# Sequence Alignment (chapter 6)

- p The biological problem
- p Global alignment
- p Local alignment
- p Multiple alignment

# Local alignment: rationale

 $\, {\rm p} \,$  Otherwise dissimilar proteins may have local regions of similarity

-> Proteins may share a function

Human bone morphogenic protein receptor type II precursor (left) has a 300 aa region that resembles 291 aa region in TGF-ß receptor (right). The shared function here is protein kinase.



Local alignment: rationale
 A
 B
 Clobal alignment would be inadequate
 Problem: find the highest scoring *local* alignment between two sequences
 Previous algorithm with minor modifications solves this problem (Smith & Waterman 1981)



# $\begin{array}{l} \textbf{Scoring local alignments} \\ \textbf{A} = \textbf{a}_1\textbf{a}_2\textbf{a}_3...\textbf{a}_n, \textbf{B} = \textbf{b}_1\textbf{b}_2\textbf{b}_3...\textbf{b}_m \\ \textbf{Let I and J be intervals (substrings) of A and B, respectively: \\ I \subset A \quad J \subset B \\ \end{array}$ $\begin{array}{l} \textbf{Best local alignment score:} \\ M(A,B) = \max\{S(I,J): I \subset A, J \subset B\} \\ \textbf{where S(I, J) is the alignment score for substrings I and J. \end{array}$







Local alignment: example													
M			0	1	2	3	4	5	6	7	8	9	10
M <sub>i,j</sub> = max { M <sub>i-1,i-1</sub> + s(a <sub>i</sub> ,			-	G	G	c	Т	С	A	А	Т	С	A
b <sub>i</sub> ), Μ δ	0	-	Ο,	0	0	0	0	0	0	0	0	0	0
$M_{i-1,j} = 0,$ $M_{i+1} = \delta,$	1	А	0 -	o									
0	2	С	0										
}	3	С	0										
	4	Т	0										
	5	А	0										
Scoring (for example)	6	А	0										
Match: +2	7	G	0										
Indel: -2	8	G	0										
208													

Local align	m	en	it:	ех	(91	mp	ole	ý					
M – max (			0	1	2	3	4	5	6	7	8	9	10
$M_{i,j} = Max (M_{i-1,i-1} + s(a_i))$			-	G	G	c	Т	С	А	A	Т	c	А
b <sub>i</sub> ),	0	-	0	0	0	0	0	0、	0	0	0	0	0
$M_{i-1,j} = 0,$ $M_{i-1,j} = \delta.$	1	А	0	0	0	0	0	0-	2				
0	2	С	0										
}	3	С	0										
	4	Т	0										
	5	Α	0										
Scoring (for example)	6	Α	0										
Match: +2	7	G	0										
Mismatch: -1	8	G	0										

Local alignment: example													
Optimal local			0	1	2	3	4	5	6	7	8	9	10
alignment:			-	G	G	c	Т	С	A	А	Т	c	А
CT-AA	0	-	0	0	0	0	0	0、	0、	0	0	0、	0
СТСАА	1	А	0	0	0、	0	0、	0	2、	2、	0、	0	2
	2	С	0	0	0,	2、	0、	2、	0	1	1、	2	0
	3	С	0	0	0	2	1	2	1	Ο,	0	3	1
	4	Т	0	0	0	0、	4	2.	1、	0	2、	1	2
	5	А	0	0	0	0	2	3、	4	3-	1	1、	3
Searing (for example)	6	А	0、	0	0	0	0	1	5	6	4 -	2	3
Match: +2	7	G	0、	2、	2	0	0	0	3	4	5	3-	1
Mismatch: -1	8	G	0	2	4 -	2	0	0	1	2	3	4-	2
Indel: -2													
210													







# Gaps in alignment

p Gap is a succession of indels in alignment

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

- p Previous model scored a length k gap as  $w(k) = -k\delta$
- P Replication processes may produce longer stretches of insertions or deletions
  - n In coding regions, insertions or deletions of codons may preserve functionality

# Gap open and extension penalties (2)

P We can design a score that allows the penalty opening gap to be larger than extending the gap:

 $w(k) = -\alpha - \beta(k-1)$ 

p Gap open cost α, Gap extension cost β
 p Alignment algorithms can be extended to use w(k) (not discussed on this course)

# Amino acid sequences

- ${\ensuremath{{\bf p}}}$  We have discussed mainly DNA sequences
- p Amino acid sequences can be aligned as well
- P However, the design of the substitution matrix is more involved because of the larger alphabet
- p More on the topic in the course Biological sequence analysis

# Demonstration of the EBI web site

- p European Bioinformatics Institute (EBI) offers many biological databases and bioinformatics tools at http://www.ebi.ac.uk/
  - n Sequence alignment: Tools -> Sequence Analysis -> Align

# Sequence Alignment (chapter 6)

- p The biological problem
- p Global alignment
- p Local alignment
- p Multiple alignment

# Multiple alignment

- P Consider a set of n sequences on the right
   n Orthologous sequences from
- different organisms n Paralogs from multiple duplications
- P How can we study
- relationships between these sequences?

aggcgagctgcgagtgcta cgttagattgacgctgac ttccggctgcgac gacacggcgacgga agtgtgcccgacgaggaggac gcgggctgtgagcgcta aagcggcctgtgtgcccta atgctgctgccgagtgc agtcgagcccgagtgc agtccgagtcc actcggtgc

# Optimal alignment of three sequences

- ${\sf p}$  Alignment of A =  $a_1a_2...a_i$  and B =  $b_1b_2...b_j$  can end either in (-,  $b_j$ ),  $(a_i, \, b_j)$  or  $(a_i,$  -)
- $p 2^2 1 = 3$  alternatives
- $\begin{array}{l} \mbox{$\mathsf{p}$} \mbox{Alignment of A, B and C = c_1c_2...c_k can end in 2^3 1 ways: (a_i, -, -), (-, b_j, -), (-, -, c_k), (-, b_j, c_k), (a_i, -, c_k), (a_i, b_j, -) or (a_i, b_j, c_k) \end{array} } \end{array}$
- P Solve the recursion using three-dimensional dynamic programming matrix: O(n<sup>3</sup>) time and space
- p Generalizes to n sequences but impractical with even a moderate number of sequences

# Multiple alignment in practice

- p In practice, real-world multiple alignment problems are usually solved with heuristics
- Progressive multiple alignment
- n Choose two sequences and align them n Choose third sequence w.r.t. two previous sequences
- and align the third against them
- n Repeat until all sequences have been aligned
- n Different options how to choose sequences and score alignments
- n Note the similarity to Overlap-Layout-Consensus

# Multiple alignment in practice

- Profile-based progressive multiple alignment: CLUSTALW
  - n Construct a distance matrix of all pairs of sequences using dynamic programming
  - Progressively align pairs in order of decreasing similarity
  - n CLUSTALW uses various heuristics to contribute to accuracy

# Additional material

- p R. Durbin, S. Eddy, A. Krogh, G. Mitchison: Biological sequence analysis
- p N. C. Jones, P. A. Pevzner: An introduction to bioinformatics algorithms
- p Course Biological sequence analysis in period II, 2008

# Rapid alignment methods: FASTA and BLAST

- p The biological problem
- p Search strategies
- p FASTA
- p BLAST

# <section-header> Description Global and local alignment algorithms are slow in practice Consider the scenario of aligning a query sequence against a large database of sequences New sequence with unknown function Mean sequence with sequences (61 132 599 sequences) abases (62 853 668 564 bases)

### Problem with large amount of sequences

- p Exponential growth in both number and total length of sequences
- P Possible solution: Compare against model organisms only
- p With large amount of sequences, chances are that matches occur by random
   n Need for statistical analysis

# Rapid alignment methods: FASTA and BLAST

- p The biological problem
- P Search strategies
- p FASTA
- p BLAST

# FASTA

- FASTA is a multistep algorithm for sequence alignment (Wilbur and Lipman, 1983)
- P The sequence file format used by the FASTA software is widely used by other sequence analysis software
- p Main idea:
  - Choose regions of the two sequences (query and database) that look promising (have some degree of similarity)
  - n Compute local alignment using dynamic programming in these regions

# FASTA outline

### p FASTA algorithm has five steps:

- n 1. Identify common k-words between I and J n 2. Score diagonals with k-word matches,
- identify 10 best diagonals n 3. Rescore initial regions with a substitution
- score matrix
- n 4. Join initial regions using gaps, penalise for gaps
- n 5. Perform dynamic programming to find final alignments



- p How to speed up the computation?
  n Find ways to limit the number of pairwise comparisons
- p Compare the sequences at word level to find out common words
  - n Word means here a k-tuple (or a k-word), a substring of length k



# Analyzing the word content

- p There are n-k+1 k-words in a string of length n
- p If at least one word of I is not found from another string J, we know that I differs from J
- p Need to consider statistical significance: I and J might share words by chance only
- $_{\rm p}$  Let n=|I| and m=|J|



# Common k-words

- p Number of common k-words in I and J can be computed using  $L_{w}(I)$  and  $L_{w}(J)$
- ${\sf p}$  For each word w in I, there are  $|L_w(J)|$  occurences in J
- p Therefore I and J have  $\sum_{w} |L_w(I)| |L_w(J)|$  common words
- p This can be computed in  $O(n + m + 4^k)$  time
  - n O(n + m) time to build the lists
  - $n O(4^k)$  time to calculate the sum (in DNA
  - strings)

Common k-words									
p I = GCATC p J = CCATC	G <mark>GC</mark> GCCATCG								
L <sub>w</sub> (I) AT: 3 CA: 2 CG: 5	L <sub>w</sub> (J) AT: 3, 9 CA: 2, 8 CC: 1, 7 CG: 5, 11	Common words 2 2 0 2							
GC: 1, 7 GG: 6 TC: 4	GC: 6 TC: 4, 10	2 0 2 10 in total							
235									

# Properties of the common word list

- p Exact matches can be found using binary search (e.g., where TCGT occurs in I?)
   n O(log 4<sup>k</sup>) time
- P For large k, the table size is too large to compute the common word count in the previous fashion
- Instead, an approach based on merge sort can be utilised (details skipped)
- P The common k-word technique can be combined with the local alignment algorithm to yield a rapid alignment approach

# FASTA outline

- p FASTA algorithm has five steps:
  - n 1. Identify common k-words between I and J n 2. Score diagonals with k-word matches,
  - identify 10 best diagonals
  - n 3. Rescore initial regions with a substitution score matrix
  - n 4. Join initial regions using gaps, penalise for gaps
  - n 5. Perform dynamic programming to find final alignments





















# FASTA outline

- p FASTA algorithm has five steps:
  - n 1. Identify common k-words between I and J
  - n 2. Score diagonals with k-word matches, identify 10 best diagonals
  - n 3. Rescore initial regions with a substitution score matrix
  - n 4. Join initial regions using gaps, penalise for gaps
  - n 5. Perform dynamic programming to find final alignments

### Rescoring initial regions

- Each high-scoring diagonal chosen in the previous step is rescored according to a score matrix
- ${\rm p}\,$  This is done to find subregions with identities shorter than k
- p Non-matching ends of the diagonal are trimmed

# $\begin{array}{c} \text{I: } C \ C \ A \ T \ C \ G \ C \ A \ T \ C \ G \\ \text{J: } C \ C \ A \ \textbf{A} \ C \ G \ C \ \textbf{A} \ T \ C \ \textbf{A} \\ \text{I': } C \ C \ A \ \textbf{A} \ C \ G \ C \ \textbf{A} \ T \ C \ \textbf{A} \\ \text{J': } C \ \textbf{C} \ \textbf{A} \ \textbf{T} \ \textbf{C} \\ \text{J': } \ A \ \textbf{C} \ \textbf{A} \ \textbf{T} \ \textbf{C} \ \textbf{A} \\ \text{A} \ A \ T \ \textbf{A} \ \textbf{A} \\ \text{A} \ \textbf{A} \ \textbf{A} \ \textbf{T} \\ \text{A} \ \textbf{A} \\ \text{A} \ \textbf{A} \ \textbf{T} \\ \text{A} \ \textbf{A} \\ \text{A} \ \textbf{A} \\ \text{T} \\ \text{A} \ \textbf{A} \\ \text{A} \ \textbf{A} \\ \text{T} \\ \text{A} \\ \text{C} \ \textbf{A} \\ \text{T} \\ \text{C} \\ \text{A} \ \textbf{A} \\ \text{T} \\ \text{C} \\ \text{A} \ \textbf{A} \\ \text{T} \\ \text{C} \\ \text{A} \ \textbf{A} \\ \text{T} \\ \text{C} \\ \text{C} \\ \text{C} \ \textbf{A} \\ \text{T} \\ \text{C} \\ \text{C} \\ \text{C} \ \textbf{A} \\ \text{C} \ \text{C} \$



# Local alignment in the highest-scoring region

- Last step of FASTA: perform local alignment using dynamic programming around the highestscoring
- Region to be aligned covers -w and +w offset diagonal to the highestscoring diagonals
- With long sequences, this region is typically very small compared to the whole n x m matrix





# **Properties of FASTA**

- p FASTA looks for initial exact matches to query sequence
  - n Two proteins can have very different amino acid sequences and still be biologically similar
  - n This may lead into a lack of sensitivity with diverged sequences

# Demonstration of FASTA at EBI

- p http://www.ebi.ac.uk/fasta/
- P Note that parameter ktup in the software corresponds to parameter k in lectures