Dynamic Programming (Local & Global Alignment)

Global alignment the gap

1	AGGATTGGAATGCTCAGAA	GCAGCTAAAGCGI	GTATGCAGGATT	GGAATTAAAGA	GGAGGTAG	ACCG 67
	111111111111111					
1	AGGATTGGAATGCTACAGA	AGCAGCTAAAGCG	GTGTATGCAGGAT	TGGAATTAAAG	GAGGAGGTA	GACCG 68

The alignment is much better when one gap is introduced

1	AGGATTGGAATGCT.CAGAAGCAGCTAAAGCGTGTATGCAGGATTGGAATTAAAGAGGAGGTAGACCG	67
1	AGGATTGGAATGCTACAGAAGCAGCTAAAGCGTGTATGCAGGATTGGAATTAAAGAGGAGGTAGACCG	68

Parameters for sequence alignment

Gap penalties

Opening:The cost to introduce a gapExtension:The cost to extend a gap

Scoring systems

Every symbol pairing is assigned with a numerical value that is based on a "symbol comparison" *or* "replacement" table/matrix

Why gap penalties ?

The optimal alignment of two similar sequences usually

- maximizes the number of matches and
- **minimizes** the number of gaps.

Permitting the insertion of arbitrarily many gaps might lead to high scoring alignments of non-homologous sequences.

Penalizing gaps forces alignments to have relatively few gaps.

Gap penalties increase the quality of an alignment – non-homologous sequences are not aligned

Gap penalties

Linear gap penalty score:

$$\gamma(g) = -gd$$

Affine gap penalty score:

$$\gamma(g) = -d - (g - 1) e$$

- γ (g) = gap penalty score of a gap of length g
 - d = gap opening penalty
 - e = gap extension penalty
 - g = gap length

Scoring insertions and deletions



Gap parameters: d = 3 (gap opening) e = 0.1 (gap extension) g = 3 (gap length) $\gamma(g) = -3 - (3 - 1) 0.1 = -3.2$

match	= 1
mismatch	= 0

Calculating alignments: Global vs. Local alignment

• For optimal GLOBAL alignment, we want best score in the final row or final column

GLOBAL - best alignment of entirety of both sequences (possibly at expense of great local similarity)

• For optimal LOCAL alignment, we want best score anywhere in matrix

LOCAL - best alignment of segments, without regard to rest of two sequences (at the expense of the overall score)

Important Points in Pairwise Sequence Alignment

Significance of Similarity

- Dependent on PID (Percent Identical Positions in Alignment)
- -Similarity/Disimilarity score
- Significance of score depend on length of alignment
- -Significance Score (Z) whether score significant
- -Expected Value (E), Chances that non-related sequence may have that score



Why we do multiple alignments?

- Multiple nucleotide or amino sequence alignment techniques are usually performed to fit one of the following scopes :
- In order to characterize protein families, identify shared regions of homology in a multiple sequence alignment; (this happens generally when a sequence search revealed homologies to several sequences)
- Determination of the consensus sequence of several aligned sequences.
- Help prediction of the secondary and tertiary structures of new sequences;
- Preliminary step in molecular evolution analysis using Phylogenetic methods for constructing phylogenetic trees

An example of Multiple Alignment VTIS**C**TGSSSNIGAG-NHVK**W**YQ**QLPG** VTISCTGTSSNIGS--ITVNWYQQLPG LRLSCSSGFIFSS--YAMYWVRQAPG LSLTCTVSGTSFDD--YYSTWVRQPPG PEVTCVVVDVSHEDPOVKFNWYVDG--ATLVCLISDFYPGA--VTVAWKADS--AALGCLVKDYFPEP--VTVSWNSG---VSLTCLVKGFYPSD--IAVEWWSNG--

Alignment of Multiple Sequences

Extending Dynamic Programming to more sequences

- -Dynamic programming can be extended for more than two
- -In practice it requires CPU and Memory (Murata et al 1985)
- MSA, Limited only up to 8-10 sequences (1989)
- -DCA (Divide and Conquer; Stoye et al., 1997), 20-25 sequences
- -OMA (Optimal Multiple Alignment; Reinert et al., 2000)
- -COSA (Althaus et al., 2002)

Progressive or Tree or Hierarchical Methods (CLUSTAL-W)

- -Practical approach for multiple alignment
- -Compare all sequences pair wise
- -Perform cluster analysis
- -Generate a hierarchy for alignment
- -first aligning the most similar pair of sequences
- -Align alignment with next similar alignment or sequence

Alignment of Multiple Sequences

Iterative Alignment Techniques

•Deterministic (Non Stochastic) methods

- -They are similar to Progressive alignment
- -Rectify the mistake in alignment by iteration
- -Iterations are performed till no further improvement
- -AMPS (Barton & Sternberg; 1987)
- -PRRP (Gotoh, 1996), Most successful
- -Praline, IterAlign
- Stochastic Methods
 - SA (Simulated Annealing; 1994), alignment is randomly modified only acceptable alignment kept for further process. Process goes until converged
 - Genetic Algorithm alternate to SA (SAGA, Notredame & Higgins, 1996)
 - -COFFEE extension of SAGA
 - -Gibbs Sampler
 - -Bayesian Based Algorithm (HMM; HMMER; SAM)

-They are only suitable for refinement not for producing *ab initio* alignment. Good for profile generation. Very slow.

Alignment of Multiple Sequences

Progress in Commonly used Techniques (Progressive)

Clustal-W (1.8) (Thompson et al., 1994) Automatic substitution matrix Automatic gap penalty adjustment Delaying of distantly related sequences Portability and interface excellent T-COFFEE (Notredame et al., 2000) Improvement in Clustal-W by iteration Pair-Wise alignment (Global + Local) Most accurate method but slow MAFFT (Katoh et al., 2002) Utilize the FFT for pair-wise alignment Fastest method Accuracy nearly equal to T-COFFEE



(A) Pairwise Alignment



(B) Multiple alignment following the tree from A.



Multiple Alignment Method

- The steps are summarized as follows:
- Compare all sequences pairwise.
- Perform cluster analysis on the pairwise data
- Generate a hierarchy for alignment
 - Binary tree or a simple ordering
- First align the most similar pair of sequences
- Then the next most similar pair and so on.
- Once an alignment of two sequences has been made, then this is fixed.
- Thus for a set of sequences A, B, C, D having aligned
- A with C and B with D
- Alignment of A, B, C, D is obtained by comparing the alignments of A and C with that of B and D

- using averaged scores at each aligned position.

ClustalW- for multiple alignment

- ClustaW is a multiple alignment program for DNA or proteins.
- Developed by Julie D. Thompson, Toby Gibson at EMBL/EBI
- ClustalW: Improving the sensitivity of multiple sequence alignment
 - sequence weighting
 - positions-specific gap penalties
 - weight matrix choice
 - Nucleic Acids Research, 22:4673-4680
- Manipulate existing alignments
- do profile analysis
- create phylogentic trees.
- Alignment can be done by 2 methods:
 - slow/accurate
 - fast/approximate

Running ClustalW

[~]% clustalw

- 1. Sequence Input From Disc
- 2. Multiple Alignments
- 3. Profile / Structure Alignments
- 4. Phylogenetic trees
- S. Execute a system command
- H. HELP
- X. EXIT (leave program)

Your choice:

Using ClustalW

***** MULTIPLE ALIGNMENT MENU *****

- 1. Do complete multiple alignment now (Slow/Accurate)
- 2. Produce guide tree file only
- 3. Do alignment using old guide tree file
- 4. Toggle Slow/Fast pairwise alignments = SLOW
- 5. Pairwise alignment parameters
- 6. Multiple alignment parameters
- 7. Reset gaps between alignments? = OFF
- 8. Toggle screen display = ON
- 9. Output format options
- S. Execute a system command
- H. HELP

or press [RETURN] to go back to main menu

Your choice:

Output of ClustalW

CLUSTAL W (1.7) multiple sequence alignment

HSTNFR	GGGAAGAGTTCCCCAGGGACCTCTCTCTAATCAGCCCTCTGGCCCAG(
SYNTNFTRP	GGGAAGAGTTCCCCAGGGACCTCTCTCTAATCAGCCCTCTGGCCCAG(
CFTNFA	TGTCCAGi
CATTNFAA	GGGAAGAGCTCCCACATGGCCTGCAACTAATCAACCCTCTGCCCCAGi
RABTNFM	AGGAGGAAGAGTCCCCAAACAACCTCCATCTAGTCAACCCTGTGGCCCAGATGGTC
RNTNFAA	AGGAGGAGAAGTTCCCAAATGGGCTCCCTCTCATCAGTTCCATGGCCCAGACCCTC
OATNFA1	GGGAAGAGCAGTCCCCAGCTGGCCCCTCCTTCAACAGGCCTCTGGTTCAG
OATNFAR	GGGAAGAGCAGTCCCCAGCTGGCCCCTCCTTCAACAGGCCTCTGGTTCAGi
BSPTNFA	GGGAAGAGCAGTCCCCAGGTGGCCCCTCCATCAACAGCCCTCTGGTTCAA
CEU14683	GGGAAGAGCAATCCCCAACTGGCCTCTCCATCAACAGCCCTCTGGTTCAGi
	* *

ClustalW options

Your choice: 5

******** PAIRWISE ALIGNMENT PARAMETERS ******** Slow/Accurate alignments:

Gap Open Penalty :15.00
 Gap Extension Penalty :6.66
 Protein weight matrix :BLOSUM30
 DNA weight matrix :IUB

Fast/Approximate alignments:

5. Gap penalty :5
6. K-tuple (word) size :2
7. No. of top diagonals :4
8. Window size :4

9. Toggle Slow/Fast pairwise alignments = SLOW

H. HELP Enter number (or [RETURN] to exit):

ClustalW options

Your choice: 6

******** MULTIPLE ALIGNMENT PARAMETERS ********

- 1. Gap Opening Penalty:15.002. Gap Extension Penalty:6.66
- 3. Delay divergent sequences :40 %
- 4. DNA Transitions Weight :0.50
- 5. Protein weight matrix :BLOSUM series
- 6. DNA weight matrix :IUB
- 7. Use negative matrix :OFF
- 8. Protein Gap Parameters

H. HELP

Enter number (or [RETURN] to exit):

ClustalX - Multiple Sequence Alignment Program

- ClustalX provides a new window-based user interface to the ClustalW program.
- It uses the Vibrant multi-platform user interface development library, developed by the National Center for Biotechnology Information (Bldg 38A, NIH 8600 Rockville Pike,Bethesda, MD 20894) as part of their NCBI SOFTWARE DEVELOPEMENT TOOLKIT.

Multiple Ali	nment Mode 🗆	Font Size: 10 -	
2			
1			

Load Sequences	Faut Circu to	-	
Load Profile 1	Font Size: 10		
Load Profile 2			
Append Sequences			
Save Sequences as			
Save Profile 1 as			
Save Profile 2 as			
Write Sequences to PS			
Write Profile 1 to PS			
Write Profile 2 to PS			
Quit			
Π			

irectories Files irectories Files outputs/ coldna.par outputs/ colprint.par outputs/ameba colprot.par outputs/temp databanks.txt outputs/temp eh_galectin.pep eh_galectin.xnu eh_galectin.xnu	irectories Files irectories Files outputs/ coldna.par outputs/ colprint.par outputs/ameba colprot.par outputs/temp databanks.txt outputs/temp eh_galectin.pep eh_galectin.xnu eh_galectin.xnu election users1/people/bfbecker/outputs/	/users1/people/bf	becker/outputs/¥
/outputs/ coldna.par /outputs/ colprint.par /outputs/ameba colprot.par /outputs/others databanks.txt /outputs/temp dotplot eh_galectin.pep eh_galectin.seg eh_galectin.xnu eh_galectin.pep	/outputs/ coldna.par /outputs/ colprint.par /outputs/ameba colprot.par /outputs/temp databanks.txt /outputs/temp eh_galectin.pep eh_galectin.xnu eh_galectin.xnu election ////////////////////////////////////	irectories	Files
election	election /users1/people/bfbecker/outputs/	'outputs/. 'outputs/ 'outputs/ameba 'outputs/others 'outputs/temp	coldna.par colprint.par colprot.par databanks.txt dotplot eh_galectin.pep eh_galectin.seg eh_galectin.xnu
7 47 1 2 2 1 2 2 2 2	/users1/people/bfbecker/outputs/	election	





