

Multiple sequence alignments

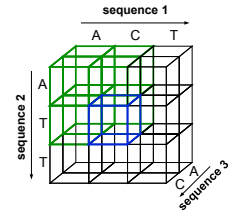


Jacques van Helden
 Jacques.van-Helden@univ-amu.fr
 Aix-Marseille Université (AMU), France
 Lab. Technological Advances for Genomics and Clinics
 (TAGC, INSERM Unit U1090)
<http://tagc.univ-mrs.fr>
 FORMER ADDRESS (1999-2011)
 Université Libre de Bruxelles, Belgique
 Bioinformatique des Génomes et des Réseaux (BIGRe lab)
<http://www.bigre.ulb.ac.be/>



Dynamical programming - multiple alignment

- Dynamical programming can be extended to treat a set of 3 sequences
 - build a 3-dimensional matrix
 - the best score of each cell is calculated on the basis of the preceding cells in the 3 directions, and a scoring scheme (substitution matrix + gap cost)
- Can be extended to n sequences by using a n-dimensional hyper-cube
- Problem: matrix size and execution time increase exponentially with the number of sequences
 - 2 sequences: $L_1 \times L_2$
 - 3 sequences: $L_1 \times L_2 \times L_3$
 - 4 sequences: $L_1 \times L_2 \times L_3 \times L_4$
 - n sequences: $L_1 \times L_2 \times \dots \times L_n$
- Aligning n sequences with dynamical programming requires $O(L^n)$ operations, which becomes thus very rapidly impractical.
- The efficiency can be improved by only considering a subspace of the n-dimensional matrix. However, even with this kind of algorithmic improvement, the number of sequences that can be aligned is still restricted (~8 sequences maximum).

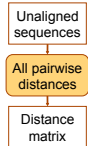


Progressive alignment

- Another approach to align multiple sequences is to perform a progressive alignment. The algorithm proceeds in several steps:
 - Calculate a **distance matrix**, representing the distance between each pair of sequences.
 - From this matrix, build a **guide tree** regrouping the closest sequences first, and the more distant sequences later.
 - Use this tree as guide to progressively align the sequences.
- This is a heuristics
 - it is a practically tractable approach, but it cannot guarantee to return the optimal solution

Progressive alignment – step 1: compute distance matrix

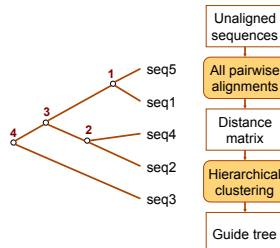
- Perform a pairwise alignment between each pair of sequences (dynamical programming or faster heuristic algorithm).
 - For n sequences: $n*(n-1)/2$ pairwise alignments.
- From each pairwise alignment, calculate the distance between the two sequences.
 - $d_{ij} = s_{ij} / L_{ij}$
 - d_{ij} : distance between sequences i and j
 - L_{ij} : length of the alignment
 - s_{ij} : number of substitutions
- Remarks
 - Gaps are not taken into account in the distance metric.
 - The matrix is symmetrical: $d_{ij} = d_{ji}$
 - Diagonal elements are null: $d_{ii} = 0$



	seq 1	seq 2	...	seq n
seq 1	d1,1	d1,2	...	d1,n
seq 2	d2,1	d2,2	...	d2,n
...
seq n	dn,1	dn,2	...	dn,n

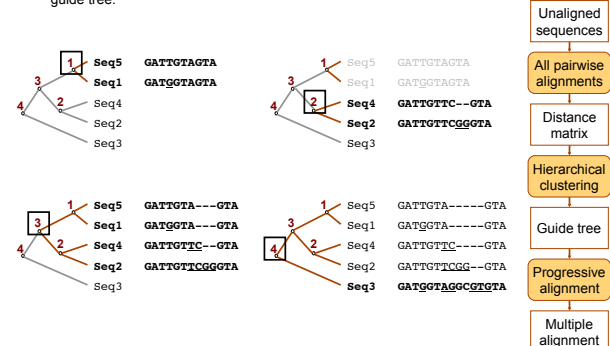
Progressive alignment – step 2: build guide tree

- A tree can be computed from the matrix distance by hierarchical clustering.
 - first regroup the two closest sequences (cluster 1)
 - next, progressively regroup the closest remaining clusters
 - the two closest sequences (cluster 2)
 - one cluster with another cluster (cluster 3)
 - one sequence with a previous cluster (cluster 4)
- This tree will then be used as **guide** to determine the order of incorporation of the sequences in the multiple alignment.
- Beware ! The guide tree should not be interpreted as a phylogenetic tree.
 - Its only purpose will be to identify the closest similarities between sequences in order to build a multiple alignment.



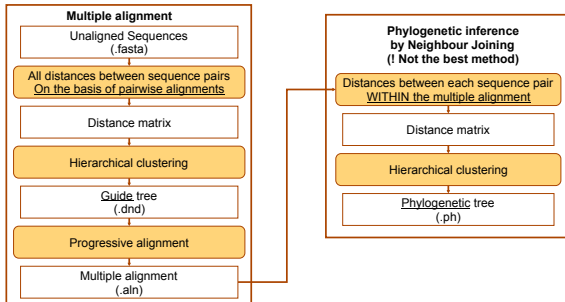
Progressive alignment – step 3: multiple alignment

- Build a multiple alignment, by progressively incorporating the sequences according to the guide tree.

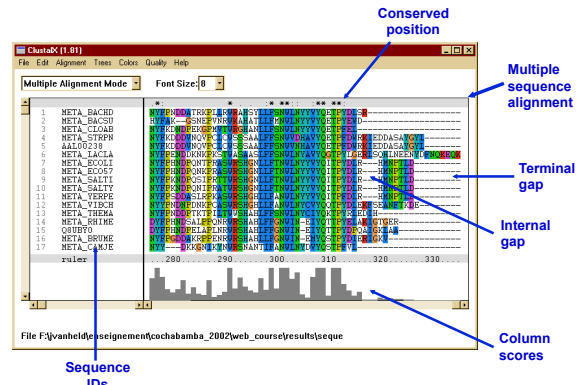


Progressive alignment and Neighbour-Joining (NJ) tree with clustalX

- Attention ! The guide tree is not a phylogenetic tree
 - Its only role is to propose an order of incorporation of the sequences for building the multiple alignment.
 - It does not aim at predicting the evolutionary history of the divergences between sequences.
- In a second time, it is possible to infer a phylogenetic tree from the multiple alignment, using the **Neighbour Joining (NJ)** method.
 - However, this method is sub-optimal for phylogenetic inference.

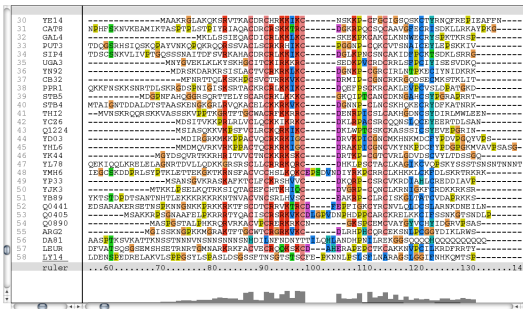


Global multiple alignment : Homoserine-O-dehydrogenase



Alignment of proteins containing a Zinc cluster domain

- The alignment of yeast Zn(2)Cys(6) binuclear cluster proteins is a difficult case.
 - The conserved region is restricted to the Zinc cluster domain.
 - This domain is not contiguous, it contains conserved and variable positions.
 - The alignment highlights 5 of the 6 characteristic cysteines.



Local multiple alignment

Progressive alignment - summary

- Processing time
 - Building the tree: proportional to $n \times n$
 - Aligning sequences: linear with number of sequences
- Heuristic method
 - cannot guarantee to return the optimal alignment.
- clustalX is a window-based environment for clustalw, which provides additional functionalities
 - Mark low scoring segments
 - The alignment can be refined manually
 - Realign selected sequences
 - Realign selected positions

Bioinformatics

Sequence motifs

Scoring a sequence with a profile matrix

Position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
A	-1.72	0.19	-1.72	-0.39	-1.72	-1.72	-1.72	-1.72	-1.72	-1.72	-1.72	-1.72	-1.72	-1.72	0.48	0.63	-1.72	0.82	0.59	-1.72
C	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26
D	-0.70	-0.70	-0.70	1.21	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70	-0.70
E	-0.20	-0.20	-0.20	1.92	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20	-0.20
F	0.42	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	1.18	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90	-0.90
G	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26	1.26
H	-1.32	-1.32	-1.32	0.60	0.98	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32	-1.32
I	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	0.21	-1.11	1.19	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	0.21
K	-1.11	-1.11	1.27	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11	-1.11
L	-1.48	-0.15	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	0.97	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	0.97
M	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	0.93	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	-0.78	1.01
N	-1.04	-1.04	-1.04	-1.04	1.11	-1.04	-1.04	-1.04	-1.04	0.74	-1.04	-1.04	-1.04	-1.04	-1.04	-1.04	-1.04	-1.04	-1.04	-1.04
P	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60	0.60
Q	1.26	1.26	0.36	0.97	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	-1.26	1.19	-1.26	-1.26	-1.26
R	-0.48	-0.48	0.85	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48	-0.48
S	-1.48	0.78	-1.48	0.14	-1.48	-1.48	-1.48	1.03	0.14	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48	-1.48
T	-1.57	0.42	-0.48	0.66	-1.57	-1.57	0.73	0.97	0.91	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57	-1.57
V	-1.52	-1.52	-0.20	-1.52	-1.52	1.01	-1.52	-1.52	-1.52	-1.52	0.63	-1.52	-1.52	-1.52	-1.52	0.45	-1.52	-1.52	-1.52	0.09
W	0.98	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	1.06	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53	-1.53
Y	-0.85	-0.85	-0.85	-0.85	1.08	-0.85	-0.85	0.77	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85	-0.85
sum	-17.8	-14.4	-13.7	-16.7	-17.2	-19.1	-16.7	-17.8	-17.6	-16.3	-14.1	-17.2	-19.1	-19.1	-17.2	-16	-17.4	-19.1	-17.2	-16.3

Sequence	A	K	D	H	V	T	S	T	M	F	V	C	W	A	V	M	A	A	L	V	SUM
Score	1.72	1.11	-0.7	1E+16	-1.52	1.97	-1.48	-1.57	-0.78	-0.9	-1.52	1.26	-1.53	-1.72	1.52	-0.78	-1.72	0.617	-1.48	0.984	21.9522

PSI-BLAST

- PSI-BLAST stands for Position-Specific Iterated BLAST (Altschul et al, 1997)
 - BLAST runs a first time in normal mode.
 - Resulting sequences are aligned together (Multiple sequence alignment) and a PSSM is calculated.
 - This PSSM is used to scan the database for new matches.
 - Steps 2-3 can be iterated several times.
- The PSSM increases the sensitivity of the search.

References

- Substitution matrices
 - PAM series
 - Dayhoff, M. O., Schwartz, R. M. & Orcutt, B. (1978). A model of evolutionary change in proteins. Atlas of Protein Sequence and Structure 5, 345-352.
 - BLOSUM substitution matrices
 - Henikoff, S. & Henikoff, J. G. (1992). Amino acid substitution matrices from protein blocks. Proc Natl Acad Sci U S A 89, 10915-9.
 - Gonnet matrices, built by an iterative procedure
 - Gonnet, G. H., Cohen, M. A. & Benner, S. A. (1992). Exhaustive matching of the entire protein sequence database. Science 256, 1443-5, 1.
- Sequence alignment algorithms
 - Needleman-Wunsch (pairwise, global)
 - Needleman, S. B. & Wunsch, C. D. (1970). A general method applicable to the search for similarities in the amino acid sequence of two proteins. J Mol Biol 48, 443-53.
 - Smith-Waterman (pairwise, local)
 - Smith, T. F. & Waterman, M. S. (1981). Identification of common molecular subsequences. J Mol Biol 147, 195-7.
 - FastA (database searches, pairwise, local)
 - W. R. Pearson and D. J. Lipman. Improved tools for biological sequence comparison. Proc. Natl. Acad. Sci. USA, 85:2444-2448, 1988.
 - BLAST (database searches, pairwise, local)
 - S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. A basic local alignment search tool. J. Mol. Biol., 215:403-410, 1990.
 - S. F. Altschul, T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller, and D. J. Lipman. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs Nucleic Acids Res., 25:3389-3402, 1997.
 - Clustal (multiple, global)
 - Higgins, D. G. & Sharp, P. M. (1988). CLUSTAL: a package for performing multiple sequence alignment on a microcomputer. Gene 73, 237-44.
 - Higgins, D. G., Thompson, J. D. & Gibson, T. J. (1996). Using CLUSTAL for multiple sequence alignments. Methods Enzymol 266, 363-402.
 - Dialign (multiple, local)
 - Morgenstern, B., Frech, K., Dress, A. & Werner, T. (1998). DIALIGN: finding local similarities by multiple sequence alignment. Bioinformatics 14, 250-4.
 - MUSCLE (multiple local)