

## Getting To Know Your Protein

### Comparative Protein Analysis: Part I. Phylogenetic Trees and Multiple Sequence Alignments

Robert Latek, PhD  
Sr. Bioinformatics Scientist  
Whitehead Institute for Biomedical Research

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

2

## Meeting Your Protein *In Silico*

- Define and characterize your favorite sequence
  - Identify homologous sequences
  - Predict function
  - Examine potential mutations
  - Study in 3D
  - Make manuscript figures :-)

## Comparative Protein Analysis

### Definition

Use information regarding a group of sequences to determine the function of an undefined sequence.

Extract novel information about a protein, or a series of proteins, through comparisons with other, related sequences.

### Application

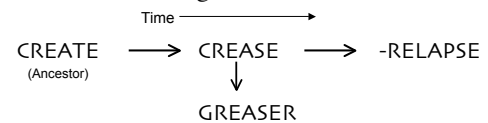
What are they?  
What are their functions?  
Why are they important?

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

3

## Comparative Protein Analysis

- Identify proteins within an organism that are related to each other and across different species
- Generate an evolutionary history of related genes
- Locate insertions, deletions, and substitutions that have occurred during evolution



WIBR Bioinformatics Courses, © Whitehead Institute, 2005

4

## Syllabus

- **Phylogenetic Trees**
- Multiple Sequence Alignments
- From Trees and MSAs to Manuscript Figures
- Exercises

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

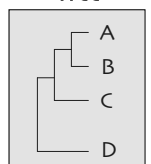
5

## Phylogenetic Trees



- A graph representing the evolutionary history of a sequence
- Relationship of one sequence to other sequences
- Dissect the order of appearance of insertions, deletions, and mutations
- Predict function, observe epidemiology, analyzing changes in viral strains
- Tree of Life  
<http://tolweb.org/tree/phylogeny.html>

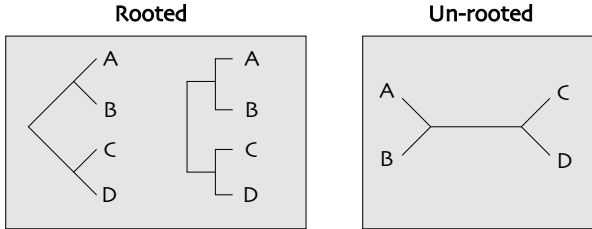
### Simple Tree



WIBR Bioinformatics Courses, © Whitehead Institute, 2005

6

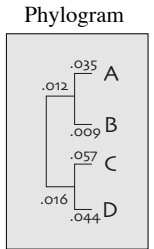
# Tree Shapes



Branches intersect at Nodes  
Leaves are the topmost branches

# Tree Characteristics

- Tree Properties**
  - Clade:** all the descendants of a common ancestor represented by a node
  - Distance:** number of changes that have taken place along a branch
- Tree Types**
  - Cladogram:** shows the branching order of nodes
  - Phylogram:** shows branching order and distances

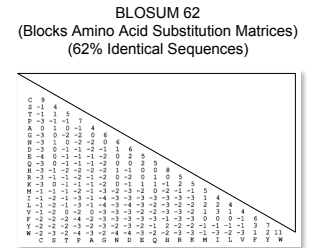


# Tree Building Algorithms

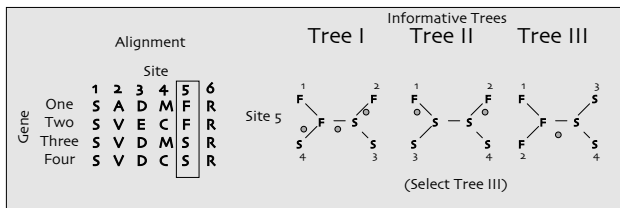
- Maximum Parsimony
- Distance Methods
  - UPGMA
  - Neighbor Joining
- Maximum Likelihood

# Amino Acid Substitutions

- Substitute one aa with another based on:
  - Chemical property
  - Size
  - Evolution (matrices)



# Maximum Parsimony



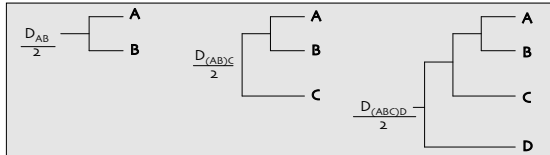
- Find the tree that changes one sequence into all of the others by the least number of steps [Focus solely on end product sequences, ignore evolutionary history]
- Only informative sites are analyzed (no gaps or conserved positions)
- Can be misleading when rates of change vary in different tree branches

# Distance Methods

- Distance** is expressed as the fraction of sites that differ between two sequences in an alignment
- Sequences with the smallest number of changes (shortest distance) are “related taxa”

## Distance Methods - UPGMA

- **UPGMA** (Unweighted Pair-Group Method with Arithmetic mean)
  - Sequentially find pair of taxa with smallest distance between them, and define branching as midpoint of two
  - Assumes the tree is additive and that rate of change is constant in all of the branches

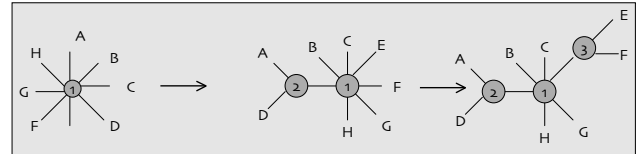


WIBR Bioinformatics Courses, © Whitehead Institute, 2005

13

## Distance Methods - NJ

- **Neighbor-Joining (NJ)**: useful when there are different rates of evolution within a tree
  - Each possible pair-wise alignment is examined. Calculate distance from each sequence to every other sequence
  - Choose the pair with the lowest distance value and join them to produce the minimal length tree
  - Update distance matrix where joined node is substituted for two original taxa and then repeat process



WIBR Bioinformatics Courses, © Whitehead Institute, 2005

14

## Maximum Likelihood

- Best accounts for variation in sequences
- Establish a **probabilistic model** with multiple solutions and determine which is most likely
- All possible trees are considered, therefore, only suitable for small number of sequences
  - Maximizes probability of finding optimal tree

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

15

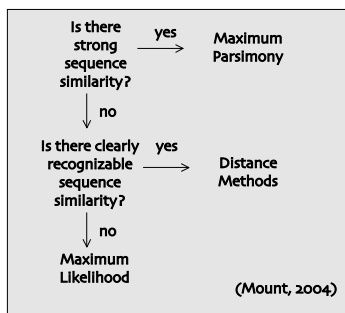
## Tree Reliability

- Probability that the members of a clade are always members of that clade
- Sample by **Bootstrapping**
  - Random sites of an alignment are randomly sampled so as to create a dataset the same size as the original. The same analysis as applied to the original data set is performed on the bootstrap dataset
  - Construct a consensus bootstrap tree and compare to the original tree

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

16

## Which Method to Use?



WIBR Bioinformatics Courses, © Whitehead Institute, 2005

17

## Syllabus

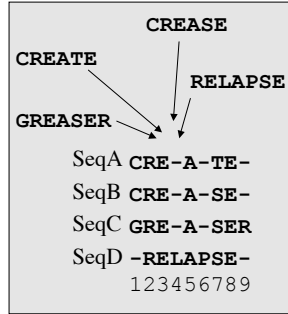
- Phylogenetic Trees
- **Multiple Sequence Alignments**
- From Trees and MSAs to Manuscript Figures
- Exercises

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

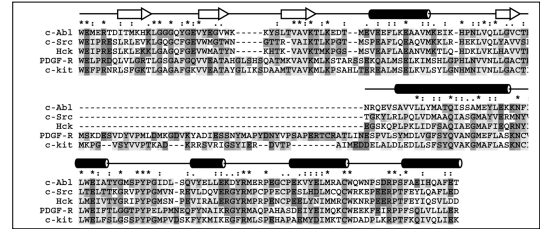
18

# Multiple Sequence Alignments

- Place residues in columns that are derived from a common ancestral residue
- MSA can reveal sequence patterns
  - Demonstration of homology between >2 sequences
  - Identification of functionally important sites
  - Protein function prediction
  - Structure prediction
  - Search for weak but significant similarities in databases
  - Design PCR primers for related gene identification
  - Genome sequencing: contig assembly



# Multiple Sequence Alignment



# Multiple Sequence Alignment

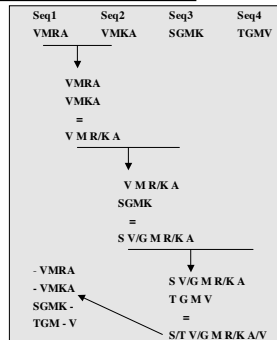
E.nidulans	IFETVEKESIS---	ENVQREWFPLINPGRIPALDITFPEQCKILFE	73
A.nidulans	VFNIRHFFPD---	IVYQREFFIN--INNGVIVATD---NTDLIWE	73
B.fuockliana	LSYEVKIDIS---	KNTQKWFPLINPGRIPALDITFPEQCKILFE	74
F.graminearum	LKFKVTFEFS---	KHREKWFPLINPGRIPALDITFPEQCKILFE	74
M.grisea	LPHITTFPHPT---	SLKQSEVLEKVNPGMSPALIDP---NTDLIWE	71
M.grisea2	LPHITTFPHPT---	SLKQSEVLEKVNPGMSPALIDP---NTDLIWE	71
N.crassa	YFDLNDICIS---	QAKSEFFVQVNGRLEIACIDP---NTDLIWE	73
Y.lipolytica	LPFNTIFLFPN---	HGRQAPFVTVINPVARVIALIDP--FNNTSIWE	127
C.albicans	LPFNTIFLFPN---	NSGRQTFPFTVINPVARVIALIDP--FNNTSIWE	170
C.glabrata	LQNTVFLEPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	181
C.maltosa	LPFNTIFLFPN---	NSGRQAPFVTVINPVARVIALIDP--FNNTSIWE	154
E.gossypii	LPFNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	180
S.lactis	MPNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	215
K.marxianus	MPNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	230
K.marxianus2	MPNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	224
S.bayanus	PHNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	171
S.cerevisiae	PHNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--FNNTSIWE	180
S.mikatae	PHNTIFLFPN---	-----GMNLSIWE	162
S.paradoxus	PHNTIFLFPN---	LGRHAPFVTVINPVARVIALIDP--GMNLSIWE	185
S.pombe	LSYEQIFVPEPK---	HPQKCKRILA--LNPGRVPELVVVK--NNQYTIWE	70
C.cinereus	QNFVFFSALIM---	LAGVYVDFPFFSISDQDFQ---MGLQVVA	66
U.maydis	ISYDIFLPHDDSD---	EKGVSAKFLINPGRVFCVSN--DSKRFVVM	71
D.cerio	LNWELRQFFP---	GLQSPFLAIPARVAVLDG---DLKLS	84
X.laevis	LEKFAASAKPFE---	NSGRQAPFVTVINPVARVIALIDP--FNNTSIWE	85
D.melanogaster	LEFNKTIINLK---	SGQMPFIFINPQHSIFELVEN---SFTIWE	48
C.elegans	VDFEYVWELLS---	SEAKS--KAKIIEAARVFFVD---SOTIWE	68
C.elegans2	IDYEVPEIDLFS---	SESKNAPVKNPAKVFVLIN---SLELTI	68
Z.maise	LDPEIVPVVLT---	TAHQKQPFALNPFQIPALVDG---DEVLPF	67
T.sectivum	AESEIVMPEV---	ASGRHFCQVCLAPFAMSPQDS---DAVLPF	67
A.thaliana	VAFETIPVDM---	KGRHQPVALQPFQFVAVVDG---DEXIFP	65
O.sativa	AESEIVPLDPS---	KGRHAPFVTVINPVARVIALIDP--FNNTSIWE	67

# Approaches

- Optimal Global Alignments** - Dynamic programming
  - Build matrices with every possible combination and search for optimal solution
  - Align 10 sequences of 100 aa length = 100<sup>10</sup> possibilities
  - Optimal in the mathematical sense
- Global Progressive Alignments** - Match most common sequences together
- Global Iterative Alignments** - Multiple re-building attempts to find best alignment
- Local alignments**
  - Profiles, Blocks, Patterns

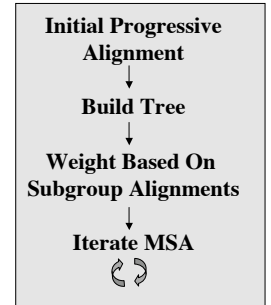
# Global Progressive Alignment

- A heuristic approach that utilizes phylogenetic information to assist in routing the alignment (clustalw/clustalx)
- Feng & Doolittle 1987, Higgins and Sharp 1988
- Most alike sequences are aligned together in order of their similarity (tree-based), a consensus is determined and then aligned to next most similar sequence



# Iterative Multiple Alignment

- “Repeatedly re-align subgroups of sequences into a global alignment to improve alignment score” (Mount, 2004)
- Start with a progressive alignment and tree
- Recalculate pair-wise scores during progressive alignment, use new scores to rebuild the tree, which is used to improve alignments



# MSA and Tree Relationship

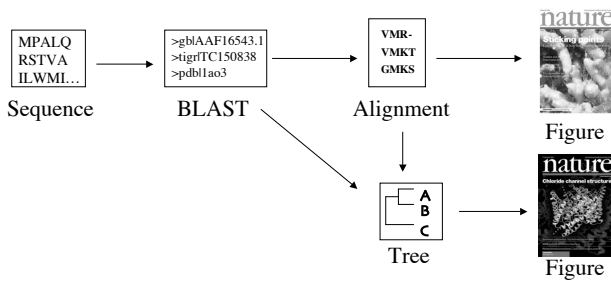
- “The optimal alignment of several sequences can be thought of as minimizing the number of mutational steps in an evolutionary tree for which the sequences are the leaves” (Mount, 2001)



# Syllabus

- Phylogenetic Trees
- Multiple Sequence Alignments
- From Trees and MSAs to Manuscript Figures
- Exercises

## Manuscript Figures 101



## 1. Find Related Sequences

- BLAST

– [www.ncbi.nih.gov/BLAST](http://www.ncbi.nih.gov/BLAST)

```
MLEICLKLVGCKSKKGLSSSSCYLEEALQRPVASFPEQGLSEARWNSKENLLAGPSENDPNLFVALY
DFVASGNTLSITKGEKLRVLGYNHNGEWCEAQTNGQGWVPSNYITPVNSLEKHSWYHGPVSRNAEYL
LSSGINGSFLVRESESPGQRSISLRYEGRVYHRYINTASDGKLYSSESRFNTLAEVLHHSTVADGLI
TTLHYAPAKRNKPTVYGVSPNYDKWEMERTDITMKHLGGQYGEVGEWKKYSLTVAVKTKEDTMEV
EEFLKEAAWKEIKHPNLVQLLGVCTREPPFYITFMTYGNLLDYLRECNRQEVNAWLLYMATQISSA
MEYLEKKNFHIRDLAARNCLVGENHLVKVADFGLSRLMTGDYTAHAGAKFFKWTAPESLAYNKFSIKS
DVAWFGVLLWEIATYQMSPPYDIDLSDVYELLEKDYRMEPEGGQPEKYEELRACWQVWVNSDRSFAEIH
QAFETMFOESSISDEVEKELGKQGVRAVSTLLOAPELPTKTRTSRRAEHRDITDVPMPHSKGGQGED
PLDHEPAVSPPLLPRKERGPPPEGLNEDERLLPKDKKTNLFSALIKKKKTAFTPPKRSSSFREMDGQPER
RGAGEEERDISNGALFPLDADPAKSPKPSNGAGVPGALRESGSGSFRPHLWKKSSSTLTSRLAT
GEEEGGSSSKRFLRSCSASCVPHGAKDTEWRSVTLPRDLQSTGROFDSSTFGGHKSEKALPRKRAGEN
RSDQVTRGTVTPPPRLVKNEEADEVFKDIMESSPPNPNTPKPLRRQVTVAPASGLPHKEAGKGS
ALGTPAAEVPVPTTSKAGSGAPGTSKGPAAEESRVRHRKHSSSESPGRDKGLSRLKPAFPVPPAASAGKA
GGKPSOSPSOEAAGEAVLGAKTATSLVDANNSDAAKPSQPEGGLKPKLPAATPKQSAKPSGTPISPAP
VPSTLPSASSALAGDQPSSTAFIPLISTRVLRKTRQPPERIASGATKGVLDSTEALCLAIRNSEQM
ASHSAVLEAGKNLYTFCVSYVDSIQQMRNKFAPREANKLENNLRELQICPATAGSGPAATQDFSKLLSS
VKEISDVQR
```

## 2. Compile & Align Sequences

- Compile Sequences into FASTA format

```
>Human
MPALGYKFSTW...
>Mouse
MDGSTDYGILQINS...
>Rat
MKKP...
>Murine_Leukemia_Virus
MTSR....
```

- Align

- PC: [www-igbmc.u-strasbg.fr/BioInfo/ClustalX/Top.html](http://www-igbmc.u-strasbg.fr/BioInfo/ClustalX/Top.html)
- OS X: [www.embl.de/~chenna/clustal/darwin/](http://www.embl.de/~chenna/clustal/darwin/)
- Web: [pir.georgetown.edu/pirwww/search/multaln.html](http://pir.georgetown.edu/pirwww/search/multaln.html)
- Jalview: [http://www.jalview.org/Web\\_Installers/install.htm](http://www.jalview.org/Web_Installers/install.htm)

## 3. Build Tree

- Create tree

– Clustalx Neighbor Joining method

- Draw tree

– TreeView:

- [taxonomy.zoology.gla.ac.uk/rod/treeview.html](http://taxonomy.zoology.gla.ac.uk/rod/treeview.html)

– Web:

- [iubio.bio.indiana.edu/treeapp/treeprint-form.html](http://iubio.bio.indiana.edu/treeapp/treeprint-form.html)

## 4. Create Figures

---

- MSAs are often multipage
  - **Convert** to PDF with **Acrobat Distiller** or open with **Ghostview** (<http://www.cs.wisc.edu/~ghost/> or <http://www.kiffe.com/macghostview.html>)
  - Extract pages individually and save as separate PDF/PS files
  - **Open** images in favorite illustration program
  - **Export** annotated alignments/trees to Powerpoint
- **Publish** paper, give award-winning presentation!

## Exercise I

---

- **BLAST your sequence**
  - NCBI BLAST
  - Collate and edit sequences in a text editor
- **Perform multiple sequence alignment**
  - Clustalx
- **Build Phylogenetic Tree**
  - Clustalx and TreeView
- **Manage Postscript Files**
  - Adobe Acrobat Distiller/Ghostview
- **Create Figure**
  - Illustrator - > Powerpoint

## References

---

- Bioinformatics: Sequence and genome Analysis. David W. Mount. CSHL Press, 2001 and 2004.
- Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins. Andreas D. Baxevanis and B.F. Francis Ouellete. Wiley Interscience, 2001.
- Bioinformatics: Sequence, structure, and databanks. Des Higgins and Willie Taylor. Oxford University Press, 2000.