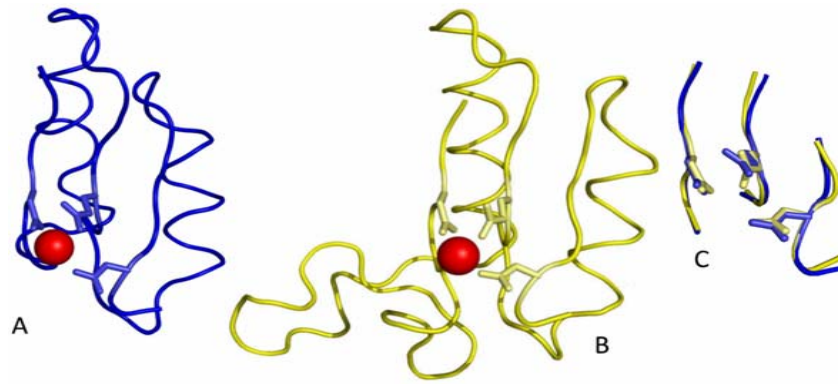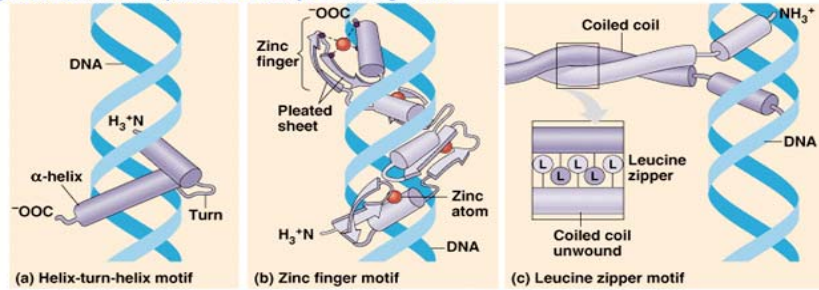# TOOLS FOR MOTIF
# AND PATTERN SEARCHING
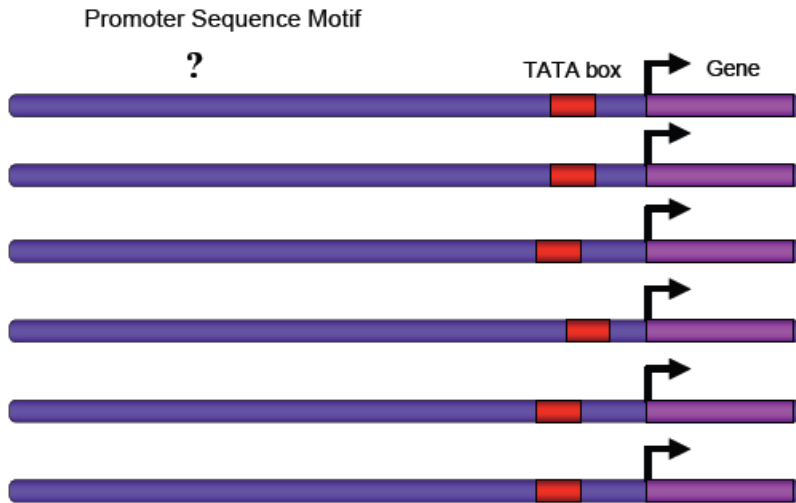
Prat Thiru

# OUTLINE

- What are motifs?
- Algorithms Used and Programs Available
- Workflow and Strategies
- MEME/MAST Demo (online and command line)

# Protein Motifs



protein motifs that help bind to DNA major & minor grooves

(a) Helix-turn-helix motif

(b) Zinc finger motif

(c) Leucine zipper motif

# DNA Motifs

# MEME Output

# Definitions

- <u>Motif</u>: Conserved regions of protein or DNA sequences

- <u>Pattern</u>: Qualitative description of a motif

  *eg. regular expression* C[AT]AAT[CG]X

- <u>Profile</u>: Quantitative description of a motif

  *eg. position weight matrix*

# Patterns

- Regular Expression Symbols
  - ➢ [ ] – OR eg. [GA] means G or A
  - ➢ { } – NOT eg. {P,V} means not P or V
  - ➢ ( ) – repeats eg. A(3) means AAA
  - ➢ X or N or "." – any

- Complex patterns representation difficult

- Loose frequency information
  eg. [AT] *vs* 20%A  80%T

# Profiles



Amino acid type *j* and gap penalties

Sequence profile position *k*

$M_{kj}$

$$M_{kj} = \log\left(\frac{p_{kj}}{p_j}\right)$$

- $p_{kj}$ = probability of amino acid j at position k in the profile
- $p_j$ = "background" probability of amino acid j in sequence

| A | .02 | .95 | .01 | .59 | .51 | .40 | .44 | .39 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|
| C | .09 | .02 | .03 | .13 | .20 | .11 | .05 | .10 |
| G | .10 | .01 | .00 | .15 | .13 | .12 | .08 | .09 |
| T | .79 | .03 | .96 | .13 | .16 | .37 | .43 | .42 |
| motif | T | A | T | A | A | [AT] | [AT] | [AT] |

# Sequence Logos

| | |
|---|---|
| HEM13 | CCCATTGTTCTC |
| HEM13 | TTTCTGGTTCTC |
| HEM13 | TCAATTGTTTAG |
| ANB1 | CTCATTGTTGTC |
| ANB1 | TCCATTGTTCTC |
| ANB1 | CCTATTGTTCTC |
| ANB1 | TCCATTGTTCGT |
| ROX1 | CCAATTGTTTTG |

**YCHATTGTTCTC**

| | |
|---|---|
| A | 002700000010 |
| C | 464100000505 |
| G | 000001800112 |
| T | 422087088261 |

# Algorithms

- Enumeration

- Probabilistic Optimization

- Deterministic Optimization

1. Identify motifs

2. Build a consensus

# Enumeration

- Exhaustive search: word counting method, count all $n$-mers and look for overrepresentation

- Less likely to get stuck in a *local* optimum

- Computationally expensive

  ➢ YMF
    http://wingless.cs.washington.edu/YMF/YMFWeb/YMFInput.pl

  ➢ Weeder
    http://159.149.109.9/weederaddons/locator.html

# Probabilistic Optimization

- Uses a Gibbs sampling approach
- One $n$-mer from each sequence is randomly picked to determine initial model.  In subsequent iterations, one sequence, $i$, is removed and the model is recalculated. Pick a new location of motif in sequence $i$ iterate until convergence
- Assumes most sequences will have the motif
- AlignAce
  http://atlas.med.harvard.edu/cgi-bin/alignace.pl
- Gibbs Motif Sampler

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

# Deterministic Optimization

- Based on expectation maximization (EM)

- EM: iteratively estimates the likelihood given the data that is present

  I.  **Expectation step:** Use current parameters (and observations) to reconstruct hidden structure

  II. **Maximization step:** Use that hidden structure (and observations) to re-estimate parameters

➢ MEME

http://meme.sdsc.edu

# Multiple EM for Motif Elicitation MEME

- Starting from a single site, EM alternates between assigning sites and updating motif model

- Performs a single iteration for each $n$-mer in target sequences, selects the best motif from this site and then iterates only that one to convergence

- Search space increases significantly with increasing number of sequence and/or sequence lengths

# Programs Available*

## Multi-purpose packages

### Motif Scanning

| | | | |
|---|---|---|---|
| **TAMO** | TAMO integrates several motif discovery programs. It includes support for motif scanning, scoring, evaluation of statistical significance, clustering, comparison, input/output, conversion between different motif representations, and visualization. http://fraenkel.mit.edu/webtamo/ | **Ahab** | The Ahab webserver allows users to scan for motifs in a set of sequences. Motifs may be user-specified or selected from a database of pre-defined matrices. http://gaspard.bio.nyu.edu/Ahab.html |
| **BEST** | BEST is a suite of four motif discovery tools integrated in a graphical user interface. BEST incorporates the BioOptimizer tool used to rank and improve the predictive power of the discovered motifs. http://webster.cs.uga.edu/~che/BEST/ | **Clover** | Clover identifies overrepresented motifs in a set of sequences, based on a pre-compiled library of motif matrices. http://zlab.bu.edu/clover/ |
| **TOUCAN2** | TOUCAN2 provides an interface to the Ensembl and EMBL databases of sequence and annotation. It incorporates tools for sequence alignment, motif discovery, and scanning. http://homes.esat.kuleuven.be/~saerts/software/toucan.php | **MAST** | MAST allows users to scan sequence databases for matches to motifs. It produces detailed annotations and figures for matches in the input sequences. http://meme.sdsc.edu/meme/intro.html |
| **Expander** | Expander is a tool for analyzing expression data. It can cluster genes, identify over-represented functional categories in clusters, and scan corresponding promoter regions for motifs. http://www.cs.tau.ac.il/~rshamir/expander/ | **Monkey** | Monkey analyzes multiple sequence alignments to identify evolutionarily conserved matches to a motif. http://rana.lbl.gov/~alan/Monkey.htm |
| **MDScan BioProspector Compare-Prospector** | MDScan uses ChIP-chip enrichment ratio data to help the motif search. BioProspector is a Gibb's sampling program. CompareProspector incorporates comparative genomics, biasing the search to regions of high conservation. http://seqmotifs.stanford.edu | **cisRED** **ORegAnno** | cisRED is a database of conserved motifs and motif patterns obtained by genome scale motif discovery. ORegAnno is a database of regulatory sites curated from the scientific literature. http://www.cisred.org/ http://www.oreganno.org/ |
| **Consensus PhyloCon** | The Consensus program finds motifs in a set of unaligned sequences. PhyloCon builds on this framework by modeling conservation across orthologous genes from multiple species. http://ural.wustl.edu/ | **UCSC Genome Browser** | Online repository of genomic sequence, multiple sequence alignments, and annotation data. The browser includes tracks for identifying conserved transcription factor binding sites. http://genome.ucsc.edu/ |
| **Weeder** | An enumerative motif discovery program that performed well in a recent comparative analysis of fourteen algorithms. http://www.pesolelab.it/ | **ENSEMBL** | Another online genomic sequence repository. Includes online tools for data mining as well as BLAST searches. http://www.ensembl.org/index.html |
| **MEME** | The popular EM-based motif discovery program. Part of the MEME/MAST system for motif discovery and search. http://meme.sdsc.edu/meme/intro.html | **TRANSFAC** | Commercial database of transcription factors, binding sites, and motifs. Includes several tools for motif scanning in sequence. http://www.gene-regulation.com/ |
| **AlignACE** | A Gibbs sampling algorithm that can identify multiple motifs in a sequence set using an iterative masking procedure. http://atlas.med.harvard.edu/ | **JASPAR** | Curated public database of transcription factor binding specificities represented as PWMs. http://jaspar.cgb.ki.se/ |

Motif Discovery Programs            Databases

*incomplete list

# Programs Available: EMBOSS Motif Searching

http://iona.wi.mit.edu/bio/tools/emboss/

- wordcount: Counts words of a specified size in a DNA sequence

- prophecy: Creates matrices/profiles from multiple alignments

- profit: Scan a sequence or database with a matrix or profile
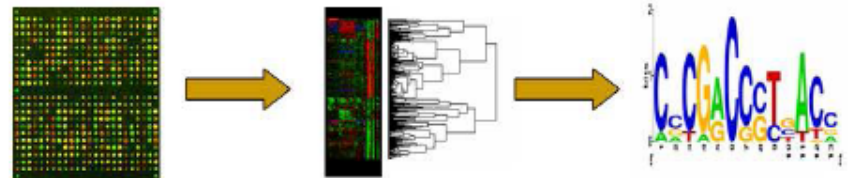
# Programs Available: EMBOSS Pattern Searching

http://iona.wi.mit.edu/bio/tools/emboss/

- fuzznuc: Nucleic acid pattern search

- fuzzpro: Protein pattern search

# Programs Available: Other

- Allegro (Expression)

http://acgt.cs.tau.ac.il/allegro/
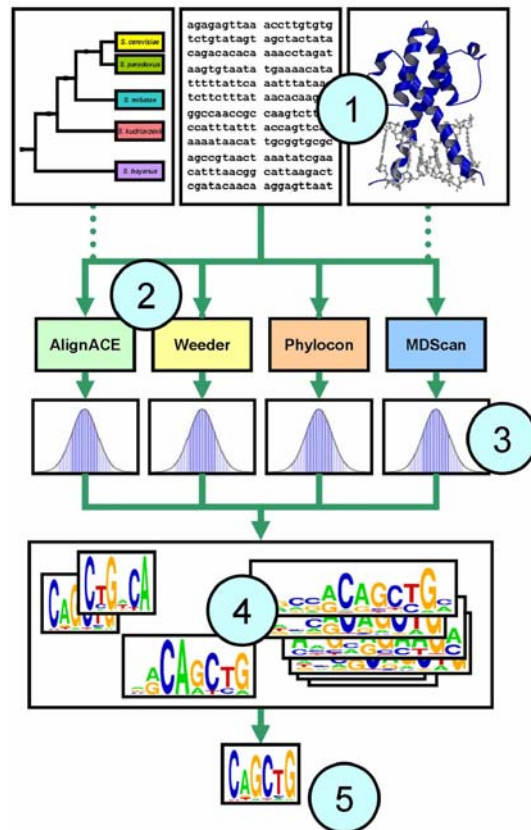


- CisGenome (ChIP-Seq)

http://www.biostat.jhsph.edu/~hji/cisgenome

# Workflow and Strategies



**Assemble input data.** Results may be improved by restricting the input to high-confidence sequences. Some algorithms achieve improved performance by using phylogenetic conservation information from orthologous sequences or information about protein DNA-binding domains. (1)

**Choose several motif discovery programs for the analysis.** For recommended programs see Figure 3. (2)

**Test the statistical significance of the resulting motifs.** Use control calculations to estimate the empirical distribution of scores produced by each program on random data. (3)

**Clustering and post-processing the motifs.** Motif discovery analyses often produce many similar motifs, which may be combined using clustering. Phylogenetic conservation information may be used to filter out statistically significant, but non-conserved motifs that are more likely to correspond to spurious sequence patterns. (4)

**Interpretation of motifs.** Algorithms exist for linking motifs to transcription factors and for combining motif discovery with expression data. (5)

**Fraenkel, E., et al.** *Practical Strategies for Discovering Regulatory DNA Sequence Motifs* **PLoS Computational Biology 2:201-210 (2006)**

# Further Reading

- **Practical Strategies for Discovering Regulatory DNA Sequence Motifs**

http://www.ploscompbiol.org/article/info:doi/10.1371/journal.pcbi.0020036

- **How does DNA sequence motif discovery work?**

http://www.nature.com/nbt/journal/v24/n8/full/nbt0806-959.html

- **MEME**

Bailey, T. L. et al. Nucl. Acids Res. 2006 34:W369-W373; doi:10.1093/nar/gkl198

# MEME/MAST Demo

## http://meme.sdsc.edu

# MEME/MAST Demo
## MEME Suite

# MEME/MAST Demo
## Command Line (on tak)

Usage: MEME (Find ungapped motifs in unaligned sequences)

eg. meme sample.fa -dna -maxw 10 -nmotifs 5 -mod zoops -pal -maxsize 1000000 -o sample_meme

meme <dataset> [optional arguments]

- <dataset> file containing sequences in FASTA format
- [-text] output in text format (default is HTML)
- [-dna] sequences use DNA alphabet
- [-protein] sequences use protein alphabet
- [-mod oops|zoops|anr] distribution of motifs
- [-nmotifs <nmotifs>] maximum number of motifs to find
- [-evt <ev>] stop if motif E-value greater than <evt>
- [-minw <minw>] minimum motif width
- [-maxw <maxw>] maximum motif width

*For complete list of options enter "meme" at the command prompt*

# MEME/MAST Demo
## Command Line (on tak)

**Usage: MAST (Searches a sequence database for occurrences of known motifs )**

**eg. mast motifs.txt -d data.fa**

**mast <mfile> <database> [ optional arguments ... ]**

- **<mfile>   file containing motifs to use; may be a MEME output file or similar file**
- **[-d <database> | -stdin] search sequences in <database> with motifs**

*For complete list of options enter "mast" at the command prompt*