

# A Phylogenetic Gibbs Recursive Sampler for Locating Transcription Factor Binding Sites

Sean P. Conlan<sup>1</sup>   Lee Ann McCue<sup>2</sup>   Lee A. Newberg<sup>1,3</sup>  
Thomas M. Smith<sup>3</sup>   William Thompson<sup>4</sup>   Charles E. Lawrence<sup>4</sup>

<sup>1</sup>Wadsworth Center, New York State Department of Health

<sup>2</sup>Pacific Northwest National Laboratory

<sup>3</sup>Department of Computer Science, Rensselaer Polytechnic Institute

<sup>4</sup>Department of Applied Mathematics, Brown University

Systems Biology: Global Regulation of Gene Expression  
2006

## Take-Home Points

- Phylogenetic modeling (Felsenstein's Algorithm) helps
- Use of ensemble centroid helps

## What We're Looking For

- Seeking elements that are short: 6–30 bp
- Only partial conserved
- Isolated elements or multiple elements per module
- Single or multiple intergenic regions per genome
- **Alignable and unalignable sequence data across genomes**

## Measures of Success

- **Sensitivity** — minimize false negatives
- **Selectivity** — minimize false positives

## Previous Work

### Non-Phylogenetic Algorithms

Many good algorithms including

- Gibbs Recursive Sampler (Thompson *et al.*, 2003)

But need to be better when analyzing closely related species.

### Phylogenetic Algorithms

Several good algorithms

- Non-statistical and/or two-species only
- PhyloGibbs (Siddharthan *et al.*, 2005). Uses successive star-topology approximations, maximum likelihood

But improvement is possible with full **Felsenstein's Algorithm** and with an **ensemble centroid**

# Gibbs Sampling

## Gibbs Sampling Overview

Move from proposed solution to proposed solution via Gibbs Sampling.

- From any proposed set of sites
  - Re-choose sites in one multi-sequence<sup>a</sup>, with probability conditioned on sites in remaining multi-sequences
- Iterate to explore parameter space.
  - Explores each proposed set of sites with probability proportional to its likelihood.

---

<sup>a</sup>An unalignable sequence or a set of aligned sequences

# Probability Conditioned on Remaining Sites

A slight oversimplification . . .

## Probability Calculation

- Current Iteration has a *position-weight matrix*, which gives current motif description, and is built from counts from current sites & pseudocounts.
- A position's weights parameterize a Dirichlet distribution, which is used to draw an equilibrium distribution.
- The equilibrium is used to parameterize a nucleotide substitution model (e.g., HKY85, HB98, New05).
- The substitution model is used to evaluate all positions attributed to it, via Felsenstein's Algorithm.

# Exact Probabilities via Felsenstein's Algorithm

A linear-time, phylogenetic-tree traversal algorithm.

## Inputs

- The nucleotide from each species in a multiple-alignment sequence position
- A phylogenetic tree with branch lengths
- A nucleotide substitution model

## Output

- The probability of the observed data for that multiple-alignment sequence position

See Felsenstein (1981), or many good textbooks, for more information.

# Ensemble Centroid

## Computing the Ensemble Centroid

With each sample from the Gibbs Sampler (after “burn in” iterations)

- For each sequence position record a “1” if it is part of a *cis*-element, record “0” otherwise.
- The vector of 0’s and 1’s is the corner of a hypercube

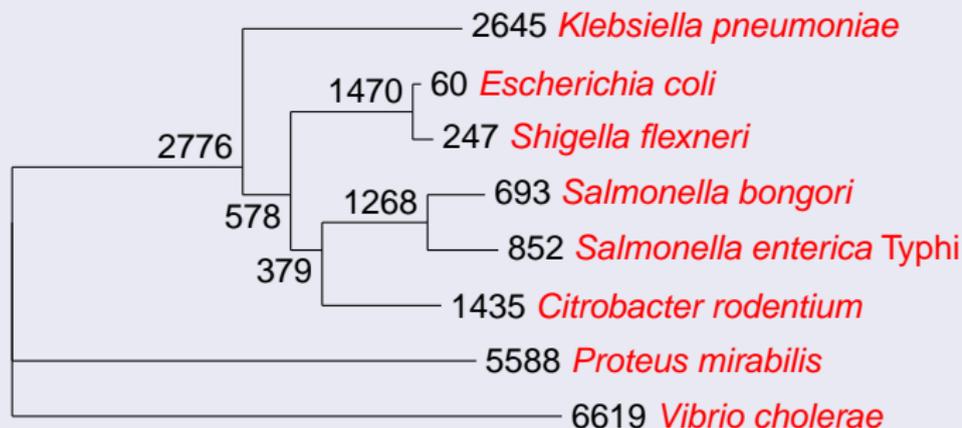
Ensemble centroid = **corner nearest to the center of mass of the collected samples**

## Advantages of Ensemble Centroid

- Expensive *a posteriori* probability calculation not needed
  - Star-topology approximation unnecessary
- Gives “entropic” solutions their due

# Synthetic Test Data

A data set: Eight species' 500 bp sequences



Expected number of substitutions  $\times 10^4$

- Gapless sequence data generated according to tree
- *P. mirabilis* and *V. cholerae* subsequently treated as not alignable

# Synthetic Test Data

## Five Collections of Data Sets

- Five collections of data sets:  $k \in \{0, 1, 2, 3, 4\}$
- 100 data sets in each collection
- A data set is 8 sequences
  - one for each species
  - each of length 500 bp
  - each with  $k$  planted *Escherichia coli* Crp binding sites
  - related by phylogenetic tree

## Data Analysis

- Each data set run separately — 500 runs total
- Accumulate results across data sets in each collection.

# Sensitivity & Selectivity

*E.g.*, entry in red shows 116 *E. coli* sites found across 61 data sets, where PhyloGibbs finds 13 *E. coli* sites across 8 data sets.

## Our Algorithm (top) and PhyloGibbs (bottom)

Data Collection	#0	#1	#2	#3	#4
Sites Found (True Positives)		17/17 0/0	116/61 13/8	154/82 54/26	176/93 75/35
False Sites (False Positives)	3/3 47/46	5/4 60/51	2/2 63/44	0/0 40/30	0/0 30/24
Sites Missed (False Negatives)		83/83 100/100	84/45 187/95	146/100 246/89	224/100 325/97

“BRASS” implementation of our algorithm (Smith, 2006), configured to find **up to two sites** per multi-sequence

# A New Nucleotide Substitution Model for Use with Felsenstein's Algorithm

Lee A. Newberg<sup>1,2</sup>

<sup>1</sup>Wadsworth Center, New York State Department of Health

<sup>2</sup>Department of Computer Science, Rensselaer Polytechnic Institute

Systems Biology: Global Regulation of Gene Expression  
2006

# Features of a New Nucleotide Substitution Model

## Features

- we explain away the apparent, nonsensical simultaneity of mutations in non-adjacent sequence positions
- we explain the failure of pooled data to behave according to some averaged model
- we permit polymorphisms, yielding a higher expected number of mismatches in conserved sequence positions

# Modeling How Nucleotides Evolve

## Existing Models

- Arbitrary equilibria
- Transition/transversion rate ratio
- Mutation rate variation within a genome
- Selection effects **via scaled fixation rates** (Halpern & Bruno, 1998)
- Context sensitive: Di- and tri-nucleotide models
- Indel support, though difficult with Felsenstein's Algorithm

## A New Model for Selection Effects

Newberg (2005) **allows that SNPs are not improbable.** (*i.e.*, without the specious fixation on species fixation.)

# Traditional Nucleotide Substitution Model

## Traditional Mutation (without Selection)

For example,

$$M_x = \begin{pmatrix} \Pr[A|A] & \Pr[C|A] & \Pr[G|A] & \Pr[T|A] \\ \Pr[A|C] & \Pr[C|C] & \Pr[G|C] & \Pr[T|C] \\ \Pr[A|G] & \Pr[C|G] & \Pr[G|G] & \Pr[T|G] \\ \Pr[A|T] & \Pr[C|T] & \Pr[G|T] & \Pr[T|T] \end{pmatrix}$$
$$= \begin{pmatrix} 0.96 & 0.01 & 0.02 & 0.01 \\ 0.01 & 0.96 & 0.01 & 0.02 \\ 0.02 & 0.01 & 0.96 & 0.01 \\ 0.01 & 0.02 & 0.01 & 0.96 \end{pmatrix}$$

Each row sums to 1.0.

# A New Nucleotide Substitution Model

## Population Model for Selection (without Mutation)

For example,

$$M_x = \begin{pmatrix} 1.1 & 0 & 0 & 0 \\ 0 & 1.0 & 0 & 0 \\ 0 & 0 & 1.0 & 0 \\ 0 & 0 & 0 & 1.0 \end{pmatrix}$$

Each row no longer sums to 1.0 but, starting with 100 organisms of each type ...

$$\begin{aligned} (100, 100, 100, 100)M_x &= (110, 100, 100, 100) \\ \frac{(110, 100, 100, 100)}{410} &= (0.268, 0.244, 0.244, 0.244) \end{aligned}$$

# A New Nucleotide Substitution Model

Combining the two ...

## Mutation *and* Selection

$$\begin{aligned} M_x &= \begin{pmatrix} \Pr[A|A] & \Pr[C|A] & \Pr[G|A] & \Pr[T|A] \\ \Pr[A|C] & \Pr[C|C] & \Pr[G|C] & \Pr[T|C] \\ \Pr[A|G] & \Pr[C|G] & \Pr[G|G] & \Pr[T|G] \\ \Pr[A|T] & \Pr[C|T] & \Pr[G|T] & \Pr[T|T] \end{pmatrix} \\ &= \begin{pmatrix} 1.056 & 0.01 & 0.02 & 0.01 \\ 0.011 & 0.96 & 0.01 & 0.02 \\ 0.022 & 0.01 & 0.96 & 0.01 \\ 0.011 & 0.02 & 0.01 & 0.96 \end{pmatrix} \end{aligned}$$

# A New Nucleotide Substitution Model

## Some Details

- Parameterized by background model and desired equilibrium
- Each generation is mutation (according to background model) followed by selection.
- Instantaneous rate formalism  $M_x = \exp(xR)$  still applies, so generation length need not be known.
- 2x invocations of Felsenstein's Algorithm, because each row no longer sums to 1.0.
- Easily computed correspondence between nucleotide equilibria  $\vec{\theta}$  and diagonal selection matrix

## Contact Information

[lnewberg@wadsworth.org](mailto:lnewberg@wadsworth.org)

[www.rpi.edu/~newbel/](http://www.rpi.edu/~newbel/)

→ Felsenstein, J. (1981) PubMed 7288891.

→ Halpern, A. L. & Bruno, W. J. (1998) PubMed 9656490.

→

Hasegawa, M., Kishino, H. & Yano, T. (1985) PubMed 3934395.

→ Newberg, L. A. (2005).

<http://www.cs.rpi.edu/research/pdf/05-08.pdf>.

→

Siddharthan, R., Siggia, E. D. & van Nimwegen, E. (2005) PubMed

→ Smith, T. M. (2006). PhD thesis, Rensselaer Polytechnic  
Institute Troy, NY. In preparation.

→

Thompson, W., Rouchka, E. C. & Lawrence, C. E. (2003) PubMed