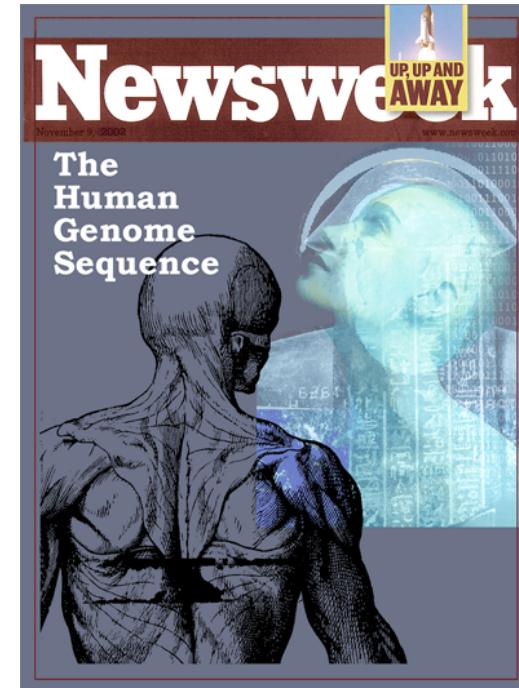
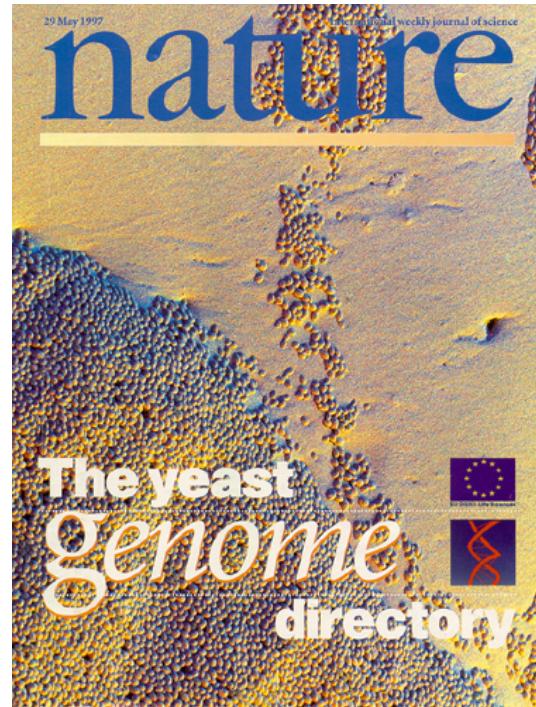
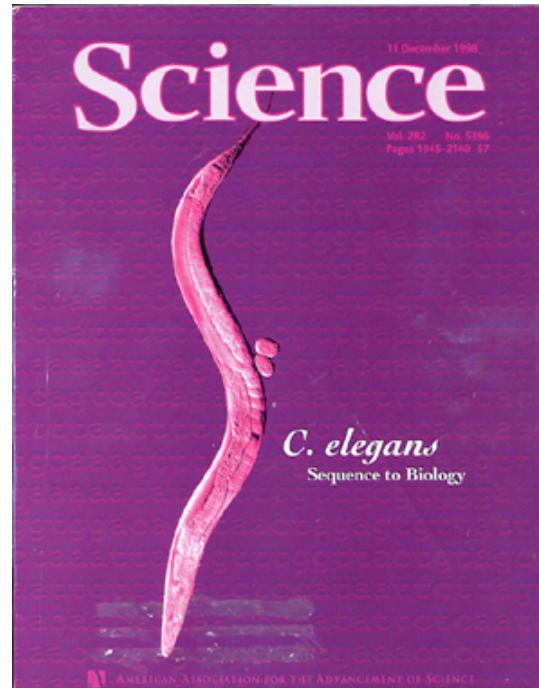


Comparative Genomics:

Surveys of a Finite Parts List

Mark Gerstein

Genomes highlight the Finiteness of the World of Sequences



1995

Bacteria, 1.6 Mb, ~1600 genes [Science 269: 496]

1997

Eukaryote, 13 Mb, ~6K genes [Nature 387: 1]

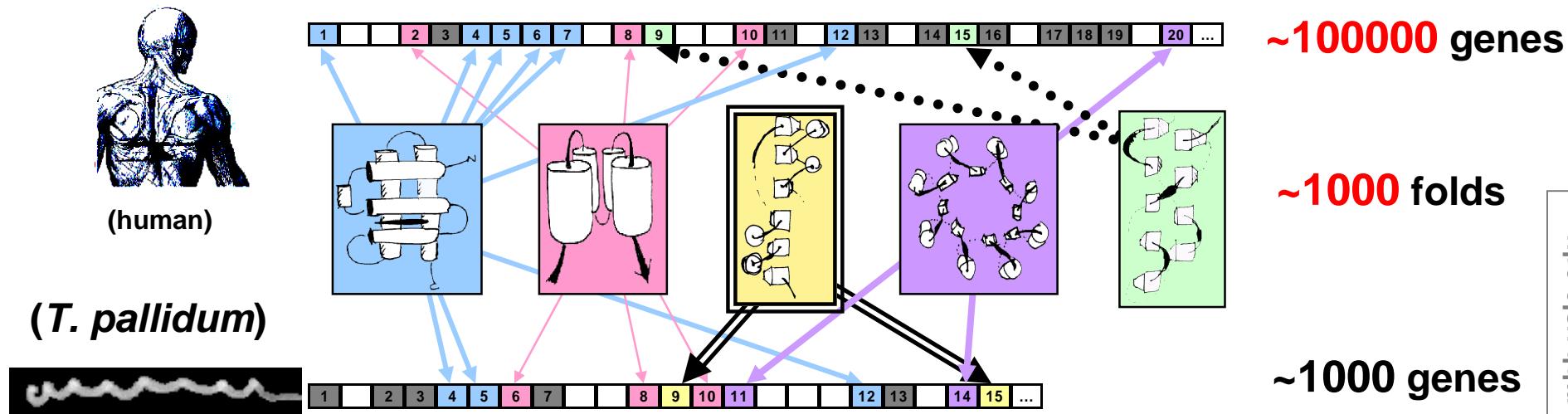
1998

Animal, ~100 Mb, ~20K genes [Science 282: 1945]

2000?

Human, ~3 Gb, ~100K genes [???

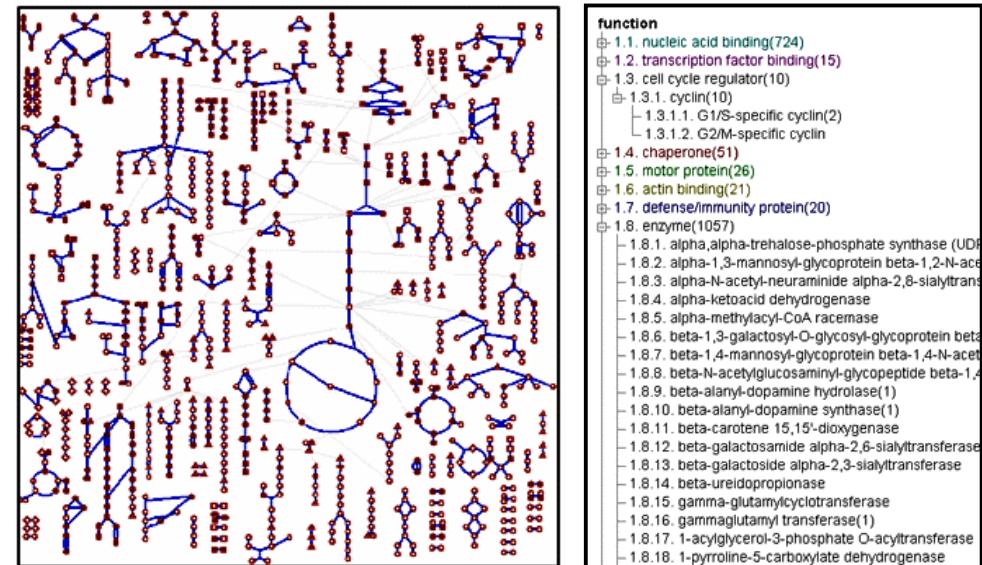
World of Structures is even more Finite, providing a valuable simplification



Likewise, for pathways,
functions, sequence
families, blocks, motifs....

Cross-referencing and interrelating among
a Different Types of Parts

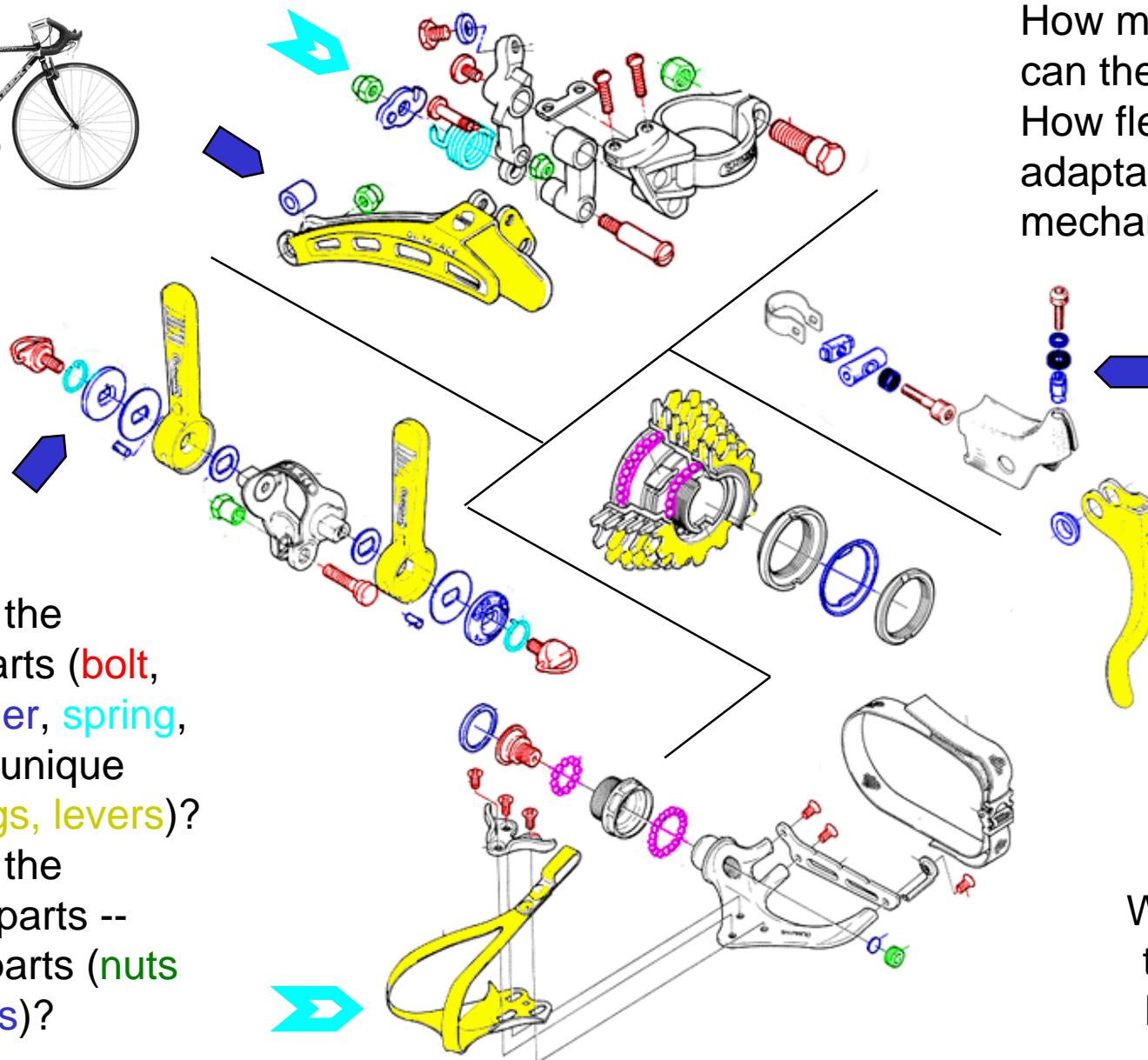
Functions picture from www.fruitfly.org/~suzi (Ashburner);
Pathways picture from ecocyc.pangeasystems.com/ecocyc
(Karp, Riley). Related resources: COGS, ProDom, Pfam, Blocks,
Domo, WIT, CATH, Scop....



A Parts List Approach to Bike Maintenance



A Parts List Approach to Bike Maintenance



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

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

Where are the parts located?

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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

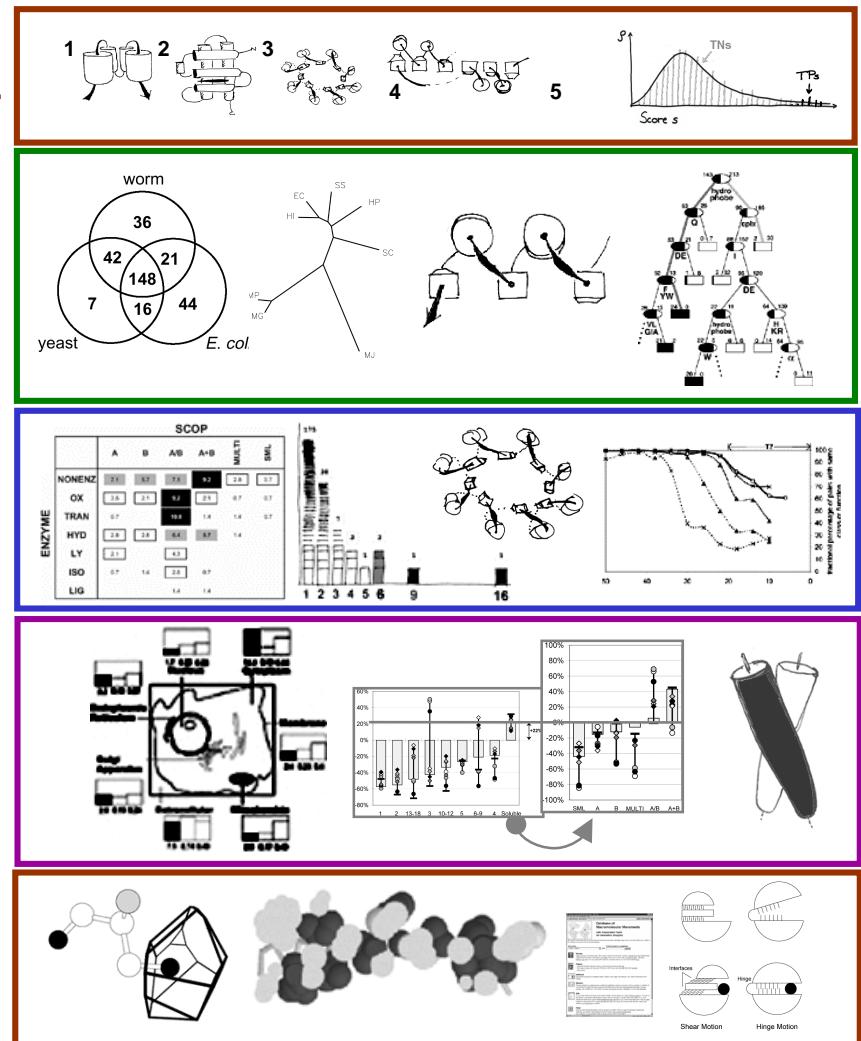
4 Folds in the Transcriptome

(Common parts? Where are parts?)

Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis, cyt. Depleted ↓ : NS, long, TM folds, transport, transcription, Leu-zip, nuc. Bayesian Localizer, phenotypes clustering

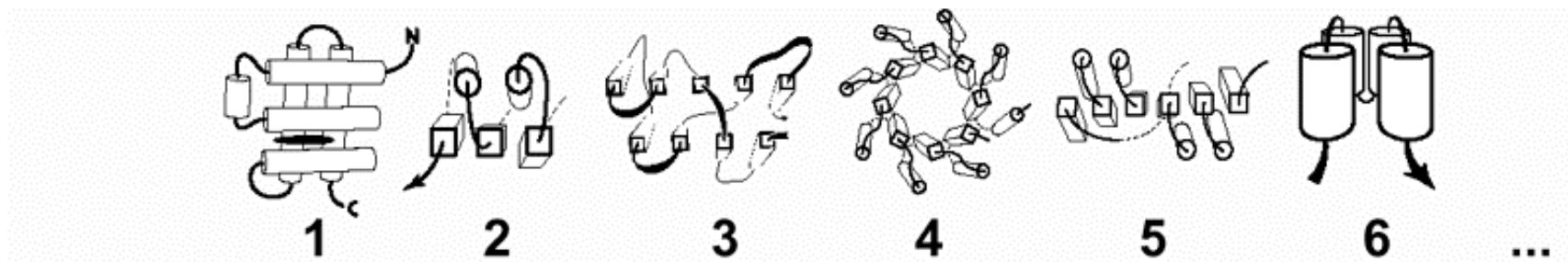
5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson, R Das, H Hegyi, J Lin, Y Kluger, C Arrowsmith, A Edwards, L Regan, S Balasubramanian, A Drawid, D Greenbaum, M Snyder, R Jansen



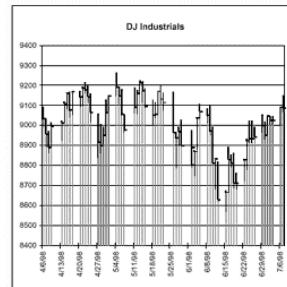
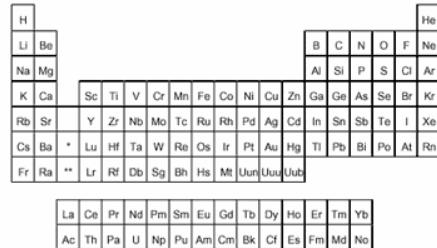
Fold Library vs. Other Fundamental Data structures

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



Folds in Molecular Biology 1000-10000

const.	mant.	exp.	unit
e	1.60	e	8 C
F	9.65	e	4 C/mol
ε_0	8.85	e	-12 F/m
μ_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
λ_C	2.43	e	-12 m
c	3.00	e	-19 m/s
G	6.67	e	-11 $m^3/kg \cdot s^2$
N_A	6.02	e	23 mol ⁻¹



10

100

**1000
-10000**

>1000000

Physics

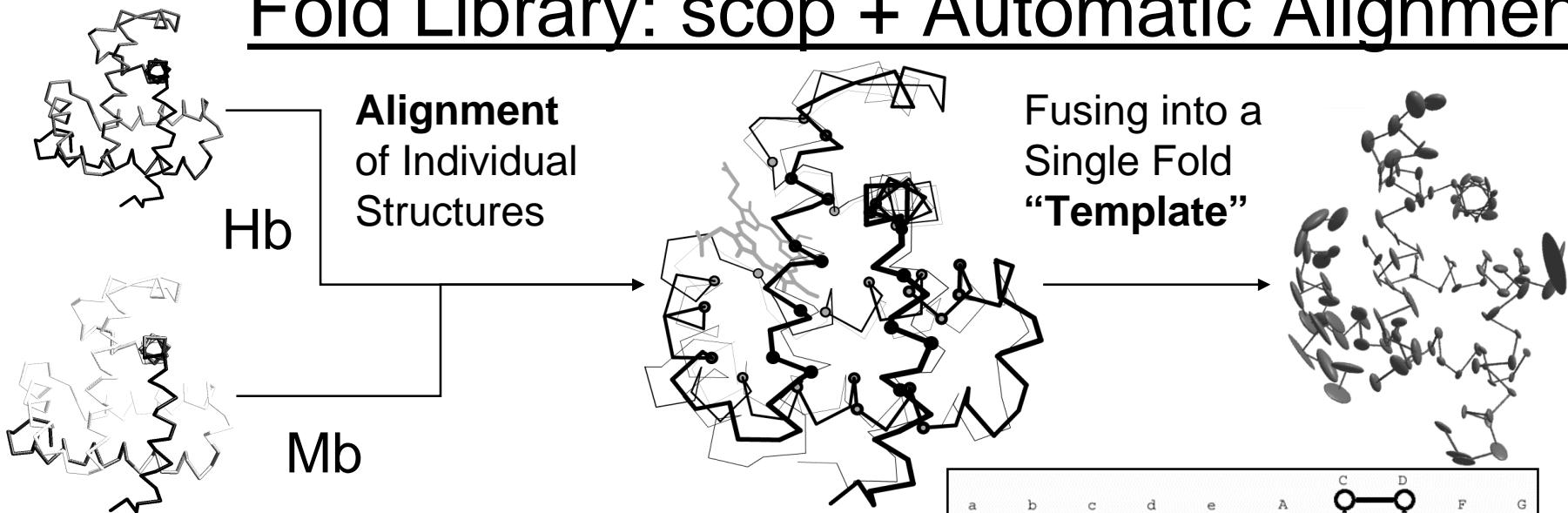
Chemistry

Finance

Politics

(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.

Fold Library: scop + Automatic Alignments



Iterative Dynamic Programming

Like repeated sequence alignment; 1 cycle doesn't converge, violates key assumption of D.P.

ACSQRP--LRV-SH	-R	SENCV
A-SNKPQLVKLMTH	VK	DFCV-

Derived from Program of G Cohen (Align, Satow et al. 1986); scop (Murzin et al., 1995). Much Previous work: Remington, Matthews '80; **Taylor, Orengo '89, '94**; Artymiuk, Rice, Willett '89; Sali, Blundell, '90; Vriend, Sander '91; Russell, Barton '92; **Holm, Sander '93**; Godzik, Skolnick '94; Gibrat, Madej, Bryant '96; Falcov, F Cohen, '96; Feng, Sippl '96; G Cohen '97; Singh & Brutlag, '98

The figure shows three stages of iterative dynamic programming:

- Initial Equivalences:** A set of initial sequence equivalences between two sets of residues (a-e and A-G).
- First Iteration:** Shows the first step of aligning sequences. The score is 57, Nbrk is 2, and RMS is 1.96.
- Second Iteration:** Shows the second step of aligning sequences. The score is 91, Nbrk is 1, and RMS is 0.65.
- Third Iteration:** Shows the final step of aligning sequences. The score is 100, Nbrk is 1, and RMS is 0.23.

Score matrices for each iteration:

	A	B	C	D	E	F	G
a	7	5	9	2	1	0	0
b	2	9	12	9	7	2	0
c	1	2	2	10	12	8	2
d	0	1	1	2	2	13	7
e	0	0	0	0	1	2	13

	A	B	C	D	E	F	G
a	19	4	4	1	1	0	0
b	4	16	16	4	4	1	0
c	1	4	4	14	18	4	1
d	0	1	1	4	4	19	4
e	0	0	0	1	1	4	19

	A	B	C	D	E	F	G
a	20	4	3	1	1	0	0
b	4	20	12	4	4	1	0
c	1	4	4	11	20	4	1
d	0	1	1	4	4	20	4
e	0	0	0	1	1	4	20

Some Similarities are Readily Apparent others are more Subtle

Easy:

Globins

125 res.,
~1.5 Å



Tricky:

Ig C & V

85 res.,
~3 Å



Very Subtle: G3P-dehydrogenase, C-term. Domain
>5 Å



Some Similarities are Readily Apparent others are more Subtle

Easy:

Globins

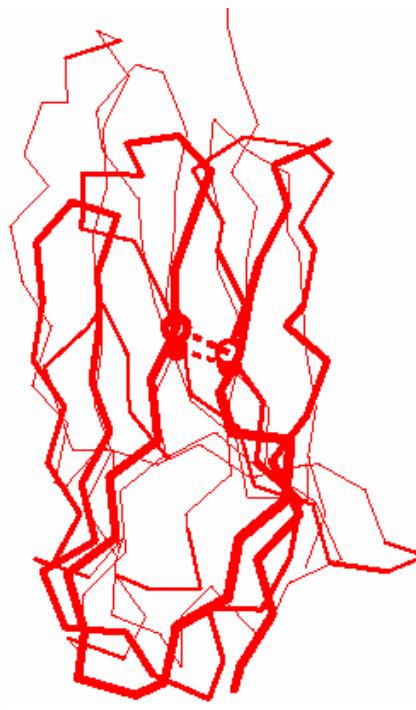
125 res.,
~1.5 Å



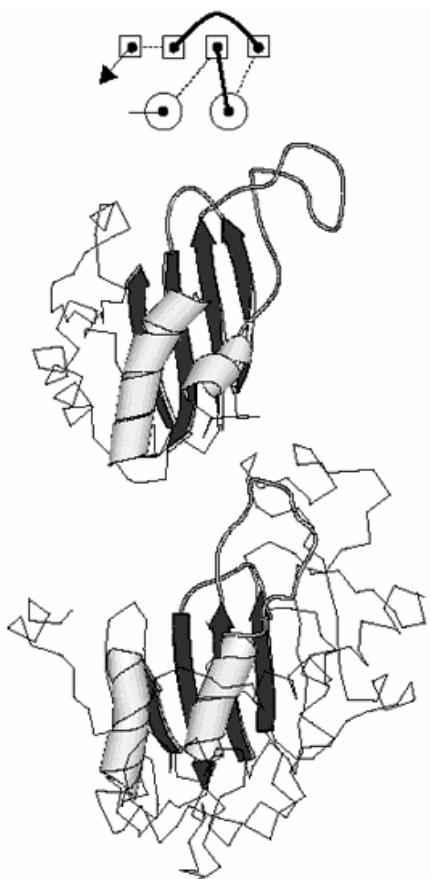
Tricky:

Ig C & V

85 res.,
~3 Å



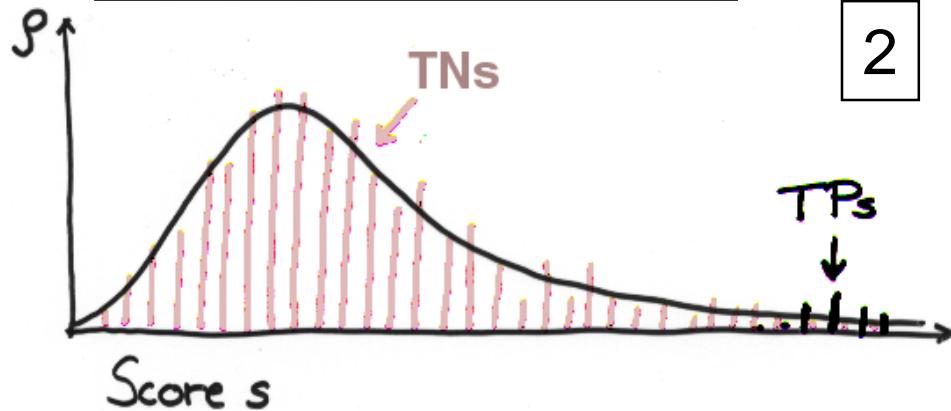
Very Subtle: G3P-dehydrogenase, C-term. Domain
>5 Å



* d1ahn					
* d1wab		0.7 TP			
* d5tima		4.9	7.1		
d4tima					
* d1htia		5.1	6.2	1.2 TP	
d1igtb1		5.1	6.0	2.1 FP	4.4
		*	*	*	
		d1ahn	d1wab	d5tima	d4tima
					d1htia
					d1igtb1

P-values

1



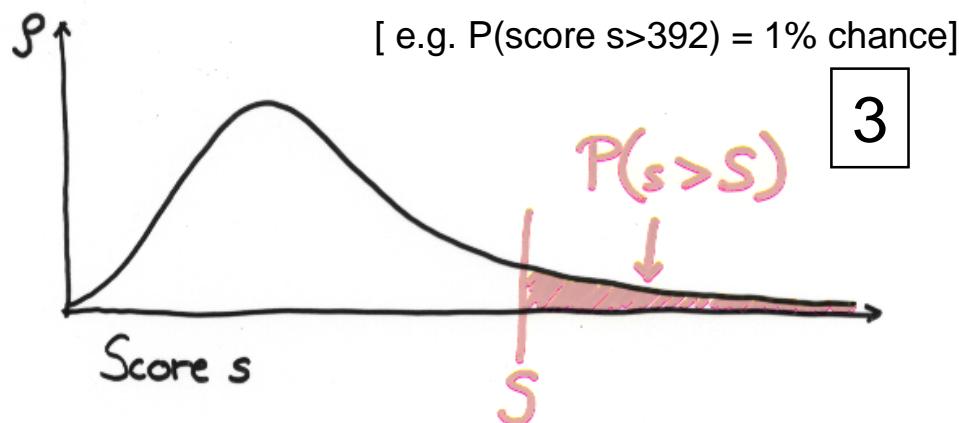
2

- Significance Statistics

- ◊ For sequences, originally used in Blast (Karlin-Altschul). Then in FASTA, &c.
- ◊ Extrapolated Percentile Rank: How does a Score Rank Relative to all Other Scores?

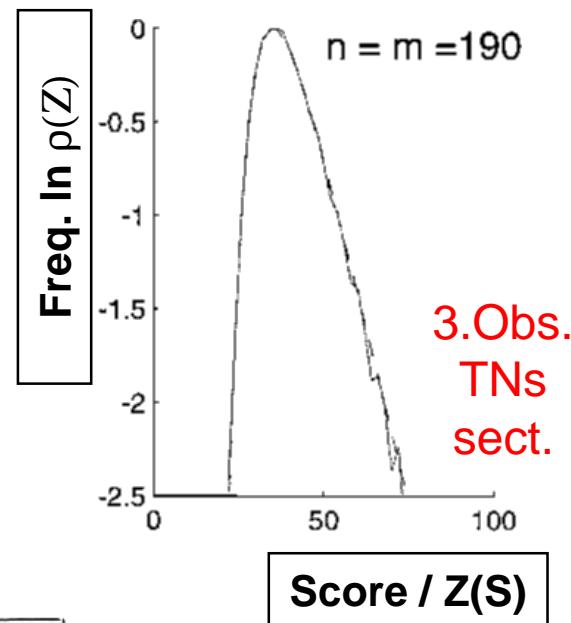
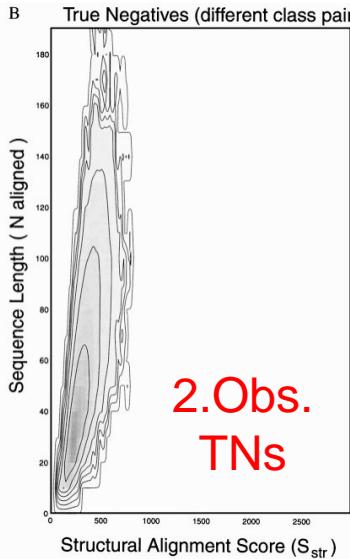
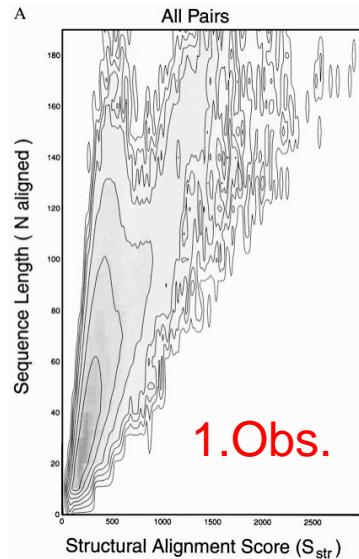
- Our Strategy: Fit to Empirical Distributions

- 1) **All-vs-All** comparison
- 2) **Graph Distribution of Scores**
in 2D (N dependence); 1K x 1K families -> ~1M scores; ~2K included TPs
- 3) **Fit a function $\rho(s)$**
to TN distribution (TNs from scop); Integrating ρ gives $P(s>S)$, chance of getting a score better than threshold S randomly
- 4) **Same formalism for sequence & structure**



EVD Fits

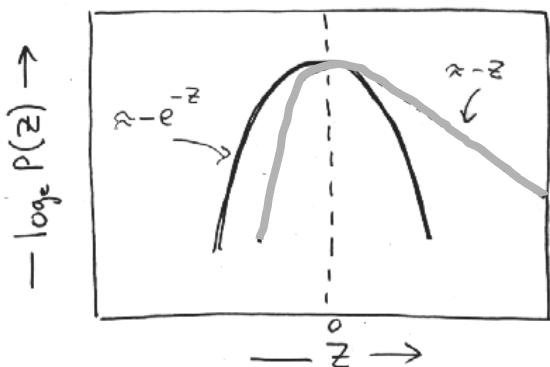
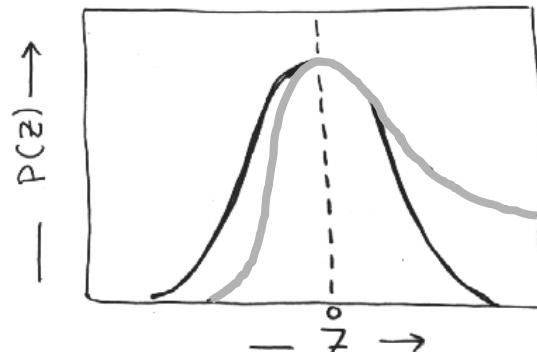
3-parm fit works;
Reasonable as
Dynamic
Programming
maximizes
over pseudo-
random variables



5. Extreme Value v. Normal Distributions:

EVD ~
Max(indep.
random
variables),

Normal ~
Sum(indep.
random
variables)
[$\exp(-z^2)$]



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

4.

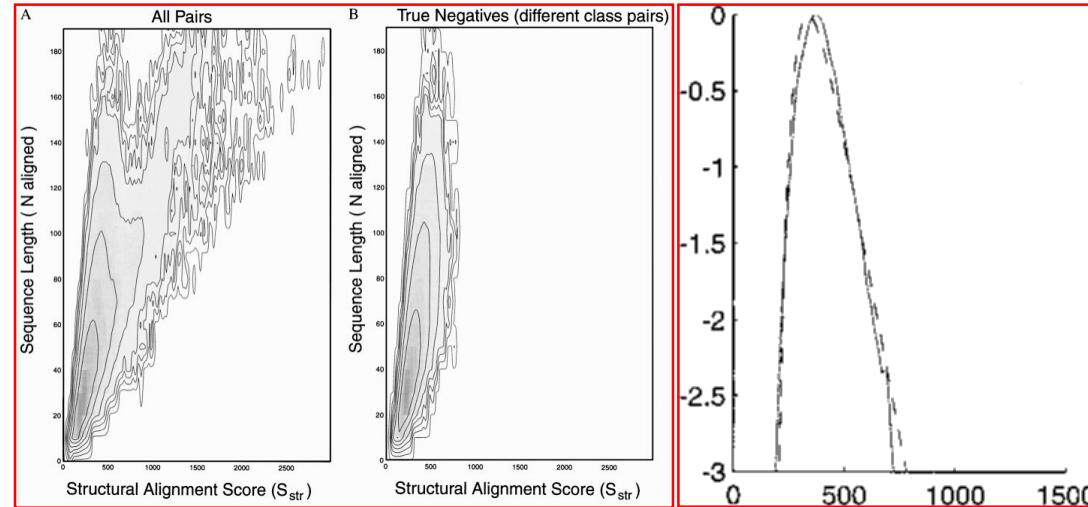
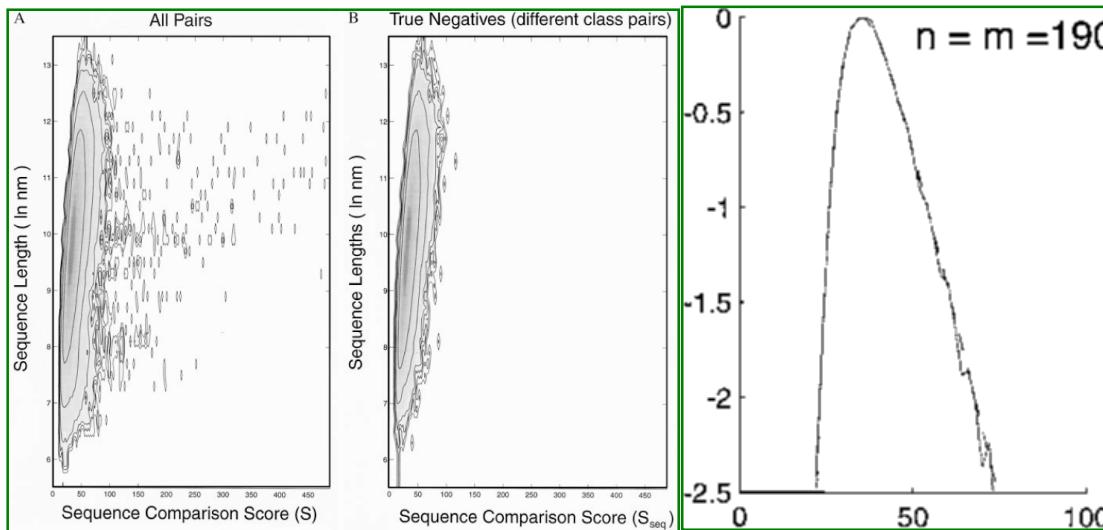
$$Z = \frac{S - (a \ln N + b)}{\sigma}$$

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

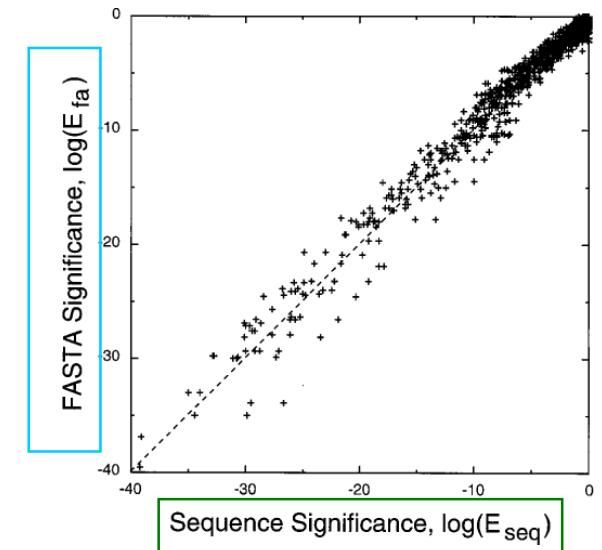
3 Free Parm. fit to EVD matches distributions best
Fit involves: a, b, σ . N = number of residues
matched; G = total gap penalty; $M(i,j)$ = similarity
matrix (Blossum for seq. or $M_{str}(i,j)$, struc.)

Same EVD Results for Sequence & Struct.

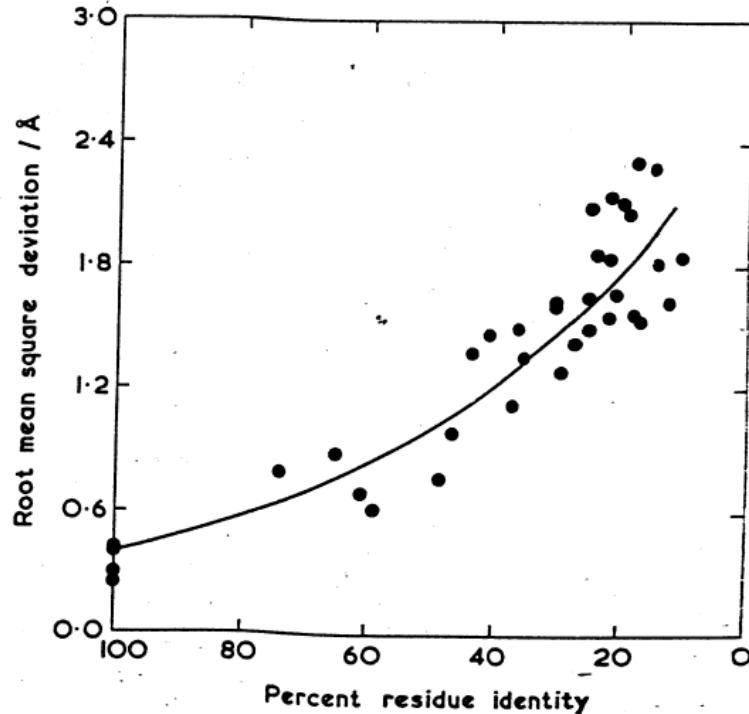
Sequence



- Use Sequence Scores to Validate
- Sequence P-value perfectly tracks FASTA e-value



Significance computation can be applied to **any** existing sequence or structure alignment; added Benefit: allows computation of an e-value without doing a db run



Chothia & Lesk,
1986 revisited,
32 → 16K pts.

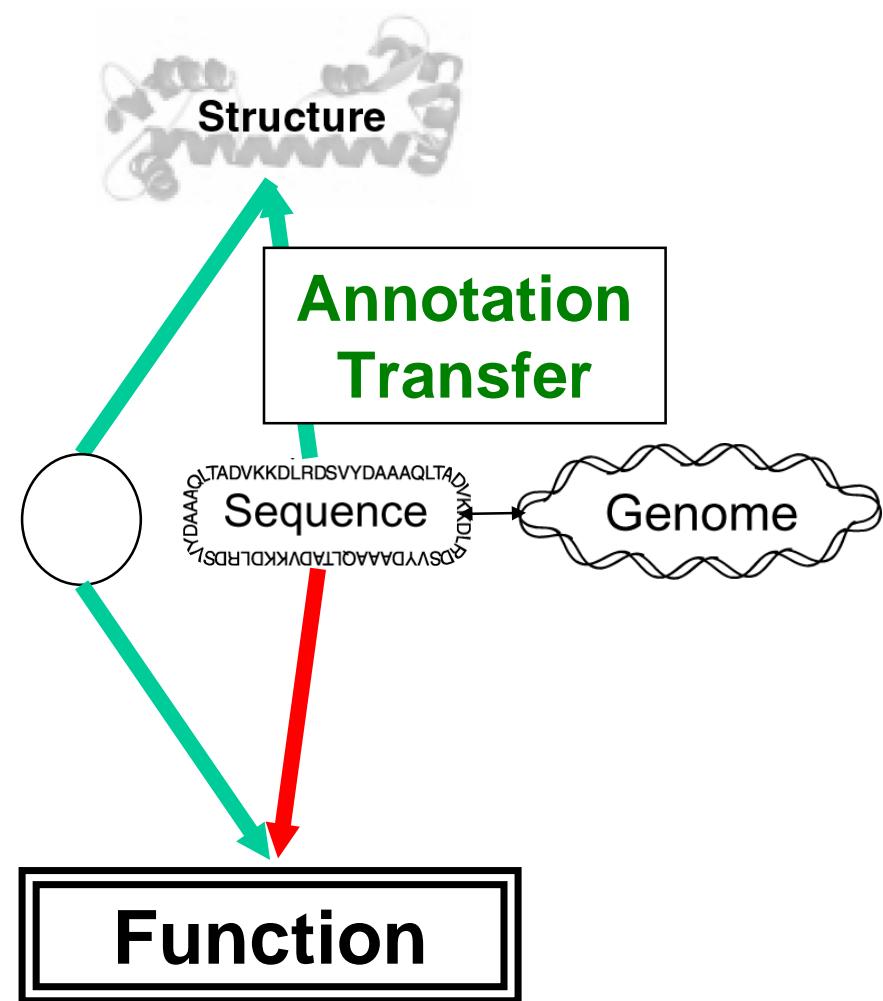
C&L $\Delta = .4 \exp(1.9 H)$

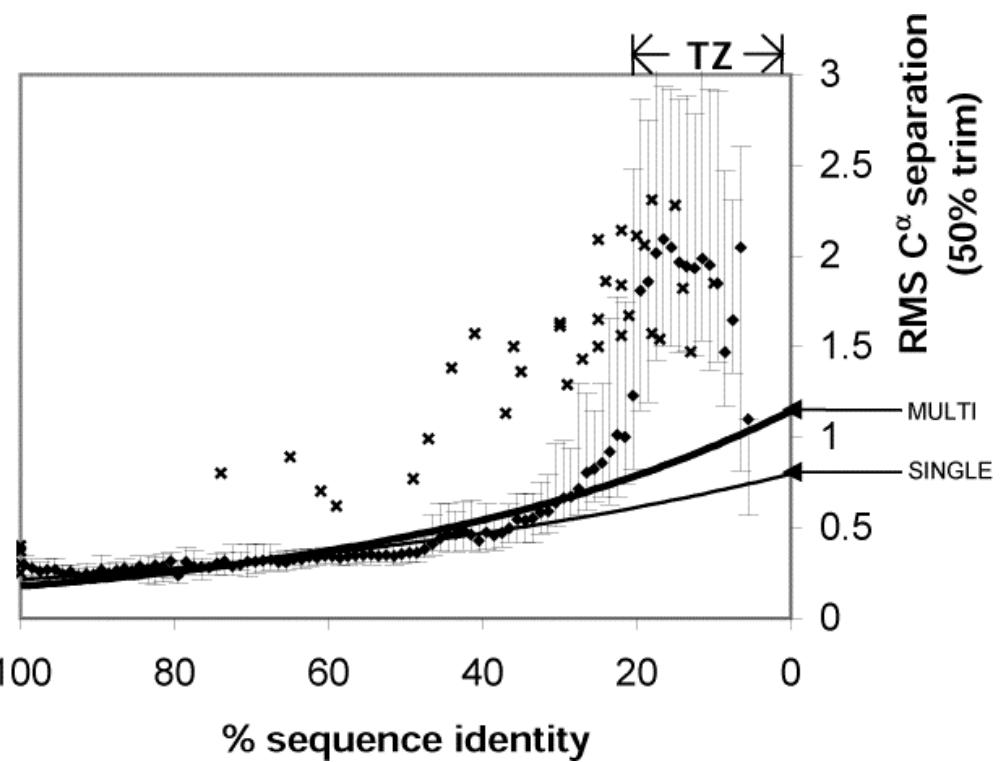
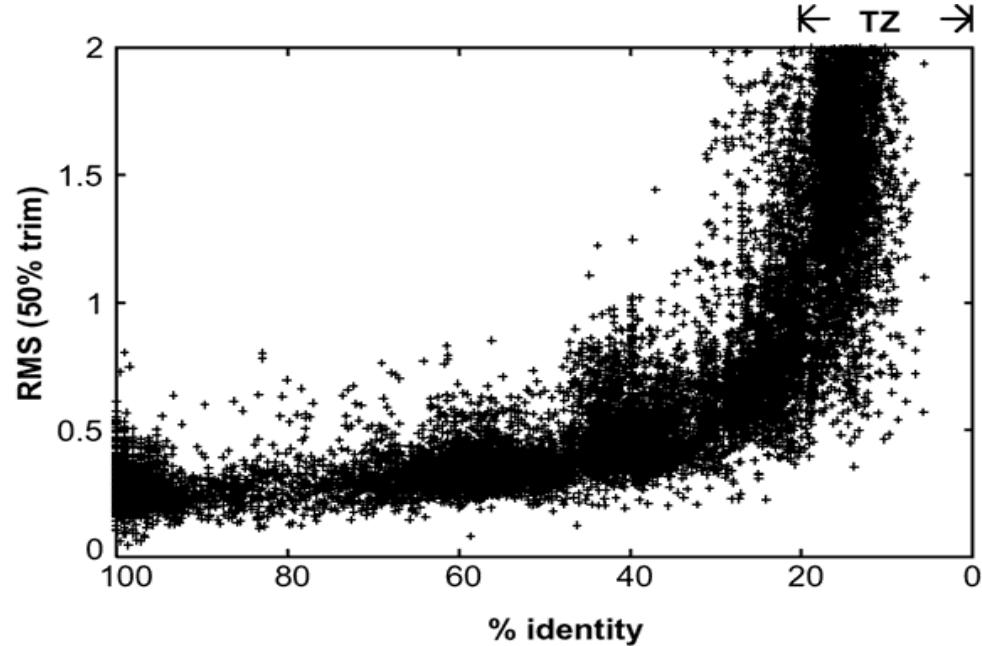
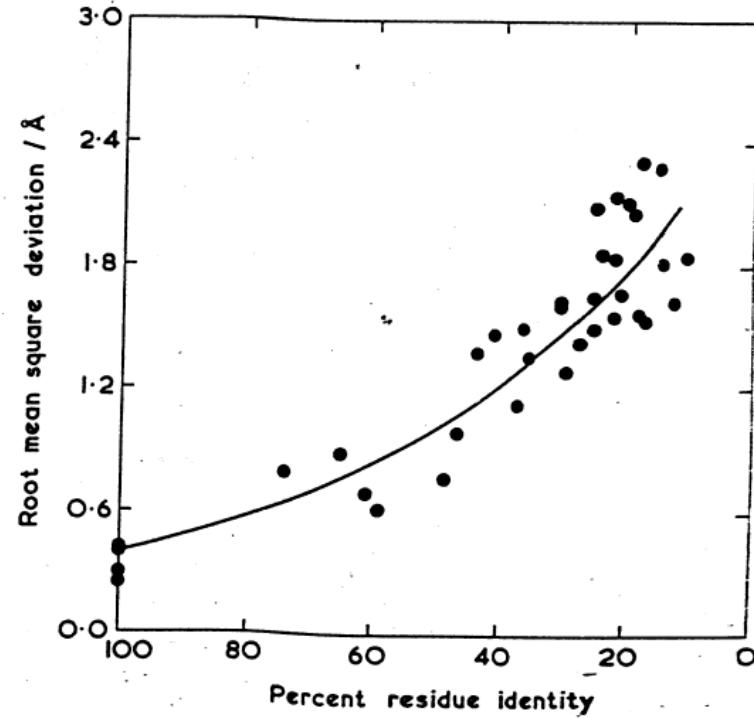
Multi: $\Delta = .2 \exp(1.3 H)$

Single: $\Delta = .2 \exp(1.9 H)$

Related work done by Thornton, Pearson, Brenner

What's the use of sequence-structure scoring schemes?
Precise Annotation Transfer for Genomics





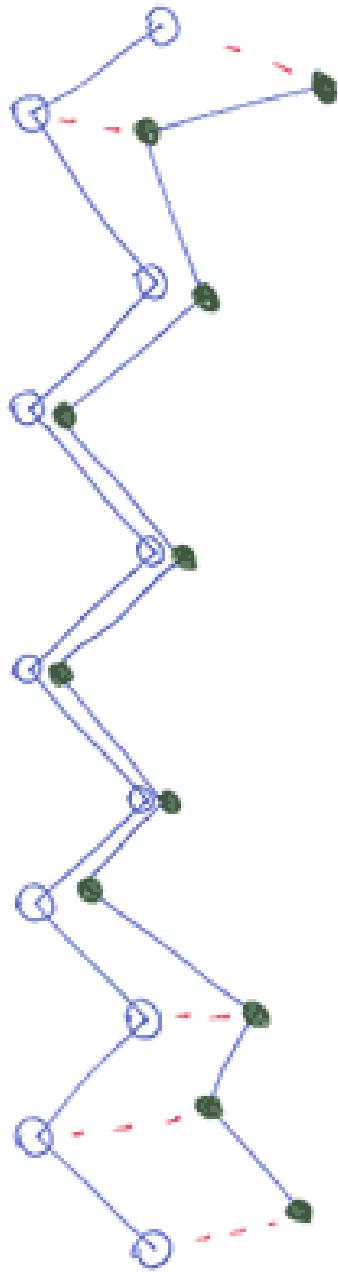
Chothia & Lesk, 1986 revisited, $32 \rightarrow 16K$ pts.

C&L $\Delta = .4 \exp(1.9 H)$

Multi: $\Delta = .2 \exp(1.3 H)$

Single: $\Delta = .2 \exp(1.9 H)$

Related
work done
by Thornton,
Pearson,
Brenner



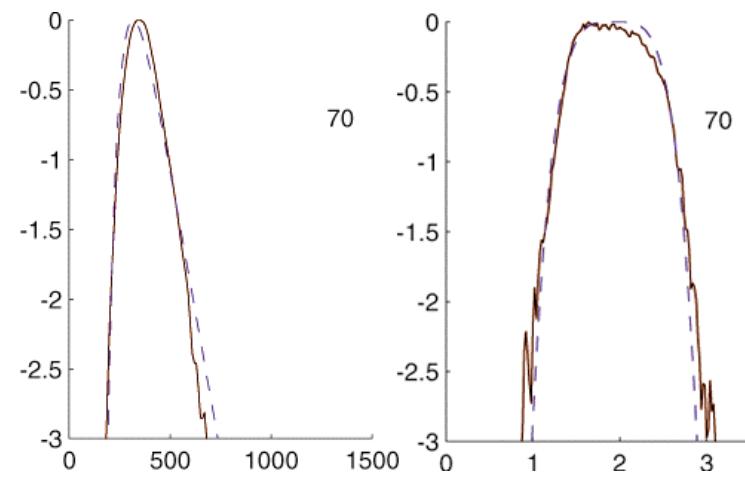
Problems with RMS and %ID

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

S_{str} RMS

$$\sum \frac{100}{5 + \mathbf{d}_i^2} vs \quad \sqrt{\sum \mathbf{d}_i^2}$$

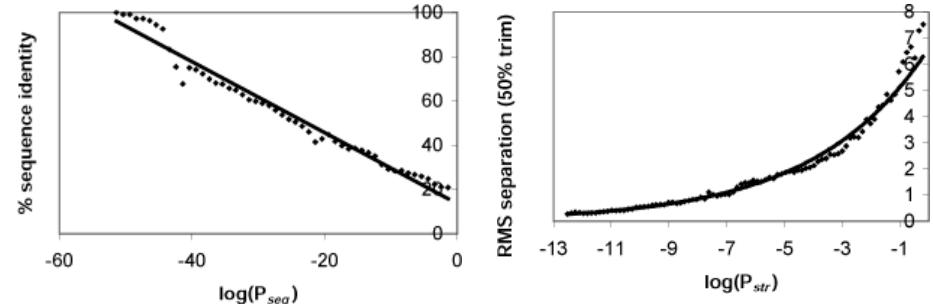
%ID problem:
“Bunching up”
between 20%
and 0%
identity



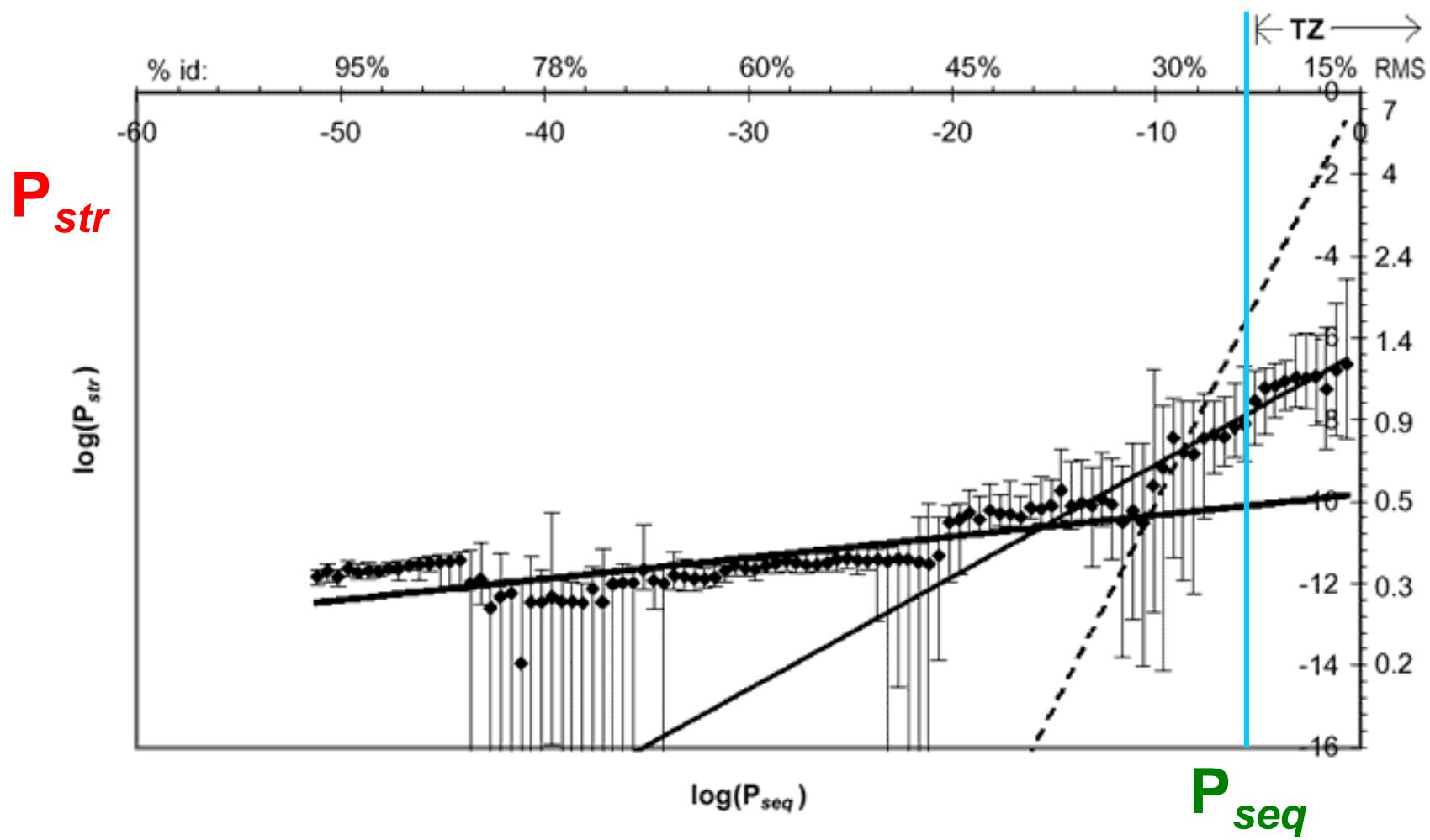
Modern statistical language

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

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



overcomes length dependency



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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

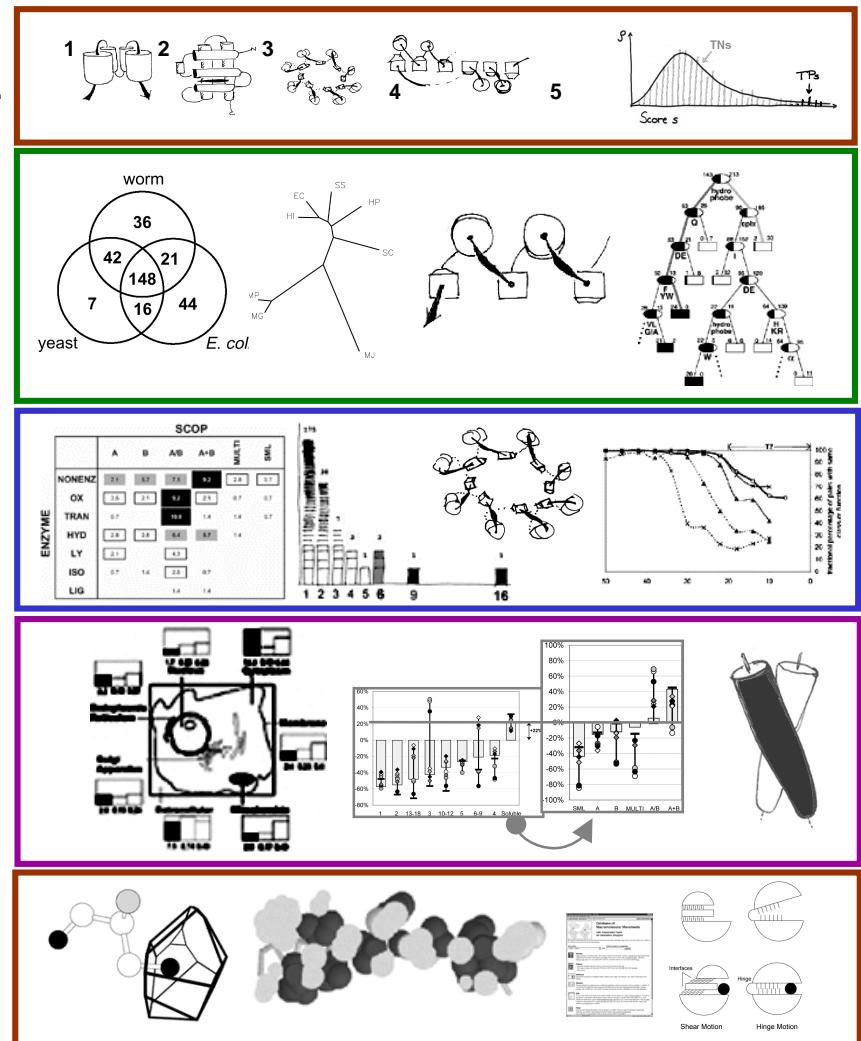
4 Folds in the Transcriptome

(Common parts? Where are parts?)

Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis, cyt. Depleted ↓ : NS, long, TM folds, transport, transcription, Leu-zip, nuc. Bayesian Localizer, phenotypes clustering

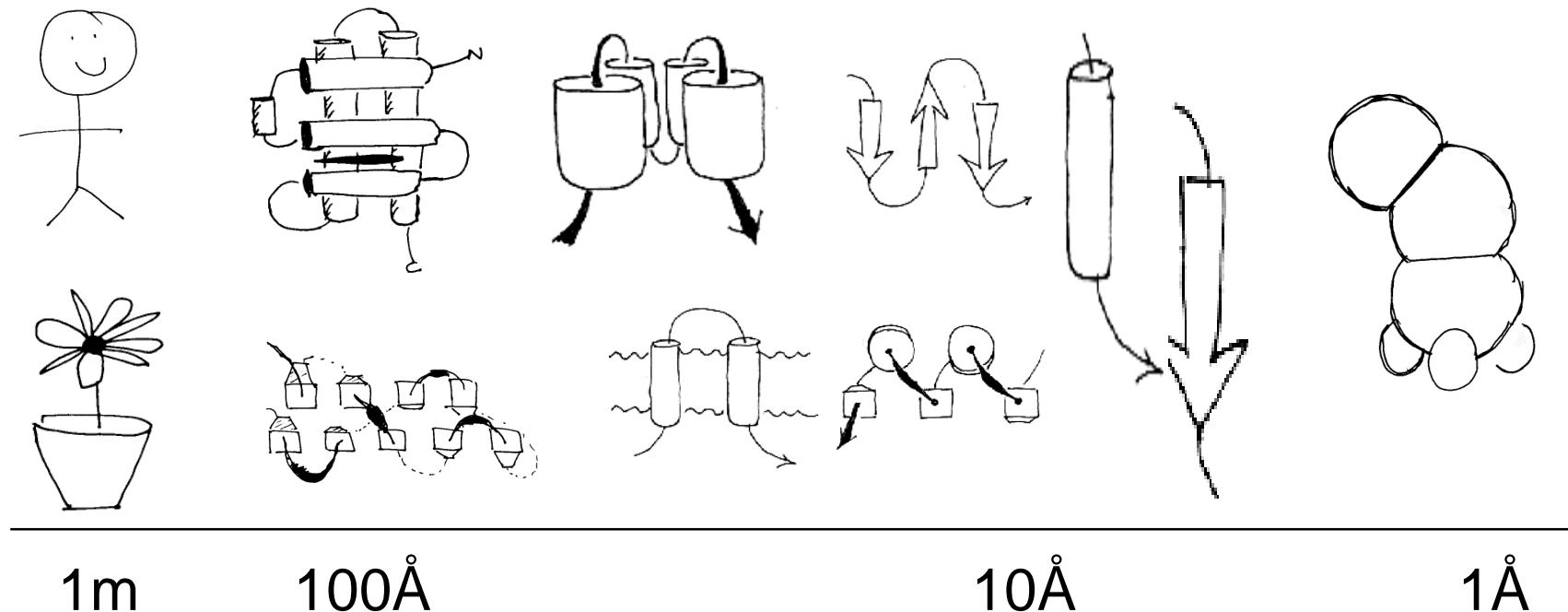
5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson, R Das, H Hegyi, J Lin, Y Kluger, C Arrowsmith, A Edwards, L Regan, S Balasubramanian, A Drawid, D Greenbaum, M Snyder, R Jansen

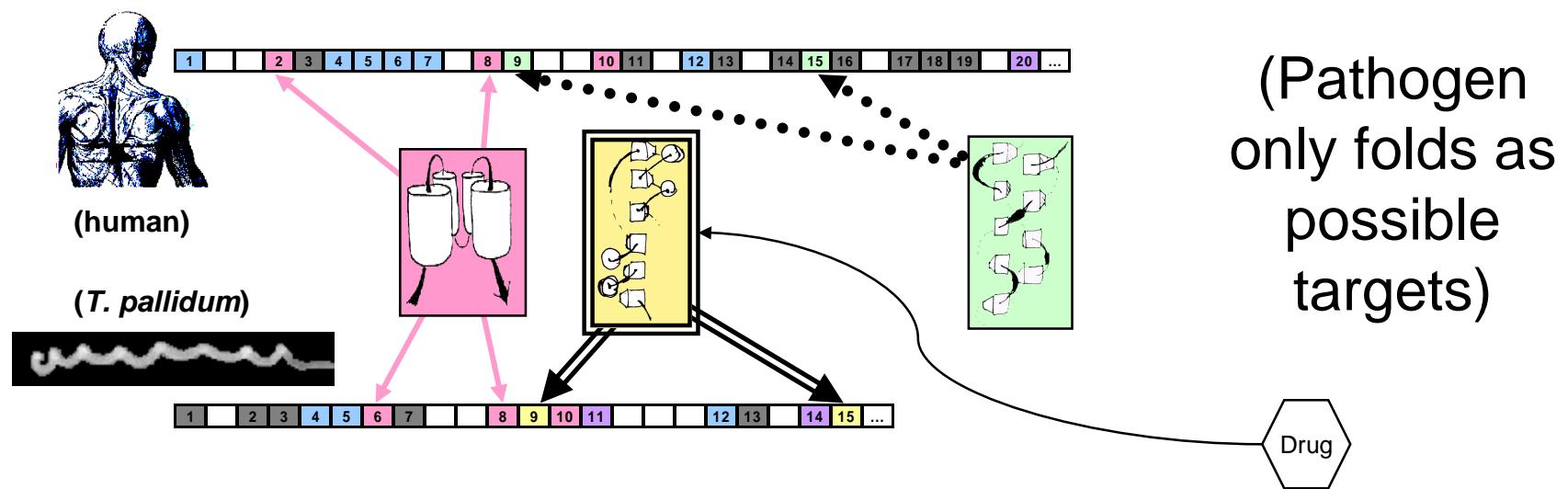


At What Structural Resolution Are Organisms Different?

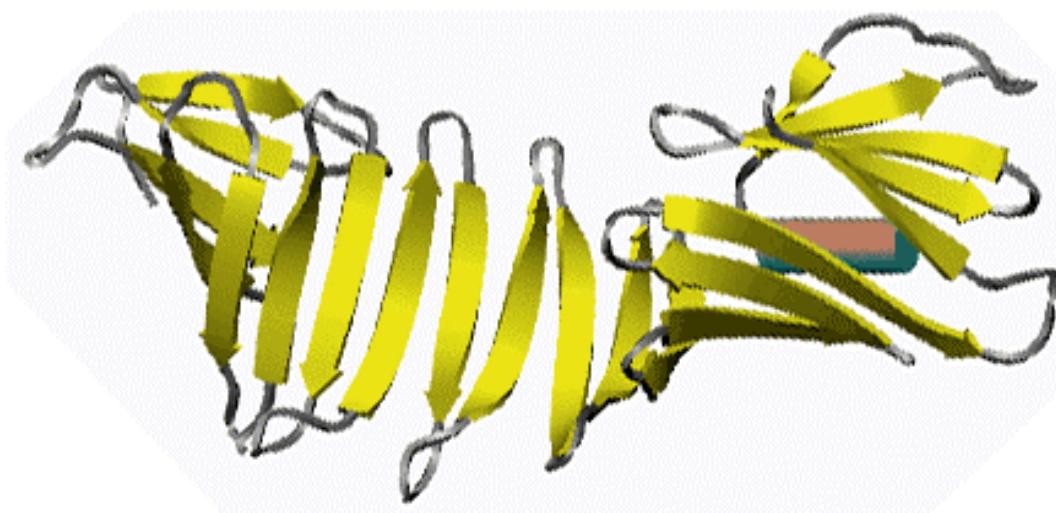
person	protein	super-secondary	helix	individual
plant	fold (Ig)	structure ($\beta\beta$,TM– TM, $\alpha\beta\alpha\beta,\alpha\alpha\alpha$)	strand	atom (C,H,O...)



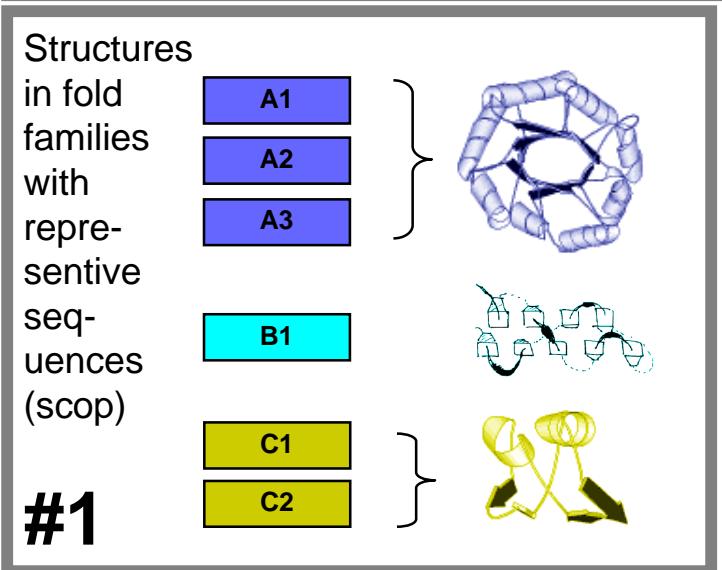
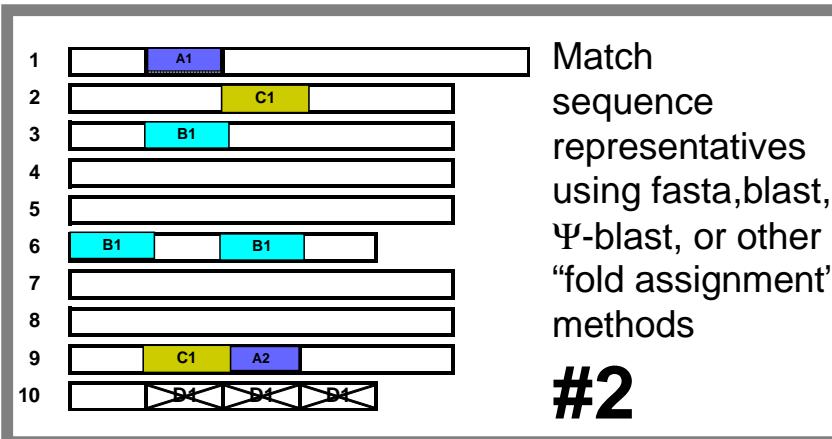
Practical Relevance of Structural Genomics



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

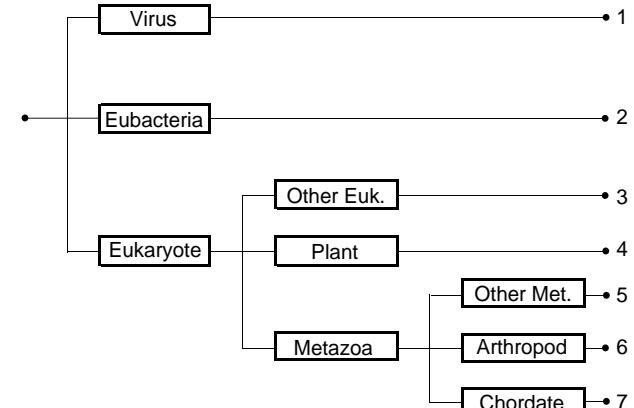


Cross-Reference: Folds→Sequences → Organisms



Organize Sequences by Genome or Taxon

#3



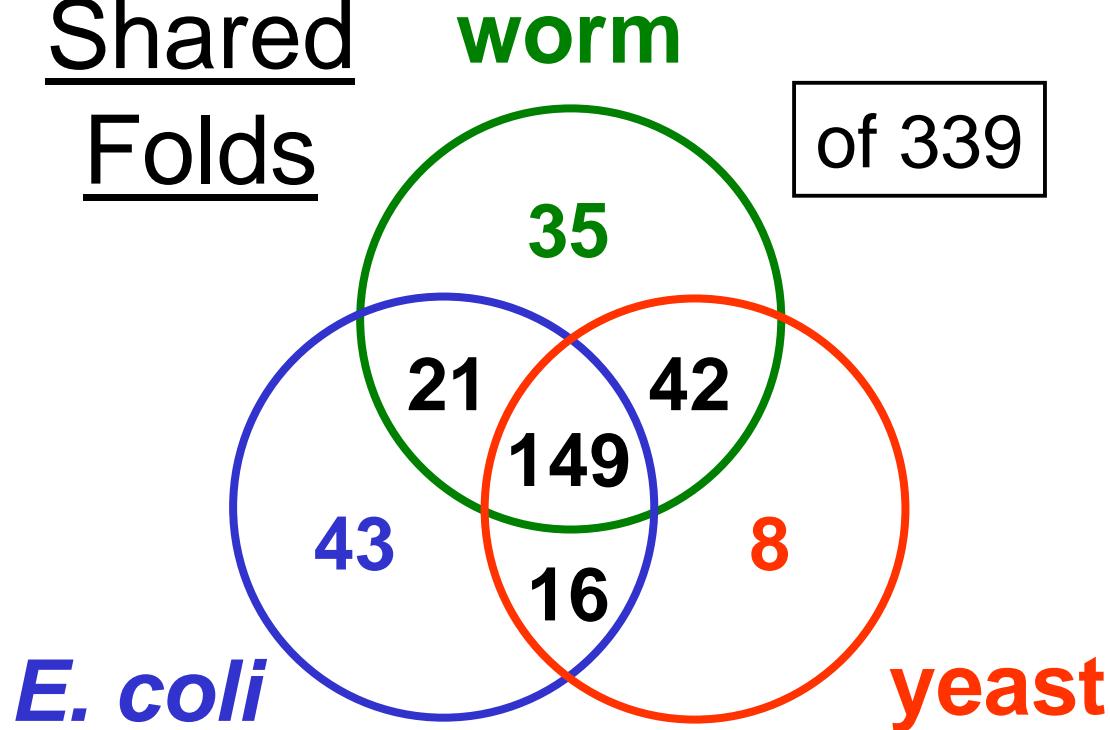
Abbrev.	Kingdom (subgroup)	Genome	Num. ORFs	Reference
EC	Bacteria (gram negative)	<i>Escherichia coli</i>	4290	Blattner et al.
HI	Bacteria (gram negative)	<i>Haemophilus influenzae</i>	1680	TIGR
HP	Bacteria (gram negative)	<i>Helicobacter pylori</i>	1577	TIGR
MG	Bacteria (gram positive)	<i>Mycoplasma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

Tabulate Results in Database

#4

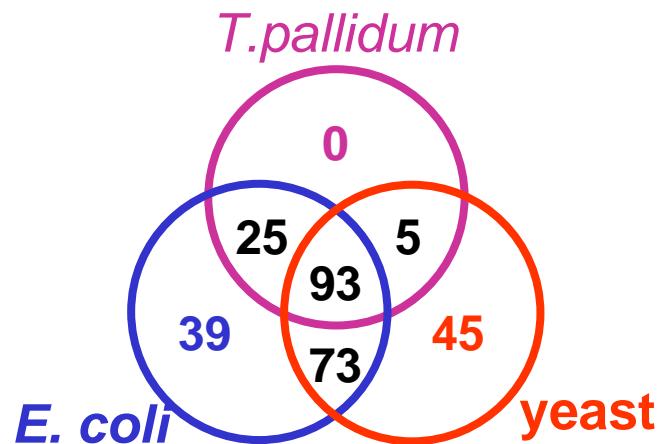
class	Fold#	EC	SC	HI	SS	HP	MJ	MP	MG	total	Fam.	PDB	Rep.	Struc.	Name
α/β	18	60	46	23	40	19	7	4	3	202	16	183	1xel	-	NAD(P)-binding Rossmann Fold
α/β	24	20	69	17	19	17	16	10	11	179	13	132	1gky	-	P-loop Containing NTP Hydrolases
α+β	31	37	28	18	16	12	40	3	3	157	23	160	1fxd	-	like Ferrodoxin
α/β	01	45	36	13	22	11	10	5	4	146	37	399	1byb	-	TIM-barrel
α/β	23	18	17	7	9	4	8	2	2	67	5	36	1pyd	a:2-181	Thiamin-binding
α/β	04	15	11	7	10	1	9	5	5	63	13	132	2tmd	a:490-645	FAD/NAD(P)-binding
α+β	55	8	9	7	8	9	3	6	6	56	4	23	1sry	a:111-421	Class-II aaRS/Biotin Synthetase
β	27	7	10	8	8	4	4	3	3	47	5	19	1fnb	19-154	Reductase/Elongation Factor D
β	24	13	7	4	3	3	3	3	3	39	18	177	1snc	-	OB-fold
α+β	11	10	8	4	8	2	2	2	1	37	11	48	1lgd	-	beta-Grasp
β	55	9	10	5	5	2	2	2	2	37	7	19	1bdo	-	Barrel-sandwich hybrid
α/β	15	5	5	4	4	5	6	3	3	35	3	22	2ts1	1-217	ATP pyrophosphatases

Shared Folds

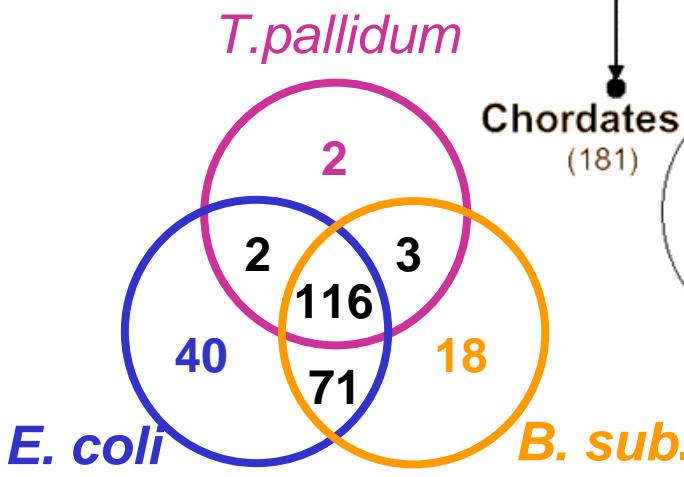


E. coli

yeast

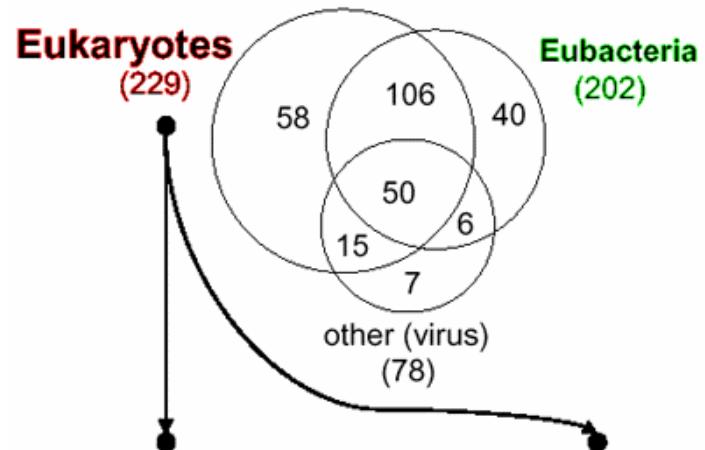


E. coli



E. coli

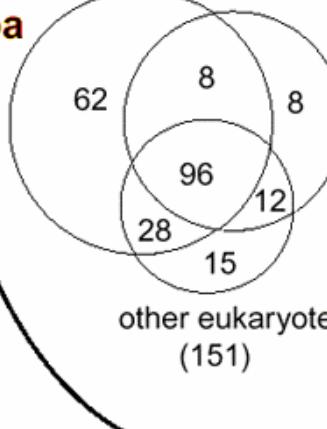
B. subtilis



Metazoa

(194)

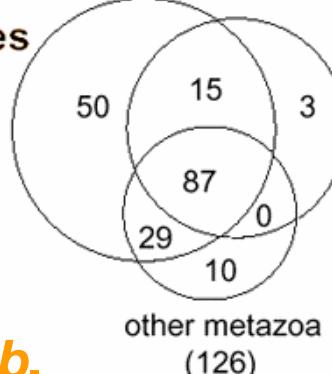
Plants
(124)



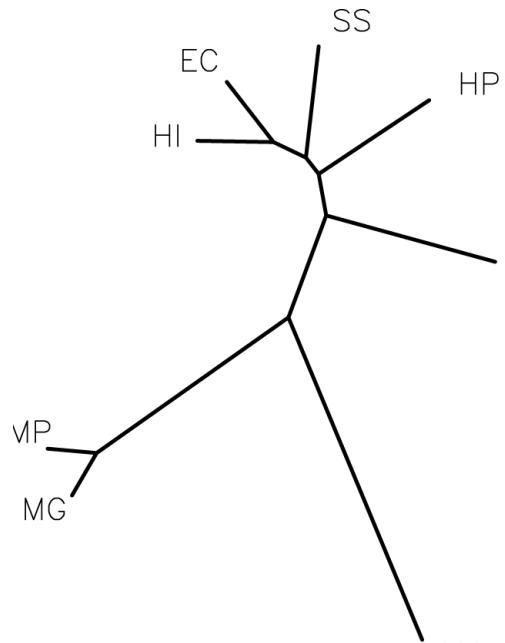
Chordates

(181)

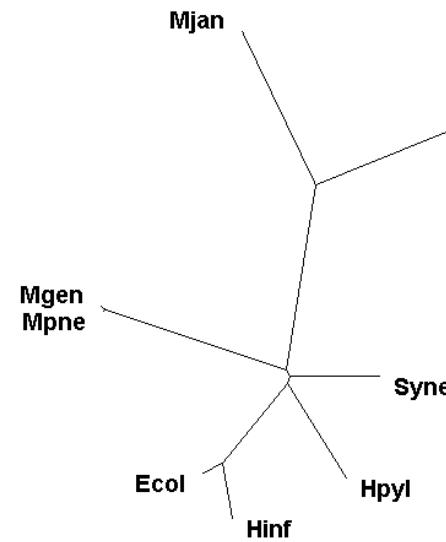
Arthropods
(105)



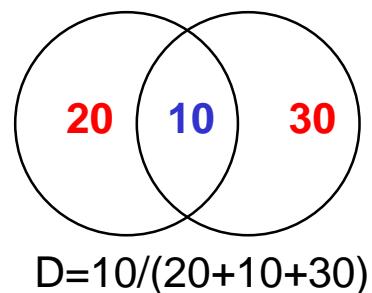
Cluster Trees Grouping Initial Genomes on Basis of Shared Folds



Fold Tree



“Classic” Tree

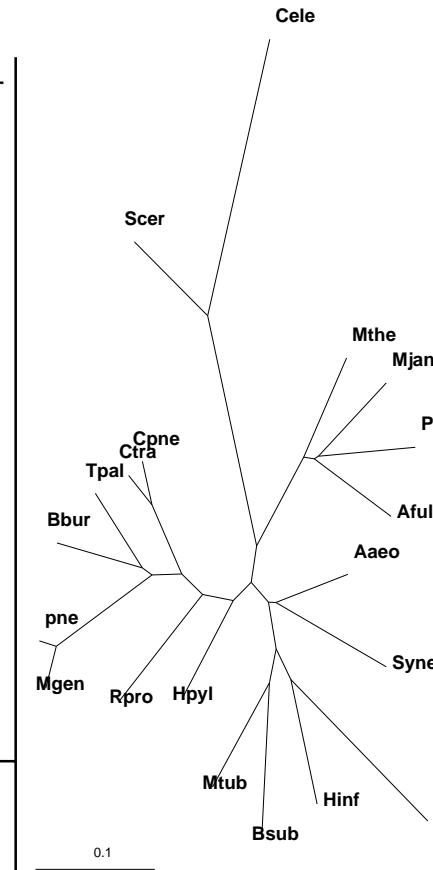


$$D = S/T$$

S = # shared folds

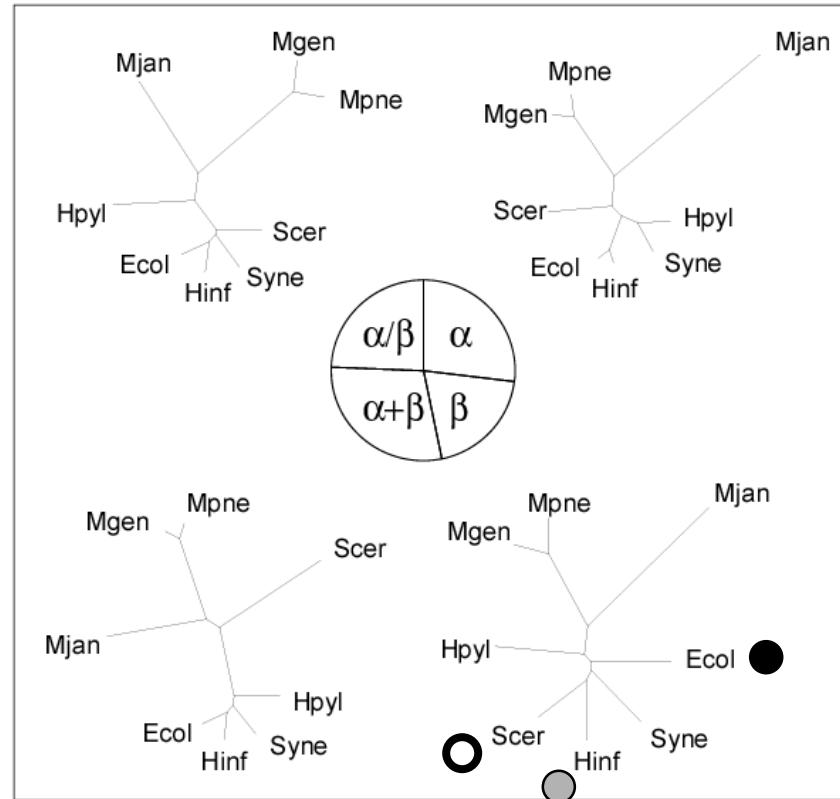
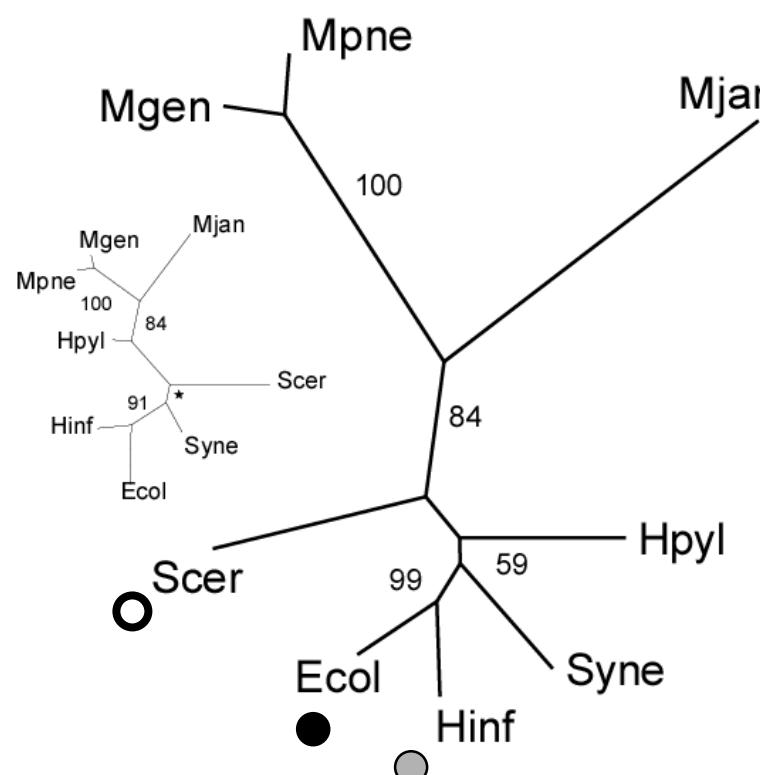
D = shared fold dist.
betw. 2 genomes

T = total #
folds in both



20 Genomes

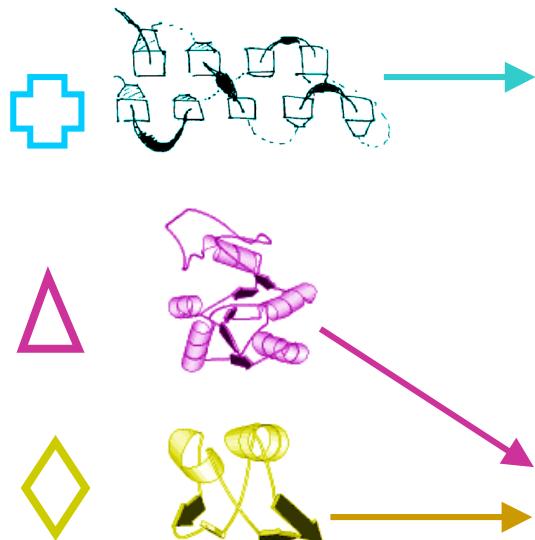
Distribution of Folds in Various Classes



Unusual distribution of all-beta folds

Common Folds in Genome, Varies Betw. Genomes

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



			num. matches in worm genome (N)	frac. all worm dom. (F)	in EC?	in SC?
	Ig	class	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%		
	Ferrodoxin	A+B	207	0.4%		

Rank	M.gen	B.sub	E.col
1	P-loop hydrolase	60	P-loop hydrolase
2	SAM methyl-transferase	16	Rossmann domain
3	Rossmann domain	13	Phosphate-binding barrel
4	Class I synthetase	12	PLP-transferase
5	Class II synthetase	11	CheY-like domain
6	Nucleic acid binding dom.	11	SAM methyl-transferase
Total ORFs	479	4268	4268
with Common Superfamilies	(22%)	(11%)	(11%)

Eubacteria

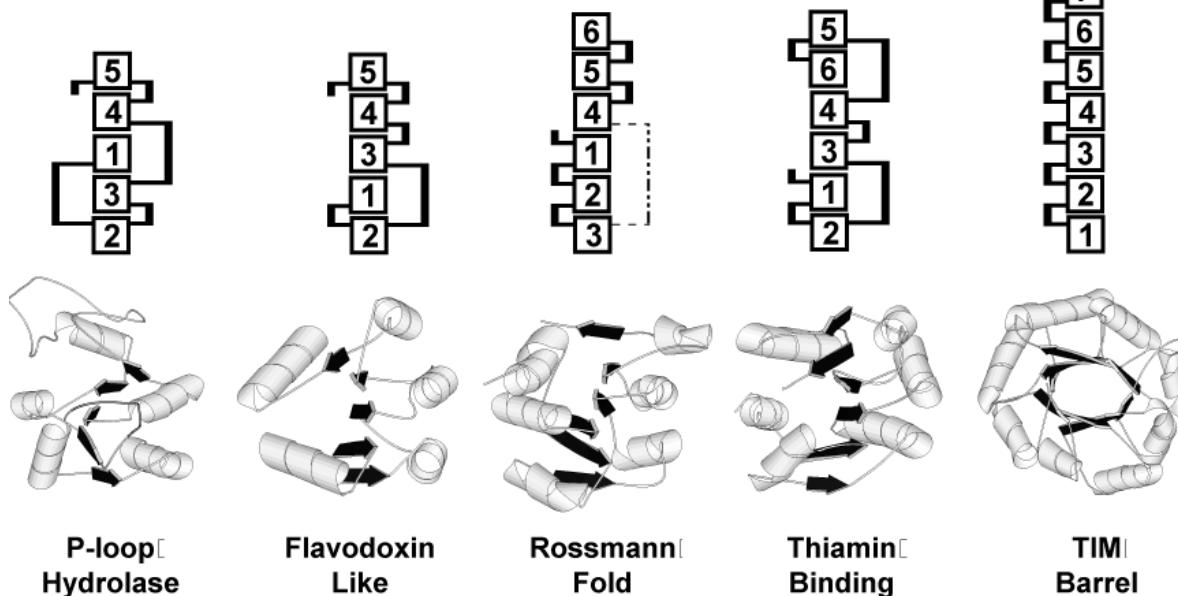
Rank	M.the	A.ful
1	P-loop hydrolase	P-loop hydrolase
2	Phosphate-binding barrel	Rossmann domain
3	Rossmann domains	Phosphate-binding barrel
4	Ferrodoxins	Ferrodoxins
5	SAM methyl-transferase	SAM methyl-transferase
6	PLP-transferases	PLP-transferases
Total ORFs	1869	2409
with Common Superfamilies	(14%)	(13%)

Archaea

Rank	S.cer
1	P-loop hydrolase
2	X Protein kinase
3	Rossmann domain
4	RNA-binding domain
5	SAM methyl-transferase
6	Ribonuclease H-like
Total ORFs	6218
with Common Superfamilies	(9%)

Yeast

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



HI, MJ, SC vs scop 1.32

42

ARTICLES

NATURE VOL. 336 3 NOVEMBER 1998

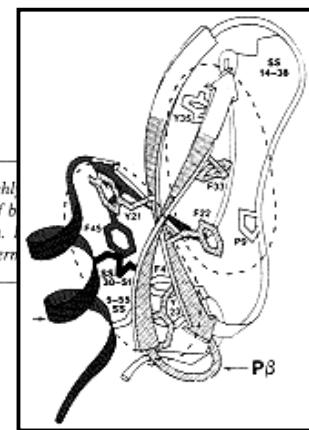
A peptide model of a protein folding intermediate

Terrence G. Oas & Peter S. Kim

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

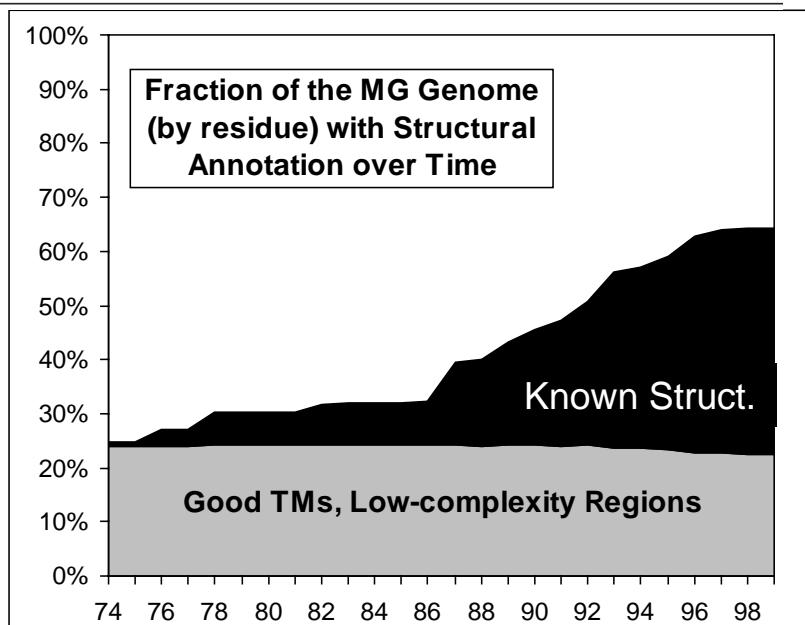
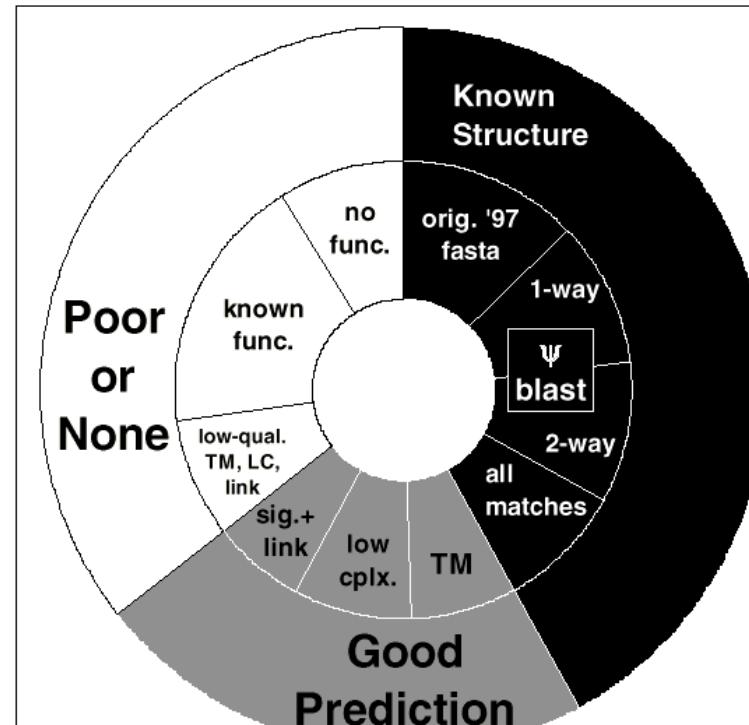
It is difficult to determine the structures of protein folding intermediates because folding is a highly disulphide-bonded peptide pair, designed to mimic the first crucial intermediate in the folding of b inhibitor, contains secondary and tertiary structure similar to that found in the native protein.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

336: 42

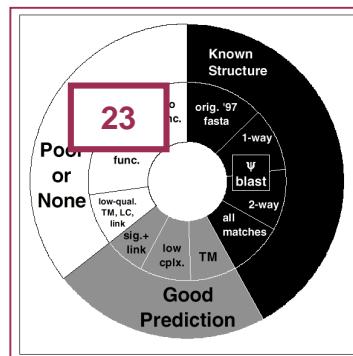
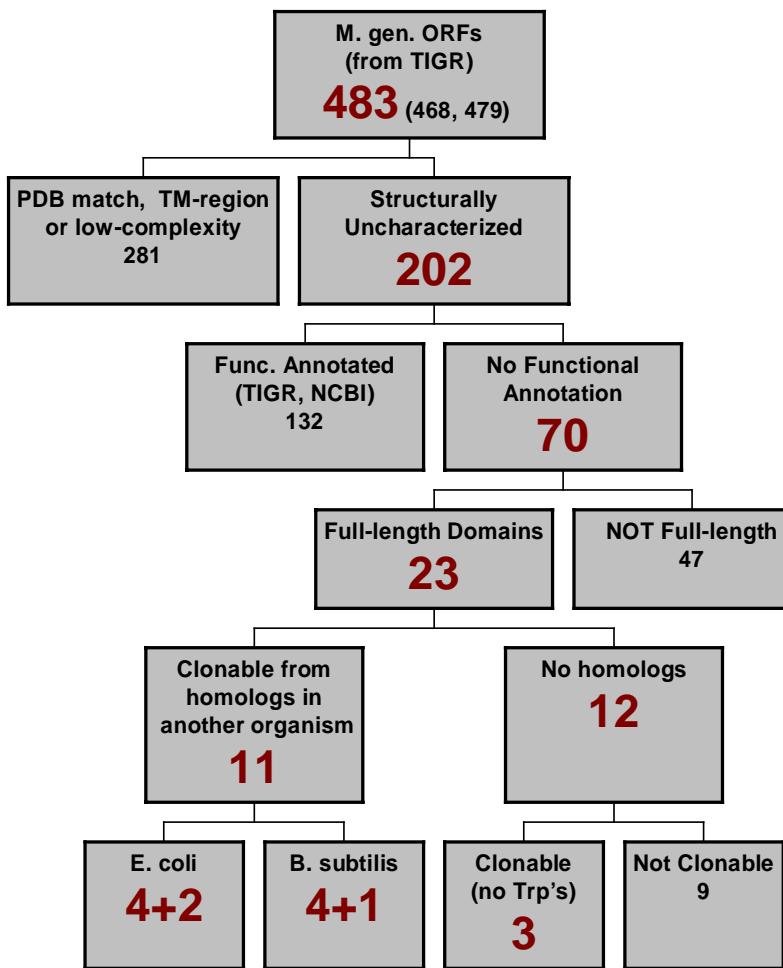


Bias Problem → Prediction, Expts.

- Known Structures are Incomplete, Biased Sample from Genome, so...
 - Resample
 - Expt. Structural Genomics...
 - Predict Structures...

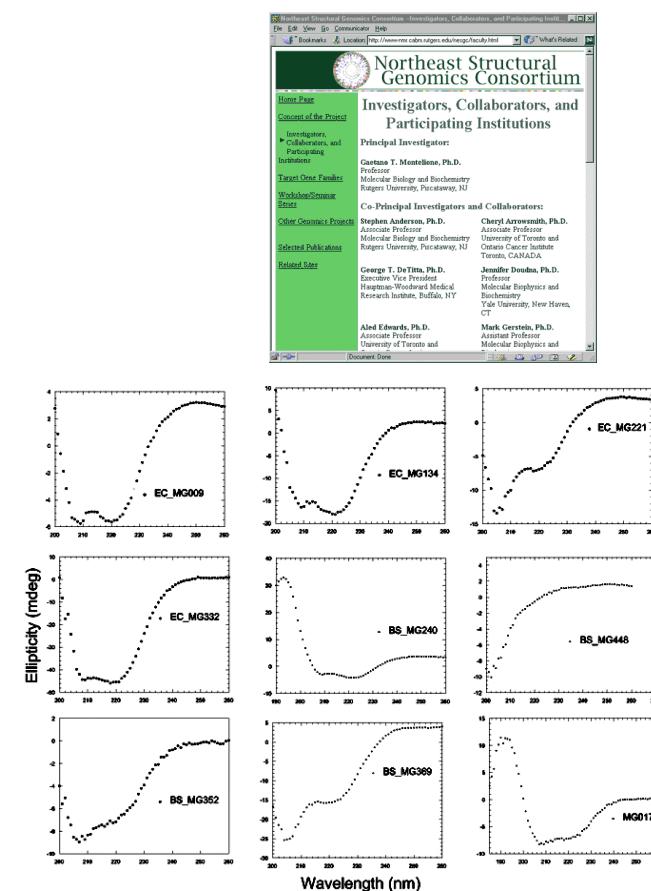


Finding Unusual Proteins for Expt. Structural Genomics

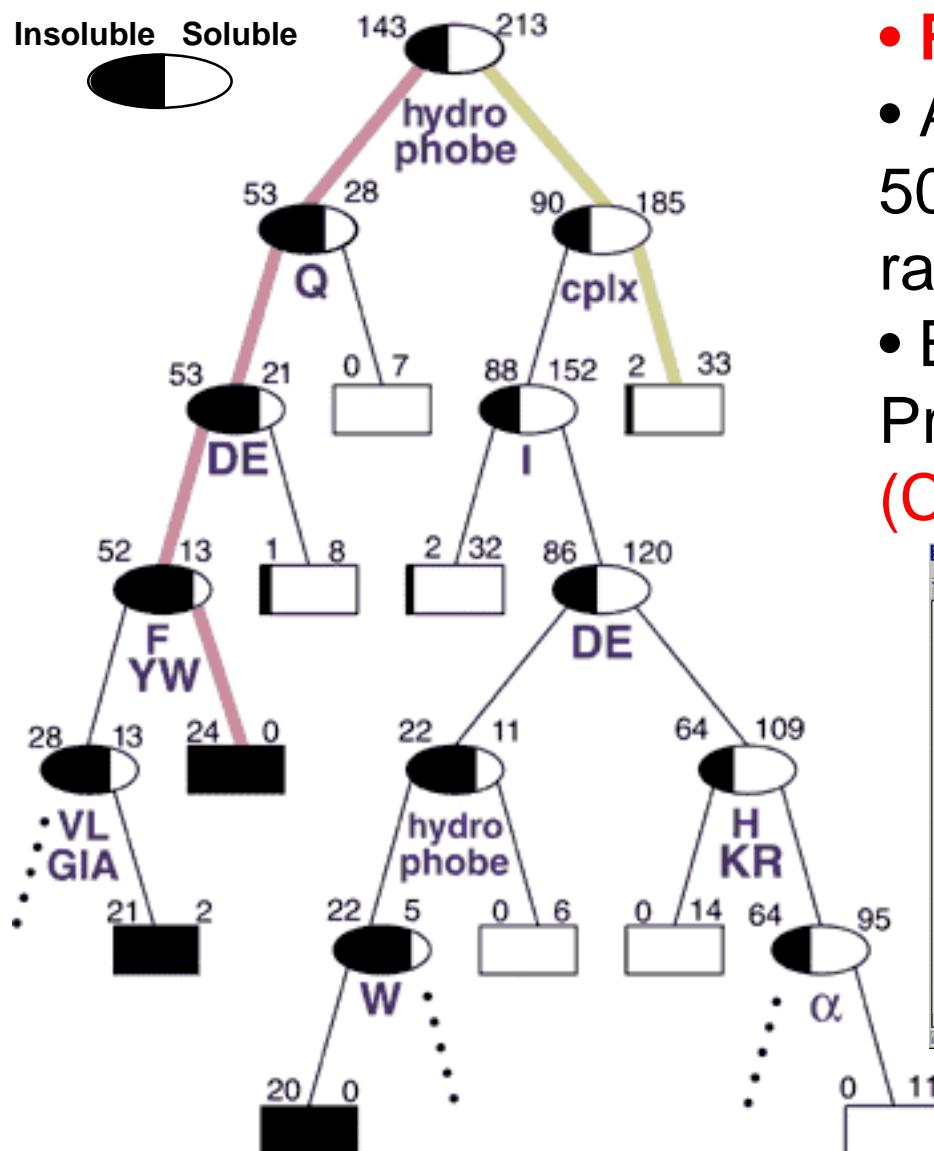


In collaboration with L Regan, S Balasubramanian

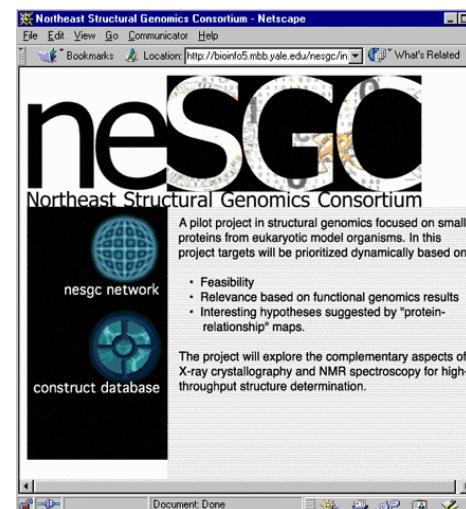
- Prospective Target Selection
- Identify Proteins in *M. genitalium* that are most atypical structurally (hardest)
- Characterize biophysically by CD (do they fold normally?)



Characterizing the Low-hanging Fruit for Experimental Structural Genomics



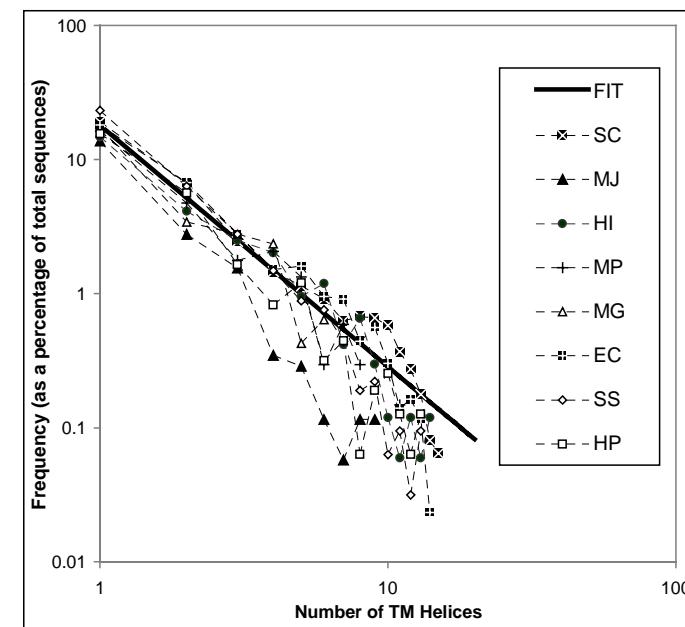
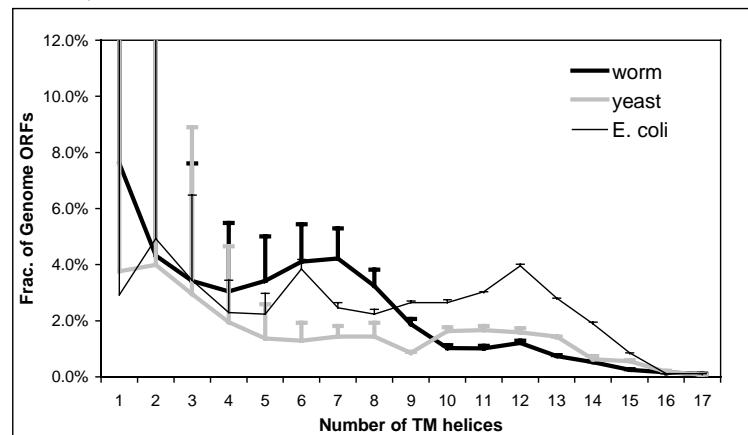
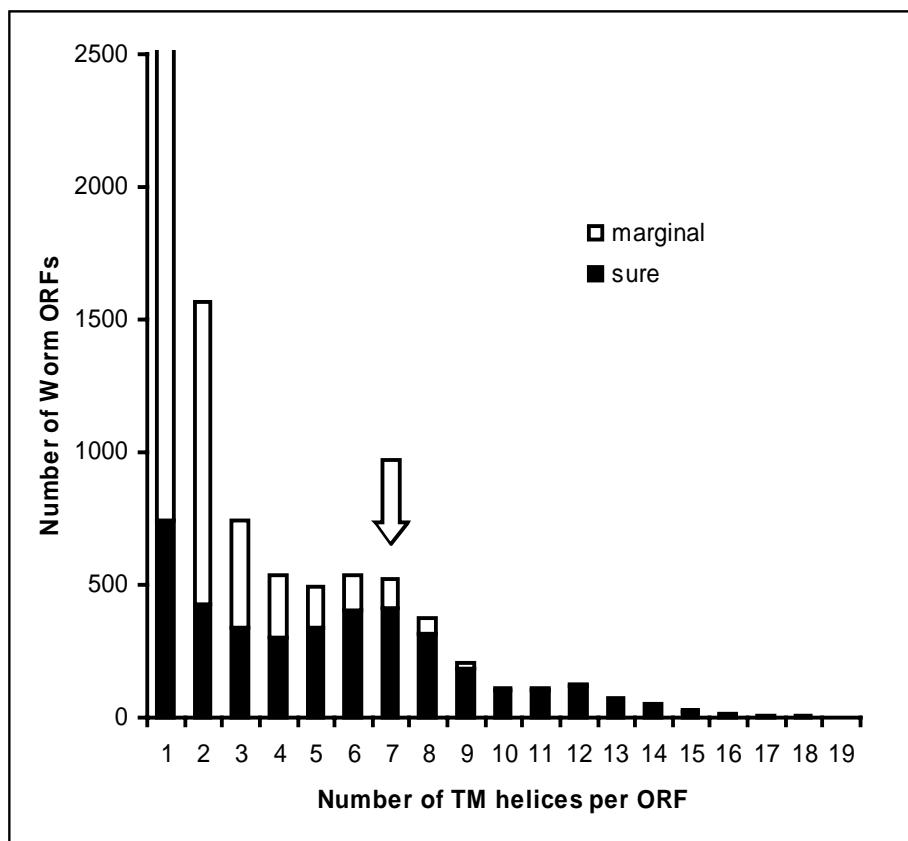
- **Retrospective Decision-Tree**
 - Analysis of the Suitability of 500 M. thermo. proteins for X-ray/NMR work
 - Based on results of Toronto Proteomics Group
(C Arrowsmith, A Edwards)



For example, proteins that fulfill the following sequence of four rules are likely to be insoluble: (1) have a hydrophobic stretch -- a long region (>20 residues) with average hydrophobicity less than -0.85 kcal/mole (on the GES scale); (2) Gln composition <4%; (3) Asp+Glu composition <17%; and (4) aromatic composition >7.5%. Conversely, proteins that do not have a hydrophobic stretch and have less than 27% of their residues in "low-complexity" regions are very likely to be soluble.

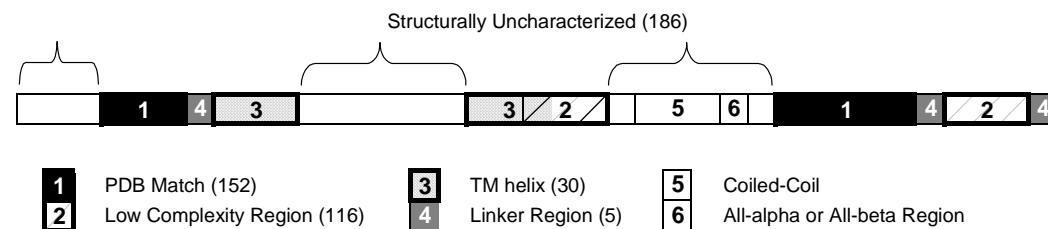
TM-helix “prediction”

- 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)



2^o Structure Prediction

- Bulk prediction of 2^o 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

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.....

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

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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

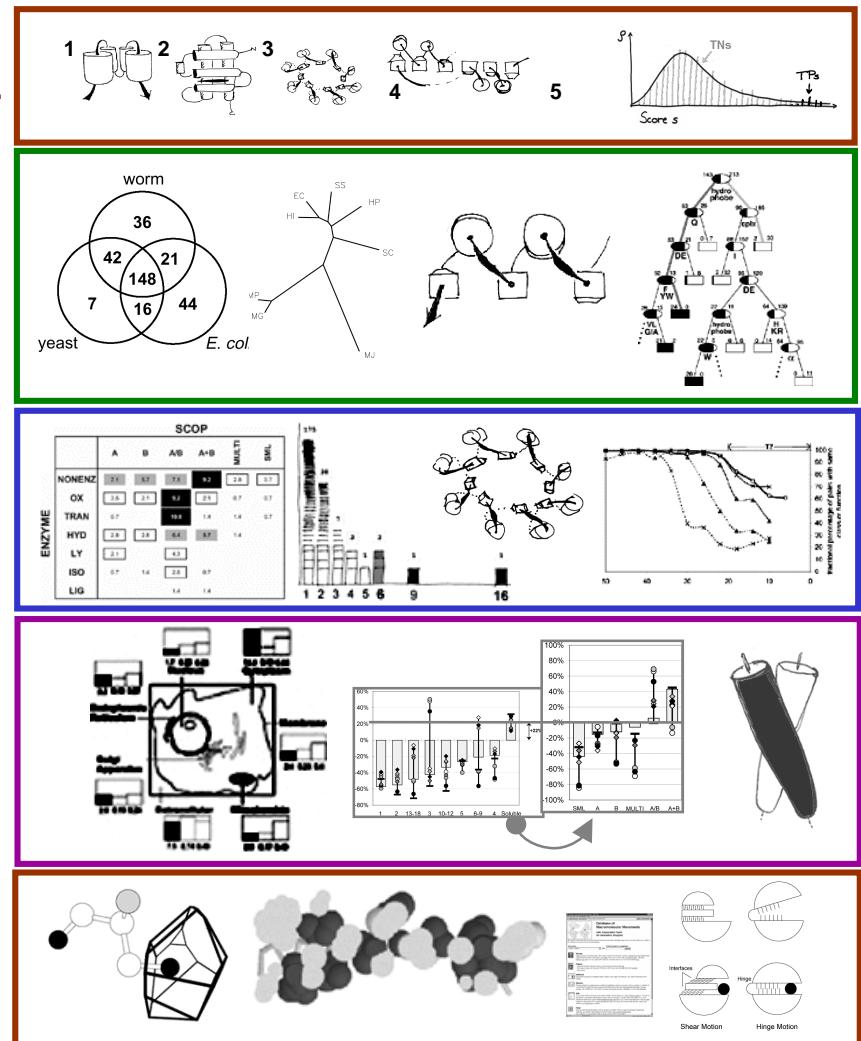
4 Folds in the Transcriptome

(Common parts? Where are parts?)

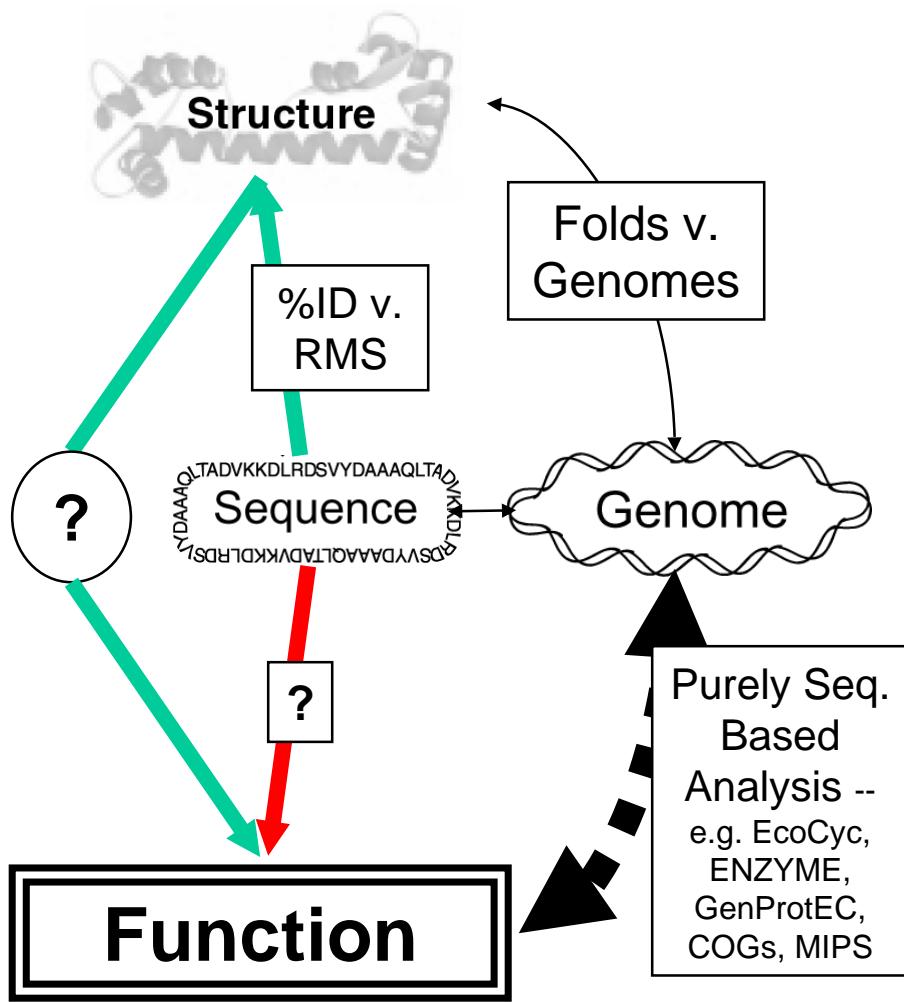
Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis, cyt. Depleted ↓ : NS, long, TM folds, transport, transcription, Leu-zip, nuc. Bayesian Localizer, phenotypes clustering

5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson, R Das, H Hegyi, J Lin, Y Kluger, C Arrowsmith, A Edwards, L Regan, S Balasubramanian, A Drawid, D Greenbaum, M Snyder, R Jansen

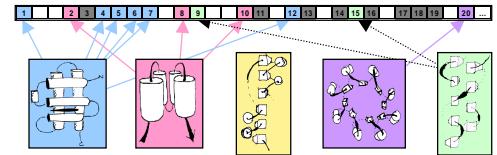


Adding Structure to Functional Genomics, Function to Structural Genomics

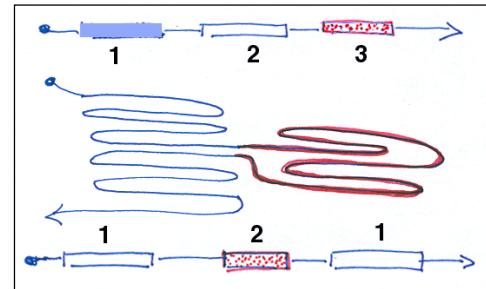


Why Structure? Do we really need it?

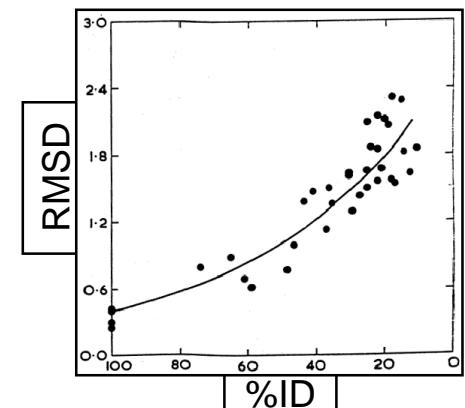
1 Most Highly Conserved



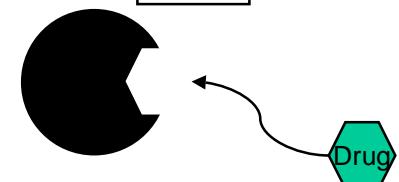
2 Precisely Defined Modules



3 Seq. \leftrightarrow Struc.
Clearer than Seq. \leftrightarrow Func.



4 Link to Chemistry, Drugs



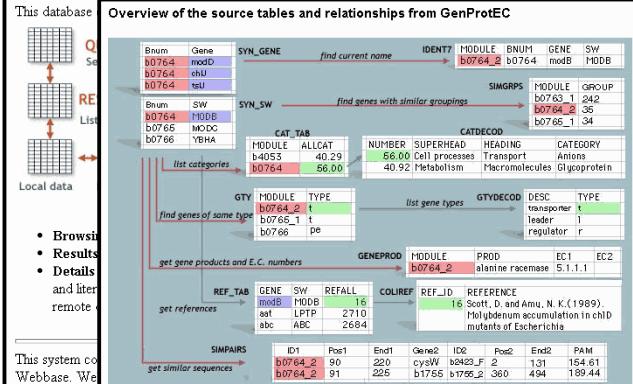
Functional Classification

GenProtEC E. coli genome and proteome database

Search Main
Click here to start searching

[Back to mainpage](#)

Schema of GenProtEC



This system contains Webbase. We have the tables in template files and then tables and them

Click a table to dump the first 10 rows.

- The primary key is SYN_GENE and identical module numbers.
- IDENT? contains folded into the SYN_GENE.
- SIMGRPS, GENES, CAT_TAB and GTY_TAB.
- GTY_TAB is a link by GENE and SW.
- REF_TAB is a link by GENE and SW.

3. GenProtEC (E. coli only)

<http://genprotec.mbl.edu/start>

Welcome to MIPS - Netscape

File Edit View Go Communicator Help

Bookmarks Location: <http://www.mips.biochem.mpg.de/> What's Related

mips

munich information center for protein sequences

MIPS Saccharomyces cerevisiae - Functional Categories - Netscape

Yeast Arabidopsis Human PIR-International Protam Pedant Sequence DB Access About MIPS Publications

METABOLISM (1046 ORFs)

- amino-acid metabolism (205 ORFs)
 - amino-acid biosynthesis (119 ORFs)
 - regulation of amino-acid metabolism (33 ORFs)
 - amino-acid transport (23 ORFs)
 - amino-acid degradation (catabolism) (35 ORFs)
 - other amino-acid metabolism activities (5 ORFs)
 - nitrogen and sulphur metabolism (75 ORFs)
 - nitrogen and sulphur utilization (38 ORFs)
 - regulation of nitrogen and sulphur utilization (29 ORFs)
 - nitrogen and sulphur transport (8 ORFs)
 - nucleotide metabolism (740 ORFs)
 - purine-nucleotide metabolism (45 ORFs)
 - pyrimidine-nucleotide metabolism (22 ORFs)
 - deoxyribonucleotide metabolism (12 ORFs)
 - metabolism of cyclic and unusual nucleotides (8 ORFs)
 - regulation of nucleotide metabolism (13 ORFs)
 - poly nucleotide degradation (23 ORFs)
 - nucleotide transport (13 ORFs)
 - other nucleotide metabolism activities (7 ORFs)
 - phosphate metabolism (31 ORFs)
 - phosphate utilization (13 ORFs)

2. MIPS/PEDANT (yeast only)

MIPS Saccharomyces cerevisiae - Functional Categories - Netscape

Yeast Arabidopsis Human PIR-International Protam Pedant Sequence DB Access About MIPS Publications

LIPID METABOLISM (1046 ORFs)

- lipid and fatty-acid binding (8 ORFs)
 - other lipid, fatty-acid and isoprenoid metabolism activities (13 ORFs)
- metabolism of vitamins, cofactors, and prosthetic groups (79 ORFs)
 - biosynthesis of vitamins, cofactors, and prosthetic groups (59 ORFs)
 - utilization of vitamins, cofactors, and prosthetic groups (7 ORFs)
 - regulation of vitamins, cofactors, and prosthetic groups (3 ORFs)
 - transport of vitamins, cofactors, and prosthetic groups (3 ORFs)
 - other vitamin, cofactor, and prosthetic group activities (7 ORFs)
 - secondary metabolism (4 ORFs)
 - biosynthesis of secondary products derived from primary amino acids (4 ORFs)
- ENERGY (246 ORFs)
 - gluconeogenesis and glycogenolysis (65 ORFs)
 - pentose-phosphate pathway (20 ORFs)
 - tricarboxylic acid pathway (24 ORFs)
 - respiration (65 ORFs)
 - fermentation (34 ORFs)
 - metabolism of energy reserves (glycogen, trehalose) (37 ORFs)

Gene Ontology

function cellular_component process

function

- 1.1. nucleic acid binding(724)
 - 1.1.1. DNA binding(350)
 - 1.1.1.1. ATP dependent DNA helicase(19)
 - 1.1.1.2. mitochondrial DNA helicase(1)
 - 1.1.1.3. AT DNA binding(2)
 - 1.1.1.4. bent DNA binding
 - 1.1.1.5. chromatin binding(24)
 - 1.1.1.6. damaged DNA binding
 - 1.1.1.7. DNA repair protein(11)
 - 1.1.1.7.1. DNA repair enzyme

cellular_component

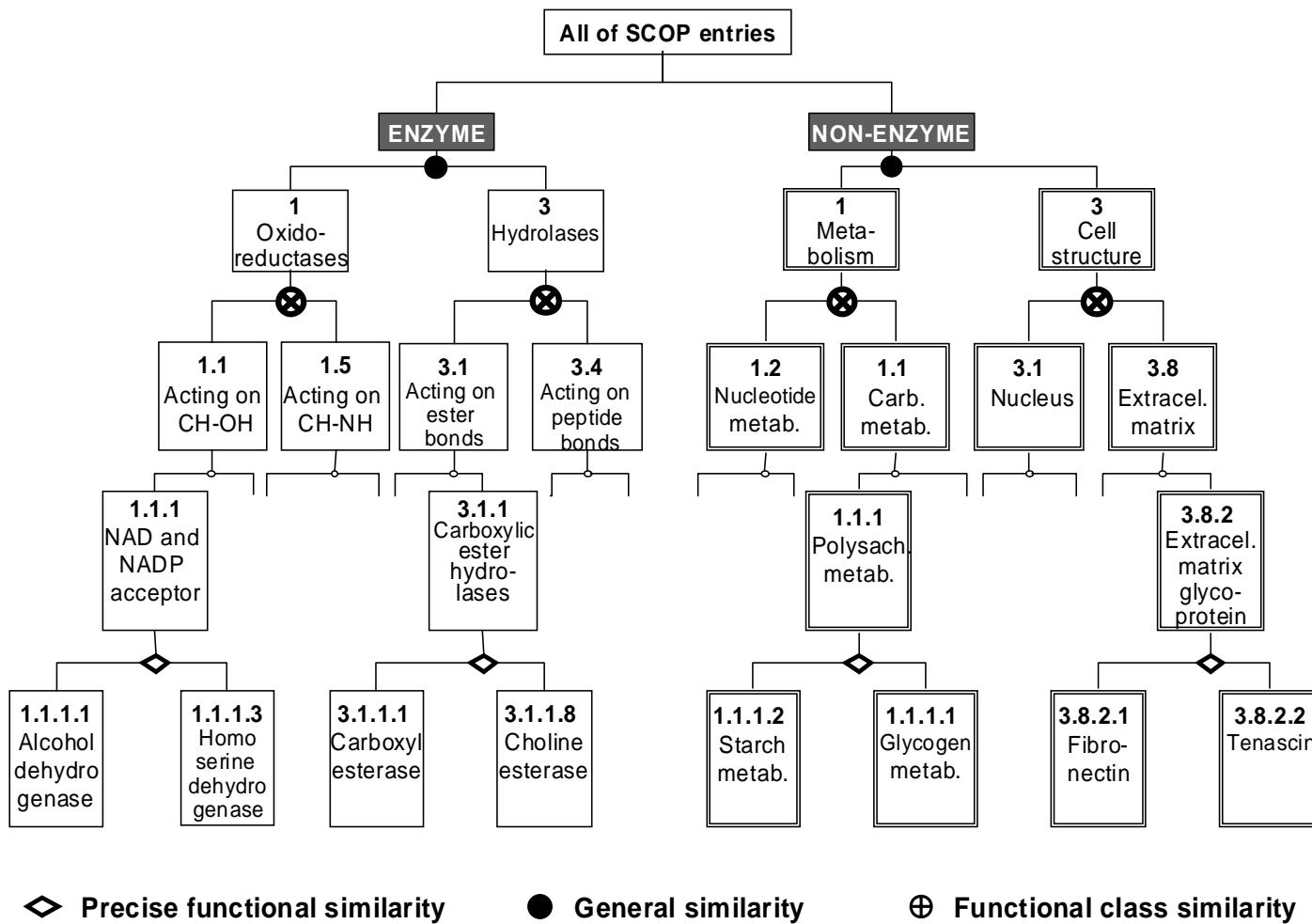
- 1.1. extracellular(71)
 - 1.1.1. fibrinogen(2)
 - 1.1.1.1. fibrinogen alpha chain
 - 1.1.1.2. fibrinogen gamma chain
 - 1.1.2. extracellular matrix(77)
 - 1.1.2.1. fibrinogen(2)
 - 1.1.2.1.1. membrane attack complex
 - 1.1.2.1.2. membrane attack complex
 - 1.1.2.2. collagen(4)
 - 1.1.2.2.1. collagen type XV
 - 1.1.2.2.2. fibrillar collagen(3)

4. FlyBase/Ashburner (fly only), extended to GO (cross-organism)

1. ENZYME (cross-organism, but just enzymes)

Also:
COGs
(cross organism, conserved)
WIT, KEGG (pathways)
TIGR, EGAD

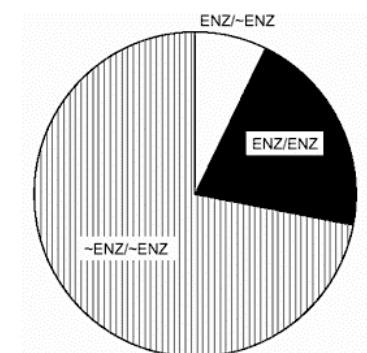
A Simple Scheme for Functionally Classifying Protein Structures



Focus on Pairs with Precise (1.1.1.*) and Broad (1.*) Similarity

A Combined Scheme, merging ENZYME + FLYBASE, with manual additions for proteins not in either (e.g. Ig's)

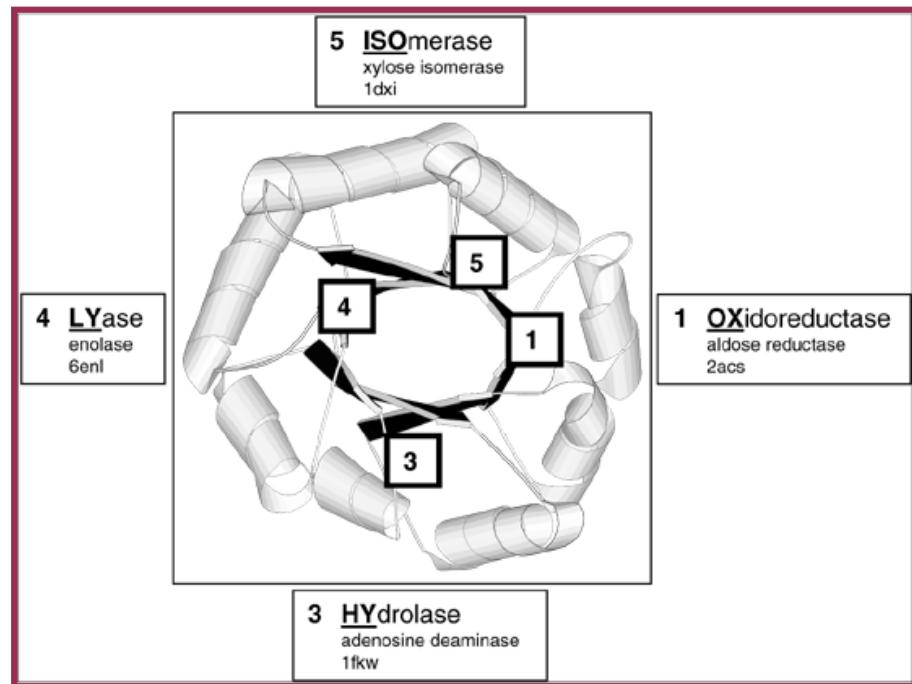
Also: MIPS, GenProtEC, GOGs



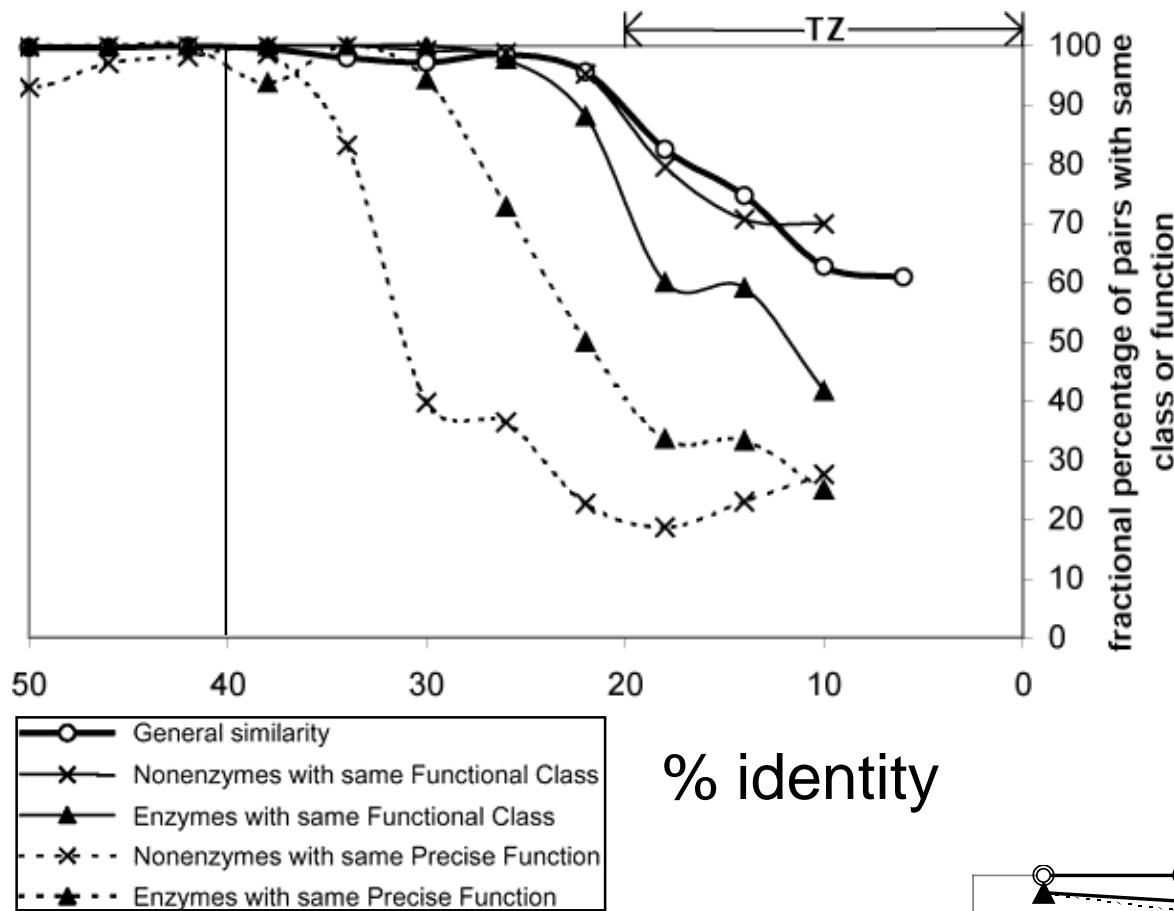
Fold-Function Combinations

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

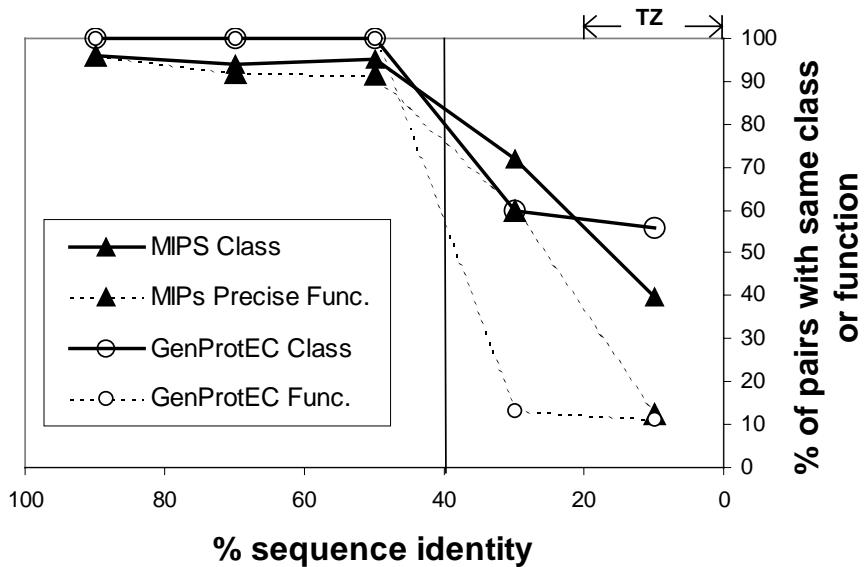
Have the sequences diverged beyond point of homology? Or does this occur at a certain sequence threshold?

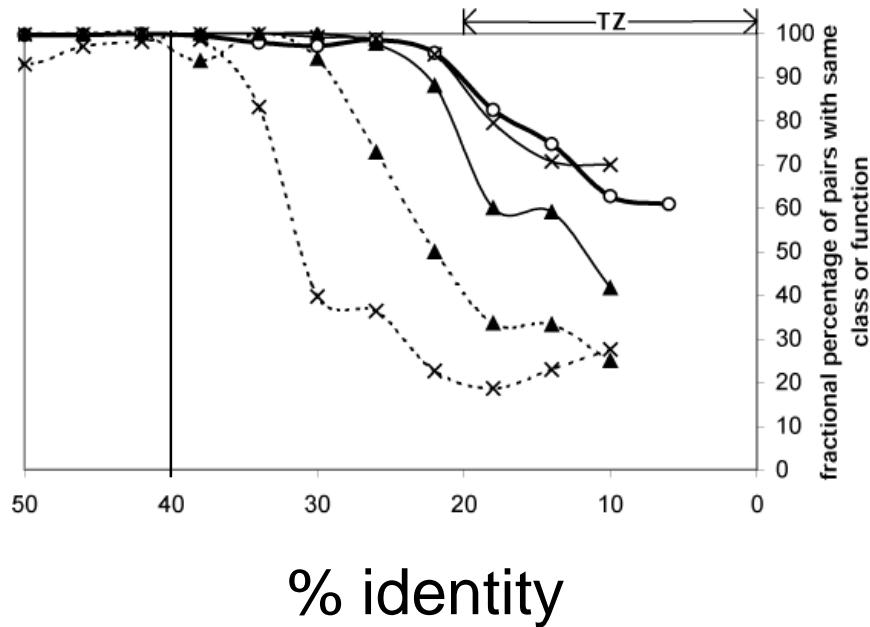


Relationship of Similarity in Sequence & Structure to that in Function



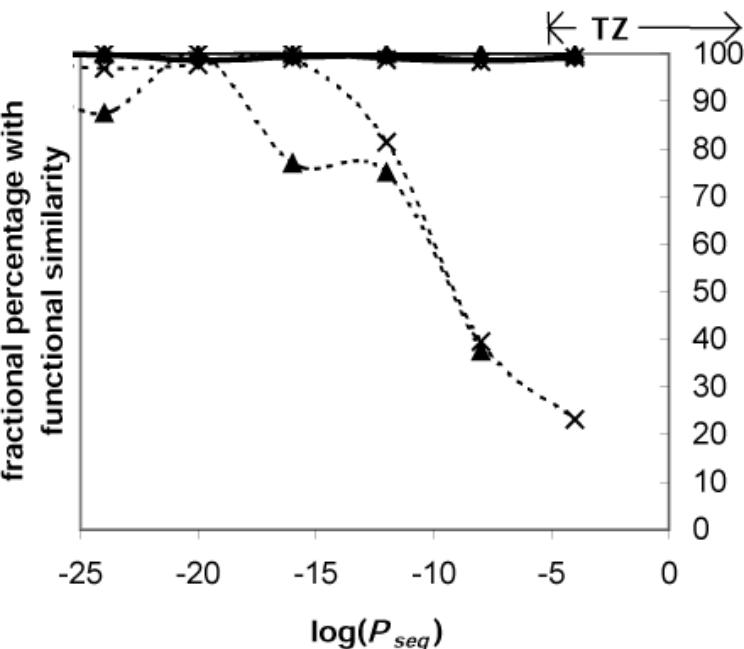
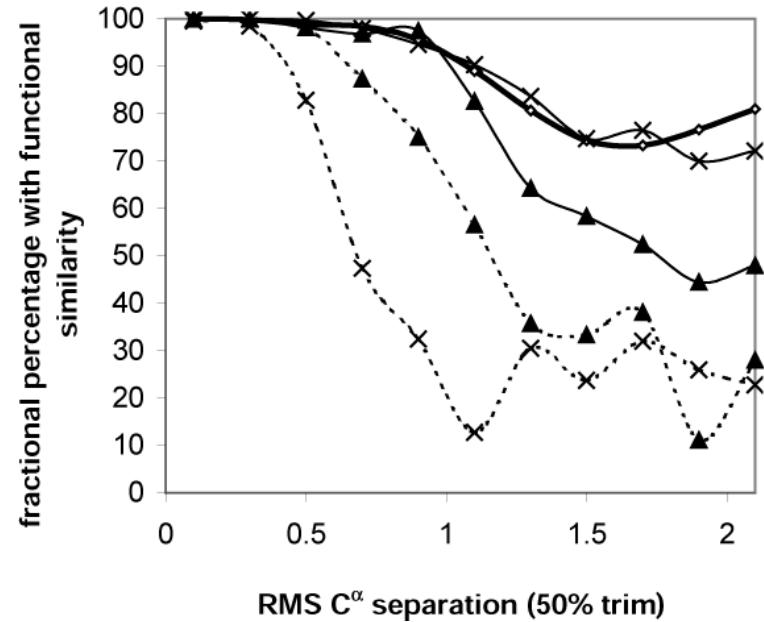
See at what %ID have diff. function (both broad & precise). Use 4 func. classifications -- ENZYME, FLYBASE (+extra), MIPS, GenProtEC





Relationship of Similarity in Sequence & Structure to that in Function II

Percent identity quite successful vs. structure sim. or statistical scores

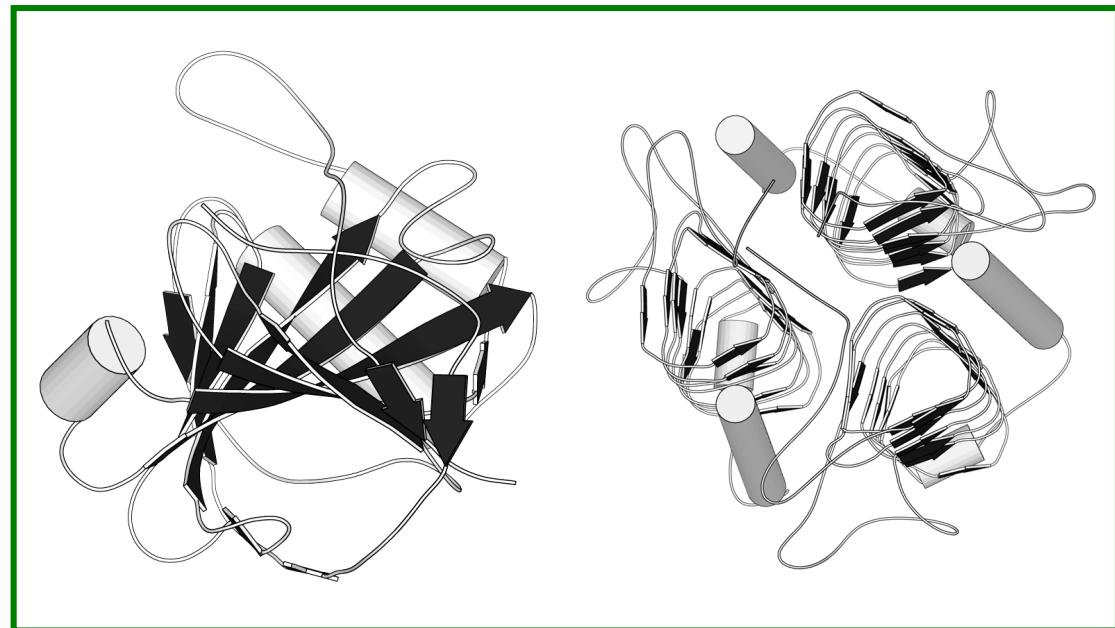
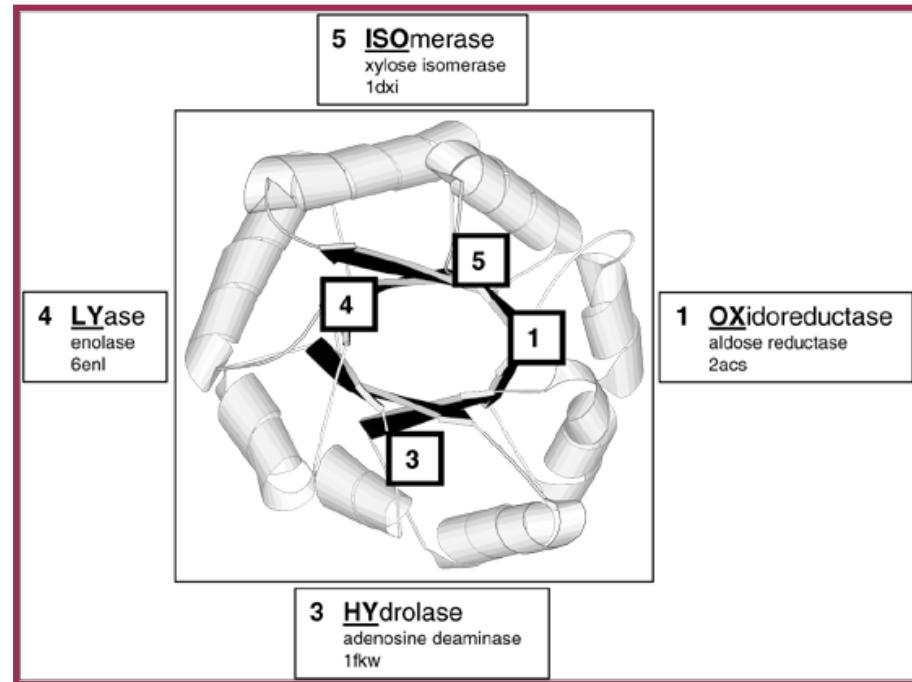


Fold-Function Combinations #2

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

Have the sequences diverged beyond point of homology? Or does this occur at a certain sequence threshold?

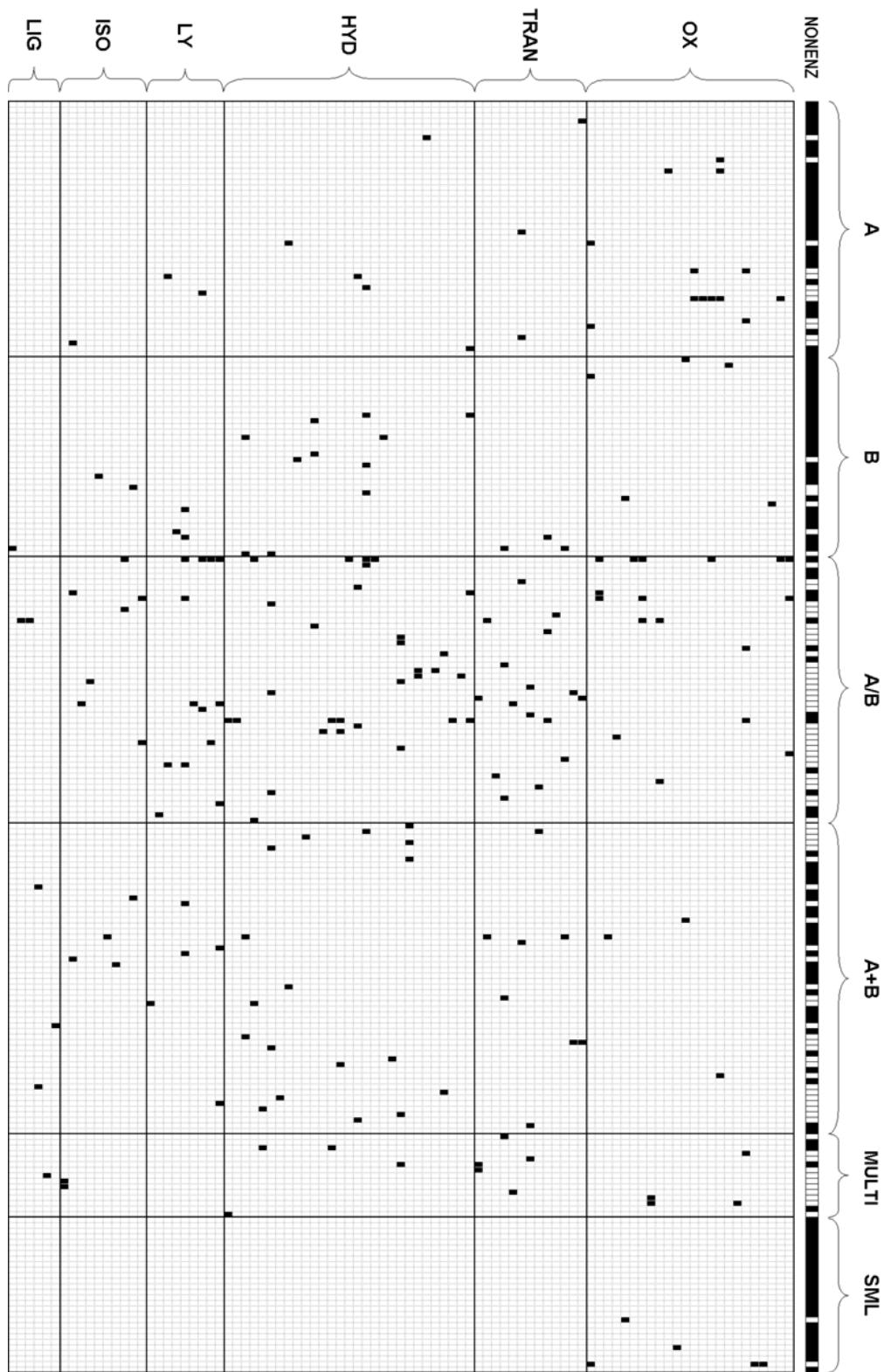
Two Different Folds Catalyze the Same Reaction -- e.g. Carbonic Anhydrases (4.2.1.1)



Fold-Function Combinations

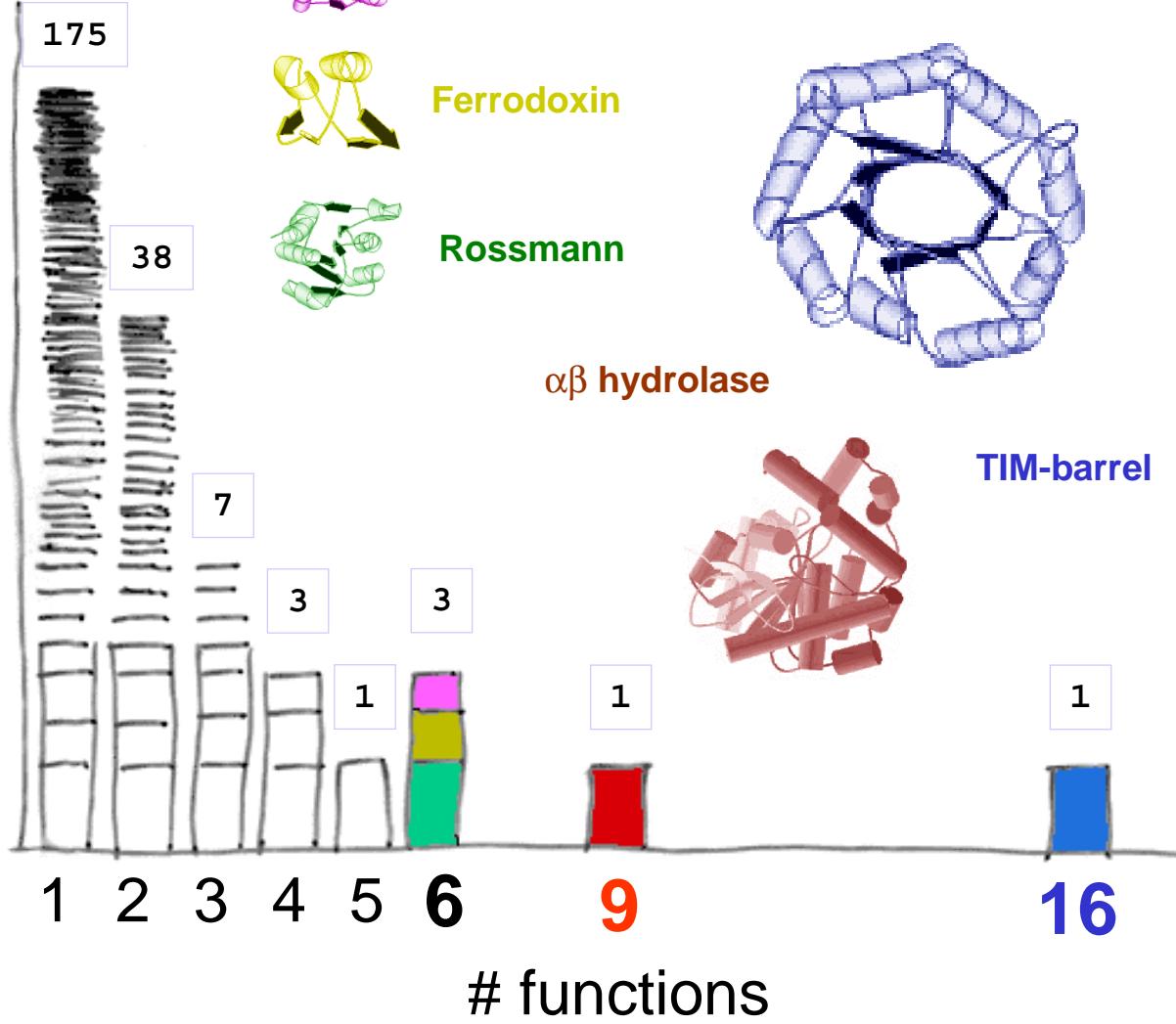
~20K (=92x229) Possible,
331 Observed

229 Folds



Most Versatile Folds & Functions

folds with certain # of functions



Top Multifold Functions →

	A	B	A+B	MULTI
NONENZ	1.1.1	2.1	3	4
OX	1.1.1.1	1.1.1.2	1.1.1.3	1.1.1.4
TRAN	2.1.1	2.1.2	2.1.3	2.1.4
HYD	3.1.1	3.1.2	3.1.3	3.1.4
LY	3.2.1	3.2.2	3.2.3	3.2.4
ISO	3.3.1	3.3.2	3.3.3	3.3.4
LIG	4.1.1	4.1.2	4.1.3	4.1.4

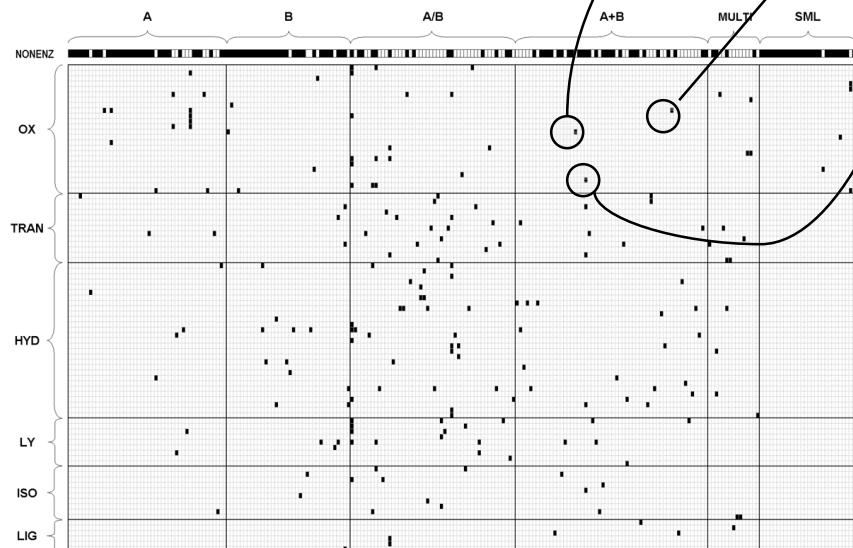
Top-4 Most Versatile Functions:

Glycosidases, carboxy-lyases, phosphoric monoester hydrolases, linear monoester hydrolases (3.2.1, 4.2.1 3.1.3, 3.5.1)

Top Multifunctional Folds →

	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
NONE	1.1.1	2.1	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
OX	1.1.1.1	1.1.1.2	1.1.1.3	1.1.1.4	1.1.1.5	1.1.1.6	1.1.1.7	1.1.1.8	1.1.1.9	1.1.1.10	1.1.1.11	1.1.1.12	1.1.1.13	1.1.1.14	1.1.1.15	1.1.1.16	1.1.1.17	1.1.1.18	1.1.1.19	1.1.1.20	1.1.1.21	1.1.1.22	1.1.1.23	1.1.1.24	1.1.1.25	1.1.1.26	1.1.1.27	1.1.1.28	1.1.1.29	1.1.1.30	1.1.1.31	1.1.1.32	1.1.1.33	1.1.1.34	1.1.1.35	1.1.1.36	1.1.1.37	1.1.1.38	1.1.1.39	1.1.1.40	1.1.1.41	1.1.1.42	1.1.1.43	1.1.1.44	1.1.1.45	1.1.1.46	1.1.1.47	1.1.1.48	1.1.1.49	1.1.1.50	1.1.1.51	1.1.1.52	1.1.1.53	1.1.1.54	1.1.1.55	1.1.1.56	1.1.1.57	1.1.1.58	1.1.1.59	1.1.1.60	1.1.1.61	1.1.1.62	1.1.1.63	1.1.1.64	1.1.1.65	1.1.1.66	1.1.1.67	1.1.1.68	1.1.1.69	1.1.1.70	1.1.1.71	1.1.1.72	1.1.1.73	1.1.1.74	1.1.1.75	1.1.1.76	1.1.1.77	1.1.1.78	1.1.1.79	1.1.1.80	1.1.1.81	1.1.1.82	1.1.1.83	1.1.1.84	1.1.1.85	1.1.1.86	1.1.1.87	1.1.1.88	1.1.1.89	1.1.1.90	1.1.1.91	1.1.1.92	1.1.1.93	1.1.1.94	1.1.1.95	1.1.1.96	1.1.1.97	1.1.1.98	1.1.1.99	1.1.1.100
TRAN	2.1.1.1	2.1.1.2	2.1.1.3	2.1.1.4	2.1.1.5	2.1.1.6	2.1.1.7	2.1.1.8	2.1.1.9	2.1.1.10	2.1.1.11	2.1.1.12	2.1.1.13	2.1.1.14	2.1.1.15	2.1.1.16	2.1.1.17	2.1.1.18	2.1.1.19	2.1.1.20	2.1.1.21	2.1.1.22	2.1.1.23	2.1.1.24	2.1.1.25	2.1.1.26	2.1.1.27	2.1.1.28	2.1.1.29	2.1.1.30	2.1.1.31	2.1.1.32	2.1.1.33	2.1.1.34	2.1.1.35	2.1.1.36	2.1.1.37	2.1.1.38	2.1.1.39	2.1.1.40	2.1.1.41	2.1.1.42	2.1.1.43	2.1.1.44	2.1.1.45	2.1.1.46	2.1.1.47	2.1.1.48	2.1.1.49	2.1.1.50	2.1.1.51	2.1.1.52	2.1.1.53	2.1.1.54	2.1.1.55	2.1.1.56	2.1.1.57	2.1.1.58	2.1.1.59	2.1.1.60	2.1.1.61	2.1.1.62	2.1.1.63	2.1.1.64	2.1.1.65	2.1.1.66	2.1.1.67	2.1.1.68	2.1.1.69	2.1.1.70	2.1.1.71	2.1.1.72	2.1.1.73	2.1.1.74	2.1.1.75	2.1.1.76	2.1.1.77	2.1.1.78	2.1.1.79	2.1.1.80	2.1.1.81	2.1.1.82	2.1.1.83	2.1.1.84	2.1.1.85	2.1.1.86	2.1.1.87	2.1.1.88	2.1.1.89	2.1.1.90	2.1.1.91	2.1.1.92	2.1.1.93	2.1.1.94	2.1.1.95	2.1.1.96	2.1.1.97	2.1.1.98	2.1.1.99	2.1.1.100
HYD	3.1.1.1	3.1.1.2	3.1.1.3	3.1.1.4	3.1.1.5	3.1.1.6	3.1.1.7	3.1.1.8	3.1.1.9	3.1.1.10	3.1.1.11	3.1.1.12	3.1.1.13	3.1.1.14	3.1.1.15	3.1.1.16	3.1.1.17	3.1.1.18	3.1.1.19	3.1.1.20	3.1.1.21	3.1.1.22	3.1.1.23	3.1.1.24	3.1.1.25	3.1.1.26	3.1.1.27	3.1.1.28	3.1.1.29	3.1.1.30	3.1.1.31	3.1.1.32	3.1.1.33	3.1.1.34	3.1.1.35	3.1.1.36	3.1.1.37	3.1.1.38	3.1.1.39	3.1.1.40	3.1.1.41	3.1.1.42	3.1.1.43	3.1.1.44	3.1.1.45	3.1.1.46	3.1.1.47	3.1.1.48	3.1.1.49	3.1.1.50	3.1.1.51	3.1.1.52	3.1.1.53	3.1.1.54	3.1.1.55	3.1.1.56	3.1.1.57	3.1.1.58	3.1.1.59	3.1.1.60	3.1.1.61	3.1.1.62	3.1.1.63	3.1.1.64	3.1.1.65	3.1.1.66	3.1.1.67	3.1.1.68	3.1.1.69	3.1.1.70	3.1.1.71	3.1.1.72	3.1.1.73	3.1.1.74	3.1.1.75	3.1.1.76	3.1.1.77	3.1.1.78	3.1.1.79	3.1.1.80	3.1.1.81	3.1.1.82	3.1.1.83	3.1.1.84	3.1.1.85	3.1.1.86	3.1.1.87	3.1.1.88	3.1.1.89	3.1.1.90	3.1.1.91	3.1.1.92	3.1.1.93	3.1.1.94	3.1.1.95	3.1.1.96	3.1.1.97	3.1.1.98	3.1.1.99	3.1.1.100
LY	3.2.1.1	3.2.1.2	3.2.1.3	3.2.1.4	3.2.1.5	3.2.1.6	3.2.1.7	3.2.1.8	3.2.1.9	3.2.1.10	3.2.1.11	3.2.1.12	3.2.1.13	3.2.1.14	3.2.1.15	3.2.1.16	3.2.1.17	3.2.1.18	3.2.1.19	3.2.1.20	3.2.1.21	3.2.1.22	3.2.1.23	3.2.1.24	3.2.1.25	3.2.1.26	3.2.1.27	3.2.1.28	3.2.1.29	3.2.1.30	3.2.1.31	3.2.1.32	3.2.1.33	3.2.1.34	3.2.1.35	3.2.1.36	3.2.1.37	3.2.1.38	3.2.1.39	3.2.1.40	3.2.1.41	3.2.1.42	3.2.1.43	3.2.1.44	3.2.1.45	3.2.1.46	3.2.1.47	3.2.1.48	3.2.1.49	3.2.1.50	3.2.1.51	3.2.1.52	3.2.1.53	3.2.1.54	3.2.1.55	3.2.1.56	3.2.1.57	3.2.1.58	3.2.1.59	3.2.1.60	3.2.1.61	3.2.1.62	3.2.1.63	3.2.1.64	3.2.1.65	3.2.1.66	3.2.1.67	3.2.1.68	3.2.1.69	3.2.1.70	3.2.1.71	3.2.1.72	3.2.1.73	3.2.1.74	3.2.1.75	3.2.1.76	3.2.1.77	3.2.1.78	3.2.1.79	3.2.1.80	3.2.1.81	3.2.1.82	3.2.1.83	3.2.1.84	3.2.1.85	3.2.1.86	3.2.1.87	3.2.1.88	3.2.1.89	3.2.1.90	3.2.1.91	3.2.1.92	3.2.1.93	3.2.1.94	3.2.1.95	3.2.1.96	3.2.1.97	3.2.1.98	3.2.1.99	3.2.1.100
ISO	3.3.1.1	3.3.1.2	3.3.1.3	3.3.1.4	3.3.1.5	3.3.1.6	3.3.1.7	3.3.1.8	3.3.1.9	3.3.1.10	3.3.1.11	3.3.1.12	3.3.1.13	3.3.1.14	3.3.1.15	3.3.1.16	3.3.1.17	3.3.1.18	3.3.1.19	3.3.1.20	3.3.1.21	3.3.1.22	3.3.1.23	3.3.1.24	3.3.1.25	3.3.1.26	3.3.1.27	3.3.1.28	3.3.1.29	3.3.1.30	3.3.1.31	3.3.1.32	3.3.1.33	3.3.1.34	3.3.1.35	3.3.1.36	3.3.1.37	3.3.1.38	3.3.1.39	3.3.1.40	3.3.1.41	3.3.1.42	3.3.1.43	3.3.1.44	3.3.1.45	3.3.1.46	3.3.1.47	3.3.1.48	3.3.1.49	3.3.1.50	3.3.1.51	3.3.1.52	3.3.1.53	3.3.1.54	3.3.1.55	3.3.1.56	3.3.1.57	3.3.1.58	3.3.1.59	3.3.1.60	3.3.1.61	3.3.1.62	3.3.1.63	3.3.1.64	3.3.1.65	3.3.1.66	3.3.1.67	3.3.1.68	3.3.1.69	3.3.1.70	3.3.1.71	3.3.1.72	3.3.1.73	3.3.1.74	3.3.1.75	3.3.1.76	3.3.1.77	3.3.1.78	3.3.1.79	3.3.1.80	3.3.1.81	3.3.1.82	3.3.1.83	3.3.1.84	3.3.1.85	3.3.1.86	3.3.1.87	3.3.1.88	3.3.1.89	3.3.1.90	3.3.1.91	3.3.1.92	3.3.1.93	3.3.1.94	3.3.1.95	3.3.1.96	3.3.1.97	3.3.1.98	3.3.1.99	3.3.1.100
LIG	3.4.1.1	3.4.1.2	3.4.1.3	3.4.1.4	3.4.1.5	3.4.1.6	3.4.1.7	3.4.1.8	3.4.1.9	3.4.1.10	3.4.1.11	3.4.1.12	3.4.1.13	3.4.1.14	3.4.1.15	3.4.1.16	3.4.1.17	3.4.1.18	3.4.1.19	3.4.1.20	3.4.1.21	3.4.1.22	3.4.1.23	3.4.1.24	3.4.1.25	3.4.1.26	3.4.1.27	3.4.1.28	3.4.1.29	3.4.1.30	3.4.1.31	3.4.1.32	3.4.1.33	3.4.1.34	3.4.1.35	3.4.1.36	3.4.1.37	3.4.1.38	3.4.1.39	3.4.1.40	3.4.1.41	3.4.1.42	3.4.1.43	3.4.1.44	3.4.1.45	3.4.1.46	3.4.1.47	3.4.1.48	3.4.1.49	3.4.1.50	3.4.1.51	3.4.1.52	3.4.1.53	3.4.1.54	3.4.1.55	3.4.1.56	3.4.1.57	3.4.1.58	3.4.1.59	3.4.1.60	3.4.1.61	3.4.1.62	3.4.1.63	3.4.1.64	3.4.1.65	3.4.1.66	3.4.1.67	3.4.1.68	3.4.1.69	3.4.1.70	3.4.1.71	3.4.1.72	3.4.1.73	3.4.1.74	3.4.1.75	3.4.1.76	3.4.1.77	3.4.1.78	3.4.1.79	3.4.1.80	3.4.1.81	3.4.1.82	3.4.1.83	3.4.1.84	3.4.1.85	3.4.1.86	3.4.1.87	3.4.1.88	3.4.1.89	3.4.1.90	3.4.1.91	3.4.1.92	3.4.1.93	3.4.1.94	3.4.1.95	3.4.1.96	3.4.1.97	3.4.1.98	3.4.1.99	3.4.1.100

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	9	5		35
HYD	4	11	30	18			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

3

	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
TRAN	0.7			10.6	1.4	1.4
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]

Compare Classifications and Genomes

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

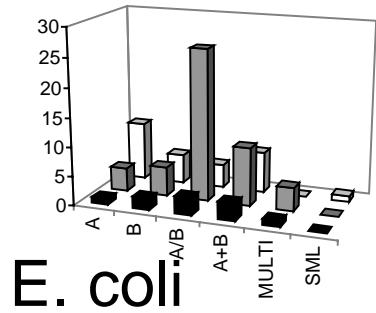
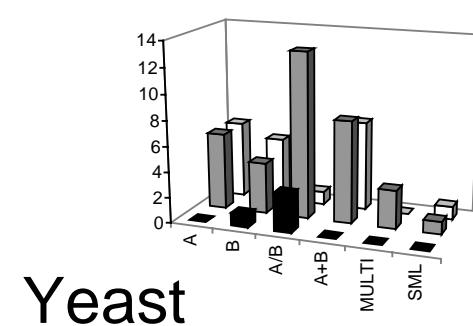
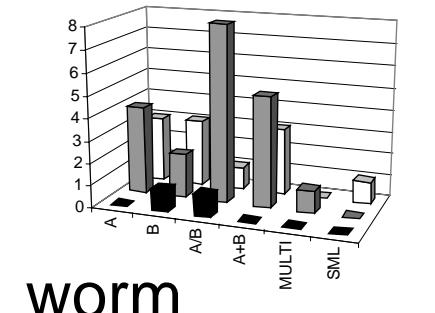
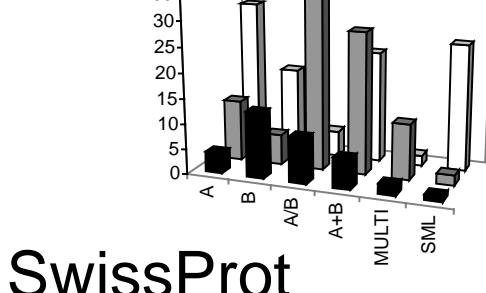
		SCOP					
		A	B	A/B	A+B	MULTI	SML
ENZYME	NONENZ	7.1	5.7	7.1	9.2	2.8	0.7
	OX	3.5	2.1	9.2	2.1	0.7	0.7
	TRAN	0.7		10.6	1.4	1.4	0.7
	HYD	2.8	2.8	6.4	5.7	1.4	
	LY	2.1		4.3			
	ISO	0.7	1.4	2.8	0.7		
	LIG			1.4	1.4		

CATH (Thornton)

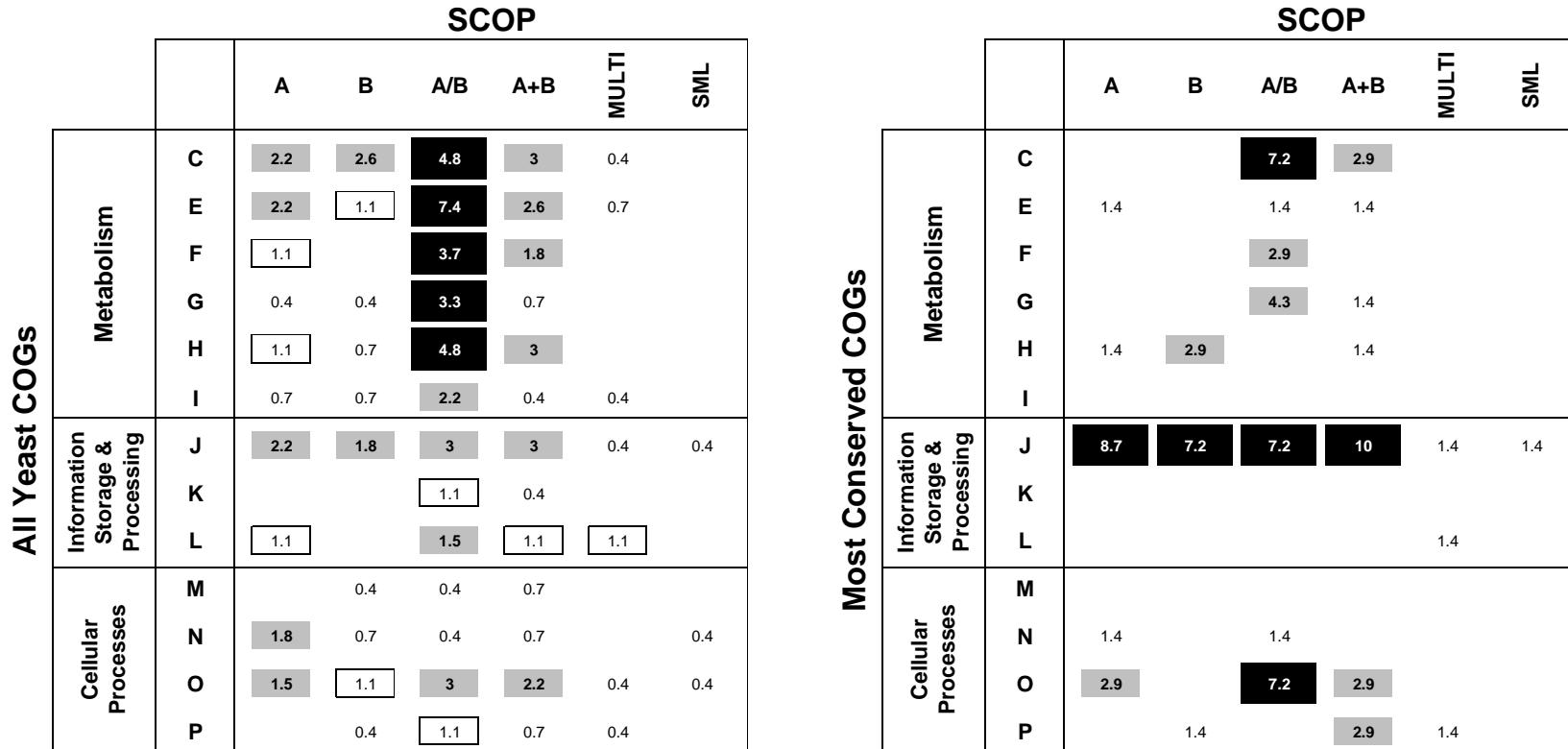
		CATH		
		A	B	AB
ENZYME	NONENZ	10	9.0	15
	OX	5.1	5.1	10
	TRAN		1.3	13
	HYD	2.6	1.3	14
	LY		2.6	1.3
	ISO	1.3	1.3	5.1
	LIG			1.3

MIPS YFC (Mewes)

		SCOP					
		A	B	A/B	A+B	MULTI	SML
MIPS Functional Cat.	metabolism	3.5	2.3	10	4.5	1.3	0.8
	energy	1.1	1.2	5	1.5	0.3	0.2
	growth, div., DNA syn.	4.9	3.6	4	4.5	1.8	1.2
	transcription	1.5	1.3	2.2	1.5	0.5	0.8
	protein synthesis	1	0.9	0.7	1.3	0.3	0.2
	protein targeting	1.2	1.7	2	1.6	0.5	0.3
	transport facilitation	0.9	0.5	0.7	0.6	0.4	
	intracellular transport	1.8	2.1	1.6	0.6	1	
	cellular biogenesis	0.9	0.7	1.2	0.3	0.3	0.1
	signal transduction	1	1	1.1	0.3	0.7	0.3
	cell rescue, defense...	1.5	1	2.6	1.9	0.7	0.5
	ionic homeostasis	0.5	0.3	0.4	0.4	0.2	



Different Structure Function Relationships for Most Ancient Proteins



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

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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

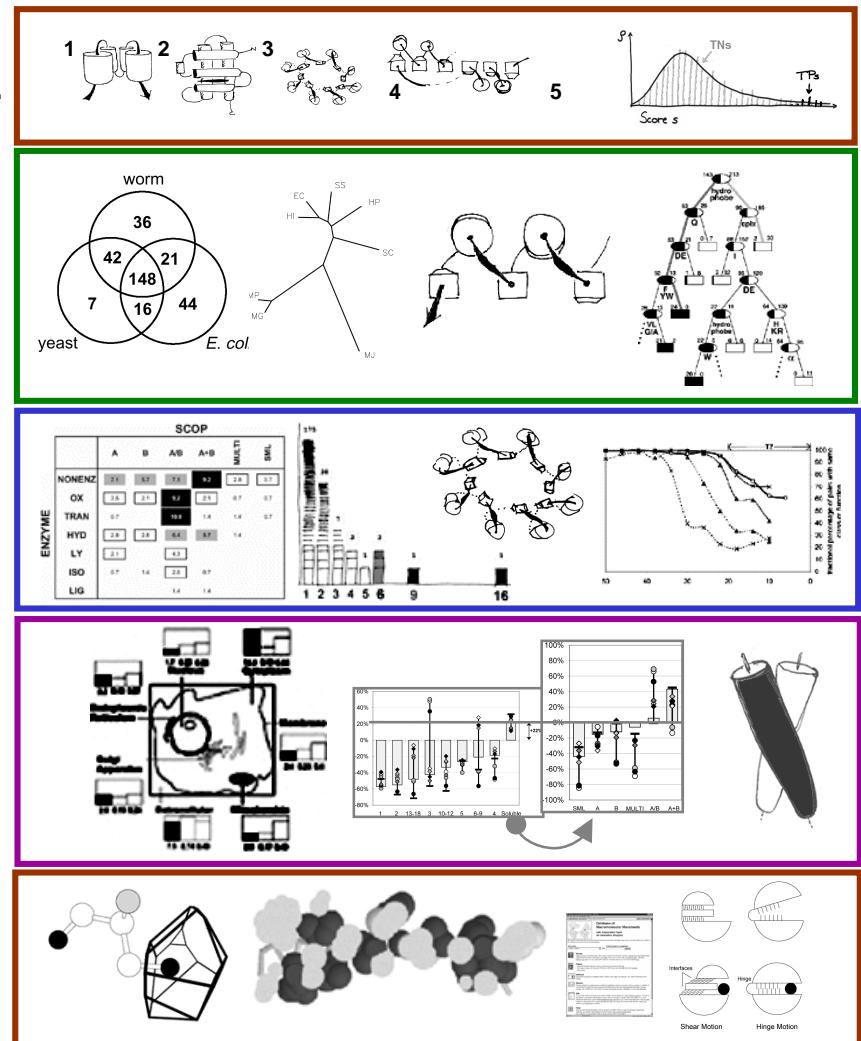
4 Folds in the Transcriptome

(Common parts? Where are parts?)

Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis, cyt. Depleted ↓ : NS, long, TM folds, transport, transcription, Leu-zip, nuc. Bayesian Localizer, phenotypes clustering

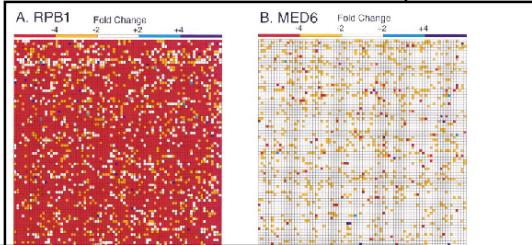
5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson, R Das, H Hegyi, J Lin, Y Kluger, C Arrowsmith, A Edwards, L Regan, S Balasubramanian, A Drawid, D Greenbaum, M Snyder, R Jansen



Dissecting the Regulatory Circuitry of a Eukaryotic Genome

Frank C. P. Holstege,* Ezra G. Jennings,*[†]
 John J. Wyrick,*[†] Tong Ihn Lee,*[†]
 Christoph J. Hengartner,*[†] Michael R. Green,[‡]
 Todd R. Golub,*[§] Eric S. Lander,*[†]
 and Richard A. Young^{†,||}
 *Whitehead Institute for Biomedical Research
 Cambridge, Massachusetts 02142
[†]Department of Biology
 Massachusetts Institute of Technology
 Cambridge, Massachusetts 02139
[‡]Howard Hughes Medical Institute
 Program in Molecular Medicine
 University of Massachusetts Medical Center



Young/Lander Affymetrix GeneChips Abs. Exp.

regulation which is superimposed on that due to specific transcription factors, a novel mechanism that may explain the specific rate of gene expression.

<http://brown.stanford.edu/brown/>

Figure 2. Genome-Wide Expression Data for Selected Components of the RNA Polymerase II Holoenzyme. Change in mRNA levels when a mutant is compared to its isogenic wild-type counterpart is presented in a grid format. In each grid square represents the left-most gene on chromosome I, and the squares to its right represent adjacent genes, in fashion through chromosome I, then II, then III, etc., until the last gene on the right arm of chromosome XVI is reached grid. The results are shown for (A) Rpb1, (B) Med6, (C) Srb15, and (D) Srb2.

The Brown Lab

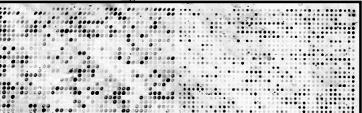
Stanford University Department of Biochemistry

The MGuide

The Complete Guide to MicroArrays
Build your own arrayer and scanner!

The transcriptional program in the response of human fibroblasts to serum

The web supplement to Iyer V.R. et al. (1999) Science 283:63-67



Brown, μarrays, Rel. Exp. over Timecourse

Also: SAGE;
Samson and
Church, Chips;
Aebersold,
Protein
Abundance

Gene Expression Datasets: the Transcriptome

Yeast Expression Data in Academia:
levels for all 6000 genes!

X-ref. with other genome data: protein fold
features common in Transcriptome....

Proc. Natl. Acad. Sci. USA
Vol. 94, pp. 190-195, January 1997
Genetics

A multipurpose transposon system for analyzing protein production, localization, and function in *Saccharomyces cerevisiae*

PETRA ROSS-MACDONALD, AMY SHEEHAN, G. SHIRLEEN ROEDER, and MICHAEL SNYDER*

Department of Biology, Yale University, P.O. Box 208103, New Haven,

Communicated by Gerald R. Fink, Whitehead Institute, Cambridge

ABSTRACT Analysis of the function of a particular product typically involves determining the expression pattern of the gene, the subcellular location of the protein, and phenotype of a null strain lacking the protein. Conditional alleles of the gene are often created as an additional tool to have developed a multifunctional, transposon-based system that simultaneously generates constructs for all the analyses and is suitable for mutagenesis of any given *Saccharomyces cerevisiae* gene. Depending on the transposon used, yeast gene is fused to a coding region for β-galactosidase or green fluorescent protein. Gene expression can therefore be monitored by chemical or fluorescence assays. The transposons create insertion mutations in the target gene, allowing phenotypic analysis. The transposon can be reduced by cre site-specific recombination to a smaller element that leaves the epitope tag inserted in the encoded protein. In addition to utility for a variety of immunodetection purposes, the ep

was mutagenizing the *E. coli* by shuttle mutagenesis. DNA containing the transposon was excised from the plasmid and inserted into yeast, where it replaced the chromosomal locus by homologous recombination. With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The transposons create insertion mutations in the target gene, allowing phenotypic analysis. The transposon can be reduced by cre site-specific recombination to a smaller element that leaves the epitope tag inserted in the encoded protein. In addition to utility for a variety of immunodetection purposes, the ep

was mutagenizing the *E. coli* by shuttle mutagenesis. DNA containing the transposon was excised from the plasmid and inserted into yeast, where it replaced the chromosomal locus by homologous recombination.

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

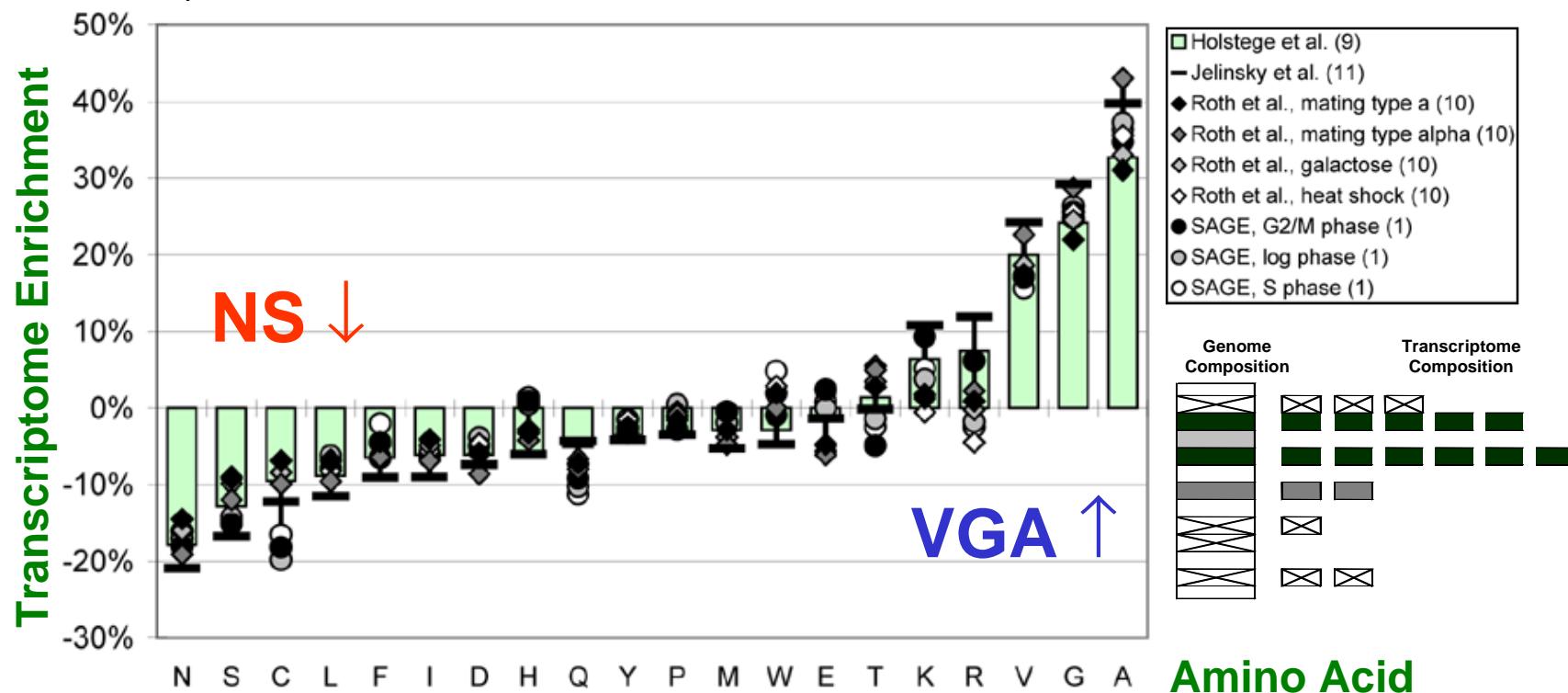
With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

With both mTrn-3HA/lacZ and mTrn-4HA/lacZ, transposons contain the *URA3* and *lacZ* genes for selection in *S. cerevisiae* and *E. coli*, respectively. The mTrn-3HA/GFP contains the coding region for GFP mutant p1(1). In each case, these are flanked by *lacZ* and *Tet* terminal repeats (TR).

Composition of Genome vs. Transcriptome

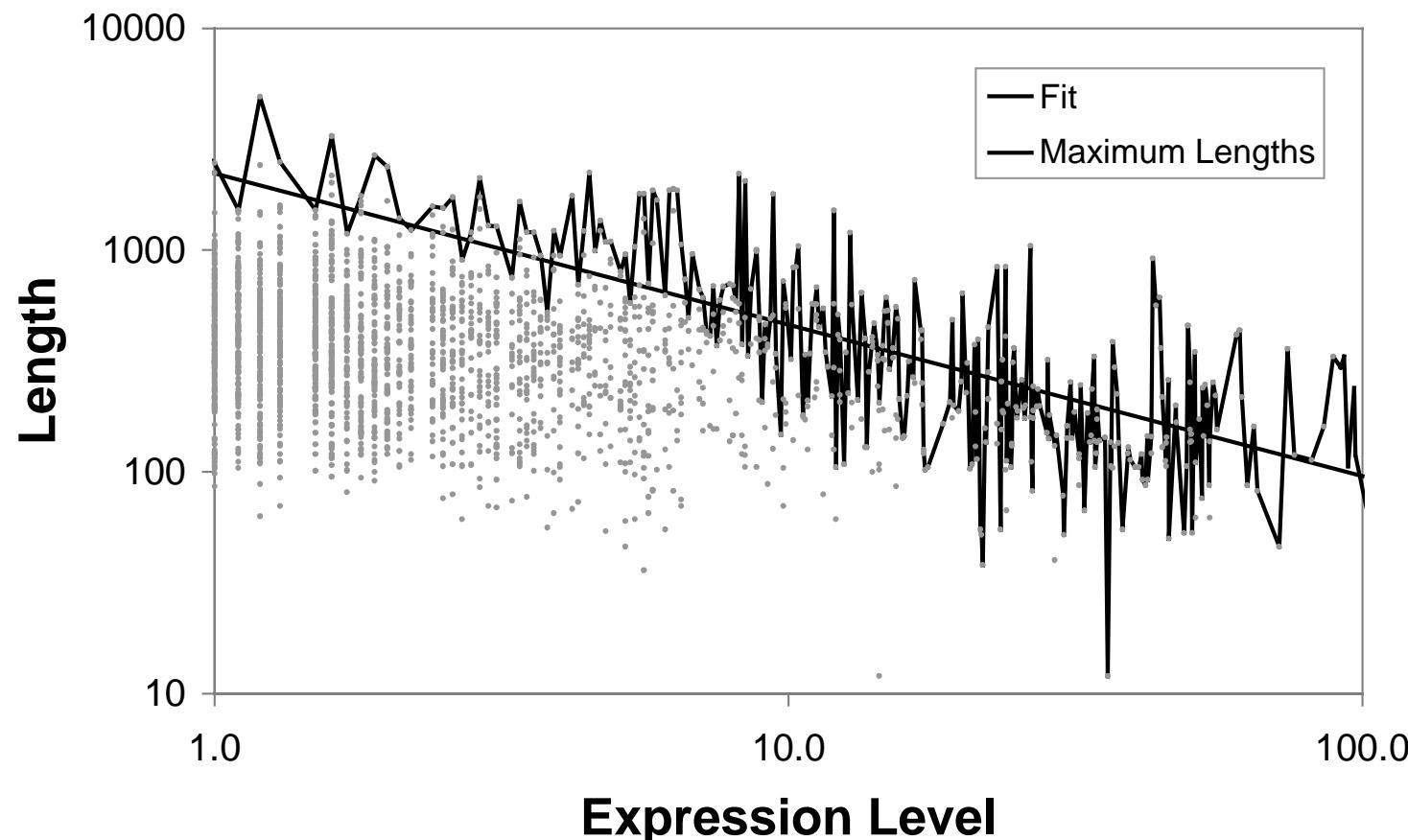
	$\sum_{\text{orf } i} n_i(F)$	$\sum_F \sum_{\text{orf } i} n_i(F)$	$G(F)$	$\sum_{\text{orf } i} e_i n_i(F)$	$\sum_F \sum_{\text{orf } i} e_i n_i(F)$	$T(F)$	$D(F)$
Feature F is Amino acids, in particular Ala	Number of Ala in yeast	Number of amino acids in yeast	Genome composition of Ala in yeast	Number of Ala weighted by expression	Number of amino acids weighted by expression	Transcriptome composition of Ala in yeast	Relative enrichment of Ala in transcriptome
Spec. Num.	141890	2574876	5.5%	347807	4758441	7.3%	32.7%
Feature F is Folds, in particular the TIM-barrel (3.1)	Number of TIM-barrel fold matches in yeast genome	Number of matches with all folds in yeast genome	Genome composition of TIM-barrel fold matches	Number of TIM-barrel fold matches weighted by expression	Number of matches with all folds weighted by expression	Transcriptome composition of TIM-barrel fold matches	Relative enrichment of TIM-barrel matches in transcriptome
Spec. Num.	65	1560	4.2%	389	4709	8.3%	97.8%



Relation between Length & Expression

Max Expression (e.g. transcripts/cell) \sim (Length) $^{-2/3}$

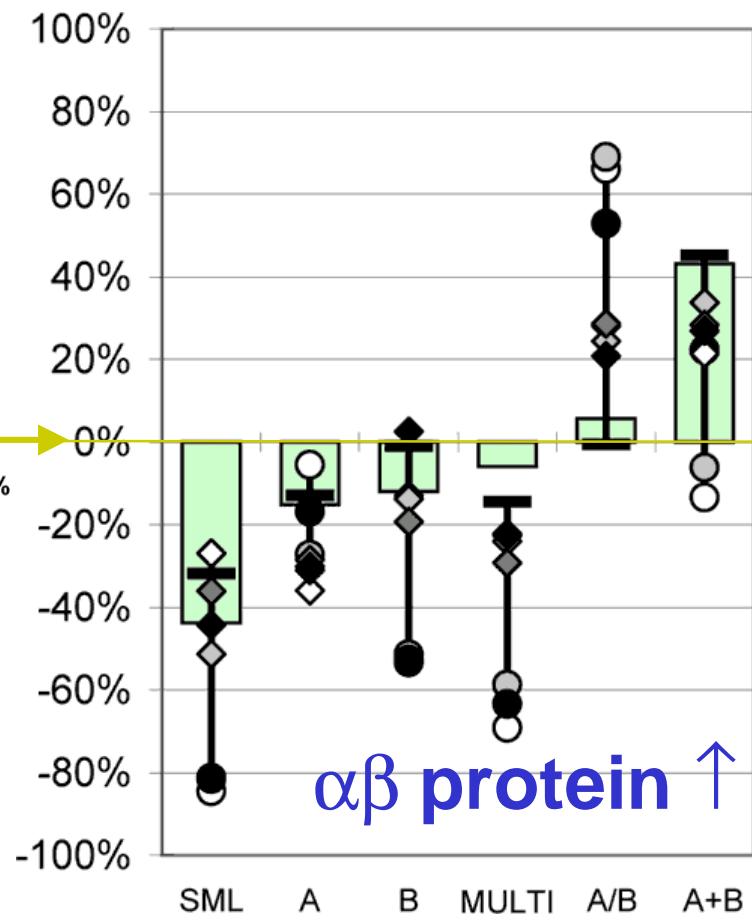
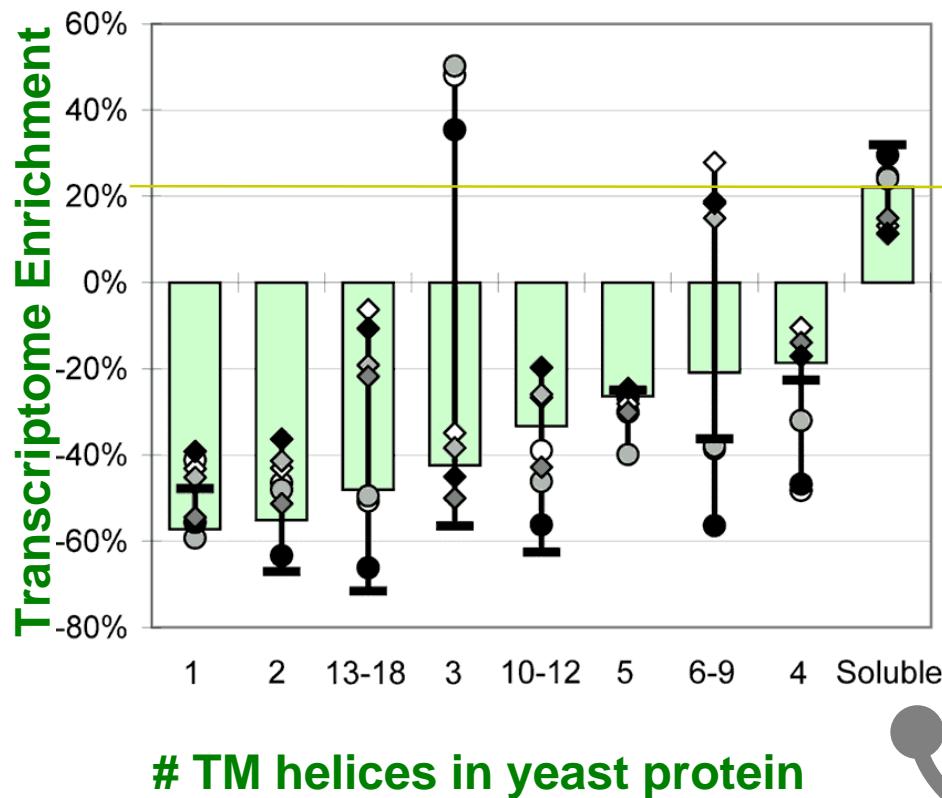
Shorter proteins can be more highly expressed



Composition of Transcriptome in terms of Broad Structural Classes

- Holstege et al. (9)
- Jelinsky et al. (11)
- ◆ Roth et al., mating type a (10)
- ◆ Roth et al., mating type alpha (10)
- ◆ Roth et al., galactose (10)
- ◆ Roth et al., heat shock (10)
- SAGE, G2/M phase (1)
- SAGE, log phase (1)
- SAGE, S phase (1)

Membrane (TM) Protein ↓



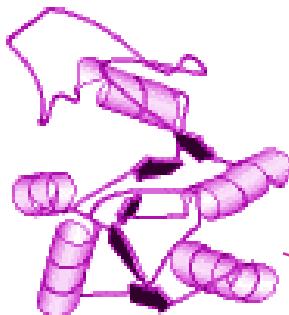
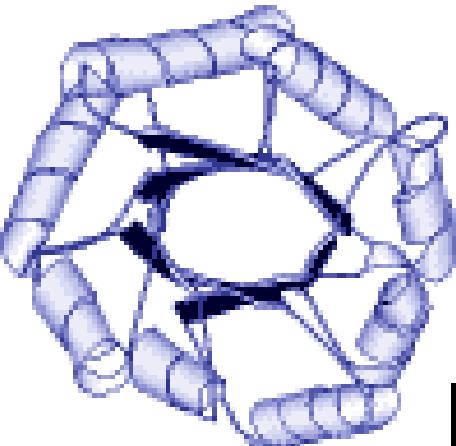
$\alpha\beta$ protein ↑

Fold Class of Soluble Proteins

Which Protein Folds

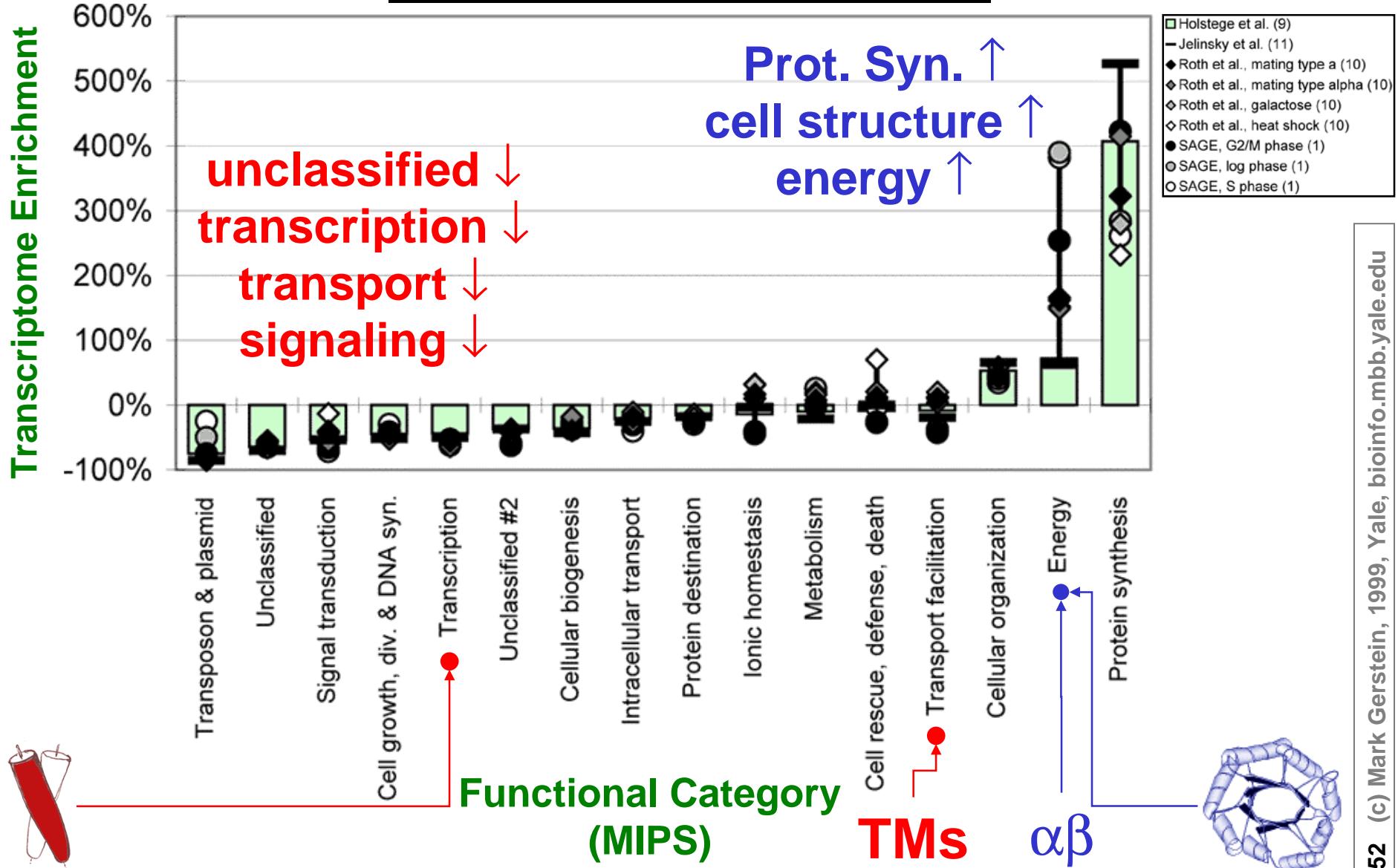
are Highly Expressed?

Top-10 folds in genome and transcriptome

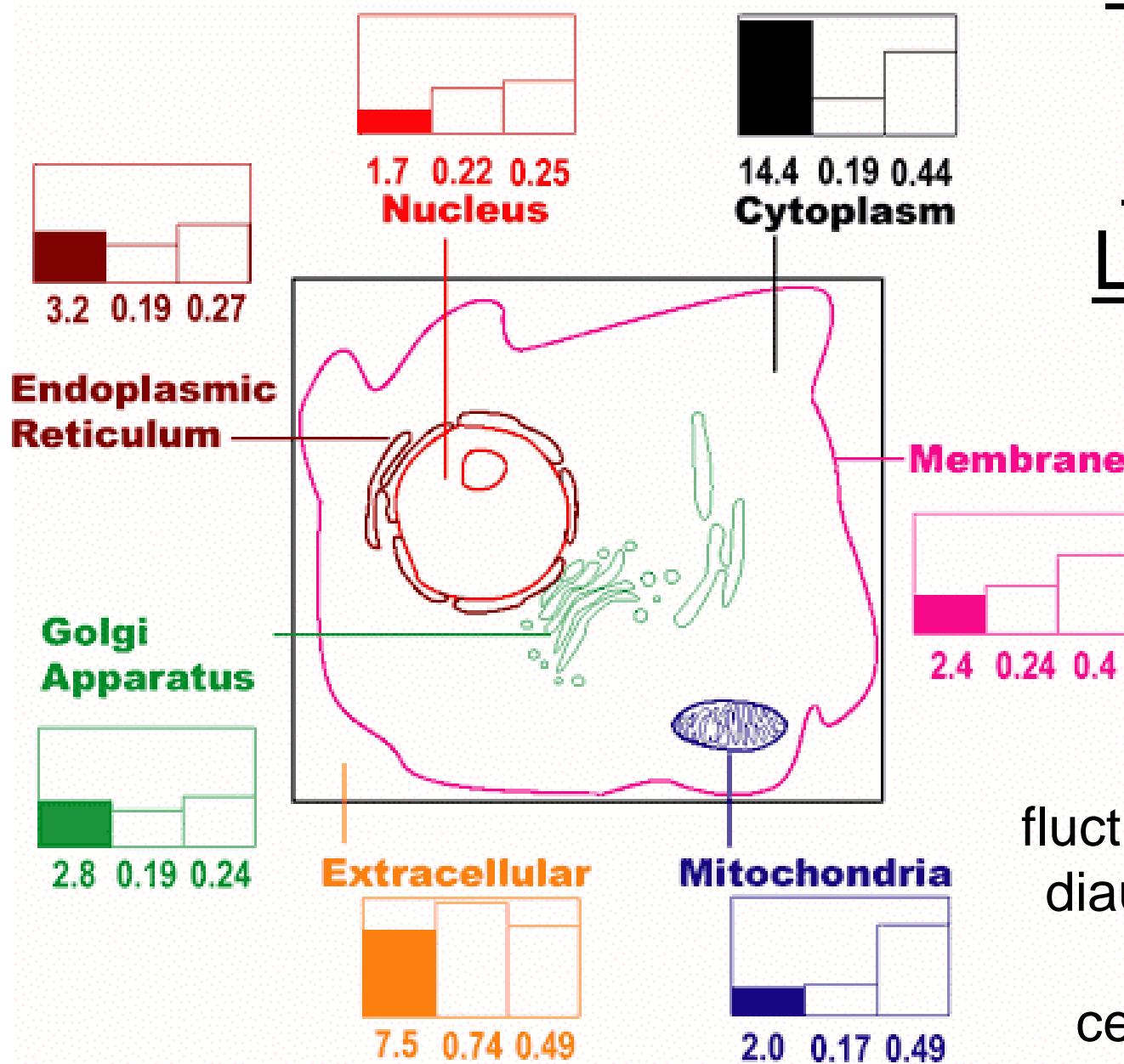


Fold	Fold Class	Rep. PDB	Composition			Rank										
			Genome [%]	Transcriptome [%]	Rel. Diff. [%]											
TIM barrel	α/β	1byb	4.2	8.3	+98	+	5	1	1	1	1	1	1	1		
P-loop NTP hydrolases	α/β	1gky	5.8	5.2	-11	•	3	2	2	4	4	4	5	6	7	
Ferredoxin like	α/β	1fxd	3.9	3.4	-14	•	6	3	7	11	9	8	10	4	10	11
Rossmann fold	α/β	1xel	3.3	3.3	0	•	8	4	3	3	3	2	2	19	15	9
7-bladed beta-propeller	β	1mda*	6.4	2.9	-55	-	2	5	4	5	6	6	7	9	9	16
alpha-alpha superhelix	α	2bct	4.4	2.7	-37	-	4	6	11	15	16	12	12	8	5	8
Thioredoxin fold	α/β	2trx	1.7	2.7	+63	+	14	7	6	8	2	5	4	11	10	6
G3P dehydrogenase-like	α/β	1drwt	0.2	2.7	+1316	+	81	8	12	2	5	3	3	35	19	30
beta grasp	α/β	1lgd	0.6	2.6	+348	+	36	9	10	21	9	18	21	82	122	120
HSP70 C-term. fragment	multi	1dky	0.8	2.6	+231	+	31	10	16	17	11	16	12	48	25	56
long helices oligomers	α	1zta	3.8	2.1	-46	-	7	15	8	14	21	15	19	21	20	33
Protein kinases (cat. core)	multi	1hcl	6.8	1.6	-77	-	1	18	19	9	16	11	15	13	16	17
alpha/beta hydrolases	α/β	2ace	2.2	0.9	-62	-	10	32	31	25	26	21	23	26	26	26
Zn2/C6 DNA-bind. dom.	sml	1aw6	2.6	0.3	-89	-	9	75	94	27	50	32	40	48	39	50

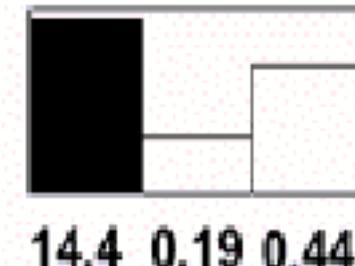
Composition of Transcriptome in terms of Functional Classes



Expression Level is Related to Localization



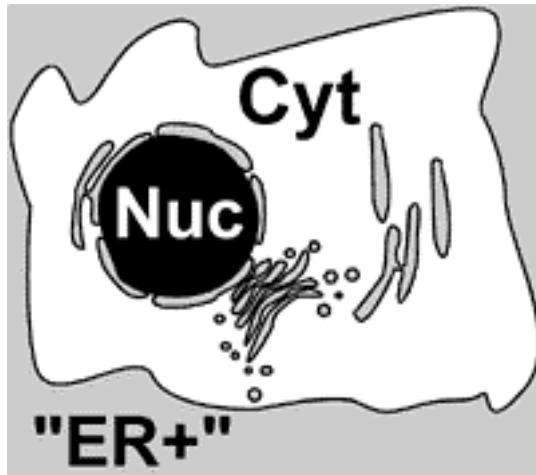
abs. exp.
(GeneChip)



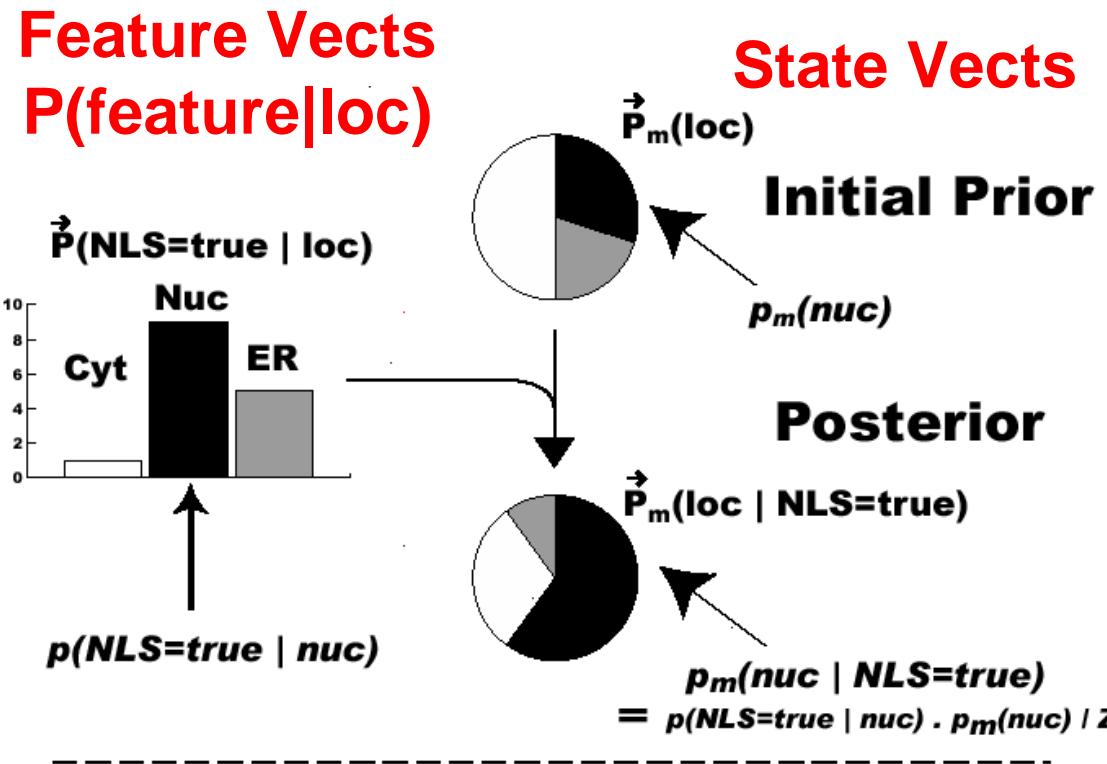
fluctuations in
diauxic shift
or
cell-cycle

Bayesian System for Localizing Proteins

loc=



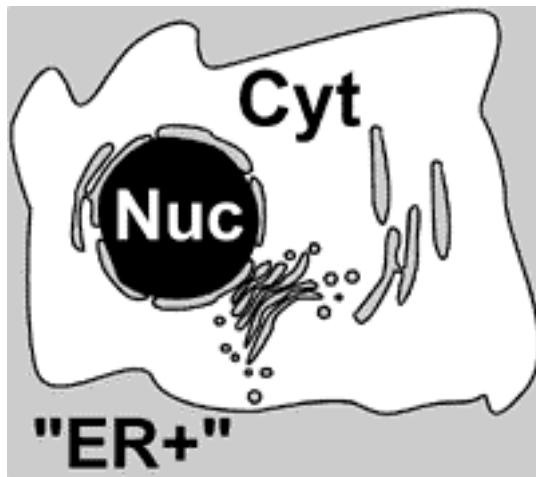
Represent localization of each protein by the state vector $\mathbf{P}(\text{loc})$ and each feature by the feature vector $\mathbf{P}(\text{feature}|\text{loc})$. Use Bayes rule to update.



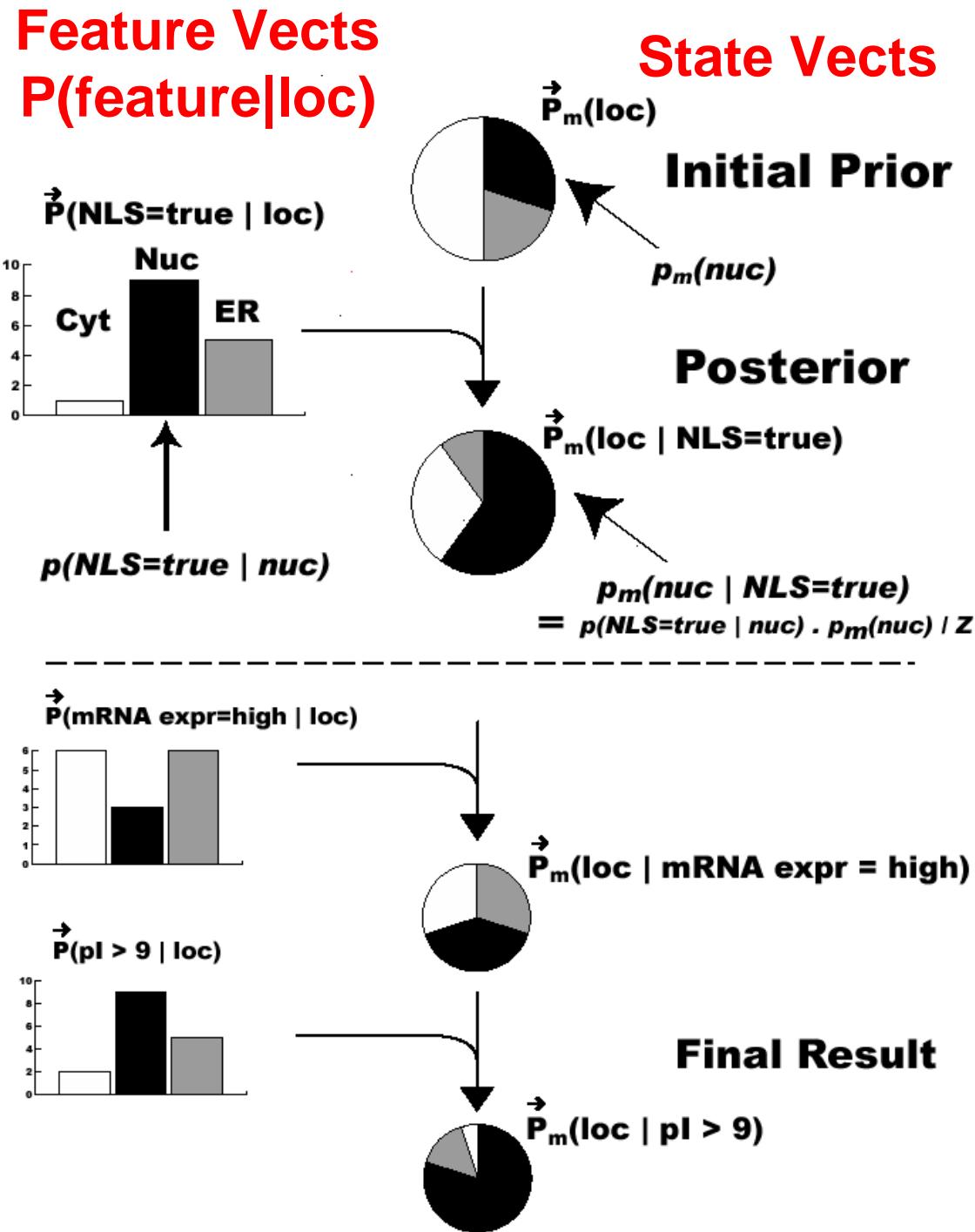
18 Features: Expression Level
(absolute and fluctuations), signal
seq., KDEL, NLS, Essential?, aa
composition

Bayesian System for Localizing Proteins

loc=



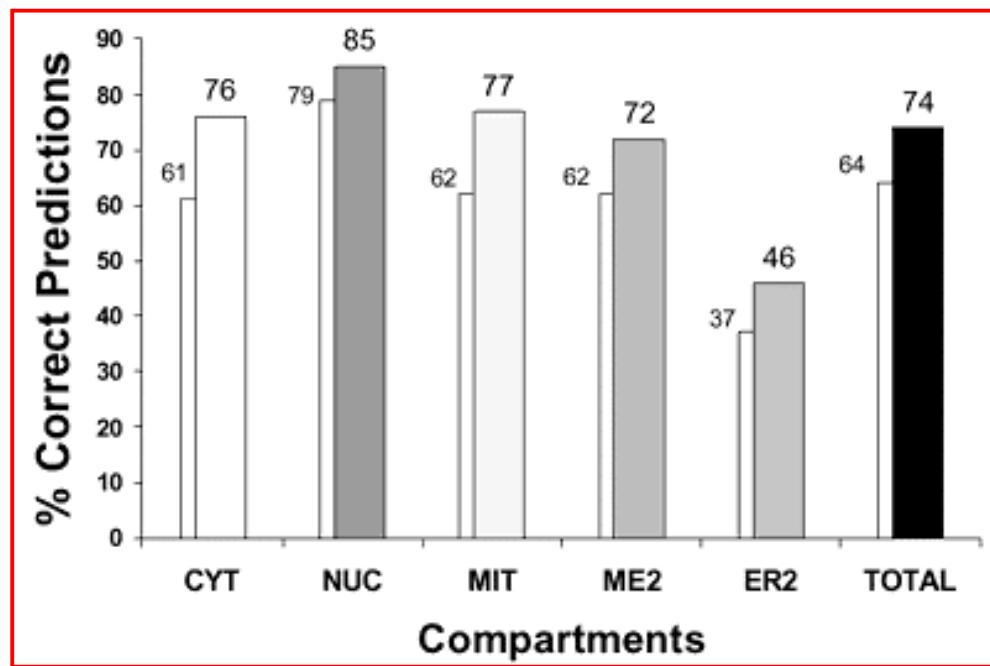
Represent localization of each protein by the state vector $\mathbf{P}(\text{loc})$ and each feature by the feature vector $\mathbf{P}(\text{feature}|\text{loc})$. Use Bayes rule to update.



Results on Testing Data

Individual
proteins:
74% with
cross-
validation

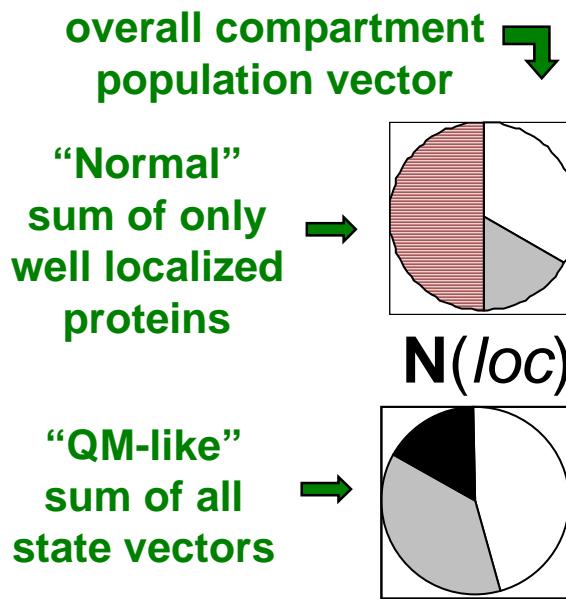
Testing, training
data, Priors: ~2000
proteins from YPD,
MIPS, SwissProt,
Snyder Lab



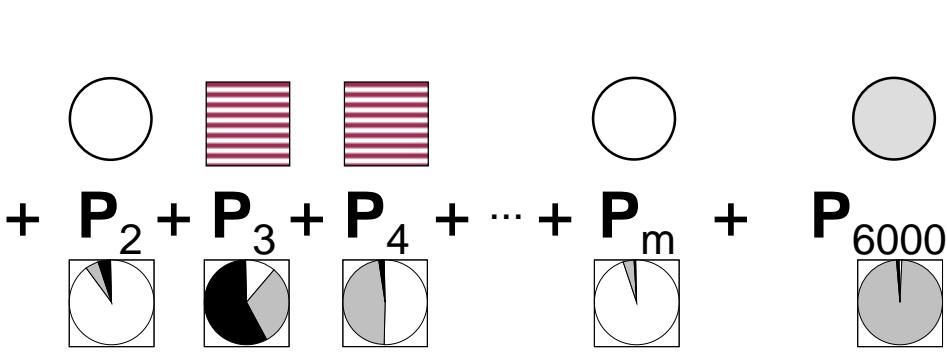
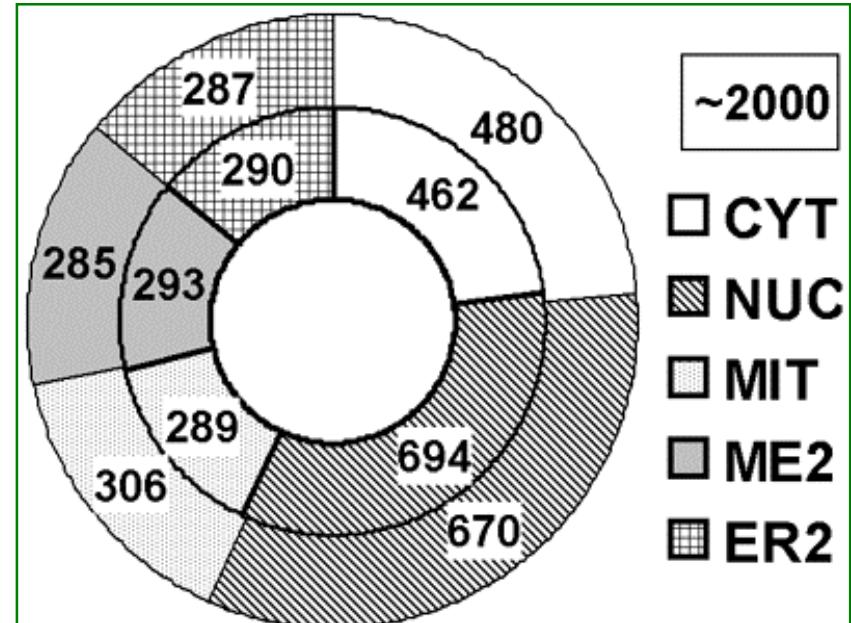
Results on Testing Data #2

Compartment Populations.

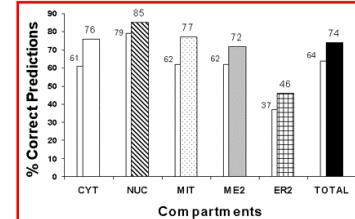
directly sum state vectors to get population. Gives **96%** pop. similarity.



state vector of protein 1



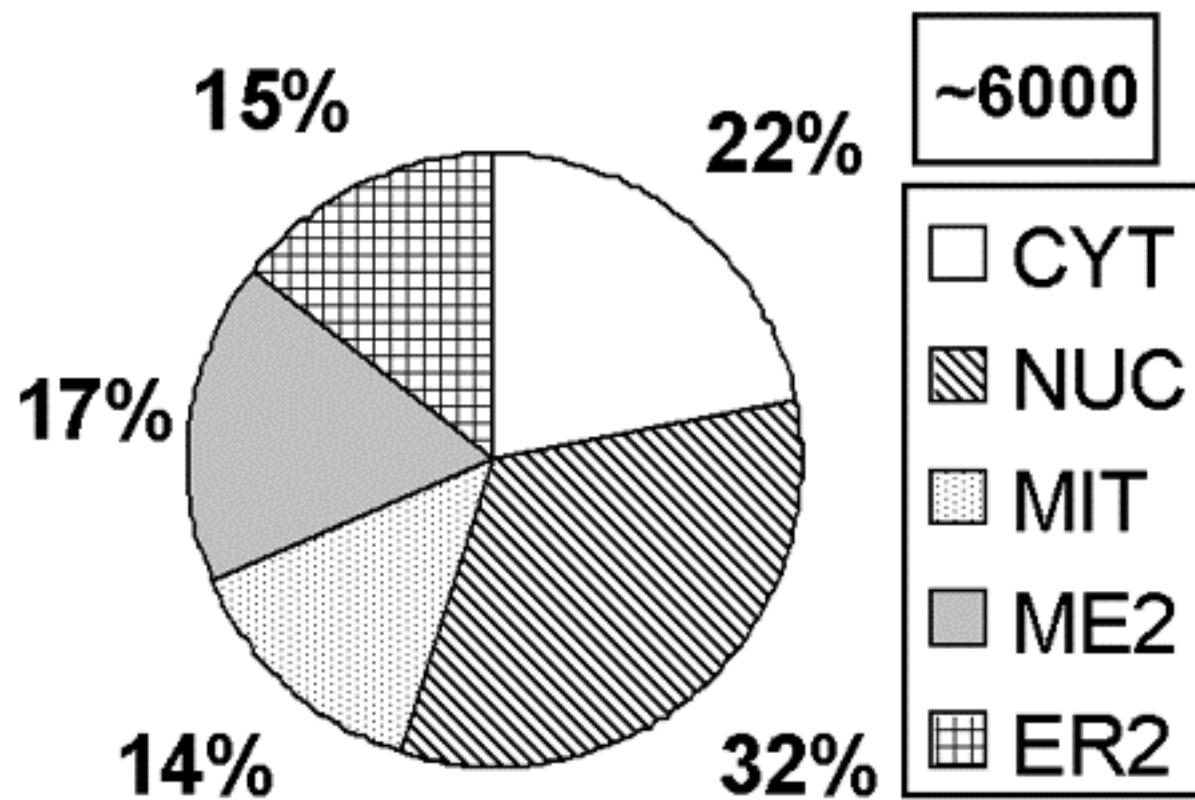
Individual proteins: 74% with cross-validation



Extrapolation to Compartment

Populations of Whole Yeast Genome:

~4000 predicted + ~2000 known

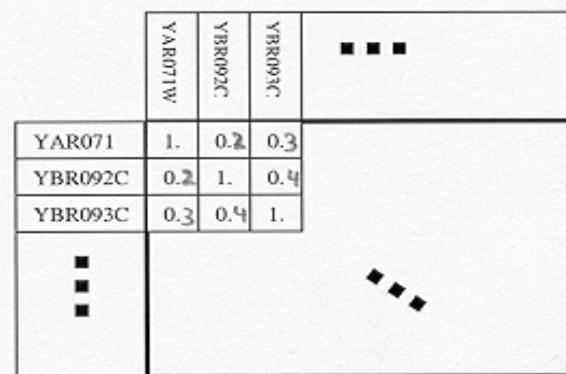
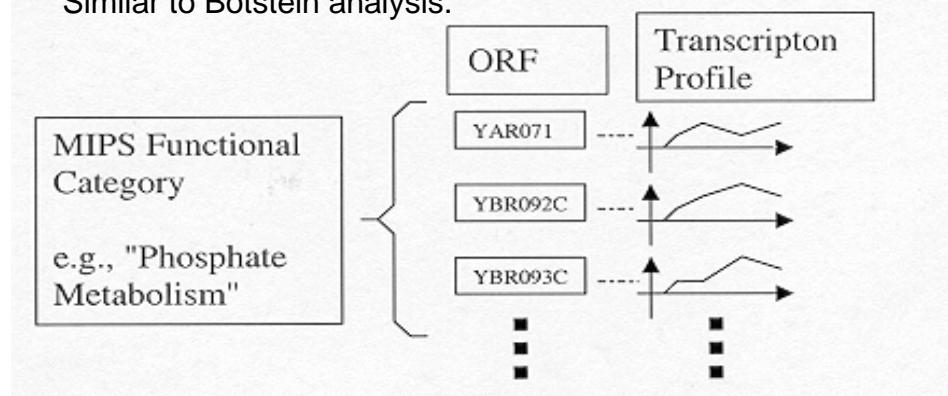


Functional category number	Function	Average correlation	# ORFs
01	METABOLISM	0.1001	1005
01.01	amino-acid metabolism	0.1488	199
01.01.01	amino-acid biosynthesis	0.239	114
01.01.04	regulation of amino-acid metabolism	0.23	32

MIPS YFC: 66 bottom classes, 10 top classes

Average correlation of uncharacterized genes is 0.16

Similar to Botstein analysis.

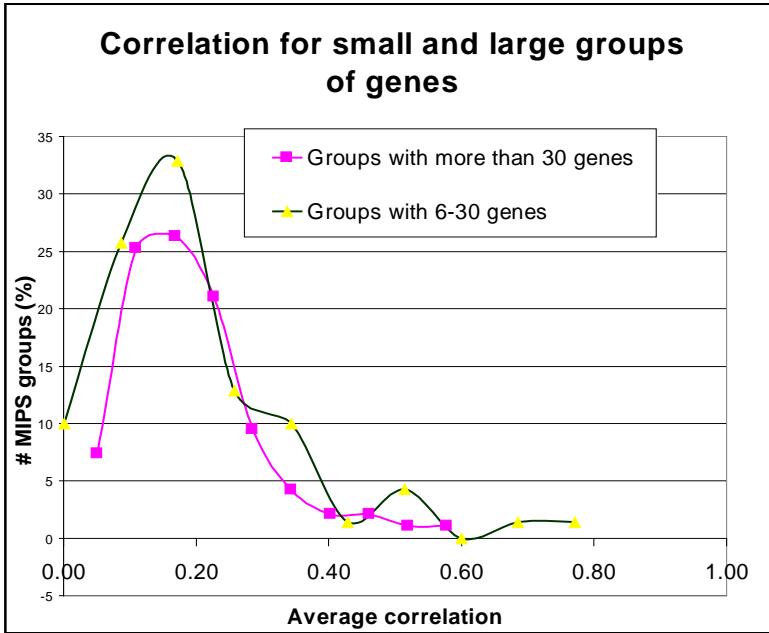


Correlation Coefficient Matrix (Pearson Coefficient)

Average Correlation Coefficient
for Group of Genes

Cluster Expression Data, Relate to MIPS Functional Category

Functional category number	Function	Average correlation	# ORFs
01	METABOLISM	0.1001	1005
01.01	amino-acid metabolism	0.1488	199
01.01.01	amino-acid biosynthesis	0.239	114
01.01.04	regulation of amino-acid metabolism	0.23	32
01.01.07	amino-acid transport	0.1198	23
01.01.10	amino-acid degradation	0.0524	42
01.01.99	other amino-acid metabolism activities	0.2205	28
01.02	nitrogen and sulphur metabolism	0.1869	6
01.02.01	nitrogen and sulphur utilization	0.0726	37
01.02.04	regulation of nitrogen and sulphur utilization	0.3715	18
01.02.07	nitrogen and sulphur transport	0.2829	10
01.03	nucleotide metabolism	0.1708	12
01.03.01	purine-ribonucleotide metabolism	0.3639	42
01.03.04	pyrimidine-ribonucleotide metabolism	0.176	28
01.03.07	deoxyribonucleotide metabolism	0.1095	12
01.03.10	metabolism of cyclic and unusual nucleotides	0.2848	8
01.03.13	regulation of nucleotide metabolism	0.2696	13
01.03.16	polynucleotide degradation	0.2461	20
01.03.19	nucleotide transport	0.1183	12
01.03.99	other nucleotide-metabolism activities	-0.037	7
01.04	phosphate metabolism	0.1363	31
01.04.01	phosphate utilization	0.142	13
01.04.04	regulation of phosphate utilization	0.0599	8
01.04.07	phosphate transport	0.0724	10
01.05	carbohydrate metabolism	0.0779	409
01.05.01	carbohydrate utilization	0.075	256
01.05.04	regulation of carbohydrate utilization	0.1174	120

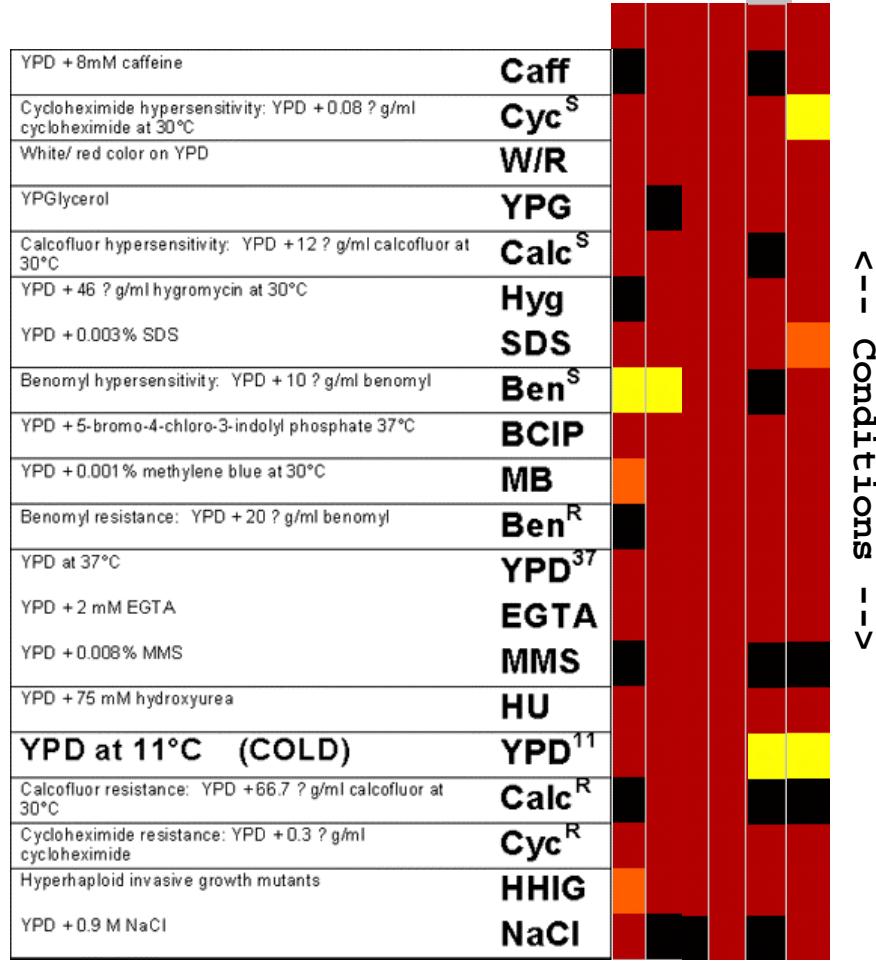


- Many groups of genes categorized by MIPS do not have higher correlation than random ORFs
- Smaller groups tend to have a slightly higher correlation

Correlation of Functional Class and Expression – not that strong

Highest Correlation

Functional category number	Function	Average correlation	# ORFs
10.04.11	key kinases	0.9403	2
10.04.13	key phosphatases	0.9283	2
11.11	ageing	0.8634	2
02.22	glyoxylate cycle	0.8136	6
10.02.07	G-proteins	0.8122	3
04.03.99	other tRNA-transcription activities	0.6932	4
09.08	biogenesis of Golgi	0.6647	2
09.19	peroxisomal biogenesis	0.6512	2
08.10	peroxisomal transport	0.646	12
04.01.04	rRNA processing	0.6074	53
01.20	secondary metabolism	0.5921	4
01.20.05	amines metabolism	0.5921	4
10.05.11	key kinases	0.5549	4
90	RETROTRANSPOSONS AND PLASMID PROTEINS	0.5299	7
02.10	tricarboxylic-acid pathway	0.5236	22
04.07	RNA transport	0.5111	27



M Snyder

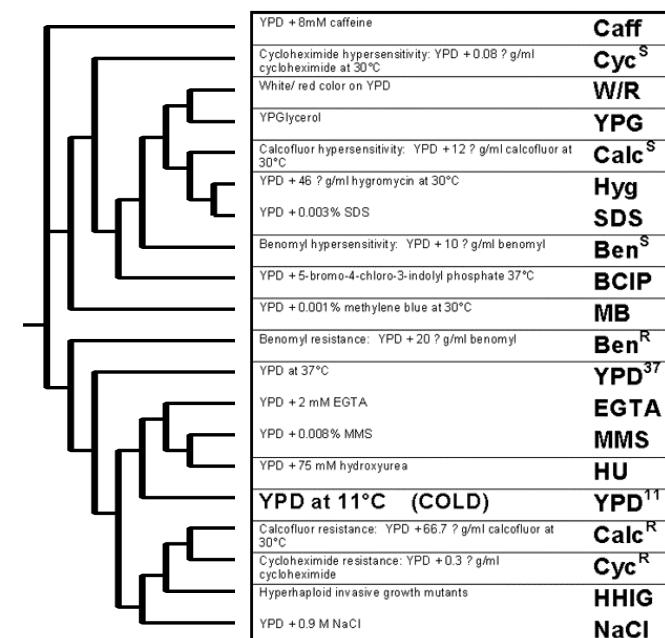
Affected by Another Condition

WT

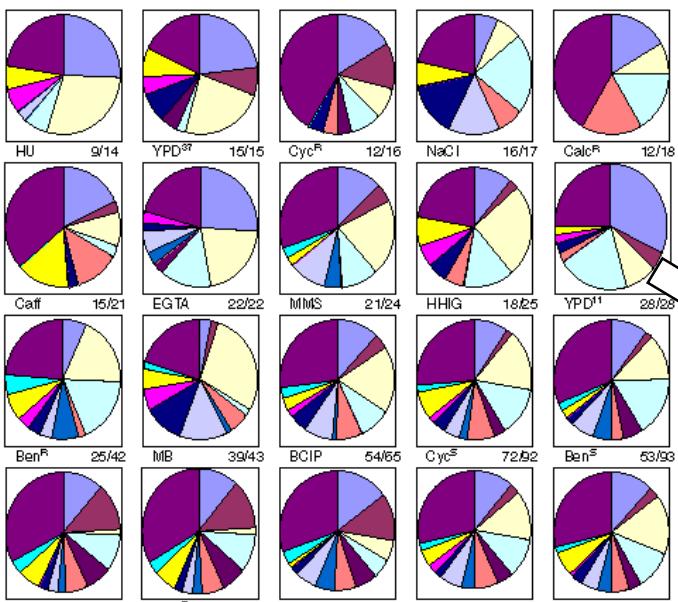
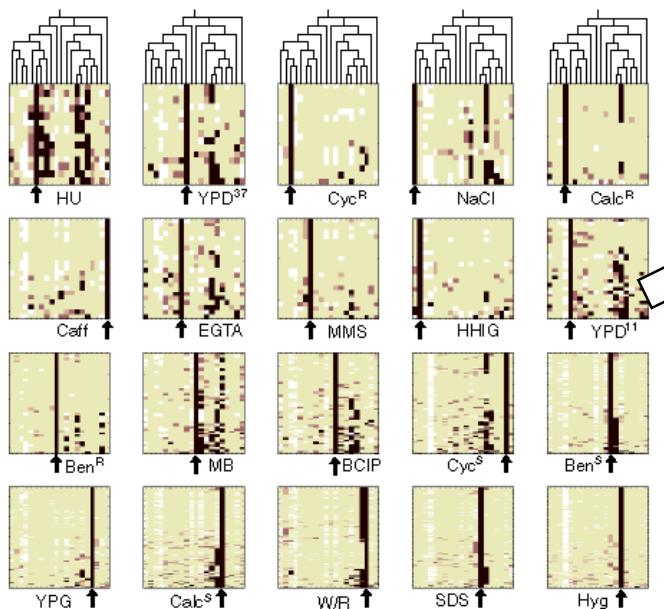
Affected by Cold

Whole Genome Phenotype Profiles

Transposon insertions into (almost) each yeast gene to see how yeast is affected in 20 conditions. Generates a phenotype pattern vector, which can be treated **similarly to expression data**



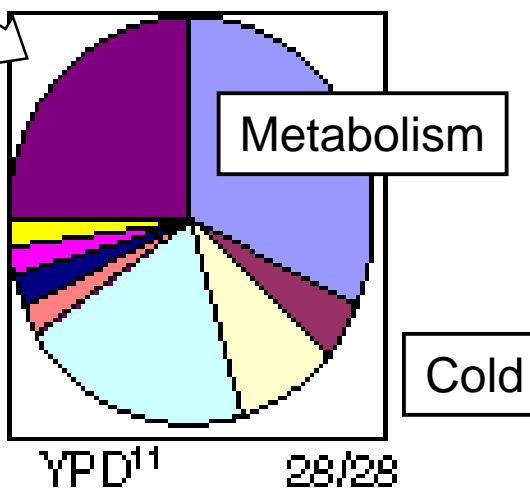
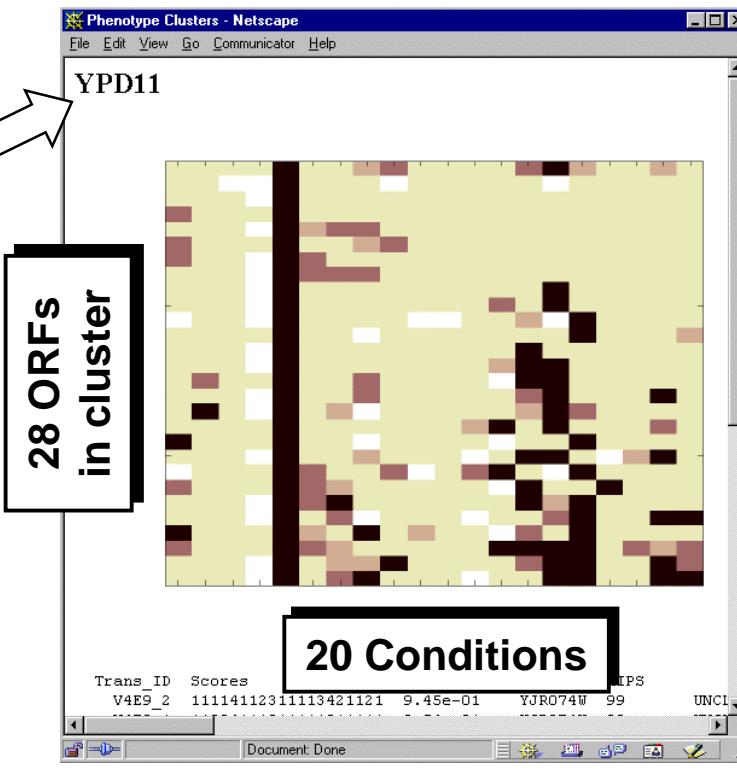
Clustering Conditions



Legend:

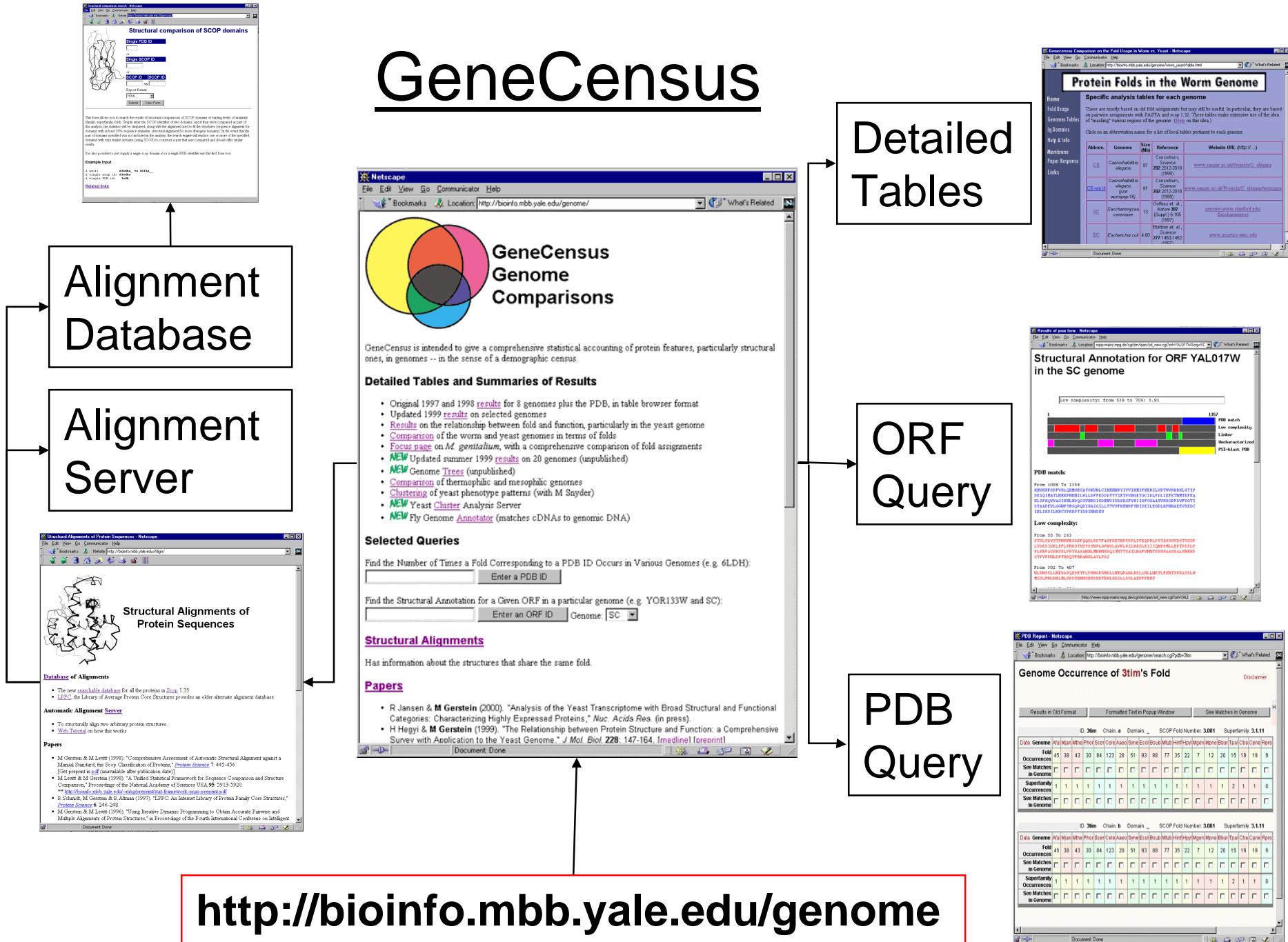
- METABOLISM
- CELL GROWTH, DIVISION AND DNA SYNTHESIS
- PROTEIN SYNTHESIS
- TRANSPORT FACILITATION
- CELLULAR BIOGENESIS
- CELL RESCUE, DEFENSE, CELL DEATH AND AGEING
- CELLULAR ORGANIZATION
- ENERGY
- TRANSCRIPTION
- PROTEIN DESTINATION
- INTRACELLULAR TRANSPORT
- SIGNAL TRANSDUCTION
- IONIC HOMEOSTASIS

Phenotype ORF Clustering



k-means clustering of ORFs based on “phenotype patterns,” cross-ref. to MIPs Functional Classes

Cluster showing cold phenotype (containing genes most necessary in cold) is enriched in metabolic functions



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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

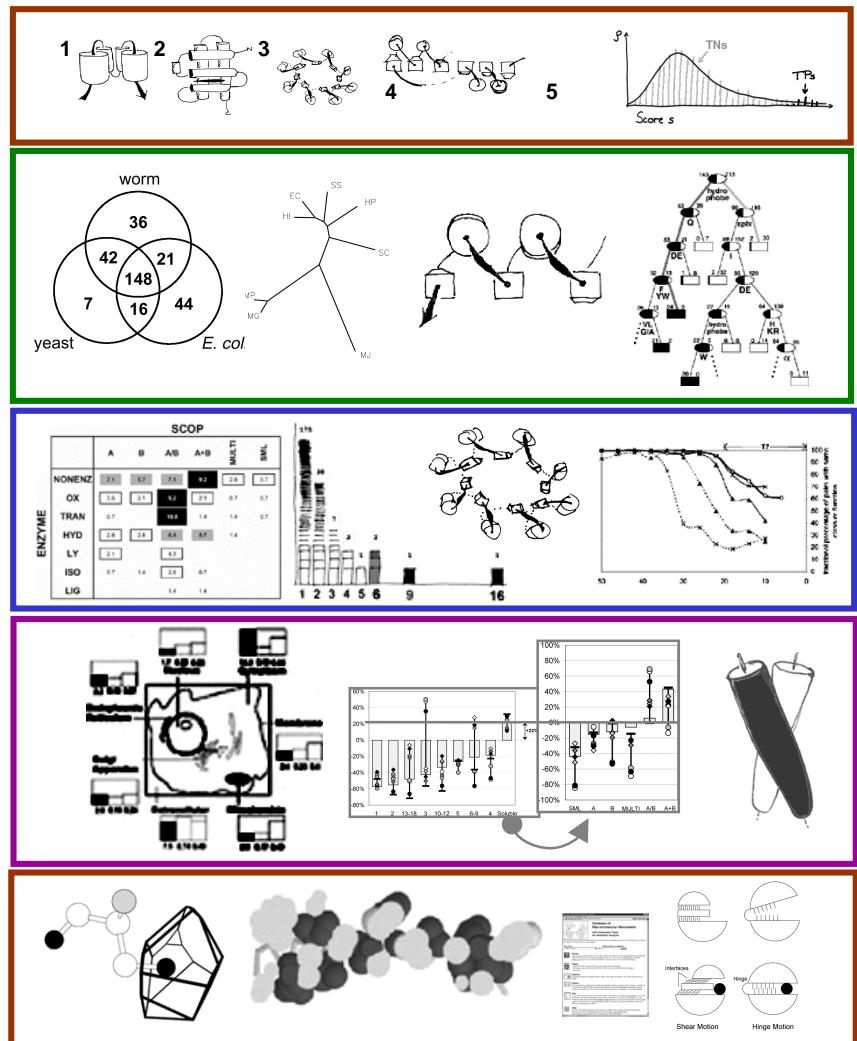
4 Folds in the Transcriptome

(Common parts? Where are parts?)

Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis, cyt. Depleted ↓ : NS, long, TM folds, transport, transcription, Leu-zip, nuc. Bayesian Localizer, phenotypes clustering

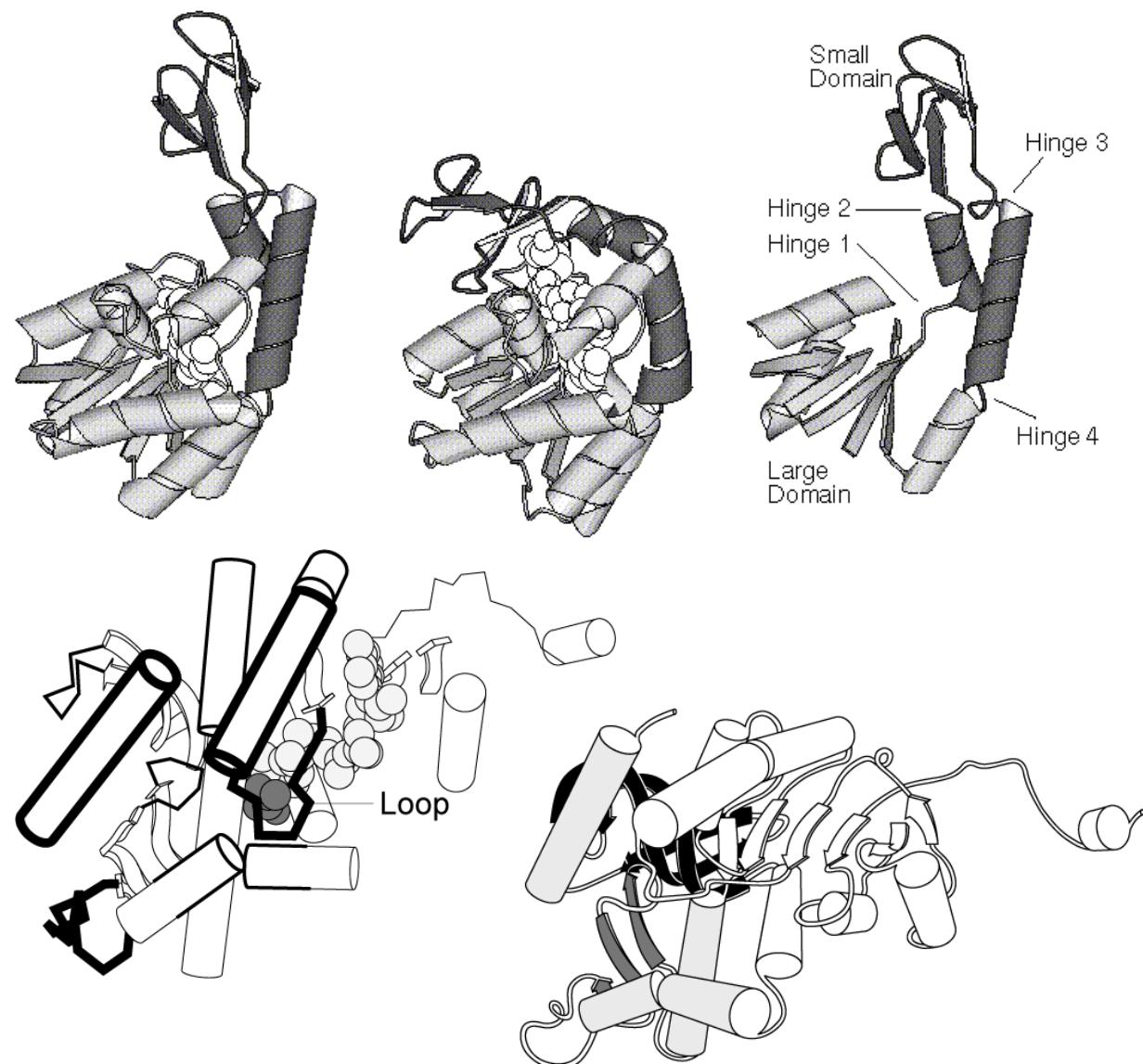
5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson, R Das, H Hegyi, J Lin, Y Kluger, C Arrowsmith, A Edwards, L Regan, S Balasubramanian, A Drawid, D Greenbaum, M Snyder, R Jansen



- What are they?
 - ◊ Proteins, Nucleic Acids (Hammerhead)
 - ◊ Sidechains (trivial), Loops (LDH), Domains (ADK), Subunits (Hb)
 - ◊ When a Ligand Binds: Open, Closed
- Essential link between structure and function
 - ◊ catalysis, regulation, transport, formation of assemblies, and cellular locomotion
- A complicated biological phenomena that can be studied in quantitative detail
 - ◊ changes in thousands of atomic coordinates

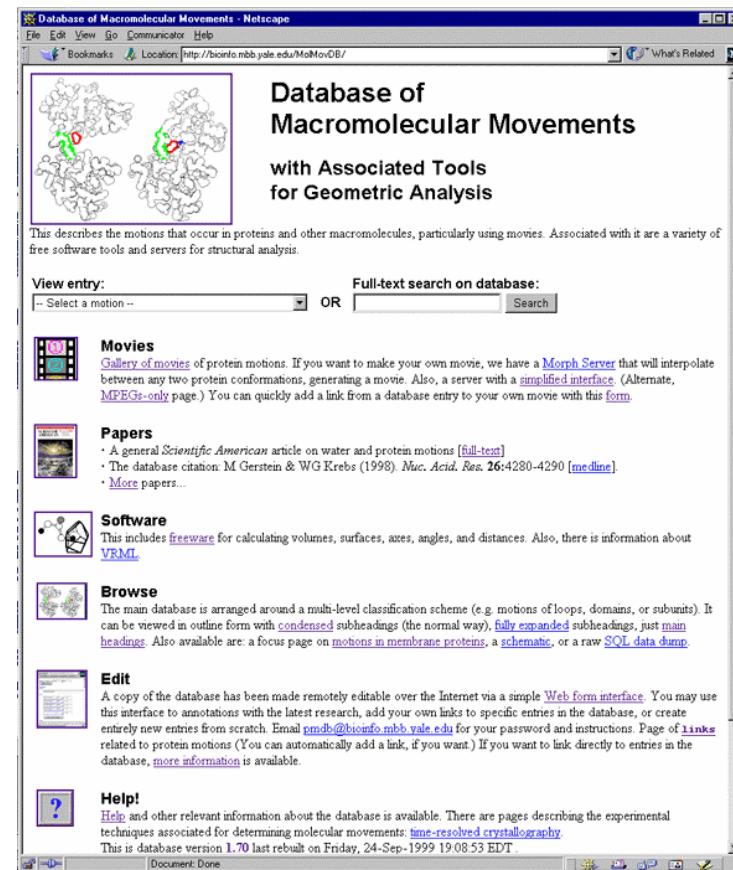
Macromolecular Motions



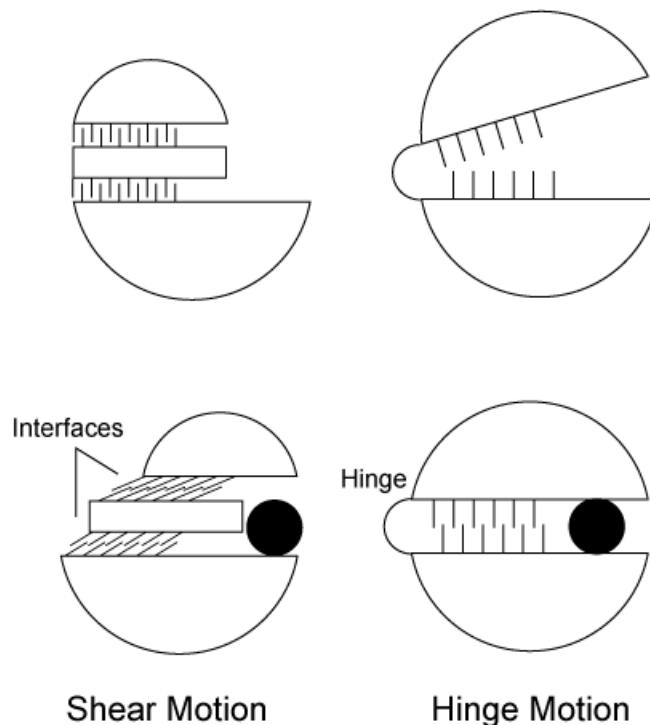
Number Known Forms	Size of Motion	Mechanism of Motion	Examples	#
Motion	Fragment	Hinge Shear Unclassifiable	TIM, LDH, TGL Insulin MS2 Coat	14 3 3
		Hinge Shear Refold Special Unclassifiable	LF, ADK, CM CS, TrpR, AAT Serpin, RT Ig elbow TBP, EF-tu	16 8 3 1 3
		Allosteric Non-allosteric Unclassifiable	PFK, Hb, GP Ig VL-VH	4 2
	Domain	Hinge Shear Unclassifiable	bR	1
		Refold Hinge Shear Special Unclassifiable	LF~TF, SBP HK~PGK, HSP	10 4
		Allosteric Non-allosteric Unclassifiable	Myosin	4
		Allosteric Non-allosteric Unclassifiable	PCNA, GroEL	3
	Subunit	Allosteric Non-allosteric Unclassifiable		

bioinfo.mbb.yale.edu/MoIMovDB

Motions DB: Information, Size, then Packing Based Classification

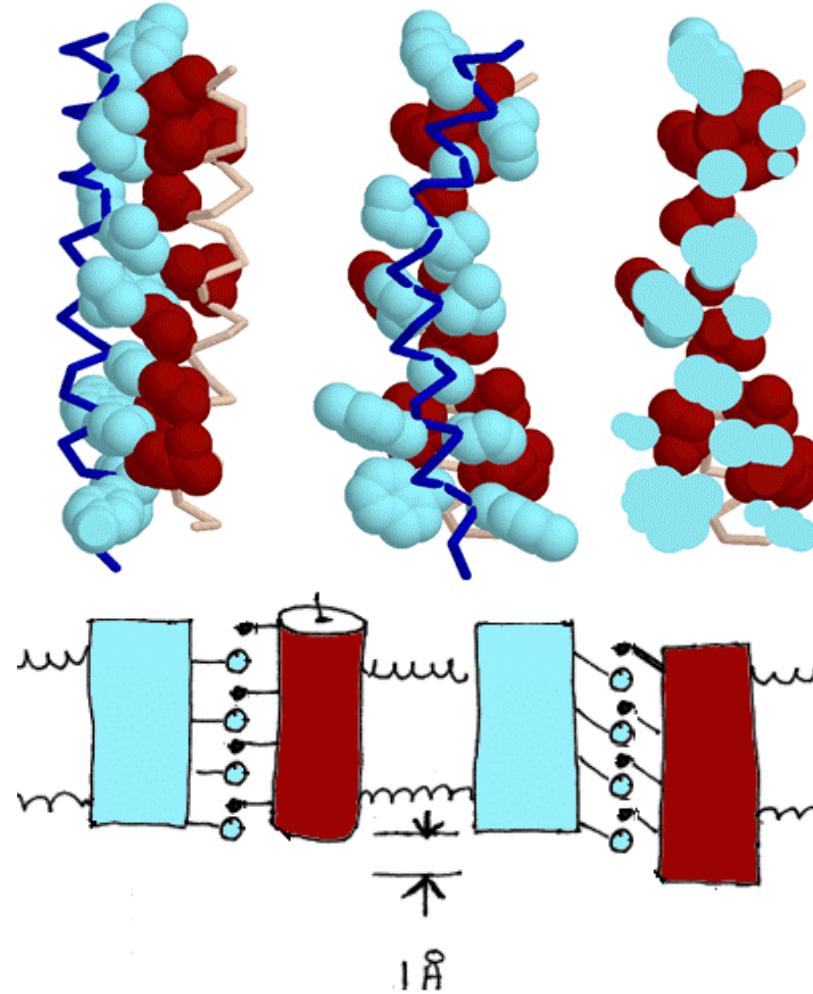


Interface Packing and Motions

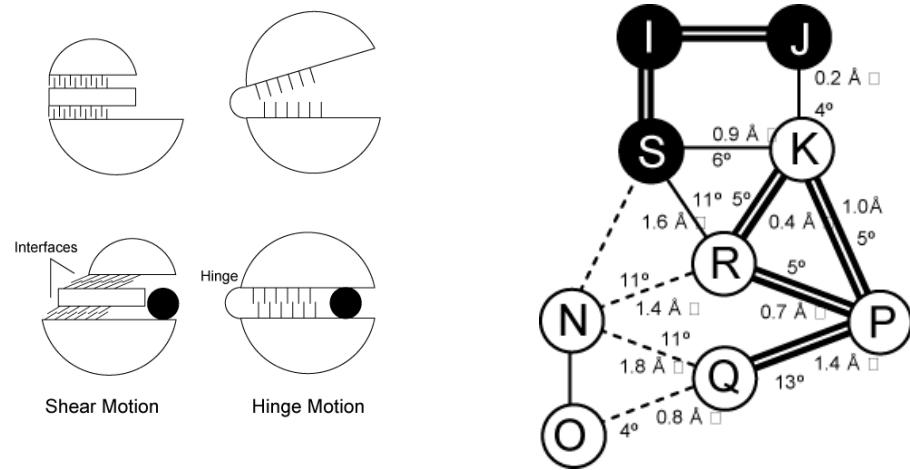
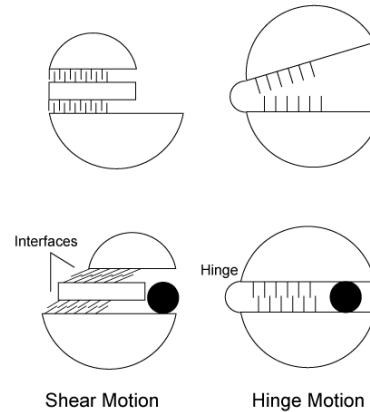
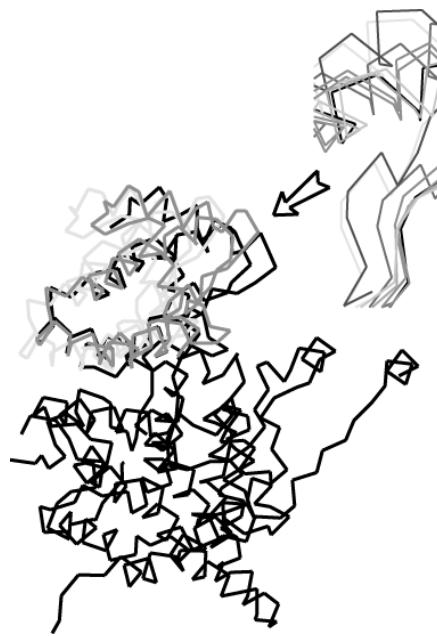
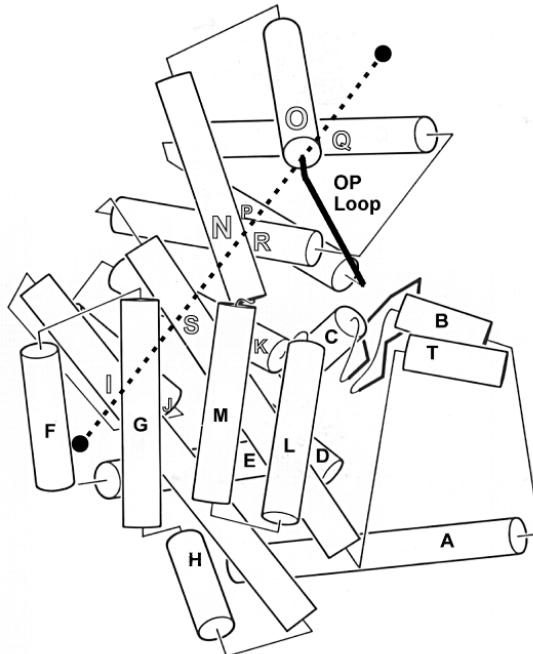


Not applicable to allosteric motions, which are much slower (~1 ms) and do involve repacking interfaces

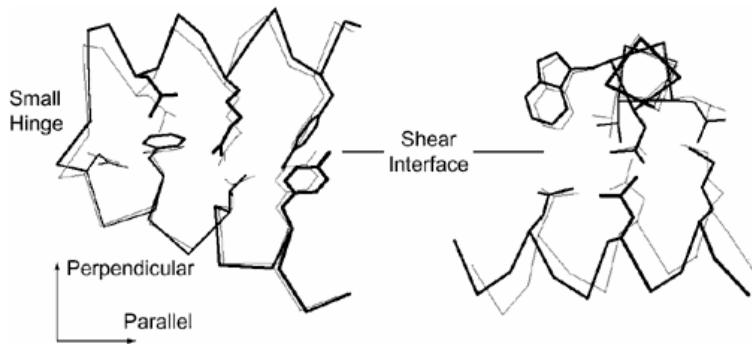
- Intercalating Interface, Knobs into Holes
- Packing is a strong constraint on motions
 - ◊ Domain or loop motions have to be fast (~10 ps – 100 ns)
 - ◊ Can't cross big energy barriers involved in repacking an interface



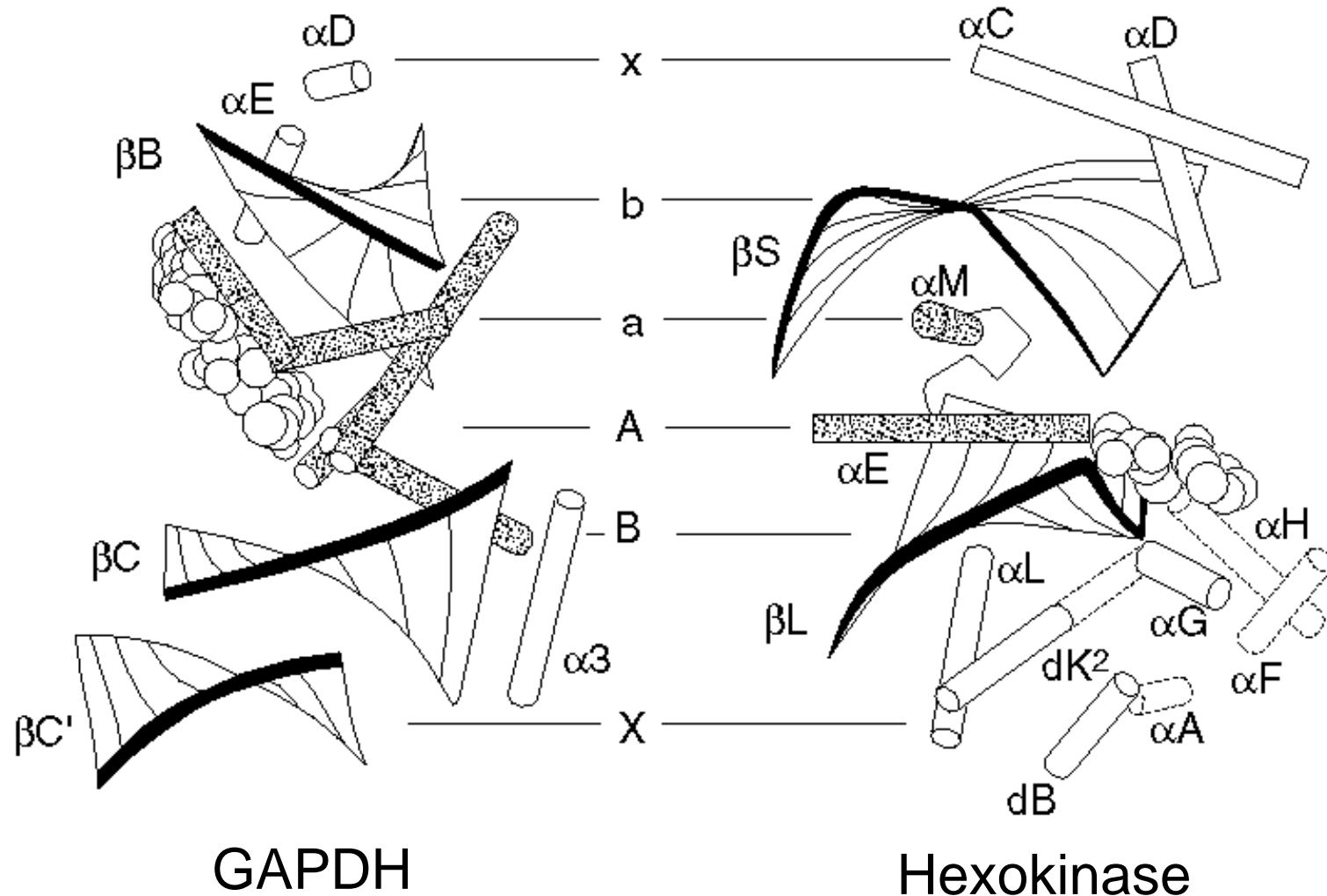
Packing Based Classification: Hinge v Shear



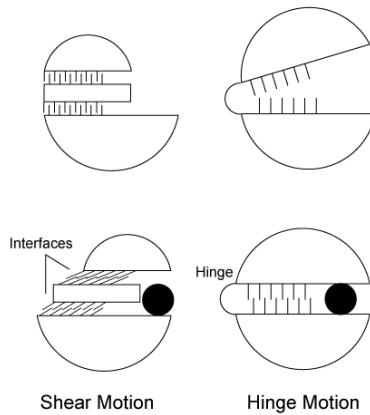
Shear Mechanism
Involves Many Small
Motions across a
Continuously
Maintained Interface



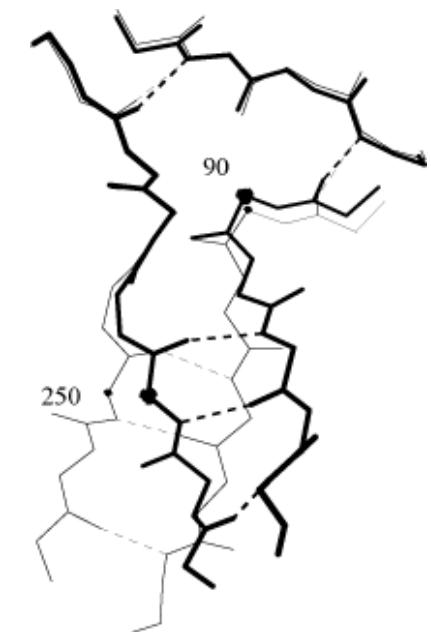
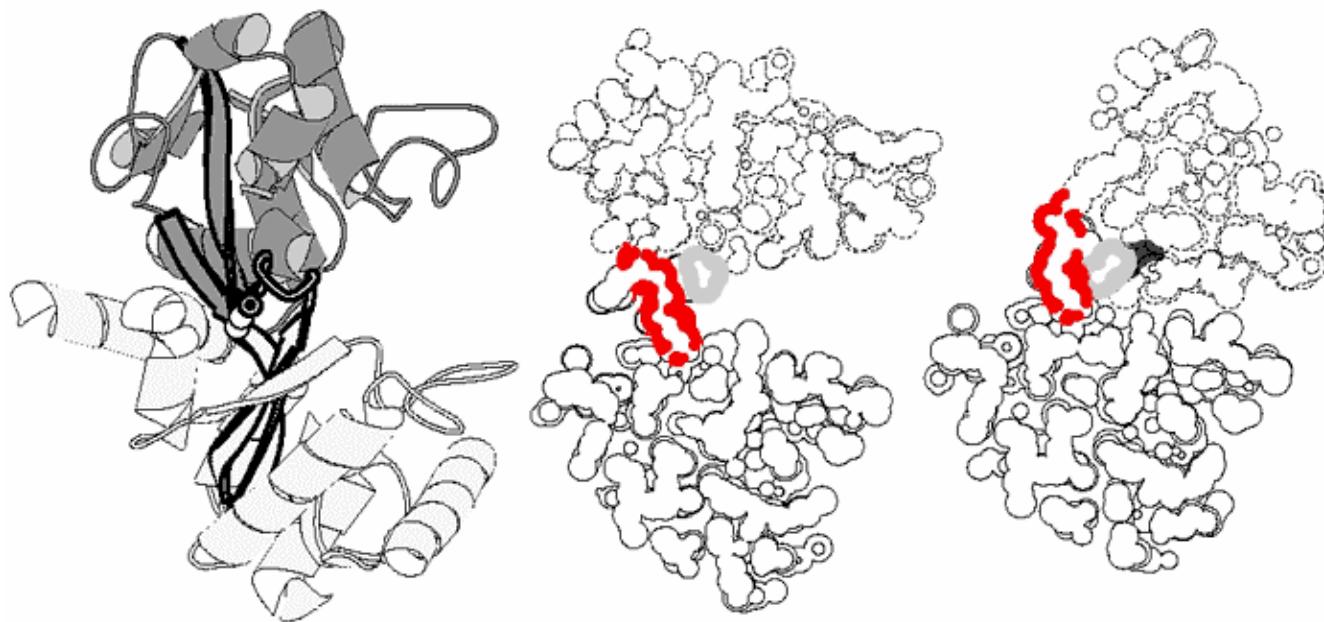
Proteins With Shear Motions are Often Divided into Layers



Packing Based Classification: Hinge v Shear



Hinge Mechanism involves absence of steric constraints (continuously maintained interface), esp. at hinge



Voronoi Volumes, the Natural Way to Measure Packing

- Each atom surrounded by a single convex polyhedron and allocated space within it
 - ◊ Allocation of all space (large V implies cavities)

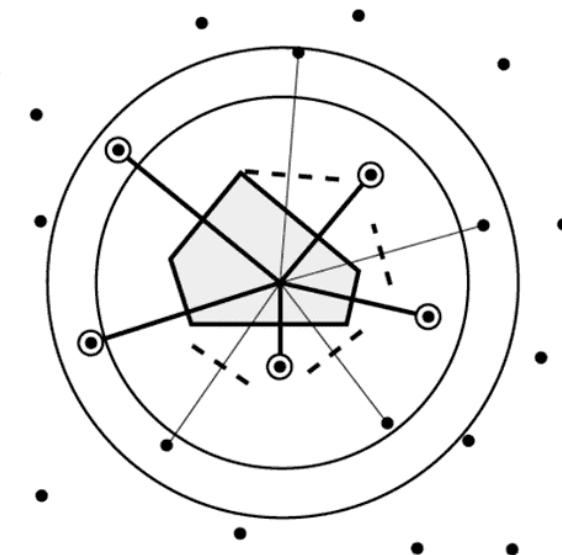
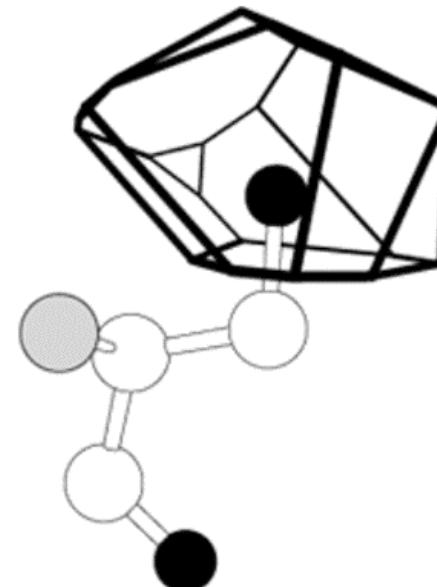
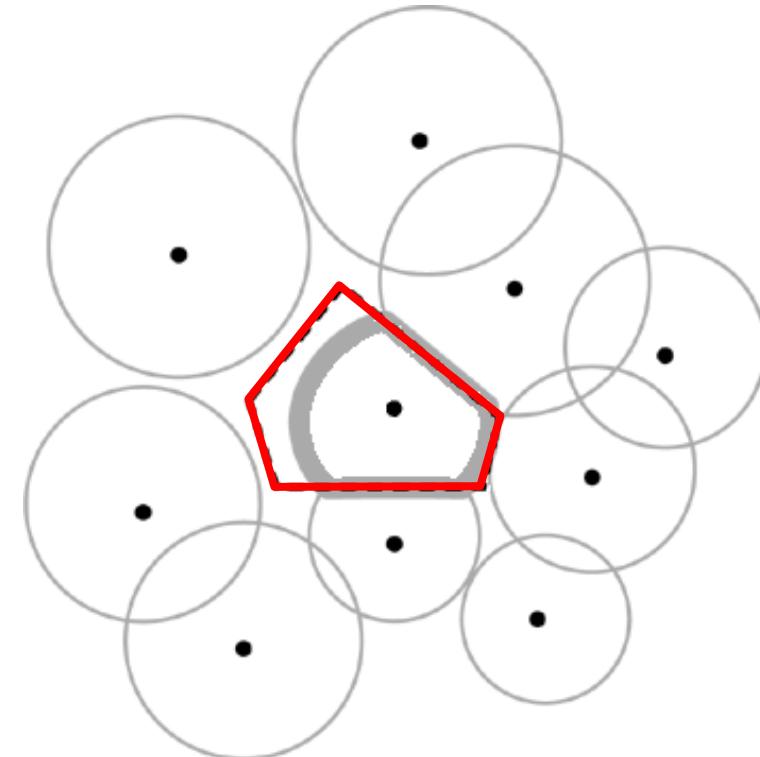
Packing Efficiency

$$= \frac{\text{Volume-of-Object}}{\text{Space-it-occupies}}$$

Space-it-occupies

$$= V(\text{VDW}) / V(\text{Voronoi})$$

Absolute v relative eff.
 V_1 / V_2



Unified Atoms		
atom	radii	volume
C3H0b	1.61	9.70
C3H0s	1.61	8.72
C3H1b	1.76	21.28
C3H1s	1.76	20.44
C4H1b	1.88	14.35
C4H1s	1.88	13.17
C4H2b	1.88	24.26
C4H2s	1.88	23.19
C4H3u	1.88	36.73
N3H0u	1.64	8.65
N3H1b	1.64	15.72
N3H1s	1.64	13.62
N3H2u	1.64	22.69
N4H3u	1.64	21.41
O1H0u	1.42	15.91
O2H1u	1.46	17.98
S2H0u	1.77	29.17
S2H1u	1.77	36.75

Residues		Parameters used in Protor Volume Derivation				
aa	volume	Typing Scheme	Hybrid chemical and numerical typing with 18 basic types			
Gly	63.8	Radii Set	ProtOr Radii, Tsai et al. (1999)			
Ala	89.3	Plane-Positioning Method	Ratio			
Val	138.2					
Leu	163.1		Atom Selection Criteria			
Ile	163.0		BL+			
Met	165.8	Structure Set	SCOP (87 structures)			
Pro	121.6					
His	157.5					
Phe	190.8					
Tyr	194.6					
Trp	226.4					
Cyh	112.8					
Cys	102.5					
Ser	94.2					
Thr	119.6					
Asn	112.4					
Gln	146.9					
Asp	114.4					
Glu	138.8					
Lys	165.1					
Arg	190.3					

ProtOr

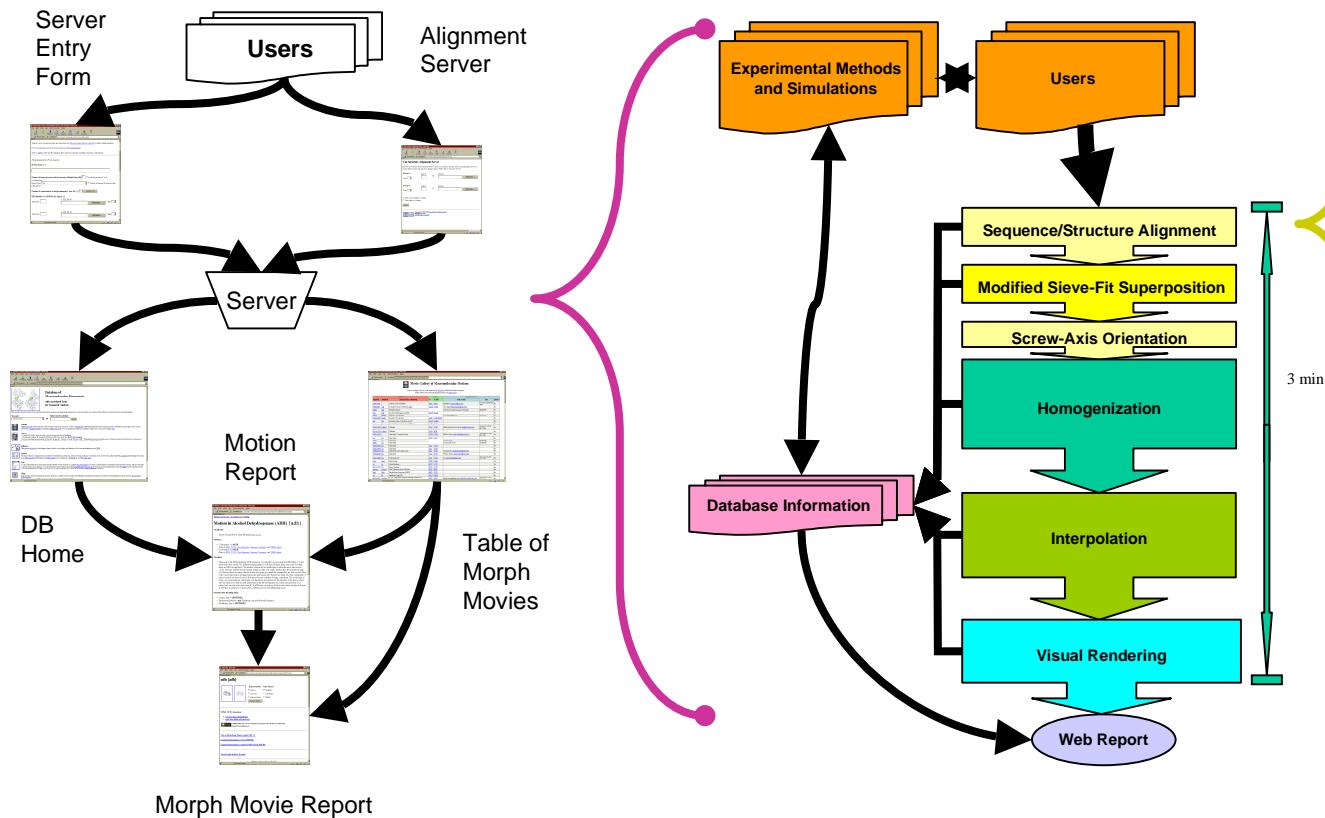
Parameter Set:

Standard Radii &

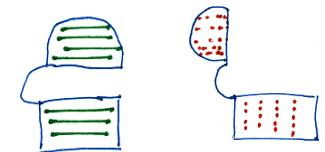
Volumes

- Consistent Radii, Typing, and Volumes for Packing Calculations
- For comparison

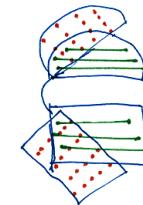
Motion Analysis Server, Morph Movies



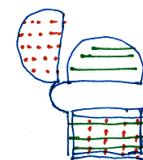
Struc-1 Struc-2



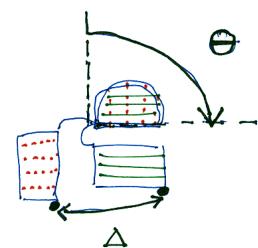
Overall Fit



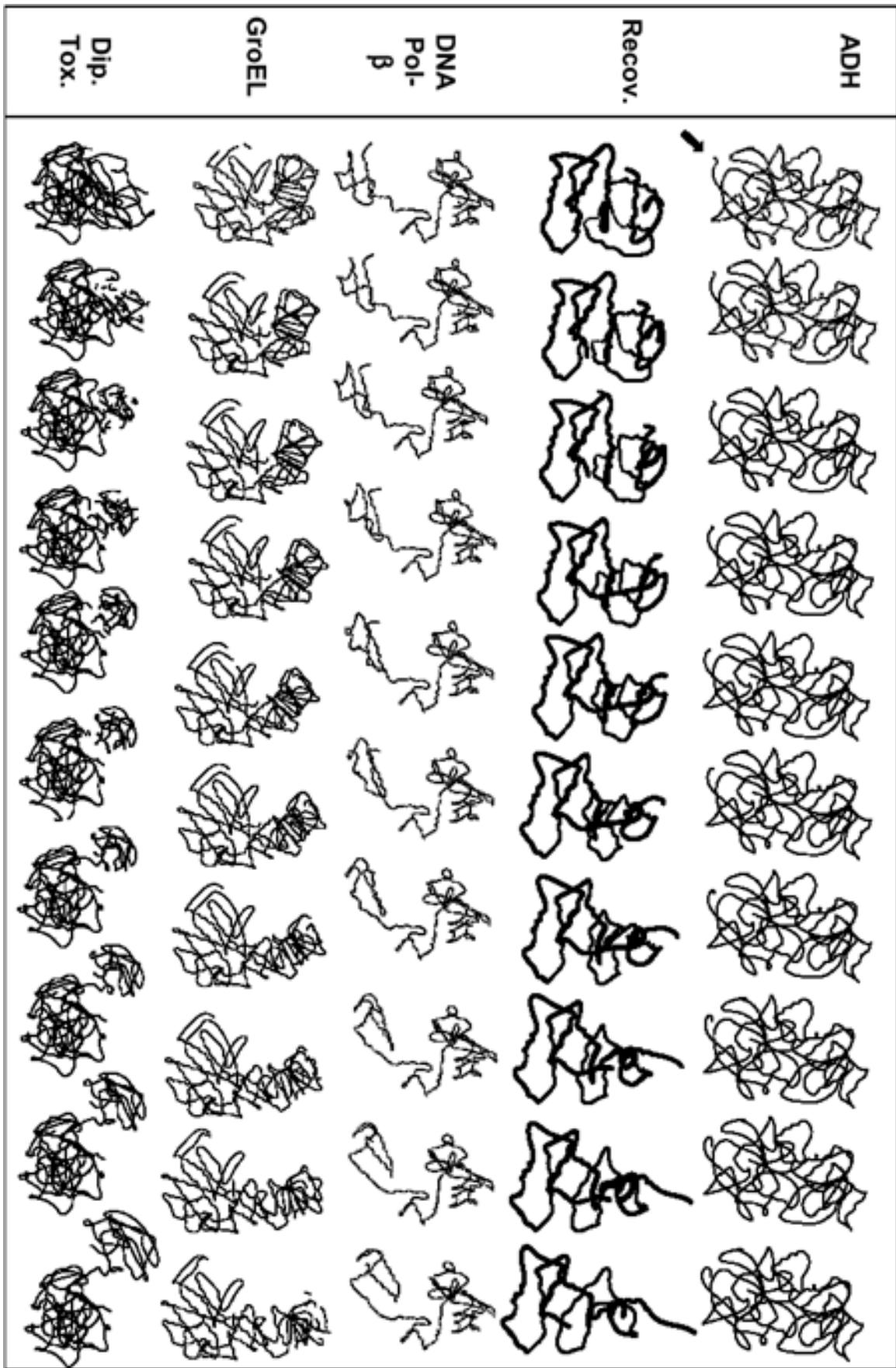
Core-1 Fit



Core-2 Fit



5 Movies Generated by the Server



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

bioinfo.mbb.yale.edu

Structures ("Classic")

(now) Structural Genomics

(now) Func. Genomics

Arrays (future)

Structures ("Classic")

1 Fold Library (A parts list.) Structural Alignment, EVD P-value, Seq. Struc. diverg.

2 Folds in Genomes (Shared, common, and/or unique parts?) Known Folds. Fold Tree, Top-10. $\beta\alpha\beta$. Biases. MG fold assignment extent. MG Target Selection, MT retrospective decision tree.

3 Folds & Functions (Roles/part?) How many folds /function? Mostly 1, but TIM versatile. Seq. diverg. vs. Func. diverg.

4 Folds in the Transcriptome

(Common parts? Where are parts?)

Enriched ↑ : VGA, TIM, $\alpha\beta$ folds, energy, synthesis. Depleted ↓ : NS, very long, TM folds, transport, transcription, Leu-zip fold. Bayesian Localizer, phenotypes clustering

5 Fold Flexibility (How adaptable is a part?). Motions DB, morph server, interface packing, Voronoi Volumes

W Krebs, J Tsai, M Levitt, C Wilson,
Y Kluger, S Balasubramanian,
L Regan, R Das, C Arrowsmith, A
Edwards, H Hegyi, J Lin, A Drawid,
D Greenbaum, M Snyder, R Jansen

