

Sequence Analysis (part II)

BBSI 2006: Lecture #($\chi+2$)

Takis Benos (2006)



BBSI 2006 30-MAY-2006

© 2006 P. Benos 1

Outline

- Sequence variation
- Distance measures
- Scoring matrices
- Pairwise alignments (global, local)
- Database searches (BLAST, FastA)
- Multiple sequence alignments



BBSI 2006 30-MAY-2006

© 2006 P. Benos 2

Sequence Variations



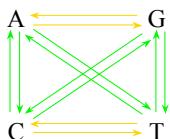
BBSI 2006 30-MAY-2006

© 2006 P. Benos 3

Sequence variation

- Base mutations: the source of sequence variation

Purines



Pyrimidines

Transitions
Transversions



BBSI 2006 30-MAY-2006

© 2006 P. Benos 4

Sequence variation (cntd)

```
tggagctAtt attgctaagt Acatattacc ccctgaagtt aatgGatcaa tcaagagaga 120
tgtggctgtt aatgaaTcgC Cttattgaat taacagggtt gatcgttctt gtcgttcaag 180
    M   R   L   I   E   L
tcattttctt tgccgtggc agtcacattg acaactatca gcacactgaa cagagtgtt 240
cggtacaaca caagtaagct ctgcaactgt ggagcgacat gtcgtccgtc cggtgcattg 300
```

silent missense nonsense



```
tggagctGtt attgctaagt Tacatattacc ccoggaagtt aatgAatcaa tcaagagaga 120
tgtggctgtt aatgaaCcgt Cttattgaat Aaacagggtt gatcgttctt gtcgttcaag 180
    M   R   V   I   E
tcattttctt tgccgtggc agtcacattg acaactatca gcacactgaa cagagtgtt 240
cggtacaaca caagtaagct ctgcaactgt ggagcgacat gtcgtccgtc cggtgcattg 300
```



BBSI 2006 30-MAY-2006

© 2006 P. Benos 5

Sequence variation (cntd)

```
tggagctAtt attgctaagt Acatattacc ccctgaagtt aatgGatcaa tcaagagaga 120
tgtggctgtt aatgaaTcgC Cttattgaat taacagggtt gatcgttctt gtcgttcaag 180
    M   R   L   I   E   L
tcattttctt tgccgtggc agtcacattg acaactatca gcacactgaa cagagtgtt 240
cggtacaaca caagtaagct ctgcaactgt ggagcgacat gtcgtccgtc cggtgcattg 300
```



deletion

```
tggagctGtt attgctaagt Tacatattacc ccctgaagtt aatgAatcaa tcaagagaga 120
tgtggctgtt aatgaaCcgt Cttattgaat taacagggtt gatcgttctt gtcgttcaag 180
    M   R   V   I   E
tcattttctt tgccgtggc agtcacattg acaactatca gcacactgaa cagagtgtt 240
cggtacaaca caagtaagct ctgcaactgt ggagcgacat gtcgtccgtc cggtgcattg 300
```



BBSI 2006 30-MAY-2006

© 2006 P. Benos 6

Sequence variation (cntd)



Figure 2. Average rates of substitution in different parts of genes and in pseudogenes.

BBSI 2006 30-MAY-2006

© 2006 P. Benos 7

Source: Li & Graur "Fundamentals of Molecular Evolution", 1991, Sinauer Assoc.

Distance measures



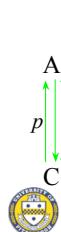
BBSI 2006 30-MAY-2006

© 2006 P. Benos 8

Nucleic acid distances

- No selection - no correction:

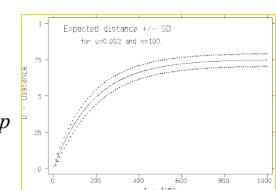
$$D = k / N$$



BBSI 2006 30-MAY-2006

© 2006 P. Benos 9

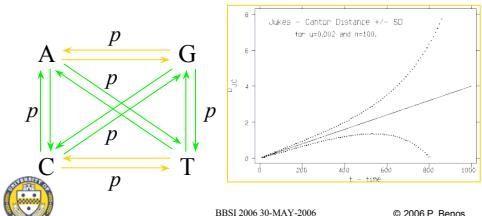
Source: <http://hebc.biology.memaster.ca/721/distance/distance.html>



Nucleic acid distances (cntd)

- Jukes-Cantor correction:

$$D_{JC} = -0.75 \ln (1 - D/0.75)$$



BBSI 2006 30-MAY-2006

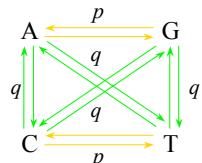
© 2006 P. Benos 10

Source: <http://helix.biology.msu.edu/cat72/distance/distance.html>

Nucleic acid distances (cntd)

- Kimura's 2-parameter model:

$$D_{K2P} = -0.5 \ln (1 - 2P - 2Q) - 0.25 \ln (1 - 2Q)$$



BBSI 2006 30-MAY-2006

© 2006 P. Benos 11

Source: <http://helix.biology.msu.edu/cat72/distance/distance.html>

Scoring matrices



BBSI 2006 30-MAY-2006

© 2006 P. Benos 12

Nucleic acid distances (cntd)

- Nucleotide substitution matrices.

A	T	C	G
A	1	0	0
T	0	1	0
C	0	0	1
G	0	0	0

Identity

A	T	C	G
A	5	-4	-4
T	-4	5	-4
C	-4	-4	5
G	-4	-4	5

BLAST

A	T	C	G
A	0	5	5
T	5	0	1
C	5	1	0
G	1	5	5

Transition/
Transversion



BBSI 2006 30-MAY-2006

© 2006 P. Benos 13

Amino acid distances: PAM

- Percent Accepted Mutations (PAM) matrices:
 - Frequency substitution matrix from aligned sequences (Dayhoff, 1978).
 - $M(i,j)$: no. of a.a. i to j mutations
 - 71 groups of closely related proteins (*why?*); 1,572 changes.
 - PAM_n : the aligned sequences have n a.a. substitutions per 100 residues.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 14

Amino acid distances: PAM (cntd)

- Assumptions of the PAM model:
 - Replacement at any site depends only on the a.a. on that site, given the mutability table.
 - Sequences in the training set (and those compared) have average a.a. composition.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 15

Amino acid distances: PAM (cntd)

Table 1 - The log odds matrix for 250 PAMs (multiplied by 10)																							
	A	C	D	E	F	G	H	I	K	L	M	N	P	R	S	T	V	W	Y				
A	2	-1.5	-0.5	-4	1	-3	-1.5	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
C	-1.5	2	-0.5	-4	1	-3	-1.5	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
D	-4	-0.5	2	-1.5	1	-3	-1.5	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
E	1	-3	-1.5	2	-0.5	-4	1	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
F	-3	-1.5	-0.5	-4	2	1	-1.5	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
G	-1.5	0	-4	1	-3	-0.5	2	-1	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
H	0	-4	-0.5	-3	1	-2	-1.5	2	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
I	-3	-0.5	-4	1	-2	-1.5	-1	2	-1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
K	-1.5	-1	-2	-1	-1	2	-0.5	-1	1	-2	-1	0	1	-1	-1	-1	-1	-1	-1				
L	-2	-1	-1	-1	-1	-0.5	-1	-1	2	-0.5	-1	0	1	-1	-1	-1	-1	-1	-1				
M	-1	-1	-2	-1	-1	-1	-0.5	-1	-1	2	-0.5	0	1	-1	-1	-1	-1	-1	-1				
N	-0.5	-1	-1	-1	-1	-1	-1	-1	-1	-0.5	2	0	1	-1	-1	-1	-1	-1	-1				
P	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	0	1	-1	-1	-1	-1	-1	-1				
R	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-0.5	2	0	-1	-1	-1	-1	-1	-1				
S	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	0	-1	-1	-1	-1	-1	-1				
T	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-0.5	2	0	-1	-1	-1	-1	-1				
V	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	0	-1	-1	-1	-1	-1				
W	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-0.5	2	0	-1	-1	-1	-1				
Y	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	-1	2	0	-1	-1	-1	-1				

$$\text{Score}(i,j) = \log_{10} M(i,j)/f(i)$$

BBSI 2006 30-MAY-2006

© 2006 P. Benos 16

Source: <http://hdbiology.msu.edu/721/distance/node9.html>

Amino acid distances: PAM (cntd)

- Sources of error in the PAM model:
 - Many proteins depart from the average a.a. composition.
 - The a.a. composition can vary even within a protein (e.g., transmembrane proteins).
 - A.a. positions are not “mutated” equally probably; especially in long evolutionary distances.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 17

Amino acid distances: PAM (cntd)

- Sources of error in the PAM model (cntd):
 - Rare replacements are observed too infrequently and...
 - ...errors in PAM1 are magnified in PAM250.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 18

A.a. distances: BLOSUM

- **Blocks Substitution Matrices (BLOSUM):**
 - Log-likelihood matrix (Henikoff & Henikoff, 1992)
 - BLOCKS database of aligned sequences used as primary source set.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 19

A.a. distances: BLOSUM (cntd)

AKAGDA---GGCDA
DRALDAFG--GSSDA
GKLGDAI---GSSAF
AKAGGA---GGTAG
CRIGFRC-DGTTDH
AKAKDA--DHSSCI

$$\text{Score}(i,j) = 2 \log_2 q_{ij} / e_{ij}$$

$$e_{ij} = p_i^2 \quad \text{for } i=j$$
$$e_{ij} = 2 p_i p_j \quad \text{for } i \neq j$$

$$p_i = 0.5 (q_{ii} + \sum q_{ij})$$



BBSI 2006 30-MAY-2006

© 2006 P. Benos 20

A.a. distances: BLOSUM (cntd)

- Weighted contribution of similar(*) sequences in order to reduce redundancy.
- BLOSUM62 is more closely related to PAM120.

(*) $n\%$ similar; the n in BLOSUM n



BBSI 2006 30-MAY-2006

© 2006 P. Benos 21

A.a. distances: BLOSUM (cntd)

Table 2 - The log odds matrix for BLOSUM 62																								
A	C	G	T	R	I	V	S	H	I	K	L	M	N	P	R	S	T	V	W	Y				
C	-9	-5	-4	-3	-2	-1	0	1	-1	-1	-1	-2	-1	-1	-1	-1	1	0	-3	-2				
G	9	-5	-4	-3	-2	-1	0	1	-1	-1	-1	-2	-1	-1	-1	-1	1	0	-3	-2				
T	5	9	-5	-4	-3	-2	-1	0	1	-1	-1	-2	-1	-1	-1	-1	1	0	-3	-2				
R	4	5	9	-5	-4	-3	-2	-1	0	1	-1	-1	-2	-1	-1	-1	1	0	-3	-2				
I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
S	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
H	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
K	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
L	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
M	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
N	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
P	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
R	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
S	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
T	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0				

Source: <http://helix.biology.mcmaster.ca/72/distance/index.html>



BBSI 2006 30-MAY-2006

© 2006 P. Benos 22

Substitution matrices: comparison

- PAM vs BLOSUM

BLOSUM 80	BLOSUM 62	BLOSUM 45
PAM 1	PAM 120	PAM 250
Less divergent	← → More divergent	

Source: <http://www.ncbi.nlm.nih.gov/Education/BLASTInfo/Scoring2.html>

- Matrices of choice:

- BLOSUM62: the all-weather matrix
- PAM250: for distant relatives



BBSI 2006 30-MAY-2006

© 2006 P. Benos 23

Substit. matrices: comparison (cntd)

- PAM vs BLOSUM (cntd)

- Lower PAM/higher BLOSUM matrices identify shorter local alignments of highly similar sequences
- Higher PAM/lower BLOSUM matrices identify longer local alignments of more distant sequences



BBSI 2006 30-MAY-2006

© 2006 P. Benos 24

Substit. matrices: comparison (cntd)

PAM10



BBSI 2006 30-MAY-2006

© 2006 P. Benos 25

Substit. matrices: comparison (cntd)

PAM250



BLOSUM62



Substit. matrices: comparison (cntd)

A	4
R	-1 5
N	-2 0 6
D	-2 -2 1 6
C	0 -3 -3 -3 9
Q	-1 1 0 0 -3 5
E	-1 0 2 -4 2 5
G	0 -2 0 -1 -3 -2 6
H	-0 1 -1 -3 0 0 -2 8
I	-1 -3 -3 -3 -1 -3 -3 -4 -3 4
L	-1 -2 -3 -4 -1 -2 -3 -4 -3 2 4
K	-1 2 0 -1 -3 1 1 -2 -1 -3 -2 5
M	-1 -2 -3 -3 -1 0 -2 -3 -2 1 2 -1 5
F	-2 -3 -3 -3 -2 -3 -3 -1 0 0 -3 0 6
P	-1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -1 -2 -4 7
S	1 -1 0 -1 0 0 0 -1 -2 -2 0 -1 -2 -1 4
T	0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -2 -1 1 5
W	-3 -3 -4 -4 -2 -2 -3 -2 -3 -2 -3 -1 1 -4 -3 -2 11
Y	-2 -2 -2 -3 -2 -2 -3 2 -1 -1 -2 -1 3 -3 -2 -2 7
V	0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4
A R N D C Q E G H I L K M F P S T W Y V	

Pairwise alignments



BBSI 2006 30-MAY-2006

© 2006 P. Benos 28

Alignment: the problem

Given two sequences, S and T , and a scoring matrix find their relative arrangement with the highest “score”.

Seq. #1: G A A T T C A G T T A
Seq. #2: G G A T C G A



BBSI 2006 30-MAY-2006

© 2006 P. Benos 29

Alignment: the problem (cntd)

G A A T T C A G T T A
| |
G G A T C G A

G A A T T C A G T T A
| | | |
G G A T C G A

G A A T T C - A G T T A
| | | | |
G G A - T C G A



BBSI 2006 30-MAY-2006

© 2006 P. Benos 30

Alignment: the problem (cntd)

- Scoring schemes: three possible situations...
 - Match **REWARD!!**
 - Mismatch **Penalise???**
 - Gap
 - Gap initiation
 - Gap extension

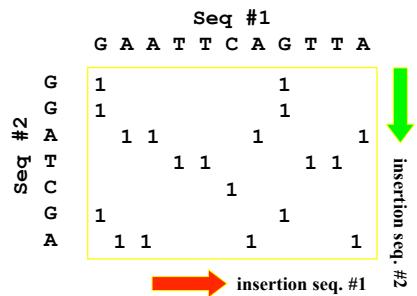
How much??



BBSI 2006 30-MAY-2006

© 2006 P. Benos 31

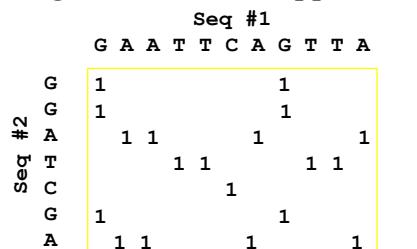
Alignment: a naïve approach



BBSI 2006 30-MAY-2006

© 2006 P. Benos 32

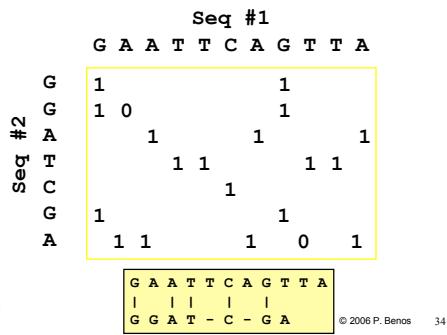
Alignment: a naïve approach



G A A T T C - A G T T A
| | | | | |
G G A - T - C G A

© 2006 P. Benos 33

Alignment: a naïve approach (cntd)



Alignment: adding scores

The formula:

$$\begin{aligned}
 M_{ij} = & \text{MAXIMUM} \{ \\
 & M_{i-1, j-1} + S_{ij} \text{ (match/mismatch in the diagonal),} \\
 & M_{i-1, j} + w \text{ (gap in sequence #1),} \\
 & M_{i, j-1} + w \text{ (gap in sequence #2)} \\
 \}
 \end{aligned}$$

- In the following example, the score for match is 1 and for mismatch and gap is 0.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 35

Alignment: adding scores (cntd)

- In each step we need to keep track only the scores of the (i,j) position and its immediate neighbours: $(i-1,j-1)$, $(i-1,j)$ and $(i,j-1)$.
- We backtrack from the right-down corner to find the actual alignment.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 36

Alignment: adding scores (cntd)											
	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	+	1								
A	0										
T	0										
C	0										
G	0										
A	0										

G A A T T T C G T T T A											
	G	A	A	T	T	T	C	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	1									
A	0	+	1								
T	0		+	1							
C	0			+	1						
G	0				+	1					
A	0					+	1				

Alignment: adding scores (cntd)										
	G	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	1	1
G	0	1	1	1						
A	0	1	1							
T	0	1								
C	0	1								
G	0	1								
A	0	1								

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

Alignment: adding scores (cntd)										
	G	A	A	T	T	C	A	G	T	A
G	0	0	1	1	1	1	1	1	1	1
G	0	0	1	1	1	1	1	2	2	2
A	0	1	1	2	2	2	2	2	2	3
T	0	1	2	2	3	3	3	3	3	3
C	0	1	2	2	3	4	4	4	4	4
G	0	1	2	2	3	3	4	5	5	5
A	0	1	2	3	3	3	4	5	5	6
(Seq #1)										
Alignment:										
(Seq #2)										
	G	A	A	T	T	C	A	G	T	A
G	0	0	0	0	0	0	0	0	0	0
G	0	1	1	1	1	1	1	1	2	2
A	0	1	2	2	2	2	2	2	2	2
T	0	1	2	2	3	3	3	3	3	3
C	0	1	2	2	3	3	4	4	4	4
G	0	1	2	2	3	3	4	4	5	5
A	0	1	2	3	3	3	4	5	5	6

Alignment: adding scores (cntd)

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

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

Alignment:

(Seq #1) G A A T T C A G T T A

(Seq #2) G G A - T C G - - A

Alignment: another example

Alignment: another example (cntd)

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	2									
G	0	2									
A	0										
T	0										
C	0										
G	0										
A	0										



Source:
http://www.sbc.su.se/~per/molbioinfo2001/dynprog/adv_dynamic.html
 BBSI 2006 30-MAY-2006 © 2006 P. Benos 43

Alignment: another example (cntd)

	G	A	A	T	T	C	A	G	T	T	A
G	0	0	0	0	0	0	0	0	0	0	0
G	0	2	1	-1							
A	0	0	4								
T	0	-1	2								
C	0	-1	0								
G	0	2	0								
A	0	0	4								



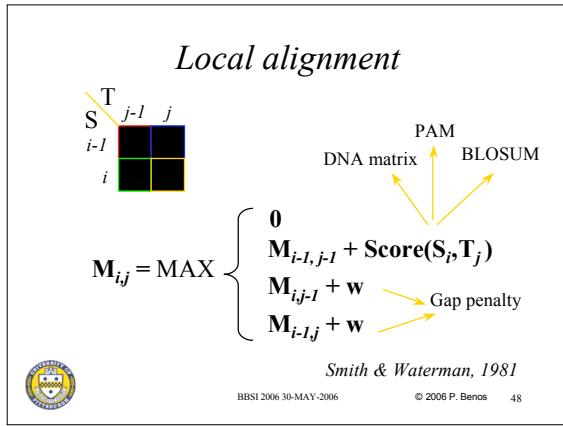
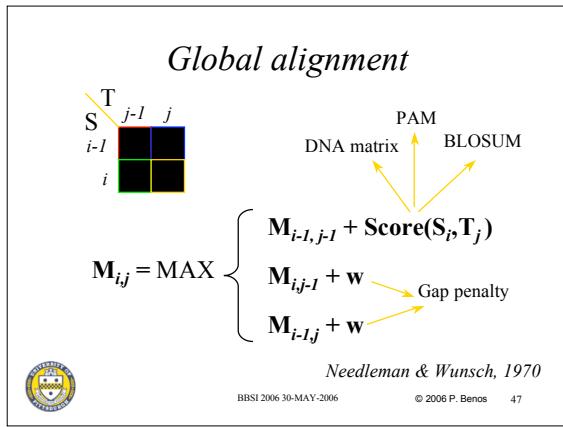
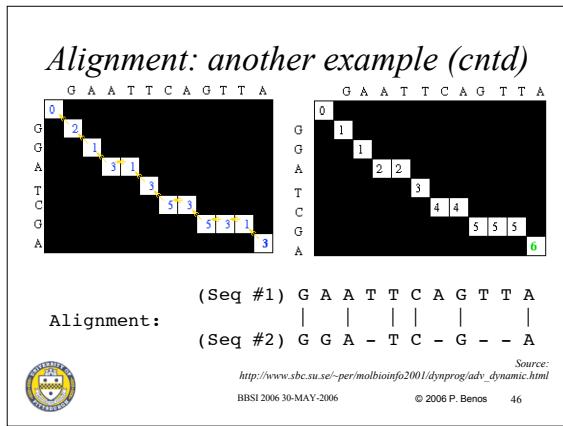
Source:
http://www.sbc.su.se/~per/molbioinfo2001/dynprog/adv_dynamic.html
 BBSI 2006 30-MAY-2006 © 2006 P. Benos 44

Alignment: another example (cntd)

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



Source:
http://www.sbc.su.se/~per/molbioinfo2001/dynprog/adv_dynamic.html
 BBSI 2006 30-MAY-2006 © 2006 P. Benos 45



Local alignment

Given two sequences, S and T , find two subsequences, s and t , whose alignment has the highest “score” amongst all subsequence pairs.

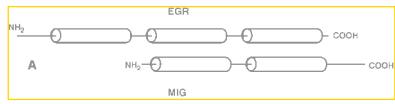
Why do we need local alignment,
if we have the global one?



BBSI 2006 30-MAY-2006

© 2006 P. Benos 49

Local alignment: an example



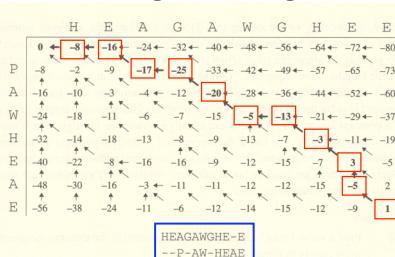
		EGR	
	NH ₂	—	—
EGR4_HUMAN	R	[FACPVESQRVSFAKSSDELNHHLR[1]	TGKEP [PQCNICLQAFNSR[DBL-TSIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR4_RAT	R	[FACPVESQRVTFAKSSDELNHHLR[1]	TGKEP [PQCNICLQAFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR3_HUMAN	R	[GACVAGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR3_MOUSE	R	[GACVAGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR1_HUMAN	R	[YACPVESCDRPFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR1_MOUSE	R	[YACPVESCDRPFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR1_CHICKEN	R	[YACPVESCDRPFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR1_BRARE	R	[YACPVESCDRPFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR2_RAT	R	[YPCPAEGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR2_MOUSE	R	[YPCPAEGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR2_CHICKEN	R	[YPCPAEGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR2_MOUSE	R	[YPCPAEGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
EGR2_HUMAN	R	[YPCPAEGCDDRFVFSKSSDELTHHLR[1]	TGKEP [PQCNICMGSFNSR[DBL-TTIVRTH] TGKEP [FACDV--CGRFAM[DEE[RSVHV]
MIG1_KLULA	--	--	--
MIG1_KLUMA	--	--	--
MIG1_YEAST	--	--	--



BBSI 2006 30-MAY-2006

© 2006 P. Benos 50

Local vs. global alignment



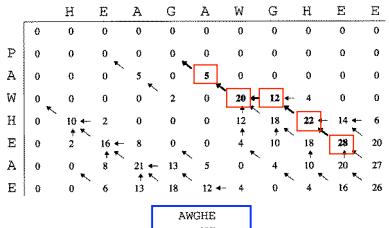
Source: Durbin et al "Biological Sequence Analysis", 1998,
Cambridge University Press



BBSI 2006 30-MAY-2006

© 2006 P. Benos 51

Local vs. global alignment (cntd)



Source: Durbin et al "Biological Sequence Analysis", 1998,
Cambridge University Press

BBSI 2006 30-MAY-2006

© 2006 P. Benos 52



Local alignment (cntd)

- Characteristics of local alignments:
 - The alignment can start/end at any point in the matrix.
 - No negative scores.
 - The mean value of the scoring matrix (e.g. PAM, BLOSUM) should be negative.
 - There should be positive scores in the scoring matrix.



BBSI 2006 30-MAY-2006

© 2006 P. Benos 53