

Bioinformatics

Supplementary Slides



UNIVERSITY OF
CAMBRIDGE

Computer Laboratory

Computer Science Tripos Part II

Pietro Lio

Michaelmas 2012

Bioinformatics Examples

Some of the images in these lectures are from several sources that have welcomed instructors to use them for educational purposes (Felsenstein, Moran, Gerstein, Yang, Arkin, Leibler, Batzoglou, Pevzner, Nussinov).

Gap penalty=-1; match=+2; mismatch=-1		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0						
C 1							
A 2							
T 3							
G 4							
T 5							

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1					
C 1							
A 2							
T 3							
G 4							
T 5							



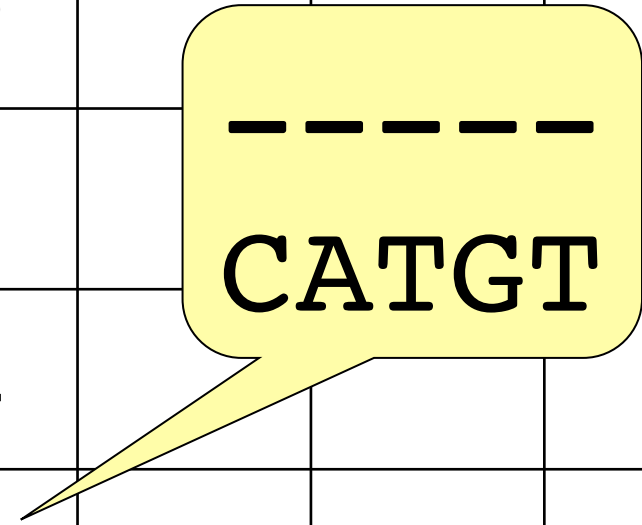
A
-

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1							
A 2							
T 3							
G 4							
T 5							



ACGCTG

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1						
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						



		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1					
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

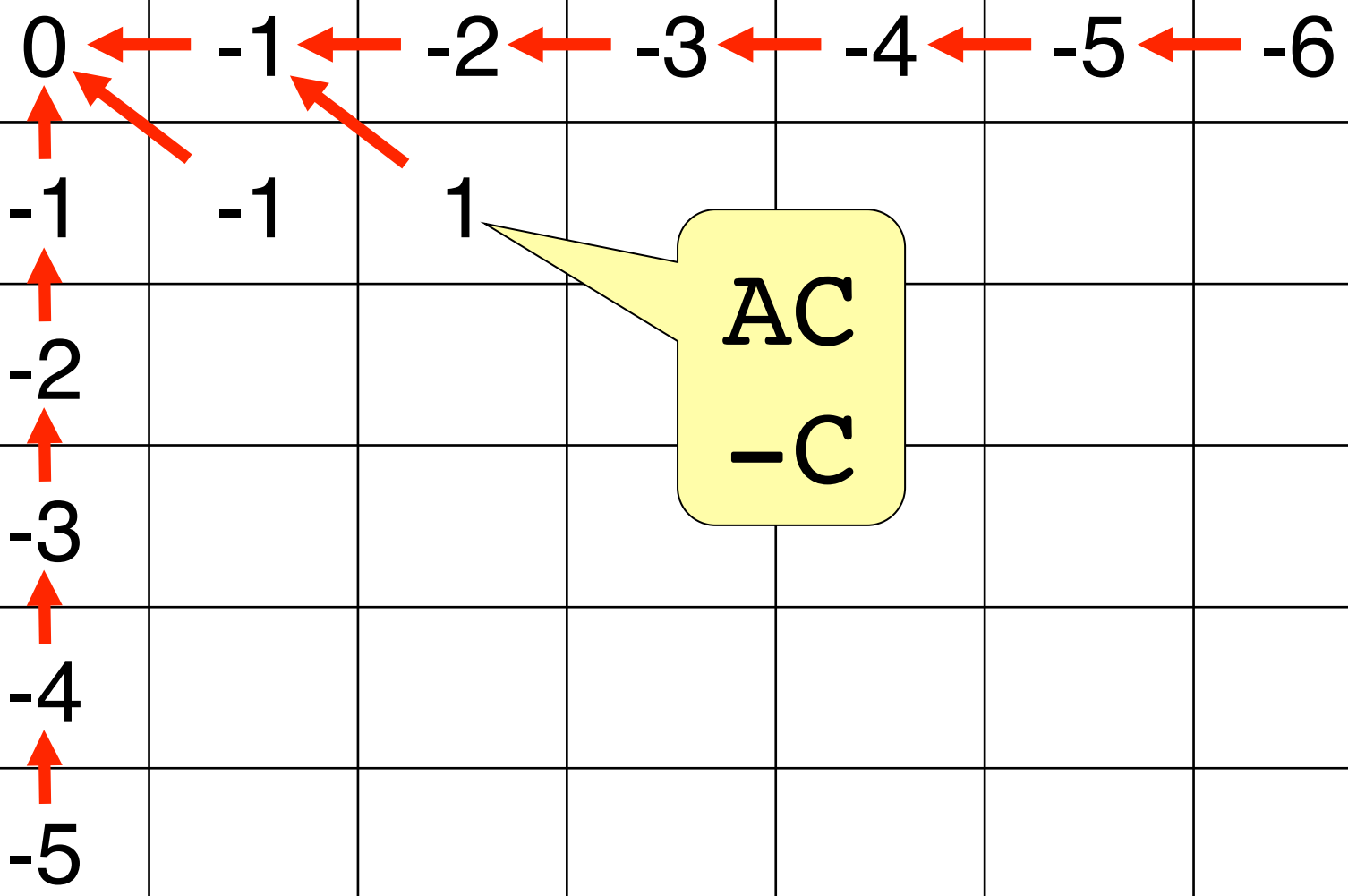
← -1 ← -2 ← -3 ← -4 ← -5 ← -6

↑
↑
↑
↑
↑

A
C

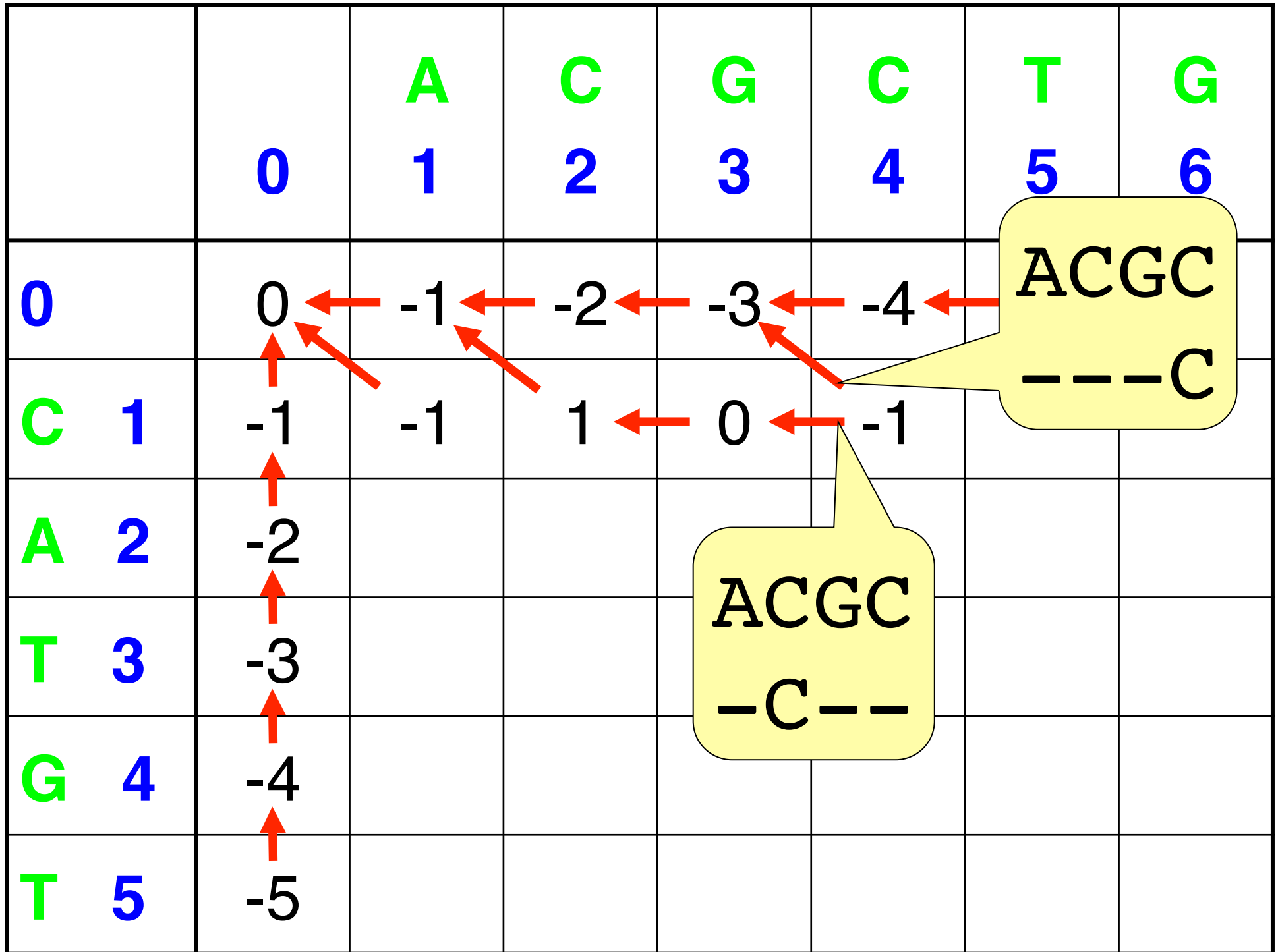
Gap penalty=-1; match=+2; mismatch=-1		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1				
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

AC
-C



		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0			
A 2	-2						
T 3	-3						
G 4	-4						
T 5	-5						

ACG
-C-



		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0			
T 3	-3						
G 4	-4						
T 5	-5						

ACG
-CA

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0	-1	-2	-3
T 3	-3	0	0	-1	-1	1	0
G 4	-4	-1	-1	2	1	0	3
T 5	-5	-2	-2	1	1	3	2

The table shows the following alignment path indicated by red arrows:

- Start at (5,6) with value 2.
- Move left to (5,5) with value 3.
- Move up to (4,5) with value 0.
- Move left to (4,4) with value 1.
- Move up to (3,4) with value 2.
- Move left to (3,3) with value 1.
- Move up to (2,3) with value -1.
- Move left to (2,2) with value -1.
- Move up to (1,2) with value 1.
- Move left to (1,1) with value -1.
- Move up to (0,1) with value -1.
- Move left to (0,0) with value 0.

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1	-2	-3	-4	-5	-6
C 1	-1	-1	1	0	-1	-2	-3
A 2	-2	1	0	0	-1	-2	-3
T 3	-3	0	0	-1	-1	1	0
G 4	-4	-1	-1	2	1	0	3
T 5	-5	-2	-2	1	1	3	2

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

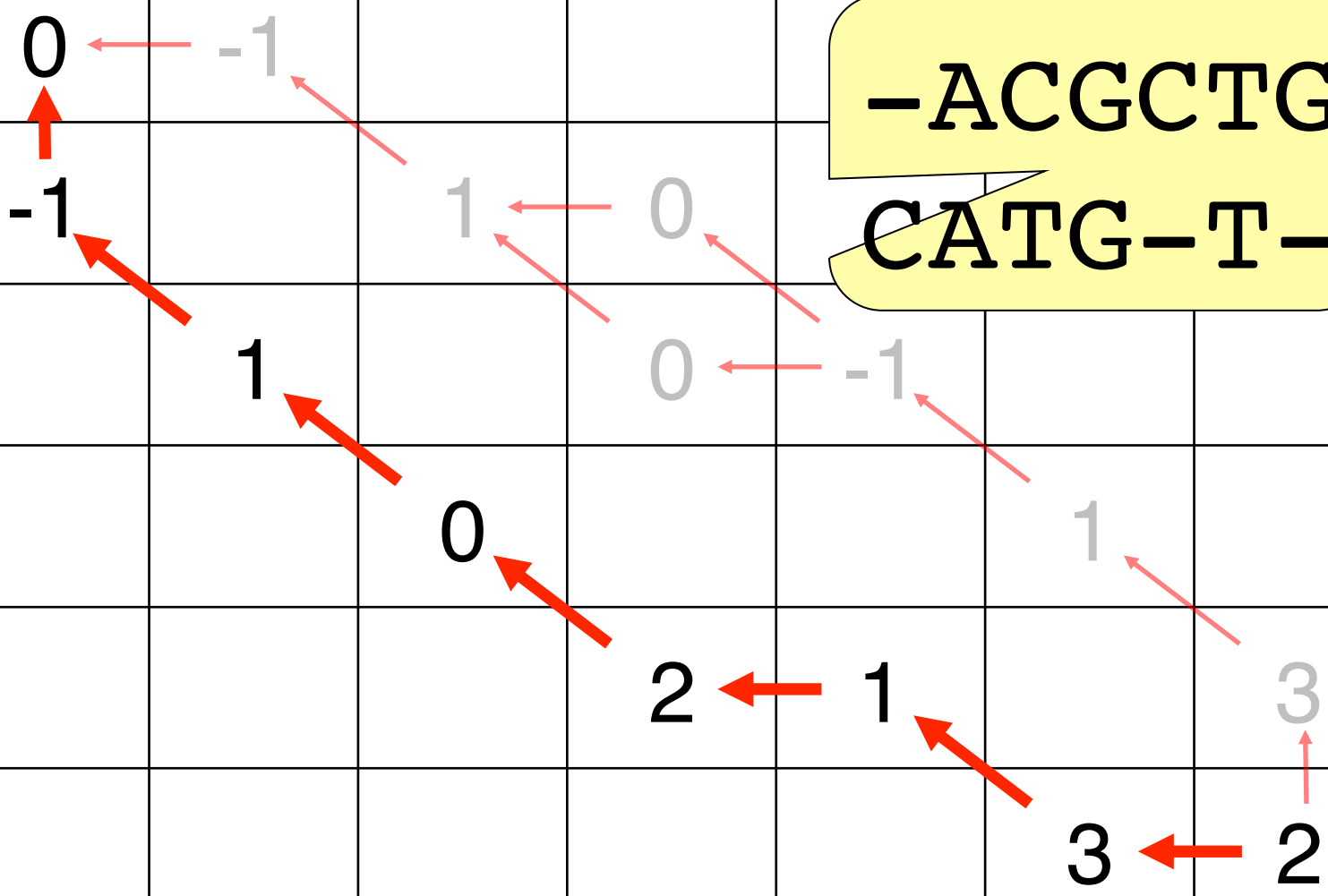
ACGCTG-
-C-ATGT

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

ACGCTG-
-CA-TGT

		A	C	G	C	T	G
	0	1	2	3	4	5	6
0	0	-1					
C 1	-1		1	0			
A 2		1		0	-1		
T 3			0			1	
G 4				2	1		3
T 5						3	2

-ACGCTG
CATG-T-



Local Alignment Example

$y = \text{TAATA}$
 $x = \text{TACTAA}$

$y \backslash x$		A	T	C	T	A	A
0	0	0	0	0	0	0	0
T 1	0						
A 2	0						
A 3	0						
T 4	0						
A 5	0						

Local Alignment Example

$y = \text{TAATA}$
 $x = \text{TACTAA}$

$y \backslash x$		T	A	C	T	A	A
	0	1	2	3	4	5	6
0	0	0	0	0	0	0	0
T 1	0	1	0	0	1	0	0
A 2	0	0	2	0	0	2	1
A 3	0						
T 4	0						
A 5	0						

Local Alignment Example

$y = \text{TAATA}$
 $x = \text{TACTAA}$

$y \backslash x$		T	A	C	T	A	A
	0	1	2	3	4	5	6
0	0	0	0	0	0	0	0
T 1	0	1	0	0	1	0	0
A 2	0	0	2	0	0	2	1
A 3	0	0	1	1	0	1	3
T 4	0	0	0	0	2	0	1
A 5	0	0	1	0	0	3	1

Local Alignment Example

$y = \text{TAATA}$
 $x = \text{TACTAA}$

$y \backslash x$	0	T	A	C	T	A	A
0	0	0	0	0	0	0	0
T 1	0	1	0	0	1	0	0
A 2	0	0	2	0	0	2	1
A 3	0	0	1	1	0	1	3
T 4	0	0	0	0	2	0	1
A 5	0	0	1	0	0	3	1

The table shows a dynamic programming matrix for local sequence alignment. The sequence $y = \text{TAATA}$ is aligned with $x = \text{TACTAA}$. The value 3 in the cell at row 5, column 6 is circled in red, indicating the maximum local alignment score. Red arrows point to the cells (1,1), (2,2), (3,3), (4,4), and (5,6), showing the path of the local alignment. Black arrows indicate the backpointers for each cell, showing the direction of the maximum score transition.

Local Alignment Example

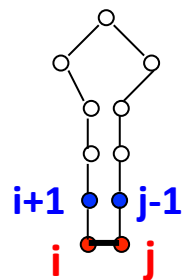
$y =$ TAATA
 $x =$ TACTAA

$y \backslash x$		T	A	C	T	A	A
	0	1	2	3	4	5	6
0	0	0	0	0	0	0	0
T 1	0	1	0	0	1	0	0
A 2	0	0	2	0	0	2	1
A 3	0	0	1	1	0	1	3
T 4	0	0	0	0	2	0	1
A 5	0	0	1	0	0	3	1

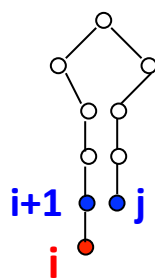
RNA Secondary Structure

- Secondary Structure :
 - Set of paired positions on interval $[i,j]$
 - This tells which bases are paired in the subsequence from x_i to x_j
- Every **optimal structure** can be built by extending **optimal substructures**.
- Suppose we know all **optimal substructures** of length less than $j-i+1$.
 The optimal substructure for $[i,j]$ must be formed in one of four ways:
 1. i,j paired
 2. i unpaired
 3. j unpaired
 4. combining two substructures

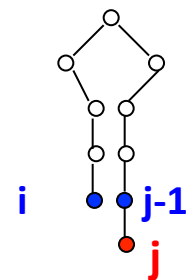
Note that each of these consists of extending or joining substructures of length less than $j-i+1$.



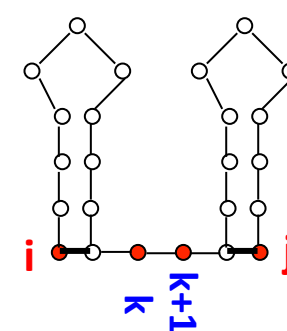
i,j pair



i unpaired



j unpaired
23



bifurcation

The Nussinov Folding Algorithm

Example: GGGAAUCC

$\gamma(i,j)$ is the maximum number of base pairs in segment $[i,j]$

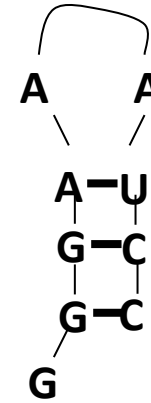
Initialisation $\gamma(i, i-1) = 0$ & $\gamma(i, i) = 0$

Starting with all subsequences of length 2, to length L :

$\gamma(i, j) =$

$$\max \left\{ \begin{array}{l} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{array} \right.$$

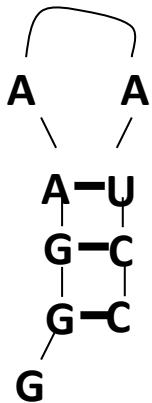
Where $\delta(i,j) = 1$ if x_i and x_j are a complementary base pair, and $\delta(i,j) = 0$, otherwise.



		j →								
		G	G	G	A	A	A	U	C	C
GGGAAUCC ↓ i	0									
	0	0								
		0	0							
			0	0						
				0	0					
					0	0				
						0	0			
							0	0		
								0	0	

Nussinov Folding Algorithm: After scores for subsequences of length 2

$$\gamma(i, j) = \max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



j →

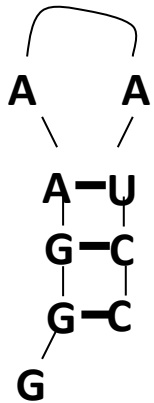
	G	G	G	A	A	A	U	C	C
G	0	0							
G	0	0	0						
G		0	0	0					
G			0	0	0				
A				0	0	1			
A					0	0	0		
A						0	0	0	
U							0	0	0
C								0	0

i ↓

Nussinov Folding Algorithm: After scores for subsequences of length 3

$$\gamma(i, j) =$$

$$\max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



j →

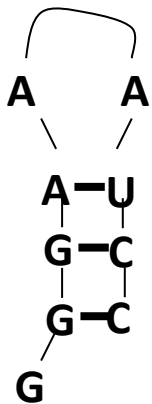
	G	G	G	A	A	A	U	C	C
G	0	0	0						
G	0	0	0	0					
G		0	0	0	0				
A			0	0	0	1			
A				0	0	1	0		
A					0	0	0	0	
U						0	0	0	
U							0	0	
C								0	0
C									0

i ↓

Nussinov Folding Algorithm

After scores for subsequences of length 4

$$\gamma(i, j) = \max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



j →

	G	G	G	A	A	A	U	C	C
G G G A A U C	0	0	0	0					
	0	0	0	0	0				
		0	0	0	0	0			
			0	0	0	0	1		
				0	0	0	1	1	
					0	0	1	1	1
						0	0	0	0
							0	0	0
								0	0

i ↓

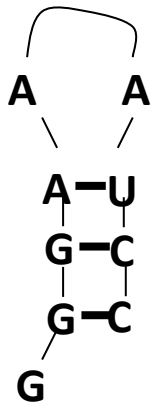
Two optimal substructures for same subsequence

Nussinov Folding Algorithm

After scores for subsequences of length 5

$$\gamma(i, j) =$$

$$\max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



j →

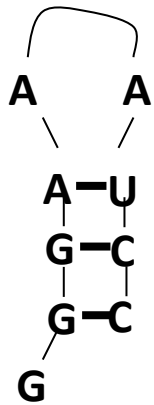
	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0				
G	0	0	0	0	0	0			
G		0	0	0	0	0	1		
G			0	0	0	0	1	1	
A				0	0	0	1	1	1
A					0	0	1	1	1
A						0	0	0	0
U							0	0	0
C								0	0
C									0

i ↓

Nussinov Folding Algorithm

After scores for subsequences of length 6

$$\gamma(i, j) = \max \left\{ \begin{array}{l} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{array} \right.$$

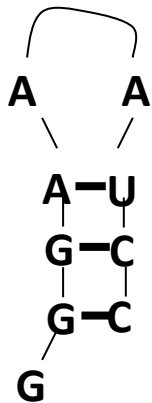


		j →								
		G	G	G	A	A	A	U	C	C
GGGAUCC ↓ i	0	0	0	0	0	0	0			
	0	0	0	0	0	0	0	1		
			0	0	0	0	0	1	2	
				0	0	0	0	1	1	1
					0	0	0	1	1	1
						0	0	1	1	1
							0	0	0	0
								0	0	0
									0	0

Nussinov Folding Algorithm

After scores for subsequences of length 7

$$\gamma(i, j) = \max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



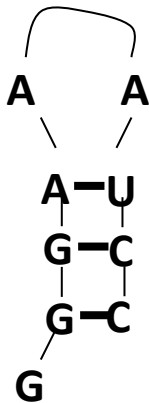
		j →								
		G	G	G	A	A	A	U	C	C
GGGAUUC ↓ i	0	0	0	0	0	0	0	1		
	0	0	0	0	0	0	0	1	2	
			0	0	0	0	0	1	2	2
				0	0	0	0	1	1	1
					0	0	0	1	1	1
						0	0	1	1	1
							0	0	0	0
								0	0	0
									0	0

Nussinov Folding Algorithm

After scores for subsequences of length 8

$$\gamma(i, j) =$$

$$\max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$

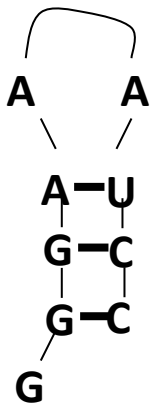


		j →								
		G	G	G	A	A	A	U	C	C
GGGA AU C ↓ i	0	0	0	0	0	0	1	2		
	0	0	0	0	0	0	1	2	3	
		0	0	0	0	0	1	2	2	
			0	0	0	0	1	1	1	
				0	0	0	1	1	1	
					0	0	1	1	1	
						0	0	0	0	
							0	0	0	
								0	0	

Nussinov Folding Algorithm

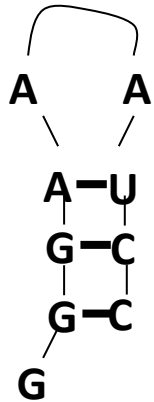
After scores for subsequences of length 9

$$\gamma(i, j) = \max \begin{cases} \gamma(i+1, j) \\ \gamma(i, j-1) \\ \gamma(i+1, j-1) + \delta(i, j) \\ \max_{i < k < j} [\gamma(i, k) + \gamma(k+1, j)] \end{cases}$$



		j →								
		G G G A A A U C C								
i ↓	G	0	0	0	0	0	0	1	2	3
	G	0	0	0	0	0	0	1	2	3
	G		0	0	0	0	0	1	<u>2</u>	2
	A			0	0	0	0	<u>1</u>	1	1
	A				0	0	<u>0</u>	1	1	1
	A					0	0	1	1	1
	U						0	0	0	0
	C							0	0	0
	C								0	0

Nussinov Folding Algorithm Traceback



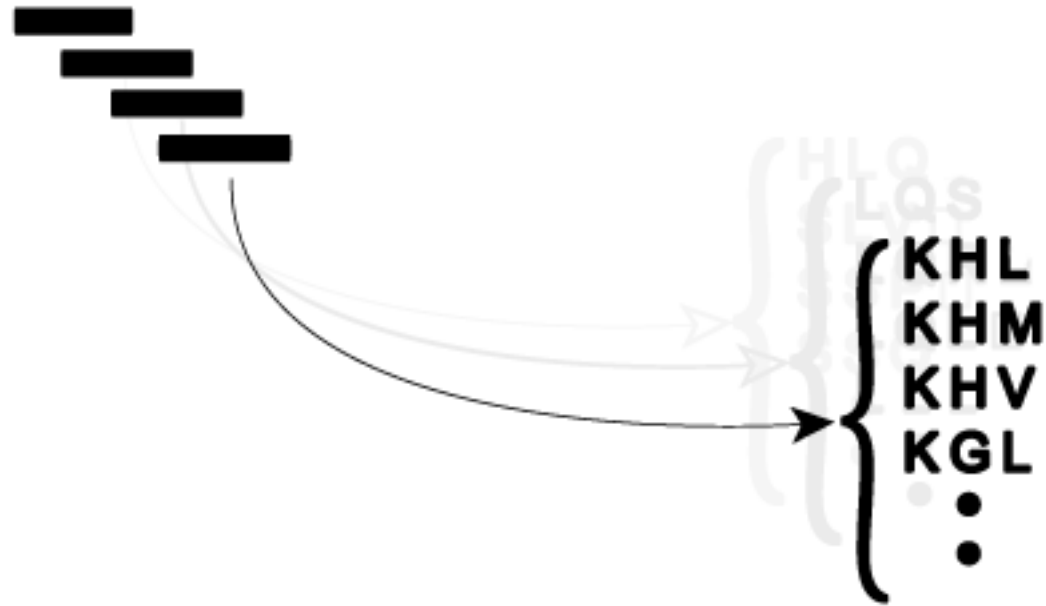
j →

	G	G	G	A	A	A	U	C	C
G	0	0	0	0	0	0	1	2	3
G	0	0	0	0	0	0	1	2	3
G		0	0	0	0	0	1	2	2
A			0	0	0	0	1	1	1
A				0	0	0	1	1	1
A					0	0	1	1	1
U						0	0	0	0
U							0	0	0
C								0	0
C									0

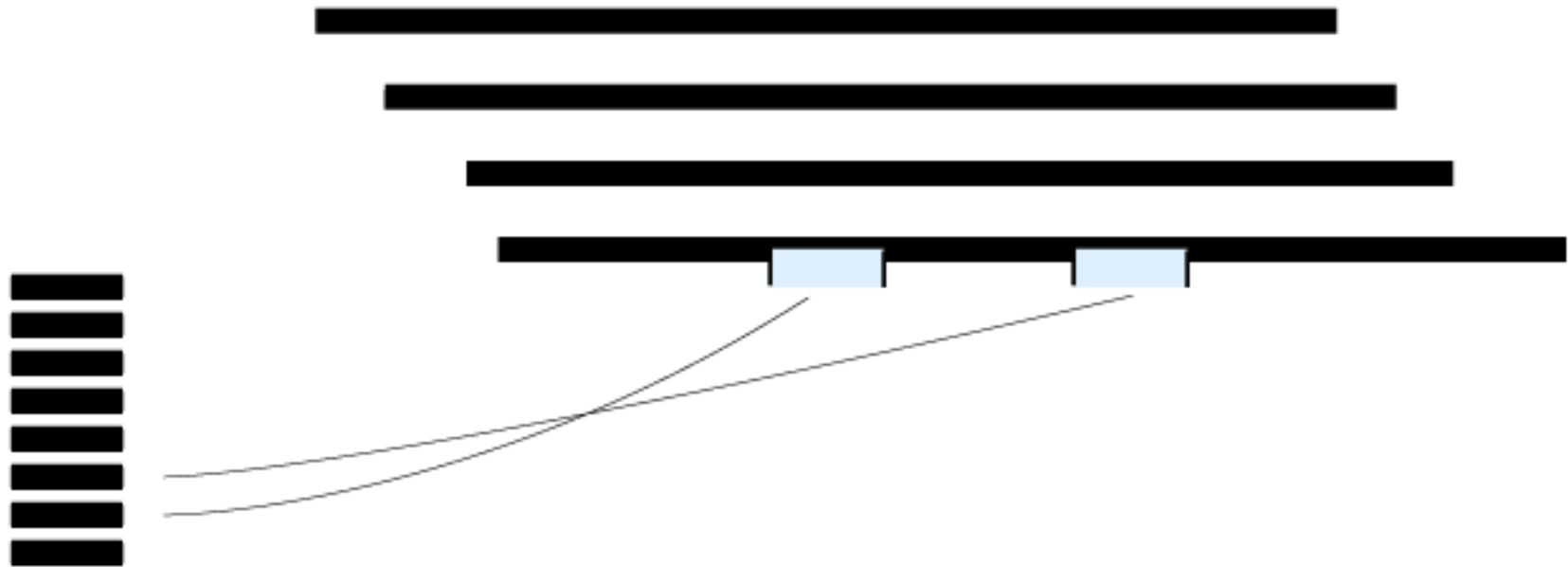
i ↓

BLAST

KHLQSLVARGVGGSLGGGAGDDA



BLAST

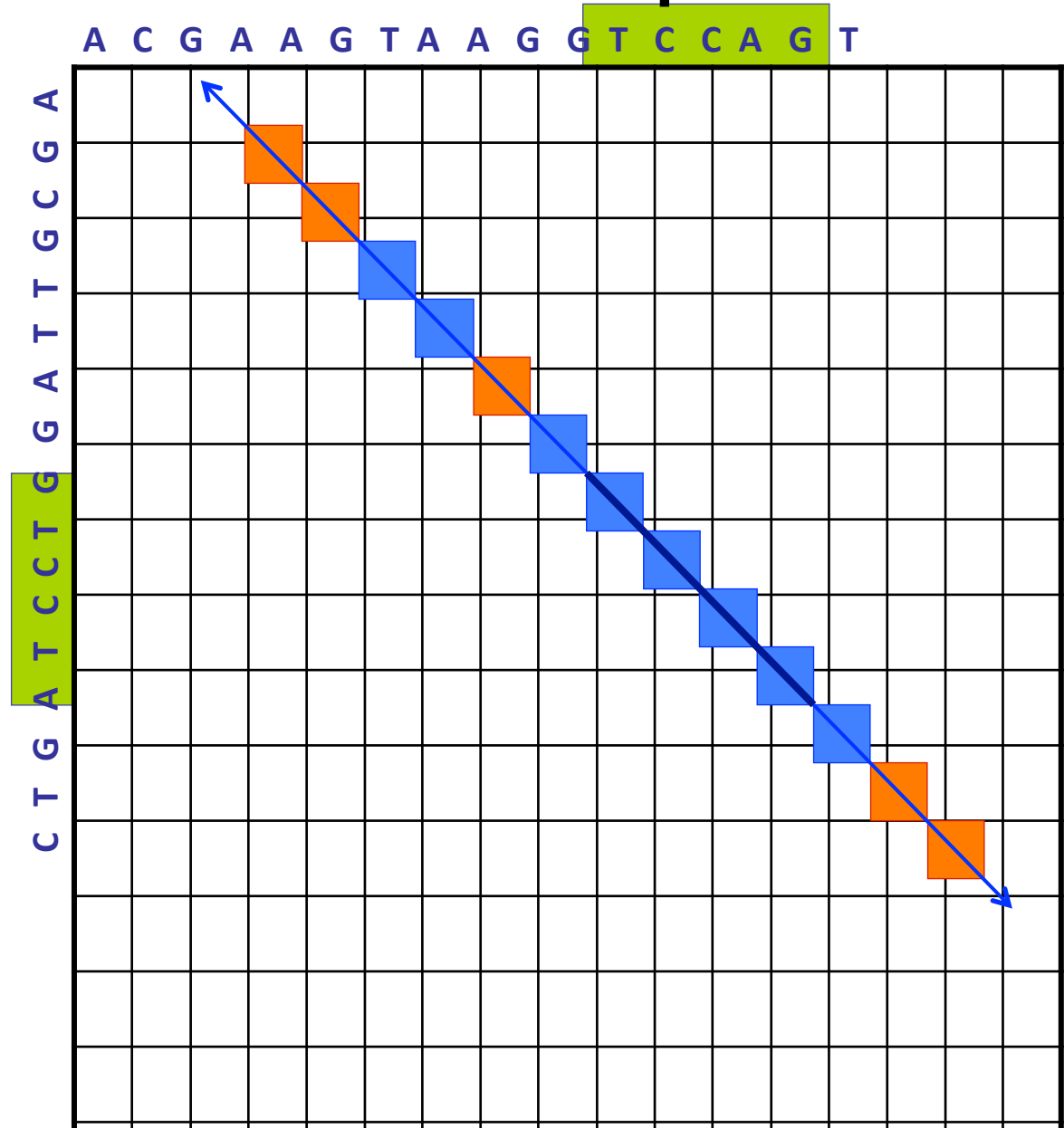


BLAST



Original BLAST: Example

- $w = 4$
- Exact keyword match of **GGTC**
- Extend diagonals with mismatches until score is under 50%
- Output result
GTAAGGTCC
GTTAGGTCC



BLAST finds a “hit” and then extends

```
GCNTACACGTCACCATCTGTGCCACCACNCATGTCTCTAGTGATCCCTCATAAGTTCCAACAAAGTTTGC
|| |||| | ||| ||| || | ||||| ||||| | ||||| | | |||||
GCCTACACACCGCCAGTTGTG-TTCCTGCTATGTCTCTAGTGATCCCTGAAAAGTTCCAGCGTATTTTGC

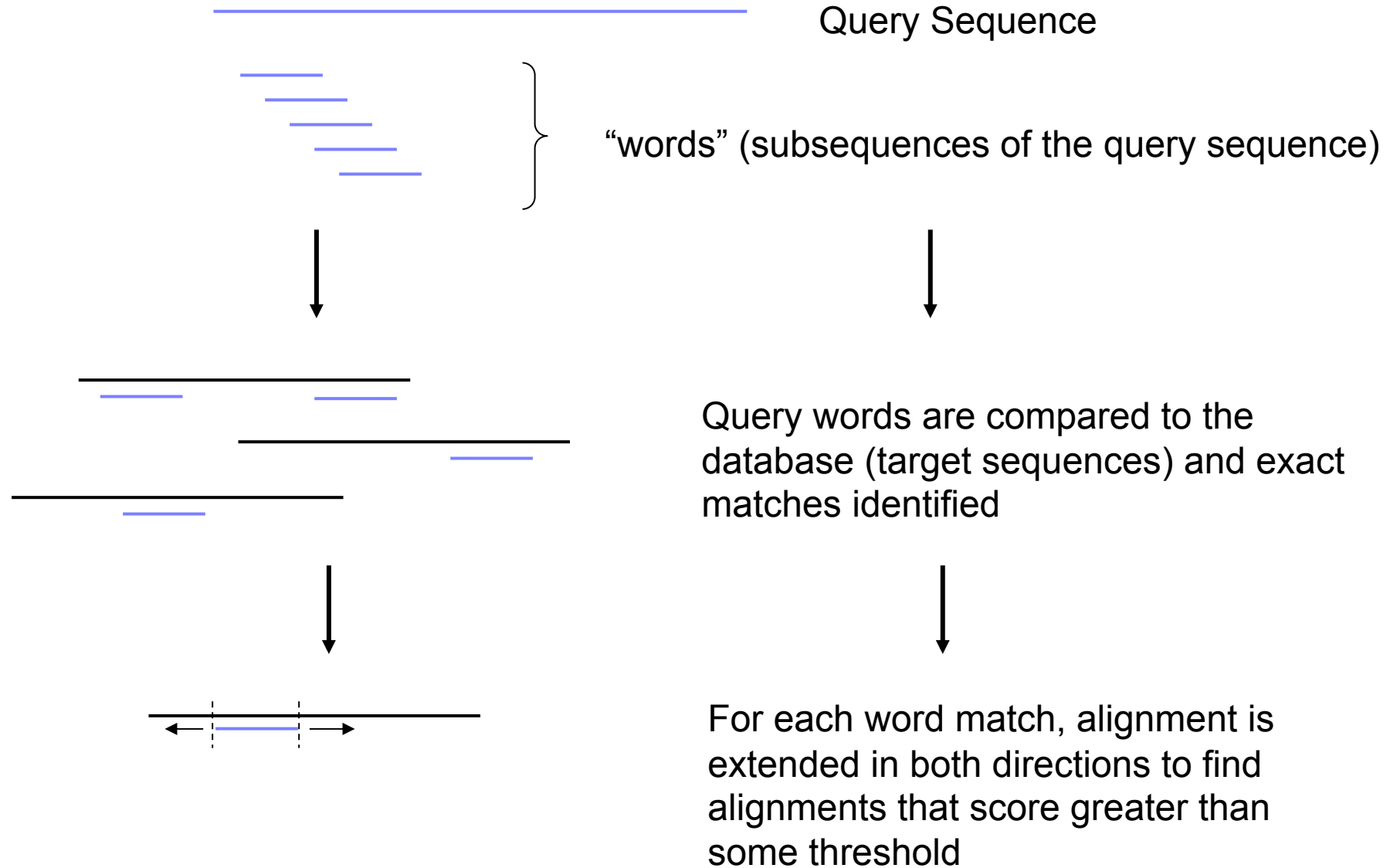
GAGTACTCAACACCAACATTGATGGGCAATGGAAAATAGCCTTCGCCATCACACCATTAAGGGTGA----
|| ||||| ||||| | |||| | ||||| || ||||| | | | | |
GAATACTCAACAGCAACATCAACGGGCAGCAGAAAATAGGCTTTGCCATCACTGCCATTAAGGATGTGGG

-----TGTTGAGGAAAGCAGACATTGACCTCACCGAGAGGGCAGGCGAGCTCAGGTA
|| ||||| ||||| || ||||| || ||||| || |||| |
TTGACAGTACACTCATAGTGTGAGGAAAGCTGACGTTGACCTCACCAAGTGGGCAGGAGAACTCACTGA

GGATGAGGTGGAGCATATGATCACCATCATAACAGAACTCAC-----CAAGATTCCAGACTGGTTCTTG
|| |||| | || | |||| |||| | |||| | | ||||| ||||| ||||| ||||| ||||| |||||
GGATGAGATGGAACGTGTGATGACCATTATGCAGAAATCCATGCCAGTACAAGATCCAGACTGGTTCTTG
```

Seed match = hit

BLAST



BLAST: example of missing a target

- Fail:

GAGTACTCAACACCAACATTAGTGGGCAATGGAAAAT

|| ||| ||| ||| ||| | ||| ||| ||| |||

GAATACTCAACAGCAACATCAATGGGCAGCAGAAAAT

- Dilemma

- Sensitivity – needs shorter seeds

- the success rate of finding a homology

- Speed – needs longer seeds

- Mega-BLAST uses seeds of length 28.

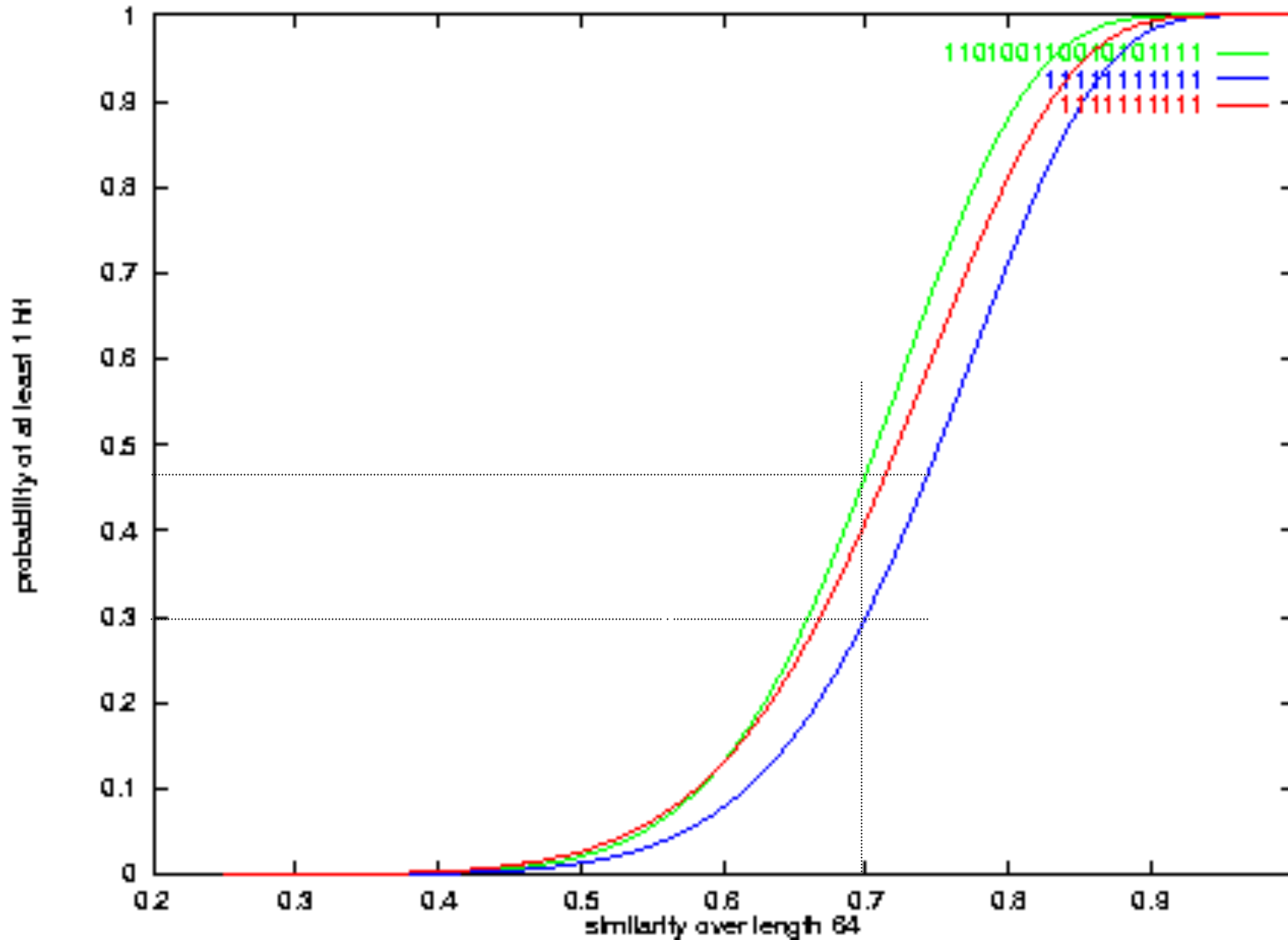
PatternHunter uses “spaced seeds”

- 111010010100110111 (called a model)
 - Eleven required matches (weight=11)
 - Seven “don’t care” positions

```
GAGTACTCAACACCAACATTAGTGGCAATGGAAAAT...
|| ||||| |||| | | |||| | |||||
GAATACTCAACAGCAACACTAATGGCAGCAGAAAAT...
      111010010100110111
```

- Hit = all the required matches are satisfied.
- BLAST seed model = 11111111111

Simulated sensitivity curves



PH2:

(homology identity = 0.7, homology length=64)

111011001011010111,

1111000100010011010111,

1100110100101000110111,

1110100011110010001101,

.....

.....

A multiple alignment

```

                10         20         30         40         50         60
                .         .         .         .         .         .
Hbb_Human.pep  -----VHLTPEEKSAVTALWGKVN--VDEVGGEALGRLLVVYPWTQRFFESFGDLST
Hbb_Horse.pep  -----VQLSGEEKAAVLALWDKVN--EEEVGGGEALGRLLVVYPWTQRFFDSFGDLSN
Hba_Human.pep  -----VLSPADKTNVKAAWGKVGAGHAGEYGAEALERMFLSFPTTKTYFPHFDLS--
Hba_Horse.pep  -----VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLS--
Myg_Phyca.pep  -----VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRLFKHLKT
Glb5_Petma.pep PIVDTGSVAPLSAAEKTKIRSAWAPVYSTYETSGVDILVKFFTSTPAAQEFFPKFKGLTT
Lgb2_Luplu.pep -----GALTESQAALVKSSWEEFNANIPKHTHRFFILVLEIAPAAKDLFSFLKGTSE
                * .         *         *         *         *

```

```

Hbb_Human.pep  PDAVMGNPKVKAHGKKVLGAFSDGLAHLA-----NLKGTFAATLSELHCDKLHVDPENFRL
Hbb_Horse.pep  PGAVMGNPKVKAHGKKVLHLSFGEGVHHLD-----NLKGTFAALSELHCDKLHVDPENFRL
Hba_Human.pep  ----HGSAQVKGHGKKVADALTNAVAHV-----DMPNALSALSDLHAHKL RVDPVNFKL
Hba_Horse.pep  ----HGSAQVKAHGKKVGDALTLAVGHLD-----DLPGALSNLSDLHAHKL RVDPVNFKL
Myg_Phyca.pep  EAEMKASEDLKKHGVTVLTALGAILKKKG-----HHEAELKPLAQSHATKHKIPIKYLEF
Glb5_Petma.pep ADQLKKSADVRWHAERIINAVNDAVASMDDT--EKMSMKLRDLSGKHAKSFQVDPQYFKV
Lgb2_Luplu.pep VP--QNNPELQAHAGKVFKLVYEAAIQVLTGTVVTDATLKNLGSVHVSKG-VADAHFPV
                . . *         .         *         *

```

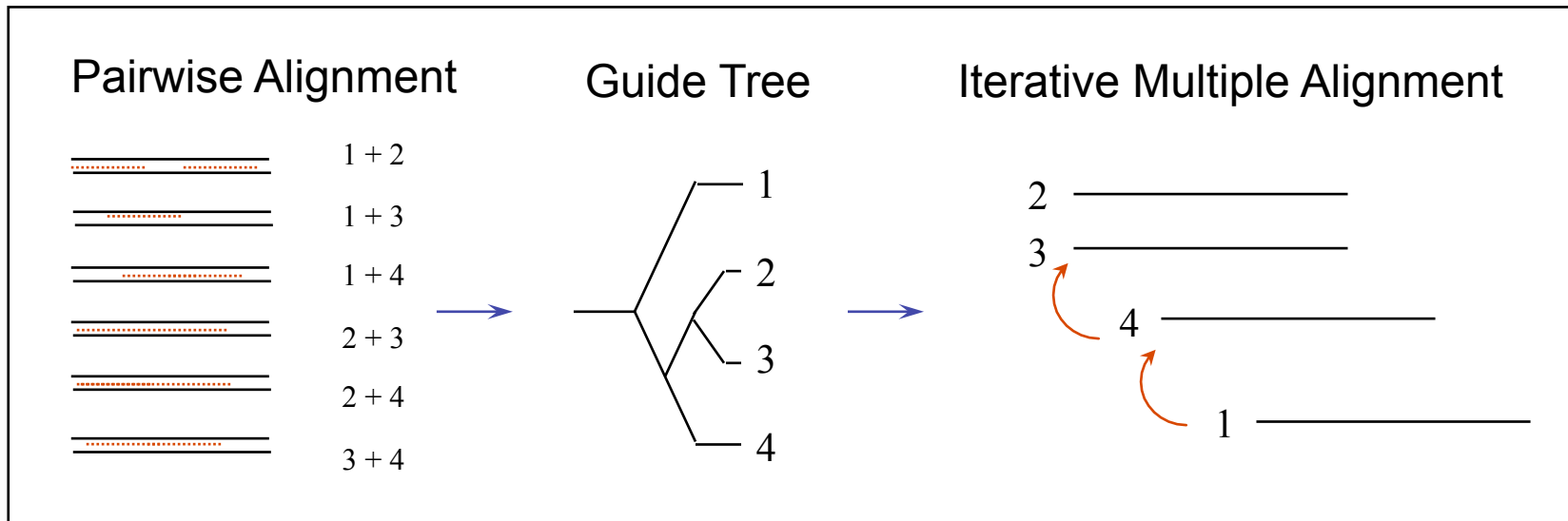
```

Hbb_Human.pep  LGNVLVCVLAHFFGKEFTPPVQAAYQKVVAGVANALAHKYH-----
Hbb_Horse.pep  LGNVLVVVLARHFGKDFTPPELQASYQKVVAGVANALAHKYH-----
Hba_Human.pep  LSHCLLVTLAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR-----
Hba_Horse.pep  LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTSKYR-----
Myg_Phyca.pep  ISEAIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG
Glb5_Petma.pep LAAVIADTVAAG-----DAGFEKLMSMICILLRSAY-----
Lgb2_Luplu.pep VKEAILKTIKEVVGAKWSEELNSAWTIAYDELAIVIKKEMNDAA---
                .         .         .

```

Progressive Alignment

Principle:



1



2



3

Step 1-pairwise alignments

Compare each sequence with each other and calculate a Similarity matrix.

Different sequences

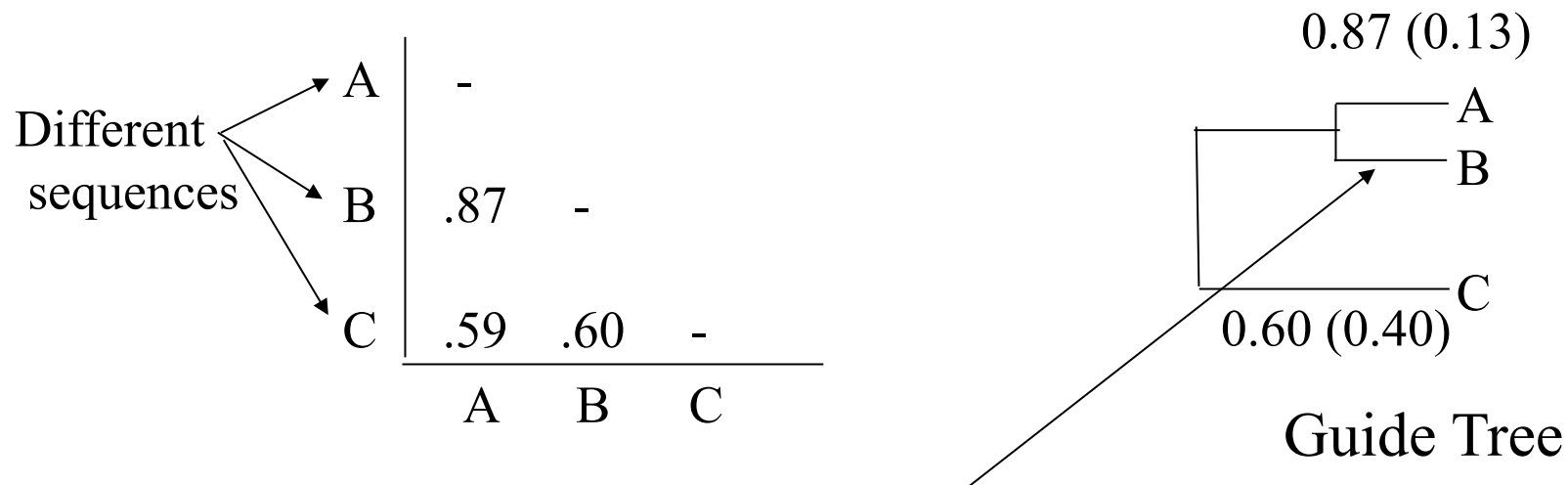
A	-		
B	.87	-	
C	.59	.60	-
	A	B	C

Each number represents the number of exact matches divided by the sequence length (ignoring gaps). Thus, the higher the number the more closely related the two sequences are.

In this similarity (distance) matrix sequence A is 87% identical to sequence B

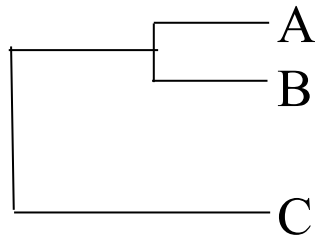
Step 2-Create Guide Tree

Use the Similarity Matrix to create a Guide Tree to determine the “order” of the sequences.



Branch length proportional to estimated divergence between A and B (0.13)

Step 3-Progressive Alignment



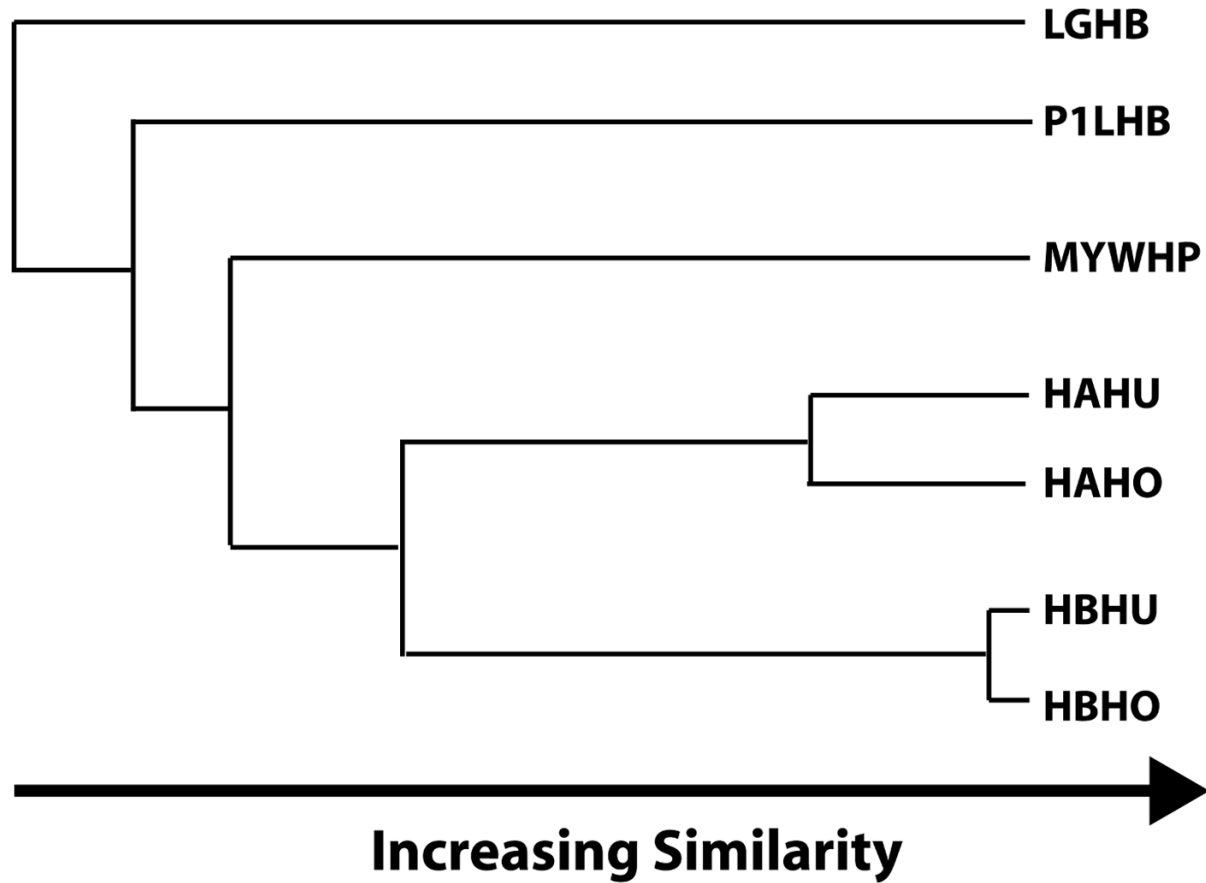
Guide Tree

Align A and B first. Then add sequence C to the previous alignment. In the closely aligned sequences and gaps are given a heavier weight than more divergent sequences.

How does it work?

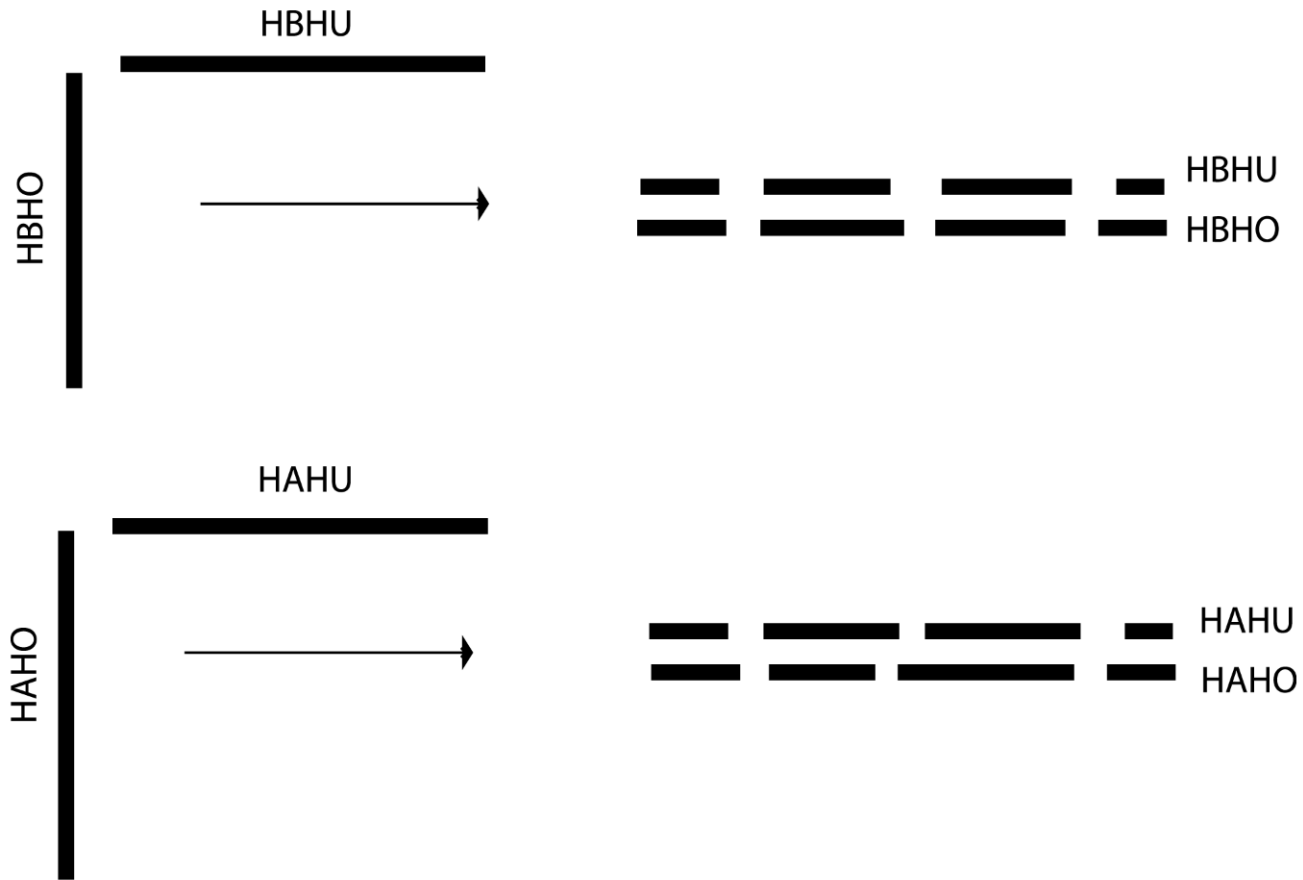
- Starting with a group of 7 sequences from different species
- Do pairwise alignments between all 7 sequences
- Score given for similarity, higher score indicates more similar

	HAHU	HBHU	HAHO	HBHO	MYWHP	P1LHB	LGHB
HAHU							
HBHU	21.1						
HAHO	32.9	19.7					
HBHO	20.7	39.0	20.4				
MYWHP	11.0	9.8	10.3	9.7			
P1LHB	9.3	8.6	9.6	8.4	7.0		
LGHB	7.1	7.3	7.5	7.4	7.3	4.3	

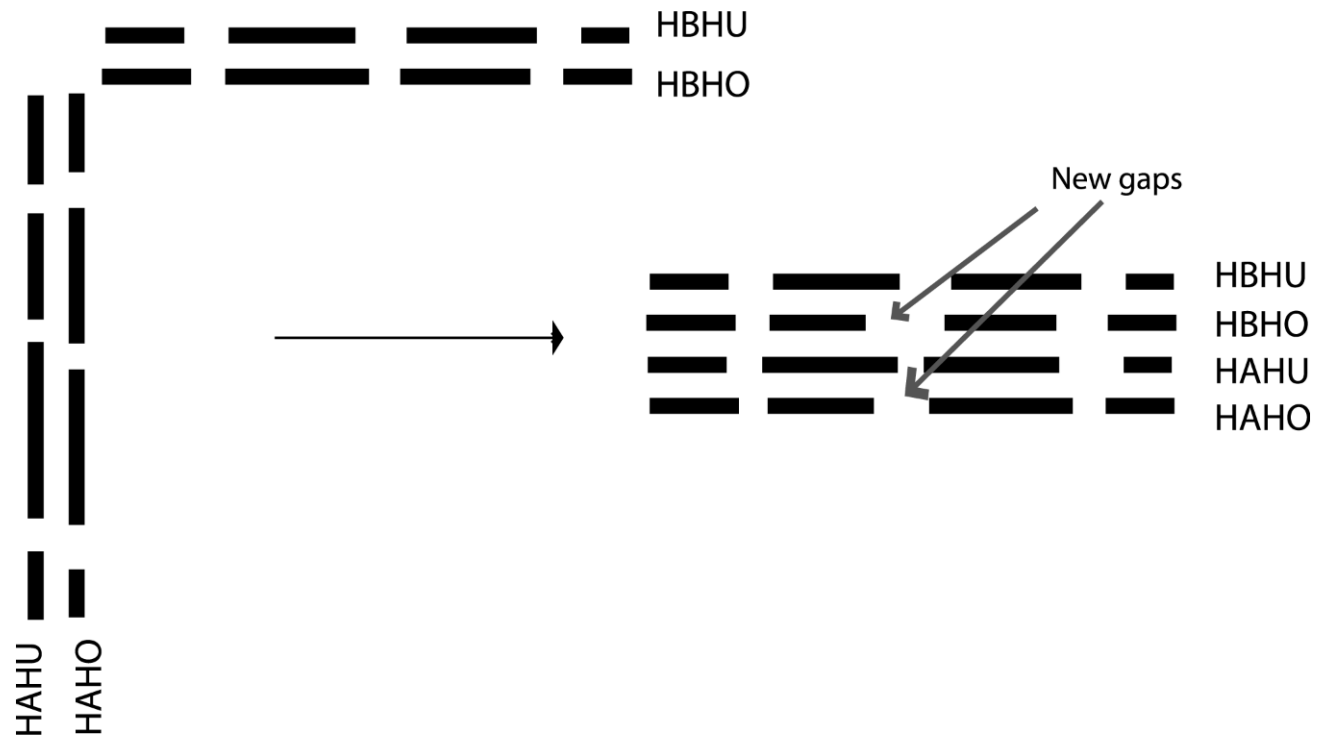


- Cluster the sequences by similarity to create a guide tree
- Branch length is proportional to estimated divergence between the two sequences

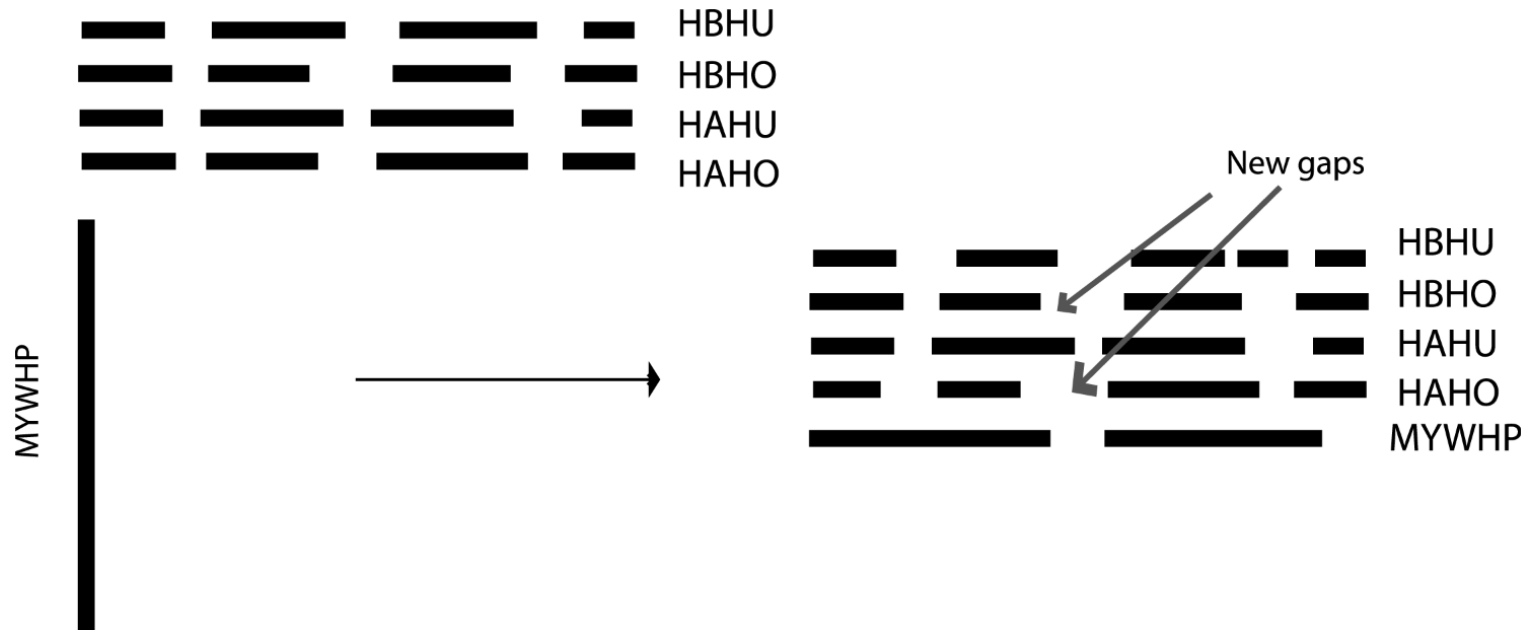
Progressive Alignment



Progressive Alignment



Progressive Alignment



More Burrows-Wheeler

Input

SIX.MIXED.PIXIES.SIFT.SIXTY.PIXIE.DUST.BOXES

Burrows-Wheeler Output

TEXYDST.E.IXIXIXSSMPPS.B..E.S.EUSFXDIIIOIIT

Repeated characters mean that it is easier to compress

More Burrows-Wheeler

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATGATACGGCGACCACCGAGATCTA

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATGATACGGCGACCACCGAGATC **TA**

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATGATACGGCGACCACCGAGATCTA

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATGATACGGCGAC **CACCGAGATCTA**

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATGATACGGCGACCCACCGAGATCTA

More Burrows-Wheeler

Reference



BWT(Reference)

Query:

AATGATACGGCGACCACCGAGATCTA

More Burrows-Wheeler

Reference



BWT(Reference)



Query:

AATG TACGGCGACCACCGAGATCTA

More Burrows-Wheeler

Reference



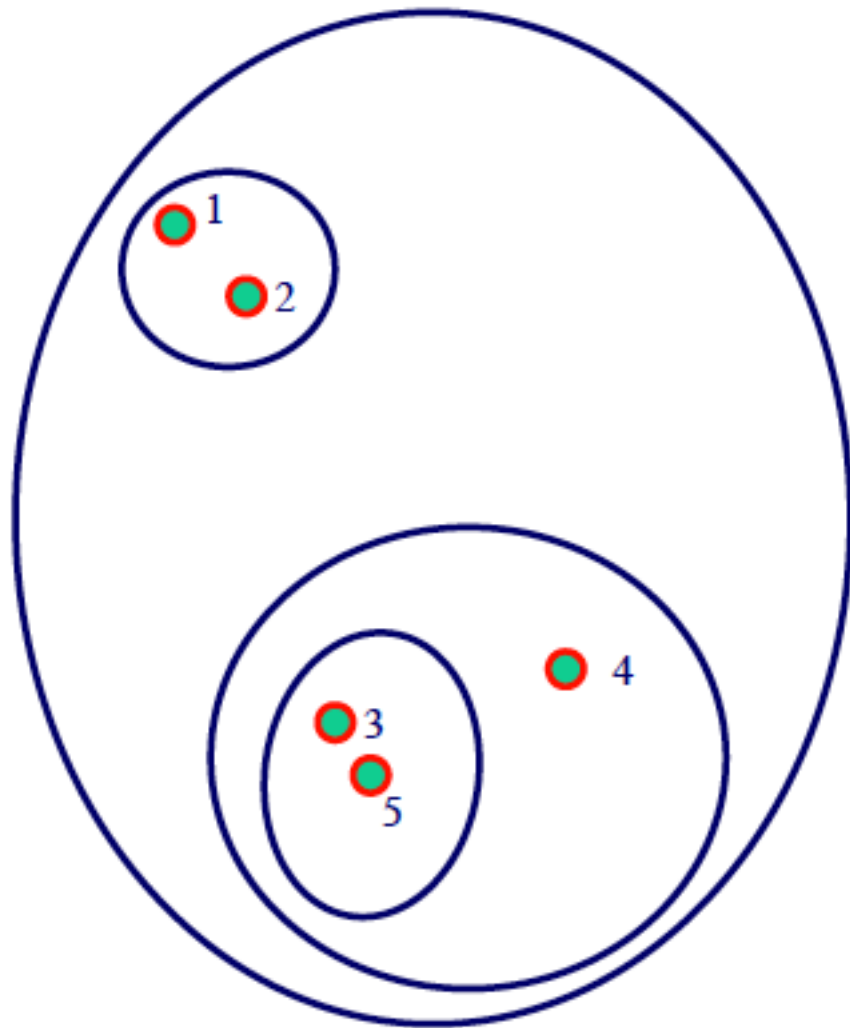
BWT(Reference)



Query:

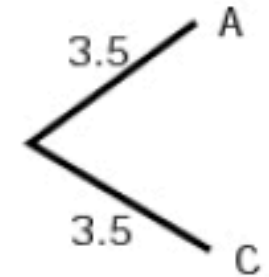
AATGTTACGGCGACCACCGAGATCTA

UPGMA



UPGMA

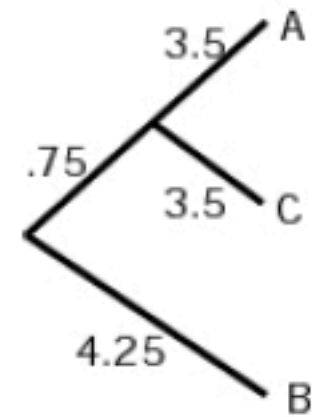
	A	B	C	D
A	0			
B	8	0		
C	7	9	0	
D	12	14	11	0



$$M_{B(AC)} = (M_{BA} + M_{BC})/2 = (8+9)/2=8.5$$

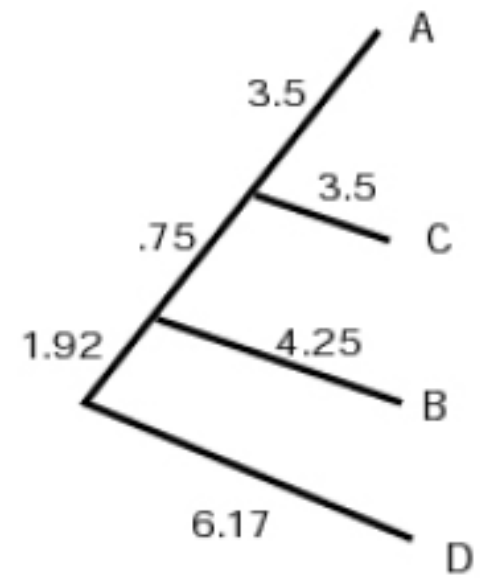
$$M_{D(AC)} = (M_{DA} + M_{DC})/2 = (12+11)/2=11.5$$

	AC	B	D
AC	0		
B	8.5	0	
D	11.5	14	0



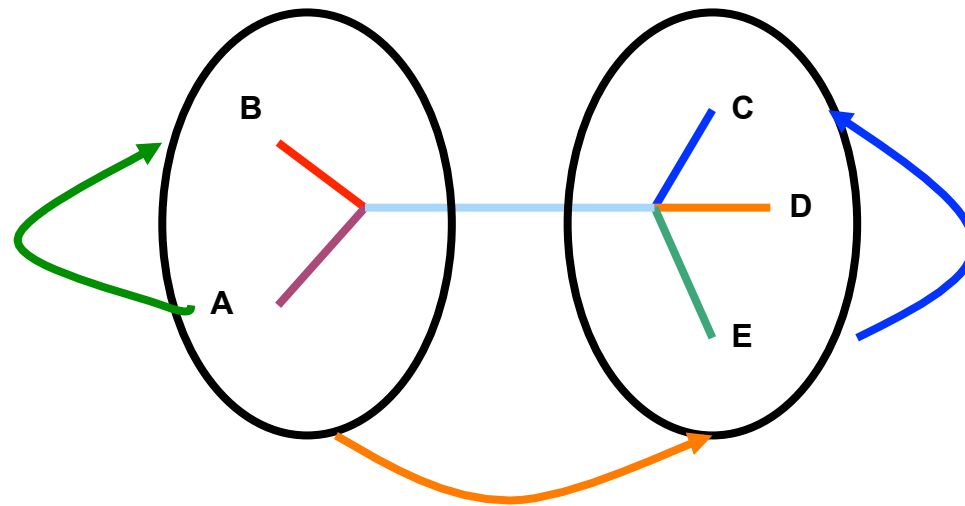
$$M_{(ABC)D} = (M_{AD} + M_{BD} + M_{CD})/3 = (12+14+11)/3$$

	ABC	D
ABC	0	
D	12.33	0

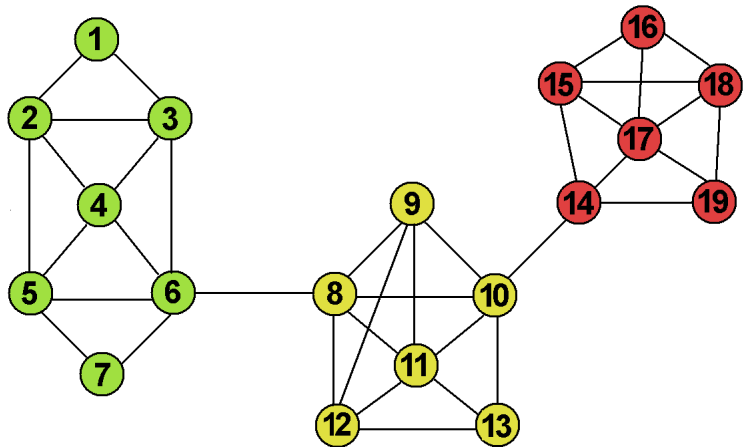


Neighbor Joining

- If A and B are joined:



$$S_{mn} = \frac{\sum (d_{im} + d_{in})}{2(N-2)} + \frac{d_{mn}}{2} + \frac{\sum d_{ij}}{N-2}$$



0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	1	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	1	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	1	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0
0	0	1	1	1	0	1	1	1	0	0	0	0	0	0	0	0	0	0
0	0	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0
0	0	0	0	0	0	1	0	0	1	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	0	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	1	1	0	1	0	1	1	0	0	0	0
0	0	0	0	0	0	0	0	1	1	1	0	1	1	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	1	1	1	0	0	0	0	0	0
0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	1	0	1
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1

MCL animation

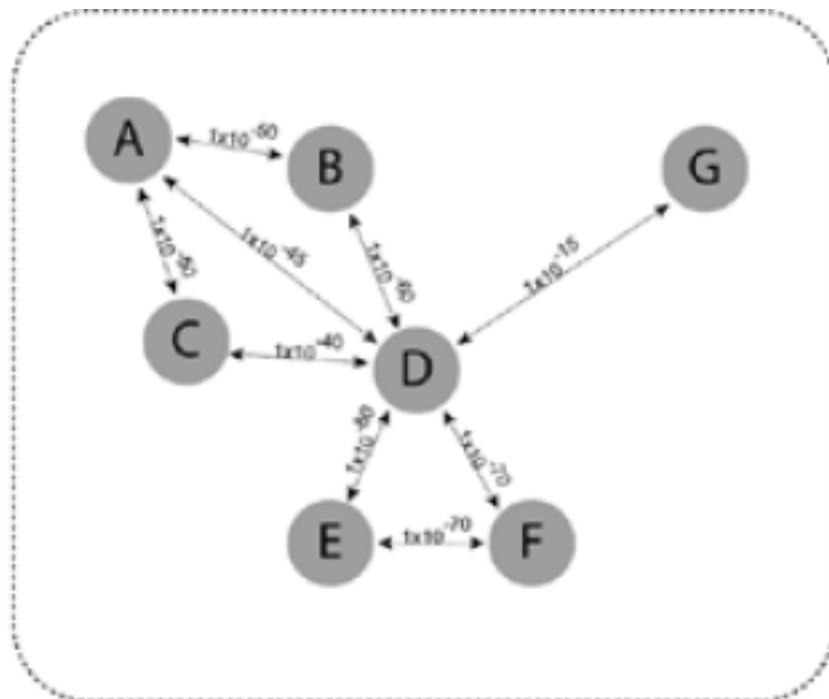
- Found at

www.micans.org/mcl/ani/mcl-animation.html

An efficient algorithm for large-scale detection of protein families

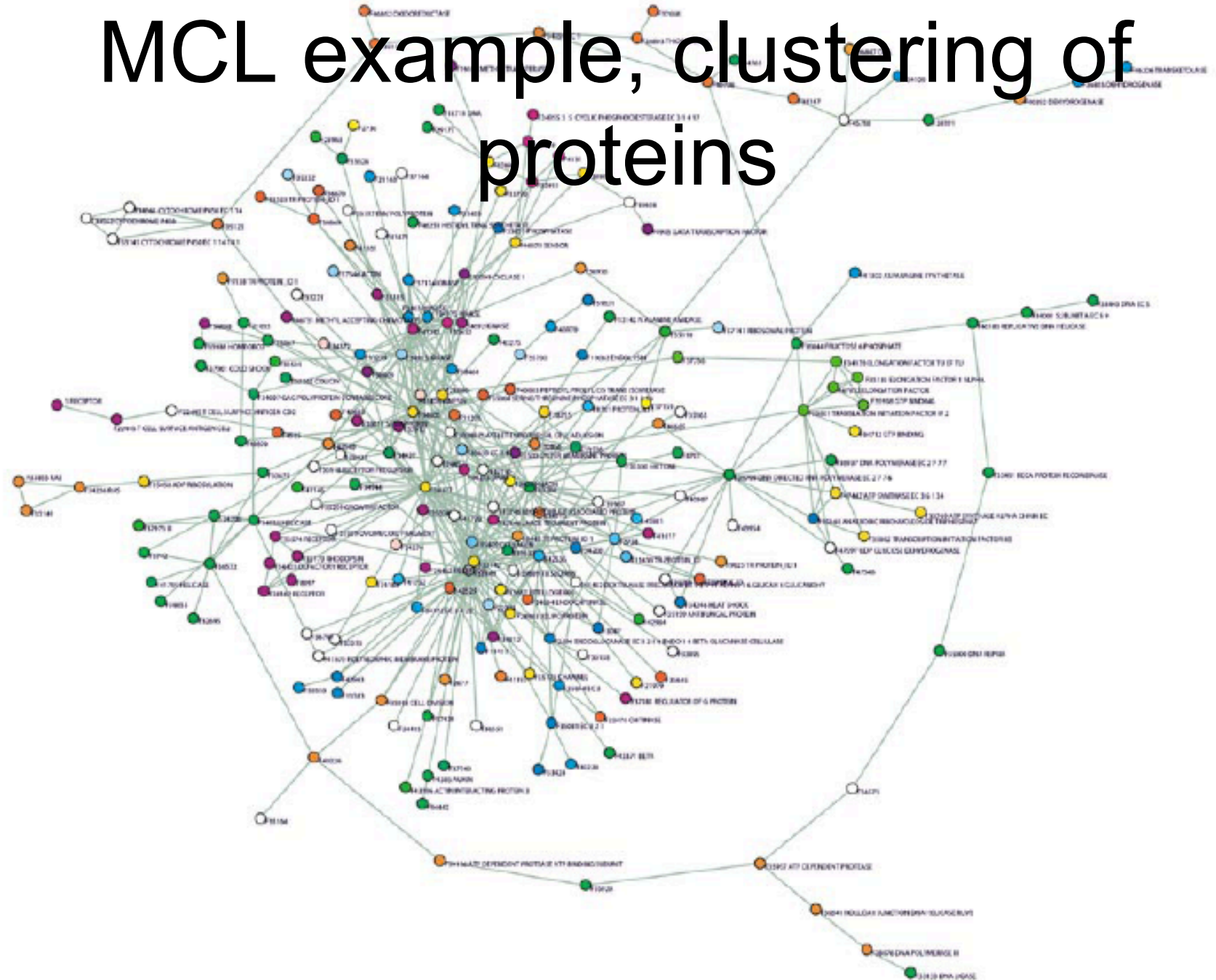
A. J. Enright*, S. Van Dongen¹ and C. A. Ouzounis

Computational Genomics Group, The European Bioinformatics Institute, EMBL Cambridge Outstation, Cambridge CB10 1SD, UK and ¹Centrum voor Wiskunde en Informatica, Kruislaan 413, NL-1098 SJ Amsterdam, The Netherlands

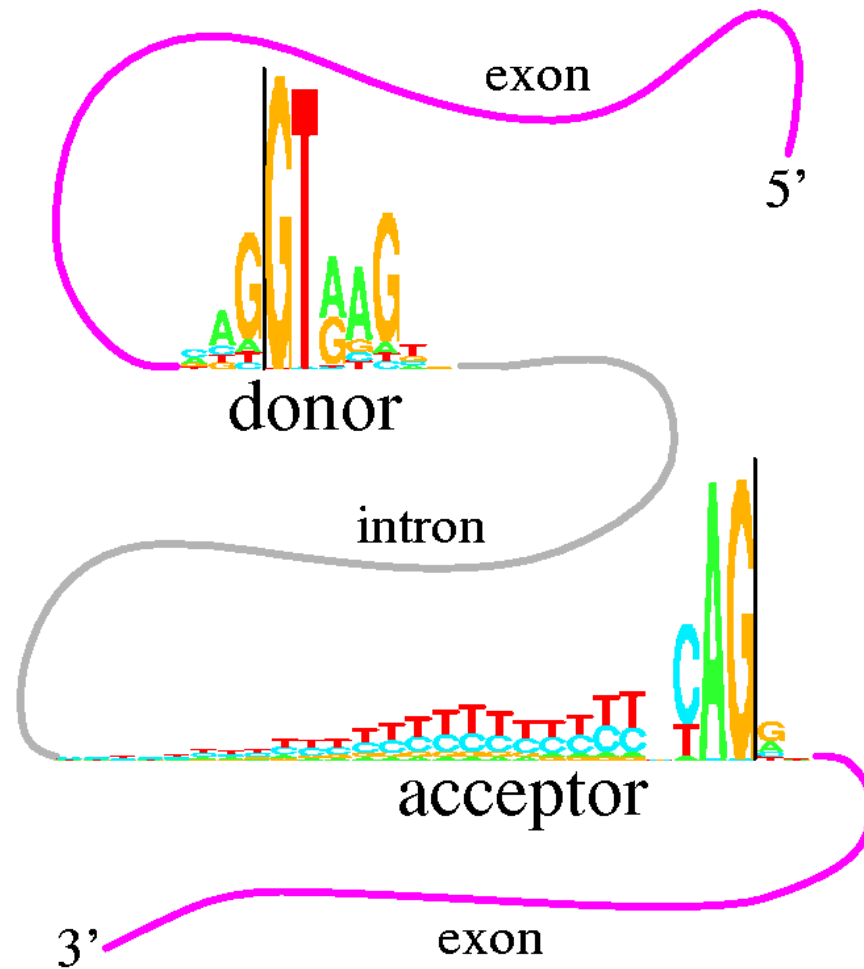


	A	B	C	D	E	F	G
A	0.42	0.24	0.20	0.11	0.00	0.00	0.00
B	0.20	0.48	0.24	0.15	0.00	0.00	0.00
C	0.20	0.00	0.40	0.10	0.00	0.00	0.00
D	0.18	0.28	0.16	0.24	0.32	0.29	0.13
E	0.00	0.00	0.00	0.19	0.40	0.29	0.00
F	0.00	0.00	0.00	0.17	0.28	0.42	0.00
G	0.00	0.00	0.00	0.04	0.00	0.00	0.87

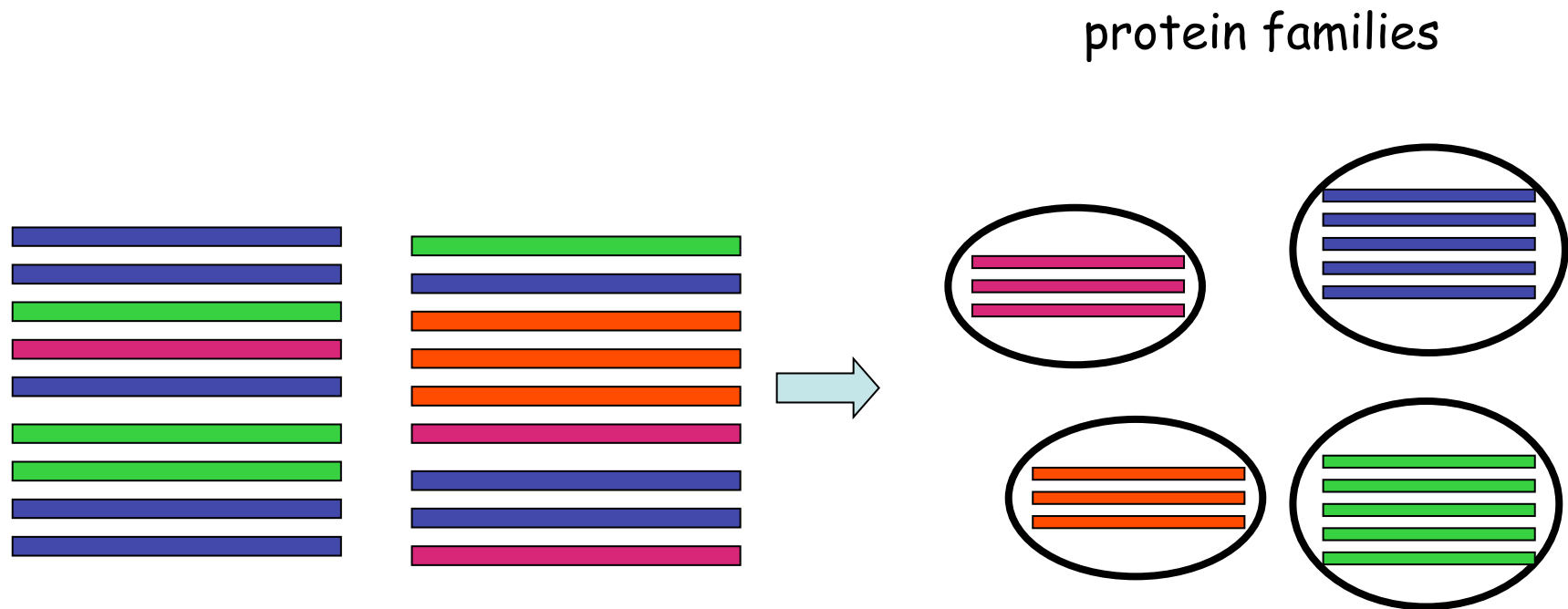
MCL example, clustering of proteins



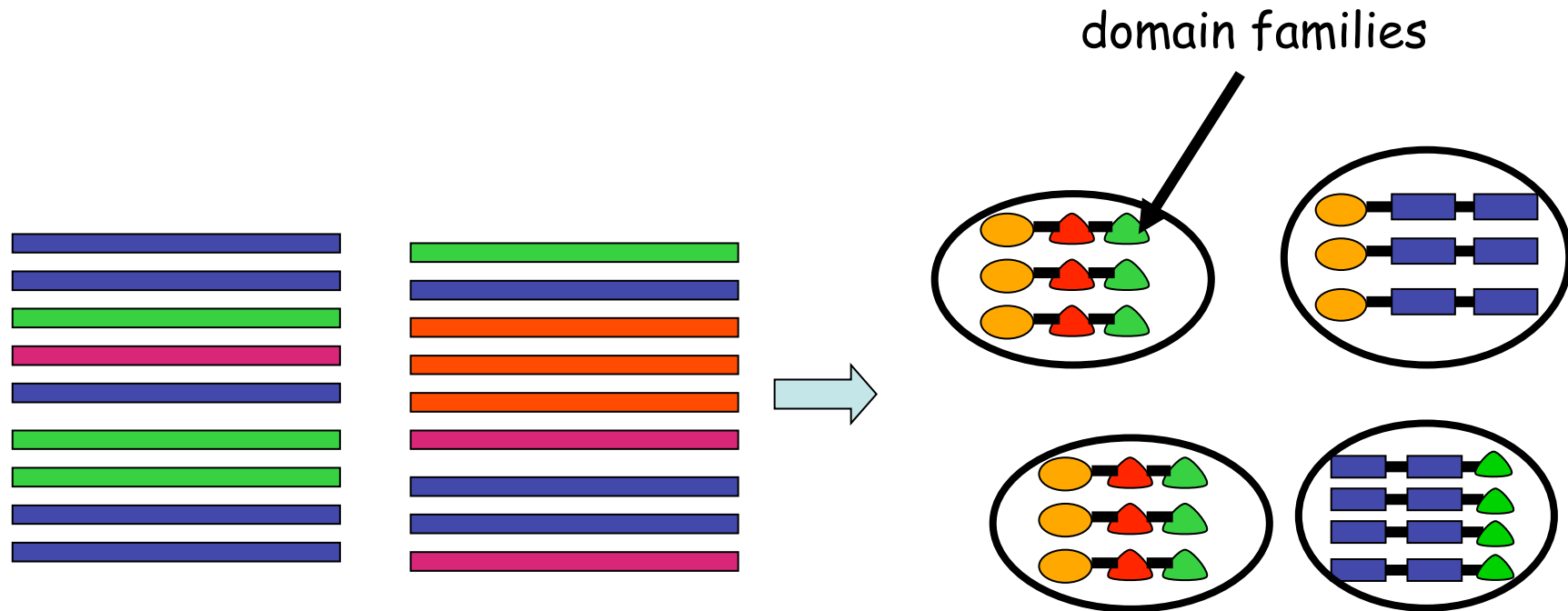
Applications of HMM



Applications of HMM: Identify protein families in the genomes



Applications of HMM: Identify domains and protein families in the genomes



The Probability that they are all Fair die rolls

Let the sequence of rolls be:

$x = 1, 2, 1, 5, 6, 2, 1, 5, 2, 4$

Then, what is the likelihood of

$\pi = \mathbf{F, F, \dots, F}$?

(say initial probs $a_{0\text{Fair}} = 1/2, a_{0\text{Loaded}} = 1/2$)

$1/2 \times P(1 | \text{Fair}) P(\text{Fair} | \text{Fair}) P(2 | \text{Fair}) P(\text{Fair} | \text{Fair}) \dots P(4 | \text{Fair}) =$

$$1/2 \times (1/6)^{10} \times (0.95)^9 = .00000000521158647211 \approx 0.5 \times 10^{-9}$$

The $1/2$ comes from the initial probability. Each value has an emission probability of $(1/6)$. There 9 transitions from Fair to Fair each with a probability of .95.

Sequence of rolls have many 6's

Let the sequence of rolls be:

$$x = 1, 6, 6, 5, 6, 2, 6, 6, 3, 6$$

Now, what is the likelihood $\pi = F, F, \dots, F$?

$$\frac{1}{2} \times (1/6)^{10} \times (0.95)^9 = 0.5 \times 10^{-9}, \text{ same as before}$$

What is the likelihood

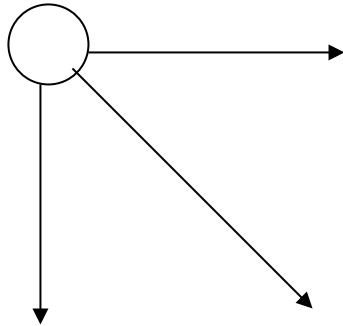
$$\pi = L, L, \dots, L?$$

$$\frac{1}{2} \times (1/10)^4 \times (1/2)^6 (0.95)^9 = .00000049238235134735 \approx 0.5 \times 10^{-7}$$

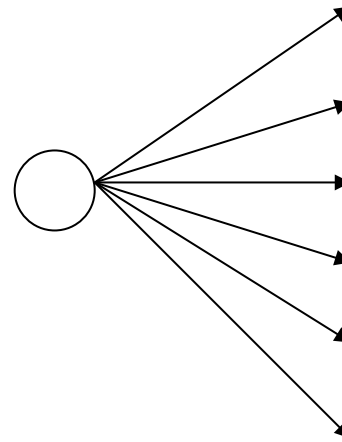
There are 4 non-6 values with a emission probability of 1/10 and 6 values of 6 each with an emission probability of (1/2) thus (1/2)^6. There are 9 state transitions of loaded die with a probability of .95 per transition.

So, it is 100 times more likely the die is loaded. However, this is not the only possible parse because some could be fair and some could be loaded.

Decoding Problem vs. Alignment Problem



Valid directions in the
alignment problem.



Valid directions in the
decoding problem.

Viterbi Example

Let x be a long sequence with a portion of $\sim 1/6$ 6's,
followed by a portion of $\sim 1/2$ 6's...

$x = 123456123456\dots 12345 \text{ 6626364656}\dots 1626364656$

Then, it is not hard to show that optimal parse is (exercise):

FFF.....F LLL.....L

The optimal parse is a set of Fairs in the beginning and a set of Loads at the end and the optimal cutting position will be at a non-6 followed by a 6. You would prefer to go one step forward and take a non-6 with a Fair state than a 6. And similarly, you would prefer to go one step back and take a 6 with a loaded die rather than a Fair state die because of the transition probabilities.

6 characters "123456" parsed as F, contribute $.95^6 \times (1/6)^6 = 1.6 \times 10^{-5}$
 parsed as L, contribute $.95^6 \times (1/2)^1 \times (1/10)^5 = 0.4 \times 10^{-5}$

"162636" parsed as F, contribute $.95^6 \times (1/6)^6 = 1.6 \times 10^{-5}$
 parsed as L, contribute $.95^6 \times (1/2)^3 \times (1/10)^3 = 9.0 \times 10^{-5}$

Underflows and Logs

Underflows are a significant problem

$$P[x_1, \dots, x_i, \pi_1, \dots, \pi_i] = a_{0\pi_1} a_{\pi_1\pi_2} \dots a_{\pi_i} e_{\pi_1}(x_1) \dots e_{\pi_i}(x_i)$$

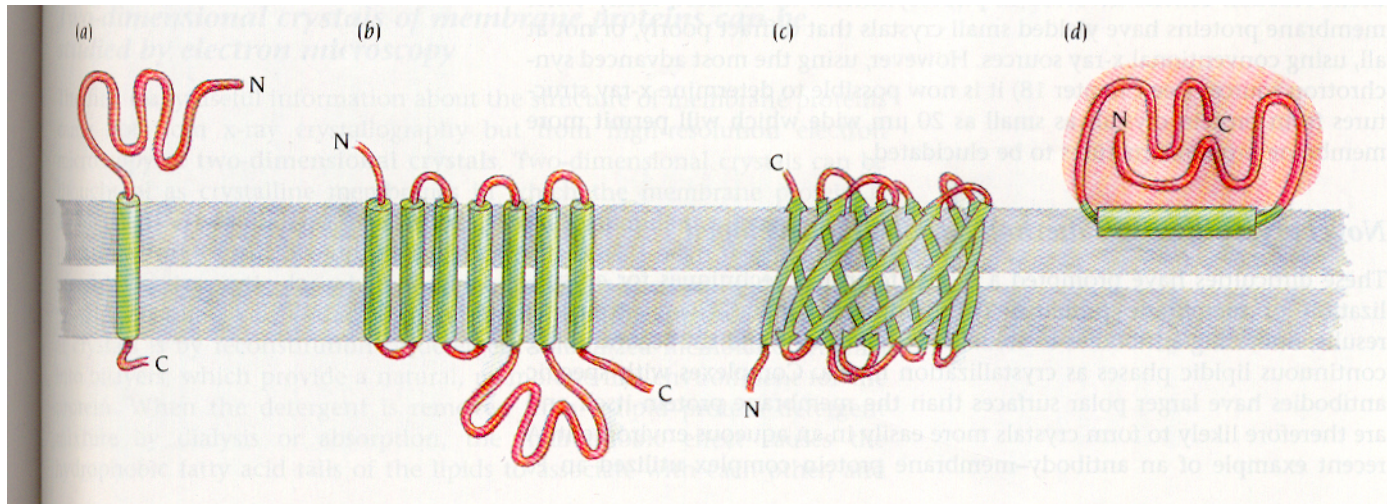
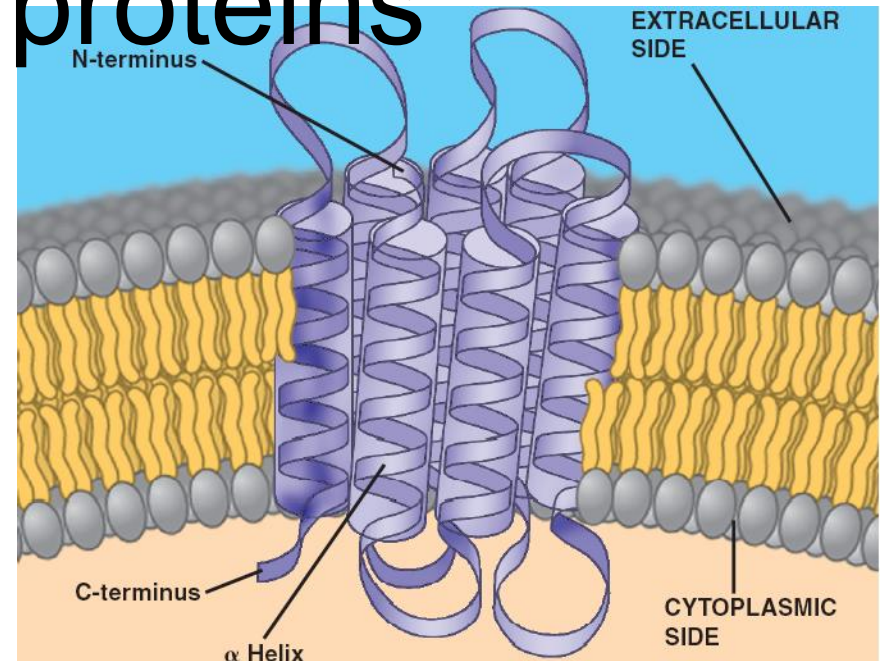
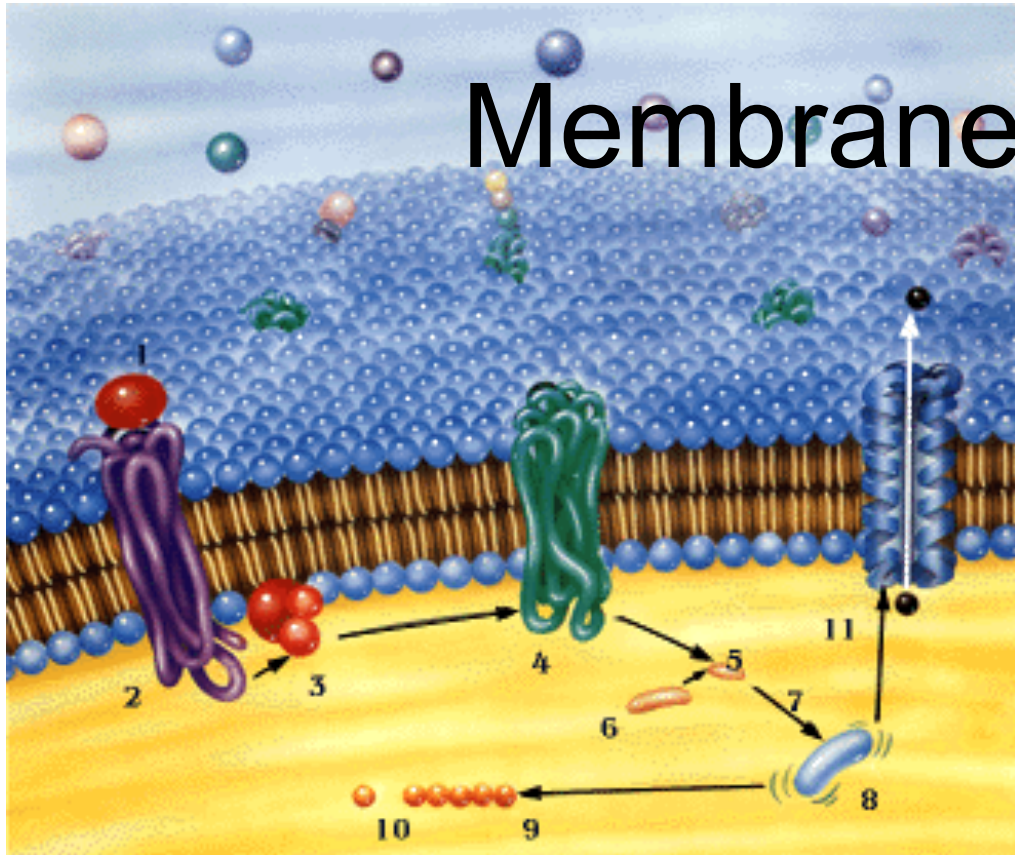
These numbers become extremely small probabilities and underflows the computer stored value.

Solution: Take the logs of all values

$$V_i(i) = \log e_k(x_i) + \max_k [V_k(i-1) + \log a_{kj}]$$

Should initial $V_0(0)$ to 0 instead of 1.

Membrane proteins



Example for TMHMM

www.cbs.dtu.dk/services/TMHMM/

>gi|218694017|ref|YP_002401684.1| membrane protein; channel [Escherichia coli 55989]

MQDLISQVEDLAGIEIDHTTSMVMIFGIIFLTAVVVHILHWVVLRTFEKRAIASS

RLWLQIITQNKLFH

RLAFTLQGIIVNIQAVFWLQKGTAAADILTTCAQLWIMMYALLSVFSLLDVILNL

AQKFPAASQLPLKGI

FQGIKLIGAILVGILMISLLIGQSPAILISGLGAMAAVLMLVFKDPILGLVAGIQLS

ANDMLKLGDWLEM

PKYGADGAVIDIGLTTVKVRNWDNTITTIPTWSLVSDSFKNWMSGMSASGGRR

IKRSISIDVTSIRFLDED

EMQRLNKAHLLKPYLTSRHQEINEWNRQQGSTESILNLRRMTNIGTFRAYLN

EYLRNHPRIRKDMTLMVR

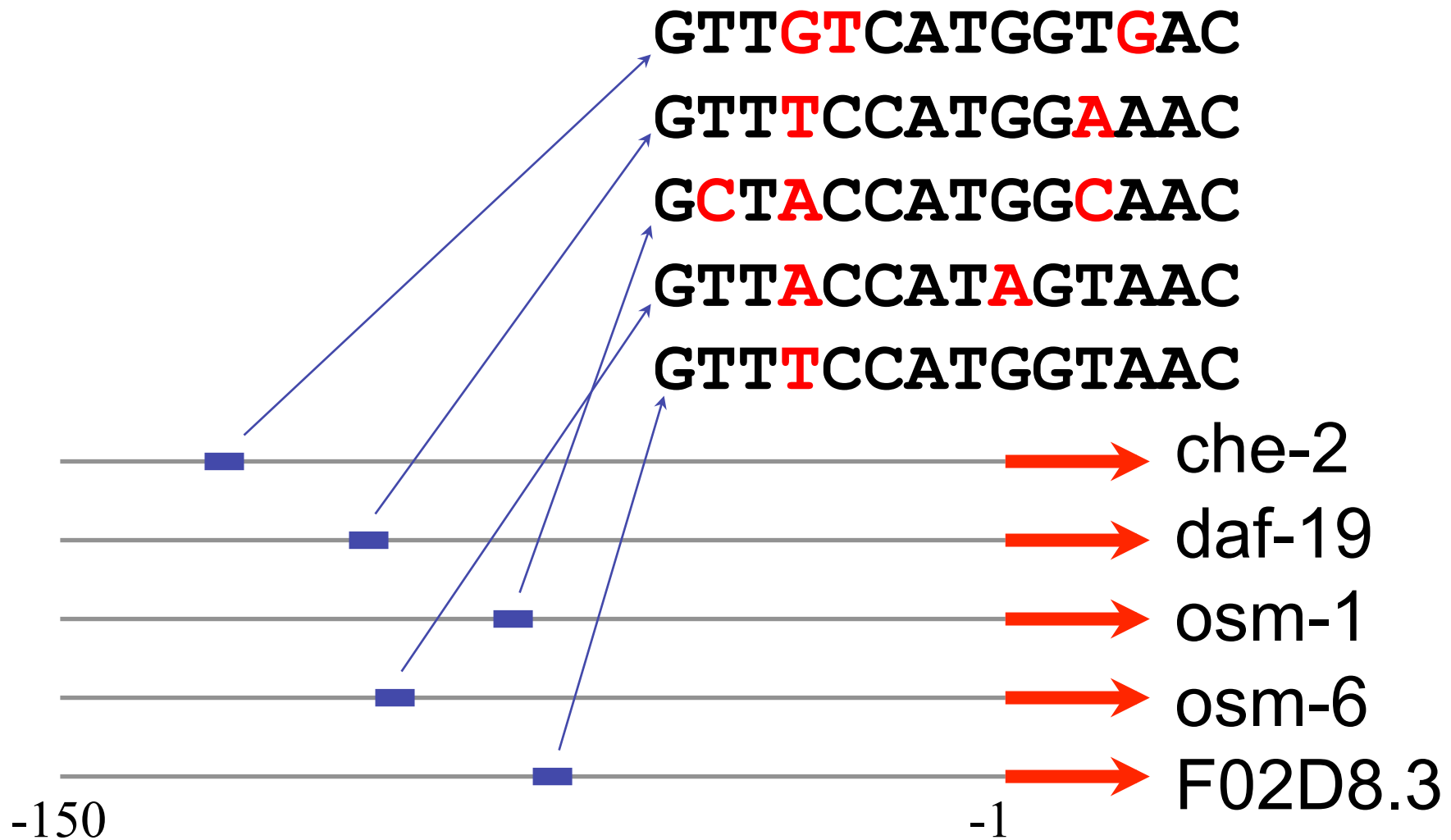
QLAPGDNGLPLEIYAFTNTVWWEYESIQADIFDHIFAIVEEFGLRLHQSPGTGN

DIRSLAGAFKQ

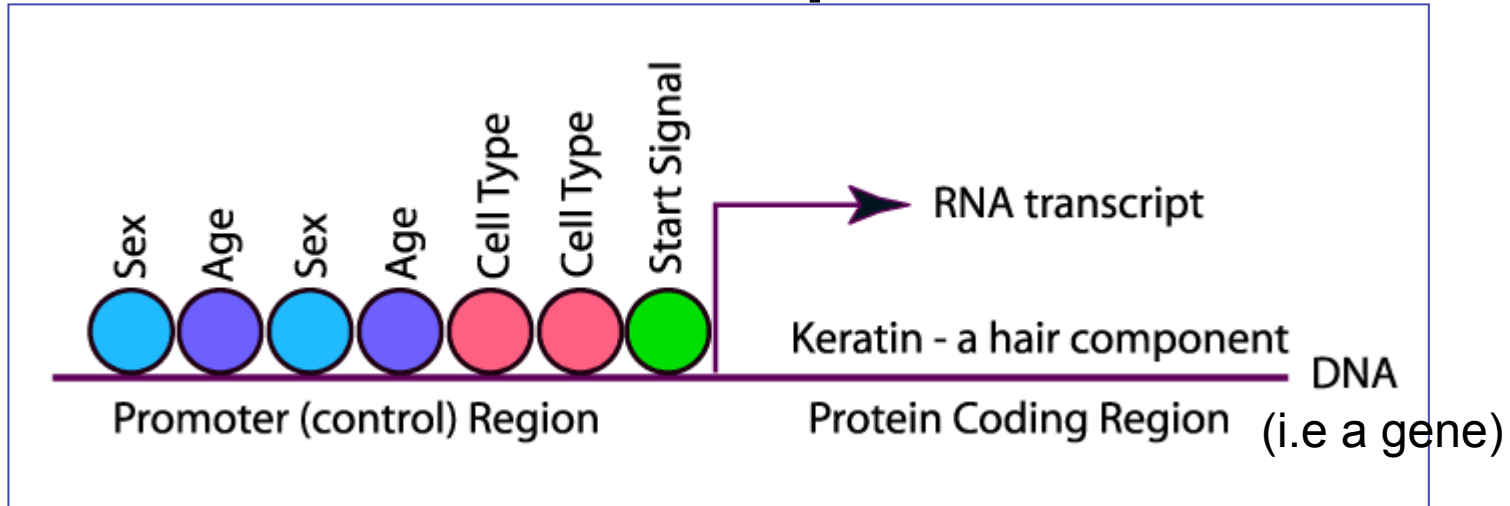
Gibbs Sampling server

http://bayesweb.wadsworth.org/cgi-bin/gibbs.8.pl?data_type=DNA

Examples of regulatory sites

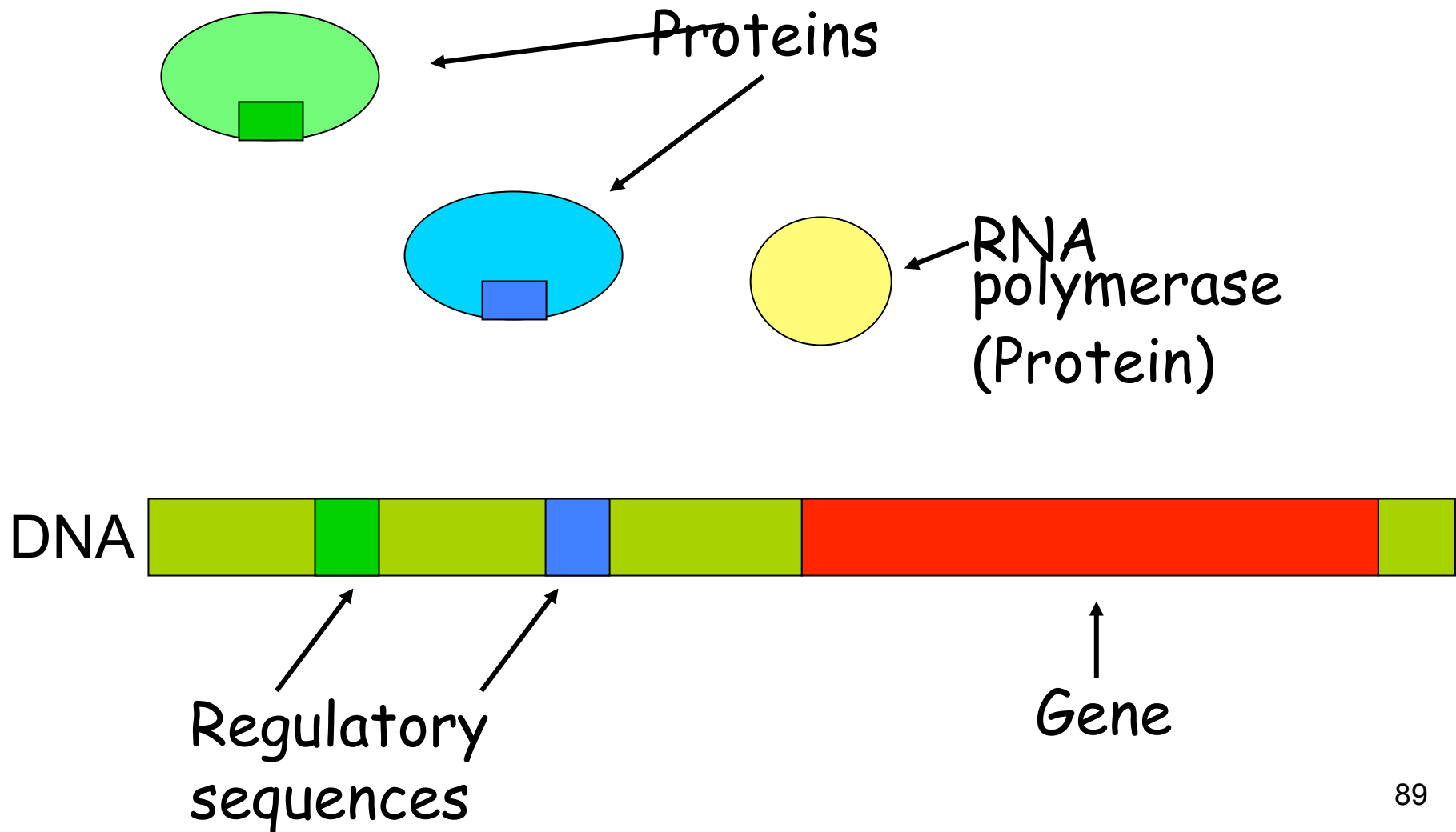


The Cell is a Computer in Soup

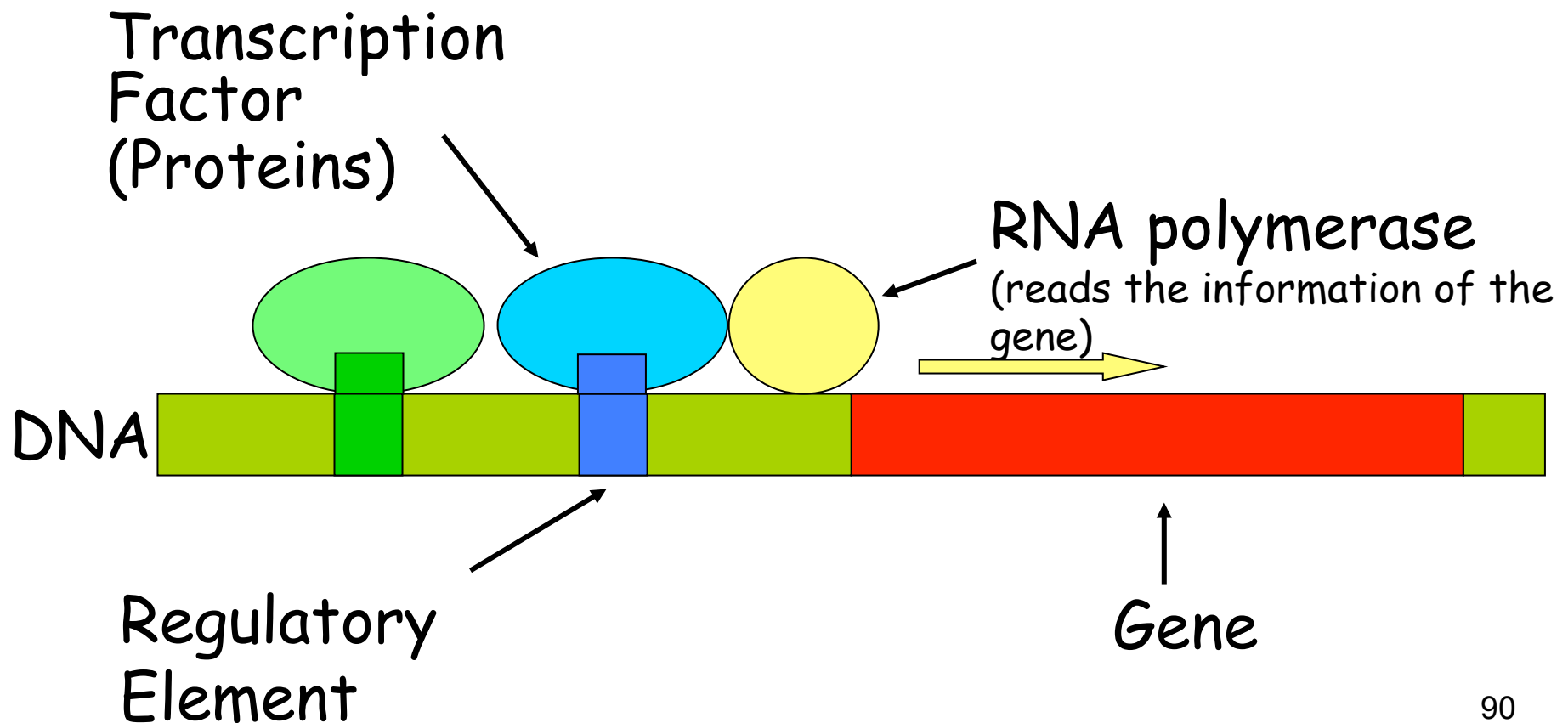


Idealized promoter for a gene involved in making hair. Proteins that bind to specific DNA sequences in the promoter region together turn a gene on or off. These proteins are themselves regulated by their own promoters leading to a gene regulatory network with many of the same properties as a neural network.

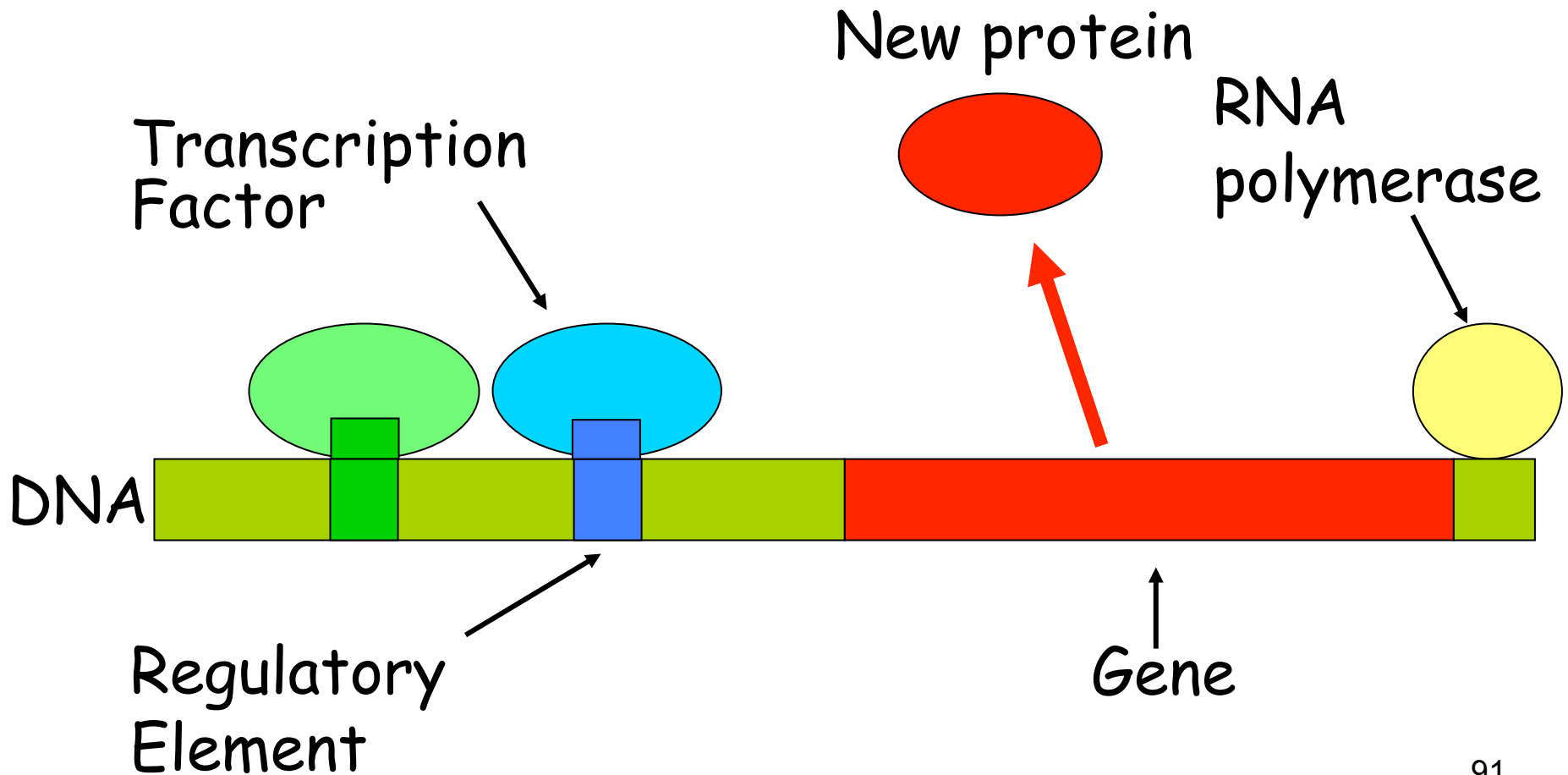
Regulation of Genes



Regulation of Genes



Regulation of Genes



Random Sample

atgaccgggatactgataccgtatTTGGCCTAGGCgtacacattagataaacgtatgaagtacgttagactcggcgccgccg
accctatTTTTTgagcagatttagtgacctggaaaaaaaaatttgagtacaaaactTTTCCgaatactgggcataaggtaca
tgagtatccctgggatgactTTTGGgaacactatagtgctctcccgattTTTgaatatgtaggatcattcgccagggtccga
gctgagaattggatgaccttgtaagtgtTTTCCacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
tcctTTTgCGgtaatgtGCCgggaggctggttacgtagggaagccctaacggacttaatggcccacttagtccacttatag
gtcaatcatgttcttTgtaatggattTTTaaactgagggcatagaccgcttggcgcacccaaattcagtgtgggCGagcgcaa
cggtTTTggcccttTtagaggccccgtactgatggaaactTTcaattatgagagagctaattctatcgCGtgCGtgttcat
aacttgagttggTTTcgaaaatgctctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcatttcaacgtatgccgaaccgaaaggggaag
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttctgggtactgatagca

Implanting Motif AAAAAAAGGGGGGGG

atgaccgggatactgatAAAAAAAGGGGGGGGggcgtacacattagataaacgtatgaagtacgtagactcggcgccgccg
accctatTTTTTgagcagatttagtgacctggaAAAAAatttgagtacaaactTTTccgaataAAAAAAAGGGGGGGG
tgagtatccctgggatgacttAAAAAAAGGGGGGGGtgctctcccgatTTTTgaatatgtaggatcattcggcagggtccga
gctgagaattggatgAAAAAAAGGGGGGGGtccacgcaatcgcgaaaccaacgcggaccCAAaggcaagaccgataaaggaga
tcctTTTgcgtaatgtgccgggaggctggttacgtaggaagccctaacggacttaataAAAAAAAGGGGGGGGcttatag
gtcaatcatgttcttTgtgaatggatttAAAAAAAGGGGGGGGgaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
cggtTTTggcccttTgttagaggccccgtAAAAAAAGGGGGGGGcaattatgagagagctaattctatcgcgtgcgtgttcat
aacttgagttAAAAAAAGGGGGGGGctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
ttggcccattggctaaaagcccaactgacaaatggaagatagaatccttTgcatAAAAAAAGGGGGGGGaccgaaaggggaag
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttAAAAAAAGGGGGGGG

Where is the Implanted Motif?

atgaccgggatactgataaaaaaagggggggggtacacattagataaacgtatgaagtacgtagactcggcgccgccg
accctatTTTTTgagcagatttagtgacctggaaaaaaatttgagtacaaaactTTTccgaataaaaaaaaggggggga
tgagtatccctgggatgacttaaaaaaagggggggtgctctcccgattTTTgaatatgtaggatcattcgccagggtccga
gctgagaattggatgaaaaaaagggggggtccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
tccTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaataaaaaaagggggggcttatag
gtcaatcatgttcttTgtaattgatttaaaaaaaggggggggaccgcttggcgcacccaaattcagtgtgggcgagcgcaa
cggTTTTggcccttTtagaggccccgtaaaaaaagggggggcaattatgagagagctaattctatcgcgtgcgtgttcat
aacttgagttaaaaaaaggggggctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcataaaaaaagggggggaccgaaaggggaag
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttaaaaaaaggggggga

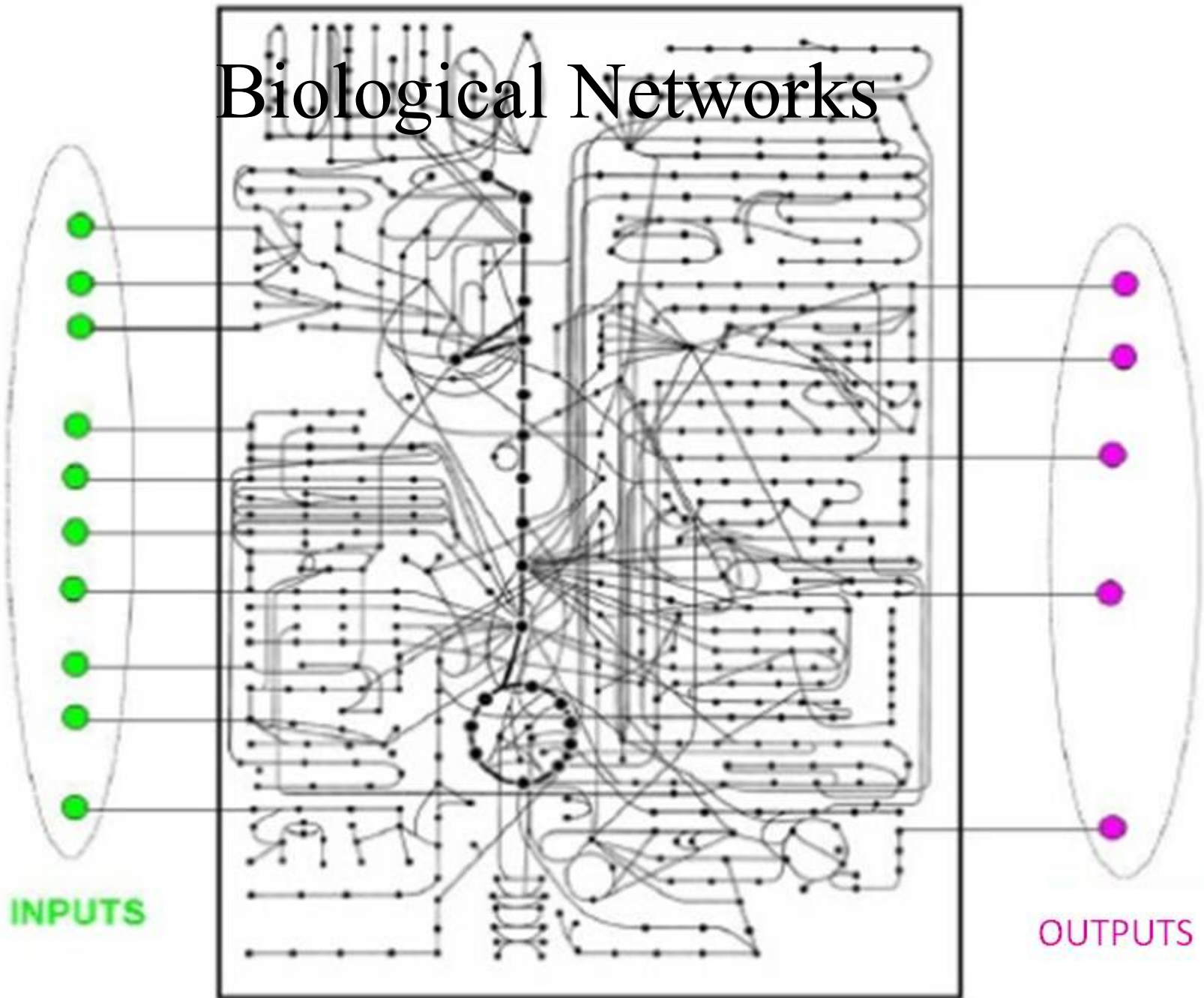
Implanting Motif AAAAAAGGGGGG with Four Mutations

atgaccgggatactgat**AgAAgAAAGGttGGG**ggcgtacacattagataaacgtatgaagtacgtagactcggcgccg
accctatTTTTTgagcagatttagtgacctggaAAAAAatttgagtacaaaactTTTccgaata**CAAtAAAACGGcGGG**a
tgagtatccctgggatgactt**AAAAtAATGGaGtGG**tgctctcccgattTTTgaatatgtaggatcattcgccagggtccga
gctgagaattggatg**CAAAAAAGGGattG**tccacgcaatcgcgaaccaacgcggaccCAAaggcaagaccgataaaggaga
tcctTTTgcgtaatgtgccgggaggctggttacgtagggaagccctaacggacttaat**AtAAtAAAGGaaGGG**cttatag
gtcaatcatgttcttTgtgaatggattt**AACAAtAAGGGctGG**gaccgcttggcgcacccaaattcagtgtgggcgagc
cgTTTTggcccttTtagaggccccgt**AtAAACAAAGGaGGGc**caattatgagagagctaattctatcgcgtgcgtgttcat
aacttgagtt**AAAAAAtAGGGaGcc**ctggggcacatacaagaggagtcttcttatcagttaatgctgtatgacactatgta
ttggcccattggctaaaagcccaacttgacaaatggaagatagaatccttgcat**ActAAAAGGaGcGG**accgaaaggggaag
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagctt**ActAAAAGGaGcGGa**

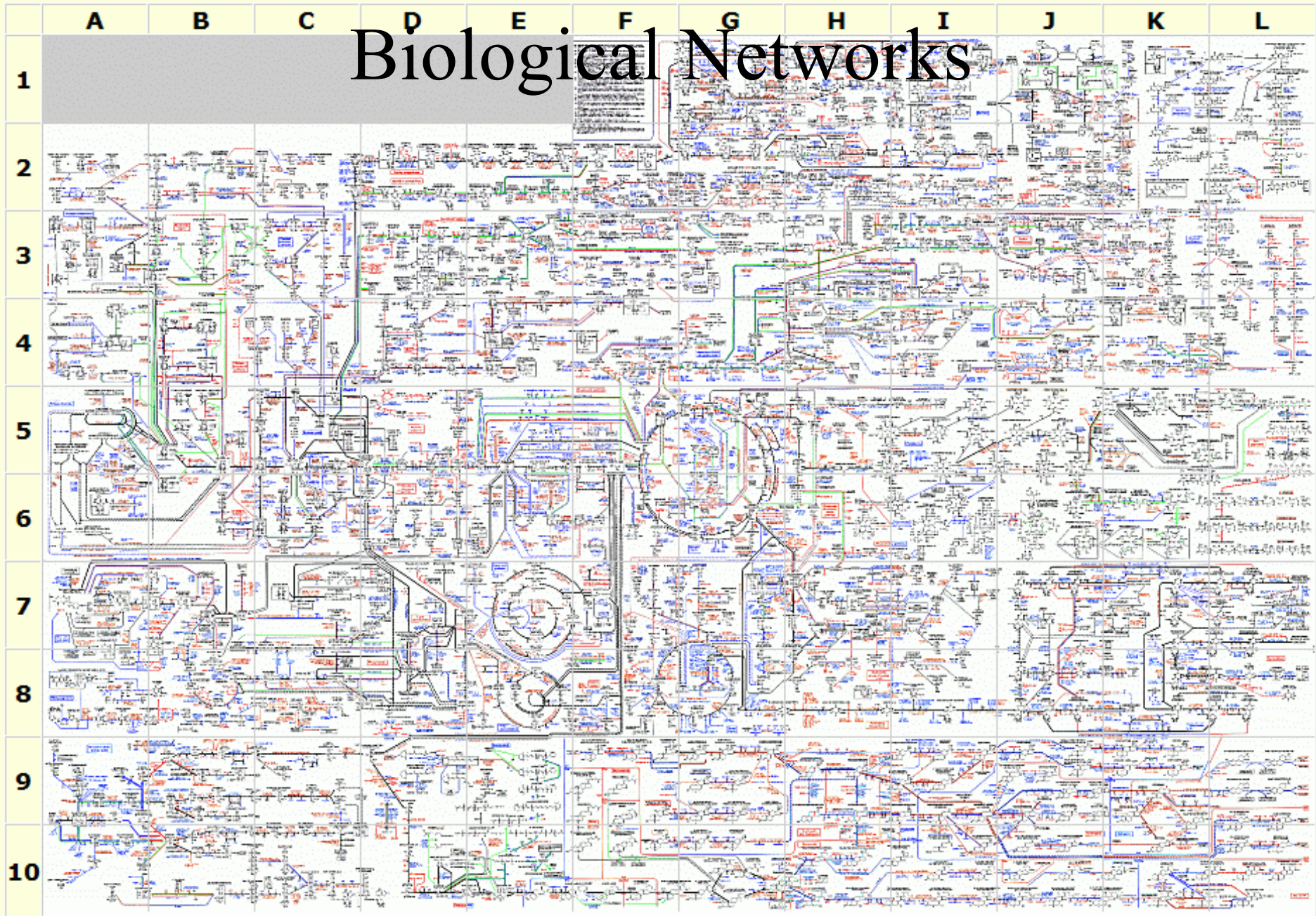
Where is the Motif???

atgaccgggatactgatagaagaaagggtggggcggtacacattagataaacgtatgaagtacgtagactcggcgccgcccg
accctatTTTTTgagcagatttagtgacctggaaaaaaaaatttgagtacaaaactTTTccgaatacaataaaacggcgggga
tgagtatccctgggatgacttaaataatggagtggtgctctcccgatTTTTgaatatgtaggatcattcgccagggtccga
gctgagaattggatgcaaaaaagggttgccacgcaatcgcgaaaccaacgcggacccaaaggcaagaccgataaaggaga
tccTTTTgcggtaatgtgccgggaggctggttacgtagggaagccctaacggacttaataataaaggaagggttatag
gtcaatcatgttcttgtgaatggatttaacaataagggtgggaccgcttggcgcacccaaattcagtggtggcgagcgcaa
cggTTTTggcccttggtagaggccccgtataaacaaggaggccaattatgagagagctaattctatcgcgtgcgtgttcat
aacttgagttaaaaataggagccctggggcacatacaagaggagtcttccttatcagttaatgctgtatgacactatgta
ttggcccattggctaaaagcccaactgacaaatggaagatagaatccttgcatactaaaaaggagcggaccgaaagggaag
ctggtgagcaacgacagattcttacgtgcattagctcgcttccggggatctaatagcacgaagcttactaaaaaggagcgga

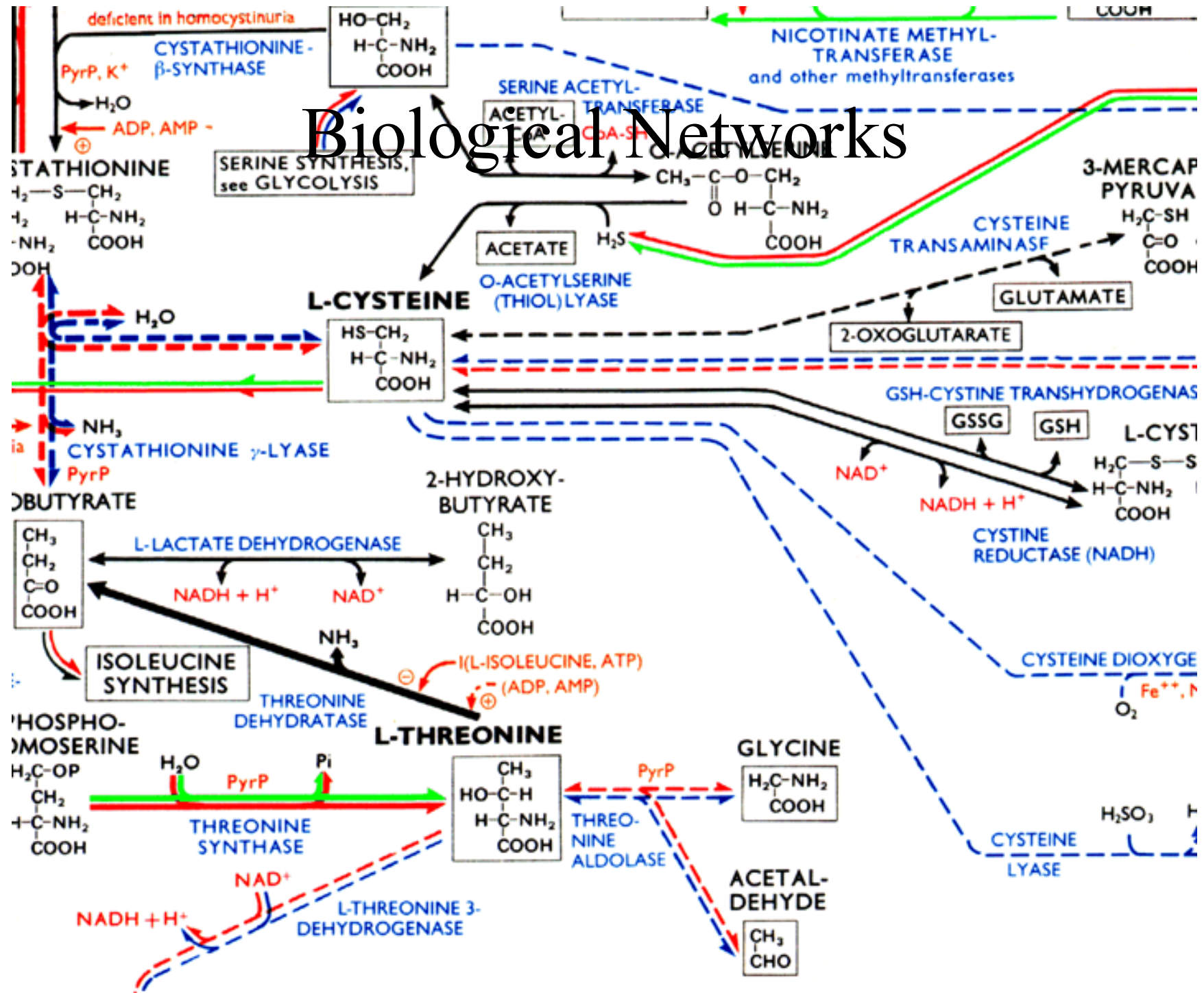
Biological Networks



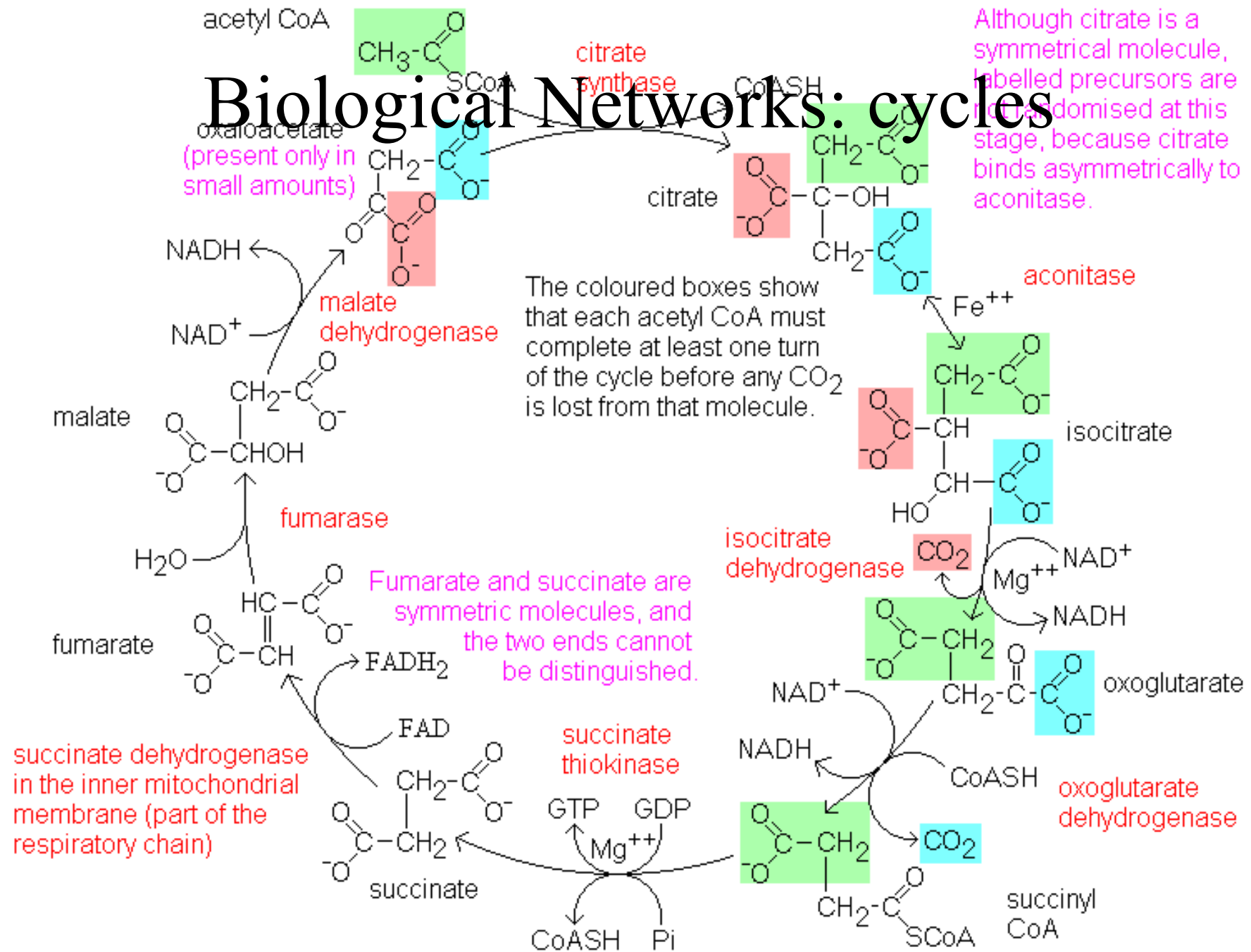
Biological Networks



Biological Networks

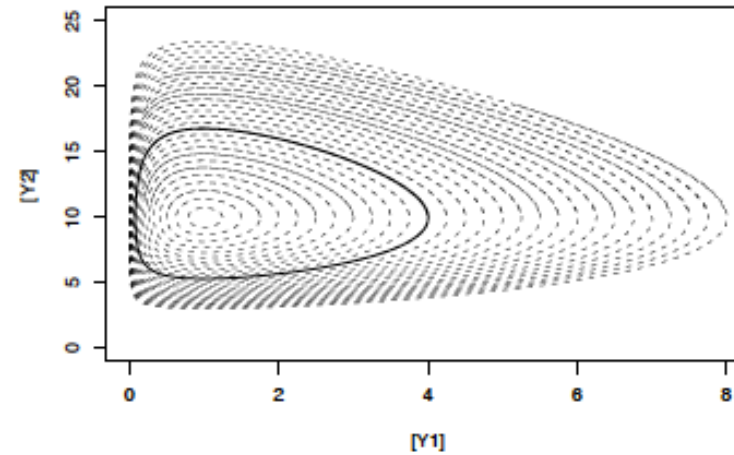
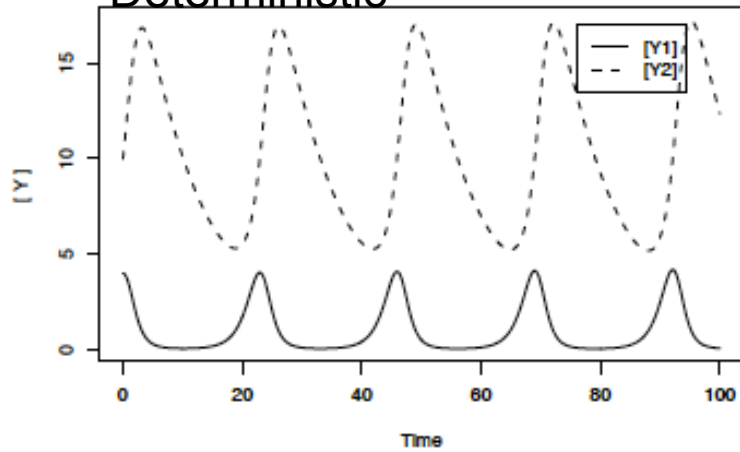


Biological Networks: cycles



Stochastic versus Deterministic

Deterministic



Stochastic

