

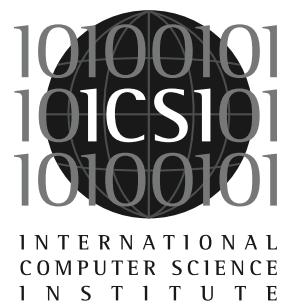
Hierarchical Dirichlet Process Hidden Markov Trees for Multiscale Image Analysis

Jyri Kivinen

International Computer Science Institute, Berkeley, USA

Helsinki University of Technology, Espoo, Finland

The University of Edinburgh, UK

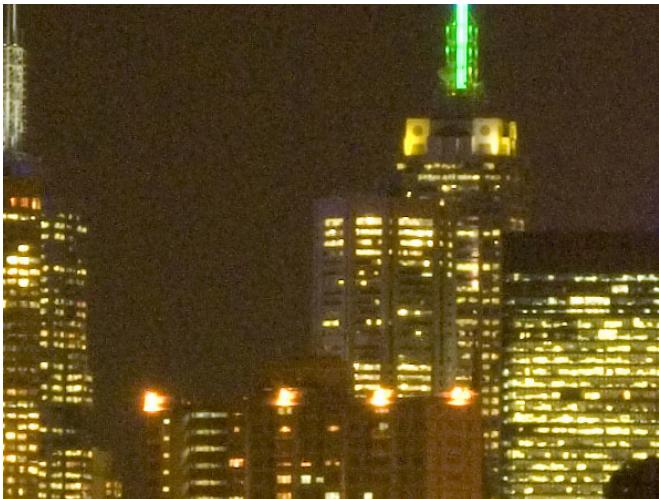


Joint work with

Erik Sudderth, Brown University
Michael Jordan, UC Berkeley



Low-level Image Analysis



Noise Removal



Deblurring



Inpainting & Restoration

Goals:

- Accurately model the statistics of *natural images*
- Exploit the availability of large digital *image collections*
 - Suggests use of data-driven, *nonparametric* models

Natural Scene Categorization



Coast

Forest

Open Country

Street

Tall Building

Goals:

- Visually *recognize* natural scene categories
- Accurately model the statistics of *natural scenes*
- Learn *global* statistical scene models

Outline

Multiscale Models for Natural Images

- Nonparametric Hidden Markov Trees (HDP-HMTs)
- Learning with Monte Carlo methods
- Truncated representations for efficient learning from large datasets

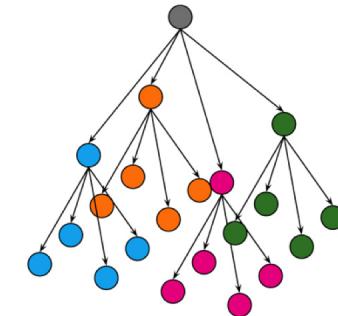
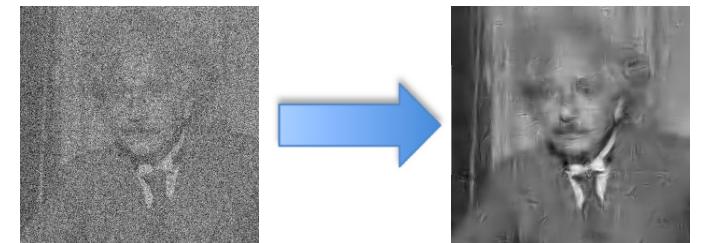


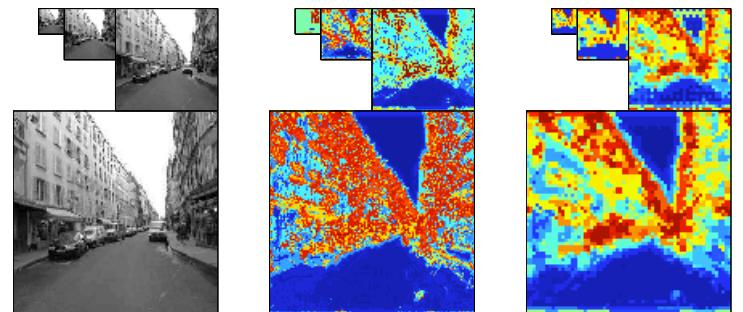
Image Denoising

- Transfer natural image statistics for making robust predictions

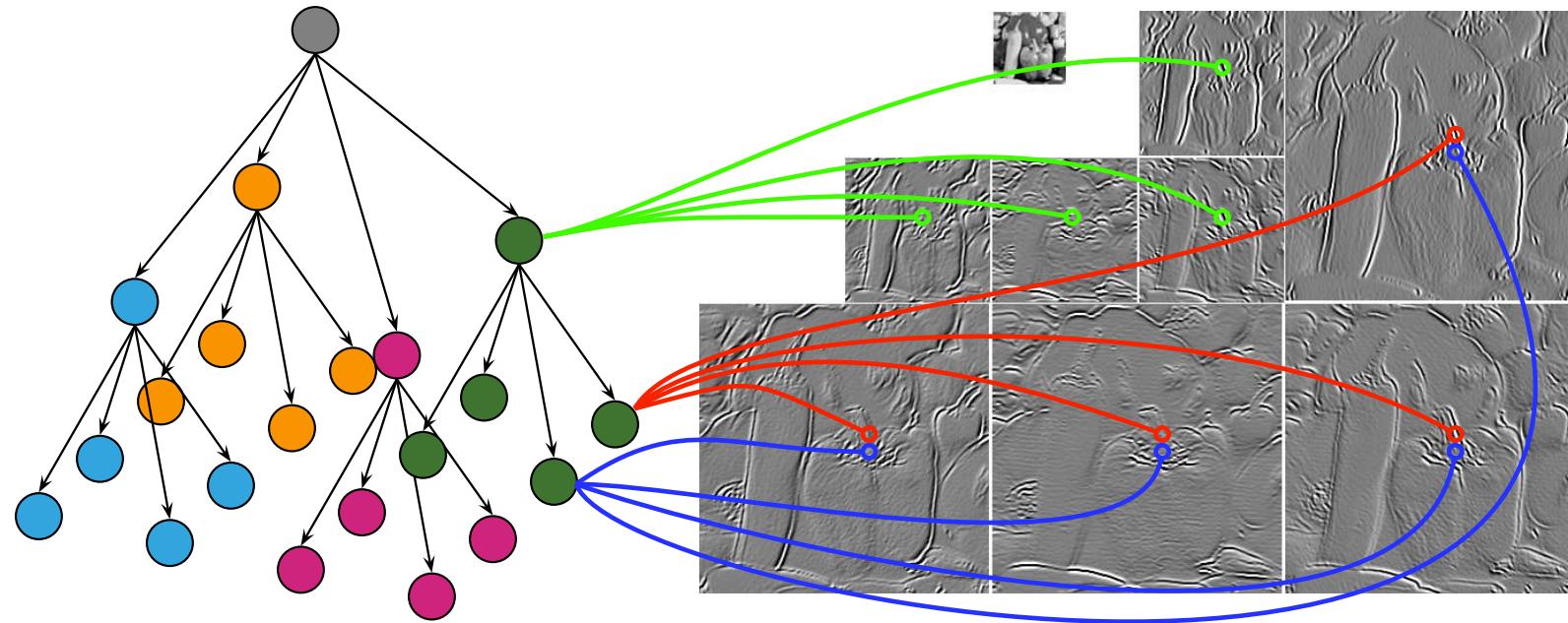


Natural Scene Analysis

- Global, data-driven scene models via HDP-HMT
- Fast categorization via Belief Propagation methods



Multiscale Models for Natural Images

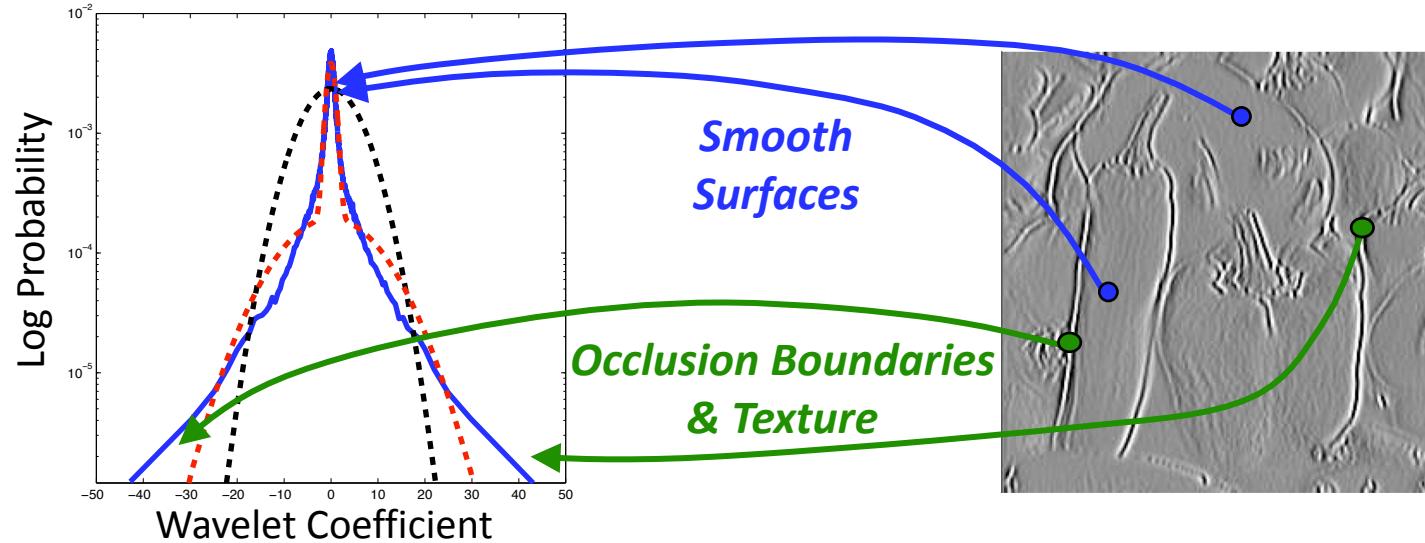


Goal: Accurately model the statistics of natural images

Approach:

- Capture multiscale dependencies using a *tree* of latent variables
- Automatically adapt to data complexity via nonparametric,
Dirichlet process priors

Mixture Models for Heavy-Tailed Wavelet Marginals

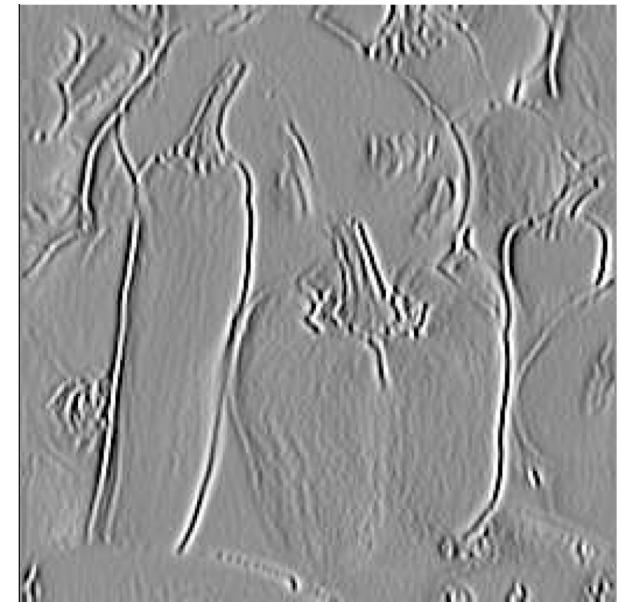
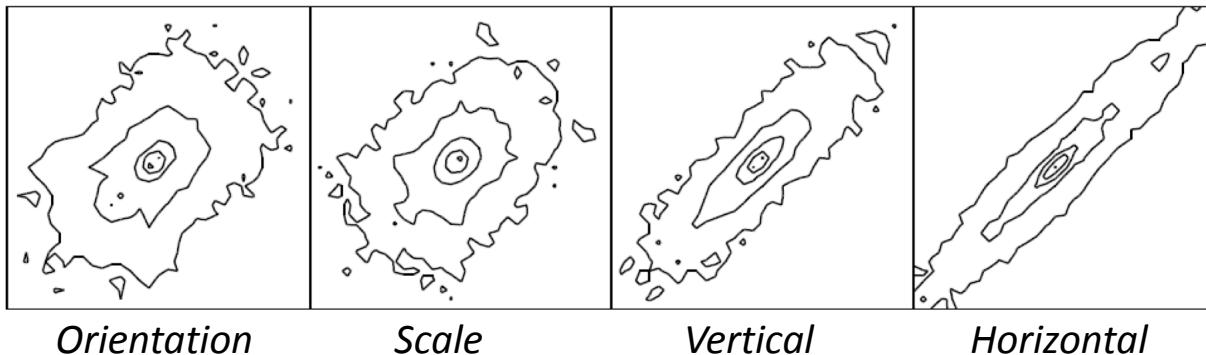


- Extreme coefficient values resultant from edges and texture occur more frequently than with a Gaussian
- Gaussian scale mixtures provide good matches for the highly kurtotic, heavy tailed distributions
$$x_{ti} = v_{ti} u_{ti}; v_{ti} \geq 0, u_{ti} \sim \mathcal{N}(0, \Lambda)$$
- Discrete mixtures easier to work with, reasonable denoising results even with binary mixtures:

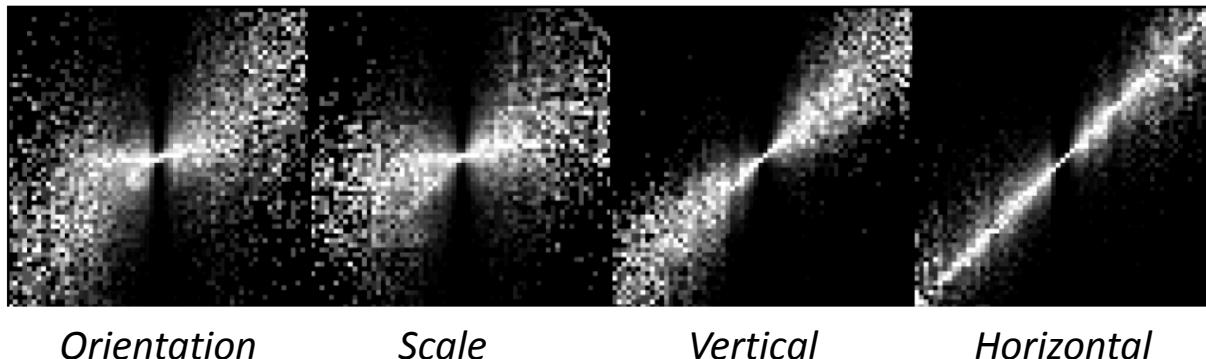
$$x_{ti} \sim \pi \mathcal{N}(0, \Lambda_0) + (1 - \pi) \mathcal{N}(0, \Lambda_1)$$

Joint Statistics of Wavelet coefficients

Pairwise Joint Histograms:



Pairwise Conditional Histograms:

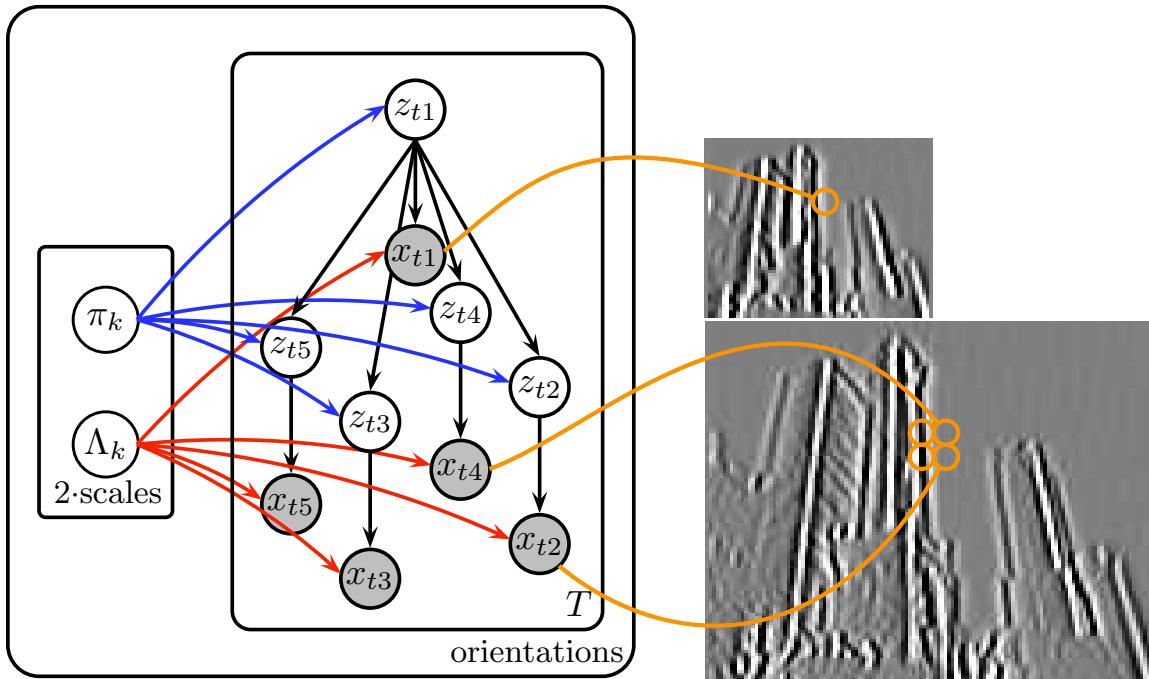


Large magnitude wavelet coefficients...

- *Persist* across multiple scales
- *Cluster* at adjacent spatial locations

Binary Hidden Markov Trees

Crouse, Nowak, & Baraniuk, 1998



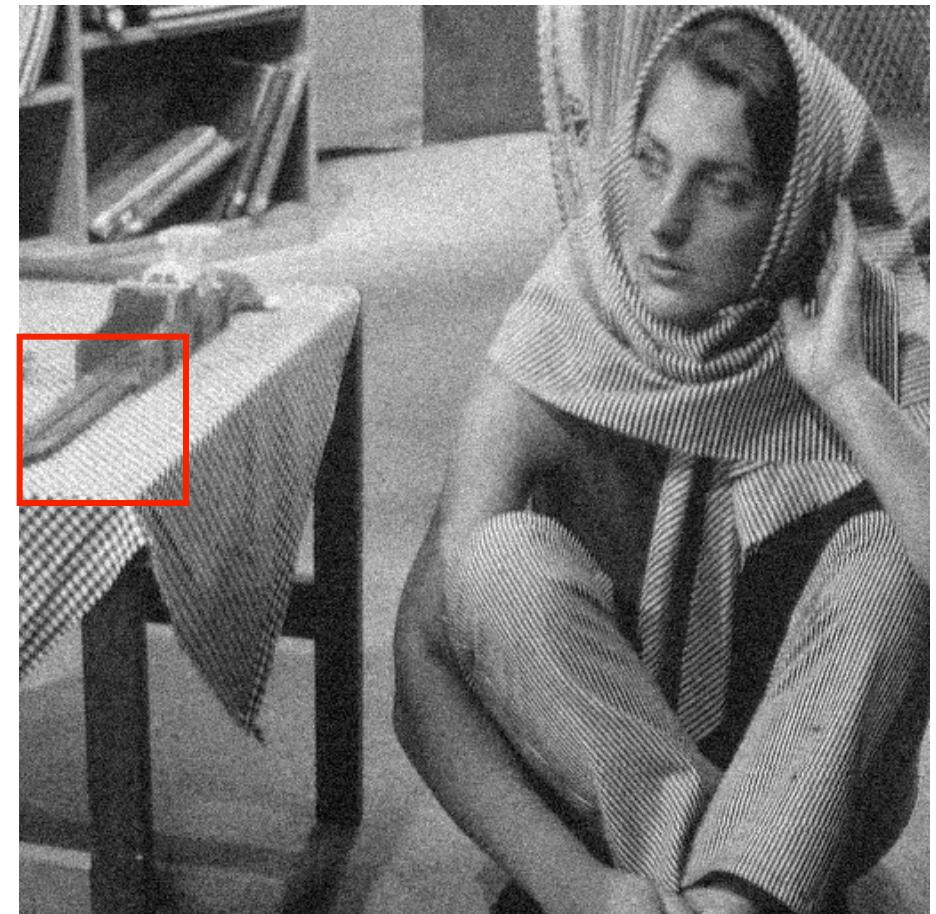
- $\pi_k \rightarrow$ state *transition* distributions
- $z_{ti} \sim \pi_{z_{Pa(ti)}}$
- $\Lambda_k \rightarrow$ state-specific *emission* covariances
- $x_{ti} \sim \mathcal{N}(0, \Lambda_{z_{ti}})$
- $z_{ti} \rightarrow$ hidden *state* or cluster assignment
- $z_{ti} \in \{0, 1\}$
- $x_{ti} \rightarrow$ *observed* wavelet coefficient

- Wavelet coefficients marginally distributed as mixtures of two Gaussians
- Markov dependencies between hidden states capture persistence of image contours across locations and scales
- Models each scale and orientation independently

Validation : Image Denoising

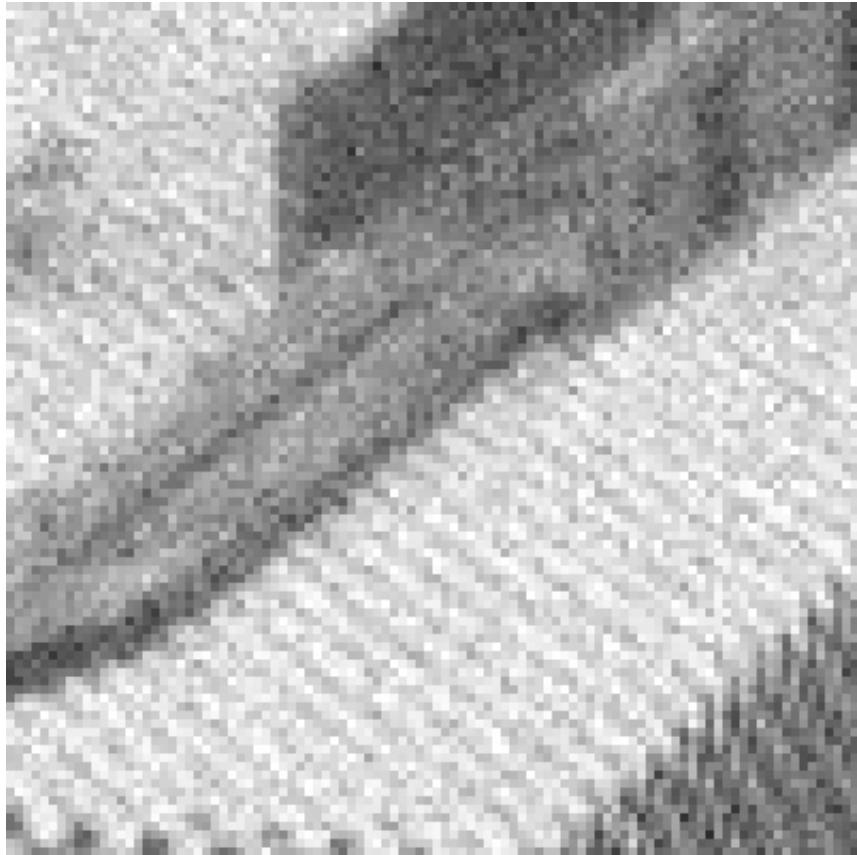


Original

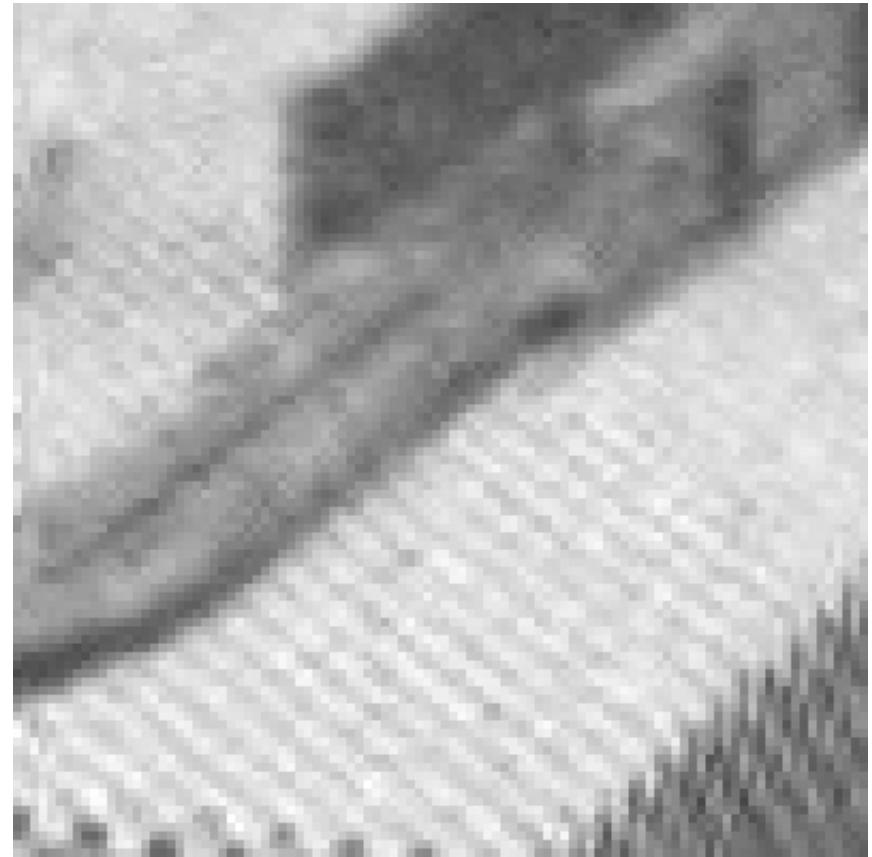


Corrupted by Additive
White Gaussian Noise
(PSNR = 24.61 dB)

Denoising with Binary HMTs



Noisy Input



Denoised (EM algorithm)

- Is two states per scale sufficient? How many is enough?
- Should states be shared the same way for all images, or for all wavelet decompositions?

Dirichlet Process Mixtures

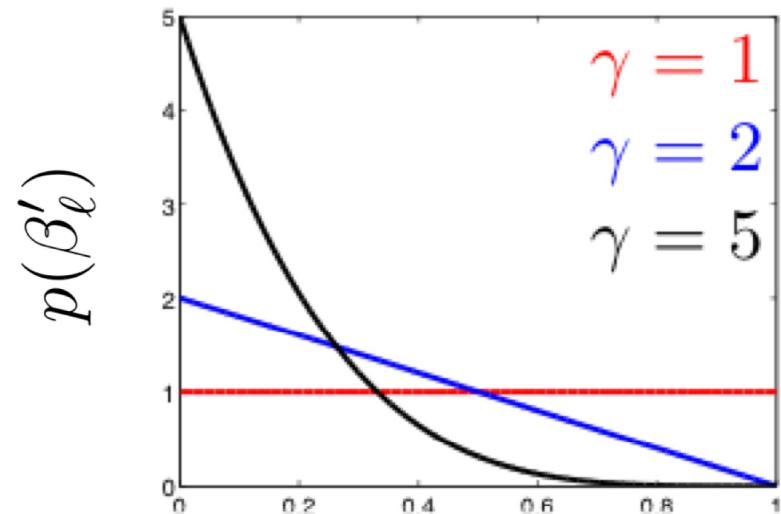
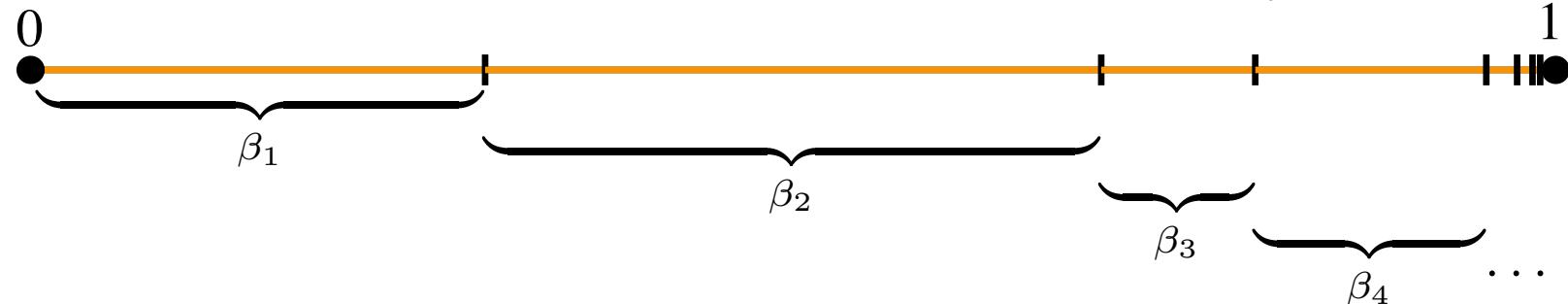
$$p(x_{ti} | \beta, \Lambda_1, \Lambda_2, \dots) = \sum_{k=1}^{\infty} \beta_k \mathcal{N}(x_{ti}; 0, \Lambda_k)$$

Stick-breaking prior for mixture weights controls complexity:

$$\beta_k = \beta'_k \prod_{\ell=1}^{k-1} (1 - \beta'_{\ell})$$

$$\beta'_\ell \sim \text{Beta}(1, \gamma)$$

$\gamma \rightarrow$ Concentration parameter

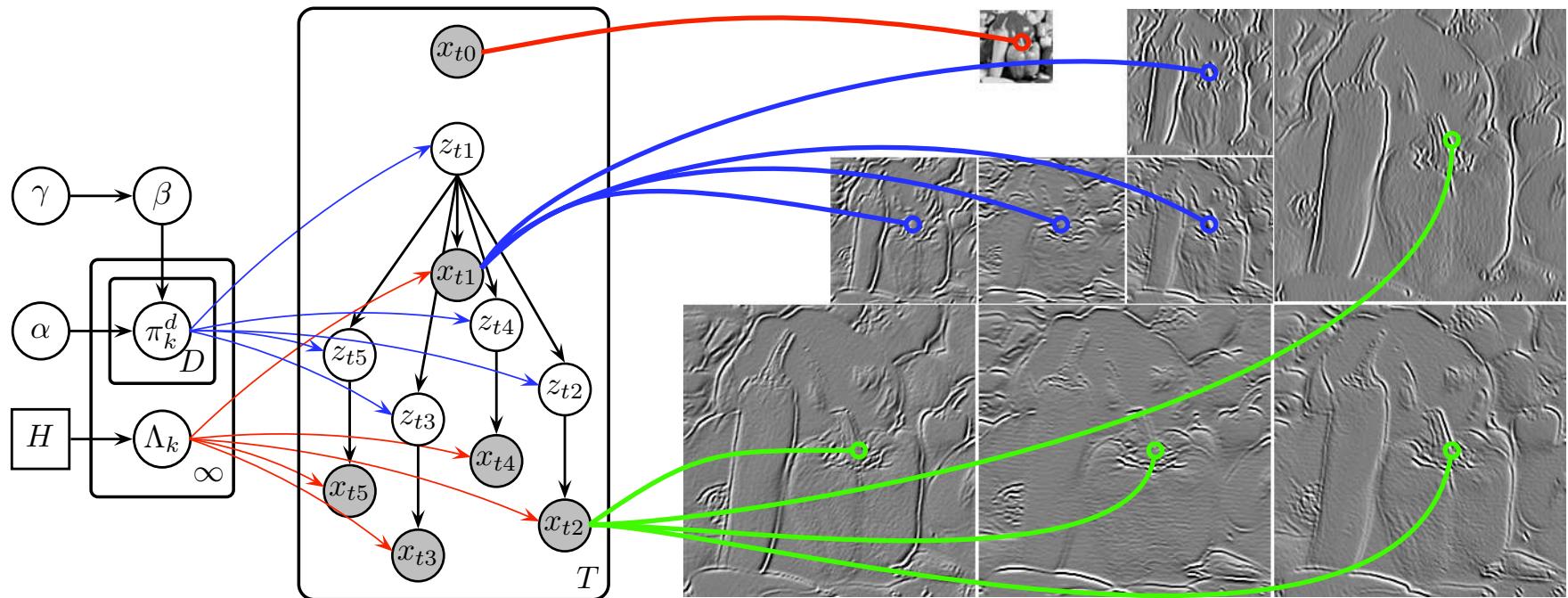


Why the Dirichlet Process ?

$$p(x) = \sum_{k=1}^{\infty} \beta_k f(x \mid \Lambda_k) \quad \begin{aligned} \beta &\sim \text{Stick}(\gamma) \\ \Lambda_k &\sim H \end{aligned}$$

- Basis for *nonparametric* models whose complexity grows as additional data is observed
 - Makes simple predictions given few observations
 - Low-weight clusters capture details of very large datasets
- Attractive *asymptotic guarantees*
 - Posterior consistency of DP mixture density estimators
 - Convergence to finite mixture parameters of any order
- Leads to simple, effective *computational methods*
 - Growing literature on Monte Carlo and variational methods
 - Integrated, efficient handling of models of varying orders

Hierarchical Dirichlet Process Hidden Markov Trees



$z_{ti} \rightarrow$ indexes *infinite* set of hidden states
 $z_{ti} \in \{1, 2, 3, \dots\}$

$\pi_k \rightarrow$ infinite set of state *transition* distributions
 $z_{ti} \sim \pi_{z_{Pa(ti)}}^{d_{ti}}$

$x_{ti} \rightarrow$ observed *vector* of wavelet coefficients

$\Lambda_k \rightarrow$ state-specific *emission* covariances
 $x_{ti} \sim \mathcal{N}(0, \Lambda_{z_{ti}})$
 $\Lambda_k \sim H$

Why a Hierarchical DP ?

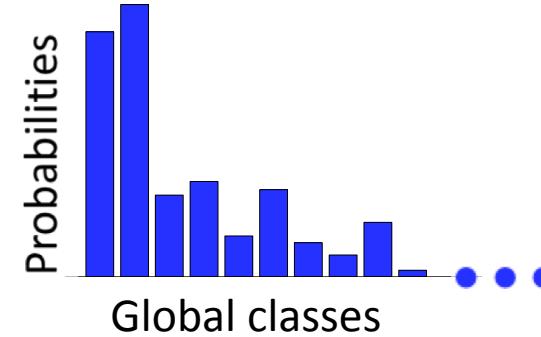
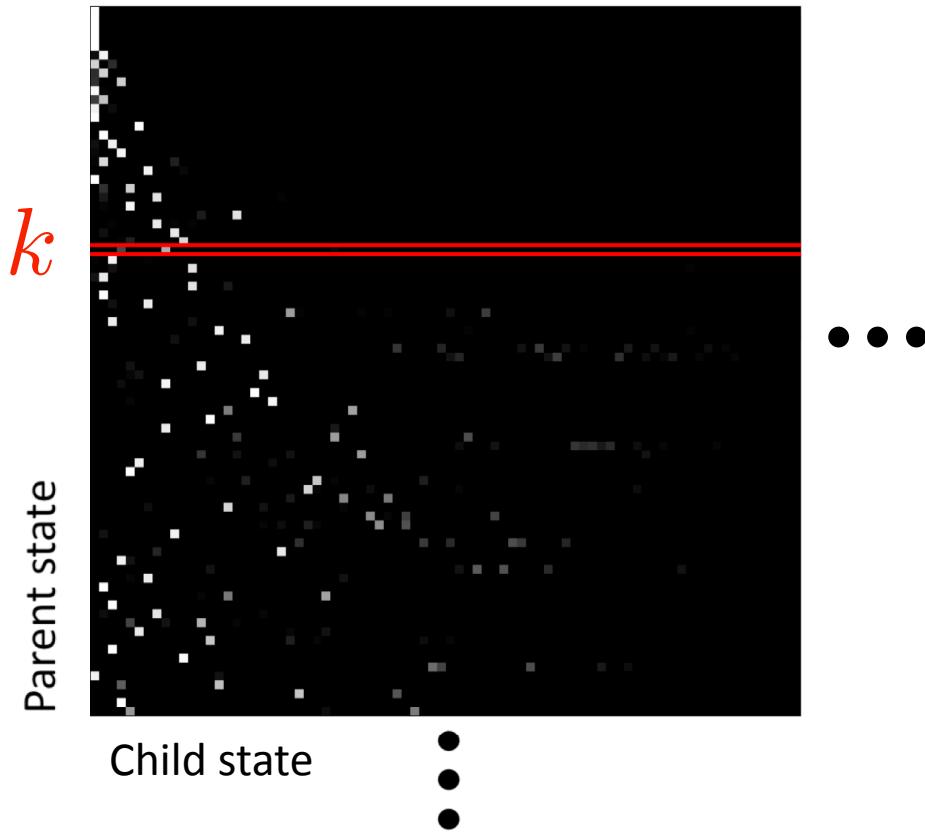
(Teh et. al. 2004)

- Hierarchical DP prior allows us to learn a potentially infinite set of *appearance patterns* from natural images
- Hierarchical coupling ensures, with high probability, that a common set of *child* states are reachable from each *parent*

$$\pi_k^{d_{ti}}(\ell) = \Pr [z_{ti} = \ell | z_{\text{Pa}(ti)}]$$

$$\beta \sim \text{Stick}(\gamma)$$

Average state frequencies



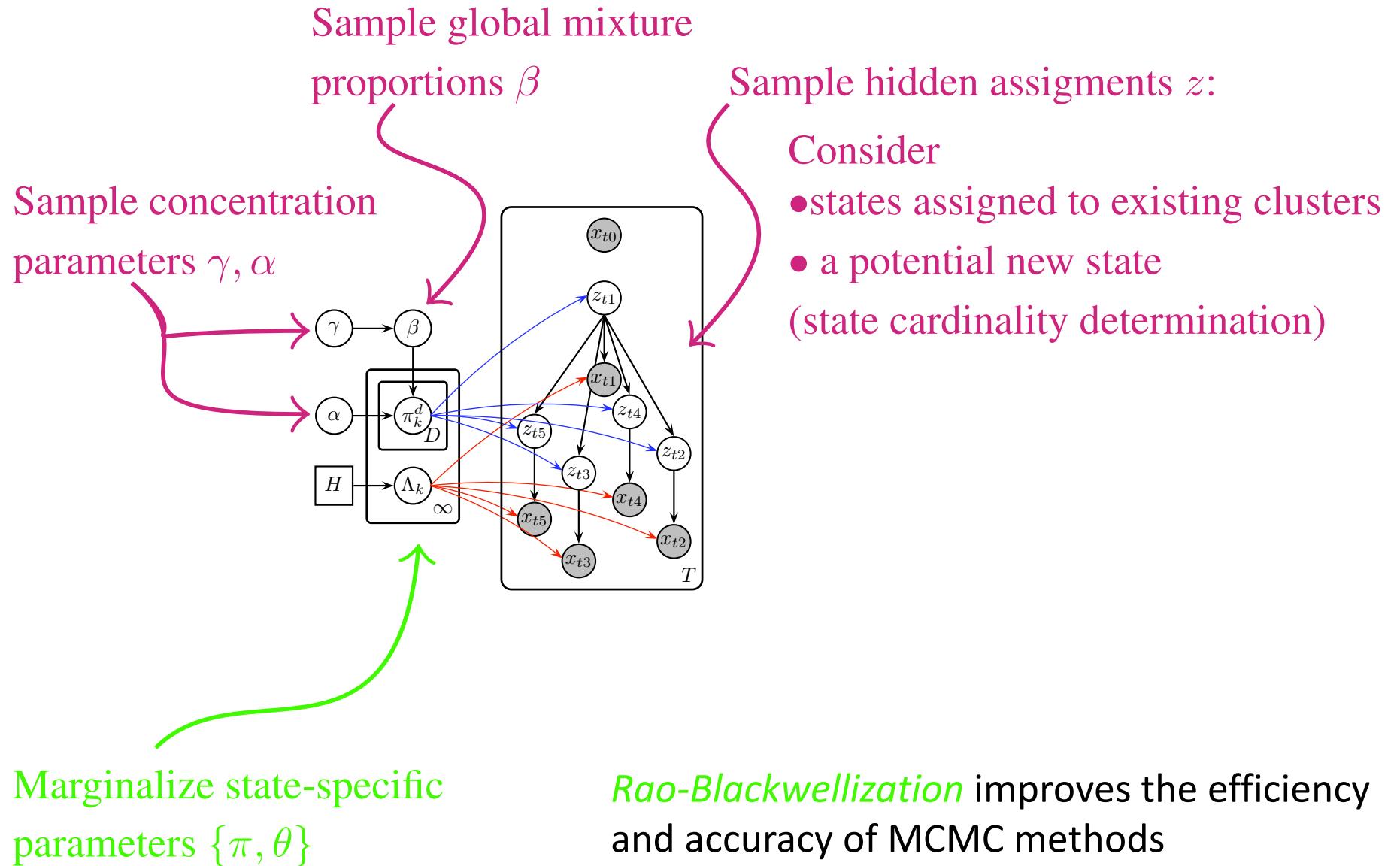
$$\pi_k^d \sim \text{DP}(\alpha, \beta)$$

Transition distributions

$$\mathbb{E} [\pi_k^d] = \beta$$

$\alpha \rightarrow$ *Sparsity & variability of transition distributions*

Learning HDP-HMT Models with a Collapsed Gibbs Sampler



Learning with a Truncated Gibbs Sampler

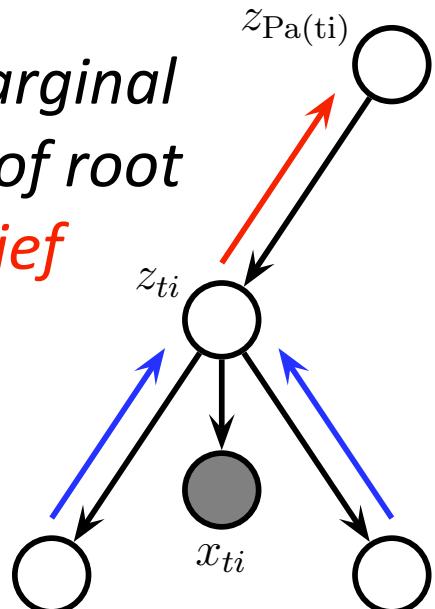
- Weak limit approximation uses high probability *upper bounds* on the number of states underlying a finite dataset:

$$\beta = (\beta_1, \dots, \beta_K) \sim \text{Dir}(\gamma/K, \dots, \gamma/K)$$

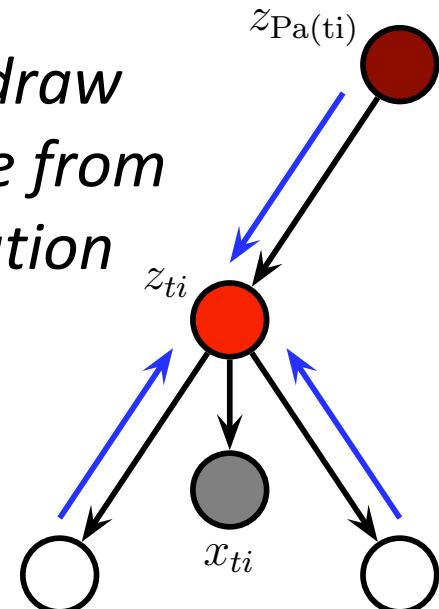
$$\pi_t = (\pi_{t1}, \dots, \pi_{tK}) \sim \text{Dir}(\alpha\beta_1, \dots, \alpha\beta_K)$$

- Predictions from *truncated* model converge in distribution to HDP as $K \rightarrow \infty$, and allow efficient *blocked sampling*:

Compute marginal distribution of root state via belief propagation

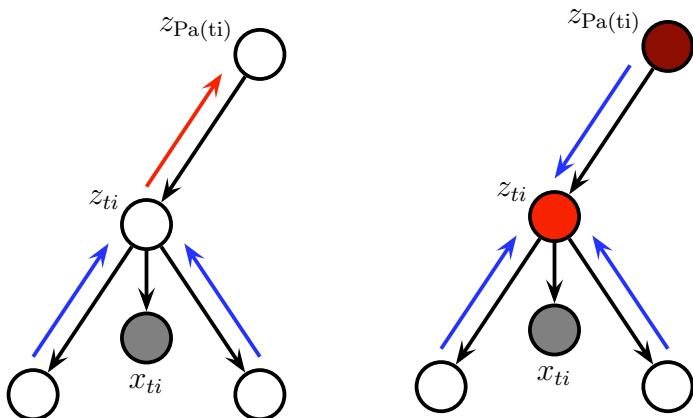


Recursively draw exact sample from joint distribution of all states



Learning with a Truncated Gibbs Sampler

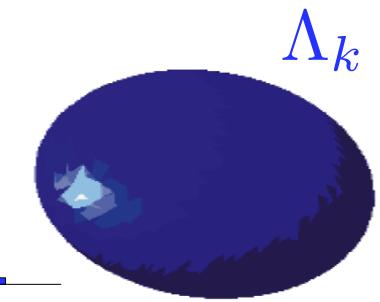
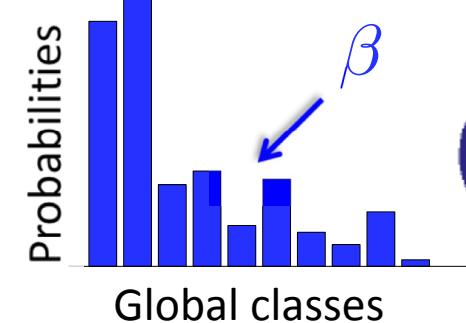
Sample hidden state assignments *jointly* using belief propagation:



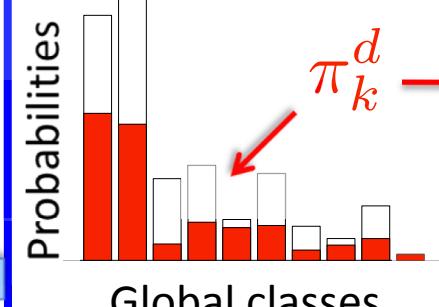
Compute messages from leaves to root

Sample assignments from root to leaves

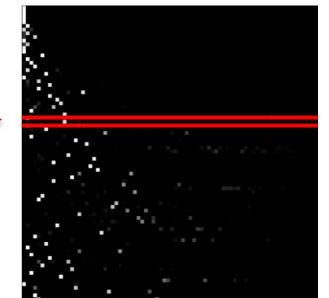
Sample parameters:



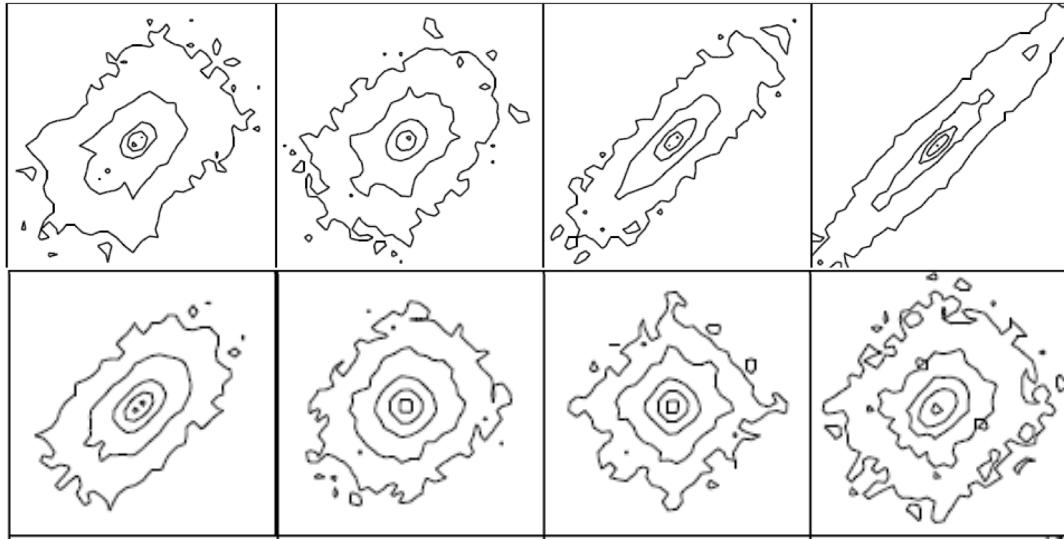
Feature distributions



Transition probabilities

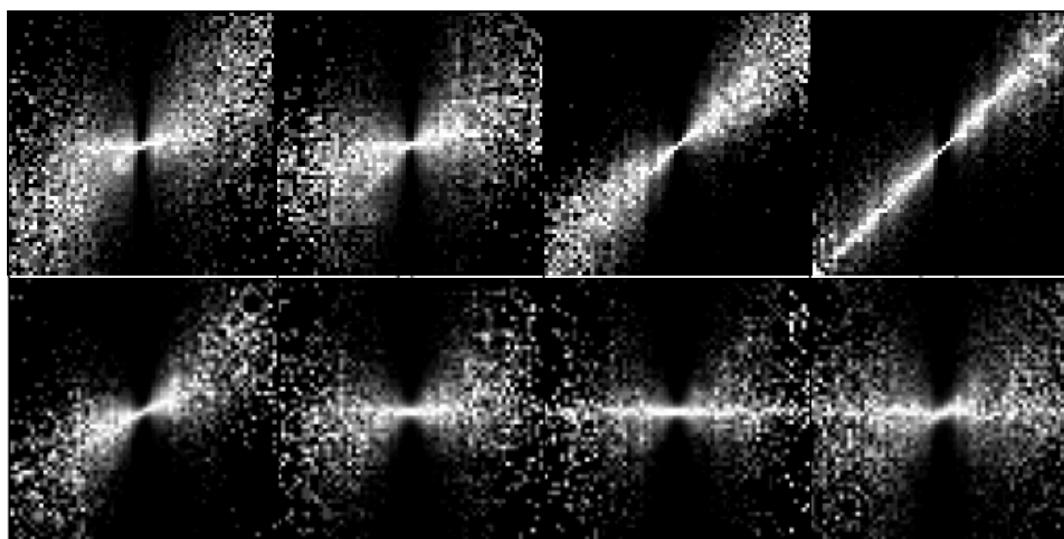


Pairwise Wavelet Histograms

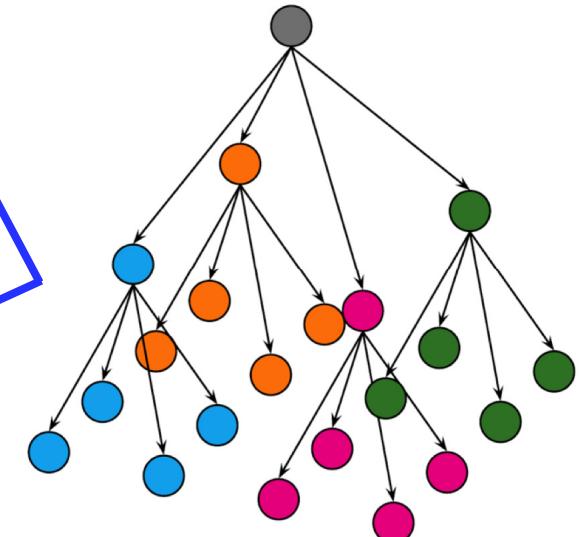


Orientation *Scale* *Vertical* *Horizontal*

Image



Model



Outline

Multiscale Models for Natural Images

- Nonparametric Hidden Markov Trees (HDP-HMTs)
- Learning with Monte Carlo methods
- Truncated representations for efficient learning from large datasets

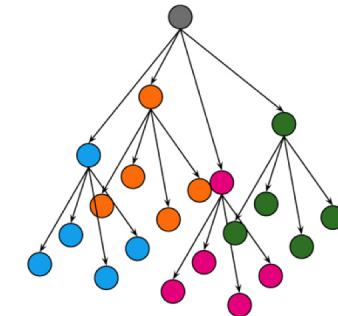
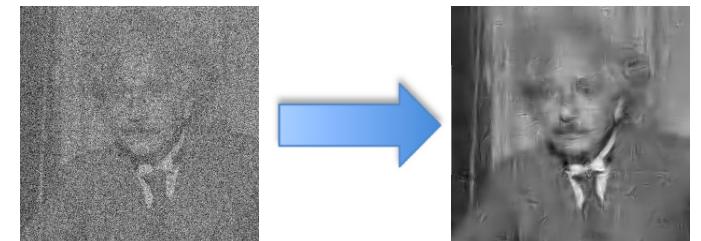


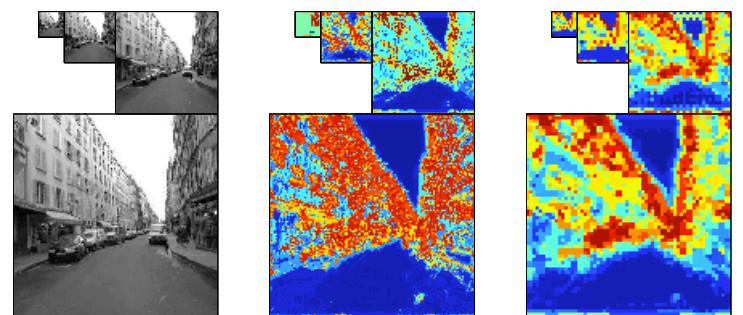
Image Denoising

- Transfer natural image statistics for making robust predictions



Natural Scene Analysis

- Global, data-driven scene models via HDP-HMT
- Fast categorization via Belief Propagation methods



Denoising Images in Wavelet-Domain

Noisy



*HDP-HMT
(Emp. Bayes)*



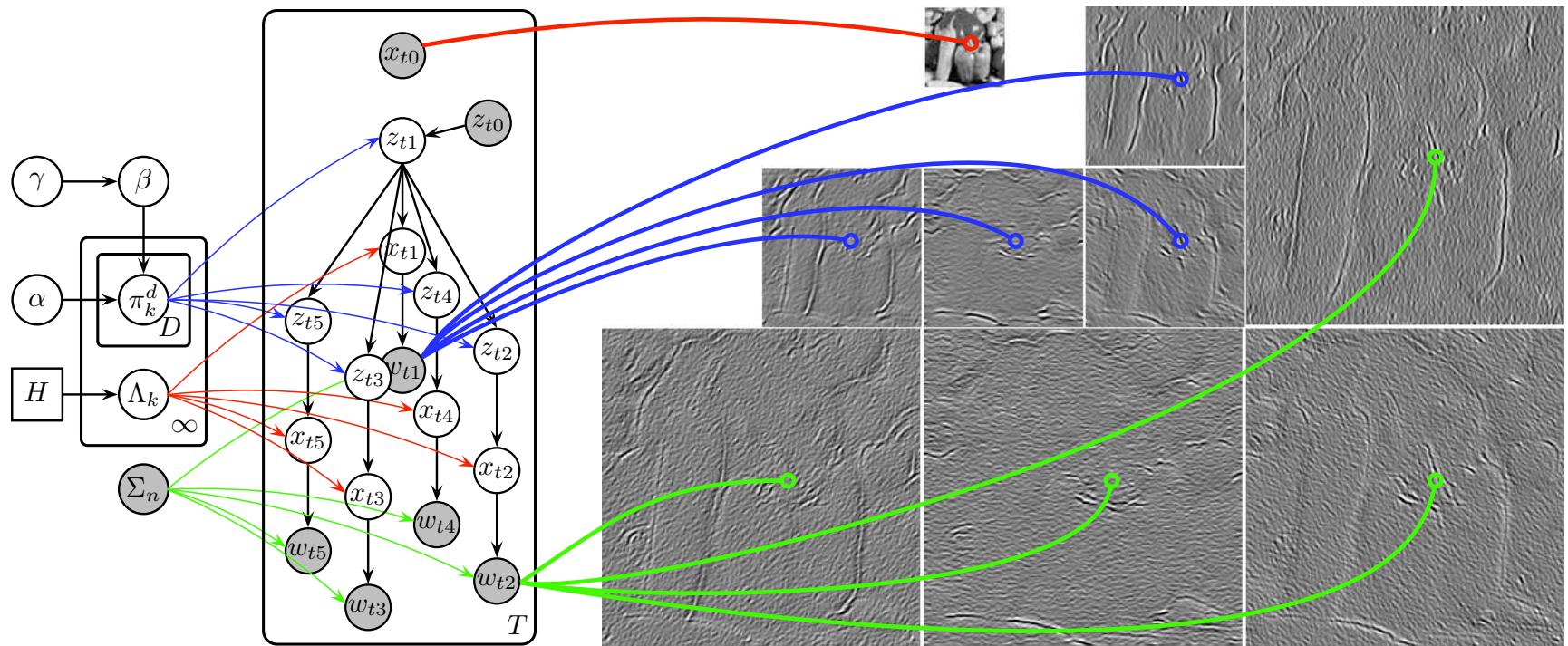
*HDP-HMT
(Transfer)*



Apply learned statistics to **denoise** images

- Using a global model increases **robustness**
- Exploit availability of large image databases to develop efficient **transfer** denoising algorithms
- Improve performance by **reusing** statistics of **clean** images

HDP-HMT for noisy data



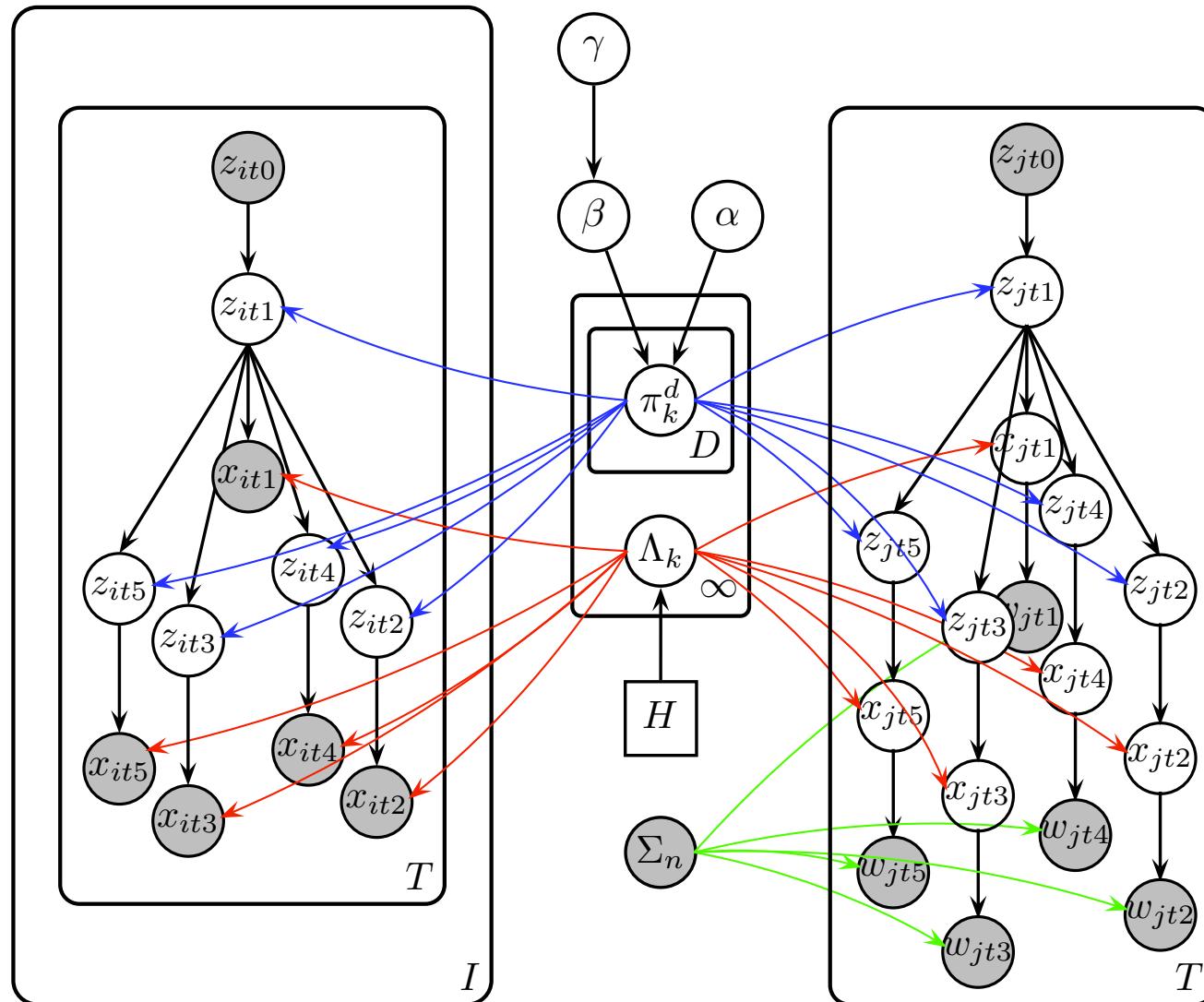
$x_{ti} \rightarrow$ unobserved vector of *clean* wavelet coefficients

$w_{ti} \rightarrow$ observed vector of *noisy* wavelet coefficients

$\Sigma_n \rightarrow$ noise variance

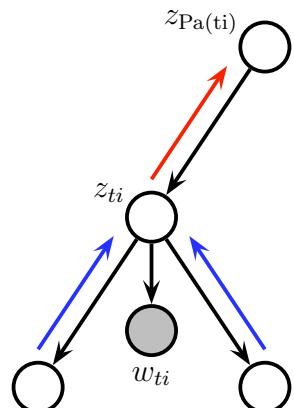
$$w_{ti} \sim \mathcal{N}(x_{ti}, \Sigma_n)$$

... and for clean data as well

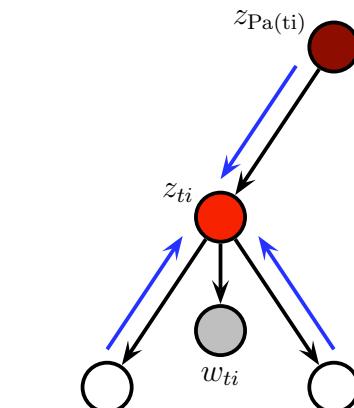


Learning via Gibbs Sampling

*Sample hidden states **jointly** using BP:*



Compute msgs from leaves to root



Sample states from root to leaves

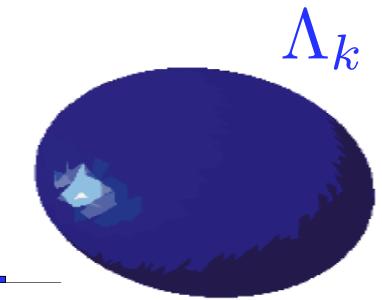
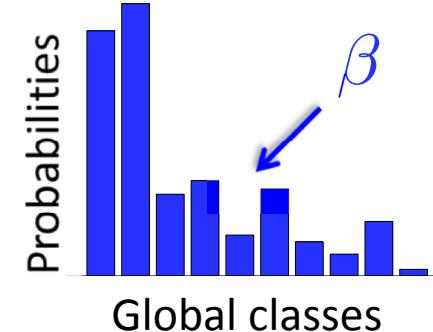
For noisy images, sample clean wavelet coefficients:

$$x_{ti} \sim \mathcal{N}(\mu, \Sigma)$$

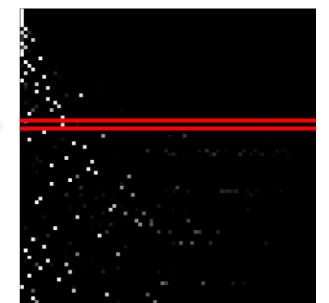
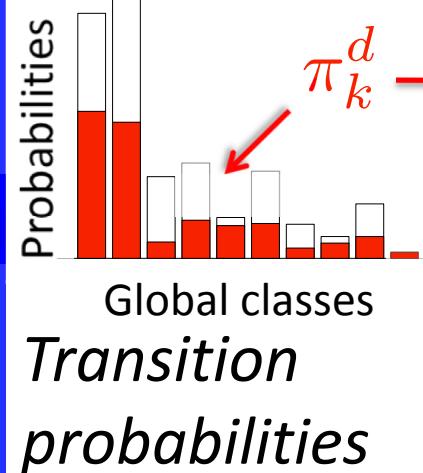
$$\mu = (\Lambda_{z_{ti}}^{-1} + \Sigma_n^{-1})^{-1} \Sigma_n^{-1} w_{ti}$$

$$\Sigma = (\Lambda_{z_{ti}}^{-1} + \Sigma_n^{-1})^{-1}$$

Sample parameters:

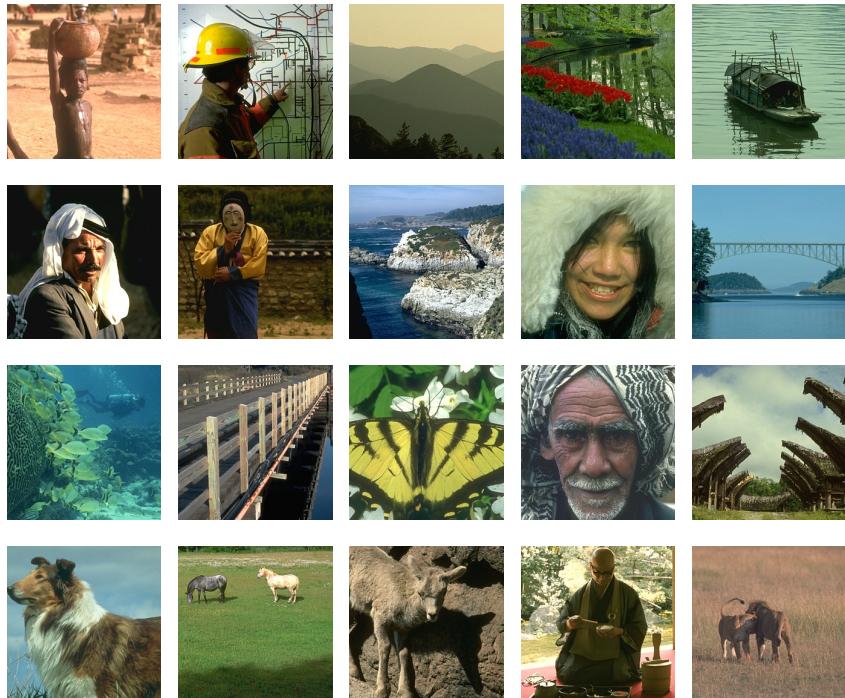


Feature distributions

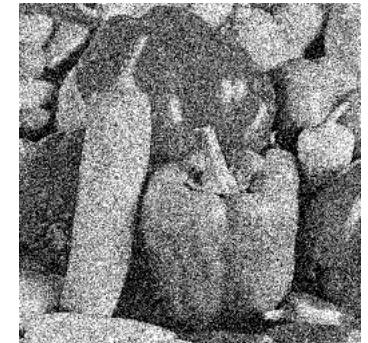


Transition probabilities

Estimating Clean Images



Empirical Bayesian approach estimates model parameters from the noisy image



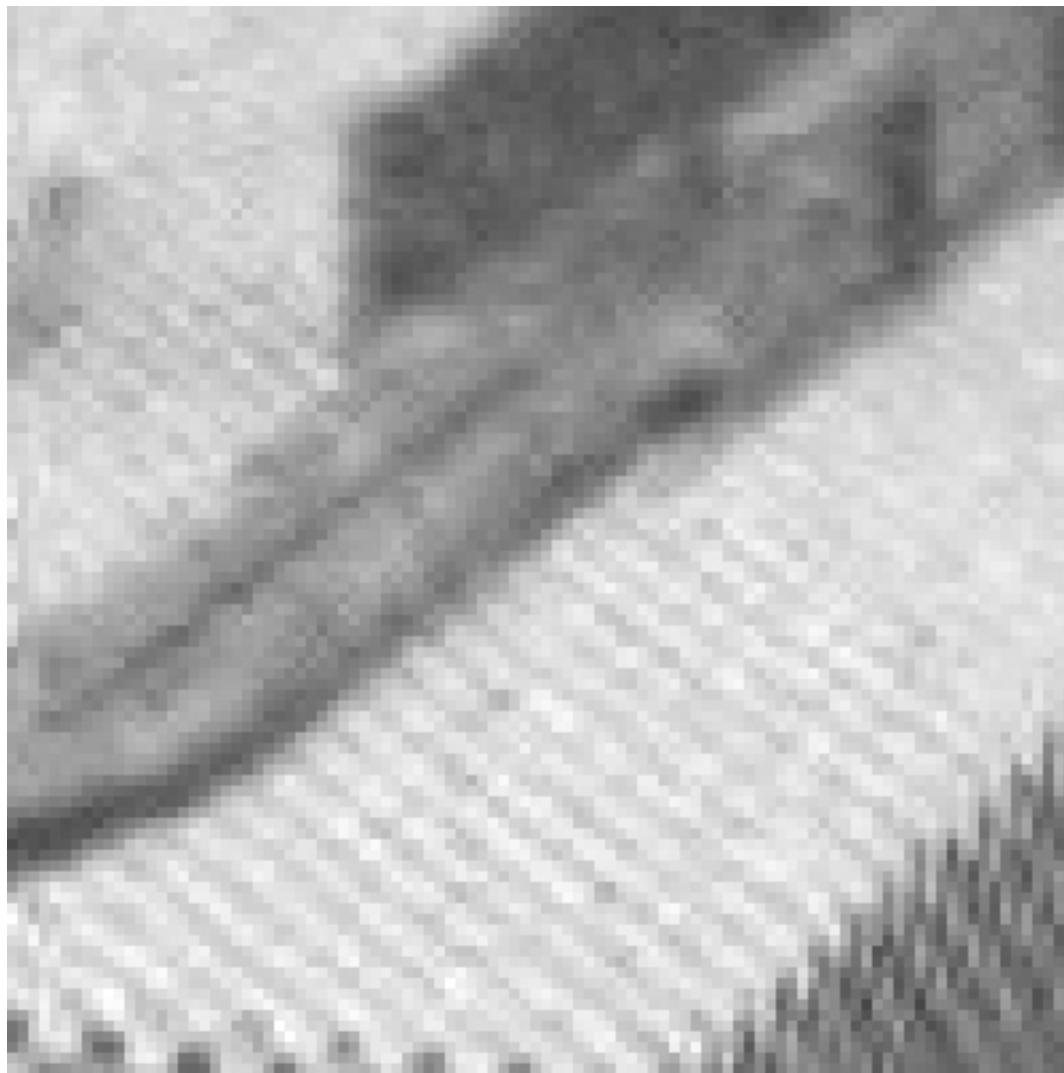
Transfer denoising approach **reuses** multiscale hidden state patterns of **clean** images for making robust predictions



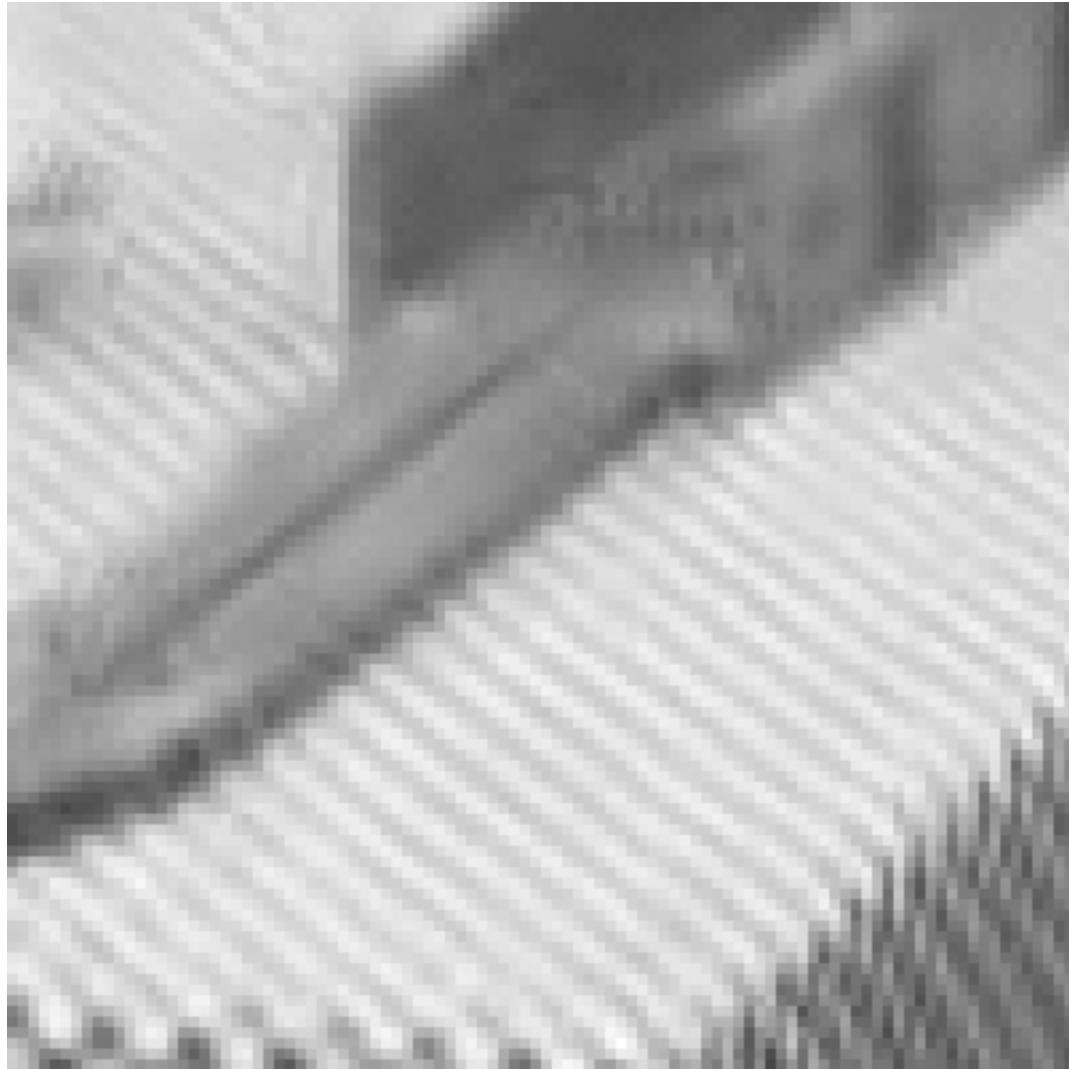
$$\mathbb{E}[x_{ti} \mid \mathbf{w}, \theta^{(s)}] = \sum_{k=1}^{K_s} p(z_{ti} = k \mid \mathbf{w}, \theta^{(s)}) \mathbb{E}[x_{ti} \mid w_{ti}, \Lambda_k^{(s)}]$$

From belief propagation *Linear least-squares smoothing*

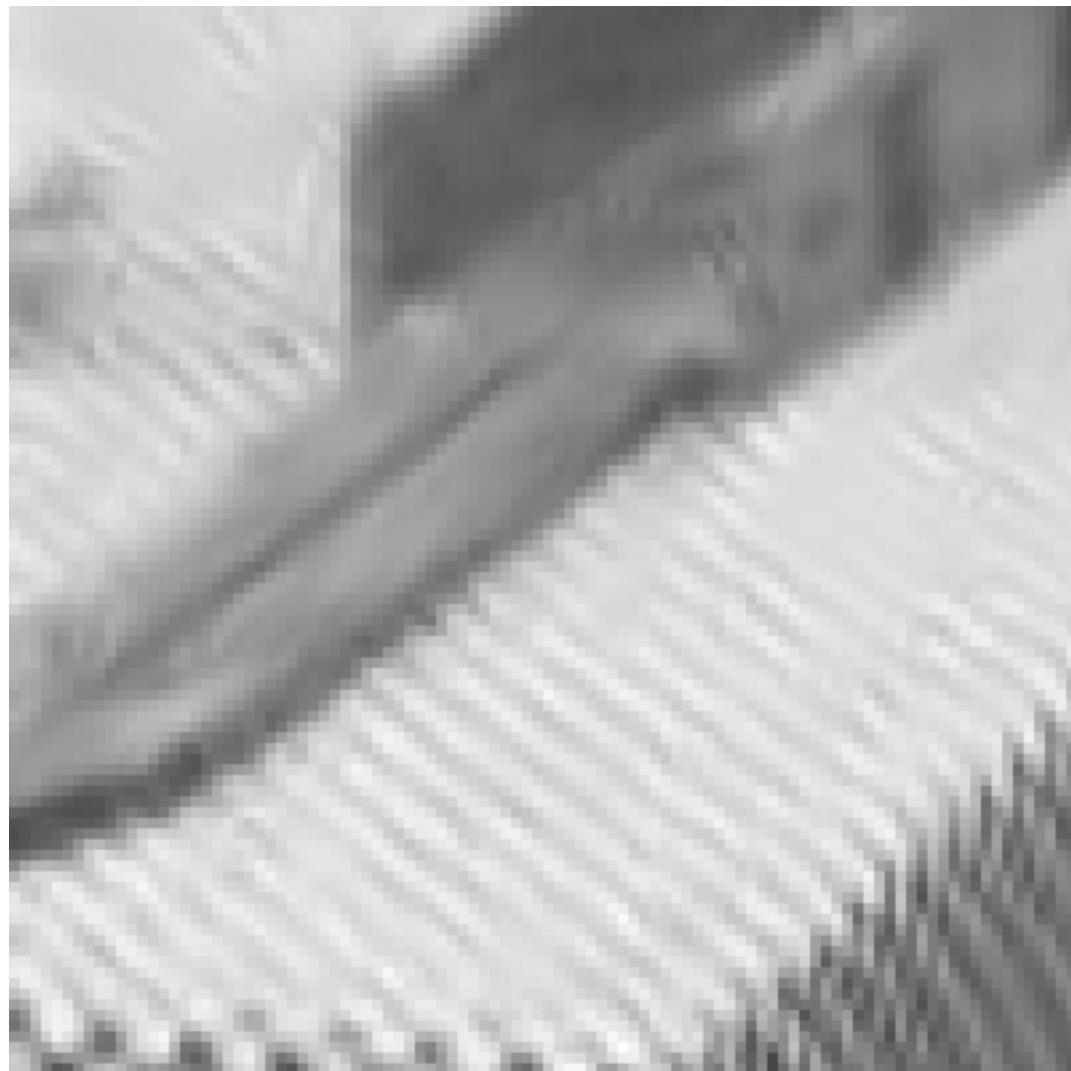
Denoising: Binary HMT



Denoising: HDP-HMT (Emp. Bayes)

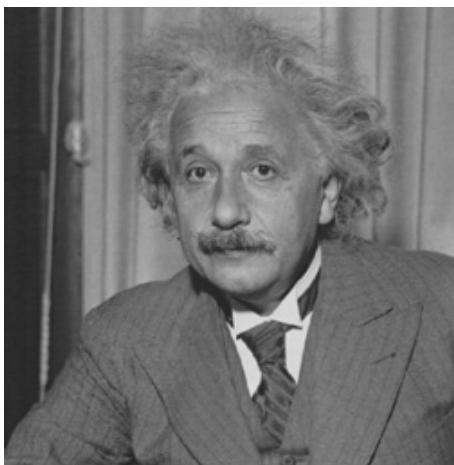


Denoising: Local GSM



Portilla, Strela, Wainwright, & Simoncelli, 2003

Denoising Einstein

<i>Noisy</i> $10.60 \text{ dB}, 0.057$	<i>HDP-HMT</i> (Emp. Bayes) $25.64 \text{ dB}, 0.564$	<i>HDP-HMT</i> (Transfer) $26.80 \text{ dB}, 0.664$
		
<i>Original</i> 	<i>BLS-GSM</i> $26.38 \text{ dB}, 0.647$ 	<i>BM3D</i> $26.49 \text{ dB}, 0.659$ 

Denoising Hill

Noisy

$8.12 \text{ dB}, 0.038$



Original



HDP-HMT

(*Emp. Bayes*)

$24.56 \text{ dB}, 0.540$



HDP-HMT

(*Transfer*)

$24.74 \text{ dB}, 0.568$



BLS-GSM

$24.54 \text{ dB}, 0.544$



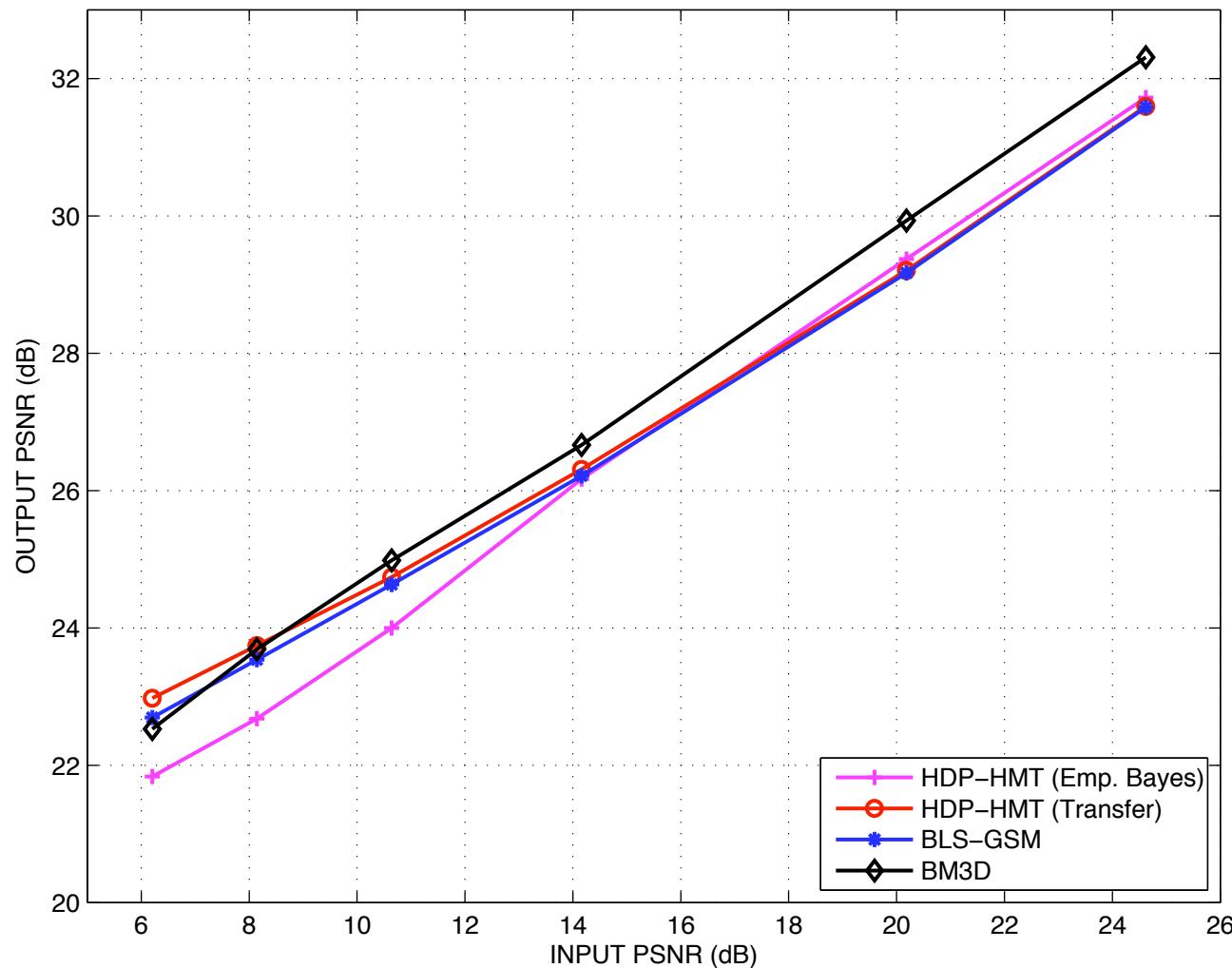
BM3D

$24.39 \text{ dB}, 0.548$



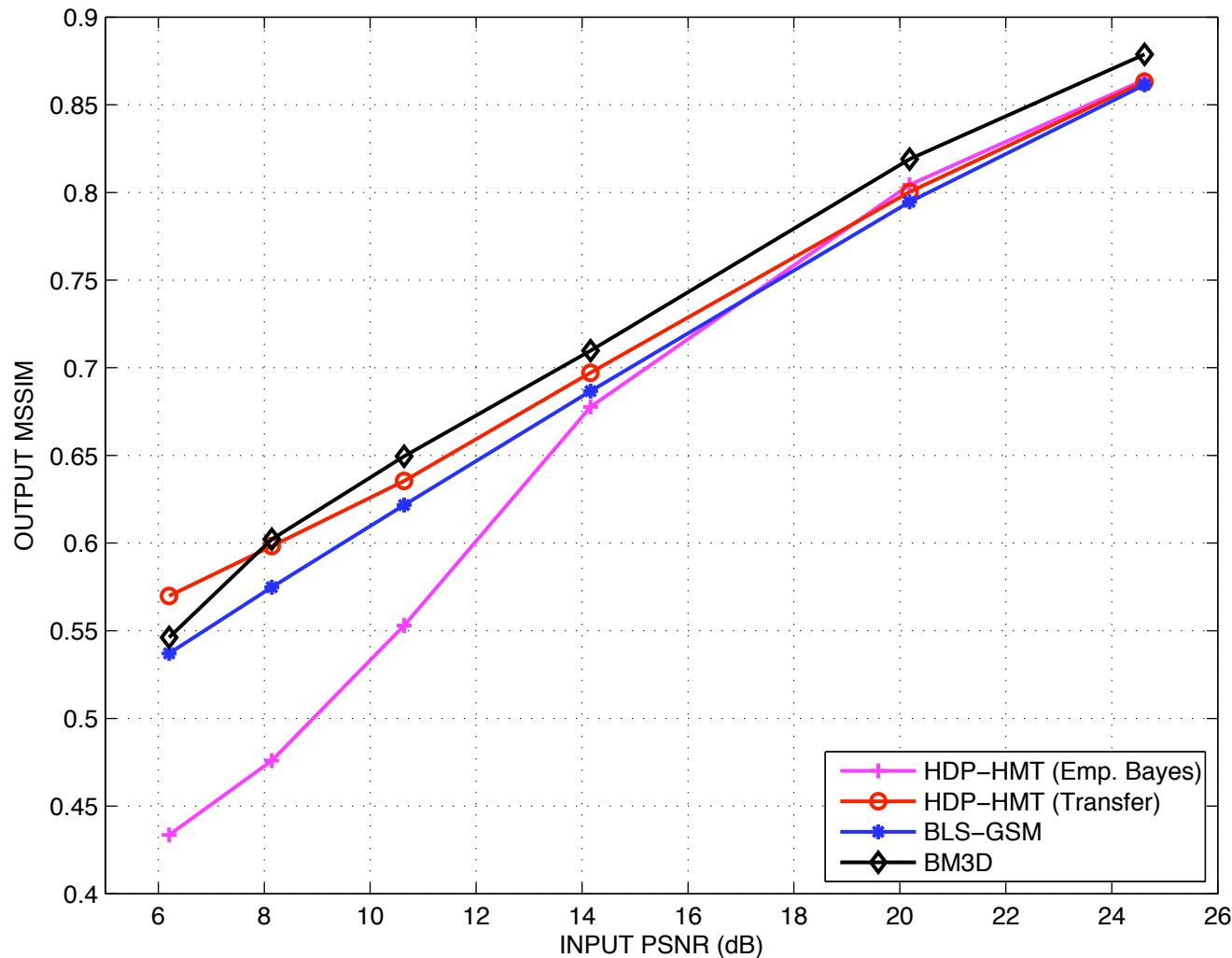
Average Denoising Performance

Peak signal-to-noise ratio (PSNR)



Average Denoising Performance

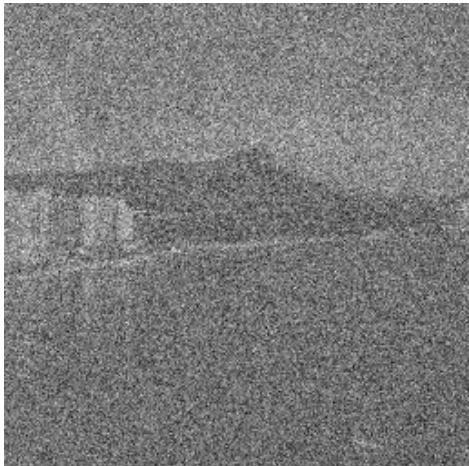
Mean structural similarity index (SSIM)



Natural Scene Denoising

Noisy

$8.14 \text{ dB}, 0.033$



Original



HDP-HMT

(*Emp. Bayes*)

$24.24 \text{ dB}, 0.519$



HDP-HMT

(*Transfer*)

$26.50 \text{ dB}, 0.794$



BLS-GSM

$25.59 \text{ dB}, 0.726$

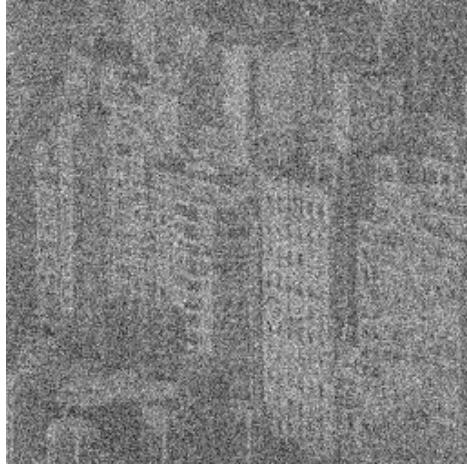


BM3D

$25.74 \text{ dB}, 0.751$



Natural Scene Denoising

<i>Noisy</i> $8.14 \text{ dB}, 0.177$	<i>HDP-HMT</i> <i>(Emp. Bayes)</i> $18.55 \text{ dB}, 0.484$	<i>HDP-HMT</i> <i>(Transfer)</i> $18.77 \text{ dB}, 0.486$
		
<i>Original</i> 	<i>BLS-GSM</i> $18.59 \text{ dB}, 0.454$ 	<i>BM3D</i> $18.65 \text{ dB}, 0.470$ 

Natural Scene Denoising

HDP-HMT (Transfer)

$23.28 \text{ dB}, 0.653$



BLS-GSM

$23.14 \text{ dB}, 0.617$



BM3D

$23.23 \text{ dB}, 0.651$



Outline

Multiscale Models for Natural Images

- Nonparametric Hidden Markov Trees (HDP-HMTs)
- Learning with Monte Carlo methods
- Truncated representations for efficient learning from large datasets

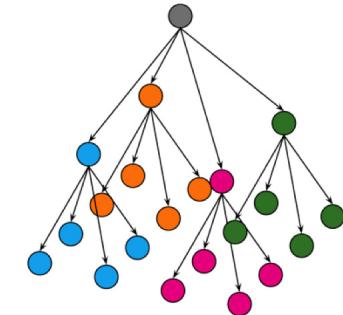
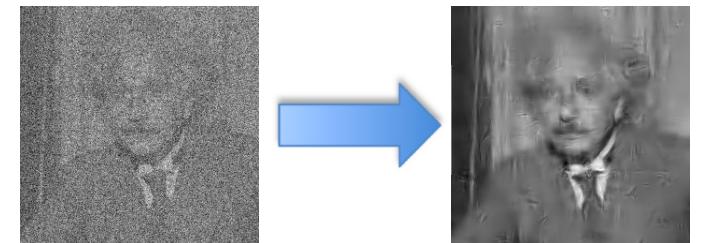


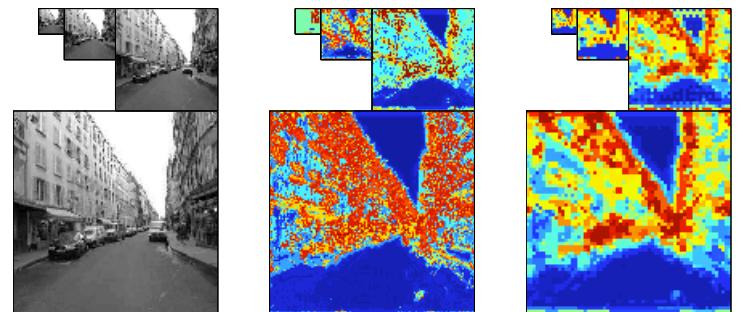
Image Denoising

- Transfer natural image statistics for making robust predictions

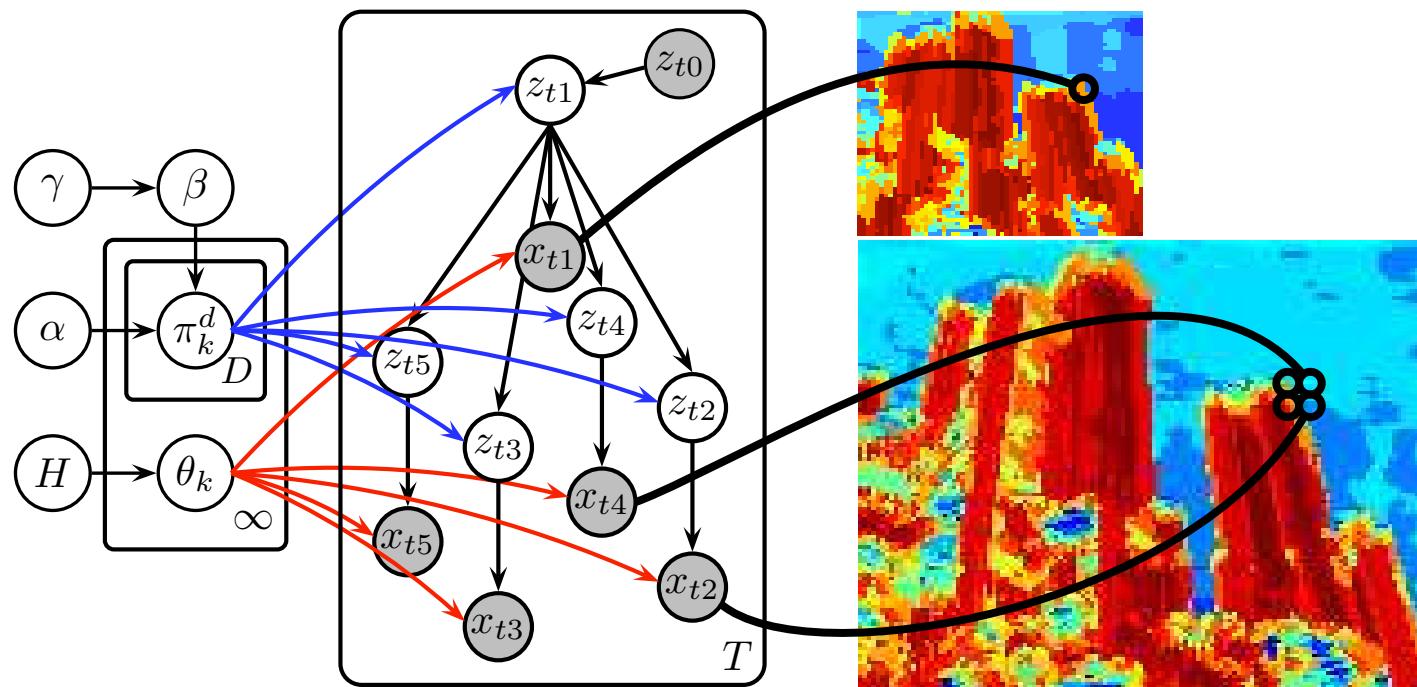


Natural Scene Analysis

- Global, data-driven scene models via HDP-HMT
- Fast categorization via Belief Propagation methods



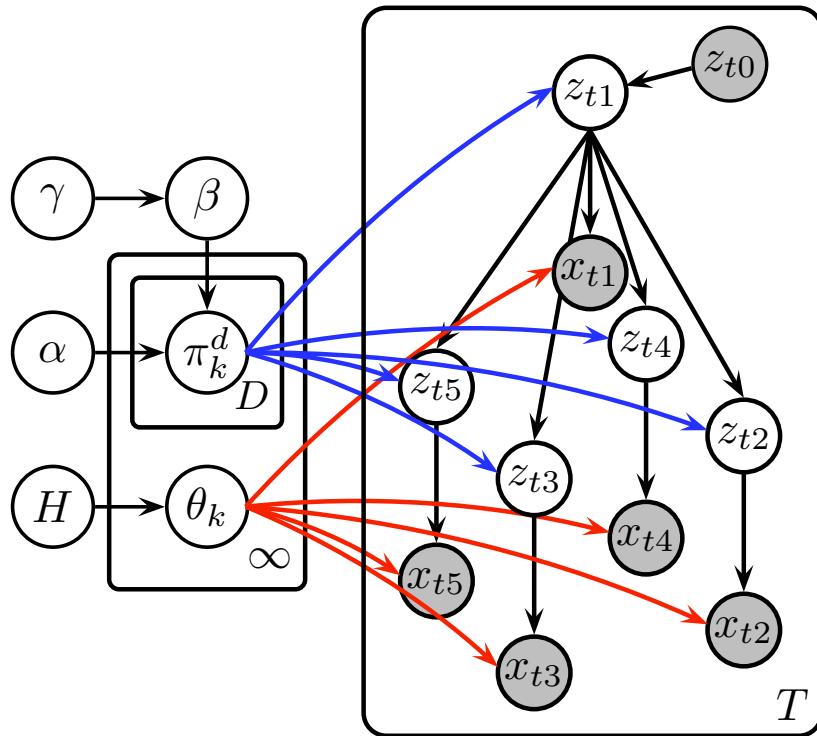
HDP-HMT Scene Model



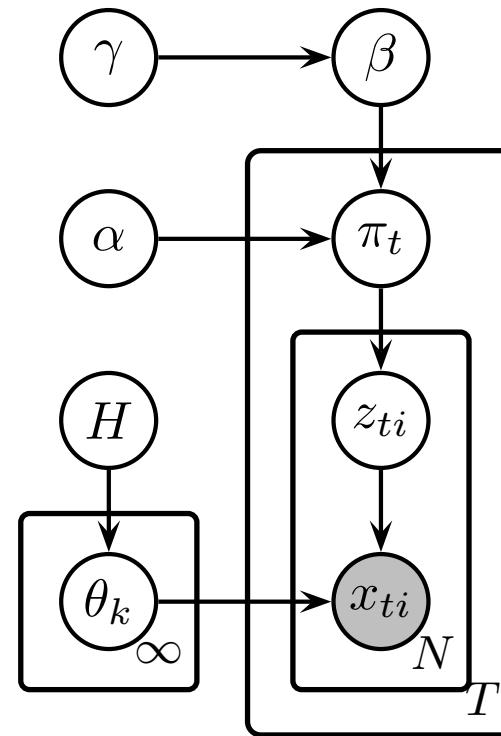
- Hidden states z_{ti} generate vectors of clean wavelet coefficients x_{ti} at multiple orientations or **SIFT-descriptors**

... versus baseline HDP-BOF

HDP-HMT



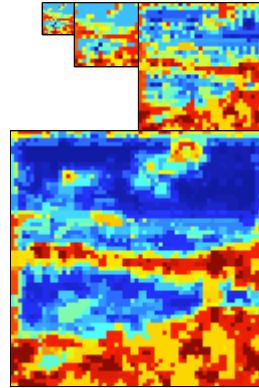
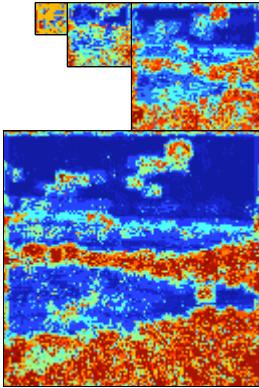
HDP-BOF



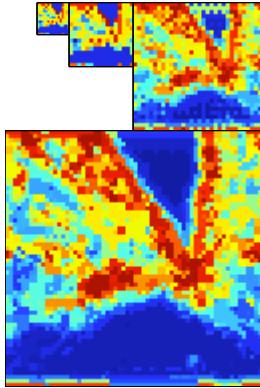
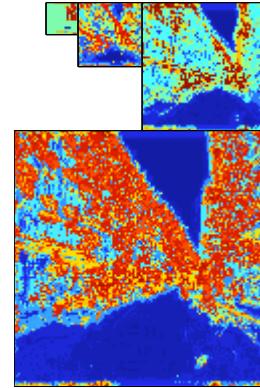
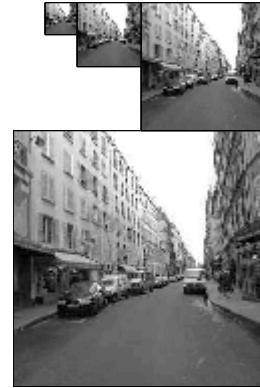
- A nonparametric Bayesian extension of the LDA-based model for scene categorization by Fei-Fei and Perona (2005), which ignores spatial dependencies in the appearance of locally extracted image features
- The HDP-HMT further extends this by incorporating dependencies in feature appearance across location and scale

MAP assignments (sp5)

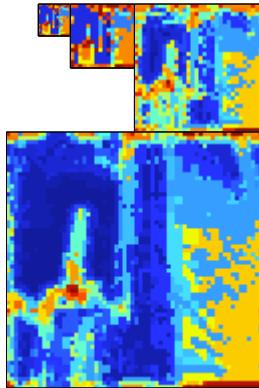
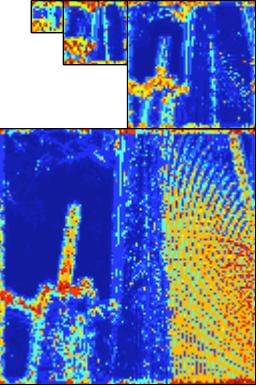
coast



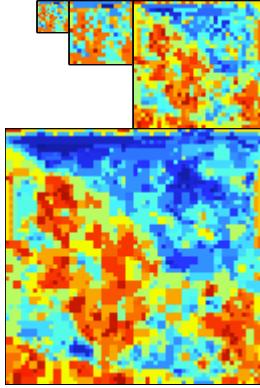
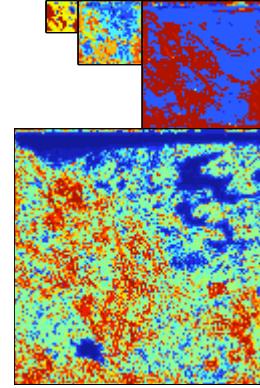
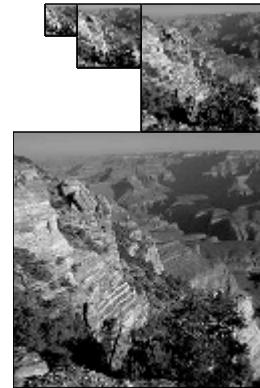
street



tall building



mountain



Input Image

HDP-BOF

HDP-HMT

Input Image

HDP-BOF

HDP-HMT

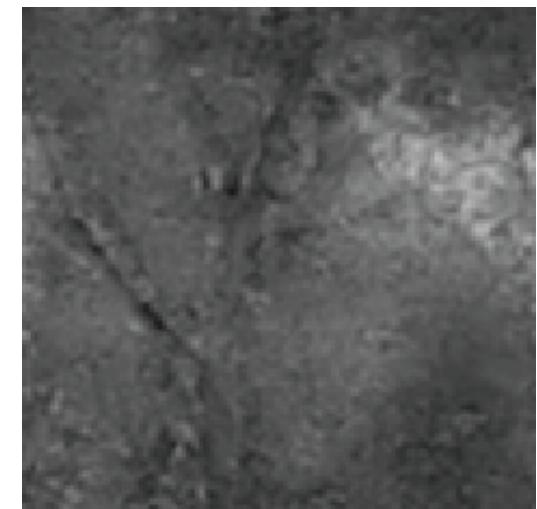
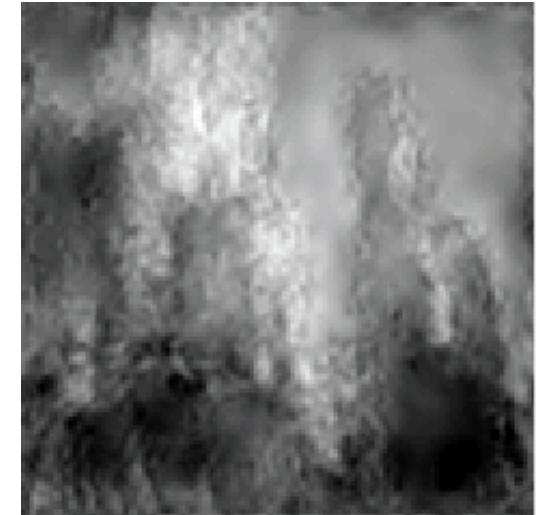
Samples given MAP states



Input Image



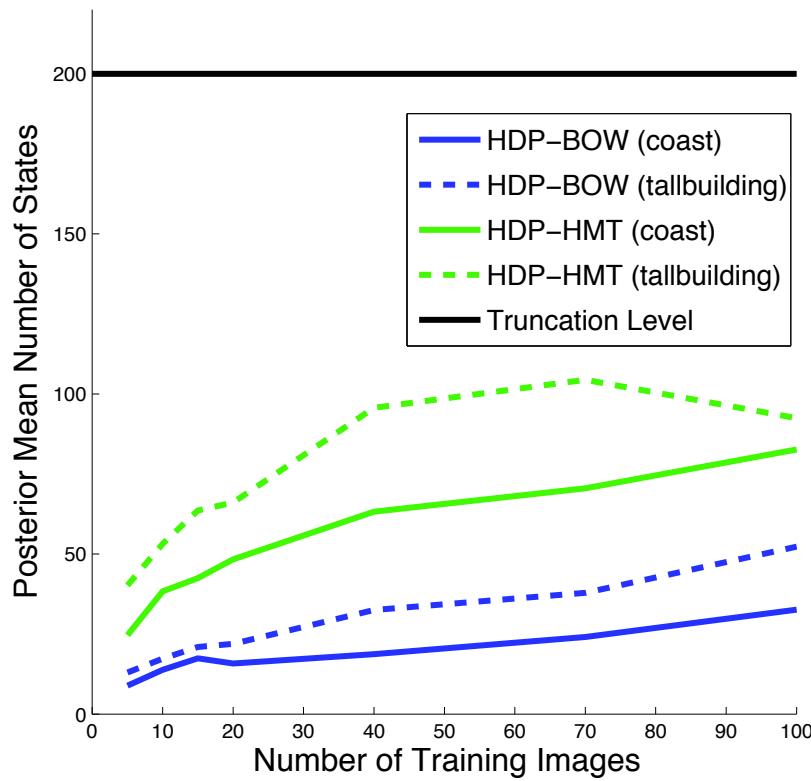
**HDP Hidden
Markov Tree**



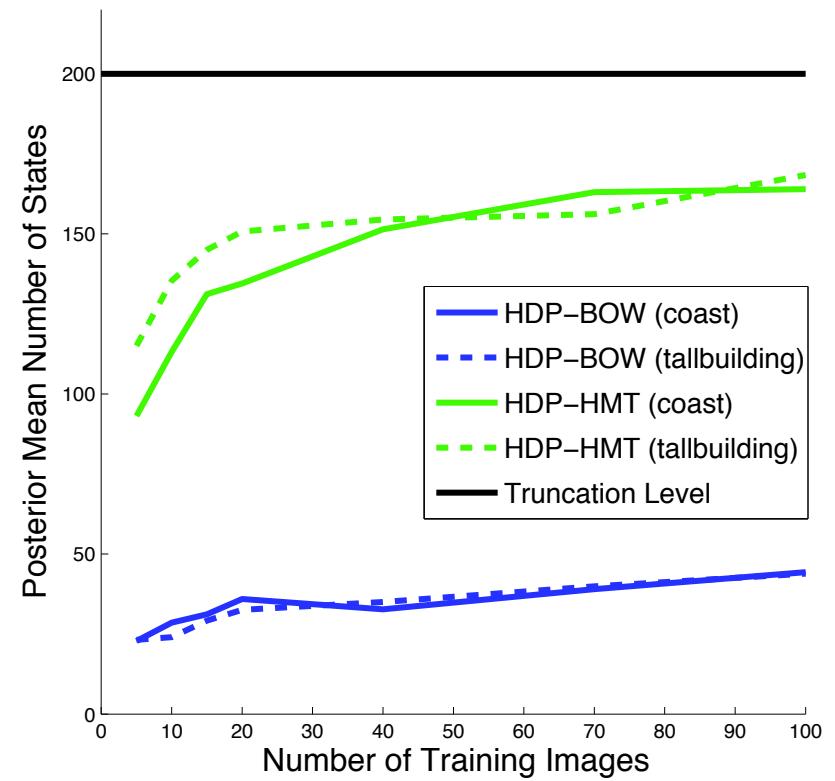
HDP Bag of Features

Number of States

Wavelet (sp5)



SIFT



Categorizing Natural Scenes

Wavelet (sfp7)

coast	77.5	0.6	10.0	0.0	0.6	10.6	0.6	0.0
forest	0.0	91.4	0.0	0.0	5.5	0.8	2.3	0.0
highway	3.3	0.0	75.0	0.0	10.0	10.0	1.7	0.0
inside city	0.9	0.9	2.8	77.8	0.0	3.7	9.3	4.6
mountain	0.6	13.8	4.6	0.6	63.2	9.2	8.0	0.0
open country	8.6	10.0	3.3	0.5	11.0	61.9	4.8	0.0
street	0.0	1.1	5.4	2.2	7.6	0.0	81.5	2.2
tall building	0.0	0.0	2.6	13.5	0.6	0.6	8.3	74.4

SIFT

coast	90.0	0.6	1.2	0.0	1.9	6.2	0.0	0.0
forest	0.0	87.5	0.0	0.0	7.8	4.7	0.0	0.0
highway	6.7	0.0	80.0	1.7	1.7	5.0	5.0	0.0
inside city	0.0	0.0	1.9	87.0	0.0	0.0	9.3	1.9
mountain	1.1	0.6	0.6	0.0	90.2	5.7	0.6	1.1
open country	11.0	1.9	1.0	0.0	5.7	80.0	0.5	0.0
street	0.0	0.0	4.3	2.2	2.2	0.0	91.3	0.0
tall building	0.0	0.0	0.0	9.0	0.6	0.0	4.5	85.9

HDP-BOF [75.3 %]

coast	90.6	0.6	4.4	0.0	1.9	2.5	0.0	0.0
forest	0.0	85.2	0.8	0.0	8.6	3.1	2.3	0.0
highway	8.3	0.0	80.0	0.0	6.7	0.0	3.3	1.7
inside city	1.9	0.9	5.6	75.0	0.9	0.9	10.2	4.6
mountain	1.7	1.1	2.9	0.0	91.4	1.1	1.7	0.0
open country	18.1	4.3	3.3	1.0	13.3	59.5	0.5	0.0
street	0.0	0.0	8.7	1.1	7.6	0.0	81.5	1.1
tall building	0.0	0.0	1.9	12.2	0.0	0.0	3.8	82.1

HDP-BOF [82.4 %]

coast	86.2	1.2	4.4	0.0	0.6	7.5	0.0	0.0
forest	0.0	91.4	0.0	0.0	4.7	3.1	0.8	0.0
highway	6.7	0.0	75.0	1.7	3.3	6.7	6.7	0.0
inside city	0.0	0.9	3.7	82.4	0.0	0.9	10.2	1.9
mountain	0.6	4.0	3.4	0.0	81.0	8.0	2.3	0.6
open country	11.0	5.2	2.9	0.0	7.6	72.9	0.5	0.0
street	0.0	0.0	6.5	2.2	1.1	0.0	89.1	1.1
tall building	0.0	0.0	0.6	7.1	1.3	0.0	10.3	80.8

HDP-HMT [80.7 %]

HDP-HMT [86.5 %]

Average Categorization Performance

	Wavelet (sfp7)		SIFT	
Man-made	82.9	85.4	86.4	89.7
Natural	78.6	83.5	85.7	87.7
Eight	75.3	80.7	82.4	86.5
Thirteen			75.9	81.8
Fifteen			69.7	77.1
	HDP-BOW	HDP-HMT	HDP-BOW	HDP-HMT

Summary and Conclusions

- Presented a hierarchical nonparametric Bayesian model for multiscale datasets with complex non-local dependencies
- Presented MCMC methods for learning HDP-HMT parameters from clean and noisy images
- Truncated representations of the DP to allow efficient blocked sampling algorithms and learning from large datasets
- The HDP-HMT captures complex natural image structures and leads to effective learning algorithms for denoising and categorization
- Robust restoration with natural image statistics transfer

Pairwise Statistics of Wavelets

