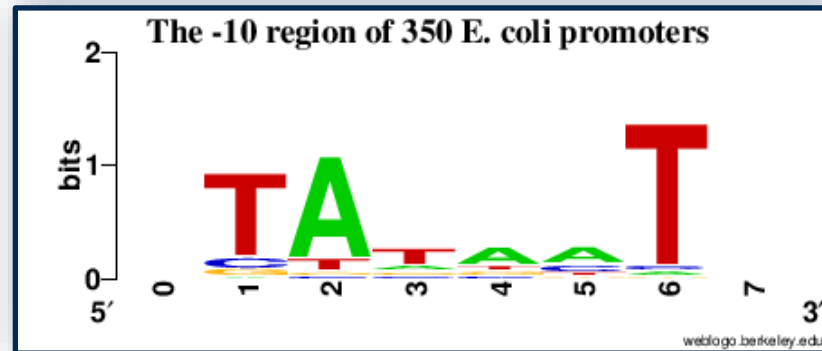# Motif Discovery in DNA and Protein Sequences

Word based and Expectation Maximization based Methods

# Multiple EM for Motif Elicitation

- MEME
  - Is a tool for discovering motifs in a group of related nucleotide or peptide sequences.
  - A MEME motif is a sequence pattern that occurs repeatedly in one or more sequences in the input group.
  - Can be used to discover novel patterns, as it bases its discoveries only on the input sequences, not on any prior knowledge (such as databases of known motifs).
  - MEME motifs allow errors (mutations) at any position in the pattern, but individual MEME motifs may <u>not</u> contain gaps (insertions or deletions).
  - Splits patterns that contain gaps into multiple motifs.
- Motifs may appear in any order, multiple times or not at all in any given sequence.
  - Input
    - a set of unaligned sequences of the same type (peptide or nucleotide) also called training set.
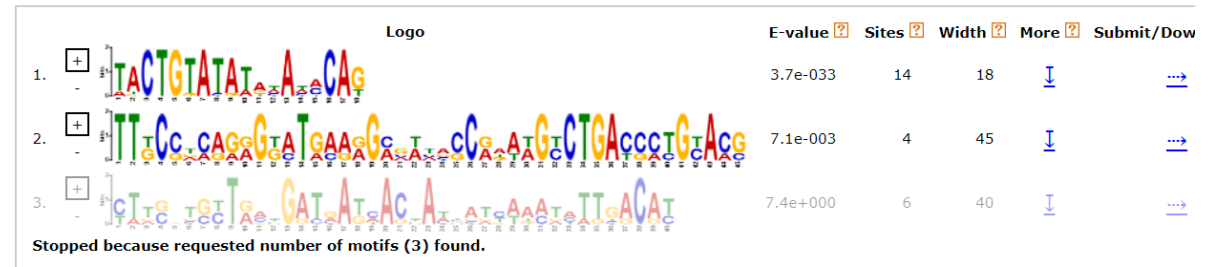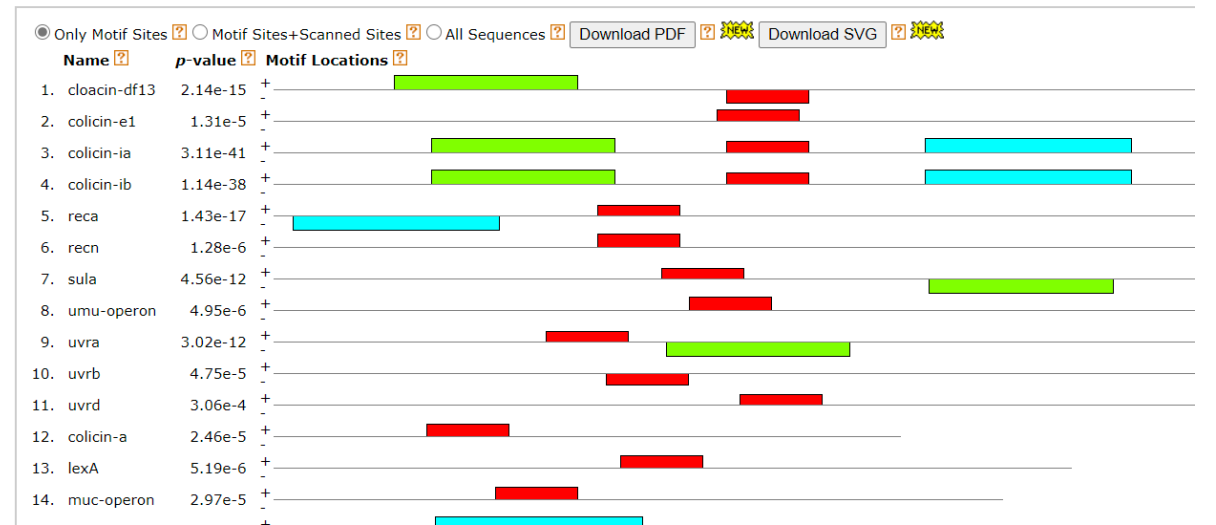
# Multiple EM for Motif Elicitation



- For each motif, MEME reports:

  o Discovered motifs;

  o Motif locations;

# Discovering  motifs in a set of Peptides sequences

meme filename.fasta -nmotifs 4 -o run1

meme filename.fasta -nmotifs 4 -minw 4 -maxw 10 -o run2

| oops | One Occurrence Per Sequence |
|------|------------------------------|
| zoops | Zero or One Occurrence Per Sequence MEME |
| anr | Any Number of Repetitions (This option can also be used to discover repeats within a single sequence) |

Allowing repeated motifs

meme filename.fasta -nmotifs 4 -minw 4 -maxw 10 -mod anr -o run.anr

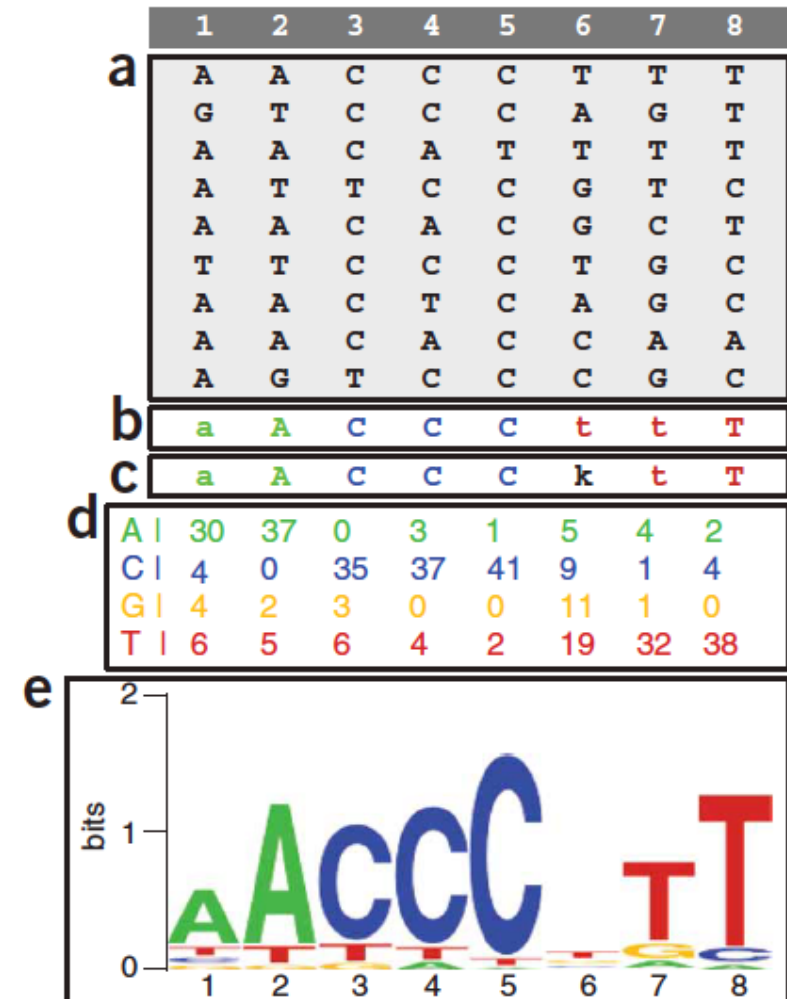http://meme-suite.org/doc/meme.html?man_type=web

# Word based Motif Discovery with RSAT

- The regulatory sequences analysis tools (RSAT) are a suite of specialized programs for detecting regulatory elements.

  o available at: http://rsat.sb-roscoff.fr/

- The oligo-analysis tool uses an exhaustive approach by scanning all the oligomers of a given size (min and max length can be defined) counting the respective occurrences in a set of sequences. It then uses statistical analysis to detect overrepresented and significant oligonucleotides.

# Motif Representations

- Representations of the binding specificity for the Kruppel transcription factor of *Drosophila melanogaster*:

  a) Kruppel site sequences;

  b) Consensus of the above sites;

  c) Degenerate consensus;

  d) Position-specific scoring matrix (PSSM);

  e) Sequence Logo obtained using WebLogo;

# Motif Discovery with oligo-analysis tool in DNA sequences

http://rsat.ulb.ac.be/rsat/
Left Panel > Motif Discovery > oligo-analysis
Select:
- Sequence type
- Oligomer Lengths
- Select Background Model

Select GO
In the results select "string-based pattern matching"

Select GO
Feature Map

Select GO

# Results from oligo-analysis

Header with parameters used for the analysis

Table with predicted sites. Each row corresponds to a predicted site, defined by its sequence, its coordinates on the input sequence and a series of scores.
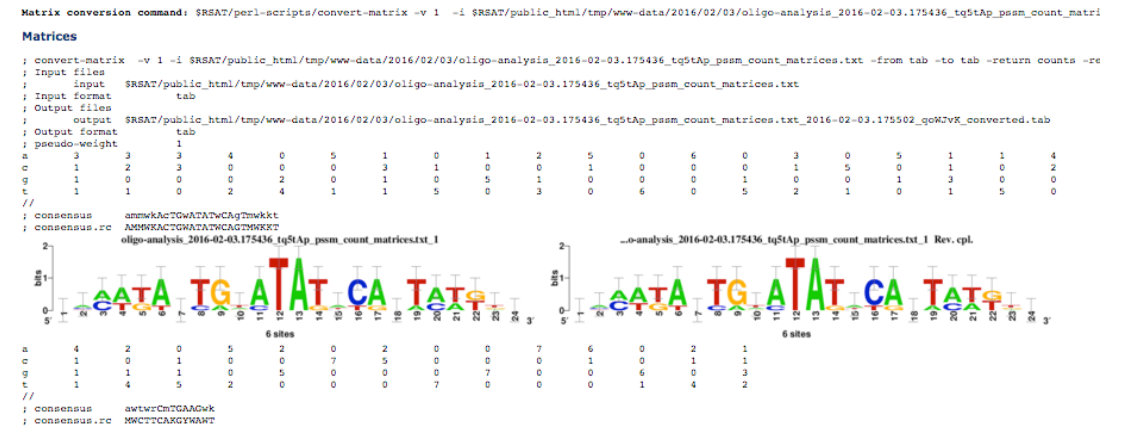
```
Oligomer length              6
Input file                   $RSAT/public_html/tmp/www-data/2016/02/03/tmp_sequence_2016-02-03.175436_yM1BMB.fasta.purged
Input format                 fasta
Output file                  $RSAT/public_html/tmp/www-data/2016/02/03/oligo-analysis_2016-02-03.175436_tq5tAp_6nt.tab
Discard overlapping matches
Counted on both strands
        grouped by pairs of reverse complements
Background model             upstream
Organism                     Saccharomyces_cerevisiae
Background estimation method Frequency file
Expected frequency file      $RSAT/public_html/data/genomes/Saccharomyces_cerevisiae/oligo-frequencies/6nt_upstream_Saccharomyces_cerevisiae-noov-2str.freq
Pseudo-frequency             0.01
Pseudo-frequency per oligo   4.80769230769231e-06
Sequence type                DNA
Nb of sequences              16
Sum of sequence lengths      3067
discarded residues           NA (quick mode)  (other letters than ACGT)
discarded occurrences        NA (quick mode)  (contain discarded residues)
nb possible positions        NA (quick mode)
total oligo occurrences      2794
total overlapping occurrences 28
total non overlapping occ    2766
alphabet size                4
nb possible oligomers        2080
oligomers tested for significance    2080
Sequences:
        cloacin-df13    200
        colicin-e1      200
        colicin-ia      200
        colicin-ib      200
        reca    200
        recn    200
        sula    200
        umu-operon      200
        uvra    200
        uvrb    200
        uvrd    200
        colicin-a       136
        lexA    173
        muc-operon      158
        hima    200
        uvrc    200

column headers
        1       seq             oligomer sequence
        2       identifier      oligomer identifier
        3       exp_freq        expected relative frequency
        4       occ             observed occurrences
        5       exp_occ         expected occurrences
        6       occ_P           occurrence probability (binomial)
        7       occ_E           E-value for occurrences (binomial)
        8       occ_sig         occurrence significance (binomial)
        9       rank            rank
        10      ovl_occ         number of overlapping occurrences (discarded from the count)
        11      forbocc         forbidden positions (to avoid self-overlap)
```

| seq | identifier | exp_freq | occ | exp_occ | occ_P | occ_E | occ_sig | rank | ovl_occ | forbocc |
|---|---|---|---|---|---|---|---|---|---|---|
| atacag | atacag|ctgtat | 0.0005908841550 | 17 | 1.65 | 2.9e-12 | 6.0e-09 | 8.22 | 1 | 0 | 85 |
| actgta | actgta|tacagt | 0.0005223348560 | 15 | 1.46 | 5.5e-11 | 1.1e-07 | 6.94 | 2 | 0 | 75 |
| acagta | acagta|tactgt | 0.0006245574949 | 16 | 1.75 | 6.6e-11 | 1.4e-07 | 6.86 | 3 | 0 | 80 |
| atactg | atactg|cagtat | 0.0005824658200 | 12 | 1.63 | 1.6e-07 | 3.3e-04 | 3.48 | 4 | 0 | 60 |
| tataca | tataca|tgtata | 0.0009981712180 | 13 | 2.79 | 7.5e-06 | 1.6e-02 | 1.81 | 5 | 1 | 65 |
| cctgaa | cctgaa|ttcagg | 0.0004525829377 | 8 | 1.26 | 5.3e-05 | 1.1e-01 | 0.96 | 6 | 0 | 40 |
| gcctga | gcctga|tcaggc | 0.0002300983709 | 6 | 0.64 | 5.6e-05 | 1.2e-01 | 0.93 | 7 | 0 | 30 |
| agcctg | agcctg|caggct | 0.0002431267464 | 6 | 0.68 | 7.6e-05 | 1.6e-01 | 0.80 | 8 | 0 | 30 |
| gctacc | gctacc|ggtagc | 0.0002569568681 | 6 | 0.72 | 0.00010 | 2.1e-01 | 0.67 | 9 | 0 | 30 |
| ctccgc | ctccgc|gcggag | 0.0001719717722 | 5 | 0.48 | 0.00014 | 3.0e-01 | 0.53 | 10 | 0 | 25 |

# Results from oligo-analysis

**Pattern Assembly**

# Additional information on Motifs

- Several high-quality transcription factor binding profile database exist:

  - JASPAR

    - Vertebrate, nematode, insects, plants, fungi, structural classes

  - TRANSFAC

    - eukaryotic **transcription factors**, their experimentally-proven binding sites, consensus binding sequences (positional weight matrices) and regulated genes.

  - CollecTF

    - database of transcription factor binding sites (**TFBS**) in the Bacteria domain.

      http://www.collectf.org/
      Search for TFs
                  select: TFS, species and experimental techniques

- Compare results from the different techniques

# Motif Scanning

FIMO
Find Individual Motif Occurences

FIMO scans a sequence database for individual matches to each of the motifs provided.

MAST
Motif Alignment & Search Tool

MAST searches sequences for matches to a set of motifs and sorts the sequences by the best combined match to all motifs.

Go to: http://meme-suite.org/tools/fimo
o   Select the consensus from oligo-analysis
o   Select the species fot the input sequence

Analyze the results. How many sequences have matches?

# Motif Discovery in DNA and Protein Sequences

Word based and Expectation Maximization based Methods

Material partially adapted from:
*Discovering Novel Sequence Motifs with MEME*, Current Protocols in Bioinformatics 2002, TL Bailey
*Using RSAT to scan genome sequences for transcription factor binding sites and cis-regulatory modules*, Nat. Protocols 2008, Turatsinze et al.