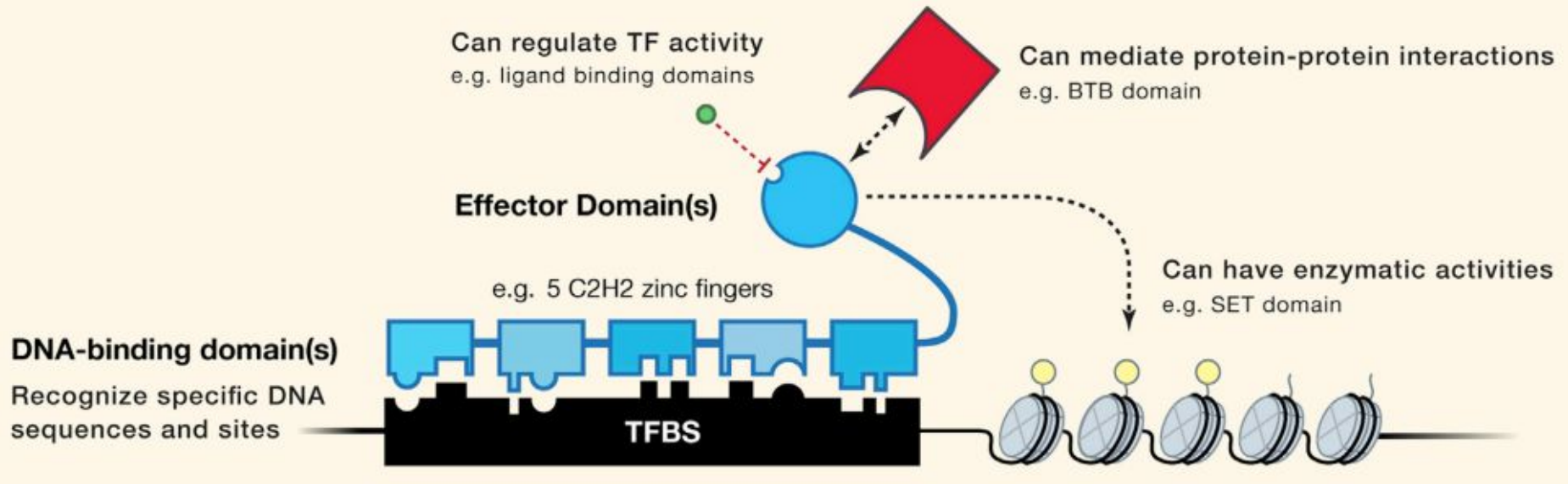# Lecture 7: Regulatory genomics

- DNA-binding sites/motifs
  - ChIP-seq
  - Position-weight matrices
  - Motif-finding
    - Expectation-Maximization
    - Gibbs Sampling

# Transcriptional regulation by TFs

# Consensus sequence of DNA-binding sites

EcoRI binds to the 6-mer GAATTC (palindrome).

- occurs once every $4^6$ (= 4,096) bp in a random DNA sequence.

HindII bind to GTYRAC.

- occur once per $4^4 \times 2^2$ (= 1,024) bp.



| | | |
|---|---|---|
| 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 |

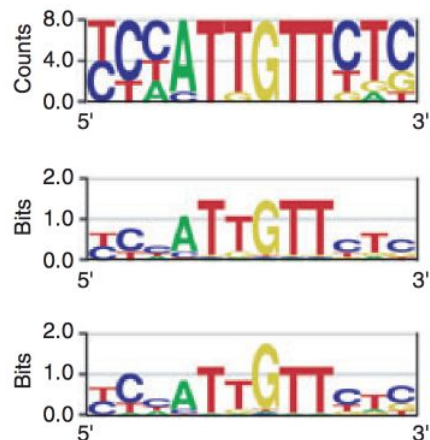Position frequency matrix

Sequence logo

# Consensus sequence of DNA-binding sites

```
A  002700000010
C  464100000505
G  000001800112
T  422087088261
```



$$I_i = 2 + \sum_b f_{b,i} \log_2 f_{b,i}$$

Scaling sequence logos based on 'information content' than frequency.

- $f_{b,i}$ : frequency of base b at position i.
- Perfectly conserved: 2 bits of information.
- Two of the four bases occur 50% of the time each: 1 bit.
- All four bases occur equally often: no information.

HindII bind to GTYRAC.

- What is its information content?

D'haeseleer (2016) Nat. Biotech.

# Consensus sequence of DNA-binding sites

```
A  002700000010
C  464100000505
G  000001800112
T  422087088261
```



$$I_{seq}(i) = -\sum_b f_{b,i} \log_2 \frac{f_{b,i}}{P_b}$$

Relative entropy (a.k.a. Kullback-Leibler distance) to correct for background nucleotide frequencies.

$$W(b,i) = \log_2 \frac{f_{b,i}}{P_b}$$

Position weight matrix (PWM).

# Consensus sequence of DNA-binding sites



```
A  002700000010
C  464100000505
G  000001800112
T  422087088261
```

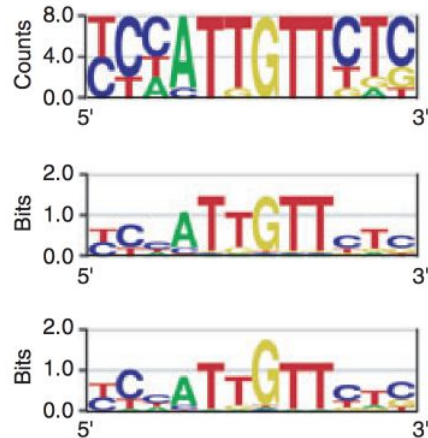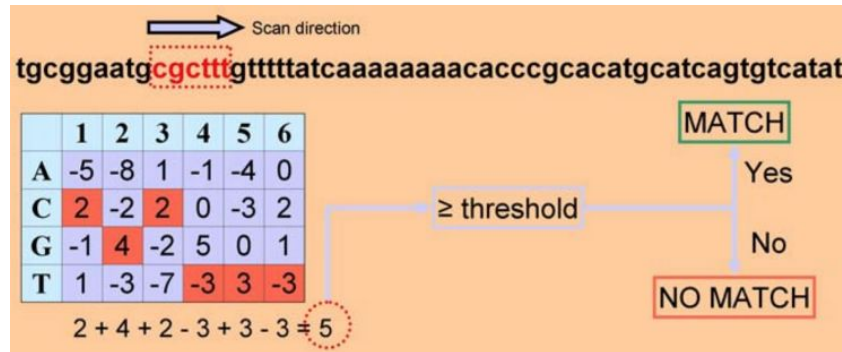$$I_{seq}(i) = -\sum_b f_{b,i} \log_2 \frac{f_{b,i}}{P_b}$$

Relative entropy (a.k.a. Kullback-Leibler distance) to correct for background nucleotide frequencies.

$$W(b,i) = \log_2 \frac{f_{b,i}}{P_b}$$

Position weight matrix (PWM).



D'haeseleer (2016) Nat. Biotech.

# Mapping of regulatory elements using ChIP-chip and ChIP-seq



Visel (2009) Nature

# Mapping of regulatory elements using ChIP-chip and ChIP-seq

# Mapping of regulatory elements using ChIP-chip and ChIP-seq



Park (2009) Nat. Rev. Genet.

Sequences are not aligned.

We don't know what the motif looks like.

The motif model learning task:

- Given: a set of sequences that are thought to contain occurrences of an unknown motif of interest
- Do:
    - infer a model (PWM) of the motif, and
    - predict the locations of the motif occurrences in the given sequences.

# Expectation-Maximization algorithm (EM)

**a** Maximum likelihood



5 sets, 10 tosses per set

|  | Coin A | Coin B |
|---|---|---|
| H T T T H H T H T H |  | 5 H, 5 T |
| H H H H T H H H H H | 9 H, 1 T |  |
| H T H H H H H T H H | 8 H, 2 T |  |
| H T H T T T H H T T |  | 4 H, 6 T |
| T H H H T H H H T H | 7 H, 3 T |  |
| | 24 H, 6 T | 9 H, 11 T |

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$$

A coin-flipping experiment

- $\theta_A$ & $\theta_B$ are the biases of two coins A & B.

- Goal: estimate $\theta = (\theta_A, \theta_B)$ by repeating the following procedure five times:

  - Randomly choose one of the two coins (with equal probability), and perform ten independent coin tosses with the selected coin.

  - Total of 50 coin tosses.

$x = (x_1, x_2, ..., x_5) \mid x_i \in \{0,1,...,10\}$ is the no. of heads observed during the ith set of tosses.

$z = (z_1, z_2,..., z_5) \mid z_i \in \{A,B\}$ is the identity of the coin used during the ith set of tosses.

Maximum likelihood estimation: statistical model that has the highest probability of generating the observed data – $\theta$ that maximizes $\log P(x,z;\theta)$.

Do & Batzoglou (2008) Nat. Biotech.

# Expectation-Maximization algorithm (EM)

**a** Maximum likelihood



| | Coin A | Coin B |
|---|---|---|
| H T T T H H T H T H | | 5 H, 5 T |
| H H H H T H H H H H | 9 H, 1 T | |
| H T H H H H H T H H | 8 H, 2 T | |
| H T H T T T H H T T | | 4 H, 6 T |
| T H H H T H H H T H | 7 H, 3 T | |
| | 24 H, 6 T | 9 H, 11 T |

5 sets, 10 tosses per set

$$\hat{\theta}_A = \frac{24}{24 + 6} = 0.80$$

$$\hat{\theta}_B = \frac{9}{9 + 11} = 0.45$$

$x = (x_1, x_2, …, x_5) \mid x_i \in \{0,1,…,10\}$ is the no. of heads observed during the ith set of tosses.

$z = (z_1, z_2,…, z_5) \mid z_i \in \{A,B\}$ is the identity of the coin used during the ith set of tosses. [**Hidden variables / Latent factors**]
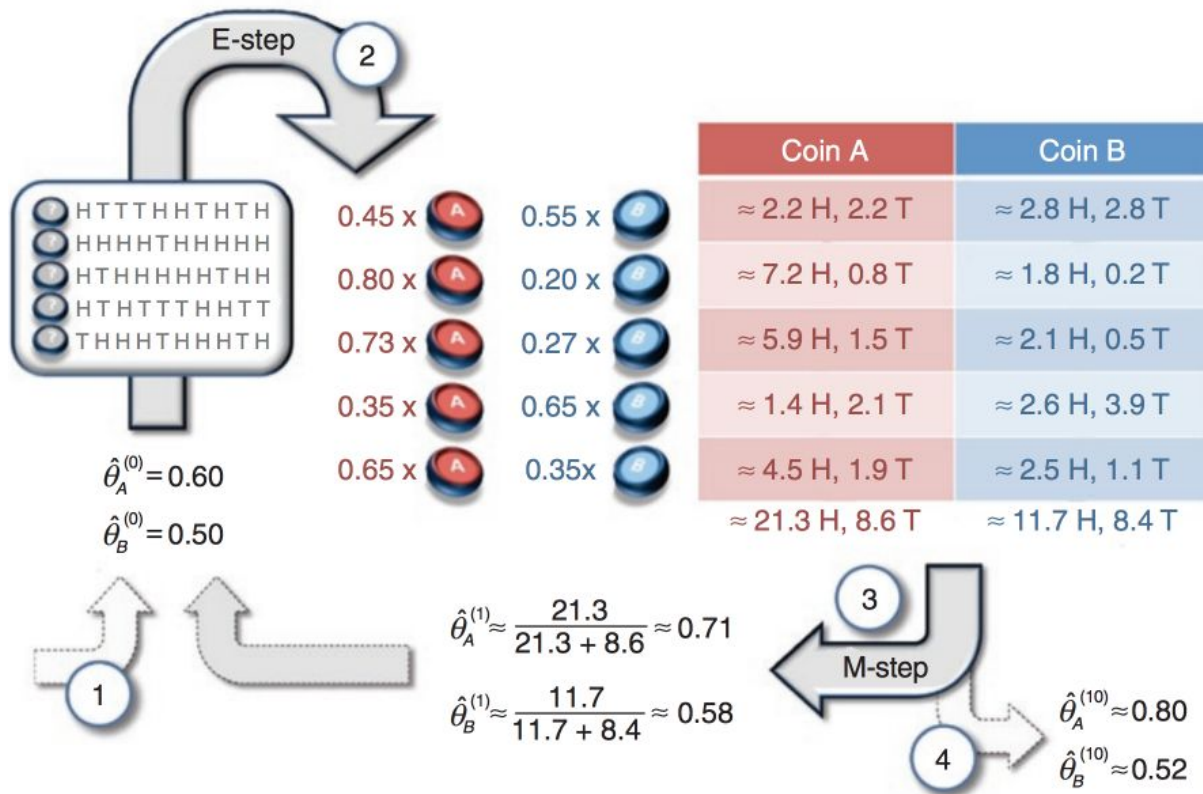
A coin-flipping experiment

- $\theta_A$ & $\theta_B$ are the biases of two coins A & B.

- Goal: estimate $\theta = (\theta_A, \theta_B)$ by repeating the following procedure five times:

  - Randomly choose one of the two coins (with equal probability), and perform ten independent coin tosses with the selected coin.

- **Not told which coin was chosen.**

Do & Batzoglou (2008) Nat. Biotech.

# Expectation-Maximization algorithm (EM)



**b** Expectation maximization

| | Coin A | Coin B |
|---|---|---|
| 0.45 x A   0.55 x B | ≈ 2.2 H, 2.2 T | ≈ 2.8 H, 2.8 T |
| 0.80 x A   0.20 x B | ≈ 7.2 H, 0.8 T | ≈ 1.8 H, 0.2 T |
| 0.73 x A   0.27 x B | ≈ 5.9 H, 1.5 T | ≈ 2.1 H, 0.5 T |
| 0.35 x A   0.65 x B | ≈ 1.4 H, 2.1 T | ≈ 2.6 H, 3.9 T |
| 0.65 x A   0.35x B | ≈ 4.5 H, 1.9 T | ≈ 2.5 H, 1.1 T |
| | ≈ 21.3 H, 8.6 T | ≈ 11.7 H, 8.4 T |

HTTTHHTHTH
HHHHTHHHHH
HTHHHHHTHH
HTHTTTHHTT
THHHTHHHTH

$\hat{\theta}_A^{(0)} = 0.60$

$\hat{\theta}_B^{(0)} = 0.50$

$\hat{\theta}_A^{(1)} \approx \dfrac{21.3}{21.3 + 8.6} \approx 0.71$

$\hat{\theta}_B^{(1)} \approx \dfrac{11.7}{11.7 + 8.4} \approx 0.58$

$\hat{\theta}_A^{(10)} \approx 0.80$

$\hat{\theta}_B^{(10)} \approx 0.52$

E-step:

- Estimate $P(x_i, z_i | \theta^{(t)})$ and the expected values of the hidden variables.
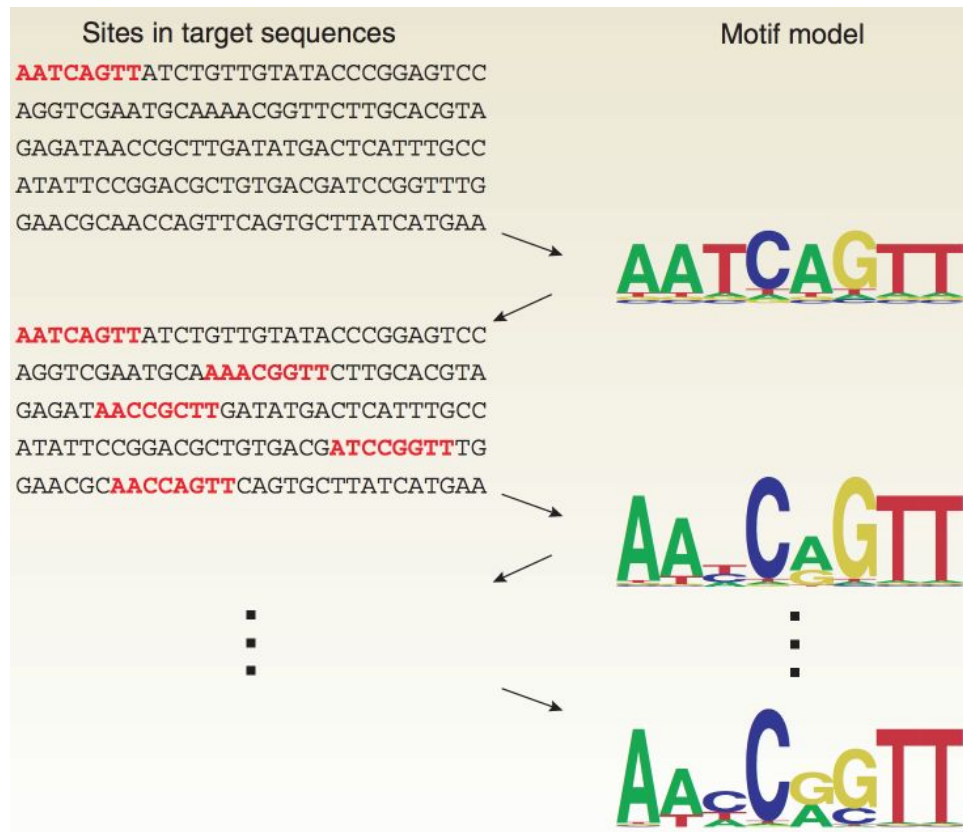
M-step:

- Estimate new parameters $\theta^{(t+1)}$ given current estimates of hidden variables & parameters.

Repeat until convergence.

$P(x_i, z_i | \theta^{(t)})$: Likelihood function, from here on also going to be written as $P(X, Z | \theta)$.

Do & Batzoglou (2008) Nat. Biotech.

# Expectation-Maximization algorithm (EM)

1. Define the probabilistic model and the likelihood function P(X | θ).

2. Identify the hidden variables (Z).

   a. Here, they are the locations of the motifs in each sequence.

3. Write the E step.

   a. Compute the expected values of the hidden variables given current parameter values.

4. Write the M step.

   a. Determine new parameters given the expected values of the hidden variables.

5. Repeat until convergence.



D'haeseleer (2016) Nat. Biotech.
Gitter @ U. Wisconsin

# Motif-finding using MEME

- MEME: Multiple EM for Motif Elicitation

- A motif is:

  - assumed to have a fixed width, **W**

  - represented by a matrix of probabilities: $p_{c,k}$ (probability of character **c** in column **k**).

- The "background" (i.e. sequence outside the motif) is given by $p_{c,0}$ (probability of character **c** in the background).

- Data is a collection of sequences, denoted **X**.

- Motif starting positions are represented by a matrix indicator variables (0/1) $Z_{i,j}$.

A motif model of length 3

$$p = \begin{array}{c} \\ \text{A} \\ \text{C} \\ \text{G} \\ \text{T} \end{array} \begin{array}{cccc} 0 & 1 & 2 & 3 \\ 0.25 & 0.1 & 0.5 & 0.2 \\ 0.25 & 0.4 & 0.2 & 0.1 \\ 0.25 & 0.3 & 0.1 & 0.6 \\ 0.25 & 0.2 & 0.2 & 0.1 \end{array}$$

background        motif positions

Given sequences L = 6.
Possible starting positions m = L − W + 1

$$Z =$$

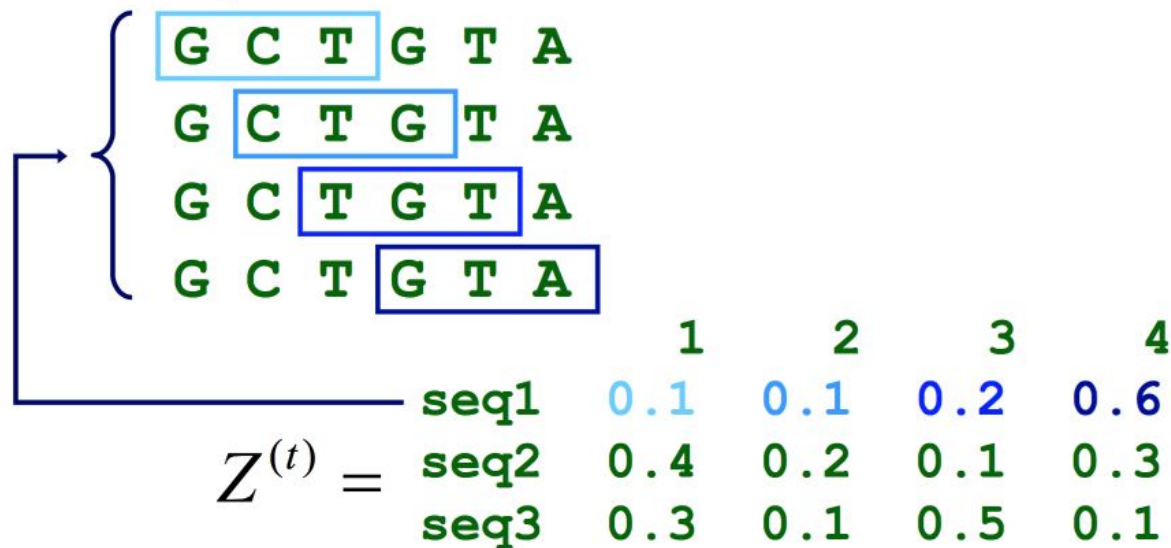| | | | | | | | | | | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| G | T | C | A | G | G | | seq1 | | | 0 | 0 | 1 | 0 |
| G | A | G | A | G | T | | seq2 | | | 1 | 0 | 0 | 0 |
| A | C | G | G | A | G | | seq3 | | | 0 | 0 | 0 | 1 |
| C | C | A | G | T | C | | seq4 | | | 0 | 1 | 0 | 0 |

# Motif-finding using MEME

1.  Define the probabilistic model and the likelihood function $P(X \mid \theta)$.

2.  Identify the hidden variables (Z).

    a.  Here, they are the locations of the motifs in each sequence.

3.  Write the E step.

    a.  Compute the expected values of the hidden variables given current parameter values.

4.  Write the M step.

    a.  Determine new parameters given the expected values of the hidden variables.

5.  Repeat until convergence.

given: length parameter **W**, set of sequences

  t=0

  set initial values for $p^{(0)}$

  do

    ++t

    re-estimate $Z^{(t)}$ from $p^{(t-1)}$ (E-step)

    re-estimate $p^{(t)}$ from $Z^{(t)}$ (M-step)

  until change in $p^{(t)} < \varepsilon$

return: $p^{(t)}$, $Z^{(t)}$

# Motif-finding using MEME

- **E-step**: compute the expected values of Z given X and $p^{(t-1)}$

- Expected values: $Z^{(t)} = E[ Z | X, p^{(t-1)} ]$

- For example:



$$Z^{(t)} = \begin{array}{c|cccc} & 1 & 2 & 3 & 4 \\ \hline seq1 & 0.1 & 0.1 & 0.2 & 0.6 \\ seq2 & 0.4 & 0.2 & 0.1 & 0.3 \\ seq3 & 0.3 & 0.1 & 0.5 & 0.1 \end{array}$$

given: length parameter **W**, set of sequences

  t=0

  set initial values for $p^{(0)}$

  do

    ++t

    re-estimate $Z^{(t)}$ from $p^{(t-1)}$ (E-step)

    re-estimate $p^{(t)}$ from $Z^{(t)}$ (M-step)

  until change in $p^{(t)} < \varepsilon$

return: $p^{(t)}, Z^{(t)}$

Gitter @ U. Wisconsin

# Motif-finding using MEME

- **E-step**: compute the expected values of Z given X and $p^{(t-1)}$

- Expected values: $Z^{(t)} \square = E[\, Z \mid X, p^{(t\square-1)}\,]$

- Applying Bayes rule to: $P(Z_{i,j} = 1 \mid X_i, p^{(t-1)})$

$$Z_{i,j}^{(t)} = \frac{P(X_i \mid Z_{i,j} = 1, p^{(t-1)})P(Z_{i,j} = 1)}{\sum_{k=1}^{m} P(X_i \mid Z_{i,k} = 1, p^{(t-1)})P(Z_{i,k} = 1)}$$

given: length parameter **W**, set of sequences

   t=0

   set initial values for $p^{(0)}$

   do

     ++t

     re-estimate $Z^{(t)}$ from $p^{(t-1)}$ (E-step)

     re-estimate $p^{(t)}$ from $Z^{(t)}$ (M-step)

   until change in $p^{(t)} < \varepsilon$

return: $p^{(t)}$, $Z^{(t)}$

$$Z_{i,j}^{(t)} = \frac{P(X_i \mid Z_{i,j} = 1, p^{(t-1)})}{\sum_{k=1}^{m} P(X_i \mid Z_{i,k} = 1, p^{(t-1)})}$$

Assuming that it is equally likely that the motif will start in any position

$$P(Z_{i,j} = 1) = \tfrac{1}{m}$$

Gitter @ U. Wisconsin

# Motif-finding using MEME

Probability of a Sequence Given a Motif Starting Position



$$P(X_i \mid Z_{i,j} = 1, p) = \prod_{k=1}^{j-1} p_{c_k, 0} \prod_{k=j}^{j+W-1} p_{c_k, k-j+1} \prod_{k=j+W}^{L} p_{c_k, 0}$$

| Before motif | Motif | After motif |
|---|---|---|

- $X_i$ is the i th sequence
- $Z_{i,j}$ is 1 if motif starts at position j in sequence i
- $c_k$ is the base at position k in sequence i

# Motif-finding using MEME

Probability of a Sequence Given a Motif Starting Position



$$P(X_i \mid Z_{i,j} = 1, p) = \prod_{k=1}^{j-1} p_{c_k,0} \prod_{k=j}^{j+W-1} p_{c_k, k-j+1} \prod_{k=j+W}^{L} p_{c_k,0}$$

| Before motif | Motif | After motif |

$$X_i = \text{G C } \boxed{\text{T G T}} \text{ A G}$$

$$p = \begin{array}{c|cccc} & 0 & 1 & 2 & 3 \\ \text{A} & 0.25 & 0.1 & 0.5 & 0.2 \\ \text{C} & 0.25 & 0.4 & 0.2 & 0.1 \\ \text{G} & 0.25 & 0.3 & 0.1 & 0.6 \\ \text{T} & 0.25 & 0.2 & 0.2 & 0.1 \end{array}$$

- $X_i$ is the i th sequence
- $Z_{i,j}$ is 1 if motif starts at position j in sequence i
- $c_k$ is the base at position k in sequence i

$$P(X_i \mid Z_{i,3} = 1, p) =$$

$$p_{G,0} \times p_{C,0} \times p_{T,1} \times p_{G,2} \times p_{T,3} \times p_{A,0} \times p_{G,0} =$$

$$0.25 \times 0.25 \times 0.2 \times 0.1 \times 0.1 \times 0.25 \times 0.25$$

$$P(X_i \mid Z_{i,1} = 1, p^{(t-1)}) \quad ?$$

# Motif-finding using MEME

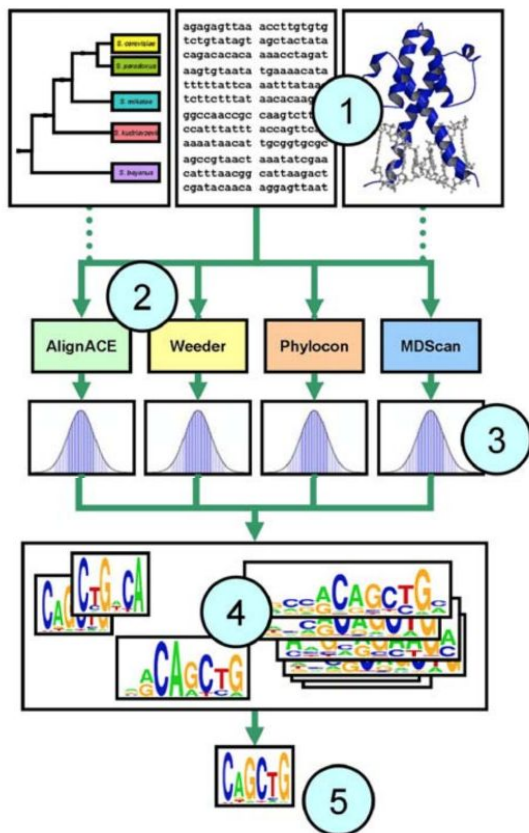- **M-step**: Estimate $p^{(t)}$ given X and $Z^{(t)}$.

- $p_{c,k}$ represents the probability of character c in position k.

- k=0 represents the background.

$$p_{c,k}^{(t)} = \frac{n_{c,k} + d_{c,k}}{\sum_{b \in \{A,C,G,T\}} (n_{b,k} + d_{b,k})}$$

$$n_{c,k} = \begin{cases} \sum_{i} \sum_{\{j|X_{i,j+k-1}=c\}} Z_{i,j}^{(t)} & k > 0 \\ \\ n_c - \sum_{j=1}^{W} n_{c,j} & k = 0 \end{cases}$$

total # c's in the dataset

sum over positions where c appers

**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)

MacIsaac & Fraenkel (2006) PLoS Comp. Biol.