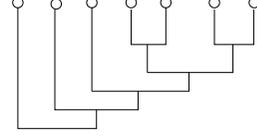
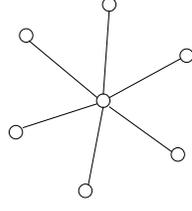


**Building a multiple alignment**  
A circle represents a sequence.

**Simultaneous**



**Progressive (tree)**



**Progressive (star)**

**Algorithms for multiple alignment**

- Simultaneous multiple alignment (example: the MSA program).  
Produces the 'true' optimal alignment.
  - Progressive multiple alignment (examples: the ClustalW and T-Coffee programs).  
Very time and space consuming: feasible only for few sequences which are highly similar, or short.
  - Manually edit a multiple alignment (very boring).  
A heuristic method that may or may not produce the true optimal alignment.  
But it is much faster than simultaneous multiple alignment.
- Example: <http://www.jalview.org/>  
Can be downloaded and installed locally.

**The MSA program**

MSA (Multiple Sequence Alignment, Kececioglu et al, 1989, 1995) implements simultaneous multiple alignment. The 1995 version of MSA has many algorithmic improvements over naïve dynamic programming. MSA fills in only a small part of the  $K$ -dimensional  $F$  matrix, a 'band' around the 'diagonal' (like Blast). The size of this 'band' is determined dynamically, by a form of branch-and-bound optimization. Therefore MSA may easily handle six closely related sequences of length  $\geq 200$ . But still, adding a single more distantly related sequence can make it hopelessly slow. This is because the 'band' needs to be much 'wider' when the sequences are not closely related. In fact the *problem* of optimal multiple sequence alignment is *NP-hard*. It is very unlikely that there is an acceptably efficient general algorithm that solves it.

**Simultaneous multiple alignment**

In principle it is easy: use dynamic programming, filling in a matrix  $F$ , just as for pairwise alignment. In pairwise alignment (two sequences), the  $F$  matrix is two-dimensional. In a multiple alignment of  $K$  sequences, the  $F$  matrix will be  $K$ -dimensional. Problem: Assume we have  $K$  sequences of length  $\ell_1, \dots, \ell_K$ . Then the size of the  $F$  matrix will be  $\ell_1 \cdot \ell_2 \cdot \dots \cdot \ell_K$ . When  $K = 6$  and  $\ell_1 = \ell_2 = \dots = \ell_K = 200$ , this is  $64,000,000,000,000$  numbers. This requires 256,000,000 MB storage, corresponding to 1 million PCs. Far too large for present-day and foreseeable computers.

**The practical solution: Progressive alignment (Feng and Doolittle 1987)**

Given  $K$  sequences to align (globally), do not align all sequences at the same time. Instead:

- Compute the  $\frac{K(K-1)}{2}$  pairwise distances between all pairs  $(s^k, s^l)$  of sequences.
- Using the distances, build a tree (as in phylogeny) to guide the progressive multiple alignment process.
- To build a multiple alignment, start with closely related sequences, add more distantly related sequences later.

Problem: the result is sensitive to the order in which sequences are added.

The order should follow the phylogeny, but the phylogeny is not known until the alignment is finished ...

To align the two first sequences, use ordinary pairwise alignment by dynamic programming (explained last week). To align a new sequence  $s$  to an existing alignment  $a$ , compute the pairwise alignment between  $a$  and every sequence in  $a$ . Then choose the best of these alignments, inserting gaps in the other sequences of  $a$  as necessary. Always preserve existing gaps in  $a$ .

To align an alignment  $a_1$  to another alignment  $a_2$ , compute a pairwise alignment between every sequence in  $a_1$  and every sequence in  $a_2$ . Then choose the best of these alignments, inserting gaps in the other sequences of  $a_1$  and  $a_2$  as necessary. Always preserve existing gaps in  $a_1$  and  $a_2$ .

Number of seqs	Time/s	Space/MB	Comments
4	< 1		4 closely related
5	8		5 closely related
6	30		6 closely related
6+1	> 4380		6 closely related + 1 more remotely (program had to be killed)

**Computer experiments with MSA: time and space consumption**

Computing multiple alignment of some bcll class B beta-lacamasases; length approximately 260 aa.

**ClustalW/ClustalX (1994): Further refinements**

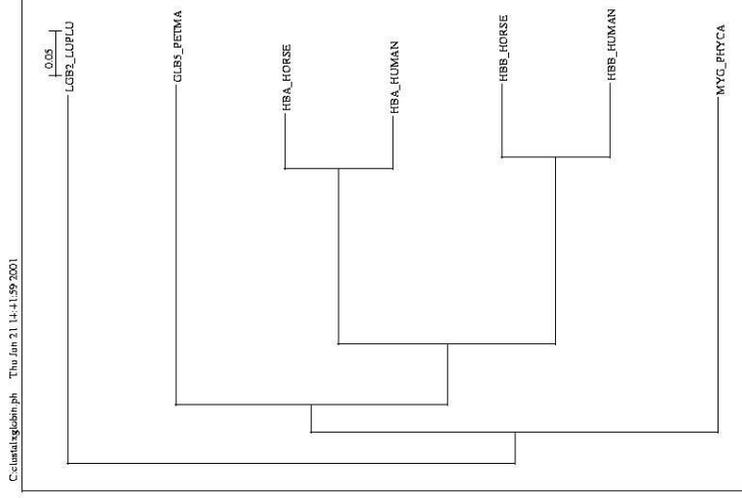
The pairwise distances used to build the guide tree are weighted by sequence similarity, so that closely related sequences are given a lower weight. This is to compensate for sequence selection bias.

The guide tree is built using Neighbour Joining to build the guide tree; an older version of ClustalW used UPGMA. Neighbour Joining and UPGMA ('unweighted pair group method using arithmetic averages') are two methods for building a clustering/phylogenetic tree, given a distance matrix. Neighbour Joining produces better results.

Use 'harder' substitution matrices for the early alignments (closely related sequences), e.g. GONNET120. Use 'softer' substitution matrices for the later alignments (distantly related sequences), e.g. GONNET350.

Gap opening cost is increased near existing gaps (to encourage fewer, longer gaps), and near hydrophobic amino acids (because they are usually inside the protein, where insertions and deletions are less likely to survive). Gap extension cost is decreased when aligning a new sequence to a gapped stretch of an existing alignment.

The graphical user interface ClustalX permits some visualization and very limited sequence editing.



- Servers, and other multiple alignment tools**
- ClustalW is available for download at <http://ftp-ibmc.u-strasbg.fr/pub/ClustalX/>
  - There is a ClustalW server at EBI <http://www.ebi.ac.uk/cluster/>
  - To edit ClustalW trees, and to include them in Word documents, use [Treeview](http://taxonomy.zoology.gla.ac.uk/rod/treeview.html):
  - T-Coffee server at <http://igs-server.cnrs-mrs.fr/Tcofee/>
  - Multalin, looks similar to ClustalW, at <http://www.toulouse.inra.fr/multalin.html>
  - AMPS (Alignment of Multiple Protein Sequences) = [multalign](http://www.compbio.dundee.ac.uk/Software/Amps/amps.html)
  - Pieup/GCG from Genetics Computer Group (commercial)
  - Many more tools, including manual editing tools, are listed at <http://www.expsy.org/tools/>
  - Example: [jalview](http://www.jalview.org/), <http://www.jalview.org/>
  - Example: [CINEMA](http://bioinf.man.ac.uk/dbbrowser/CINEMA2.1/), <http://bioinf.man.ac.uk/dbbrowser/CINEMA2.1/>

**BLOCKS motifs**

A BLOCKS motif is an ungapped multiple alignment of related sequences.

BLOCKS motif for class B beta-lactamases (pattern IPB001018A):

```

BLAB_ERHNY|P26918 ( 87) PLELVINTNTHYHTRDRAAGNAYWKSIGAKVSTRQTRDL 34
BLAB_BACPR|P25910 ( 91) KVTTFIPNHHGDCIGLGLYLRKRGVQSYANQMTIDL 100
BLAB_ERMA|P52699 ( 87) KIKGSISSHSHFSDSTGIEIWLNSRSIPTYSLELTEL 31
BLA1_XANMA|P52700 (154) ANADRIVMDSGEVITLVGGIVFTAHFMAGHTPSTAWTM 72
BLA2_BACE|P04190 (108) RVTDVIIITHAHADRIGGKIKTKERGIKASHSTALTAEL 35
BLAB_BACCE|P14488 (107) RVTDVIIITHAHADRIGGKTKERGIKASHSTALTAEL 43
BLA2_BACBP|P10425 (108) RVTDVIIITHAHADRIGGKTKERGIKASHSTALTAEL 41
051899 (154) ASADRIMDGEVVTGGAFTAHFMGHTPSTAWTM 59
067103 ( 88) PVIYAIVTTHYHLDHMHVYGAKFKKAKVIAHHRKIKKEF 72
067667 ( 90) PIRFLVVTTHYHLDHMHVYGAKAFARFEGAEVIAHHEWMAFDY 89
068919 ( 76) PVRLLVNTTHFHGDHSFNGNIGIKDAVIVAAHRRRTTEM 82
Q02057 ( 62) PGRIVVNTTHFHGDHAFNGQVFAFGPRTRIIAHEDMRSAM 100
    
```

The number to the right is the weight of the sequence (100 = most distant from consensus).

The BLOSUM substitution matrices were created from BLOCKS motifs.

**T-Coffee (2000): More precise multiple alignments**

In ClustalW, an early mistake in progressive alignment cannot be corrected later. Consider these sequences:

```

1 GARFIELD THE LAST FAT CAT
2 GARFIELD THE FAST CAT
3 GARFIELD THE VERY FAST CAT
4 THE FAT CAT
    
```

The sequences will be aligned (1,2), then ((1,2),3), then (((1,2),3),4), giving the wrong (left) alignment:

```

GARFIELD THE LAST FA-T CAT
GARFIELD THE FAST CA-T ---
GARFIELD THE VERY FAST CAT
----- THE ----FA-T CAT
----- THE ---- FA-T CAT
    
```

T-Coffee (Notredame, Higgins, Heringa 2000) improves on ClustalW/ClustalX as follows:

- Build a 'library' of all pairwise local and all pairwise global alignments between the sequences.
- For each matching pair of amino acids from two sequences, record a weighted constraint.
- For instance, C in sequence 2 is likely to match C in sequence 1, sequence 3 and sequence 4.
- Use the information from the library when doing the progressive multiple alignment.

In the exercises you will see that this makes a difference in real life.

**Motifs and multiple alignments**

A motif is a section of a multiple alignment where many sequences agree.

There are many ways to represent motifs:

Database	Representation	URL
BLOCKS	Blocks	<a href="http://www.blocks.fhcrc.org/">www.blocks.fhcrc.org/</a>
PROSITE patterns	Regular expressions	<a href="http://www.expasy.org/prosite/">www.expasy.org/prosite/</a>
PROSITE profiles	Profiles (position-specific scores)	<a href="http://www.expasy.org/prosite/">www.expasy.org/prosite/</a>
Pfam	Hidden Markov Models	<a href="http://www.sanger.ac.uk/Pfam/">www.sanger.ac.uk/Pfam/</a>
Interpro	Integration of other databases	<a href="http://www.ebi.ac.uk/interpro/">http://www.ebi.ac.uk/interpro/</a>

Usually based on high-quality sequences from SwissProt.

A BLOCKS block, or a PROSITE pattern or profile, or a Pfam HMM, (partially) characterizes a protein family.

Given a PROSITE pattern or profile, one can find protein sequences that match the profile.

Given a sequence, one can find a matching PROSITE pattern or profile, and hence the protein family.





**Phylip: Classic phylogenetic tools by Joseph Felsenstein**

Download from <http://evolution.genetics.washington.edu/phylip.html>  
 Or use <http://bioweb.pasteur.fr/seganal/phylogeny/phylip-uk.html>

- Clustering based on pairwise evolutionary distance: DNADIST or PROTDIST, followed by NEIGHBOR or FITCH or KITSCH
- Parsimony: DNAPARS or PROTPARS
- Maximum likelihood: fastDNAmI

After computing a phylogeny (and a treefile), use Phylip to draw the trees:

- Use DRAWGRAM to draw a rooted tree
- Use DRAWTREE to draw an unrooted tree

**Cladistics**

Cladistics is a particular school of phylogeny founded by Willi Hennig: <http://www.cladistics.org/>

A *cladogram* is a tree with organisms only at the leaves, and each branch is a clade. A *clade* is a monophyletic taxon: a group of organisms which includes the most recent common ancestor of all of its members and all of the descendants of that most recent common ancestor. Cladistics determines the evolutionary relationships based on *derived* similarity. There seems to be some disagreement (or perhaps, a religious wars) around phylogeny, and in particular cladistics. For instance, Felsenstein's Phylip webpage says: *Inspiration to keep going came from Cladistics magazine, for endless negative reviews and disparaging comments. If they liked it I'd worry.*

**Bootstrap: How reliable is the tree?**

Usually some branches in a phylogenetic tree are poorly determined. Two estimate which are and which aren't, one uses bootstrap.

- Create many (100) new synthetic multiple alignments as follows: Randomly select columns (with replacement) from the original multiple alignment.
- From each synthetic multiple alignment, create a synthetic phylogenetic tree.
- For each branch in the original tree, count how many times it appears in the synthetic trees.

The synthetic alignment to the right is created from columns 5 6 1 2 3 9 7 3 6 of the left one:

ESFGDLSSTPD	DDLSEFPSPFL
DSFGDLSNPG	DDLDSFPSPFL
PHF-DLS---	DDLPHF-SFL
PHF-DLS---	DDLPHF-SFL
PKFKGLTTAD	GGLPKFATFL
DRFKHLKTEA	HHLDRFEKFL
SFLKGTSEVP	GGTSFLVSLT

The program's pseudo-random number generator needs a seed, so you must choose a random number.

**The phylogenetic utility of different kinds of sequences**

- Protein (amino acid) sequences can be used to determine distant evolutionary events. The short-term non-functional changes in the third codon position are hidden.
- Protein-coding DNA sequences can be used to determine recent evolutionary events for closely related organisms. The short-term non-functional changes in the third codon position are visible.
- Some RNA-coding DNA sequences can be used to study the phylogeny of ancient organisms. For instance, 16 S rRNA is part of the protein-synthesizing machinery (the ribosome). This machinery is ancient, universal and essential, so the 16 S rRNA sequence is well conserved across long evolutionary distances.
- 16 S rRNA was used by Woese et al. to establish the domains of Archaea, Eukaryota and Eubacteria (1977). Mitochondrial DNA (in humans) is inherited only from the mother and therefore permits 'single inheritance' phylogeny despite sexual reproduction.