

Bioinformatics

Four Multiple Sequence Alignment

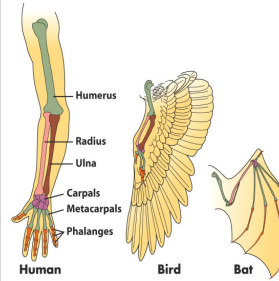
PKKLGCCDNALRAV
DCRLIIVVYQQASE
DCRLIPVVFATKNVS
DKKLLIIVYQLAVI
DCRTHIIVYAMKAMN
EGCVKQLDFHRQTTE
TGQIDQLSYAQRKAD

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

PKKLGCCDNALRAV
DCRLIIVVYQQASE
DCRLIPVVFATKNVS
DKKLLIIVYQLAVI
DCRTHIIVYAMKAMN
EGCVKQLDFHRQTTE
TGQIDQLSYAQRKAD

©2012 Sami Khuri

Multiple Sequence Alignment



- ❖ Progressive Alignment
- ❖ Guide Tree
- ❖ ClustalW
- ❖ Toffee
- ❖ Muscle
- ❖ MAFFT

©2012 Sami Khuri

Aligning BRCA1 Sequences (I)

```

Wombat : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Opossum : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Armadillo : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Sloth : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Dugong : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 74
Hyrax : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 74
Aardvark : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 74
Tenrec : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 74
Rhinoceros : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Pig : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Hedgehog : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Human : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Rat : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
Hare : AAAGTTAAGTGGTGTATCCAGAAGTATGACACTTTAGCTCTGATACTCCAGGTTAGGAGCCATGACAGACGGCAGA : 83
    
```

Part of the alignment of the DNA sequences of the BRCA1 gene

From "Bioinformatics and Molecular Evolution" by Paul Higgs and Teresa Attwood

©2012 Sami Khuri

Aligning BRCA1 Sequences (II)

```

Wombat : KVNEWLRSRSDILASDNSNGRSHEQSAEVPFSALEDGHPDTAEGNSVSEKID : 52
Opossum : KVNEWLRSNDVLPFYDYSVRSRHEQNAEATNALEYGHVET-DGNSSIASEKID : 51
Armadillo : KVNEWFRSDDILTSDDSHDRGSELNAEVAGALKV--SKEVDYSSSEKID : 50
Sloth : KVNEWFRSDDILTSDDSHDRGSELNAEVAGALKV--PNEVDYSSSEKID : 50
Dugong : KVNEWFRSDGL---DDLHDKGSENAEVAGALEV--PEEVHGYSSSEKID : 47
Hyrax : KVNEWFRSDNL---SDSPFEGSENLGKVPVKL--PGEVHRYSPFENID : 47
Aardvark : KVNEWFRSDGL---DGHDEGSENAEIGGALEV--SNEVHYSYSSSEKID : 47
Tenrec : KVNEWFRSKSHGL---GDSRDGRPESEADVAVAVEV--PDEACESYSSPEKID : 47
Rhinoceros : KVNEWFRSDELLTSDDSHDRGSENAEIVAGALEV--QNEVDYSSSEKID : 50
Pig : KVNEWFRSDEMLTSDDSDRRSENAEIVAGALEV--PNEADGHLGSSEKID : 50
Hedgehog : KVNEWFRSDELLTSDDSDRRSENAEIVAGALEV--PNEADGHLGSSEKID : 50
Human : KVNEWFRSDELLTSDDSHDRGSENAEIVAGALEV--LNEVDYSSSEKID : 50
Rat : KVNEWFRSDEMLTSDDSDRRSENAEIVAGALEV--SNEVDGCFSSSKID : 50
Hare : KVNEWFRSDEMLTSDNSDRRSENAEIVAGALEV--PKEVDYSSSEKID : 50
    
```

Alignment of BRCA1 protein sequences for the same region on the gene

From "Bioinformatics and Molecular Evolution" by Paul Higgs and Teresa Attwood

©2012 Sami Khuri

Conserved Regions in Genes in Divergent Species

- Species that are very different from one another have similar genes that generally perform identical or similar functions.
 - **Example:** Marsupial vs. Placental
- Sometimes these genes undergo mutations due to natural selection, thus altering their function.

©2012 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.

©2012 Sami Khuri

Multiple Sequence Alignment

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

©2012 Sami Khuri

Importance of MSA (I)

- If protein X with unknown function, has domains that are similar to domains of annotated proteins, then we can infer that protein X has a similar structure or function to the annotated proteins.
- A **Multiple Sequence Alignment** generally reveals more information than the analysis of a sequence by itself or even the analysis obtained from a Pairwise Sequence Alignment.

©2012 Sami Khuri

Aligning Kinases: An Example

```

p110β      SYVLGIG-----DRHSDNINVKKTGQLFHIDFGHILGNFKSKFGIKRERVFFILT
p110δ      TYVLGIG-----DRHSDNIMIRESGQLFHIDFGHFLGNFKTKFGINRERVFFILT
p110α      TFILGIG-----DRHNSNIMVKDDGQLFHIDFGHFLDHKKKFGYKRRVPPVLT
p110γ      TFVLGIG-----DRHNDNIMITETGQLFHIDFGHILGNYSFLGINKERVPPVLT
p110_dicti TYVLGIG-----DRHNDNLMVTKGGRLFHIDFGHFLGNYSKFKGKERAPFVFT
cAMP-kinase QIVLTFEYLSLDLIYRDLKPENLLIQQGVIQVTDDFGFAKRVKGRTWXLCG--TPEYLA
  
```

Multiple sequence alignment between a cAMP-kinase and 5 PI-3 kinases. Green indicates total conservation (identical residues), while blue indicates physicochemically conserved residues (belonging to the same partition of amino acids).

©2012 Sami Khuri

Pairwise vs. Multiple Alignment

```

p110α      TFILGIGDRHNSNIMVKDDG-QLFHIDFGHFLDHKKKFGYKRRVPPVLT--QDFLIVI
cAMP-kinase QIVLTFEYLSLDLIYRDLKPENLLIQQGVIQVTDDFGFAKRVKGRTWXLCGTPEYLAPE

p110β      SYVLGIG-----DRHSDNINVKKTGQLFHIDFGHILGNFKSKFGIKRERVFFILT
p110δ      TYVLGIG-----DRHSDNIMIRESGQLFHIDFGHFLGNFKTKFGINRERVFFILT
p110α      TFILGIG-----DRHNSNIMVKDDGQLFHIDFGHFLDHKKKFGYKRRVPPVLT
p110γ      TFVLGIG-----DRHNDNIMITETGQLFHIDFGHILGNYSFLGINKERVPPVLT
p110_dicti TYVLGIG-----DRHNDNLMVTKGGRLFHIDFGHFLGNYSKFKGKERAPFVFT
cAMP-kinase QIVLTFEYLSLDLIYRDLKPENLLIQQGVIQVTDDFGFAKRVKGRTWXLCG--TPEYLA
  
```

Top Figure: The pairwise alignment of the two homologous kinases does not align the important active-site residues and the DFG motif (in green).

Bottom Figure: The multiple sequence alignment of 5 homologous kinases forces the best-conserved regions to be matched.

©2012 Sami Khuri

Importance of MSA (II)

Given a group of sequences:

- Are they homologous?
 - MSA will reveal the relationship between them.
- Do they contain conserved regions?
 - Similar regions may reveal similar functions, eg. active sites.
- Can we build a family profile?
 - The profile can be used to search and fish out members of that family in databases.
- Can we build a consensus sequence?
 - The consensus sequence can be used for further analysis

©2012 Sami Khuri

Importance of MSA (III)

- MSA can help in the prediction of secondary and tertiary structures of new sequences.
- Homology Modeling:
 - MSA can be used for protein modeling programs.
- MSA's are used as input for constructing phylogenetic trees
 - Especially for distance-based algorithms such as UPGMA and Neighbor-Joining.

©2012 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.

©2012 Sami Khuri

Heuristic Algorithms

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

©2012 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.

©2012 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.

©2012 Sami Khuri

Step One of Clustal: Pairwise Alignments

1) Perform pairwise alignments of all sequences Compare each sequence with each other calculate a **distance matrix**.

A	-		
B	.87	-	
C	.59	.60	-
	A	B	C

Distance = Number of exact matches divided by the sequence length (ignoring gaps).

Distance Matrix

Note that .87 means 87% identical.

©2012 Sami Khuri

Step Two of Clustal: Create Guide Tree

2) Use the results of the Distance Matrix to create a **Guide Tree** to help determine in what order the sequences are aligned.

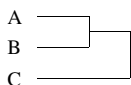


The **Guide Tree**, or Dendrogram has no phylogenetic meaning. It cannot be used to show evolutionary relationships.

©2012 Sami Khuri

Step Three of Clustal: Progressive Alignment

3) Use the Guide Tree to align the sequences



- Align A and B first
- Then add sequence C to the previous alignment

Align the most closely related sequences first, then add in the most distantly related ones and align them to the existing alignment, inserting gaps if necessary.

©2012 Sami Khuri

Multiple Alignment Problems

- Does the quality of the **guide tree** matter?
 - Not for very closely related sequences, but perhaps for distantly related ones.
- **Local minimum** problem
 - If the initial alignments have a problem, they cannot be removed during subsequent steps.

©2012 Sami Khuri

Which Comparison Table?

- **Single Parameter problem**
 - You are using one weight matrix, and one set of penalties for all the sequences.
 - The best set of parameters for one part of the alignment may not be the best for another part.
- **Do we use**
 - **BLOSUM 35** to best align the distant sequences
 - **BLOSUM 90** to align the very closely related sequences, or
 - **BLOSUM 62** as an average?

©2012 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.

©2012 Sami Khuri

Steps of ClustalW



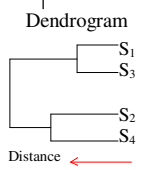
- Multiple Alignment Step:
1. Aligning S₁ and S₃
 2. Aligning S₂ and S₄
 3. Aligning (S₁,S₃) with (S₂,S₄).

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



©2012 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.



©2012 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.



©2012 Sami Khuri

When Not To Use Clustal

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

©2012 Sami Khuri

Alignment Problems

- Final result sometimes depends on the **order** that sequences were analyzed.
- **Gaps** can make alignment unrealistically long.
- Sequences of **different lengths** can cause problems.
- **Non-homologous** regions can dilute homologous areas.
 - Only need to align the shared domain.
 - So trim away any excess sequence and realign.

©2012 Sami Khuri

DNA or Protein Alignment

- If we are comparing two or more sequences, is it better to align the **DNA**, or **Protein**?

It depends on what we want to compare.

- If **protein function**, then look at the amino acids
- If **genetic changes**, then look at the DNA

- The **initial mutations** take place at the DNA level, but the **evolutionary pressure** occurs at the protein level.

©2012 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.

©2012 Sami Khuri

Alternatives to CLUSTALW (I)

- **TCoffee**: A collection of tools for Computing, Evaluating and Manipulating Multiple Alignments of DNA, RNA, Protein Sequences and Structures.
 - Good for distantly related sequences too.
 - www.tcoffee.org
- **MUSCLE**: Multiple Sequence Comparison by Log-Expectation
 - www.drive5.com/muscle

©2012 Sami Khuri

Alternatives to CLUSTALW (II)

- **MAFFT**: Multiple Alignment using Fast Fourier Transform.
 - A good balance between accuracy and speed.
 - align.genome.jp/mafft
- **PRRN**: A web-based multiple sequence alignment package.
 - align.genome.jp/prrn

©2012 Sami Khuri

Alternatives to CLUSTALW (III)

- **Praline**: Multiple sequence alignment toolkit with several strategies to optimize alignment quality.
 - Has an option for “transmembrane structure prediction”.
 - www.ibi.vu.nl/programs/pralinewww
- **Blocks**: Blocks Multiple Alignment Processor
 - Performs a local alignment (finds conserved blocks) blocks.fhrc.org/blocks/process_blocks.html

©2012 Sami Khuri

Alternatives to CLUSTALW (IV)

- **Meme**: Multiple Em for Motif Elicitation
 - Performs local multiple alignment, searching for motifs.
 - meme.sdsc.edu/meme/cgi-bin/meme.cgi
- **SAM**: Sequence Alignment and Modeling System
 - collection of flexible software tools for creating, refining, and using linear hidden Markov models for biological sequence analysis
 - compbio.soe.ucsc.edu/sam.html

©2012 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.

©2012 Sami Khuri

MSA Editor and Formatter Programs

- Multiple Sequence Alignment programs:
 - CINEMA (Color Interactive Editor for Multiple Alignments)
 - GDE (Genetic Data Environment)
 - GeneDoc
 - MACAW
- Multiple Sequence Alignment programs:
 - Boxshade
 - CLUSTALX

©2012 Sami Khuri

Appropriate Approaches

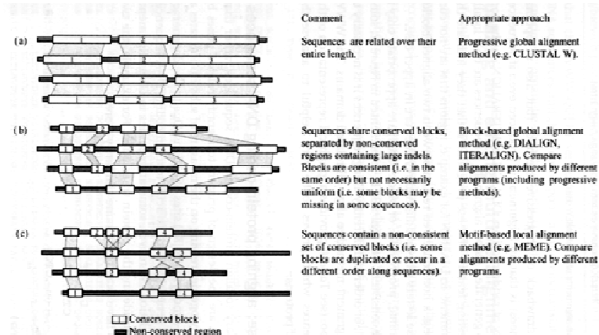


Figure 4 Choice of multiple alignment methods according of the nature of the sequence set.

©2012 Sami Khuri