

Integrative Genomics:

Surveys of a Finite Parts List

Mark Gerstein

*H Hegyi, J Lin, N Echols,
P Harrison, A Drawid, M Levitt, R Jansen, D
Greenbaum, M Snyder, S Teichmann, P Bertone,
B Stenger, V Alexandrov, J Tsai,
C Wilson, J Qian, W Krebs*

Talk at Stanford U, 16 April 2001

1995

Bacteria,
1.6 Mb,
~1600 genes
[*Science* 269: 466]



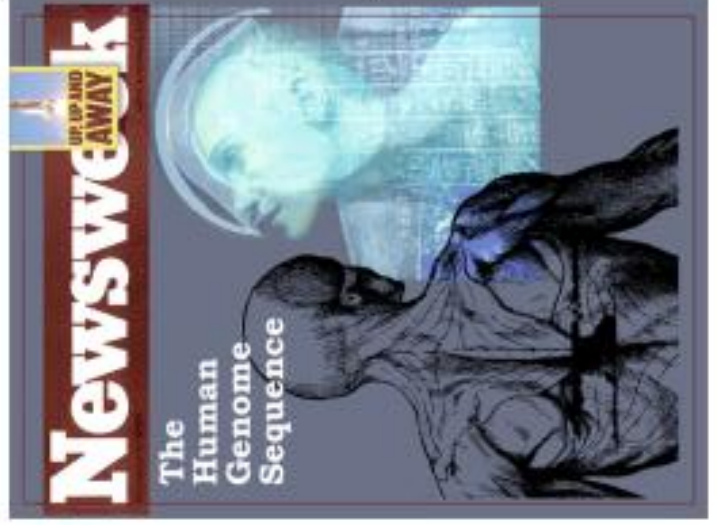
1997

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



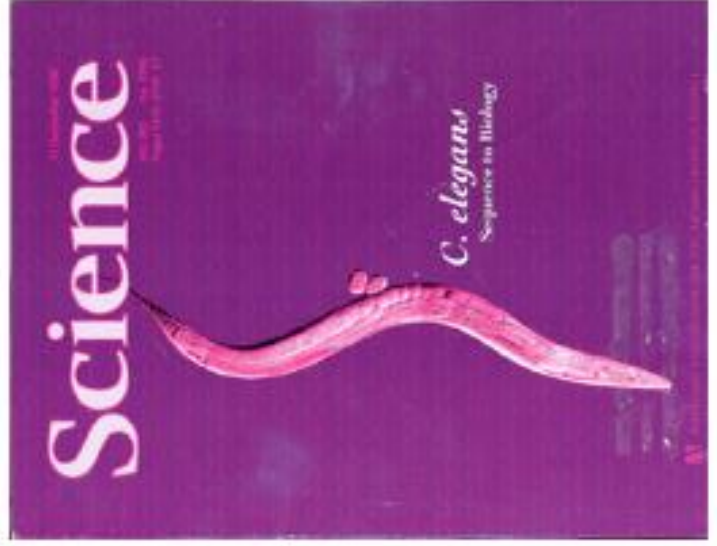
1998

Animal,
~100 Mb,
~20K genes
[*Science* 282:
1946]



2000

Human,
~3 Gb,
~30K genes



Genomes
highlight
the
Finiteness
of the
"Parts" in
Biology

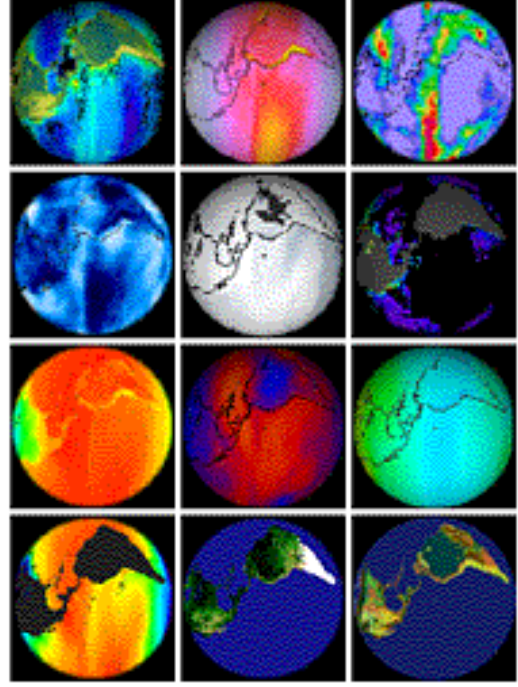
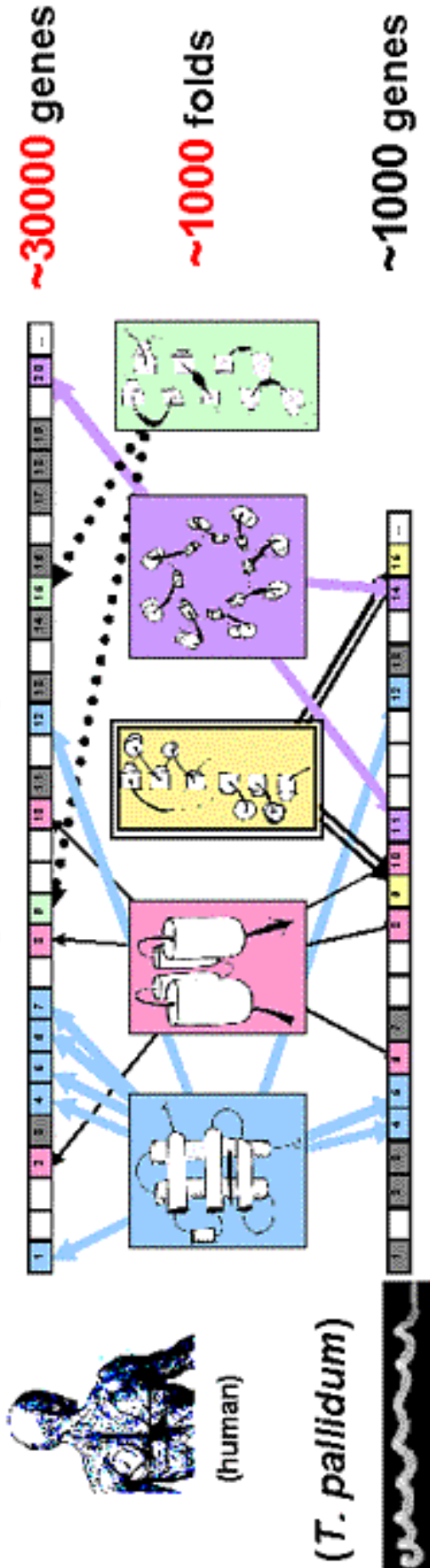


'98 spoof

real thing, Apr '00

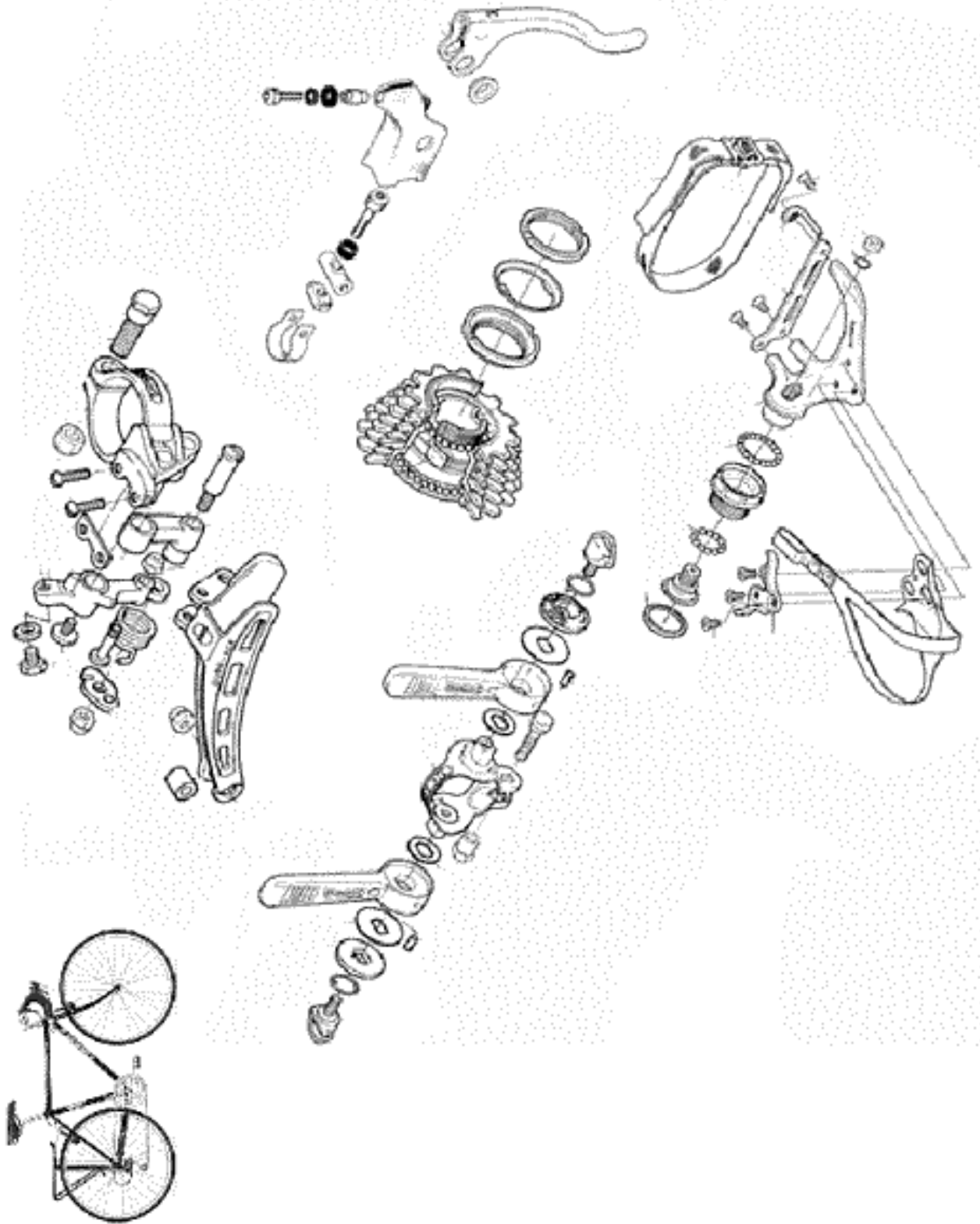
Global Surveys of a Finite Set of Parts from

Many Perspectives

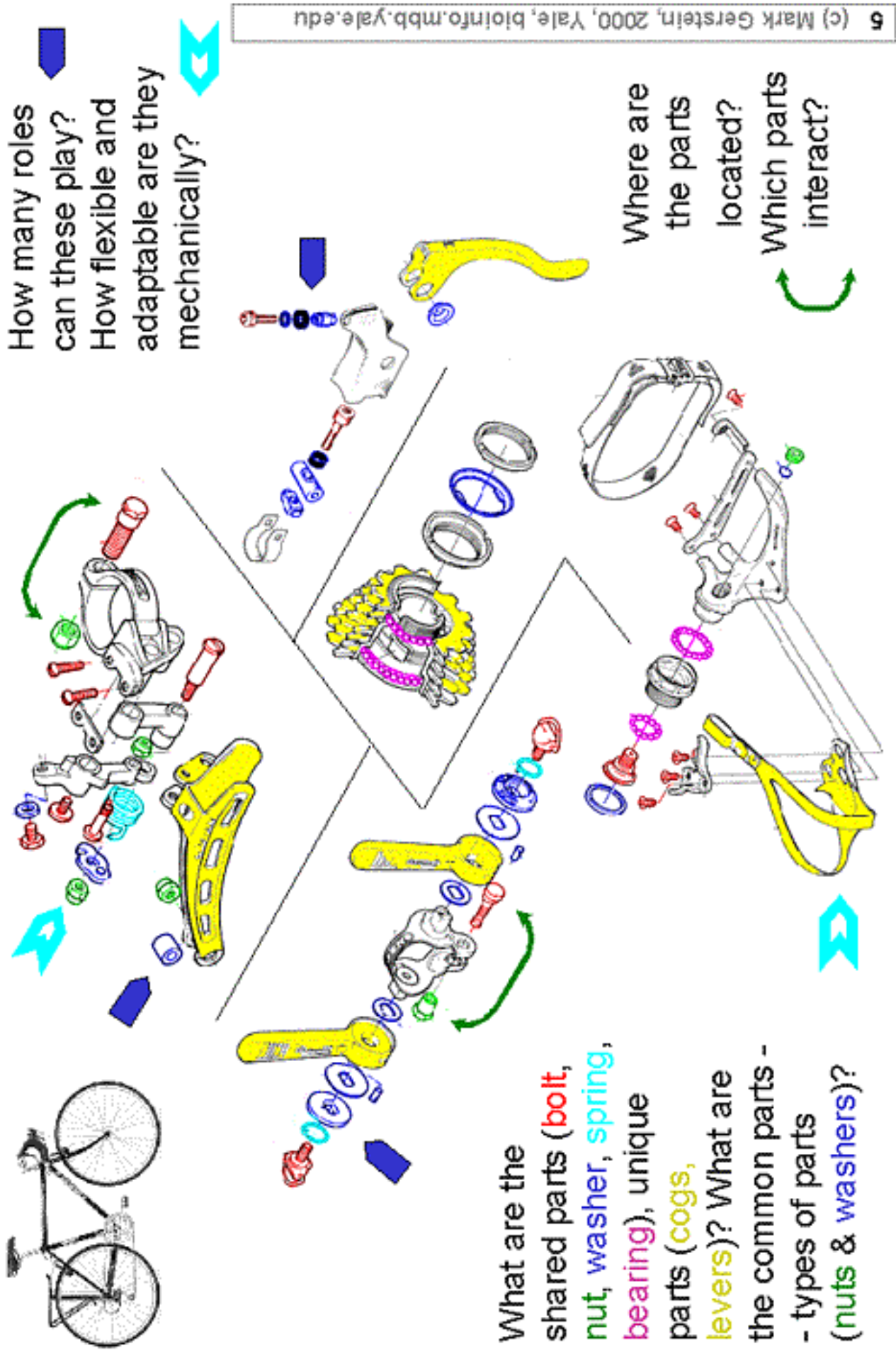


World of structures is even more finite than that of sequences, providing simplification...

A Parts List Approach to Bike Maintenance



A Parts List Approach to Bike Maintenance



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

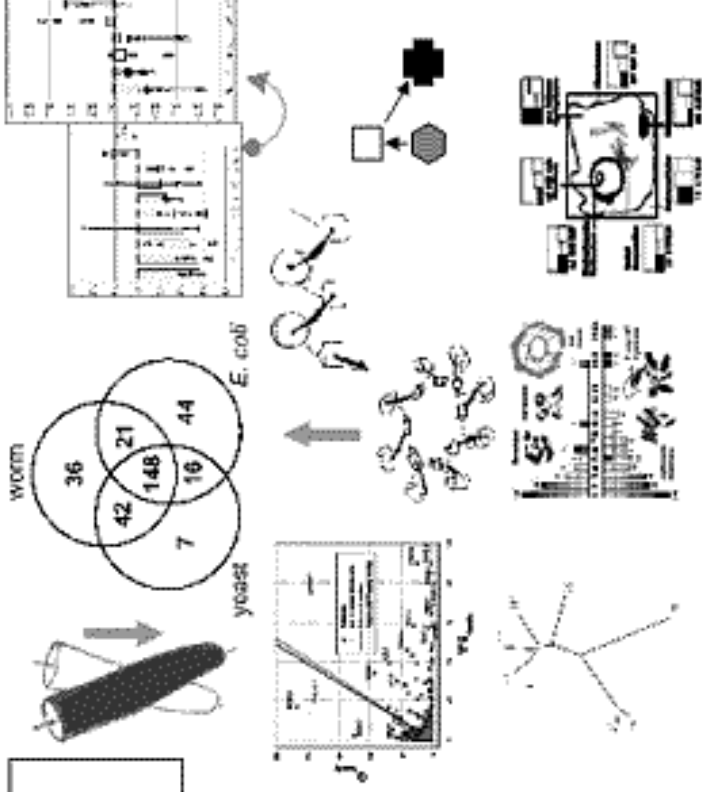
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

What is a Proteomic “Part” or Category?

Result of **grouping together**
sequences & structures

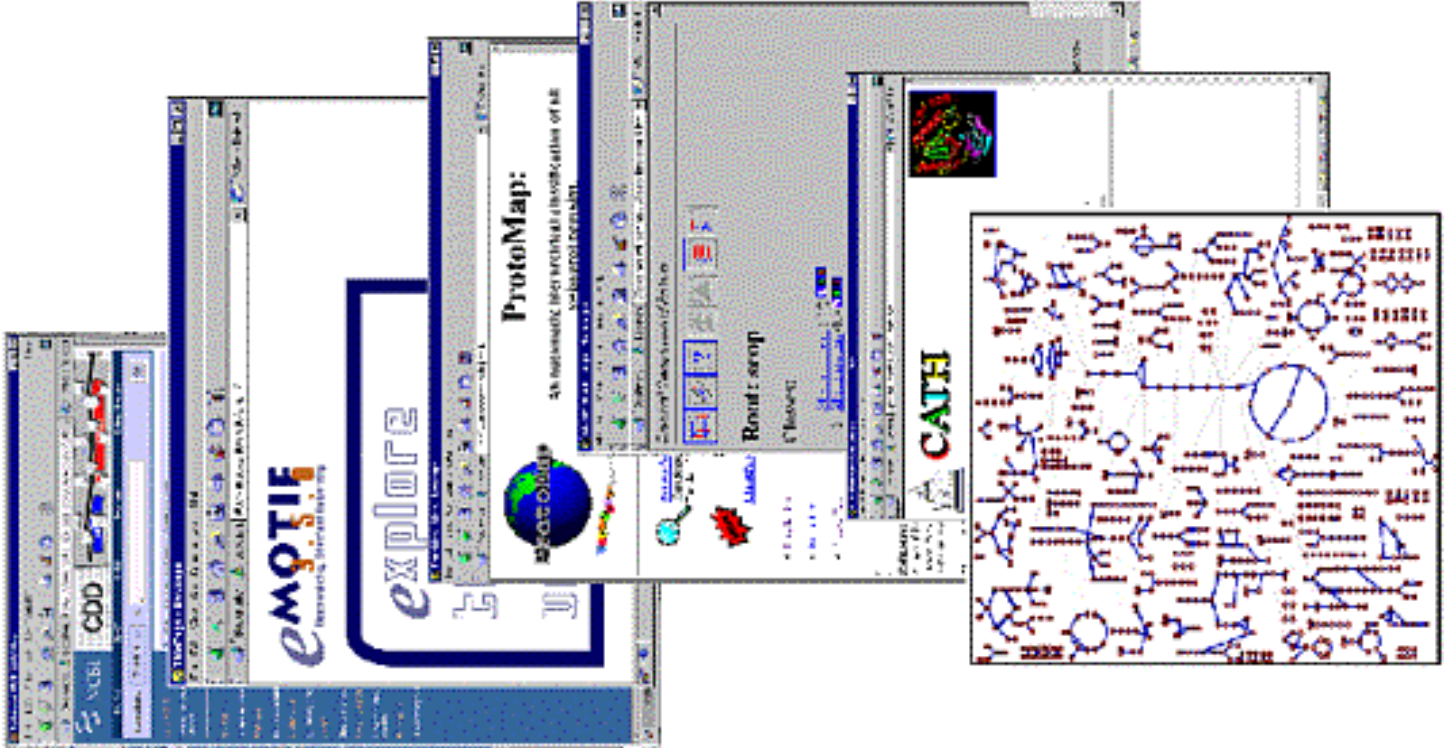
Many **non-unique but self-**
consistent definitions

Folds: most intuitive & fundamental
parts (focus here)

Same logic for seq. families
(homologs), orthologs, blocks,
motifs.... even pathways,
functional systems

Much previous work...

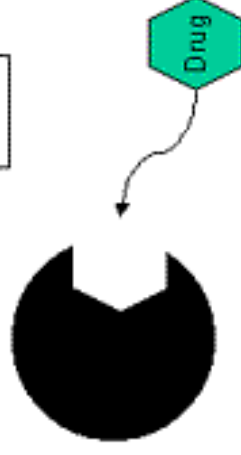
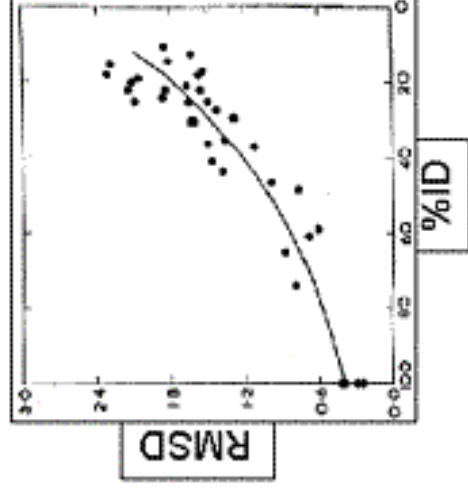
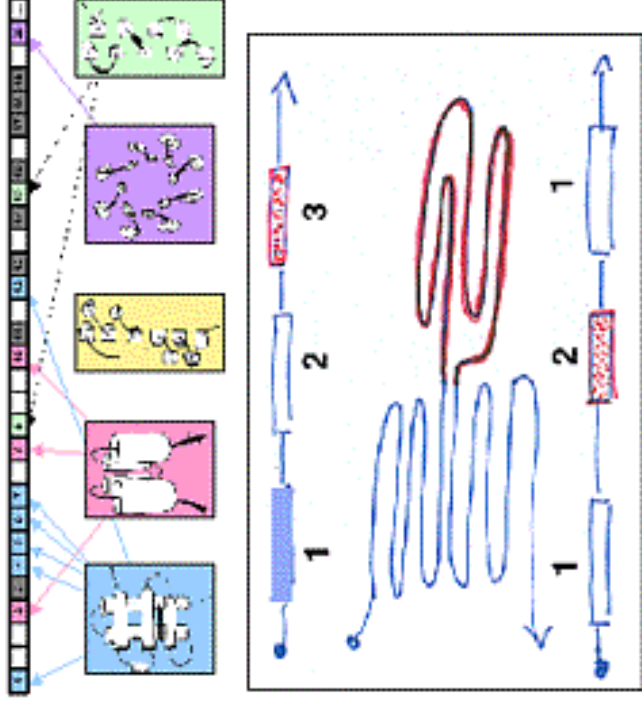
(CATH, scop, Blocks, FSSP, Interpro, eMotif, Prosite, CDD,
Pfam, ProtoMap, Prints, VAST, TOGA... Remington, Matthews '80;
Taylor, Orengo '89, '94; Thornton, CATH; Artymnik, Rice, Willett '89; Sali, Blundell, '90;
Vriend, Sander '91; Russell, Barton '92; Holm, Sander '93+ (FSSP); Godzik, Skolnick '94;
Gibrat, Bryant '96 (VAST); F Cohen, '96; Feng, Sippl '96; G Cohen '97; Singh & Brutlag,
'98) (Functions picture from www.fruitfly.org/~suzi (Ashburner); Pathways picture from,
ecocyc.pangeasystems.com/ecocyc (Karp, Riley).)

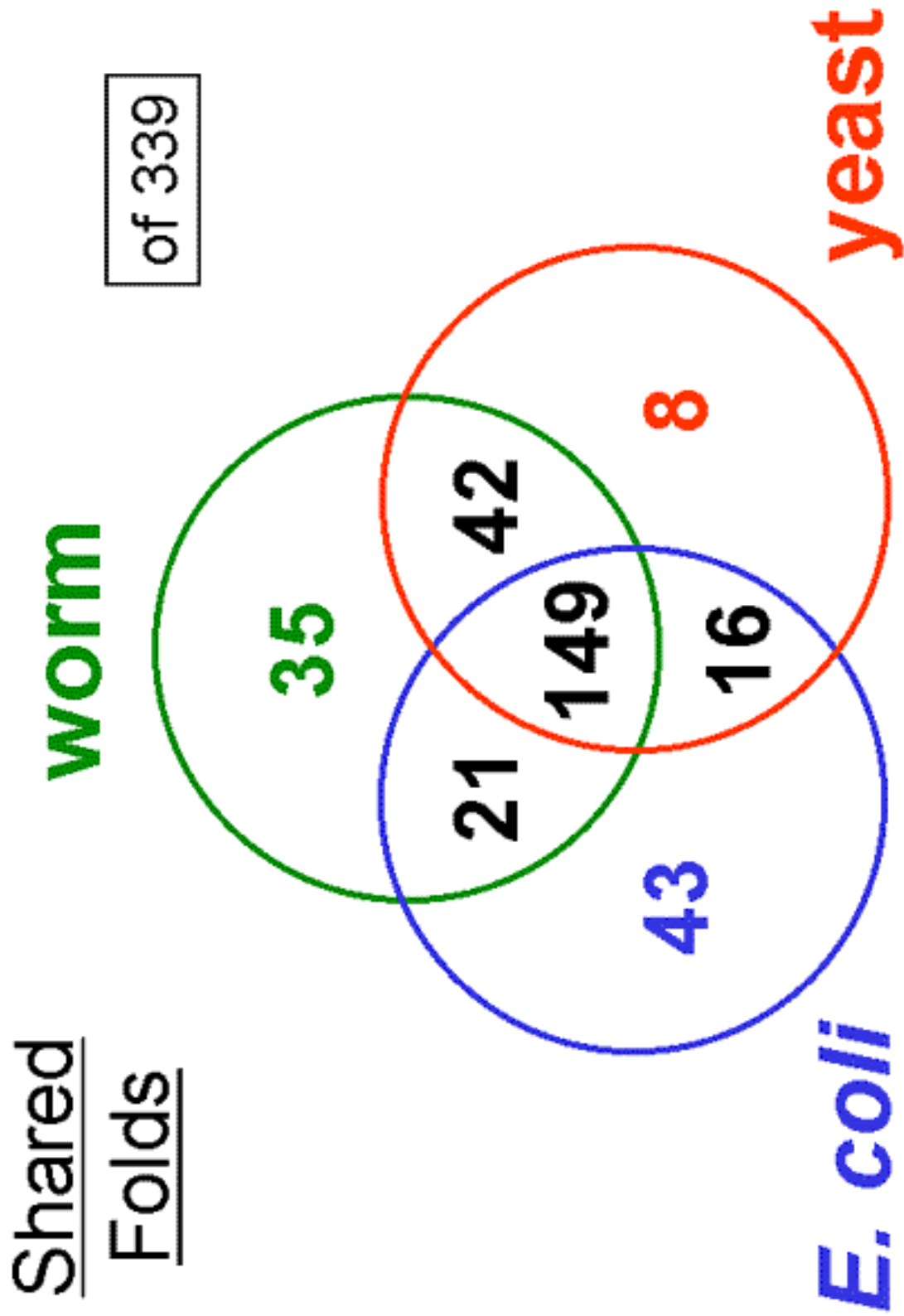


Folds as the best parts, the final annotation for the human genome. Why?

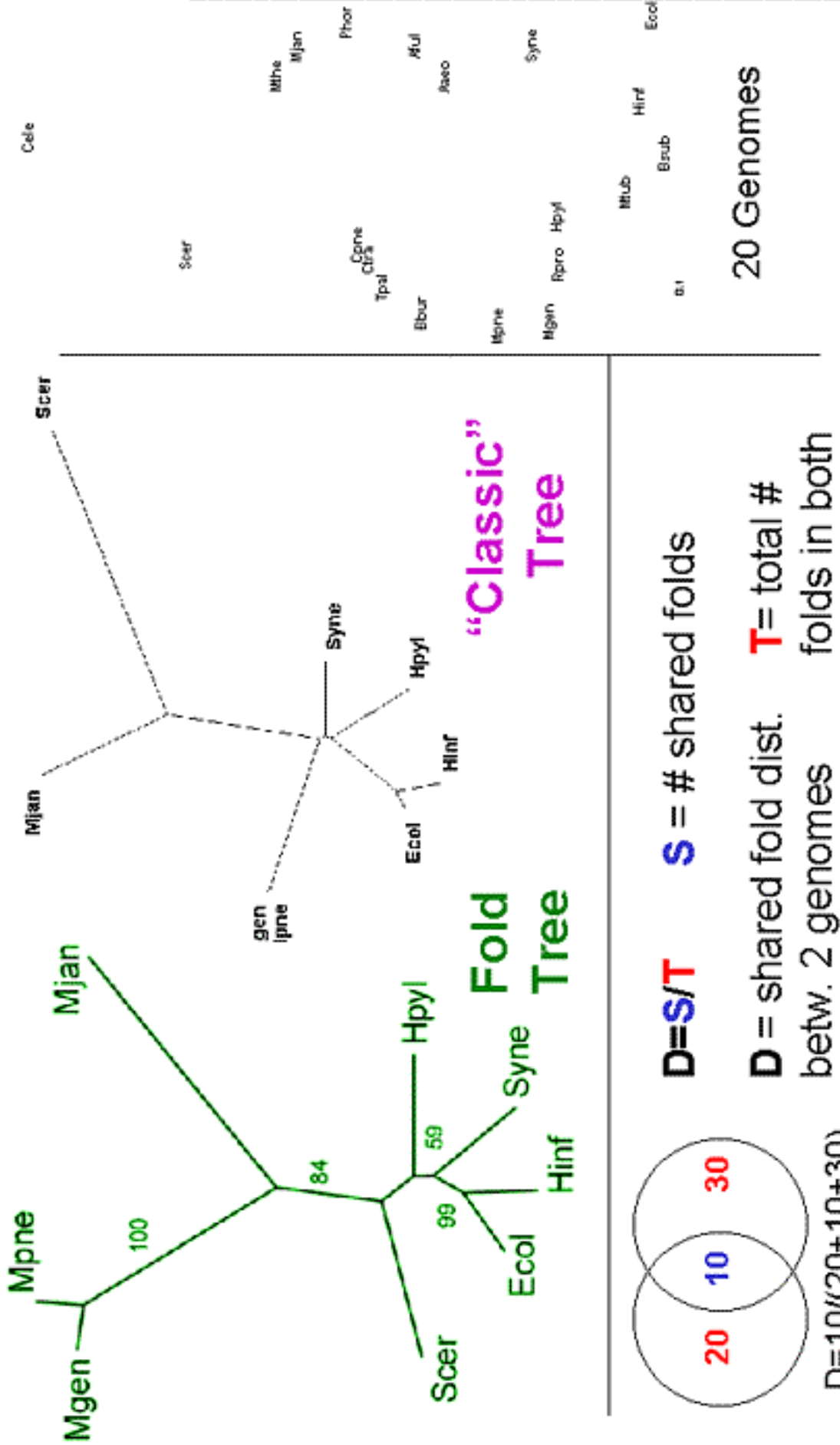
- 1 Most Highly Conserved
- 2 Precisely Defined Modules
- 3 Seq. \Leftrightarrow Struc. Clearer than Seq. \Leftrightarrow Func.
- 4 Link to Chemistry, biochemical function, drugs

Here mostly
scop+astral+auto-alignments
(Chothia, Levitt, Brenner, Murzin)



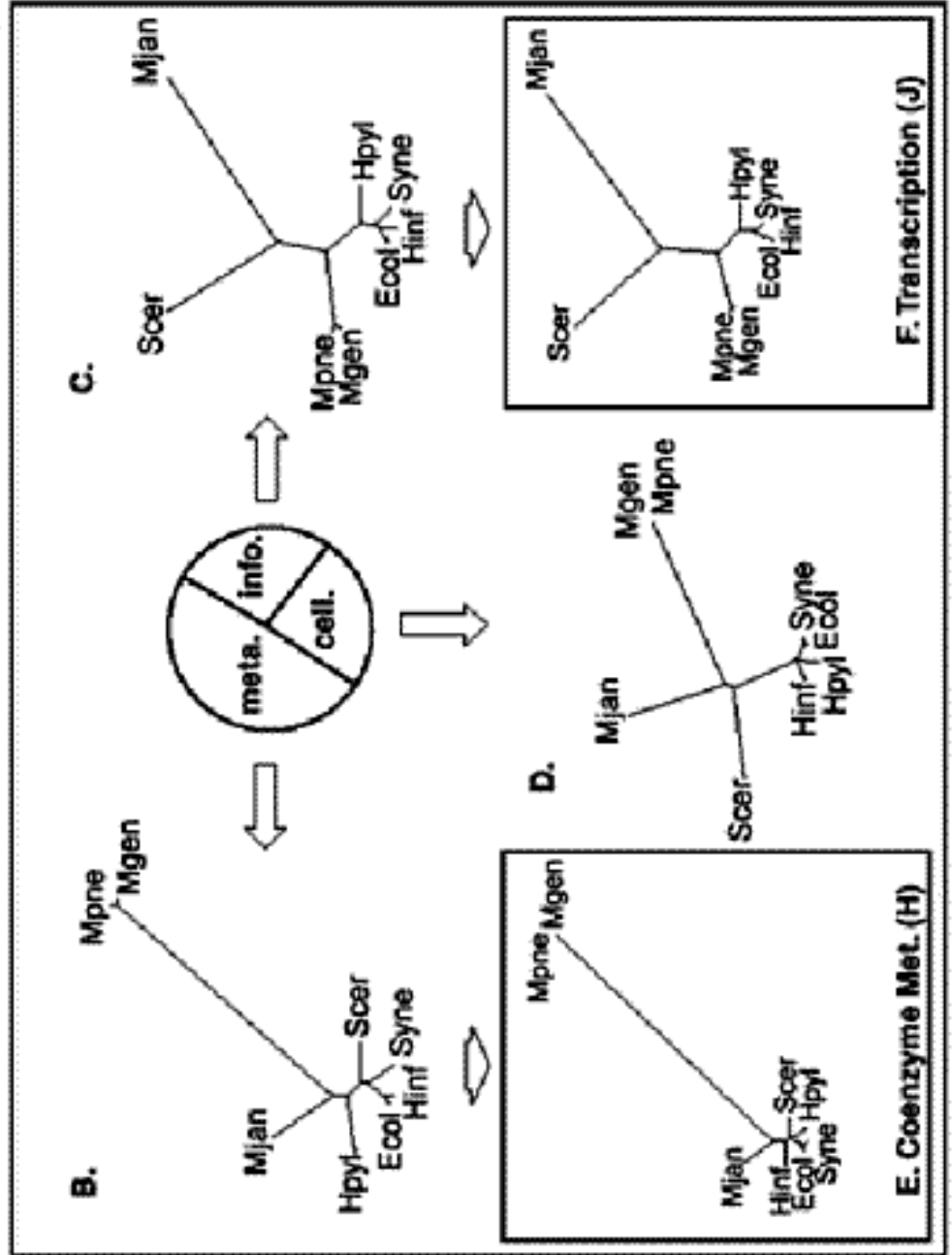
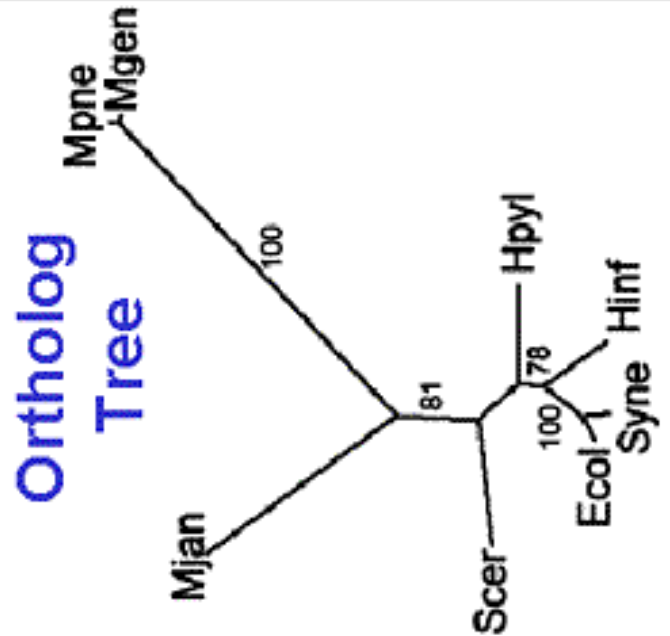
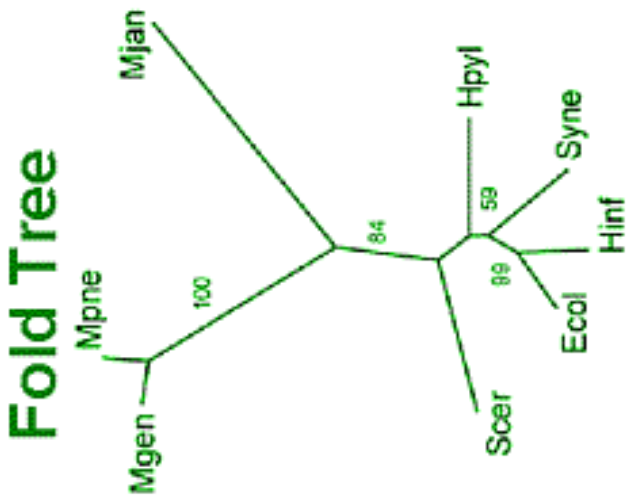


Cluster Trees on Basis of Shared Initial Genomes on Basis of Shared Folds



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

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



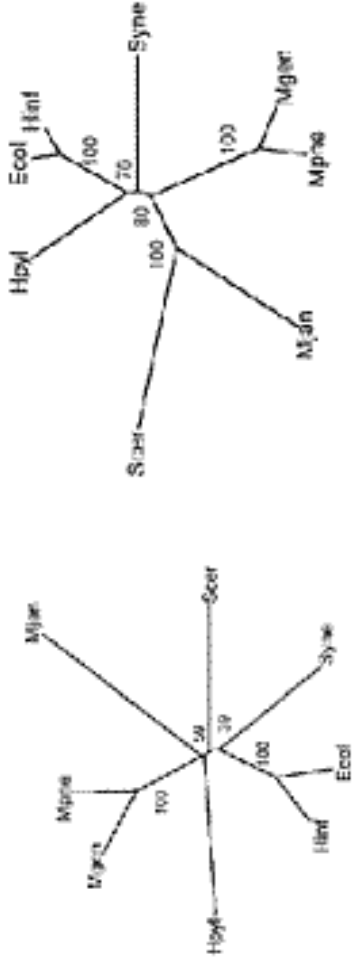
Compare with trees

on spectrum of

“levels”: single-gene

trees, whole-genome

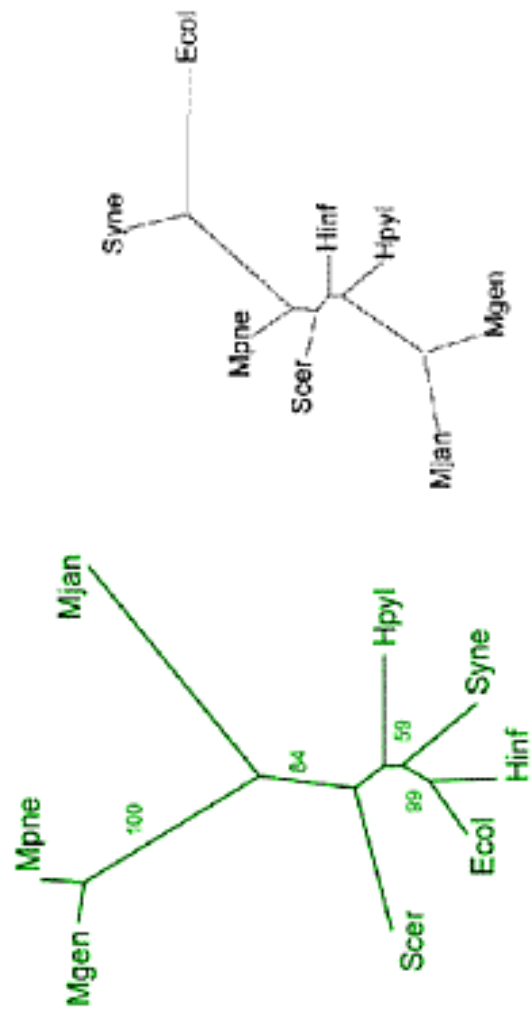
composition trees



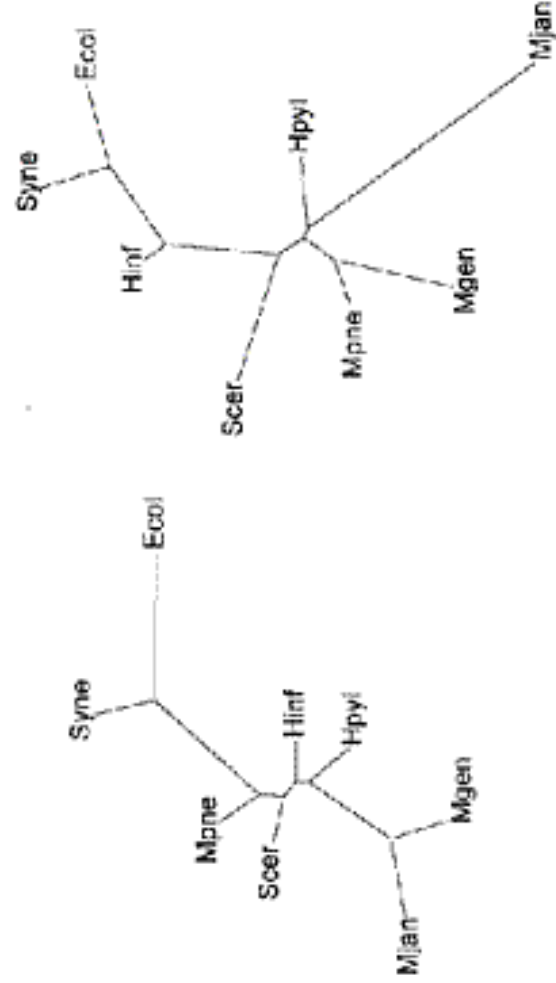
ribosomal protein

TIM

Single-gene Trees



Fold Tree



Ortholog Tree

“Classic” Tree

A. Di-Nucleotide

B. Amino Acid

AA & di-NT Composition Trees

(S Karlin)

Common Folds

⇐ **Worm**

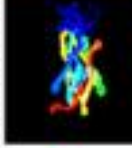

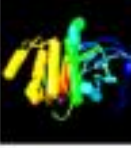


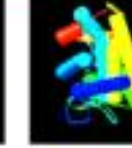

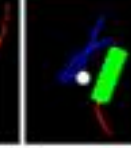
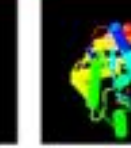
Extracellular, signaling and trans. factor folds (lg, kinase, lectin, nuc. receptor, ZnF)







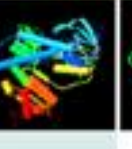



E. coli ⇨

Metabolic Folds (TIM, Ferredoxin, Rossmann, P-loop hydrolase, FAD-binding)

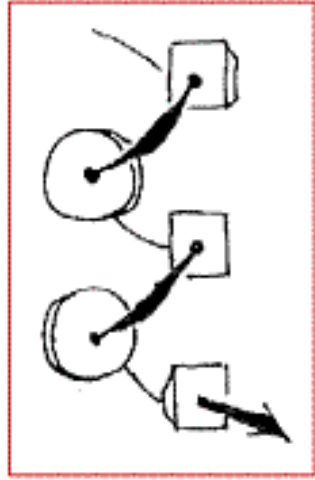
(partslist.org)

Ranked folds (click on arrows to re-rank)

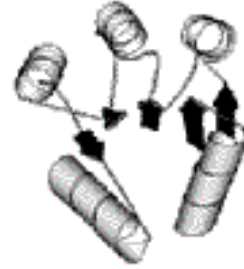
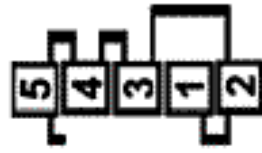
829		Immunoglobulin-like beta-sandwich d1ajw :2.1 (1.39) B
556		Knottins (Small inhibitors) d1dec :7.3 (1.39) S
454		Protein kinases (PK) d3ick :5.1 (1.39) M
322		C-type lectin-like d1tsg :4.105 (1.39) A+B
284		Glucocorticoid receptor-like (DNA-binding domain) d1zfo :7.33 (1.39) S
257		Ligand-binding domain of nuclear receptor d2lbd :1.95 (1.39) A
250		alpha-alpha superhelix d1a17 :1.91 (1.39) A
238		Classic zinc finger d1sp2 :7.31 (1.39) S
223		P-loop containing nucleotide triphosphate

93		beta/alpha (TIM)-barrel d1a12 :3.1 (1.39) A/B
82		Ferredoxin-like d2aw0 :4.34 (1.39) A+B
79		NAD(P)-binding Rossmann-fold domains d1eny :3.22 (1.39) A/B
72		P-loop containing nucleotide triphosphate hydrolases d1d6i :3.29 (1.39) A/B
67		Flavodoxin-like d1wab :3.14 (1.39) A/B
50		Ribonuclease H-like motif d2ing :3.47 (1.39) A/B
38		FAD/NAD(P)-binding domain d1oit :3.4 (1.39) A/B
36		Periplasmic binding protein-like II d1a8e :3.83 (1.39) A/B
36		PLP-dependent transferases d2dkb :3.54 (1.39) A/B
		S-adenosyl-L-methionine-dependent

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



P-loop
hydrolase



Flavodoxin
like

42

ARTICLES

NATURE VOL 393 NOVEMBER 1998

A peptide model of a protein folding intermediate

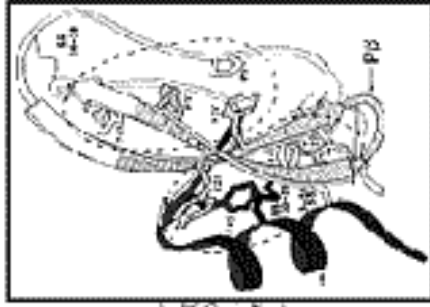
Terrence G. Oas & Peter S. Kim

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

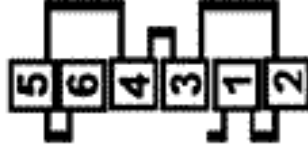
It is difficult to determine the structures of protein folding intermediates because folding is a highly disordered process. To overcome this, we designed a model of a protein folding intermediate, consisting of a beta-strand and an alpha-helix, and used NMR to determine its structure. We found that the structure of the folding intermediate is similar to that of the native protein.

336: 42

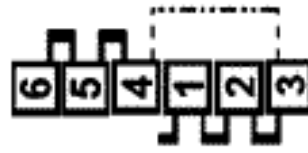
HI, MJ, SC vs scop 1.32. All share α/β structure with repeated R.H. β regions connecting adjacent strands or nearly so (18+4+2 of 24)



TIM-
barrel



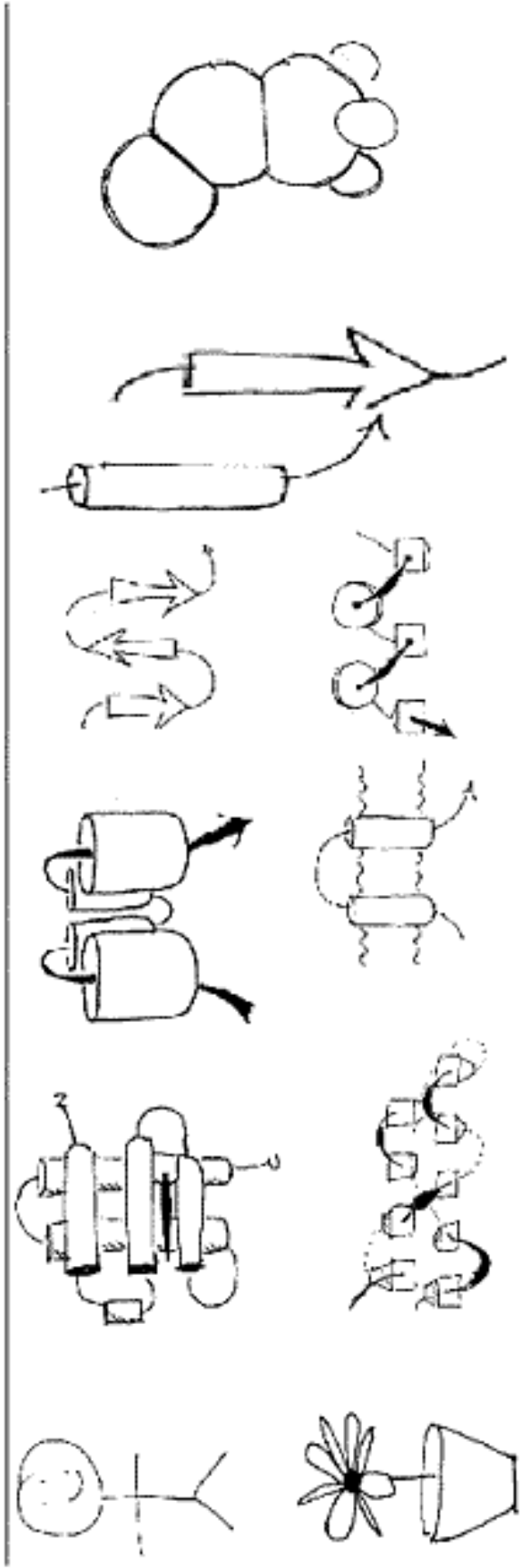
Thiamin
Binding



Rossmann
Fold

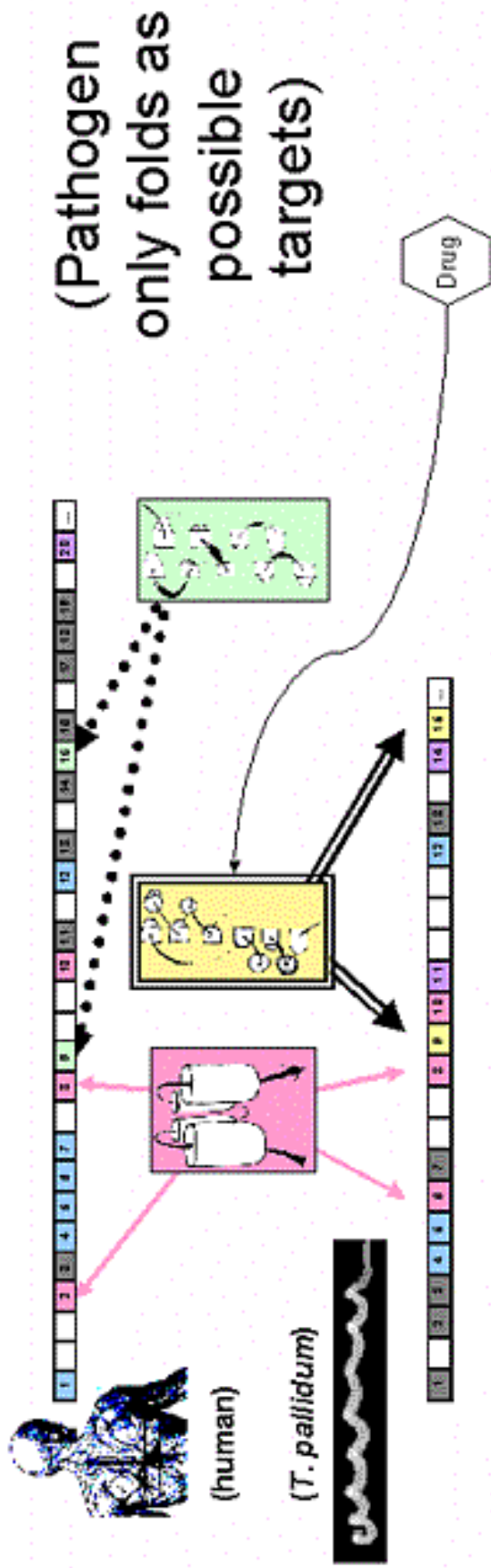
At What Structural Resolution Are Organisms Different?

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



1m 100A 10A 1A

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)



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

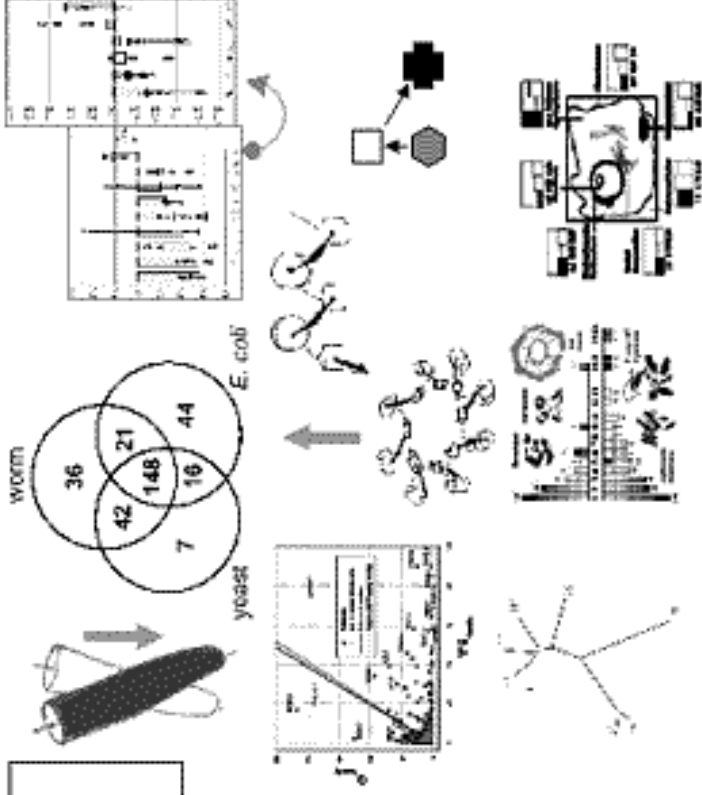
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

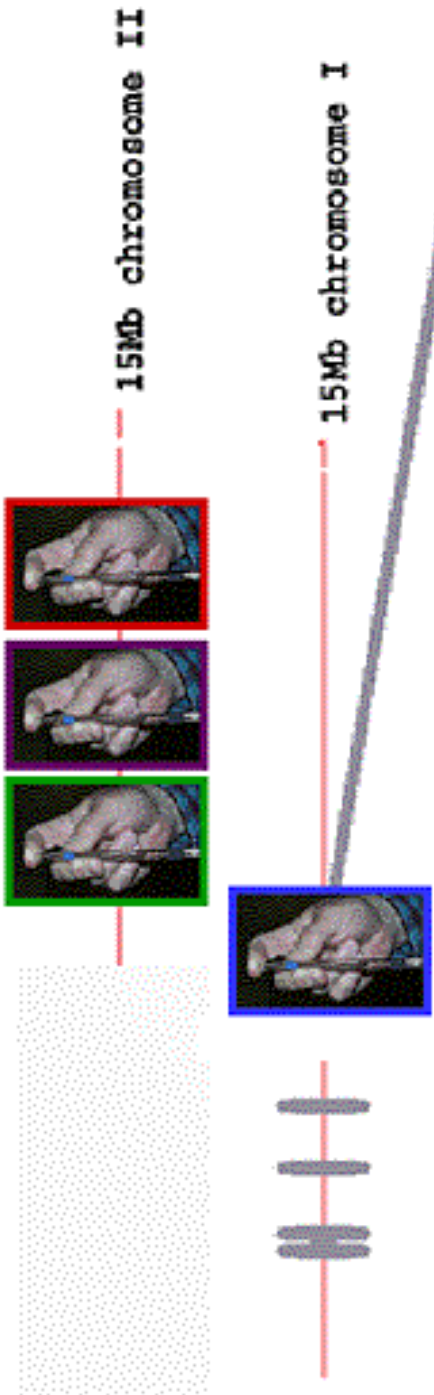
Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

Pseudogenomics: Surveying "Dead" Parts



MTAPMDVDNLSRLLNVGMSGGRLTTSVNEQELOQTCCAVAKSVFASQASLLEVEPP IVC
 GDIHGQYSDLLRIFDKNGFPDVFNLFLGDYVDRGRQNIETICLMLCFKIKYPENFFMLR
 GNHECPA INRVYGFYEECNRRYKSTRLWSIFQDTFNWMLCGLIGSRILCMHGGLSPHLQ
 TLDO LROLPRPODPPNPSIGIDLLWADPDOWVKGWQANTRGVSYVFGODVVADVCSRLDI
DIVARAHOVODGYEEFFASKKMVTIFSAPHYCGGOFDNSAATMKVDENMVCTFVMYKPTPK
 SMRRG *

Pseudogenomics: Surveying "Dead" Parts

15Mb chromosome II

15Mb chromosome I

Example of a potential Ψ G with frameshift in mid-domain

pseudogene fragment on worm chromosome II

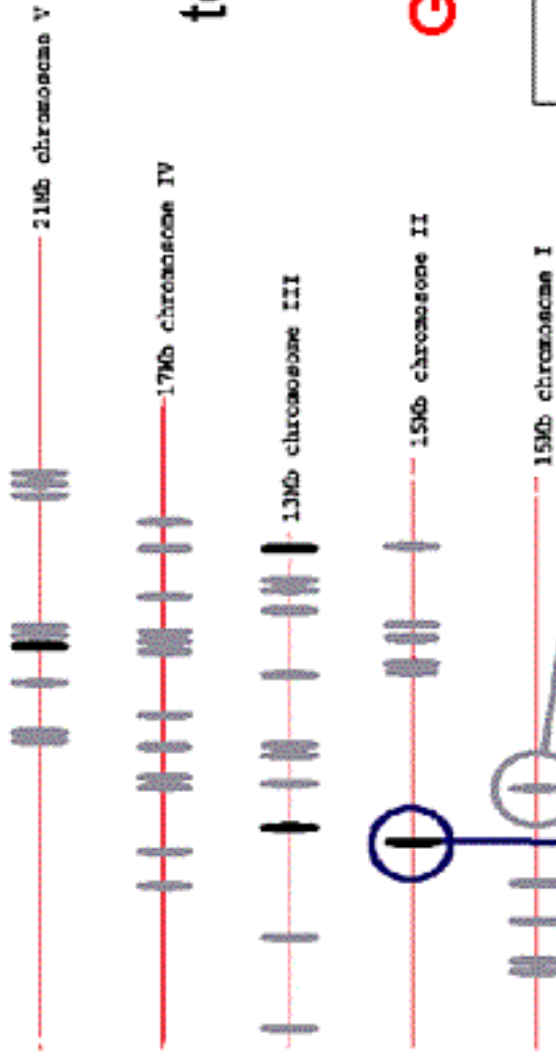
a paralog with the homologous segment highlighted (from chromosome I)
(W09C3.6, serine/threonine protein phosphatase PPI)

TKRTSNGFGQDVVDLFSILDSGLVARAHXVLQDIFEEFAS
KKMVTIFS#APHSPHAPHYCAQFDNSAATVKV

MTAPMDVDNLSRLLNVGMSGGRLTTSVNEQEQLQCCAVAKSVFASQASLLEVEPPIVC
GDIHGQYSDLLRIFDKXNGFPDVFNLFLGDYVDRGRQNIETICLMLCFKIKYPENFFMLR
GNHECPA INRVYGFYEENRRYKSTRLWSIPQDTFNWMLCGLIGSRILCMHGGLSPHLQ
TLDOLRQLPRPODPPNPSIGIDLLWADPDQWVKGMQANTRGVS YVFGODVVADVCSRLDI
DIVARAHOVODGYEEFASKKMVTIFSAPHYCGGFDNSAATMKVDENMVCTFVMYKPTPK
SMRRG*

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

Folds in Pseudogenes



pseudogene fragment on worm chromosome II

```

TKSTNGFGDEVVVDLPSLDSGLVARAKXVLQDIFEPAE
EAMTIFE#AFESFESHAPHYCHAGFENSAATYKY
MIAFNDVDELSRLALNVGSGGLTITSNEQLQTCACAAKSVFASQASLLEVEPPPVC
GEIMGQYSDLLRIFDNGFFPDVNFLLGQVYDGRONITICLMLCFKIKYDNPFLER
GHECPALNRYVIGFEECMREKSTRLMSIFQDTNNMPLQGLIGSSILCMBGGLSPHLQ
TLQLEQLPFPQDFPHFSEGIDLDAADPEQVYKQVAFKGVVFGDDVADGCEERLH
DQKARAHGKVVGGVGHFHSASLKHPTFHSLSHYGGGGRSANNKNDENNVCIFVNIATFPE
SKRGG*
    
```

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

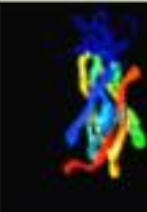



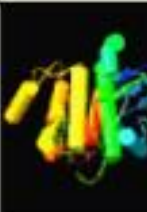


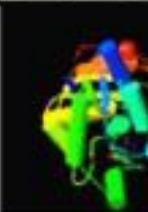


Example of a potential Ψ G with frameshift in mid-domain

Ψ G identification pipeline to Summary of Pseudogenes in worm

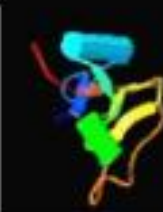

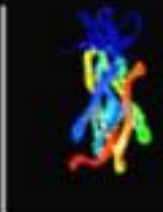

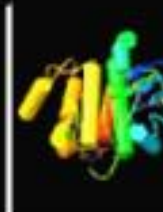

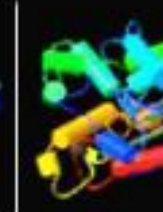



G=19K **G_E=8K** **Ψ G=4K (2K)**

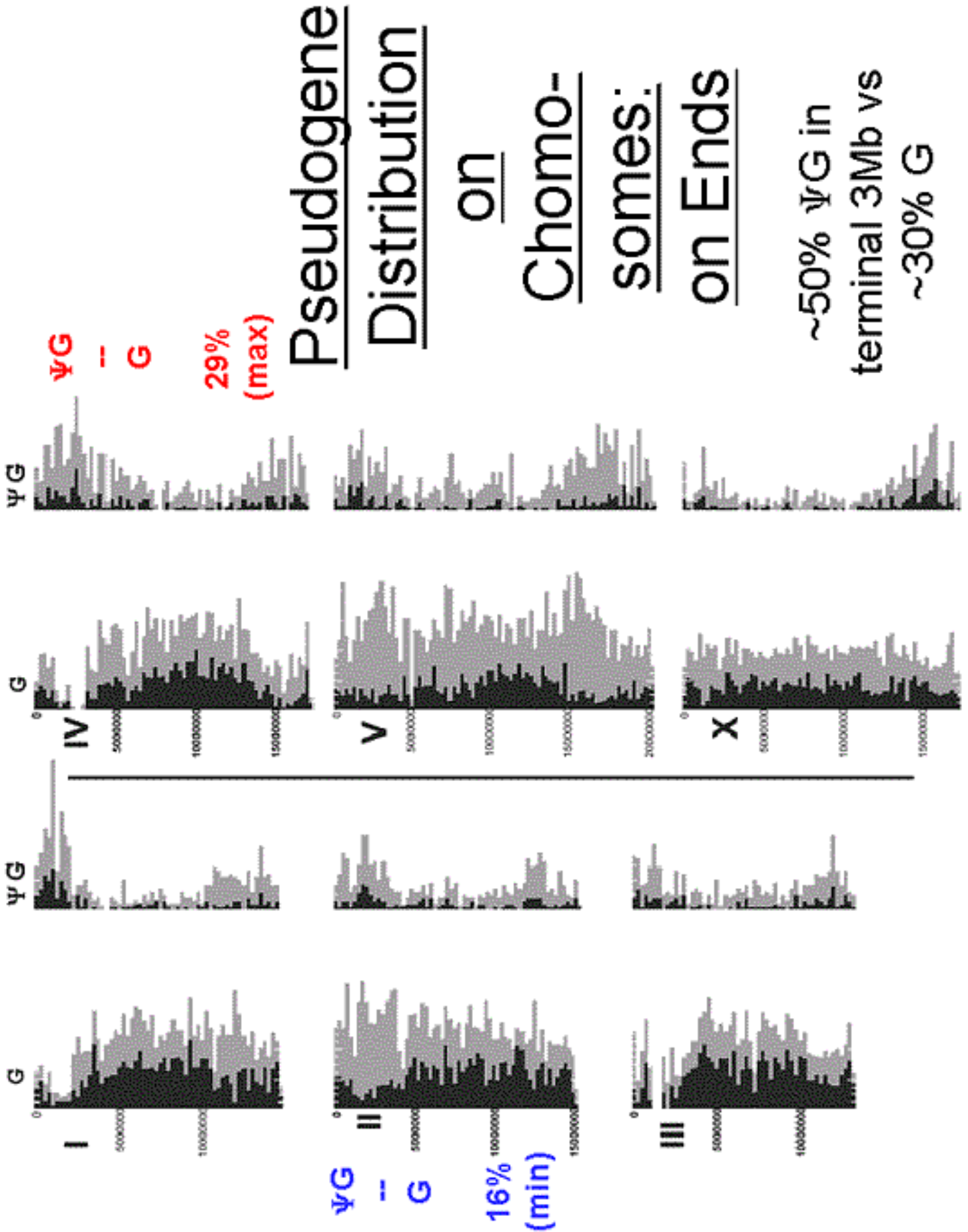
Category	Total number	Number for genes with EST match	Genes with EST match as percentage of Category	Number for genes in paralog families with EST match	Genes in paralog families with EST match as percentage of Category
Genes	18,576 (G)	7,829 (G _E)	42%	13,417 (G _p)	72%
Pseudogenes and pseudogene fragments	3,814 (Ψ G)	2,788 (Ψ G _E)	47%	2,729 (Ψ G _p)	72%
Singletons	637 (17% of Ψ G)	233	36%	---	---
Intronic pseudogenes*	1,155 (30% of Ψ G)	351	30%	704	61%

Most Common Worm "Pseudofolds" #1

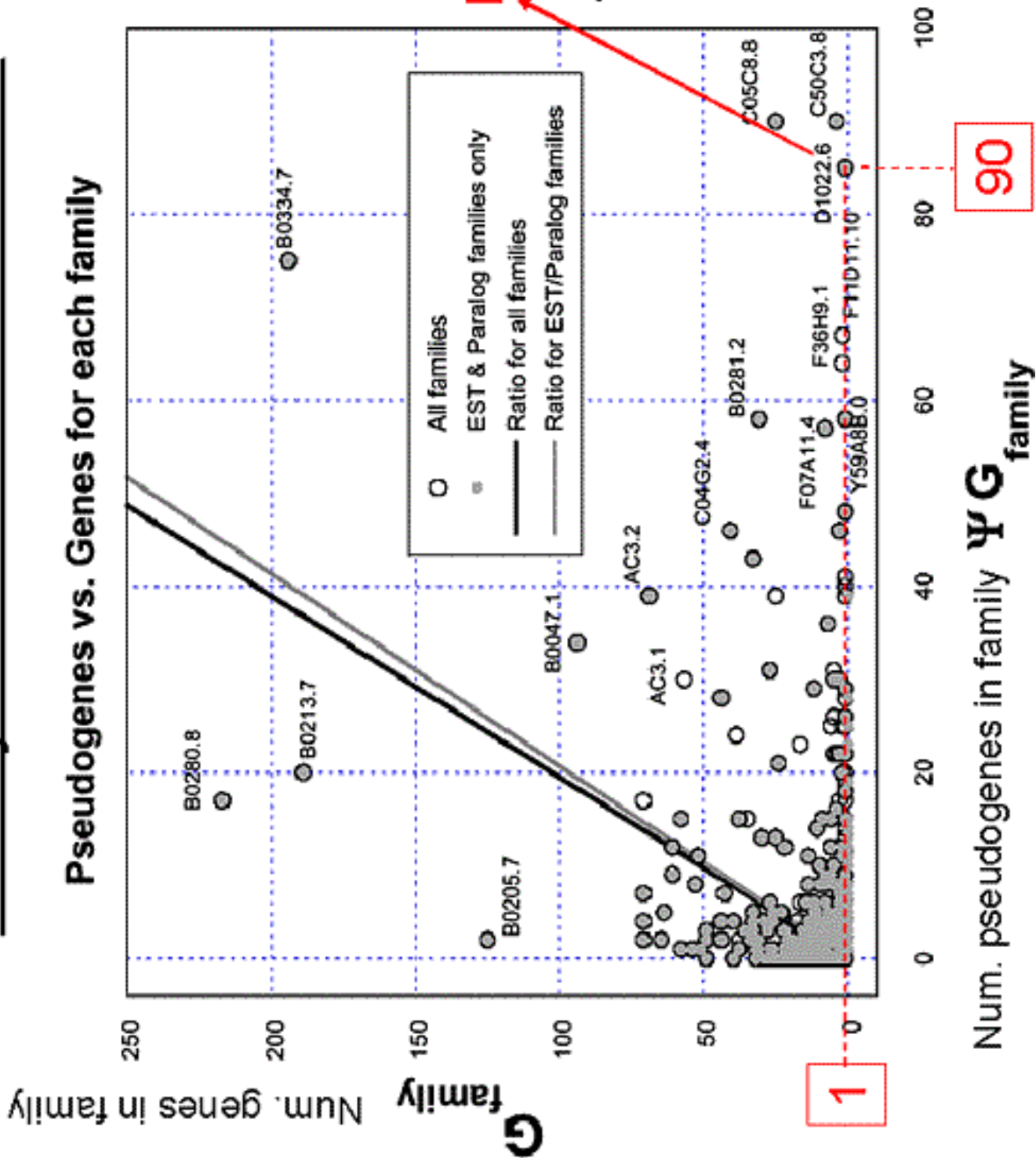
G Rank (Number matches)	Ψ G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description	G Rank (Number matches)	Ψ G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description
1 (769)	2		d1ajw_ 2.1 Immunoglobulin	6 (246)	8		d2lbd_ 1.95 Nuc. receptor ligand-binding domain
2 (555)	6		d1dec_ 7.3 Knottin	7 (243)	34		d1a17_ 1.91 Alpha/alpha superhelix
3 (434)	3		d3lck_ 5.1 Protein kinase	8 (227)	17		d1sp2_ 7.31 Classic zinc finger
4 (302)	1		d1tsg_ 4.105 C-type lectin	9 (215)	20		d1dai_ 3.29 P-loop NTP hydrolase
5 (274)	7		d1zfo_ 7.33 Glucocorticoid receptor DNA-binding dom.	10 (197)	13		d2aw0_ 4.34 Ferredoxin

Most Common Worm "Pseudofolds" #2

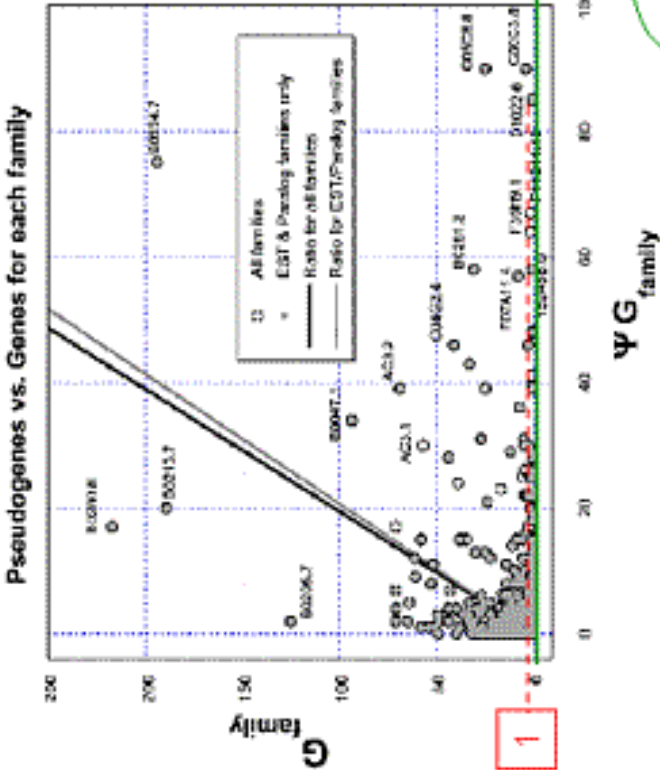
Ψ G Rank (Number matches)	G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description	Ψ G Rank (Number matches)	G Rank	Fold	Representative Domain, SCOP 1.39 Number, Description
1 (39)	4		d1tsg___ 4.105 C-type lectin	6 (18)	2		d1dec___ 7.3 Knottin
2 (32)	1		d1ajw___ 2.1 Immunoglobulin	7 (17)	5		d1zfo___ 7.33 Glucocorticoid receptor DNA-binding dom.
3 (27)	3		d3lck___ 5.1 Protein kinase	8 (15)	6		d2lbd___ 1.95 Nuc. receptor ligand-binding domain
4 (25)	11		d1cvl___ 3.56 Alpha/beta-hydrolase	9 (13)	58		d1bus___ 7.14 Ovomucoid PCI inhibitor fold
5 (23)	63		d1ako___ 4.93 DNase-I fold	9 (13)	19		d2bnh___ 3.7 Leu-rich, right-handed β/α superhelix



Decayed Lines of Genes?

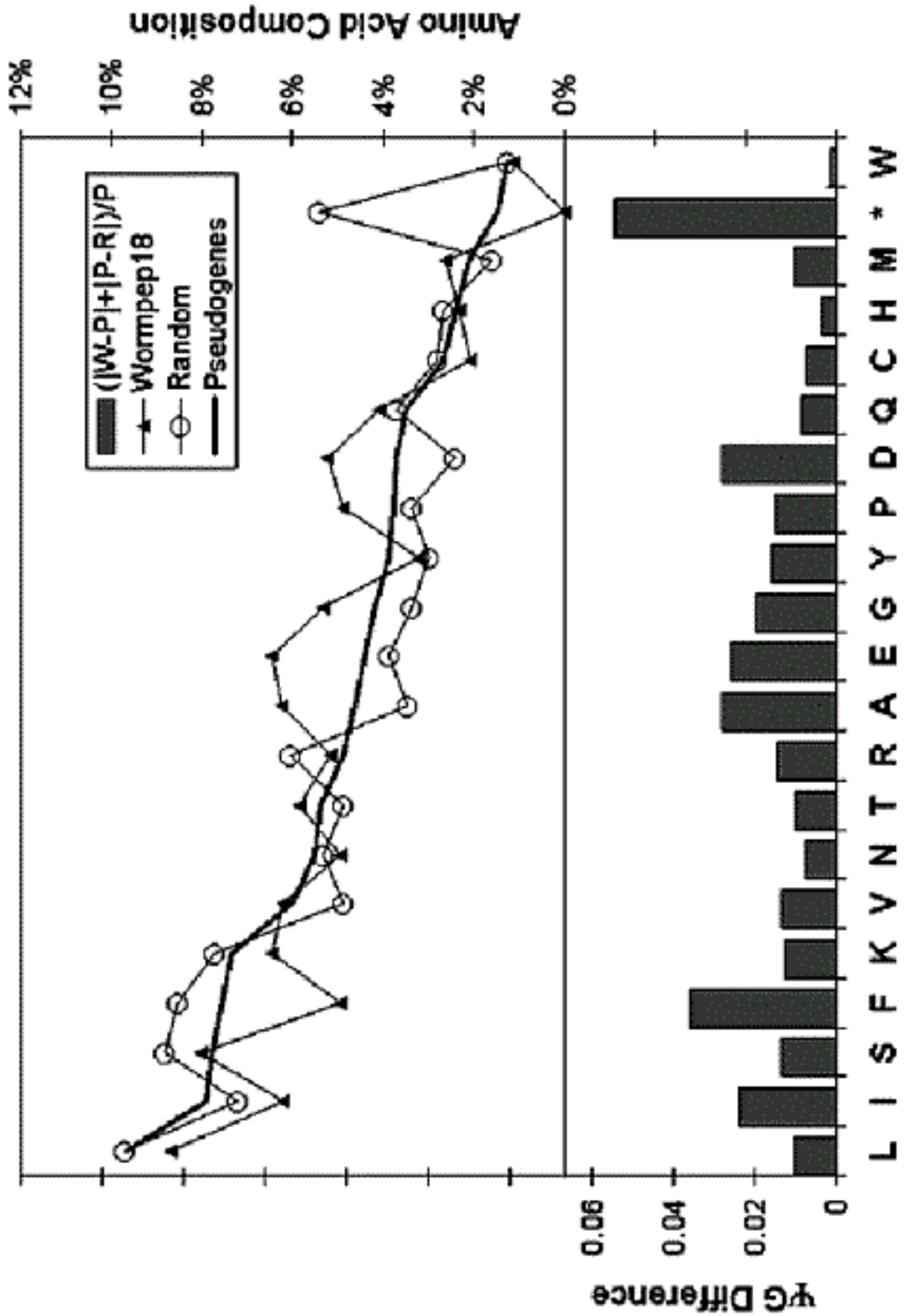


Completely Dead Families



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	Acridone repressor

Amino Acid Composition of Pseudogenes is Midway between Proteins and Random



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

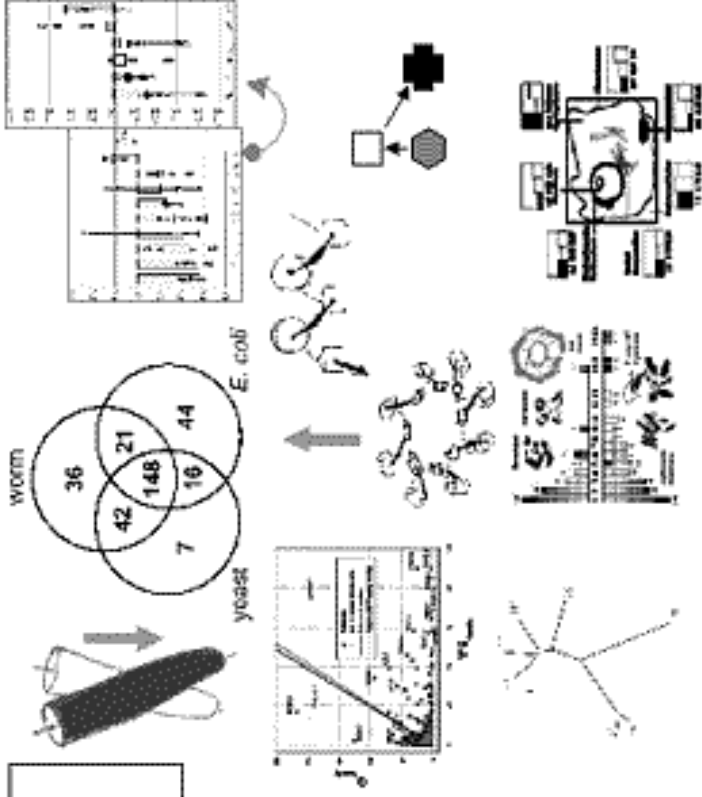
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

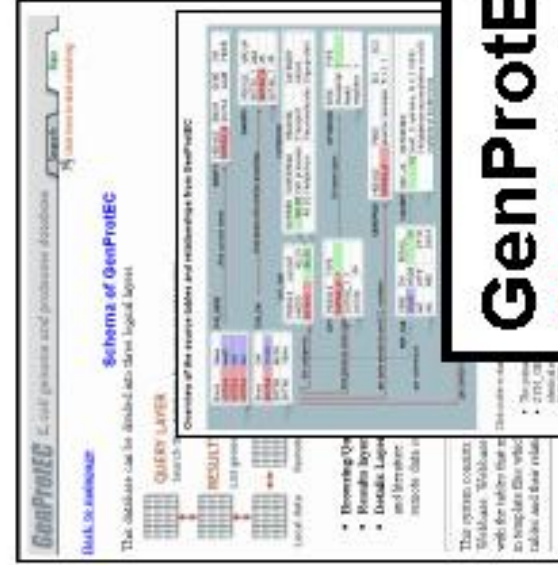
Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

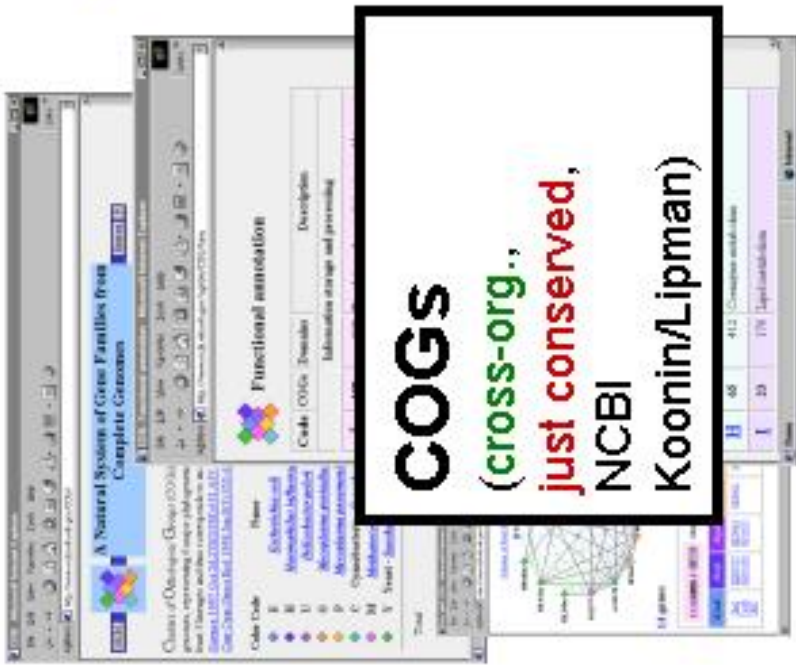
Functional Classification



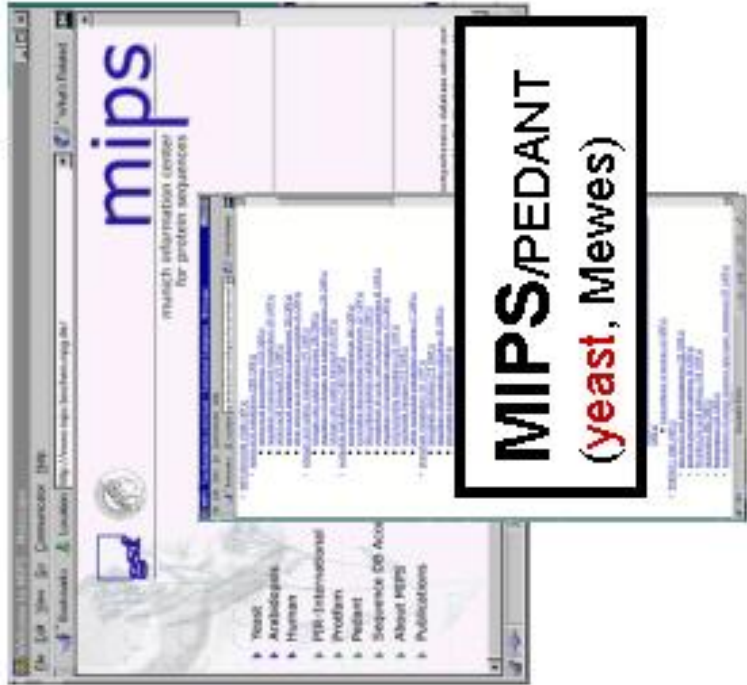
GenProtEC
(*E. coli*, Riley)

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

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



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



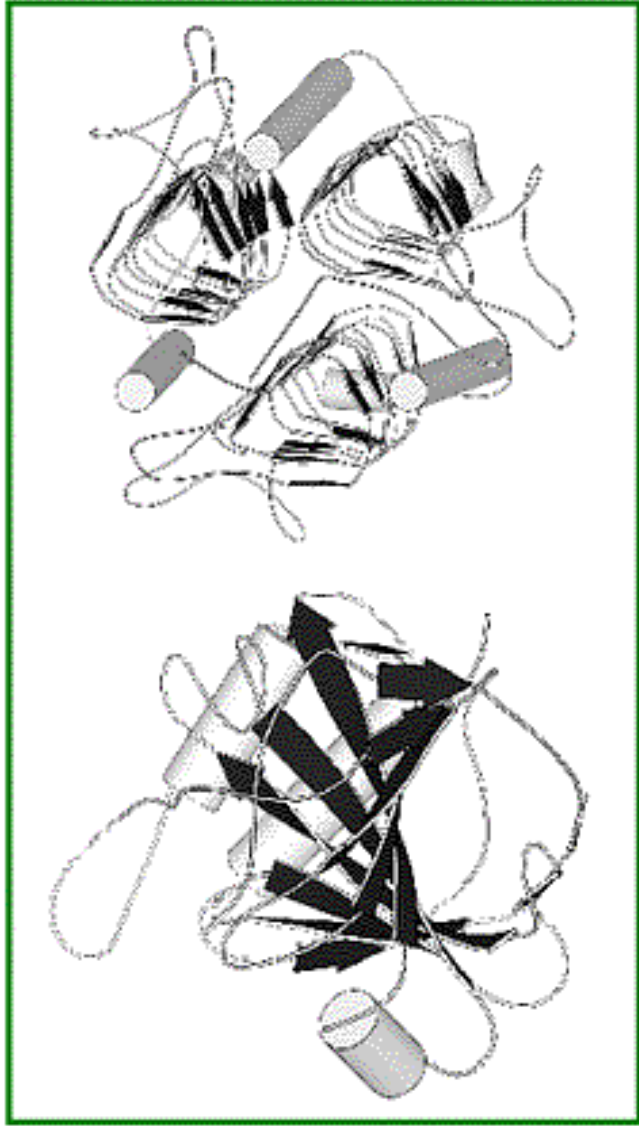
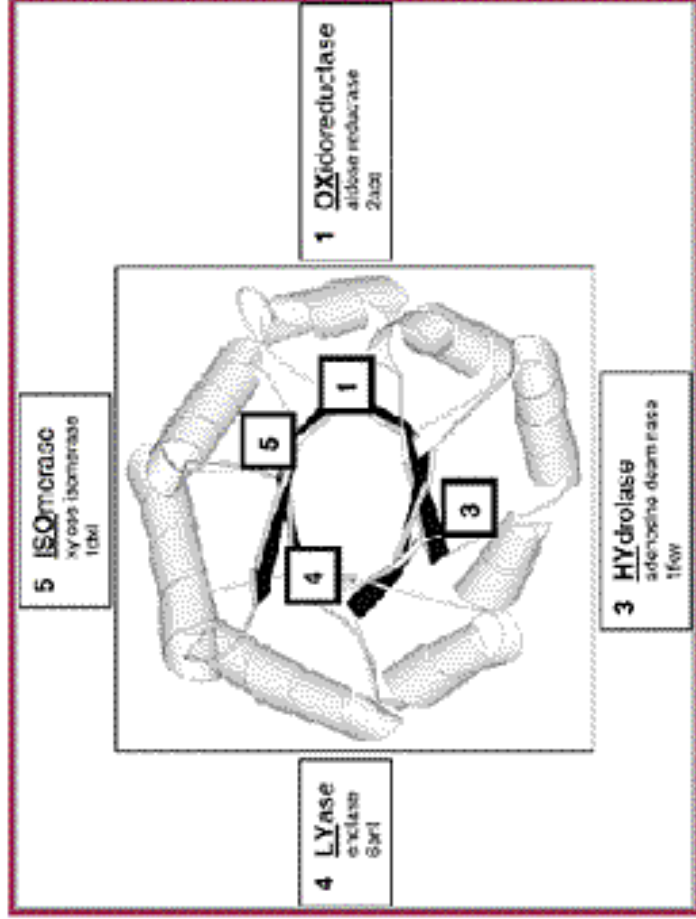
MIPS/PEDANT
(yeast, Mewes)



"FLY"
(fly, Ashburner)
now extended to
GO (cross-org.)

Fold-Function Combinations

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

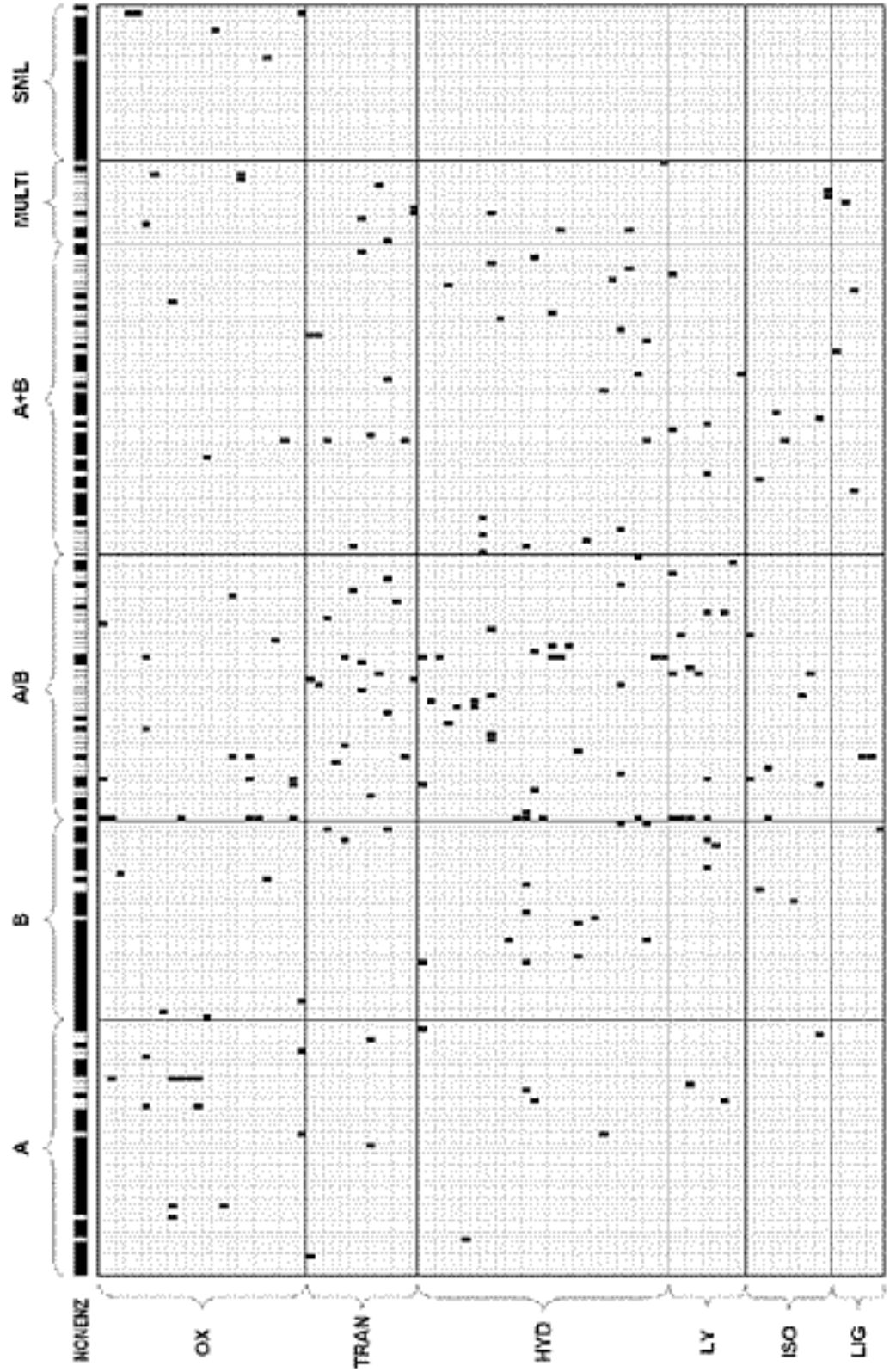


Two Different Folds
Catalyze the Same
Reaction -- e.g.
Carbonic Anhydrases
(4.2.1.1)

Fold-Function Combinations

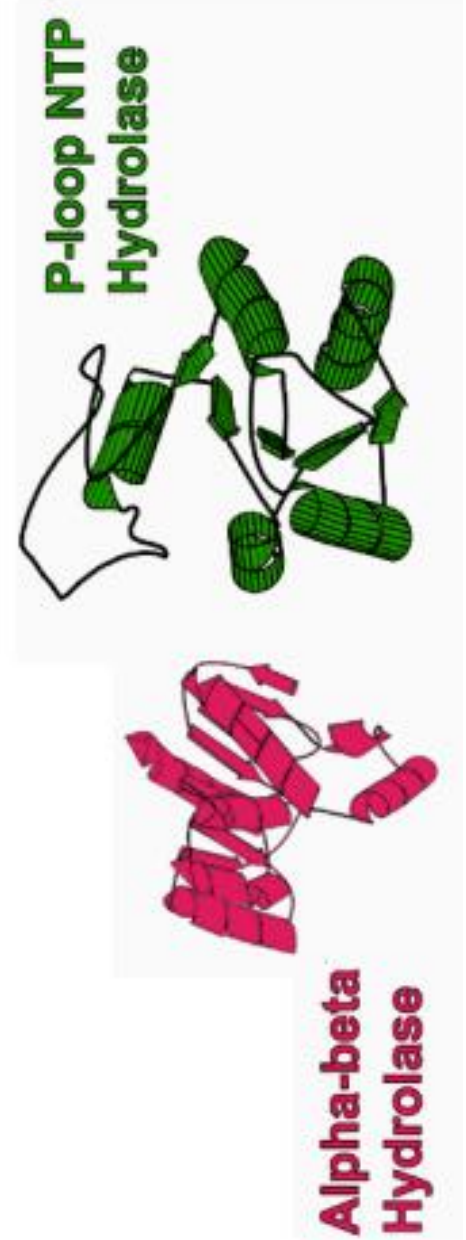
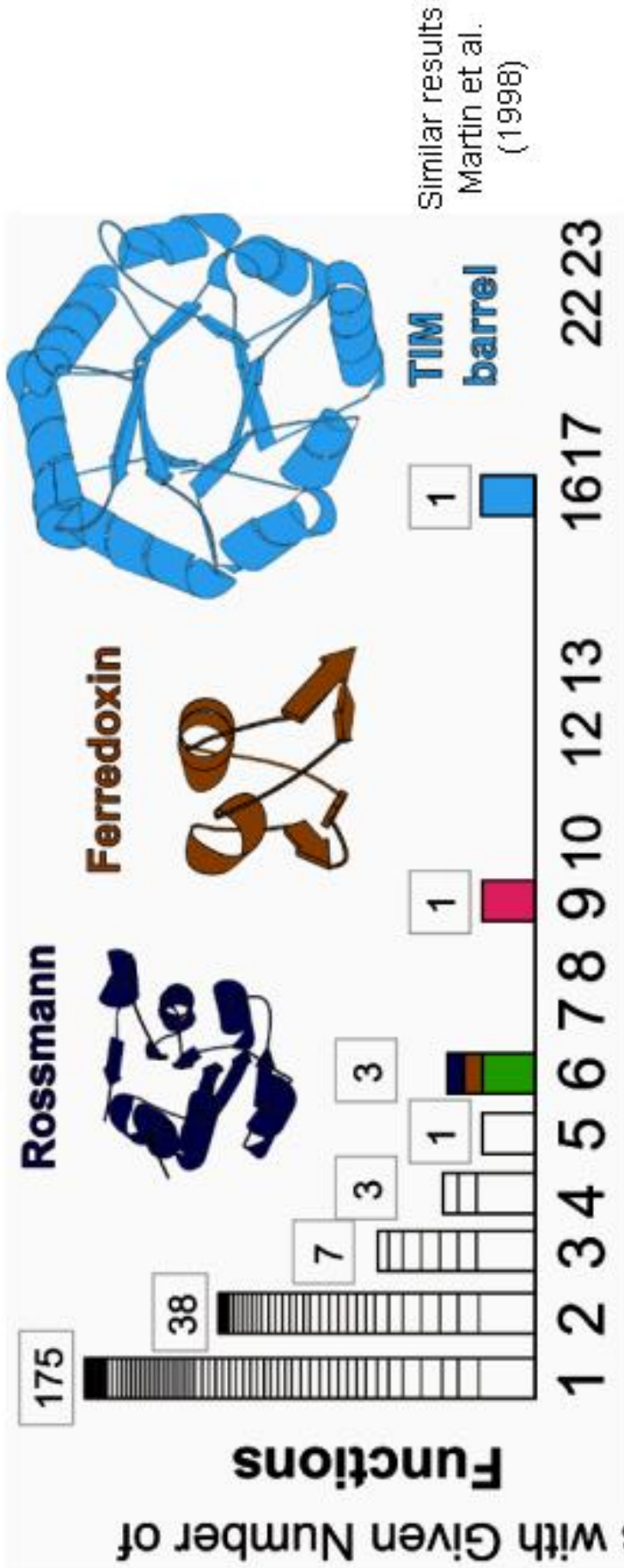
~20K (=92x229) Possible,
331 Observed

229 Folds

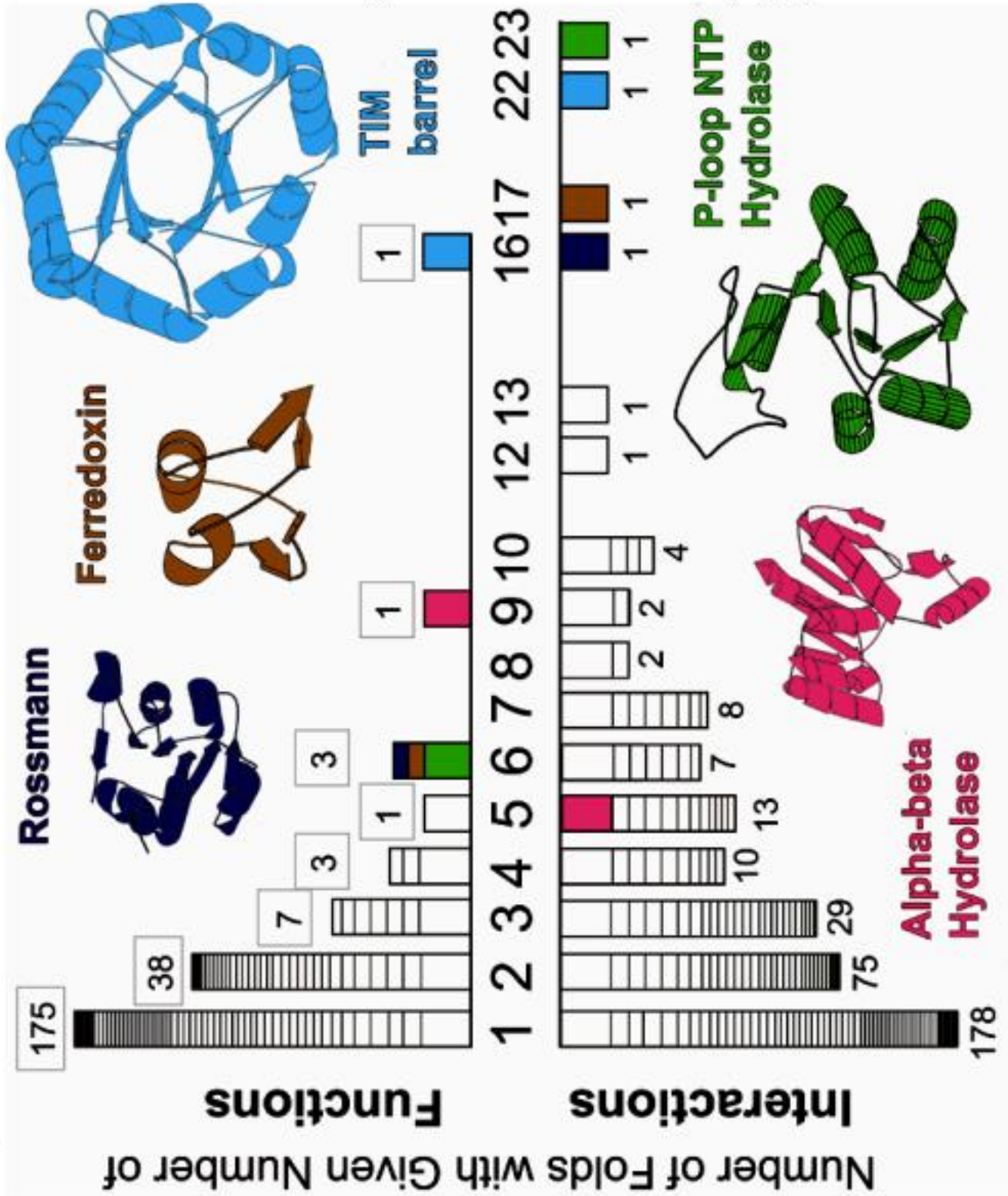


91 Enzymatic Functions
+ Non-Enzyme

Most Versatile Folds



Most Versatile Folds – Relation to Interactions



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

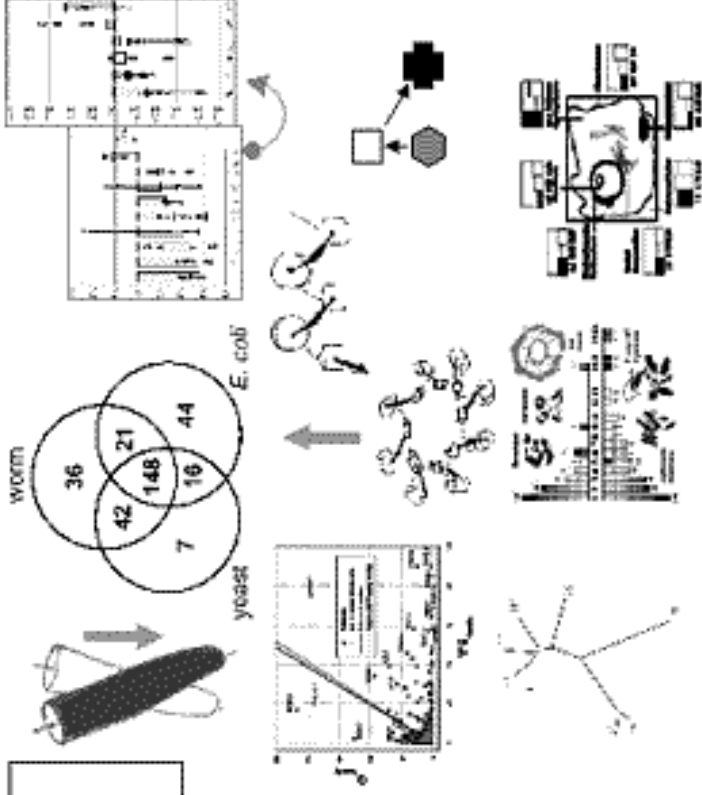
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

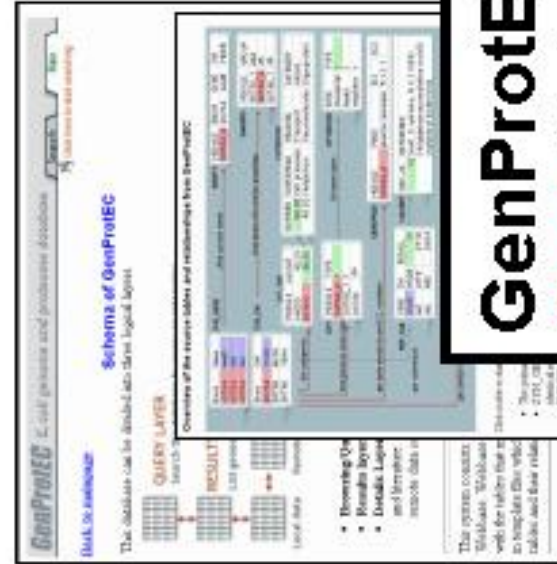
Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

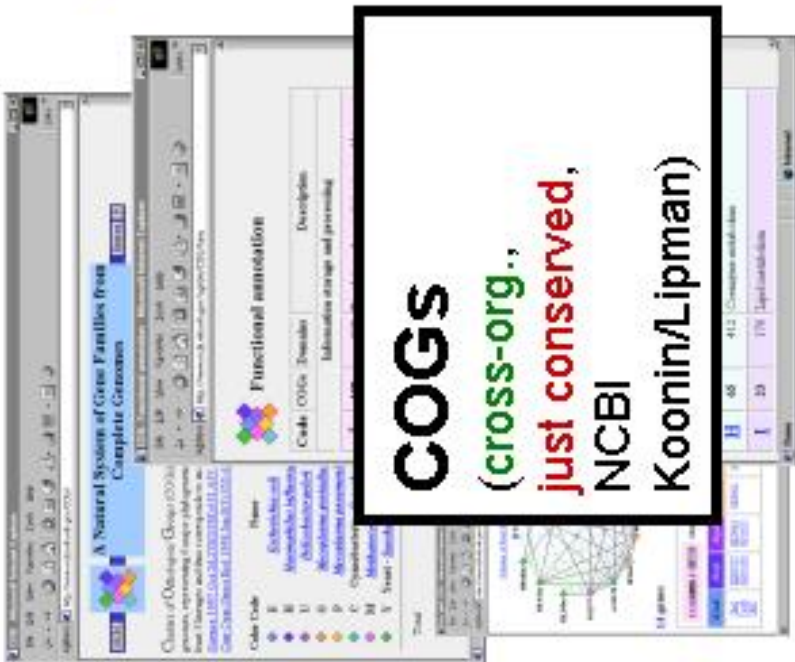
Functional Classification



GenProtEC
(*E. coli*, Riley)

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

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



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

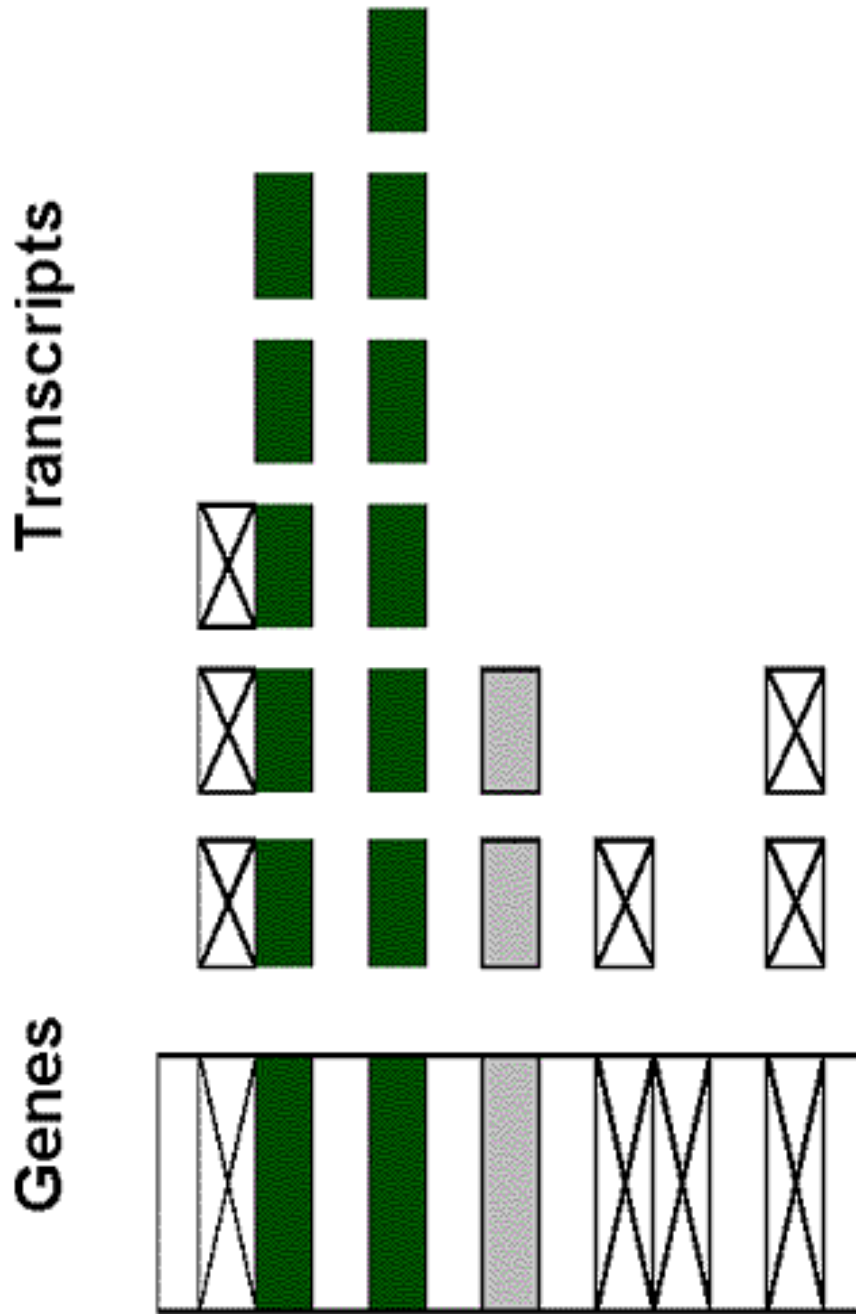


MIPS/PEDANT
(yeast, Mewes)



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

Transcriptome Composition



Common Folds in the Transcriptome

Yeast Genome Rank

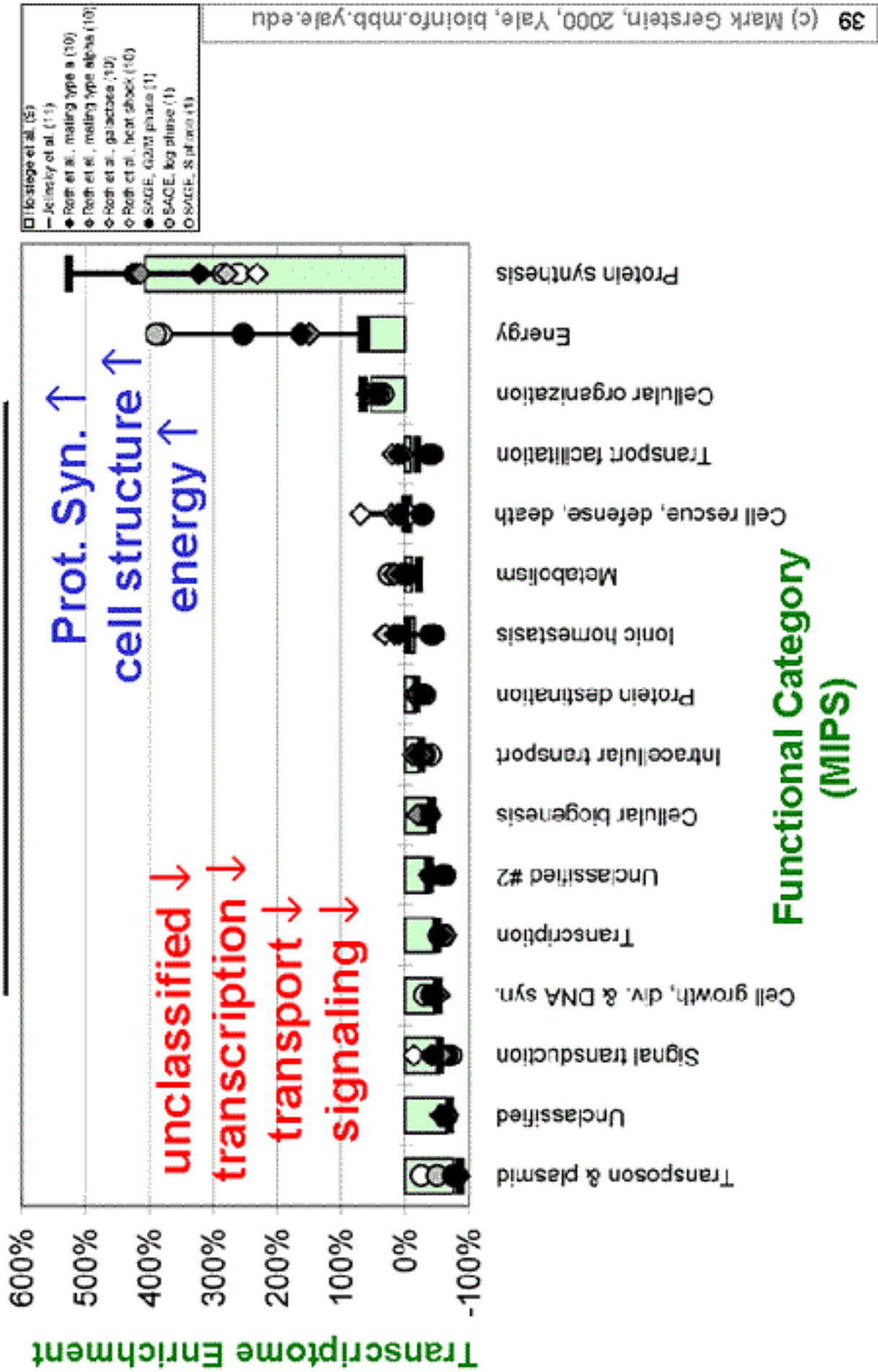
Transcriptome Rank

(partslist.org)

beta/alpha (TIM)-barrel d1aj2 :3.1 (1.39) A/B	5	1
P-loop containing nucleotide triphosphate hydrolases d1dai :3.39 (1.39) A/B	3	2
Ferredoxin-like d2aw0 :4.34 (1.39) A+B	6	3
NAD(P)-binding Rossmann-fold domains d1eny :3.22 (1.39) A/B	8	4
7-bladed beta-propeller d1gotb :2.51 (1.39) B	2	5
alpha-alpha superhelix d1a17 :1.91 (1.39) A	4	6
Thioredoxin fold d1mek :3.38 (1.39) A/B	14	7
Glyceraldehyde-3-phosphate dehydrogenase-like d1drw :2.42 (1.39) A+B	78	8
beta-Grasp d1tif :4.11 (1.39) A+B	36	9
Heat shock protein 70kD (HSP70) d1dkra :5.17 (1.39) M	31	10

Reductase/isomerase/elongation factor common domain d1dfr :1.232 (1.39) B	60	11
Serpins d1psl :5.2 (1.39) M	78	12
N-terminal nucleophile aminohydrolases (Ntn hydrolases) d1neda :4.95 (1.39) A+B	18	13
Metallothionein d1aoo :7.38 (1.39) S	96	14
Oligomers of long helices (LEU-Zip) d2ifo :1.105 (1.39) A	7	15
DNA/RNA polymerases d1har :5.9 (1.39) M	25	16
Ribonuclease H-like motif d2lqg :3.47 (1.39) A/B	15	17
Protein kinases (PK) d3lck :5.1 (1.39) M	1	18
Class II aaRS and biotin synthetases d1py8 :4.61 (1.39) A+B	21	19
PLP-dependent transferases d2lkb :3.54 (1.39) A/B	12	20
DNA/RNA-binding 3-helical bundle d1a5j :1.14 (1.39) A	15	21

Composition of Transcriptome in terms of Functional Classes



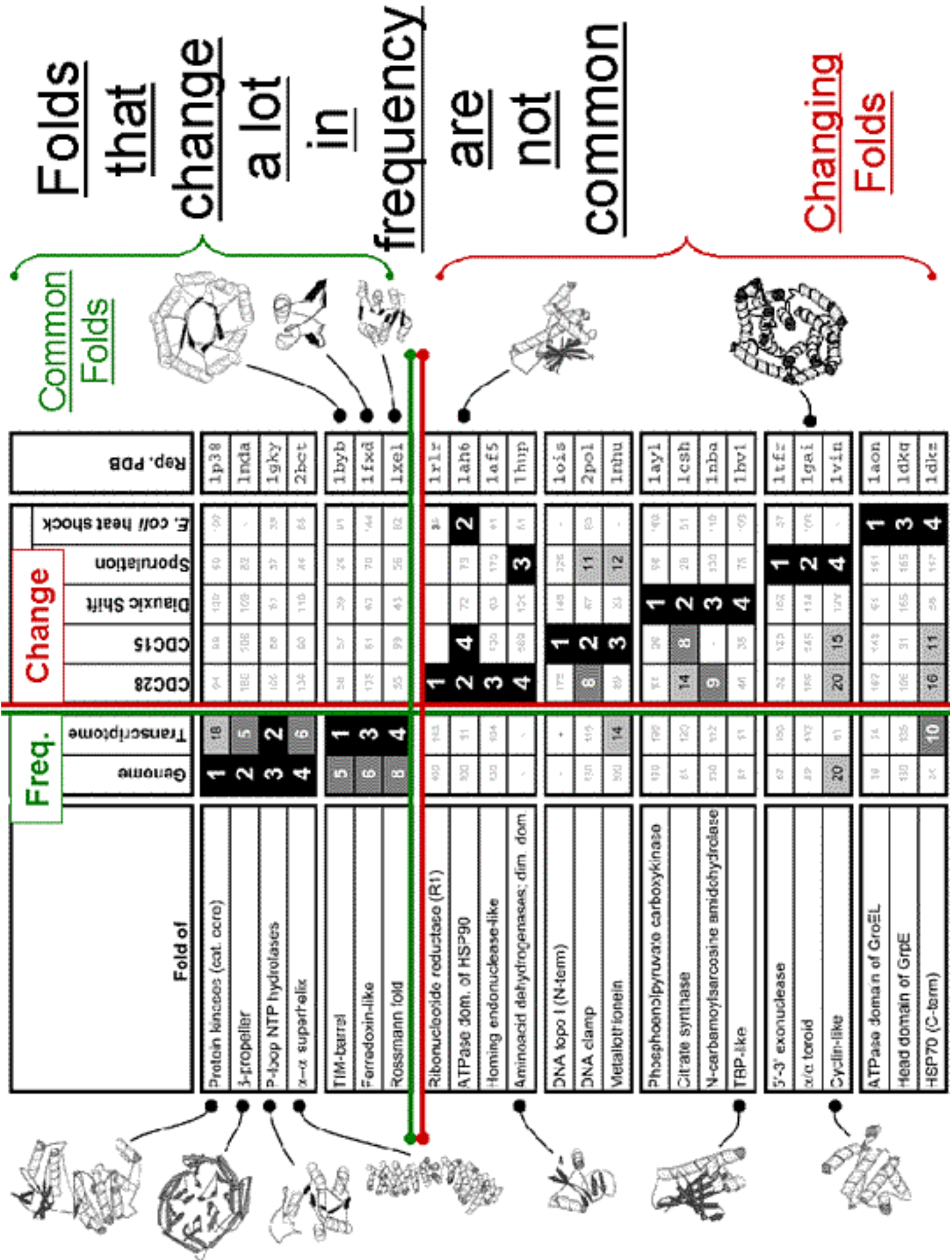
Common Folds in the Transcriptome #2

(partslist.org)

Yeast Genome Rank
Transcriptome Rank

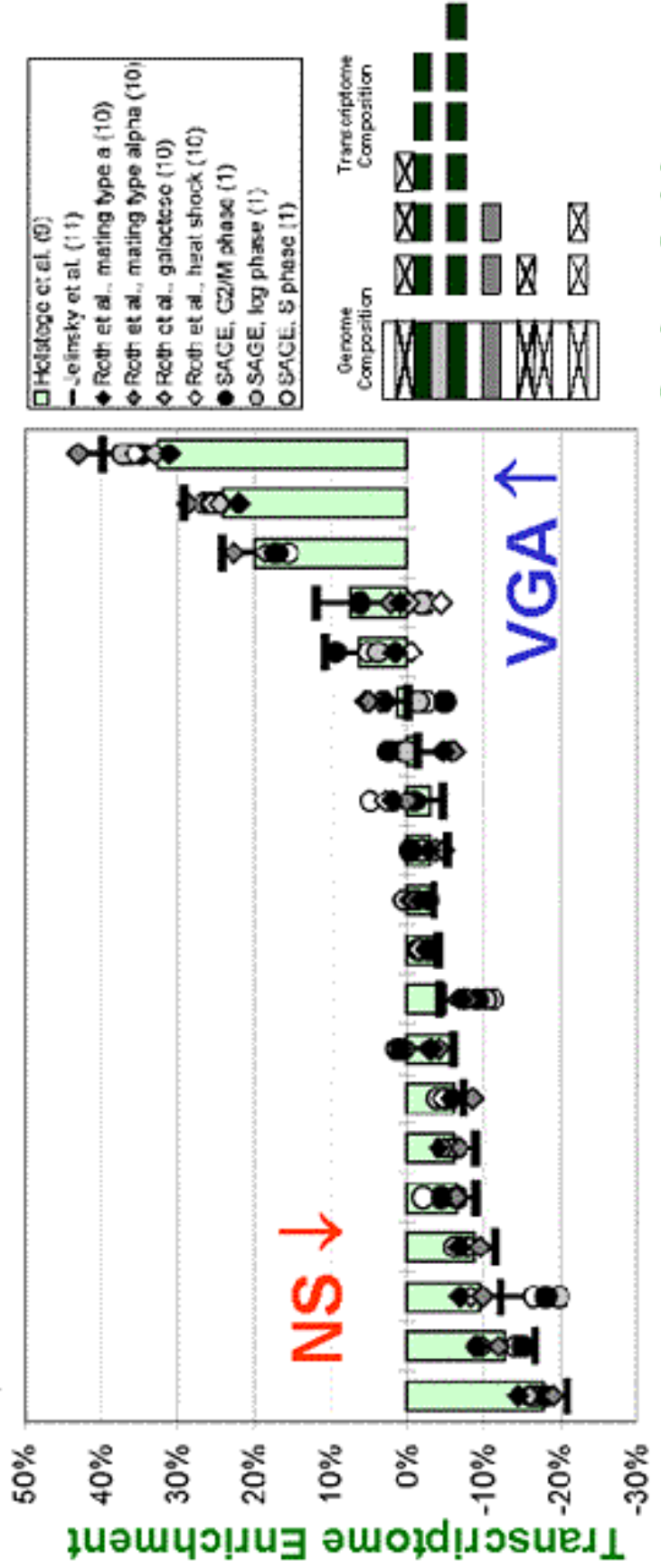
beta/alpha (TIM)-barrel d1aj2 :3.1 (1.39) A/B	5	1	<Met.↑
P-loop containing nucleotide triphosphate hydrolases d1dai :3.39 (1.39) A/B	3	2	<Met.↑
Ferredoxin-like d2aw0 :4.34 (1.39) A+B	6	3	<Met.↑
NAD(P)-binding Rossmann-fold domains d1eny :3.22 (1.39) A/B	8	4	<Met.↑
7-bladed beta-propeller d1gotb :2.51 (1.39) B	2	5	↓Tfac.>
alpha-alpha superhelix d1a17 :1.91 (1.39) A	4	6	
Thioredoxin fold d1mek :3.38 (1.39) A/B	14	7	
Glyceraldehyde-3-phosphate dehydrogenase-like d1drw :2.42 (1.39) A+B	78	8	↓Sign.> <Met.↑
beta-Grasp d1tif :4.11 (1.39) A+B	36	9	
Heat shock protein 70kD (HSP70) d1dkra :5.17 (1.39) M	31	10	↓Tfac.>

Reductase/isomerase/elongation factor common domain d1dfr :2.32 (1.39) B	60	11
Serpins d1psl :5.2 (1.39) M	78	12
N-terminal nucleophile aminohydrolases (Ntn hydrolases) d1neda :4.95 (1.39) A+B	18	13
Metallothionein d1aoo :7.38 (1.39) S	96	14
Oligomers of long helices d2ifo :1.105 (1.39) A (LEU-zip)	7	15
DNA/RNA polymerases d1har :5.9 (1.39) M	25	16
Ribonuclease H-like motif d2lrg :3.47 (1.39) A/B	15	17
Protein kinases (PK) d3lck :5.1 (1.39) M	1	18
Class II aaRS and biotin synthetases d1pyr :4.61 (1.39) A+B	21	19
PLP-dependent transferases d2lkb :3.54 (1.39) A/B	12	20
DNA/RNA-binding 3-helical bundle d1a5j :1.14 (1.39) A	15	21



Composition of Genome vs. Transcriptome

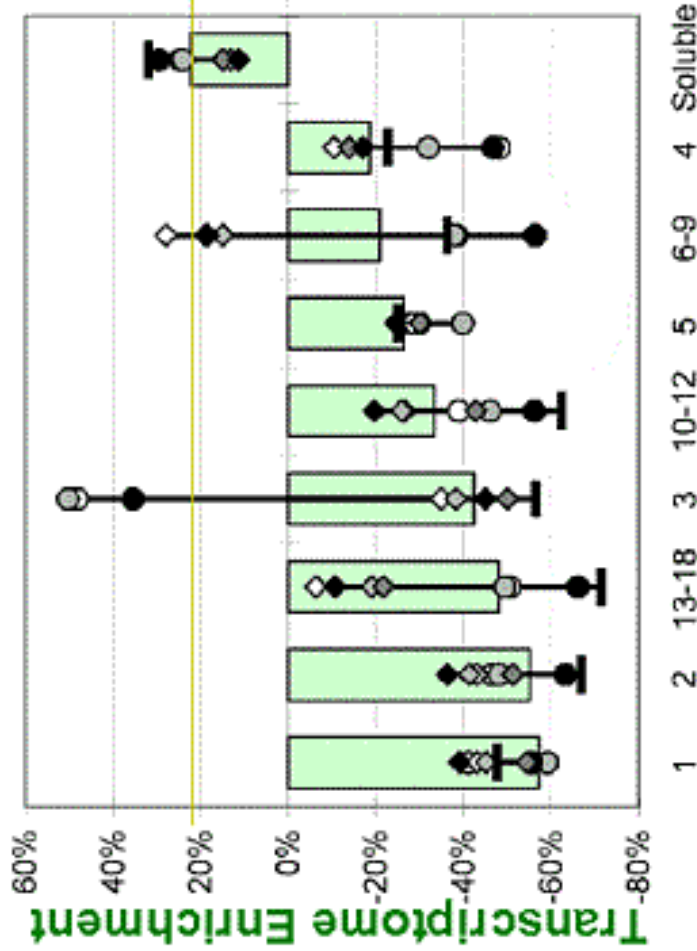
	$\sum_{\text{set } i} n_i(F)$	$\sum_{i=1}^n n_i(F)$	$G(F)$	$\sum_{i=1}^n e_i n_i(F)$	$\sum_{i=1}^n e_i n_i(F)$	$T(F)$	$D(F)$
Feature F is Amino acids, in particular Ala	Number of Ala in yeast	Number of amino acids in yeast	Genome composition of Ala in yeast	Number of Ala weighted by expression	Number of amino acids weighted by expression	Transcriptome composition of Ala in yeast	Relative enrichment of Ala in <u>transcriptome</u>
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 <u>transcriptome</u>
Spec. Num.	65	1560	4.2%	389	4709	8.3%	97.8%



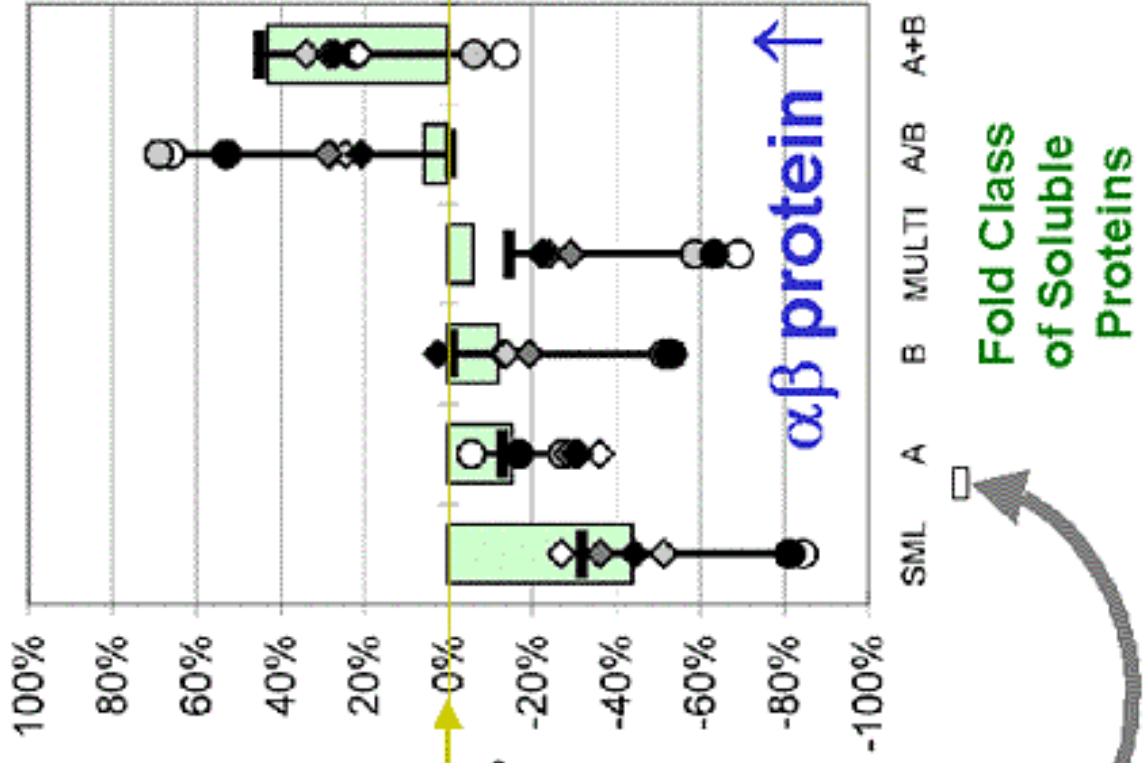
Composition of Transcriptome in terms of Broad Structural Classes

- ▣ Holstege et al. (8)
- Jaffe et al. (5)
- ◆ 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 ↓



TM helices in yeast protein



Fold Class of Soluble Proteins

Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

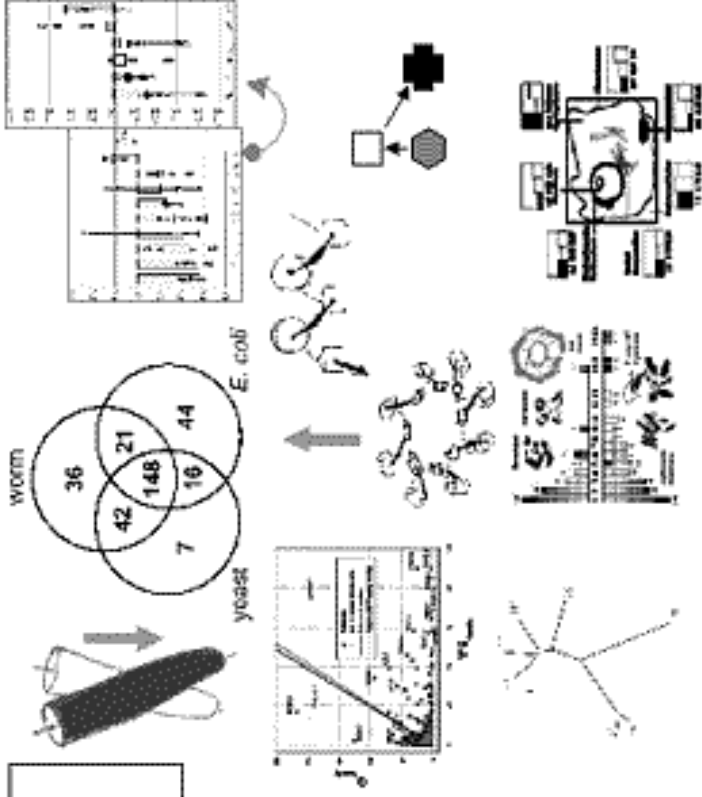
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

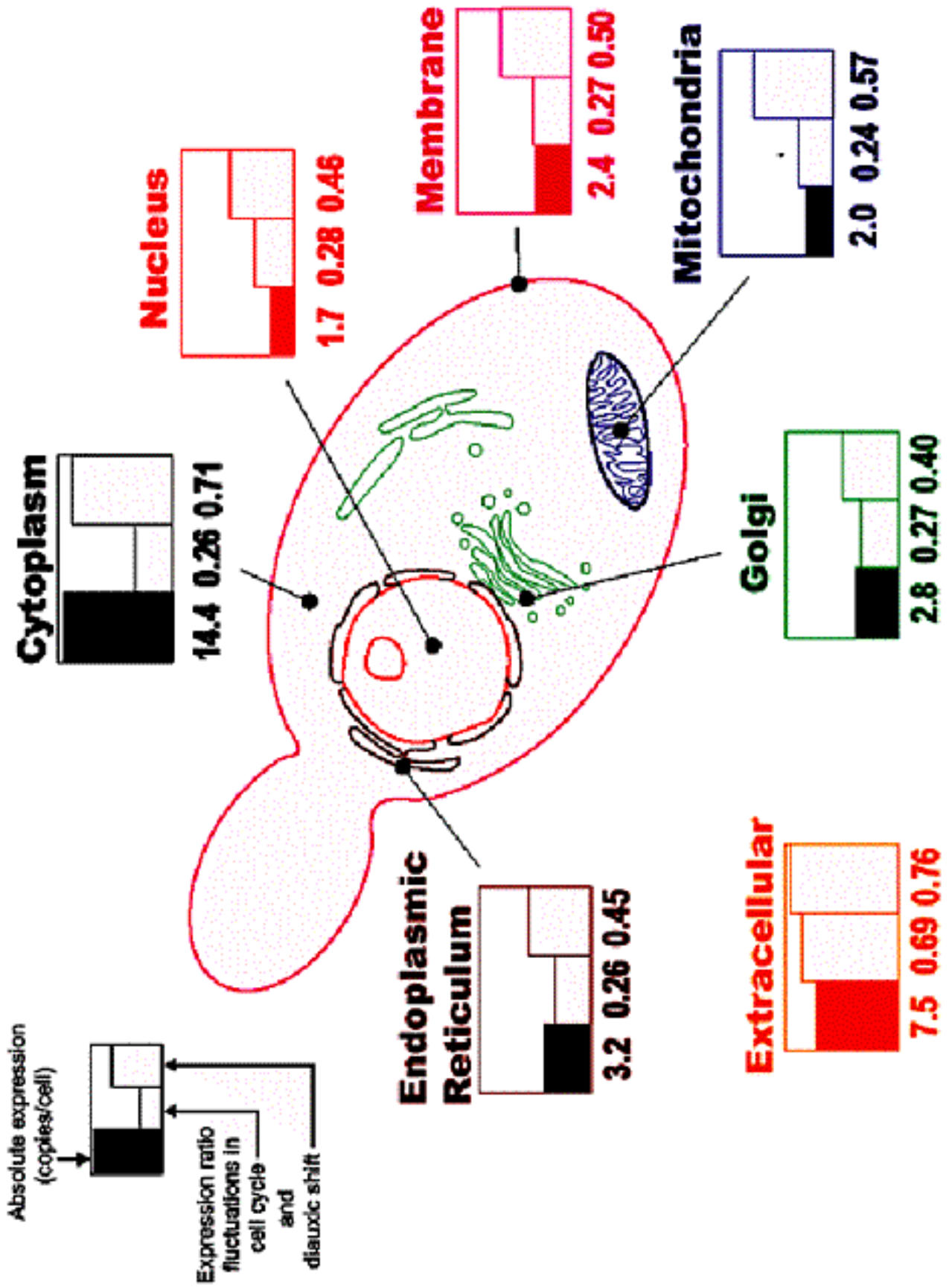
Integrated Views based on Parts



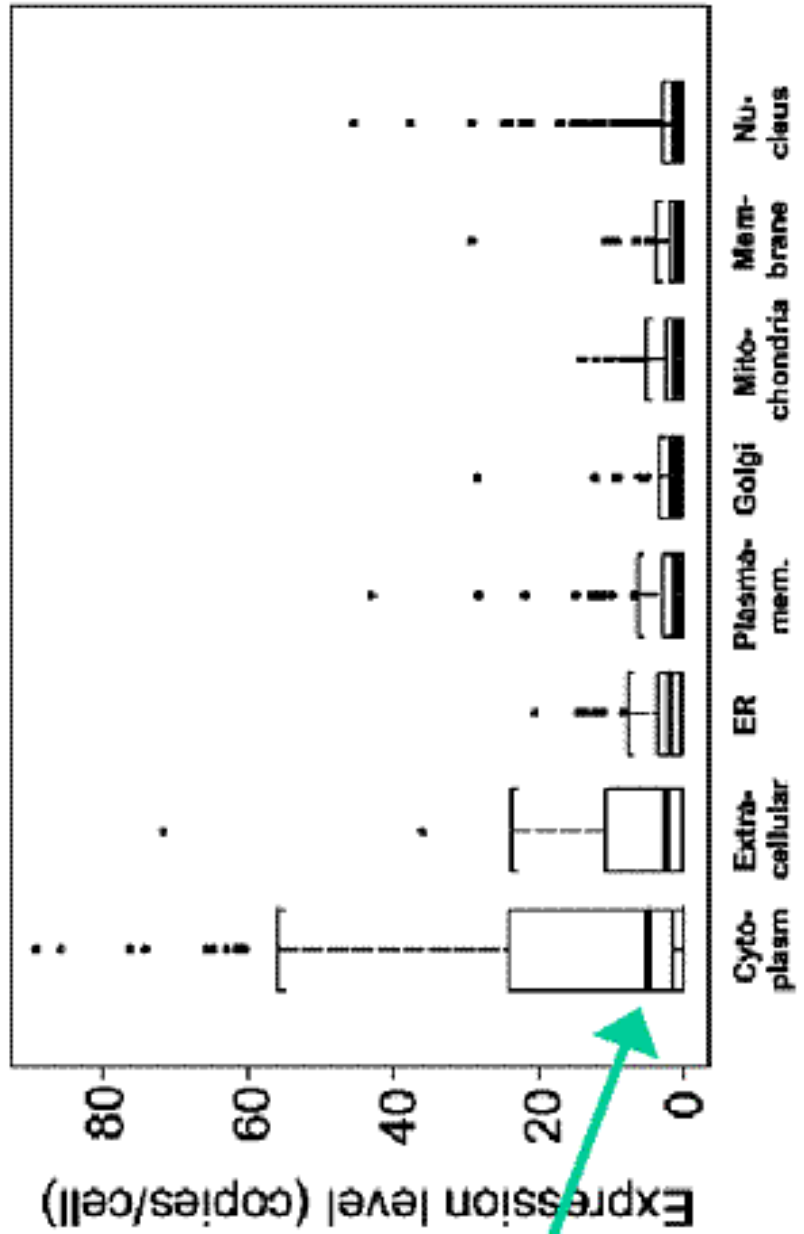
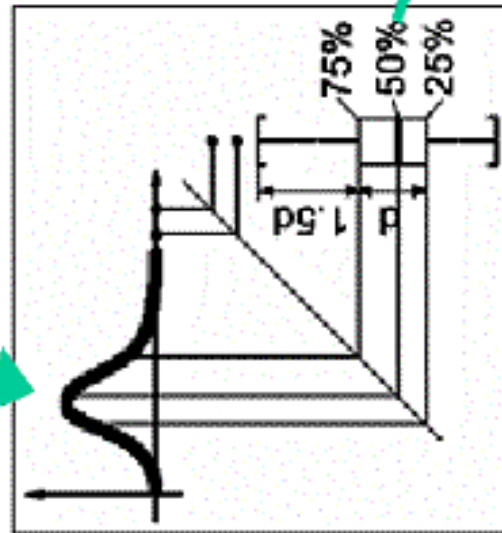
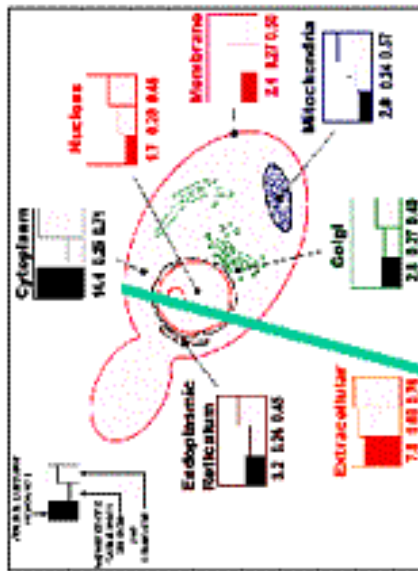
*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

Expression Level is Related to Localization



Distributions of Expression Levels

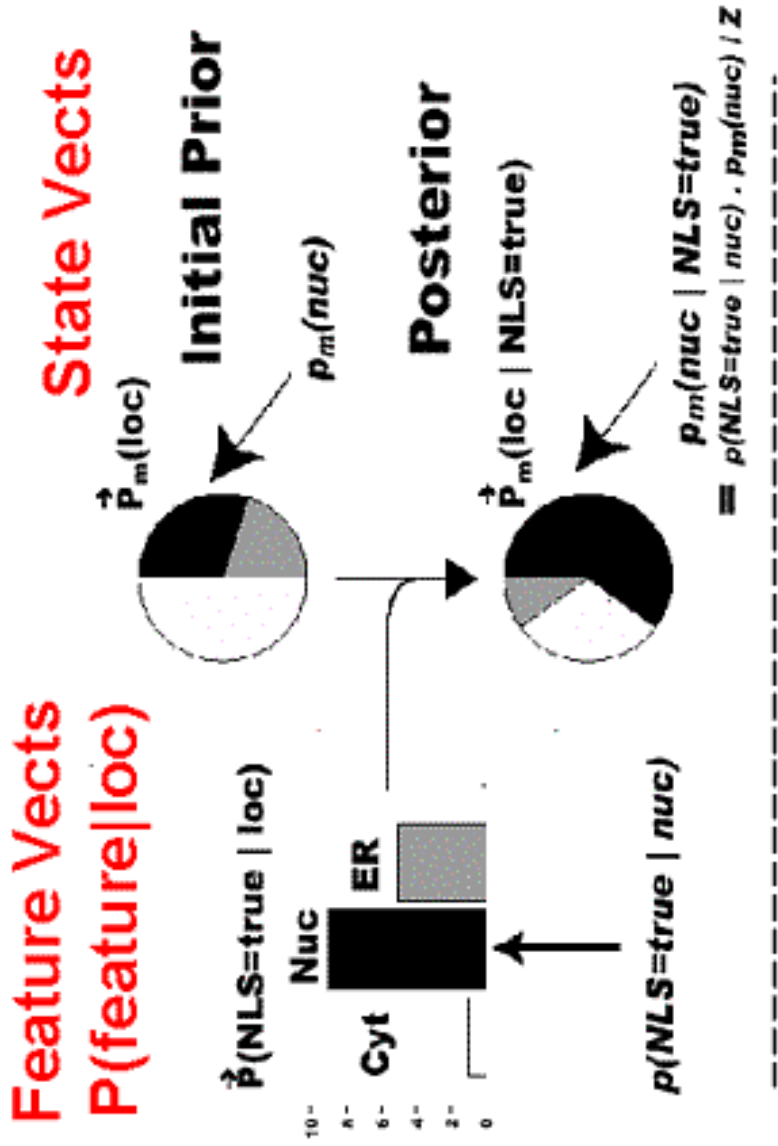


Subcellular localization

Bayesian System for Localizing Proteins



~6000 yeast genes
with expression levels but only
~2000 with known localization



Bayesian
System for
Localizing
Proteins

loc=



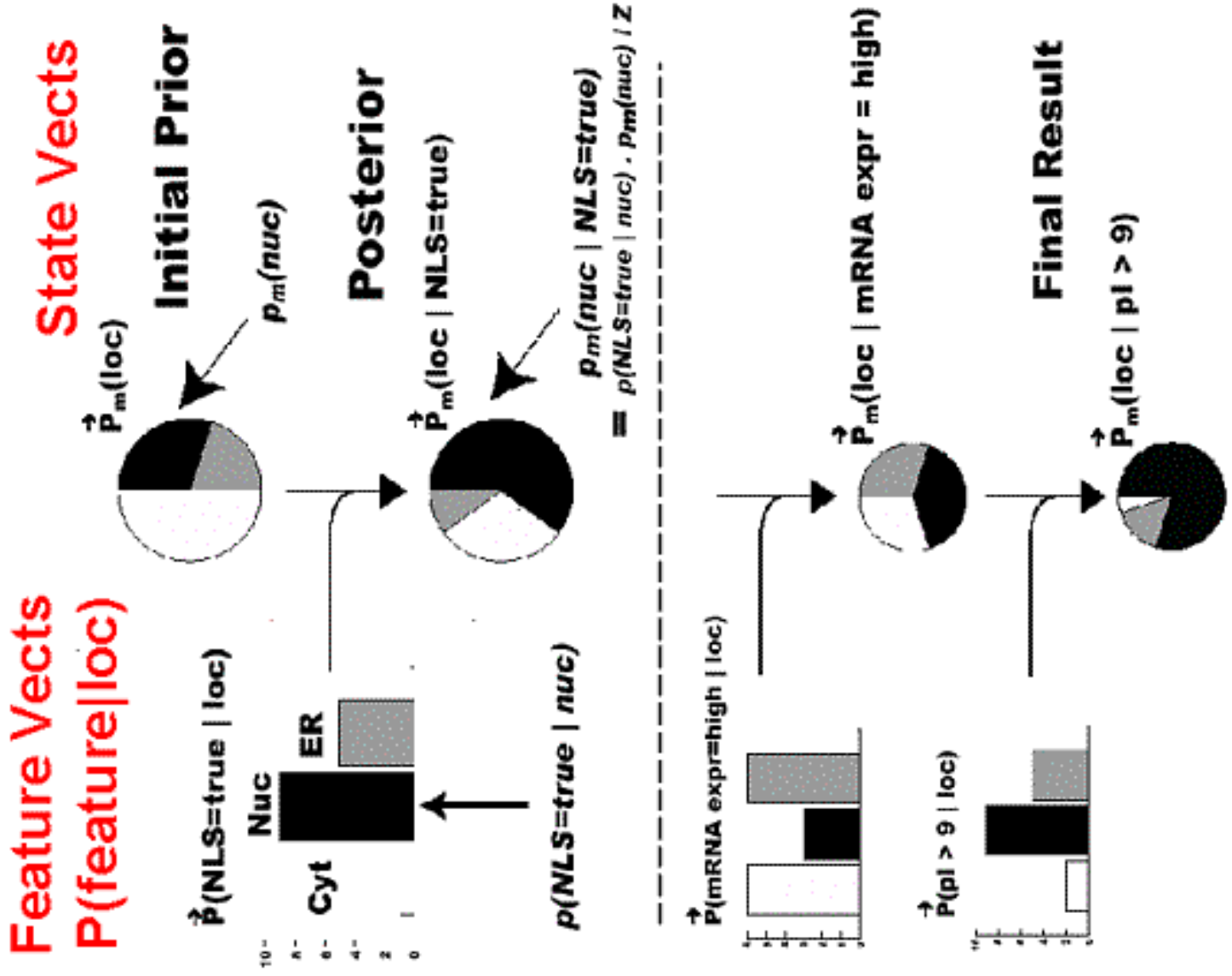
18 Features: Expression Level
 (abs. & fluct.), signal seq.,
 KDEL, NLS, Essential?,
 composition

Bayesian System for Localizing Proteins

loc=

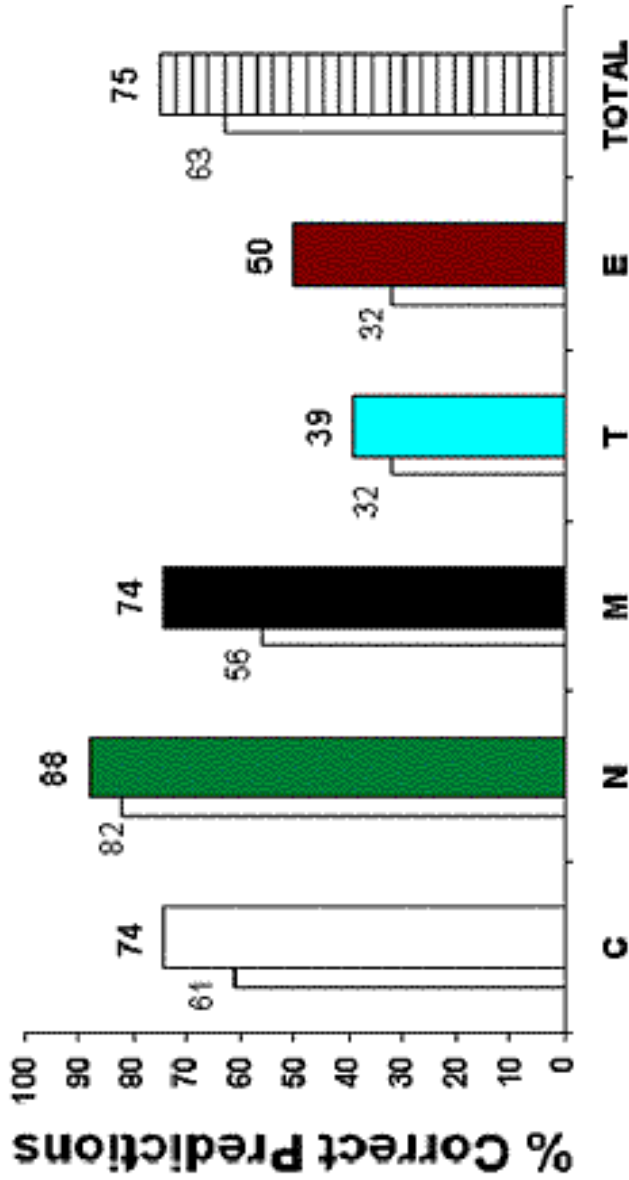
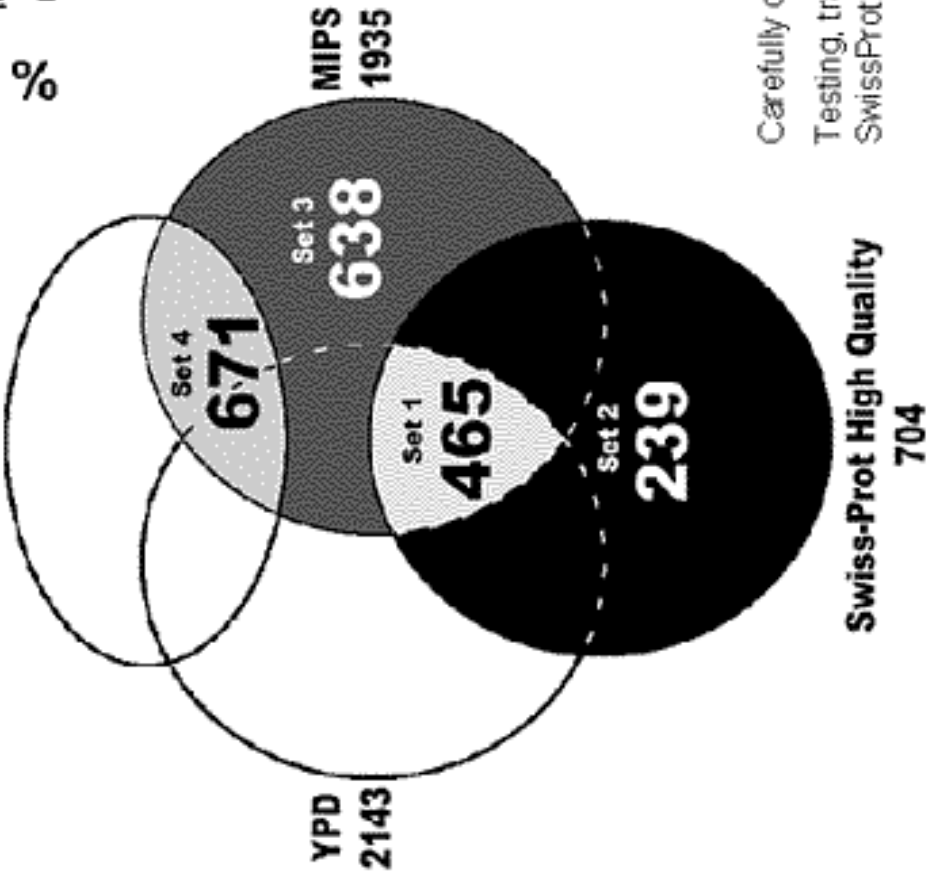


Represent localization of each protein by the state vector $\mathbf{P}(\text{loc})$ and each feature by the feature vector $\mathbf{P}(\text{feature}|\text{loc})$. Use Bayes rule to update.



Results on Testing Data

Localization Annotated as Predicted



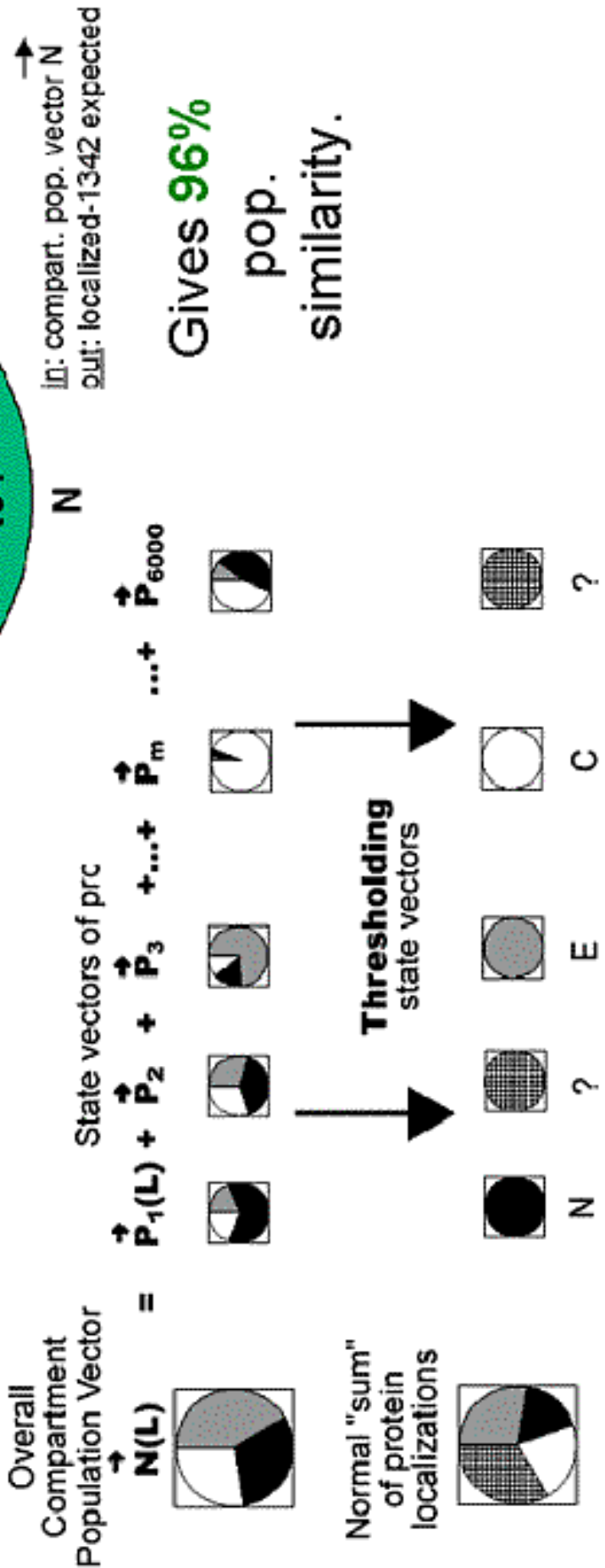
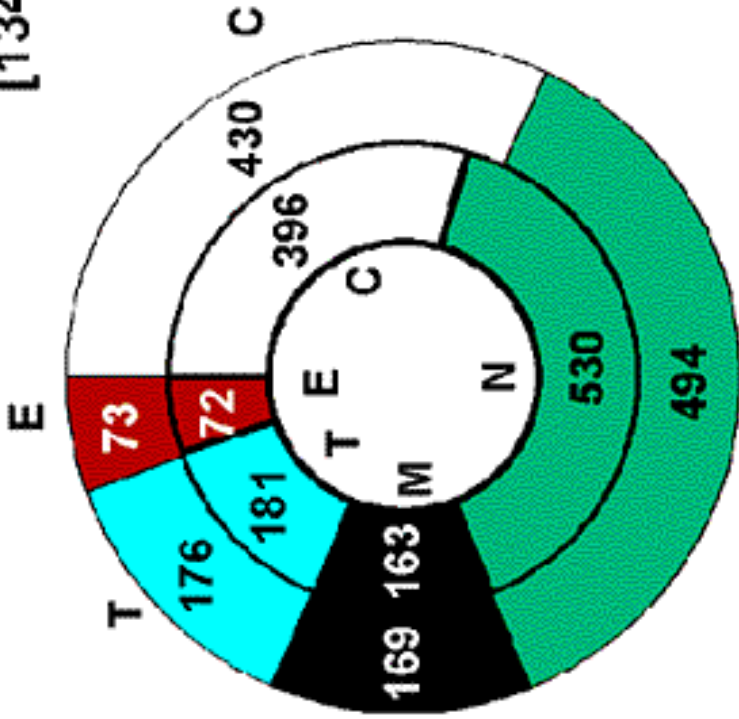
Individual proteins: 75% with cross-validation

Carefully clean training dataset to **avoid circular logic**
 Testing, training data, Priors: ~2000 proteins from YPD, MIPS, SwissProt, Snyder Lab

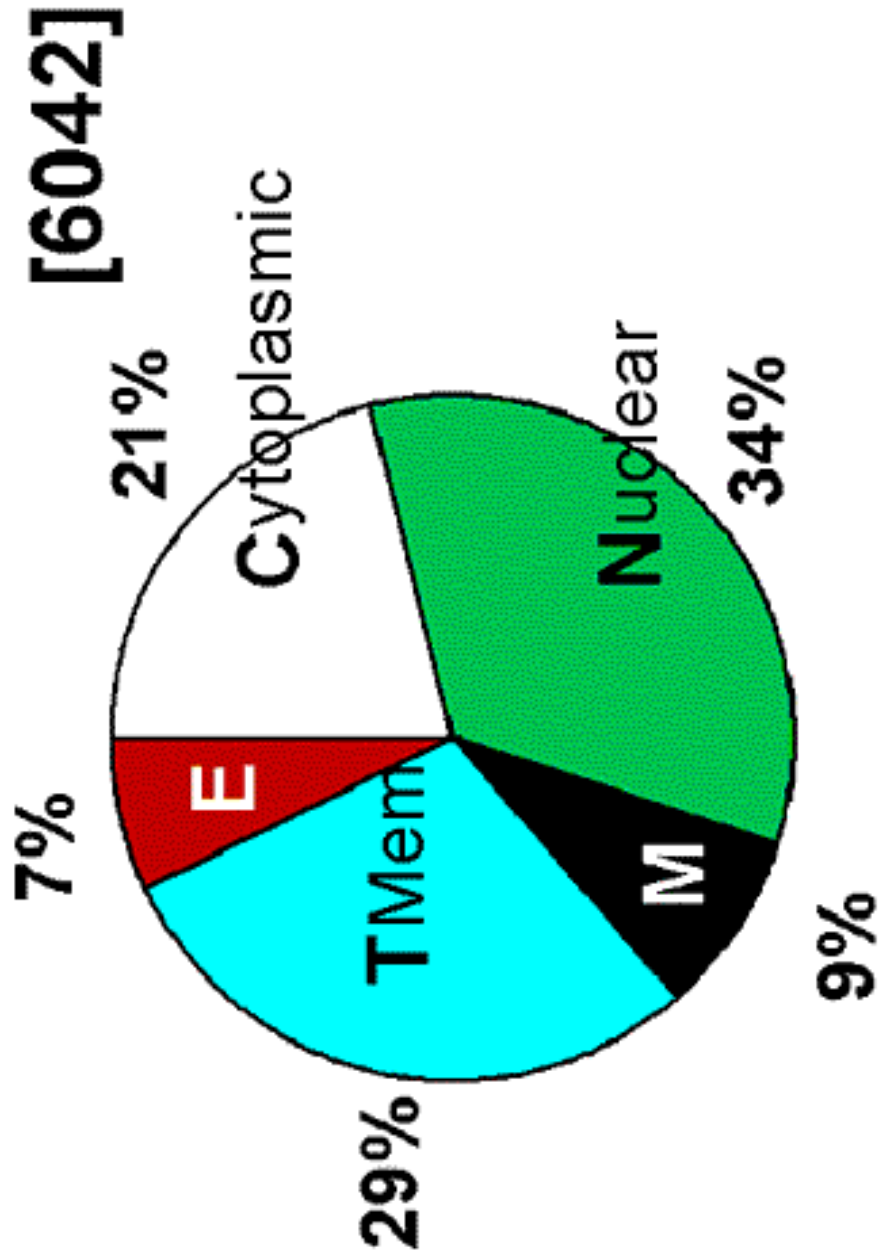
Results on Testing Data #2

Compartment Populations. Like QM, directly sum state vectors to get population.

[1342]



Extrapolation to Compartment
Populations of Whole Yeast Genome:
~4000 predicted + ~2000 known



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

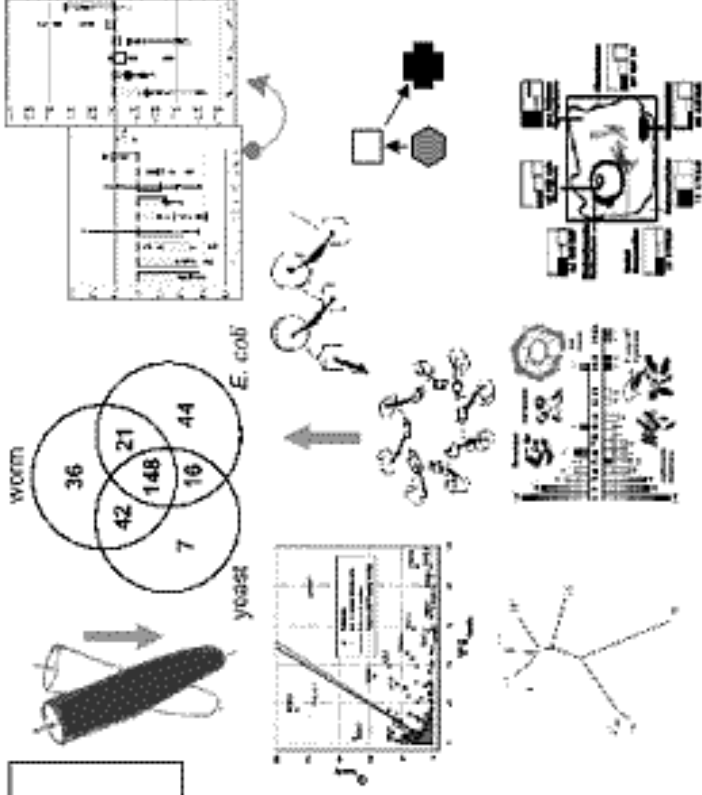
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

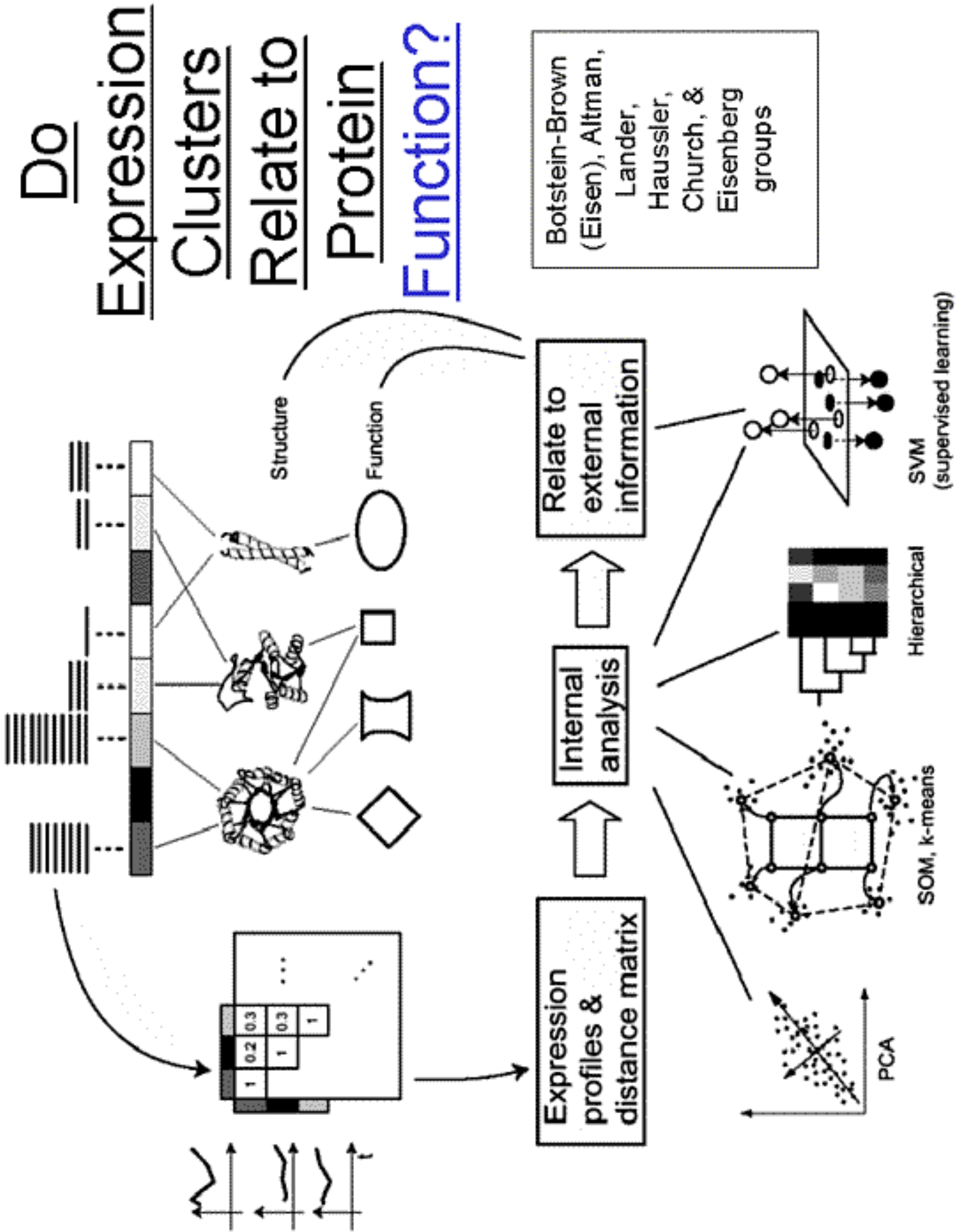
Interactions: Permanent cplx. vs other types

Integrated Views based on Parts

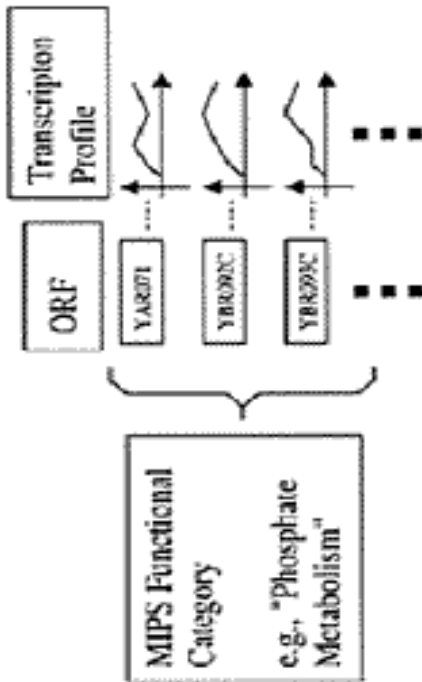


*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu



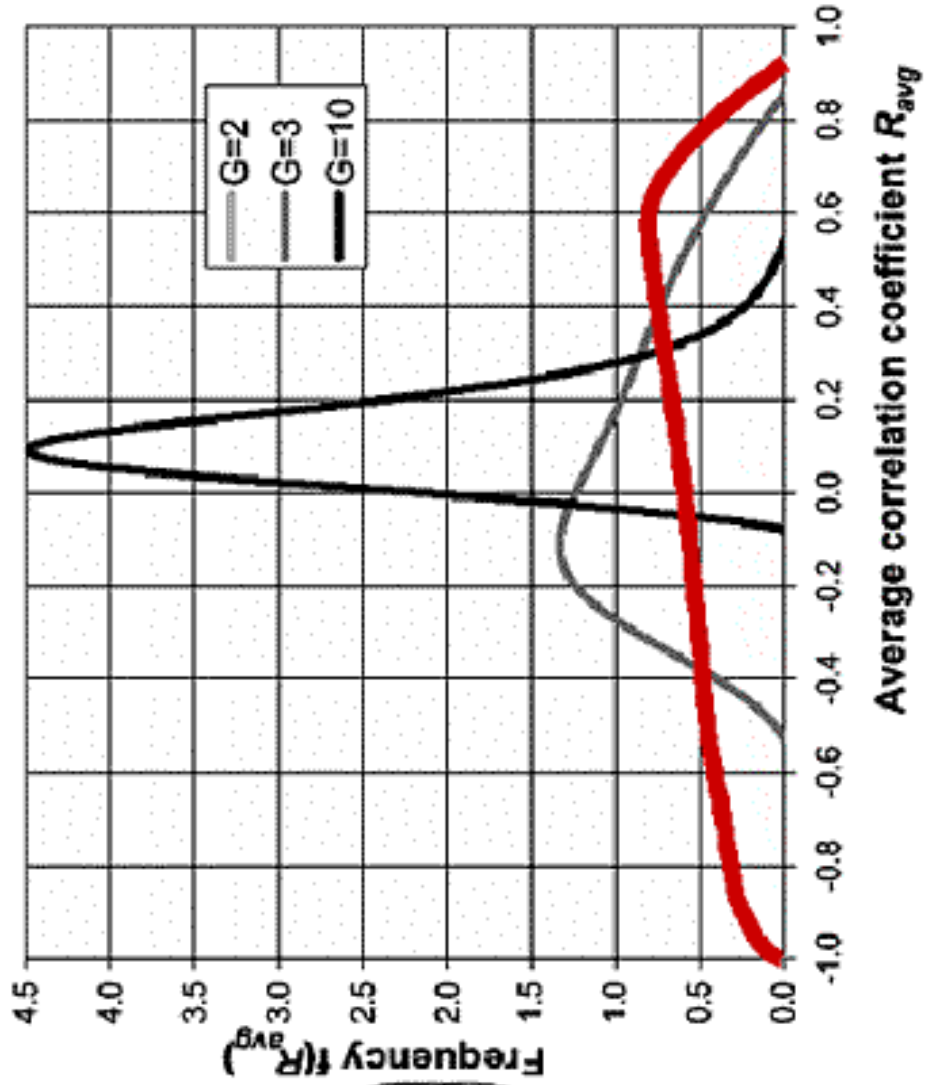
Distributions of Gene Expression Correlations, for All Possible Groupings



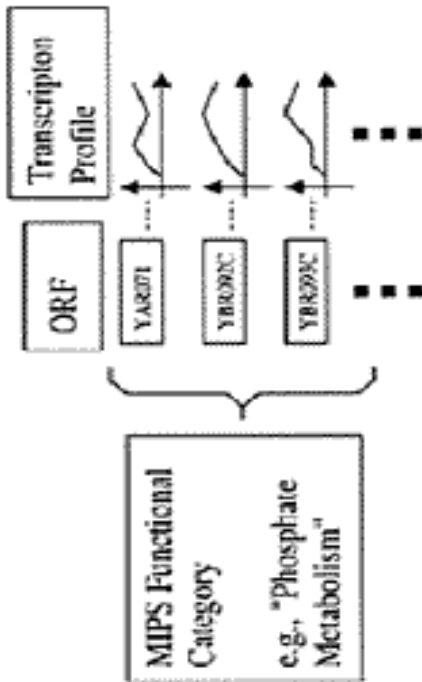
YAR01	1	0.2	0.3	...
YBR09C	0.2	1	0.4	...
YBR09C	0.3	0.4	1	...
...

Correlation Coefficient Matrix (Pearson Coefficient)

Average Correlation Coefficient for Group of Genes



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



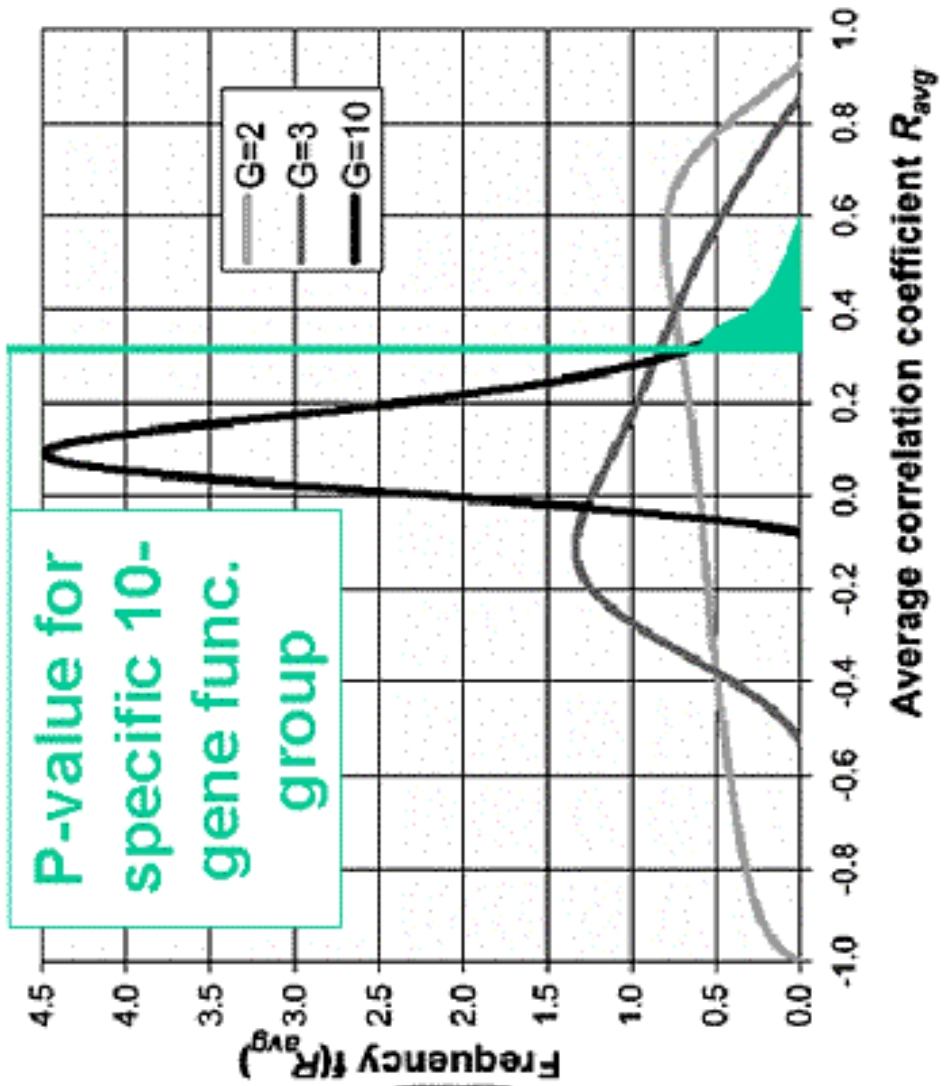
	YAR01	YBR092C	YBR092C	...
YAR01	1	0.2	0.3	...
YBR092C	0.2	1	0.4	...
YBR092C	0.3	0.4	1	...
...

Correlation Coefficient Matrix (Pearson Coefficient)

Average Correlation Coefficient for Group of Genes

Sample for Diauxic shift Expt. (Brown).

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



Correlation:

Always Significant
Sometimes Significant (depends on expt.)
Never Significant

MIPS category	Experiment			
	Cell Cycle (CDC28)	Cell cycle (CDC15)	Diauxic shift	Spo- rulation
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

MIPS category	Experiment			
	Cell Cycle (CDC28)	Cell cycle (CDC15)	Diauxic shift	Spo- rulation
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
Glucosogenesis	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

Based on Distributions,

Correlation of

Established Functional

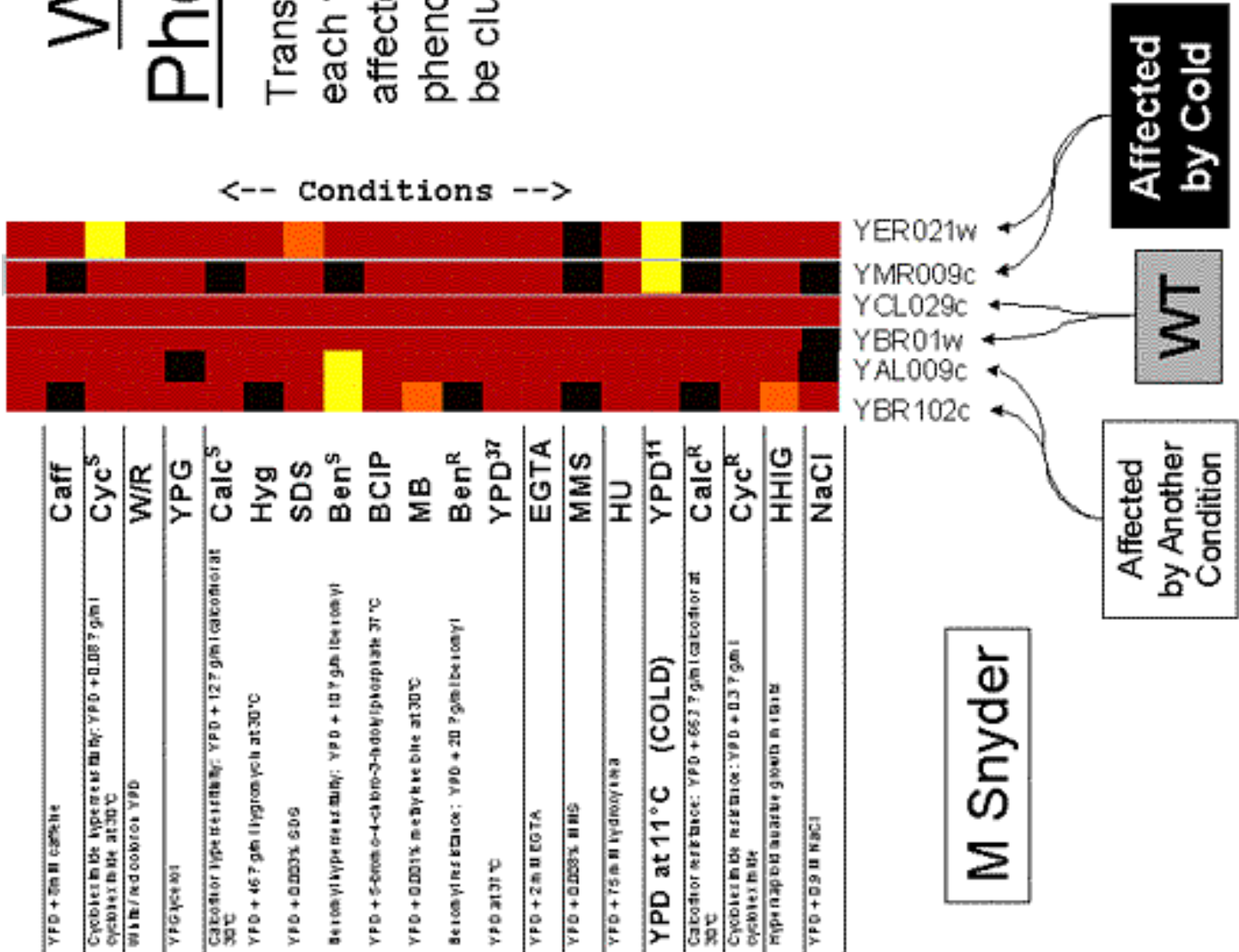
Categories, Computer

Clusterings

	Fraction of significant groups		Total # groups
	CDC28	Diauxic Sporulation	
MIPS 1	63%	19%	16
MIPS 2	50%	17%	102
MIPS 3	23%	5%	73
"Energy" (2 nd level)	40%	20%	10
SOM	93%	-	30
Clustering	-	80%	25

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 clustered similarly to expression



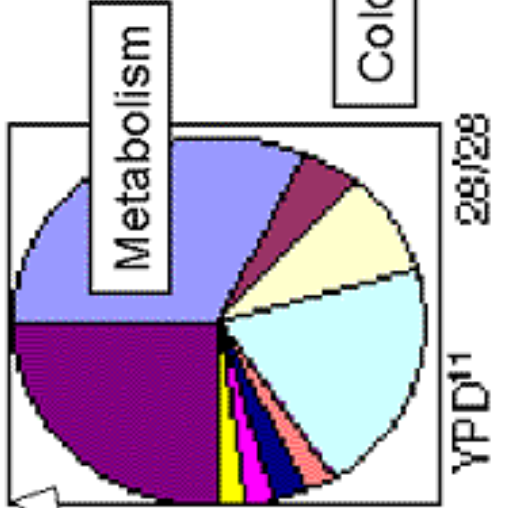
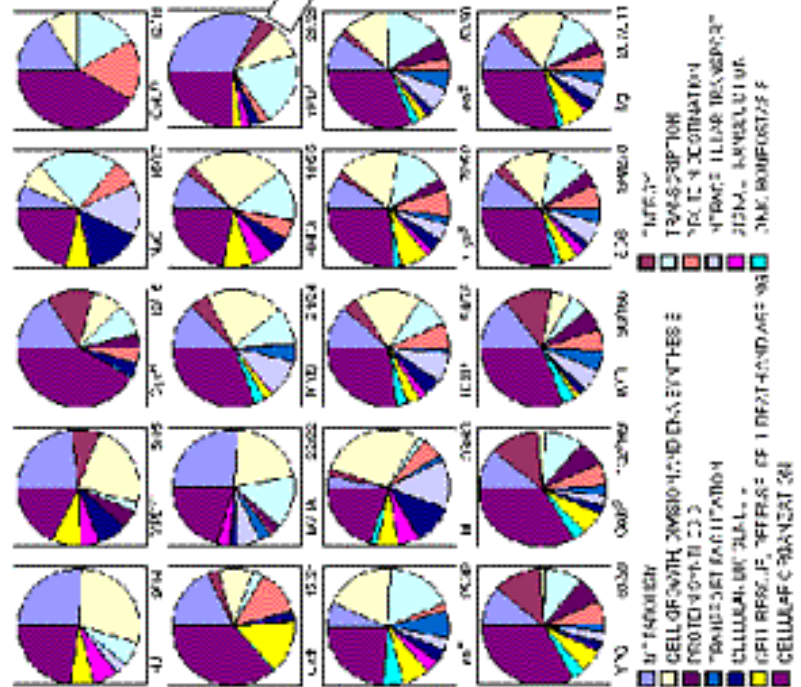
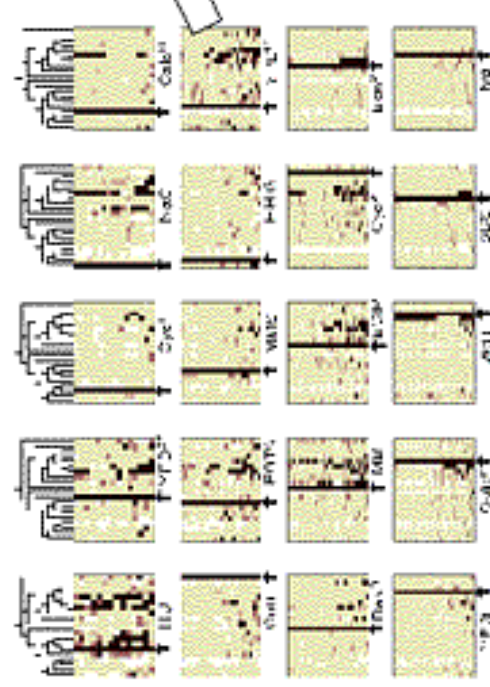
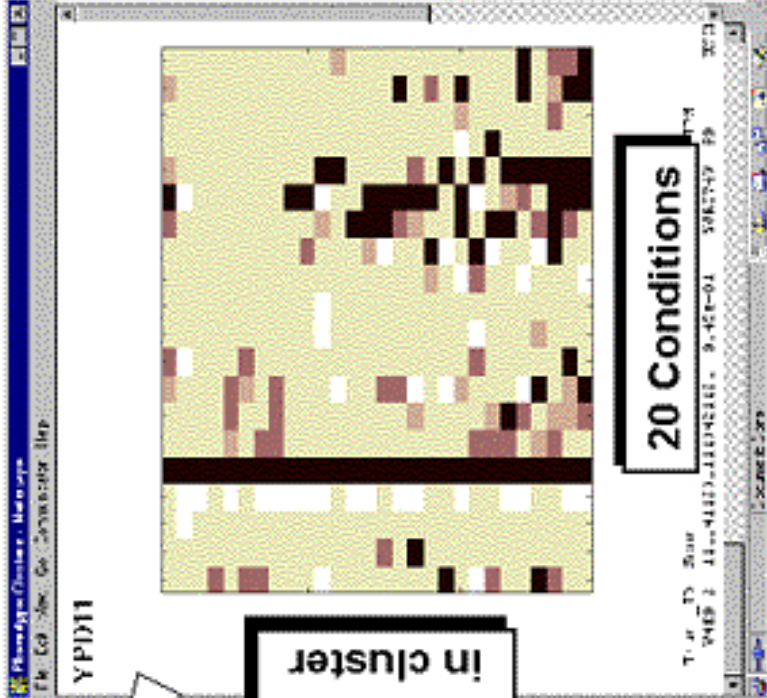
Phenotype ORF Clusters from Transposon Expt.

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....



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

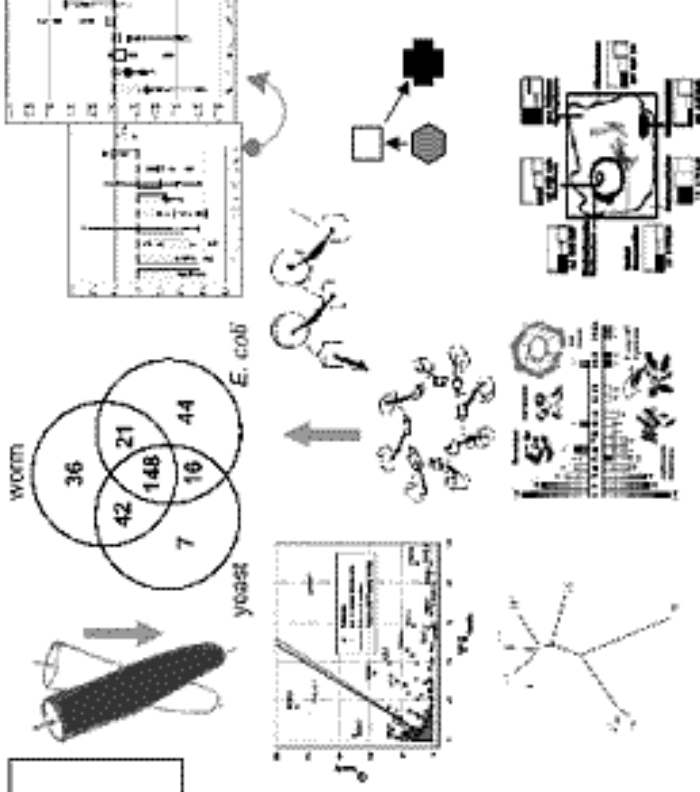
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

Can FUNCTION be defined well enough to relate to expression?

Problems defining function:

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

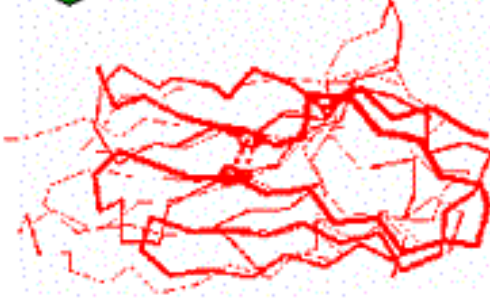
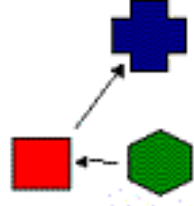
Conflating of Roles: molecular action, cellular role, phenotypic manifestation.

Non-systematic Terminology:

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

More clearly defined attributes of proteins

Fold, Localization, Interactions & Regulation



VS.

Functional Classification

COGS
(cross-org, just conserved, NCBI Koonin/Lipman)

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

GenProtEC
(E. coli, Riley)

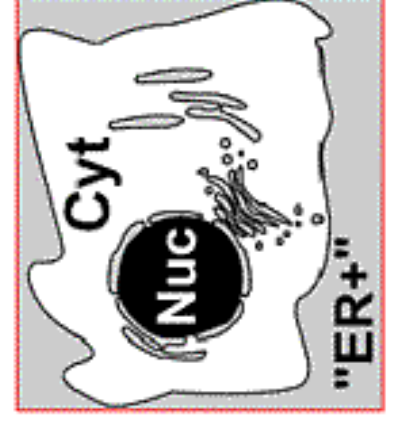
MIPS
(yeast, Mewes)

"Fly"
(fly, Ashburner) now extended to GO (cross-org.)

MIPS/PEDANT
(yeast, Mewes)

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

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



Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

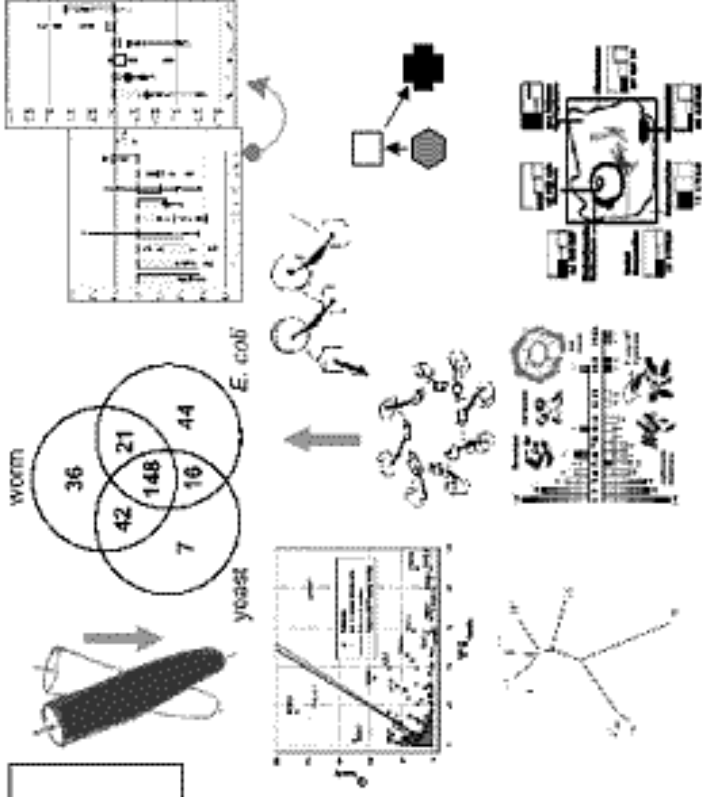
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

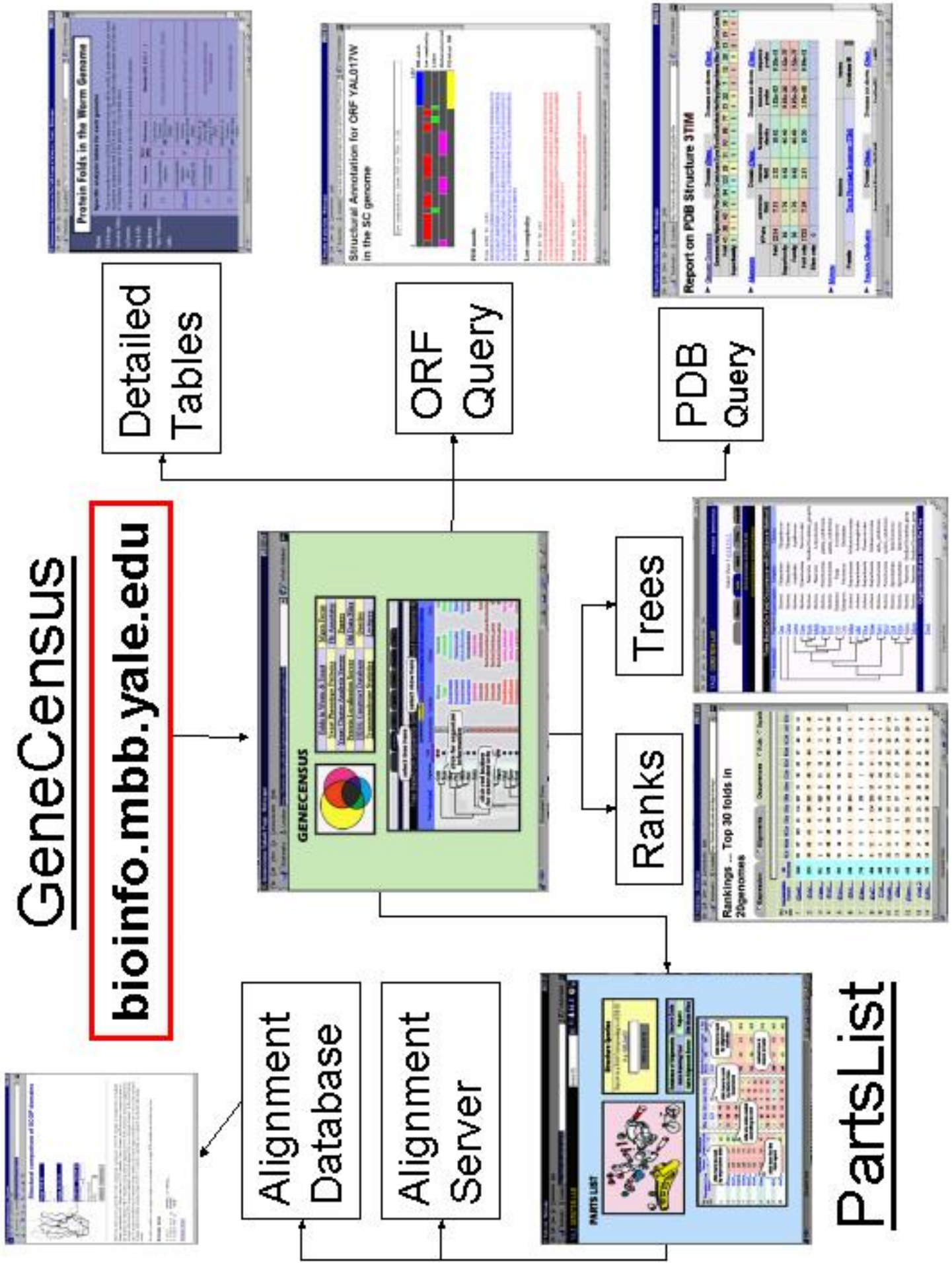
Interactions: Permanent cplx. vs other types

Integrated Views based on Parts



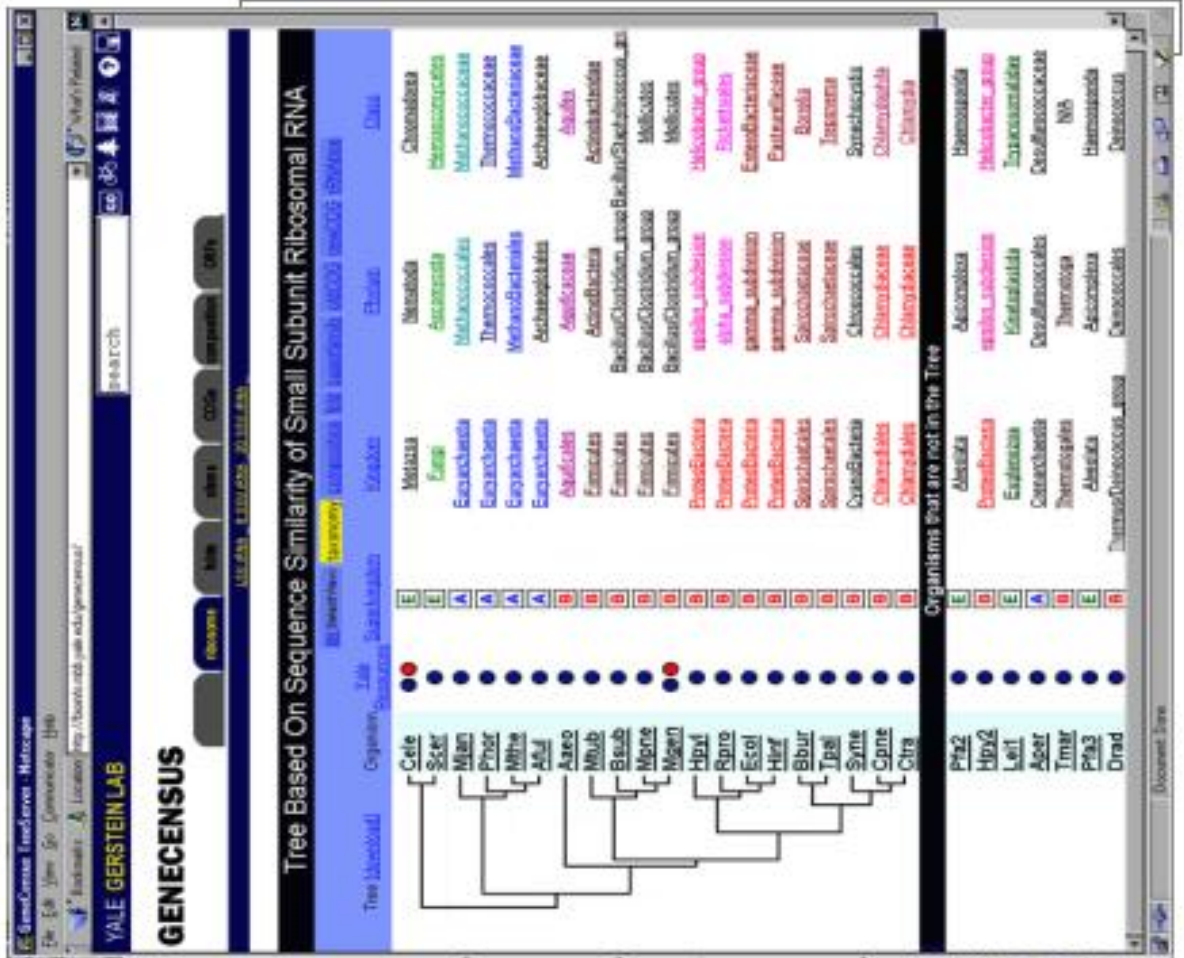
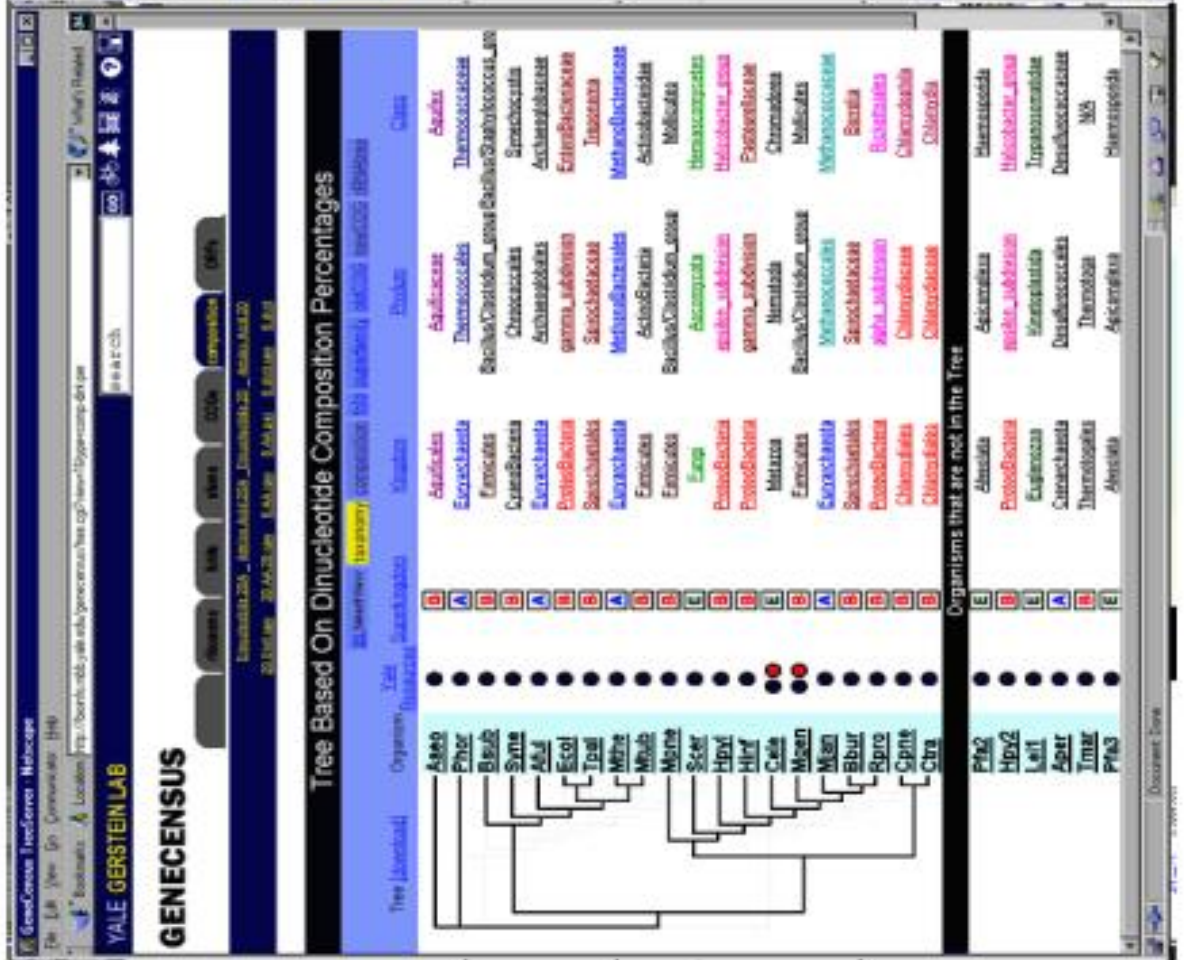
*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu

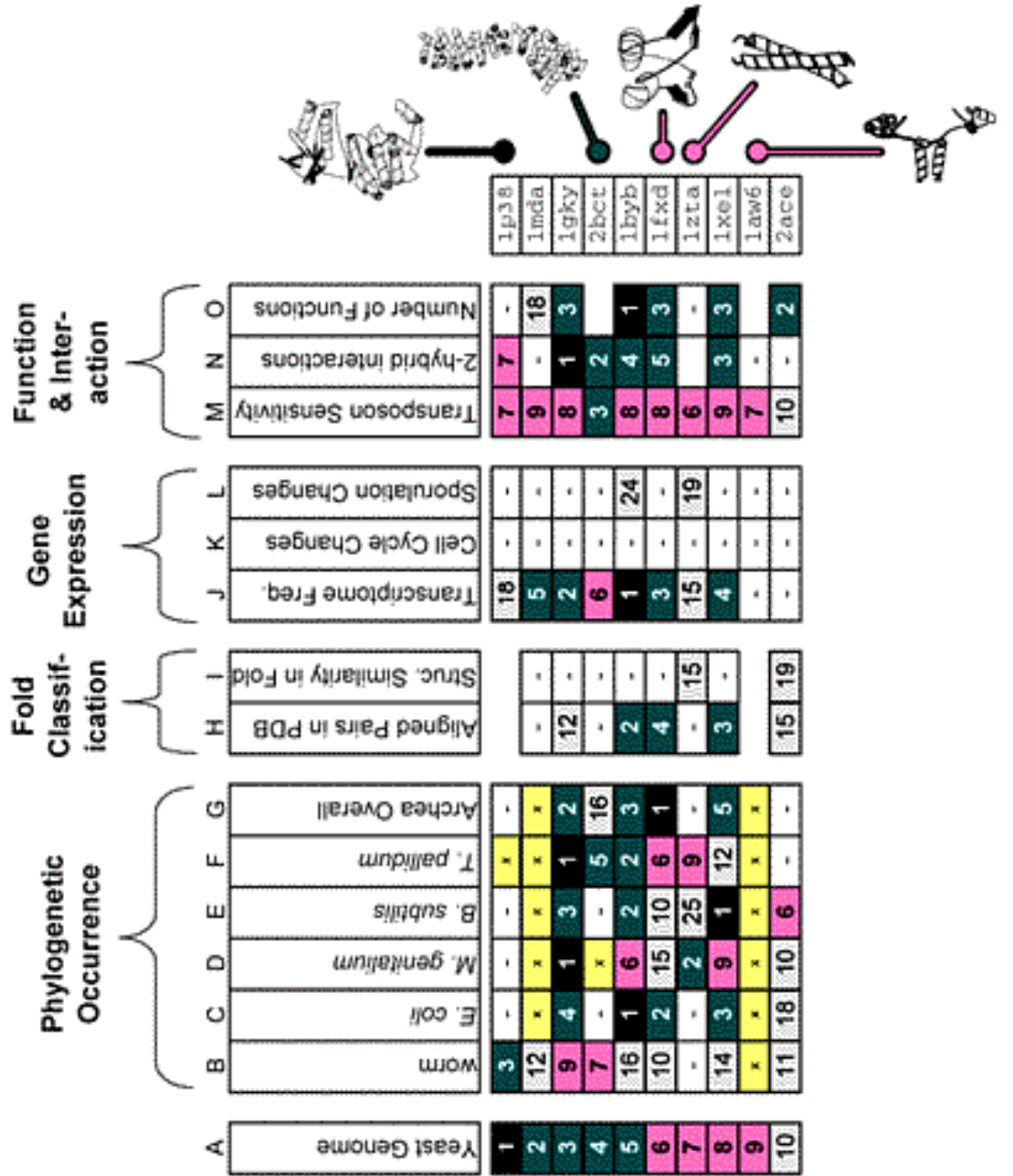


GeneCensus Dynamic Tree Viewers

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



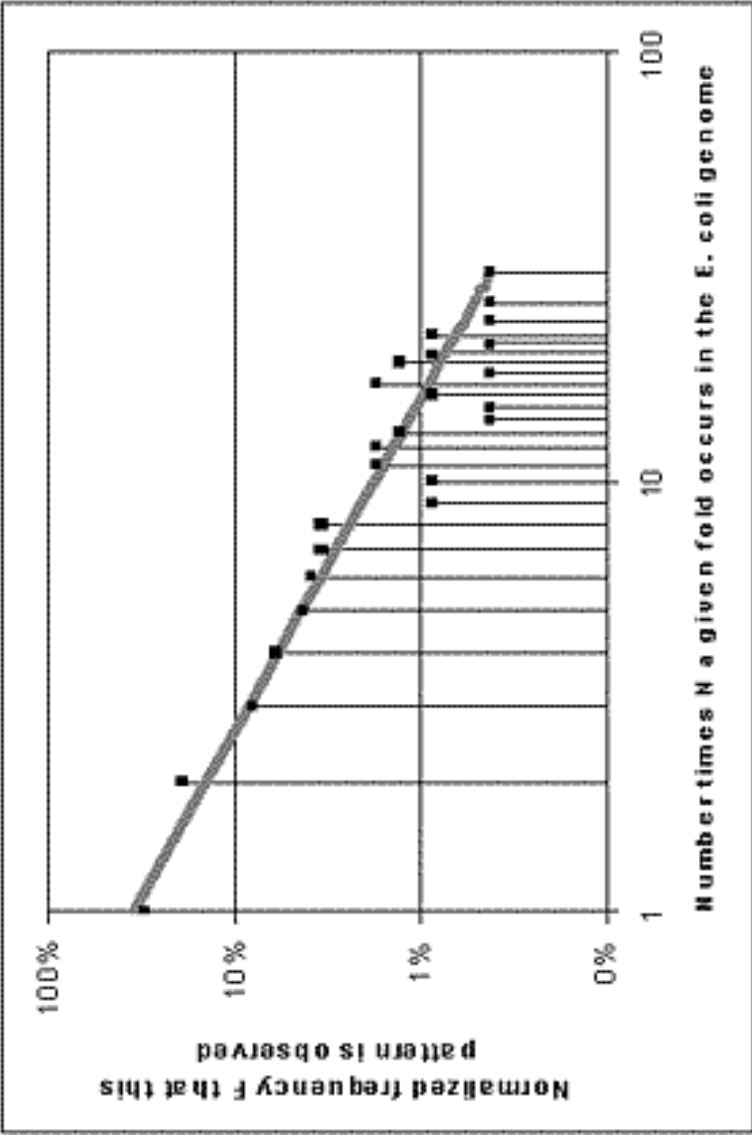
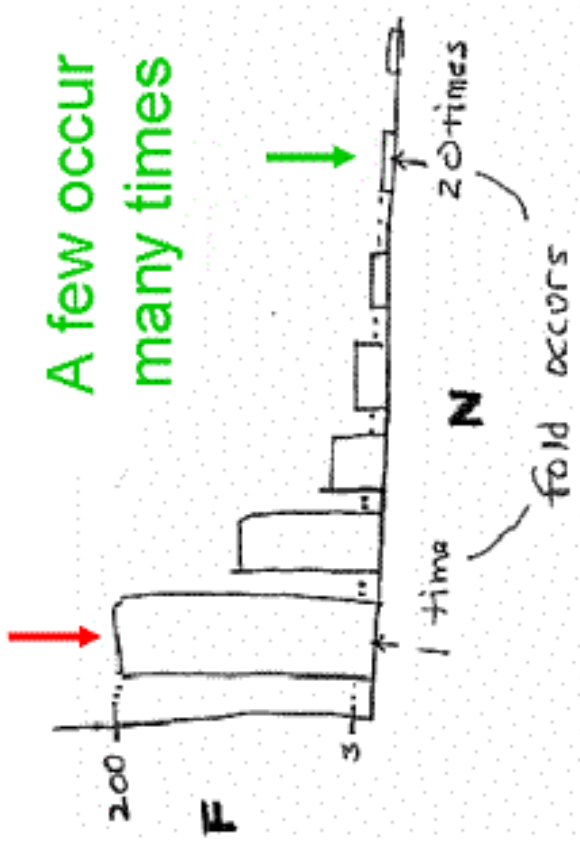
Surveying a Finite PartsList from Many Perspectives



Insert Animations

Part
Occurrences
Frequencies
Follow a
Power Law

Most occur
 just once

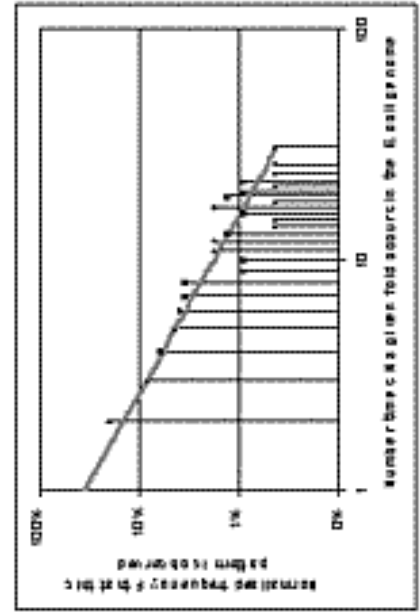
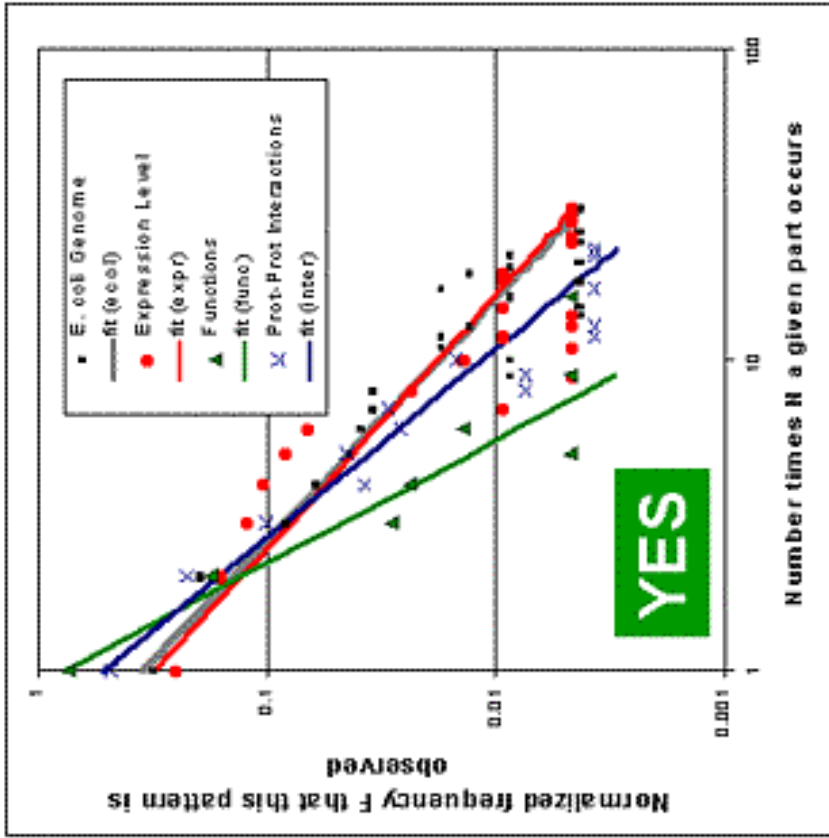
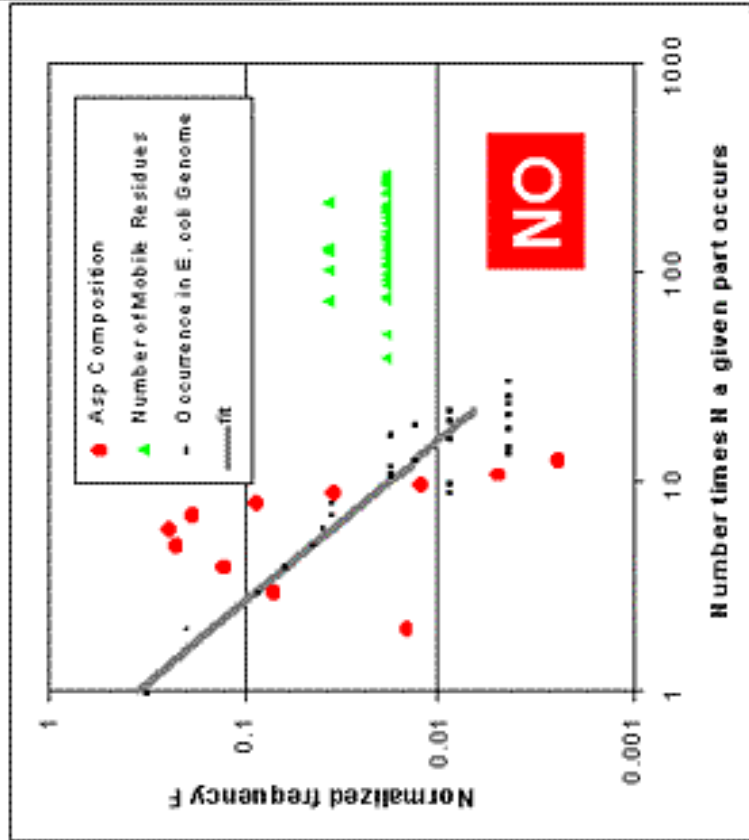


Overall falloff follows
 a Power Law
 (Zipf law, "80-20" rule)

$$F = AN^{-B}$$

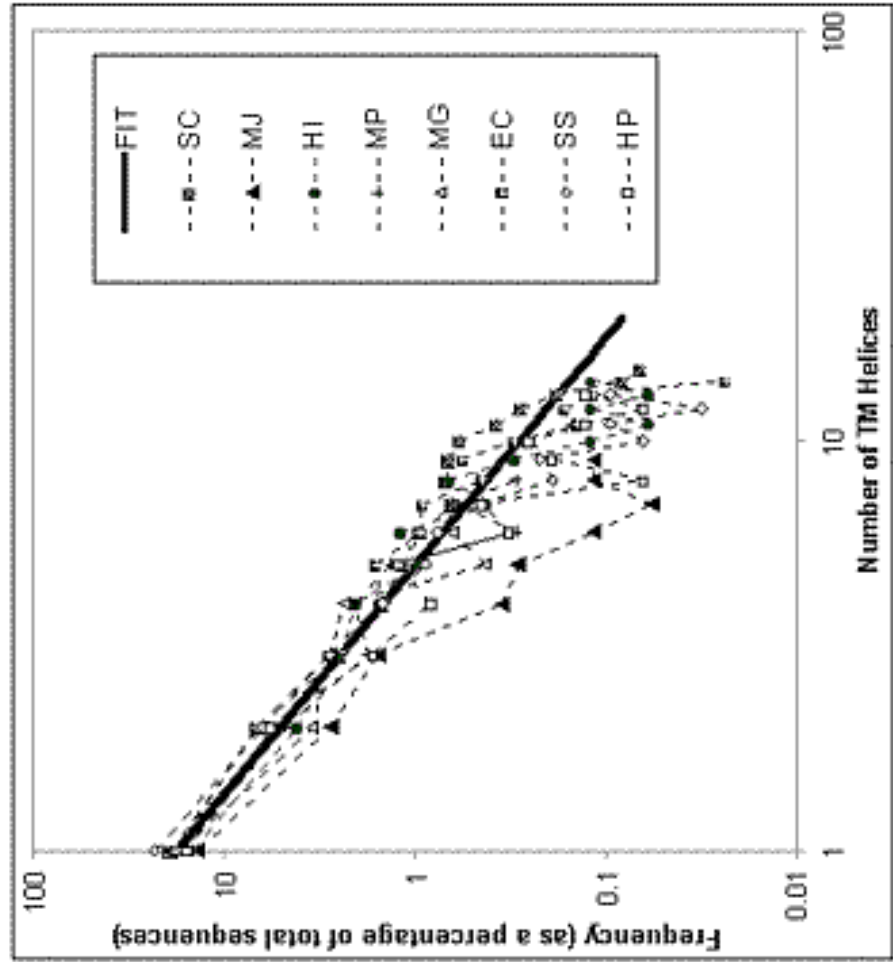
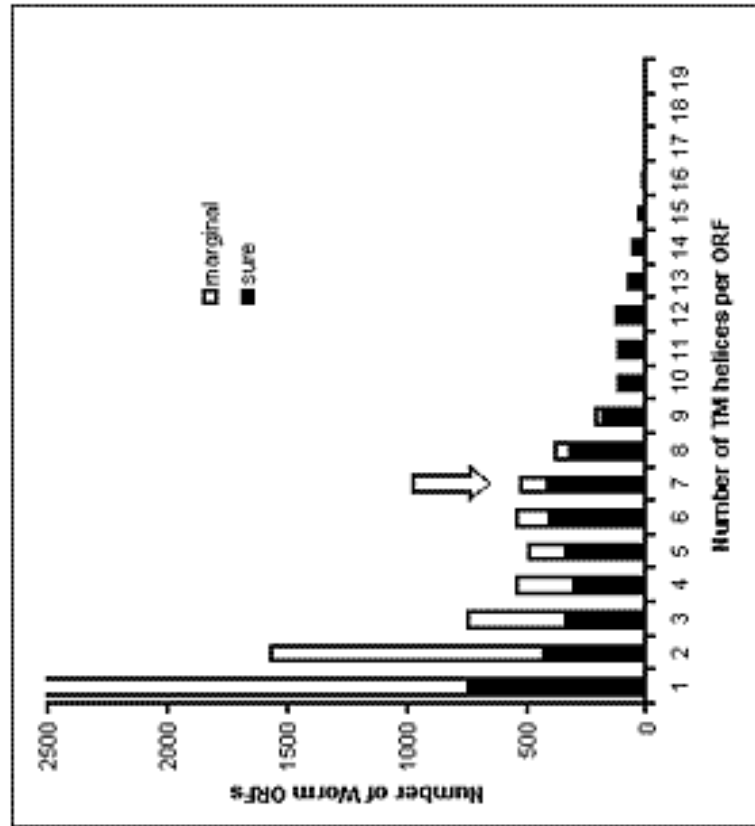
(A=.35, B=1.2 above)

Many (but not all) Part Quantities Follow Power Laws



Occurrence of TM-helices follow Zipf Law

- TM identification (KD, GES).
- Count number with 2 peaks, 3 peaks, &c.
- Similar conclusions to others: von Heijne, Rost, Jones, &c.
- Divide into sure and marginal
- (Boyd & Beckwith's criteria)
- $F = 1/[5N^2]$



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

Integrative Genomics: Surveys of a Finite Parts List

Using Parts to Interpret Genomes

Shared & Common parts: Venn Diag.

Whole-genome trees, top-10 with $\beta\alpha\beta$.

Ψ -genes

Folds/func? A few versatile scaffolds (TIM).

Using Parts & Categories to Mine Expression Data

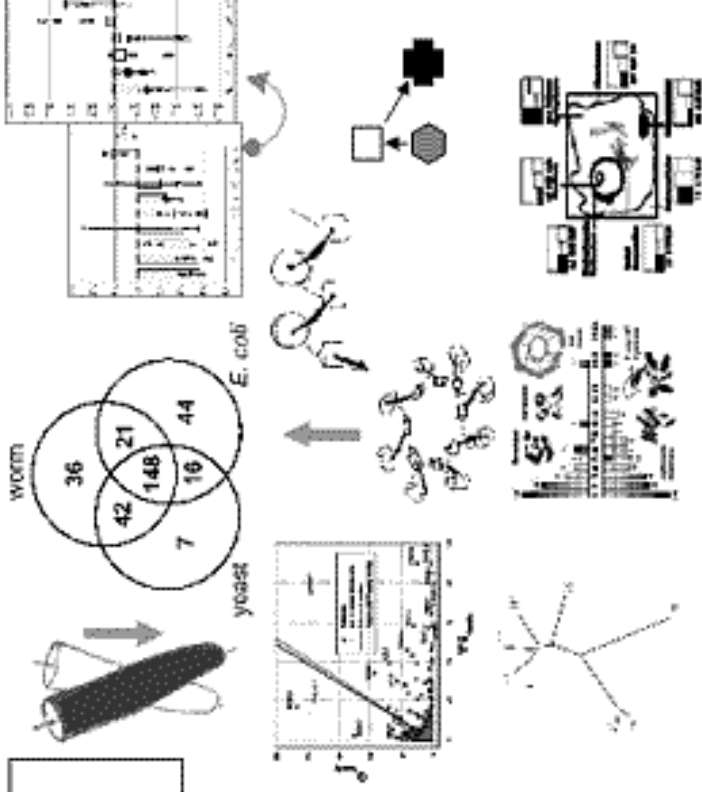
Folds: Top-10 in expression (TIM)

Localization: Bayesian framework

Function: Is there a relation?

Interactions: Permanent cplx. vs other types

Integrated Views based on Parts



*H Hegyi, J Lin, N Echols,
P Harrison, M Levitt, C Wilson,
R Das, A Drawid, R Jansen,
D Greenbaum, M Snyder,
S Teichmann, P Bertone,
B Stenger, J Tsai, C Wilson,
V Alexandrov, J Qian,
W Krebs, M Snyder*

bioinfo.mbb.yale.edu