

# Discovering binding site motifs

Yaron Orenstein

Simons Institute Bootcamp  
Gene regulation mini-course

January 21<sup>st</sup> 2016

P value	HS binding sequence	MM binding sequence
1.0E-45		
1.0E-29		
1.0E-136		
1.0E-27		

# Part I: de novo motif discovery using EM

Slides sources:

Computational genomics, TAU

Prof. Roded Sharan, Prof. Ron Shamir

ARG81		PhyloCon
ARO80		PhyloCon
ARR1		Converge
ASH1		Converge

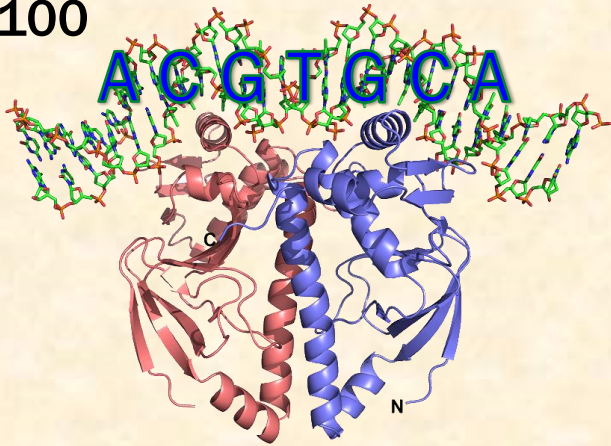
# Gene expression regulation

- Transcription is regulated mainly by transcription factors (**TFs**) - proteins that bind to DNA subsequences, called binding sites (**BSs**).
- TFBSs are located mainly in the gene's promoter – the DNA sequence upstream the gene's transcription start site.

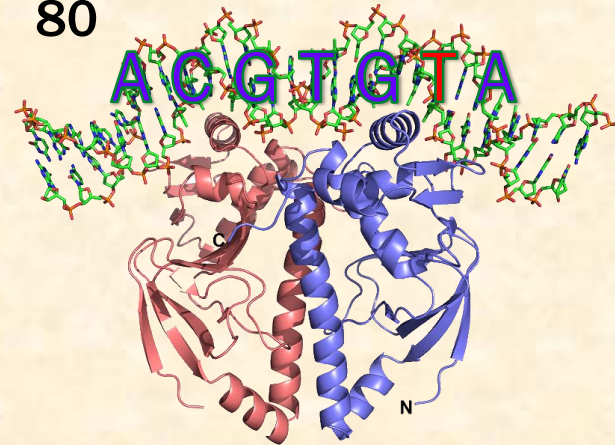


# Binding preferences

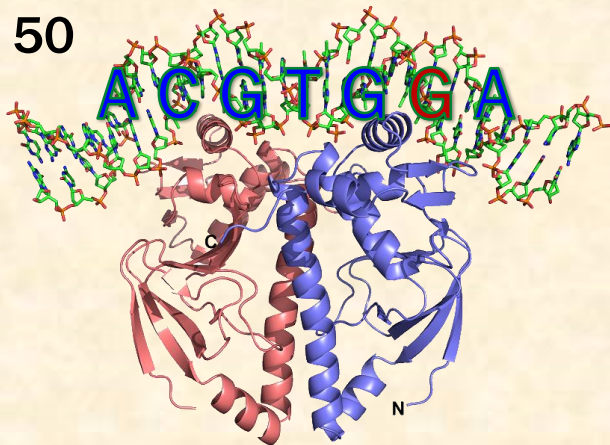
100



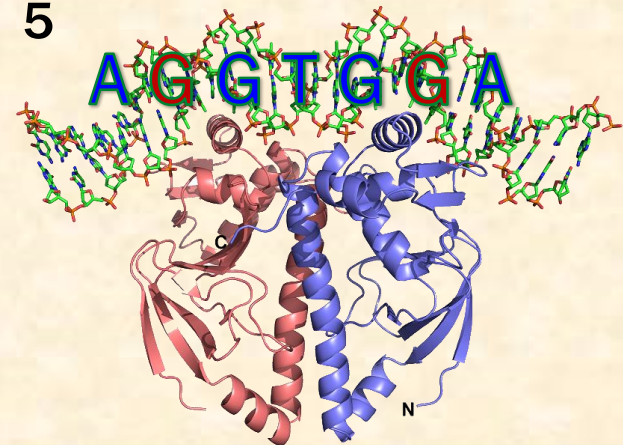
80



50



5



# TFBS models

The BSs of a particular TF share a common pattern, or **motif**, which is often modeled using:

- **Consensus string**

**GGAATT**

with up to  $d$  mismatches

$d=1$

ATC**GGAATT**CTGCAG  
G**GCAATT**CG**GGAATG**  
A**GGTATT**CTCAGATTA

- **Degenerate string**

**GGWATB** ( $W=\{A,T\}$ ,  $B=\{C,G,T\}$ )

ATC**GGAATT**CTGCAG  
GGCAATT**CGGGAATG**  
A**GGTATT**CTCAGATTA

# TFBS models

- PWM = Position weight matrix**

	1	2	3	4	5	6
A	0.1	0.8	0	0.7	0.2	0
C	0	0.1	0.5	0.1	0.4	0.6
G	0	0	0.5	0.1	0.4	0.1
T	0.9	0.1	0	0.1	0	0.3

➤ Cutoff = 0.009

AGCTACACCCATTAT 0.06  
 AGTAGAGCCTTCGTG 0.06  
 CGATTCTACAATATGA 0.01



- K-mer model**

```

ATGATTAATTGC 1.00
GCAATTAATCAT 0.96
GTAATTAATCAT 0.94
ATGATTAATCAT 0.91
GTCATTAATCAT 0.91
ATGATTAATGAC 0.91
AATGATTAATGA 0.88
ATGATTAATGGC 0.88
AATGATTAATGG 0.86
AATGATTAATTG 0.86
    
```

- other variants:**

**DNA shape features**  
**di-nucleotides**

# Motif profile

- Given the binding sites, the motif profile is easily constructed.

Alignment

```

a G g t a c T t
C c A t a c g t
a c g t T A g t
a c g t C c A t
C c g t a c g G
  
```

Profile

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

Consensus **A C G T A C G T**

- Line up the patterns by their start indexes

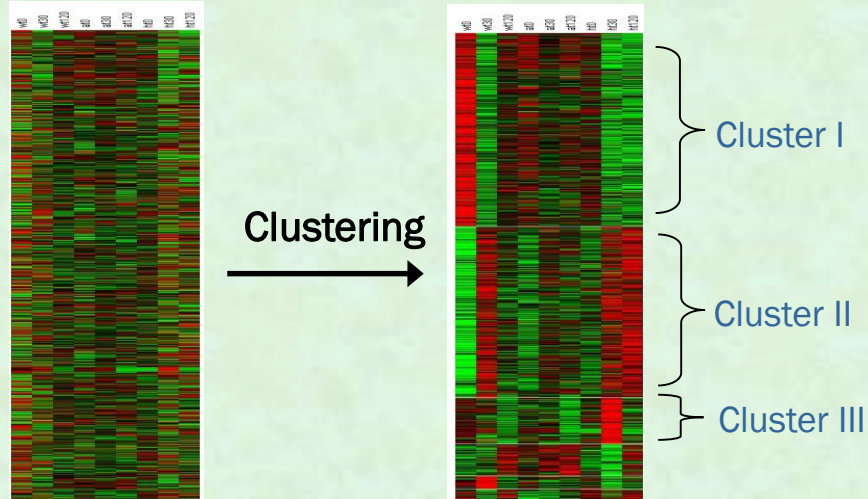
$$s = (s_1, s_2, \dots, s_t)$$

- Construct matrix profile with frequencies of each nucleotide in columns

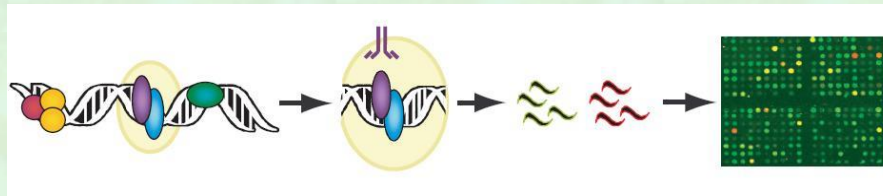
# Motif discovery

## Co-regulated gene set

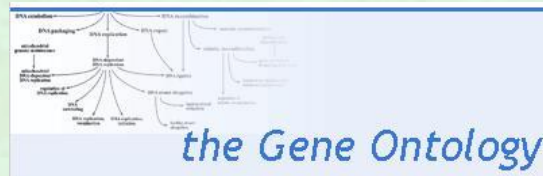
Gene expression microarrays



Location analysis  
(ChIP-chip, ...)



Functional group  
(e.g., GO term)



## Promoter/3'UTR sequences



Motif discovery



1 GG 2 3 4 5 TT 6 7 8 9 10 CC

The image shows a sequence logo for a motif. The sequence is GGTTCC, with positions 1 through 10 indicated below the letters. The letters are colored: G (yellow), G (yellow), T (red), T (red), C (blue), and C (blue). The positions are numbered 1 through 10.



# Motif discovery: *Goals and challenges*

- **Goal:** Given a set of co-regulated genes, find a recurrent motif in their promoter regions.
- **Challenges:**
  - **BSs** are short and degenerate (**non-specific**)
  - **Promoters** are long + complex (**hard to model**)
  - **Search space** is huge (**motif and sequence**)
  - **Data** is noisy

# An example: Implanting Motif AAAAAAAAAGGGGGGG

atgaccgggatactgatAAAAAAAAAGGGGGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaAAAAAatttgagtacaaaactTTTccgaataAAAAAAAAAGGGGGGGa  
tgagtatccctgggatgacttAAAAAAAAAGGGGGGGtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgAAAAAAAAAGGGGGGGtccacgcaatcgcaaccaacgcgaccCAAaggcaagaccgataaaggaga  
tcctTTTgcggaatgtgccgggaggctggttacgtagggaagccctaacggacttaatAAAAAAAAAGGGGGGGcTTatag  
gtcaatcatgTTcttTgtaatggatttAAAAAAAAAGGGGGGGgaccgcttggcgcacccaattcagtgtgggCGagcgcaa  
cgTTTTggcccttTgtagaggccccgtAAAAAAAAAGGGGGGGcaattatgagagagctaatttatcgCGtgCGtTtcat  
aacttgagttAAAAAAAAAGGGGGGGctggggcacatacaagaggagtcttCttatcagTTaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcattAAAAAAAAAGGGGGGGaccgaaaggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttAAAAAAAAAGGGGGGGa

# Where is the Implanted Motif? (\*)

atgaccgggatactgataaaaaaagggggggggtacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaaggggggga  
tgagtatccctgggatgacttaaaaaaagggggggtgctctcccgatTTTTgaatatgtaggatcattcgccaggggtccga  
gctgagaattggatgaaaaaaagggggggtccacgcaatcggaaccaacgcgacccaaaggcaagaccgataaaggaga  
tcctTTTgcggaatgtgccgggaggctggttacgtagggaagccctaacggacttaataaaaaaagggggggttatag  
gtcaatcatgttcttTggaatggatttaaaaaaaggggggggaccgcttggcgcacccaaattcagtgtgggagcgcaa  
cggTTTTggccttTtagaggccccgtaaaaaaagggggggcaattatgagagagctaatttatcgcgTgctgttcat  
aacttgagttaaaaaaagggggggtggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaagggggggaccgaaaggggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttaaaaaaaggggggga

# Implanting Motif AAAAAAGGGGGG with Four Mutations

atgaccgggatactgatAgAAgAAAGGttGGGggcgtacacattagataaacgtatgaagtacgttagactcggcgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataCAAtAAACGGcGGGga  
tgagtaccctgggatgacttAAAAtAAtGGaGtGGtgctctcccgattTTTTgaatatgtaggatcattcgccaggggccga  
gctgagaattggatgCAAAAAAAGGGattGtccacgcaatcgcgaaaccaacgcggaaccaaggcaagaccgataaaggaga  
tccTTTTgCGgtaatgtgCCgggaggctggttacgtaggaagccctaacggacttaAtAAtAAAGGaaGGGcttatag  
gtcaatcatgttcttgtgaatggattAACAAtAAGGGctGGgaccgcttggcgacccaaattcagtgtggcgagcgcaa  
cggTTTTggcccttgttagaggccccgtAtAAACAAAGGaGGGccaattatgagagagctaatactatcgctgctgttcat  
aacttgagttAAAAAtAGGGaGccctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta  
ttggccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataActAAAAAGGaGcGGgaccgaaaggaag  
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttActAAAAAGGaGcGGa

# Where is the Motif???

atgaccgggatactgatagaagaaagggttgggggctacacattagataaacgtatgaagtacgttagactcggcgccgccg  
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaatacaataaaacggcgga  
tgagtatccctgggatgacttaaataatggagtggtgctctcccgatTTTTgaatatgtaggatcattcgccagggTccga  
gctgagaattggatgcaaaaaagggttgtccacgcaatcgcgaaccaacgCGGaccCAAaggcaagaccgataaaggaga  
tcctTTTgCGGtaatgtGCCgggaggctggttacgtagggaagccCTaacggacttaataataaaaggaagggttatag  
gtcaatcatgTtcttgTgaatggatttaacaataagggtgggaccgcttggcgcacCCaaattcagtgTgggCGagCGcaa  
CGTTTTgGCcttgTtagaggCCCCgtataaacaaggaggccaattatgagagagctaatttatCGcgtGcgtgTtcat  
aacttgagTtaaaaaataggagccctggggcacatacaagaggagtcttCcttatcagTtaatgctgTatgacactatgta  
ttggcccattggctaaaagCCcaacttgacaaatggaagatagaatccttgcatactaaaaaggagCGGaccgaaagggaag  
ctggTgagcaacgacagattcttacgTgcattagctcGcttccggggatctaatagcacgaagcttactaaaaaggagCGga

# MEME

## Multiple EM for Motif Elicitation

[Bailey, Elkan ISMB '94]

**Goal:** Given a set of sequences, find a motif (PWM) that maximizes the expected likelihood of the data

**Technique:** EM (Expectation Maximization)  
(based on [Lawrence, Reilly '90])

# EM reminder: the goal

- Input: **incomplete** data originating from a probability distribution with some unknown parameters
- Want to find the parameter values that maximize the likelihood
- EM – approach that helps when maximum likelihood solution cannot be directly computed.
- Seeks a local maximum by iteratively solving two easier subproblems

# EM-reminder: the setting

---

Input: data  $X$  coming from a probabilistic model with hidden information  $y$

Goal: Learn the model's parameters  $\theta$  so that the likelihood is maximized.

---



# EM-reminder: the algorithm

Main component:

$$Q(\theta | \theta^t) = \sum_y P(y | x, \theta^t) \cdot \log P(x, y | \theta)$$

$\log P(x, y | \theta)$  is called the **complete log likelihood** function

→  $Q$  is the expectation of the complete log likelihood over the distribution of  $y$  given the current parameters  $\theta^t$

The algorithm:

repeat

- **E-step:** Calculate the  $Q$  function
- **M-step:** Maximize  $Q(\theta | \theta^t)$  with respect to  $\theta$
- Stopping criterion: improvement in log likelihood  $\leq \epsilon$

Note: local optimum guaranteed to be reached, not global.  
Starting point matters! Try many..

# MEME: The Mixture Model

Data:  $X = (X_1, \dots, X_n)$  :

all (overlapping)  $l$ -mers in the input sequences

Assume  $X_i$ 's were generated by a two-component mixture model -  $\theta = (\theta_1, \theta_2)$  :

Model #1:  $\theta_1 =$  motif model:

$f_{i,b}$  = prob. of base  $b$  at pos  $i$  in motif,  $1 \leq i \leq l$

Model #2:  $\theta_2 =$  background (BG) model:

$f_{0,b}$  = prob. of base  $b$

Mixing parameter:  $\lambda = (\lambda_1, \lambda_2)$

$\lambda_j$  = prob. that model # $j$  is used ( $\lambda_1 + \lambda_2 = 1$ )

Assume independence between  $l$ -mers

# Log Likelihood

Missing data:  $Z = (Z_1, \dots, Z_n)$  :

$Z_i = (Z_{i1}, Z_{i2})$ ;  $Z_{ij} = 1$  if  $X_i$  from model # $j$  ; 0 o/w

Complete Likelihood of model given data:

$$L(\theta, \lambda \mid X, Z) = p(X, Z \mid \theta, \lambda)$$

$$= \prod_{i=1 \dots n} p(X_i, Z_i \mid \theta, \lambda)$$

$$p(X_i, Z_i \mid \theta, \lambda) = p(X_i \mid Z_i, \theta, \lambda) p(Z_i \mid \theta, \lambda) =$$

$$= \lambda_1 p(X_i \mid \theta_1) \text{ if } Z_{i1}=1; \lambda_2 p(X_i \mid \theta_2) \text{ if } Z_{i2}=1$$

$$\rightarrow \log L = \sum_{i=1 \dots n} \sum_{j=1,2} Z_{ij} \log(\lambda_j p(X_i \mid \theta_j))$$

# MEME: Algorithm

Goal: Maximize  $E[\log L]$

Outline of EM algorithm:

- Choose starting  $\theta^0, \lambda^0$
- Repeat until convergence of  $\theta$ :
  - E-step: Re-estimate  $Z$  from  $\theta^t, \lambda^t, X$
  - M-step: Re-estimate  $\theta^{t+1}, \lambda^{t+1}$  from  $X, Z$
- Repeat all of the above for various  $\theta, \lambda \dots$

# E-step

Compute expectation of  $\log L$  over  $Z$ :

$$E_Z[\log L] = \sum_{i=1 \dots n} \sum_{j=1,2} Z'_{ij} \log (\lambda_j p(X_i | \theta_j))$$

where:

$$\begin{aligned} Z'_{ij} &= p(Z_{ij}=1 | \theta^t, \lambda^t, X_i) = \\ &= p(Z_{ij}=1, X_i | \theta^t, \lambda^t) / p(X_i | \theta^t, \lambda^t) = \\ &= p(Z_{ij}=1, X_i | \theta^t, \lambda^t) / \sum_{k=1,2} p(Z_{ik}=1, X_i | \theta^t, \lambda^t) = \\ &= \lambda_j^t p(X_i | \theta_j^t) / \sum_{k=1,2} \lambda_k^t p(X_i | \theta_k^t) \end{aligned}$$

# M-step

Find  $\theta, \lambda$  that maximize  $E[\log L] = Q(\theta, \lambda | \theta^t, \lambda^t)$ :

$$E[\log L] = \sum_{i=1 \dots n} \sum_{j=1,2} Z'_{ij} \log (\lambda_j p(X_i | \theta_j))$$

Finding  $\lambda$ :

Suffices to maximize  $L_1 = \sum_{i=1 \dots n} \sum_{j=1,2} Z'_{ij} \log \lambda_j$

$$\lambda_1 + \lambda_2 = 1 \rightarrow L_1 = \sum_{i=1 \dots n} (Z'_{i1} \log \lambda_1 + Z'_{i2} \log (1 - \lambda_1))$$

$$dL_1/d\lambda_1 = \sum_{i=1 \dots n} (Z'_{i1} / \lambda_1 - Z'_{i2} / (1 - \lambda_1))$$

# MEME: Algorithm

M-step (cont.):

$$dL_1/d\lambda_1 = \sum_{i=1\dots n} (Z'_{i1} / \lambda_1 - Z'_{i2} / (1-\lambda_1)) = 0$$

$$\rightarrow \lambda_1 \sum_{i=1\dots n} Z'_{i2} = (1-\lambda_1) \sum_{i=1\dots n} Z'_{i1}$$

$$\rightarrow \lambda_1 (\sum_{i=1\dots n} (Z'_{i1} + Z'_{i2})) = \sum_{i=1\dots n} Z'_{i1}$$

$$\rightarrow \lambda_1 = (\sum_{i=1\dots n} Z'_{i1}) / n$$

$$\lambda_2 = 1 - \lambda_1 = (\sum_{i=1\dots n} Z'_{i2}) / n$$

Expected #times that letter  $k$   
is produced at column  $j$

Finding  $\theta$ :  $\theta_{jk} = C_{jk} / \sum_{k \in \{A,C,G,T\}} C_{jk}$

$$C_{jk} = \sum_{i=1\dots n} Z'_{i1} I(k, X_{ij})$$

# Multiple passes

- To find multiple motifs: After finding one, erase (mask) its positions from all sequences and rerun the alg.
- By doing several passes, a drop in the log likelihood score may indicate which are true motifs



## MEME Suite Menu

- Submit A Job
- Documentation
- Downloads
- User Support
- Alternate Servers
- Authors
- Citing

New! Postdoc Available



# MEME

## Multiple Em for Motif Elicitation

Version 4.3.0

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.

### 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 **60000 characters** total in any of a large number of **formats**.

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

 Browse... 

or  
the **actual sequences** here (**Sample Protein Input Sequences**):

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 ( $\geq 2$ )

**Maximum** width ( $\leq 300$ )

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 ( $\geq 2$ )

**Maximum** sites ( $\leq 300$ )

**Shuffle** sequence letters

Enter the name of a file containing a **background Markov model**:

 Browse... 

#### DNA-ONLY OPTIONS

(Ignored for protein searches)

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

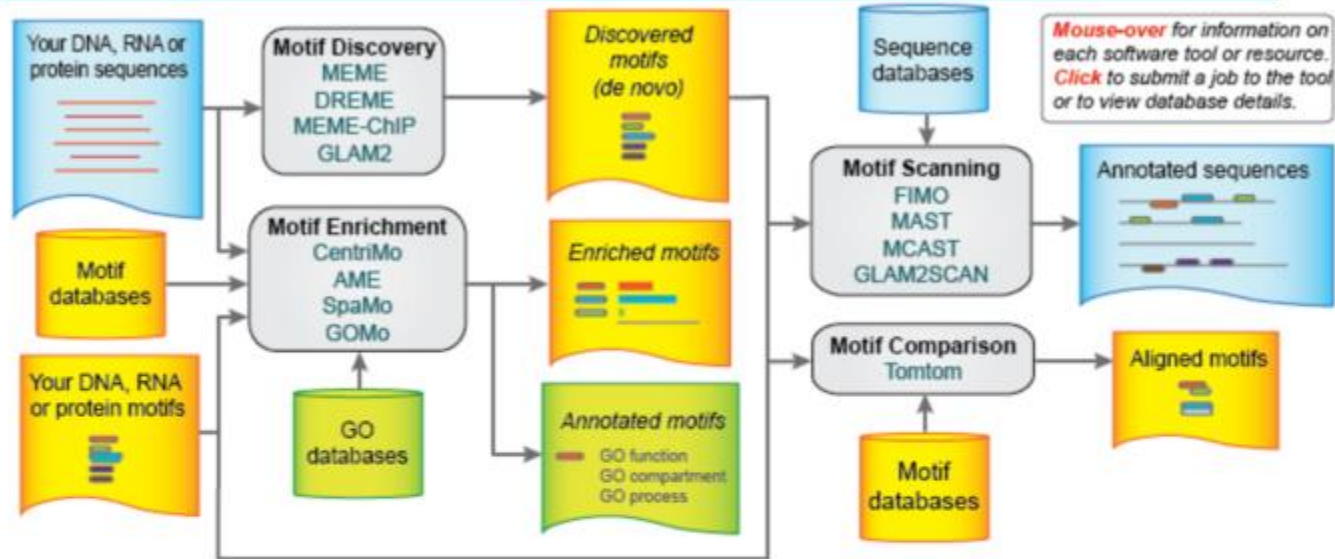
Version 4.3.0

Please send comments and questions to: [meme@nbc.net](mailto:meme@nbc.net)

Powered by Opal

# The MEME Suite

Motif-based sequence analysis tools



 **MEME**  
Multiple Em for Motif Elicitation

 **CentriMo**  
Local Motif Enrichment Analysis

 **FIMO**  
Find Individual Motif Occurrences

 **DREME**  
Discriminative Regular Expression Motif Elicitation

 **AME**  
Analysis of Motif Enrichment

 **MAST**  
Motif Alignment & Search Tool

 **MEME-ChIP**  
Motif Analysis of Large Nucleotide Datasets

 **SpaMo**  
SpaceM Motif Analysis Tool

 **MCAST**  
Motif Cluster Alignment and Search Tool

 **GLAM2**  
Gapped Local Alignment of Motifs

 **GOMo**  
Gene Ontology for Motifs

 **GLAM2Scan**  
Scanning with Gapped Motifs

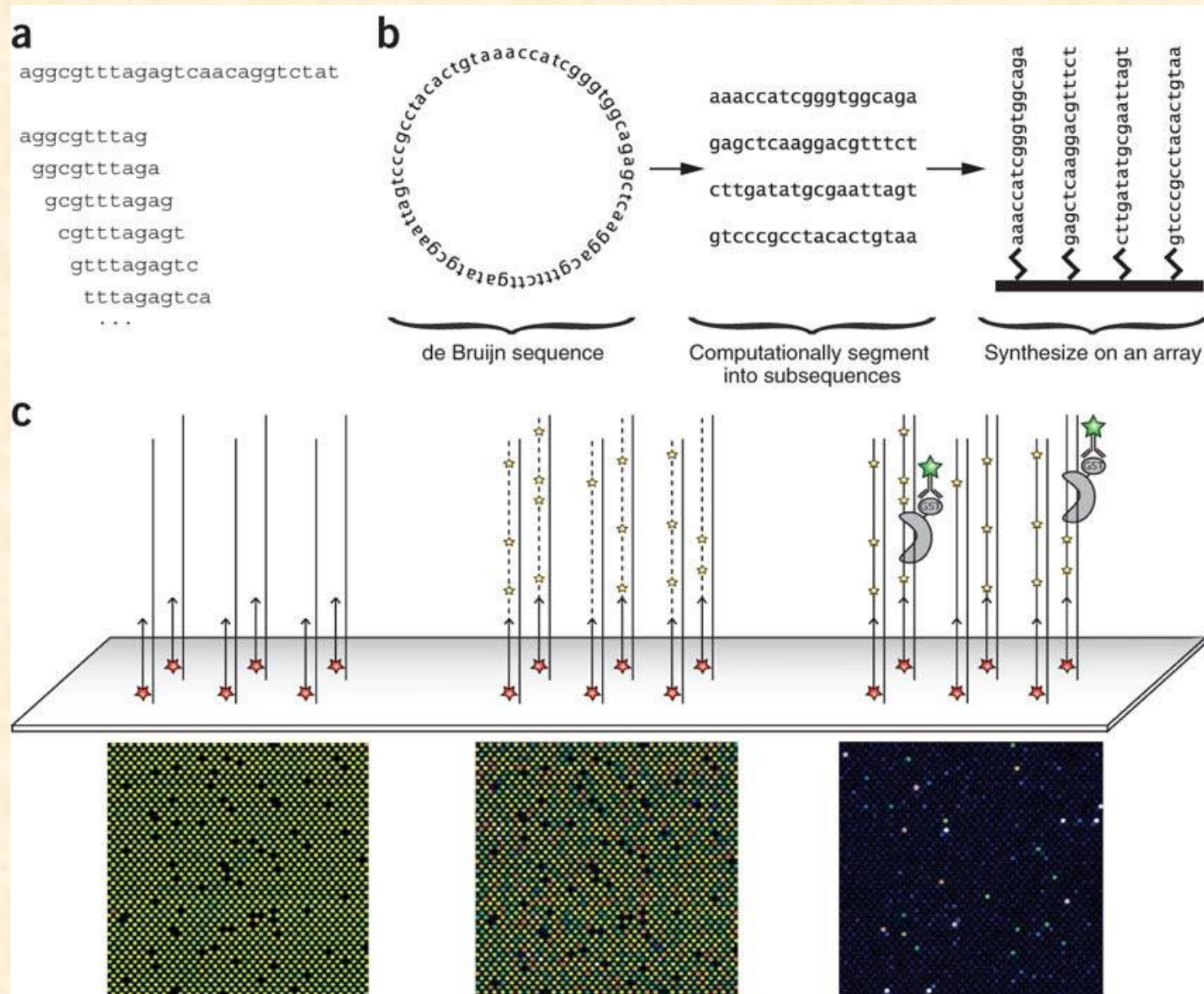
 **Tomtom**  
Motif Comparison Tool

 **GT-Scan**  
Identifying Unique Genomic Targets



**Part II:**  
**Discovering motifs**  
**in high-throughput**  
***in vitro* data**

# Protein binding microarrays



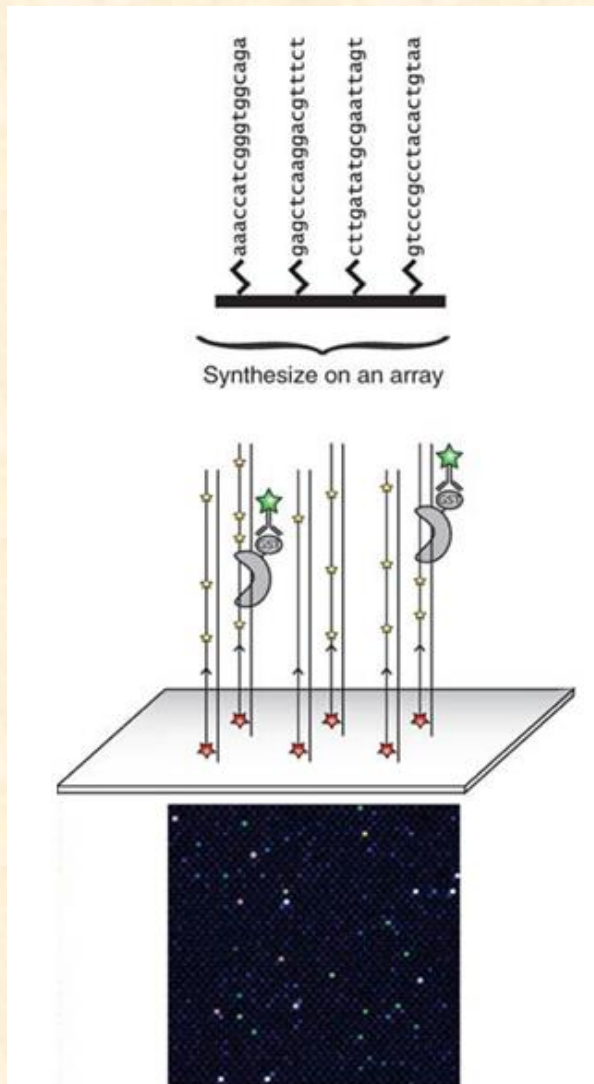
# The data

- Around 41,000 probes, each 36bp-long.
- Plus, corresponding binding intensities.

```
936656.772613291 AAATTTGATACAACCTGCGTCTTGAGCTCTAATTCAC
908113.295168776 TTTTGGATTTAGCCGGGTCAATCCTCCAATATCCTA
897768.671357285 TTTTGGATAAGAATAATTACATAAGAGTGAAAACAT
889610.449301355 ATATATATATTGGTAAAGAGCAATTAATCTTCGCTG
788786.81899799 GGTTTGATACACTCAGTTGTTATAAGAGCGTAGGCA
758147.585093369 GAAGCTTAATTAATCCTGCACTACGCCTGTTTTCA
748029.092176931 TAGTTTGATTGACTGCAATTTGTTTCCCCTGCCTT
740301.561188748 GTTTGGATAGGGGTATGCAGAGTCTCTTTCTCCGTT
...
...
979.99239669736 GAAGGAGCAGCCCCTGAGGGCGGCTCGCATAACTGG
973.610994073221 TCTCTAGGAGTCAGGGCTGTCCTCGCATACCGTGCC
964.88923536786 GAGCGTCAGCGCGGGTCACTTCATGTATGGTCGCC
577.174517674586 AGGGATTCTTGGCTGCGGCCTGTATCGCTGTTGGAA
552.401117184654 GGCGTGGCTTGTTACCTCATGCAGAAGACCGGTAAG
-218.900301178986 CGTGTACTTCAAACGCGGAGTTGGCATATGCGGTCC
-288.772718876298 GGAGGTCAACCACGGTCAGTCCATGTGAAGTCCATA
-2830.13880248726 ACCGGTCGCAGTTGGCCTTCCGCTGATAAGCATTGC
```

(UniPROBE database, Robasky and Bulyk, *NAR* 11)

# Computational challenge



experiment

**AGCTCGAGTAG** 208  
**GTTGATGCATT** 134  
**GCTAGCCAGGA** 10  
**CGGCTAGCGCC** 45  
**GCATCGCTGCG** 65  
**CGATCGCTCGA** 544  
**AGCTCGAGGCG** 414  
**GTTGATGCAGC** 13  
**GCTAGCCAGGC** 51  
**CGGCTAGCGAG** 62  
**GCATCGCTGTT** 102  
**CGATCGCTCGA** 54

Computational analysis



	1	2	3	4	5
A	0.1	0.8	0	0	0
C	0.9	0.1	0.9	0.5	0.1
G	0	0	0.1	0.4	0
T	0	0.1	0	0.1	0.9

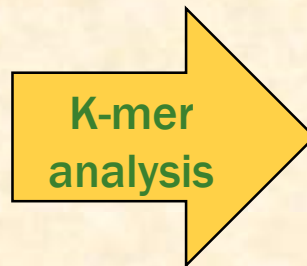
(Berger et al., NBT 2006)

## RAP algorithm

Input: a list of sequences and scores.

Output: a binding model.

<b>AGCTCGAGTAG</b>	<b>10</b>
<b>GTTGATGCATT</b>	<b>5</b>
<b>GCTAGCCAGGA</b>	<b>20</b>
<b>CGGCTAGCGCC</b>	<b>25</b>
<b>GCATCGCTGCG</b>	<b>35</b>
<b>CGATCGCTCGA</b>	<b>50</b>
<b>AGCTCGAGGCG</b>	<b>70</b>
<b>GTTGATGCAGC</b>	<b>30</b>
<b>GCTAGCCAGGC</b>	<b>45</b>
<b>CGGCTAGCGAG</b>	<b>60</b>
<b>GCATCGCTGTT</b>	<b>100</b>
<b>CGATCGCTCGA</b>	<b>40</b>



Ranked 4-mer list

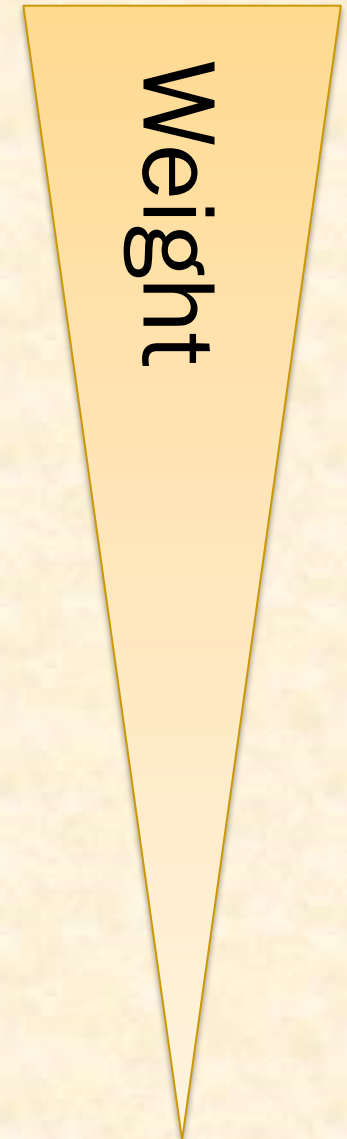
<b>TCGA</b>	<b>42.5</b>
<b>GCTC</b>	<b>42.5</b>
<b>CCAG</b>	<b>32.5</b>
<b>AGCT</b>	<b>40</b>

...

TAATGCGA	150
ATAATGGGG	132
AAATGGCGAAC	98
GATTAATGACGC	95
TAATGTGC	85
AAATGGGGCC	76
TAATGTAC	74
ATATGTAACC	71



				0.95			
					0.95		
			0.9				
		0.85					

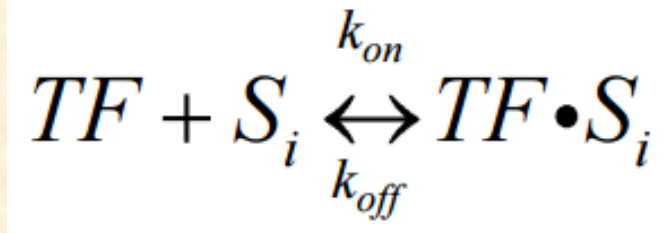




# BEEML-PBM

(Binding Energy Estimation by Maximum Likelihood for PBMs)

- Follows a biomechanical model of TF binding to sequence  $S_i$ :



- $K_{on}$  and  $K_{off}$  = on- and off-rates.

$$K_d = \frac{K_{off}}{K_{on}} = \frac{[S_i][TF]}{[TF \bullet S_i]}$$

# BEEML-PBM model

- Probability of binding in equilibrium is:

$$P(S_i) = \frac{[TF \cdot S_i]}{[S_i] + [TF \cdot S_i]} =$$

- $E_i$  = free energy difference between  $S_i$  and  $S_{ref}$ .

$$\mu = \frac{K_d(S_{ref})}{[TF]}$$

# BEEML-PBM model

- Energy contribution by PWM model  $\varepsilon$ :

$$E_i = \prod_{k=1}^l \varepsilon(S_i(k), k)$$

- $l$  = motif length.
- $S_i(k)$  = letter  $k$  of  $S_i$
- $\varepsilon(b, k)$  = energy contribution of  $b$  in position  $k$ .

# BEEML-PBM model

- Binding probability to double-stranded  $S_i$  :

$$F(S_i) = P(S_i) + (1 - P(S_i))P(S_i')$$

- $S_i'$  = reverse complement of  $S_i$ .

- Binding probability to sequence  $T$ :

$$B(T) = \sum_{i=1}^{|T|-l+1} F(T_{i:i+l-1})$$

- $T_{i:i+l-1}$  = subseq of  $T$  of length  $l$  starting at  $i$ .

# BEEML-PBM model parameters

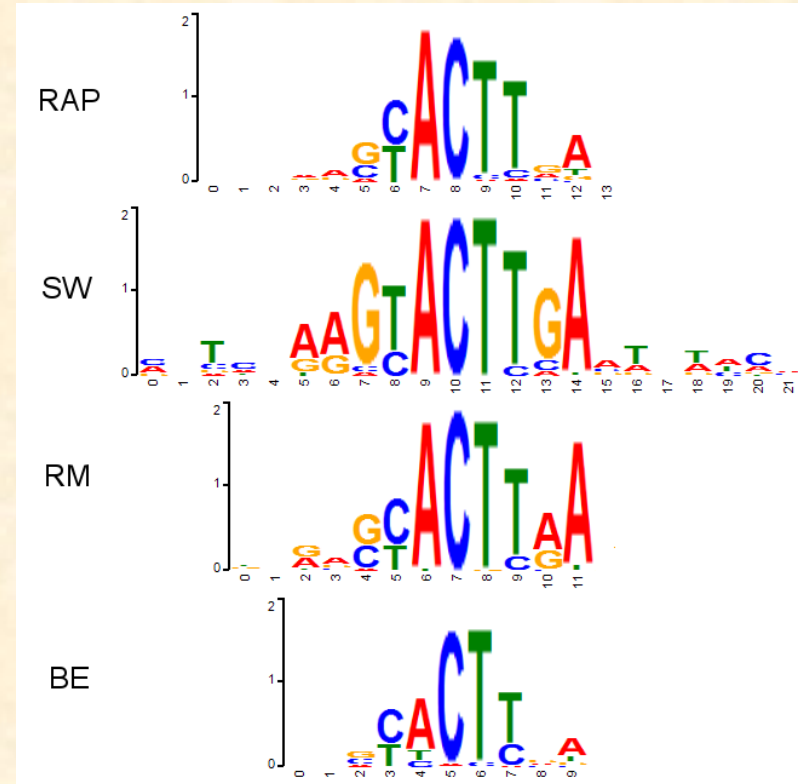
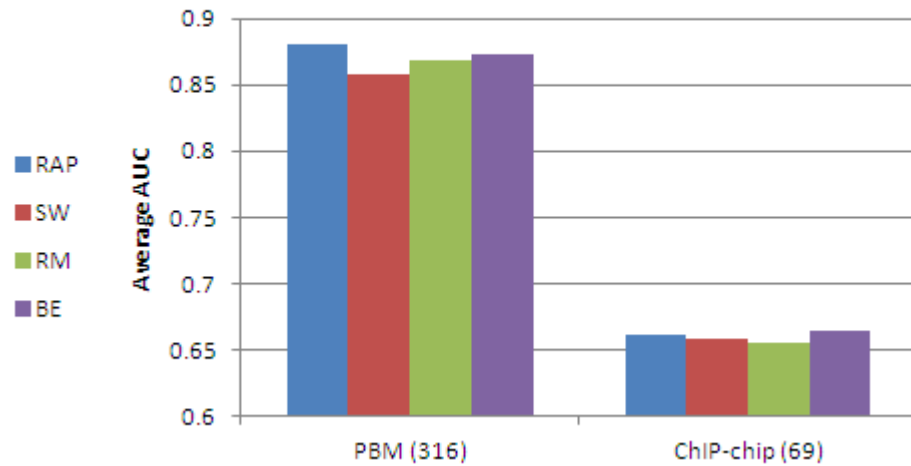
- $\varepsilon = \text{PWM}$ ,  $\mu = \frac{K_d(S_{ref})}{[TF]}$
- Objective function:

$$O(\varepsilon, \mu) = \sum_i (Y_i - a - cB(T_i))^2 + \lambda \sum_{b=A}^T \sum_{k=1}^l \varepsilon(b, k)^2$$

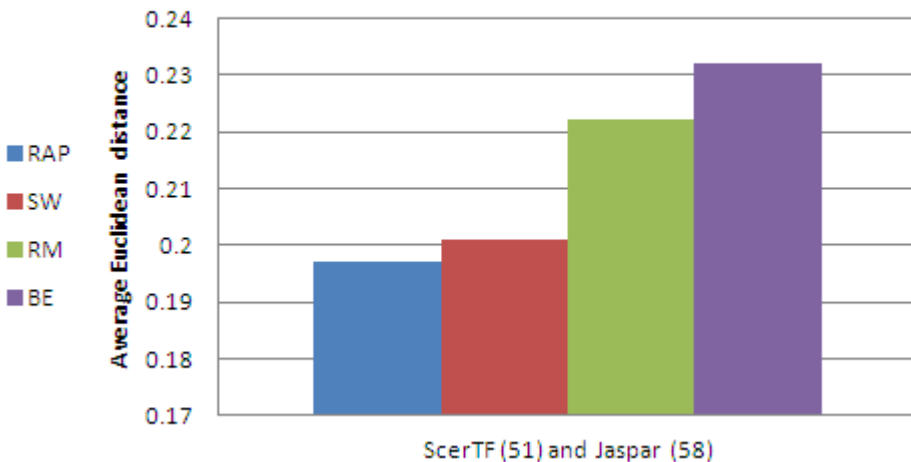
- $Y_i$  = observed binding intensity.
- $T_i$  = sequence of probe  $i$ .
- Minimization by Levenberg-Marquardt algorithm.

# Results

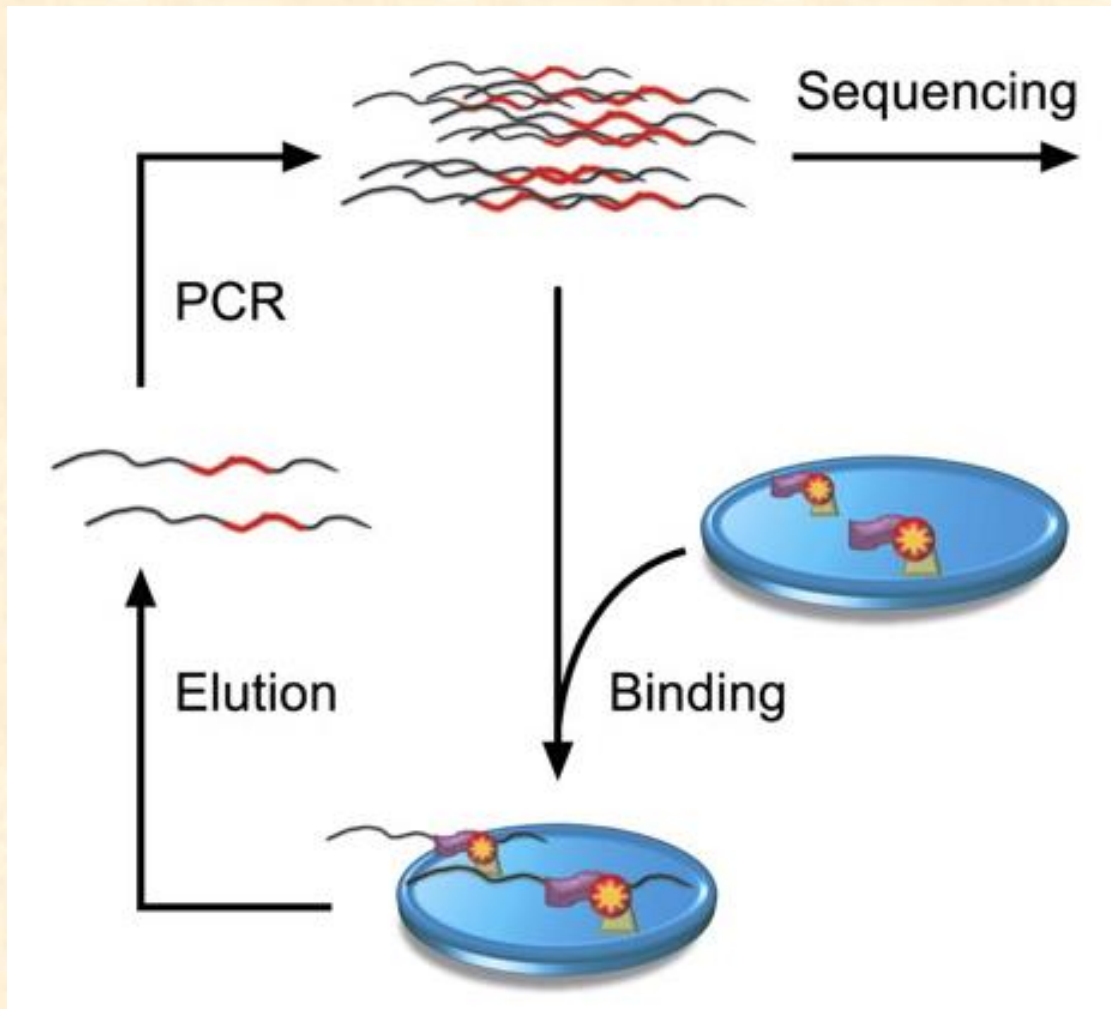
Predicting *in vitro* and *in vivo* binding



Similarity to known motifs



# The next step: High-throughput SELEX



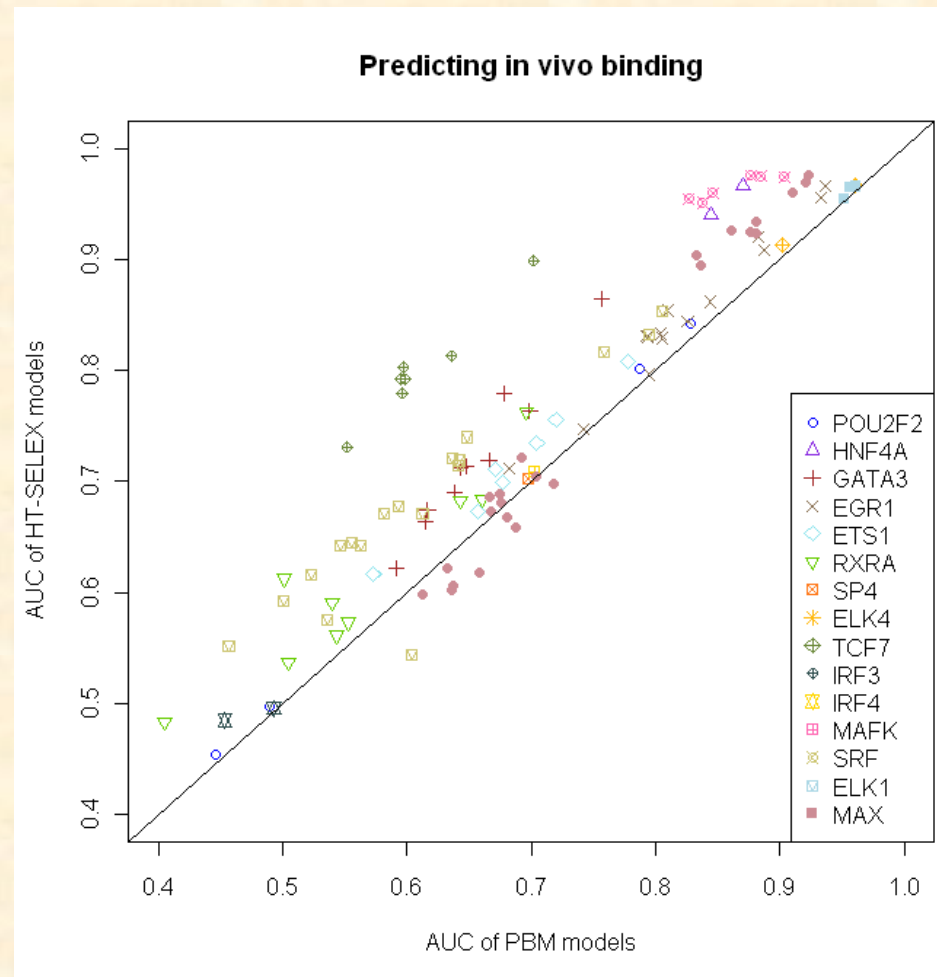
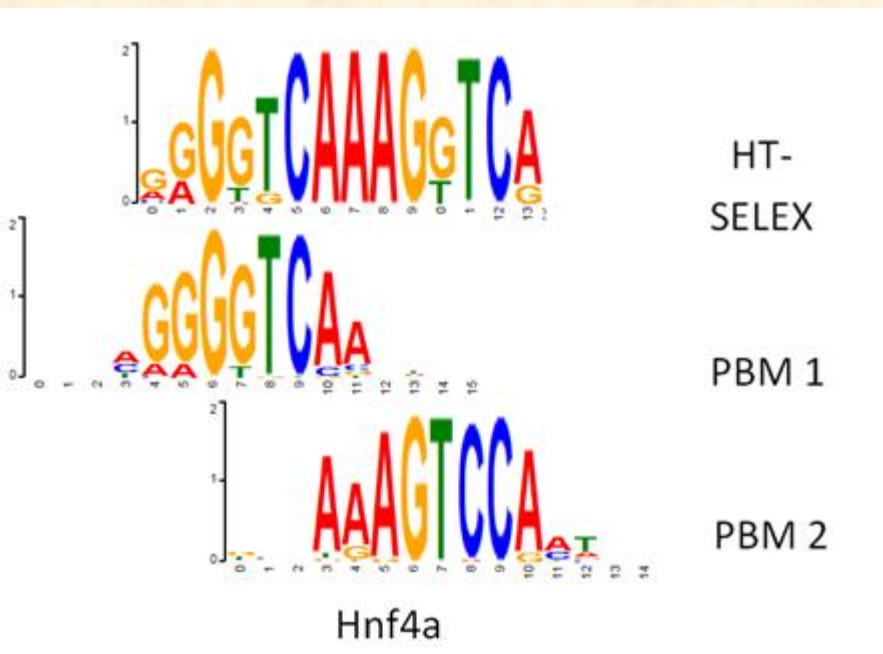
# What happens through the cycles...

Cycle 0	Cycle 1	Cycle 2	Cycle 3
GTCAGTGTTCAGTGACGTG	GTCAGTGTTCAGTGACGTG	GTCAGTGTTCAGTGACGTG	<b>TTC</b> ACTCTCGAGGCTAGC
GCGTGACGTGAAGTCGTA	GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA
GTGCAACGTTGTTTGAGG	GTGCAACGTTGTTTGAGG	GTGCAACGTTGTTTGAGG	GTGCAACGTTGTTTGAGG
GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA
AAGGTGACCGTACGAGTC	AAGGTGACCGTACGAGTC	AAGGTGACCGTACGAGTC	AAGGTGACCGTACGAGTC
GTCCATTCACTGGTGAGT	GTCCATTCACTGGTGAGT	GTCCATTCACTGGTGAGT	GTCCATTCACTGGTGAGT
GTGCAGCGTGAACGTTGG	GTGCAGCGTGAACGTTGG	<b>TTC</b> ACTCTCGAGGCTAGC	<b>TTC</b> ACTCTCGAGGCTAGC
GTGATTTGGAACACCCAG	GTGATTTGGAACACCCAG	GTGATTTGGAACACCCAG	GTGATTTGGAACACCCAG
CCCCACCCACCGGCCTGC	CCCCACCCACCGGCCTGC	CCCCACCCACCGGCCTGC	<b>TTC</b> ACTCTCGAGGCTAGC
ACAGTCAGCCCTAGCACG	<b>TTT</b> GACTCTGCTACGCAT	<b>TTT</b> GACTCTGCTACGCAT	<b>TTT</b> GACTCTGCTACGCAT
CACATACGCTGACTCGTA	CACATACGCTGACTCGTA	GCGA <b>TTAACTCT</b> GTGAAA	GCGA <b>TTAACTCT</b> GTGAAA
<b>TTT</b> GACTCTGCTACGCAT	<b>TTT</b> GACTCTGCTACGCAT	<b>TTT</b> GACTCTGCTACGCAT	<b>TTT</b> GACTCTGCTACGCAT
CGATCGATCAGGCTAGCT	CGATCGATCAGGCTAGCT	CGATCGATCAGGCTAGCT	GCGA <b>TTAACTCT</b> GTGAAA
...	...	...	...
...	...	...	...
...	...	...	...
<b>TTC</b> ACTCTCGAGGCTAGC	<b>TTC</b> ACTCTCGAGGCTAGC	<b>TTC</b> ACTCTCGAGGCTAGC	<b>TTC</b> ACTCTCGAGGCTAGC



# Comparison of HT-SELEX to PBM

(Orenstein and Shamir 2014)



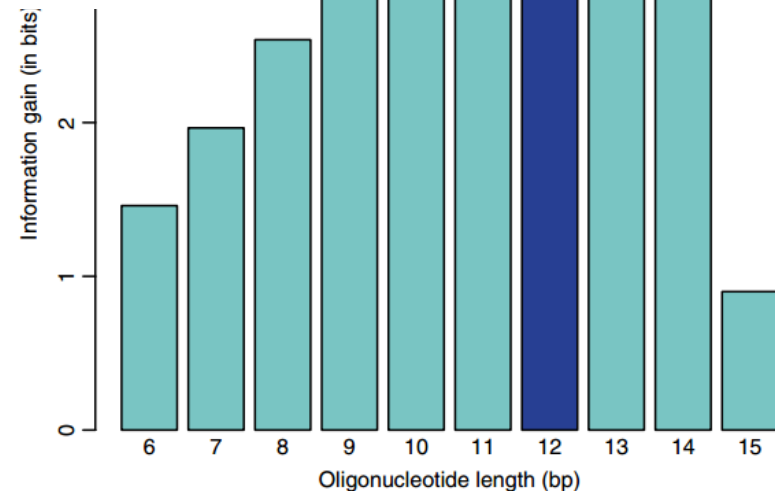
# SELEX-seq

(Riley et al. 2014)

- Find  $k$  by KL-divergence score:

$$D_{CGKL}(K) = \sum_{i \in S_{100}(K)} \left( P_2(w_i) \log \frac{P_2(w_i)}{P_{MM}(w_i)} \right) + \left[ 1 - \sum_{i \in S_{100}(K)} P_2(w_i) \right] \log \left( \frac{1 - \sum_{i \in S_{100}(K)} P_2(w_i)}{1 - \sum_{i \in S_{100}(K)} P_{MM}(w_i)} \right)$$

- $S_{100}(K)$ : count  $\geq 100$

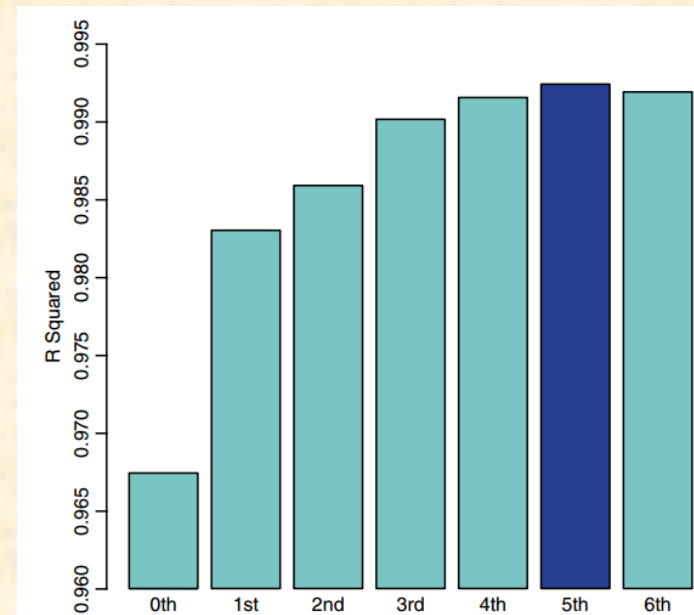


# SELEX-seq

- Scoring k-mers:

$$score_i(w) = \sqrt[i]{freq_i(w)/freq_0(w)}$$

- Estimated frequency of cycle 0 by 5-th order Markov model.



# Open challenges

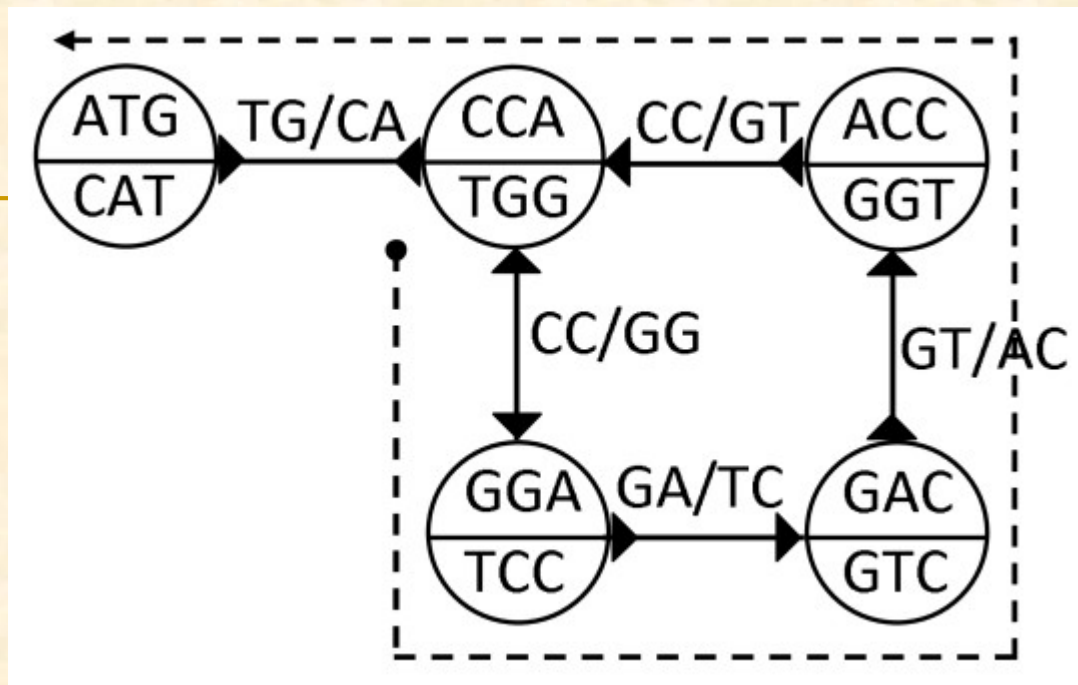
- HT-SELEX
  - Choice of cycle.
  - Utilization of several cycles together.
  - Technological biases.
- *In vivo* binding prediction:
  - Improve prediction using *in vitro* models.

# Summary

- PWM model for protein-DNA binding.
- MEME: find motifs in genomic sequences.
- RAP, BEEML-PBM: infer motifs from PBM data.
- New computational challenges with HT-SELEX.



# Bi-directed de Bruijn graphs for assembly and design



# De Bruijn sequence of order $k$

- A cyclic sequence s.t. every possible  $k$ -mer appears in it exactly once.
- The most compact sequence to cover all  $k$ -mers.
- Length:  $|\Sigma|^k$
- DNA alphabet = {A, C, G, T}
- dB sequence of order  $k=2$ :  
ACGTAGAATTGGCCTC



# De Bruijn graph of order k

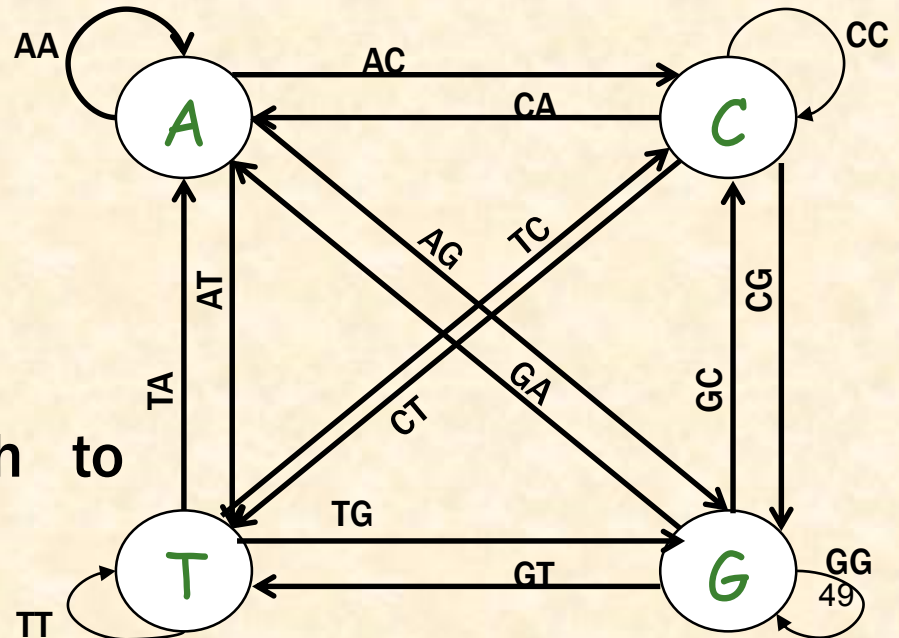
- Vertex = k-mer  $(x_1, \dots, x_k)$      $|V| = |\Sigma|^k$
- Directed edge for overlap of (k-1):  
 $(x_1, x_2, \dots, x_k) \rightarrow (x_2, \dots, x_k, x_{k+1})$      $|E| = |\Sigma|^{k+1}$
- $\forall v \ d^+(v) = d^-(v) = |\Sigma|$

outdegree

indegree

The graph is **strongly connected**  
and each vertex is **balanced**

→ **Euler tour** in a de Bruijn graph to  
generate de Bruijn sequence.



# Assembly by de Bruijn graph

- Construct de Bruijn graph according to k-mer reads.
- [Edge weights by read scores.]
- Find an Euler tour in the graph.
- Is it always possible?

# Chinese Postman Problem

- Find the shortest path that traverses each edge at least once.
- Solved by minimum-cost flow:
  - $V^- = \{v \mid d^-(v) > d^+(v)\}, V^+ = \{\dots\}$
  - $c(s,v) = d^-(v) - d^+(v), c(v,t) = d^+(v) - d^-(v), c(u,v) = \infty$
  - $a(u,v) = w(u,v), a(s,v) = a(v,t) = 0$
- Add  $f(u,v)$  copies of edge  $(u,v)$ .

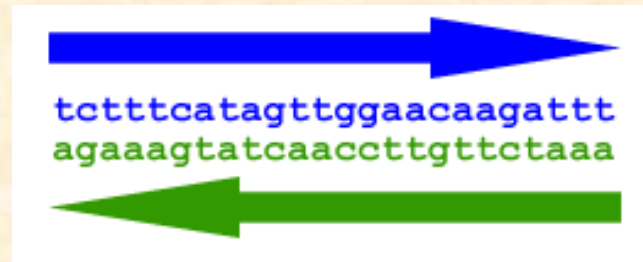
# Running time of CPP

- Successive shortest path:
  - Continue till no augmenting path exists:
    - Send flow through the cheapest path.
    - Update residual network (negative costs).
  - #iterations  $\leq \max F = \frac{1}{2} \sum_v |d^+(v) - d^-(v)|$
  - Finding cheapest augmenting path:
    - Bellman-Ford:  $O(|V| |E|)$
    - Dijkstra (Johnson's):  $O(|V| \log |V|)$
    - Integral and bounded costs:  $O(|V| + |E|)$

# Room for improvement

- **Reverse Complementary** of double-stranded DNA:

A ⇔ T      C ⇔ G



- In a double stranded sequence, when a k-mer is present, so is its RC.

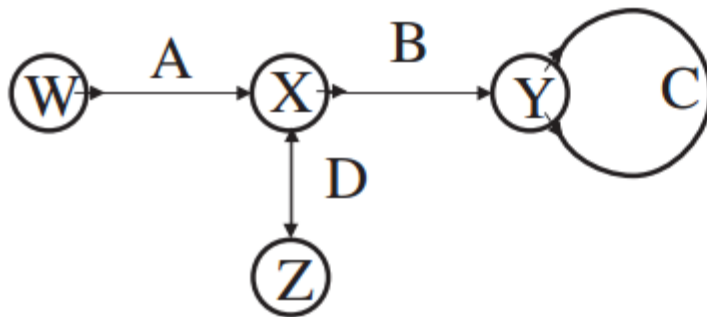
Goal: find the shortest sequence  
to cover all the k-mers.

# Outline

- Bi-directed de Bruijn graphs.
- Solution to assembly: LP and bipartite.
- Solution to design: RCdB.
- Summary and open problems.

# Bi-directed graphs

- Each edge endpoint has an orientation.
- Incidence matrix:  $V \times E \rightarrow \{-2, -1, 0, 1, 2\}$
- $d^-(v) = -\sum_{\{e \in E \mid I(v,e) < 0\}} I(v,e)$ ,  $d^+(v) = +\sum \dots$
- $\text{bal}(v) = d^+(v) - d^-(v) = \sum_e I(v,e)$



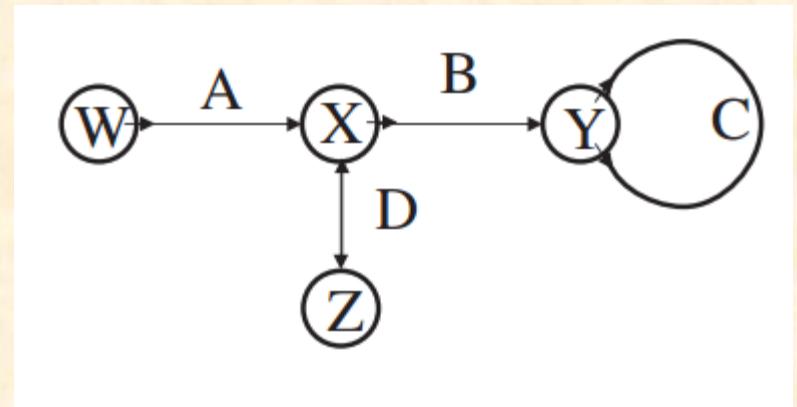
A is negative-incident to X

B is positive-incident to X

	A	B	C	D
W	1			
X	-1	1		-1
Y		-1	2	
Z				-1

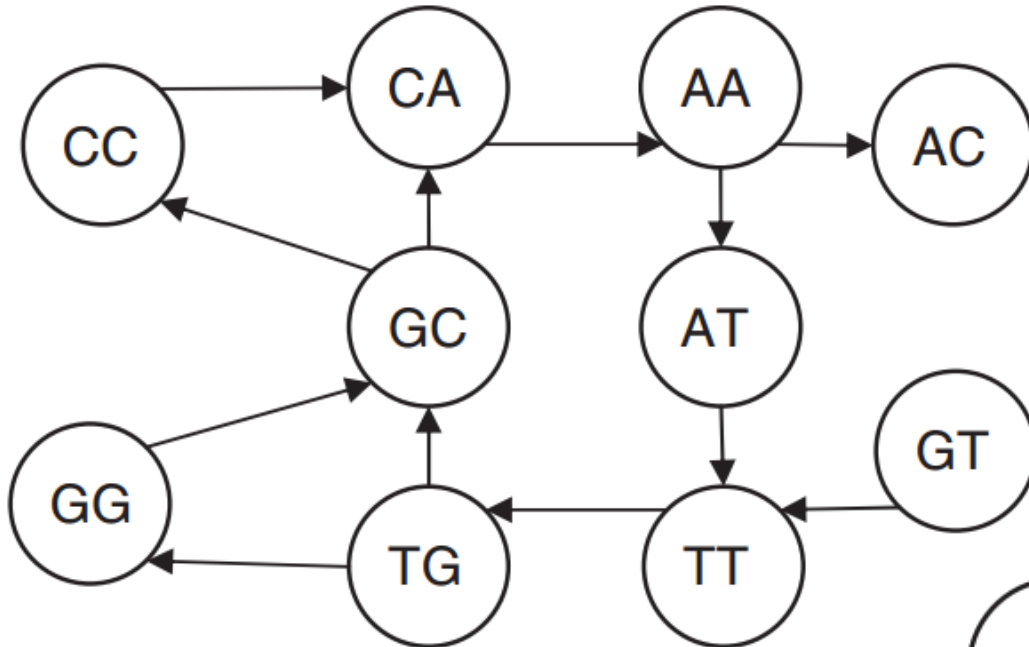
# Walks in bi-directed graphs

- $v_1, e_1, \dots, v_{k-1}, e_{k-1}, v_k$ 
  - $e_i$  incident to  $x_i$  and  $x_{i+1}$
  - $e_{i-1}$  and  $e_i$  have opposite orientation at  $x_i$
- $W, A, X, B, Y, C, Y, B, X, D, Z.$
- $W, A, X, D, Z.$





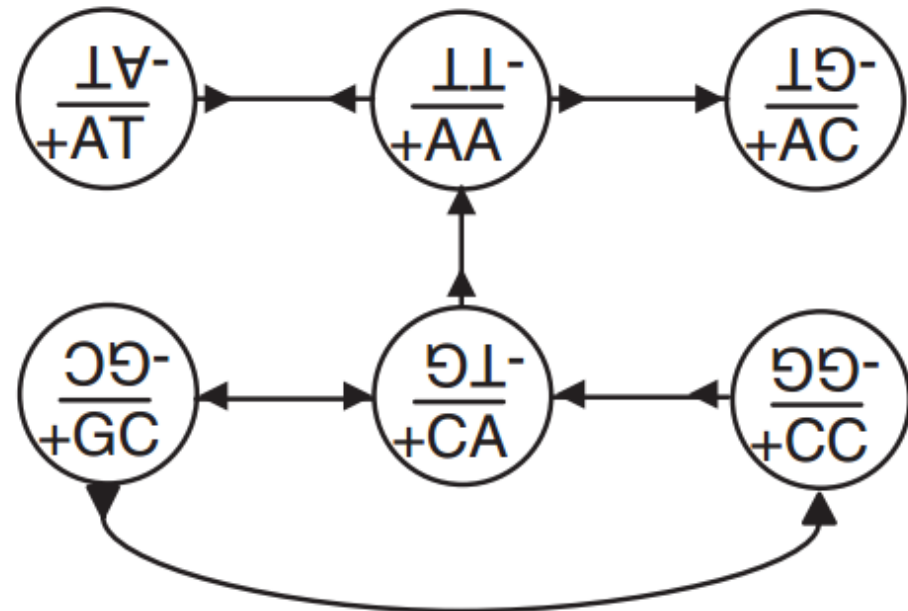
# Bi-directed de Bruijn graphs



Bi-directed overlap:

- e pos to x, neg to y, p(x) overlaps p(y)
- e pos to x, pos to y, p(x) overlaps n(y)

- e neg to x, neg to y, n(x) overlaps p(y)
- e neg to x, pos to y, n(x) overlaps n(y)



# Linear Programming to solve Chinese postman problem

Find  $f: E \rightarrow \mathbb{N}$  (#copies to add of each edge)

$$\text{minimize } \sum w(e)f(e)$$

$$\text{subject to } \sum_e I^G(x, e)f(e) = -bal^G(x) \text{ for each vertex } x$$

$$f(e) \geq 0 \quad \text{for each edge } e$$

# Solving the LP

- LP solutions to binet matrices are half-integral (Appa and Kotnyek, 2006).
- A reduction to undirected graphs (totally unimodular matrices) has integral solutions (Hochbaum, 2004).
- The solution is a 2-approximation.

# Shortest path in bidirected

- Modified Dijkstra:
  - updates neighbors consistent with a walk.
  - $O((|V| + |E|) \log |V|)$
- Shortest paths between unbalanced:
  - $V^- = \{v \mid d^-(v) > d^+(v)\}, V^+ = \{\dots\}$
  - Run  $P = \min\{|V^+|, |V^-|\}$  times.
- Now, how can we solve assembly problem?

# Solution by bipartite matching

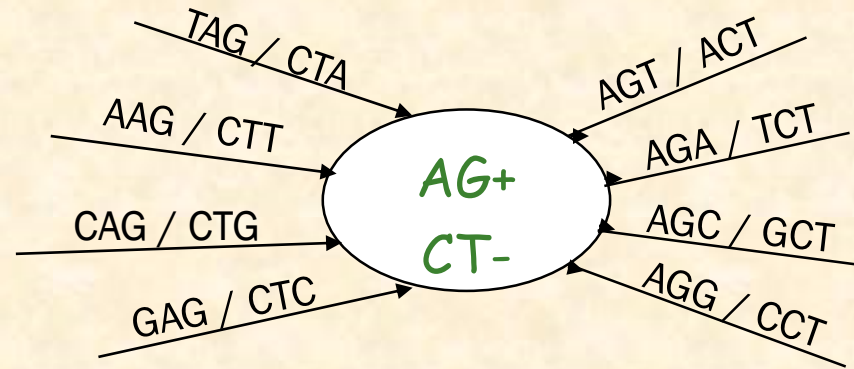
- Construct bipartite graph:
  - $U^- = U_{v \in V^-} \{v^{(1)}, \dots, v^{(d^-(v) > d^+(v))}\}$ ,  $U^+ = \{\dots\}$
  - $w(u^{(i)}, v^{(j)}) = \text{shortest\_path}(u, v)$
- For each edge in the matching:
  - Add edges of the shortest path.
- Find Euler tour in original graph.

# Running time

- Modified Dijkstra:
  - $O((|V| + |E|) \log |V|)$
  - Run  $P = \min \{|V^+|, |V^-|\}$  times.
- Bipartite matching:
  - Hungarian method
  - $O(|V|^3) = O((\sum_v |d^+(v) - d^-(v)|)^3)$

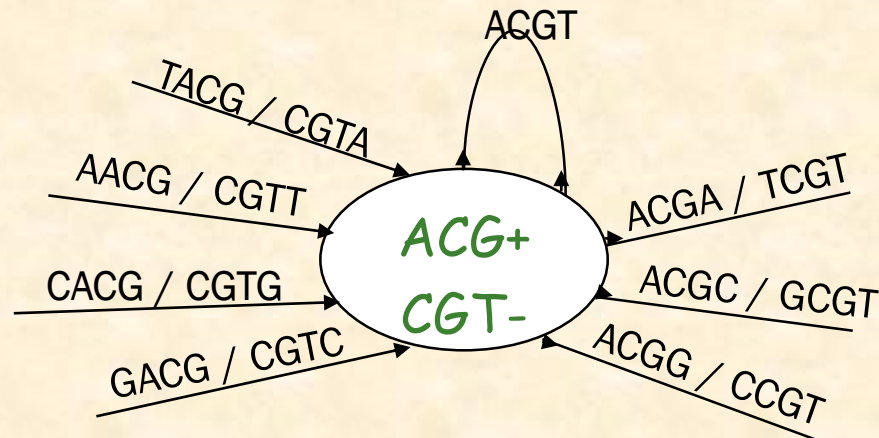
# Reverse complement de Bruijn sequence

- For odd k, the graph is balanced.



- For even k, the graph is unbalanced.

e pos to x, pos to y,  
p(x) overlaps n(y)



# Solutions

- Can be solved by linear programming.
- Or by maximum weighted matching:
  - Shortest paths = label overlaps.
  - #vertices =  $4^{k/2}$
  - Run time =  $O(4^{3k/2})$
- Integer weights => improved run-time.



# Computational results

## Reverse complement de Bruijn sequence

K	2	4	6	8	10	12	14
De Bruijn	16	256	4,096	65,536	1,048,576	16,777,216	268,435,456
Optimal	10	142	2,140	33,262	526,816	8,400,772	134,274,844
Lower bound	10	136	2,080	32,896	524,800	8,390,656	134,225,920
Saving factor	1.6	1.8	1.91	1.97	1.99	1.997	1.999

# Open questions

1. A sequence with improved coverage of gapped k-mers.

2. What is the number of sequences?

$$\frac{(k!)^{k^{n-1}}}{k^n}$$

3. Probe design that minimizes the number of probes.

4. Linear time algorithm.

5. Closed formula for length.

# References

1. Maximum likelihood of genome assembly.  
(Medvedev and Brudno, *Journal of Computational Biology* 2009)
2. An efficient algorithm for Chinese postman walk on bi-directed de Bruijn graphs.  
(Kundet, Rajasekaran and Dinh, *Discrete Math. Algorithm. Appl.* 2012)
3. Design of shortest double-stranded DNA sequence covering all k-mers.  
(Orenstein and Shamir, *Bioinformatics* 2013)