

# Comparing Genomes in terms of Protein Structure:

## **Surveys of a Finite Parts List**

Mark Gerstein

# Genomes highlight the Finiteness of Biology

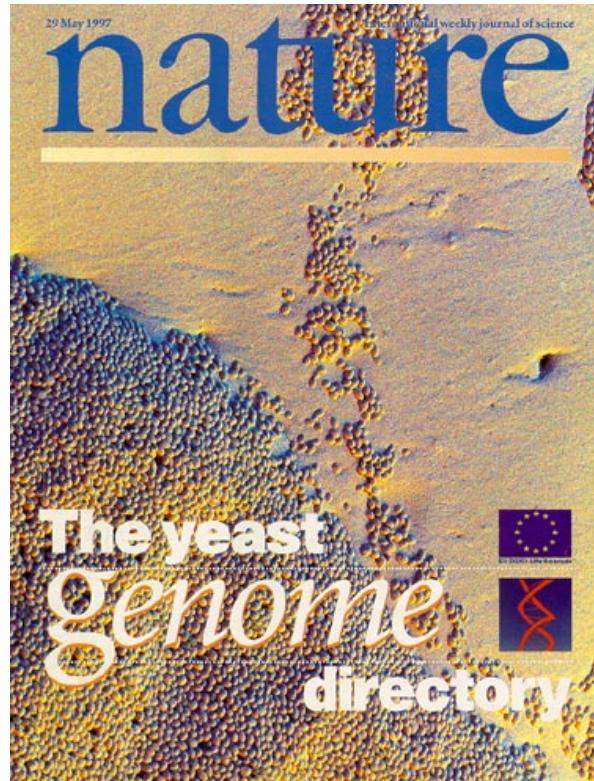
1995



Bacteria  
1.6 Mb, ~1600 genes

[Fleischmann *et al.* (1995). "Whole-genome random sequencing and assembly of *Haemophilus influenzae* rd." *Science* **269**: 496-512.]

1997



Eukaryote  
13 Mb, ~6000 genes

1998....

Microbial Genomes  
>15 completed,  
~40 underway

The Worm:  
75% of 100 Mb  
done, with ~13 K  
genes so far)

The Human:  
3 Gb & 100 K  
genes, 2003? <sub>2</sub>.

# Comparing Genomes in terms of Protein Structure: Surveys of a Finite Parts List

## 1 Library of Known Folds

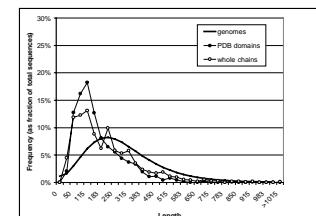
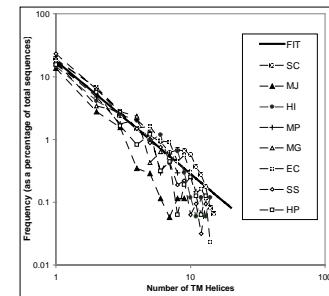
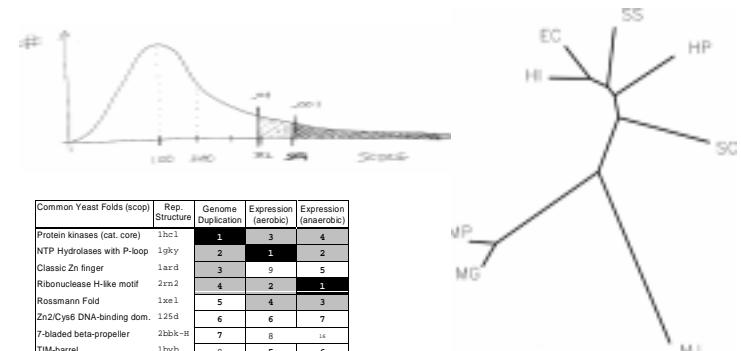
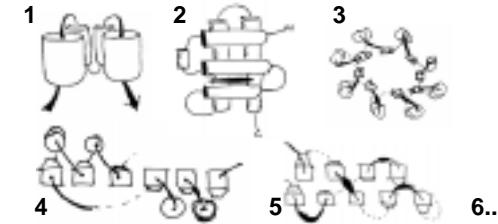
Importance of statistics. Scop auto-alignments. Significance follows EVD stats, same as sequences.

## 2 A Census of Known Folds

Which folds in which Organisms: Plants v. People, E coli v. yeast? Shared Fold Tree. Top-10 by duplication/expression, repeated  $\beta\alpha\beta$  supersecondary struc.

## 3 Prediction of Unknown Folds

Zipf law for TM's but no 7-TM's. Same 2<sup>o</sup> comp. but different a.a. comp. Can extrapolate from known structures to genomes?

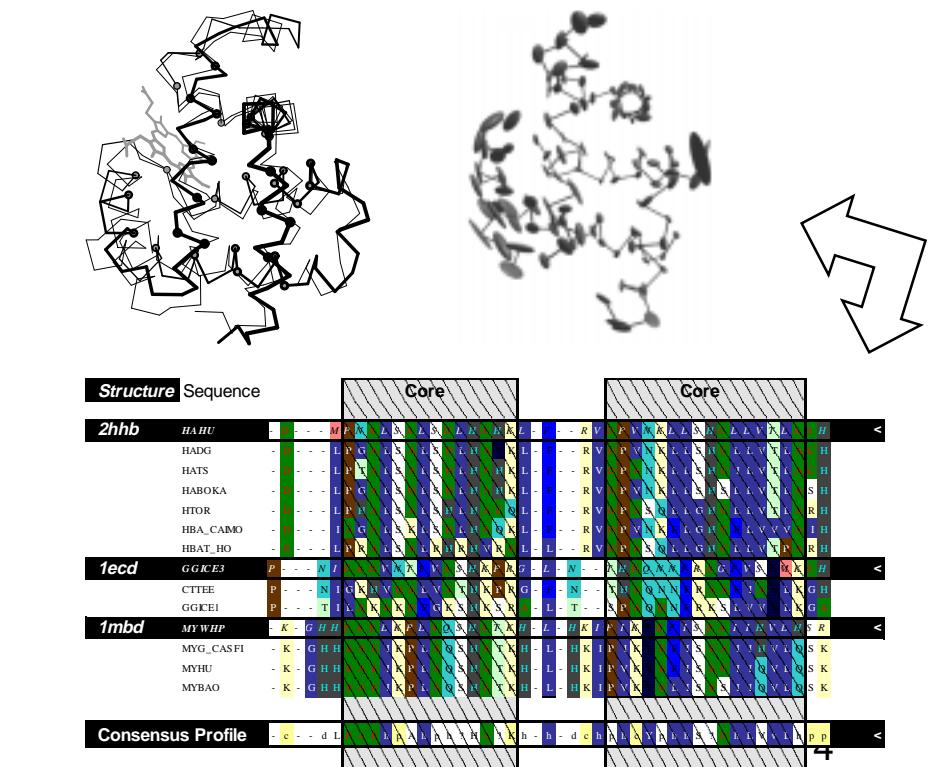
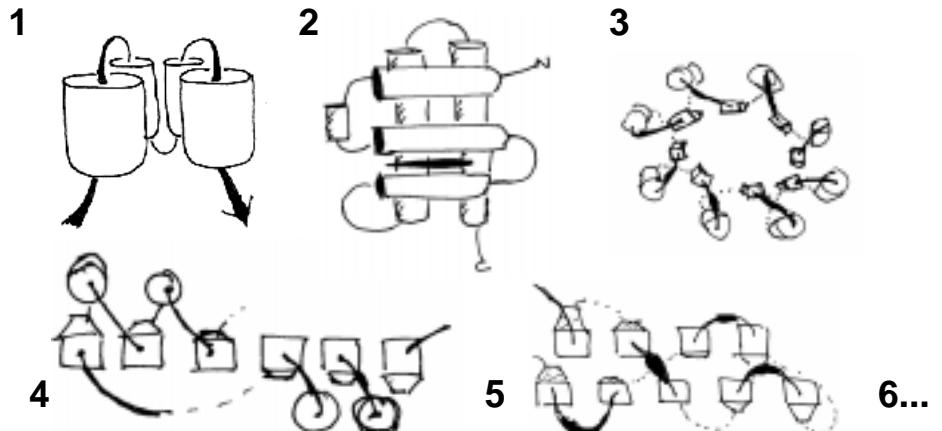


<http://bioinfo.mbb.yale.edu/genome>

**Acknowledgements:** M Levitt, scop  
(Murzin, Brenner, Ailey, Hubbard, Chothia)

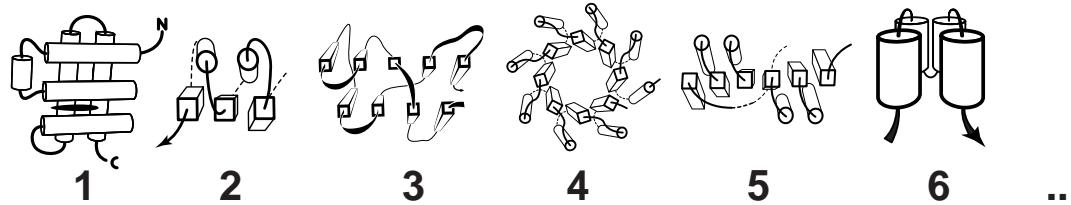
# Fold Library

- Primary way to interpret Genome Sequences in terms of Structure
- Very Limited Number of Folds (~1K-10K, Chothia)
- Elements: Domain definitions; Aligned structures, collecting together Non-homologous Sequences; Core region annotation
- Many approaches to building Library
  - ◊ Automatic: FSSP-HSSP (Sander), Entrez-MMDB (Bryant)
  - ◊ Semi-automatic: CATH (Thornton), HOMALDB (Sali)
  - ◊ **Manual (scop, Murzin)**
  - ◊ Start with Sequences: Pfam (Durbin, Eddy), COGs (Koonin, Lipman), Blocks (Henikoff), ProSite (Bairoch)



# Scale Fold Library vs. Other Fundamental Data structures

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



Folds in Molecular Biology    1000-10000

const.	mant.	exp.	unit
e	1.60	e	8 C
F	9.65	e	4 C/mol
$\epsilon_0$	8.85	e	-12 F/m
$\mu_0$	1.26	e	-6 H/m
h	6.63	e	-34 J*s
k	1.38	e	-23 J/K
$m_e$	9.11	e	-31 kg
$m_p$	1.67	e	-27 kg
$m_n$	1.68	e	-27 kg
$a_0$	5.29	e	-11 m
$\lambda_C$	2.43	e	-12 m
c	3.00	e	-19 m/s
G	6.67	e	-11 m <sup>3</sup> /kg*s <sup>2</sup>
$N_A$	6.02	e	23 mol <sup>-1</sup>

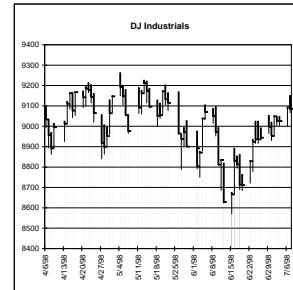
10

Physics

H	He
Li	Be
Na	Mg
K	Ca
Rb	Sr
Cs	Ba
Fr	Ra
La	Ce
Ac	Th
Sc	Ti
Lu	Rf
Y	Zr
Hf	Ta
Nb	W
Nb	Os
Mo	Ir
Tc	Pt
Ru	Au
Rh	Hg
Pd	Tl
Ag	Pb
Cd	Bi
In	Po
Zn	At
Ga	Rn
Ge	
As	
Se	
Br	
Kr	
B	C
C	N
N	O
O	F
F	Ne
Ne	
Al	Si
Si	P
P	S
S	Cl
Cl	Ar

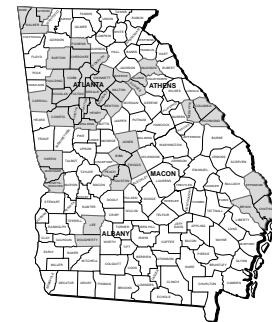
100

Chemistry



1000  
-10000

Finance



>1000000

Politics

5

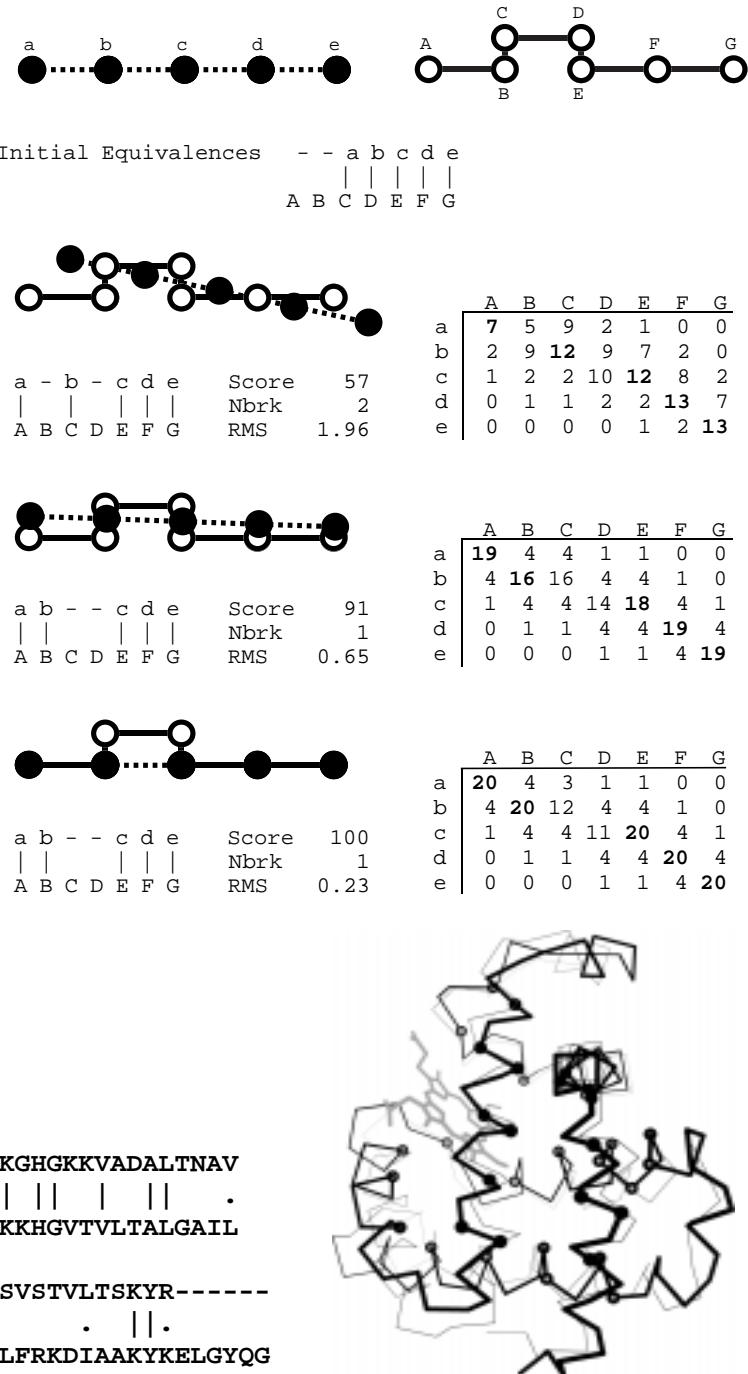
(Large than physics and chemistry, Similar to Finance (Exact Finite Number of Objects (3,056 on NYSE by 1/98), descrip. by Standardized Statistics (even abbrevs, INTC) and groups (sectors))  
Smaller than Social Surveys, Indefinite Number of People, Not Well Defined Vocabulary and statistics.

# Automatic Alignments of Scop, Focussing on Statistics of Relationships

- Our Approach
  - ◊ Iterative Dynamic Programming, like repeated sequence alignment
  - ◊ Derived from Program of G Cohen (Align)
  - ◊ **Score** =  $M_{str}(i,j) = \sum 100 / (5 + d^2)$
- Numerous other approaches to struc. alignment: SAP (Taylor), VAST (Bryant), Artymiuk, Sali, Sippl, Sander, Cohen, STAMP (Barton)

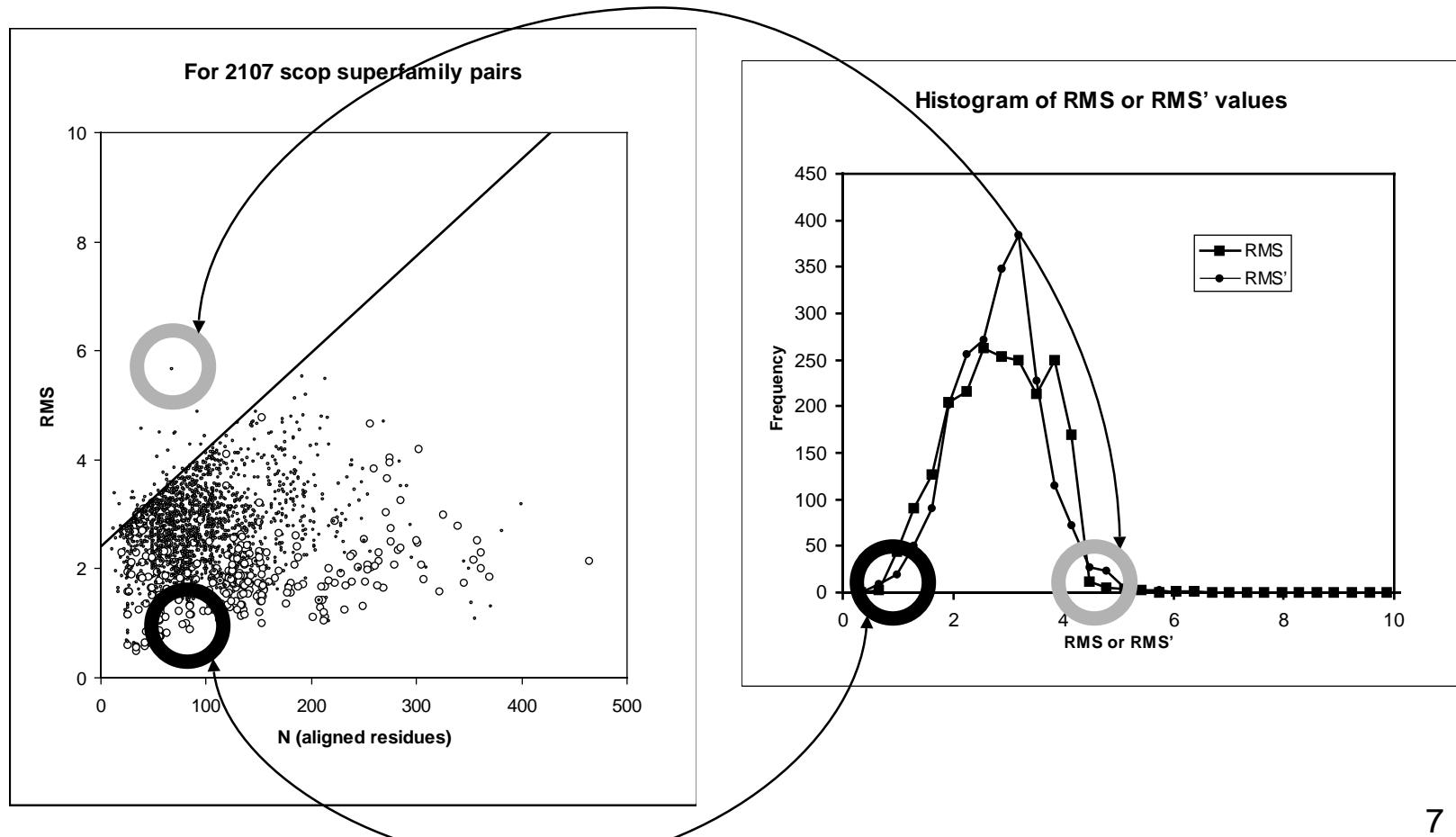
Hb VLSPADKTNVKAAGWKVGAHAGEYGAEARLMFLSFPTTKTYFPHF-DLS-----HGSAQVKGHGKKVADALTN  
 Mb VLSEGEWQLVLHVWAKVEADVAGHGQDILIRLFKSHPETLEKFDRFKHLKTEAEMKASEDLKKHGVTVLTALGAIL

Hb AHVD-DMPNALSALSDLHAHKLRDPVNFKLLSHCLLVTAAHLPAEFTPASLDKFLASVSTVLTSKYR-----  
 Mb KK-KGHHEAELKPLAQSHATKHKIPIKYLEFISEAIIHVLHSRHPGDFGADAQGAMNKALELFRKDIAAKYKELGYQG



# Statistics on Range of Similarities

For 2107 pairs, only 2% Outliers (with subtle similarity)

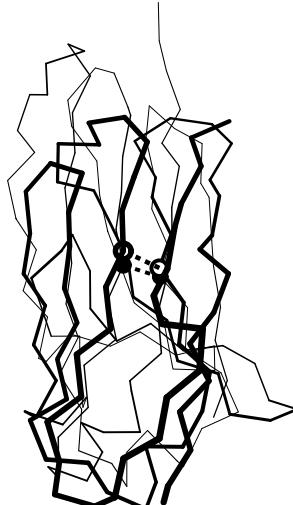
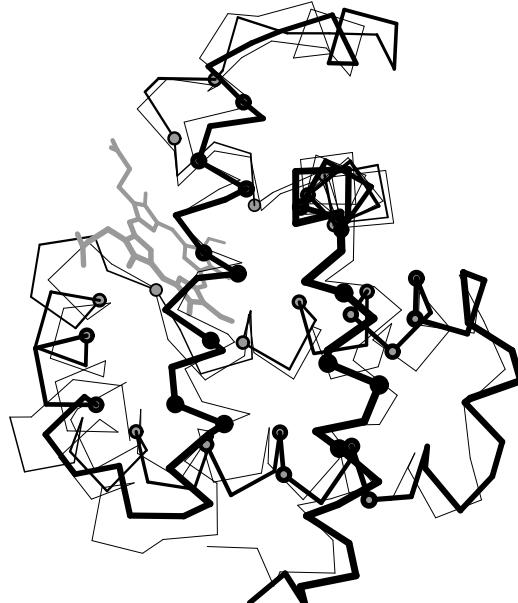


# Some Similarities are Readily Apparent others are more Subtle

Easy:  
Globins

Tricky:  
Ig C,  
Ig V,

Very Subtle: G3P-dehydro-  
genase, C-term. domain

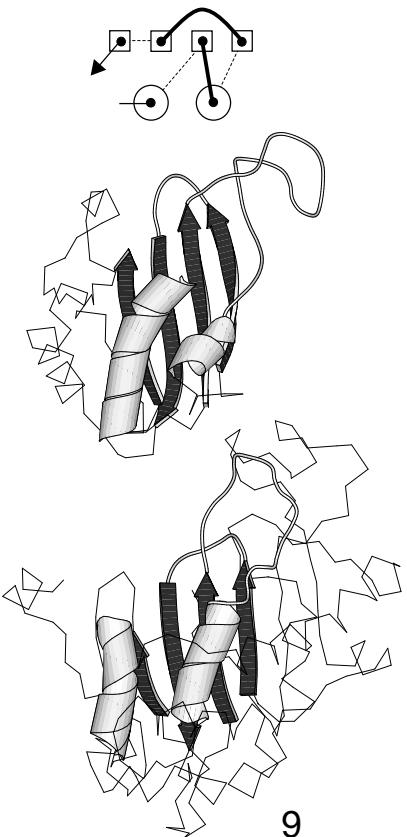
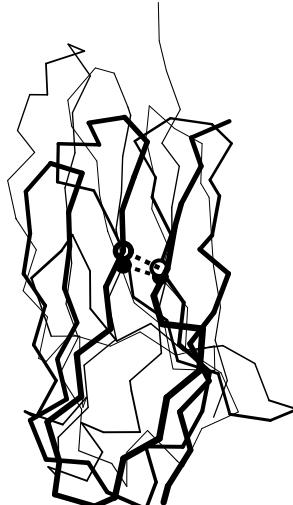
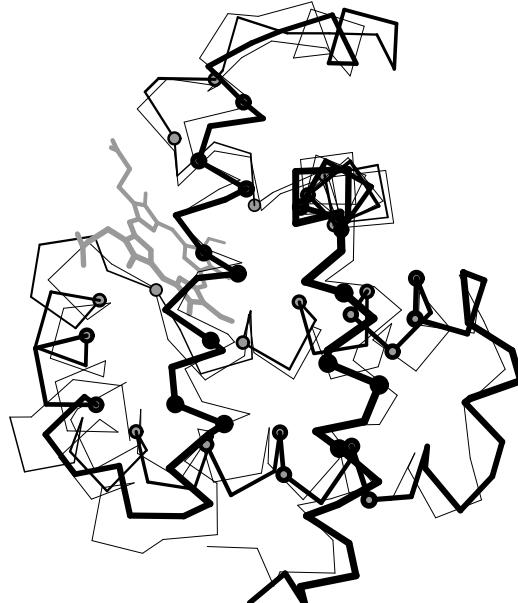


# Some Similarities are Readily Apparent others are more Subtle

Easy:  
Globins

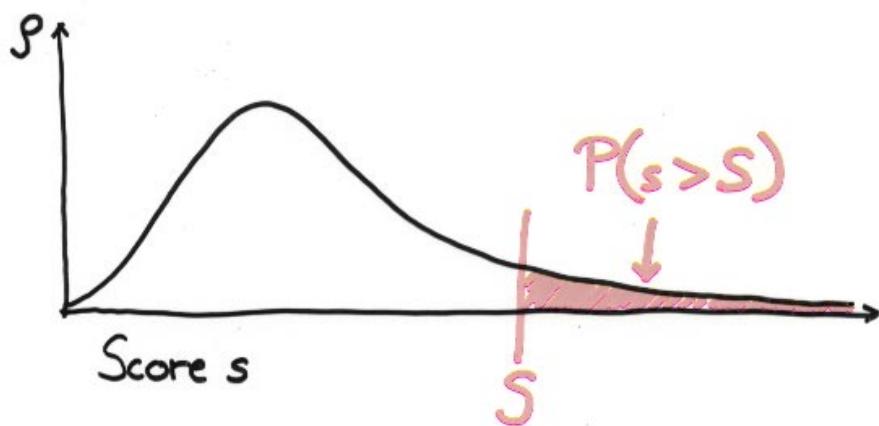
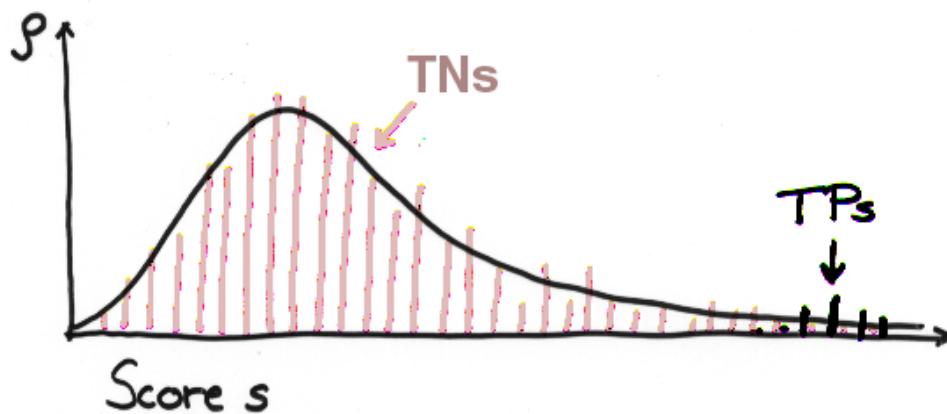
Tricky:  
Ig C,  
Ig V,

Very Subtle: G3P-dehydro-  
genase, C-term. domain



9

# Strategy: same formalism for sequence & structure



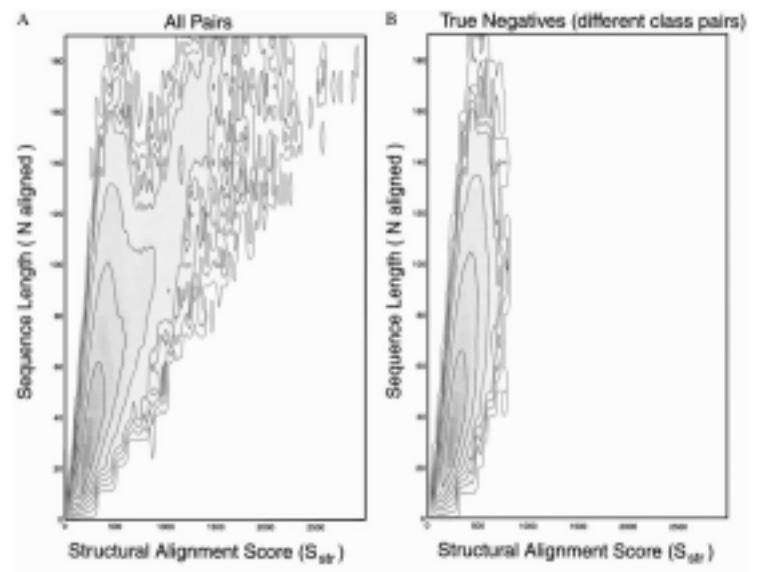
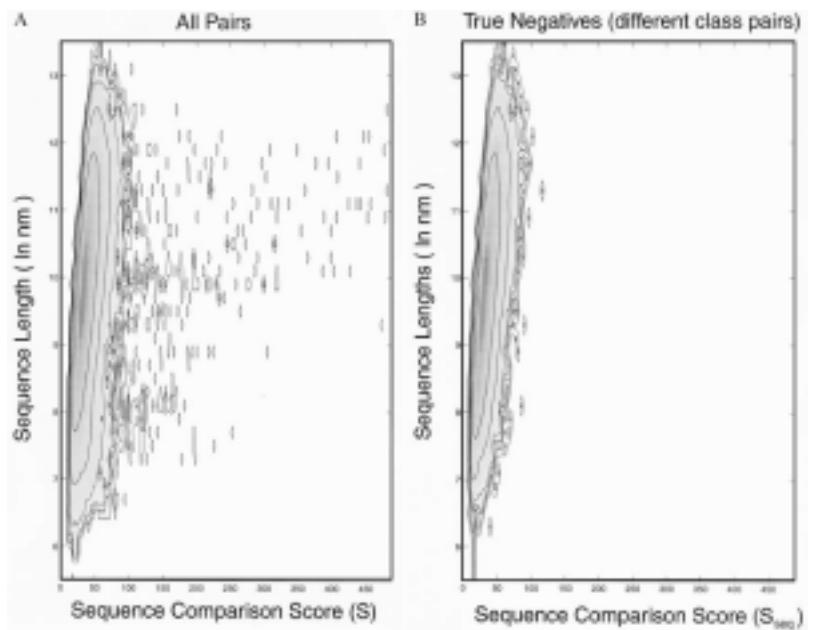
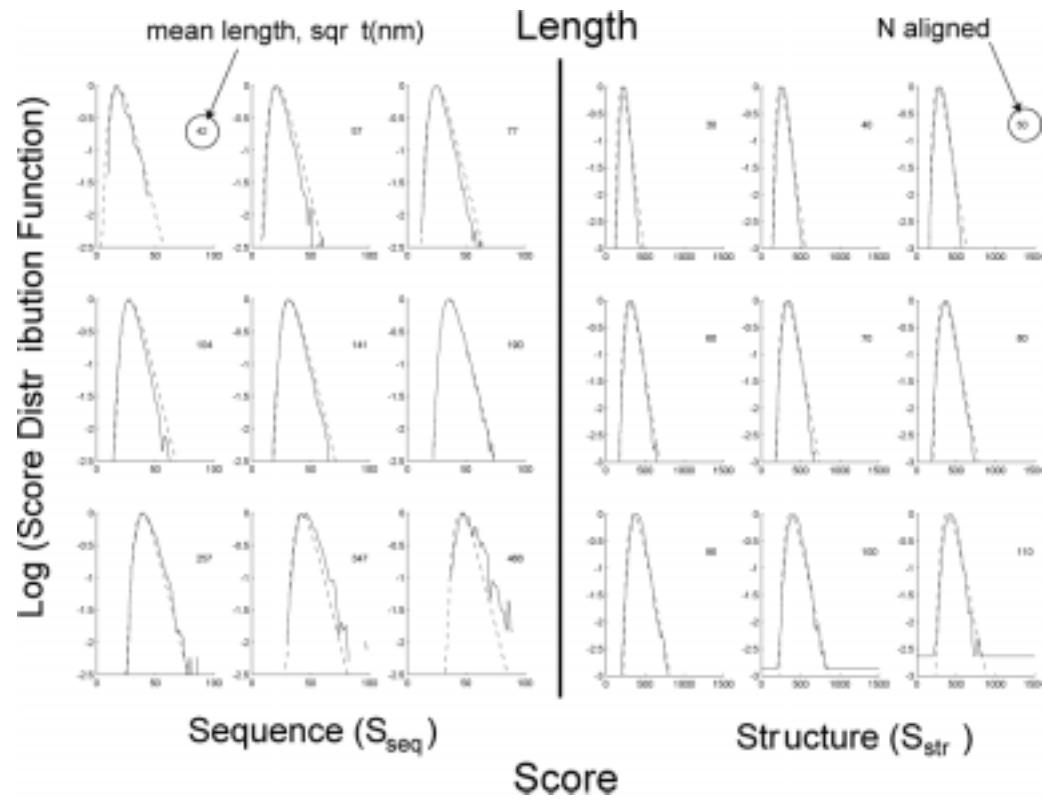
[ e.g.  $P(\text{score } s>392) = 1\% \text{ chance}$ ]

- Fit to Observed Distribution
  - ◊ All-vs-All comparison
  - ◊ Graph Distribution of Scores in 2D (N dependence)
  - ◊  $1K \times 1K$  families  $\rightarrow \sim 1M$  scores;  $\sim 2K$  TP embedded in this
  - ◊ Fit a function  $p(S)$  to distribution of true negatives (as determined by scop)

- A P-value for Significance

- ◊ Integration of  $p$  (ie the CDF) gives  $P(s>S)$ , the chance of getting a score better than threshold  $S$  randomly
- ◊ Extrapolated Percentile Rank: How does a Score Rank Relative to all Other Scores?
- ◊ For sequences, originally used in Blast (Karlin-Altschul). Then in FASTA, &c.

# Observed Distributions

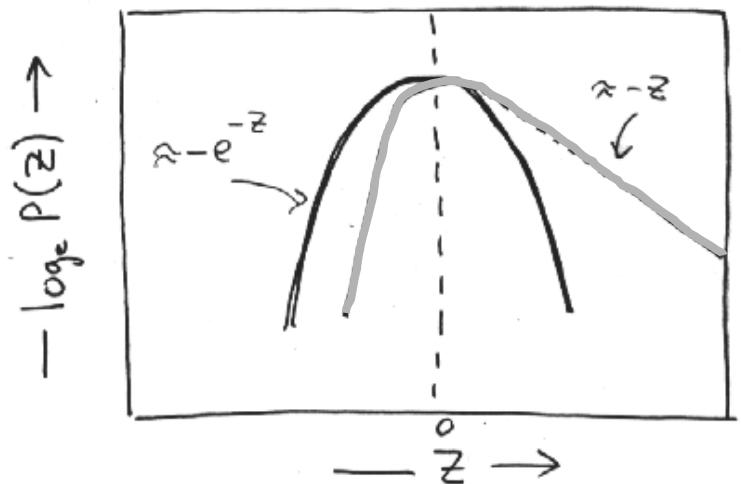
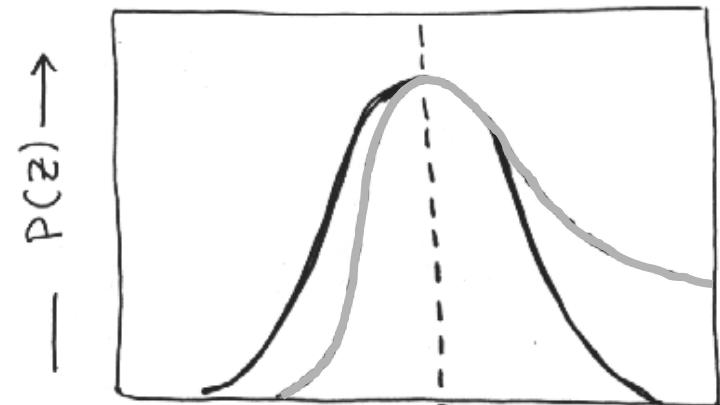
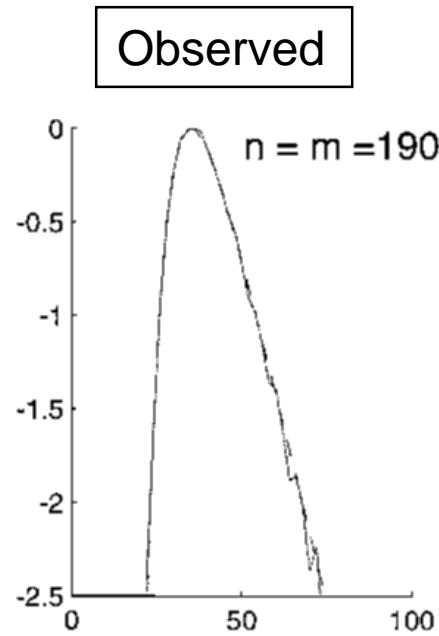


# EVD Fits

$$\rho(z) = \exp(-z - e^{-z})$$

$$(\ln \rho(z) = -z - e^{-z})$$

- Reasonable as Dyn. Prog. maximizes over pseudo-random variables
- EVD is **Max**(indep. random variables);
- Normal is **Sum**(indep. random variables)
  - ◊  $\rho(z) = \exp(-z^2)$ ,  $\ln \rho(z) = -z^2$



Extreme Value Distribution (EVD, long-tailed) fits the observed distributions best. The corresponding formula for the P-value:

$$P(z > Z) = \int \rho(z) dz = 1 - \exp(-e^{-Z}) \quad 12$$

# Same Results for Sequence & Structure

3 Free Parm. fit to EVD involving:  $a, b, \sigma$ .  
These are the only difference betw. sequence and structure.

$$Z = \frac{S - (a \ln N + b)}{\sigma}$$

$$S = \sum_{i,j} M(i, j) - G$$

$$\rho(z) = \exp(-z - e^{-z})$$

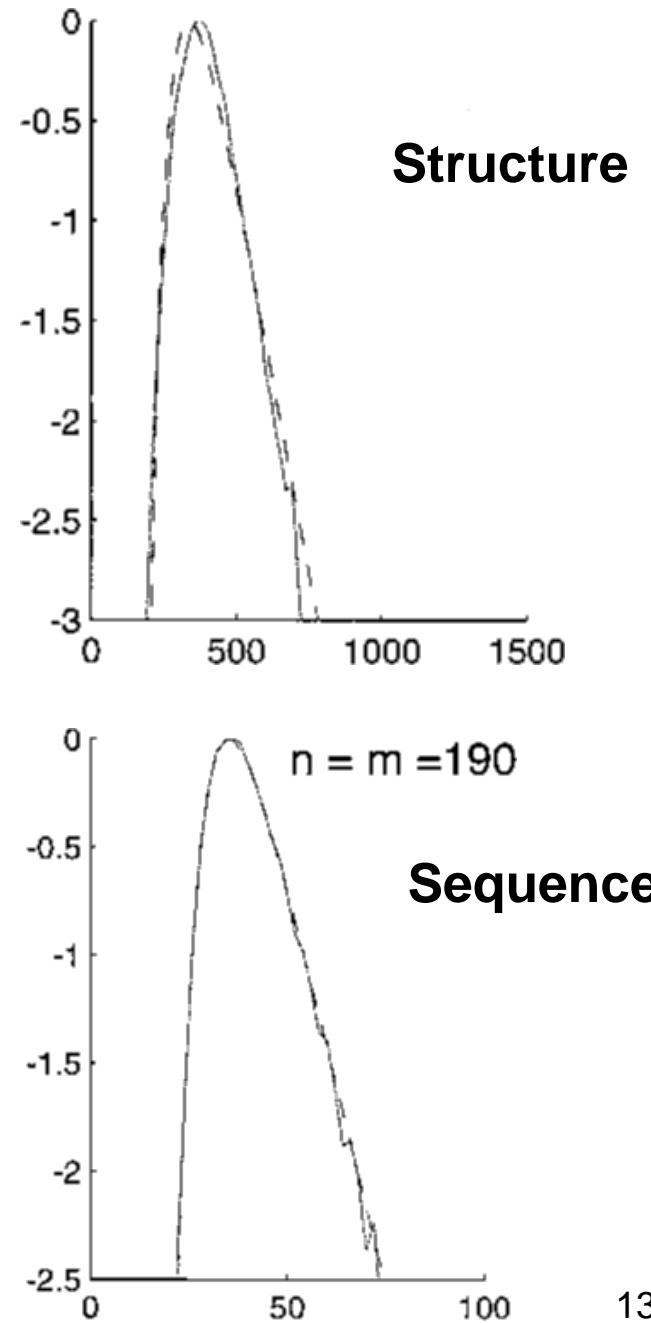
$N, G, M$  also defined differently for sequence and structure.

$N$  = number of residues matched.

$G$  = total gap penalty.

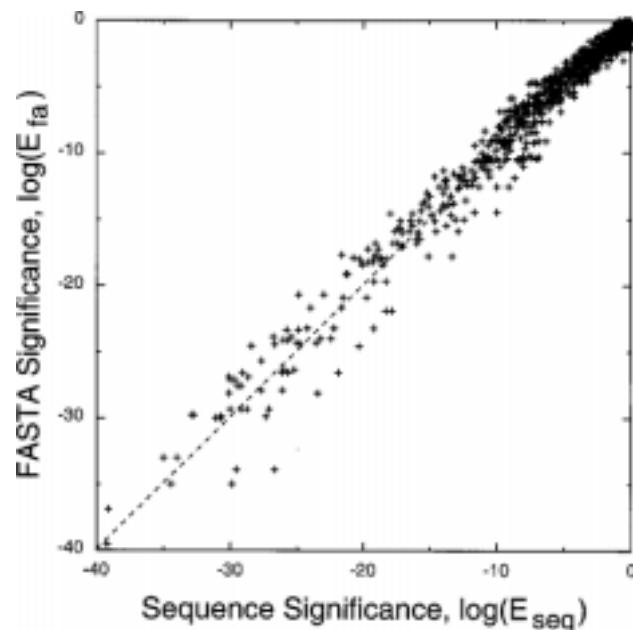
$M(i,j)$  = similarity matrix

(Blossum for seq. or  $M_{\text{str}}(i,j)$ , struc.)

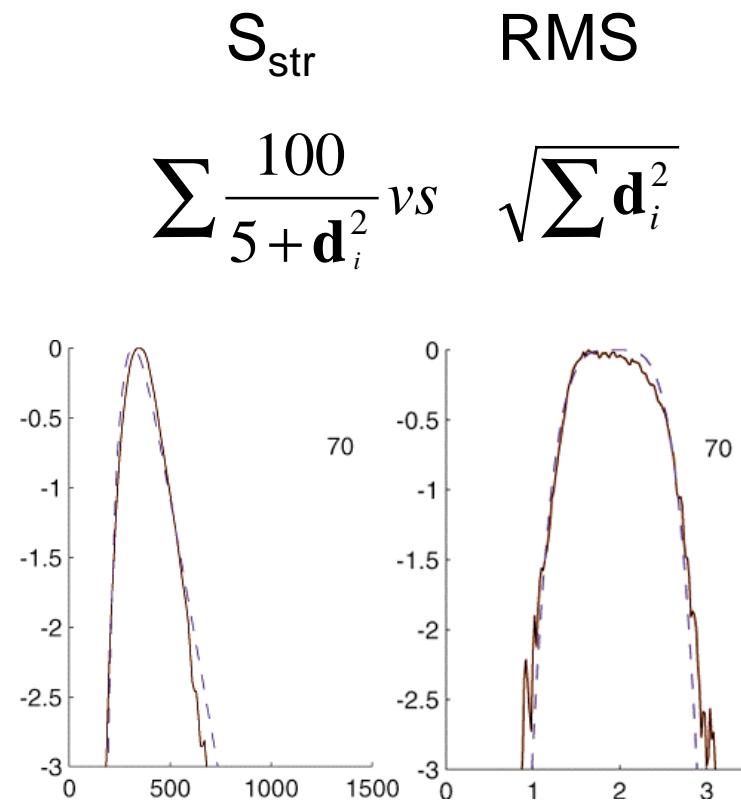


# Use Sequence Scores to Validate

- Sequence P-value perfectly tracks FASTA e-value
  - ◊ Validates approach
  - ◊ Added Benefit: allows computation of an e-value without doing a db run
- Significance computation can be applied to **any** existing sequence or structure alignment
- Also, RMS doesn't work instead of structural alignment (no EVD fit)
  - ◊ RMS penalizes worst fitting atoms, easily skewed



(c) M Gerstein (<http://bioinfo.mbb.yale.edu>)



# Comparing Genomes in terms of Protein Structure: Surveys of a Finite Parts List

## 1 Library of Known Folds

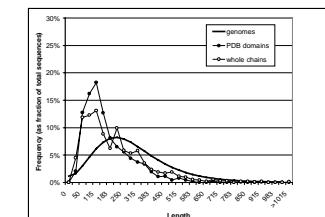
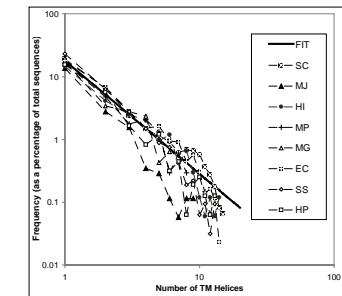
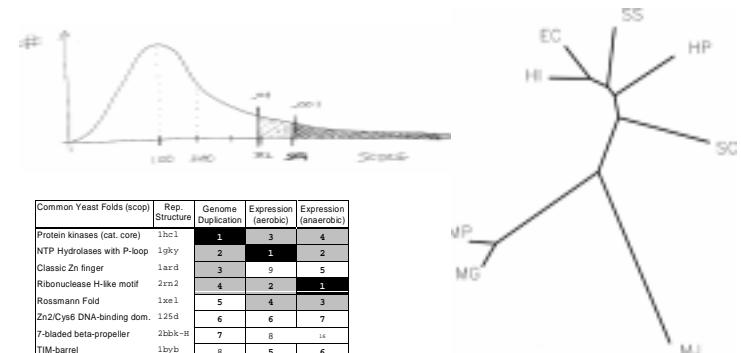
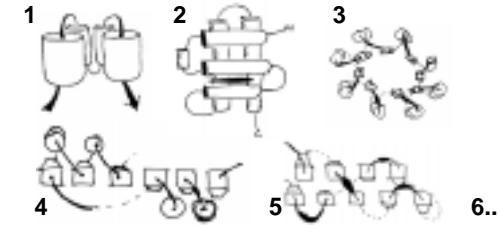
Importance of statistics. Scop auto-alignments. Significance follows EVD stats, same as sequences.

## 2 A Census of Known Folds

Which folds in which Organisms: Plants v. People, E coli v. yeast? Shared Fold Tree. Top-10 by duplication/expression, repeated  $\beta\alpha\beta$  supersecondary struc.

## 3 Prediction of Unknown Folds

Zipf law for TM's but no 7-TM's. Same 2<sup>o</sup> comp. but different a.a. comp. Can extrapolate from known structures to genomes?

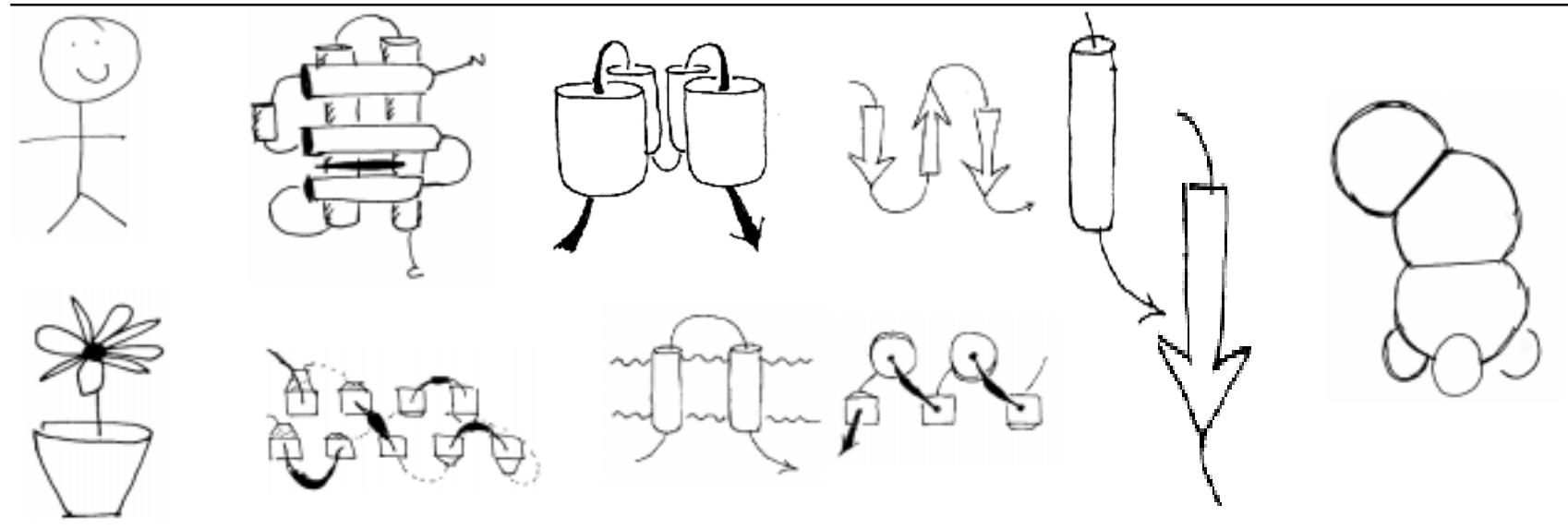


<http://bioinfo.mbb.yale.edu/genome>

**Acknowledgements:** M Levitt, scop  
(Murzin, Brenner, Ailey, Hubbard, Chothia)

# At What Structural Resolution Are Organisms Different?

person	protein	super-secondary	helix	individual
plant	fold (Ig)	structure ( $\beta\beta$ ,TM– TM, $\alpha\beta\alpha\beta,\alpha\alpha\alpha$ )	strand	atom (C,H,O...)



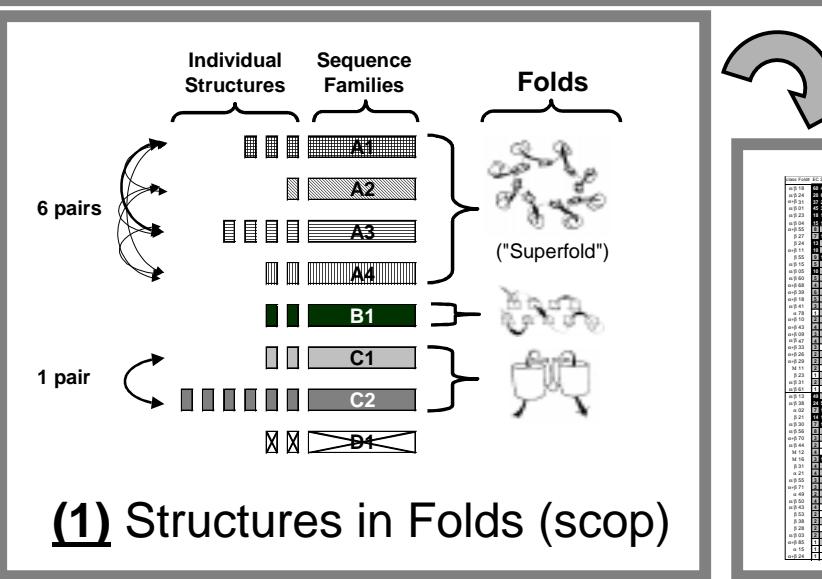
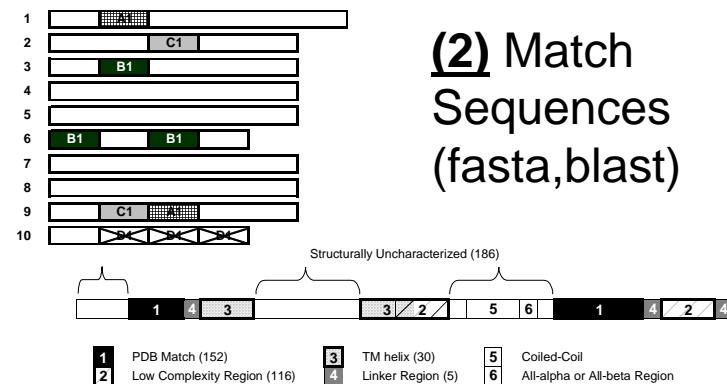
1m

100Å

10Å

1Å

# Cross-Reference: Folds→Sequences → Organisms



**(3) Organize Sequences by Genome or Taxon**

3 + 5

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

**(4) Results in “Fold Table”**

class	Fold#	EC	SC	HI	SS	HP	MJ	MP	MG	total	Fam.	PDB	Rep.	Struc.	Name
$\alpha/\beta$	18	<b>60</b>	<b>46</b>	<b>23</b>	<b>40</b>	<b>19</b>	<b>7</b>	<b>4</b>	<b>3</b>	202	<b>16</b>	183	1xel	-	NAD(P) bindii
$\alpha/\beta$	24	<b>20</b>	<b>69</b>	<b>17</b>	<b>19</b>	<b>17</b>	<b>16</b>	<b>10</b>	<b>11</b>	179	<b>13</b>	132	1gky	-	P-loop Contai
$\alpha+\beta$	31	<b>37</b>	<b>28</b>	<b>18</b>	<b>16</b>	<b>12</b>	<b>40</b>	<b>3</b>	<b>3</b>	157	<b>23</b>	160	1fxd	-	like Fe/Fodoxi
$\alpha/\beta$	01	<b>45</b>	<b>36</b>	<b>13</b>	<b>22</b>	<b>11</b>	<b>10</b>	<b>5</b>	<b>4</b>	146	<b>37</b>	399	1byb	-	TIM-barrel
$\alpha/\beta$	23	<b>18</b>	<b>17</b>	<b>7</b>	<b>9</b>	<b>4</b>	<b>8</b>	<b>2</b>	<b>2</b>	67	<b>5</b>	36	1pyd	a:2-181	Thiamin bindi
$\alpha/\beta$	04	<b>15</b>	<b>11</b>	<b>7</b>	<b>10</b>	<b>1</b>	<b>9</b>	<b>5</b>	<b>5</b>	63	<b>13</b>	132	2tmd	a:490-645	FAD/NAD(P)-
$\alpha+\beta$	55	<b>8</b>	<b>9</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>3</b>	<b>6</b>	<b>6</b>	56	<b>4</b>	23	1sry	a:11-421	Class-I aaRS
$\beta$	27	<b>7</b>	<b>10</b>	<b>8</b>	<b>8</b>	<b>4</b>	<b>4</b>	<b>3</b>	<b>3</b>	47	<b>5</b>	19	1fnb	19-154	Reductase/EI
$\beta$	24	<b>13</b>	<b>7</b>	<b>4</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	<b>3</b>	39	<b>18</b>	177	1snc	-	OB-fold
$\alpha+\beta$	11	<b>10</b>	<b>8</b>	<b>4</b>	<b>8</b>	<b>2</b>	<b>2</b>	<b>2</b>	<b>1</b>	37	<b>11</b>	48	1lgd	-	beta-Grasp

# Shared Folds in OWL and Initial Genomes

Venn Diagrams

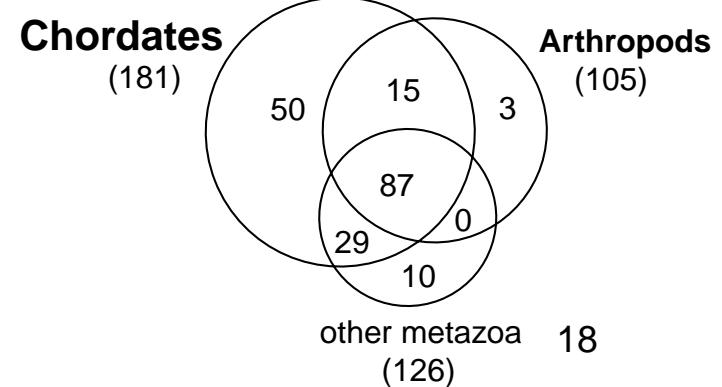
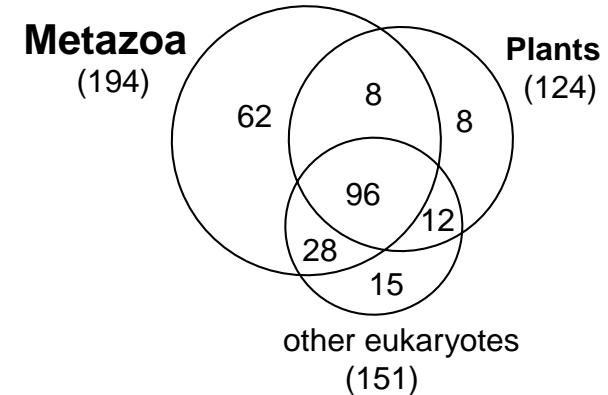
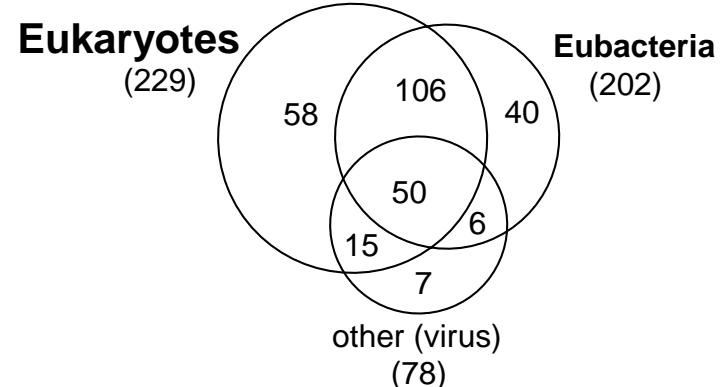
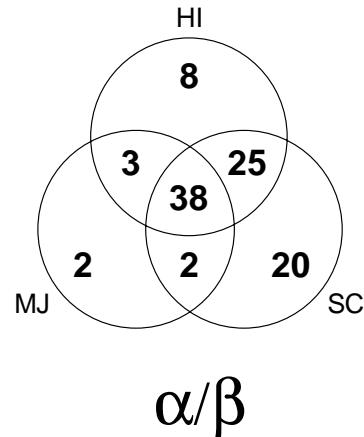
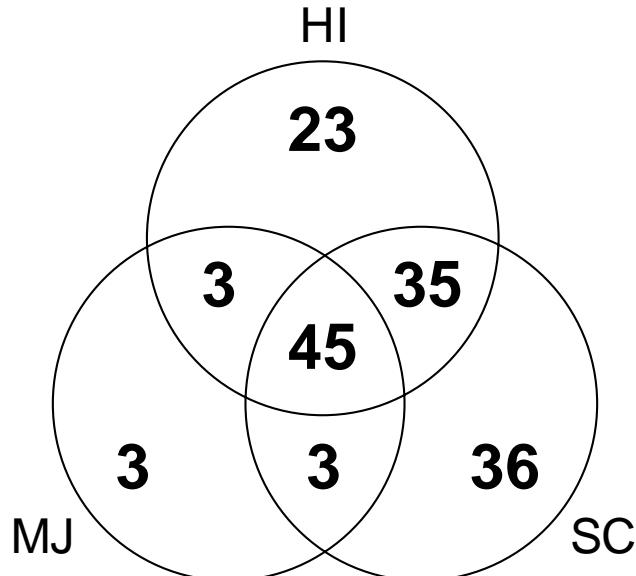
~300 folds (282 folds in scop 1.32 ['96])

~120K sequences in OWL 27.1

7 groups of organisms

3 of the first genomes

HI (bacteria), MJ (archeon), SC (eukaryote)



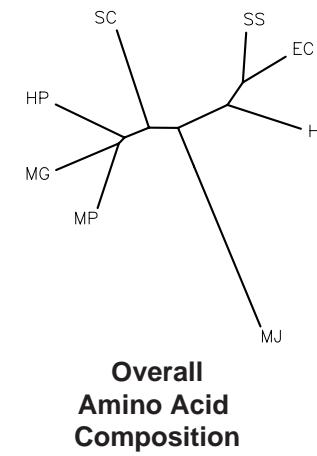
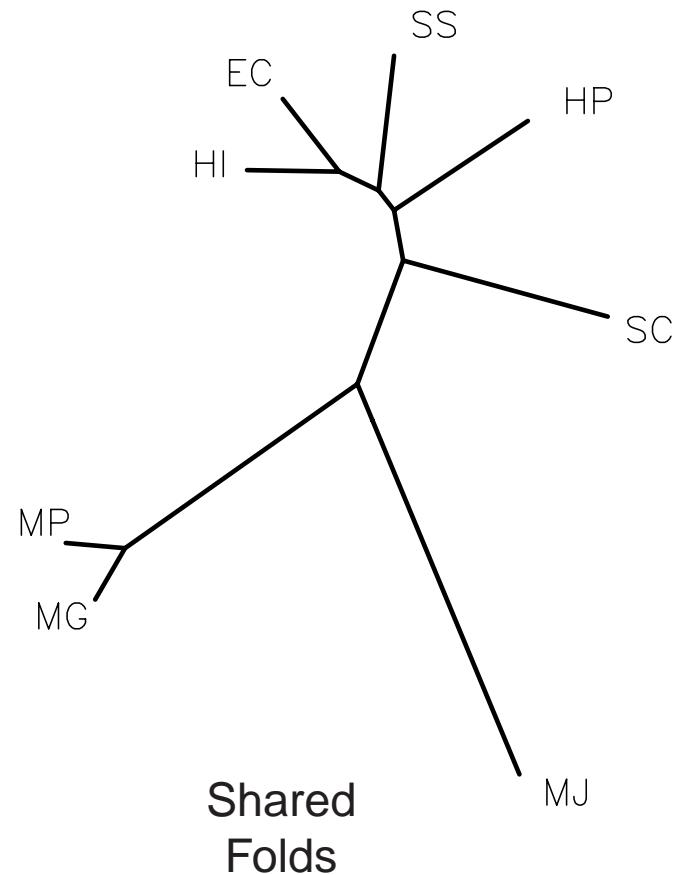
# Cluster Trees Grouping 8 Initial Genomes on Basis of Shared Folds

**D=S/T**

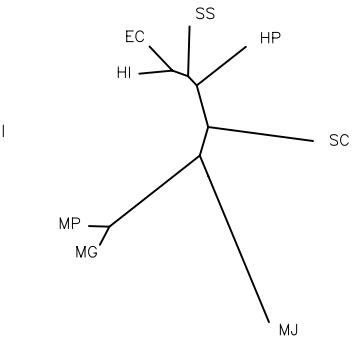
**D** = shared fold dist. betw.  
2 genomes

**S** = # shared folds

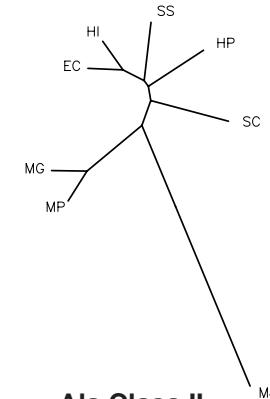
**T**= total # folds in both



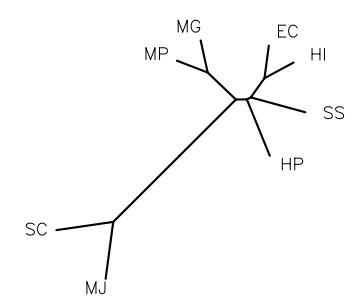
Overall  
Amino Acid  
Composition



Shared Sequence  
Families of  
Known Structure



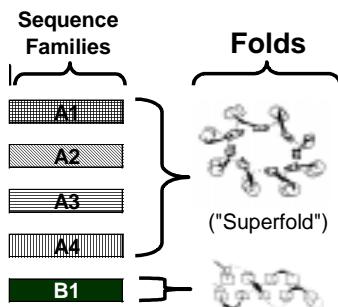
Ala Class II  
synthetase  
similarity



s17 ribosomal  
protein similarity  
(OB Fold)

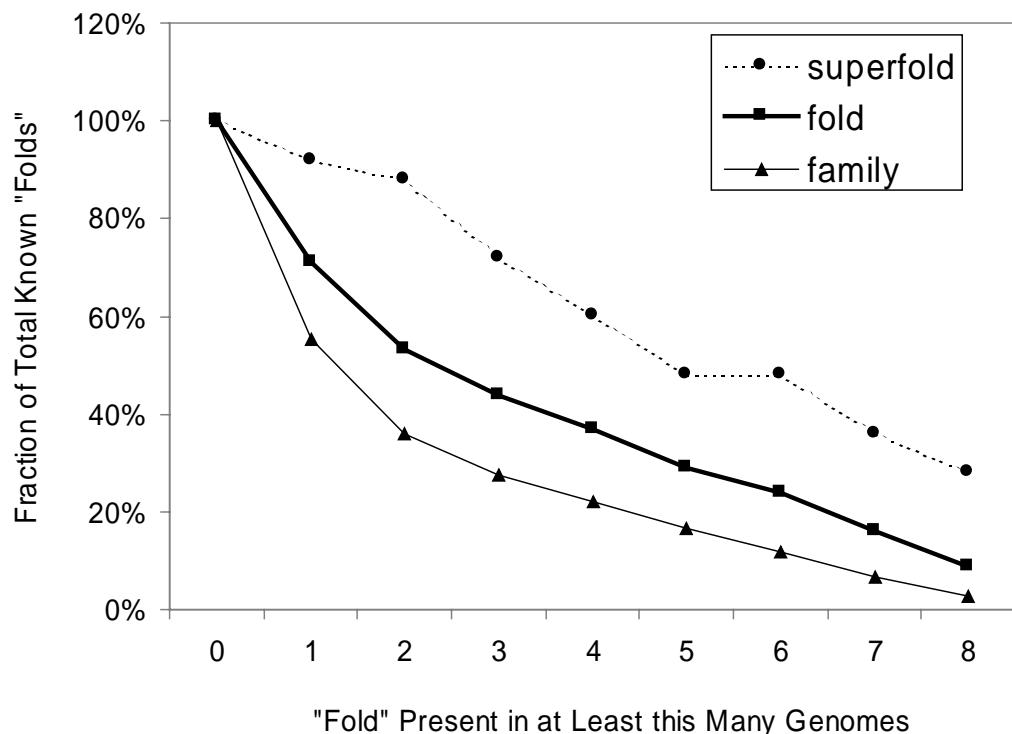
# Patterns of Folds Usage in 8 Genomes

	fold	fam.	super fold
total in PDB	338	990	25
in at least one of 8 genomes	240	547	23
present in this many genomes			
1	60	192	1
2	32	82	4
3	23	54	3
4	27	53	3
5	17	50	0
6	27	49	3
7	24	41	2
8	30	26	7



Superfold = fold that allows many non-homologous seq.  
(Thornton/Orengo)

| <b>ESHSHMM</b> (##) |
|---------------------|---------------------|---------------------|---------------------|---------------------|
| <b>CCISPJP</b> G    |
11111111 (30)	.1..... (23)	1..... (19)	11111..11 (16)	111111.. (16)
1111.... (09)	11111... (08)	1.1..... (08)	1.111..11 (06)	11..... (06)
...1.... (06)	1.11.... (05)	.1.1.... (05)	1.111... (04)	11.1.... (04)
.1....1.. (04)	..1.... (04)	111111.1 (03)	1111111. (03)	1111..11 (03)
1111.1.. (03)	.....1.. (03)	1111.111 (02)	111...11 (02)	111..11.. (02)
1..11.1.. (02)	..111... (02)	.1.11... (02)	1..1.1.. (02)	1..1..1.. (02)
111..... (02)	.11..... (02)	.....1. (02)	....1... (02)	111..111 (01)
111..111 (01)	1..111..1 (01)	1..1111.. (01)	.1.1..11 (01)	.1..11..1 (01)
.11..1..1 (01)	1.....111 (01)	1..111.. (01)	1..1...11 (01)	1..1..11.. (01)
11.....11 (01)	11..1..1.. (01)	11..11... (01)	111..1.. (01)	111..11.. (01)
.11....1.. (01)	1.....111 (01)	1...11.. (01)	1..1..1... (01)	....111 (01)
....1..1.. (01)	....1..1.. (01)	....11.. (01)	..1..1... (01)	.1....1.. (01)
1.....1.. (01)	.....1..1 (01)	.....1.. (01)	.....1.. (01)	.....1.. (01)



# Top-10 Folds in 8 First Genomes

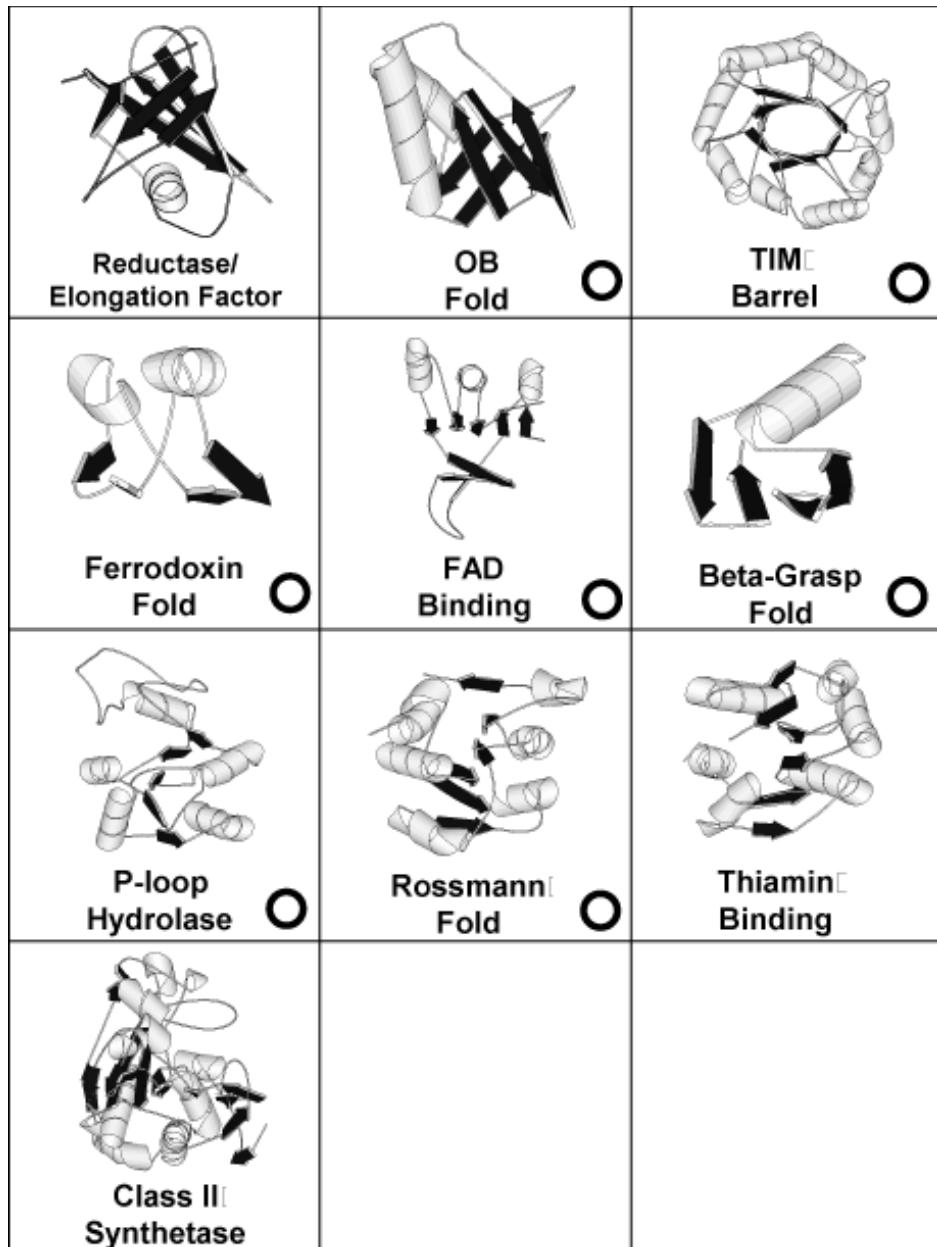
Num in genome	Class	Fold Name	Representative Structure (PDB selection)
<b>Top-10 in a eukaryotic genome (SC)</b>			
84	$\alpha\beta$	Protein kinases (catalytic core)	1irk
49	$\alpha\beta$	P-loop containing NTP hydrolases	1gky
35	$\alpha\beta$	Rossmann Fold	2ohx A:175-324
31	$\alpha\beta$	TIM barrel	1tim A:
25	$\alpha\beta$	Ribonuclease H-like	2rn2
18	S	Classic zinc finger	1zaa C:
14	$\alpha\beta$	Ubiquitin conjugating enzyme	1aaK
12	$\beta$	GroES-like	1acy L:109-211
10	$\alpha\beta$	Thioredoxin-like	1txx
9	$\alpha\beta$	Thiamin-binding Fold	1pvd A:2-181
5x8	...	...	...
7	$\alpha\beta$	Flavodoxin-like	3chy
<b>Top-11 in a eubacterial genome (H)</b>			
18	$\alpha\beta$	Rossmann Fold	2ohx A:175-324
13	$\alpha\beta$	P-loop containing NTP hydrolases	1gky
12	$\alpha\beta$	Flavodoxin-like	3chy
10	$\alpha\beta$	TIM barrel	1tim A:
10	$\alpha\beta$	Ferrodoxin-like	1fxd
10	$\alpha\beta$	Ribonuclease H-like	2rn2
6	$\alpha\beta$	Periplasmic binding protein-like II	1sbp
5	$\alpha\beta$	Periplasmic binding protein-like I	2drj
5	$\alpha\beta$	Like Class II arRS synthetases	1sry A:111-421
4	$\beta$	OB-fold	1pyp
4	$\alpha\beta$	Thiamin-binding Fold	1pvd A:2-181
<b>Top-11 in an archaeal genome (M)</b>			
19	$\alpha\beta$	Ferrodoxin-like	1fxd
10	$\alpha\beta$	P-loop containing NTP hydrolases	1gky
7	$\alpha\beta$	TIM barrel	1tim A:
6	$\alpha\beta$	Rossmann Fold	2ohx A:175-324
5	$\alpha$	Histone-fold	1ntx
4	$\alpha\beta$	Thiamin-binding Fold	1pvd A:2-181
4	$\alpha\beta$	Flavodoxin-like	3chy
4	$\beta$	Reductase/elongation factor common	1efg A:283-403
3	$\alpha\beta$	ATP-grasp	1bnc A:115-330
3	$\alpha\beta$	PLP-dependent transferases	1dkx
3	$\alpha\beta$	ATP pyrophosphatases	1gpm A:208-404

What are  
the most  
common  
folds?

How many  
shared?

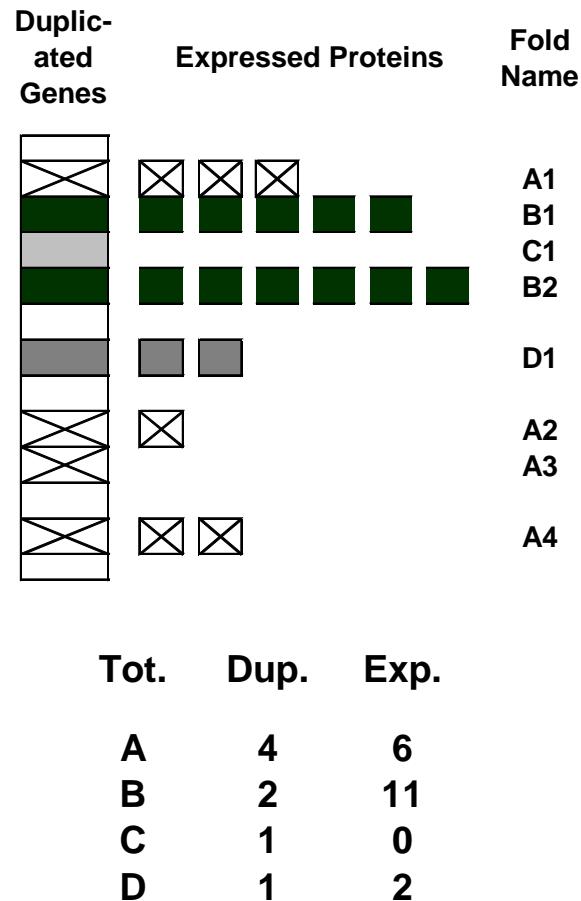
Many  
Super-  
folds

Mostly  
 $\alpha/\beta$





# Top-10 Folds according to Expression

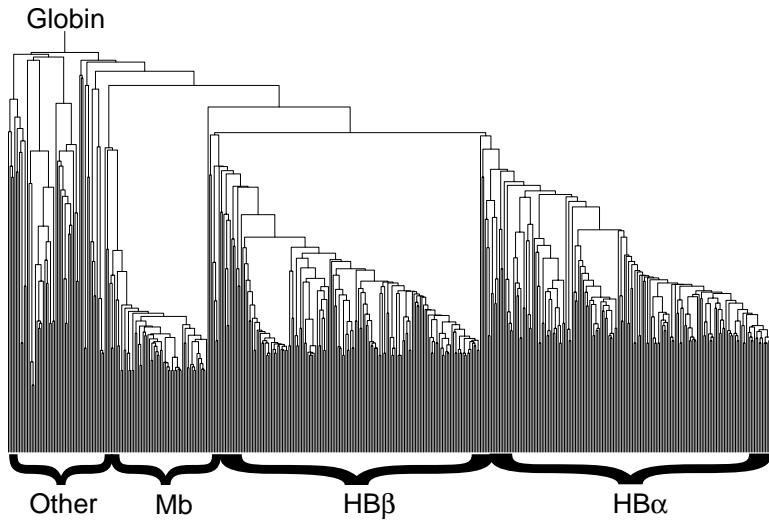


- Previous top-10 measures duplication
- Now weight by expression using data from Brown et al.

Common Yeast Folds (scop)	Rep. Structure	Genome Duplication	Expression (aerobic)	Expression (anaerobic)
Protein kinases (cat. core)	1hcl	1	3	4
NTP Hydrolases with P-loop	1gky	2	1	2
Classic Zn finger	1ard	3	9	5
Ribonuclease H-like motif	2rn2	4	2	1
Rossmann Fold	1xel	5	4	3
Zn2/Cys6 DNA-binding dom.	125d	6	6	7
7-bladed beta-propeller	2bbk-H	7	8	16
TIM-barrel	1byb	8	5	6
like Ferrodoxin	1fxd	9	7	10
DNA-binding 3-helix bundle	1enh	10	30	36
...	...	...	...	...
GroES-like	1lep-A	17	10	9
...	...	...	...	...
like HSP70, Ct-dom.	1dkz-A	22	11	8



# An Issue with Fold Counting: Biases in the Databanks

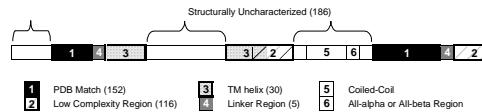


Example Structure (PDB)	Fold Name	Percentage of known folds in genome	Rank in eubacterial Top-10
<b>Top-10 in a bacterial genome (<i>H. influenzae</i>)</b>			
2HSD-A	Rossmann Fold (NAD binding)	9.6	1
1AKE-A	NTP Hydrolases containing P-loop	5.7	3
1RCF	Flavodoxin-like	5.1	4
6TIM-B	TIM-barrel	4.5	2
1FXD	Ferredoxin-like	4.2	5
2RN2	like Ribonuclease H	3.0	16
1SBP	like Periplasmic binding protein (class II)	3.0	11
2DRI	like Periplasmic binding protein (class I)	3.0	19
1SRY-*	Class II aaRS and biotin synthetases	2.7	50
1PYP	OB-fold	2.7	9

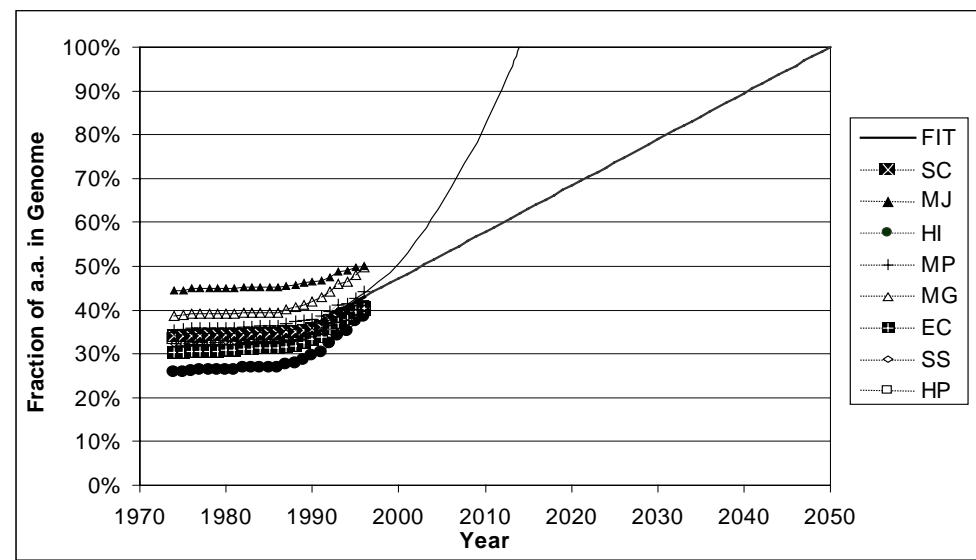
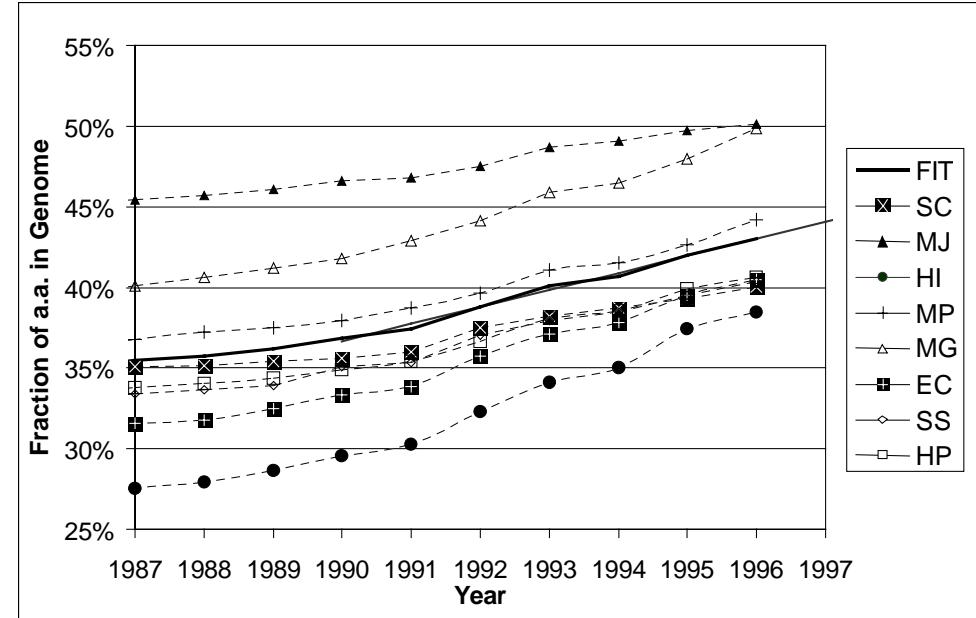
- Over-representation of certain species and functions in the databanks (e.g. human v. plant globins, Ig's)
  - Nevertheless HI top-10 like eubacterial top-10
- PDB is small, biased sample of genome (6-12%)
- Selection of structures in PDB is biased, anecdotal
- Different numbers with different comparison sensitivity: FASTA, HMM, &c
- Some Correction with Seq. Weighting, Diff. Sampling

# The Problem: All Folds in Genome Not Known until 20?? → Prediction

- Separate TM, LC, linkers
- How many residues in genome matched by known folds, in 1975, '76, '77... '00... '50



(c) M Gerstein (<http://bioinfo.mbb.yale.edu>)



# Comparing Genomes in terms of Protein Structure: Surveys of a Finite Parts List

## 1 Library of Known Folds

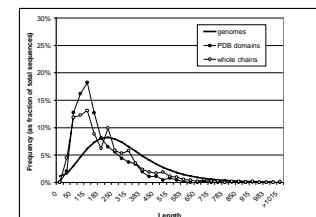
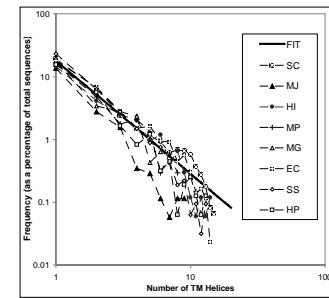
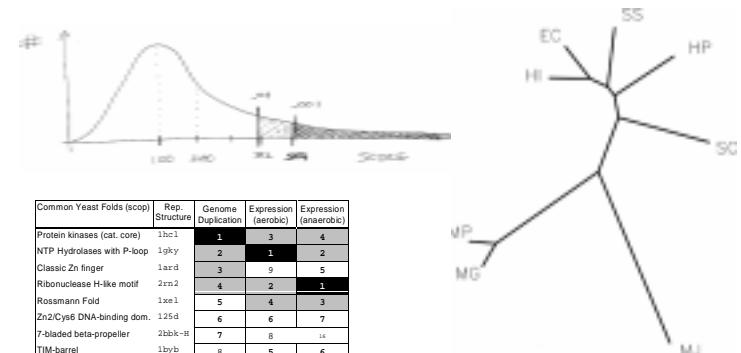
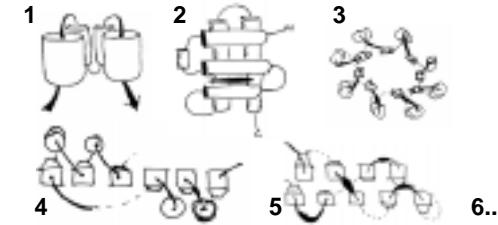
Importance of statistics. Scop auto-alignments. Significance follows EVD stats, same as sequences.

## 2 A Census of Known Folds

Which folds in which Organisms: Plants v. People, E coli v. yeast? Shared Fold Tree. Top-10 by duplication/expression, repeated  $\beta\alpha\beta$  supersecondary struc.

## 3 Prediction of Unknown Folds

Zipf law for TM's but no 7-TM's. Same 2<sup>o</sup> comp. but different a.a. comp. Can extrapolate from known structures to genomes?

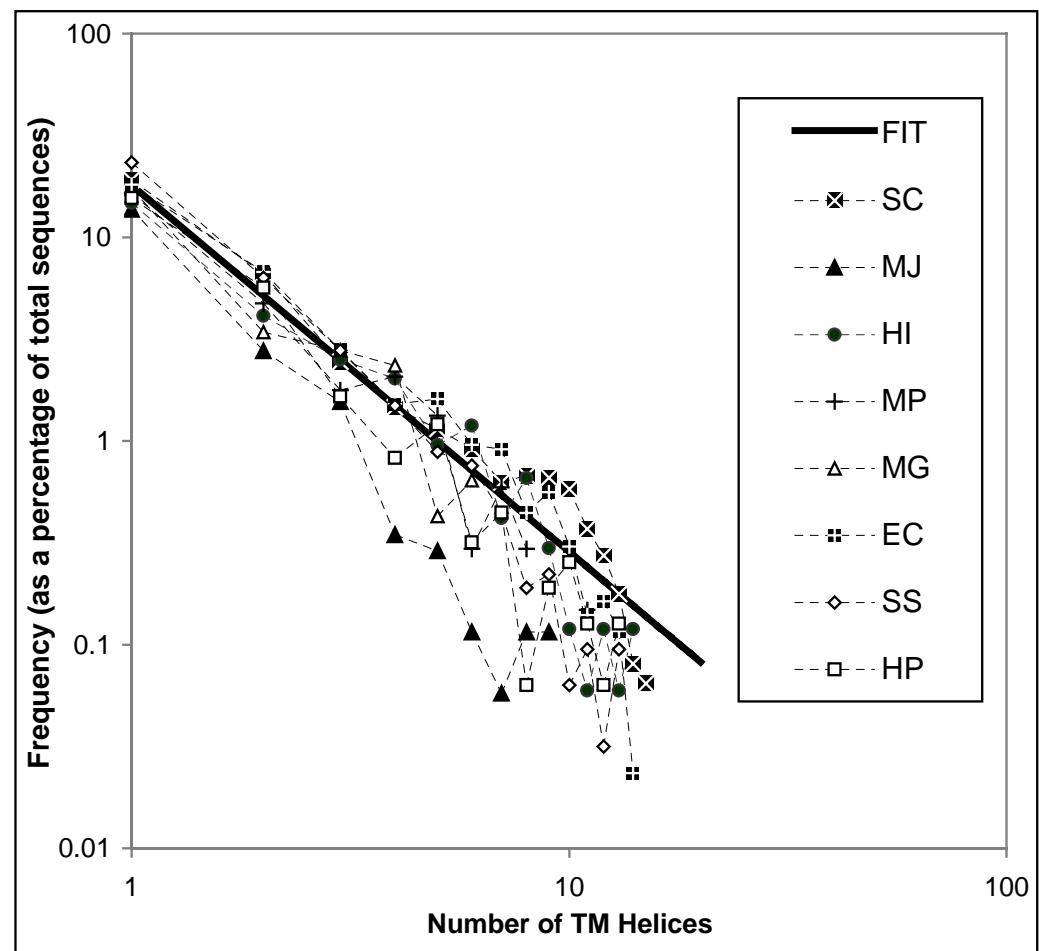
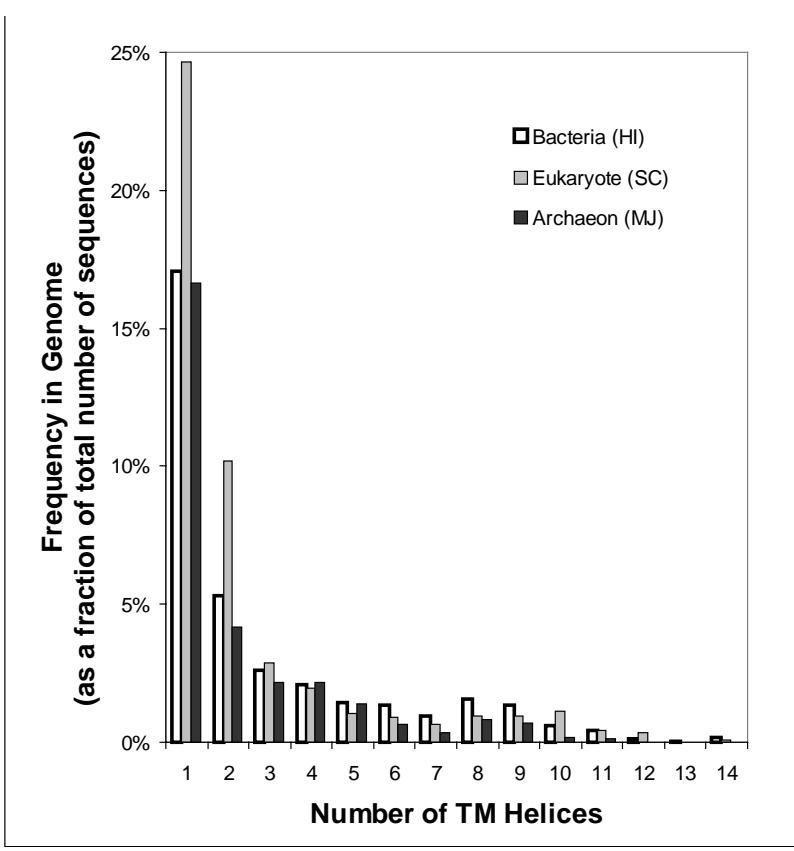


<http://bioinfo.mbb.yale.edu/genome>

**Acknowledgements:** M Levitt, scop  
(Murzin, Brenner, Ailey, Hubbard, Chothia)

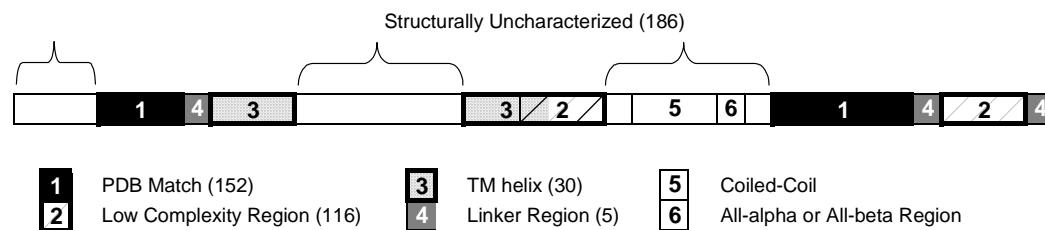
# TM-helix “prediction”

- TM prediction (KD, GES). Count number with 2 peaks, 3 peaks, &c.
- Yeast has more mem. prots., esp. 2-TMs
- Similar conclusions to others: von Heijne, Rost, Jones, &c.
- No preference for particular supersecondary structures: 7-TM's
- Freq. of Number of TM helices follows a Zipf-like law:  $F=1/[5n^2]$



# 2<sup>o</sup> Structure Prediction

- Bulk prediction of 2<sup>o</sup> struc. in genomes
- Same fraction of  $\alpha$  and  $\beta$  (by element, half each)
- Both overall and only for unknown soluble proteins.



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

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%

Not expected  
since.....





# Genomes Sequences are longer than those in Known Structures

Assess 2<sup>o</sup>,TM predictions

- (+) comprehensive, statistical
- (-) predictions inaccurate  
(~65%)

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

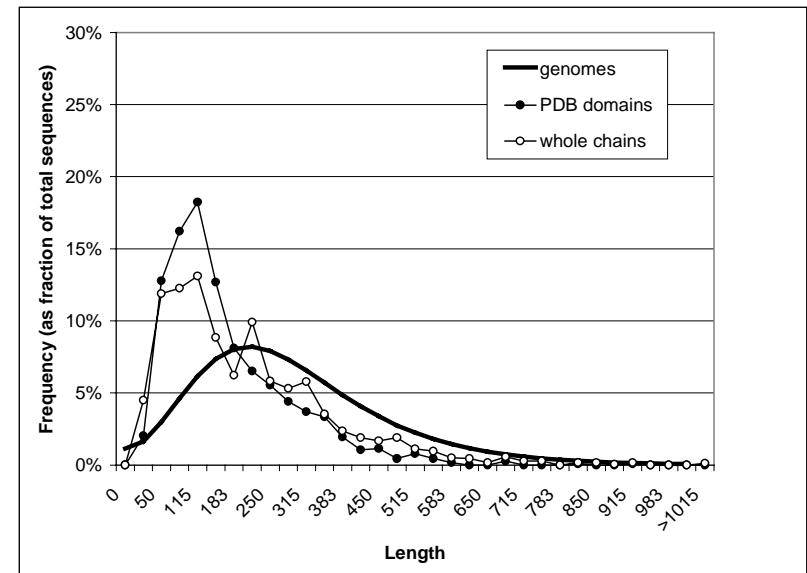
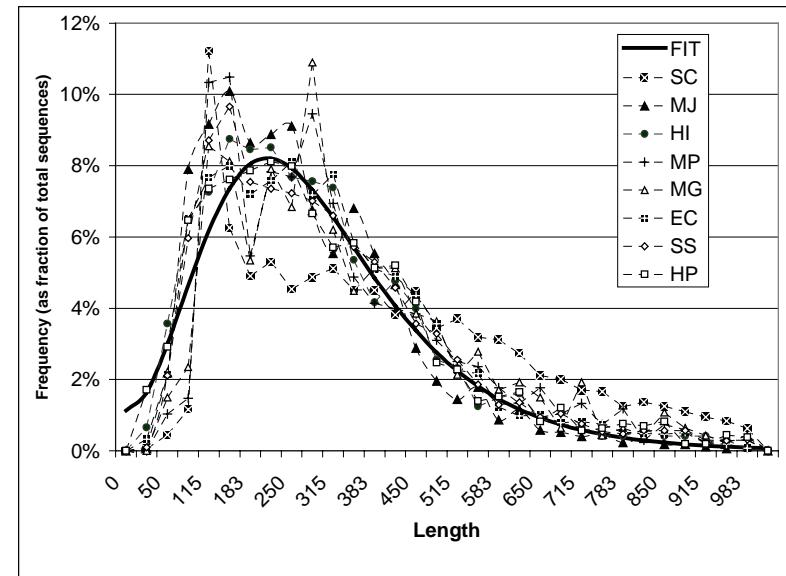
How Representative are the Known Structures of the Proteins in Complete Genome? Is prediction (extrapolation) based on known structures justified?

340 aa for avg. genome seq.

(470 aa for yeast)

205 aa for PDB chain

170 aa for PDB domain





# Comparing Genomes in terms of Protein Structure: Surveys of a Finite Parts List

## 1 Library of Known Folds

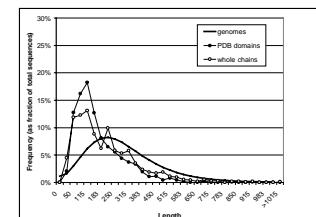
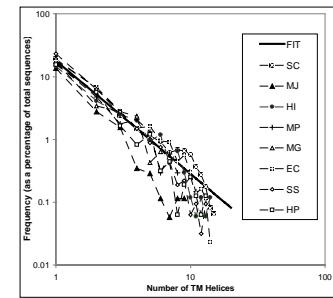
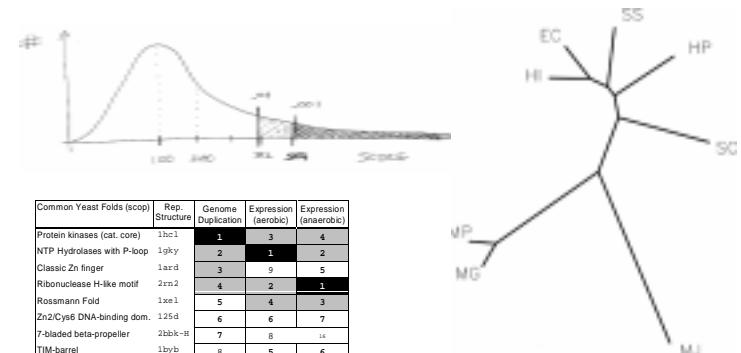
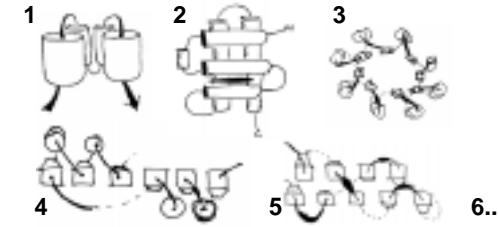
Importance of statistics. Scop auto-alignments. Significance follows EVD stats, same as sequences.

## 2 A Census of Known Folds

Which folds in which Organisms: Plants v. People, E coli v. yeast? Shared Fold Tree. Top-10 by duplication/expression, repeated  $\beta\alpha\beta$  supersecondary struc.

## 3 Prediction of Unknown Folds

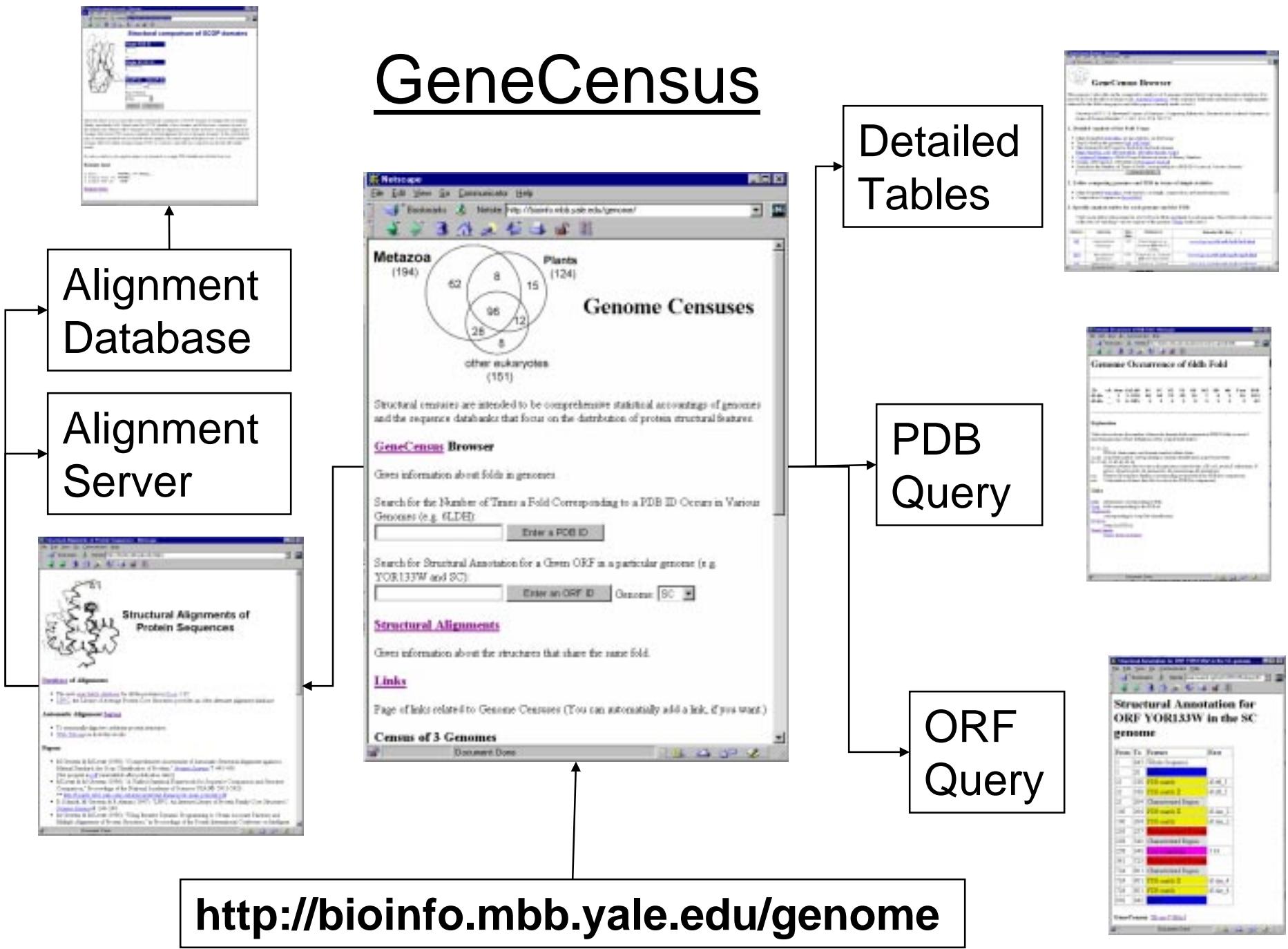
Zipf law for TM's but no 7-TM's. Same 2<sup>o</sup> comp. but different a.a. comp. Can extrapolate from known structures to genomes?



<http://bioinfo.mbb.yale.edu/genome>

**Acknowledgements:** M Levitt, scop  
(Murzin, Brenner, Ailey, Hubbard, Chothia)

# GeneCensus



# Comparing Genomes in terms of Protein Structure: Surveys of a Finite Parts List

## 1 Library of Known Folds

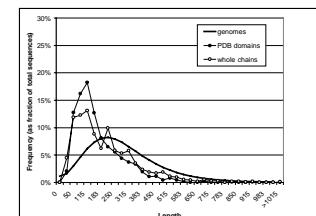
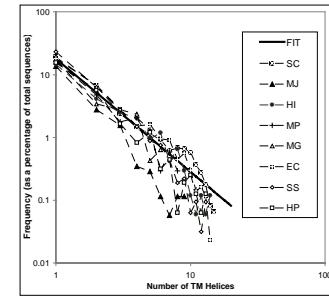
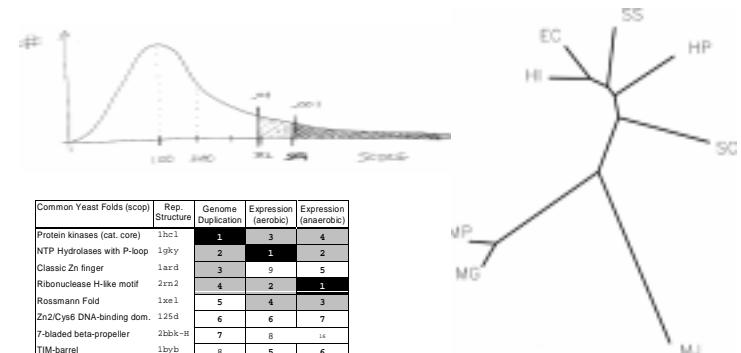
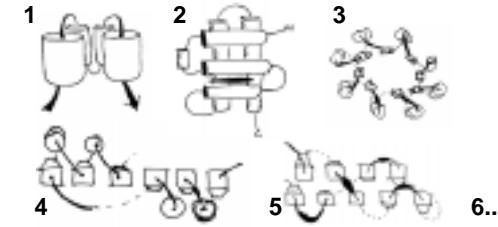
Importance of statistics. Scop auto-alignments. Significance follows EVD stats, same as sequences.

## 2 A Census of Known Folds

Which folds in which Organisms: Plants v. People, E coli v. yeast? Shared Fold Tree. Top-10 by duplication/expression, repeated  $\beta\alpha\beta$  supersecondary struc.

## 3 Prediction of Unknown Folds

Zipf law for TM's but no 7-TM's. Same 2<sup>o</sup> comp. but different a.a. comp. Can extrapolate from known structures to genomes?



<http://bioinfo.mbb.yale.edu/genome>

**Acknowledgements:** M Levitt, scop  
(Murzin, Brenner, Ailey, Hubbard, Chothia)

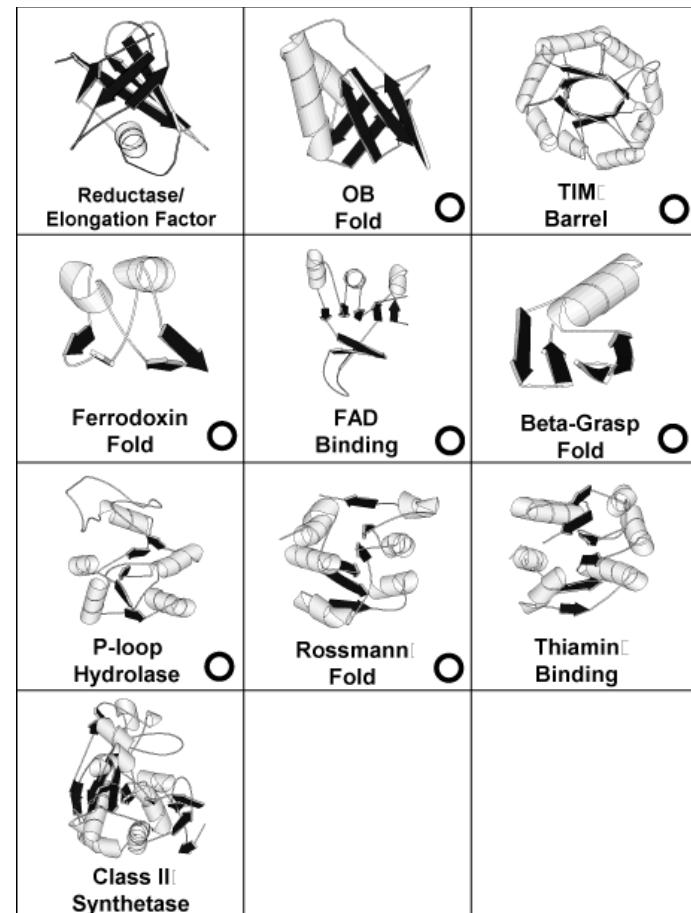
# Issues with Top-10

Depend on:  
Database (version) used, Biases  
Which Genomes  
Comparison Prog. (FASTA/BLAST)  
Threshold ( $1e-3$ ,  $1e-4$ )

**Want a Fair Census with  
Uniform & Consistent Sampling**

Common Yeast Folds	Rep. Struct.	GeneCensus		Sacch3D, SGD		diff.
		count	rank	count	rank	
Protein kinases (cat. core)	1hcl	110	1	109	1	
NTP Hydrolases with P-loop	1gky	69	2	52	2	
Classic Zn finger	1ard	55	3	34	7	.
Ribonuclease H-like motif	2rn2	54	4	30	8	.
Rossmann Fold	1xel	46	5	41	5	
Zn2/Cys6 DNA-binding dom.	125d	46	6	30	9	.
7-bladed beta-propeller	2bbk-H	46	7	0	-	<
TIM-barrel	1byb	36	8	39	6	.
Ferrodoxin-like	1fxd	28	9	43	4	.
DNA-binding 3-helix bundle	1enh	22	10	22	10	
Long Helix Oligomers (coils)	1zta	1	-	47	3	<

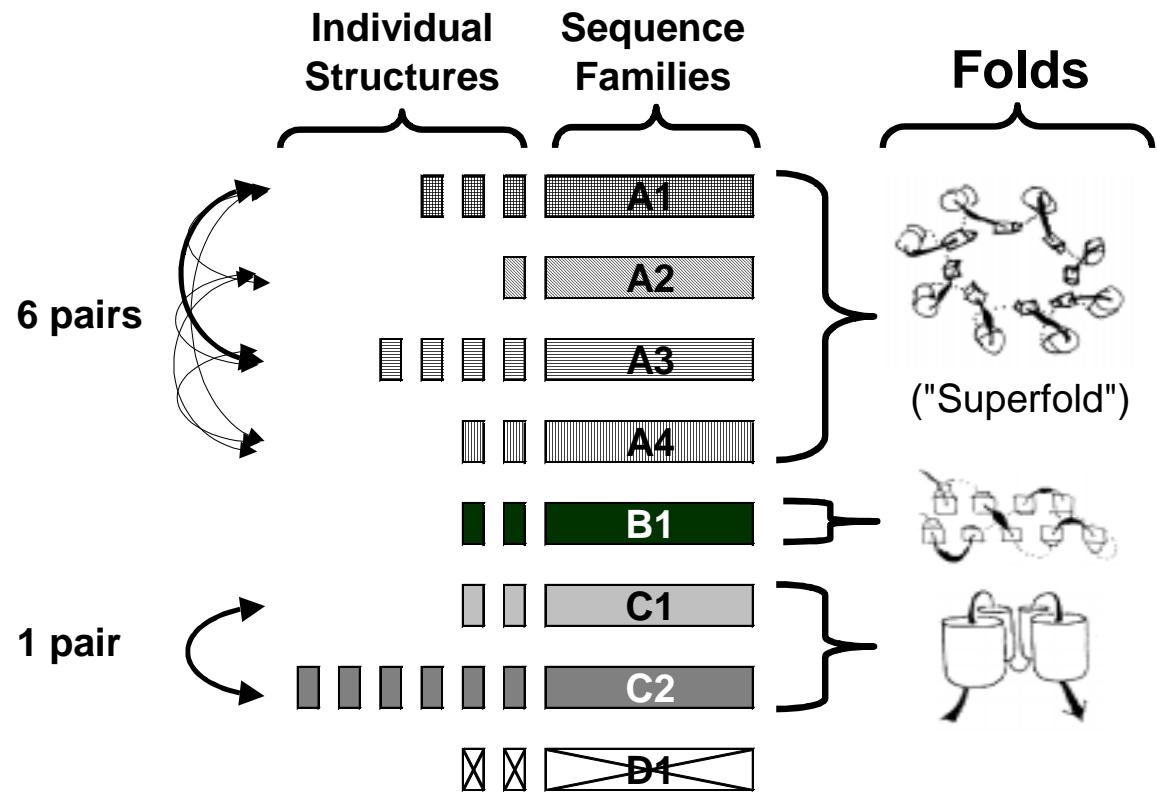
Top-10 in 8 Genomes



Sacc3D  
(Chervitz, Cherry)  
vs GeneCensus

# Clustering the PDB, assess comparison methods

- Structure of “scop”
  - ◊ ~10K structure **domains**
  - ◊ ~1K sequence **families**
  - ◊ ~350 **folds**
  - ◊ 25 **superfolds** (>10 fam)
  - ◊ ~2K + ~2K true positive pairs (1.32)
- Different Sequence Thresholds for Clustering
  - ◊ high threshold, struc. sim.
  - ◊ Low threshold, seq. sim.
  - ◊ .01 e-value for all comparisons



- Brenner et al. used scop pairs to calibrate the statistics of comparison methods
  - ◊ How many of the TP pairs can be found with FASTA?
  - ◊ e-value ~ epq
- Credits: scop, Murzin, Brenner, Ailey, Hubbard, Chothia; Thornton, superfolds; Pearson, fasta