

The problem: Grappling with Function on a Genome Scale?



~1,200 protein-coding genes

(~950 pseudogenes)



3 - Lectures.GersteinLab.org

(c) '09

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

- Fundamental complexities
 - Role Conflation: molecular, cellular, phenotypic
 - Often >2 proteins/function
 - Also Multi-functionality:
 - 2 functions/protein
 - phenotypically e.g. Pleiotropic effects such as human PKU being involved in retardation & eczema
 - cellular role e.g. Depending on the molecule it interacts with HSP70 is involved with protein folding, translocation of proteins into mitochondia, biogenesis of certain subunits..

<u></u>

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

- Fundamental complexities
 - Role Conflation: molecular, cellular, phenotypic
 - Often >2 proteins/function
 - Also Multi-functionality:2 functions/protein
 - phenotypically e.g. Pleiotropic effects such as human PKU being involved in retardation & eczema
 - cellular role e.g. Depending on the molecule it interacts with HSP70 is involved with protein folding, translocation of proteins into mitochondia, biogenesis of certain subunits..
- Fun terms... but do they scale?....
 Starry night (P Adler, '94)



Hierarchies & DAGs of controlled-vocab terms but still have issues...



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

8 - Lectures.GersteinLab.org (a) 00

Networks as a universal language



Key comparisons

- Electronic circuits share similar design principles compared to biological networks (synthetic biology)
- Social Systems provide useful intuition
- Biological systems and computer operating systems (The genome has often been called the OS for a living organism)

 \Diamond Execute information processing tasks

 Adaptive systems shaped by changing environments (Natural vs man-made systems)



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



- Why Networks?
- Generating Networks
 - Scanning for Targets of Modular Domains
 - Propagating Known Information
- Central Network Points
 - Hubs & Bottlenecks (yeast ppi & reg. net)
- Networks & Variation (human ppi & miRNA-targ. net)
- Social Network Comparisons (reg. net. in many organisms)
 - in rel. to social hierarchy
 - scaling in rel. to partnerships
- Computer OS Comparisons

(E. coli reg. net)

<u>Outline: Molecular</u> <u>Networks</u>



Example: yeast PPI network

Actual size:

- \diamond ~6,000 nodes
 - → Computational cost: ~18M pairs
- ♦ Estimated ~15,000 edges
 → Sparseness: 0.08% of all pairs (Yu et al., 2008)

Known interactions:

- $\Diamond\,$ Small-scale experiments: accurate but few
 - \rightarrow Overfitting: ~5,000 in BioGRID, involving
 - ~2,300 proteins
- Large-scale experiments: abundant but
 noisy

 \rightarrow Noise: false +ve/-ve for yeast two-hybrid data up to

45% and 90% (Huang et al., 2007)



Different Types of Molecular Networks



Protein-protein Interaction networks





Undirected



Generating Networks

How do we construct large molecular networks. From connecting sequence patterns for modular interaction domains to matching sequences







Determining the interaction specificity

- Experimental methods to identify domain interaction specificity
 - \Diamond phage display experiment (e.g. SH3)
 - \Diamond the peptide library screen (e.g. S/T Kinase)
- Scan the target proteome with normalized PWMs



Binding Specificity Experiment

Α R Ρ

Each spot is a

one of the 20

amino acids

fixed at that

Example for

S/T Kinase

position

mixture of



Integrating structural and conservation features

- Calculate features for potential interaction sites:
 - \Diamond Surface accessibility
 - \Diamond Protein disorder
 - \Diamond Sequence Conservation
- Build a naïve Bayes classifier based on a validated training set to integrate the features



Generating Networks #2

How do we construct large molecular networks? From extrapolating correlations between functional genomics data with fairly small sets of known interactions, making best use of the known training data.



Training sets



4

1

?

?

?

Network prediction: features

• Example 1: gene expression



Gasch et al., 2000

Network prediction: features

• Example 2: sub-cellular localization



Data integration & Similarity Matrix



Learning methods

An endless list:

- Docking (e.g. Schoichet and Kuntz 1991)
- Evolutionary (e.g. Ramani and Marcotte, 2003)
- Topological (e.g. Yu et al., 2006)
- Bayesian (e.g. Jansen et al., 2003)
- Kernel methods
 - \Diamond Global modeling:
 - em (Tsuda et al., 2003)
 - kCCA (Yamanishi et al., 2004)
 - kML (Vert and Yamanishi, 2005)
 - Pairwise kernel (Pkernel) (Ben-Hur and Noble, 2005)
 - \Diamond Local modeling:
 - Local modeling (Bleakley et al., 2007)

Let's compare in a public challenge!

(DREAM: Dialogue for Reverse Engineering Assessment and Methods)

Our work: efficiently propagating known information

Training set expansion Local model 1 → Local model 2 • Motivation: lack of training examples Expand training sets horizontally Multi-level learning PPI predictions • Motivation: hierarchical nature of interaction ↓ • Expand training sets vertically DDI predictions • Expand training sets vertically ↓ • RRI predictions ↓

reconstruction challenge

Protein interaction



Yeast NADP-dependent alcohol dehydrogenase 6 (PDB: 1piw)

Protein-level features for interaction prediction: functional genomic information

Domain interaction



Pfam domains: PF00107 (inner) and PF08240 (outer)

Domain-level features for interaction prediction: evolutionary information

Residue interaction



Interacting residues: 283 (yellow) with 287 (cyan), and 285 (purple) with 285

Residue-level features for interaction prediction: physical-chemical information

[Yip and Gerstein, BMC Bioinfo. ('09, press)]

Combining the three problems



Empirical results (AUCs)

	Ind. levels	Unidirectional flow			Bidirectional flow			
Level		PD	PR	DR	PD	PR	DR	PDR
Proteins	71.68				72.23	72.50		72.82
Domains	53.18	61.51			71.71		68.94	71.20
Residues	57.36		54.89	53.81		72.26	63.16	77.86



- Highest accuracy by bidirectional flow
- Additive effect: 2 vs. 3 levels

Finding Central Points in Networks: Hubs & Bottlenecks

Where are key points networks ? How do we locate them ?



Global topological measures

Indicate the gross topological structure of the network



[Barabasi]



Regulatory and metabolic networks are *directed*

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









Non-hub-bottleneck **node**

Hub-non-bottleneck node

Non-hub-non-bottleneck **node**

Bottlenecks are what matters in regulatory networks



Networks & Variation

Which parts of the network vary most in sequence? Which are under selection, either positive or negative?



METHODOLOGY: MAP SNP AND CNV DATA ONTO ENSEMBL GENES, AND THEN MAP ENSEMBL GENES TO THE KNOWN INTERACTOME

ILLUSTRATIVE



* From Nielsen et al. *PLoS Biol.* (2005) and Bustamante et al. *Nature* (2005)

ADAPTIVE EVOLUTION CAN BE SEEN ON TWO DIFFERENT LEVELS



POSITIVE SELECTION LARGELY TAKES PLACE AT THE NETWORK PERIPHERY



Positive selection in the human interactome

Source: Nielsen et al. PLoS Biol. (2005), HPRD, and Kim et al. PNAS (2007)

CENTRAL PROTEINS ARE LESS LIKELY TO BE UNDER POSITIVE SELECTION

Degree vs. Positive Selection



 Peripheral genes are likely to under positive selection, whereas hubs aren't

Reasoning

Hubs

- This is likely due to the following reasons:
 - Hubs have stronger structural constraints, the network periphery doesn't
 - Most recently evolved functions (e.g. "environmental interaction genes" such as sensory perception genes etc.) would probably lie in the network periphery
- Effect is independent of any bias due to gene expression differences

* With a probability of over 80% to be positively selected as determined by Ka/Ks. Other tests of positive selection (McDonald Kreitmann and LDD) corroborate this result.

Source: Nielsen et al. PLoS Biol. (2005), Bustamante et al. Nature (2005), HPRD, Rual et al. Nature (2005), and Kim et al. PNAS (2007)

CENTRAL NODES ARE LESS LIKELY TO LIE INSIDE OF SDs



* Specifically, a number of the SDs are likely not fixed, but rather common CNVs in the reference genome Source: Database of genetic variation, HPRD, Rual et al. *Nature* (2005), and Kim et al. PNAS (2007)

Networks & Variation 2

Variation in the miRNA network



Analyze Regulation in microRNAtarget Network

- Relationship between target in degree (number of micro-RNAs that regulate gene) & evolutionary rate of gene?
 - \Diamond In deg. related 3' UTR size
- Expectation: more regulation, more constraint

Relationship between microRNA regulation and protein evolution



Human vs.	Number of genes	Correlation	P-value
chimpanzee	11326	-0.11	2.E-32
mouse	13280	-0.21	7.E-128
rat	12270	-0.20	4.E-107
COW	11683	-0.21	8.E-115
chicken	8061	-0.18	1.E-57

Important genes are regulated more intensively regulated by the microRNAs

[Cheng et al., BMC Genomics, 2009 (in press)]



For non-housekeeping genes, functionally critical genes are intensively regulated by miRNAs and prefer long 3'UTR.

housekeeping genes, however conserved, are selected to have shorter 3'UTRs to avoid miRNA regulation.



Social 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)]

S. cerevisiae

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

Social 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 _{Net-collab}	0	0.91	0.71





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



Collaborative Nature of the Levels





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.



Co-regulation Instantiates a Multi-Input Motif





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

Network Comparisons #3 Relating the size of co-regulation in partnership networks with the scale of the regulated





- Readily seen in many commonplace social contexts.
- An academic institution (say a high school), multiple teachers supervise the same set of students and have partnership interactions amongst themselves.

Building and Analysis of

Networks

-Edge placed if two regulators co-regulate

Network type	Species	Number of	Number of	Number of
		regulators	targets	interactions
Transcription	E. coli	160	1,420	3,123
Transcription	Yeast	157	4,410	12,873
Transcription	Mouse	144	1,092	2,403
Transcription	Rat	91	461	1,092
Transcription	Human	156	3,032	6,896
Phosphorylation	Yeast	87	1,337	4,083
Modification	Human	518	1,218	2,782

[Bhardwaj et al., PLoS Comp Biol (2010), in press]

<mark>6</mark>

 $\overset{w}{\supset}$

XV

Scaling of Regulators with Targets



Linear in *E. coli* (Due to operons) Exponential Saturation in others

[Bhardwaj et al., PLoS Comp Biol (2010), in press]

<u>Comparison to Social Networks: Partnership networks</u> <u>effectively saturate with increasingly complex output</u>



Software Network Comparison Comparing the structure and evolution of biological regulatory networks and software call graphs


E. Coli Transcriptional regulatory network vs Linux kernel call graph



		<i>E. coli</i> transcriptional regulatory network	Linux call graph
Basic properties of systems	Nodes	Genes (TFs & targets)	Functions (subroutines)
	Edges	Transcriptional regulation	Function calls
	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



Comparison: hierarchical organization





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 Spearman correlation r=0.25 P<10⁻⁷⁵ P<10⁻⁷⁵ Vin 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

- Why Networks?
- Generating Networks
 - Scanning for Targets of Modular Domains
 - Propagating Known Information
- Central Network Points
 - Hubs & Bottlenecks (yeast ppi & reg. net)
- Networks & Variation (human ppi & miRNA-targ. net)
- Social Network Comparisons (reg. net. in many organisms)
 - in rel. to social hierarchy
 - scaling in rel. to partnerships
- Computer OS Comparisons

(E. coli reg. net)

Outline: Molecular <u>Networks</u>



<u>Conclusions on Networks:</u> <u>Generation</u>



- Predicting Networks
 - Scanning for sequence motifs recognized by modular protein domains (motips)
 - Extrapolating from the Training Set
 - Principled ways of using known information in the fullest possible fashion
 - Multi-level learning

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



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

Conclusions: Connecting Networks & Variation



- Positive selection (adaptive evolution) at the network periphery
 - On a sequence level, it can be seen as positive selection of peripheral nodes
 - On a structural level, it can be seen as the pattern of SDs that display significantly higher allele frequencies in non-central genes
- miRNA network
 - Ø More highly regulated genes are under more constraint in miRNAtarget networks
 - $\Diamond\,$ Exception for housekeeping genes

<u>Conclusions: Comparison to</u> <u>Social Hierarchies</u>

- Regulatory Network Hierarchies
 - ◊ Middle managers dominate, sitting at info. flow bottlenecks
 - \Diamond Paradox of influence & essentiality

<u>Conclusions: Comparison to</u> <u>Social and Regulatory Hierarchies</u>

- Regulatory Network Hierarchies
 - ◊ Democratic v Autocratic
 - Ollaborative (locally democratic) fraction of networks increases with organism complexity
 - \Diamond Middle managers most collaborative
 - Ø Most interaction occur between two middle managers (as seen in efficient corporate hierarchies)
- Number of collaborative partners saturates even while scale of targets governed increases

 \Diamond Also seen in social networks



		<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)







- an automated web tool

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

Ele Edit View Favorites Iools Help			
🔾 Back 🔹 🕤 👻 😰 🏠 🔎 Search 🤺 Favorites 🤕 🍰 🛬 ک 🔯 🔹 🛄			
Agdress 👔 http://networks.gersteinlab.org:8080/tyna/index.jsp?networkOrder=id8categoryOrder=id8cview=ADVANCED_VIEW88istType=owned8iistNetworkType=18iistNetwetw 🐑 🐑 💿 🛛 Links 🦇 🐑 🔹			(WAZZA)
tY	NA	1	
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 💌 networks with (Attribute name) 📼	= (Attribute value) List		
Workspace manager	Networks in database (<u>upload</u> <u>download</u>)		
Load an existing network 🛛	ID Name Creation date		
Load 14. Uetz 2000 yeast two 💌	14 Uetz 2000 yeast two hybrid kevin 21-Feb-06 Delete		
Into workspace 0 💌	15 Ito 2001 yeast two hybrid kevin 21-Feb-06 <u>Delete</u>		TUNK
Categorized by Nil 💌	16 Ho 2002 pull down kevin 21-Feb-06 Delete		Display options:
Load	18 Jansen 2003 PIT kevin 21-Feb-06 Delete		Default colors:
	19 MIPS yeast PPI kevin 21-Feb-06 Delete		Node: blue 🔽 Edge: lightgrey 🔽 Text: white 🔽
Current working networks in your workspaces:	21 BIND yeast data kevin 21-Feb-06 Delete		C None
Workspace D: statFilter(degrees, geq, 1, value, neighbors=false, intersection(22 DIP yeast data kevin 21-Feb-06 Delete		Color gradient: Degree 🔽 of Original network 🔽 from green 🔽 to red
"Uetz 2000 yeast two hybrid",	23 Kim 2006 structural interaction kevin 21-Feb-06 Delete		C Color class: Class name: V white V
"ito 2001 yeast two hybrid"))	24 Han 2004 FYI data kevin 21-Feb-06 Delete		Redraw
Workspace 1: (empty)	25 Luscombe 2004 regulatory kevin 21-Feb-06 Delete		
Workspace 2: (empty)			Statistics:
vvorkspace J: (empty)	Categories in database (<u>upload</u> <u>download</u>)		Category Node Edge Connected Degrees @ Clustering Coefficients Coefficients Betweenness @
Multiple extensel evolution	ID Name Creator Creation date		Counts Count Count Avg. S.D. Min. Max.
Interpretence of the second se		-	Whole 276 187 109 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
ê 🗧	🔮 Internet		

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]

Acknowledgements

H Yu P Kim K Yip C Cheng N Bhardwaj K-K Yan H Lam

R Alexander G Fang M Seringhaus Y Xia J Korbel E Franzosa B Turk J Mok M Snyder



Networks.GersteinLab.org

Job opportunities currently for postdocs & students

More Information on this Talk

SUBJECT: Networks

DESCRIPTION:

```
Brown Applied Math, Providence, RI; 2010.04.09, 16:00-17:00;
[I:BROWNMATH] (Long networks talk, derived from [I:MBINETS],
including callgraph*, coregscaling*, reghier*, & motips* for 1st time.
Whole talk took 2 hrs. with questions.)
```

(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, 2008. 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 .