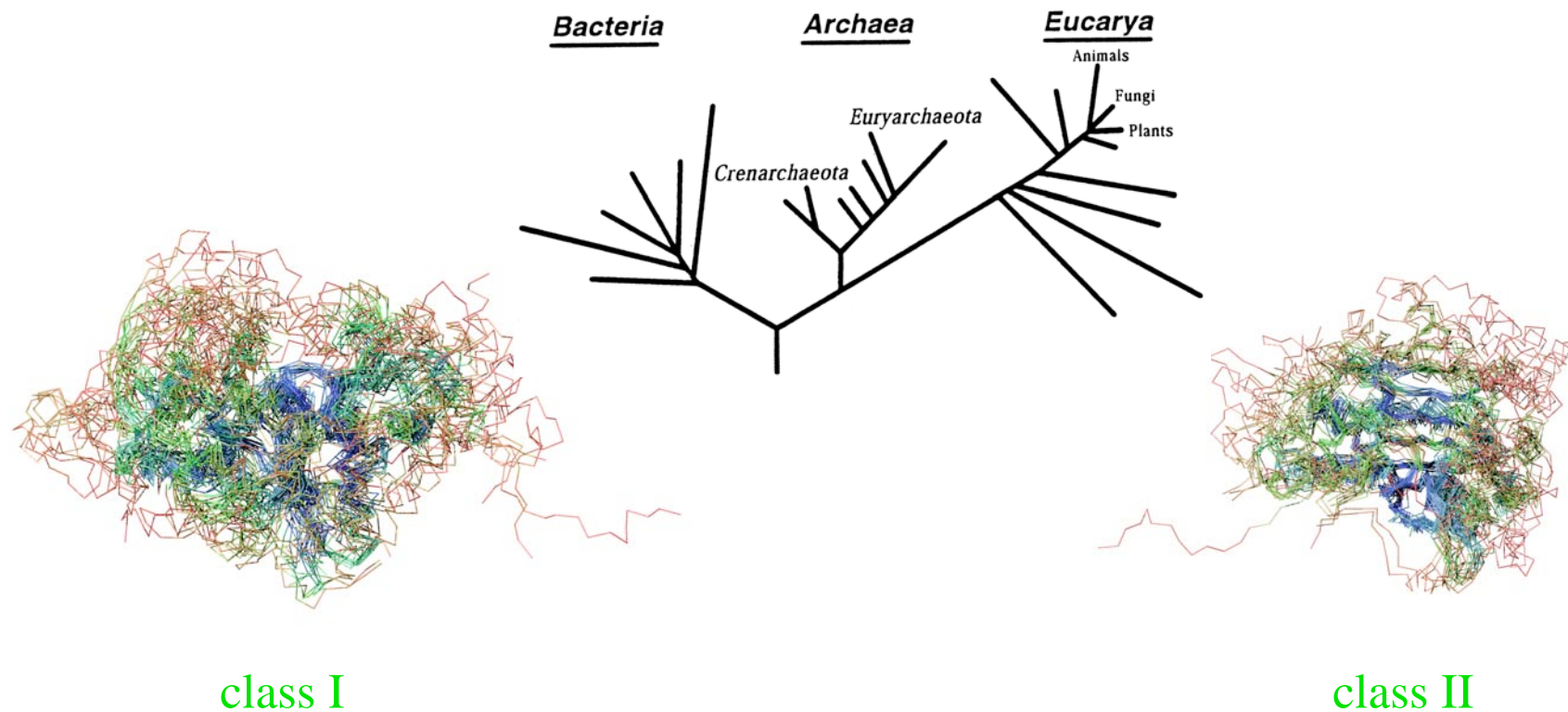


Evolution of Protein Structure in the Aminoacyl-tRNA Synthetases



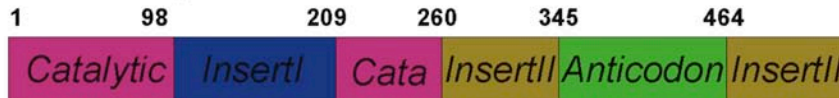
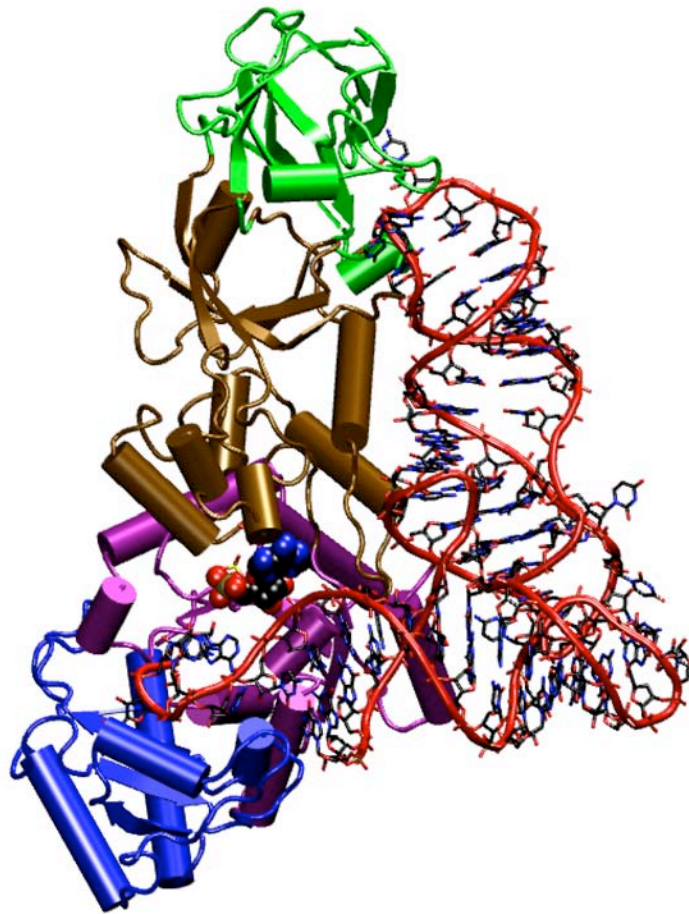
P. O'Donoghue and Z. Luthey-Schulten*

Department of Chemistry, Beckman Institute,
Center for Biophysics and Computational Biology
University of Illinois at Urbana-Champaign

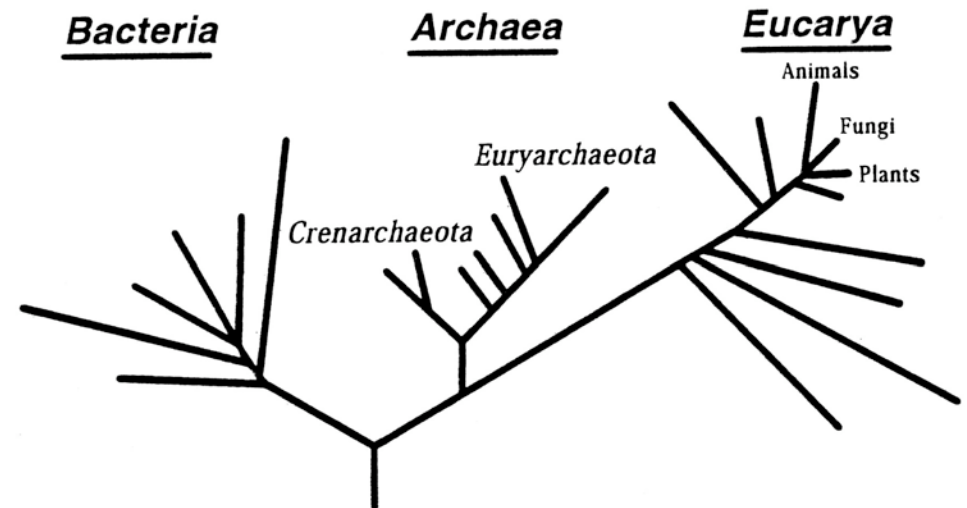
What can be learned from AARS?

- “The aminoacyl-tRNA synthetases, perhaps better than any other molecules in the cell, optimize the current situation and help to understand (the effects) of HGT” Woese (PNAS, 2000; MMBR 2000)

Aminoacyl-tRNA synthetases

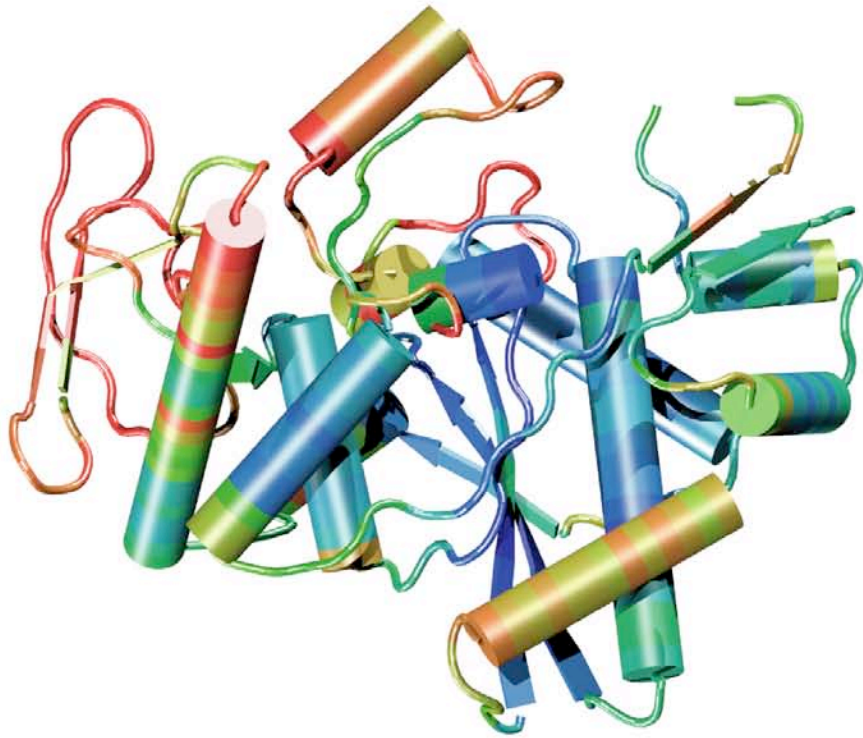


Universal Tree of Life

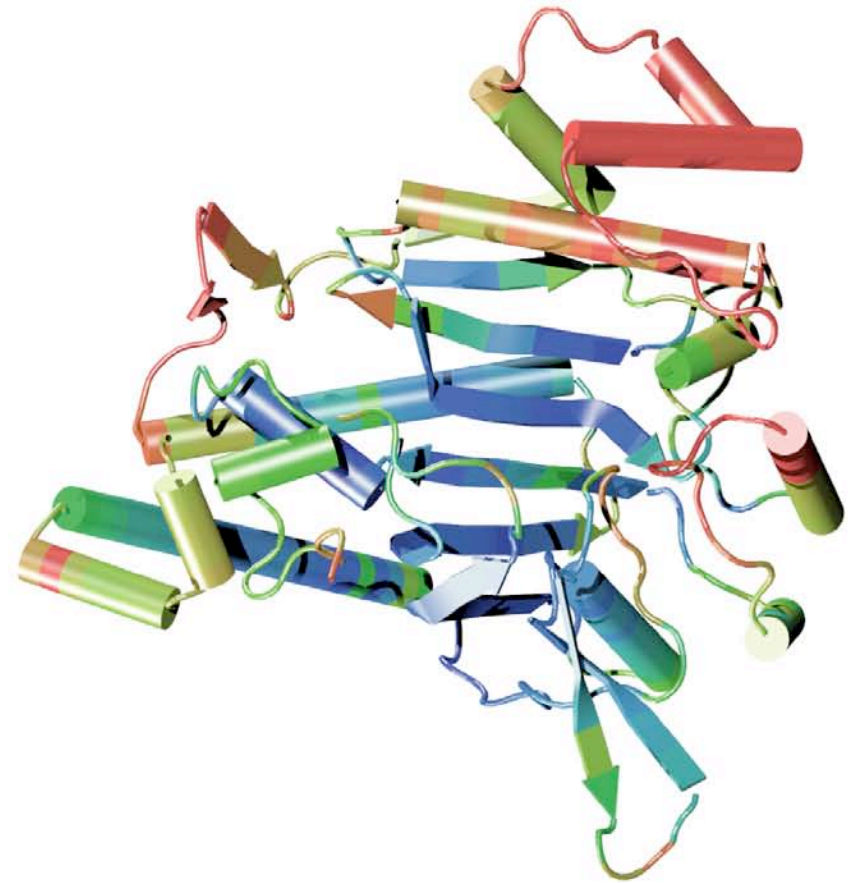


Woese *PNAS* 1990, 2002.

Structural Conservation in the Catalytic Domain of the AARSs



Class I Lysyl-tRNA Synthetase

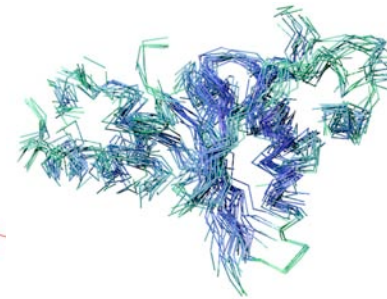
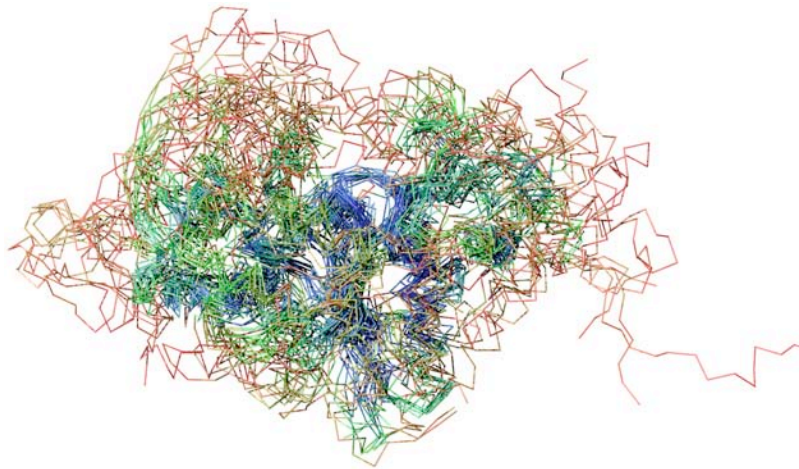


Class II Lysyl-tRNA Synthetase

Why Study the Evolution of Protein Structure?

1. Important for Homology Modeling

Better profiles improve database searches and give better alignments of distant homologs.
Allows mixing of sequence and structure information systematically.



13% sequence id
in the core (blue)

2. Learn how evolutionary dynamics changed protein shape.

Mapping a protein of unknown structure onto a homologous protein of known structure is equivalent to defining the evolutionary pathway connecting the two proteins

3. Impact on protein structure prediction, folding, and function

Evolutionary profiles increase the signal to noise ratio - Evolution is the foundation of bioinformatics.

Outline

1. Summarize evolutionary theory of the universal phylogenetic tree.

Methods

2. Introduce a structure-based metric which accounts for gaps, and show that evolutionary information is encoded in protein structure.
3. Introduce multidimensional QR factorization for computing non-redundant representative multiple alignments in sequence or structure.

Applications

4. Non-redundant multiple alignments which well represent the evolutionary history of a protein group provide better profiles for database searching.

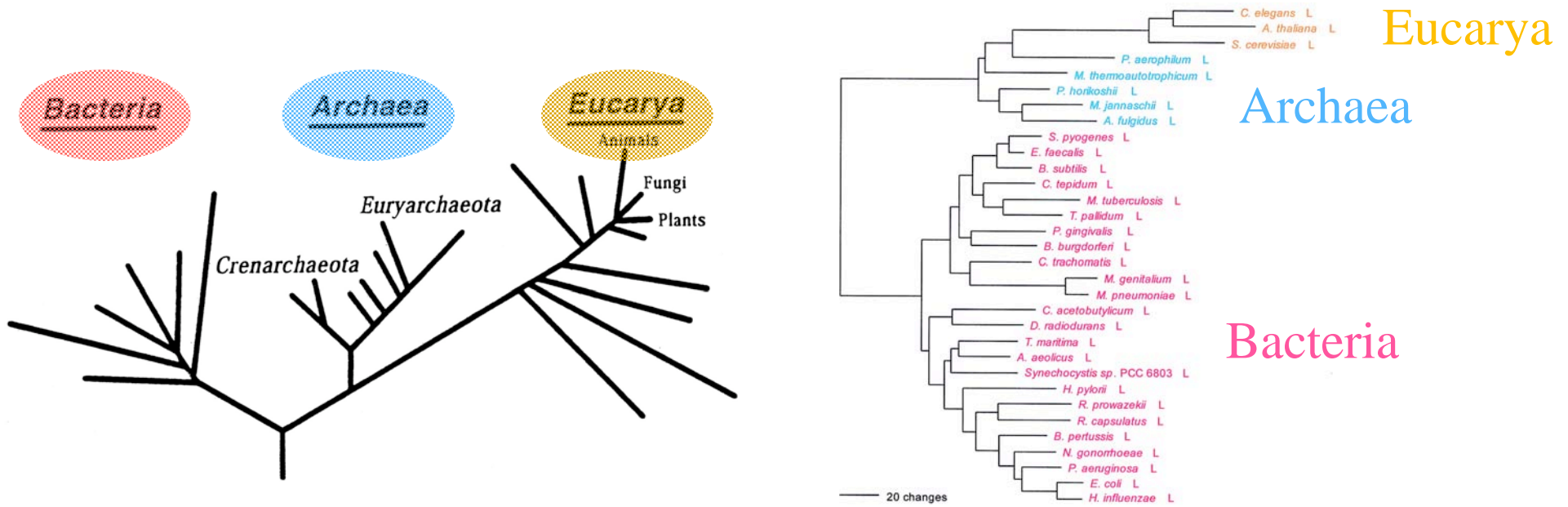
Eliminate bias inherited from structure or sequence databases.

Important for bioinformatic analysis (substitution matrices, knowledge based potentials structure pred., genome annotation) and evolutionary analysis.

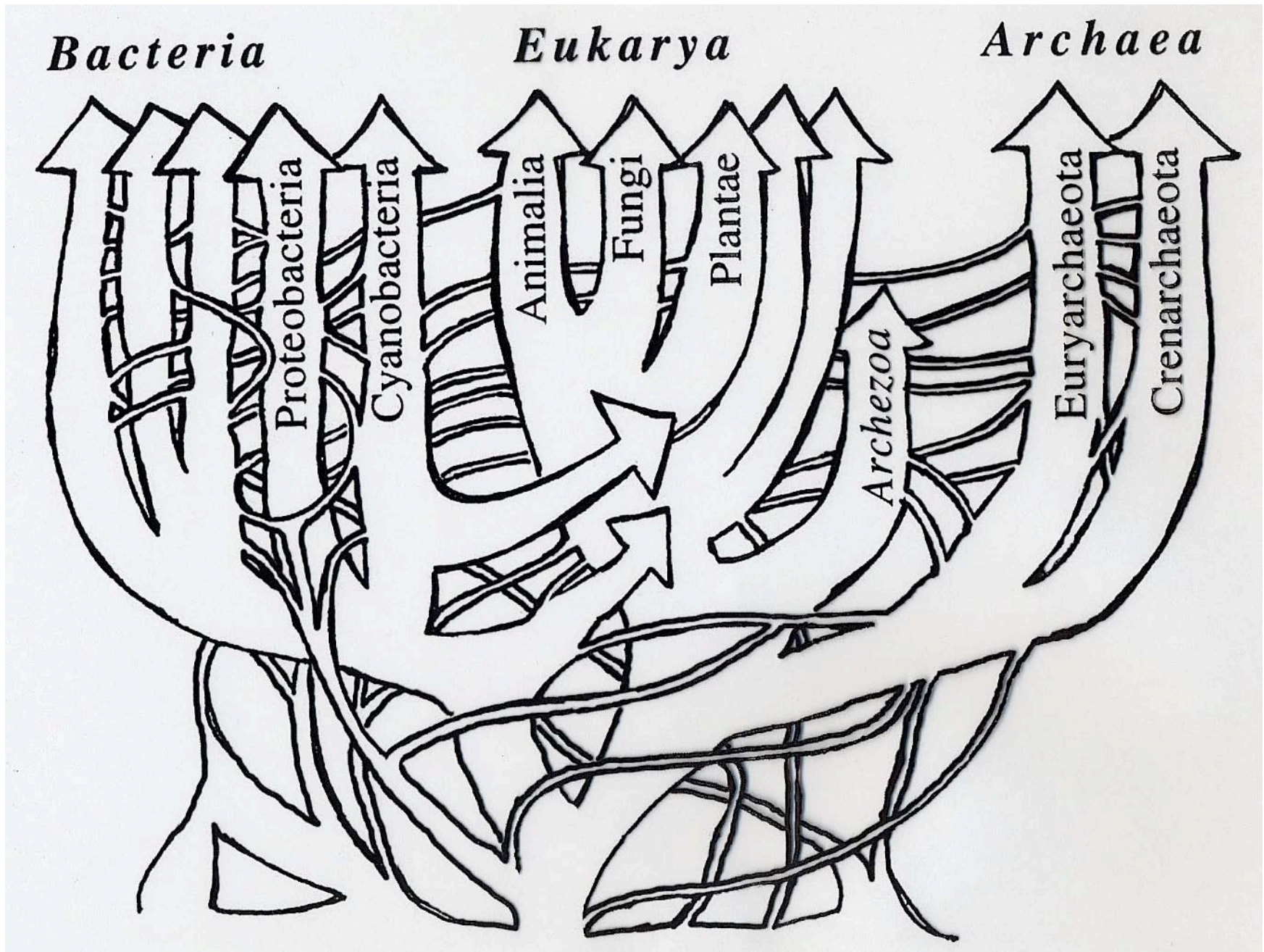
5. Depict the evolution of structure and function in Aspartyl-tRNA synthetase.

Universal Phylogenetic Tree

three domains of life



Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution.



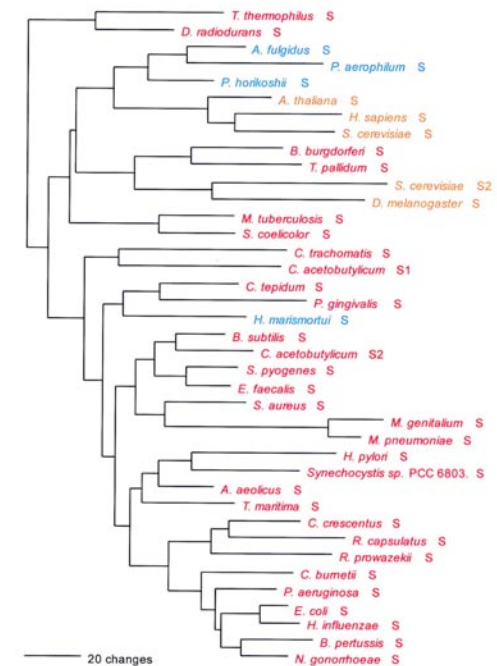
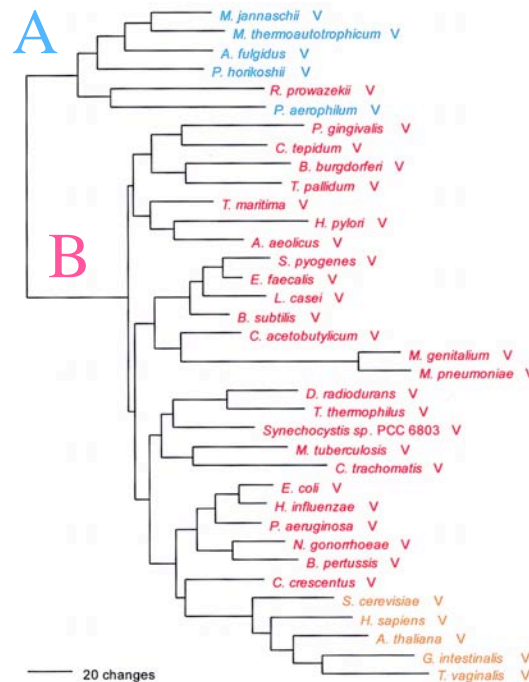
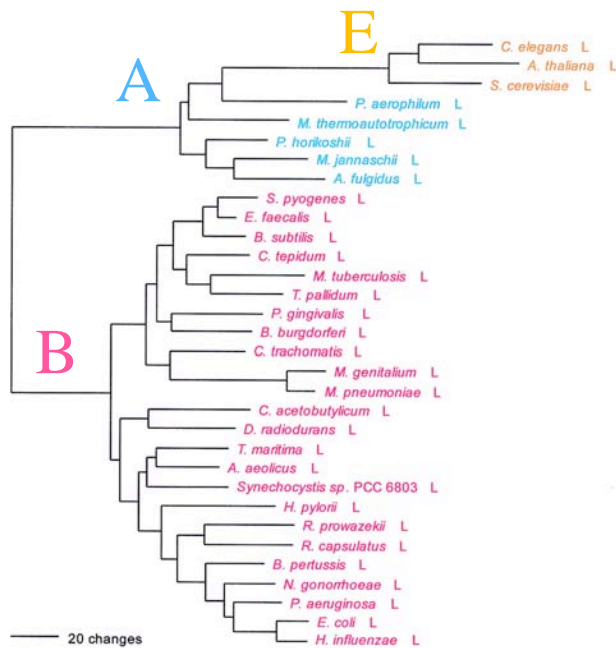
After W. Doolittle, modified by G. Olsen

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

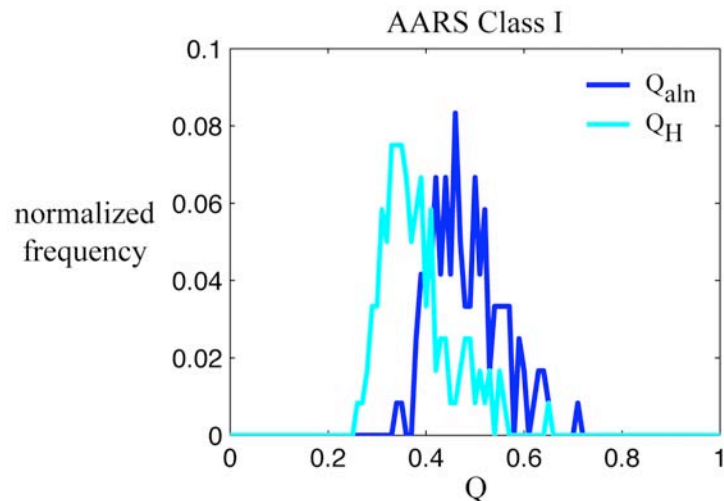
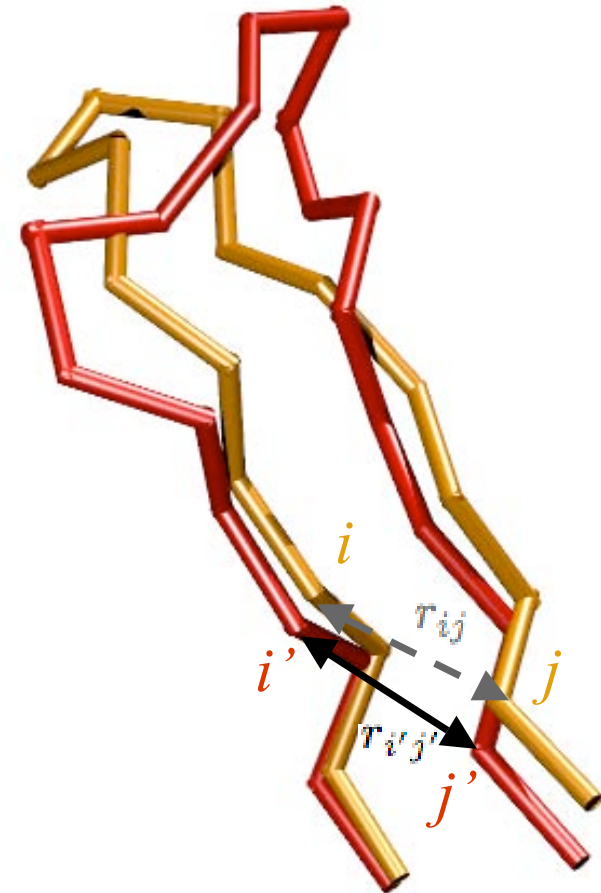
Protein Structure Similarity Measure

Q_H Structural Homology

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

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

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

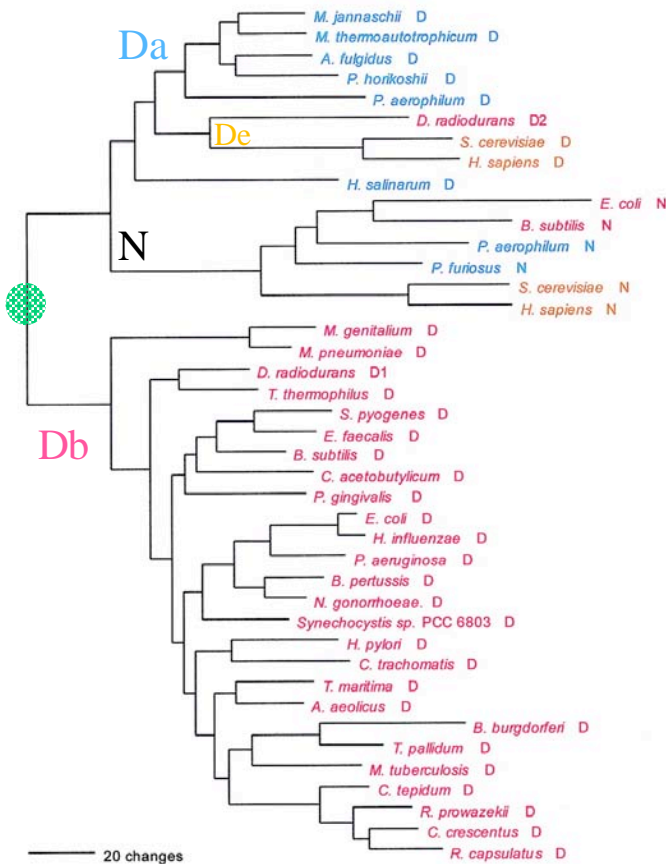


“Gaps should count as a character
but not dominate” C. Woese

Protein structure encodes evolutionary information

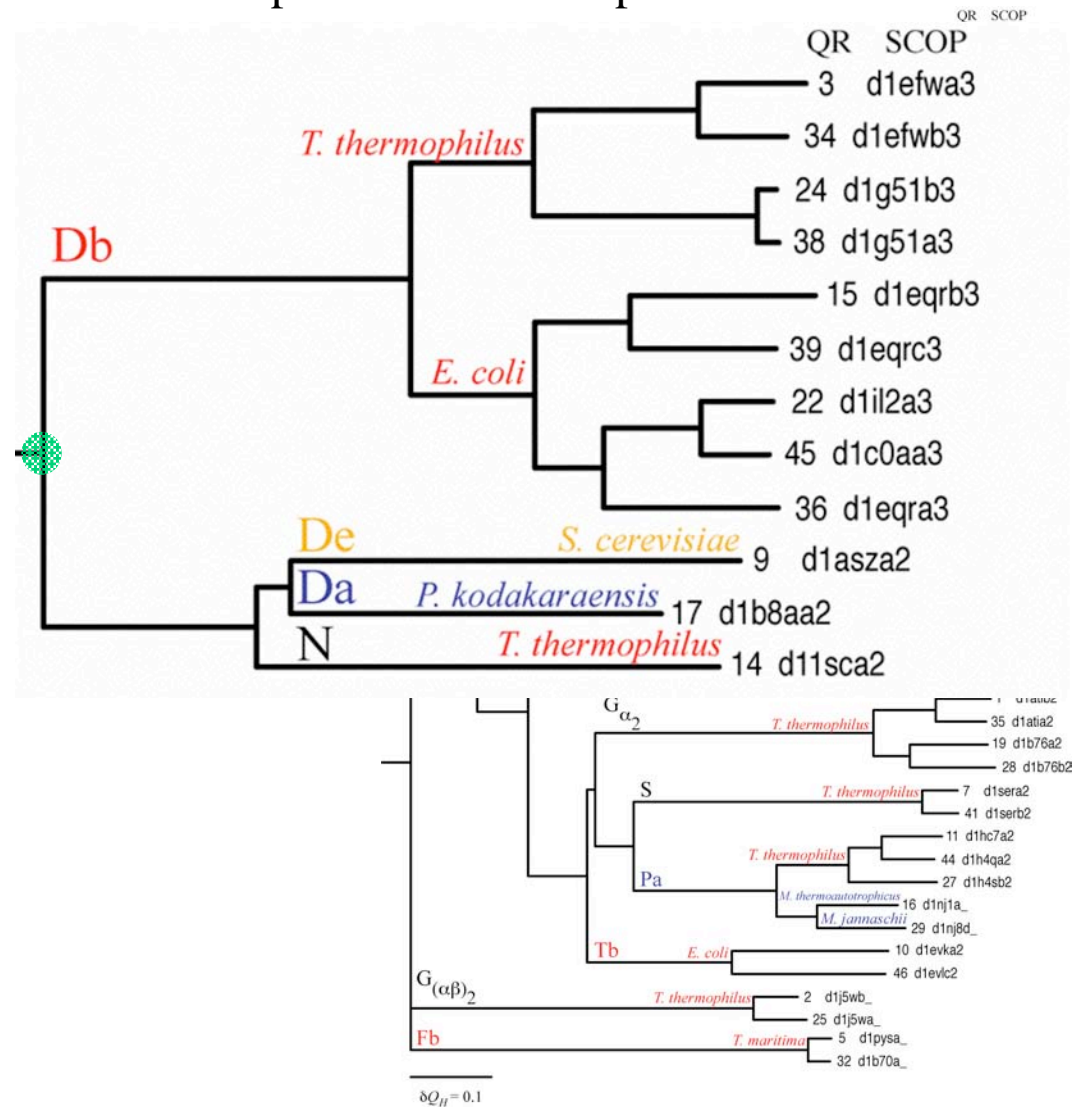
Sequence Phylogeny

AspRS-AsnRS Group



Structure Phylogeny

AspRS-AsnRS Group Class II AARSs



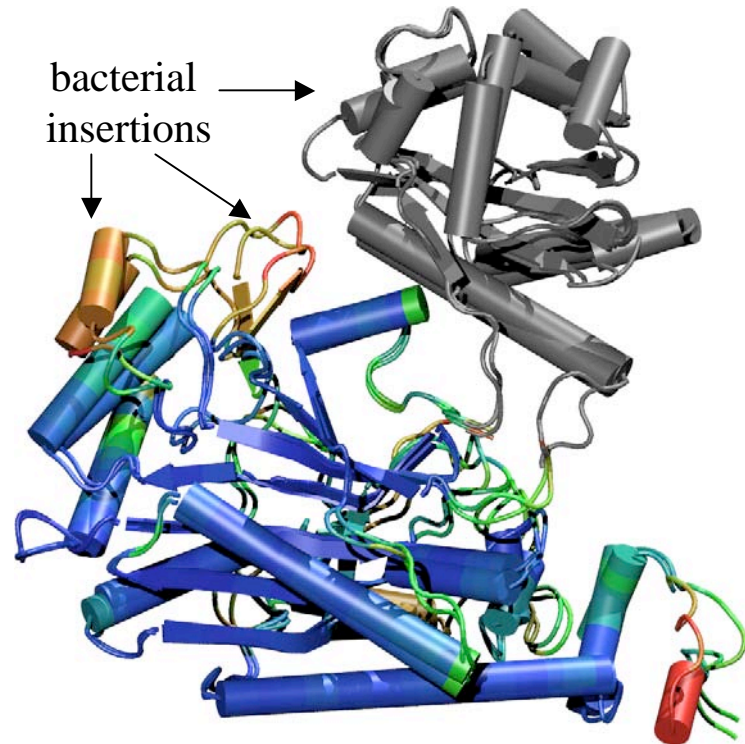
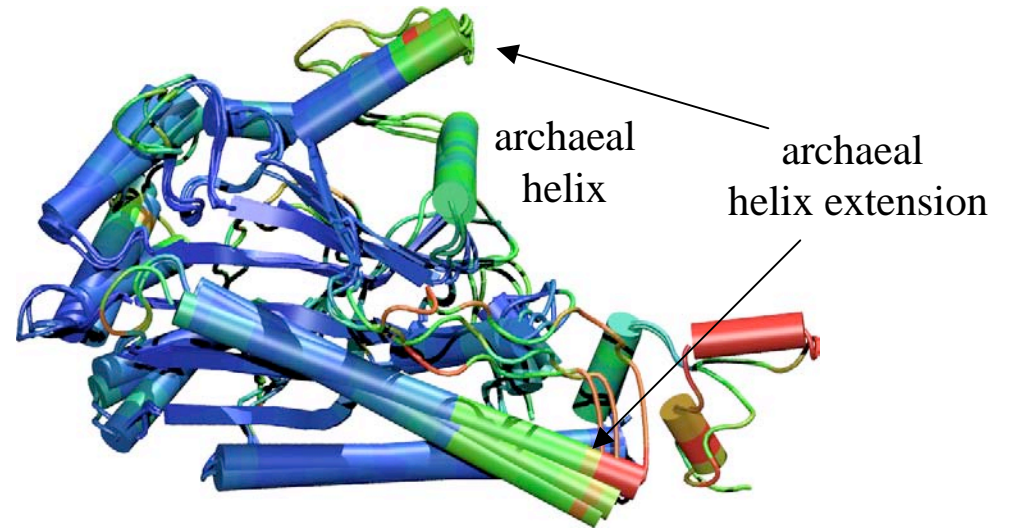
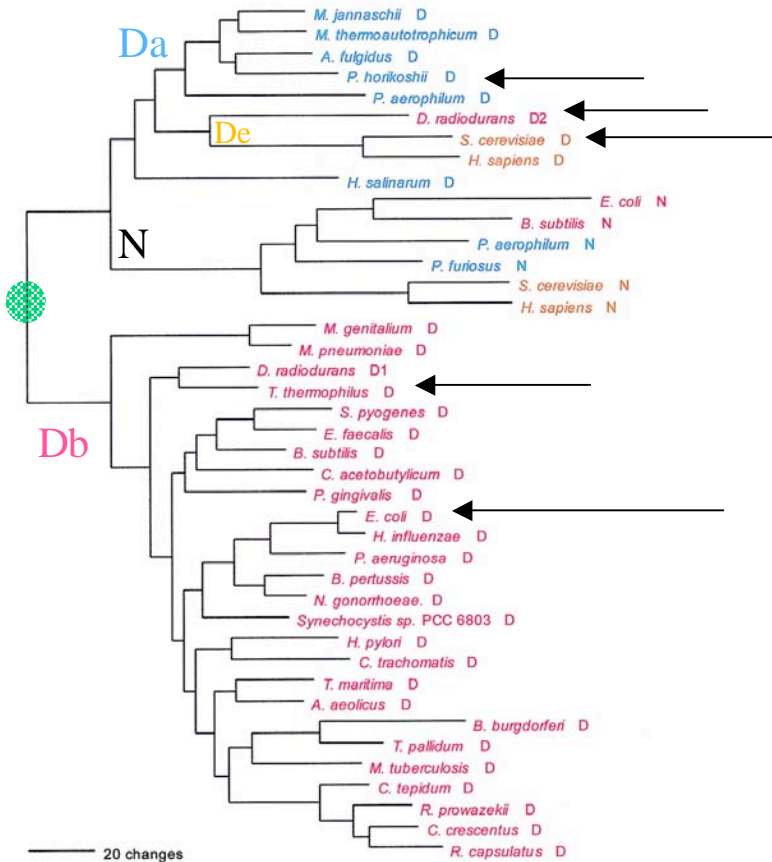
Woese, Olsen, Ibba, Soll *MMBR* 2000

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

Horizontal Gene Transfer in Protein Structure

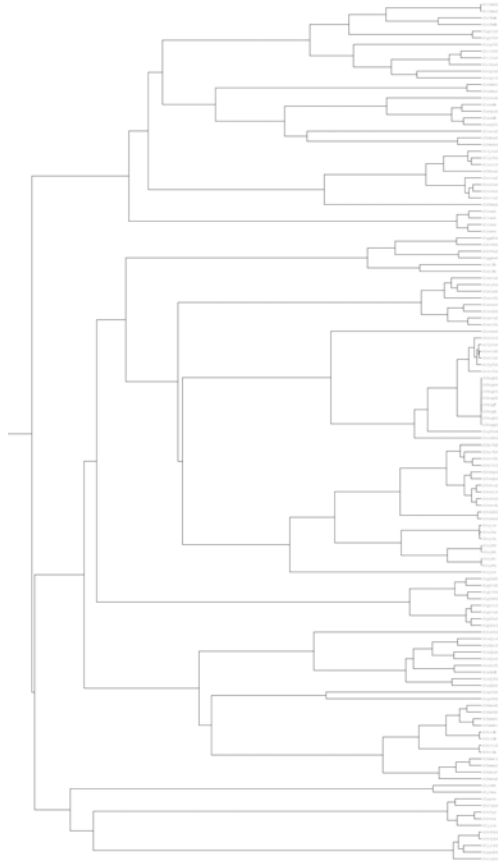
Sequence Phylogeny

AspRS-AsnRS Group



Non-redundant Representative Sets

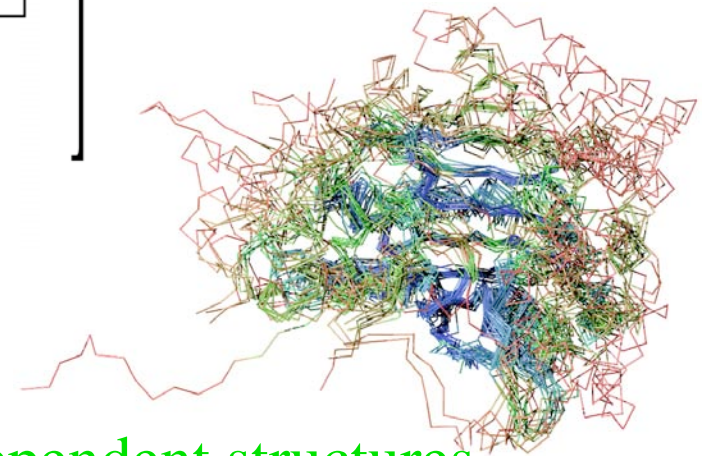
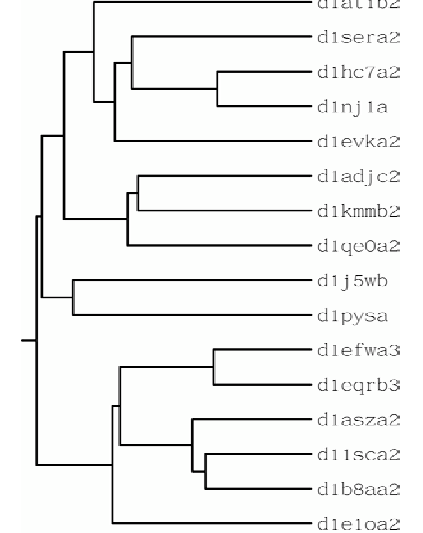
Too much information
129 Structures



Multidimensional QR
factorization
of alignment matrix, A .

$$A = \left[\begin{array}{c} \nearrow d=4 \\ \downarrow l_{aln} \\ \xrightarrow{k_{proteins}} \end{array} \begin{array}{c} G \\ Z \\ Y \\ X \end{array} \right]$$

Economy of information
16 representatives



QR computes a set of maximal linearly independent structures.

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)

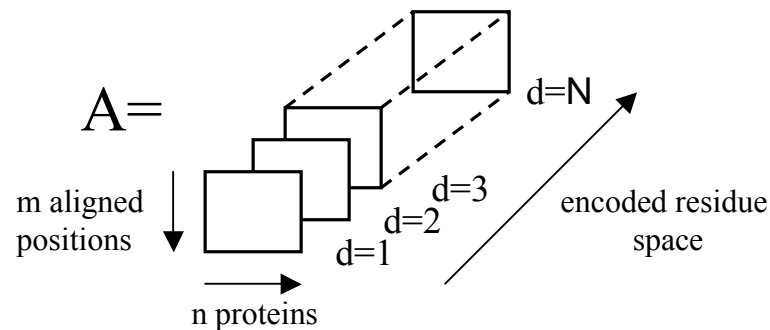
B = (0,1,0)

C = (0,0,1,0)

...

GAP = (0,1)

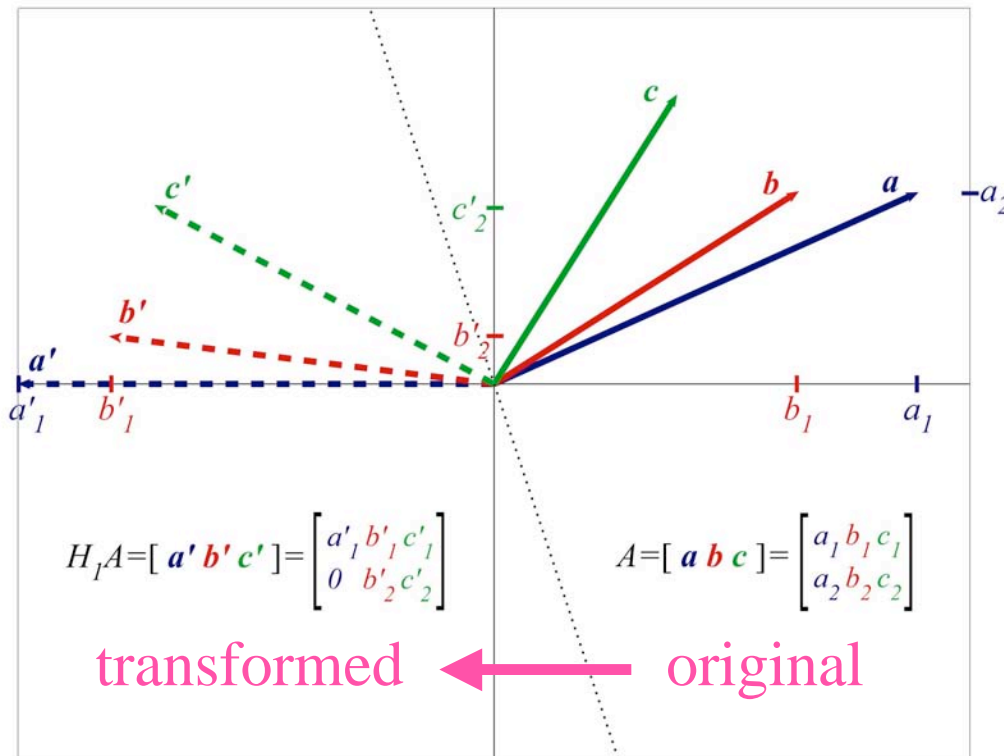
Alignment Matrix



The QR establishes an **order** of linear dependence

by applying Householder transformations and permutations

$$Q^T = H_n \dots H_1$$



Three 1-D (2 residue) proteins **a b c**.

a is our measuring stick, reference frame.

The transformation reveals that **b** is more linearly dependent on **a**, so the permutation swaps **b'** with **c'**.

Given **a**, **c** adds more information to the system than **b**.

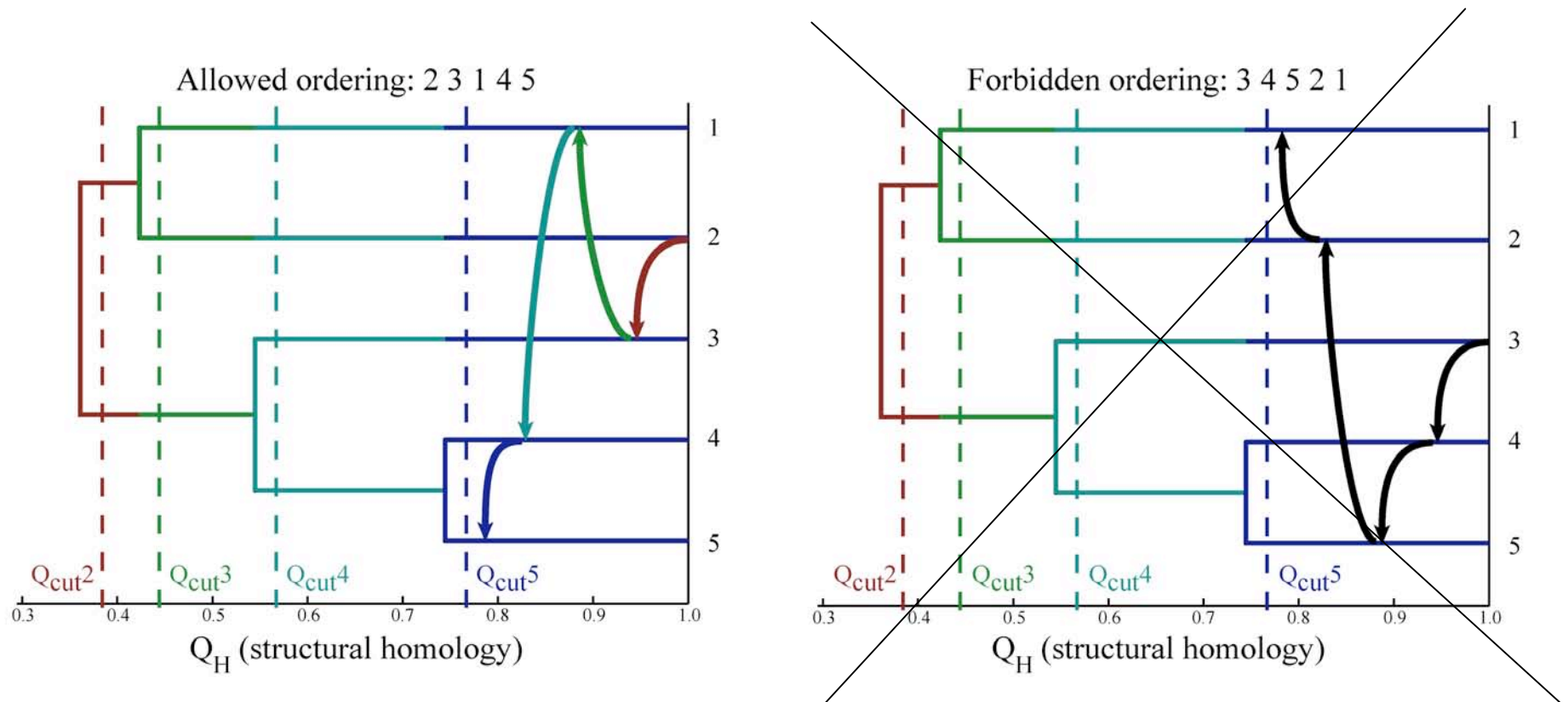
Multiply aligned proteins exist in a higher dimensional space, so this magnitude is computed with a matrix p-norm:

$$\|a_j\|_{F_p} = \left(\sum_{d=1}^4 \sum_{i=k}^{m_{a1n}} |a_{ijd}|^p \right)^{1/p}$$

adjustable
parameter

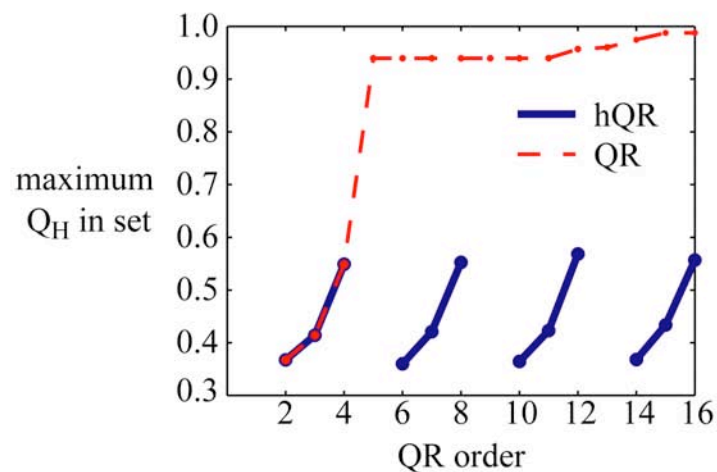
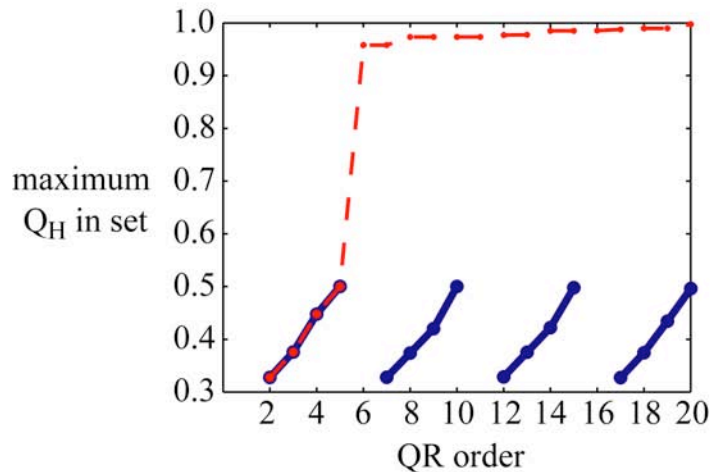
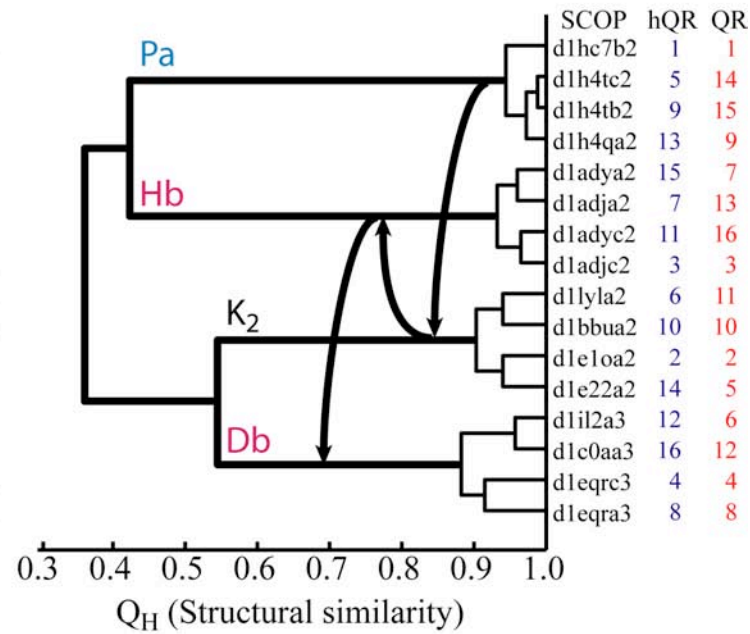
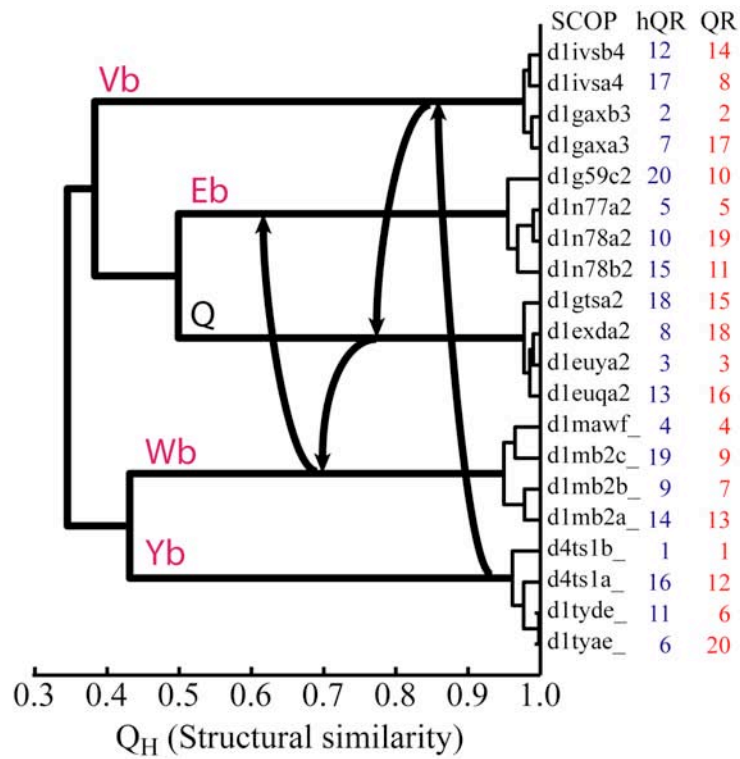
What are the constraints on the parameters?

Must maintain the evolutionary history of the protein group.



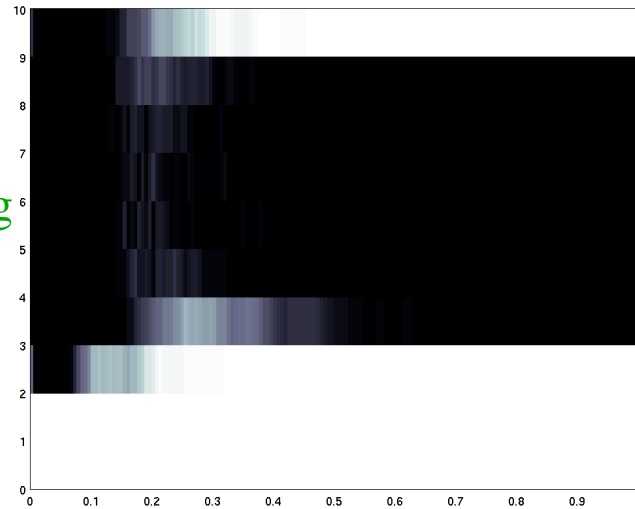
This rule is used to determine the value of two adjustable parameters in our implementation of the QR.

Hierarchical Multidimensional QR -



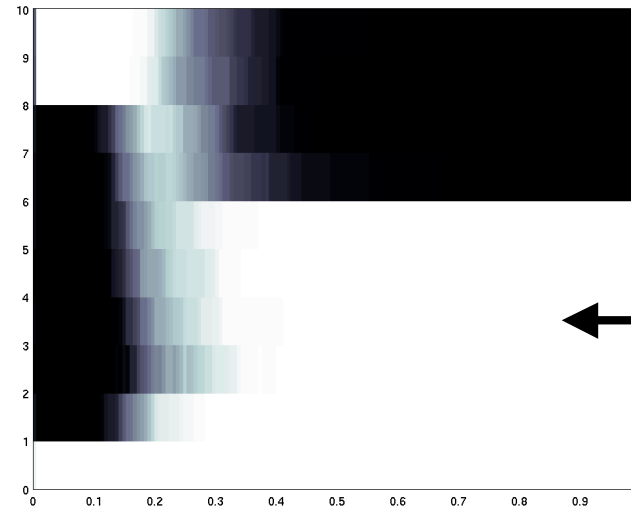
Parameters Define the Measure of Linear Dependence

AARS class I, Rossman fold



γ (normalized)

AARS class II, Novel fold



γ (normalized)

ordering
p-norm

← forbidden

← allowed

ordering norm

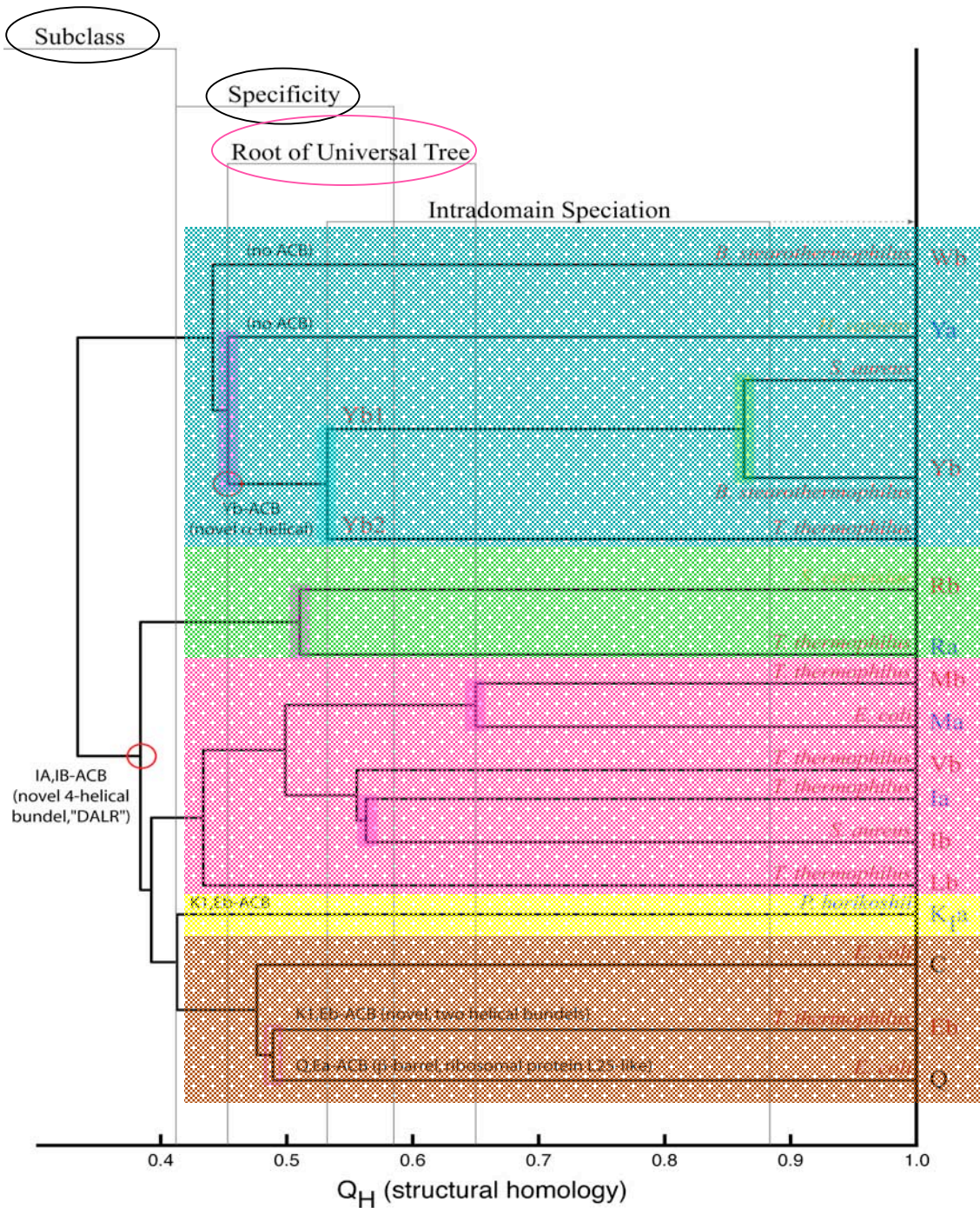
$$\max_{j=k, \dots, n_{\text{proteins}}} (\|a_j\|_{F_p})$$

$$\|a_j\|_{F_p} = \left(\sum_{d=1}^4 \sum_{i=k}^{m_{aln}} |a_{ijd}|^p \right)^{1/p}$$

gap scale

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

Class I AARSs evolutionary events

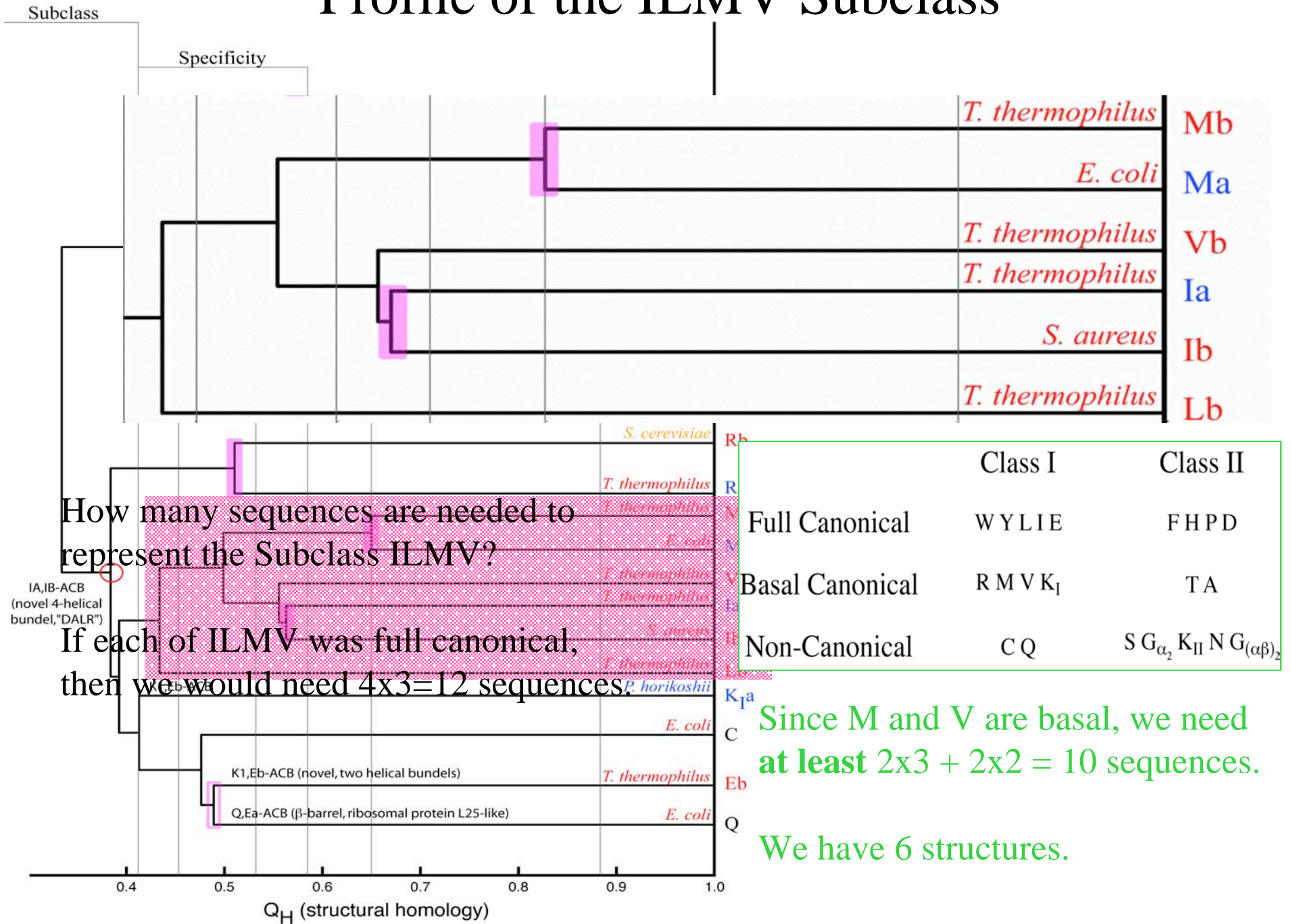


5 Subclasses

Specificity – 11 Amino acids

Domain of life A, B, E

Profile of the ILMV Subclass



How many sequences are needed to represent the Subclass ILMV?

If each of ILMV was full canonical, then we would need $4 \times 3 = 12$ sequences.

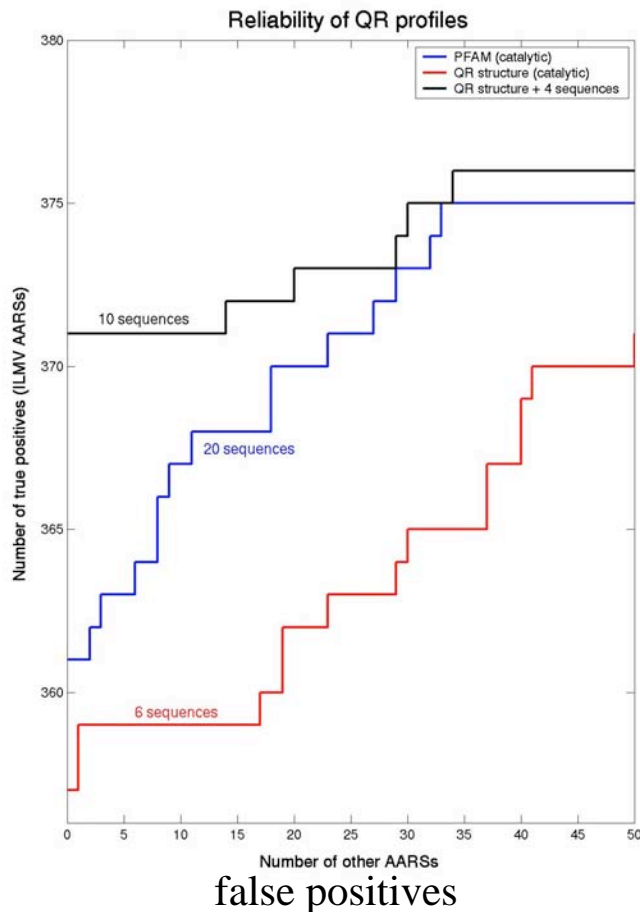
Since M and V are basal, we need at least $2 \times 3 + 2 \times 2 = 10$ sequences.

We have 6 structures.

Non-Redundant Profiles for Database Searching

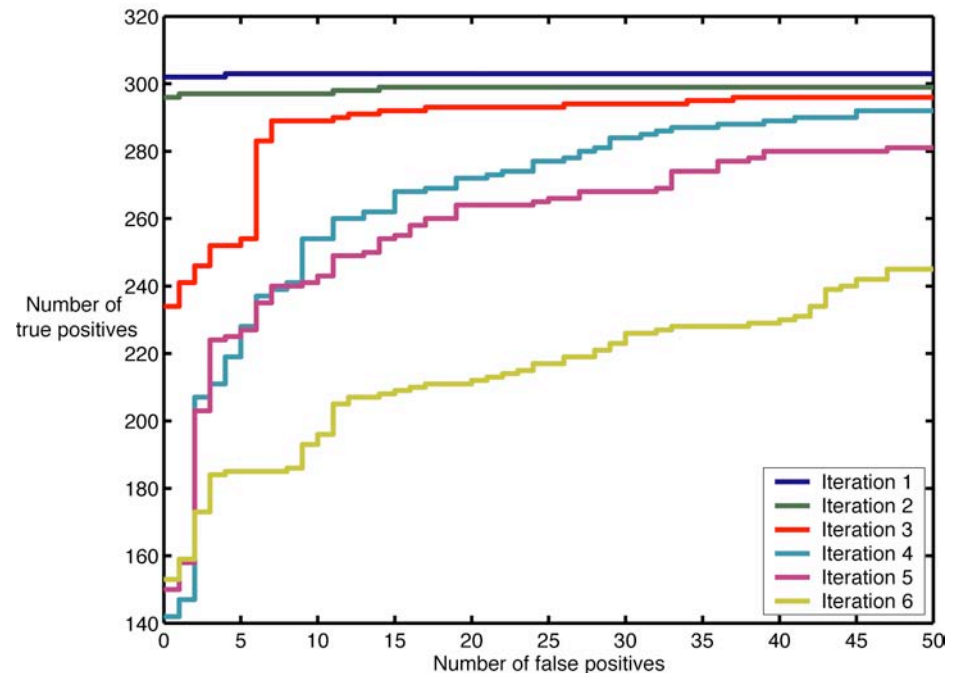
AARS Subclass ILMV

HMMER



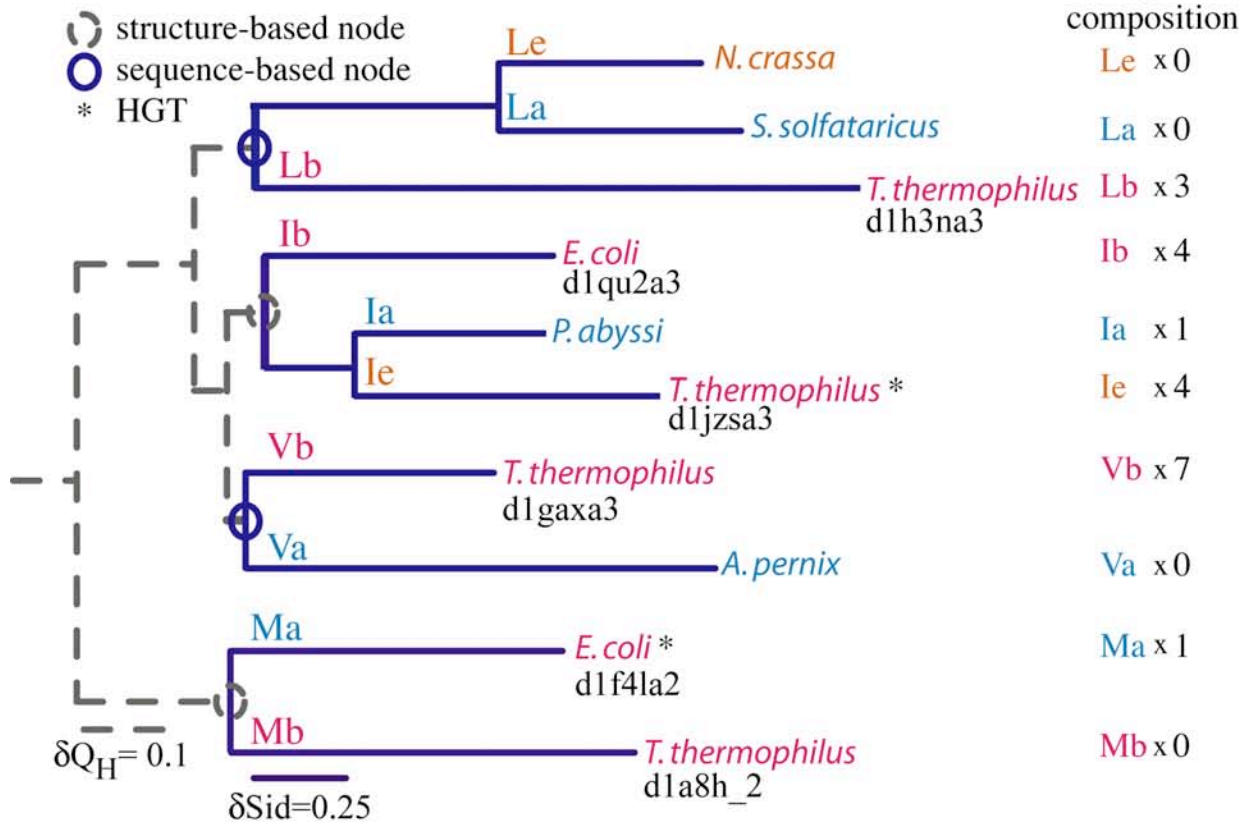
Choosing the right 10 sequence makes all the difference.

Psi-Blast



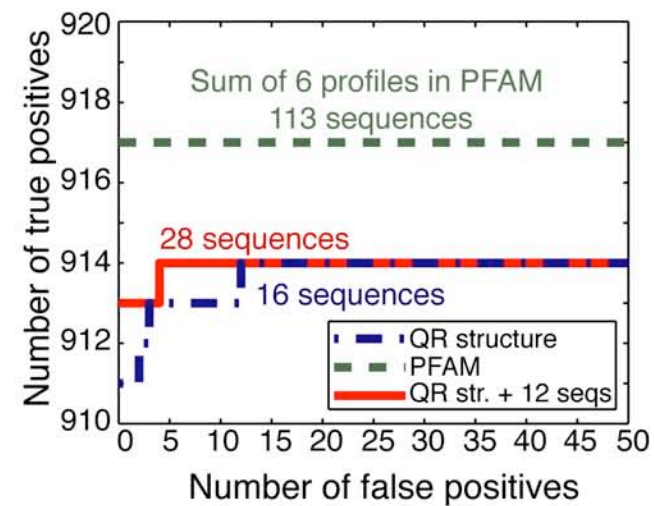
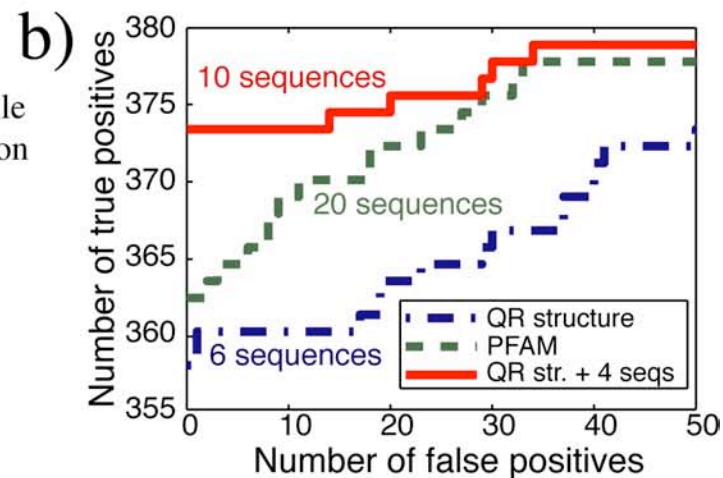
Starting with a non-redundant profile, accuracy diminishes with Psi-blast iterations which add in bias. Repair with QR filter.

a) Combined Structure-Sequence Phylogeny
an evolutionary profile of subclass IA AARs



Pfam profile composition

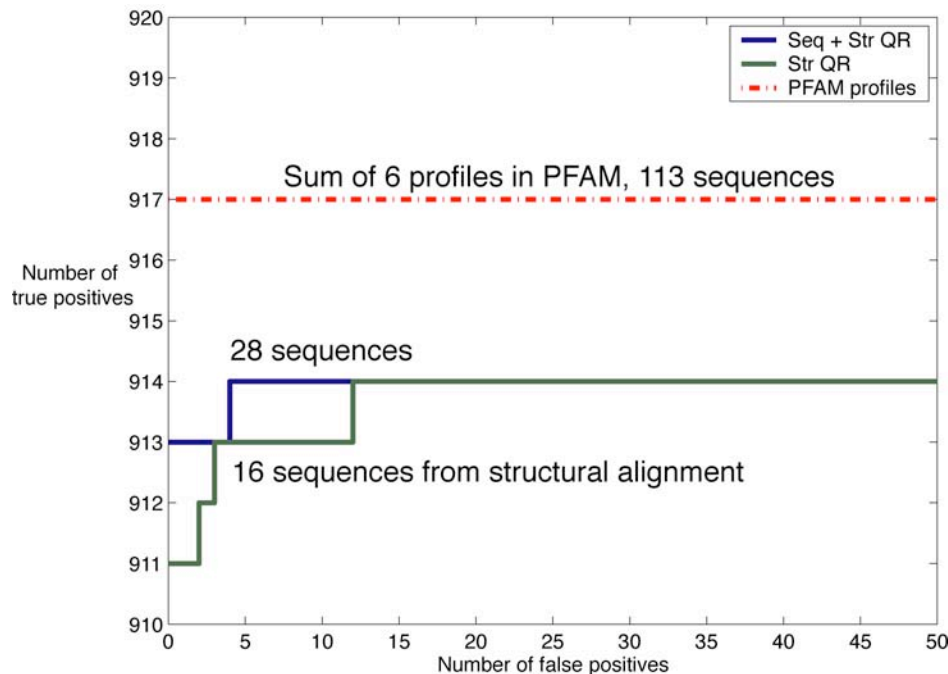
Le x 0
 La x 0
 Lb x 3
 Ib x 4
 Ia x 1
 Ie x 4
 Vb x 7
 Va x 0
 Ma x 1
 Mb x 0



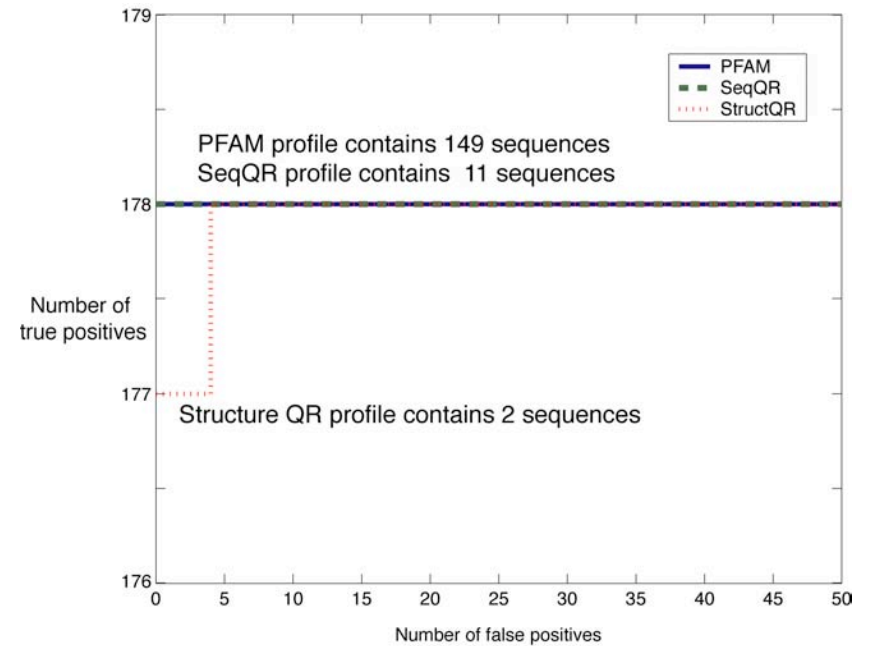
The Economy of Information

How many sequence are needed for profiles?

A single profile
for class I AARSs



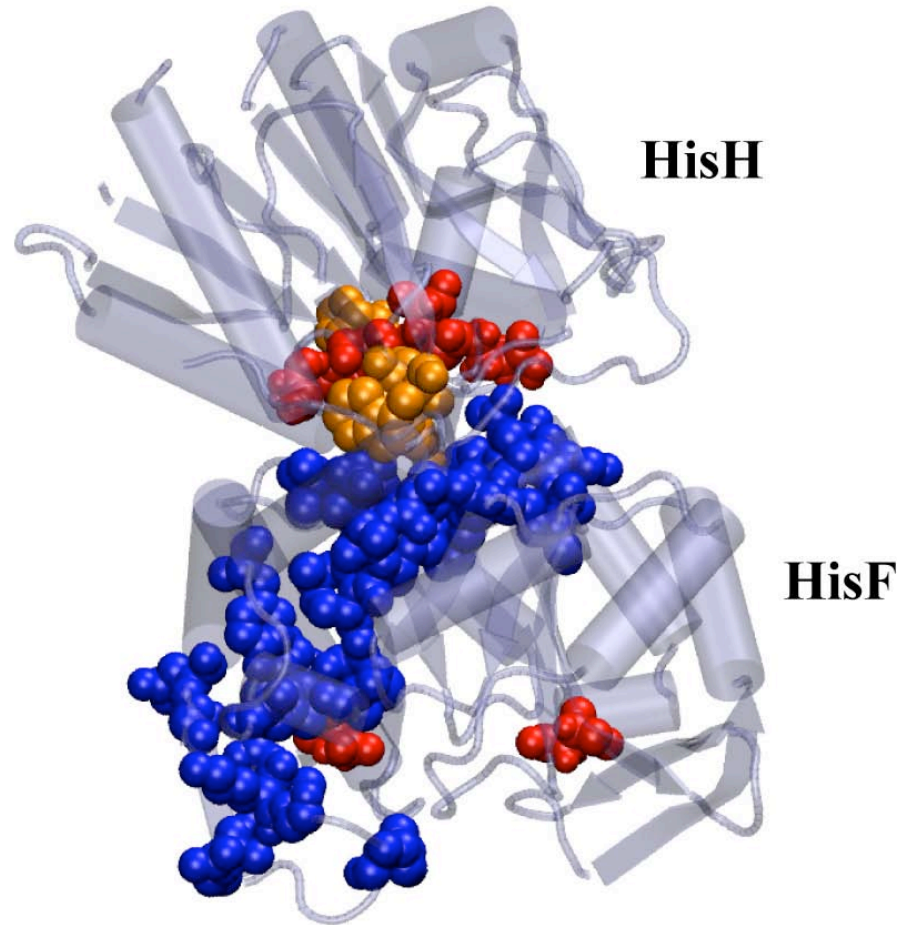
HisA and HisF Protein Family
TIM Barrel fold



PFAM profile of 113 sequences finds 3 additional sequence fragments compared to the non-redundant profile of 28 sequences.

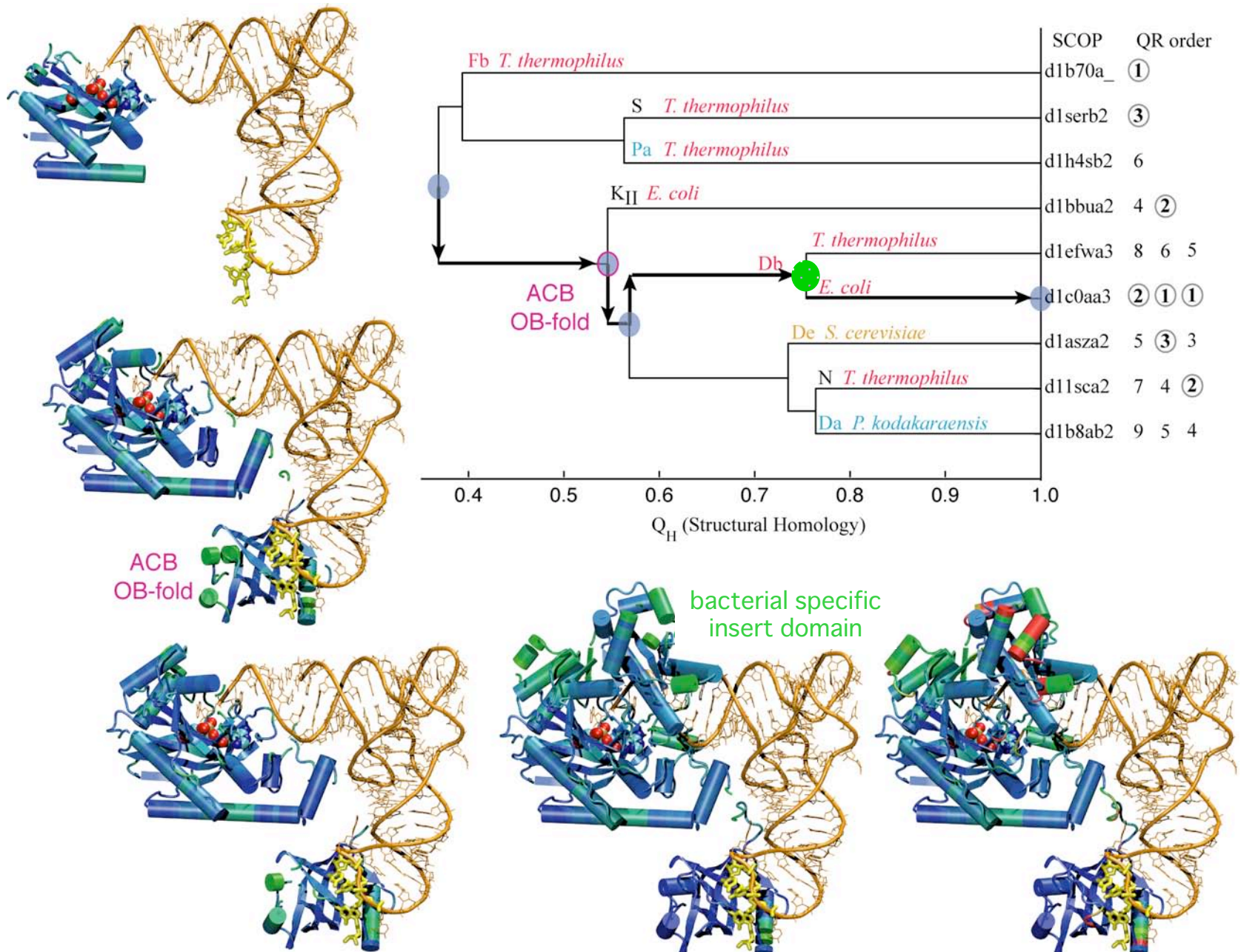
If the sequences well represent the evolutionary history of the protein family, a factor of 10 to 100 less information is required.

Evolutionary Structure/Sequence Profiles Suggest Reaction Pathway



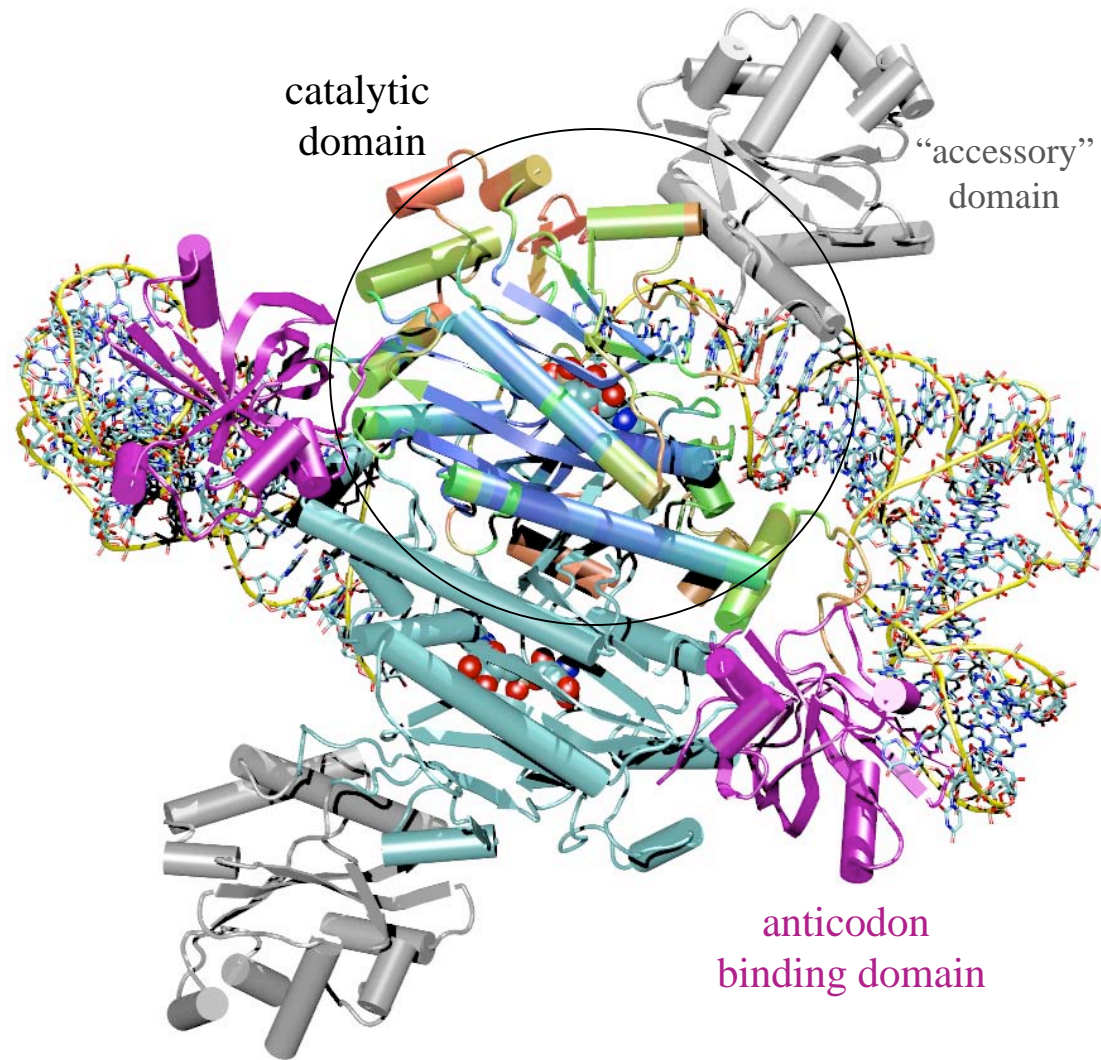
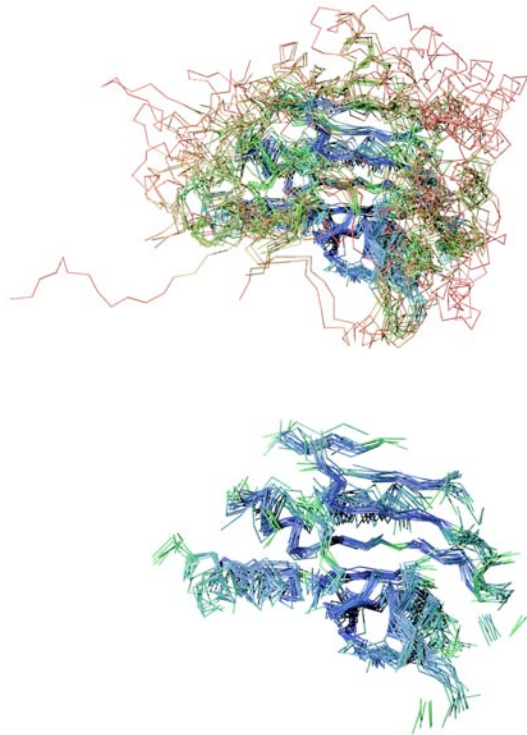
R. Amaro and Z. Schulten, *MD Simulations of Substrate Channeling*, Chemical Physics Special Issue, 2004 (in press). *FE Landscapes of Ammonia Channeling*, PNAS 2003

Evolution of Structure and Function in AspRS



AARS domains have different Evolutionary Histories

catalytic domains
AARSs II



bacterial type aspartyl-tRNA synthetase
E. coli, homodimer

Summary

Evolutionary information is encoded in protein structure.

Protein structure can be used to investigate early evolutionary events.

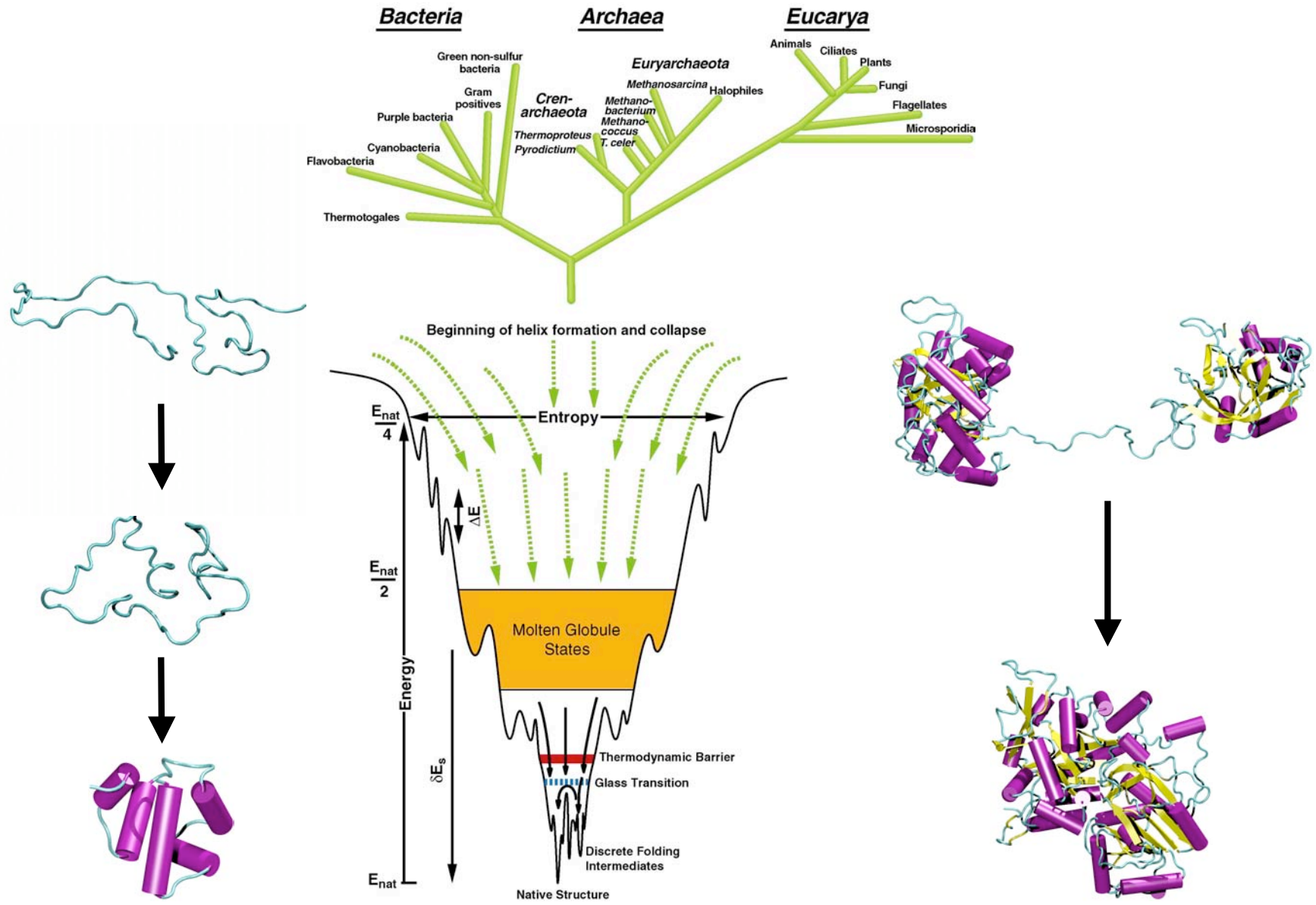
Accounting for gaps is important for comparing homologous structures - structure metric

Multidimensional QR factorization computes non-redundant sets from multiple sequence or structure alignments which well represent the evolutionary history of the group as expressed in phylogenetic tree

Structure databases are limited, but multiple structural alignments provide accurate alignments, especially in the case of distant homologies

Supplement the structures with an appropriate number and type of sequences (in accord with the phylogenetic topology) to produce minimal representative profiles. Search profiles for foldons!!

Evolution of Protein Structure



VMD Multiple Sequence Display with Evolution Analysis Algorithms

The screenshot displays the VMD 1.8.3a2 OpenGL Display interface. The main window shows a 3D protein structure with various colored helices and loops. An 'Extensions' menu is open, listing several plugins: sequence, autoimd, apbsrun, imd, contactmap, pdbtool, ramaplot, rmsd, solvate, timeline, multiseq, tkcon, and vmdmovie. A 'treeWindow' is open, showing a phylogenetic tree with the following entries: d1efwa3.ent Thermus thermophilus B, d1c0aa3.ent Escherichia coli B, d1n9wb1.ent d1n9wb1.ent, d1asza2.ent Saccharomyces cerevisiae E, and d1b8aa2.ent Pyrococcus kodakaraensis A. A scale bar indicates a distance of 0.56. Below the tree, a 'Sequence Display' window shows a multiple sequence alignment of the protein sequences, with specific residues highlighted in yellow.

Extensions

- sequence
- autoimd
- apbsrun
- imd
- contactmap
- pdbtool
- ramaplot
- rmsd
- solvate
- timeline
- multiseq
- tkcon
- vmdmovie

treeWindow

Tree

d1efwa3.ent Thermus thermophilus B
d1c0aa3.ent Escherichia coli B
d1n9wb1.ent d1n9wb1.ent
d1asza2.ent Saccharomyces cerevisiae E
d1b8aa2.ent Pyrococcus kodakaraensis A

0.56

Sequence Display

```
d1b8aa2.ent  IDTEGERLLGKYM--MENENAPLYFLYQYPS-----EAKPFYIMKYDN-----K--PEICRAFDFLEYRGV  
d1asza2.ent  DLSTENEKFLGKLV--RDKYDTDFYILDKFPL-----EIRPFYTMPDPA-----N--PKYSNSYDFFMERGE  
d1n9wb1.ent  DLSEEAERLLGEYA--KERWGSDFVTRYP-----SVRPFYTYP--EE-----DGTTRSFDLLFRGL  
d1c0aa3.ent  ---GSD-KP-DLRDE---SKWAPLWVIDFPMFE--DDGEGGLTAMHHPPTSPPK--DMTAAELKAAPENAVANAYDMVINGY  
d1efwa3.ent  ---GSD-KP-DL-RR---EGFRFLWVDFPLLEWDEEEEAWTYMHHPPTSPPHED--LPLLEKDPGRVRALAYDLVINGV
```


Acknowledgements

Patrick O'Donoghue

Rommie Amaro

Anurag Sethi

John Eargle

Corey Hardin

Michael Baym

Michael Januszyk

Felix Autenrieth

Taras Pogorelov

Brijeet Dhaliwal

Funding: NSF, NIH, NIH Resource for Macromolecular Modeling and Bioinformatics, NRAC NSF Supercomputer Centers

Graphics Programmers VMD

John Stone, Dan Wright, John Eargle

<http://www.ks.uiuc.edu/Research/vmd/alpha/zs04/>

Collaborators

Evolutionary Studies

Gary Olsen, Carl Woese (UIUC)

Algorithms

Mike Heath (UIUC)

Rob Russell (EMBL) **STAMP**

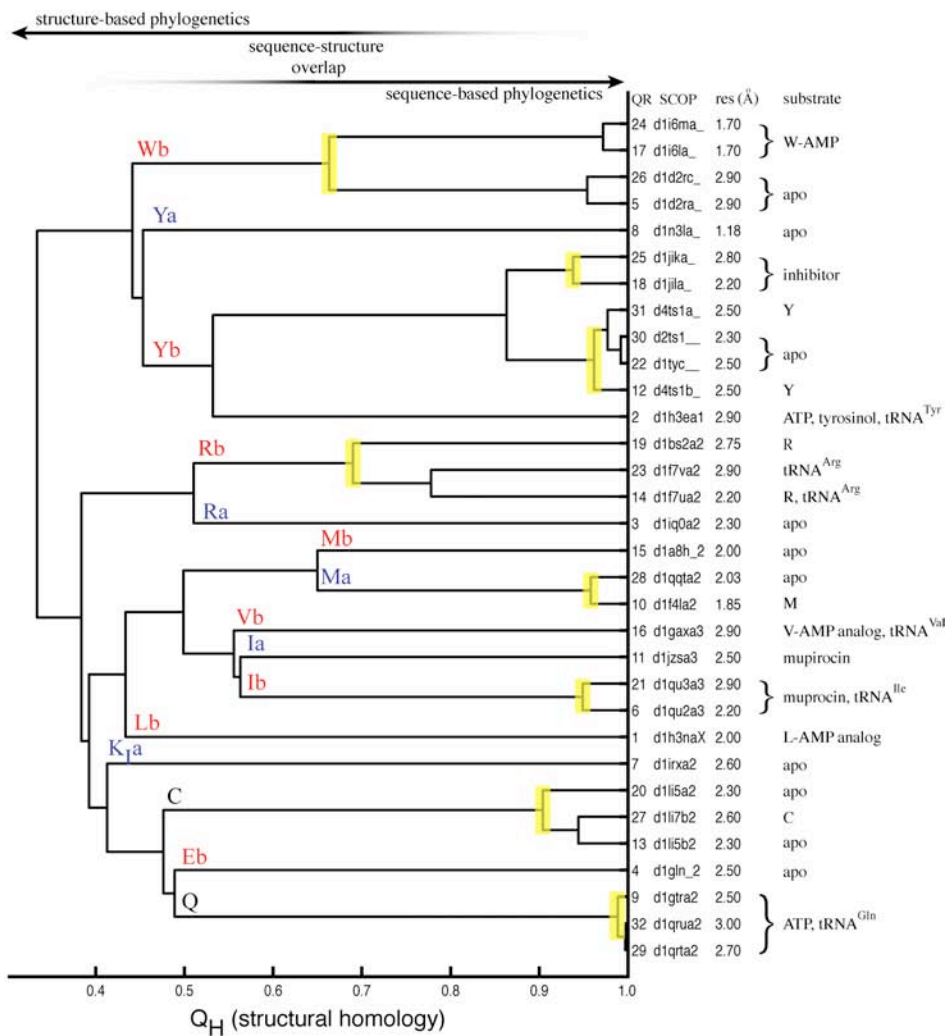
Protein Structure Prediction

Peter Wolynes, Jose Onuchic,

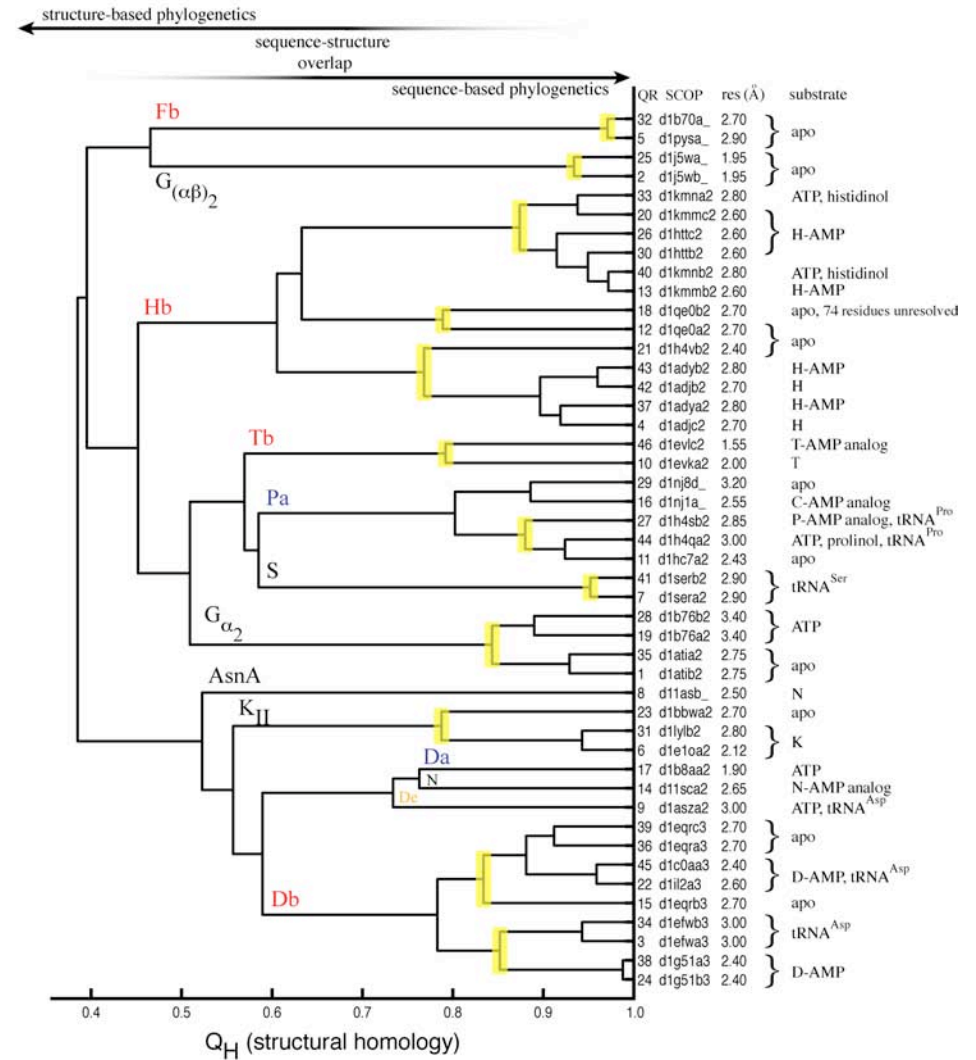
Ken Suslick

extra slides

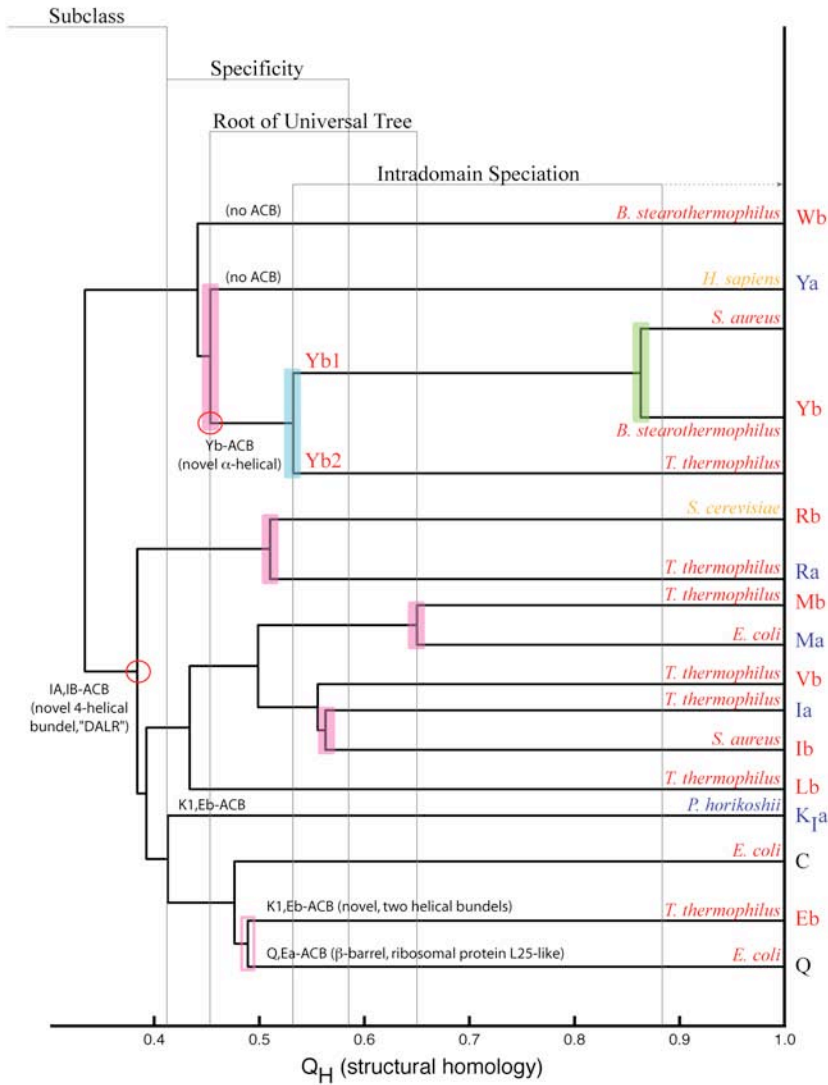
Structure Phylogeny Class I AARSs



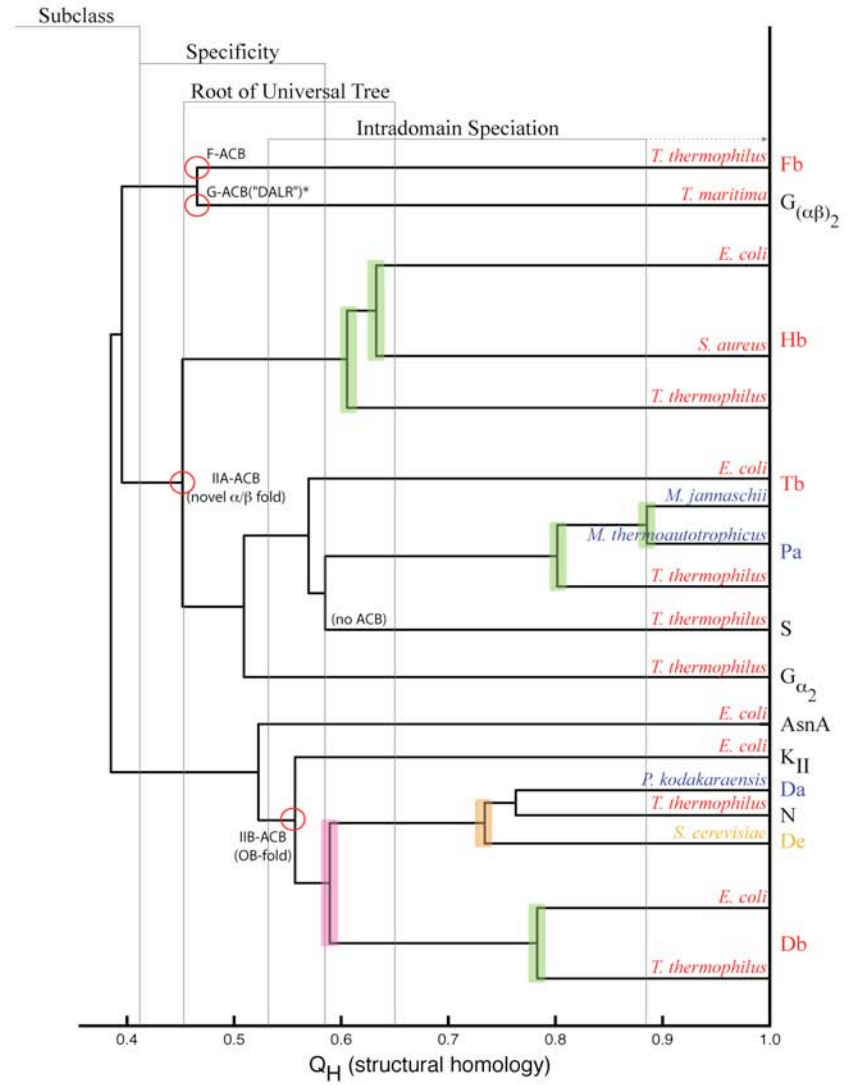
Structure Phylogeny Class II AARSs



Structure Phylogeny Class I AARSs

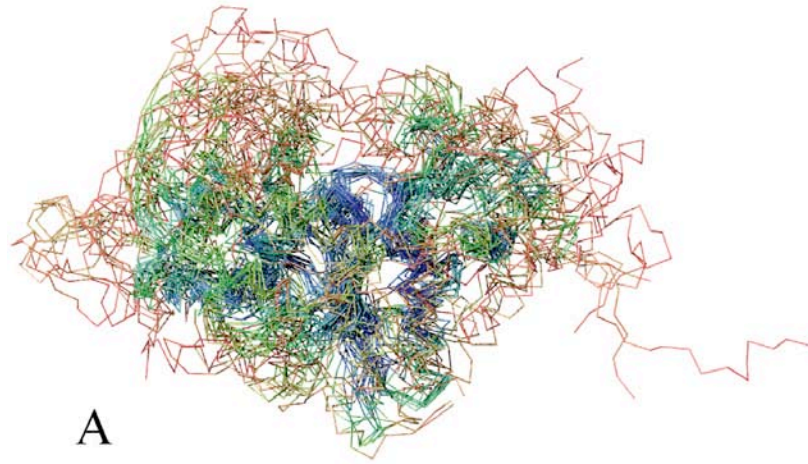


Structure Phylogeny Class II AARSs

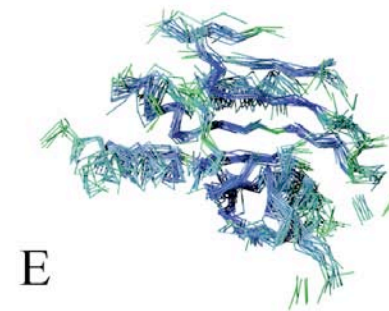
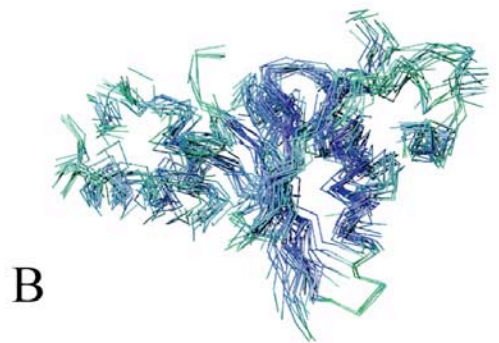
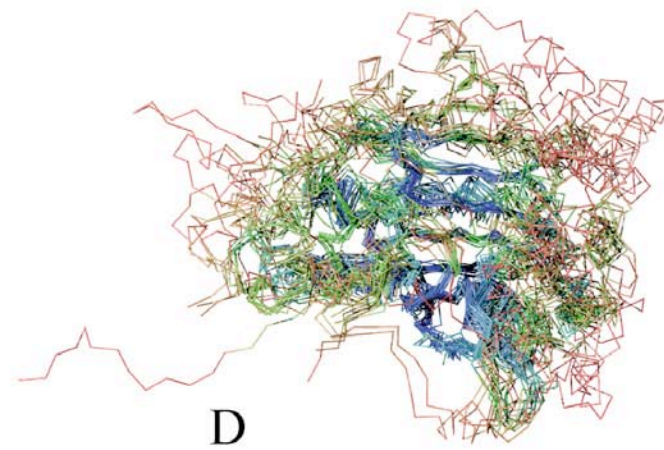


Structural Overlap of the AARSs

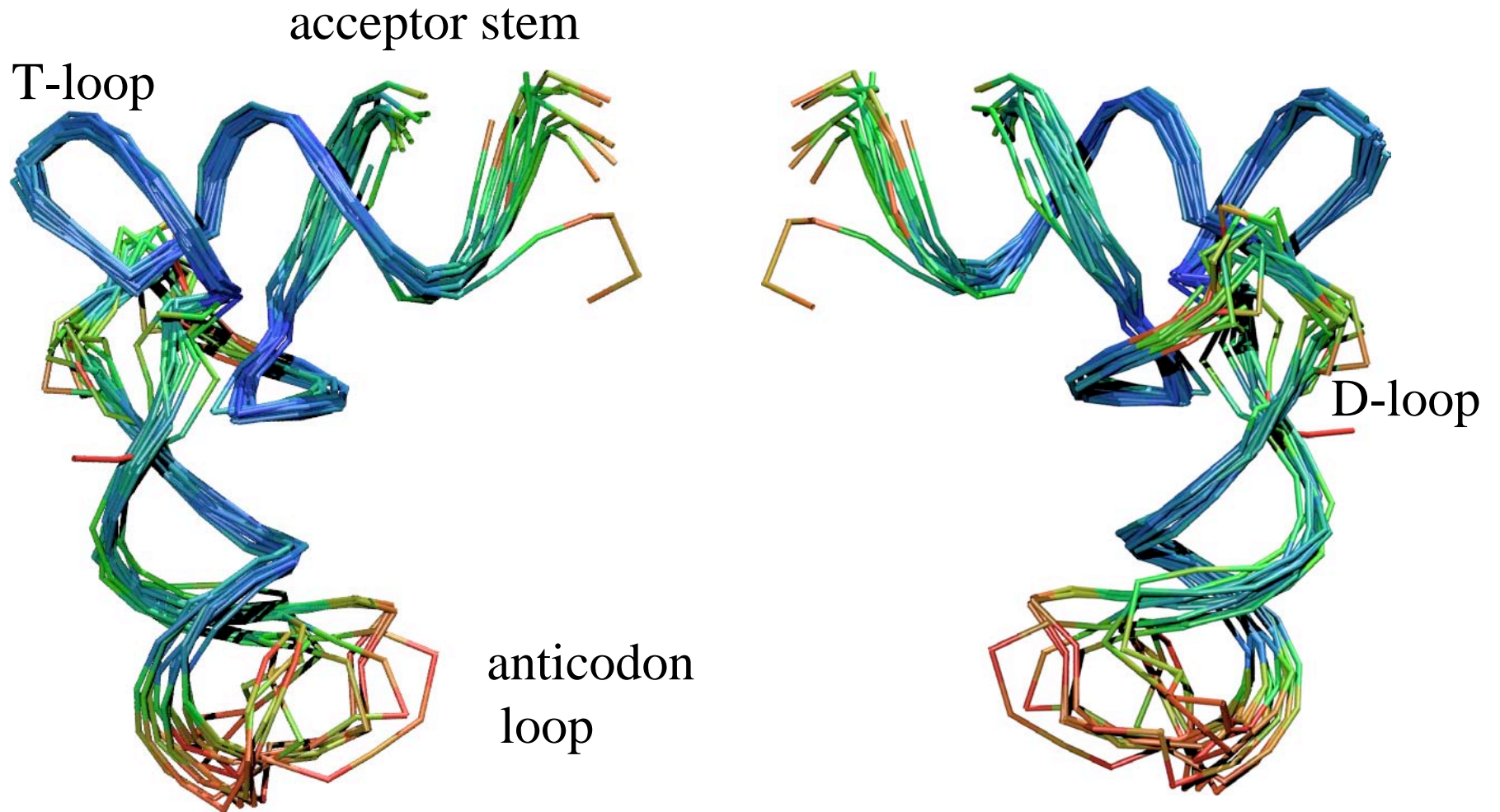
Class I



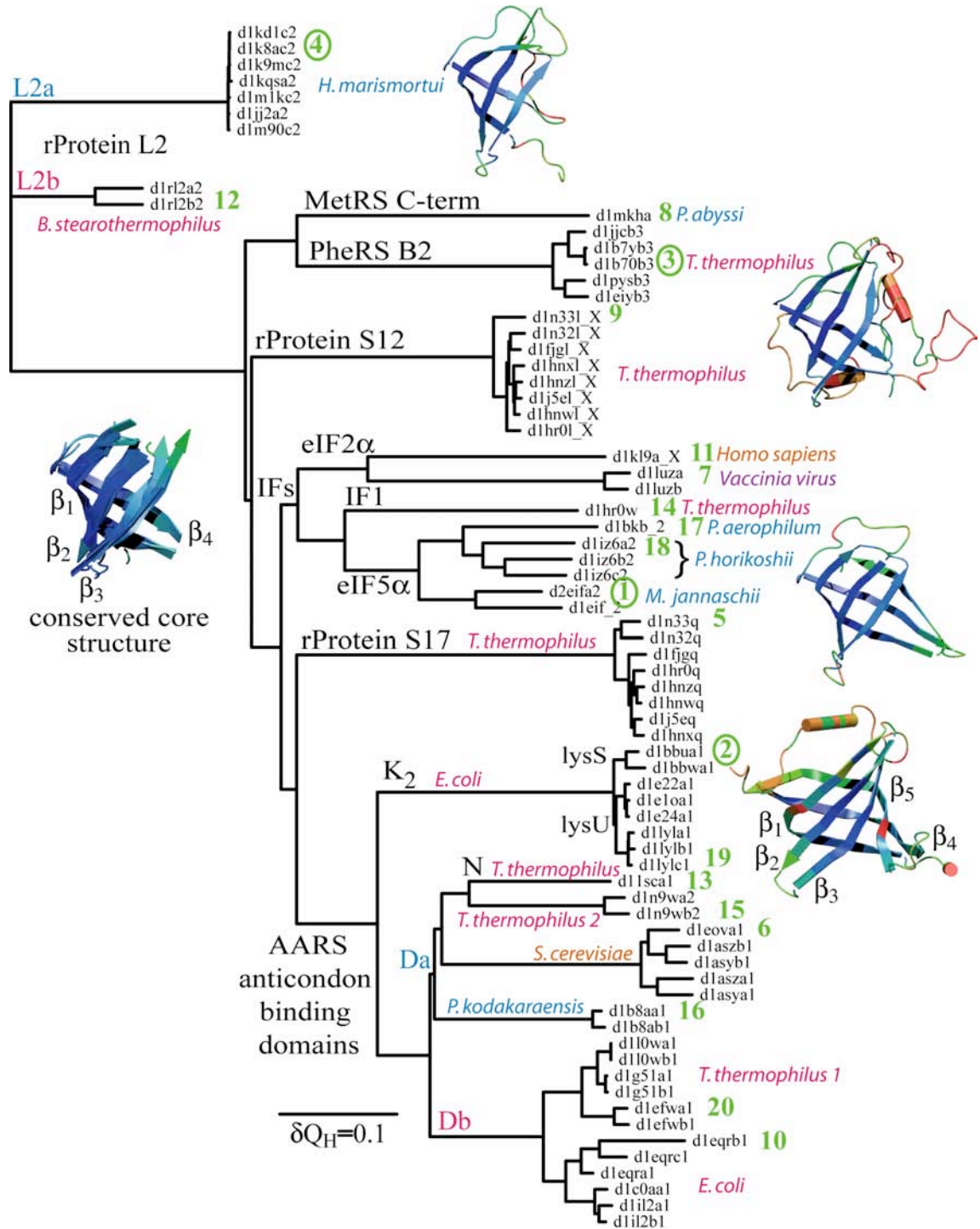
Class II

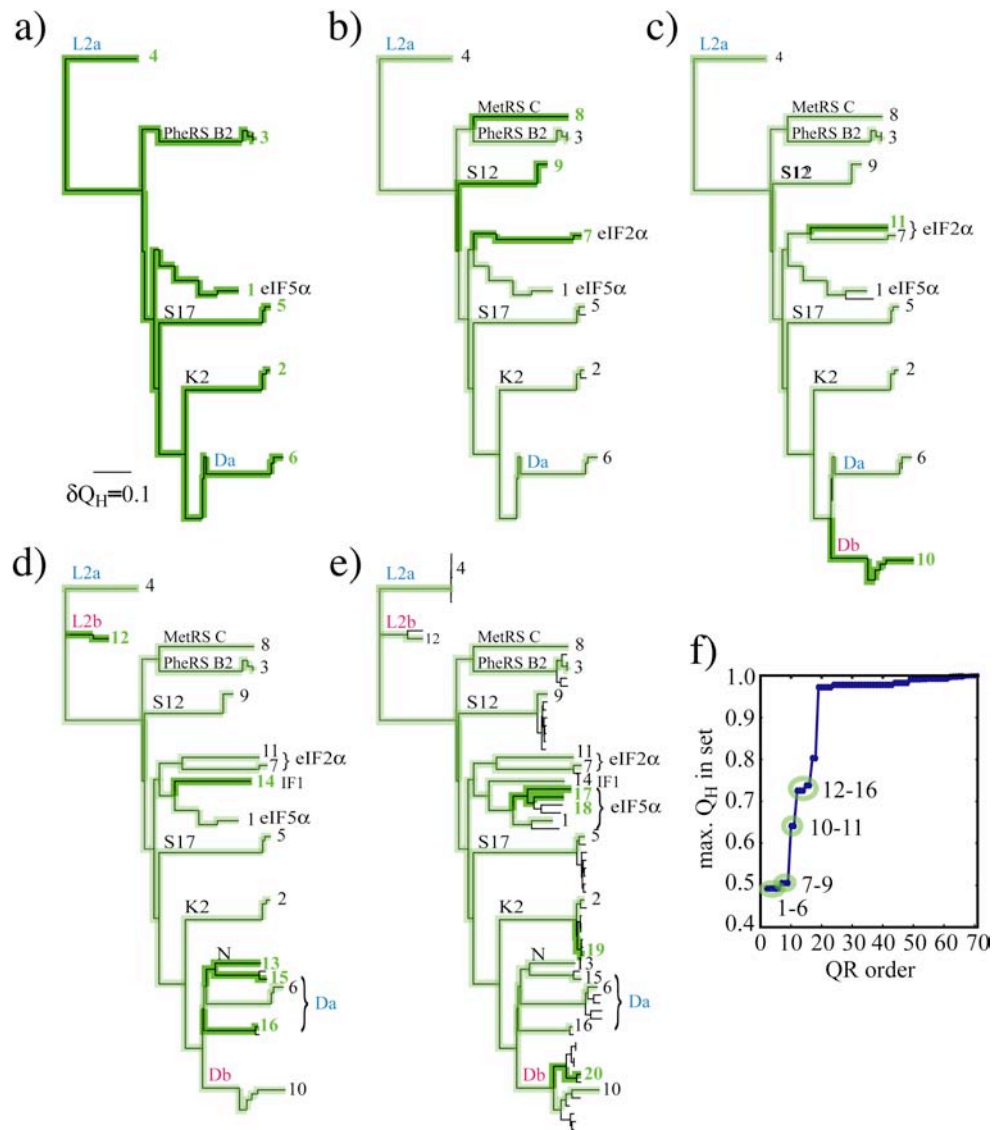


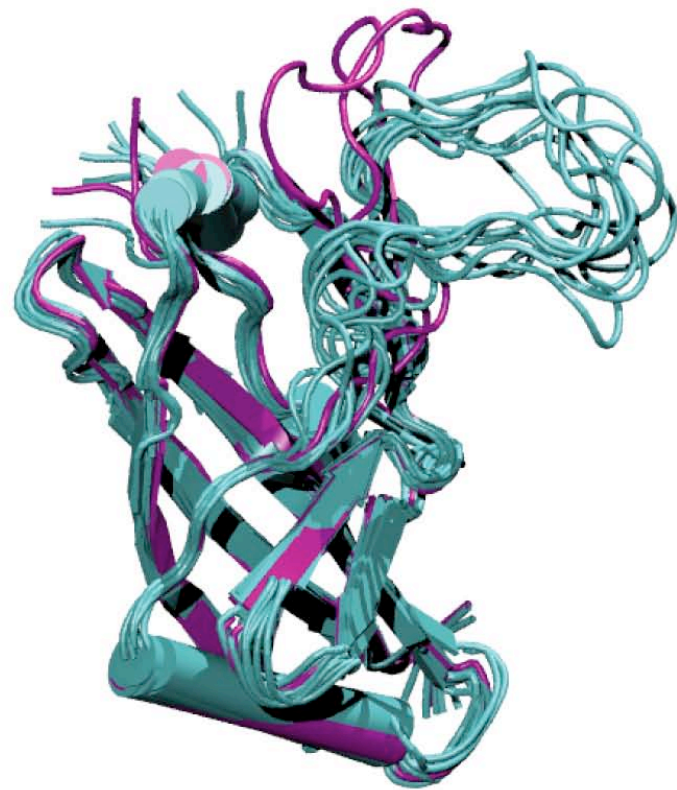
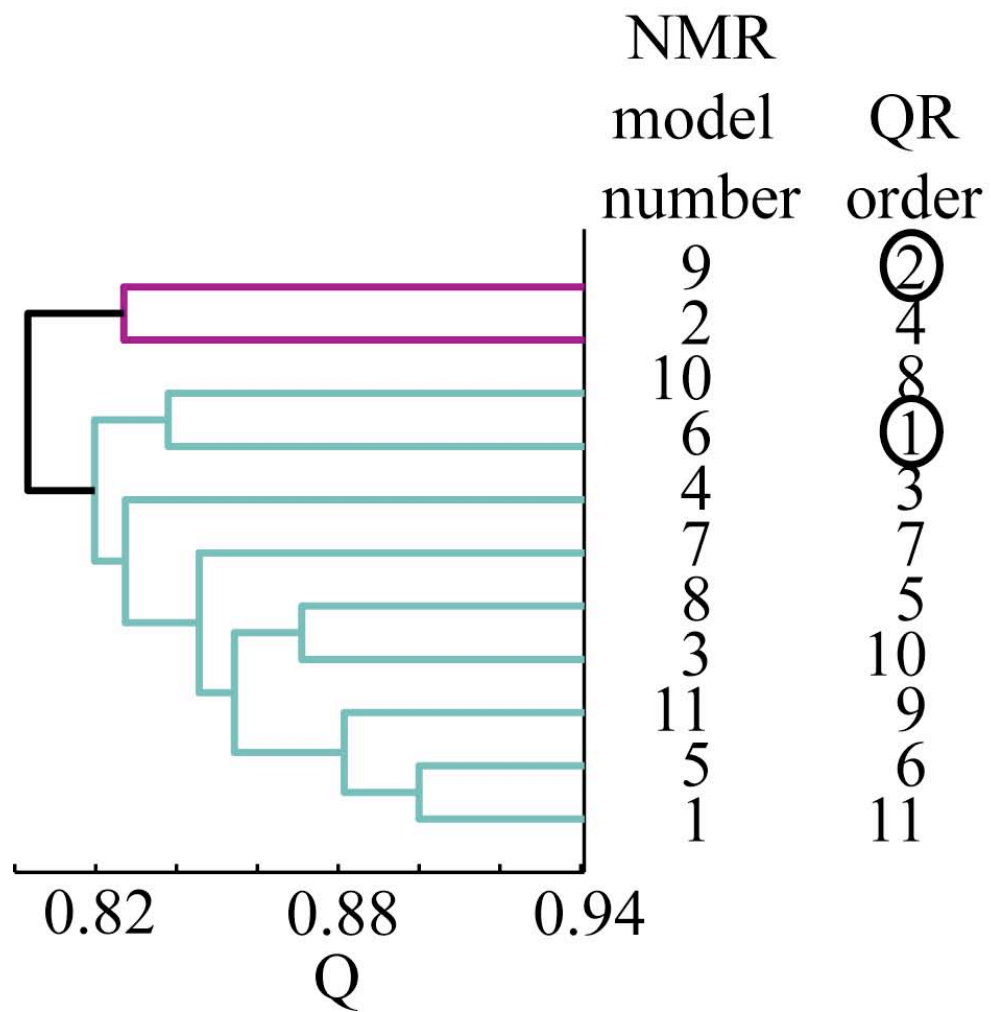
Structural Conservation in tRNA



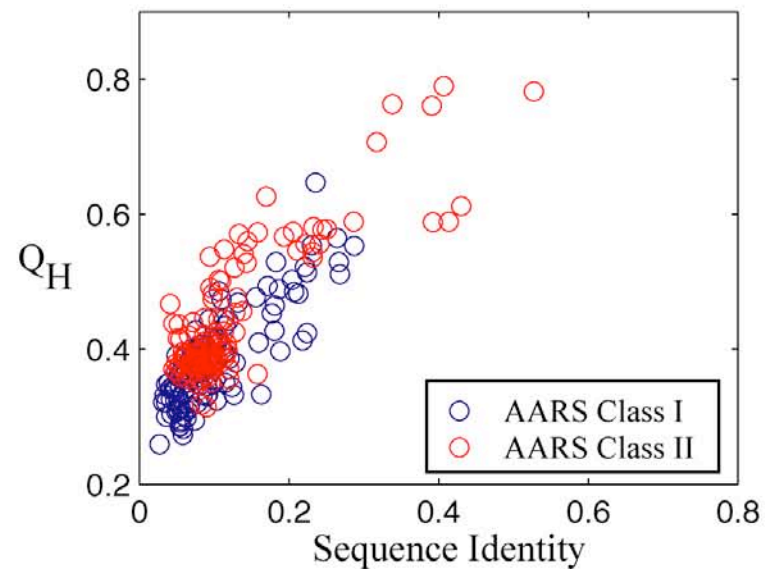
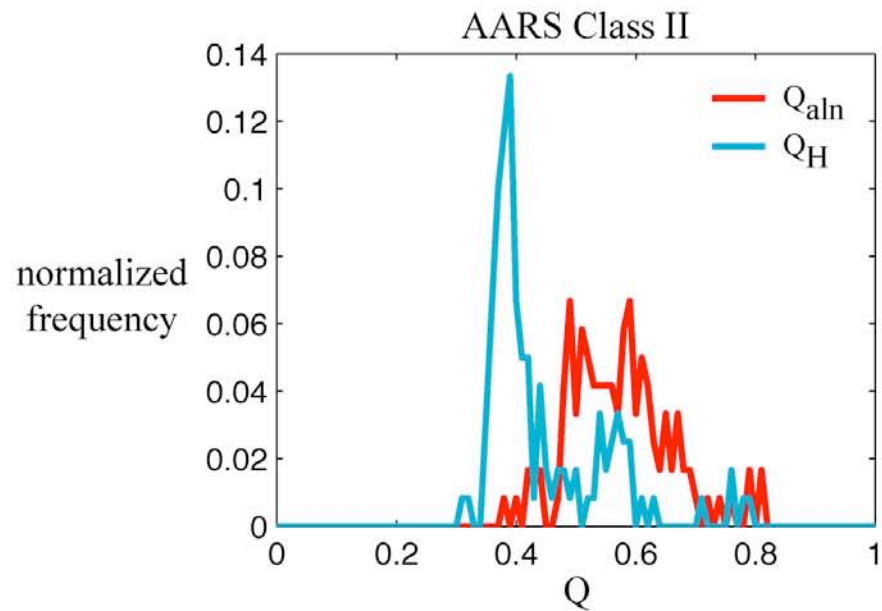
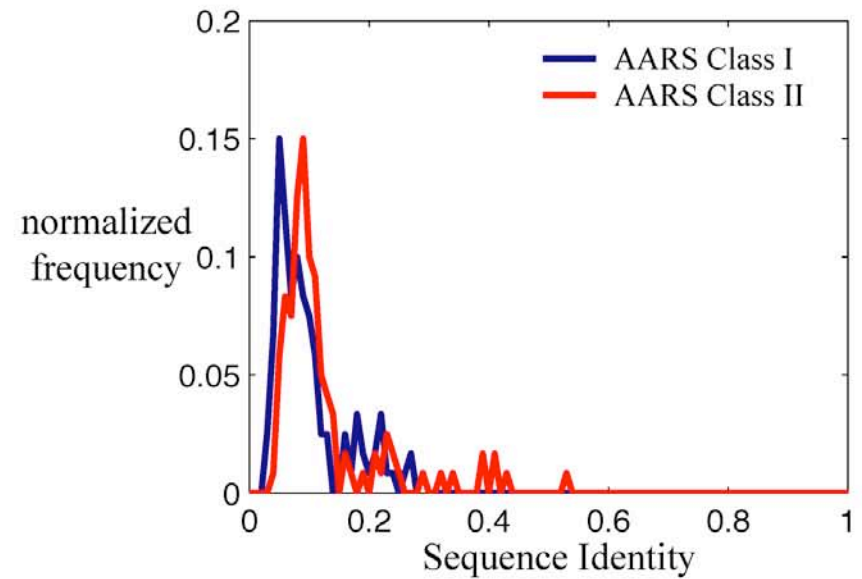
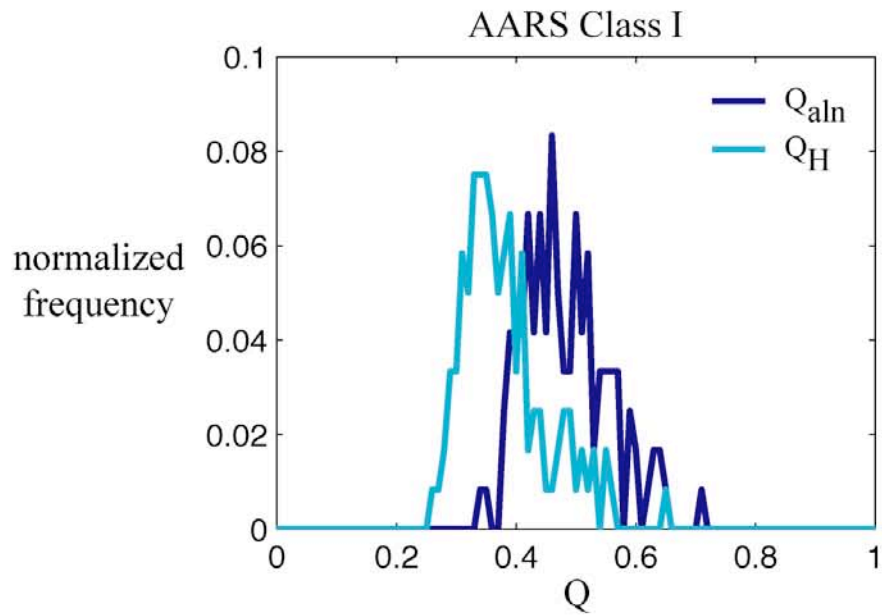
Representative set of
OB folds involved in
translation ?





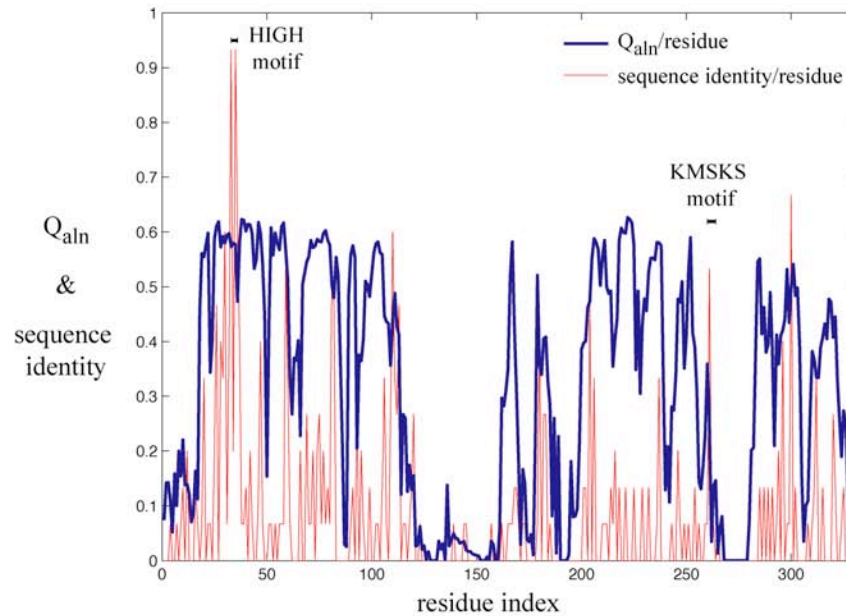


Only structure can reveal distant evolutionary relationships

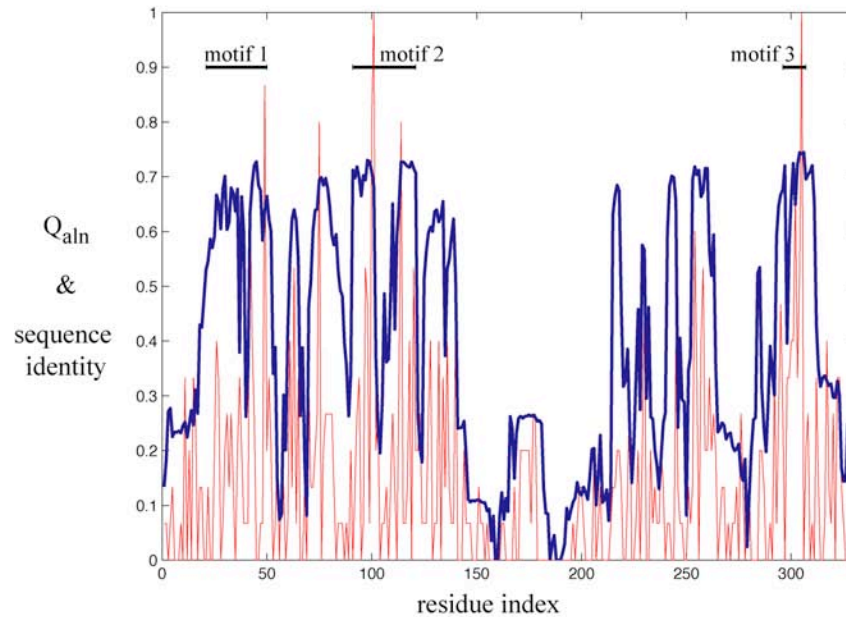


Conservation of Sequence and Structure

GlnRS

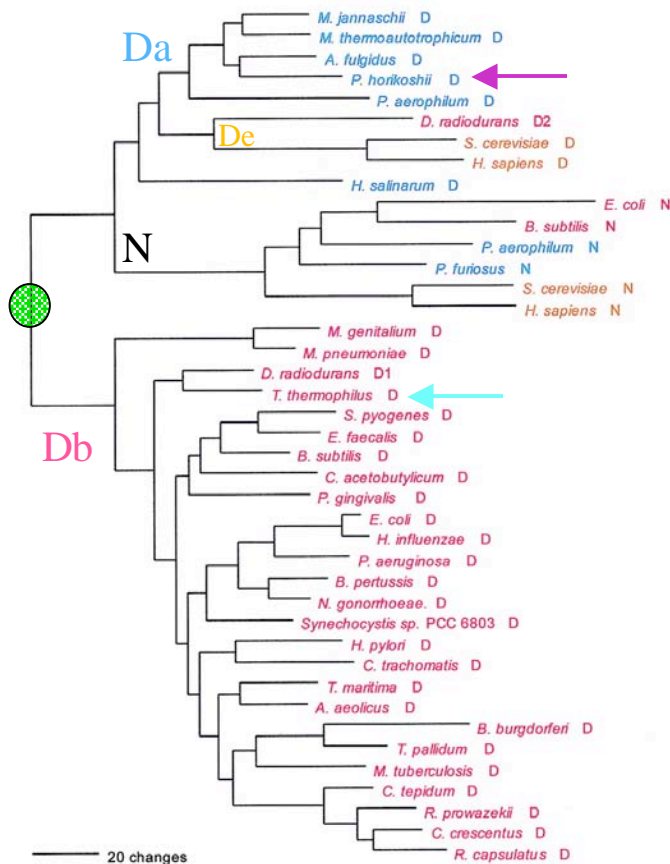


AsnRS

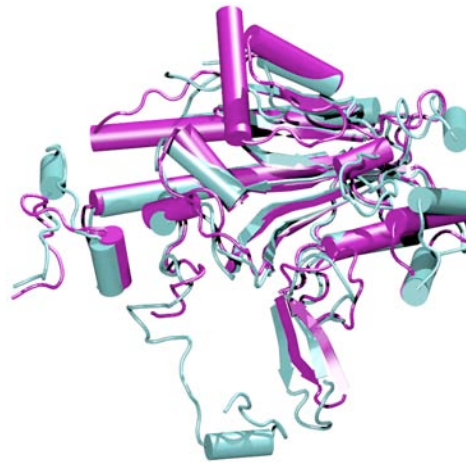


Protein structure encodes evolutionary information

Sequence Phylogeny
Woese et al. 2000

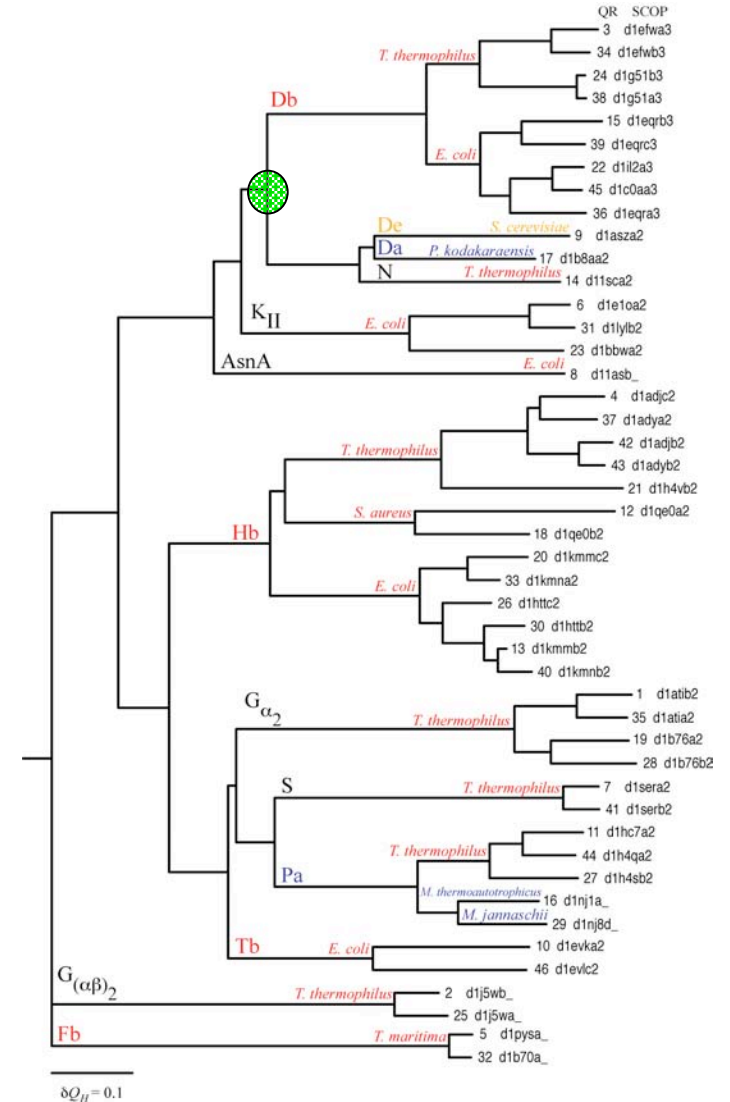


Structural Overlap



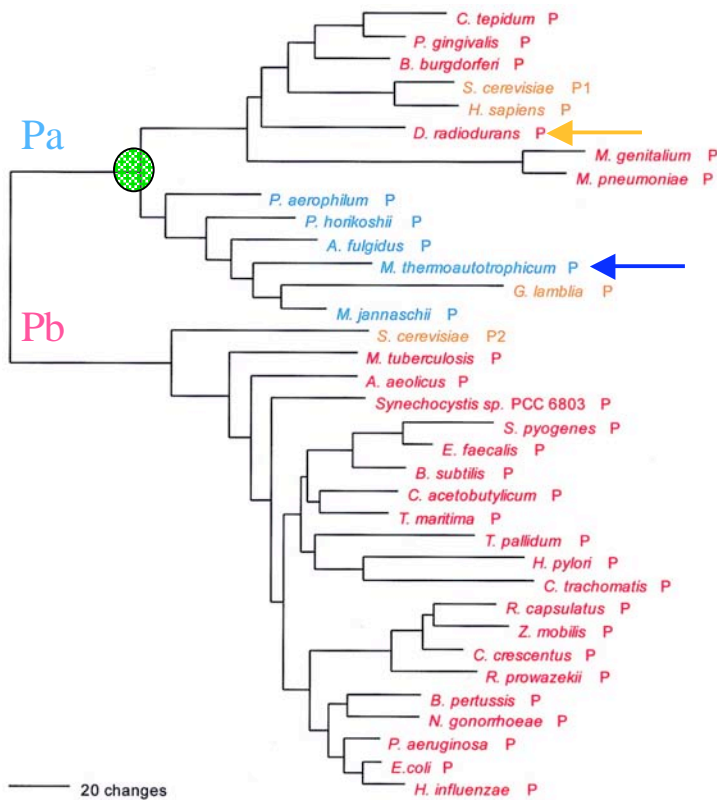
T. thermophilus
P. kodakaraensis

Structure Phylogeny
Class II AARSs

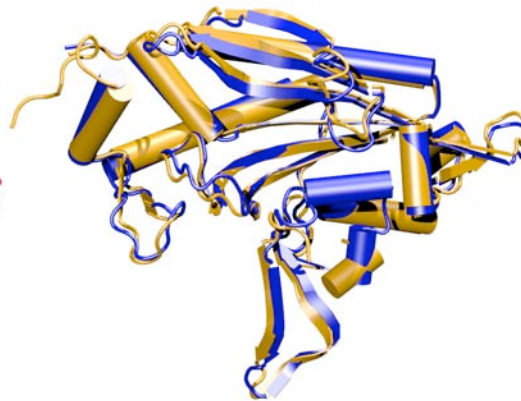


Horizontal Gene Transfer and Protein Structure in ProRS

Sequence Phylogeny
Woese et al. 2000

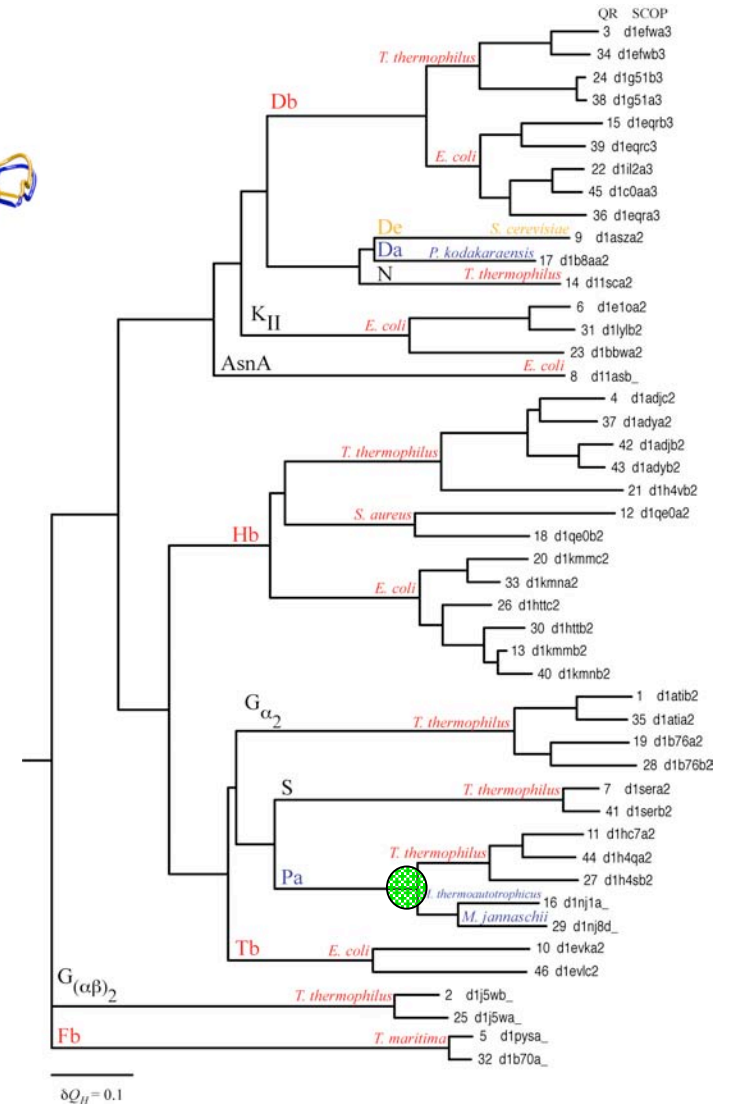


Structural Overlap

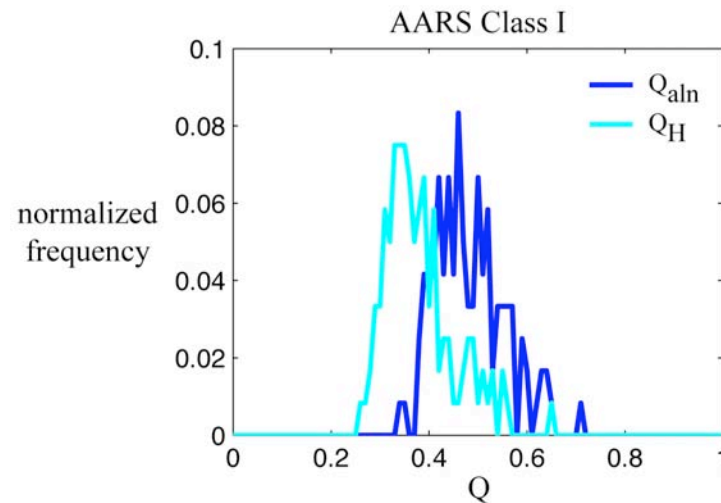
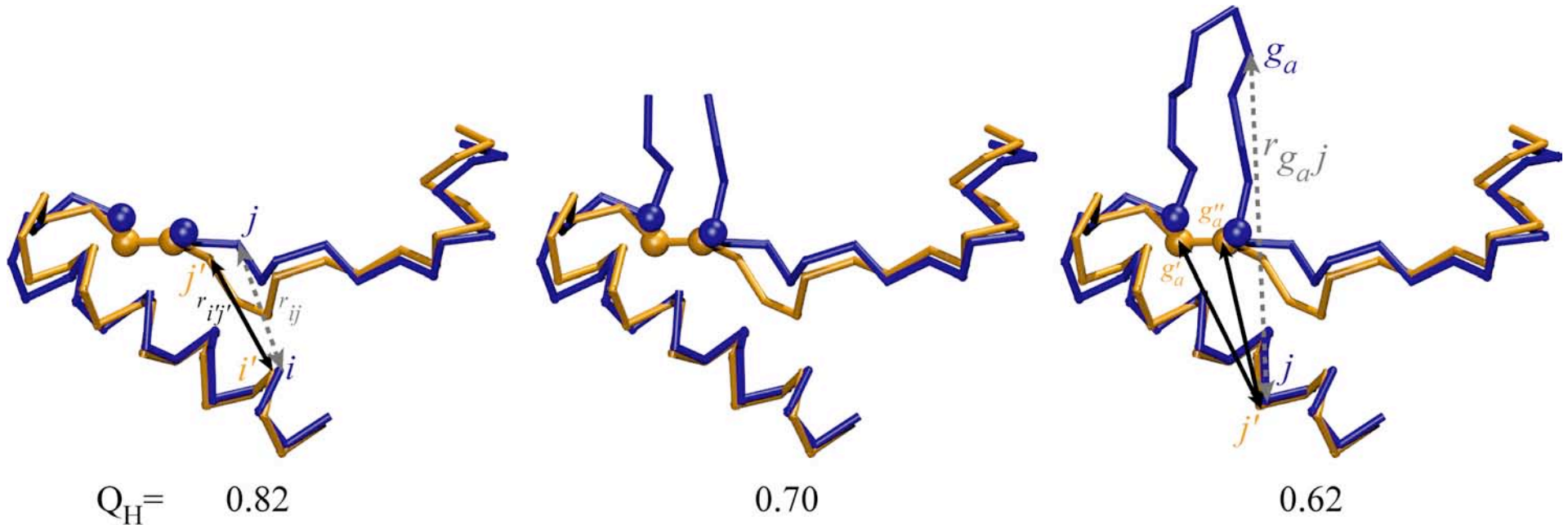


T. thermophilus
M. thermoautotrophicum

Structure Phylogeny
Class II AARs



Structural Homology Measure the effect of insertions



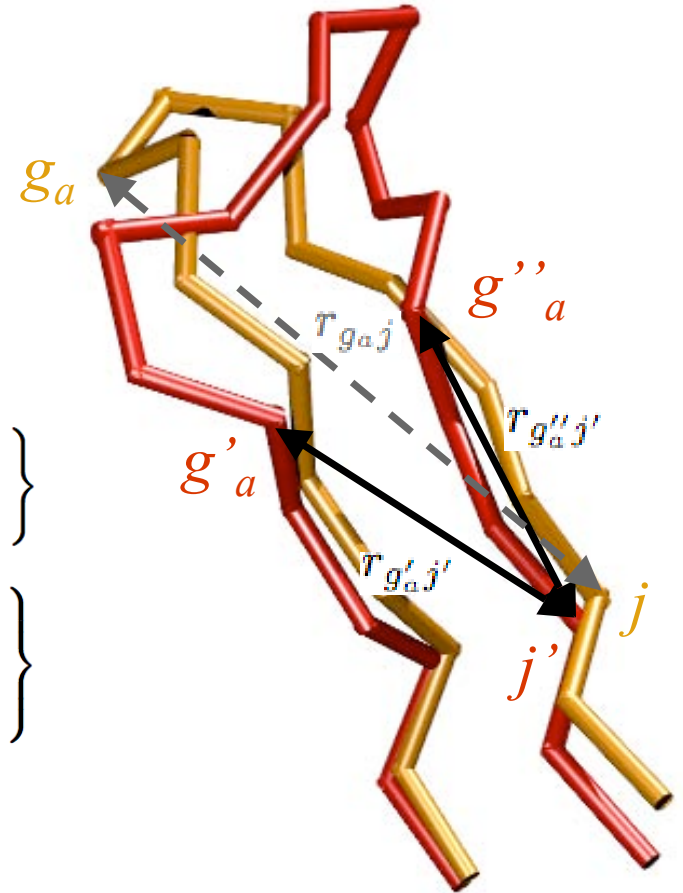
“Gaps influence the analysis
But should not dominate it” CW

Structural Homology Measure

compare inserted residues to gap edges

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

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



QR Factorization

Solve the least squares problem $Ax = b$

by triangularizing A with an orthogonal transformation.

$$Q^T A = \begin{bmatrix} R \\ 0 \end{bmatrix} \quad Q^T (Ax) = Q^T (b)$$

The system is now solved by back substitution,

$$\begin{bmatrix} R \\ 0 \end{bmatrix} x = \begin{bmatrix} c_1 \\ c_2 \end{bmatrix} \quad Rx = c_1$$

with a minimum residual of

$$\|r\|_2 = \|b - Ax\|_2 = \|c_2\|_2$$

Multi-Dimensional QR

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

N-dimensional QR = N one-dimensional QRs.

Permutation matrix is constant for each dimension, ordering norm is Frobenius-like matrix p-norm.

$$\max_{j=k,\dots,n_{proteins}} (\|a_j\|_{F_p}) \quad \|a_j\|_{F_p} = \left(\sum_{d=1}^4 \sum_{i=k}^{m_{aln}} |a_{ijd}|^p \right)^{1/p}$$

Encoding Structure

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

Gap “residues”: $(0, 0, 0, g)$

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

Encoding Sequence

Orthogonal Encoding = 24-space
23 amino acids symbols (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)

L. Heck, J. Olkin, and K. Nagshineh (1998) *J. Vibration Acoustics* **120**:663.

P. O’Donoghue and Z. Luthey-Schulten (2003) *Micro. Mol. Biol. Rev.* **67**:550-571.

QR Factorization with Column Pivoting

1. Calculate column norm of column i and all columns to the right.

$$A^{(k-1)} = H_{k-1} \dots H_1 A P_1 \dots P_{k-1}$$

$$\max_{j=k, \dots, n} (s_j^{(k)}) \quad s_j^{(k)} = \left(\sum_{i=k}^m a_{ij}^2 \right)^{1/2} \quad \leftarrow \text{Ordering Norm}$$

2. Swap column i with column to the right of maximum norm and record column permutation.

$$H_{k-1} \dots H_1 A P_1 \dots P_{k-1} P_k$$

3. Construct and apply H_k

$$A^{(k)} = H_k H_{k-1} \dots H_1 A P_1 \dots P_{k-1} P_k$$

$$\tilde{A} = A P_1 \dots P_n = A P$$

Original matrix, A , columns ordered by **increasing linear dependence.**

Protein Structure Prediction

1-D protein sequence

Ab Initio protein folding

SISSIRVKS KRIQLG...

Threading/Profile Alignment

$$E_{AM} = - \sum_{\mu=1}^{N_{\mu}} \sum_{i,j} \left\{ \gamma_{AM} [P_i, P_j, P_i^{\mu}, P_j^{\mu}] \right\}$$

$$X \exp \left[\frac{-(r_{ij} - r_{ij}^{\mu})^2}{2\sigma_{ij}^2} \right]$$

$$E = E_{match} + E_{gap}$$

Target Sequence

SISSRVKSKRIQLGLNQAELAQKV-----GTTQ...
 QFANEFKVRRIKLGYTQ-----TNVGEALAAVHGS...

Known structure(s)

3-D protein structure



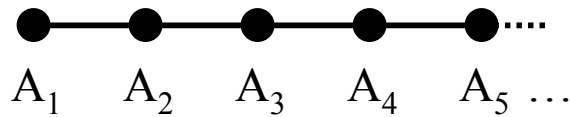
Eastwood, Hardin, Luthey-Schulten, Wolynes (2001)

IBM. J.RES.&DEV.45:475-497

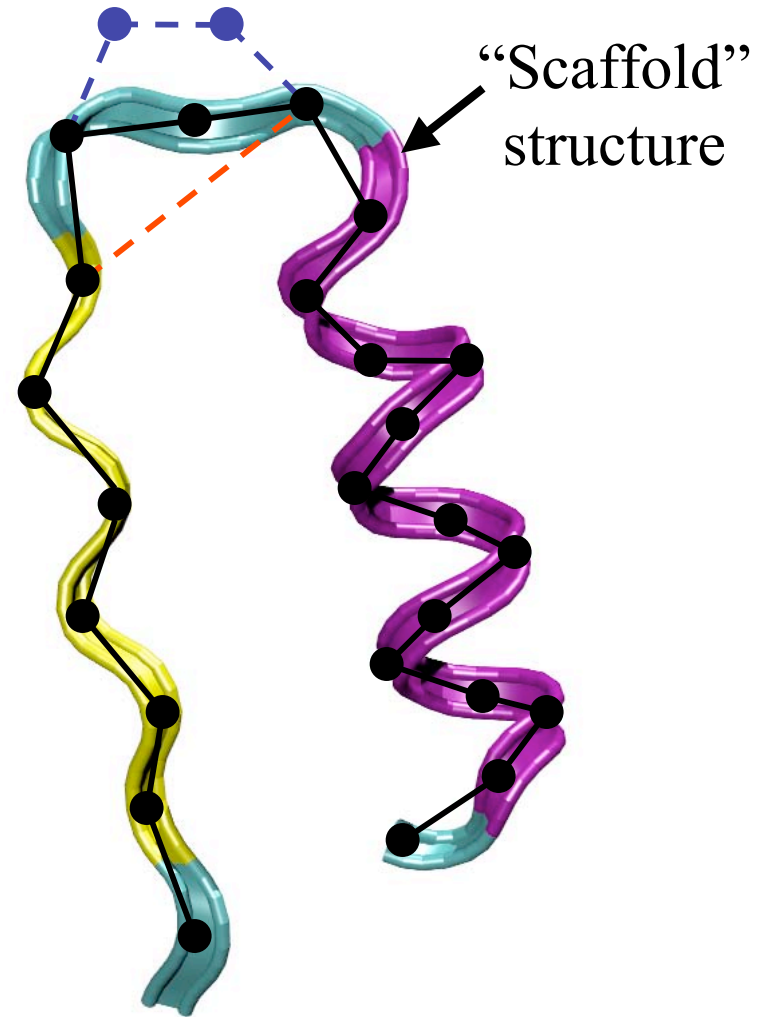
Papoian, et.al. PNAS (2004)

Sequence-Structure Alignment

Target sequence



Alignment between
target(s) and scaffold(s)



1. Energy Based Threading*

$$H = E_{contact} + E_{profile} + E_{H-bonds} + E_{gap}$$

$$E_{profile} = \sum_i^n \gamma^{(p)}(A_i, SS_i, SA_i)$$

$$E_{contact} = \sum_{i,j} \sum_{k=1}^2 \gamma_k^{(ct)}(A_i, A_j) * U(r_k - r_{ij})$$

2. Sequence – Structure Profile Alignments

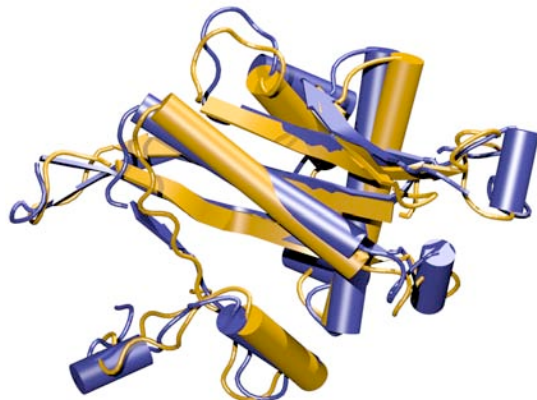
Clustal, Hidden Markov (HMMER, PSSM)
with position dependent gap penalties

*R. Goldstein, Z. Luthey-Schulten, P. Wolynes (1992, PNAS), K. Koretke et.al. (1996, Proteins)

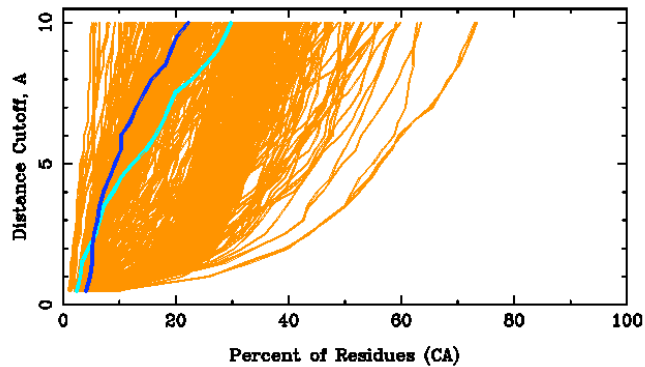
CASP5

Fold Recognition/Threading

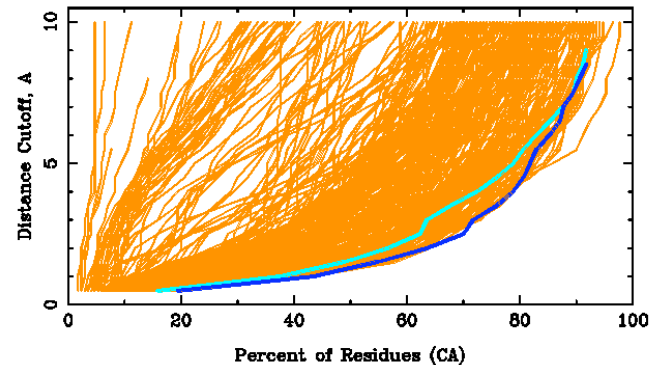
Schulten-Wolynes Group



T017ZTS093_1



T019ZTS093_1

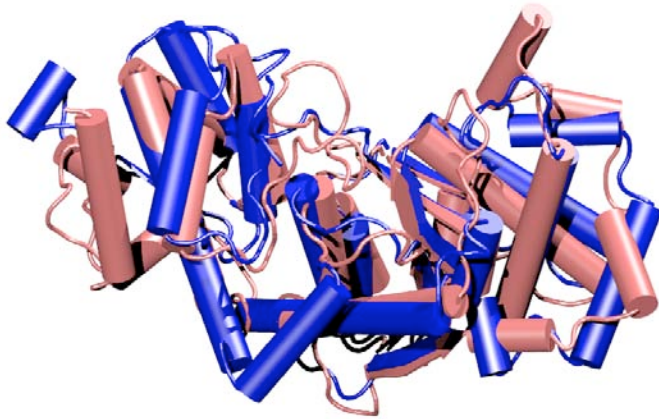


The prediction is never better than the scaffold.

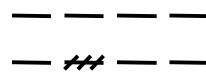
Threading energy function/profiles requires improvement.

Why Study the Evolution of Protein Structure?

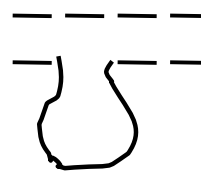
In what specific ways has the evolutionary dynamic changed protein shape over time?



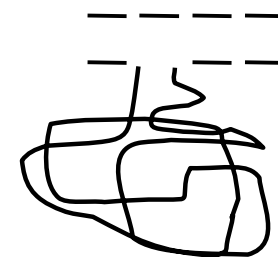
Substitution



Indel



Domain Insertion



implications for protein structure prediction,
protein design

What can studying the change in protein shape over time tell us about the evolutionary process?

How did translation evolve?

When, with respect to the root of the universal phylogenetic tree, was translation established in its modern form?

What was the role of the AARSs in the evolution of the translation mechanism, development of the genetic code?