# Overview

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

- Week 4
  - 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
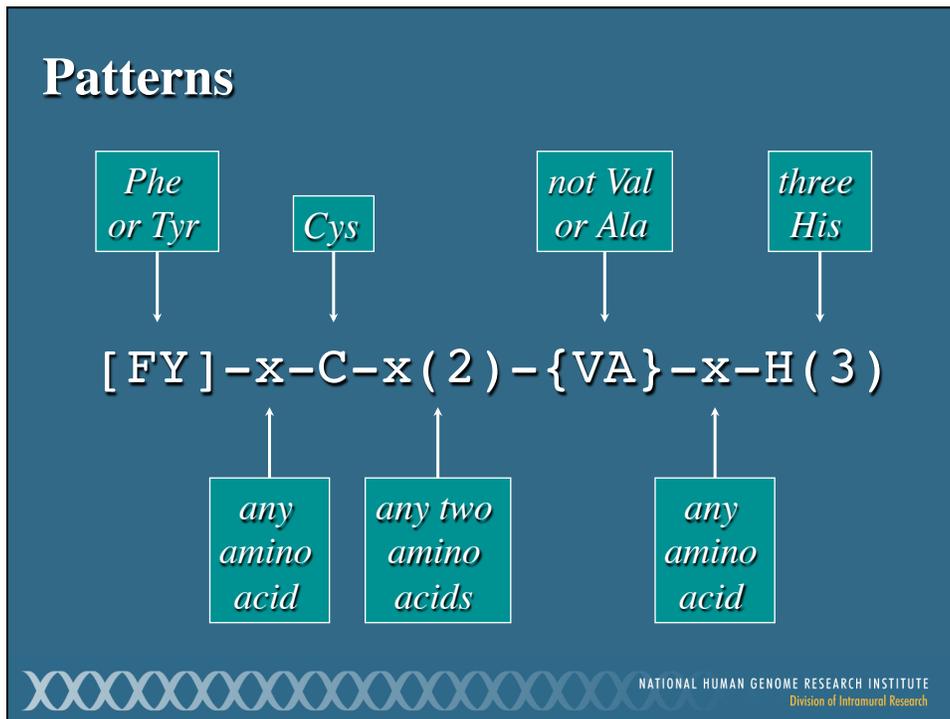- Allow for the analysis of distantly-related proteins

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 | 13 | 0 | -12 | 13 | 0 | 8 | -5 | -5 | -1 | -2 | | 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 | 6 | 7 | 6 | -11 | 19 | 11 | -5 | 6 | -10 | -11 | | | 89 | 17 | 17 | 24 | 22 | 9 | -50 | -48 | 12 |
| G | 70 | 60 | 20 | 70 | 30 | -60 | 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 |

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

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

# Sequences Used in Examples

*http://research.nhgri.nih.gov/ teaching/seq_analysis.shtml*



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Further Reading

*Current Protocols in Bioinformatics*
*Unit 2.5*
*Identifying Protein Domains with the*
*Pfam Database*

*Current Protocols in Bioinformatics*
*Unit 2.7*
*The InterPro Database and Tools for*
*Protein Domain Analysis*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Conserved Domain Database (CDD)

- Identify conserved domains in a protein sequence

- "Secondary database"
  - Pfam A (not Pfam B)
  - Simple Modular Architecture Research Tool (SMART)
  - COG (orthologous prokaryotic protein families)
  - KOG (eukaryotic equivalent of COG)
  - PRK ("protein clusters" of related protein RefSeq entries)
  - 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

# 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

# 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
  - Ongoing curation by EBI staff and *external experts*
  - Expert annotation includes editing/updates of

    | | |
    |---|---|
    | **KW** | Keyword lines |
    | **CC** | Comment lines |
    | **FT** | Feature table |

  - Distinct accession series
    **[OPQ]12345**

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

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

# 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*



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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, VAST provides a simple and fast first answer to the question of structural similarity

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research

# Overview

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

- Week 4
  - 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

- Provide basis for:
  - Predicting secondary structure
  - Performing phylogenetic analyses
  - Generating position-specific scoring matrices for use with sensitive sequence search methods

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# Overarching 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

- 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

**Find sequences to align through database searches satisfying a reasonable E-value cutoff**

↓

**Run the multiple sequence alignment program**

↓

**Inspect and assess the quality of the alignment**

↓

**Remove sequences that seriously disrupt the alignment, then realign**

↓

**Add back remaining sequences, based on key residues in the alignment**

↓

**Interpret the alignment**

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

---

# Selecting the Sequences

- 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

## Selecting the Sequences

- Sequences should be of about the same length

- Trim sequences down, so as to only use regions that have been deemed similar by either:
    - Pairwise search methods (*e.g.*, BLAST)
    - Profile-based search methods (*e.g.*, PSI-BLAST)

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

## Selecting the Sequences

- 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 dissimilar"

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

## Inspection: An Iterative Process

- Perform alignment on small set of sequences
- Examine the quality of the alignment, looking for:
  - Conservation of residues across alignment
  - Conservation of physicochemical properties
  - Relatively neat block-type structure
  - Excessive numbers of gaps

- 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

## Inspection: An Iterative Process

- Use visualization tools to identify "key residues" and "problem regions" (*e.g.*, JalView)

- Cross-check against "expertly created" multiple sequence alignments available online

- Use any available information from solved X-ray or NMR structures to nail down structurally important regions and to assess where gaps can (or cannot) be tolerated

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

- Allows for automatic multiple alignment of nucleotide or amino acid sequences
- Can align data sets quickly and easily
- Uses scoring matrices as a series
- Can bias the location of gaps, based on known structural information
- Works with Jalview, Java applet for viewing and manipulating results

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
  - Generally 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 |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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



NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

## Progressive Alignment

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

A
B

C
D

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

## 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"

A
B

C
D

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# 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, in ClustalW alignment format

    *Alternative formats available:*
    *GCG*
    *PHYLIP*
    *NEXUS*
    *NBRF/PIR*
    *GDE*
    *FASTA*

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

# ClustalW2 Output

- Cladogram
    - Tree that is 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 a phylogeny
    - Branches are *not* of equal length
    - 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  
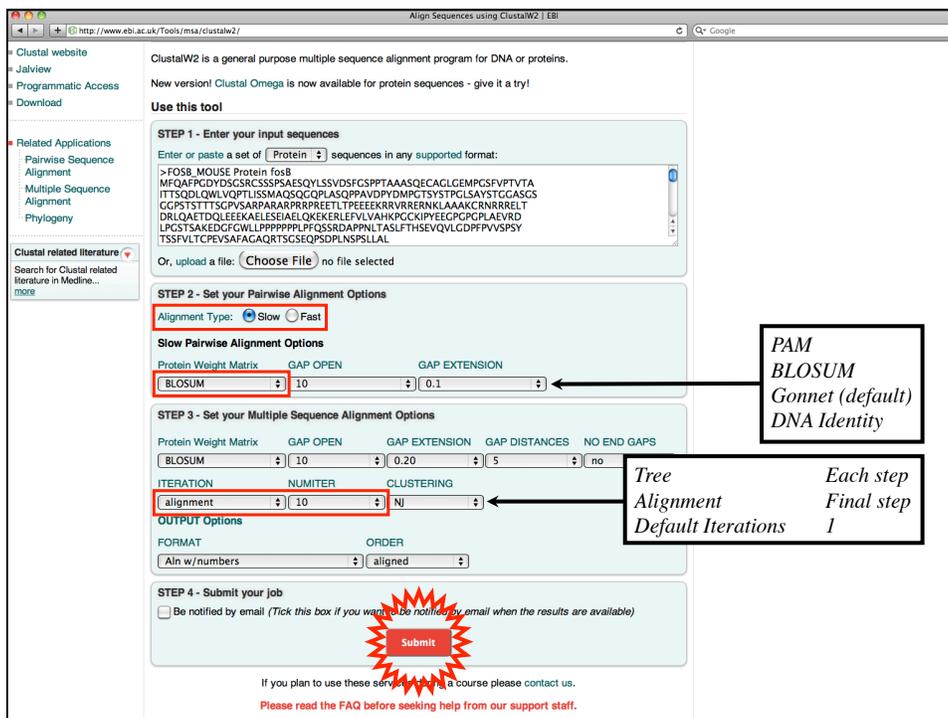Division of Intramural Research

# ClustalW2 Conservation Patterns
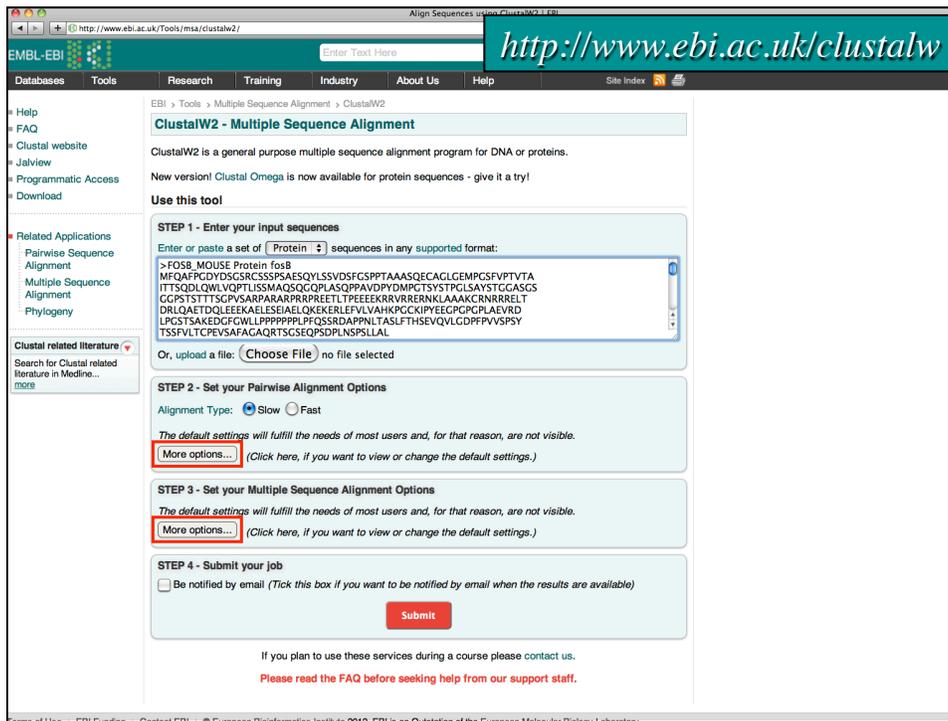
*Interpretation is empirical — there is no parallel to the E-values seen in BLAST searches to assess "significance"*

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

**:**     "conserved"  
(strongly similar properties)

**.**     "semi-conserved"  
(weakly similar properties)

NATIONAL HUMAN GENOME RESEARCH INSTITUTE  
Division of Intramural Research

| 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 |

# 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
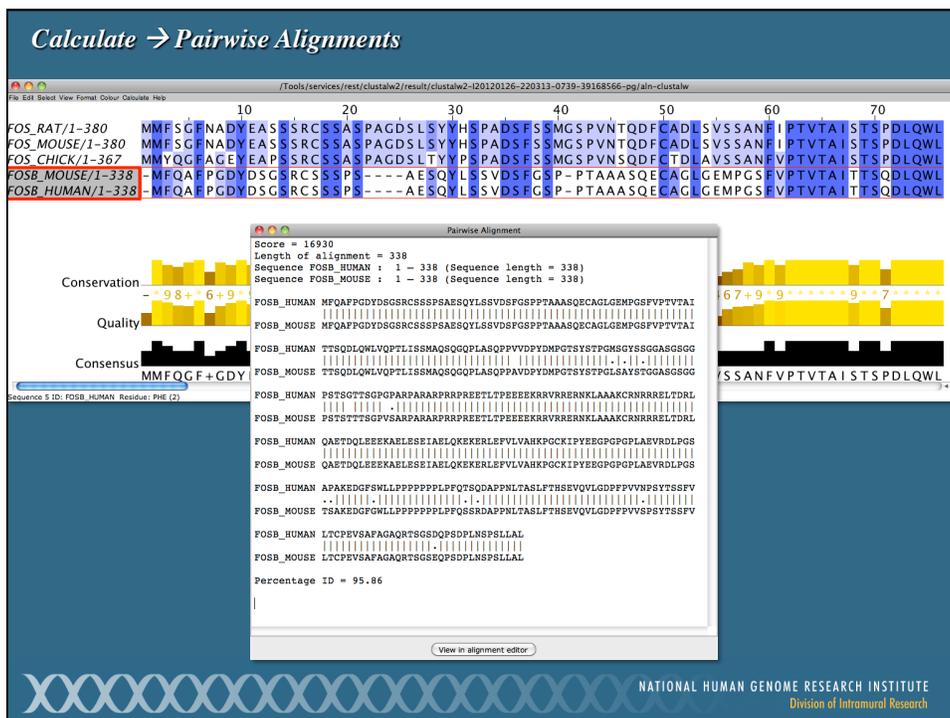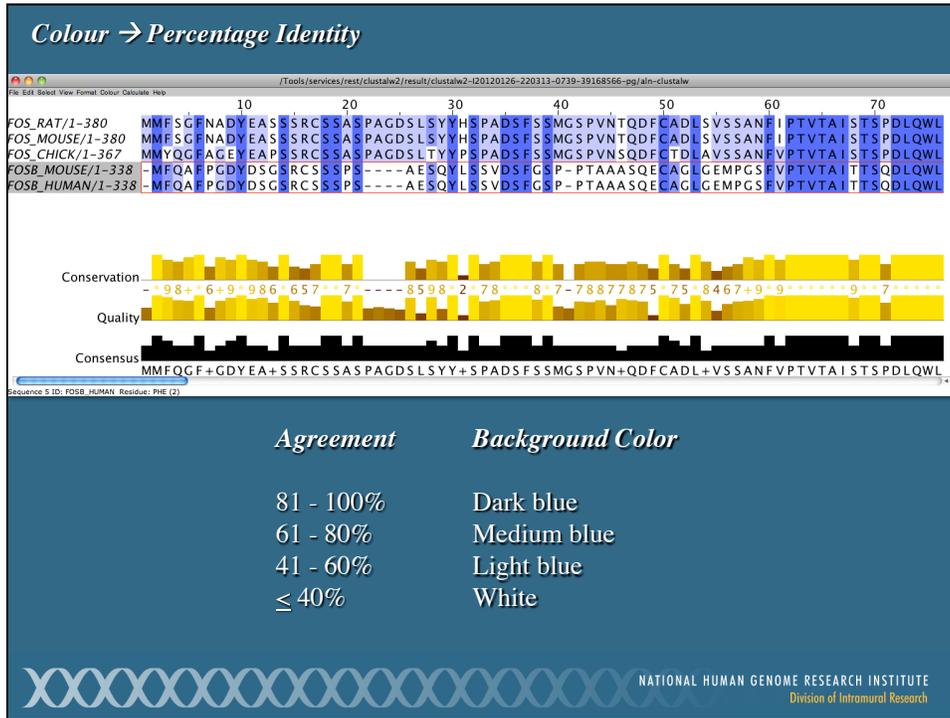- Calculation of phylogenetic trees

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 BLOSUM scores | |
| *Consensus* | Based on percent identity | |

NATIONAL HUMAN GENOME RESEARCH INSTITUTE
Division of Intramural Research

*Colour → Percentage Identity*

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



*Calculate → Pairwise Alignments*