

MCMC for Sequence Motif Search

Jianlin Cheng, PhD

Department of Computer Science
University of Missouri, Columbia



Fall, 2013

Free for academic use. Copyright @ Jianlin Cheng & original sources for some materials

Sequence Motif Search Problem

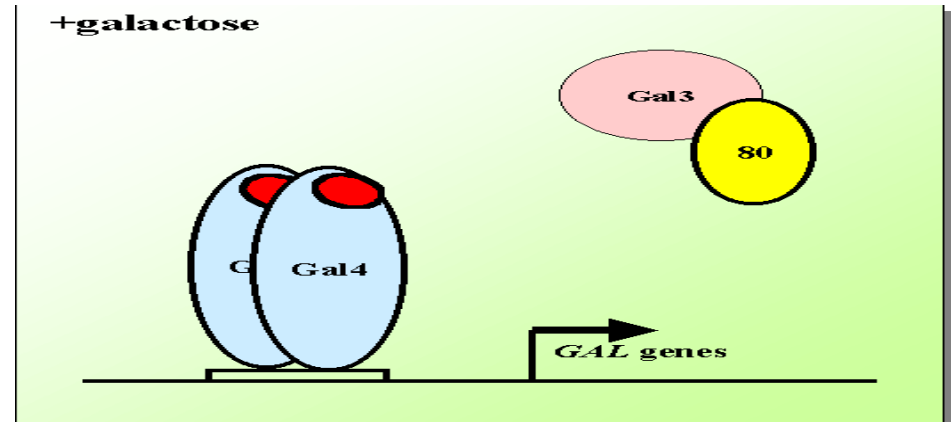
- Find a set of sub-sequences of multiple sequences whose alignment has maximum alignment score (highest similarity).
- It is NP-hard.
- Biologically find highly conserved regions (motifs) of related genes or a protein family

A Motif Example

```
0      5      10     15     20     25     30     35     40     45
TCTCATCCGGTGGGAATCACTGCCGCATTTGGAGCATAAACAATGGGGGGG
TACGAAGGACAAACACTTTAGAGGTAATGGAAACACAACCGGCGCATAAA
ATACAAACGAAAGCGAGAAGCTCGCAGAAGCATGGGAGTGTAATAAGTG
GGCGCCTCATTCTCGGTTATAAGCCAAAACCTTGTCGAGGCAACTGTCA
TCAAATGATGCTAGCCGTCGGAATCTGGCGAGTGCATAAAAAGAGTCAAC
```

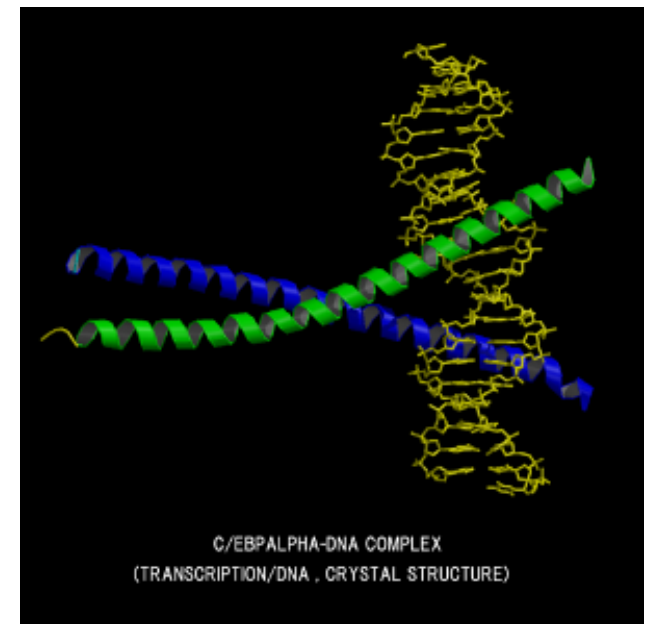
Examples: Transcription Factors

- yeast: Gal4
- drosophila
- mammal



1: actcgtcggggcggtacgtacgtaacgtacgta**CGGACA**ACTGTTGACCG
2: cggagcactgttgagcgacaagta**CGGAGCA**CTGTTGAGCGgtacgtac
3: ccccgtagg**CGGCGCA**CTCTCGCCCGggcgtacgtacgtaacgtacgta
4: agggcgcggtacgtaccgtcgacgtcg**CGCGCCGCA**CTGCTCCGacgct

Kathrina Kechris, 2005



Motif Model

Data: Upstream sequences from co-regulated/co-expressed genes.

Assumption: Binding site occurs in most sequences

1: actcgtcggggcggtacgtacgtaacgtacgtacggacaactggttgaccg
 2: cggagcactggttgagcgcacaagtacggagcactggttgagcgggtacgtac
 3: ccccgtagggcggcgcaactctcgcccgggcggtacgtacgtaacgtacgta
 4: agggcgcggtacgctaccgctcgacgctcgcgcgccgcactactccacact

Goals: 1) Estimate motif
 2) Predict motif locations



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
A	0	0	0	$\frac{3}{4}$	0	$\frac{1}{4}$	$\frac{1}{2}$	0	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$	0	0	0
C	$\frac{4}{4}$	0	0	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{4}$	$\frac{1}{2}$	0	$\frac{1}{4}$	0	$\frac{1}{4}$	0	$\frac{1}{4}$	$\frac{3}{4}$	$\frac{4}{4}$	0
G	0	$\frac{4}{4}$	$\frac{4}{4}$	0	$\frac{3}{4}$	$\frac{1}{2}$	$\frac{1}{4}$	0	0	$\frac{3}{4}$	0	0	$\frac{3}{4}$	0	$\frac{1}{4}$	0	$\frac{4}{4}$
T	0	0	0	0	0	0	0	$\frac{1}{2}$	$\frac{3}{4}$	0	$\frac{3}{4}$	$\frac{1}{2}$	0	$\frac{1}{4}$	0	0	0



1: actcgtcggggcggtacgtacgtaacgtacgta**CGGACAAC**TGTTGACCG
 2: cggagcactggttgagcgcacaagta**CGGAGCACT**GTTGAGCGgtacgtac
 3: ccccgtagg**CGGCGCACTCTCGCCC**Gggcggtacgtacgtaacgtacgta
 4: agggcgcggtacgctaccgctcgacgctcg**CGCGCCGCACTGCTCCG**acgct

Initialization of Locations (variables)

actcgtcggggcggtacgtacgttaa|cgtacgta**CGGACA**ACTGTTGACCG
cggagcactgttgagcgacaagta**CGGAGCA**CTGTTGAGCGgtacgtac
ccccgtagg**CGGCGCA**CTCTCGCCCGggcgta|cgtacgtaacgtacgtac
agggcgcgta|cgtacgtaccgtcgacgtcg**CGCGCCGCA**CTGCTCCGacgct

Motif Size: 17

Construct a probability matrix (profile)

[illegible]

Conditional Probability of a Position

1

actcgtcggggcggtacgtacgtaacgtacgta**CGGACA**ACTGTTGACCG
cggagcactgttgagcgcacaagta**CGGAGCA**CTGTTGAGCGgtacgtac
ccccgtagg**CGGCGCACTCTCGCCCCG**ggcgtacgtacgtaacgtacgtacgta
agggcgcgtacgtctaccgtcgacgtcg**CGCGCCGCACTGCTCCG**acgct

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
A	1/4	1/4											
C	2/4	1/4															
G	0	2/4															
T	1/4	0															

*Find the best position in each sequence
That maximize product of probability*

$$\text{Prob (posi_i)} = 2/4 * 2/4 * \dots$$

Position Resampling

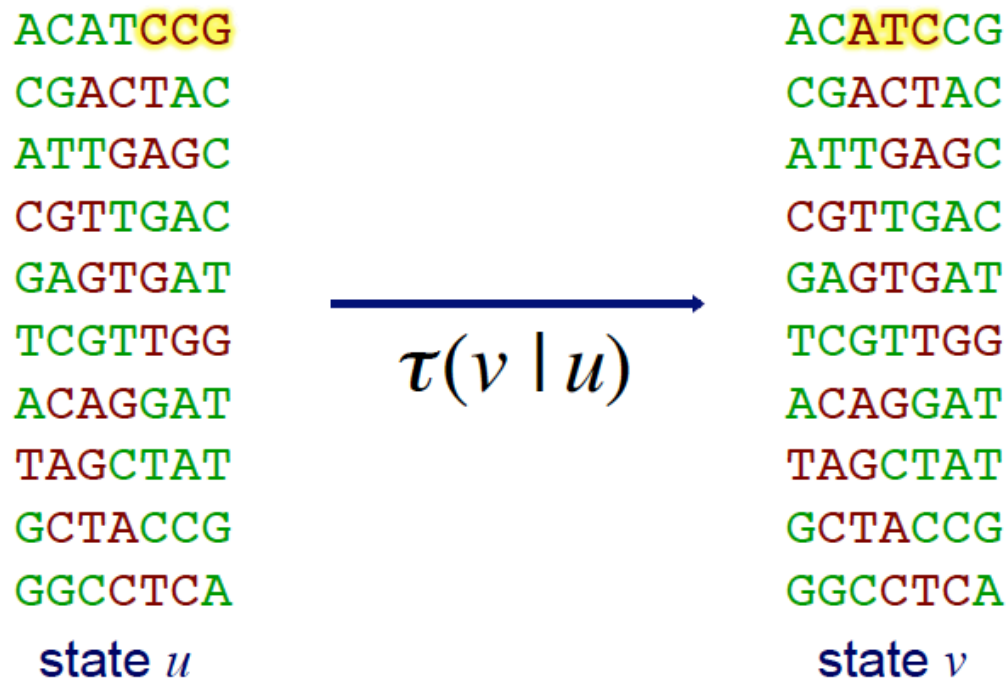
actcgtcggggcggtacgtacgtaacgtacgtac**CGGACA**ACTGTTGACCG
cggagcactgttgagcgacaagta**CGGAGC**ACTGTTGAGCGgtacgtac
ccccgtagg**CGGCGC**ACTCTCGCCCGggcggtacgtacgtaacgtacgtac
agggcgcggtacgtacgtaccgtcgacgtcg**CGCGCCGC**ACTGCTCCGacgct

Construct a probability matrix (profile) from new positions

[illegible]

Markov Chain Monte Carlo (MCMC)

- we can view the motif finding approach in terms of a Markov chain
- each state represents a configuration of the starting positions (a_i values for a set of random variables $A_1 \dots A_n$)
- transitions correspond to changing selected starting positions (and hence moving to a new state)



Markov Chain Monte Carlo

- for the motif-finding task, the number of states is enormous
- key idea: construct Markov chain with stationary distribution equal to distribution of interest; use sampling to find most probable states
- detailed balance:

$$P(u)\tau(v | u) = P(v)\tau(u | v)$$

probability of state u probability of transition $u \rightarrow v$

- when detailed balance holds:

$$\frac{1}{N} \lim_{N \rightarrow \infty} \text{count}(u) = P(u)$$

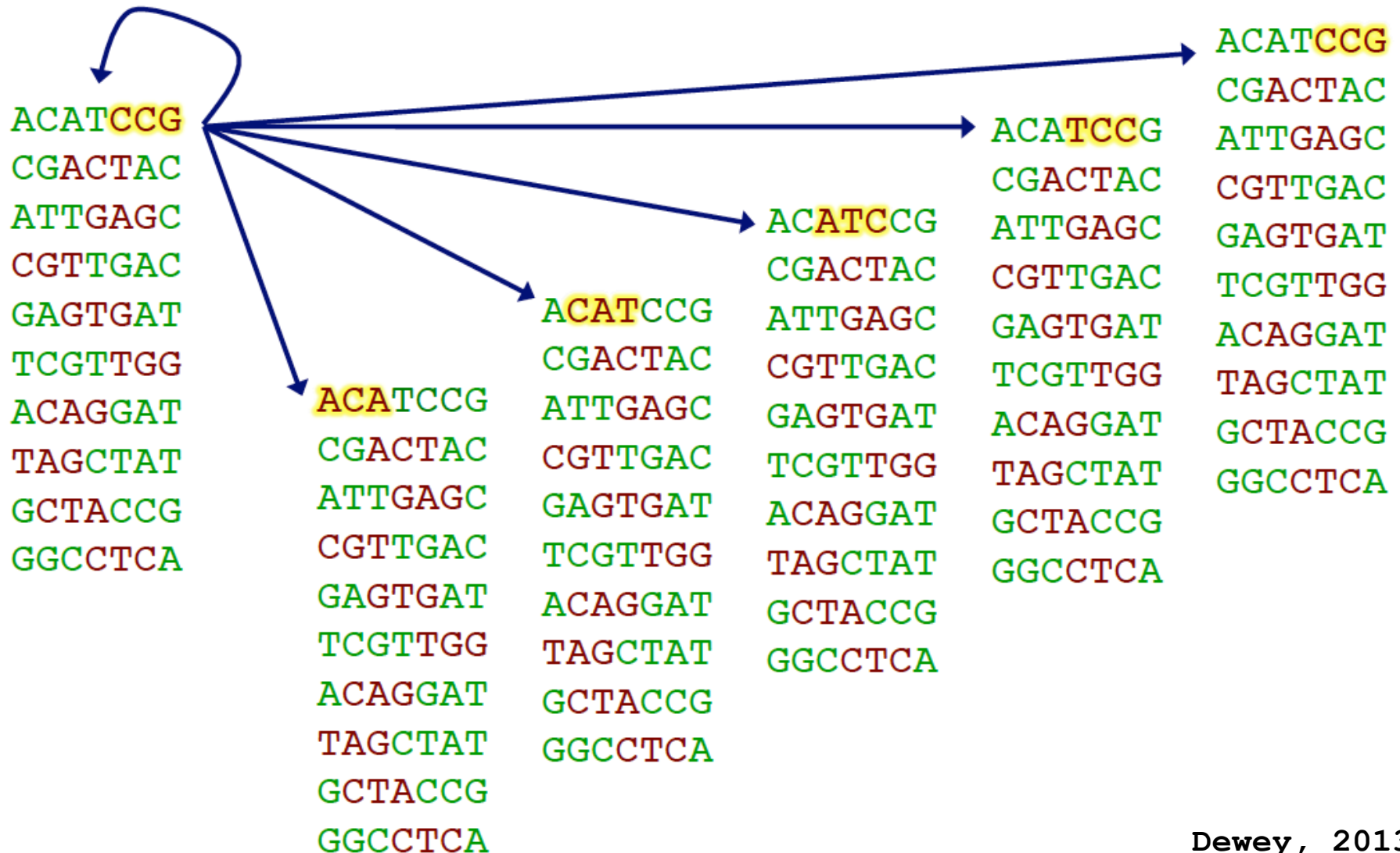
MCMC with Gibbs Sampling

Gibbs sampling is a special case of MCMC in which

- Markov chain transitions involve changing one variable at a time
- transition probability is conditional probability of the changed variable given all others
- i.e. we sample the joint distribution of a set of random variables $P(A_1 \dots A_n)$ by iteratively sampling from $P(A_i \mid A_1 \dots A_{i-1}, A_{i+1} \dots A_n)$

Gibbs Sampling Approach

- possible state transitions when first sequence is selected



Project I

- **Objective:** Design and develop a MCMC method to find a motif from a group of DNA sequence
- **Other tools and data:**
<http://biowhat.ucsd.edu/homer/motif/> ; Download the package to find the data in one of sub directories?
- MEME tool: <http://meme.nbcr.net/meme/>
- Motif Visualization tool (weblogo):

`http://weblogo.berkeley.edu/logo.cgi`

Project Discussion

Group 1: Kishore, Sairam, Abhimanyu, Maruthi, Shravya, Rajni

Group 2: Xiaokai, Xiao, Rui, Linfei, Pei, Zhiluo

Group 3:

Group 4

Others: participating in discussion

Items to Discuss

Group Discussion (25 minutes)

- Problem definition
- Algorithm
- Implementation
- Evaluation
- Visualization
- Task Assignment
- Select Coordinator

• Informal Presentation (5 minutes per group)

• To do:

Present your plan (PPT): 15 minutes per group
(Wednesday)

Create group accounts on server (email notification)

Gibbs Sampling Algorithm for Motif Finding

given: length parameter W , training set of sequences

choose random positions for a

do

pick a sequence X_i

estimate p given current motif positions a

(using all sequences but X_i) (predictive update step)

sample a new motif position a_i for X_i (sampling step)

until convergence

return: p, a

Gibbs Sampling Algorithm II

Assumption: size of motif is fixed

Initialization:

Make an initial guess of the motif locations and compute a probability matrix

Repeat:

Select one sequence randomly

Use the matrix to evaluate the probabilities of all positions in the sequence (product of probability)

Select (or sample) a position in the sequence according to their probability

Recalculate the motif probability matrix with the new position

Until matrix converges.

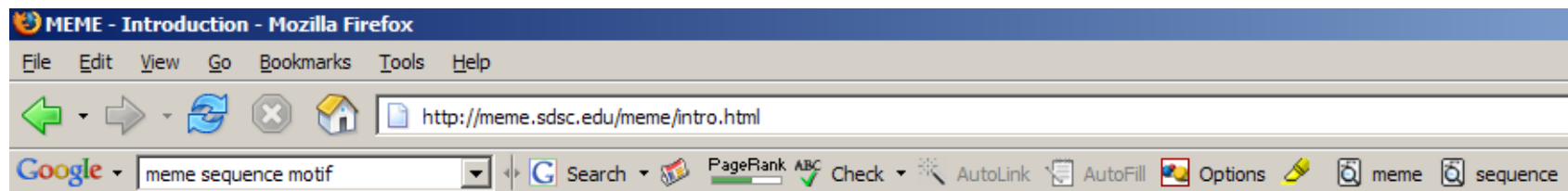
Sample a position according to probability

actcgtcggggcggtacgtacgtaacgtacgtaⁱ**CGGACA**ACTGTTGACCG
 cgggagcactgttgagcgacaagta**CGGAGCACTGTTGAGCG**gtacgtac
 ccccgtagg**CGGCGCACTCTCGCCCCG**ggcggtacgtacgtaacgtacgtacgta
 agggcgcggtacgtaccggtcgacgtcg**CGCGCCGCACTGCTCCG**acgct

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17
A	1/4	1/4											
C	2/4	1/4															
G	0	2/4															
T	1/4	0															

Compute $P_i = 2/4 * 2/4 * \dots 1 \leq i \leq n$
 Select a position according to its
 Normalized probability.

Sample probability of $i = \frac{p_i}{\sum_{i=1}^n p_i}$



Menu

- + Submit A Job
- + Resources
- + Alternate Servers
- + Other Tools



MEME

Multiple Em for Motif Elicitation



MAST

Motif Alignment & Search Tool

THE MEME/MAST SYSTEM

Motif Discovery and Search

Version 3.5.3

The MEME/MAST system allows you to

1. [discover](#) motifs (highly conserved regions) in groups of related DNA or protein sequences using [MEME](#) and,
2. [search](#) sequence databases using motifs using [MAST](#).

- ◆ The MEME/MAST system was developed by [Timothy Bailey](#), [Charles Elkan](#), and [Bill Noble](#) at the [UCSD Computer Science and Engineering](#) department with input from [Michael Gribskov](#) at [Purdue University](#).
- ◆ MEME and MAST are described in detail in the [papers](#) available here.
- ◆ Answers to Frequently Asked Questions about MEME and MAST are given in the [GENERAL FAQ](#).
- ◆ Visit the [MEME user forum](#) for online discussions with the MEME support team members and other MEME users.
- ◆ You can see [sample MEME output](#) or [sample MAST output](#).
- ◆ Differences between the current release of the MEME/MAST system and earlier releases are described in the [release notes](#).
- ◆ You can [download](#) the MEME/MAST software and install it on your own computer. This will allow you to use many features that are not available with the interactive versions of MEME and MAST.
- ◆ [Meta-MEME](#) combines motif models from MEME into a hidden Markov model framework for use in searching sequence databases.
- ◆ MEME and MAST are copyrighted software and can be [licensed](#) for commercial use.

Menu

- Submit A Job
- Resources
- Alternate Servers
- Other Tools



MEME

Multiple Em for Motif Elicitation

Version 3.5.3

Use this form to submit DNA or protein sequences to MEME. MEME will analyze your sequences for similarities among them and produce a description (**motif**) for each pattern it discovers. Your results will be sent to you by e-mail.

Data Submission Form

Required

Your **e-mail address**:

Re-enter e-mail address:

Please enter the **sequences** which you believe share one or more motifs. The sequences may contain no more than **60,000 characters** total in any of a large number of **formats**.

Enter the **name of a file** containing the sequences here:

or

The **actual sequences** here (Sample Input Sequences):

```
>seq1
actcgtcggggcgtacgtacgtaacgtacgtaCGGACAACTGTTGA
CCG
>seq2
cggagcactgttgagcgacaagtaCGGAGCACTGTTGAGCGgt
```

How do you think the occurrences of a single motif are **distributed** among the sequences?

- ☒ **One** per sequence
☐ **Zero or one** per sequence
☐ **Any number** of repetitions

MEME will find the optimum **width** of each motif within the limits you specify here:

 Minimum width (≥ 2) **Maximum** width (≤ 300) **Maximum number of motifs** to find

Optional

Description of your sequences:

MEME will find the optimum **number of sites** for each motif within the limits you specify here:

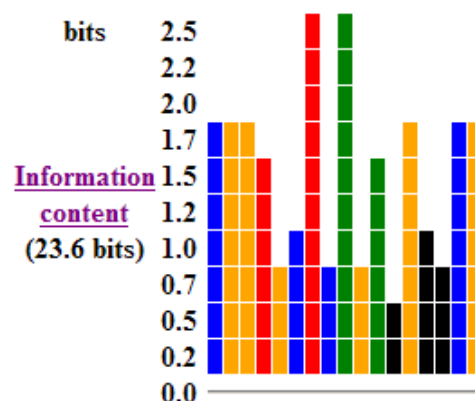
 Minimum sites (≥ 2) **Maximum** sites (≤ 300)

- ☐ **Text** output format
☐ **Shuffle** sequence letters

For DNA sequences only:

- ☐ Search given **strand** only
☐ Look for **palindromes** only

Simplified A : : : 8 : 3 a : : : : 5 : :
pos.-specific C a : : 3 3 8 : 8 : 3 3 3 : 5 5 a :
probability G : a a : 8 : : 3 : 8 : 3 a : 5 : a
matrix T : : : : : : : a : 8 5 : : : :



Multilevel CGGAGCACTGTTGACCG
consensus CCA G CCC CG
sequence G

NAME STRAND START P-VALUE

seq2 + 1 2.44e-10
 seq1 + 33 5.18e-09
 seq4 - 28 1.08e-07
 seq3 + 10 1.08e-07

SITES

CGGAGCACTGTTGAGCG ACAAGTACGG
 AACGTACGTA CGGACAACCTGTTGACCG
 AGCGT CGGAGCAGTGCGGCGCG CGACGTGCGAC
 CCCCCTAGG CGGCGCACTCTCGCCCG GGCGTACGTA

Motif 2 block diagrams

Gibbs Motif Sampler

<http://bayesweb.wadsworth.org/gibbs/gibbs.html>

The Gibbs Motif Sampler

(for DNA)

Show advanced options [How to enter data?](#)

Email Address:

Please enter the data sequence: ([FASTA](#) format) *

Prokaryotic Defaults

Sampler Mode: ☐ Site Sampler

No. of different motifs (patterns):

Motif Width(s):*

Eukaryotic Defaults

☐ Motif Sampler ☒ Recursive Sampler

Max sites per seq: (recursive sampler)

Est. total sites for each motif type:

Gibbs Motif Sampler

<http://bayesweb.wadsworth.org/gibbs/gibbs.html>

Email Address:

Please enter the data sequence: ([FASTA](#) format) *

```
>seq1
actcgtcggggcggtacgtacgtaacgtacgtacCGGACAACGTGTTGACCG
>seq2
cggagcactggttgagcgacaagtaCGGAGCACTGTTGAGCGgtacgtac
>seq3
ccccgtaggCGGCGCACTCTCGCCCGggcggtacgtacgtaacgtacgtac
>seq4
agggcgcggtacgtacaccgtcgacgtcgCGCGCCGCACTGCTCCGacgct
```

Prokaryotic Defaults

Sampler Mode:

☒ Site Sampler

No. of different motifs
(patterns):

Motif Width(s):*

Eukaryotic Defaults

☐ Motif Sampler

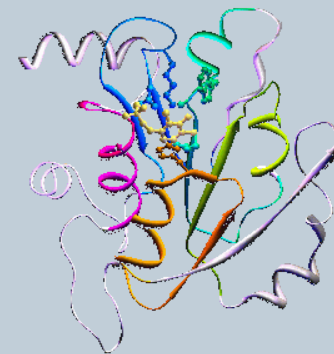
☐ Recursive Sampler

Max sites per seq:
(recursive sampler)

Est. total sites for each
motif type:

The Gibbs Motif Sampler

(for DNA)



[Browse the Gibbs Motif Sampler Manual](#)

Output of Gibbs Sampler

actcgtcggggcgtagctacgtaacgtacgtaCGGACAACGTGTTGACCG
cggagcactgttgagcgacaagtaCGGAGCACTGTTGAGCGgtacgtac
ccccgtaggCGGCGCACTCTCGCCCGggcgtagctacgtaacgtacgta
agggcgcgtagctaccgtcgacgtcgCGCGCCGCACTGCTCCGacgct

Motif probability model

Pos. #	a	t	c	g
1	0.014	0.013	0.949	0.024
2	0.014	0.013	0.023	0.950
3	0.014	0.013	0.023	0.950
4	0.755	0.013	0.209	0.024
5	0.014	0.013	0.209	0.765
6	0.199	0.013	0.764	0.024
7	0.940	0.013	0.023	0.024
8	0.014	0.013	0.764	0.209
9	0.014	0.939	0.023	0.024
10	0.014	0.013	0.209	0.765
11	0.014	0.754	0.209	0.024
12	0.014	0.568	0.209	0.209
13	0.014	0.013	0.023	0.950
14	0.570	0.013	0.394	0.024
15	0.014	0.013	0.394	0.579
16	0.014	0.013	0.949	0.024
17	0.014	0.013	0.023	0.950

Prob Matrix

Confidence

Motif

End pos

Start pos

Background probability model
0.225 0.189 0.279 0.306

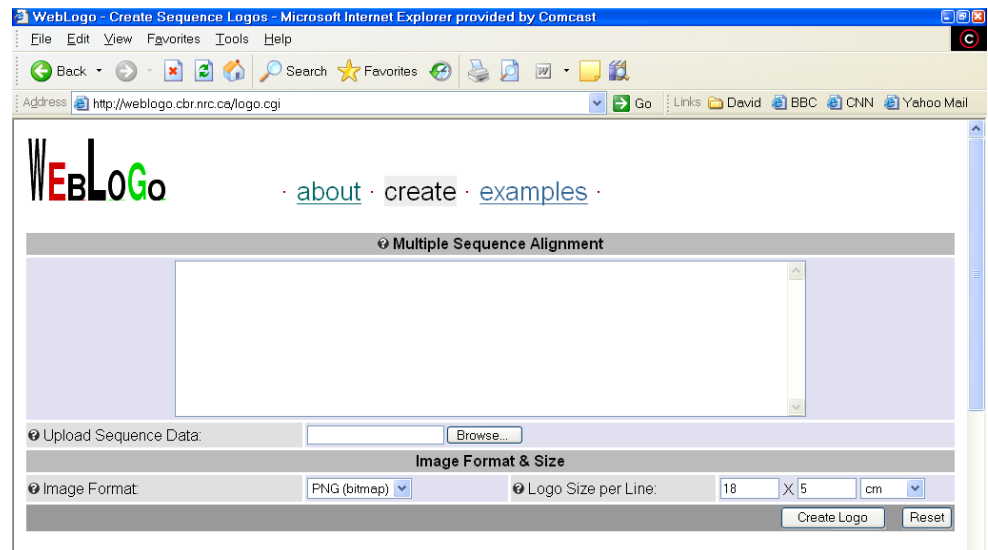
17 columns
Num Motifs: 5

1, 1	33	acgta	CGGACAAC	TGTTGACCG	49	1.00	F	seq1	
2, 1	1		CGGAGCACT	TGTTGAGCG	17	0.49	F	seq2	
2, 2	25	aagta	CGGAGCACT	TGTTGAGCG	gtacg	41	0.51	F	seq2
3, 1	10	gtagg	CGGCGCACT	CTCGCCCG	ggcgt	26	1.00	F	seq3
4, 1	44	agcgt	CGGAGCAGT	GCGGCGCG	cgacg	28	1.00	R	seq4



- Graphical representation of nucleotide base (or amino acid) conservation in a motif (or alignment)
- Information theory $2 + \sum_{b=\{A,C,G,T\}} p(b) \log_2 p(b)$
- Height of letters represents relative frequency of nucleotide bases

<http://weblogo.berkeley.edu/>



Kathrina Kechris, 2005

Entropy and Information

Visualization goals

- (1) The height of the position is proportional to the information contained at the position
- (2) The height of a letter is proportional to the probability of the letter appearing at the position

Two new concepts related to probability matrix:

Entropy

Information

- Entropy is a measure of uncertainty of a distribution $\sum_i -p_i \log_2 p_i$

	A	C	G	T
1	1/4	1/4	1/4	1/4
2	0	1	0	0
3	1/2	1/2	0	0
4				
⋮				

What is the entropy
Of positions 1,2,3?

- Information is the opposite of entropy. It measures the certainty of a distribution
- Information = maximum entropy – the entropy of a position (or distribution)

Maximum entropy for n characters is the Entropy when n characters are uniformly Distributed. $\log_2 n$

Info. Of pos 1 = $2 - 2 = 0$

Info. Of pos 2 = $2 - 0 = 2$

Info. Of pos 3 = $2 - 1 = 1$

Multiple Sequence Alignment

```
>seq1
CGGACAACCTGTTGACCG
>seq2
CGGAGCACTGTTGAGCG
>seq3
CGGCGCACTCTCGCCCG
>seq4
CGGCGCGCACTGCTCCG
```

Upload Sequence Data:

Image Format & Size

Image Format:

Logo Size per Line:

X

Advanced Logo Options

Sequence Type:

☐ amino acid ☐ DNA / RNA ☒ Automatic Detection

First Position Number:

Logo Range:

-

Small Sample Correction:

☒

Frequency Plot:

☐

