

#### Biological Network Analysis

Mark B Gerstein

Yale

slides at

#### Lectures.GersteinLab.org

(See Last Slide for References & More Info.)

# Networks occupy a midway point in terms of level of understanding



 $\begin{array}{c} & & & \\ & & & \\ & & \\ & & & \\ & & \\ & & & & \\ & & & \\ & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \\ &$ 

1D: Complete Genetic Partslist ~2D: Bio-molecular Network Wiring Diagram

3D and 4D: Detailed structural understanding of cellular machinery (e.g. ribosome in different functional states)

### Networks as a universal language



### <u>Combining networks forms an ideal way</u> of integrating diverse information



- Why Networks?
- Network Comparisons

(reg. net. in many organisms)

- in rel. to social hierarchy
- scaling in rel. to partnerships
- Computer OS Comparisons
- Network Dynamics Across Environments

(prokaryote metab. pathways)

- Metabolic Pathways
- Entry pts. (Mem. Proteins)

### Outline: Molecular Networks



### Network Comparison #1 Comparing the Yeast Regulatory Network to a Governmental Hierarchy





### Determination of "Level" in Regulatory Network Hierarchy with Breadth-first Search

I. Example network with all 4 motifs



III. Finding mid-level nodes (Green)



II. Finding terminal nodes (Red)





### **Regulatory Networks have similar** <u>hierarchical structures</u>







[Yu et al., Proc Natl Acad Sci U S A (2006)]

### Yeast Regulatory Hierarchy: the Middle-managers Rule



10 - Lectures.GersteinLab.org

### Yeast Network Similar in Structure to Government Hierarchy with Respect to Middle-managers



### <u>Characteristics of Regulatory Hierarchy:</u> <u>Middle Managers are Information Flow</u> Bottlenecks



### Network Comparison #2 Broadening the comparison to different types of hierarchies & different types of biological networks



### **Different kinds of Hierarchies**



- Well-defined levels and a clear chain of command
- A military hierarchy



- Without well-defined levels & with more coregulatory partnerships
- A club or a scientific collaboration network



Intermediate

- High degree of coregulation and can be organized into hierarchies
- A law firm

	Autocratic	Democratic	Intermediate
Betweenness 🛆	1.03	3.6	3.3
Betweenness (	4.1	1.08	3.4
Var. Betw. (triangles)	2.1	0.58	1.74
Var. Betw. (all)	2.9	1.4	1.9
D <sub>Net-collab</sub>	0	0.91	0.71





[Bhardwaj et al., PNAS (2010), in press]

# Higher species are more show more collaborative nodes (more democratic)



[Bhardwaj et al., PNAS (2010), in press]

### **Collaborative Nature of the Levels**



[Bhardwaj et al., PNAS (2010), in press]



### **Collaboration Between Levels**



$$D_{betw-level-collab}^{L,M} = \frac{\sum_{A \in L} \sum_{B \in M} \frac{G_A \cap G_B}{G_A \cup G_B}}{\left|L\right| \bullet \left|M\right|}$$

[Bhardwaj et al., PNAS (2010), in press]



### Middle Managers Interact the Most in Efficient Corporate Settings

- Floyd, S. W. et al (1992)
  Middle management involvement in strategy and its association with strategic type Strategic Management Journal 13, 153-167.
- Woodward, J. (1982) Industrial Organization: Theory and Practice (Oxford University Press, Oxford).
- Floyd, S. W. et al (1993)
  Dinosaurs or Dynamos?
  Recognizing Middle
  Management's Strategic Role
  The Academy of Management Executive 8, 47-57.
- Floyd, S. W. et al (1997)
  Middle management's strategic influence and organizational performance

Journal of Management Studies 34, 465-485.



Network Comparison #3: Comparing the structure and evolution of biological regulatory networks and software call graphs



### **E. Coli Transcriptional regulatory network vs Linux kernel call graph**

0	
	XE

		<i>E. coli</i> transcriptional regulatory network	Linux call graph
	Nodes	Genes (TFs & targets)	Functions (subroutines)
Basic properties of	Edges	Transcriptional regulation	Function calls
systems	External constraints	Natural environment	Hardware architecture, customer requirements
	Origin of evolutionary changes	Random mutation & natural selection	Designers' fine tuning



	<i>E. coli</i> transcriptional	Linux call graph
	regulatory network	
Number of nodes	1378	12391
Number of persistent nodes	72* (5%)	5120 (41%)
Number of edges	2967	33553
Number of modules	64	3665
Number of comparative	200 bacterial genomes	24 versions of kernels
references		
Years of evolution	Billions years	20 years



[Yan et al., PNAS (2010), in press]

### **Comparison: hierarchical organization**

% in E. coli % in Linux regulatory call graph network Pyramidal vs Top-heavy master 4.6 29.6 regulator middle 5.1 58.2 manager workhorse 90.2 12.3 10<sup>0</sup> 10<sup>0</sup> -----out--deg ---out-deg ←in-deg ↔ in-deg 10 10 Probability Distribution 10<sup>-2</sup> Degree distribution Roles of hubs 10<sup>-2</sup> 10<sup>-3</sup>  $10^{-4}$ ' 10<sup>-3</sup>⊧ out-degree hubs 10<sup>-5</sup> in-degree hubs e.g. "crp" e.g. "printk" 10 \_\_\_\_ 10<sup>\_6</sup> 10<sup>4</sup> 10<sup>0</sup> Degree [Yan et al., PNAS (2010), in press] 10<sup>0</sup> 10<sup>2</sup> 10<sup>2</sup> 10<sup>4</sup>

#### **Comparison: organization of modules**



### **Comparison of persistent components**

 Persistent genes (preserve among different genomes) vs persistent functions (preserve among different releases)



specialized proteins are preserved across genomes

- Building of the hierarchy:
  - TRN: Bottom up. Regulatory changes are the main driving forces of evolution
  - ♦ Call graph: top down

### **Evolutionary rate of persistent functions**



### Why and so what?

The difference can be explained by the nature of hubs evolution: tinkering vs design Getweenness Centrality Kim et.al. PNAS 2007

- Independent modules:
  - robust
  - costly: the system needs a variety of tools for different tasks
- Overlap modules (reuse):
  - Less robust:
    - Breakdown of a generic component is harmful to the whole system
    - Fragile in the sense any change in a module may require compensating changes in a generic function
  - cost effective: components can be used by need to be fine-tuned

# Network Dynamics Across Environments: Metabolic Pathways

How do molecular networks change across environments? What pathways are used more ? Used as a biosensor ?



#### What is Metagenomics?

#### **Traditional Genomics**



#### **Metagenomics**



#### Sorcerer II Global Ocean Survey



Sorcerer II journey August 2003- January 200

Sample approximately every 200 miles

#### Sorcerer II Global Ocean Survey



#### **Extracting Environmental Data from Other Sources**

Sample Depth:	1 meter	
Water Depth:	32 meters	
Chlorophyll:	4.0 ug/kg	
Salinity:	31 psu	
Temperature:	11 C	
Location: 41°5'28'	'N, 71°36'8''W	
		Unites States Unites States Mexico Hondusos Costa Rica e Parama Galapagos Islands 10 sates

Annual Phosphate [umol/I] at the surface



World Ocean Atlas 2005 NOAA/NODC

Nutrient Features Extracted: Phosphate Silicate Nitrate Apparent Oxygen Utilization Dissolved Oxygen 35.5

34.5

33.5

#### 40% of Oceans are Impacted by Humans

\* Resolution is 1 km square

\* Value of a activity at a particular location is determined by the type of ecosystem present:

Impact = ∑ Features \* Ecosystem \* impact weight



### Anthropogenic Features Extracted:

**Ultraviolet radiation** 

Shipping

Pollution

**Climate Change** 

**Ocean Acidification** 







### Expressing data as matrices indexed by site, env. var., and pathway usage

[Rusch et. al., (2007) PLOS Biology; Gianoulis et al., PNAS (in press, 2009]



#### Canonical Correlation Analysis: Simultaneous weighting



#### Canonical Correlation Analysis: Simultaneous weighting



### CCA: Finding Variables with Large Projections in "Correlation Circle"



The goal of this technique is to interpret cross-variance matrices We do this by defining a change of basis.



## Strength of Pathway co-variation with environment



Environmentally Environmentally invariant variant



Gianoulis et al., PNAS 2009

Conclusion #1: energy conversion strategy, temp and depth



Gianoulis et al., PNAS 2009

#### Conclusion #2: Outer Membrane components vary with the environment







Membrane proteins interact with the environment, transporting available nutrients, sensing environmental signals, and responding to changes

> Gianoulis et al., *PNAS* 2009 Patel et al. *Genome Research* 2010

## Network Dynamics Across Environments: Membrane Proteins (Pathway Entry Points)



### Membrane Proteins: Sensing and Responding the Environment



- 2.3 million predicted membrane proteins
- 1.2 million unique
- 850,000 mapped to 151 membrane protein COGs

**107** variant membrane protein families

44 invariant membrane protein families

20% have NO KNOWN FUNCTION

Patel and Gianoulis et al., (in press) Genome Research

#### Membrane Proteins co-vary more than Metabolic Pathways

Median absolute structural Correlation Coefficient

Membrane Proteins = 0.3





Patel and Gianoulis et al., (in press) Genome Research

### **CCA** Limitations



**Dimension 1** 

- Four Major Obstacles (1) Strength and directionality of relationships not intuitive
- (2) Relative weights of features are difficult to visualize and compare.
- (3) No real means of quantifying covariation between specific sets of features.
- (4) Difficult to visualize or compare results in more than 2 dimensions.

#### Protein Families and Environmental Features Network (PEN)



Distance: Dot product between 1st and 2nd Dimension of CCA

# Protein Families and Environmental Features Network (PEN)



"Bi-modules": groups of environmental features and membrane proteins families that are associated

UV, dissolved oxygen, apparent oxygen utilization, sample depth, and water depth are not in the network





#### Bi-module 2: Iron Transporters/Pollution/Shipping



#### **Bi-module 2: Iron Transporters/Pollution/Shipping**



Rigwell A. J. (2002) Phil. Trans. R. Soc. Lond.

#### **Bi-module 2: Iron Transporters/Pollution/Shipping**



-Negative correlation between COG4558 and COG0609 and dust/pollution values (p-value <0.01)

- Searching the BRENDA database for enzymes using iron as a cofactor reveal that an increase in these two COGs negatively correlated to the amount of enzymes present that required iron.

### **Biosensors: 4 logs in 4 years Beyond Canaries in a Coal Mine**

(Moore's law) 1.5x/yr for electronics vs 10x/yr for DNA Sequencing

\$1000 Human genome ~ \$1 E. coli \$100 Human genome ~\$.10 E. coli





Carr and Church, Nat Biotech 2009

- Why Networks?
- Network Comparisons

(reg. net. in many organisms)

- in rel. to social hierarchy
- scaling in rel. to partnerships
- Computer OS Comparisons
- Network Dynamics Across Environments

(prokaryote metab. pathways)

- Metabolic Pathways
- Entry pts. (Mem. Proteins)

### Outline: Molecular Networks



### **Conclusions: Comparison of Social and Regulatory Hierarchies**

- Middle managers dominate, sitting at info. flow bottlenecks
- Democratic v Autocratic
- Collaborative (locally democratic) fraction of networks increases with organism complexity
- Middle managers most collaborative
- Most interaction occur between 2 middle managers (as seen in efficient corporate hierarchies)



		<i>E. coli</i> transcriptional regulatory network	Linux call graph
Hierarchical organization	Structure	Pyramidal	Top-heavy
	Characteristic hubs	Upper-level TFs with high out-degree	Generic workhorse functions with high in-degree
Organization of modules	Downstream modules as labeled by	Master TFs responsible for sensing environmental signals	High-level starting functions which initiate execution for specific tasks
	Node reuse	Low	High
	Overlap between modules	Low	High
Persistent nodes	Characteristics	Specialized (non- generic) workhorses	Generic or reusable functions
	Location in hierarchy	Mostly bottom	Mostly top
	Evolutionary rate	Mostly conservative (e.g. dnaA)	Conservative (e.g. strlen) & adaptive (e.g. mempool_alloc)
Design principles	Building of hierarchy	Bottom up	Top down
	Optimal solution favors	Robustness	Cost-effectiveness (reuse of components)



### Conclusions: Network Dynamics Across Environments

- Developed approach to connect quantitative features of environment to usage of pathways & families
  - CCA + PEN
- Applied to available aquatic datasets, identified footprints predictive of environment (potentially useful as biosensor)
- Integration of geospatial data can highlight unexpected trends as anthropogenic factors seem to be reflected in microbial function

- Specific Conclusions
  - Strong correlation exists between a community's energy conversion strategies & env. parameters (e.g. temperature & chlorophyll)
  - Relation between Fe and P transporters & amt. of chemical in environment
    - For Fe illustrates impact of pollution & shipping

N Bhardwaj K-K Yan P Patel T Gianoulis H Yu

A Paccanaro K Yip R Bjornson G Fang Y Xia J Korbel J Raes P Bork D Engelman

M Snyder



Acknowledgements

> Job opportunities currently for postdocs & students

### Networks.GersteinLab.org

### **Default Theme**

- Default Outline Level 1
  - Level 2



### **More Information on this Talk**

SUBJECT: Networks

**DESCRIPTION:** 

New York Academy of Sciences, A Look at the Tools and Comparative Approaches of Systems Biology, 2010.06.10, 17:00-17:40; [I:NYSYSBIO] (Medium-length networks talk, derived from [I:BROWNMATH] with metagenomics updated.)

NOTES:

This PPT should work on mac & PC. Paper references in the talk were mostly from Papers.GersteinLab.org.

PERMISSIONS: This Presentation is copyright Mark Gerstein, Yale University, 2010. Please read permissions statement at http://www.gersteinlab.org/misc/permissions.html . Feel free to use images in the talk with PROPER acknowledgement (via citation to relevant papers or link to gersteinlab.org).

PHOTOS & IMAGES. For thoughts on the source and permissions of many of the photos and clipped images in this presentation see http://streams.gerstein.info . In particular, many of the images have particular EXIF tags, such as kwpotppt , that can be easily queried from flickr, viz: http://www.flickr.com/photos/mbgmbg/tags/kwpotppt .