*Current Topics in Genome Analysis*
*Spring 2010*

*Week 3: Biological Sequence Analysis II*

*Andy Baxevanis, Ph.D.*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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 | 10 | 0 | 15 | 0 | 0 | -12 | 15 | 0 | 8 | -5 | -5 | -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 | -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 | 0 | 7 | 0 | 0 | -11 | 19 | 11 | -9 | 0 | -10 | -11 | | 89 | 17 | 17 | 24 | 22 | 9 | -50 | -48 | 12 |
| G | 70 | 00 | 20 | 70 | 30 | -30 | 150 | -20 | -30 | -10 | -50 | -30 | 40 | 30 | 20 | -30 | 60 | 40 | 20 | -100 | -70 | 30 |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

---

# 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 |
|:---:|:---:|:---:|

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

## Family: *p450* (PF00067)

152 architectures  18883 sequences  2 interactions  1392 species  516 structures

### Summary

**Cytochrome P450**  Add annotation

Cytochrome P450s are haem-thiolate proteins [6] involved in the oxidative degradation of various compounds. They are particularly well known for their role in the degradation of environmental toxins and mutagens. They can be divided into 4 classes, according to the method by which electrons from NAD(P)H are delivered to the catalytic site. Sequence conservation is relatively low within the family - there are only 3 absolutely conserved residues - but their general topography and structural fold are highly conserved. The conserved core is composed of a coil termed the 'meander', a four-helix bundle, helices J and K, and two sets of beta-sheets. These constitute the haem-binding loop (with an absolutely conserved cysteine that serves as the 5th ligand for the haem iron), the proton-transfer groove and the absolutely conserved EXXR motif in helix K. While prokaryotic P450s are soluble proteins, most eukaryotic P450s are associated with microsomal membranes. their general enzymatic function is to catalyse regiospecific and stereospecific oxidation of non-activated hydrocarbons at physiological temperatures [6].

**Literature references**

1. Graham-Lorence S, Amarneh B, White RE, Peterson JA, Simpson ER; , Protein Sci 1995;4:1065-1080.: A three-dimensional model of aromatase cytochrome P450. PUBMED:7549871

2. Degtyarenko KN, Archakov AI; , FEBS Lett 1993;332:1-8.: Molecular evolution of P450 superfamily and P450-containing monooxygenase systems. PUBMED:8405421

3. Nelson DR, Kamataki T, Waxman DJ, Guengerich FP, Estabrook RW, Feyereisen R, Gonzalez FJ, Coon MJ, Gunsalus IC, Gotoh O, et al; , DNA Cell Biol 1993;12:1-51.: The P450 superfamily: update on new sequences, gene mapping, accession numbers, early trivial names of enzymes, and nomenclature. PUBMED:7678494

4. Guengerich FP; , J Biol Chem 1991;266:10019-10022.: Reactions and significance of cytochrome P-450 enzymes. PUBMED:2037557

5. Nebert DW, Gonzalez FJ; , Annu Rev Biochem 1987;56:945-993.: P450 genes: structure, evolution, and regulation. PUBMED:3304150

6. Werck-Reichhart D, Feyereisen R; , Genome Biol 2000;1:REVIEWS3003.: Cytochromes P450: a success story. PUBMED:11178272

**InterPro entry IPR001128**

Cytochrome P450 enzymes are a superfamily of haem-containing mono-oxygenases that are found in all kingdoms of life, and which show extraordinary diversity in their reaction chemistry. In mammals, these proteins are found primarily in microsomes of hepatocytes and other cell types, where they oxidise steroids, fatty acids and xenobiotics, and are important for the detoxification and clearance of various compounds, as well as for hormone synthesis and breakdown, cholesterol synthesis and vitamin D metabolism. In plants, these proteins are important for the biosynthesis of several compounds such as hormones, defensive compounds and fatty acids. In bacteria, they are important for several metabolic processes, such as the biosynthesis of antibiotic erythromycin in Saccharopolyspora erythraea (Streptomyces erythraeus).

**Example structure**  
PDB entry 2a1n: Crystal structure of ferrous dioxygen complex of D251N cytochrome P450cam  
View a different structure: 2a1n

---

## Family: *p450* (PF00067)

152 architectures  18883 sequences  2 interactions  1392 species  516 structures

### Domain organisation

Below is a listing of the unique domain organisations or architectures in which this domain is found. More...

**There are 16131 sequences with the following architecture: p450**  
AVNA_ASPPA [Aspergillus parasiticus] Averantin oxidoreductase EC=1.14.-.- (495 residues)  
Show all sequences with this architecture.

**There are 1087 sequences with the following architecture: p450 x 2**  
C1331_XYLFA [Xylella fastidiosa] Putative cytochrome P450 1338I EC=1.14.-.- (402 residues)  
Show all sequences with this architecture.

**There are 137 sequences with the following architecture: p450, Flavodoxin_1, FAD_binding_1, NAD_binding_1**  
CSO5_FUSOX [Fusarium oxysporum] Bifunctional P-450:NADPH-P450 reductase Cytochrome P450 505 NADPH--cytochrome P450 reductase EC=1.14.14.1 EC=1.6.2.4 (1066 residues)  
Show all sequences with this architecture.

**There are 54 sequences with the following architecture: An_peroxidase, p450**  
Q4W941_ASPFU [Aspergillus fumigatus (Sartorya fumigata)] Fatty acid oxygenase, putative EC=1.-.-.- (1136 residues)  
Show all sequences with this architecture.

**There are 38 sequences with the following architecture: p450, FAD_binding_6, NAD_binding_1, Fer2**  
A1UZ98_BURMS [Burkholderia mallei (strain SAVP1)] Cytochrome P450 (784 residues)  
Show all sequences with this architecture.

**There are 33 sequences with the following architecture: p450 x 3**  
Q93N82_STRLA [Streptomyces lavendulae] P450-related oxidase (397 residues)  
Show all sequences with this architecture.

**There are 17 sequences with the following architecture: p450, KR**  
Q629N7_BURMA [Burkholderia mallei (Pseudomonas mallei)] Cytochrome P450-related protein (1373 residues)  
Show all sequences with this architecture.

**There are 15 sequences with the following architecture: An_peroxidase x 2, p450**  
Q0CZ99_ASPTN [Aspergillus terreus (strain NIH 2624 / FGSC A1156)] Putative uncharacterized protein (1045 residues)

*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*

# Further Reading

*Current Protocols in Bioinformatics*
*Unit 2.5*
*Pfam*

*Current Protocols in Bioinformatics*
*Unit 2.7*
*InterPro*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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
  - PRK
  - TIGRFAM

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Conserved Domain Database (CDD)

- 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Swiss-Prot

- *Goal:* Provide a single reference sequence for each protein sequence
- Distinguishing Features
  - Non-redundancy
  - Integration with other databases (`db_xref`)
  - Ongoing curation by EBI staff and *external experts*
  - Expert annotation includes editing/updates of
    - `CC`     Comment lines
    - `FT`     Feature table
  - Distinct accession series
    - `[OPQ]12345`

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research



*Default = 500*

*Default = 10*

*Default = 0.005*

*Change cutoffs to show hits "below the line"*

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

---

*Cn3D Viewer*

*Rendering: Tubes*

*Coloring: Identity*

  *Red – matches*
  *Blue – mismatches*

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

**Cn3D Viewer**

**Rendering: Tubes**

**Coloring: Identity**
   *Red – matches*
   *Blue – mismatches*

A → V



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

# Further Reading

*Current Protocols in Bioinformatics*
*Unit 1.3*
*Entrez and Cn3D*

*Current Protocols in Bioinformatics*
*Unit 5.1*
*An Introduction to Modeling Protein*
*Structure from Sequence*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

---

# Overview

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

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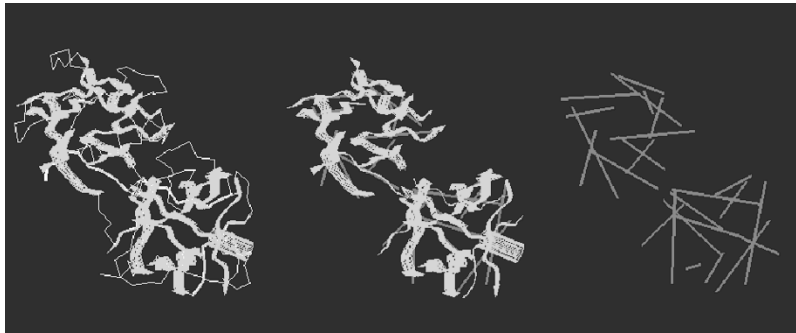  - Structures: VAST, Cn3D, and *de novo* Prediction
  - Multiple Sequence Alignment

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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"

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Interpretation

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

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research
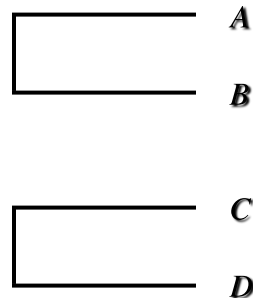
# Progressive Alignment

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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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 |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

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

# Progressive Alignment

- Align "sequence" AB with "sequence" CD

- 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# ClustalW2 Output

- Pairwise scores

- Multiple sequence alignment (`.aln`)

    - Alternative formats available:

        GCG
        Phylip
        PIR
        GDE

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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 |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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"

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# ClustalW Colors

| AVFPMILW | RED | Small (small+ hydrophobic (incl.aromatic -Y)) |
|----------|-----|-----------------------------------------------|
| DE | BLUE | Acidic |
| RK | MAGENTA | Basic - H |
| STYHCNGQ | GREEN | Hydroxyl + sulfhydryl + amine + G |
| Others | Grey | Unusual amino/imino acids etc |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

**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 |

**Calculate → Pairwise Alignments**



**Calculate → Calculate Tree → Neighbour Joining Using BLOSUM62**

# Further Reading

*Current Protocols in Bioinformatics*
*Unit 2.3*
*ClustalW*

*Current Protocols in Bioinformatics*
*Unit 3.8*
*T-Coffee*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Understanding Analyses

*Sequence* → *Results*

*Inspection*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research