### Using a Finite Parts List to Understand Complexity:

### Analysis of Whole-Genome Expression Patterns in terms of Protein Folds

Mark Gerstein



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

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]



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

## Using a Finite Parts List to Understand Complexity: Analysis of Whole-Genome Expression Patterns in terms of Protein Folds

### 1 Past: Comparing Genomes in terms of Protein Folds

Fold Library, Shared and Unique Folds (Venn, fold tree), Common Folds (top-10), Repeated  $\beta\alpha\beta$ , Relation of Structural Class to Functional Class, Bias Problem, Prediction (esp. TMs)

### 2 Preliminary:

## Integrating Expression Data into the Analysis

Brown lab yeast data, Top fold by expression, 12-TMs exp. Change, Functional Class Correlated with Expression Change





## Fold Library vs. Other Fundamental Data structures

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



(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.









### **Shared Folds in Initial Genomes**



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











# $\frac{\text{Bias Problem}}{\text{Prediction}}$

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









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



## Using a Finite Parts List to Understand Complexity: Analysis of Whole-Genome Expression Patterns in terms of Protein Folds

### 1 Past: Comparing Genomes in terms of Protein Folds

Fold Library, Shared and Unique Folds (Venn, fold tree), Common Folds (top-10), Repeated  $\beta\alpha\beta$ , Relation of Structural Class to Functional Class, Bias Problem, Prediction (esp. TMs)

### 2 Preliminary:

## Integrating Expression Data into the Analysis

Brown lab yeast data, Top fold by expression, 12-TMs exp. Change, Functional Class Correlated with Expression Change







## Integrate Gene Expression Data into Folds in Genome Analysis



Yeast Expression Data Principally from Brown lab site (1st available!).

Also: SAGE data, Church lab data, Snyder lab transposon data, Young lab data

> Data from: "Exploring the Metabolic and Genetic Control of Gene Expression on a Genomic Scale". DeRisi et al. Science 278:680. Also, Chu *et al. Science* 1998 October 23; 282: 699-705

<u>Caveat:</u> With GeneChips, only rough ORFto-ORF comparison possible, 1-ORF comparisons along timecourse better. SAGE and transposon data better for this « regard. Mostly we just aggregate the data. <u>Top-10 Folds</u> <u>according to</u> <u>Expression</u>



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

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

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

## Expression not related to Overall Fold Class or Overall Composition

- Fold class composition weighted by transcript frequency does **not** change during differential expression of genes.
- Amino acid composition weighted by transcript frequency does **not** change during differential expression of genes.





## <u>Different Classes of</u> <u>Membrane Proteins</u> <u>Have Different</u> <u>Changes in Expression</u> Level (esp. 12 TMs)



Column gives the expression in aerobic conditions (high sugar, second time-series data point in DeRisi et al.), and other column, in anaerobic conditions (low sugar, high ethanol, last time-series data ( point in DeRisi et al.). 9 hexose permeases, 1 lactate transporter.

Most Expressed TMs

Most Expressed TMs

n aerobic condi	tions	
ORF	TMs	
YHR078W	4	
YGL008C	6	
YBR012W-B	2	
YLR340W	2	
YPL131W	2	
YHR099W	2	
YMR205C	2	
YHR216W	2	
YLR432W	2	
	5	

<u>in ana</u>	aerodic cor	naitions
ORF		TMs
YPR1	149W	4
YDR3	343C	9
YDR3	342C	9
YKL2	217W	7
YHRO	)96C	9
YBR1	116C	2
YIL08	38C	6
YBRO	)12W-B	2
YBR	)54W	7
YBR2	218C	2





Time

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

MIPS YFC: 66 bottom classes, 10 top classes Average correlation of uncharacterized genes is 0.16 Similar to Botstein analysis.





**Correlation Coefficient Matrix (Pearson Coefficient)** 

Average Correlation Coefficient for Group of Genes

<sup>°</sup> <u>Correlate with</u> Expression Level with Functional

**Category** 

Functional category number	Function	Average correlation	# ORFs
01	METABOLISM	0.1001	1005
01.01	amino-acid metabolism	0.1488	199
01.01.01	amino-acid biosynthesis	0.239	114
01.01.04	regulation of amino-acid metabolism	0.23	32
01.01.07	amino-acid transport	0.1198	23
01.01.10	amino-acid degradation	0.0524	36
01.01.99	other amino-acid metabolism activities	0.2205	4
01.02	nitrogen and sulphur metabolism	0.1869	73
01.02.01	nitrogen and sulphur utilization	0.0726	37
01.02.04	regulation of nitrogen and sulphur utilization	0.3715	28
01.02.07	nitrogen and sulphur transport	0.2829	8
01.03	nucleotide metabolism	0.1708	134
01.03.01	purine-ribonucleotide metabolism	0.3639	42
01.03.04	pyrimidine-ribonucleotide metabolism	0.176	28
01.03.07	deoxyribonucleotide metabolism	0.1095	1:
01.03.10	metabolism of cyclic and unusual nucleotides	0.2848	
01.03.13	regulation of nucleotide metabolism	0.2696	1;
01.03.16	polynucleotide degradation	0.2461	
01.03.19	nucleotide transport	0.1187	1:
01.03.99	other nucleotide-metabolism activities	-0.0328	
01.04	phosphate metabolism	0.1348	3'
01.04.01	phosphate utilization	0.16	1:
01.04.04	regulation of phosphate utilization	P 3599	8
01.04.07	phosphate transport	0.0724	1(
01.05	carbohydrate metabolism	0.0779	409
01.05.01	carbohydrate utilization	0.075	256
01.05.04	regulation of carbohydrate utilization	0.1174	120



## <u>Results from Analysis</u> of Correlation of Functional Class and <u>Expression</u>

#### **Highest Correlations**

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

Functional category number	Function	Average correlation	# ORFs
10.04.11	key kinases	0.9403	2
10.04.13	key phosphatases	0.9283	2
11.11	ageing	0.8634	2
02.22	glyoxylate cycle	0.8136	6
10.02.07	G-proteins	0.8122	3
04.03.99	other tRNA-transcription activities	0.6932	4
09.08	biogenesis of Golgi	0.6647	2
09.19	peroxisomal biogenesis	0.6512	2
08.10	peroxisomal transport	0.646	12
04.01.04	rRNA processing	0.6074	53
01.20	secondary metabolism	0.5921	4
01.20.05	amines metabolism	0.5921	4
10.05.11	key kinases	0.0049	4
90	RETROTRANSPOSONS AND PLASMID PROTEINS	0.5299	7
02.10	tricarboxylic-acid pathway	0.5236	22
04.07	RNA transport	0.5111	27

## Using a Finite Parts List to Understand Complexity: Analysis of Whole-Genome Expression Patterns in terms of Protein Folds

### 1 Past: Comparing Genomes in terms of Protein Folds

Fold Library, Shared and Unique Folds (Venn, fold tree), Common Folds (top-10), Repeated  $\beta\alpha\beta$ , Relation of Structural Class to Functional Class, Bias Problem, Prediction (esp. TMs)

### 2 Preliminary:

## Integrating Expression Data into the Analysis

Brown lab yeast data, Top fold by expression, 12-TMs exp. Change, Functional Class Correlated with Expression Change









### http://bioinfo.mbb.yale.edu/MolMovDB

Server Produces Semi-realistic Minimized Interpolation (as MPEG, VRML, &c) between Any 2 Aligned Conformations, Analyzes the Motion

### scop (Murzin et al.) Acknowledgements: Ronald Jansen **S** Teichmann **M** Levitt H Hegyi **C** Wilson **C** Chothia W Krebs **J** Lin **ONR**

(Bright)

PhRMA NSF