

# Biological Network Analysis



Mark B Gerstein

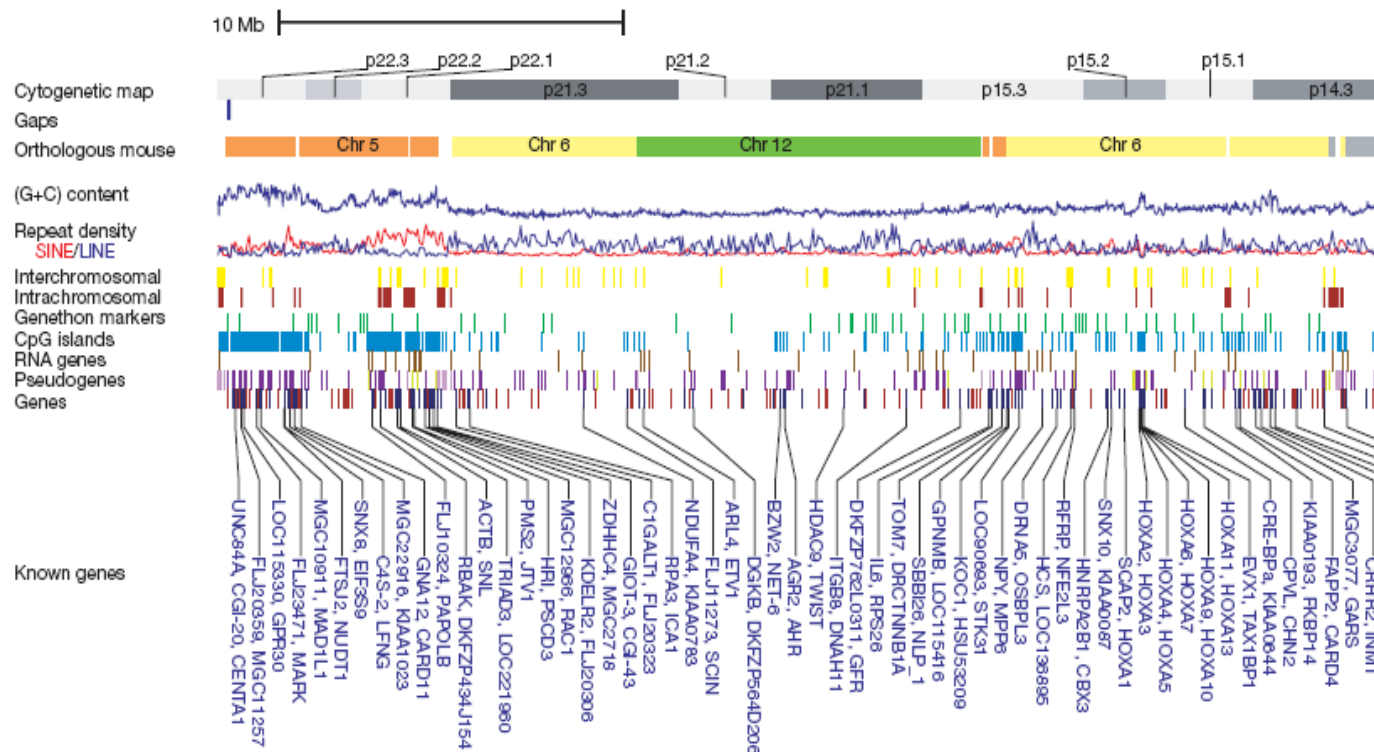
Yale

slides at  
**Lectures.GersteinLab.org**

(See Last Slide for References  
& More Info.)

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

## sequence of human chr. 7



**~1,200 protein-coding genes**  
(~950 pseudogenes)

[Hillier et al, Nature, 424, 157]

# Traditional single molecule way to integrate evidence & describe function

EF2\_YEAST

**Descriptive Name:**  
Elongation Factor 2

**Lots of references**  
to papers

**Summary sentence describing function:**  
This protein promotes the GTP-dependent translocation of the nascent protein chain from the A-site to the P-site of the ribosome.

UniProt

Basic UniProt Protein Viewer - UniProt [the Universal Protein Resource] - Microsoft Internet Explorer

Home > Database > UniProt Protein Viewer

Text Search UniProt Knowledgebase

Home About UniProt Getting Started Searches/Tools Databases Support/Documentation

General information about the UniProt/Swiss-Prot entry	
Entry name	EF2_YEAST
Primary accession number	P32324
Entered in Swiss-Prot	Release 27, 01-OCT-1993
Sequence was last modified	Release 27, 01-OCT-1993
Annotations were last modified	Release 47, 01-MAY-2005

Protein description	
Protein name	Elongation factor 2
Synonyms	EF-2

References	
[1]	NUCLEOTIDE SEQUENCE (EFT1 AND EFT2). MEDLINE=92112760; PubMed=1730643; [NCBI, ExPASy, EBI, Israel, Japan] Perentesis J.P., Phan L.D., Laporte D.C., Livingston D.M., Bodley J.W.; "Saccharomyces cerevisiae elongation factor 2. Genetic cloning, characterization of expression, and G-domain modeling.";

Comments	
FUNCTION	This protein promotes the GTP-dependent translocation of the nascent protein chain from the A-site to the P-site of the ribosome.
SUBCELLULAR LOCATION	Cytoplasmic.

DIR Δ41778 Δ41778

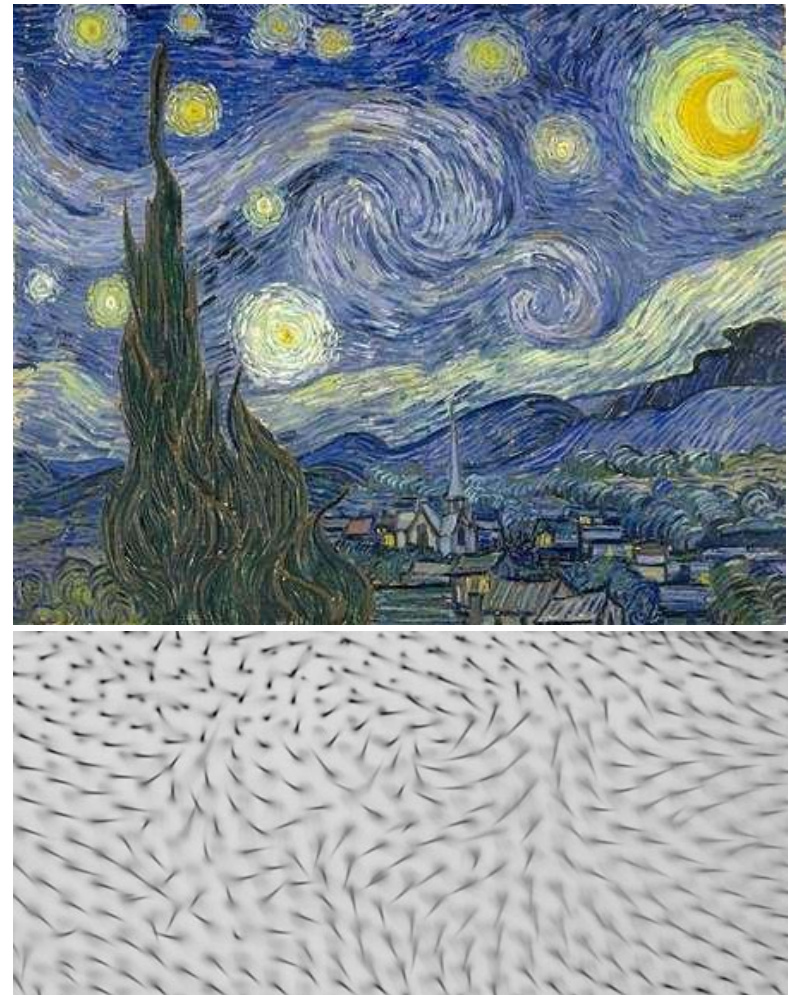
# Some obvious issues in scaling single molecule definition to a genomic scale

- 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 mitochondria, biogenesis of certain subunits..



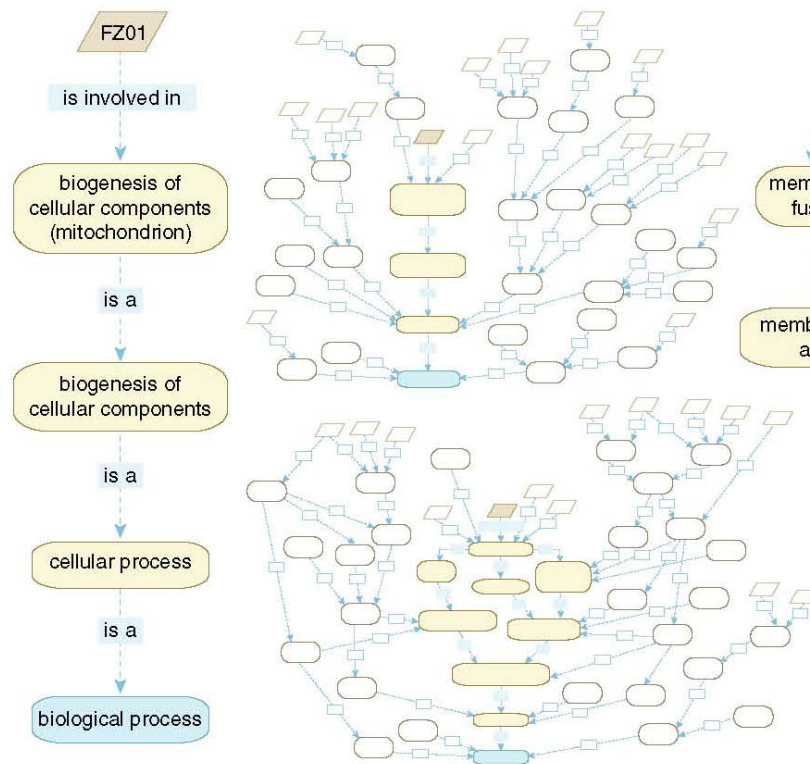
# Some obvious issues in scaling single molecule definition to a genomic scale

- 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 mitochondria, biogenesis of certain subunits..
- Fun terms... but do they scale?....
  - **Starry night** (P Adler, '94)

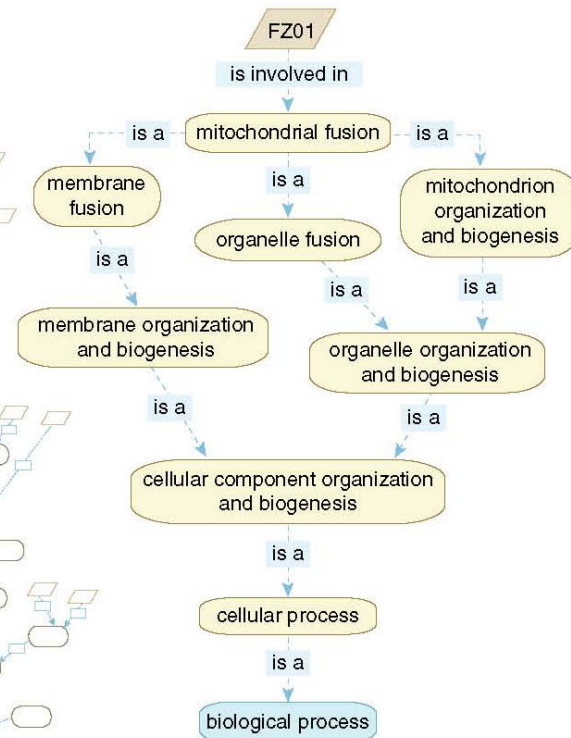


[HSP from Craig et al, Rev Physiol Biochem Pharmacol (2006) 156:1 ; Terms from Seringhaus et al. GenomeBiology (2008)]

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

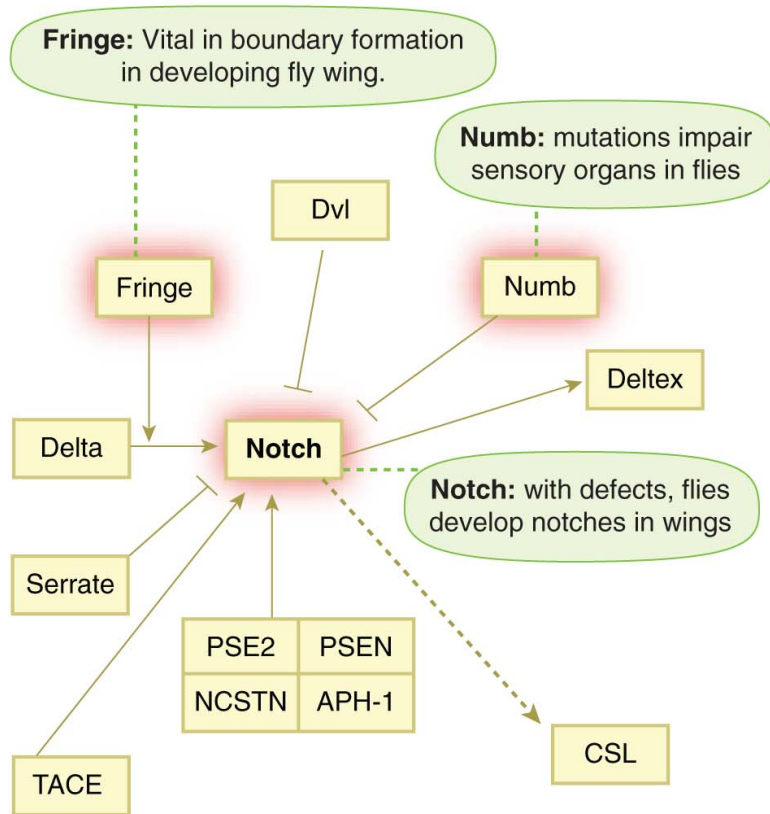


**MIPS (Mewes et al.)**

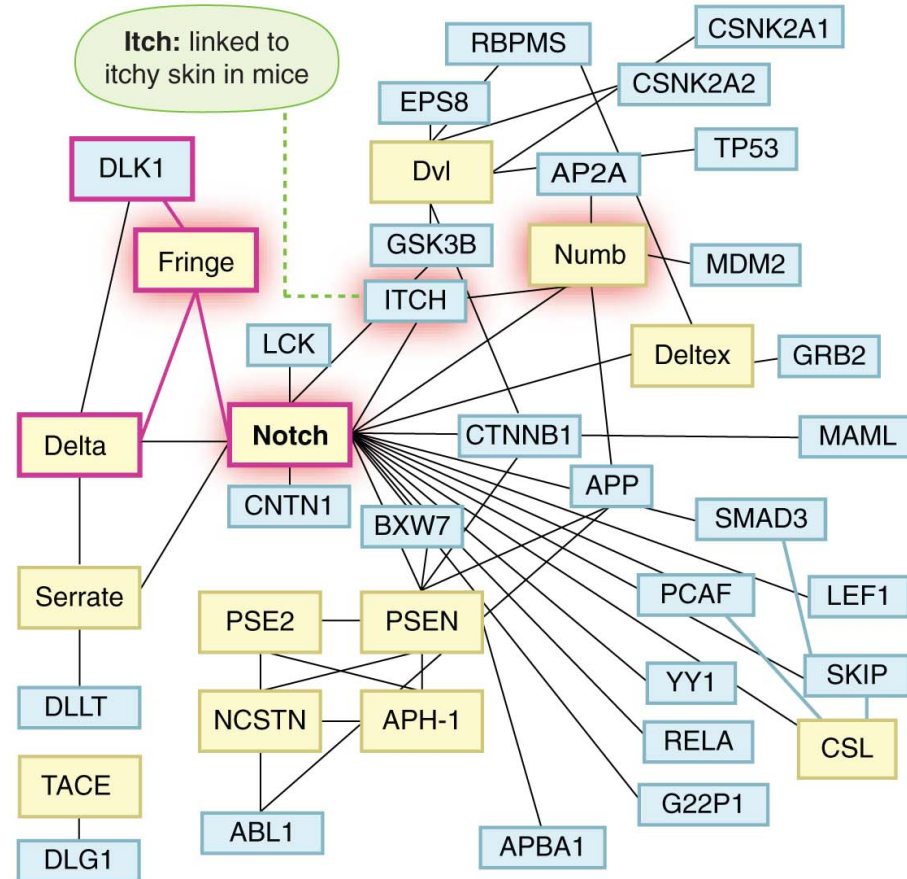


**GO (Ashburner et al.)**

# Networks (Old & New)



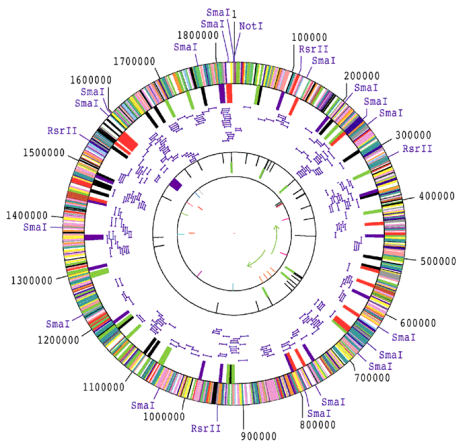
Classical KEGG pathway



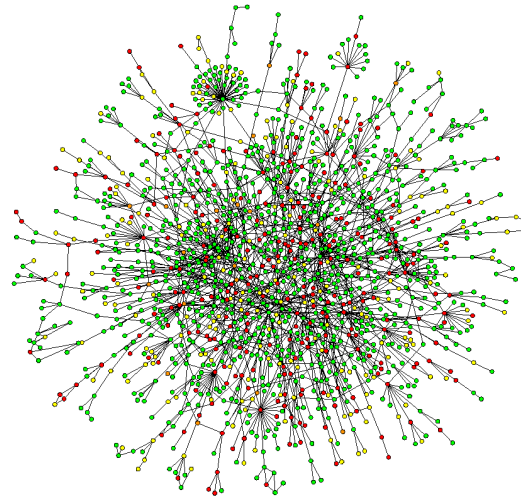
Same Genes in High-throughput Network



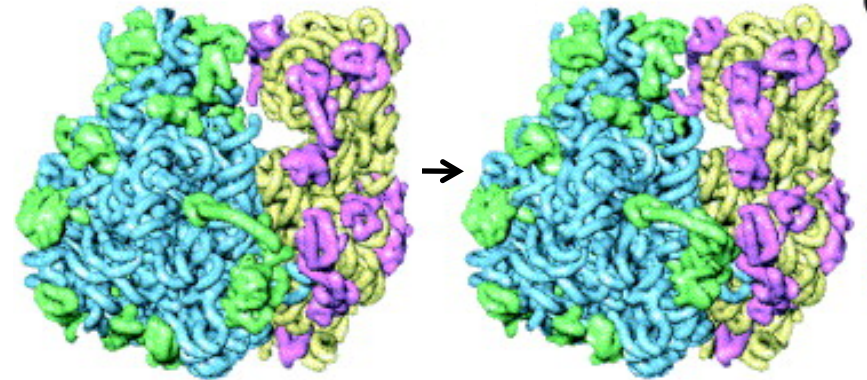
# Networks occupy a midway point in terms of level of understanding



1D: Complete  
Genetic Partslist

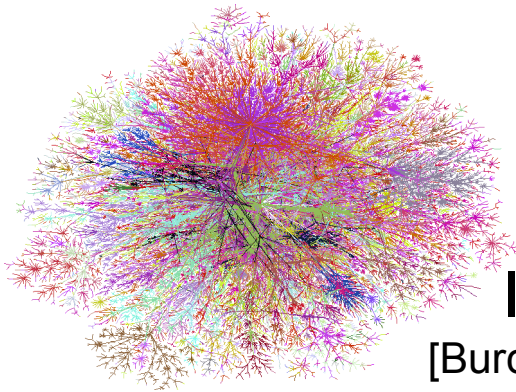


~2D: Bio-molecular  
Network  
Wiring Diagram

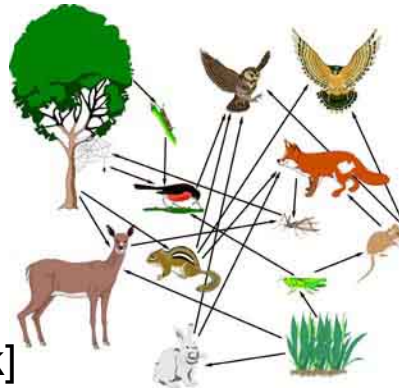


3D and 4D:  
Detailed structural understanding  
of cellular machinery  
(e.g. ribosome in different  
functional states)

# Networks as a universal language



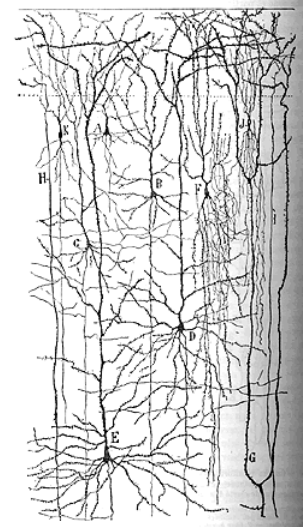
Internet  
[Burch & Cheswick]



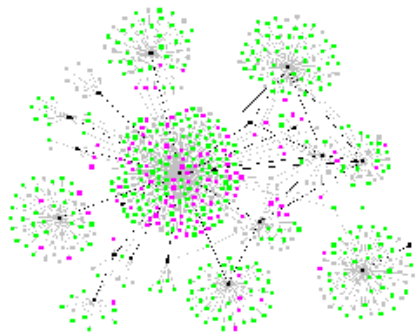
Food Web



Electronic  
Circuit



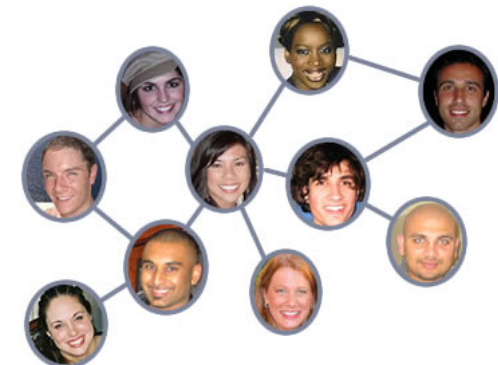
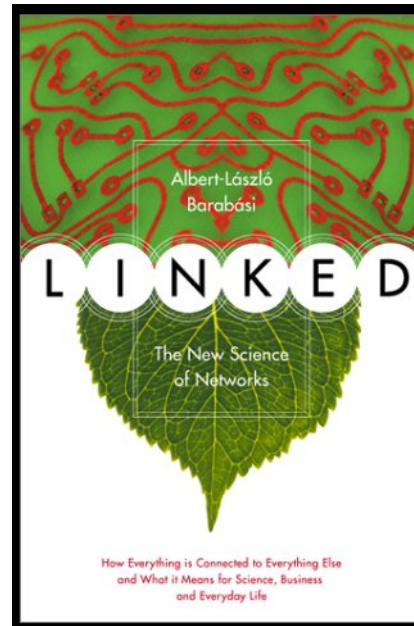
Neural Network  
[Cajal]



Disease  
Spread  
[Krebs]



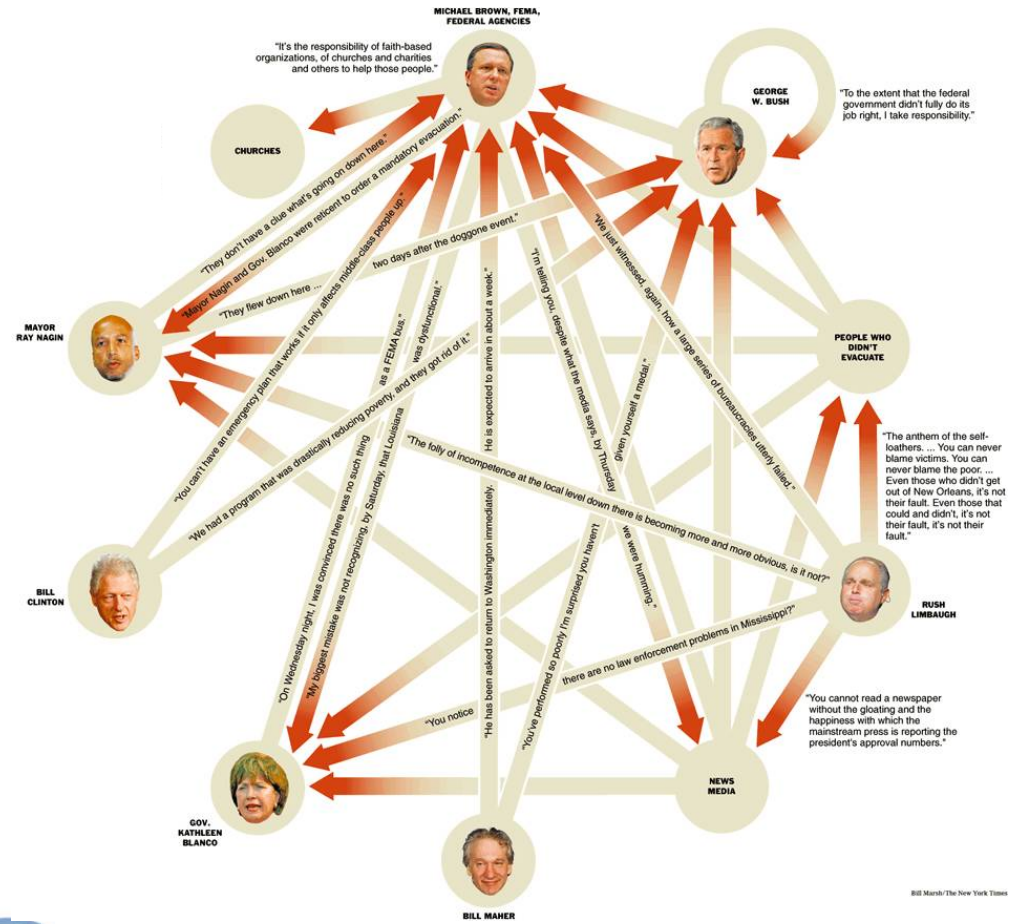
Protein  
Interactions  
[Barabasi]



Social Network



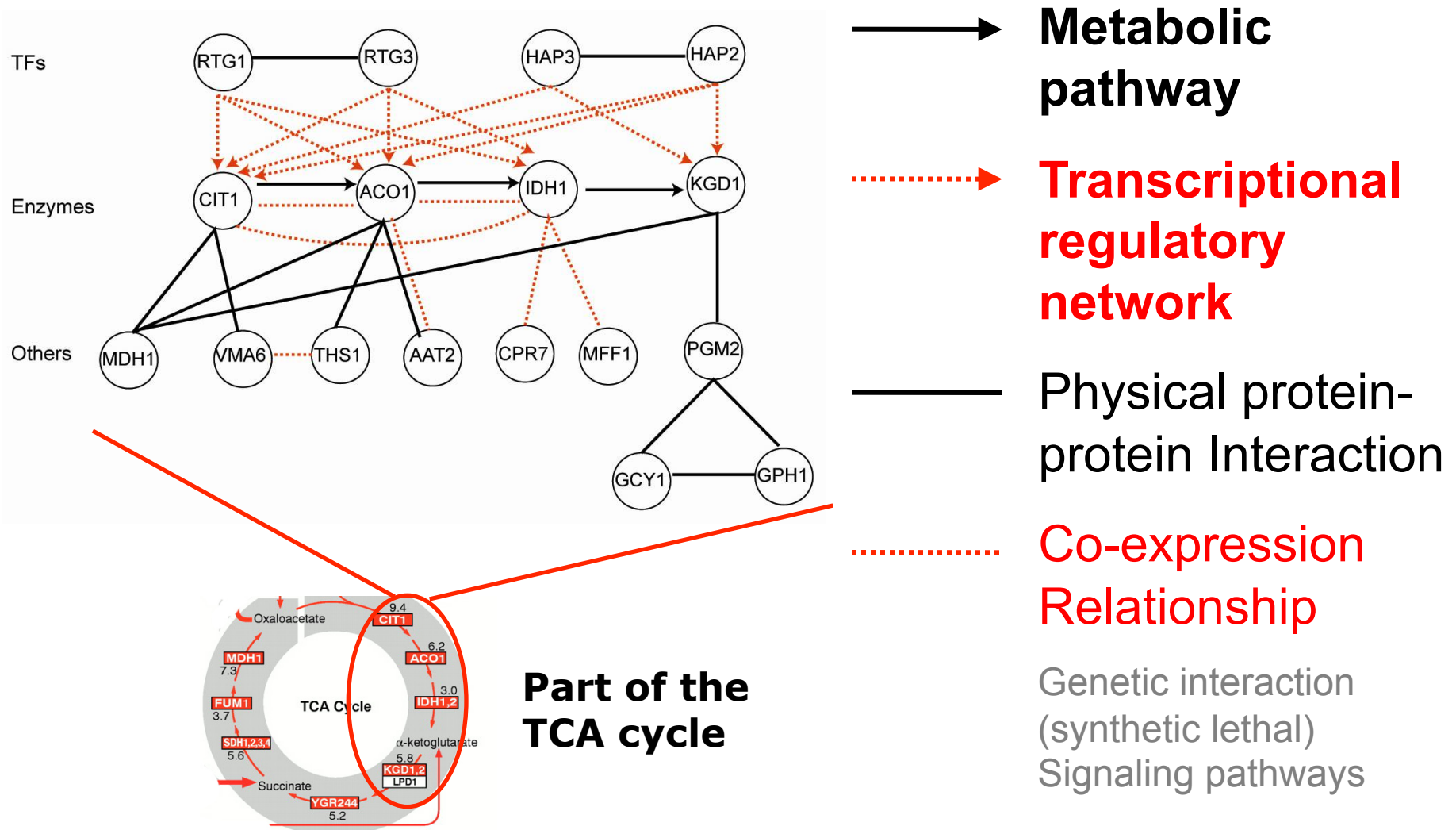
## Guilt by association



## Finding the causal regulator (the "Blame Game")

[NY Times, 2-Oct-05, 9-Dec-08]

# Combining networks forms an ideal way of integrating diverse information



- Why Networks?
- Background:  
Central Network Points
- Networks & Variation  
(human ppi)
- 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)
- Network Dynamics Across Environments  
(prokaryote metab. pathways)
  - Metabolic Pathways
  - Entry pts. (Mem. Proteins)

## Outline: Molecular Networks



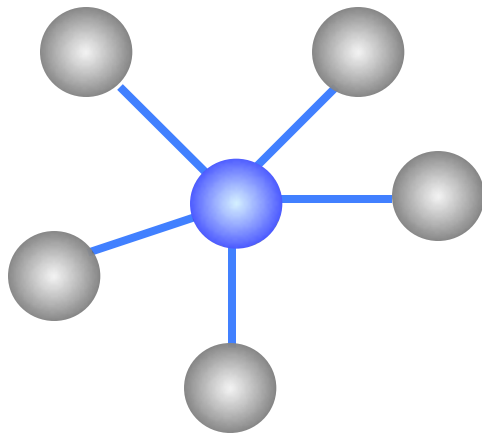


# Background: Finding Central Points in Networks

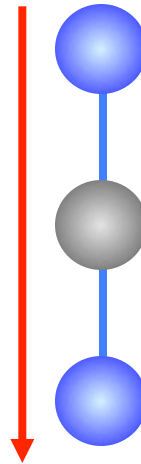


# Global topological measures

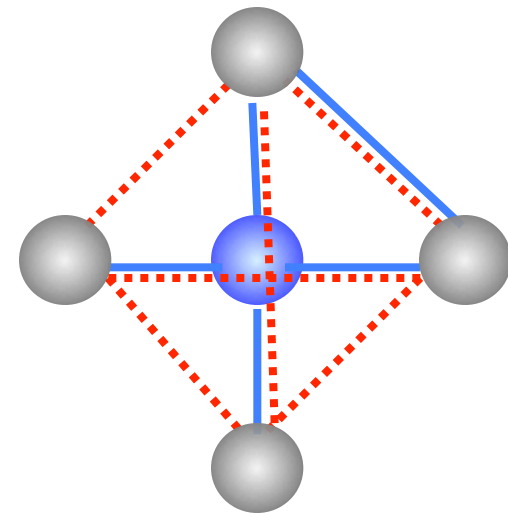
Indicate the gross topological structure of the network



Degree ( $K$ )  
5



Path length ( $L$ )  
2



Clustering coefficient ( $C$ )  
 $1/6$

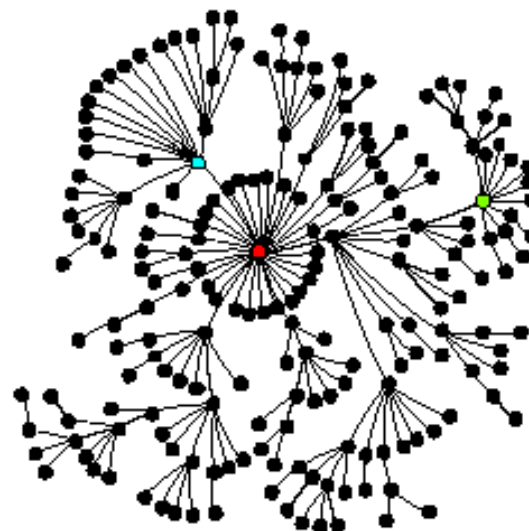
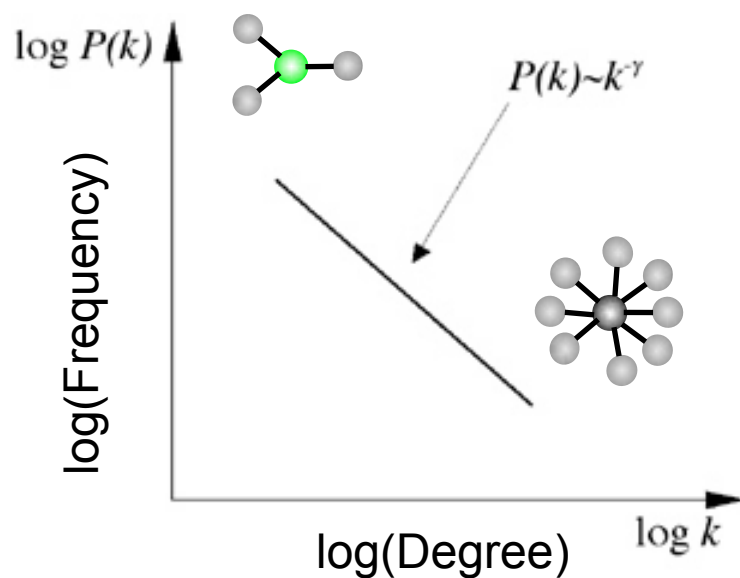
Interaction and expression networks are ***undirected***

[Barabasi]



# Scale-free networks

Power-law distribution



**Hubs** dictate the structure of the network

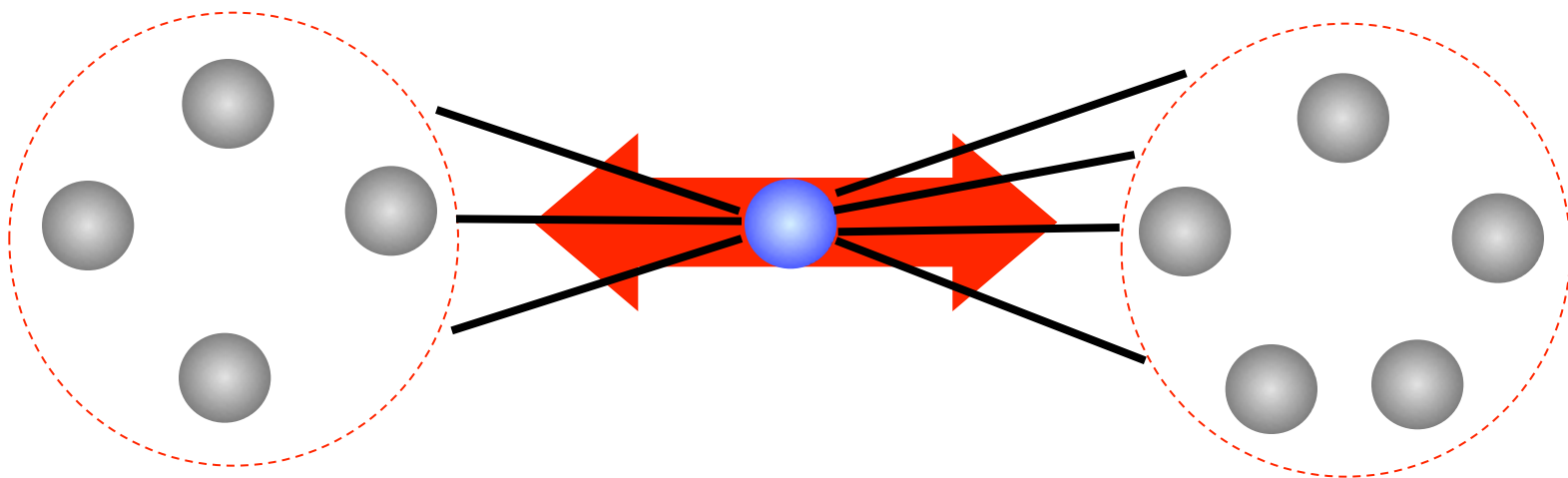
[Barabasi]

## 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.**

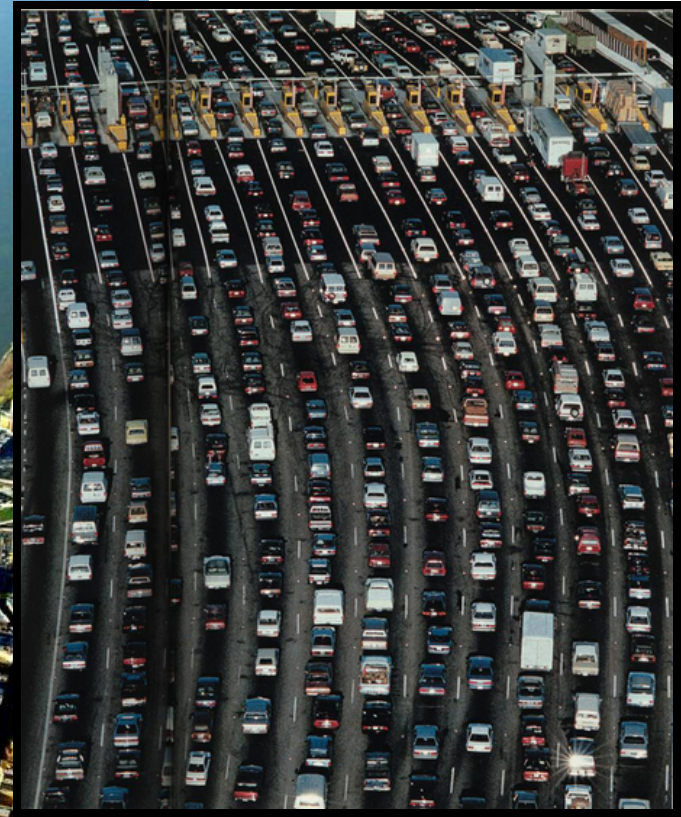


# Betweenness centrality -- Bottlenecks

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



George Washington  
Bridge





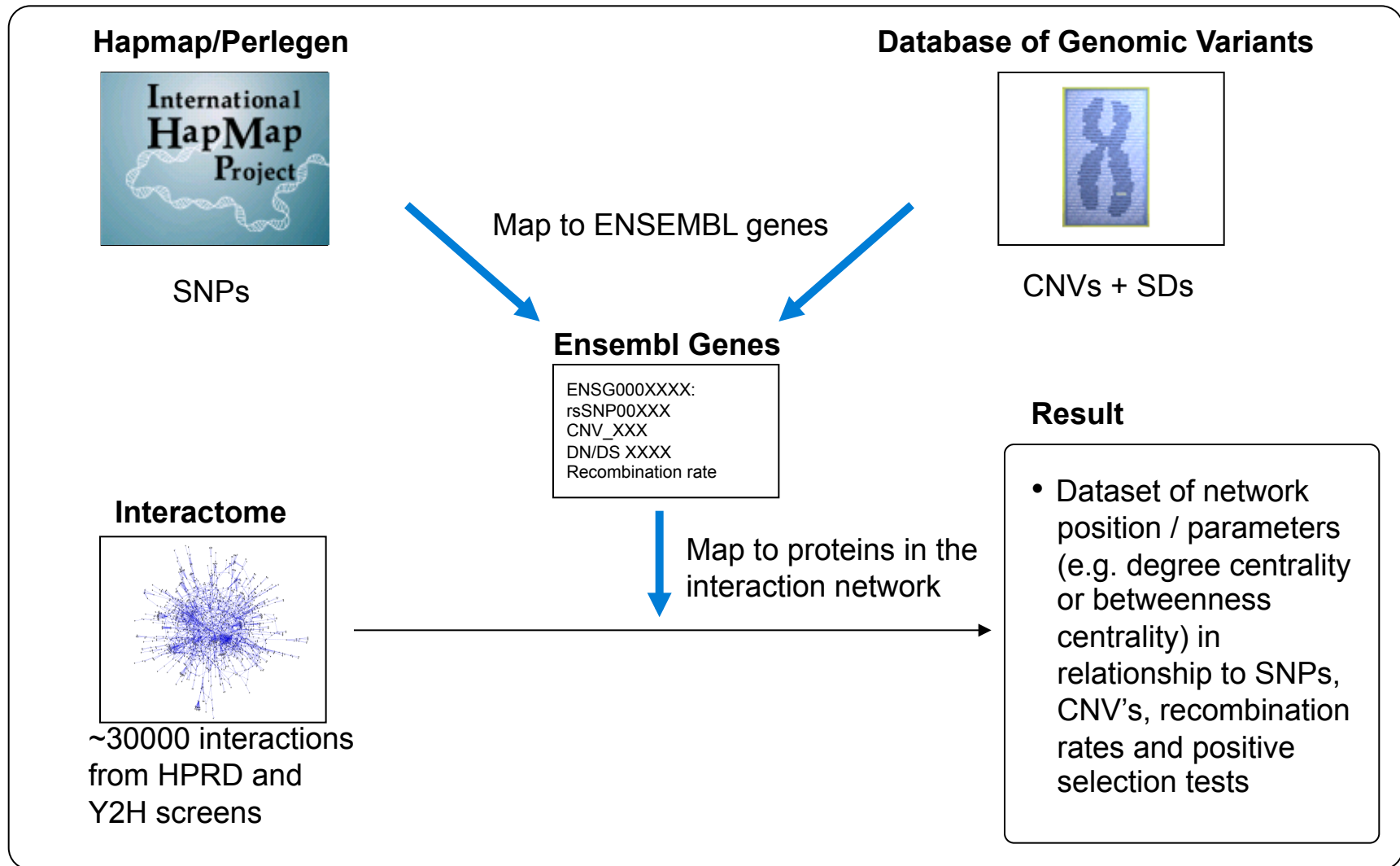
# Networks & Variation

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



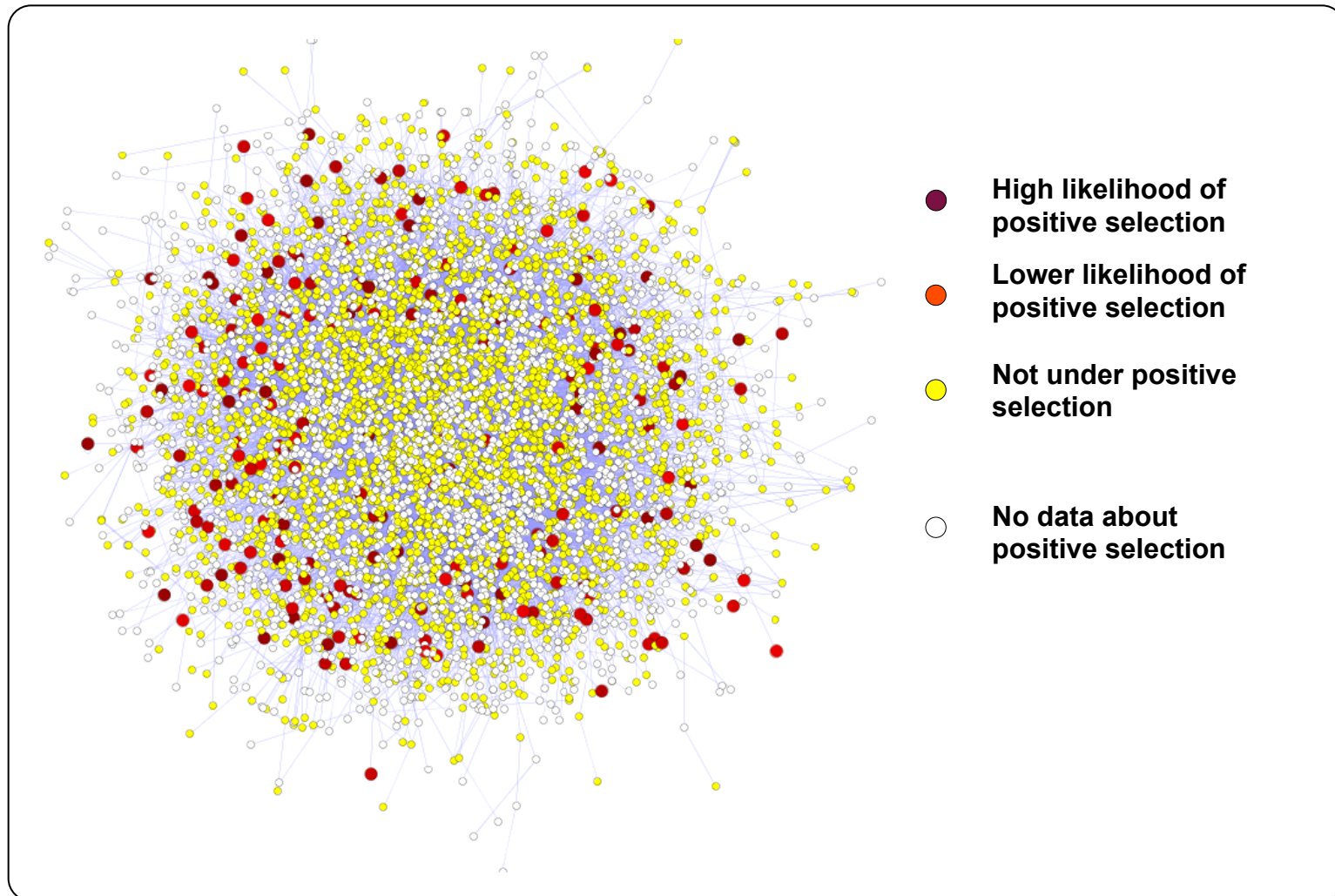
\* From Nielsen et al. *PLoS Biol.* (2005) and Bustamante et al. *Nature* (2005)

Source: PMK



# POSITIVE SELECTION LARGELY TAKES PLACE AT THE NETWORK PERIPHERY

Positive selection in the human interactome

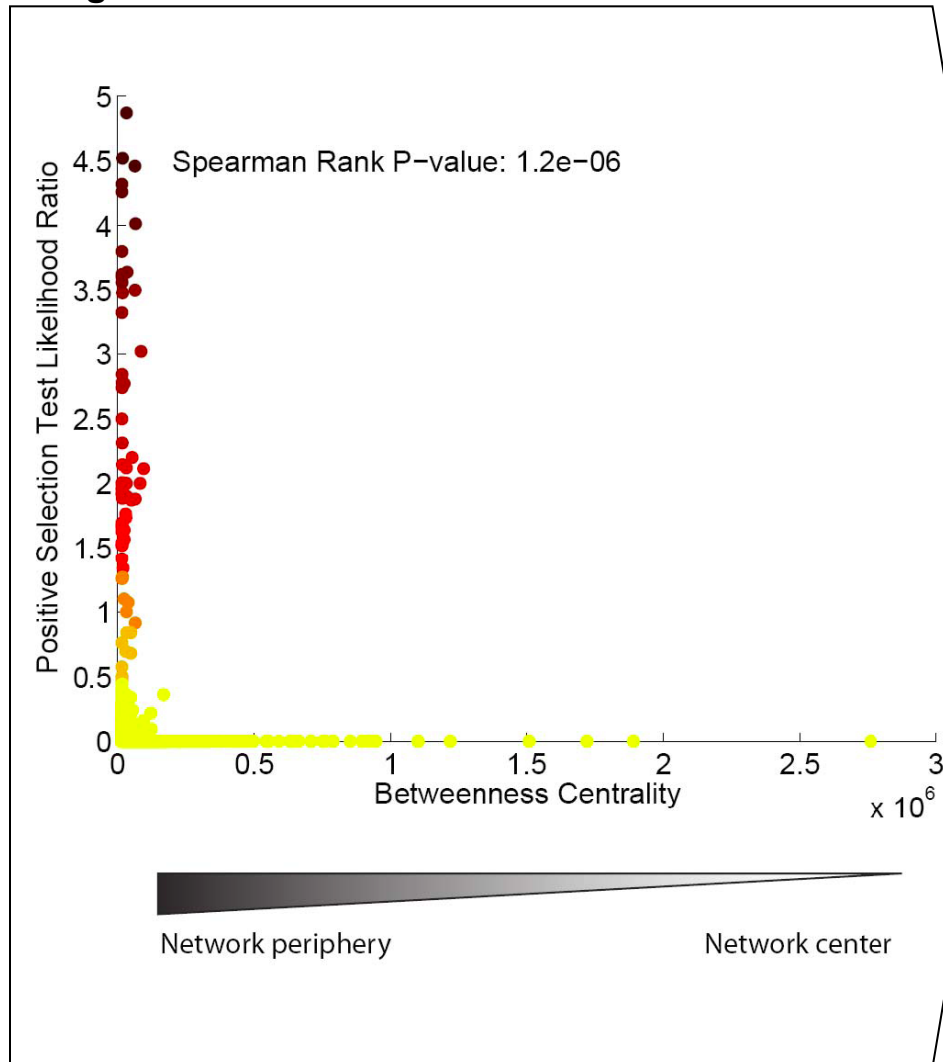


Source: Nielsen et al. *PLoS Biol.* (2005), HPRD, and Kim et al. *PNAS* (2007)

# CENTRAL PROTEINS ARE LESS LIKELY TO BE UNDER POSITIVE SELECTION

▢ Hubs

Degree vs. Positive Selection



## Reasoning

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

# **Social Network Comparison #1**

## **Comparing the Yeast Regulatory Network to a Governmental Hierarchy**

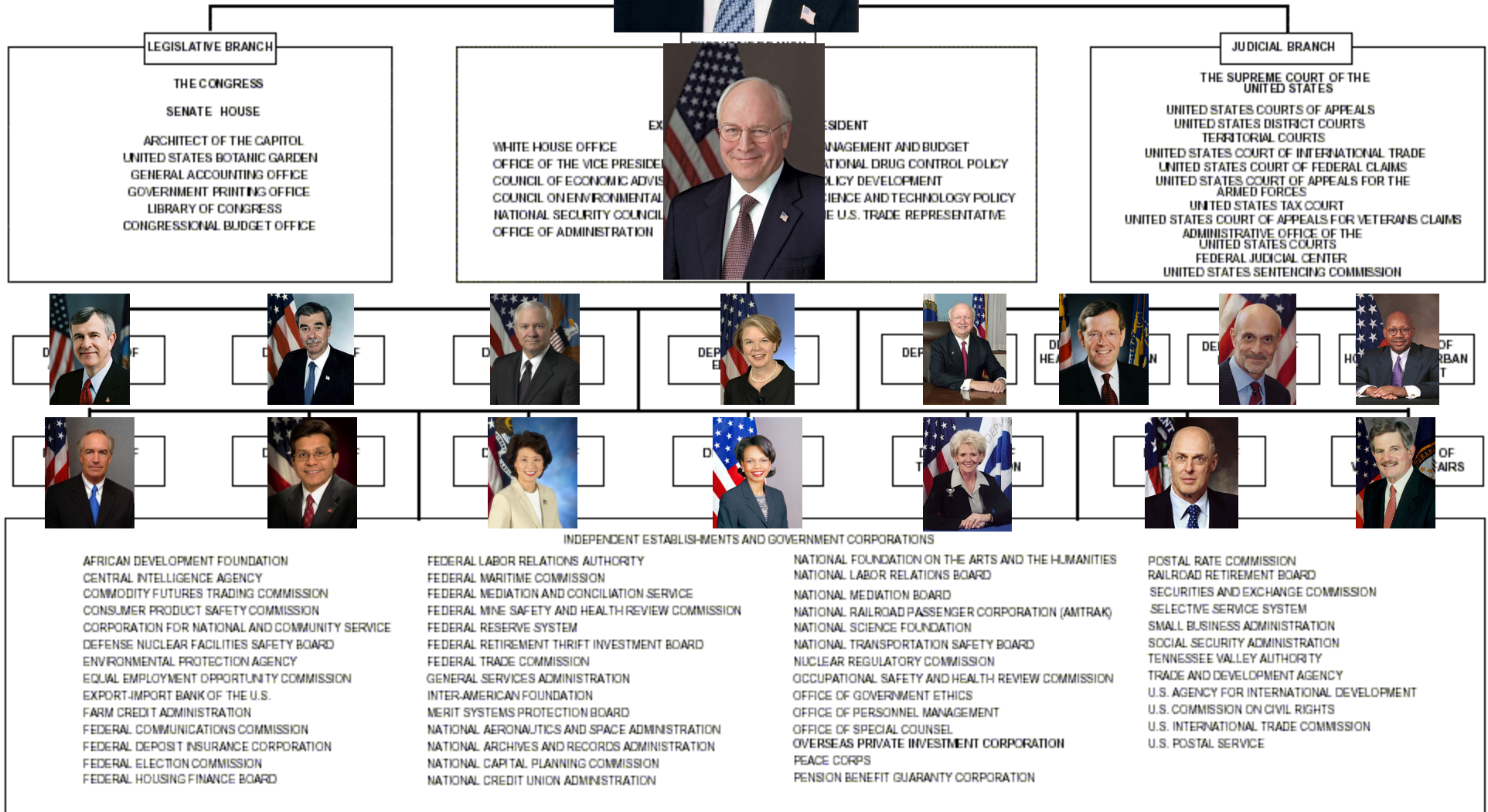
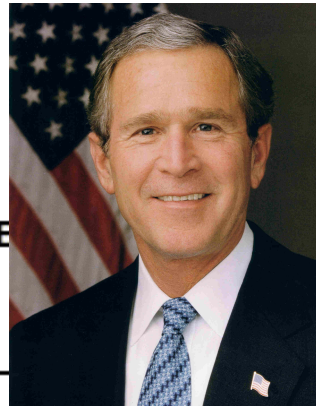




# Social Hierarchy

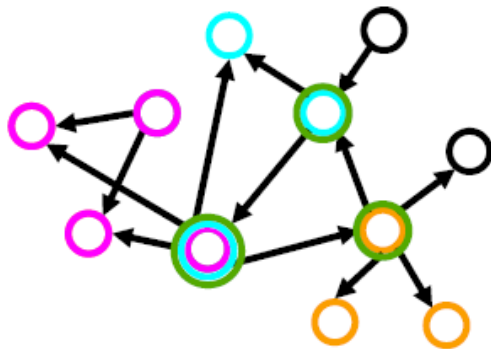
THE GOVERNMENT

UNITED STATES

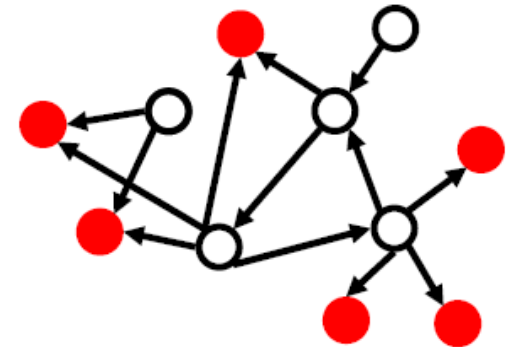


# Determination of "Level" in Regulatory Network Hierarchy with Breadth-first Search

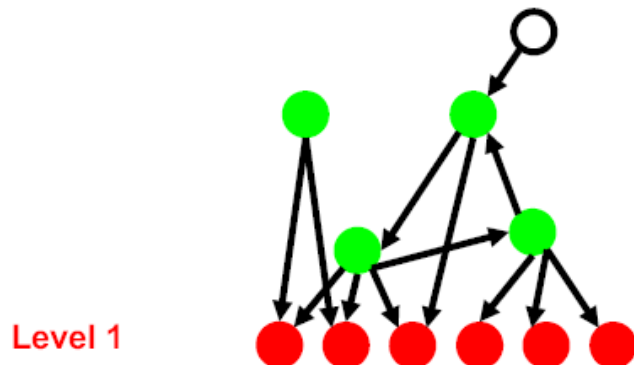
I. Example network with all 4 motifs



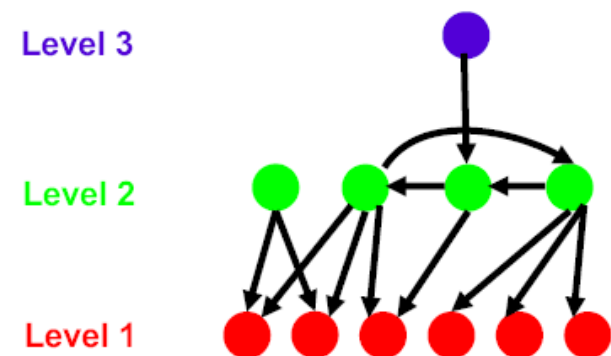
II. Finding terminal nodes (Red)



III. Finding mid-level nodes (Green)



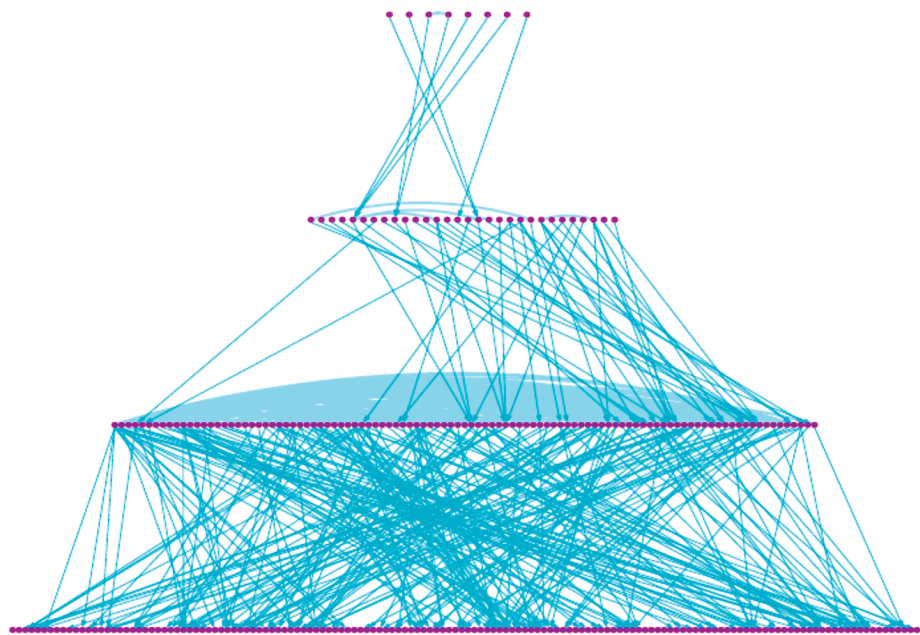
IV. Finding top-most nodes (Blue)



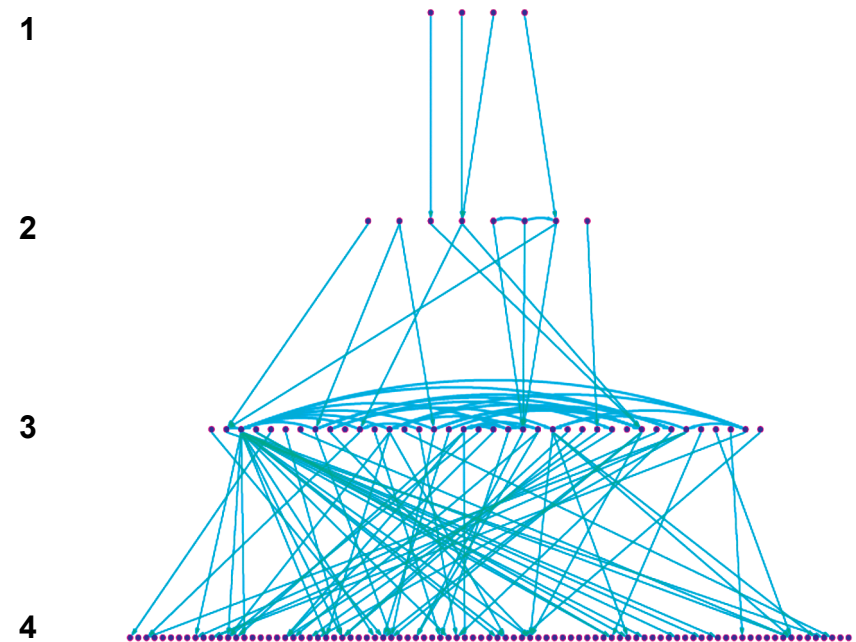
[Yu et al., PNAS (2006)]



# Regulatory Networks have similar hierarchical structures



*S. cerevisiae*

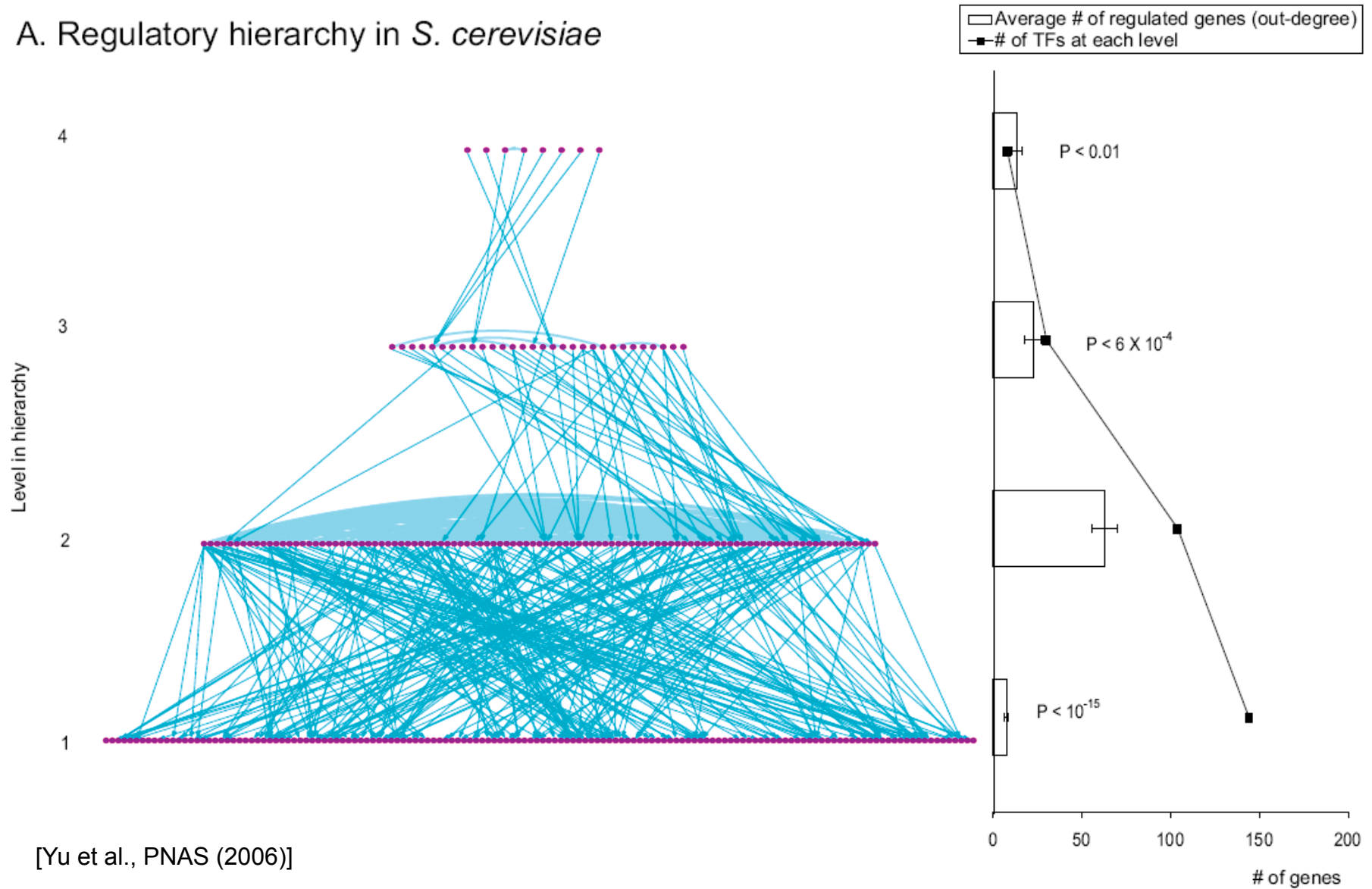


*E. coli*

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

# Yeast Regulatory Hierarchy: the Middle-managers Rule

## A. Regulatory hierarchy in *S. cerevisiae*

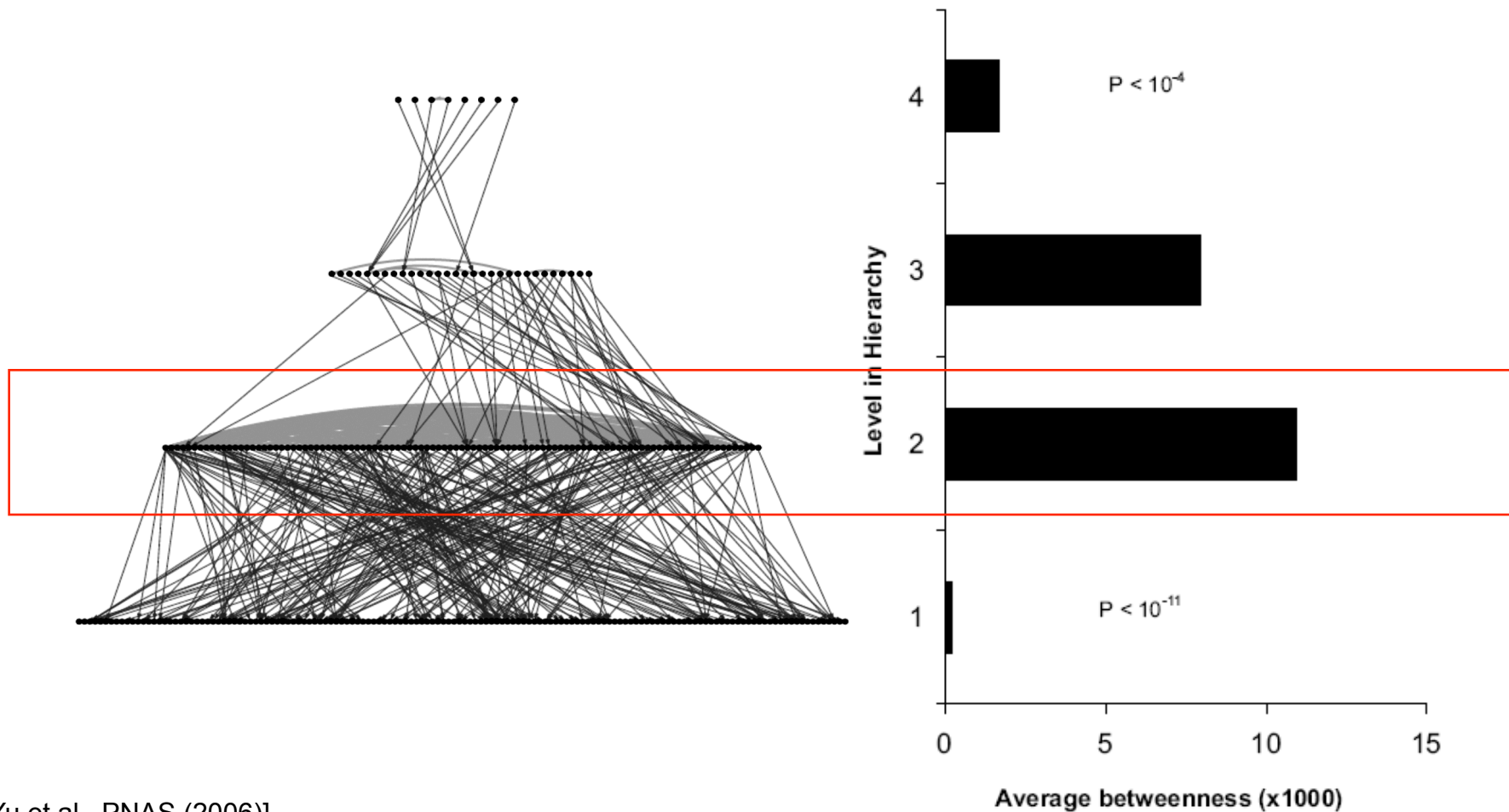


# Yeast Network Similar in Structure to Government Hierarchy with Respect to Middle-managers



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

Average betweenness at each level



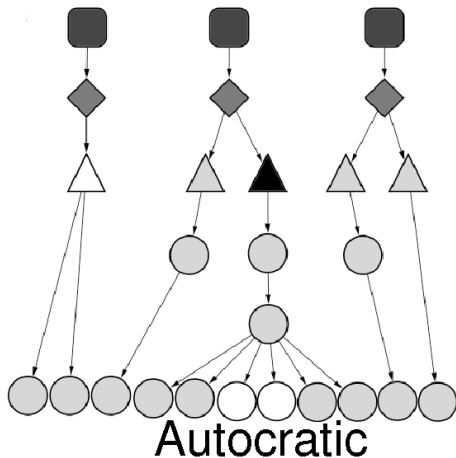
[Yu et al., PNAS (2006)]

# **Social Network Comparison #2**

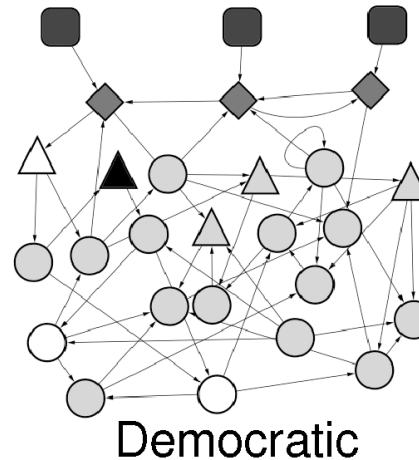
## **Broadening the comparison to different types of hierarchies & different types of biological networks**



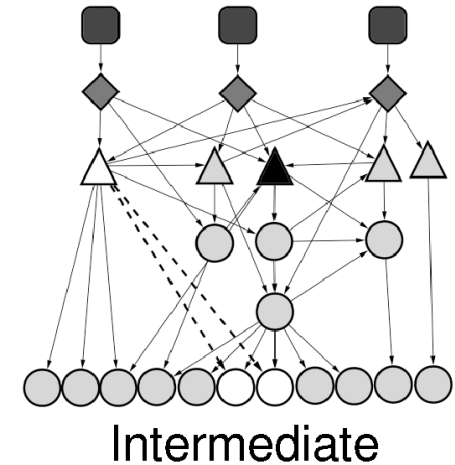
# Different kinds of Hierarchies



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

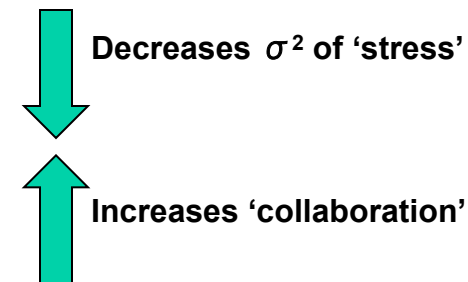


- Without well-defined levels & with more co-regulatory partnerships
- A club or a scientific collaboration network



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

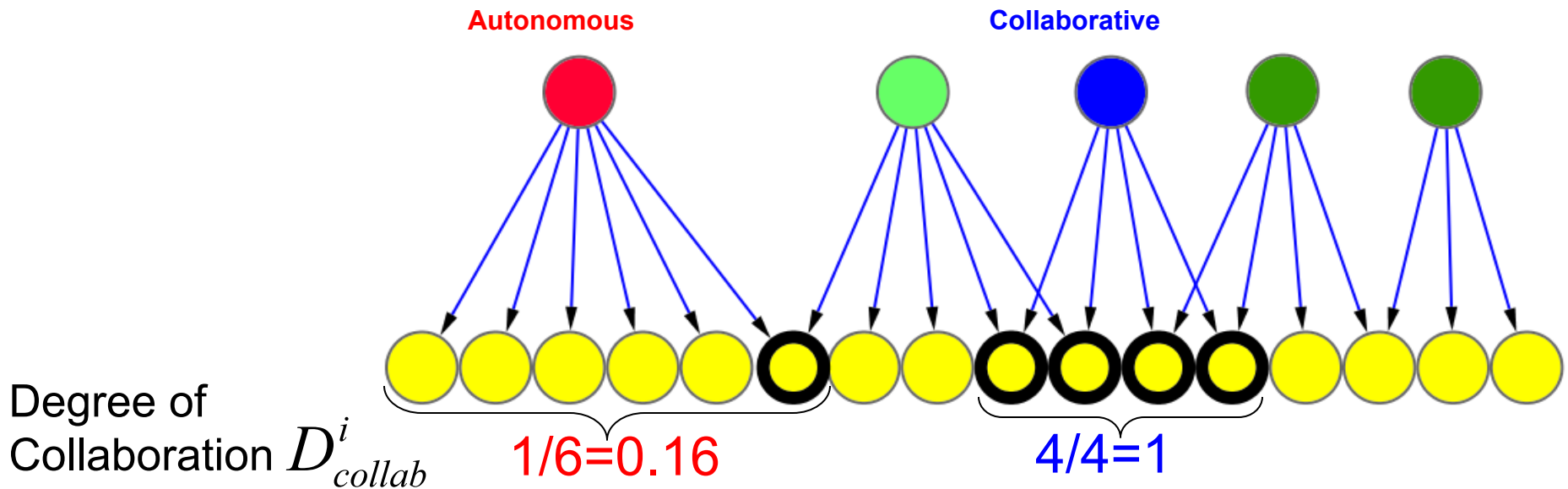
	Autocratic	Democratic	Intermediate
Betweenness ( $\triangle$ )	1.03	3.6	3.3
Betweenness ( $\blacktriangle$ )	4.1	1.08	3.4
Var. Betw. (triangles)	2.1	0.58	1.74
Var. Betw. (all)	2.9	1.4	1.9
$D_{Net-collab}$	0	0.91	0.71



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

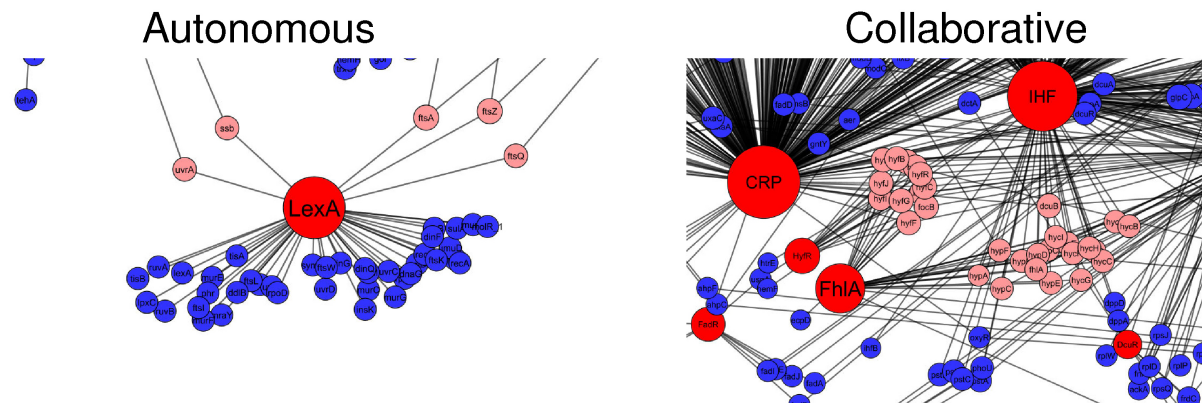


# Collaborative Nature of the Nodes



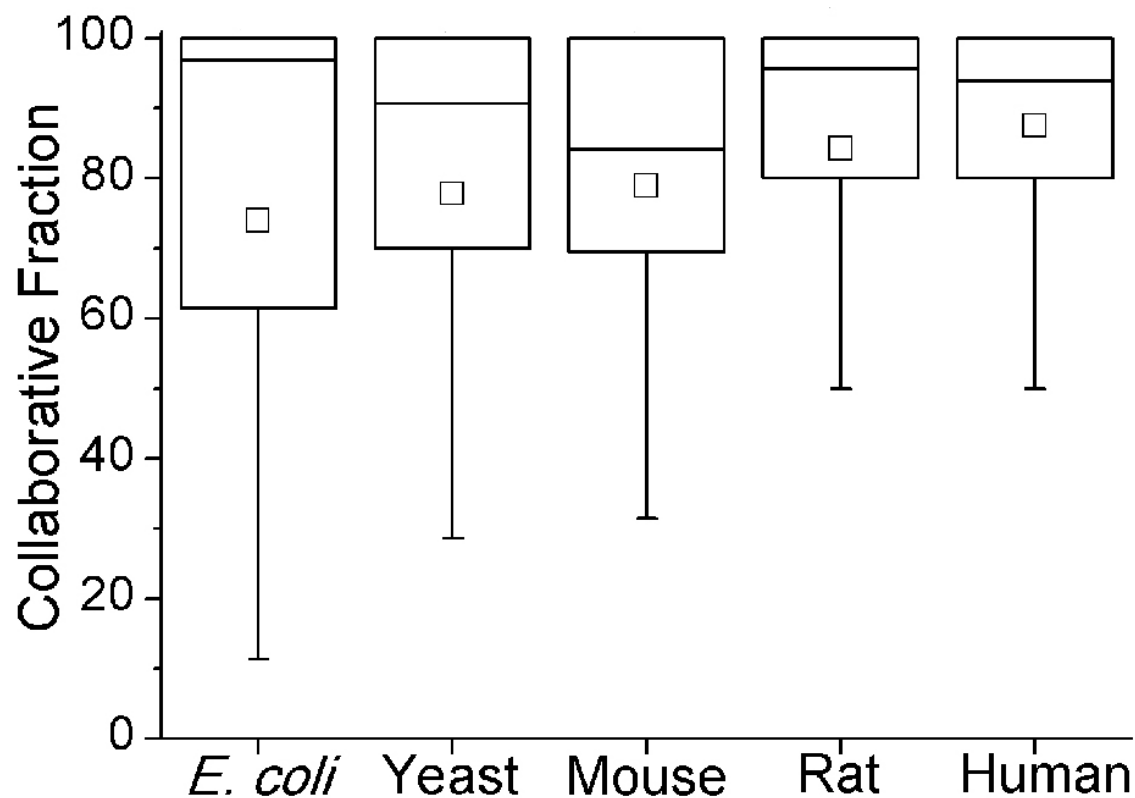
**More Collaborative:  
Democratic**

**More Autonomous:  
Autocratic**



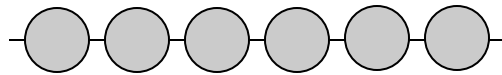
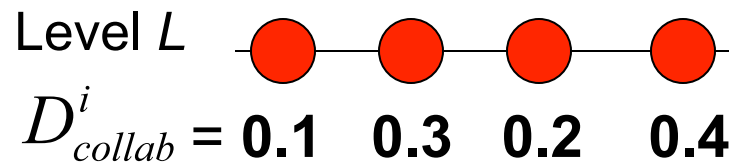
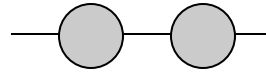
[Bhardwaj et al., PNAS (2010), in press]

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



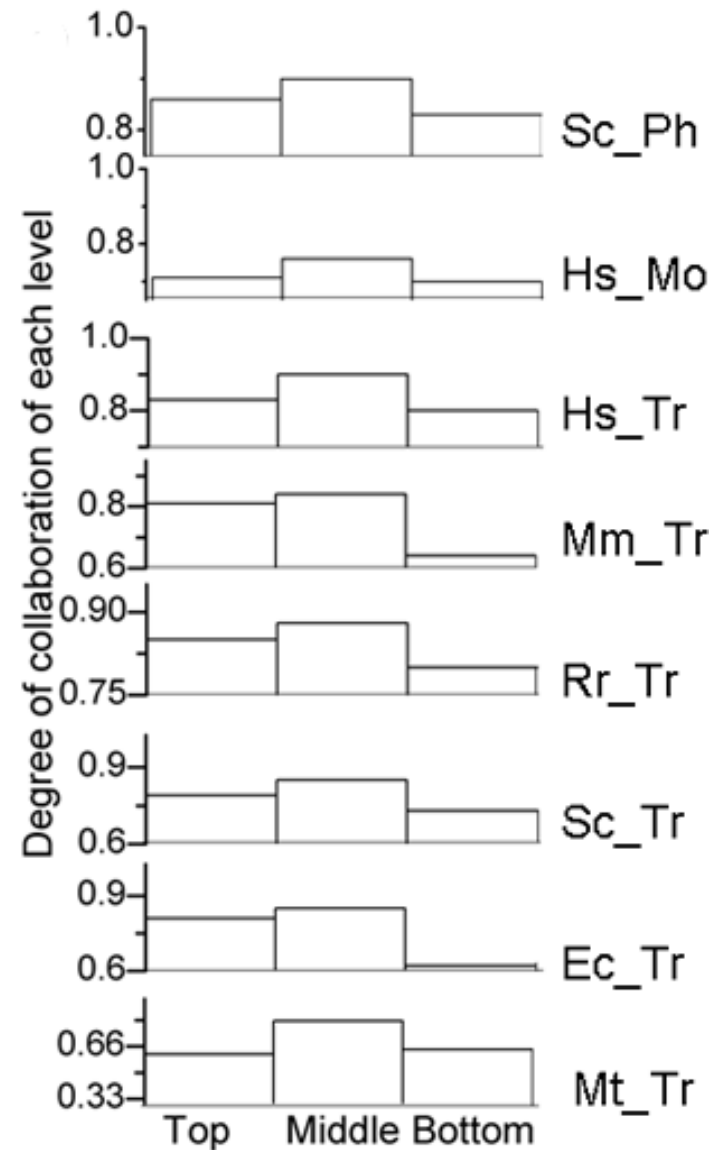
[Bhardwaj et al., PNAS (2010), in press]

# Collaborative Nature of the Levels



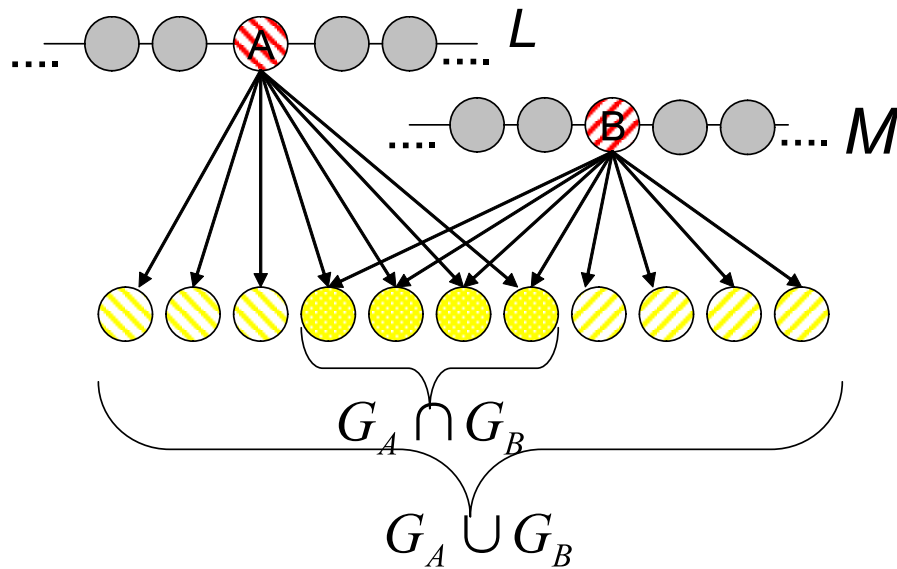
$$D_{Level-collab}^L = \frac{0.1 + 0.3 + 0.2 + 0.4}{4} = 0.25$$

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



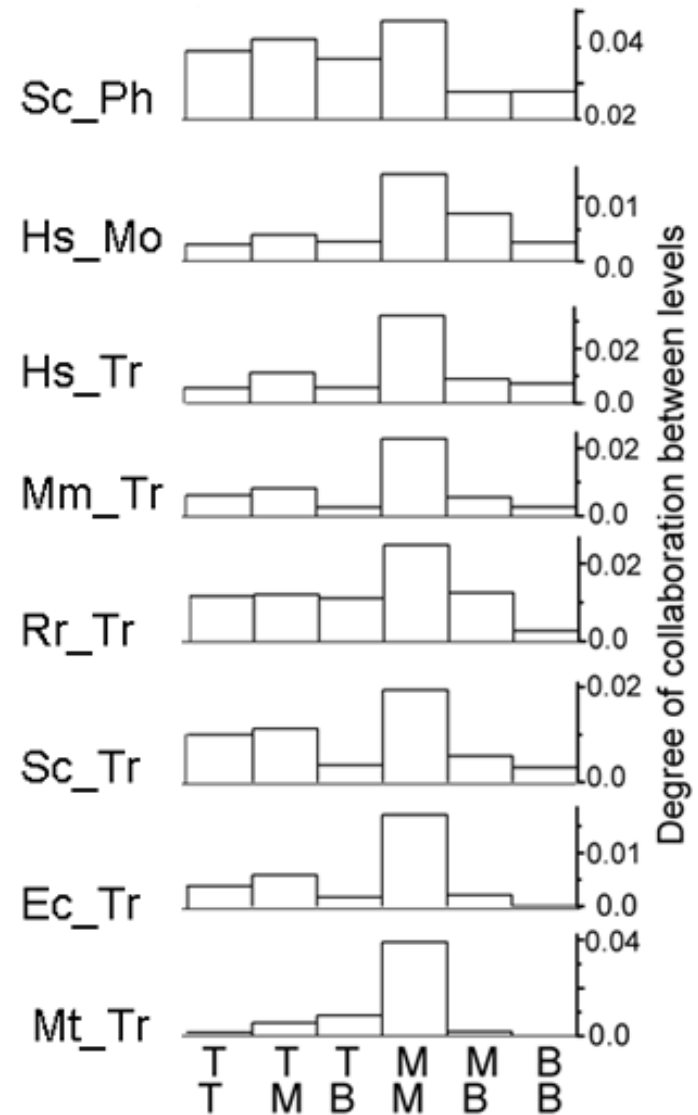


# Collaboration Between Levels



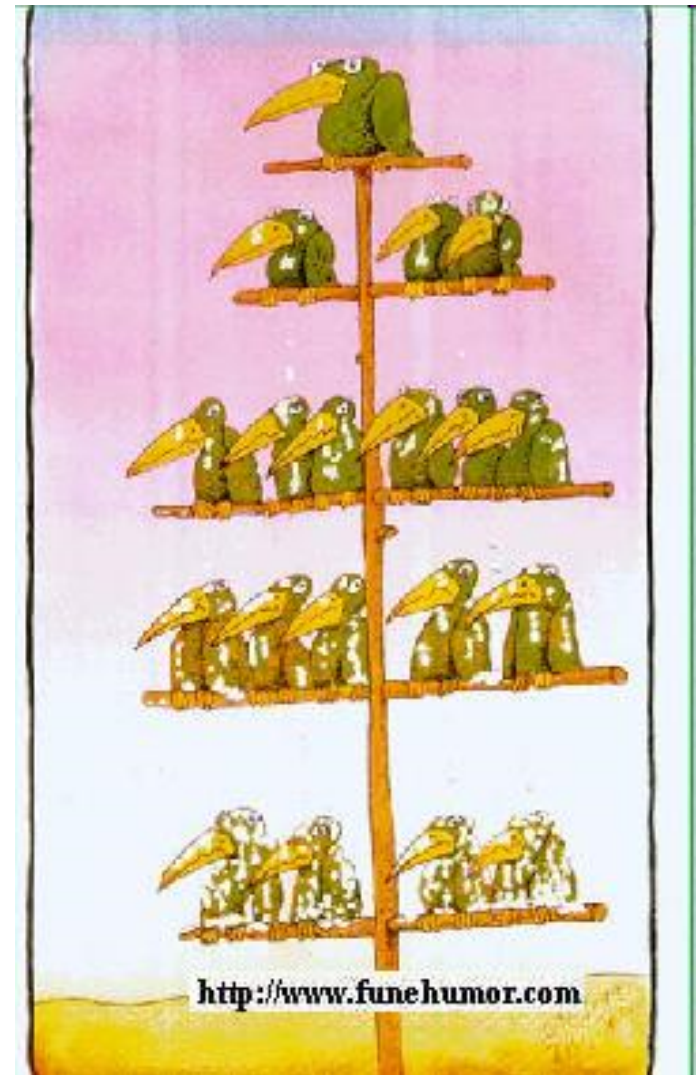
$$D_{\text{betw-level-collab}}^{L,M} = \frac{\sum_{A \in L} \sum_{B \in M} \frac{|G_A \cap G_B|}{|G_A \cup G_B|}}{|L| \cdot |M|}$$

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



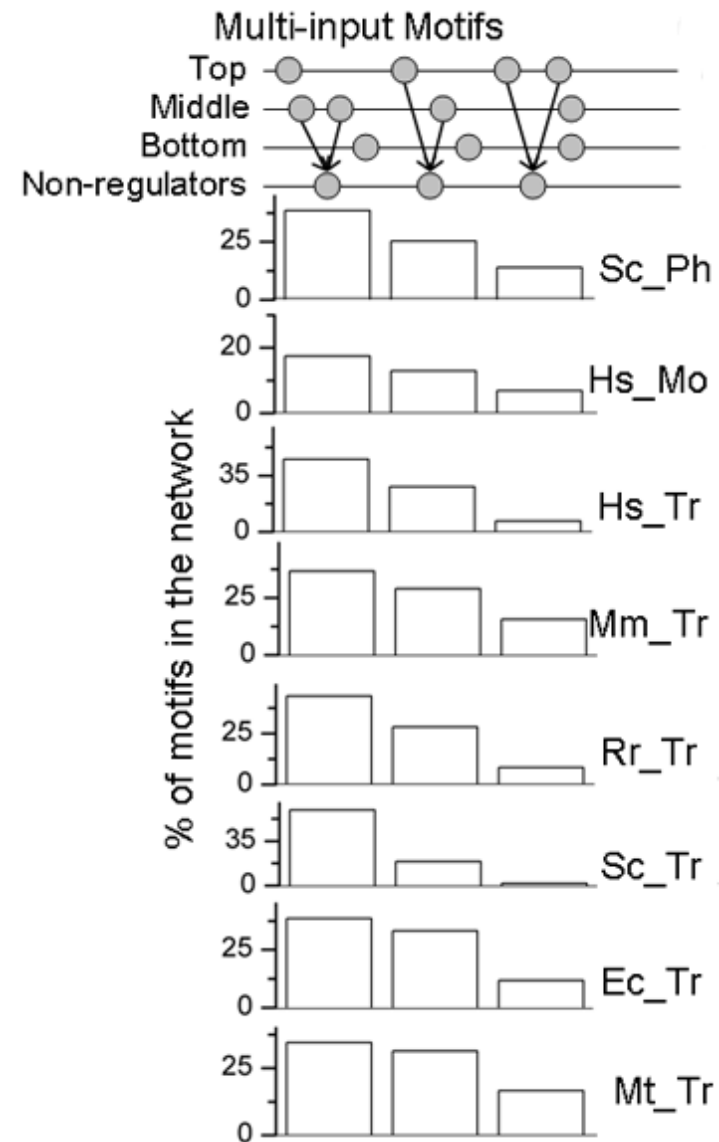
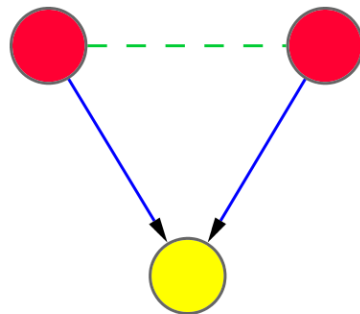
# Middle Managers Interact the Most in Efficient Corporate Settings

- 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.



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

# Co-regulation Instantiates a Multi-Input Motif



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

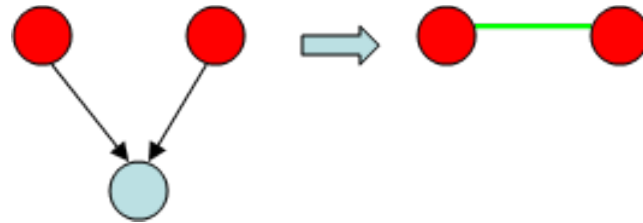


# **Network Comparisons #3**

## **Relating the size of co-regulation in partnership networks with the scale of the regulated**



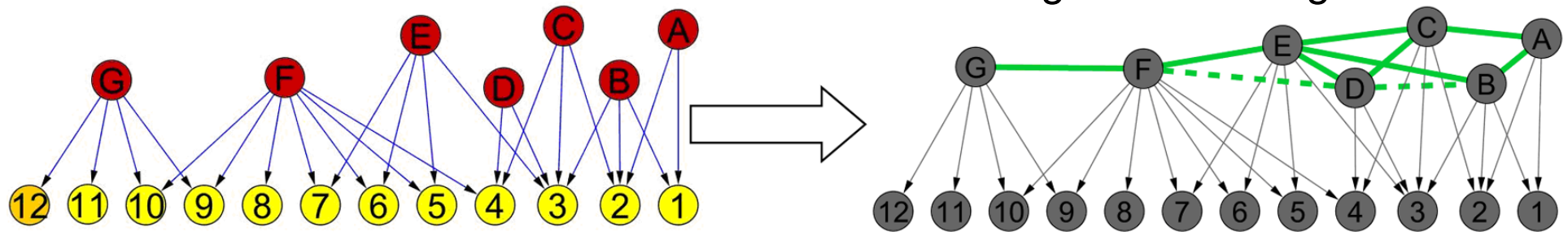
## Co-regulation Partnerships



- 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.

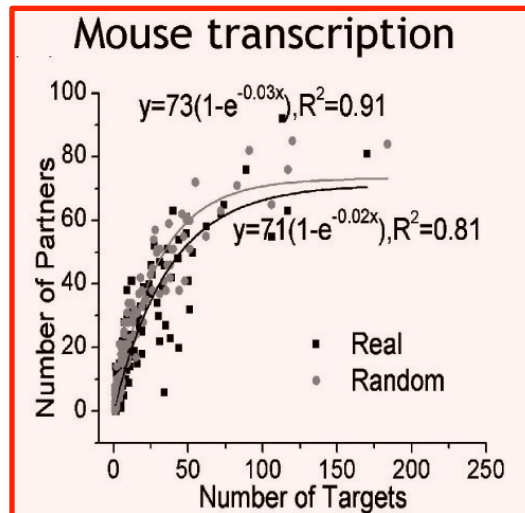
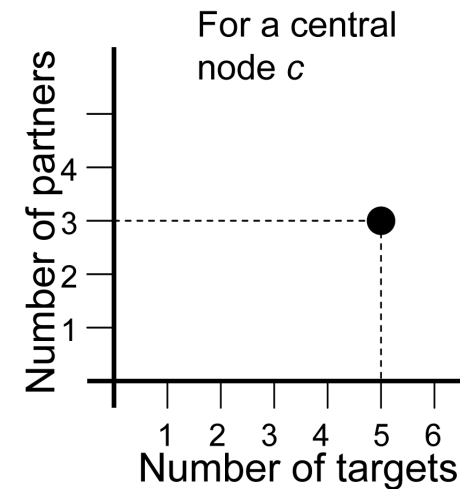
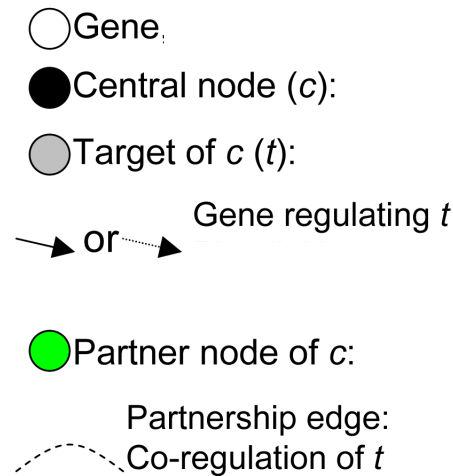
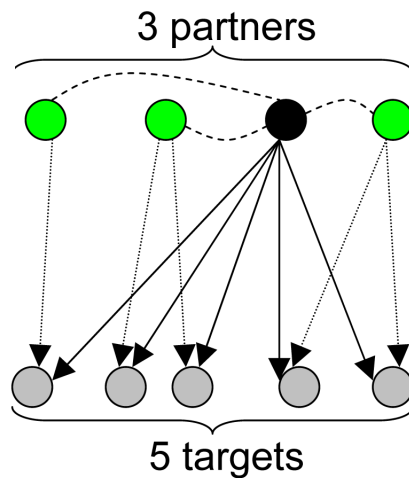
[Bhardwaj et al., PLoS Comp Biol (2010), in press]

# Building and Analysis of Networks



Network type	Species	Number of regulators	Number of targets	Number of interactions
Transcription	<i>E. coli</i>	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

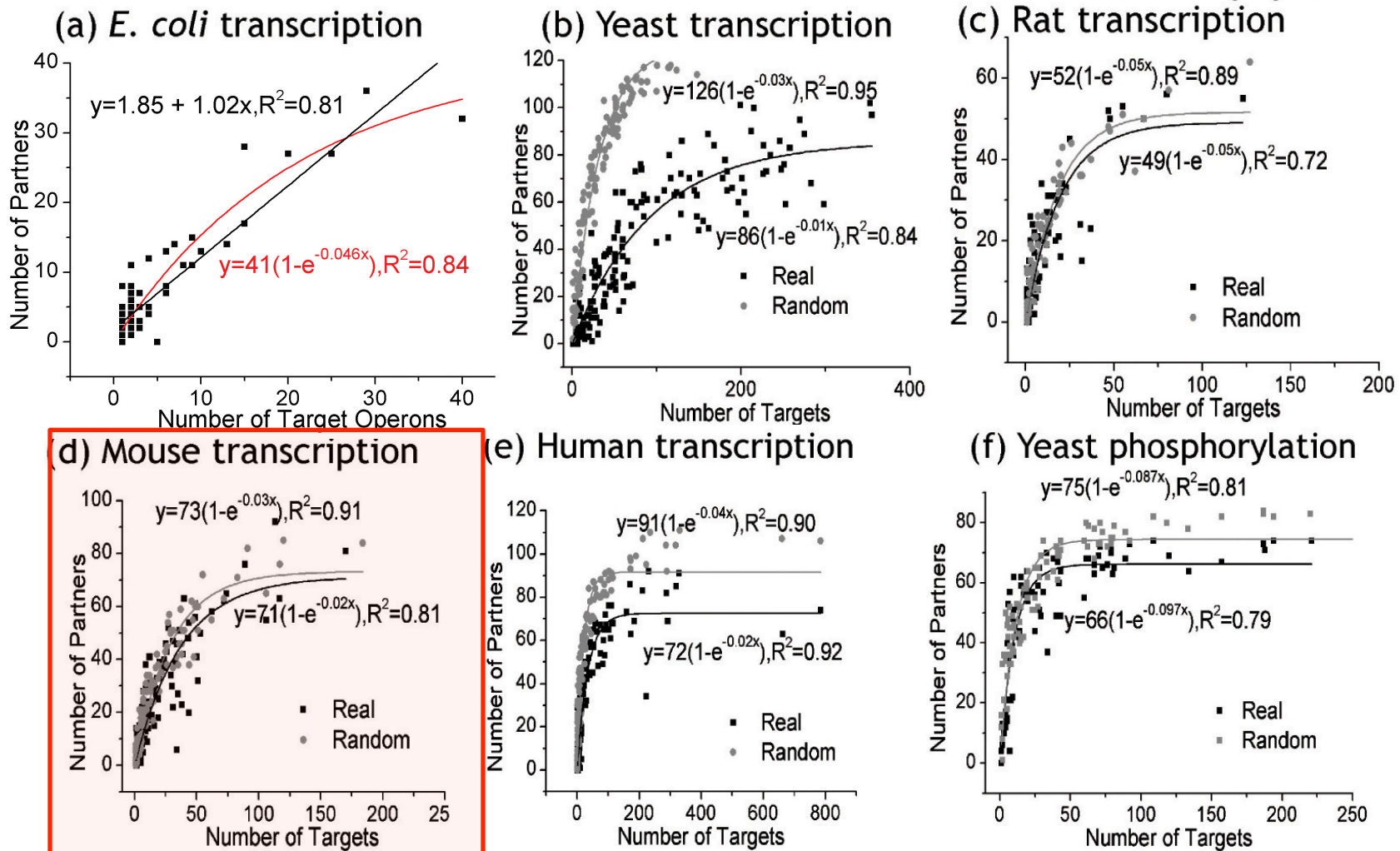
[Bhardwaj et al., PLoS Comp Biol (2010), in press]



## Scaling of Regulatory Partnerships with Targets



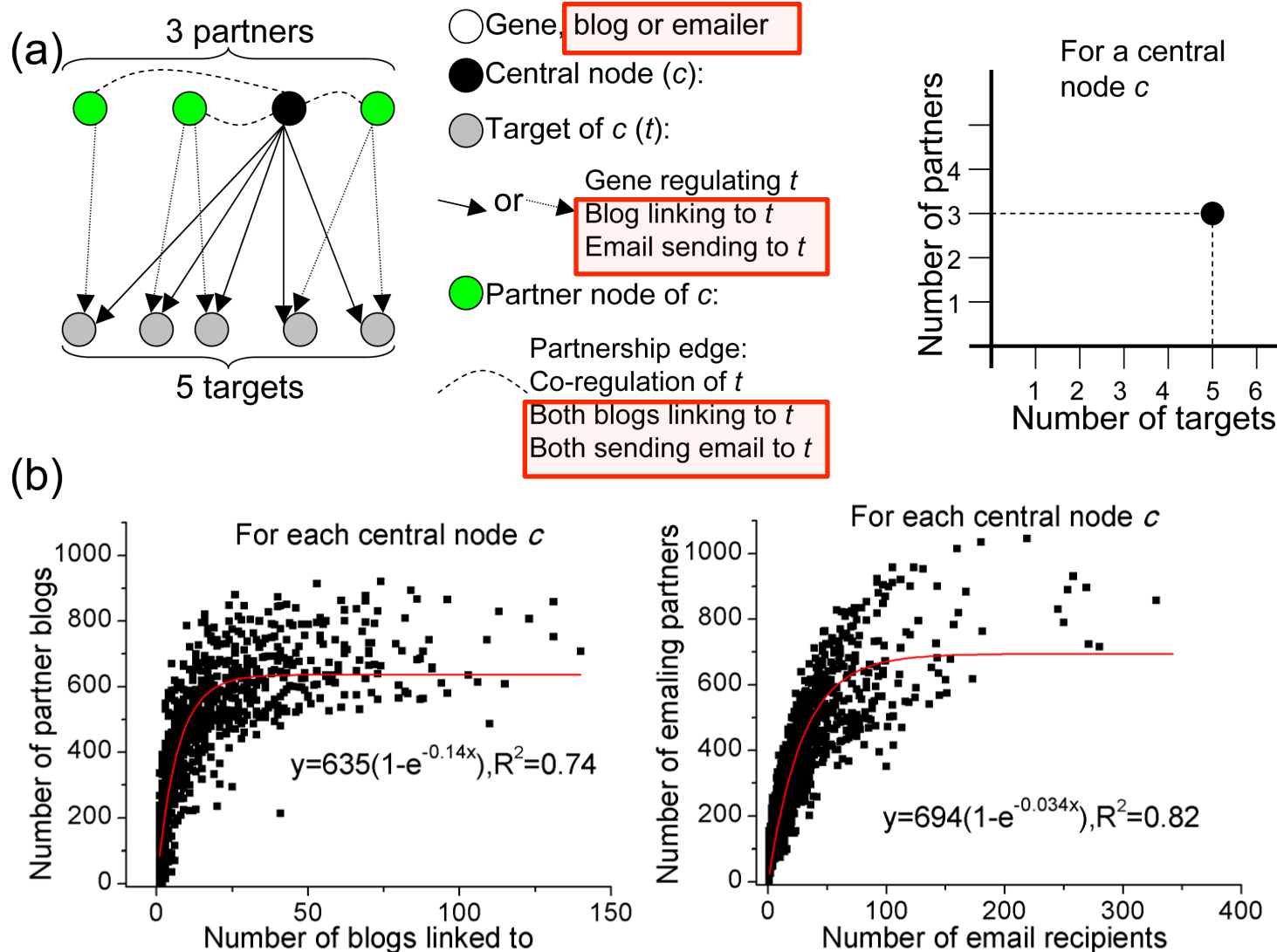
# Scaling across many networks



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

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

# Comparison to Social Networks: Partnership networks effectively saturate with increasingly complex output

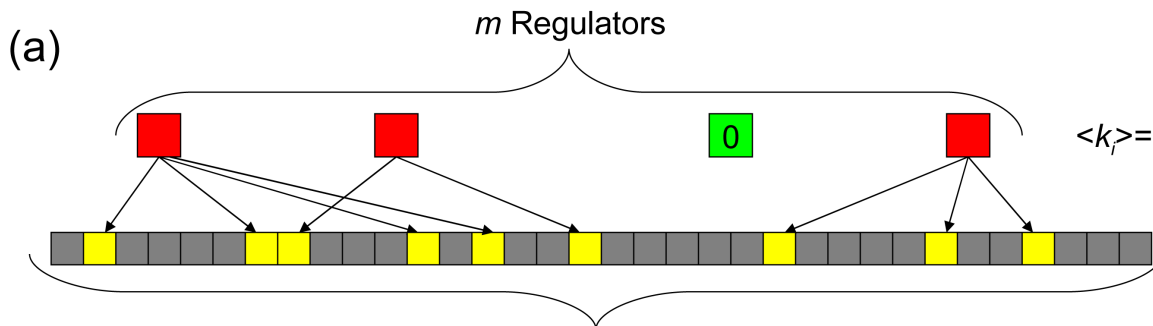


# A Simple Theoretical Model

On average, each regulator has  $n$  targets

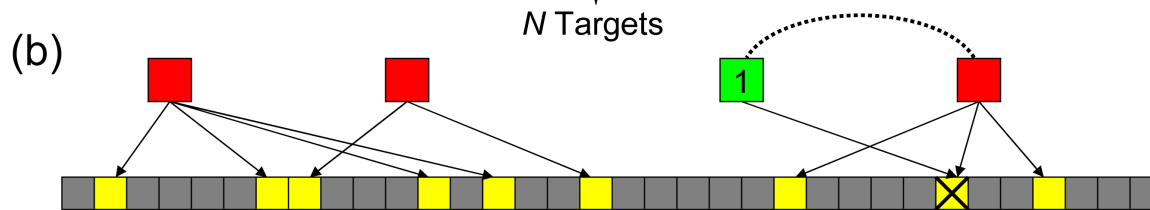
$f_i$  : the number of partners

$k_i$  : the number of targets



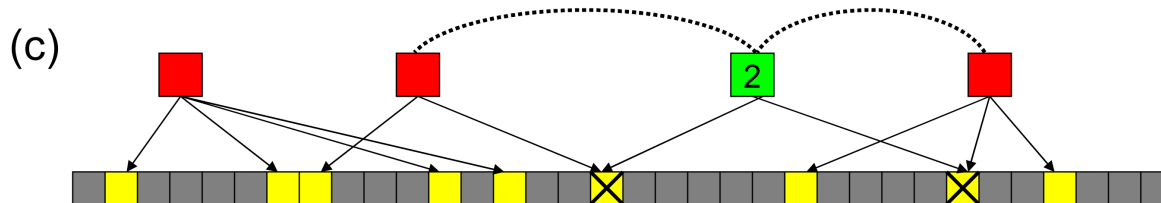
When  $f=0$ ,

$$\left. \frac{\Delta f}{\Delta k} \right|_{f=0} = \frac{nm}{N}$$



When  $f=1$ ,

$$\left. \frac{\Delta f}{\Delta k} \right|_{f=1} = \frac{(m-1)n}{N}$$



Generalizing,

$$\frac{\partial f}{\partial k} = \frac{(m-f)n}{N}$$

Integrating this, we get:

$$f = m(1 - e^{-n/Nk}) = a(1 - e^{-bk})$$

43

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

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

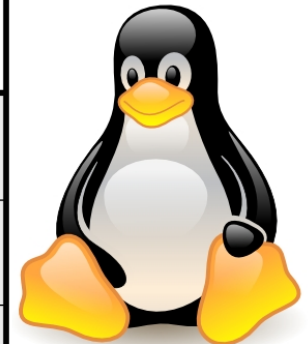




# *E. Coli* Transcriptional regulatory network vs Linux kernel call graph

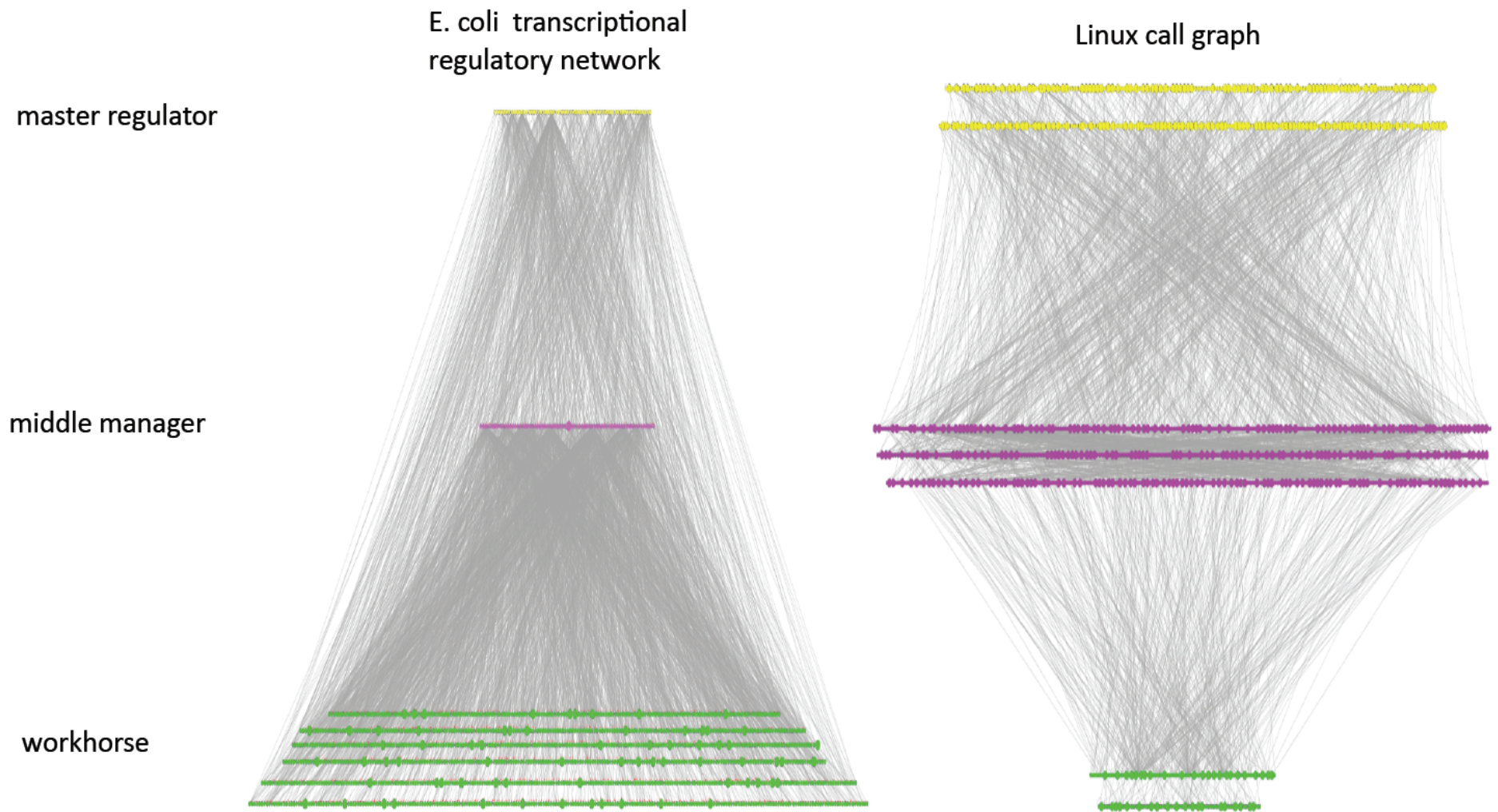


		<i>E. coli</i> transcriptional regulatory network	Linux call graph
<b>Basic properties of systems</b>	Nodes	Genes (TFs & targets)	Functions (subroutines)
	Edges	Transcriptional regulation	Function calls
	External constraints	Natural environment	Hardware architecture, customer requirements
	Origin of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



	<i>E. coli</i> transcriptional regulatory network	Linux call graph
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 references	200 bacterial genomes	24 versions of kernels
Years of evolution	Billions years	20 years

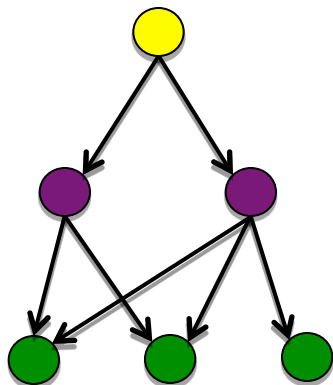
[Yan et al., PNAS (2010), in press]



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

# Comparison: hierarchical organization

Pyramidal vs  
Top-heavy



% in *E. coli*  
regulatory  
network

% in Linux  
call graph

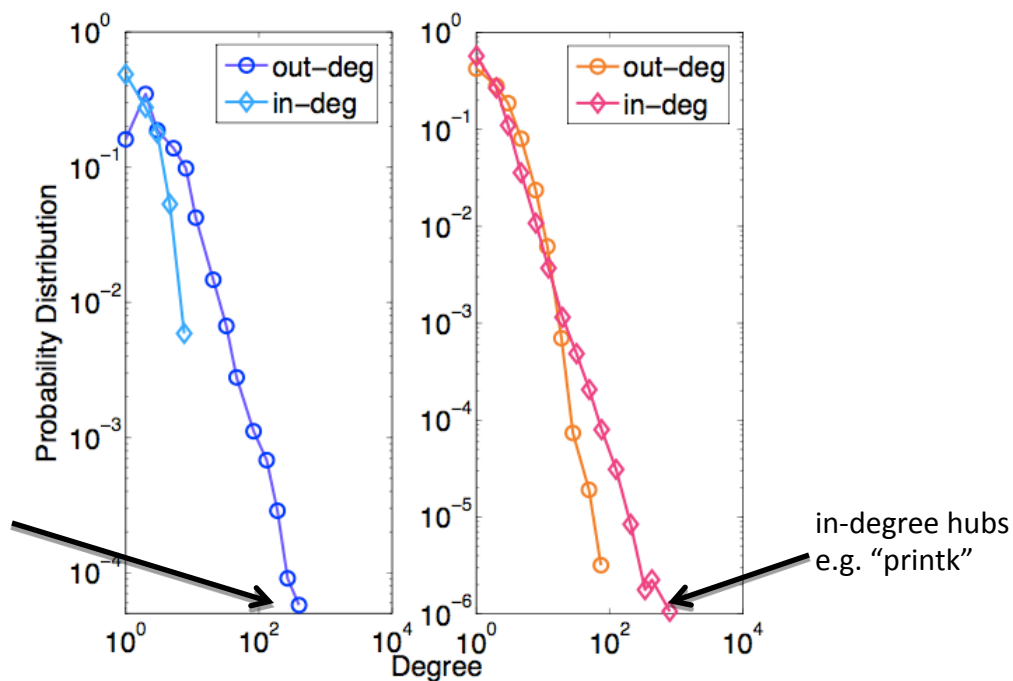
master  
regulator  
  
middle  
manager  
  
workhorse

master regulator	4.6	29.6
middle manager	5.1	58.2
workhorse	90.2	12.3

Degree distribution  
Roles of hubs

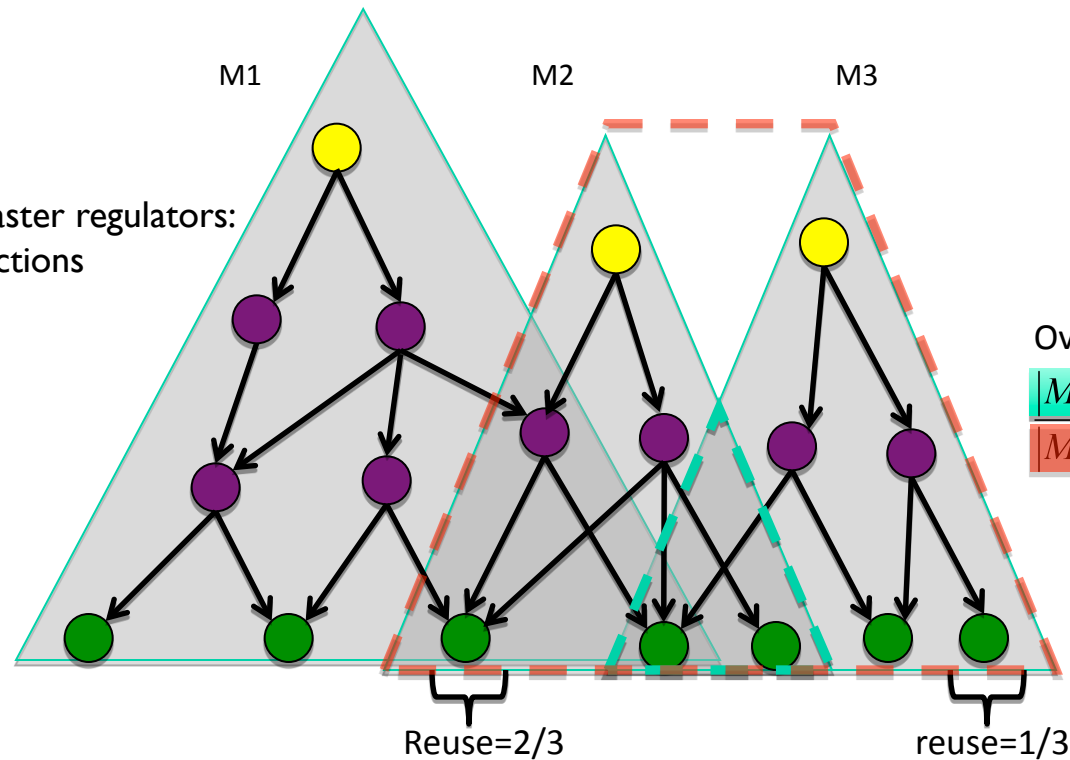
out-degree hubs  
e.g. "crp"

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



# Comparison: organization of modules

Modules are labeled by master regulators:  
TFs, high-level starting functions



TRN:  
modules overlap little,  
components are  
less generic: “ompF”

	E. Coli TRN	Linux call graph
# of modules	64	3665
Average overlap	4.3%	80.7%
Maximum node reuse	15.6%	87.5%
Average node reuse	3.5%	8.4%

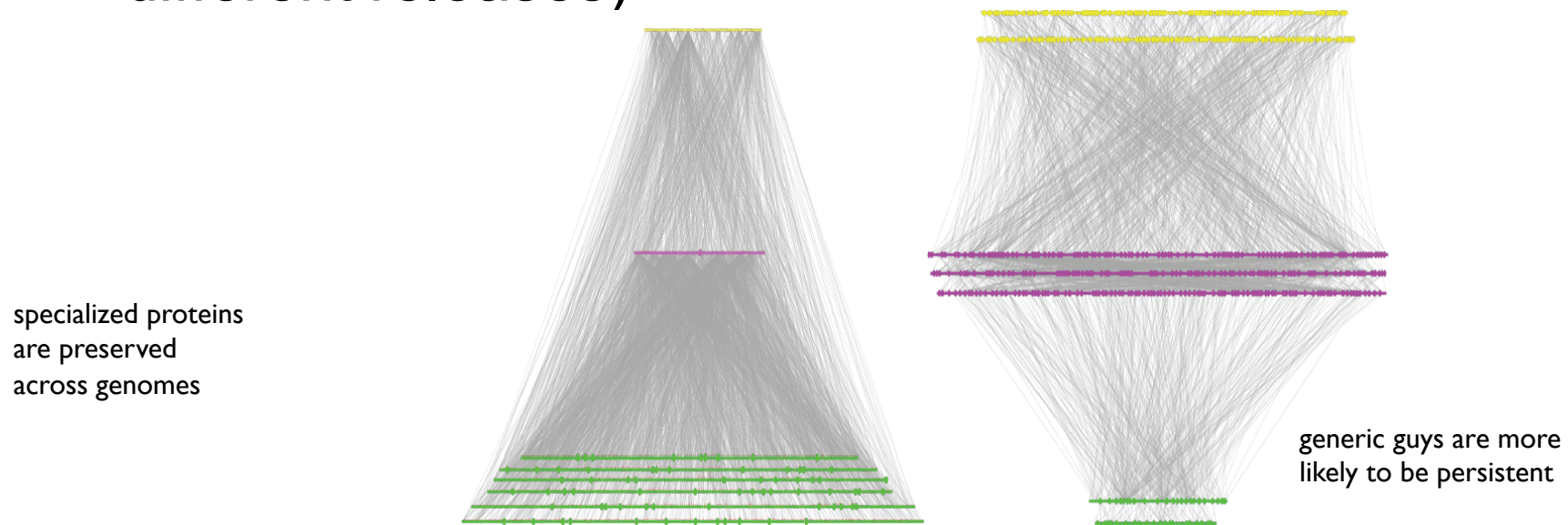
Call graph:  
modules overlap,  
Functions are highly  
reused (generic):  
“printf”

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



# Comparison of persistent components

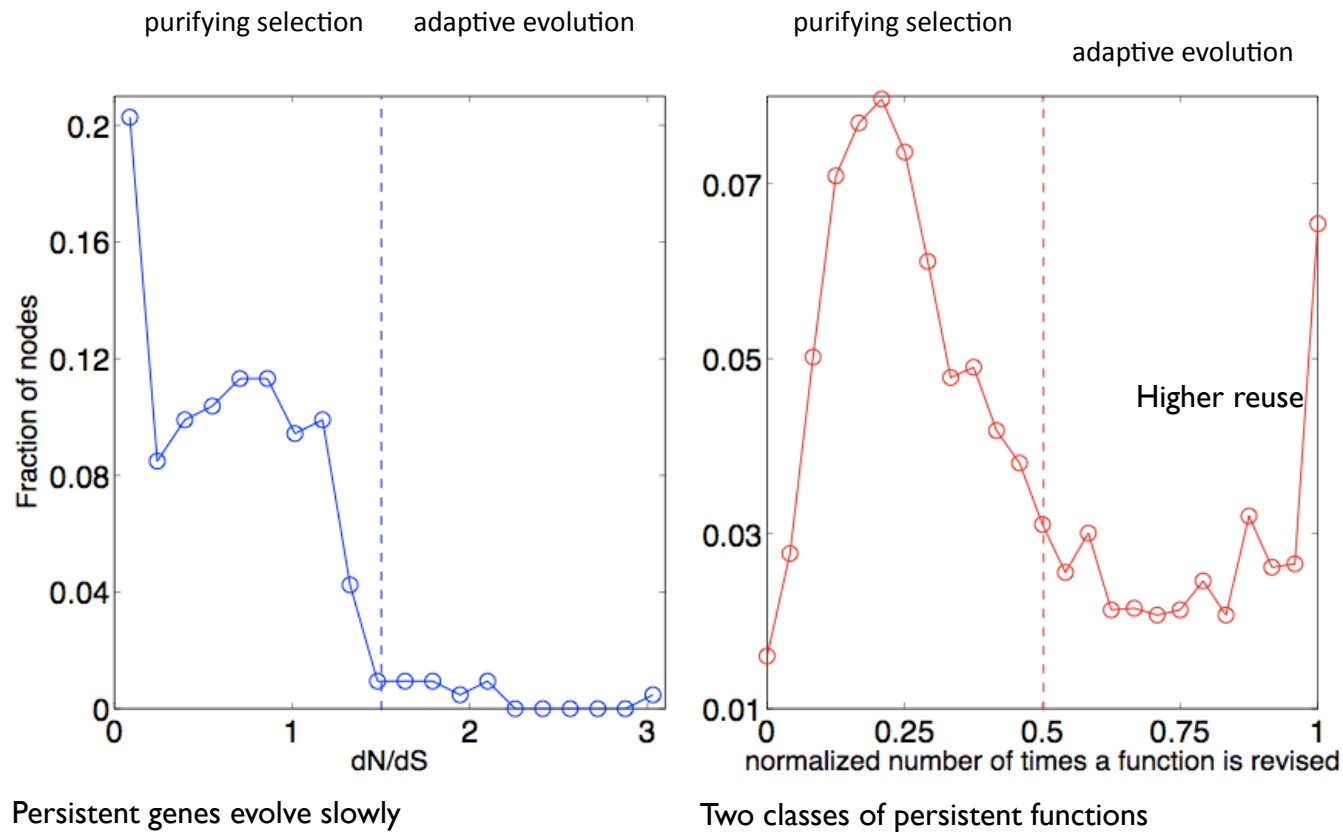
- Persistent genes (preserve among different genomes) vs persistent functions (preserve among different releases)



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

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

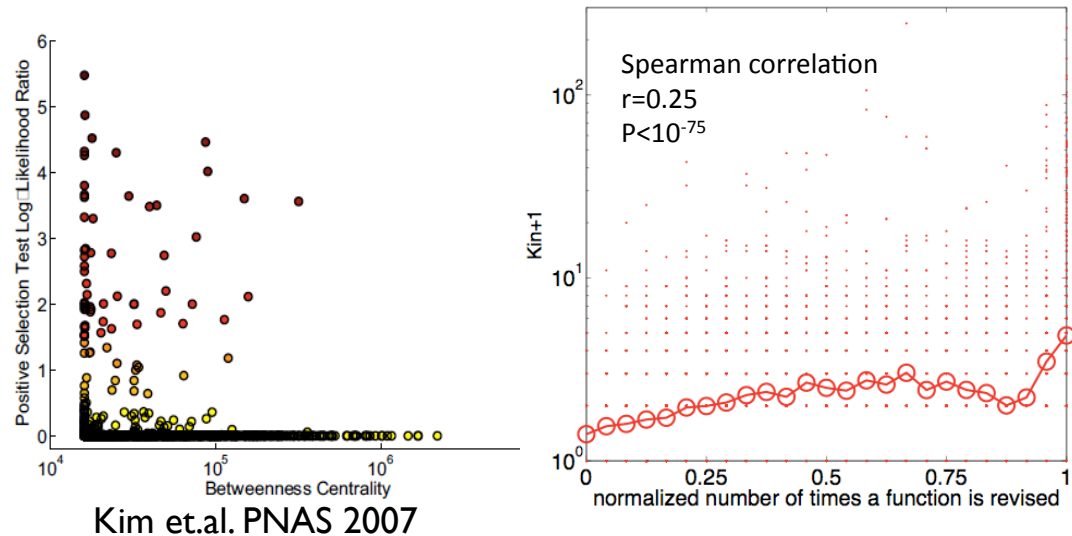
# Evolutionary rate of persistent functions



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

# Why and so what?

The difference can be explained by the nature of hubs evolution: tinkering vs design



- ▶ 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

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

# **Network Dynamics Across Environments: Metabolic Pathways**

**How do molecular networks change across environments?**

**What pathways are used more ?**

**Used as a biosensor ?**

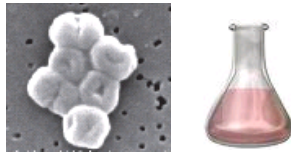




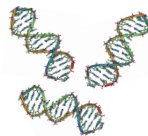
# What is Metagenomics?

## Traditional Genomics

Select organism and culture

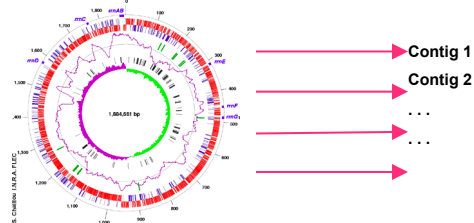


Extract DNA and sequence



atgctcgatctcg  
atcgatctcgctg  
atgccgatctaa

Assemble and annotate



Estimated that less than 1% of microbes can be cultured

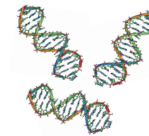


## Metagenomics

Collect sample from environment

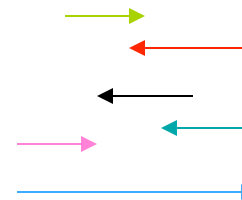


Extract DNA and sequence

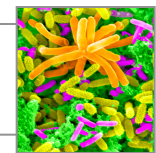


atgctcgatctcg  
atcgatctcgctg  
atgccgatctaa

Assemble and annotate



Lose information about which gene belongs to which microbe



# Sorcerer II Global Ocean Survey

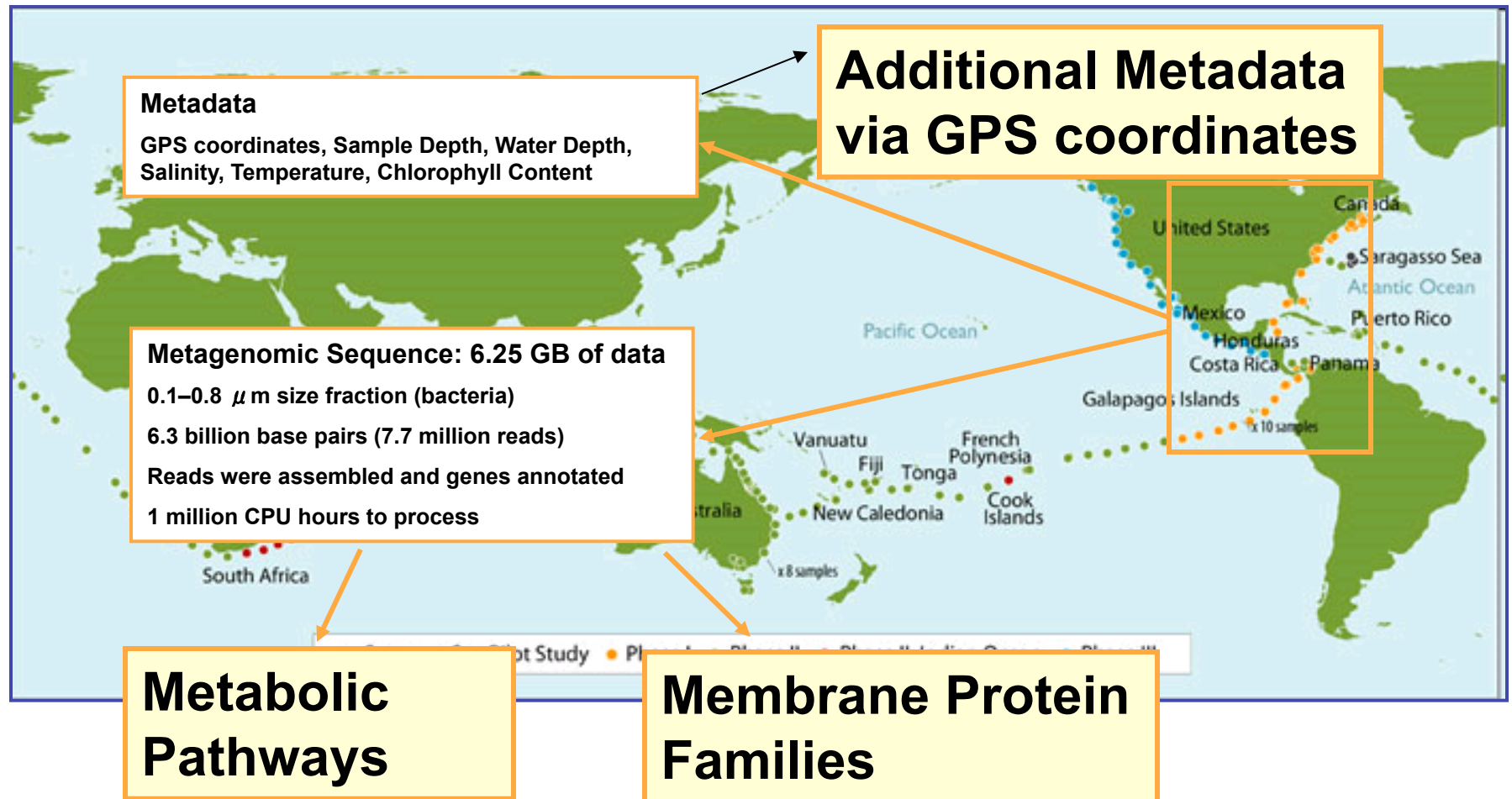


**Sorcerer II journey August 2003- January 2006**

**Sample approximately every 200 miles**



# Sorcerer II Global Ocean Survey

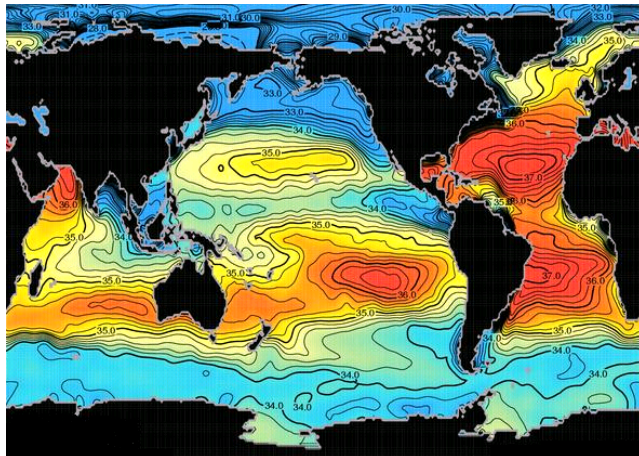


# Extracting environmental data using GPS Coordinates



**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

**GPS coordinates allow us to extract information from other sources:**



- \* **World Ocean Atlas**
- \* **National Center for Ecological Analysis and Synthesis**

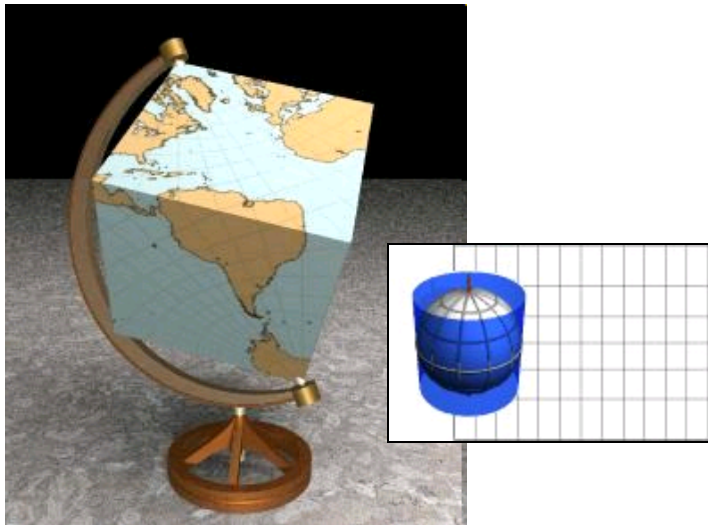


# World Ocean Atlas 2005

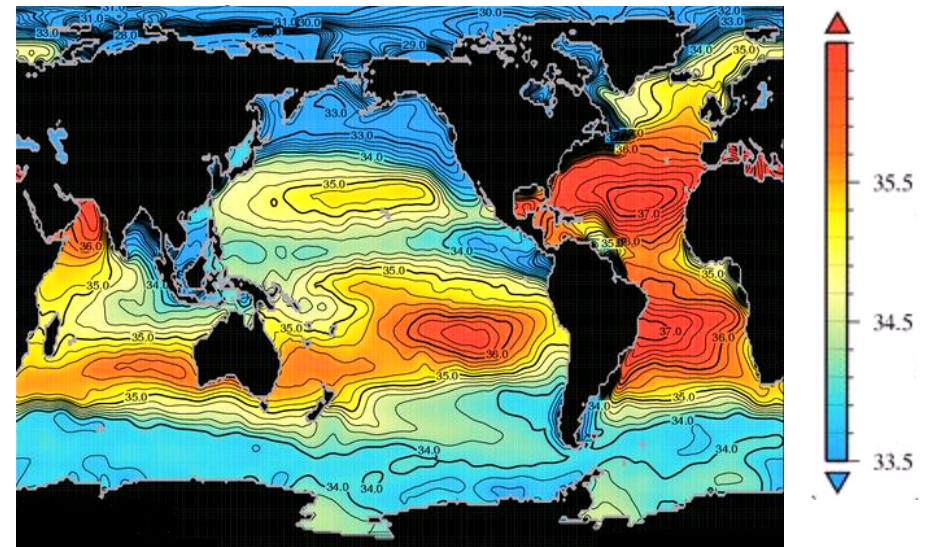
NOAA (National Oceanic and Atmospheric Administration) and  
NODC (National Oceanographic Data Center)

- \* Cumulative annual data at the ocean surface
- \* Resolution is 1 degree latitude/longitude

... no simple geometric shape matches the Earth



Annual Phosphate [ $\mu\text{mol/l}$ ] at the



## Nutrient Features Extracted:

Phosphate

Silicate

Nitrate

Apparent Oxygen Utilization

Dissolved Oxygen

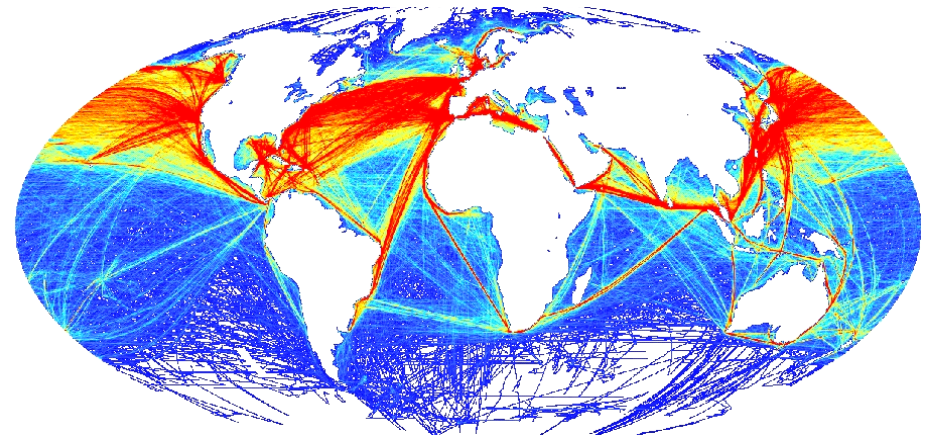
# National Center for Ecological Analysis and Synthesis (NCEAS)

- \* Resolution is 1 km square
- \* Value of a activity at a particular location is determined by the type of ecosystem present:

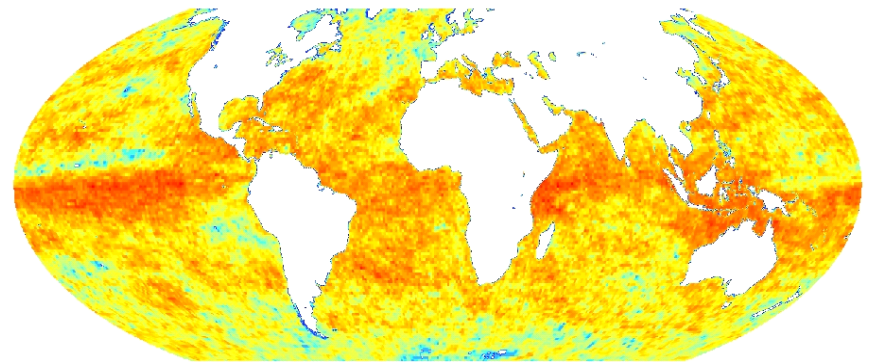
$$\text{Impact} = \sum \text{Features} * \text{Ecosystem} * \text{impact weight}$$

## Anthropogenic Features Extracted:

Ultraviolet radiation  
Shipping  
Pollution  
Climate Change  
Ocean Acidification



Shipping



Climate Change

READS → PROTEIN FAMILIES → PATHWAYS

CCGTGAGCACGATGCGC-----  
 ATGCTCATGCT-----  
 ATCGTGACGCGATGC-----  
 CCGTGAGCACGATGCGC-----  
 ATGCTCATGCT-----  
 ATCGTGACGCGATGC-----  
 ATGCTCATGCT-----  
 GCGATCGATCGATCGTAGC-----  
 TGCTGCTAGCATGCT-----  
 GCGATCGATCGATCGTAGC-----  
 TGCTGCTAGCATGCT-----  
 CCGTGAGCACGATGCGC-----  
 GTATCGTAGCATGCTT-----  
 CCGTGAGCACGATGCGC-----  
 GCGATCGATCGATCGTAGC-----



$$P_1 = f_1 + f_2 + f_3$$

$$P_2 = f_4 + f_5 + f_6$$

Mapping Raw  
Metagenomic  
Reads to a  
Matrix of  
Families or  
Pathways for  
each Site

PATHWAYS



SITES

$$P_{1,1} = 2 + 1 + 3$$

$$P_{2,1} = 2 + 4 + 3$$

$$P_{1,2} = 5 + 2 + 6$$

$$P_{2,1} = 5 + 7 + 6$$



	Fam 1	Fam 2	...	...	Fam 151
Site 1	.01	.02			
Site 2	0	.01			
...					
Site 29					

**Families Matrix**

# counts Fam 2  
#total protein counts at site 2

Pathway Sequences  
(Community Function)

Environmental  
Features

Metabolic Pathways

Sites

	P1	P2	P3		
B1	3800	1400	1000		
B2	2200	100	400		
↓	---	---	---		



Environmental

Metadata

Sites

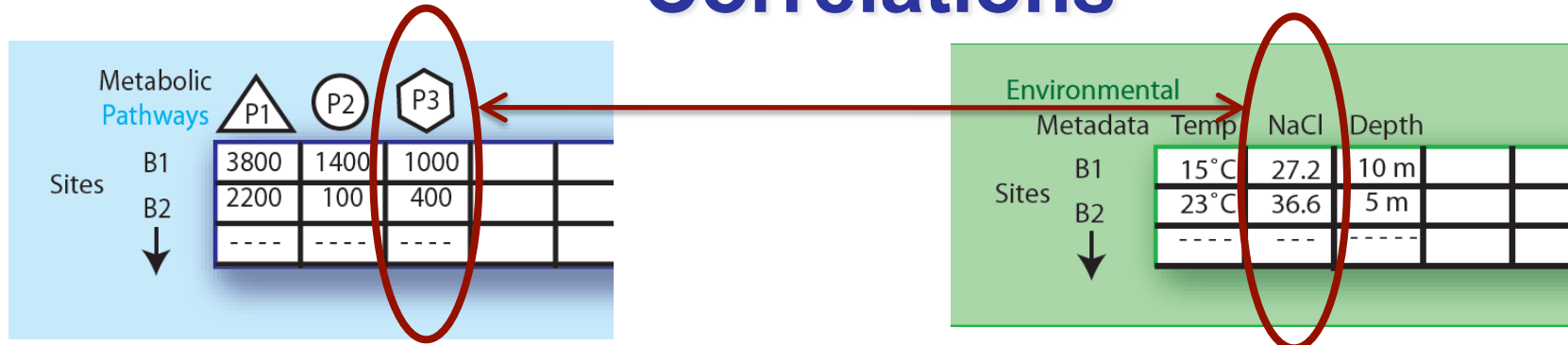
	Temp	NaCl	Depth		
B1	15°C	27.2	10 m		
B2	23°C	36.6	5 m		
↓	---	---	---		

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

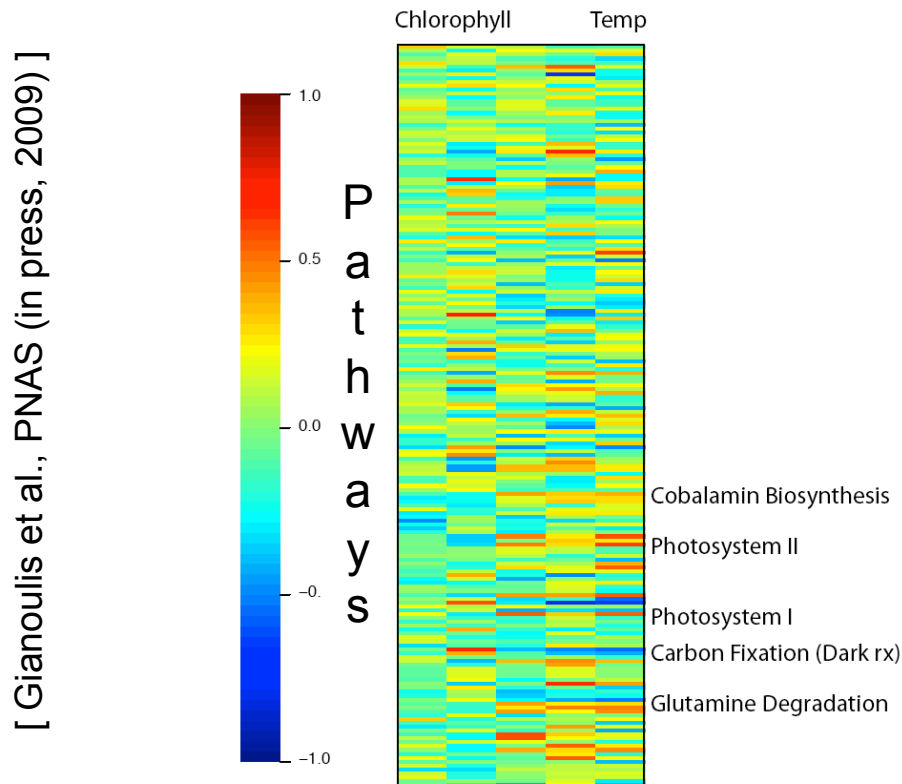
[Rusch et. al., (2007) PLOS Biology;  
Gianoulis et al., PNAS (in press, 2009)]




# Simple Relationships: Pairwise Correlations







Environmental Features







# Canonical Correlation Analysis: Simultaneous weighting


Fitness Test	# km run/week
5K run	





Fitness Text Index	Exercise Index
<div>5K run</div> <div>Heart Rate</div> <div>Lifting max</div>	   

$$FTI = a \text{ 5K run } + b \text{ Lifting max } + c \text{ Heart Rate }$$

$$EI = a \text{  } + b \text{  } + c \text{  } + d \text{  }$$

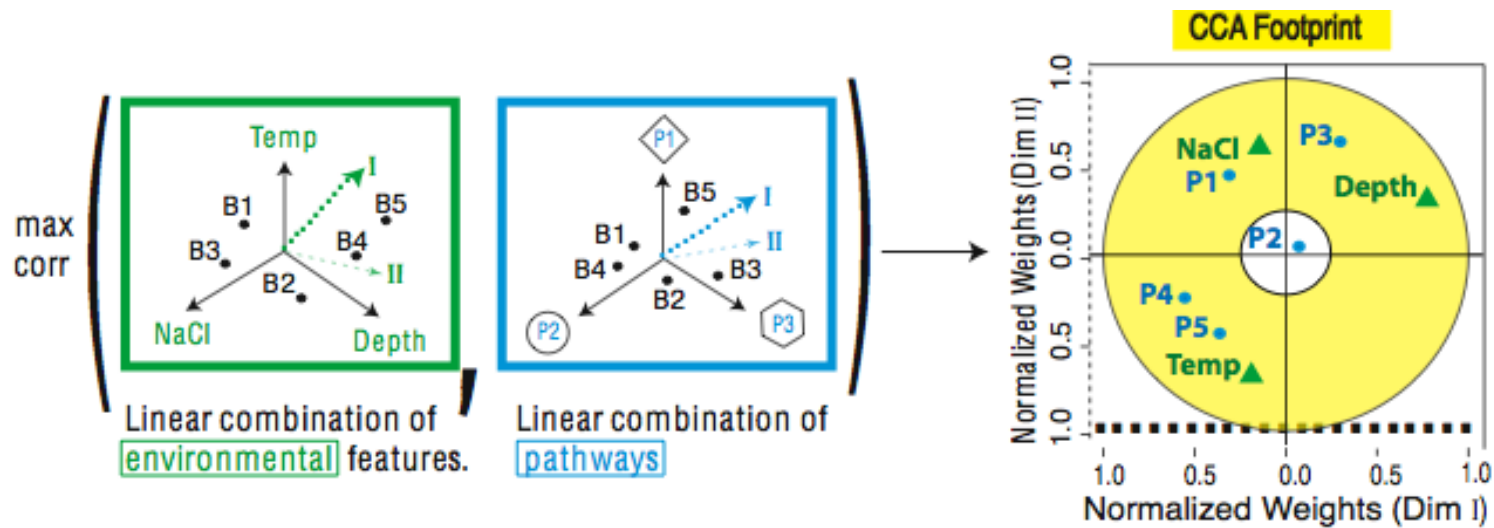
# Canonical Correlation Analysis: Simultaneous weighting

Fitness Test	# km run/week
5K run	

Fitness Text Index	Exercise Index
<div>5K run</div> <div>Heart Rate</div> <div>Lifting max</div>	   

Environmental Features	Metabolic Pathways/ Protein Families
Temp      etc	Photosynthesis      etc
Chlorophyll	Lipid Metabolism

# Environmental-Metabolic Space



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

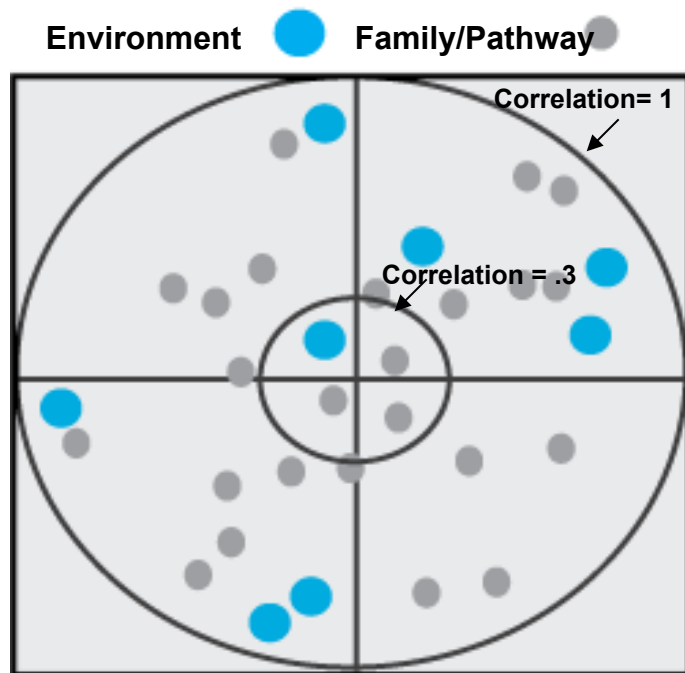
Given  $X = \{x_1, x_2, \dots, x_n\}$  and  $Y = \{y_1, y_2, \dots, y_m\}$

$$C = \begin{matrix} \sum_X & \sum_{X,Y} \\ \sum_Y & \sum_{Y,X} \end{matrix} \quad \max_{a,b} \text{Corr}(U,V) = \frac{a' \sum_{12} b}{\sqrt{a' \sum_{11} a} \sqrt{b' \sum_{22} b}}$$

Gianoulis et al., PNAS 2009



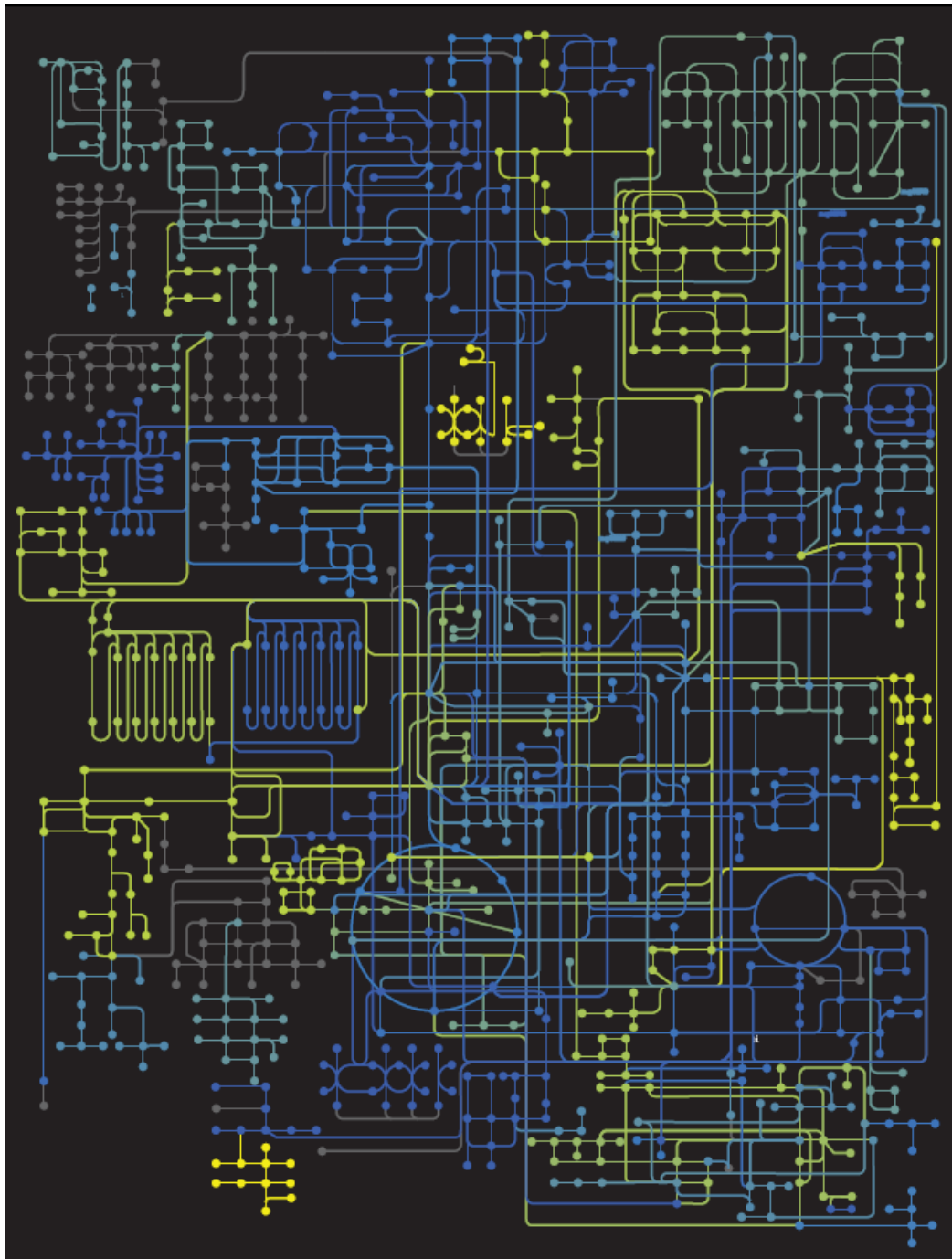
# CCA Example



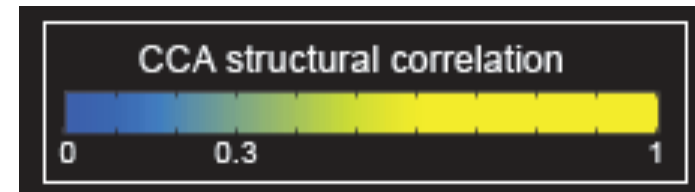
**This plot shows the correlations in the first and second dimensions**

**Correlation Circle: The closer the point is to the outer circle, the higher the correlation**

**Variables projected in the same direction are correlated**

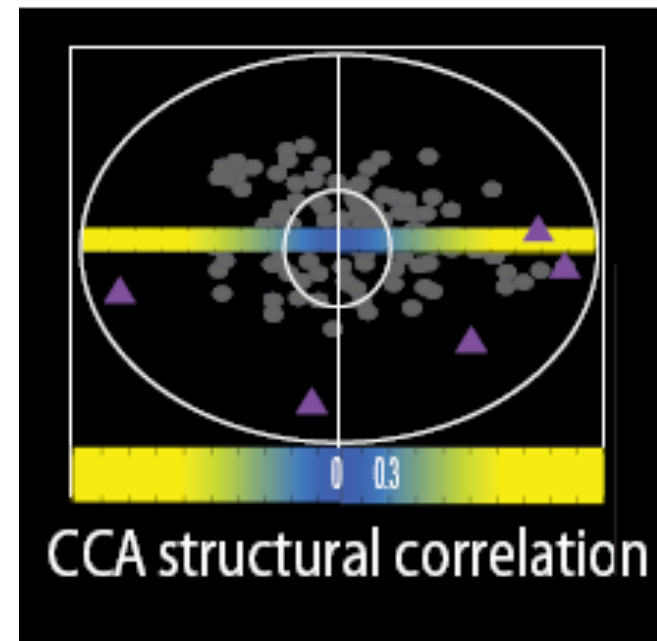


## Strength of Pathway co-variation with environment



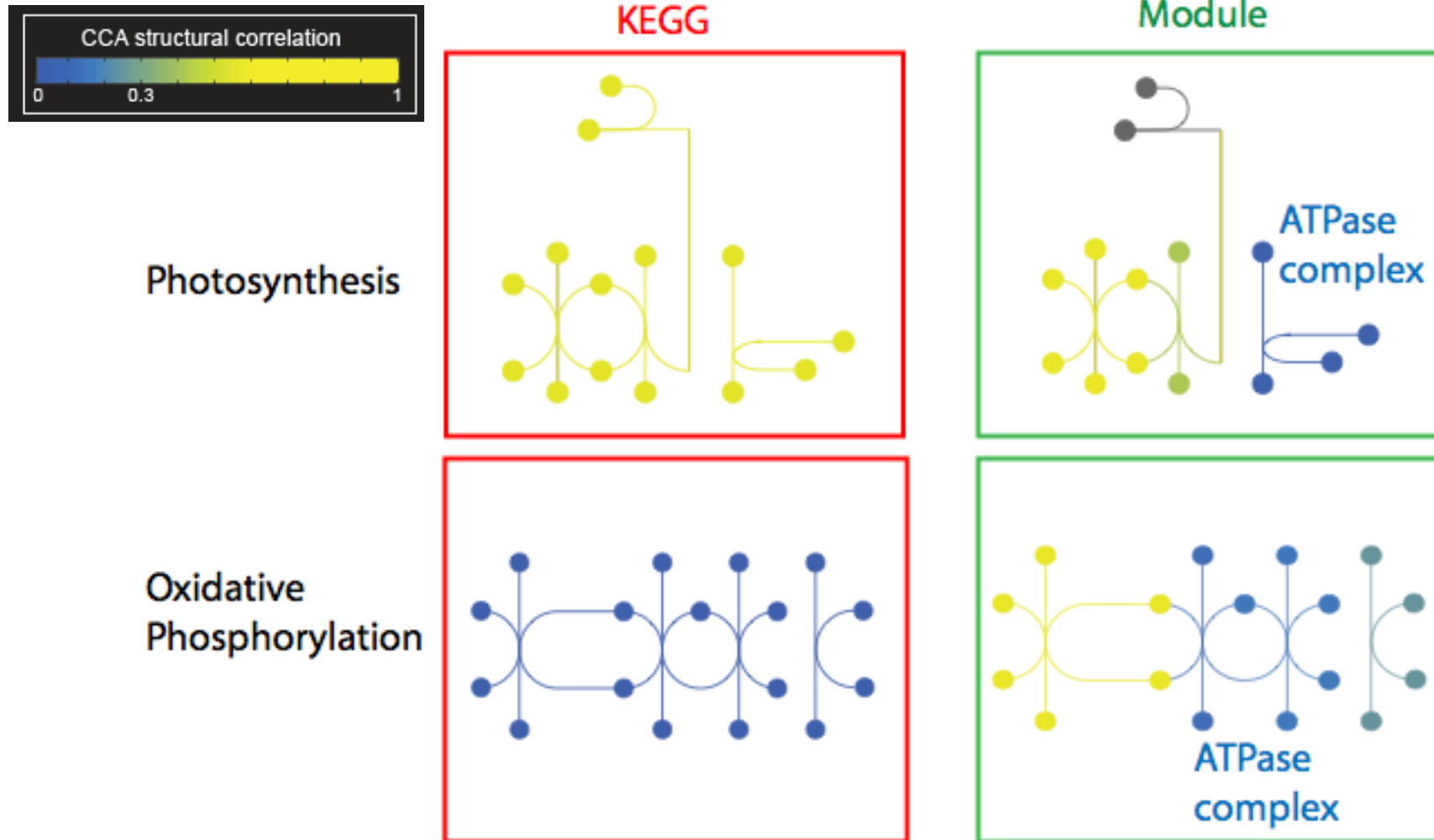
Environmentally  
invariant

Environmentally  
variant

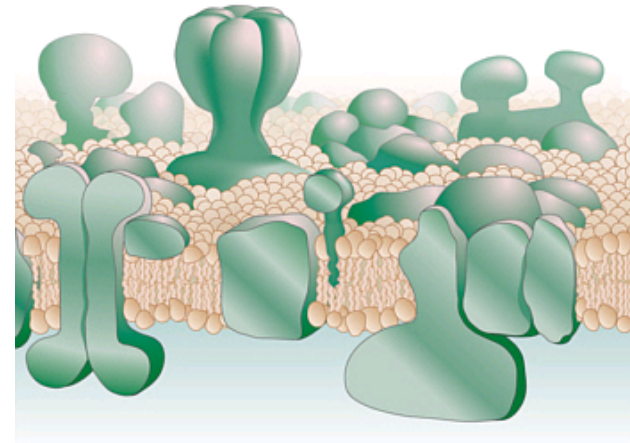
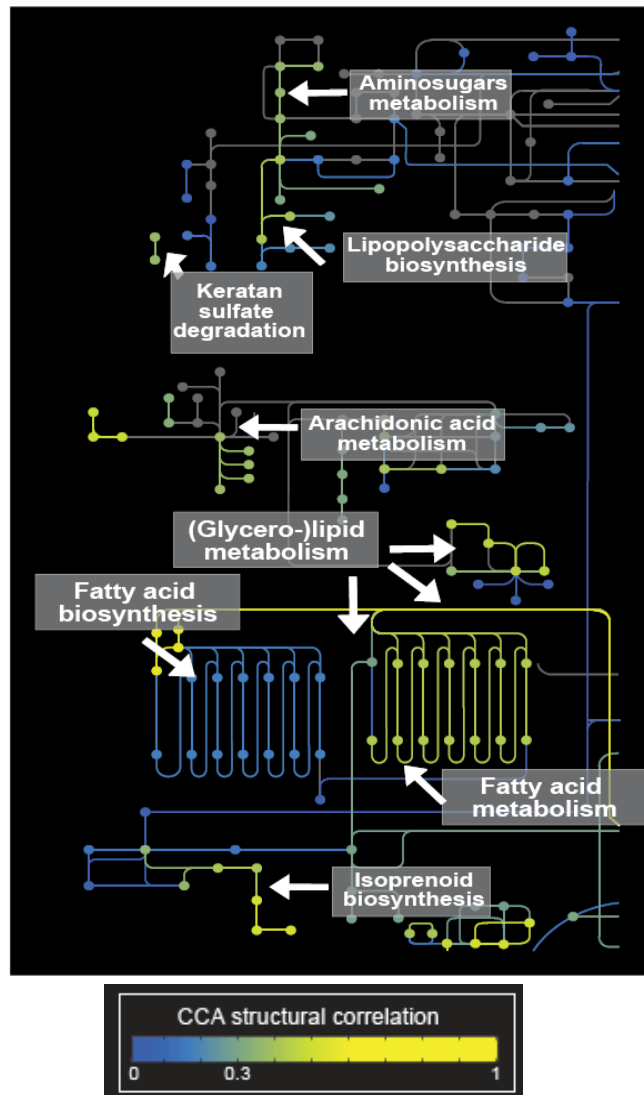


Gianoulis et al., *PNAS* 2009

# Conclusion #1: energy conversion strategy, temp and depth



## 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**

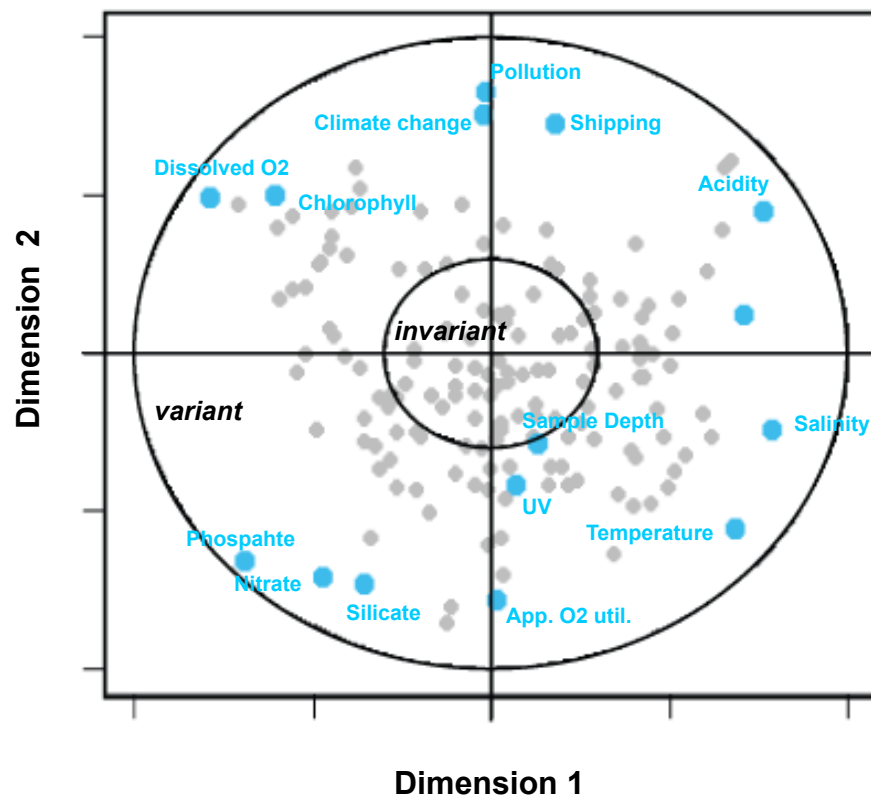
Gianoulis et al., *PNAS* 2009  
Patel et al. *Genome Research* 2010



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



# CCA results for Membrane Proteins



**107 variant membrane protein families**

**44 invariant membrane protein families**

**Difficult to see the strength and directionality of a relationship**

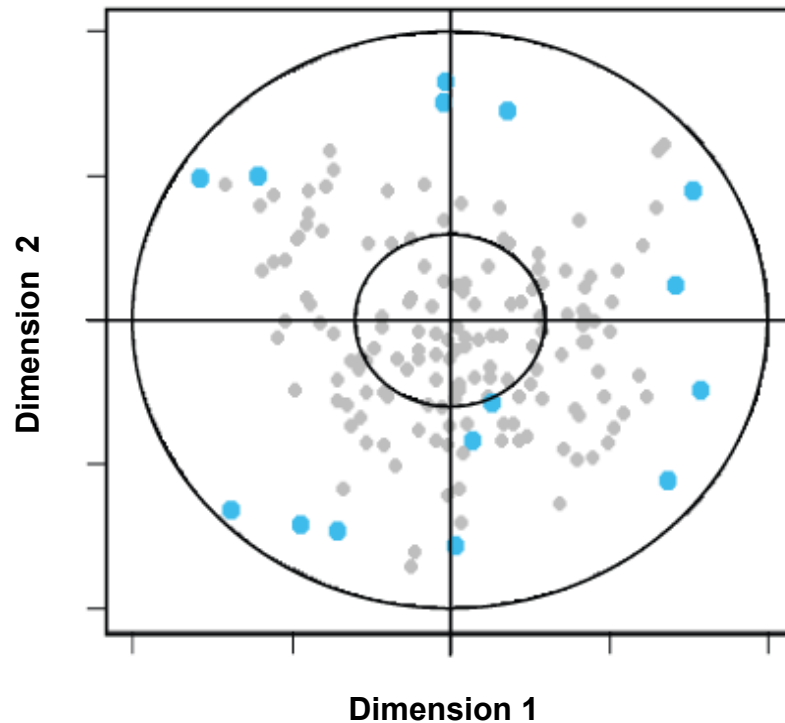
**Weights of the features are difficult to visualize and compare**

**There is no means of quantifying the variation between sets of features**

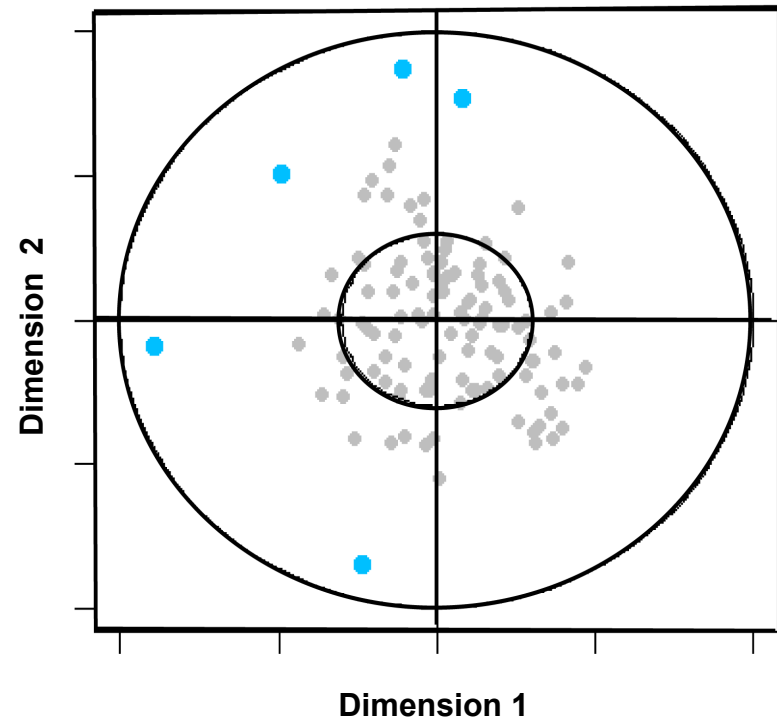
# Membrane Proteins vary more than Metabolic Pathways

Median absolute structural Correlation Coefficient

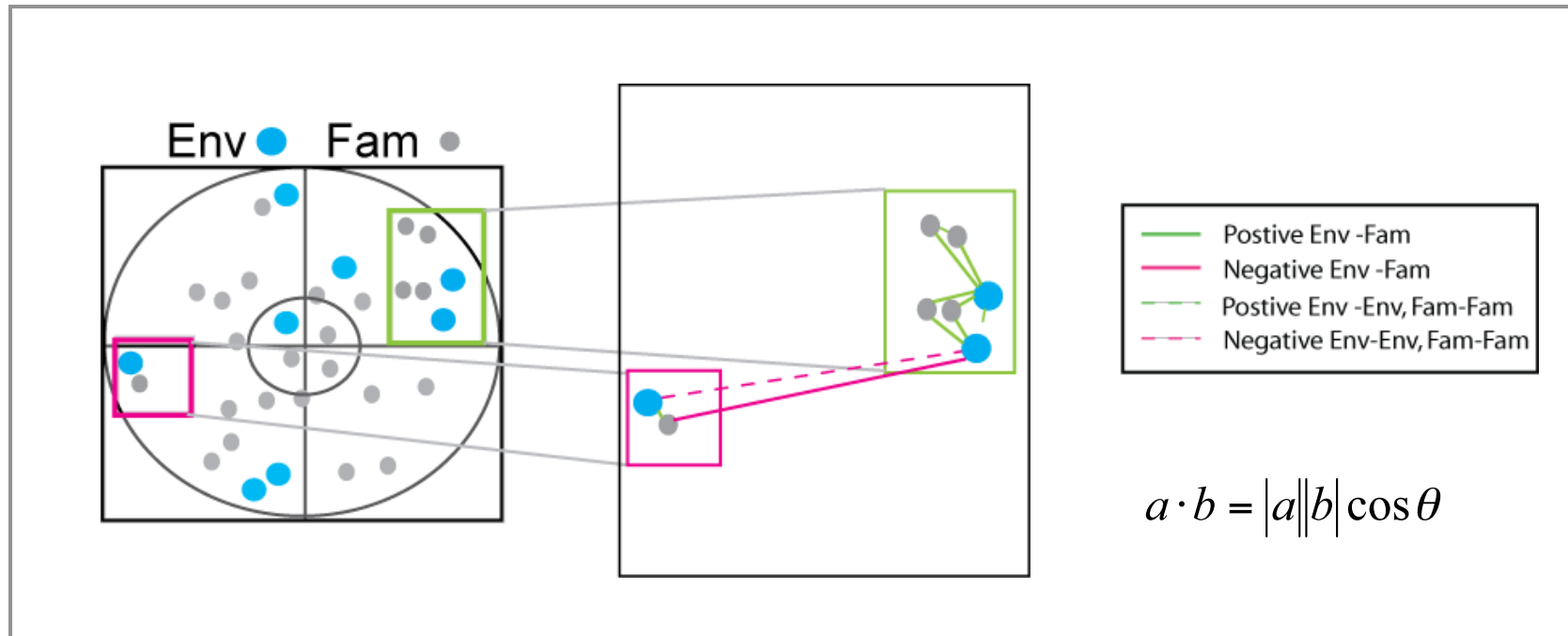
**Membrane Proteins = 0.3**



**Metabolic Pathways = 0.17**



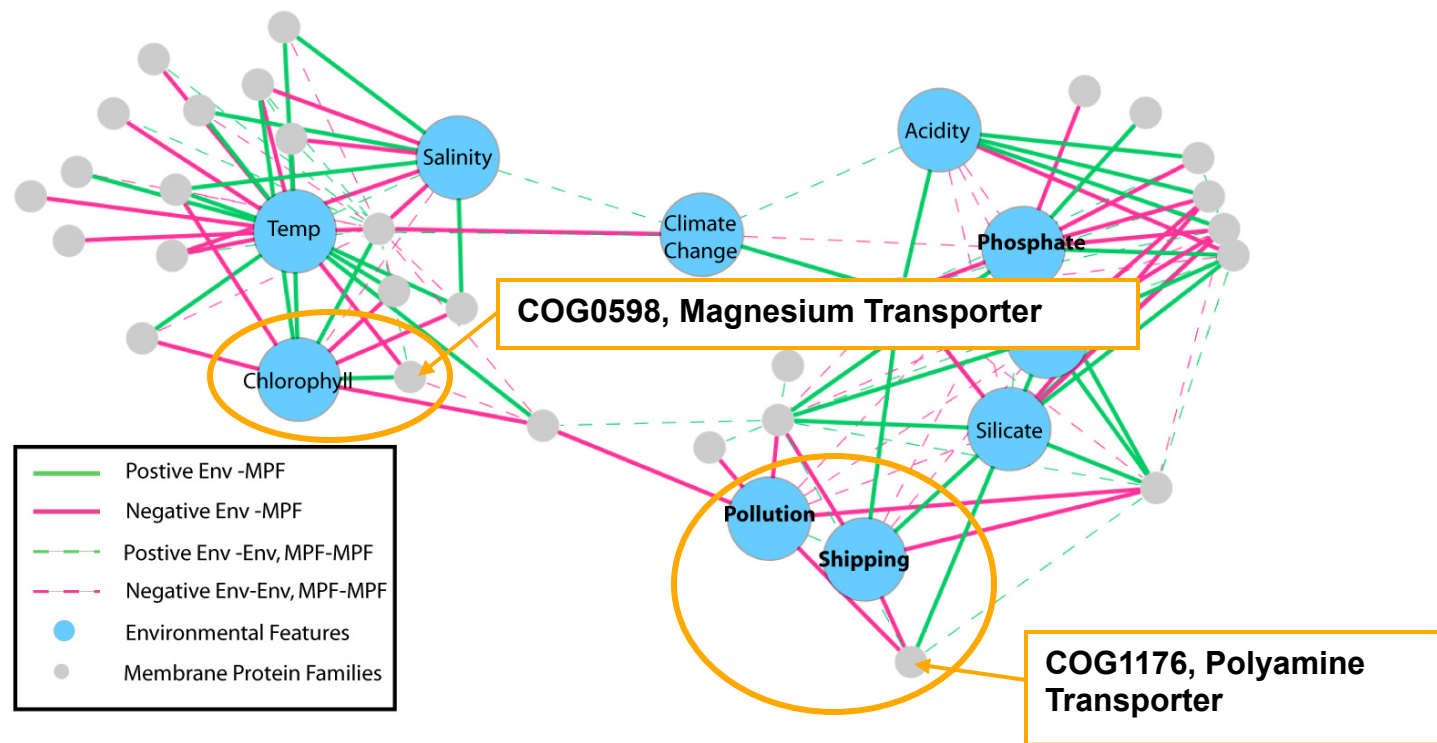
# Protein Families and Environmental Features Network (PEN)



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



# Protein Families and Environmental Features Network (PEN)

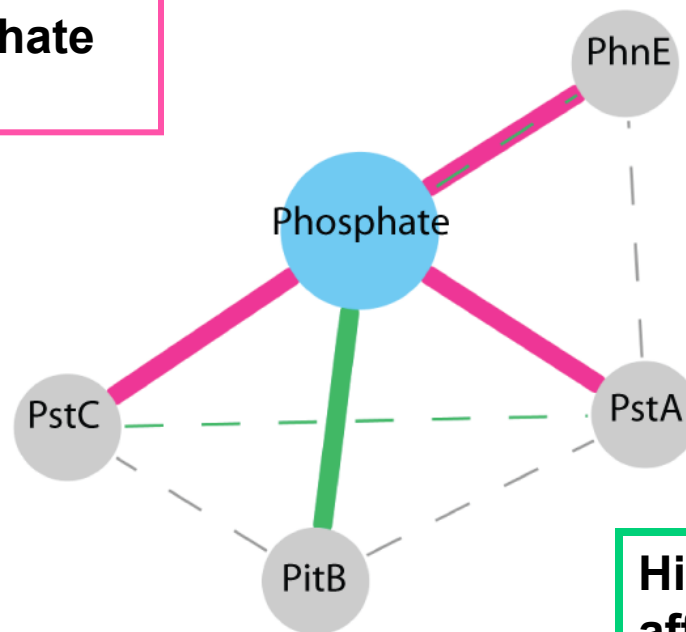


“Bi-modules”: groups of environmental features and membrane proteins families that are associated

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

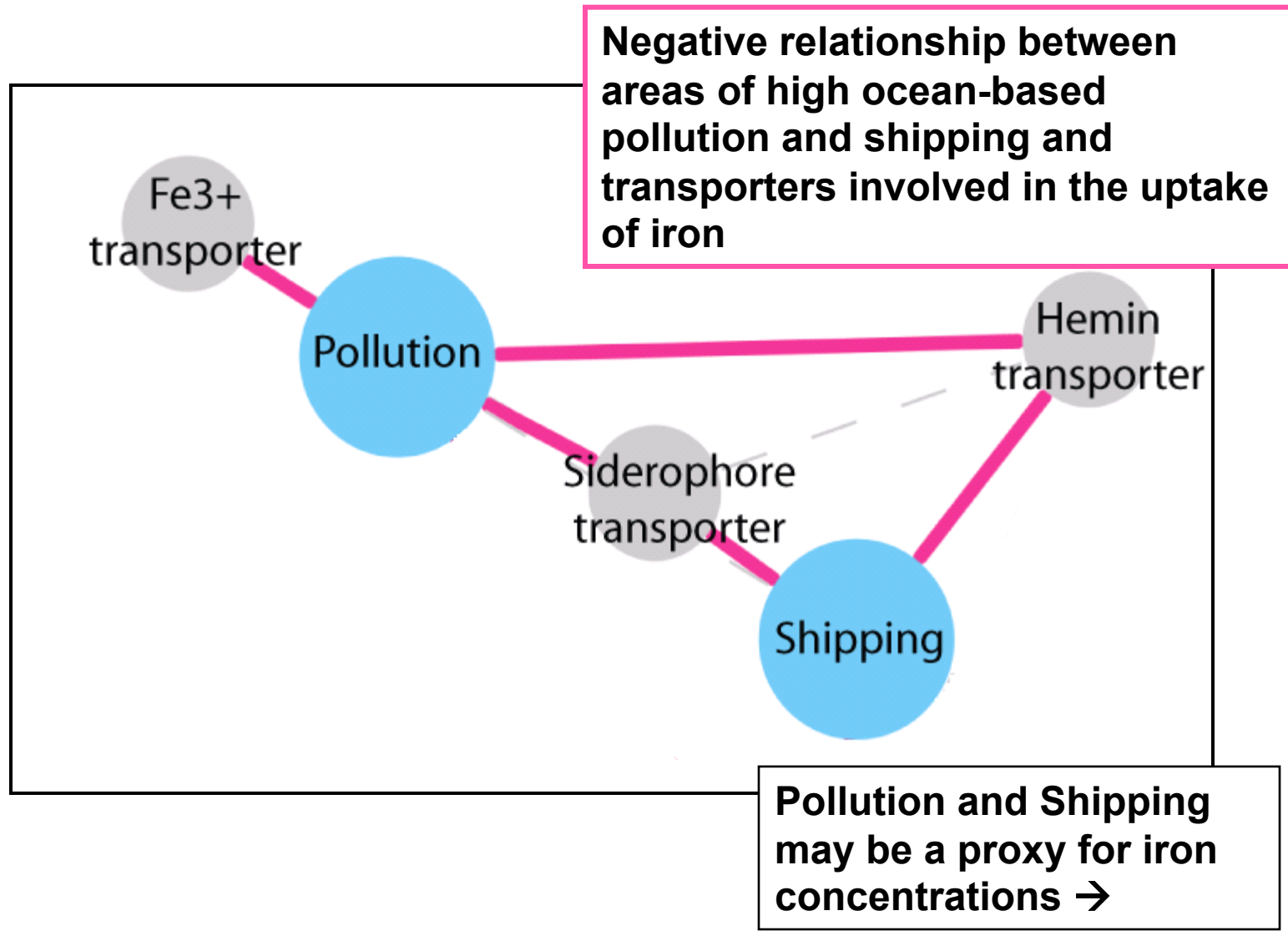
## Bi-module 1: Phosphate/Phosphate Transporters

**Low Phosphate, high affinity phosphate transporters which are induced during phosphate limitation**

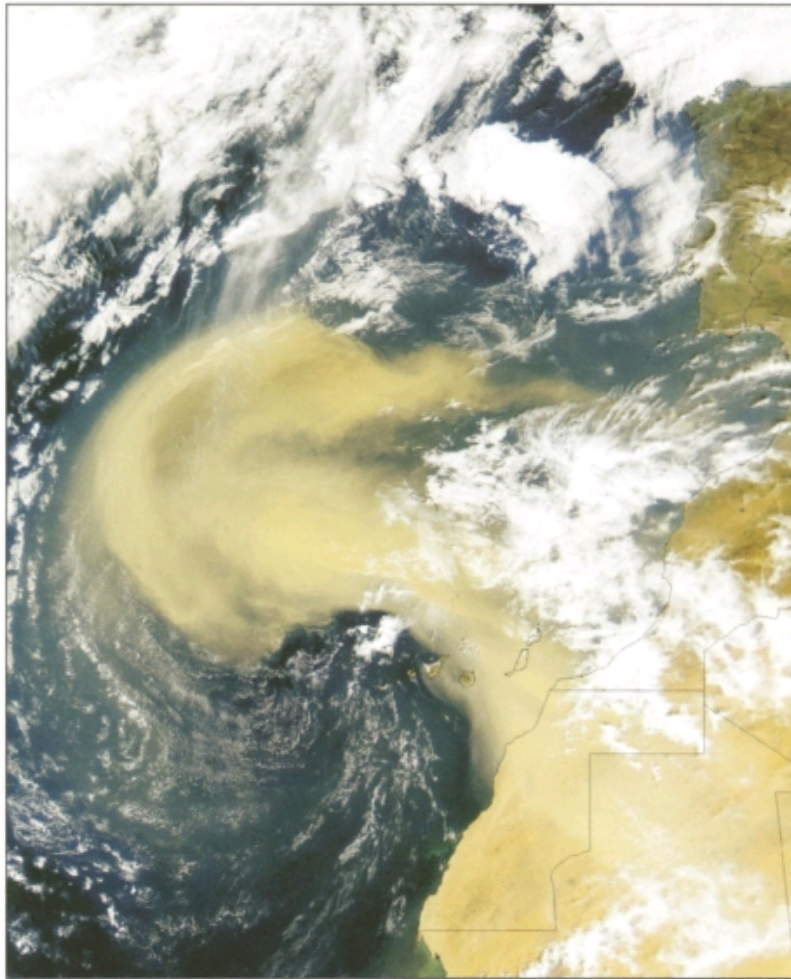


**High Phosphate, low affinity inorganic phosphate ion transporter which are constitutively expressed**

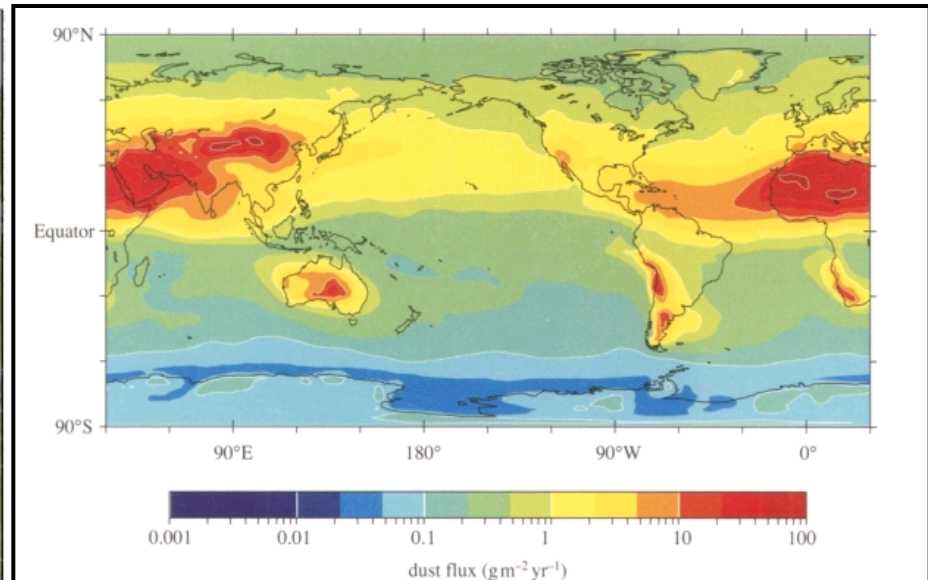
## Bi-module 2: Iron Transporters/Pollution/Shipping



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



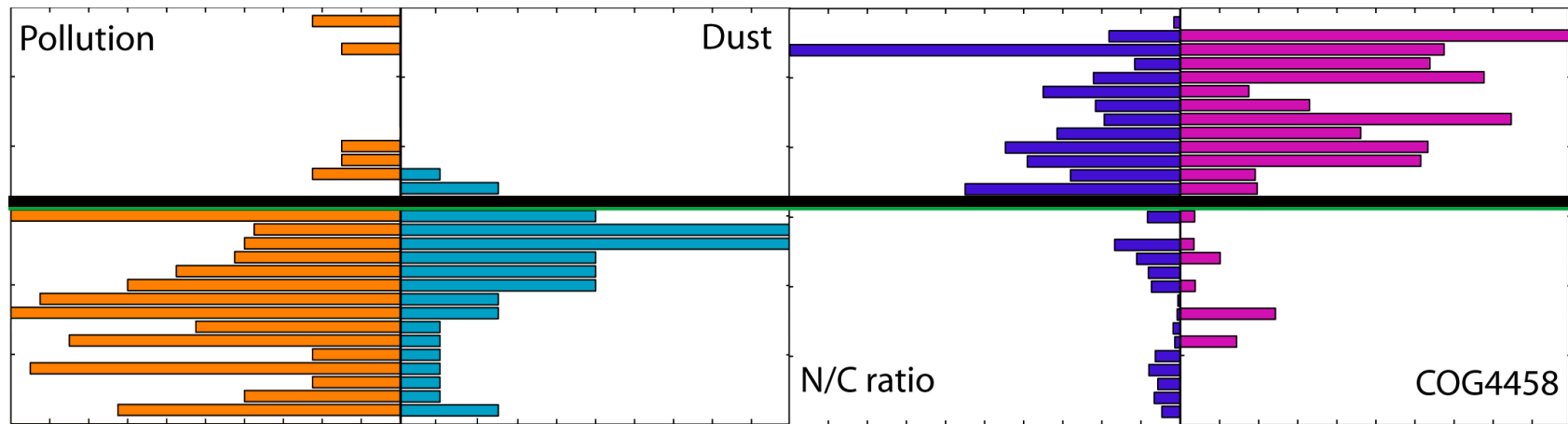
Rigwell A. J. (2002) *Phil. Trans. R. Soc. Lond.*



**Iron is usually limiting in oceans: High Nitrate-Nutrient/Low Chlorophyll regions**  
**Delivery of iron to is usually by:**

- terrestrial input
- fluvial (rivers) input
- upwelling from the ocean floor
- aeolian dust from land

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



Pollution and Dust ↑

N/C and Iron Transporters ↓

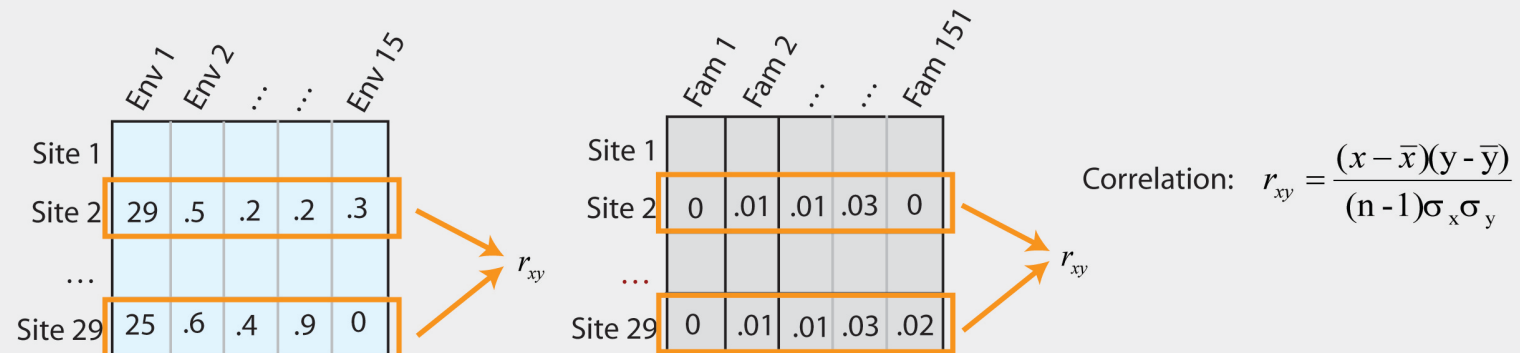
-Negative correlation between COG4558 and COG0609 and dust/pollution values (p-value <0.01)

- 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.

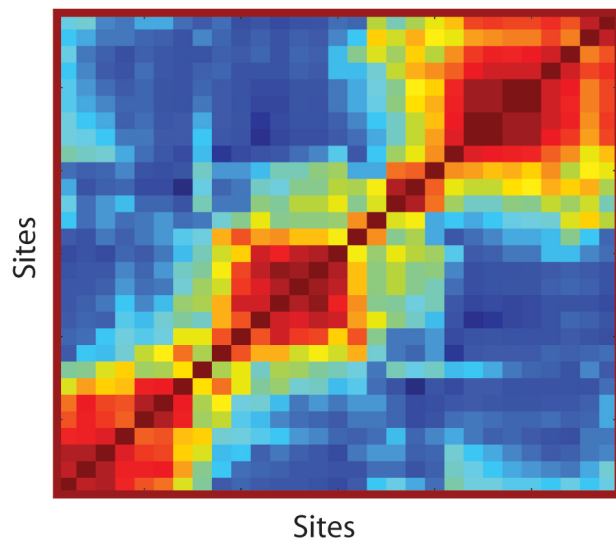


# How Similar are the Sites to each other?

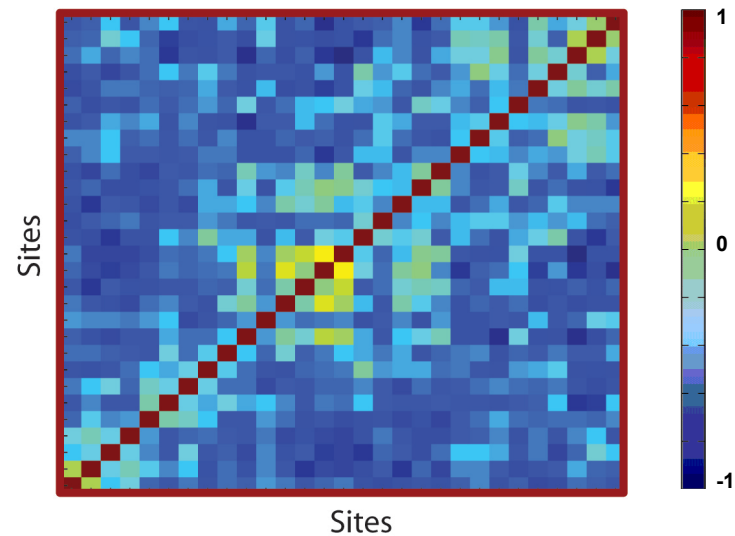
## Site-Site Correlations (Fam, Env)



Environmental Features Site Correlations

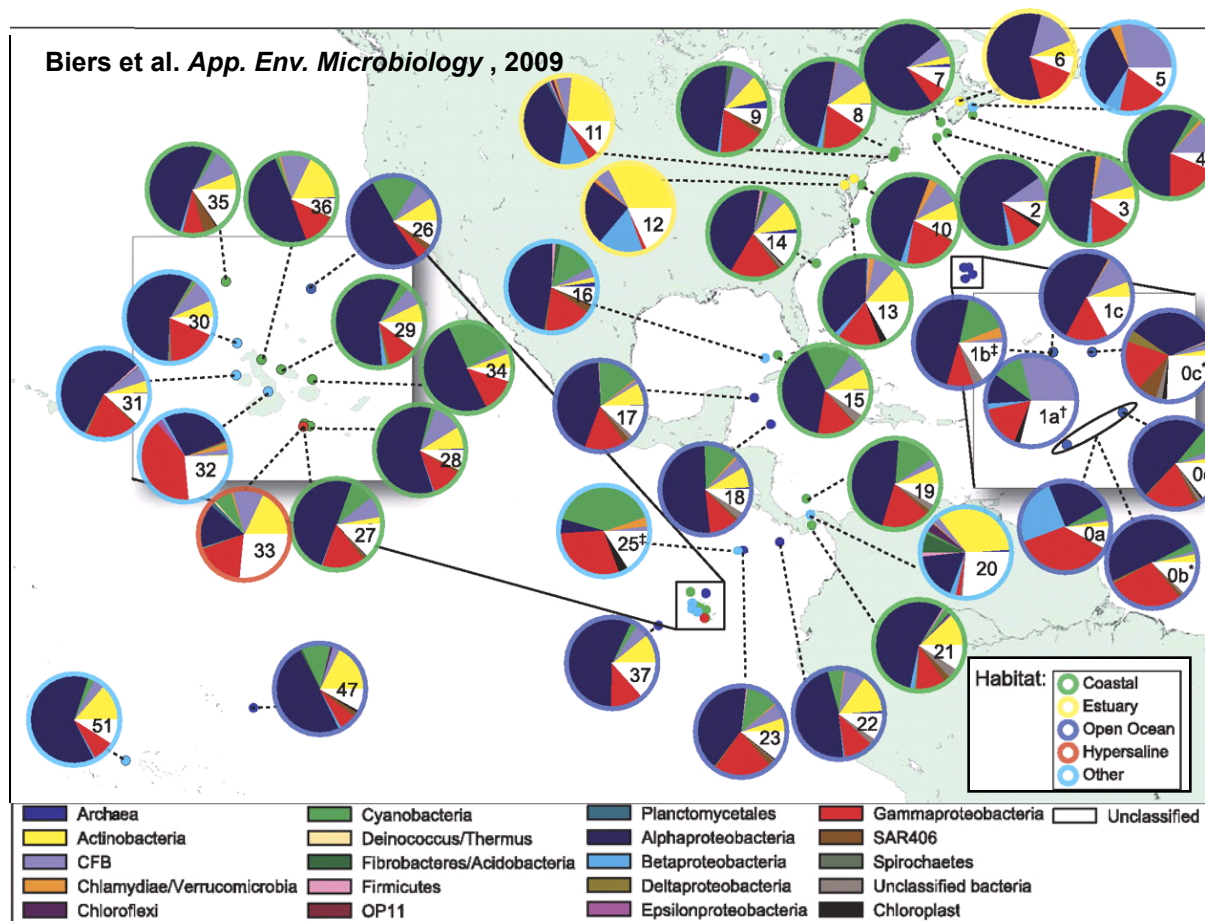


Membrane Protein Families Site Correlations



# Species Distribution

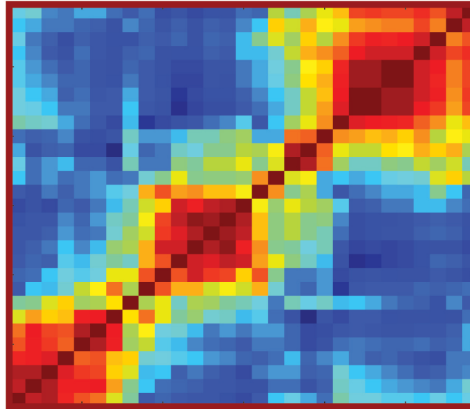
- The 16S rRNA gene is a component of the small prokaryotic ribosomal subunit
- Bacteria with 16S rRNA gene sequences more similar than 97% are considered the same 'species'
- 10,025 16S genes found and classified



	Cyanobacteria	Alphaproteobacteria	Gammaproteobacteria	...
Site 1	15	3		
Site 2	20	7		
...				
Site 29				

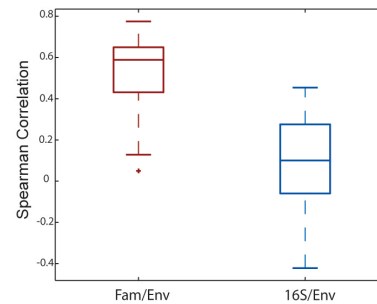
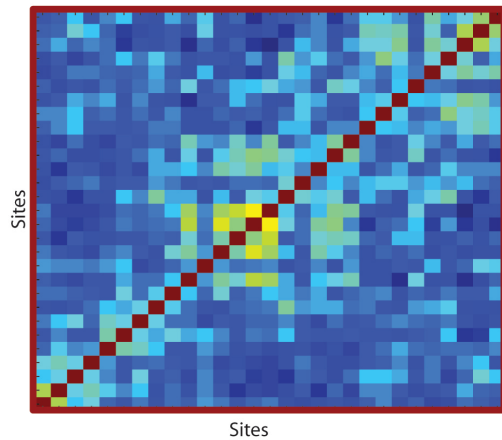
20% level, "phylum"

Environmental Features Site Correlations

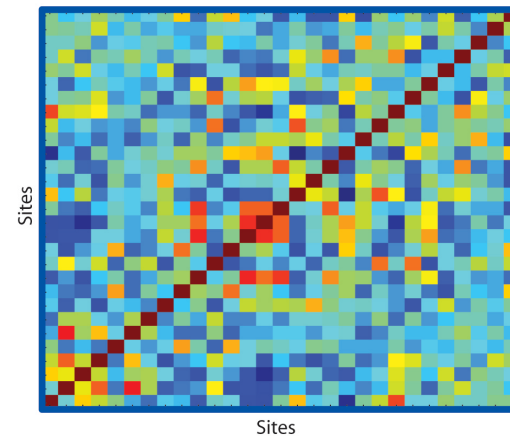


This suggests that the observed membrane protein variation is more a function of the measured environmental features, than phylogenetic diversity.

Membrane Protein Families Site Correlations



16S Site Correlations



Method: For each site, we correlated the EF profile distances and its MPF frequency profile distances and 16S profile distances

Environmental Features Site Correlations

	Site 1	Site 2	Site 3	Site 4	...					
Site 1	1	.8	.5	.7						
Site 2	.8	1	.3	.4						
Site 3	.5	.3	1	.7						
...	.7	.4	.7	1						

Spearman rank correlation:

1	2	4	3	...						
---	---	---	---	-----	--	--	--	--	--	--

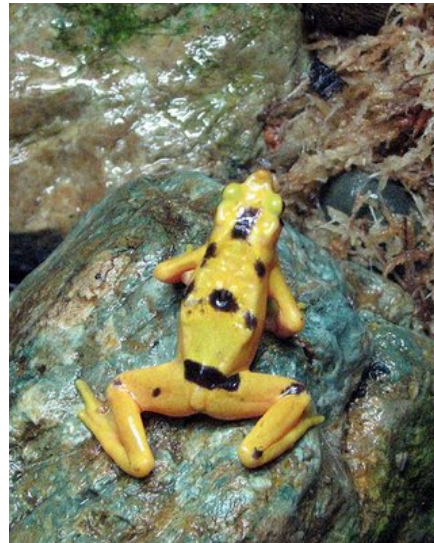
Membrane Protein Families Site Correlations

	Site 1	Site 2	Site 3	Site 4	...					
Site 1	1	.2	.6	.3						
Site 2	.2	1	.9	.8						
Site 3	.6	.9	1	.6						
...	.3	.8	.6	1						

1	4	2	3	...						
---	---	---	---	-----	--	--	--	--	--	--

$$r_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{(n-1)\sigma_x\sigma_y}$$

# Biosensors: Beyond Canaries in a Coal Mine



- Why Networks?
- Background:  
Central Network Points
- Networks & Variation  
(human ppi)
- 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)
- Network Dynamics Across Environments  
(prokaryote metab. pathways)
  - Metabolic Pathways
  - Entry pts. (Mem. Proteins)

## Outline: Molecular Networks





# Conclusions: Networks & Variation



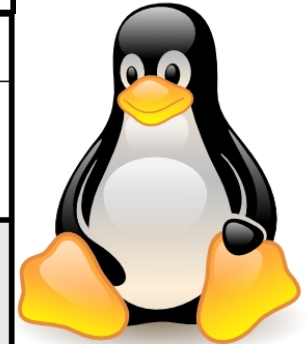
- Positive selection (adaptive evolution) at the network periphery
  - On a sequence level, it can be seen as positive selection of peripheral nodes
  - On a structural level, it can be seen as the pattern of SDs that display significantly higher allele frequencies in non-central genes

# Conclusions: Comparison to Social and Regulatory Hierarchies

- Regulatory Network 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)
- Number of collaborative partners saturates even while scale of targets governed increases
  - Also seen in social networks



		<i>E. coli</i> transcriptional regulatory network	Linux call graph
<b>Hierarchical organization</b>	Structure	Pyramidal	Top-heavy
	Characteristic hubs	Upper-level TFs with high out-degree	Generic workhorse functions with high in-degree
<b>Organization of modules</b>	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
<b>Persistent nodes</b>	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)
<b>Design principles</b>	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



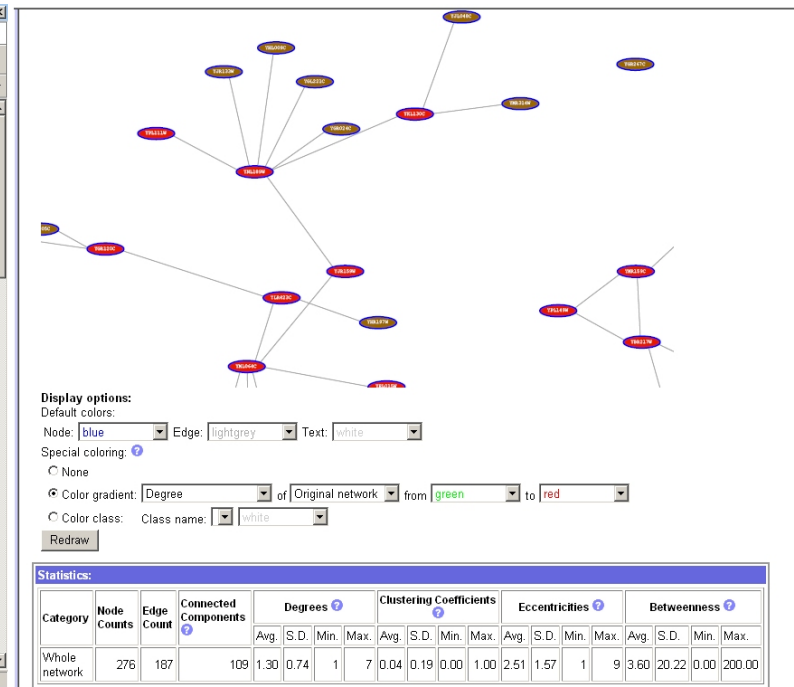
- an automated web tool

**tYNA**

(vers. 2 :

**"TopNet-like**

**Yale Network Analyzer"**)



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

**P Kim**

**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



**Networks.GersteinLab.org**

Job opportunities currently for postdocs & students

# More Information on this Talk

**SUBJECT:** Networks

**DESCRIPTION:**

Institut de recherches cliniques de Montreal (IRCM), Montreal, Quebec; 2010.05.03, 11:30-12:30; [**I:IRCM**] (Long networks talk, derived from [**I:BROWNMATH**], including **metamembrane\*** for 1st time. Takes 55' without callgraph sect.)

(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, 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> .