

# Statistical Modeling and Conceptualization of Visual Patterns

Song-Chun Zhu

**Abstract**—Natural images contain an overwhelming number of visual patterns generated by diverse stochastic processes. Defining and modeling these patterns is of fundamental importance for generic vision tasks, such as perceptual organization, segmentation, and recognition. The objective of this epistemological paper is to summarize various threads of research in the literature and to pursue a unified framework for conceptualization, modeling, learning, and computing visual patterns. This paper starts with reviewing four research streams: 1) the study of image statistics, 2) the analysis of image components, 3) the grouping of image elements, and 4) the modeling of visual patterns. The models from these research streams are then divided into four categories according to their semantic structures: 1) descriptive models, i.e., Markov random fields (MRF) or Gibbs, 2) variants of descriptive models (causal MRF and “pseudodescriptive” models), 3) generative models, and 4) discriminative models. The objectives, principles, theories, and typical models are reviewed in each category and the relationships between the four types of models are studied. Two central themes emerge from the relationship studies. 1) In representation, the integration of descriptive and generative models is the future direction for statistical modeling and should lead to richer and more advanced classes of vision models. 2) To make visual models computationally tractable, discriminative models are used as computational heuristics for inferring generative models. Thus, the roles of four types of models are clarified. The paper also addresses the issue of conceptualizing visual patterns and their components (vocabularies) from the perspective of statistical mechanics. Under this unified framework, a visual pattern is equalized to a statistical ensemble, and, furthermore, statistical models for various visual patterns form a “continuous” spectrum in the sense that they belong to a series of nested probability families in the space of attributed graphs.

**Index Terms**—Perceptual organization, descriptive models, generative models, causal Markov models, discriminative methods, minimax entropy learning, mixed Markov models.



## 1 INTRODUCTION

### 1.1 Quest for a Common Framework of Visual Knowledge Representation

NATURAL images consist of an overwhelming number of visual patterns generated by very diverse stochastic processes in nature. The objective of image analysis is to parse generic images into their constituent patterns. For example, Fig. 1a shows an image of a football scene which is parsed into: Fig. 1b a point process for the music band, Fig. 1c a line and curve process for the field marks, Fig. 1d a uniform region for the ground, Fig. 1e two texture regions for the spectators, and Fig. 1f two objects—words and human face. Depending on the types of patterns that a task is interested in, the image parsing problem is respectively called 1) *perceptual grouping* for point, line, and curve processes, 2) *image segmentation* for region process, and 3) *object recognition* for high level objects. In other words, grouping, segmentation, and recognition are subtasks of the image parsing problem and, thus, they ought to be solved in a unified way. This requests a common and mathematically sound framework for representing visual knowledge, and the visual knowledge includes two parts.

1. Mathematical definitions and models of various visual patterns.
2. Computational heuristics for effective inference of the visual patterns and models.

• The author is with the Departments of Statistics and Computer Science, University of California, 8130 Math Sciences Bldg., Box 951554, Los Angeles, Los Angeles, CA 90095. E-mail: sczhu@stat.ucla.edu.

Manuscript received 1 Jan. 2002; revised 2 Aug. 2002; accepted 29 Jan. 2003. Recommended for acceptance by D. Jacobs and M. Lindenbaum. For information on obtaining reprints of this article, please send e-mail to: tpami@computer.org, and reference IEEECS Log Number 118002.

The objective of this epistemological paper is to pursue such a unified framework. More specifically, it should address the following four problems.

**Conceptualization of visual patterns.** What is a quantitative definition for a visual pattern? For example, what is a “texture” and what is a “human face?” The concept of a pattern is an abstraction of some properties decided by certain “vision purposes.” These properties are feature statistics computing from either raw signals or some hidden descriptions inferred from raw signals. In both ways, a visual pattern is equalized to a set of observable signals governed by a statistical model—which we call an ensemble. In other words, each instance in the set is assigned a probability. For homogeneous patterns, such as texture, on large lattice this probability is a uniform distribution and the visual pattern is an equivalence class of images that satisfy certain descriptions. The paper should review some theoretical background in statistical mechanics and typical physics ensembles, from which a consistent framework is derived for defining various visual patterns.

**Statistical modeling of visual patterns.** First of all, why are the statistical models much needed in vision? In other words, what is the origin of these models? Some argued that probabilities are involved because of noise and distortion in images. This is truly a misunderstanding! With high quality digital cameras, there is rarely noise or distortion in images anymore. Probabilities are associated with the definitions of patterns and are even derived from deterministic definitions. In fact, the statistical models are intrinsic representation of visual knowledge and image regularities. Second, what are the mathematical space for patterns and models? Patterns are represented by attributed graphs and, thus, models are defined in the space of attributed graphs. The paper should

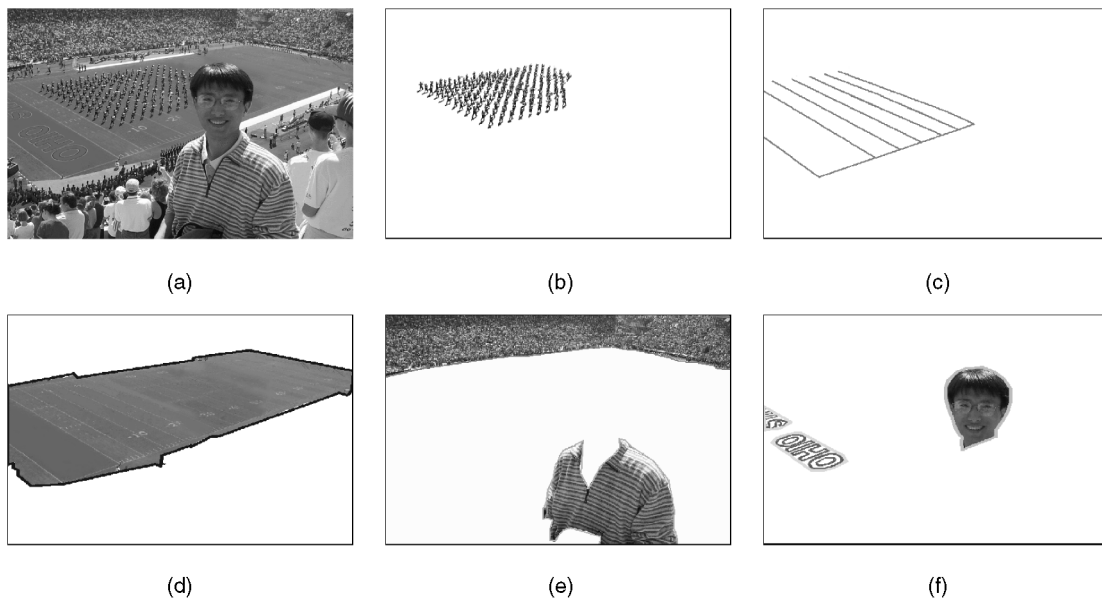


Fig. 1. Parsing an image into its constituent patterns. (a) An input image. (b) A point process. (c) A line/curve process. (d) A uniform region. (e) Two texture regions. (f) Objects: face and words. Courtesy of Tu and Zhu [83].

review two classes of models. One is descriptive model that are Markov random fields (or Gibbs) and its variants (including causal Markov models). The other is generative models which engage hidden variables for generating images in a top-down manner. It is shown that the two classes of models should be integrated. In the literature, a generative model often has a trivial descriptive component and a descriptive model usually has a trivial generative component. As a result of this integration, the models for various visual patterns, ranging from textures to geometric shapes, should form a “continuous spectrum” in the sense that they are from a series of nested probability families in this space.

**Learning a visual vocabulary.** What is the hierarchy of visual descriptions for general visual patterns? Can this vocabulary of visual description be defined quantitatively and learned from the ensemble of natural images? Compared with the large vocabulary in speech and language (such as phonemes, words, phrases, and sentences), and the rich structures in physics (such as electrons, atoms, molecules, and polymers), the current visual vocabulary is far from being enough for visual pattern representation. This paper reviews some progress in learning image bases and textons as visual dictionaries. These dictionaries are associated with generative models as parameters and are learned from natural images through model fitting.

**Computational tractability.** Besides the representational knowledge (definitions, models, and vocabularies), there is also computational knowledge. The latter are computational heuristics for effective inference of visual patterns, i.e., inferring hidden variables from raw images. These heuristics are the discriminative models that are approximations to posterior probability or ratios of posterior probabilities. The approximative posteriors are computed through local image features, in contrast to the real posterior computed by the Bayes rule following generative models. Then, it is natural to ask what are the intrinsic relationships between representational and computational models? Generally speaking, the generative models are expressed as top-down probabilities and the hidden variables have to be inferred from posterior

probabilities following the Bayes rule, by Markov chain Monte Carlo techniques, in general, such as the Metropolis-Hastings method. In contrast, the discriminative models approximate the posterior in a bottom-up and speedy fashion. These discriminative probabilities are used as proposal probabilities that drive the Markov chain search for fast convergence and mixing.

The questions raised above have motivated long threads of research from many disciplines, for example, applied mathematics, statistics, computer vision, image coding, psychology, and computational neurosciences. Recently, a uniform mathematical framework emerges from the interactions between the research streams and, experimentally, a large number of visual patterns can be modeled realistically. This inspires the author to write an *epistemology* paper to summarize the progress in the field. The objective of the paper is to facilitate communications between different fields and provide a road map for the pursuit of a common mathematical theory for visual pattern representation and computation.

## 1.2 Plan of the Paper

The paper consists of the following five parts.

*Part 1 Literature survey.* The paper starts with a survey of the literature in Section 2 to set the background. We divide the literature in four research streams:

1. the study of natural image statistics,
2. the analysis of natural image components
3. the grouping of natural image elements, and
4. the modeling of visual patterns.

These streams develop four types of models:

1. descriptive model (Markov random fields or Gibbs),
2. variants of descriptive models (causal MRF and “pseudodescriptive” models),
3. generative models, and
4. discriminative model.

The relationships of the models will be studied.

*Part 2: A common framework for learning models.* Section 3 presents a common maximum-likelihood formulation for modeling visual patterns. Then, it leads to the choice of two families of the probability models: descriptive models (and its variants) and generative models. Then, the paper presents the descriptive and generative models in parallel.

*Part 3: Descriptive models and its variants.* This includes two sections. First, the paper presents, in Section 4, the basic assumptions and the minimax entropy principles for learning descriptive models and seven typical examples from low-level image pattern to high-level human face patterns in the literature. Second, in Section 8, the paper discusses a few variants to the descriptive models, including causal Markov models and the pseudodescriptive models.

*Part 4: Generative models.* In parallel, Section 6 presents the basic assumptions, methods, and five typical examples for learning generative models.

*Part 5: Conceptualization of visual patterns.* This includes two sections. First, in Section 5.2, it addresses the issue of conceptualization from the perspective of descriptive models. It presents the statistical physics foundation of descriptive models and three types of ensembles: the microcanonical, canonical, and grand-canonical ensembles. Then, it conceptualizes a visual pattern to an ensemble of physical states. In Section 7, the paper revisits the conceptualization of patterns from the perspectives of generative models and states that the visual vocabulary can be learned as parameters in the generative models.

*Part 6: Discriminative models.* Then, the paper turns to computational issues in Section 9. It reviews how discriminative models can be used for inferring hidden structures in generative models and presents maximum mutual information principle for selecting informative features for discriminations.

Finally, Section 10 concludes the paper by raising some challenging issues in model selection and the balance between descriptive and generative models.

## 2 LITERATURE SURVEY—A GLOBAL PICTURE

In this section, we briefly review four research streams and summarize four types of probabilistic models to represent a global picture of the field.

### 2.1 Four Research Streams

#### 2.1.1 Stream 1: The Study of Natural Image Statistics

Any generic vision systems, biologic or machine, must account for image regularities. Thus, it is of fundamental importance to study the statistical properties of natural images. Most of the early work studied natural image statistics from the perspective of image coding and redundancy reduction and often used them to predict/explain the neuron responses.

Historically, Attneave [3], Barlow [5], and Gibson [35] were among the earliest who argued for the ecologic influence on vision perception. Kersten [49], did perhaps, the first experiment measuring the conditional entropy of the intensity at a pixel given the intensities of its neighboring pixels, in a spirit similar to Shannon's [76] experiment of measuring the entropy of English words. Clearly, the strong correlation of intensities between adjacent pixels results in low entropy. Further study of the intensity correlation in natural images leads to an interesting rediscovery of a  $1/f$  power law by Field

[28].<sup>1</sup> By doing a Fourier transform on natural images, the amplitude of the Fourier coefficients at frequency  $f$  (averaged over orientations) fall off in a  $1/f$ -curve (see Fig. 4a). The power may not be exactly  $1/f$  and vary in different image ensembles [72]. This inspired a large body of work in biologic vision and computational neurosciences which study the correlations of not only pixel intensities but responses of various filters at adjacent locations. These works also expand from gray-level static images to color and motion images (see [2], [78] for more references).

Meanwhile, the study on natural image statistics extends from correlations to histograms of filter responses, for example, using Gabor filters.<sup>2</sup> This leads to two interesting observations. First, the histograms of Gabor type filter responses on natural images have high kurtosis [29]. This reveals that natural images have high order (non-Gaussian) structures. Second, it was reported independently by [72], [94] that the histograms of gradient filtered images are consistent over a range of scales (see Fig. 5). The scale invariance experiment is repeated by several teams [13], [38]. Further studies along this direction include investigations on joint histograms and low-dimensional manifolds in high-dimensional spaces. For example, the density on a 7D unit sphere for all  $3 \times 3$  pixel patches of natural images [52], [53]. Going beyond pixel statistics, some most recent work measured the statistics of object shapes [96], contours [32], and the size of regions and objects in natural images [4].

#### 2.1.2 Stream 2: The Analysis of Natural Image Components

The high kurtosis in image statistics observed in stream 1 is only a marginal evidence for hidden structures in natural scenes. A direct way for discovering structures and reducing image redundancy is to transform an image into a superposition of image components. For example, Fourier transform, wavelet transforms [16], [57], and various image pyramids [77] for generic images, and principal component analysis for some specific ensembles of images.

The transforms from image pixels to bases achieve two desirable properties. The first is variable decoupling. The coefficients of these bases are less correlated or become independent in ideal cases. The second is dimension reduction. The number of bases for approximately reconstructing an image is often much smaller than the number of pixels.

If one treats an image as a continuous function, then a mathematical tool for decomposing images is *harmonic analysis* (see [25], [59], [60]). Harmonic analysis is concerned with decomposing various classes of functions (i.e., mathematic spaces) by different bases. Further development along this vein includes the wedgelets, ridgelet, edgelets, curvelets [10], [101].

Obviously, the ensemble of natural images is quite different from those functional classes. Therefore, the image components must be adapted to natural images. This leads to inspiring ideas in recent literature—*sparse coding with overcomplete basis or dictionary* [67]. With overcomplete basis, an image may be reconstructed by a

1. The spectra power-law was first reported in [23] in studying television signals and rediscovered by Cohen et al. [15] in photographic analysis, and then by Burton and Moorhead [100] in optics study. It was Fields' work that brought it to attention of the broad vision communities.

2. Correlations only measures second order moments while histograms include all the high order information, such as skewness (third order) and kurtosis (fourth order).

small (sparse) number of bases in the dictionary. This often leads to 10-100 folds of dimension reduction. For example, an image of  $200 \times 200$  pixels can be reconstructed approximately by about 100 – 500 base images. Olshausen and Field then learned the overcomplete dictionary from natural images. Fig. 13 shows some of the bases. Added to this development is the independent component analysis (ICA) [17], [84]. It is shown in harmonic analysis that the Fourier, wavelet, and ridgelet bases are independent components for various ensembles of mathematical functions (see [25] and references therein). But, for the ensemble of natural images, it is not possible to have an independent basis and one can only compute a basis that maximize some measure of independence. Going beyond the image bases, recently, Zhu et al. [99] proposed the texton representation with each texton consisting of a number of image bases at various geometric, photometric, and dynamic configurations. If we compare the image bases to phonemes in speech, then the textons are larger structures corresponding to words.

### 2.1.3 Stream 3: The Grouping of Natural Image Elements

The third research stream originated from Gestalt psychology [51]. Human visual perception has strong tendency (bias) toward forming global percept (“whole” or pattern) by grouping local elements (“parts”). For example, human vision completes illusory figures [47], and perceives hallucinatory structures from totally random dot patterns [79]. In contrast to research streams 1 and 2, early work in stream 3 focused on *computational procedures and algorithms* that seem to demonstrate performance similar to human perception. This includes work on illusory figure completion and grouping from local edge elements (e.g., Guy and Medioni [41]).

While the Gestalt laws are quite successful in many artificial illusory figures, their applicability in real-world images was haunted by ambiguities. A pair of edge elements may be grouped in one image but separated in the other image, depending on information that may have to be propagated from distant edge elements. So, the Gestalt laws are not really *deterministic laws* but rather *heuristics* or *importance hypotheses* which are better used with probabilities.

Lowe [56] was the first who computed the likelihoods (probabilities) for grouping a pair of line segments based on proximity, colinearity, or parallelism, respectively. Considering a number of line segments that are independently and uniformly distributed in terms of lengths, locations, and orientations in a unit square, Lowe estimated the expected number for a pair of line segments at a certain configuration that are formed *accidentally* according to this uniform distribution. Lowe conjectured that the likelihood of grouping a pair of line segments in real images should be proportional to the inverse of this expected number—which he called *nonaccidental property*. In a similar method, Jacobs [43] calculated the likelihood for grouping a convex figure from a set of line segments. In a similar way, Moisan et al. [61] compute the likelihoods for “meaningful alignments.” More advanced work includes Sarkar and Boyer [74] and Dickinson et al. [24] for generic object grouping and recognition (see Fig. 20). Bienenstock et al. [7] proposed a compositional vision approach for grouping of handwritten characters.

Generally speaking, the probabilities for grouping are, or can be reformulated to, posterior probability ratios of

“grouping” versus “ungrouping” or “on” versus “off” an object. The probabilities are computed based on local features. Recently, people started learning these probabilities and ratios from natural images with supervised input. For example, Geisler et al. [32] computed the *likelihood ratio* for the probability that a pair of edge elements appear in the same curve (grouped manually) against the probability that they appear in different curves (not grouped). Konishi et al. [50] computes probability ratio for a pixel on versus off the edge (object boundary) from some manually segmented images.

### 2.1.4 Stream 4: The Modeling of Natural Image Patterns

The fourth stream of research follows the Bayesian framework and develops explicit models for visual patterns. In the literature, Grenander [36], Cooper [18], and Fu [31] were the pioneers using statistical models for various visual patterns. In the late 1980s and early 1990s, image models become popular and indispensable when people realized that vision problems, typically the shape-from-X problems, are fundamentally ill-posed. Extra information is needed to account for regularities in real-world scenes and the models represent our visual knowledge. Early models all assumed simple *smoothness* (sometimes piecewisely) of surfaces or image regions, and they were developed from different perspectives. For example, physically-based models [8], [81], regularization theory [69], and energy functionals [63]. Later, these concepts all converged to statistical models that prevailed due to two pieces of influential work. The first work is the Markov random field (MRF) modeling [6], [19] introduced from statistical physics. The second work is the Geman and Geman [33] paper which showed that vision inference can be done rigorously by Gibbs sampler under the Bayesian framework. There was extensive literature on Markov random fields and Gibbs sampling in the late 1980s. This trend went down in the early 1990s for two practical reasons: 1) Most of those Markov random field models are based on pair cliques and, thus, do not realistically characterize natural image patterns. 2) The Gibbs sampler is computationally very demanding on such problems.

Other probability models of visual patterns include deformable templates for objects, such as human face [90] and hands [37]. In contrast to the homogeneous MRF models for texture and smoothness, deformable templates are inhomogeneous MRF on small graphs whose nodes are labeled. We should return to more recent MRF models in later section.

## 2.2 Four Categories of Statistical Models

The interactions of the research streams produce four categories of probability models. In the following sections, we briefly review the four types of models to set background for a mathematical framework that unifies them.

### 2.2.1 Category 1: Descriptive Models

First, the integration of stream 1 and stream 4 yields a class of models that we call “descriptive models.” Given an image ensemble and its statistics properties, such as the  $1/f$ -power law, scale invariant gradient histograms, studied in stream 1, one can always construct a probability model which produces the same statistical properties as observed in the image ensemble. The probability is of the Gibbs (MRF) form following a maximum entropy principle [44]. By maximum entropy, the model minimizes the bias while

it satisfies the statistical descriptions. We call such models the descriptive models because they are constructed based on statistical descriptions of the image ensembles.

The descriptive model is attractive because a single probability model can integrate all statistical measures of different image features. For example, a Gibbs model of texture [93] can account for the statistics extracted by a bank of filters, and a Gibbs model of shapes (2D simple curves) can integrate the statistics of various Gestalt properties: proximity, colinearity, parallelism [96]. Such integration is not a simple product of the likelihoods or marginals on different features (like the projection pursuit method) but uses sophisticated energy functions to account for the dependency of these features. This provides a way to exactly measure the “nonaccidental statistics” sought after by Lowe [56]. We shall deliberate on this point in latter section.

The descriptive models are all built on certain graph structures including lattices. There are two types descriptive models in the literature: 1) *Homogeneous models* where the statistics are assumed to be the same for all elements (vertices) in the graph. The random variables are the attributes of vertices, such as texture models. 2) *Inhomogeneous model* where the elements (vertices) of the graph are labeled and different features and statistics are used at different sites, for example, deformable models of human faces.

### 2.2.2 Category 2: Variants of Descriptive Models and Energy Approximations

The descriptive models are often computationally expensive, due to the difficulty of computing the partition (normalizing) functions. This problem becomes prominent when the descriptive models have large image structures and account for high order image statistics. In the literature, there are a few variants to the descriptive models and approximative methods.

The first is causal Markov models. A causal MRF model approximates a descriptive model by imposing a partial (or even linear) order among the vertices of the graph such that the joint probability can be factorized as a product of conditional probabilities. The latter have lower dimensions and, thus, are much easier to learn and to compute. The Causal MRF models are still maximum entropy distributions subject to, sometimes, the same set of statistical constraints as the descriptive models. But, the entropy is maximized in a limited probability space. Examples include texture synthesis in [26], [70] and the recent cut-and-paste work [27], [54].

The second is called pseudodescriptive model. Typical examples include texture synthesis methods by Heeger and Bergen [42] and DeBonet and Viola [21]. They draw independent samples in the feature space, for example, filter responses at multiple scales and orientations at each pixel from the marginal or joint histograms. Though the sampled filter responses satisfy the statistical description in an observed image, there is no image that can produce all these filter responses, as the latter are conflicting with each other. Then, an image is synthesized by a pseudoinverse method. Sampling in the feature space and the pseudoinverse are often computationally convenient but the whole method does not follow a rigorous model.

The other approximative approach for computing the descriptive model introduces a *belief* at each vertex. These beliefs are only normalized at a single site or a pair of sites and they do not necessarily form a legitimate (well normalized)

joint probability for the whole graph. Thus, it avoids computing the partition functions. This technique, originated in statistical physics, includes the mean field approximation, the Bethe and Kikuchi approximations (see Yedidia et al. [89] and Yuille [92]).

### 2.2.3 Category 3: Generative Models

The principled way for tackling the computational complexity of descriptive models (no “hacks” or approximations) is to introduce hidden variables that can “explain away” the strong dependency in observed images. For example, the sparse coding scheme [67] is a typical generative model which assumes an image being generated by a small number of bases. Other models include [20], [30]. The computation becomes less intensive because of the reduced dimensions and the partially decoupling of hidden variables. The generative model must engage some vocabulary of visual descriptions. For example, an overcomplete dictionary for image coding. The elements in the vocabulary specify how images are generated from hidden variables.

The generative models are not separable from descriptive models because the hidden variables must be characterized by a descriptive model, though in the literature, the latter may often be a trivial iid Gaussian model or a causal Markov model. For example, the sparse coding scheme is a two layer generative model and assumes that the image bases are iid hidden variables. Hidden Markov models in speech and motion are also two layer models whose hidden layer is a Markov chain (causal MRF model with linear order).

So, descriptive and generative models must be integrated for developing richer and computationally tractable models. We should deliberate on this in latter sections. Thus, we have a unified family of models for the descriptive (its variants) and generative models. These models are representational.

### 2.2.4 Category 4: Discriminative Models

In contrast to the representational models (descriptive plus generative), some probabilities are better considered computational heuristics viewed from the general task of image parsing—the discriminative models used in stream 3 belong to this category.

In comparison, descriptive models and generative models are used as prior probabilities and likelihoods in the Bayesian framework, while discriminative models approximate the posterior probabilities of hidden variables (often individually) based on local features. As we shall show in later sections, they are *importance proposal probabilities* which drive the stochastic Markov chain search for fast convergence. It was shown, through simple case, that the better the proposal probability approximates the posterior, the faster the algorithm converges [58].

The interaction between discriminative and generative models has not gone very far in the literature. Recent work include the data driven Markov chain Monte Carlo (DDMCMC) algorithms for image segmentation, parsing, and object recognition [82], [83], [98].

### 2.2.5 Summary and Justification of Terminology

To clarify the terminology used above, Fig. 2 shows a trivial example of the four models for a desk object. A desk consists of four legs and a top, denoted, respectively, by variables  $d, l_1, l_2, l_3, l_4, t$  for their attributes (vector valued). Fig. 2a shows the undirected graph for a descriptive model

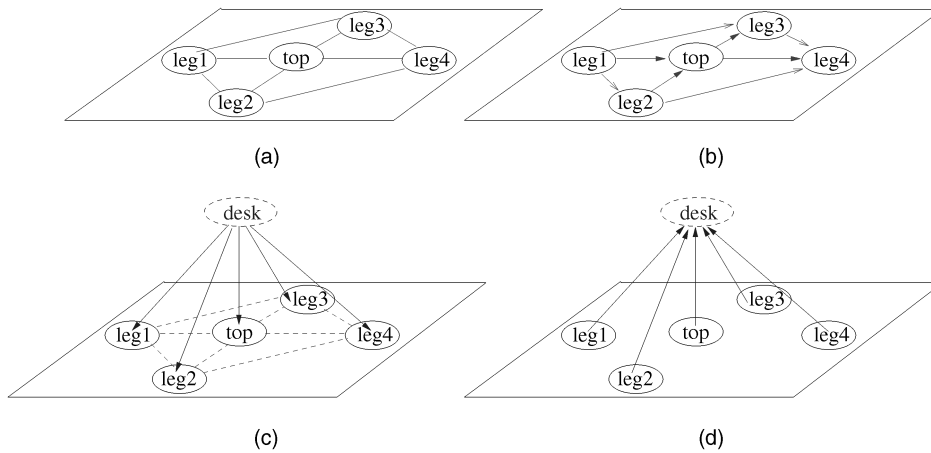


Fig. 2. Four types of models for a simple desk object. (a) Descriptive (MRF), (b) causal MRF, (c) generative + descriptive, and (d) discriminative.

$p(l_1, l_2, l_3, l_4, t)$ . It is in the Gibbs form with a number of energy terms to account for the spatial arrangement of the five pieces. The potential functions of the Gibbs assign low energies and, thus, high probabilities, to more general configurations. This descriptive model accounts for the phenomenological probability that the five pieces occur together without “understanding” a hidden concept of “desk”—denoted by hidden variable  $d$ . The causal MRF model assumes a directed graph in Fig. 2b. Thus, it simplifies the descriptive model as  $p(l_2)p(t|l_1, l_2)p(l_3|t, l_1) p(l_4|t, l_2, l_3)$ . Fig. 2c is a two level generative model which involves a hidden variable  $d$  for the “whole” desk. The desk generates the five pieces by a model  $p(l_1, l_2, l_3, l_4, t|d)$ .  $d$  contains global attributes of the desk which controls the positions of the five parts. If we assume that the five pieces are conditionally independent, then it becomes a context free grammar (without the dashed lines). In general, we still need a descriptive model to characterize the spatial deformation by a descriptive model (see the dashed links). But, this new descriptive model  $p(l_1, l_2, l_3, l_4, t|d)$  is much less complicated than  $p(l_1, l_2, l_3, l_4, t)$  in Fig. 2a. For example, if there are five types of desks, the descriptive model  $p(l_1, l_2, l_3, l_4, t)$  must have complicated energy function so that it has five distinct modes (maxima). But, if  $d$  contains a variable for the desk type, then  $p(l_1, l_2, l_3, l_4, t|d)$  has a single mode for each type of desk and its potential is quite easy to compute. Finally, Fig. 2d is a discriminative model, the links are pointed from parts to whole (reversing the generative arrows). It tries to compute a number of posterior probabilities  $p(d|t)$ ,  $p(d|l_i)$ ,  $i = 1, 2, 3, 4$ . These probabilities are often treated as “votes” that are then summed up in a generalized Hough transform.

Syntactically, the generative, causal Markov, and discriminative models can all be called Bayesian (causal, belief) networks as long as there are no loops in the graphs. But, this terminology is very confusing in the literature. Our terminology for the four types of models is from a semantic perspective. We call it a generative model if the links are directed downwards in the conceptual hierarchy. We call it a discriminative model if the links are upward. For example, the Bayes networks used by [24], [74], [75] (see Fig. 20) are discriminative models. We call it a causal Markov model if the links are pointed to variables at the same conceptual level (also see Fig. 17). For example, we consider hidden Markov

models in motion or speech as two layer generative models where the hidden variables is governed by a causal Markov (descriptive) model because they belong to the same semantic level. When a generative model is integrated with descriptive model, the integrated model can still be called generative model—a slight abuse of terminology.

It is worth noting that not all hidden variables are used in generative models. The mixed Markov model, as a variant of descriptive model, uses hidden variables to specify the neighborhood for variables at the same semantic level. These hidden variables are called “address variables” [65]. In contrast, the hidden variables in generative models represent entities of large structures.

### 3 PROBLEM FORMULATION

Now, we start with a general formulation of visual modeling, from which we derive the descriptive and generative models for visual knowledge representation.

Let  $\mathcal{E}$  denote the ensemble of natural images in our environment. As the number of natural images is so large, it makes sense to talk about a frequency  $f(\mathbf{I})$  for images  $\mathbf{I} \in \mathcal{E}$ .  $f(\mathbf{I})$  is intrinsic to our environment and our sensory system. For example,  $f(\mathbf{I})$  would be different for fish living in a deep ocean or rabbits living in a prairie, or if our vision is 100 times more acute. The general goal of visual modeling is to estimate the frequency  $f(\mathbf{I})$  by a probabilistic model  $p(\mathbf{I})$  based on a set of observations  $\{\mathbf{I}_1^{\text{obs}}, \dots, \mathbf{I}_M^{\text{obs}}\} \sim f(\mathbf{I})$ .  $p(\mathbf{I})$  represents, exclusively, our understandings of image regularities and, thus, all of our representational knowledge for vision.<sup>3</sup>

It may sound quite ridiculous to estimate a density like  $f(\mathbf{I})$  which is often in a  $256 \times 256$  space. But as we shall show in the rest of the paper, this is possible because of the strong regularities in natural images, and easy access to a very large number of images. For example, if a child sees 20 images per second, and opens eyes 16 hours a day, then by the age of 10, he/she has seen three billion images. The probability model  $p(\mathbf{I})$  should approach  $f(\mathbf{I})$  by minimizing a Kullback-Leibler divergence  $KL(f||p)$  from  $f$  to  $p$ ,

3. A frequency  $f(\mathbf{I})$  is an objective probability for the ensemble  $\mathcal{E}$ , while a model  $p(\mathbf{I})$  is subjective and biased by the finite data observation and choice of model families.

$$KL(f||p) = \int f(\mathbf{I}) \log \frac{f(\mathbf{I})}{p(\mathbf{I})} d\mathbf{I} = E_f[\log f(\mathbf{I})] - E_f[\log p(\mathbf{I})]. \quad (1)$$

Approximating the expectation  $E_f[\log p(\mathbf{I})]$  by a sample average leads to the standard maximum-likelihood estimator (MLE),

$$p^* = \arg \min_{p \in \Omega_p} KL(f||p) \approx \arg \max_{p \in \Omega_p} \sum_{m=1}^M \log p(\mathbf{I}_m^{\text{obs}}), \quad (2)$$

where  $\Omega_p$  is a family of distributions where  $p^*$  is searched for. One general procedure is to search for  $p$  in a sequence of nested probability families,

$$\Omega_0 \subset \Omega_1 \subset \dots \subset \Omega_K \rightarrow \Omega_f \ni f,$$

where  $K$  indexes the dimensionality of the space, e.g., the number of free parameters. As  $K$  increases, the probability family should be general enough to approach  $f$  to an arbitrary predefined precision.

There are two choices for the families  $\Omega_p$  in the literature.

The first choice is the descriptive model. They are called exponential or log-linear models in statistics, and Gibbs models in physics. We denote them by

$$\Omega_1^d \subset \Omega_2^d \subset \dots \subset \Omega_K^d \rightarrow \Omega_f \ni f. \quad (3)$$

The dimension of the space  $\Omega_1^d$  is augmented by increasing the number of *feature statistics* of  $\mathbf{I}$ .

The second choice is the generative model, or mixture models in statistics, denoted by

$$\Omega_1^g \subset \Omega_2^g \subset \dots \subset \Omega_K^g \rightarrow \Omega_f \ni f. \quad (4)$$

The dimension of  $\Omega_p$  is augmented by introducing hidden variables for the underlying image structures in  $\mathbf{I}$ .

Both families are general enough for approximating any distribution  $f$ . In the following sections, we deliberate on the descriptive and generative models and learning methods and then discuss their unification and the philosophy of model selection.

## 4 DESCRIPTIVE MODELING

In this section, we review the basic principle of descriptive modeling and show a spectrum of seven examples for modeling visual patterns from low to high levels.

### 4.1 The Basic Principle of Descriptive Modeling

The basic idea of descriptive modeling is shown in Fig. 3. Let  $\mathbf{s} = (s_1, \dots, s_n)$  be a representation of a visual pattern. For example,  $\mathbf{s} = \mathbf{I}$  could be an image with  $n$  pixels and, in general,  $\mathbf{s}$  could be a list of attributes for vertices in a random graph representation. An observable data ensemble is illustrated by a cloud of points in an  $n$ -space and each point is an instance of the visual pattern. A descriptive method extracts a set of  $K$  **features** as *deterministic transforms* of  $\mathbf{s}$ , denoted by  $\phi_k(\mathbf{s}), k = 1, \dots, K$ . For example,  $\phi_k(\mathbf{I}) = \langle F, \mathbf{I} \rangle$  is a projection of image  $\mathbf{I}$  on a linear filter (say Gabor)  $F$ . These features (such as  $F$ ) are illustrated by axes in Fig. 3. In general, the axes don't have to be straight lines and could be more than one-dimensional. Along these axes, we can compute the projected histograms of the ensemble (the right side of Fig. 3). We denote these

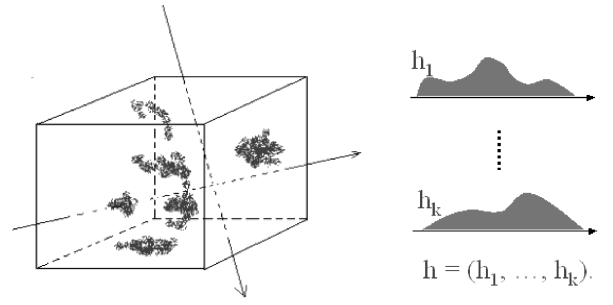


Fig. 3. Descriptive modeling: estimating a high-dimensional frequency  $f$  by a maximum entropy model  $p$  that matches the low-dimensional (marginal) projections of  $f$ . The projection axes could be nonlinear.

histograms as  $\mathbf{h}_k^{\text{obs}}$  for features  $\phi_k(\mathbf{s}), k = 1, 2, \dots, K$ . They are estimates to the marginal statistics of  $f(\mathbf{s})$ .

A model  $p$  must match the marginal statistics  $\mathbf{h}_k^{\text{obs}}, k = 1, \dots, K$  if it is to estimate  $f(\mathbf{s})$ . Thus, we have descriptive constraints:

$$E_p[h(\phi_k(\mathbf{s}))] = \mathbf{h}_k^{\text{obs}} \approx E_f[h(\phi_k(\mathbf{s}))], \quad k = 1, \dots, K. \quad (5)$$

The least biased model that satisfies the above constraints is obtained by maximum entropy [44] and this leads to the FRAME model [93],

$$p_{\text{des}}(\mathbf{s}; \beta) = \frac{1}{Z(\beta)} \exp \left\{ - \sum_{k=1}^K \langle \lambda_k, h(\phi_k(\mathbf{s})) \rangle \right\}. \quad (6)$$

The parameters  $\beta = (\lambda_1, \dots, \lambda_K)$  are Lagrange multipliers and they are computed by solving the constraint equations (5).  $\lambda_k$  is a vector whose length is equal to the number of bins in the histogram  $h(\phi_k(\mathbf{s}))$ . As the features  $\phi_k(\mathbf{s}), k = 1, 2, \dots, K$  are often correlated, the parameters  $\beta$  are learned to weight these features. Thus,  $p_{\text{des}}(\mathbf{s}; \beta)$  integrates all the observed statistics.<sup>4</sup>

The selection of features in  $p_{\text{des}}$  is guided by a minimum entropy principle. For any new feature  $\phi^+$ , we can define its nonaccidental statistics following Zhu et al. [93].

**Definition 1 (Nonaccidental Statistics).** Let  $\mathbf{h}_f^+$  be the observed statistics for a novel feature  $\phi^+$  computed from the ensemble, i.e.,  $\mathbf{h}_f^+ \approx E_f[h(\phi_+(\mathbf{s}))]$  and  $\mathbf{h}_p^+ = E_{p_{\text{des}}}[h(\phi_+(\mathbf{s}))]$  its expected statistics according to a current model  $p_{\text{des}}$ . Then, the nonaccidental statistics of  $\phi^+$ , with its correlations to the previous  $K$  features removed, is a quadratic distance  $d(\mathbf{h}_f^+, \mathbf{h}_p^+)$ .

$d(\mathbf{h}_f^+, \mathbf{h}_p^+)$  measures the statistics discrepancy of  $\phi^+$  which are not captured by the previous  $K$  features. Let  $p_{\text{des}}^+$  be an augmented descriptive model with the  $K$  statistics in  $p_{\text{des}}$  plus the feature  $\phi^+$ , then the following theorem is observed in Zhu et al. [93].

**Theorem 1 (Feature Pursuit).** In the above notation, the nonaccidental statistics of feature  $\phi^+$  is equal to the entropy deduction,

$$d(\mathbf{h}_f^+, \mathbf{h}_p^+) = KL(f||p_{\text{des}}) - KL(f||p_{\text{des}}^+) = \text{entropy}(p_{\text{des}}) - \text{entropy}(p_{\text{des}}^+), \quad (7)$$

4. In natural language processing, such Gibbs model was also used in modeling the distribution of English letters [22].

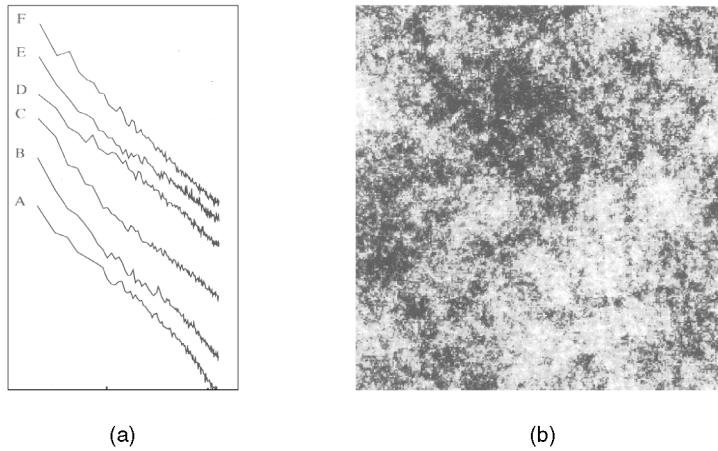


Fig. 4. (a) The log-Fourier-amplitude of natural images are plotted against  $\log f$ , courtesy of Field [28]. (b) A randomly sampled image with  $1/f$  Fourier amplitude, courtesy of Mumford [65].

where  $d(\mathbf{h}_f^+, \mathbf{h}_p^+)$  is a quadratic distance between the two histograms.

As entropy is the logarithmic volume of the ensemble governed by  $p_{\text{des}}$ , the higher the nonaccidental statistics, the more informative feature  $\phi^+$  is for the visual pattern in terms of reducing uncertainty. Thus, a feature  $\phi^+$  is selected sequentially for maximum entropy reduction following (7).

The Cramer and Wold theorem states that the descriptive model  $p_{\text{des}}$  can approximate any densities  $f$  using linear axes only (also see [95]).

**Theorem 1 (Cramer and Wold).** *Let  $f$  be a continuous density, then  $f$  is a linear combination of  $\mathbf{h}$ , the latter are the marginal distributions on the linear filter response  $F^{(\xi)} * \mathbf{s}$ , and  $f$  can be reconstructed by  $p_{\text{des}}$ .*

## 4.2 A Spectrum of Descriptive Models for Visual Patterns

In the past few years, the descriptive models have successfully accounted for the observed natural image statistics (stream 1) and modeled a broad spectrum of visual patterns displayed in Fig. 1. In this section, we show seven examples.

### 4.2.1 Model D1: Descriptive Model for $1/f$ -power Law of Natural Images

An important discovery in studying the statistics of natural images is the  $1/f$  power-law (see review in stream 1). Let  $\mathbf{I}$  be a natural image and  $\hat{\mathbf{I}}(\xi, \eta)$  its Fourier transform. Let  $A(f)$  be the Fourier amplitude  $|\hat{\mathbf{I}}(\xi, \eta)|$  at frequency  $f = \sqrt{\xi^2 + \eta^2}$  averaged over all orientations, then  $A(f)$  falls off in a  $1/f$ -curve.

$$A(f) \propto 1/f, \quad \text{or} \quad \log A(f) = \text{const} - \log f.$$

Fig. 4a is a result in logarithmic scale by Field [28] for six natural images. The curves are fit well by straight lines in log-plot. This observation reveals that natural images contain equal Fourier power at each frequency band—scale invariance. That is,

$$\iint_{f^2 \leq \xi^2 + \eta^2 \leq (2f)^2} |\hat{\mathbf{I}}(\xi, \eta)| d\xi d\eta = 2\pi \int_{f^2}^{4f^2} \frac{1}{f^2} df^2 = \text{const.}, \quad \forall f.$$

The descriptive model that accounts for such statistical regularity is surprisingly simple. It was showed by

Mumford [65] that a Gaussian Markov random field (GMRF) model below has exactly  $1/f$ -Fourier amplitude.

$$p_{1/f}(\mathbf{I}; \beta) = \frac{1}{Z} \exp \left\{ - \sum_{x,y} \beta |\nabla \mathbf{I}(x, y)|^2 \right\}, \quad (8)$$

where  $|\nabla \mathbf{I}(x, y)|^2 = (\nabla_x \mathbf{I}(x, y))^2 + (\nabla_y \mathbf{I}(x, y))^2$ .  $\nabla_x$  and  $\nabla_y$  are the gradients. As the Gibbs energy is of a quadratic form and its matrix is real symmetric circulant, by a spectral analysis (see [68]) its eigenvectors are the Fourier bases and its eigenvalues are the spectra.

This simply demonstrates that the much celebrated  $1/f$ -power law is nothing more than a second order moment constraint in the maximum entropy construction,

$$E_p [|\nabla \mathbf{I}(x, y)|^2] = \frac{1}{2\beta} \approx E_f [|\nabla \mathbf{I}(x, y)|^2], \quad \forall x, y. \quad (9)$$

This is equivalent to a  $1/f$  constraint in the Fourier amplitude.

Since  $p_{1/f}(\mathbf{I}; \beta)$  is a Gaussian model, one can easily draw a random sample  $\mathbf{I} \sim p_{1/f}(\mathbf{I}; \beta)$ . Fig. 4b shows a typical sample image by Mumford [65]. It has very little structure in it! We will revisit the case in the generative model.

### 4.2.2 Model D2: Descriptive Model for Natural Images with Scale-Invariant Histograms

The second important discovery of natural image statistics is the scale-invariance of gradient histograms [72], [94]. Take a natural image  $\mathbf{I}$  and build a pyramid with a number of  $n$  scales,  $\mathbf{I} = \mathbf{I}^{(0)}, \mathbf{I}^{(1)}, \dots, \mathbf{I}^{(n)}$ .  $\mathbf{I}^{(s+1)}$  is obtained by an average of  $2 \times 2$  pixels in  $\mathbf{I}^{(s)}$ . The histograms  $\mathbf{h}^{(s)}$  of gradients  $\nabla_x \mathbf{I}^{(s)}(x, y)$  (or  $\nabla_y \mathbf{I}^{(s)}(x, y)$ ) are plotted in Fig. 5a for three scales  $s = 0, 1, 2$ . Fig. 5b shows the logarithm of the histograms averaged over a number of images.

These histograms demonstrate high kurtosis and are amazingly consistent over a range of scales. Let  $\mathbf{h}^{\text{obs}}$  be the normalized histogram averaged over three scales and impose constraints that a model  $p$  should produce the same histograms (marginal distributions),

$$E_p [h(\nabla_x \mathbf{I}^{(s)})] = E_p [h(\nabla_y \mathbf{I}^{(s)})] = \mathbf{h}^{\text{obs}}, \quad s = 0, 1, 2, 3. \quad (10)$$



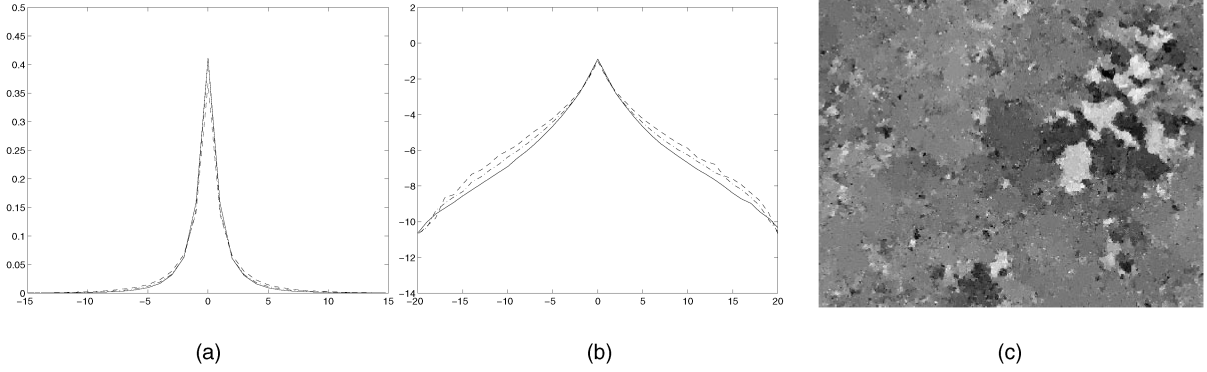


Fig. 5. (a) Gradient histograms over three scales. (b) Logarithm of histograms. (c) A randomly sampled images from a descriptive model  $p_{\text{inv}}(\mathbf{I}; \beta)$ . Courtesy of Zhu and Mumford [94].

Zhu and Mumford [94] derived a descriptive model,

$$p_{\text{inv}}(\mathbf{I}; \beta) = \frac{1}{Z} \exp \left\{ - \sum_{s=0}^3 \sum_{(x,y) \in \Lambda^{(s)}} \lambda_x^{(s)}(\nabla_x \mathbf{I}^{(s)}(x,y)) + \lambda_y^{(s)}(\nabla_y \mathbf{I}^{(s)}(x,y)) \right\}. \quad (11)$$

$\Lambda^{(s)}$  is the image lattice at scale  $s$ .  $\beta = (\lambda_x^{(0)}(), \lambda_y^{(0)}(), \dots, \lambda_x^{(3)}(), \lambda_y^{(3)}())$  are the parameters and each  $\lambda_x^{(s)}()$  is a 1D potential function quantized by a vector.

Fig. 5c shows a typical image sampled from this model by a Gibbs sampler that was used in [33]. This image has the scale-invariant histograms shown in Figs. 5a and 5b. Clearly, the sampled image demonstrates some piecewise smoothness and consists of microstructures of various sizes.

To make connection with other models, we remark on two aspects of  $p_{\text{inv}}(\mathbf{I}; \beta)$ .

First, by choosing only one scale  $s = 0$ , the constraints in (10) is a superset of the constraints in (9), as the histogram includes the variance. Therefore,  $p_{\text{inv}}$  also observes the  $1/f$ -power law but with much more structures.

Second, with only one scale,  $p_{\text{inv}}$  reduces to the general smoothness models widely used in shape-from-X and denoising (see review of stream 4). The learned potential functions  $\lambda_x()$  and  $\lambda_y()$  match pretty close to the manually selected energy functions. This bridges the learning of Gibbs model with PDEs in image processing (see details in [94]).

#### 4.2.3 Model D3: Descriptive Model for Textures

The third descriptive model accounts for interesting psychophysical observations in texture study that histograms of a set of Gabor filters may be *sufficient statistics* in texture perception, i.e., two textures cannot be told apart in early vision if they share the same histograms of Gabor filters [14].

Let  $F_1, \dots, F_K$  be a set of linear filters (such as Laplacian of Gaussian, Gabors), and  $h(F_k * \mathbf{I})$  the histograms of filtered image  $F_k * \mathbf{I}$  for  $k = 1, 2, \dots, K$ . Each  $F_k$  corresponds to an axis and  $h(F_k * \mathbf{I})$  a 1D marginal distribution in Fig. 3. From an observed image, a set of histograms  $\mathbf{h}_k^{\text{obs}}, k = 1, 2, \dots, K$  are extracted. By imposing the descriptive constraints

$$E_p[h(F_k * \mathbf{I})] = \mathbf{h}_k^{\text{obs}}, \quad \forall k = 1, 2, \dots, K. \quad (12)$$

A FRAME model [93], [95] is obtained through maximum entropy.

$$p_{\text{tex}}(\mathbf{I}; \beta) = \frac{1}{Z} \exp \left\{ - \sum_{(x,y) \in \Lambda} \sum_{k=1}^K \lambda_k(F_k * \mathbf{I}(x,y)) \right\}. \quad (13)$$

where  $\beta = (\lambda_1(), \lambda_2(), \dots, \lambda_K())$  are potential functions with each function  $\lambda_i()$  being approximated by a vector.  $p_{\text{tex}}(\mathbf{I}; \beta)$  extends traditional Markov random field models [6], [19] by replacing pairwise cliques with Gabor filters and by upgrading the quadratic energy to nonparametric potential functions which account for high order statistics.

Fig. 6 illustrates the modeling of a texture pattern. As texture is homogeneous, it uses spatial average in a single input image in Fig. 6a to estimate the ensemble average  $\mathbf{h}_k^{\text{obs}}, k = 1, 2, \dots, K$ . With  $K = 0$  constraints,  $p_{\text{tex}}(\mathbf{I}; \Theta)$  is a uniform distribution and a *typical* random sample is a noise image shown in Fig. 6b. With  $K = 1, 2, 7$  histogram constraints, the randomly sampled images from the learned Gibbs models  $p_{\text{tex}}(\mathbf{I}; \Theta)$ , are shown in Figs. 6c, 6d, and 6e, respectively. The samples are drawn by Gibbs sampler [33] from  $p_{\text{tex}}(\mathbf{I}; \beta)$  and the selection of filters are governed by a minimax entropy principle [93]. A wide variety of textures are modeled in this way. In a similar way, one can put other statistics, such as filter correlations, in the model (See [71]).

#### 4.2.4 Model D4: Descriptive Model for Texton (Attributed Point) Process

The descriptive models  $p_{1/f}$ ,  $p_{\text{inv}}$ , and  $p_{\text{tex}}$  are all based on lattice and pixel intensities. Now, we review a fourth model for texton (attributed point process) that extends lattices to graphs and extends pixel intensity to attributes. Texton processes are very important in perceptual organization. For example, Fig. 1 shows a point process for the music band and Fig. 7a shows a wood pattern where a texton represents a segment of the tree trunk.

Suppose a texton  $t$  has attributes  $x, y, \sigma, \theta, c$  for its location, scale, orientation, and photometric contrast, respectively. A texton pattern with an unknown number of  $n$  textons is represented by,

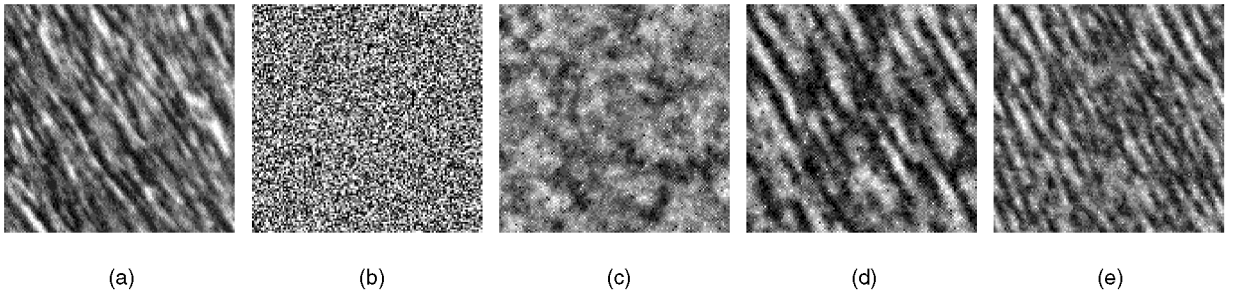


Fig. 6. Learning a sequence of descriptive models for a fur texture: (a) The observed texture image, (b), (c), (d), and (e) are the synthesized images as random samples from  $p_{\text{tex}}(\mathbf{I}; \beta)$  using  $K = 0, 1, 2, 7$  filter histograms, respectively. The images are obtained by Gibbs sampler. Courtesy of Zhu et al. [95].

$$\mathbf{T} = (n, \{ t_j = (x_j, y_j, s_j, \theta_j, c_j), j = 1, \dots, n \}).$$

Each texton  $t$  has a neighborhood  $\partial t$  defined by spatial proximity, good continuity, parallelism or other Gestalt properties. It can be decided deterministically or stochastically. Once a neighborhood graph is decided, one can extract a set of features  $\phi_k(t|\partial t), k = 1, 2, \dots, K$  at each  $t$  measuring some Gestalt properties between  $t$  and its neighbors in  $\partial t$ . If the point patterns are homogeneous, then through constraints on the histograms, a descriptive model is obtained to capture the spatial organization of textons [40],

$$p_{\text{txn}}(\mathbf{T}; \beta_o, \beta) = \frac{1}{Z} \exp \left\{ -\beta_o n - \sum_{j=1}^n \sum_{k=1}^K \lambda_k(\phi_k(t_j|\partial t_j)) \right\}, \quad (14)$$

$p_{\text{txn}}$  is distinct from previous descriptive models in two respects. 1) The number of elements varies, thus a death-birth process must be used in simulating the model. 2) Unlike the static lattice, the spatial neighborhood of each element can change dynamically during the simulation.

Fig. 7a shows an example of a wood pattern with  $\mathbf{T}$  given, from which a texton model  $p_{\text{txn}}$  is learned. Figs. 7b, 7c, and 7d show three stages of the MCMC sampling process of  $p_{\text{txn}}$  at  $t = 1, 30, 332$  sweeps, respectively. This example demonstrates that global pattern arises through simple local interactions in  $p_{\text{txn}}$ . More point patterns are referred to [40].

#### 4.2.5 Model D5: Descriptive Models for 2D Open Curves: Snake and Elastica

Moving up the hierarchy from point and textons to curves, we see that many existing curve models are descriptive.

Let  $C(s) s \in [a, b]$  be an open curve, there are two curve models in the literature. One is the prior term used in the popular SNAKE or active contour model [48].

$$p_{\text{snk}}(C; \alpha, \beta) = \frac{1}{Z} \exp \left\{ - \int_a^b \alpha |\nabla C(s)|^2 + \beta |\nabla^2 C(s)|^2 ds \right\},$$

where  $\nabla C(s)$  and  $\nabla^2 C(s)$  are the first and second derivatives.

The other is an Elastica model [62] simulating a Ulenbeck process of a moving particle with friction, let  $\kappa(s)$  be the curvature, then

$$p_{\text{els}}(C; \beta) = \frac{1}{Z} \exp \left\{ - \int_a^b [\alpha + \beta \kappa^2(s)] ds \right\}.$$

$\alpha$  controls the curve length as a decay probability for terminating the curve, like  $\beta_o$  in  $p_{\text{txn}}$ .

Figs. 8a and 8b show two sets of randomly sampled curves each starting from an initial point and orientation, the curves show general smoothness like the images in Fig. 5c. Williams and Jacobs [85] adopted the Elastica model for curve completion. They define the so-called ‘‘stochastic completion field’’ between two oriented line segments (a source and a sink). Suppose a particle is simulated by a random walk, it starts from the source and ends at the sink. The completion fields shown in Figs. 8c and 8d display the probability that the particle passing a point  $(x, y)$  in the lattice (dark means high probability). This was used as a model for illusory contours.

#### 4.2.6 Model D6: Descriptive Models for 2D Closed Curves

The next descriptive model generalizes the smoothness curve model to 2D shape models with both contour and region-based features. Let  $\Gamma(s), s \in [0, 1]$  be a simple closed curve of normalized length. One can always represent a curve by polygon with a large enough number of vertices. Some edges can be added on the polygon for spatial proximity, parallelism, and symmetry. Thus, a random graph structure is established, and some Gestalt properties  $\phi_k(), k = 1, 2, \dots, K$  can be extracted at each vertex and its

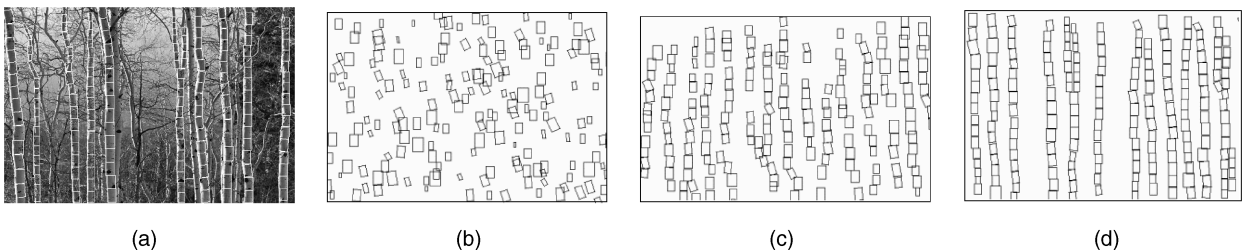


Fig. 7. Different stages of simulating a wood pattern with local spatial interactions of textons. Each texton is represented by a small rectangle. (a) observed, (b)  $t = 1$ , (c)  $t = 30$ , and (d)  $t = 332$ . After Guo et al. [40].

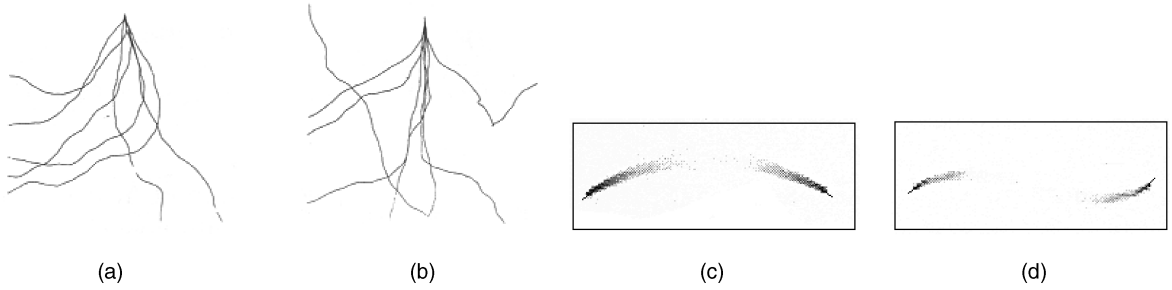


Fig. 8. (a) and (b) Two sets of random sampled curves from the Elastica model. After Mumford [62]. (c) and (d) The stochastic completion fields. After Williams and Jacobs [85].

neighbors, such as colinearity, cocircularity, proximity, parallelism, etc. Through constraints on the histograms of such features, a descriptive model is obtained in [96],

$$p_{\text{shp}}(\Gamma; \beta) = \frac{1}{Z} \exp \left\{ \sum_{k=1}^K \int_0^1 \lambda_k(\phi_k(s)) ds \right\}. \quad (15)$$

This model is invariant to translation, rotation, and scaling. By choosing features  $\phi_k(s)$  to be  $\nabla, \nabla^2, \kappa(s)$ , this model is a nonparametric extension of the SNAKE and Elastica models on open curves.

Fig. 9 shows a sequence of shapes randomly sampled from  $p_{\text{shp}}(\Gamma; \beta)$ . The training ensemble includes contours of animals and tree leaves. The sampled shapes at  $K = 0$  (i.e., no features) are very irregular (sampled by Markov chain random walk under the hard constraint that the curve is closed and has no self-intersection; the MC starts with a circle) and become smooth at  $K = 2$  which integrates two features: colinearity and cocircularity measured by the curvature and derivative of curvature  $\kappa(s)$  and  $\nabla\kappa(s)$ , respectively. Elongated and symmetric “limbs” appear at  $K = 5$  when we integrate crossing region proximity, parallelism, etc.

#### 4.2.7 Model D7: Descriptive Models for 2D Human Face

Moving up to high-level patterns, descriptive models were used for modeling human faces [90] and hand [37], but early deformable models were manually designed, though in principle, they could be reformulated in the maximum entropy form. Recently, a descriptive face model is learned from data by [55] following the minimax entropy scheme.

A face is represented by a list of  $n$  (e.g.,  $n = 83$ ) key points which are manually decided. Connecting these points forms the sketch shown in Fig. 10. Thus, each face is a point in a 166-space. After normalization in location,

rotation and scaling, it has 162 dimensions. Fig. 10a shows four of example faces from the data ensemble.

Unlike the previous homogeneous descriptive models where all elements in a graph (or lattice) are subject to the same statistical constraints, these key points on the face are labeled and, thus, different statistical constraints are imposed at each location.

Suppose we extract  $K$  features  $\phi_k(V), k = 1, 2, \dots, K$  on the graph  $V$ , then a descriptive model is,

$$p_{\text{fac}}(V; \beta) = \frac{1}{Z} \exp \left\{ - \sum_{k=1}^K \lambda_k(\phi_k(V)) \right\}. \quad (16)$$

Liu et al. did a PCA to reduce the dimension first and, therefore, the features  $\phi_k(V)$  are extracted on the PCA coefficients. Fig. 10b shows four sampled faces from a uniform model in the PCA-coefficient space bounded by the covariances. The sampled faces in Figs. 10c and 10d become more pleasant as the number of features increases. When  $K = 17$ , the synthesized faces are no longer distinguishable from faces in the observed ensemble.

#### 4.2.8 Summary: A Continuous Spectrum of Models on the Space of Random Graphs

To summarize this section, visual patterns, ranging from generic natural images, textures, textons, curves, 2D shapes, and objects, can all be represented on attributed random graphs. All the descriptive models reviewed in this section are focused on different subspaces of a huge space of random graphs. Thus, these models are examples in a “continuous” spectrum in the graph space (see (3))! Though the general ideas of defining probability on random graphs were discussed in Grenander’s pattern theory [36], it will be a long

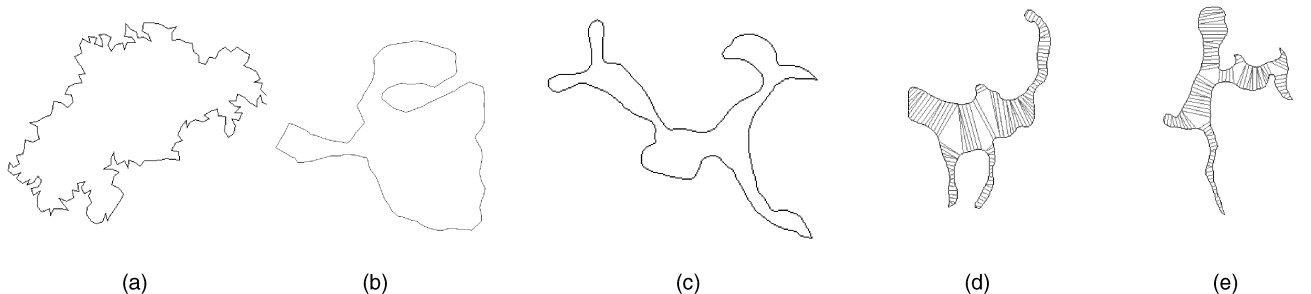


Fig. 9. Learning a sequence of models  $p_{\text{shp}}(\Gamma; \beta)$  for silhouettes of animals and plants, such as cats, dogs, fish, and leaves. (a), (b), (c), and (e) are typical samples from  $p_{\text{shp}}$  with  $K = 0, 2, 5, 5, 5$ , respectively. The line segments show the medial axis features. Courtesy of Zhu [96].

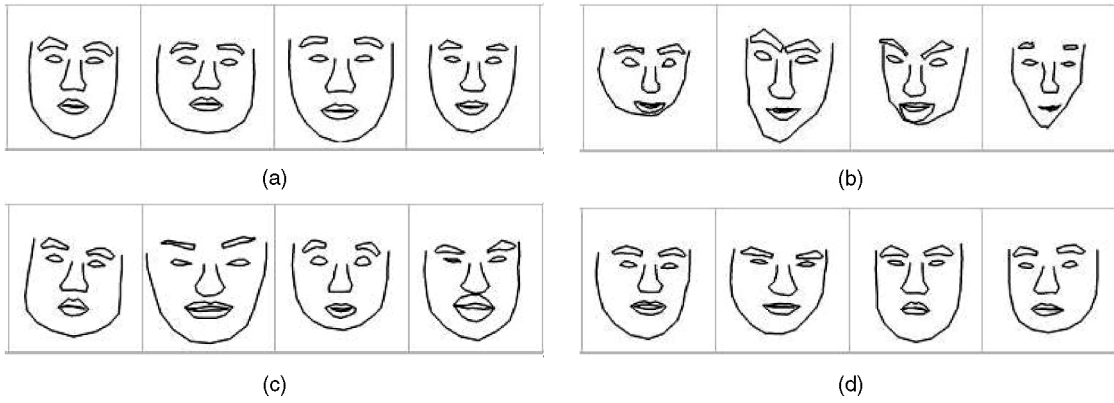


Fig. 10. Learning a sequence of face models  $p_{fac}(V; \beta)$ . (a) Four of the observed faces from a training data set. (b), (c), and (d) Four of the stochastically sampled faces with  $K = 0, 4, 17$  statistics, respectively. Courtesy of Liu et al. [55].

way for developing such models as well as discovering a sufficient set of features and statistics on various graphs.

## 5 CONCEPTUALIZATION OF VISUAL PATTERNS AND STATISTICS PHYSICS

Now, we study an important theoretical issue associated with visual modeling: How do we define a visual pattern mathematically? For example, what is the definition of a human face, or a texture? In mathematics, a concept is equalized to a set. However, a visual pattern is characterized by a probabilistic model as the previous section showed. The connection between a deterministic set and a statistical model was established in modern statistical physics.

### 5.1 Background: Statistical Physics and Ensembles

Modern statistical physics is a subject studying macroscopic properties of a system involving massive amounts of elements [12]. Fig. 11 illustrates three types of physical systems that are interesting to us.

*Microcanonical ensembles.* Fig. 11a is an insulated system of  $N$  elements. The elements could be atoms, molecules, and electrons in systems such as gas, ferro-magnetic material, fluid, etc.  $N$  is really big, say  $N = 10^{23}$  and is considered infinity. The system is decided by a configuration or state  $s = (\mathbf{x}^N, \mathbf{m}^N)$ , where  $\mathbf{x}^N$  describes the coordinates of the  $N$  elements and  $\mathbf{m}^N$  their momenta [12]. It is impractical to study the  $6N$  vector  $s$  and, in fact, these microscopic states are less relevant, and people are more interested in the macroscopic properties of the system as a whole, say the number of elements  $N$ , the total energy  $E(s)$ , and total volume  $V$ . Other derivative properties are temperature and pressure, etc.

If we denote by  $\mathbf{h}(s) = (N, E, V)$  the macroscopic properties, at thermodynamic equilibrium all microscopic states that satisfy this property is called a *microcanonical ensemble*,

$$\Omega_{mce}(\mathbf{h}_o) = \{s = (\mathbf{x}^N, \mathbf{m}^N) : \mathbf{h}(s) = \mathbf{h}_o = (N, V, E)\}.$$

$s$  is an **instance** and  $\mathbf{h}(s)$  is a **summary** of the system state for practical purposes. Obviously,  $\Omega_{mce}$  is a deterministic set or an equivalence class for all states that satisfy a descriptive constraints  $\mathbf{h}(s) = \mathbf{h}_o$ .

An essential assumption in statistical physics is, as a first principle,

*“all microscopic states are equally likely at thermodynamic equilibrium.”*

This is simply a maximum entropy assumption. Let  $\Omega \ni s$  be the space of all possible states, then  $\Omega_{mce} \subset \Omega$  is associated with a uniform probability,

$$p_{unif}(s; \mathbf{h}_o) = \begin{cases} 1/|\Omega_{mce}(\mathbf{h}_o)| & \text{for } s \in \Omega_{mce}(\mathbf{h}_o), \\ 0 & \text{for } s \in \Omega/\Omega_{mce}(\mathbf{h}_o). \end{cases}$$

*Canonical ensembles.* The canonical ensemble refers to a small system (with fixed volume  $V_1$  and elements  $N_1$ ) embedded in a microcanonical ensemble, see Fig. 11b. The canonical ensemble can exchange energy with the rest system (called heat bath or reservoir). The system is relatively small, e.g.,  $N_1 = 10^{10}$ , so that the bath can be considered a microcanonical ensemble itself.

At thermodynamic equilibrium, the microscopic state  $s_1$  for the small system follows a Gibbs distribution,

$$p_{Gib}(s_1; \beta) = \frac{1}{Z} \exp\{-\beta E(s_1)\}.$$

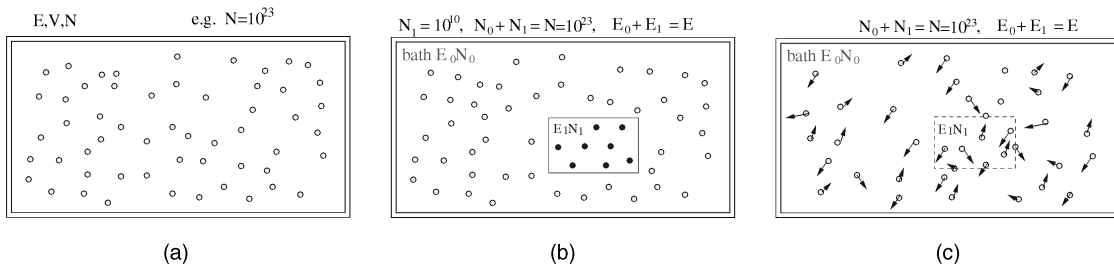


Fig. 11. Three typical ensembles in statistical mechanics. (a) Microcanonical ensemble, (b) canonical ensemble, and (c) grand-canonical ensemble.

The conclusion was stated as a general theorem by Gibbs [34]: “If a system of a great number of degrees of freedom is micro-canonically distributed in phase, any very small part of it may be regarded as canonically distributed.”

Basically, this theorem states that the Gibbs model  $p_{\text{Gib}}$  is a conditional probability of the uniform model  $p_{\text{unif}}$ . This conclusion is extremely important because it bridges a deterministic set  $\Omega_{mce}$  with a descriptive model  $p_{\text{Gib}}$ . We consider this as a true origin of probability for modeling visual patterns. Some detailed deduction of this conclusion in vision models can be found in [87].

*Grand-Canonical ensembles.* When the small system with a fixed volume of  $V_1$  can also exchange elements with the bath as in liquid and gas materials, then it is called a grand-canonical ensemble, see Fig. 11c. The grand-canonical ensemble follows a distribution,

$$p_{\text{gce}}(s_1; \beta_o, \beta) = \frac{1}{Z} \exp\{-\beta_o N_1 - \beta E(s_1)\},$$

where an extra parameter  $\beta_o$  controls the number of elements  $N_1$  in the ensemble.

## 5.2 Conceptualization of Visual Patterns

In statistical mechanics, one is concerned with macroscopic properties for practical purposes and ignores the differences between the enormous number of microscopic states. Similarly, our concept of a visual pattern must be defined for a purpose. The purpose is reflected in the selection of some “sufficient” statistics  $\mathbf{h}(s)$ . That is, depending on a visual task, we are only interested in some global (macro)properties  $\mathbf{h}(s)$ , and ignore the differences between image instances within the set. This was clearly the case in Julesz psychophysics experiments in the 1960s-1970s on texture discrimination [46]. Thus, we define a visual concept in the same way as the microcanonical ensemble.

**Definition 2 (Homogeneous Visual Patterns).** For any homogeneous visual pattern  $v$  defined on a lattice or graph  $\Lambda$ , let  $s$  be the visual representation (e.g.,  $s = \mathbf{I}$ ) and  $\mathbf{h}(s)$  a list of sufficient feature statistics, then a pattern  $v$  is equal to a maximum set (or equivalence class), as  $\Lambda$  goes to infinity in the von Hove sense,

$$\text{A pattern } v = \Omega(\mathbf{h}_o) = \{s_\Lambda : \mathbf{h}(s) = \mathbf{h}_o, \Lambda \rightarrow \infty\}. \quad (17)$$

As  $\Lambda$  goes to infinity and the pattern is homogeneous, the statistical fluctuations and the boundary condition effects both diminish. It makes sense to put the constraints  $\mathbf{h}(s) = \mathbf{h}_o$ .

In the literature, a texture pattern was first defined as a Julesz ensemble by [97]. This can be easily extended to any patterns, including, generic images, texture, smooth surfaces, texton process, etc.

The connections between the three physical ensembles also reveals an important duality between a *descriptive constraints*  $\mathbf{h}(s) = \mathbf{h}_o$  in the deterministic set  $\Omega_{\text{mch}}(\mathbf{h}_o)$  and the parameters  $\beta$  in Gibbs model  $p_{\text{Gib}}$ . In vision, the duality is between the image statistics  $\mathbf{h}_o = (\mathbf{h}_1^{\text{obs}}, \dots, \mathbf{h}_K^{\text{obs}})$  in (6) and the parameters of the descriptive models  $\beta = (\lambda_1, \dots, \lambda_K)$ .

The connection between set  $\Omega(\mathbf{h}_o)$  and the descriptive model  $p(s; \beta)$  is restated by the theorem below [87].

**Theorem 3 (Ensemble Equivalence).** For visual signals  $s_\Lambda \in \Omega(\mathbf{h}_o)$  on large (or infinity) lattice (or graph)  $\Lambda$ , then on any small lattice  $\Lambda_o \subset \Lambda$ , the signal  $s_{\Lambda_o}$  given its neighborhood  $s_{\partial\Lambda_o}$  follows a descriptive model  $p(s_{\Lambda_o} | s_{\partial\Lambda_o}; \beta)$ .

The duality between  $\beta$  and  $\mathbf{h}_o$  is reflected by the maximum entropy constraints  $E_{p(s; \beta)}[\mathbf{h}(s)] = \mathbf{h}_o$ . More precisely, it is stated in the following theorem [87].

**Theorem 4 (Model and Concept Duality).** Let  $p(s_\Lambda; \beta)$  be a descriptive model of a pattern  $v$ , and  $\Omega_\Lambda(\mathbf{h})$  the set for pattern  $v$ , and let  $\psi(\mathbf{h})$  and  $\rho(\beta)$  be the entropy function and pressure defined as

$$\psi(\mathbf{h}) = \lim_{\Lambda \rightarrow \infty} \frac{1}{|\Lambda|} \log |\Omega_\Lambda(\mathbf{h})|, \quad \text{and} \quad \rho(\beta) = \lim_{\Lambda \rightarrow \infty} \frac{1}{|\Lambda|} \log Z(\beta).$$

If  $\mathbf{h}_o$  and  $\beta_o$  correspond to each other, then

$$\phi'(\mathbf{h}_o) = \beta_o, \quad \text{and} \quad \rho'(\beta_o) = \mathbf{h}_o,$$

in the absence of phase transition.

For visual patterns on finite graphs, such as a human face or a 2D shape of animal, the definition of pattern is given below.

**Definition 3 (Finite Patterns).** For visual pattern  $v$  on a finite lattice or graph  $\Lambda$ , let  $s$  be the representation and  $\mathbf{h}(s)$  its sufficient statistics, and the visual concept is an ensemble governed by a maximum entropy probability  $p(s; \beta)$ ,

$$\text{pattern } v = \Omega(\mathbf{h}_o) = \{(s, p(s; \beta)) : E_p[\mathbf{h}(s)] = \mathbf{h}_o\}. \quad (18)$$

Each pattern instance  $s$  is associated with a probability  $p(s; \beta)$ .

The ensemble is a set with each instance assigned a probability. Obviously, (17) is a special case of (18). That is, when  $\Lambda \rightarrow \infty$ , one homogeneous signal is enough to compute the expectation, i.e.,  $E_p[\mathbf{h}(s)] = \mathbf{h}(s)$ . The limit of  $p(s; \beta)$  is the uniform probability  $p_{\text{unif}}(s; \mathbf{h}_o)$  as  $\Lambda \rightarrow \infty$ .

The probabilistic notion in defining finite visual signal is the root for errors in recognition, segmentation, and grouping. On any finite graph, the ensembles for two different patterns will overlap and the ability of distinguishing two patterns is limited by the Chernoff information that measures the distances of the two distributions. Some in-depth discussions on the relationship between performance bounds and models are referred to the order parameter theory [91].

To conclude this section, we have the following equivalence for conceptualization of visual pattern.

$$\text{A visual pattern } v \longleftrightarrow \mathbf{h} \longleftrightarrow \beta \in \Omega_K^d.$$

## 6 GENERATIVE MODELING

In this section, we revisit the general MLE learning formulated in (2), (3), and (4) and review some progress in generative models of visual patterns and the integration with descriptive models.

### 6.1 The Basic Principle of Generative Modeling

Descriptive models are built on features and statistics extracted from the signal and use complex potential functions to characterize visual patterns. In contrast, generative models introduce hidden (latent) variables to account for the generating process of large image structures.

For simplicity of notation, we assume  $L$ -levels of hidden variables which generate image  $\mathbf{I}$  in a linear order. At each level,  $W_i$  generates  $W_{i-1}$  with a dictionary (vocabulary)

$\mathcal{D}_i, i = 1, \dots, L$ . The dictionary is a set of description, such as image bases, textons, parts, templates, lighting functions, etc.

$$W_L \xrightarrow{\mathcal{D}_L} W_{L-1} \xrightarrow{\mathcal{D}_{L-1}} \dots \xrightarrow{\mathcal{D}_2} W_1 \xrightarrow{\mathcal{D}_1} \mathbf{I}. \quad (19)$$

Let  $p(W_{i-1}|W_i; \mathcal{D}_i, \beta_{i-1})$  denote the conditional distribution for pattern  $W_{i-1}$  given  $W_i$ , with  $\beta_{i-1}$  being the parameter of the model. Then, by summing over the hidden variables, we have an image model,

$$p(\mathbf{I}; \Theta) = \sum_{W_L} \dots \sum_{W_1} p(\mathbf{I}|W_1; \mathcal{D}_1, \beta_0) p(W_1|W_2; \mathcal{D}_1, \beta_1) \dots p(W_{L-1}|W_L; \mathcal{D}_L, \beta_{L-1}). \quad (20)$$

$\Theta = (\mathcal{D}_1, \dots, \mathcal{D}_L; \beta_0, \dots, \beta_{L-1})$  are the parameters, and each conditional probability is often a descriptive model specified by  $\beta_i$ .

By analogy to speech, the observable image  $\mathbf{I}$  is like the speech wave form. Then, the first-level dictionary  $\mathcal{D}_1$  is like the set of *phonemes* and  $\beta_1$  parameterizes the transition probability between phonemes. In image model,  $\mathcal{D}_1$  is a set of image bases like Gabor wavelets. The second-level dictionary  $\mathcal{D}_2$  is like the set of *words*, each being a short sequences of phonemes in  $\mathcal{D}_1$ , and  $\beta_2$  parameterizes the transition probability between words. In image models,  $\mathcal{D}_2$  is the set of textons. Going up the hierarchy, we need dictionaries like the *grammatical reproduction rules* for phrases and sentences in language and probabilities for how frequently each reproduction rule is used, etc.

A hidden variable  $W_i$  is fundamentally different from an image feature  $\phi_i$  in descriptive models, though they may be closely related.  $W_i$  is a *random variable* that should be inferred from images, while  $\phi_i$  is a *deterministic transform* of images.

Following the ML-estimate in (2), one can learn the parameters  $\Theta$  in  $p(\mathbf{I}; \Theta)$  by EM-type algorithm, like stochastic gradients [39]. Take derivative of the log-likelihood with respect to  $\Theta$ , and set  $\frac{d \log p(\mathbf{I}; \Theta)}{d\Theta} = 0$ , one gets

$$0 = \sum_{W_L} \dots \sum_{W_1} \left[ \frac{\partial \log p(\mathbf{I}|W_1; \mathcal{D}_1, \beta_0)}{\partial (\mathcal{D}_1, \beta_0)} + \dots + \frac{\partial \log p(W_{L-1}|W_L; \mathcal{D}_L, \beta_{L-1})}{\partial (\mathcal{D}_L, \beta_{L-1})} \right] \times p(W_1|\mathbf{I}; \mathcal{D}_1, \beta_0) \dots p(W_L|W_{L-1}; \mathcal{D}_L, \beta_{L-1}). \quad (21)$$

In theory, these equations can be solved with global optimum by iterating two steps [39]:

1. The E-type step. Making inferences about the hidden variables by sampling from a sequence of posteriors,

$$\begin{aligned} W_1 &\sim p(W_1|\mathbf{I}; \mathcal{D}_1, \beta_0), \quad \dots, \\ W_L &\sim p(W_L|W_{L-1}; \mathcal{D}_L, \beta_{L-1}). \end{aligned} \quad (22)$$

Then, we can approximate the summation (integration) by importance sampling.

2. The M-type step. Given the samples, one optimizes the parameters  $\Theta$ . The learning results in  $\Theta$  includes the visual dictionaries  $\mathcal{D}_1, \dots, \mathcal{D}_L$  and the descriptive models  $\beta_0, \dots, \beta_{L-1}$  that govern their spatial layouts of the hidden structures. It is beyond this review to discuss the algorithm.

## 6.2 Some Examples of Generative Models

Now, we review a spectrum of generative image models, starting again with a model for the  $1/f$ -power law.

### 6.2.1 Model G1: A Generative Model for the $1/f$ -power Law of Natural Images

The  $1/f$ -law of the Fourier amplitude in natural images was analytically modeled by a Gaussian MRF  $p_{1/f}$  (see (8)). We transform (8) into the Fourier domain, thus

$$p_{1/f}(\mathbf{I}; \beta) = \frac{1}{Z} \exp \left\{ - \sum_{\xi, \eta} \beta (\xi^2 + \eta^2) |\hat{\mathbf{I}}(\xi, \eta)|^2 \right\}. \quad (23)$$

The Fourier bases are the independent components for the Gaussian ensemble governed by  $p_{1/f}$ . From the above Gaussian model, one obtains a two-layer generative model [65],

$$\mathbf{I}(x, y) = \sum_{\xi} \sum_{\eta} \frac{1}{2\beta(\xi^2 + \eta^2)} a(\xi, \eta) e^{2\pi i \frac{x\xi + y\eta}{N}}, \quad a(\xi, \eta) \sim N(0, 1). \quad (24)$$

The dictionary  $\mathcal{D}_1$  is the Fourier basis, and the hidden variables are the Fourier coefficients  $a(\xi, \eta) \forall \xi, \eta$  which are iid normal distributed. Only two parameters are used in  $\beta = (0, 1)$  for specifying the normal density. Therefore,

$$W_1 = \{a(\xi, \eta) : \forall \xi, \eta\} \text{ and } \mathcal{D}_1 = \{ \mathbf{b}(\mathbf{I}; \xi, \eta) = e^{2\pi i \frac{x\xi + y\eta}{N}} : \forall \xi, \eta \}.$$

One can sample a random image  $\mathbf{I} \sim p_{1/f}(\mathbf{I}; \beta)$ , according to (24) it's:

1. drawing the iid Fourier coefficients and
2. generating the synthesis image  $\mathbf{I}$  by linear superposition of the Fourier bases.

A result is displayed in Fig. 4b.

To the author's knowledge, this is the only image model whose descriptive and generative versions are analytically transferable. Such happy endings perhaps only occur in the Gaussian family!

In the literature, Ruderman [73] explains the  $1/f$ -law by an occlusion model. It assumes that image  $\mathbf{I}$  is generated by a number of independent "objects" (rectangles) of size subject to a cubic law  $1/r^3$ . A synthesis image is shown in Fig. 12a.

### 6.2.2 Model G2: A Generative Model for Scale-Invariant Gradient Histograms

The scale-invariance of gradient histograms in natural images inspired a number of research for generative models in parallel with the descriptive model  $p_{inv}$ . The objective is to search for some "laws" that governs the distribution of objects in natural scenes.

The first is the random collage model [53], which is also called the "dead leaves" model (see Stoyan et al. [80]). It assumes that an image is generated by a number of  $n$  opaque disks. Each disk is represented by hidden variables  $x, y, r, \alpha$  for center, radius, and intensity, respectively.

$$\begin{aligned} W_1 &= (n, \{x_i, y_i, r_i, \alpha_i\} : i = 1, 2, \dots, n), \\ \mathcal{D}_1 &= \{disk(\mathbf{I}; x, y, r) : \forall (x, y) \in \Lambda, r \in [r_{min}, r_{max}]\}. \end{aligned}$$

The dictionary  $\mathcal{D}_1$  includes disk templates at all possible sizes and locations. Therefore, let  $a \oslash b$  denote that  $a$  occludes  $b$ , the image is generated by

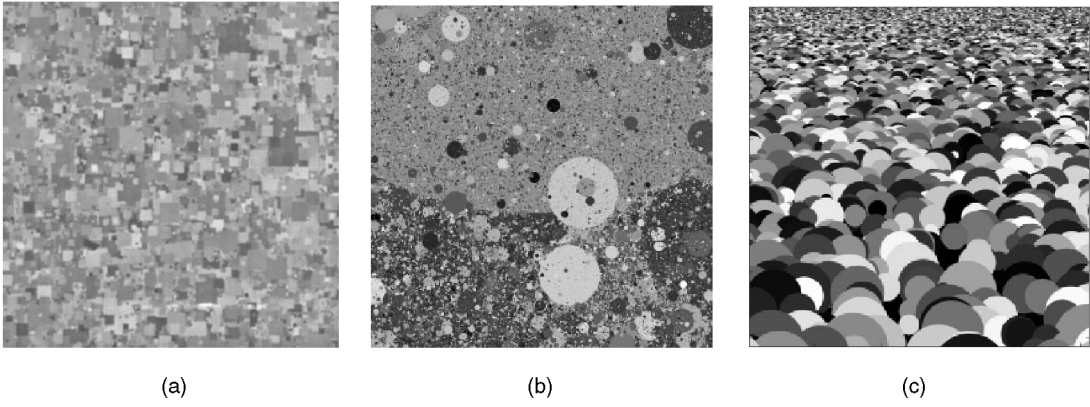


Fig. 12. Synthesized images from three generative models. (a) Ruderman [73], (b) Lee et al. [53], and (c) Chi [13]. See text for explanations.

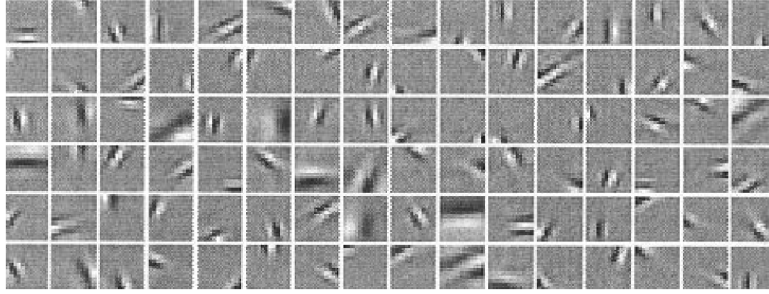


Fig. 13. Some of the linear bases (dictionary) learned from natural images by Olshausen and Field [67].

$$\mathbf{I} = \text{disk}(x_n, y_n, r_n, \alpha_n) \circ \text{disk}(x_{n-1}, y_{n-1}, r_{n-1}, \alpha_{n-1}) \circ \cdots \circ \text{disk}(x_1, y_1, r_1, \alpha_1). \quad (25)$$

Lee et al. [53] showed that, if  $p(n)$  is Poisson distributed, and the disk location  $(x, y)$  and intensity  $\alpha$  are iid uniform distributed, and the radius  $r_i$  subject to a  $1/r^3$ -law,

$$p(r) = c/r^3, \quad \text{for } r \in [r_{\min}, r_{\max}]. \quad (26)$$

Then, the generative model  $p(\mathbf{I}; \Theta)$  has scale invariance gradient histograms. Fig. 12b shows a typical image sampled from this model.

The second model is studied by Chi [13]. This offers a beautiful 3D generative explanation. It assumes that the disks (objects) are sitting vertically on a 2D plane (the ground) facing the viewer. The sizes of the disks are iid uniformly distributed and they have proven that the 2D projected (by perspective projection) sizes of the objects then follow the  $1/r^3$  law in (26). The locations and intensities are iid uniformly distributed like the random collage model. A typical image sampled from this model is shown in Fig. 12c. More rigorous studies and laws along this vein are in [66]. These results put a reasonable explanation for the origin of scale invariance in natural images. Nevertheless, these models are all biased by the object elements they choose, as they are not maximum entropy models, in comparison with  $p_{\text{inv}}(\mathbf{I})$ .

### 6.2.3 Model G3: Generative Model for Sparse Coding: Learning the Dictionary

In research stream 2 (image coding, wavelets, image pyramids, ICA, etc.) discussed in Section 2.1, a linear additive model is widely assumed and an image is a superposition of some local image bases from a dictionary plus a Gaussian noise image  $\mathbf{n}$ .

$$\mathbf{I} = \sum_i^n \alpha_i \cdot \psi_{\ell_i, x_i, y_i, \tau_i, \sigma_i} + \mathbf{n}, \quad \psi_i \in \mathcal{D}, \forall i. \quad (27)$$

$\psi_\ell$  is a base function, for example, Gabor, Laplacian of Gaussian, etc. It is specified by hidden variables  $x_i, y_i, \tau_i, \sigma_i$  for position, orientation, and scale. Thus, a base is indexed by hidden variables  $b_i = (\ell_i, \alpha_i, x_i, y_i, \tau_i, \sigma_i)$ . The hidden variables and dictionary are

$$W_1 = (n, \{b_i : i = 1, 2, \dots, n\}; \mathbf{n}), \\ \mathcal{D}_1 = \{\psi_\ell(x, y, \tau, \sigma) : \forall x, y, \tau, \sigma, \ell\}.$$

$x_i, y_i, \tau_i, \sigma_i$  are assumed iid uniformly distributed, and the coefficients  $\alpha_i \sim p(\alpha)$ ,  $\forall i$  follow an iid Laplacian or mixture of Gaussian for sparse coding,

$$p(\alpha) \sim \exp\{-|\alpha|/c\} \quad \text{or} \quad p(\alpha) = \sum_{j=1}^2 \omega_j N(\alpha, \sigma_j). \quad (28)$$

According to the theory of generative model (Section 6.1), one can learn the dictionary from raw images in the  $M$ -step. Olshausen and Field [67] used the sparse coding prior  $p(\alpha)$  learned a set of  $144 = 12 \times 12$  pixels bases, some of which are shown in Fig. 13. Such bases capture some image structures and are believed to bear resemblance to the responses of simple cells in V1 of primates.

### 6.2.4 Model G4: A Generative Model for Texton and Texture

In the previous three generative models, the hidden variables are assumed to be iid distributed. Such distributions can be viewed as degenerated descriptive models. But obviously these variables and objects are not iid, and sophisticated descriptive models are needed for the spatial relationships between the image bases or objects.

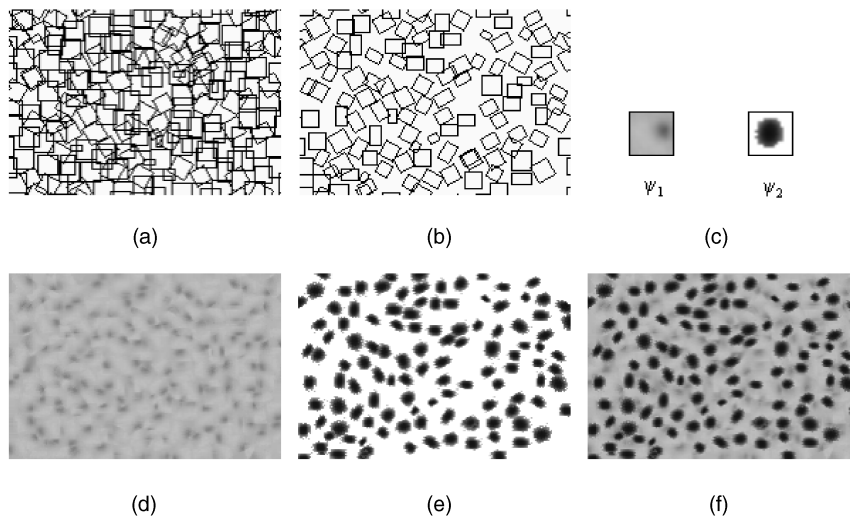


Fig. 14. An example of integrating descriptive texton model and a generative model for a cheetah skin pattern. (a) Sampled texton map  $T_1$ . (b) Sampled texton map  $T_2$ . (c) Templates. (d) Layer I  $I(T_1, \psi_1)$ . (e) Layer II  $I(T_2, \psi_2)$ . (f) Synthesized image. After Guo et al. [40].

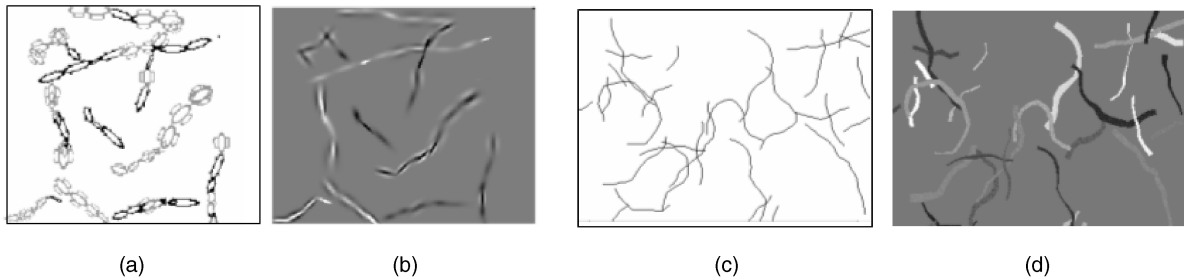


Fig. 15. (a) and (b) A rope model is a Markov chain of knots and each knot has 1-3 image bases shown by the ellipses. (c) and (d) The smooth curve model on intensity. After Tu and Zhu [83].

The first work that integrates the descriptive and generative model was presented in [40] for texture modeling. It assumes that a texture image is generated by two levels (a foreground and a background) of hidden texton processes plus a Gaussian noise. Fig. 14 show an example of cheetah skin pattern. Figs. 14a and 14b shows two texton patterns  $T_1$ ,  $T_2$ , which are sampled from descriptive textons models  $p_{\text{txn}}(\mathbf{T}; \beta_{o,1}, \beta_1)$  and  $p_{\text{txn}}(\mathbf{T}; \beta_{o,2}, \beta_2)$ , respectively. The models are learned from an observed cheetah skin (raw pixel) image. Each texton is symbolically illustrated by an oriented window. Then, two base functions  $\psi_1, \psi_2$  are learned from images and shown in Fig. 14c. The two image layers are shown in Figs. 14d and 14e. The superposition (with occlusion) of the two layers renders the synthesized image in Fig. 14f. More examples and discussions are referred to in [40].

### 6.2.5 Model G5: A Generative Rope Model of Curve Processes

A three-layer generative model for curve, called a ‘‘rope model,’’ was studied by Tu and Zhu [83]. The model extends the descriptive model for SNAKE and Elastica  $p_{\text{snk}}$  and  $p_{\text{els}}$  by integrating it with base and intensity representation.

Fig. 15a shows a sketch of the rope model that is a Markov chain of knots. Each knot  $\zeta$  has 1-3 linear bases, for example, difference of Gaussian (DoG), and difference of offset Gaussians (DooG) at various orientations and scales

$$W_2 = (n, \zeta_1, \zeta_2, \dots, \zeta_n), \text{ with}$$

$$\zeta_i = (\alpha_{ij}, \ell_{ij}, x_{ij}, y_{ij}, \tau_{ij}, \sigma_{ij})_{j=1}^k, \quad k \leq 3,$$

$$W_1 = (N, \{b_{ij} : i = 1, 2, \dots, n; j = 1, \dots, 3\}).$$

Fig. 15b shows a number of random curves (image not pure geometry) sampled from the rope model. The image  $I$  is the linear sum of the bases in  $W_1$ .

This additive model is insufficient for occlusion, etc. Figs. 15c and 15d show a occlusion type curve model. Each curve is a SNAKE/Elastica type Markov chain model with width and intensity at each point. Fig. 15c is the sampled curve skeleton and Fig. 15d is the image. Smoothness are assumed for both geometry, width, and intensity.

### 6.2.6 Summary

The generative models used in vision are still preliminary and they often assume a degenerated descriptive model for the hidden variables. To develop richer generative models, one needs to integrate generative and descriptive models.

## 7 CONCEPTUALIZATION OF PATTERNS AND THEIR PARTS: REVISITED

With generative models, we now revisit the conceptualization of visual patterns in a more general setting.

In Section 5.2, a visual pattern  $v$  with representation  $s$  is equalized to a statistical ensemble governed by a model  $p(s; \beta)$  or, equivalently, a statistical description  $h_o$ . In reality, the representation  $s$  is given in a supervised way and is not



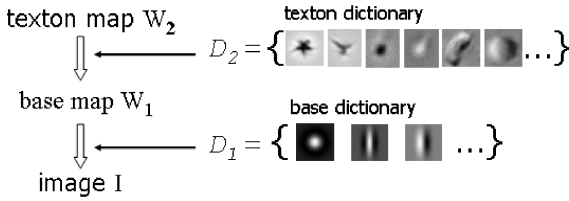


Fig. 16. A three-level generative image model with textons. Modified from Zhu et al. [99].

observable unless  $s$  is an image. Thus, we need to define visual concepts based on images so that they can be learned and verified from observable data.

Following the notation in Section 6.1, we have the following definition extending from Definition 3.

**Definition 4 (Visual Pattern).** A visual pattern  $v$  is a statistical ensemble of image  $\mathbf{I}$  governed by a generative model  $p(\mathbf{I}; \Theta^v)$  with  $L$  layers,

$$\text{pattern } v = \Omega(\Theta^v) = \{ (\mathbf{I}, p(\mathbf{I}; \Theta^v)) : \Theta^v \in \Omega_K^g \},$$

where  $p(\mathbf{I}; \Theta^v)$  is defined in (20).

In this definition, a pattern  $v$  is identified by a vector of parameters in the generative family  $\Omega_K^g$ , which include the  $L$  dictionaries and  $L$  descriptive models,

A visual pattern  $v \longleftrightarrow \Theta^v = (\mathcal{D}_1^v, \dots, \mathcal{D}_L^v, \beta_0^v, \dots, \beta_{L-1}^v) \in \Omega_K^g$ .

By analogy to speech,  $\Theta^v$  defines the whole language system, say  $v = \text{English}$  or  $v = \text{Chinese}$ , and it includes all the hierarchic descriptions from waveforms to phonemes, and to sentences—both the vocabulary and models.

Therefore, the definition of many intuitive but vague concepts, such as textons, meaningful parts of shape, etc., must be defined in the context of a generative model  $\Theta$ . It is meaningless to talk about a texton or part without a generative image model.

**Definition 5 (Visual Vocabulary).** A visual vocabulary, such as textons, meaningful parts of shape, etc. are defined as an element in the dictionaries  $\mathcal{D}_i, i = 1, \dots, L$  associated with the generative model of natural images  $p(\mathbf{I}; \Theta)$ .

To show some recent progress, we show a three-level generative model for textons in Fig. 16. It assumes that an image  $\mathbf{I}$  is generated by a linear superposition of bases  $W_1$  in (28). These bases are, in turn, generated by a smaller number of textons  $W_2$ . Each texton is a deformable template consisting of a few bases in a graph structure. The dictionary  $\mathcal{D}_1$  includes a number of base functions, such as Laplacian of Gaussian, Gabor, etc. They like the phonemes in speech. The dictionary  $\mathcal{D}_2$  includes a larger number of texton templates. Each texton in  $\mathcal{D}_2$  represents a small iconic object at distance, such as stars, birds, cheetah blobs, snowflakes, beans, etc. It is expected that natural images have levels of vocabularies with sizes  $|\mathcal{D}_1| = O(10)$  and  $|\mathcal{D}_2| = O(10^3)$ . These must be learned from natural images.

## 8 VARIANTS OF DESCRIPTIVE MODELS

In this section, we review the third category of models that are two variants of descriptive models—causal MRF and pseudodescriptive models. These variants are most popular due to

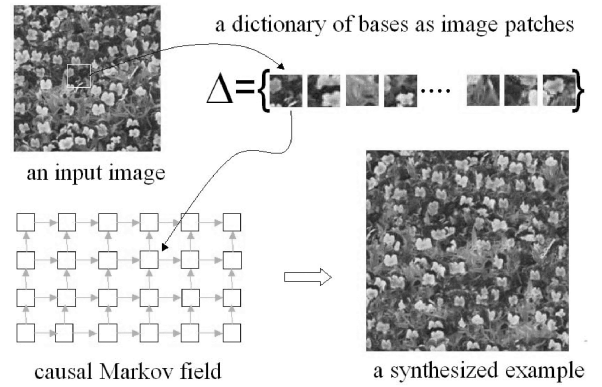


Fig. 17. A causal MRF model for example-based texture synthesis [26], [27], [54].

their computational convenience. However, people should be aware of their limitations and use them with caution.

### 8.1 Causal Markov Models

Let  $s = (s_1, \dots, s_n)$  be the representation of a pattern. As Fig. 2b illustrates, a causal Markov model imposes a partial order in the vertices and, thus, factorizes the joint probability into a product of conditional probabilities,

$$p_{\text{cau}}(s; \beta) = \prod_{i=1}^n p(s_i | \text{parent}(s_i); \beta_i). \quad (29)$$

$\text{parent}(s_i)$  is the set of parent vertices which point to  $s_i$ . Though the graph is directed in syntax, this is not a generative model because the variables are at the same semantic level.  $p_{\text{cau}}(s)$  can be derived from the maximum entropy learning scheme in Section 4.1.

$$p_{\text{cau}}^* = \arg \max_s - \sum_s p_{\text{cau}}(s) \log p_{\text{cau}}(s).$$

Thus,  $p_{\text{cau}}(s; \beta)$  is a special class of descriptive model. When the dimension of  $p(s_i | \text{parent}(s_i))$  is not high, (e.g.,  $|\text{parent}(s_i)| + 1 \leq 4$ ), the conditional probability is often estimated by a nonparametric Parzen window.

There are many causal Markov models for texture in the 1980s and early 1990s (See Popat and Picard [70] and references therein). In the following, we review two pieces of interesting work that appeared recently.

One is the work on example-based texture synthesis by Efros and Leung [26], Liang et al. [54], and Efros and Freeman [27]. Hundreds of realistic textures can be synthesized by a patching technique. Fig. 17 reformulates the idea in a causal Markov model. An example texture image is first chopped into a number of image patches of a predefined size. These patches form a vocabulary  $\mathcal{D}_1 = \Delta$  of image “bases” specific to this texture. Then, a causal Markov field is set up with each element being chosen from  $\Delta$  conditional on two other previous patches (left and below). The patches are pasted one by one in a linear order by sampling from a nonparametric conditional distribution. A synthesized image is shown to the lower-right side. The vocabulary  $\Delta$  greatly reduces the search space and, thus, the causal model can be simulated extremely fast. The model is biased by the dictionary and the causality assumption.

Another causal Markov model was proposed by [88]. Wu et al. represent an image by a number of bases from a

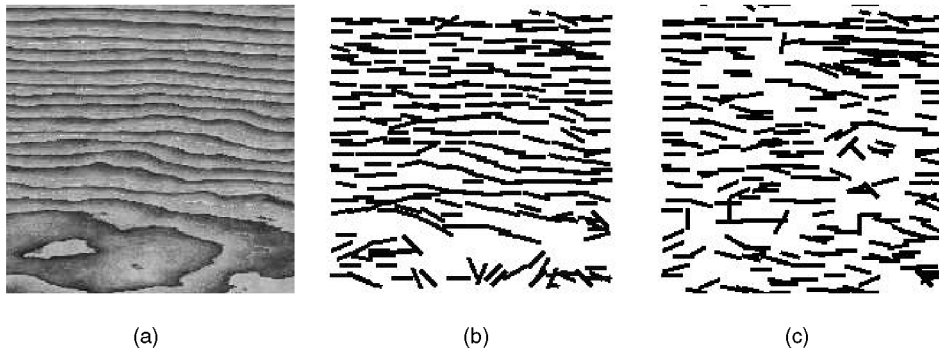


Fig. 18. A causal Markov model for texture sketch. (a) Input, (b) image sketch, and (c) a synthesized sketch. After [88].

generic base dictionary (Log. DoG. DooG) as in sparse coding model. Each base is then symbolically represented by a line segment, as Figs. 18a and 18b show. This forms a base map similar to the texton (attributed point) pattern in Fig. 7. Then, a causal model is learned based on Fig. 18b for the base map. The graph structure is more flexible than the grid in Fig. 17. A random sample is drawn from the model and shown in Fig. 18c.

## 8.2 Pseudodescriptive Models

While causal Markov models approximate the Gibbs distributions  $p_{\text{des}}$  and have sound probabilities  $p_{\text{cau}}$ , the second variant, called pseudodescriptive model in this paper, approximates the Julesz ensemble.

For example, the texture synthesis work by Heeger and Bergen [42] and De Bonet and Viola [21] belong to this family. Given an observed image  $\mathbf{I}^{\text{obs}}$  on a large lattice  $\Lambda$ , suppose a number of  $K$  filters  $F_1, F_2, \dots, F_K$  are chosen, say Gabors at various scales and orientations. Convolution of the filters with image  $\mathbf{I}^{\text{obs}}$ , one obtains a set of filter responses

$$S^{\text{obs}} = \{F_i^{\text{obs}}(x, y) = F_i * \mathbf{I}^{\text{obs}}(x, y) : i = 1, 2, \dots, K, (x, y) \in \Lambda\}.$$

Usually,  $K > 30$  and, thus,  $S^{\text{obs}}$  is a very redundant representation of  $\mathbf{I}^{\text{obs}}$ . In practice, to reduce the dimensionality and computation, these filter responses are organized in a pyramid representation with low-frequency filters subsampled (see Fig. 19).

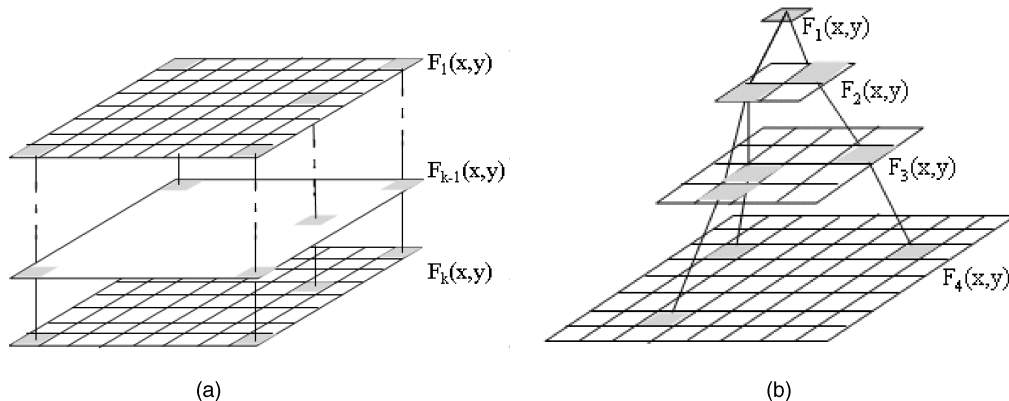


Fig 19. (a) Extracting feature vectors ( $F_1(x, y), \dots, F_K(x, y)$ ) for every pixels in a lattice and, thus, obtain  $K|\Lambda|$  filter responses. (b) Extracting the feature vectors in a pyramid. See Heeger and Bergen [42] and De Bonet and Viola [21].

Let  $\mathbf{h}^{\text{obs}} = \mathbf{h}(\mathbf{I}^{\text{obs}}) = (\mathbf{h}_1^{\text{obs}}, \dots, \mathbf{h}_K^{\text{obs}})$  be the  $K$  marginal histograms of the filter responses. A Julesz ensemble (or texture) is defined by  $\Omega(\mathbf{h}^{\text{obs}}) = \{\mathbf{I} : \mathbf{h}(\mathbf{I}) = \mathbf{h}^{\text{obs}}\}$ . Heeger and Bergen [42] sampled the  $K \cdot |\Lambda|$  filter responses independently according to  $\mathbf{h}^{\text{obs}}$ , which is computationally very convenient. Obviously, the sampled filter responses  $F_i(x, y), i = 1, \dots, K, (x, y) \in \Lambda$  produce histograms  $\mathbf{h}_o$  (or very closely), but these filter responses are inconsistent as they are sampled independently. There is no image  $\mathbf{I}$  that can produce these filter responses. Usually, one finds an image  $\mathbf{I}$  that has least-square error by pseudoinverse. In fact, this employs an image model,

$$p_{\text{psdes}}(\mathbf{I}) \propto \exp \left\{ - \sum_{i=1}^K \sum_{(x,y) \in \Lambda} (F_i^{\text{syn}}(x, y) - F_i * \mathbf{I}(x, y))^2 / \sigma^2 \right\},$$

$$F_i^{\text{syn}}(x, y) \stackrel{\text{iid}}{\sim} \mathbf{h}_i^{\text{obs}}, \forall i, \forall (x, y). \quad (30)$$

Of course, the image computed by pseudoinverse usually does not satisfy  $\mathbf{h}(\mathbf{I}) = \mathbf{h}^{\text{obs}}$ . So, we call it a “pseudodescriptive” model. The work by De Bonet and Viola [21] was done in the same principle, but it used a  $K$ -dimensional joint histogram for  $\mathbf{h}(\mathbf{I})$ . As  $K$  is very high in their work (say,  $K = 128$ ), sampling the joint histogram is almost equal to shuffling the observed image.

In a descriptive model or Julesz ensemble, the number of constraints in  $\mathbf{h}(\mathbf{I}) = \mathbf{h}^{\text{obs}}$  is much lower than the image

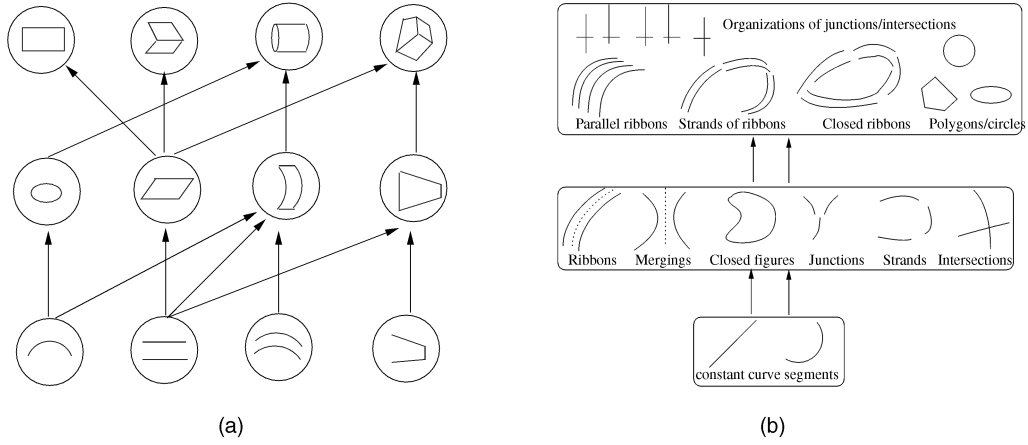


Fig. 20. Hierarchic perceptual grouping. (a) After Dickinson et al. [24]. (b) After Sarkar and Boyer [75].

pixels  $\Lambda$ . In contrast, a pseudodescriptive model puts  $K|\Lambda|$  constraints and produces an empty set.

## 9 DISCRIMINATIVE MODELS

Many perceptual grouping work (research stream 3) fall in category 4—discriminative models. In this section, we briefly mention some typical work and then focus on the theoretical connections between discriminative models to the descriptive and generative models. A good survey of grouping literature is given in [9].

### 9.1 Some Typical Discriminative Models

The objective of perceptual grouping is to compose image elements into larger and larger structures in a hierarchy. Fig. 20 shows two influential works in the literature. Dickinson et al. [24] adopted a hierarchic Bayesian network for grouping short line and curve segments into generic object facets, and the latter are further grouped into 2D views of 3D object parts. Sarkar and Boyer [75] used the Bayesian network for grouping edge elements into hierarchic structures in aerial images. More recent work is Amir and Lindenbaum [1].

If we represent the hierarchic representation by a linear order for ease of discussion, the grouping proceeds in the inverse order of the generative model (see (19), Fig. 2).

$$\mathbf{I} \rightarrow W_1 \rightarrow W_2 \rightarrow \dots \rightarrow W_L. \quad (31)$$

As the grouping must be done probabilistically, both Dickinson et al. [24] and Sarkar and Boyer [75] adopted a list of conditional probabilities in their Bayesian networks. Reformulated in the above notation, they are,

$$q(W_1|\mathbf{I}), \quad q(W_2|W_1), \quad \dots, \quad q(W_L|W_{L-1}).$$

Again, we use linear order here for clarity. There may be expressways for computing objects from edge elements directly, such as generalized Hough transform. In the literature, most of these conditional probabilities are manually estimated or calculated in a similar way to [56].

### 9.2 The Computational Role of Discriminative Models

The discriminative models are effective and useful in vision and pattern recognition. However, there are a number of conceptual problems suggesting that they should perhaps

not be considered *representational models*, instead they are *computational heuristics*. In the desk example of Fig. 2, the presence of a leg may, as a piece of evident, “suggests” the presence of a desk but it does not “cause” a desk. A leg can also suggest chairs and a dozen other types of furniture that have legs. It is the desk concept that causes four legs and a top at various configurations in the generative model.

What is wrong with the inverted arrows in discriminative models? A key point associated with Bayes (causal, belief) networks is the idea of “explaining-away” or “lateral inhibition” in a neuroscience term. If there are multiple competing causes for a symptom, then the recognition of one cause will suppress the other causes. In a generative model, if a leg is recognized as belonging to a desk during computation, then the probability of a chair at the same location is reduced drastically. But, in a discriminative model, it appears that the four legs are competing causes for the desk, then one leg should drive away the other three legs in explanation! This is not true. Without the guidance of generative model, the discriminative methods could create combinatorial explosions.

In fact, the discriminative models are approximations to the posteriors,

$$\begin{aligned} q(W_1|\mathbf{I}) &\sim p(W_1|\mathbf{I}; \mathcal{D}_1, \beta_0), \quad \dots, \\ q(W_L|W_{L-1}) &\sim p(W_L|W_{L-1}; \mathcal{D}_L, \beta_{L-1}). \end{aligned} \quad (32)$$

Like most pattern recognition methods, the approximative posteriors  $q(\cdot)$ s use only local deterministic features at each level for computational convenience. For example, suppose  $W_1$  is an edge map, then it is usually assumed that  $q(W_1|\mathbf{I}) = q(W_1|\Phi_1(\mathbf{I}))$  with  $\Phi_1(\mathbf{I})$  being some local edge measures [50]. For the other levels,  $q(W_{i+1}|W_i) = q(W_{i+1}|\Phi_i(W_i))$  with  $\Phi_i(W_i)$  being some *compatibility functions and metrics* [75], [7].

For ease of notation, we only consider one level of approximation:  $q(W|\mathbf{I}) = q(W|\Phi(\mathbf{I})) \approx p(W|\mathbf{I}; \mathcal{D}, \beta)$ . By using local and deterministic features, information is lost in each approximation. The amount of information loss is measured by the Kullback-Leibler divergence. Therefore, the best set of features is chosen to minimize the loss.

$$\begin{aligned} \Phi^* &= \arg \min_{\Phi \in \text{Bank}} KL(p||q) \\ &= \arg \min_{\Phi \in \text{Bank}} \sum_W p(W|\mathbf{I}; \mathcal{D}, \beta) \log \frac{p(W|\mathbf{I}; \mathcal{D}, \beta)}{q(W|\Phi(\mathbf{I}))}. \end{aligned}$$

Now, we have the following theorem for what are most discriminative features.<sup>5</sup>

**Theorem 5.** For linear features  $\Phi$ , the divergence  $KL(p \parallel q)$  is equal to the mutual information between variables  $W$  and image  $\mathbf{I}$  minus the mutual information between  $W$  and  $\Phi(\mathbf{I})$ .

$$KL(p(W|\mathbf{I}; \mathcal{D}, \beta) \parallel q(W|\Phi(\mathbf{I}))) = MI(W, \mathbf{I}) - MI(W, \Phi(\mathbf{I})).$$

$MI(W, \mathbf{I}) = MI(W, \Phi(\mathbf{I}))$  if and only if  $\Phi(\mathbf{I})$  is the sufficient statistics for  $W$ .

This theorem leads to a *maximum mutual information principle* for discriminative feature selection and it is different from the most informative feature for descriptive models.

$$\begin{aligned} \Phi^* &= \arg \max_{\Phi \in \text{Bank}} MI(W, \Phi(\mathbf{I})) \\ &= \arg \min_{\Phi \in \text{Bank}} KL(p(W|\mathbf{I}; \mathcal{D}, \beta) \parallel q(W|\Phi(\mathbf{I}))). \end{aligned}$$

The main problem with the discriminative models is that they do not pool global and top-down information in inference. In our opinion, the discriminative models are *importance proposal probabilities* for sampling the true posterior and inferring the hidden variables. Thus, they are crucial in computation for both Bayesian inference and for learning generative models (see the E-step in (22)). In both tasks, we need to draw samples from the posteriors through Markov chain Monte Carlo (MCMC) techniques. The latter need to design some proposal probabilities  $q()$ s to suggest the Markov chain moves.

The convergence of MCMC critically depend on how well  $q()$  approximate  $p()$ . This is stated in the theorem below by Mengersen and Tweedie [58].

**Theorem 6.** Sampling a target density  $p(x)$  by the independence Metropolis-Hastings algorithm with proposal probability  $q(x)$ . Let  $P^n(x_o, y)$  be the probability of a random walk to reach point  $y$  at  $n$  steps from an initial point  $x_o$ . If there exists  $\rho > 0$  such that,

$$\frac{q(x)}{p(x)} \geq \rho, \quad \forall x,$$

then the convergence measured by a  $L_1$  norm distance

$$\|P^n(x_o, \cdot) - p\| \leq (1 - \rho)^n.$$

This theorem, though on a simple case, states the computational role of discriminative model. The idea of using discriminative models, such as edge detection, clustering, Hough transforms, are used in a data-driven Markov chain Monte Carlo (DDMCMC) framework for generic image segmentation, grouping, and recognition [98], Tu and Zhu [83].

## 10 DISCUSSION

The modeling of visual patterns is to pursue a probability model  $p()$  to estimate an ensemble frequency  $f()$  in a sequence of nested probability families which integrate both descriptive and generative models. These models are adapted and augmented in four aspects:

1. learning the parameters of the descriptive models,

5. This proof was given in a unpublished note by Wu and Zhu. A similar conclusion was also given by a variational approach by Wolf and George [86], who sent an unpublished manuscript to Zhu.

2. pursuing informative features and statistics in descriptive models.
3. selecting address variables and neighborhood configurations for the descriptive model, and
4. introducing hidden variables in the generative models.

The main challenge in modeling visual patterns is the choice of models that cannot be answered unless we understand the different purposes of vision.

**What is the ultimate goal of learning? Where does it end?** Our ultimate goal is to find the "best" generative model. Starting from the raw images, each time when we add a new layer of hidden variables, we make progress in discovering the hidden structures. At the end of this pursuit, suppose we dig out all the hidden variables, then we will have a physically-based model which is the ultimate generative model denoted by  $p_{\text{gen}}^*$ . This model cannot be further compressed and we reach the Komogorov complexity of the image ensemble.

For example, the chemical diffusion-reaction equations with a few parameters may be the most parsimonious model for rendering some textures. But, obviously, this is not a model used in human vision. Why didn't human vision pursue such ultimate model? This leads to the second question below.

**How do you choose a generative model from many possible explanations?** There are two extremes of models. At one extreme, Theorem 2 states that the pure descriptive model  $p_{\text{des}}^*$  on raw pixels, i.e., no hidden variables at all, can approximate the ensemble frequency  $f(\mathbf{I})$  as long as we put a huge number of features statistics. At the other extreme end, we have the ultimate generative model  $p_{\text{gen}}^*$  mentioned above. In graphics, there is also a spectrum of models, ranging from image-based rendering to physically-based ray tracing. Certainly, our brains choose a model somewhere between  $p_{\text{des}}^*$  and  $p_{\text{gen}}^*$ .

We believe that the choice of generative models is decided by two aspects. The first is the different purposes of vision for navigation, grasping not just for coding. Thus, it makes little sense to justify models by a simple minimum description length principle or other statistics principles, such as AIC/BIC. The second is the computational effectiveness. It is hopeless to have a quantitative formulation for vision purposes at present. We only have some understanding on the second issue.

A descriptive model uses features  $\Phi()$  which is deterministic and, thus, easy to compute (filtering) in a bottom-up fashion. But, it is very difficult to do synthesis using features. For example, sampling the descriptive model (such as FRAME) is expensive. In contrast, the generative model uses hidden variables  $W$  which has to be inferred stochastically and, thus, expensive to compute (analysis). But, it is easier to do top-down synthesis using the hidden variables. For the two extreme models,  $p_{\text{des}}^*$  is infeasible to sample (synthesis) and  $p_{\text{gen}}^*$  is infeasible to infer (analysis). For example, it is infeasible to infer parameters of a reaction-diffusion equation from observed texture images. The choice of generative model in the brain should make both analysis and synthesis convenient. As vision can be used for many diverse purposes, there will be many models coexist.

**Where do features and hidden variables (i.e., visual vocabulary) come from?** The mathematical principles (minimax entropy or maximum mutual information) can choose "optimal" features and variables from predefined sets, but the creation of these candidate sets often come from three sources: 1) observations in human vision, such as psychology and neuroscience, thus related to purposes of vision, 2) physics

models, or 3) artist models. For example, the Gabor filters and Gestalt laws are found to be very helpful in visual modeling. At present, the visual vocabulary is still far from being enough.

This may sound ad hoc to someone who likes analytic solutions! Unfortunately, we may never be able to justify such vocabulary mathematically, just as physicists cannot explain why they have to use forces or basic particles and why there are space and time. Any elegant theory starts from some creative assumptions. In this sense, we have to accept that *The far end of modeling is art.*

## ACKNOWLEDGMENTS

This work is supported by an US National Science Foundation grant IIS-00-92664 and an US Office of Naval Research grant N-000140-110-535. The author would like to thank David Mumford, Yingnian Wu, and Alan Yuille for extensive discussions that lead to the development of this paper, and also thank Zhuowen Tu, and Cheng-en Guo for their assistance.

## REFERENCES

- [1] A. Amir and M. Lindenbaum, "Ground from Figure Discrimination," *Computer Vision and Image Understanding*, vol. 76, no. 1, pp. 7-18, 1999.
- [2] J.J. Atick and A.N. Redlich, "What Does the Retina Know about Natural Scenes?" *Neural Computation*, vol. 4, pp. 196-210, 1992.
- [3] F. Attneave, "Some Informational Aspects of Visual Perception," *Psychological Rev.*, vol. 61, pp. 183-193, 1954.
- [4] L. Alvarez, Y. Gousseau, and J.-M. Morel, "The Size of Objects in Natural and Artificial Images," *Advances in Imaging and Electron Physics*, J.-M. Morel, ed., vol. 111, 1999.
- [5] H.B. Barlow, "Possible Principles Underlying the Transformation of Sensory Messages," *Sensory Communication*, W.A. Rosenblith, ed. pp. 217-234, Cambridge, Mass.: MIT Press, 1961.
- [6] J. Besag, "Spatial Interaction and the Statistical Analysis of Lattice Systems (with discussion)," *J. Royal Statistical Soc., B*, vol. 36, pp. 192-236, 1973.
- [7] E. Bienenstock, S. Geman, and D. Potter, "Compositionality, MDL Priors, and Object Recognition," *Proc. Neural Information Processing Systems*, 1997.
- [8] A. Blake and A. Zisserman, *Visual Reconstruction*. Cambridge, Mass.: MIT Press, 1987.
- [9] K.L. Boyer and S. Sarkar, "Perceptual Organization in Computer Vision: Status, Challenges, and Potentials," *Computer Vision and Image Understanding*, vol. 76, no. 1, pp. 1-5, 1999.
- [10] E.J. Candes and D.L. Donoho, "Ridgelets: A Key to Higher Dimensional Intermitency?" *Philosophical Trans. Royal Soc. London, A*, vol 357, no. 1760, pp. 2495-2509, 1999.
- [11] C.R. Carlson, "Thresholds for Perceived Image Sharpness," *Photographic Science and Eng.*, vol. 22, pp. 69-71, 1978.
- [12] D. Chandler, *Introduction to Modern Statistical Mechanics*. Oxford Univ. Press, 1987.
- [13] Z.Y. Chi, "Probabilistic Models for Complex Systems," doctoral dissertation with S. Geman, Division of Applied Math, Brown Univ., 1998.
- [14] C. Chubb and M.S. Landy, "Othogonal Distribution Analysis: A New Approach to the Study of Texture Perception," *Comp. Models of Visual Processing*, M.S. Landy, ed. Cambridge, Mass.: MIT Press, 1991.
- [15] R.W. Cohen, I. Gorog, and C.R. Carlson, "Image Descriptors for Displays," Technical Report contract no. N00014-74-C-0184, Office of Navy Research, 1975.
- [16] R.R. Coifman and M.V. Wickerhauser, "Entropy Based Algorithms for Best Basis Selection," *IEEE Trans. Information Theory*, vol. 38, pp. 713-718, 1992.
- [17] P. Common, "Independent Component Analysis—A New Concept?" *Signal Processing*, vol. 36, pp. 287-314, 1994.
- [18] D. Cooper, "Maximum Likelihood Estimation of Markov Process Blob Boundaries in Noisy Images," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 1, pp. 372-384, 1979.
- [19] G.R. Cross and A.K. Jain, "Markov Random Field Texture Models," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 5, pp. 25-39, 1983.
- [20] P. Dayan, G.E. Hinton, R. Neal, and R.S. Zemel, "The Helmholtz Machine," *Neural Computation*, vol. 7, pp. 1022-1037, 1995.
- [21] J.S. De Bonet and P. Viola, "A Non-Parametric Multi-Scale Statistical Model for Natural Images," *Advances in Neural Information Processing*, vol. 10, 1997.
- [22] S. Della Pietra, V. Della Pietra, and J. Lafferty, "Inducing Features of Random Fields," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 4, Apr. 1997.
- [23] N.G. Deriugin, "The Power Spectrum and the Correlation Function of the Television Signal," *Telecomm.*, vol. 1, no. 7, pp. 1-12, 1957.
- [24] S.J. Dickinson, A.P. Pentland, and A. Rosenfeld, "From Volumes to Views: An Approach to 3D Object Recognition," *CVGIP: Image Understanding*, vol. 55, no. 2, pp. 130-154, Mar. 1992.
- [25] D.L. Donoho, M. Vetterli, R.A. DeVore, and I. Daubechic, "Data Compression and Harmonic Analysis," *IEEE Trans. Information Theory*, vol. 6, pp. 2435-2476, 1998.
- [26] A. Efros and T. Leung, "Texture Synthesis by Non-Parametric Sampling," *Proc. Int'l Conf. Computer Vision*, 1999.
- [27] A. Efros and W.T. Freeman, "Image Quilting for Texture Synthesis and Transfer," *Proc. SIGGRAPH*, 2001.
- [28] D.J. Field, "Relations between the Statistics and Natural Images and the Responses Properties of Cortical Cells," *J. Optical Soc. Am. A*, vol. 4, pp. 2379-2394, 1987.
- [29] D.J. Field, "What Is the Goal of Sensory Coding?" *Neural Computation*, vol 6, pp. 559-601, 1994.
- [30] B. Frey and N. Jojic, "Transformed Component Analysis: Joint Estimation of Spatial Transforms and Image Components," *Proc. Int'l Conf. Computer Vision*, 1999.
- [31] K.S. Fu, *Syntactic Pattern Recognition*. Prentice-Hall, 1982.
- [32] W.S. Geisler, J.S. Perry, B.J. Super, and D.P. Gallogly, "Edge Co-occurrence in Natural Images Predicts Contour Grouping Performance," *Vision Research*, vol. 41, pp. 711-724, 2001.
- [33] S. Geman and D. Geman, "Stochastic Relaxation, Gibbs Distributions and the Bayesian Restoration of Images," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 6, pp 721-741, 1984.
- [34] J.W. Gibbs, *Elementary Principles of Statistical Mechanics*. Yale Univ. Press, 1902.
- [35] J.J. Gibson, *The Perception of the Visual World*. Boston: Houghton Mifflin, 1966.
- [36] U. Grenander, *Lectures in Pattern Theory I, II, and III*. Springer, 1976-1981.
- [37] U. Grenander, Y. Chow, and K.M. Keenan, *Hands: A Pattern Theoretical Study of Biological Shapes*. New York: Springer-Verlag, 1991.
- [38] U. Grenander and A. Srivastava, "Probability Models for Clutter in Natural Images," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 23, no. 4, Apr. 2001.
- [39] M.G. Gu and F.H. Kong, "A Stochastic Approximation Algorithm with MCMC Method for Incomplete Data Estimation Problems," *Proc. Nat'l Academy of Sciences*, vol. 95, pp 7270-7274, 1998.
- [40] C.E. Guo, S.C. Zhu, and Y.N. Wu, "Visual Learning by Integrating Descriptive and Generative Methods," *Proc. Int'l Conf. Computer Vision*, 2001.
- [41] G. Guy and G. Medioni, "Inferring Global Perceptual Contours from Local Features," *Int'l J. Computer Vision*, vol. 20, pp. 113-133, 1996.
- [42] D.J. Heeger and J.R. Bergen, "Pyramid-Based Texture Analysis/Synthesis," *Proc. SIGGRAPH*, 1995.
- [43] D.W. Jacobs, "Recognizing 3D Objects Using 2D Images," doctoral dissertation, MIT AI Laboratory, 1993.
- [44] E.T. Jaynes, "Information Theory and Statistical Mechanics," *Physical Rev.*, vol. 106, pp. 620-630, 1957.
- [45] B. Julesz, "Textons, the Elements of Texture Perception and Their Interactions," *Nature*, vol. 290, pp. 91-97, 1981.
- [46] B. Julesz, *Dislogues on Perception*. Cambridge, Mass.: MIT Press, 1995.
- [47] G. Kanizsa, *Organization in Vision*. New York: Praeger, 1979.
- [48] M. Kass, A. Witkin, and D. Terzopoulos, "Snakes: Active Contour Models," *Proc. Int'l Conf. Computer Vision*, 1987.
- [49] D. Kersten, "Predictability and Redundancy of Natural Images," *J. Optical Soc. Am. A*, vol. 4, no. 12, pp. 2395-2400, 1987.
- [50] S.M. Konishi, J.M. Coughlan, A.L. Yuille, and S.C. Zhu, "Fundamental Bounds on Edge Detection: An Information Theoretic Evaluation of Different Edge Cues," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 25, no. 1, Jan. 2003.

- [51] K. Koffka, *Principles of Gestalt Psychology*. New York: Harcourt, Brace and Co., 1935.
- [52] A. Koloydenko, "Modeling Natural Microimage Statistics," PhD Thesis, Dept. of Math and Statistics, Univ. of Massachusetts, Amherst, 2000.
- [53] A.B. Lee, J.G. Huang, and D.B. Mumford, "Random Collage Model for Natural Images," *Int'l J. Computer Vision*, Oct. 2000.
- [54] L. Liang, X.W. Liu, Y. Xu, B.N. Guo, and H.Y. Shum, "Real-Time Texture Synthesis by Patch-Based Sampling," Technical Report MSR-TR-2001-40, Mar. 2001.
- [55] C. Liu, S.C. Zhu, and H.Y. Shum, "Learning Inhomogeneous Gibbs Model of Face by Minimax Entropy," *Proc. Int'l Conf. Computer Vision*, 2001.
- [56] L.D. Lowe, *Perceptual Organization and Visual Recognition*. Kluwer Academic Publishers, 1985.
- [57] S.G. Mallat, "A Theory for Multiresolution Signal Decomposition: The Wavelet Representation," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 11, no. 7, pp. 674-693, July 1989.
- [58] K.L. Mengersen and R.L. Tweedie, "Rates of Convergence of the Hastings and Metropolis Algorithms," *Annals of Statistics*, vol. 24, pp. 101-121, 1994.
- [59] Y. Meyer, "Principe d'Incertitude, Bases Hilbertiennes et Algebres d'Operateurs," *Bourbaki Seminar*, no. 662, 1985-1986.
- [60] Y. Meyer, *Ondelettes et Operateurs*. Hermann, 1988.
- [61] L. Moisan, A. Desolneux, and J.-M. Morel, "Meaningful Alignments," *Int'l J. Computer Vision*, vol. 40, no. 1, pp. 7-23, 2000.
- [62] D.B. Mumford, "Elastic and Computer Vision," *Algebraic Geometry and Its Applications*, C.L. Bajaj, ed. New York: Springer-Verlag, 1994.
- [63] D.B. Mumford and J. Shah, "Optimal Approximations of Piecewise Smooth Functions and Associated Variational Problems," *Comm. Pure and Applied Math.*, vol. 42, 1989.
- [64] D.B. Mumford, "Pattern Theory: A Unifying Perspective," *Proc. First European Congress of Math.*, 1994.
- [65] D.B. Mumford, "The Statistical Description of Visual Signals," *Proc. Third Int'l Congress on Industrial and Applied Math.*, K. Kirchgassner, O. Mahrenholtz, and R. Mennicken, eds., 1996.
- [66] D.B. Mumford and B. Gidas, "Stochastic Models for Generic Images," *Quarterly of Applied Math.*, vol. LIX, no. 1, pp. 85-111, 2001.
- [67] B.A. Olshausen and D.J. Field, "Sparse Coding with an Overcomplete Basis Set: A Strategy Employed by V1?" *Vision Research*, vol. 37, pp. 3311-3325, 1997.
- [68] M.B. Priestley, *Spectral Analysis and Time Series*. London: Academic Press, 1981.
- [69] T. Poggio, V. Torre, and C. Koch, "Computational Vision and Regularization Theory," *Nature*, vol. 317, pp. 314-319, 1985.
- [70] K. Popat and R.W. Picard, "Novel Cluster-Based Probability Model for Texture Synthesis, Classification, and Compression," *Proc. SPIE Visual Comm. and Image*, pp. 756-768, 1993.
- [71] J. Portilla and E.P. Simoncelli, "A Parametric Texture Model Based on Joint Statistics of Complex Wavelet Coefficients," *Int'l J. Computer Vision*, vol. 40, no. 1, pp. 49-71, 2000.
- [72] D.L. Ruderman, "The Statistics of Natural Images," *Network: Computation in Neural Systems*, vol. 5, pp. 517-548, 1994.
- [73] D.L. Ruderman, "Origins of Scaling in Natural Images," *Vision Research*, vol. 37, pp. 3385-3398, Dec. 1997.
- [74] S. Sarkar and K.L. Boyer, "Integration, Inference, and Management of Spatial Information Using Bayesian Networks: Perceptual Organization," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 15, no. 3, Mar. 1993.
- [75] S. Sarkar and K.L. Boyer, *Computing Perceptual Organization in Computer Vision*. Singapore: World Scientific, 1994.
- [76] C. Shannon, "A Mathematical Theory of Communication," *Bell System Technical J.*, vol. 27, 1948.
- [77] E.P. Simoncelli, W.T. Freeman, E.H. Adelson, and D.J. Heeger, "Shiftable Multiscale Transforms," *IEEE Trans. Information Theory*, vol. 38, no. 2, pp. 587-607, 1992.
- [78] E.P. Simoncelli and B.A. Olshausen, "Natural Image Statistics and Neural Representation," *Ann. Rev. Neuroscience*, vol. 24, pp. 1193-1216, 2001.
- [79] B.J. Smith, "Perceptual Organization in a Random Stimulus," *Human and Machine Vision*, A. Rosenfeld, ed. San Diego, Calif.: Academic Press, 1986.
- [80] D. Stoyan, W.S. Kendall, and J. Mecke, *Stochastic Geometry and Its Applications*. John Wiley and Sons, 1987.
- [81] D. Terzopoulos, "Multilevel Computational Process for Visual Surface Reconstruction," *Computer Vision, Graphics, and Image Processing*, vol. 24, pp. 52-96, 1983.
- [82] Z.W. Tu and S.C. Zhu, "Image Segmentation by Data Driven Markov Chain Monte Carlo," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 24, no. 5 May 2002.
- [83] Z.W. Tu and S.C. Zhu, "Parsing Images into Region and Curve Processes," *Proc. European Conf. Computer Vision*, 2002.
- [84] J.H. Van Hateren and D.L. Ruderman, "Independent Component Analysis of Natural Image Sequences Yields Spatiotemporal Filters Similar to Simple Cells in Primary Visual Cortex," *Proc. Royal Soc. London*, vol. 265, 1998.
- [85] L.R. Williams and D.W. Jacobs, "Stochastic Completion Fields: A Neural Model of Illusory Contour Shape and Salience," *Neural Computation*, vol. 9, pp. 837-858, 1997.
- [86] D.R. Wolf and E.I. George, "Maximally Informative Statistics," unpublished manuscript, 1999.
- [87] Y.N. Wu and S.C. Zhu, "Equivalence of Julesz and Gibbs Ensembles," *Proc. Int'l Conf. Computer Vision*, 1999.
- [88] Y.N. Wu, S.C. Zhu, and C.E. Guo, "Statistical Modeling of Image Sketch," *Proc. European Conf. Computer Vision*, 2002.
- [89] J.S. Yedidia, W.T. Freeman, and Y. Weiss, "Generalized Belief Propagation," TR-2000-26, Mitsubishi Electric Research Lab., 2000.
- [90] A.L. Yuille, "Deformable Templates for Face Recognition," *J. Cognitive Neuroscience*, vol. 3, no. 1, 1991.
- [91] A.L. Yuille, J.M. Coughlan, Y.N. Wu, and S.C. Zhu, "Order Parameter for Detecting Target Curves in Images: How Does High Level Knowledge Helps?" *Int'l J. Computer Vision*, vol. 41, no. 1/2, pp. 9-33, 2001.
- [92] A.L. Yuille, "CCCP Algorithms to Minimize the Bethe and Kikuchi Free Energies: Convergent Alternatives to Belief Propagation," *Neural Computation*, 2001.
- [93] S.C. Zhu, Y.N. Wu, and D.B. Mumford, "Minimax Entropy Principle and Its Application to Texture Modeling," *Neural Computation*, vol. 9, no. 8, pp. 1627-1660, Nov. 1997.
- [94] S.C. Zhu and D.B. Mumford, "Prior Learning and Gibbs Reaction-Diffusion," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 19, no. 11, pp. 1236-1250, Nov. 1997.
- [95] S.C. Zhu, Y.N. Wu, and D.B. Mumford, "Filters, Random Fields, and Maximum Entropy (FRAME): Towards a Unified Theory for Texture Modeling," *Int'l J. Computer Vision*, vol. 27, no. 2, pp. 1-20, 1998.
- [96] S.C. Zhu, "Embedding Gestalt Laws in Markov Random Fields," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 21, no. 11, pp. 1170-1187, Nov. 1999.
- [97] S.C. Zhu, X.W. Liu, and Y.N. Wu, "Exploring Julesz Texture Ensemble by Effective Markov Chain Monte Carlo," *IEEE Trans. Pattern Analysis and Machine Intelligence*, vol. 22, no. 6, June 2000.
- [98] S.C. Zhu, R. Zhang, and Z.W. Tu, "Integrating Top-Down/Bottom-Up for Object Recognition by DDMCMC," *Proc. Computer Vision and Pattern Recognition*, 2000.
- [99] S.C. Zhu, C.E. Guo, Y.N. Wu, and Y.Z. Wang, "What Are Textons," *Proc. European Conf. Computer Vision*, 2002.
- [100] G.J. Burton and J.R. Moorehead, "Color and Spatial Structures in Natural Scenes," *Applied Optics*, vol. 26, no. 1, pp. 157-170, 1987.
- [101] D.L. Donoho, "Wedgelets: Nearly Minimax Estimation of Edges," *Annals of Statistics*, vol. 27, no. 3, pp. 859-897, 1999.



**Song-Chun Zhu** received the BS degree from the University of Science and Technology of China in 1991, and the MS and PhD degrees from Harvard University in 1994 and 1996, respectively. All degrees are in computer science. He is currently an associate professor jointly with the Departments of Statistics and Computer Science at the University of California, Los Angeles (UCLA). He is a codirector of the UCLA Center for Image and Vision Science.

Before joining UCLA, he worked at Brown University (applied math), Stanford University (computer science), and Ohio State University (computer science). His research is focused on computer vision and learning, statistical modeling, and stochastic computing. He has published more than 50 articles and received a number of honors, including a David Marr prize honorary nomination, the Sloan fellow in computer science, the US National Science Foundation Career Award, and an Office of Naval Research Young Investigator Award.