# Scalable, Dynamic Analysis and Visualization for Genomic Datasets

Grant Wallace, Matthew Hibbs, Maitreya Dunham, Rachel Sealfon, Kai Li, and Olga Troyanskaya

> Olga Troyanskaya Assistant Professor Department of Computer Science & Lewis-Sigler Institute for Integrative Genomics Princeton University



Laboratory for BIOINFORMATICS and FUNCTIONAL GENOMICS

# Science has become data-centric

- Enormous amounts of data being observed
- Data then analyzed by experts in the field or computer scientists/statisticians
- Most disciplines critical to have human insight during data analysis => visualization

#### Data growing in genomics: systems-level studies

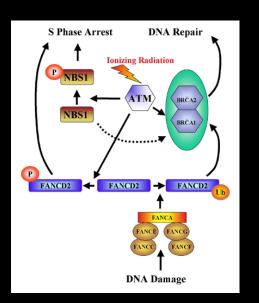
- Cells & organisms are complex systems
- Many biological processes & diseases result from multiple changes on molecular level
- To understand and cure cancer & other diseases, need to observe and model cellular processes on a systems level



Laboratory for BIOINFORMATICS and FUNCTIONAL GENOMICS

## Data-Knowledge gap in biology





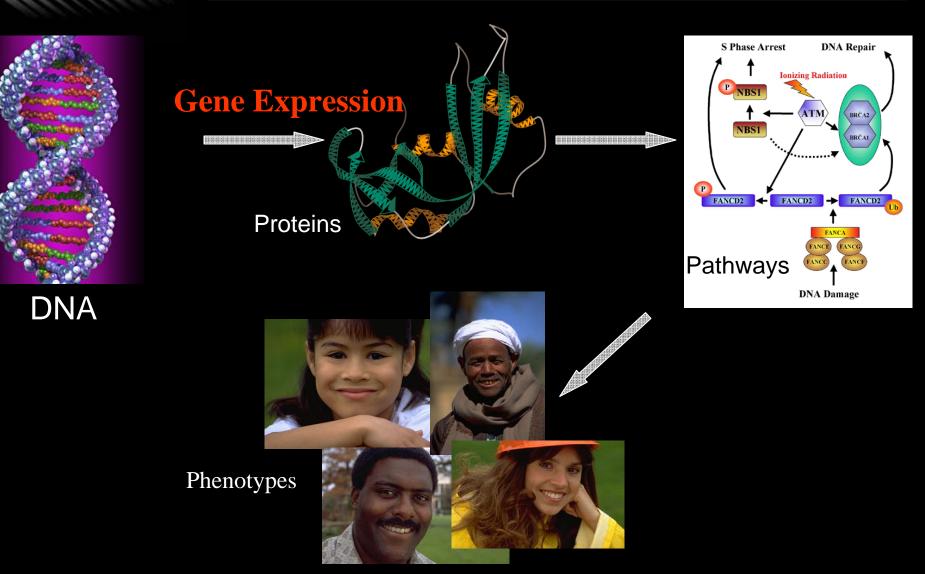
Explosion of genomic **DATA** 

**KNOWLEDGE** of components and inter-relationships that lead to function

#### Outline

- Background
- Visualization and Analysis Solutions
  - ForestView
  - SPELL
  - GOLEM
- Integration of Solutions

# Isn't genomics "done" with the human genome?



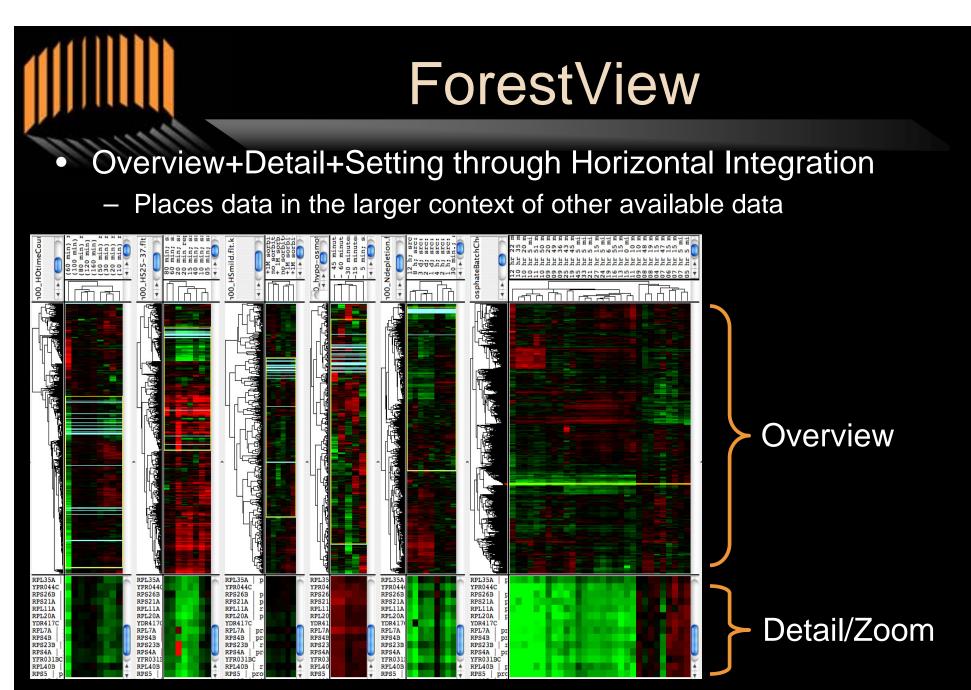
## Microarray Results

Raw Image from Spellman et al., 98

#### Challenges of Microarray Analysis

#### LOTS of data

- Hundreds of millions of measurements exploring cancer
- Strong need for collaboration
- Biologist's insight critical analysis => visualizationbased analysis is critical
- Limited by effective visualization capacities
  - Screen space at a premium
  - Need to identify relevant data subsets
- Most analysis small-scale due to these issues
- Future of discovery will require larger-scale, integrated analysis and visualization

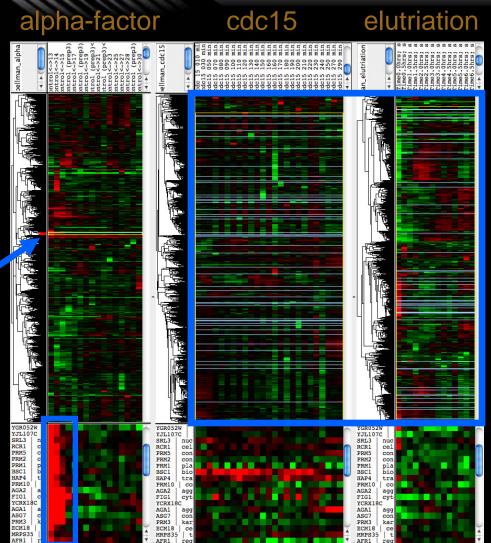


Six diverse datasets linked together

#### **ForestView Overview+Detail+Setting through Horizontal Integration** Places data in the larger context of other available data TEEEEEEEEEEEEEEEEEEEEEEE EEEEEEE 45 60 30 15 Nden ADDDALL HS25-XX A Selections 000014000 10H\_00r G+G++ 8 linked across datasets Overview **Details aligned** horizontally RPL35A RPL35A RPL357 RPL35A RPL3 RPL357 RPS26B RPS261 RPS26 RPSZIA RPL11A RPL11A RPL11 RPL1 RPL11 RPL11A RPL20A RPL20A RPT.20A RPL20A RPL20 RPT.207 Detail/Zoom YDR41 YDR4170 YDR417 YDR417C YDR41 YDR417C RPL7A RPL7A RPL7A RPL7 RPL7A RPL7A RPS4B RPS4B RPS4B RPS4P RPS4B RPS4B RPS23B RPS23B RPS23B RPS23 RPS23E RPS23B RPS4A RPS4A RPS4A RPS4P RPS4A RPS4A p YFR031E YFR03 YFR031BC YFR0: YFR03 YFR031BC RPL40B RPL40 RPL40B RPL4 RPL40 RPL40B RPS5

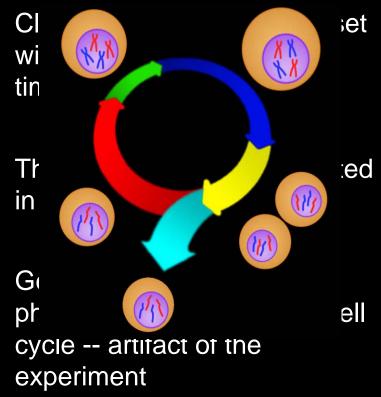
Six diverse datasets linked together

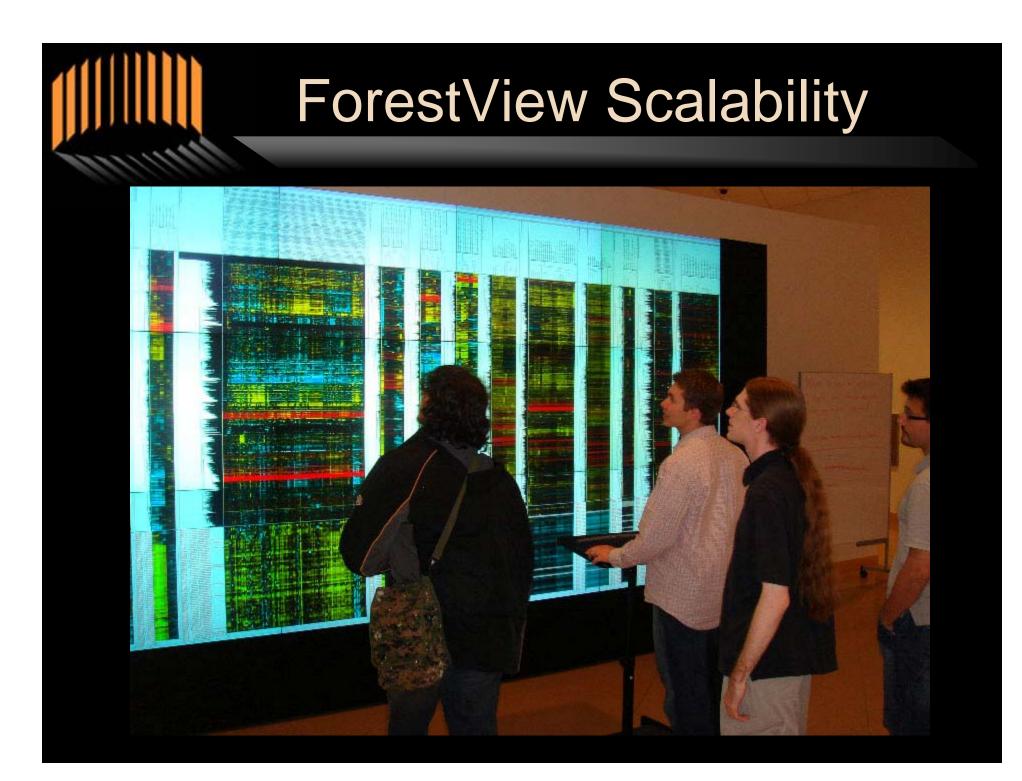
#### **ForestView Observation**



Spellman cell cycle datasets

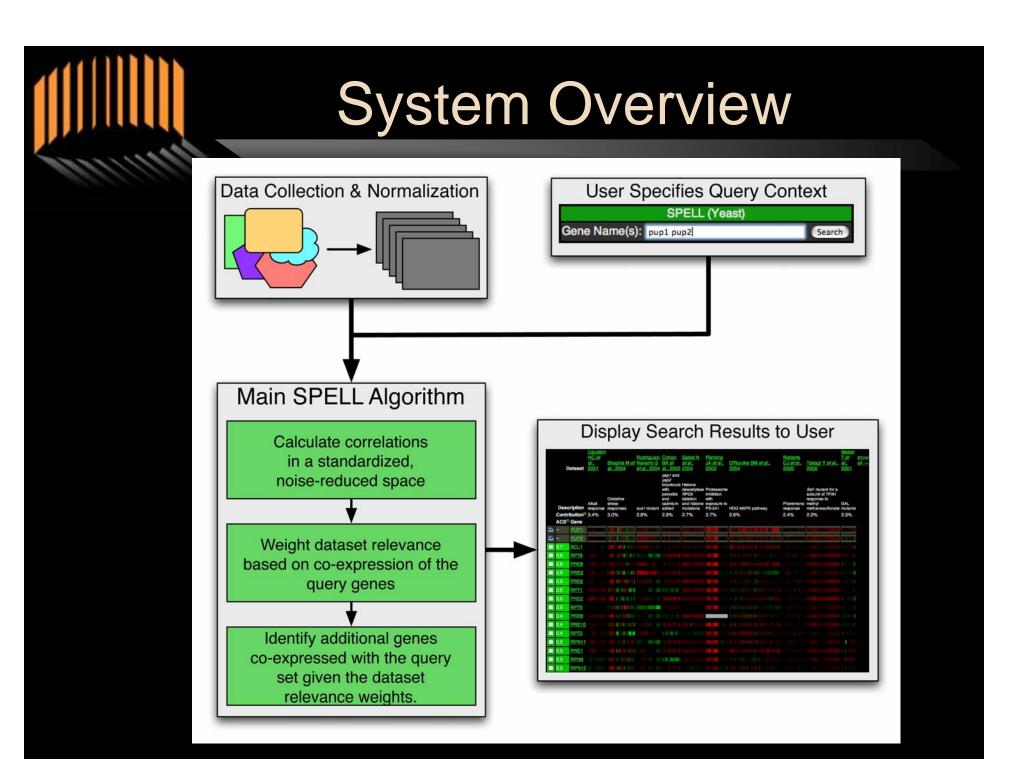
Three datasets studying cell cycle phases -- cancer can be related to cell cycle defects

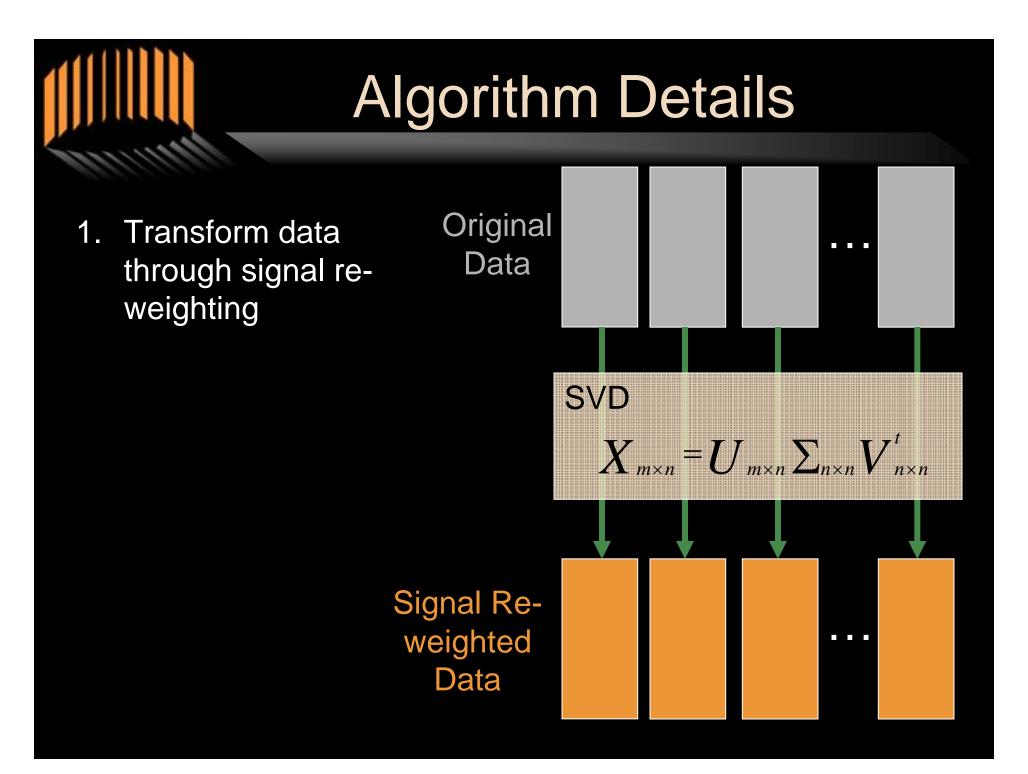




#### Analysis Across Multiple Datasets

- ForestView increases amount of data simultaneously viewed...
- But still much more data available
- Need to identify relevant datasets, and identify novel genes related to functions
- Our solution:
- SPELL (Serial Patterns of Expression Levels Locator) – a similarity search and visualization engine





### **Algorithm Details**

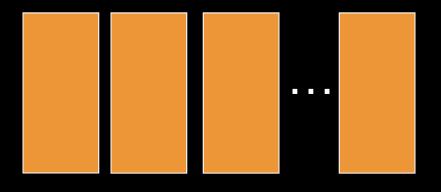
- Transform data through signal reweighting
- Given a query, determine dataset weights

#### User specifies N query genes

SPELL (Yeast)
Gene Name(s): ARP8, INO80
Search

Avg z-score between query pairs is dataset weight

$$w_{d} = \left(\frac{2}{N(N-1)}\right)_{i=1}^{N-1} \sum_{j=i+1}^{N} (z'_{q_{i},q_{j}})$$



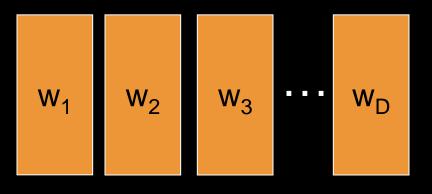
## Algorithm Details

- Transform data through signal reweighting
- Given a query, determine dataset weights
- Identify additional co-expressed genes



Scores for all other genes are weighted average z-scores





## **Algorithm Details**

- Transform data through signal reweighting
- 2. Given a query, determine dataset weights
- Identify additional co-expressed genes
- 4. Display results

	Dataset		Diamide treatment	<u>et al.,</u> 2005	AJ et al., 2004 Limitation		<u>Bro C el</u> <u>al.,</u> 2003
			time course	Genotoxic stress	by Phosphate		Lithium response
		ibution		4.1%	3.8%		3.2%
	ACS			,.	0.070		
⊠		CTR9					
		MED2					
	1.4	BIR1					
	1.4	LSM3					
	1.4	<u>TAF12</u>					
	1.4	ENT1					
	1.4	TFG1					
	1.4	RSC58					
	1.4	YNG2					
	1.3	NUP2					
	1.3	ENT5					
	1.3	ARP8					
	1.3	DJP1					
	1.3	PRP3					
	1.3	VPS16					
	1.3	FHL1					
	1.3	SRB4					
	1.3	SET2					

all →

## Search from Cell Cycle Cluster



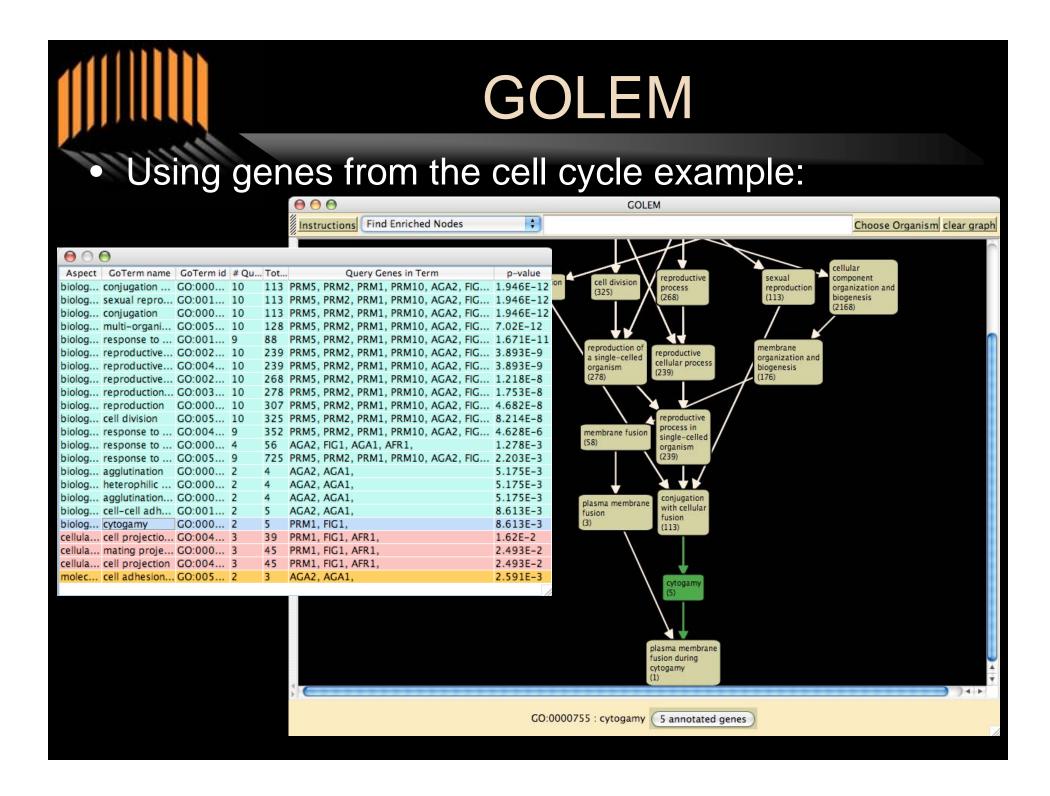
#### SPELL - S. cerevisiae

#### Search results<sup>II</sup> for DEM1, CDC28

	Refi	ine Sea	arch											
New Search				Spellman				Gasch AP et	Saldanha	Gasch AP et		Yoshimoto	Fellmann DD	show
Dataset Listing			Dataset	<u>PT et al.,</u> 1998		PT et al., 1998	Pitkanen JP et al., 2004	al., 2000	2004	al., 2000	<u>AJ et al.,</u> 2004	<u>H et al.,</u> 2002	Eriksson PR et al., 2005	all →
				Cell cycle,		0-1-1-	<b>D</b> harabaan a ina ana ana ana ana ana ana ana an	Linet Oberly		Dithiothrietol			spt10 global	
Show Expression evels				cdc15 block-	GAL		Phosphomannose isomerase pmi40 deletion strain	from various		exposure time			transcription regulator null	
-04015	Description Contribution					response to excess mannose 3.3%	temp to 37C 3.0%	by Leucine 2.9%	course (y13) 2.6%	by Uracil 2.5%	exposure 2.5%	mutant 2.5%		
		ACS		4.1 70	0.070	3.370	3.376	3.076	2.370	2.070	2.0%	2.070	2.370	
	₫-	-	CDC28			-								
		-	DEM1											
	1	1.9	POP7		1111									
		1.8	MSG5		1.0.00						I BH S			
		1.7	UBS1		III II									
		1.7	AME1											
	-	1.6	MCM2		1111									
		1.6	ARL1						10.001			1111		
		1.6	PRP9			10110								
		1.5	PPS1											
		1.5	FTH1									11 11 101		
	-	1.5	TSC10											
		1.5	ALG14											
		1.5	YBR255W											
	-	1.5	TAH1											
		1.5	SWD3									1111		
	-	1.5	CSH1											
		1.5	TOS1											
	-	1.5	ERF2											
		1.5	PBP2									Inter II		
	- 1	1.4	YBR259W											
		1.4	HSL7											

### GOLEM: going beyond genomic data

- Also need analysis tools for verification and interpretation of search or visual analysis results
- Gene Ontology (GO) curated hierarchical structure that represents known biology
- GOLEM tool for finding GO enrichment, viewing ontology structure



#### The future: dynamic, integrated search, analysis and visualization SPELL query used to identify datasets relevant to a gene set SPELL (Yeast) Datasets explored in ForestView, Gene Name(s): CDC28, CDC15, SWI4 Search related genes selected Results used to further explore current data, or refine SPELL search $\Theta \odot \Theta$ Aspect GoTerminame GoTermid # Ou., Tot., **Ouery Genes in Term** n-value piolog... conjugation ... GO:000... 10 113 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 1,946E-1 113 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 1,946E-12 piolog... sexual repro... GO:001... 10 piolog... conjugation GO:000... 10 113 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 1,946E-1 128 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 7,02E-12 piolog... multi-organi... GO:005... 10 olog... response to ... GO:001... 9 88 PRM5, PRM2, PRM1, PRM10, AGA2, FIG... 1.671E-1 piolog... reproductive... GO:002... 10 239 PRM5, PRM2, PRM1, PRM10, AGA2, FIG... 3.893E-9 reproductive... GO:004... 10 239 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 3,893E-9 iolog... reproductive... GO:002... 10 268 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 1,218E-8 piolog... reproduction... GO:003... 10 278 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 1,753E-8 iolog... reproduction GO:000... 10 307 PRM5, PRM2, PRM1, PRM10, AGA2, FIG., 4,682E-8 325 PRM5 PRM2 PRM1 PRM10 AGA2 FIG 8 214E-8 iolog cell division CO:005 10 352 PRM5, PRM2, PRM1, PRM10, AGA2, FIG... plog... response to ... GO:004... 9 4.628E-6 piolog., response to ... GO:000... 4 56 AGA2, FIG1, AGA1, AFR1, 1.278E-3 response to ... GO:005. 725 PRM5, PRM2, PRM1, PRM10, AGA2, FIG ... 2.203E-3 iolog... applutination GO:000... 2 AGA2, AGA1, 5.175E-3 heterophilic .... GO:000... AGA2, AGA1, 5.175E-3 piolog... applutination... GO:000... 2 AGA2, AGA1. 5.175E-3 AGA2, AGA1, iolog... cell-cell adh... GO:001... 2 8.613E-3 iolog... cytogamy GO:000... PRM1. FIG1. 8.613E-3 cellula... cell proiectio... GO:004... 3 PRM1, FIG1, AFR1, 1.62E-2 39 45 ellula... mating proje... GO:000... 3 PRM1, FIG1, AFR1 2.493Ecellula... cell projection GO:004... 3 45 PRM1, FIG1, AFR1, 2.493E-2 cell adhesion... GO:005 AGA2, AGA1 2.591E-3 YJL107 SRL3 YJL1070 SRL3 RCR1 PRM5 PRM2 PRM1 BSC1 HAP4 PRM10 AGA2 FIG1 YCRX1 AGA1 ASG7 PRM3 ECM18 GOLEM used to identify RCR1 PRM5 PRM2 PRM1 BSC1 HAP4 PRM1 AGA2 FIG1 YCRX AGA1 ASG7 PRM3 ECM1 PRM5 PRM2 PRM1 SSC1 IAP4 PRM1 VGA2 FIG1 VCRX VGA1 VGA1 VGA1 VGA3 PRM3 functional enrichments of clusters

MRPS3

### Conclusions

- Scalable, dynamic visualization-based analysis enables novel biological discoveries
- Input from biology researchers critical in analysis
- Integration of multiple data sets is important in modeling and analyzing biological data
- Integration of visualization and analysis/search is important in genomics

### Acknowledgements

#### **Research Staff**

\* Camelia Chiriac Postdoctoral Fellows

- \* Florian Markowetz
- \* Edo Airoldi
- \* David Hess

Graduate Students

- \* Chad Myers
- \* Matthew Hibbs
- \* Curtis Huttenhower
- \* Patrick Bradley
- \* Maria Chikina
- \* Yuanfang Guan
- \* Zafer Barutcuoglu

\* Lars Bongo Undergraduate Students

- \* Drew Robson
- \* Daniel Barrett
- \* Adam Wible
- \* Rachel Sealfon





Laboratory for BIOINFORMATICS and FUNCTIONAL GENOMICS http://function.princeton.edu



#### Collaborator:

#### Kai Li





#### **Central Dogma**



- Common cellular mechanisms to create proteins
- Understanding the function, coordination of proteins key to understand disease
- Gene expression microarrays give us a picture of cellular states

