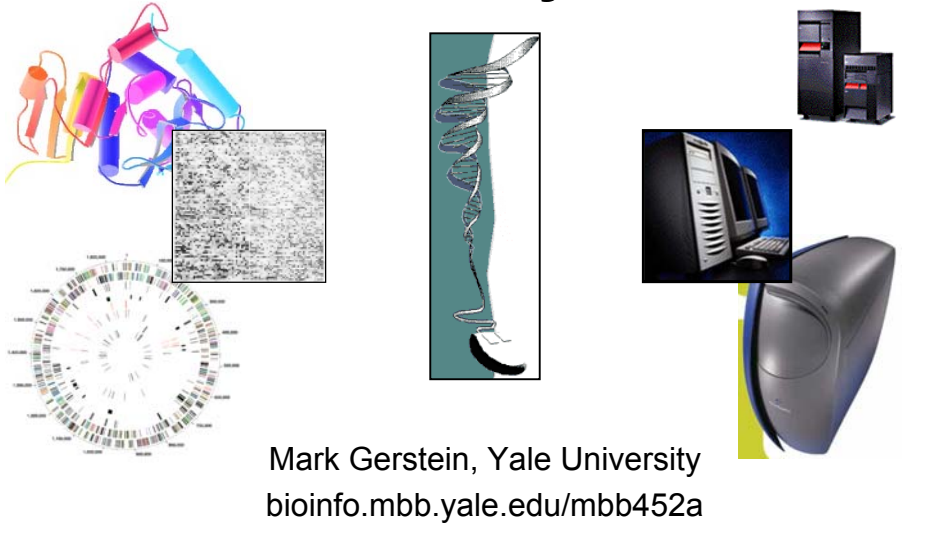


BIOINFORMATICS Surveys

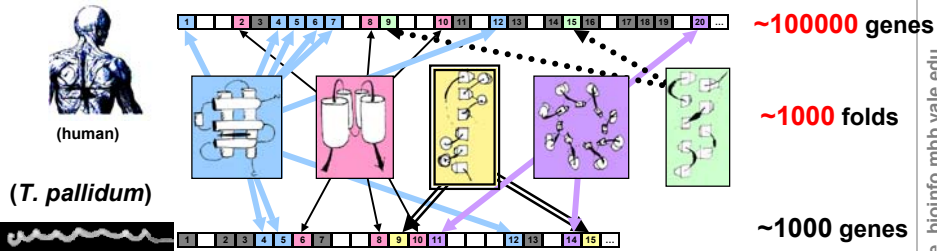


Mark Gerstein, Yale University
bioinfo.mbb.yale.edu/mbb452a

Large-scale Database Surveys (contents)

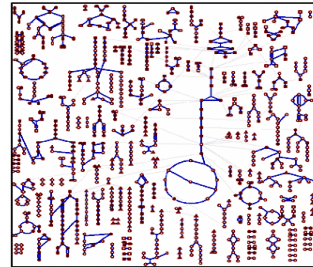
- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Function Classification
- Cross-tabulation, folds and functions
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection

Simplifying the Complexity of Genomes: Global Surveys of a Finite Set of Parts from Many Perspectives



Same logic for sequence families, blocks, orthologs, motifs, pathways, functions....

Functions picture from www.fruitfly.org/~suzi (Ashburner); Pathways picture from ecocyc.pangeasystems.com/ecocyc (Karp, Riley).
Related resources: COGS, ProDom, Pfam, Blocks, Domo, WIT, CATH, Scop....



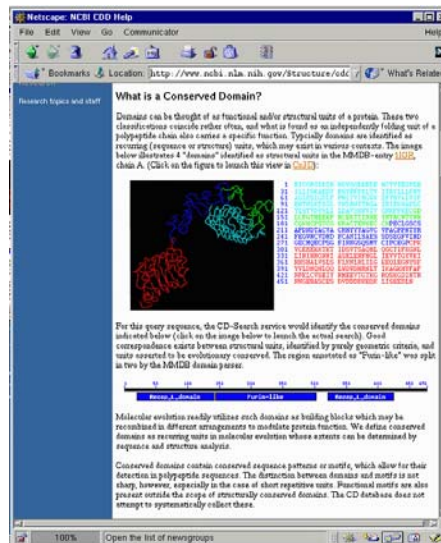
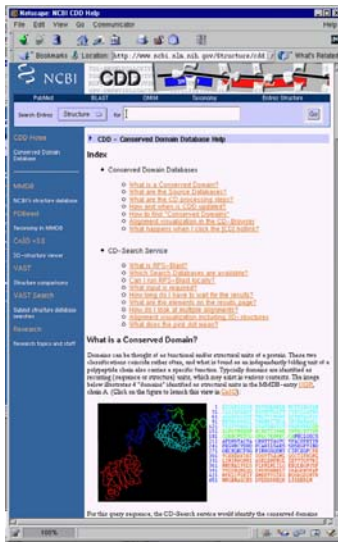
Part = Homolog

Part = Motif



5 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Part = Conserved Domains



6 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Part = Ortholog COGs - Orthologs

Ortholog ~ gene with precise same role in diff. organism, directly related by descent from a common ancestor

vs
Paralog

Ortholog,
homolog,
fold

(Lipman, Koonin, NCBI)

A Natural System of Gene Families from Complete Genomes

Clusters of Orthologous Groups (COGs) were delineated by comparing protein sequences recorded in 8 complete genomes, representing 6 major phylogenetic lineages. Each COG consists of individual proteins or groups of paralogous genes at least 2 domains and thus corresponds to an ancient conserved domain.

Science 1997 Oct 24;278(5318):611-617
 doi:10.1126/science.1255432

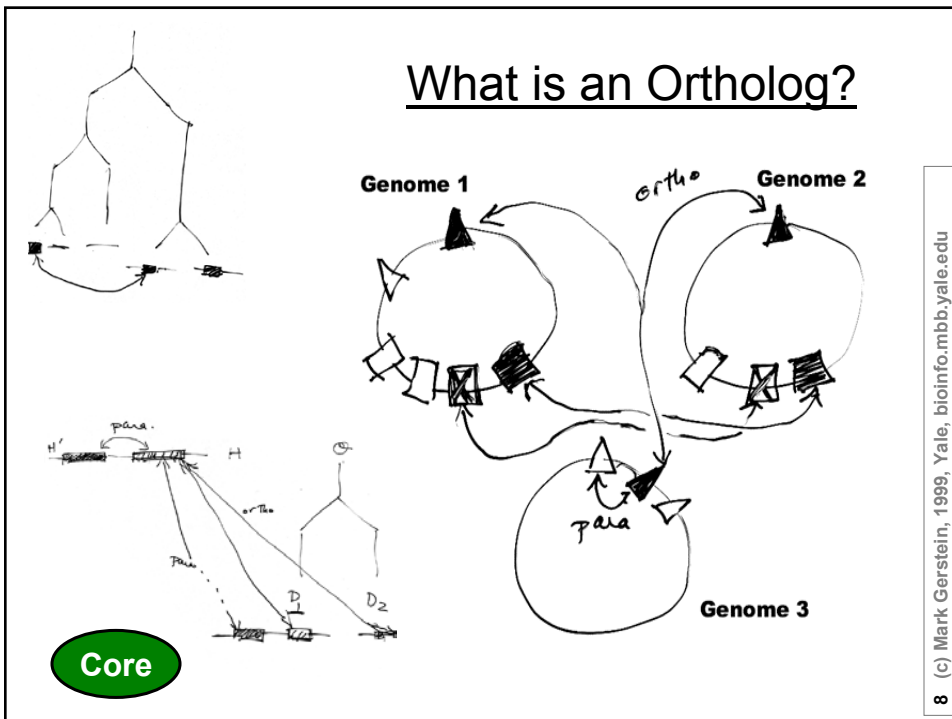
| Color Code | Name | Genome size | Proteins |
|------------|---|--------------|----------|
| E | <i>Escherichia coli</i> | 4,453,231 bp | 4283 |
| H | <i>Haemophilus influenzae</i> | 1,830,260 bp | 1703 |
| U | <i>Mycobacterium parvum</i> | 1,667,867 bp | 1564 |
| G | <i>Mycobacterium genitalium</i> | 580,073 bp | 468 |
| F | <i>Mycobacterium fortuitum</i> | 816,384 bp | 677 |
| C | <i>Cryptosporidium parvum</i> | 3,573,470 bp | 3168 |
| M | <i>Mollicoccus janssonii</i> | 1,779,534 bp | 1736 |
| Y | Yeast - <i>Saccharomyces cerevisiae</i> | 12,060 Kbp | 5932 |

Functional annotation

| Code | COG(s) | Domain(s) | Description |
|------|--------|-----------|---|
| J | 288 | 943 | Information storage and processing |
| K | 18 | 205 | Transcription |
| L | 64 | 676 | Translation, repair, recombination |
| | | | Cellular processes |
| O | 32 | 405 | Motility |
| M | 47 | 469 | Other membrane, cell wall, lipoproteins |
| N | 38 | 127 | Structure and biophysics |
| P | 30 | 225 | Transport, catabolism and metabolism |
| | | | Metabolism |
| C | 77 | 711 | Energy production and conversion |
| G | 33 | 301 | Carbohydrate metabolism and transport |
| E | 87 | 3742 | Amino acid metabolism and transport |
| F | 87 | 381 | Nucleotide metabolism and transport |
| H | 45 | 412 | Cofactor metabolism |
| I | 33 | 176 | Lipid metabolism |

7 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

What is an Ortholog?



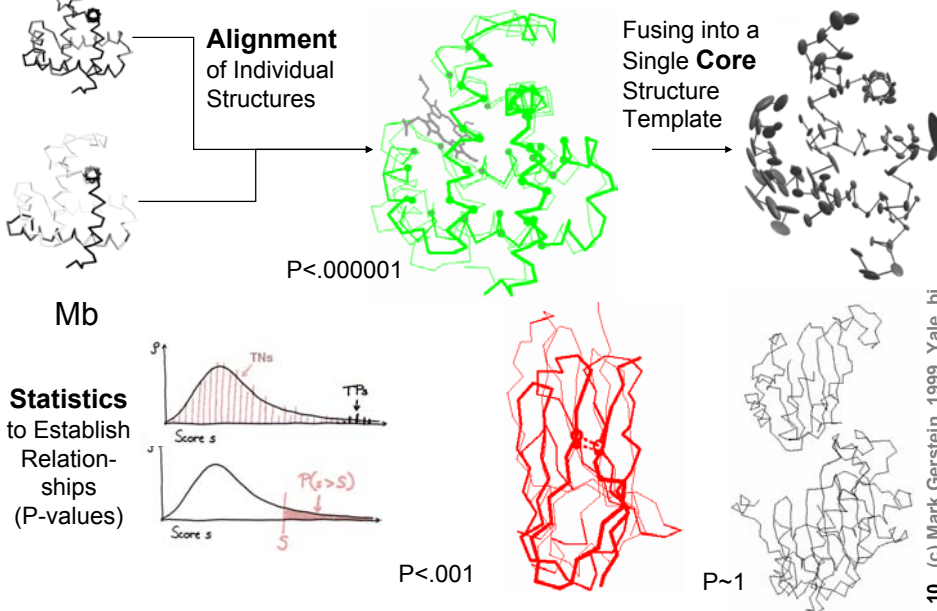
8 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Function Classification
- Cross-tabulation, folds and functions
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection

9 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

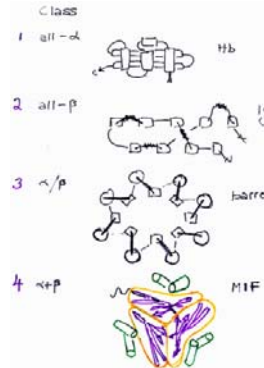
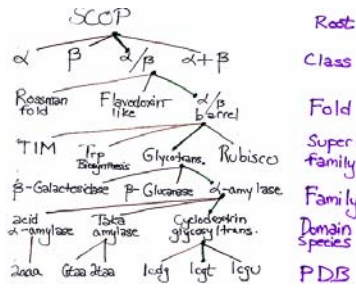
Hb The Parts List: A Library of Known Folds



10 (c) Mark Gerstein, 1999, Yale, bi

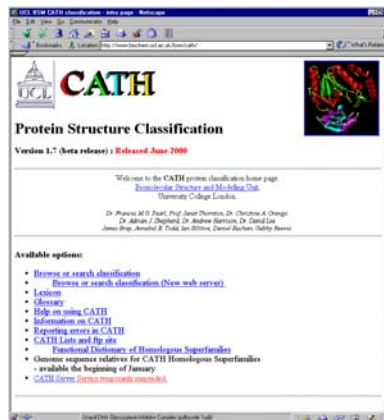
Fold Classifications

- Scop
 - ◊ Chothia, Murzin (Cambridge)
 - ◊ Manual classification, auto-alignments available
 - ◊ Evolutionary clusters
- Cath
 - ◊ Thornton (London)
 - ◊ semi-automatic classification with alignments
 - ◊ class, arch, topo., homol.
- FSSP
 - ◊ Sander, Holm (Cambridge)
 - ◊ totally automatic with DALI
 - ◊ objective but not always interpretable clusters
- VAST



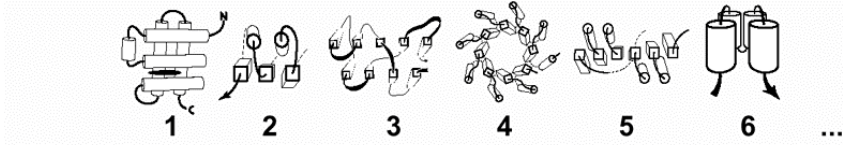
11 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Part = Fold



Fold Library vs. Other Fundamental Data structures

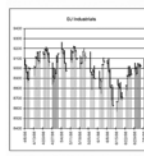
Parts List **Database: Statistical**, rather than mathematical relationships and conclusions



Folds in Molecular Biology **1000-10000**

| Element | Symbol | Atomic Weight | Boiling Point |
|---------|--------|---------------|---------------|
| H | H | 1.00 | 8°C |
| F | F | 18.99 | -4°C/mol |
| Cl | Cl | 35.45 | -12°C/mol |
| Br | Br | 79.90 | -108°C/mol |
| I | I | 126.90 | -184°C/mol |
| Na | Na | 22.99 | 883°C/mol |
| K | K | 39.09 | 770°C/mol |
| Ca | Ca | 40.08 | 1485°C/mol |
| Mg | Mg | 24.31 | 900°C/mol |
| Al | Al | 26.98 | 2450°C/mol |
| Si | Si | 28.09 | 2355°C/mol |
| P | P | 30.97 | 2537°C/mol |
| S | S | 32.06 | 2445°C/mol |
| Fe | Fe | 55.85 | 2750°C/mol |
| Cu | Cu | 63.55 | 2562°C/mol |
| Zn | Zn | 65.38 | 2537°C/mol |
| Ni | Ni | 58.71 | 2730°C/mol |
| Co | Co | 58.93 | 2700°C/mol |
| N | N | 14.01 | -196°C/mol |
| O | O | 15.99 | -183°C/mol |
| Ne | Ne | 20.18 | -196°C/mol |

| Element | Symbol | Atomic Weight | Boiling Point |
|--------------|--------|---------------|---------------|
| H | H | 1.00 | 8°C |
| He | He | 4.00 | -269°C/mol |
| Li | Li | 6.94 | 1330°C/mol |
| Be | Be | 9.01 | 2970°C/mol |
| B | B | 10.81 | 2550°C/mol |
| C | C | 12.01 | 3642°C/mol |
| N | N | 14.01 | -196°C/mol |
| O | O | 15.99 | -183°C/mol |
| F | F | 18.99 | -188°C/mol |
| Ne | Ne | 20.18 | -196°C/mol |
| Na | Na | 22.99 | 883°C/mol |
| Mg | Mg | 24.31 | 900°C/mol |
| Al | Al | 26.98 | 2450°C/mol |
| Si | Si | 28.09 | 2355°C/mol |
| P | P | 30.97 | 2537°C/mol |
| S | S | 32.06 | 2445°C/mol |
| Cl | Cl | 35.45 | -34°C/mol |
| Ar | Ar | 39.95 | -186°C/mol |
| K | K | 39.09 | 770°C/mol |
| Ca | Ca | 40.08 | 1485°C/mol |
| Sc | Sc | 44.96 | 2835°C/mol |
| Ti | Ti | 47.88 | 3287°C/mol |
| V | V | 50.94 | 3407°C/mol |
| Cr | Cr | 51.99 | 2672°C/mol |
| Mn | Mn | 54.94 | 2071°C/mol |
| Fe | Fe | 55.85 | 2750°C/mol |
| Co | Co | 58.93 | 2700°C/mol |
| Ni | Ni | 58.71 | 2730°C/mol |
| Cu | Cu | 63.55 | 2562°C/mol |
| Zn | Zn | 65.38 | 2537°C/mol |
| Ga | Ga | 69.72 | 2400°C/mol |
| Ge | Ge | 72.64 | 2833°C/mol |
| As | As | 74.92 | 2110°C/mol |
| Se | Se | 78.96 | 2179°C/mol |
| Br | Br | 79.90 | -108°C/mol |
| Kr | Kr | 83.80 | -153°C/mol |
| Rb | Rb | 85.47 | 390°C/mol |
| Sr | Sr | 87.62 | 1362°C/mol |
| Y | Y | 88.91 | 2790°C/mol |
| Zr | Zr | 91.22 | 3552°C/mol |
| Nb | Nb | 92.91 | 2471°C/mol |
| Mo | Mo | 95.94 | 2623°C/mol |
| Tc | Tc | 98.91 | 2537°C/mol |
| Ru | Ru | 101.07 | 2631°C/mol |
| Rh | Rh | 102.91 | 2673°C/mol |
| Pd | Pd | 106.42 | 2698°C/mol |
| Ag | Ag | 107.87 | 2162°C/mol |
| Cd | Cd | 112.41 | 2042°C/mol |
| In | In | 114.82 | 2019°C/mol |
| Sn | Sn | 118.71 | 2270°C/mol |
| Sb | Sb | 121.76 | 1780°C/mol |
| Te | Te | 127.60 | 1769°C/mol |
| I | I | 126.90 | -184°C/mol |
| Xe | Xe | 131.29 | -108°C/mol |
| Ba | Ba | 137.33 | 1290°C/mol |
| La | La | 138.91 | 3233°C/mol |
| Ce | Ce | 140.12 | 3442°C/mol |
| Pr | Pr | 140.91 | 3273°C/mol |
| Nd | Nd | 144.24 | 3288°C/mol |
| Pm | Pm | 144.91 | 3273°C/mol |
| Sm | Sm | 150.36 | 3273°C/mol |
| Eu | Eu | 151.96 | 3243°C/mol |
| Gd | Gd | 157.25 | 3273°C/mol |
| Tb | Tb | 158.93 | 3273°C/mol |
| Dy | Dy | 162.50 | 3273°C/mol |
| Ho | Ho | 164.93 | 3273°C/mol |
| Er | Er | 167.26 | 3273°C/mol |
| Tm | Tm | 168.93 | 3273°C/mol |
| Yb | Yb | 173.05 | 3273°C/mol |
| Lu | Lu | 174.97 | 3273°C/mol |
| Hf | Hf | 178.49 | 3537°C/mol |
| Ta | Ta | 180.95 | 3290°C/mol |
| W | W | 183.84 | 3422°C/mol |
| Re | Re | 186.21 | 3180°C/mol |
| Os | Os | 190.23 | 3047°C/mol |
| Ir | Ir | 192.22 | 2709°C/mol |
| Pt | Pt | 195.08 | 2637°C/mol |
| Au | Au | 196.97 | 2537°C/mol |
| Hg | Hg | 200.59 | 2019°C/mol |
| Tl | Tl | 204.38 | 1723°C/mol |
| Pb | Pb | 207.2 | 1749°C/mol |
| Bi | Bi | 208.98 | 1564°C/mol |
| Po | Po | 209 | 1564°C/mol |
| At | At | 210 | 1564°C/mol |
| Rn | Rn | 222 | 1564°C/mol |
| Fr | Fr | 223 | 1564°C/mol |
| Ra | Ra | 226 | 1564°C/mol |
| Ac | Ac | 227 | 1564°C/mol |
| Th | Th | 232.04 | 3273°C/mol |
| Pa | Pa | 231.04 | 3273°C/mol |
| U | U | 238.03 | 3273°C/mol |
| Np | Np | 237.05 | 3273°C/mol |
| Pu | Pu | 244.06 | 3273°C/mol |
| Am | Am | 243.06 | 3273°C/mol |
| Cm | Cm | 247.07 | 3273°C/mol |
| Bk | Bk | 247.07 | 3273°C/mol |
| Cf | Cf | 251.08 | 3273°C/mol |
| Es | Es | 252.08 | 3273°C/mol |
| Fm | Fm | 257.09 | 3273°C/mol |
| Mendelevium | Md | 258.10 | 3273°C/mol |
| Nobelium | Nb | 262.11 | 3273°C/mol |
| Lanthanum | La | 138.91 | 3233°C/mol |
| Cerium | Ce | 140.12 | 3442°C/mol |
| Praseodymium | Pr | 140.91 | 3273°C/mol |
| Neodymium | Nd | 144.24 | 3288°C/mol |
| Europium | Eu | 151.96 | 3243°C/mol |
| Gadolinium | Gd | 157.25 | 3273°C/mol |
| Terbium | Tb | 158.93 | 3273°C/mol |
| Dysprosium | Dy | 162.50 | 3273°C/mol |
| Ytterbium | Yb | 173.05 | 3273°C/mol |
| Lutetium | Lu | 174.97 | 3273°C/mol |



10
Physics

100
Chemistry

1000
-10000
Finance

>1000000
Politics

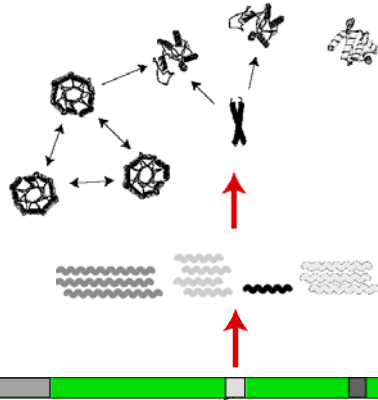
(Larger than physics and chemistry, similar to Finance (Exact Finite Number of Objects (3,056 on NYSE by 1/98), description by Standardized Statistics (even abbrevs, INTC) and groups (sectors)) Smaller than Social Surveys, Indefinite Number of People, Not Well Defined Vocabulary and statistics.

The Next Step: Post-genomic Challenges

↑ #1: Understanding Protein Function on a Genomic Scale

Large-scale Biochemistry:
Expression, Structural Genomics, Protein Interactions

Initial Step: genome sequence & genes



Extra

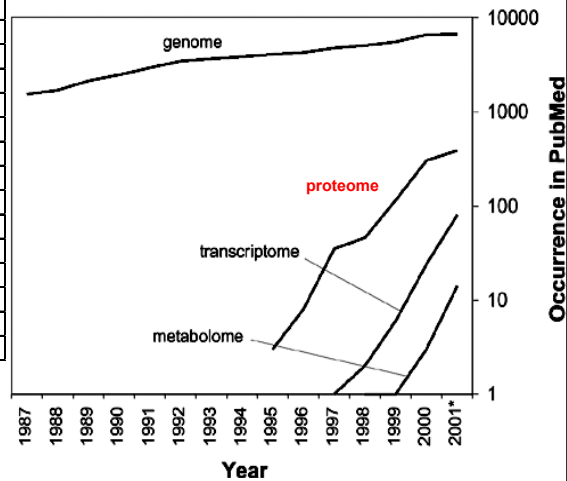
■ #2: Understanding the Meaning of Intergenic Regions

Evolutionary Implications as a Graveyard: Pseudogenes, Regulatory Regions, Repeats.

| Term | Google Hits | PubMed Hits | 1st PubMed Hit Year |
|-----------------|-------------|-------------|---------------------|
| Genome | ~1880000 | 66171 | 1932 ** |
| Proteome | ~63,000 | 703 | 1995 |
| Transcriptome | 3520 | 72 | 1997 |
| Physiome | 2980 | 15 | 1997 |
| Metabolome | 349 | 12 | 1998 |
| Phenome | 4980 | 6 | 1995 |
| Morphome | 238 | 2 | 1996 |
| Interactome | 56 | 2 | 1999 |
| Glycome | 46 | 1 | 2000 |
| Secretome | 21 | 1 | 2000 |
| Ribonome | 1 | 1 | 2000 |
| Orfeome | 42 | - | - |
| Regulome | 18 | - | - |
| Cellome | 17 | - | - |
| Operome | 8 | - | - |
| Transportome | 1 | - | - |
| Functome | 1 | - | - |

An 'Omic Language to Describe the Next Steps

Extra



proteomics



Integrating heterogeneous 'omic information through proteins: families, folds, locations, functions, interactions, pseudogenes...

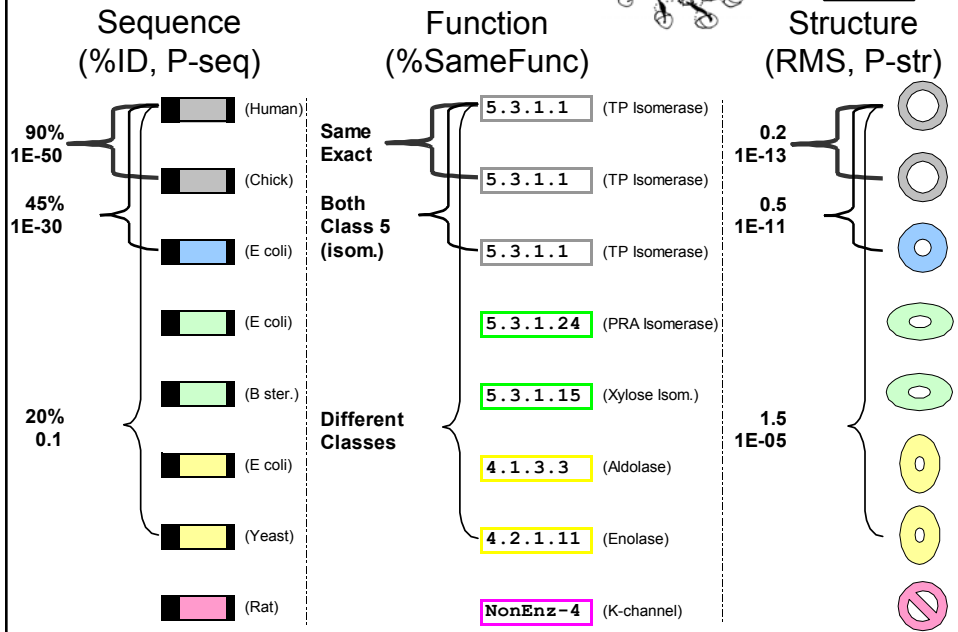
Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Function Classification
- Cross-tabulation, folds and functions
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection

Annotation Transfer: TIM ex.



Extra



Chothia & Lesk, 1986 -- 32 points

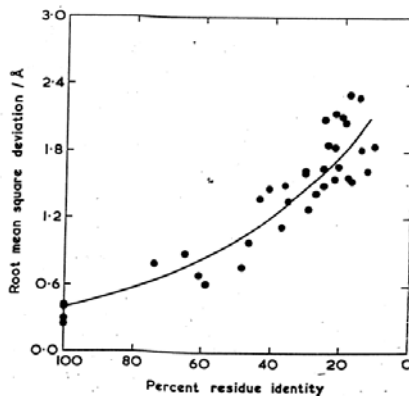


Fig. 2. The relation of residue identity and the r.m.s. deviation of the backbone atoms of the common cores of 32 pairs of homologous proteins (see Table II).

EMBO J 4: 823 (1986)

“The relation between the divergence of sequence and structure in proteins”

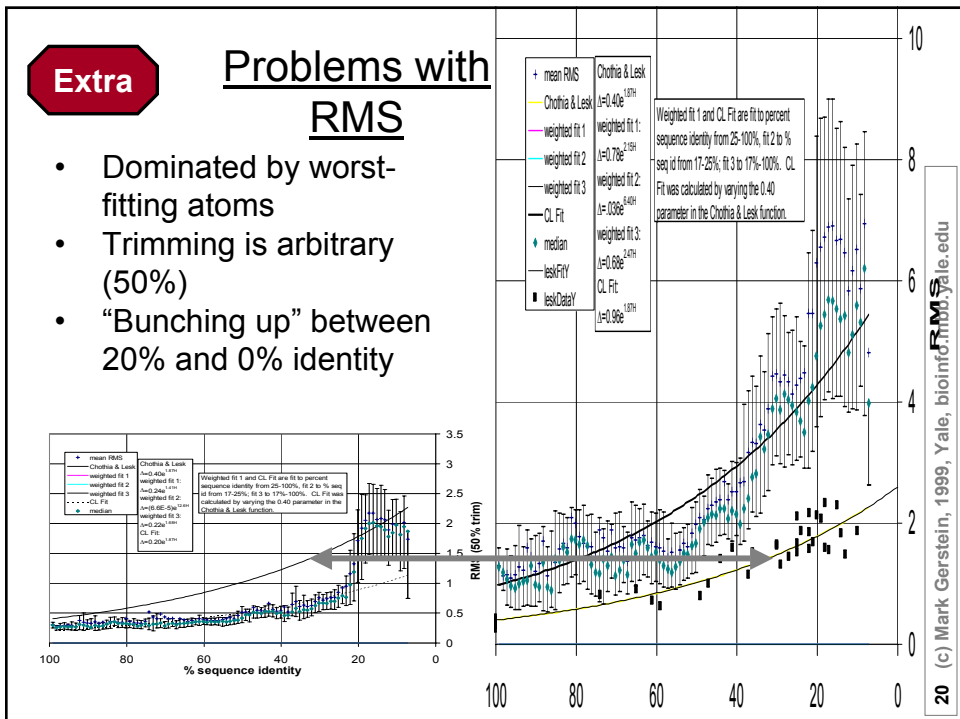
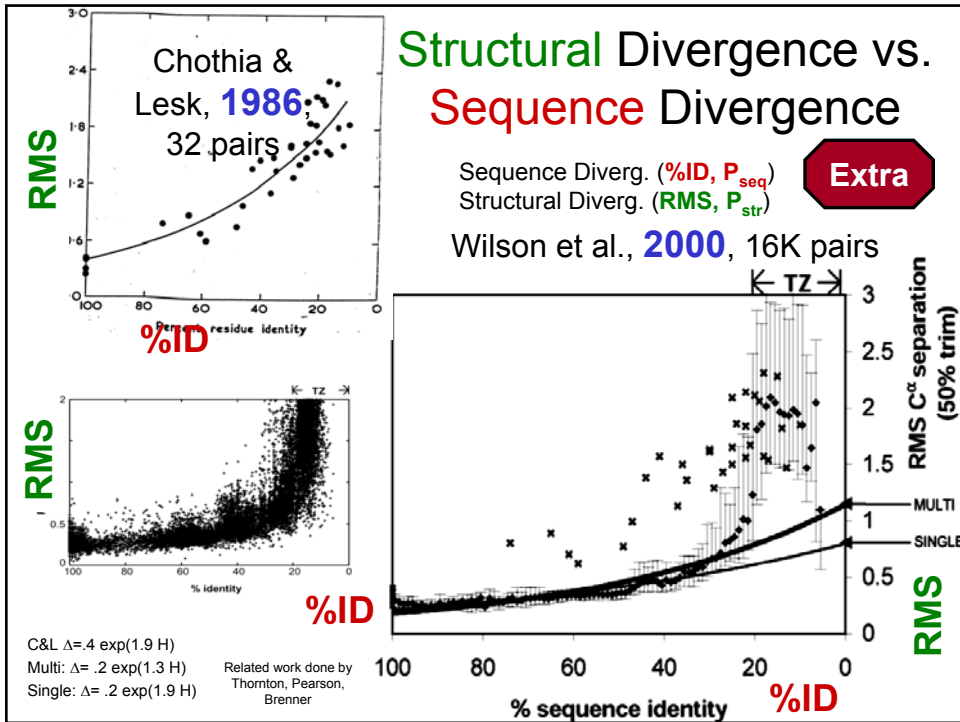
32 pairs of homologous proteins

RMS, percent identity

$$\Delta = 0.40 e^{1.87H}$$

Now redo with >16,000 pairs in scop + auto-alignments (pdb95d)....

Extra



Problems with RMS and %ID

- Difference not similarity, NO EVD fit
- Dominated by worst-fitting atoms, easily skewed
- Trimming is arbitrary (50%)

$$S_{str} = \sum \frac{100}{5 + d_i^2}$$

$$RMS = \sqrt{\sum d_i^2}$$

%ID problem:
"Bunching up"
between 20%
and 0%
identity

21 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Structural Comp. Score vs. Smith-Waterman Score

overcomes zero bunching, trimming problem

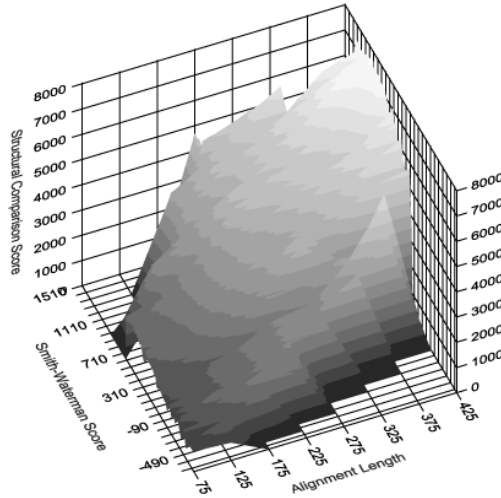
$S_{str} = 100(21 - 11 \exp(-0.0054 S_{seq}))$

Extra

22 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Problems with Structural Alignment Score

Different Lengths give different scores.
 Scores follow equation of the form:
 $y = An + Mx + B$

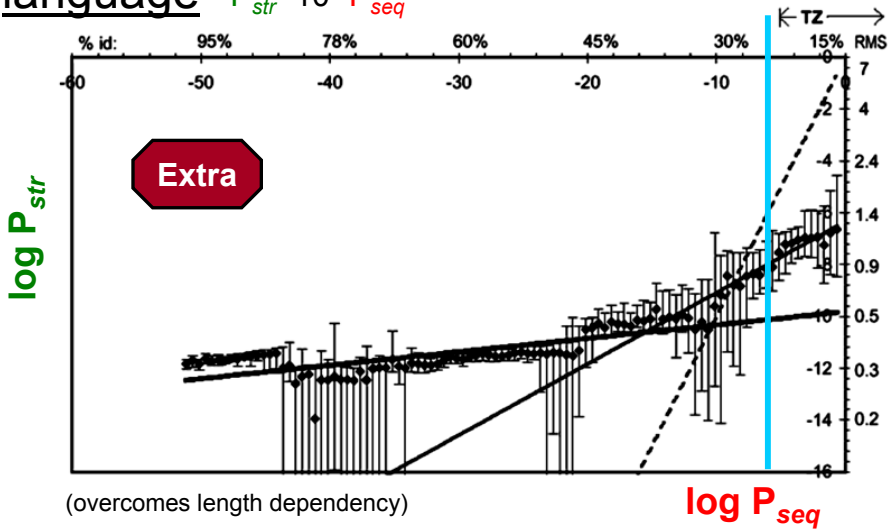
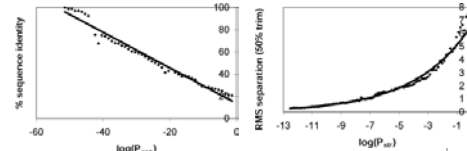


Extra

Modern statistical language

Not in TZ
 $P_{str} = 10^{-10} P_{seq}^{.05}$

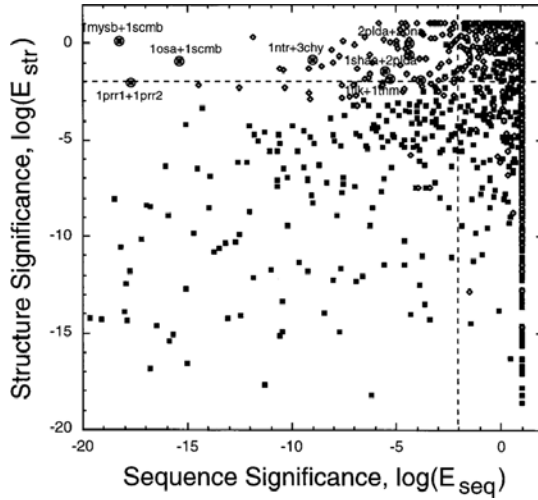
in TZ
 $P_{str} = 10^{-6} P_{seq}^{.274}$



(overcomes length dependency)

Focus on Twilight Zone

- Sequence Sig. without structure signif.
 - ◇ Protein motions
 - ◇ small proteins
 - ◇ low-res, NMR
- Struc. Sig. without Seq. signif.
 - ◇ More in bottom-right than top-left



Extra

Relationship of Similarity in Sequence & Structure - Summary

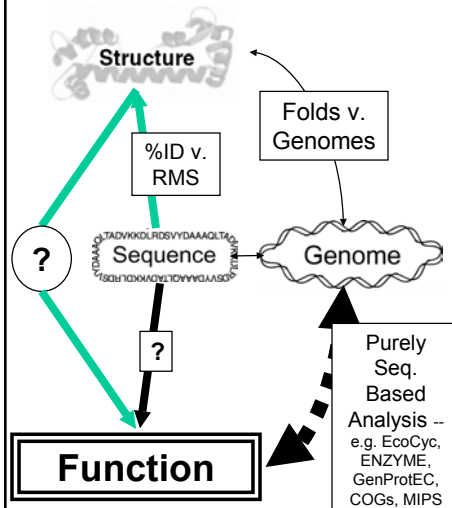
Extra

| | Sequence Similarity | Structural Similarity | Features | Limitations |
|-----------------------------|---------------------------|-------------------------------|---|--|
| Traditional Scores | Percent sequence identity | RMS C ^α separation | Well understood, in use | RMS depends most highly on worst matches, requiring arbitrary trimming |
| Alignment Similarity Scores | S _{seq} | S _{str} | Analogous similarity scores, S _{str} depends most highly on best matches | Dependence on alignment length |
| Modern Probabilistic Scores | P _{seq} | P _{str} | Statistical significance, unified framework for different comparisons | Not as familiar as RMS and percent identity, some residual length-dependency |

Large-scale Database Surveys (contents)

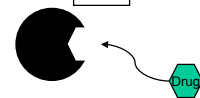
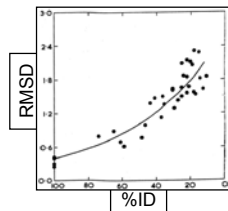
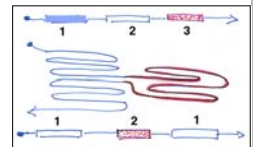
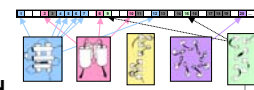
- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Function Classification
- Cross-tabulation, folds and functions
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection

Adding Structure to Functional Genomics, Function to Structural Genomics



Why Structure? Do we really need it?

- 1 Most Highly Conserved
- 2 Precisely Defined Modules
- 3 Seq. \leftrightarrow Struc. Clearer than Seq. \leftrightarrow Func.
- 4 Link to Chemistry, Drugs



Functional Classification

COGs
(*cross-org.*,
just *conserved*,
NCBI
Koonin/Lipman)

GenProtEC
(*E. coli*, Riley)

ENZYME
(SwissProt
Bairoch/
Apweiler,
just *enzymes*,
cross-org.)

“Fly”
(*fly*, Ashburner)
now extended to
GO (*cross-org.*)

MIPS/PEDANT
(*yeast*, Mewes)

Also:
Other
SwissProt
Annotation
WIT, KEGG
(just *pathways*)
TIGR EGAD
(*human ESTs*)

GenProtEC - Functional Classification

the *E. coli* database
<http://genprotec.mbl.edu/start>

GenProtEC *E. coli* genome and proteome database

Back to mainpage

Schema of GenProtEC

This database can be divided into three logical layers.

- QUERY LAYER**: Search the database by query or categories.
- RESULTS LAYER**: List genes within the search criteria.
- DETAIL LAYER**: Display all available information about that gene.

Local data ↔ Remote data

- Browsing/Query layer** - Data used to categorize genes by gene type or physiological role
- Results layer** Gene Selection - gene and SWISSPROT names and synonyms
- Details Layer** - Specific information about a gene including sequence similarity, groupings, and literature. Data is available from local data sources as GenProtEC as well as to other remote data sources where URL's can be created with embedded data.

This system consists of tables originating in FoxPro 2.6 and served to the web with Expertelligence's Webase. Webase uses an ODBC data source and embedded SQL statements to communicate with the tables that make up the GenProtEC system. SQL and Webase expressions are embedded in template files which also contain the HTML markup used to display the data. An overview of the tables and their relationships is shown below.

Overview of the source tables and relationships from GenProtEC

| SYN_GENE | IDENTY | MODULE | SWISSPROT | GENE | SWISS |
|----------|--------|--------|-----------|------|-------|
| 80764 | 1600 | 80764 | 80764 | 1600 | 1600 |
| 80764 | 1600 | 80764 | 80764 | 1600 | 1600 |
| 80764 | 1600 | 80764 | 80764 | 1600 | 1600 |

| SYN_SW | SYN_SW | SYN_SW | SYN_SW | SYN_SW | SYN_SW |
|--------|--------|--------|--------|--------|--------|
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |

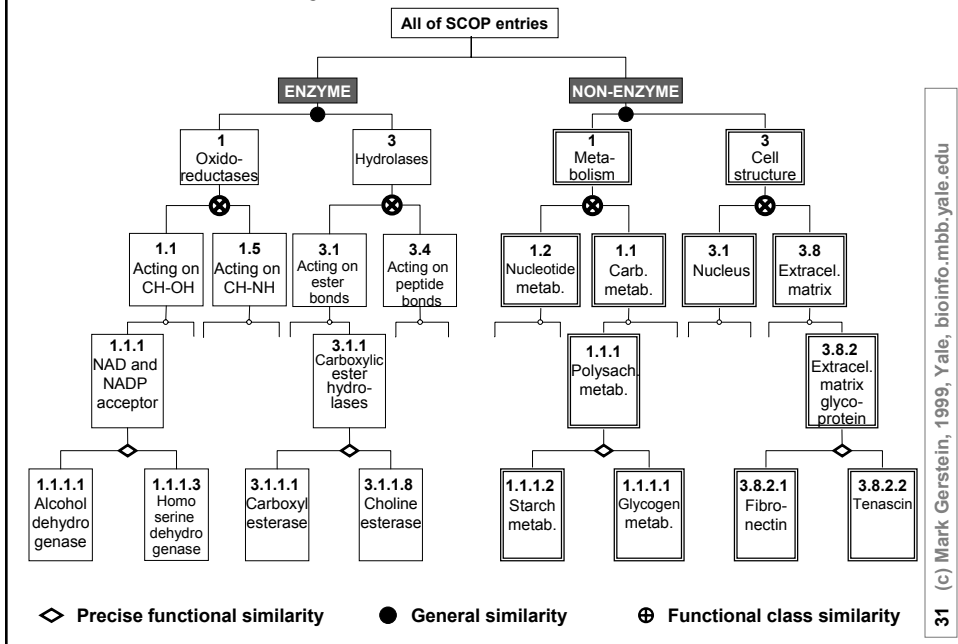
| GENE | SWISSPROT | SYN_GENE | SYN_SW | SYN_SW | SYN_SW |
|-------|-----------|----------|--------|--------|--------|
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |

| REF_TAB | GENE | SWISSPROT | SYN_GENE | SYN_SW | SYN_SW |
|---------|------|-----------|----------|--------|--------|
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |
| 80764 | 1600 | 80764 | 1600 | 80764 | 1600 |

Click a table to dump the first 20 records

- The primary key is the MODULE field which is derived from the ENRUM field
- SYN_GENE and SYN_SW contain all known genes and SWISSPROT names. Rows with identical module numbers are synonyms.
- IDENTY contains the authority gene and SWISSPROT name. Eventually this should be folded into the SYN tables.
- SIMGRPS, GENEROD, SIMPAIRS are keyed to MODULE.
- CAT_TAB and QTY are keyed to MODULE and are the link tables to CATDECOD and GTYDECOD.
- REF_TAB is a link table between REFERENCES and the SYN tables but is currently keyed by GENE and SW and not by MODULE.

Hierarchy of Protein Functions



31 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Can we define FUNCTION?

Problems defining function:

Multi-functionality: 2 functions/protein (also 2 proteins/function)

Conflating of Roles: molecular action, cellular role, phenotypic manifestation.

Non-systematic Terminology:

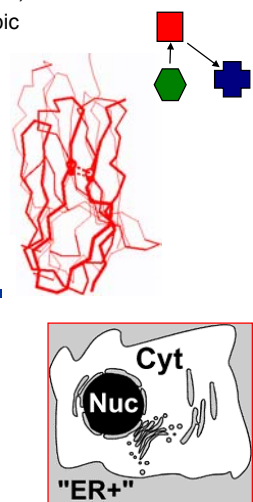
'suppressor-of-white-apricot' & 'darkener-of-apricot'

Fold, Localization, Interactions & Regulation are attributes of proteins that are much more clearly defined

Functional Classification

COGs (cross-org., just conserved, NCBI, Koonin/Lipman)
GenProtEC (*E. coli*, Riley)
ENZYME (SwissProt, Bairoch/ Apweiler, just enzymes, cross-org.)
"Fly" (fly, Ashburner) now extended to **GO** (cross-org.)
MIPS PEDANT (yeast, Mewes)
Other: SwissProt Annotation, WIT, KEGG (just pathways), TIGR EGAD (human ESTs)

VS.



32 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

They have been called aloof, apathetic, uncoordinated, strung out, two-faced, puny and spastic—and by some of the most respected scientists in the world. Who are recognizing such "stupid" names, when their discoverers are free to name as they wish.

Stange Gene Names

Many of the most outlandish labels come from the study of *Drosophila*, the fruit flies. **Yippee**, a novel fly gene, which was described in the *Journal of Molecular Biology*, is one example.

It was named for the reaction of Katarina Roxström-Lindquist, a graduate student at the University of Stockholm, upon cloning yippee. If she has a good result, Katarina has a habit of writing "yippee" in the margin of her notebook.

In such cases, the appellation says more about the scientist than the gene. Star Trek aficionados surface with names like **vulcan** and **klington**.

For an axon bundle that describes what flies see under the microscope, a clear way to describe a mutant that dies during development, usually in the second larval stage. Even liquor preferences make their way out of the cabinet and into the literature with genes such as **grappa**.

But few names are so closely tied to popular tastes or culture. Scientists more often rely on scriptural, literary or historical sources to choose a colorful name that describes what flies see under the microscope.

A gene that affects female fertility was dubbed **sarah** after the wife of Abraham who was infertile for many years but eventually bore a child. (One mutation that causes fly embryos to remain headless affects the *copra* gene, named after a famous slave in ancient Europe who was hounded as a result of his faith. In a recent issue of the *Journal of Cell Biology*, David St. Johnston and colleagues reported on

barentsz, which they named after a Dutch explorer who froze to death in the ice near the North Pole. Why? Because the mutant blocks the movement of a key messenger RNA, causing it to get stuck in the wrong place. **Agoraphobic** refers to a mutant for which the larvae look normal but never crawl out of the egg shell.

Stephen Trevis, a professor of biochemistry at the University of North Carolina at Chapel Hill, reports. "No gene is named in a clever manner, then that probably helps you remember." He followed this

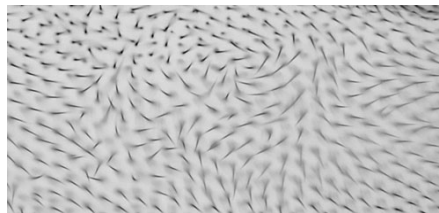
preception in 1987 when he named a *Drosophila* gene **single-minded** after the visual effect of the mutant morphology. Flies with mutations in this gene possess a single bundle of axons in their nervous systems instead of two. He had also considered using simple-minded but abandoned that label because

the name could have been taken as offensive, especially if the function of the equivalent gene in people proved similar related to a mental disability. His concern was well founded. Recent studies have implicated one of the two human single-minded genes in Down syndrome.

Political correctness has indeed been an issue in the past. In 1963 a mutant fly gene was discovered that caused males to court other males. The original gene name of *fruity* was eventually changed to *fruity* after much public disapproval. A similar situation arose more recently, when scientists at Princeton University found mutations in flies that caused them to be leprosy-infected, or, in the vernacular of the investigators, "nigger out." They therefore named the corresponding genes after vegetable-cabbage, rutabaga, radish and turnip—while some scientists found objectionable.



Strange Gene Names 2



how beautiful it was," says Adler. Under the microscope, you can see swirling patterns, rather than all the hairs pointing in the same direction as in wildtype. It immediately brought to mind **Starry Night**, the painting by Van Gogh. So when he isolated a similar gene that same year, he naturally enough named it **Van Gogh**.

Adler asserts that with a descriptive designation, the connection becomes more personal and subtle between the name of the gene and its visual effect.

For many such results study genes, in part because most genes do not affect the behavior or appearance of the organism so directly. When obvious clues to the action of the gene are lacking,

geneticists often pick a name based on the inferred function of the gene product. **Redtape** is the most recent in a series of designations given to genes which, when mutated, block transport along axons.

The predecessors of redtape include **roadblock**, **gridlock** and **Sunday driver**, all emanating from the lab of Lawrence Goldstein at the University of California, San Diego. However, names based on

humor are often not particularly creative—sometimes they are even misleading. For example, one human gene, first described in 1992, is aryl hydrocarbon receptor nuclear translocator, a mouthful that is hard to remember except by its acronym, **ARNT**. Worse, recent studies have shown that ARNT might not act as an aryl hydrocarbon receptor translocator at all, suggesting it might soon be due for rechristening.

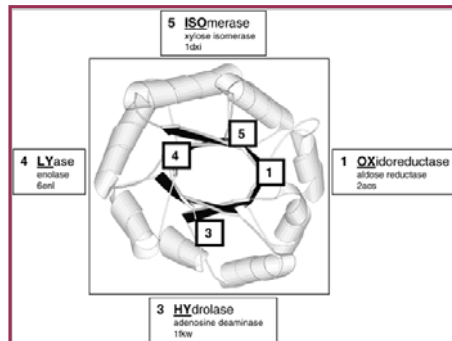
Perhaps the new name will be shorter and more memorable. Or maybe not. Even the author of *Starry Night* and Van Gogh concedes that not all genes can be designated so creatively. As Adler says, "There are a lot of genes, and only so many names you can come up with." Maria Vack

End of class 2002, 11.20 (Bioinfo-12) [starting in databases]

35 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Fold-Function Combinations #1

Many Functions on the
Same Fold
-- e.g. the TIM-barrel



Two Different Folds
Catalyze the Same
Reaction -- e.g.
Carbonic Anhydrases
(4.2.1.1)



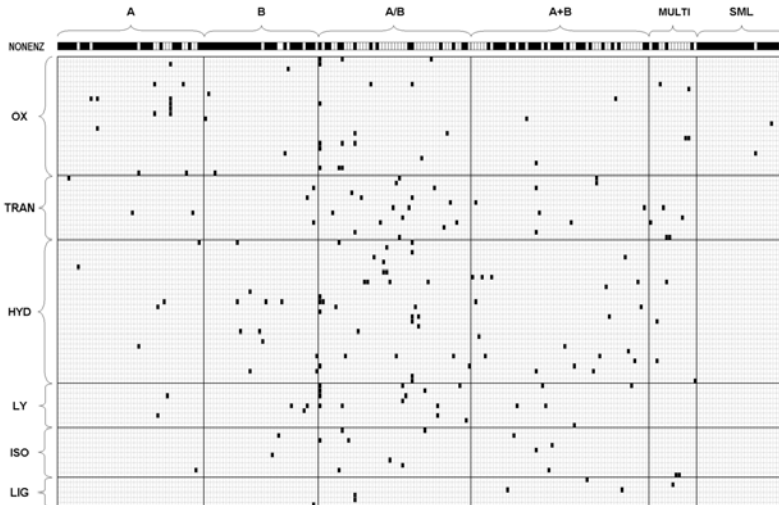
36 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Fold-Function Combinations Cross-Tabulation

~20K (=92x229) Possible,
331 Observed

229 Folds

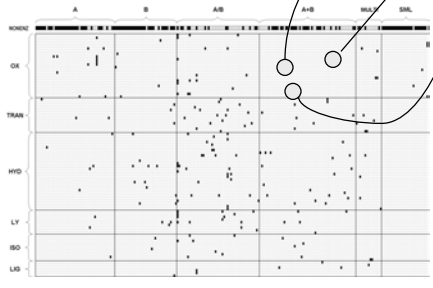
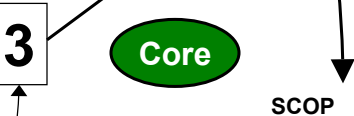
91 Enzymatic Functions
+ Non-Enzyme



37 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Fold-Function Combinations Cross- Tabulation Summary Diagram

| | A | B | A/B | A+B | MULTI | SML | sum |
|--------|----|----|-----|-----|-------|-----|-----|
| NONENZ | 34 | 30 | 14 | 28 | 4 | 26 | 136 |
| OX | 13 | 5 | 17 | 3 | 4 | 5 | 47 |
| TRAN | 3 | 3 | 16 | 3 | 5 | | 35 |
| HYD | 4 | 11 | 30 | 18 | 4 | | 67 |
| LY | 2 | 3 | 13 | 5 | | | 23 |
| ISO | 1 | 2 | 7 | 4 | 2 | | 16 |
| LIG | | 1 | 2 | 3 | 1 | | 7 |
| sum | 57 | 55 | 99 | 69 | 20 | 31 | 331 |



| | A | B | A/B | A+B | MULTI | SML |
|--------|-----|-----|------|-----|-------|-----|
| NONENZ | 7.1 | 5.7 | 7.1 | 9.2 | 2.8 | 0.7 |
| OX | 3.5 | 2.1 | 9.2 | 2.1 | 0.7 | 0.7 |
| TRAN | 0.7 | | 10.6 | 1.4 | 1.4 | 0.7 |
| HYD | 2.8 | 2.8 | 6.4 | 5.7 | 1.4 | |
| LY | 2.1 | | 4.3 | | | |
| ISO | 0.7 | 1.4 | 2.8 | 0.7 | | |
| LIG | | | 1.4 | 1.4 | | |

[Similar analysis in Martin et al. (1998), Structure 6: 875]

The Most Versatile Folds, Versatile Functions

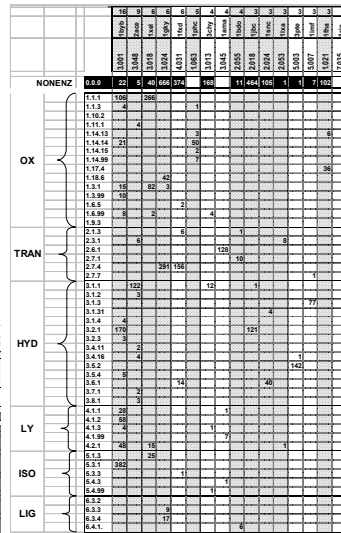
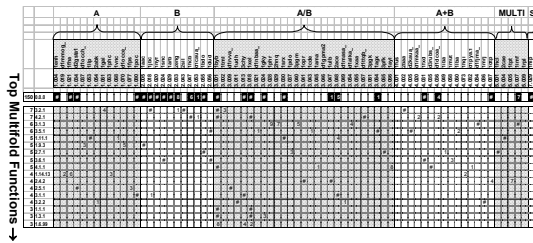
Top Multifunctional Folds →

Top-4 Functions:

Glycosidases, carboxylases, phosphoric monoester hydrolases, linear monoester hydrolases (3.2.1, 4.2.1, 3.1.3, 3.5.1)

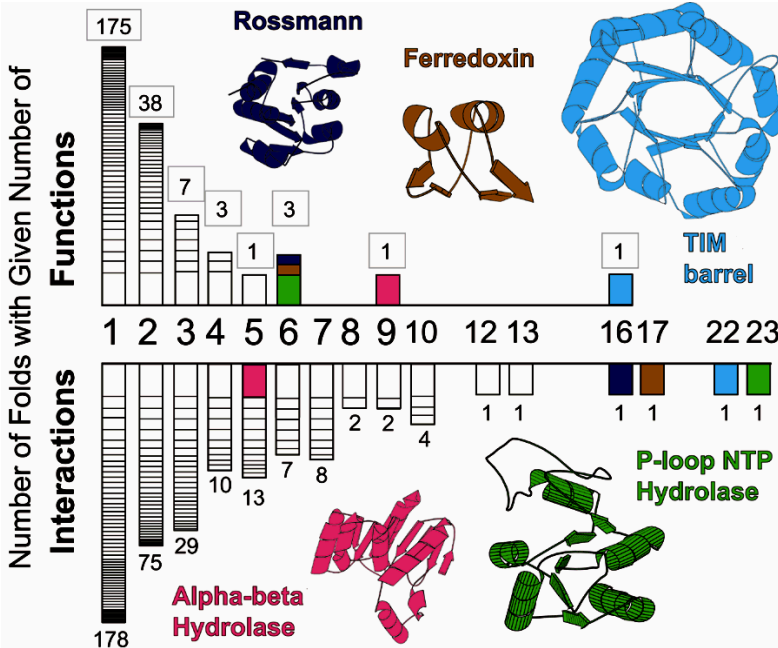
Top-5 Folds:

TIM-barrel (16), alpha-beta hydrolase fold (9), Rossmann fold (6), P-loop NTP hydrolase fold (6), Ferredoxin fold (6)



39 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Most Versatile Folds – Relation to Interactions



Similar results
Martin et al.
(1998)

The number of interactions for each fold = the number of other folds it is found to contact in the PDB

40 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Compare Classifications and Genomes

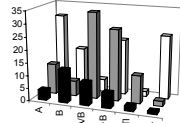
Compare 1 Structure-Function Cross-Tab for Different Genomes and Different Functional & Structural Classifications for the Yeast Genome



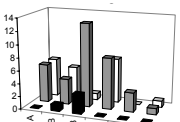
| | | SCOP | | | | | |
|--------|--------|------|-----|------|-----|-------|-----|
| | | A | B | A/B | A+B | MULTI | SML |
| ENZYME | NONENZ | 7.1 | 5.7 | 7.1 | 9.2 | 2.8 | 0.7 |
| | OX | 3.5 | 2.1 | 9.2 | 2.1 | 0.7 | 0.7 |
| | TRAN | 0.7 | | 10.6 | 1.4 | 1.4 | 0.7 |
| | HYD | 2.8 | 2.8 | 6.4 | 5.7 | | 1.4 |
| | LY | 2.1 | | 4.3 | | | |
| | ISO | 0.7 | 1.4 | 2.8 | 0.7 | | |
| | LIG | | | 1.4 | 1.4 | | |

CATH (Thornton)

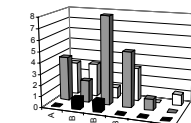
| | | CATH | | |
|--------|--------|------|-----|-----|
| | | A | B | AB |
| ENZYME | NONENZ | 10 | 9.8 | 15 |
| | OX | 5.1 | 5.1 | 10 |
| | TRAN | | 1.3 | 10 |
| | HYD | 2.8 | 1.3 | 14 |
| | LY | 2.8 | | 1.3 |
| | ISO | 1.3 | 1.3 | 5.1 |
| | LIG | | | 1.3 |



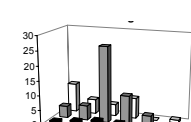
SwissProt



Yeast



worm



E. coli

MIPS YFC (Mewes)

| | | SCOP | | | | | | |
|---------------------|------------------------|------|-----|-----|-----|-------|-----|-----|
| | | A | B | A/B | A+B | MULTI | SML | |
| MIPS Functional Cnt | transcription | 1 | 32 | 23 | 7 | 45 | 1 | 1 |
| | transp. in. class. sp. | 2 | 1 | 1 | 1 | 1 | 1 | 1 |
| | transp. in. class. sp. | 3 | 45 | 34 | 4 | 45 | 14 | 12 |
| | transcription | 4 | 11 | 13 | 22 | 15 | 15 | 15 |
| | protein synthesis | 5 | 1 | 0.5 | 0.7 | 1.3 | 0.3 | 0.2 |
| | protein synthesis | 6 | 12 | 1 | 2 | 1.6 | 0.5 | 0.3 |
| | transp. in. class. sp. | 7 | 59 | 0.5 | 0.7 | 0.6 | 0.4 | |
| | transp. in. class. sp. | 8 | 11 | 24 | 14 | 0 | 1 | |
| | transp. in. class. sp. | 9 | 0.5 | 0.7 | 1 | 0.3 | 0.3 | 0.1 |
| | transp. in. class. sp. | 10 | 1 | 1 | 1.1 | 0.3 | 0.7 | 0.3 |
| | transp. in. class. sp. | 11 | 1 | 1 | 2.6 | 1.6 | 0.7 | 0.5 |
| | transp. in. class. sp. | 13 | 0.5 | 0.3 | 0.4 | 0.4 | 0.2 | |

41 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

COGs vs SCOP: Different Structure Function Relationships for Most Conserved Proteins



| | | SCOP | | | | | | |
|--------------------|----------------------------------|------|-----|-----|-----|-------|-----|-----|
| | | A | B | A/B | A+B | MULTI | SML | |
| All Yeast COGs | Metabolism | C | 2.2 | 2.6 | 4.8 | 3 | 0.4 | |
| | | E | 2.2 | 1.1 | 7.4 | 2.6 | 0.7 | |
| | | F | 1.1 | | 3.7 | 1.8 | | |
| | | G | 0.4 | 0.4 | 3.3 | 0.7 | | |
| | | H | 1.1 | 0.7 | 4.8 | 3 | | |
| | I | 0.7 | 0.7 | 2.2 | 0.4 | 0.4 | | |
| | Information Storage & Processing | J | 2.2 | 1.8 | 3 | 3 | 0.4 | 0.4 |
| | | K | | | 1.1 | 0.4 | | |
| | | L | 1.1 | | 1.8 | 1.1 | 1.1 | |
| M | | | 0.4 | 0.4 | 0.7 | | | |
| Cellular Processes | N | 1.8 | 0.7 | 0.4 | 0.7 | | 0.4 | |
| | O | 1.5 | 1.1 | 3 | 2.2 | 0.4 | 0.4 | |
| | P | | 0.4 | 1.1 | 0.7 | 0.4 | | |

| | | SCOP | | | | | | |
|---------------------|----------------------------------|------|-----|-----|-----|-------|-----|-----|
| | | A | B | A/B | A+B | MULTI | SML | |
| Most Conserved COGs | Metabolism | C | | | 7.2 | 2.9 | | |
| | | E | 1.4 | | 1.4 | 1.4 | | |
| | | F | | | 2.9 | | | |
| | | G | | | 4.3 | 1.4 | | |
| | | H | 1.4 | 2.9 | | 1.4 | | |
| | Information Storage & Processing | J | 8.7 | 7.2 | 7.2 | 10 | 1.4 | 1.4 |
| | | K | | | | | | |
| | | L | | | | | 1.4 | |
| | | M | | | | | | |
| Cellular Processes | N | 1.4 | | 1.4 | | | | |
| | O | 2.9 | | 7.2 | 2.9 | | | |
| | P | | 1.4 | | 2.9 | 1.4 | | |

(Scop, Murzin, Ailey, Brenner, Hubbard, Chothia; COGs, Tatusov, Koonin, Lipman)

42 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

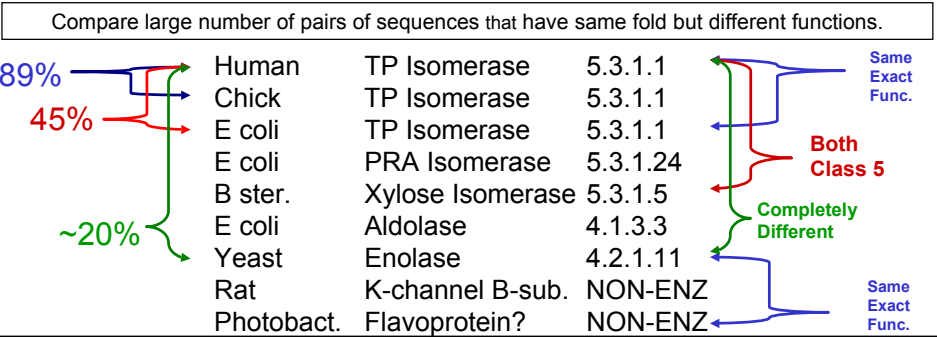
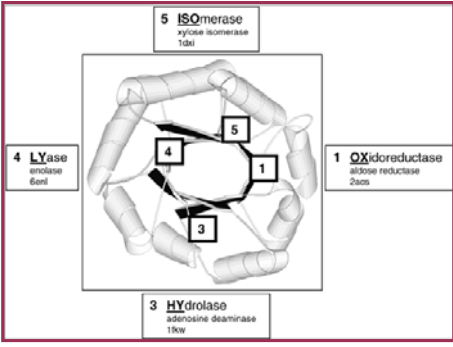
Extra

From here to end of Surveys all is "extra" unless otherwise marked.

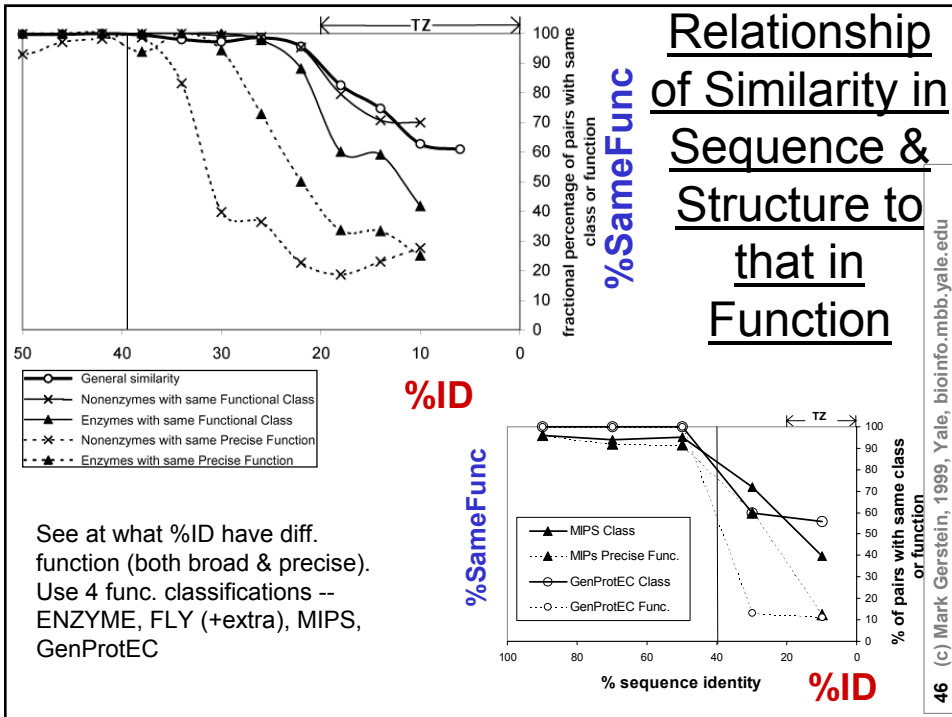
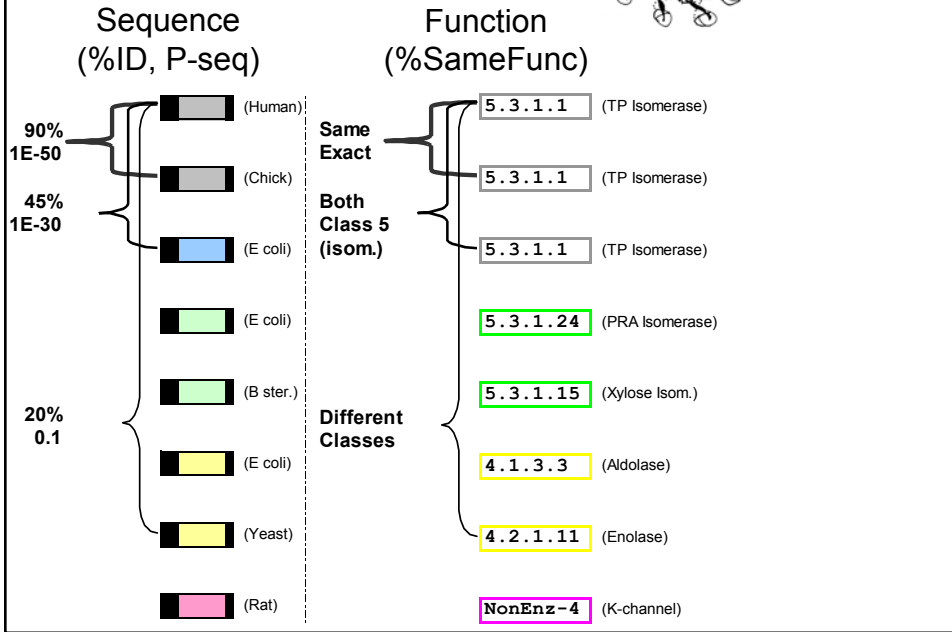
Fold-Function Combinations #2

Many Functions on the Same Fold
 -- e.g. the TIM-barrel

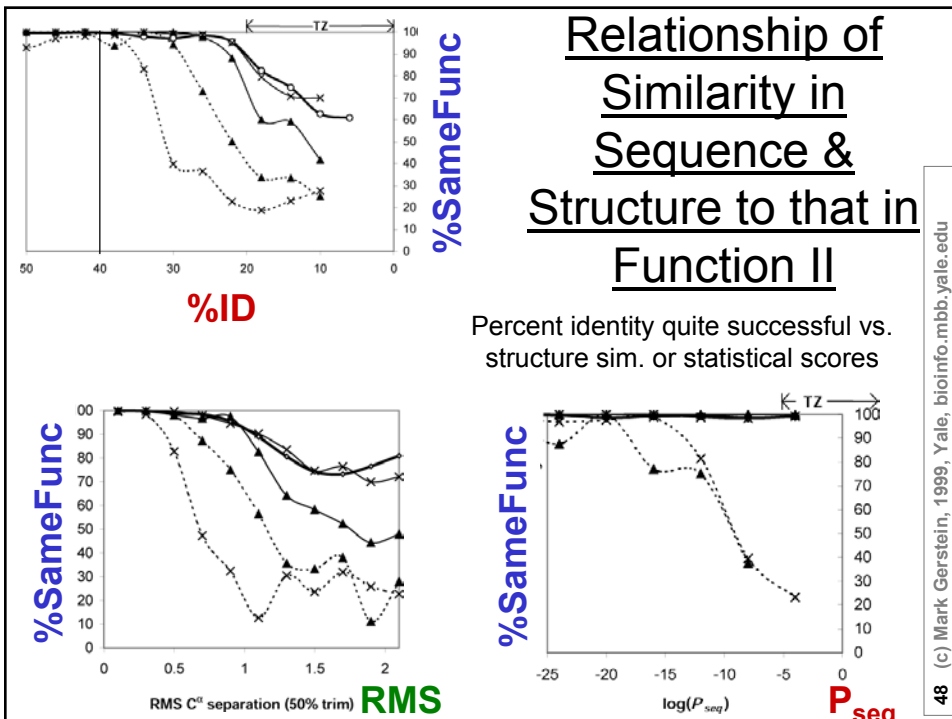
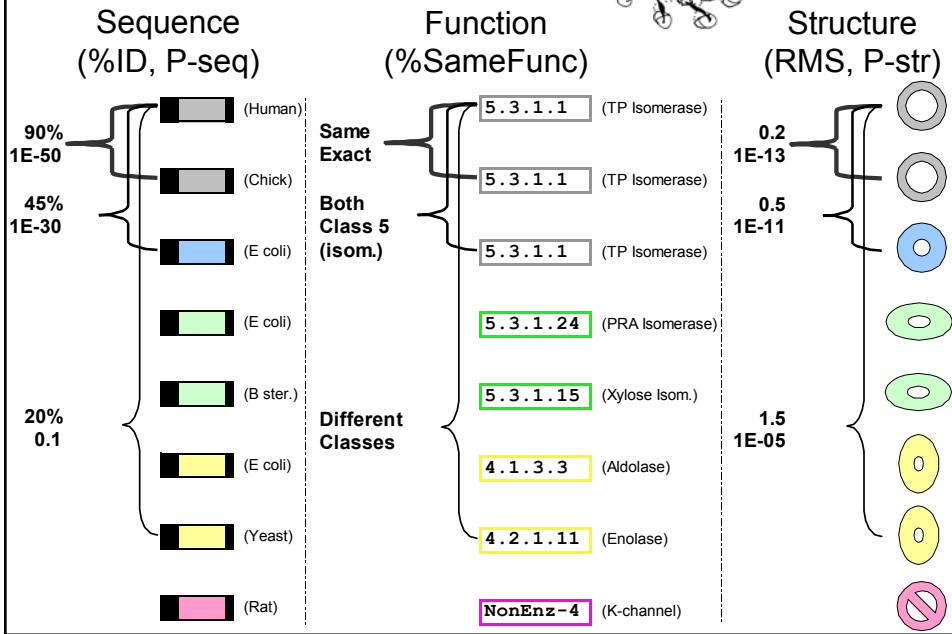
at what degree of divergence?
 Sequence Diverg. (%ID, P_{seq})
 Structural Diverg. (RMS, P_{str})
 Functional Diverg. (%SameFunc)



Annotation Transfer: TIM ex.



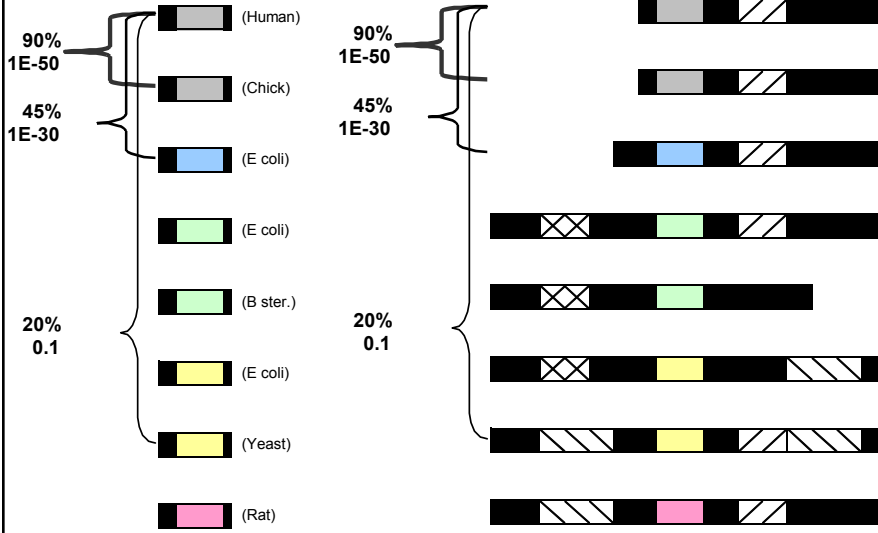
Annotation Transfer: TIM ex.



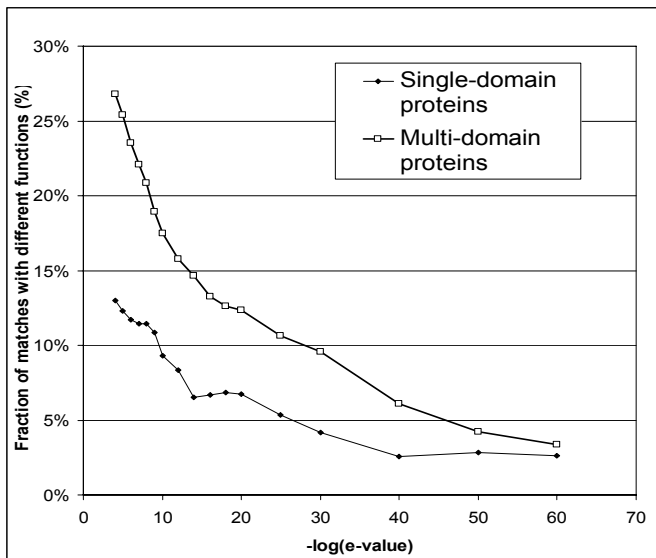
Sequence Divergence of Multidomain Proteins

Divergence in Single Domain Sequences (%ID, P-seq)

Divergence in Multi-domain Sequences



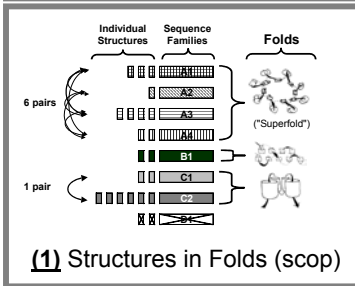
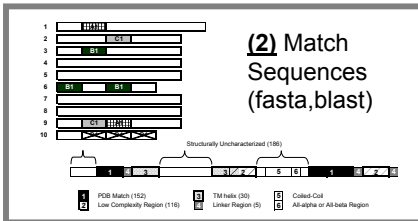
Multi-domain proteins have greater divergence in function with sequence



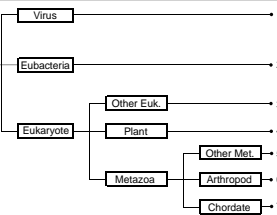
Large-scale Example: Census DB

- 9 Genome Comparison
- 1437 Relational Tables
- 442 Mb
- Simple ASCII Layout

Cross-Reference: Folds → Sequences → Organisms



(3) Organize Sequences by Genome or Taxon



3
+

| Abbrev. | Kingdom (subgroup) | Genome | Num. ORFs | Reference |
|---------|--------------------------|---------------------------------|-----------|-------------------|
| EC | Bacteria (gram negative) | <i>Escherichia coli</i> | 4290 | Blattner et al. |
| HI | Bacteria (gram negative) | <i>Haemophilus influenzae</i> | 1680 | TIGR |
| HP | Bacteria (gram negative) | <i>Helicobacter pylori</i> | 1577 | TIGR |
| MG | Bacteria (gram positive) | <i>Mycoplasma genitalium</i> | 468 | TIGR |
| MJ | Archaea (Euryarchaeota) | <i>Methanococcus jannaschii</i> | 1735 | TIGR |
| MP | Bacteria (gram positive) | <i>Mycoplasma pneumoniae</i> | 677 | Himmelsich et al. |
| SC | Eukarya (fungi) | <i>Saccharomyces cerevisiae</i> | 6218 | Goffeau et al. |
| SS | Bacteria (Cyanobacteria) | <i>Synechocystis sp.</i> | 3168 | Kaneko et al. |

(4) Results in "Fold Table"

| Class | Fold# | EC | SC | HI | SS | HP | MJ | MP | MG | total | Fam. | Rep. | Struc. | Name |
|----------------|-------|----|----|----|----|----|----|----|----|-------|------|------|----------------|---------------|
| α/β | 18 | 60 | 46 | 23 | 40 | 19 | 7 | 4 | 3 | 202 | 16 | 183 | 1xe1- | NAD(P) bindi |
| α/β | 24 | 20 | 69 | 17 | 19 | 17 | 16 | 10 | 11 | 179 | 13 | 132 | lgky- | P-loop conta |
| α/β | 31 | 37 | 28 | 18 | 16 | 12 | 40 | 3 | 3 | 157 | 23 | 160 | 1fxd- | like Fe. pdox |
| α/β | 01 | 45 | 36 | 13 | 22 | 11 | 10 | 5 | 4 | 146 | 37 | 399 | lbyb- | TM-ba at |
| α/β | 23 | 18 | 17 | 7 | 9 | 4 | 8 | 2 | 2 | 67 | 5 | 36 | lpyd a:2-181 | Thiamin bindi |
| α/β | 04 | 15 | 11 | 7 | 10 | 1 | 9 | 5 | 5 | 63 | 13 | 132 | 2tmd a:490-645 | FAD/Ni(D/P) |
| α/β | 55 | 8 | 9 | 7 | 8 | 9 | 3 | 6 | 6 | 56 | 4 | 23 | lszy a:111-421 | Class-aaRS |
| β | 27 | 7 | 10 | 8 | 8 | 4 | 4 | 3 | 3 | 47 | 5 | 19 | 1fnb 19-154 | Reducase/EI |
| β | 24 | 13 | 7 | 4 | 3 | 3 | 3 | 3 | 3 | 39 | 18 | 177 | lsnc- | OB-fold |
| α/β | 11 | 10 | 8 | 4 | 8 | 2 | 2 | 2 | 1 | 37 | 11 | 48 | l1gd- | beta-Grasp |

Integrated Analysis System: X-ref Parts with Genomes

Folds: scop+automatic
 Orthologs: COGs
 "Families": homebrew, ProtoMap

finding parts in genome sequences
blast,
 ψ -blast,
fasta,
 TM, low-complexity
 &
 (Altschul, Pearson, Wootton)



part occurrence profiles

One approach of many...
 Much previous work on
 Sequence & Structure Clustering
 CATH, Blocks, FSSP,
 Interpro, eMotif, Prosite,
 CDD, Pfam, Prints, VAST,
 TOGA...

Remington, Matthews '80; **Taylor, Orengo '89, '94**; Thornton, CATH; Artymiuk, Rice, Willett '89; Sali, Blundell, '90; Vriend, Sander '91; Russell, Barton '92; **Holm, Sander '93+ (FSSP)**; Godzik, Skolnick '94; **Gibrat, Bryant '96 (VAST)**; F Cohen, '96; Feng, Sippl '96; G Cohen '97; Singh & Brutlag, '98

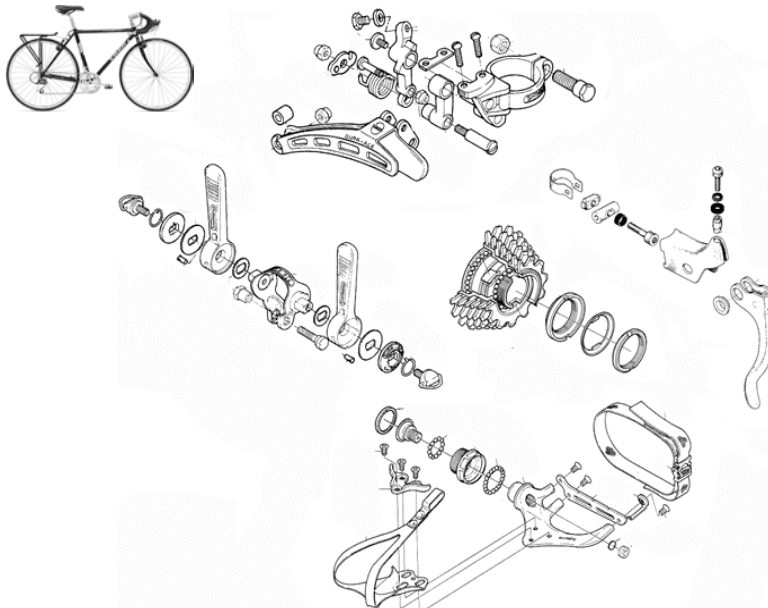
PDB Report - Netscape

Genome Occurrence of 3tim's Fold

Results in Old Format Formatted Text in Popup Window See Matches in Genome

| ID | 3tim | Chain | a | Domain | SCOP Fold Number | 3.001 | Superfamily | 3.1.11 | | | | | | | | | | | | |
|-------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|--------------------------|
| Data Genome | Mu | Msr | Mm | Mph | Mse | Mca | Mae | Mco | Mbs | Mhr | Msp | Mgn | Mno | Mou | Mtr | Mcr | Mpr | | | |
| Fold Occurrences | 44 | 38 | 43 | 30 | 84 | 123 | 28 | 51 | 63 | 88 | 77 | 36 | 22 | 7 | 12 | 20 | 15 | 19 | 18 | 9 |
| See Matches in Genome | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |
| Superfamily Occurrences | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 2 | 1 | 1 | 1 | 0 |
| See Matches in Genome | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> | <input type="checkbox"/> |

A Parts List Approach to Bike Maintenance



A Parts List Approach to Bike Maintenance

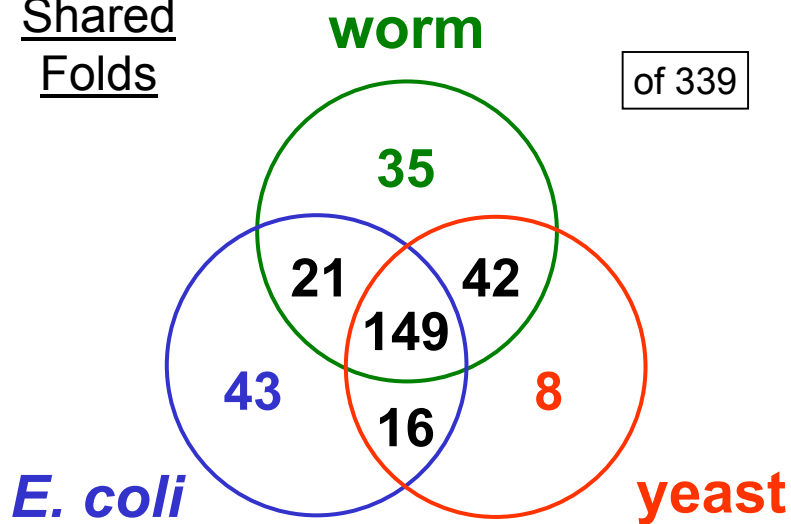
How many roles can these play?
How flexible and adaptable are they mechanically?

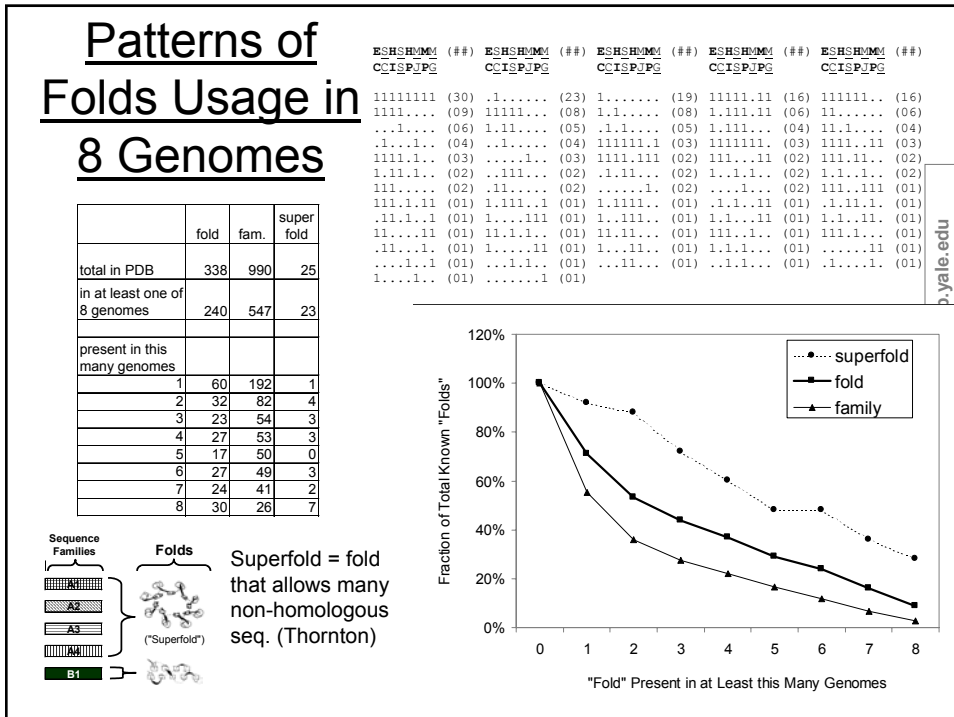
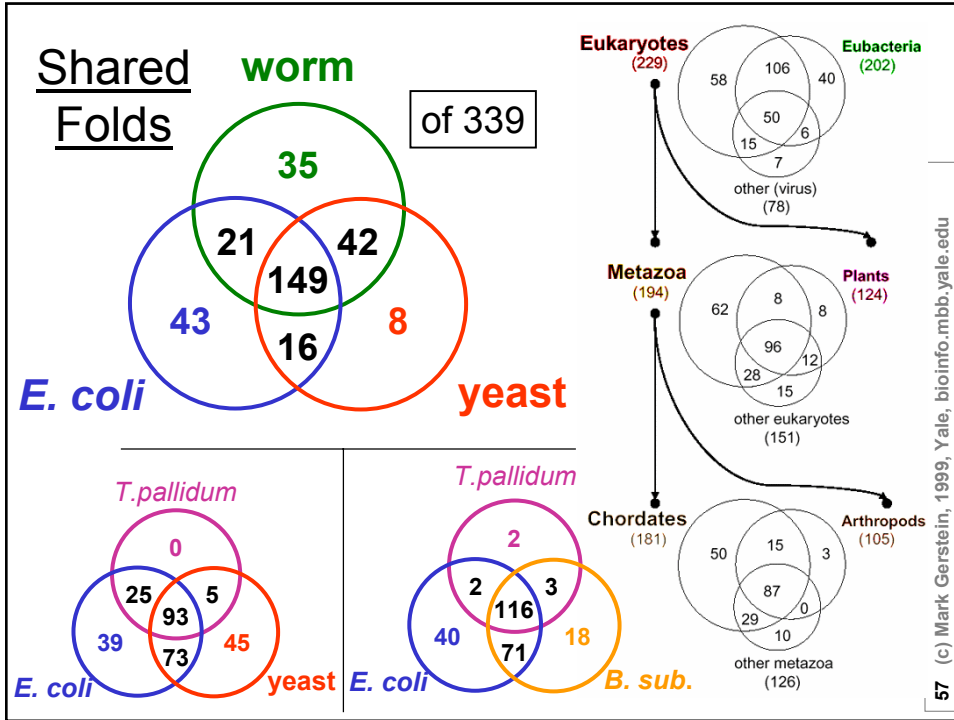
What are the shared parts (**bolt**, **nut**, **washer**, **spring**, **bearing**), unique parts (**cogs**, **levers**)? What are the common parts - types of parts - (nuts & washers)?

Where are the parts located?
Which parts interact?

55 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Shared Folds





Whole Genome Trees

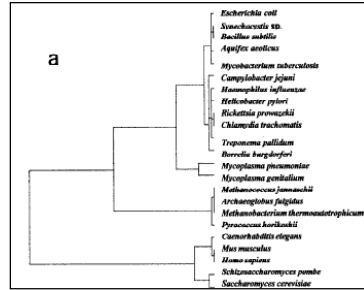
The Genomic Tree as Revealed from Whole Proteome Comparisons

Fredj Tekala,^{1,3} Antonio Lazcano,² and Bernard Dujon¹

¹Unité de Génétique Moléculaire des Levures [URA1300 Centre National de la Recherche Scientifique (CNRS) and UFR927 University Pierre and Marie Curie], Institut Pasteur, 75224 Paris Cedex 15, France; ²Facultad de Ciencias, UNAM, Apdo. Cd. Universitaria, 04510 Mexico City, Mexico

The availability of a number of complete cellular genome sequences allows the development of organisms' classification, taking into account their genome content, the loss or acquisition of genes, and overall gene similarities as signatures of common ancestry. On the basis of correspondence analysis and hierarchical classification methods, a methodological framework is introduced here for the classification of the available 20 completely sequenced genomes and partial information for *Schizosaccharomyces pombe*, *Homo sapiens*, and *Mus musculus*. The outcome of such an analysis leads to a classification of genomes that we call a genomic tree. Although these trees are phenograms, they carry with them strong phylogenetic signatures and are remarkably similar to 16S-like rRNA-based phylogenies. Our results suggest that duplication and deletion events that took place through evolutionary time were globally similar in related organisms. The genomic trees presented here place the Archaea in the proximity of the Bacteria when the whole gene content of each organism is considered, and when ancestral gene duplications are eliminated. Genomic trees represent an additional approach for the understanding of evolution at the genomic level and may contribute to the proper assessment of the evolutionary relationships between extant species.

The determination of complete genome sequences from ≥ 20 organisms offers an unprecedented opportunity to study horizontal transfer (which may have been more intense during early cellular evolution, Woese 1998), un-



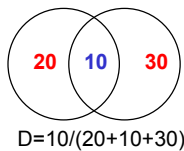
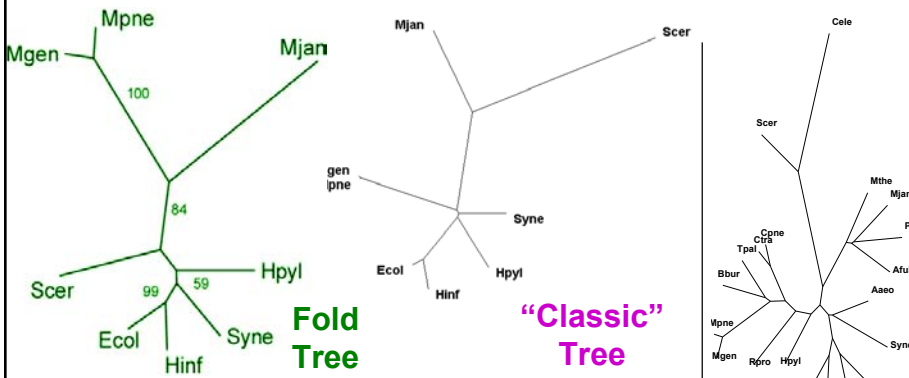
550 Genome Research
www.genome.org

9550-557 ©1999 by Cold Spring Harbor Laboratory Press ISSN 1054-9803/99 \$5.00 www.genome.org

orthologs, homologs, folds, motifs

59 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Cluster Trees Grouping Initial Genomes on Basis of Shared Folds

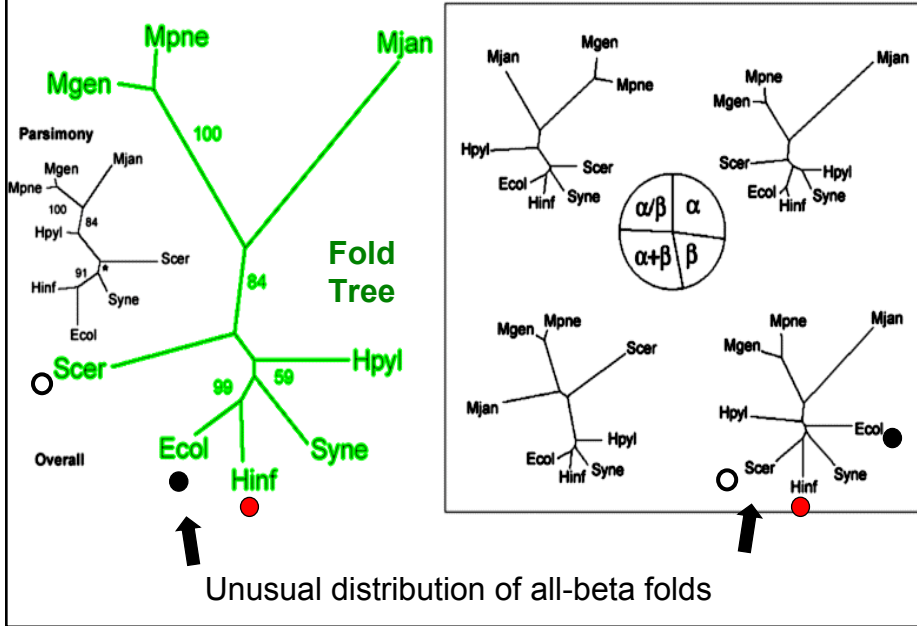


$D=S/T$ $S = \#$ shared folds
 $D =$ shared fold dist. betw. 2 genomes $T =$ total # folds in both

20 Genomes

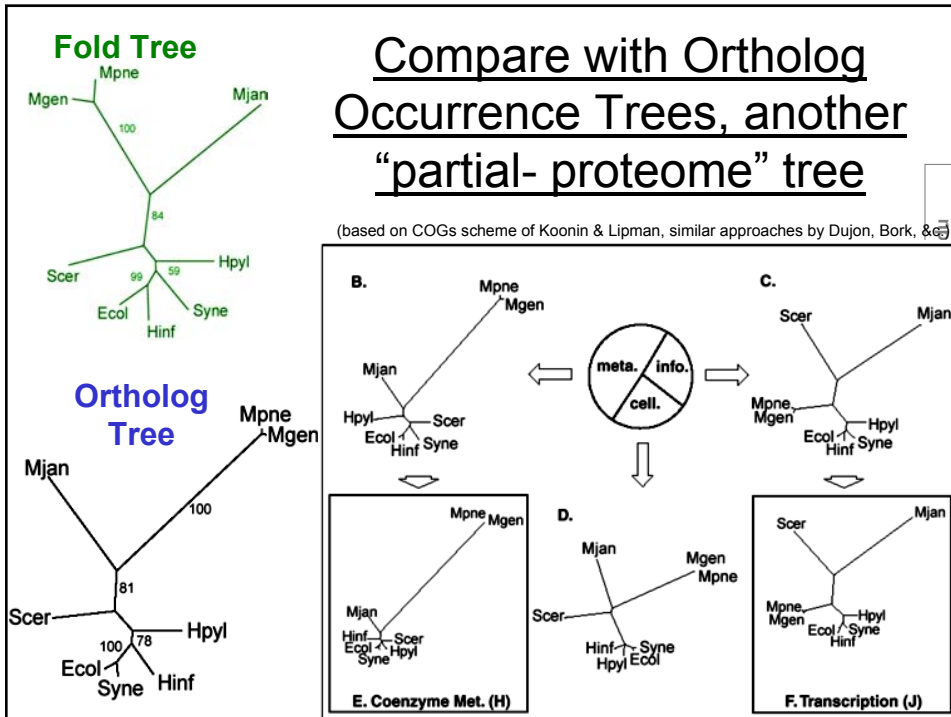
60 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Distribution of Folds in Various Classes

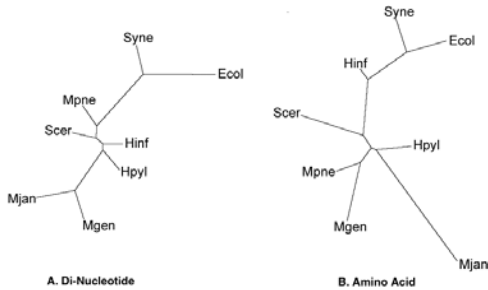
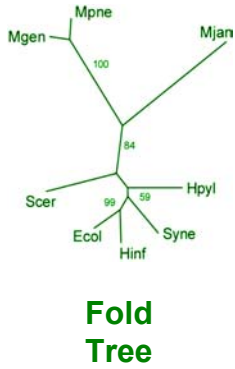
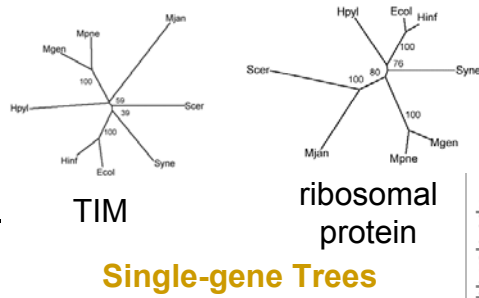


Compare with Ortholog Occurrence Trees, another "partial- proteome" tree

(based on COGs scheme of Koonin & Lipman, similar approaches by Dujon, Bork, & Ge)



Compare with trees on spectrum of "levels": single-gene trees, whole-genome composition trees



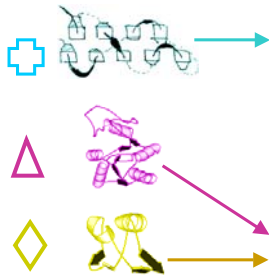
AA & di-NT Composition Trees
(S Karlin)

Ortholog Tree

"Classic" Tree

Common Folds in Genome, Varies Betw. Genomes

Depends on comparison method, DB, sfams v folds, & (new top superfamilies via v-Blast, Intersection of top-10 to get shared and common)



| Top-10 Worm Folds | class | num. matches in worm genome (N) | frac. all worm dom. (F) | in EC? | in SC? |
|----------------------------------|-------|---------------------------------|-------------------------|--------|--------|
| Ig | B | 830 | 1.7% | | |
| Knottins | SML | 565 | 1.1% | | |
| Protein kinases (cat. core) | MULT | 472 | 0.9% | | |
| C-type lectin-like | A+B | 322 | 0.6% | | |
| corticoid recep. (DNA-bind dom.) | SML | 276 | 0.5% | | |
| Ligand-bind dom. nuc. receptor | A | 257 | 0.5% | | |
| alpha-alpha superhelix | A | 247 | 0.5% | | |
| C2H2 Zn finger | SML | 239 | 0.5% | | |
| P-loop NTP Hydrolase | A/B | 235 | 0.5% | | |
| Ferredoxin | A+B | 207 | 0.4% | | |

| | <i>M. genitalium</i> | <i>S. subtilis</i> | <i>E. coli</i> |
|---------------------------|------------------------------|-----------------------------|-----------------------------|
| Rank | Superfamily # | Superfamily # | Superfamily # |
| 1 | P-loop hydrolase 60 | P-loop hydrolase 173 | P-loop hydrolase 191 |
| 2 | SAM methyl-transferase 16 | Rossmann domain 165 | Rossmann domain 158 |
| 3 | Rossmann domain 13 | Phosphate-binding barrel 79 | Phosphate-binding barrel 64 |
| 4 | Class I synthetase 12 | PLP-transferase 44 | PLP-transferase 38 |
| 5 | Class II synthetase 11 | CheY-like domain 36 | CheY-like domain 36 |
| 6 | Nucleic acid binding dom. 11 | SAM methyl-transferase 30 | Ferredoxins 35 |
| Total ORFs | 479 | 4268 | 4268 |
| with Common Superfamilies | 105 (22%) | 465 (11%) | 458 (11%) |

Eubacteria

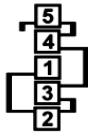
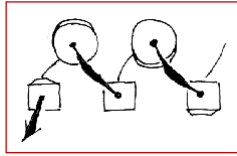
| | <i>M. thermoautotrophicum</i> | <i>A. fulgidus</i> |
|---------------------------|-------------------------------|-----------------------------|
| Rank | Superfamily # | Superfamily # |
| 1 | P-loop hydrolase 93 | P-loop hydrolase 118 |
| 2 | Phosphate-binding barrel 54 | Rossmann domain 104 |
| 3 | Rossmann domain 53 | Phosphate-binding barrel 56 |
| 4 | Ferredoxins 48 | Ferredoxins 49 |
| 5 | SAM methyl-transferase 17 | SAM methyl-transferase 24 |
| 6 | PLP-transferase 15 | PLP-transferase 18 |
| Total ORFs | 1869 | 2409 |
| with Common Superfamilies | 252 (14%) | 309 (13%) |

Archaea

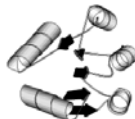
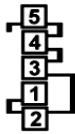
| | <i>S. cerevisiae</i> |
|---------------------------|---------------------------|
| Rank | Superfamily # |
| 1 | P-loop hydrolase 249 |
| 2 | Protein kinase 128 |
| 3 | Rossmann domain 90 |
| 4 | RNA-binding domain 75 |
| 5 | SAM methyl-transferase 63 |
| 6 | Ribonuclease H-like 57 |
| Total ORFs | 6218 |
| with Common Superfamilies | 560 (9%) |

Yeast

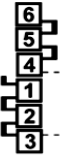
Common, Shared Folds: $\beta\alpha\beta$ structure



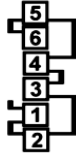
P-loop hydrolase



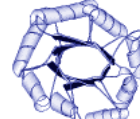
Flavodoxin like



Rossmann Fold



Thiamin Binding



TIM-barrel

A peptide model of a protein folding intermediate

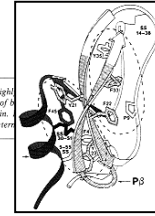
Terrence G. Oas & Peter S. Kim

Whitehead Institute for Biomedical Research, New Cambridge Center, Cambridge, Massachusetts 02142, USA
Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

It is difficult to determine the structures of protein folding intermediates because folding is a highly disulphide-bonded peptide pair, designed to mimic the first crucial intermediate in the folding of B inhibitor, contains secondary and tertiary structure similar to that found in the native protein. circumvent the problem of cooperativity and permit characterization of structure of folding intermediates.

All share $\alpha\beta$ structure with repeated R.H. $\beta\alpha\beta$ units connecting adjacent strands or nearly so (18+4+2 of 24)

336: 42



HI, MJ, SC vs scop 1.32

65 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

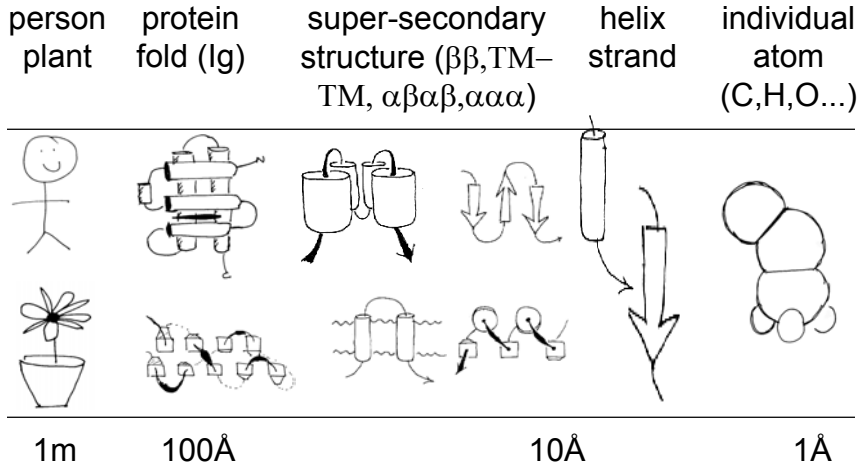
What are the most common folds: Overall? In plants? In animals?

| EC Class | EC Class | Fold Name | Num. Sequ. | Families | Num. of Sequences | | | | | | |
|-------------------------------|---------------|-----------------------------------|------------|----------|-------------------|-------|--------|---------|-----------|-------|------|
| | | | | | Total | Plant | Animal | Microbe | Eukaryote | Other | |
| Totals | | | | | 719 | 37796 | 5169 | 3029 | 4900 | 19316 | 1828 |
| Overall Top-10 | | | | | | | | | | | |
| 3J82-A | β | Immunoglobulin-like | 22 | 19 | 6 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6T3K-B | $\alpha\beta$ | TIM-barrel | 29 | 6 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1ATP-E | O | Protein Kinases (catalytic core) | 1 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1JPD | O | Ferredoxin-like | 17 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1JXK-A | $\alpha\beta$ | NTP Hydrolases containing P-loop | 9 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1J8D-C | α | DNA-binding 3-helical bundle | 13 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2R5D-A | $\alpha\beta$ | Rossmann Fold (NAD binding) | 11 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1J8D | α | Globin-like | 3 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2R82 | $\alpha\beta$ | like Ribonuclease H | 15 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1Z8P | S | Classic Zinc Finger | 3 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| Sequence Family Top-11 | | | | | | | | | | | |
| 1J82-A | β | Immunoglobulin-like | 22 | 19 | 6 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6T3K-B | $\alpha\beta$ | TIM-barrel | 29 | 6 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1JPD | O | Ferredoxin-like | 17 | 4 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2R82 | $\alpha\beta$ | like Ribonuclease H | 15 | 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1Z8P | β | OS-fold | 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1PTX | S | Small inhibitors, toxins, lectins | 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2T8V-C | β | Viral coat and capsid proteins | 14 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1J8D-C | α | DNA-binding 3-helical bundle | 13 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2R5D-A | $\alpha\beta$ | Rossmann Fold (NAD binding) | 11 | 3 | 2 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1J8P | $\alpha\beta$ | Flavodoxin-like | 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1J82B | α | 4-helical cytokines | 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

| Fold Name | Number | Percent of Sequences | | | | |
|--|--------|----------------------|-----------|-------|---------|-------|
| | | Virus | Eukaryote | Plant | Microbe | Other |
| Plant Top-10 | | | | | | |
| $\alpha\beta$ TIM-barrel | 29 | 6 | 7 | 20 | 4 | 13 |
| O like Ferredoxin | 17 | 4 | 2 | 13 | 6 | 6 |
| $\alpha\beta$ NTP Hydrolases containing P-loop | 9 | 3 | 0 | 0 | 0 | 0 |
| O Protein Kinases (catalytic core) | 1 | 4 | 3 | 0 | 0 | 0 |
| S Small inhibitors, toxins, lectins | 14 | 0 | 0 | 0 | 0 | 0 |
| $\alpha\beta$ Rossmann Fold (NAD binding) | 11 | 3 | 0 | 3 | 1 | 3 |
| O RuBisCO (small subunit) | 1 | 0 | 0 | 0 | 0 | 0 |
| β like Concavein A | 6 | 0 | 0 | 0 | 0 | 0 |
| $\alpha\beta$ like Hydrophobic Seed Protein | 2 | 0 | 0 | 0 | 0 | 0 |
| $\alpha\beta$ like Ribonuclease H | 15 | 2 | 0 | 0 | 0 | 0 |
| Metazoan Top-10 | | | | | | |
| β like Immunoglobulin | 32 | 13 | 1 | 1 | 26 | 1 |
| O Protein Kinases (catalytic core) | 1 | 4 | 3 | 3 | 6 | 6 |
| α DNA-binding 3-helical bundle | 13 | 3 | 2 | 5 | 0 | 0 |
| α like Globin | 3 | 2 | 0 | 0 | 0 | 0 |
| S Classic Zinc Finger | 2 | 1 | 0 | 0 | 0 | 0 |
| $\alpha\beta$ NTP Hydrolases containing P-loop | 9 | 3 | 0 | 0 | 0 | 0 |
| β Trypsin-like serine proteases | 4 | 1 | 1 | 0 | 0 | 0 |
| α Cytochrome P450 | 1 | 1 | 0 | 0 | 0 | 0 |
| S like Glucocort. receptor (DNA-binding) | 4 | 1 | 0 | 0 | 0 | 0 |
| EF-hand | 3 | 1 | 0 | 0 | 0 | 0 |

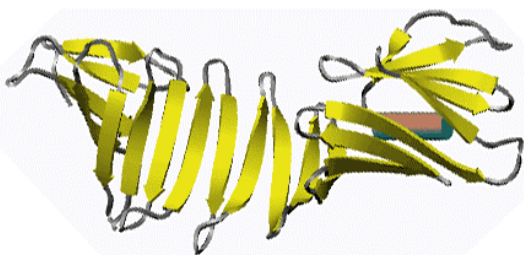
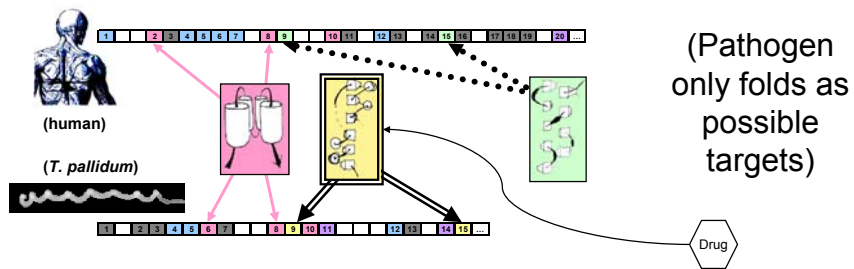
66 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

At What Structural Resolution Are Organisms Different?



67 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Practical Relevance of Structural Genomics



- OspA protein
 - ◇ in Lyme-disease spirochete *B. burgdorferi*
 - ◇ previously identified as the antigen for vaccine
 - ◇ has novel fold (C Lawson)

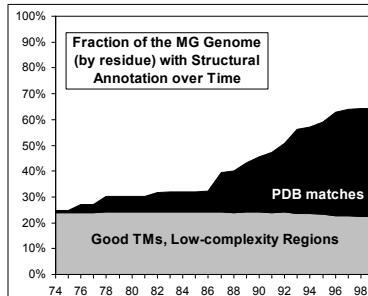
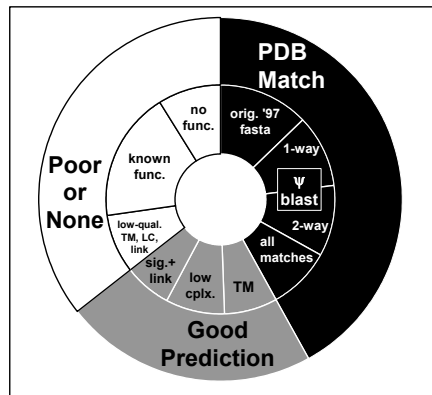
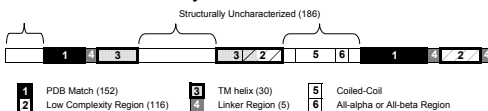
68 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection
- Function Classification
- Cross-tabulation, folds and functions

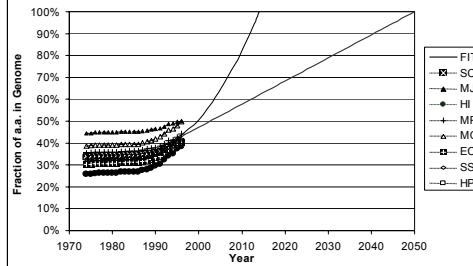
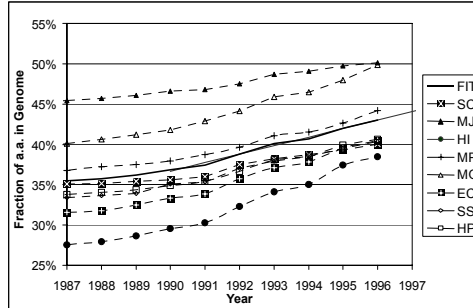
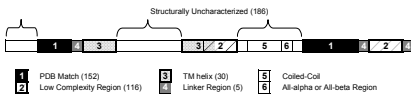
Know All Folds in a Genome: How are we doing on MG?

- MG smallest genome with 479 ORFs
- Separate PDB Match, TMs, LC (SEG), linkers
- How many residues in genome matched by known folds, in 1975, '76, '77...'00...'50
- The impact of PSI-blast in comparison to pairwise methods
 - ◊ Two way PSI-blast gives an improvement (genome vs PDB, PDB vs. genome)
- Union of many sets of PDB matches finds >40% of a.a. and more than half the ORFs (242/479)
 - ◊ (Eisenberg, Godzik, Bork, Koonin, Frishman)
- ~65% structurally characterized



Know All Folds in Genome: MG Optimistic → Prediction

- Just use one pairwise method for matching
- Multiple, big genomes (e.g. SC)

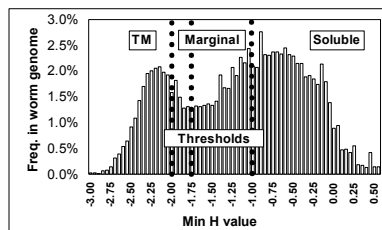
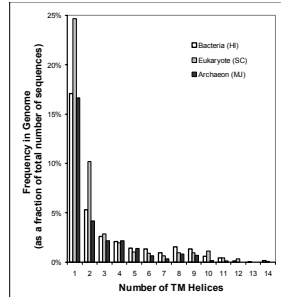
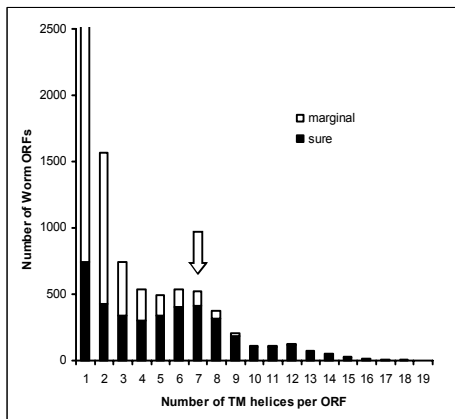


71 (c) M2

TM-helix "prediction"

- TM prediction (KD, GES). Count number with 2 peaks, 3 peaks, &c.
- Similar conclusions to others: von Heijne, Rost, Jones, &c.

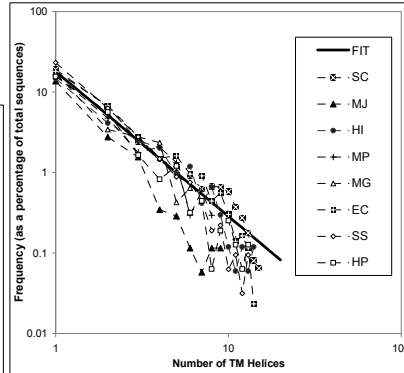
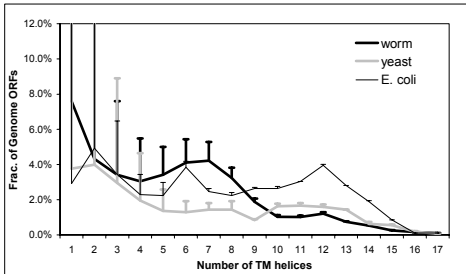
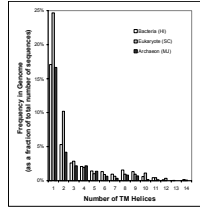
- Divide Predictions into sure and marginal (Boyd & Beckwith's criteria)



72 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Comparative Genomics of Membrane Proteins

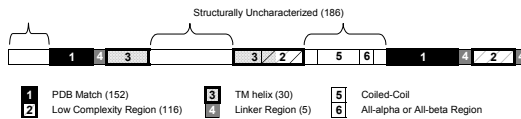
- Yeast has more mem. prots., esp. 2-TMs
- Similar conclusions to others: von Heijne, Rost, Jones, &c.



- Overall, no strong preference for particular supersecondary structures
 - ◊ Freq. of Number of TM helices follows a Zipf-like law: $F=1/[5n^2]$
- In detail, worm has a peak for 7-TMs and *E. coli* for 12-TMs

2° Structure Prediction

- Bulk prediction of 2° struc. in genomes
- Same fraction of α and β (by element, half each)
- Both overall and only for unknown soluble proteins.



- Diff From PDB: 31% helical and 21% strand.
- Related results: Frishman

Not expected since.....

| Fraction of residues Predicted to be in... | strand | helix |
|--|------------|------------|
| Avg | 17% | 39% |
| SD | 1% | 2% |
| EC | 17% | 39% |
| HI | 16% | 41% |
| HP | 15% | 42% |
| MG | 17% | 39% |
| MJ | 19% | 37% |
| MP | 17% | 39% |
| SC | 17% | 34% |
| SS | 16% | 38% |

Different Amino Acid Composition Should Give Different 2° Structure

Each a.a. has different propensity for local structure
 -> Different Compositions (K from 4.4 in EC to 10.4 in MJ, Q too)
 -> Different Local Structure (but compensation?)
 Propensities from Regan (beta) and Baldwin (alpha)

| | Amino Acid Composition | | | | | | | | Propensity (kcal/mole) | | |
|------------------|------------------------|-------|-------|-------|-------|-------|-------|-------|------------------------|-------|--------|
| | EC | HI | SS | SC | HP | MP | MG | MJ | TM-hlx | helix | strand |
| K | 4.4 | 6.3 | 4.2 | 7.3 | 8.9 | 8.6 | 9.5 | 10.4 | 8.8 | -1.5 | -0.4 |
| C | 1.2 | 1.0 | 1.0 | 1.3 | 1.1 | .8 | .8 | 1.3 | -2 | -1.1 | -0.8 |
| R | 5.5 | 4.5 | 5.1 | 4.5 | 3.5 | 3.5 | 3.1 | 3.8 | 12.3 | -1.9 | -0.4 |
| N | 4.0 | 4.9 | 4.0 | 6.1 | 5.9 | 6.2 | 7.5 | 5.3 | 4.8 | -1 | -0.5 |
| Q | 4.4 | 4.6 | 5.6 | 3.9 | 3.7 | 5.4 | 4.7 | 1.5 | 4.1 | -1.3 | -0.4 |
| A | 9.5 | 8.2 | 8.5 | 5.5 | 6.8 | 6.7 | 5.6 | 5.5 | -1.6 | -1.9 | 0 |
| I | 6.0 | 7.1 | 6.3 | 6.6 | 7.2 | 6.6 | 8.2 | 10.5 | -3.1 | -1.2 | -1.3 |
| H | 2.3 | 2.1 | 1.9 | 2.2 | 2.1 | 1.8 | 1.6 | 1.4 | 3 | -1.1 | -0.4 |
| S | 5.8 | 5.8 | 5.8 | 9.0 | 6.8 | 6.5 | 6.6 | 4.5 | -0.6 | -1.1 | -0.9 |
| M | 2.8 | 2.4 | 2.0 | 2.1 | 2.2 | 1.6 | 1.5 | 2.2 | -3.4 | -1.4 | -0.9 |
| P | 4.4 | 3.7 | 5.1 | 4.3 | 3.3 | 3.5 | 3.0 | 3.4 | 0.2 | 3 | >3.0 |
| G | 7.4 | 6.6 | 7.4 | 5.0 | 5.8 | 5.5 | 4.6 | 6.3 | -1 | 0 | 1.2 |
| F | 3.9 | 4.5 | 4.0 | 4.5 | 5.4 | 5.6 | 6.1 | 4.2 | -3.7 | -1 | -1.1 |
| E | 5.7 | 6.5 | 6.0 | 6.5 | 6.9 | 5.7 | 5.7 | 8.7 | 8.2 | -1.2 | -0.2 |
| Y | 2.9 | 3.1 | 2.9 | 3.4 | 3.7 | 3.2 | 3.2 | 4.4 | 0.7 | -1.2 | -1.6 |
| V | 7.1 | 6.7 | 6.7 | 5.6 | 5.6 | 6.5 | 6.1 | 6.9 | -2.6 | -0.8 | -0.9 |
| T | 5.4 | 5.2 | 5.5 | 5.9 | 4.4 | 6.0 | 5.4 | 4.0 | -1.2 | -0.6 | -1.4 |
| D | 5.1 | 5.0 | 5.0 | 5.8 | 4.8 | 5.0 | 4.9 | 5.5 | 9.2 | -1 | 0.9 |
| L | 10.6 | 10.5 | 11.4 | 9.6 | 11.2 | 10.3 | 10.7 | 9.5 | -2.8 | -1.6 | -0.5 |
| W | 1.5 | 1.1 | 1.6 | 1.0 | .7 | 1.2 | 1.0 | .7 | -1.9 | -1.1 | -1 |
| total propensity | | | | | | | | | | | |
| α | -1.00 | -1.02 | -0.96 | -1.00 | -1.05 | -1.03 | -1.05 | -1.01 | | | |
| β | -0.27 | -0.33 | -0.26 | -0.36 | -0.37 | -0.38 | -0.42 | -0.36 | | | |

75 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Supersecondary structure words

- Look at super-secondary patterns ("words" such as $\alpha\alpha$ or $\beta\alpha\beta$) in predictions
- Compare observed freq. with expected freq.
 $odds = f(\alpha\beta)/f(\alpha)f(\beta)$
 (Freq. Words, Karlin)
- Do have differences between genomes (and PDB) here

HI more $\alpha\alpha$, $\alpha\alpha\alpha$, $\alpha\alpha\alpha\alpha$...



SC more $\beta\beta$, $\beta\beta\beta$, $\beta\beta\beta\beta$...



MJ more $\alpha\beta\alpha\beta$, $\beta\alpha\beta\alpha$...



| Super-Secondary Structure "Word" | Maximum Difference between 3 Genomes | Relative Abundance (Odds Ratio) | | | |
|----------------------------------|--------------------------------------|---------------------------------|-------------|-------------|------|
| | | HI | MJ | SC | PDB |
| $\beta\beta$ | 26% | 0.96 | 1.06 | 1.24 | 1.22 |
| $\alpha\alpha$ | 15% | 0.97 | 0.85 | 0.83 | 0.85 |
| $\alpha\beta$ | 10% | 1.09 | 1.09 | 0.99 | 0.95 |
| $\beta\alpha$ | 7% | 0.98 | 1.00 | 0.93 | 0.99 |
| $\beta\beta\beta$ | 41% | 0.96 | 1.15 | 1.46 | 1.62 |
| $\alpha\alpha\alpha$ | 19% | 1.01 | 0.83 | 0.84 | 0.92 |
| $\alpha\beta\alpha$ | 18% | 1.04 | 1.03 | 0.87 | 1.16 |
| $\alpha\alpha\beta$ | 15% | 1.03 | 0.97 | 0.89 | 0.70 |
| $\beta\alpha\beta$ | 12% | 1.15 | 1.24 | 1.10 | 1.19 |
| $\beta\alpha\alpha$ | 11% | 0.93 | 0.87 | 0.83 | 0.78 |
| $\beta\beta\alpha$ | 9% | 0.90 | 0.94 | 0.99 | 0.82 |
| $\alpha\beta\beta$ | 6% | 0.97 | 0.98 | 1.03 | 0.80 |
| $\beta\beta\beta\beta$ | 54% | 1.03 | 1.35 | 1.78 | 2.28 |
| $\alpha\alpha\alpha\alpha$ | 29% | 1.10 | 0.82 | 0.89 | 1.18 |
| $\beta\beta\beta\alpha$ | 25% | 0.85 | 0.94 | 1.10 | 0.98 |
| $\beta\alpha\beta\alpha$ | 23% | 1.11 | 1.18 | 0.94 | 1.48 |
| $\alpha\beta\alpha\beta$ | 21% | 1.21 | 1.23 | 0.99 | 1.39 |
| $\alpha\beta\alpha\alpha$ | 21% | 1.00 | 0.95 | 0.81 | 1.00 |
| ... | ... | ... | ... | ... | ... |

76 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

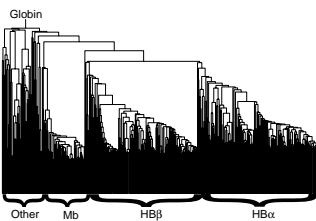
Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection
- Function Classification
- Cross-tabulation, folds and functions

An Issue with Fold Counting: Biases in the Databanks

- Over-representation of certain species and functions in the databanks (e.g. human v. plant globins, Ig's)
 - Nevertheless HI top-10 like eubacterial top-10
- PDB small, biased sample of genome (6-12%)
- Diff. numbers with diff. comparison sensitivity
 - FASTA, HMM, &c
 - Some Correction with Seq. Weighting, Diff. Sampling
 - Uniform sampling is better than high sensitivity for some and low for others (ψ -blast problem)
 - Best to avoid FPs than FNs for Venn

| Example Structure (PDB) | Fold Name | Percentage of known folds in genome | Rank in eubacterial Top-10 |
|---|---|-------------------------------------|----------------------------|
| Top-10 in a bacterial genome (H. influenzae) | | | |
| 2HSB-A | Rossmann Fold (NAD binding) | 9.6 | 1 |
| 1AKE-A | NTP Hydrolases containing P-loop | 5.7 | 2 |
| 1RCP | Flavodoxin-like | 5.1 | 3 |
| 6TDM-B | TIM-barrel | 4.5 | 4 |
| 1EXD | Ferredoxin-like | 4.2 | 5 |
| 2EN2 | like Ribonuclease H | 3.0 | 16 |
| 1SBP | like Periplasmic binding protein (class II) | 3.0 | 17 |
| 2DK1 | like Periplasmic binding protein (class I) | 3.0 | 19 |
| 1SRV-Y | Class II aaRS and biotin synthetases | 2.7 | 50 |
| 1PYP | OB-fold | 2.7 | 51 |



Same Issues with Real US Census!!
Sampling

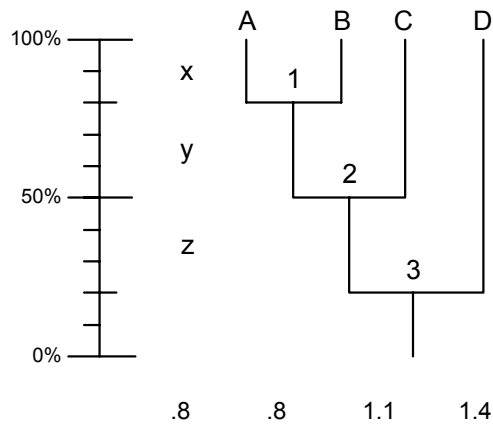
- Databank has biases.
- Assuming "fair" distribution spreads sequences uniformly through "space", want to weight sequences:

- ◊ over-represented, down (mammal)
- ◊ under-represented, up (plant & NV)

- Weights derived from a tree

- ◊ Length of an unshared branch is allotted directly to sequence
- ◊ Length of a shared branch is divided proportionally among sequences

Using a Tree to Correct for Biases



Other schemes (Argos, Sander)

Class Notes

Paper topics due by end of week
Brief email to MG, DG, JS
(1 sentence to 1 paragraph)
We'll respond with a thumbs up or down

Probably won't get to simulation

New datamining notes

Different Perspectives on Protein Thermostability

In depth focus on single molecule vs. broad view of many (all?) proteins. Anecdotal vs. Comprehensive (the genomic perspective)

Change in entropy of unfolded state in engineering of TLP (disulfides)

Proc. Natl. Acad. Sci. USA
Vol. 96, pp. 2088-2092, March 1999
Biochemistry

Engineering an enzyme to resist boiling

BRUNUS VAN DEN BURG^{1,2}, GERT VERHEIJ², OENE R. VELDMAN², GERARD VENEEMA², AND VINCENT G. H. EBRINK³

¹Department of Genetic Engineering, Biotechnology Research and Development, University of Leiden, Leiden, 3720 ZB, The Netherlands; ²Nationaal Natuurhistorisch Laboratorium, 6525 XD, Wageningen, The Netherlands; ³Department of Biotechnology, University of Twente, Enschede, 7500, The Netherlands

Communicated by Bruce W. Matthews, University of Oregon, Eugene, OR, December 23, 1997 (received for review October 23, 1997)

ABSTRACT In recent years, many efforts have been made to isolate enzymes from extremophilic organisms in the hope to reveal the structural basis for thermostability and to obtain hyperstable biocatalysts. Here we show how a moderately stable enzyme in thermophilic proteins from *Methanocaldococcus* (TLP) can be made hyperstable by a limited number of mutations. The mutational strategy included replacing residues in TLP not by residues found at equivalent positions in naturally occurring more thermostable variants, as well as rationally designed mutations. Thus, an extremely stable *Kid* mutant enzyme was obtained that was able to function at 100°C and in the presence of denaturing agents. This *Kid* mutant contained a relatively large number of mutations whose stabilizing effect is generally

items (13, 14), the introduction of a salt bridge (15), and the introduction of a disulfide bridge (16). These latter mutations were not based on comparison of the sequence and structure of naturally occurring TLPs, but on rational design.

In the present report we describe the construction of a TLP-size variant in which five TLP-size → thermophilin mutations (MT, TSA, GSA, TSE, LAPP) were combined with SGP and a disulfide bridge between residues 60 and 8 in the catalytic β-strays (E1). It is shown that this engineered enzyme resembles proteins isolated from thermophilic *Archaea* and *Eubacteria* in terms of its resistance to high temperatures and denaturing agents. It is also shown that conferring extreme stability to TLP-size did not result in major changes in catalytic properties. These findings are discussed in

Proc. Natl. Acad. Sci. USA
Vol. 96, pp. 12306-12310, October 1999
Biochemistry

Protein thermostability above 100°C: A key role for ionic interactions

Superthermophilic, two-pyruvate carboxylase (phosphoenolpyruvate carboxylase)

CHRISTOPHER WILSON¹, DEBRA L. MANNING¹, NADIA TRILAKSHY¹, KAREN S. E. YIP¹, THOMAS J. WILLIAMS¹, K. LINDA BOSTIAN², DAVID W. ROY², HEATH H. KILPATRICK², AND FRANK T. BEGG¹

¹Department of Biochemistry, University of Western Australia, Perth, 6009, Australia; ²Department of Biochemistry, University of Queensland, St. Louis, 4072, Australia

Submitted by the F. Aron, Medical Research Council, Cambridge, United Kingdom and approved August 19, 1999 (received for review July 22, 1999)

ABSTRACT The discovery of superthermophilic microorganisms has indicated that significant increases of stability may be achieved in proteins from mesophiles by the inclusion of "rigidifying" mutations. However, we believe that the protein from the hyperthermophile has been evolved against edge values of this work. Our work has focused on hydrogen bond and ionic interactions in the hyperthermophilic enzyme. We have indicated that significant increases of stability may be achieved in proteins from mesophiles by the inclusion of "rigidifying" mutations. However, we believe that the protein from the hyperthermophile has been evolved against edge values of this work. Our work has focused on hydrogen bond and ionic interactions in the hyperthermophilic enzyme.

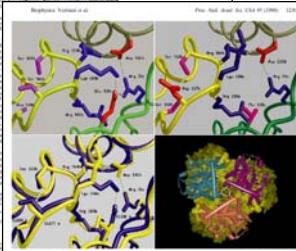
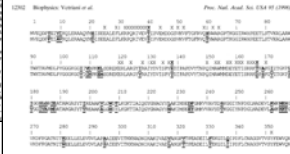
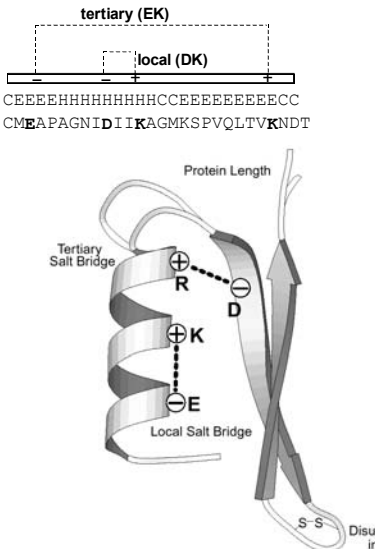


Fig. 1. The structure of the hyperthermophilic enzyme. The structure is shown as a ribbon diagram. The structure is shown as a ribbon diagram. The structure is shown as a ribbon diagram. The structure is shown as a ribbon diagram.

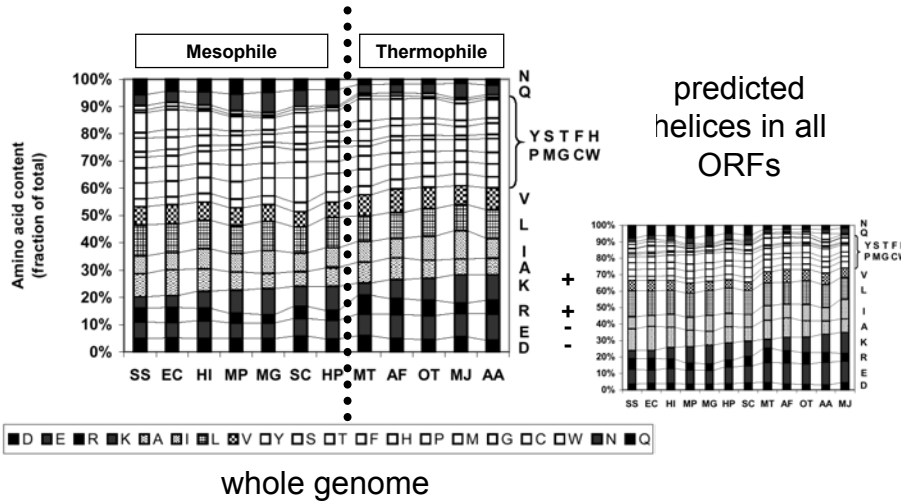
Thermostability: Analyzing a few Factors with Genome Comparison



| Organism | Category | Genome Abbreviation | # of Proteins | Physiological condition |
|--|---------------------------|---------------------|---------------|-------------------------|
| <i>Pyrococcus horikoshii</i> (Strain OT3) (Kawarayashi et al., 1998) | archaea | OT | 2061 | 98°C anaerobe |
| <i>Aquifex aeolicus</i> (Decker et al., 1998) | eubacteria, gram negative | AA | 1522 | 95°C |
| <i>Methanococcus janaschii</i> (Bull et al., 1996) | archaea | MJ | 1735 | 85°C anaerobe |
| <i>Archaeoglobus fulgidus</i> (Klenk et al., 1997) | archaea | AF | 2409 | 83°C anaerobe |
| <i>Methanobacterium thermoautotrophicum</i> (Smith et al., 1997) | archaea | MT | 1869 | 65°C anaerobe |
| <i>Haemophilus influenzae</i> (Fleischmann et al., 1995) | eubacteria, gram negative | HI | 1680 | mesophilic temp. |
| <i>Mycoplasma genitalium</i> (Fraser et al., 1995) | eubacteria, gram positive | MG | 470 | mesophilic temp. |
| <i>Mycoplasma pneumoniae</i> (Himmelreich et al., 1996) | eubacteria, gram positive | MP | 677 | mesophilic temp. |
| <i>Helicobacter pylori</i> (Tomb et al., 1997) | eubacteria, gram negative | HI | 1590 | mesophilic temp. |
| <i>Escherichia coli</i> (Blattner et al., 1997) | eubacteria, gram negative | EC | 4288 | mesophilic temp. |
| <i>Synechocystis</i> sp. (Kuroki et al., 1995) | cyanobacteria | SS | 3168 | mesophilic temp. |
| <i>Saccharomyces cerevisiae</i> (Goffeau et al., 1997) | eukaryote, fungus | SC | 6218 | mesophilic temp. |

Composition Analysis of the Proteome

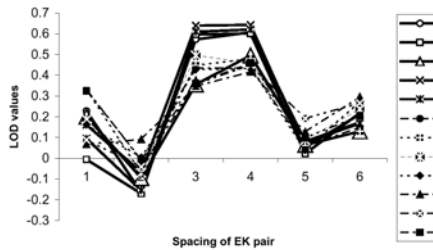
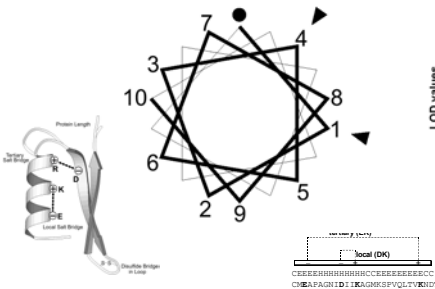
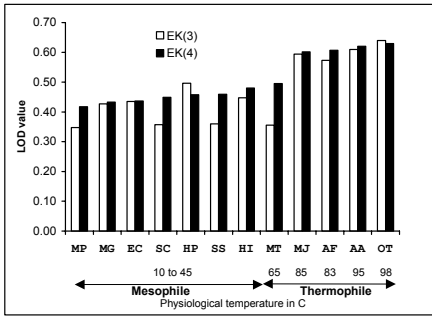
More Charged Residues in Thermophiles, Suggestive of Salt Bridges



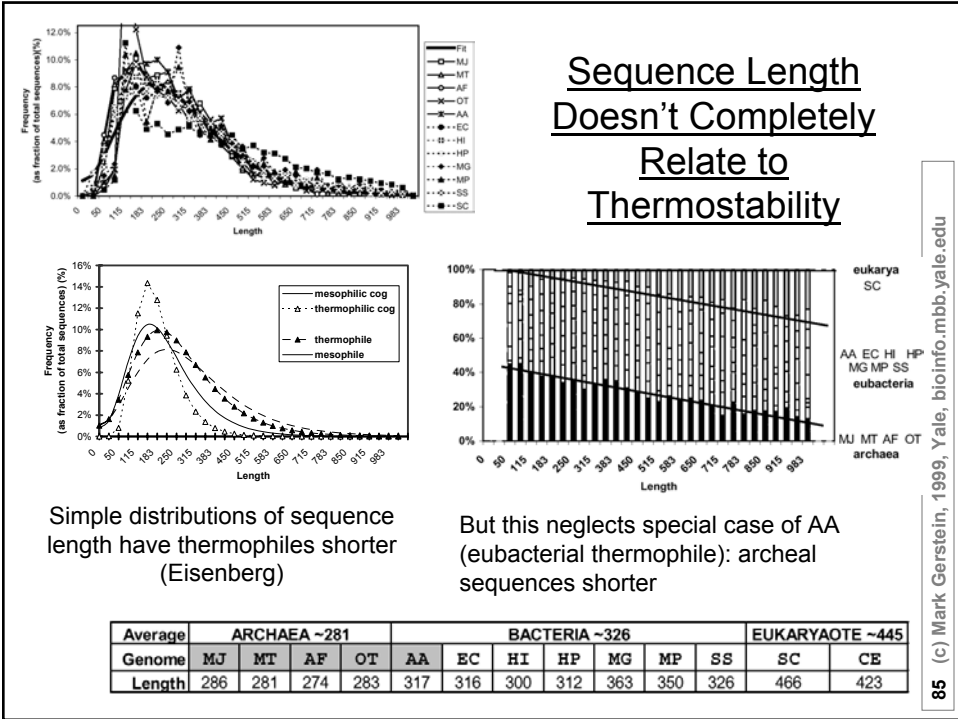
83 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

1-4 Spacing of Charged Residues More than Expected in Thermophile Helices ⇒ Salt Bridges

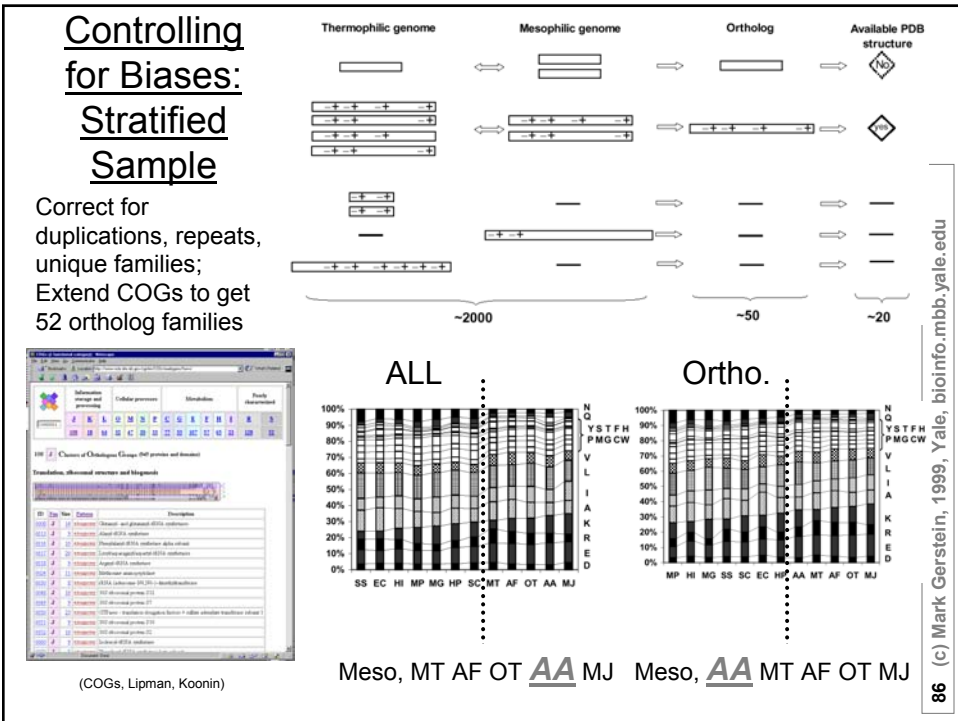
Quantify with LOD score
 $LOD = \log(\text{observed/expected})$
 For inst.,
 $\text{expected}[EK(4)] \sim f(E) \cdot f(K)$
 $LOD > 0$, greater than expected



84 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu



85 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

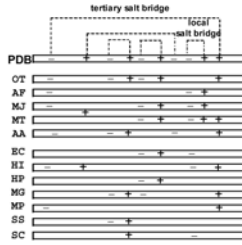


86 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Controls II: Known Structures, Random Genomes

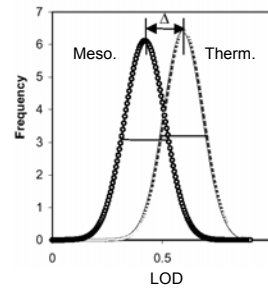
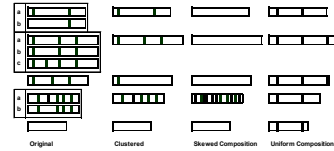
3D Structures

For orthologs of known structure: map tertiary salt bridges onto multiple alignment and look at conservation in Therm. vs. Meso.



| COG | | Cat. | PDB | Therm. Avg. SB | Meso Avg. SB | Diff | |
|-----|---|------------|------|----------------|--------------|------|---|
| 49 | J | ribosomal | 1fss | 5.6 | 3.1 | 3 | + |
| 80 | J | ribosomal | 1aci | 0.8 | 0.7 | 0.1 | |
| 81 | J | ribosomal | 1ac2 | 6.4 | 4.3 | 2.1 | + |
| 91 | J | ribosomal | 1axe | 1.8 | 0.9 | 0.9 | |
| 92 | J | ribosomal | 1awi | 3 | 1.9 | 1.1 | + |
| 96 | J | ribosomal | 1asa | 2 | 2.1 | -0.1 | |
| 98 | J | ribosomal | 1pkp | 0.6 | 1.7 | -1.1 | - |
| 184 | J | ribosomal | 1a32 | 1.8 | 1.9 | -0.1 | |
| 186 | J | ribosomal | 1fip | 0.4 | 0.9 | -0.5 | |
| 16 | J | synthetase | 1pys | 7.6 | 2.6 | 5 | + |
| 124 | J | synthetase | 1ady | 9.6 | 6.1 | 3.5 | + |
| 162 | J | synthetase | 2z1 | 3.8 | 3.3 | 0.5 | |
| 30 | J | other | 1yub | 5 | 5.3 | -0.3 | |
| 125 | F | other | 1fmk | 0.8 | 0.4 | 0.4 | |
| 149 | C | other | 1bim | 3 | 4.3 | -1.3 | - |
| 541 | N | other | 1fs | 3.6 | 3.4 | 0.2 | |
| 112 | E | other | 1cg0 | 6.2 | 4.6 | 1.6 | + |
| 552 | N | other | 1fh | 4.2 | 4.6 | -0.4 | |

Random Sampling: Make up random thermo. and meso. genomes, see what distribution of each statistic is



87 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

End of class on 11.27

88 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

How Representative are the Known Structures of the Proteins in a Complete Genome? The issue of Bias

Assess 2° TM predictions

(+) comprehensive, statistical

(-) predictions inaccurate

(~65%)

(-) extrapolate from PDB (esp. TM), domain problem

Is prediction (extrapolation) based on known structures justified?

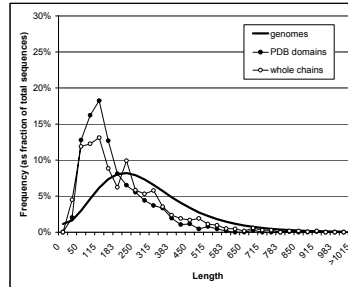
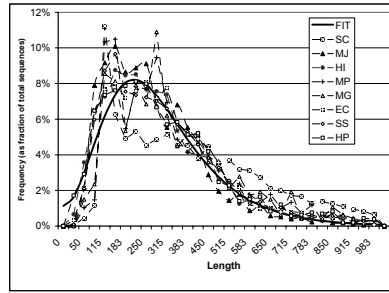
Length: Genomes Sequences are longer than those in Known Structures

340 aa for avg. genome seq.

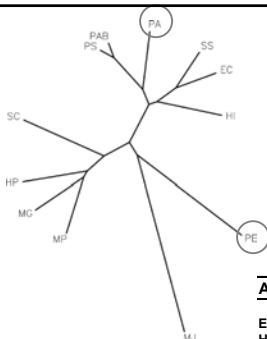
(470 aa for yeast)

205 aa for PDB chain

~160 aa for PDB domain



89 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu



Amino Acid Composition

How Representative are the Known Structures of the Proteins in Complete Genome?

| Name | Soluble PDB | = all-β | + all-α |
|------|-------------|---------|---------|
| A | 8.40% | 6.8% | 9.2% |
| C | 1.72% | 1.6% | 1.4% |
| D | 5.91% | 5.9% | 5.8% |
| E | 6.29% | 5.2% | 7.3% |
| F | 3.94% | 4.2% | 4.2% |
| G | 7.79% | 8.4% | 6.4% |
| H | 2.19% | 2.1% | 2.2% |
| I | 5.54% | 5.4% | 5.1% |
| K | 6.02% | 5.6% | 6.5% |
| L | 8.37% | 7.3% | 9.6% |
| M | 2.15% | 1.7% | 2.4% |
| N | 4.57% | 5.3% | 4.4% |
| P | 4.70% | 5.1% | 4.4% |
| Q | 3.73% | 3.6% | 4.2% |
| R | 4.78% | 4.2% | 5.4% |
| S | 5.97% | 7.2% | 5.7% |
| T | 5.87% | 7.2% | 5.2% |
| V | 6.96% | 7.6% | 5.7% |
| W | 1.46% | 1.7% | 1.5% |
| Y | 3.64% | 3.8% | 3.5% |

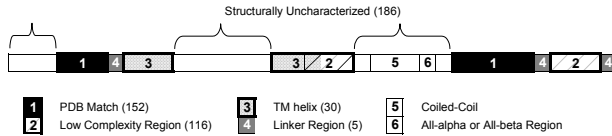
| Abs. | rms | K | I | C | Q | W | N | F | L | G | A | P | S | R | H | M | E | D | T | Y | V |
|------|-----|------|------|-----|-----|-----|-----|-----|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| EC | | 4.4 | 6.0 | 1.2 | 4.4 | 1.6 | 4.0 | 3.9 | 10.6 | 7.4 | 9.5 | 4.4 | 5.8 | 5.6 | 2.3 | 2.8 | 5.7 | 5.1 | 5.4 | 2.9 | 7.1 |
| HI | | 6.3 | 7.1 | 1.0 | 4.6 | 1.1 | 4.9 | 4.5 | 10.5 | 6.6 | 8.2 | 3.7 | 5.8 | 4.5 | 2.1 | 2.4 | 6.5 | 5.0 | 5.2 | 3.1 | 6.7 |
| SS | | 4.2 | 6.3 | 1.0 | 5.6 | 1.6 | 4.0 | 4.0 | 11.4 | 7.4 | 8.5 | 5.1 | 5.8 | 5.1 | 1.9 | 2.0 | 6.0 | 5.0 | 5.5 | 2.9 | 6.7 |
| SC | | 7.3 | 6.6 | 1.3 | 3.9 | 1.0 | 6.1 | 4.5 | 9.6 | 5.0 | 5.5 | 4.3 | 9.0 | 4.5 | 2.2 | 2.1 | 6.5 | 5.8 | 5.9 | 3.4 | 5.6 |
| HP | | 8.9 | 7.2 | 1.1 | 3.7 | 7 | 5.9 | 5.4 | 11.2 | 5.8 | 6.8 | 3.3 | 6.8 | 3.5 | 2.1 | 2.2 | 6.9 | 4.8 | 4.4 | 3.7 | 5.6 |
| MP | | 8.6 | 6.6 | 1.8 | 5.4 | 1.2 | 6.2 | 5.6 | 10.3 | 5.5 | 6.7 | 3.5 | 6.5 | 3.5 | 1.8 | 1.6 | 5.7 | 5.0 | 6.0 | 3.2 | 6.5 |
| MG | | 9.5 | 8.2 | 1.8 | 4.7 | 1.0 | 7.5 | 6.1 | 10.7 | 4.6 | 5.6 | 3.0 | 6.6 | 3.1 | 1.6 | 1.5 | 5.7 | 4.9 | 5.4 | 3.2 | 6.1 |
| MJ | | 10.4 | 10.5 | 1.3 | 1.5 | 7 | 5.3 | 4.2 | 9.5 | 6.3 | 5.5 | 3.4 | 4.5 | 3.8 | 1.4 | 2.2 | 8.7 | 5.5 | 4.0 | 4.4 | 6.9 |
| AVG | | 7.5 | 7.3 | 1.1 | 4.2 | 1.1 | 5.5 | 4.8 | 10.5 | 6.1 | 7.0 | 3.8 | 6.4 | 4.2 | 1.9 | 2.1 | 6.5 | 5.1 | 5.2 | 3.3 | 6.4 |
| SD | | 2.3 | 1.4 | 2 | 1.3 | 3 | 1.2 | 1.8 | 7 | 1.0 | 1.5 | 1.7 | 1.3 | 1.9 | 3 | 4 | 1.0 | 3 | 7 | 5 | 6 |

| Diff. | EC | HI | SS | SC | HP | MP | MG | MJ | AVG | RMS | | | | | | | | | | |
|-------|-----|----|-----|-----|-----|-----|----|----|-----|-----|-----|-----|-----|-----|-----|----|-----|-----|-----|-----|
| 16 | -25 | 8 | -29 | 19 | 7 | -15 | -2 | 28 | -6 | 13 | -5 | -3 | 16 | 3 | 28 | -7 | -14 | -7 | -22 | 1 |
| 17 | 8 | 27 | -38 | 24 | -21 | 6 | 12 | 26 | -15 | -2 | -20 | -2 | -6 | -7 | 10 | 5 | -17 | -11 | -14 | -4 |
| 20 | -29 | 13 | -39 | 49 | 9 | -13 | 1 | 37 | -6 | 1 | 11 | -3 | 6 | -15 | -8 | -2 | -16 | -6 | -20 | -4 |
| 21 | 24 | 18 | -21 | 5 | -27 | 31 | 14 | 15 | -36 | -34 | -7 | 51 | -7 | -2 | 4 | 5 | -4 | 0 | -8 | -20 |
| 27 | 52 | 29 | -34 | 0 | -61 | 27 | 36 | 34 | -26 | -18 | -29 | 14 | -28 | -4 | 2 | 11 | -20 | -25 | 1 | -20 |
| 28 | 45 | 18 | -56 | 44 | -17 | 35 | 41 | 24 | -29 | -20 | -25 | 8 | -27 | -18 | -28 | -8 | -17 | 2 | -11 | -7 |
| 36 | 61 | 48 | -50 | 27 | -32 | 62 | 53 | 28 | -41 | -33 | -36 | 11 | -35 | -28 | -30 | -8 | -18 | -8 | -11 | -12 |
| 38 | 77 | 88 | -23 | -61 | -49 | 14 | 6 | 14 | -19 | -35 | -28 | -25 | -20 | -35 | 1 | 40 | -3 | -31 | 20 | -2 |
| AVG | 26 | 31 | -36 | 13 | -23 | 19 | 20 | 26 | -22 | -16 | -17 | 6 | -13 | -13 | -4 | 4 | -14 | -11 | -8 | -9 |
| RMS | 45 | 39 | 38 | 35 | 31 | 30 | 28 | 27 | 25 | 24 | 23 | 21 | 21 | 18 | 18 | 16 | 15 | 15 | 15 | 11 |

90 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

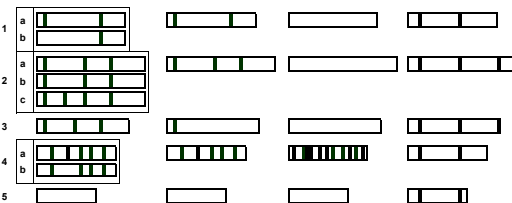
Composition of Different Regions of Genomes

- Are composition differences uniform?
- Resampling
- Non-globular regions differ most in occurrence and composition
- Remove Repetitive Regions (SEG)



| Statistics for Amino Acids | AVG | SD | EC | HI | HP | MG | MJ | MP | SC | SS |
|----------------------------|--------|-------------|---------|--------|--------|--------|--------|--------|---------|---------|
| Total Number | 775998 | | 1358465 | 505279 | 500616 | 170400 | 497968 | 237905 | 2900670 | 1033450 |
| Fraction Masked by... | | | | | | | | | | |
| PDB Match | 8.7% | 3.7% | 11.1% | 13.7% | 8.8% | 12.9% | 7.1% | 9.7% | 6.2% | 9.0% |
| Non-globular Region | 21.7% | 6.9% | 16.7% | 13.9% | 22.2% | 28.2% | 35.1% | 24.7% | 23.9% | 20.5% |
| TM-helix | 4.9% | 1.4% | 7.3% | 6.1% | 4.8% | 3.8% | 2.9% | 4.5% | 5.2% | 5.9% |
| Linker Region | 5.1% | 0.4% | 5.3% | 4.8% | 4.8% | 5.0% | 5.0% | 5.2% | 4.6% | 5.1% |
| Fraction Remaining | | | | | | | | | | |
| Uncharacterized | 59.7% | 8.9% | 59.6% | 61.5% | 59.4% | 50.2% | 49.9% | 55.8% | 60.0% | 59.6% |

| | AVG | SD | EC | HI | HP | MG | MJ | MP | SC | SS |
|------------------------|------------|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Overall | 23% | 10% | 16% | 17% | 27% | 36% | 38% | 28% | 21% | 20% |
| PDB Match | 18% | 9% | 12% | 14% | 24% | 27% | 34% | 20% | 12% | 15% |
| Non Globular Region | 36% | 13% | 32% | 33% | 39% | 50% | 52% | 40% | 42% | 35% |
| TM-helix | 49% | 15% | 55% | 53% | 55% | 67% | 55% | 56% | 56% | 51% |
| Linker Region | 27% | 10% | 22% | 24% | 29% | 39% | 33% | 35% | 21% | 25% |
| Uncharacterized Region | 23% | 6% | 15% | 17% | 26% | 34% | 32% | 27% | 20% | 19% |



| PDB | Select | length | class | name |
|------|---------|--------|-------|-------------------|
| lsty | - | 137 | β | Staph nuclease |
| lcbp | a:9-137 | 129 | β | CAP |
| lbgh | - | 85 | β | Gene V protein |
| lpht | - | 83 | β | SH3 domain |
| ltpf | a: | 250 | α/β | TIM |
| lway | a: | 248 | α/β | Trp Synthase |
| 8dfr | - | 186 | α/β | DHFR |
| 2zn2 | - | 155 | α/β | Ribonuclease H |
| lhrs | d: | 87 | α/β | Barstar |
| lgsb | - | 185 | α+β | Hen Lysozyme |
| 1191 | - | 162 | α+β | T4 lysozyme |
| 1931 | - | 129 | α+β | alpha-Lactalbumin |
| 7rea | - | 124 | α+β | RNAse A |
| lbrn | l: | 108 | α+β | Barnase |
| lfkd | - | 107 | α+β | FK506 |
| 9rnt | - | 104 | α+β | RNAse T1 |
| lsha | a: | 103 | α+β | SH2 domain |
| lubi | - | 76 | α+β | Ubiquitin |
| lcse | l: | 63 | α+β | Cl-2 inhibitor |
| ligd | - | 61 | α+β | B1 domain |
| lmbd | - | 153 | α | Globin |
| lhrc | - | 105 | α | Cytochrome c |
| 2wzp | r: | 104 | α | Trp Repressor |
| 1111 | a: | 89 | α | Cro Repressor |
| lcbp | d: | 66 | α | Lambda Repressor |
| lcpo | - | 61 | α | ROP |
| lmyk | a: | 47 | α | Arc Repressor |
| 2zta | a: | 31 | α | GCN4 zipper |
| lbt1 | - | 263 | M | beta-Lactamase |
| lbp1 | - | 58 | S | BPTI |
| AVG | | 116 | | |

| Name | Hydroph. Polar | Soluble PDB | biophys. proteins | Rel. Diff. |
|------|----------------|-------------|-------------------|------------|
| | | PS | BP | BP/PS -1 |
| P | H | 4.7% | 3.7% | -21% |
| F | H | 4.0% | 3.2% | -19% |
| M | H | 2.1% | 1.8% | -16% |
| D | P | 6.0% | 5.1% | -16% |
| V | H | 7.0% | 6.2% | -12% |
| C | H | 1.7% | 1.5% | -9% |
| S | P | 6.0% | 5.7% | -5% |
| G | . | 7.8% | 7.7% | -1% |
| I | H | 5.6% | 5.5% | -1% |
| N | P | 4.6% | 4.6% | 0% |
| W | H | 1.4% | 1.5% | 1% |
| T | P | 5.8% | 6.0% | 2% |
| L | H | 8.4% | 8.7% | 5% |
| A | . | 8.4% | 8.8% | 6% |
| Y | . | 3.7% | 3.9% | 6% |
| H | P | 2.2% | 2.4% | 6% |
| Q | P | 3.7% | 4.0% | 6% |
| R | P | 4.8% | 5.2% | 9% |
| E | P | 6.2% | 7.0% | 13% |
| K | P | 5.9% | 7.7% | 30% |

Biophysical Proteins

Proteins that inform our view of the folding process -- as compared to the PDB.

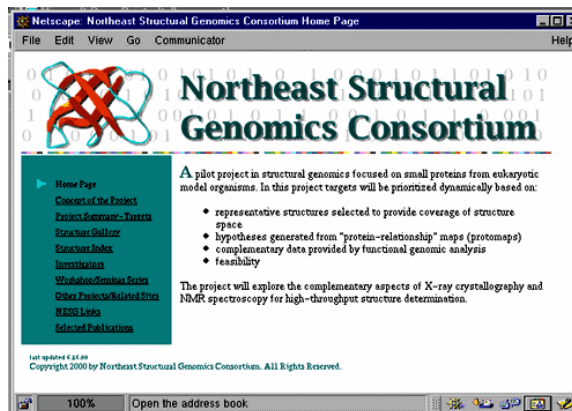
Shorter (116 v 161)

Fewer hydrophobes

Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection
- Function Classification
- Cross-tabulation, folds and functions

nesg.org

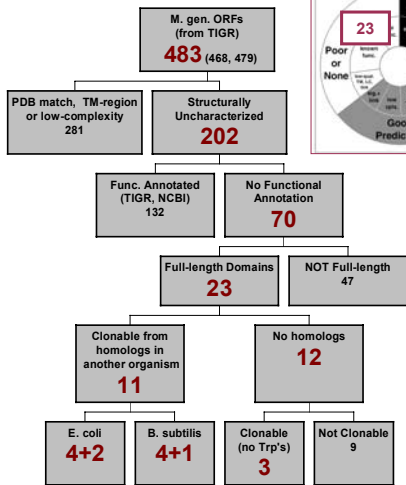


G Montelione

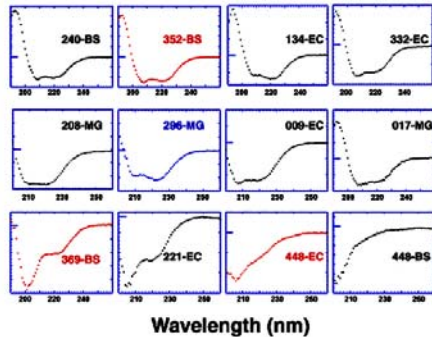
Finding Unusual Proteins for Expt. Structural Genomics

Prospective Target Selection

- Identify Proteins in *M. genitalium* that are most atypical structurally (hardest)
- Characterize biophysically by CD (do they fold normally?)



L Regan



95 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Tracking Database

| ORF ID | PDB ID | Organism | Target ORF | Clonal | Expression | Purified | RSDC spectrum | CD long-term | Crystal | X-ray data | X-ray structure | NMR long-term | NMR structure | Electrostatics |
|--------|--------|----------|------------|--------|------------|----------|---------------|--------------|---------|------------|-----------------|---------------|---------------|----------------|
| WH0415 | | Coli | | | | | | | | | | | | |
| TT11 | | Mbe | MF7110 | | | | | | | | | | | |
| TT2 | | Mbe | MF7102 | | | | | | | | | | | |
| TT3 | | Mbe | MF7109 | | | | | | | | | | | |
| TT4 | | Mbe | MF7100 | | | | | | | | | | | |
| TT5 | | Mbe | MF7101 | | | | | | | | | | | |
| TT2 | 1ak | Mbe | mf1501 | | | | | | | | | | | |
| TT10 | 1ak | Mbe | mf1512 | | | | | | | | | | | |
| TT11 | | Mbe | mf1500 | | | | | | | | | | | |
| TT12 | | Mbe | mf1502 | | | | | | | | | | | |
| TT13 | | Mbe | mf1504 | | | | | | | | | | | |

| ORF ID | Construct | Metadata |
|--------|-----------|----------|
| WH0415 | WH0415 | WH0415 |
| TT11 | TT11 | TT11 |
| TT2 | TT2 | TT2 |
| TT3 | TT3 | TT3 |
| TT4 | TT4 | TT4 |
| TT5 | TT5 | TT5 |
| TT2 | TT2 | TT2 |
| TT10 | TT10 | TT10 |
| TT11 | TT11 | TT11 |
| TT12 | TT12 | TT12 |
| TT13 | TT13 | TT13 |

96 (c) Mark Gerstein, 1999, Yale, bioinfo.mbb.yale.edu

Large-scale Database Surveys (contents)

- Fold Library
- Parts Lists: homologs, motifs, orthologs, folds
- Overall Sequence-structure Relationships, Annotation Transfer
- Parts in Genomes, shared & common folds
- Genome Trees
- Extent of Fold Assignment: the Bias Problem
- Bulk Structure Prediction
- The Genomic vs. Single-molecule Perspective
- Understanding Biases in Sampling
- Relationship to experiment: LIMS, target selection
- Function Classification
- Cross-tabulation, folds and functions