

What is a multiple alignment

- Alignment of 3 or more DNA or AA sequences
 - EVERY nucleotide position is assumed homologous

Example multiple alignment

What can you do with multiple alignments?

- Identify conserved and divergent regions
 - Functional inferences
 - Phylogenies
 - Relationships between organisms
 - Define motifs or domains
 - Identify additional family members; motif searching more sensitive than Blast.
 - Secondary structure prediction
 - Protein structure

ClustalW

- Align each pair of sequences and calculate distance matrix.
 - Calculate an neighbor joining phylogenetic tree.
 - Progressively align sequences, starting with most similar sequences based on tree.
 - **Retain information from closely related sequences that are easy to align.**

Distances

- Percent divergence of sequences.
 - Sequences aligned, number identities/number of residues, and subtract from 1

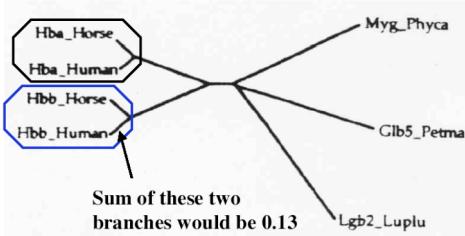
AWYSTRGHKIL
AWYTTRGHKIL

$$\begin{aligned} \frac{9}{10} &= 0.9 \\ 1 - 0.9 &= 0.10 \end{aligned}$$

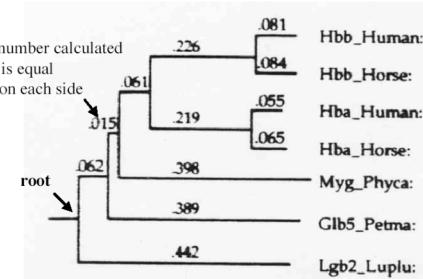
Make a matrix of all pair-wise distances.
Sequences that have the lowest divergence are
the most closely related and are grouped together.

fbb_Human	1	-				
fbb_Horse	2	17	-			
fba_Human	3	.59	.60	-		
fba_Horse	4	.59	.59	13	-	
Mvg_Phyc	5	.77	.77	.75	.75	-
Gibb_Petra	6	.81	.82	.73	.74	.80
gb2_Luplu	7	.87	.86	.86	.88	.93
	1	2	3	4	5	6

Unrooted tree shows the grouping of the most closely related sequences. Length of “branch” is the distance. Make a “rooted tree” by midpoint method

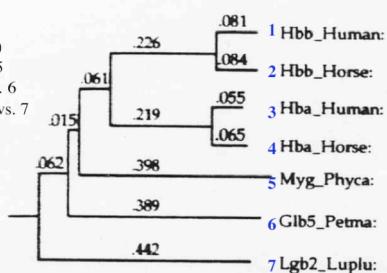


Make a “rooted” tree using midpoint method.
Longest branch considered root.



Progressive alignment: retain info as you align sequences

- 3 vs. 4
- 1 vs. 2
- (3,4) vs. (1,2)
- (3,4,1,2) vs. 5
- (3,4,1,2,5) vs. 6
- (3,4,1,2,5,6) vs. 7



Sequence weighting

- Way to down-weight duplicate information
 - Gives more weight to rare, divergent sequences

Sequence weighting: distance to root, but share distances with common branches

Phylogenetic tree with bootstrap values:

- Root (0.081+)
- Node 062: .398
- Node 061: .015, .061
- Node 015: .219
- Node 065: .055
- Node 084: .226
- Node 081: .081
- Leaf 1: Hbb_Human: 0.221
- Leaf 2: Hbb_Horse: 0.225
- Leaf 3: Hba_Human: 0.194
- Leaf 4: Hba_Horse: 0.203
- Leaf 5: Myg_Phyc: 0.411
- Leaf 6: Glb5_Petma: 0.398
- Leaf 7: Lgb2_Luplu: 0.442

Sequence weighting: distance to root, but share distances with common branches

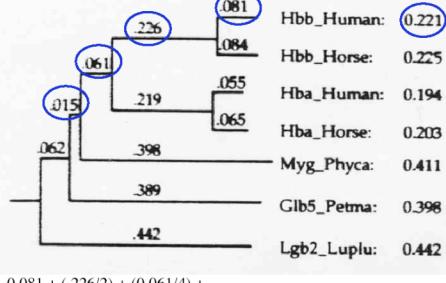
Phylogenetic tree with bootstrap values at nodes:

- Root: .015
- Node 061: .226
- Node 084: .084
- Node 055: .219
- Node 065: .065
- Node 398: .398
- Node 389: .389
- Node 442: .442

Species names and their corresponding bootstrap values:

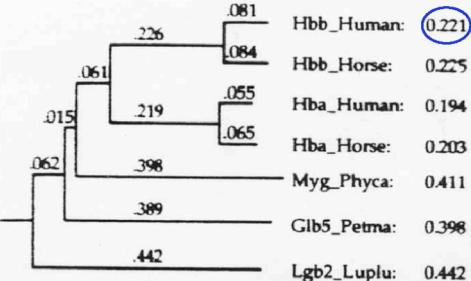
- Hbb_Human: 0.081
- Hbb_Horse: 0.225
- Hba_Human: 0.194
- Hba_Horse: 0.203
- Myg_Phycal: 0.411
- Glb5_Petma: 0.398
- Lgb2_Luplu: 0.442

Sequence weighting: distance to root, but share distances with common branches



$$0.081 + (.226/2) + (0.061/4) +$$

Sequence weighting: distance to root, but share distances with common branches



$$0.081 + (.226/2) + (0.061/4) + (0.015/5) + (0.062/6)$$

Gap penalties

- Not just gap open and gap extend penalty.
- Use information in alignment to guide choice of gap penalties.
- Position specific gap penalties.

Initial Gap penalty

- Just like Blast
 - Gap open
 - Gap extend

Sequence similarity

- Closely related sequences less likely to have gaps, reduce the penalty.
- Gap open penalty decreases on linear scale for more divergent sequences

Matrix used

- Different matrices have different optimal gap open penalties.
 - Similar to sequence divergence
- Take average score for mismatched residues as a scaling factor.

Lowered gap penalties for existing alignments

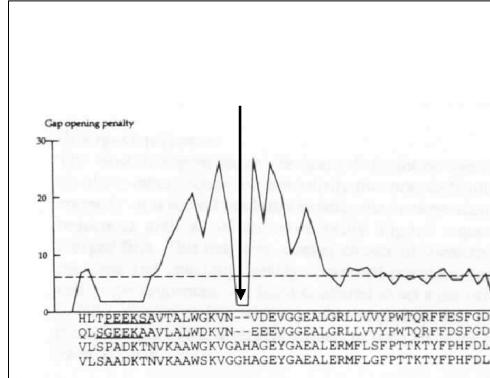
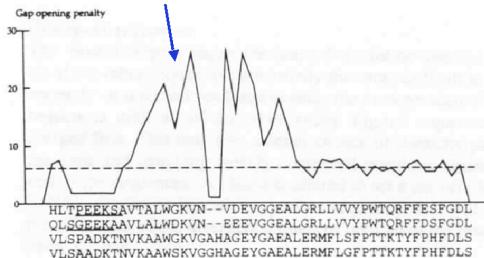
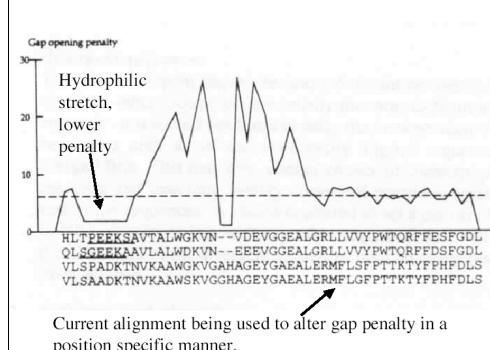
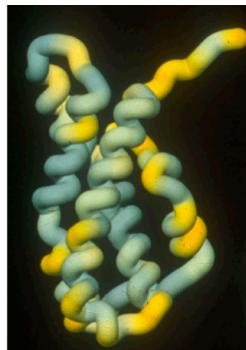
- If a gap already occurs in the progressive alignment, it is likely that it will occur in the same position in subsequent alignments
- Reduced based upon proportion of sequences with a gap in that position

Increased penalty near existing gaps.

- Force gaps to occur in same position in the alignment
- Increase gap penalties in region 8 residues from existing gaps

Hydrophilic stretches

- Reduce gap in a run of 5 hydrophilic residues
- Hydrophilic residues tend to be in loops, where most gaps occur.



Clustal using different matrices during alignment (remember difference between BLOSUM62 and BLOSUM 45 in BLAST)

- Matrix varied throughout alignment dependent upon %ID
 - BLOSUM or PAM (BLOSUM default)
- Sequences 80% - 100% identical
 - BLOSUM80
- Sequences 60% - 80% identical
 - BLOSUM62 (Blast default)
- Sequences 30% - 60% identical
 - BLOSUM45
- Sequences 0% - 30% identical
 - BLOSUM30

Multiple Fasta formatted unaligned sequences

```
>red18
MRSILVLLCVLLMAICAAADKKSTVSKENAAMKVAMIKFLDSRTDRFKKRIE
KIGYPITPPQYTTLLYYNRERLMDWCHNVEVSKKIIILGGNKLNNKNFARMGRIIGWKN
QWILKRRQHMHMVRVMRKYKASAIAKKIVAMKVADLPCN
>white18
MRSILVLLCVLLMAICAAADKKTTVSKENAAAMKIAMIKFLDARAAGKFKKRVE
NNGYPYITPPQWTTLLYYNRQRMLMECHTYVEFSKKIILMGGNKLNNKNFARMGRIIGWKN
QWILKRRQHMHMVRVMRKYKSTAIAKKIVAMKVADLPCN
>asimi18
MRSILVLLCVLLMAICAAADKKTTSVSKNEEAAMKVAMMKFLDMKAGVFKEIIIE
DMGYPYITPPQWTTLLYYNRERLIEFCRSFLASKKIIILGGNKLNNKANFARMGRILGWKS
QWAVRQRQHMHMVRVSRRHTSTAIAKRIIVAMKVADLPCN
>pink18
MRFLLLLCLVLMAGNVSGAVCKRKPNVNGKIVVKEKNKAAMKIGFMEYLDALKLVFKRHWLV
GANWLQKPFETDEMRYLAIRLKVCHGTYIWQSRLIMLYRPLNEKYFKKVGRLAWRN
YLIVFRMRMIGVLKKNLKRSEITKPMQKLDDTKDGECPCPVRKING
>green18
MRSILVLLCVLMAGCVAFDVVVVSRRQEYSVYQRGMVNFLDEEMHKLVKRFR
DMRWNLQPGFVFLLKKVNRRERMMRKYCMYARYSKKILQLKHLPVNKKTLTKMGRFVGYN
YGVIRELYADYFVFRDVQGFRGPKMTAAMRKYSSKDPGTPCKNEKRRG
```

Clustal Demo

ClustalW Parameters: for most uses the default works fine

- Output format
 - Phylip, PIR, etc.
- Gap open/extend
- Matrix
 - BLOSUM is best
- Ktup
 - Word size

ClustalW will give you a guide tree. There are better methods for constructing phylogenies, and we will go over them later

Your guide tree:
phy440_209301-329072.faa

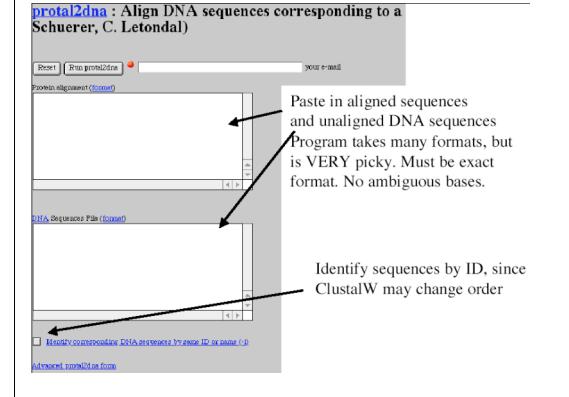
```
:
&lt;t tip_id=0 weight=0.27951;
&lt;t tip_id=1 weight=0.05723;
&lt;t tip_id=2 weight=0.05204;
&lt;t tip_id=3 weight=0.12081;
&lt;t tip_id=4 weight=0.36159;
&lt;t tip_id=5 weight=0.37268;
&lt;t tip_id=6 weight=0.04741;
```



Aligning Protein coding regions

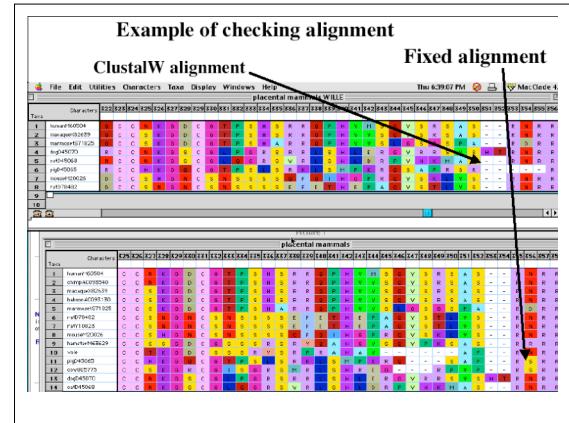
- Align protein sequence
 - Gaps are easier to deal with; wobble
 - AA information regarding substitutions
- Use aligned protein sequence to align DNA
 - Use program Protal2DNA

Protal2DNA Demo



You must check alignments by eye

- Each position must be homologous.
- Focus on regions around gaps.
- Look for conservative AA changes, as indicated by sequence editor color.
- If doing protein coding region, check both protein and DNA alignment (silent changes).
- Do not over edit. If sequence is so far out of alignment you should try changing gap parameters in Clustal.



Use Se-Al to check alignment

- Look at Protein
- Look at DNA (minimize silent changes)
- Export to different file formats
- Color coded amino acids
- Has some problems, occasionally crashes the computer (but it is free). Save work as you go.

Se-al example

