

Identifying and representing DNA binding sites

Biol4230

Tues, April 17, 2018

Bill Pearson wrp@virginia.edu

4-2818 Pinn 6-057

- Looking for functional sites: promoters, regulatory elements, modification sites
- Products of convergent, not divergent evolution
- Weak spacing constraints
- Usually represented as a consensus sequence
- If alignment is given, consensus is obvious
- If consensus is given, alignment is obvious
- Search for consensus and alignment together
- consensus, meme, gibbs

fasta.bioch.virginia.edu/biol4230

1

To learn more:

- Overview of multiple alignment and motif finding – Mount (2001) Chapter 4
- Schneider et al. (1986) J. Mol. Biol. Information content of binding sites on nucleotide sequences 188:415-431
- Stormo and Hartzell (1989) Identifying protein binding sites from unaligned DNA fragments
- Lawrence and Reilly (1990) An expectation maximization (EM) algorithm for the identification and characterization of common sites in unaligned biopolymer sequences PROTEINS 7:41-51
- Lawrence et al. (1993) Detecting subtle sequence signals: a Gibbs sampling strategy for multiple alignment Science 262:208-214

fasta.bioch.virginia.edu/biol4230

2

Regulation of transcription: RNA polymerase, promoters, enhancers, transcription factors, DNA binding proteins

- Gene expression is regulated at many levels:
 - production of RNA (transcription)
 - promoter occupancy, transcription rate, termination
 - transcript RNA splicing
 - mRNA stability
 - mRNA translation
 - post-translational processing/stability

fasta.bioch.virginia.edu/biol4230

3

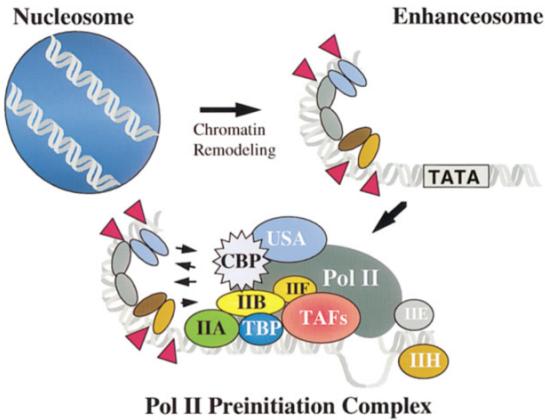
Regulation of transcription: RNA polymerase, promoters, enhancers, transcription factors, DNA binding proteins

- Transcription factors interact with DNA to increase/decrease transcription levels
 - lacR (bacterial lac Repressor)
 - hundreds of transcription factors modify expression in response to signals
 - FOS/JUN
 - nuclear receptors
 - homeobox proteins
 - Myc
 - etc.,etc.

fasta.bioch.virginia.edu/biol4230

4

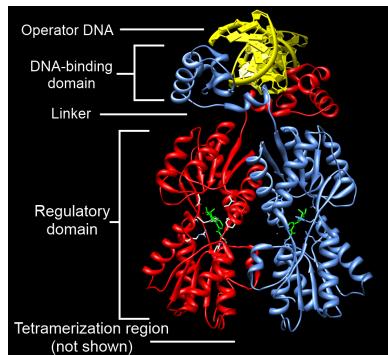
Regulation of transcription: RNA polymerase, promoters, enhancers, transcription factors, DNA binding proteins



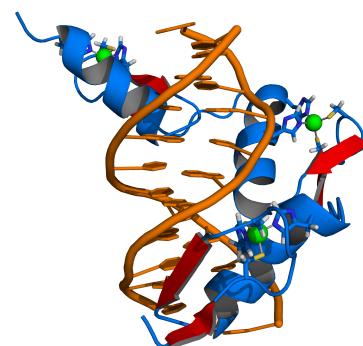
fasta.bioch.virginia.edu/biol4230

5

Regulation of transcription: RNA polymerase, promoters, enhancers, transcription factors, DNA binding proteins



Lac repressor

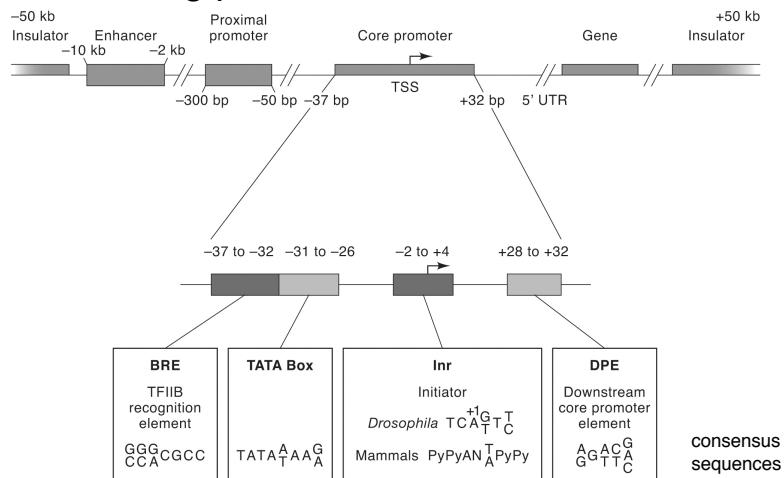


Zn finger

fasta.bioch.virginia.edu/biol4230

6

Regulation of transcription: RNA polymerase, promoters, enhancers, transcription factors, DNA binding proteins



Transcriptional regulation in eukaryotes, concepts, strategies, and techniques. CSHL Press 2009

7

The "first" promoter motifs: *E. coli* -10, -35

Harley, C. B. & Reynolds, R.
P. Nucleic Acids Res 15, 2343–2361 (1987).

Oliphant, A. R. & Struhl, K.
Nucleic Acids Res 16, 7673–7683 (1988).

Alignment of <i>E. coli</i> Promoter Sequences												
SEQUENCE	TYPE	-35		-10		SP	PHI	DISREP.	TS	REF		
(a)	(b)	(c)	(d)	(e)	(f)	(g)	(h)	(i)	(j)	(k)	(l)	(m)
aceEF	b	ACCTAGCCTTGTATTGACCTTC		CCCGAGAC TTCAAG CGGGAGCTTACAG		17	-4.3	-4.4	4	24		
ads	b	AGATGTTTGTTT TTTCTT GATGCTCA		CCCCGAGC TTCAAG CGGGAGCTTACAG		17	-5.5	-3.4	-4.6	4	25,26	
alaS	b	AACCCATACGCGT TTACG TTCCAGTC		AGGAAACT TAATCT ATTCCGCTTTTGAATG		18	-3.1			9		
ampC	b	TCCGTTGCTGAGG TTCTCA CGCTTAT		GCGTGTCT TACATG CTAGGGCTGGCGGAG		16	-1.5			9		
ampC/C16	b	TCGTTGCTGAGG TTCTCA CGCTTAT		GCGTGTCT TACATG CTAGGGCTGGCGGAG		16	-1.5			9		
arcA	b	TTCACCGGTTAGG GTCAGT TTTCGCG		CGGGGGCT TGTGGT GTTGTGGTGTGGGGCT		16	-1.6			1,3	25	
arcC	b	CGGAAATGATGCG AGTACT TTTCGCG		GTGATGAGA GAGCTT TGTGGTGTGGGGCT		17	-3.6			9		
arcE	b	CGGTTTCCAG CGAAC CGGGGAGA		GTCGGTCAAG TATTTC TTACAGCTTCTACAG		17	-3.2			4	28	
arcA(I)	b	ACGGGGATGAC CGGGG CCTTTAT		CGGAGCTC TTACATG TTCTGAGCGGGGCT		16	-4.3			4	29	
arcA(I)(X)c	b	ACGGGGATGAC CGGGG CCTTTAT		CGGAGCTC TTACATG TTCTGAGCGGGGCT		18	-3.8			4	29	
arcB	b	TTCGTTTTCGATG TTGAG GAGCTTG		TCATGAGA TATCAA TATTCGCGGCTTAT		18	-2.4			9		
arcB(P-1/-6)	b	TTCGTTTTCGATG TTGAG GAGCTTG		GTGCTATA TATCAT CAACTTCTTCGCTTG		15	-2.0			30		
arcB(P-1/L)	b	TTCGTTTTCGATG TTGAG GAGCTTG		GTGCTATA TATCAT CAACTTCTTCGCTTG		15	-2.0			30		
argP-1	b	TTACGGGCGTGGG TTTCAT TGGCGCA		ACCTGGGG TATTTC TTATTCGGAACTACGAA		17	-2.6			4	31	
argP-2	b	CGGGGGGGGGGGG TGGGGT GGAGGT		GGGGGGGG TGGGGT GGAGGT GGGGGGGGGGG		17	-3.1			3.9	4	31
argP-1/113	b	CGGGGGGGGGGGG TGGGGT GGAGGT		GGGGGGGG TGGGGT GGAGGT GGGGGGGGGGG		17	-3.3			3.1		
argP	b	ATTCGAAATGGGG TTCCCA ATGAAATGAA		TTGGCGAGA TTAAGT GAATTTTAATTCGAA		17	-1.7			4	31,32	
argI	b	AGC TTCCCA ATGAAATGAA		TTGGCGAGA TTAAGT GAATTTTAATTCGAA		17	-1.5			4	31	
argR	b	TGGGGGGGGGGG TTGGGG GAGGAGG		CTTCATCAA TAAATG GAATTTTAATTCGAA		17	-3.2			-5.9	2,4	31
arcF	b	TACGAAAATGAGG TTGAA ACTTGCT		TTGGGGGG TGGGGT GGAGGT GGGGGGGGG		16	-1.9			2,4	33	
arcG	b	AGCTTAAACGGGG TTGAA CATTGCTCA		CGGAGAGA TAGATT CGGAGCTTTCACATCA		17	-1.6			2,4	33	
arcH	b	GTACTGAGGAGCTT GTCAT TACGTAT		TTTTTGGG TACATG CCTGGGGGGGGGGGG		16	-3.1			9		
bioA	b	GGCTGTGCGAAC GGGTTT TTTCGTT		AACTGGGG TAGACT TGCTGGGGGGGGGG		18	-3.8	-3.4				
bioB	b	TGCGACGATGAGC TTGAA ACCAAAT		GAAGAGAT TAGGT TACAGGGTTCACCGGA		17	-2.2					

A.														
G	1	4	42	2	2	8	...	3	1	11	7	6	2	G
C	0	4	3	15	32	11	...	15	47	13	21	12	2	C
A	0	3	1	33	11	16	...	1	47	13	21	12	2	A
T	57	47	12	8	13	23	...	33	2	14	16	10	47	T

B.												
15	16	17	18	19	20	21	22	23	24	25	26	27

Figure 5: Matrix for the consensus of *E. coli* promoter elements. The number of times each base occurs in each position of the -35 and -10 sequence elements selected from random DNA (A), and the number of those elements separated by a spacing of from 15 to 19 base pairs (B).

fasta.bioch.virginia.edu/biol4230

8

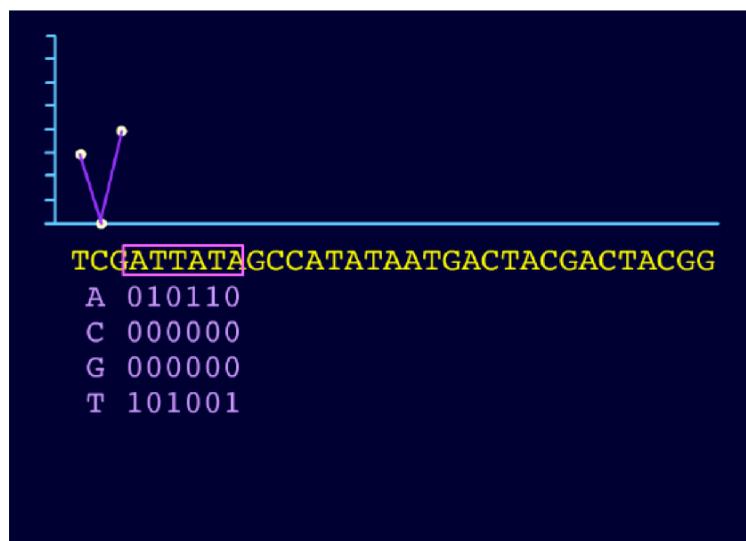
Consensus Patterns and Motifs

- How to represent motifs
 - consensus patterns
 - weight matrices – Position Specific Scoring Matrices
- How to “weight” positions?
 - frequency
 - information content ($\Sigma \log \text{odds}$)
- How to search for motifs
 - Heuristic greedy (consensus, wconsensus)
 - Expectation-Maximization (MEME)
 - Gibbs sampling

fasta.bioch.virginia.edu/biol4230

9

Scanning a sequence with PATSER

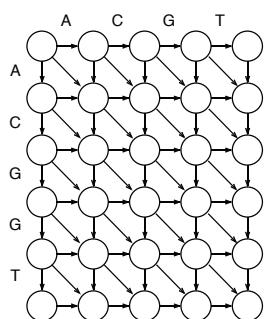


Scanning a sequence with PATSER (PWM, PSSM)

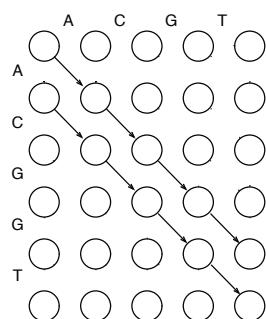


Pattern scanning vs Alignment Computational complexity

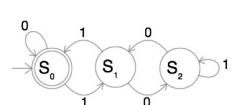
Pairwise alignment
PSSM with gaps
 $O(n^2)$



Global alignment
PSSM/PWM no gaps
 $O(n^2)$



Exact match
No gaps
 $O(n)$



[wikipedia.org/wiki/
Deterministic_finite_automaton](https://en.wikipedia.org/wiki/Deterministic_finite_automaton)

Consensus Patterns and Motifs

- How to represent motifs

- consensus patterns
- weight matrices

- do not require identity
- include consensus (1's, 0's)
- allow mismatches

Position Specific
Scoring Matrix

$$S = \log \left(\frac{f_{pos,b}}{p_b} \right)$$

fasta.bioch.virginia.edu/biol4230

13

How to represent motifs – information content

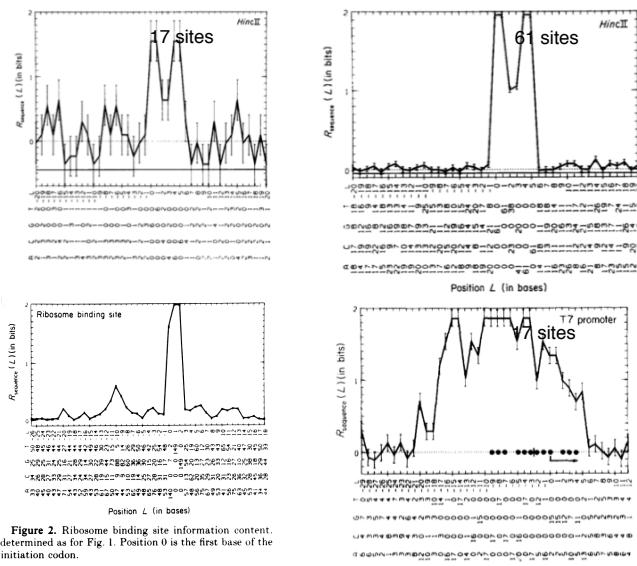


Figure 2. Ribosome binding site information content, determined as for Fig. 1. Position 0 is the first base of the initiation codon.

Schneider TD, Stormo GD, Gold L, Ehrenfeucht A. *J Mol Biol.* (1986) **188**:415-431
Information content of binding sites on nucleotide sequences.

14

Representing Consensus Sequences ("TATAA" box)

present/absent	<table border="1"> <tr><td>A</td><td>0</td><td>1</td><td>0</td><td>1</td><td>1</td><td>0</td></tr> <tr><td>C</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>G</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td><td>0</td></tr> <tr><td>T</td><td>1</td><td>0</td><td>1</td><td>0</td><td>0</td><td>1</td></tr> </table>	A	0	1	0	1	1	0	C	0	0	0	0	0	0	G	0	0	0	0	0	0	T	1	0	1	0	0	1
A	0	1	0	1	1	0																							
C	0	0	0	0	0	0																							
G	0	0	0	0	0	0																							
T	1	0	1	0	0	1																							
counts	<table border="1"> <tr><td>A</td><td>0</td><td>11</td><td>4</td><td>7</td><td>9</td><td>0</td></tr> <tr><td>C</td><td>1</td><td>0</td><td>1</td><td>2</td><td>1</td><td>0</td></tr> <tr><td>G</td><td>1</td><td>0</td><td>1</td><td>2</td><td>1</td><td>0</td></tr> <tr><td>T</td><td>10</td><td>1</td><td>6</td><td>1</td><td>1</td><td>12</td></tr> </table>	A	0	11	4	7	9	0	C	1	0	1	2	1	0	G	1	0	1	2	1	0	T	10	1	6	1	1	12
A	0	11	4	7	9	0																							
C	1	0	1	2	1	0																							
G	1	0	1	2	1	0																							
T	10	1	6	1	1	12																							
percent	<table border="1"> <tr><td>A</td><td>2</td><td>95</td><td>26</td><td>59</td><td>51</td><td>1</td></tr> <tr><td>C</td><td>9</td><td>2</td><td>14</td><td>13</td><td>20</td><td>3</td></tr> <tr><td>G</td><td>10</td><td>1</td><td>16</td><td>15</td><td>13</td><td>0</td></tr> <tr><td>T</td><td>79</td><td>3</td><td>44</td><td>13</td><td>17</td><td>96</td></tr> </table>	A	2	95	26	59	51	1	C	9	2	14	13	20	3	G	10	1	16	15	13	0	T	79	3	44	13	17	96
A	2	95	26	59	51	1																							
C	9	2	14	13	20	3																							
G	10	1	16	15	13	0																							
T	79	3	44	13	17	96																							
log-odds	<table border="1"> <tr><td>A</td><td>-38</td><td>19</td><td>1</td><td>12</td><td>10</td><td>-48</td></tr> <tr><td>C</td><td>-15</td><td>-38</td><td>-8</td><td>-10</td><td>-3</td><td>-32</td></tr> <tr><td>G</td><td>-13</td><td>-48</td><td>-6</td><td>-7</td><td>-10</td><td>-48</td></tr> <tr><td>T</td><td>17</td><td>-32</td><td>8</td><td>-9</td><td>-6</td><td>19</td></tr> </table>	A	-38	19	1	12	10	-48	C	-15	-38	-8	-10	-3	-32	G	-13	-48	-6	-7	-10	-48	T	17	-32	8	-9	-6	19
A	-38	19	1	12	10	-48																							
C	-15	-38	-8	-10	-3	-32																							
G	-13	-48	-6	-7	-10	-48																							
T	17	-32	8	-9	-6	19																							

$$S = f_{pos,b}$$

$$S = \log\left(\frac{f_{pos,b}}{p_b}\right)$$

15

Consensus Patterns and Motifs

- How to represent motifs
 - consensus patterns
 - weight matrices
- How to “weight” positions?
 - frequency $S = f_{pos,b}$
 - information content Σ (log odds) $S = \log\left(\frac{f_{pos,b}}{p_b}\right)$
- How to search for motifs
 - heuristic greedy
 - Gibbs sampling

Representing Consensus Sequences

counts	<table border="1" style="border-collapse: collapse; width: 100%;"> <tr><td>A</td><td>9</td><td>214</td><td>63</td><td>142</td><td>118</td><td>8</td></tr> <tr><td>C</td><td>22</td><td>7</td><td>26</td><td>31</td><td>52</td><td>13</td></tr> <tr><td>G</td><td>18</td><td>2</td><td>29</td><td>38</td><td>29</td><td>5</td></tr> <tr><td>T</td><td>193</td><td>19</td><td>124</td><td>31</td><td>43</td><td>216</td></tr> </table>	A	9	214	63	142	118	8	C	22	7	26	31	52	13	G	18	2	29	38	29	5	T	193	19	124	31	43	216
A	9	214	63	142	118	8																							
C	22	7	26	31	52	13																							
G	18	2	29	38	29	5																							
T	193	19	124	31	43	216																							

frequency	<table border="1" style="border-collapse: collapse; width: 100%;"> <tr><td>A</td><td>0.04</td><td>0.88</td><td>0.26</td><td>0.59</td><td>0.49</td><td>0.03</td></tr> <tr><td>C</td><td>0.09</td><td>0.03</td><td>0.11</td><td>0.13</td><td>0.22</td><td>0.05</td></tr> <tr><td>G</td><td>0.07</td><td>0.01</td><td>0.12</td><td>0.16</td><td>0.12</td><td>0.02</td></tr> <tr><td>T</td><td>0.80</td><td>0.08</td><td>0.51</td><td>0.13</td><td>0.18</td><td>0.89</td></tr> </table>	A	0.04	0.88	0.26	0.59	0.49	0.03	C	0.09	0.03	0.11	0.13	0.22	0.05	G	0.07	0.01	0.12	0.16	0.12	0.02	T	0.80	0.08	0.51	0.13	0.18	0.89
A	0.04	0.88	0.26	0.59	0.49	0.03																							
C	0.09	0.03	0.11	0.13	0.22	0.05																							
G	0.07	0.01	0.12	0.16	0.12	0.02																							
T	0.80	0.08	0.51	0.13	0.18	0.89																							

log(2)-odds	<table border="1" style="border-collapse: collapse; width: 100%;"> <tr><td>A</td><td>-2.76</td><td>1.82</td><td>0.06</td><td>1.23</td><td>0.96</td><td>-2.92</td></tr> <tr><td>C</td><td>-1.46</td><td>-3.11</td><td>-1.22</td><td>-1.00</td><td>-0.22</td><td>-2.21</td></tr> <tr><td>G</td><td>-1.76</td><td>-5.00</td><td>-1.06</td><td>-0.67</td><td>-1.06</td><td>-3.58</td></tr> <tr><td>T</td><td>1.67</td><td>-1.66</td><td>1.04</td><td>-1.00</td><td>-0.49</td><td>1.84</td></tr> </table>	A	-2.76	1.82	0.06	1.23	0.96	-2.92	C	-1.46	-3.11	-1.22	-1.00	-0.22	-2.21	G	-1.76	-5.00	-1.06	-0.67	-1.06	-3.58	T	1.67	-1.66	1.04	-1.00	-0.49	1.84
A	-2.76	1.82	0.06	1.23	0.96	-2.92																							
C	-1.46	-3.11	-1.22	-1.00	-0.22	-2.21																							
G	-1.76	-5.00	-1.06	-0.67	-1.06	-3.58																							
T	1.67	-1.66	1.04	-1.00	-0.49	1.84																							

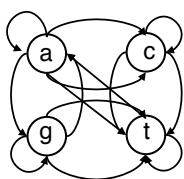
position-specific
information content
(bits)

$$I_p = \sum_b^{A,T} f_b \log_2 \left(\frac{f_b}{p_b} \right)$$

I	0.98	1.33	0.28	0.38	0.22	1.35
---	------	------	------	------	------	------

17

Where do scoring matrices come from?



$$q_{ij} = M^{20} = \text{PAM20}$$

$$\begin{aligned} &\{0.828, 0.019, 0.133, 0.019\}, \\ &\{0.019, 0.828, 0.019, 0.133\}, \\ &\{0.133, 0.019, 0.828, 0.019\}, \\ &\{0.019, 0.133, 0.019, 0.828\} \end{aligned}$$

$$\lambda S = \log \left(\frac{q_{ij}}{p_j} \right)$$

probability of mutation

probability of alignment by chance

$$p_i(a,c,g,t) =$$

$$p_j = 0.25$$

$$\lambda S = 10 \log \left(\frac{q_{a,c}}{p_c} \right)$$

$$= 10 \log \left(\frac{0.019}{0.25} \right) = -11.2$$

$$\lambda_2 = \frac{\log(2)}{10} = 0.33$$

fasta.bioch.virginia.edu/biol4230

18

Ranking criterion: Information Content

$$\text{Information content} = \frac{1}{N} \log \text{Likelihood Ratio}$$

$$I_{total} = \sum_{pos=1}^n \sum_b^{A,T} f_{pos,b} \log_2 \frac{f_{pos,b}}{p_b}$$

Position independent
Scoring Matrix

$$S = \log \left(\frac{f_{ij}}{p_j} \right)$$

Position Specific
Scoring Matrix (PSSM)

$$S = \log \left(\frac{f_{pos,b}}{p_b} \right)$$

fasta.bioch.virginia.edu/biol4230

19

Information content and binding energy

If we denote H_i the *binding energy* for a DNA site S_i , then the probability that the protein would be bound to S_i (at equilibrium) is given by the Boltzman distribution:

$$P_i = \frac{e^{-H_i}}{Z}$$

where Z is the *partition function* and is defined as the sum of the e^{-H_x} over all possible sites S_x :

$$Z = \sum_x e^{-H_x}$$

The average binding energy for this protein over all sites S_i would be:

$$\langle H \rangle = \sum_i P_i H_i = -\sum_i (P_i \ln P_i) - \ln Z$$

The sum on the right is called the entropy of the probability distribution.

A useful measure of difference between two probability distributions is the relative entropy, which is defined as:

$$H(P, Q) \equiv \sum_i P_i \ln \frac{P_i}{Q_i}$$

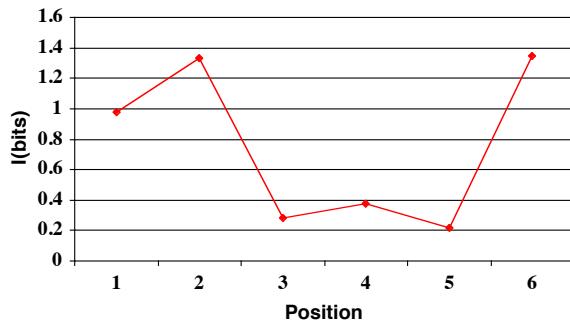
Benos, P. V., Lapedes, A. S. & Stormo, G. D. *Bioessays* **24**, 466–475 (2002). fasta.bioch.virginia.edu/biol4230

20

Consensus Information Content

$$I_p = \sum_b^{A..T} f_b \log\left(\frac{f_b}{p_b}\right)$$

I	0.98	1.33	0.28	0.38	0.22	1.35
---	------	------	------	------	------	------



fasta.bioch.virginia.edu/biol4230

21

A	T	C	G	T	A	T	C	G	T	A	T	C	G	T	A	T	C	G	T	A	T	C	G
Col E1 site 2	T	T	T	T	G	G	G	G	T	A	T	A	C	T	G	T	A	C	G	A	A	A	A
Col E1 site 1	A	T	A	T	A	T	G	T	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
att site 1	A	T	A	T	A	T	G	T	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
bsf R mut	T	A	T	A	C	T	G	T	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
CRP	A	G	A	T	A	C	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
cys	A	G	A	T	A	C	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
deo P2 site 2	A	T	G	A	T	G	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
deo P2 site 1	A	T	G	A	T	G	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
deo	A	T	G	A	T	G	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
ivv B	A	T	G	A	T	G	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lac	T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lacZ	T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lacZ E	T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lacZ K	T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lacZ M	T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
lacZ T	A	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A	A
tms A	A	G	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	A	A	A	A
tms B	T	A	G	A	T	A	T	G	T	A	T	C	A	T	G	T	A	C	G	A	A	A	A
pBR 74	C	G	C	G	C	T	G	T	G	A	T	C	A	T	G	T	A	C	G	T	A	T	T
cba site 2	A	T	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	T	A	T	T
cba site 1	A	T	A	T	A	T	G	T	G	A	T	C	A	T	G	T	A	C	G	T	A	T	T
ido	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-

frequency:

A	T	C	G
0.48	0.48	0.39	0.04
0.04	0.04	0.04	0.92
0.29	0.29	0.13	0.09

bits:

A	T	C	G
0.94	0.94	0.64	-2.64
-2.64	0.94	-2.64	0.94
-0.94	-0.94	-0.94	-0.94
0.04	0.04	0.04	0.04
0.04	0.04	0.04	0.04
0.13	0.13	0.04	0.38
0.63	0.63	0.38	0.22
0.35	0.35	0.12	0.12
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09
0.22	0.22	0.09	0.09

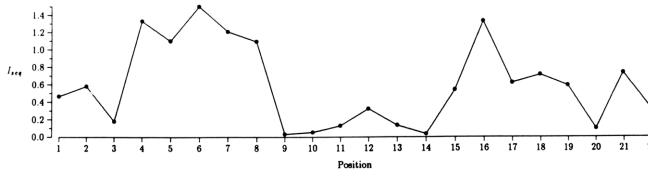


FIG. 1. CRP binding sites. (A) Twenty-three sites identified as binding to CRP (8, 9). (B) The frequency with which each base occurs at each position in the CRP binding sites. (C) The specificity matrix for the protein, based on the binding sites, which is calculated as $\log(f_b/p_b)$, where f_b is the observed frequency of each base (from the matrix above) and p_b is the a priori probability of obtaining base b . In this example, $p_b = 0.25$ for all b , approximating the *E. coli* genomic composition. At positions for which $f_b = 0$, an estimated frequency of 0.5/23 is used in the calculation. (D) The "information content" at each position of the CRP binding sites is plotted (10). The sum of all positions is 13.06 bits.

Stormo GD, Hartzell GW (1989) Identifying protein-binding sites from unaligned DNA fragments. Proc Natl Acad Sci U S A. 86:1183-1187.

Consensus Patterns and Motifs

- How to represent motifs
 - consensus patterns
 - weight matrices
- How to “weight” positions?
 - frequency
 - information content ($\Sigma \log$ odds)
- How to search for motifs
 - heuristic greedy (consensus, wconsensus)
 - Expectation-Maximization (MEME)
 - Gibbs sampling

fasta.bioch.virginia.edu/biol4230

23

Identifying Functional Domains in Biological Sequences

A problem in:
Feature Detection
or
Multiple Alignment

Two parts to the problem:

1. Can't look at all possible alignments,
pick a subset likely to contain the answer
2. Need criterion for ranking alignments
that is reasonable and efficient

fasta.bioch.virginia.edu/biol4230

24

E. coli CRP binding regions

```
CE1CG      taatgtttgtctgggtttgtggcatcggcgagaatacgccgttgtgaaagactgttttgcgtttcacaaaaatggaaaggccacagtcttgacag  
ECOARABOP  gacaaaacgcgtacaaaagggtctataatcaccgcggaaaaggccatgttgcacggcgtacacttgcgtatgcctagcatttatccataag  
ECOBGLR1   acaaaatccaaataacttaatttggatttttatataacttataatccctaaaattacacaaatgttgcgtatgcctagcatttatccataat  
ECOCRP    cacaaggcggaaatgtctaaaacaggcaggatgtcactacatgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOCYA    acgggtctacacttgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECODEOP2   agtgaattttgaaccacatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOGALE    ggcataaaaaacgcgttaattttgtgttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOILVBPR  gtcgcgggggtttttttttatctgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOLAC     aacgcattaaatgtgttagtgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALBA   acattaccgcatttcgttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALBA   ggaggaggcggggaggatgagaacacggcttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALT    gatcgcgttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOOMPA    gctgacaaaaagattaaacatccatatacagactttttcatatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOTNAA    tttttaaacattaaattttatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOUXU1    cccatgagatgttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttt  
PBR322     ctggcttaactatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
TRN9CAT    ctgtgacggaaagatcaatccgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
TDC        gattttataactttatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc
```

fasta.bioch.virginia.edu/biol4230

25

E. coli CRP binding sites

```
CE1CG      taatgtttgtctgggtttgtggcatcggcgagaatacgccgttgtgaaagactgtttTTTGATCCTTCACAAAatggaaaggccacagtcttgacag  
ECOARABOP  gacaaaacgcgtacaaaagggtctataatcaccgcggaaaaggccatgttgcacggcgtacacttgcgtatgcctagcatttatccataag  
ECOBGLR1   acaaaatccaaataacttaatttggatttttatataacttataatccctaaaattacacaaatgttgcgtatgcgtatgcgtatgcgtatgc  
ECOCRP    cacaaggcggaaatgtctaaaacaggcaggatgtcactacatgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOCYA    acgggtctacacttgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECODEOP2   agtgaattTTGAAACAGATCGCATTCAGTGTatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOGALE    ggcataaaaaacgcgttaattttgtgttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOILVBPR  gtcgcgggggtttttttttatctgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOLAC     aacgcattaaTGTAACTTGCATCAGTGTatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALBA   acattaccgcatttcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALBA   ggaggaggcggggaggatgagaacacggcttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOMALT    gatcgcgttgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOOMPA    gctgacaaaaagattaaacatccatatacagactttttcatatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOTNAA    tttttaaacattaaattttatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
ECOUXU1    cccatgagatgttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttttt  
PBR322     ctggcttaactatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
TRN9CAT    ctgtgacggaaagatcaatccgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc  
TDC        gattttataactttatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgcgtatgc
```

fasta.bioch.virginia.edu/biol4230

26

Ranking criterion: Information Content

Information content = $\frac{1}{N} \text{LogLikelihood Ratio}$

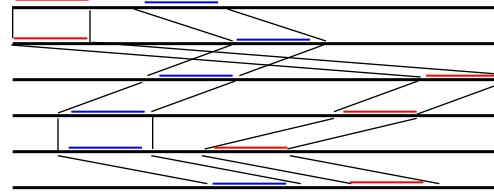
$$I_{total} = \sum_{pos=1}^n \sum_b^{A,T} f_{pos,b} \log_2 \frac{f_{pos,b}}{p_b}$$

Consensus Patterns and Motifs

- How to represent motifs
 - consensus patterns
 - weight matrices
- How to “weight” positions?
 - frequency
 - information content ($\Sigma \log$ odds)
- How to search for motifs
 - heuristic greedy (consensus, wconsensus)
 - Expectation-Maximization (MEME)
 - Gibbs sampling

Finding a consensus sequence: consensus

A C T G A A T
A G C G T C C
C T T G C C G



fasta.bioch.virginia.edu/biol4230

31

Finding a consensus sequence: consensus

A C T G A A T
A G C G T C C
C T T G C C G

I=12.0	A C T G A A T
A	1 0 0 0 1 1
C	0 1 0 0 0 0
G	0 0 0 1 0 0
T	0 0 1 0 0 0

I=12.0	C T G A A T
A	0 0 0 1 1 0
C	1 0 0 0 0 0
G	0 0 1 0 0 0
T	0 1 0 0 0 1

I=8.0	A C T G A A T
A	1 0 0 0 1 1
G	0 1 1 0 0 1
C	0 1 0 2 0 0
T	0 0 1 0 1 0

I=7.0	A C T G A A T
A	1 0 0 0 1 1
G	0 2 0 0 1 0
C	0 1 1 0 0 1
T	0 0 1 1 0 0

I=6.0	C T G A A T
A	1 0 0 1 1 0
C	1 0 1 0 0 1
G	0 1 1 1 0 0
T	0 1 0 0 1 1

I=7.0	C T G A A T
A	0 0 0 1 1 0
C	1 1 0 0 1 1
G	1 0 2 0 0 0
T	0 1 0 1 0 1

I=6.1	A C T G A A T
A	2 0 0 0 1 1
G	1 1 1 0 1 2
C	0 1 0 3 0 0
T	0 1 2 0 1 0

I=3.8	A C T G A A T
A	2 0 0 0 1 1
G	0 1 1 1 1 1
C	0 1 1 2 0 1
T	1 1 1 0 1 0

I=5.8	C T G A A T
A	0 0 0 1 1 0
C	2 1 0 0 2 2
G	1 0 2 1 0 0
T	0 2 1 1 0 1

I=5.4	C T G A A T
A	0 0 0 1 1 0
C	1 1 0 1 2 1
G	1 0 3 0 0 1
T	1 2 0 1 0 1

fasta.bioch.virginia.edu/biol4230

32

E. coli CRP binding sites

```

CE1CG      taatgtttgtctgggatcgccgagaatgcgcgtgtgaaagactgtttTGATGTTTACAAAatggaagtccacagtcttacag
ECOARABOP   gacaaaacgcgtacaacaaaagtgtataatcacgcggaaaaggccatgttaaTTCGACGGCGTCACACTTtgctatgccatagcattttatccataag
ECOBGLR1    acaaaatccaataacttaattatggatttttatataaaatcttataaaattccatataaTGTGAGCATGGTCATATTtttatcaat
ECOCRP     cacaagcgaaagatgtctaaaacagtcaggatgtcaactatgtatgtactgcTGCAAGGACGTCACATTaccgtgcgtacatggtagc
ECOCYA     acggtgatcacatgttagcgcacccggatcgtcaaggTGTAAATTGATCAGTTttagaccattttcgtgtgaactaaaaaaac
E CODEOP2    agtgaattTGTGACCAGATCCATTacatgtatgtcaacttgtttagatgttttcttatttgtgtgtatcgaagtgtggagtagatgtttagata
ECOGALE    ggcataaaaaaacggctaatttttgtgtaaaacggatccactaattttatccTCTCACACTTTCCATCTtttgcattatgttgcattaccataaggcc
ECOILVBPR   gctccgggggtttttttaTCTGCAATTCAGTACAACactgtatcacccctcttttgcgtgaaaaattttccatttgcgttgc
ECOLAC      aacgcattaaTGTGAGTTAGCTCACTCattaggcaccccgggttacactttatgttcccgctgtatgtgtggattgtgagcataaatttcac
ECOMALBA    acattaccgcattTGTAAACAGAGATCACAaaggcgtttggcgatggggcaaggaggatggcgtataaaaaactagatgtcggttta
ECOMALBA    ggaggaggcgagggatggaaacgggttgcatttttgtgaattCCTGACGTTGCTGCAAatcgtgggattttatgtggc
ECOMALT     gatcgcgtcggttttaggtgttataaaaggattttgaatTCTGACACAGTGCAAATTcacacataaaaaacgtcatcgcgttgcatt
ECOOMPA     gctgacaaaaaaaatttttaaaatccatgctatatataatttttttttgcattatCCTGACGTTGCTGCAAatcgtgggattttatgtggc
ECOTNAA     ttttttaacattaaaatccatgctatatataatttttttttgcattttttggcgcacgtTGTGATTGATCACATttaaacatttcga
ECOUXU1     cccatgaggtaattgtTGTGATGTGTTAACCAatttttagatcggatttttgatttttacccaaaggtagaaattatcgc
PBR322      ctggotttaatcatgcggcatcagagatgtactggatggaccatttttgTGTGAAATACCCACAGAtgctaaggaaaaatccgc
TRN9CAT     ctgtgacggaaatgcacccgctatattttatcgttgcctgtatccggggccccccattttggggaaaaaTGAGACGTTGATGGCAcg
TDC         gatttttttataacttgtattttaaaggatttttattttatgttgcacggatTGTGAGTGTTGCACATAtctgtt

```

fasta.bioch.virginia.edu/biol4230

33

Aligned CRP binding sites

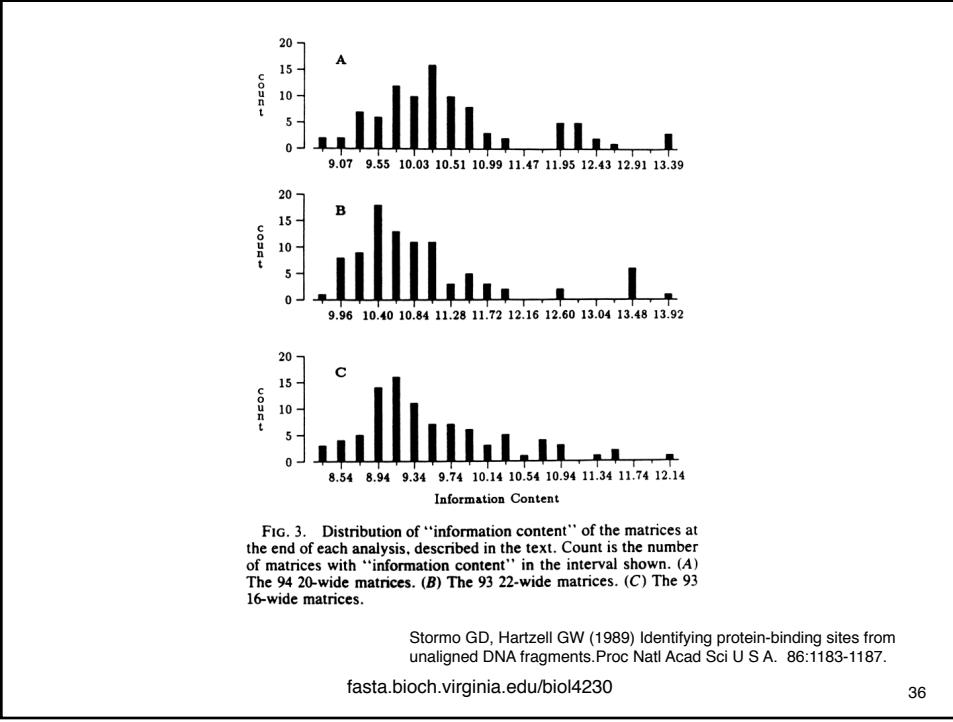
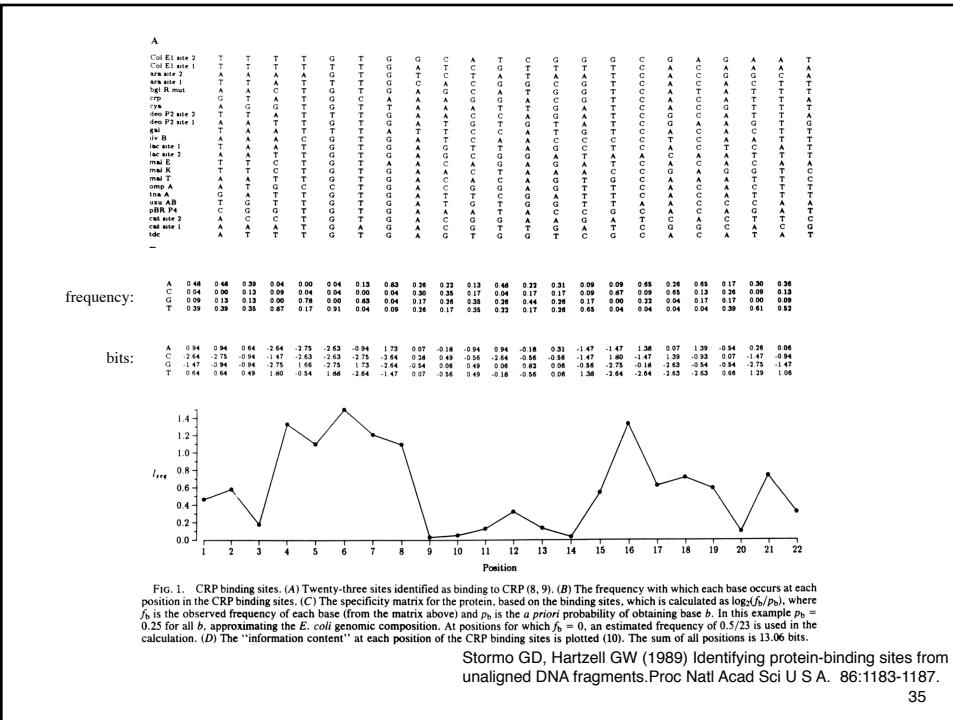
```

CE1CG      cgcgtgtgaaagactgtttTGATGTTTACAAAatggaagtccacagtcttgcacag
ECOARABOP  gcagaaaacgcacatgttttaTTCGACGGCGTCACACTTtgctatgccatagcattttatccataag
ECOBGLR1    taaaatcacacaaaatgttaaTGTGAGCATGGTCATATTtttatcaat
ECOCRP     aatacattgtgtactgcTGCAAGGACGTCACATTaccgtgcgtacatggtagc
ECOCYA     ttcttacggatcgtcaaggTGTAAATTGATCAGTTttagaccattttcgtgtaaactaaa
E CODEOP2    agtgaattTGTGACCAGATCCATTacatgtatgtcaacttgtttagatgttttcttatttgtgtgtatcgaagtgtggagtagatgtttagata
ECOGALE    aacgatccactattttccaTCTCACACTTTCCATCTttttgcattatgttgcattaccata
ECOILVBPR   gctccgggggtttttttaTCTGCAATTCAGTACAACactgtatcacccctcttttgcgtgaaaaattttccatttgcgttgc
ECOLAC      aacgcattaaTGTGAGTTAGCTCACTCattaggcaccccgggttacactttatgttcccggttacactttatgttcccggt
ECOMALBA    acattaccgcattTGTAAACAGAGATCACAaaggcgtttggcgatggggcaaggaggatggcgtataaaaaactagatgtcggttta
ECOMALBA    aaccgggtatggattCCTGATGTGTTGCAAAatcgtggattttatgtggc
ECOMALT     gatgttaataaaatggattTGTGACACAGTGCAAttttcacacataaaaaaacgtcatcgcgttgcattaaaaaa
ECOOOMPA   tatacaagacttttttcatatCCTGACGGAGTTGACACTTttaaaccattttcaactacgttgtagactttatcgcc
ECOT       attaatattgtcccccgaacgtTGTGATTGATCACATttaaacatttcaga
ECOUXU1     cccatgaggtaattgtTGTGATGTGTTAACCAattttttagattttcgggattttgatttttacccaaaggta
PBR322      gtactgaggatgcacccatttttgTGTGAAATACCCACAGAtgctaaggaaaaatccgc
TRN9CAT     ccctggccaactttggcaaaaTGAGACGTTGATGGCAcg
TDC         tctggaaattttttgtattttaaaggatttttattttatgttgcacggatTGTGAGTGTTGCACATAtctgtt

```

fasta.bioch.virginia.edu/biol4230

34



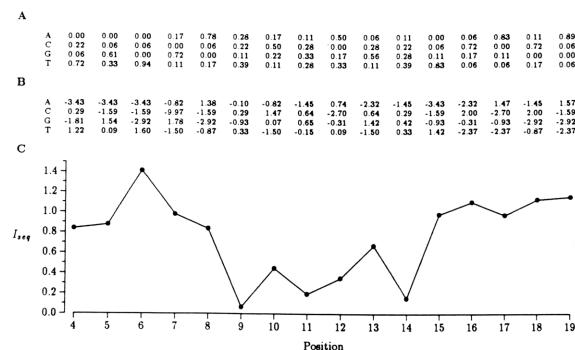


FIG. 4. The best 16-wide matrix. The positions are numbered 4–19, corresponding to the central positions of Fig. 1. (A) The frequency of each base at the sites included in the best matrix, as in Fig. 1B. (B) The specificity matrix determined from the frequency matrix, as in Fig. 1C. In this case the ρ values were varied. The analyses in Fig. 2; $\rho_A = \rho_C = \rho_G = \rho_T = 0.2$; $\rho_A = \rho_C = \rho_G = 0.21$; and $\rho_T = 0.31$. The analyses were also performed using $\rho_b = 0.25$ for all b . In that case, essentially the same matrix remained the best, but the distribution had more high-scoring matrices due to the high probability of adenine and thymine matches. The specificity matrix values for positions with $f_b = 0$ were estimated using $f_b = 0.5/18$ from the 18 sequences in the data set. (C) The “information content” at each position of the matrix is plotted. The sum from all positions is 12.15 bits.

Stormo GD, Hartzell GW (1989) Identifying protein-binding sites from unaligned DNA fragments. Proc Natl Acad Sci U S A. 86:1183-1187.

fasta.bioch.virginia.edu/biol4230

37

consensus -L 16 -q 1000 -c0 -pr2 -pt 4 -pf 4
L-mer Width: 16

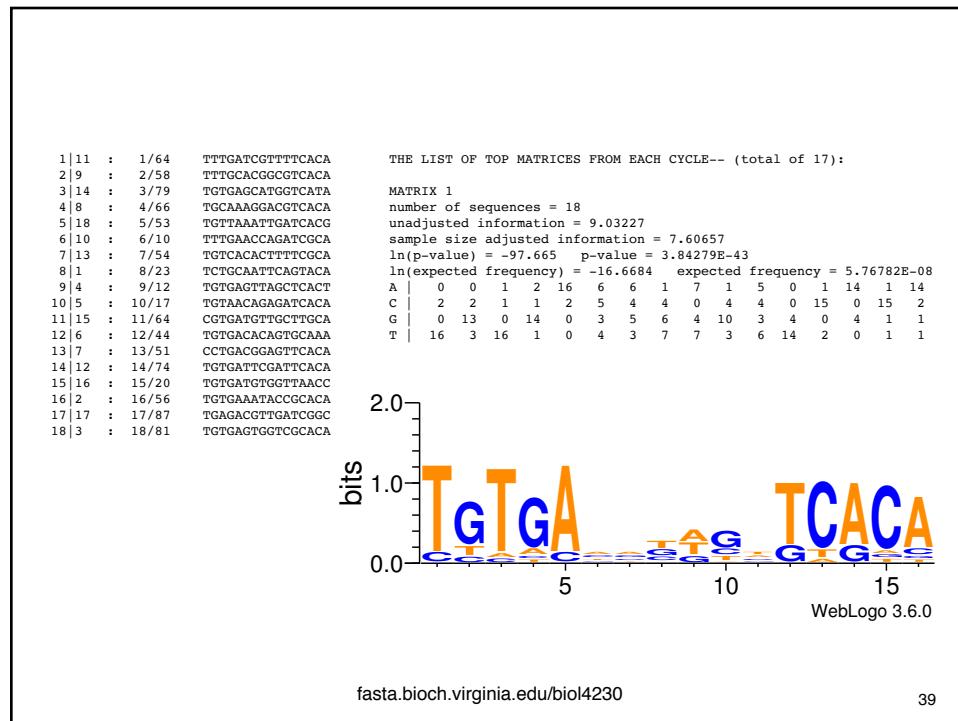
Top Matrices saved from each cycle: 4
Matrices Saved from the last cycle: 4
sequence 1: CEB1CG fragments: 1-105
sequence 18: TDC fragments: 1-105
Total number of sequences: 18.
Total number of sequence fragments: 18.

Consensus (CRP)

```
# Observed frequency and occurrence of each letter.  
#number of letters in the input sequences = 1890  
A 0.302646; observed occurrence = 572 (letter 1)  
C 0.182540; observed occurrence = 345 (letter 2)  
G 0.208995; observed occurrence = 395 (letter 3)  
T 0.305820; observed occurrence = 578 (letter 4)  
  
PRIOR FREQUENCIES DETERMINED BY OBSERVED FREQUENCIES.  
* Information for the alphabet from the command line.  
letter 1: A (complement: T) prior freq = 0.302646 CYCLE [] total number information ln top ln expected frequency []  
letter 2: C (complement: G) prior freq = 0.182540 []-----[]-----[]-----[]-----[]-----[]-----[]-----[]-----[]-----[]  
letter 3: G (complement: C) prior freq = 0.208995 1 [] 1620 2.9492 0.0000 [] 7.3902 []  
letter 4: T (complement: A) prior freq = 0.305820 2 [] 748 6.8222 -13.5825 [] 0.4475 []  
3 [] 849 8.2051 -20.1461 [] 0.0577 []  
4 [] 832 8.7882 -26.3882 [] -0.3628 []  
5 [] 848 8.9275 -31.8728 [] -0.3179 []  
6 [] 857 9.1860 -38.8902 [] -2.0624 []  
7 [] 877 9.2908 -45.6680 [] -3.8014 []  
8 [] 864 9.1929 -51.3100 [] -4.6251 []  
9 [] 876 9.1152 -57.2498 [] -5.9597 []  
10 [] 879 9.0644 -63.5596 [] -7.8751 []  
11 [] 864 8.8973 -68.7012 [] -8.8353 []  
12 [] 854 8.8738 -75.4184 [] -11.5917 []  
13 [] 850 8.6817 -80.0738 [] -12.5205 []  
14 [] 873 8.5267 -85.0113 [] -13.9878 []  
15 [] 852 8.2955 -88.6578 [] -14.4562 []  
16 [] 865 8.1066 -92.6288 [] -15.6014 []  
17 [] 857 7.8793 -95.7075 [] -16.3204 []  
18 [] 936 7.6066 -97.6650 [] -16.6684 []
```

fasta.bioch.virginia.edu/biol4230

38



```

consensus -L 16 -q 1000 -c0 -pr2 -pt 4 -pf 4
L-mer Width: 16

Top Matrices saved from each cycle: 4
Matrices Saved from the last cycle: 4

sequence 1: CEB1CG      fragments: 1-105
sequence 18: TDC        fragments: 1-105
Total number of sequences: 18.
Total number of sequence fragments: 18.

# Observed frequency and occurrence of each letter.
#number of letters in the input sequences = 1890
A 0.302646; observed occurrence = 572 (letter 1)
C 0.182540; observed occurrence = 345 (letter 2)
G 0.208995; observed occurrence = 395 (letter 3)
T 0.305820; observed occurrence = 578 (letter 4)

PRIOR FREQUENCIES DETERMINED BY OBSERVED FREQUENCIES.
* Information for the alphabet from the command line.
letter 1: A (complement: T) prior freq = 0.302646 CYCLE [] number [ ] information [ ] ln top [ ] ln expected [ ] p-value [ ] frequency [ ]
letter 2: C (complement: G) prior freq = 0.182540 [ ]-----[ ]-----[ ]-----[ ]-----[ ]
letter 3: G (complement: C) prior freq = 0.208995 1 [] 1620 | 2.9492 | 0.0000 [ ] 7.3902 [ ]
letter 4: T (complement: A) prior freq = 0.305820 2 [] 748 | 6.8222 | -13.5825 [ ] 0.4475 [ ]
letter 5: C (complement: G) prior freq = 0.208995 3 [] 849 | 8.2051 | -20.1461 [ ] 0.0577 [ ]
letter 6: G (complement: C) prior freq = 0.182540 4 [] 832 | 8.7882 | -26.3882 [ ] -0.3628 [ ]
letter 7: T (complement: A) prior freq = 0.305820 5 [] 848 | 8.9275 | -31.8728 [ ] -0.3179 [ ]
letter 8: C (complement: G) prior freq = 0.208995 6 [] 857 | 9.1860 | -38.8902 [ ] -2.0624 [ ]
letter 9: G (complement: C) prior freq = 0.182540 7 [] 877 | 9.2908 | -45.6680 [ ] -3.8014 [ ]
letter 10: T (complement: A) prior freq = 0.305820 8 [] 864 | 9.1929 | -51.3100 [ ] -4.6251 [ ]
letter 11: C (complement: G) prior freq = 0.208995 9 [] 876 | 9.1152 | -57.2498 [ ] -5.9597 [ ]
letter 12: G (complement: C) prior freq = 0.182540 10 [] 879 | 9.0644 | -63.5596 [ ] -7.8751 [ ]
letter 13: T (complement: A) prior freq = 0.305820 11 [] 864 | 8.8973 | -68.7012 [ ] -8.8353 [ ]
letter 14: C (complement: G) prior freq = 0.208995 12 [] 854 | 8.8738 | -75.4184 [ ] -11.5917 [ ]
letter 15: G (complement: C) prior freq = 0.182540 13 [] 850 | 8.6817 | -80.0738 [ ] -12.5205 [ ]
letter 16: T (complement: A) prior freq = 0.305820 14 [] 873 | 8.5267 | -85.0113 [ ] -13.9878 [ ]
letter 17: C (complement: G) prior freq = 0.208995 15 [] 852 | 8.2955 | -88.6578 [ ] -14.4562 [ ]
letter 18: G (complement: C) prior freq = 0.182540 16 [] 865 | 8.1066 | -92.6288 [ ] -15.6014 [ ]
letter 19: T (complement: A) prior freq = 0.305820 17 [] 857 | 7.8793 | -95.7075 [ ] -16.3204 [ ]
letter 20: A (complement: T) prior freq = 0.302646 18 [] 936 | 7.6066 | -97.6650 [ ] -16.6684 [ ]

INFORMATION CONTENT IS CALCULATED USING NATURAL
LOGARITHMS (i.e. BASE e). DIVIDE BY  $\ln(2)$  = 0.693 TO
CONVERT TO BASE 2, WHICH WAS USED IN PREVIOUS VERSIONS
OF THIS PROGRAM.

```

fasta.bioch.virginia.edu/biol4230

40

Consensus from random sequences

```

[]          MATRICES SAVED FOR NEXT CYCLE []
[]-
[] total    top adjusted      ln top      [] ln expected []
CYCLE [] number   information   p-value   [] frequency []
-----[]-----[]-----[]-----[]
1 [] 2040 | 1.4133 | 0.0000 [] 7.6207 []
2 [] 648 | 6.5863 | -14.0510 [] 0.4547 []
3 [] 866 | 8.2103 | -21.0667 [] -0.1259 []
4 [] 819 | 8.6626 | -26.3960 [] 0.6457 []
5 [] 849 | 8.6417 | -30.6613 [] 2.2093 []
6 [] 871 | 8.7530 | -36.3388 [] 2.1271 []
7 [] 872 | 8.7389 | -41.7409 [] 2.1122 []
8 [] 878 | 8.6112 | -46.5558 [] 2.4937 []
9 [] 875 | 8.4596 | -51.2305 [] 2.8370 []
20 [] 937 | 6.4446 | -88.4984 [] 9.6257 []
21 [] 943 | 6.3042 | -91.2183 [] 9.6902 []
22 [] 955 | 6.1507 | -93.4632 [] 9.8955 []
23 [] 947 | 5.9713 | -94.9290 [] 10.4300 []
24 [] 969 | 5.8093 | -96.4855 [] 10.1381 []

```

MATRIX 1 number of sequences = 3
unadjusted information = 18.0219
sample size adjusted information = 8.21031
ln(p-value) = -21.0667 p-value = 7.09314E-10
ln(expected frequency) = -0.125938 expected frequency = 0.88167

A	0	0	0	3	1	0	0	3	0	0	0	0	3	0	1	0
C	0	0	3	0	0	2	0	0	0	3	1	0	0	1	2	0
G	0	3	0	0	0	0	1	0	3	0	1	0	0	2	0	3
T	3	0	0	0	2	1	2	0	0	0	1	3	0	0	0	0

41

consensus greedy search for information content

- No longer used (good for teaching)
- Measure for success – information content
- Greedy strategy:
 - consider all windows of length(n) vs all windows in second sequence
 - take best pairs of windows (most information content), continue to next sequence, repeat
 - order dependent
 - time kN^2 rather than N^k

fasta.bioch.virginia.edu/biol4230

42

Statistical Strategies for Consensus Alignment - EM and Gibbs

- A problem of estimation with hidden data - the positions are easy to find if the consensus is known, and the consensus is easy to find if the positions are known
- Start with random positions, build a consensus estimate
- Apply consensus to sequences, assign probability of being a consensus, repeat
- Gibbs is similar, but a target sequence is left out and scanned at each stage

fasta.bioch.virginia.edu/biol4230

43

Expectation Maximization for Consensus Alignment

(1) Begin with a set of sequences:

```
CE1CG      taatgtttgtgctgggttttgtggc
ECOARABOP  gacaaaaacgcgtaacaaaagtgtc
ECOBGLR1   acaaatcccaataacttaatttattg
ECOCRP     cacaaggcgaaagctatgtctaaac
ECOCYA     acggtgctacacttgttatgtcgcc
ECODEOP2   agtgaattttttgaaccagatcgca
ECOGALE    ggcataaaaaacgcgtaaattctt
ECOILVBPR  gctccgggggtttttgttatctt
```

(3) Use the consensus to build a matrix

CE1CG	GTTTG	A	5	6	3	3	3
ECOARABOP	TAACA	C	0	1	0	2	2
ECOBGLR1	AATAA	G	1	0	1	1	1
		T	2	1	4	2	2
ECOCRP	AAAAC						
ECOCYA	ACGGT	f _{ns}	f ₁	f ₂	f ₃	f ₄	f ₅
ECODEOP2	ATTAA	0.3	A	0.6	0.8	0.4	0.4
ECOGALE	AAAC	0.2	C	0.0	0.1	0.0	0.2
ECOILVBPR	TATCT	0.2	G	0.2	0.0	0.1	0.1
		0.3	T	0.2	0.1	0.5	0.3

(2) Select consensus sites at random:

```
CE1CG      taatGTTGtgctgggttttgtggc
ECOARABOP  gacaaaaacgcgTAACAaaaagtgtc
ECOBGLR1   acaaatccAATAActtaatttattg
ECOCRP     cacaaggcgaaAAAACgtatgtctaaac
ECOCYA     ACGGTgtacacttgttatgtcgcc
ECODEOP2   agtATTAAtttgaaccagatcgca
ECOGALE    ggcataaaAAACggctaaattctt
ECOILVBPR  gctccgggggtttttgtTATCT
```

(4) Use the consensus to weight each position

$$p_i = \sum_{j=0..5} f_{sb} + \sum_{non-site} f_{ns}$$

(5) Use the weighted site to build a new consensus

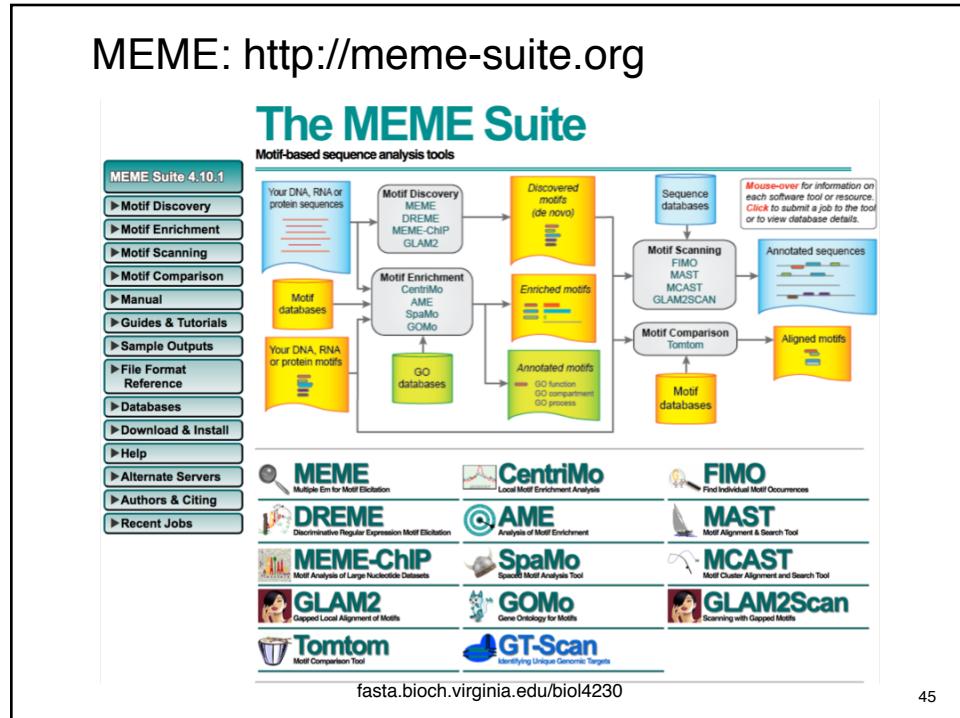
(6) Repeat (3..5)

Lawrence, C. E. & Reilly, A. A.
PROTEINS 7, 41–51 (1990).

44

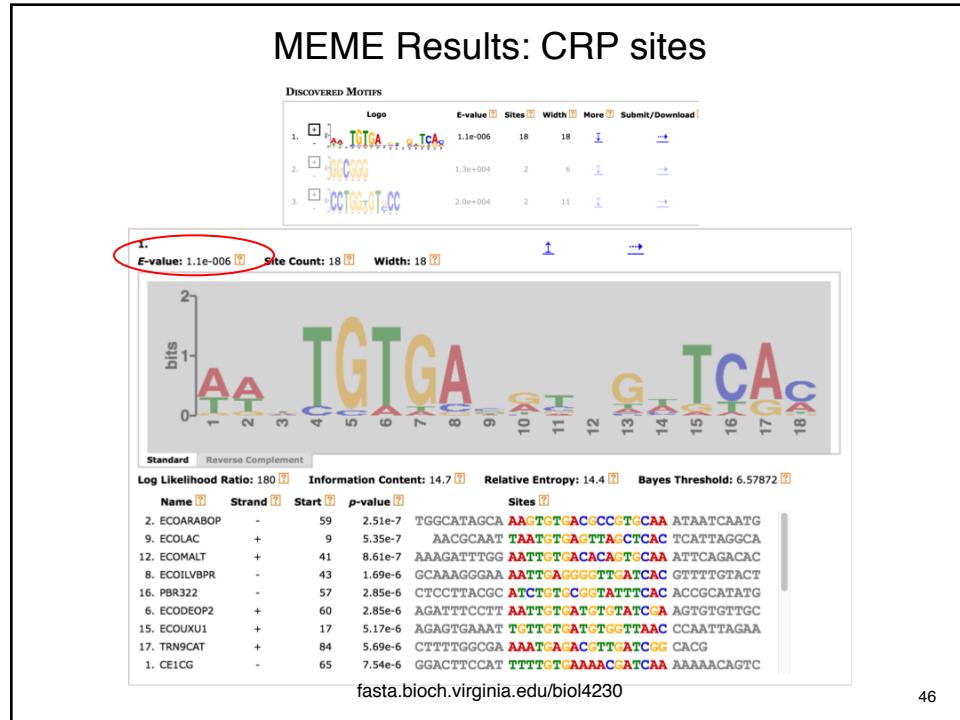
fasta.bioch.virginia.edu/biol4230

MEME: <http://meme-suite.org>



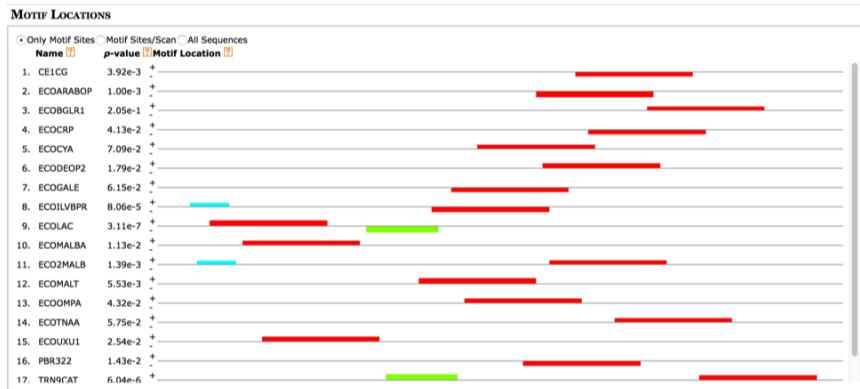
45

MEME Results: CRP sites



46

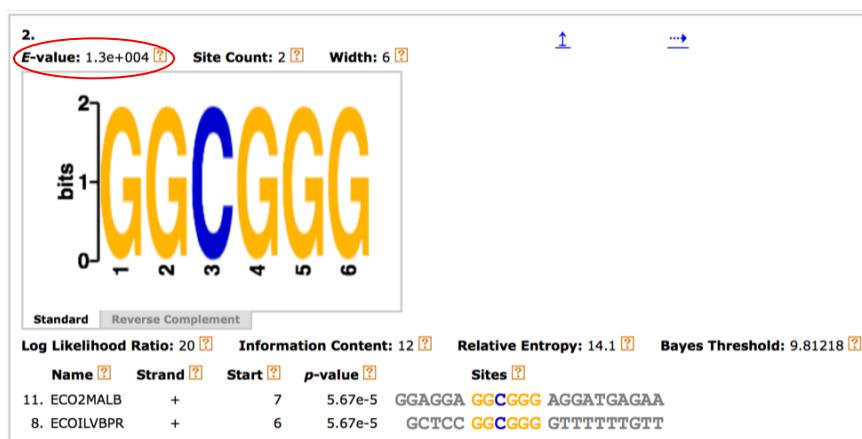
MEME Results: CRP site positions



fasta.bioch.virginia.edu/biol4230

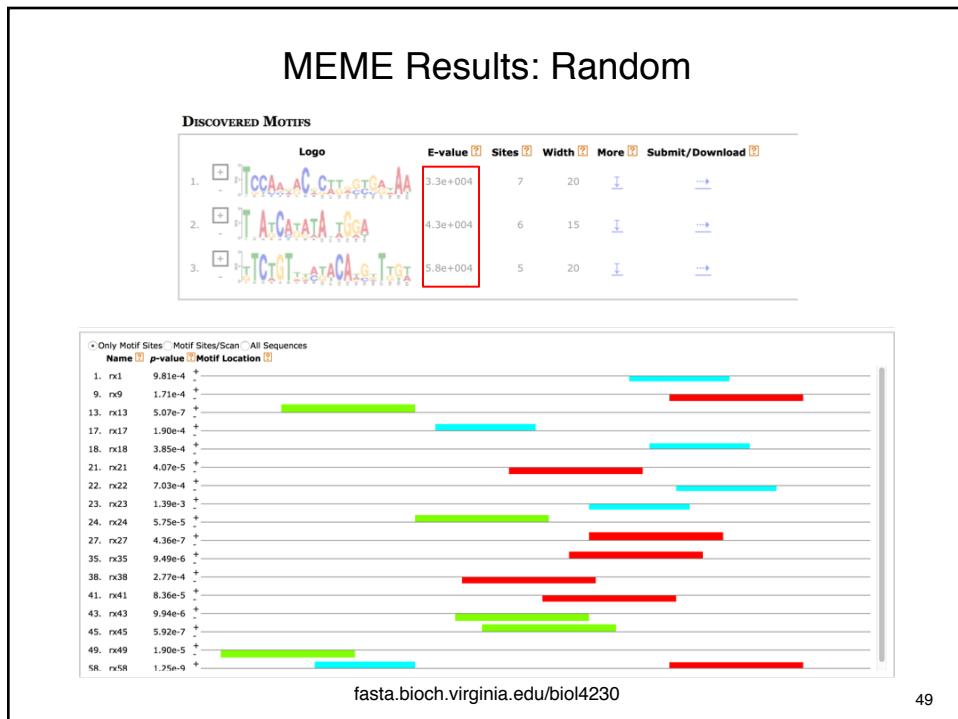
47

MEME Results: CRP - Motif 2



fasta.bioch.virginia.edu/biol4230

48



MEME sensitivity (recall) and specificity (precision)

dataset name	Output of MEME+		Analysis of discovered motifs					
	pass	log likelihood of discovered motif (d)	motif name	recall	precision	log likelihood of known motif (k)	difference ($d - k$)	
lipocalin	1 2	-55013 -55057	lipA lipB	1.000 0.400	0.357 0.200	-55090 -55092	77 35	
hth	1	-496332	hth	0.933	0.571	-496346	14	
farn	1 2 3	-92518 -92585 -92569	farnL farnB farnA	0.917 0.615 0.733	0.880 0.842 0.647	-92525 -92517 -92566	7 -68 -3	
crp	1	-60547	crp	0.792	0.905	-60590	43	
lexa	1	-109155	lexa	0.842	0.615	-109147	-8	
cplexa	1 2	-169923 -170048	lexa crp	0.842 0.667	0.696 0.471	-169918 -170116	-5 68	

Table 2: Overview of results of MEME+ on test datasets. MEME+ was run with W set to the values shown in Table 1 and $\beta = 0.01$. The *log likelihood* values are base-2 logarithms.

fasta.bioch.virginia.edu/biol4230

50

Bailey and Elkan, (1994) "Fitting a mixture model by expectation maximization to discover motifs in biopolymers", *Proceedings of the Second International Conference on Intelligent Systems for Molecular Biology*, pp. 28-36

Gibb's sampling for Consensus Alignment

- (1) Begin with a set of sequences with randomly located motifs:

```

CE1CG      taatGTTTGtgtgggttttgtgc
ECOARABOP  gacaaaacgcgTAACAaagtgtc
ECOBGLR1   acaaattcccATAActaatttattg
ECOCRP    cacaaggcgaaagctatgtAAAC
ECOCYVA   ACGGTgctacacttgatgtgcgc
ECODEOP2   agtgaATTAtttgacccatgcga
ECOGALE   ggcataaaAAACggcttaattctt
ECOILVBPR  gtcggcggggttttgtTATCT
```

- (3) Using the probability matrix from the included sequences, calculate the probability of each site on the excluded sequence

```
ECOCRP      cacaaggcgaaagctatgtctaaac
```

- (4) Select a site at random, using weights from the probabilities in (3)

- (5) Repeat steps (2) - (4)

- (2) Exclude one of the sequences at random, and build a consensus matrix from the other motifs

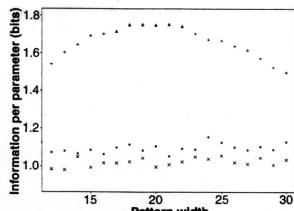
```

CE1CG      taatGTTTGtgtgggttttgtgc
ECOARABOP  gacaaaacgcgTAACAaagtgtc
ECOBGLR1   acaaattcccATAActaatttattg
ECOCRP    cacaaggcgaaagctatgtAAAC
ECOCYVA   ACGGTgctacacttgatgtgcgc
ECODEOP2   agtgaATTAtttgacccatgcga
ECOGALE   ggcataaaAAACggcttaattctt
ECOILVBPR  gtcggcggggttttgtTATCT
```

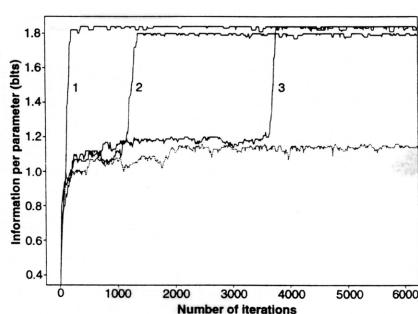
fasta.bioch.virginia.edu/biol4230

51

g. 2. Information per parameter as the criterion pattern width for helix-turn-helix (HTH) proteins. The points indicate the maximum values of information per parameter found by the algorithm. The upper points (\blacktriangle and \blacktriangledown) used the complete sequences of the 30 HTH proteins listed in Fig. 1A. (\blacktriangle) All of the sequences in the data set were aligned in the correct register (as Fig. 1A). (+) One or more of the sequences in the data set were incorrectly aligned. All completely correct alignments in the width range from 17 to 22 residues gave greater values of information per parameter than any incorrect alignments outside this width range. (●) The nonsites' sequence data of the 30 HTH proteins, constructed by deleting the 18 residues of the H pattern itself (Fig. 1A) from each of the sequences. (\times) A shuffled data set (46) of the 30 HTH sequences. The alignments from the nonsites background of the HTH proteins give values slightly greater than random expectation.



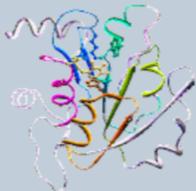
i. 3. Convergence behavior of the Gibbs sampling algorithm. Because the Gibbs sampler, when run for a long time, is a heuristic rather than a rigorous optimization procedure, one cannot guarantee the optimality of its results if it produces a local minimum. Therefore, the best solution found in a series of runs will be called "maximal." A single pattern of width 18 residues was sought in the data of 30 HTH proteins shown in Fig. 1A. Solid lines show the course of three independent runs with different random seeds. Evolution



Lawrence CE, Altschul SF, Boguski MS, Liu JS, Neuvald AF, Wootton JC. (1993) Detecting subtle local protein homologies using a Gibbs sampling strategy for multiple alignment. *Science* 262:208-214.

Gibbs: ccmbweb.ccv.brown.edu/gibbs/gibbs.html

The Gibbs Motif Sampler Homepage



Welcome to the Gibbs Motif Sampler Homepage.

The Gibbs Motif Sampler will allow you to identify motifs, conserved regions, in DNA or protein sequences. This software was developed by Eric C. Rouchka and Bill Thompson based on work by C. E. Lawrence, J. S. Liu, L. A. McCue, A. F. Neuwald, L. A. Newberg and others (References).

→ Gibbs version 3.1 source and binaries for Linux, MS Windows (using Cygwin), Solaris, Solaris.x86 and MAC OS-X are available [here](#).

Gibbs is described in:

- Thompson WA, Newberg LA, Conlan S, McCue LA, and Lawrence CE. (2007) The Gibbs Centroid Sampler. *Nucleic Acids Res.* PubMed: 17483517, doi: 10.1093/nar/gpm265.
- Newberg LA, Thompson WA, Conlan S, Smith TM, McCue LA, and Lawrence CE. (2007) A phylogenetic Gibbs sampler that yields centroid solutions for *cis* regulatory site prediction. *Bioinformatics*. PubMed: 17488758, doi: 10.1093/bioinformatics/btm241.
- Thompson W, Rouchka EC, and Lawrence CE. (2003) Gibbs Recursive Sampler: finding transcription factor binding sites. *Nucleic Acids Res.* 31(13):3580-3585. PubMed: 12824370, doi: 10.1093/nar/gkx608.
- Thompson W, Palumbo MJ, Wasserman WW, Liu JS, and Lawrence CE. (2004) Decoding human regulatory circuits. *Genome Res.* 14(10A):1967-1974. PubMed: 15466295, doi: 10.1101/gr.2589004.
- Supplementary data for these papers are available [here](#).

fasta.bioch.virginia.edu/biol4230

53

GIBBS sampler of CRP sites: (1000 iterations)

```
16 columns
Num Motifs: 17
  1,   1      61 actgt TTTTTTGATCGTTTCACAAAA atgga      82  0.93 F CE1CG
  2,   1      55 attga TTATTGACGGCGTCACACTT tgcta      76  0.99 F ECOARABOP
  3,   1      76 ttaat AACCTGTGAGCATGGTCATATTT ttatc      97  0.97 F ECOBGLR1
  4,   1      63 tgcac GTATGCAAAGGACGTCACATTA cccgt      84  0.99 F ECOCRP
  5,   1      50 cagca AGGTCTTAAATTGATCACGTTT tagac      71  0.89 F ECOCYPA
  6,   1      7 gtgaa TTATTGAAACCAGATCGCATTAA cagtg      28  1.00 F ECODEOP2
  7,   1      42 tccac TAATTATTCCATGTGCACACTT ttcgc      63  0.72 F ECOGALE
  8,   1      20 ttttg TTATCTGCAATTCACTGACAAAA cgtga      41  0.74 F ECOILVBPR
  9,   1      9 gcaat TAATGTGAGTTAGCTCACTCAT taggc      30  1.00 F ECOLAC
 10,   1     14 gccaa TTCTGTAACAGAGATCACACAA agcga      35  1.00 F ECOMALBA
 11,   1     61 aggaa TTTCGTGATGTTGCTGCAAAA atcgt      82  0.93 F ECO2MALB
 12,   1     41 ttggg AATTGTGACACAGTGCAAATTC agaca      62  0.97 F ECOMALT
 13,   1     48 ttcat ATGCCCTGACGGAGTTCACACTT gtaag      69  0.98 F ECOOMPA
 14,   1     71 cgaac GATTGTGATTGCGATTCACATTT aaaca      92  1.00 F ECOTNAA
 15,   1     17 gaaat TGTTGTGATGTGGTTAACCCAA ttaga      38  0.50 F ECOUXU1
 16,   1     53 statg CGGTCTGAAATACCCACAGAT gcgta      74  0.91 F PBR322
 18,   1     78 agtta ATTGTGAGTGGTCGGCACATAT cctgt      99  1.00 F TDC
***** *****
```

fasta.bioch.virginia.edu/biol4230

54

GIBBS sampler of CRP sites: (5000 iterations)

16 columns

Num Motifs: 19

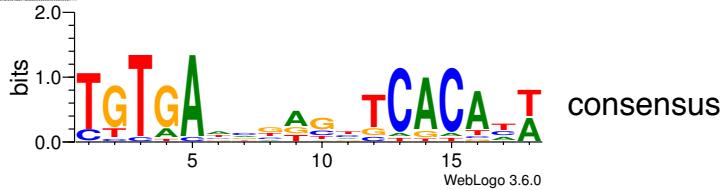
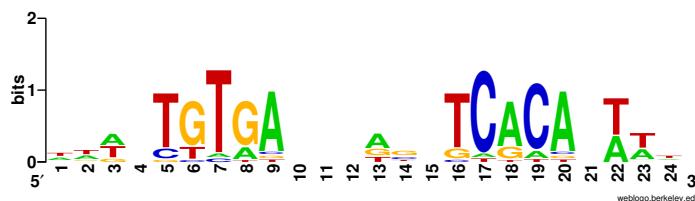
1,	1	60	gactg	TTTTTTGATCGTTTCACAAAAAA	tggaa	83	0.94	F	CE1CG
2,	1	54	cattg	ATTATTGCACGGCGTCACACTT	gttat	77	1.00	F	ECOARABOP
3,	1	75	gttaa	TAACTGTGAGCATGGTCATATT	tatca	98	0.98	F	ECOBGLR1
4,	1	62	ctgca	TGTATGCAAAGGACGTCACATT	cgtgc	85	0.98	F	ECOCRP
5,	1	49	tcagc	AAGGTGTTAAATTGATCACGTTT	agacc	72	0.94	F	ECOCYA
6,	1	6	agtga	ATTATTGAAACCAGATCGCATT	agtga	29	0.95	F	ECODEOP2
6,	2	59	tccct	TAATTGTGATGTCATCGAAGTGT	gttgc	82	0.56	F	ECODEOP2
7,	1	41	tccca	CTAATTATTCCATGTCACACTT	tcgca	64	0.45	F	ECOGALE
8,	1	38	agtac	AAAACGTGATCAACCCCTCAATT	tccct	61	0.73	F	ECOILVBPR
9,	1	8	cgcac	TTAACGTGAGTTAGCTCACTCATT	aggca	31	1.00	F	ECOLAC
10,	1	13	cgcac	ATTCTGTAACAGAGATCACACAAA	gcgcac	36	1.00	F	ECOMALBA
11,	1	60	aaggc	ATTTCGTGATGTTGCTGCAAAAAA	tgcgtg	83	0.97	F	ECO2MALB
12,	1	40	atttg	GAATTGTGACACAGTCACATTCA	gacac	63	0.94	F	ECOMALT
13,	1	47	tttca	TATGCCTGACGGAGTTCACACTTG	taagt	70	0.95	F	ECOOMP
14,	1	70	ccgaa	CGATTGTGATTGATTCACATT	aacaa	93	1.00	F	ECOTNAA
15,	1	16	tgaaa	TTGTTGTGATGTTAACCAAT	tagaa	39	0.55	F	ECOUXU1
16,	1	52	catat	GCGCTGTGAAATACCGCACAGATG	cgtaa	75	0.88	F	PBR322
17,	1	4	ctg	TGACGGAAGATCACTTCGAGAAT	aaata	27	0.05	F	TRN9CAT
18,	1	77	aagtt	AATTGTGAGTGGTCGCCACATATC	ctgtt	100	1.00	F	TDC

*** * * * * ***

fasta.bioch.virginia.edu/biol4230

55

GIBBS sampler of CRP sites: (5000 iterations)



fasta.bioch.virginia.edu/biol4230

56

Finding consensus regions in unaligned sequences

- Some introduction: regulation of transcription
- Looking for functional sites: promoters, regulatory elements, modification sites
- Products of convergent, not divergent evolution
- Weak spacing constraints
- Usually represented as a consensus sequence
- If alignment is given, consensus is obvious
- If consensus is given, alignment is obvious
- Search for consensus and alignment together
- **consensus, meme, gibbs**