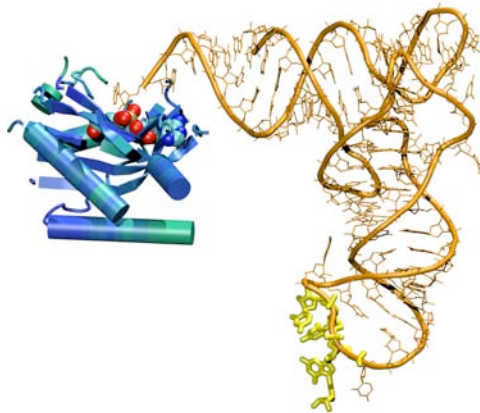


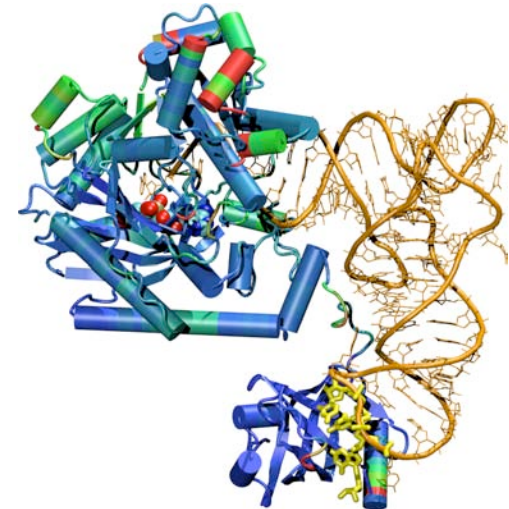
Evolutionary concepts in bioinformatics

evolution of protein structure
in the aminoacyl-tRNA synthetases



| | | Second position | | | | | |
|-----|-----|-----------------|-----|-----|-----|-----|----------------|
| | | U | C | A | G | | |
| U | UUU | Phe | UCU | UAU | UGU | U | Third position |
| | UUC | | | | | | |
| | UUA | Leu | UCA | UAA | UGA | A | |
| | UUG | | | UCG | UAG | UGG | |
| CUU | Leu | CCU | CAU | CGU | U | | |
| CUC | | | | | | CCC | CAC |
| CUA | Pro | CCA | CAA | CGA | A | | |
| CUG | | | | | | CCG | CAG |
| A | AUU | Ile | ACU | AAU | AGU | U | |
| | AUC | | | | | | ACC |
| | AUA | Thr | ACA | AAA | AGA | A | |
| | AUG | | | | | | ACG |
| GUU | Val | GCU | GAU | GGU | U | | |
| GUC | | | | | | GCC | GAC |
| GUA | Ala | GCA | GAA | GGA | A | | |
| GUG | | | | | | GCG | GAG |

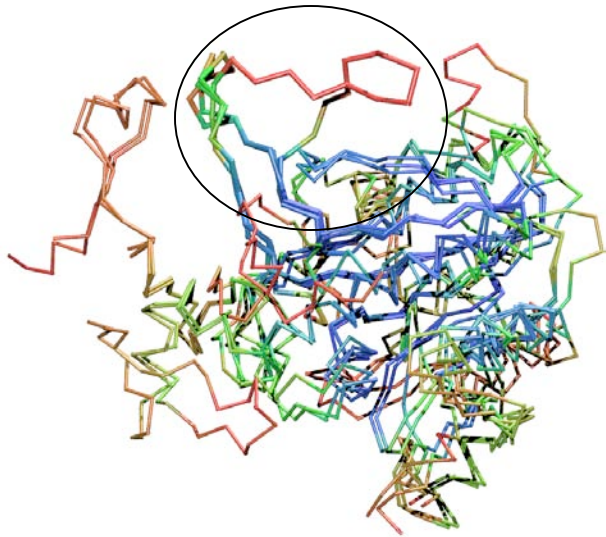
Copyright © 2004 Pearson Education, Inc., publishing as Benjamin Cummings.



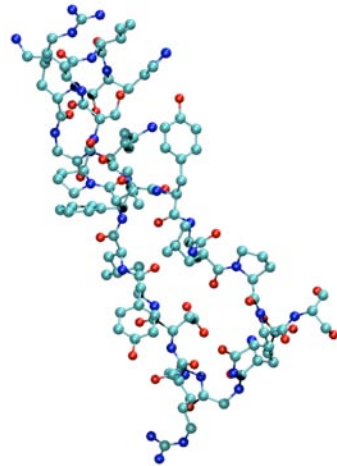
Patrick O'Donoghue* and Zan Luthey-Schulten

Department of Chemistry
University of Illinois at Urbana-Champaign

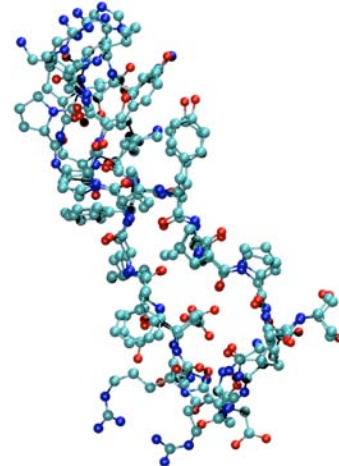
Protein Homology in Structure and Sequence



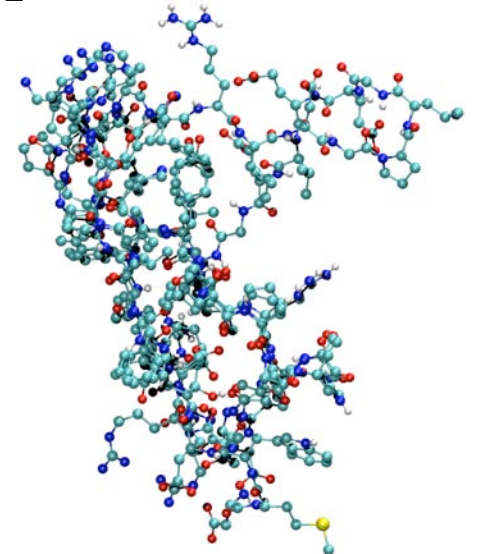
3 homologous structures;
2 closely related (Db1, Db2),
1 more distant (Fb).



Db1



Db1, Db2



Db1, Db2, Fb



Db1, Db2, Fb
backbone only

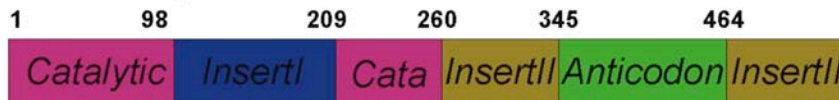
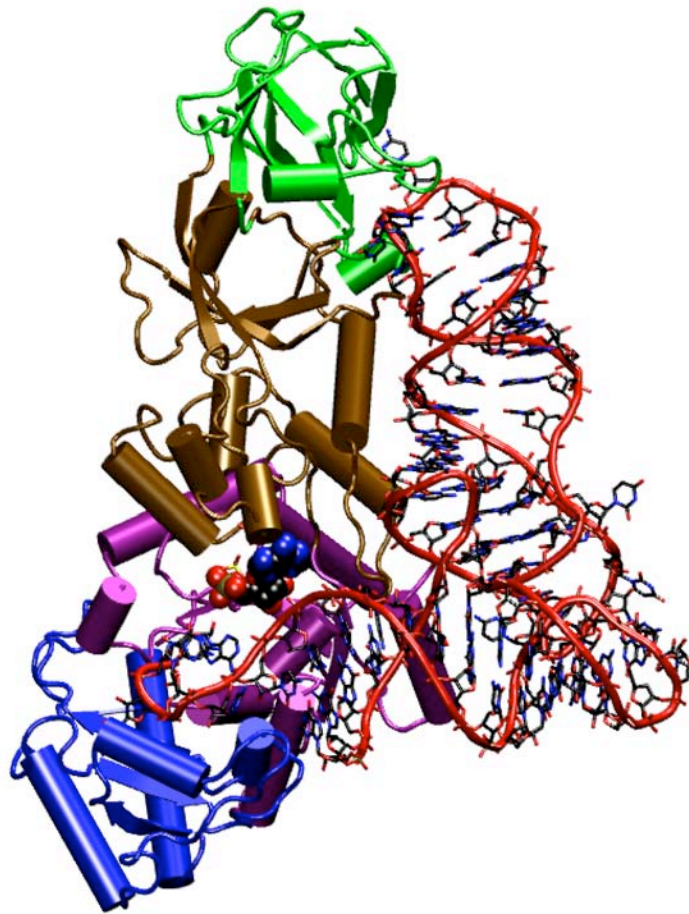
| | |
|-----|--|
| Db1 | E---GARDFLV-PYRHE-----PGLFYALPQS |
| Db2 | -E--GARDYLV-PSRVH-----KGKFYALPQS |
| Fb | ---DMWDTFWLT-GE--GFRLEGPLGEEVEGRLLLRTH |

What can be learned from AARSs?

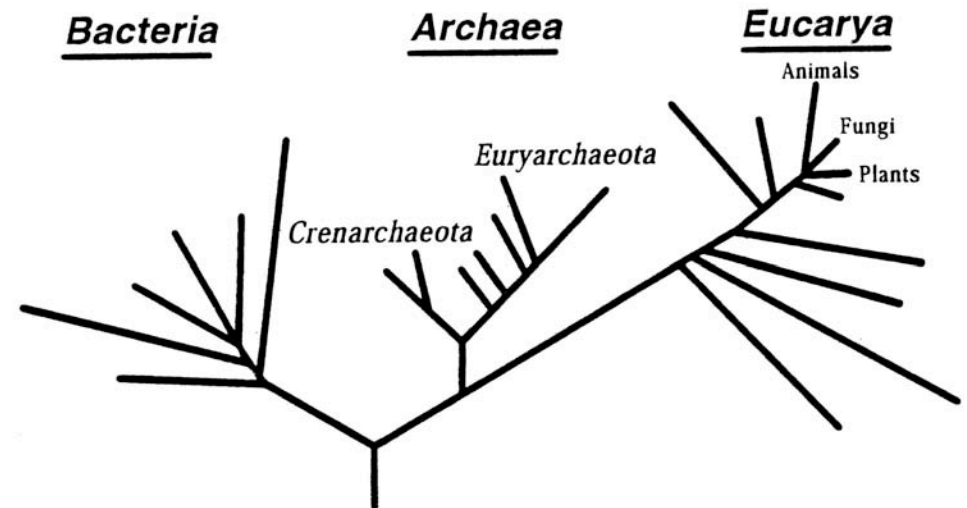
“The aminoacyl-tRNA synthetases, perhaps better than any other molecules in the cell, optimize the current situation and help to understand (the effects) of Horizontal Gene Transfer (HGT).”

Carl Woese (PNAS, 2000; MMBR 2000)

Aminoacyl-tRNA synthetases



Universal Tree of Life

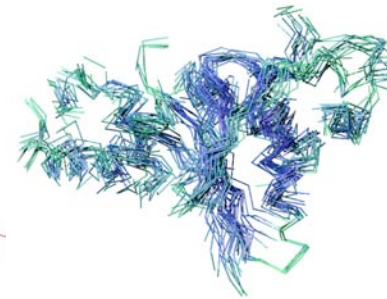
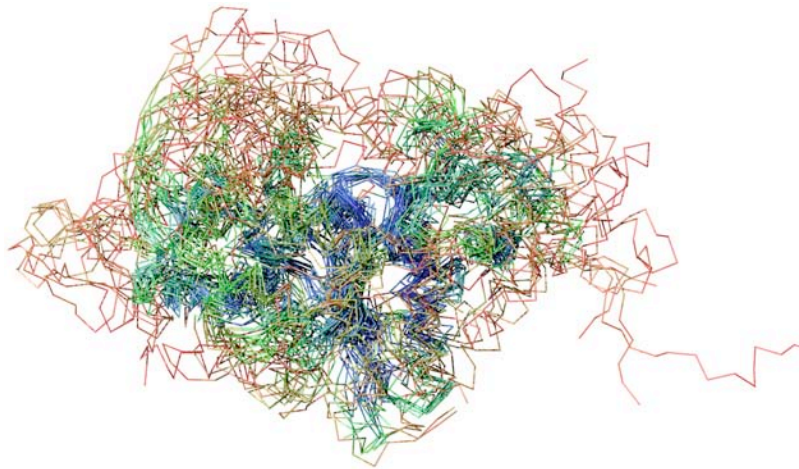


Woese *PNAS* 1990, 2002.

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.

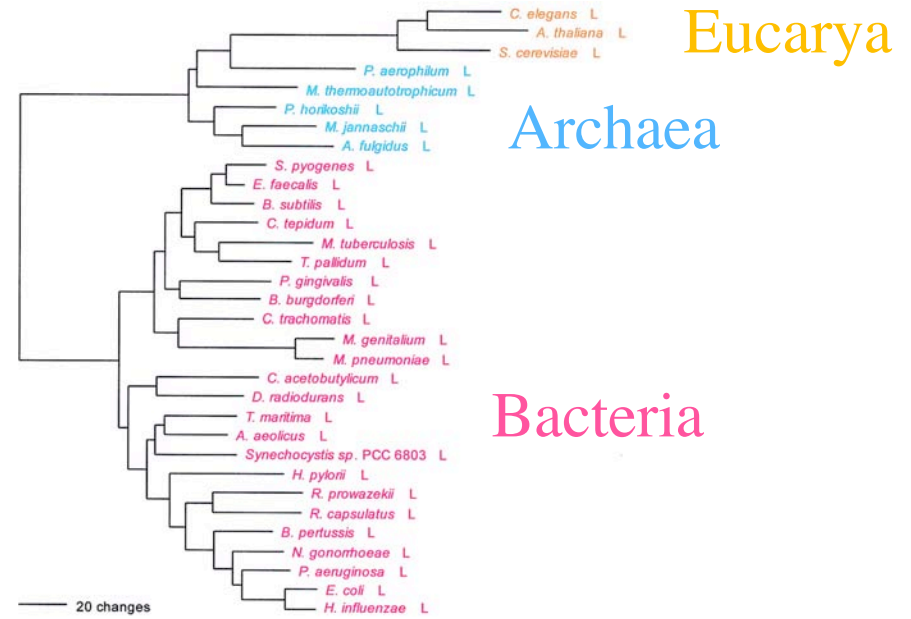
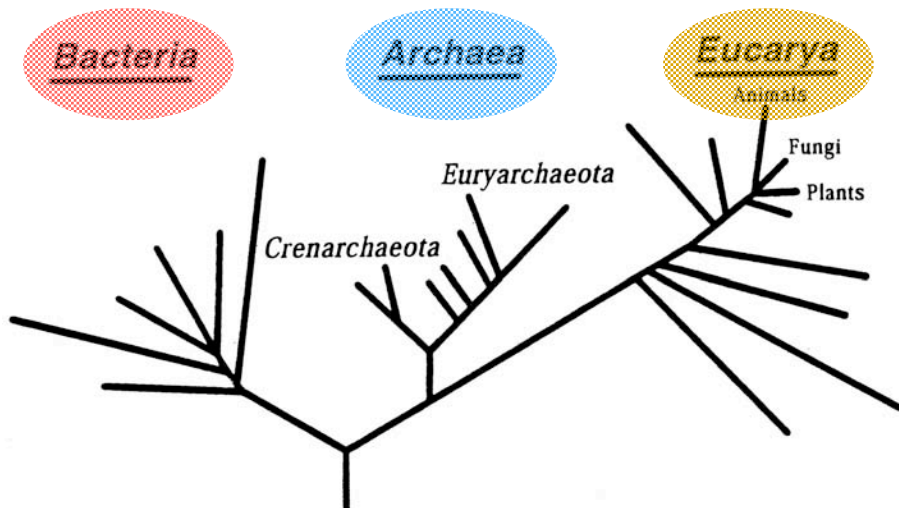
Eliminate bias inherited from structure or sequence databases.

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

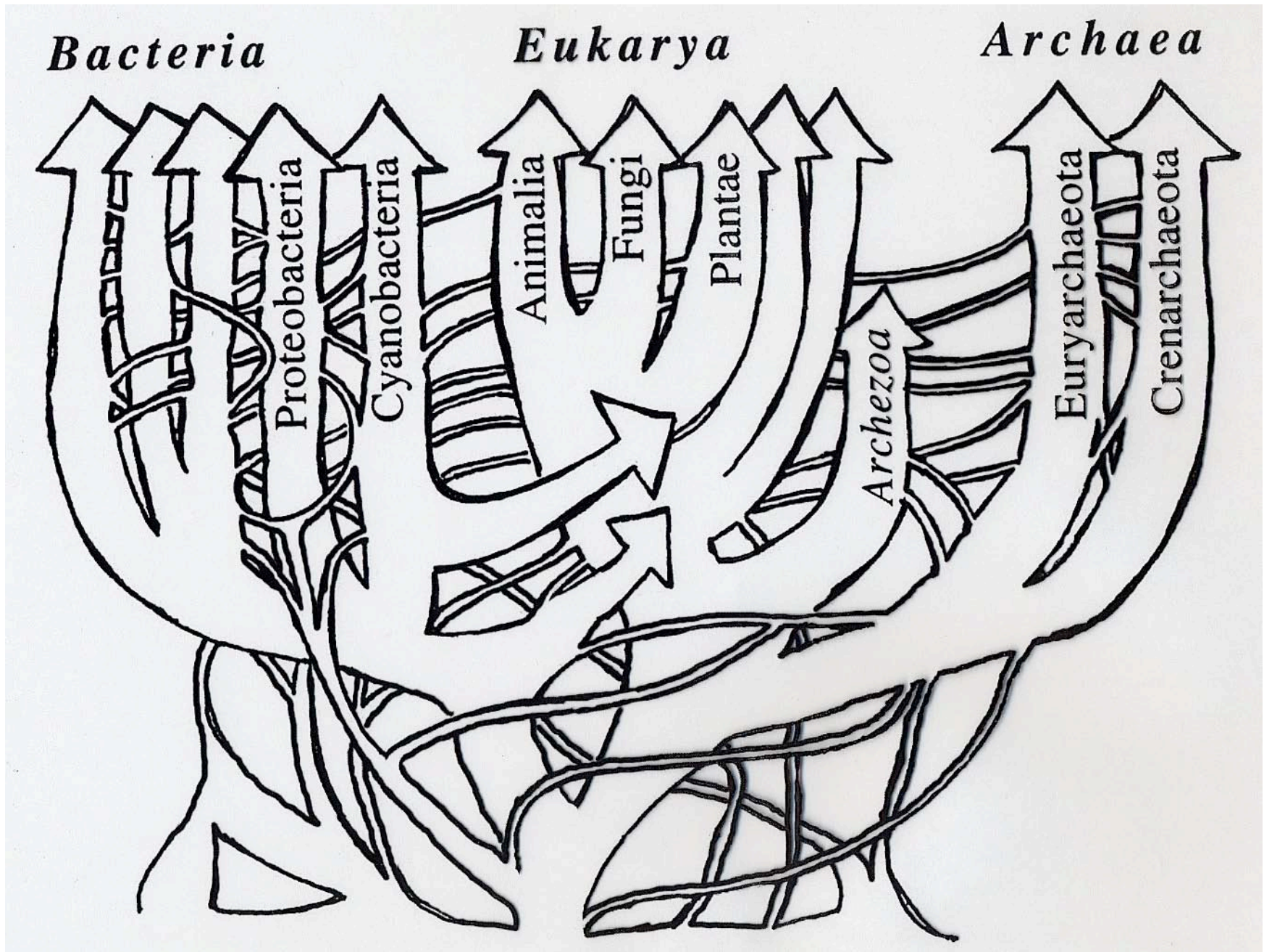
5. Depict the evolution of structure and function in Aspartyl-tRNA synthetase.

Universal Phylogenetic Tree

three domains of life



Leucyl-tRNA synthetase displays the full canonical phylogenetic distribution.



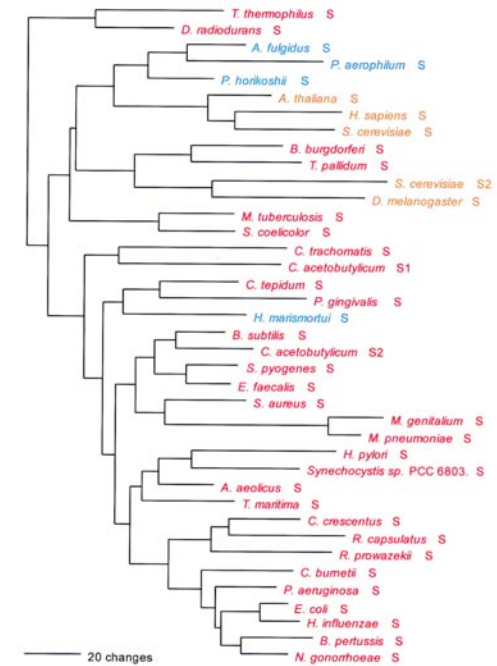
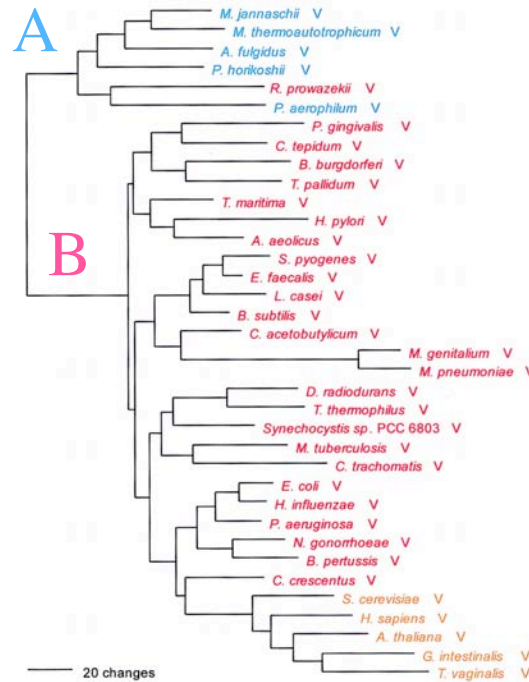
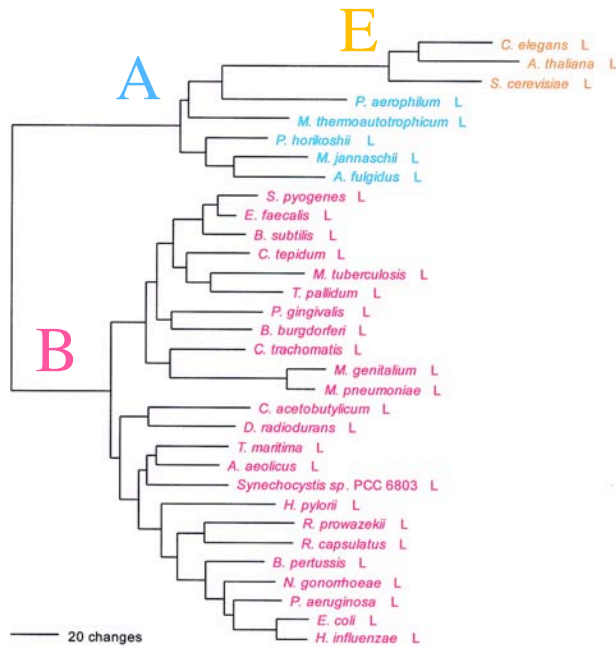
After W. Doolittle, modified by G. Olsen

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

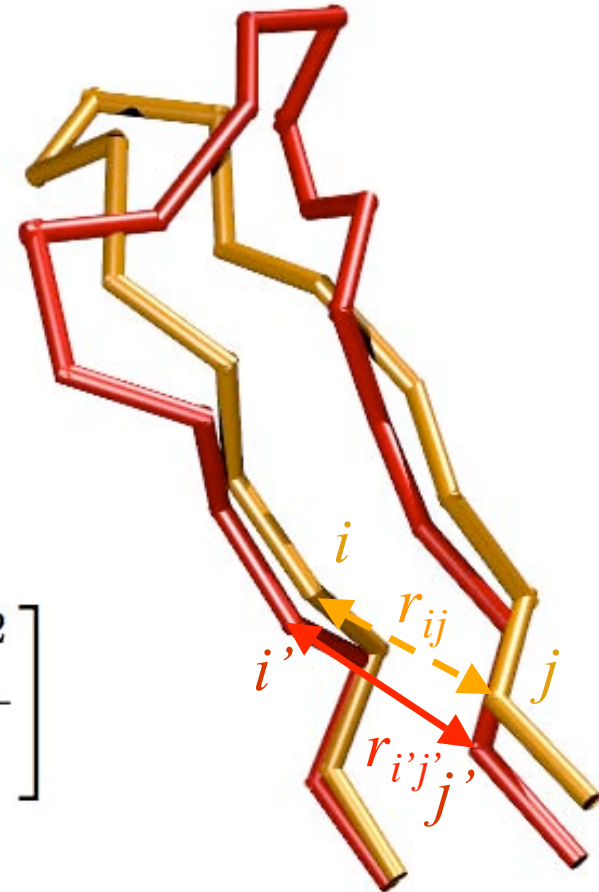
Protein Structure Similarity Measure

Q_H Structural Homology

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

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

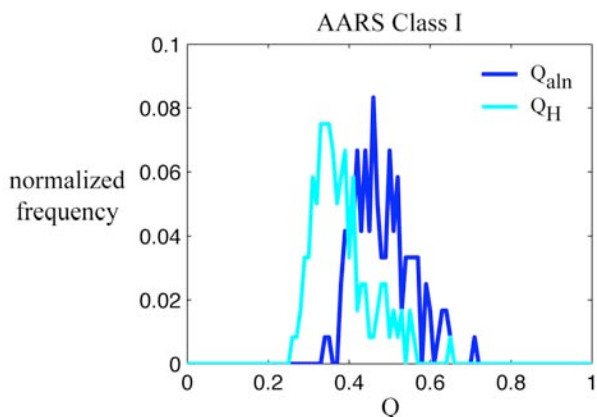
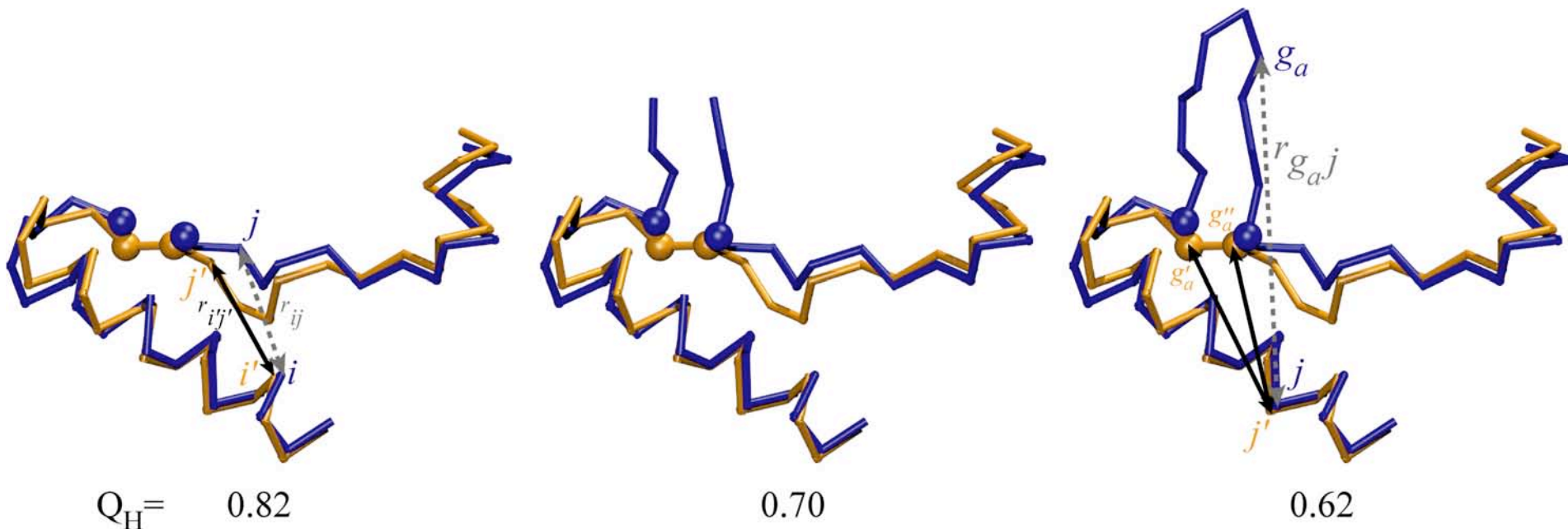
$$q_{aln} = \sum_{i < j-2} \exp \left[-\frac{(r_{ij} - r_{i'j'})^2}{2\sigma_{ij}^2} \right]$$



Structural Similarity Measure

the effect of insertions

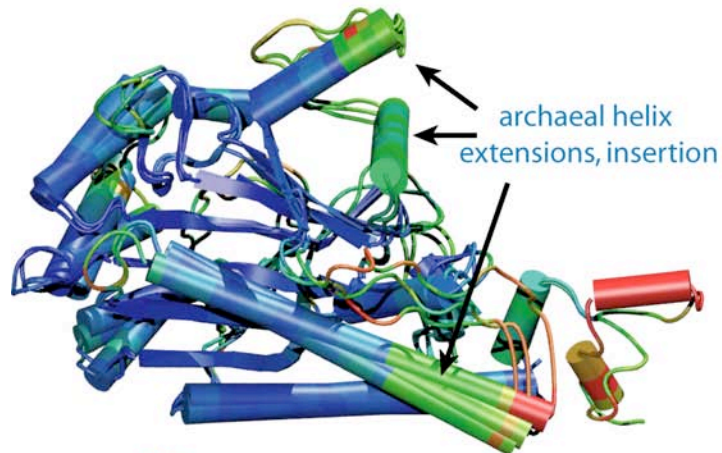
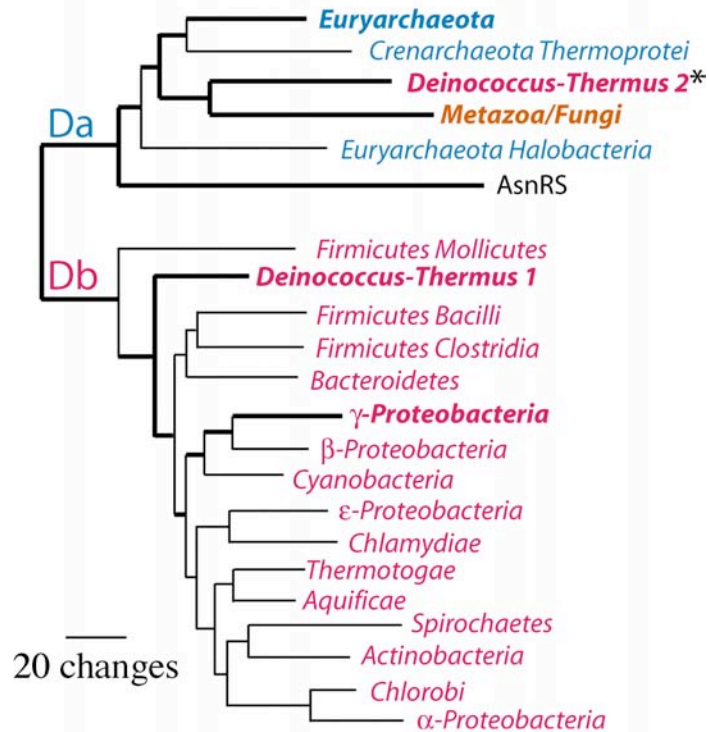
“Gaps should count as a character but not dominate” C. Woese



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

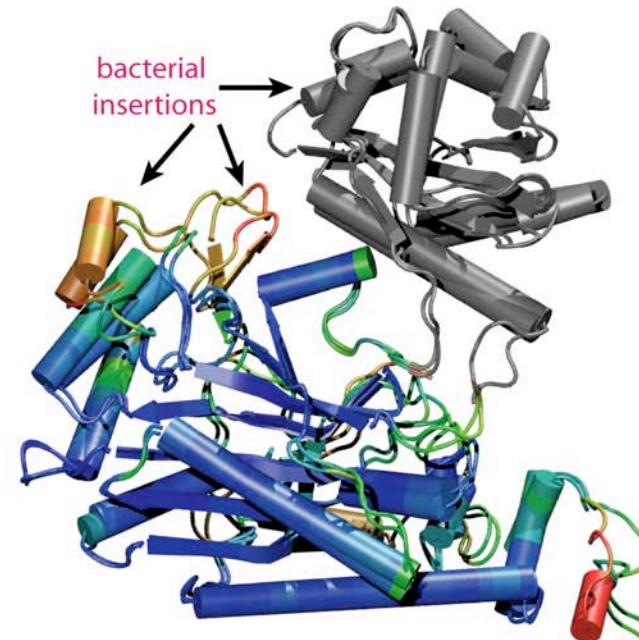
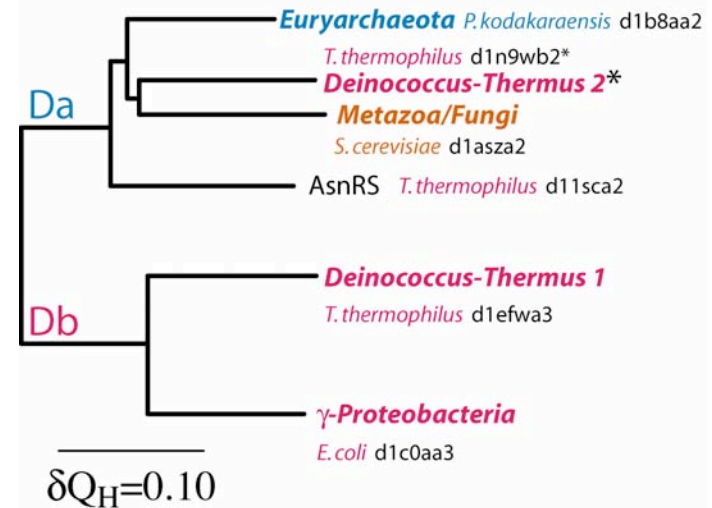
Protein structure encodes evolutionary information

sequence-based phylogeny



Da - AspRS archaeal genre

structure-based phylogeny

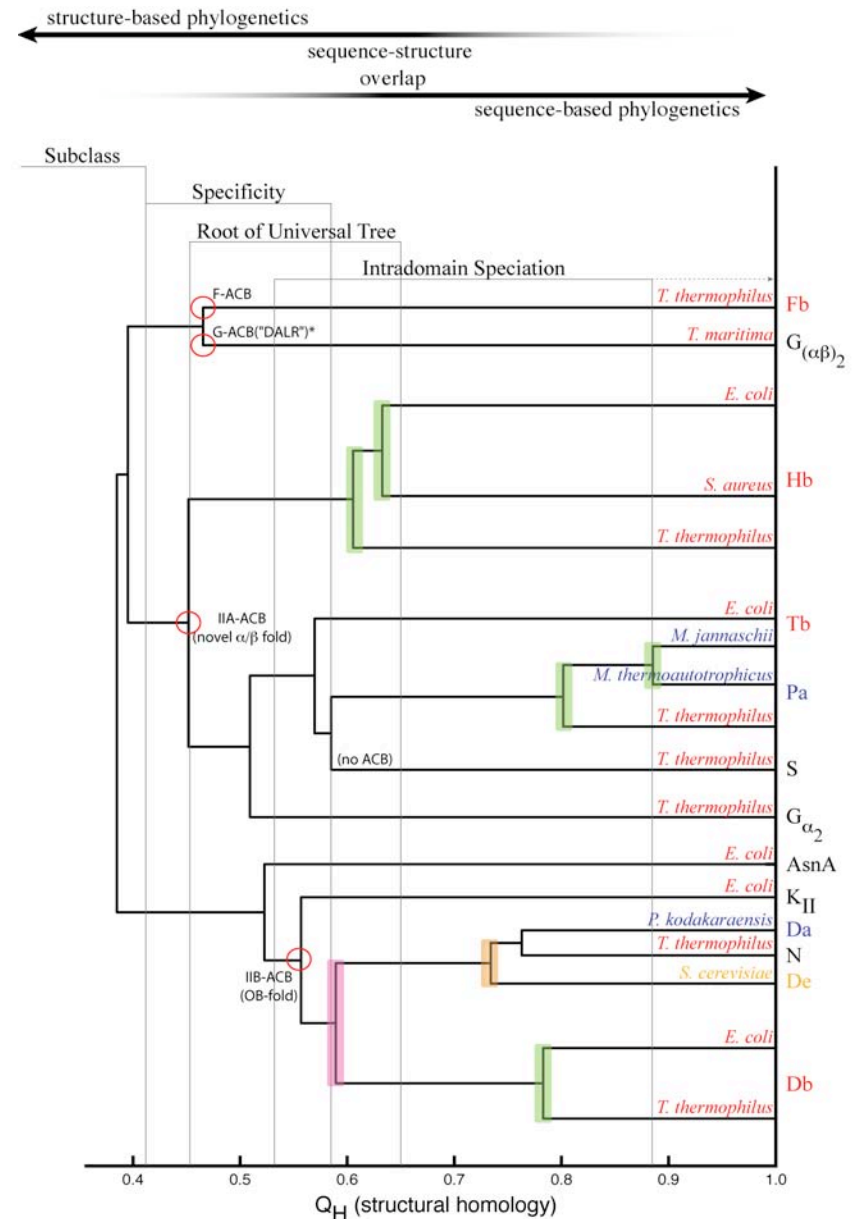
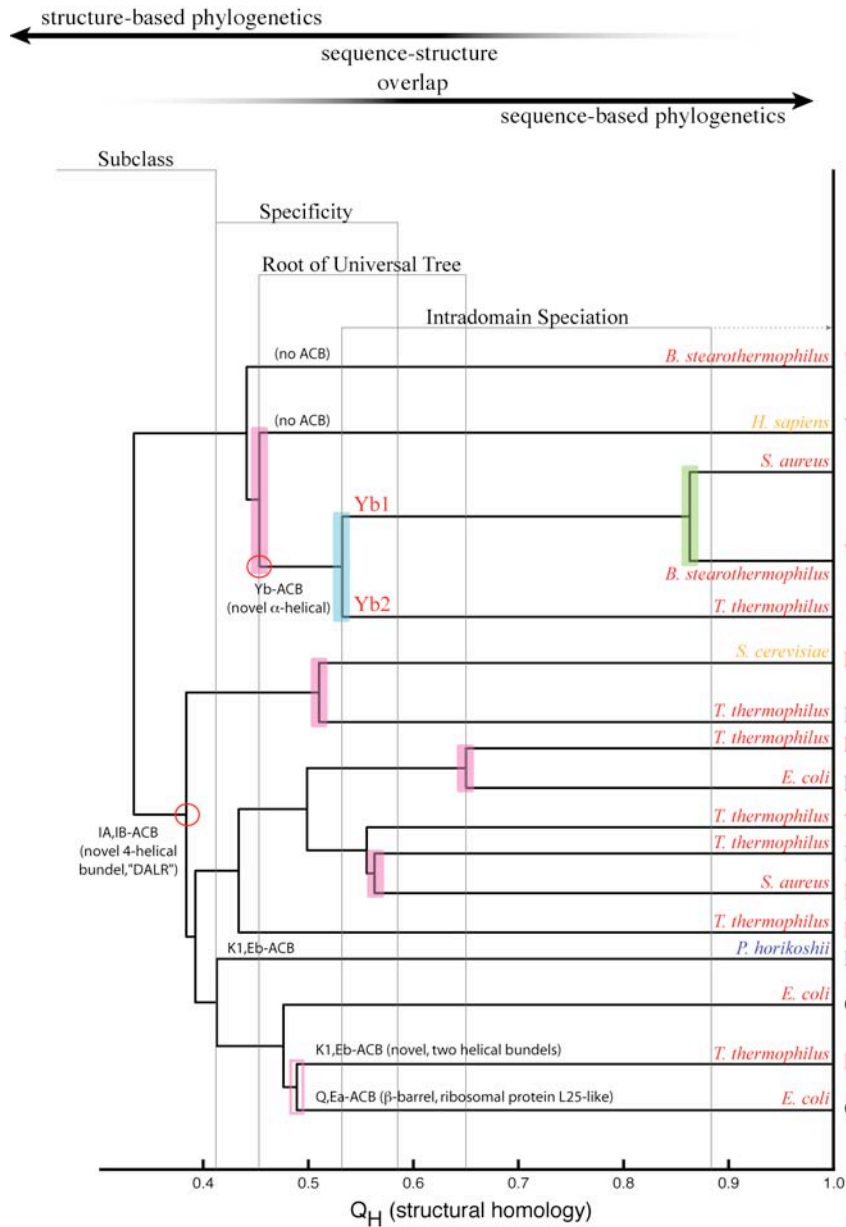


Db - AspRS bacterial genre

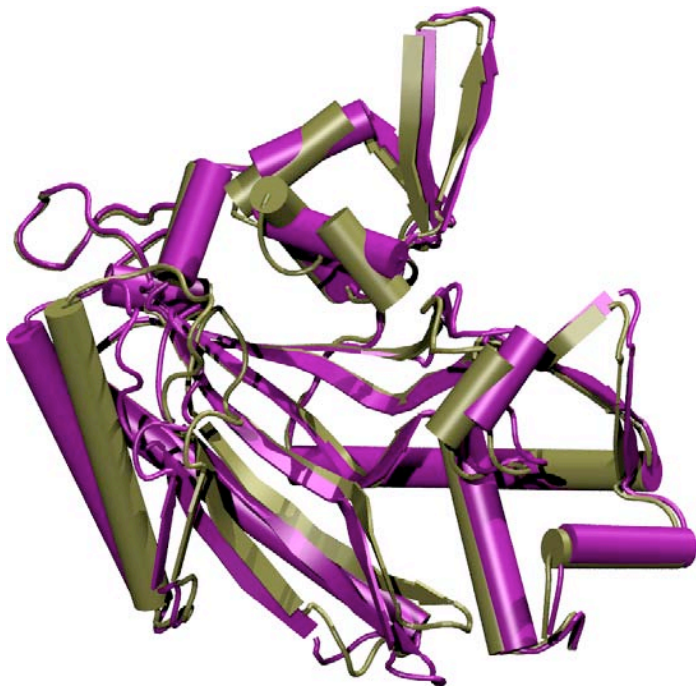
Protein structure reveals distant evolutionary events

Class I AARs

Class II AARs



Sequences define more recent evolutionary events



Conformational changes
in the same protein.

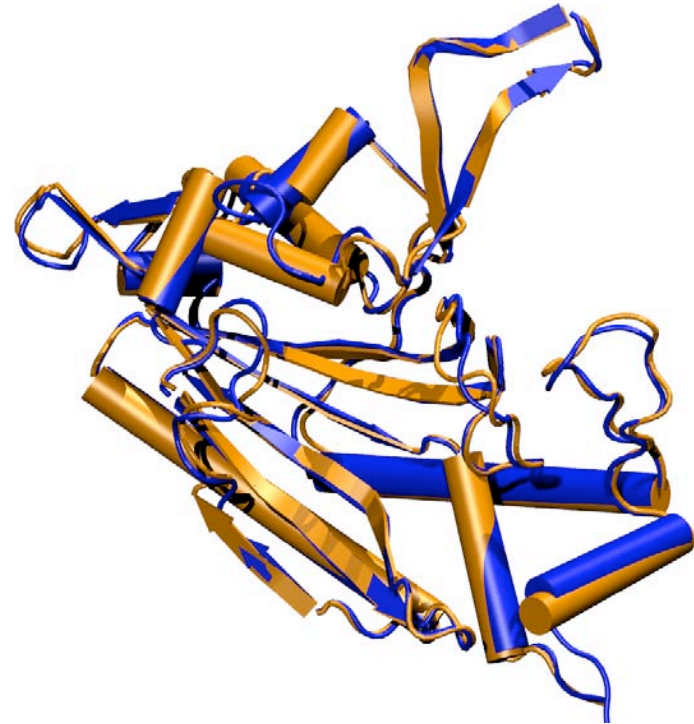
ThrRS

T-AMP analog, 1.55 Å.

T, 2.00 Å.

$Q_H = 0.80$

Sequence identity = 1.00



Structures for two
different species.

ProRS

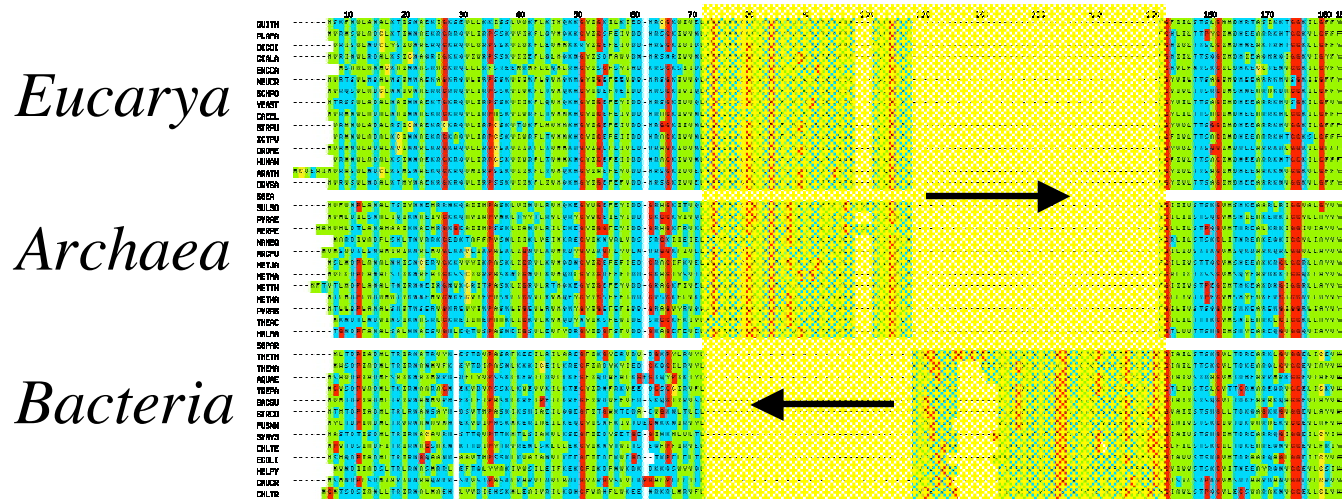
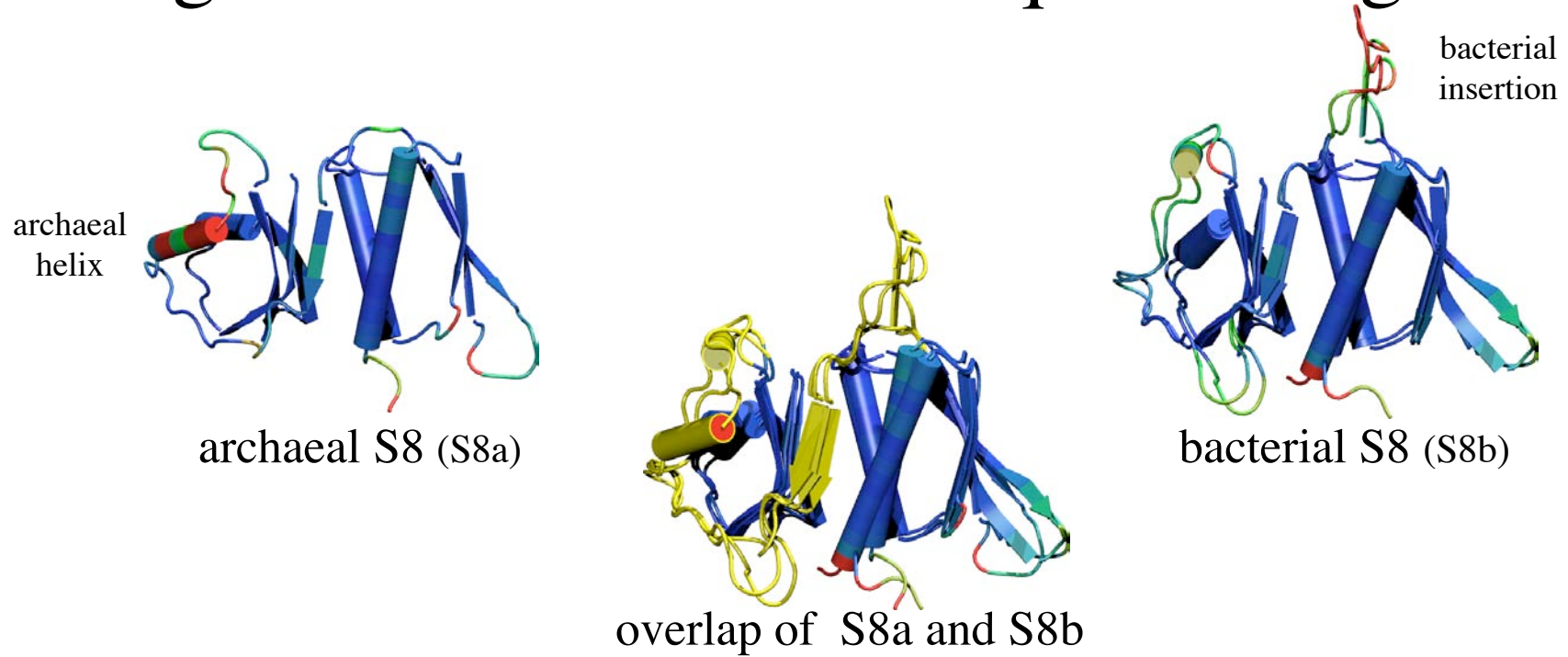
M. jannaschii, 2.55 Å.

M. thermoautotrophicus, 3.20 Å.

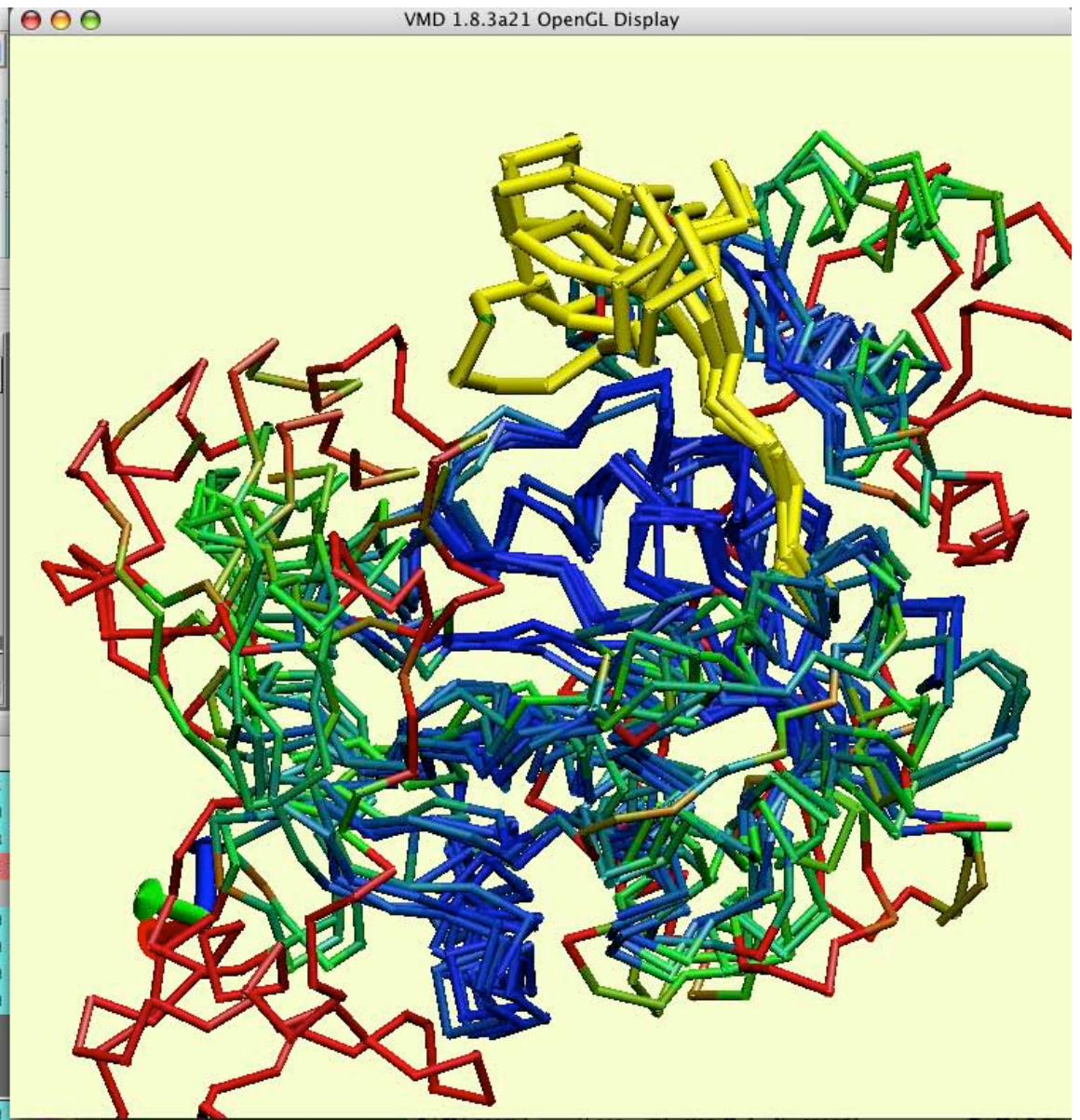
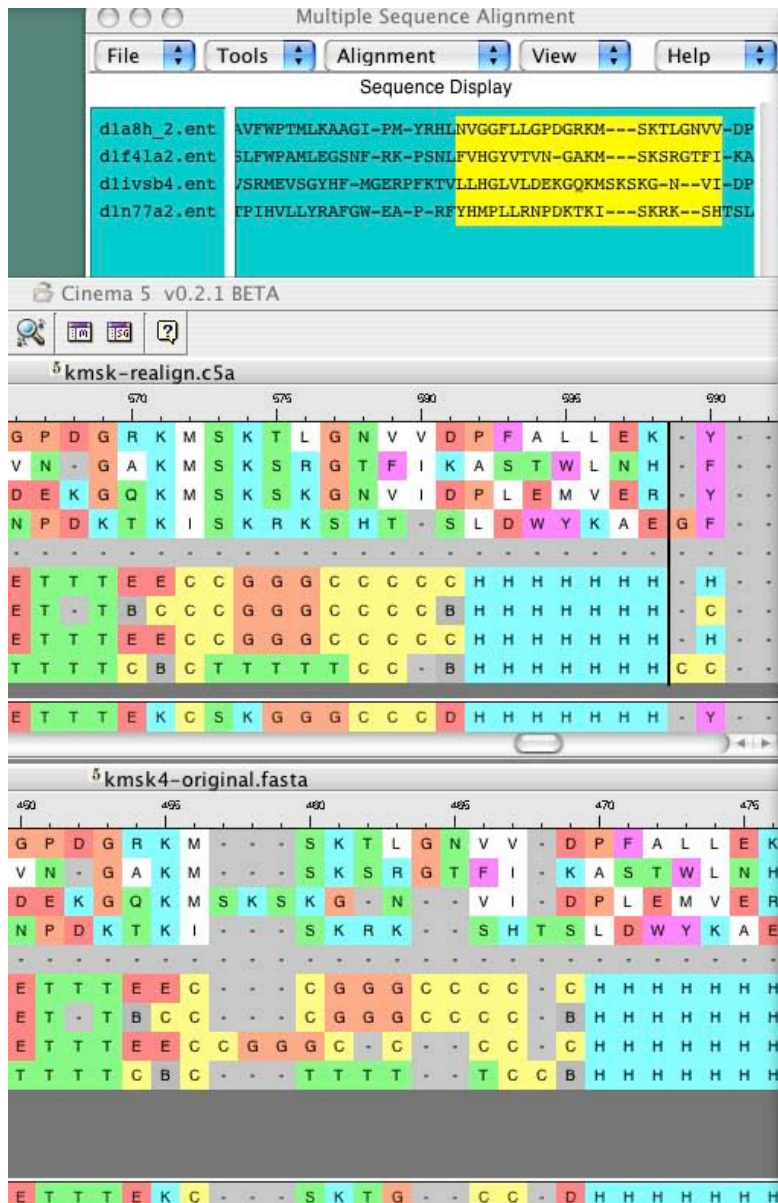
$Q_H = 0.89$

Sequence identity = 0.69

Using structure to correct sequence alignments

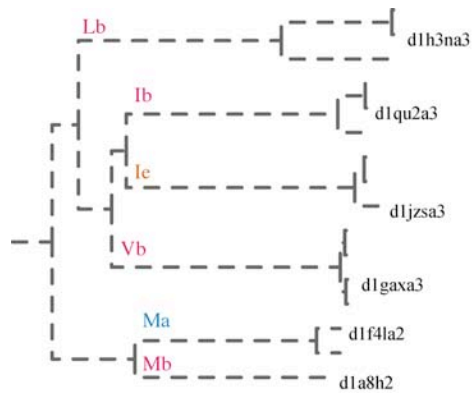


Conformational versus evolutionary change

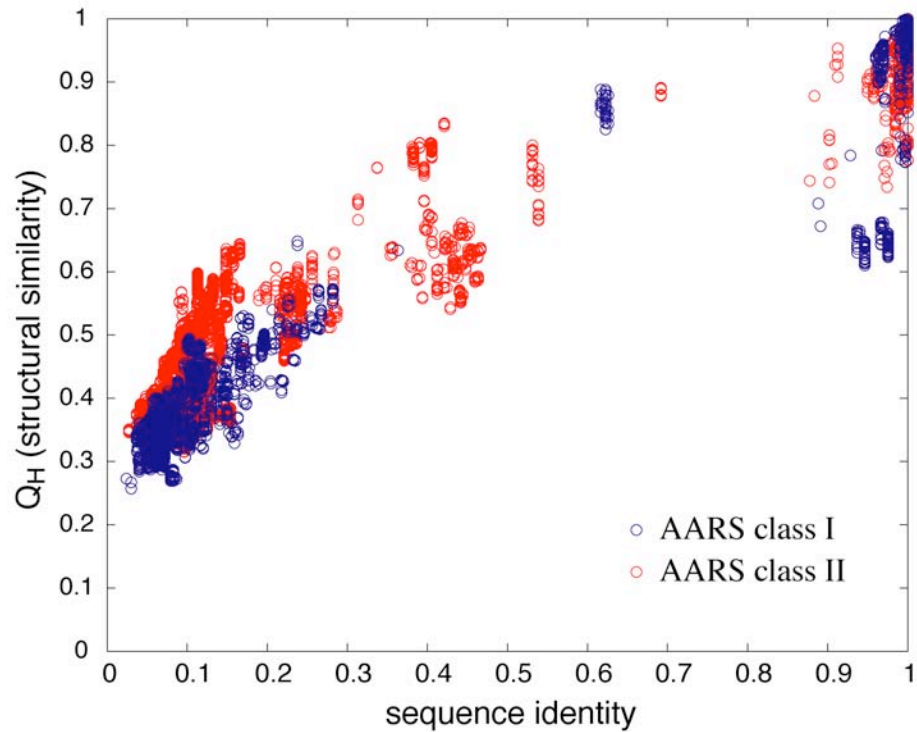
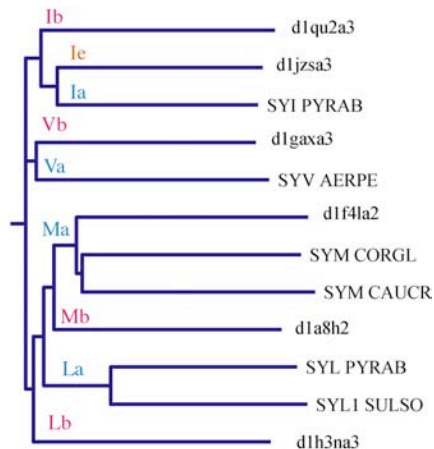


Towards a unified phylogenetic framework in sequence and structure

Structure-based tree

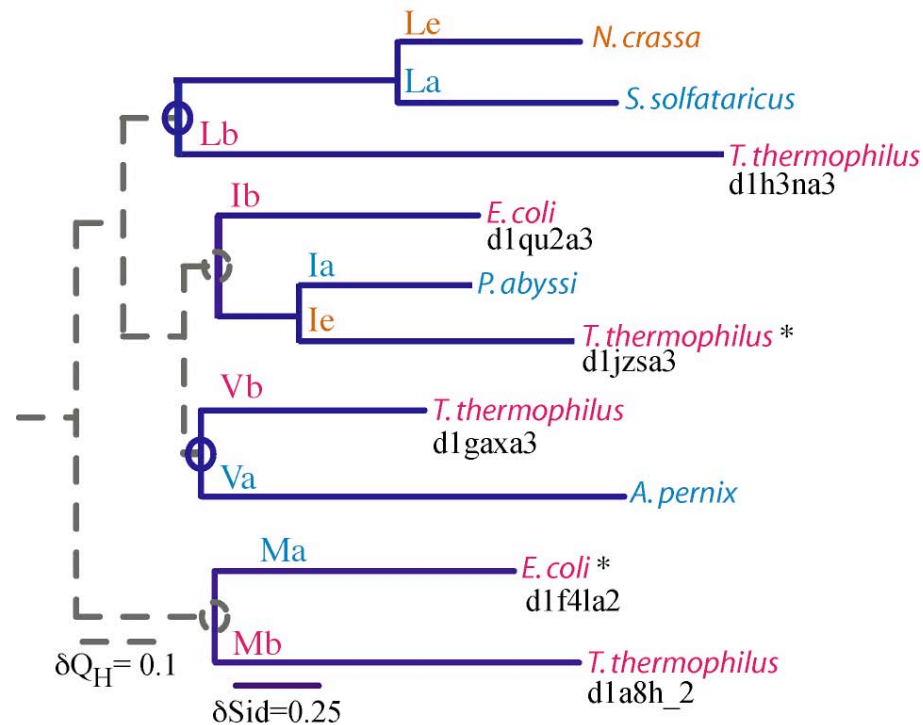


Sequence-based tree



Towards a unified phylogenetic framework in sequence and structure

Combined sequence-structure tree

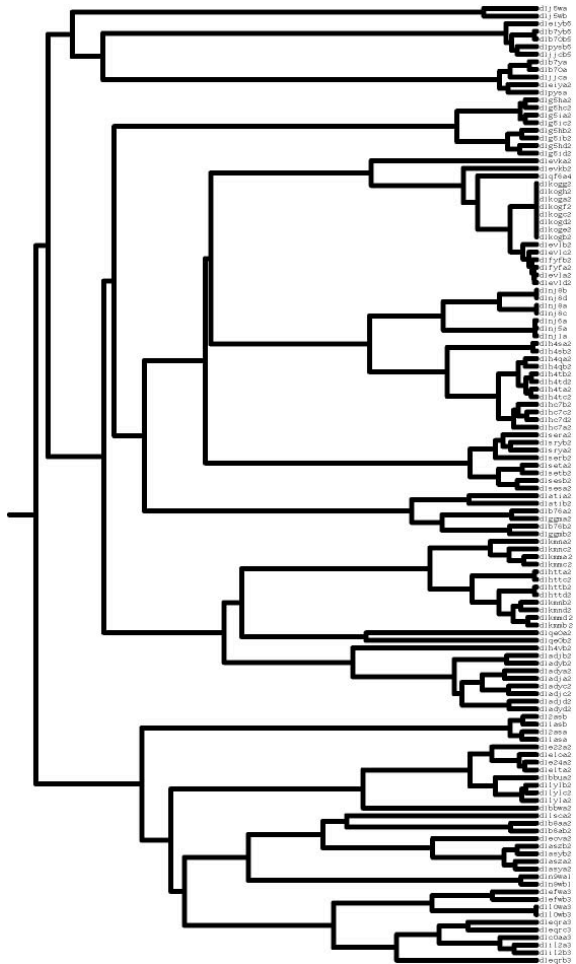


Structure is used to infer distant evolutionary events, i.e., the development of basic structures and functions.

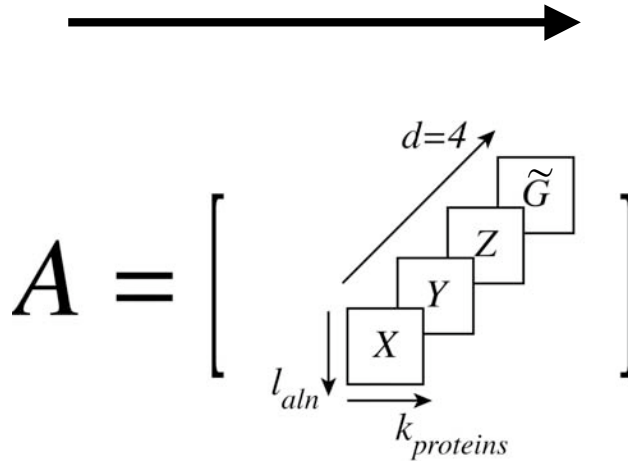
Sequences supplement the missing structure data, and define more recent evolutionary events, i.e., speciation.

Non-redundant representative sets

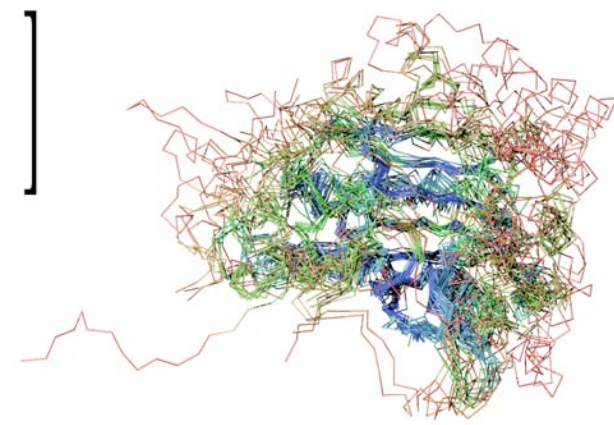
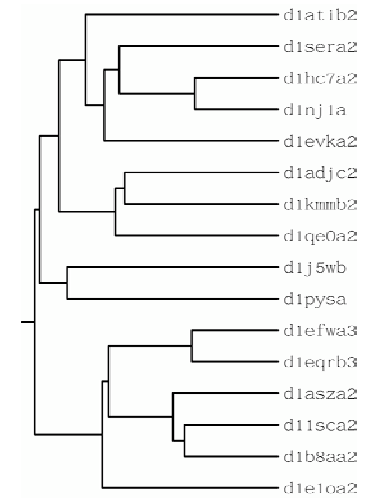
Too much information
129 Structures



Multidimensional QR
factorization
of alignment matrix, A .



Economy of information
16 representatives



QR computes a set of minimal linearly dependent structures.

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

P. O'Donoghue and Z. Luthey-Schulten (2004) *J. Mol. Biol.*, in press.

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 $\tilde{G} = gG \quad (0, 0, 0, g)$

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)

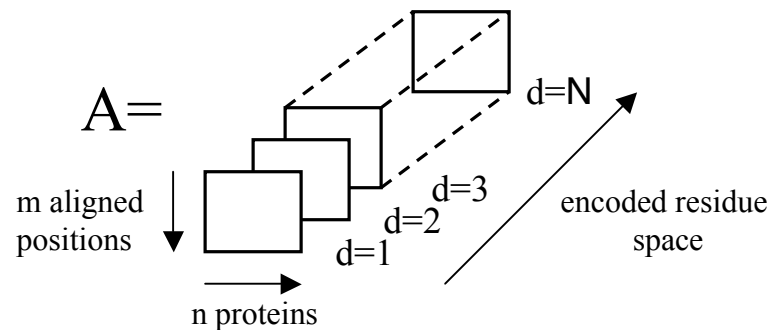
B = (0,1,0)

C = (0,0,1,0)

...

GAP = (0,1)

Alignment Matrix



A multiple alignment is a matrix with linearly dependent columns
redundancy is equivalent to linear dependence

QR factorization

Re-orders the columns of A, segregating the linearly independent columns from the dependent ones without scrambling the information in A. **SVD not an option.**

$$Q^T A P = \tilde{R}$$

$$\tilde{A} = A P$$

Q^T – orthogonal matrix of product of Householder transformations.

P – permutation matrix encodes column pivoting which exchanges columns of A and puts the redundant or similar proteins to the right hand side.

Multidimensional QR

N simultaneous QR factorizations, one for each d-dimension.

$$Q_{(d)}^T A_{(d)} P = Q_{(d)}^T \left[\begin{array}{c} \\ \\ \\ \end{array} \right] P = \tilde{R}_{(d)}$$

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

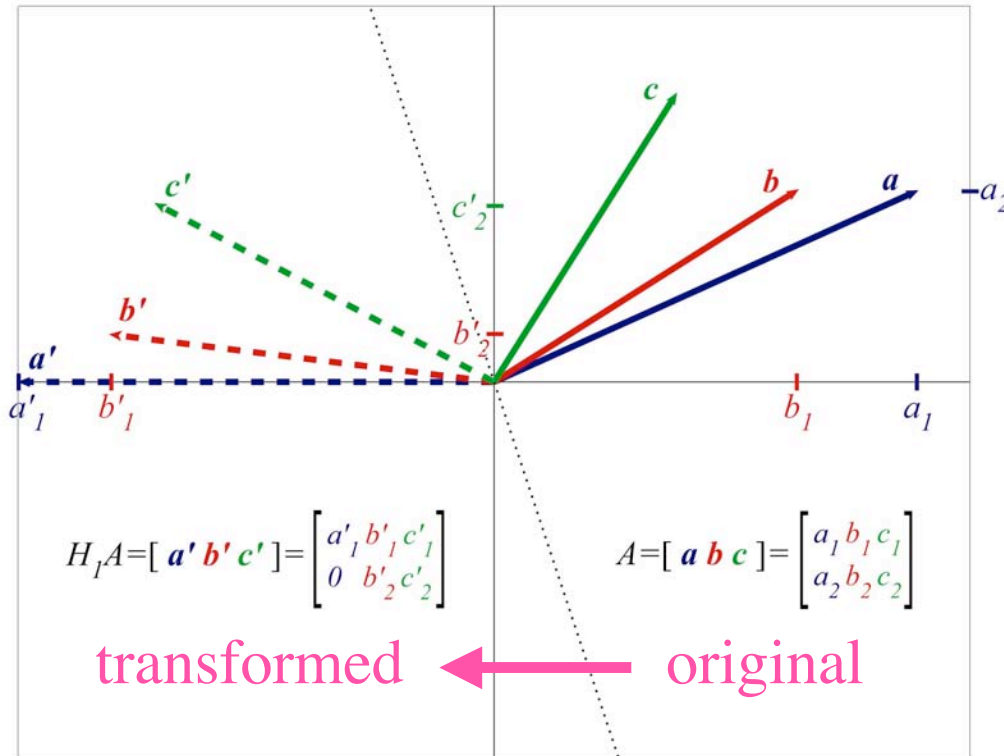
L. Heck, J. Olkin, and K. Nagshineh (1998) *J. Vibration Acoustics* **120**:663.

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

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**.

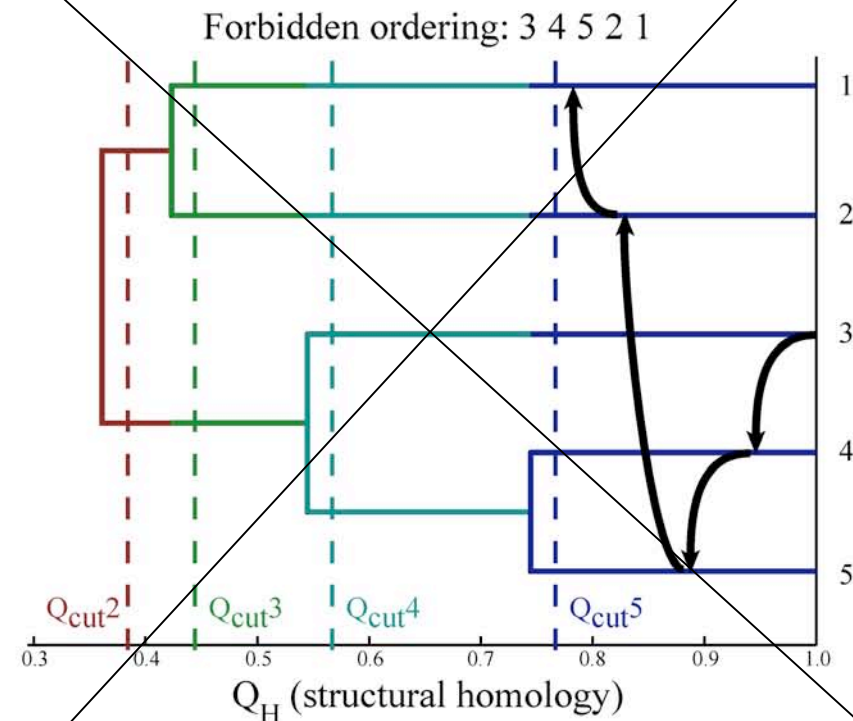
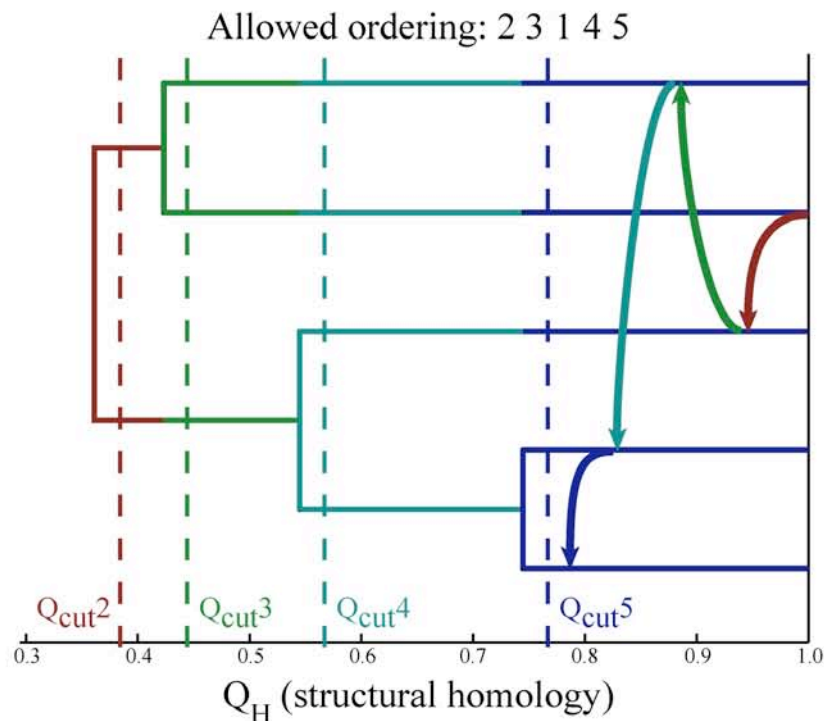
Multiply aligned proteins exist in a higher dimensional space, so this magnitude is computed with a matrix p-norm:

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

adjustable
parameter

What are the constraints on the parameters?

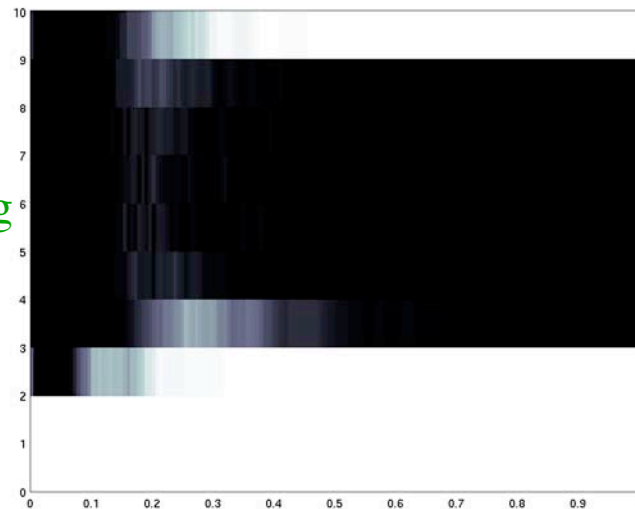
Represent the evolutionary history of the protein group with a spanning set of structures.



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

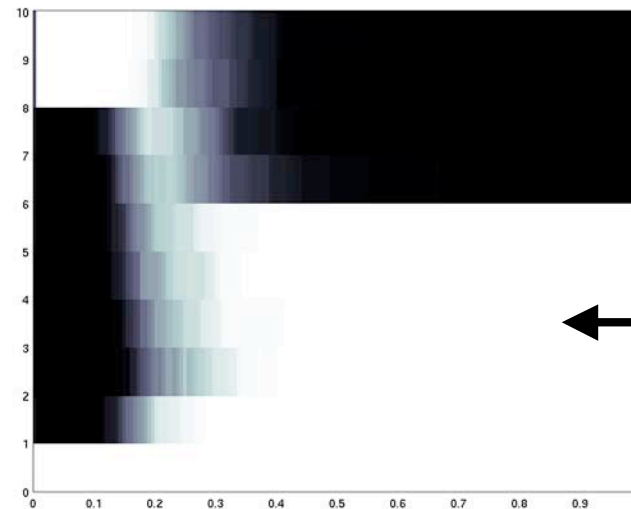
Parameters Define the Measure of Linear Dependence

AARS class I, Rossman fold



γ (normalized)

AARS class II, Novel fold



γ (normalized)

ordering
p-norm

ordering norm

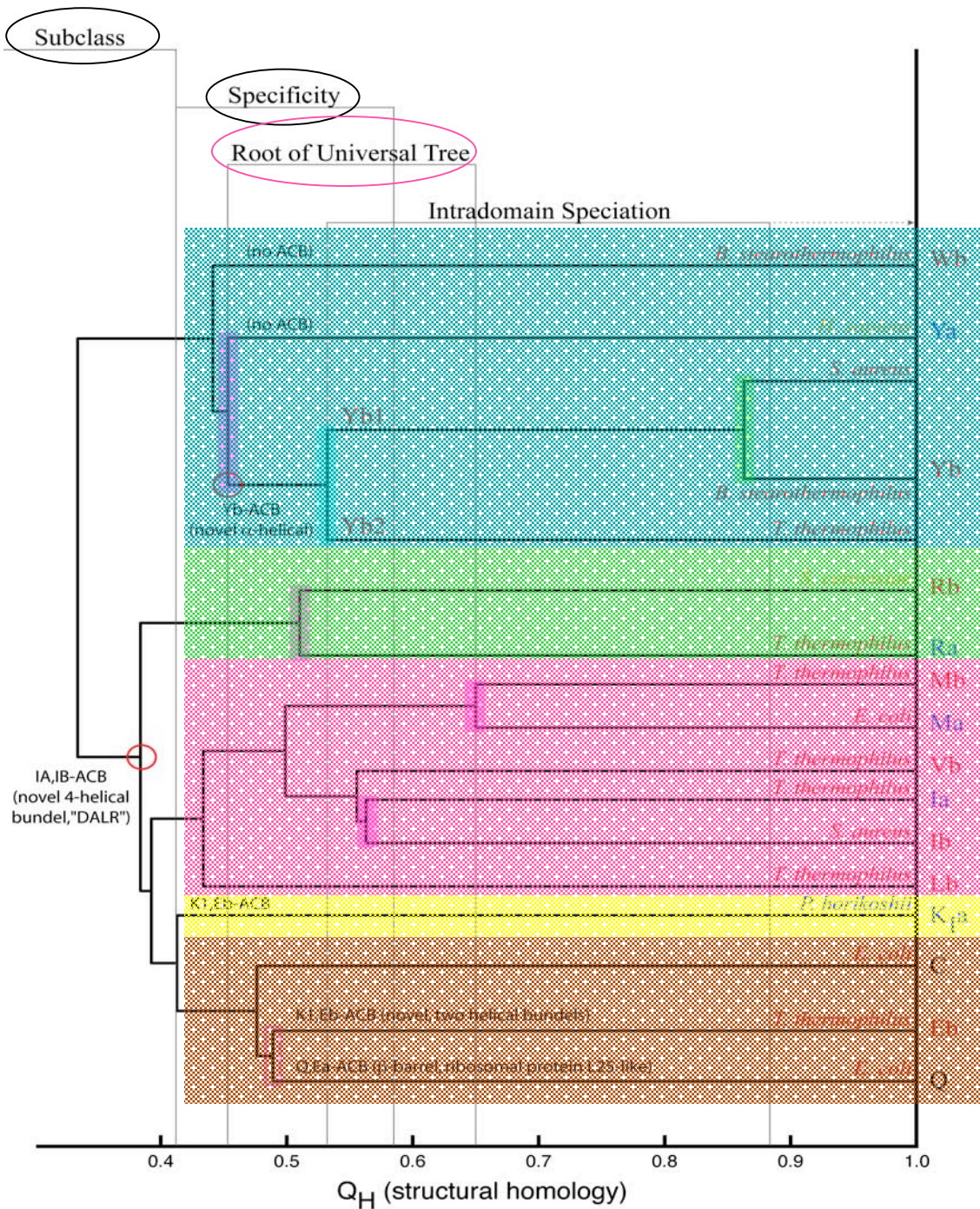
$$\max_{j=k, \dots, n_{\text{proteins}}} (\|a_j\|_{F_p})$$

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

gap scale

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

Class I AARSs evolutionary events

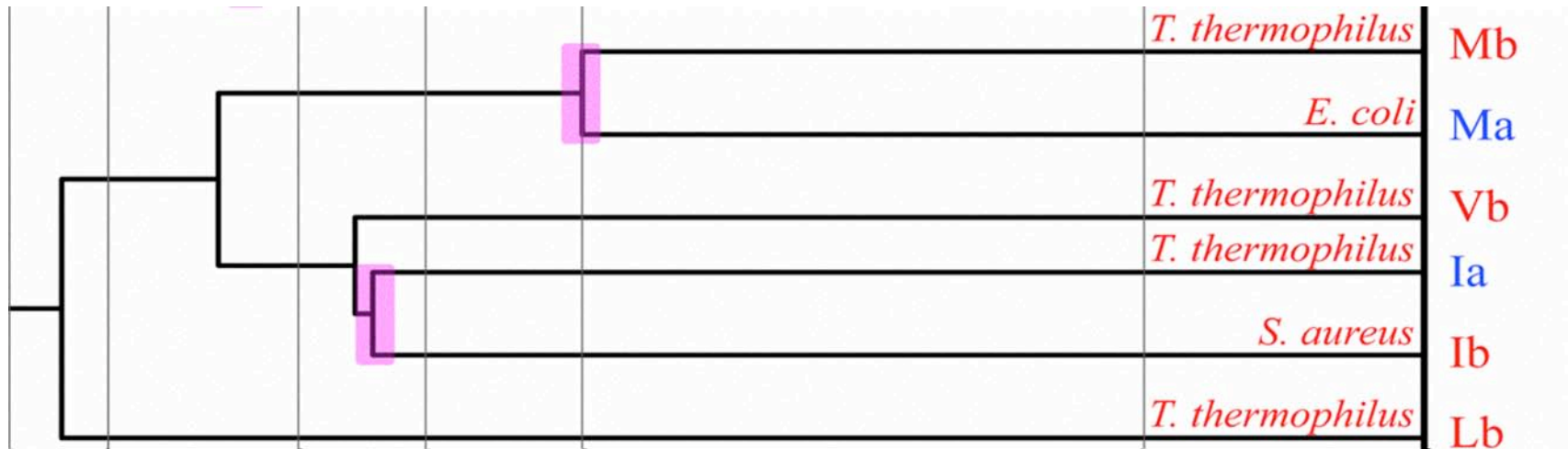


5 Subclasses

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.

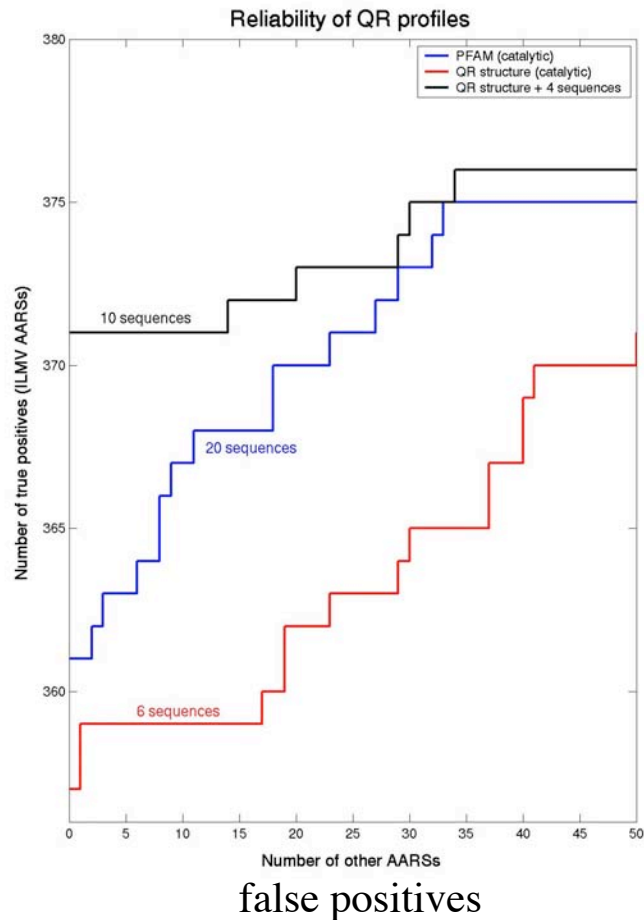
| | Class I | Class II |
|-----------------|----------------------|--|
| Full Canonical | W Y L I E | F H P D |
| Basal Canonical | R M V K _I | T A |
| Non-Canonical | C Q | S G _{α₂} K _{II} N G _{(αβ)₂} |

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

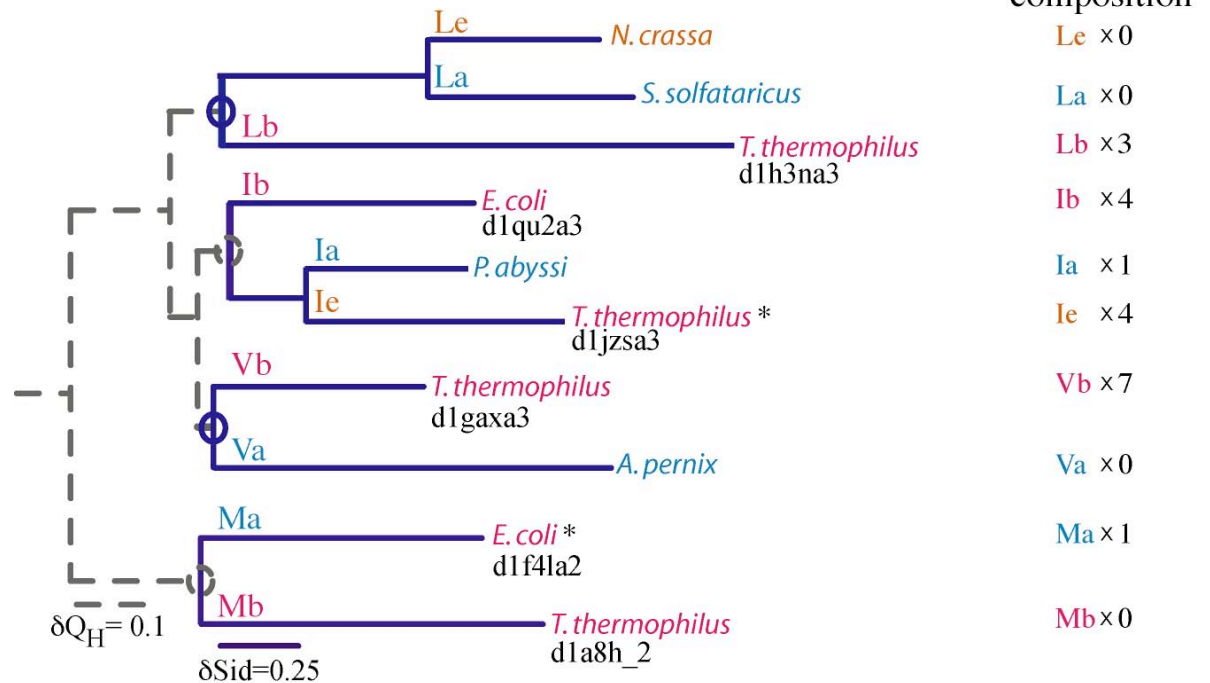
We have 6 structures.

Evolutionary Profiles for Homology Recognition

AARS Subclass ILMV



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

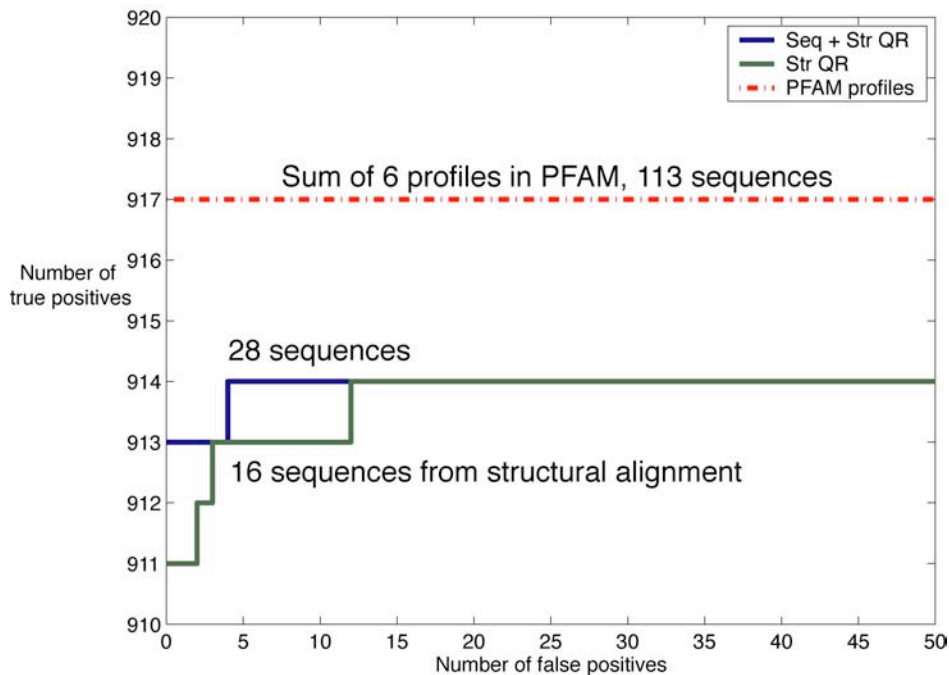


The composition of the profile matters.
 Choosing the right 10 sequence makes all the difference.

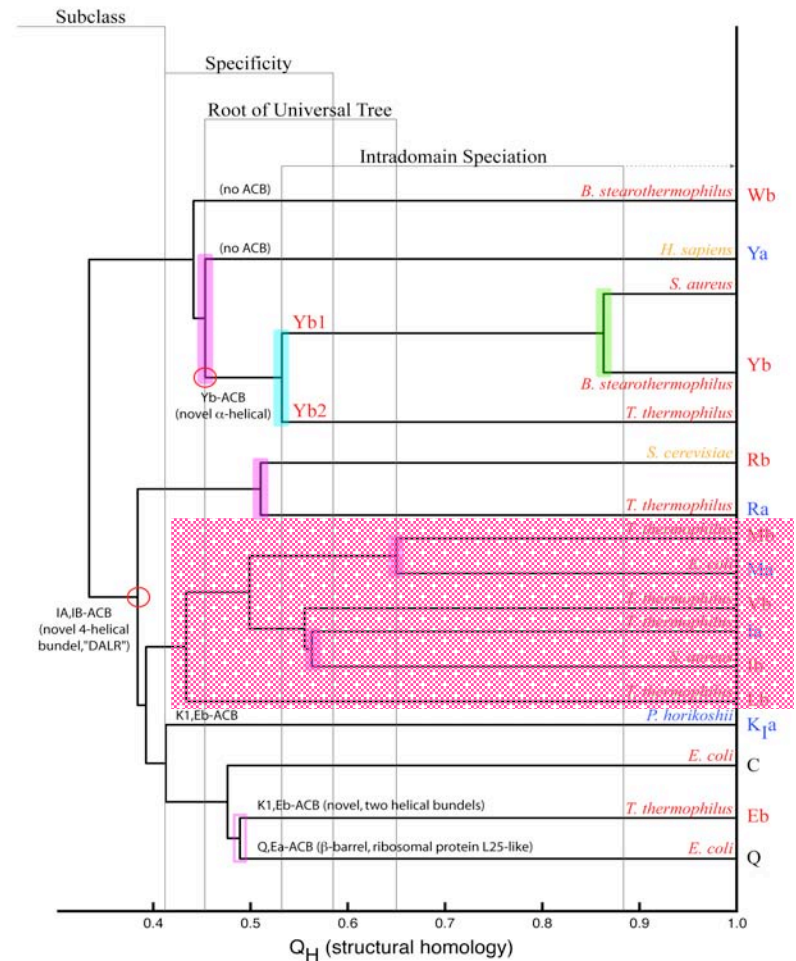
The Economy of Information

How many sequence are needed for profiles?

A single profile
for class I AARSs

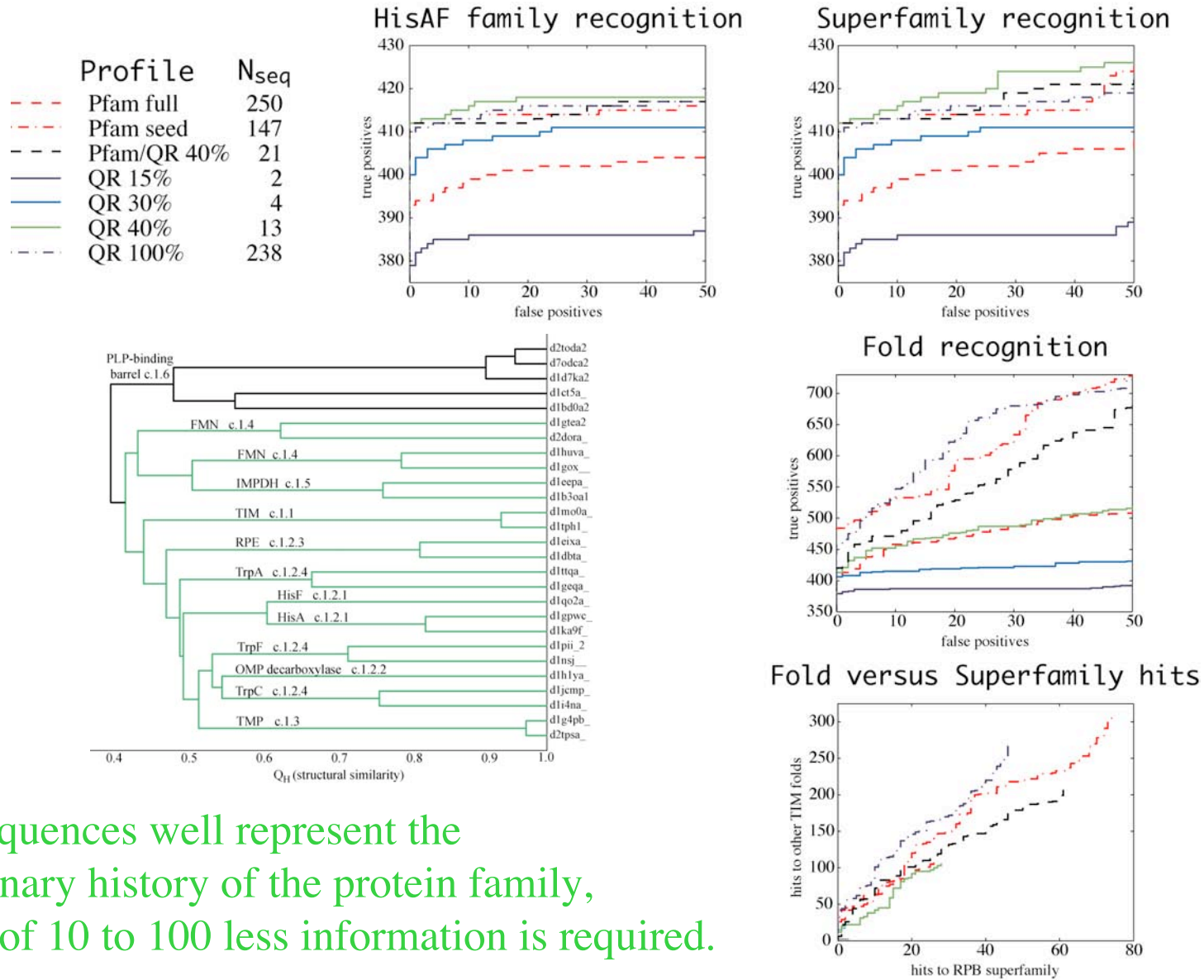


PFAM profile of 113 sequences finds 3 additional sequence fragments compared to the non-redundant profile of 28 sequences.



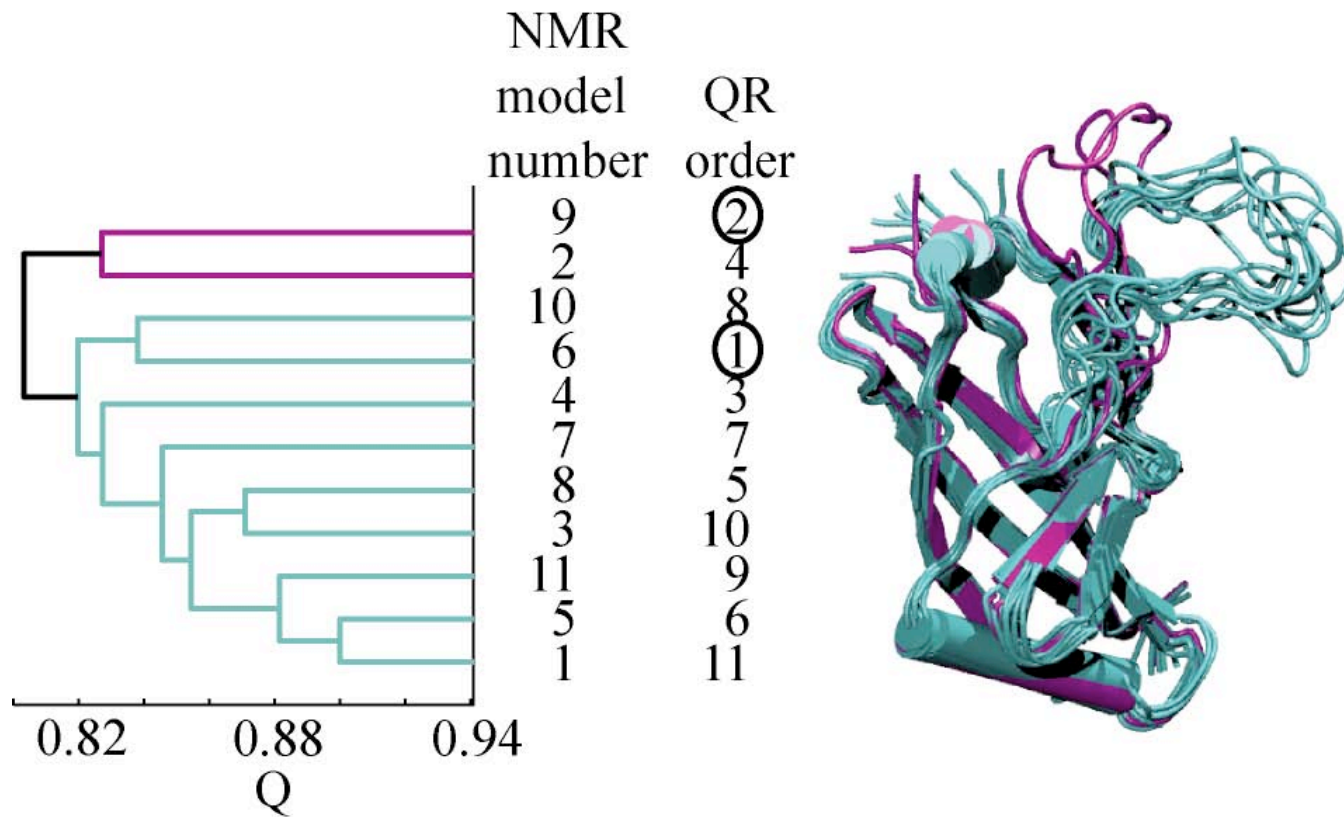
Economy of Information

How many sequences are needed for profiles?



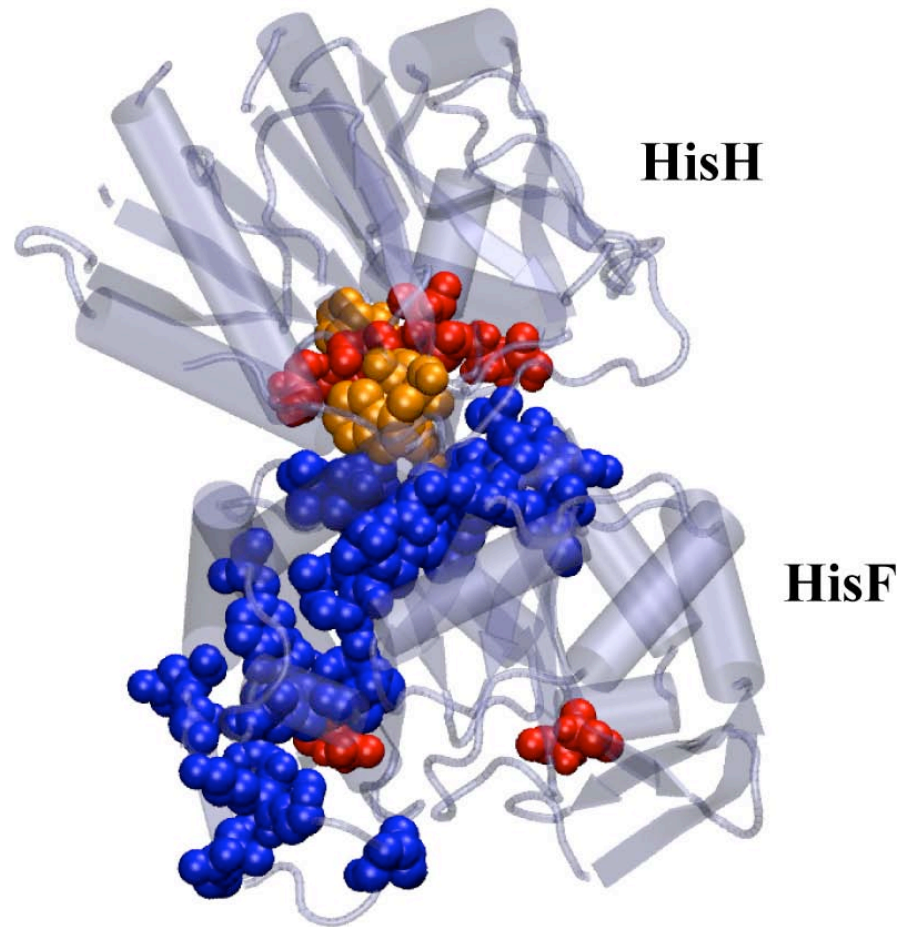
If the sequences well represent the evolutionary history of the protein family, a factor of 10 to 100 less information is required.

QR factorization of an ensemble of NMR structures



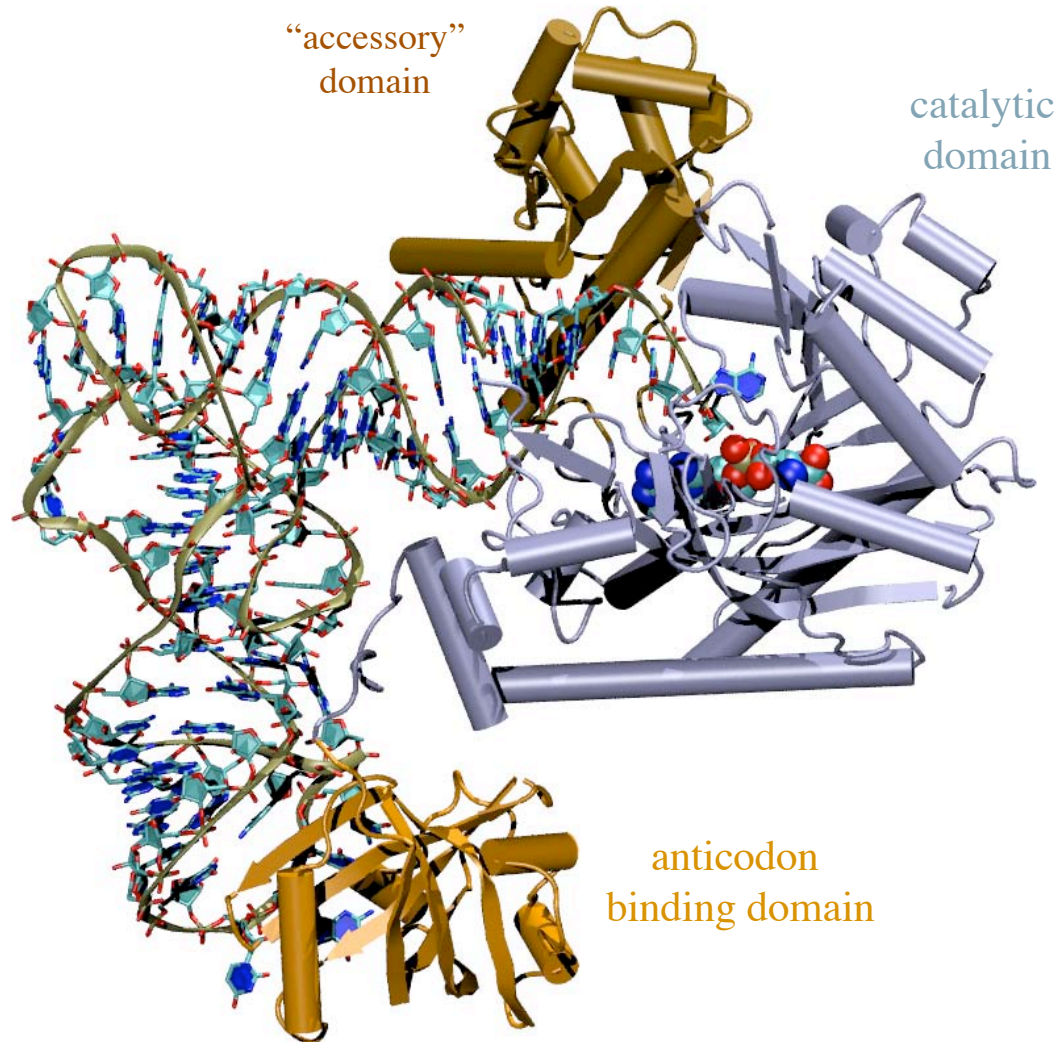
The QR algorithm can be applied to conformational, evolutionary ensembles or both simultaneously.

Evolutionary Structure/Sequence Profiles Suggest Reaction Pathway

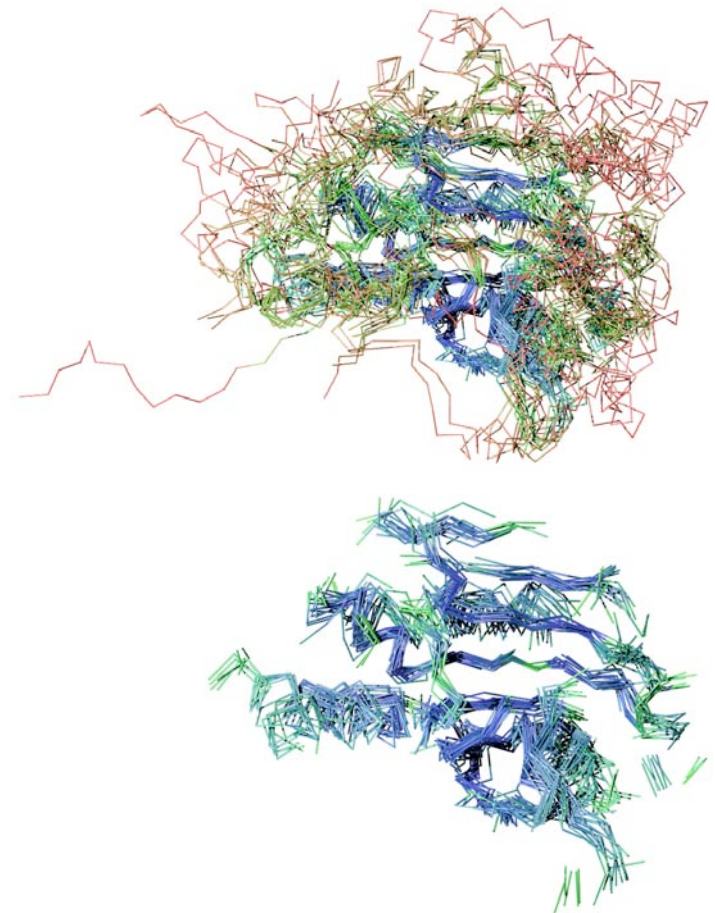


R. Amaro and Z. Schulten, *MD Simulations of Substrate Channeling*, Chemical Physics Special Issue, 2004 (in press). *FE Landscapes of Ammonia Channeling*, PNAS 2003

Domain Structure in AspRS



catalytic domains
class II AARSs

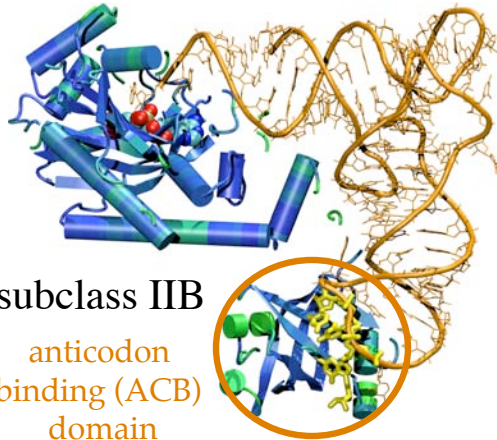


bacterial type aspartyl-tRNA synthetase
E. coli, homodimer

Evolution of Structure and Function in AspRS

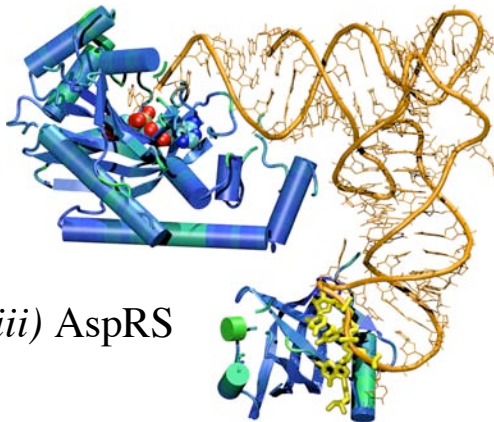


i) class II

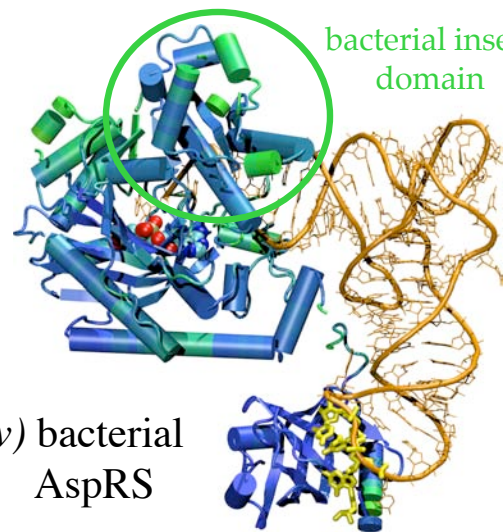


ii) subclass IIB

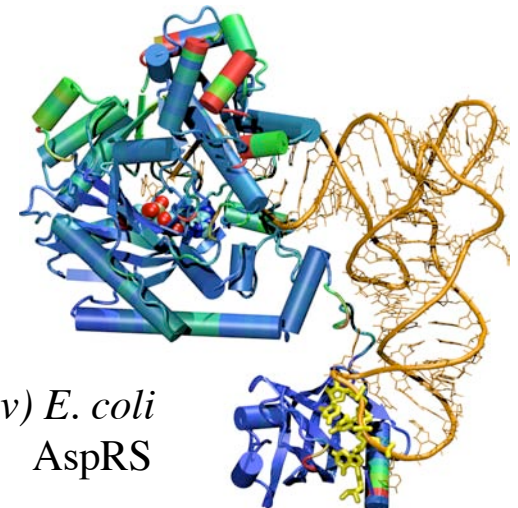
anticodon
binding (ACB)
domain



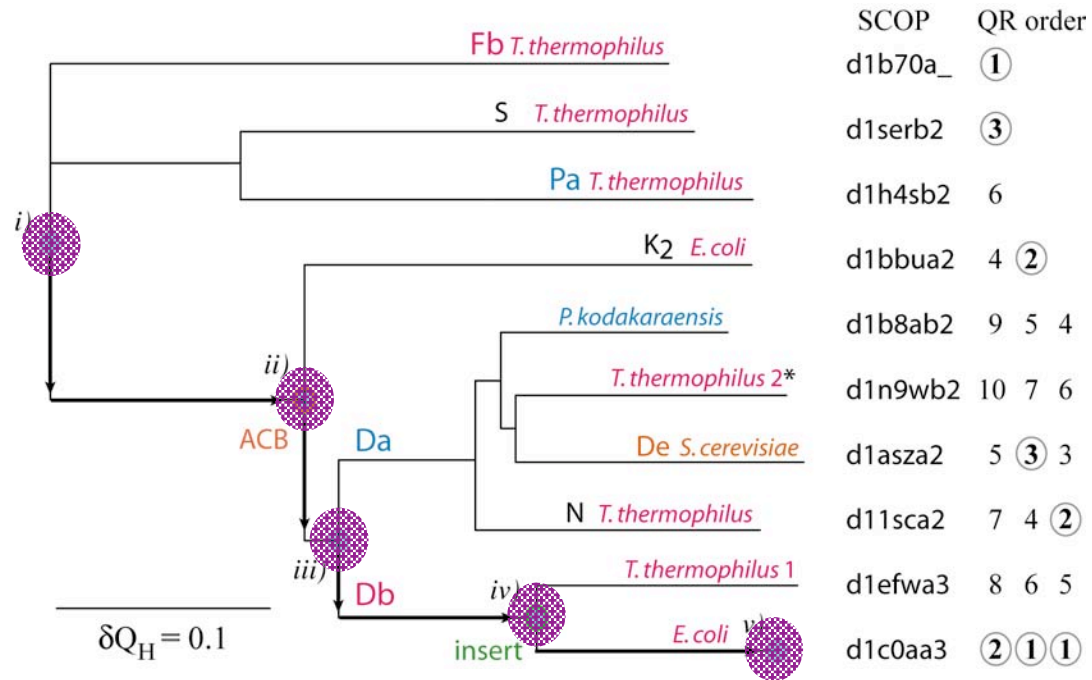
iii) AspRS



iv) bacterial
AspRS



v) *E. coli*
AspRS



$\delta Q_H = 0.1$

bacterial insert
domain

Summary

Evolutionary information is encoded in protein structure.

Protein structure allows investigation of evolutionary events that pre-date the origin of species.

Accounting for gaps is critical for comparing homologous structures.

Sequence and structure can be combined to give a unified phylogenetic framework.

The QR factorization provides evolutionary profiles (EPs).

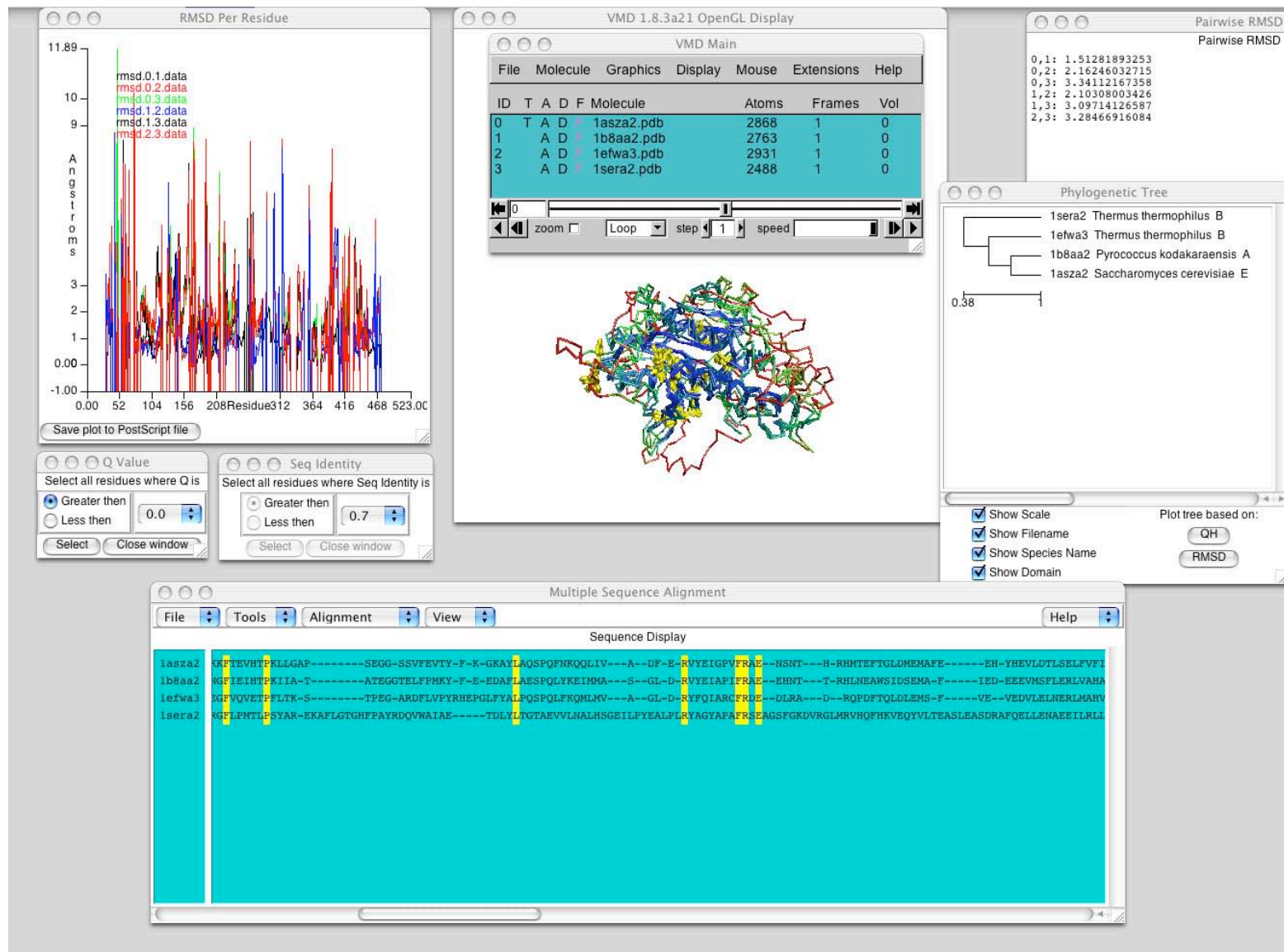
By spanning the evolutionary space with a small number of representative sequences EPs outperform traditional profiles.

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.

The QR algorithm can be applied to conformational, evolutionary ensembles or both simultaneously.

Multiseq in VMD: Merging the sequence and structure worlds



Brijet Dhaliwal, John Eargle, John Stone, Dan Wright

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