# Evolution of Protein Structure



# Luthey-Schulten Group

Department of Chemistry, Biophysics, and Beckman Institute
University of Illinois at Urbana-Champaign

# Universal Phylogenetic Tree
## three domains of life



**Bacteria**  **Archaea**  **Eucarya**
Animals

Fungi

Plants

Euryarchaeota

Crenarchaeota

Eucarya

Archaea

Bacteria

C. elegans  L
A. thaliana  L
S. cerevisiae  L
P. aerophilum  L
M. thermoautotrophicum  L
P. horikoshii  L
M. jannaschii  L
A. fulgidus  L
S. pyogenes  L
E. faecalis  L
B. subtilis  L
C. tepidum  L
M. tuberculosis  L
T. pallidum  L
P. gingivalis  L
B. burgdorferi  L
C. trachomatis  L
M. genitalium  L
M. pneumoniae  L
C. acetobutylicum  L
D. radiodurans  L
T. maritima  L
A. aeolicus  L
Synechocystis sp. PCC 6803  L
H. pylorii  L
R. prowazekii  L
R. capsulatus  L
B. pertussis  L
N. gonorrhoeae  L
P. aeruginosa  L
E. coli  L
H. influenzae  L

20 changes

## Based on 16S rRNA

Leucyl-tRNA synthetase displays the
full canonical phylogenetic distribution.

for review see Woese *PNAS* 2000

Woese, Olsen, Ibba, Soll *MMBR* 2000

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

# Protein Structure Similarity Measure

## $Q_H$ Structural Homology

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

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

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



O'Donoghue & Luthey-Schulten  *MMBR* 2003.

# 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

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

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

# Protein structure encodes evolutionary information



sequence-based phylogeny

structure-based phylogeny

Da

Db

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

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

Da

Db

*Euryarchaeota* P. kodakaraensis d1b8aa2
T. thermophilus d1n9wb2*
**Deinococcus-Thermus 2***
**Metazoa/Fungi**
S. cerevisiae d1asza2
AsnRS T. thermophilus d1l1sca2

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

γ-**Proteobacteria**
E. coli d1c0aa3

$\delta Q_H = 0.10$

bacterial insertions

archaeal helix extensions, insertion

Da - AspRS archaeal genre

Db - AspRS bacterial genre

JMB 2005
MMBR 2003

# Protein structure reveals distant evolutionary events



Class I AARSs

Class II AARSs

# 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

# Non-redundant Representative Sets

Too much information
129 Structures

Economy of information
16 representatives

Multidimensional QR
factorization
of alignment matrix, $A$.



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

d1atib2
d1scra2
d1hc7a2
d1nj1a
d1evka2
d1adjc2
d1kmmb2
d1qe0a2
d1j5wb
d1pysa
d1efwa3
d1eqrb3
d1asza2
d1lsca2
d1b8aa2
d1eloa2

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

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

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

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

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

$d=4$, $d=1$, G, Z, Y, X, $m_{aln}$, $n_{proteins}$

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

# Applications of Evolutionary Profiles

**I.** **Genome Annotation** AARS -  MJ1660

**II.** **Conserved Core -- Folding Nuclei? HD Exchange?**

**III. Functional Ancestor ?**

**IV. Classification of Protein Structures** - Superfamilies

Class I AARSs evolutionary events

5 Subclasses

Specificity – 11 Amino acids

Domain of life A,B,E

# Profile of the ILMV Subclass?

# Profile of the ILMV Subclass



*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.

|  | 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.

# Evolutionary Profiles for Homology Recognition
## AARS Subclass ILMV



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

A. Sethi, P. O'Donoghue, Z. Luthey-Schulten (2005) *JMB, PNAS*

# Genome Annotation

*M.jannaschii* genome was completely sequenced in 1996.
Genome had four missing AARSs:

AsnRS
GlnRS } Indirect Mechanism

LysRS    Class I AARS

CysRS    ?

Cysteinyl-tRNA(Cys) formation in *Methanocaldococcus jannaschii*: the mechanism is still unknown. *J. Bacteriology*, Jan. 2004, **186:**8-14.
Ruan B, Nakano H, Tanaka M, Mills JA, DeVito JA, Min B, Low KB, Battista JR, and Söll D.

| Protein | E-value |
|---|---|
| HisRS | 1.1e-10 |
| AspRS | 1.9e-10 |
| PheRS α-chain | 9.5e-10 |
| ThrRS | 6.6e-04 |
| ProRS | 9.1e-03 |
| SerRS | 9.2e-03 |
| putative CysRS | 1.6e-02 ← MJ1660 |
| AlaRS | 5.1e-02 |
| GlyRS | 0.12 |
| PheRS β-chain | 0.15 |
| DNA repair protein | 7.5 |

*M. jannaschii* genome database search using EP of class II AARS with HMMER

Sethi, et. al., PNAS, **102**, 2005

# Connections of Direct and Indirect Pathways for Cysteinyl-tRNA formation to Cysteine Biosynthesis



Sauerwald et al., Science, 307, 2005, 1969-1972.

# Genes for Cysteine Biosynthesis and Aminoacylation

| | Cys coding | | Cys biosynthesis | | | Cys biosynthesis/coding | |
|---|---|---|---|---|---|---|---|
| | CysRS | CysE | CysK/M | CBS | CGL | SepRS | SepCysS |
| **Crenarchaea** | | | | | | | |
| Aeropyrum pernix | NP_148045 | - | NP_148041 | NP_147802 | NP_147803 | - | - |
| Sulfolobus solfataricus | NP_343652 | - | (NP_341900) | (NP_341900) | (NP_343729) | - | - |
| Sulfolobus tokodaii | NP_378245 | - | (NP_377338) | (NP_377338) | (NP_376392) | - | - |
| Pyrobaculum aerophilum | NP_558873 | (NP_559322) | (NP_559045) | (NP_559045) | (NP_559999) | - | - |
| **Euryarchaea** | | | | | | | |
| Haloarcula marismortui | YP_135935 | YP_135755 | YP_134915 | (YP_135866) | (YP_136993) | - | - |
| Halobacterium sp. | NP_280014 | NP_280304 | NP_280167 | NP_279635 | (NP_279780) | - | - |
| Methanothermobacter thermautotrophicus | - | - | - | - | - | NP_276615 | NP_276195 |
| Methanocaldococcus jannaschii | - | - | - | - | - | NP_248670 | NP_248688 |
| Methanococcus maripaludis | NP_988180 | - | - | - | - | NP_987808 | NP_988360 |
| Methanopyrus kandleri | - | - | - | - | - | NP_613724 | NP_613516 |
| Methanosarcina acetivorans | NP_615709 | NP_617620 | NP_617619 | - | (NP_617435) | NP_615064 | NP_615682 |
| Methanosarcina barkeri | AAF18751 | 40160510* | AAF07039 | - | - | ZP_00298242 | ZP_00297376 |
| Methanosarcina mazei | NP_633935 | NP_635293 | - | - | NP_635109 | NP_633407 | NP_633905 |
| Methanosarcina thermophila | ? | AAG01805 | AAG01804 | ? | ? | ? | ? |
| Methanococcoides burtonii | ? | ZP_00149388 | ZP_00149387 | ? | ? | ZP_00147576 | ZP_00148017 ZP_00148733 |
| Methanospirillum hungatei | 401798240* | 401798540* | 401798280* | ? | ? | 40179880* | 401798260* |
| Methanogenium frigidum | ? | ? | Contig384.gene842** | ? | ? | Contig1085.gene108** | Contig1260.gene378** |
| Pyrococcus abyssi | NP_127080 | NP_126842 | (NP_126065) | (NP_126065) | (NP_126586) | - | - |
| Pyrococcus furiosus | NP_578753 | NP_578497 | (NP_578587) | (NP_578587) | NP_578995 | - | - |
| Pyrococcus horikoshii | NP_142595 | - | - | - | NP_142999 | - | - |
| Ferroplasma acidarmanus | 401193730* | ? | ZP_0306996 | ? | ? | ? | ? |
| Thermoplasma acidophilum | NP_394604 | - | (NP_394010) | (NP_394010) | NP_393559 | - | - |
| Thermoplasma volcanium | NP_111763 | - | (NP_111108) | (NP_111108) | (NP_110693) | - | - |
| Picrophilus torridus | YP_022862 | - | YP_022929 | (YP_023731) | (YP_023880) | - | - |
| Archaeoglobus fulgidus | NP_069247 | - | - | - | - | NP_068951 | NP_068869 NP_069020 |
| **Nanoarchaea** | | | | | | | |
| Nanoarchaeum equitans | NP_069247 | - | - | - | - | - | - |

*gene object identifiers from Integrated Microbial Genomes database at JGI.
**M. frigidum draft genome sequence, Saunders et al. (2003) Gen. Res. 13, 1580–1588.
All other codes are NCBI-NR database gene identifiers. - absence of gene. ? absence of gene in incomplete genome.

Evolutionary history of SepRS

Same pattern as euryarchaeal portion of rRNA tree.

Was present in LUCAS.

P. O' Donoghue, A. Sethi, C. Woese, and Z. Luthey-Schulten, PNAS, 2005.

# Evolutionary history of class I CysRS



Bacterial groupings from UPT tree also seen in CysRS phylogeny.

Multiple HGT events of the direct route to archaeal organisms.

The direct route to cysteine aminoacylation was also present in the LUCAS.

# Evolution of Structure and Function in AspRS



*i)* class II

*ii)* subclass IIB

anticodon binding (ACB) 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 |

ACB  Da

$\delta Q_H = 0.1$

Db  insert

*iii)* AspRS

*iv)* bacterial AspRS

bacterial insert domain

*v)* E. coli AspRS

# Evolutionary profile for HisA-HisF family



EP outperforms popular profile methods with an economy of information.

Sethi, et. al., PNAS, 2005.

# Economy of Information

## How many sequences are needed for profiles?



A. Sethi, P. O'Donoghue, ZLS, PNAS **102**, 2005

# Reclassification of TIM barrel Superfamilies ?



PLP-binding barrel c.1.6

d2toda2
d7odca2
d1d7ka2
d1ct5a_
d1bd0a2

FMN   c.1.4
d1gtea2
d2dora_

FMN   c.1.4
d1huva_
d1gox__

IMPDH   c.1.5
d1eepa_
d1b3oa1

TIM   c.1.1
d1mo0a_
d1tph1_

OMP decarboxylase   c.1.2.3
d1eixa_
d1dbta_

TrpA   c.1.2.4
d1ttqa_
d1geqa_

HisA   c.1.2.1
d1qo2a_

HisF   c.1.2.1
d1gpwc_
d1ka9f_

TrpF   c.1.2.4
d1pii_2
d1nsj__

RPE   c.1.2.2
d1h1ya_

TrpC   c.1.2.4
d1jcmp_
d1i4na_

TPS   c.1.3
d1g4pb_
d2tpsa_

0.4   0.5   0.6   0.7   0.8   0.9   1.0

$Q_H$ (structural similarity)

# Unifying the Worlds of Sequence and Structure

# Multiseq in VMD : Merging the sequence and structure worlds



Version 1.83

# 2006 MultiSeq: New Features

## Analyze the Evolution of Sequence and Structure



## Eliminate Redundancy



## Plus More Functions

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

# Acknowledgements

Patrick O'Donoghue

Anurag Sethi

Rommie Amaro
Felix Autenrieth
Alexis Black
**John Eargle**
Corey Hardin
Taras Pogorelov
**Elijah Roberts**
**Dan Wright**

## Graphics Programmers VMD

Elijah Roberts, Dan Wright, John Eargle

John Stone

## Collaborators

Evolutionary Studies
Gary Olsen, Carl Woese (UIUC)
QR Algorithms
Mike Heath (UIUC)
Protein Structure Prediction
Peter Wolynes, Jose Onuchic (UCSD)
Ken Suslick (UIUC)

## Funding