

Seeing as Statistical Inference

Song-Chun Zhu

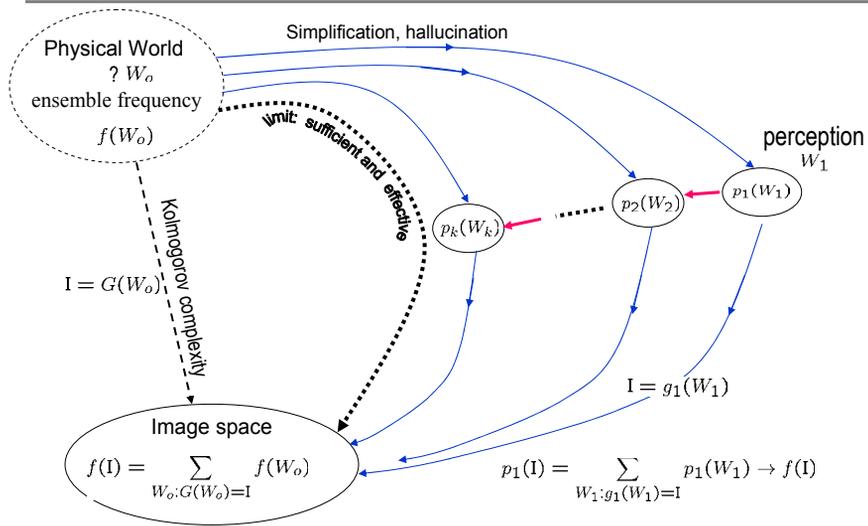
Center for Image and Vision Science
University of California, Los Angeles

An overview of joint work with many students and colleagues

MSRI, January, 2005

Song-Chun Zhu

Pursuit of image models

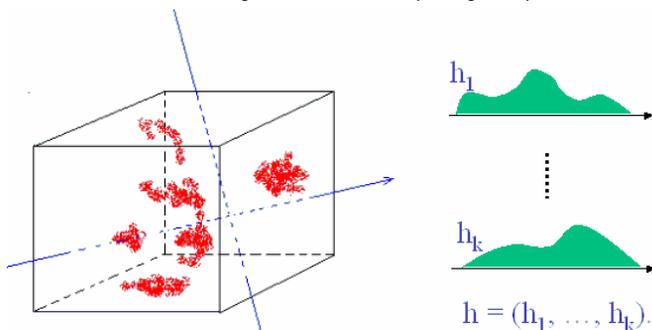


MSRI, January, 2005

Song-Chun Zhu

Family I: Descriptive Modeling

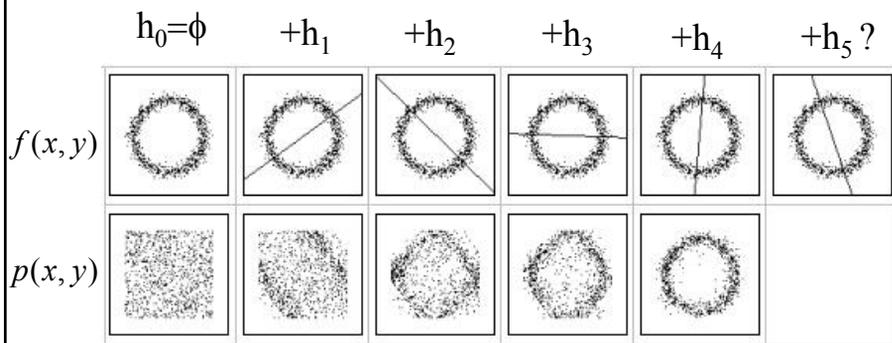
The extracted statistics are marginal distributions (histograms).



1. Given the observed statistics, a maximum entropy model is learned to reproduce the observed statistics.
2. The most informative statistics (features) is selected by minimizing the entropy of the max. ent. model

Leads to Markov random Field and Gibbs models

Toy Example: Estimating 1D manifold embedded in 2D



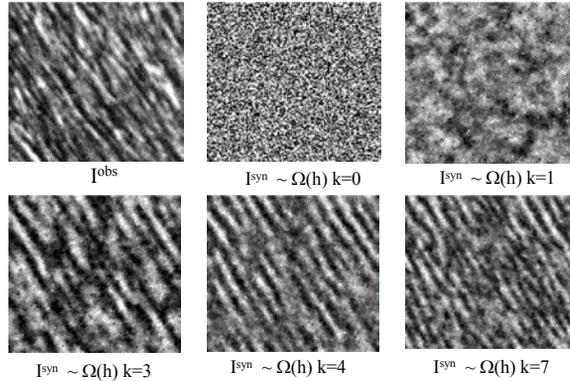
By Ce Liu, 2001.

Augmenting the model by minimax entropy

(Zhu, Wu, Mumford 95-97)

a texture = $\Omega(h_c) = \{I : h(F_i * I) = h_{c,i}, i = 1, 2, \dots, K, \Lambda \rightarrow Z^2\}$

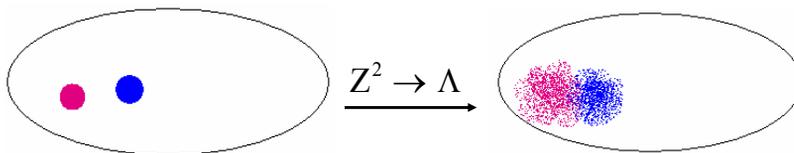
$h_c = (h_1, \dots, h_K)$ are histograms of Gabor filters, i.e. marginal distributions of $f(I)$



MSRI, January, 2005

Song-Chun Zhu

Prob. Model derived from deterministic ensemble



texture ensembles :

$$\Omega(h_c) = \{I : h(I) = h_c\}$$

texture models :

$$p(I_\Lambda | I_{\partial\Lambda}; \beta)$$

Markov random fields and FRAME models on finite lattice (Zhu, Wu, Mumford, 1997):

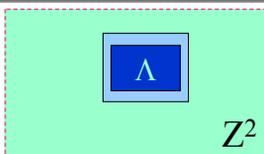
$$p(I_\Lambda | I_{\partial\Lambda}; \beta) = \frac{1}{Z(\beta)} \exp\left\{-\sum_{j=1}^k \beta_j h_j(I_\Lambda | I_{\partial\Lambda})\right\}$$

MSRI, January, 2005

Song-Chun Zhu

Prob. Model derived from deterministic ensemble

Wu and Zhu '99



Theorem

For a very large image from the Julesz ensemble

$$I \sim \Omega(h_c) = \{ I : h(I) = h_c \}$$

any local patch of the image I_Λ given its neighborhood follows a conditional distribution specified by a FRAME model $p(I_\Lambda | I_{\partial\Lambda} : \beta)$

Theorem

As the image lattice goes to infinity, $f(I; h_c)$ is the limit of the FRAME model $p(I_\Lambda | I_{\partial\Lambda} : \beta)$, in the absence of phase transition.

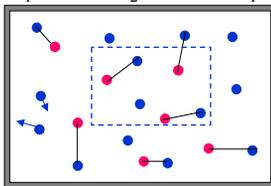


MSRI, January, 2005

Song-Chun Zhu

Correspondence to ensembles in Stat. Physics

$$N_1 = 10^{23}, N_2 = 10^{18} \quad N_1 + N_2 = 10^{23}$$



$$\text{Micro-canonical Ensemble} = \Omega(N, E, V) = \{ s : h(S) = (N, V, E) \}$$

A large system with fix number of elements N , volume V , and energy E .

What are the basic elements in the ensemble of visual patterns?

The minimax entropy principle does not tell us about it !

MSRI, January, 2005

Song-Chun Zhu

Family II: layered generative modeling

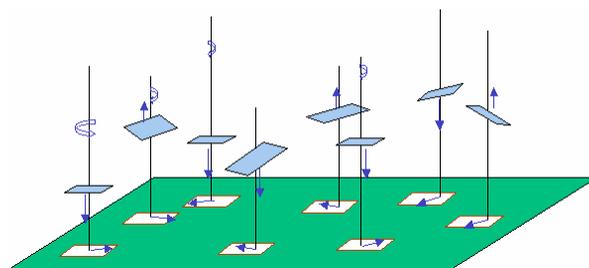
Seeking Fundamental Image Elements (Isolated)

Sparse coding, Olshausen and Felds, 95

Transformed Component Analysis, Frey and Jojic 00

Textons, Leung and Malik 99, Guo, Zhu and Wu, 01,02, (Dated back to Julesz 70s)

Image primitives, Guo, Zhu and Wu, ICCV, 03 (Dated back to Marr)



MSRI, January, 2005

Song-Chun Zhu

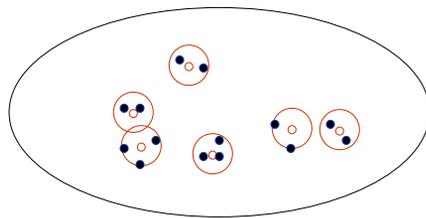
Minimizing the Kolmogorov entropy

The learning problem is to pursue the *best dictionary* so as to minimize the the *Kolmogorov entropy*.

Suppose we have a set of signals (in contrast to assuming a underlying density f)

$$\Omega = \{I_1, I_2, \dots\}$$

which lie on a low-dimensional manifold with unknown dimension H . Suppose the ensemble is covered by at least $N(\epsilon)$ -balls with radius ϵ .



$$N(\epsilon) = (1/\epsilon)^H$$

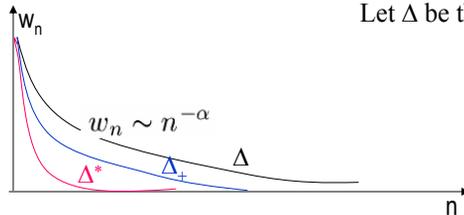
$$H = \log_{1/\epsilon}^{N(\epsilon)} \rightarrow \log |\Omega|$$

MSRI, January, 2005

Song-Chun Zhu

Augmenting generative models

Therefore, the goal is to pursue the optimal bests so that the decreasing rate is fast, which corresponds to minimizing the Kolmogorov entropy.



Let Δ be the dictionary.

When the dictionary is orthogonal, there is a close relation between the Kolmogorov entropy H (or the dimension of signal) and the optimal decreasing rate of the coefficients.

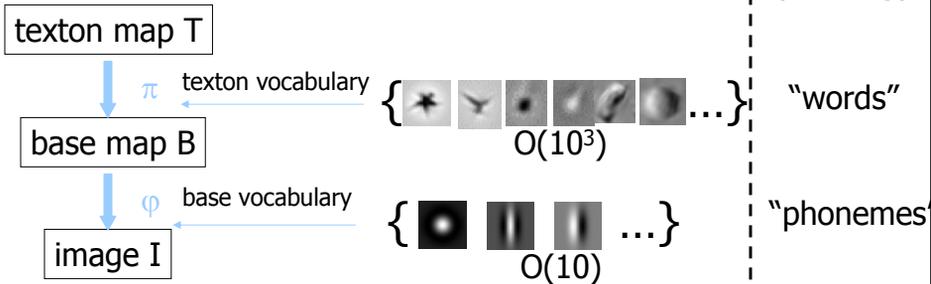
$$\alpha^* = \frac{1}{H} - \frac{1}{2} \quad (\text{Donoho, 98})$$

A Three Level Generative image model

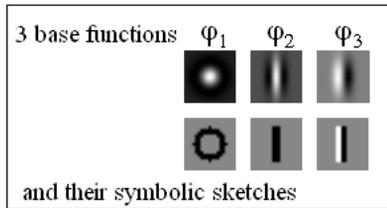
(Zhu et al 02)

Textons are defined as a vocabulary associated with a generative image model.

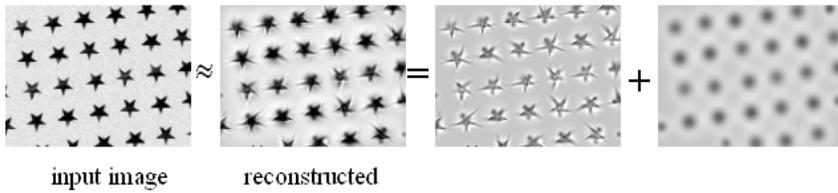
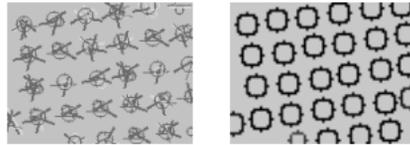
A two level image model:



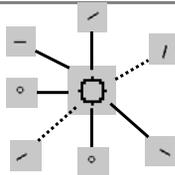
Learning textons



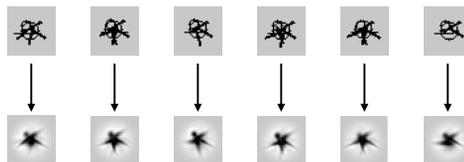
symbolic sketch of
electronic bases and nucleus bases



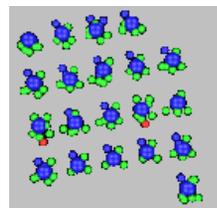
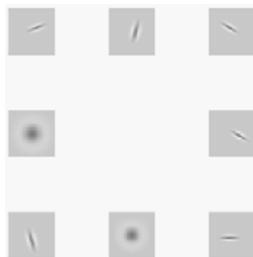
Texton with geometric variations



texton template



a sample of texton instances



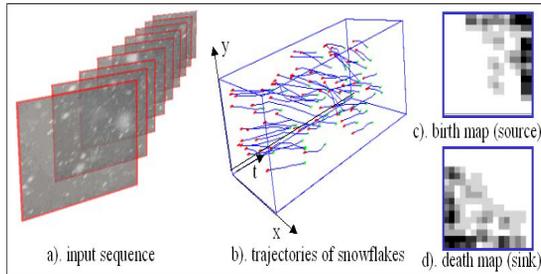
“atomic” model

"Motons" --- Moving textons

Observed Sequence



Synthesized Sequence

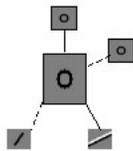


MSRI, January, 2005

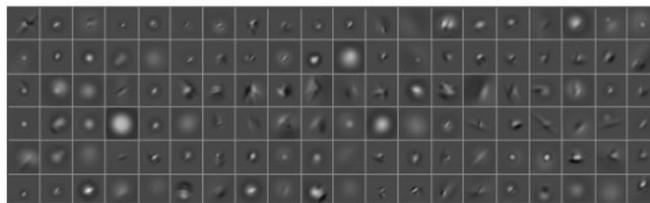
Song-Chun Zhu

"Motons" --- Moving textons

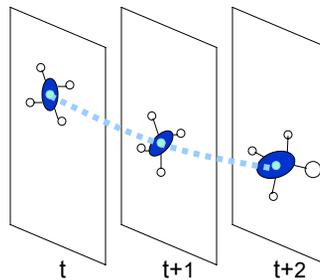
(Wang and Zhu 02,03)



a texton template π



many texton instances randomly sampled from π



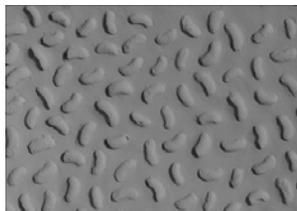
MSRI, January, 2005

Song-Chun Zhu

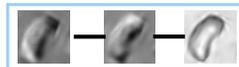
Lightons: textons with lighting variations

(Zhu et al 02)

photometric stereo images

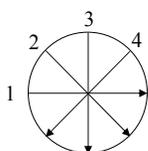


A lighton is a triplet

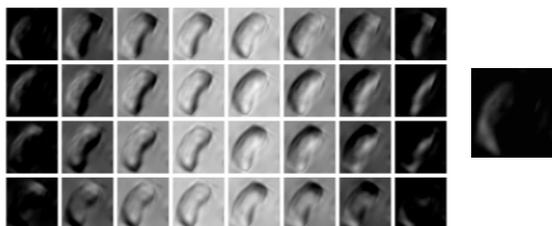


$$B = \alpha_1 b_1 + \alpha_2 b_2 + \alpha_3 b_3$$

Sampling the 3D elements under varying lighting directions



4 lighting directions



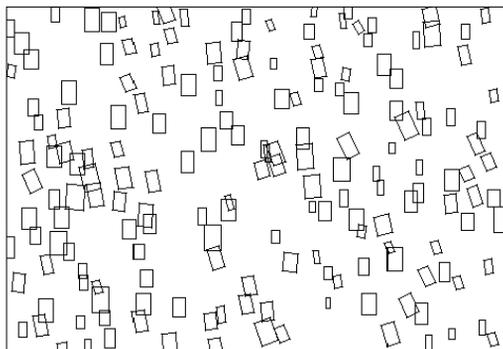
MSRI, January, 2005

Song-Chun Zhu

MCMC simulation of the "texton" process

Guo, Zhu, and Wu 01

The spatial relationship between textons are characterized by the MRF model (family I)



The textons form dynamic neighborhood (Mumford called mixed random fields)
We realized that textons should not be modeled in isolation, and must pursue global structures

MSRI, January, 2005

Song-Chun Zhu

Primal Sketch Model

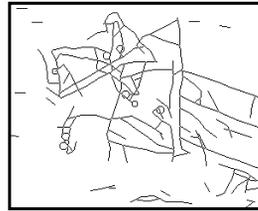
(Guo, Zhu and Wu, 2003)



org image



sketching pursuit process



sketches



syn image



synthesized textures



sketch image

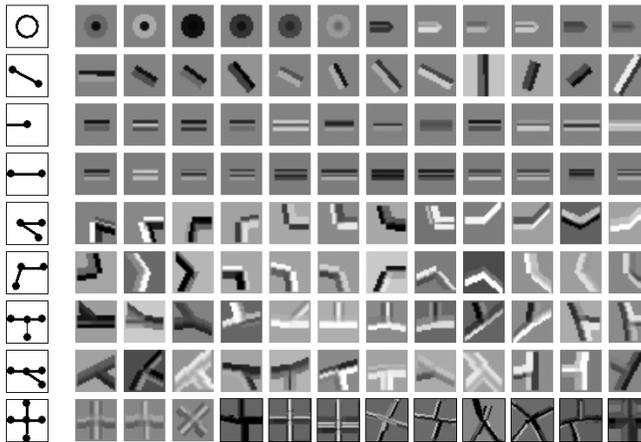
MSRI, January, 2005

Song-Chun Zhu

Simple examples of the image primitive

Learned texton dictionary

(Guo, Zhu and Wu, 2003)



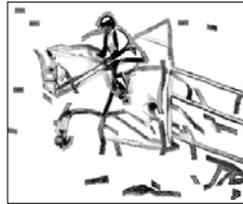
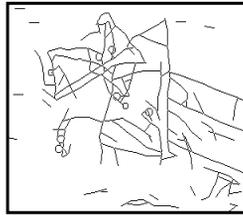
(a)

(b)

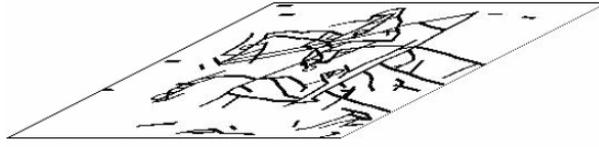
MSRI, January, 2005

Song-Chun Zhu

Primal Sketch: two-level model



Spatial MRF



dictionary Δ \Downarrow generative



Texture MRF

MSRI, January, 2005

Song-Chun Zhu

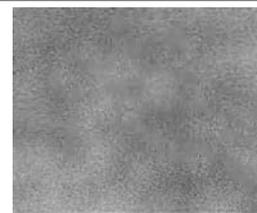
Ix Kurtosis is attributed to structures



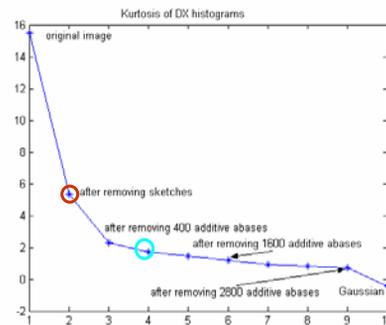
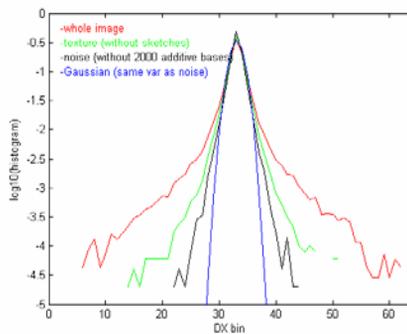
Input image



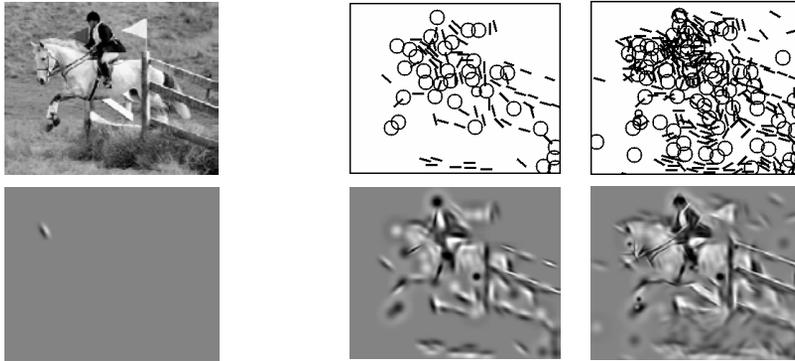
sketchable pixels removed



removing 800 additive bases

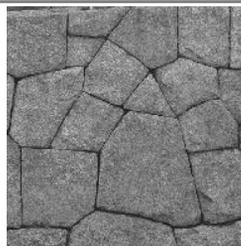


Comparison with linear additive bases

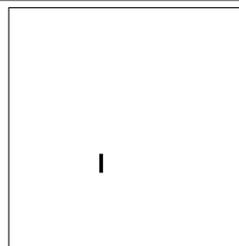


1. Missing the semantics structures
2. Not sparse enough!

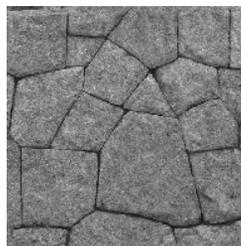
More Example



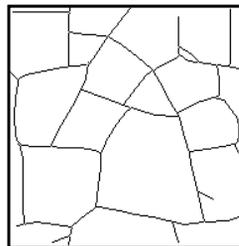
original image



sketching pursuit process



synthesized image

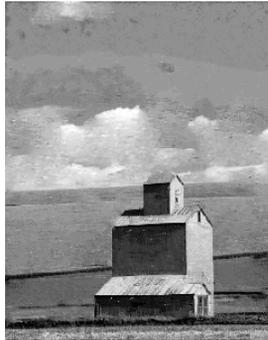


sketches

More example



original image



synthesized image



sketching pursuit process

MSRI, January, 2005

Song-Chun Zhu



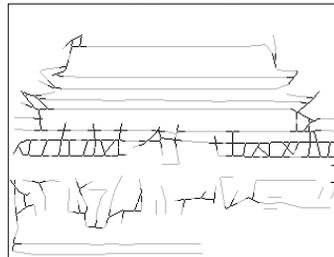
original image



sketching pursuit process



synthesized image



sketches

MSRI, January, 2005

Song-Chun Zhu

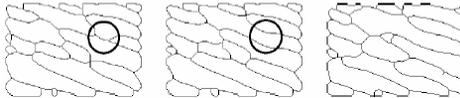
primal sketch over time: topological changes

(Wang and Zhu 2004)

River sequence



Sketch sequence
shows topologic changes



Water sketch over time



Fire sketch over time



MSRI, January, 2005

Song-Chun Zhu

Summary: Generative and Descriptive manifolds

1. Textons are atomic structures in natural images.
Each texton is specified by 4 types of intrinsic dimensions

geometric, photometric, topological, and dynamic

generative manifold : $\Omega_D = \{I : I = g(W), W \in \Omega_d\}$.

2. Textures can also be viewed as manifolds

descriptive manifold : $\Omega(h_c) = \{I : h(I) = h_c\}$.

Theorems show: (1) works for low entropy regimes and (2) works for high entropy regimes

MSRI, January, 2005

Song-Chun Zhu

Climbing up the hierarchy of representation

Sketches of human figure

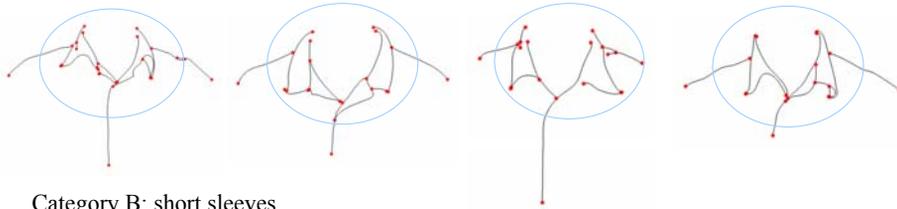


MSRI , January, 2005

Song-Chun Zhu

Supervised Learning of sub-graphs for parts

Category A: Collars



Category B: short sleeves



MSRI , January, 2005

Song-Chun Zhu

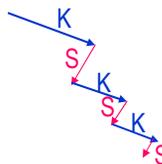
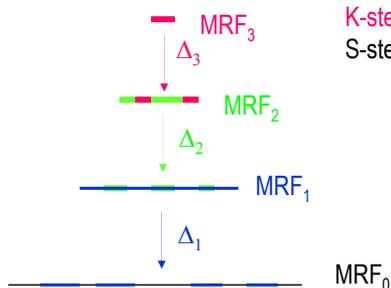
Principle: Integrating the two methods

Now it is not hard to see that the modeling and learning process is to pursue hierarchic models by minimizing the Shannon entropy and Kolmogorov entropy in turns,

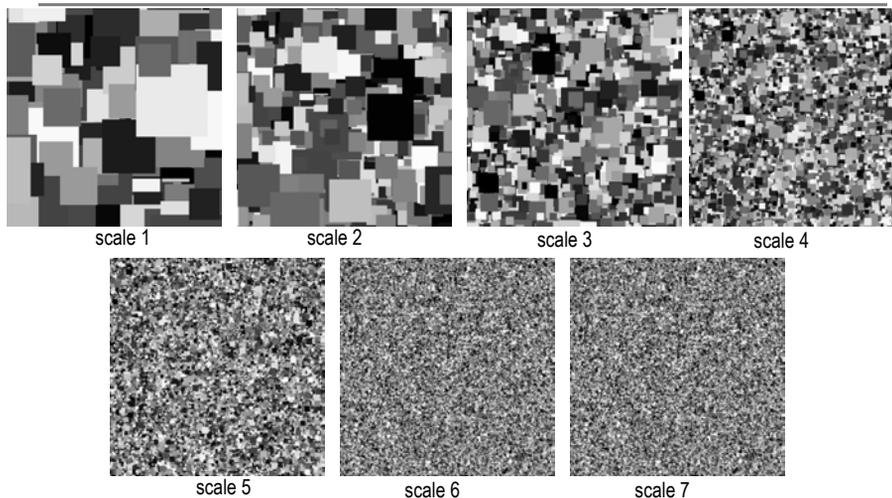
At each level, we run two steps
in reducing the entropy

K-step --- Transform the structures to higher layer

S-step --- Put a MRF/Gibbs model on the remainings

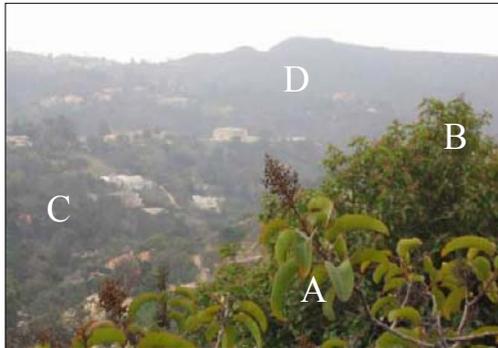


Challenges: model transitions in scale space



We need a seamless transition between the MRF theory and the generative analysis

Nature image contains objects at a range of scales



This picture contains trees/leaves at four ranges of distance, over which our perception changes.

- A: see individual leaves with sharp edge/boundary (occlusion model)
- B: see leaves but blurry edge (additive model)
- C: see a texture impression (MRF)
- D: see constant area (iid Gaussian)

Three types of changes over scales

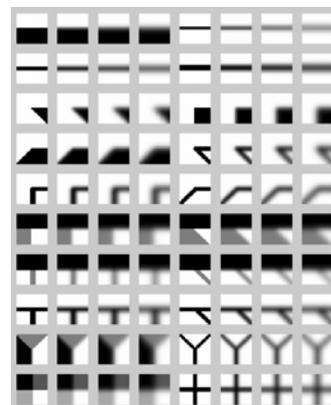
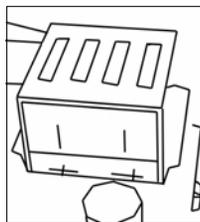
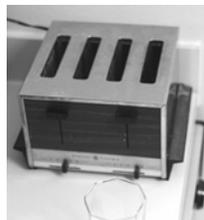
1. Catastrophic (texture to texton explosion),
2. graph grammatical splitting,

3. boundary sharpening

image I

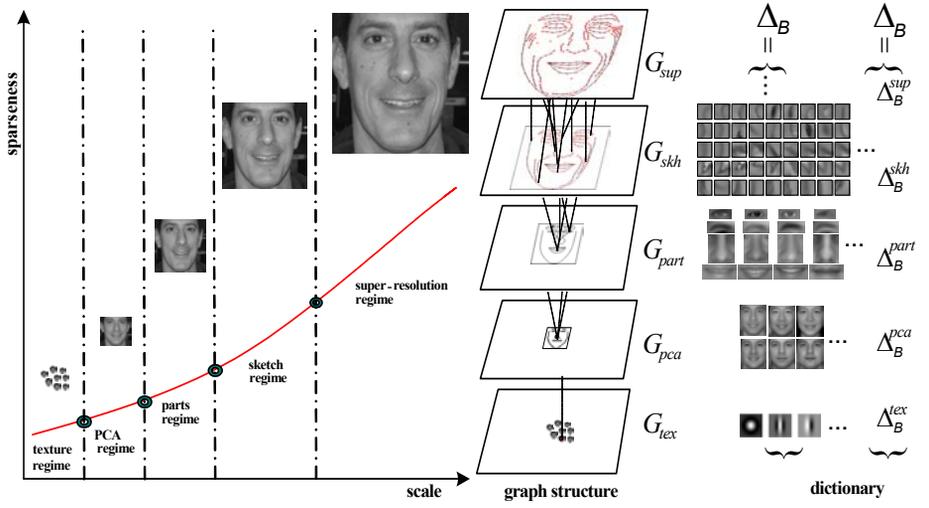


primal sketch G



(Guo, Zhu and Wu, 2004)

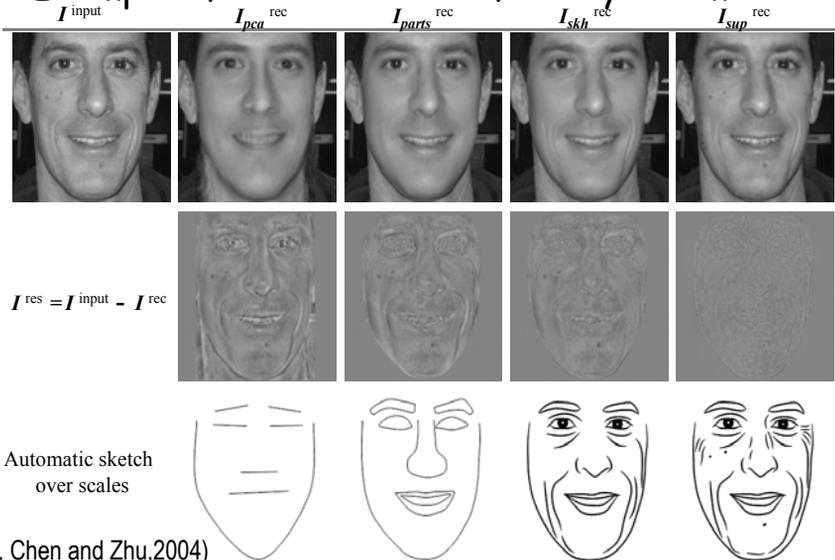
"Scale - Regime" diagram (Xu, Chen and Zhu,2004)



MSRI , January, 2005

Song-Chun Zhu

Example of reconstructed face by our model

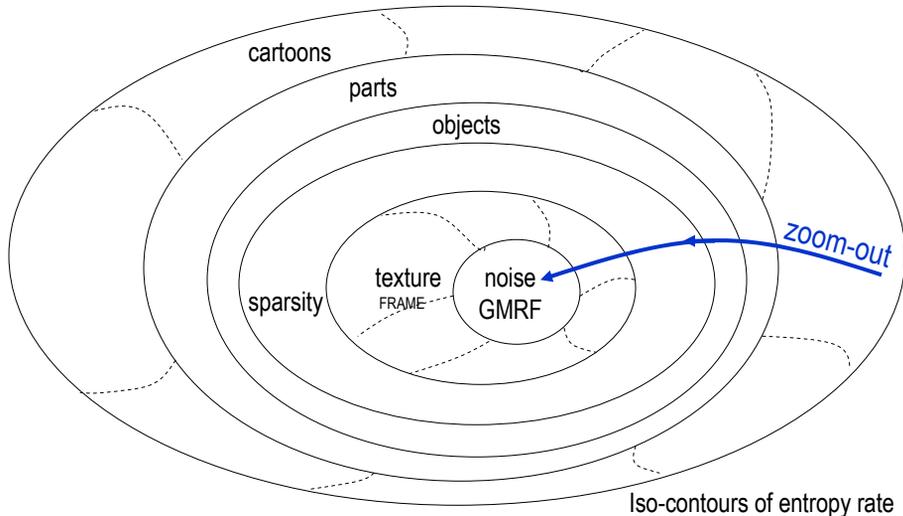


(Xu, Chen and Zhu,2004)

MSRI , January, 2005

Song-Chun Zhu

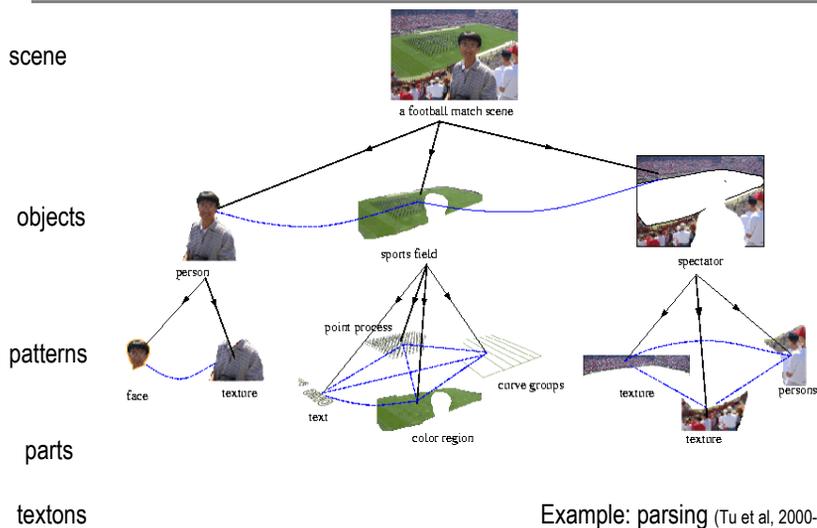
Mapping the Image Universes at different entropy rate



MSRI, January, 2005

Song-Chun Zhu

Part 2: Generic Images parsing



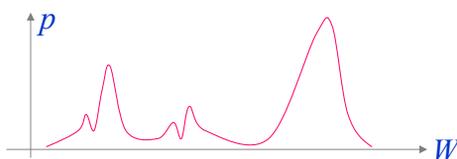
MSRI, January, 2005

Song-Chun Zhu

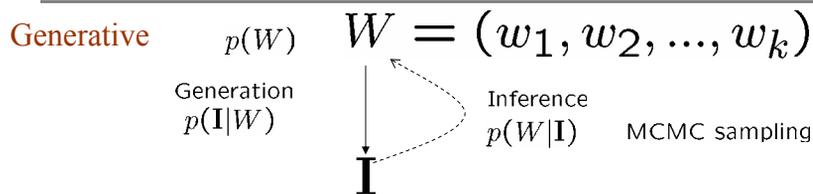
Part 2: statistical Computing

In statistics, we need to sample from the joint posterior probability.

$$(W_1, W_2, \dots, W_k) \sim p(W | I) \text{ or } p(I | W)p(W)$$

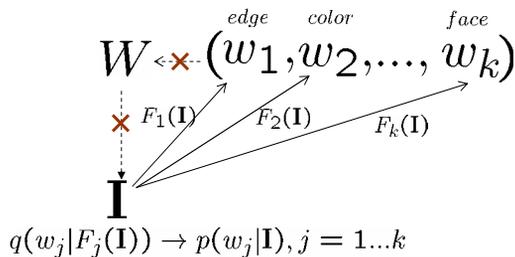


Generative vs. Discriminative Algorithms

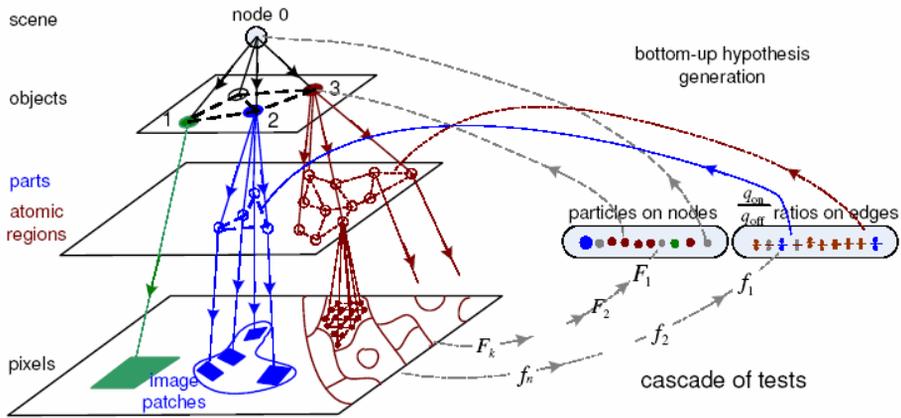


$$W^* = \arg \max p(W|I) = \arg \max p(I|W)p(W)$$

Discriminative



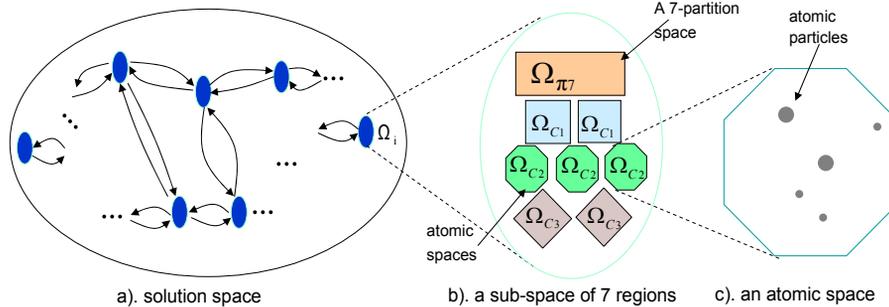
Integrating generative and discriminative



MSRI, January, 2005

Song-Chun Zhu

The Search/state Space



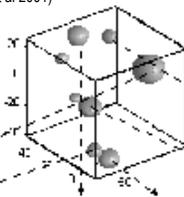
MSRI, January, 2005

Song-Chun Zhu

Example: Clustering in Color Space

Using Mean-shift clustering (Cheng, 1995, Meer et al 2001)

$$q(\theta|\mathbf{I}) = \sum_{i=1}^K \omega_i g(\theta - \theta_i)$$



Input



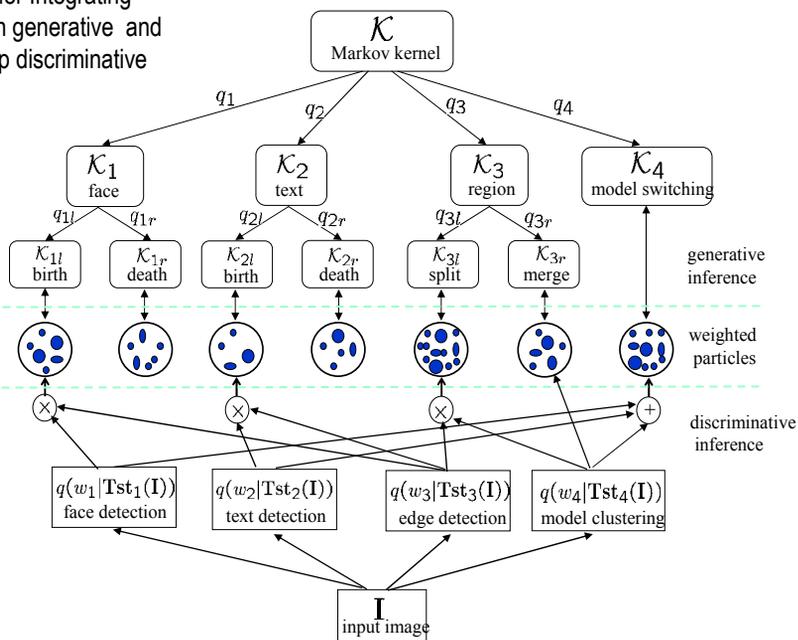
saliency maps 1 2 3 4 5 6

The brightness represents how likely a pixel belongs to a cluster.

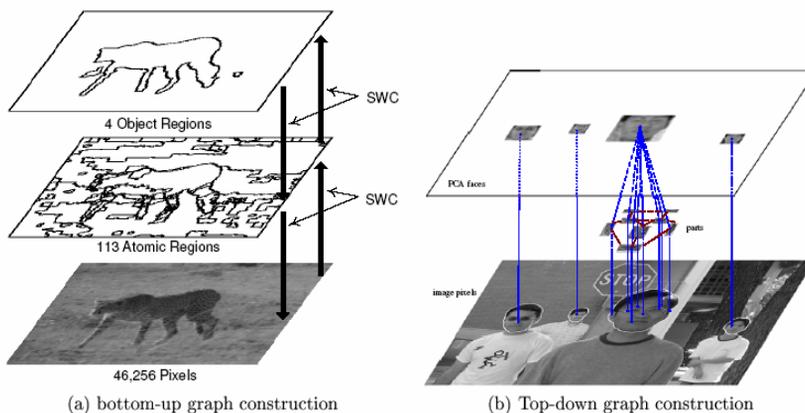
MSRI, January, 2005

Song-Chun Zhu

Diagram for Integrating
Top-down generative and
Bottom-up discriminative
Methods.



Two Computing Mechanisms



Alternating Bottom-up and Top-Down

Measuring the power of a discriminative Test

$$\begin{aligned} \delta(w|F_+) &= KL(p(w|\mathbf{I})||q(w|Tst_t(\mathbf{I}))) - KL(p(w|\mathbf{I})||q(w|Tst_t(\mathbf{I}), F_+)) \\ &= MI(w|Tst_t(\mathbf{I}, F_+)) - MI(w|Tst_t(\mathbf{I})) = KL(q(w|Tst_t(\mathbf{I}), F_+)||q(w|Tst_t(\mathbf{I}))) \end{aligned}$$

Measuring the power of sub-kernels

$$\begin{aligned} W_t \sim \mu_t(W) &= \nu(W_0) \circ K_{a(1)} \circ K_{a(2)} \circ \dots \circ K_{a(t)} \\ \delta_{a(t)} &\stackrel{def}{=} KL(p(W|\mathbf{I})||\mu_t(W)) - KL(p(W|\mathbf{I})||\mu_{t+1}(W)) = KL(K_{a(t)}(W_t|W_{t+1})||p_{MC}(W_t|W_{t+1})) \end{aligned}$$

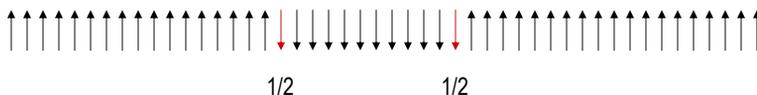
Gibbs sampler with Ising / Potts models

$$p(I) = \frac{1}{Z} \exp\{-\beta \sum_{\langle s,t \rangle} 1(I_s = I_t)\} = \frac{1}{Z} \prod_{\langle s,t \rangle} \exp\{-\beta \cdot 1(I_s = I_t)\}, \quad \langle s,t \rangle \in E_o$$

For example, in a 1D string of spins, suppose we use a Gibbs sampler to flip one spin at a time. It has a $p=1/2$ probability for flipping the spin at the boundary. Flipping a string of length n will need on average

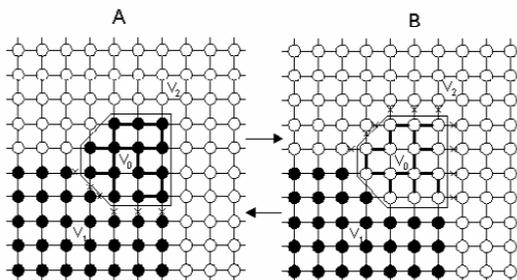
$$t = 1/p^n = 2^n \text{ steps!}$$

This is exponential waiting time.



SW with Ising / Potts models

Swendsen-Wang (1987) is an extremely smart idea that flips a patch at a time. There are multiple interpretations. We explain it from the Metropolis-Hastings method.



Each edge in the lattice $e=\langle s,t \rangle$ is associated with a constant probability q .

If s and t have different labels at the current state, e is turned off.

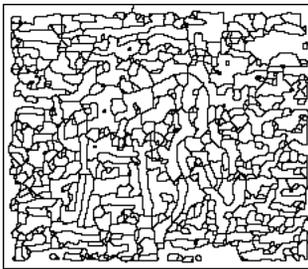
If s and t have the same label, e is turned off with probability q .

Thus each object is broken into a number of connected components (subgraph).

Convergence comparison: in sweep#

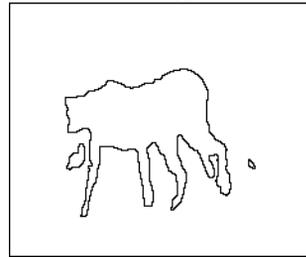


Input image



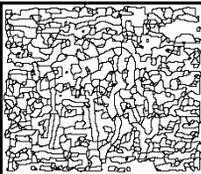
Atomic regions

Segment results



MSRI, January, 2005

Song-Chun Zhu



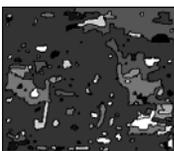
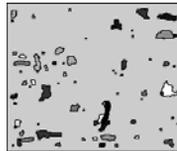
Graph partition/clustering: Using the discriminative model in the partition space

T=1

T=2

T=4

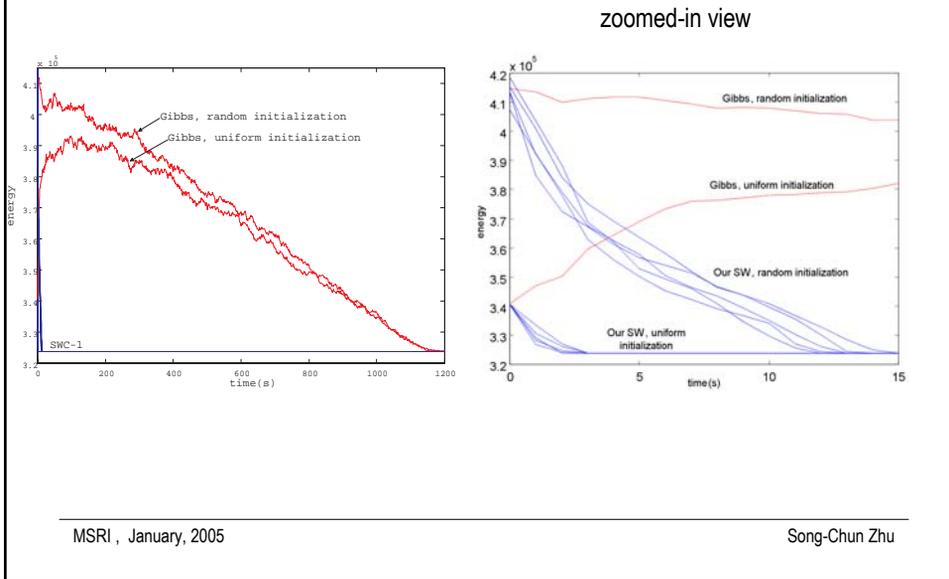
T=8



MSRI, January, 2005

Song-Chun Zhu

Convergence comparison: in cpu time



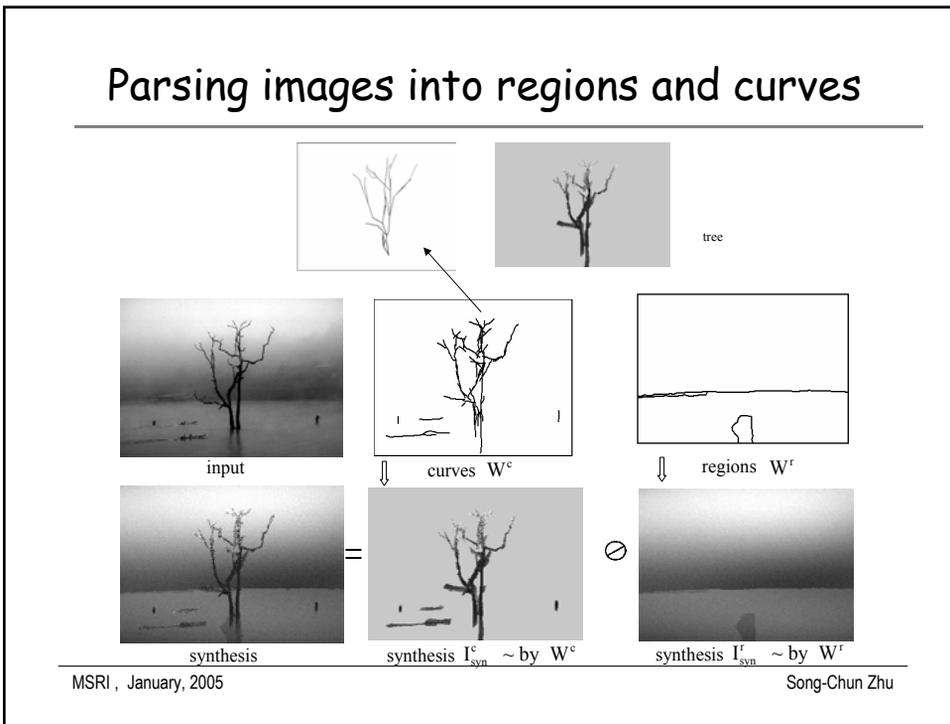
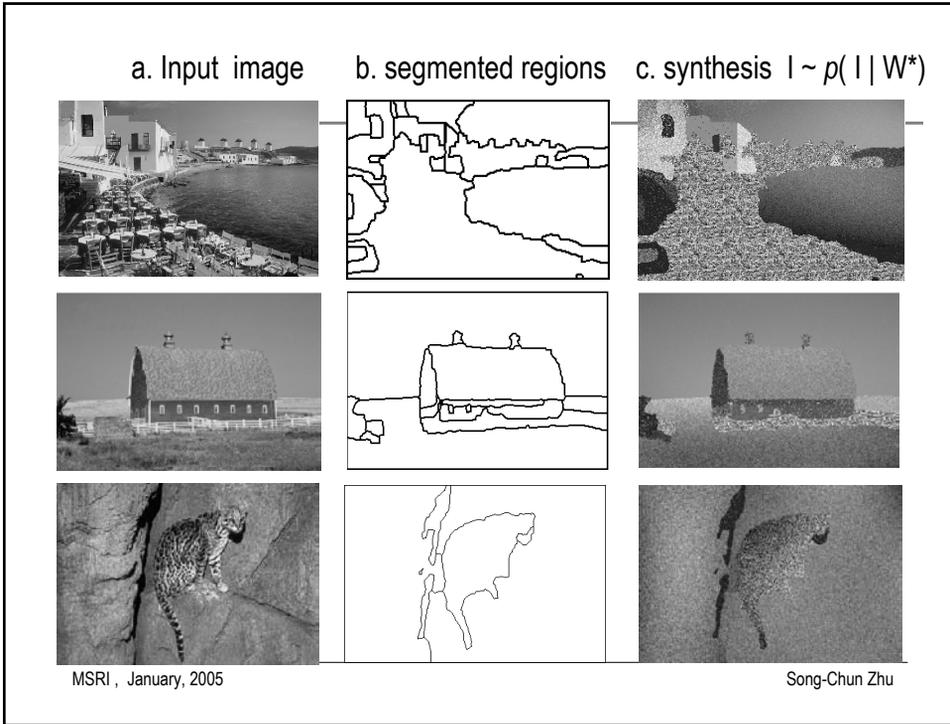
The Berkeley Benchmark Study

(David Martin et al, 2001)

test images	DDMCMC	manual segment	"error" measure
			0.1083
			0.3082
			0.5627

MSRI , January, 2005

Song-Chun Zhu



from image parsing to 3D

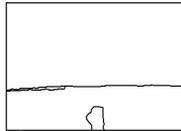
Example I: 3D reconstruction from a Single Image (Han and Zhu, 2003)



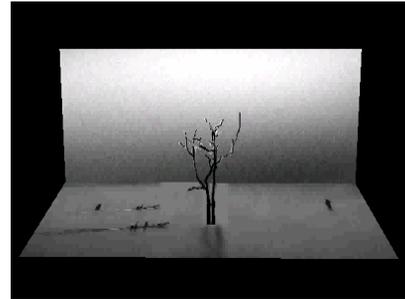
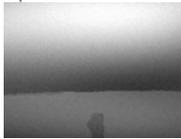
input I



curve & tree layer



region layer



3D reconstruction and rendering

from image parsing to 3D

3D reconstruction (Han and Zhu, 2003)

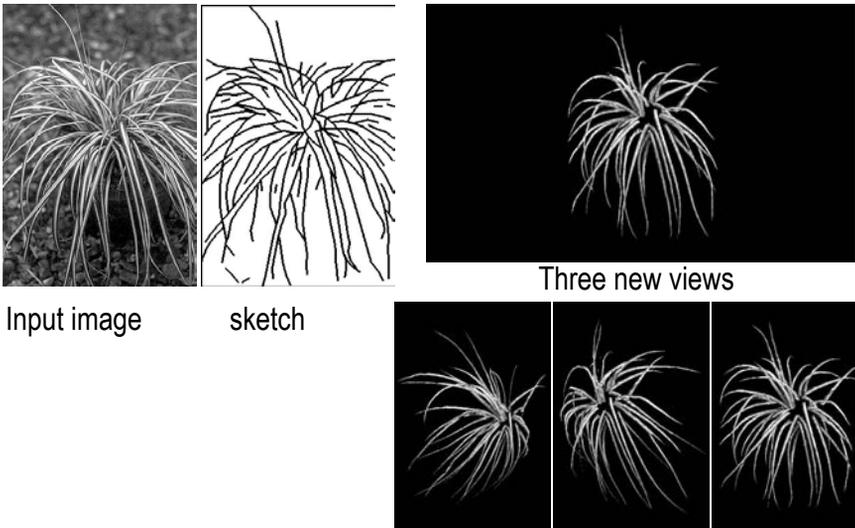


input image



3D reconstruction from a single image

from image parsing to 3D

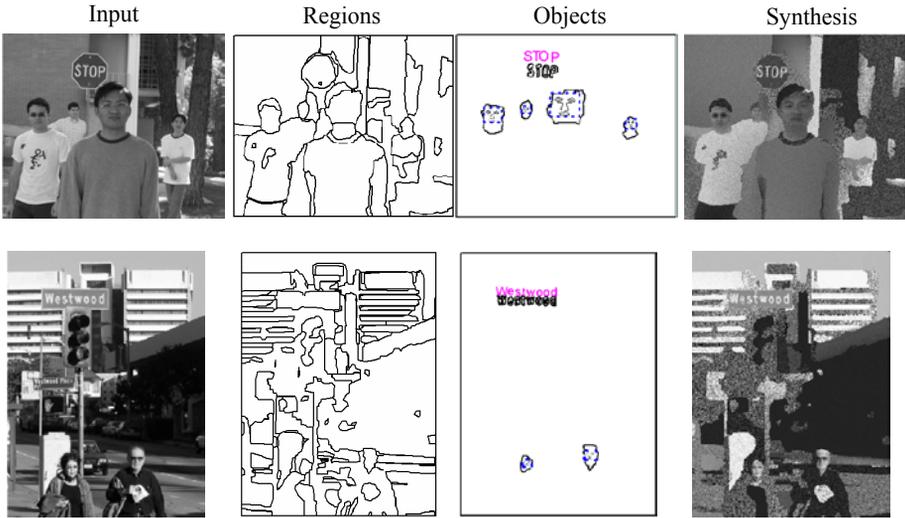


MSRI , January, 2005

Song-Chun Zhu

Image Parsing Results

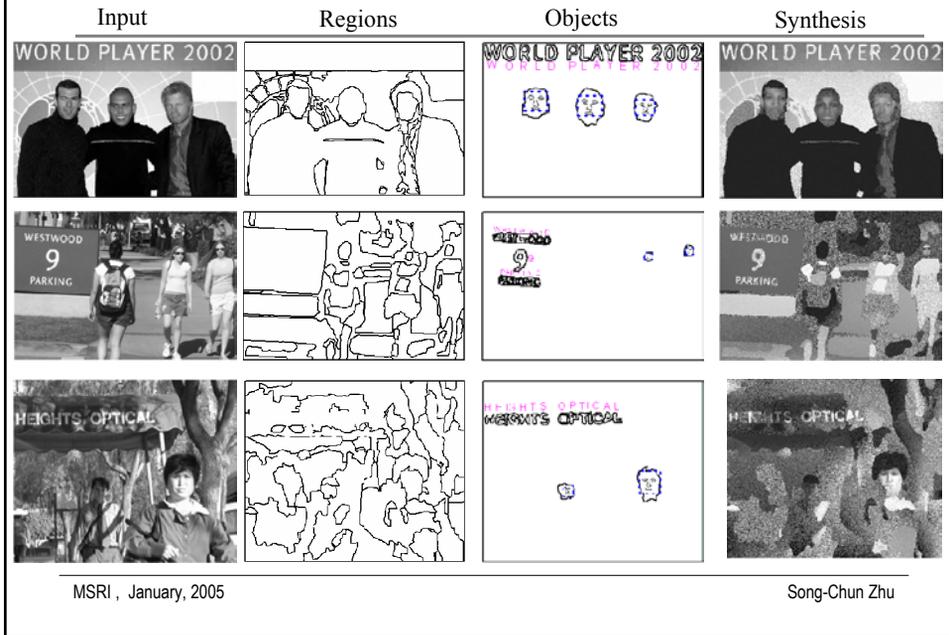
Tu, Chen, Yuille, and Zhu, iccv2003



MSRI , January, 2005

Song-Chun Zhu

Image Parsing Results



Summary: We see three methods and we need to work on the cracks between them

Generative methods

Hierarchic models,
Harmonic analysis/wavelet...

Heuristic Search,
MCMC

Descriptive methods

Markov fields/networks
Graphical models, Stat. Physics

Relaxation, Gibbs sampler,
Swendsen-Wang, Belief prop.

Discriminative methods

Clustering/detection,
Machine learning
Adaboosting