NATIONAL HUMAN GENOME RESEARCH INSTITUTE   Division of Intramural Research
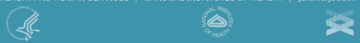
*Current Topics in Genome Analysis*
*Spring 2008*

*Week 3: Biological Sequence Analysis II*

*Andy Baxevanis, Ph.D.*

U.S. DEPARTMENT OF HEALTH AND HUMAN SERVICES  |  NATIONAL INSTITUTES OF HEALTH  |  genome.gov/DIR

## Overview

- Week 2
    - Similarity *vs.* Homology
    - Global *vs.* Local Alignments
    - Scoring Matrices
    - BLAST
    - BLAT

- Week 3
    - Profiles, Patterns, Motifs, and Domains
    - Structures: VAST, Cn3D, and *de novo* Prediction
    - Multiple Sequence Alignment

## Sequence Comparisons

- Homology searches
  - Usually "one-against-one"          *BLAST, FASTA*
  - Allows for comparison of individual sequences against databases comprised of individual sequences

- Profile searches
  - Uses collective characteristics of a family of proteins
  - Search can be "one-against-many"          *Pfam, InterPro, CDD*

    or "many-against-one"          *PSI-BLAST*

## Profiles

- Numerical representations of multiple sequence alignments
- Depend upon *patterns* or *motifs* containing conserved residues
- Represent the common characteristics of a protein family
- Can find similarities between sequences with little or no sequence identity
- Allow for the analysis of distantly-related proteins

## Profile Construction

```
APHIIVATPG
GCEIVIATPG
GVEICIATPG
GVDILIGTTG
RPHIIVATPG
KPHIIIATPG
KVQLIIATPG
RPDIVIATPG
APHIIVGTPG
APHIIVGTPG
GCHVVIATPG
NQDIVVATTG
```

• *Which residues are seen at each position?*
• *What is the frequency of observed residues?*
• *Which positions are conserved?*
• *Where can gaps be introduced?*

*Position-Specific Scoring Table*

| Cons | A | B | C | D | E | F | G | H | I | K | L | M | N | P | Q | R | S | T | V | W | Y | Z |
|------|----|----|----|----|----|----|-----|----|----|----|----|----|----|----|----|----|----|-----|----|-----|-----|----|
| G | 17 | 18 | 0 | 19 | 14 | -22 | 31 | 0 | -9 | 12 | -15 | -5 | 15 | 10 | 9 | 6 | 18 | 14 | 1 | -15 | -22 | 11 |
| P | 18 | 8 | 13 | 8 | 8 | -12 | 13 | 8 | 8 | 8 | 8 | -1 | | 23 | 2 | -2 | 12 | 11 | 17 | -31 | -8 | 1 |
| H | 5 | 24 | -12 | 29 | 25 | -20 | 8 | 32 | -9 | 9 | -10 | -9 | 22 | 7 | 30 | 10 | 0 | 4 | -8 | -20 | -7 | 27 |
| I | -1 | -12 | 6 | -13 | -11 | 33 | -12 | -13 | 63 | -11 | 40 | 29 | -15 | -9 | -14 | -15 | -6 | 7 | 50 | -17 | 8 | -11 |
| V | 3 | -11 | 1 | -11 | -9 | 22 | -3 | -11 | 46 | -9 | 37 | 30 | -13 | -3 | -9 | -13 | -6 | 6 | 50 | -19 | 2 | -8 |
| V | 5 | -9 | 9 | -9 | 19 | -1 | -13 | 57 | -9 | 35 | 26 | -13 | -2 | -11 | -13 | -4 | 9 | 58 | -29 | 0 | -9 | |
| A | 54 | 15 | 12 | 20 | 17 | -24 | 44 | -6 | -4 | -1 | -11 | -5 | 12 | 19 | 9 | -13 | 21 | 19 | 9 | -39 | -20 | 10 |
| T | 40 | 20 | 20 | 20 | 20 | -30 | 40 | -10 | 20 | 20 | -10 | 0 | 20 | 30 | -10 | -10 | 30 | 150 | 20 | -60 | -30 | 10 |
| P | 31 | 6 | 7 | 6 | 6 | 11 | 10 | 11 | 2 | 6 | 16 | 11 | | 89 | 17 | 17 | 24 | 22 | 9 | -50 | -48 | 12 |
| G | 70 | 60 | 20 | 70 | 50 | 5 | 150 | -20 | -30 | -10 | -50 | -30 | 40 | 30 | 20 | -30 | 60 | 40 | 20 | -100 | -70 | 30 |

## Patterns

| Phe *or* Tyr | Cys | | *not* Val *or* Ala | three His |

$$[FY]-x-C-x(2)-\{VA\}-x-H(3)$$

| any amino acid | any two amino acids | | any amino acid | |

## Pfam

- Collection of multiple alignments of protein domains and conserved protein regions (regions which probably have structural or functional importance)

- Each Pfam entry contains:
  - Multiple sequence alignment of family members
  - Protein domain architectures
  - Species distribution of family members
  - Information on known protein structures
  - Links to other protein family databases

## Pfam

- Pfam A
  - Based on *curated* multiple alignments ("seed alignment")
  - Hidden Markov models (HMMs) used to find all detectable protein sequences belonging to the family
  - Given the method used to construct the alignments, hits are highly likely to be true positives

- Pfam B
  - Automatically generated from database searches
  - Deemed "lower quality", but can be useful when no Pfam A family is identified

[FW]-[SGNH]-x-[GD]-{F}-[RKHPT]-{P}-C-[LIVMFAP]-[GAD]

**Parent-Child Relationships (Subfamilies)**

*Child entries are more specific than the parent*
*A match to the child entry implies a match to the parent*
*Signatures for the parent and child entries must overlap*



| | |
|---|---|
| Center | Tree root |
| Inner circles | Tree nodes |
| Outer circles | Representative model organisms |

*There is no significance to the placement of individual nodes on the circles*

# Conserved Domain Database (CDD)

- Identify conserved domains in a protein sequence

- "Secondary database"
  - Pfam A and B
  - Simple Modular Architecture Research Tool (SMART)
  - Clusters of Orthologous Groups

- Search performed using RPS-BLAST
  - Query sequence is used to search a database of precalculated position-specific scoring tables
  - *Not* the same method used by Pfam or InterPro



http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml

# PSI-BLAST

- Position-Specific Iterated BLAST search
- Easy-to-use version of a profile-based search
  - Perform BLAST search against protein database
  - Use results to calculate a position-specific scoring matrix
  - PSSM replaces query for next round of searches
  - May be iterated until no new significant alignments are found
    - Convergence – all related sequences deemed found
    - Divergence – query is too broad, make cutoffs more stringent



*http://www.ncbi.nlm.nih.gov/BLAST*

## Overview

- Week 2
  - Similarity *vs.* Homology
  - Global *vs.* Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 3
  - Profiles, Patterns, Motifs, and Domains
  - Structures: VAST, Cn3D, and *de novo* Prediction
  - Multiple Sequence Alignment

## Predicting Tertiary Structure

- Sequence specifies conformation, *but* conformation does *not* specify sequence

- Structure is conserved to a much greater extent than sequence

- Similarities between proteins may not necessarily be detected through "traditional" methods

VAST Structure Comparison

*Step 1:* Construct vectors for secondary structure elements



VAST Structure Comparison

*Step 2:* Optimally align structure element vectors

*Protein 1*                *Protein 2*

*Alignment 1*     *Alignment 2*     *Alignment 3*     *Alignment 4*

## VAST Shortcomings

- Not the best method for determining structural similarities

- Reducing a structure to a series of vectors necessarily results in a loss of information (less confidence in prediction)

- Regardless of the "simplicity" of the method, provides a simple and fast first answer to the question of structural similarity

| Worms | **Rendering** | Spacefill |
| Secondary Structure | **Coloring** | Charge |



# Current Protocols in Bioinformatics

*CPBI Unit 1.3*  
*Entrez and Cn3D*

*http://nihlibrary.nih.gov*  
*Search "Online Journals" for "Current Protocols in Bioinformatics"*

# SWISS-MODEL

- Automated comparative protein modelling server
- Web front-end at *http://www.expasy.org/swissmod*
- Results returned by E-mail

BLAST search to find similarities in PDB *by sequence*

Select templates with sequence identity > 25% and projected model size > 20 amino acids

Generate models

Do energy minimization

Generate PDB file for new protein model



```
21DJH.pdb: 42.77 % identity
21DJG.pdb: 42.77 % identity
11DJG.pdb: 42.22 % identity
11QAS.pdb: 44.17 % identity
11QAT.pdb: 43.52 % identity
21QAT.pdb: 43.52 % identity
21QAS.pdb: 43.52 % identity

Target:     |=============================================|
21DJH.pdb   |           ——
21DJG.pdb   |           ——
11DJG.pdb   |           ——
11QAS.pdb   |           ——
11QAT.pdb   |           ——
21QAT.pdb   |           ——
21QAS.pdb   |           ——
```

```
ATOM    1   H1   SER   1    24.219  22.954
ATOM    2   H2   SER   1    24.770  21.435
ATOM    3   N    SER   1    24.355  22.187
ATOM    4   H3   SER   1    23.466  21.925
ATOM    5   CA   SER   1    25.266  22.675
ATOM    6   CB   SER   1    24.826  24.072
ATOM    7   OG   SER   1    24.857  25.006
ATOM    8   HG   SER   1    24.717  25.929 -55.233  1.00 99.00
ATOM    9   C    SER   1    25.471  21.750 -53.751  1.00 25.00
ATOM   10   O    SER   1    25.923  22.169 -52.684  1.00 25.00
ATOM   11   N    LYS   2    25.227  20.460 -53.972  1.00 25.00
ATOM   12   H    LYS   2    24.961  20.142 -54.878  1.00 99.00
ATOM   13   CA   LYS   2    25.366  19.408 -52.943  1.00 25.00
ATOM   14   CB   LYS   2    24.003  18.772 -52.622  1.00 25.00
```

## Structural Modeling Software

- Modeller  *http://www.salilab.org/modeller/*

- DeepView  *http://us.expasy.org/spdbv/*

- WHAT IF  *http://swift.cmbi.kun.nl*

## Current Topics in Genome Analysis

Week 14  
Tuesday, April 15, 2008

**Protein Structure Analysis and  
Protein-Protein Interactions**

*David Wishart, Ph.D.*  
*Departments of Computing Science and  
Biological Sciences  
University of Alberta*

## Overview

- Week 2
  - Similarity *vs*. Homology
  - Global *vs*. Local Alignments
  - Scoring Matrices
  - BLAST
  - BLAT

- Week 3
  - Profiles, Patterns, Motifs, and Domains
  - Structures: VAST, Cn3D, and *de novo* Prediction
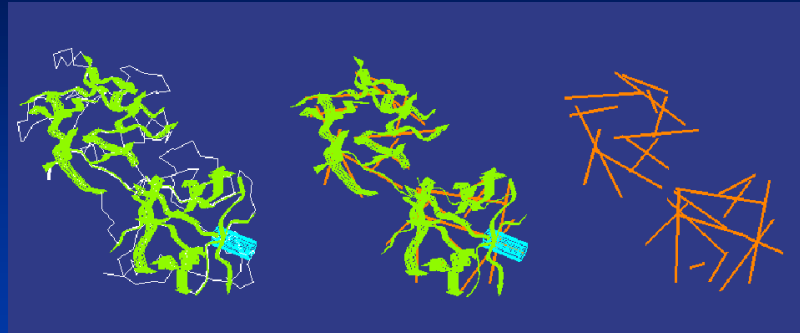  - Multiple Sequence Alignment

## Why do multiple sequence alignments?

- Identify conserved regions, patterns, and domains
  - Experimental design
  - Predicting structure and function
  - Identifying new members of protein families
- Perform phylogenetic analysis
- Generate position-specific scoring matrices for subsequent searches ("many-against-one" or "one against many")
- Bolster confidence in secondary structure predictions

## Considerations

- Absolute sequence similarity
  *Create the alignment by lining up as many common characters as possible*

- Conservation
  *Take into account residues that can substitute for one another and not adversely affect the function of the protein*

- Structural similarity
  *Knowledge of the secondary or tertiary structure of the proteins being aligned can be used to fine-tune the alignment*

## General Guidelines

- As with most analyses, concentrate on the protein level rather than on the nucleotide level
  - More informative
  - Less prone to inaccurate alignment ("20 *vs.* 4")
  - Can "translate back" to nucleotide sequences *after* doing the alignment

## General Guidelines

- Use a reasonable number of sequences to avoid technical difficulties
  - *Global* alignment method: compute time increases exponentially as sequences are added to the set
  - Most alignment algorithms are ineffective on huge data sets (and may yield inaccurate alignments)
  - Phylogenetic studies resulting from inordinately large data sets are almost impossible
  - Good starting point: 10-15 sequences
  - Ballpark upper limit: 50 sequences

## General Guidelines

- Selecting sequences for alignment
  - Sequences should be of about the same length
  - Use closely-related sequences to determine "required" amino acids
  - Use more divergent sequences to study evolutionary relationships
  - Good starting point: use sequences that are 30-70% similar to most of the other sequences in the data set
  - The most informative alignments result when the sequences in the data set are not "too similar", but also not "too different"

## General Guidelines

- Iterative process
  - Perform alignment on small set of sequences
  - Examine the quality of the alignment
  - If alignment good, can add new sequences to data set, then realign
  - If alignment not good, remove any sequences that result in the inclusion of long gaps, then realign

## Interpretation

- Absolutely-conserved positions are *required* for proper structure and function

- Relatively well-conserved positions are able to tolerate limited amounts of change and not adversely affect the structure or function of the protein

- Non-conserved positions may "mutate freely," and these mutations can possibly give rise to proteins with new functions

# Interpretation

- Gap-free blocks probably correspond to regions of secondary structure

- Gap-rich blocks probably correspond to unstructured or loop regions

# ClustalW2

- Automatic multiple alignment of nucleotide or amino acid sequences

- Implementations
  - Client versions
    *command-line text menu system, all platforms*
  - Web-based version
    *http://www.ebi.ac.uk/clustalw2*

## Progressive Alignment

- Align two sequences at a time

- Gradually build up the multiple sequence alignment by merging larger and larger sub-alignments, clustering on the basis of similarity

- Uses protein scoring matrices and gap penalties to calculate alignments having the best score

- Major advantages of method
  - Very fast
  - Alignments generally of high quality

## Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLST
>sequence B
VQLSGEEKAAVLALWDKVNEEEVGGEALGRLLVVYPWTQRFFDSFGDSLN
>sequence C
VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSH
>sequence D
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSH
```

# Progressive Alignment

1. Calculate a similarity score (percent identity) between every pair of sequences to drive the alignment

   For $N$ sequences, this requires the calculation of $[N \times (N - 1)] / 2$ pairwise alignments

   | Sequences | Alignments |
   |-----------|------------|
   | 4         | 6          |
   | 10        | 45         |
   | 25        | 300        |
   | 50        | 1,225      |
   | 100       | 4,950      |

# Progressive Alignment

```
>sequence A
VHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLST
>sequence B
VQLSGEEKAAVLALWDKVNEEEVGGEALGRLLVVYPWTQRFFDSFGDSLN
>sequence C
VLSPADKTNVKAAWGKVGAHAGEYGAEALERMFLSFPTTKTYFPHFDLSH
>sequence D
VLSAADKTNVKAAWSKVGGHAGEYGAEALERMFLGFPTTKTYFPHFDLSH
```
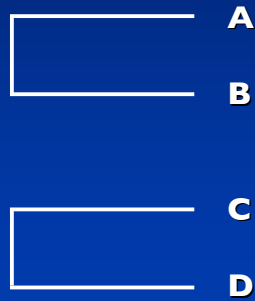
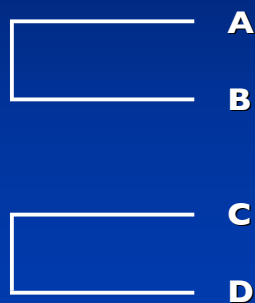| %ID | A   | B   | C  | D   |
|-----|-----|-----|----|-----|
| A   | 100 |     |    |     |
| B   | 80  | 100 |    |     |
| C   | 44  | 40  | 100 |    |
| D   | 40  | 40  | 92 | 100 |

## Progressive Alignment

2. Derive a dendrogram (guide tree) based on the pairwise comparisons (.dnd file)

    Can infer from tree that A and B share greater similarity with each other than with C or D
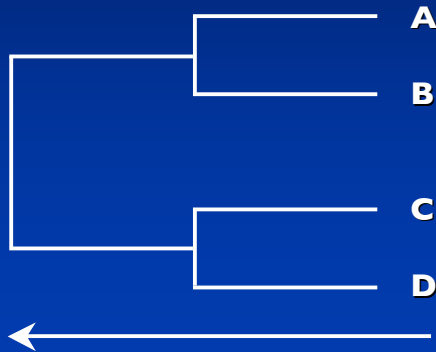
    A
    B

    C
    D

## Progressive Alignment

3. Align A with B → alignment AB (fixed)
4. Align C with D → alignment CD (fixed)
5. Represent alignments AB and CD as *single sequences*

    A
    B

    C
    D

## Progressive Alignment

6. Align "sequence" AB with "sequence" CD

7. Continue following the branching order of the tree, from the tips to the root, merging each new pair of "sequences"



## Progressive Alignment: Advantages

- Do "easier" alignments between highly-related sequences first

- Use information regarding conservation at each position to help with more difficult alignments between more distantly-related sequences later on in process

## Progressive Alignment: Disadvantages

- If initial alignments are made on distantly related sequences, there may be errors in the initial alignments

- Once an alignment is "fixed", it is not reconsidered, so any errors in the early alignments may propagate through subsequent alignments

- New version of ClustalW2 does provide a "remove first" iteration scheme to attempt to improve alignments

## ClustalW2 Output

- Pairwise scores

- Multiple sequence alignment (`.aln`)
  - Alternative formats available: GCG, Phylip, PIR, GDE

# ClustalW2 Output

- Cladogram
  - Tree assumed to be an estimate of a phylogeny
  - Branches are of equal length
  - Cladograms show common ancestry, but do not provide an indication of the amount of "evolutionary time" separating taxa

- Phylogram
  - Tree that is assumed to be an estimate of phylogeny
  - Branch lengths proportional to the amount of inferred evolutionary change

# ClustalW2 Conservation Patterns

- Conservation patterns in multiple sequence alignments usually follow the following rules:

  | | |
  |---|---|
  | [WYF] | Aromatics |
  | [KRH] | Basic side chains (+) |
  | [DE] | Acidic side chains (−) |
  | | |
  | [GP] | Ends of helices |
  | [HS] | Catalytic sites |
  | [C] | Cysteine cross-bridges |

## ClustalW2 Conservation Patterns

- Interpretation is *empirical* — there is no parallel to the *E*-values seen in BLAST searches to assess "significance"

  \*     entirely conserved column  
          (want in at least 10% of positions)

  :     "conserved"  
          (according to color table)

  .     "semi-conserved"

## ClustalW Colors

| | | |
|---|---|---|
| AVFPMILW | **Red** | Small |
| DE | **Blue** | Acidic |
| RK | **Magenta** | Basic |
| STYHCNGQ | **Green** | |

# Jalview

- Java applet available within ClustalW2 results
- Used to manually edit ClustalW2 alignments
- Color residues based on various properties
- Pairwise alignment of selected sequences
- Consensus sequence calculations
- Removal of redundant sequences
- Calculation of phylogenetic trees
- Color PostScript output

*Default view*

**Conservation**    *Conservation of total alignment (indication of percent identity)*

**Quality**    *Alignment quality, based on BLOSUM62 scores*

**Consensus**    *Based on percent identity*



*Colour → Percentage Identity*

| Agreement | Background Color |
|-----------|------------------|
| 81 - 100% | Dark blue |
| 61 - 80% | Medium blue |
| 41 - 60% | Light blue |
| ≤ 40% | White |

# A User's Guide to the Human Genome II

*http://www.nature.com/
ng/supplements/*

*Commentary:
Keeping Biology in Mind*



# Current Topics in Genome Analysis

### Next Lecture:

## Mining Data from Genome Browsers

*Tyra Wolfsberg, Ph.D.*
*National Human Genome Research Institute*
*National Institutes of Health*