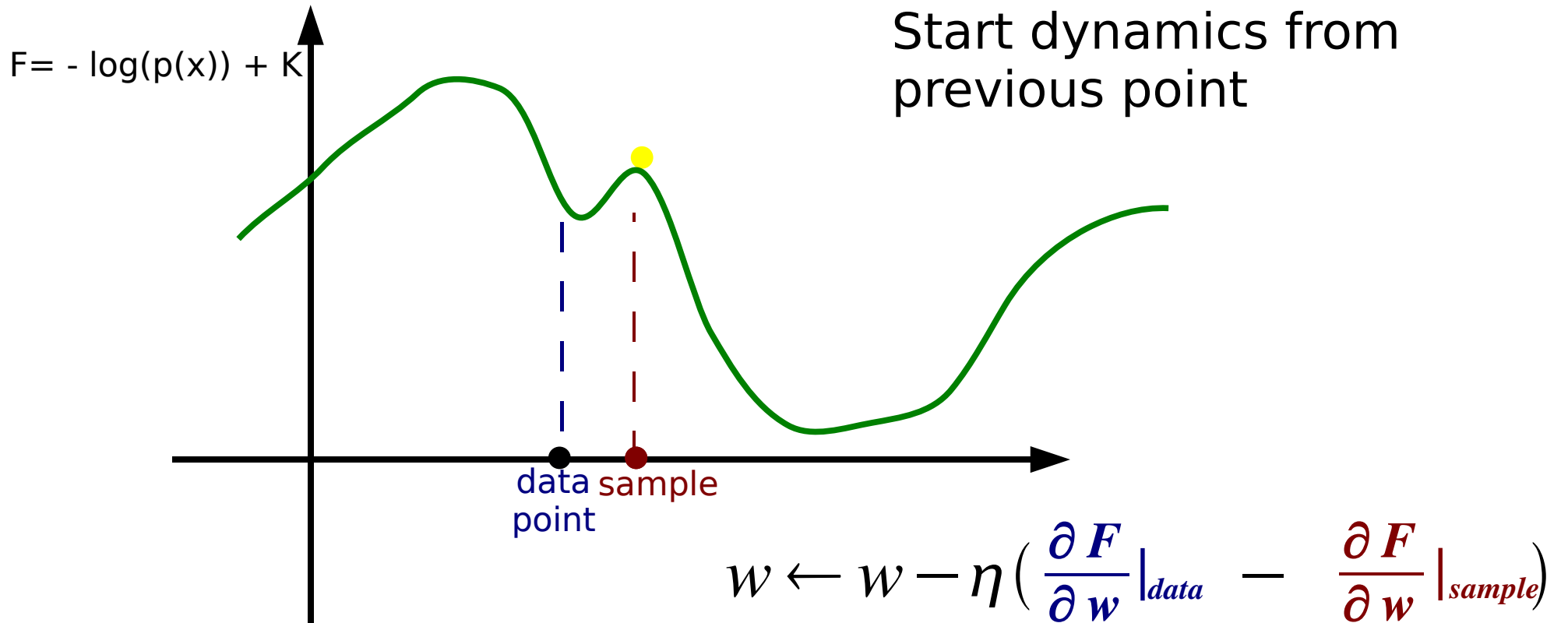
Learning



Learning



Learning



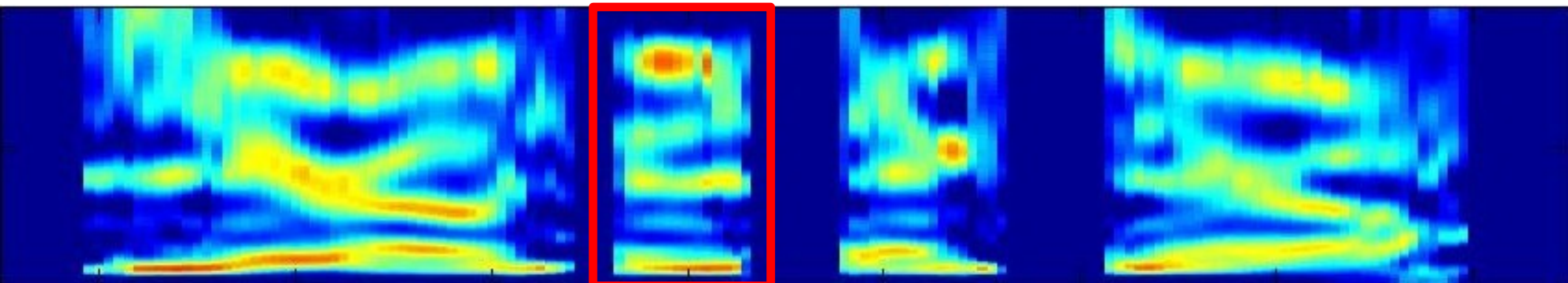
Outline

- mathematical formulation of the model
- training
- **learning acoustic features for speech recognition**
- generation of natural images
- recognition of facial expression under occlusion
- conclusion

Speech Recognition on TIMIT

INPUT

- frame: 25ms Hamming window
- 10ms overlap between frames
- 39 log-filterbank outputs + energy per frame
- 15 frames
- PCA whitening → 384 components
- no deltas, no delta-deltas...

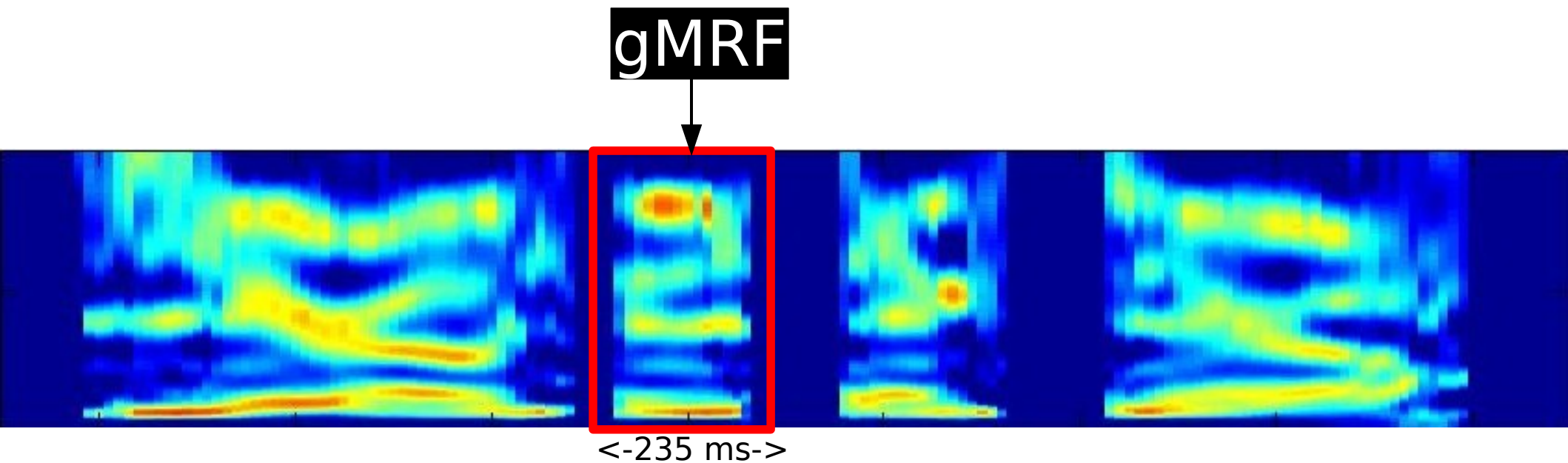


<-235 ms->

Dahl, Ranzato, Mohamed, Hinton, NIPS 2010

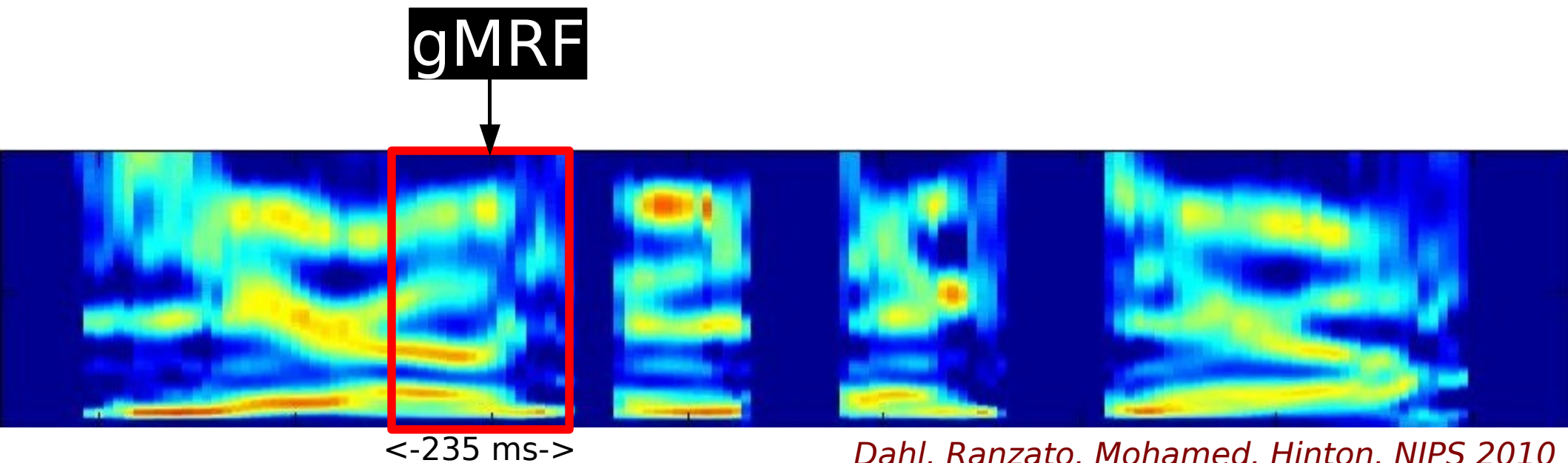
Speech Recognition on TIMIT

gated MRF with 1024 precision and 512 mean units
trained on independent windows



Speech Recognition on TIMIT

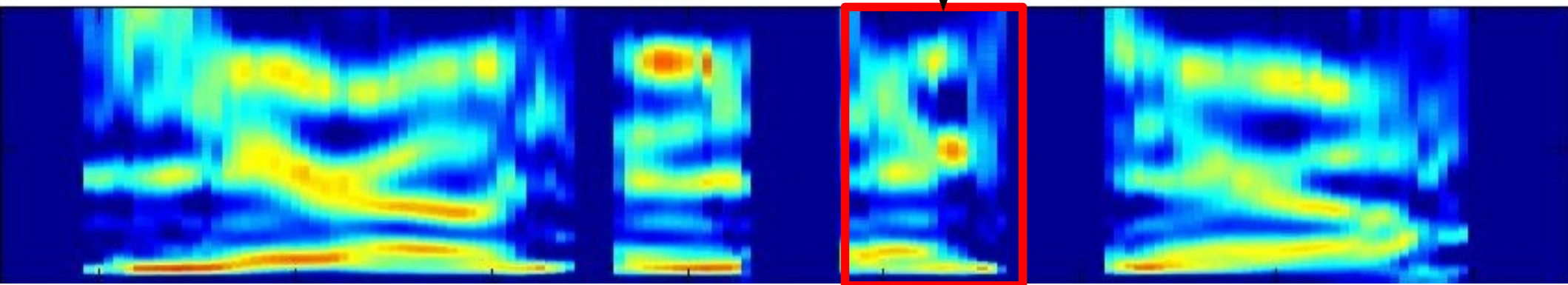
gated MRF with 1024 precision and 512 mean units
trained on independent windows



Speech Recognition on TIMIT

gated MRF with 1024 precision and 512 mean units
trained on independent windows

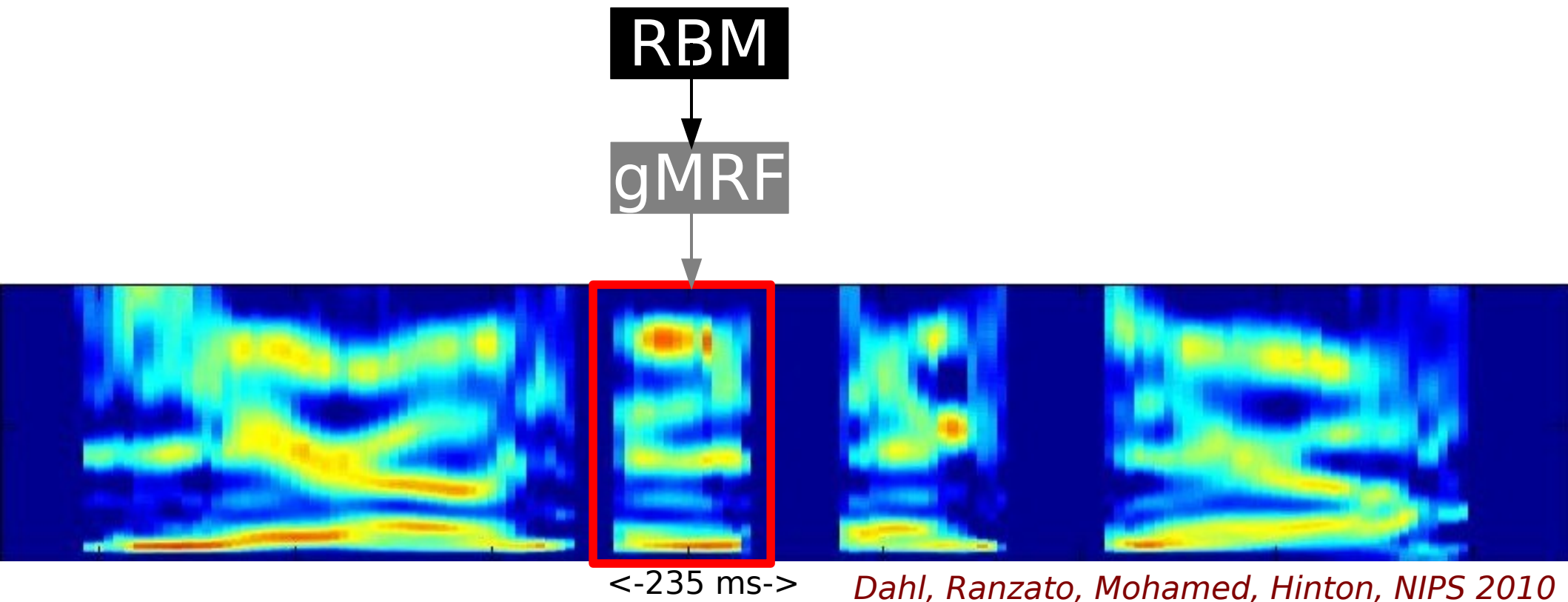
gMRF



Dahl, Ranzato, Mohamed, Hinton, NIPS 2010

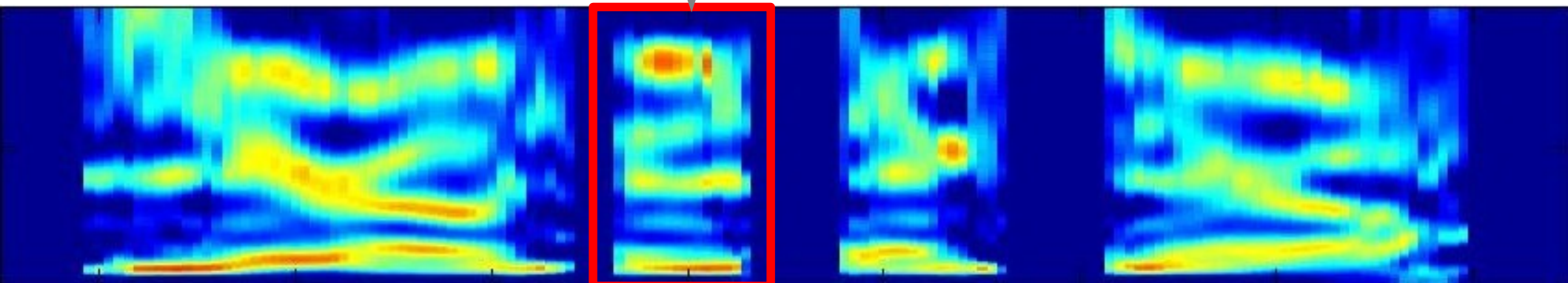
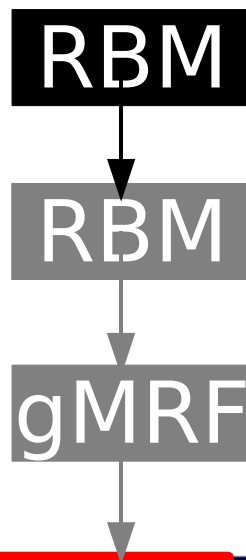
Speech Recognition on TIMIT

- For each training sample, infer latent variables of gated MRF
- Train 2nd layer binary-binary RBM (2048 units)



Speech Recognition on TIMIT

- Train 1st layer gated MRF
- Train 2nd layer binary-binary RBM (2048 units)
- Train 3rd layer binary-binary RBM (2048 units)



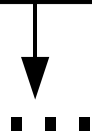
<-235 ms->

Dahl, Ranzato, Mohamed, Hinton, NIPS 2010

Speech Recognition on TIMIT

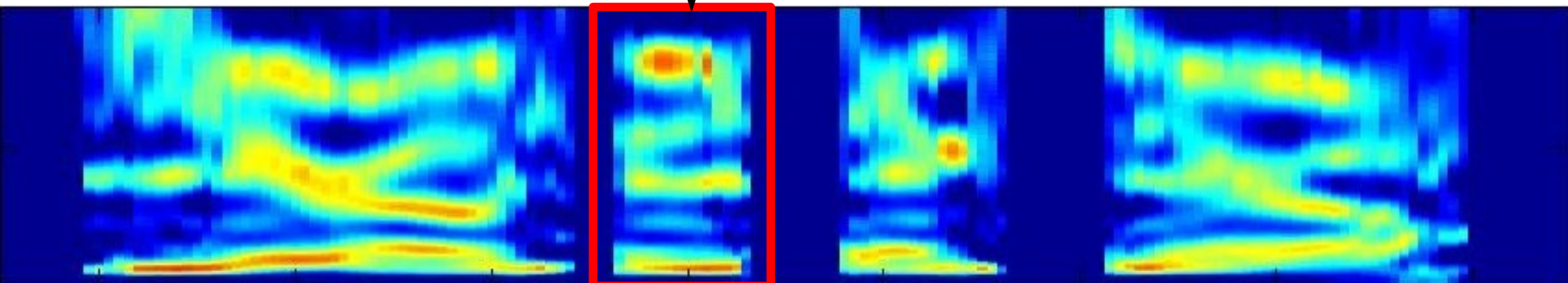
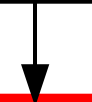
8 layer (deep) model

RBM



RBM

gMRF

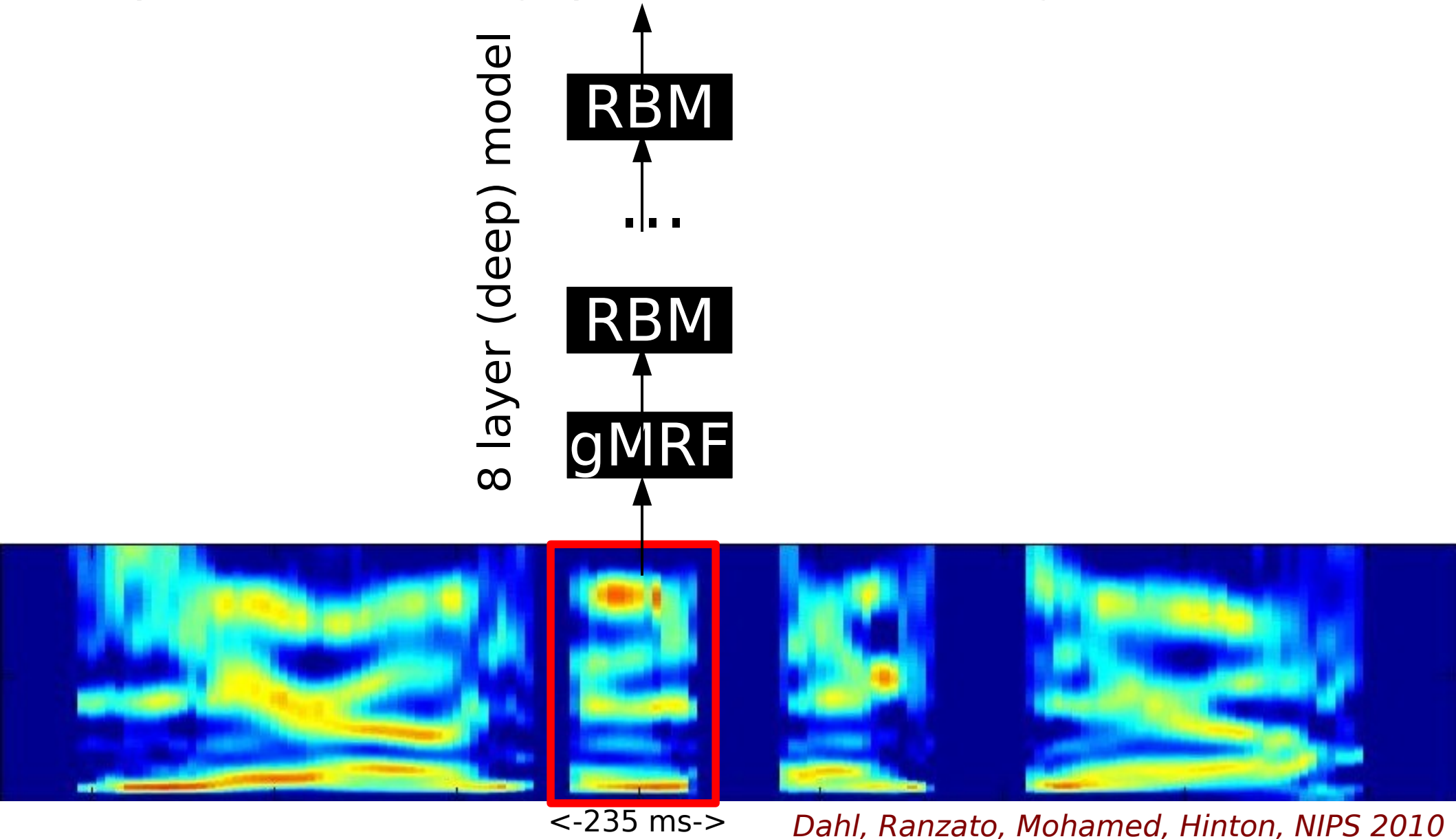


<-235 ms->

Dahl, Ranzato, Mohamed, Hinton, NIPS 2010

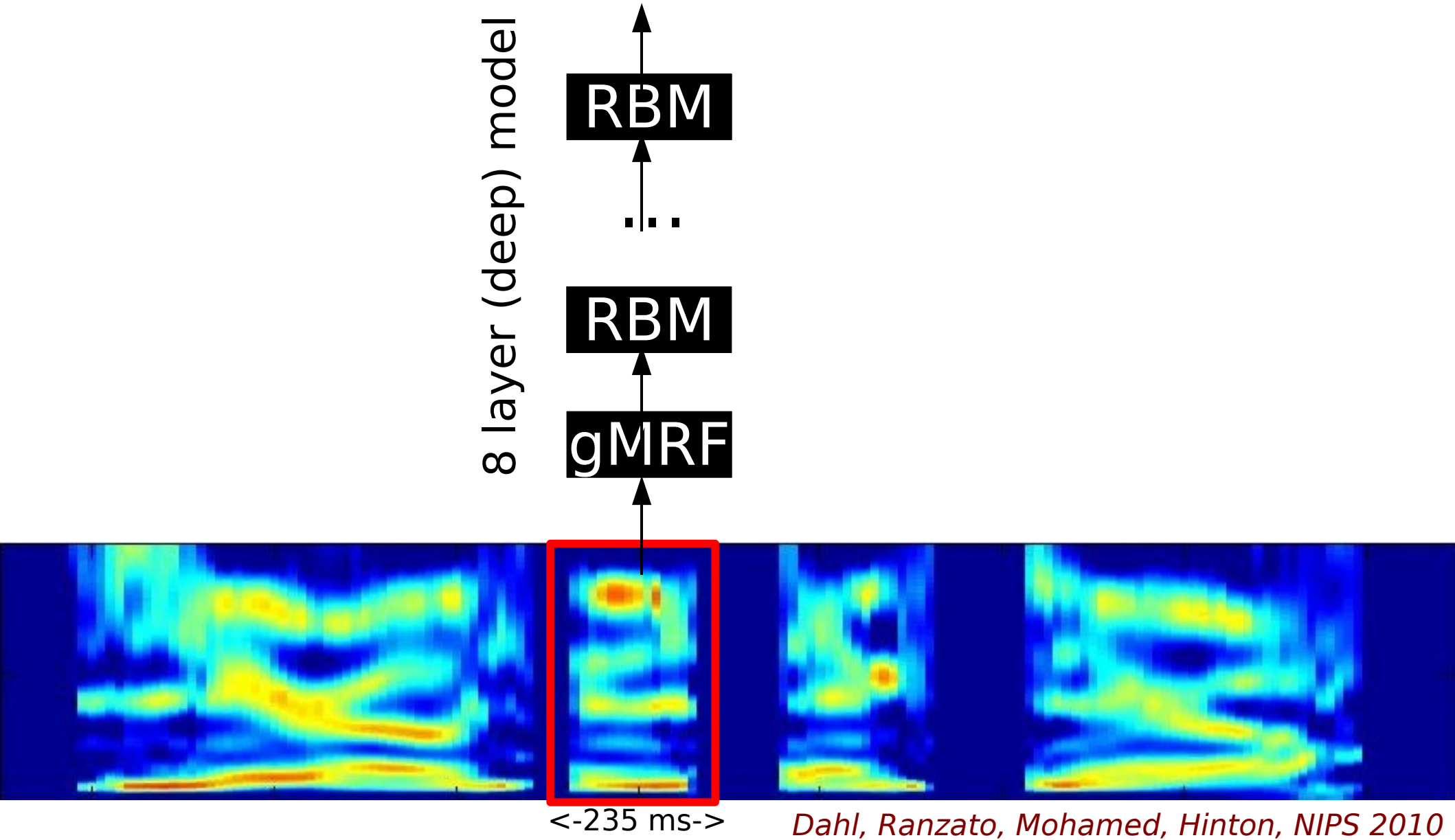
Speech Recognition on TIMIT

Supervised training: predict states of aligned HMM



Speech Recognition on TIMIT

Test: predict states of HMM → decoding (bigram language model)



Speech Recognition on TIMIT

METHOD	PER
CRF	34.8%
Large-Margin GMM	33.0%
CD-HMM	27.3%
Augmented CRF	26.6%
RNN	26.1%
Bayesian Triphone HMM	25.6%
Triphone HMM discrim. trained	22.7%
DBN with gated MRF	20.5%

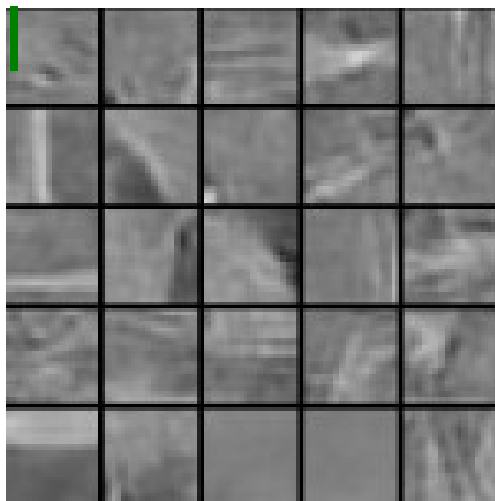
Speech Recognition on TIMIT

METHOD	PER	Year
CRF	34.8%	2008
Large-Margin GMM	33.0%	2006
CD-HMM	27.3%	2009
Augmented CRF	26.6%	2009
RNN	26.1%	1994
Bayesian Triphone HMM	25.6%	1998
Triphone HMM discrim. trained	22.7%	2009
DBN with gated MRF	20.5%	2010

Outline

- mathematical formulation of the model
- training
- learning acoustic features for speech recognition
- **generation of natural images**
- recognition of facial expression under occlusion
- conclusion

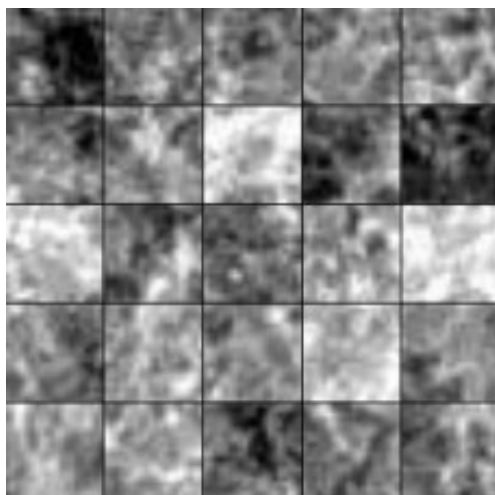
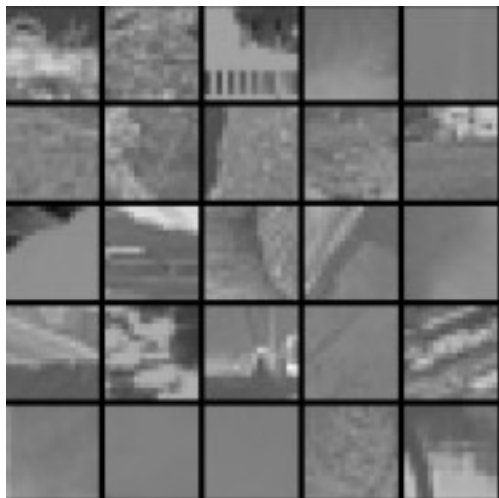
Generation natural
image patches



mcRBM

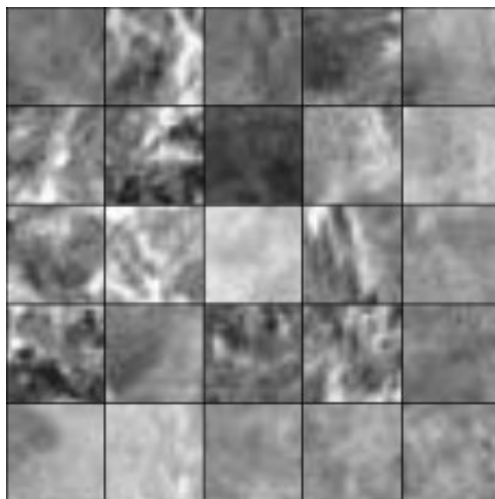
Ranzato and Hinton CVPR 2010

Natural images



GRBM

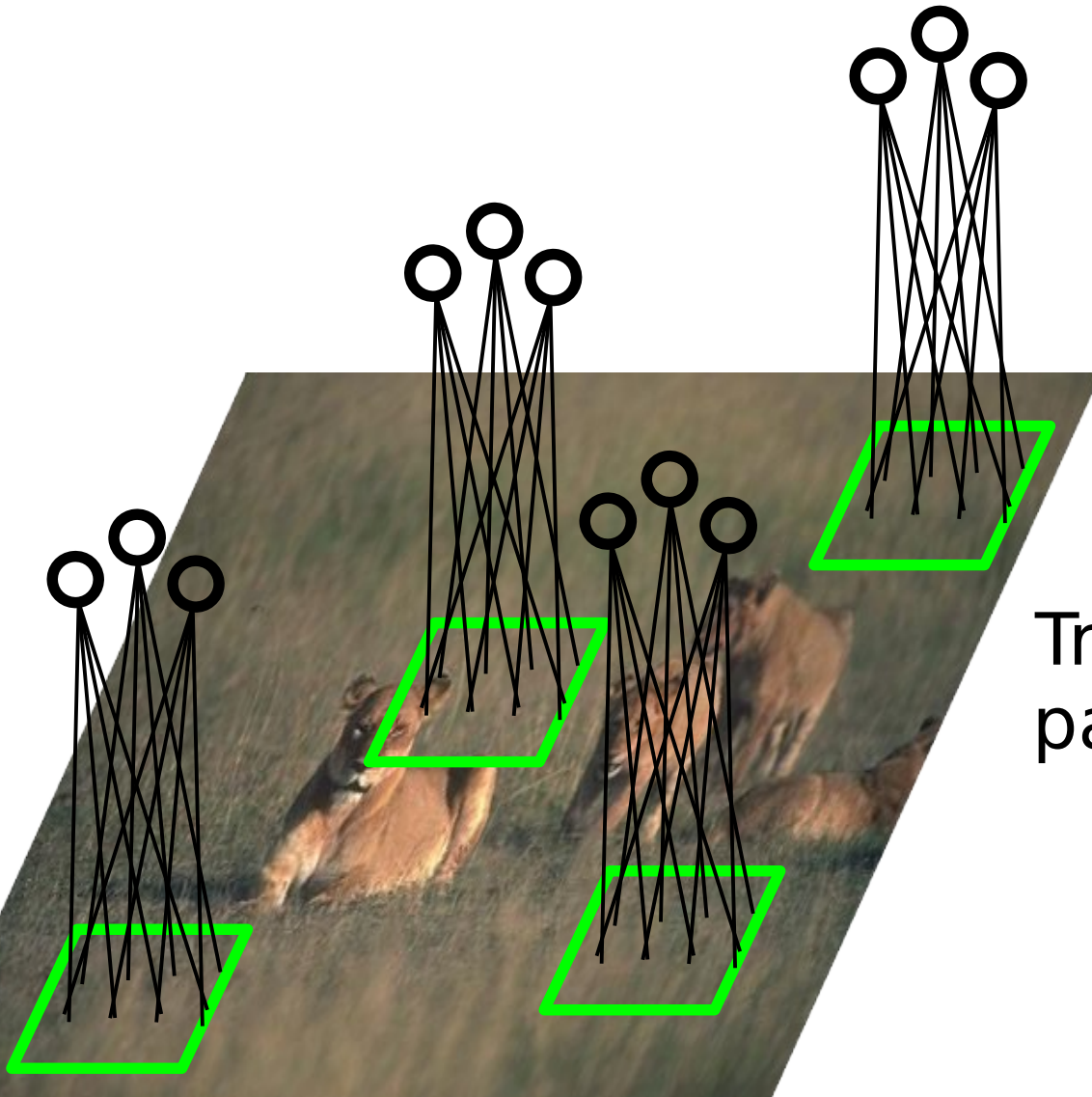
from Osindero and Hinton NIPS 2008



S-RBM + DBN

from Osindero and Hinton NIPS 2008

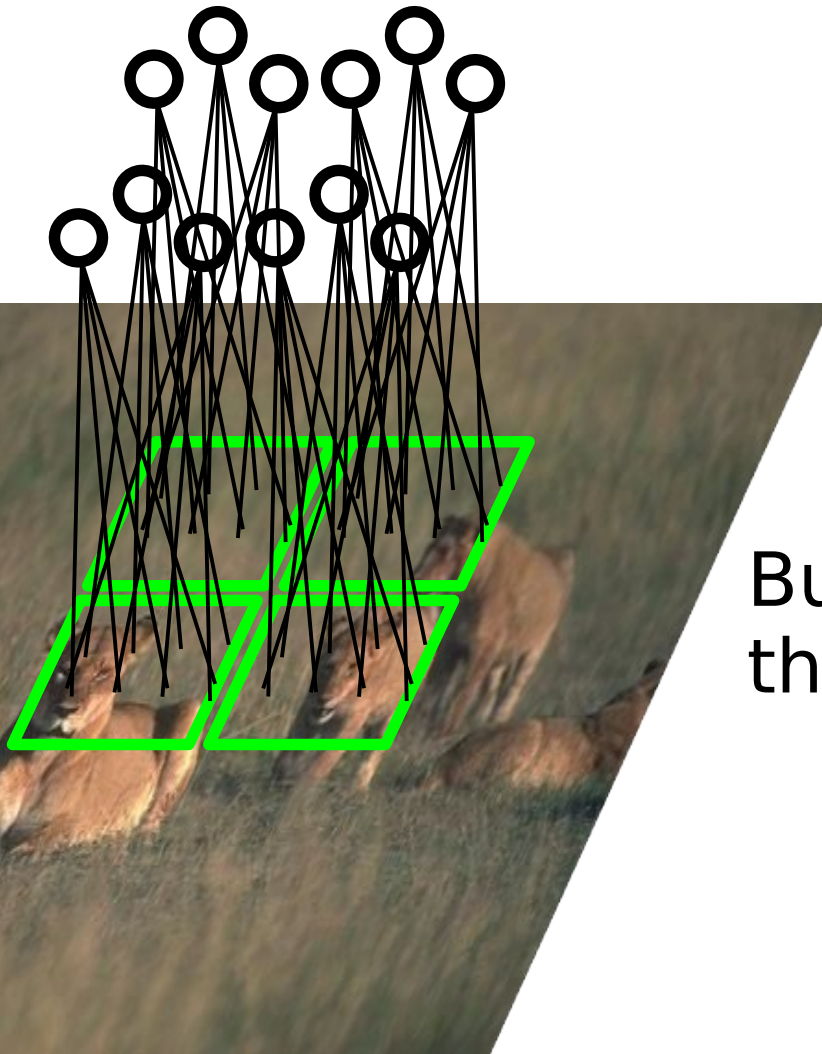
From Patches to High-Resolution Images



Training by picking patches at random

From Patches to High-Resolution Images

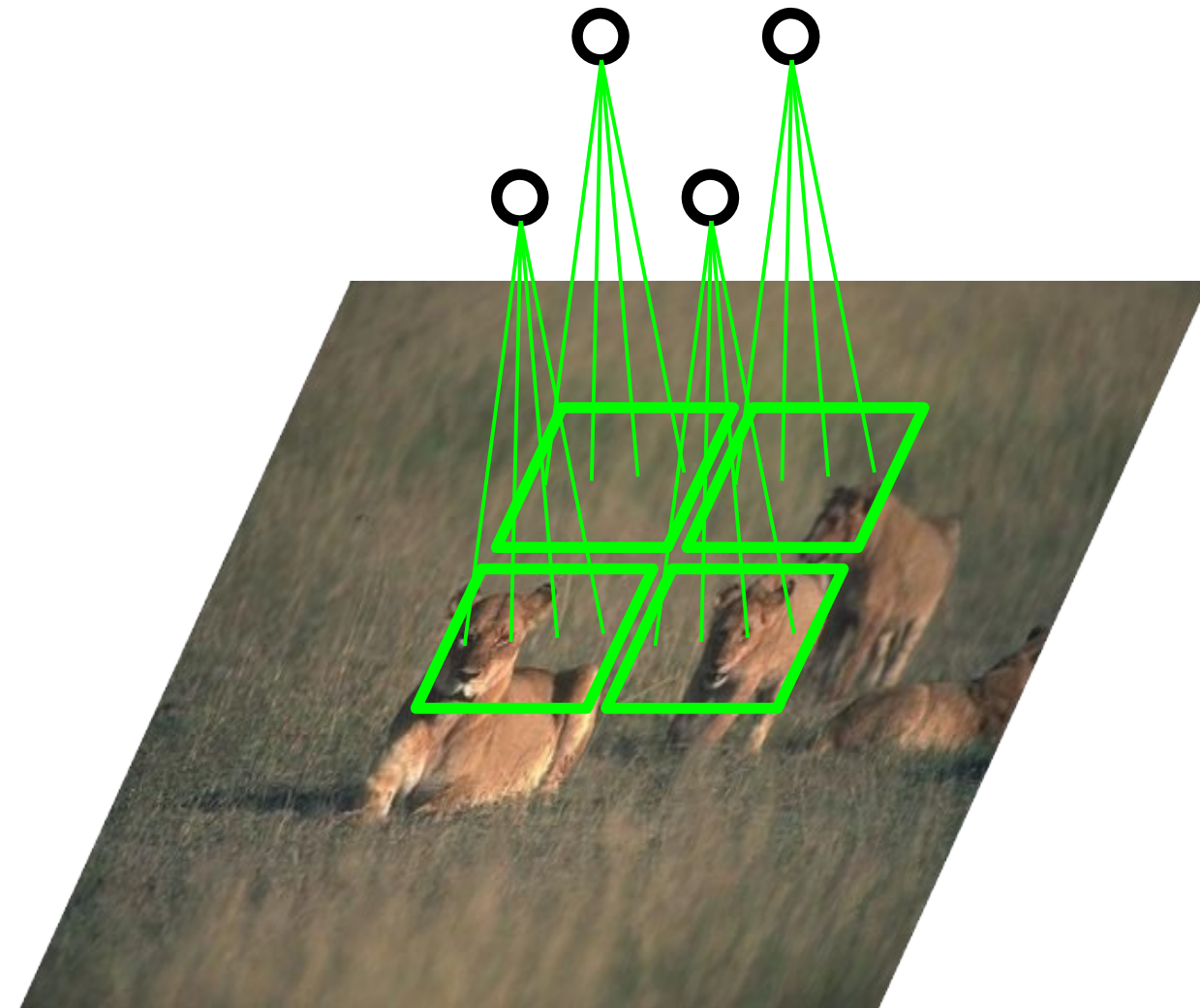
This is not a good way to extend the model to big images: block artifacts



But we could also take them from a grid

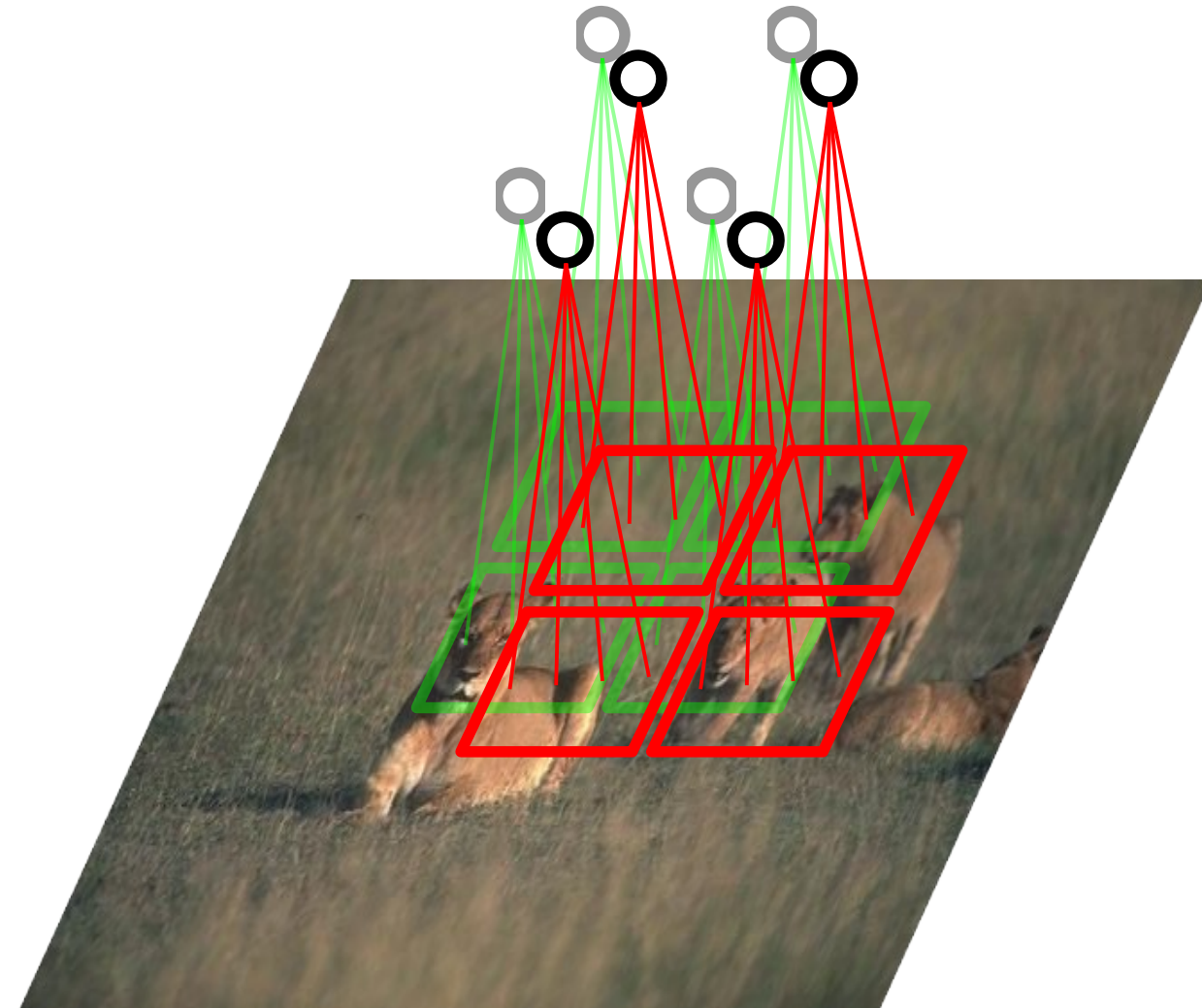
From Patches to High-Resolution Images

IDEA: have one subset of filters applied to these locations,



From Patches to High-Resolution Images

IDEA: have one subset of filters applied to these locations, another subset to these locations

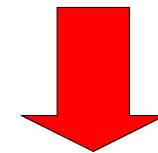


From Patches to High-Resolution Images

IDEA: have one subset of filters applied to these locations, another subset to these locations, etc.



Train jointly all parameters.

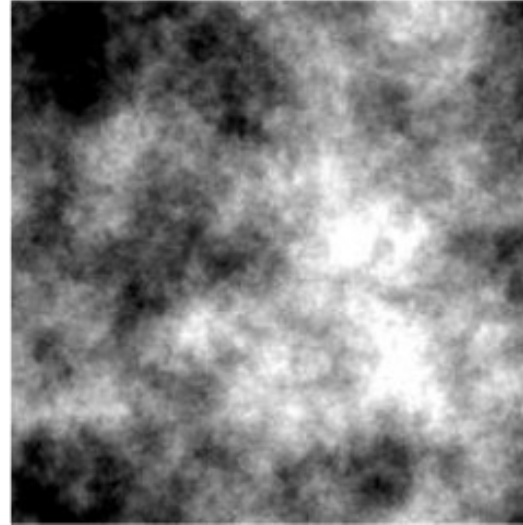


No block artifacts
Reduced redundancy

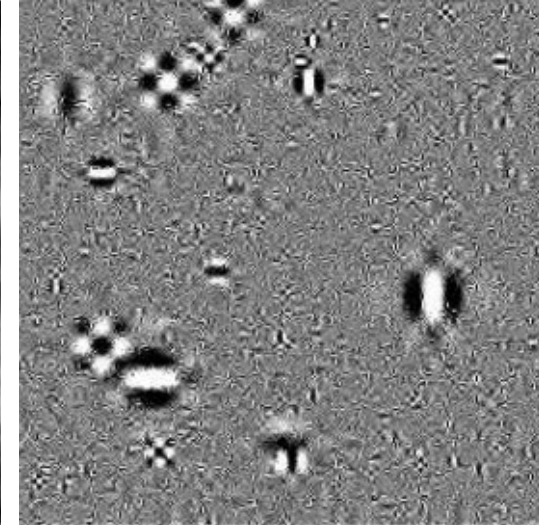
*Gregor LeCun arXiv 2010
Ranzato, Mnih, Hinton NIPS 2010*

Sampling High-Resolution Images

Gaussian model



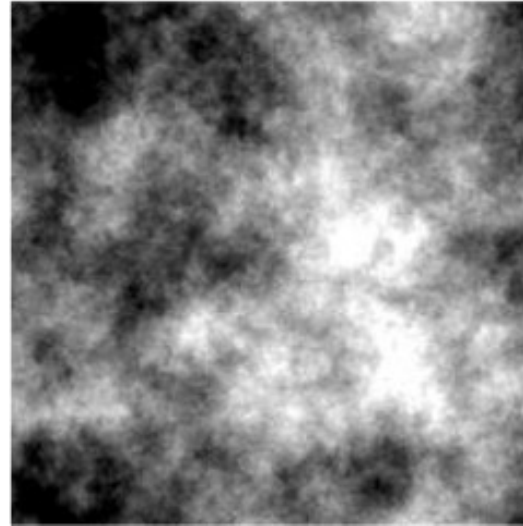
marginal wavelet



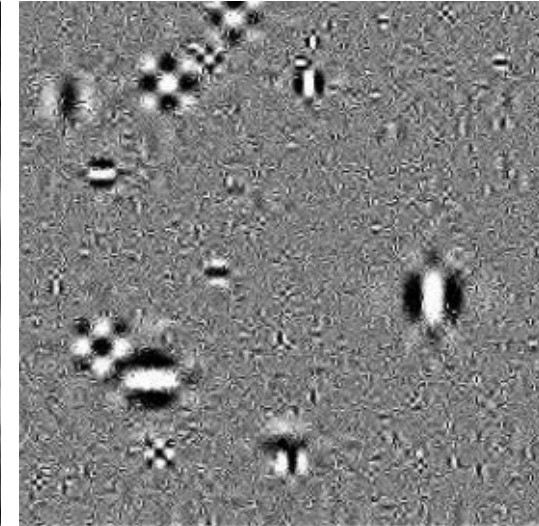
from Simoncelli 2005

Sampling High-Resolution Images

Gaussian model

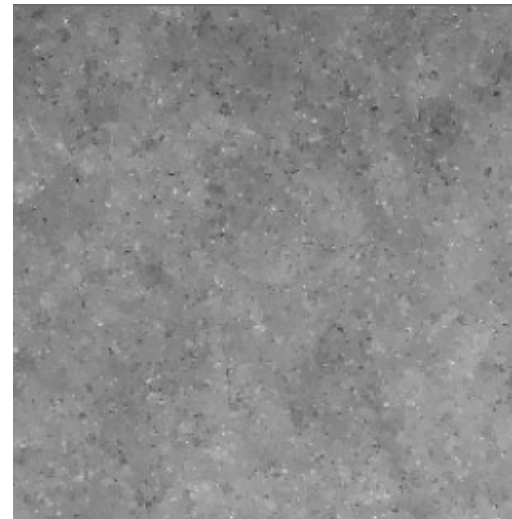


marginal wavelet

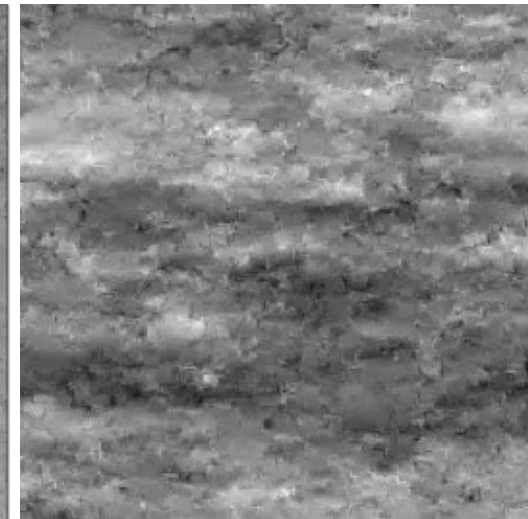


from Simoncelli 2005

Pair-wise MRF



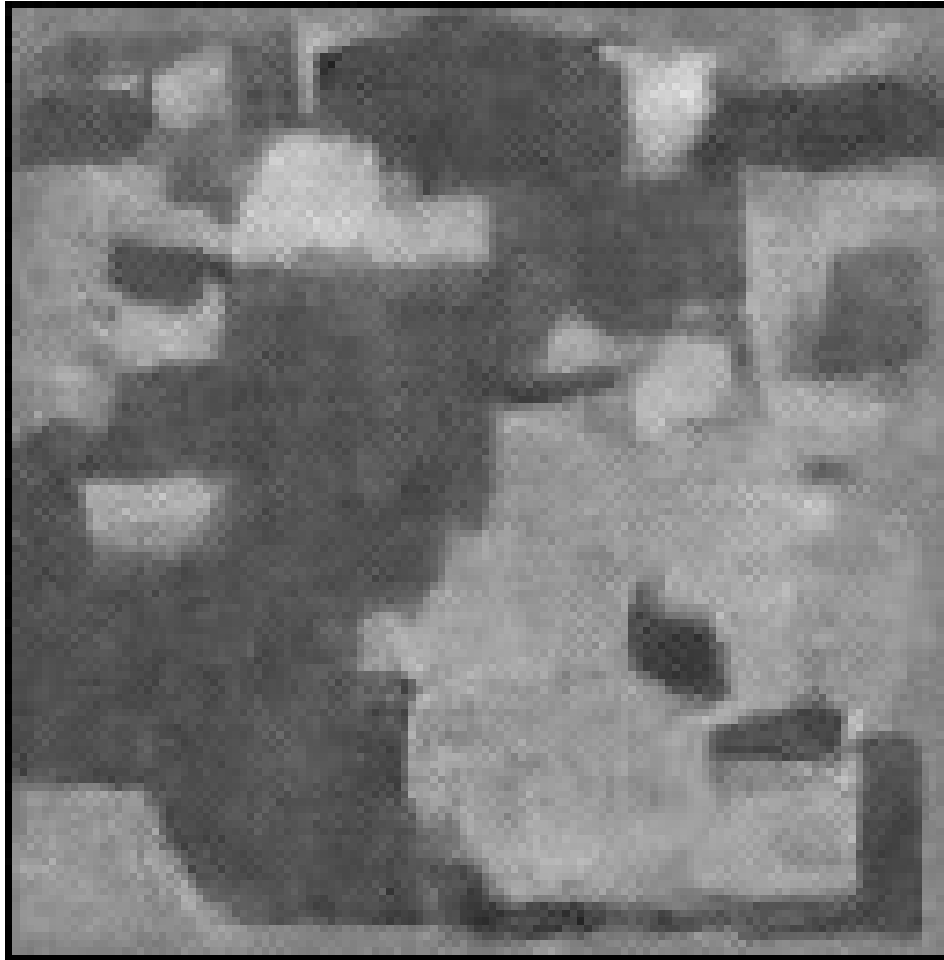
FoE



from Schmidt, Gao, Roth CVPR 2010

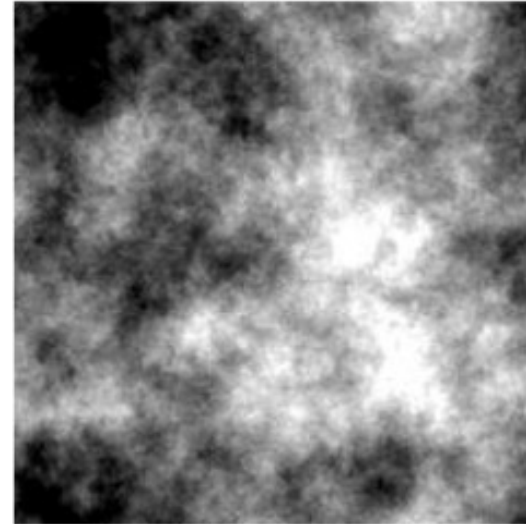
Sampling High-Resolution Images

Mean Covariance Model



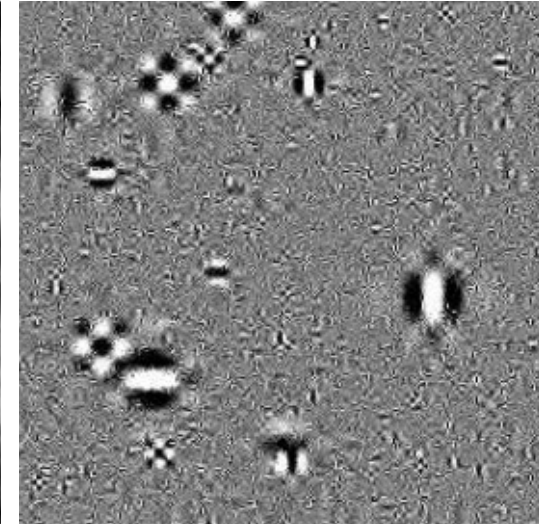
Ranzato, Mnih, Hinton NIPS 2010

Gaussian model

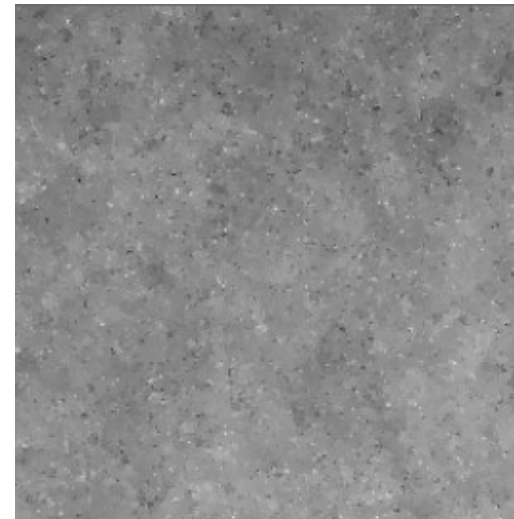


from Simoncelli 2005

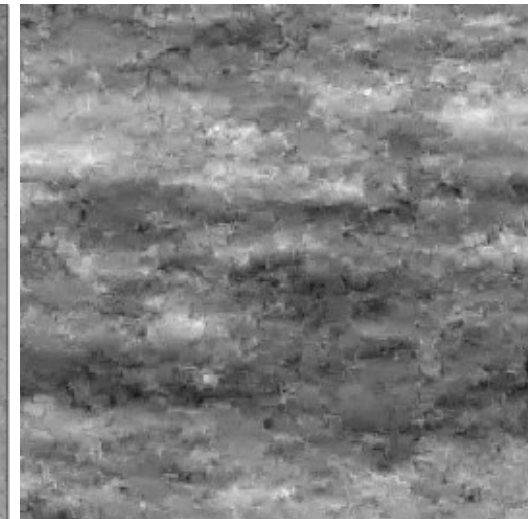
marginal wavelet



Pair-wise MRF



FoE



from Schmidt, Gao, Roth CVPR 2010

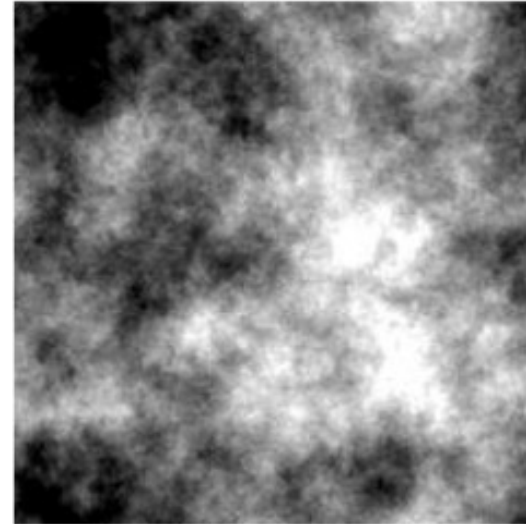
Sampling High-Resolution Images

Mean Covariance Model



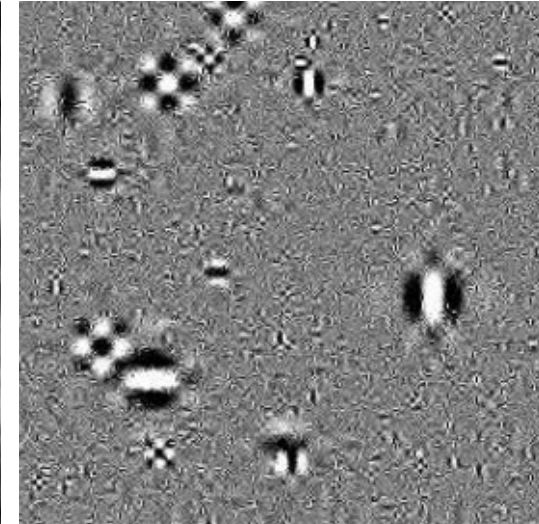
Ranzato, Mnih, Hinton NIPS 2010

Gaussian model

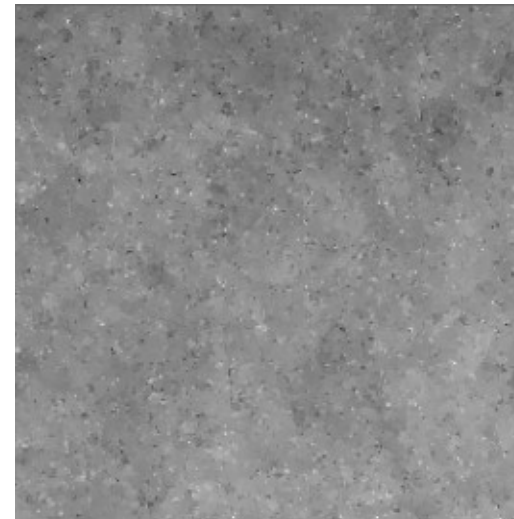


from Simoncelli 2005

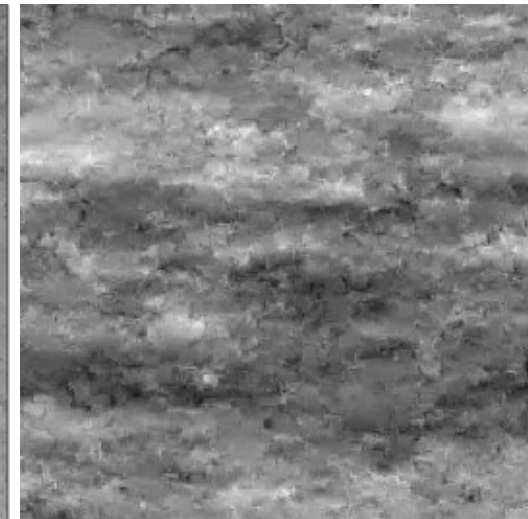
marginal wavelet



Pair-wise MRF



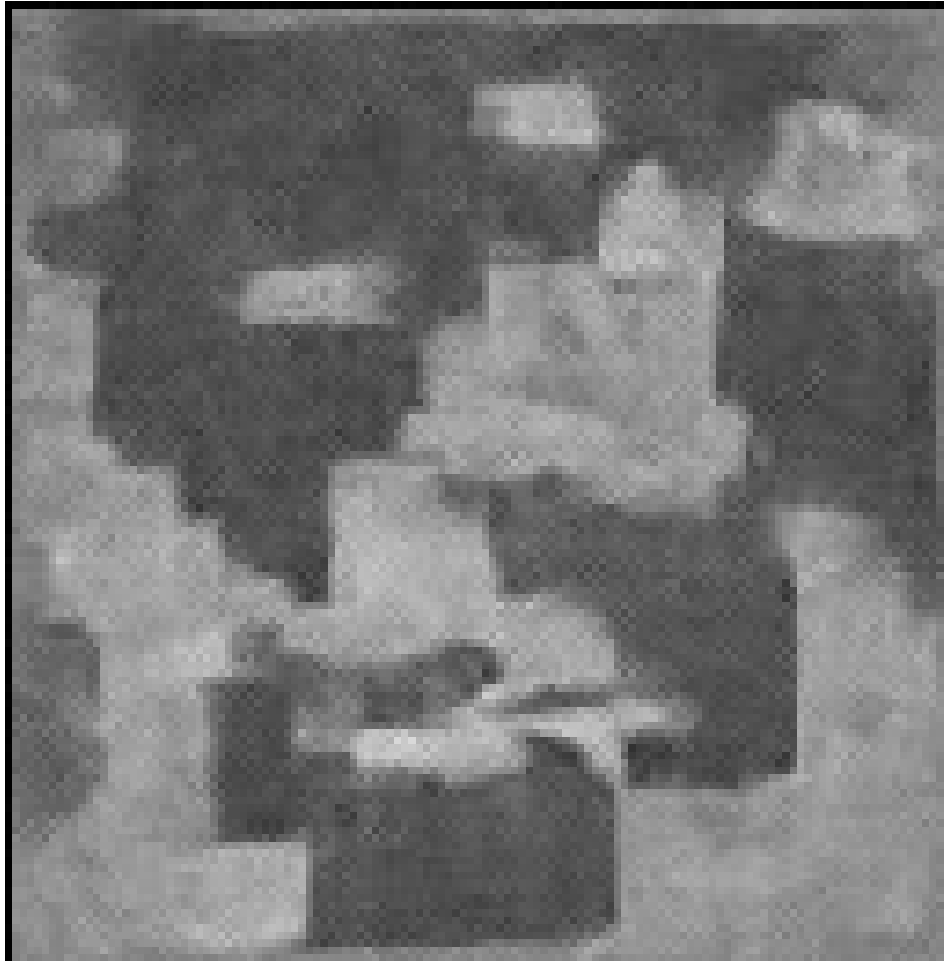
FoE



from Schmidt, Gao, Roth CVPR 2010

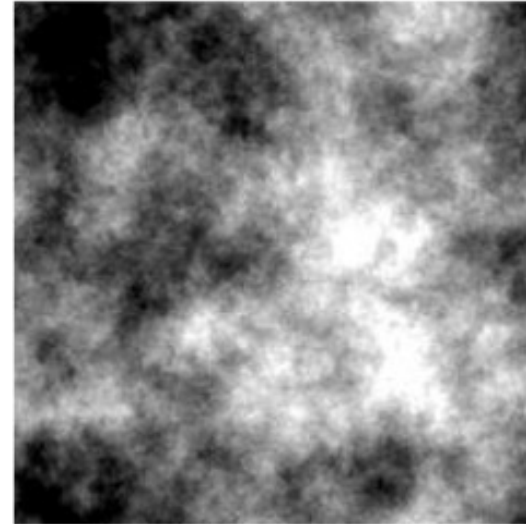
Sampling High-Resolution Images

Mean Covariance Model



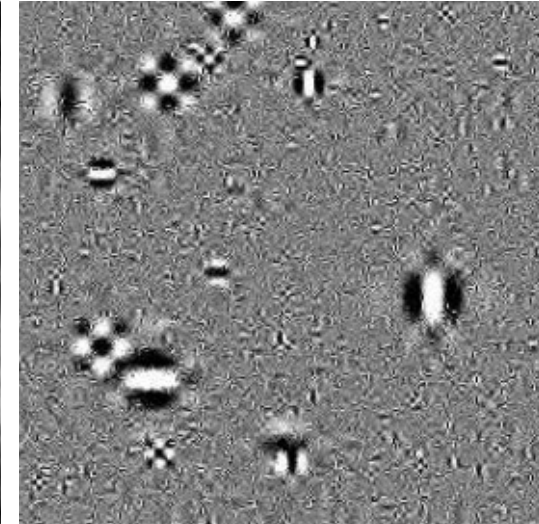
Ranzato, Mnih, Hinton NIPS 2010

Gaussian model

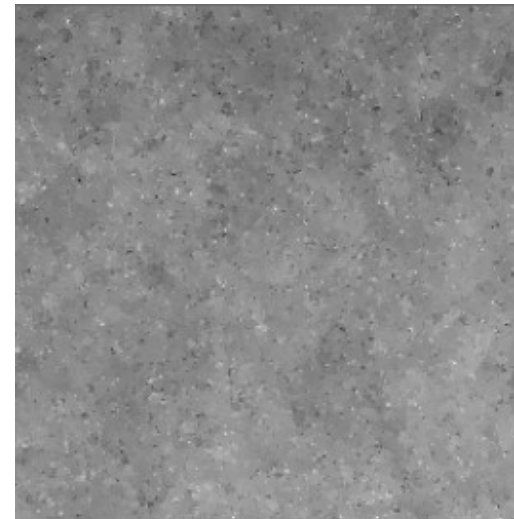


from Simoncelli 2005

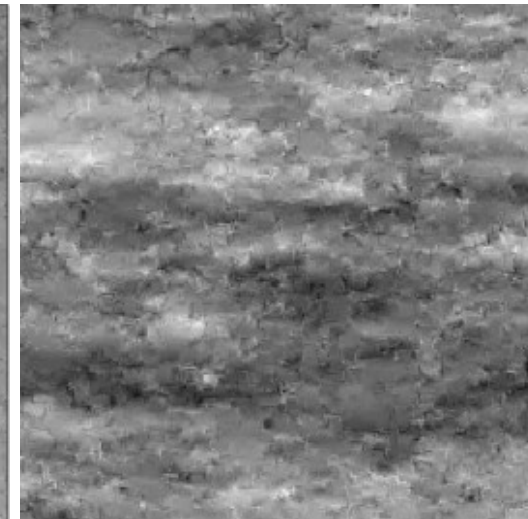
marginal wavelet



Pair-wise MRF



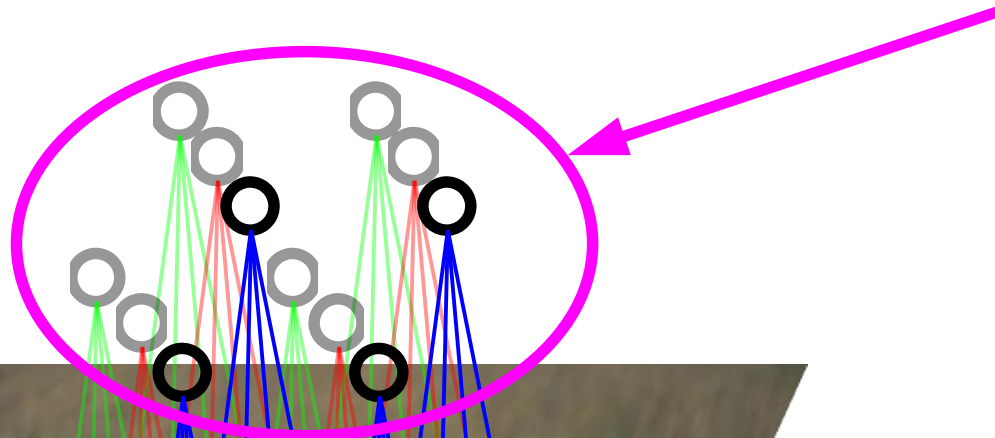
FoE



from Schmidt, Gao, Roth CVPR 2010

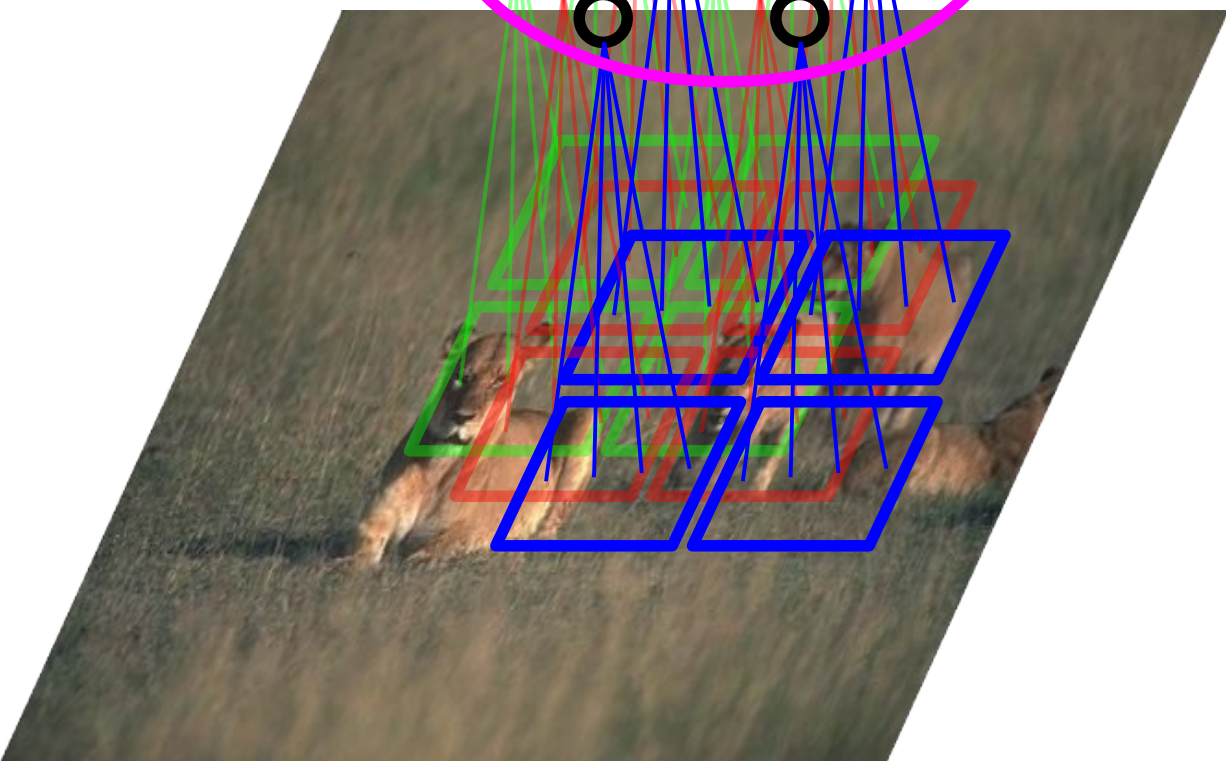
Making the model.. “DEEPER “

Treat these units as data
to train a similar model on the top



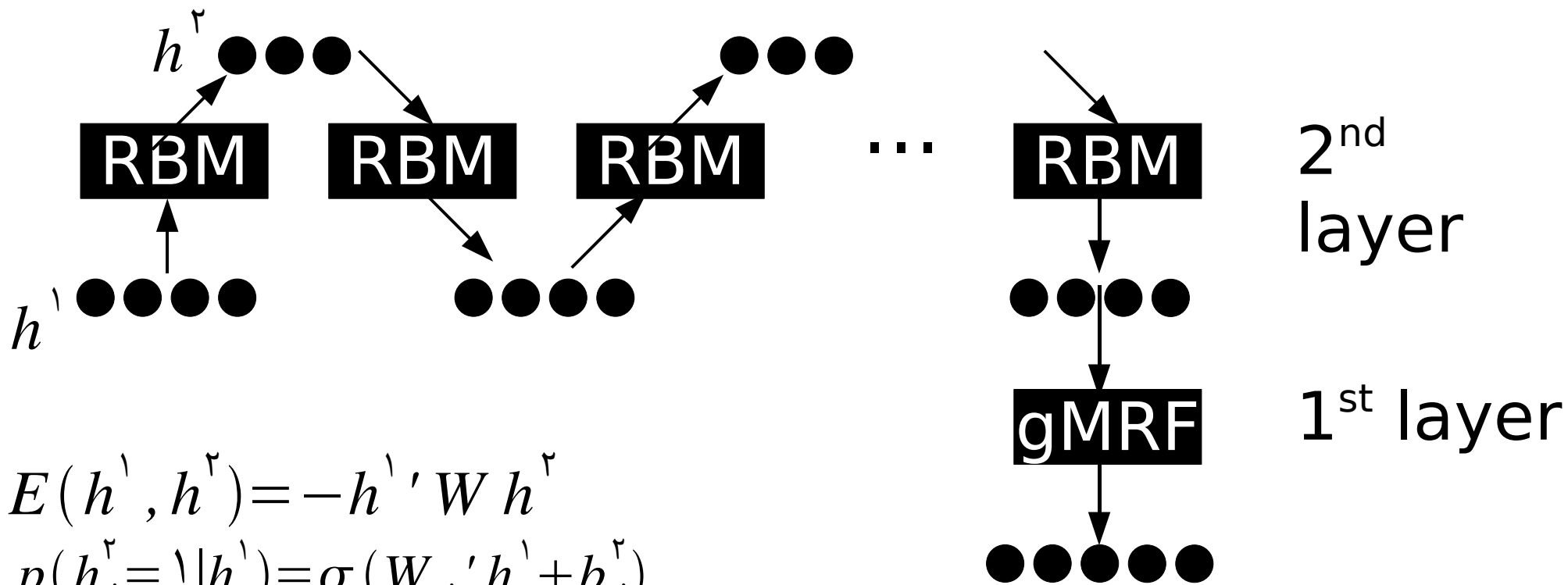
SECOND STAGE

Field of binary RBM's.
Each hidden unit has a
receptive field of 30x30
pixels in input space.



Sampling from the DEEPER model

- Sample from 2nd layer
- project sample in image space using 1st layer $p(x|h)$



$$E(h^l, h^r) = -h^{l'} W h^r$$

$$p(h_j^r = \lambda | h^l) = \sigma(W_j^r h^l + b_j^r)$$

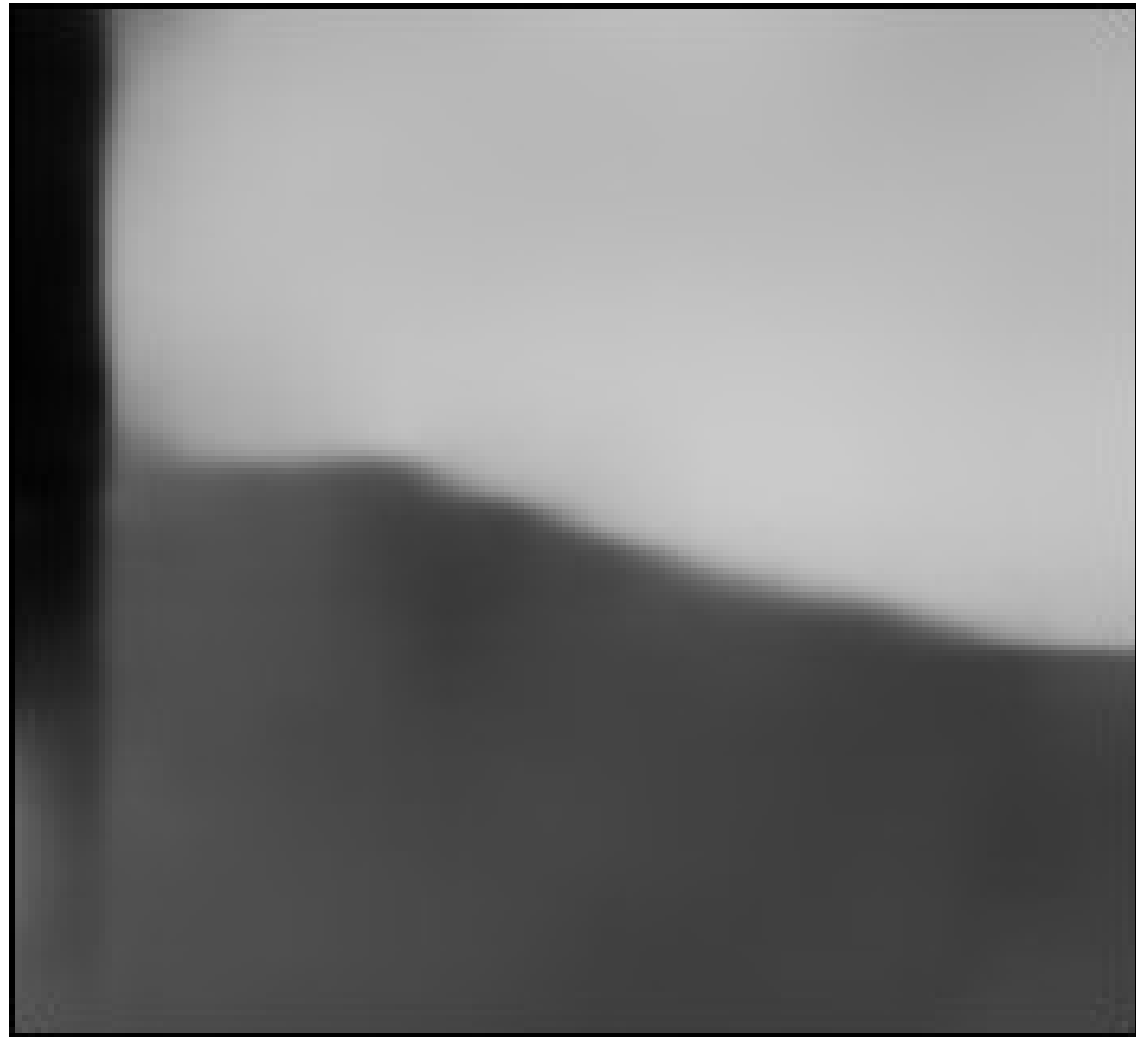
$$p(h_k^l = \lambda | h^r) = \sigma(W_k^l h^r + b_k^l)$$

Samples from Deep Generative Model

1st stage model



3rd stage model



Samples from Deep Generative Model

1st stage model



3rd stage model



Samples from Deep Generative Model

1st stage model



3rd stage model

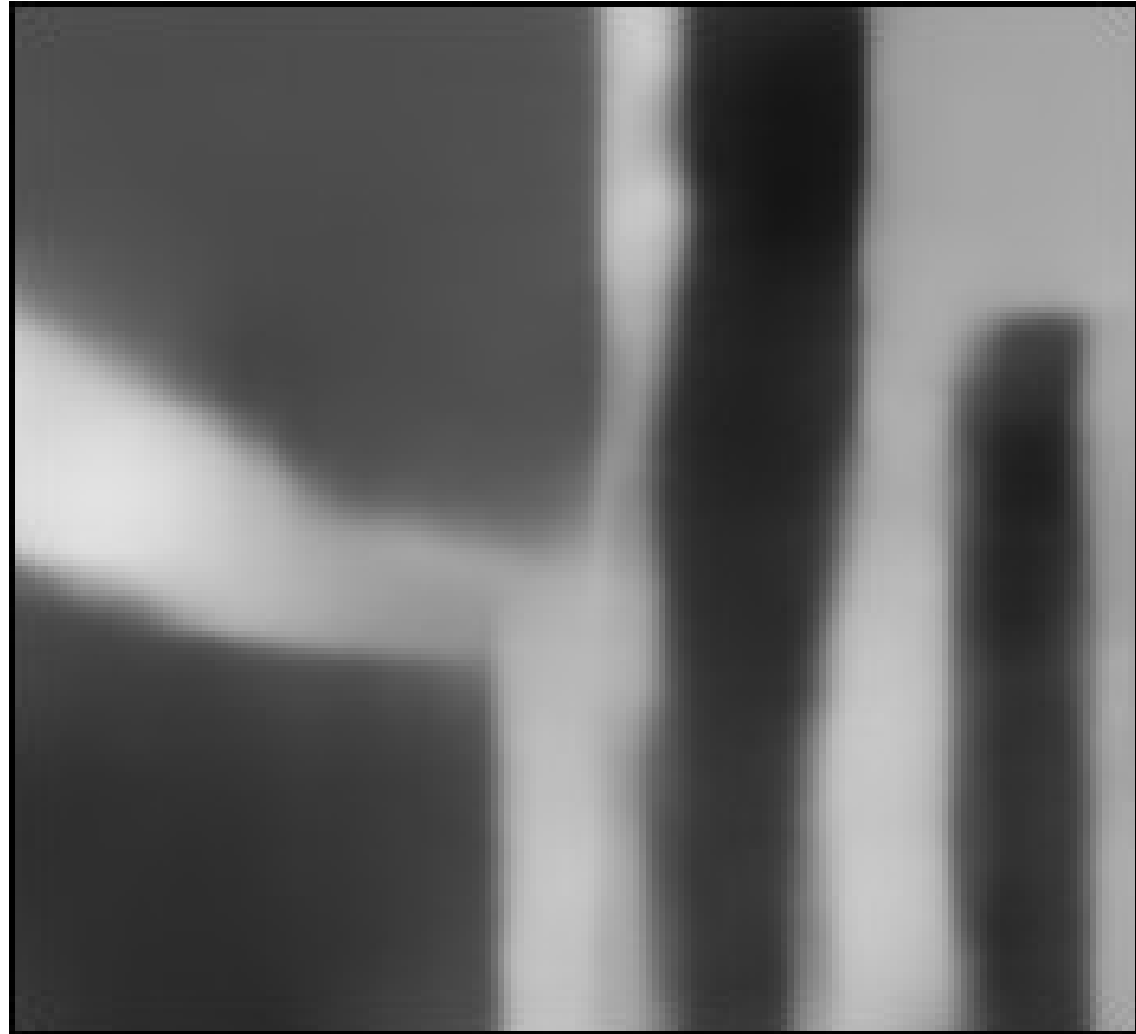


Samples from Deep Generative Model

1st stage model



3rd stage model



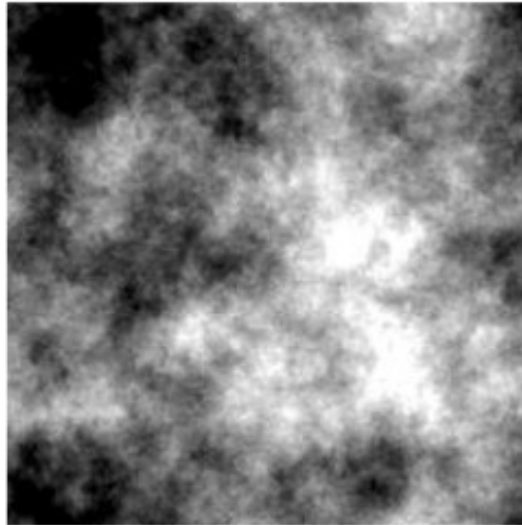
Sampling High-Resolution Images

FoE



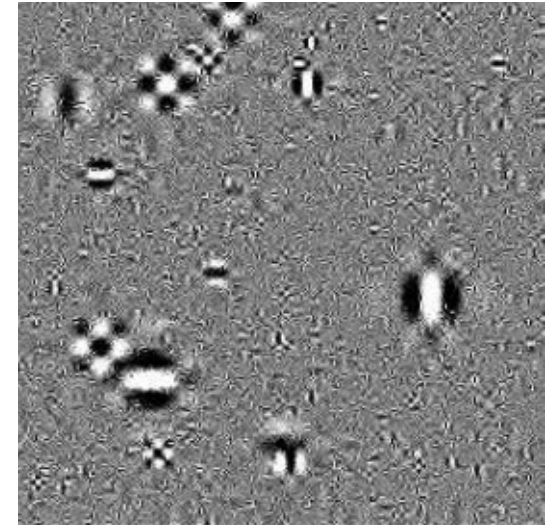
from Schmidt, et al CVPR 2010

Gaussian model



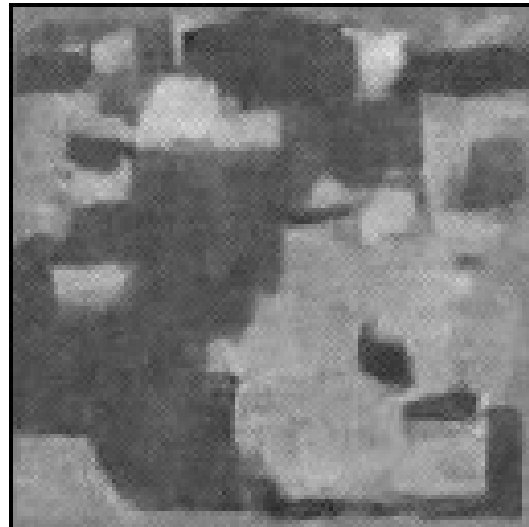
from Simoncelli 2005

marginal wavelet



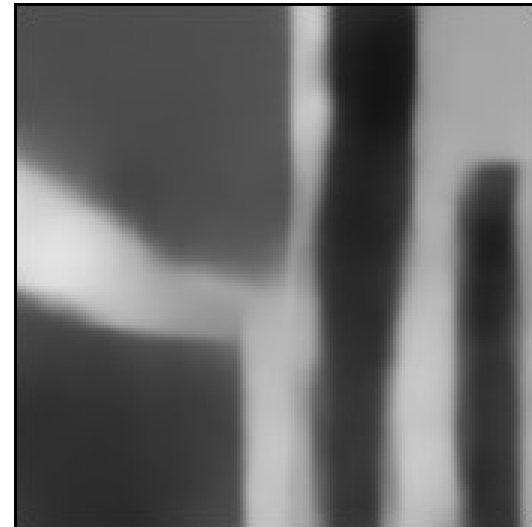
from Simoncelli 2005

Deep - 1



Ranzato, Mnih, Hinton NIPS 2010

Deep - 3 layers



Ranzato, et al. CVPR 2011

Outline

- mathematical formulation of the model
- training
- learning acoustic features for speech recognition
- generation of natural images
- **recognition of facial expression under occlusion**
- conclusion

Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

anger



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

disgust



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

fear



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

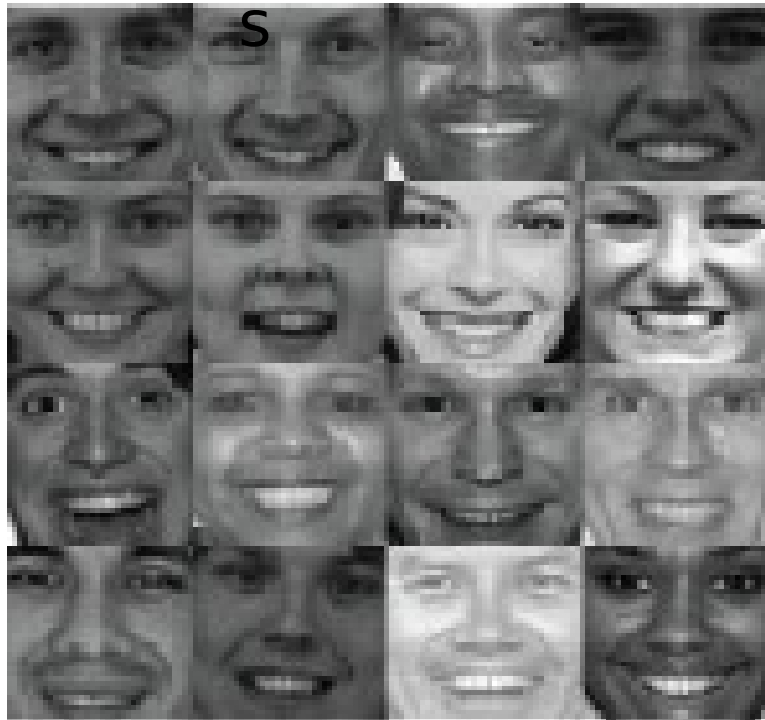
~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

happines



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

neutral



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

sadness



Facial Expression Recognition

Toronto Face Dataset (J. Susskind et al. 2010)

~ 100K unlabeled faces from different sources

~ 4K labeled images

Resolution: 48x48 pixels

7 facial expressions

surprise



Facial Expression Recognition

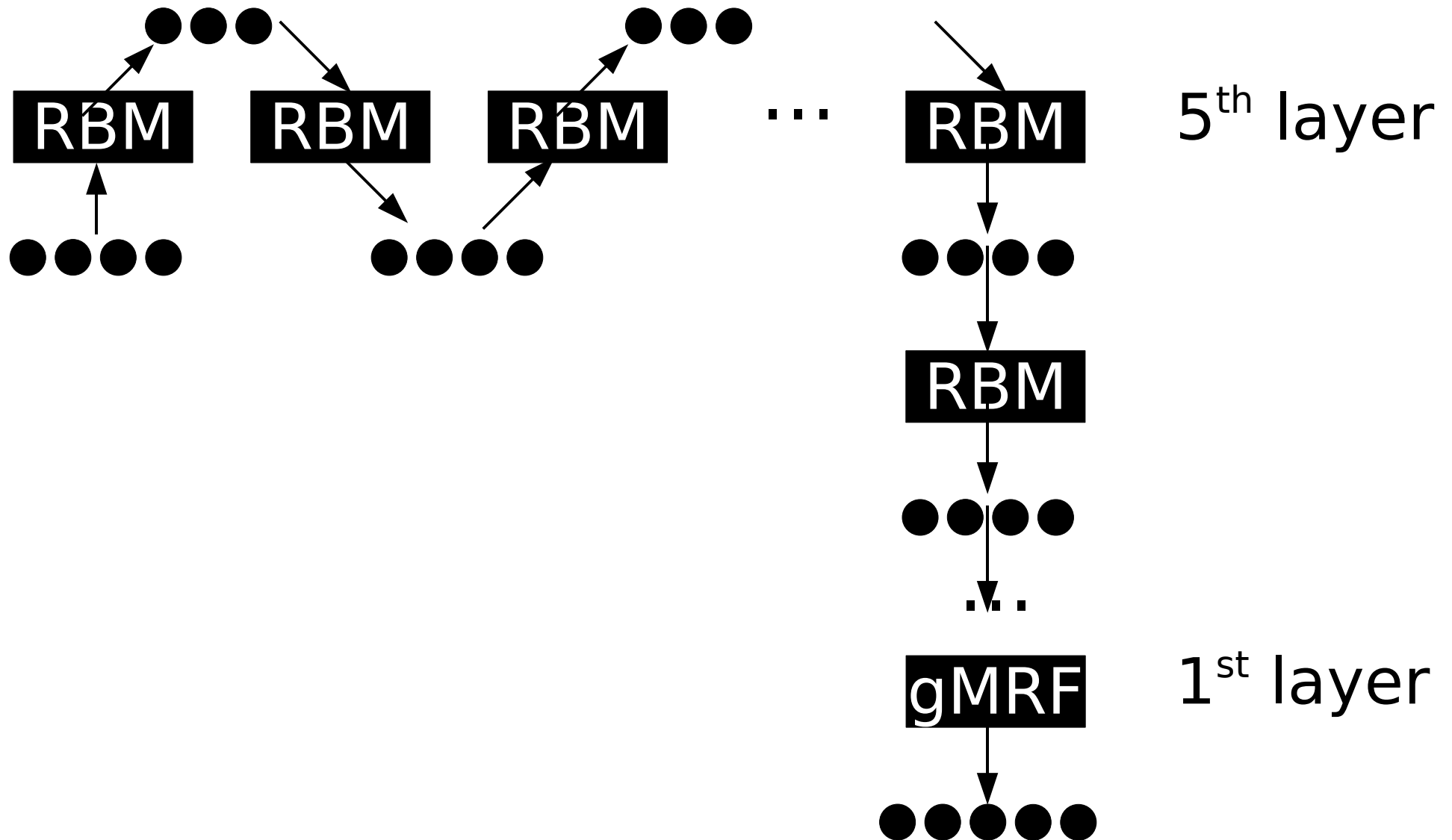
- 1st layer using local (not shared) connectivity
- layers above are fully connected
- 5 layers in total

- Result

- | | |
|--|-------|
| - Linear Classifier on raw pixels | 71.5% |
| - Gaussian RBF SVM on raw pixels | 76.2% |
| - Gabor + PCA + linear classifier
<i>Dailey et al. J. Cog. Science 2002</i> | 80.1% |
| - Sparse coding
<i>Wright et al. PAMI 2008</i> | 74.6% |
| - DEEP model (3 layers): | 82.5% |

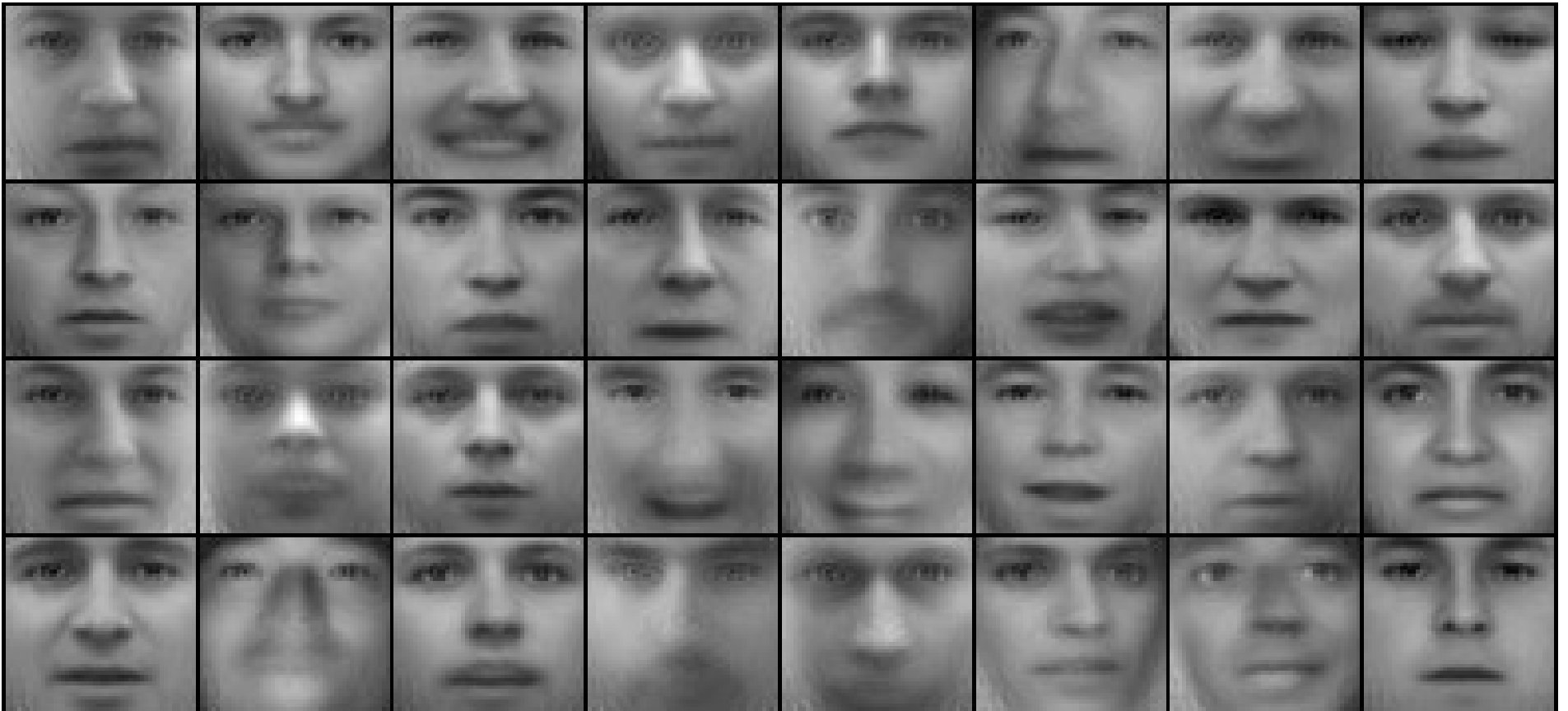
Facial Expression Recognition

- We can draw samples from the model (5th layer with 128 hiddens)



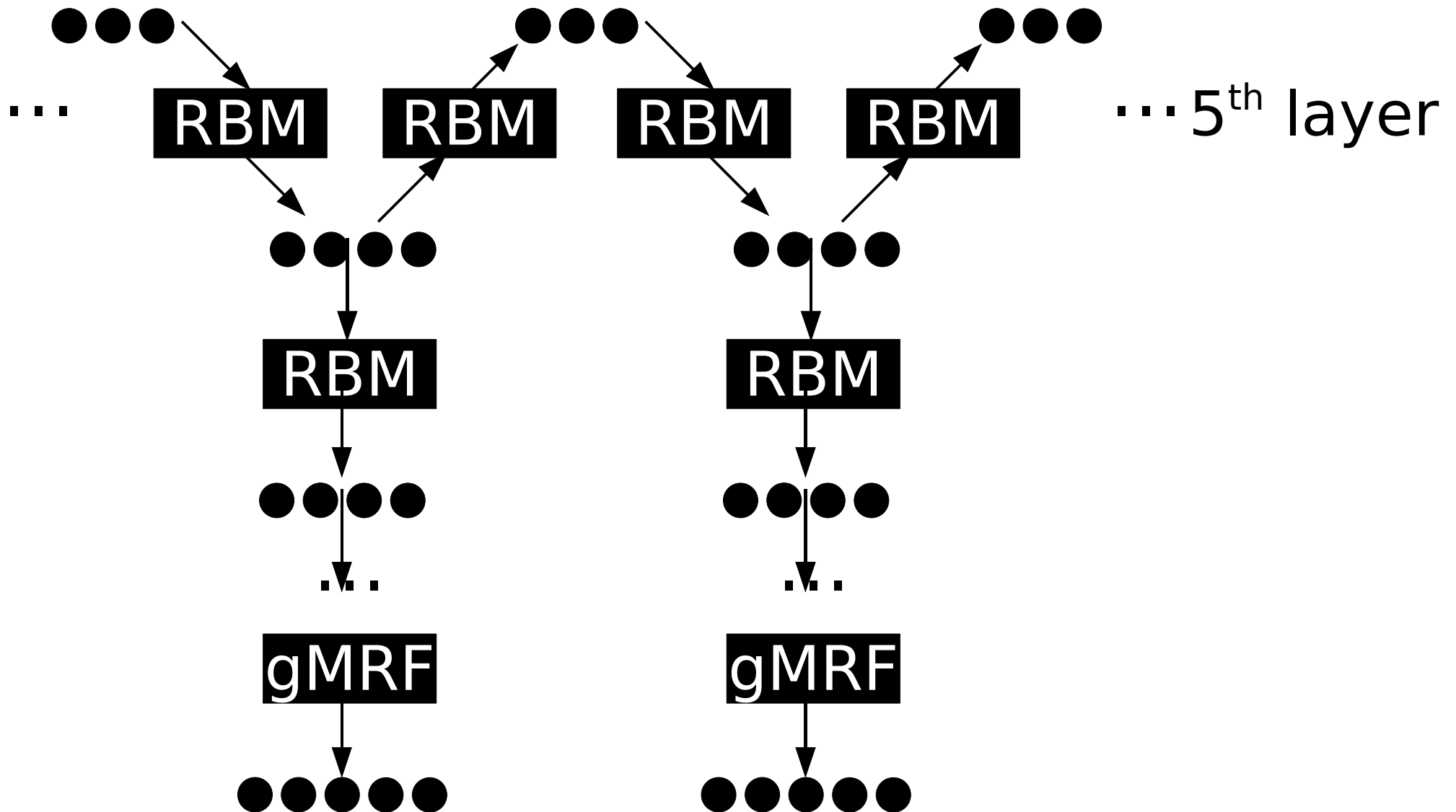
Facial Expression Recognition

- We can draw samples from the model (5th layer with 128 hiddens)



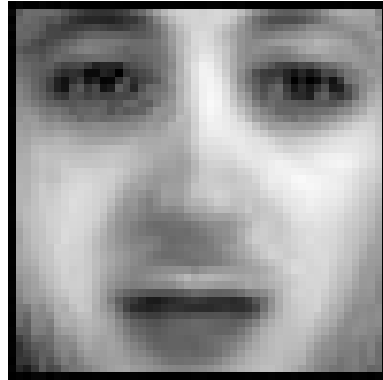
Facial Expression Recognition

- We can draw samples from the model (5th layer with 128 hiddens)



Facial Expression Recognition

- We can draw samples from the model (5th layer with 128 hiddens)



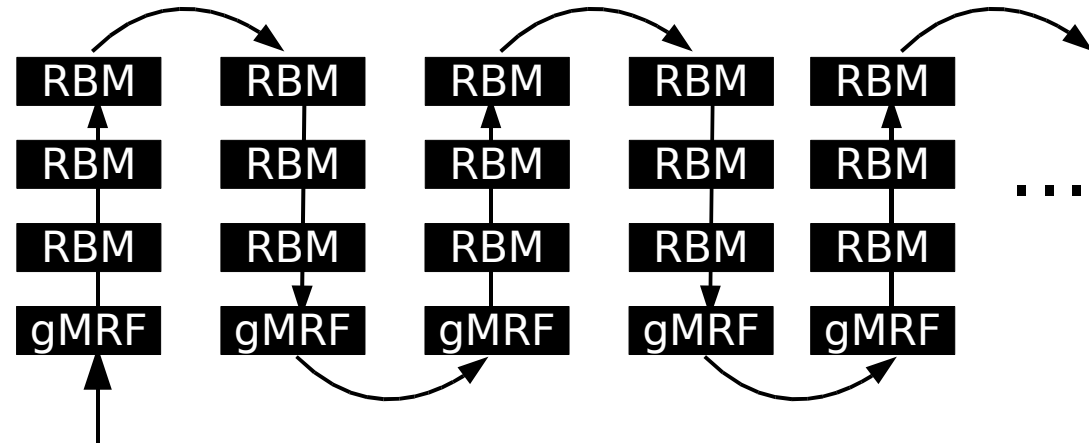
Facial Expression Recognition

- We can draw samples from the model (5th layer with 128 hiddens)



Facial Expression Recognition

- 7 synthetic occlusions
- use generative model to fill-in (conditional on the known pixels)



Facial Expression Recognition

originals



Type 1 occlusion: eyes



Restored images

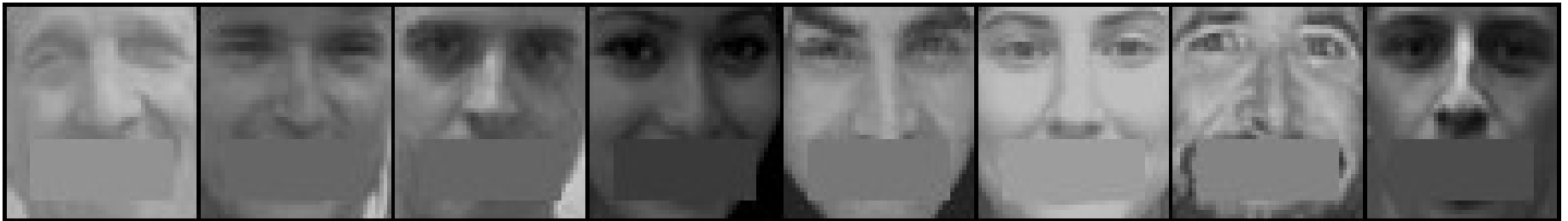


Facial Expression Recognition

originals



Type 2 occlusion: mouth



Restored images

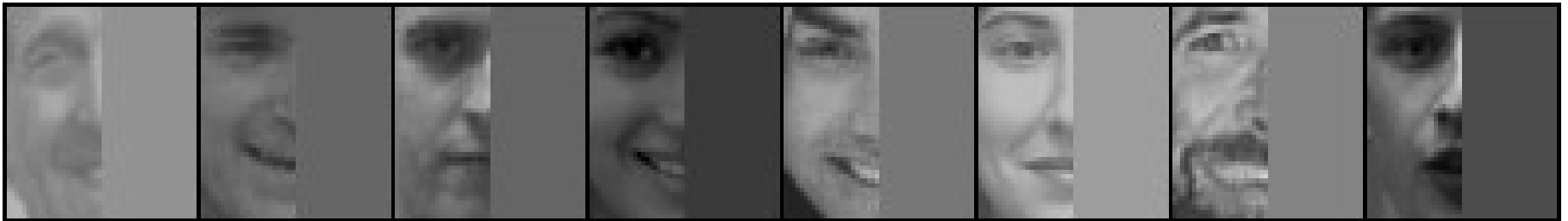


Facial Expression Recognition

originals



Type 3 occlusion: right half



Restored images



Facial Expression Recognition

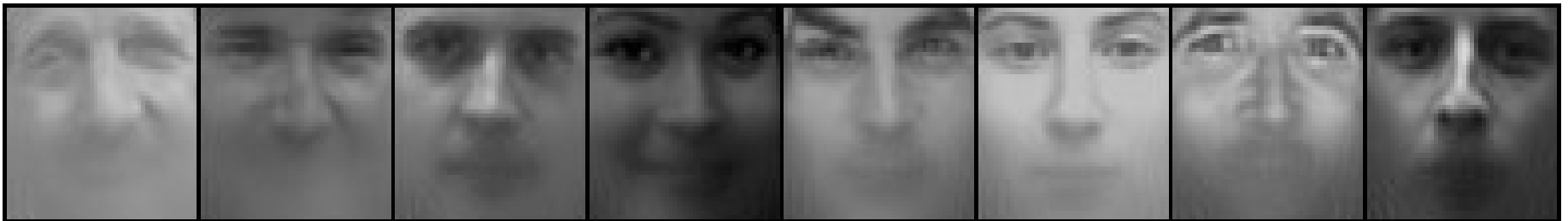
originals



Type 4 occlusion: bottom half



Restored images



Facial Expression Recognition

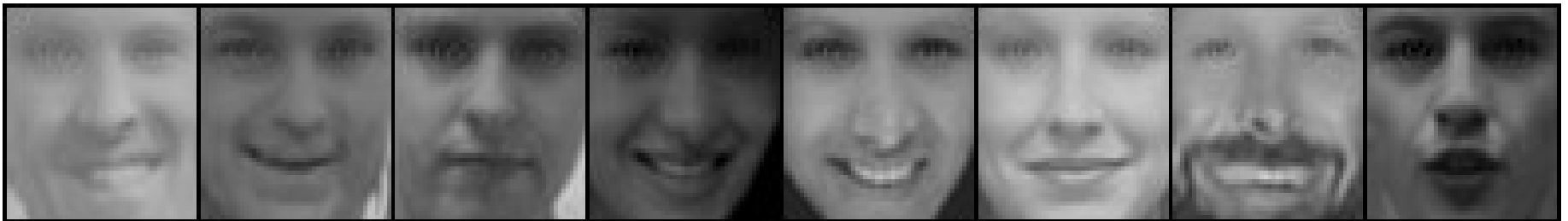
originals



Type 5 occlusion: top half



Restored images



Facial Expression Recognition

originals



Type 6 occlusion: nose



Restored images

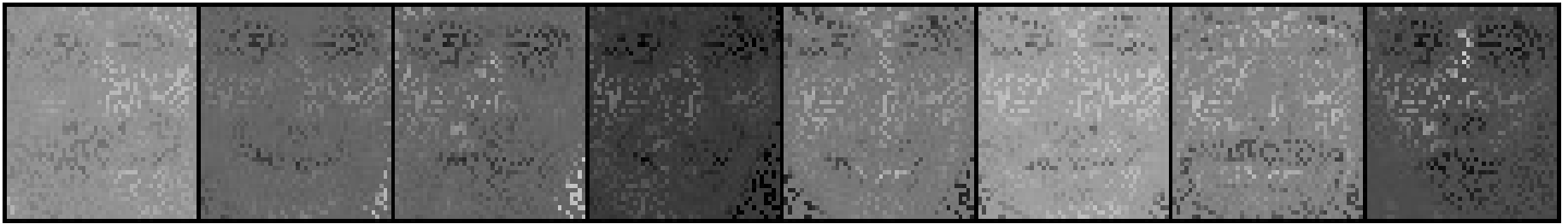


Facial Expression Recognition

originals



Type 7 occlusion: 70% of pixels at random



Restored images



Facial Expression Recognition

Original



Input



Facial Expression Recognition

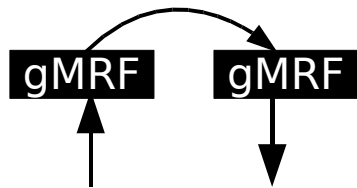
Original



Input



1st



Facial Expression Recognition

Original



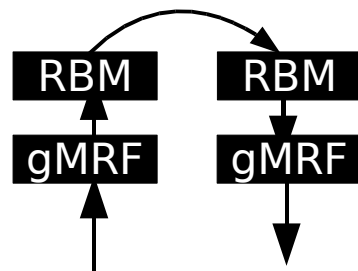
Input



1st



2nd



Facial Expression Recognition

Original



Input



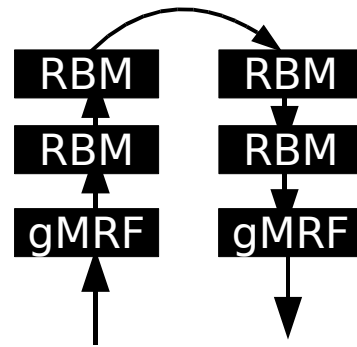
1st



2nd



3rd



Facial Expression Recognition

Original



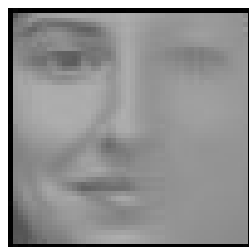
Input



1st



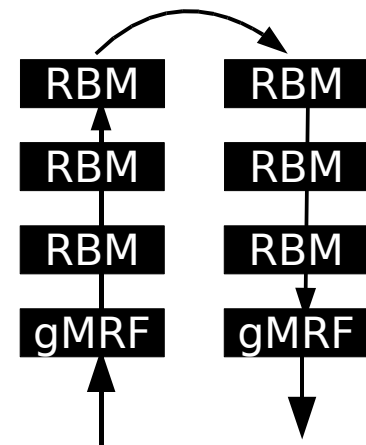
2nd



3rd



4th



Facial Expression Recognition

Original



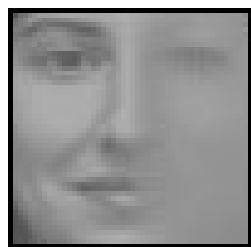
Input



1st



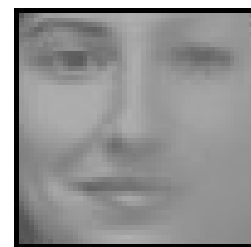
2nd



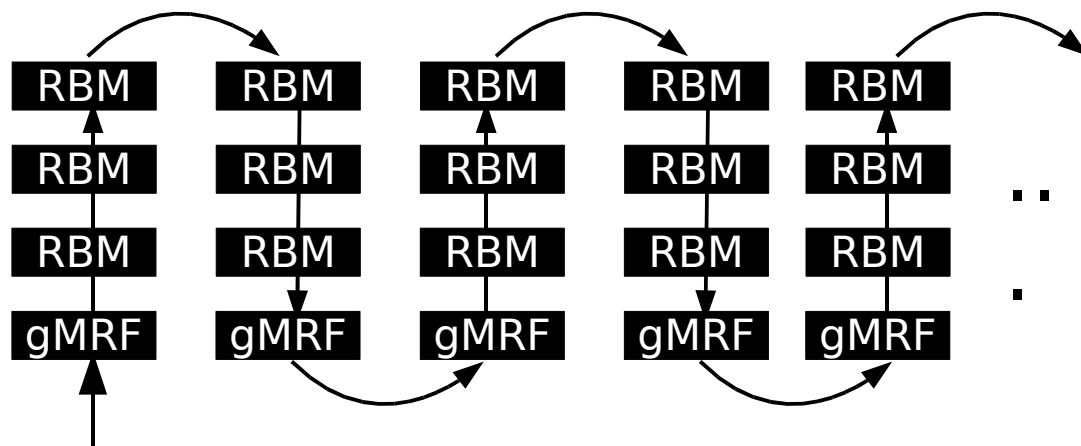
3rd



4th

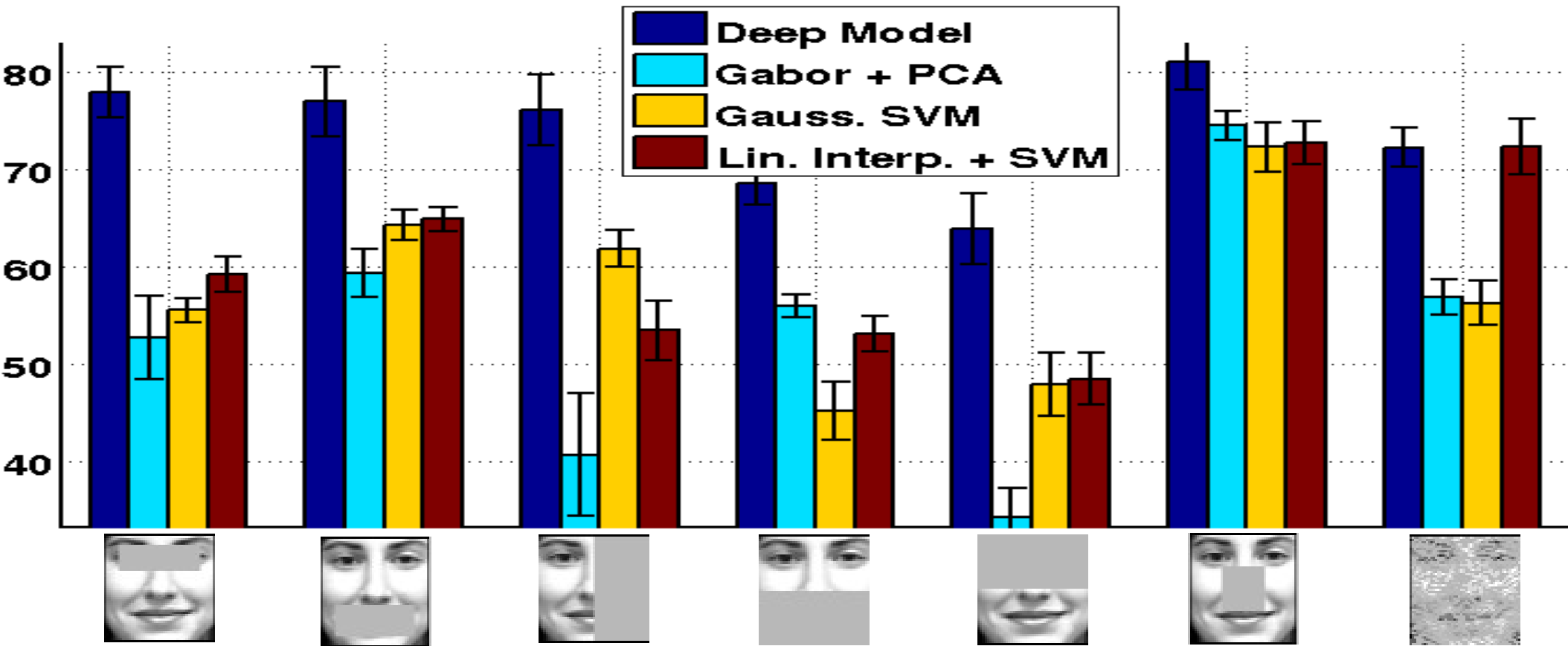


10 times



Facial Expression Recognition

CASE 1: original images for training, occluded for test



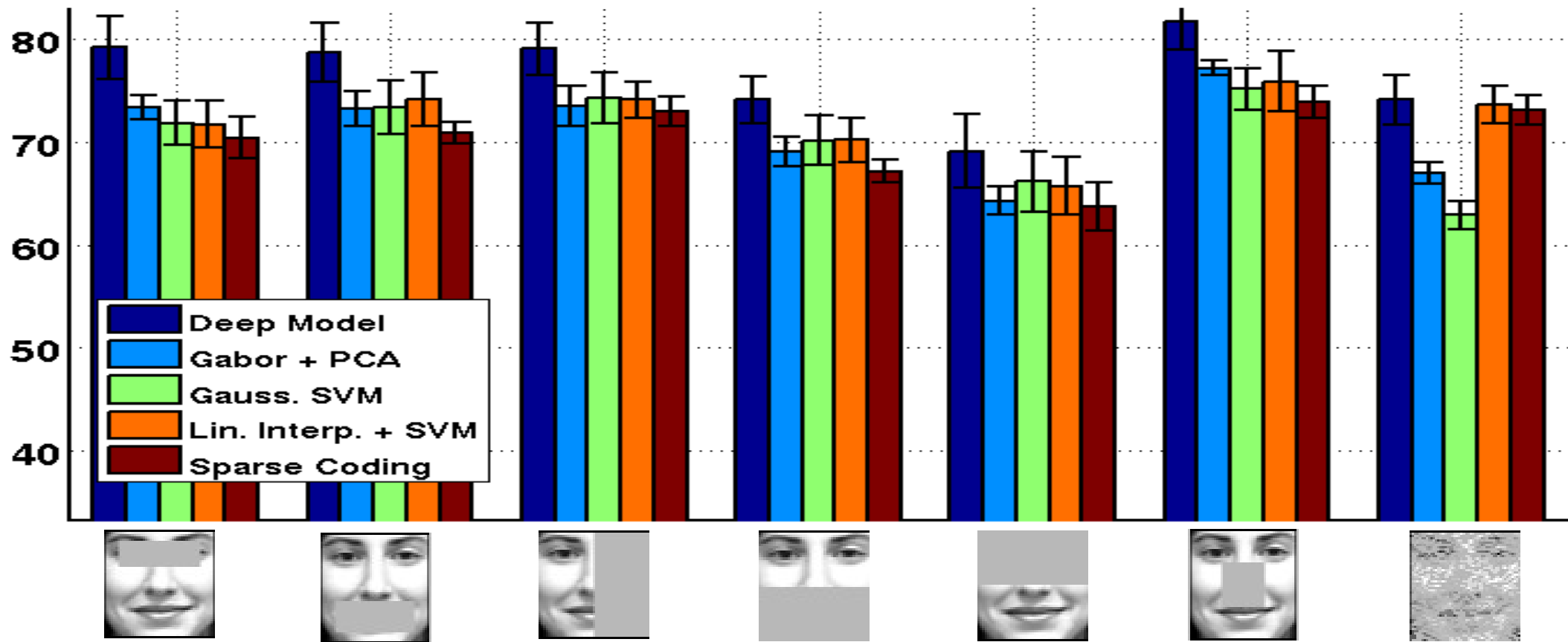
Dailey, et al. *J. Cog. Neuros.* 2003

Wright, et al. *PAMI* 2008

Ranzato, et al. *CVPR* 2011

Facial Expression Recognition

CASE 2: occluded images for both training and test



Dailey, et al. J. Cog. Neuros. 2003

Wright, et al. PAMI 2008

Ranzato, et al. CVPR 2011

Outline

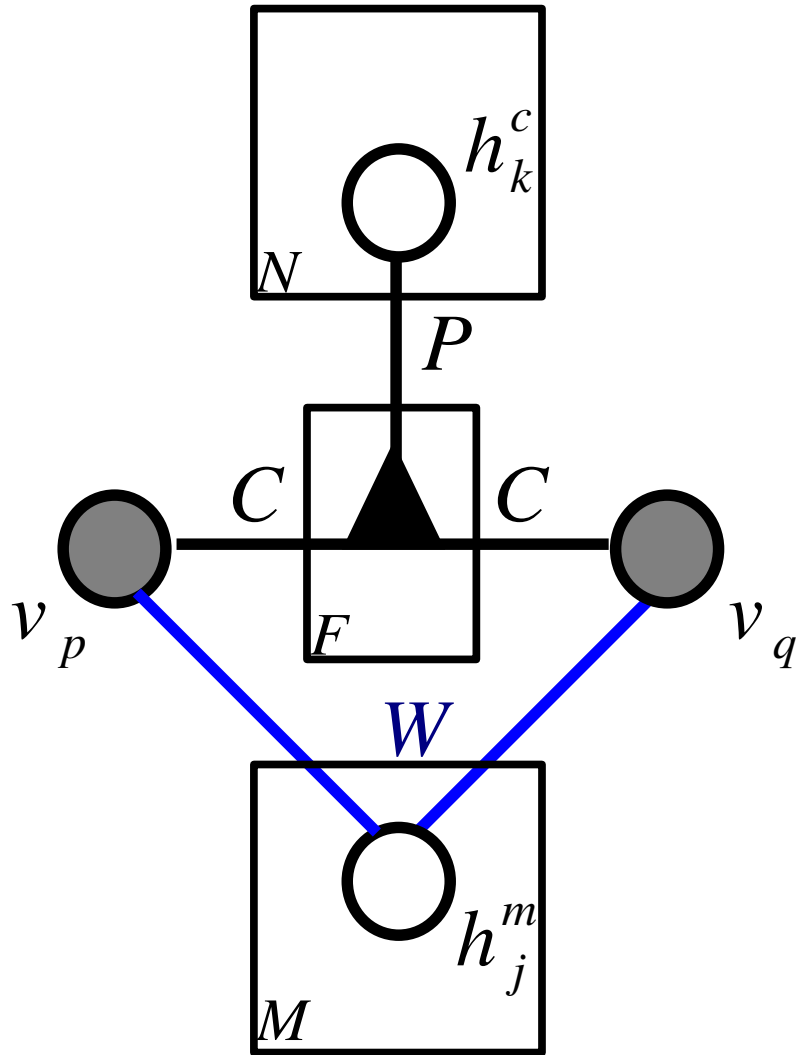
- mathematical formulation of the model
- training
- learning acoustic features for speech recognition
- generation of natural images
- recognition of facial expression under occlusion
- **conclusion**

Summary

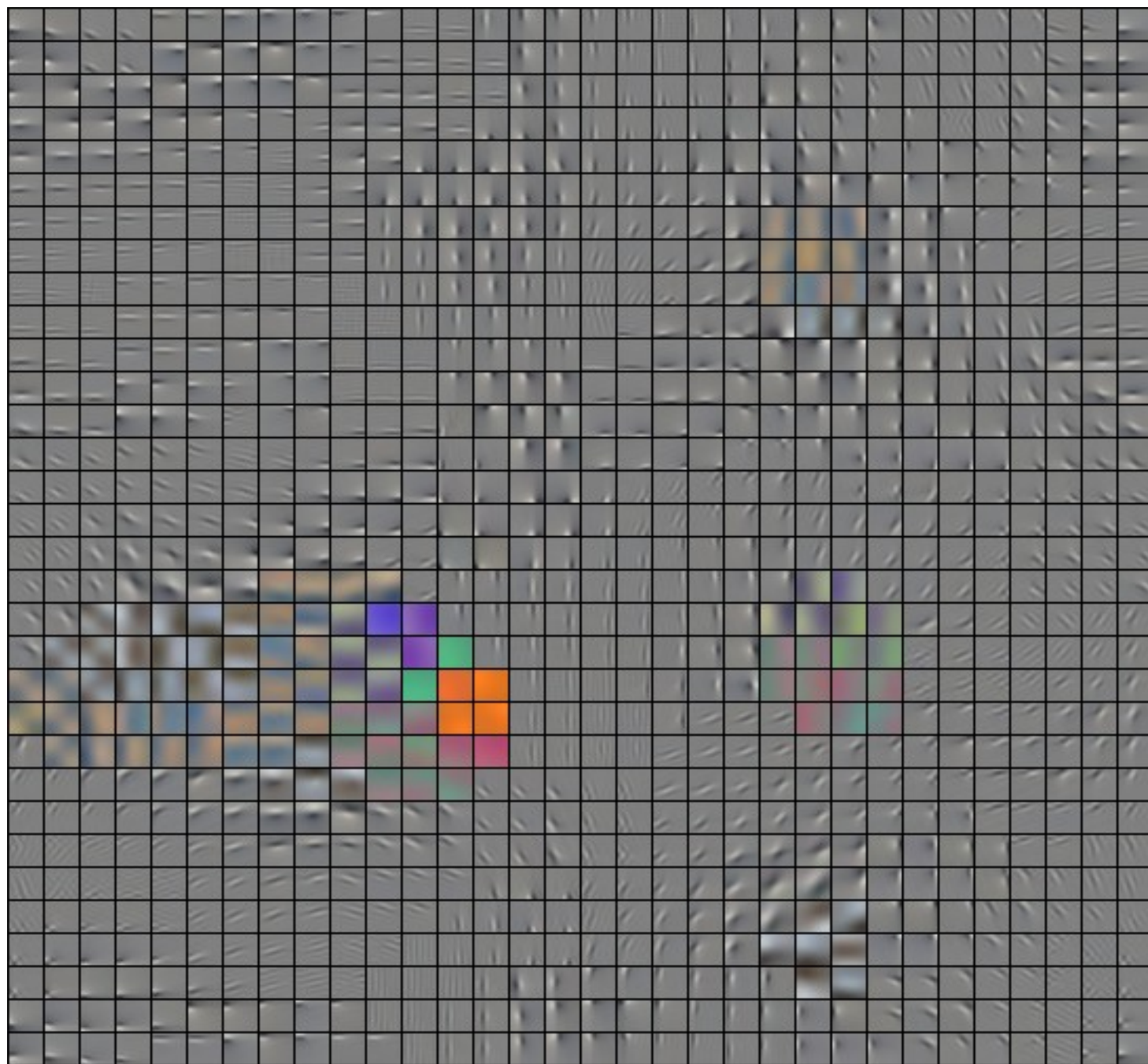
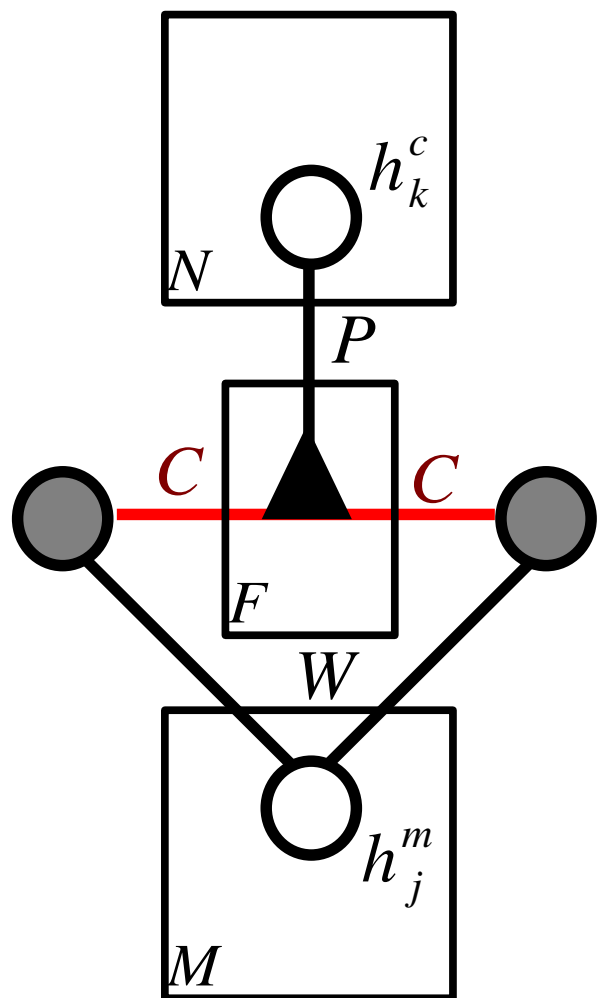
- Unsupervised Learning: regularize & learn representations
- Deep Generative Model
 - 1st layer: gated MRF
 - Higher layers: binary RBM's
 - fast inference
- Realistic generation: natural images
- Applications:
 - facial expression recognition robust to occlusion
 - speech recognition

THANK YOU

Learned Filters



Learned Filters



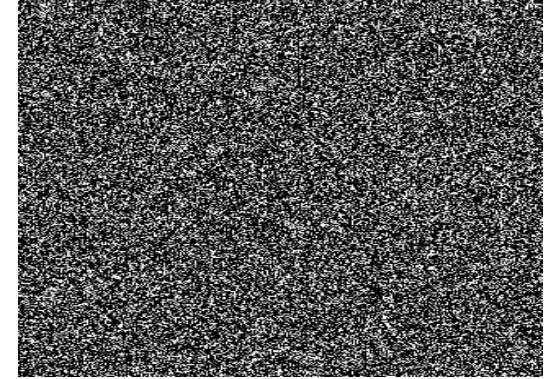
Using -Energy to Score Images

less likely

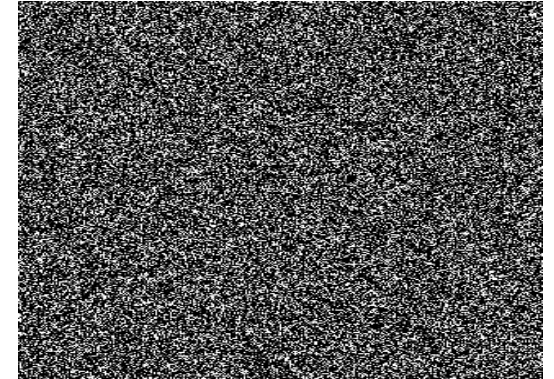
test images



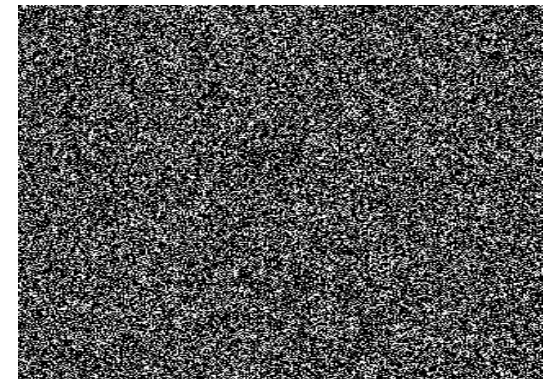
>



>

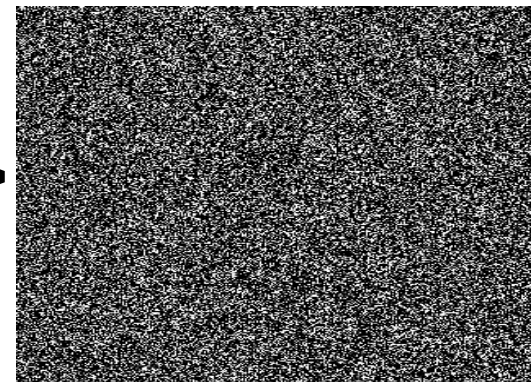
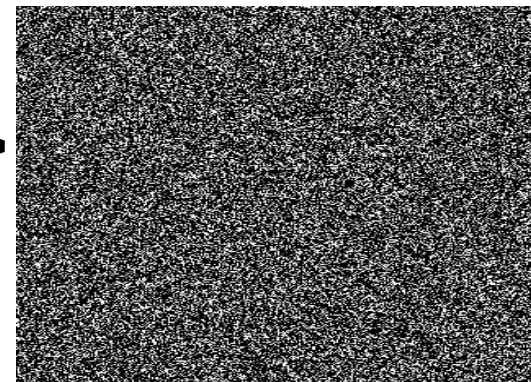
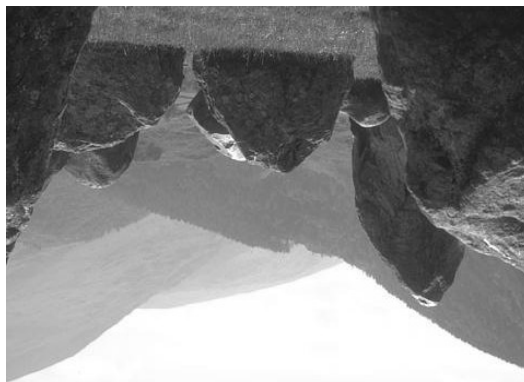
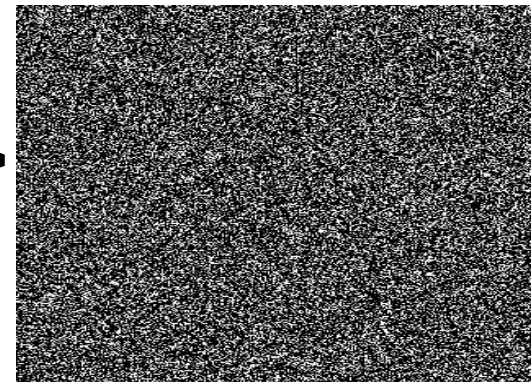


>



Using Energy to Score Images

Upside-down images



Using Energy to Score Images

Average of those images for which
difference of energy is higher

