

# Structural Genomics:

# Surveys of a Finite Parts List

Mark Gerstein

*H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison*

Talk at DIMACS, 15 Dec. 2000

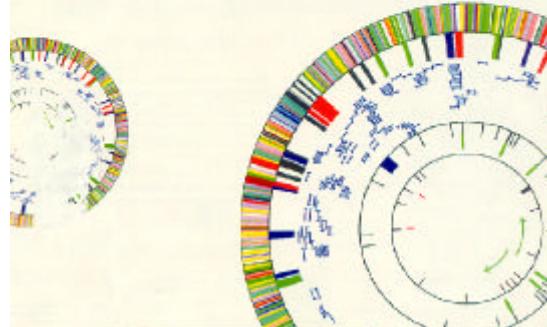
**1995**

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



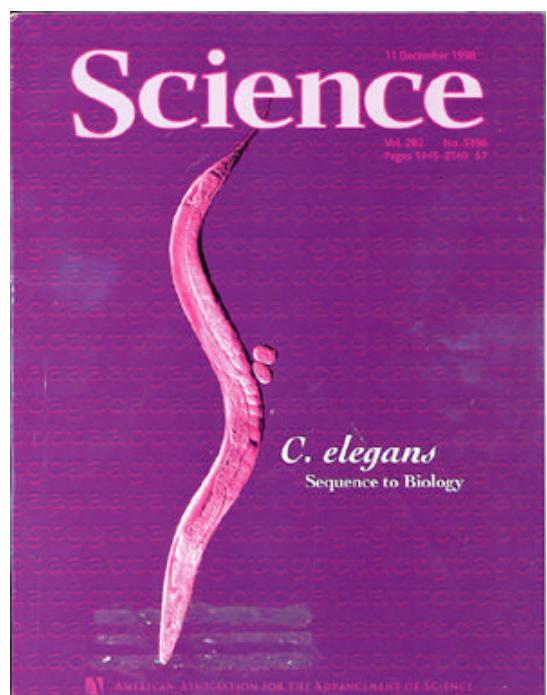
**1997**

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



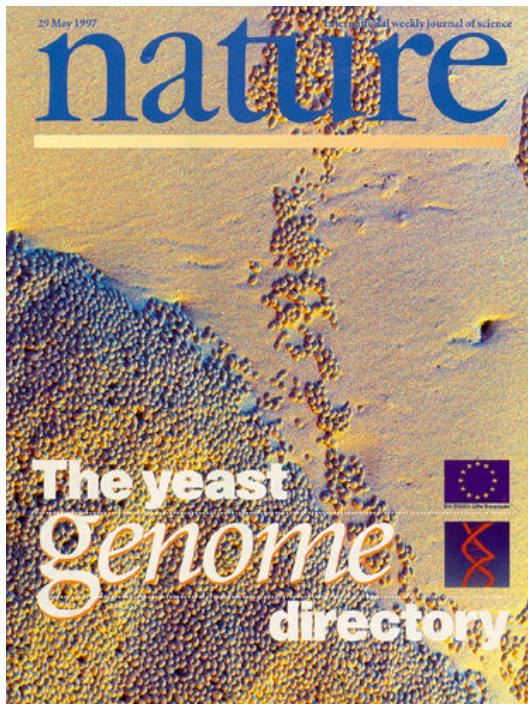
**1998**

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



**2000?**

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



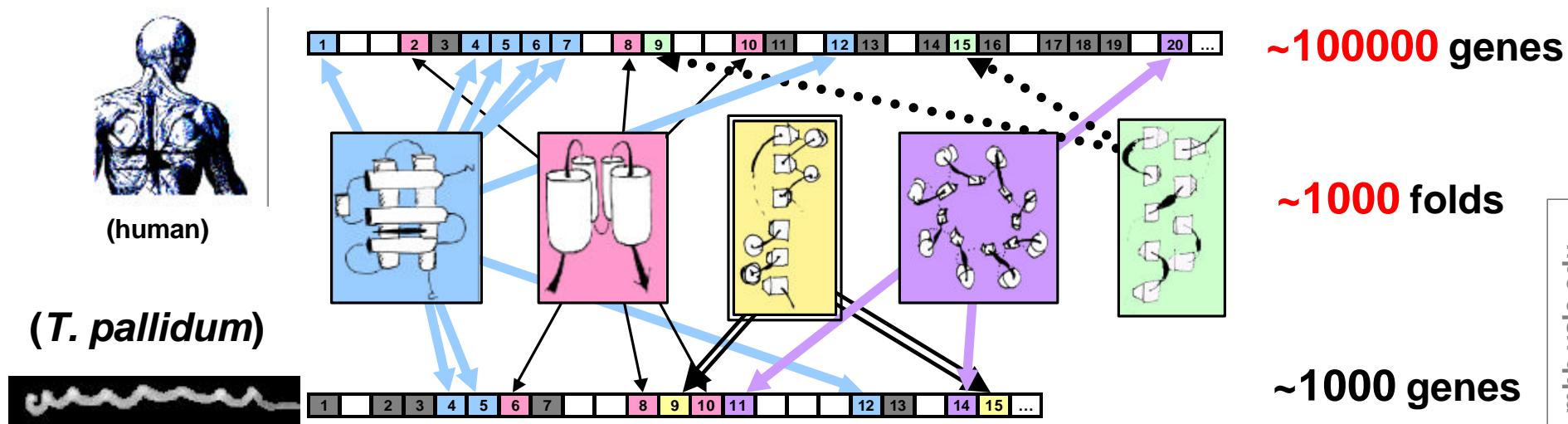
Genomes  
highlight  
the  
**Finiteness**  
of the  
**"Parts"** in  
Biology



'98 spoof

real thing, Apr '00

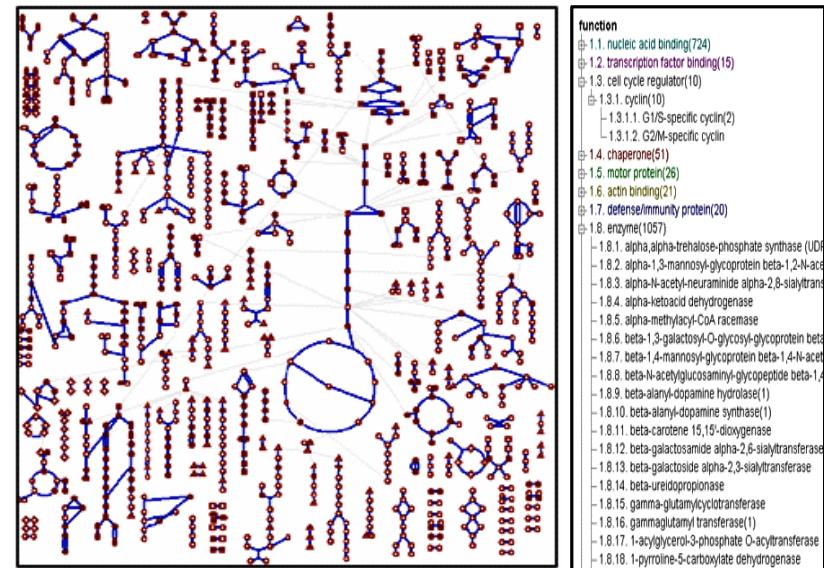
# World of Structures is even more Finite, providing a valuable simplification



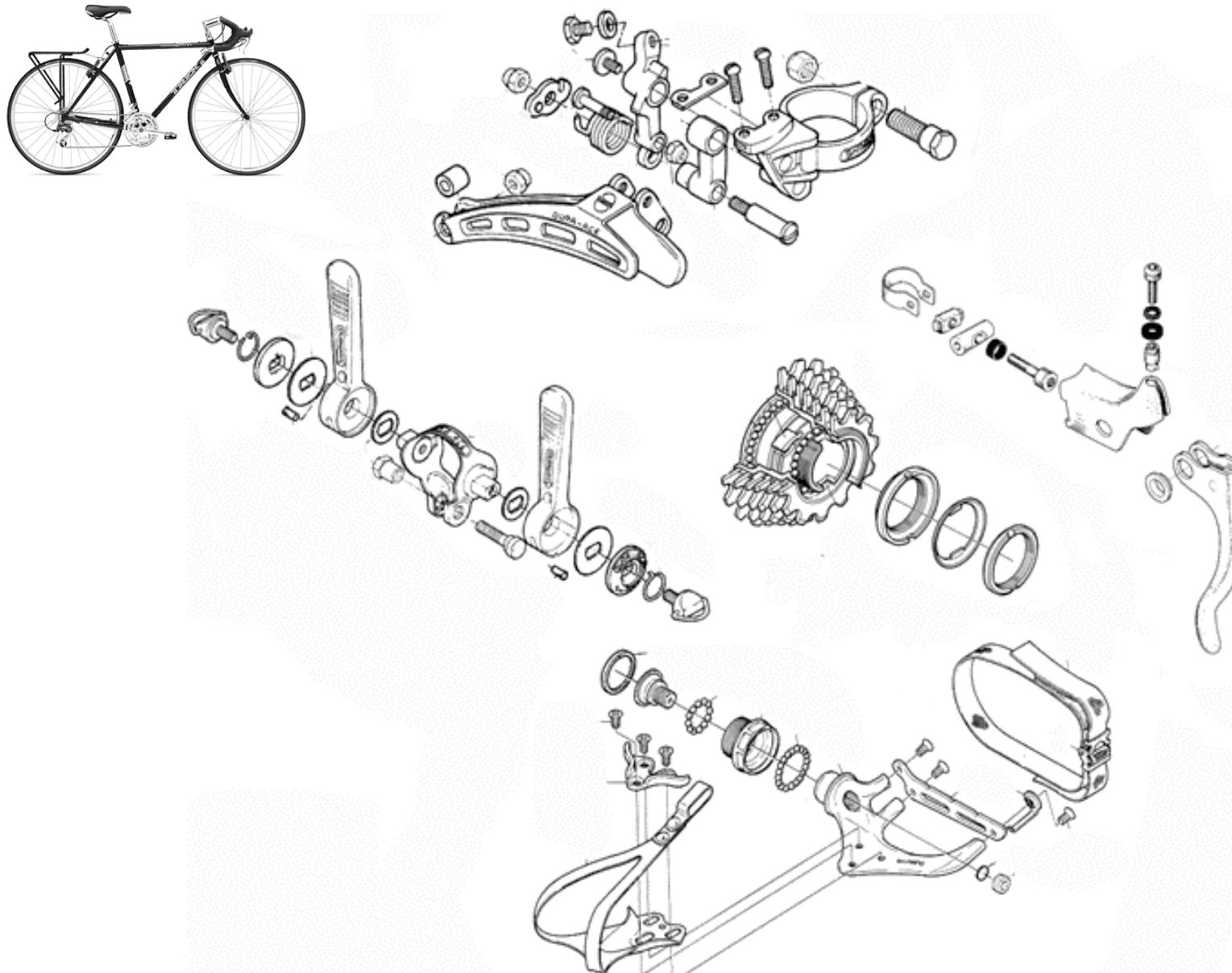
Same logic for pathways, functions, sequence families, blocks, motifs....

**Global Surveys of a  
Finite Set of Parts from  
Many Perspectives**

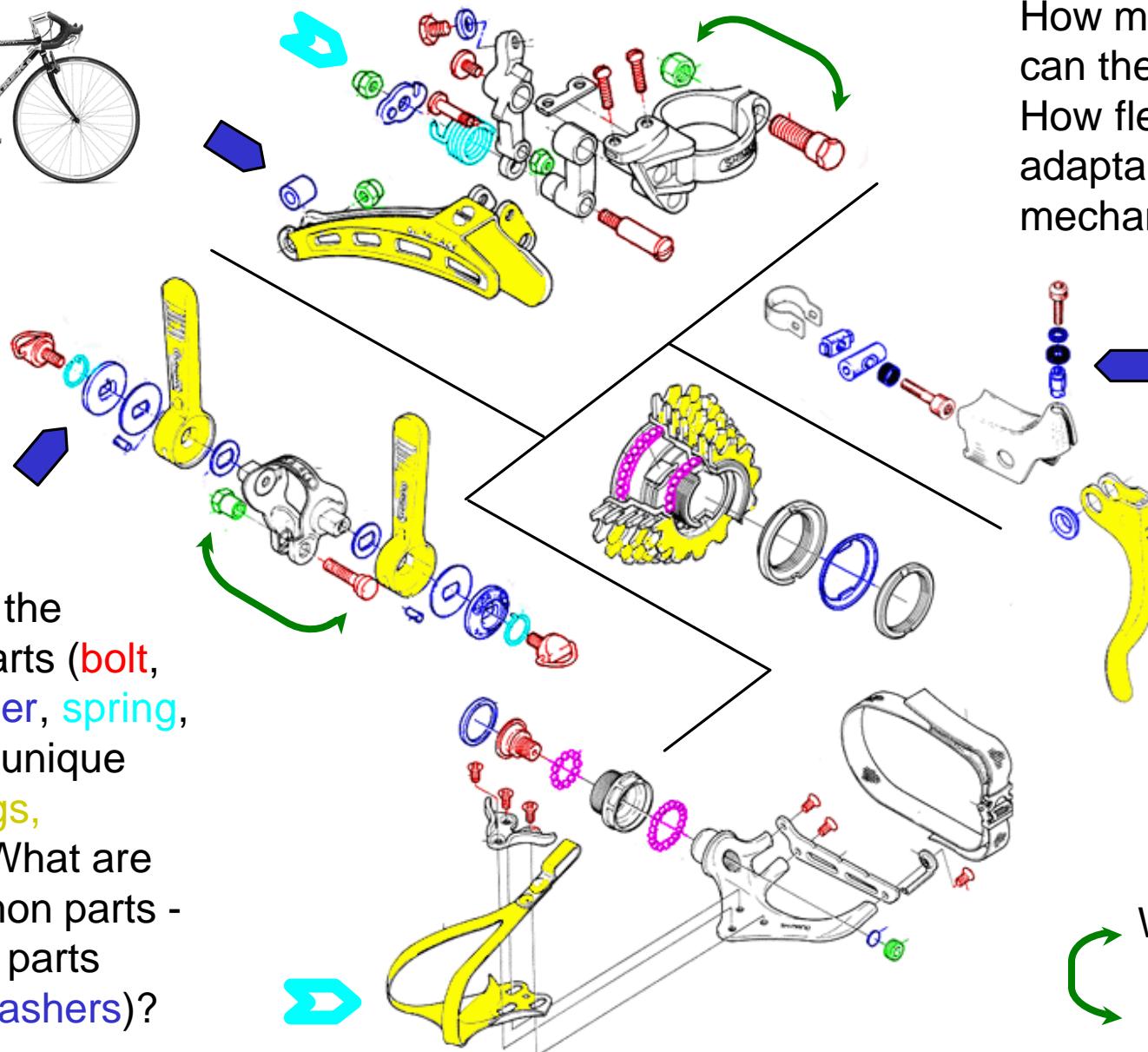
Functions picture from [www.fruitfly.org/~suzi](http://www.fruitfly.org/~suzi) (Ashburner); Pathways picture from, [ecocyc.pangeasystems.com/ecocyc](http://ecocyc.pangeasystems.com/ecocyc) (Karp, Riley). Related resources: COGS, ProDom, Pfam, Blocks, Domo, WIT, CATH, Scop....



# 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? Which parts interact?

# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

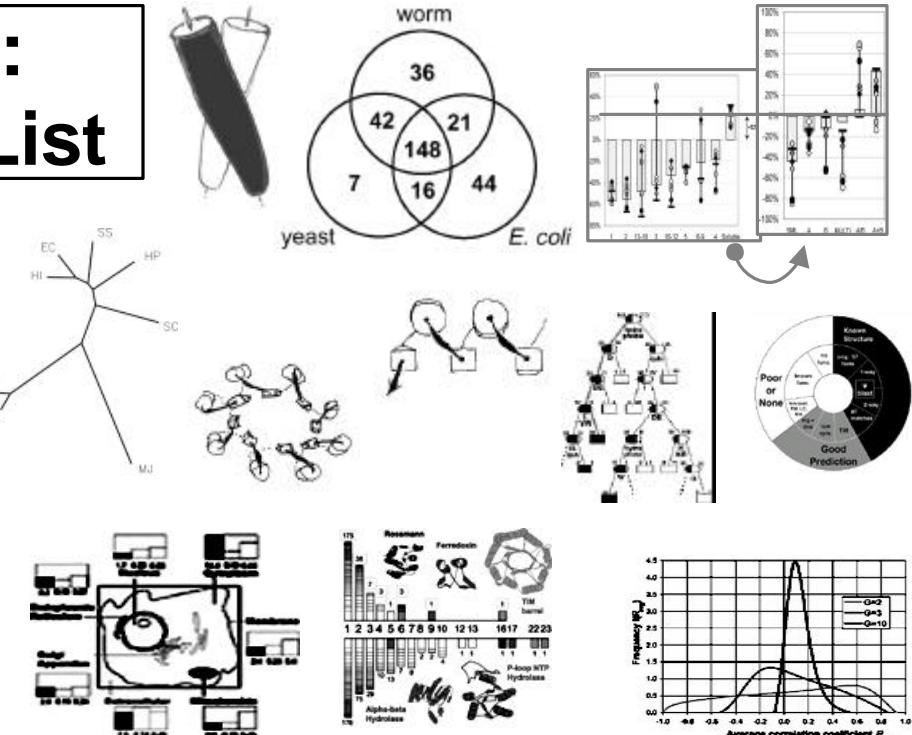
**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

**2 Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

**3 Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

**4 Using Folds to Interpret Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

**5 Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.

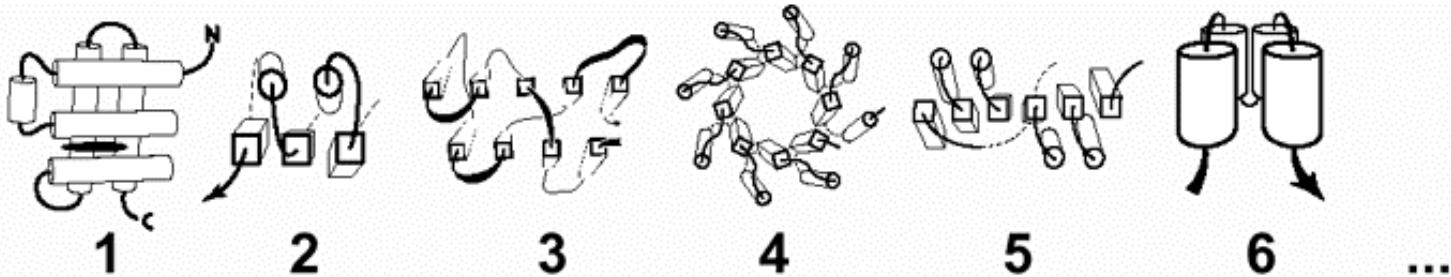


**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

**bioinfo.mbb.yale.edu**

# 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	0	8°C
F	9.65	0	4 C/mol
$\epsilon_0$	8.85	-12	F/m
$\mu_0$	1.26	0	6 H/m
$\hbar$	6.63	0	34 J·s
k	1.38	0	23 J/K
$m_e$	9.11	0	-31 kg
$m_p$	1.67	0	-27 kg
$m_n$	1.68	0	-27 kg
$a_0$	5.29	0	-11 m
$\lambda_C$	2.43	0	-12 m
c	3.00	0	-19 m/s
G	6.67	0	-11 m <sup>2</sup> /kg·s <sup>2</sup>
N <sub>A</sub>	6.02	0	23 mol <sup>-1</sup>

10

Physics

100

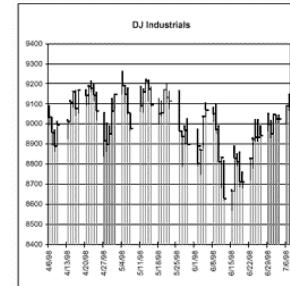
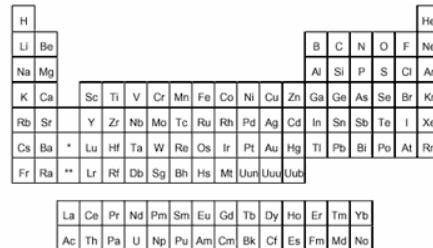
Chemistry

1000  
-10000

Finance

>1000000

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.

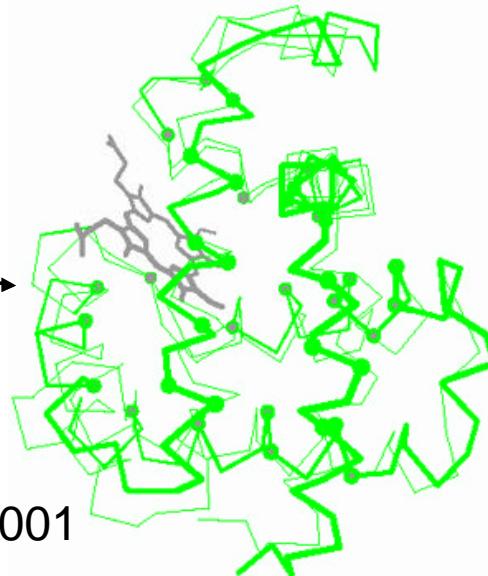
Hb

# The Parts List: A Library of Known Folds



Alignment  
of Individual  
Structures

P<.000001

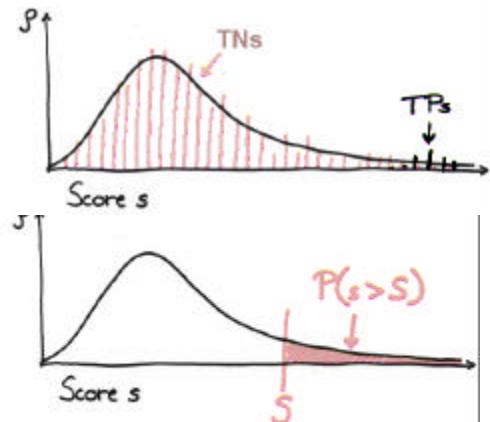


Fusing into a  
Single **Core**  
Structure  
Template

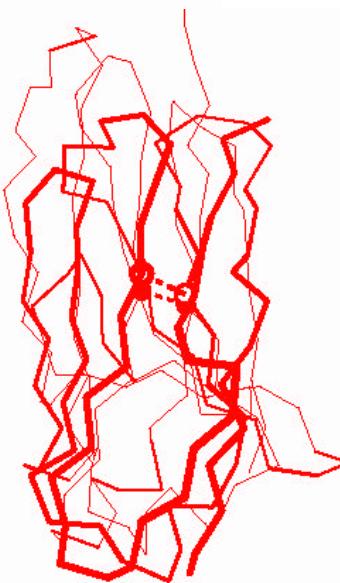


Mb

Statistics  
to Establish  
Relation-  
ships  
(P-values)



P<.001



P~1



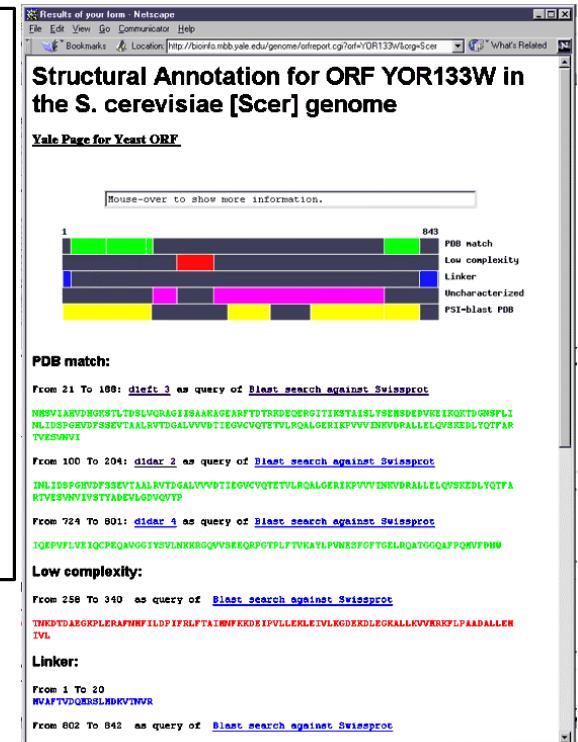
# Integrated Analysis System: X-ref Parts with Genomes

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

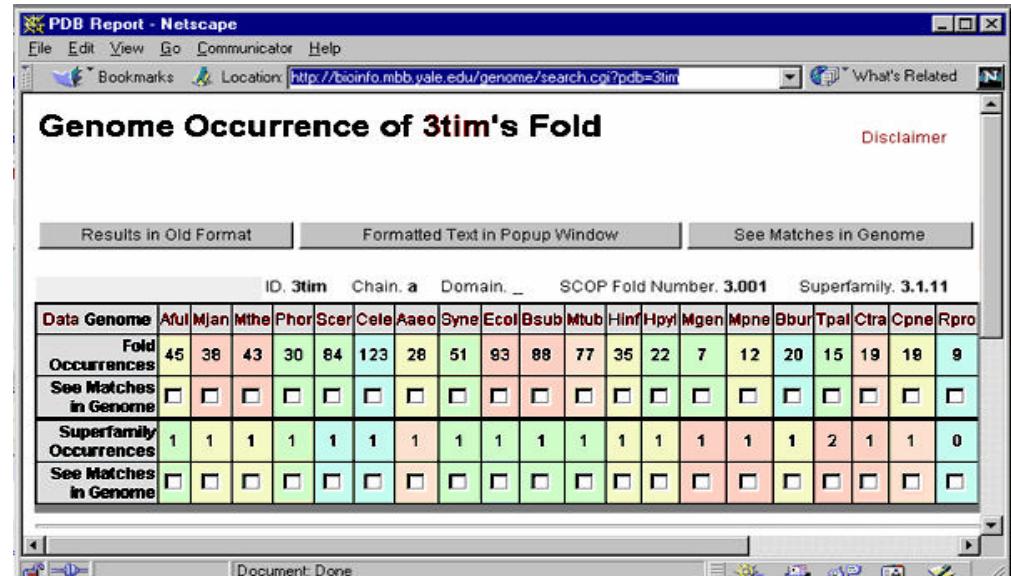
Remington, Matthews '80; **Taylor, Orengo '89, '94; Thornton, CATH**; Artymiuk, Rice, Willett '89; **Sali, Blundell, '90; Vriend, Sander '91; Russell, Barton '92; Holm, Sander '93+ (FSSP)**; Godzik, Skolnick '94; **Gibrat, Bryant '96 (VAST)**; F Cohen, '96; Feng, Sippl '96; G Cohen '97; Singh & Brutlag,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

Folds: scop+automatic  
Orthologs: COGs  
“Families”: homebrew,  
ProtoMap

finding parts in genome sequences  
**blast**,  
 $\psi$ -blast,  
**fasta**,  
TM, low-complexity, &c  
(Altschul, Pearson, Wooton)

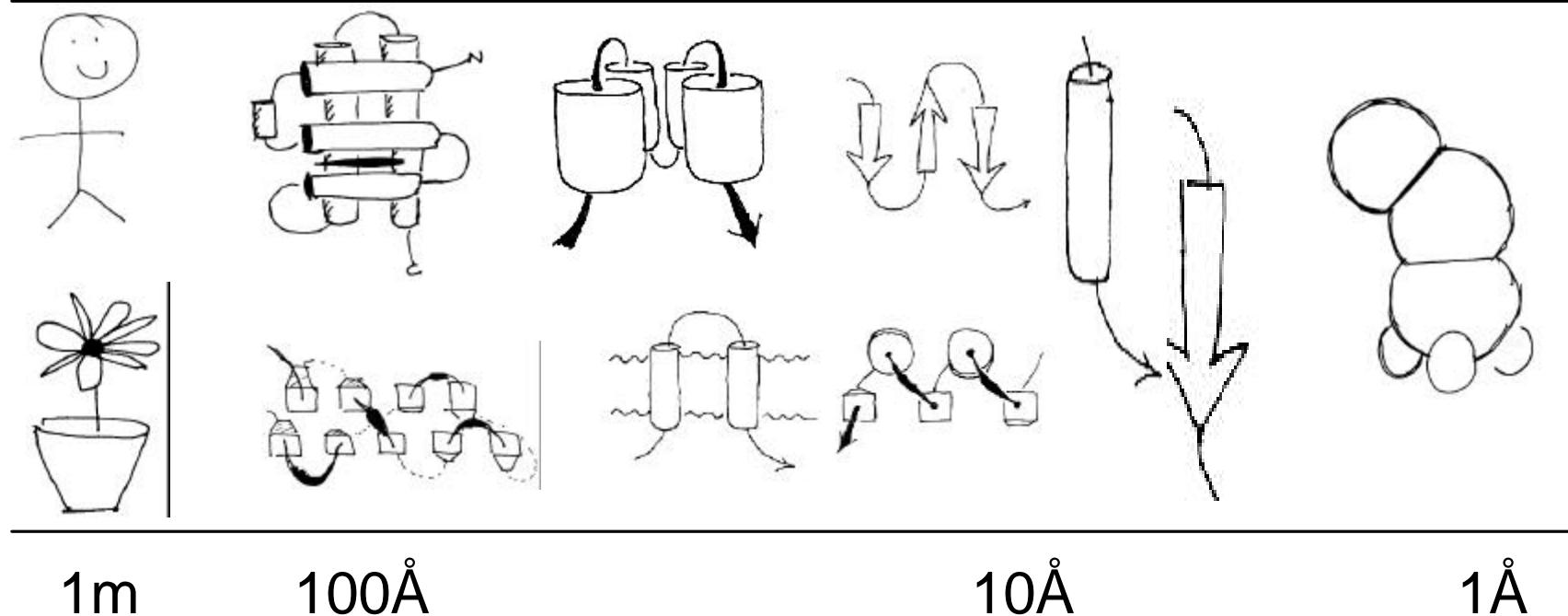


part occurrence profiles

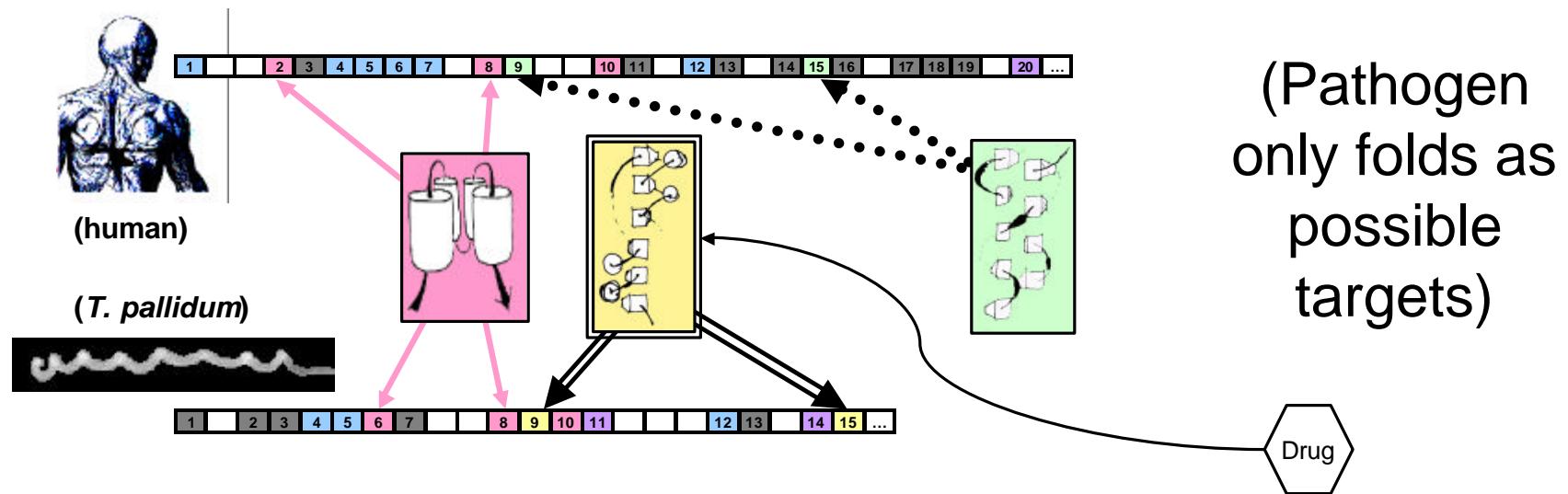


# At What Structural Resolution Are Organisms Different?

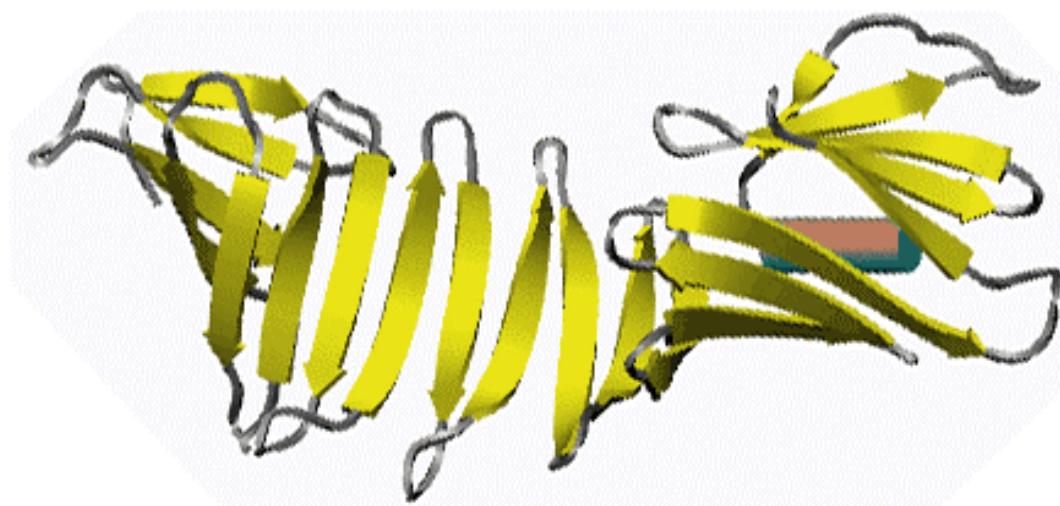
person	protein	super-secondary	helix	individual
plant	fold (Ig)	structure ( $\beta\beta$ ,TM– TM, $\alpha\beta\alpha\beta,\alpha\alpha\alpha$ )	strand	atom (C,H,O...)



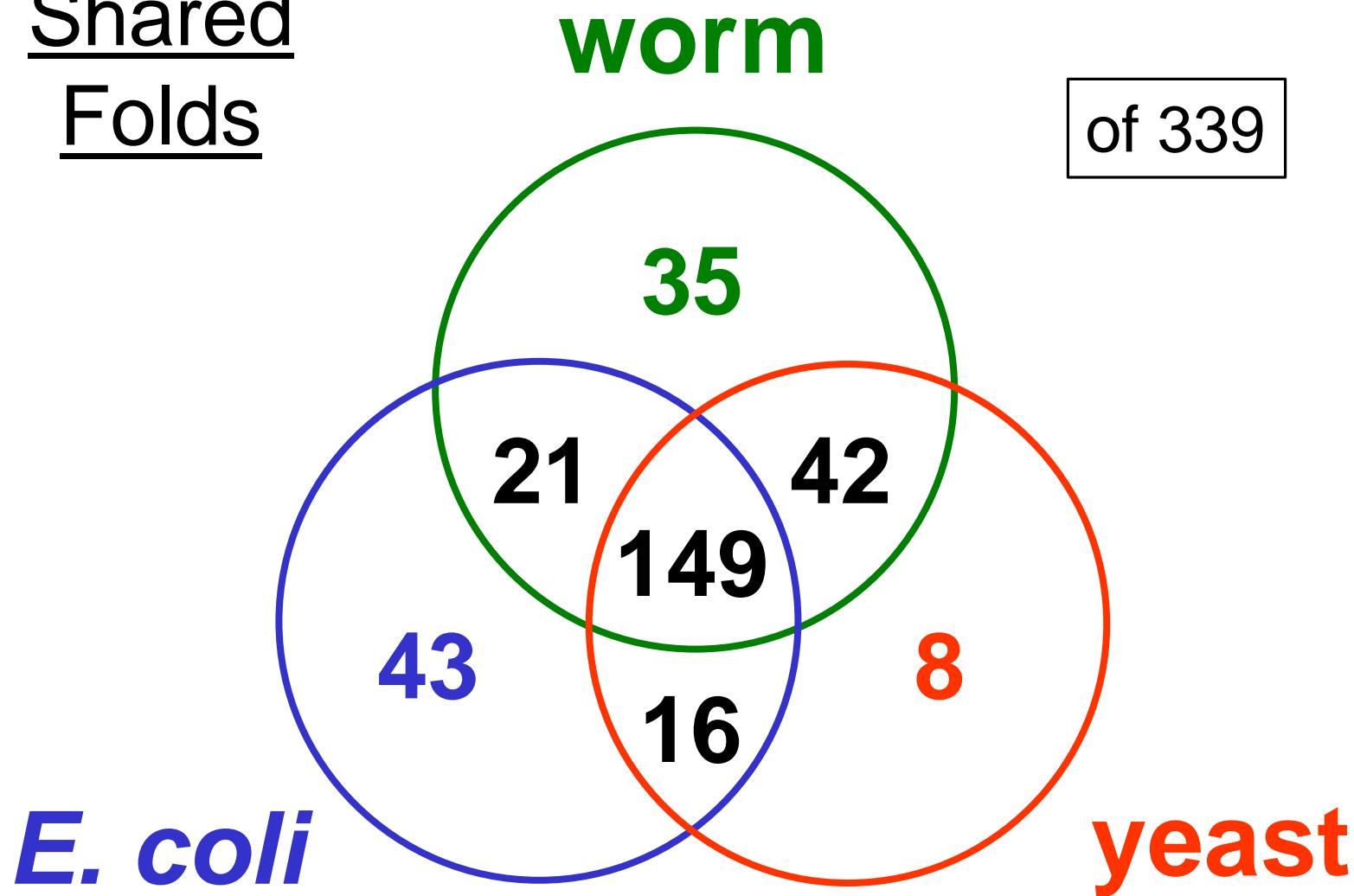
# Practical Relevance of Structural Genomics



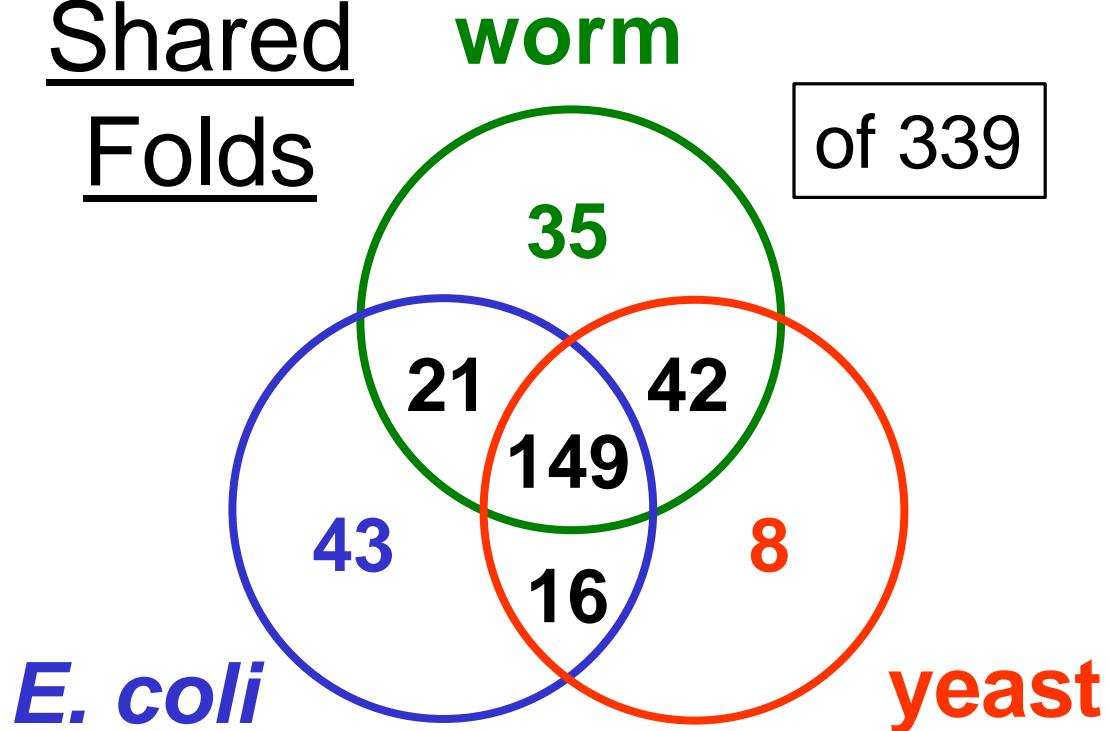
- 
- OspA protein
    - ◊ in Lyme-disease spirochete *B. burgdorferi*
    - ◊ previously identified as the antigen for vaccine
    - ◊ has novel fold (C Lawson)



## Shared Folds



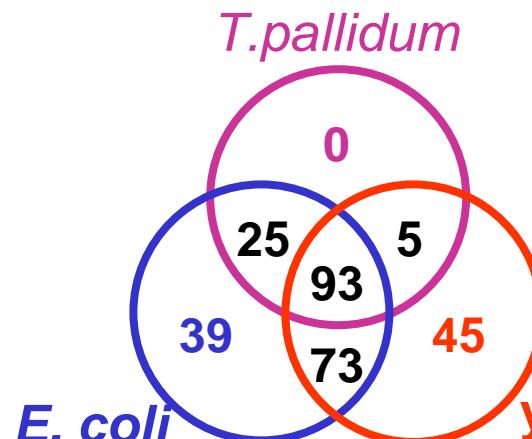
# Shared Folds



*E. coli*

*yeast*

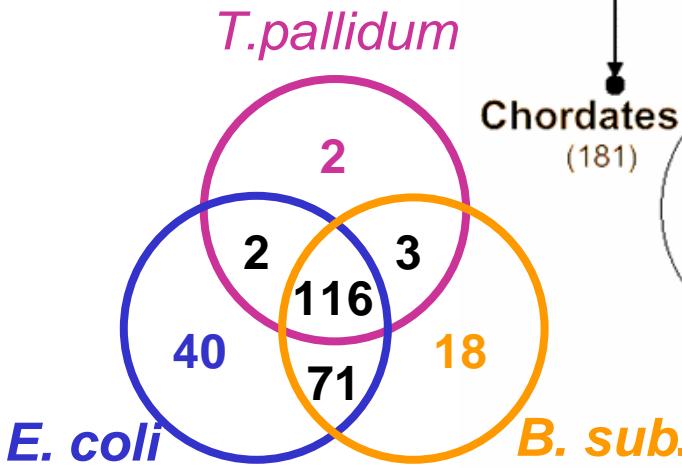
*T. pallidum*



*E. coli*

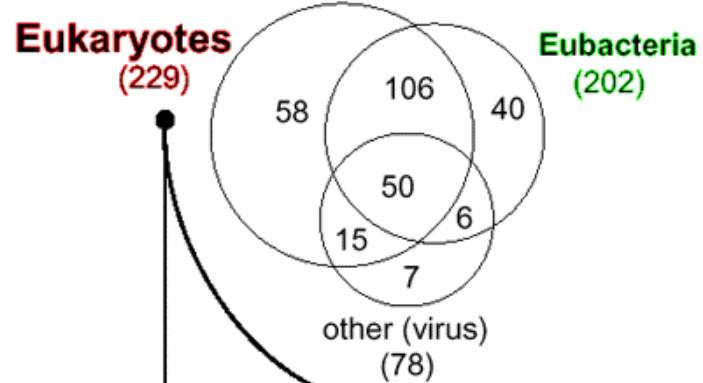
*yeast*

*T. pallidum*



*E. coli*

*B. sub.*

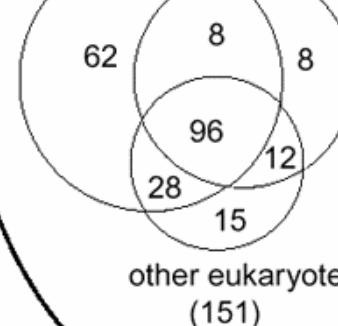


other (virus)  
(78)

Metazoa

(194)

Plants  
(124)

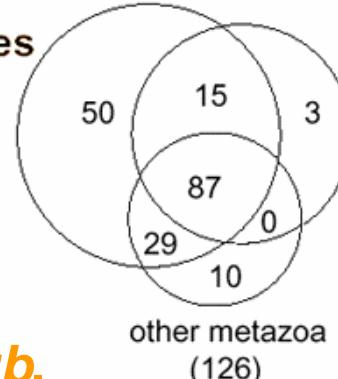


other eukaryotes  
(151)

Chordates

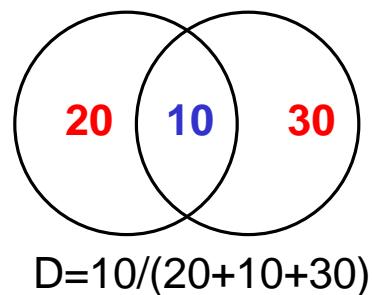
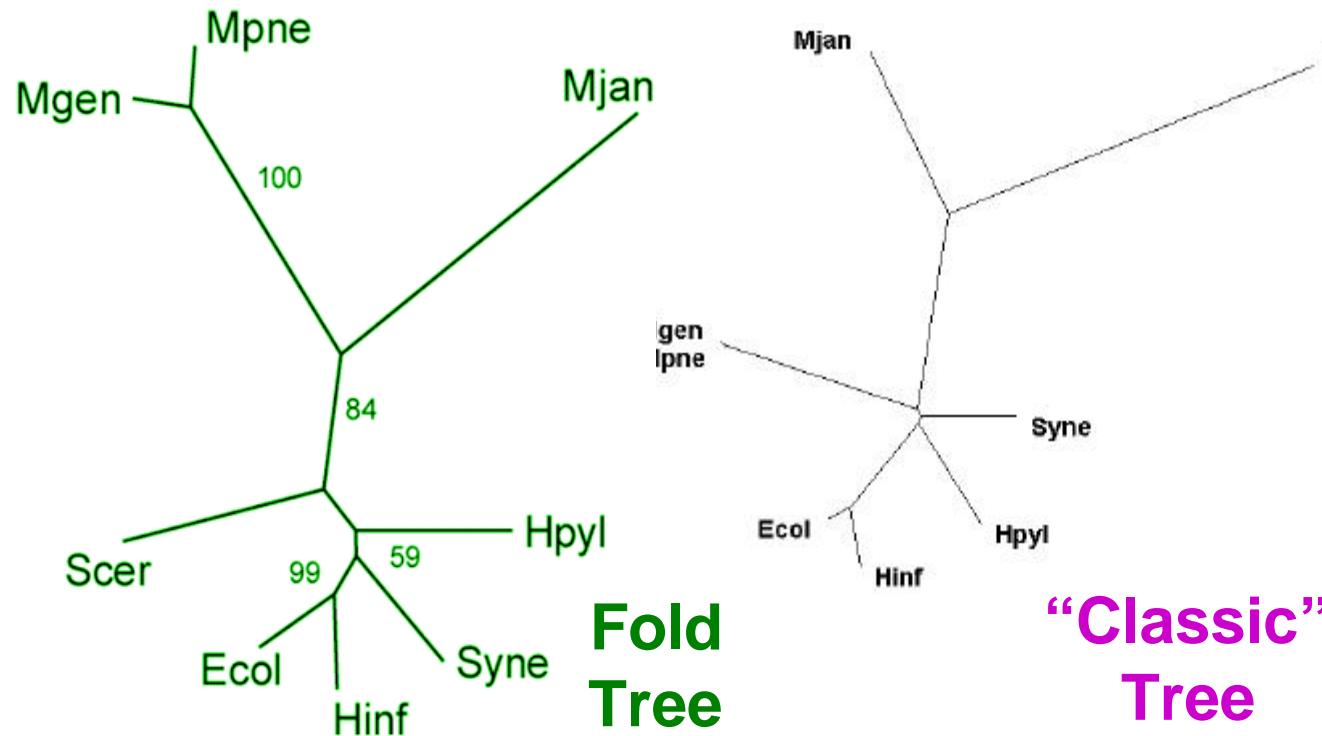
(181)

Arthropods  
(105)



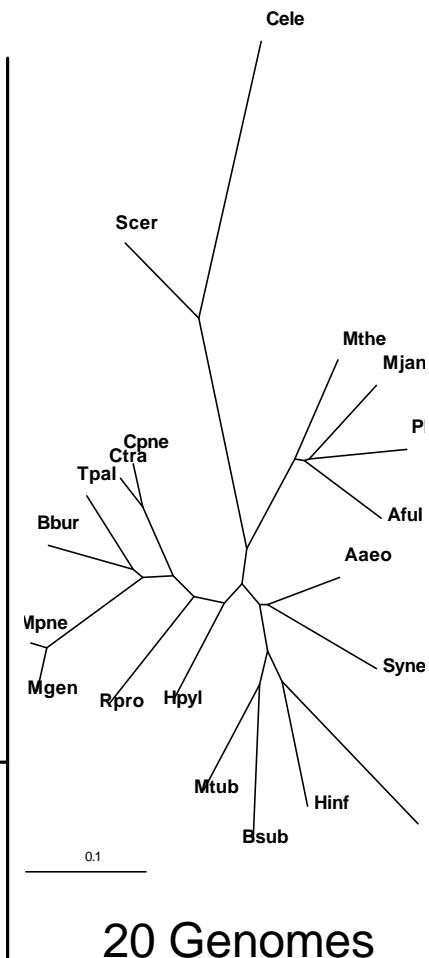
other metazoa  
(126)

# Cluster Trees Grouping Initial Genomes on Basis of Shared Folds

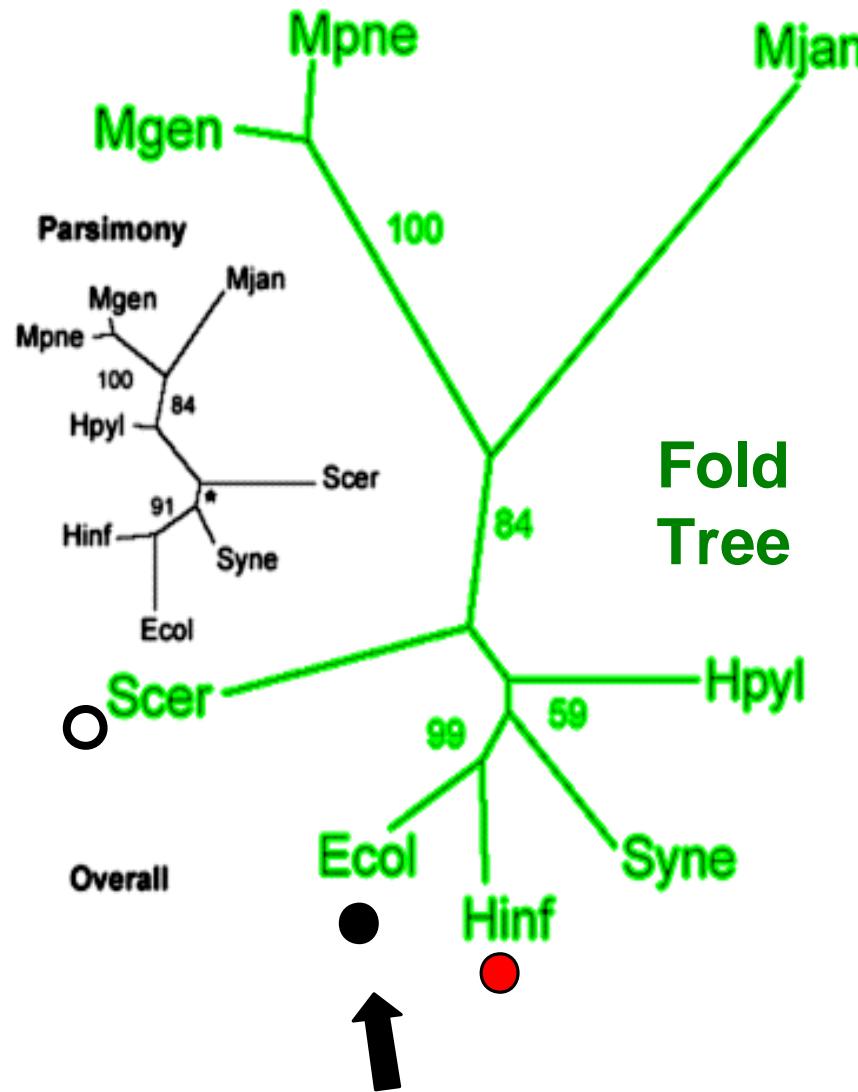


$D = S/T$        $S = \# \text{ shared folds}$   
 $D = \text{shared fold dist.}$        $T = \text{total \#}$   
 betw. 2 genomes      folds in both

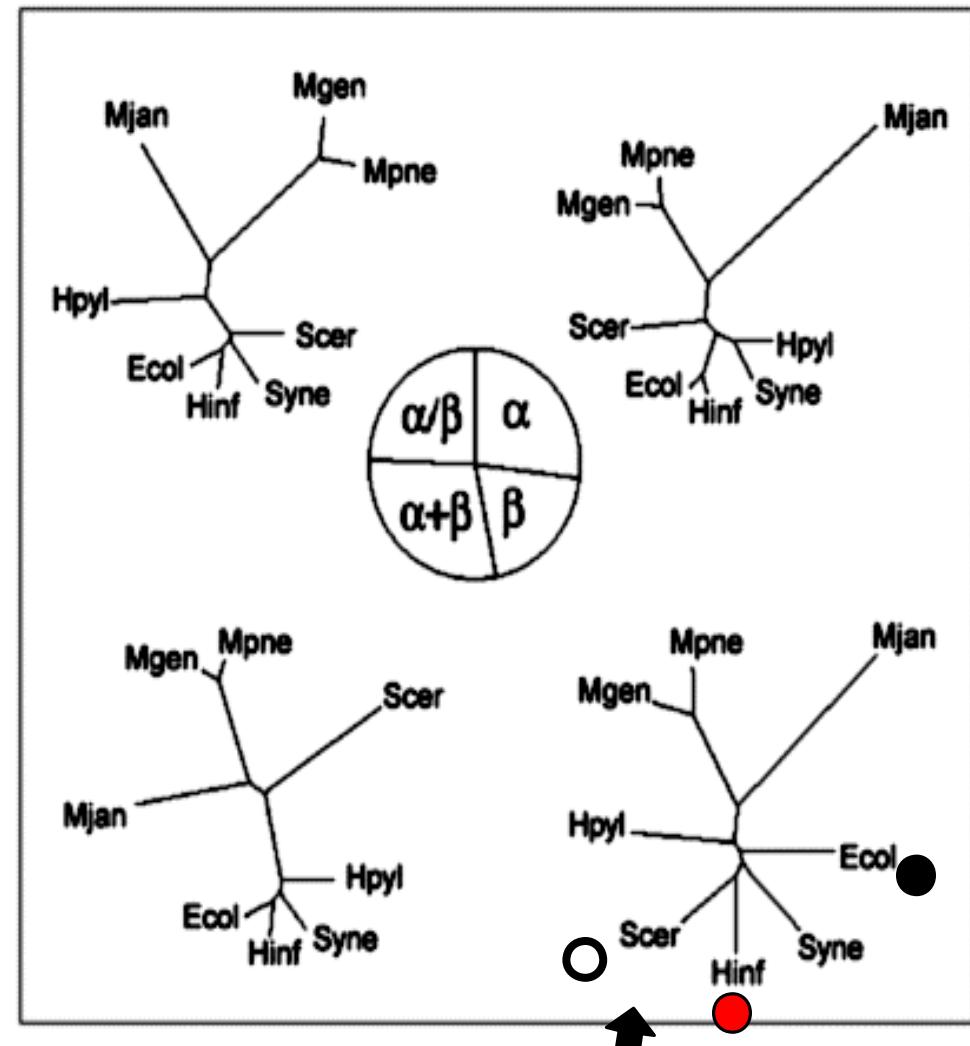
“Classic”  
Tree



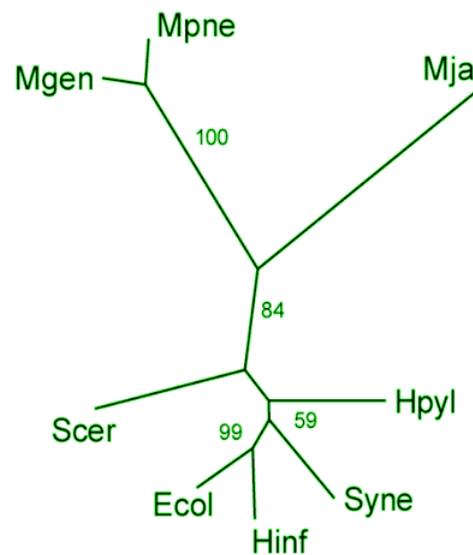
# Distribution of Folds in Various Classes



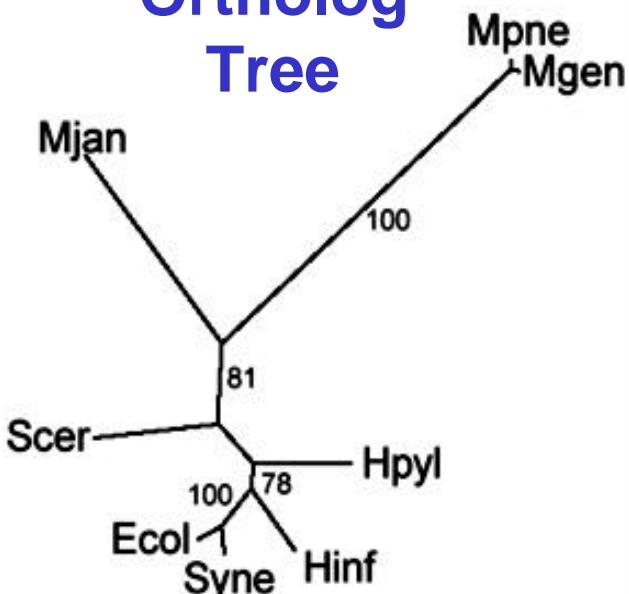
Unusual distribution of all-beta folds



## Fold Tree

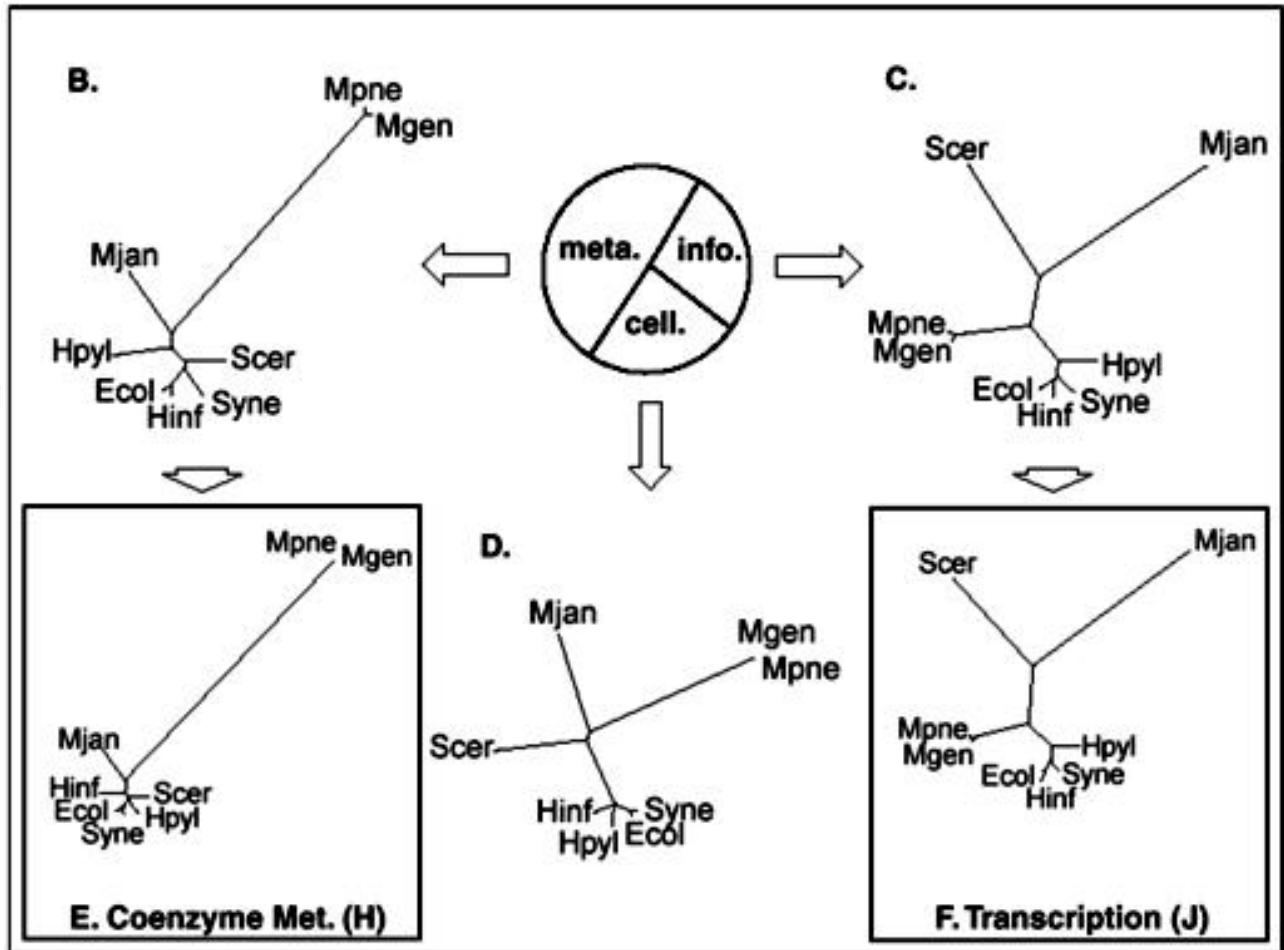


## Ortholog Tree

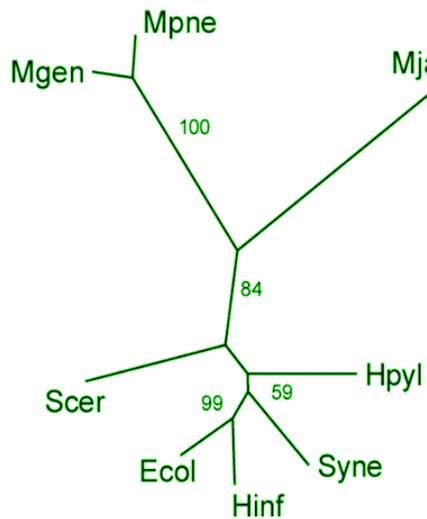


## Compare with Ortholog Occurrence Trees, another “partial- proteome” tree

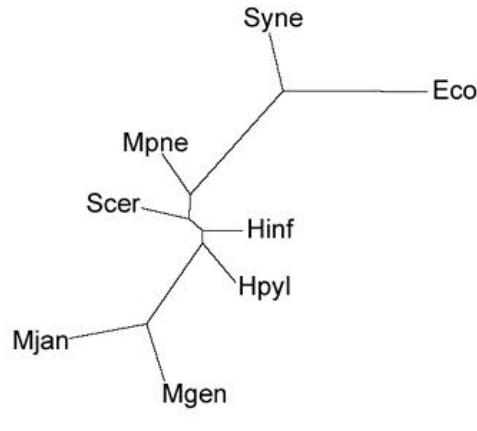
(based on COGs scheme of Koonin & Lipman, similar approaches by Dujon, Bork, &c.)



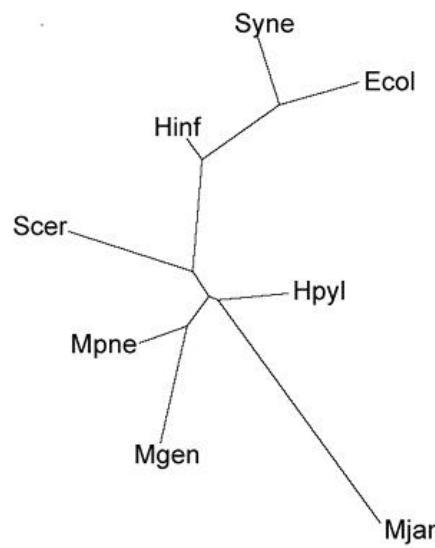
# Compare with trees on spectrum of “levels”: single-gene trees, whole-genome composition trees



Fold  
Tree

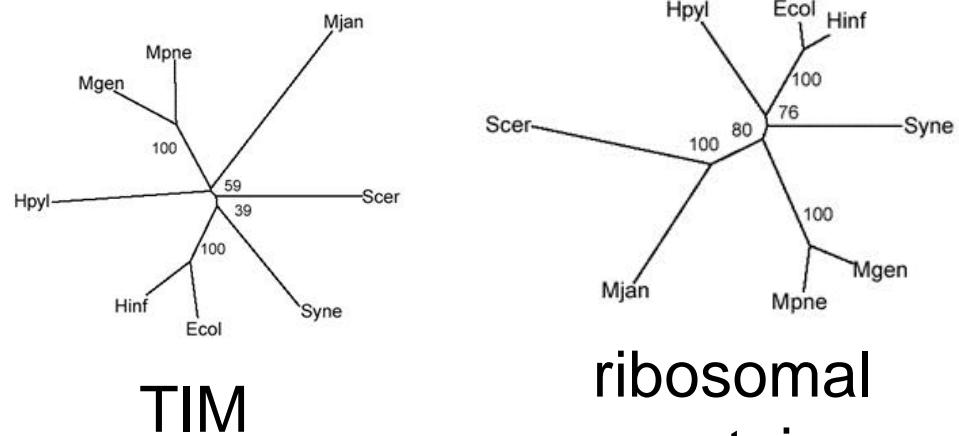


A. Di-Nucleotide



B. Amino Acid

AA & di-NT Composition Trees  
(S Karlin)



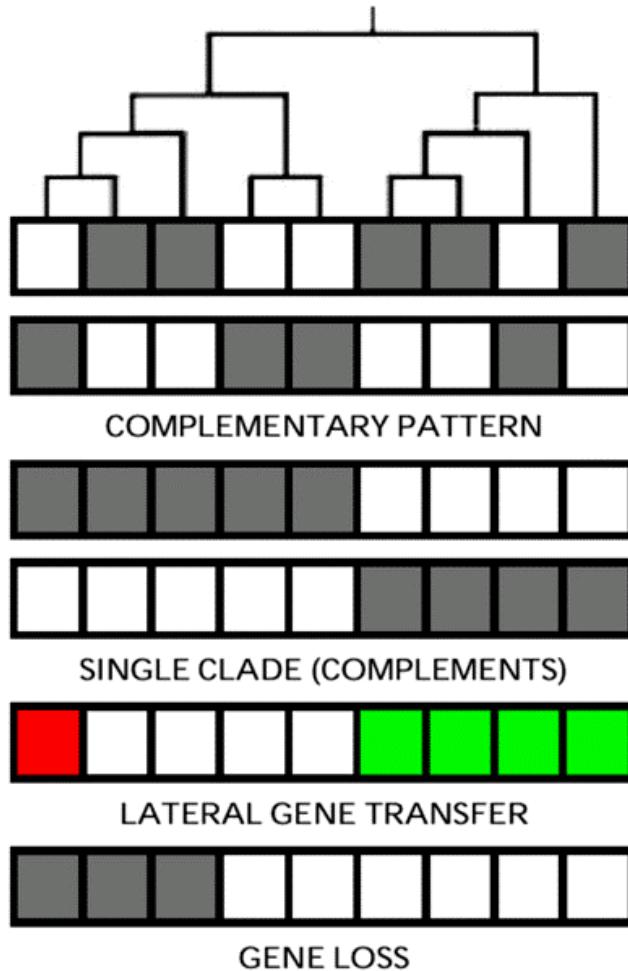
TIM

ribosomal  
protein

## Single-gene Trees

## Ortholog Tree

## “Classic” Tree

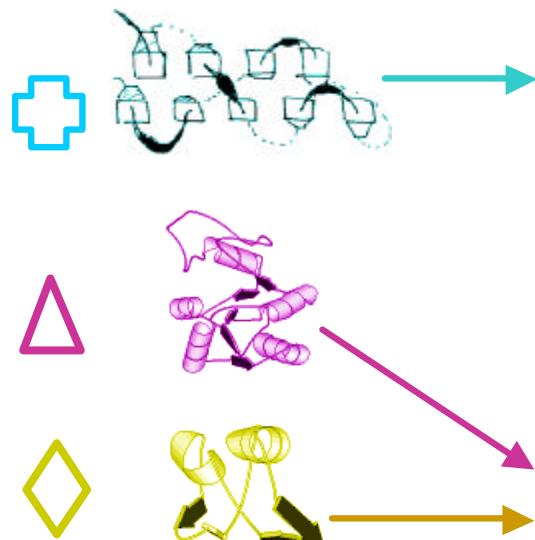


# Specific Patterns of Fold Usage: Horizontal Transfer

afu	mjan	mthe	phor	scer	cele	aaoe	syne	ecol	mtsub	mtub	hifn	hpyl	mgcn	mpne	bbur	tpal	ctra	cpne	ppro	Sfam	domain	SCOP Function	Swissprot	Swissprot Function
●	●	●	●																	1.86.1	d1aora1	Aldehyde FerOR C' domain	AOR_PYRFU	ALDEHYDE:FERREDOXIN OXIDOREDUCTASE
●	●	●	●																	4.94.1	d1aora2	Aldehyde FerOR N' domain	AOR_PYRFU	ALDEHYDE:FERREDOXIN OXIDOREDUCTASE
●	●	●	●																	3.1.10	d5ruba1	RuBisCo, C' domain	RBL_NITVU	RUBISCO LARGE SUBUNIT
																				1.101.1	d5csma_	Chorismate mutase II	CHMU_ARATH	CHORISMATE MUTASE (EC 5.4.99.5)
																				1.37.1	d1rec_	EF-hand	TPC2_DROME	TROPONIN C
																				2.1.5	d1suh_	Cadherin	CAD5_HUMAN	VASCULAR ENDOTHELIAL-CADHERIN
																				2.45.1	d1eal_	Lipocalins	PGHD_HUMAN	PROSTAGLANDIN-H2
																				3.7.1	d2bnh_	Leucine-rich repeats	RINI_PIG	RIBONUCLEASE INHIBITOR
																				4.70.1	d1axx_	Cytochrome b5	NIA1_MAIZE	NITRATE REDUCTASE (EC 1.6.6.1)
																				4.112.1	d1t0h_	Tyrosine hydroxylase	TY3H_HUMAN	TYROSINE 3-HYDROXYLASE (EC 1.14.16.2)

# 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?
	class				
Ig	B	830	1.7%		
Knottins	SML	565	1.1%		
Protein kinases (cat. core)	MULT	472	0.9%		
C-type lectin-like	A+B	322	0.6%		
corticoid recep. (DNA-bind dom.)	SML	276	0.5%		
Ligand-bind dom. nuc. receptor	A	257	0.5%		
alpha-alpha superhelix	A	247	0.5%		
C2H2 Zn finger	SML	239	0.5%		
P-loop NTP Hydrolase	A/B	235	0.5%		
Ferredoxin	A+B	207	0.4%		

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

Eubacteria

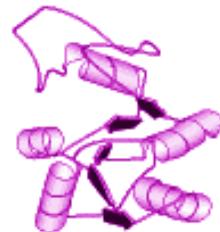
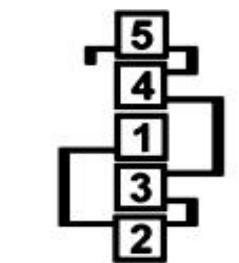
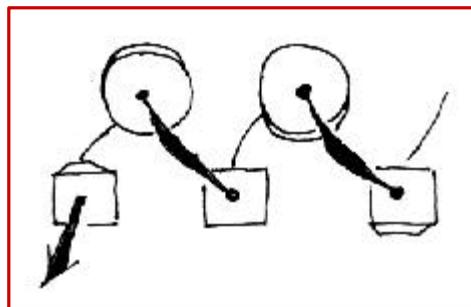
Rank	<i>M. thermo-autotrophicum</i>		<i>A. fulgidus</i>			
	Superfamily	#	Superfamily	#		
1		P-loop hydrolase	93		P-loop hydrolase	118
2		Phosphate-binding barrel	54		Rossmann domain	104
3		Rossmann domains	53		Phosphate-binding barrel	56
4		Ferredoxins	48		Ferredoxins	49
5		SAM methyl-transferase	17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

Archaea

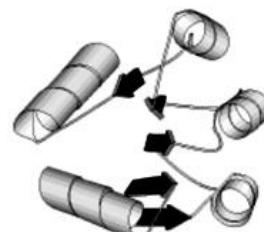
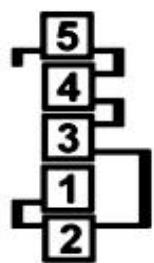
Rank	<i>S. cerevisiae</i>		
	Superfamily	#	
1		P-loop hydrolase	249
2		Protein kinase	123
3		Rossmann domain	90
4		RNA-binding domain	75
5		SAM methyl-transferase	63
6		Ribonuclease H-like	57
Total ORFs		6218	
with Common Superfamilies		560 (9%)	

Yeast

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



P-loop  
hydrolase



Flavodoxin  
like

42

ARTICLES

NATURE VOL. 316 3 NOVEMBER 1985

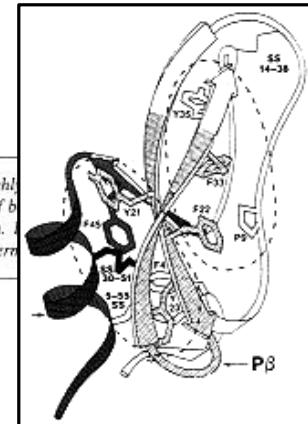
## A peptide model of a protein folding intermediate

Terrence G. Oas & Peter S. Kim

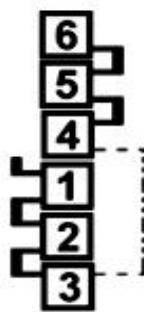
Whitehead Institute for Biomedical Research, Nine Cambridge Center, Cambridge, Massachusetts 02142, USA  
Department of Biology, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

*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 biotin inhibitor, contains secondary and tertiary structure similar to that found in the native protein. It circumvents the problem of cooperativity and permit characterization of structures of folding intermediates.*

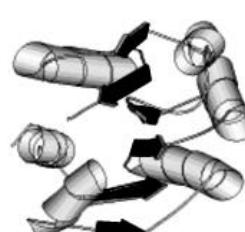
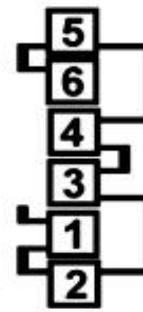
336: 42



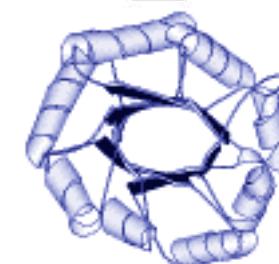
All share  $\alpha/\beta$  structure with repeated R.H.  $\beta\alpha\beta$  units connecting adjacent strands or nearly so (18+4+2 of 24)



Rossmann  
Fold



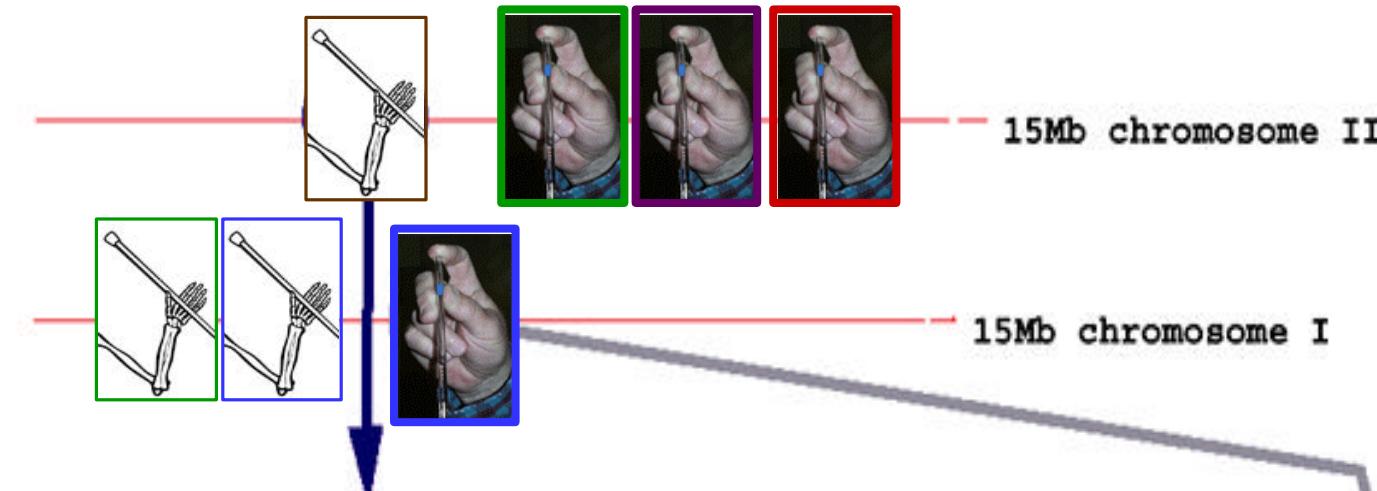
Thiamin  
Binding



TIM-  
barrel

HI, MJ, SC  
vs scop  
1.32

# Pseudogenomics: Surveying “Dead” Parts



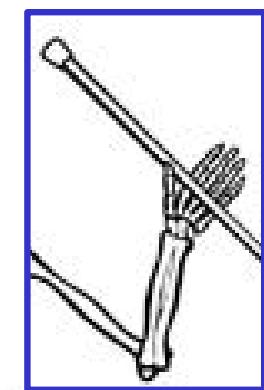
pseudogene fragment on worm chromosome II

TKRTSNGFGQDVVVDLFSILDGLVARAHXVLQDIFEFFAS  
KKMVTIFS#APHSPHSAPHYCAQFDNSAATVKV

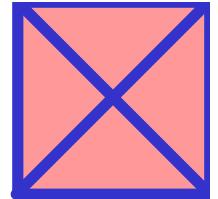
a paralog with the homologous segment highlighted (from chromosome I)  
(W09C3.6, serine/threonine protein phosphatase PP1)

M**TAPMDVDNLMSRLLNVGMSGGRLLT**SNEQELQTCCAVAKSVFASQASLLEVEPPIIVC  
GDIHGQYS**DLLRIFDKNGFPFDVNFLFLG**DYVDRGRQNIETICLMLCFKIKYPENFFMLR  
GNHECPAINRVYGFYEECNRRYKSTR**LWSIFQDTFNWMPLCG**LIGSRILCMHGGLSPHLQ  
TLDOLROLPRPQDPNPSIGIDLLWADPDOWVKGWOANTRGVSYVFGODVVADVCSRLDI  
**DLVARAHOVVODGYEFFASKMVTIFSAPHYCGOFDNSAATMKVDENMVCTFV**MYKPTPK  
SMRRG\*

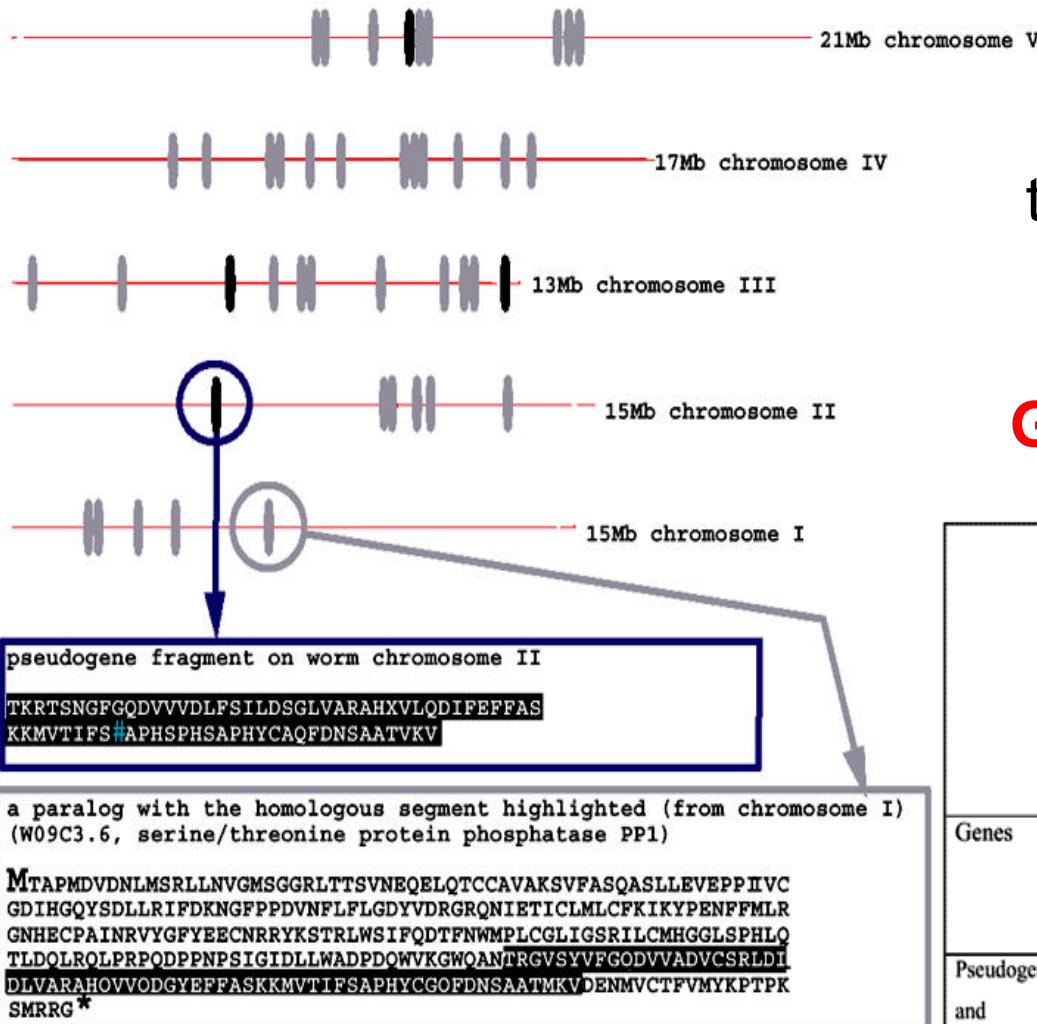
Example of  
a potential  
ΨG with  
frameshift in  
mid-domain



(Our def'n: ΨG = obvious homolog to known protein with frameshift or stop in mid-domain)



# Folds in Pseudogenes



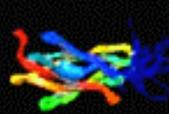
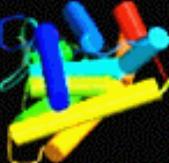
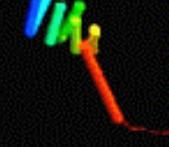
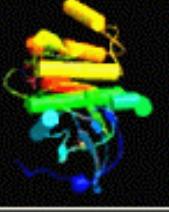
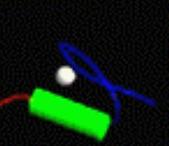
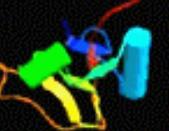
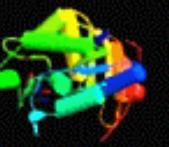
YG identification pipeline  
to Summary of Pseudogenes  
in worm

**G=19K G<sub>E</sub>=8K YG=4K (2K)**

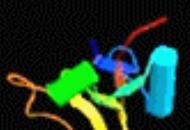
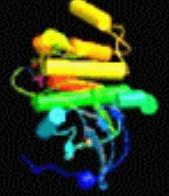
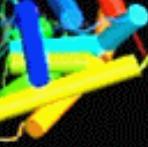
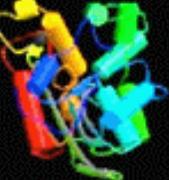
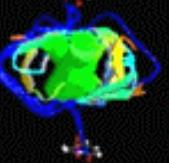
	Category	Total number	Number for genes with EST match	Genes with EST match as percentage of Category	Number for genes in paralog families with EST match	Genes in paralog families with EST match as percentage of Category
Genes	Total	18,576 (G)	7,829 (G <sub>E</sub> )	42%	13,417 (G <sub>P</sub> )	72%
	Singletons	5,913	2,788	47%	---	---
Pseudogenes and pseudogene fragments	Total	3,814 ( $\Psi$ G)	997 ( $\Psi$ G <sub>E</sub> )	26%	2,729 ( $\Psi$ G <sub>P</sub> )	72%
	Singletons	637 (17% of $\Psi$ G)	233	36%	---	---
	Intronic pseudogenes *	1,155 (30% of $\Psi$ G)	351	30%	704	61%

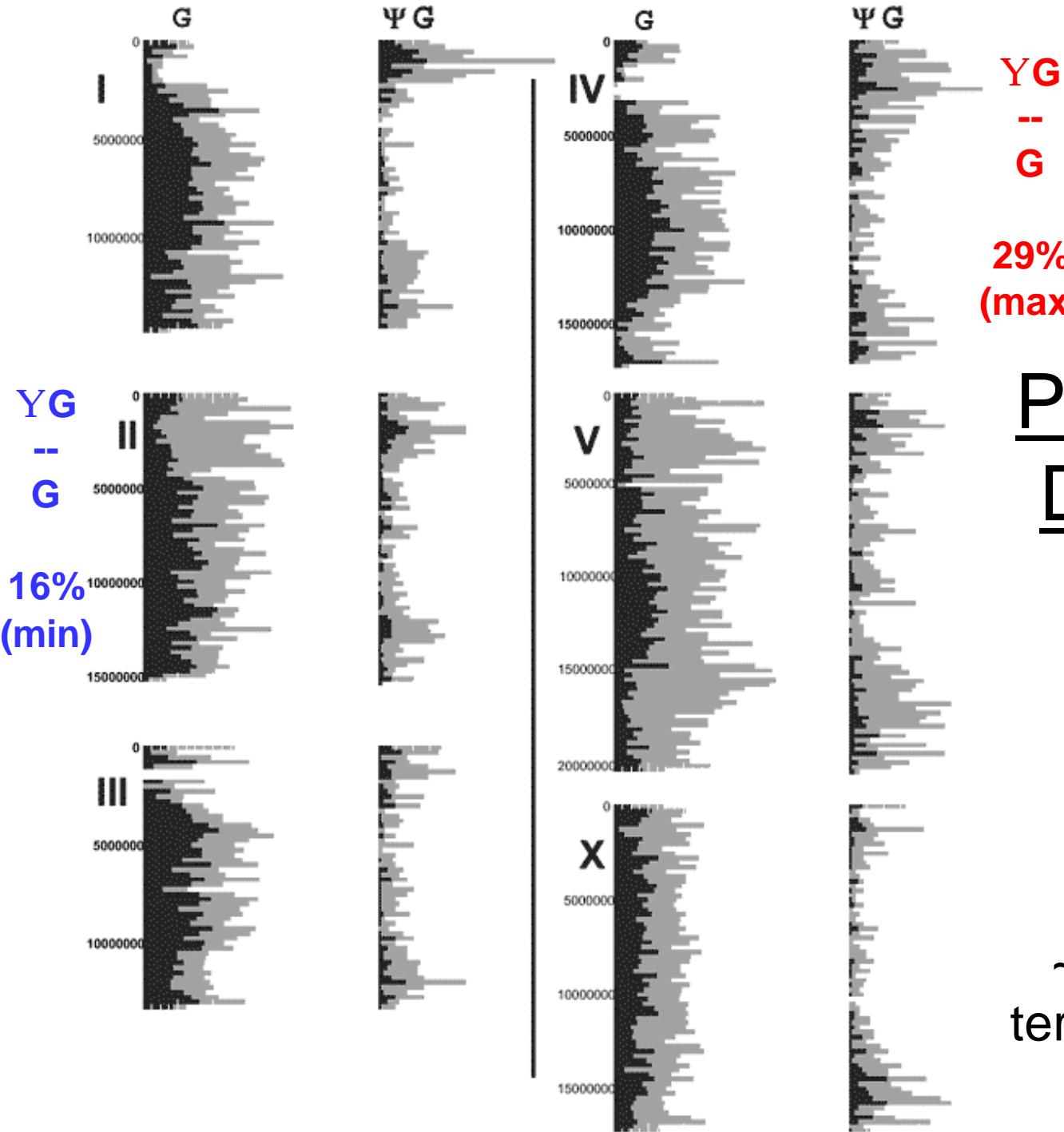
Example of a potential  $\Psi$ G with frameshift in mid-domain

# Most Common Worm “Pseudofolds” #1

G Rank (Number matches)	$\Psi$ G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description	G Rank (Number matches)	$\Psi$ G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description
<b>1</b> (769)	<b>2</b>		d1ajw_— 2.1 Immuno-globulin	<b>6</b> (246)	<b>8</b>		d21bd_— 1.95 Nuc. receptor ligand-binding domain
<b>2</b> (555)	<b>6</b>		d1dec_— 7.3 Knottin	<b>7</b> (243)	34		d1a17_— 1.91 Alpha/alpha superhelix
<b>3</b> (434)	<b>3</b>		d3lck_— 5.1 Protein kinase	<b>8</b> (227)	17		disp2_— 7.31 Classic zinc finger
<b>4</b> (302)	<b>1</b>		d1tsg_— 4.105 C-type lectin	<b>9</b> (215)	20		d1dai_— 3.29 P-loop NTP hydrolase
<b>5</b> (274)	<b>7</b>		d1zfo_— 7.33 Glucocorticoid receptor DNA- binding dom.	<b>10</b> (197)	13		d2aw0_— 4.34 Ferredoxin

# Most Common Worm “Pseudofolds” #2

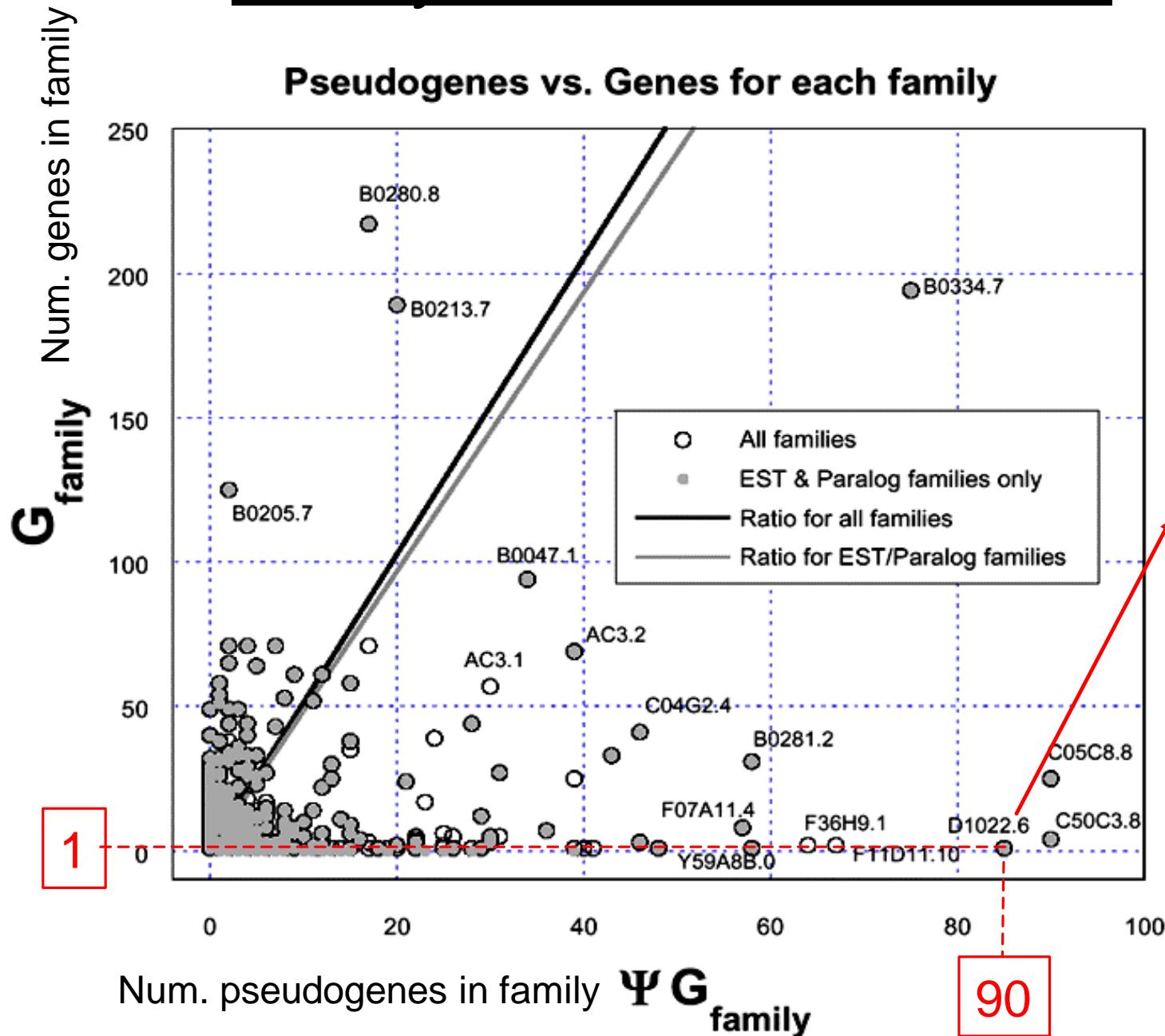
<b>ΨG Rank (Number matches)</b>	<b>G Rank</b>	<b>Fold</b>	<b>Representative Domain, SCOP 1.39 Number, Description</b>	<b>ΨG Rank (Number matches)</b>	<b>G Rank</b>	<b>Fold</b>	<b>Representative Domain, SCOP 1.39 Number, Description</b>
<b>1</b> (39)	<b>4</b>		d1tsg_— 4.105 C-type lectin	<b>6</b> (18)	<b>2</b>		d1dec_— 7.3 Knottin
<b>2</b> (32)	<b>1</b>		d1ajw_— 2.1 Immuno-globulin	<b>7</b> (17)	<b>5</b>		d1zfo_— 7.33 Glucocorticoid receptor DNA-binding dom.
<b>3</b> (27)	<b>3</b>		d3lck_— 5.1 Protein kinase	<b>8</b> (15)	<b>6</b>		d2lbd_— 1.95 Nuc. receptor ligand-binding domain
<b>4</b> (25)	11		d1cvl_— 3.56 Alpha/beta-hydrolase	<b>9</b> (13)	58		d1bus_— 7.14 Ovomucoid PCI inhibitor fold
<b>5</b> (23)	63		d1ako_— 4.93 DNAse-I fold	<b>9</b> (13)	19		d2bnh_— 3.7 Leu-rich, right-handed β/α superhelix



## Pseudogene Distribution on Chromo- somes

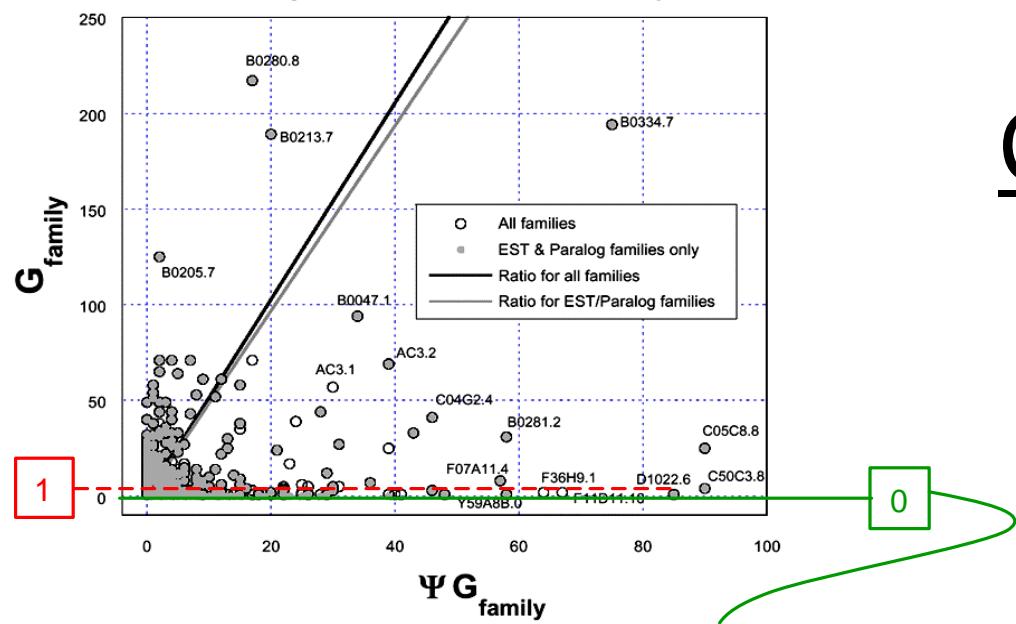
~50%  $\Psi G$  in  
terminal 3Mb vs  
~30% G

# Decayed Lines of Genes?



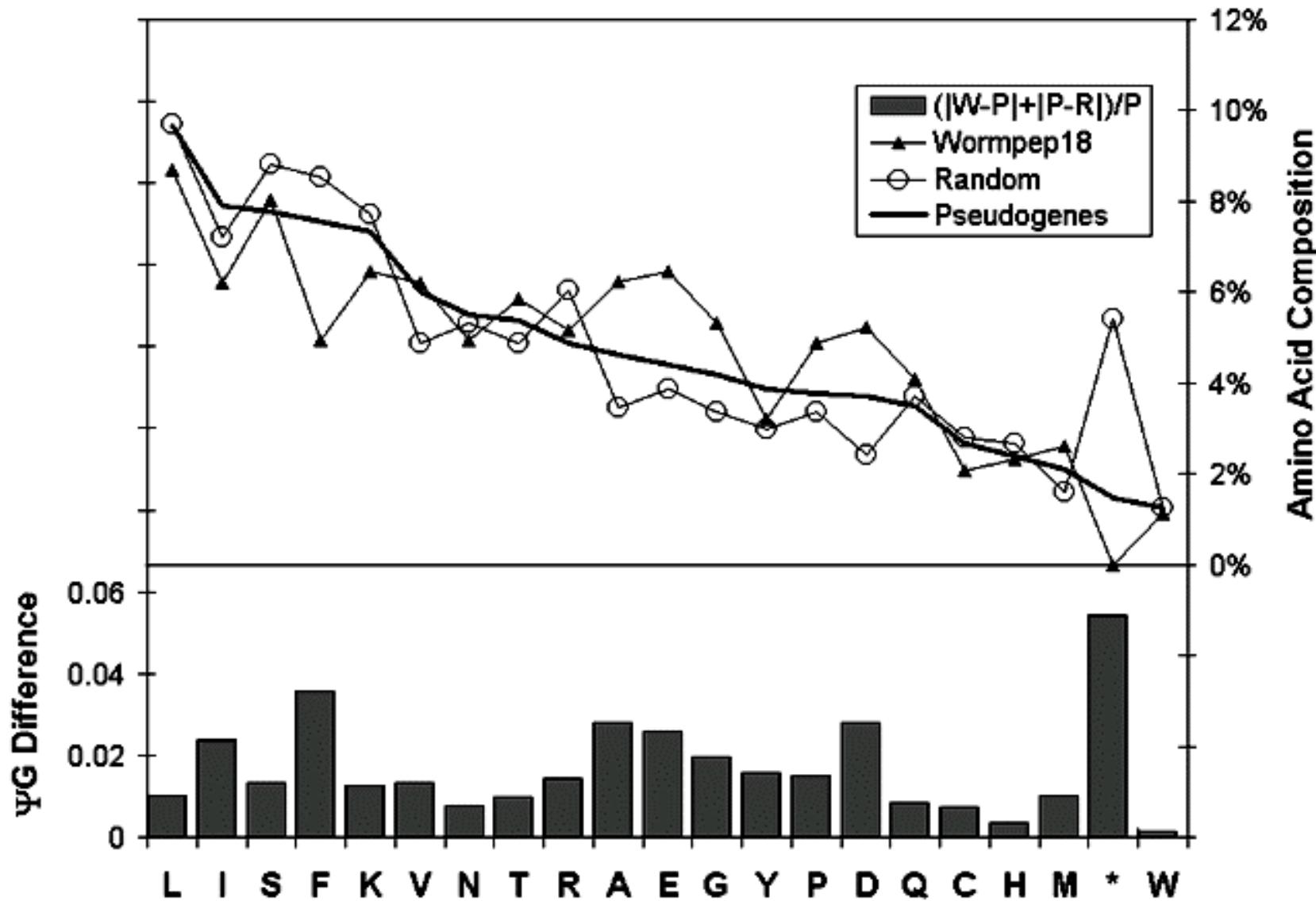
D1022.6 has  
90 dead  
fragments of  
itself – a  
disused line  
of chemo-  
receptors?

# Completely Dead Families



Rank	Number matches	Organism of closest match*	PROTOMAP family representative	Notes on representative
#1	7 *****	Yeast	YJA7_YEAST	Hypothetical protein in yeast
#2 =	5 ****	Human	XPD_MOUSE	Xeroderma pigmentosum group D complementing protein
#2 =	5 ****	Cow	CPSA_BOVIN	Cleavage and polyadenylation specificity factor
#4 =	4 ****	Frog	THB_RANCA	Thyroid hormone receptor beta
#4 =	4 ****	Human	SEX_HUMAN	SEX gene
#4 =	4 ****	Fly	MDR1_RAT	Multidrug resistance protein 1
#7 =	3 ***	Vaccinia virus	YVFB_VACCC	Hypothetical vaccinia virus protein
#7 =	3 ***	Fly	VHRP_VACCC	Host range protein from vaccinia
#7 =	3 ***	Human	IF4V_TOBAC	Eukaryotic initiation factor 4A
#7 =	3 ***	<i>E. coli</i>	ACRR_ECOLI	AcrAB operon repressor

# Amino Acid Composition of Pseudogenes is Midway between Proteins and Random



# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

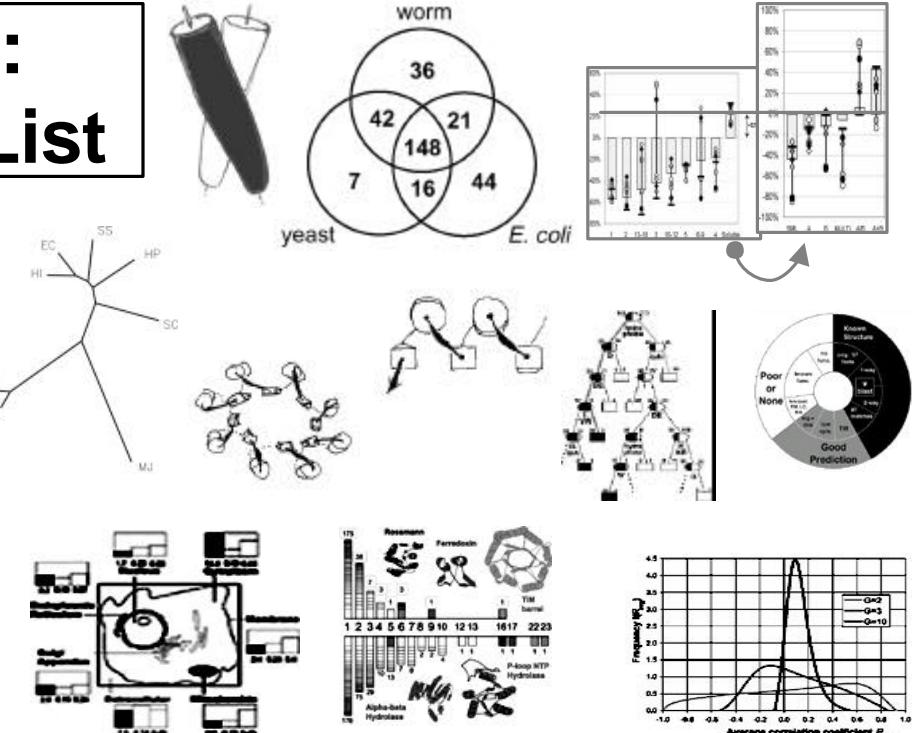
- **Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

- **Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

## 4 Using Folds to Interpret Expression Data

**Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

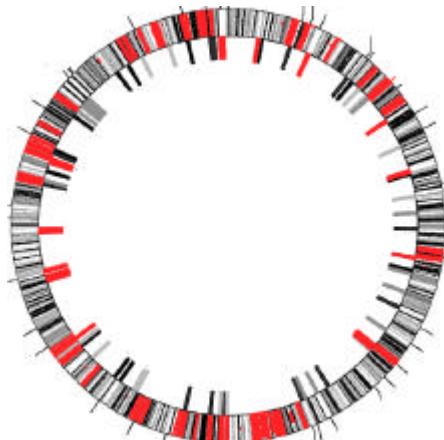
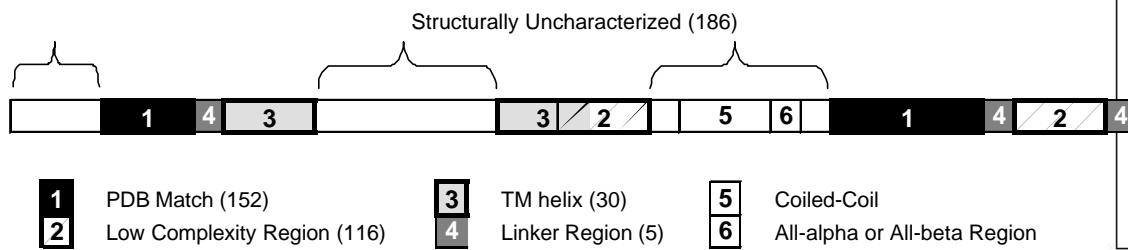
- **Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.



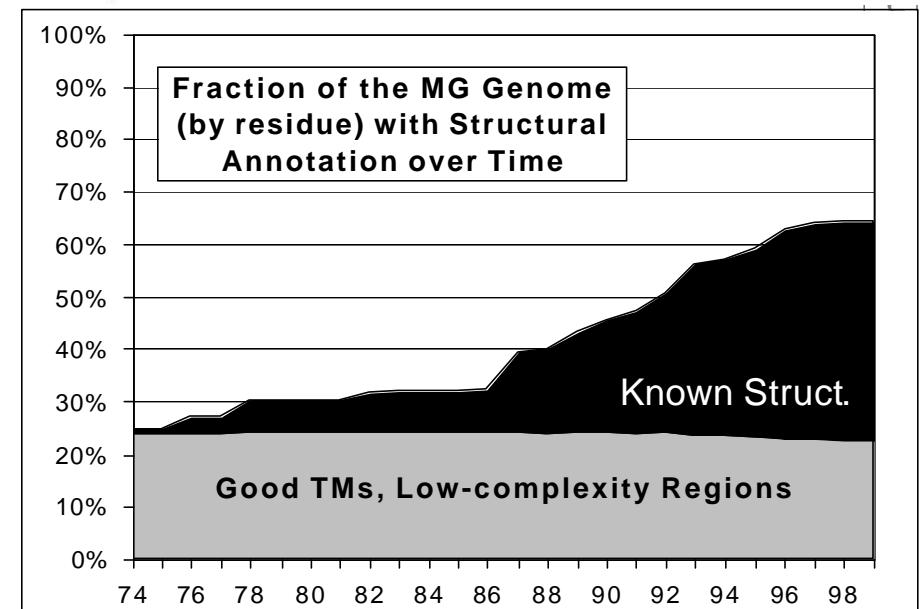
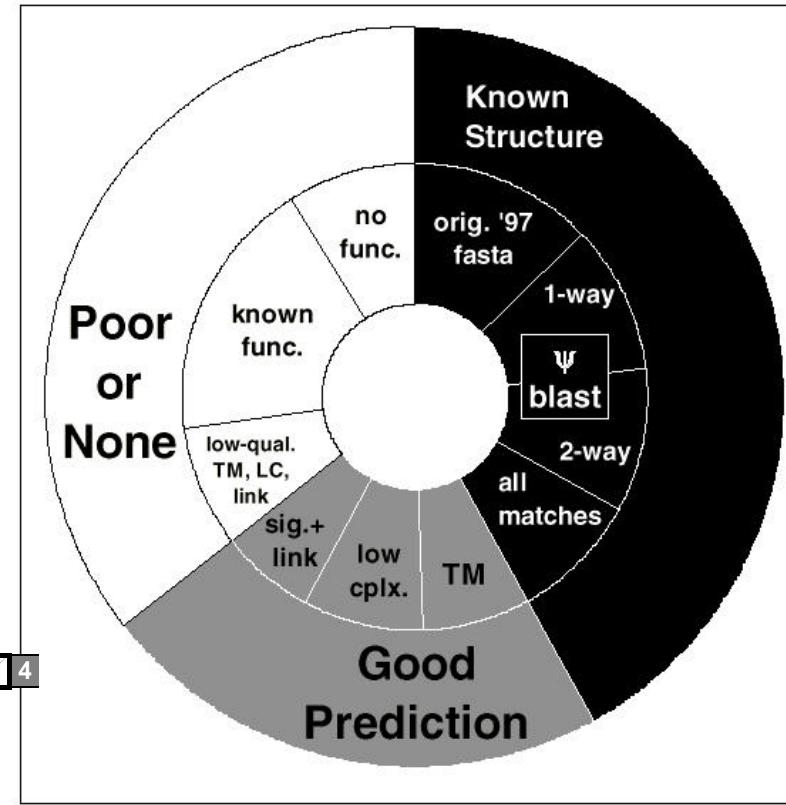
**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

**bioinfo.mbb.yale.edu**

# The Obvious Problem: Incomplete Coverage of the Genome by Known Folds



Status of  
M. gen.  
at end  
end of  
'98



# Bias Problem → Prediction, Experimental Structural Genomics

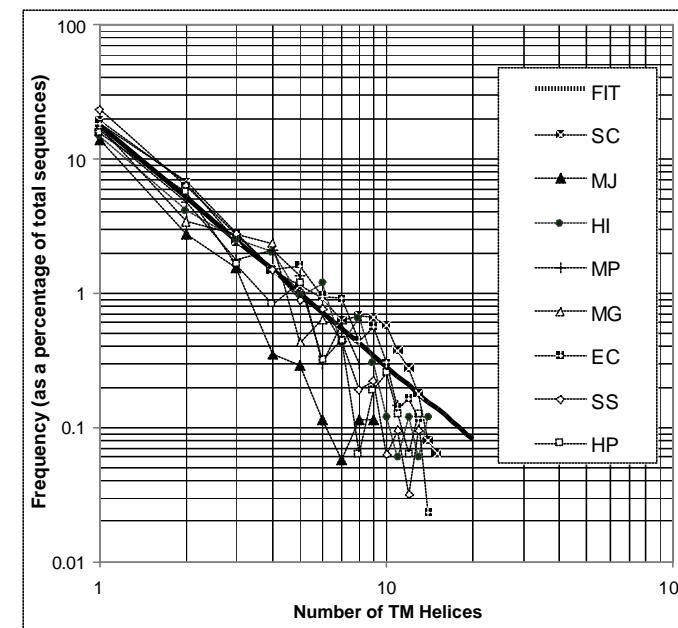
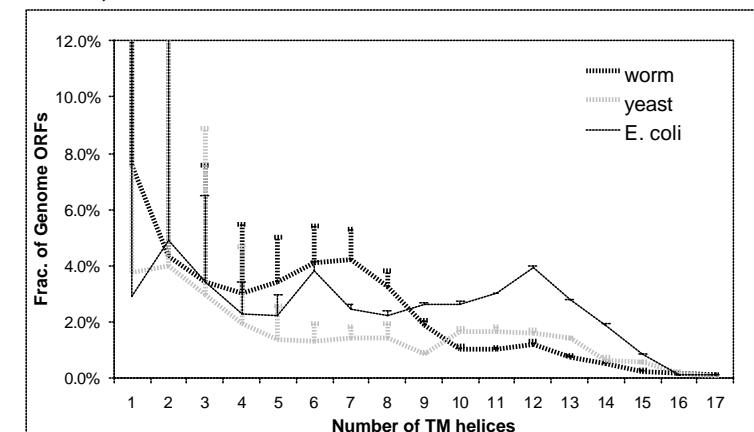
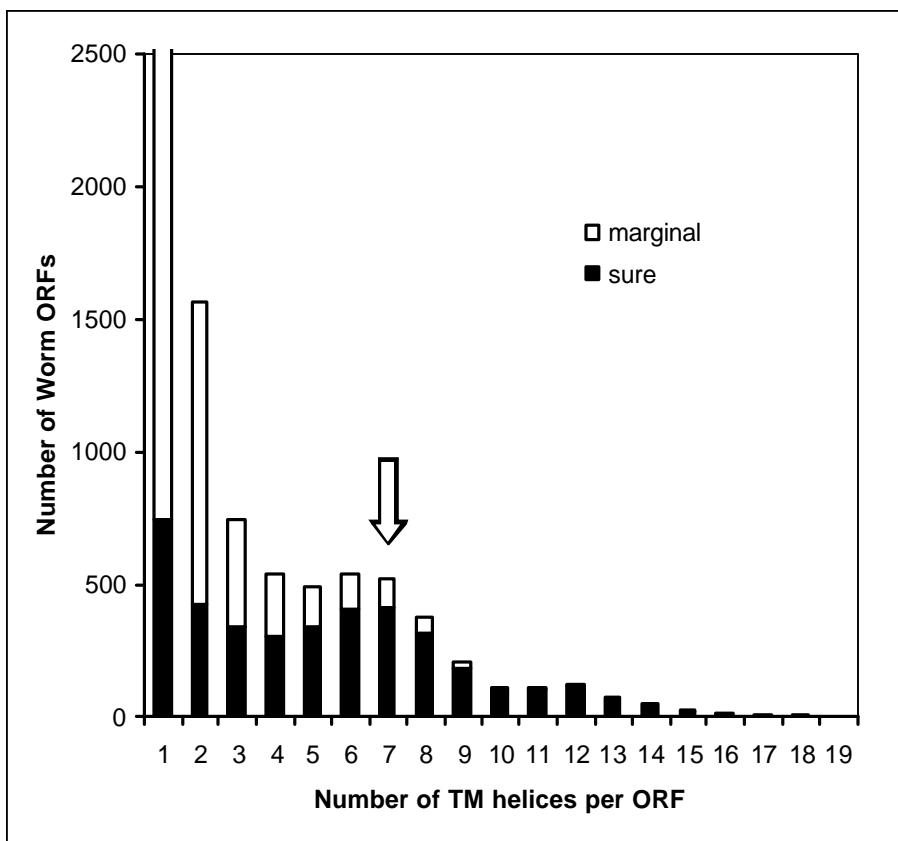
- Known Structures are Incomplete,  
Biased Sample from Genome, so...
  - ◊ Resample
  - ◊ Solve Structures
  - ◊ Predict Structures



Same  
Sampling  
Issues with  
US Census!!

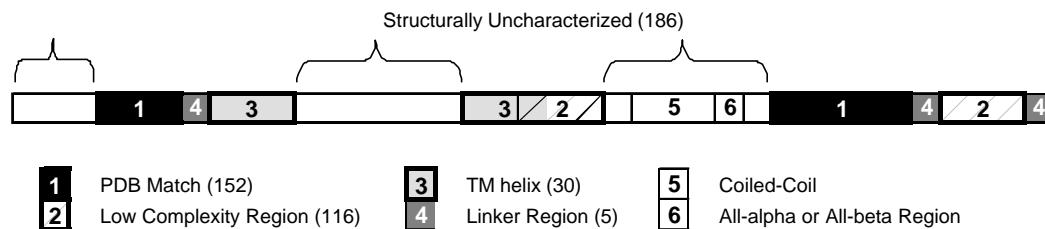
# TM-helix “prediction”

- 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<sup>o</sup> Structure Prediction

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



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

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

Not expected since.....

# 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)

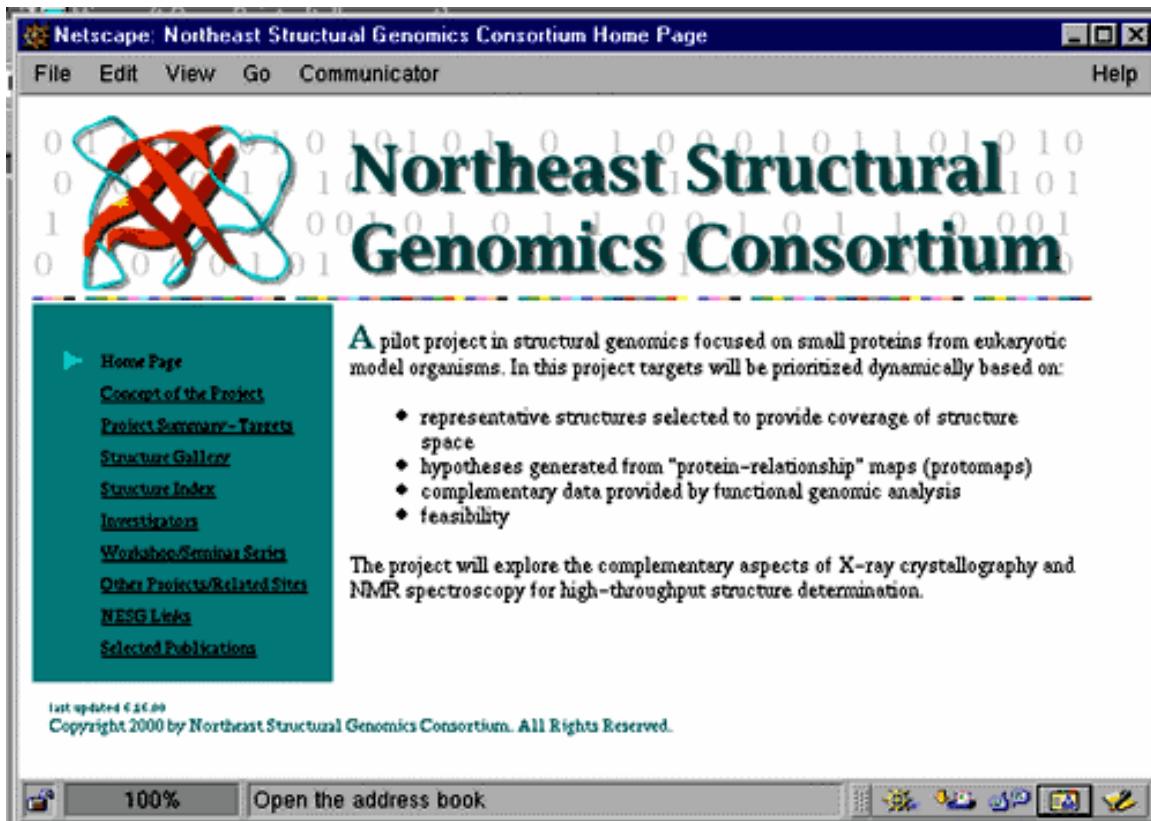
	Amino Acid Composition								Propensity (kcal/mole)		
	EC	HI	SS	SC	HP	MP	MG	MJ	TM-hlx	helix	strand
K	4.4	6.3	4.2	7.3	8.9	8.6	9.5	10.4	8.8	-1.5	-0.4
C	1.2	1.0	1.0	1.3	1.1	.8	.8	1.3	-2	-1.1	-0.8
R	5.5	4.5	5.1	4.5	3.5	3.5	3.1	3.8	12.3	-1.9	-0.4
N	4.0	4.9	4.0	6.1	5.9	6.2	7.5	5.3	4.8	-1	-0.5
Q	4.4	4.6	5.6	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

$\alpha$  -1.00 -1.02 -0.96 -1.00 -1.05 -1.03 -1.05 -1.01

$\beta$  -0.27 -0.33 -0.26 -0.36 -0.37 -0.38 -0.42 -0.36

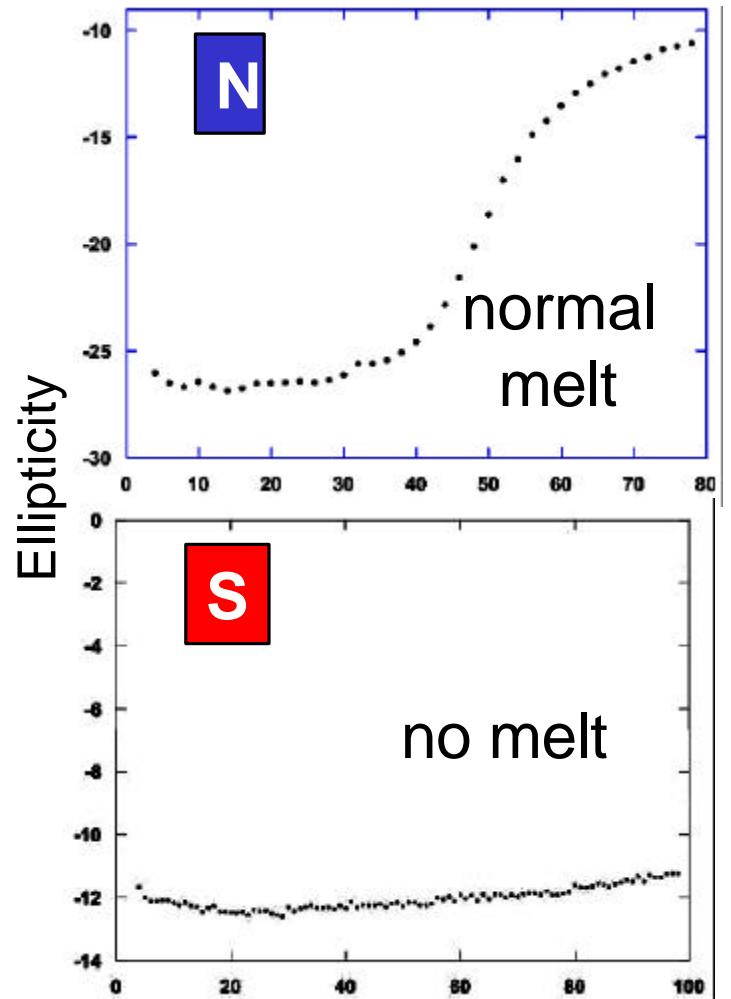
nesq.org



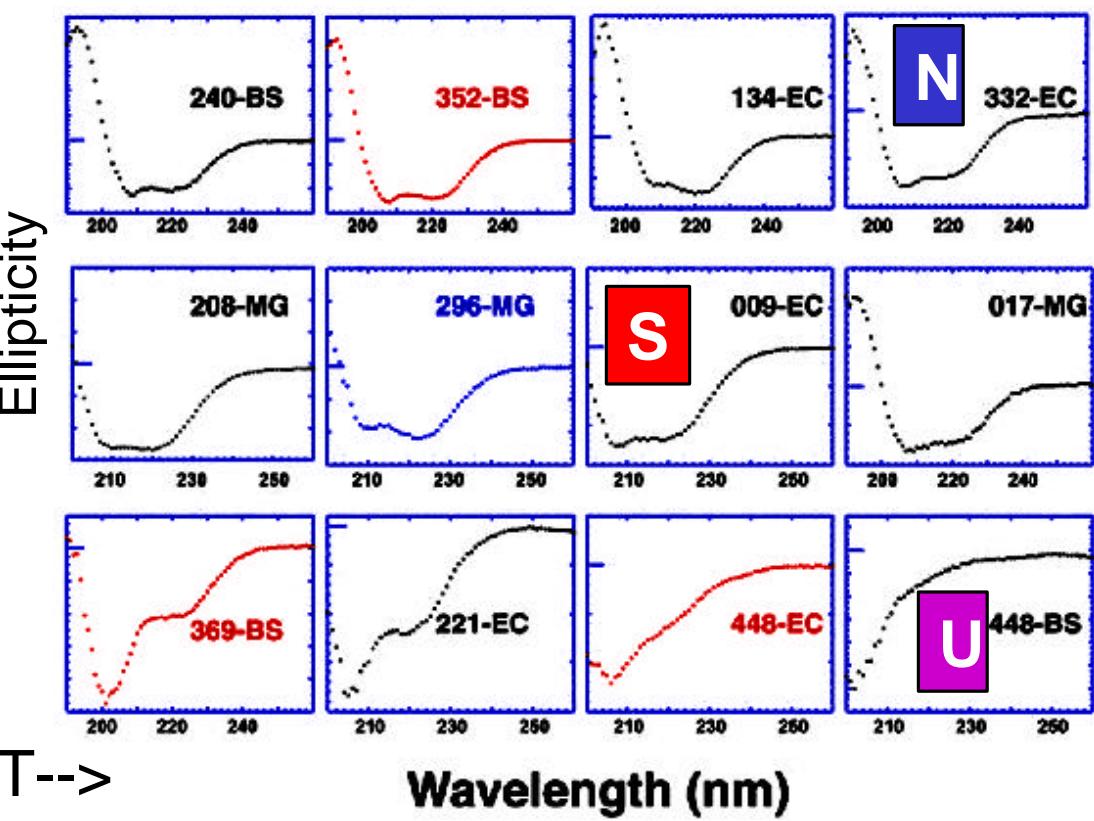
# G Montelione

Columbia, Cornell, Yale, Rutgers, Toronto....

Strange S		ORFID	Length	Source	Cloned	Expression	Fraction at 37°C	Purification	2 <sup>b</sup> Structure	Thermal melts	Minimal	Essential
9	262	EC	Yes	++	IB	Sol (0.1mM IPTG)	H	No	Yes	No		
17	176	MG	Yes	++	IB		H	No	No	No		
56	277	EC	Yes	-	-		-	-	-	Yes	Yes	
134	100	EC	Yes	+	Sol	Sol	H	Yes	No	Yes		
208	196	MG	Yes	+	IB+Sol	Sol	H	Yes	No	No		
221	154	EC	Yes	++	IB+Sol	Sol <sup>1</sup>	H+C	No <sup>2</sup>	Yes	Yes		
240	292	BS	Yes	++	IB+Sol	Sol	H	Yes	No	Yes		
296	129	MG	Yes	++	IB	Sol (25°C)	H	Yes	No	No		
332	239	EC	Yes	++	IB+Sol	Sol	H	Yes	Yes	Yes		
352	166	BS	Yes	++	IB	Sol (30°C)	H	Yes	No	NO		
369	557	BS	Yes	++	IB+Sol	Sol	H+C	Yes	No	Yes		
448	150	BS	Yes	++	IB+Sol	Sol	C	-	Yes	Yes		
448	150	EC	Yes	++	IB+Sol	Sol	C	-	Yes	Yes		
461	425	BS	No	-	-	-	-	-	-	No	Yes	



## M.gen. CD results



# Progress Database

**Project Progress Summary**

Select the database entries to display: [Advanced Form](#)

Target Organism: Any Organism Attribute: Tertiary Structure

Homolog Organism: Any Organism

Laboratory: Any Laboratory

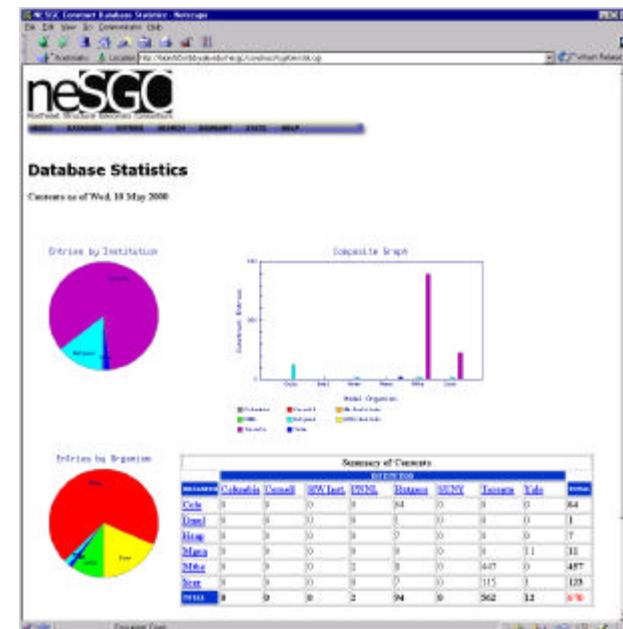
Institution: Any Institution

[Create Table](#) Display: Reduced

There are 14 entries which have either X-ray or NMR tertiary structure, excluding restricted access and test entries.

[Download Sequences](#)

NESGC ID	PDB ID	Organism	Target ORF	Cloned	Expressed	Purified	HSQC Spectrum	CD Assignments	Crystal	X-ray Data	X-ray Structure	NMR Assignments	NMR Structure	Biochemical Function
WR64Tc	-	Cele	-	*	*	*	*	-	-	-	-	*	*	*
TT1	-	Mthe	MT 0152	*	*	*	-	-	*	*	*	-	-	*
TT2	-	Mthe	MT0129	*	*	*	-	-	*	*	*	-	-	*
TT3	-	Mthe	MT1790	*	*	*	-	-	*	*	*	-	-	*
TT4	-	Mthe	MT0150	*	*	*	-	-	*	*	*	-	-	*
TT5	-	Mthe	MT1791	*	*	*	-	-	*	*	*	-	-	*
TT9	Leik	Mthe	mt1048	*	*	*	*	-	-	-	-	*	*	*
TT10	Leik	Mthe	mt1615	*	*	*	*	-	-	-	-	*	*	*
TT11	-	Mthe	mth0040	*	*	*	*	-	-	-	-	*	*	*
TT12	-	Mthe	mt1699	*	*	*	*	-	-	-	-	*	*	*
TT13	-	Mthe	mt1184	*	*	*	*	-	-	-	-	*	*	*



**Construct Database Record for TT15**

[Edit](#) [Advanced Sequence](#)

GENERAL	
NESGC ID:	TT15
Created:	04-21-2000
Modified:	04-21-2000
Investigator:	Chenyl
Laboratory:	AnneArach
Institution:	Toronto

TARGET	
ORF ID:	Mthe0035
Target Organism:	Thermotolerant (Mtc)
Homolog ORF ID:	
Homolog Organism:	NA

PLASMID	
Insert P.A. Protein:	Start: 0 End: 0
Vector:	pET28b
Sequence Verified:	NA
Date Cloned:	04-21-2000
Cloning Comments:	

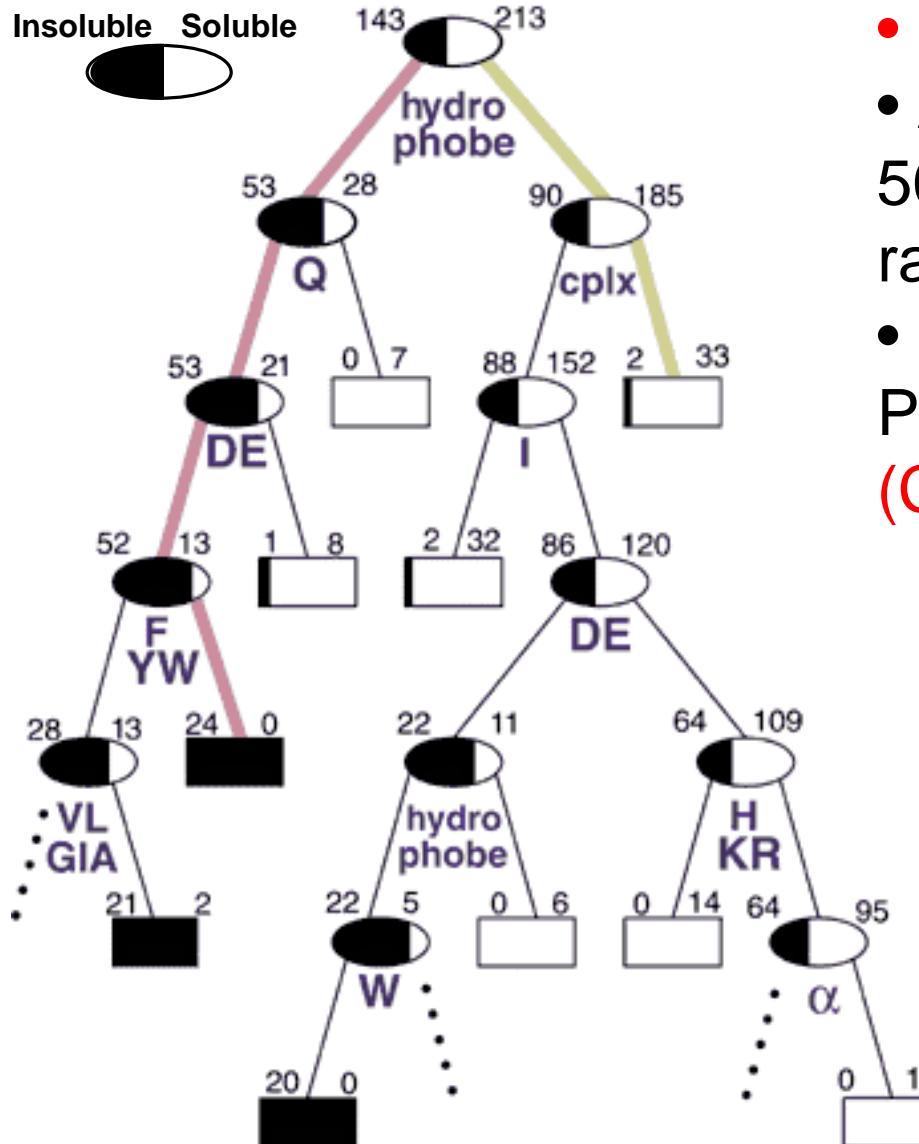
  

CONDITIONS	
Expression Level:	NA
Expansion Comments:	

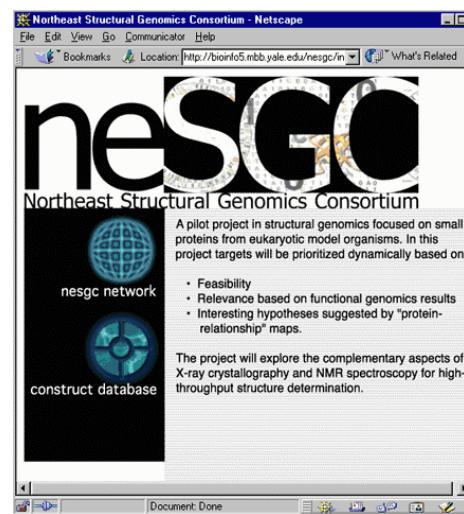
  

SOLUBILITY	
Phase:	NA
Concentrability:	5
Solubility of Cell Extract:	5
Solubility of Purified Protein:	NA
Solubility Comments:	

# 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.

# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

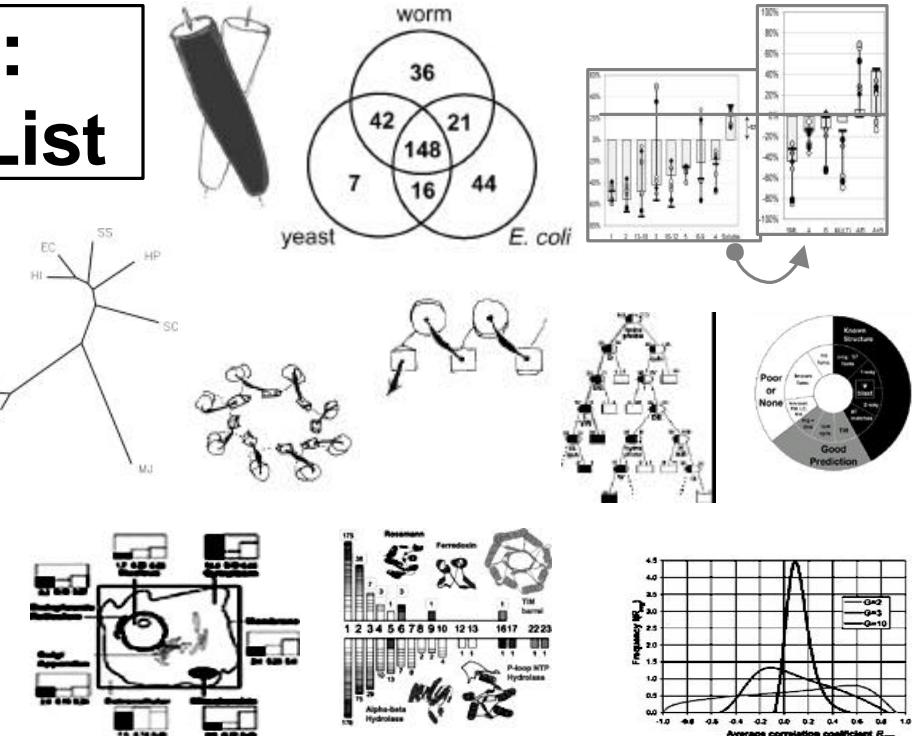
**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

**2 Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

**3 Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

**4 Using Folds to Interpret Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

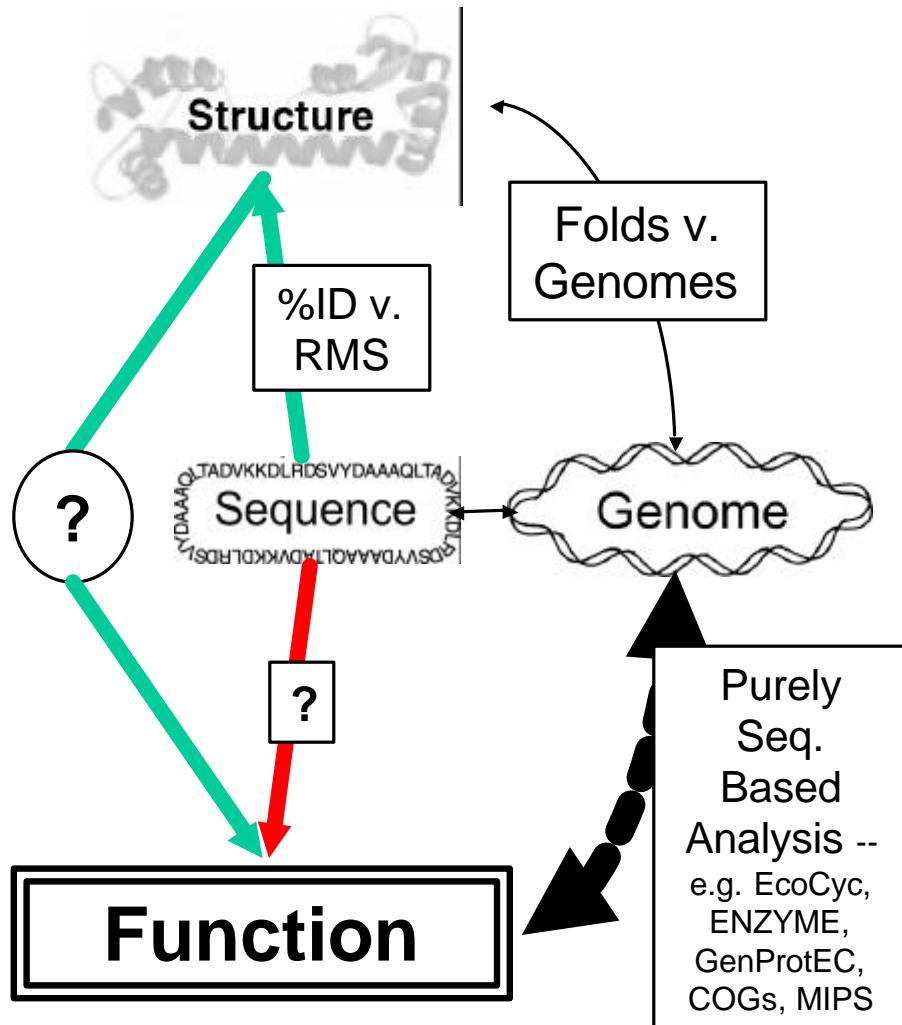
**5 Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.



**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
S Teichmann, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

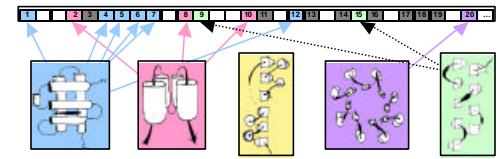
**bioinfo.mbb.yale.edu**

# Adding Structure to Functional Genomics, Function to Structural Genomics

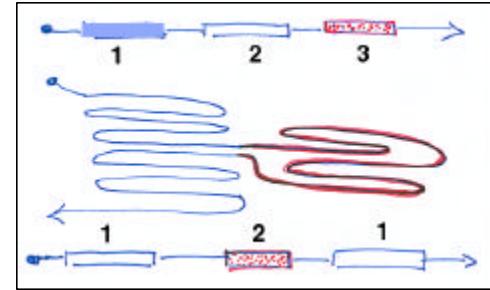


## Why Structure? Do we really need it?

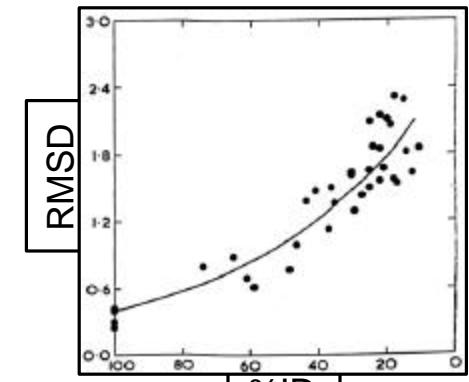
1 Most Highly Conserved



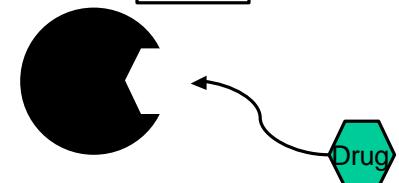
2 Precisely Defined Modules



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



4 Link to Chemistry,  
Drugs



# Functional Classification

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

**GenProtEC**  
(*E. coli*, Riley)

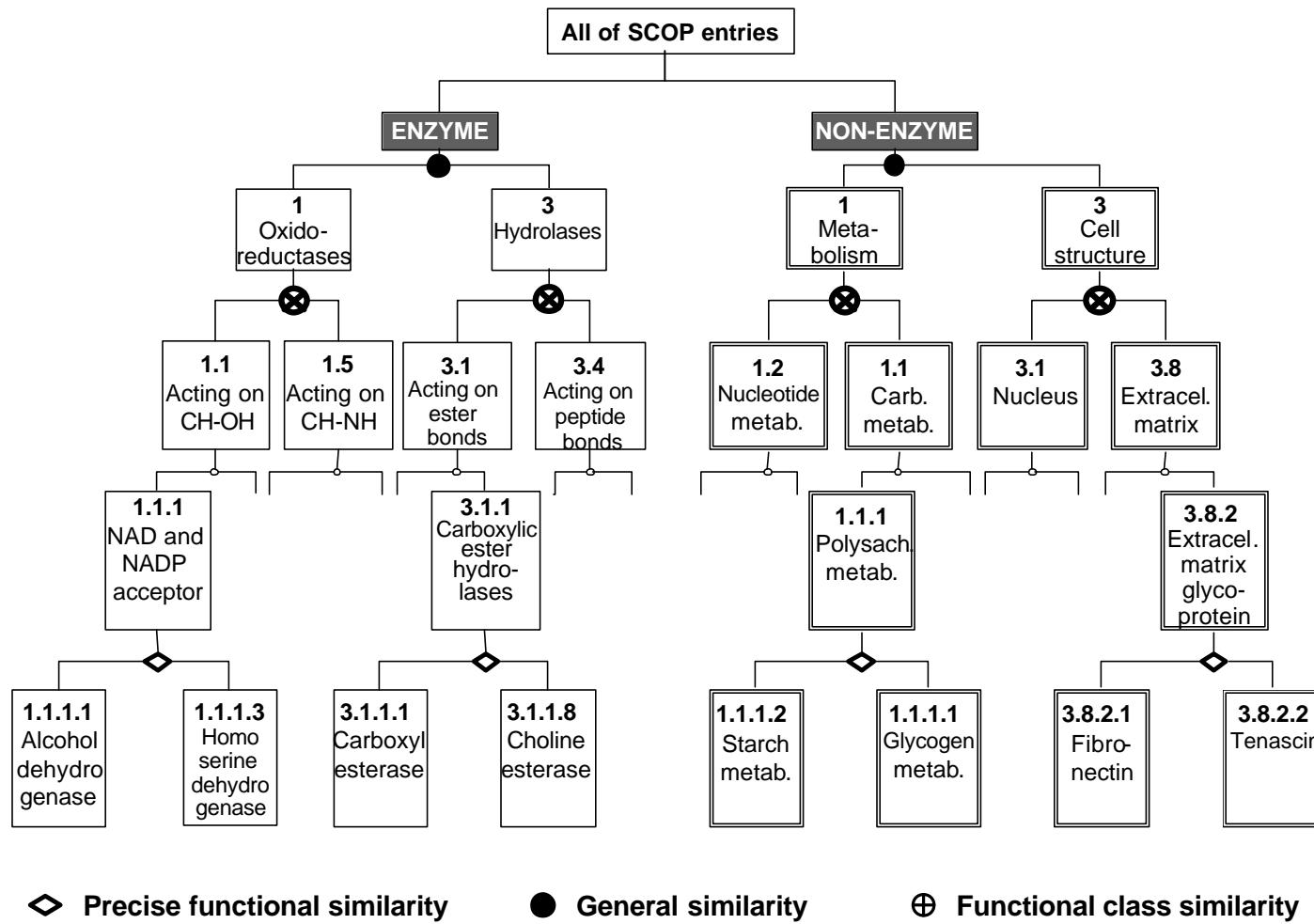
**“Fly”**  
(fly, Ashburner)  
now extended to  
**GO** (cross-org.)

**MIPS/PEDANT**  
(yeast, Mewes)

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

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

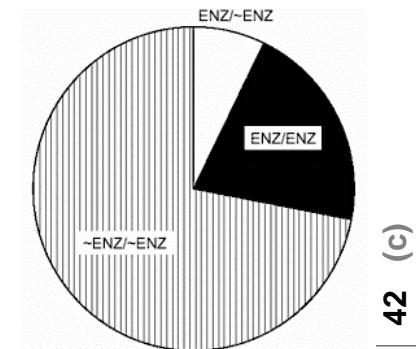
# A Simple Scheme for Functionally Classifying Protein Structures



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

A Combined Scheme, merging ENZYME + FLY, with some annotation from MIPS, GenProtEC, & manual additions for proteins not in either (e.g. Ig's)

Also: COGs



# Fold-Function Combinations #1

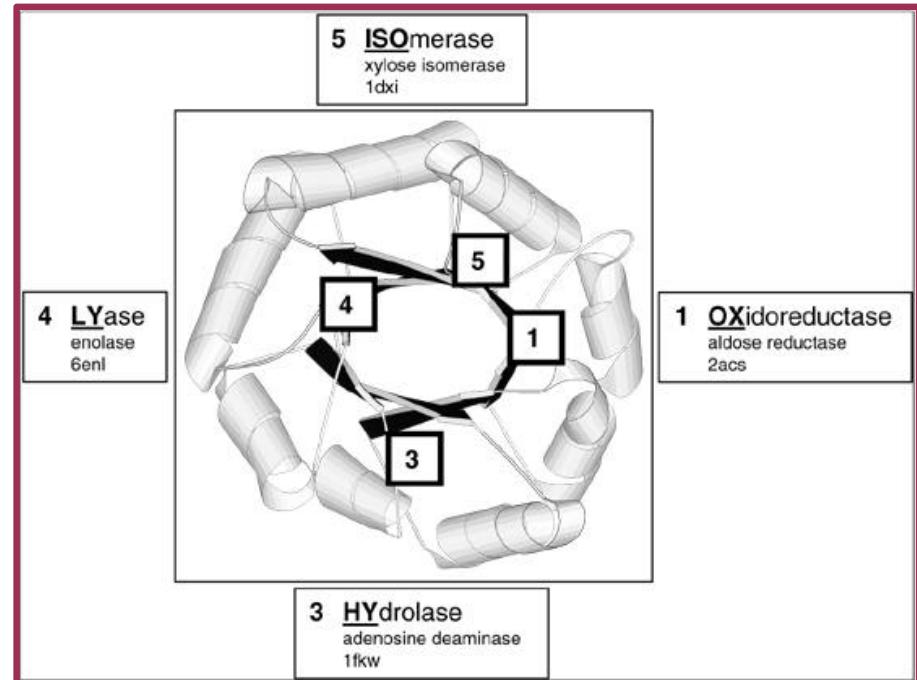
**Many Functions on the Same Fold**  
-- e.g. the TIM-barrel

at what degree of divergence?

Sequence Diverg. (%ID,  $P_{seq}$ )

Structural Diverg. (RMS,  $P_{str}$ )

Functional Diverg. (%SameFunc)

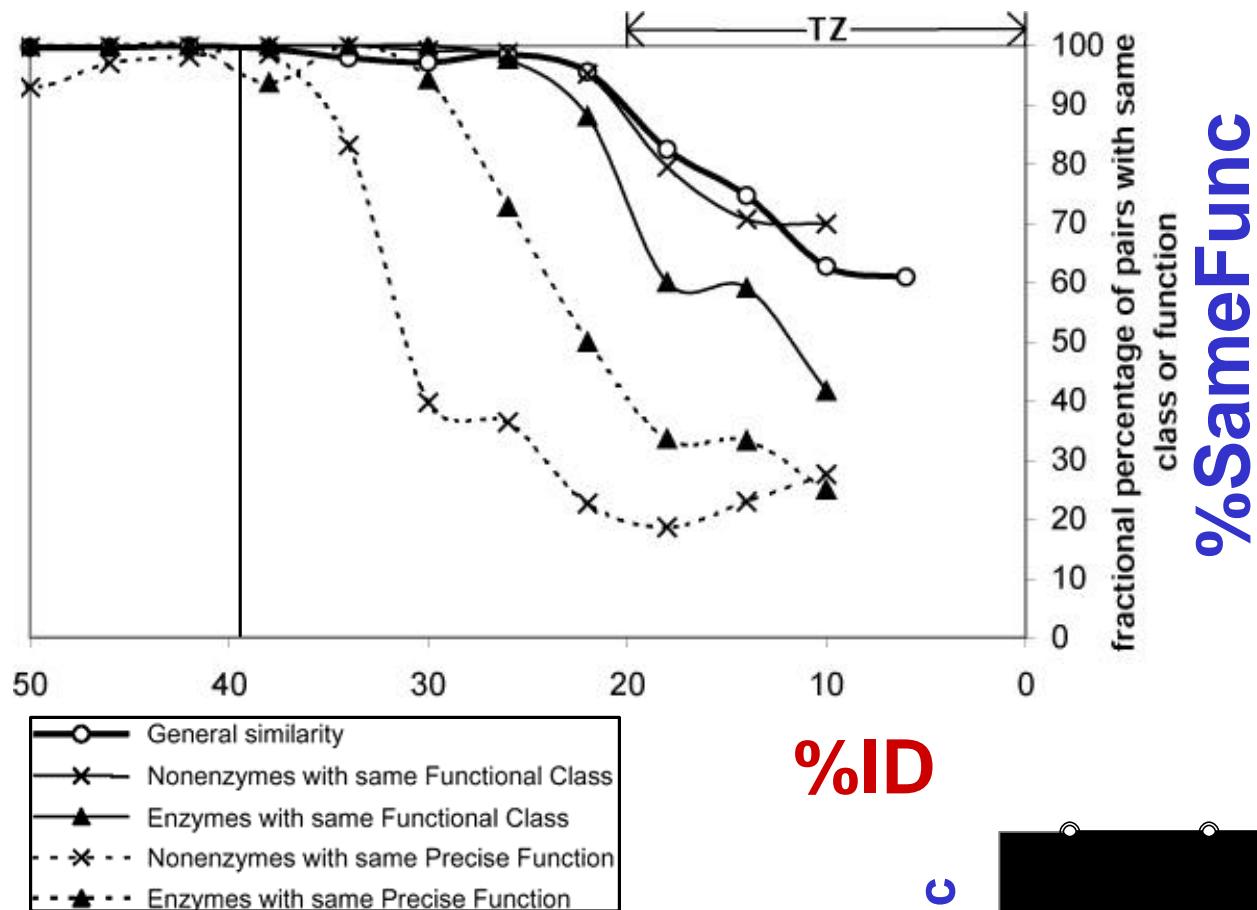


Compare large number of pairs of sequences that have same fold but different functions.

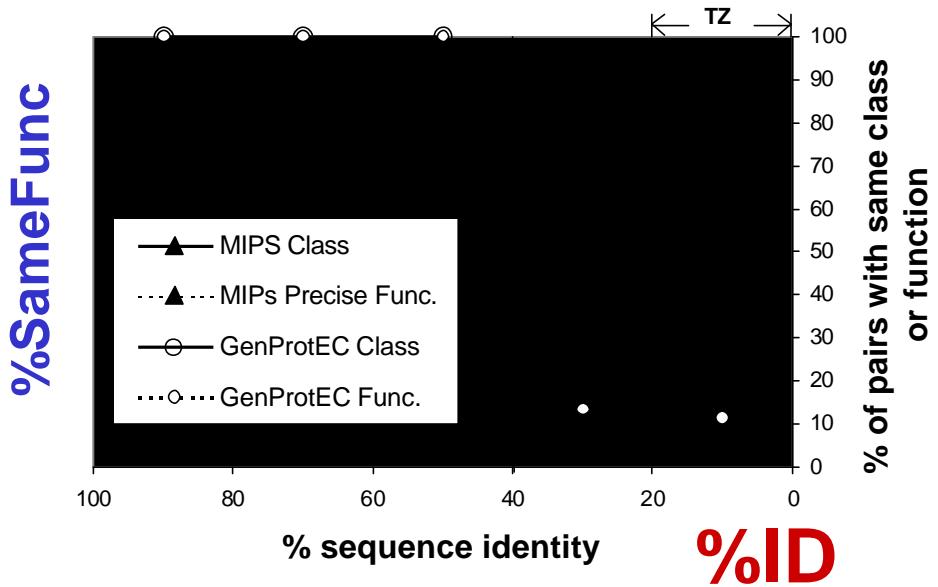
89%	Human	TP Isomerase	5.3.1.1	Same Exact Func.
45%	Chick	TP Isomerase	5.3.1.1	
	E coli	TP Isomerase	5.3.1.1	
	E coli	PRA Isomerase	5.3.1.24	
	B ster.	Xylose Isomerase	5.3.1.5	
	E coli	Aldolase	4.1.3.3	
	Yeast	Enolase	4.2.1.11	Both Class 5
	Rat	K-channel B-sub.	NON-ENZ	Completely Different
	Photobact.	Flavoprotein?	NON-ENZ	Same Exact Func.

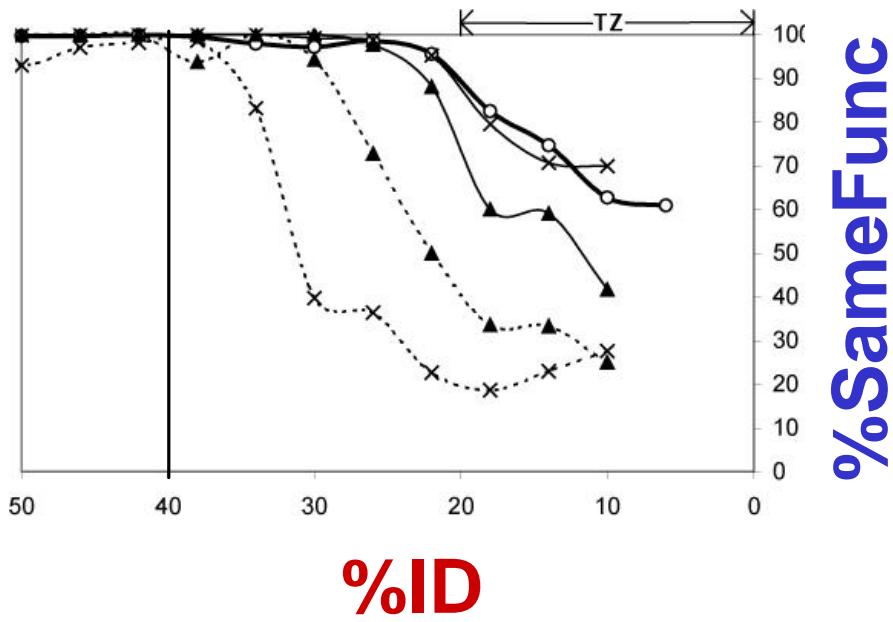
45% (red arrows) and ~20% (green arrows) indicate sequence divergence levels between the compared organisms.

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



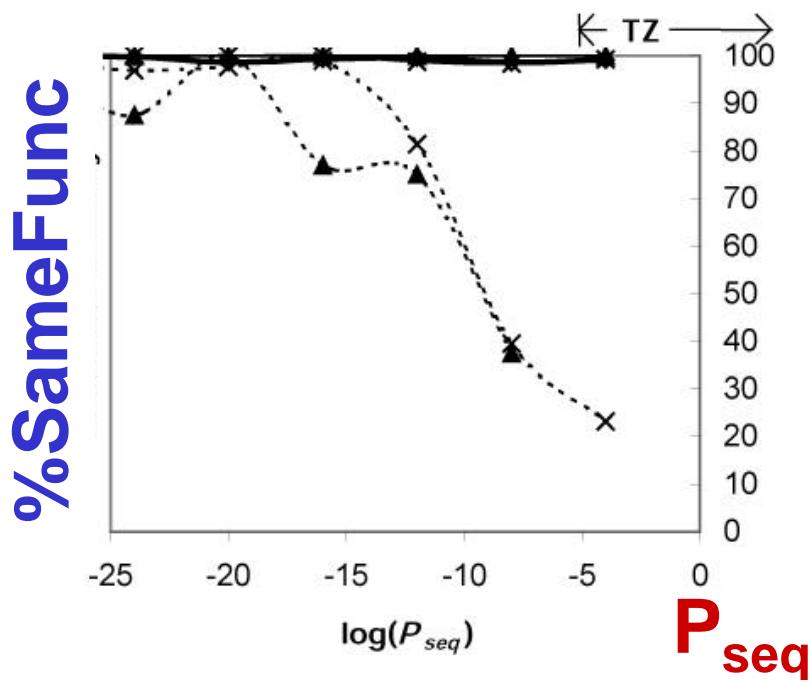
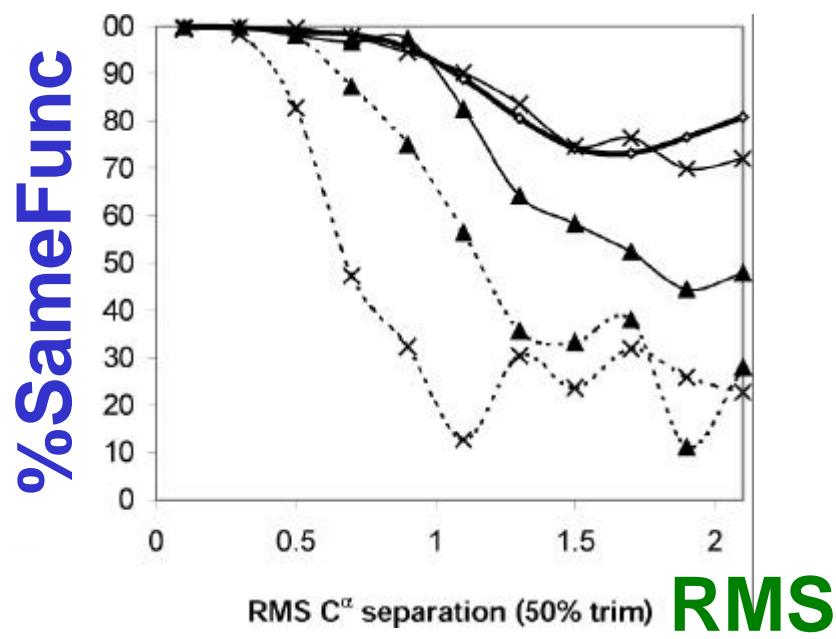
See at what %ID have diff. function (both broad & precise). Use 4 func. classifications -- ENZYME, FLY (+extra), MIPS, GenProtEC





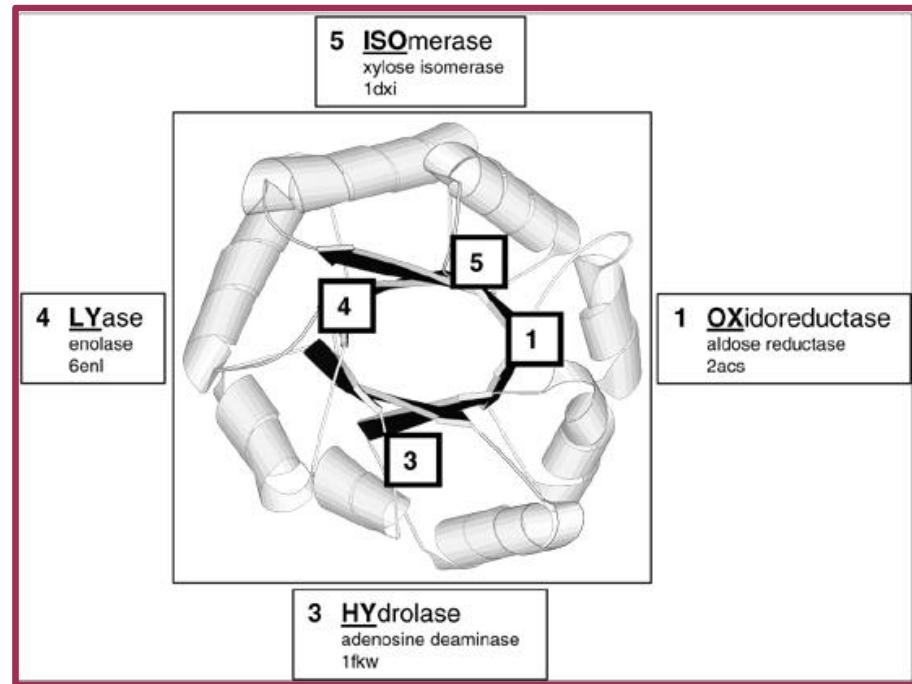
## Relationship of Similarity in Sequence & Structure to that in Function II

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

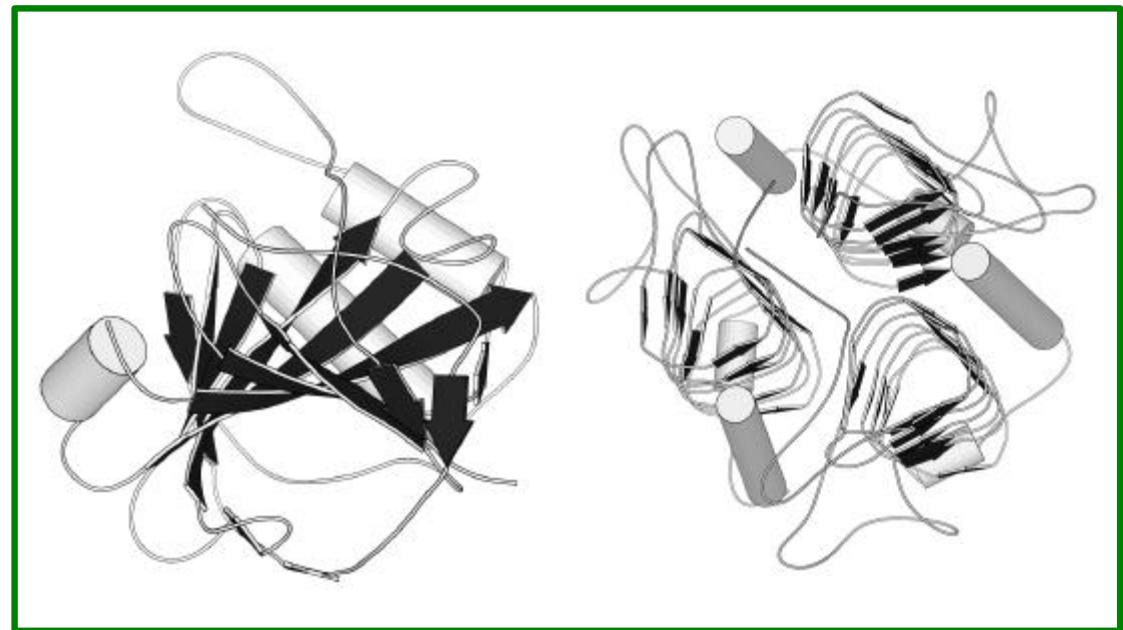


# Fold-Function Combinations #2

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



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

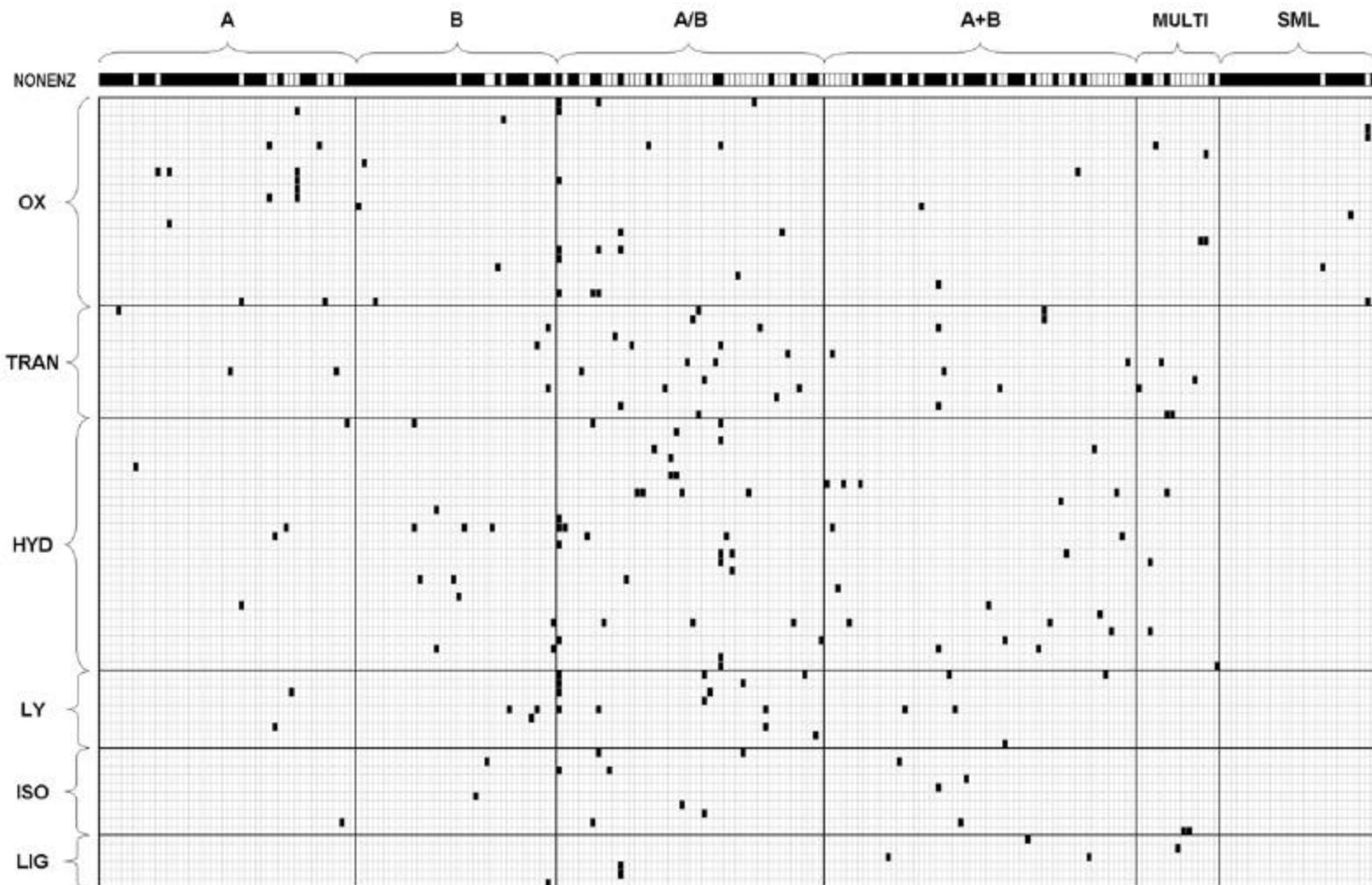


# Fold-Function Combinations

91 Enzymatic Functions  
+ Non-Enzyme

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

229 Folds



# Most Versatile Functions

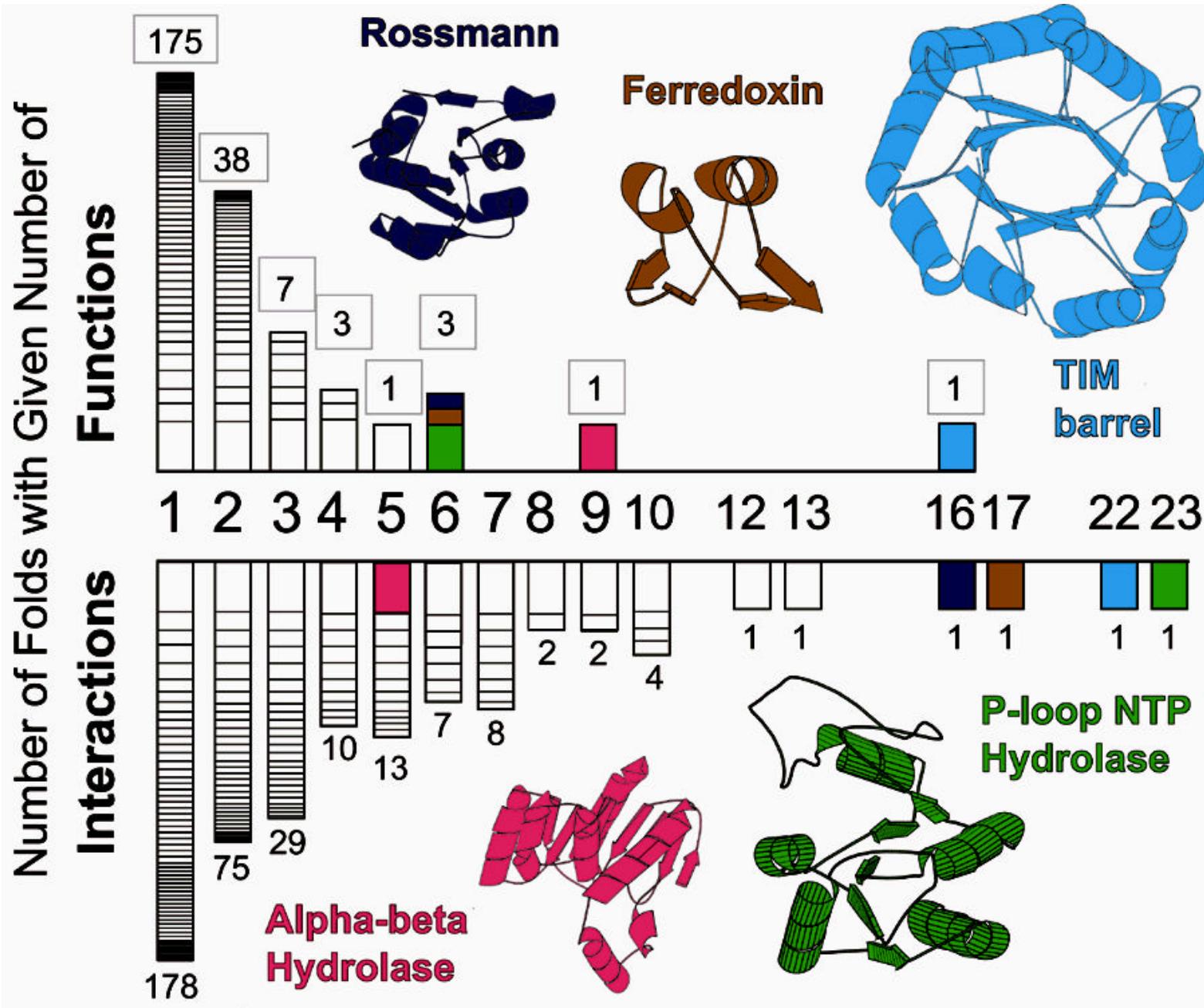
Sort table rows....

		A	B	A/B	A+B	MULTI S
1.004	tenh	d1mmog				
1.019	d1ffa	1ffa				
1.021	d1gsta1					
1.034	d1occh_					
1.037						
1.053	1llp					
1.054	2abk					
1.061	1gal					
1.063	1phc					
1.068	1vnc					
1.070	d1occe_					
1.077	1fps					
1.080	1poc					
2.005	iaac					
2.018	1jbc					
2.020	1nyf					
2.024	1snc					
2.029	1arb					
2.033	Zeng					
2.043	2sil					
2.047	1hcb					
2.053	d1caua_					
2.055	1bdo					
2.056	1duu					
3.001	1byb					
3.002	1tmi					
3.009	d1rvva_					
3.011	1udh					
3.013	3chy					
3.018	1xel					
3.021	d1nbaa_					
3.024	19ky					
3.026	1phr					
3.029	2hng					
3.030	1srx					
3.037	1pao					
3.040	3pgm					
3.041	1opr					
3.043	1ode					
3.045	1ama					
3.046	d1gpmaz					
3.047	1ub					
3.048	2ace					
3.049	d1masa_					
3.054	d1alka_					
3.055	1xaa					
3.057	d1ttqp_					
3.061	3pgk					
3.064	1agx					
3.065	3pkf					
3.066	1ayl					
4.001	1fus					
4.002	2baa					
4.005	d2kaaua_					
4.020	d1imkaa_					
4.031	1fxd					
4.035	d3rubs_					
4.036	d1dcoca_					
4.049	1iba					
4.058	1mut					
4.060	1lba					
4.073	1friqi					
4.082	d1pxa.1					
4.084	d1rima_					
4.086	1mrj					
4.087	1dp					
5.001	1hcl					
5.004	2cae					
5.005	1tpf					
5.007	1imf					
5.009	1rp1					
7.029	1hip					

## 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)

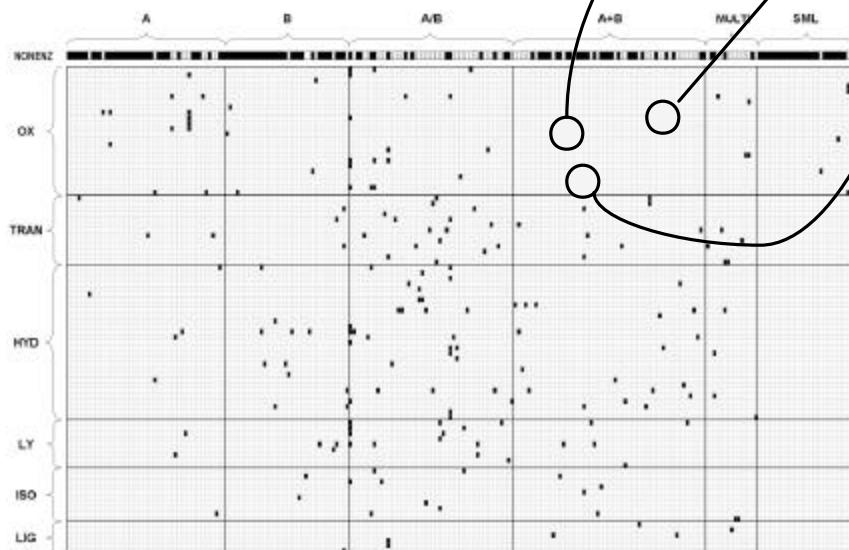
# Most Versatile Folds – Relation to Interactions



Similar results  
Martin et al.  
(1998)

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

# 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	8	5	5	35
HYD	4	11	30	18	4	4	67
LY	2	3	13	5			23
ISO	1	2	7	4	2		16
LIG	1	2	3	1	1		7
sum	57	55	99	69	20	31	331

3

SCOP

	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	

ENZYME

[ 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

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

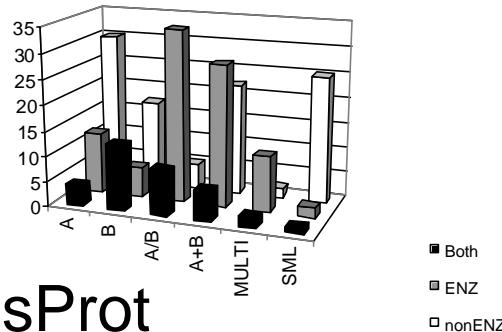
CATH (Thornton)

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

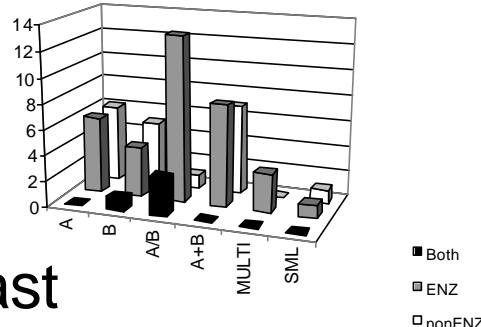
MIPS YFC (Mewes)

	SCOP						
	A	B	A/B	A+B	MULTI	SML	
metabolism	1	3.5	2.3	10	4.5	1.3	0.8
energy	2	1.1	1.2	5	1.5	-0.3	0.2
growth, div., DNA syn.	3	4.3	3.6	4	4.5	1.3	1.2
transcription	4	1.5	1.9	2.2	1.5	0.5	0.8
protein synthesis	5	1	0.9	0.7	1.3	0.3	0.2
protein targeting	6	1.2	1.7	2	1.6	0.5	0.3
transport facilitation	7	0.9	0.5	0.7	0.6	0.4	
intracellular transport	8	1.8	2.1	1.6	0.6	1	
cellular homeostasis	9	0.9	0.7	1.2	0.3	0.3	0.1
signal transduction	10	1	1	1.1	0.3	0.7	0.5
cell process, defense...	11	1.5	1	2.6	1.8	0.7	0.5
ionic homeostasis	13	0.5	0.3	0.4	0.4	0.2	

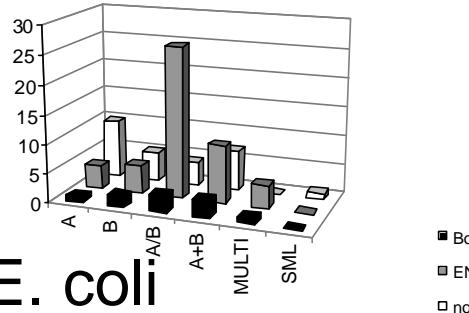
SwissProt



Yeast

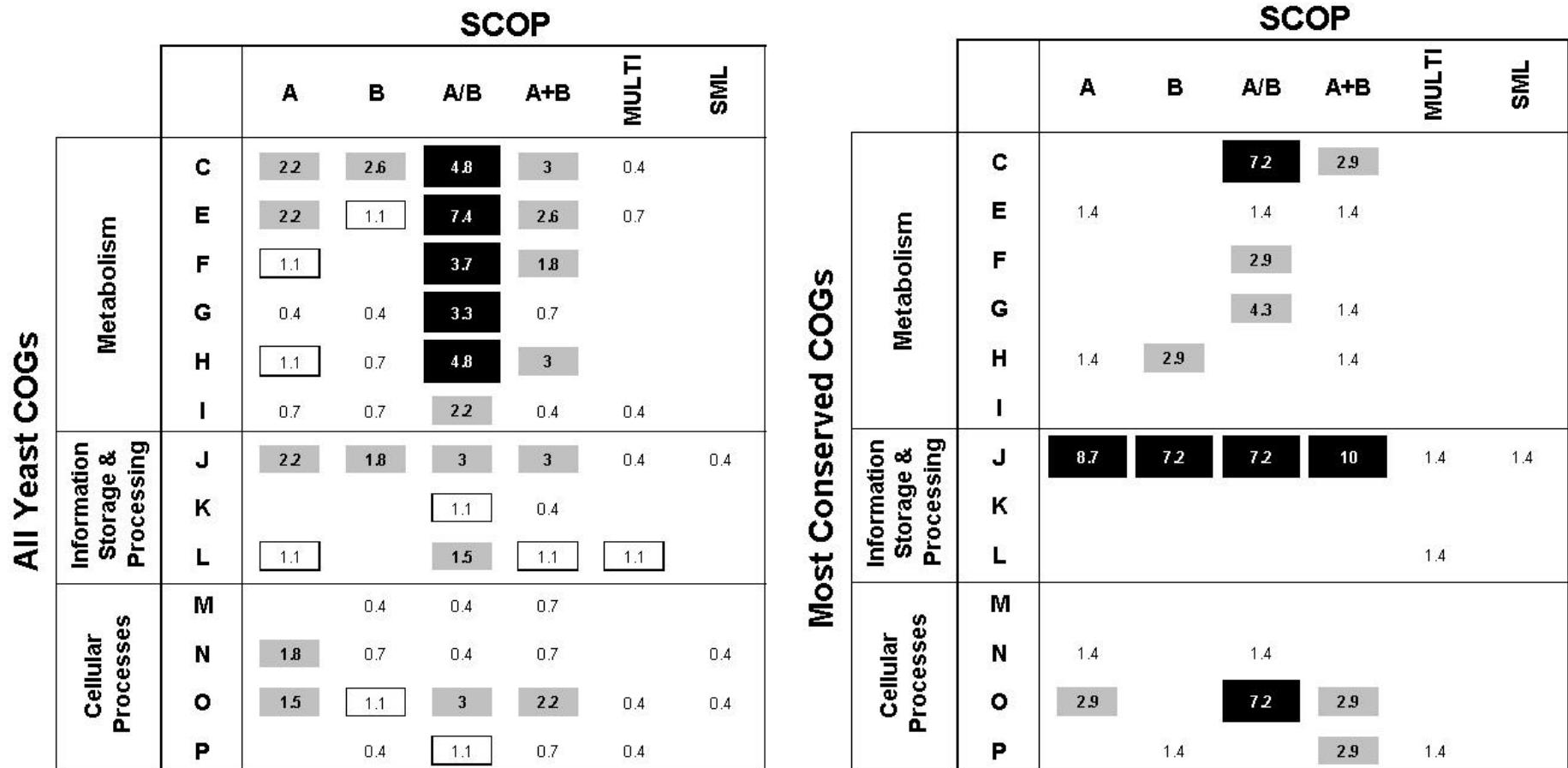


worm



E. coli

# Different Structure Function Relationships for Most Ancient Proteins



(Scop, Murzin/Chothia; COGs, Koonin/Lipman)

# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

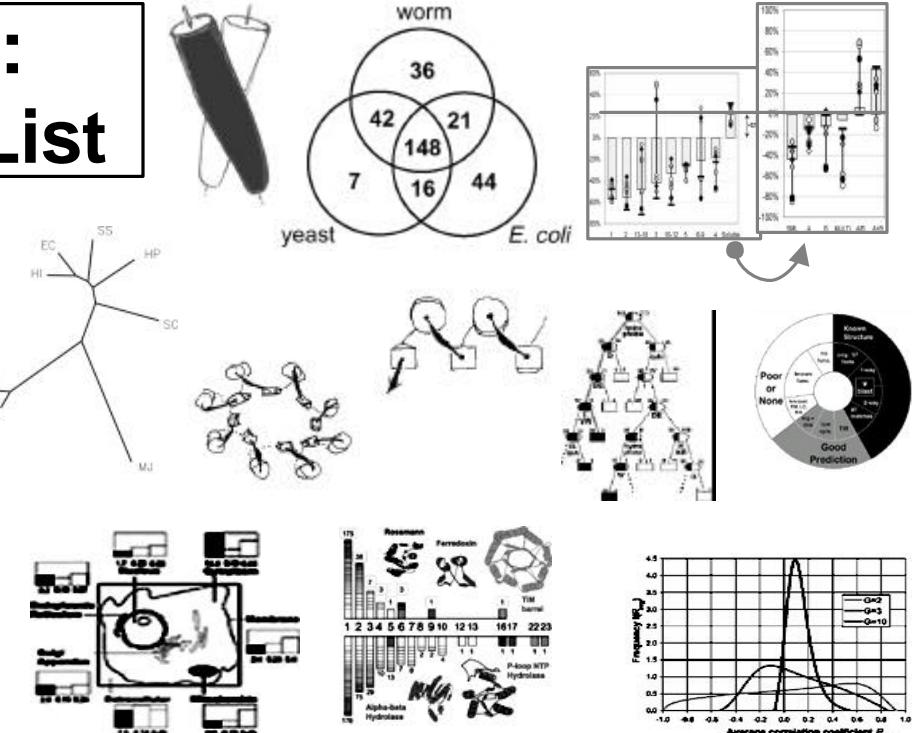
- **Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

- **Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

## 4 Using Folds to Interpret Expression Data

**Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

- **Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.



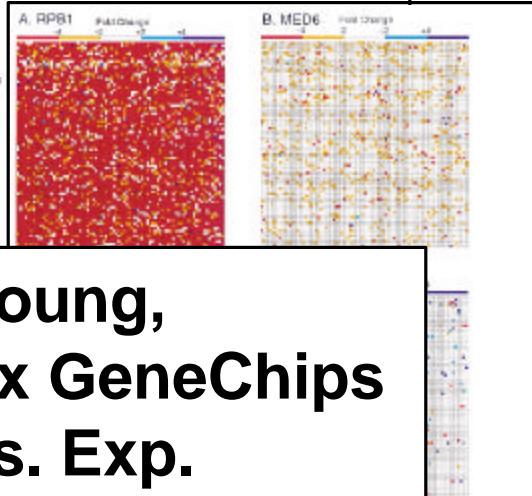
**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

**bioinfo.mbb.yale.edu**

## Dissecting the Regulatory Circuitry of a Eukaryotic Genome

Frank C. P. Holstage,<sup>a</sup> Ezra G. Jennings,<sup>a</sup>  
 John J. Wyrick,<sup>a</sup> Tong Ibin Lee,<sup>b</sup>  
 Christoph J. Hengartner,<sup>a</sup> Michael R. Green,<sup>a</sup>  
 Todd R. Golub,<sup>b</sup> Eric S. Lander,<sup>b</sup>  
 and Richard A. Young<sup>a,b</sup>

and Michael N. Young.<sup>1</sup>  
<sup>1</sup>Whitehead Institute for Biomedical Research  
Cambridge, Massachusetts 02142  
<sup>†</sup>Department of Biology  
Massachusetts Institute of Technology  
Cambridge, Massachusetts 02139  
<sup>‡</sup>Howard Hughes Medical Institute  
Program in Molecular Medicine  
University of Massachusetts Medical Center



# Young, Affymetrix GeneChips Abs. Exp.

regulation which is superimposed on that due to specific transcription factors, a novel mechanism accounts for regulation of oncogene expression.

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 article-supplement to Iyer-Hegde et al., (2005) *Biochem. Biophys. Res. Commun.* 380, 81-87

The Transcription  
of Sporulation in  
The Web Companion

# Brown, marrays, Rel. Exp. over Timecourse

Also:  
SAGE (mRNA);  
2D gels for  
Protein  
Abundance  
(Aebersold,  
Futcher)

# 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 BOSS-MACDONALD, AMY SHERMAN, G. SEBASTIEN ROEDER, AND MICHAEL SNYDER

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

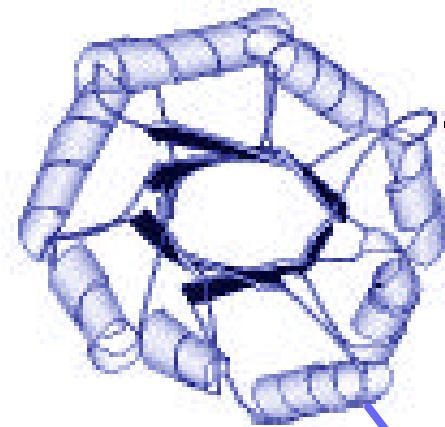
**ABSTRACT** Analysis of the function of a particular product typically involves determining the expression pattern of the gene, the subcellular location of the protein, and the phenotype of a null strain lacking the protein. Conditionally expressed genes are often created as an additional tool that simultaneously generates constructs for all the analyses and is suitable for mutagenesis of any given *S. cerevisiae* gene. Depending on the transposon used, the yeast gene is fused to a coding region for  $\beta$ -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 physical analysis. The transposon can be reduced by site-specific recombination to a smaller element that leaves a cleavage tag inserted in the encoded protein. In addition, the ability for insertion of immunodetection purposes, the expression cassette can be modified.

yeast gene is fused to a coding region for  $\beta$ -galactosidase green fluorescent protein. Gene expression can therefore be monitored by chemical or fluorescence assays. The T1 positions create insertion mutations in the target gene, the T2 positions create recombination mutations, and the T3 positions create point mutations. The transposon can be reduced by cassette-specific recombination to a smaller element that leaves a single base pair insertion in the untagged protein. In addition, it is feasible to tag a variety of immunodetected proteins, the entire genome, and certain cell lines for selection of transposon-containing clones. The modified construct containing the T1-T3 cassette can be used to generate stable T-1, T-2, and T-3 cell populations. Transposition of the BAC-PAC vector into the genome of a cell line creates a transposon insertion, which induces a frameshift mutation in the targeted gene, whereas insertion into the T1-T3 cassette does not. Thus, these three types of recombination allow for detection of a transposon insertion in a cell line. When a transposon insertion is detected, the T-1, T-2, and T-3 cell lines can be used to map the insertion site.

# Snyder, Transposons, Protein Abundance

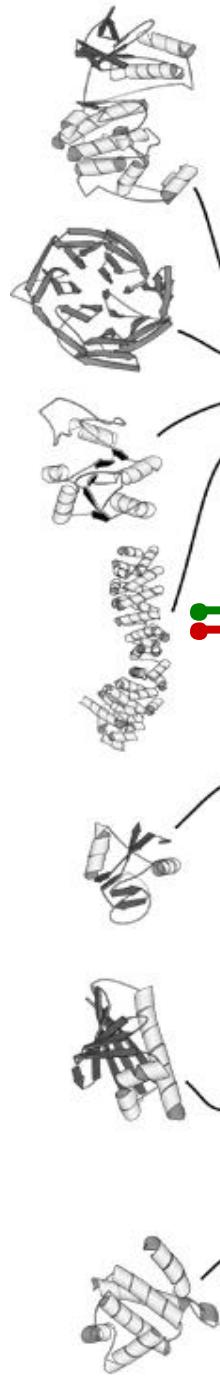
intergrating the A8799 and S851 genes (Fig. 2). A8799 is an essential gene whose encoded protein localizes to the apical pole body region (N. Rupin and M.S., unpublished data). The S851 gene encodes a karyopherin-like protein.

# Common Parts: the Transcriptome



Fold	Fold Class	Rep. PDB	Composition		Rel. Diff. [%]	Rank	
			Genome [%]	Transcriptome [%]		Genome	Young
TIM barrel	$\alpha/\beta$	1byb	4.2	8.3	+98	5	1
P-loop NTP hydrolases	$\alpha/\beta$	1gky	5.8	5.2	-11	3	2
Ferredoxin like	$\alpha/\beta$	1fxd	3.9	3.4	-14	6	3
Rossmann fold	$\alpha/\beta$	1xel	3.3	3.3	0	8	4
7-bladed beta-propeller	$\beta$	1mda*	6.4	2.9	-55	2	5
alpha-alpha superhelix	$\alpha$	2bct	4.4	2.7	-37	4	6
Thioredoxin fold	$\alpha/\beta$	2trx	1.7	2.7	+63	14	7
G3P dehydrogenase-like	$\alpha/\beta$	1drwt	0.2	2.7	+1316	81	8
beta grasp	$\alpha/\beta$	1igd	0.6	2.6	+348	36	9
HSP70 C-term. fragment	multi	1dky	0.8	2.6	+231	31	10
Leu-zipper	$\alpha$	1zta	3.8	2.1	-46	7	15
Protein kinases (cat. core)	multi	1hcl	6.8	1.6	-77	1	18
alpha/beta hydrolases	$\alpha/\beta$	2ace	2.2	0.9	-62	10	32
Zn2/C6 DNA-bind. dom.	sml	1aw6	2.6	0.3	-89	9	75

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%



Fold of	Freq.		Change					Rep. PDB
	Genome	Transcriptome	CDC28	CDC15	Diauxic Shift	Sporulation	E. coli heat shock	
Protein kinases (cat. core)	1	18	94	98	139	60	100	1p38
$\beta$ -propeller	2	5	160	108	109	82	-	1mda
P-loop NTP hydrolases	3	2	100	88	91	57	39	1gky
$\alpha$ - $\alpha$ superhelix	4	6	136	90	110	44	55	2bct
TIM-barrel	5	1	58	57	39	24	91	1byb
Ferredoxin-like	6	3	135	61	63	70	144	1fxd
Rossmann fold	8	4	55	99	43	56	92	1xel
Ribonucleotide reductase (R1)	100	143	1	-	-	-	35	1rlr
ATPase dom. of HSP90	100	91	2	4	72	73	2	1ah6
Homing endonuclease-like	130	164	3	136	85	175	41	1af5
Aminoacid dehydrogenases; dim. dom.	-	-	4	169	121	3	51	1hup
DNA topo I (N-term)	-	-	175	1	148	126	-	1ois
DNA clamp	130	115	8	2	87	11	60	2pol
Metallothionein	100	14	89	3	33	12	-	1mhu
Phosphoenolpyruvate carboxykinase	130	190	51	28	1	96	169	1ayl
Citrate synthase	81	120	14	8	2	28	51	1csh
N-carbamoylsarcosine amidohydrolase	130	112	9	-	3	138	118	1nba
TBP-like	81	91	46	38	4	75	100	1bvl
5'-3' exonuclease	67	150	32	125	162	1	157	1tfr
$\alpha$ / $\alpha$ toroid	62	132	169	145	114	2	100	1gai
Cyclin-like	20	61	20	15	129	4	-	1vin
ATPase domain of GroEL	36	34	183	143	61	151	1	1aon
Head domain of GrpE	130	135	196	31	165	165	3	1dkg
HSP70 (C-term)	31	10	16	11	56	117	4	1dkz

Common Folds



Folds  
that  
change  
a lot  
in  
frequency

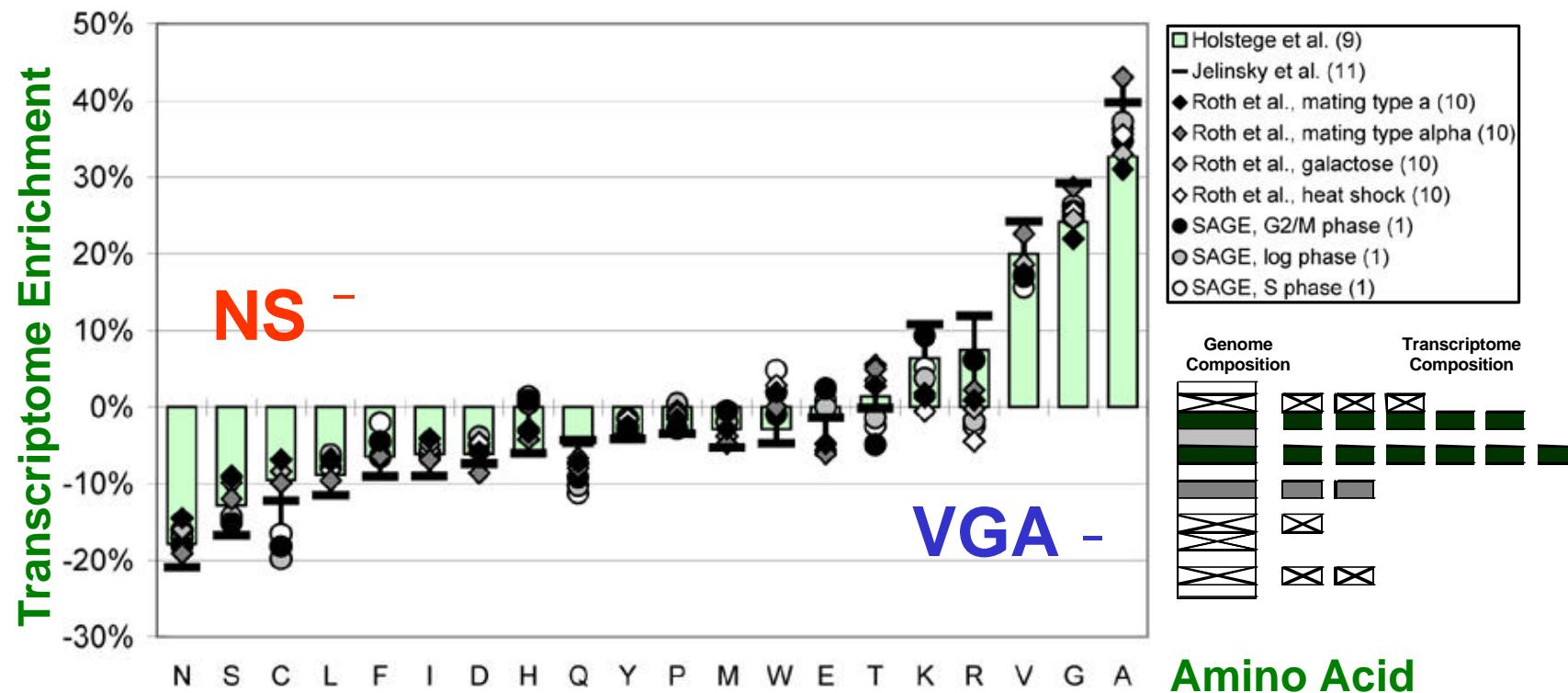
are  
not  
common



Changing Folds

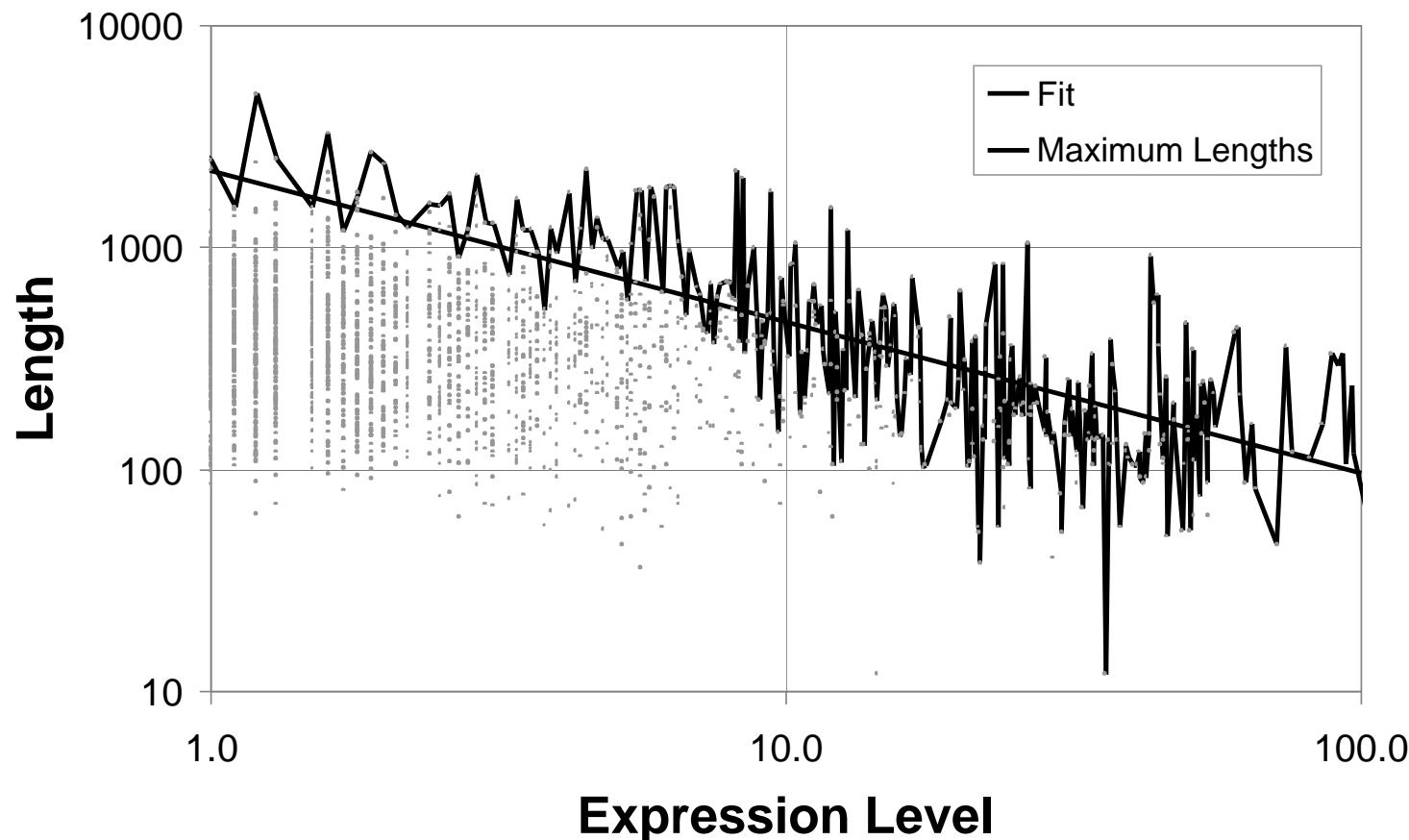
# Composition of Genome vs. Transcriptome

	$\sum_{\text{orf } i} n_i(F)$	$\sum_F \sum_{\text{orf } i} n_i(F)$	$G(F)$	$\sum_{\text{orf } i} e_i n_i(F)$	$\sum_F \sum_{\text{orf } i} e_i n_i(F)$	$T(F)$	$D(F)$
<b>Feature F is Amino acids, in particular Ala</b>	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
<b>Spec. Num.</b>	<b>141890</b>	<b>2574876</b>	<b>5.5%</b>	<b>347807</b>	<b>4758441</b>	<b>7.3%</b>	<b>32.7%</b>
<b>Feature F is Folds, in particular the TIM-barrel (3.1)</b>	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
<b>Spec. Num.</b>	<b>65</b>	<b>1560</b>	<b>4.2%</b>	<b>389</b>	<b>4709</b>	<b>8.3%</b>	<b>97.8%</b>

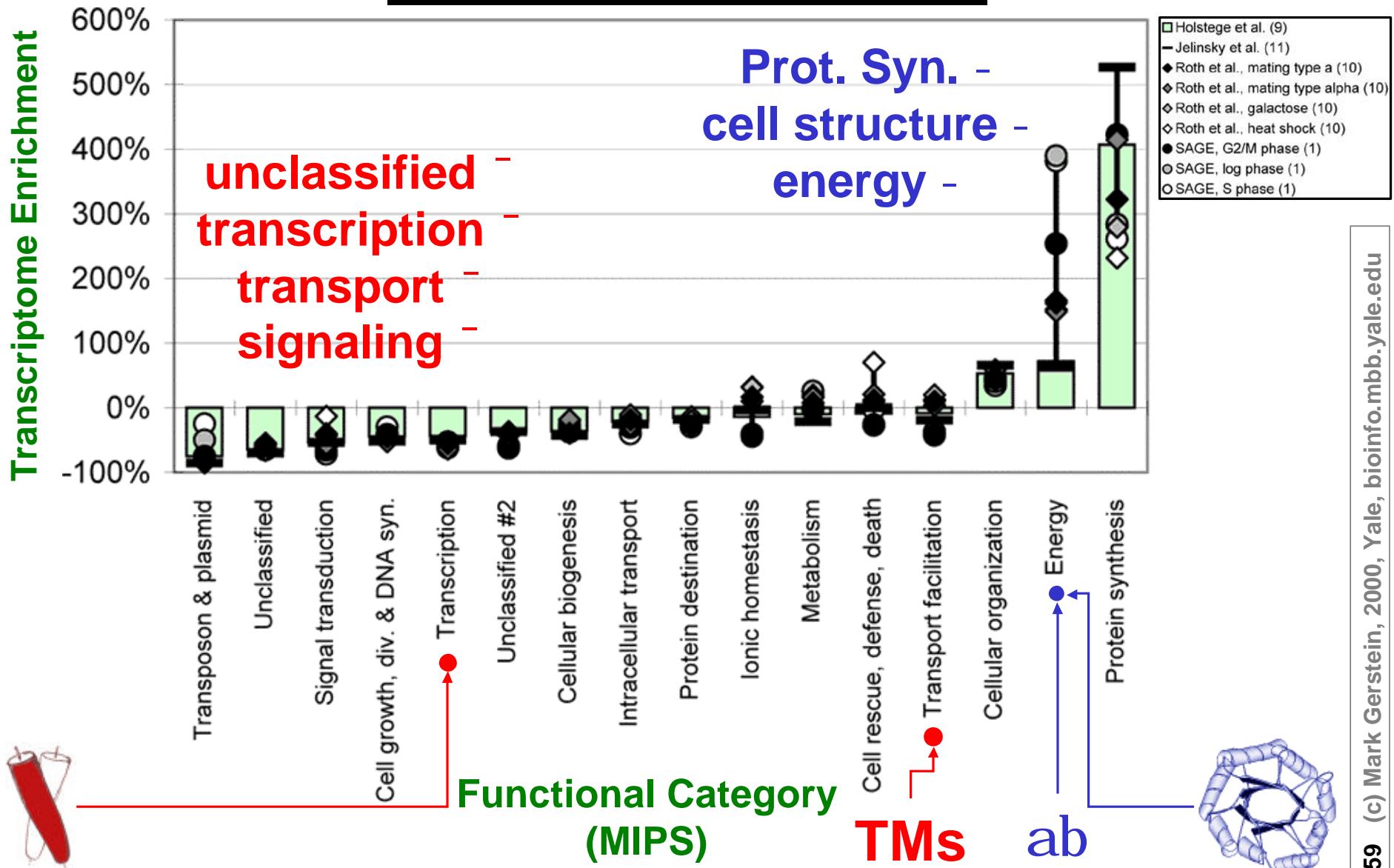


# Relation between Length & Expression

Max Expression (e.g. transcripts/cell)  $\sim$  (Length) $^{-2/3}$   
Shorter proteins can be more highly expressed



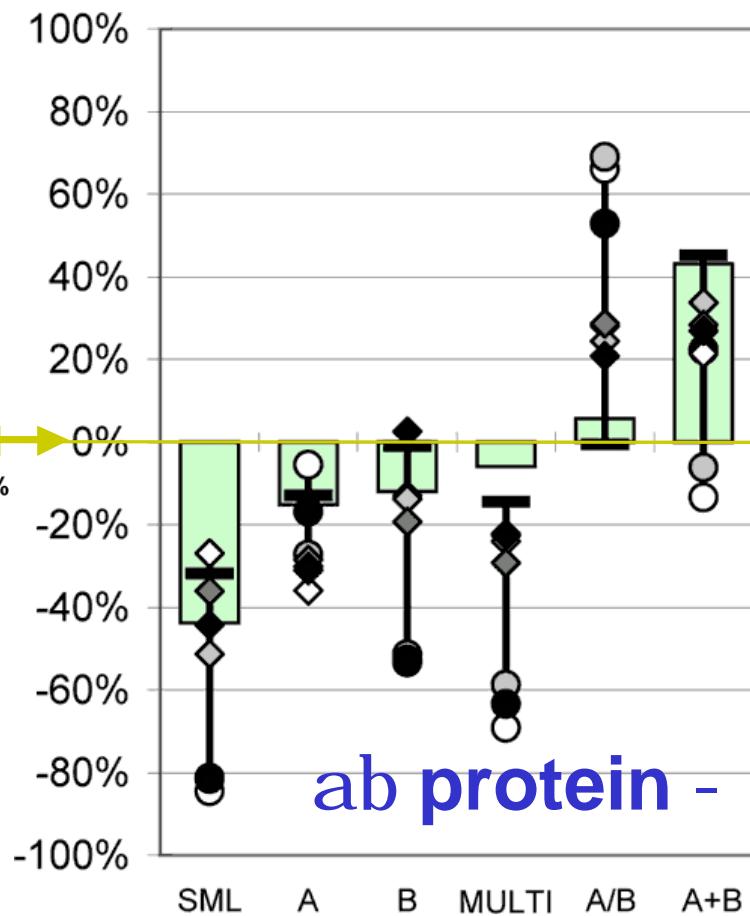
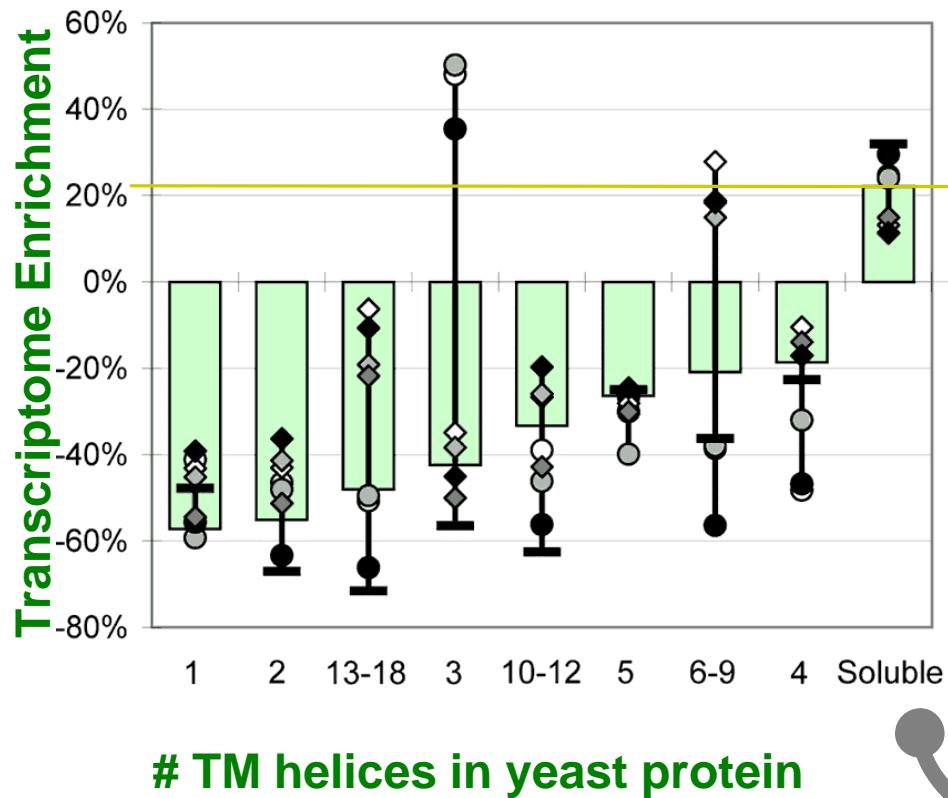
# Composition of Transcriptome in terms of Functional Classes



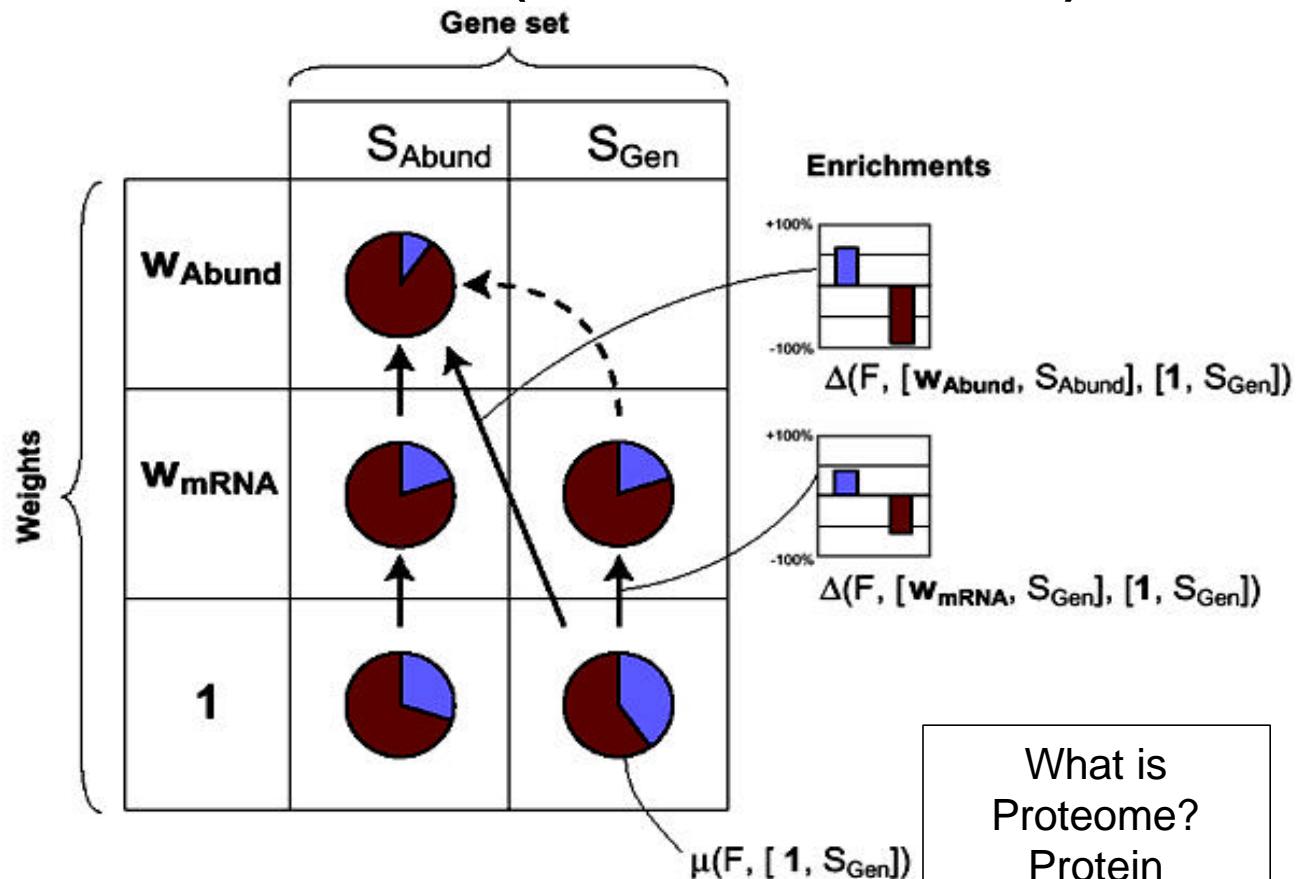
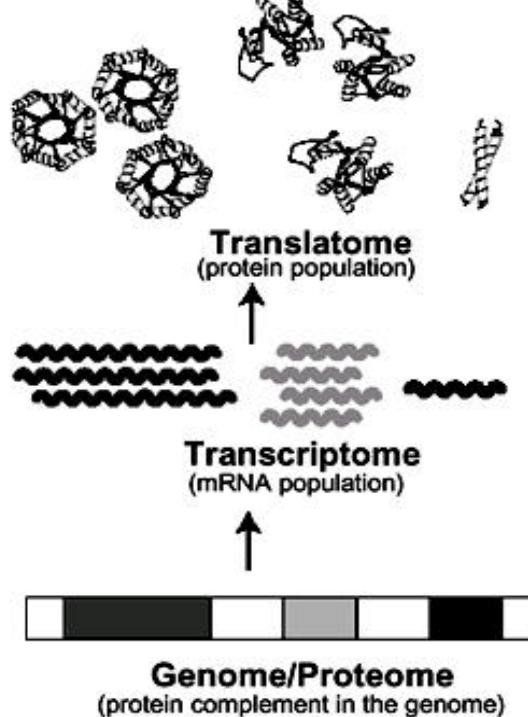
# Composition of Transcriptome in terms of Broad Structural Classes

- 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 -



# Relating the Transcriptome to Cellular Protein Abundance (Translatome)

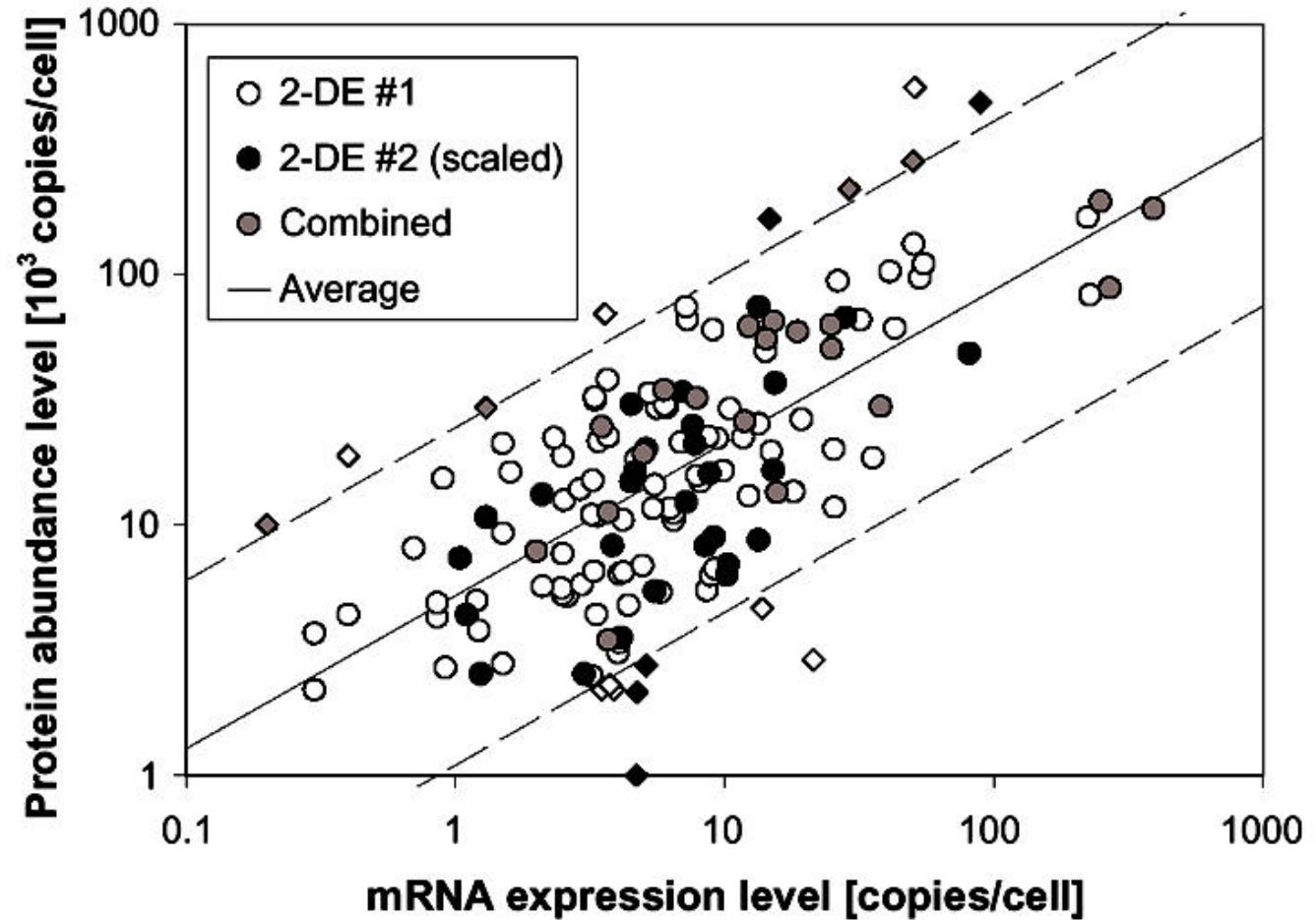


2D-gel electrophoresis Data sets: Futcher (71), Aebersold (156), scaled set with 171 proteins  
New effect is dealing with gene selection bias

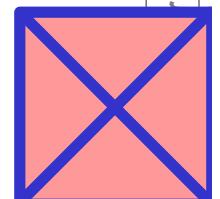
What is Proteome?  
Protein complement in genome or cellular protein population?

# mRNA and protein abundance related, roughly

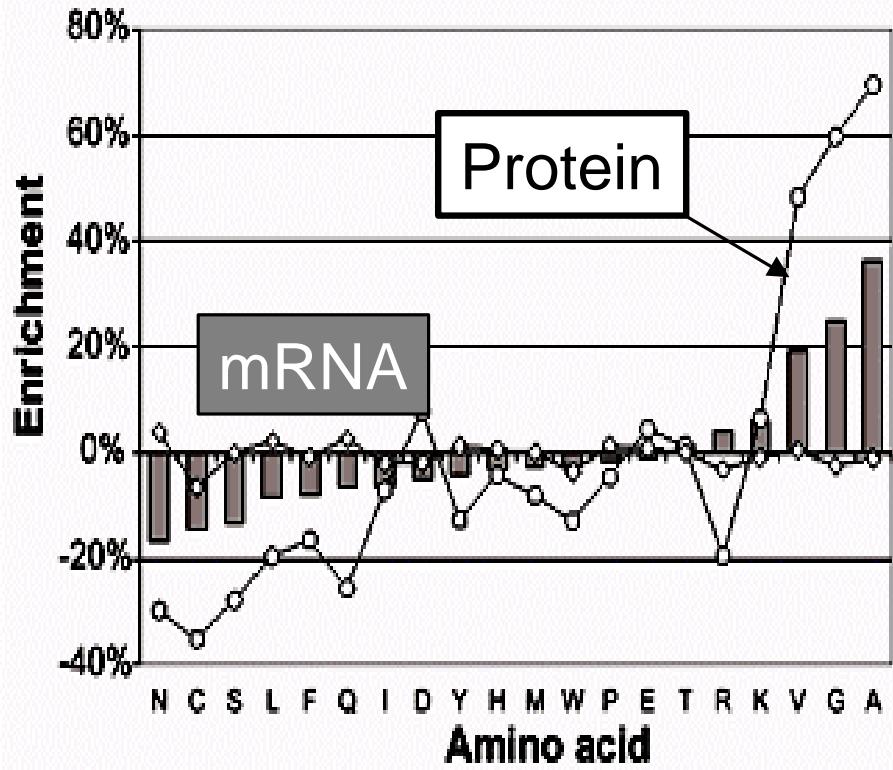
~150 protein abundance values from merging results of 2D gel expts. of Aebersold & Futcher



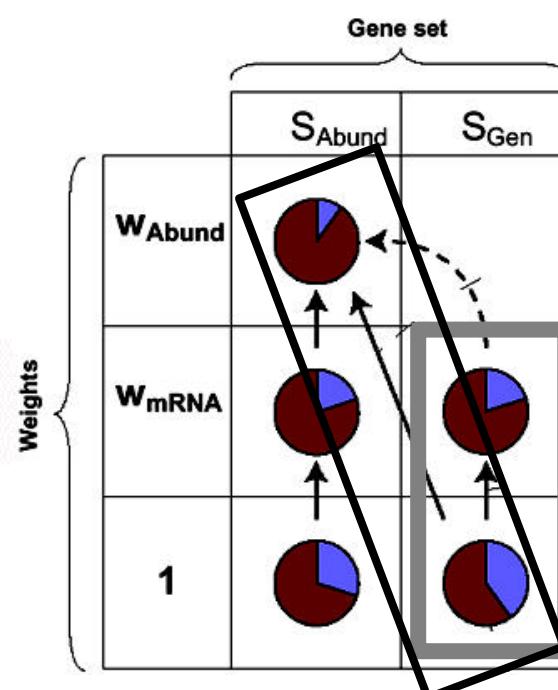
mRNA values for same 150 genes from merging and scaling 6 yeast expressions



# Amino Acid Enrichment

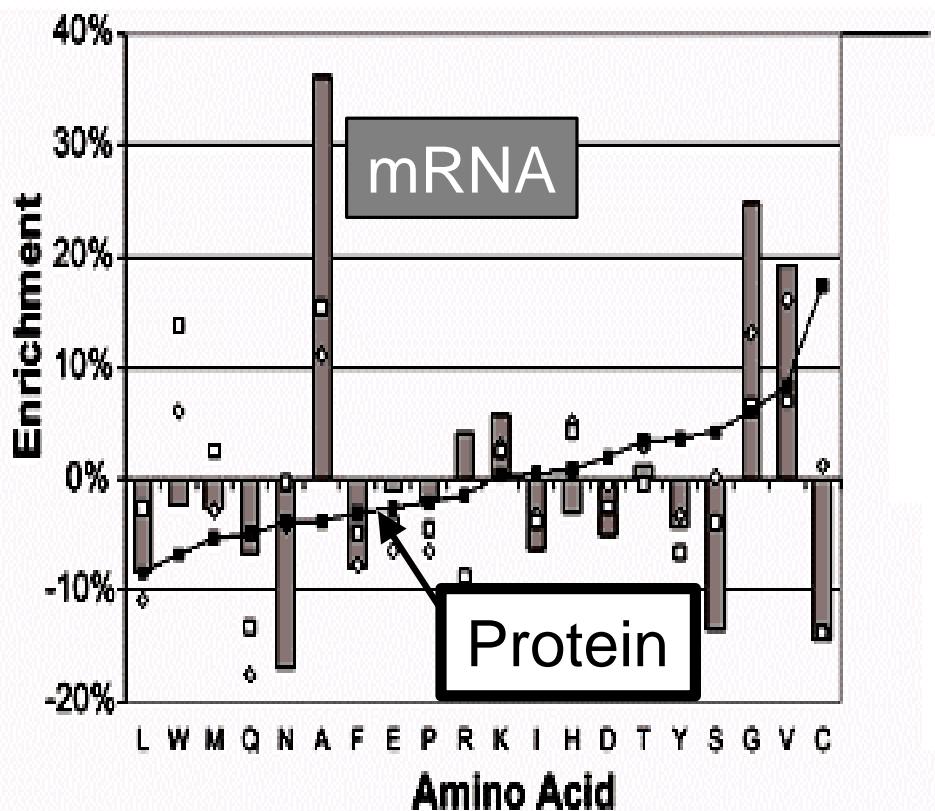


Legend:  
■  $\Delta(\text{aa}, [W_{\text{mRNA}}, S_{\text{Gen}}], [1, S_{\text{Gen}}])$   
○  $\Delta(\text{aa}, [W_{\text{Abund}}, S_{\text{Gen}}], [1, S_{\text{Gen}}])$   
◊  $\Delta(\text{aa}, [W_{\beta\text{-Gal}}, S_{\beta\text{-Gal}}], [1, S_{\text{Gen}}])$

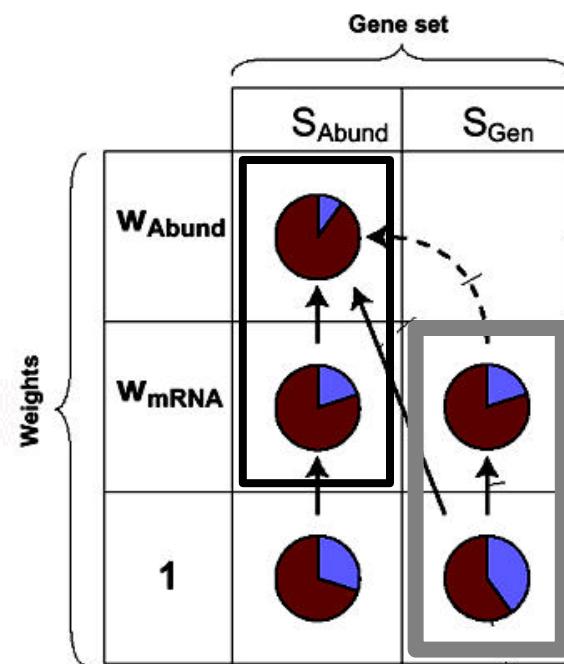


Simple story is translatome is enriched in same way as transcriptome

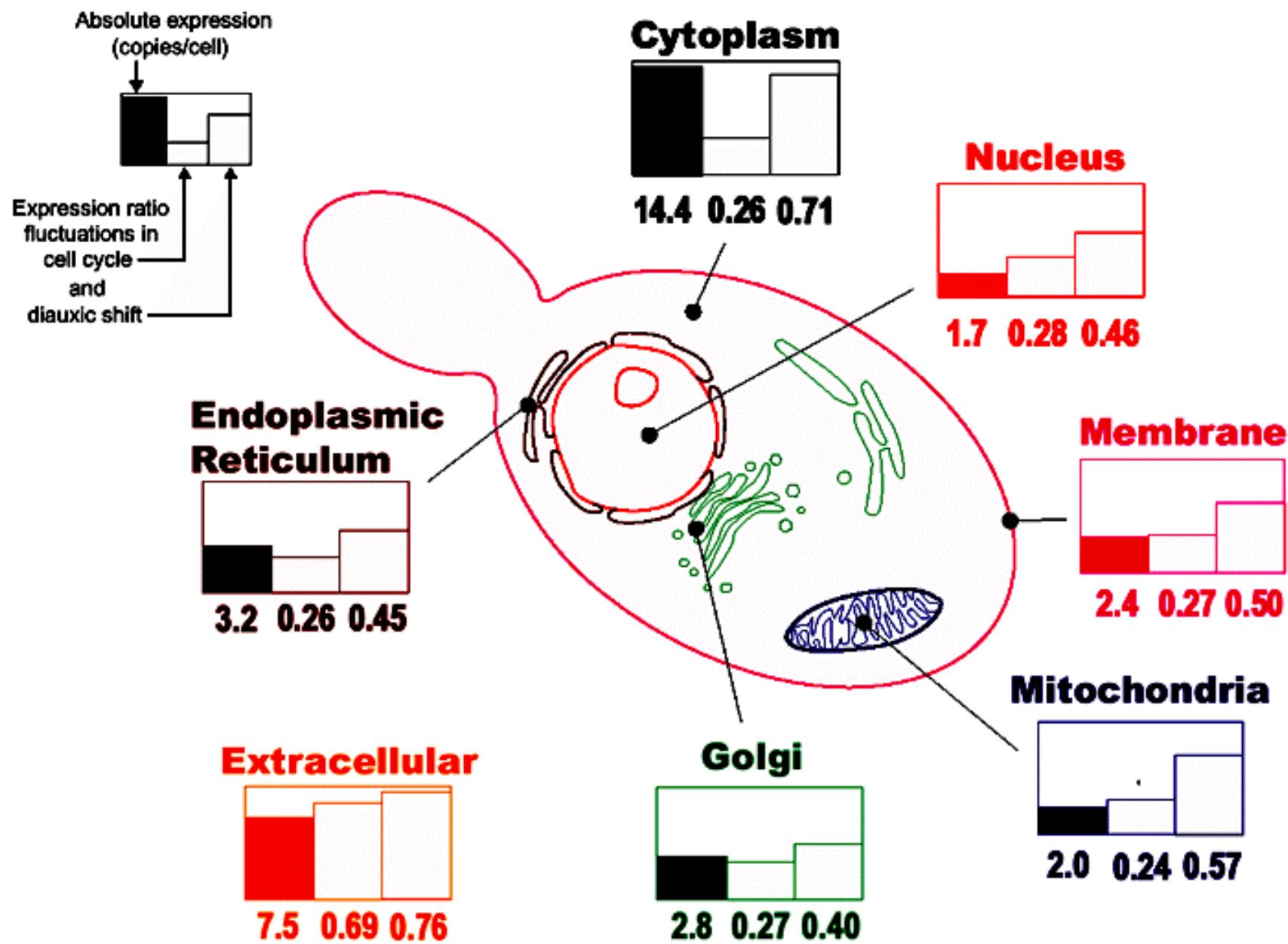
# Amino Acid Enrichment – Complexities

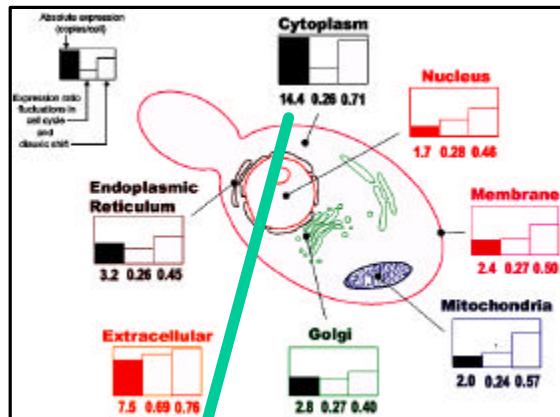


- $\Delta(\text{aa}, [W_{\text{mRNA}}, S_{\text{Gen}}], [1, S_{\text{Gen}}])$
- $\Delta(\text{aa}, [W_{\text{Abund}}, S_{\text{Abund}}], [1, S_{\text{Abund}}])$
- $\Delta(\text{aa}, [W_{\text{mRNA}}, S_{\text{Abund}}], [1, S_{\text{Abund}}])$
- $\Delta(\text{aa}, [W_{\text{Abund}}, S_{\text{Abund}}], [W_{\text{mRNA}}, S_{\text{Abund}}])$

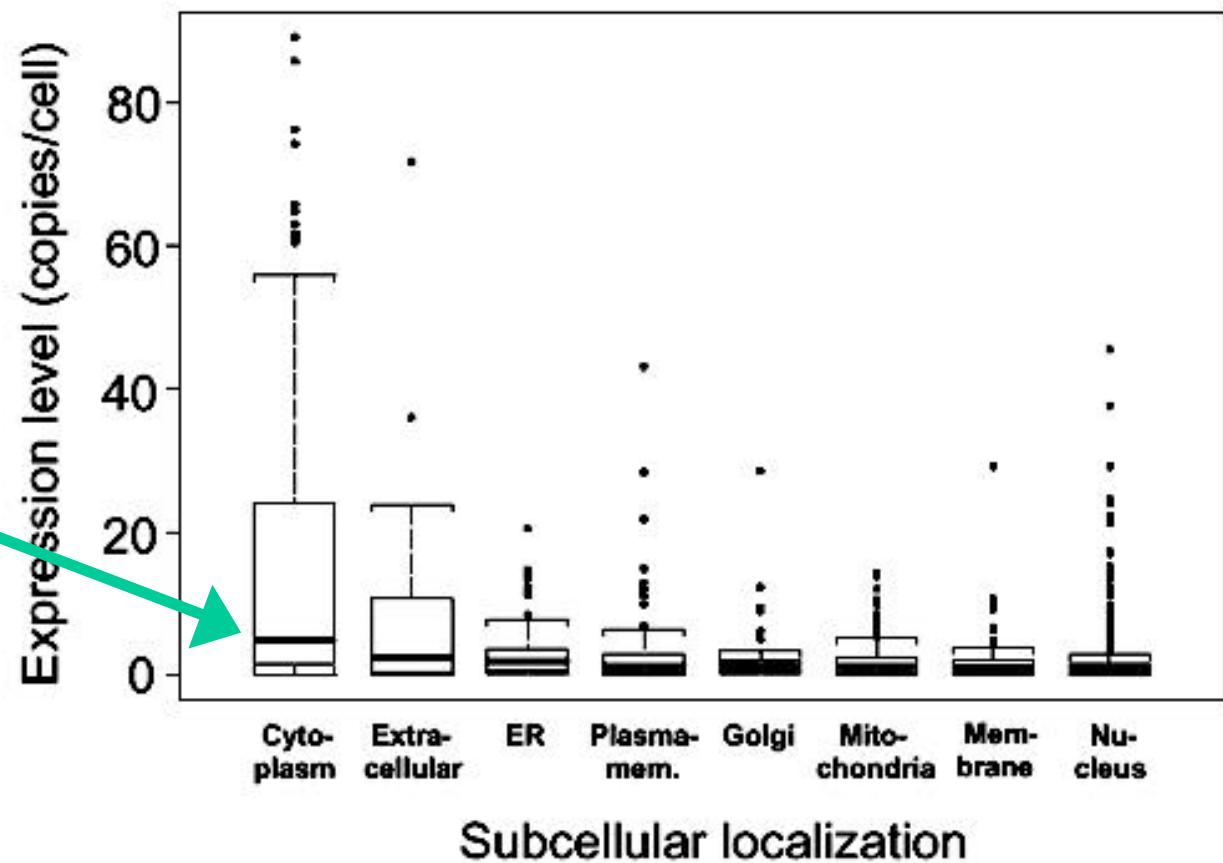
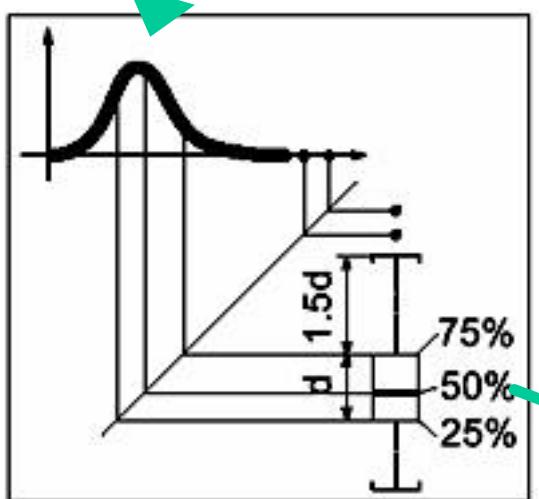


# Expression Level is Related to Localization



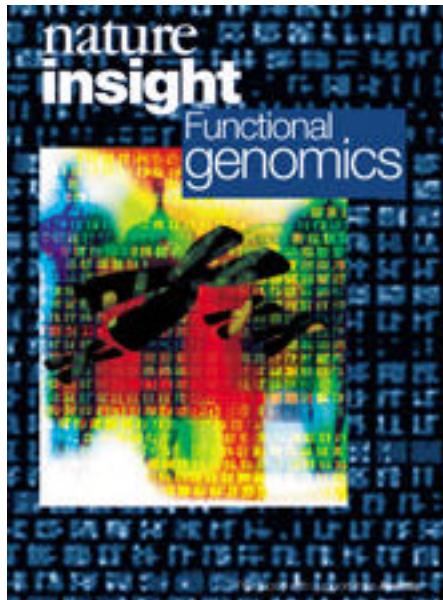


# Distributions of Expression Levels



~6000 yeast genes  
with expression levels

but only ~2000 with localization....



insight review articles

## Genomics, gene expression and DNA arrays

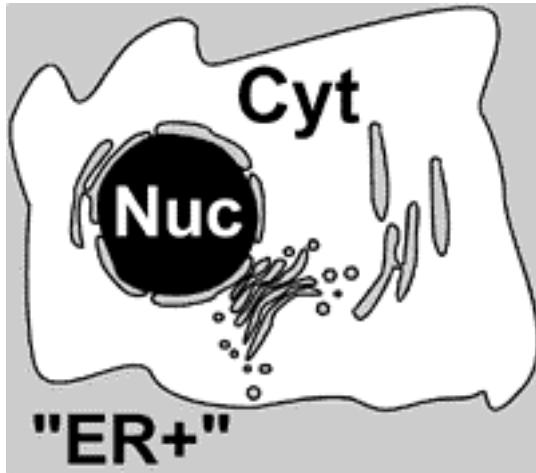
David J. Lockhart & Elizabeth A. Winzeler

*Genomics Institute of the Novartis Research Foundation, 3115 Merryfield Row, San Diego, California 92121, USA*

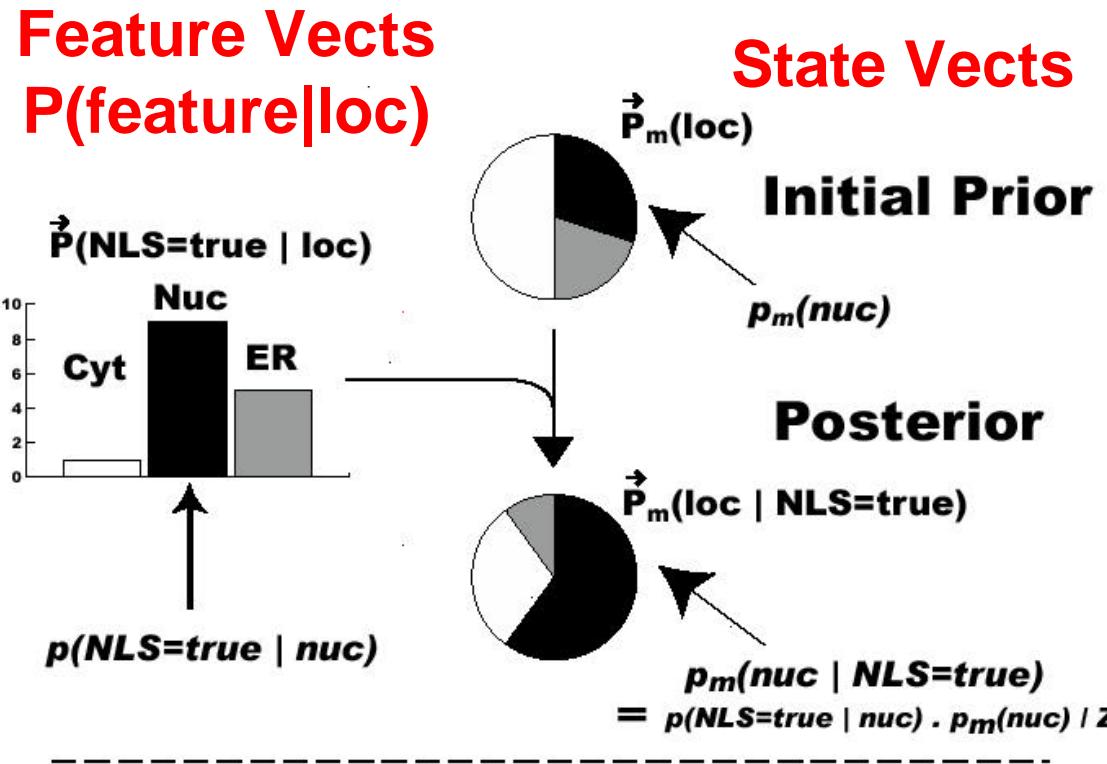
Experimental genomics in combination with the growing body of sequence information promise to revolutionize the way cells and cellular processes are studied. Information on genomic sequence can be used experimentally with high-density DNA arrays that allow complex mixtures of RNA and DNA to be interrogated in a parallel and quantitative fashion. DNA arrays can be used for many different purposes, most prominently to measure levels of gene expression (messenger RNA abundance) for tens of thousands of genes simultaneously. Measurements of gene expression and other applications of arrays embody much of what is implied by the term (genomics); they are broad in scope, large in scale, and take advantage of all available sequence information for experimental design and data interpretation in pursuit of biological understanding.

# Bayesian System for Localizing Proteins

loc=



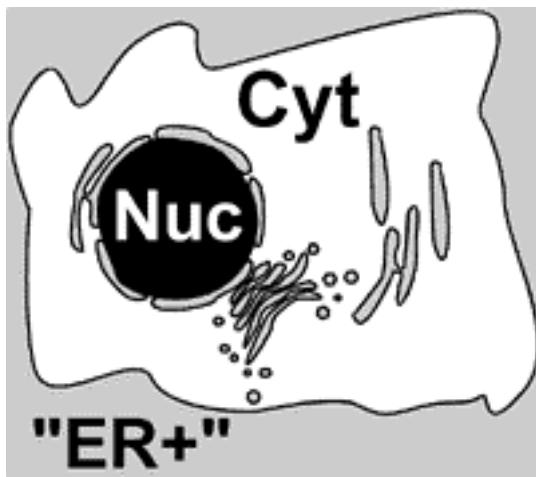
Represent localization of each protein by the state vector  $P(\text{loc})$  and each feature by the feature vector  $P(\text{feature}|\text{loc})$ . Use Bayes rule to update.



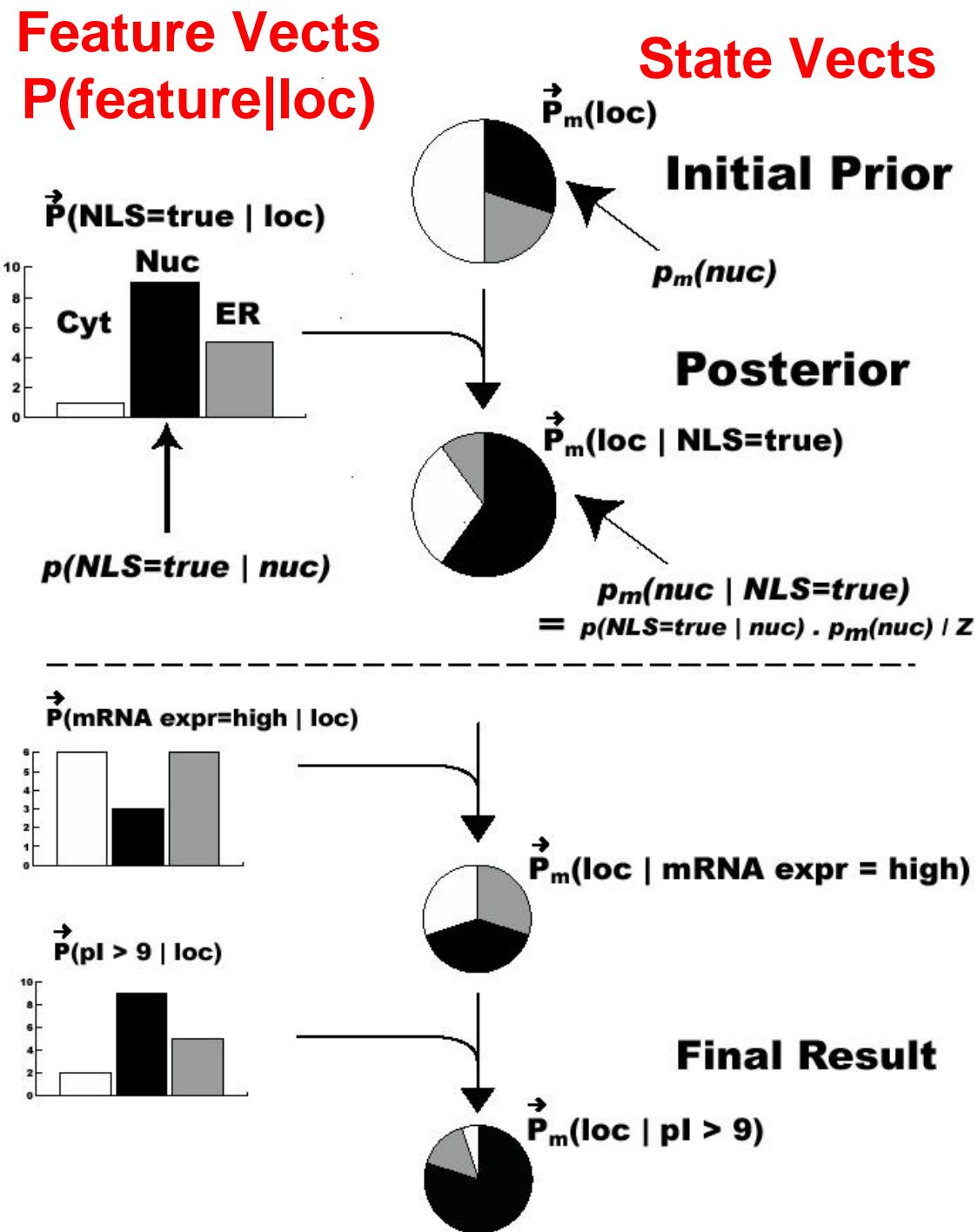
18 Features: Expression Level  
(absolute and fluctuations), signal  
seq., KDEL, NLS, Essential?, aa  
composition

# Bayesian System for Localizing Proteins

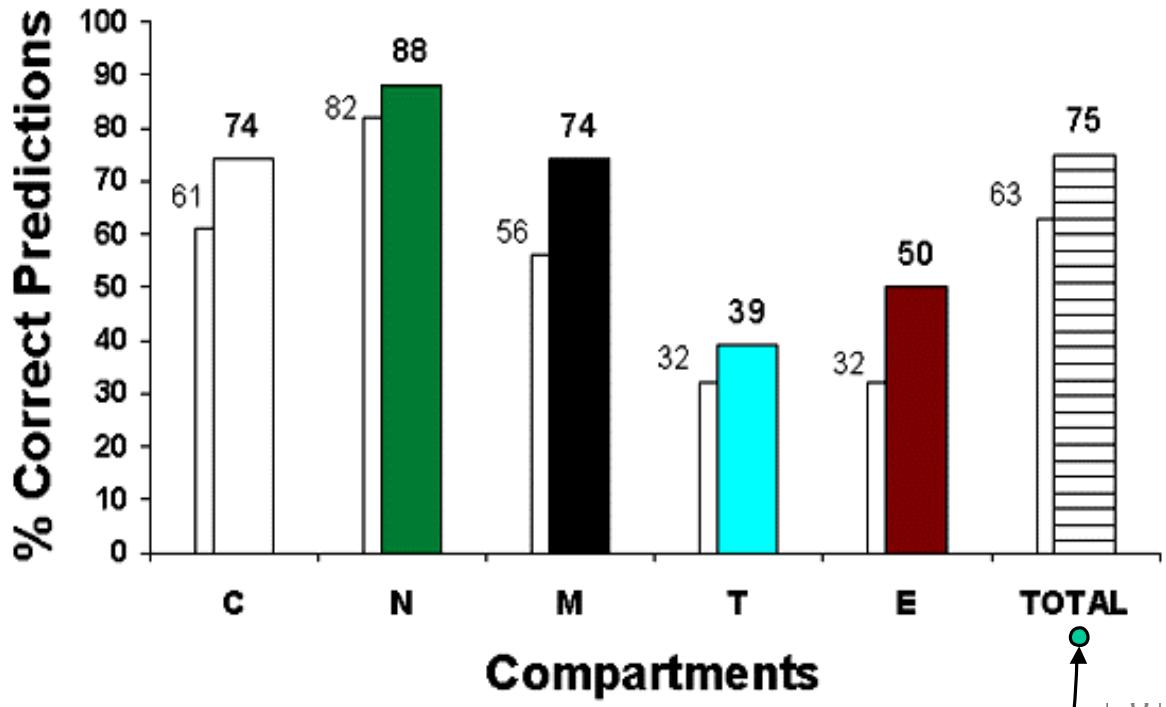
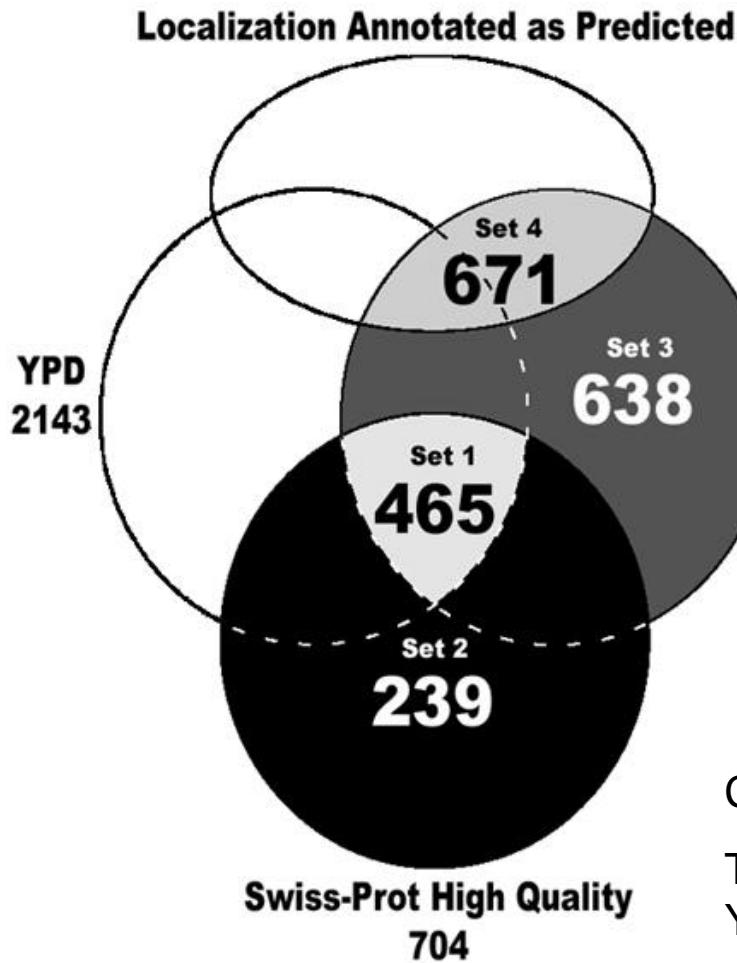
**loc=**



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: 75% with cross-validation**

Carefully clean training dataset to **avoid circular logic**

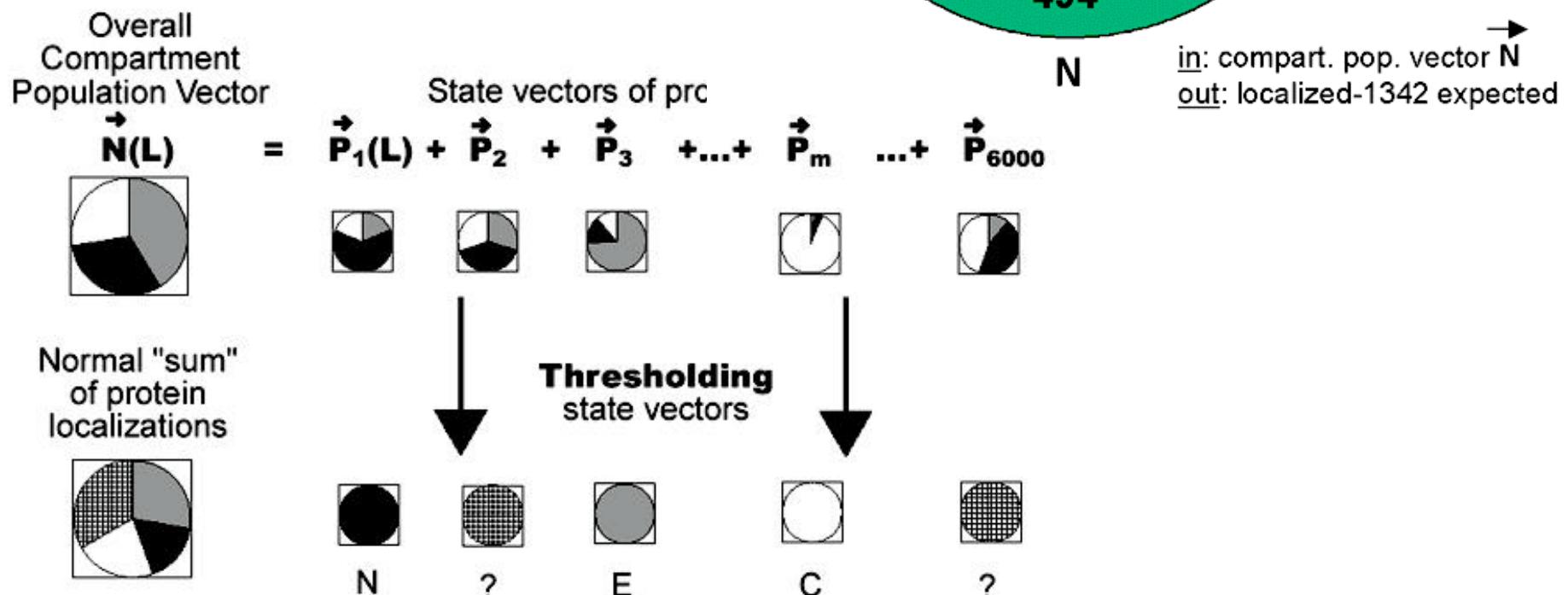
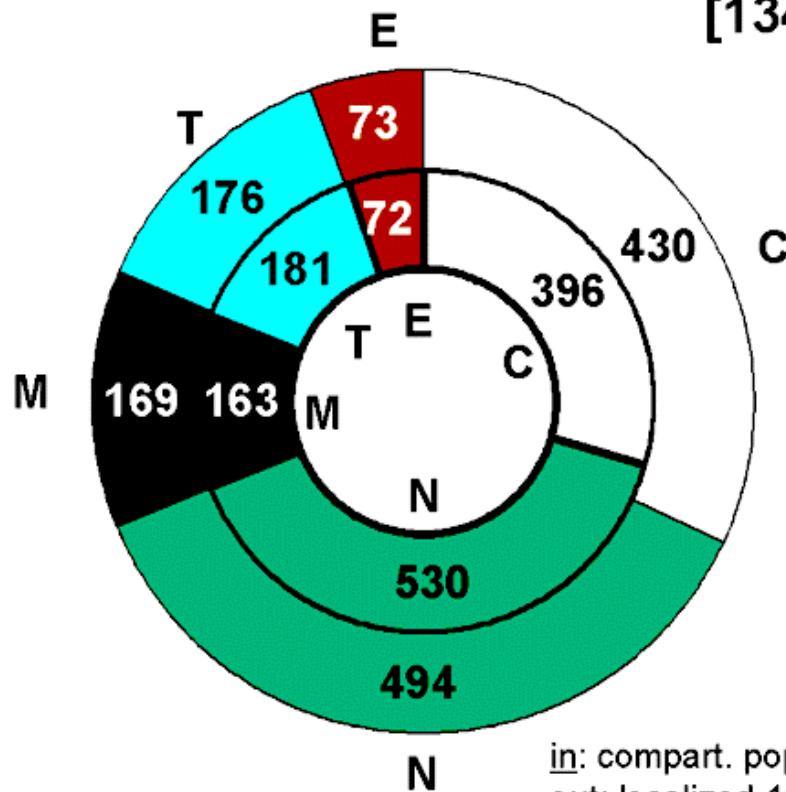
Testing, training data, Priors: ~2000 proteins from YPD, MIPS, SwissProt, Snyder Lab

[1342]

# Results on Testing Data #2

## Compartment

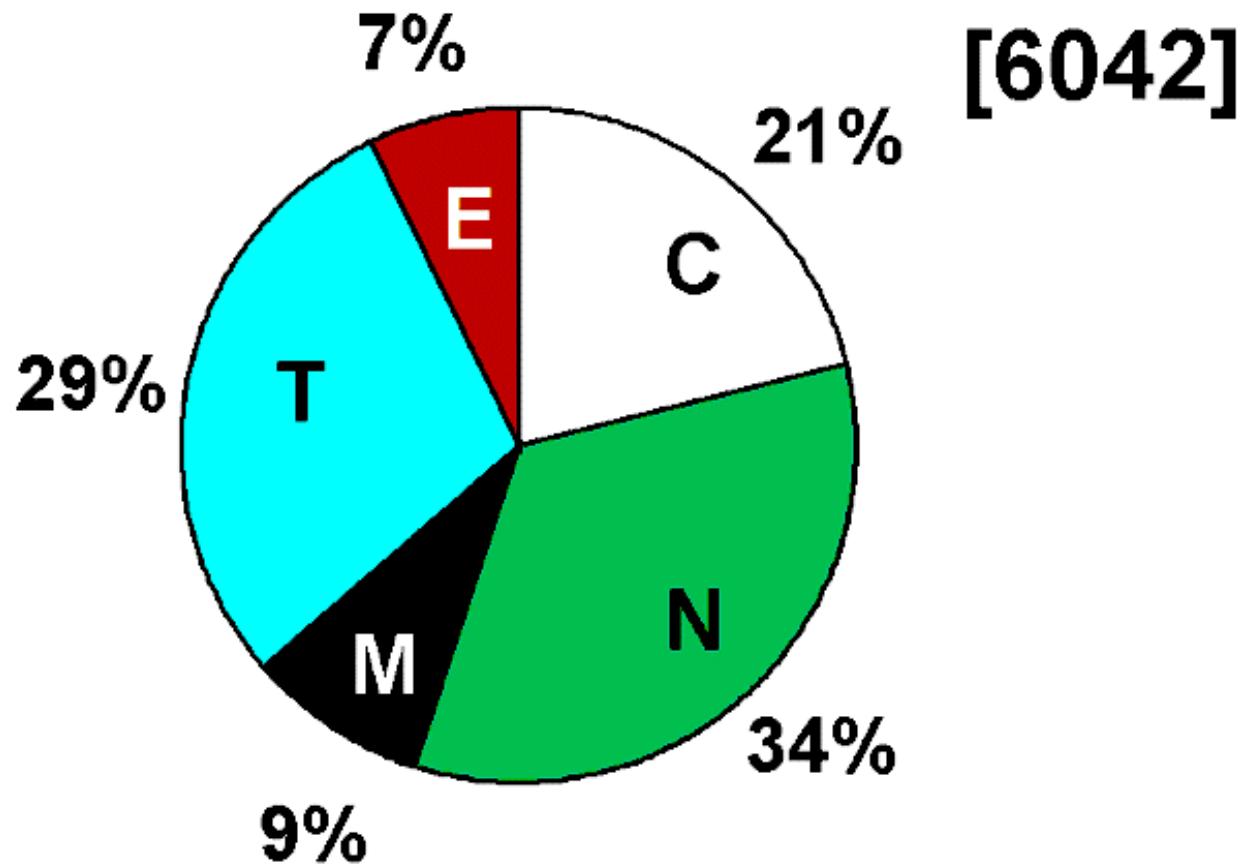
Populations. Like QM,  
directly sum state vectors  
to get population. Gives  
**96%** pop. similarity.



# Extrapolation to Compartment

## Populations of Whole Yeast Genome:

~4000 predicted + ~2000 known



# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

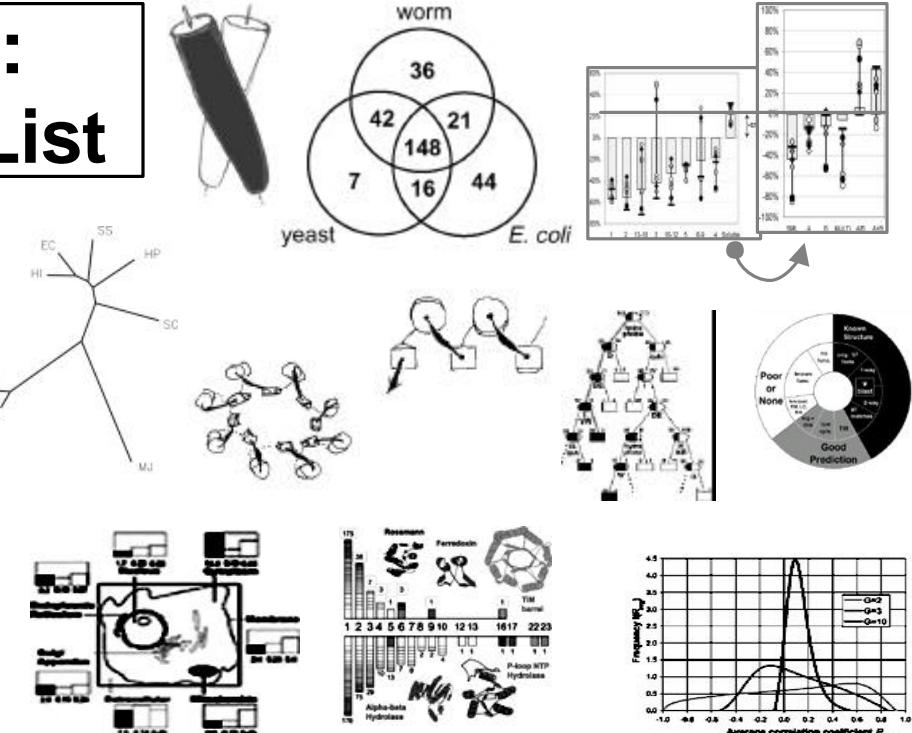
- **Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

- **Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

## 4 Using Folds to Interpret Expression Data

**Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

- **Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.

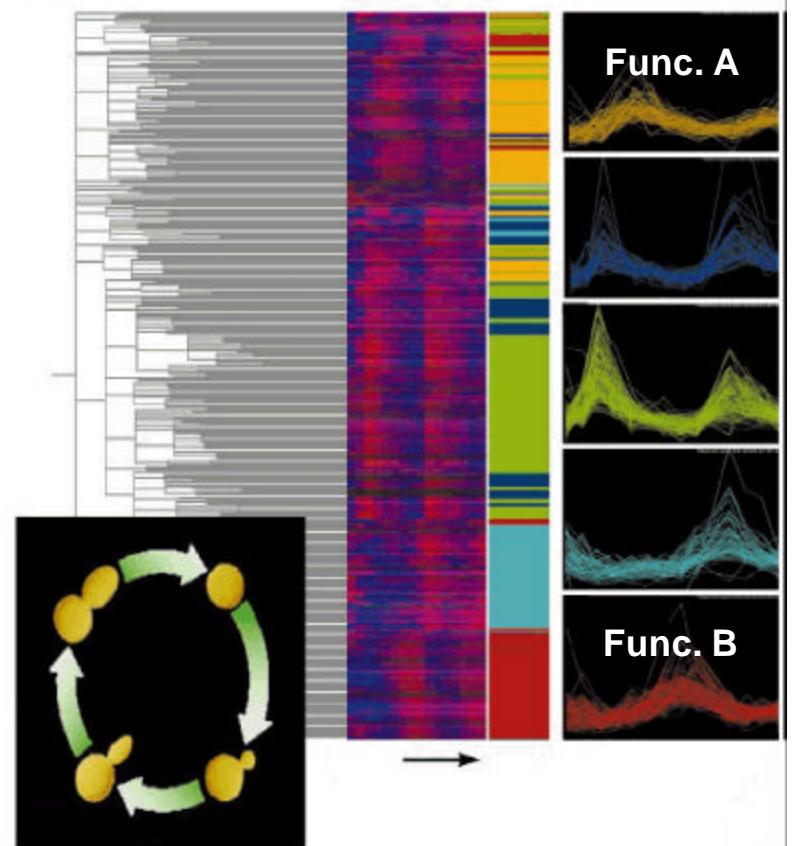
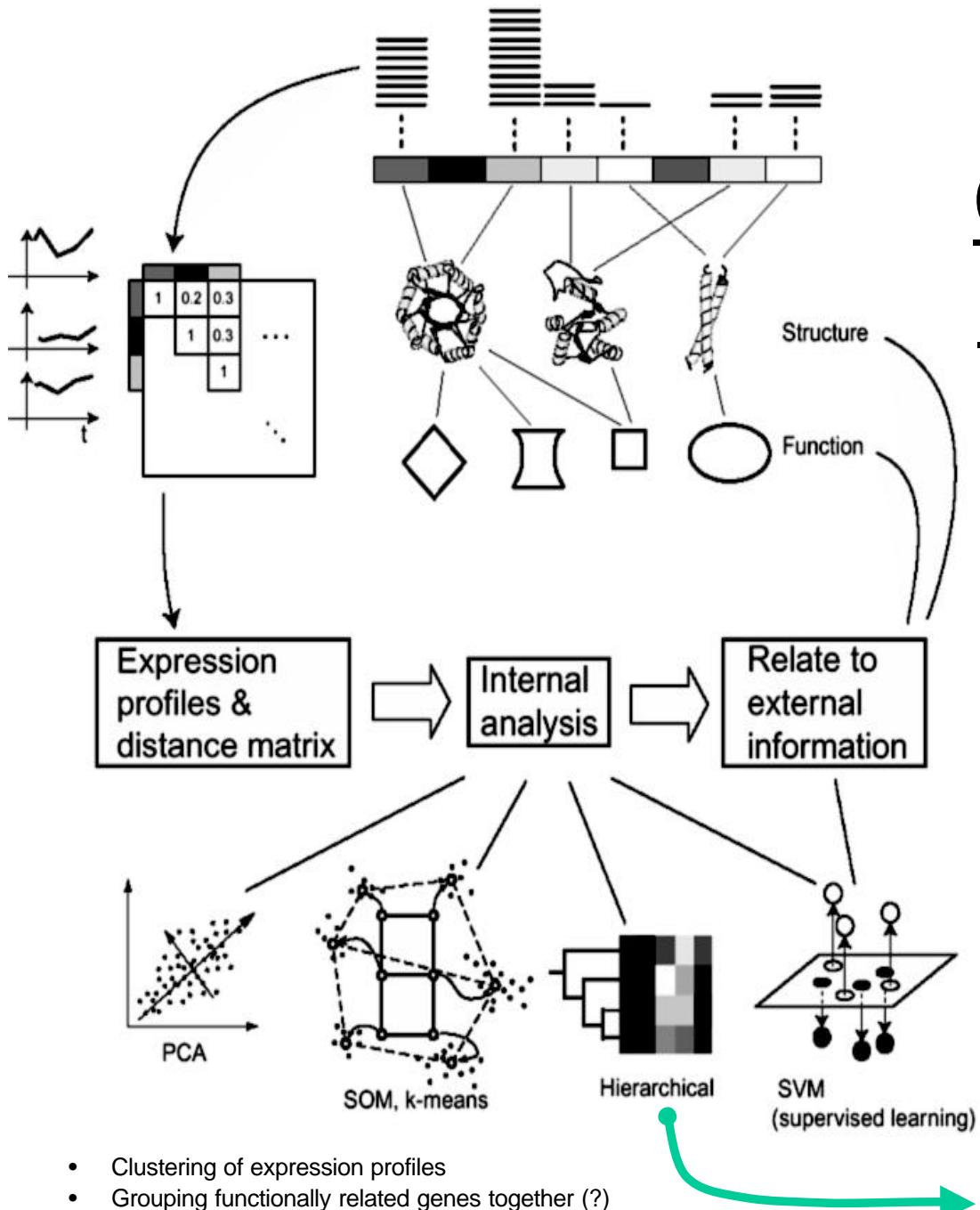


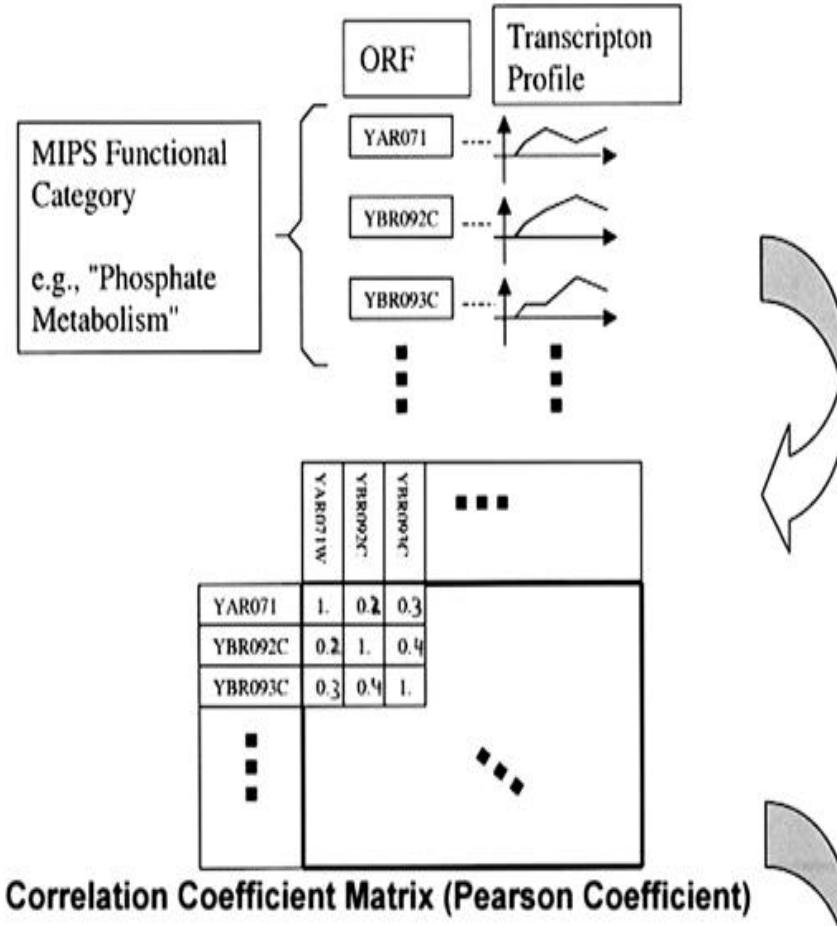
**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

**bioinfo.mbb.yale.edu**

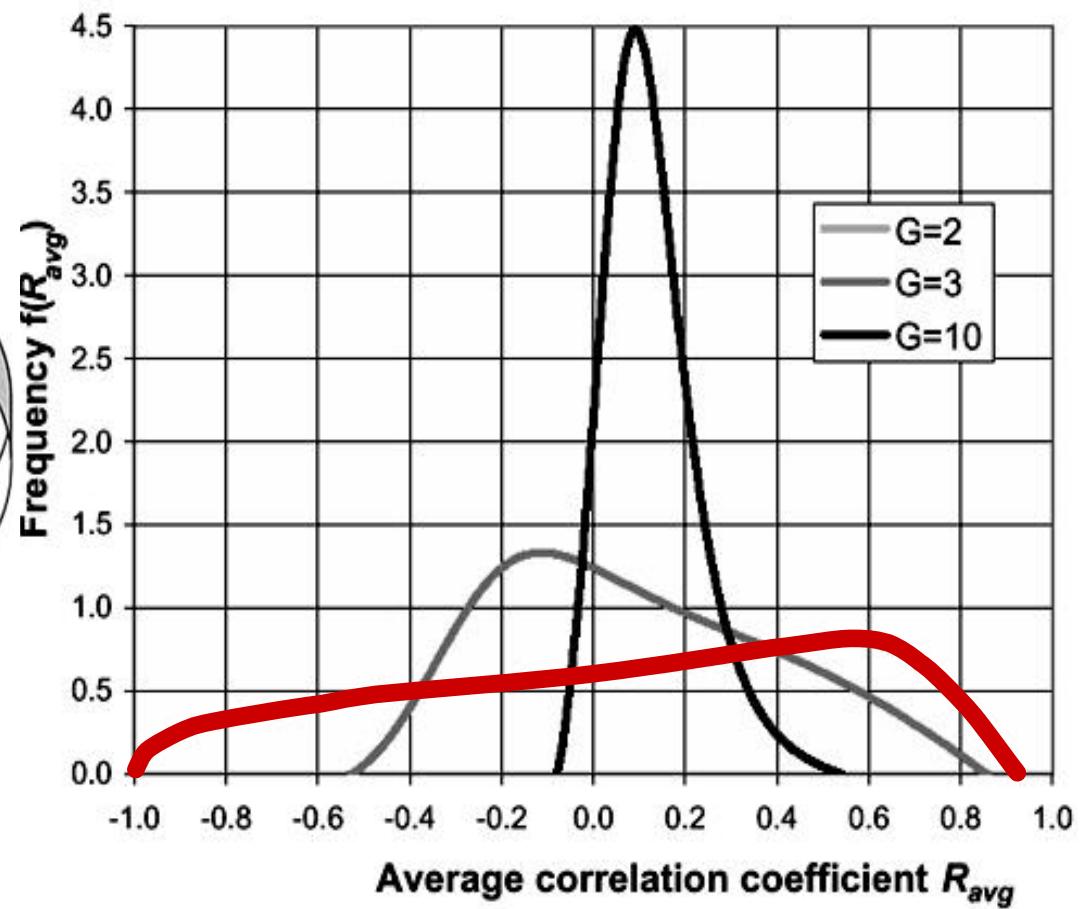
# Do Expression Clusters Relate to Protein Function?

Can they predict functions?

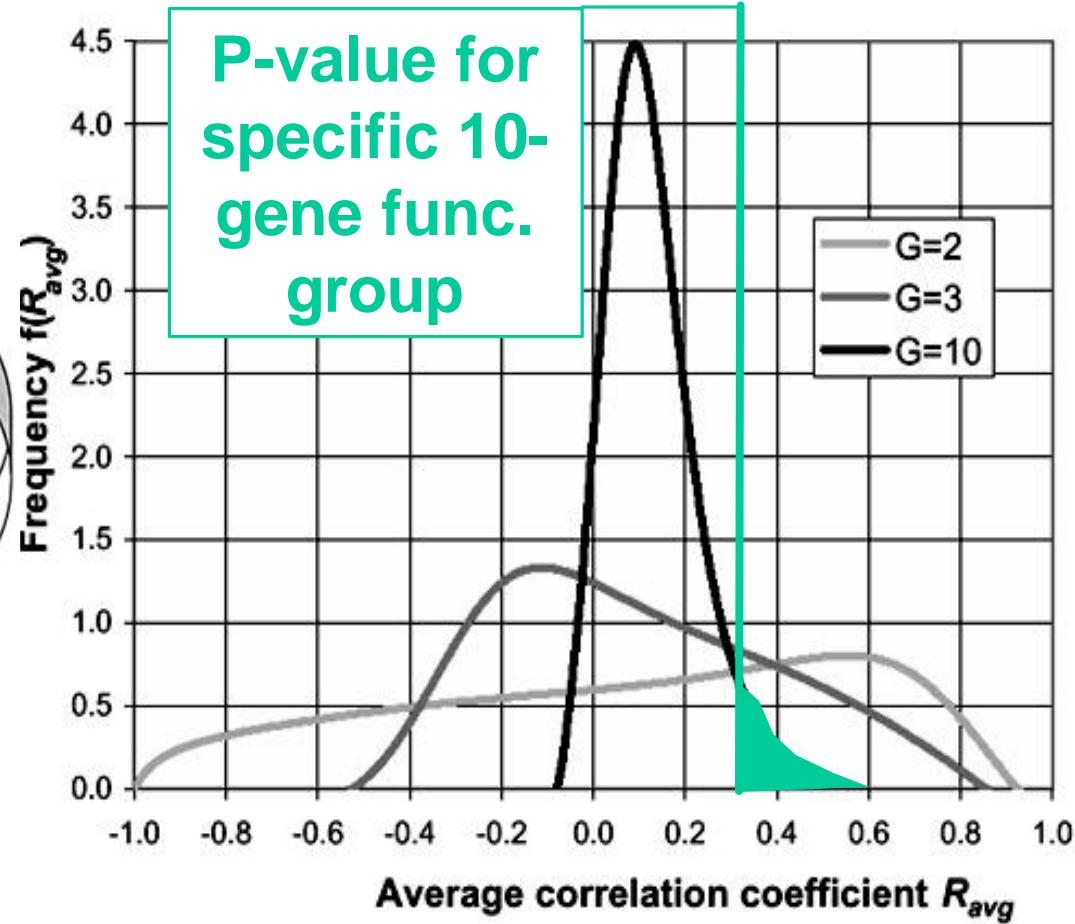
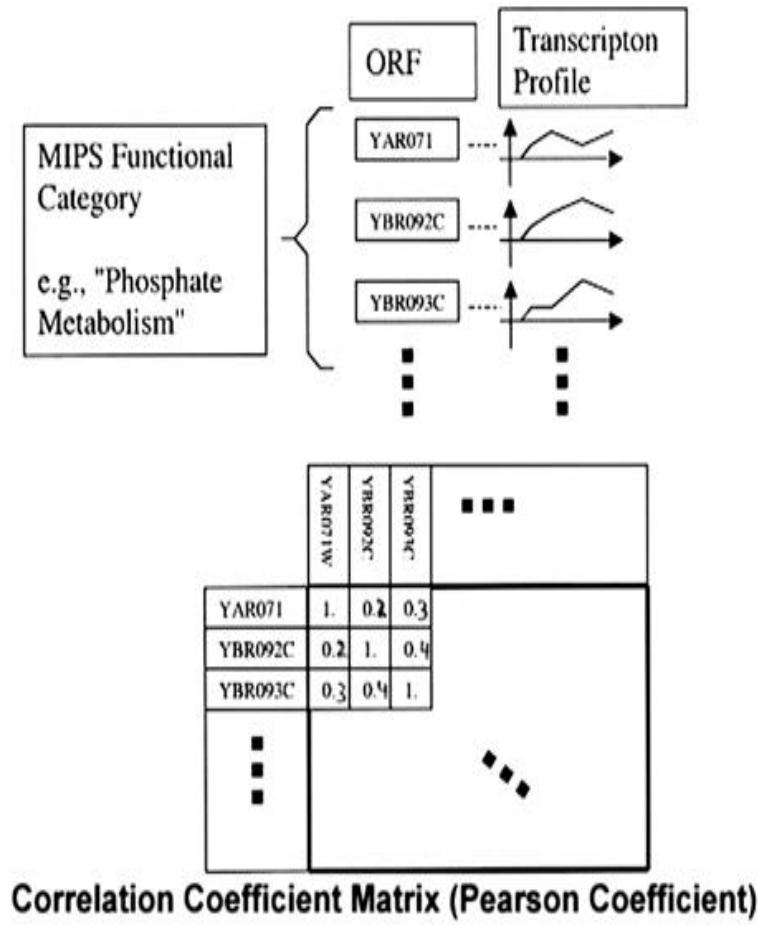




# Distributions of Gene Expression Correlations, for All Possible Gene Groupings



# Distributions of Gene Expression Correlations, for All Possible Gene Groupings 2



Average Correlation Coefficient for Group of Genes

Sample for Diauxic shift Expt. (Brown),

$$\text{Ex. } R_{avg,G=3} = \frac{[ R(\text{gene-1,gene-3}) + R(\text{gene-1,gene-4}) + R(\text{gene-5,gene-7}) ] / 3}{3}$$

MIPS category	Experiment			
	Cell Cycle (CDC28)	Cell cycle (CDC15)	Diauxic shift	Sporulation
Cell growth, division & DNA syn.	>4	>4	>4	>4
Protein synthesis	>4	>4	>4	>4
Transcription	>4	>4	>4	1.6
Cellular organization	>4	>4	0.3	0.3
Energy	>4	>4	0.1	0.9
Cell rescue, defense, death	>4	>4	0	0
Intracellular transport	>4	>4	0	0
Ionic homeostasis	>4	>4	0	0.8
Metabolism	>4	>4	0	0
Transport facilitation	>4	>4	0	0
Signal transduction	2.5	1.6	0.1	0.6
Unclassified	2.3	>4	0	0
Cellular biogenesis	2.0	>4	0.4	0.2
Protein destination	0.3	>4	0.2	0.6
Retrotransposon & plasmid	0	2.8	1.9	1.0

	Fraction of significant groups				Total # groups
	CDC28	CDC15	Diauxic Shift	Sporulation	
MIPS 1	63%	81%	19%	13%	16
MIPS 2	50%	63%	17%	13%	102
MIPS 3	23%	33%	5%	4%	73
"Energy" (2 <sup>nd</sup> level)	40%	60%	20%	0%	10
SOM	93%	-	-	-	30
Hierarch. Clustering		80%			25

MIPS category	Experiment			
	Cell Cycle (CDC28)	Cell cycle (CDC15)	Diauxic shift	Sporulation
Respiration	>4	>4	>4	3.4
TCA pathway	>4	>4	>4	0.6
Glycogen, trehalose metabolism	>4	>4	1.2	0.7
Glycolysis	>4	>4	0.9	2.1
Gluconeogenesis	3.7	>4	0.1	1.7
Glyoxylate cycle	1.6	0.7	3.0	2.3
Pentose-phosphate pathway	1.5	0.8	0	0.6
Fermentation	1.3	>4	0	2.2
Other energy generation activities	0.7	0.1	0.1	0.2
Beta-oxidation of fatty acids	0.5	0.4	0.4	0.2

Correlation:

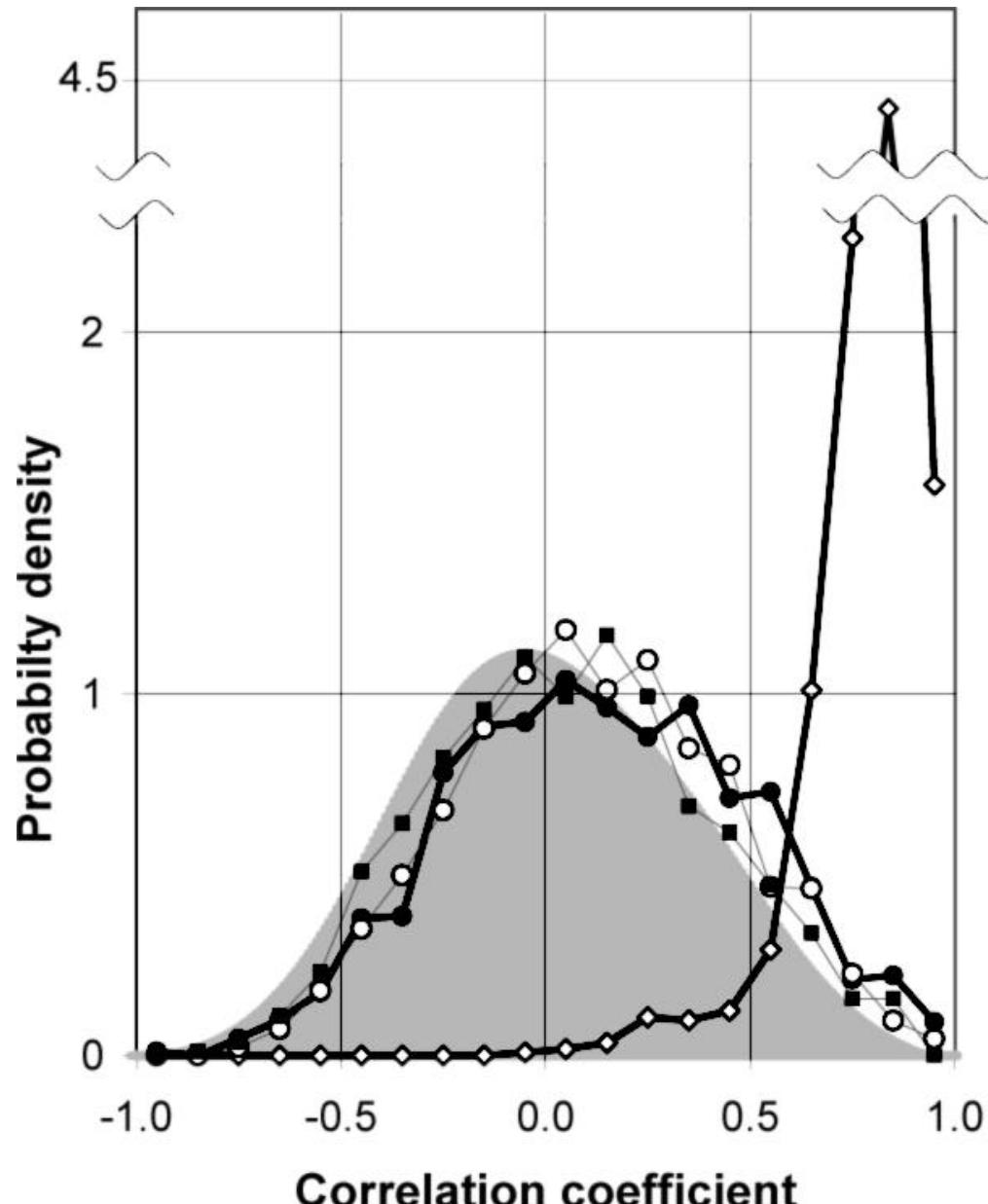
Always Significant

Sometimes Significant (depends on expt.)

Never Significant

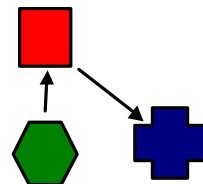
Based on Distributions,  
Correlation of  
Established Functional  
Categories, Computer  
Clusterings

# Protein-Protein Interactions & Expression



between selected expression  
timecourses in CDC28 expt. (Davis)

Use same formalism to assess  
how closely related expression  
timecourses to sets of known p-p  
interactions

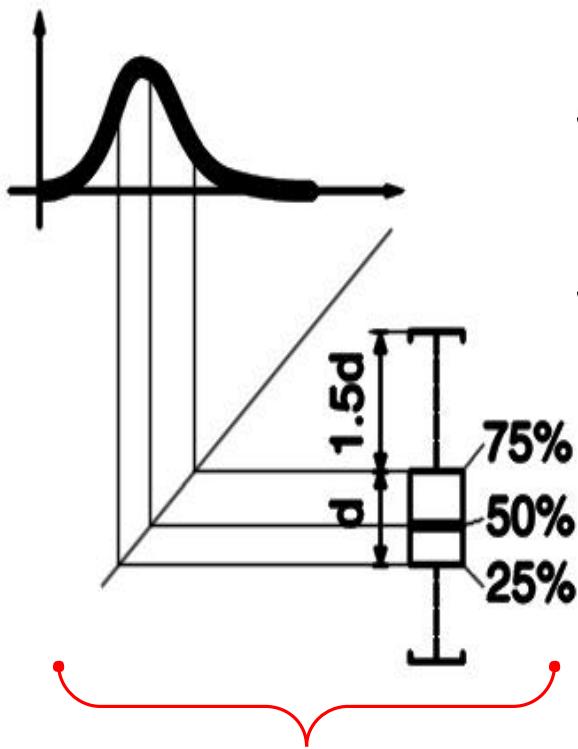


## Sets of interactions

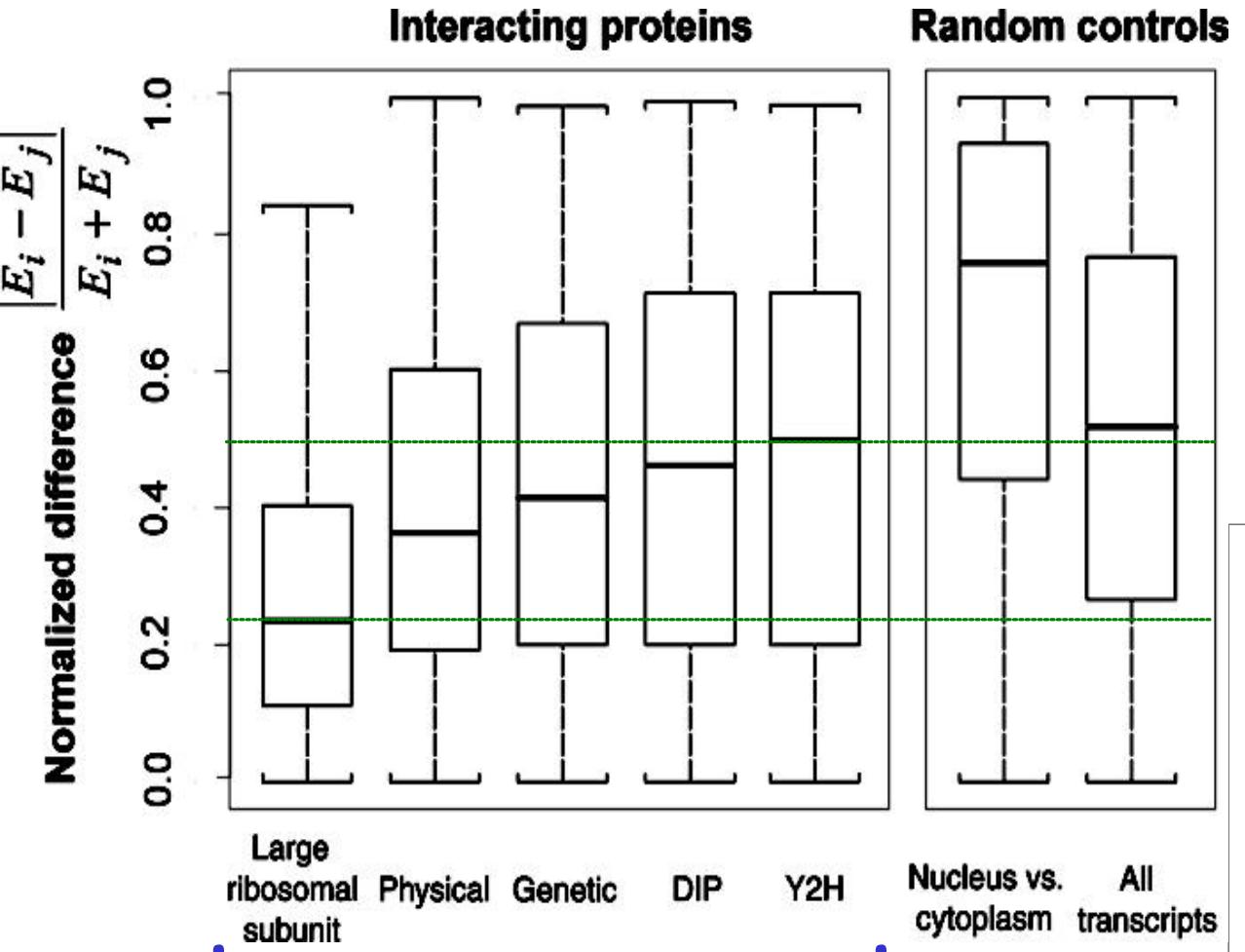
Random (cell cycle CDC28) (all pairs)  
(control)

physical (from MIPS)  
genetic (Uetz et al.)  
Y2H

Large ribosomal subunit  
(strong interaction, clearly diff.)

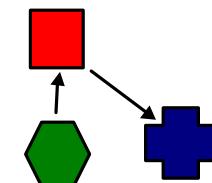


**Distribution of  
Normalized  
Expression Levels**



**Relation of P-P Interactions  
to Abs. Expression Level**

for



# Can we define FUNCTION well enough to relate to expression?

Problems defining function:

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

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

**Non-systematic Terminology:**

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

**Functional Classification**

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

**GenProtEC**  
(*E. coli*, Riley)

**mips**

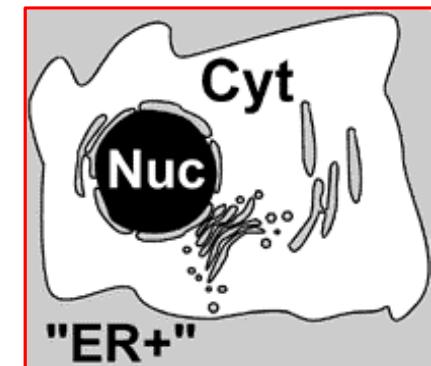
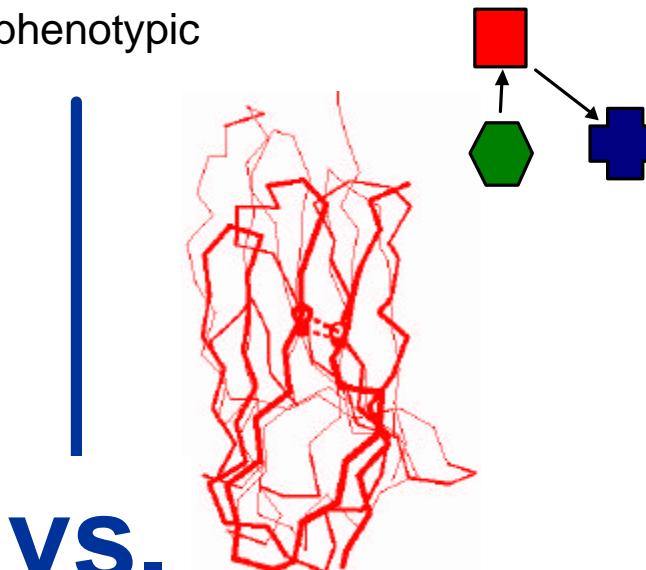
**MIPS/PEDANT**  
(yeast, Mewes)

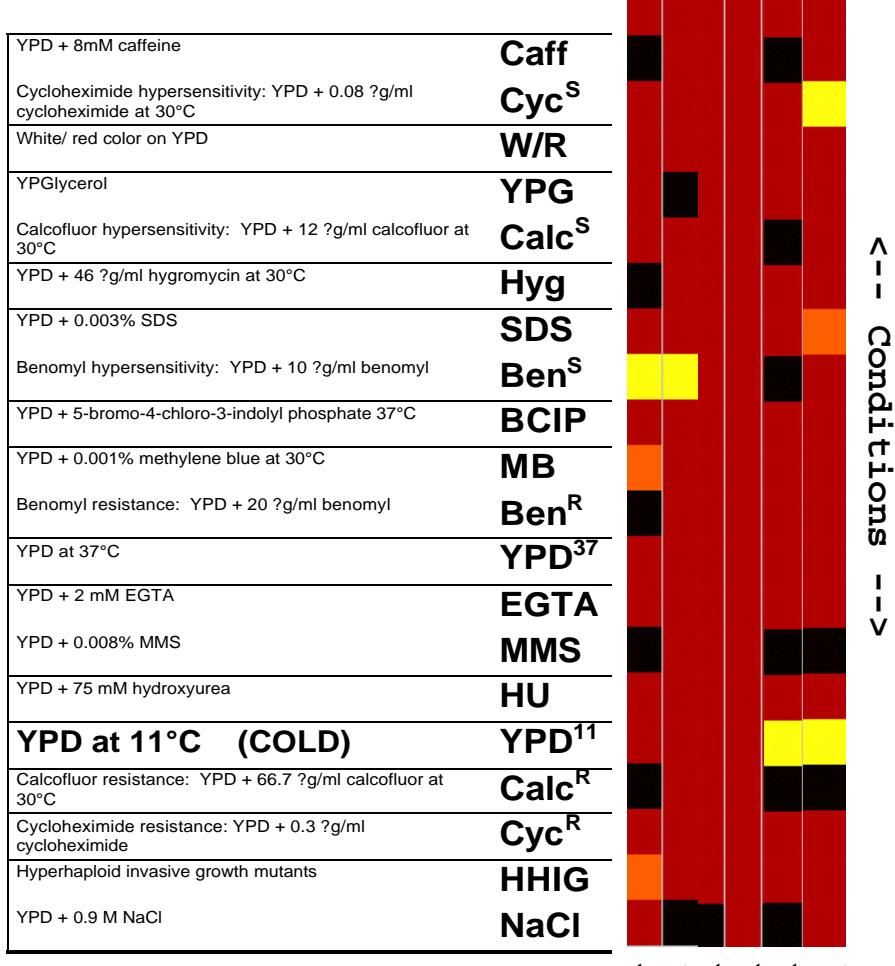
**“Fly”**  
(fly, Ashburner) now extended to **GO** (cross-org.)

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

24 (c) Mark Gerstein, 2000, Yale, bioinfo.mbb.yale.edu

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





M Snyder

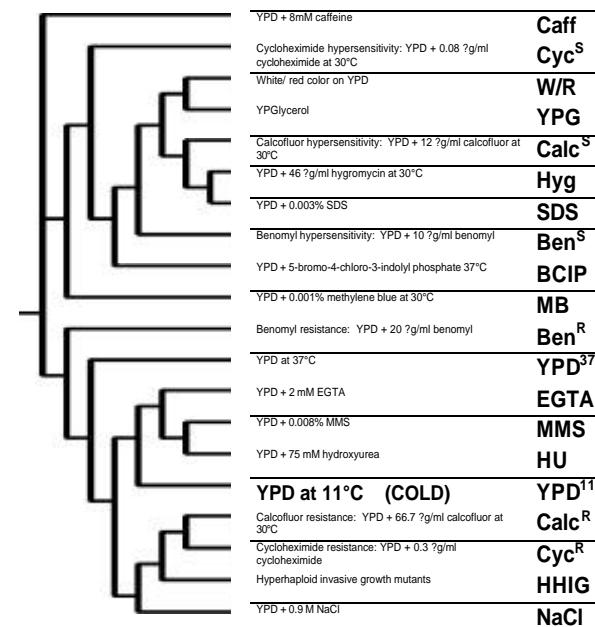
Affected  
by Another  
Condition

WT

Affected  
by Cold

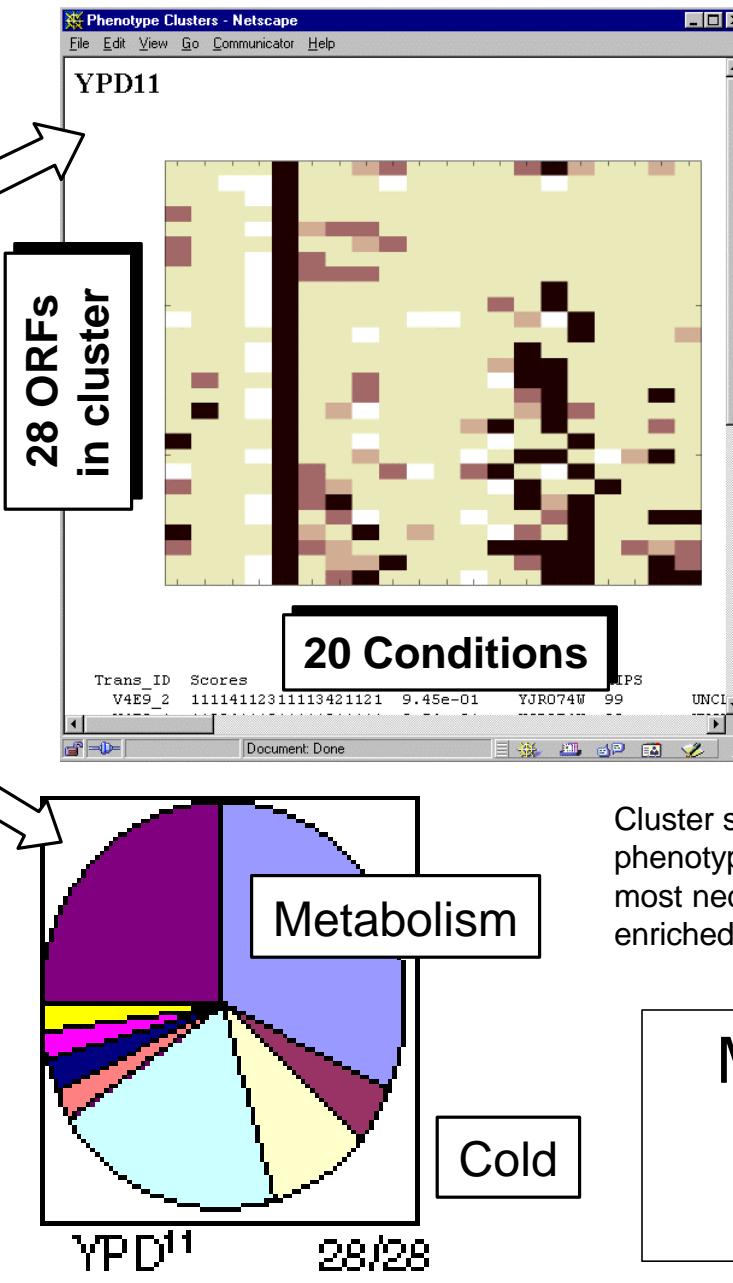
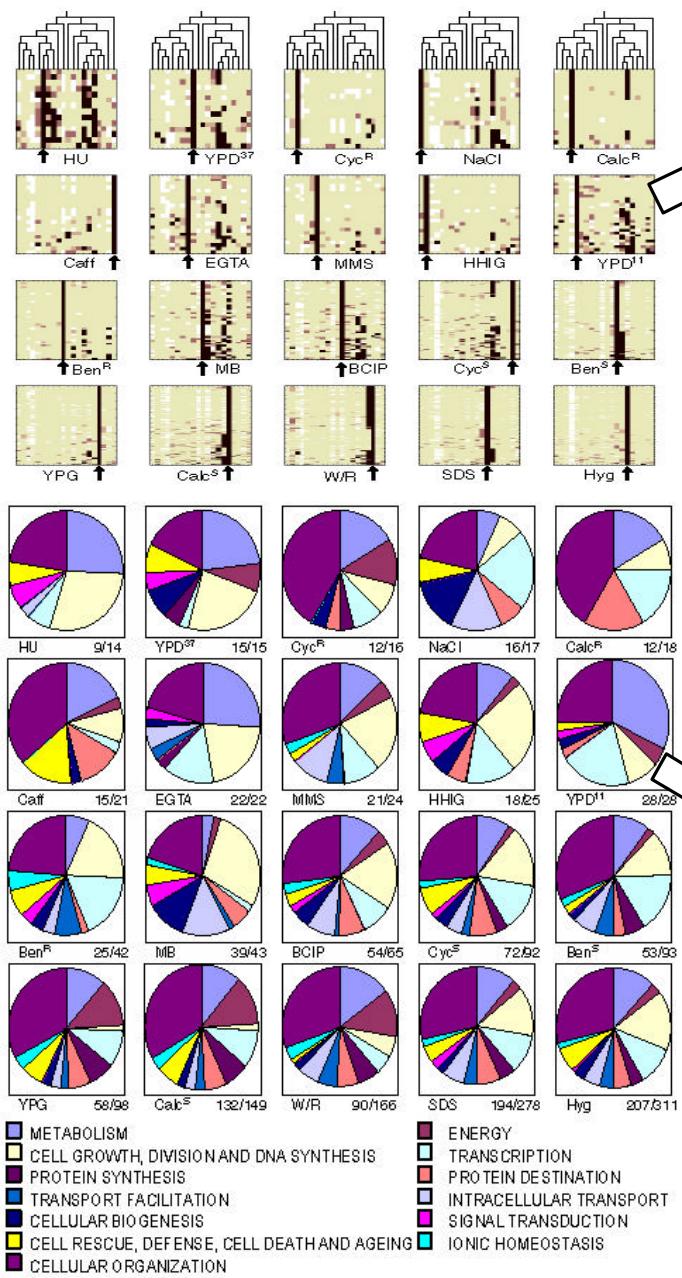
# Whole Genome Phenotype Profiles

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

# Phenotype ORF Clusters from Transposon Expt.



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**

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

M Snyder,  
A Kumar,  
et al....

# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

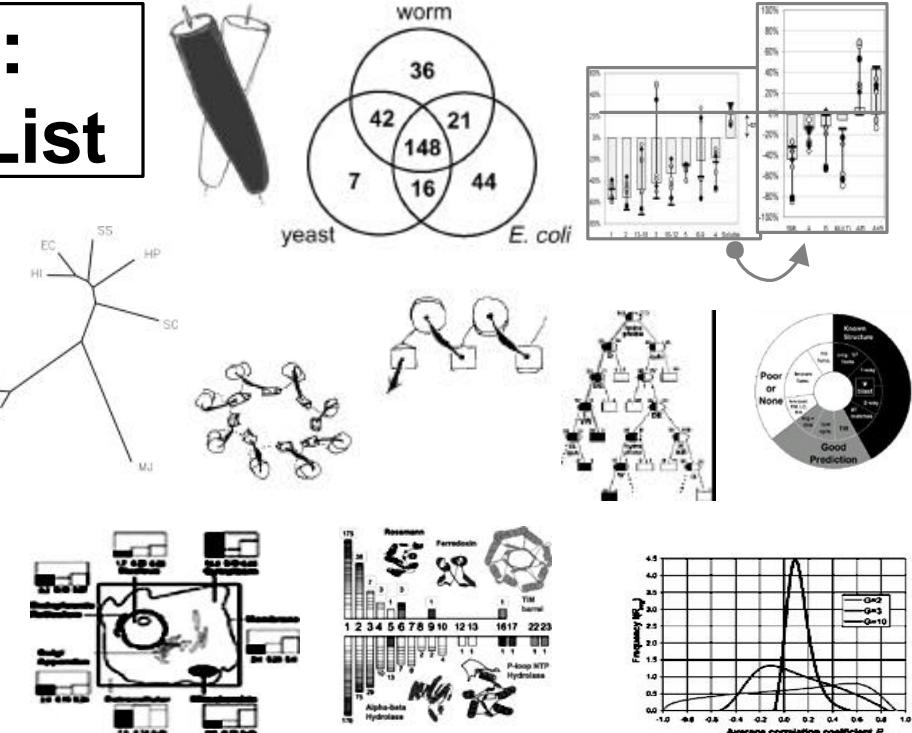
- **Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

- **Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

## 4 Using Folds to Interpret Expression Data

**Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

- **Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.

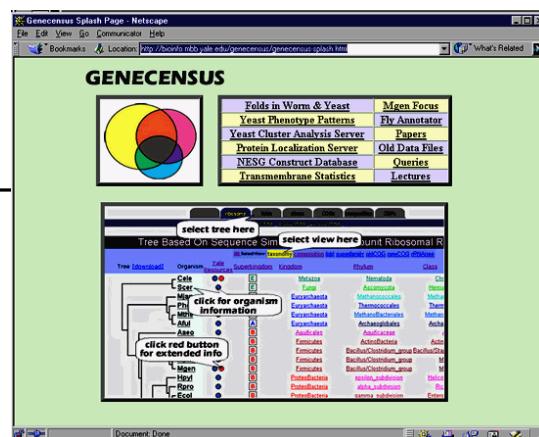
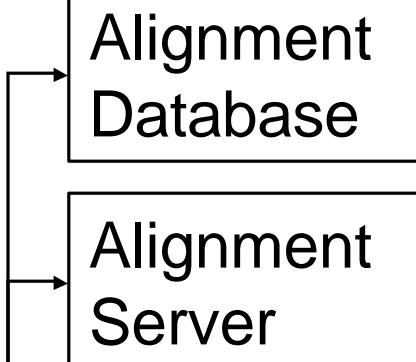


**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

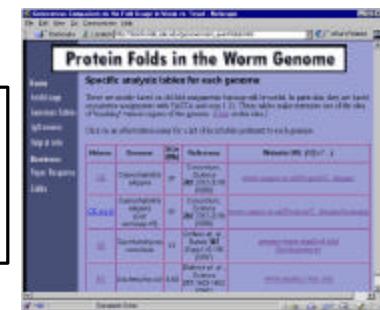
**bioinfo.mbb.yale.edu**

# GeneCensus

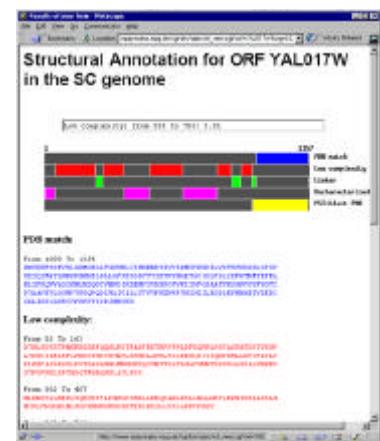
**bioinfo.mbb.yale.edu**



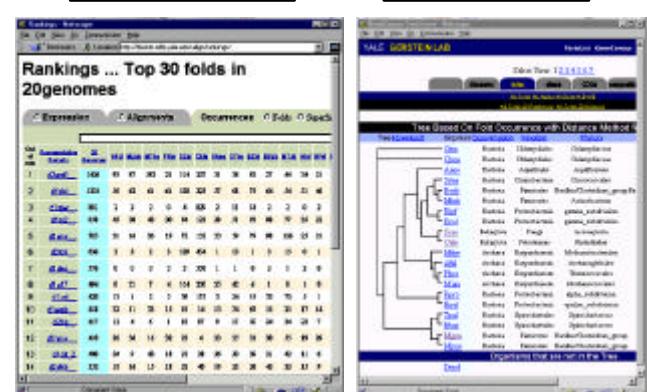
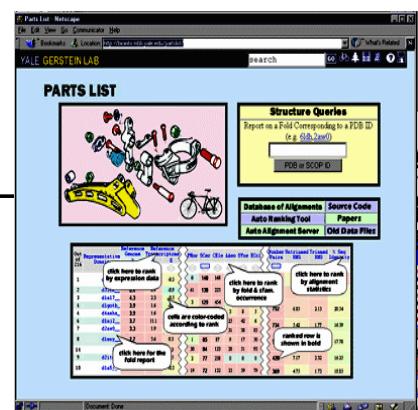
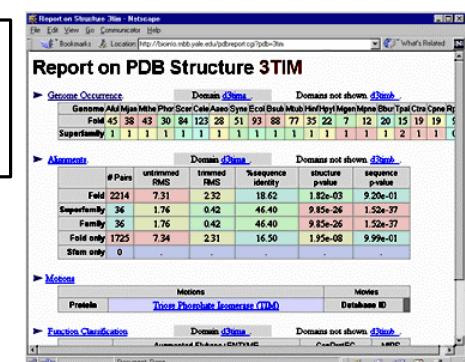
Detailed  
Tables



ORF  
Query



PDB  
Query



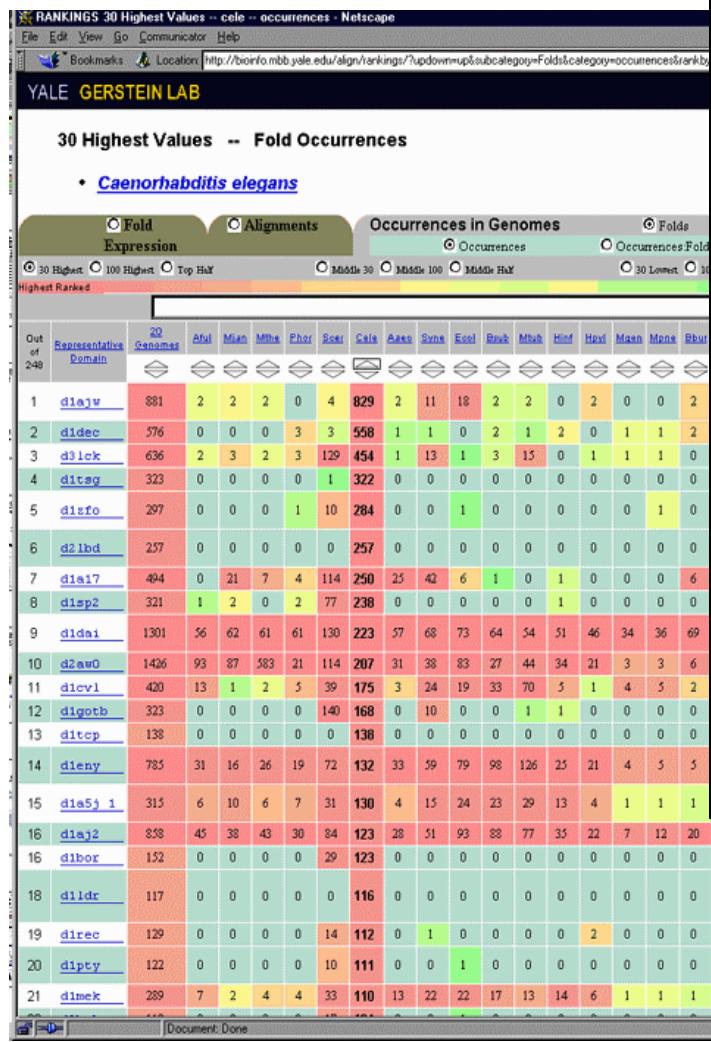
Ranks

Trees

# PartsList

## Ranking

### Viewers



**Rankings2 - Netscape**

**YALE GERSTEIN LAB**

**Rankings<sup>2</sup>**

[View the first 30 folds](#)   [View the entire table](#)

	fold occurrence in Scer	fold occurrence in Cele	fold percentage in reference genome	fold percentage in reference transcriptome
Max of all	140	829	11.3	11.1
Min of all	0	0	0.1	0
Average	6.9	22.9	0.5	0.5
Non zero hits	215	247	213	134
Rank again				
<a href="#">d1ajw</a> :1.002.001 Immunoglobulin-like beta-sandwich Class: All beta proteins	4	829	0.2	0.2
<a href="#">d1dec</a> :1.007.003 Knottins (Small inhibitors) Class: Small Proteins	3	556	-	-
<a href="#">d3lck</a> :1.005.001 Protein kinases (PK) Class: Multi-domain (alpha and beta) proteins	129	454	6.4	0.9
<a href="#">d1tsg</a> :1.004.105 C-type lectin-like Class: Alpha plus beta proteins	1	322	-	-
<a href="#">d1zfo</a> :1.007.033 Glucocorticoid receptor-like (DNA-binding domain) Class: Small Proteins	10	284	0.4	0.1
<a href="#">d21bd</a> :1.001.093 Ligand-binding domain of nuclear receptor Class: All alpha proteins	0	257	-	-
<a href="#">d1al7</a> :1.001.091 alpha-alpha superhelix Class: All alpha proteins	114	250	4.3	2.3
<a href="#">d1sp2</a> :1.007.031 Classic zinc finger Class: Small Proteins	77	238	1.4	0.2

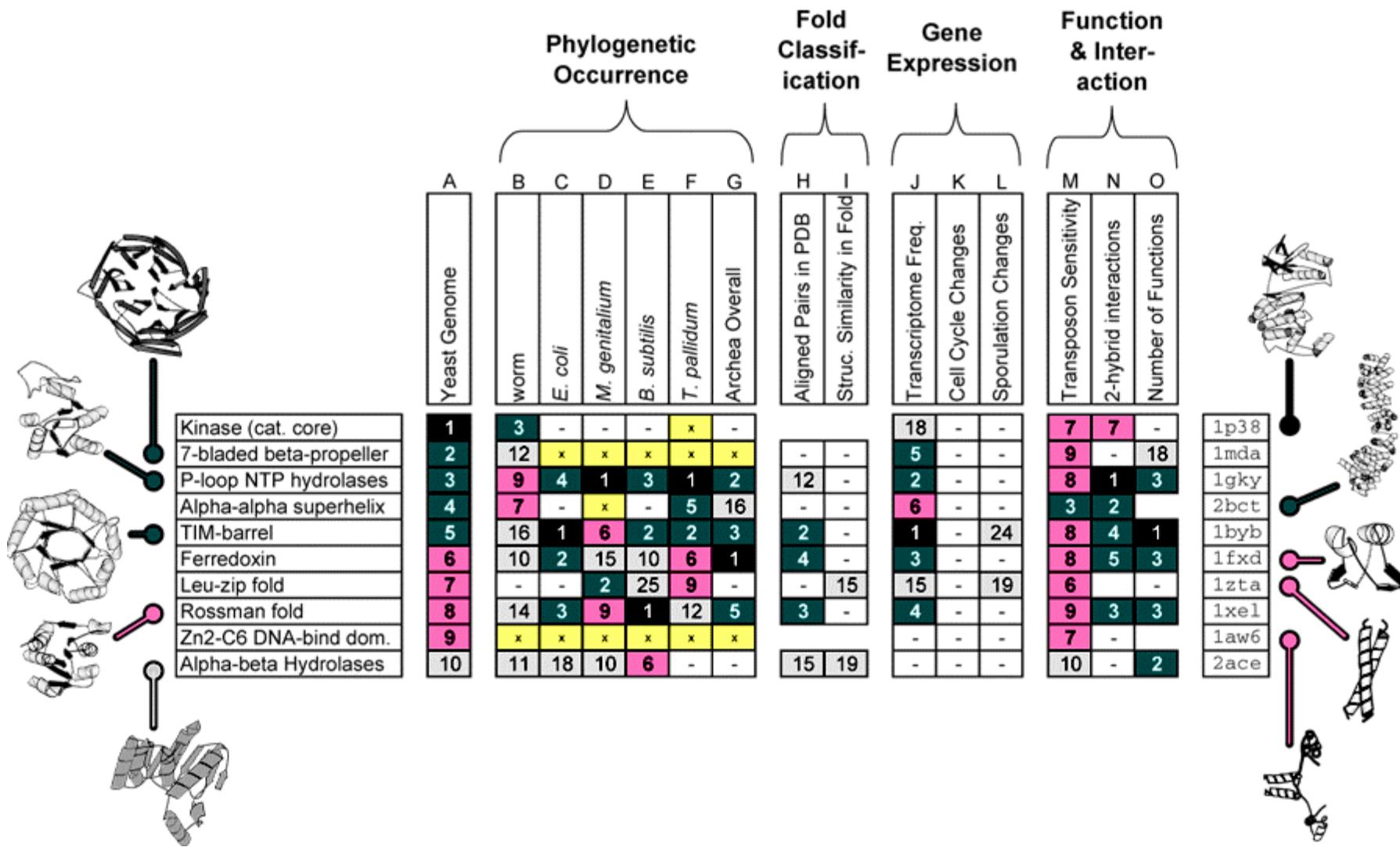
Document: Done

Rank Folds by Genome  
Occurrence, Expression, Fold  
Clustering, Length, &c

J Qian,  
B Stenger,  
J Lin....



# Surveying a Finite PartsList from Many Perspective



# GeneCensus Dynamic Tree Viewers

Recluster organisms based on folds, composition, &c and compare to traditional taxonomy

**Tree Based On Dinucleotide Composition Percentages**

Tree [download]	Organism	Yale Resources	Superkingdom	Kingdom	Phylum	Class
Aaeo		B		Aquificales	Aquificae	Aquifex
Phor		A		Euryarchaeota	Thermococcales	
Bsub		B		Firmicutes	Bacillus/Clostridium_group	
Syne		B		CyanoBacteria	Chroococcales	Synechocystis
Aful		A		Euryarchaeota	Archaeoglobales	Archaeoglobaceae
Ecol		B		ProteoBacteria	gamma_subdivision	Enterobacteriaceae
Tpal		B		Spirochaetales	Spirochaetaceae	Treponema
Mthe		A		Euryarchaeota	Methanobacterales	Methanobacteriaceae
Mtub		B		Firmicutes	Actinobacteria	Actinobacteridae
Mgne		B		Firmicutes	Bacillus/Clostridium_group	Mollicutes
Scer		E		Fungi	Ascomycota	Hemiascomycetes
H pyl		B		ProteoBacteria	epsilon_subdivision	Helicobacter_group
Hinf		B		ProteoBacteria	gamma_subdivision	Pasteurellaceae
Cele		E		Metazoa	Nematoda	Chromadorea
Mgen		B		Firmicutes	Bacillus/Clostridium_group	Hemiscomycetes
Mian		A		Euryarchaeota	Methanococcales	Methanococcaceae
Bbur		B		Spirochaetales	Spirochaetaceae	Thermococcaceae
Rpro		B		ProteoBacteria	alpha_subdivision	Methanobacteriales
Cpne		B		Chlamydiales	Chlamydiae	Methanobacteriaceae
Ctra		B		Chlamydiales	Chlamydiae	Archaeoglobaceae

**Organisms that are not in the Tree**

Pfa2		E	Alveolata	Apicomplexa	Haemosporida
Hpy2		B	ProteoBacteria	epsilon_subdivision	Helicobacter_group
Lei1		E	Euglenozoa	Kinetoplastida	Trypanosomatidae
Aper		A	Crenarchaeota	Desulfurococcales	Desulfurococcaceae
Tmar		B	Thermogales	Thermotoga	N/A
Pfa3		E	Alveolata	Apicomplexa	Haemosporida

**Tree Based On Sequence Similarity of Small Subunit Ribosomal RNA**

Tree [download]	Organism	Yale Resources	Superkingdom	Kingdom	Phylum	Class
Cele		E		Metazoa	Nematoda	Chromadorea
Scer		E		Fungi	Ascomycota	Hemiscomycetes
Mjan		A		Euryarchaeota	Methanococcales	Methanococcaceae
Phor		A		Euryarchaeota	Thermococcales	Thermococcaceae
Mthe		A		Euryarchaeota	Methanobacterales	Methanobacteriales
Aful		A		Euryarchaeota	Archaeoglobales	Archaeoglobaceae
Aaeo		B		Aquificales	Aquificae	Aquifex
Mtub		B		Firmicutes	Actinobacteria	Actinobacteridae
Bsub		B		Firmicutes	Bacillus/Clostridium_group	Mollicutes
Mgne		B		Firmicutes	Bacillus/Clostridium_group	Mollicutes
H pyl		B		ProteoBacteria	epsilon_subdivision	Helicobacter_group
Rpro		B		ProteoBacteria	alpha_subdivision	Rickettsiales
Ecol		B		ProteoBacteria	gamma_subdivision	Enterobacteriaceae
Hinf		B		ProteoBacteria	gamma_subdivision	Pasteurellaceae
Bbur		B		Spirochaetales	Spirochaetaceae	Borrelia
Tpal		B		Spirochaetales	Spirochaetaceae	Treponema
Syne		B		CyanoBacteria	Chroococcales	Synechocystis
Cpne		B		Chlamydiales	Chlamydiae	Chlamydophila
Ctra		B		Chlamydiales	Chlamydiae	Chlamydia

**Organisms that are not in the Tree**

Pfa2		E	Alveolata	Apicomplexa	Haemosporida
Hpy2		B	ProteoBacteria	epsilon_subdivision	Helicobacter_group
Lei1		E	Euglenozoa	Kinetoplastida	Trypanosomatidae
Aper		A	Crenarchaeota	Desulfurococcales	Desulfurococcaceae
Tmar		B	Thermogales	Thermotoga	N/A
Pfa3		E	Alveolata	Apicomplexa	Haemosporida
Drad		B	Thermus/Deinococcus_group	Deinococcales	Deinococcus

# Comparative Genomics: Surveys of a Finite Parts List

## 1 Using Folds to Interpret Genomes

**Genomes.** Fold Library background. Shared and/or unique parts. Venn Diagram, Fold tree with all- $\beta$  diff. Ortholog tree. Horizontal Transfer. Common Parts: Top-10 folds with  $\beta\alpha\beta$ . Common  $\Psi$ -fold.

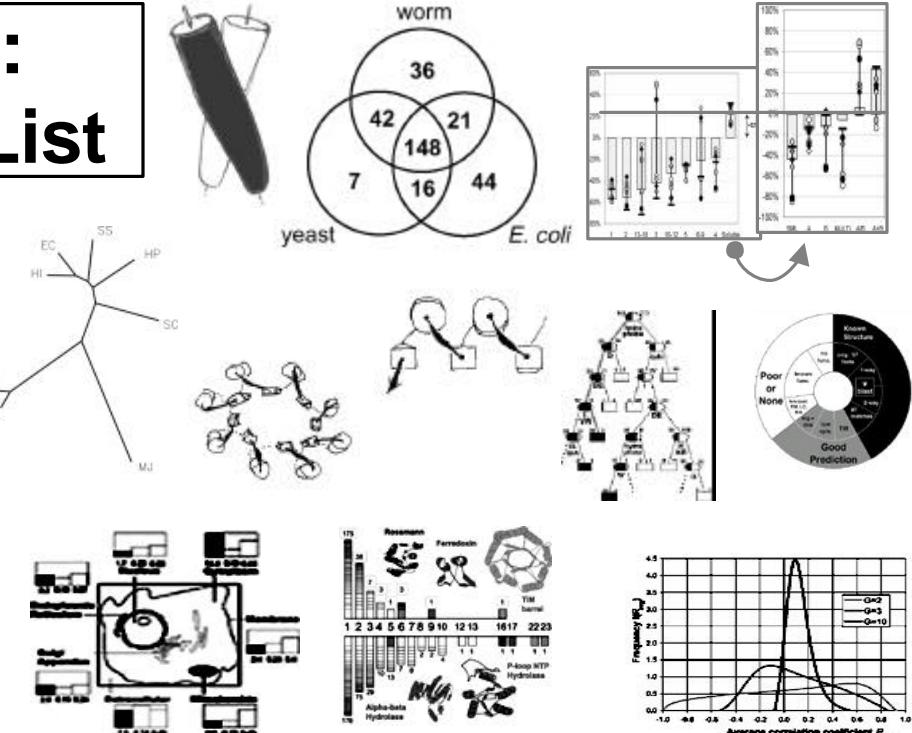
- **Tricky Issues → Expt. SG.** Extent of fold assignment (MG, 65%, 98). Predictions. Biases. NESG.org Construct DB. Datamining this w/ decision trees. Selecting weird MG CD targets.

- **Folds & Functions.** Roles/part? How many folds/func? 331 of ~20K combinations. Mostly 1 func/fold, but some versatile scaffolds -- TIM most versatile. Similar for interactions. Func. Divergence vs. Seq. & Struc. Diverg.

## 4 Using Folds to Interpret Expression Data

**Expression Data.** Top-10 parts in other terms. Enriched in transcriptome: VGA,  $\alpha\beta$  folds, energy, synthesis, Cytoplasmic, TIM fold. Depleted: NS, long, TMs, transport, transcription, Nuclear, Leu-zip fold. Bayesian localizer.

- **Tricky Issues: Relating Expression to Function.** Expression relates to structure & localization but to function, globally? Weak relation to protein-protein interactions.



**H Hegyi, J Lin, B Stenger,  
N Echols, P Bertone, J Qian,  
L Regan, S Balasubramanian,  
V Alexandrov, G Montelione,  
A Edwards, C Wilson, Y Kluger,  
C Arrowsmith, A Drawid,  
R Jansen, D Greenbaum,  
S Teichmann, P Harrison**

**bioinfo.mbb.yale.edu**