

Multiple Sequence Alignment: An Introduction

	G	E	N	E	T	I	C	S
G	60	40	30	20	20	0	10	0
E	40	50	30	30	20	0	10	0
N	30	30	40	20	20	0	10	0
E	20	20	20	30	20	10	10	0
S	20	20	20	20	20	0	10	10
I	10	10	10	10	10	20	10	0
S	0	0	0	0	0	0	0	10

Lecture 1.1

1

Resources for this lecture

- Web page
<http://www.carleton.ca/~jcheetha/topics>
- Lecture Slides
- Lab
- Notes
- Papers
- Programs

Lecture 1.1

2

Alignments tell us about...

- Function or activity of a new gene/protein
- Structure or shape of a new protein
- Location or preferred location of a protein
- Stability of a gene or protein
- Origin of a gene or protein
- Origin or phylogeny of an organelle
- Origin or phylogeny of an organism

Lecture 1.1

3

Factoid:

*Sequence comparisons
lie at the heart of all
bioinformatics*

Lecture 1.1

4

Similarity versus Homology

- Similarity refers to the likeness or % identity between 2 sequences
- Similarity means sharing a statistically significant number of bases or amino acids
- **Similarity does not imply homology**
- Homology refers to shared ancestry
- Two sequences are homologous if they are derived from a common ancestral sequence
- **Homology usually implies similarity**

Lecture 1.1

5

Similarity versus Homology

- **Similarity can be quantified**
- It is correct to say that two sequences are X% identical
- It is correct to say that two sequences have a similarity score of Z
- It is generally **incorrect** to say that two sequences are X% *similar*

Lecture 1.1

6

Similarity versus Homology

- Homology cannot be quantified
- If two sequences have a high % identity it is OK to say they are homologous
- It is **incorrect** to say two sequences have a homology score of Z
- It is **incorrect** to say two sequences are X% homologous

Lecture 1.1

7

Sequence Complexity

MCDEFGHIKLAN.... High Complexity

ACTGTCACTGAT.... Mid Complexity

NNNNTTTTTNNN.... Low Complexity

Translate those DNA sequences!!!

Lecture 1.1

8

Assessing Sequence Similarity

THESTORYOFGENESIS
THISBOOKONGENETICS **Two Character Strings**

THESTORYOFGENESI-S
THISBOOKONGENETICS **Character Comparison**

THE STORY OF GENESIS
THIS BOOK ON GENETICS **Context Comparison**

Lecture 1.1

9

Is This Alignment Significant?

Gelsolin	89	G	N	E	L	S	D	E	S	G	A	A	A	I	F	T	V	Q	L	108		
Annexin	82	P	S	A	L	K	S	A	L	S	G	H	L	E	T	V	I	L	G	L	101	
	154	E	K	D	I	I	S	D	T	S	G	D	F	R	K	L	M	V	A	L	173	
	240	E	-	S	I	K	E	V	K	G	D	L	E	N	A	P	L	N	L	258		
	314	E	Y	Y	I	Q	D	T	K	G	D	Y	Q	K	A	L	L	Y	L	333		
Consensus		L	x	P	x	x	x	P	D	x	S	G	x	h	x	x	h	x	V	L	L	

Lecture 1.1

10

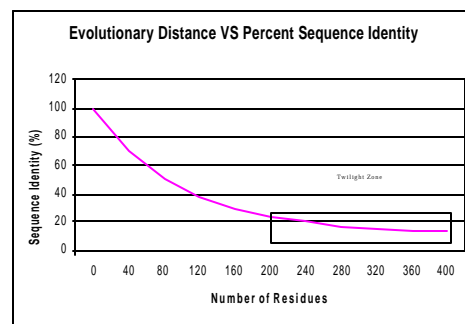
Some Simple Rules

- If two sequence are > 100 residues and > 25% identical, they are likely related
- If two sequences are 15-25% identical they **may** be related, but more tests are needed
- If two sequences are < 15% identical they are probably not related
- If you need more than 1 gap for every 20 residues the alignment is suspicious

Lecture 1.1

11

Doolittle's Rules of Thumb



Lecture 1.1

12

Multiply Aligned Proteins

- Ideally, the amino acids in each column of the alignment occupy a similar 3D structural position. All aa's in that position descend from a common ancestral amino acid, i.e. are homologous.
- But, protein sequences and structures diverge, so it is not always possible to know if an alignment is "correct," i.e. represents homology.
- See the structural example next.

Lecture 1.1

19

Superposition of Polypeptide Backbones of Aligned Proteins



Lecture 1.1

20

Why Do Multiple Alignments?

- To represent sequence similarities and differences within a group of related sequences.
- To show evolutionary pattern of nature's successful experiments in generating genetic diversity.
- To indicate which residues can be changed without destroying adaptiveness.
- For DNA: To find regions for PCR primers.
- For proteins: To predict regions of conserved 2D and 3D structures.
- First step in molecular phylogenetics analysis.

Lecture 1.1

21

Molecular Phylogenetics



Lecture 1.1

22

Purposes of Multiple Alignments of Proteins

- Discover conserved regions.
- Uncover patterns of α -helix, β -sheet.
- Uncover patterns of hydrophobicity and hydrophilicity.
- Find "gappy" regions of surface loops, hypervariability.
- To search for family members and remote homologs in databases.
- Do phylogenetic analysis, tree drawing.

Lecture 1.1

23

Databases of Protein Alignments

- Pfam
- Prosite
- Prints
- Blocks

ID	p99.1.4782; BLOCK
AC	BP04782A; distance from previous block=(0,11)
DE	PROTEIN SEX-DETERMINING REGION Y TESTIS-DETERMI
BL	DPA; width=17; seqs=8; 99.5%=837; strength=1238
SRY_CALJA P51501	(8) MLRVFNSDEYNPAALQN 100
SRY_BISBO Q27949	(1) MFRVLNDDVYSPAVVQQ 55
SRY_BOVIN Q03255	(1) MFRVLNDDVYSPAVVQQ 55
SRY_CAPHI Q03256	(12) MFRVLKDDVYSPAVVQQ 57
SRY_PIG P36393	(9) MFRVLKDDVYSPAVVQQ 73
SRY_SHEEP Q03257	(12) MFRVLKDDVYSPAVVQQ 57
SRY_GORGO P48046	(8) MLSVFNDDVYSPAVVQQ 77
SRY_HUMAN Q05066	(8) MLSVFNDDVYSPAVVQQ 85



Lecture 1.1

24

Dynamic Programming

	G	E	N	E	T	I	C	S		G	E	N	E	T	I	C	S
G	10	0	0	0	0	0	0	0	G	6	40	30	20	20	0	10	0
E	0	10	0	10	0	0	0	0	E	40	6	30	30	20	0	10	0
N	0	0	10	0	0	0	0	0	N	30	30	6	20	20	0	10	0
E	0	0	0	10	0	10	0	0	E	20	20	20	6	20	10	10	0
S	0	0	0	0	0	0	10	0	S	20	20	20	20	6	0	10	10
I	0	0	0	0	0	0	10	0	I	10	10	10	10	10	6	0	0
S	0	0	0	0	0	0	0	10	S	0	0	0	0	0	0	6	0

G E N E T I C S
 | | | | * | |
 G E N E S I S

Lecture 1.1

25

Dynamic Programming

- Developed by Needleman & Wunsch (1970)
- Refined by Smith & Waterman (1981)
- Ideal for quantitative assessment
- Guaranteed to be mathematically optimal
- Slow N^2 algorithm
- Performed in 2 stages
 - Prepare a scoring matrix using recursive function
 - Scan matrix diagonally using traceback protocol

Lecture 1.1

26

The Recursive Function

$$S_{ij} = S_{ij} + \max \begin{cases} S_{i-1,j-1} & \text{or} \\ \max_{2 < x < i} S_{i-x,j-1} + w_{x-1} & \text{or} \\ \max_{2 < y < i} S_{i-1,j-y} + w_{y-1} \end{cases}$$



W = gap penalty
S = alignment score

27

Identity Scoring Matrix (S_{ij})

	A	R	N	S	C	G	E	G	H	L	K	M	F	P	S	T	W	Y	V
A	1																		
R	0	1																	
N	0	0	1																
S	0	0	0	1															
C	0	0	0	0	1														
G	0	0	0	0	0	1													
H	0	0	0	0	0	0	1												
L	0	0	0	0	0	0	0	1											
K	0	0	0	0	0	0	0	0	1										
M	0	0	0	0	0	0	0	0	0	1									
F	0	0	0	0	0	0	0	0	0	0	1								
P	0	0	0	0	0	0	0	0	0	0	0	1							
S	0	0	0	0	0	0	0	0	0	0	0	0	1						
T	0	0	0	0	0	0	0	0	0	0	0	0	0	1					
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1				
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1			
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1		

Lecture 1.1

28

A Simple Example...

AATVD	AATVD	AATVD
A 1	A 1 1	A 1 1 0 0 0
V	V	V
V	V	V
D	D	D

AATVD	AATVD	AATVD
A 1 1 0 0 0	A 1 1 0 0 0	A 1 1 0 0 0
V 0	V 0 1 1	V 0 1 1 2
V	V	V
D	D	D

Lecture 1.1

29

A Simple Example...

AATVD	AATVD	AATVD
A 1 1 0 0 0	A 1 1 0 0 0	A 1 1 0 0 0
V 0 1 1 2 1	V 0 1 1 2 1	V 0 1 1 2 1
V	V 0 1 1 2 2	V 0 1 1 2 2
D	D 0 1 1 1 3	D 0 1 1 1 3

AATVD	AATVD	AATVD
AV - VD	A - VVD	AVVD

Lecture 1.1

30

Could We Do Better?

- Key to the performance of Dynamic Programming is the scoring function
- Dynamic Programming always gives the mathematically correct answer
- Dynamic Programming does not always give the biologically correct answer
- The weakest link -- The Scoring Matrix

Lecture 1.1

31

Scoring Matrices

- An empirical model of evolution, biology and chemistry all wrapped up in a 20 X 20 table of integers
- Structurally or chemically similar residues should ideally have high diagonal or off-diagonal numbers
- Structurally or chemically dissimilar residues should ideally have low diagonal or off-diagonal numbers

Lecture 1.1

32

A Better Matrix - PAM250

	A	R	N	D	C	G	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	2																			
R	-2	6																		
N	0	0	2																	
D	0	-1	2	4																
C	-2	-4	-4	-6	4															
G	0	1	1	2	3	4														
E	0	-1	1	3	3	2	4													
G	1	-3	0	1	3	3	1	0	3											
H	-1	2	2	1	3	3	1	2	0	3										
I	-1	2	2	2	3	2	2	2	2	3										
L	-2	-3	-3	-4	-4	0	2	3	3	2	3									
K	-1	3	1	1	0	3	1	0	2	0	-2	-3	5							
M	-1	0	2	3	3	1	2	3	2	2	2	0	0	6						
F	2	3	2	3	3	3	3	2	2	2	2	3	3	0	0					
P	1	0	1	1	0	1	0	1	1	-1	-2	0	-2	-3	1	3				
S	1	-1	1	0	0	2	1	0	0	0	0	-1	-2	0	1	3				
T	-1	-1	0	0	0	1	0	0	0	0	-3	0	-1	-2	0	1	3			
W	-6	2	4	7	8	9	7	7	3	-5	-2	-3	-4	0	6	-2	-5	17		
Y	-3	-4	-2	-4	-4	-4	-4	-4	-3	0	-1	-1	-4	-2	7	-5	-3	0	10	
V	0	-2	-2	-2	-2	-2	-2	-2	-2	-2	-2	-1	-1	-1	0	-6	-2	-4		

Lecture 1.1

33

PAM Matrices

- Developed by M.O. Dayhoff (1978)
- PAM = Point Accepted Mutation
- Matrix assembled by looking at patterns of substitutions in closely related proteins
- 1 PAM corresponds to 1 amino acid change per 100 residues
- 1 PAM = 1% divergence or 1 million years in evolutionary history

Lecture 1.1

34

Using PAM250...

A T V D
A 2
T 1 3
V 0 0 4
D 0 0 -2 4

Gap
Penalty = -1

A A T V D	A A T V D	A A T V D
A 2	A 2 2	A 2 1 1 0 0
V	V	V
V	V	V
D	D	D
A A T V D	A A T V D	A A T V D
A 2 1 1 0 0	A 2 1 1 0 0	A 2 1 1 0 0
V 0 2	V 0 2 1	V 0 2 1 5
V	V	V
D	D	D

Lecture 1.1

35

Using PAM250...

A T V D
A 2
T 1 3
V 0 0 4
D 0 0 -2 4

Gap
Penalty = -1

A A T V D	A A T V D	A A T V D
A 2 1 1 0 0	A 2 1 1 0 0	A 2 1 1 0 0
V 0 2 1 5 -1	V 0 2 1 5 -1	V 0 2 1 5 -1
V	V 0 1 2 5 3	V 0 1 2 5 3
D	D 0 1 1 0 9	D 0 1 1 0 9
A A T V D		
A V - V D		

Lecture 1.1

36

Alignment Methods

- By hand: requires expert knowledge of protein evolution and structure; also tedious.
- Automatic: done by a computer program; requires a scoring method that can be used to measure how good a particular alignment is. ClustalW is this type of method.
- Refinement: after automatic alignment, improve with advanced methods (software) or by hand (requires an editor).



Lecture 1.1

43

Methods: Multidimensional Dynamic Programming

- Guarantees an optimal solution.
- Memory and computational time required increase exponentially with the product of the sequence lengths.
- **MSA** algorithm is the foremost implementation; can align 5-7 proteins of 200-300 residues. But servers have limitations:
 - Runs > 10 minutes are terminated.
 - If > 4 sequences, limit lengths to 100 aas.

Lecture 1.1

44

Methods: Progressive Alignment

- General method is to construct a succession of pairwise alignments.
- Most common approach
- Fast and efficient due to use of heuristics
- **Feng-Doolittle** method: construct a "guide tree" of similarities; pairwise align most similar sequences first.
- **ClustalW**: adds more sophisticated methods, but still works by adding one sequence at a time

Lecture 1.1

45

Clustal Algorithm



1. Compare n sequences in pairs. Do all possible pairs $[n \times (n-1) / 2]$.
2. Construct a similarity matrix $[n \times n]$ of pairwise alignments.
3. Rank sequences according to how well they align pairwise.
4. Do progressive alignment: Start with best pair of sequences and continue aligning, until all sequences are aligned.

Lecture 1.1

46

Overview of ClustalW Procedure



Lecture 1.1

47

Clustal Algorithm II

4. Do progressive alignment:
 - 1) start with best pair of sequences
 - 2) align next most similar sequence to first pair
 - 3) continue aligning next most similar, until all sequences are aligned.
5. No iteration is done to refine the final alignment.
6. ClustalW may not give best possible MSA. Usually need to refine MSA by hand.



Lecture 1.1

48

ClustalW Features



- Global alignment
- Progressive alignment
- Fast, heuristic, but not guaranteed to find the best alignment
 - Should adjust MSA by hand
- Good for similar length sequences
- Good for protein families

Lecture 1.1

49

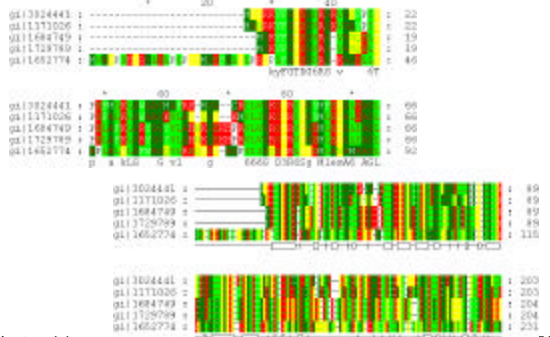
Alignment Editors

- Programs that allow the user to manipulate a multiple sequence alignment that has been produced by software.
- Sections of sequences can be moved to improve the alignment (possibly).
- Most importantly, sections of the MSA can be shaded and highlighted to show which aas are the most similar (conserved), and what the chemical properties are.
- The GeneDoc and JalView programs are examples.

Lecture 1.1

50

GeneDoc Screens

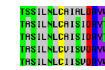


Lecture 1.1

51

Progressive Alignments

- True for ClustalW and all progressive alignment methods: heuristic, not exhaustive, may not find "best" MSA.
- Mathematical process, not connected with biological reality.
- Output MSA can usually be improved by hand. Inspect MSA closely, adjust.
- Use viewer with color-coding.
- Make judgments based on biology.
- Emphasize most reliable regions of MSA, based on convincing positional homology.



Lecture 1.1

52

Various Methods of Multiple Alignment

- MSA program: if few, short sequences
- ClustalW: usually best
- ClustalX: graphical
- PILEUP: in GCG package
- DIALIGN: local gap-free alignments first
- Match-Box: good for conserved domains
- SAGA: "genetic algorithm"
- MEME: iterative realignment



Lecture 1.1

53

Test of Alignment Method

1. Align related protein sequences that have known 3D structures.
2. Match structures, identify common domains and structural motifs.
3. Check if sequence alignment has correctly aligned the common motifs.
4. Evaluate alignment methods by comparing with each other.

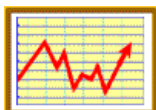


Lecture 1.1

54

Still Improvements to Seek

- Better algorithms: faster, more efficient
- Better scoring: reasonable biologically
- Able to deal with repeats
- Able to deal with domain swaps
- Incorporate external information
 - eg. active sites, critical residues



Lecture 1.1

55

Web Sites



- ClustalW servers:
 - <http://www.clustalw.genome.ad.jp/>;
 - <http://dot.imgen.bcm.tmc.edu:9331/multi-align/multi-align.html>;
 - <http://www.genebee.msu.su/genebee.html>;
 - <http://www.bionavigator.com>;also others
- multiple sequence analysis tools:
 - <http://www.techfak.uni-bielefeld.de/bcd/Curric/MulAli/welcome.html>

Lecture 1.1

56

More Web Sites



- Jalview: <http://www.ebi.ac.uk/~michele/jalview/>
- BoxShade viewer:
http://www.ch.embnnet.org/software/BOX_form.html
- GeneDoc viewer:
<http://www.psc.edu/biomed/genedoc/>
- Phylip phylogenetic analysis:
<http://evolution.genetics.washington.edu/phylip.html>

Lecture 1.1

57

The problem: Sequence comparison

- How to compare one sequence (target) to many sequences (database search)
- How to compare more than two sequences simultaneously

Lecture 1.1

58

What is a Multiple Sequence Alignment (MSA) ?

- MSA is the alignment of N sequences (Protein or DNA) simultaneously, where $N > 2$.
- Let S_i denote a sequence in the Global Multiple Sequence Alignment of $N > 2$ sequences $S = \{S_1, \dots, S_N\}$ is obtained by inserting gaps denoted by “ - ” possibly at the beginning or end, positions.
- The new set of N sequences denoted by $S' = \{S'_1, \dots, S'_N\}$ will all have length L .

Lecture 1.1

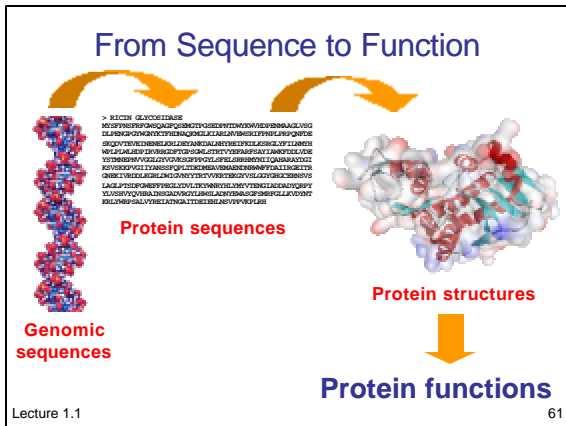
59

Reasons for aligning sets of sequences

- Infer *phylogenetic* trees from *homologous* sequences
- Highlight conserved sites/regions
 - Residues conserved during evolution play an important role
- Highlight variable sites/regions
- Uncover changes in gene/protein structure
- Prediction of protein structure and function
 - Proteins which are very similar in sequence generally have similar 3D structure and function as well (*Homology modeling*)
 - By searching a sequence of unknown structure against a database of known proteins the structure and/or function can in many cases be predicted

Lecture 1.1

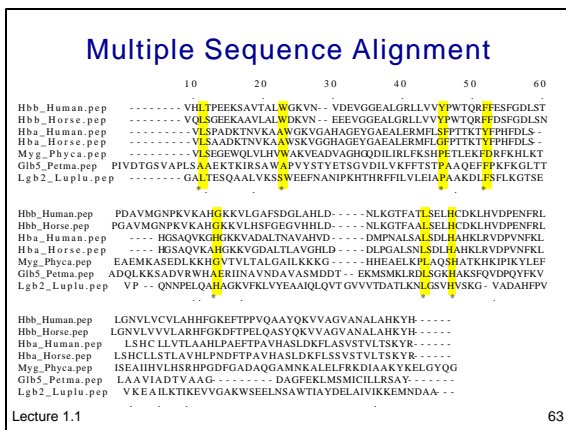
60



Global vs. Local Alignments

- **Global alignment** algorithms start at the beginning of two sequences and add gaps to each until the end of one is reached.
- **Local alignment** algorithms finds the region (or regions) of highest similarity between two sequences and build the alignment outward from there.

Lecture 1.1



The challenge: pairwise to multiple alignment

- in principle pairwise approach applicable to multiple alignment - but not practical
- CPU time proportional to sequence length
e.g. if aligning 2 sequences of 300 positions = 1 second
... 3 sequences = 300 seconds
and 10 sequences = 300⁸ seconds
= 951 years

Lecture 1.1 64

Progressive Multiple Alignment (Clustal)

- Any Exact Method would be **TOO SLOW**
- We use a **Heuristic Algorithm**.
- **Progressive Alignment Algorithms** are the most Popular
 - -Fast
 - -Clustal
 - -Greedy Heuristic (No Guarranty).

Lecture 1.1 65

Progressive alignment

Scerevisiae	[1]				
C.elegans	[2]	0.640			
Drosophila	[3]	0.634	0.327		
Human	[4]	0.630	0.408	0.420	
Mouse	[5]	0.619	0.405	0.469	0.289

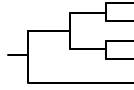
1. Do pairwise alignment of all sequences and calculate distance matrix.
2. Create a guide tree based on this pairwise distance matrix.
3. Align progressively following guide tree.
 - start by aligning most closely related pairs of sequences
 - (keep gaps that appeared in sequence pairs fixed)
 - at each step align two sequences or one to an existing subalignment

Lecture 1.1 66

Parallel Clustal

•**Pairwise (PW) alignment matrix**
average alignment calculation spends most of its time here
can parallelize as all elements are independent

•**Guide tree calculation**
Calculation of closest sequences (branch) can be parallelized. Together with PW matrix calculation Clustal W is ~85-92% parallel



•**Progressive alignment**
Remaining ~5-10% of the code can be parallelized at this stage by calculating profile scores in parallel. As a result the whole application can be ~93-98% parallel depending on a size of a problem



Lecture 1.1

67



Lecture 1.1

68

Parallel Clustal: extensions

Optimization of input parameters - scoring matrices, gap penalties - requires many repetitive Clustal W calculations with various input parameters.

Minimum Vertex Cover – use k-minimum vertex cover to remove erroneous sequences, and identify clusters of highly similar sequences.

Lecture 1.1

69

Minimum Vertex Cover

Minimum vertex cover (classic problem)

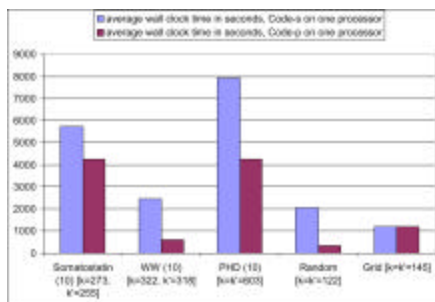
Definition: A set of vertices in an undirected graph where every edge connects at least one vertex. The minimum vertex cover problem is to find a minimum size set and is NP-complete.

Cheetham, J.J., Dehne, F., Rau-Chaplin, A., Stege, U. and Taillon, P.J. (2003) **Solving Large FPT Problems on Course Grained Parallel Machines**, *J. Computer and Systems Science* (in press).

Lecture 1.1

70

Parallel FPT Minimum Vertex Cover



Lecture 1.1

71

Clustal XP

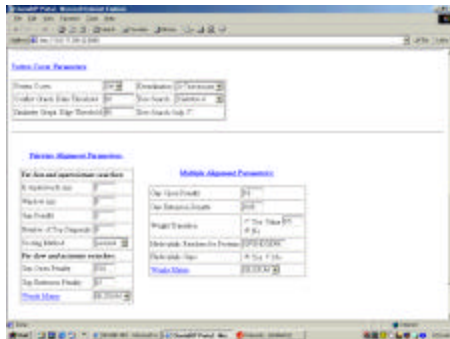


<http://134.117.206.42:8000/>

Lecture 1.1

72

Clustal XP



Lecture 1.1

73

Conclusions and Summary

- Parallel computing can greatly speed up multiple sequence alignments
- Minimum vertex cover can be very useful in multiple sequence alignment to detect misplaced sequences and detect groups of very similar sequences
- A fixed parameter tractable algorithm can be used to solve minimum k-vertex cover problems for large values of k.

Lecture 1.1

74