# Pursuing Explicit and Implicit Manifolds by Information Projection

Song-Chun Zhu

University of California, Los Angeles, USA

Lotus Hill Research Institute, China （中国莲花山研究院）

Joint work with Yingnian Wu, Kent Shi, ….

---

# 1, Background on visual (appearance) manifolds

Image patches from a single object category
are often found to form low dimensional manifolds.

e.g.  ISOMAP, LLE:
Saul and Roweis, 2000.

But, people found that image patches of generic
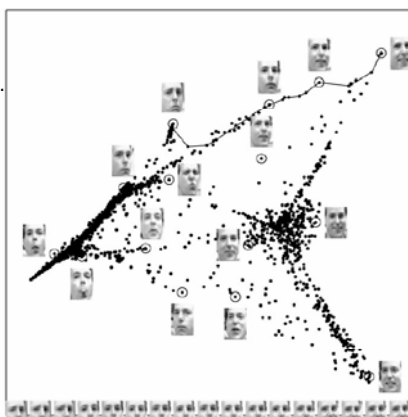natural images do not follow this observation.

Fig. 3. Images of faces (11) mapped into the embedding space described by the first two coordinates of LLE. Representative faces are shown next to circled points in different parts of the space. The bottom images correspond to points along the top-right path (linked by solid line), illustrating one particular mode of variability in pose and expression.

# Looking at local, generic natural image statistics
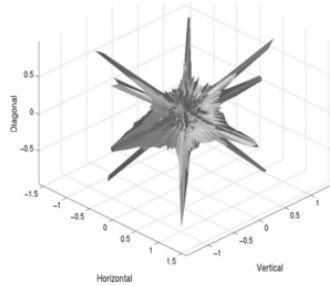
....
Ruderman and Bialek 87, 94
Fields 87, 94
Zhu and Mumford 95-96
Chi and Geman 97-98
Huang, 2000
Simoncelli etc 98-03
.....

Here is an example of how real world data can be truly complex – non-Gaussian and highly kurtotic. This is an iso-density contour for a 3D histogram of log(range) images (2x2 patches minus their means) (Brown range image database, thesis of James Huang)
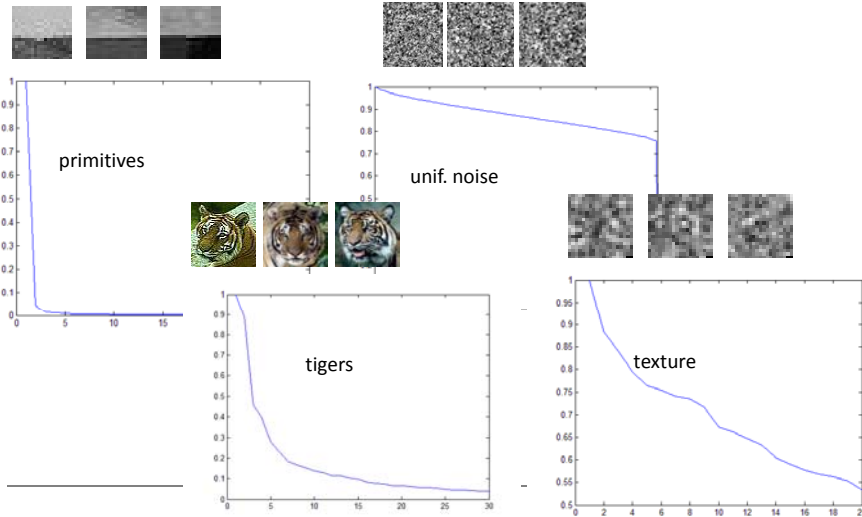
# A wide spectrum of categories from low to high entropy

| Edge | Bar | Two Parallel Lines | Cat | Dog | Lion | Tiger | Fur | Carpet | Grass | Noise |
|---|---|---|---|---|---|---|---|---|---|---|

Entropy ~ Dimension ~ Log volume( manifold )

# Visual manifolds have varying dimensions

Take 16x16 image patches (256-space), run PCA for each
category, and plot the eigen-values in decreasing order.



primitives

unif. noise

tigers

texture

# By analogy: pictures of our universe



entropy ( temperature ) regimes.

compositional structures.

How do we learn these manifolds?
Can we do it by K-mean clustering?
3 modeling theories in vision:
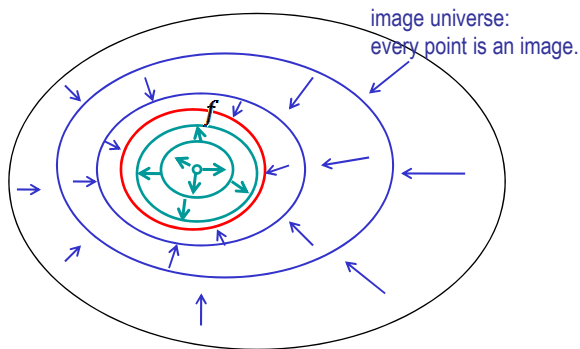    (1) Markov random fields, (2) Sparse coding, (3) Grammar and Composition

# 2, Manifold pursuit in the universe of image patches

$f$ : target distribution;    $p$: our model;    $q$: initial model

$$q = p_0 \rightarrow p_1 \rightarrow \cdots \rightarrow p_k \quad to \quad f$$

image universe:
every point is an image.

1,  $q = unif()$

2,  $q = \delta()$

model ~ image set ~ manifold ~ cluster

---

# Intuitive idea:  a professor grading an exam

The full score (like dimension in our case) is 100. You have two ways:

For top students (high dimensional manifolds), you start from 100 and deduct points :

$$100 - 2 - 0 - 0 - 3 - 0 - 2 - 0 - 0 - 0 - 0 - 0 - 1 = 92$$

For bottom students (low dimensional manifolds), you start from 0 and add points
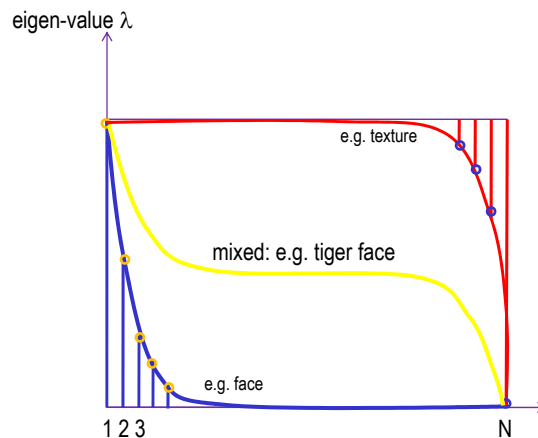
$$0 + 8 + 0 + 0 + 3 + 0 + 2 + 0 + 0 + 5 + 0 + 0 + 1 = 19$$

In reality, suppose the exam is very long (just like the large image has >1M pixels), a student may have mixed performance, e.g. doing excellent in the 1st half and doing poorly in the 2nd half. Thus a most effective way is to use the two methods for different sections of the exam.

$$(50 - 2 - 0 - 0 - 3 - 0) + (0 + 5 + 3 + 0 + 0 + 2) = 45 + 10 = 55$$

In fact, most of the object categories are middle entropy manifolds and have mixed structures.

# Manifold pursuit in the image universe

In a simple case: $f$ is a Gaussian distribution



# Manifold pursuit by information projection

Given only positive examples from a class c

$$\Omega_c^+ = \{I_i^{obs}; \quad i = 1, 2, \ldots, M^+\} \sim f(I)$$

We pursue a series of models $p$ to approach a underlying "true" probability $f$

$$q = p_0 \rightarrow p_1 \rightarrow \cdots \rightarrow p_k \quad to \quad f$$

At each step, we augment the current model $p$ to a new model $p_+$

$$h_+^* = \arg\max \ KL(f \mid p) - KL(f \mid p_+)$$
$$= \arg\max KL \ (p_+ \mid p)$$

Subject to a projection constraint:

$$E_{p_+}[h_+(I)] = E_f[h_+(I)] \cong \bar{h}_+$$

$h_+(I)$ is a feature statistics of image $I$

# Manifold pursuit by information projection

Solving the constrained optimization problem leads to the Euler-Lagrange equation

$$p_+^* = \arg\min \int p_+(I) \log \frac{p_+(I)}{p(I)} \, dI + \lambda_+ [\int p_+(I) h_+(I) dI - \bar{h}_+] + \lambda_o \int p_+(I) dI - 1]$$

$$p_k(I; \Theta_k) = \frac{1}{z_{+,k}} p_{k-1}(I; \Theta_{k-1}) e^{-\lambda_k h_k(I)}$$

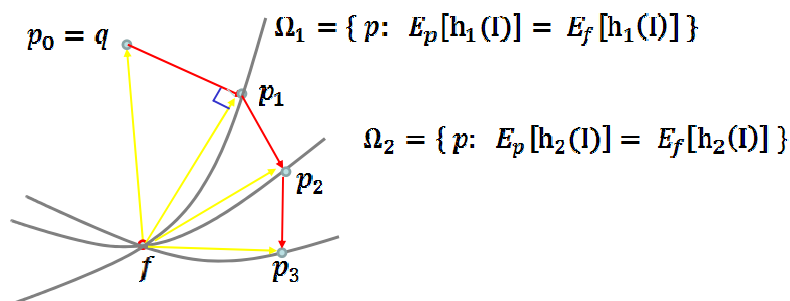$$= \frac{1}{z_k} q(I) \exp \{ -\sum_{i=1}^{k} \lambda_l h_l(I) \}$$

where

$$Z_k = Z_{+,1} \cdot Z_{+,2} \cdots Z_{+,k} \qquad \Theta_k = (\lambda_1, \lambda_2, \cdots \lambda_k)$$

For q being a uniform distribution, we have $\qquad q(I) = \frac{1}{z_o}$

---

# Information projection

DellaPietra, DellaPietra, Lafferty, 97
Zhu, Wu, Mumford, 97



$p_0 = q$

$$\Omega_1 = \{ p: \ E_p[h_1(I)] = E_f[h_1(I)] \}$$

$p_1$

$$\Omega_2 = \{ p: \ E_p[h_2(I)] = E_f[h_2(I)] \}$$

$p_2$

$f \qquad p_3$

$$KL(f \mid p) = KL(f \mid p_+) + KL(p_+ \mid p)$$

So the KL-divergence decreases monotonically.

# A Maximin Learning Principle

max-step: choosing a distinct feature and statistics

$$h_+^* = \arg\max KL\ (p_+ \mid p)$$

min-step: given the selected feature constraint, computing the parameter
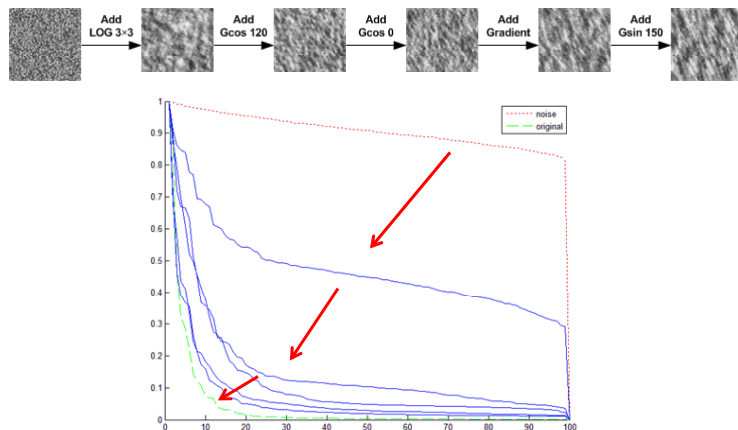
$$\lambda_+^* = \arg\min KL(p_+ \mid p)$$

Claim: this learning procedure unifies almost all we know in visual modeling
PCA, sparse coding,
MRF, Gibbs, FRAME,
Adaboost (when h() is binary),
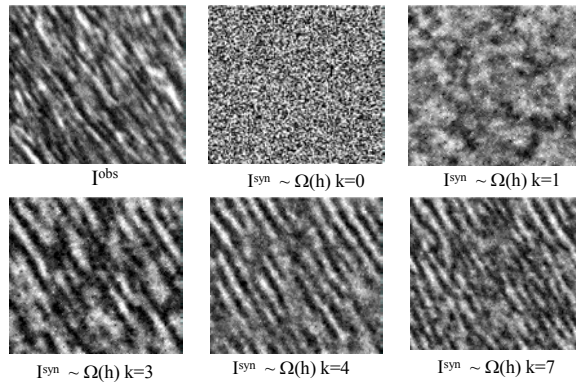Stochastic grammar

---

# 3, Case studies:

Case 1: Pursuing texture models by compression from white noise

# A texture pattern is an "implicit manifold"

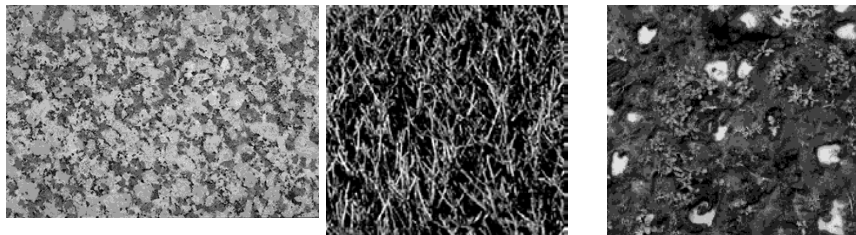$$\text{a texture} = \Omega(h_c) = \{ I : \ h_i(I) = h_{c,i} \ , i = 1,2,...,K \}$$

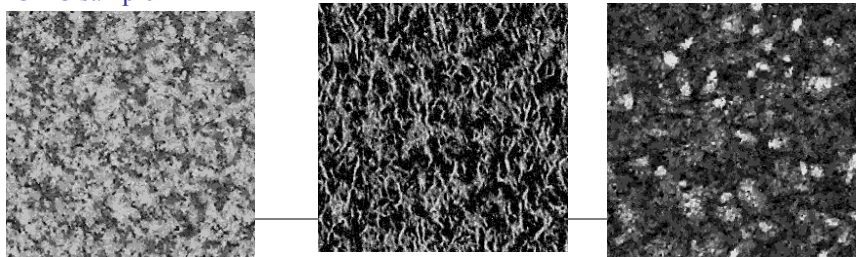$H_c$ are histograms of Gabor filters, i.e. marginal distributions of $f(I)$



$I^{obs}$      $I^{syn} \sim \Omega(h)$ k=0      $I^{syn} \sim \Omega(h)$ k=1

$I^{syn} \sim \Omega(h)$ k=3      $I^{syn} \sim \Omega(h)$ k=4      $I^{syn} \sim \Omega(h)$ k=7

(Zhu,Wu, Mumford 97,99,00)

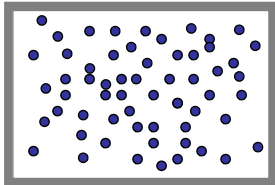# More examples of the texture manifold (implicit)

Observed



MCMC sample

## This is originally from statistical physics !

Statistical physics studies macroscopic properties of systems
that consist of massive elements with microscopic interactions.

e.g.: a tank of insulated gas or ferro-magnetic material

$N = 10^{23}$

A state of the system is specified by the position of the
N elements $X^N$ and their momenta $p^N$

$$S = (x^N, \ p^N)$$

But we only care about some global properties
Energy E, Volume V, Pressure, ....

Micro-canonical Ensemble

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

---

## Equivalence of Julesz ensemble and FRAME/MRF models

$\Lambda$

Zhu, Wu, Mumford, 1997
Wu and Zhu, 1999

Theorem 1
   For a very large image from the Julesz ensemble $I \sim f(I; h_c)$ any
   local patch of the image $I_\Lambda$ given its neighborhood follows a conditional
   distribution specified by a FRAME model $p(I_\Lambda \,|\, I_{\partial\Lambda} : \beta)$

Theorem 2
   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.

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

# Case 2: A car pattern is an "explicit manifold"

## Learning active basis as deformable template

A basis is an image space spanned by a number of vectors (e.g. Gabor/primitives)

$$B = (B_1, B_2, \ldots, B_k)$$

$$A\ car = \Omega = \{I: \ I = \sum_i \gamma_i \, B_{i,\delta} \}$$
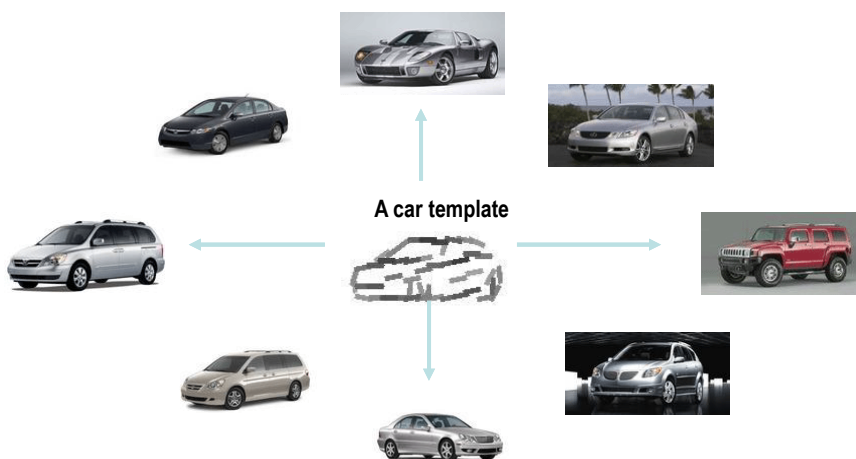
A car template



(Gabor elements represented by bar)

An incoming car image:



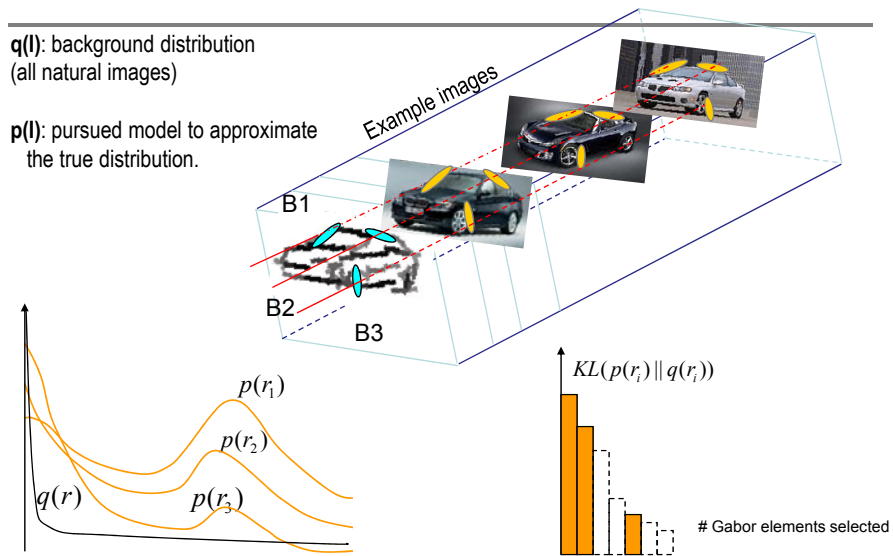With slight modification, this model can handle multi-views        Wu, Si. Gong, Zhu, 2008

---

# Deformed to fit many car instances



**A car template**

# Pursuing the active basis model (explicit manifold)

**q(I)**: background distribution
(all natural images)

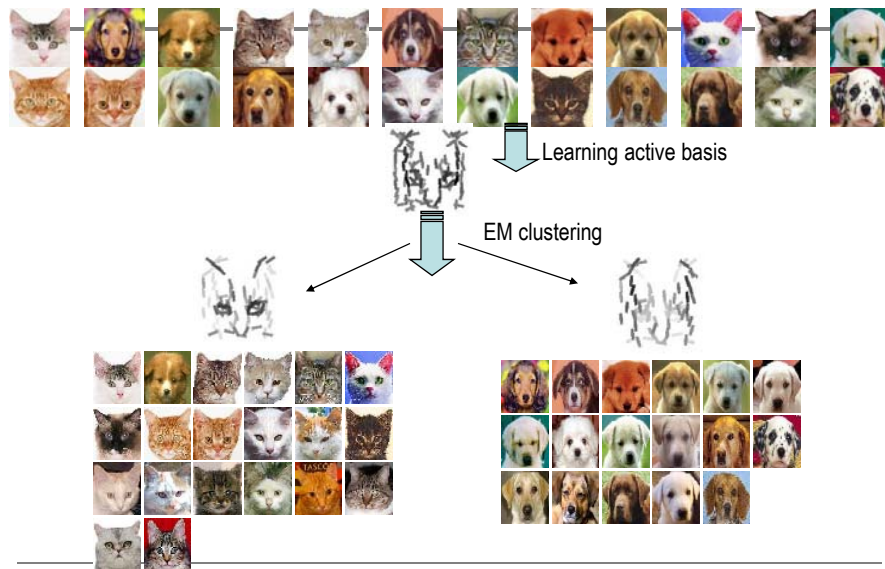**p(I)**: pursued model to approximate
the true distribution.

Example images

B1

B2

B3



$p(r_1)$

$p(r_2)$

$q(r)$    $p(r_3)$

$KL(p(r_i) \| q(r_i))$

\# Gabor elements selected

---

# A running example



A car template consisting of
48 Gabor elements

Car instances

11

# Experiment : learning and clustering



Learning active basis

EM clustering

# Experiment : learning and detection



AUC = 0.929

AUC = 0.860

true positive rate

false positive rate

active basis
adaboost

Wu,Si,Fleming,Zhu,07

vs: Viola, Jones, 04

## Template detection experiment



Wu,Si,Fleming,Zhu,07

## Summary: two pure manifolds

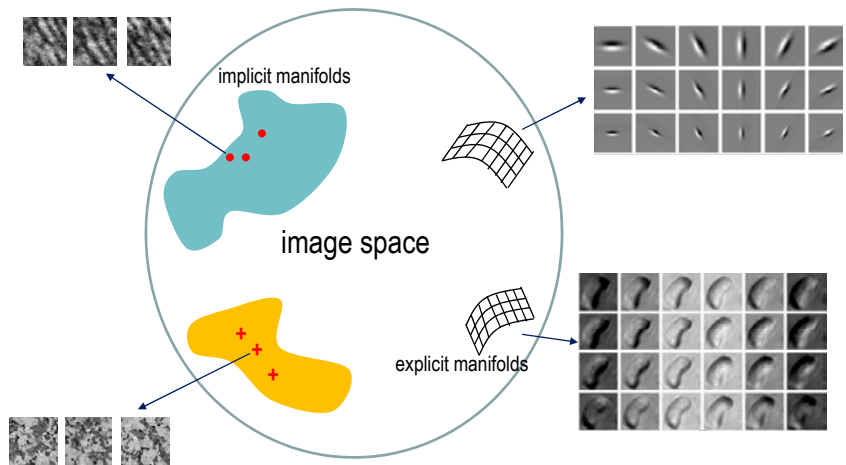*implicit* vs. *explicit*

$$\Omega = \{I: \quad h(I) = h_o \}$$

h(I) is some image feature/statistics

$$\Omega = \{ I: \quad I = g(w; \Delta) \}$$

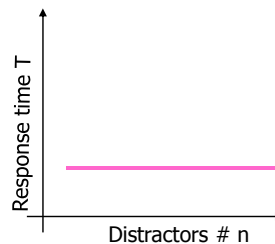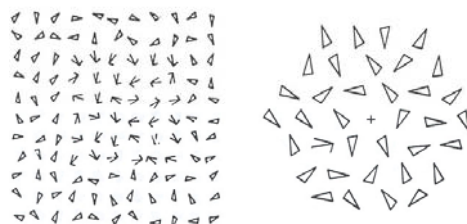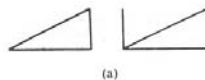g is a generation function,
w is intrinsic dimension
Δ is a dictionary

# Summary: a second look at the space of image patches



implicit manifolds

image space

explicit manifolds

# 4, Relations to the literature: psychophysics

(1) textures vs textons   (Julesz, 60-70s)

textons



(a)

Response time T

Distractors # n

# Textons vs. Textures

textures



Response time T

Distractors # n
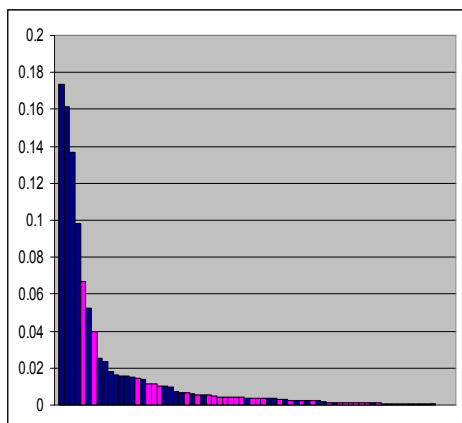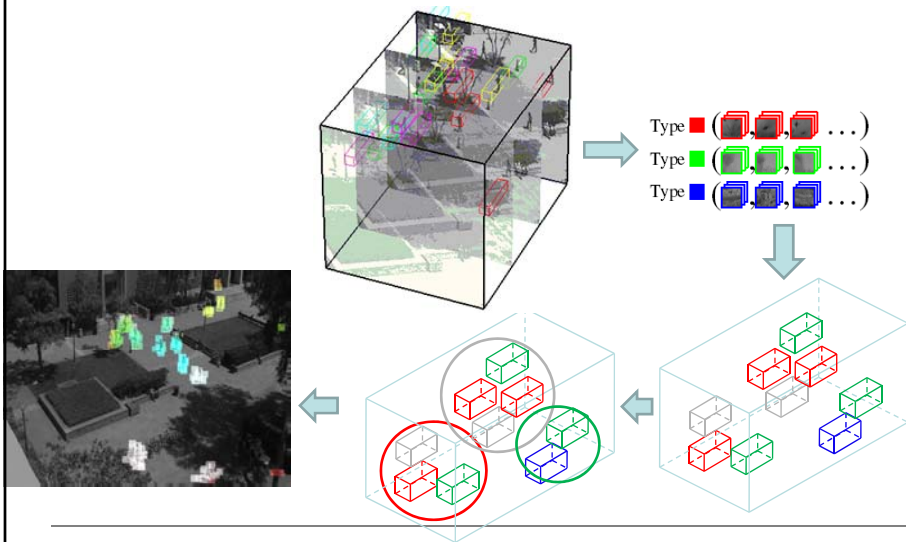
---

# Frequency plot of the ex/implicit manifolds in natural images

implicit texture clusters (blue),
explicit primitive clusters (pink).



| | cluster centers | | instances in each cluster | | | | | |
|---|---|---|---|---|---|---|---|---|
| 1 | | | | | | | | sky, wall, floor |
| 2 | | | | | | | | dry wall, ceiling |
| 3 | | | | | | | | carpet, ceiling, thick clouds |
| 4 | | | | | | | | step edge |
| 5 | | | | | | | | concrete floor, wood wall |
| 6 | | | | | | | | L-junction |
| 7 | | | | | | | | ridge/bar |
| 8 | | | | | | | | carpet, wall |
| 9 | | | | | | | | L-junction centered at 165° |
| 10 | | | | | | | | water |
| 11 | | | | | | | | lawn grass |
| 12 | | | | | | | | terminator |
| 13 | | | | | | | | wild grass, roof |
| 14 | | | | | | | | L-junction at 130° |
| 15 | | | | | | | | plants from far distance |
| 16 | | | | | | | | sand |
| 17 | | | | | | | | close-up of concrete |
| 18 | | | | | | | | wood grain |
| 19 | | | | | | | | T-junction at 90° |
| 20 | | | | | | | | Y-junction |

## Clustering in video



## Examples in video

| explicit | implicit |
| --- | --- |

# 6, Primal sketch: integrating the two regimes



org image     sketching pursuit process     sketches

syn image     synthesized textures     sketch image
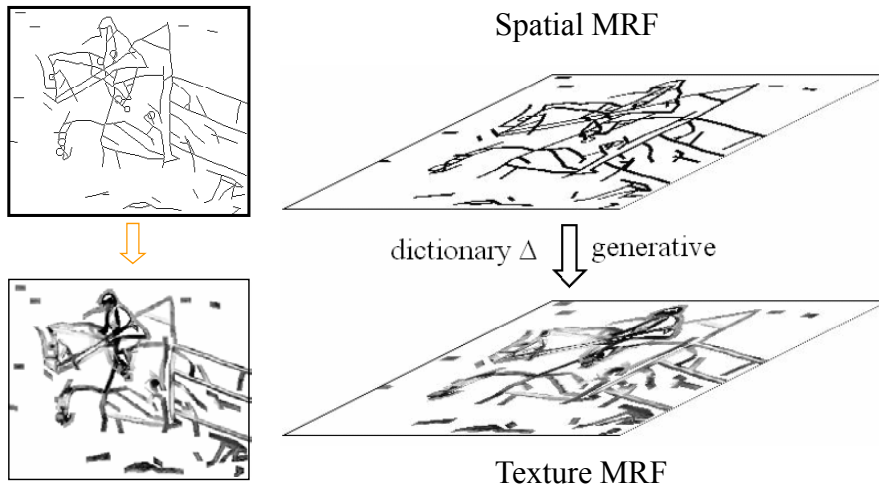
(Guo,Zhu,Wu, 2003-05)
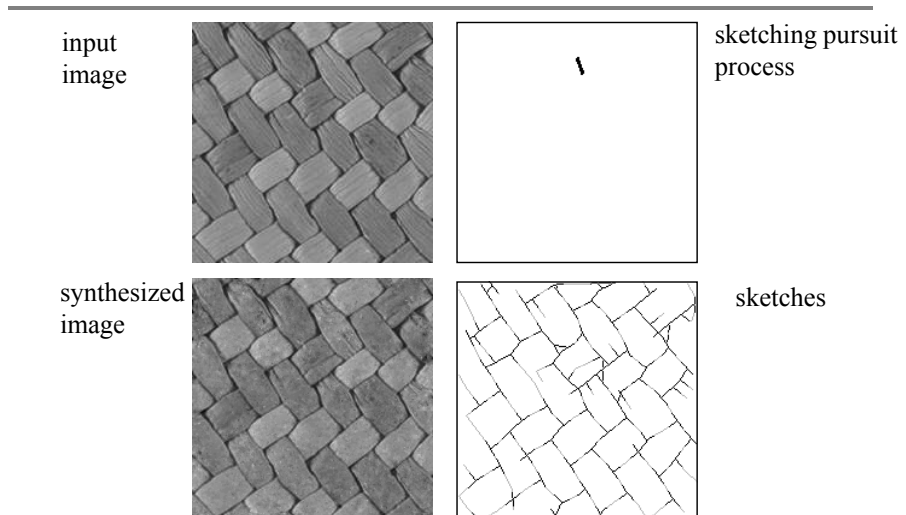
# manifolds of image primitives

Learned texton/primitive dictionary with some landmarks that transform and warp the patches



(a)          (b)
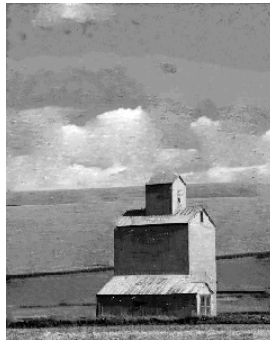
# Primal Sketch is a two-level MRF model



Spatial MRF

dictionary $\Delta$  generative

Texture MRF

# Primal sketch example

input
image

sketching pursuit
process

synthesized
image

sketches

## Primal sketch example



original image      synthesized image      sketching pursuit process

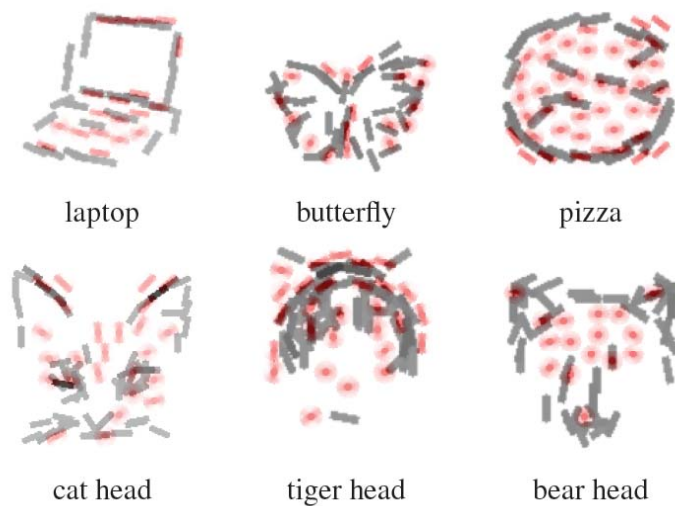## 7, deformable template: mixing the im/explicit manifolds



space of orientation histograms
(a projection space)

$H(I) = \mathbf{h} + \varepsilon$

texture patch

space of image patches

sketch patch

$I = cB + \varepsilon$

(a)

Si et al 2008

## The two types of models compete in learning the templates



head and shoulder

hedgehog

water patches

## Some examples of learn object categories



laptop

butterfly

pizza

cat head

tiger head

bear head

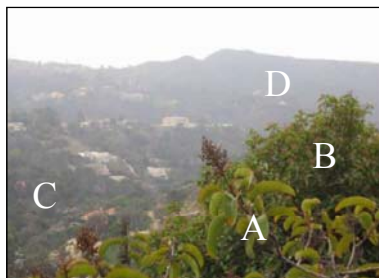# 8, Information scaling leads to manifold transitions !



Scaling (zoom-out) increases the image entropy (dimensions)



Wu, Zhu, Guo, 04,07

# Transition of the manifolds through info. scaling

How are these manifolds related to each other ?



perceptual scale space theory (Wang and Zhu 2005)

# Summary: understanding the "ingredients of our herbs" !

2 type manifolds, pursuit, integration, mixing, and transition