

LOGOS: A MODULAR BAYESIAN MODEL FOR *DE NOVO* MOTIF DETECTION

ERIC P. XING*

*Computer Science Division, University of California
Berkeley, CA 94720, USA
epxing@cs.berkeley.edu*

WEI WU

*Life Sciences Division, Lawrence Berkeley National Laboratory
Berkeley, CA 94720, USA
wwu@lbl.gov*

MICHAEL I. JORDAN

*Computer Science and Statistics, University of California
Berkeley, CA 94720, USA
jordan@cs.berkeley.edu*

RICHARD M. KARP

*Computer Science Division, University of California
Berkeley, CA 94720, USA
karp@cs.berkeley.edu*

Received 19 August 2003

Revised 11 November 2003

Accepted 7 December 2003

The complexity of the global organization and internal structure of motifs in higher eukaryotic organisms raises significant challenges for motif detection techniques. To achieve successful *de novo* motif detection, it is necessary to model the complex dependencies within and among motifs and to incorporate biological prior knowledge. In this paper, we present **LOGOS**, an integrated **L**Ocal and **G**IObal motif **S**equences model for biopolymer sequences, which provides a principled framework for developing, modularizing, extending and computing expressive motif models for complex biopolymer sequence analysis. **LOGOS** consists of two interacting submodels: HMDM, a local alignment model capturing biological prior knowledge and positional dependency within the motif local structure; and HMM, a global motif distribution model modeling frequencies and dependencies of motif occurrences. Model parameters can be fit using training motifs within an empirical Bayesian framework. A variational EM algorithm is developed for *de novo* motif detection. **LOGOS** improves over existing models that ignore biological priors and dependencies in motif structures and motif occurrences, and demonstrates

*To whom correspondence should be addressed.

superior performance on both semi-realistic test data and *cis*-regulatory sequences from yeast and *Drosophila* genomes with regard to sensitivity, specificity, flexibility and extensibility.

Keywords: Cis-regulatory system; Bayesian model; Dirichlet prior; hidden Markov model; variational inference.

1. Introduction

The identification of motif structures within biopolymer sequences such as DNA and protein is an important task in computational biology and is essential in advancing our knowledge about biological systems. It is known that only a small fraction of the genomic sequences in multi-cellular higher organisms constitute the protein coding information of the genes (e.g. only 1.5% for human genomes¹), whereas the rest of the genome, besides playing purely structural roles such as forming the centromeres and telomeres of the chromosomes, contains a large number of short sequence motifs that make up an immensely rich codebook of the gene regulation program, known as the **cis-regulatory system**. It is believed that this regulatory program determines the level, location and chronology of gene expression, which significantly, if not predominantly, contributes to the developmental, morphological and behavioral diversity of complex organisms.⁷

The problem of *de novo* motif detection^a has been widely studied. Numerous algorithmic approaches have been proposed, most of which use probabilistic generative models to model motifs as stochastic string patterns randomly embedded in a simple background. In such a setting, motif detection can be formulated as a standard missing-value inference and parameter estimation problem (for motif locations and position weight matrices, respectively), and standard methods such as EM and Gibbs sampling can be applied. This literature is too large to survey here, but some relevant examples include MEME,² BioProspector,¹⁹ and AlignACE.²¹ A different framework based on word segmentation and dictionary construction was proposed in the MobyDick algorithm,⁵ which pointed out the importance of combinatorial analysis of a large set of potential motifs jointly, so that some dependencies among motifs can be captured. A similar “word-enumeration” idea also appeared in Ref. 20. Recently, Gupta and Liu extended the dictionary model to a stochastic dictionary (SD) model by replacing the words in the dictionary with *probabilistic word matrices*,¹² allowing stochasticity of motif instances to be modeled. Many of these methods have been widely used and have been successful empirically for motif detection in well curated bacterial and yeast gene regulatory sequences. However, generalization of these successful results to longer, more complex and weakly characterized input sequences such as those from higher eukaryotic genomes seems less immediate. A recent survey by Eisen raises concerns over the inability of some

^aNot to be confused with model-based *motif scan*, the task of searching known motifs based on given position weight matrices, as addressed by Frith *et al.*¹⁰ and Huang *et al.*¹⁵

contemporary motif models to incorporate biological knowledge of global motif distribution, motif structures and motif sequence composition.⁹

Several recent studies have tried to address these concerns from different perspectives. For example, some authors have proposed better objective functions for motif detection, by scoring motifs based on the statistical significance of the information content,¹⁴ and by considering cooperative motif binding between multiple transcription factors.¹¹ Van Helden *et al.* recently suggested using a signature conservation pattern to constrain the motif patterns.¹³ Bussemaker *et al.* proposed incorporating gene expression data from microarrays for motif detection.⁶ Frith *et al.* used an HMM in their motif scanner to model the possible presence of clustered motif occurrences in complex *cis*-regulatory sequences.¹⁰ Although these attempts head in the direction of more expressive motif models, it is not clear whether these ideas can be integrated to assemble a powerful yet transparent and computationally efficient motif detection algorithm.

We are interested in developing a principled general framework for motif modeling, which is expressive (in terms of being able to describe internal structures, inter-motif relations, motif abundances, etc., and readily incorporates prior knowledge from experimental biology), yet mathematically and algorithmically transparent and well-structured, hence simplifying model construction, computation and extension. In a recent methodological paper, we briefly laid out a theoretical foundation for modular motif models where we made explicit the decomposition of a full motif model into the following two components: the *global distribution model*, which models the frequencies of different motifs and the dependencies between motif occurrences in a sequence; and the *local alignment model*, which captures the intrinsic properties within motifs, including characteristic position weight matrices (PWMs) and site dependencies.²⁶ Based on this framework, we extended the conventional motif-alignment model into a very expressive hierarchical Bayesian Markovian model, called a hidden Markov Dirichlet-Multinomial (HMDM) model, for local alignment, which successfully captures internal motif structure and incorporates prior knowledge from biologically known motifs using a structured Bayesian prior model for the PWMs of motifs. In the current paper, we integrate the HMDM model into a general framework for the modeling of motif-containing biopolymer sequences and present a *de novo* motif detector developed based on this framework. This framework uses the HMDM model as the local alignment submodel and uses a newly designed HMM that we describe here for the global submodel. A *variational EM* (VEM) algorithm is developed for efficient Bayesian learning and prediction. We call our framework **LOGOS**, for integrated **L**ocal and **G**lobal motif **S**equences model.^b

^bNot to be confused with “*logo*,” a graphic representation of an aligned set of biopolymer sequences first introduced by Tom Schneider²² to help visualizing the consensus and the entropy (or “information”) patterns of monomer frequencies. A *logo* is not a motif finding algorithm, but is often used as a way to present motifs visually.

2. LOGOS: A Modular Generative Framework for Motif Sequences

2.1. Preliminaries

Motifs are short stochastic string patterns scattered in biopolymer sequences such as DNA and proteins. The characteristic sequence patterns of motifs and their locations often relate to potentially important biological functions such as serving as the *cis*-elements for gene regulation or as the catalytic sites for protein activities. Numerous biological studies have revealed rich architecture in the global organization and the internal structures of motifs in higher eukaryotic organisms. Taking DNA motifs as an example, it is well known that the *cis*-regulatory elements often occur in clusters (referred to as *cis*-modules), possibly for eliciting synergistic or more robust regulatory signals.⁷ The biophysical mechanisms of DNA-protein interactions at the motif-binding sites further suggest that the sites within the DNA motifs are not necessarily uniformly conserved.^{24,25} Rather, the conservation pattern may be subject to a constraint imposed by the structure of the binding protein, resulting in the so-called ‘shape’ bias (Fig. 1). These *meta-sequence features* of the motif structure raise significant challenges to conventional motif-finding algorithms, which primarily rely on simplifying independence assumptions that decouple (potential) associations among sites within each single motif and among multiple instances of motifs. (For example, the conventional product multinomial model to be described shortly assigns equal probability to both the original motif and its permuted version in Fig. 1.)

In the following paragraph, we introduce the necessary notation for our presentation. Note that to simplify the presentation, we use DNA motifs as a running example, but it should be clear that our technique is readily applicable to protein motifs.

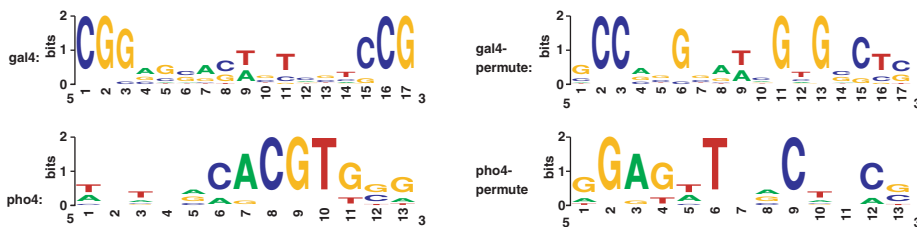


Fig. 1. An illustration of the *shape bias*. On the left-hand side are two genuine motif patterns. On the right are artificial patterns resulted from a column permutation of the original motifs. It is believed that for plausible biological motifs the conserved sites are more likely to occur consecutively and possibly followed (or preceded) by heterogeneous sites that are also consecutive (rather than interspersed). Such characteristic conservation patterns of the sites in a motif are often reflected in the ‘contour shape’ of the motif logo (e.g. *U*- or *bell-shaped*, as exhibited by motifs *gal4* and *pho4*, respectively), which reflects the *spatial* pattern of the information content over all sites. It is important to note that ‘shape’ is only associated with the conservation pattern of a motif PWM, but **not** with any specific consensus sequences of the motif.

We denote a regulatory DNA sequence by a character string $y = (y_1, \dots, y_T) \in \{A, T, C, G\}$. An indicator string x signals the locations of the motif occurrences. Following biological convention, we denote the *multi-alignment* of M instances of a motif of length L by an $M \times L$ matrix \mathbf{A} , of which each *column* corresponds to a *position* or *site* in the motif. The multi-alignment of all instances of motif k specified by the indicator string x in sequence y is denoted by $\mathbf{A}^{(k)}(x, y)$. We define a *counting matrix* $h(\mathbf{A}^{(k)})$ (or $h^{(k)}(x, y)$) for each motif alignment, where each column $h_l = [h_{l1}, \dots, h_{l4}]'$ is an integer vector with four elements, specifying the number of occurrences of each nucleotide (nt) at position l of the motif. (Similarly we define the *counting vector* h_{bk} for the background sequence $y - \mathbf{A}$, where the somewhat abusive use of the minus sign means excluding all motif sub-sequences in \mathbf{A} from y .) We assume that the nucleotides at position l of motif k admit a *position-specific multinomial distribution*, $\theta_l^{(k)} = [\theta_{l1}^{(k)}, \dots, \theta_{l4}^{(k)}]'$. The ordered set of position-specific multinomial parameters of all positions of motif k , $\theta^{(k)} = \{\theta_1^{(k)}, \dots, \theta_{L^{(k)}}^{(k)}\}$, is referred to as a *position weight matrix* (PWM). It is clear that the counting matrix $h^{(k)}$ corresponds to the *sufficient statistics* of PWM $\theta^{(k)}$. Formally, the problem of motif detection is that of inferring $\mathbf{x} = \{x^{(1)}, \dots, x^{(N)}\}$ and estimating $\theta = \{\theta^{(1)}, \dots, \theta^{(K)}\}$, given a set of sequences $\mathbf{y} = \{y^{(1)}, \dots, y^{(N)}\}$. For simplicity, we omit the superscript k (motif type index) of variable θ and the superscript n (sequence index) of variable x and y in wherever it is clear from the context that we are focusing on a generic motif type or a generic sequence.

2.2. The modular motif model

Without loss of generality, assume that the occurrences of motifs in a DNA sequence, as indicated by x , are governed by a **global distribution model** $p(x|\Theta_g, \mathcal{M}_g)$, and for each type of motif, the nucleotide sequence pattern shared by all its instances admits a **local alignment model** $p(\mathbf{A}(x, y)|x, \Theta_l, \mathcal{M}_l)$. We further assume that the background non-motif sequences are modeled by a simple conditional model, $p(y - \mathbf{A}(y, x)|x, \Theta_{bk})$, where the background nt-distribution parameters Θ_{bk} are assumed to be estimated *a priori* from the entire sequence. The symbols Θ_g , Θ_l , \mathcal{M}_g and \mathcal{M}_l stand for the parameters (e.g. the PWMs) and model classes (e.g. a product multinomial model) in the respective submodels. Thus, the likelihood of a regulatory sequence y is:

$$\begin{aligned} p(y|\Theta, \mathcal{M}) &= \sum_x p(x|\Theta_g, \mathcal{M}_g)p(y|x, \Theta_l, \mathcal{M}_l) \\ &= \sum_x p(x|\Theta_g, \mathcal{M}_g)p(\mathbf{A}|x, \Theta_l, \mathcal{M}_l)p(y - \mathbf{A}|x, \Theta_{bk}), \end{aligned} \quad (1)$$

where $\mathbf{A} \triangleq \mathbf{A}(x, y)$. Note that Θ_l here is not necessarily equivalent to the PWMs (θ) of the motifs, but is a generic symbol for the parameters of a more general model of the aligned motif instances. (For example in the HMDM model to be defined shortly, Θ_l refers to the hyperparameters that describe a distribution of PWMs.)

Equation (1) makes explicit the modular structure of the **LOGOS** framework for generic motif models. The submodel $p(x|\Theta_g, \mathcal{M}_g)$ captures properties such as the frequencies of different motifs and the dependencies between motif occurrences. On the other hand, the submodel $p(\mathbf{A}|x, \Theta_l, \mathcal{M}_l)$ captures the intrinsic properties within motifs that can help to improve sensitivity and specificity to genuine motif patterns. Depending on the value of the latent indicator x_t (e.g., motif or not) at each position t , y_t admits different probabilistic distributions, such as a particular nucleotide distribution inside a motif or a background distribution.

For example, the conventional *uniform and independent* (UI) model for motif start-positions used in many motif finding algorithms is an instance of a simple global model, where the motif instances are assumed to occur independently with uniform probability at all possible locations in a sequence. Therefore, $p(x) = \prod_{m=1}^M p(x_m)$, where $p(x_m = t)$ is the marginal probability of the m -th motif at location t , which in this case is a uniform distribution over all t , and the same for all M instances. Note that there is no *model constraint* to prohibit overlapping motif instances.^c The UI model does not appear to be problematic in *de novo* motif finding tasks involving bacterial or even simple yeast sequence sets, in which the input sequences are usually small in size and homogeneous in content (e.g. pre-screened according to mRNA co-expression) and the motif occurrences tend to be sparse. But some recent studies including our own experiments suggest that the correctness of motif finding based on the UI assumption starts to break down for less well pre-screened input sequences or for those with clustered motif occurrences, such as the *Drosophila* gene regulatory sequences.⁴

An example of the local model is the standard product multinomial (PM) model, where the position-specific nt-distributions within a motif are assumed to be independent.¹⁹ Thus the likelihood of a multi-alignment \mathbf{A} is: $p(\mathbf{A}|\Theta) = \prod_{l=1}^L \prod_{j=1}^4 [\theta_{lj}]^{h_{lj}}$. Although a popular model for many *de novo* motif finders, PM nevertheless is sensitive to noise and random or trivial recurrent patterns (e.g. poly-N or repetitions of short k -mers such as CpG islands), and is unable to capture potential site dependencies inside the motifs. Various pattern-driven approaches (e.g. using a fragmentation model,¹⁸ splitting a “two-block” motif into two coupled sub-motifs,^{2,19} or imposing explicit “shape”¹³ or entropy constraints¹⁷), have been developed to handle special patterns such as the *U-shaped* motifs, but generalization to other “shapes” seen from known motifs is not very straightforward. The mixture of PWMs and tree-based Bayesian network models recently developed by Barash *et al.* can capture positional dependencies within motifs.³ But these models are motif specific and do not incorporate prior knowledges about typical dependency patterns implied in the actual biologically identified motifs (we will

^cHeuristics are generally employed — such as throwing away overlapping sampled motifs (in the Gibbs sampler) or rescaling the joint posterior of x (in MEME) — to enforce the *non-overlapping constraint*. Nevertheless, this results in inconsistencies between the computed motif distribution and the one defined by the model, and incurs a sizable overhead due to wasteful computations.

elaborate this point in later discussions). Dirichlet priors for θ have been used in the PM setting,^{2,18} but they are primarily used for smoothing rather than for explicitly incorporating prior knowledge about motifs.

Recently, Xing *et al.*²⁶ developed the HMDM model for motif alignment, which captures site dependencies inside the motifs and incorporates prior knowledge of nt-distributions of all motif sites from biologically known motifs. It shows improved sensitivity (compared to PM) to true biological motifs in the presence of synthetic false motifs in the motif detection setting. Frith *et al.*¹⁰ proposed an HMM model for *cis*-element clusters in higher eukaryotic DNA, which shows promising performance in motif scanning (for which the PWMs are given). Our goal in this paper is to develop an expressive modular motif model that builds on these previous lines of research.

We present a *de novo* motif detection algorithm using an HMM as the global distribution model and an HMDM as the local alignment model. The resulting composite **LOGOS** model is capable of: (1) performing formal and efficient inference of global motif occurrences under a flexible setting that allows clustered motif instances, multiple motif types, and motifs on reverse complementary sequences; (2) correctly enforcing the non-overlapping constraint; (3) capturing site dependencies inside the motifs so as to bias prediction toward more biologically plausible motifs while remaining flexible with regards to motif shapes and lengths; and (4) incorporating prior knowledge of nt composition at each motif site to provide smoothed and robust Bayesian estimation of the PWMs.

2.3. The local model: An HMDM for motif alignment

The local alignment model is crucial for identifying the correct motif patterns in a noisy background. As mentioned before, many motifs are not uniformly well-conserved at all their sites^d (e.g. *gal4* in Fig. 1). Biological evidence shows that conserved sites are likely to occur consecutively.⁹ This is called *site clustering*, one of the main motivations for the HMDM model. Obviously the PM model can not model such patterns: given a length L motif for which only $\frac{L}{2}$ positions are conserved, PM would assign the same probability regardless of the locations of the conserved sites.

In the HMDM model (Fig. 2), we assume that there are I underlying latent nt-distribution prototypes,^e according to which position-specific multinomial distributions of nt are determined, and that each prototype is represented by a Dirichlet distribution, which defines a probability distribution over the simplex of

^dA possible reason could be that a binding protein only interacts with a DNA target through a few highly specific aa-nt interactions, but is tolerant of variations in other sites.

^eWe can roughly imagine that the set of prototypes should include prototypes corresponding to four possible conserved nt-distributions (i.e. those having most of the probability mass at A, C, G, T, respectively), as well as other prototypes corresponding to distributions that are less conserved or even heterogeneous in different ways.

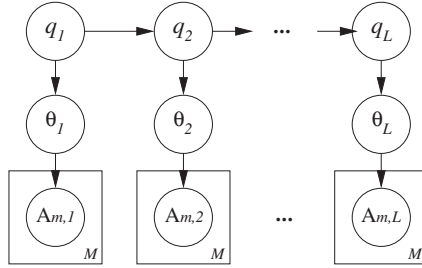


Fig. 2. The HMDM model for motif instances specified by a given x . The circles are random variables and the boxes are plates representing replicates (i.e. M instances of a motif).

multinomial parameters (see Appendix A). Furthermore, the sequence of prototypes at consecutive positions in the motif is governed by a first-order Markov process, which defines a structural prior for the position weight matrices.

More precisely, a multi-alignment \mathbf{A} containing M motif instances is generated by the following process. (1) We first sample a sequence of prototype indicators $q = (q_1, \dots, q_L)$ from a first-order Markov chain with initial distribution π and transition matrix B . Then we repeat the following for each column $l \in \{1, \dots, L\}$: (1) A component from a Dirichlet mixture $\alpha = \{\alpha_1, \dots, \alpha_I\}$, where each $\alpha_i = [\alpha_{i1}, \dots, \alpha_{i4}]'$ is the parameter vector for a Dirichlet distribution, is picked (deterministically) according to indicator q_l . Say, $q_l = i$, thus we picked α_i . (2) A multinomial distribution θ_l is sampled according to $p(\theta|\alpha_i)$, the probability defined by Dirichlet component i . (3) All the nucleotides in column l are generated *i.i.d.* according to the multinomial distribution parametrized by θ_l .

The complete likelihood of motif alignment $\mathbf{A}_{M \times L}$ characterized by a counting matrix h is:

$$p(\mathbf{A}, q, \theta | x, \Theta_l, \mathcal{M}_l) = p(\mathbf{A} | x, \theta) p(\theta | q, \alpha) p(q | \pi, B), \quad (2)$$

where (using the update properties of the Dirichlet distribution and denoting $q_l^i = 1$ if q_l is in state i and 0 otherwise):

$$p(\mathbf{A} | x, \theta) p(\theta | q, \alpha) = \prod_{l=1}^L \prod_{i=1}^I \text{Dir}(\alpha_i + h_l) q_l^i, \quad (3)$$

$$p(q | \pi, B) = \prod_{i=1}^I [\pi_i]^{q_1^i} \prod_{l=1}^{L-1} \prod_{i,j=1}^I [B_{i,j}]^{q_l^i q_{l+1}^j}. \quad (4)$$

The major role of the HMDM model is to impose dynamic priors for modeling data whose distributions exhibit spatial dependencies.

As Fig. 2 makes clear, this model is *not* a simple HMM for discrete sequences. In an HMM model the transitions would be between the emission models (i.e. multinomials) themselves, and the output at each time would be a single data instance in the sequence. In HMDM, the transitions are between different priors for

the emission models, and the direct output of the HMM is the parameter vector of a generative model, which will be sampled multiple times at each position to generate random instances. This approach is especially useful when we have either empirical or learned prior knowledge (e.g., from training motifs) about motif properties such as *site clustering* or other positional dependencies that can be captured by a first order Markov chain (Eq. (4)). We will see an example of this point in Sec. 4.1.

2.4. The global model: An HMM for motif indicators

The HMDM generative process only creates aligned multiple instances of a motif, but does not complete the generation of the observed sequence set. We need a model for the background sequences and another process that generates the positions of the motif instances. For this we need a global model for the indicator variable sequence x that can specify the locations of all motif instances.

Let x be the indicator variable sequence specifying whether each y_t in a DNA sequence is in the background or in a motif, and if in a motif, which motif and where in the motif: $x = (x_1, \dots, x_T)$, where $x_t \in \mathbb{S}$. The indicator state space \mathbb{S} includes all possible identity labels of a monomer (nt) in a sequence: $\mathbb{S} = \mathbb{M} \cup \mathbb{M}' \cup \{b^0, b^1, \dots, b^k, d\}$, where $\mathbb{M} = \{1^{(1)} \dots L_1^{(1)}, 1^{(2)} \dots L_2^{(2)}, \dots, 1^{(K)} \dots L_k^{(K)}\}$ is the set of all possible sites within a motif on the forward strand (i.e., states $1^{(1)}$ to $L_1^{(1)}$ correspond to the sites in motif type 1 on the forward strand, and so on); \mathbb{M}' is the set of all possible sites within a motif on the reverse complementary strand; b^0 corresponds to the inter-cluster background state; $b^k, k \neq 0$ corresponds to the intra-cluster background states; and d represents dummy states. We model the distribution of x with the first-order Markov process depicted in Fig. 3.

The motivation for this Markov model is that we expect to see occasional motif clusters in a large ocean of global background sequences (represented by state b^0), and each motif instance in a cluster is embedded in a corresponding sea of intra-cluster background sequences (b^i). The model assumes that the distance between clusters is geometrically distributed with mean $1/(1 - \beta_{0,0})$, and the distance between motif instances within cluster k is also geometrically distributed with mean $1/(1 - \beta_{k,0})$. As shown in Fig. 3, with equal probability $\beta_{k,k}/2$, an intra-background state b^k reaches the start states $1^{(k)}$ and $L_k^{(k')}$ of motif k on the forward or reverse strand, deterministically passes through all internal sites of motif k (thus avoiding motif overlapping), and transitions back to the same background state b^k , thereby stochastically generating a cluster of occurrences of motif k^f ; b^k also has a small probability $\beta_{k,i}/2$ of transitioning to the start state of another motif i , which terminates cluster k and leads into cluster i ; all intra-background states also have probability $\alpha_{k,0}$ of returning to the global background state. These parameters can in principle be fitted using a training set, or just specified empirically

^fNote that such a scheme does not imply that the actual nucleotide-sites involved in protein binding are all on the same strand, but merely means that we *represent* the motif using the sequence pattern on one of the two complementary strands.

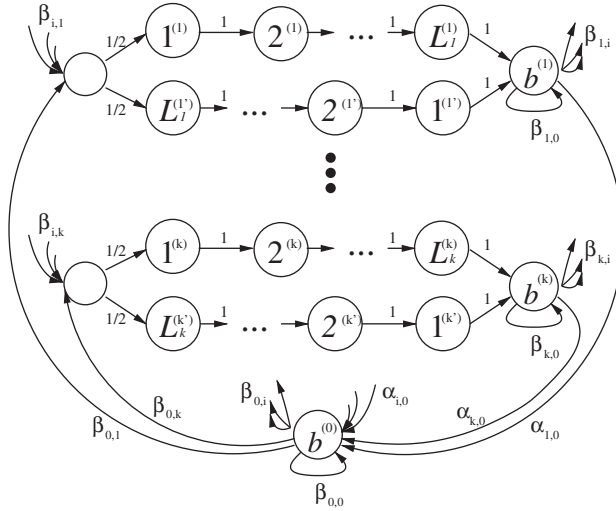


Fig. 3. The global HMM. Labeled circles represent the functional states in DNA sequences. Unlabeled circles are dummy states; arrows between nodes represent state-transitions that have non-zero probabilities; numbers on the arrows represent transition probabilities, and so do the parameter symbols accompanied the arrows (with the parameter subscripts denoting the source and target of the transitions). A path that follows the arrows stochastically leads to a state sequence that is reminiscent of a regulatory DNA sequence with motifs embedded in the background.

based on a rough estimation of the motif or *cis*-module frequencies.[§] Note that these parameters do not impose rigid constraints on the number of motif instances or modules; the actual number of instances is determined by the posterior distribution of the indicator sequence $p(x|y)$.

In accordance with the above framework, we introduce multinomial parameters $\theta_{b_0} = [\theta_{b_0,1}, \dots, \theta_{b_0,4}]'$ and $\theta_{b_1} = [\theta_{b_1,1}, \dots, \theta_{b_1,4}]'$ for the inter- and intra-cluster background nt-distributions, respectively (assuming all intra-cluster backgrounds use the same multinomial model). Thus, given the PWMs Θ of motifs and the background parameters $\Theta_{bk} = \{\theta_{b_0}, \theta_{b_1}\}$, we have the usual joint probability for a conditional HMM model of a motif-containing sequence:

$$p(x, y|\Theta, \Theta_{bk}, \Theta_g, \mathcal{M}_g) = p(x_1) \prod_{t=2}^T p(x_t|x_{t-1}) \prod_{t=1}^T p(y_t|x_t, \Theta, \Theta_{bk}) \quad (5)$$

$$\text{where } p(y_t|x_t, \Theta, \Theta_{bk}) = \prod_{i \in \mathbb{S}} \prod_{j=1}^4 [\theta_{ij}]^{\delta(y_t, j) \delta(x_t, i)}.$$

[§]When no strong knowledge about modular dependencies is available, it is better to just set all bk-to-mt transitions $\beta_{i,j}, i, j \neq 0$, to the same small constant reflecting motif frequency, and similarly for $\beta_{0,k}$ and $\alpha_{k,0}$ reflecting cluster frequency, to avoid overfitting. In our experiment, we parametrize our HMM model in such a fashion. This reduced model is very similar to the one used in Cister,¹⁰ but with unknown PWMs in our case.

The locations of all motif instances encoded in x can be inferred from the global model using Bayes rule.

The HMM model we proposed is not meant to capture fine details of the global motif dependencies, because without a sufficiently large and well-characterized training data set, we could risk overfitting to hypothetical structures and fail to generalize to sequences bearing unknown (and possibly simpler) structures. But within the **LOGOS** framework, if so desired, we can easily generalize to more elaborate models, such as one that models higher order dependencies, or one which uses a more complex background (e.g. a higher-order Markov model) in a principled way. All that is needed is to simply expand the state space \mathbb{S} , and either train or empirically parametrize a more expressive initial and transition model in the global HMM.

3. Inference and Learning Algorithm

3.1. Variational Bayesian learning

In order to do Bayesian estimation of the motif parameter θ , and to predict the locations of motif instances via the indicator sequence x , we need to be able to compute the posterior distribution $p(\theta|y)$, which is infeasible in closed form for a complex motif model (because we have to marginalize out q and x in the joint posterior $p(\theta, q, x|y, \mathcal{M})$). A possible approach is to use a Markov Chain Monte Carlo (MCMC) method, such as a Gibbs sampler, which performs “asymptotically exact inference.” However, concerns over likely slow mixing and difficulties in detecting convergence motivate us to use *variational Bayesian inference*, which has a more deterministic flavor similar to that of EM and is computationally more efficient.

The variational Bayesian inference method developed in Ref. 26 for the HMDM model is a special instance of the *generalized mean field* (GMF) algorithm.²⁷ Briefly, in the GMF framework, a complex joint distribution p , such as the joint posterior $p(\theta, q, x|y, \mathcal{M})$, is approximated with a simpler distribution Q defined by the product of inter-dependent local marginals over disjoint subsets of all domain variables, e.g. $Q(\theta, q, x) = Q_l(\theta, q)Q_g(x)$. The optimal form of each local marginal can be obtained via minimizing the Kullback–Leibler (KL) divergence between Q and p with respect to free distributions Q_l and Q_g .²⁷ Omitting mathematical details, this optimization results in the following coupled updates:

$$Q_g(x) = p(x|\mathcal{M}_g)p(y|x, \bar{\phi}(\theta), \mathcal{M}_g) \quad (6)$$

$$Q_l(\theta, q) = p(q|\mathcal{M}_l)p(\theta|q, \alpha, \bar{h}(y), \mathcal{M}_l) \quad (7)$$

where, E_D denotes expectation with respect to the distribution D , and $\bar{h}(y) = E_{Q_g}[h(x, y)]$, $\bar{\phi}(\theta) = E_{Q_l}[\ln \theta]$,^h which are referred to as the *generalized mean field*

^h $\ln \theta$, where $\ln(\cdot)$ is a componentwise operation, is called the *natural parameterization* of a multinomial.

messages exchanged between submodels conducting probabilistic influences of the respective submodel.

A key property revealed in Eqs. (6) and (7) is the isomorphism of their right-hand sides to those of the Eqs. (2) and (5). Essentially, the variational marginals $Q_g(x)$ and $Q_l(\theta, q)$ recover exactly the same form of the original global and local submodels, except that the motif parameters θ on which the global submodel is conditioned are replaced by their Bayesian estimates (in the natural parameter form), and the sufficient statistics h propagated from the global submodel to the local submodel are replaced by their posterior expectations. This means that the locality of inference and marginalization in the composite **LOGOS** model is preserved in both local and global submodels. We can easily obtain the optimal approximate posterior distribution of θ by marginalizing $Q_l(\theta, q)$ over q , and that of x , using $Q_g(x)$. It can be further proved that the coupled updates (6) and (7) actually optimize a lower bound of the likelihood $p(y|\Theta, \mathcal{M})$ and are guaranteed to converge to a local maximum (as in standard EM).¹⁶ In the following section we summarize the computation procedure involved in **LOGOS**, which we call a “*variational EM*” algorithm (VEM), after its operational resemblance to conventional EM.

3.2. The variational EM algorithm

Due to the locality of variational Bayesian inference, we can perform inference in the local alignment model HMDM as if we have “observations” \bar{h} (to obtain a distribution $Q_l(\theta, q)$ that approximates the marginalized conditional $p(\theta, q|y)$), and in the global HMM model as if the position-specific multinomial distribution of a motif $\bar{\phi}(\theta)$ is given (to obtain $Q_g(x)$ that approximates $p(x|y)$). Therefore, Bayesian estimates of the multinomial parameters can be obtained via fixed-point iteration through the following EM-like procedure:

3.2.1. Variational E step

Compute the expected sufficient statistics, the count matrix $\bar{h} = E_{Q_g(x)}[h]$, via inference in the global motif model given $\bar{\phi}(\theta)$ and sequence set \mathbf{y} :

$$\bar{h} = \sum_{n=1}^N \sum_{t=1}^{T_n-L+1} h(y_{t:t+L-1}^{(n)}) p(x_t^{(n)} = 1 | y^{(n)}, \bar{\phi}(\theta), \Theta_{bk}), \quad (8)$$

where superscript n indicates the n th DNA sequence; $p(x_t^{(n)} = 1 | y^{(n)}, \bar{\phi}(\theta), \Theta_{bk})$ is the posterior probability of position t in sequence n being the start site of a motif given sequence $y^{(n)}$, Bayesian estimate of motif PWMs, and the background, which can be computed using the standard forward-backward algorithm for HMMs on $Q_g(x)$.

3.2.2. Variational M step

Compute the posterior mean of the natural parameter, $\bar{\phi}(\theta) = E_{Q_i(\theta, q)}[\phi(\theta)]$, via inference in the local motif alignment model given \bar{h} :

$$\bar{\phi}(\theta_{i,j}) = \sum_{i=1}^I p(q_l = i | \bar{h}) (\Psi(\alpha_{ij} + \bar{h}_{lj}) - \Psi(|\alpha_i| + |\bar{h}_l|)), \quad (9)$$

where $\Psi(x) = \frac{\partial \log \Gamma(x)}{\partial x} = \frac{\Gamma'(x)}{\Gamma(x)}$ is the digamma function; $p(q_l = i | \bar{h})$ is the posterior probability of hidden state q given “observation” \bar{h} , which can be computed using the forward-backward algorithm on $Q_i(\theta, q)$.

This *modular inference* procedure provides a framework that scales readily to more complex models. For example, the motif distribution model $p(x)$ can be made more sophisticated so as to model complex properties of multiple motifs such as motif-level dependencies (e.g. co-occurrence, overlaps and concentration within regulatory modules) without complicating the inference in the local alignment model. Similarly, the motif alignment model can also be more expressive (e.g. a mixture of HMDMs) without interfering with inference in the motif distribution model.

The Dirichlet parameters and HMM transition matrix of the HMDM are fitted from a training dataset via empirical Bayes estimation⁸ (see Appendix B for details).

4. Experiments

In a prior work,²⁶ we systematically examined the performance of the HMDM model by implementing a prototype motif detector using HMDM as the local model and testing it on semi-realistic datasets in which biologically identified motifs are planted in a random background, possibly in the presence of artificially produced “false motifs” as decoys. The major advantage of using such a test system is that we know the ground truth, i.e. the true locations and PWMs of the motifs to be detected, and hence can reliably compare performance of different models. We showed that HMDM has a notably higher specificity (than PM) to the genuine motifs in the presence of an artificial decoy, and significantly out-performs the PM-based MEME algorithm in the one-motif-per-sequence scenario.

The **LOGOS** model developed in the current paper integrates HMDM as a sub-component, which models the motif alignments, accompanied with an expressive HMM model, which models the global distribution of motifs in a biologically more realistic way than the UI model. In the following sections, we examine the performance of **LOGOS** using both semi-realistic datasets and real genomic sequences from yeast. All yeast motif sequences are obtained from the *Promoter Database of Saccharomyces cerevisiae* (SCPD), 15 of which are used to fit the hyperparameters of the HMDM, and others (independent of the training set) are used for testing. We compare three variants of **LOGOS**, ordered with decreasing model expressiveness, HMDM+HMM (**LOGOS_{hh}**), PM+HMM (**LOGOS_{ph}**) and PM+UI (**LOGOS_{pu}**), as well as the MEME and AlignACE program (both of which are essentially the

same as **LOGOS**_{pu} in terms of model assumptions, but are enhanced by additional pattern-driven submodels, i.e. gapped motifs, and a more sophisticated implementation).

4.1. Learning the HMDM parameters

We learn our HMDM model using a motif collection from the SCPD. Our dataset contains 15 training motifs. Each has 6 to 32 instances all of which have been identified via biological experiments.

We begin with an experiment showing how HMDM can capture intrinsic properties of the motifs. The prior distribution of the position-specific multinomial parameters θ , reflected in the parameters of the Dirichlet mixtures learned from data, can reveal the nt-distribution patterns of the motifs. Examining the transition probabilities between different Dirichlet components further tells us about dependencies between adjacent positions (which indirectly reveals the “shape” information). We set the total number of Dirichlet components to be eight based on an empirical model selection decision to strike a balance between the expressiveness and complexity of the HMDM model. Intuitively, it can be understood as a prior choice of the size of the collection of nt-distribution prototypes needed to capture all conserved patterns (i.e., prototypes centered at a specific nucleotide), the non-conserved patterns (i.e., uniform nt-distribution prototypes), and the intermediate patterns. Figure 4a shows the Dirichlet parameters fitted from the dataset via empirical Bayes estimation. Among the eight Dirichlet components, components 1–4 favor highly concentrated multinomial distributions centered at each possible single nucleotides A, T, G, and C, respectively, suggesting they correspond to “homogeneous” prototypes, whereas components 7 and 8 favor a near uniform distribution of all 4 nt-types, hence “heterogeneous” prototypes. Components 5 and 6 are somewhat in between. Such patterns are consistent with the biologically possible nt-distributions anticipated for motif sites, and suggest that

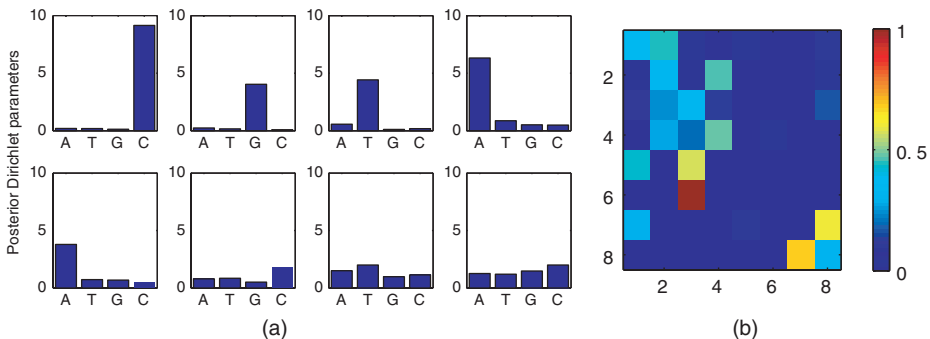


Fig. 4. (a) Dirichlet hyperparameters. Each of the eight panels represents the 4-dimensional parameter vector of a Dirichlet component (the height of the bar represents the magnitude of the corresponding element in the vector); (b) Markov transition matrix. Each element in the transition matrix B specifies the color of a rectilinear patch in the image.

the choice of eight components offer sufficient expressive power to accommodate the necessary prototypes (as the homogeneous prototypes are already somewhat redundant). Interestingly, from the learned transition model of the HMM (Fig. 4b), it can be seen that the transition probability from a homogeneous prototype to a heterogeneous prototype is significantly less than that between two homogeneous or two heterogeneous prototypes, confirming an empirical speculation in biology that motifs have the so-called *site clustering* property. Of course, with an HMM, we can only capture first-order dependencies. To model higher-order dependencies, more complex models, such as tree models, are needed.

4.2. Performance on semi-realistic sequence data

4.2.1. Single motif, and multiple instances per sequence

Under a realistic motif detection condition, the number of motif instances is unknown. Rather than trying all possible numbers of occurrences suggested by the user or decided by the algorithm and reporting a heuristically determined plausible number, **LOGOS** uses the global HMM model to describe a posterior distribution of motif instances, which depends on both the prespecified indicator state transition probabilities and the actual sequence y to be analyzed. Currently, the transition probabilities are empirically set at a default value to reflect our rough estimates of motif frequencies (i.e. 5%). But as more training data of annotated regulatory sequences are collected, we plan to fit these parameters in a genome-specific fashion. Due to modularity of variational inference in **LOGOS**, the locations of all instances, which are specified by the indicator sequence x , can be efficiently inferred from the variational marginal distribution $Q(x)$, a standard HMM, using *posterior decoding*, which computes the posterior expectation of x .

Table 1 summarizes the performance of three variants of **LOGOS** for single motif detection, with an unknown number of instances per sequence. We present the median false positive (FP) and false negative (FN) rates (in terms of finding each instance of the motifs within an offset of 3 bp) of motif detection experiments over 20 test datasets. Each test dataset consists of 20 sequences, each generated by planting (uniformly at random) 0–7 instances of a motif (real sites from SCPD),

Table 1. Performance of **LOGOS** for single motif detection, with unknown number of instances per sequence.

Motif name	LOGOS _{hh}		LOGOS _{ph}		LOGOS _{pu}	
	FP	FN	FP	FN	FP	FN
abf1	0.3115	0.2116	0.6774	0.1957	0.7917	0.9123
gal4	0.1569	0.1569	0.1895	0.1534	0.2917	0.7939
gcn4	0.1820	0.2355	0.6142	0.2821	0	0.9594
gcr1	0.1962	0.2134	0.3371	0.2038	0.3333	0.9437
mat	0.0723	0.0337	0.3563	0	0.5000	0.9643
mcb	0.3734	0.0910	0.3628	0.0792	0.3333	0.9431
mig1	0.0774	0	0.0854	0	0.9764	0.1000
crp	0.3768	0.3398	0.2727	0.5294	0	0.9487

together with its permuted “decoy,” in a 300–400 bp random background sequence. As Table 1 shows, **LOGOS**_{pu} yields the weakest results, losing in all eight motif detections (in terms of $(FP + FN)/2$), suggesting that the conventional PM + UI model, which is used in MEME, and with slight variation, in AlignACE and BioProspector, is not powerful enough to handle non-trivial detection tasks as posed by our testset. **LOGOS**_{ph} improves significantly over **LOGOS**_{pu}, even yielding the best performance in one case (for *mcb*), suggesting that the HMM global model we introduced indeed strengthens the motif detector. Finally, as hoped, **LOGOS**_{hh} yields the strongest results, performing best on 7 of the 8 motifs, convincingly showing that capturing the internal structures of motifs and making use of prior knowledge from known motifs, combined with the use of the HMM global model, can yield substantially improved performance. Our results are reasonably robust under different choices of the global HMM parameters.

4.2.2. Simultaneous detection of multiple motifs

Detecting multiple motifs simultaneously is arguably a better strategy than detecting one at a time followed by deleting or masking the detected motifs, especially when motif concentrations are high, because the latter strategy mistakenly treats the other motifs as background, causing potentially suboptimal estimation of both motif and background parameters. The global HMM model we propose readily handles simultaneous multiple motif detection: we only need to encode all motif states into the state space \mathbb{S} of the motif indicator x , and perform standard HMM inference. The locations of all motifs can be directly read off from the state configuration of x . Table 2 summarizes the results on 20 testsets each containing 20 sequences harboring motifs *abf1*, *gal4* and *mig1* (0–6 total instances/seq). The upper panels show the predictive performance based on the optimal (in terms of maximal log-likelihood of y from 50 independent runs of the VEM) posterior expectation of x . Note that with a HMDM local model, **LOGOS**_{hh} exhibits better performance. In the lower panels, we show the best result out of the top three predictions made by **LOGOS** (note the “ k -at-a-time” prediction yields a total of $3k$ possibly redundant

Table 2. Simultaneous multiple motif detection (median FP-FN rate over 20 testsets containing three motifs).

	LOGOS _{hh}		LOGOS _{ph}	
	FP	FN	FP	FN
abf1	0.3591	0.3274	0.7778	0.7434
gal4	0.1259	0.1714	0.3751	0.1491
mig1	0.3849	0.2243	0.3481	0
abf1	0.3841	0.2400	0.4721	0.3972
gal4	0.0926	0.0986	0.2609	0.1255
mig1	0.1250	0.0333	0.2318	0

Table 3. Simultaneous detection of three motifs, with lengths improperly specified (18, 22, and 20 bp, respectively, instead of the actual 13, 17, and 11 bp).

	LOGOS _{hh}		LOGOS _{ph}	
	FP	FN	FP	FN
abf1	0.7295	0.6667	0.8021	0.7680
gal4	0.1167	0.2042	0.2357	0.1325
mig1	0.4183	0.2128	0.8150	0.8381
abf1	0.3310	0.2804	0.5742	0.4821
gal4	0.0955	0.1222	0.1882	0.1250
mig1	0.2124	0.1327	0.3218	0.1623

motif patterns). This is close to the stochastic dictionary scenario where the predicted motif is to be identified from the optimal dictionary of the patterns resulting from the motif detection program.¹² It is expected that a human observer could easily pick out the biologically more plausible motifs when given a visual presentation of the most likely motifs suggested by a motif finder.

4.2.3. Detecting motifs of uncertain lengths

A useful property of the HMDM submodel is that it actually does not need to know the exact lengths of the motifs to be detected, since HMDM allows a motif to start (and end) with consecutive heterogeneous sites. Thus, a blurred motif boundary is permissible, especially when the resulting window is large enough to cover at least the entire length of the motif. As a result, we do not have to know the exact length of the motif, but just need to roughly guess it conservatively, during *de novo* motif detection. This is another appealing feature of **LOGOS**, which extends its flexibility. As shown in Table 3, even in simultaneous multiple motif detection, with improperly specified motif lengths, HMDM + HMM performs nearly as well as when motif lengths are precisely specified, whereas PM + HMM is not as good.

4.3. Performance on real genomic sequence data

4.3.1. Motif detection in yeast promoter regions

In this section we report a performance comparison of **LOGOS** (HMM + HMDM) with two popular motif detection programs, MEME and AlignACE, on 12 yeast genomic sequence sets gathered from the SCPD database (the selection is based on having at least a total of 5 motif instances in all sequences and the motif being independent of our training set). Each sequence set consists of multiple yeast promoter regions each about 500 bp long and containing on both strands an unknown number of occurrences of a predominant motif (but also possibly other minor motifs) as specified by the name of the dataset (Table 4, where the rightmost column gives the number of sequences in each dataset). Note that both the relatively large sizes of the

Table 4. Comparison of motif detectors on yeast promoter sequences.

Set name	LOGOS		MEME		AlignACE		Seq no.
	FP	FN	FP	FN	FP	FN	
abf1	0.7949	0.6522	1.0000	1.0000	0.5294	0.6087	20
csre	0.4444	0.1667	0.7778	0.5000	0.8000	0.5000	4
gal4	0.1333	0.0714	0.1667	0.2857	0.3333	0.1429	6
gcn4	0.3529	0.1852	1.0000	1.0000	0.3333	0.5556	9
gcr1	0.2859	0.6154	1.0000	1.0000	0.4545	0.4615	6
hstf	0.8571	0.5556	0.6000	0.5556	0.8500	0.6667	6
mat	0.4194	0	0.3750	0.5625	0.2500	0.2500	7
mcb	0.4706	0.2500	0.2000	0.3333	0.2500	0.2500	6
mig1	0.8077	0.2857	1.0000	1.0000	0.8333	0.7857	22
pho2	0.9024	0.5000	1.0000	1.0000	1.0000	1.0000	3
swi5	0.7647	0.5000	1.0000	1.0000	0.9412	0.7500	2
uash	0.8250	0.6818	1.0000	1.0000	0.9231	0.9545	18

input sequences and the possible presence of motifs other than what has been annotated make the motif finding task significantly more difficult than a semi-realistic test data or a small, well curated real test data. We use the following command to run MEME: “meme \$file -p 2 -dna -mod tcm -revcomp -nmotifs 1.” In practice, this means that we search for a DNA sequence on both strands for at most one motif, which can occur zero or more times in any given sequence. AlignACE is run with default command-line arguments nearly identical to those for MEME, with the only difference that AlignACE can return multiple predicted motifs (of which we select the best match from the top five MAP predictions). **LOGOS** is set in the multiple-detection mode and is used to make two motif predictions simultaneously. As shown in Table 4, for this non-trivial *de novo* motif detection task, **LOGOS** outperforms the other two programs by a significant margin.

4.3.2. Motif detection in *Drosophila* regulatory DNAs

In this section we report on a preliminary *de novo* motif discovery analysis of the regulatory regions of the 9 *Drosophila* genes involved in body segmentation. The input data consists of 9 DNA sequences ranging from 512 to 5218 bp, as described in Berman *et al.*⁴ Biologically identified motifs include *bcd*, *cad*, *hb*, *kni* and *kr*. For comparison, we provide the PWMs postulated by Berman *et al.* for these five motifs, which were used in their *motif scan* analysis (Fig. 5). The sources of all PWMs are biologically identified sequence segments from the literature (which are unaligned, ranging from 5 to 93 instances per motif, and about 20 ~ 40 bases in length). The PWMs are derived from an alignment of all these identified motif sequences.

We apply **LOGOS** (which is set to identify four motifs at a time) to the *Drosophila* dataset and Fig. 6 gives a partial list of the top-scoring motif patterns (of the top three runs out of a total of 50 runs, evaluated by the likelihood under the

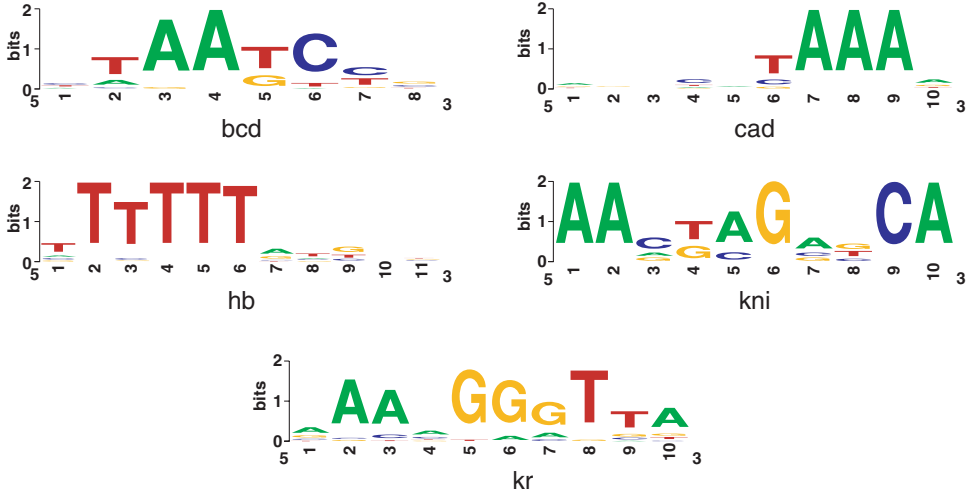


Fig. 5. Berman *et al.*'s *Drosophila* motif patterns derived from multi-alignments of biologically identified motif instances.

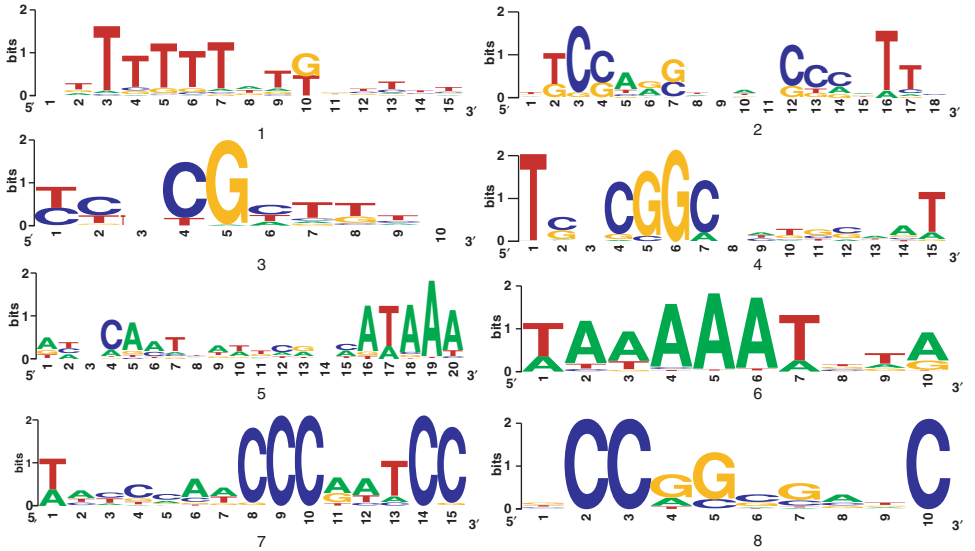


Fig. 6. Motif patterns detected by LOGOS in the regulatory regions of 9 *Drosophila* genes.

LOGOS model at convergence). Note that the *logos* shown here are not the conventional sequence *logos* based on counts of aligned nucleotides; instead we use the *logo* visualization software to graphically present the **Bayesian estimate** of the position-specific multinomial parameters θ of each motif, so they are not necessarily equal to the usual nt frequencies of aligned sequences, but represent a more robust

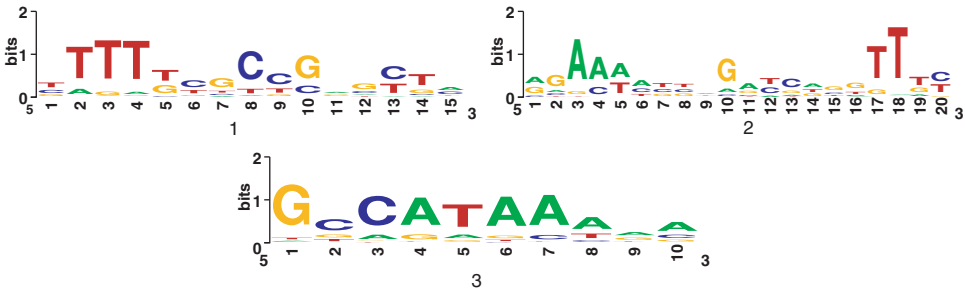


Fig. 7. Motif patterns detected by MEME in the regulatory regions of the *Drosophila* *eve*-skipped gene.

probabilistic model of the motif sequences. A visual inspection reveals that patterns 1 and 5 correspond to the *hb* and *cad* binding sites, respectively (as confirmed by the matching of the locations of our results and the sequence annotations). Part of pattern 2 agrees with the reverse complement of the *kr* motif (containing -CCCxTT-), but this motif seem to be actually a “two-block” motif because the pattern we detected under a longer estimated motif length contains an additional co-occurring conserved pattern a few bases upstream. Part of pattern 7 is close to the *bcd* motif (containing -AATCC-) but also contains additional sites (i.e., the three highly conserved C’s upstream), which turned out to be resulted from a number of false positive substrings picked up together with the true *bcd* motifs. A careful examination of pattern 6 suggests that it may be actually derived from putative motif subsequences that correspond to the *kni* binding site. This is not obvious at first because it appears quite different from the *kni* logo in Fig. 5. But after seeing an example *kni* site in stripe2/7: 5’agaaaactagatca3’, starting at position 35, we realized that the answer might be plausible. The discrepancy is likely due to the artifacts in the original generation of the alignment data supporting the *kni* logo: only 5 biologically identified instances were used and they are quite diverse; the resulting multiple alignment is visually sub-optimal in that homogeneous sites are severely interspersed with heterogeneous sites. Patterns 3, 4, and 8 are putative motifs not annotated in the input sequences. We also ran the same dataset through MEME (also four patterns to be found a time) and the output is in general weaker and harder to interpret. Figure 7 shows the best three patterns, from which one could recognize a *hb* (pattern 1) and a *cad* (pattern 3). Note that the motif logos given in Fig. 5 are based on the nucleotide-frequency profiles of biologically identified instances from many sources. Thus it is not surprising that some of the patterns we found are similar but do not match the logos in Fig. 5 exactly since our logos are derived from Bayesian estimates of the motif parameters and our data source consists of a small number of regulatory regions of the *Drosophila* genome, which might be smaller and less representative compared to the data source underlying Fig. 5 (except for *kni*).

5. Conclusions

We have presented a generative probabilistic framework for modeling motifs in biopolymer sequences. A modular architecture is proposed, which consists of a local submodel of motif alignment, and a global submodel of motif distribution.

We use an HMDM model for local motif alignment, which captures site dependencies inside motifs and incorporates learnable prior knowledge from known motifs for Bayesian estimation of the PWMs of novel motifs in unseen sequences. We use an HMM model for the global motif distribution, which introduces simple dependencies among motif instances and allows efficient and consistent inference of motif locations. A deterministic algorithm, variational EM, is developed to solve the complex missing value and Bayesian inference problems associated with our model. VEM allows probabilistic inference in the local alignment and the global distribution submodels to be carried out virtually separately with a proper Bayesian interface connecting the two processes. This *divide and conquer* strategy makes it much easier to develop more sophisticated models for various aspects of motif analysis without being overburdened by the daunting complexity of the full motif problem.

As discussed at length in a previous paper,²⁶ HMDM is a hierarchical Bayesian model that describes a structured prior distribution of PWMs, so that dependencies between sites within a motif can be modeled in the finite space of “nt-distribution prototypes.” A recent paper by Barash *et al.* proposed several expressive Bayesian network representations (e.g. tree network, mixture of trees, etc.) for motifs, which also aimed at modeling dependencies between motif sites.³ An important difference between these two approaches is that, in Barash’s Bayesian network representations, the site-dependencies are modeled directly at the level of site-specific nt-distributions in a “sequence-context dependent” way; whereas in the HMDM model, the site-dependencies are modeled at the level of the **prior distributions** of the site-specific nt-distributions in a “conservation-context dependent” way. Thus, Barash’s motif models have one-to-one correspondence with particular motif consensus patterns, and need to be trained on an one-model-per-motif basis. On the other hand, the HMDM model corresponds to a generic signature structure at the meta-sequence level, and is not meant to commit to any specific consensus motif sequence, but aims at generalizing across different motifs bearing similar conservation structures (e.g. a bell-shape). In terms of the resulting computational task in *de novo* motif detection, Barash’s model needs to be estimated in an *unsupervised* fashion and makes no use of the biologically identified motifs in the database, whereas HMDM helps to turn the model estimation task into a *semi-unsupervised* learning problem that draws connection between novel motifs to be found and the biologically identified motifs via a shared Bayesian prior, so that the patterns to be found are biased toward biologically more plausible motifs. It is interesting to note that these two approaches are complementary in that Barash’s models provide a more expressive likelihood model of the motif instances, and the HMDM model can be straightforwardly generalized to define a prior distribution

for these more expressive models (e.g. replacing the Markov chain for the prototype sequence in HMDM with a tree model and/or introducing Dirichlet mixture priors for the parameters of Barash’s models).

Due to the functional diversity of the DNA motifs, we expect that there could exist more complex dependencies and regularities in the structures of motifs, thus further investigations into these properties and more powerful local models for motifs (e.g. the combination of HMDM to expressive motif representations proposed above) are needed. Similarly, the HMM global model we propose is only a first step beyond the conventional UI model, and is only able to capture dependencies between motifs and motif clusters at a very limited level (e.g. it cannot model higher order dependencies such as hierarchical structures and long-distance influence between motifs). More expressive models are needed to achieve these goals. Nevertheless, under the **LOGOS** architecture, extensions from baseline models are modular and the probabilistic calculations involved can also be handled in a *divide-and-conquer* fashion via generalized mean-field inference. We are in the process of developing more expressive versions of **LOGOS**. In particular, recent work by Gupta *et al.*¹² have motivated us to pursue combination of the dictionary-based models with our approach to capture richer motif properties in complex sequences. We are optimistic that **LOGOS** can serve as a flexible framework for motif analysis in biopolymer sequences.

Acknowledgments

We thank Prof. Michael Eisen for helpful discussions on motif structures, and two anonymous reviewers for careful examination of the manuscript and for many valuable comments and suggestions.

Appendix A. Multinomial Distributions and Dirichlet Priors

To model a categorical random variable z , which can take J possible discrete values (e.g. all four possible nucleotides, A, C, G and T, in a DNA sequence), a standard distribution is the **multinomial distribution**: $p(z = j|\theta) = \theta_j$, $|\theta| = \sum_{j=1}^J \theta_j = 1$, $\theta_j > 0, \forall j$, where j represents one of the J possible values. The (column) vector $\theta = [\theta_1, \dots, \theta_J]'$ is called the multinomial parameters.ⁱ For a set of M *i.i.d.* samples of z , $\mathbf{z} = (z_1, \dots, z_M)$ (e.g. a whole column of nucleotides in a multi-alignment **A**), the sufficient statistics are the counts of each possible value: $h_j = \sum_{m=1}^M \delta(z_m, j)$,

ⁱNote that for simplicity, in this section, we reuse the symbol θ (and also h and α in the sequel) to denote a single column vector, whose elements are singly subscripted (e.g. θ_j); whereas in the main text and the next section, these symbols each denotes a two-dimensional array consisting of a sequence of columns vectors, whose elements are consequently doubly subscripted (e.g. θ_{lj}).

where $\delta(a, b) = 1$ if $a = b$ and 0 otherwise. Under a multinomial distribution, the likelihood of a single sample z_m is:

$$p(z_m|\theta) = \prod_{j=1}^J [\theta_j]^{\delta(z_m, j)}, \tag{A.1}$$

and the joint likelihood of the *i.i.d.* sample set \mathbf{z} is:

$$p(\mathbf{z}|\theta) = \prod_{m=1}^M \prod_{j=1}^J [\theta_j]^{\delta(z_m, j)} = \prod_{j=1}^J [\theta_j]^{h_j}. \tag{A.2}$$

To model uncertainty about the multinomial parameters, we can treat θ as a multivariate continuous random variable, and use a **Dirichlet Density** to define a prior distribution $\text{Dir}(\alpha)$ for θ :

$$p(\theta|\alpha) = C(\alpha) \prod_{j=1}^J [\theta_j]^{\alpha_j - 1}, \tag{A.3}$$

where the hyperparameters $\alpha = [\alpha_1, \dots, \alpha_J]'$, $\alpha_j > 0, \forall j$ are called the Dirichlet parameters, and $C(\alpha)$ is the normalizing constant which can be computed analytically:

$$C(\alpha) = \frac{\Gamma(|\alpha|)}{\prod_{j=1}^J \Gamma(\alpha_j)}. \tag{A.4}$$

Now we can calculate the joint probability $p(\theta, \mathbf{z}|\alpha)$:

$$p(\theta, \mathbf{z}|\alpha) = p(\mathbf{z}|\theta)p(\theta|\alpha) = C(\alpha) \prod_{j=1}^J [\theta_j]^{\alpha_j + h_j - 1}. \tag{A.5}$$

Integrating Eq. (A.5) over θ , we obtain the marginal likelihood:

$$p(\mathbf{z}|\alpha) = \int p(\theta, \mathbf{z}|\alpha) d\theta = \frac{\Gamma(|\alpha|)}{\Gamma(|\alpha| + |h|)} \prod_{j=1}^J \frac{\Gamma(\alpha_j + h_j)}{\Gamma(\alpha_j)}. \tag{A.6}$$

From Eq. (A.6) we can see that the quantity $\alpha_j - 1$ represents the number of imaginary counts that event ($z = j$) has already occurred. Furthermore, we have posterior distribution $p(\theta|\mathbf{z}, \alpha) = p(\theta, \mathbf{z}|\alpha)/p(\mathbf{z}|\alpha) = \text{Dir}(\alpha + h)$, which is isomorphic to the prior distribution, thus analytically integrable. This isomorphism between the prior and posterior is called *conjugacy* and priors of such nature are called *conjugate priors*.

Appendix B. Estimating Hyper-Parameters in HMDM

We can compute the maximum likelihood estimation of the hyper-parameters $\Theta = \{\pi, B, \alpha\}$ of the HMDM model from a training dataset of known motifs using an EM algorithm. This approach is often referred to as empirical Bayes parameter estimation.

Following Sjölander *et al.*,²³ for a given set of multi-alignment matrices $\{\mathbf{A}^{(1)}, \dots, \mathbf{A}^{(K)}\}$, where each $\mathbf{A}^{(k)}$ represents a multiple alignment of M_k biologically identified instances of motif k of length L_k , the likelihood of the count vector $h_l^{(k)}$ summarizing the column of aligned nucleotides at site l of motif k , under Dirichlet prior α_i , is

$$p(h_l^{(k)}|\alpha_i) = \frac{\Gamma(|h_l^{(k)}| + 1)\Gamma(|\alpha_i|)}{\Gamma(|h_l^{(k)}| + |\alpha_i|)} \prod_{j=1}^4 \frac{\Gamma(h_{lj}^{(k)} + \alpha_{ij})}{\Gamma(h_{lj}^{(k)} + 1)\Gamma(\alpha_{ij})}. \quad (\text{B.1})$$

Note that this formula is slightly different from Eq. (A.6) because $h_l^{(k)}$ can be resulted from $\frac{\Gamma(|h_l^{(k)}| + 1)}{\prod_{j=1}^4 \Gamma(h_{lj}^{(k)} + 1)}$ distinct permutations of the M_k nucleotides. Since no particular ordering of the motif instances in multi-alignment matrices is assumed for the training data, it is more appropriate to model the probability of the count matrices h resulted from \mathbf{A} than that of \mathbf{A} themselves.²³

Thus, the complete log likelihood of the count matrices $h^{(k)} = \{h_1^{(k)}, \dots, h_{L_k}^{(k)}\}$, $\forall k$, and the latent HMDM state sequences $q^{(k)} = \{q_1^{(k)}, \dots, q_{L_k}^{(k)}\}$, $\forall k$, can be obtained by replacing the $\mathbf{A}^{(k)}$'s in Eq. (2) with $h^{(k)}$'s, integrating over each $\theta^{(k)}$ (which results in a term like Eq. (B.1) for each count vector), and taking the logarithm of the resulting marginal:

$$\begin{aligned} l_c(\Theta) &= \log p(h^{(1)}, \dots, h^{(K)}, q^{(1)}, \dots, q^{(K)}|\Theta) \\ &= \log \left\{ \prod_{k=1}^K \left[p(q_1^{(k)}|\Theta) \prod_{l=1}^{L_k-1} p(q_{l+1}^{(k)}|q_l^{(k)}, \Theta) \prod_{l=1}^{L_k} p(h_l^{(k)}|q_l^{(k)}, \Theta) \right] \right\} \\ &= \sum_{k=1}^K \sum_{i=1}^I \delta(q_1^{(k)}, i) \log \pi_i + \sum_{k=1}^K \sum_{l=1}^{L_k-1} \sum_{i,i'=1}^I \delta(q_l^{(k)}, i) \delta(q_{l+1}^{(k)}, i') \log B_{i,i'} \\ &\quad + \sum_{k=1}^K \sum_{l=1}^{L_k} \sum_{i=1}^I \delta(q_l^{(k)}, i) \left(\log \frac{\Gamma(|h_l^{(k)}| + 1)\Gamma(|\alpha_i|)}{\Gamma(|h_l^{(k)}| + |\alpha_i|)} + \sum_{j=1}^4 \log \frac{\Gamma(h_{lj}^{(k)} + \alpha_{ij})}{\Gamma(h_{lj}^{(k)} + 1)\Gamma(\alpha_{ij})} \right). \end{aligned} \quad (\text{B.2})$$

The EM algorithm is essentially a coordinate ascent procedure that maximizes the expected complete log likelihood $E_{Q(q)}(l_c(\Theta))$ (also written as $\langle l_c(\Theta) \rangle_Q$) over the distribution $Q(q)$ and the parameter Θ . In the E step, we seek $Q(q) = \arg \max_Q \langle l_c(\Theta) \rangle_Q$, which turns out to be $Q(q) = p(q|h, \Theta) = \prod_k p(q^{(k)}|h^{(k)}, \Theta)$. Thus the E step is equivalent to computing $\langle l_c(\Theta) \rangle_{p(q|h, \Theta)}$, which reduces to replacing the sufficient statistics dependent on $q^{(k)}$ in Eq. (B.2) by their expectations with respect to $p(q^{(k)}|h^{(k)}, \Theta)$. In the M step, we compute $\Theta = \arg \max_{\Theta} \langle l_c(\Theta) \rangle_Q$. Specifically, we iterate between the following two steps until convergence:

E step:

- Compute the posterior probabilities $p(q_l^{(k)}|h^{(k)})$ of the hidden states, and the matrix of co-occurrence probabilities $p(q_l^{(k)}, q_{l+1}^{(k)}|h^{(k)})$ of each motif k , using

the *forward-backward* algorithm in a hidden Markov model with initial and transition probabilities defined by $\{\pi, B\}$ and emission probabilities defined by $p(h_l^{(k)} | q_l^{(k)} = i) = p(h_l^{(k)} | \alpha_i)$ (i.e., Eq. (B.1)).

M step:

- Baum–Welch update for the HMM parameters $\{\pi, B\}$ based on expected sufficient statistics computed from all the $p(q_l^{(k)} | h^{(k)})$ and $p(q_l^{(k)}, q_{l+1}^{(k)} | h^{(k)})$:

$$\pi_i = \frac{\sum_{k,l} p(q_l^{(k)} = i | h^{(k)})}{\sum_k L_k} \quad (\text{B.3})$$

$$B_{i,j} = \frac{\sum_{k,l} p(q_l^{(k)} = i, q_{l+1}^{(k)} = j | h^{(k)})}{\sum_{k,l} \sum_j p(q_l^{(k)} = i, q_{l+1}^{(k)} = j | h^{(k)})} \quad (\text{B.4})$$

- Gradient ascent for the Dirichlet parameters: (To force the Dirichlet parameters to be positive, we reparameterize the Dirichlet parameters as $\alpha_{ij} = e^{w_{ij}}, \forall i, j$, as described by Sjölander *et al.*²³)

$$w_{ij} = w_{ij} + \eta \frac{\partial \langle l_c(\Theta) \rangle}{\partial w_{ij}} \quad (\text{B.5})$$

where

$$\begin{aligned} \frac{\partial \langle l_c(\Theta) \rangle}{\partial w_{ij}} &= \frac{\partial \langle l_c(\Theta) \rangle}{\partial \alpha_{ij}} \frac{\partial \alpha_{ij}}{\partial w_{ij}} \\ &= \sum_{k=1}^K \sum_{l=1}^{L_k} \alpha_{ij} p(q_l^{(k)} = i | h^{(k)}) \\ &\quad \times (\Psi(|\alpha_i|) - \Psi(|h_l^{(k)}| + |\alpha_i|) + \Psi(h_{l_j}^{(k)} + \alpha_{ij}) - \Psi(\alpha_{ij})) \end{aligned}$$

and η is the learning rate, usually set to be a small constant.

References

1. B. Alberts, A. Johnson, J. Lewis, M. Raff, K. Roberts and P. Walter, *Molecular Biology of the Cell, 4th Edition*. Taylor and Francis (2002).
2. T. L. Bailey and C. Elkan, “Unsupervised learning of multiple motifs in biopolymers using EM,” *Machine Learning* **21**, 51–80 (1995).
3. Y. Barash, G. Elidan, N. Friedman and T. Kaplan, “Modeling dependencies in protein-DNA binding sites,” In *RECOMB* (2003).
4. B. P. Berman, Y. Nibu, B. D. Pfeiffer, P. Tomancak, S. E. Celniker, M. Levine, G. M. Rubin and M. B. Eisen, “Exploiting transcription factor binding site clustering to identify cis-regulatory modules involved in pattern formation in the *Drosophila* genome,” *Proc. Natl. Acad. Sci. USA* **99**, 757–762 (2002).
5. H. Bussemaker, H. Li and E. Siggia, “Building a dictionary for genomes: Identification of presumptive regulatory sites by statistical analysis,” *Proc. Natl. Acad. Sci. USA* **97** (2000).
6. H. J. Bussemaker, H. Li and E. D. Siggia, “Regulatory element detection using correlation with expression,” *Nat. Genet.* **27**(2), 167–171 (2001).
7. E. H. Davidson, *Genomic Regulatory Systems*. Academic Press (2001).

8. B. Efron, "Empirical Bayes methods for combining likelihoods (with discussion)," *J. Amer. Statist. Assoc.* **91**, 538–565 (1996).
9. M. Eisen, "Structural properties of transcription factor-DNA interactions and the inference of sequence specificity," submitted (2003).
10. M. C. Frith, U. Hansen and Z. Weng, "Detection of cis-element clusters in higher eukaryotic DNA," *Bioinform.* **17**, 878–889 (2001).
11. D. GuhaThakurta and G. D. Stormo, "Identifying target sites for cooperatively binding factors," *Bioinform.* **17**, 608–621 (2001).
12. M. Gupta and J. Liu, "Discovery of conserved sequence patterns using a stochastic dictionary model," *J. Amer. Statist. Assoc.* **98** (2003).
13. J. V. Helden, A. Rios and J. Collado-Vides, "Discovering regulatory elements in non-coding sequences by analysis of spaced dyads," *Nucleic Acids Res.* **28**, 1808–1818 (2000).
14. G. Z. Hertz and G. D. Stormo, "Identifying DNA and protein patterns with statistically significant alignments of multiple sequences," *Bioinform.* **15**, 563–577 (1999).
15. H. Huang, M. Kao, X. Zhou, J. S. Liu and W. H. Wong, "Determination of local statistical significance of patterns in Markov with application to promoter element identification," *Journal of Computational Biology* (in press).
16. M. I. Jordan, Z. Ghahramani, T. S. Jaakkola and L. K. Saul, "An introduction to variational methods for graphical models," in M. I. Jordan (ed.) *Learning in Graphical Models*. MIT Press, Cambridge (1999).
17. K. Kechris, E. van Zwet, P. Bickel and M. Eisen, "Detecting DNA regulatory motifs by incorporating position-specific base conservation," *submitted* (2003).
18. J. Liu, A. Neuwald and C. Lawrence, "Bayesian models for multiple local sequence alignment and Gibbs sampling strategies," *J. Amer. Statist. Assoc.* **90**, 1156–1169 (1995).
19. X. Liu, D. Brutlag and J. Liu, "Bioprospector: Discovering conserved DNA motifs in upstream regulatory regions of co-expressed genes," *Proc. of Pac. Symp. Biocomput.*, pp. 127–138 (2001).
20. X. S. Liu, D. L. Brutlag and J. S. Liu, "An algorithm for finding protein-DNA binding sites with applications to chromatin immunoprecipitation microarray experiments," *Nat. Biotechnol.* **20**(8), 835–839 (2002).
21. F. P. Roth, J. D. Hughes, P. W. Estep and G. M. Church, "Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation," *Nat. Biotechnol.* **16**(10), 939–945 (1998).
22. T. D. Schneider and R. M. Stephens, "Sequence logos: A new way to display consensus sequences," *Nucl. Acids Res.* **18**, 6097–6100 (1990).
23. K. Sjölander, K. Karplus, M. Brown, R. Hughey, A. Krogh, I. Mian and D. Haussler, "Dirichlet mixtures: A method for improving detection of weak but significant protein sequence homology," *Computer Applications in the Biosciences* **12** (1996).
24. G. D. Stormo and D. S. Fields, "Specificity, free energy and information content in protein-DNA interactions," *Trends in Biochemical Sciences* **23**, 109–113 (1998).
25. L. Stryer, *Biochemistry* (4th edition). W. H. Freeman and Company (1995).
26. E. P. Xing, M. I. Jordan, R. M. Karp and S. Russell, "A hierarchical Bayesian Markovian model for motifs in biopolymer sequences," *Proc. of Advances in Neural Information Processing Systems* **16** (2003).
27. E. P. Xing, M. I. Jordan and S. Russell, "A generalized mean field algorithm for variational inference in exponential families," *Proceedings of the 19th Annual Conference on Uncertainty in AI* (2003).



Eric Xing received his B.S. with honors in Physics and Biology from Tsinghua University, his Ph.D. in Molecular Biology and Biochemistry from Rutgers University and will complete his Ph.D. in Computer Science at UC Berkeley in May 2004. His early work in molecular biology focused on the genetic mechanisms of human carcinogenesis and the mutational spectrum of tumor suppressor genes. Then he moved into the field of statistical machine learning and has focussed on probabilistic graphical models, approximate inference and pattern recognition. He is interested in studying biological problems in systems biology, genetics and evolution using statistical learning approaches, and in theory and application of graphical models, nonparametric Bayesian analysis and semi-supervised learning.

Wei Wu is currently a post doctoral fellow at the Lawrence Berkeley National Laboratory under Dr. I. Saira Mian and Dr. Mina Bissell. She received her Ph.D. in computational molecular biology from the joint graduate program of Rutgers University and University of Medicine and Dentistry of New Jersey in 2000. Later, she received her M.Sc. in computer science from University of California, Santa Cruz in 2002 under the supervision of Prof. David Haussler. Her areas of research are Bioinformatics and Computational Biology.



Michael Jordan is Professor in the Department of Electrical Engineering and Computer Science and the Department of Statistics at the University of California at Berkeley. He received his Masters from Arizona State University, and earned his Ph.D. from the University of California, San Diego. He was a professor at the Massachusetts Institute of Technology from 1988 to 1998. His research has spanned a number of areas in computer science and statistics, and he has published over 200 research papers in these fields. In recent years he has focused on algorithms for approximate probabilistic inference in graphical models, on kernel machines, and on applications of machine learning to problems in bioinformatics, information retrieval, and signal processing.



Richard M. Karp received the Ph.d. from Harvard University in 1959. From 1959 to 1968 he was a member of the Mathematical Sciences Department at IBM Research. From 1968 to 1994 and from 1999 to the present he has been a Professor at the University of California, Berkeley, where he held the Class of 1939 Chair and is currently a University Professor. From 1988 to 1995 and 1999 to the present he has been a Research Scientist at the International Computer

Science Institute in Berkeley. From 1995 to 1999 he was a Professor at the University of Washington.

His research interests include combinatorial, parallel and randomized algorithms, computational complexity and NP-completeness, probabilistic combinatorics and algorithmic problems in genomics and computer networking.

His honors and awards include: U.S. National Medal of Science, Turing Award, Fulkerson Prize, Harvey Prize (Technion), Centennial Medal (Harvard), Lanchester Prize, Von Neumann Theory Prize, Von Neumann Lectureship, Distinguished Teaching Award (Berkeley), Faculty Research Lecturer (Berkeley), Miller Research Professor (Berkeley), Babbage Prize and seven honorary degrees. He is a member of the U.S. National Academies of Sciences and Engineering, the American Philosophical Society and the French Academy of Sciences, and a Fellow of the American Academy of Arts and Sciences, the American Association for the Advancement of Science, the Association for Computing Machinery and the Institute for Operations Research and Management Science.