#### **Biological Networks**

Gavin Conant 163B ASRC <u>conantg@missouri.edu</u> 882-2931

# Types of Network

- Regulatory
- Protein-interaction
- Metabolic
- Signaling
- Co-expressing

#### **General principle**



Gene/protein/enzyme

#### Example

- We can create a network where actors and actresses are the nodes
- Two actors are joined if they co-stared in a film
- Kevin Bacon game













#### Example

- We can create a network where actors and actresses are the nodes
- Two actors are joined if they co-stared in a film
- Kevin Bacon game
  - On average, two actors can be linked by 3.65 films
  - Christopher Lee is actually more highly connected than Kevin Bacon

Watts, D. J., and S. H. Strogatz. 1998. Collective dynamics of 'small-world' networks. Nature **393**:440-442.

# Example—Protein interaction networks

- Many or most proteins bind to other proteins
  - Generally not covalent
  - But can still be
    - Long term (multi-peptide enzyme complexes)
    - Short term (signal transduction)
  - Generically treat this as protein interactions

# Methods to identify protein interactions

- Lots of small-scale experiments
- We will discuss the "genome-scale" methods
- Earliest: yeast two-hybrid
  - Uses yeast as a tool—can be done for proteins from any organism
  - But has been done at genome scale in yeast
     ③



 Originally, the transcription factor induced transcription of the gene and turned the colonies blue



- TF now functions if the two pieces are "close"
- Note however that the two pieces do not "stick together" on their own



- We are testing if the bait and prey interact
- If they do, they will bring the two halves of the TF close enough to turn on LacZ
- If the colonies turn blue, the two proteins interaction

# Yeast two-hybrid results

- > 60% of the proteins in the yeast genome have been tested for some interactions
- Question, if there are 5000 unique proteins in yeast, what is the possible number of interactions?
- Problems with two-hybrid method
  - We've altered the proteins—false negatives and false positives

# Yeast two-hybrid results

- Problems with two-hybrid method (cont.)
  - Probably won't be able to sample all possible interactions
  - Doesn't account for time of protein expression or location
    - We could infer that two proteins interact when they are in fact never at the same place at the same time

#### Mass-spec. methods

• Start by adding a tag to the "bait" protein of interest:



 Now grow cells with this construct: copies of the bait protein will be tagged

### Interactions vs. complexes

- Two-hybrid methods find pairwise protein interactions
- Here, we are looking at larger groups of proteins, aka complexes
- Complex is a somewhat vague term
  - Length of residue?
  - Functional?

# Mass spec. continued

- Extract the cell proteins without disturbing the protein complexes
- From those proteins, extract any complexes with a member having the tag

Several steps

- Uses an antibody to the tag for identification
- Result is complexes that have the bait protein as a member

#### Actual mass spec

- We now have a group of complexes which the bait protein is a member of
- What other proteins are present?



# Separation of fragments



- Magnetic field separates fragments based on:
  - Charge
  - Mass
- Result is a list of ions with mass and change that are present

# Identifying the ions

- 1. The mass spectrometer gives us a list of peptide fragment masses and charges
- 2. You search the genome for all possible peptide fragments and calculate their mass and charge
- 3. Match 1 to 2!
  - Obviously this is a computational challenge

### Mass-spec protein complexes

- We put the bait-centered protein complexes into the mass-spec and identify the peptides present
- Look those up in the genome to identify the proteins present
- This gives us a list of the proteins in the complex for that bait protein



Gavin, A. C. et al., 2002. Functional organization of the yeast proteome by systematic analysis of protein complexes. Nature **415**:141-147.

# In yeast

- ~590 proteins used as bait
- ~ 232 complexes found among these
  - ~130 proteins used as baits did not return a complex
  - But complexes can also interact with each other

# **Building networks**

- Protein interaction data has lots of uses
- One is to study biological networks
- Networks are an abstraction of the relationship among entities!



Jeong, H., S. P. Mason, A.-L. Barabási, and Z. N. Oltvai. 2001. Lethality and centrality in protein networks. Nature **411**:41-42.

### Network Measurements

- Simple: *Degree* distribution
  - How many edges?
  - Expressed as a distribution

# **Network Measurements**

- Simple: *Degree* distribution
  - How many edges?
  - Expressed as a distribution
- Common degree distributions:
  - Normal: Mean number of edges is  $\mu$
  - Exponential
    - P(x)≈e<sup>x</sup>
- Power-law
  - $P(x) \approx x^a$

# Comparing exponential and "scale-free" networks

- Exponential
  - P(x)≈e<sup>x</sup>
  - Most nodes have similar numbers of edges
  - Sensitive to node loss
- Scale-free
  - P(x) ≈ x<sup>a</sup>
  - Few nodes with many edges (aka hubs)
  - Less sensitive to random node loss



Figure 1 Visual illustration of the difference between an exponential and a scale-free network. **a**, The exponential network is homogeneous: most nodes have approximately the same number of links. **b**, The scale-free network is inhomogeneous: the majority of the nodes have one or two links but a few nodes have a large number of links, guaranteeing that the system is fully connected. Red, the five nodes with the highest number of links; green, their first neighbours. Although in the exponential network only 27% of the nodes are reached by the five most connected nodes, in the scale-free network more than 60% are reached, demonstrating the importance of the connected nodes in the scale-free network. Both networks contain 130 nodes and 215 links  $\langle k \rangle = 3.3$ ). The network visualization was done using the Pajek program for large network analysis: {http://vlado.tmf.uni-lj.si/pub/networks/pajek/pajekman.htm}.

Albert, R., H. Jeong, and A. L. Barabasi. 2000. *Nature* **406**:378-382.



Figure 2 Changes in the diameter d of the network as a function of the fraction f of the removed nodes. a, Comparison between the exponential (E) and scale-free (SF) network models, each containing N = 10,000 nodes and 20,000 links (that is,  $\langle k \rangle = 4$ ). The blue symbols correspond to the diameter of the exponential (triangles) and the scale-free (squares) networks when a fraction / of the nodes are removed randomly (error tolerance). Red symbols show the response of the exponential (diamonds) and the scale-free (circles) networks to attacks, when the most connected nodes are removed. We determined the r dependence of the diameter for different system sizes (N = 1,000; 5,000; 20,000) and found that the obtained curves, apart from a logarithmic size correction, overlap with those shown in a, indicating that the results are independent of the size of the system. We note that the diameter of the unperturbed (f = 0) scale-free network is smaller than that of the exponential network, indicating that scale-free networks use the links available to them more efficiently, generating a more interconnected web. b, The changes in the diameter of the Internet under random failures (squares) or attacks (circles). We used the topological map of the Internet, containing 6,209 nodes and 12,200 links  $\langle k \rangle = 3.4$ , collected by the National Laboratory for Applied Network Research (http://moat.nlanr.net/ Routing/rawdata/). c, Error (squares) and attack (circles) sun/vability of the World-Wide Web, measured on a sample containing 325,729 nodes and 1,498,353 links<sup>3</sup>, such that (k) = 4.59.

## Other statistics

- Mean path length
- Longest path (aka diameter)
- Clustering coefficient
  - Number of connections between your "neighbors" over the possible number of connections
- Number of components

# Features of protein interaction networks

- "Small world"
- There are a few proteins with very many interactions (hubs)
- Proteins are "cliquish":
  - If you and I interact, and I interact with Susan
  - It is more likely you and Susan will interact
- Most proteins "talk" to every other protein in very few steps

## Uses of interaction networks

- Predict disease-related genes
  - The density of protein interactions that a protein shares with other proteins can be used to predict whether it is likely to influence a particular disease
- Follow cell-signaling?

# Predictions from interaction networks

- The scale-free character of biological networks suggests that hub nodes may be more evolutionarily important
- Jeong et al., tested this by asking if genes for yeast protein interaction hubs were more likely to kill the cell when "knocked out" than other genes
  - Jeong, H., S. P. Mason, A.-L. Barabási, and Z. N. Oltvai. 2001.*Nature* **411**:41-42.



Characteristics of the yeast proteome. **a**, Map of protein-grotein interactions. The largest cluster, which contains -73% of all processes, is shown. The colour of a node signifies the phenotypic effect of removing the corresponding protein (ed., lethal: green, non-lethat; orange, slow growth; yellow, unknown). **b**, Connectivity distribution F(k) of interacting yeast proteins, giving the probability that a given protein interacts with *k* other proteins. The exponential cut-off<sup>5</sup> indicates that the number of proteins with more than 20 interactions is slightly less than expected for pure scale-free networks. In the absence of data on the link directions, all interactions have been considered as bidirectional. The parameter controlling the short-length scale correction has value  $k_i = 1$ . **c**. The fraction of essential proteins with exactly *k* links versus their connectivity, *k*, in the yeast proteome. The list of 1,572 mutants with known phenotypic profile was obtained from the Proteome database<sup>13</sup>. Detailed statistical analysis, including r = 0.75 for Pearson's linear correlation coefficient, demonstrates a positive correlation between lethality and connectivity. For additional details, see http://www.nd.edu/~networks/cell.

19.1

10.0

No. of firsts

# Other kinds of networks

- Regulatory networks
- Metabolic Networks
- Neural networks
- Gene expression networks





Genetics: A Conceptual Approach, Third Edition © 2009 W. H. Freeman and Company



© 2009 W. H. Freeman and Company

# Why study gene expression?

- Identify genes associated with pathways or diseases
- Diagnostics of disease
- Reconstruct the gene regulatory network

# Inferring the yeast regulatory network



Fig. 1. Systematic genome-wide location analysis for yeast transcription regulators. (A) Methodology. Yeast transcriptional regulators were tagged by introducing the coding sequence for a c-myc epitope tag into the normal genomic locus for each regulator. Of the yeast strains constructed in this fashion, 106 contained a single epitope-tagged regulator whose expression could be detected in rich growth conditions. Chromatin immunoprecipitation (ChIP) was performed on each of these

106 strains. Promoter regions enriched through the ChIP procedure were identified by hybridization to microarrays containing a genome-wide set of yeast promoter regions. (B) Effect of P value threshold. The sum of all regulator-promoter region interactions is displayed as a function of varying P value thresholds applied to the entire location data set for the 106 regulators. More stringent P values reduce the number of interactions reported but decrease the likelihood of false-positive results.

Lee, T. I. et al., Science 298, 799 (2002).



# Network "Motifs"

- We've looked at the "macro" structure of networks—what about their "micro" structure"
- Milo et al., defined the set of "motifs" for a network
  - Aka "subgraphs"
  - Milo, R., S. Shen-Orr, S. Itzkovitz, N. Kashtan,
    D. Chklovskii., and U. Alon. 2002. *Science* 298:824-827.

Fig. 1. (A) Examples of interactions represented by directed edges between nodes in some of the networks used for the present study. These networks go from the scale of biomolecules. (transcription factor protein X binds regulatory DNA regions of a gene to regulate the production rate of. protein. YÌ. through cells (neuron X is synaptically connected to neuron Y), organisms. (X to.



feeds on Y). (B) All 13 types of three-node connected subgraphs.

Network	Nodes	Edges	$N_{\rm real}$	$N_{\rm mand} \pm { m SD}$	Ziscore	$N_{\rm real}$	$N_{\rm mand}\pm {\rm SD}$	Ziscore	$N_{\rm real}$	$N_{\rm mand} \pm SI$	Z score
Gone regulation				X	Feed-	-X	X	Bi-fan			
(transcription)				¥.	forward		S.				
				V = 0	Joop	- <del>2</del> 2 Z	-mr W				
			b	ż		n n					
E, coli	424	\$19	40	$7 \pm 3$	10	203	$47 \pm 12$	13	i i		
S. cereviriae*	685	1,052	70	$11 \pm 4$	14	1812	$300 \pm 40$	41			
Neurons				X	Feed-	L K	X	Bi-fan	$\mathbb{R}^{\mathbf{X}}$	34	Bi-
				Y V	forward house	1 12	S.		v		parallel
				ŵ	roop	z	W		- 24	$\mu^{a}$	
			$\Rightarrow$ z						W	f	
C. depres†	252	509	125	$90 \pm 10$	3.7	127	$55 \pm 13$	5.3	227	$35 \pm 10$	20
Food webs				X	Three	17 X	<b>1</b> 4	34-			
				¥ v	chain	- 190 - 1	-94 	parallel.			
				$\tilde{\Psi}$		×24 -	$\mu^{\kappa}$				
				z			W				
Little Rock	92	984	3219	$-3120\pm 50$	2.1	7295	$2220\pm240$	25	l i		
Yiban	83	391	1182	$-1020\pm20$	7.2	1357	$230\pm50$	23			
St. Martin	42	205	409	$450 \pm 10$	NS	362	$130 \pm 20$	12			
Chesapeake	31	40	30	82 ± 4 - 435 ± 14	248 3.6	29	$2 \pm 2$ $20 \pm 20$	3			
Shima ith	22	199	184	222 ± 12 190 ± T	3.0	297	$80 \pm 20$ $80 \pm 25$	13			
B. Breek	25	104	181	$130 \pm 7$	7.4	267	$30 \pm 7$	32			
Electronic circuits				Σ.	Feed-	X	Y	Ni-fan	∠ X	м	Bi-
(forward logic chips)				Y.	forward	L D	<u></u>		Y.	z	parallel
			Ŵ		loop	¥2≤ -30¥ Z W		- 24	¥		
			- 5	z					· · ·	r -	
x15850	10,383	14,240	424	$2 \pm 2$	285	1040	1 = 1	1200	480	$2 \pm 1$	335
s38584	20,717	34,204	413	$10 \pm 3$	120	1739	$6 \pm 2$	800	711	$9\pm2$	320
s38417	23,843	33,661	612	$3 \pm 2$	400	2404	1 ± 1	2550	531	$2 \pm 2$	340
x9234 x13267	3,844	8,197	211	$2 \pm 1$ $2 \pm 1$	224	754	1#1	1050	209	2 + 1	200
Electronic cit	enits.		3		Three	v	v	Ridan	x-	⇒v	East.
(digital fractional multipliers)			A node		node			「本	1	node	
BB			7 34		feedback	<i>\₩</i> 2	Sale -			- V	feedback
			$Y \leqslant -$	— z	loop	Z	W		$z \in$	-w	loop
1216	122	199	10	1+1			1+1	3.8		$1 \pm 1$	
\$420	252	399	20	1±1	18	10	1±1	10	ŭ –	1±1	11
s#3#‡	512	8.19	40	$1 \pm 1$	38	22	$1 \pm 1$	20	23	$1 \pm 1$	25
World Wide Web			P₹ I		Feedback	X Fully		AX.	b.,	Uplinked	
				¥	with two	2	$\mathbb{Z}_{4}$	connected	7	7	motual
			6		matual	$y \iff z$ triad		$\dot{\mathbf{x}} \longleftrightarrow$	> Z	dyad	
				ź	ayaas						
nd.edu§	325,729	1.46e5	1.165	$2e3 \pm 1e2$	800	6.8c6	5e4±4e2	15,000	1.2e6	$1.04 \pm 2.0$	2 5000

Table 1. Network motifs found in biological and technological networks. The numbers of nodes and edges for each network are shown. For each motif, the numbers of appearances in the real network (N<sub>mal</sub>) and in the randomized networks (N<sub>cavd</sub> ± SD, all values rounded) (17, 18) are shown. The P value of all motifs is P < 0.01, as determined by comparison to 1000 randomized networks (100 in the case of the World Wide Web). As a qualitative measure of statistical significance, the Z score = (N<sub>real</sub> - N<sub>read</sub>)/SD is shown. NS, not significant. Shown are motifs that occur at least U = 4 times with completely different sets of nodes. The networks are as follows (18): transcription interactions between regulatory proteins and genes in the bacterium E. coli (11) and the yeast S. cerevisiae (20); synaptic connections between neurons in C. elegans, including neurons connected by at least five synapses (24); trophic interactions in ecological. food webs (22), representing pelagic and benthic species (Little Rock Lake), birds, fishes, invertebrates (Ythan Estuary), primarily larger fishes (Chesapeake Bay), lizards (St. Martin Island), primarily invertebrates (Skipwith Pond), pelagic lake species (Bridge Brook Lake), and diverse desert taxa (Coachella Valley); electronic sequential logic circuits parsed from the ISCAS89 benchmark set (7, 25), where nodes represent logic gates and flip-flops (presented are all five partial scans of forward-logic chips and three digital fractional multipliers in the benchmark set); and World Wide Web hyperlinks between Web pages in a single domain (4) (only three-node motifs are shown). e, multiplied by the power of 10 (e.g., 1.46e6  $= 1.46 \times 10^{6}$ ).

### UNDERSTANDING REGULATORY NETWORKS

#### Genes have "switches"

- Transcription is controlled at a number of levels
- Conceptually, there are proteins (transcription factors) that bind to specific non-coding DNA near transcribed genes
- These proteins can turn *on* transcription or turn *off* transcription
- The transcription factors themselves respond to various cellular signals

### Metabolic networks

- Two types of entities
  - Metabolites: Compounds the cell degrades or synthesizes
  - Enzymes: Proteins that do this
- First work: Complete catalog of all enzymes in yeast and *E. coli*
  - Still don't have effective methods to detect all metabolites







### Human metabolic network

- Two versions
- ~1475 to ~2269 genes coding for enzymes
- 2478-3188 metabolites named
- 1052-3732 reactions

Duarte et al., *PNAS*. **104**:1777 Ma et al.,*Molecular Systems Biology* **3**:135.



Fig. 1. Human metabolic knowledge landscape. Colors represent the percentage of reactions within a pathway that have a confidence score of 3 (biochemical or genetic evidence), 2 (physiological data or evidence from a nonhuman mammalian cell), 1 (modeling evidence), or 0 (unevaluated). Metabolic pathways (primarily defined by the Kyoto Encyclopedia of Genes and Genomes LIGAND database) were classified into three categories based on their level of characterization as detailed in the text.



Figure 3 False prediction percentages for genes in particular cellular compartments (A) and particular metabolic subsystems ( $\theta$ ). The overall error rate is the percentage of false predictions out of all of the predictions. The false-negative (FN) rate is the percentage of FN predictions out of all predictions in which the experimental data show normal growth. The false-positive (FP) rate is the percentage of FP predictions out of all predictions in which the experimental data show retarded growth. Genes that participate in transport functions between compartments are classified according to Table 2. Compartments with at least 10 genes and metabolic subsystems with at least 15 genes are included.