

Introduction to Evolutionary Concepts and VMD/MultiSeq - Part I

Zaida (Zan) Luthey-Schulten

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

NIH Workshop 2010

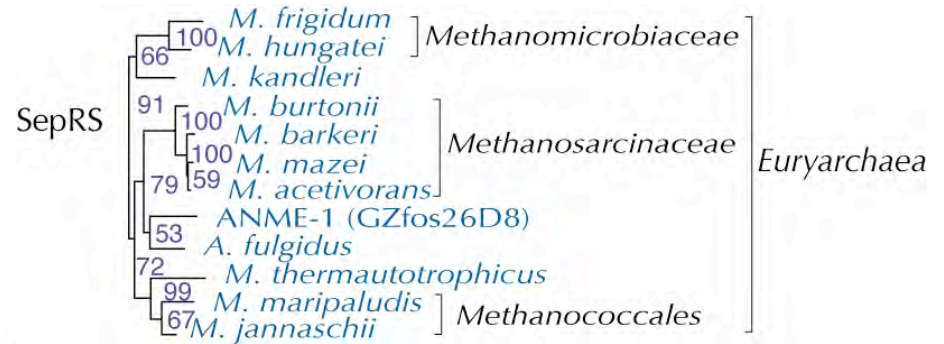
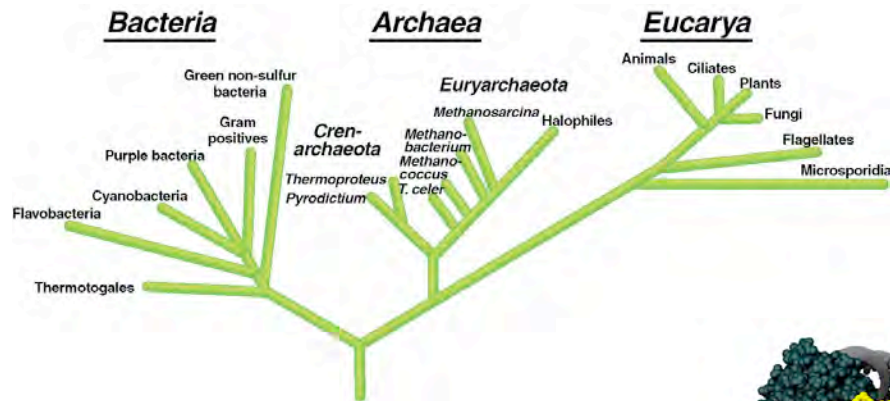


ILLINOIS

UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

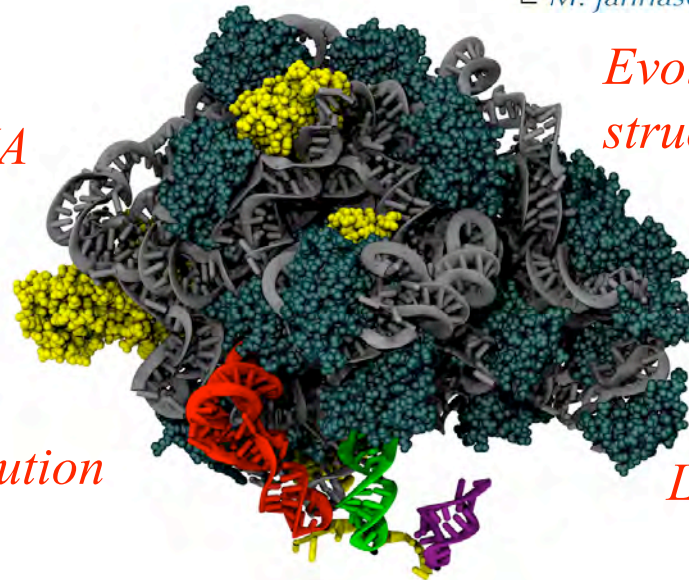
VMD/MultiSeq - “A Tool to Think”

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.*”



UPT - Woese 16S rRNA

Evolutionary profiles for protein structure & function prediction



Signatures ribosomal evolution

LSU (23S rRNA + rproteins)

New Tools in VMD/MultiSeq

Protein / RNA
Sequence Data

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

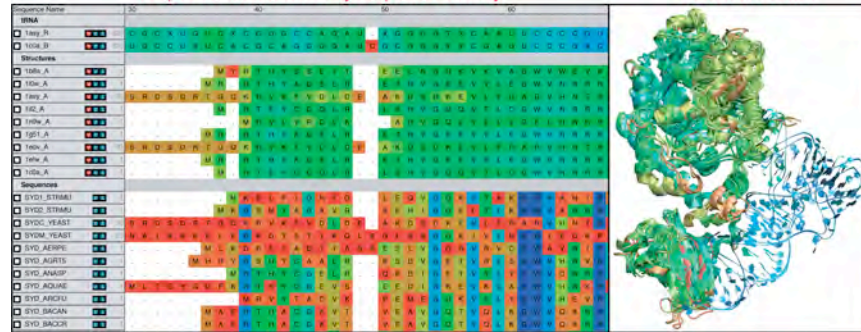
Metadata Information,
Clustal &
Phylogenetic Trees

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

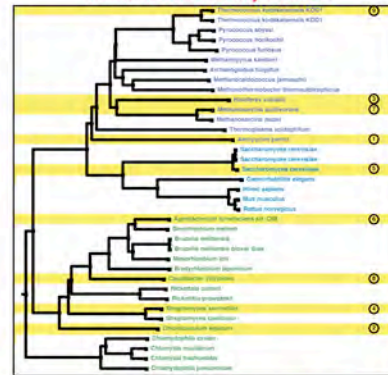
Sequence /Structure
Alignment

Protein & RNA
secondary structure

Group data by taxonomic classification

Sequence Name	50
Eukaryota:Fungi	
1asy_A	S R R D S D R T G Q K R V K F V D
1aov_A	S R R D S D R T G Q K R V K F V D
SYDC_YEAST	S R R D S D R T G Q K R V K F V D
Eukaryota:Metazoa	
SYD_CAEL	G L V N S K E K K V L N F L K V
SYD_HUMAN	S M I Q S Q E K P D R V L V R V
SYD_MOUSE	S M I Q S Q E K P D R V L V R V
Archaea:Crenarcha	
SYD_AERPE M L K D R F I A D I
Archaea:Euryarchaeota	
1n9w_A M R V L V R D
1b8a_A M Y R T H Y S S E
SYD_METMA M S L A N L R T H Y T A D
SYD_HALN1 M L E R T Y I E D
SYD_THEAC M P R T Y I D T
SYD_PYRHO M L R T H Y S N E
Bacteria:Proteobacteria	
110w_A M R R R T H Y A G S
112_A M R T E Y C G G O

View sequence or structure phylogenies and
eliminate redundancy with QR



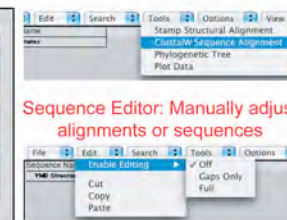
QR non-redundant
seq / str sets

Cluster analysis
/ Bioinformatics
scripting

Import data directly from BLAST databases



Align sequences with Clustal

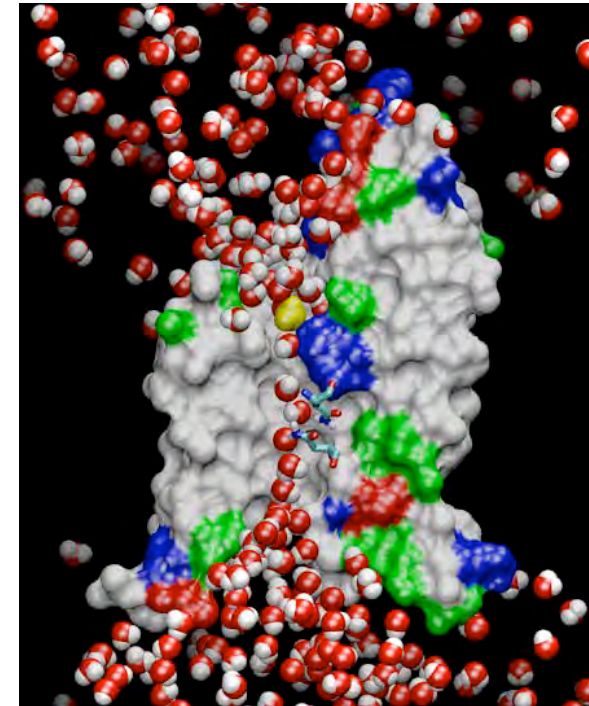
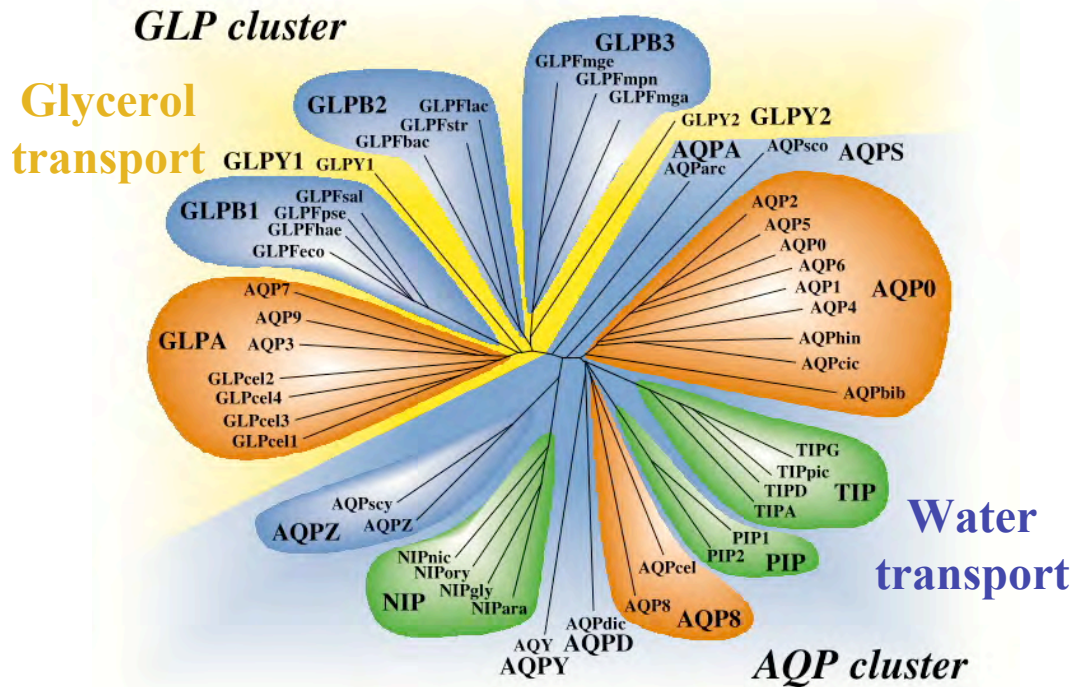


Tutorials
MultiSeq/AARS
EF-Tu/Ribosome

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)

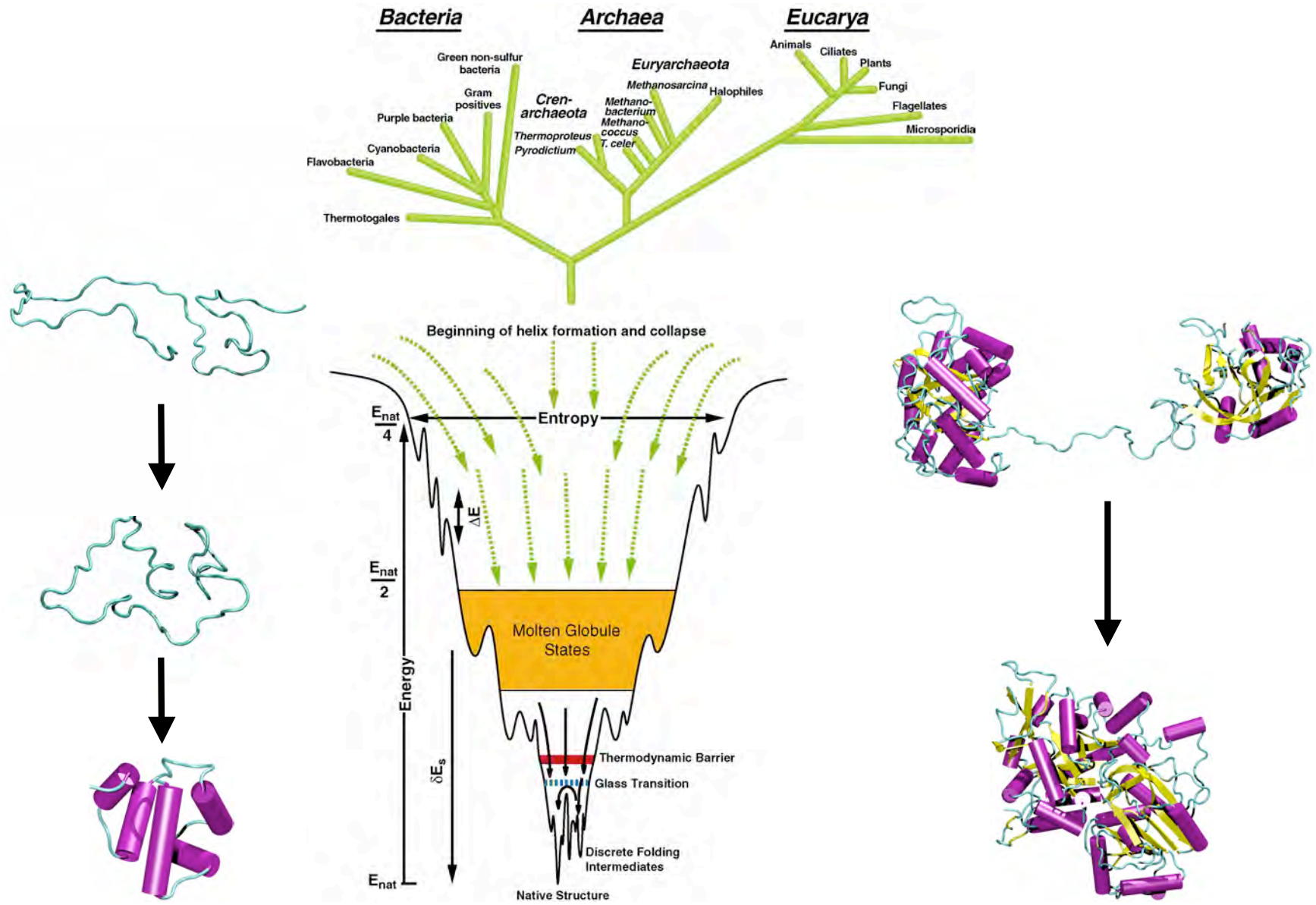
Aquaporin Superfamily: Bacterial & Eucaryal



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

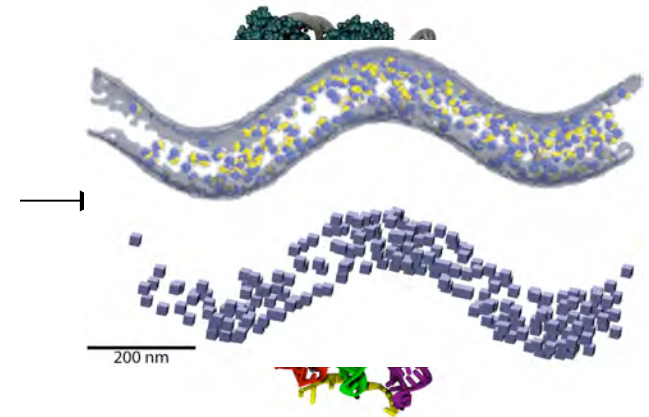
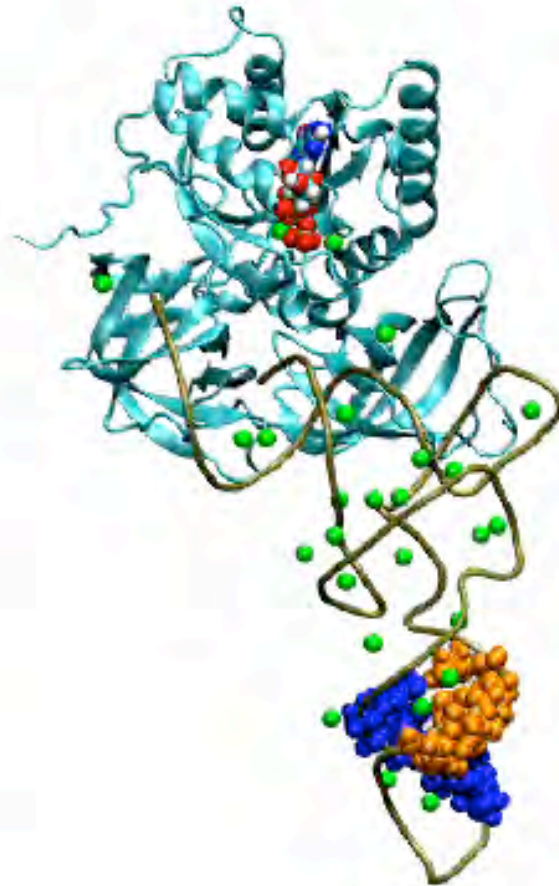
AQP0_HUMAN	---	LNTLHPAVSVGQATTVEIFLTLQFVLCIFATYDE	-RRNGQLGSVALAVGFSLALGHLFGMYT	GAGM	183
AQP1_HUMAN	---	RNDLADGVNSGQGLGIEIIGTLQVLVLCVLATDR	-RRRDLGGSAPLAIGLSVALGHLAIDYT	GCGI	191
AQP2_HUMAN	---	VNALSNSTTAGQAVTVLFLTLQVLVLCIFASTDE	-RRGENPGTPALSIGFSVALGHLGHIHYT	GCSM	183
AQP3_HUMAN		GIFATVPSGHLDMINGFFDQFIGTASLIVCVLAI	VDPNNPVPRGLEAFTVGLVVLVIGTSMGFNS	GYAV	214
AQP4_HUMAN	---	VTMVHGNLTAGHLLVELIITFQLVFTIFASCDS	-KRTDVTGSIALAIGFSVAIGHLFAINYTG	GASM	212
AQP5_HUMAN	---	VNALNNNTTQQAAMVVELILTFQLALCIFASTDS	-RRTSPVGSVALSIGLSVTLGHLVGIYFT	GCSM	184
AQP6_HUMAN	---	INVVRNSVSTGQAVAVELLLTLQVLVLCVFASTDS	-RQTS--GSPATMIGISWALGHLIGILFT	GCSM	195
AQP7_HUMAN		GIFATYLPDHMTLWRGFLNEAWLTGMLQLCLFA	ITDQENNPALPGTEALVIGILVVIIGVSL	GMNTGYAI	225
AQP8_HUMAN	-	AAFVTVQEQGQVAGALVAEIIITLLALAVCMGAIN	--EKTGKPLAPFSIGFAVTVDILAGGPV	SGGCM	209
AQP9_HUMAN		HIFATYPAPYLSLANAFADQVVATMILLIIVFA	IFDSRNLGAPRGLEPIAIGLLIIVIAS	SLGLNSGCAM	215
GLPF_ECOLI		GTFSTYPNPHINFVQAFVEMVITAILMGLLILAL	TDDGNGVPRGPLAPLLIGLLIAVIGAS	MGPLTGFM	202
ruler	...	180.....190.....200.....210.....220.....230.....240....			

Protein (RNA) Folding, Structure, & Function



Protein:RNA Complexes in Translation

Evolution, Structure, and Dynamics



Proteins/RNA
Polyribosomes
Ribosome

“*Evolution SepRS/CysRS*”, *PNAS*, 2005

“*Dynamic Signaling Network*”, *PNAS* 2009

“*Exit Strategy Charged tRNA*” *JMB* 2010

“*Dynamical Recognition Novel Amino Acids*”, *JMB* 2008

“*Signatures ribosomal evolution*”
PNAS 2008, *BMC* 2009

“*Whole cell simulations on GPUs*” *IEEE* 2009

“*Dynamics of tRNA*” *FEBS* 2010

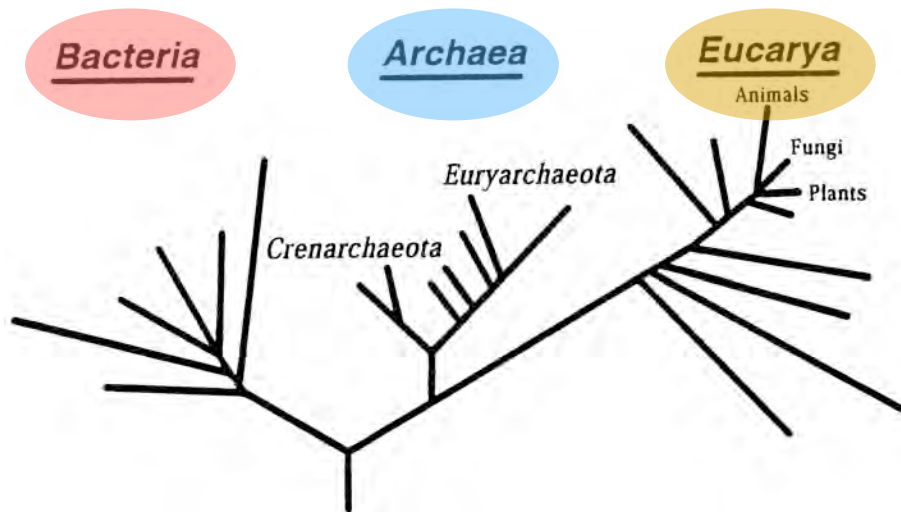
Basic principles of evolutionary analysis for proteins & RNAs

- 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

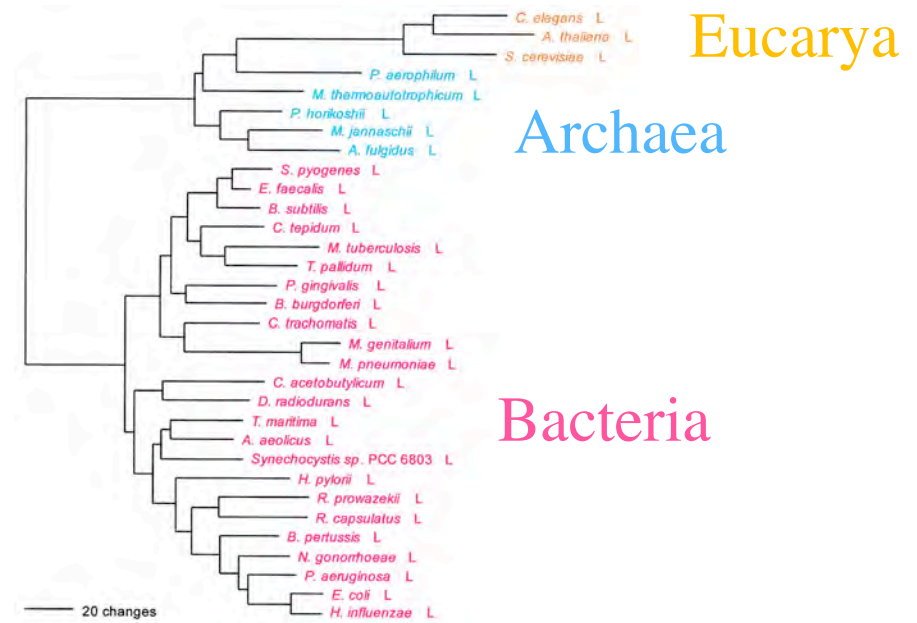
*Various models of evolutionary change

Universal Phylogenetic Tree

3 domains of life

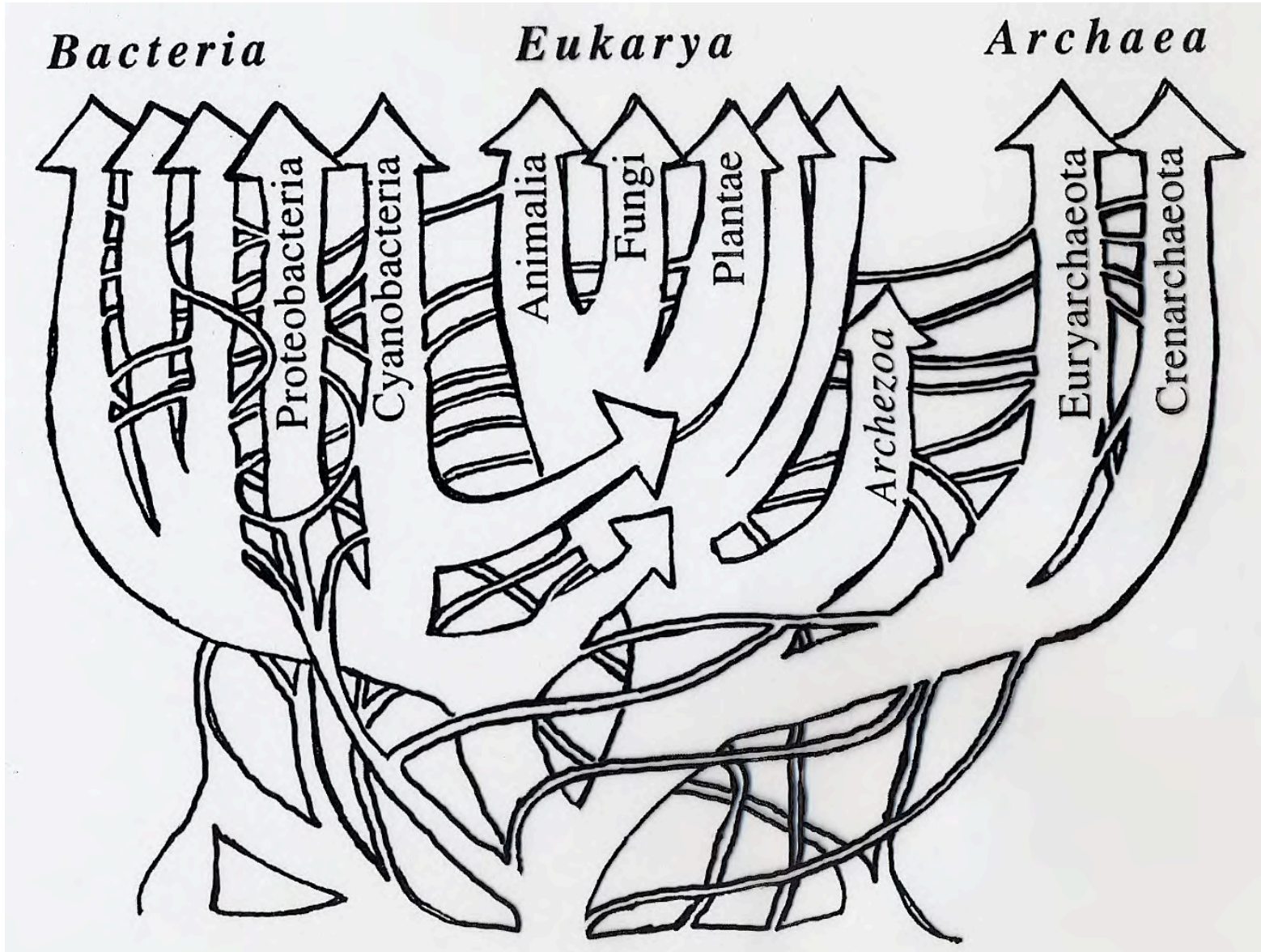


Reference 16S rRNA tree



Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution.

Look for horizontal gene transfer events



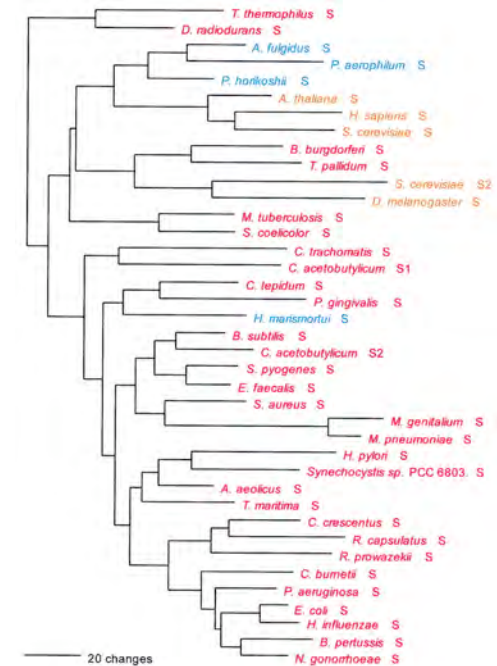
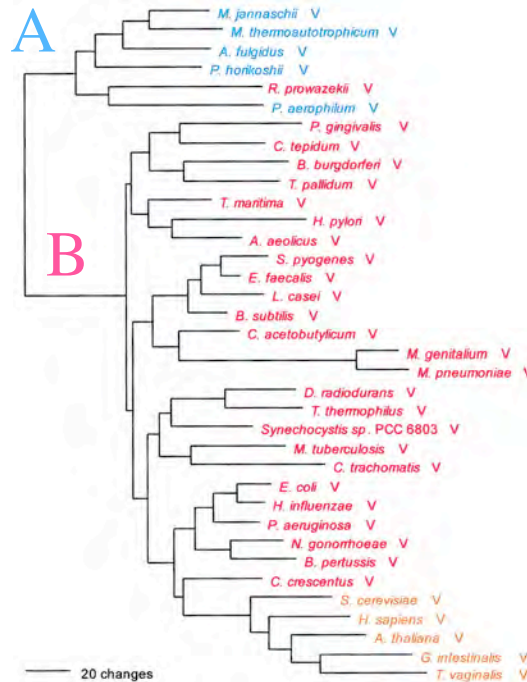
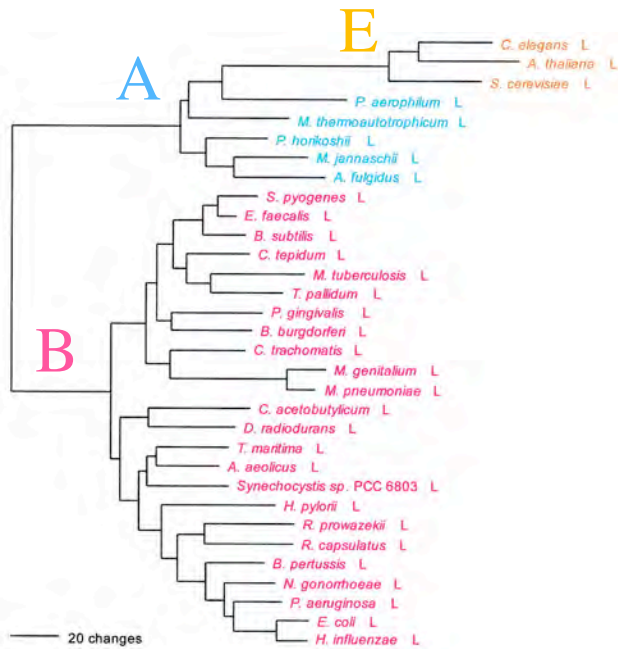
After W. Doolittle, modified by G. Olsen

Phylogenetic Distributions

Full Canonical

Basal Canonical

Non-canonical



increasing inter-domain of life Horizontal Gene Transfer

“HGT erodes the historical trace, but does not completely erase it....” G. Olsen

Woese, Olsen, Ibba, Soll *MMBR* 2000

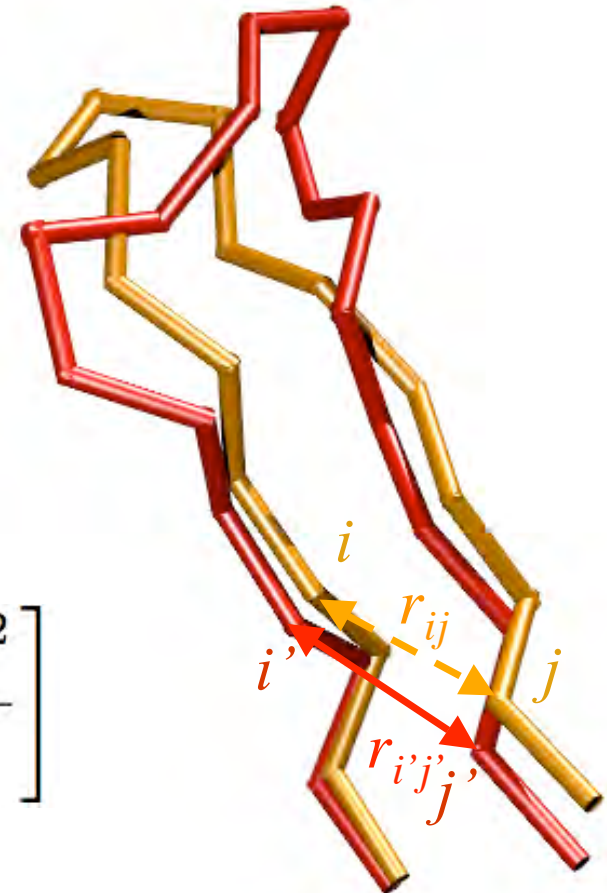
Protein Structure Similarity Measure

Q_H Structural Homology

fraction of native contacts for aligned residues +
presence and perturbation of gaps

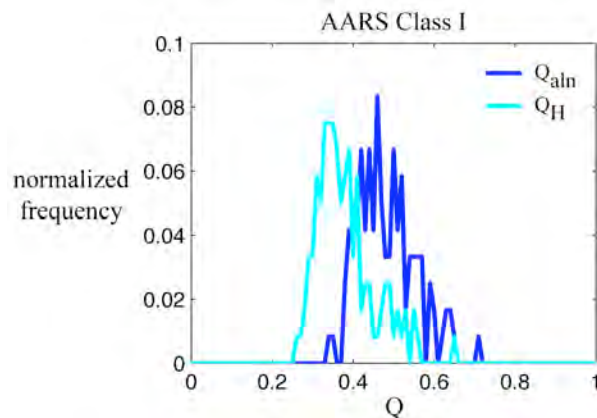
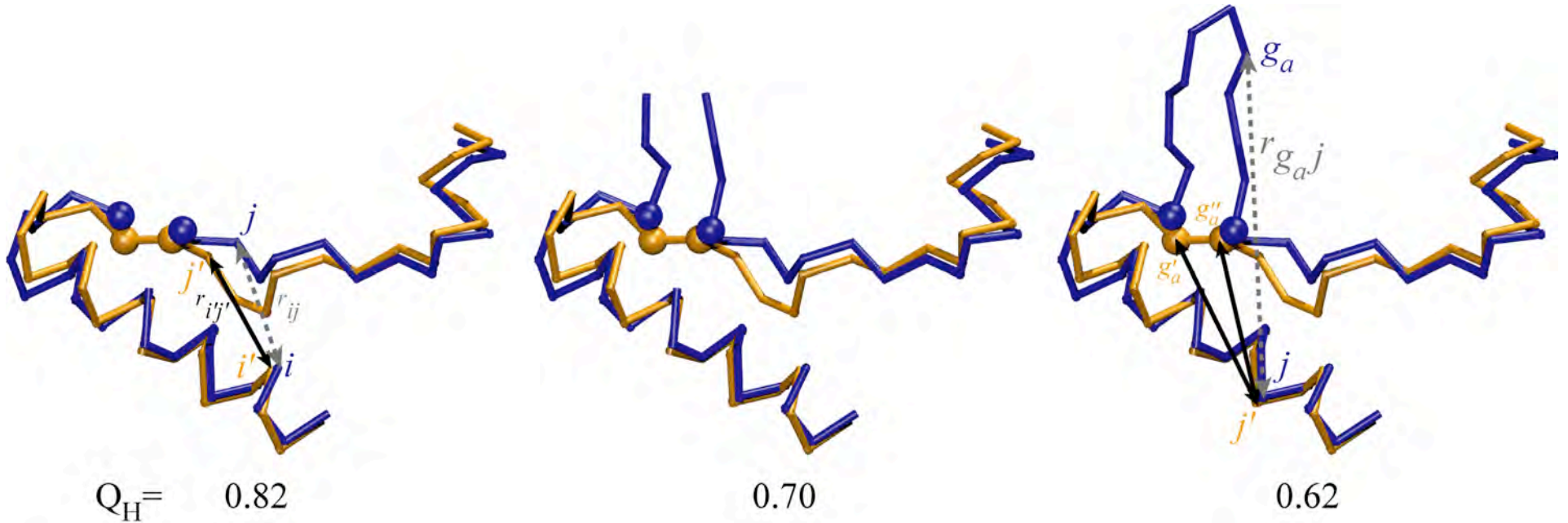
$$Q_H = N [q_{aln} + q_{gap}]$$

$$q_{aln} = \sum_{i < j-2} \exp \left[-\frac{(r_{ij} - r_{i'j'})^2}{2\sigma_{ij}^2} \right]$$



Structural Similarity Measure the effect of insertions

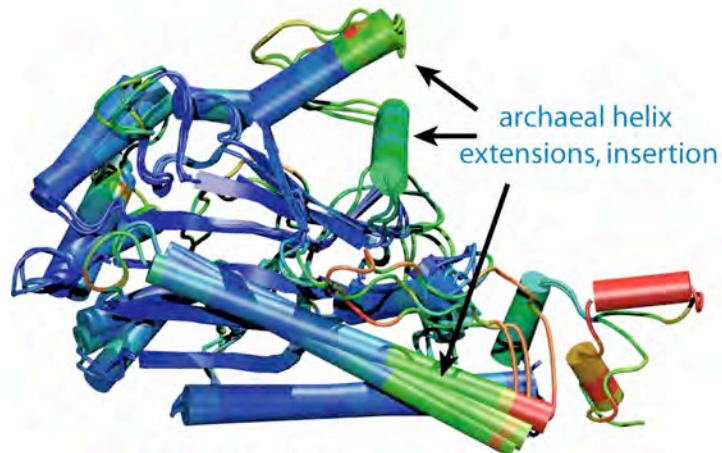
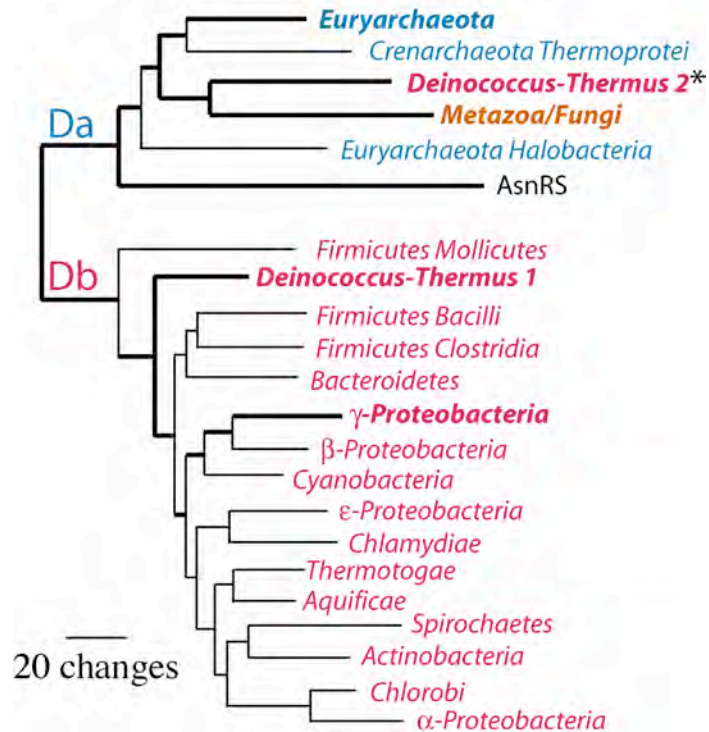
“Gaps should count as a character but not dominate” C. Woese



$$q_{gap} = \sum_{g_a} \sum_j^{N_{aln}} \max \left\{ \exp \left[-\frac{(r_{g_a j} - r_{g_a' j'})^2}{2\sigma_{g_a j}^2} \right], \exp \left[-\frac{(r_{g_a j} - r_{g_a'' j'})^2}{2\sigma_{g_a j}^2} \right] \right\} + \sum_{g_b} \sum_j^{N_{aln}} \max \left\{ \exp \left[-\frac{(r_{g_b j} - r_{g_b' j'})^2}{2\sigma_{g_b j}^2} \right], \exp \left[-\frac{(r_{g_b j} - r_{g_b'' j'})^2}{2\sigma_{g_b j}^2} \right] \right\}$$

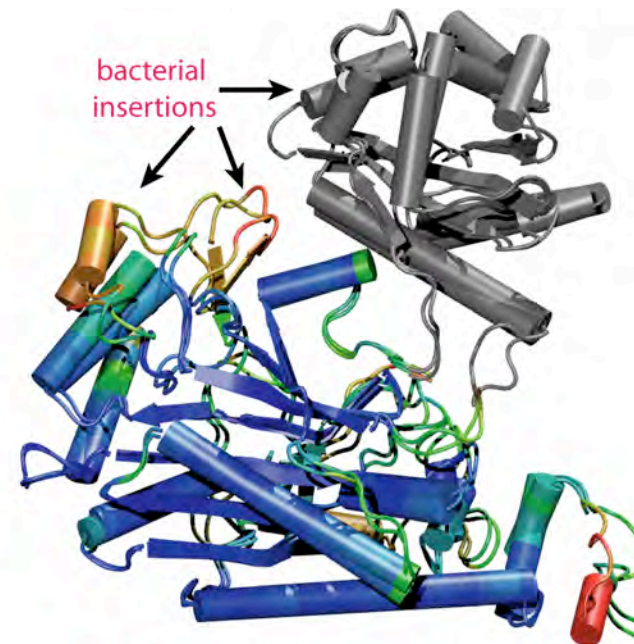
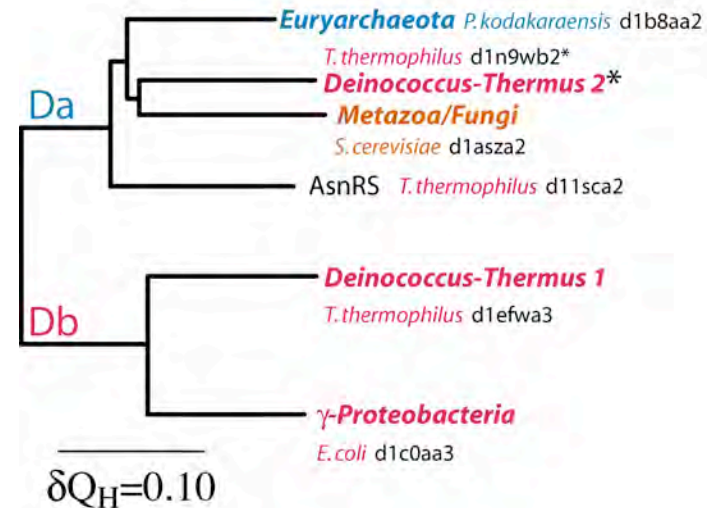
Structure encodes evolutionary information!

sequence-based phylogeny



Da - AspRS archaeal genre

structure-based phylogeny

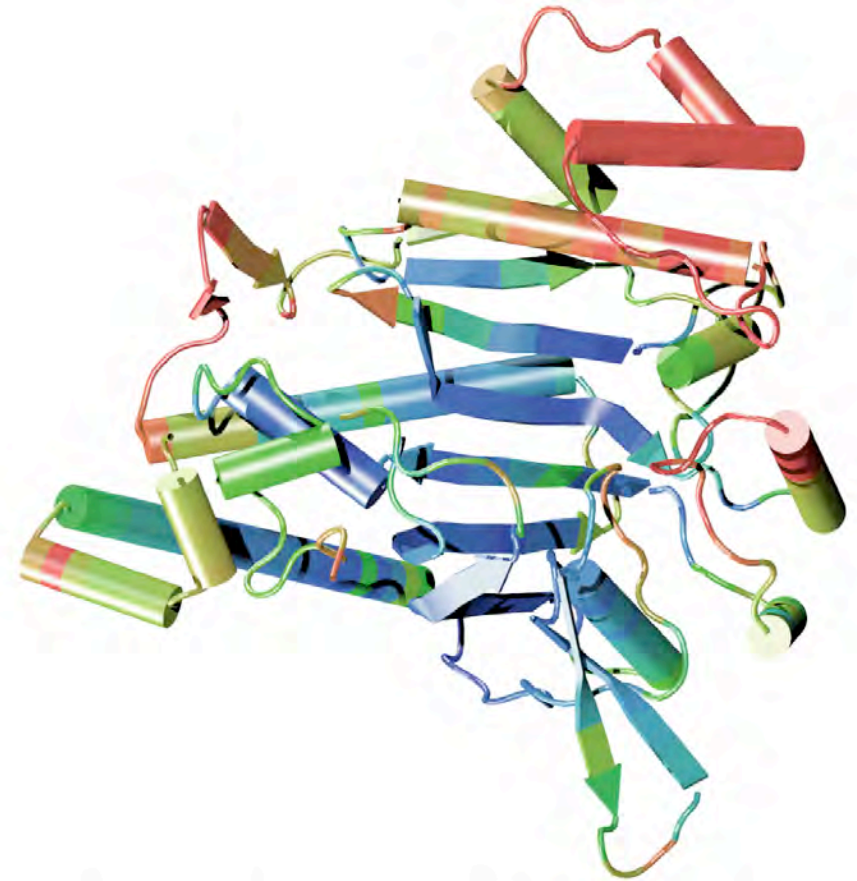
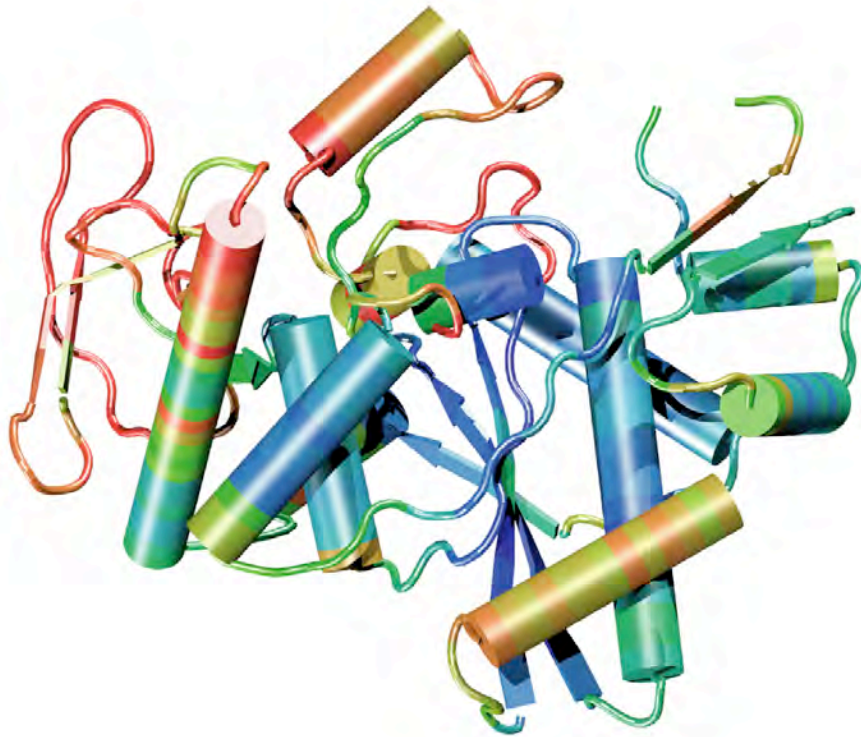
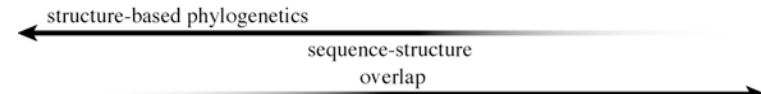
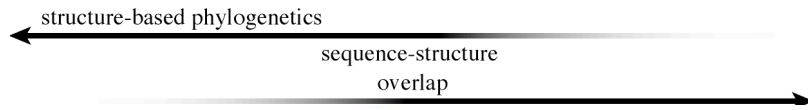


Db - AspRS bacterial genre

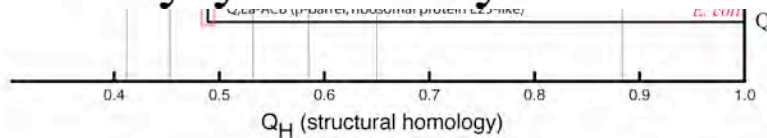
Structure reveals distant evolutionary events

Class I AARSs

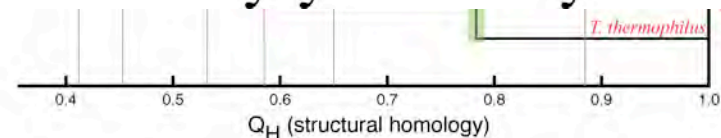
Class II AARSs



Class I Lysyl-tRNA Synthetase



Class II Lysyl-tRNA Synthetase



Sequences define more recent evolutionary events:



Conformational changes
in the same protein.

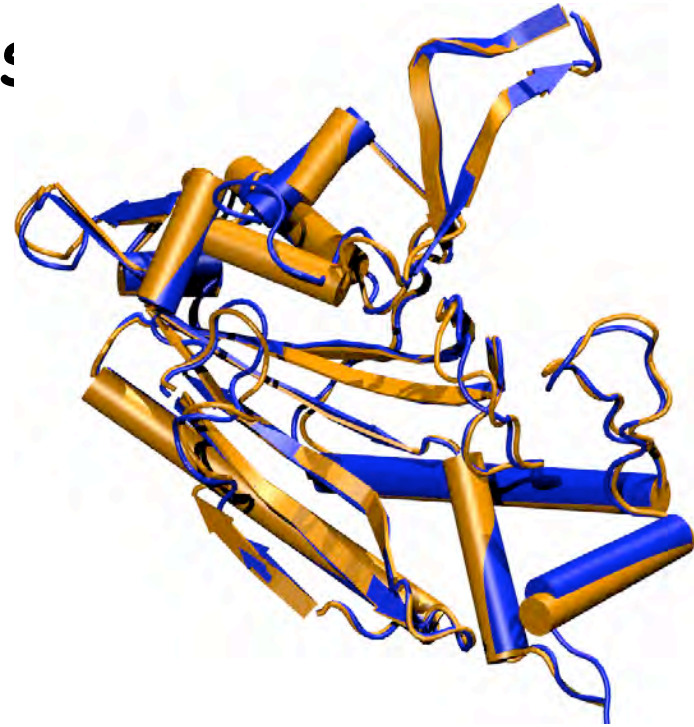
ThrRS

T-AMP analog, 1.55 Å.

T, 2.00 Å.

$Q_H = 0.80$

Sequence identity = 1.00



Structures for two
different species.

ProRS

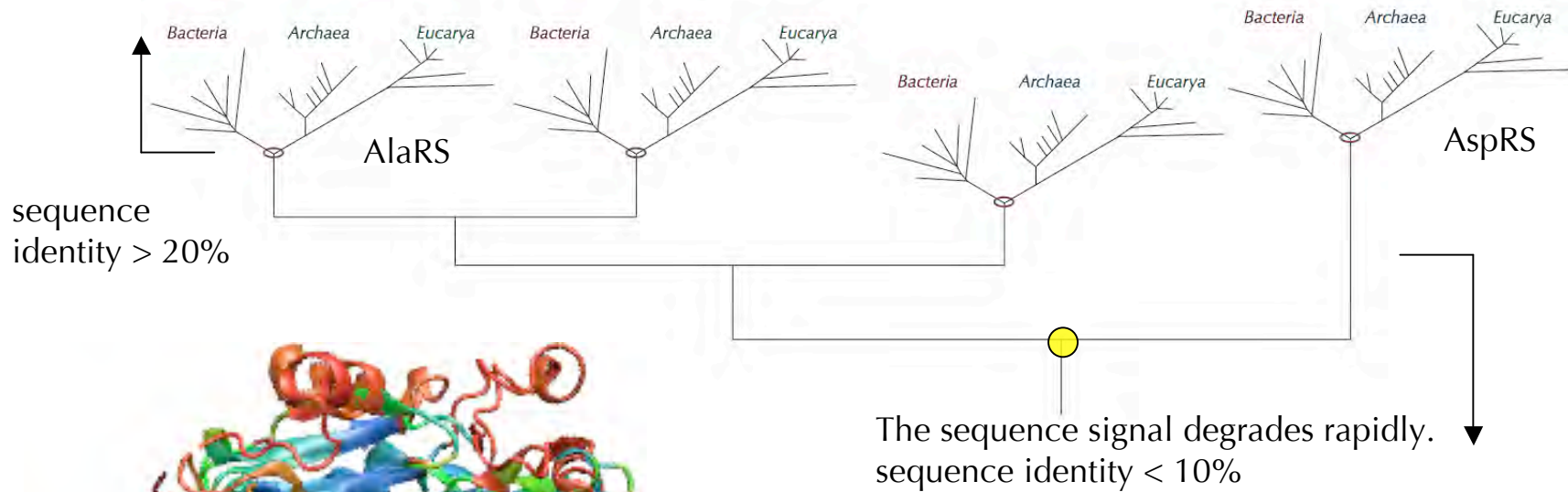
M. jannaschii, 2.55 Å.

M. thermoautotrophicus, 3.20 Å.

$Q_H = 0.89$

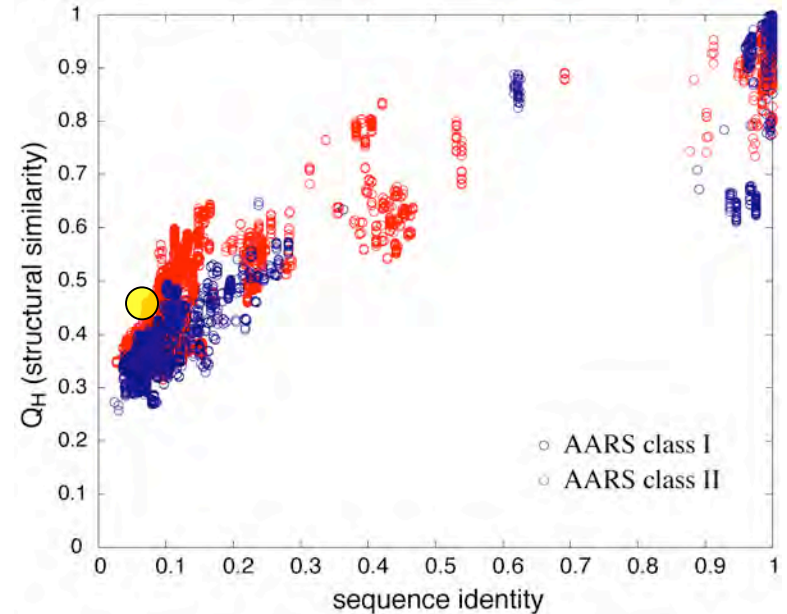
Sequence identity = 0.69

Relationship Between Sequence & Structure



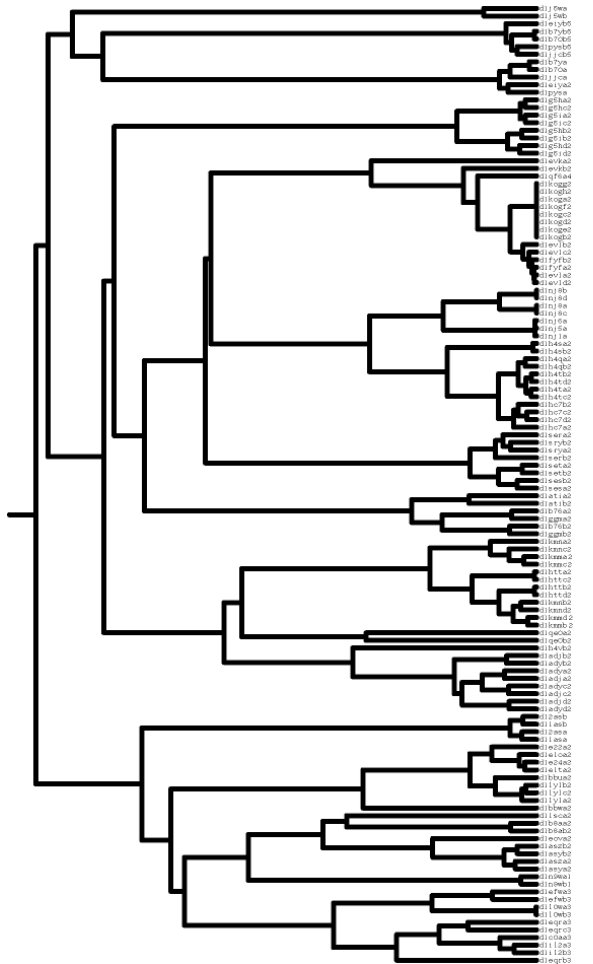
Structural superposition of AlaRS & AspRS.

● Sequence id = 0.055, $Q_H = 0.48$

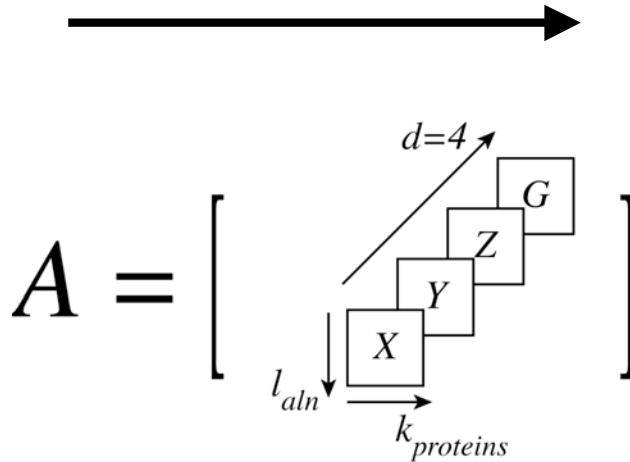


Non-redundant Representative Profiles

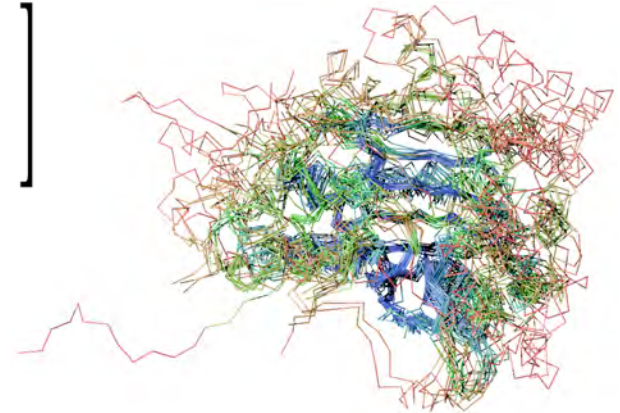
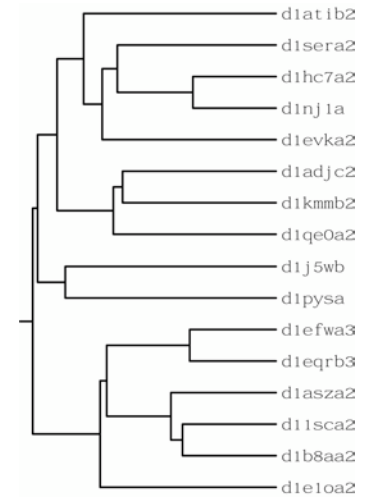
Too much information
129 Structures



Multidimensional QR
factorization
of alignment matrix, A .



Economy of information
16 representatives



QR computes a set of maximal linearly independent structures.

P. O'Donoghue and Z. Luthey-Schulten (2003) *MMBR* **67**:550-571.

P. O'Donoghue and Z. Luthey-Schulten (2005) *J. Mol. Biol.*, **346**, 875-894.

Numerical Encoding of Proteins in a Multiple Alignment

Encoding Structure

Rotated Cartesian + Gap = 4-space

Aligned position $(x_{C_\alpha}, y_{C_\alpha}, z_{C_\alpha}, 0)$

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)

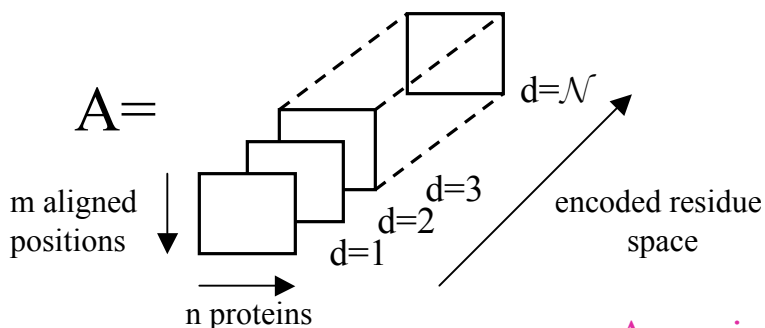
B = (0,1,0)

C = (0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0)

...

GAP = (0,1)

Alignment is a Matrix with Linearly Dependent Columns

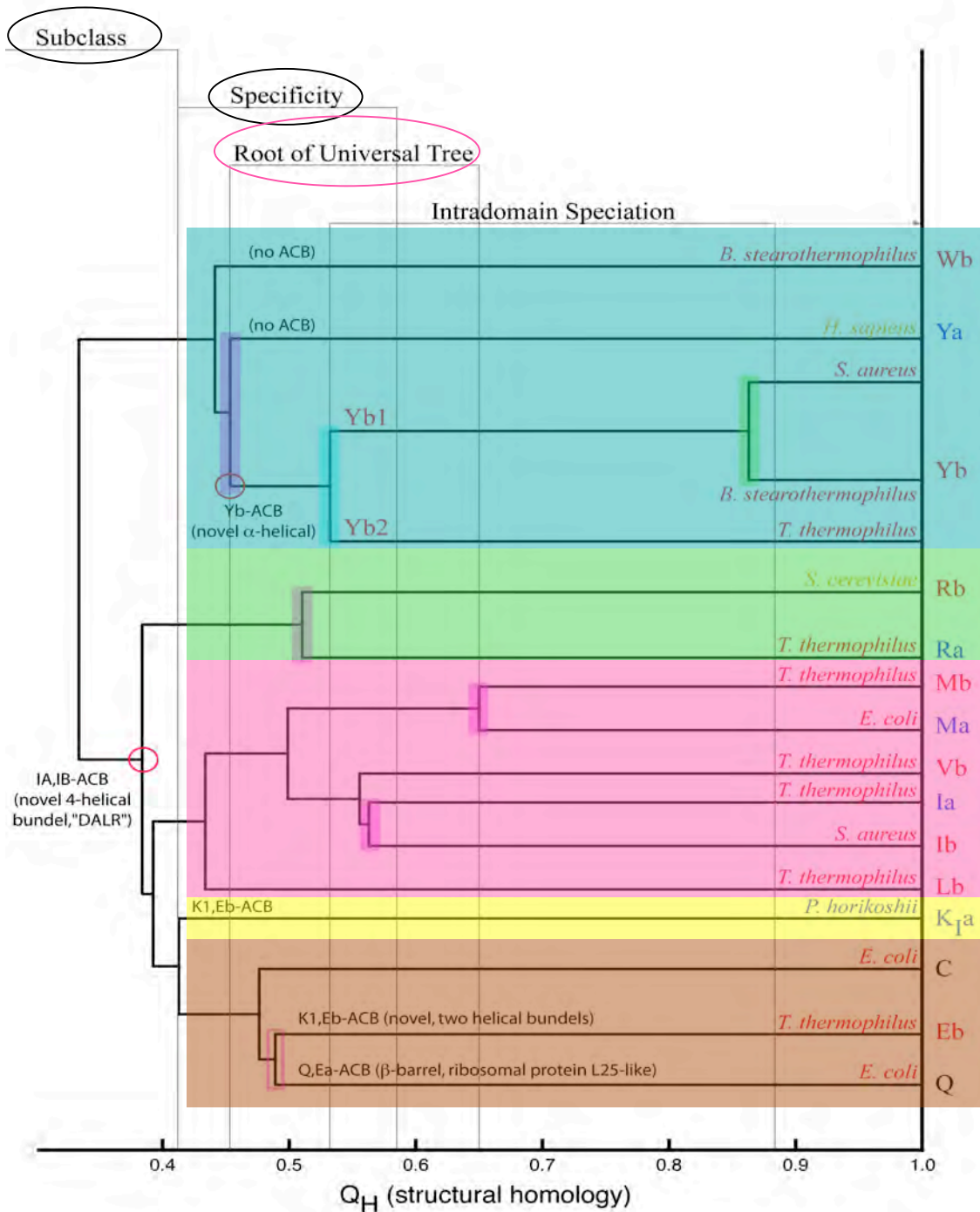


$$Q_{(d)}^T A_{(d)} P = Q_{(d)}^T \begin{bmatrix} d=1 & & & & & \\ & d=2 & & & & \\ & & d=3 & & & \\ & & & d=4 & & \\ & & & & d=5 & \\ & & & & & d=N \end{bmatrix} P = \tilde{R}_{(d)}$$

The diagram shows the matrix equation $Q_{(d)}^T A_{(d)} P = Q_{(d)}^T \dots P = \tilde{R}_{(d)}$. The matrix A is transformed into a form where the columns are linearly dependent. The axes are labeled d=1, d=2, d=3, d=4, d=5, and d=N. The matrix A is represented as a grid of boxes with labels X, Y, Z, and G. The dimensions are m_{aln} and n_{proteins}.

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

Class I AARSs evolutionary events

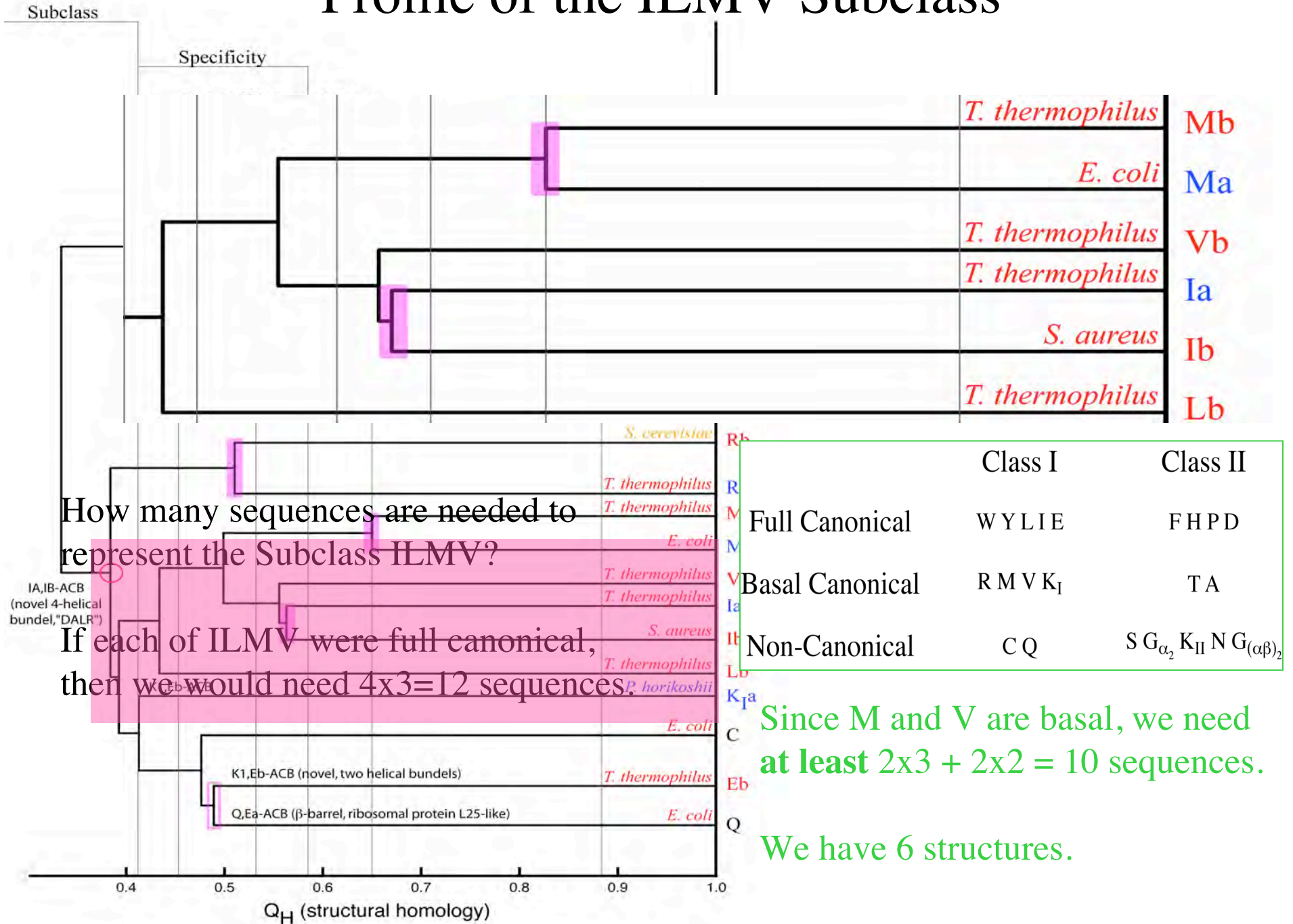


5 Subclasses

Specificity - 11 Amino acids

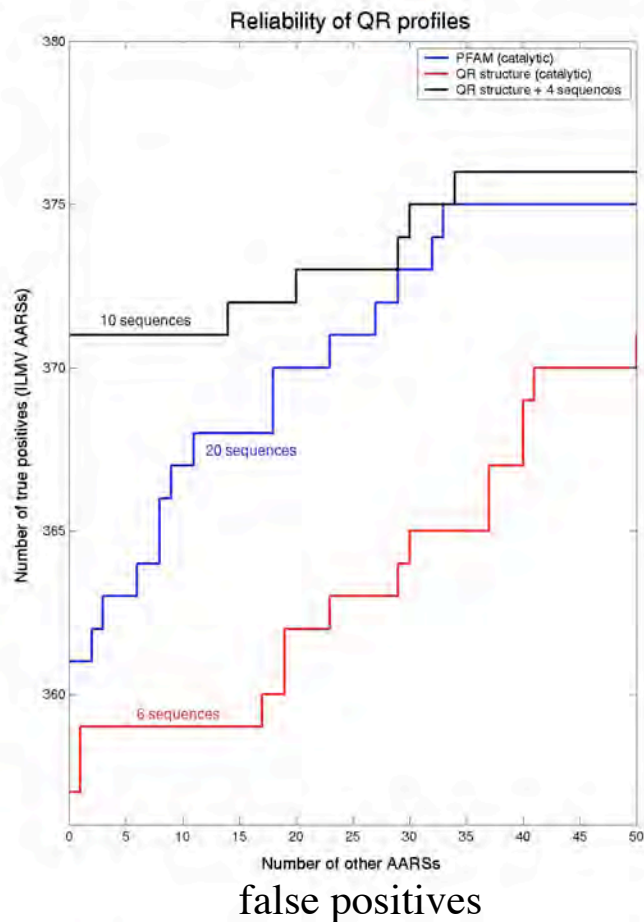
Domain of life A,B,E

Profile of the ILMV Subclass

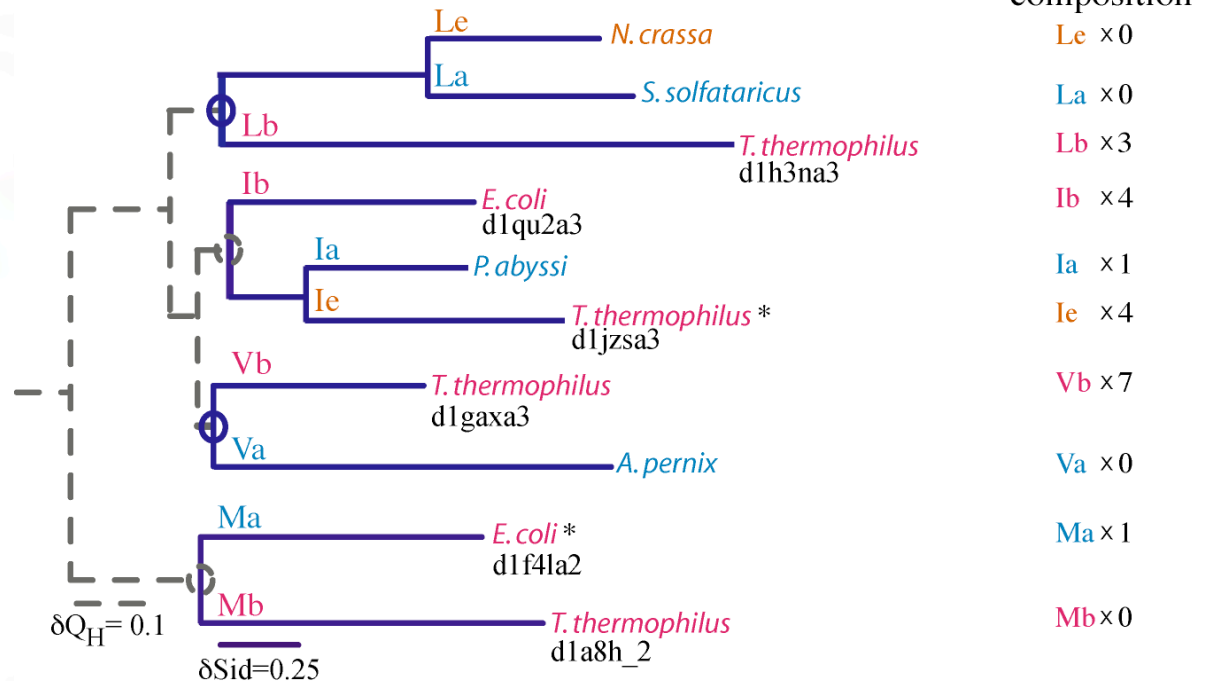


Evolutionary Profiles for Homology Recognition

AARS Subclass ILMV

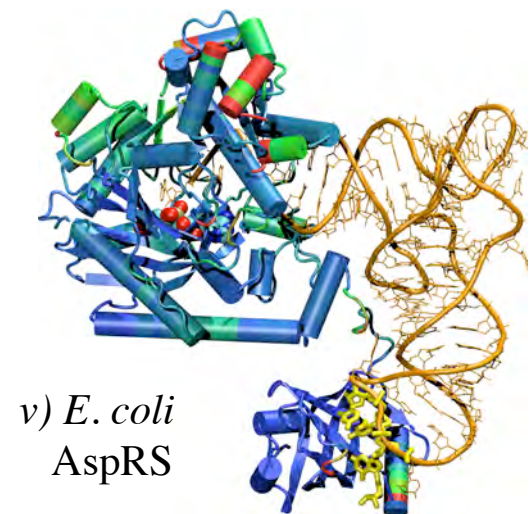
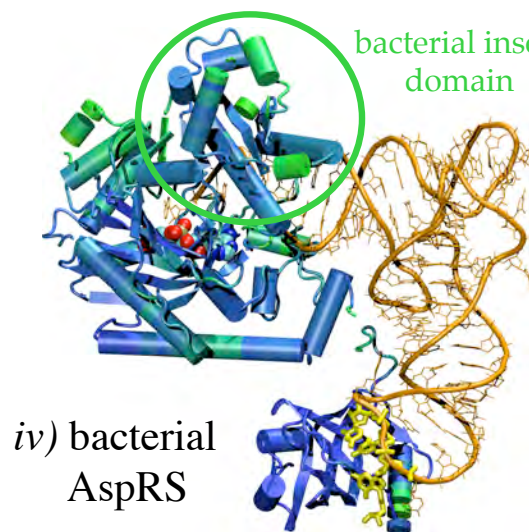
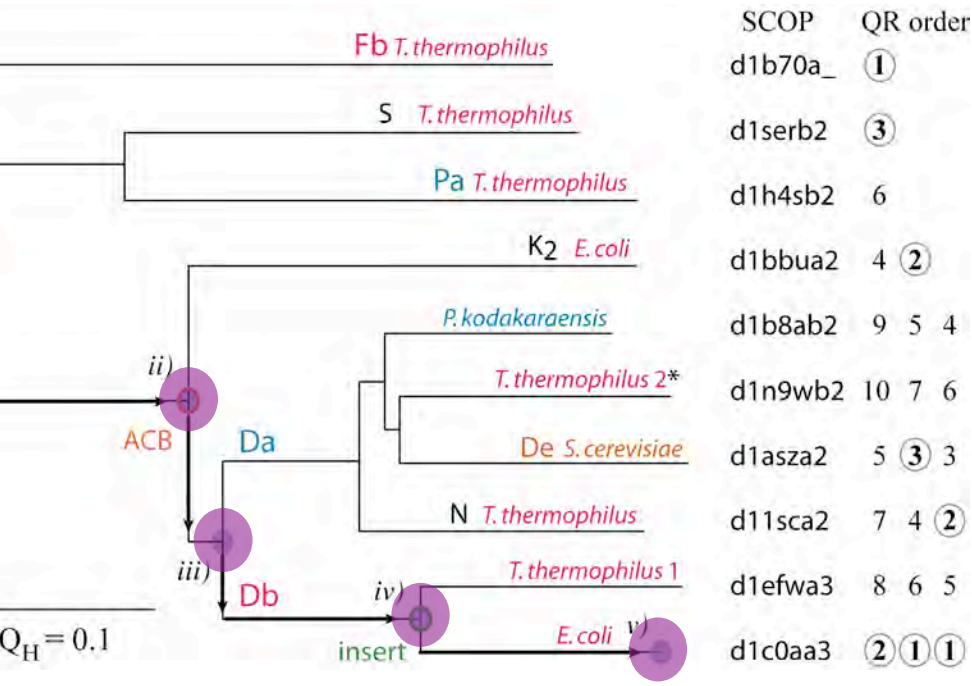
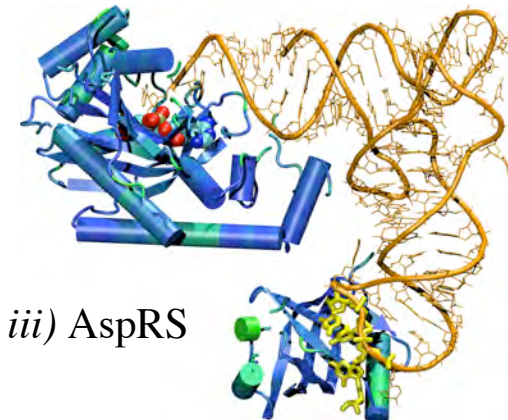
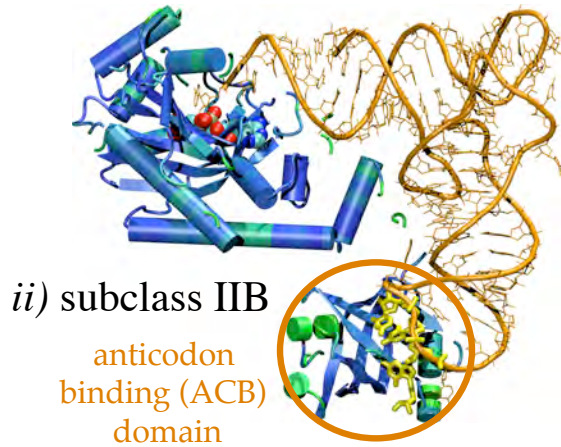
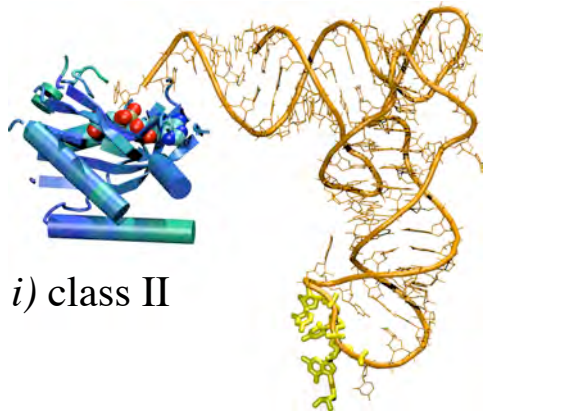


Combined Structure-Sequence Phylogeny
 an evolutionary profile of the AARS subclass IA

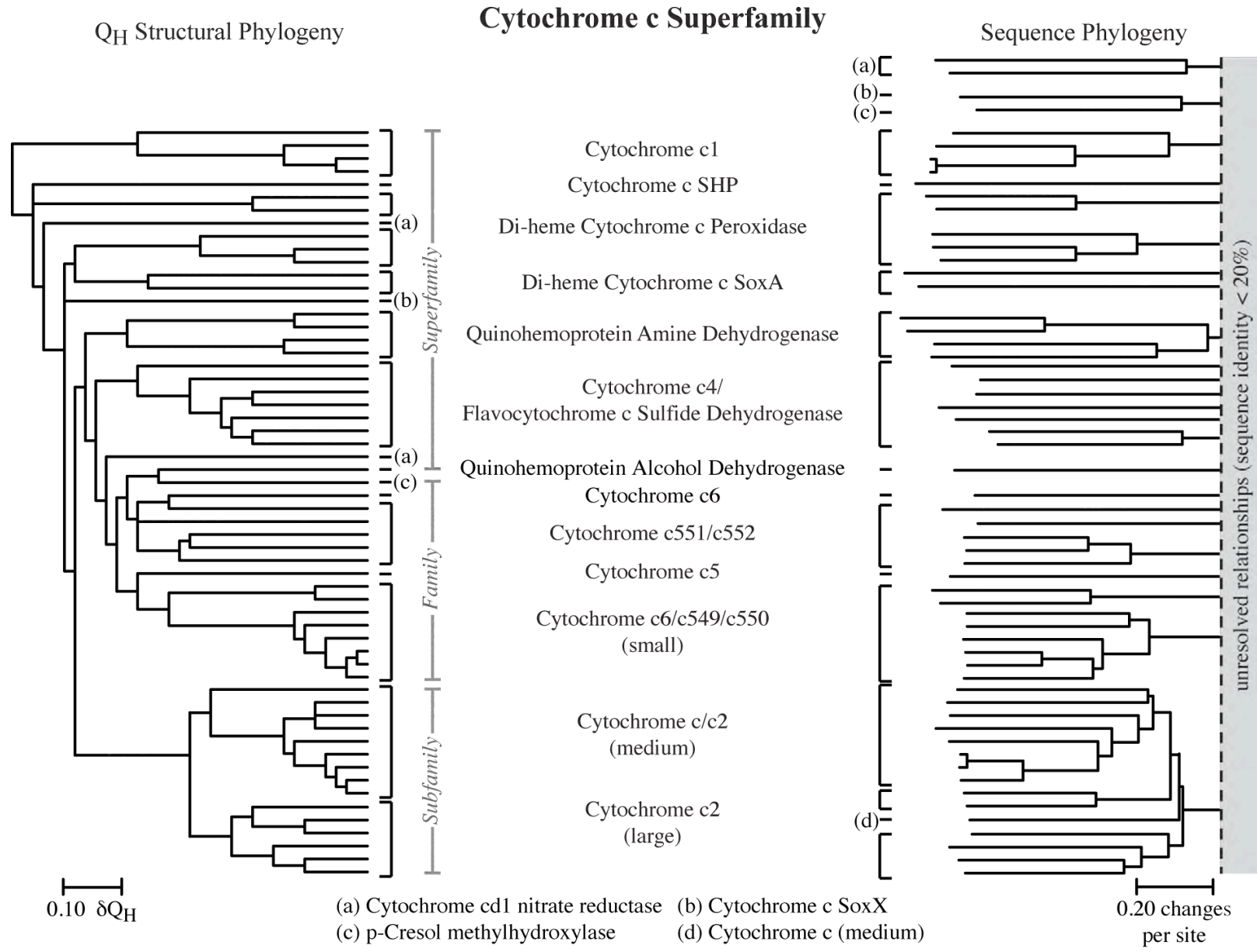


The composition of the profile matters.
 Choosing the right 10 sequence makes all the difference.

Design - Evolution of Structure and Function in Class II



Superfamily structural analysis



Structural Profiles

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

Add more structural information

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

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.

Multiple Structural Alignments

STAMP – cont'd

2. Refine Initial Alignment & Produce Multiple Structural Alignment

Alignment score:

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

$$S_P = \sum_{aln.path} P_{ij}$$

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

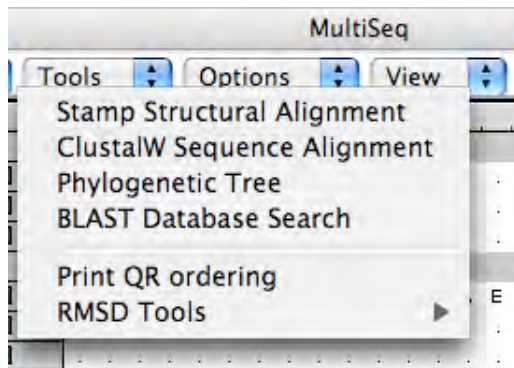
i_A, i_B -- length of gaps in A and B.

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.

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, BLAST, **STAMP, Signatures, Mutual information, QR, PT,.....**)

New Tools in VMD/MultiSeq

Protein / RNA
Sequence Data

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

Metadata Information,
Clustal &
Phylogenetic Trees

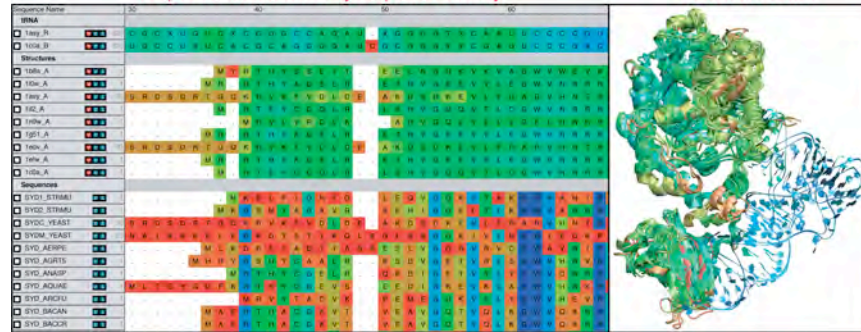
RAXml Trees,
Genomic Content,
Temperature DB

Blast & PsiBlast

Sequence Editor

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

Synchronization between
1D and 3D views



Sequence /Structure
Alignment

Protein & RNA
secondary structure

Group data by taxonomic classification

View sequence or structure phylogenies and
eliminate redundancy with QR

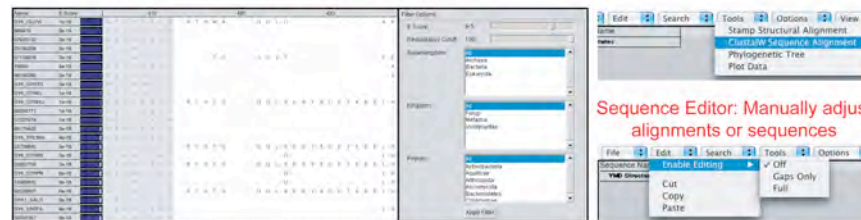


QR non-redundant
seq / str sets

Cluster analysis
/ Bioinformatics
scripting

Import data directly from BLAST databases

Align sequences with Clustal



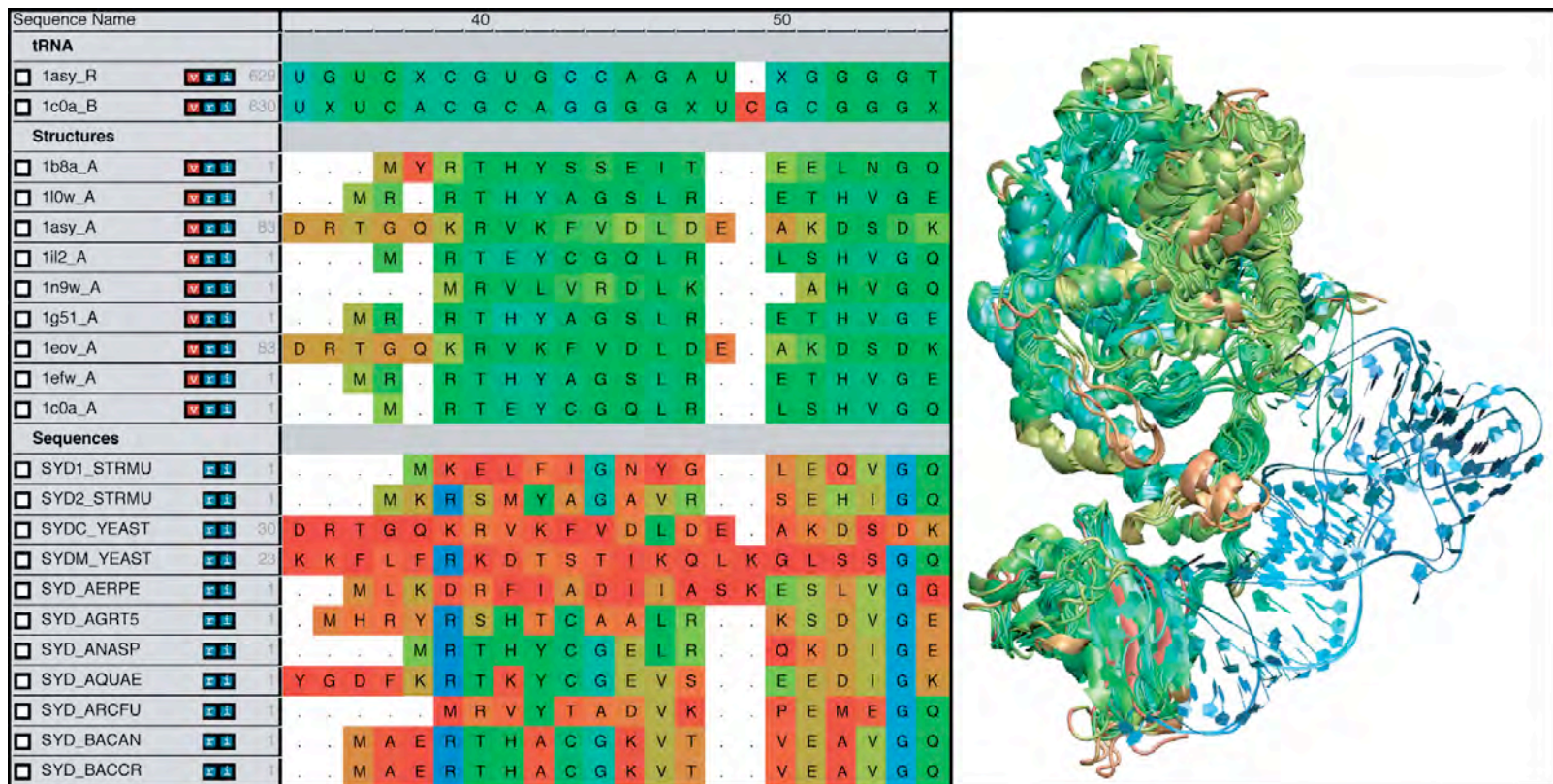
Tutorials
MultiSeq/AARS
EF-Tu/Ribosome

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



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)

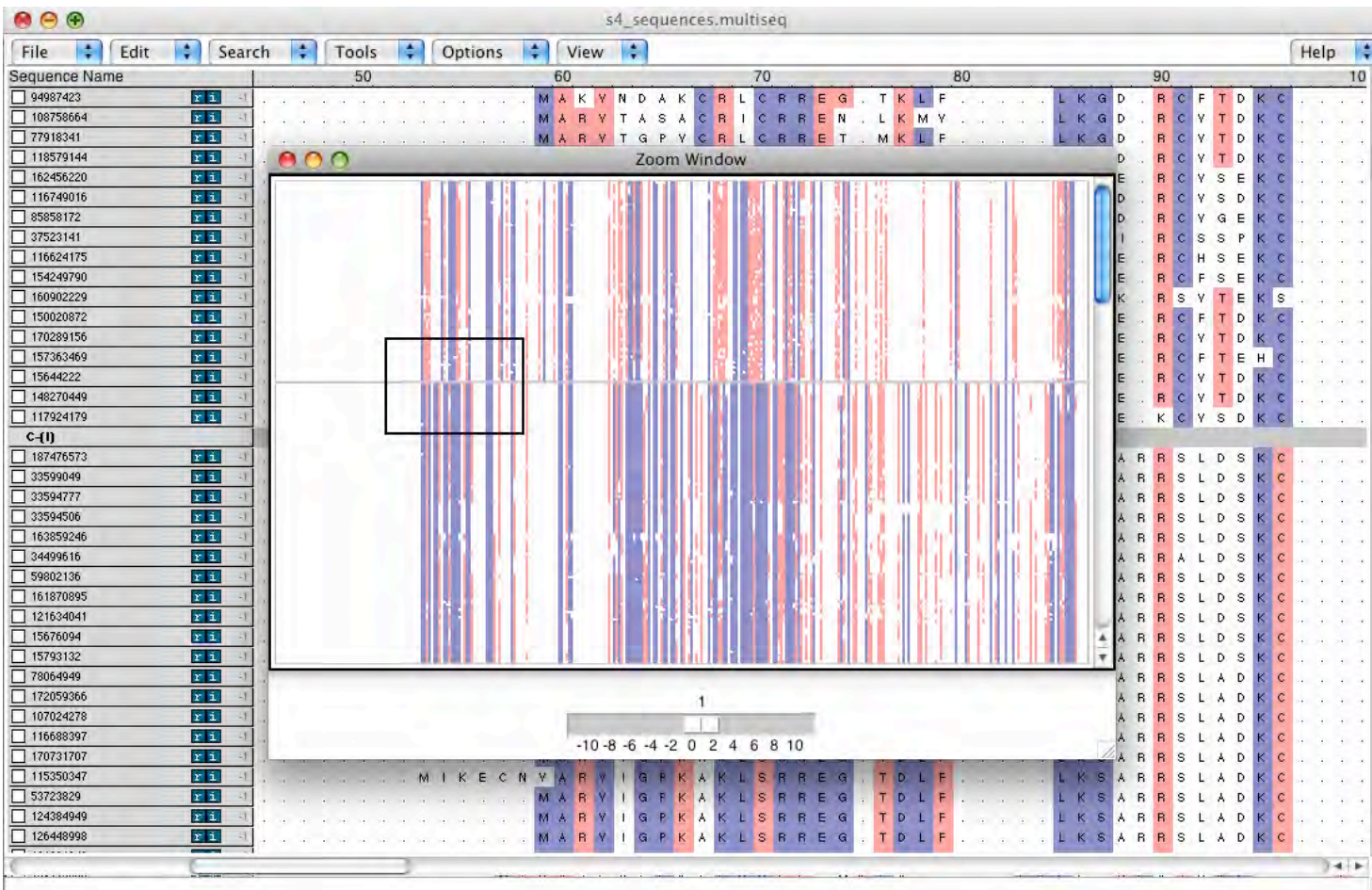
670 MB on disk

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.

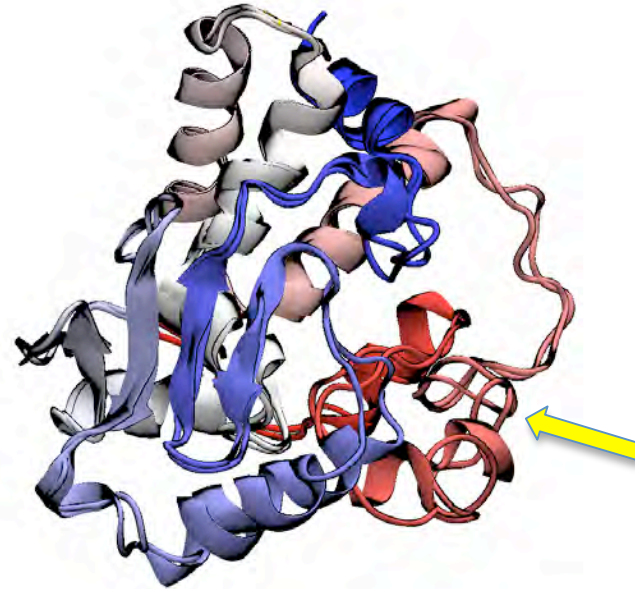
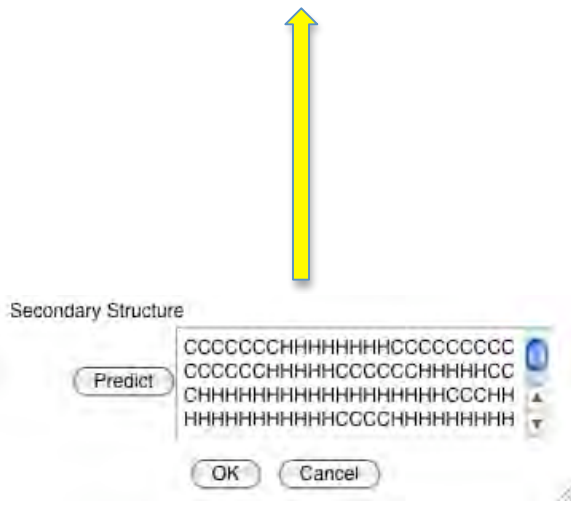
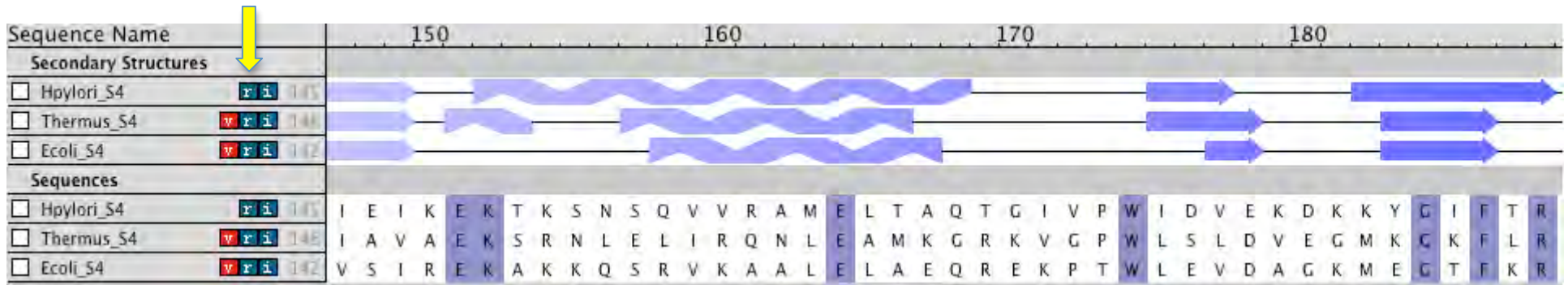


660 sequences of ribosomal protein S4 from all complete bacterial genomes*.

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

Secondary structure prediction

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



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.

A screenshot of the MultiSeq software preferences dialog, specifically the 'Software' tab. The dialog has two tabs at the top: 'Metadata' and 'Software'. The 'Software' tab is active. The dialog is divided into three main sections. The first section is 'BLAST Installation Directory' with fields for the directory path (set to '/usr/local/blast'), 'BLASTMAT' (set to 'data'), and 'BLASTDB'. The second section is 'PSIPRED Installation Directory' with fields for the directory path (set to '/Volumes/HomeRAID2/Homes/erobert3/Applications/OSX-1386/bin/'), 'PSIPREDDATA' (set to '/Volumes/HomeRAID2/Homes/erobert3/Applications/OSX-1386/share/psipred/data'), and 'PSIPREDDDB' (set to '/Volumes/Homes/Databases/psipred/psipred-sp'). The third section is 'Path to external editor'. At the bottom of the dialog are 'Close' and 'Help' buttons. A blue arrow points from the text on the left to the PSIPREDDATA field.

Export Modeller compatible alignments

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

```
>P1; Hpylori_S4
sequence:Hpylori_S4:::::0.00:0.00
MARYRGAVERLERRFVSLALKGE-RRLSGKSALDKRAYGPGQHGOR-RAKTSYGLQLK
EKQKAKMMYGISEKQFRSIFVEANRLDGNTGENLIRLIERRLDNVVYRMGFATTRSSARQ
LVTHGHVLDVGKRLDIPSYFVRSQKIEIKEKTKSNSQVVRAMELTAQTGIVPWIDVEKD
KKYGIFTRYPEREEVVVPIEERLIVELYSK*
```

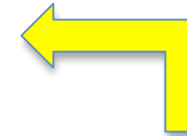
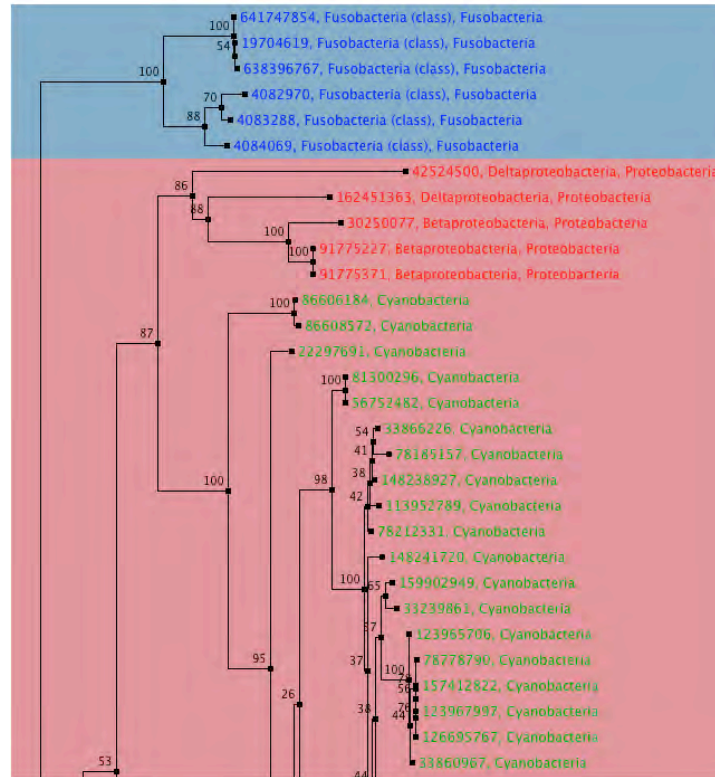
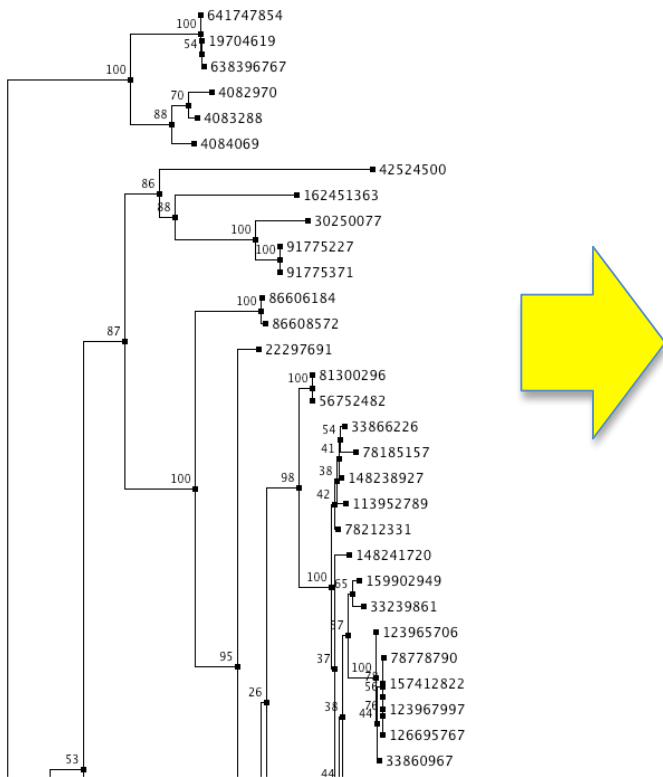
```
>P1; Thermus_S4
structureX:Thermus_S4:2:D:209:D::-1.00:-1.00
-GRYIGPVCRLCRREGVKLYLKGE-RCYSPKCAMERRPYPPGQHGQKRARRPSDYAVRLR
EKQKLRRRIYGISERQFRNLFEEASKKKGVTGVSFVLGLESRLDNVVYRLGFAVSRROARQ
LVRHGHITVNGRRVDLPSYRVRPGDEIAVAEKSRLNLELIRQNLEAMKGRKVGWLSLDVE
GMKGKFLRLPDREDLALPVNEQLVIEFYSR*
```

```
>P1; Ecoli_S4
structureX:Ecoli_S4:1:D:205:D::-1.00:-1.00
-ARYLGPKLKLSRREGTDLFLKSGVRAIDTKCKIE---QAPGQHGAR-KPRLSDYGVQLR
EKQKVRRIYGVLERQFRNYYKEAARLKGNTGENLLALLEGRLDNVVYRMGFATRAEARQ
LVSHKAIMVNGRVVNIASYQVSPNDVVSIREKAKKQSRVKAALAEQREKPTWLEVDAG
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()
```

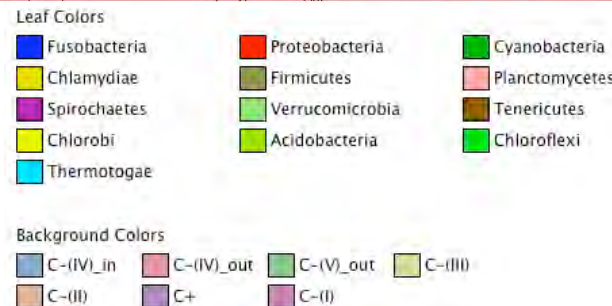
Phylogenetic tree editor

- 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.

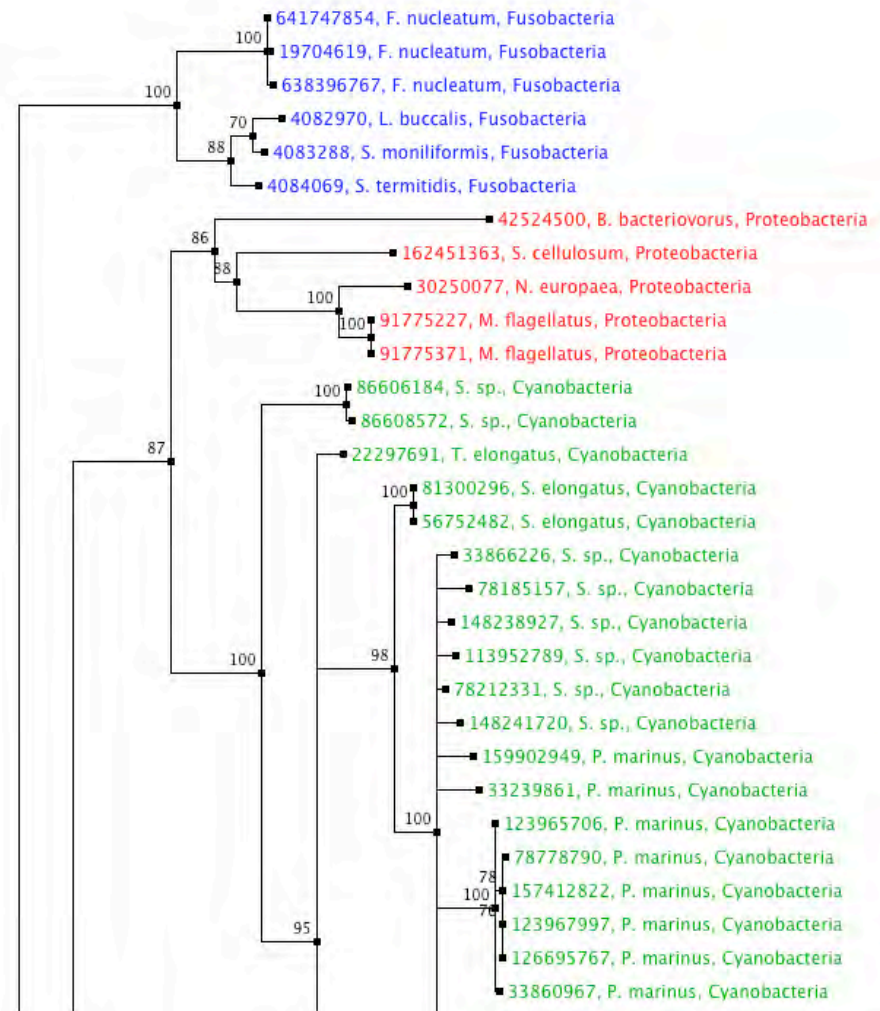
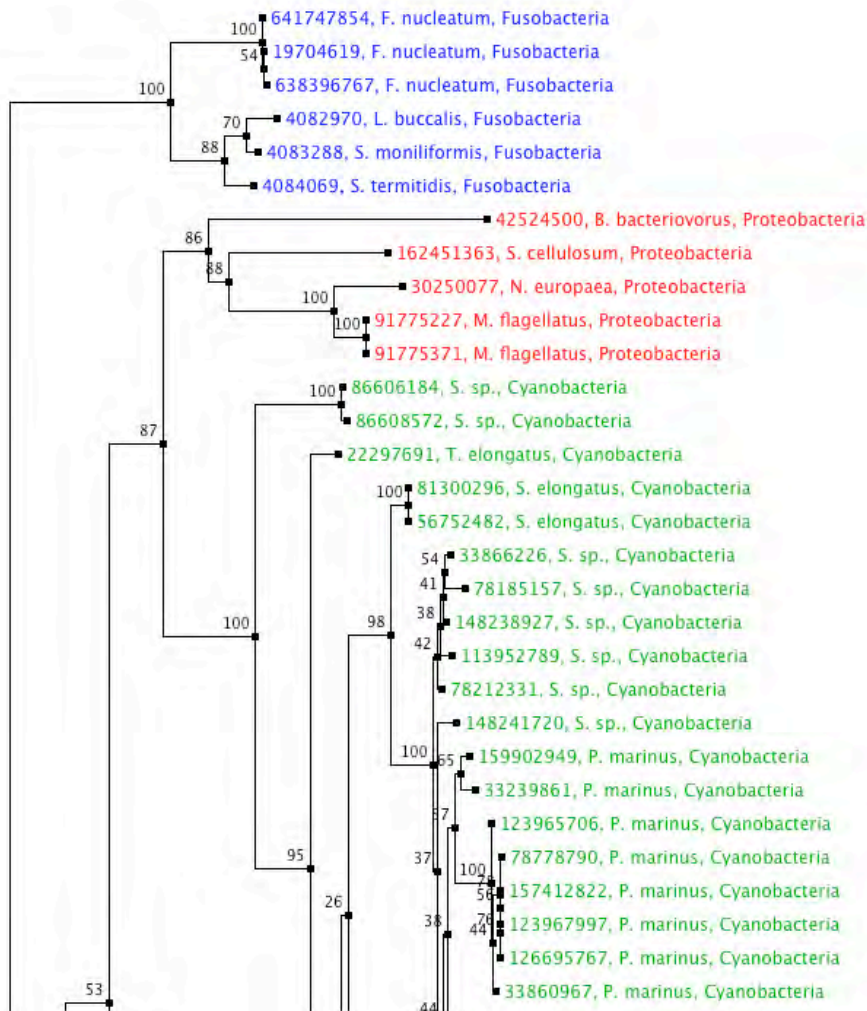
Maximum likelihood tree of 660 S4 sequences reconstructed using RAXML.



Elijah Roberts 2009

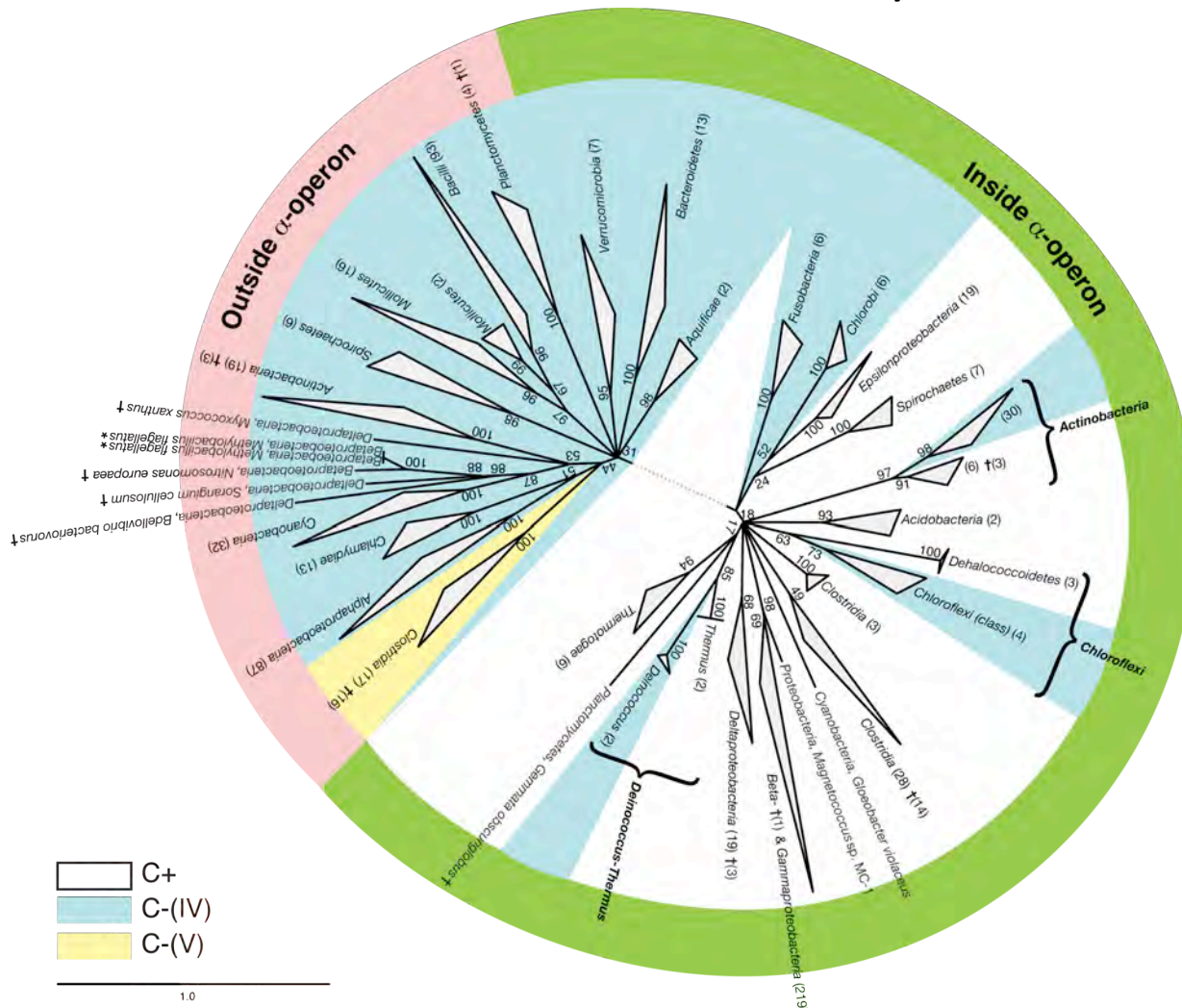
Edit the physical layout of the tree

- 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.

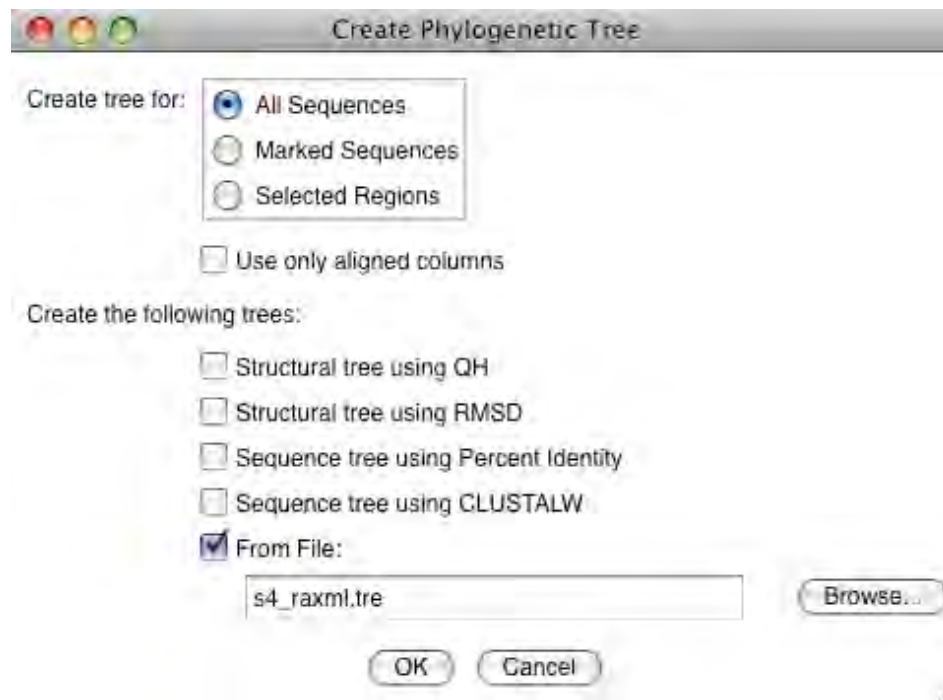


Combined phylogenetic tree and genome content analysis of ribosomal protein S4 for all complete bacterial genomes.

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

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

```

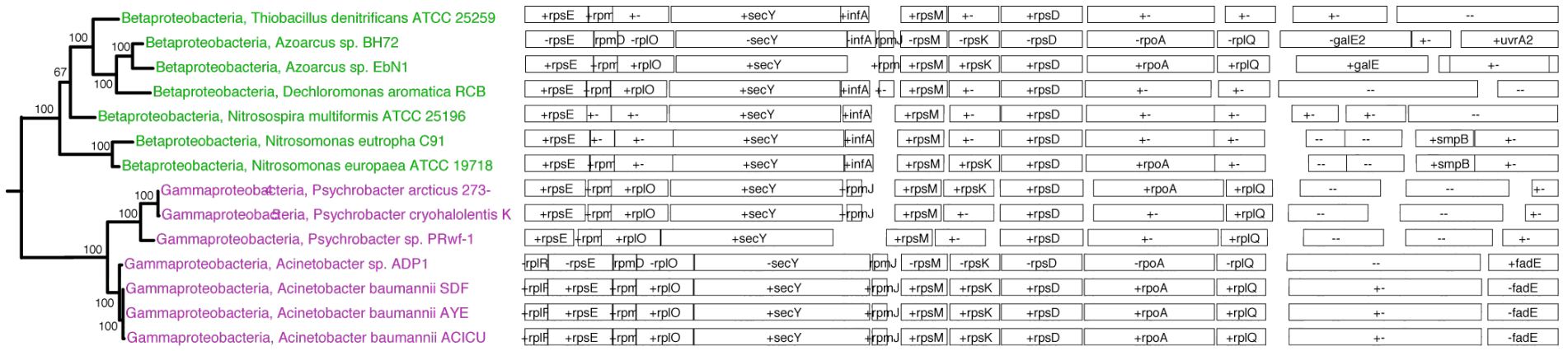
proc drawGenomicContextOfAlignment {outputFilename alignment contextDistance scaling genomeDirectory} {
    foreach sequence $alignment {
        # Make sure we have the GI number for this sequence.
        set giNumber [::SeqData::getSourceData $sequence "gi"]

        # 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
    }
}

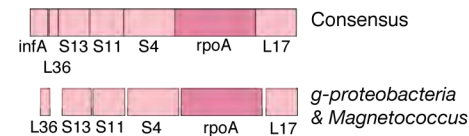
```



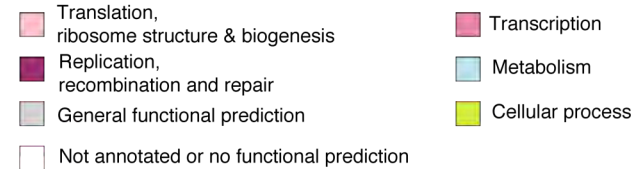
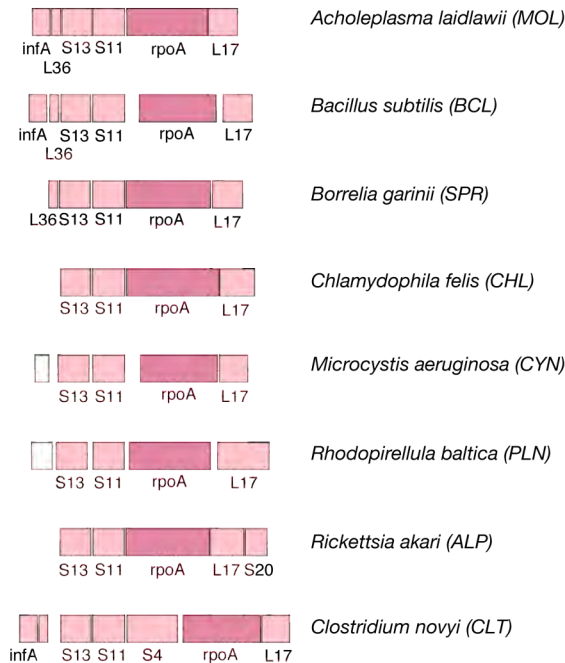
Genome content future directions

- 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.

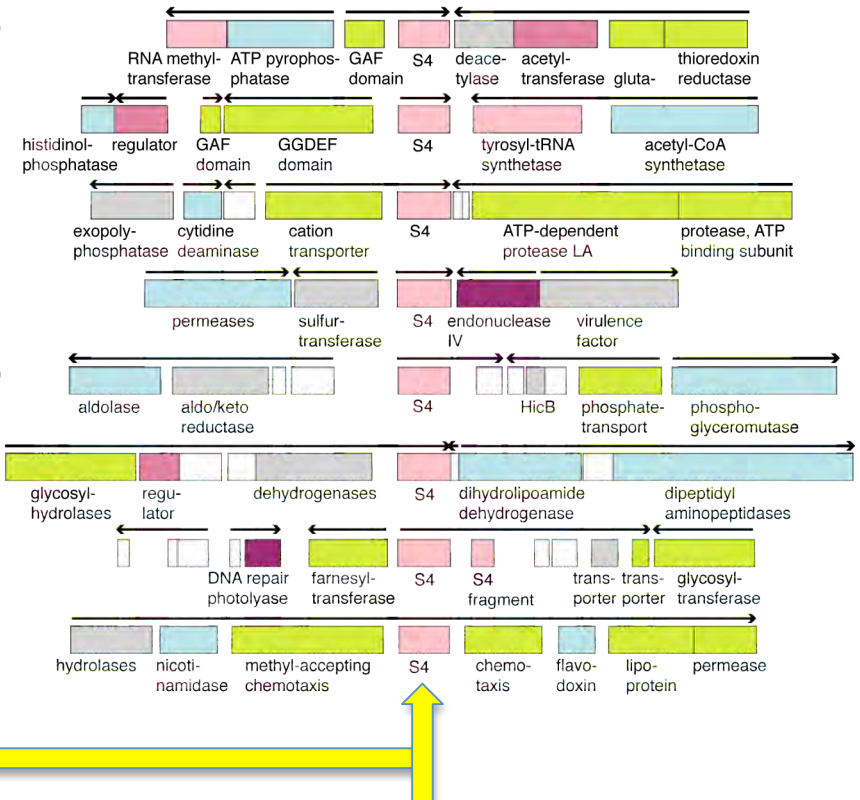
A) a-operon Organization



B) Corresponding a-operon for comparison



C) Outside-operon S4 context



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
SYK_GLOVI	1e-19	N P Y P Y R Y E	. R T H M A . . . G D L Q A K
666876	2e-19	T Q I C K I K S
67920132	2e-19	N G E E V E V D
23130228	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
1N9W	4e-19	R V L V R D L K A
46199389	5e-19	R V L V R D L K A
SYK_SYNY3	5e-19	R D L S N G E E
SYK_SYNEL	1e-18	A H L A A G E A
SYK_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	A A A L E G C E
SYK_PROMA	4e-18	P 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
SYK_STRR6	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
SYK_STRPN	6e-18	K Y A N L D K E Q L H
15900610	6e-18	K Y A N L D K E Q L H
62526807	6e-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
SYK1_SALTI	6e-18	E E L E A L N I
SYK_ENTFA	8e-18	Y D N H T K E E L S
56707357	8e-18	E L E E L D N K

Filter Options

E Score: e-5

Redundancy Cutoff: 100

Superkingdom:
Archaea
Bacteria
Eukaryota

Kingdom:
Fungi
Metazoa
Viridiplantae

Phylum:
Actinobacteria
Aquificae
Arthropoda
Ascomycota
Bacteroidetes
Chlamydiae

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,

The image shows the NCBI protein-protein BLAST search interface. At the top, the NCBI logo is on the left, and the text "protein-protein BLAST" is on the right. Below the logo, there are four tabs: "Nucleotide", "Protein", "Translations", and "Retrieve results for an RID". The "Protein" tab is selected. The main search area contains a large text input field for the query sequence. Below the input field, there are several options: "Search" (a link), "Set subsequence" (with "From:" and "To:" input fields), "Choose database" (with a dropdown menu showing "nr"), and "Do CD-Search" (with a checked checkbox). At the bottom, there are three buttons: "BLAST!", "Reset query", and "Reset all".

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;	Snapz Pro X

Final Blast Result: Sequence Alignment

 >[gi|46395801|sp|O88F17|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 VSLAFGLSIATLAQSVGHISGAHLNPAVTLGLLLSCQISVLRAIMYIIAQCVGAIVATAI 112

V+ AFGL++ T+A ++GHISG HLNPAV+ GL++ + + Y+IAQ +GAI+A +

Sbjct: 40 VAFAFGLTVLTMFAIGHISGCHLNPAVVSFGLVVGGRFPAKELLPYVIAQVIGAILAAGV 99

Query: 113 LSGITSSLP--DNSLGL--NALAP----GVNSGQGLGIEIIGTLQLVLCVLATTD RRRRD 164

+ I S + S GL N A G G G E++ T ++ ++ TD R

Sbjct: 100 IYLIASGKAGFELSAGLASNGYADHSPGGYTLGAGFVSEVVMTAMFLVVIMGATDARAP- 158

Query: 165 LGGSGPLAIGFSVALGHLLAIDYTGCGINPARSFGSSVITHNF--QDHWIFWVGPFIGAA 222

G P+AIG ++ L HL++I T +NPARS G ++ + Q W+FWV P IGAA

Sbjct: 159 -AGFAPIAIGLALTLIHLISIPVTNTSVNPARSTGPALFVGGWALQQLWLFWVAPLIGAA 217

Query: 223 LAVLIY 228

+ +Y

Sbjct: 218 IGGALY 223

Search returns approximate alignments - needing refinement!
Clustal, Muscle, MAFT, Tcoffee, pileup, Smith-Waterman, and
hand editing

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_CAEEL		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"/> 110w_A		1 M R . R T H Y A G S
<input type="checkbox"/> 1i12_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

The screenshot shows a metadata display window for the sequence SYDC_YEAST. The window is organized into several sections:

- Sequence Name:** SYDC_YEAST
- Source Organism:** Saccharomyces cerevisiae
- Common Name:** yeast
- EC Number:** 6.1.1.12
- EC Description:** Aspartate--tRNA ligase.
- Description:** Aspartyl-tRNA synthetase, cytoplasmic (EC 6.1.1.12) (Aspartate--tRNA ligase) (AspRS) - Saccharomyces cerevisiae (Baker's yeast).
- Data Sources:** sp=P04802,SYDC_YEAST; pdb=1EOV,A
- Lineage:** Eukaryota, Fungi, Ascomycota, Saccharomycotina, Saccharomyces, Saccharomycetales
- Notes:** A large empty text area for adding notes and annotations.

At the bottom of the window are two buttons: **OK** and **Cancel**.

Acknowledgements

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

