

**Gerstein Bioinformatics Group Member**

File Edit View Go Communicator Help

[MB&B] YALE UNIVERSITY Gerstein Lab

Group Member	Title
<a href="#">Mark Gerstein</a>	Principal Investigator
<a href="#">Suganthi Balasubramanian</a>	Postdoctoral Fellow
<a href="#">Paul Harrison</a>	Postdoctoral Fellow
<a href="#">Hedi Hegyi</a>	Postdoctoral Fellow
<a href="#">Jochen Junker</a>	Postdoctoral Fellow
<a href="#">Yuval Kluger</a>	Sloan Fellow
<a href="#">Ning Lan</a>	Postdoctoral Fellow
<a href="#">Nicholas Luscombe</a>	Postdoctoral Fellow
<a href="#">Jiang Qian</a>	Postdoctoral Fellow
<a href="#">Vadim Alexandrov</a>	Graduate Student
<a href="#">Paul Bertone</a>	Graduate Student
<a href="#">Joint with Michael Snyder</a>	
<a href="#">Rajdeep Das</a>	Graduate Student
<a href="#">Dov Greenbaum</a>	Graduate Student
<a href="#">Ronald Jansen</a>	Graduate Student
<a href="#">Ted C Johnson</a>	Graduate Student
<a href="#">Werner G Krebs</a>	Graduate Student
Nathaniel Echols	
<a href="#">Jimmy Lin</a>	
<a href="#">Patrick McGarvey</a>	
Christopher Titcombe	Undergraduate Student
<a href="#">Joann Delvecchio</a>	Administrator

**Who we are?**

Search  Only this site (best)  Go! Special Search

# MB&B Bioinformatics Group (Gerstein Lab)

**Scientific American: Feature Article: The Bioinformatics Gold Rush: July...**

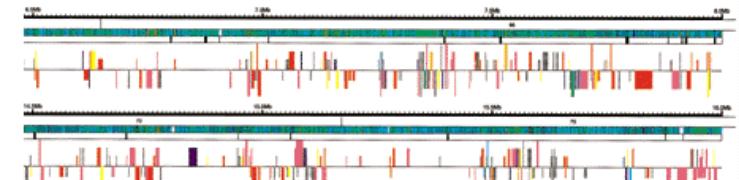
File Edit View Go Communicator Help  
Bookmarks Netsite: 0700issue/0700howard.html What's Related

**SCIENTIFIC AMERICAN**

Main Menu Interview Bookmarks Feedback  
Current Issue Explore! Ask the Experts Marketplace Search the Web

**FEATURE ARTICLES**

## The Bioinformatics Gold Rush



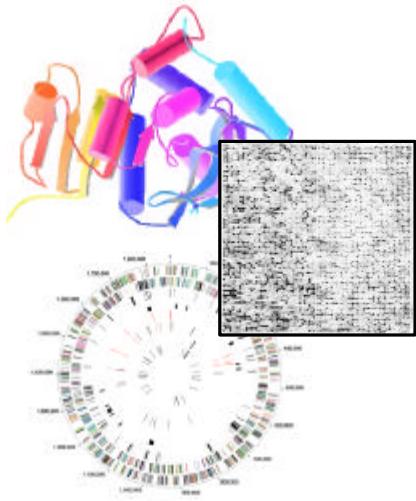
A \$300-million industry has emerged around turning raw genome data into knowledge for making new drugs

*By Ken Howard*

**SUBTOPICS:** Plastics." When a family friend whispered this word to Dustin Hoffman's character in the 1967 film *The*

Document: Done

Seems to becoming  
of interest...



## What do we do?



*(Molecular) **Bio - informatics** =*

Conceptualizing **biology in terms of molecules** (in the physical-chemical sense) and then **applying “informatics”** techniques (derived from disciplines such as CS and statistics) to **organize and analyze** the **information** associated with these molecules, on a **large-scale** and in an **integrative fashion**. Bioinformatics is a practical discipline with many **applications**.

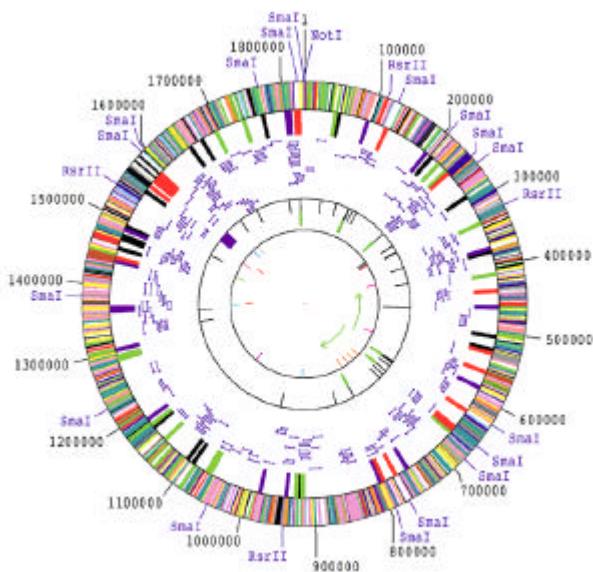
# Large-scale, Integrative Analysis of...

Structures



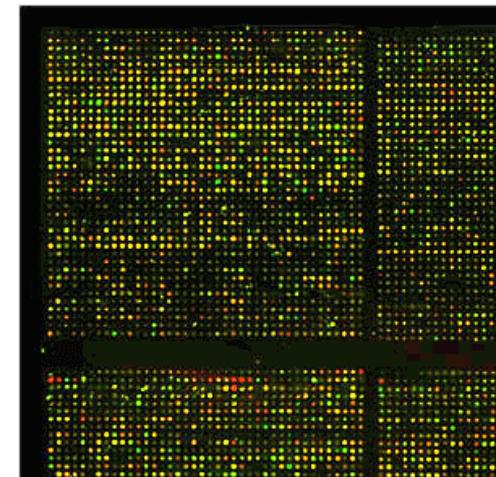
1990

Genomes



2000

Expression  
Data



2010

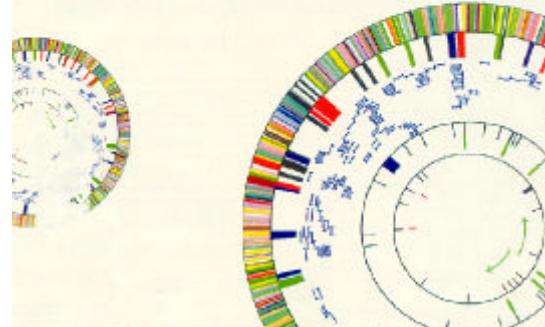
**1995**

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



**1997**

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



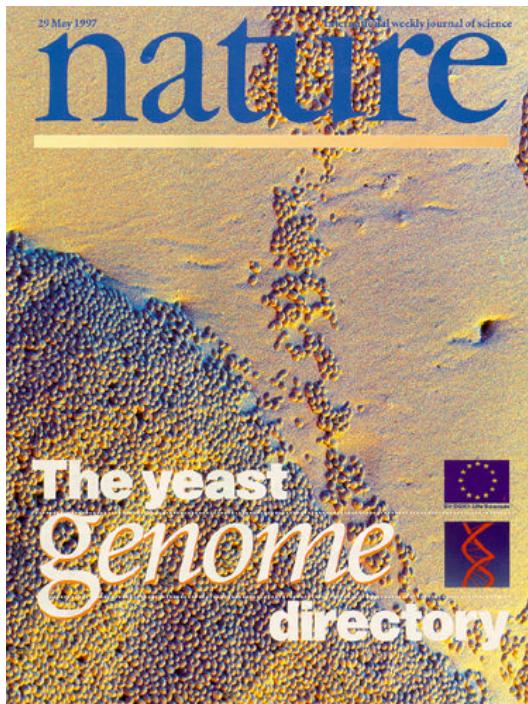
**1998**

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



**2000**

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



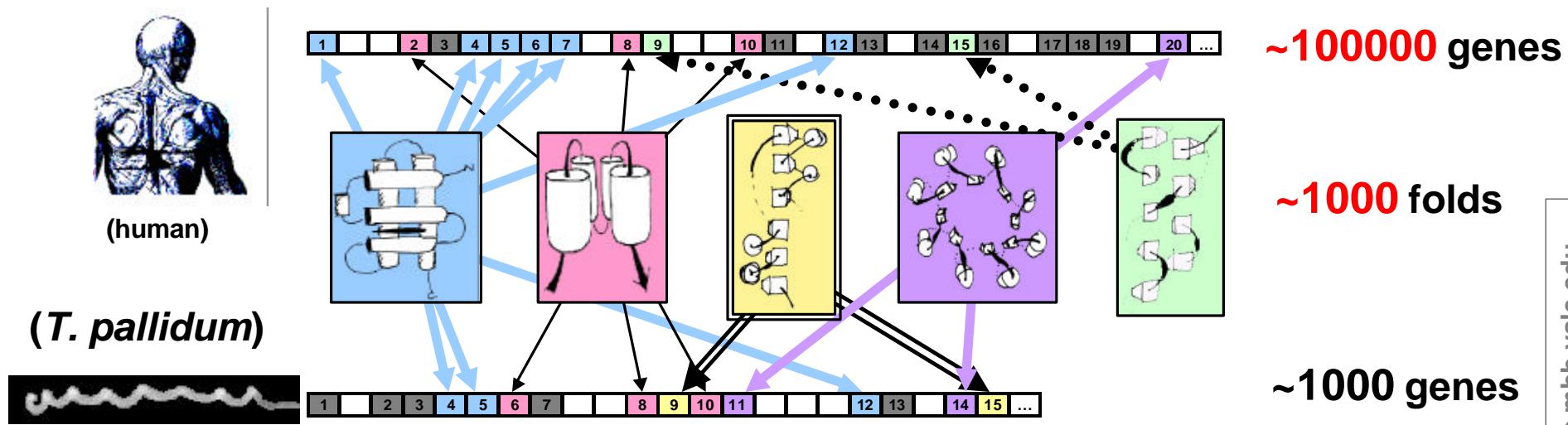
**2000**  
The Year  
of  
Genomics



'98 spoof

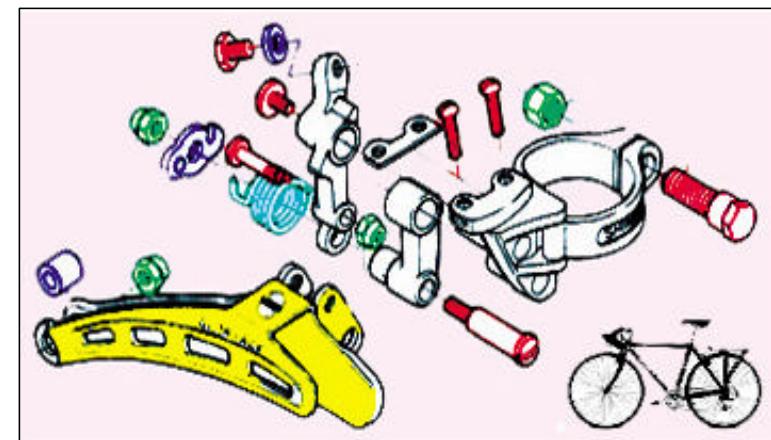
real thing, Apr '00

# The Finite Parts List as a Simplifying Theme

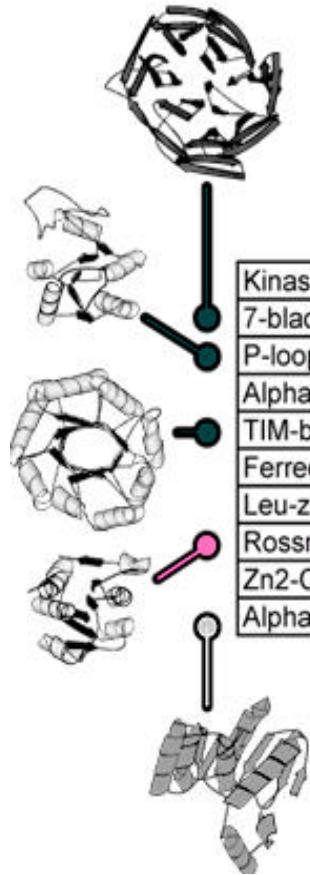


Parts = **FOLDS**, pathways, functions, sequence families, blocks, motifs....

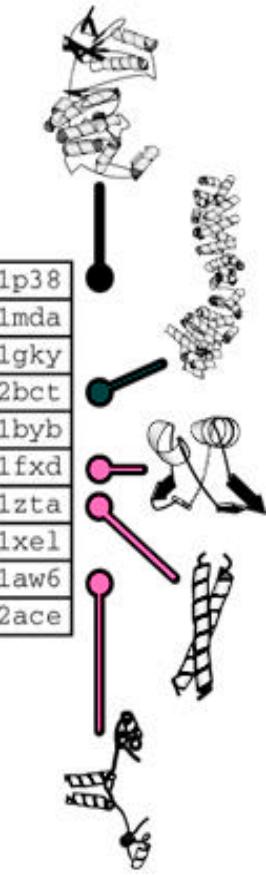
**Global Surveys of a Finite Set of Parts from Many Perspectives**



# Surveying a Finite Parts List from a Infinite Number of Perspectives



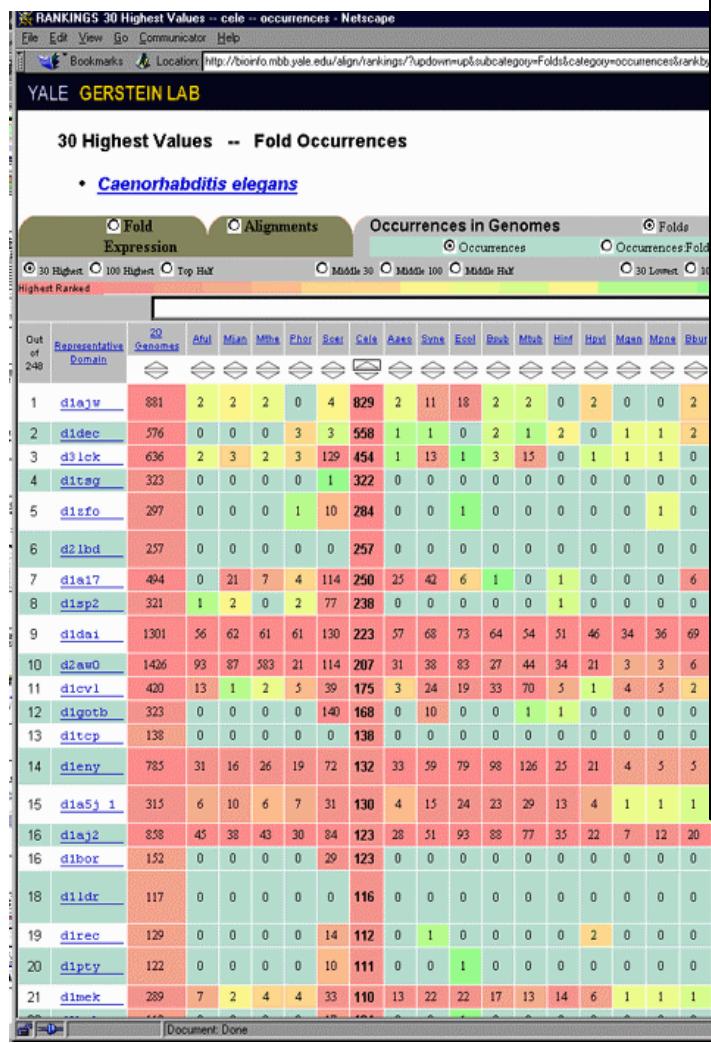
Phylogenetic Occurrence							Fold Classification			Gene Expression			Function & Interaction		
Yeast Genome	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
Kinase (cat. core)	1	3	-	-	-	x	-	18	-	-	7	7	-	1p38	
7-bladed beta-propeller	2	12	x	x	x	x	x	5	-	-	9	-	18	1mda	
P-loop NTP hydrolases	3	9	4	1	3	1	2	2	-	-	8	1	3	1gky	
Alpha-alpha superhelix	4	7	-	x	-	5	16	6	-	-	3	2	-	2bct	
TIM-barrel	5	16	1	6	2	2	3	4	-	1	-	24	1byb		
Ferredoxin	6	10	2	15	10	6	1	3	-	3	-	-	8	5	3
Leu-zip fold	7	-	-	2	25	9	-	15	-	15	-	19	6	-	-
Rossman fold	8	14	3	9	1	12	5	3	-	4	-	-	9	3	3
Zn2-C6 DNA-bind dom.	9	x	x	x	x	x	x	15	19	-	7	-	10	-	2
Alpha-beta Hydrolases	10	11	18	10	6	-	-	-	-	-	-	-	1xel	1aw6	2ace



# PartsList

## Ranking

### Viewers



**Rankings2 - Netscape**

**YALE GERSTEIN LAB**

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

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

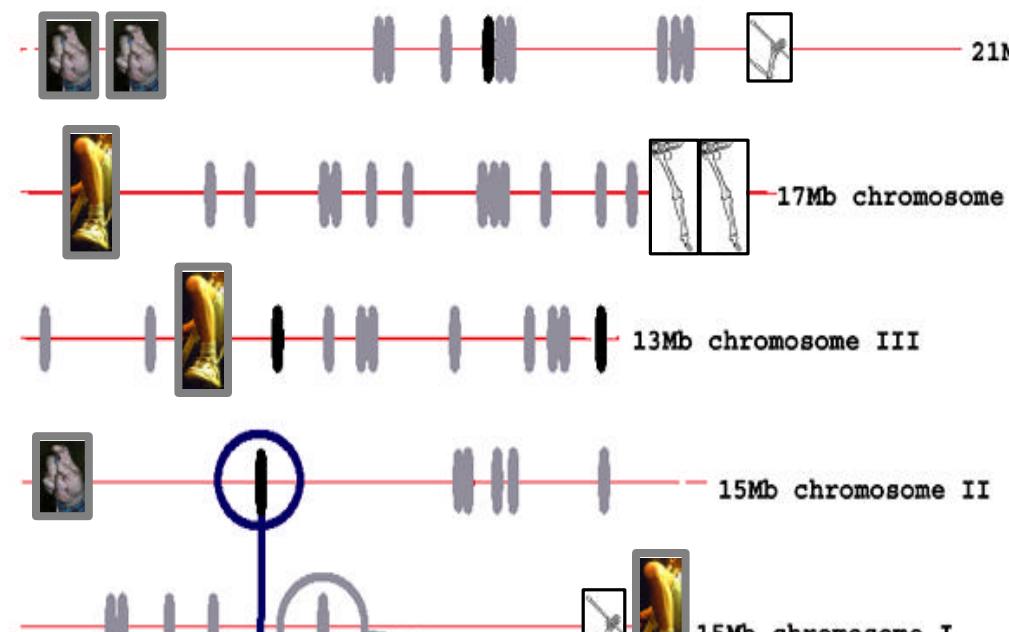
	fold occurrence in Scer	fold occurrence in Cele	fold percentage in reference genome	fold percentage in reference transcriptome
Max of all	140	829	11.3	11.1
Min of all	0	0	0.1	0
Average	6.9	22.9	0.5	0.5
Non zero hits	215	247	213	134
Rank again				
<a href="#">d1ajw</a> :1.002.001 Immunoglobulin-like beta-sandwich 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

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

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



# Pseudogenomics: Surveying “Dead” Parts

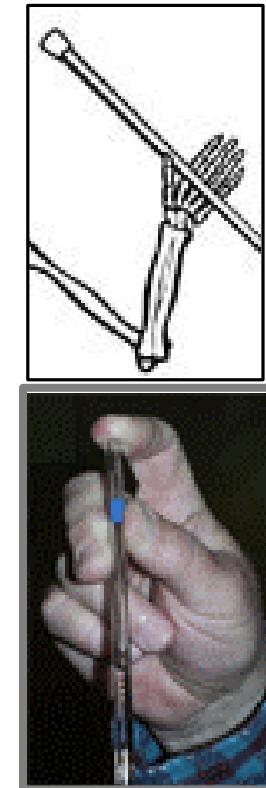


pseudogene fragment on worm chromosome II

TKRTSNGFGQDVVVDLFSILDGLVARAHXVLQDIFEFFAS  
KKMVTIFS#APHSPHSAPHYCAQFDNSAATVKV

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

MTAPMDVDNLMSRLLNVGMGGRLTTSVNEQELQTCCAVAKSVFASQASLLEVEPPIVC  
GDIHGQYS DLLRIFDKNGFPDVNFLFLGDYVDRGRQNIETICLMLCFKIKYPENFFMLR  
GNHECPAINRVYGFYEECNRRYKSTRLWSIFQDTFNWMPLCGLIGSRILCMHGGLSPHLQ  
TLDOLROLPRPQDPNPNSIGIDLLWADPDOWVKGWQANTRGVSYVFGODVVADVCSDLI  
DLVARAHOVVODGYEFFASKMVTIFSAPHYCGOFDNSAATMKV DENMVCTFVMYKPTPK  
SMRRG\*



Example of a potential ΨG with frameshift in mid-domain