# *Introduction to Evolutionary Concepts and VMD/MultiSeq - Part I*
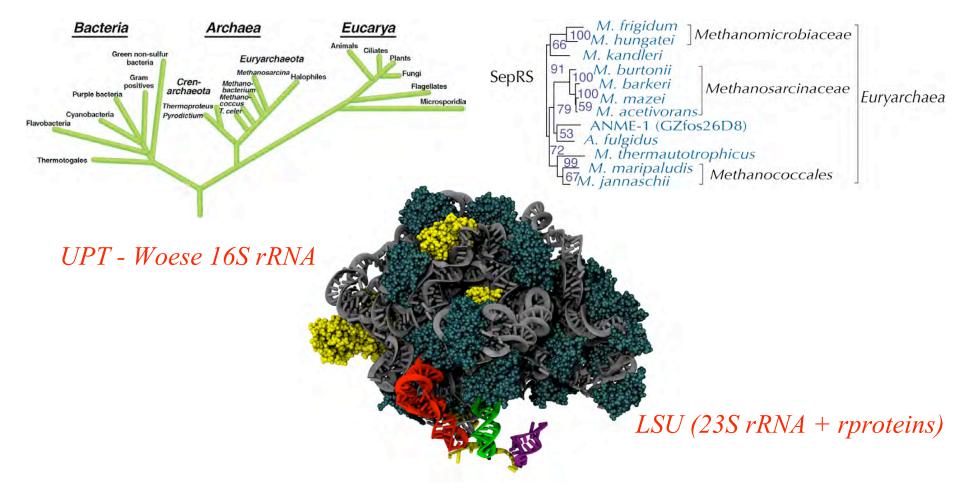
## Zaida (Zan) Luthey-Schulten

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

NIH Workshop 2009

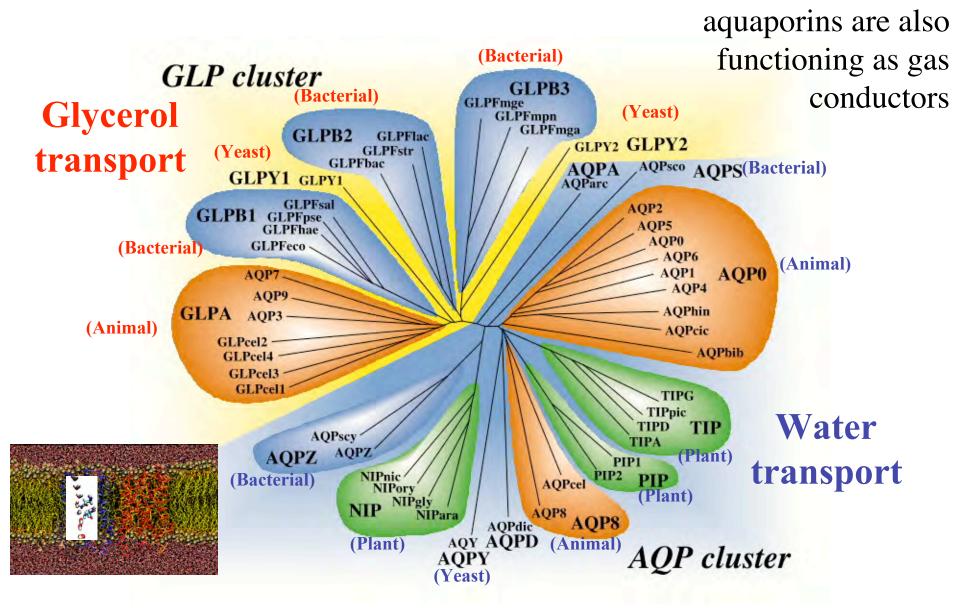ILLINOIS
UNIVERSITY OF ILLINOIS AT URBANA-CHAMPAIGN

# 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."*
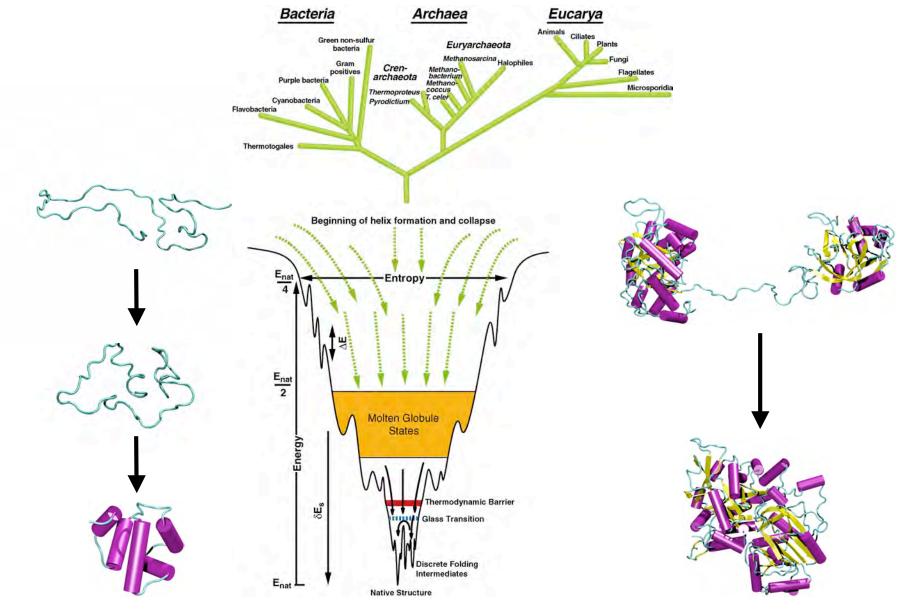


*UPT - Woese 16S rRNA*

*LSU (23S rRNA + rproteins)*

# The Aquaporin Superfamily

aquaporins are also functioning as gas conductors

Glycerol transport

GLP cluster

(Bacterial)
GLPB3
GLPFmge
GLPFmpn
GLPFmga
(Yeast)
GLPY2 GLPY2
AQPA AQPsco AQPS (Bacterial)
AQParc

(Bacterial)
GLPB2 GLPFlac
GLPFstr
GLPFbac
(Yeast)
GLPY1 GLPY1

GLPFsal
GLPB1 GLPFpse
GLPFhae
GLPFeco

(Bacterial)

AQP7
AQP9
GLPA AQP3

(Animal)

GLPcel2
GLPcel4
GLPcel3
GLPcel1

AQP2
AQP5
AQP0
AQP6
AQP1
AQP4
AQPhin
AQPcic
AQPbib

AQP0 (Animal)

TIPG
TIPpic
TIPD TIP
TIPA
(Plant)

Water transport

AQPscy
AQPZ AQPZ
(Bacterial)

NIPnic
NIPory
NIP NIPgly
NIPara
(Plant)

PIP1
PIP2 PIP
AQPcel (Plant)

AQP8
AQPdic AQP8
AQY AQPD (Animal)
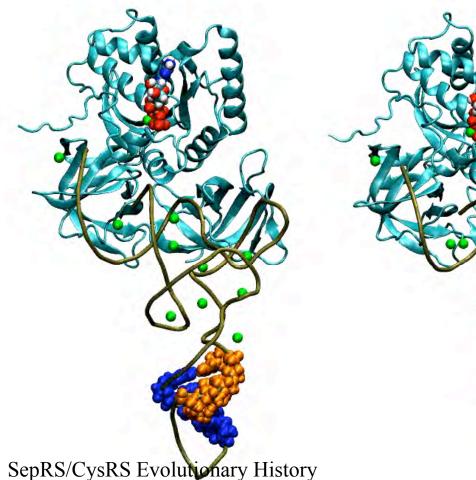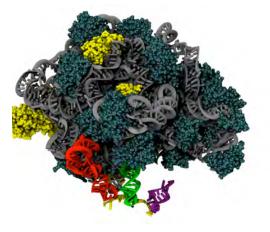AQPY AQPdic
(Yeast)

AQP cluster

Heymann and Engel *News Physiol. Sci.* **14**, 187 (1999)

# Evolution of Protein (RNA) Folding, Structure, & Function

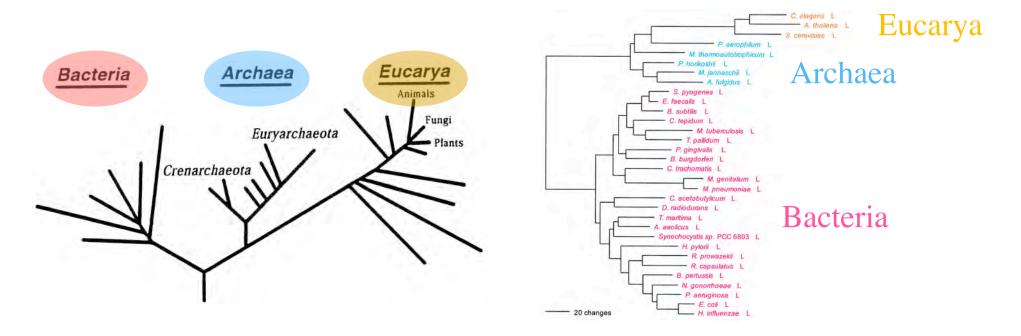# Evolution of Translational Machinery



SepRS/CysRS Evolutionary History
Sethi, O'Donoghue, ZLS, *PNAS* 2005
O'Donoghue, Sethi, Woese, ZLS,
 *PNAS* 2005, Sethi, et al. Dynamics of
Allosteric Network, **PNAS** 09

Dynamical Recognition EF-Tu:tRNA-
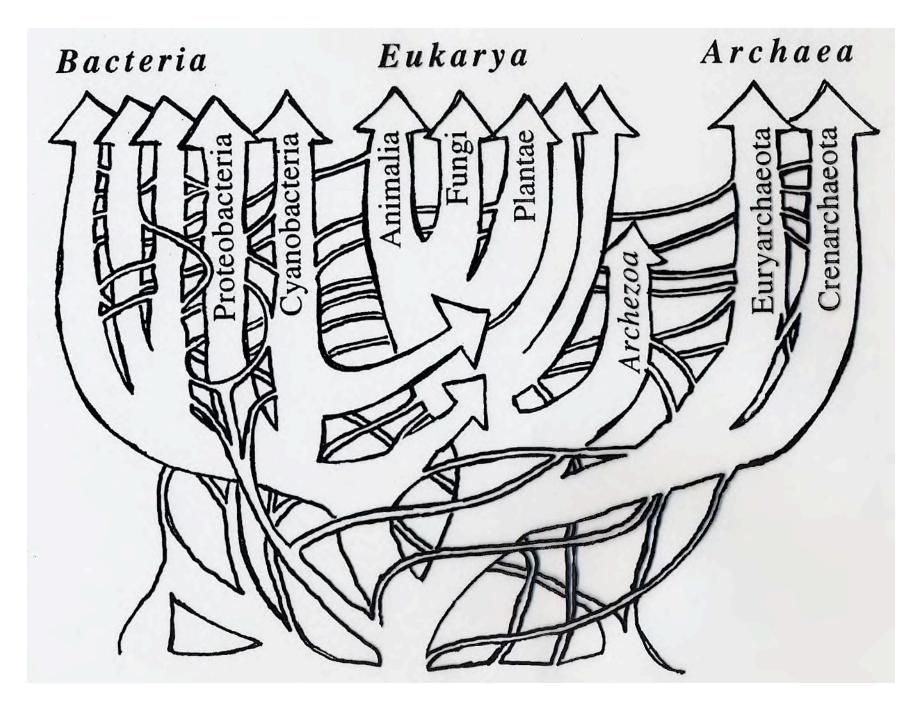Novel Amino Acids
Eargle et al., *JMB* 2008

Proteins/RNA
Ribosome

Molecular Signatures of
Ribosomal Evolution*,*
*PNAS* **2008,** Roberts,
Sethi, Woese, ZLS

# Universal Phylogenetic Tree
## three domains of life



Based on 16S rRNA



Eucarya

Archaea

Bacteria

Leucyl-tRNA synthetase displays the
full canonical phylogenetic distribution.

for review see Woese *PNAS* 2000

Woese, Olsen, Ibba, Soll *MMBR* 2000

**Bacteria** — Proteobacteria, Cyanobacteria

**Eukarya** — Animalia, Fungi, Plantae, Archezoa
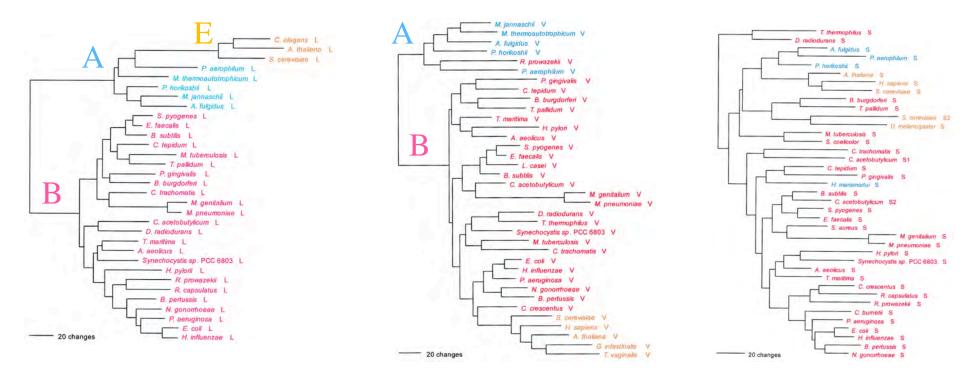
**Archaea** — Euryarchaeota, Crenarchaeota

After W. Doolittle, modified by G. Olsen

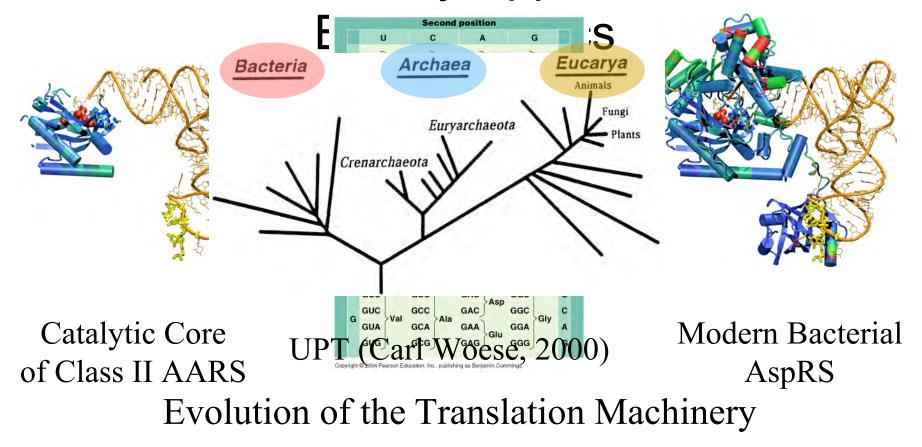# Phylogenetic Distributions



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
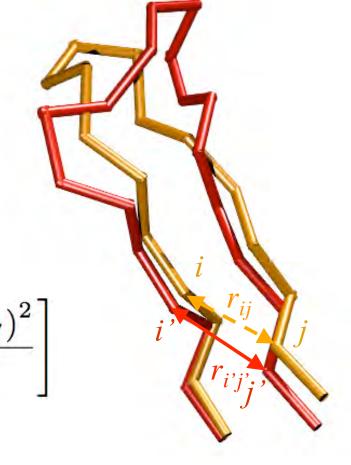
# MultiSeq in VMD

## Evolutionary Approach to Enzymes

Catalytic Core
of Class II AARS

UPT (Carl Woese, 2000)

Modern Bacterial
AspRS

Evolution of the Translation Machinery

# Protein Structure Similarity Measure

## Q<sub>H</sub> Structural Homology

_fraction of native contacts for aligned residues +_
_presence and perturbation of gaps_

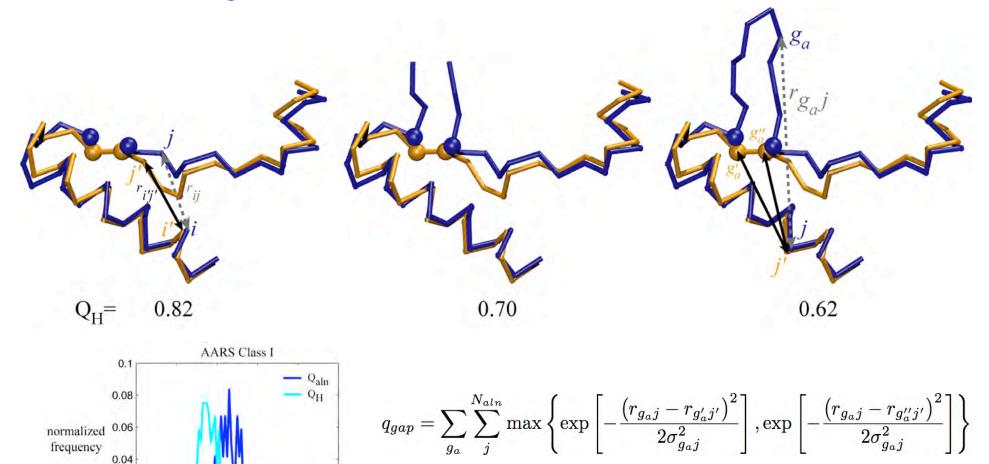$$Q_H = \aleph \left[ q_{aln} + q_{gap} \right]$$

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

# 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

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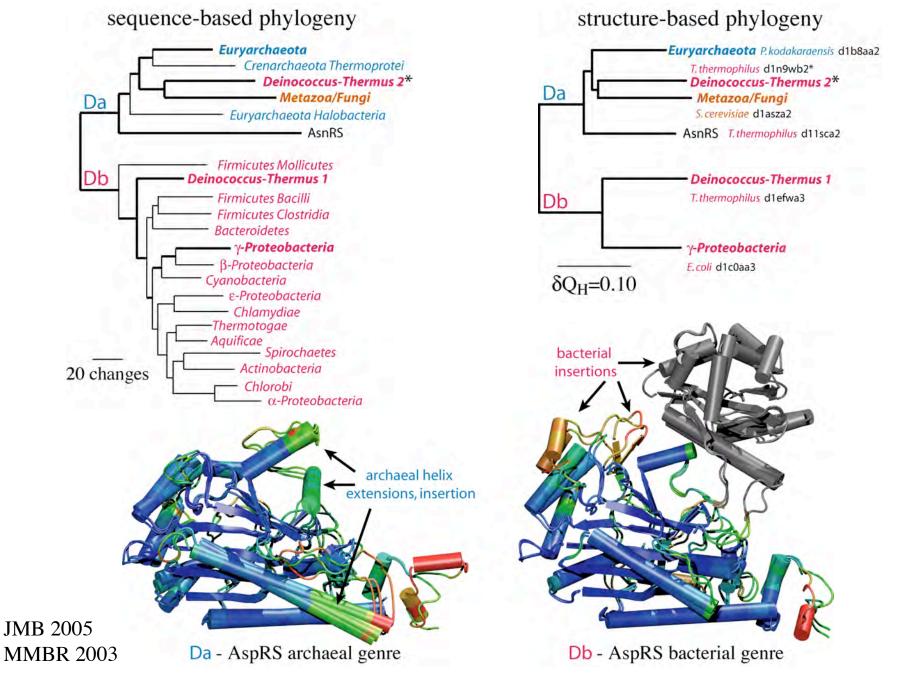
$$+ \sum_{g_b}\sum_{j}^{N_{aln}} \max\left\{ \exp\left[-\frac{(r_{g_b j} - r_{g'_b j'})^2}{2\sigma^2_{g_b j}}\right], \exp\left[-\frac{(r_{g_b j} - r_{g''_b j'})^2}{2\sigma^2_{g_b j}}\right]\right\}$$

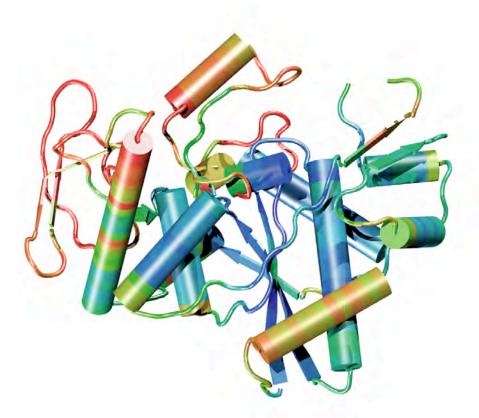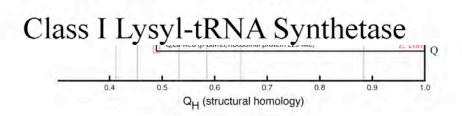# Protein structure encodes evolutionary information



sequence-based phylogeny

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

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

20 changes

structure-based phylogeny

Da
- **Euryarchaeota** *P. kodakaraensis* d1b8aa2
- *T. thermophilus* d1n9wb2*
- **Deinococcus-Thermus 2***
- **Metazoa/Fungi**
- *S. cerevisiae* d1asza2
- AsnRS *T. thermophilus* d11sca2

Db
- **Deinococcus-Thermus 1**
- *T. thermophilus* d1efwa3
- γ-**Proteobacteria**
- *E. coli* d1c0aa3

$\delta Q_H = 0.10$

archaeal helix extensions, insertion

Da - AspRS archaeal genre

bacterial insertions

Db - AspRS bacterial genre

JMB 2005
MMBR 2003

# Protein structure reveals distant evolutionary events



Class I AARSs

Class II AARSs

Class I Lysyl-tRNA Synthetase

Class II Lysyl-tRNA Synthetase
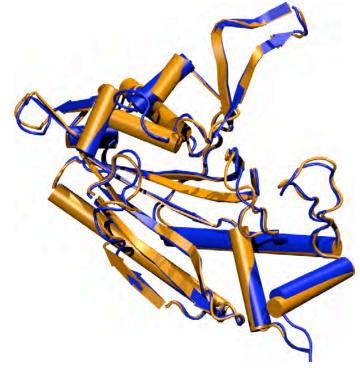
# Sequences define more recent evolutionary events



Conformational changes
in the same protein.

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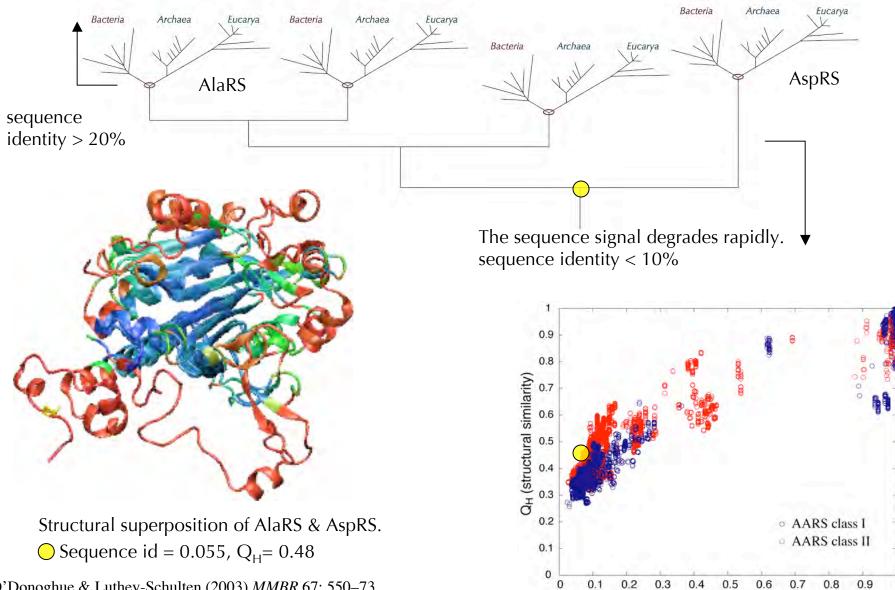T, 2.00 A.

$Q_H = 0.80$
Sequence identity = 1.00



Structures for two
different species.

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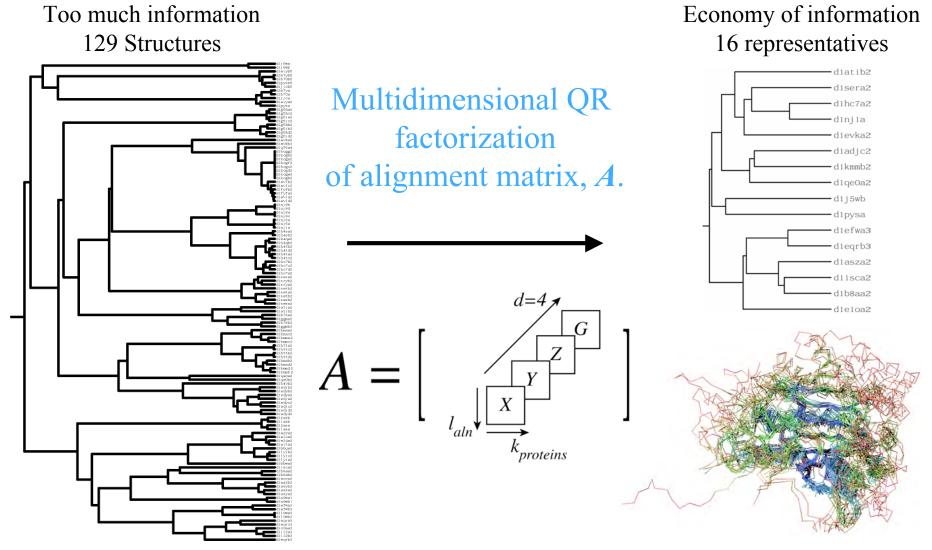*M. thermoautotrophicus*, 3.20 A.

$Q_H = 0.89$
Sequence identity = 0.69

# The Relationship Between Sequence & Structure



Bacteria  Archaea  Eucarya  Bacteria  Archaea  Eucarya  Bacteria  Archaea  Eucarya  Bacteria  Archaea  Eucarya

AlaRS

AspRS

sequence
identity > 20%

The sequence signal degrades rapidly.
sequence identity < 10%

Structural superposition of AlaRS & AspRS.

⬤ Sequence id = 0.055, $Q_H$ = 0.48

$Q_H$ (structural similarity)

AARS class I
AARS class II

sequence identity

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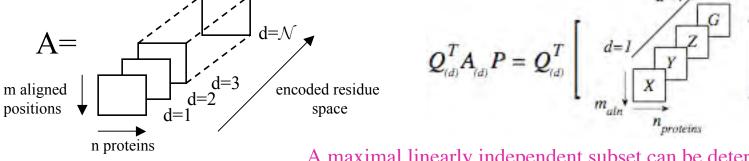
# Non-redundant Representative Sets

Too much information
129 Structures

Economy of information
16 representatives



Multidimensional QR
factorization
of alignment matrix, $A$.

$$A = \begin{bmatrix} & & & G \\ & & Z & \\ & Y & & \\ X & & & \end{bmatrix}$$

with $d=4$, $l_{aln}$, $k_{proteins}$

QR computes a set of maximal linearly independent structures.

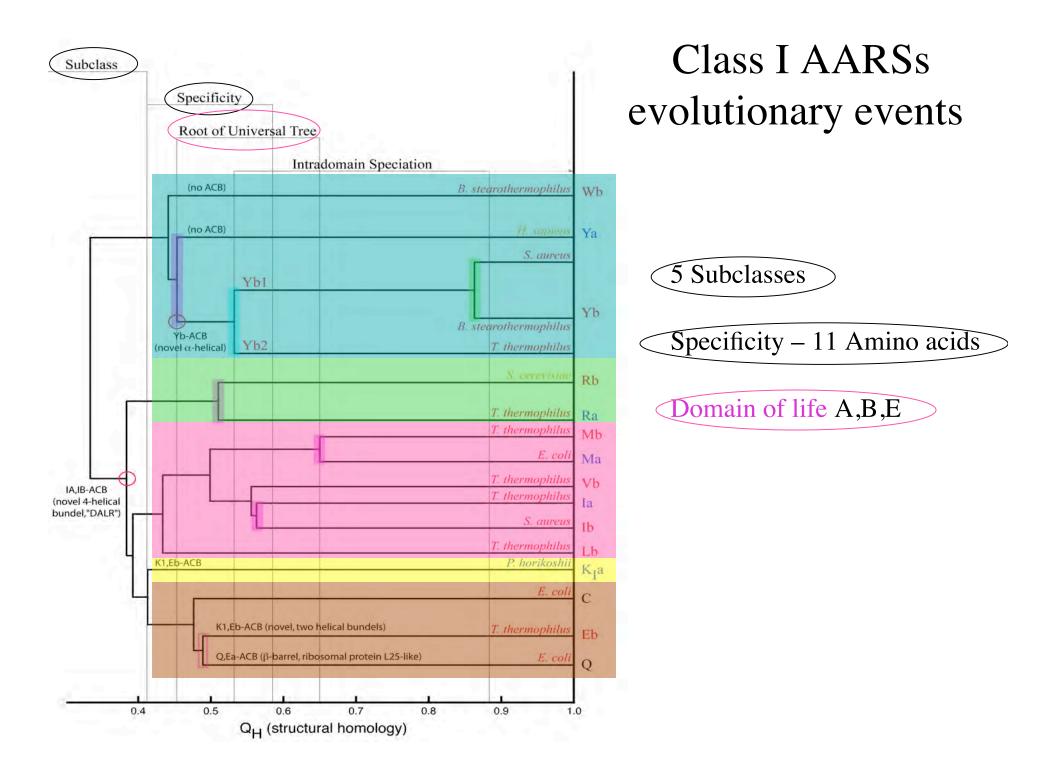P. O'Donoghue and Z. Luthey-Schulten (2003) *MMBR* **67**:550-571.
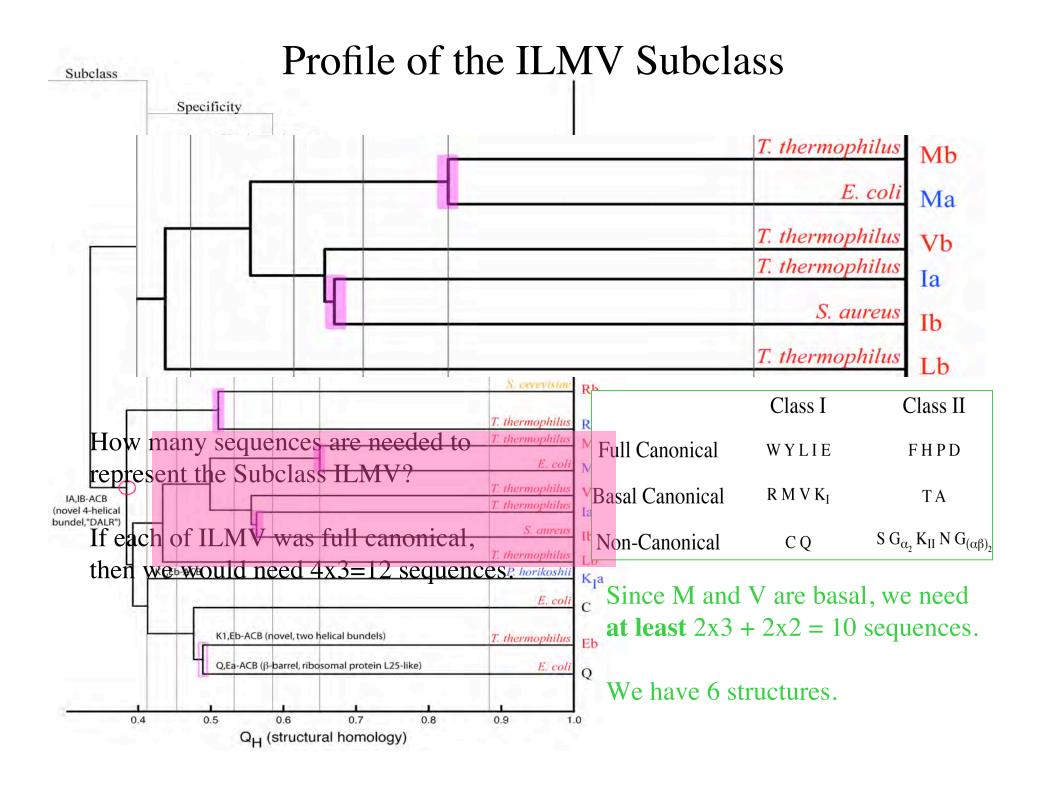
P. O'Donoghue and Z. Luthey-Schulten (2005) *J. Mol. Biol.*, **346**, 875-894.
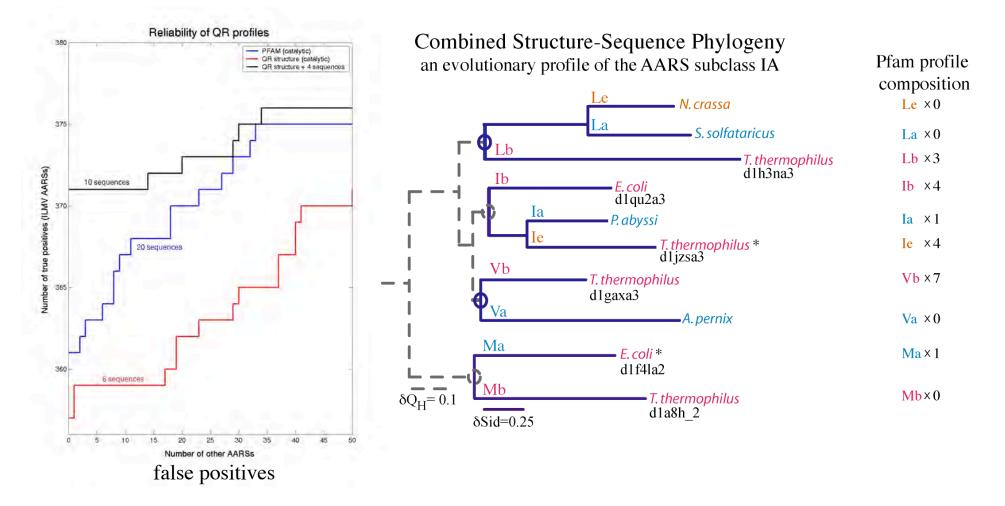
# Numerical Encoding of Proteins in a Multiple Alignment

## 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



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

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

Subclass

Specificity

| T. thermophilus | Mb |
| E. coli | Ma |
| T. thermophilus | Vb |
| T. thermophilus | Ia |
| S. aureus | Ib |
| T. thermophilus | Lb |

How many sequences are needed to represent the Subclass ILMV?

If each of ILMV was full canonical, then we would need 4x3=12 sequences.

IA,IB-ACB (novel 4-helical bundel,"DALR")

K1,Eb-ACB (novel, two helical bundels)

Q,Ea-ACB (β-barrel, ribosomal protein L25-like)

|  | Class I | Class II |
|---|---|---|
| Full Canonical | W Y L I E | F H P D |
| Basal Canonical | R M V $K_I$ | T A |
| Non-Canonical | C Q | S $G_{\alpha_2}$ $K_{II}$ N $G_{(\alpha\beta)_2}$ |

Since M and V are basal, we need **at least** 2x3 + 2x2 = 10 sequences.

We have 6 structures.

*S. cerevisiae* Rb
*T. thermophilus* R
*T. thermophilus* M
*E. coli* M
*T. thermophilus* V
*T. thermophilus* Ia
*S. aureus* Ib
*T. thermophilus* Lb
*P. horikoshii* $K_I$a
*E. coli* C
*T. thermophilus* Eb
*E. coli* Q

0.4    0.5    0.6    0.7    0.8    0.9    1.0

$Q_H$ (structural homology)

# Evolutionary Profiles for Homology Recognition
## AARS Subclass ILMV



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

# 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

ACB

Da

Db

insert

$\delta Q_H = 0.1$

| | 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 |

# 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

# Structural Domains

## Protein: Aspartyl–tRNA synthetase (AspRS) from *Escherichia coli*

## Lineage:

1. Root: scop
2. Class: All beta proteins
3. Fold: OB–fold
   *barrel, closed or partly opened n=5, S=10 or S=8; greek–key*
4. Superfamily: Nucleic acid–binding proteins
5. Family: Anticodon–binding domain
   *barrel, closed; n=5, S=10*
6. Protein: Aspartyl–tRNA synthetase (AspRS)
   *this is N–terminal domain in prokaryotic enzymes and the first "visible" domain in eukaryotic enzymes*
7. Species: *Escherichia coli*

## PDB Entry Domains:

1. 1c0a
   1. region a:1–106
2. 1il2
   *complexed with 1mg, 5mc, 5mu, amo, h2u, psu, so4*
   1. region a:1–106
   2. region b:1001–1106
3. 1eqr
   *complexed with mg*
   1. region a:1–106
   2. region b:1–106
   3. region c:1–106

# Profile - Multiple Structural Alignments

Representative Profile of AARS Family

Catalytic Domain

# 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 bewteen 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.

# 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.

# Planned Tools in MultiSeq



Protein / RNA Sequence Data

Entire SwissProt DB, 100,000+ RNA seqs

Metadata Information, Clustal & Phylogenetic Trees

Incorporate genomic content

Blast & PsiBlast

Sequence Editor

Sequence / Structure Alignment

RNA Secondary Structure

QR non-redundant seq / str sets

Cluster analysis / Bioinformatics scripting

Tutorials MultiSeq/AARS
EF-Tu/Ribosome

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

Import data directly from BLAST databases

Align sequences with Clustal

Sequence Editor: Manually adjust alignments or sequences

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)

# 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 ($Q_{res}$, $Q_H$,..) or imported experimental properties



- Integrates popular bioinformatics tools making them easier to use and reducing the barrier to performing bioinformatics analysis (ClustalW, STAMP, BLAST)

# 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

# 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)

# Protein sequence alignment
# How do I align two similar, but different sequences

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

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

*There exist web accessible tools, e.g., BLAST search:* **http://www.ncbi.nlm.nih.gov/**

# NiceProt View of Swiss-Prot:
## P47865

(Printer-friendly view) (Submit update) (Quick BlastP search)

[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**<br>**Aquaporin 1**<br>**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

**Scroll down to the sequence:**

# Final Result: Sequence Alignment

```
>gi|46395801|sp|Q88F17|AQPZ_PSEPK  [G]  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    VAFAFGLTVLTMFAIGHISGCHLNPAVSFGLVVGGRFPAKELLPYVIAQVIGAILAAGV 99

Query: 113   LSGITSSLP--DNSLGL--NALAP----GVNSGQGLGIEIIGTLQLVLCVLATTDRRRRD 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 method returns approximate alignments - needing refinement

# 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 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

# MultiSeq Tutorials