


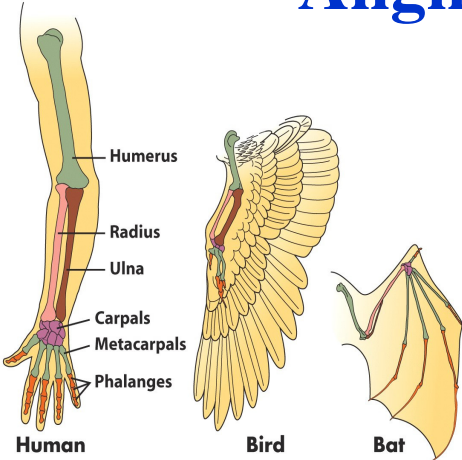
Introduction to Bioinformatics

Sami Khuri
Department of Computer Science
San José State University
San José, California, USA
khuri@cs.sjsu.edu
www.cs.sjsu.edu/faculty/khuri

©2010 Sami Khuri

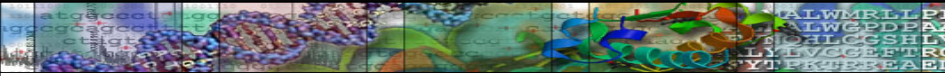


Multiple Sequence Alignment



- ❖ Progressive Alignment
- ❖ Iterative Pairwise
- ❖ Guide Tree
- ❖ ClustalW
- ❖ Co-linearity
- ❖ Multiple Sequence Alignment Editors

©2010 Sami Khuri



```


Wombat      : AAAGTTAATGAGTGGTTATCCAGAAGTAGTGACATTTTAGCCCTCTGATAACTCCAAACGGTAGGAGCCATGACGAGAGCCGAGA : 83
Opossum     : AAAGTTAATGAGTGGTTATTCCAGAAGTAGTGACGTTTTAGCCCCAGATTACTCAAGTGTTAGGAGCCATGAACAGAAATGCAGA : 83
Armadillo   : AAAGTTAACGAGTGGTTTTCCAGAAGTGTGACATATTAACTTCTGATGACTCAGACAGATAGGGGGTCTGAATTAATGTCAGA : 83
Sloth       : AAAGTTAATGAGTGGTTTTCCAGAAGTGTGACATACTAACTTCTGATGACTCAGACAAATGGGGGTCTGAATCAAATGCAGA : 83
Dugong      : AAAGTTAATGAGTGGTTTTCCAGAAGTGTGGCCTG-----GATGACTTGCATGATAAGGGGTCTGAGTCAAATGCAGA : 74
Hyrax       : AAAGTTAATGAGTGGTTTTCCAGAAGTGNACACCCTA-----AGTGATTCACTAGTGGGGGTCTGAATTAATGGAATA : 74
Aardvark    : AAAGTTAATGAGTGGTTTTCCAGAAGTGTGGCCTG-----GATGGCTCAGATGATGAAGGGTCTGAATCAAATGCAGA : 74
Tenrec      : AAGGTTAACGAGTGGTTTTCCAAAGCCACGGCCTG-----GTTGACTCTCGCATGGGGCCTGAGTCAAGCCGAGA : 74
Rhinoceros  : AAAGTTAATGAGTGGTTTTCTAGAAGCGATGAAATGTTAACTTCTGATGACTCAGATGATGGGGCCTGAATCAAATACTGA : 83
Pig         : AAAGTTAATGAGTGGTTTTCTAGAAGCGATGAAATGTTAACTTCTGACGACTCAGAGGACAGGAGGTCTGAATCAAATACTGG : 83
Hedgehog    : AAAGTGAATGAAATGGCTTTCCAGAAGTGTGAAGTGTAACTTCTGATGACTCATATGAAAGGGATCTAAATCAAATACTGA : 83
Human       : AAAGTTAATGAGTGGTTTTCCAGAAGTGTGAAGTGTAACTTCTGATGACTCAGATGATGGGGAGTCTGAATCAAATGCCAA : 83
Rat         : AAAGTGAATGAGTGGTTTTCCAGAAGTGTGAAATGTTAACTTCTGACAAATGCATGACAGGAGGGCCTGCTCAAATGCAGA : 83
Hare        : AAAGTTAACGAGTGGTCTCCAGAAGTAAATGAAATGTTAACTTCTGATGACTCACTTGACCGGGGTCTGAATCAAATGCCAA : 83
AaaGTLAatGAgTGGtTtTccAgAagt atga T gatgactca gat g gg cFga t aaatgc ga

* 20 * 40 * 60 * 80

Wombat      : GGTGCCAGTGGCTTAGAAGATGGGCATCCAGATACCGAGAGGGAAATCTAGCGTTTCTGAGAAAGACTGAC : 156
Opossum     : GGCACCAATGCTTTAGAAATATGGCATGTAGAGACA---GATGAAATCTAGCATTTCTGAAAGACTGAT : 153
Armadillo   : AGTAGCTGGTGCATTGAAGTT-----TCAAAGAAGTAGATGAAATTTCTAGTTTTCCAGAGAAATAGAC : 150
Sloth       : AGTAGCTGGTGCATTGAAGTT-----CCAAATGAAGTAGATGATATTCTGGTCTTCCAGAGAAATAGAC : 150
Dugong      : AGTAGCTGGTGCATTGAAGTT-----CCAGAAAGTAGATGATATTCTAGTTCTTCCAGAGAAATAGAC : 141
Hyrax       : AGTGGCTGGTCCAGTAAACT-----CCAGGTGAAGTAGATGATATTCTAGTTTTCCAGAGAAATAGAC : 141
Aardvark    : AATAGTGGTGGCATTGAAGTT-----TCAAATGAAGTAGATGATATTCTGGTCTTCCAGAGAAATAGAC : 141
Tenrec      : CGTAGCTGGTGCATTGAAGTT-----CCAGACGAAGCATGTGAATCTTATAGTTCTCCAGAGAAATAGAC : 141
Rhinoceros  : AGTAGCTGGTGCAGTGAAGTT-----CAAATGAAGTAGATGATATTCTGGTCTTCCAGAGAAATAGGC : 150
Pig         : AGTAGCTGGTGCAGTGAAGTT-----CAAATGAAGTAGATGATATTCTGGTCTTCCAGAGAAATAGGC : 150
Hedgehog    : AGTAACTGTAACACAGAAAGTT-----CAAATGCAATAGATGATTTTTTTGGTCTTCCAGAGAAATAGAC : 150
Human       : AGTAGCTGATGATTTGAGCGTT-----CTAAATGAGGTAGATGAAATTTCTGGTCTTCCAGAGAAATAGAC : 150
Rat         : AGCTGCTGTTGTTGTTGAAGTT-----TCAAATGAAGTGGATGATGTTTCAAGTTCTTCAAAGAAATAGAC : 150
Hare        : AGTGGCTGGTGCATTGAAGTC-----CCAAAGGAGGTAGATGATATTCTGGTCTTCCAGAGAAATAGAC : 150
gt gctg tgc t gAagtt cA a gaag a atggatatT t Gtt TtCagAGAA Atagac
    
```

Part of the alignment of the DNA sequences of the BRCA1 gene
From "Bioinformatics and Molecular Evolution" by Paul Higgs and Teresa Attwood

©2010 Sami Khuri



Aligning BRCA1 Sequences

```

* * * * *
Wombat      : KVNEWLSRSSDILASDNSNGRSHEQSAEVPSALEDGHPDTAEGNSSVSEKTD : 52
Opossum     : KVNEWLFRSNDVLPADYSSVRSHEQNAEATNALEYGHVET-DGNSSISEKTD : 51
Armadillo   : KVNEWFSRGDDILTSDDSHDRGSELNAEVAGALKV--SKEVDYSSFSSEKID : 50
Sloth       : KVNEWFSRSDDILTSDDSHNGGSESNAEVVGALKV--PNEVDGYSGSSEKID : 50
Dugong      : KVNEWFFRSDGL---DDLHDKGSESNAEVAGALEV--PEEVHGYSSSSEKID : 47
Hyrax       : KVNEWFSRSDNL---SDSPSEGSSELNGKVAGPVKL--PGEVHRYSSFPENID : 47
Aardvark    : KVNEWFSRSDGL---DGSHEGSESNAEIGGALEV--SNEVHSYSGSSEKID : 47
Tenrec      : KVNEWFSKSHGL---GDSRDGRPESGADVAVAFEV--PDEACEYSSESPEKTD : 47
Rhinoceros  : KVNEWFSRSEILTSDDSHDGGPESNTEVAGAVEV--QNEVDGYSGSSEKIG : 50
Pig         : KVNEWFSRSEMLTSDDSQDRRESNTEGVAGAAEV--PNEADGHLGSSEKID : 50
Hedgehog    : KVNEWLSRSEDELLTSDDSYDKGSKSKTEVTVTTEV--PNAIDXFFGSSEKIN : 50
Human       : KVNEWFSRSEDELLGSDSHDGESESNAKVADVLDV--LNEVDEYSGSSEKID : 50
Rat         : KVNEWFSRTGEMLTSDNASDRRPASNAEAAVVLEV--SNEVDGCFSSSKKID : 50
Hare        : KVNEWFSRSNEMLTDPDSDLRRSESNAKVAGALEV--PKEVDGYSGSTEKID : 50
KVNEWfs4   6   d   s   e   n   e   e   ki
    
```

Alignment of BRCA1 protein sequences for the same region on the gene
From "Bioinformatics and Molecular Evolution" by Paul Higgs and Teresa Attwood

©2010 Sami Khuri



What is Multiple Alignment

Most simple extension of pairwise alignment

Given:

- Set of sequences
- Match matrix
- Gap penalties

Find:

Alignment of sequences such that an optimal score is achieved.

©2010 Sami Khuri



Uses of Multiple Alignment

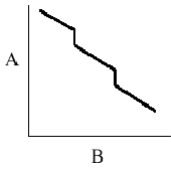
A good **alignment** is critical for further analysis

- Determine the **relationships** between a group of sequences
- Determine the **conserved** regions
- **Evolutionary Analysis**
 - Determine the phylogenetic relationships and evolution
- **Structural Analysis**
 - Determine the overall structure of the proteins

©2010 Sami Khuri

Alignment Difficulties

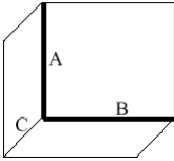
We cannot use sequence comparison algorithms (dynamic programming) and just add more sequences.



A
B

→

Add 1
sequence



A
B
C

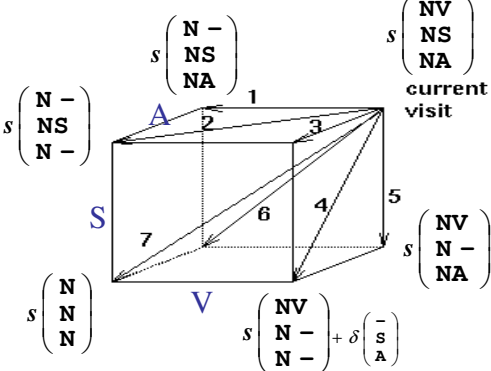
With each sequence, another dimension is added that we need to search for the optimal alignment.

©2010 Sami Khuri

Hyperlattice Computation

k=3 2^k-1=7

$$s \begin{pmatrix} NV \\ NS \\ NA \end{pmatrix} = \max \begin{matrix} s \begin{pmatrix} N \\ N \\ N \end{pmatrix} + \delta \begin{pmatrix} V \\ S \\ A \end{pmatrix} \\ s \begin{pmatrix} NV \\ N - \\ N - \end{pmatrix} + \delta \begin{pmatrix} - \\ S \\ A \end{pmatrix} \\ s \begin{pmatrix} N - \\ NS \\ N - \end{pmatrix} + \delta \begin{pmatrix} V \\ - \\ A \end{pmatrix} \\ s \begin{pmatrix} NS \\ N - \\ NA \end{pmatrix} + \delta \begin{pmatrix} - \\ V \\ A \end{pmatrix} \\ s \begin{pmatrix} N - \\ N - \\ NA \end{pmatrix} + \delta \begin{pmatrix} V \\ S \\ - \end{pmatrix} \\ s \begin{pmatrix} NS \\ NA \\ NV \end{pmatrix} + \delta \begin{pmatrix} - \\ - \\ - \end{pmatrix} \\ s \begin{pmatrix} N - \\ NA \\ NV \end{pmatrix} + \delta \begin{pmatrix} - \\ S \\ - \end{pmatrix} \\ s \begin{pmatrix} NV \\ N - \\ NS \end{pmatrix} + \delta \begin{pmatrix} - \\ - \\ A \end{pmatrix} \end{matrix}$$



current visit

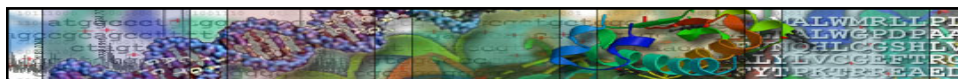
©2010 Sami Khuri



MSA: Exact vs. Heuristic

- The **exact algorithm**
 - traverses the entire search space
 - finds overall measure of alignment quality and tries to maximize this quality.
- The operation is computationally intensive.
- The largest computers can only optimally align a few sequences (7-8).
- Therefore, we have to use **heuristics**; i.e., faster algorithms, if we want to align many sequences.

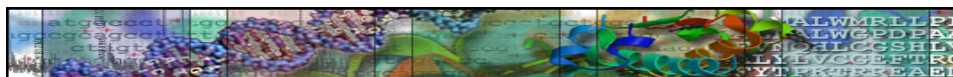
©2010 Sami Khuri



Heuristic Algorithms

- Based on a **progressive pairwise** alignment approach
 - ClustalW (**C**luster **A**lignment)
 - PileUp (GCG)
 - MACAW
- Builds a global alignment based on **local alignments**
- Builds local multiple alignments
- Based on **Hidden Markov Models**
- Based on **Genetic algorithms**.

©2010 Sami Khuri



Progressive Strategies for MSA

- A common strategy to the MSA problem is to **progressively align** pairs of sequences.
 - A starting pair of sequences is selected and aligned
 - Each subsequent sequence is aligned to the previous alignment.
- **Progressive alignment** is a greedy algorithm.


©2010 Sami Khuri



Iterative Pairwise Alignment

- The **greedy algorithm**:
 - align some pair*
 - while not done*
 - pick an unaligned string “near”*
 - some aligned one(s)*
 - align with the previously aligned group*
- There are many variants to the algorithm.


©2010 Sami Khuri



ClustalW: Package for MSA

- **ClustalW** [the **W** is from **W**eighted] is a software package for the MSA problem.
- Different weights are given to sequences and parameters in different parts of the alignment to and create an alignment that makes sense biologically.
- **Scalable Gap Penalties** for protein profile alignments
 - A gap opening next to a conserved hydrophobic residue can be penalized more heavily than a gap opening next to a hydrophilic residue.
 - A gap opening very close to another gap can be penalized more heavily than an isolated gap.

©2010 Sami Khuri



Steps of ClustalW

S₁ —————

S₂ —————

S₃ —————

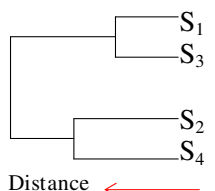
S₄ —————

↓ All Pairwise Alignments

Similarity Matrix

	S ₁	S ₂	S ₃	S ₄
S ₁		4	9	4
S ₂			4	7
S ₃				4
S ₄				

Cluster Analysis →



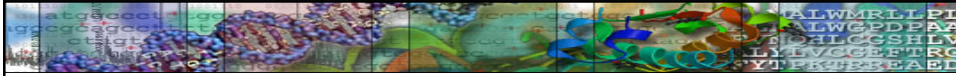
Dendrogram

Distance ←

Multiple Alignment Step:

1. Aligning S₁ and S₃
2. Aligning S₂ and S₄
3. Aligning (S₁,S₃) with (S₂,S₄).

©2010 Sami Khuri



ClustalW: An Example

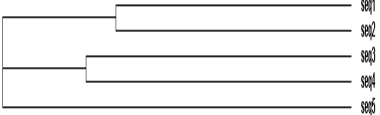
CLUSTAL W (1.82) multiple sequence alignment

```


seq3      FEGGILVEAL 10
seq4      FDG-ILVQAV 9
seq5      YEGGAVVQAL 10
seq1      YDG-GAVEAL 9
seq2      YDG-G--EAL 7
          :::*   :*:
    
```

* = identity
: = strongly conserved
. = weakly conserved

By using the same five sequences and aligning them with CLUSTALW, we get the illustrated results.



©2010 Sami Khuri



```

HEA_HUMAN  -----VLSPADKTNVKAAMGKVGAHAGEYGAELERMFLSFPTTKTYFFPHF-DLS- 49
HEA_HORSE  -----VLSAADKTNVKAAMSKVGGHAGEYGAELERMFLGFPTTKTYFFPHF-DLS- 49
HEA_CHICK  -----MVLSAADKNNVKGIFTKIAGHAEYGAETLERMFTTYPTTKTYFFPHF-DLS- 50
HBB_HUMAN  -----VHLTPEEKSAVTALWGKVNV--VDEVGGEALGRLLVVYPWTQRFESFGDLST 50
HBB_BOSMU  -----MLTAEKKAAVTAFWGKVK--VDEVGGEALGRLLVVYPWTQRFESFGDLSS 49
HBB_HORSE  -----VQLSGEKKAAVLALMDKVN--EEVGGEALGRLLVVYPWTQRFESFGDLSN 50
HBB_MACGI  -----VHLTAEKNAITSLMGKVA--IEQTGGEALGRLLVVYPWTQRFDFHFDLSN 50
MYG_PHYCA  -----VLSGEWQLVLVHVMKVEADVAGHGQDILIRLFSKHPETLEKFDKFKHLKT 51
GLB5_PETMA PIVDTGTSVAPLSAAEKTKIRSAWAPVYSIYETSGVDILVKFFTSPTAAQEFPPKFKGLIT 60
LGB2_LUPLU -----GALTESQAALVKSSWEEFMANIPKHTHFPLLVLEIAPAAKDLFSFLKGTSE 52
          * : : : : : : : : : : : : : : : : : : : : : : : : : : : :
          * : : : : : : : : : : : : : : : : : : : : : : : : : : :

HEA_HUMAN  ---HGSAQVKGHGKQVADALTNVAHVDD---MPNALSALSDPAHAKLRVDPVNFKL 100
HEA_HORSE  ---HGSAQVKAHGKQVGDALTLAVGHLD---LPGALSNSLSDPAHAKLRVDPVNFKL 100
HEA_CHICK  ---HGSAQIRGHGKQVVAALIEAANHIDD---IAGTSLKLSLSDPAHAKLRVDPVNFKL 101
HBB_HUMAN  PDAVMGNPKVKAHGKKVLDGFASDGLAHLDM---LKGTFATLSEHHDGDKLHVDPENFKL 105
HBB_BOSMU  ADAVMGNPKVKAHGKKVLDGFASDGLAHLDM---LKGTFALSEHHDGDKLHVDPENFKL 104
HBB_HORSE  PGAVMGNPKVKAHGKKVLDGFASDGLAHLDM---LKGTFALSEHHDGDKLHVDPENFKL 105
HBB_MACGI  AKAVMANPKVLAHGAKVLAIFGDAIKRLDN---LKGTFALSEHHDGDKLHVDPENFKL 105
MYG_PHYCA  EAEMKASEDLKKGHTVLTALGAILKKGKH---HEAELKPLAQSHATKHKIPIKYLEF 106
GLB5_PETMA ADQLKKSADVRWHAERLINAVMDAVASMDDT--EKMSHKLRLDLSGSHAKSPQVDPQYFKV 118
LGB2_LUPLU VP--QMNPELQAHAAGKVKLVYEAAIQLVTVGVVVVTDATLKNLGSVH*SKG-VADAHFPV 109
          . : : * : : : : : : : : : : : : : : : : : : : : : :
          . : : * : : : : : : : : : : : : : : : : : : : : : :

HEA_HUMAN  LSHCLLVTAAHLPAEFTPAVHASLDKFLASVSTVLTSKYR----- 141
HEA_HORSE  LSHCLLSTLAVHLPNDFTPAVHASLDKFLSSVSTVLTSKYR----- 141
HEA_CHICK  LGQCFLVVVAIHHPAALTPEVHASLDKFLCAVGTVLTAKYR----- 142
HBB_HUMAN  LGMVLCVLAHFGKKEFTPPVQAAYQKVVAGVANALAHKYH----- 146
HBB_BOSMU  LGMVLCVLAHFGKKEFTPPVQAAYQKVVAGVANALAHKYH----- 145
HBB_HORSE  LGMVLCVLAHFGKKEFTPELQASYQKVVAGVANALAHKYH----- 146
HBB_MACGI  LGMIIVICLAHFGKKEFTIDTQVAWQKLVAGVANALAHKYH----- 146
MYG_PHYCA  ISEAIHVLHSPHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG 153
GLB5_PETMA LAAVIADTVAAAG--DAGFEKLMSMICILLRSAY----- 149
LGB2_LUPLU VKEAIIKTKREVVGAKSWSEELNSAWTIAYDELAIVLKRKEMNDAA--- 153
          : : : : : : : : : : : : : : : : : : : : : : :
    
```

©2010 Sami Khuri

Red Blood Cells

Red blood cells contain several hundred hemoglobin molecules which transport oxygen

Oxygen binds to heme on the hemoglobin molecule

©2010 Sami Khuri

Hemoglobin Structure

- 4 subunits: 2 alpha and 2 beta (red and gold)
- Each subunit holds a heme group
- Each heme group contains an iron atom, which is responsible for the binding of oxygen

©2010 Sami Khuri



Practical Considerations

- When to use Clustal?
- Can be used to align any group of protein or nucleic acid sequences that are related to each other over their entire lengths.
- Clustal is optimized to align sets of sequences that are entirely co-linear, i.e. sequences that have the same protein domains, in the same order.



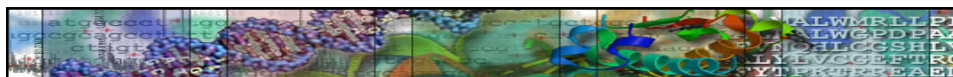
©2010 Sami Khuri



When Not To Use Clustal

- Sequences do not share common ancestry.
- Sequences are partially related.
- Sequences include short non overlapping fragments.

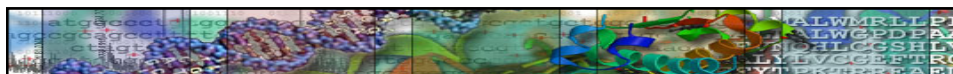
©2010 Sami Khuri



Structural Alignment

- What you really want to do is “align regions of similar function”.
- These are the areas that are evolutionarily conserved. (Folds, domains, disulfide bonds)
- **Problem**
 - The computer does not know anything about the structure or function of the proteins.
- **Solution**
 - Use computer alignment as a first step, then manually adjust the alignment to account for regions of structural similarity.

©2010 Sami Khuri



MSA Editors

- Once the multiple alignment is produced, it may be necessary to edit the sequence manually to obtain a more reasonable or expected alignment.
- Some of the considerations for an editor:
 - the use of colors to aid in the visual representation of the alignment,
 - the capability of recognizing the alignment format,
 - the ability of using the mouse to add, delete, or move sequences, thus allowing for an adequate windows interface.

©2010 Sami Khuri

Divide and Conquer Alignment

Use the **divide and conquer** technique to perform multiple sequence alignment.

MSA with the Divide & Conquer Method by J. Stoye
Gene 211(2), GC45-GC56, 1998.
(Gene-COMBIS).

<http://bibiserv.techfak.uni-bielefeld.de/dca/algorithm/>

©2010 Sami Khuri

Appropriate Approaches

	Comment	Appropriate approach
(a)	Sequences are related over their entire length.	Progressive global alignment method (e.g. CLUSTAL W).
(b)	Sequences share conserved blocks, separated by non-conserved regions containing large indels. Blocks are consistent (i.e. in the same order) but not necessarily uniform (i.e. some blocks may be missing in some sequences).	Block-based global alignment method (e.g. DIALIGN, ITERALIGN). Compare alignments produced by different programs (including progressive methods).
(c)	Sequences contain a non-consistent set of conserved blocks (i.e. some blocks are duplicated or occur in a different order along sequences).	Motif-based local alignment method (e.g. MEME). Compare alignments produced by different programs.

Conserved block
 Non-conserved region

Figure 4 Choice of multiple alignment methods according of the nature of the sequence set.

©2010 Sami Khuri