

# An Empirical Prior Improves Accuracy for Bayesian Estimation of Transcription Factor Binding Site Frequencies within Gene Promoters

Stephen A. Ramsey<sup>1,2</sup>

<sup>1</sup>Department of Biomedical Sciences, Oregon State University, Corvallis, OR, USA. <sup>2</sup>School of Electrical Engineering and Computer Science, Oregon State University, Corvallis, OR, USA.

## Supplementary Issue: Current Developments in Domestic Animal Bioinformatics

**ABSTRACT:** A Bayesian method for sampling from the distribution of matches to a precompiled transcription factor binding site (TFBS) sequence pattern (conditioned on an observed nucleotide sequence and the sequence pattern) is described. The method takes a position frequency matrix as input for a set of representative binding sites for a transcription factor and two sets of noncoding, 5' regulatory sequences for gene sets that are to be compared. An empirical prior on the frequency  $A$  (per base pair of gene-vicinal, noncoding DNA) of TFBSs is developed using data from the ENCODE project and incorporated into the method. In addition, a probabilistic model for binding site occurrences conditioned on  $\lambda$  is developed analytically, taking into account the finite-width effects of binding sites. The count of TFBS  $\beta$  (conditioned on the observed sequence) is sampled using Metropolis–Hastings with an information entropy-based move generator. The derivation of the method is presented in a step-by-step fashion, starting from specific conditional independence assumptions. Empirical results show that the newly proposed prior on  $\beta$  improves accuracy for estimating the number of TFBS within a set of promoter sequences.

**KEYWORDS:** transcription factor, binding site, Bayesian statistics, enrichment analysis, gene regulation

**SUPPLEMENT:** Current Developments in Domestic Animal Bioinformatics

**CITATION:** Ramsey. An Empirical Prior Improves Accuracy for Bayesian Estimation of Transcription Factor Binding Site Frequencies Within Gene Promoters. *Bioinformatics and Biology Insights* 2015:9(S4) 59–69 doi: 10.4137/BBI.S29330.

**TYPE:** Methodology

**RECEIVED:** May 11, 2016. **RESUBMITTED:** September 11, 2016. **ACCEPTED FOR PUBLICATION:** September 18, 2016.

**ACADEMIC EDITOR:** J. T. Efrid, Editor in Chief

**PEER REVIEW:** Eight peer reviewers contributed to the peer review report. Reviewers' reports totaled 1158 words, excluding any confidential comments to the academic editor.

**FUNDING:** This study was supported by award number K25HL098807 from the National Heart, Lung and Blood Institute, the Medical Research Foundation of Oregon (New Investigator Grant award), the National Science Foundation (award numbers 1557605-DMS and 1553728-DBI), the PhRMA Foundation (Research Starter Grant), and Oregon State University (Division of Health Sciences Interdisciplinary Research Grant award).

The author confirms that the funder had no influence over the study design, content of the article, or selection of this journal.

**COMPETING INTERESTS:** Author discloses no potential conflicts of interest.

**CORRESPONDENCE:** stephen.ramsey@oregonstate.edu

**COPYRIGHT:** © the authors, publisher and licensee Libertas Academica Limited. This is an open-access article distributed under the terms of the Creative Commons CC-BY-NC 3.0 License.

Paper subject to independent expert blind peer review. All editorial decisions made by independent academic editor. Upon submission manuscript was subject to anti-plagiarism scanning. Prior to publication all authors have given signed confirmation of agreement to article publication and compliance with all applicable ethical and legal requirements, including the accuracy of author and contributor information, disclosure of competing interests and funding sources, compliance with ethical requirements relating to human and animal study participants, and compliance with any copyright requirements of third parties. This journal is a member of the Committee on Publication Ethics (COPE).

Published by Libertas Academica. Learn more about this journal.

## Introduction

In bioinformatics, an enduring and fundamental question is how best to use an organism's genome sequence as well as prior knowledge of the DNA sequence preferences of transcription factors (TFs) in order to determine which TFs are responsible for an observed pattern gene expression differences between sample groups, such as tissues at different stages of disease and cells cultured in the presence or absence of a chemical stimulus.<sup>1–3</sup> The general approach of computationally analyzing noncoding DNA sequences within 5' (upstream) regions of differentially expressed gene sets to identify statistically overrepresented TF binding site (TFBS) sequence matches – known as TFBS enrichment analysis<sup>4–8</sup> – has proved useful for identifying the gene regulatory mechanisms from transcriptome data.<sup>9–14</sup> Databases such as MatBase,<sup>15</sup> TRANSFAC,<sup>16</sup> JASPAR,<sup>17</sup> UniPROBE,<sup>18</sup> Factorbook,<sup>19</sup> and FootprintDB<sup>20</sup> are rapidly accumulating position-nucleotide frequency matrices (PFMs) that represent the sequence preferences of individual TFs. This rapid accumulation is driven by high-throughput assays such as ChIP-seq and protein-binding microarrays and through the use of improved *in silico* structural models for

predicting TF–DNA affinities. Although a ChIP-seq assay can be used to map the binding sites of a specific TF genome-wide within a specific cell type or tissue,<sup>21</sup> only a small percentage of known TFs have been successfully assayed using this technique. In vertebrates, there have been relatively few reports of applications of ChIP-seq outside of humans and model species such as mouse.<sup>22,23</sup> Thus, the approach of computationally analyzing a set of 5' regulatory sequences to measure the enrichment of TFBS – leveraging databases of *known* TFBS sequence patterns – remains unmatched in terms of the number of TFBS sequence patterns that can be simultaneously analyzed. This discovery power is particularly important in vertebrates, for which there are ~1800 different TFs, of which hundreds can be expressed in any given cell type or tissue.<sup>24</sup>

Reflecting the importance of this problem, multiple computational approaches have been proposed for PFM-guided detection of enrichment of TFBS within gene-vicinal sequences.<sup>7,25–27</sup> For the purpose of specificity, I define *gene-vicinal* to mean within approximately 5 kbp (in either direction) of a transcription start site.<sup>28</sup> The TFBS enrichment analysis method of Frith et al.<sup>7</sup> involves the direct use of the



position-probability matrix (PPM, which is the row-normalized PFM) in order to compute a likelihood ratio of the PPM model to a nucleotide frequency-based background model, for a binding site-sized sequence window at a given position. The likelihood ratios are then averaged over all nucleotide positions within a single gene-vicinal sequence to obtain a single-gene score. For each possible subset of genes from the gene set, the product of gene-level scores is computed, and these subset-level scores are averaged.<sup>7</sup> In another approach, Ho Sui et al.<sup>25</sup> used a log-likelihood ratio approach with an empirically determined hard threshold in order to identify TFBS and then used the binomial distribution to test the enrichment of TFBS. Sinha and Tompa<sup>29</sup> used a multi-TF approach in which the weighted sum of occurrences of a specific TF's PPM was computed over binding site configurations for all TF PPMs to be analyzed. The prior on the expected number of binding sites is not treated probabilistically but is a fixed parameter value. Pavesi and Zambelli<sup>27</sup> rescaled the positional log-likelihood score in order to map the score to a compact interval and then computed the maximum of this rescaled score at all positions within a gene-vicinal sequence; this per-gene score is then averaged over all genes in the gene set. The diversity of methods for PFM-guided TFBS enrichment analysis and the significant numbers of studies (over 600 combined, for Refs. 7 and 25) that have reported using these methods underscores the importance of this problem in the field of bioinformatics.

Despite its discovery power, TFBS enrichment analysis using prior TF binding pattern information in the form of PFMs has a fundamental challenge that PFMs are highly variable in terms of their specificity for nucleotide sequences and in terms of the uncertainty of the composition of the corresponding PPMs.<sup>30,31</sup> Within databases of TFBS sequence patterns, the numbers of representative binding sites from which individual vertebrate TF PFMs have been compiled can vary by four orders of magnitude, from half a dozen to tens of thousands of representative oligomer sequences.<sup>15–17</sup> For cases of TFs with highly specific nucleotide affinity and/or very low sampling of representative binding site sequences, PFM counts of zero pose a problem in the standard PPM-based approach and necessitate the use of ad hoc pseudocounts to enable the scoring of nucleotide sequences that do not perfectly match the TFBS consensus sequence.<sup>32,33</sup> Furthermore, because the precision of the PPM that is associated with a PFM depends directly on the number of representative binding site sequences used to compile the PFM,<sup>30</sup> TFBS enrichment analysis using only the PPM (and not taking into account the uncertainty in the PPM's structure) can be a source of both type I and type II errors. Finally, in order to assess the significance of a finding that the frequency of PPM sequence matches for a TF is statistically overrepresented for 5' upstream sequences for a gene set versus for a background set of genes, it is necessary to quantify the *magnitude* of the frequency enrichment and not just statistical significance (eg, using a *P* value). In addition, it is useful to be able to estimate the uncertainty on the magnitude

of the TFBS frequency enrichment. A Bayesian approach to TFBS frequency estimation, as described below, has the potential to address the challenge of highly variable accuracy (sharpness) of known TFBS motifs. Bayesian methods have long been used for de novo motif discovery<sup>34–37</sup> and have also been proposed for TFBS recognition and demonstrated to have improved accuracy over traditional motif scanning.<sup>30</sup> In the context of PFM-guided enrichment analysis, a Bayesian approach is appealing because it could account for uncertainty in the PPM and it could provide an estimate of the TFBS frequency per base pair of noncoding DNA, while appropriately weighting high-quality and low-quality matches to the PPM. By using a Bayesian approach, an additional benefit is that an empirical prior distribution of TFBS frequencies (across many TFs) can be included in the model to improve the TFBS frequency estimation in the case of a weak (ie, degenerate) PPM.

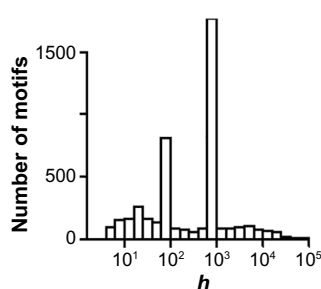
In this article, I describe a Bayesian approach to PFM-guided TFBS enrichment analysis, which produces samples from the posterior distribution of the number of TFBS for a given PFM, within a given sequence. The method incorporates an empirical prior on the per base pair TFBS frequency that is informed by the analysis of human TFBS from the ENCODE project (as opposed to the geometric prior used in a previous study<sup>30</sup>). Finally, because the method is developed from an explicit joint probability model of all of the observables and model parameters, the method could be readily extended to incorporate other types of regulatory potential scores.<sup>38,39</sup> I show empirical results from applying the new prior to estimate the number of TFBS for a synthetic set of promoter sequences in which representative TFBS sequences are introduced. The empirical results show that the new prior improves accuracy when compared to a previously proposed prior on the per-promoter number of TFBS.

### Mathematical Preliminaries and Notation

For the purpose of detecting TFs that may control a given cluster of coexpressed genes, it is simplest to consider a single TF at a time; I use “TFX” as a generic symbol for this TF. (Although this article is focused on single-TF enrichment analysis for simplicity, the pairwise TF enrichment analysis could in principle be accommodated by extensions of the general approach described herein.) Consistent with a Bayesian approach, I start by framing the problem of PFM-based TFBS enrichment analysis in terms of a set of random variables including observations, nuisance parameters, and a single parameter – the number of TFBSs within a given set of gene promoters – whose distribution (conditioned on the observations) will ultimately be sampled. To do so, I introduce a bit of mathematical notation needed to define these random variables. It is convenient to denote the set of natural DNA nucleotides by integers,  $\mathbb{D} = \{1, 2, 3, 4\}$ , corresponding to A, C, G, and T (so the complementary nucleotide for nucleotide  $d \in \mathbb{D}$  is  $5 - d$ ). For simplicity, let the promoter sequences of a cluster of differentially expressed genes be concatenated and represented as a

sequence  $s \in \mathbb{D}^L$ , where  $L$  is the total sequence length in base pair. Let the noncoding DNA sequence within the promoters of a set of randomly selected genes that are expressed (but not necessarily differentially expressed) within the same cell type or tissue, be represented by  $s'' \in \mathbb{D}^{L''}$ , where  $L''$  is the length in base pair. Finally, let  $s' \in \mathbb{D}^{L'}$  be a large DNA sequence comprising noncoding, gene-vicinal sequences selected at random and in which any known TFBS (as mapped by high-throughput ChIP-seq studies) have been excluded (here again,  $L'$  is the total length, in base pair). The *background model* sequence  $s'$  will be used to obtain a model for nucleotide frequencies in noncoding, nonbinding site DNA. The *regulatory region* sequences  $s$  and  $s''$  will be analyzed for a relative enrichment of TFBSs for a given TF, as described below.

The TFX is assumed to have a set of representative binding site sequences (numbering  $c$ ; depending on the type of assay used to compile the representative binding site sequences,  $c$  could range from 6 to 100,000, as shown in Fig. 1) obtained from the literature and/or from high-throughput protein–DNA binding measurements. The representative binding site sequences are assumed to have been multiply aligned; I denote by  $w$  the length (in base pair) of the core region of overlapping representative binding sites in the multiple alignment. The counts of nucleotides of each type at each position within a PFM will be denoted by a matrix  $\mathbf{c} \in \mathbb{C}^{w \times 4}$ , where  $\mathbb{C} = \{0, 1, 2, \dots, c\}$ . I note that  $c$  and  $w$  are specific to the transcription factor TFX, and this dependence could be denoted by  $c(\text{TFX})$  and  $w(\text{TFX})$ ; however, for the simplicity of notation, I use the more compact  $c$  and  $w$ . The height of the TF matrix,  $w$ , can vary significantly from TF to TF; across all 4528 matrices in the TRANSFAC 2015 Professional database,  $w$  varies from 5 to 30, with a median



**Figure 1.** The distribution of values of  $c$  for a collection of 4,528 vertebrate transcription factors from the TRANSFAC Professional 2013.4 and JASPAR 5.0 databases. The sharp peak in the distribution at  $c = 100$  is due to the inclusion of motif matrix information for which the original sequence alignments are not available (such as from high-throughput *in vitro* protein–DNA binding screens,<sup>40</sup> for which a default value of  $c = 100$  was selected based on consistency with the number of significant digits for the values reported in the motif matrices). The sharp peak at  $c = 998$  is due to the 2,076 structure-derived motifs that were originally obtained using the 3DTF tool<sup>41</sup> and were then incorporated into the TRANSFAC database). The long tail of  $c$  values above  $10^3$  represents motif matrices compiled from high-throughput TF location assays such as ChIP-array and ChIP-seq.

of 12. The index set for sequence positions within a binding site for TFX will be denoted by  $\mathbb{G} = (1, \dots, w)$ . The factor TFX is assumed to have an overall frequency per base pair, represented by  $\lambda$ , of binding sites in a given sequence of noncoding, gene-vicinal DNA; I represent our uncertainty about  $\lambda$  by treating it as a random variable  $\Lambda$  on  $(0, \lambda_{\max}]$ , where a fixed value  $\lambda_{\max} \in (0, 1)$  is chosen as an upper limit. An absolute physical upper limit on  $\lambda_{\max}$  would be  $1/w$ , given the requirement that binding sites for TFX not overlap. However, TFBSs in mammals are in general sparsely distributed, even within gene-vicinal regions<sup>42</sup>; thus, the value  $\lambda_{\max} = 0.001 \text{ bp}^{-1}$  is used here (and as seen in the “**Modeling  $P(\beta|\lambda)p(\lambda)$** ” section, of all the TFs in the ENCODE human ChIP-seq dataset,<sup>43</sup> none has a binding site frequency per base pair in gene-vicinal sequence that exceeds  $0.001 \text{ bp}^{-1}$ ). Because the actual locations of the TFX binding sites within a given noncoding, gene-vicinal sequence of length  $L$  are not known, I denote the binding site locations by a  $\{-1, 0, 1\}^L$ -valued random variable  $B$ . Specifically, for the outcome  $B = \beta \in \{-1, 0, 1\}^L$ , at each location  $l \in \{1, \dots, L\}$ ,

$$\beta_l = \begin{cases} 1 & \text{if there is a plus-orientation TFBS whose} \\ & \text{5' end is at location } l \\ 0 & \text{if there is no TFBS whose} \\ & \text{5' end is at location } l \\ -1 & \text{if there is a minus-orientation TFBS whose} \\ & \text{5' end is at location } l, \end{cases} \quad (1)$$

where in the above, +(plus)-orientation means that the PPM pattern derived from  $c$  matches the sequence on the forward strand, and –(minus)-orientation means that the PPM pattern matches the sequence on the reverse strand. The use of a discrete parameter to represent the presence/absence of a TFBS at a given nucleotide position (rather than a continuous parameter) makes the space of TFBS configurations more tractable to explore,<sup>30,31</sup> without sacrificing the ability to differentially weight a poor-affinity binding site from a high-affinity binding site in the analysis. The total number of binding sites in a given binding site configuration  $\beta$  is obtained by the  $L^1$  norm,  $\beta \equiv \|\beta\|_1$ . For simplicity of notation, only binding site configurations  $\beta$  for which the entire binding site is contained within the range of sequence positions  $\mathbb{L}_r$ , and for which no two binding sites are overlapping by any number of nucleotides (even if the two binding sites have opposite orientations) will be allowed. Thus, the range of  $B$  is not the entirety of  $\{-1, 0, 1\}^L$ , but a subset  $\mathbb{B} \subset \{-1, 0, 1\}^L$  defined by the above constraints. In the Bayesian approach to TFBS enrichment analysis that described below,  $B = \|\mathbf{B}\|_1$  is the integer-valued random variable whose distribution (conditioned on the observed regulatory region sequences and on the PFM for representative binding sites) is sought.

Importantly, the probability distribution on the number of binding sites will depend on the length of the DNA sequence



being analyzed (longer combined regulatory sequences will, in general, contain more binding sites of a given type), and thus, the probability distribution for  $\|\mathbf{B}\|_1$  (the number of binding sites) conditioned on the DNA sequence  $s$  cannot be directly compared to the probability distribution for  $\mathbf{B}|s''$  unless  $L = L''$ . Thus, in practice, one would compare samples of  $B/L_r|s$  with samples of  $B/L''|s''$ , with  $\Lambda$  treated as a nuisance variable. A key benefit of a Bayesian approach is that it will not require a specific value for the parameter  $\lambda$ ; all possible values (consistent with the imposed constraint  $\lambda_{\max}$ ) are considered.

A Bayesian approach to analyzing whether binding sites for TFX are enriched within sequences  $s'$  versus  $s''$  can now be succinctly described as comparing samples from the distribution of

$$B/L_r | c, s, s' \tag{2}$$

with samples from the distribution of

$$B/L'' | c, s'', s', \tag{3}$$

with  $\Lambda$  marginalized, under an explicit probability model. Thus, the technical problem to be solved here is how to accurately sample from the conditional distribution

$$B | c, r, s', \tag{4}$$

where  $r$  is an arbitrary observed set of (concatenated) promoter sequences (and in practice, one set of samples would be generated for the case  $r = s$  and one set of samples would be generated for the case  $r = s''$ ). I denote the length of the sequence  $r$  by  $L_r$  (which will have the value  $L$  or  $L''$  depending on whether we are modeling the case  $r = s$  or the case  $r = s''$ ), and the sequence of unique positions within the combined gene-vicinal DNA regions by  $\mathbb{L}_r = (1, \dots, L_r)$ . Similarly, I define the sequence of positions within  $s'$  by  $\mathbb{L}'_r = (1, \dots, L')$ . Now, we can more precisely state our goal as modeling the posterior distribution  $B|r, s', c$ . In order to be able to do this, it is convenient to define a matrix-valued random variable  $\Phi$  and a vector-valued random variable  $\Psi$ . The  $w \times 4$  matrix random variable  $\Phi$  represents the PPM that is associated with the PFM  $c$ , and it is a random variable because the true probability model will always be uncertain if the number of representative binding site sequences (ie, the number  $c$ ) is finite. In keeping with a PPM model, for each sample  $\phi$  from the random variable  $\Phi$ , each row of  $\phi$  (which I denote by  $\phi_g$  where  $g \in \mathbb{G}$ ) has unit  $L^1$  norm. This means that each row  $\phi_g$  of  $\Phi$  is a random variable whose range is the unit three-simplex  $\mathbb{H}^3$ . A central assumption that makes a Bayesian analysis of TFBS enrichment tractable is that the  $\Phi_g$  are all independent random variables. The  $\mathbb{H}^3$ -valued random variable  $\Psi$  represents the nucleotide frequencies on  $s'$ , and its distribution is generally very sharply peaked since the sequence  $s'$  from which the background model is obtained is usually hundreds to thousands of kilobase pair in length.

In the application of PFM-guided TFBS enrichment analysis, the observations  $r$ ,  $c$ , and  $s'$  are known by definition; however, it is helpful in a Bayesian approach to formally define a generative model in which we can compute the probability of these observations, conditioned on  $\Lambda$ ,  $\mathbf{B}$ ,  $\Phi$ , and  $\Psi$ . Such a generative model can be more concisely defined in terms of random variables, and thus, I refer to a  $\mathbb{D}^{L_r}$ -valued random variable  $\mathbf{R}$  for which we have the observed sequence  $r$ , and a  $\mathbb{D}^{L'}$ -valued random variable  $\mathbf{S}'$  for which we have observed  $s'$ , and a  $\{1, \dots, c\}^{w \times 4}$ -valued random variable  $\mathbf{C}$  for which we have observed  $c$ . The random variables in this model are summarized in Table 1.

In order to be able to model the conditional probability of the sequence  $r$  given a PFM  $\phi$ , the specified locations of TFBS  $\beta$ , and the background nucleotide frequency model  $\psi$ , it is necessary to define a function  $\mathcal{U}: \{-1, 0, 1\}^{L_r} \rightarrow \mathcal{P}(\mathbb{L}_r)$  that maps a configuration  $\beta$  of binding sites to the set of nucleotide positions within  $\mathbb{L}_r$  that the binding sites occupy. Thus,  $\mathcal{U}(\beta)$  is the *footprint* of the binding sites whose 5' locations are specified by  $\beta$ . Let us define the set of all pairs of binding site *footprint* positions and binding site configurations by  $\mathbb{U} \subset \mathbb{L}_r \times \mathbb{B}$ :

$$\mathbb{U} = \bigcup_{\beta \in \mathbb{B}} \mathcal{U}(\beta) \times \{\beta\}. \tag{5}$$

Given a configuration of binding sites  $\beta$ , any position  $l \in \mathcal{U}(\beta)$  within one of the binding sites will correspond to a specific binding site orientation (1 or -1), and this correspondence will be denoted by a mapping  $\mathcal{F}: \mathbb{U} \rightarrow \{-1, 1\}$ ,

$$(l, \beta) \mapsto_{\mathcal{F}} \begin{cases} 1 & \text{if } l \text{ is in a forward-orientation binding site} \\ -1 & \text{if } l \text{ is in a reverse-orientation binding site.} \end{cases} \tag{6}$$

**Table 1.** Random variables in the full probability model for TFBS enrichment analysis.

VARIABLE	SAMPLE/OBSERVATION	RANGE	MEANING
$\Lambda$	$\lambda$	$(0, \lambda_{\max}]$	Frequency (per bp) of binding sites within $r$
$\mathbf{B}$	$\beta$	$\{-1, 0, 1\}^{L_r}$	Presence/absence (and orientation) of TFBS
$\Phi$	$\phi$	$(\mathbb{H}^3)^w$	The true PPM of the TF
$\Psi$	$\psi$	$\mathbb{H}^3$	Nucleotide frequencies on non-TFBS DNA
$\mathbf{R}$	$r$	$\mathbb{D}^{L_r}$	Gene-vicinal, noncoding sequence
$\mathbf{S}'$	$s'$	$\mathbb{D}^{L'}$	Background noncoding sequence
$\mathbf{C}$	$c$	$\{1, \dots, c\}^{w \times 4}$	The PFM for representative binding sites

In the case of a reverse-orientation binding site, the PPM  $\phi$  will correspond to the reverse complement of the nucleotide sequence within the binding site, in which case it is convenient to define a conditional complementation function  $\mathcal{C}: \mathbb{D} \times \{-1,1\} \rightarrow \mathbb{D}$  by

$$\mathcal{C}(d, j) = jd + \frac{5}{2}(1-j), \quad (7)$$

which is the identity on  $d$  when  $j = 1$  and which complements  $d$  when  $j = -1$ . Similarly, any configuration  $\beta$  and any position  $l$  within a binding site will correspond to a specific row of the PFM for the TF, depending on the orientation of the binding site; I denote this correspondence by a mapping  $\mathcal{G}: \mathbb{U} \rightarrow \mathbb{G}$ ,

$$(l, \beta) \xrightarrow{\mathcal{G}} \text{the index of the row of } c \text{ corresponding to position } l \in \mathcal{U}(\beta). \quad (8)$$

Finally, in order to be able to model the joint probability of  $r$  and  $s$ , it will be necessary to count nucleotides of each type (ie, A, C, G, and T) outside of TFBS as well as at different positions within the binding sites of TFX. Outside of TFBS, I represent the nucleotide counts by the 4-vector  $f$  whose elements are defined by

$$f_d = |\{l \in \mathbb{L}_r \setminus \mathcal{U}(\beta) \mid r_l = d\}| + |\{l \in \mathbb{L}' \mid s'_l = d\}|, \quad (9)$$

for all  $d \in \mathbb{D}$ . I represent the position-nucleotide counts for the sequence within all TFBS by a  $w \times 4$  matrix  $\sigma$  whose elements are

$$\sigma_{gd} = \left| \left\{ l \in \mathcal{U}(\beta) \mid \mathcal{G}(l, \beta) = g \wedge \mathcal{C}(r_l, \mathcal{J}(l, \beta)) = d \right\} \right| \quad (10)$$

for all  $g \in \mathbb{G}$  and  $d \in \mathbb{D}$ . Because of the physical constraint that binding sites for TFX cannot overlap, it follows that

$$\sum_{d \in \mathbb{D}} \sigma_{gd} = \beta, \quad (11)$$

for all  $g \in \mathbb{G}$ . In the next section, I introduce the statistical approach by defining a joint probability model.

### Bayesian Approach to TFBS Enrichment Analysis

Having defined random variables to represent all of the observed information ( $R, C, S'$ ) and the latent variables ( $\Phi, \Psi, \Lambda$ ), and the model parameter  $B$ , the first step in a Bayesian approach<sup>44</sup> is to define a simplified model for the joint probability distribution. I choose the model

$$p(r, s', \Lambda, \beta, \phi, \psi, c) = P(r \mid \phi, \psi, \beta) P(s' \mid \psi) p(\psi) P(c \mid \phi) p(\phi) P(\beta \mid \Lambda) p(\Lambda), \quad (12)$$

where the condensed notation  $P(\lambda)$  means  $P(\Lambda = \lambda)$  and so forth for the other random variables,  $P$  denotes a probability

distribution, and  $p$  denotes a probability density. Eq. 12 can be derived from first principles based on the following independence assumptions:

$$R \perp S', C, \Lambda \mid \phi, \psi, \beta \quad (13)$$

$$S' \perp \Phi, B, C, \Lambda \mid \psi \quad (14)$$

$$C \perp \Lambda, B, \Psi \mid \lambda \quad (15)$$

$$B \perp \Phi, \Psi \mid \lambda \quad (16)$$

$$\Phi \perp \Psi, \Lambda \quad (17)$$

$$\Psi \perp \Lambda \quad (18)$$

The independence structure of Eq. 12 can be summarized in graphical model notation,<sup>45</sup> as shown in Figure 2. To make the joint probability model explicit, each of the conditional probabilities in Eq. 12 will be specified below.

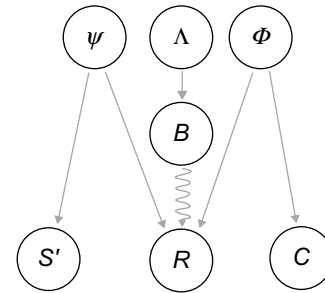
**Modeling  $P(r \mid \phi, \psi, \beta) P(s' \mid \psi) p(\psi)$ .** For PFM-guided computational recognition of TFBS, a fundamental assumption is that the probability model for the counts of nucleotides outside of TFBS is independent of the probability model for the counts of nucleotides within TFBS.<sup>32,46</sup> This means that the conditional probability of  $r$  can be expressed as the product of conditional probabilities for the subsequences of  $r$  corresponding to  $\mathcal{U}(\beta)$  (the TFBS) and corresponding to  $\mathbb{L}_r \setminus \mathcal{U}(\beta)$  (outside the TFBS). Conditioned on  $\psi$ , the nucleotide probabilities at positions outside of TFBS, which are denoted by the random variables  $\{R_l\}_{l \in \mathbb{L}_r \setminus \mathcal{U}(\beta)}$  and the random variables  $\{S'_l\}_{l \in \mathbb{L}'}$ , are assumed to satisfy

$$R_l \mid \psi \stackrel{\text{iid}}{\sim} \text{Categ}(\mathbb{D}, \psi) \quad (19)$$

for any  $l \in \mathbb{L}_r \setminus \mathcal{U}(\beta)$ , and

$$S'_l \mid \psi \stackrel{\text{iid}}{\sim} \text{Categ}(\mathbb{D}, \psi), \quad (20)$$

for any  $l \in \mathbb{L}'$  (where iid denotes independent and identically distributed). Conditioned on  $\phi$  and  $\beta$ , the nucleotide sequence



**Figure 2.** Graphical model diagram of the independence assumptions shown in Eqs. 13–18. Each arrow denotes a relationship between a *parent variable* and a *child variable*. Collectively, the variables and arrows indicate conditional independence as follows: each variable  $X$  is independent of other variables, jointly conditioned on all parents of  $X$ .



probabilities at locations within the footprint  $\mathcal{U}(\beta)$  of the binding sites specified by  $\beta$ , the nucleotide probabilities are denoted by random variables  $\{R_l\}_{l \in \mathcal{U}(\beta)}$  that are independent and distributed as follows:

$$R_l | \phi, \beta, \mathcal{G}(l, \beta) = g, \mathcal{J}(l, \beta) = 1 \sim \text{Categ}(\mathbb{D}, \phi_g) \quad (21)$$

$$5 - R_l | \phi, \beta, \mathcal{G}(l, \beta) = g, \mathcal{J}(l, \beta) = -1 \sim \text{Categ}(\mathbb{D}, \phi_g). \quad (22)$$

for any  $l \in \mathbb{L}_r \setminus \mathcal{U}(\beta)$ . Because the length of  $s'$  is assumed to be quite substantial, the distribution of  $\Psi | s'$  will be quite sharply peaked, and thus, any weak prior on  $\Psi$  will have little effect. Thus, it is reasonable to assume a uniform prior  $p(\Psi) = 3!$ . Given the uniform prior assumption for  $\Psi$  and the definition in Eqs. 9 and 10 and the assumptions in Eqs. 19–22, the conditional probability of  $R, S'$  has a compact form,

$$P(r | \phi, \Psi, \beta) P(s' | \Psi) p(\Psi) = 3! \prod_{d \in \mathbb{D}} \left[ (\psi_d)^{f_d} \left( \prod_{g \in \mathbb{G}} (\phi_{gd})^{\sigma_{gd}} \right) \right], \quad (23)$$

that will be compatible with collapsing of  $\phi$  and  $\Psi$  (as shown in the ‘‘Obtaining the distribution of  $B | r, s', c$ ’’ section).

**Modeling  $P(c | \phi) p(\phi)$ .** To account for uncertainty in the PFM  $c$  due to sampling from a finite (and in many cases, very limited) number of representative binding site sequences, the PFM is represented by a random variable  $C$ . A core assumption in the field of PFM-guided TFBS recognition is that rows of  $C$ , denoted by  $C_g$  (where  $g \in \mathbb{G}$ ), are independent and multinomial distributed with a fixed number of trials.<sup>47,48</sup> Because in some cases, some representative binding site sequences will be outside the core portion of the multiple alignment from which the PFM is tabulated, the row sums of  $c$  may in some cases be less than the count  $c$  of representative binding site sequences. Thus, to accommodate such cases, I denote by  $c_g$  the sum of the elements of row  $g$  of  $c$ . In terms of the row-specific counts  $c_g$  (for  $g \in \mathbb{G}$ ), the conditional distribution of  $C_g$  can be expressed as

$$C_g \mid \|C_g\|_1 = c_g, \Phi_g = \phi_g \sim \text{Mult}(c_g, \phi_g), \quad (24)$$

for which the formula for  $P(c | \phi)$  immediately follows

$$P(c | \phi) = \prod_{g \in \mathbb{G}} \left[ \frac{c_g!}{\prod_{d \in \mathbb{D}} c_{gd}! \prod_{d' \in \mathbb{D}} (\phi_{gd'})^{c_{gd'}}} \right]. \quad (25)$$

The most common approach for selecting the prior probability  $p(\phi)$  for the PPM is to choose a uniform prior, in which case  $p(\phi)$  is just the constant  $(3!)^w$ . Although other authors have pointed out the possibility of using an empirical prior on  $\phi$ ,<sup>30</sup> it is nontrivial to collapse  $\Phi$  by analytic integration over all  $\phi$ , in the case of a nonuniform  $p(\phi)$ , so here I assume a uniform  $p(\phi)$ .

**Modeling  $P(\beta | \lambda) p(\lambda)$ .** At a given base pair location with no binding sites nearby (and with no sequence information), I model the probability that there is a binding site – in a specific orientation – as  $\lambda/2$ . In the absence of sequence information, intuition would suggest treating the occurrence or absence of a binding site for TFX at each position in DNA as independent and identically distributed Bernoulli trials. However, because of the physical constraint that two binding sites are not allowed to overlap, each binding site (ie, each nonzero entry of  $\beta$ ) affects the probability of a binding site at nearby positions. Specifically, each binding site prevents the possibility of an overlapping binding site (in either orientation) at  $w - 1$  bp positions, and for an additional  $2(w - 1)$  flanking positions, a binding site is only possible in one orientation. Thus, the probability model consistent with the physical constraints would be

$$P(\beta | \lambda) = \mathcal{N}_1(L_r, w, \lambda) (\lambda/2)^\beta (1 - \lambda)^{L_r - w\beta - 2(w-1)\beta} \times \left( 1 - \frac{\lambda}{2} \right)^{2(w-1)\beta}, \quad (26)$$

where  $\mathcal{N}_1$  is function that is implicitly defined by the law of total probability for  $P(\beta | \lambda)$ . In the limit where  $L_r \gg w$ , and solving for  $\mathcal{N}_1$  using the law of total probability, we have the approximate result,

$$P(\beta | \lambda) \simeq \left( 1 + \frac{\lambda w \left( 1 - \frac{\lambda}{2} \right)^{2(w-1)}}{(1 - \lambda)^{3w-2}} \right)^{-\frac{L_r}{w}} \cdot \left( \frac{\lambda \left( 1 - \frac{\lambda}{2} \right)^{2(w-1)}}{2(1 - \lambda)^{3w-2}} \right)^\beta + O(w/L_r). \quad (27)$$

In the case  $w = 1$ , the above can be seen to reduce to  $\lambda^\beta (1 - \lambda)^{L_r - \beta/2}$ , which is the expected joint probability of outcome sequence  $\beta$  for  $L_r$  independent trials of the categorical distribution with outcomes  $(-1, 0, 1)$  with probabilities  $(\lambda/2, 1 - \lambda, \lambda/2)$ , in which  $\beta$  trials have a nonzero outcome.

The prior distribution  $p(\lambda)$  reflects the range and relative probability of different  $\Lambda$  values for TFX, before the sequence  $r$  has been taken into account. The prior  $p(\lambda)$  is important because for real-world applications, it can exert a significant effect on the distribution of  $\Lambda | r, c$ . For mammals, the prior  $p(\lambda)$  can be formulated empirically using binding site frequencies (per base pair of noncoding, gene-vicinal DNA sequence) for 620 human TF ChIP-seq experiments (comprising 119 distinct TFs) obtained from the ENCODE project.<sup>43</sup> For each ChIP-seq experiment, binding sites within regions of noncoding DNA within  $-1500$  to  $+500$  bp transcription start sites of VEGA transcripts (from Ensembl Release 75, GRCh37 assembly coordinates) were

mapped, using ChIP-seq peak data that were peak-called using the SPP program<sup>49</sup> and for which the data files were downloaded from the ENCODE data access page at the European Bioinformatics Institute from the June 2012 release ([http://ftp.ebi.ac.uk/pub/databases/ensembl/encode/integration\\_datajan2011/byDataType/peaks/june2012/spp/optimal/](http://ftp.ebi.ac.uk/pub/databases/ensembl/encode/integration_datajan2011/byDataType/peaks/june2012/spp/optimal/)) in narrowPeak format. The counts of binding sites within these noncoding regions were fit to a Poisson model parameterized by a binding site frequency  $\lambda$  per base pair of DNA; for each ChIP-seq experiment, a  $\lambda$  estimate was obtained using maximum likelihood. The resulting histogram of  $\lambda$  estimates is well-described by a beta distribution, as shown in Figure 3, with parameters as given in Table 2. Thus, it is convenient to adopt a prior density.

$$p(\lambda) = \frac{1}{\mathcal{B}(\lambda_{\max}; \alpha, \nu)} \lambda^{\alpha-1} (1-\lambda)^{\nu-1}, \quad (28)$$

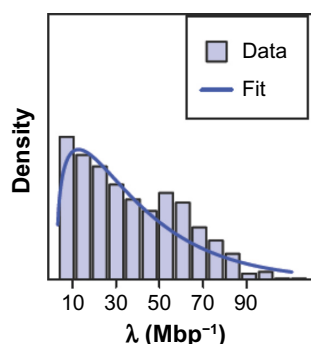
where  $\mathcal{B}$  is the incomplete beta function, and the shape hyperparameters are as given in Table 2 (recall that the range of  $\Lambda$  is  $(0, \lambda_{\max}]$ ). Combining Eqs. 27 and 28, and in the limit where  $\lambda \ll \sqrt{L_r / \tau w}$ , the product  $P(\beta | \lambda) p(\lambda)$  can be approximated

$$P(\beta | \lambda) p(\lambda) = \frac{1}{2^\beta \mathcal{B}(\lambda_{\max}; \alpha, \nu)} \lambda^{\beta+\alpha-1} (1-\lambda)^{\nu+L_r-(2w-1)\beta-1} \times (1 - O(\lambda^2)), \quad (29)$$

where with our choice of  $\lambda_{\max}$ , the second-order term can be neglected, resulting in a beta distribution-like dependence on  $\lambda$ ,

$$P(\beta | \lambda) p(\lambda) = p(\beta, \lambda) = \frac{1}{2^\beta \mathcal{B}(\lambda_{\max}; \alpha, \nu)} \lambda^{\beta+\alpha-1} \times (1-\lambda)^{\nu+L_r-(2w-1)\beta-1}, \quad (30)$$

for  $\lambda \in (0, \lambda_{\max}]$ . From this equation, and integrating  $\lambda$  over the range from  $[0, \lambda_{\max}]$ , we see that the probability model Eq. 26 corresponds to a prior



**Figure 3.** Distribution of frequencies of TFBS (per base pair of noncoding, gene-vicinal DNA sequence) for human transcription factors, based on the analysis of 620 ChIP-seq datasets from the ENCODE project.<sup>43</sup>

**Table 2.** Parameter estimate for the beta distribution model for the prior  $p(\lambda)$  on the binding site frequency per base pair, for human transcription factors.

	$\alpha$	$\nu (\times 10^4)$
Least-squares estimate	1.37	3.62
95% confidence interval	$\pm 0.15$	$\pm 0.53$

$$P(\beta) = \frac{\mathcal{B}(\lambda_{\max}; \beta + \alpha, \nu + L_r - (2w-1)\beta)}{2^\beta \mathcal{B}(\lambda_{\max}; \alpha, \nu)}. \quad (31)$$

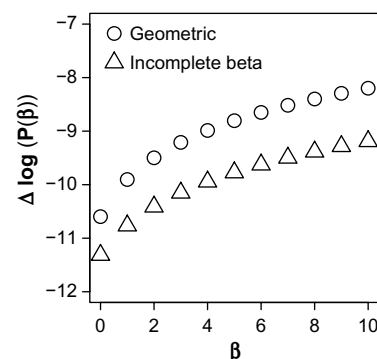
In contrast to Eq. 31, Lähdesmäki and Shmulevich<sup>50</sup> used a geometric prior on  $\beta$ ,  $P(\beta) = 1/2^{\beta+1}$ , implying

$$P(\beta) = q \frac{\beta! (L_r - \beta)!}{L_r!} \frac{1}{2^{\beta+1}}, \quad (32)$$

up to a normalization constant  $q$ . A comparison of the two priors (Fig. 4) suggests that the incomplete beta function prior (Eq. 31) may be more conservative than the geometric prior – a hypothesis that I investigate by analyzing simulated promoter sequence data in the next section.

### Obtaining the Distribution of $B|r, s', c$

The second step in a Bayesian approach<sup>44</sup> is to condition on the observed data – in this case,  $r, s'$ , and  $c$  – and then obtain the conditional distribution of the parameter(s) of interest, in this case  $B$ . In order to be able to do so, starting from the joint probability model (Eq. 12), the nuisance parameters  $\phi, \psi$ , and  $\lambda$  must be either estimated or marginalized. A key advantage of the Bayesian approach is that we can take into account the probability distribution of  $\Phi|c$ , in the process of eliminating  $\Phi$  by marginalization. The parameter  $\psi$  can be similarly marginalized, although the uncertainty in the background nucleotide



**Figure 4.** Plot of the change in  $\log P(\beta)$  with the addition of a single binding site, as a function of  $\beta$ , for the incomplete beta function-based prior (Eq. 31) and the previously proposed geometric prior (Eq. 32). The  $\Delta \log P$  values between the two priors are closer for  $\beta = 0$  but become greater with increasing  $\beta$ , indicating that the empirical prior in Eq. 31 is not simply equivalent to a rescaling of the geometric prior.



frequency is generally very small for real-world applications in which  $L'$  is large. We marginalize  $\phi$  and  $\psi$  by integration,

$$p(\mathbf{r}, \mathbf{s}', \lambda, \boldsymbol{\beta}, \mathbf{c}) = \int d^{3w} \phi \int d^3 \psi p(\mathbf{r}, \mathbf{s}', \lambda, \boldsymbol{\beta}, \phi, \psi, \mathbf{c}). \quad (33)$$

Given Eqs. 12, 23, and 25, the dependence of the integrand in Eq. 33 on  $\phi$  and  $\psi$  has the same algebraic form as the probability density function for independent Dirichlet random variables  $\{\Psi, \Phi_1, \dots, \Phi_w\}$ , as a consequence of the fact that the Dirichlet distribution is the conjugate prior for the multinomial distribution.<sup>44</sup> Thus, the two integrals in Eq. 33 can be evaluated analytically.<sup>30</sup> The desired joint conditional probability of  $\Lambda, \mathbf{B}$  follows by the definition of conditional probability,

$$p(\lambda, \boldsymbol{\beta} | \mathbf{r}, \mathbf{s}', \mathbf{c}) = \frac{p(\mathbf{r}, \mathbf{s}', \lambda, \boldsymbol{\beta}, \mathbf{c})}{P(\mathbf{r}, \mathbf{s}', \mathbf{c})}. \quad (34)$$

After performing the integrals in Eq. 33, using Eqs. 30 and 34, and then log-transforming, we have

$$\begin{aligned} & \log p(\lambda, \boldsymbol{\beta} | \mathbf{r}, \mathbf{s}', \mathbf{c}) \\ &= \log \mathcal{N}_2(\lambda_{\max}, L_r, w, \alpha, \nu, \mathbf{c}) - \log P(\mathbf{r}, \mathbf{s}', \mathbf{c}) - \beta \log 2 \\ & \quad + (\alpha + \beta - 1) \log \lambda + (\nu + L_r - (2w - 1)\beta - 1) \log(1 - \lambda) \\ & \quad + \sum_{g \in \mathbb{G}} \sum_{d \in \mathbb{D}} \log(c_{gd} + \sigma_{gd})! - \sum_{g \in \mathbb{G}} \log(c_g + \beta + 3)! \\ & \quad + \sum_{d \in \mathbb{D}} \log f_d! - \log(L_r + L' - \beta w + 3)!, \end{aligned} \quad (35)$$

where  $\mathcal{N}_2(\lambda_{\max}, L_r, w, \alpha, \nu, \mathbf{c})$  is function that will not need to be evaluated. The parameter  $\lambda$  can then be marginalized by integration, yielding

$$\begin{aligned} & \log P(\boldsymbol{\beta} | \mathbf{r}, \mathbf{s}', \mathbf{c}) \\ &= \log \mathcal{N}_2(\lambda_{\max}, L_r, w, \alpha, \nu, \mathbf{c}) - \log P(\mathbf{r}, \mathbf{s}', \mathbf{c}) - \beta \log 2 \\ & \quad + \log \mathcal{B}(\lambda_{\max}; \alpha + \beta, \nu + L_r - (2w - 1)\beta) \\ & \quad + \sum_{g \in \mathbb{G}} \sum_{d \in \mathbb{D}} \log(c_{gd} + \sigma_{gd})! - \sum_{g \in \mathbb{G}} \log(c_g + \beta + 3)! \\ & \quad + \sum_{d \in \mathbb{D}} \log f_d! - \log(L_r + L' - \beta w + 3)!, \end{aligned} \quad (36)$$

where  $\mathcal{B}$  is the incomplete beta function. As we will see below, in order to obtain samples of  $\mathbf{B} | \mathbf{r}, \mathbf{s}', \mathbf{c}$ , it will not be necessary to explicitly evaluate  $P(\mathbf{r}, \mathbf{s}', \mathbf{c})$ . Now that we have an explicit formula for  $\log[P(\boldsymbol{\beta} | \mathbf{r}, \mathbf{s}', \mathbf{c})]$  up to additive terms that do not depend on  $\boldsymbol{\beta}$ , it is possible to generate  $\boldsymbol{\beta}$  samples from this distribution using Markov Chain Monte Carlo (MCMC) sampling.

**MCMC approach.** For sampling from  $\mathbf{B} | \mathbf{r}, \mathbf{s}', \mathbf{c}$ , the Metropolis–Hastings algorithm,<sup>51</sup> in which a probabilistic proposal generator  $g(\boldsymbol{\beta}, \boldsymbol{\beta}')$  for a transition from  $\boldsymbol{\beta} \rightarrow \boldsymbol{\beta}'$  can be defined so as to optimize the acceptance rate for moves, is

convenient. For the problem of TFBS enrichment detection, following the general approach used by Lähdesmäki et al for TFBS recognition, I use a two-stage proposal generator in which a base pair position  $l \in \mathbb{L}_r$  is selected at random, and then, depending on the current state of  $\boldsymbol{\beta}$ , binding site removal or addition (in the latter case, with a randomly selected value  $j \in \{-1, 1\}$ ) is proposed (in the case of binding site addition,  $j = -1$  or  $j = 1$  is chosen with equal probability). For this approach, it will be useful to have a simplified expression for the log probability ratio for  $B_l = j$  versus  $B_l = 0$ , conditioned on  $\mathbf{r}, \mathbf{s}', \boldsymbol{\beta}_{\mathbb{L}_r \setminus \{l\}}, \mathbf{c}$ ; it is convenient to define some additional notation in order to make this conditional probability ratio explicit. Without loss of generality, let us assume that the current state for the *hypothetical* binding site configuration  $\boldsymbol{\beta} \in \mathbb{B}$ , a location  $l \in \mathbb{L}_r$  such that  $\beta_l = 0$ , and an orientation  $j \in \{-1, 1\}$  such that the configuration  $\boldsymbol{\beta}$  with  $\beta_l = j$  would not violate the TFBS physical constraints. In order to simplify notation, I define a function  $\mathcal{H}: \mathbb{D}^{L_r} \times \mathbb{L} \times \{-1, 1\} \times \mathbb{G} \rightarrow \mathbb{D}$  by

$$\mathcal{H}(\mathbf{r}, l, j, g) = \mathcal{C}(r_{l+j(g-1)}, j), \quad (37)$$

whose value represents the nucleotide at position  $g$  within the binding site for TFX that has orientation  $j$  and whose 5' most nucleotide is at location  $l \in \mathbb{L}_r$ . I also define a function  $\mathcal{I}: \mathbb{D}^{L_r} \times \mathbb{L}_r \times \{-1, 1\} \times \mathbb{D} \rightarrow \{0, 1, \dots, g\}$  by

$$\mathcal{I}(\mathbf{r}, l, j, d) = \left\{ \left\{ g \in \mathbb{G} \mid d = r_{l+j(g-1)} \right\} \right\}, \quad (38)$$

whose value is the count of nucleotides of base  $d$  within a binding site for TFX in orientation  $j$  whose 5' most nucleotide is at location  $l$ .

Applying Eq. 35 to two different binding site configurations that differ by one binding site being present/absent at a specific location  $l \in \mathbb{L}_r$ , and using the definitions of  $\mathcal{H}$  and  $\mathcal{I}$ , we obtain a closed-form expression for the log ratio of the conditional probability of there being a binding site at location  $l$  (in orientation  $j$ ), to the conditional probability that there is not a binding site at  $l$

$$\begin{aligned} & \log \left( \frac{P(B_l = j | \mathbf{r}, \mathbf{s}', \boldsymbol{\beta}_{\mathbb{L}_r \setminus \{l\}}, \mathbf{c})}{P(B_l = 0 | \mathbf{r}, \mathbf{s}', \boldsymbol{\beta}_{\mathbb{L}_r \setminus \{l\}}, \mathbf{c})} \right) \\ &= -\log 2 \\ & \quad + \log \mathcal{B}(\lambda_{\max}; \alpha + \beta + 1, \nu + L_r - (2w - 1)(\beta + 1)) \\ & \quad - \log \mathcal{B}(\lambda_{\max}; \alpha + \beta, \nu + L_r - (2w - 1)\beta) \\ & \quad + \sum_{g \in \mathbb{G}} \log(c_{g, \mathcal{H}(\mathbf{r}, l, j, g)} + \sigma_{g, \mathcal{H}(\mathbf{r}, l, j, g)} + 1) \\ & \quad - \sum_{g \in \mathbb{G}} \log(c_g + \beta + 4) \\ & \quad - \sum_{d \in \mathbb{D}} \log \{ (f_d)_{[\mathcal{I}(\mathbf{r}, l, j, d)]} \} \\ & \quad + \log \{ (L_r + L' - \beta w + 3)_{[w]} \}, \end{aligned} \quad (39)$$



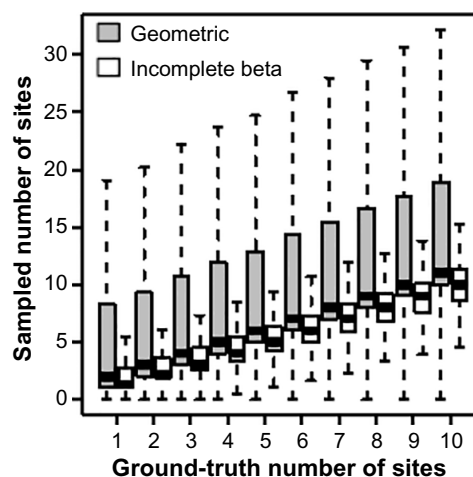
where the notation  $(x)_{[n]}$  denotes the falling factorial and  $\sigma$  is computed for  $\beta$ . Given Eq. 39, sampling from the distribution of  $B|r, s', c$  can be accomplished using the Metropolis–Hastings algorithm with the following proposal distribution:

$$g\left((\beta_1, \dots, \beta_{l-1}, j, \beta_{l+1}, \dots, \beta_{L_r}) \mid (\beta_1, \dots, \beta_{l-1}, 0, \beta_{l+1}, \dots, \beta_{L_r})\right) = \frac{(S_l)^q}{\sum_{r \in \mathbb{L}_r} (S_r)^q}, \quad (40)$$

where  $S_l$  is the Shannon entropy of the probabilities  $P(\beta_j = j | r, s', \beta_{\mathbb{L}_r \setminus \{l\}}, c)$  (for  $j \in \{-1, 0, 1\}$ ) with  $\beta_{\mathbb{L}_r \setminus \{l\}} = \mathbf{0}$ , and where the weight exponent  $q$  is tuned to achieve the desired average acceptance rate.<sup>52</sup> The reason for using a proposal distribution that weights each position by the Shannon entropy is that at most positions  $l \in \mathbb{L}_r$ , the entropy of the three conditional probabilities  $P(\beta_l = j | r, s', \beta_{\mathbb{L}_r \setminus \{l\}}, c)$  (for  $j \in \{-1, 0, 1\}$  and fixed  $l$ ) is very small, and thus, from the standpoint of optimizing the acceptance rate, it is convenient to weight the generation of proposed moves toward moves with more even odds of move acceptance. Results from empirical testing with sequence lengths  $L_r = 2 \times 10^4$  suggest that a value  $q = 0.85$  gives an acceptance rate of about 0.12, with increasing values of  $q$  increasing the acceptance rate. In this study, the Markov chain is initialized by iterating over  $l = 0$  to  $l = L_r$ , for each  $l$ , setting  $\beta_l$  to be the most probable configuration given Eq. 39, conditioned on  $\beta_{\mathbb{L}_r \setminus \{l\}} = \mathbf{0}$ . Once the Markov chain has converged, at most positions  $l \in \mathbb{L}_r$ , the entropy of the three conditional probabilities  $P(\beta_l = j | r, s', \beta_{\mathbb{L}_r \setminus \{l\}}, c)$  (for  $j \in \{-1, 0, 1\}$  and fixed  $l$ ) is very small, and thus, from the standpoint of optimizing the acceptance rate, it is convenient to weight the generation of proposed moves toward moves with more even odds of move acceptance.

## Empirical Results

Based on a direct comparison (Fig. 4) of the incomplete beta function prior (Eq. 31) and the previously proposed geometric prior (Eq. 32), it seems reasonable to suppose that, within the context of a Bayesian approach for PFM-based TFBS frequency estimation (as described in the “Obtaining the distribution of  $B|r, s', c$ ” section) the incomplete beta function prior and the geometric prior might have different effects on the conditional distribution of the number of TFBS, ie, the samples of  $B|c, r, s'$ . To test this hypothesis, I generated a synthetic dataset based on a simulated background sequence  $s'$  (with  $L' = 100,000$ ) and 120 gene promoter sequences (each with  $L_r = 20,000$ ), with uniform probabilities for each nucleotide. Into each simulated base sequence  $r$ , and for each of a fixed set of 100 TF PFMs selected at random from TRANSFAC Professional 2015,  $t \in \{1, \dots, 10\}$  TFBS were inserted into the  $r$  sequence (using representative binding site sequences from which the PFMs were computed, resulting in a modified sequence  $r_t$ ). Ten samples from the stationary distribution of  $B$  (the number of TFBS) were then generated using the MCMC



**Figure 5.** Comparing the accuracies of two MCMC implementations of the Bayesian method for estimating the number of binding sites of a TF, based on the geometric prior (Eq. 32) and the incomplete beta function-based prior (Eq. 31). For each combination of  $t$  (number of ground-truth sites) and type of prior, the bar denotes the median, the box denotes the interquartile range, and the whiskers are offset 1.5 interquartile range above or below the 75th and 25th percentiles.

approach described in the “MCMC approach” section (with 5000 burn-in steps, 100 steps per sample, and  $q = 0.85$ ), for both the geometric prior and the incomplete beta function-based prior (with  $\nu = 10,000$  and  $\alpha = 1.0$ ). For each of the two priors and for each combination of sequence  $r$ , PFM  $c$ , and number of ground-truth binding sites  $t$ , the 10  $B|c, r_t, s'$  samples were averaged, producing one *geometric prior* sample and one *incomplete beta function prior* sample for each of the 120,000 combinations of  $c, r$ , and  $t$ . The distributions of  $B|c, r_t, s'$ , organized by  $t$  and by prior, as shown in Figure 5, reveal several interesting patterns. First, across the fixed set of 100 randomly selected TFs, the MCMC method incorporating the incomplete beta function prior appears to yield samples that are more accurate than the MCMC method incorporating the geometric prior. In terms of mean-squared error, the MCMC method with the incomplete beta function-based prior is 19.6, whereas the mean-squared error with the geometric prior is 107.7. Second, the samples generated using the MCMC method with the geometric prior appear to be substantially higher-variance than the samples generated using the MCMC method with the incomplete beta function-based prior (quantitatively, the  $t$ -averaged standard deviation of the TFBS count samples obtained using the incomplete beta function prior was 4.05 versus 8.75 for the geometric prior).

## Discussion

This study demonstrates the utility of incorporating an empirical prior on the TFBS frequency per base pair within the context of a Bayesian method for PFM-based TFBS enrichment analysis, but there are several aspects in which the work raises interesting questions that could be explored in future studies. First, in this work, a two-parameter



parametric function has been fit to empirical data on the density distribution of frequencies of human TFBS per base pair of noncoding, gene-vicinal sequence. Thus, the results shown here do not reveal to what extent the estimated parameters for the distribution would generalize to TFBS frequencies in other species. At least for the mouse genome, available evidence from the modENCODE project suggests that overall, TF binding within promoter regions is highly conserved between human and mouse.<sup>53</sup> Moreover, for two TFs whose TFBS were assayed in five mammalian species by ChIP-seq, the numbers of genome-wide binding sites did not vary more than 2× between species.<sup>22</sup> Thus, it seems reasonable to expect that the  $\Lambda$  prior distribution (across TFs) would be similar, for gene-vicinal noncoding sequence. Nevertheless, in future work, it would be informative to estimate the hyperpriors  $\alpha$  and  $\nu$  for human, mouse, fruit fly, and worm to enable a cross-species comparison. Second, it would be useful to characterize how the choice of  $q$  parameter affects the empirical performance of the MCMC approach used here, ie, the acceptance ratio, the number of steps required for burn-in, and the number of steps required between samples; it may be possible to significantly improve the speed of the proposed MCMC method through tuning  $q$  and the sampling parameters. Third, a key aspect to be explored is the extent to which the accuracy improvement with the incomplete beta function-dependent prior is associated with high-count versus low-count PFMs. Intuitively, it seems reasonable to suppose that for most TFs, an increase in the accuracy of the prior would be expected to have more of an effect on the posterior distribution of  $\beta$  when the PFM count is low, since a higher count PFM would be expected to have a much bigger likelihood ratio that would, in turn, be more likely to dominate over the prior on the number of TFBS.

## Conclusions

This study presents a Bayesian approach to the bioinformatics problem of PFM-guided TFBS enrichment analysis. The method incorporates an empirical prior on the frequency distribution  $\lambda$  of binding sites for TFs that is based on genome location data from the ENCODE project. In addition, the method incorporates a probabilistic model for TFBS occurrence conditioned on the parameter  $\lambda$  that takes into account the finite width of the TFBS, in contrast to a previous approach in which the TFBS probability was assumed to have a geometric dependence with a fixed factor of 1/2.<sup>30</sup> The sampling equation for adding/removing a binding site (Eq. 39) could be easily extended to include other sources of information, such as a regulatory potential score derived from phylogenetic sequence conservation or from epigenetic measurements. The R software code implementing the MCMC method described in the “MCMC approach” section and the promoter analyses shown in Figure 5 is available at <http://github.com/ramseylab/tfbsincbeta>.

## Acknowledgments

The author thanks Jichen Yang, Tanjin Xu, Holly Arnold, and Theo Knijnenburg, for reviewing early drafts of the article and for providing helpful feedback. The author also thanks Harri Lähdesmäki for providing technical insights on the Lähdesmäki–Rust–Shmulevich method for TFBS recognition and Yuan Jiang for technical advice. Part of this work was carried out in the laboratories of Alan Aderem and Ilya Shmulevich, and their support is gratefully acknowledged.

## Author Contributions

Conceived and designed the experiments: SAR. Analyzed the data: SAR. Wrote the first draft of the manuscript: SAR. Contributed to the writing of the manuscript: SAR. Agree with manuscript results and conclusions: SAR. Jointly developed the structure and arguments for the paper: SAR. Made critical revisions and approved final version: SAR. The author reviewed and approved of the final manuscript.

## REFERENCES

- Stormo GD. Computer methods for analyzing sequence recognition of nucleic acids. *Annu Rev Biophys Biophys Chem.* 1988;17:241–63.
- Roth FP, Hughes JD, Estep PW, Church GM. Finding DNA regulatory motifs within unaligned noncoding sequences clustered by whole-genome mRNA quantitation. *Nat Biotechnol.* 1998;16:939–45.
- Wasserman WW, Sandelin A. Applied bioinformatics for the identification of regulatory elements. *Nat Rev Genet.* 2004;5:276–87.
- Rahmann S, Müller T, Vingron M. On the power of profiles for transcription factor binding site detection. *Stat Appl Genet Mol Biol.* 2003;2:Article7.
- KelAE, Gössling E, Reuter I, Cheremushkin E, Kel-Margoulis OV, Wingender E. MATCHTM: a tool for searching transcription factor binding sites in DNA sequences. *Nucleic Acids Res.* 2003;31:3576–9.
- Aerts S, Thijs G, Coessens B, Staes M, Moreau Y, De Moor B. Toucan: deciphering the cis-regulatory logic of coregulated genes. *Nucleic Acids Res.* 2003;31:1753–64.
- Frith MC, Fu Y, Yu L, Chen JF, Hansen U, Weng Z. Detection of functional DNA motifs via statistical over-representation. *Nucleic Acids Res.* 2004;32:1372–81.
- Tomba M, Li N, Bailey TL, et al. Assessing computational tools for the discovery of transcription factor binding sites. *Nat Biotechnol.* 2005;23:137–44.
- Gilchrist M, Thorsson V, Li B, et al. Systems biology approaches identify ATF3 as a negative regulator of Toll-like receptor 4. *Nature.* 2006;441:173–8.
- Tan K, Tegner J, Ravasi T. Integrated approaches to uncovering transcription regulatory networks in mammalian cells. *Genomics.* 2008;91:219–31.
- Ramsey SA, Klemm SL, Zak DE, et al. Uncovering a macrophage transcriptional program by integrating evidence from motif scanning and expression dynamics PLOS. *Comput Biol.* 2008;4:e1000021.
- Litvak V, Ratushny AV, Lampano AE, et al. A FOXO3-IRF7 gene regulatory circuit limits inflammatory sequelae of antiviral responses. *Nature.* 2012;490:421–5.
- Gold ES, Ramsey SA, Sartain MJ, et al. ATF3 protects against atherosclerosis by suppressing 25-hydroxycholesterol-induced lipid body formation. *J Exp Med.* 2012;209:807–17.
- Ramsey SA, Vengrenyuk Y, Menon P, et al. Epigenome-guided analysis of the transcriptome of plaque macrophages during atherosclerosis reveals activation of the Wnt signaling pathway. *PLoS Genet.* 2014;10:e1004828.
- Quandt K, Frech K, Karas H, Wingender E, Werner T. MatInd and MatInspector: new fast and versatile tools for detection of consensus matches in nucleotide sequence data. *Nucleic Acids Res.* 1995;23:4878–84.
- Wingender E, Chen X, Hehl R, et al. TRANSFAC: an integrated system for gene expression regulation. *Nucleic Acids Res.* 2000;28:316–9.
- Sandelin A, Alkema W, Engstrom P, Wasserman WW, Lenhard B. JASPAR: an open-access database for eukaryotic transcription factor binding profiles. *Nucleic Acids Res.* 2004;32:D91–4.
- Newburger DE, Bulyk ML. UniPROBE: an online database of protein binding microarray data on protein-DNA interactions. *Nucleic Acids Res.* 2009;37:D77–82.



19. Wang J, Zhuang J, Iyer S, et al. Sequence features and chromatin structure around the genomic regions bound by 119 human transcription factors. *Genome Res.* 2012;22:1798–812.
20. Sebastian A, Contreras-Moreira B. footprintDB: a database of transcription factors with annotated cis elements and binding interfaces. *Bioinformatics.* 2014;30:258–65.
21. Johnson DS, Mortazavi A, Myers RM, Wold B. Genome-wide mapping of in vivo protein-DNA interactions. *Science.* 2007;316:1497–502.
22. Schmidt D, Wilson MD, Ballester B, et al. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Science.* 2010;328:1036–40.
23. Villar D, Flicek P, Odom DT. Evolution of transcription factor binding in metazoans—mechanisms and functional implications. *Nat Rev Genet.* 2014;15:221–33.
24. Vaquerizas JM, Kummerfeld SK, Teichmann SA, Luscombe NM. A census of human transcription factors: function, expression and evolution. *Nat Rev Genet.* 2009;10:252–63.
25. Ho Sui SJ, Mortimer JR, Arenillas DJ, et al. oPOSSUM: identification of over-represented transcription factor binding sites in co-expressed genes. *Nucleic Acids Res.* 2005;33:3154–64.
26. Sinha S. On counting position weight matrix matches in a sequence, with application to discriminative motif finding. *Bioinformatics.* 2006;22:e454–63.
27. Pavesi G, Zambelli F. Prediction of over represented transcription factor binding sites in co-regulated genes using whole genome matching statistics. In: Masulli F, Mitra S, Pasi G, eds. *Applications of Fuzzy Sets Theory.* Berlin: Springer; 2007:651–8.
28. Cheng C, Alexander R, Min R, et al. Understanding transcriptional regulation by integrative analysis of transcription factor binding data. *Genome Res.* 2012;22:1658–67.
29. Sinha S, Tompa M. YMF: a program for discovery of novel transcription factor binding sites by statistical overrepresentation. *Nucleic Acids Res.* 2003;31:3586–8.
30. Lähdesmäki H, Rust AG, Shmulevich I. Probabilistic inference of transcription factor binding from multiple data sources. *PLoS One.* 2008;3:e1820.
31. Miller AK, Print CG, Nielsen PMF, Crampin EJ. A Bayesian search for transcriptional motifs. *PLoS One.* 2010;5:e13897.
32. Berg OG. Selection of DNA binding sites by regulatory proteins: the LexA protein and the arginine repressor use different strategies for functional specificity. *Nucleic Acids Res.* 1988;16:5089–105.
33. Nishida K, Frith MC, Nakai K. Pseudocounts for transcription factor binding sites. *Nucleic Acids Res.* 2009;37:939–44.
34. Liu JS. The collapsed Gibbs sampler in Bayesian computations with applications to a gene regulation problem. *J Am Stat Assoc.* 1994;89:958–66.
35. Thijs G, Marchal K, Lescot M, et al. A Gibbs sampling method to detect over-represented motifs in the upstream regions of coexpressed genes. *J Comput Biol.* 2002;9:447–64.
36. Xing EP, Wu W, Jordan MI, Karp RM. LOGOS: a modular Bayesian model for de novo motif detection. *Proc IEEE Comput Soc Bioinform Conf.* 2003;2:266–76.
37. Jensen ST, Liu XS, Zhou Q, Liu JS. Computational discovery of gene regulatory binding motifs: a Bayesian perspective. *Stat Sci.* 2004;19:188–204.
38. Siepel A, Bejerano G, Pedersen JS, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.* 2005;15:1034–50.
39. Taylor J, Tyekucheva S, King DC, Hardison RC, Miller W, Chiaromonte F. ESPERR: learning strong and weak signals in genomic sequence alignments to identify functional elements. *Genome Res.* 2006;16:1596–604.
40. Berger MF, Badis G, Gehrke AR, et al. Variation in homeodomain DNA binding revealed by high-resolution analysis of sequence preferences. *Cell.* 2008;133:1266–76.
41. Gabdoulline R, Eckweiler D, Kel A, Stegmaier P. 3DTF: a web server for predicting transcription factor PWMs using 3D structure-based energy calculations. *Nucleic Acids Res.* 2012;40:W180–5.
42. Muratani M, Deng N, Ooi WF, et al. Nanoscale chromatin profiling of gastric adenocarcinoma reveals cancer-associated cryptic promoters and somatically acquired regulatory elements. *Nat Comm.* 2014;5:4361.
43. Gerstein MB, Kundaje A, Hariharan M, et al. Architecture of the human regulatory network derived from ENCODE data. *Nature.* 2012;489:91–100.
44. Gelman A, Carlin JB, Stern HS, Dunson DB, Vehtari A, Rubin DB. *Bayesian Data Analysis.* 3rd ed. CRC Press, Boca Raton, FL; 2013.
45. Buntine WL. Operations for learning with graphical models. *J Artif Intel Res.* 1994;2:159–225.
46. Heumann JM, Lapedes AS, Stormo GD. Neural networks for determining protein specificity and multiple alignment of binding sites. *Proc Int Conf Intel Syst Mol Biol.* 1994;2:188–94.
47. Hertz GZ, Hartzell GW, Stormo GD. Identification of consensus patterns in unaligned DNA sequences known to be functionally related. *Comput App Biosci CABIOS.* 1990;6:81–92.
48. Jensen ST, Liu JS. Bayesian clustering of transcription factor binding motifs. *J Am Stat Assoc.* 2008;103:188–200.
49. Kharchenko P, Tolstorukov M, Park P. Design and analysis of ChIP-seq experiments for DNA-binding proteins. *Nat Biotechnol.* 2008;26(12):1351–9.
50. Lähdesmäki H, Shmulevich I. Learning the structure of dynamic Bayesian networks from time series and steady state measurements. *Mach Learn.* 2008;71:185–217.
51. Metropolis N, Rosenbluth AW, Rosenbluth MN, Teller AH, Teller E. Equation of state calculations by fast computing machines. *J Chem Phys.* 1953;21:1087–92.
52. Roberts GO, Gelman A, Gilks WR. Weak convergence and optimal scaling of random walk Metropolis algorithms. *Ann Appl Probab.* 1997;7:110–20.
53. Shen Y, Yue F, McCleary DF, et al. A map of the cis-regulatory sequences in the mouse genome. *Nature.* 2012;488:116–20.