# *Introduction to evolutionary concepts and VMD/MultiSeq - Part I*

# *Characterizing molecular systems*

## Zaida (Zan) Luthey-Schulten

Dept. Chemistry, Physics, Beckman Institute, Institute of

Genomics Biology, & Center for Biophysics

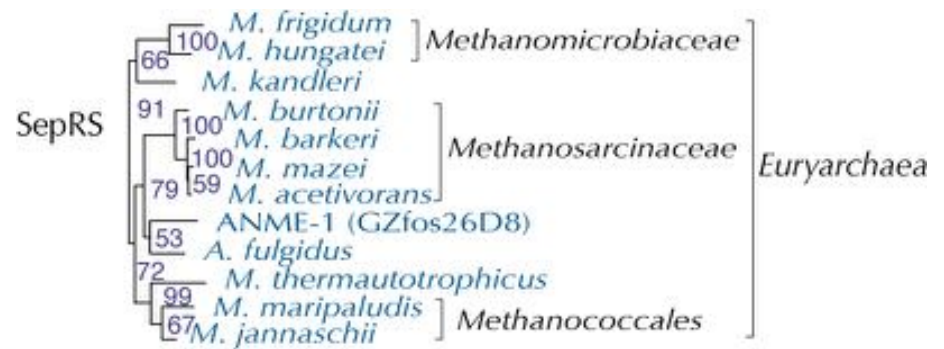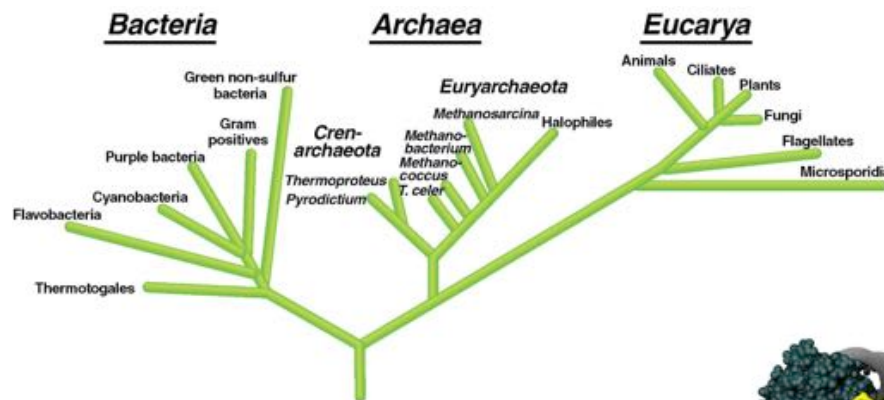ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# 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."*
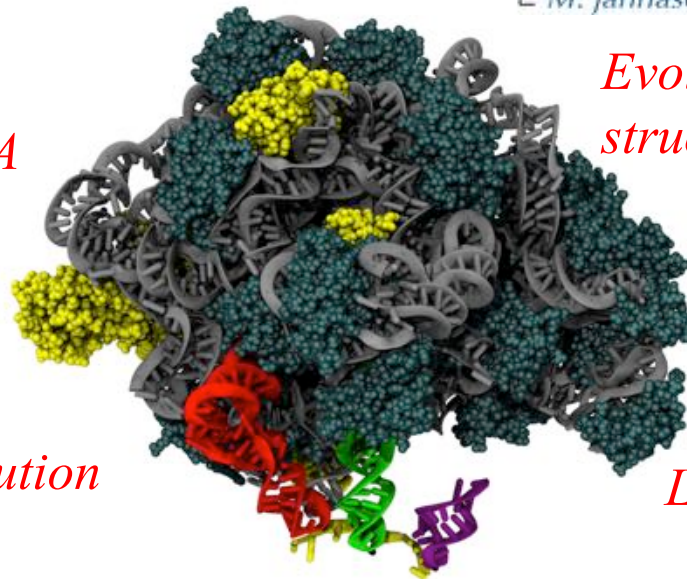


*UPT - Woese 16S rRNA*

*Evolutionary profiles for protein structure & function prediction*

*Signatures ribosomal evolution*

*LSU (23S rRNA + rproteins)*

# Why Look at More Than One Sequence?

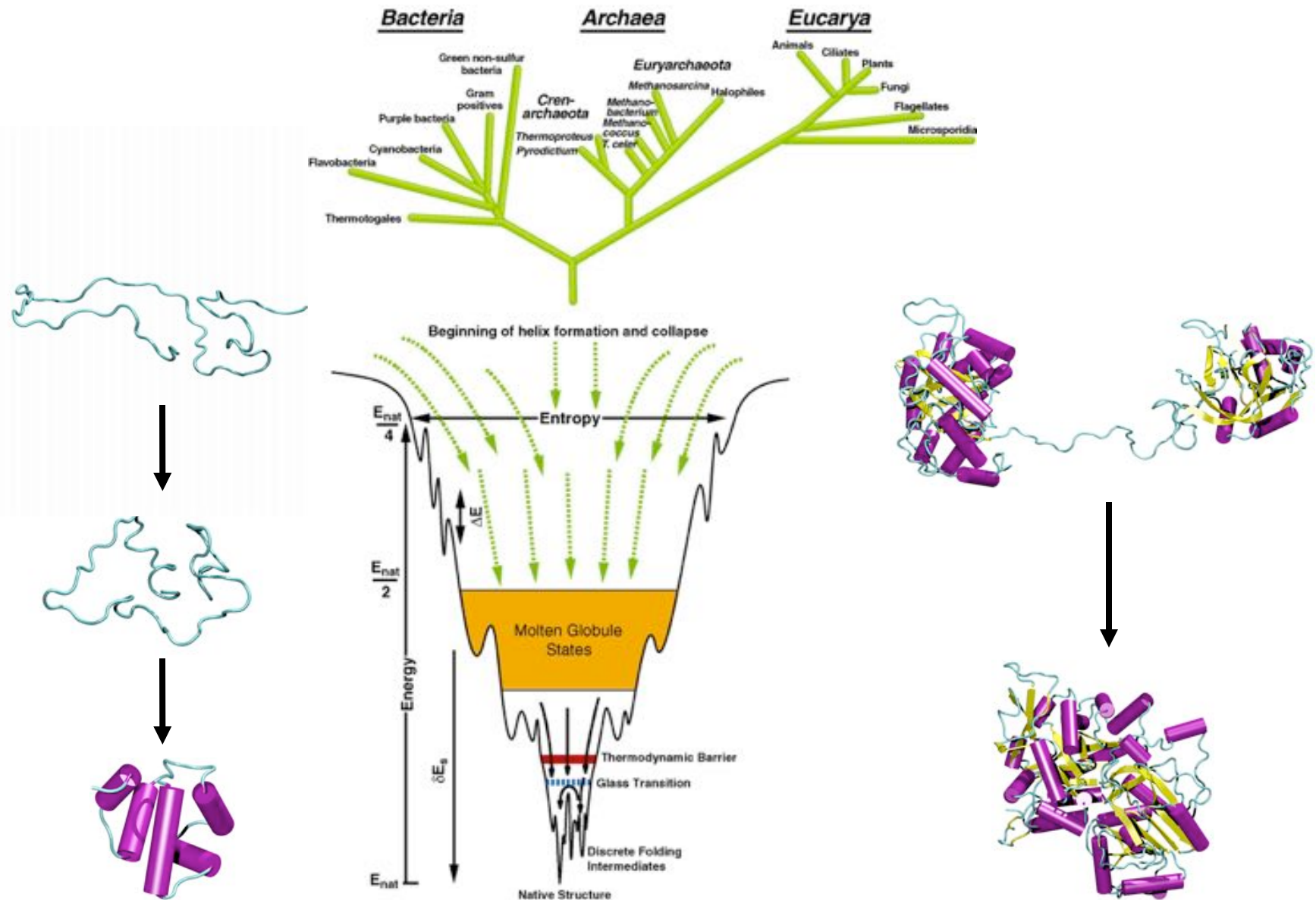## 1. Multiple Sequence Alignment shows patterns of conservation



2. Are these positions functionally important? Active sites, folding,..

3. What and how many sequences should be included?

4. Where do I find the sequences and structures for MS alignment?

5. How to generate pairwise and multiple sequence alignments?

# Protein (RNA) Folding, Structure, & Function

# New Tools in VMD/MultiSeq

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

Synchronization between 1D and 3D views

Protein / RNA
Sequence Data

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

Sequence /Structure
Alignment

Protein & RNA
secondary structure

Group data by taxonomic classification

View sequence or structure phylogenies and eliminate redundancy with QR

Metadata Information,
Clustal, MAFFT &
Phylogenetic Trees

RAXml Trees,
Genomic Content,
Temperature DB

QR non-redundant
seq / str sets

Cluster
analysis /
Bioinformatics
scripting

Import data directly from BLAST databases

Align sequences with Clustal

Blast & PsiBlast

Sequence Editor: Manually adjust alignments or sequences
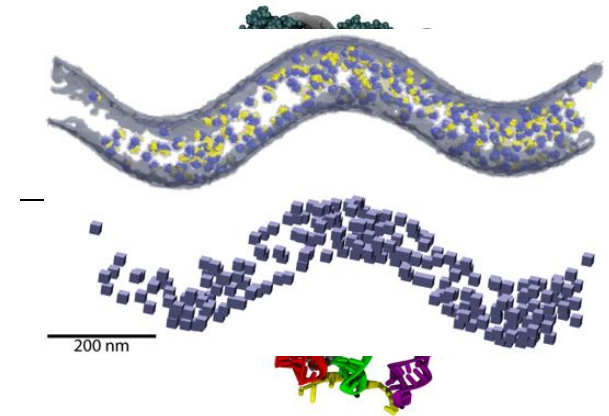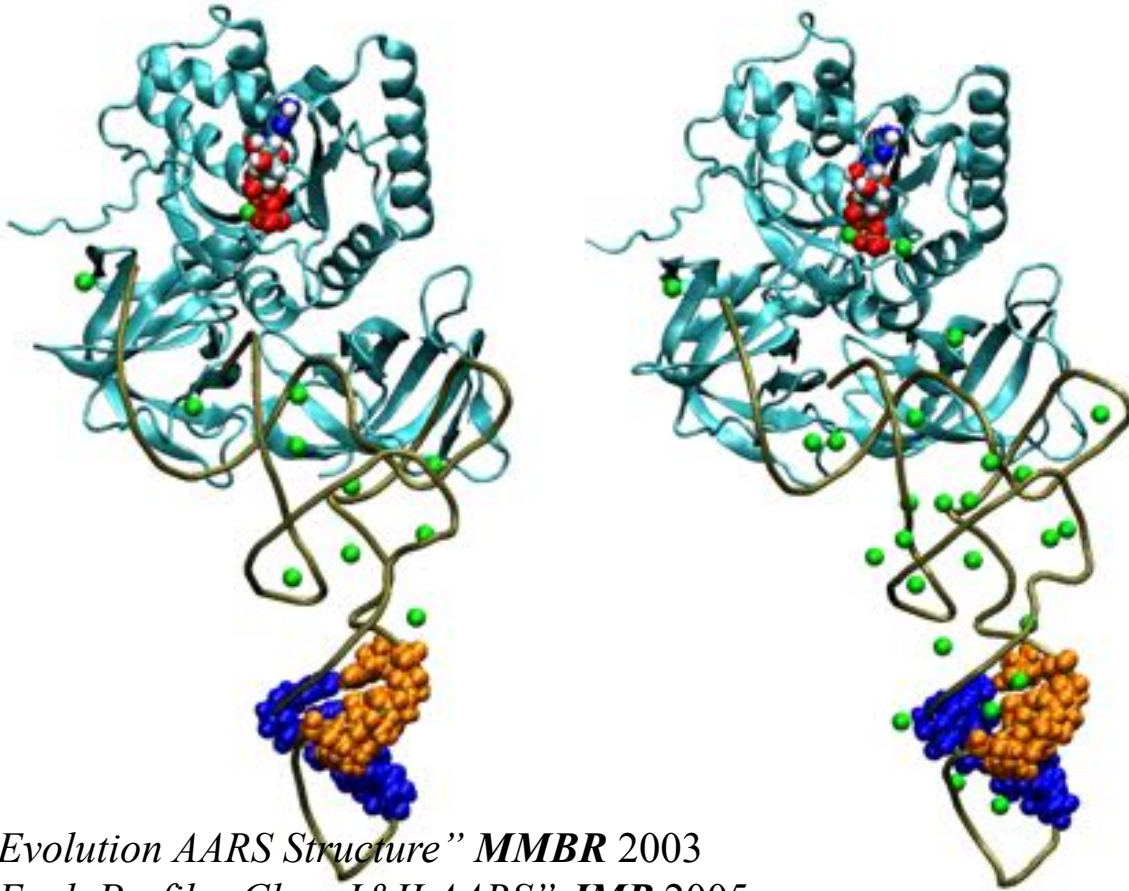
Tutorials MultiSeq/
AARS

Sequence Editor

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)

# Protein:RNA Complexes in Translation
## Evolutionary Analysis & Dynamics



**r-Proteins/r-RNA**
**Ribosome LSU**

*"Signatures ribosomal evolution"*
***PNAS*** 2008, ***BMC*** 2009, ***BJ*** 2010
*"Motion L1 Stalk:tRNA"* ***JMB*** 2010**,**
**"***Ribosome Biogenesis*" ***JPC*** 2012,3
*"Whole cell simulations on GPUs"*
***IEEE*** 2009,***Plos CB*** 2011,***PRL***2011,
***JCC*** 2013**, *PNAS*** 2013,
***PRL*** 2013**, *CSB*** 2013
***Nature* 2014, *BJ*** 2015

*"Evolution AARS Structure"* **MMBR** 2003
*"Evol. Profiles Class I&II AARS"* **JMB** 2005
*"Evolution SepRS/CysRS"* **PNAS 2005**
*"Dynamic Signaling Network"* **PNAS** 2009
*"Exit Strategy Charged tRNA"* **JMB** 2010
*"Mistransl. in Mycoplasma"* **PNAS** 2011
*"Capture & Selection of ATP"***JACS** 2013

*"Recognition & tRNA Dynamics"*
***JMB*** *2008,****FEBS*** 2010, ***RNA*** 2012
*Network Viewer,* ***Bioinf., JCTC*** 2012

# Basic principles of evolutionary analysis for proteins & RNAs

- Comparative analysis of sequences and <span style="color:red">structures</span>
- Multiple sequence alignments <span style="color:red">(gaps and editing)</span>
- Sequence and <span style="color:red">structure</span> 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

"G" scattered around gaps



"Classic"
ClustalW
alignment

~ 5 minutes

MAFFT7*
alignment

~ 30 seconds

More sequences!

"G" aligned

http://www.clustal.org/clustal2/     *  MAFFT v7.221, Katoh and Standley, Mol.Biol and Evol. 2015

# Sequence Alignment & Dynamic Programming

Seq. 1: $a_1\ a_2\ a_3$ - - $a_4\ a_5 \ldots a_n$

Seq. 2: $c_1$ - $c_2\ c_3\ c_4\ c_5$ - $\ldots c_m$

number of possible alignments:

$$= \binom{2n}{n} = 2^{2n}\left(\sqrt{n\pi}\right)^{-1}$$

Needleman-Wunsch alignment algorithm

$$H(i,j) = MAX \begin{cases} H(i-1, j-1) + S[a(i), b(j)] \\ H(i, j-k) - W(k), \\ H(i-m, j) - W(m) \end{cases}$$

$\underline{\underline{S}}$ : substitution matrix



Score Matrix H: Traceback

gap penalty W = - 6

Reference: "Biological Sequence Analysis - Probabilistic Models of Proteins and Nucleic Acids" R. Durbin, S. Eddy, A. Krogh, and G. Mitchison, Cambridge U. P.London, 1998; pp. 19-22 (see also other sections)

# Needleman-Wunsch Global Alignment

## Similarity Values

|   | M | G | K | P |
|---|---|---|---|---|
| M | 5 | -3 | -1 | -2 |
| G | -3 | 6 | -2 | -2 |
| P | -2 | -2 | -1 | 7 |
| K | -1 | -2 | 5 | -1 |
| K | -1 | -2 | 5 | -1 |
| P | -2 | -2 | -1 | 7 |

## Initialization of Gap Penalties

|   |   | M | G | K | P |
|---|---|---|---|---|---|
|   | 0 → -6 → -12 → -18 → -24 |   |   |   |   |
| M | -6 | 5 | -3 | -1 | -2 |
| G | -12 | -3 | 6 | -2 | -2 |
| P | -18 | -2 | -2 | -1 | 7 |
| K | -24 | -1 | -2 | 5 | -1 |
| K | -30 | -1 | -2 | 5 | -1 |
| P | -36 | -2 | -2 | -1 | 7 |

# Filling out the Score Matrix H

# Traceback and Alignment



The Alignment

Traceback (blue) from optimal score

# STAMP - Multiple Structural Alignments

1. **Initial Alignment Inputs**

- Multiple Sequence alignment
- Ridged Body "Scan"
- Pairwise Alignments and Hierarchical Clustering

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 with no gap penalty

R. Russell, G. Barton (1992) *Proteins* **14**: 309     R.B. Russel, T. Walsh, G. Barton, STAMP version 4.4: User Guide, 2010.

# Multiple Structural Alignments

## STAMP – cont'd

2. Refine Initial Alignment & Produce Multiple Structural Alignment

Alignment score:

$$S_c = \frac{S_p}{L_P} \frac{L_p - i_A}{L_A} \frac{L_p - i_B}{L_B}$$

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

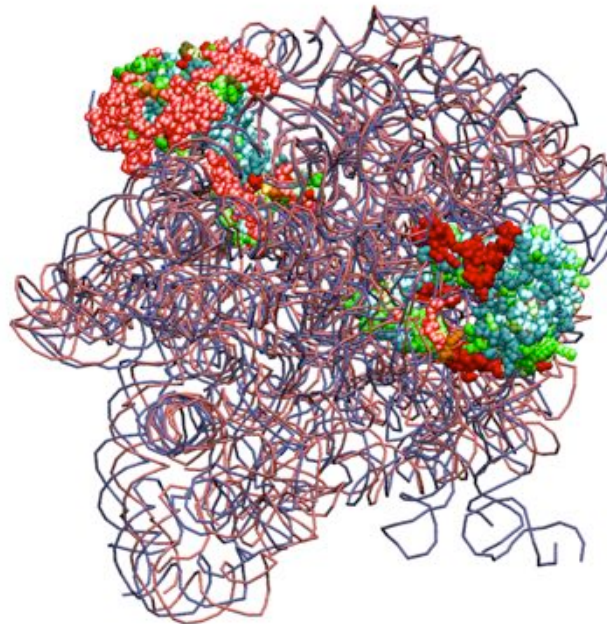$L_p, L_A, L_B$  — length of alignment, sequence A, sequence B

$i_A, i_B$  — 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.

# Structural Overlaps - STAMP

Ribosome large subunit showing ribosomal proteins L2 and L3
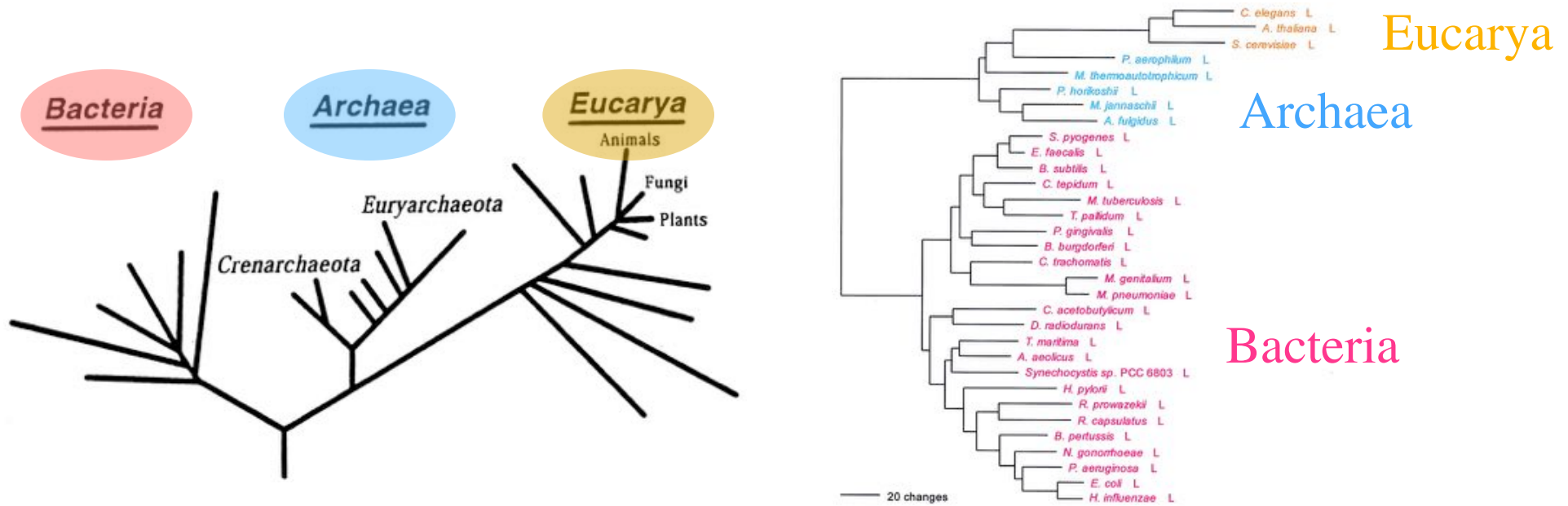180,000 atoms in 4 rRNAs and 58 proteins



E. coli                                          arismortui
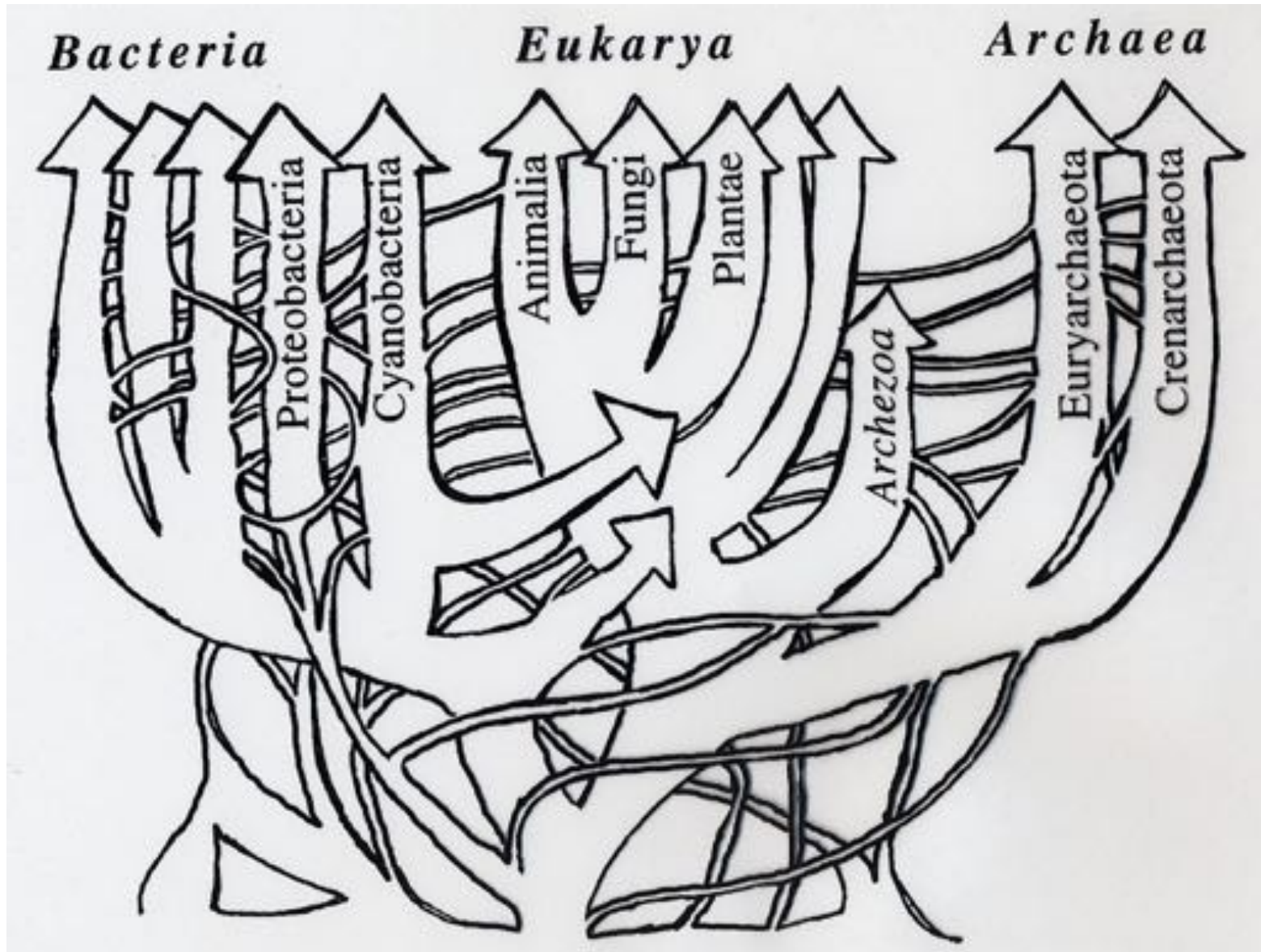
# Universal Phylogenetic Tree

## 3 domains of life



Reference 16S rRNA tree

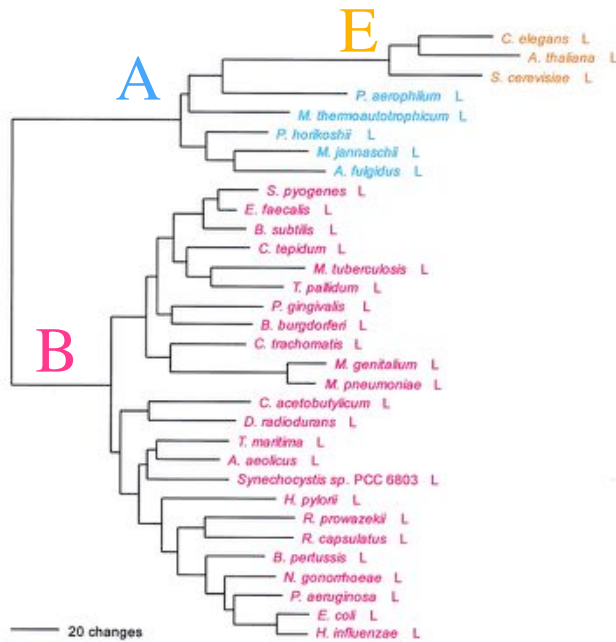Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution.

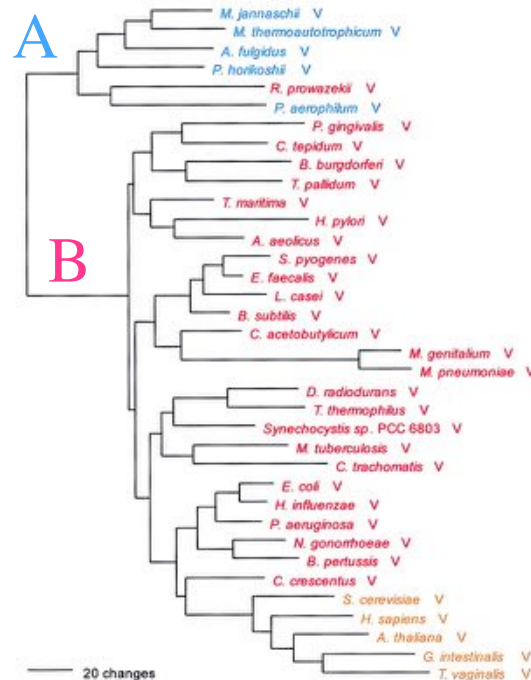# Look for horizontal gene transfer events



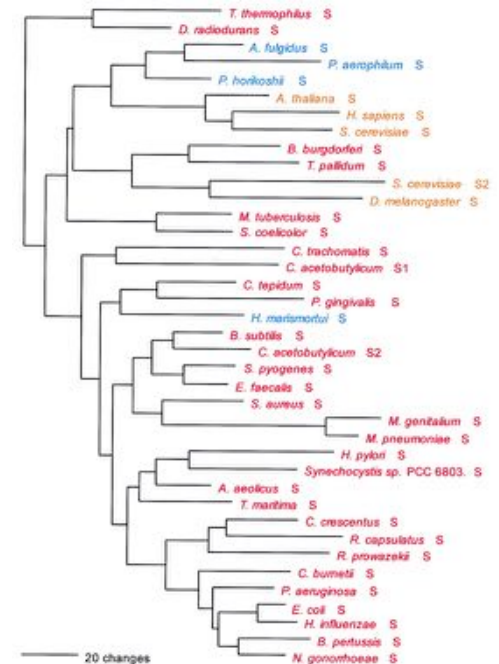After W. Doolittle, modified by G. Olsen

# 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
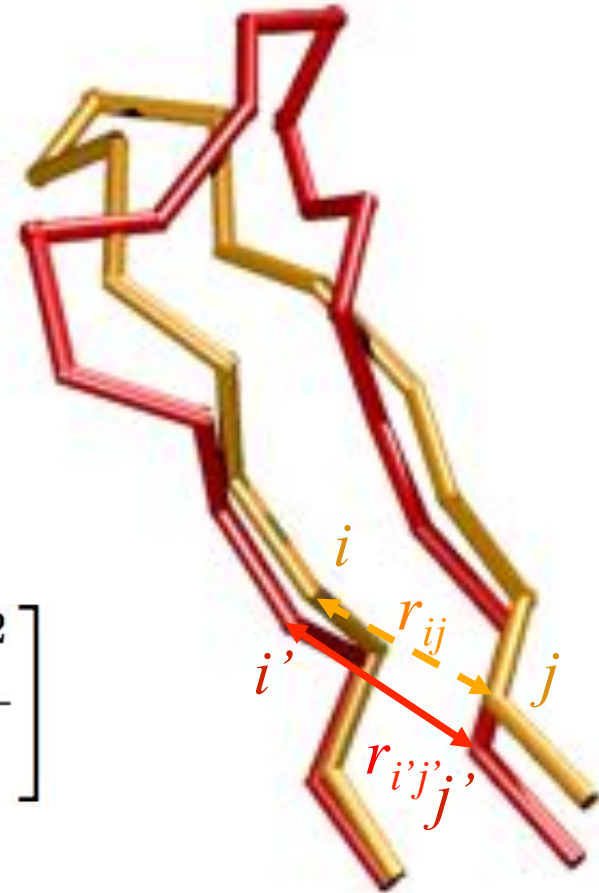
Woese, Olsen, Ibba, Soll *MMBR* 2000

# Protein Structure Similarity Measure

## $Q_H$ Structural Homology

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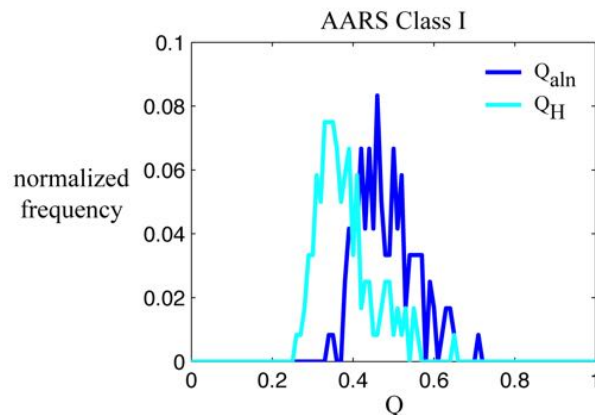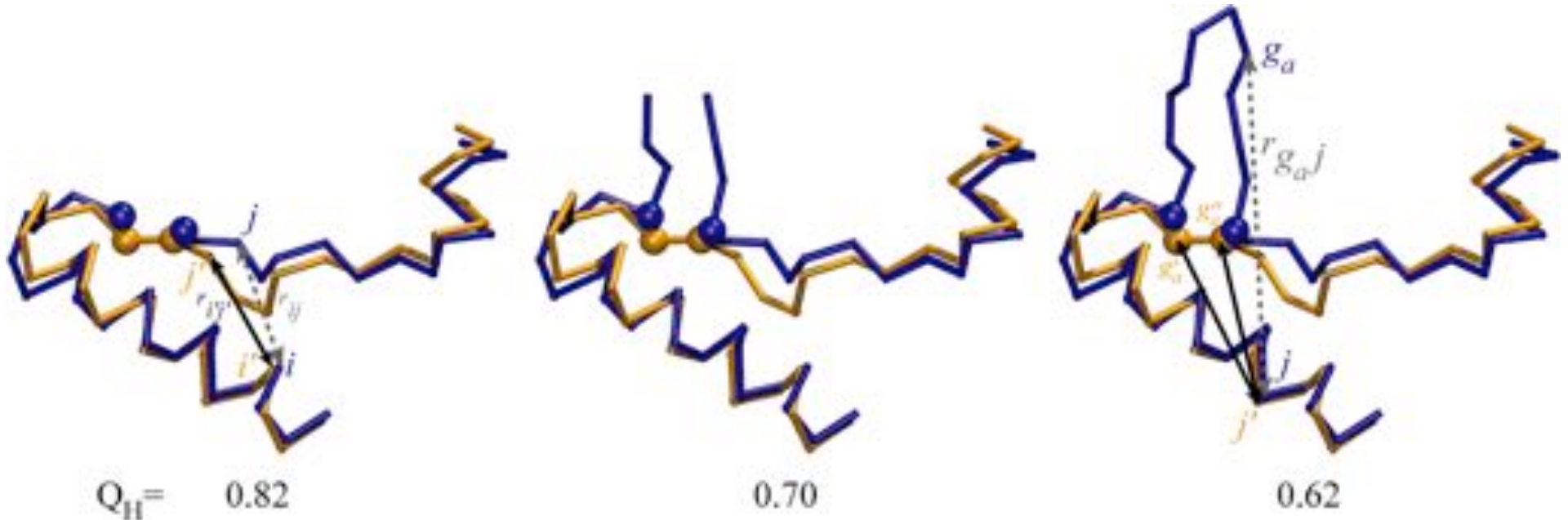presence and perturbation of gaps

$$Q_H = \aleph \left[ q_{aln} + q_{gap} \right]$$

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

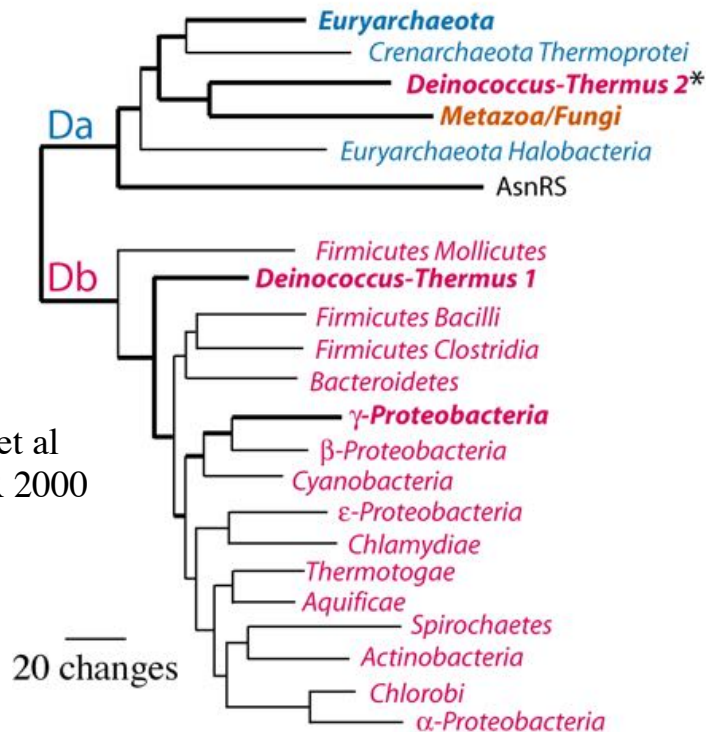# Structural Similarity Measure: The effect of insertions

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



$Q_H =$  0.82          0.70          0.62



AARS Class I

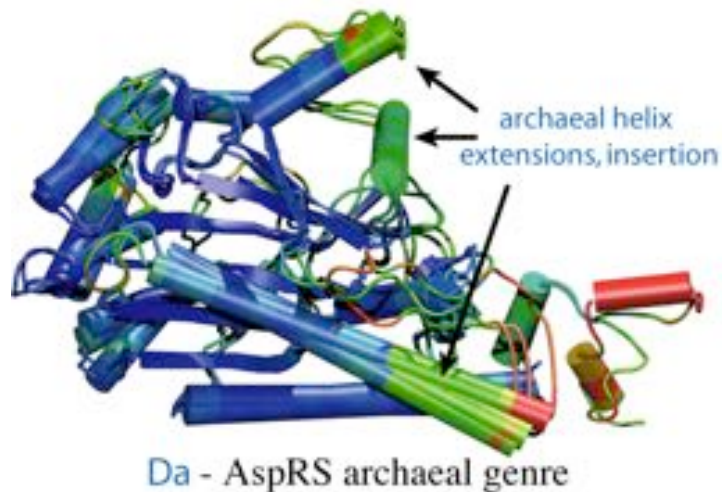$$q_{gap} = \sum_{g_a}\sum_{j}^{N_{aln}} \max\left\{ \exp\left[-\frac{(r_{g_a j} - r_{g'_a j'})^2}{2\sigma_{g_a j}^2}\right], \exp\left[-\frac{(r_{g_a j} - r_{g''_a j'})^2}{2\sigma_{g_a j}^2}\right]\right\}$$

$$+ \sum_{g_b}\sum_{j}^{N_{aln}} \max\left\{ \exp\left[-\frac{(r_{g_b j} - r_{g'_b j'})^2}{2\sigma_{g_b j}^2}\right], \exp\left[-\frac{(r_{g_b j} - r_{g''_b j'})^2}{2\sigma_{g_b j}^2}\right]\right\}$$
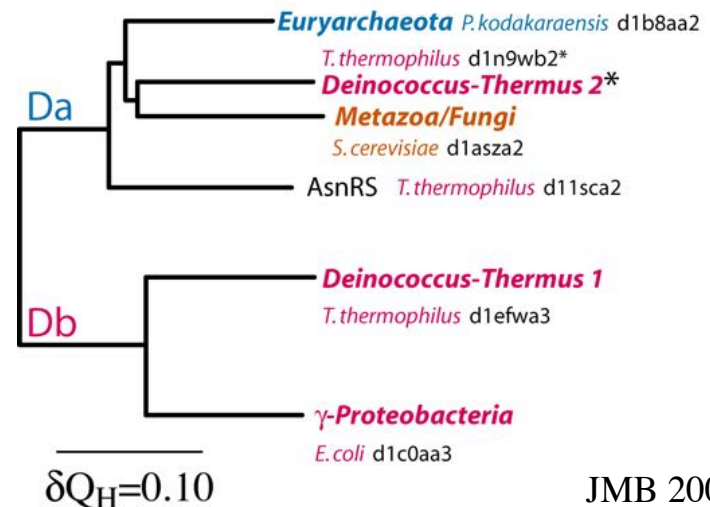
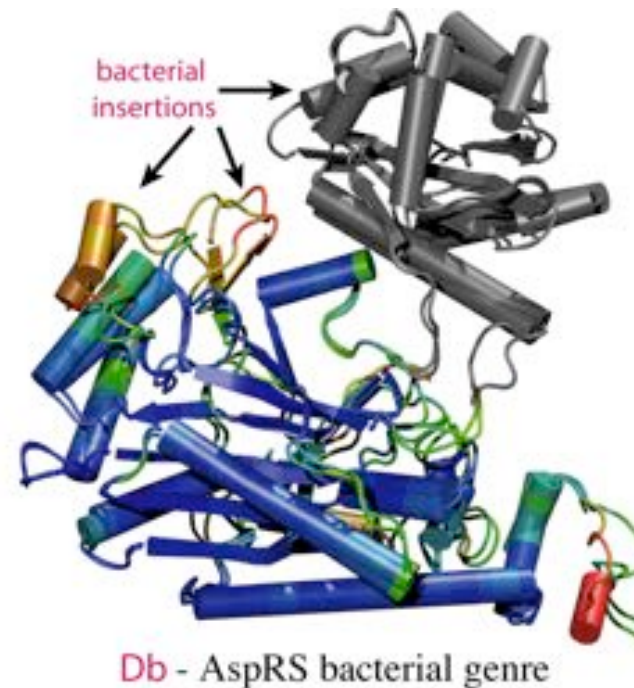# Structure encodes evolutionary information!



sequence-based phylogeny

**Euryarchaeota**
*Crenarchaeota Thermoprotei*
**Deinococcus-Thermus 2***
**Metazoa/Fungi**
*Euryarchaeota Halobacteria*
AsnRS

Da

Db

*Firmicutes Mollicutes*
**Deinococcus-Thermus 1**
*Firmicutes Bacilli*
*Firmicutes Clostridia*
*Bacteroidetes*
**γ-Proteobacteria**
*β-Proteobacteria*
*Cyanobacteria*
*ε-Proteobacteria*
*Chlamydiae*
*Thermotogae*
*Aquificae*
*Spirochaetes*
*Actinobacteria*
*Chlorobi*
*α-Proteobacteria*

Woese et al
MMBR 2000

20 changes

structure-based phylogeny

**Euryarchaeota** *P. kodakaraensis* d1b8aa2
*T. thermophilus* d1n9wb2*
**Deinococcus-Thermus 2***
**Metazoa/Fungi**
*S. cerevisiae* d1asza2
AsnRS *T. thermophilus* d1lsca2

Da

Db

**Deinococcus-Thermus 1**
*T. thermophilus* d1efwa3

**γ-Proteobacteria**
*E. coli* d1c0aa3

δQ_H=0.10

JMB 2005
MMBR 2003

archaeal helix
extensions, insertion

Da - AspRS archaeal genre

bacterial
insertions
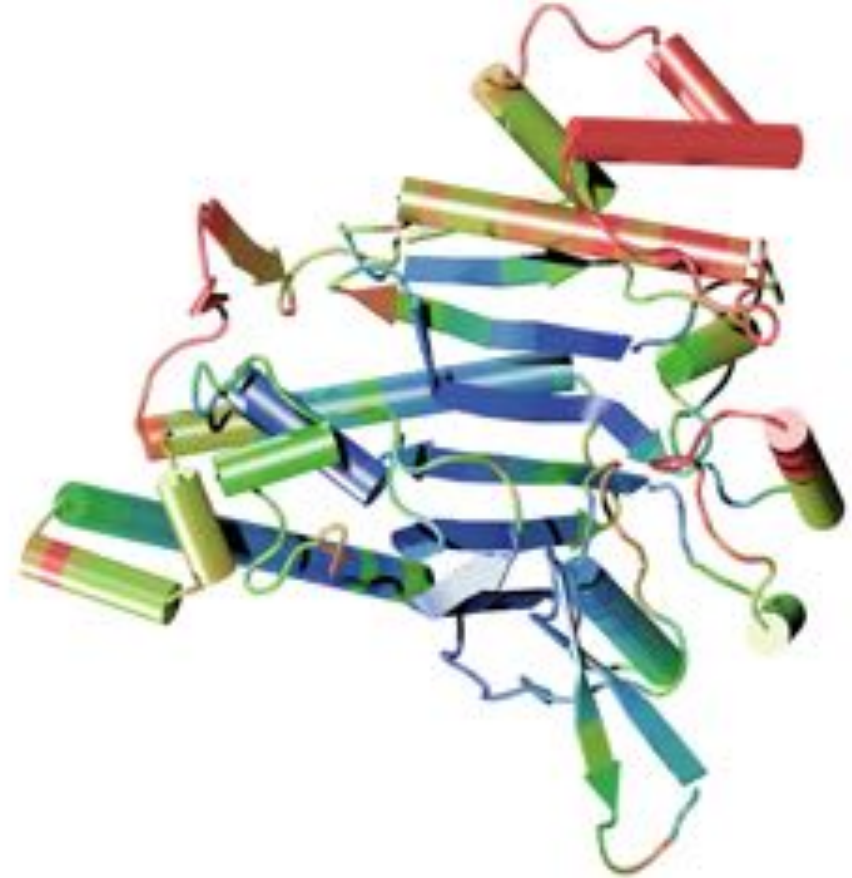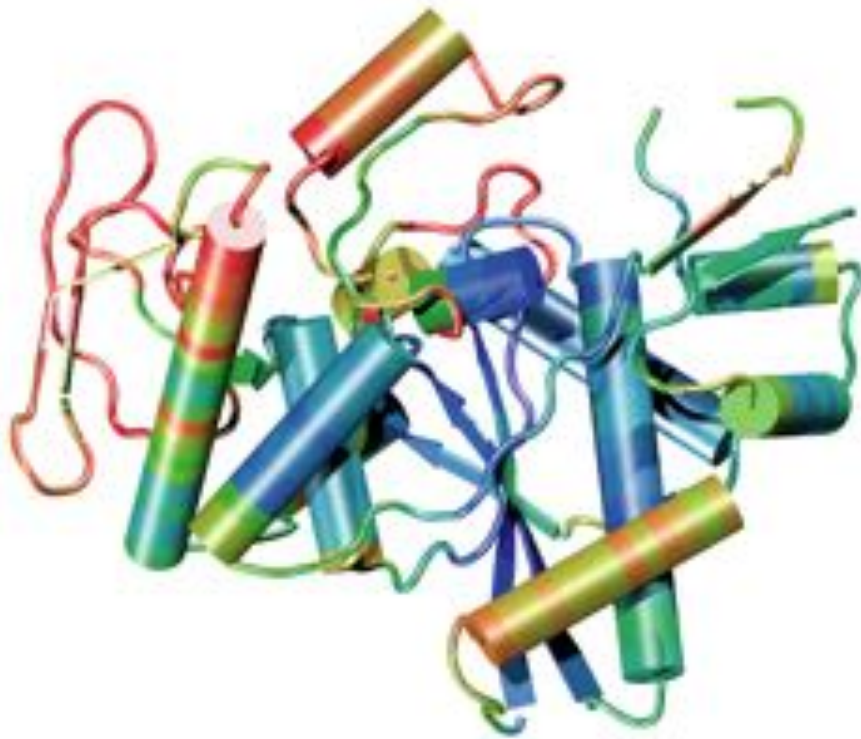
Db - AspRS bacterial genre

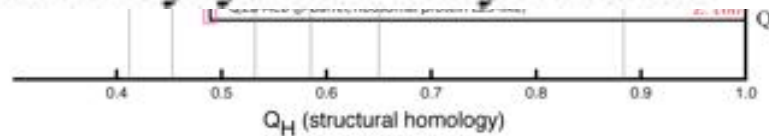# Structure reveals distant evolutionary events

## Class I AARSs

## Class II AARSs



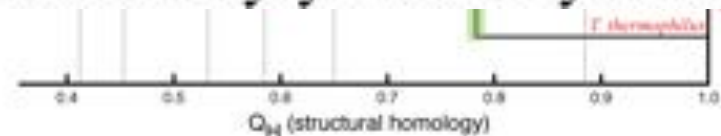structure-based phylogenetics

sequence-structure overlap

structure-based phylogenetics

sequence-structure overlap
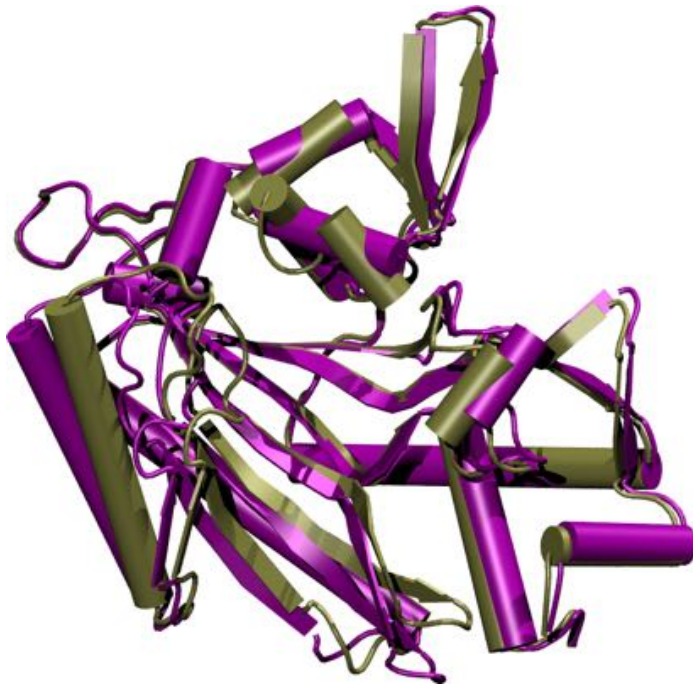
**Class I Lysyl-tRNA Synthetase**

$Q_H$ (structural homology)

**Class II Lysyl-tRNA Synthetase**

$Q_H$ (structural homology)

# Sequences define more recent evolutionary events



Conformational changes
in the same protein.



Structures for two
different species.

ThrRS
T-AMP analog, 1.55 A.
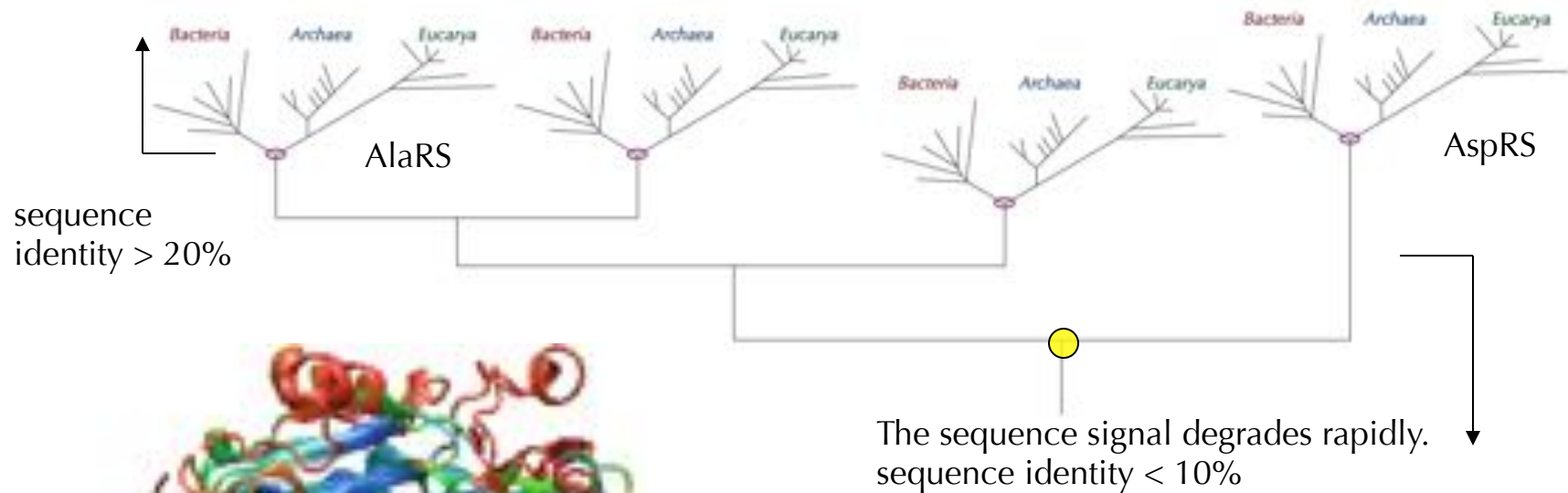T, 2.00 A.

$Q_H = 0.80$
Sequence identity = 1.00

ProRS
$M. jannaschii$, 2.55 A.
$M. thermoautotrophicus$, 3.20 A.

$Q_H = 0.89$
Sequence identity = 0.69

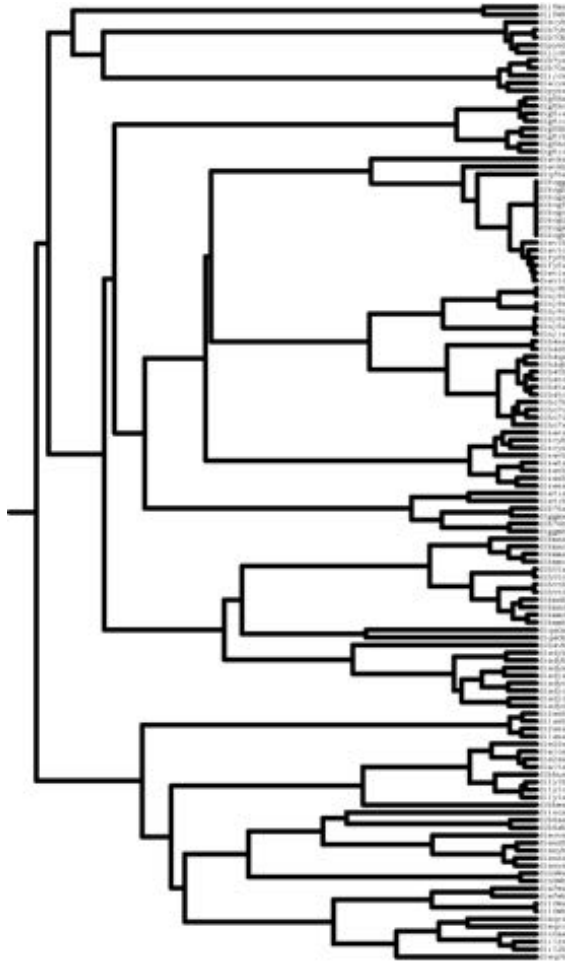# Relationship Between Sequence & Structure



sequence identity > 20%

AlaRS

AspRS

The sequence signal degrades rapidly.
sequence identity < 10%

Structural superposition of AlaRS & AspRS.
⬤ Sequence id = 0.055, $Q_H$= 0.48

O'Donoghue & Luthey-Schulten (2003) *MMBR* 67: 550–73.
Structural alignment & visualization software MultiSeq/VMD

$Q_H$ (structural similarity) vs sequence identity

○ AARS class I
○ AARS class II

# Non-redundant Representative Profiles

Too much information
129 Structures

Economy of information
16 representatives

Multidimensional QR
factorization
of alignment matrix, $A$.

$$A = \begin{bmatrix} & & & \nearrow^{d=4} \\ & & \boxed{Z} & \boxed{G} \\ l_{aln} \downarrow & \boxed{X} & \boxed{Y} & \\ & & \xrightarrow{k_{proteins}} & \end{bmatrix}$$

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 $g = \gamma \dfrac{\|X\|_{F_4} + \|Y\|_{F_4} + \|Z\|_{F_4}}{\|G\|_{F_4}}$

adjustable parameter

## Sequence Space
Orthogonal Encoding = 24-space

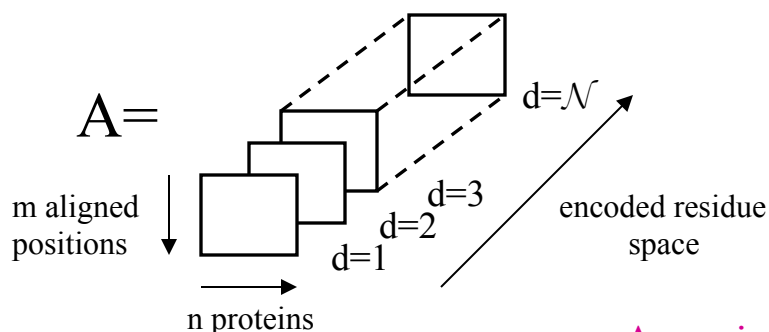23 amino acids (20 + B, X, Z) + gap

A = (1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0)
B = (0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0)
C = (0,0,1,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0)
...
GAP = (0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,0,1)

## Alignment is a Matrix with Linearly Dependent Columns

A=

m aligned positions
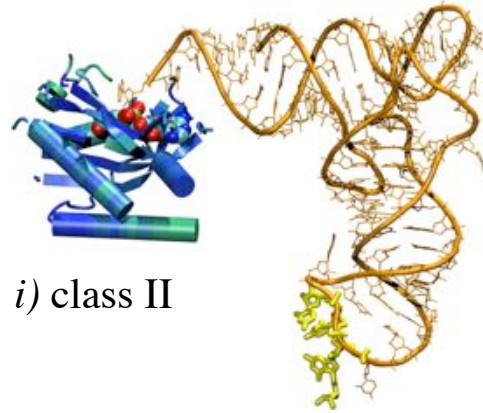
n proteins

d=1, d=2, d=3, ..., d=$\mathcal{N}$ encoded residue space

$$Q_{(d)}^T A_{(d)} P = Q_{(d)}^T \begin{bmatrix} & & G \\ & Z & \\ Y & & \\ X & & \end{bmatrix} \quad P = \tilde{R}_{(d)}$$

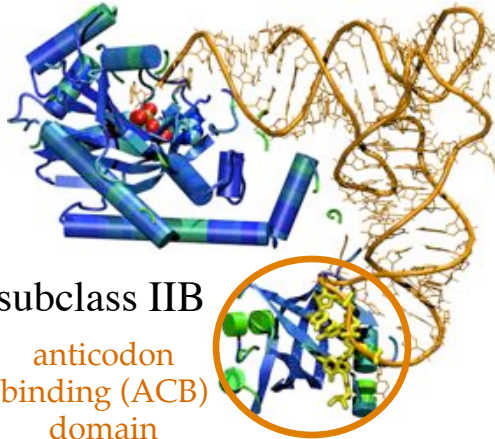d=4, d=1, $m_{aln}$, $n_{proteins}$

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

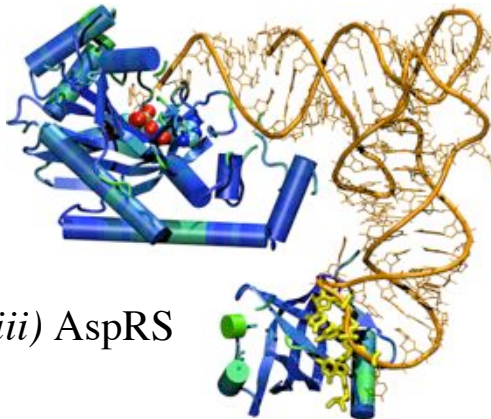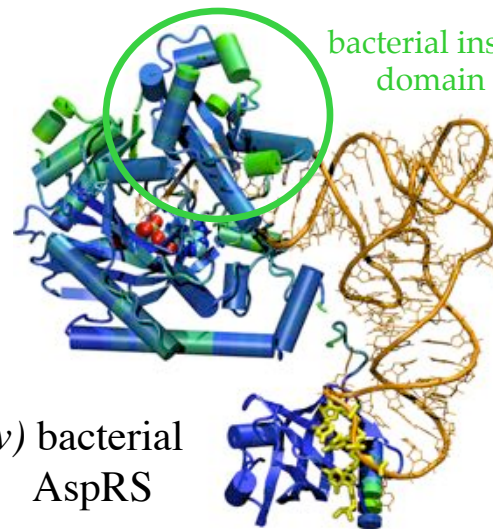# Evolution of Structure and Function in AspRS
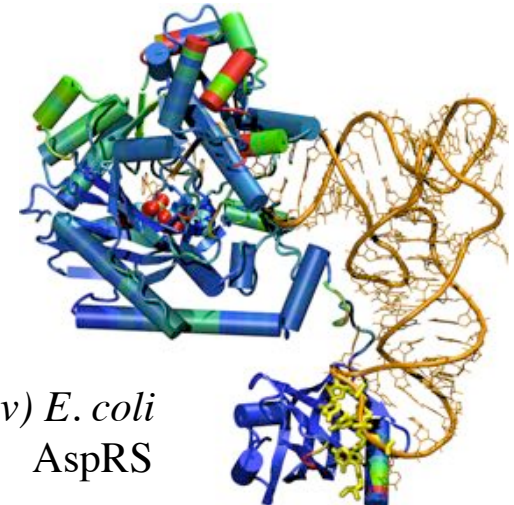


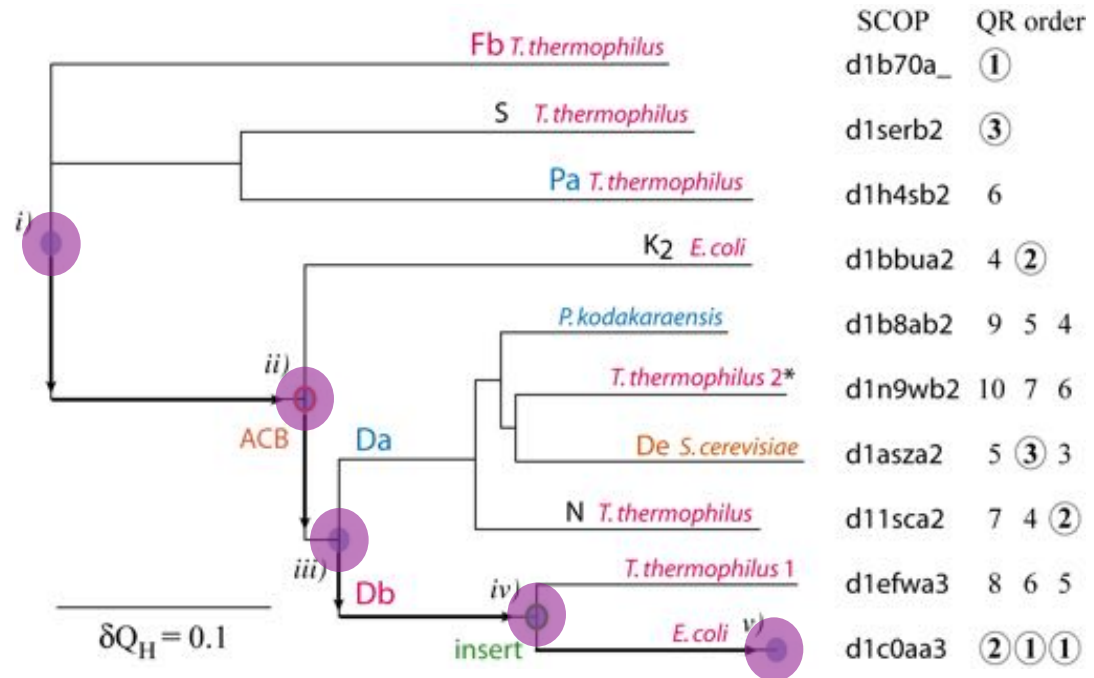*i)* class II

*ii)* subclass IIB

anticodon binding (ACB) domain

*iii)* AspRS

*iv)* bacterial AspRS

*v) E. coli* AspRS

bacterial insert domain

| SCOP | QR order |
|---|---|
| Fb *T. thermophilus* | d1b70a_ ① |
| S *T. thermophilus* | d1serb2 ③ |
| Pa *T. thermophilus* | d1h4sb2 6 |
| K₂ *E. coli* | d1bbua2 4 ② |
| *P. kodakaraensis* | d1b8ab2 9 5 4 |
| *T. thermophilus* 2* | d1n9wb2 10 7 6 |
| De *S. cerevisiae* | d1asza2 5 ③ 3 |
| N *T. thermophilus* | d1lsca2 7 4 ② |
| *T. thermophilus* 1 | d1efwa3 8 6 5 |
| *E. coli* | d1c0aa3 ② ① ① |

$\delta Q_H = 0.1$

# Summary Structural Evolutionary 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

# New Tools in VMD/MultiSeq



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

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

Sequence /Structure
Alignment

Protein & RNA
secondary structure

QR non-redundant
seq / str sets

Cluster
analysis /
Bioinformatics
scripting

Tutorials MultiSeq/
AARS
EF-Tu/Ribosome

# 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

# Load large sequence sets*

| Swiss-Prot (Proteins) | Greengenes (RNA)* |
|---|---|
| Curated sequences | Environmental 16S rRNA |
| 392,667 sequences | 90,654 entries |
| Unaligned | Aligned (7682 positions) |
| 177 MB on disk | 670 MB on disk |
| 2 minutes to load | 2.5 minutes to load * |
| 2.4 GB memory used | 4.0 GB memory used* |

*"Signatures of ribosomal evolution" with Carl Woese, PNAS (2008)
*Release May 2013 contains 1.2 million sequences – Memory??

# 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

# 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

# Scripting MultiSeq

- All MultiSeq functions can 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.

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

See also ITEP for microbial genomes, Benedict et al. **BMC Genomics 2014**



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.

# Molecular Signatures of Translation- Drug Targets



**16S rRNA**

I  II  III  IV

E. coli
T. thermophilus
H. marismortui

Ribosomal Signatures: Idiosyncrasies in rRNA and/or r-proteins characteristic of the domains of life

**69 ( 119)** & **6 (14)** in 16S ( 23S)

Number of Archaeal/Eukaryal Signatures

Drug

746  2610  2058

Number of Bacterial Signatures

**MultiSeq Zoom**

**23S rRNA**

I  II  III  IV  V  VI

E. coli
T. thermophilus
H. marismortui

E. Roberts, A. Sethi, J. Montoya, C. R. Woese & Z. Luthey-Schulten. *PNAS*
*"Molecular Signatures of Ribosomal Evolution" (2008)*

Kim,… Luthey-Schulten, Ha, and Woodson, *Nature* *"Protein-guided RNA dynamics during early ribosome assembly (2014)*

# 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

# MultiSeq: Display and Edit Metadata

- External databases are <span style="color:red">cross-referenced</span> to display <span style="color:red">metadata</span> such as taxonomy (lineage), data source (sp, <span style="color:red">Uniprot #</span>), EC, enzymatic function

- Changes to metadata should periodically be updated!!!

- <span style="color:red">Electronic Notebook</span>: Notes and annotations about a specific sequence or structure can be added – and saved



There were missing residues