### Large-scale Analysis of Molecular Networks



Mark B Gerstein Yale

slides at Lectures.GersteinLab.org

(See Last Slide for References & More Info.)

### The problem: Grappling with Function on a Genome Scale?



- 250 of ~530 originally characterized on chr. 22 [Dunham et al. Nature (1999)]
- >25K Proteins in Entire Human Genome (with alt. splicing)



3 - Lectures.GersteinLab.org

(c) '09

## Some obvious issues in scaling single molecule definition to a genomic scale

- Fundamental complexities
  - ◊ Often >2 proteins/function
  - ♦ Multi-functionality:2 functions/protein
  - Role Conflation: molecular, cellular, phenotypic

### **Networks (Old & New)**



Same Genes in High-throughput Network

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







### 1D: Complete Genetic Partslist

~2D: Bio-molecular Network Wiring Diagram 3D: Detailed structural understanding of cellular machinery <u></u>

### **Networks as a universal language**





### **Network** pathology & pharmacology



[Adapted from H Yu]

9 Lectures.GersteinLab.org (c) 2009

# 10 Lectures.GersteinLab.org (c) 2009

### **Outline: Molecular Networks**

- Why Networks?
- Network Structure: Key Positions
  - ♦ Hubs & Bottlenecks
  - $\Diamond$  Tops of a Hierarchy
  - $\Diamond \, \mathsf{RE}\text{-score} \ \mathsf{nodes}$
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment



### **Different Types of Molecular Networks**



**Protein-protein Interaction networks** 









[Toenjes, *et al*, *Mol. BioSyst.* (2008); Jeong *et al*, *Nature* (2001); [Horak, et al, Genes & Development, 16:3017-3033; DeRisi, Iyer, and Brown, Science, 278:680-686]

Undirected

#### Directed

Metabolic pathway networks

### **Global topological measures**

Indicate the gross topological structure of the network



### **Scale-free networks**

Power-law distribution



Hubs dictate the structure of the network

[Barabasi]

### Hubs tend to be Essential

Integrate gene essentiality data with protein interaction network. Perhaps hubs represent vulnerable points?

[Lauffenburger, Barabasi]



### **Relationships extends to "Marginal Essentiality"**

Marginal essentiality measures relative importance of each gene (e.g. in growth-rate and condition-specific essentiality experiments) and scales continuously with "hubbiness"



### Another measure of Centrality: Betweenness centrality

## Betweenness of a node is the number of shortest paths of pairs of vertices that run through it -- a measure of information flow.

Freeman LC (1977) Set of measures of centrality based on betweenness. Sociometry 40: 35–41.



Girvan & Newman (2002) PNAS 99: 7821.

### **Betweenness centrality -- Bottlenecks**

### Proteins with high betweenness are defined as *Bottlenecks* (top 20%), in analogy to the traffic system







[Yu et al., PLOS CB (2007)]



Non-hub-bottleneck **node** 

) Hut

Hub-non-bottleneck **node** 

Non-hub-non-bottleneck node

## Bottlenecks are what matters in regulatory networks



### Signaling transduction pathways are directed



[Xianglin Shi]

### Bottlenecks in signaling pathways are important



# 22 Lectures.GersteinLab.org (c) 2009

### **Outline: Molecular Networks**

- Why Networks?
- Network Structure: Key Positions
  - ♦ Hubs & Bottlenecks
  - $\Diamond$  Tops of a Hierarchy
  - $\Diamond \, \mathsf{RE}\text{-score} \ \mathsf{nodes}$
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment





### 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** hierarchical structures





E. coli

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

S. cerevisiae

**5** - Lectures.GersteinLab.org (a) 00 N

### **Example of Path Through Regulatory Network**



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



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



### **Characteristics of Regulatory Hierarchy:** Middle Managers are Information Flow **Bottlenecks**



Average betweenness at each level

5

Average betweenness (x1000)

0

15

10

### **Characteristics of Regulatory Hierarchy: The Paradox of Influence and Essentiality**



[Yu et al., PNAS (2006)]

# **31 Lectures.GersteinLab.org** (c) 2009

### **Outline: Molecular Networks**

- Why Networks?
- Network Structure: Key Positions
  - ♦ Hubs & Bottlenecks
  - $\Diamond$  Tops of a Hierarchy
  - $\Diamond \, \mathsf{RE}\text{-score} \ \mathsf{nodes}$
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment



- How much does a regulator influence its targets?
- For miRNA-target networks easy to calculate, as all influence is downregulation
  - ◊ target prediction via: TargetScan, PITA, PicTar, miRanda, ...
- Look at down-reg. genes in a sample & compare with targets of a specific micro-RNA
  - Ø more down-reg genes => stronger regulatory effect

### **RE-score: Another way to identify** <u>"important" network nodes</u>





**Application of RE-score to** measure changing miRNA effect in different conditions (ER- and ER+ breast cancer)

Cheng et al., Genome Biology, 2009



# 35 Lectures.GersteinLab.org (c) 2009

### **Outline: Molecular Networks**

- Why Networks?
- Network Structure: Key Positions
  - ♦ Hubs & Bottlenecks
  - $\Diamond$  Tops of a Hierarchy
  - $\Diamond \, \mathsf{RE}\text{-score} \ \mathsf{nodes}$
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment



### What is metagenomics?

#### **Genomics Approach**



#### **Metagenomics Approach**



#### Partially Assemble and Annotate



36 - Lectures.GersteinLab.org (a) w

### **Global Ocean Survey Statistics (GOS)**



6.25 GB of data7.7M Reads1 million CPU hoursto process





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]



**39** - Lectures.GersteinLab.org<sub>(0)</sub>

### Canonical Correlation Analysis: Simultaneous weighting



### Canonical Correlation Analysis: Simultaneous weighting



### **Environmental-Metabolic Space**



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

Given 
$$X = \{x_1, x_2, ..., x_n\}$$
 and  $Y = \{y_1, y_2, ..., y_m\}$   

$$C = \sum_{X} \sum_{Y} \sum_{X,Y} \max_{X,Y} Corr(U,V) = \frac{a' \sum_{12} b}{\sqrt{a' \sum_{11} a} \sqrt{b' \sum_{22} b}}$$



### Strength of Pathway co-variation with environment



Environmentally Environmentally invariant variant





### <u>Conclusion #1: energy</u> <u>conversion strategy,</u> <u>temp and depth</u>



### **<u>Conclusion #2: Outer Membrane</u> components vary the environment**





### Conclusion #3: Covariation of AA biosynthesis and Import



Why is their fluctuation in amino acid metabolism? Is there a feature(s) that underlies those that are environmentally-variant as opposed to those which are not?

### Biosensors: Beyond Canaries in a Coal Mine



# 48 Lectures.GersteinLab.org (c) 2009

### **Outline: Molecular Networks**

- Why Networks?
- Network Structure: Key Positions
  - ♦ Hubs & Bottlenecks
  - $\Diamond$  Tops of a Hierarchy
  - $\Diamond \, \mathsf{RE}\text{-score} \ \mathsf{nodes}$
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment



### <u>Conclusions:</u> Analysis of Network Structure



- Centrality Measures in Protein Network
  - $\Diamond$  Hubs & Bottlenecks
  - Importance of later in regulatory networks

### Regulatory Network Hierarchies

- Middle managers dominate, sitting at info. flow bottlenecks
- Paradox of influence and essentiality
- Output Description of the second structure of the s

### Conclusions: Points of Network Centrality



- RE-score measures degree of (down) regulation of targets vs. non-targets
- Application to miRNA network
- Different RE-score of miRNAs can be used in cancer classification

### Conclusions: Networks Dynamics across Environments



- Developed and adapted techniques to connect quantitative features of environment to metabolism.
- Applied to available aquatic datasets, we identified footprints that were predictive of their environment (potentially could be used as biosensor).
- Strong correlation exists between a community's energy conversion strategies and its environmental parameters (e.g. temperature and chlorophyll).
- Suggest that limiting amounts of cofactor can (partially) explain increased import of amino acids in nutrient-limited conditions.





- an automated web tool

OI (vers. 2 : "TopNet-like Yale Network Analyzer")

			and the second se
File Edit View Favorites Tools Help			
😮 Back + 🕑 - 💌 📓 🎧 🔎 Search 👷 Favorites 🚱 🔗 + 🍃 🖸 - 🔜 📓			
Agdress 📳 http://networks.gersteinlab.org:8080/tyna/index.jsp?networkOrder=id8categoryOrder=id8cview=ADVANCED_VIEW8iksType=owned8ikstNetworkType=18ikstNetw 👱 🔂 Go 🛛 Links 🏾 🐑 🔹			
			TELIN
<b>4V</b>			
Getting started API WSDL Download tYNA Installation guide Plugins for Cytoscape Contact Known problems			
You are logged in as kevin. <u>Logout</u>	View: Simple.	Advanced	
List Owned  Biological  Attribute name)	= (Attribute value) V List		
workspace manager	Networks in database ( <u>upload</u> <u>download</u> )		
Load an existing network 🕢	ID Name		Tuak
	ID Name Creator date		
Load 14. Uetz 2000 yeast two	14 Uetz 2000 yeast two hybrid kevin 21-Feb-06 De	elete	
Into workspace 0 💌	15 lto 2001 yeast two hybrid kevin 21-Feb-06 De	elete	TEMA
Categorized by	16 Ho 2002 pull down kevin 21-Feb-06 De	elete	
	17 Gavin 2002 pull down kevin 21-Feb-06 De	elete	Display options: Default colors:
Load	18 Jansen 2003 PIT kevin 21-Feb-06 De	elete	Node: blue 🔻 Edge: lightgrey 🔽 Text: white 🔽
	19 MIPS yeast PPI kevin 21-Feb-06 De	elete	Special coloring: 0
Markanaga D:	21 BIND yeast data kevin 21-Feb-06 <u>De</u>	<u>elete</u>	C None
intersection(	22 DIP yeast data kevin 21-Feb-06 De	elete	Color gradient: Degree 🔽 of Original network 🔽 from green 🔽 to red
"Uetz 2000 yeast two hybrid", "te 2001 yeast two hybrid")	23 Kim 2006 structural interaction kevin 21-Feb-06 De	elete	O Color class: Class name: 🔽 white 🖃
Northease 1. (anato)	24 Han 2004 FYI data kevin 21-Feb-06 De	elete	Redraw
Workspace 7. (empty)	25 Luscombe 2004 regulatory kevin 21-Feb-06 De	elete 🗸	
workspace 2: (empty)			Statistics:
vvorkspace 3: (empty)	Categories in database ( <u>upload</u> <u>download</u> )		Conteners Node Edge Connected Degrees @ Clustering Coefficients Eccentricities @ Betweenness @
	ID Name Creator Creation date		Counts Count O Avg. S.D. Min. Max.
Multiple-network analysis		-	Whole 376 197 199 1 30 0 74 1 7 0 04 0 10 0 0 1 00 3 51 1 57 1 0 0 00 30 30 30 30 30
· · · · · · · · · · · · · · · · · · ·		Internet	network 2/0 10/ 109 1.30 0.74 1 7 0.04 0.19 0.00 1.00 2.51 1.57 1 9 3.60 20.22 0.00 200.00

Normal website + Downloaded code (JAVA) + Web service (SOAP) with Cytoscape plugin

[Yu et al., NAR (2004); Yip et al. Bioinfo. (2006); Similar tools include Cytoscape.org, Idekar, Sander et al]

### H Yu P Kim K Yip T Gianoulis C Cheng

A Paccanaro P Alves T Emonet P Cayting M Seringhaus Y Xia **J** Korbel A Sboner P Patel P Bork **J** Raes E Franzosa M Snyder N Bhardwaj **R** Alexander

### Acknowledgements



### Networks.GersteinLab.org

Job opportunities currently for postdocs & students

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

**<u>TITLE</u>**: Large-scale Analysis of Molecular Networks

**SUBJECT:** Networks

**DESCRIPTION**: Laufer Center Inaugural Symposium, Stony Brook University, NY; 2009.09.25, 09:00-09:30; [I:**STONEYBROOK**] (Short networks talk. Aiming to fit into 27' w. 3' questions. Modifications of **rescore\*** compared to previous networks talk.)

**NOTES**: Orange background slides are hidden and not meant to me shown. PPT works on mac & PC and has many photos. Paper references in the talk were mostly from Papers.GersteinLab.org. The above topic list can be easily cross-referenced against this website. Each topic abbrev. which is starred is actually a papers "ID" on the site. For instance, the topic **pubnet**\* can be looked up at http://papers.gersteinlab.org/papers/pubnet )

**PERMISSIONS**: This Presentation is copyright Mark Gerstein, Yale University, 2009. 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).

<u>PHOTOS & IMAGES</u>. 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 .