

Introduction to evolutionary concepts and VMD/MultiSeq - Part I

Characterizing your systems

Zaida (Zan) Luthey-Schulten

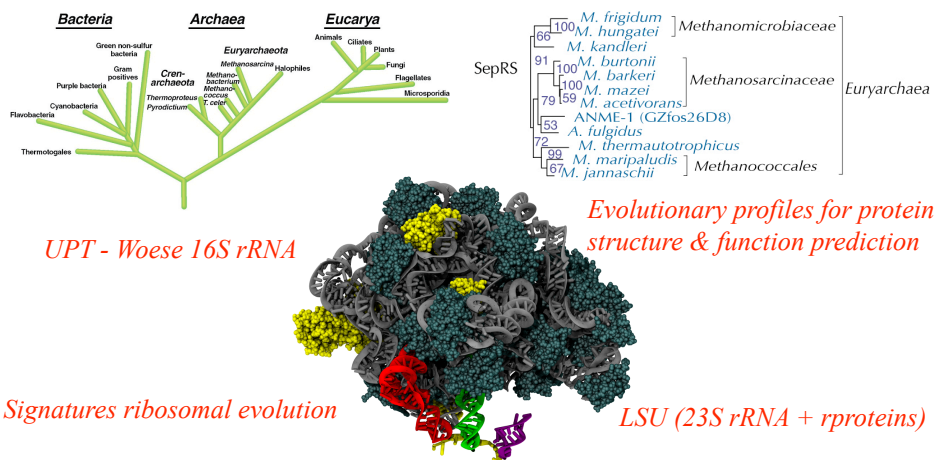
Dept. Chemistry, Beckman Institute, Biophysics, Institute of
Genomics Biology, & Physics

NIH Resource Macromolecular Modeling and Bioinformatics
Atlanta Workshop 2011



VMD/MultiSeq - “A Tool to Think”

Carl Woese - “VMD is far from a simple visualization tool for a biologist, it is a true thinking tool. Without it a whole class of biological hypotheses would simply not exist.”



New Tools in VMD/MultiSeq

Protein / RNA Sequence Data

SwissProt DB (400K), Greengenes RNA (100K) Signatures, Zoom

Metadata Information, Clustal, MAFFT & Phylogenetic Trees

RAXml Trees, Genomic Content, Temperature DB

Blast & PsiBlast

Sequence Editor

Sequence /Structure Alignment

Protein & RNA secondary structure

QR non-redundant seq / str sets

Cluster analysis / Bioinformatics scripting

Tutorials MultiSeq/AARS

EF-Tu/Ribosome

J. Eargle, D. Wright, Z. Luthey-Schulten, *Bioinformatics*, 22:504 (2006)

E. Roberts, J. Eargle, D. Wright, Z. Luthey-Schulten, *BMC Bioinformatics*, 7:382 (2006)

Aquaporin Superfamily: Bacterial & Eucaryal

Glycerol transport

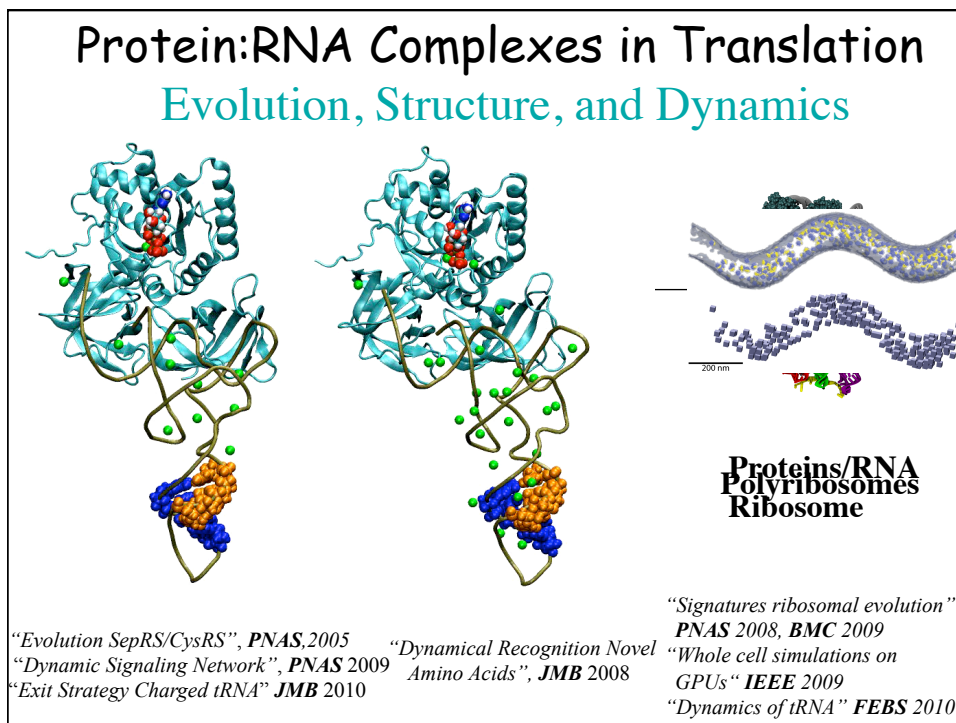
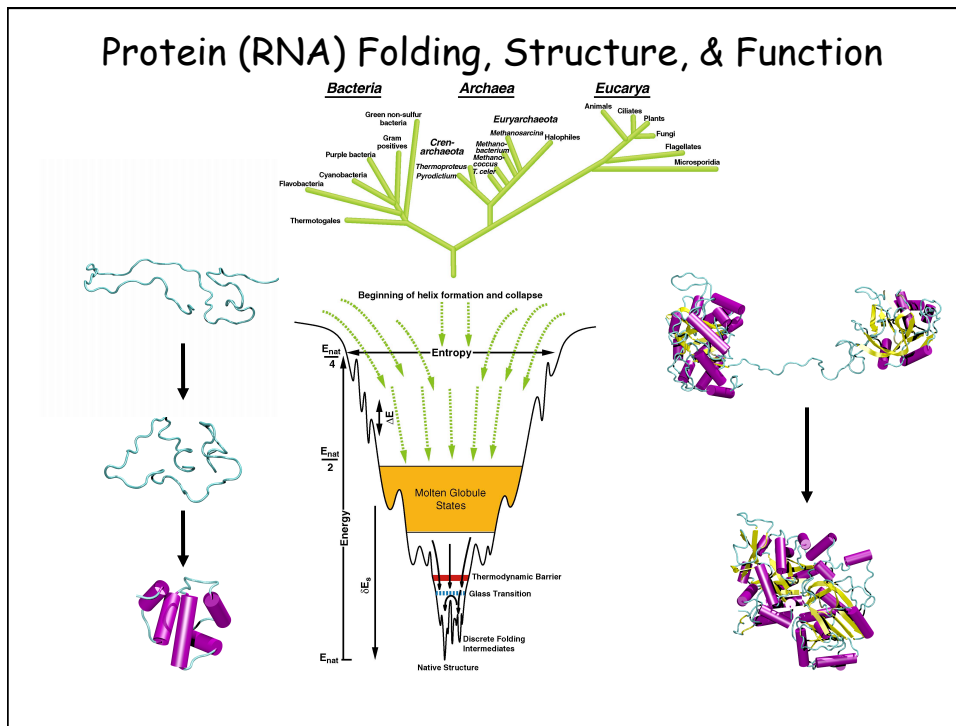
Water transport

AQP cluster

Heymann and Engel *News Physiol. Sci.* (1999) Archaeal AqpM *M. Marburgensis*, *JBC* 2003, *PNAS* 2005

```

AQP0 HUMAN ---LNTLHPAVSVGCAITVEIFLILQFVLICIFATYDE-RRNQQLGSSVALAVGFSLALGHLFGMYITGAGM 183
AQP1 HUMAN ---RNDLADGVNSGQGLGIEIIGLQLVLCVLAITDR-RRRDLGSSAPLAIGLSVALGHLAIDYTGCGI 191
AQP2 HUMAN ---VNALSNSITAGCAVTVVLFLLQLVLCIFASTDE-RRGENPGTFALSIGFSVALGHLGHIHYTGCSM 183
AQP3 HUMAN G1FATYFSGHLDMMINGFEDQFIGTASLIVCVLAITVDYNNPVPVRLGEAFIVGLVVLVIGSMGFNSGYAV 214
AQP4 HUMAN ---VTMVHGNLTAGHLLVLIIRFQLVFIIFASCDGSKRTDVTGSLIALAIGFSVAIGHLFAINYTGASM 212
AQP5 HUMAN ---VNALNNTTQGGAMVVELILFQLALCIFASTDS-RRRTSFGSSPALSIGLSVTLGHLVGIYFTGCSM 184
AQP6 HUMAN ---INVVRNSVSTGGAVAVELLQLVLCVFASTDS-RQTS--GSPATMIGISWALGHLGILFTGCSM 195
AQP7 HUMAN G1FATYLPDHMLLWRGFLNEAWLTMGLQLCLFAITDQENNPALPCEALVIGILVVIIGVSLGMNNGYAI 225
AQP8 HUMAN -AAFTYVQEQGQVAGALVAEIIITLLALAVCMGAIN--ERTKGPLAPFSIGFAVTVVDILAGGPVSGCCM 209
AQP9 HUMAN H1FATYPPAYLSLANAFADQVVAIMILLIIVPAIFDSRNLCAPRGLPEPTAIGLLIIVIASLGLNSGCAM 215
GLPF ECOLI G1FSTYFNPFINFVQAFVEMVITAILMGLLILALDDGNGVPRGPLAPLLIGLLIIVIGASMGPLTGAFM 202
ruler .....180.....200.....210.....220.....230.....240....
    
```

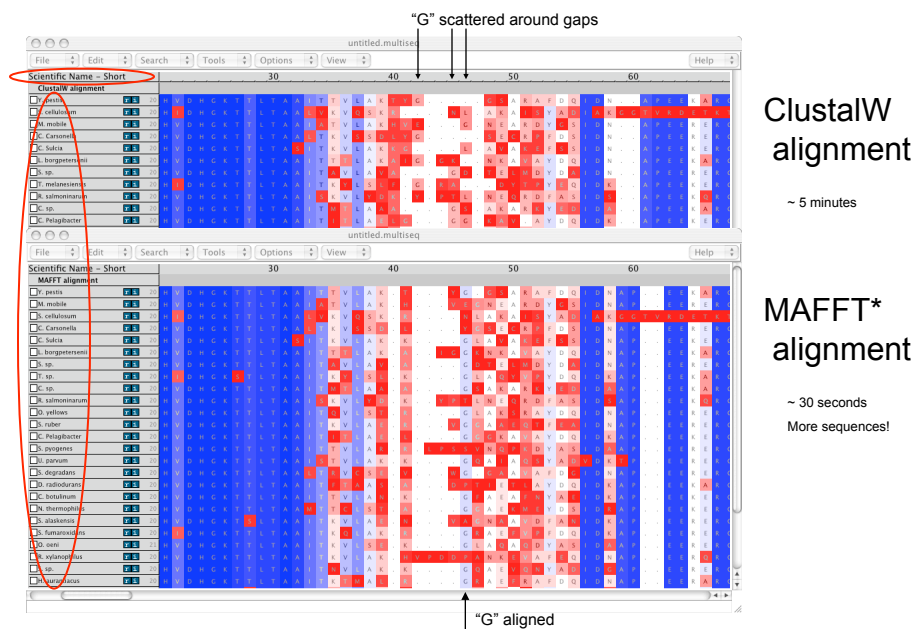


Basic principles of evolutionary analysis for proteins & RNAs

- Comparative analysis of sequences and **structures**
- Multiple sequence alignments (**gaps and editing**)
- Sequence and **structure** phylogenetic trees*
- Reference to 16S rRNA tree
- Horizontal or lateral gene transfer events
- Genomic context
- Evolutionary profiles representing diversity
- Conservation analysis of evolutionary profiles

*Various models of evolutionary change

Alignment of ~200 EF-Tu sequences in VMD/MultiSeq



* "Mafft" Katoh, Misawa, Kuma, Miyata, NAR 2002, 2005

STAMP - Multiple Structural Alignments

1. Initial Alignment Inputs

- Multiple Sequence alignment
- Ridged Body "Scan"

2. Refine Initial Alignment & Produce Multiple Structural Alignment

$$P_{ij} = \left\{ e^{-d_{ij}^2 / 2E_1} \right\} \left\{ e^{-s_{ij}^2 / 2E_2} \right\}$$

probability that residue i on structure A is equivalent to residue j on structure B.

d_{ij} -- distance between i & j

S_{ij} -- conformational similarity; function of rms between $i-1, i, i+1$ and $j-1, j, j+1$.

- Dynamic Programming (Smith-Waterman) through P matrix gives optimal set of equivalent residues.
- This set is used to re-superpose the two chains. Then iterate until alignment score is unchanged.
- This procedure is performed for all pairs.

R. Russell, G. Barton (1992) *Proteins* 14: 309.

Modified for RNA, John Eargle, VMD/MultiSeq.

Multiple Structural Alignments

STAMP – cont'd

2. Refine Initial Alignment & Produce Multiple Structural Alignment

Alignment score:

$$S_C = \frac{S_P}{L_P} \frac{L_P - i_A}{L_A} \frac{L_P - i_B}{L_B}$$

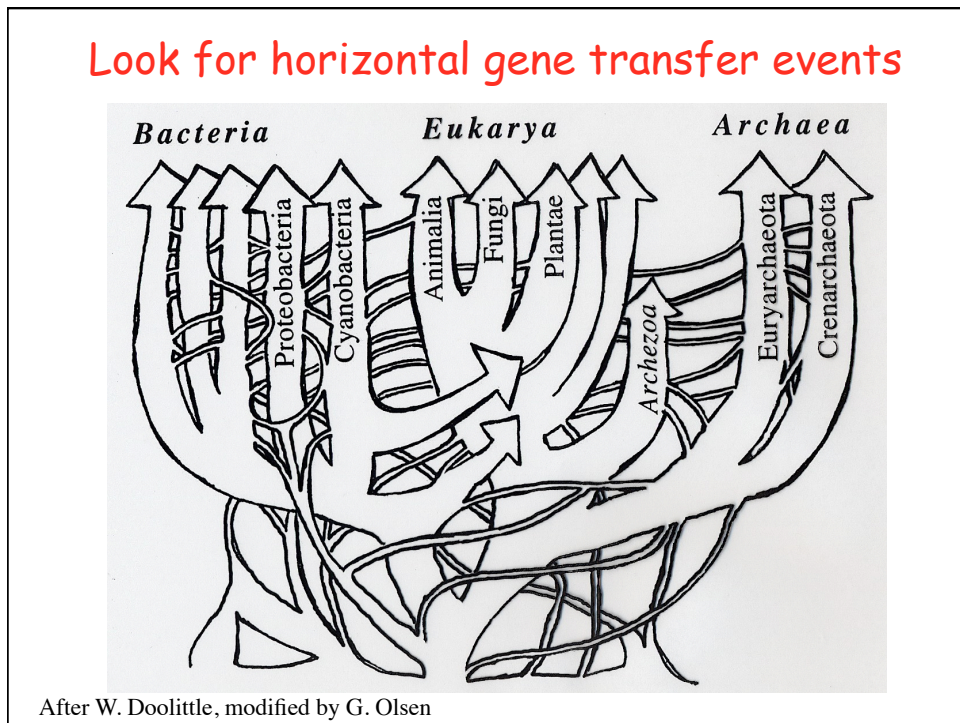
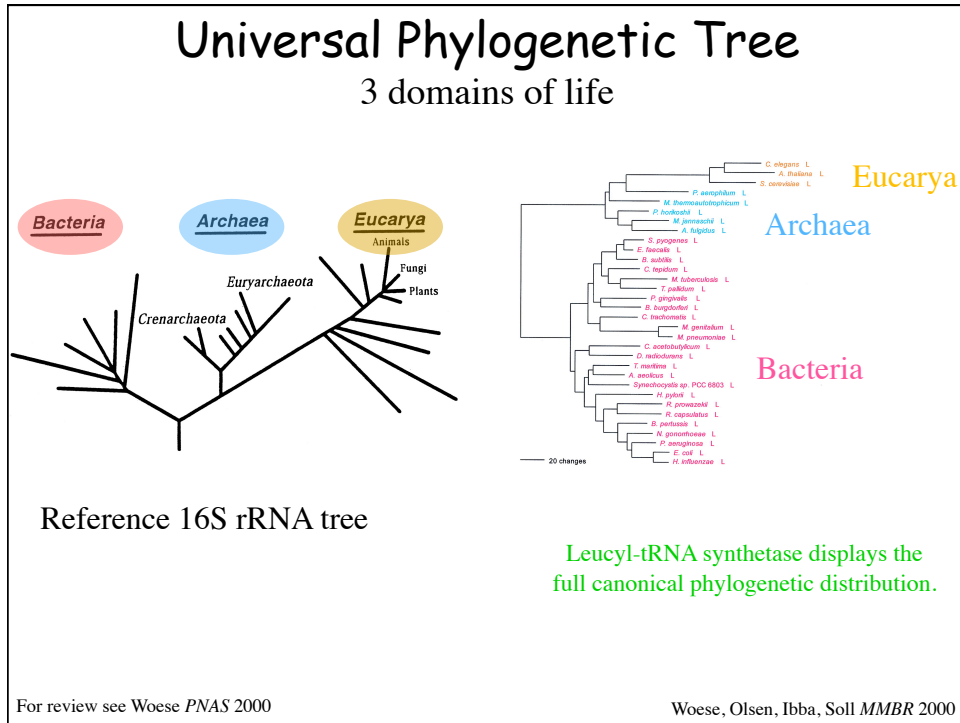
$$S_P = \sum_{\text{align. path}} P_{ij}$$

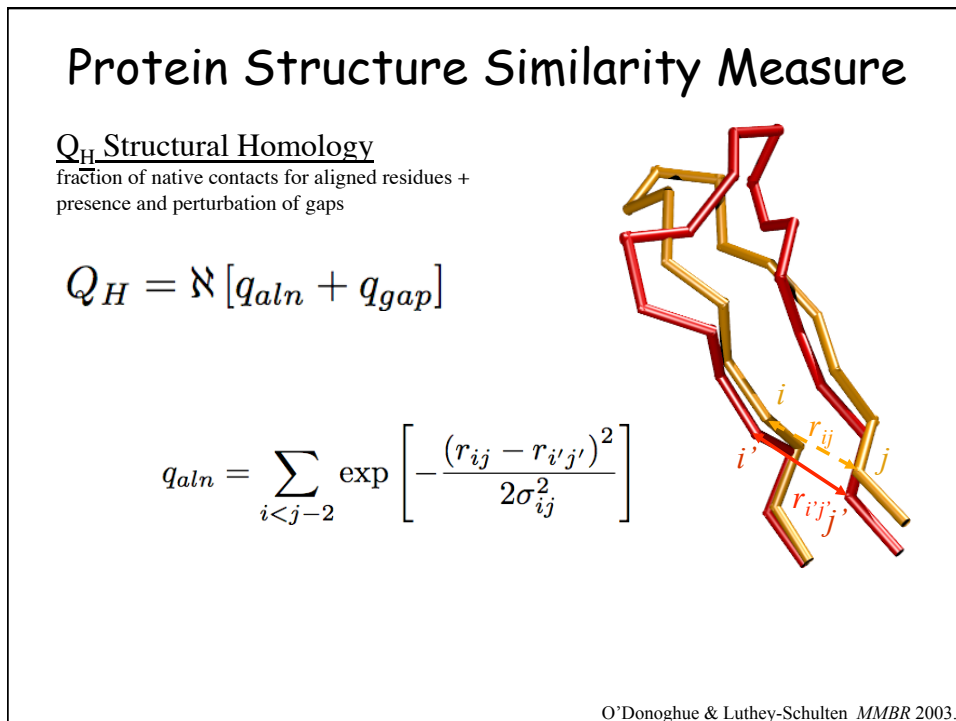
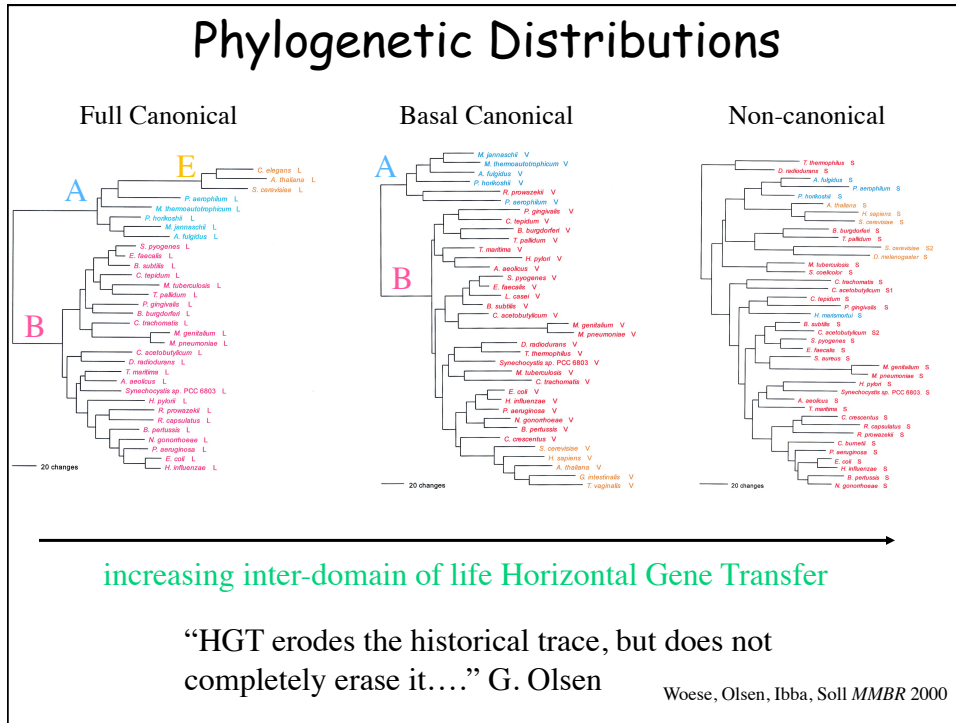
L_P, L_A, L_B -- length of alignment, sequence A, sequence B

i_A, i_B -- length of gaps in A and B.

Multiple Alignment:

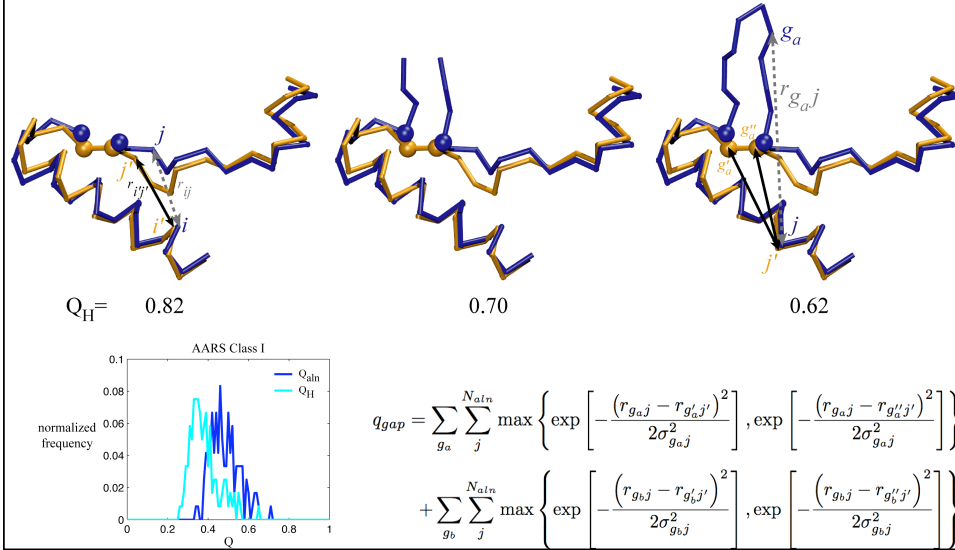
- Create a dendrogram using the alignment score.
- Successively align groups of proteins (from branch tips to root).
- When 2 or more sequences are in a group, then average coordinates are used.



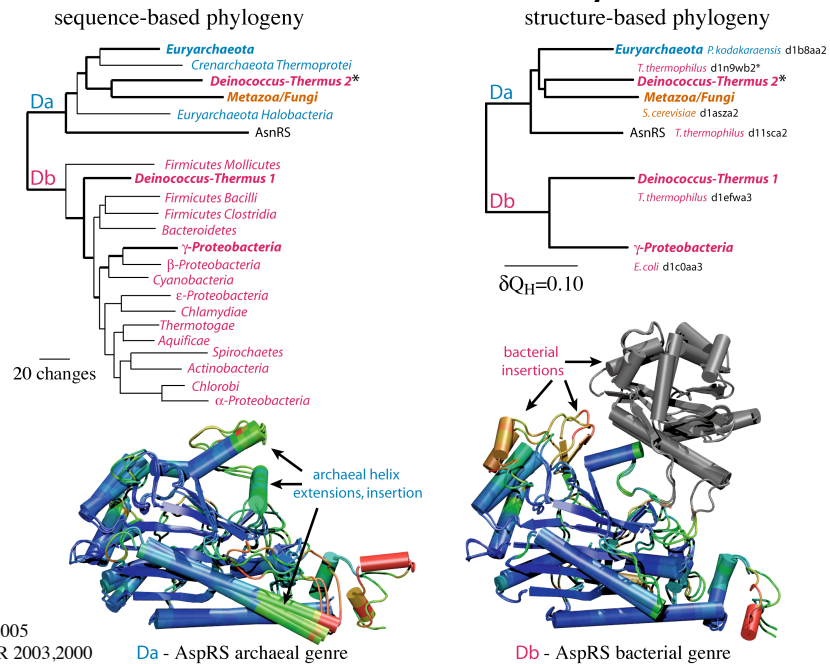


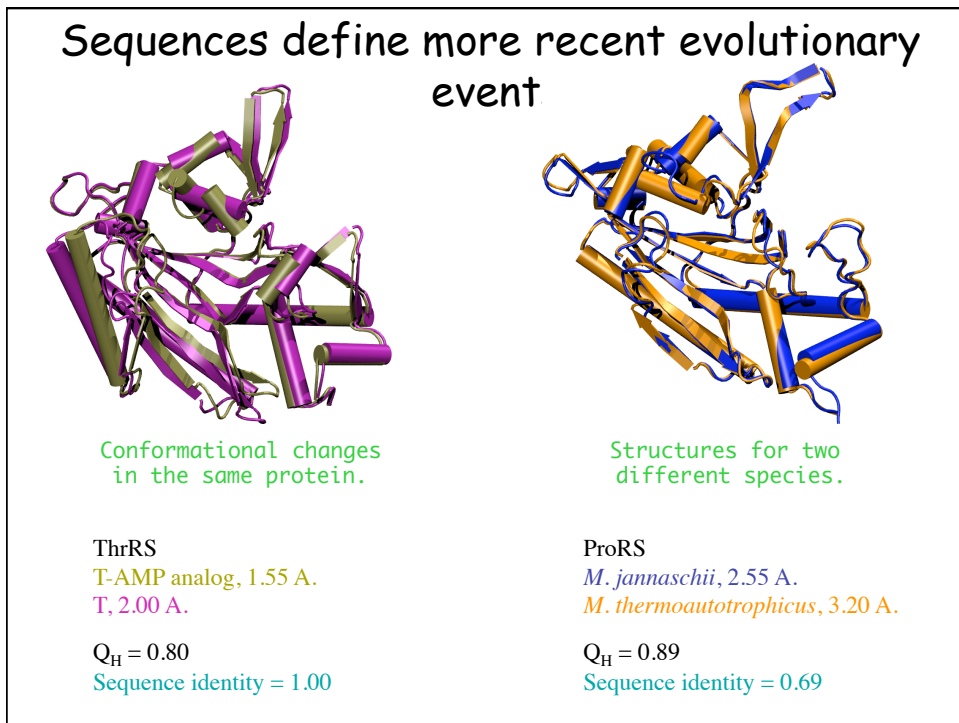
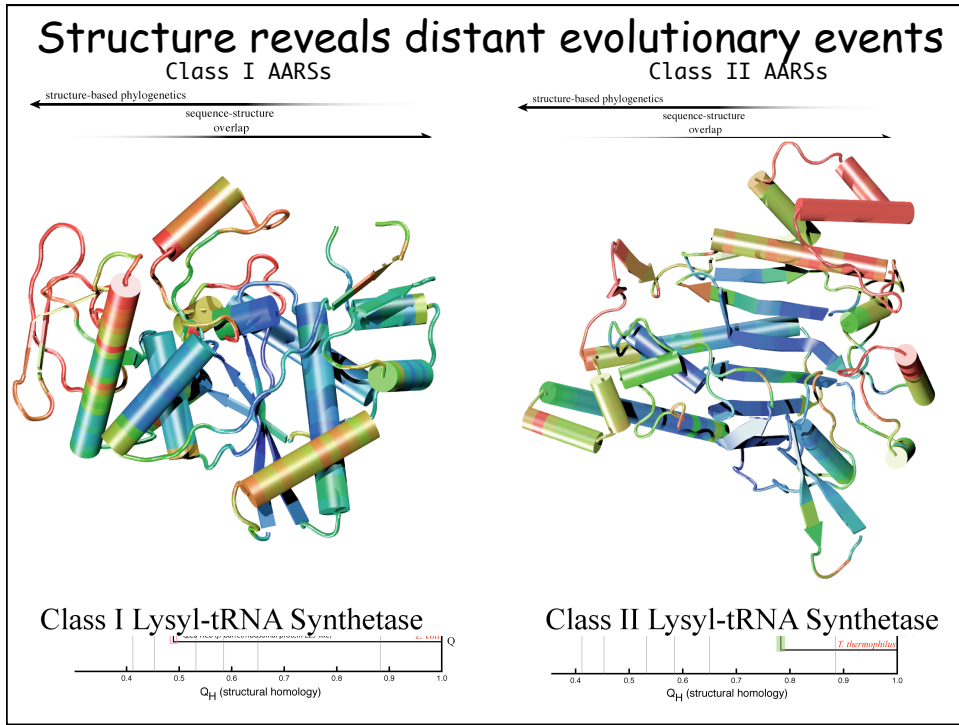
Structural Similarity Measure: The effect of insertions

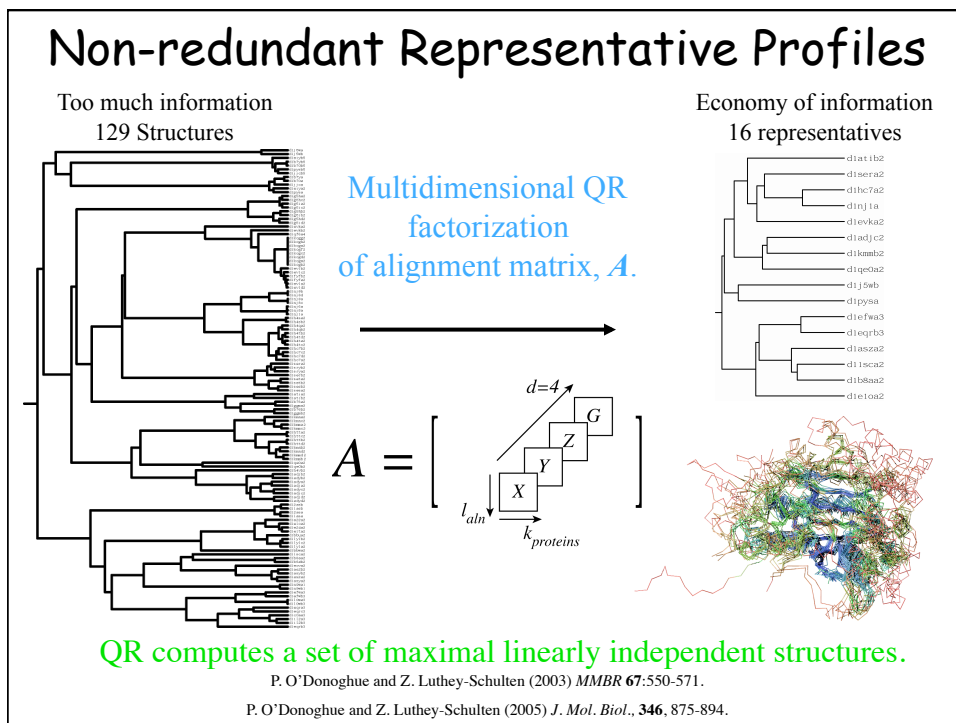
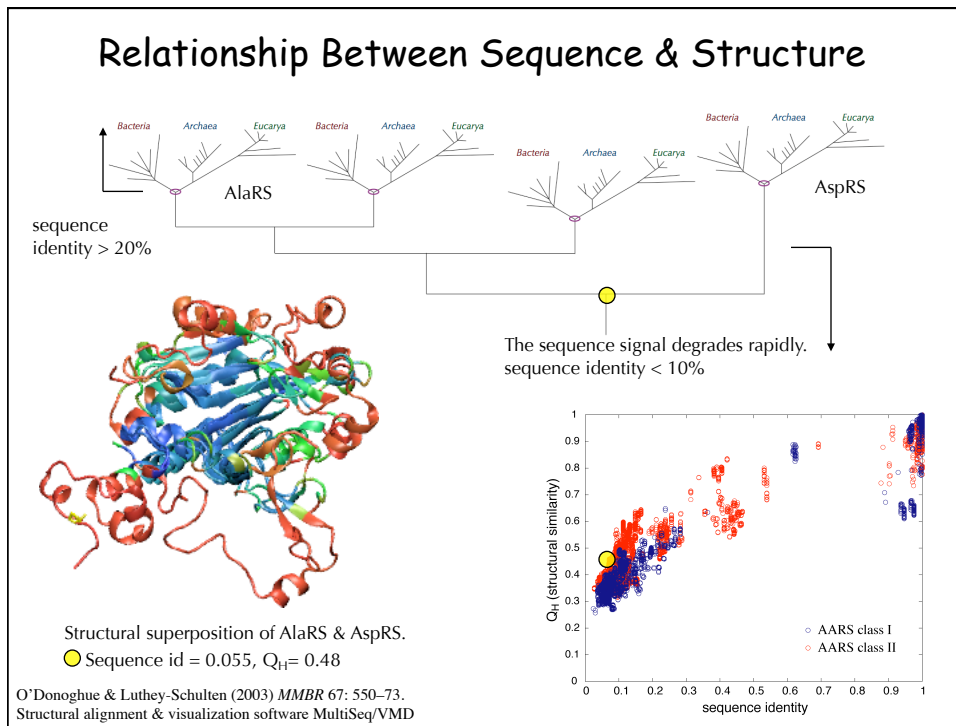
“Gaps should count as a character but not dominate” C. Woese



Structure encodes evolutionary information!







Numerical Encoding of Proteins in a Multiple Alignment

Encoding Structure

Rotated Cartesian + Gap = 4-space

Aligned position $(x_{C_a}, y_{C_a}, z_{C_a}, 0)$

Gapped position $(0, 0, 0, g)$

Gap Scaling $g = \gamma \frac{\|X\|_{F_4} + \|Y\|_{F_4} + \|Z\|_{F_4}}{\|G\|_{F_4}}$

adjustable
parameter

Sequence Space

Orthogonal Encoding = 24-space

23 amino acids (20 + B, X, Z) + gap

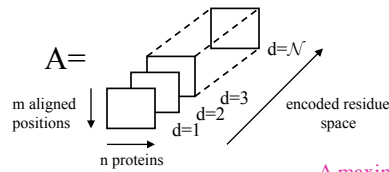
A = (1,0)

B = (0,1,0)

C = (0,0,1,0)

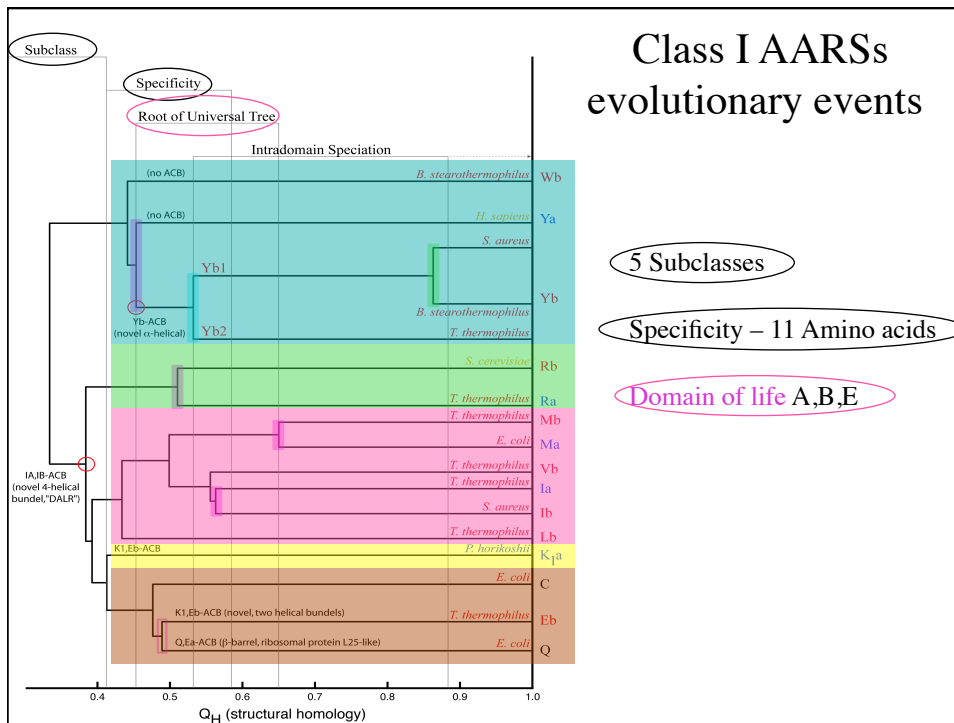
...
GAP = (0,1)

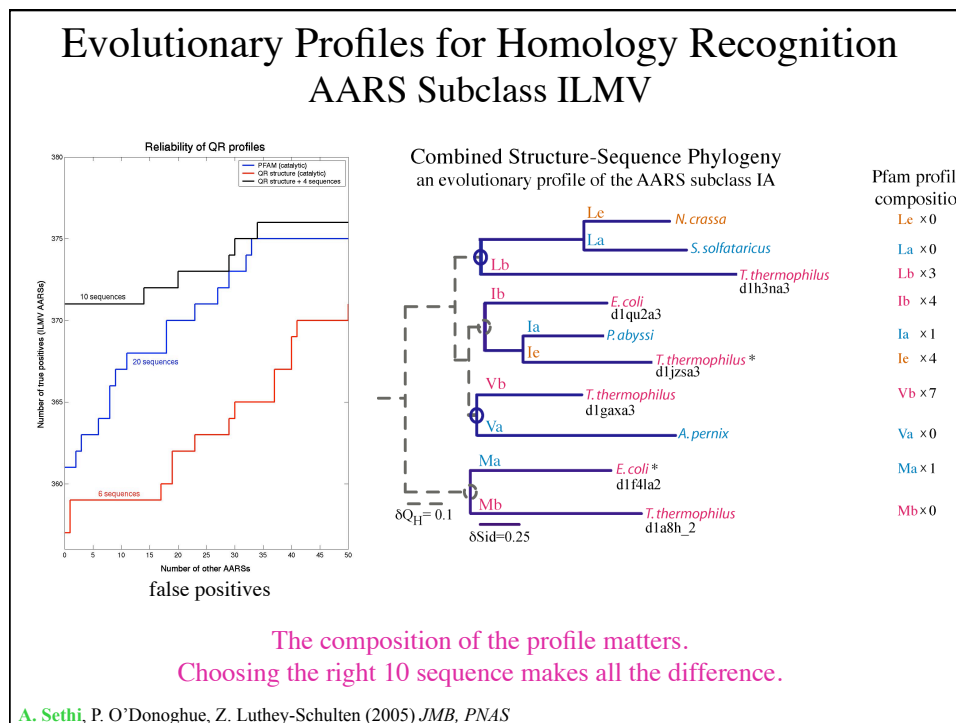
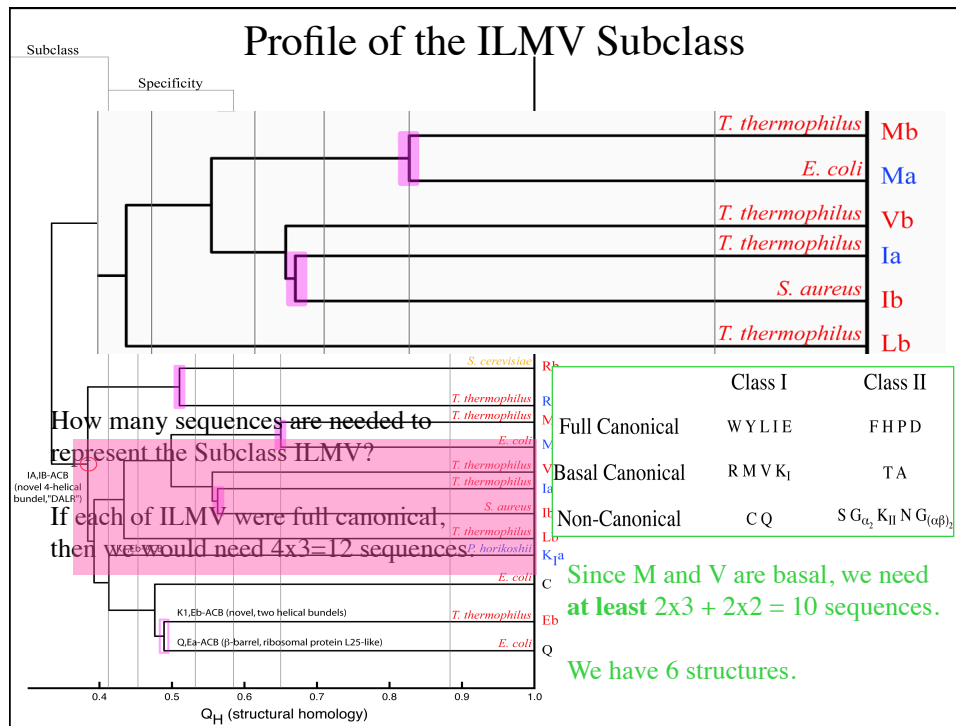
Alignment is a Matrix with Linearly Dependent Columns

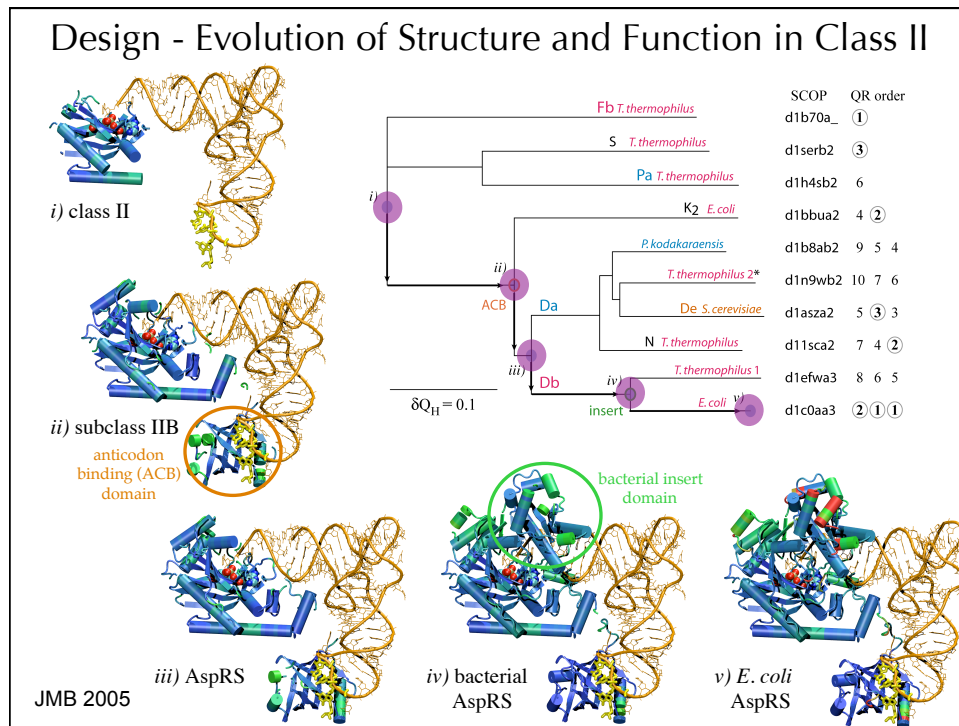


$$Q_{(d)}^T A_{(d)} P = Q_{(d)}^T \begin{bmatrix} d=4 & & & & \\ & d=3 & & & \\ & & d=2 & & \\ & & & d=1 & \\ & & & & m_{aln} \end{bmatrix} \begin{bmatrix} G \\ Z \\ Y \\ X \\ n_{proteins} \end{bmatrix} P = \tilde{R}_{(d)}$$

A maximal linearly independent subset can be determined with respect to a threshold, e.g., similarity measure threshold.







Summary Structural Profiles

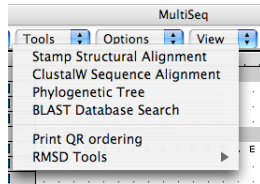
1. Structures often more conserved than sequences!! Similar structures at the Family and Superfamily levels.

Add more structural information to identify core and variable regions

2. Which structures and sequences to include? Use evolution and eliminate redundancy with QR factorization

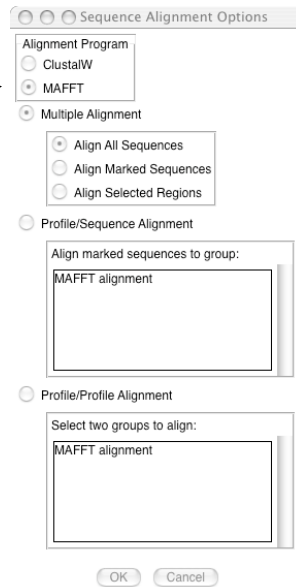
What is MultiSeq?

- MultiSeq is an extension to VMD that provides an environment to combine sequence and structure data
- A platform for performing bioinformatics analyses within the framework of evolution
- Provides software for improving the signal-to-noise ratio in an evolutionary analysis by eliminating redundancy (**StructQR, SeqQR, Evolutionary Profiles “EP”**)
- Visualizes computationally **derived metrics** (Q_{res} , $Q_{H,\dots}$) or imported experimental properties



- Integrates popular bioinformatics tools along with new algorithms (ClustalW, **MAFFT**, BLAST, **STAMP**, **Signatures**, **Mutual information**, QR, PT,....)

Choose MAFFT to perform
multiple sequence
alignment →



New Tools in VMD/MultiSeq

Protein / RNA Sequence Data

SwissProt DB (400K), Greengenes RNA (100K) Signatures, Zoom

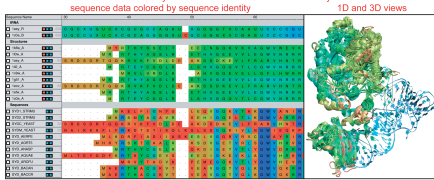
Metadata Information, Clustal & Phylogenetic Trees

RAXml Trees, Genomic Content, Temperature DB

Blast & PsiBlast

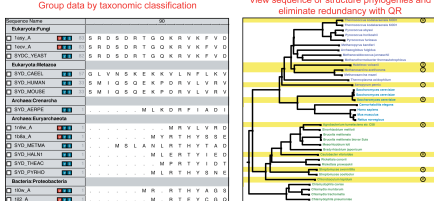
Sequence Editor

View structural data colored by structural conservation and sequence data colored by sequence identity



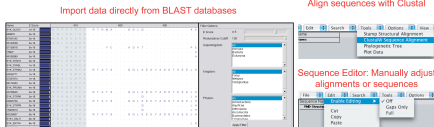
Synchronization between 1D and 3D views

Group data by taxonomic classification

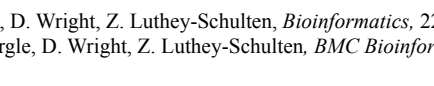


View sequence or structure phylogenies and eliminate redundancy with QR

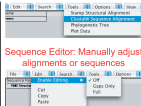
Align sequences with Clustal



Import data directly from BLAST databases



Sequence Editor: Manually adjust alignments or sequences



Sequence /Structure Alignment

Protein & RNA secondary structure

QR non-redundant seq / str sets

Cluster analysis / Bioinformatics scripting

Tutorials MultiSeq/ AARS

EF-Tu/Ribosome

J. Eargle, D. Wright, Z. Luthey-Schulten, *Bioinformatics*, 22:504 (2006)
 E. Roberts, J. Eargle, D. Wright, Z. Luthey-Schulten, *BMC Bioinformatics*, 7:382 (2006)

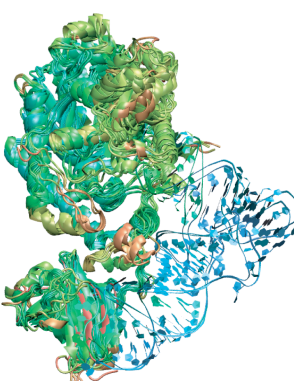
MultiSeq Combines Sequence and Structure

- Align sequences or structures; manually edit alignments
- View data colored by numerous metrics including structural conservation and sequence similarity
- Synchronized coloring between 1D and 3D views

Variation in structures

Variation in sequences

Sequence Name	40	50
IRNA		
1asy_R	U G U C X C G U G C C A G A U	X G G G G G T
1c0a_B	U X U C A C G C A G G G G X U C	G C G G G G X
Structures		
158a_A	M Y R T H Y S S E I T	E E L N G G
150w_A	M R R T H Y A G S L R	E T H V G E
1asy_A	D H T O O K R V K F V D L D E	A K D S D K
152_A	M R T E Y C G G L R	L S H V G G
159w_A	M R H V L V R D L K	L A H V G G
1551_A	M R R T H Y A G S L R	E T H V G E
159v_A	D R T D G K R W K F V D L D E	A K D S D K
159e_A	M R R T H Y A G S L R	E T H V G E
1c0a_A	M R H T E Y G G G L R	L S H V G G
Sequences		
SYD1_STRMU	M K E L F I G N Y G	L E Q V G G
SYD2_STRMU	M K S M A G A V R	S E H I G G
SYD3_YEAST	D R T G Q K R V K F Y D L D E	A K D S D K
SYD4_YEAST	K K F L F K D T S T I K O L K	G L S S G G
SYD_AERPE	M L K D R F I A D I I A S K	E S L V G G
SYD_AGR15	M H R Y R S H T C A A L R	K S D V G E
SYD_ANASP	M R T H Y C G E L R	Q K D I G E
SYD_AQUAE	Y G D F K R T K Y C G E V S	E E D I G K
SYD_ARCFU	M R V Y T A D V K	P E M E G G
SYD_BAGAN	M A E R T H A C G K V T	V E A V G G
SYD_BAGCR	M A E R T H A C G K V T	V E A V G G



Load large sequence sets

Swiss-Prot (Proteins)

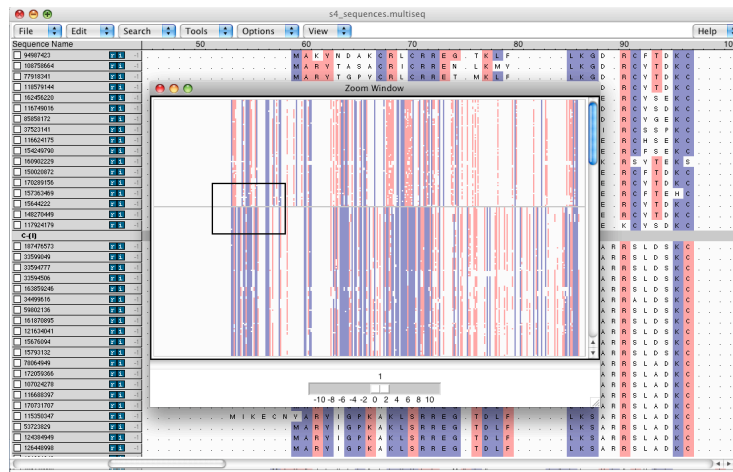
Curated sequences
392,667 sequences
Unaligned
177 MB on disk
2 minutes to load
2.4 GB memory used

Greengenes (RNA)*

Environmental 16S rRNA
90,654 entries
Aligned (7682 positions)
670 MB on disk
2.5 minutes to load *
4.0 GB memory used*

Sequence editor

- New sequence API allows editing of large alignments. Align closely related sequences by group, combine groups, and then manually correct.
- Zoom window gives an overview of the alignment, quickly move the editing window to any part of the alignment.



660 sequences
of ribosomal
protein S4 from
all complete
bacterial
genomes*.

* K. Chen, E. Roberts, Z Luthey-Schulten (2009) BMC Bioinformatics

Secondary structure prediction

- Integration with PSIPRED* to predict secondary structure of sequences.
- Compare to VMD STRIDE predictions from structures.

Sequence Name	150	160	170	180
Secondary Structures				
<input type="checkbox"/> Hpylori_S4	[Diagram showing secondary structure prediction for Hpylori_S4]			
<input type="checkbox"/> Thermus_S4	[Diagram showing secondary structure prediction for Thermus_S4]			
<input type="checkbox"/> Ecoli_S4	[Diagram showing secondary structure prediction for Ecoli_S4]			
Sequences				
<input type="checkbox"/> Hpylori_S4	I E I K	E K T K S N S Q V V R A M	E L T A Q T G I V P	W I D V E K D K K Y G I F T R
<input type="checkbox"/> Thermus_S4	I A V A	E K S R N L E L I R Q N L	E A M K G R K V G P	W L S L D V E G M K G K F L R
<input type="checkbox"/> Ecoli_S4	V S I R	E K A K K Q S R V K A A L E	L A E Q R E K P T	W L E V D A G K M E G T F K R

Secondary Structure

Predict

OK Cancel

Modeling of *Helicobacter pylori* ribosomal protein S4 using two known bacterial structures from *Thermus thermophilus* and *Escherichia coli*.

Zinc-binding site replaced by salt bridge in *H. pylori*.

* D. Jones (1999) *J Mol Biol*

PSIPRED installation

- PSIPRED is not included with VMD, must be installed locally.
- Configured in the MultiSeq software preferences dialog (File->Preferences).

Requires a sequence database filtered for problematic regions. Here using Swiss-Prot for relatively fast predictions.

Metadata Software

BLAST Installation Directory

/usr/local/blast Browse...

BLASTMAT data

BLASTDB

PSIPRED Installation Directory

/Volumes/HomeRAID2/Homes/erobert3/Applications/OSX-386/bin/ Browse...

PSIPREDDATA /AID2/Homes/erobert3/Applications/OSX-386/share/psipred/data

PSIPREDDB /Volumes/Homes/Databases/psipred/psipred-sp

Path to external editor

Browse...

Close Help

Export Modeller compatible alignments

- MultiSeq can automatically export SIF alignment files compatible with Modeller.

```
>P1; Hpylori_S4
sequence:Hpylori_S4:::::0.00:0.00
MARYRGAVERLERRRFGVSLALKEG-RRLSGKSALDKRAYGPGHGQR-RAKTSYDGLQLK
EKQKAMMYGISEKQFRSIFVEANRLDGNTEENLRLIERRLDNVVYRMGFATTRSSARQ
LVTHGHVLDGKRLDIPSYFVRSQGKIEIKETKNSQVVRAMELTAQTGIVPWIDVEKD
KKGIFTRYPEREEVVVPIERLIVELYSK*

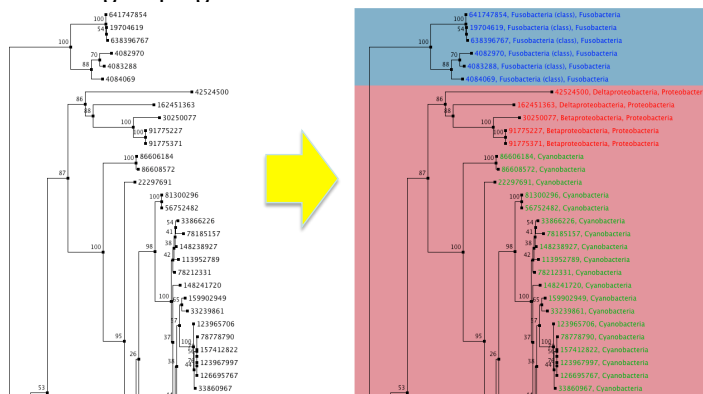
>P1; Thermus_S4
structureX:Thermus_S4:2:D:209:D::-1.00:-1.00
-GRYIGPVCRLCRREGVKLYLKEG-RCYSPKCAMERRPYPPGQHGQKRARRRPSDYAVRLR
EKQKLRRIYGISERQFRNLFEEASKKGVTVGSVFLGLESRLDNVVYRLGFVSRROARQ
LVRHGHTVNGRRVLDLPSYRVRPGDEIAVAEKSRNLELIRQNLEAMKGRKVGWPLSLDVE
GMKGFRLRPREDLALPVNEQLVIEFYSR*

>P1; Ecoli_S4
structureX:Ecoli_S4:1:D:205:D::-1.00:-1.00
-ARYLGPVKLSRREGTDLFLKSGVRAIDTKCKIE---QAPGQHGAR-KPRLSDYGVQLR
EKQKVRRIYGVLERQFRNYKAEARLKGNTGENLLALLEGRLDNVVYRMGFATRAEARQ
LVSHKAIMVNGRVVNIASVQVSPNDVVSIREKAKQSRVKAALAEQREKPTWLEVDAG
KMEGTFKRKPERSDLSADINEHLIVELYSK*
```

```
a = mymodel(env, alnfile='alignment.ali', knowns=('Ecoli_S4','Thermus_S4'), sequence='Hpylori_S4')
a.starting_model = 1
a.ending_model = 20
a.make()
```

Phylogenetic tree editor

- Automatically add annotations and colors to phylogenetic trees based on taxonomy, enzyme, temperature class, and/or MultiSeq groupings.



A cluster of five proteobacterial sequences branch near the cyanobacterial sequences. These are cases of horizontal gene transfer.

Maximum likelihood tree of 660 S4 sequences reconstructed using RAXML.

Leaf Colors

■ Fusobacteria	■ Proteobacteria	■ Cyanobacteria
■ Chlamydiae	■ Firmicutes	■ Planctomycetes
■ Spirochaetes	■ Verrucomicrobia	■ Thermicutes
■ Chlorobi	■ Acidobacteria	■ Chloroflexi
■ Thermotogae		

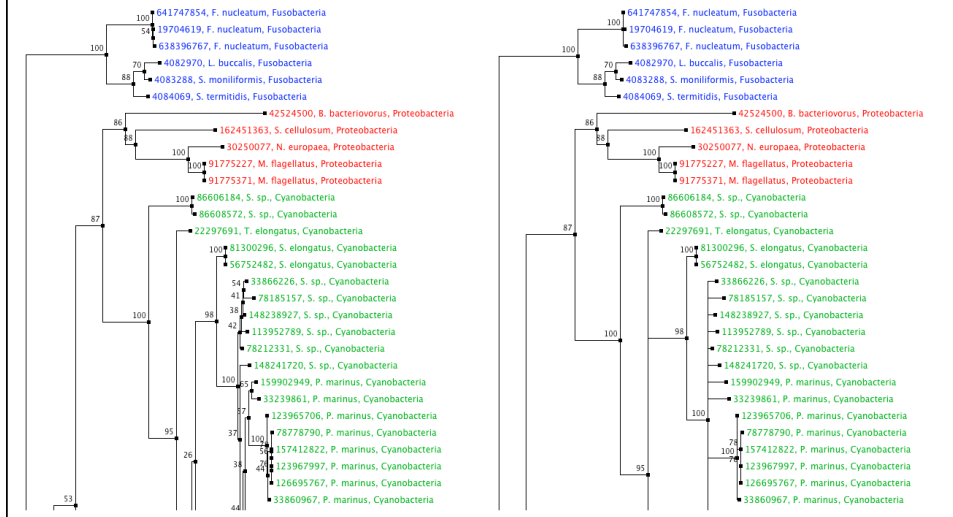
Background Colors

 C-(IV)_in	 C-(IV)_out	 C-(VI)	 C-(III)
 C-(III)	 C+	 C-(I)	

Elijah Roberts 2009

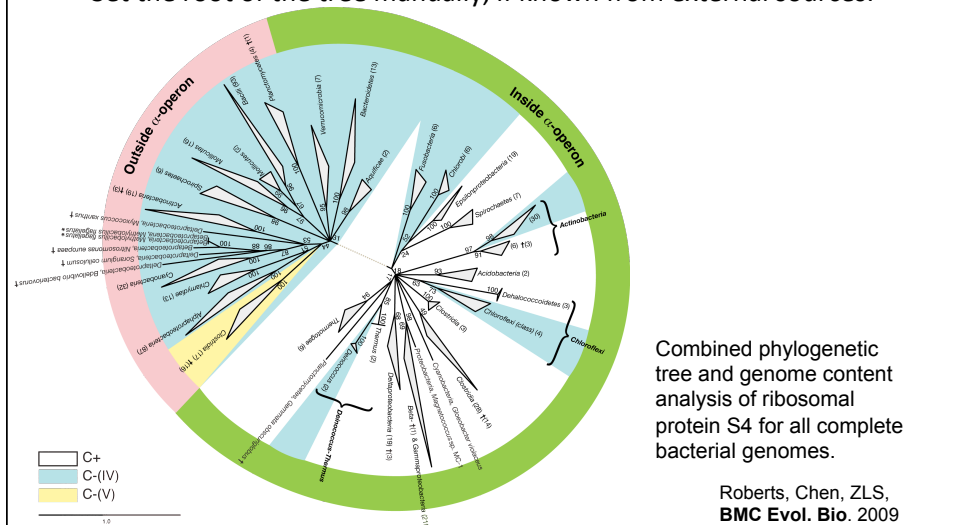
Edit the physical layout of the tree

- Nodes with low support can be removed.
- Nodes can be rotated for easier reading.



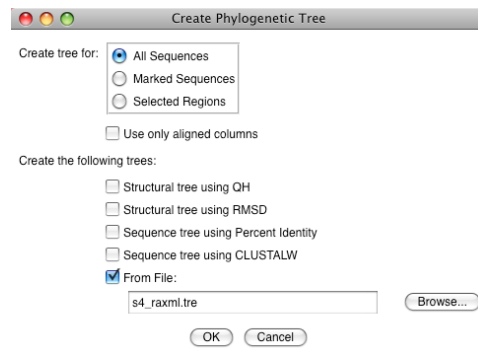
Manipulate branches to simplify the tree

- Manually collapse by node.
- Automatically collapse clades that are alike according to taxonomy, enzyme, temperature class, and/or MultiSeq grouping.
- Set the root of the tree manually, if known from external sources.



Phylogenetic tree generation

- Generate distance based trees only over well-aligned columns (no indels).
- Export alignments in Phylip format (PHY) compatible with RAxML for maximum likelihood reconstructions.
- Import Newick trees from phylogenetic reconstruction programs (including RAxML).



Scripting MultiSeq

- All MultiSeq functions can now be scripted.
- Scripting an analysis provides benefits:
 - It can be checked for correctness.
 - It can be quickly repeated by anyone.
 - It can be modified later with new functionality.
 - It can be run on a cluster in VMD text mode.
(if it can be easily broken into independent chunks)
- Many functions are too user specific and/or too complex to be turned into a GUI.
- Some examples of MultiSeq scripts...

Genome content

- When using sequence from fully sequenced genomes, additional information is available in the genome content.
- Conservation of gene ordering, neighbors, or intergenic regions can provide additional evolutionary information not contained in the sequence.
- Gene names and ordering can be obtained from the genome PTT files, want to organize the information in an evolutionarily meaningful manner.

Location	Strand	Length	PID	Gene	Synonym	Code	COG	Product
3437638..3438021	-	127	16131173	rplQ	b3294	-	COG0203J	50S ribosomal subunit protein L17
3438062..3439051	-	329	16131174	rpoA	b3295	-	COG0202K	RNA polymerase, alpha subunit
3439077..3439697	-	206	16131175	rpsD	b3296	-	COG0522J	30S ribosomal subunit protein S4
3439731..3440120	-	129	16131176	rpsK	b3297	-	COG0100J	30S ribosomal subunit protein S11
3440137..3440493	-	118	16131177	rpsM	b3298	-	COG0099J	30S ribosomal subunit protein S13
3440640..3440756	-	38	16131178	rpmJ	b3299	-	COG0257J	50S ribosomal subunit protein L36
3440788..3442119	-	443	16131179	secY	b3300	-	COG0201U	preprotein translocase membrane subunit
3442127..3442561	-	144	16131180	rplO	b3301	-	COG0200J	50S ribosomal subunit protein L15
3442565..3442744	-	59	16131181	rpmD	b3302	-	COG1841J	50S ribosomal subunit protein L30
3442748..3443251	-	167	16131182	rpsE	b3303	-	COG0098J	30S ribosomal subunit protein S5

Combined genomic context/phylogenetic tree

- Use a script to walk through a phylogenetic tree, find the genome content near the source gene, create a graphical representation of the combined data.

```

proc draw_genome_context_of_phylogeny {args} {
  # Load the sequences.
  set alignment [::SeqData::Fasta::loadSequences $alignmentFilename]

  # Load the tree
  set tree [::PhyloTree::Newick::loadTreeFile $treeFilename]

  # Reorder the alignment by the tree.
  set treeAlignment {}
  set leafNodes [::PhyloTree::Data::getLeafNodes $tree]
  foreach node $leafNodes {
    set foundNode 0
    set nodeName [::PhyloTree::Data::getNodeName $tree $node]
    foreach sequence $alignment {
      if {$nodeName == [::SeqData::getName $sequence]} {
        lappend treeAlignment $sequence
        set foundNode 1
        break
      }
    }
  }

  # Draw the genomic context.
  drawGenomicContextOfAlignment $outputFilename $treeAlignment $contextDistance $scaling $genomeDirectory
}

```

Combined genomic context/phylogenetic tree

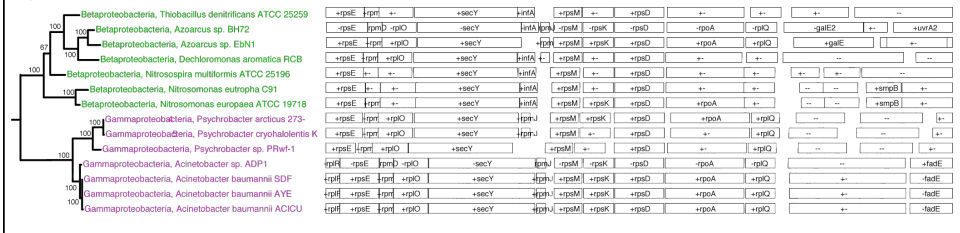
```

proc drawGenomicContextOfAlignment {outputFilename alignment contextDistance scaling genomeDirectory} {
    foreach sequence $alignment {
        # Make sure we have the GI number for this sequence.
        set giNumber [::SeqData::getSourceData $sequence "gi"]

        # Make sure we can tell which genome this sequence is from.
        set taxonomy [join [::SeqData::getLineage $sequence 1 0 1] ","]
        if (![info exists genomeTaxonomyMap($taxonomy)]) {
            error "ERROR: Unknown genome for sequence [::SeqData::getName $sequence]: $taxonomy"
        }

        # Go through each of the genome context files for the genome.
        set foundGene 0
        foreach genomeName $genomeTaxonomyMap($taxonomy) {
            ...
        }
    }

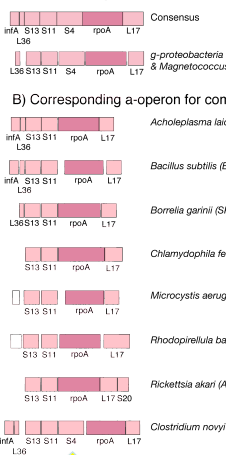
    # Draw the genomic context.
    drawMultipleGenomicContext $outputFilename $alignment $geneFiles $genePositions $geneStrands $contextDistance
}
    
```



Genome content future directions

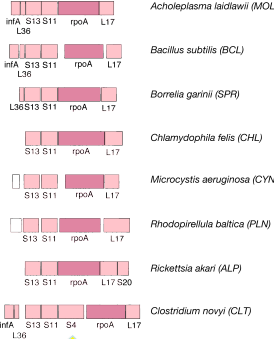
- Genome content still a work in progress.
- Good candidate for a GUI: combined phylogenetic tree/ genome content viewer.
- Can also use COG codes to color by gene function.
- Still need API for manipulating PTT files.

A) a-operon Organization

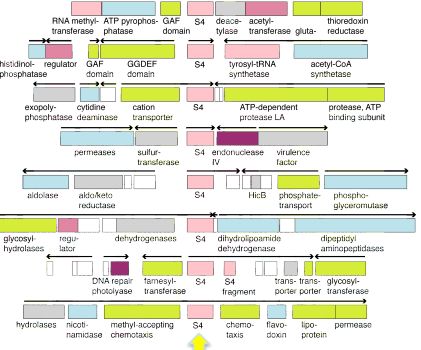


- Translation, ribosome structure & biogenesis
- Metabolism
- Replication, recombination and repair
- General functional prediction
- Cellular process
- Not annotated or no functional prediction

B) Corresponding a-operon for comparison



C) Outside-operon S4 context



Roberts, Chen, ZLS, BMC Evol. Bio. 2009

Genome content of ribosomal protein S4 by occurrence of the gene in the alpha operon.

Fifteen Clostridia genomes contain two copies of S4: one zinc-binding and one zinc-free.

BLAST DB Searching

- Import sequence data directly from BLAST databases
- Search using a single sequence or an **EP** profile
- Filter results based on taxonomy or redundancy (**QR**)

Name	E Score	410	420	430	
BYK_GLOVI	1e-19	N P Y P Y Y V E	R T H M A	G D L O	A K
666876	2e-19	I Q I C K I K S			
67520132	2e-19	N G E E V E V D			
23130288	3e-19	A D L A S G E E			
57159018	3e-19	M I D K V	Y C	A D V T	P E
1NGW	4e-19	H V L V R D L K			A
46195989	5e-19	H V L V R D L K			A
BYK_SVWYQ	5e-19	H D L S N G E E			
BYK_SVNEI	1e-18	K H L A A G E A			
BYK_STRMU	1e-18	D P F G K R F E	R T A T S	G Q L K E K Y A D K T K E E L H	
50256771	1e-18	E E V I D M P A			
57227974	1e-18	E E V I D M P A			
68179432	3e-18	I A A A L E G C E			
BYK_PHOMA	4e-18	I N G Q D R E I			
55738646	5e-18	D P F G K R F E	R T A T S	G Q L K E K Y A D K T K E E L H	
BYK_STRRG	5e-18	K Y A N L D K E		Q	L H
55820759	5e-18	D P F G K R F E	R T A T S	G Q L K E K Y A D K T K E E L H	
BYK_STRPN	6e-18	K Y A N L D K E		Q	L H
15500510	6e-18	K Y A N L D K E		Q	L H
82529807	5e-18	D P F G K R F E	R T A T S	G Q L K E K Y A D K T K E E L H	
BYK1_SALTI	6e-18	I E L E A L N I			
BYK_ENTFA	8e-18	I D N H T K E E			L S
56707357	8e-18	I L E E L D N K			

Filter Options

E Score: e-5

Redundancy Cutoff: 100

Superkingdom: **All**
Archaea
Bacteria
Eukaryota

Kingdom: **All**
Fungi
Metazoa
Viridiplantae

Phylum: **All**
Actinobacteria
Aquificae
Arthropoda
Ascomycota
Bacteroidetes
Chlamydiae

Apply Filter

Protein sequence alignment

How do I align two similar, but different sequences ?

Sequence 1: $a_1 a_2 a_3 - - a_4 a_5 \dots a_n$

Sequence 2: $c_1 - c_2 c_3 c_4 c_5 - \dots c_m$

There exist fast web tools, e.g., **BLAST** search: <http://www.ncbi.nlm.nih.gov/>
See also Blastn, Psi-Blast, ...

protein-protein BLAST

Nucleotide Protein Translations Retrieve results for an RID

[Search](#)

[Set subsequence](#) From: To:

[Choose database](#)

[Do CD-Search](#)

Now: **BLAST** or [Resequency](#) [Basicall](#)

Sequences from Swiss-Prot, NCBI, JGI,

Structures from PDB, CATH, SCOP,

[ExpASY Home page](#)
[Site Map](#)
[Search ExpASY](#)
[Contact us](#)
[Swiss-Prot](#)

Search for

NiceProt View of Swiss-Prot:

P47865

[\[Entry info\]](#)
[\[Name and origin\]](#)
[\[References\]](#)
[\[Comments\]](#)
[\[Cross-references\]](#)
[\[Keywords\]](#)
[\[Features\]](#)
[\[Sequence\]](#)
[\[Tools\]](#)

Note: most headings are clickable, even if they don't appear as links. They link to the user manual or other documents.

Entry information	
Entry name	AQP1_BOVIN
Primary accession number	P47865
Secondary accession numbers	None
Entered in Swiss-Prot in	Release 33, February 1996
Sequence was last modified in	Release 44, July 2004
Annotations were last modified in	Release 45, October 2004
Name and origin of the protein	
Protein name	Aquaporin-CHIP
Synonyms	Water channel protein for red blood cells and kidney proximal tubule Aquaporin 1 Water channel protein CHIP29
Gene name	Name: AQP1
From	Bos taurus (Bovine) [TaxID: 9913]
Taxonomy	Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Cetartiodactyla; Ruminantia; Pecora; Bovidae; Bovinae; Bos.
References	
[1] SEQUENCE FROM NUCLEIC ACID. TISSUE=Ocular ciliary epithelium;	Stearitz (Pro X)

Final Blast Result: Sequence Alignment

>gi|46395801|sp|Q88F17|AQPZ_PSEPK Aquaporin Z
 Length = 230

Score = 119 bits (299), Expect = 6e-27
 Identities = 70/186 (37%), Positives = 105/186 (56%), Gaps = 12/186 (6%)

Query: 53 VSLAFGLSIATLAQSVGHISGAHLNPAVTLGLLLSCQISVLRAIMYIIAQCVGAIIVATAI 112
 V+ AFGL++ T+A ++GHISG HLNPAV+ GL++ + + Y+IAQ +GAI+A +
 Sbjct: 40 VAFAPGLTVLTMFAIGHISGCHLNPAVSPFGLVVGGRRFPKELLPYVIAQVIGAILAAGV 99

Query: 113 LSGITSSLP--DNSLGL--NALAP---GVNSGQGLGIEIIGTLQLVLCVLATDRRRR 164
 + I S + S GL N A G G G E++ T ++ ++ TD R
 Sbjct: 100 IYLIASGKAGFELSAGLASNGYADHSPGGYTLGAGFVSEVVMAMFLVVMGATDARAP- 158







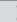



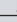









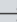

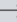

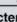
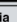




Query: 165 LGGSGPLAIGFSVALGHLLAIDYTGCGINPARSPFGSSVITHNF--QDHWIFWVGPFIGAA 222
 G P+AIG ++ L HL++I T +NPARS G ++ + Q W+FWV P IGAA
 Sbjct: 159 -AGFAPAIAGLALTLIHLISIPVTNTSVNPARSTGPALFVGGWALQQLWLFVWVAPLIGAA 217

Query: 223 LAVLIY 228
 + +Y
 Sbjct: 218 IGGALY 223

Search returns approximate alignments - needing refinement!
 Clustal, Muscle, MAFT, Tcoffee, pileup, Smith-Waterman, and
 manual editing in sequence editor

Flexible Grouping of Data

- Automatically group data by taxonomic classification to assist in evolutionary analysis (HGT) or create custom groups
- Apply metrics to groups independently, e.g bacterial signal

Sequence Name	90
Eukaryota:Fungi	
<input type="checkbox"/> 1asy_A   83	S R D S D R T G Q K R V K F V D
<input type="checkbox"/> 1eov_A   83	S R D S D R T G Q K R V K F V D
<input type="checkbox"/> SYDC_YEAST   82	S R D S D R T G Q K R V K F V D
Eukaryota:Metazoa	
<input type="checkbox"/> SYD_CAEL   57	S K . . . E K K V L N F L K V K E
<input type="checkbox"/> SYD_HUMAN   33	S Q . . . E K P D R V L V R V R D
<input type="checkbox"/> SYD_MOUSE   33	S Q . . . E K P D R V L V R V K D
Archaea:Crenarcha	
<input type="checkbox"/> SYD_AERPE   1 M L K D R F I A D
Archaea:Euryarchaeota	
<input type="checkbox"/> 1n9w_A   1 M R V L V R D
<input type="checkbox"/> 1b8a_A   1 M Y R T H Y S S E
<input type="checkbox"/> SYD_METMA   1 M S L A N L R T H Y T A D
<input type="checkbox"/> SYD_HALN1   1 M E N R T Y T A D
<input type="checkbox"/> SYD_THEAC   1 M L S I A E
<input type="checkbox"/> SYD_PYRHO   1 M I E K V Y C Q E
Bacteria:Proteobacteria	
<input type="checkbox"/> 1l0w_A   1 M R . R T H Y A G S
<input type="checkbox"/> 1l12_A   1 M . R T E Y C G Q

MultiSeq: Display and Edit Metadata

- External databases are **cross-referenced** to display **metadata** such as taxonomic information and enzymatic function
- Changes to metadata are preserved for future sessions
- **Electronic Notebook**: Notes and annotations about a specific sequence or structure can be added

Sequence Name:	SYDC_YEAST
Source Organism:	Saccharomyces cerevisiae
Common Name:	yeast
EC Number:	6.1.1.12
EC Description:	Aspartate--tRNA ligase.
Description:	Aspartate--tRNA synthetase, cytoplasmic (EC 6.1.1.12) (Aspartate--tRNA ligase) (AsPRS) - Saccharomyces cerevisiae (Baker's yeast).
Data Sources:	sp=P04802,SYDC_YEAST pdb=1EOVA
Lineage:	Eukaryota Fungi Ascomycota Saccharomycotina Saccharomycetes Saccharomycetales
Notes:	
<input type="button" value="OK"/> <input type="button" value="Cancel"/>	

Acknowledgements

- Elijah Roberts
- John Eargle
- Ke Chen
- Kirby Vandivort
- John Stone
- Michael Bach
- NIH Resource for Macromolecular Modeling and Bioinformatics

