# Understanding Protein Function on a Genome-scale through the Analysis of Molecular Networks
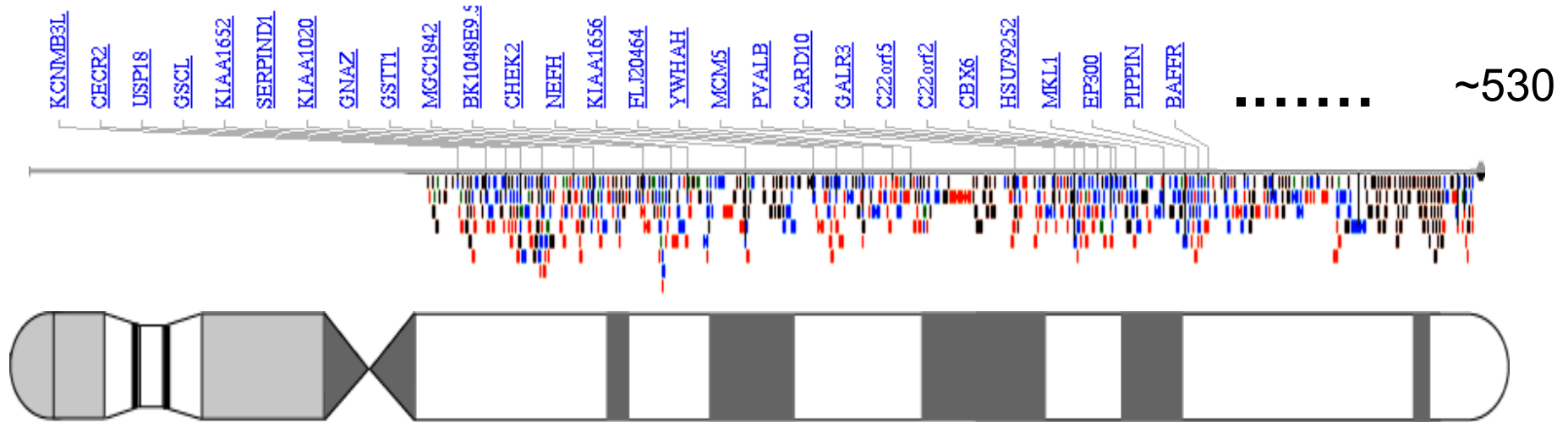
Mark B Gerstein
Yale

**Slides at**

**Lectures.GersteinLab.org**

**(See Last Slide for References & More Info.)**

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



KCNMB3L CECR2 USPI8 GSCL KIAA1652 SERPIND1 KIAA1020 GNAZ GSTT1 MGC1842 BK1048E9.5 CHEK2 NEFH KIAA1656 FLJ20464 YWHAH MCM5 PVALB CARD10 GALR3 C22orf5 C22orf2 CBX6 HSU79252 MKL1 EP300 PIPPIN BAFFR  ....... ~530

- 250 of ~530
originally characterized on chr. 22
[Dunham et al. Nature (1999)]

- >25K Proteins in Entire Human Genome
(with alt. splicing)

# 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
the universal protein knowledgebase

Home > Database > UniProt Protein Viewer

hosted by

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**

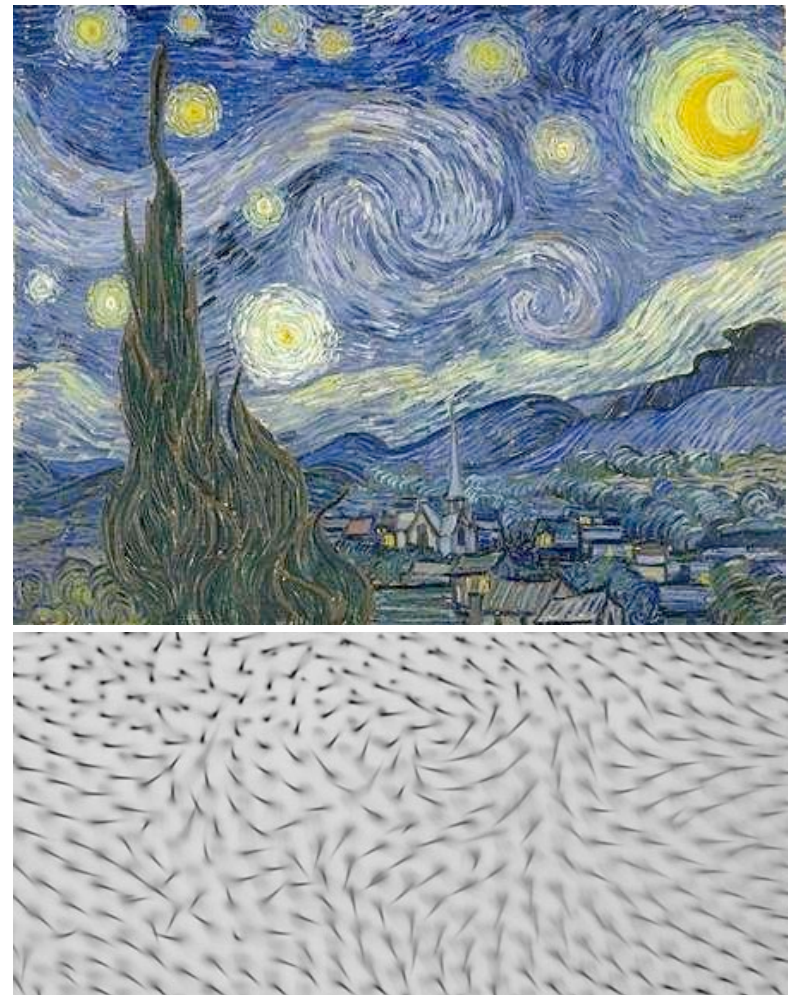| FUNCTION | This protein promotes the GTP-dependent translocation of the nascent protein chain from the A-site to the P-site of the ribosome. |
| SUBCELLULAR LOCATION | Cytoplasmic. |

# 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

# 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

- Fun terms… but do they scale?....
  - ◊ **Starry night** (P Adler, '94)

[Seringhaus et al. GenomeBiology (2008)]

# An Ontology of Naming Pathologies

**M** Explicit meaning

M-scientific    SEMA5A[a]

Not "funny"; usually acronym or concatenation of long descriptive scientific name

M-literal    drop dead[b]

Inherent meaning of words is sufficient to describe gene function in some way; no cultureal knowledge is required

M-embed

Clever reference or allusion. Cutural savvy or other knowledge required to make sense

Literary    malvolio[c]

Acronym    LOV[d]

Historical    yuri[e]

Pop culture    tribbles[f]

**~M** No explicit meaning

~M-outside    kuzbanian[g]

Some outside, non-obvious reason for name

~M-irrel    ring[h]

Irrelevant acronym; not tied to gene function

~M-nr    yippee[i]

Silly or funny names. No relevance to underlying gene function

Multi

**T** Transferred naming system

T-relation    kryptonite and superman

Naming ceases to make sense if names are shuffled among genes

T-norelation    arleekin
valiet
tungus...[k]

Names could be shuffled among genes with no loss of meaning

**P** Problematic relationships

P-clash    PKD1 and lov-1[l]

Analogous genes with very different names

P-confusion    MT-1[m]

Many genes with same name, or many names for one gene

P-defunct    BAF45 and BAF47[n]

Gene named to reflect information later shown to be inaccurate or untrue

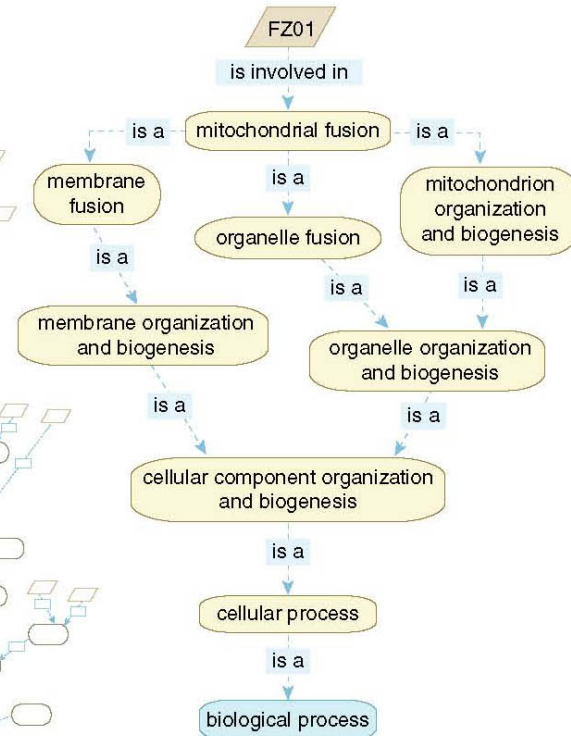**[Seringhaus et al. GenomeBiology (2008)]**

# Gene Name Skew



[Seringhaus et al. GenomeBiology (2008)]

7 Lectures.GersteinLab.org (c) 2009

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



**MIPS (Mewes et al.)**

**GO (Ashburner et al.)**

[Seringhaus & Gerstein, Am. Sci. '08]
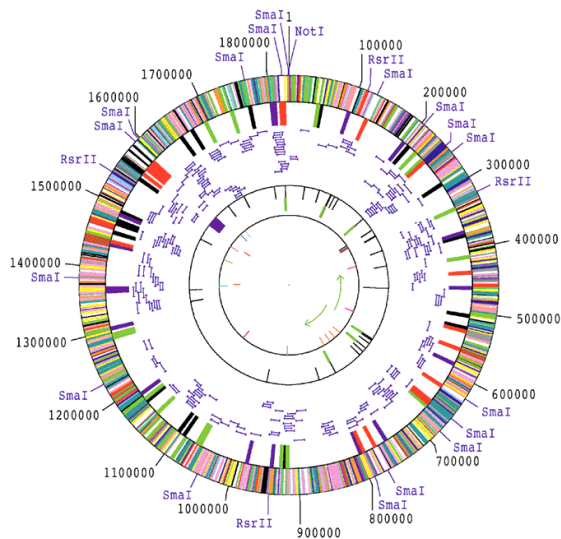
# Networks (Old & New)



Fringe: Vital in boundary formation in developing fly wing.

Numb: mutations impair sensory organs in flies

Notch: with defects, flies develop notches in wings

Itch: linked to itchy skin in mice

Classical KEGG pathway

Same Genes in High-throughput Network

[Seringhaus & Gerstein, Am. Sci. '08]

9 Lectures.GersteinLab.org (c) 2009

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



1D: Complete Genetic Partslist

[Fleischmann et al., Science, 269 :496]



~2D: Bio-molecular Network Wiring Diagram

[Jeong et al. Nature, 41:411]



3D: Detailed structural understanding of cellular machinery

# Networks as a universal language



Internet
[Burch & Cheswick]

Food Web

Electronic
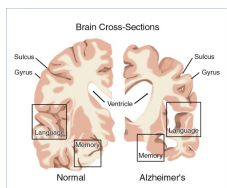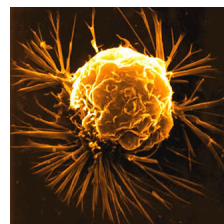Circuit

Neural Network
[Cajal]

Disease
Spread
[Krebs]

Protein
Interactions
[Barabasi]

Albert-László
Barabási

L I N K E D

The New Science
of Networks

How Everything is Connected to Everything Else
and What it Means for Science, Business
and Everyday Life

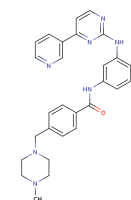Social Network

# Network pathology & pharmacology

Breast Cancer
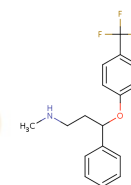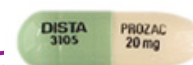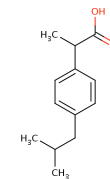
Alzheimer's Disease

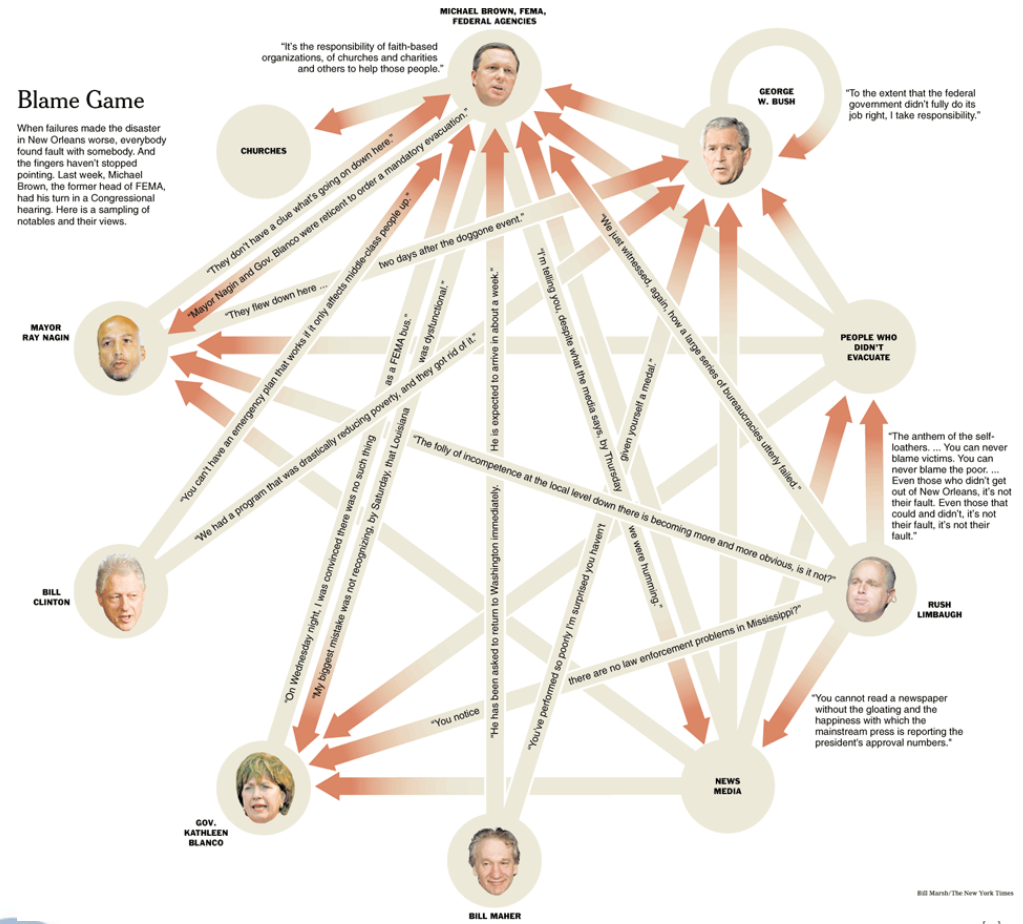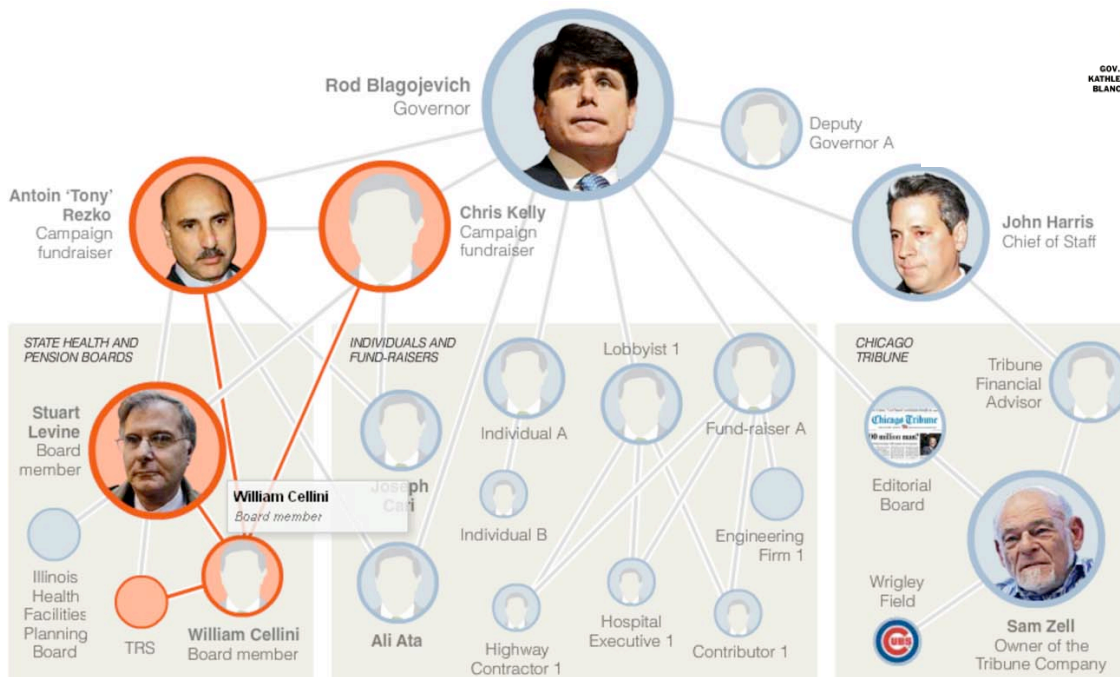Parkinson's Disease

Multiple Sclerosis

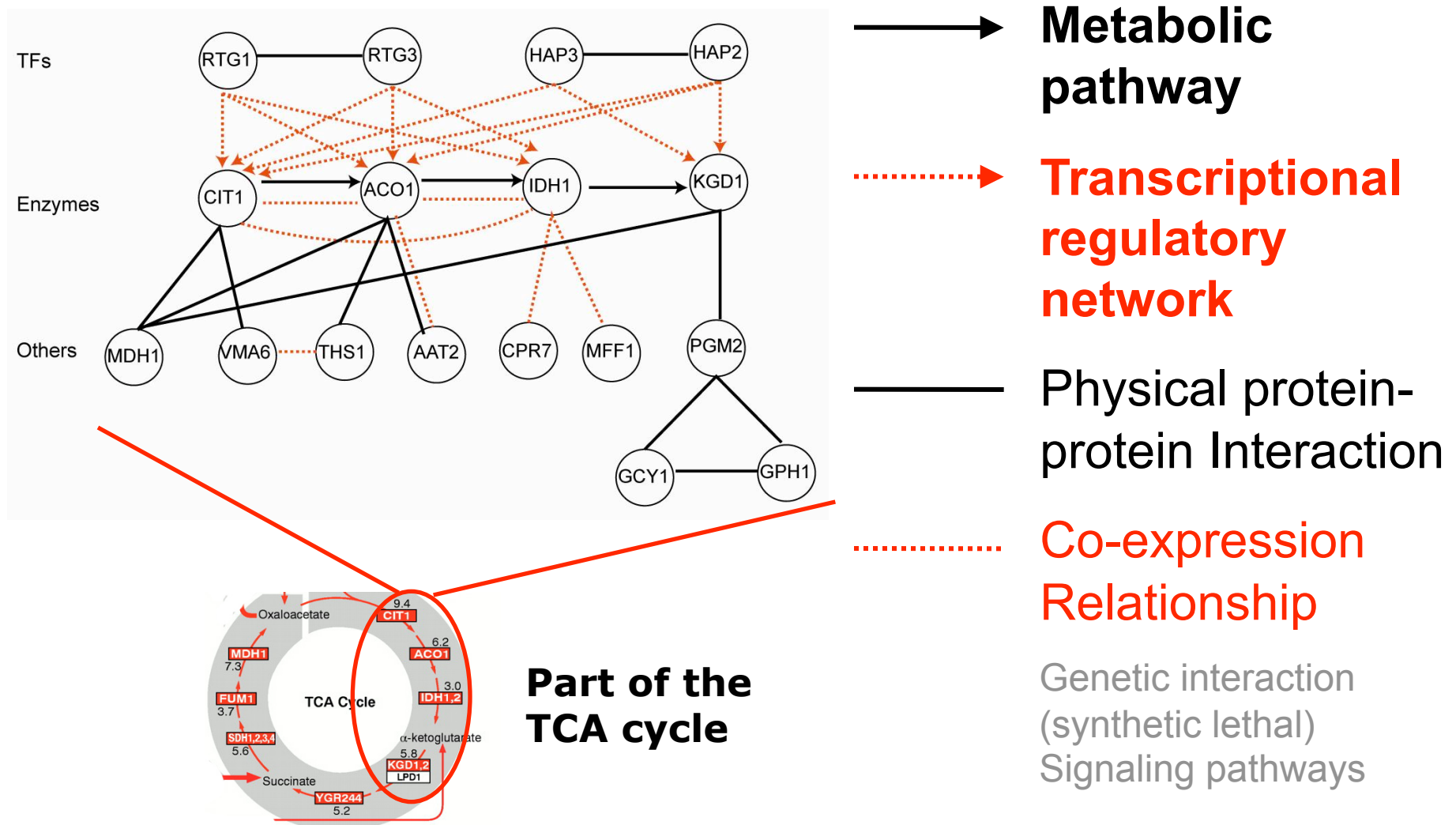**Interactome networks**

**[Adapted from H Yu]**

# Using the position in networks to describe function



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

# Combining networks forms an ideal way of integrating diverse information



Part of the TCA cycle

→ **Metabolic pathway**

⇢ **Transcriptional regulatory network**

— Physical protein-protein Interaction

⋯ Co-expression Relationship

Genetic interaction (synthetic lethal) Signaling pathways

14 Lectures.GersteinLab.org (c) 2009
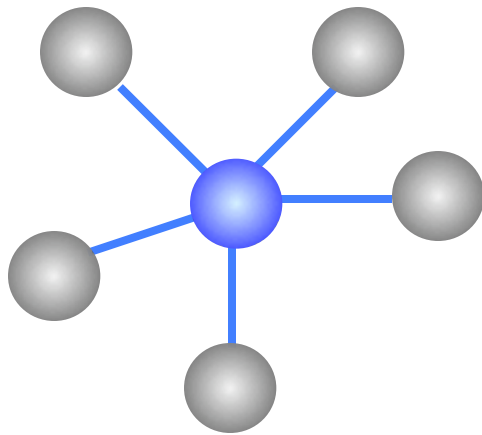
# Outline: Molecular Networks

- Why Networks?
- Network Structure: Key Positions
  - ◊ Hubs & Bottlenecks
  - ◊ Tops of a Hierachy
- Networks, Variation & the Environment
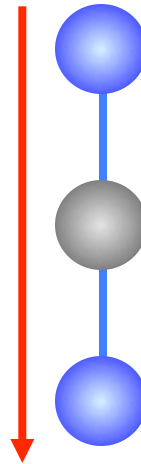  - ◊ Which pathways change most with the environment
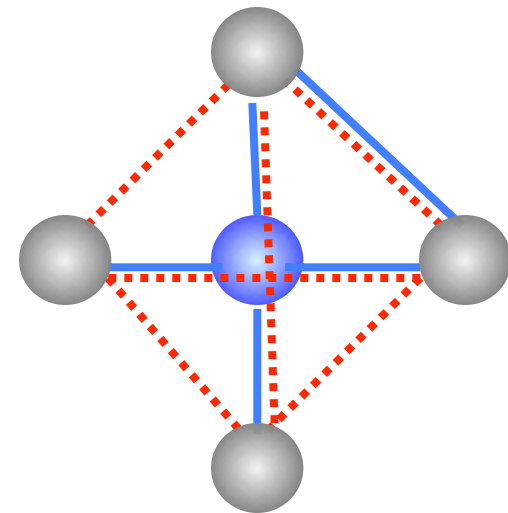
# Global topological measures

Indicate the gross topological structure of the network
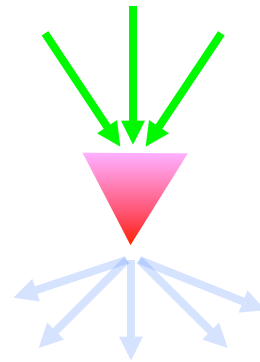


Degree ($K$)

5

Path length ($L$)

2

Clustering coefficient ($C$)

1/6
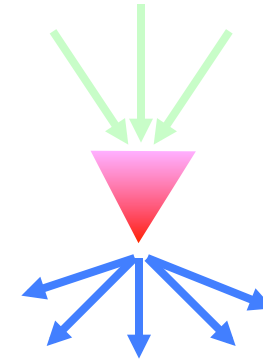
Interaction and expression networks are ***undirected***

[Barabasi]
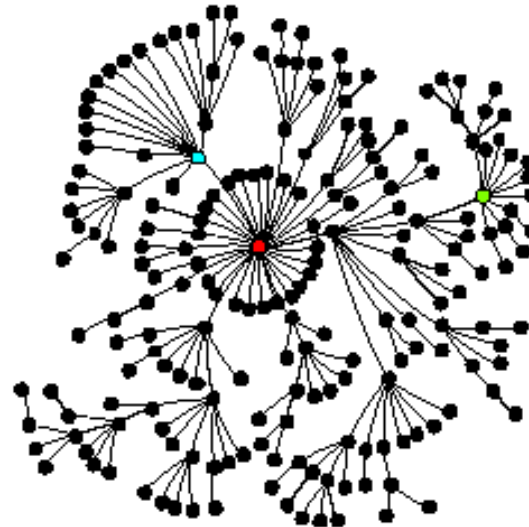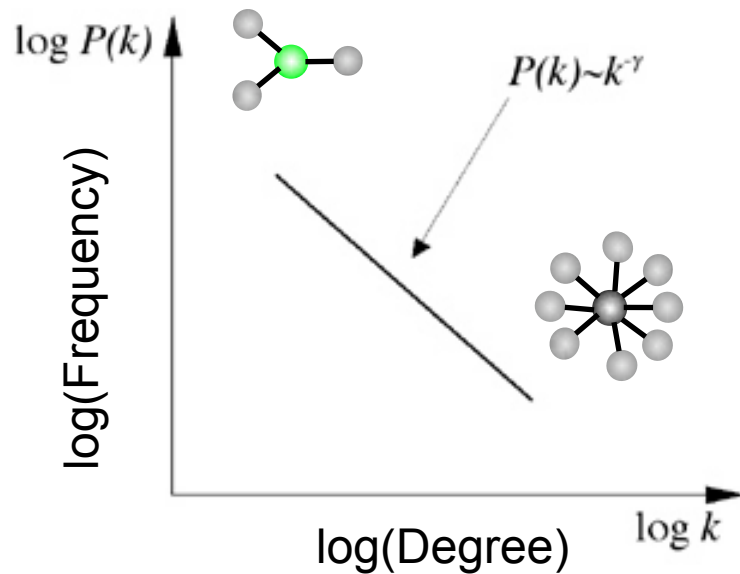
# Global topological measures for directed networks

TFs

Targets

In-degree
3

Out-degree
5

Regulatory and metabolic networks are **directed**

# Scale-free networks

Power-law distribution



$P(k) \sim k^{-\gamma}$

log P(k)

log(Frequency)

log(Degree)

log k

***Hubs*** dictate the structure of the network
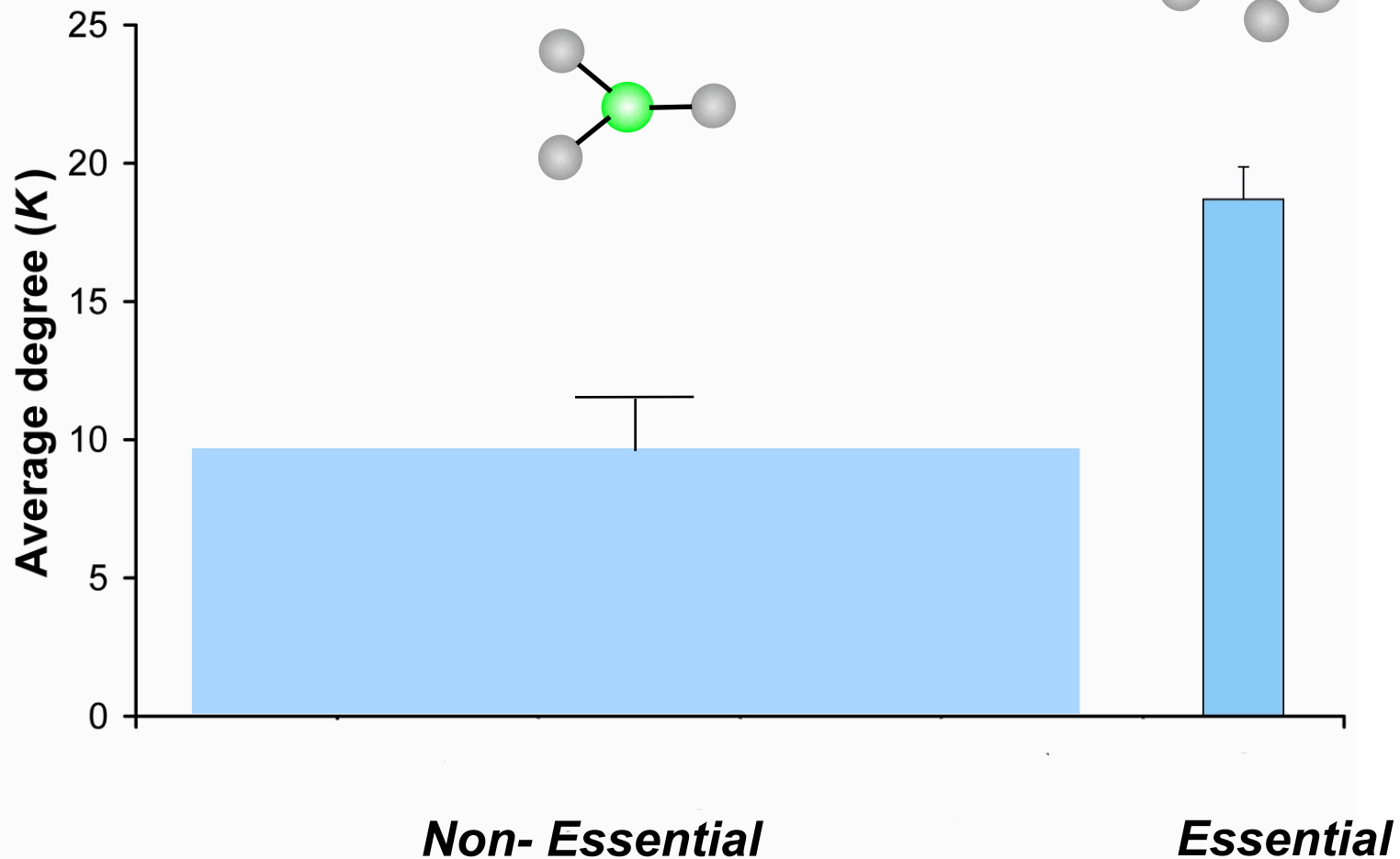
**[Barabasi]**

# Hubs tend to be Essential

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

[Lauffenburger, Barabasi]



**"hubbiness"**

Average degree ($K$)

25
20
15
10
5
0

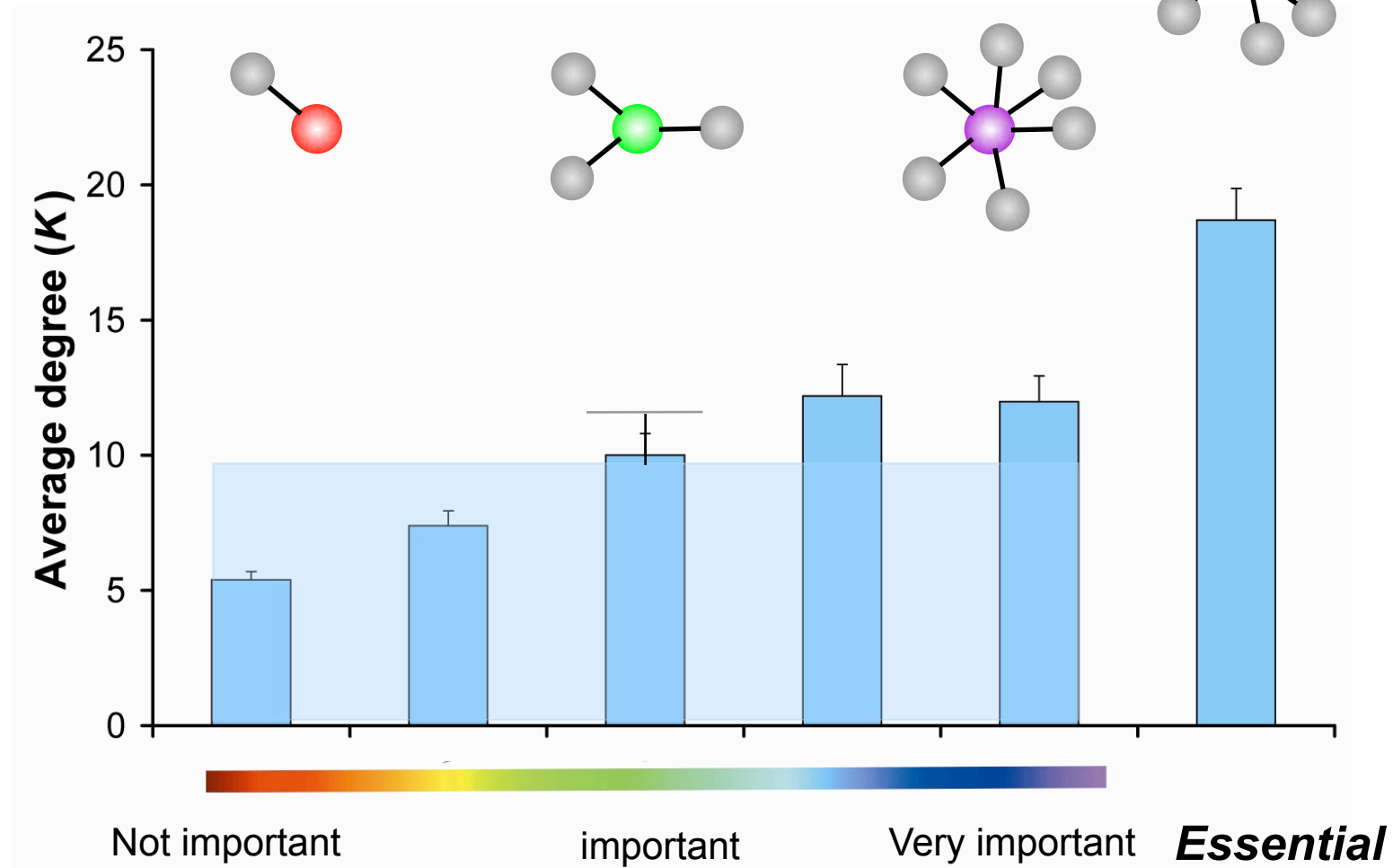*Non- Essential*          *Essential*

# 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"
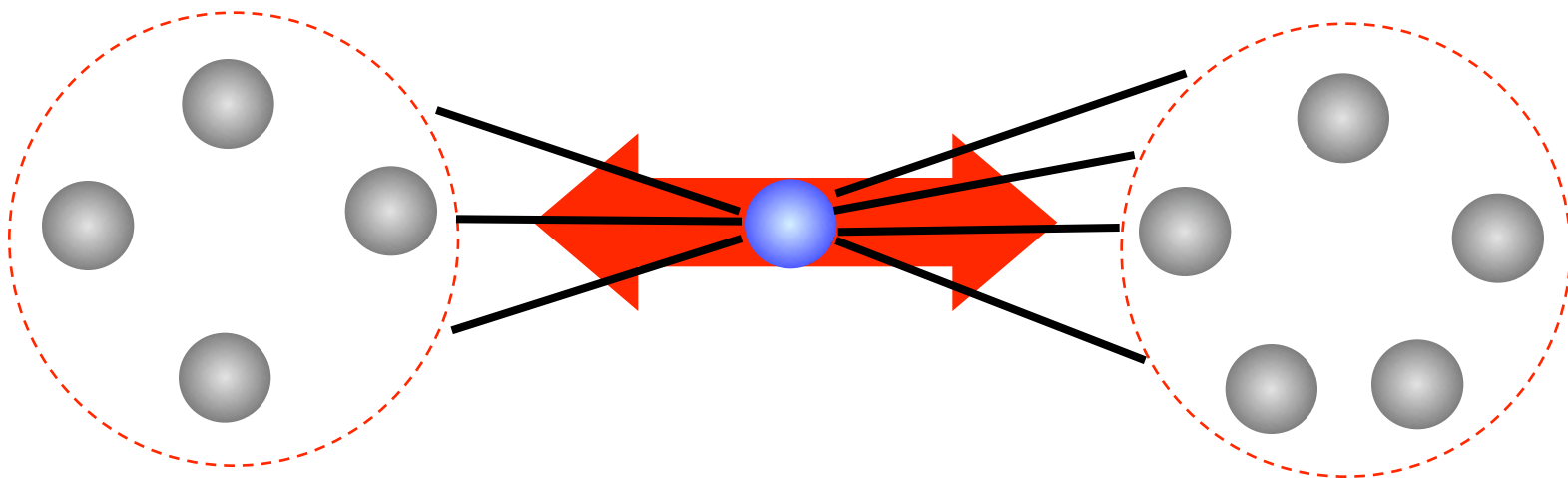
# Another measure of Centrality: Betweenness centrality

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**
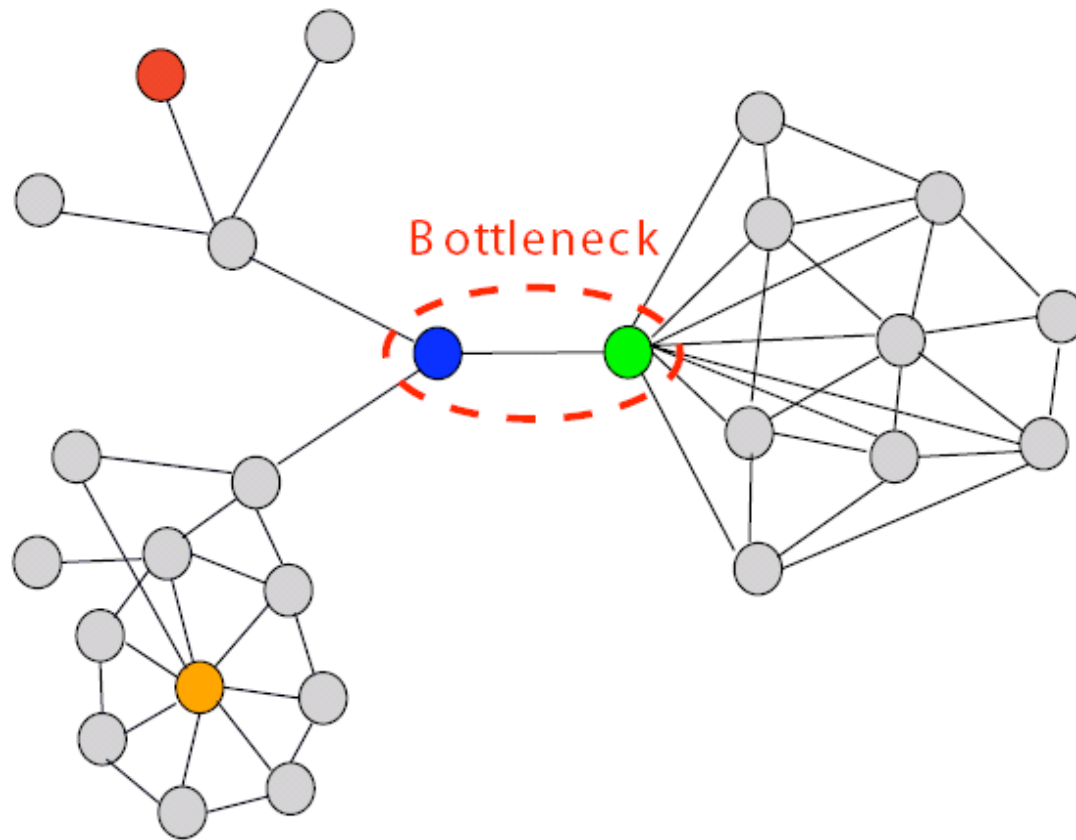


George Washington Bridge
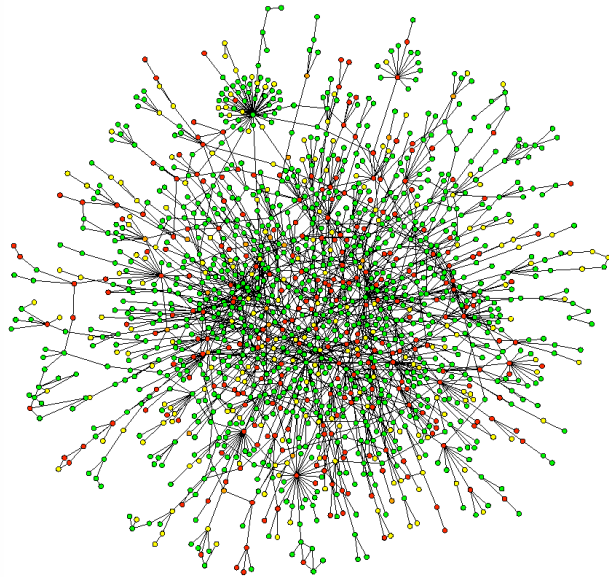
Bottleneck

# Bottlenecks & Hubs

Legend:

- 🟢 Hub-bottleneck **node**
- 🔵 Non-hub-bottleneck **node**
- 🟠 Hub-non-bottleneck **node**
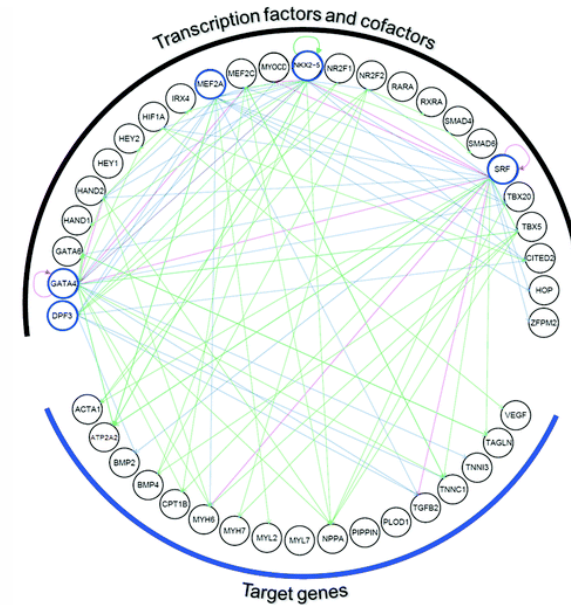- 🔴 Non-hub-non-bottleneck **node**

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

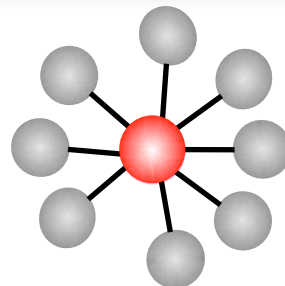# Different Interactome networks


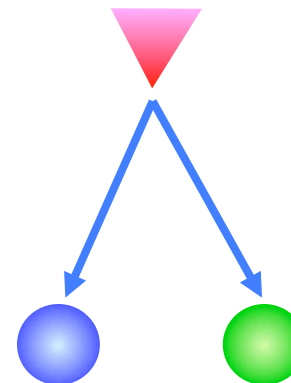
**Interaction networks**



**Regulatory networks**

**Undirected**

**Directed**
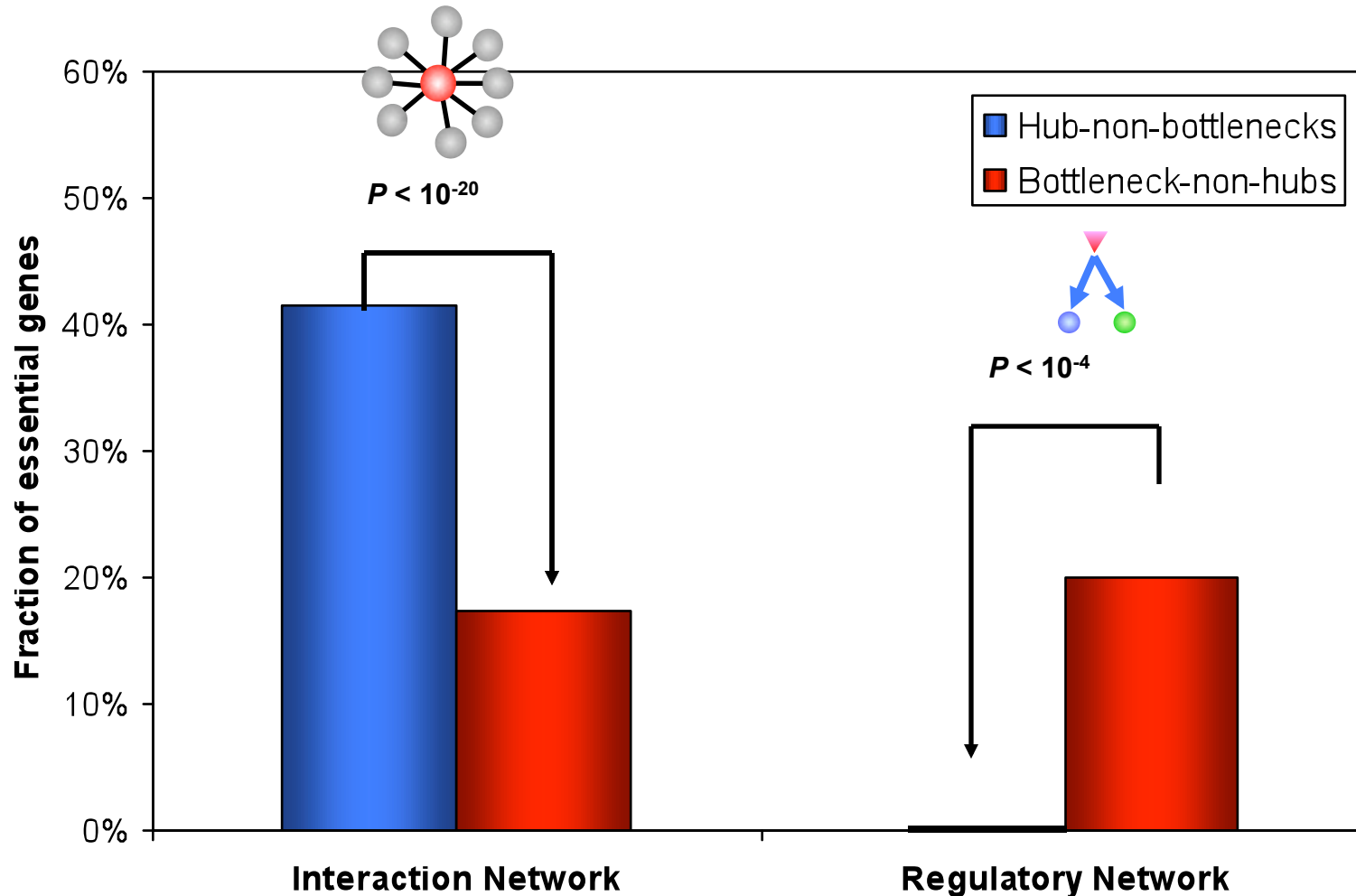
[Toenjes, *et al*, *Mol. BioSyst.* (2008)]

[Jeong *et al*, *Nature* (2001)]

# Bottlenecks are what matters in regulatory networks



[Yu *et al.*, *PLoS Comput Biol* (2007)]

# Signaling transduction pathways are directed



[Xianglin Shi ]

[Yu *et al.*, *PLoS Comput Biol* (2007]

# Bottlenecks in signaling pathways are important



P < 0.02

Fraction of essential genes

30%

20%

10%

0%

Bottlenecks in
Signal Pathways

Bottlenecks in
Interaction network

Random

TLR/IL-1
TNF
Metals, Stresses
TIR
MyD88
IRAK
RIP1
TRADD
AF6
Ub(63)
TRAF2
FADD
TAK1
NIK
p
IKKα/β/γ
p105
p100
β-catenin
κB
NF-κB
p
p
Cytokines,
Adhesion molecules,
Cox-2, iNOS,
Bcl-xl, cIAPs,
Cyclin D1,......

[Xianglin Shi ]

[Yu *et al.*, *PLoS Comput Biol* (2007)]

# Outline: Molecular Networks

- Why Networks?
- Network Structure:
  Key Positions
    ◊ Hubs & Bottlenecks
    ◊ Tops of a Hierachy
- Networks, Variation & the
  Environment
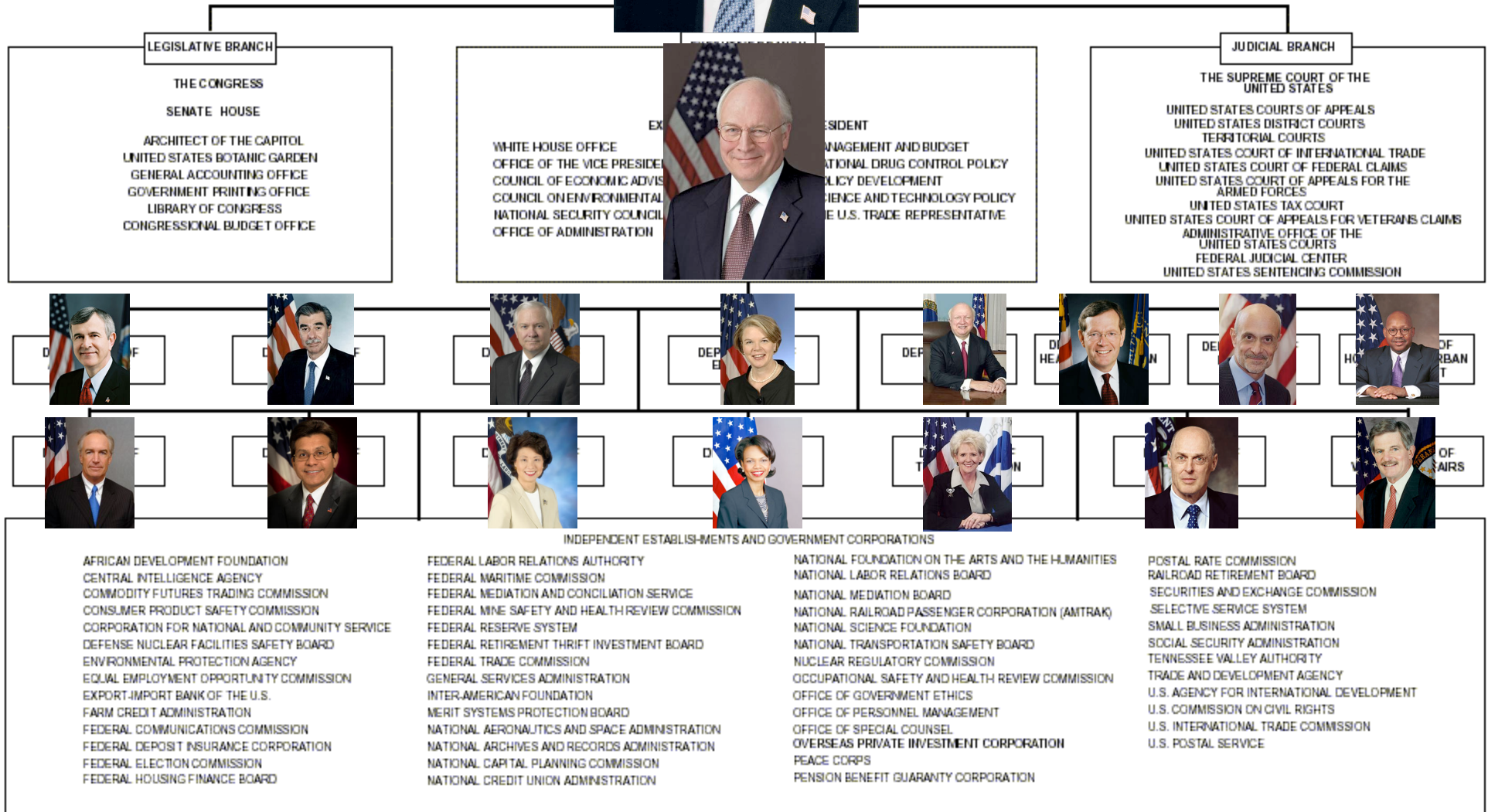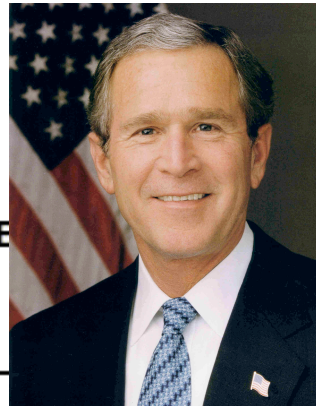    ◊ Which pathways change
      most with the
      environment

# Social Hierarchy



THE GOVE... UNITED STATES

THE CONGRESS

SENATE  HOUSE

ARCHITECT OF THE CAPITOL
UNITED STATES BOTANIC GARDEN
GENERAL ACCOUNTING OFFICE
GOVERNMENT PRINTING OFFICE
LIBRARY OF CONGRESS
CONGRESSIONAL BUDGET OFFICE

EX... ...RESIDENT

WHITE HOUSE OFFICE
OFFICE OF THE VICE PRESIDE...
COUNCIL OF ECONOMIC ADVIS...
COUNCIL ON ENVIRONMENTAL...
NATIONAL SECURITY COUNCIL
OFFICE OF ADMINISTRATION

...ANAGEMENT AND BUDGET
...ATIONAL DRUG CONTROL POLICY
...OLICY DEVELOPMENT
...IENCE AND TECHNOLOGY POLICY
...E U.S. TRADE REPRESENTATIVE

JUDICIAL BRANCH

THE SUPREME COURT OF THE
UNITED STATES

UNITED STATES COURTS OF APPEALS
UNITED STATES DISTRICT COURTS
TERRITORIAL COURTS
UNITED STATES COURT OF INTERNATIONAL TRADE
UNITED STATES COURT OF FEDERAL CLAIMS
UNITED STATES COURT OF APPEALS FOR THE
ARMED FORCES
UNITED STATES TAX COURT
UNITED STATES COURT OF APPEALS FOR VETERANS CLAIMS
ADMINISTRATIVE OFFICE OF THE
UNITED STATES COURTS
FEDERAL JUDICIAL CENTER
UNITED STATES SENTENCING COMMISSION

INDEPENDENT ESTABLISHMENTS AND GOVERNMENT CORPORATIONS

AFRICAN DEVELOPMENT FOUNDATION
CENTRAL INTELLIGENCE AGENCY
COMMODITY FUTURES TRADING COMMISSION
CONSUMER PRODUCT SAFETY COMMISSION
CORPORATION FOR NATIONAL AND COMMUNITY SERVICE
DEFENSE NUCLEAR FACILITIES SAFETY BOARD
ENVIRONMENTAL PROTECTION AGENCY
EQUAL EMPLOYMENT OPPORTUNITY COMMISSION
EXPORT-IMPORT BANK OF THE U.S.
FARM CREDIT ADMINISTRATION
FEDERAL COMMUNICATIONS COMMISSION
FEDERAL DEPOSIT INSURANCE CORPORATION
FEDERAL ELECTION COMMISSION
FEDERAL HOUSING FINANCE BOARD

FEDERAL LABOR RELATIONS AUTHORITY
FEDERAL MARITIME COMMISSION
FEDERAL MEDIATION AND CONCILIATION SERVICE
FEDERAL MINE SAFETY AND HEALTH REVIEW COMMISSION
FEDERAL RESERVE SYSTEM
FEDERAL RETIREMENT THRIFT INVESTMENT BOARD
FEDERAL TRADE COMMISSION
GENERAL SERVICES ADMINISTRATION
INTER-AMERICAN FOUNDATION
MERIT SYSTEMS PROTECTION BOARD
NATIONAL AERONAUTICS AND SPACE ADMINISTRATION
NATIONAL ARCHIVES AND RECORDS ADMINISTRATION
NATIONAL CAPITAL PLANNING COMMISSION
NATIONAL CREDIT UNION ADMINISTRATION

NATIONAL FOUNDATION ON THE ARTS AND THE HUMANITIES
NATIONAL LABOR RELATIONS BOARD
NATIONAL MEDIATION BOARD
NATIONAL RAILROAD PASSENGER CORPORATION (AMTRAK)
NATIONAL SCIENCE FOUNDATION
NATIONAL TRANSPORTATION SAFETY BOARD
NUCLEAR REGULATORY COMMISSION
OCCUPATIONAL SAFETY AND HEALTH REVIEW COMMISSION
OFFICE OF GOVERNMENT ETHICS
OFFICE OF PERSONNEL MANAGEMENT
OFFICE OF SPECIAL COUNSEL
OVERSEAS PRIVATE INVESTMENT CORPORATION
PEACE CORPS
PENSION BENEFIT GUARANTY CORPORATION

POSTAL RATE COMMISSION
RAILROAD RETIREMENT BOARD
SECURITIES AND EXCHANGE COMMISSION
SELECTIVE SERVICE SYSTEM
SMALL BUSINESS ADMINISTRATION
SOCIAL SECURITY ADMINISTRATION
TENNESSEE VALLEY AUTHORITY
TRADE AND DEVELOPMENT AGENCY
U.S. AGENCY FOR INTERNATIONAL DEVELOPMENT
U.S. COMMISSION ON CIVIL RIGHTS
U.S. INTERNATIONAL TRADE COMMISSION
U.S. POSTAL SERVICE
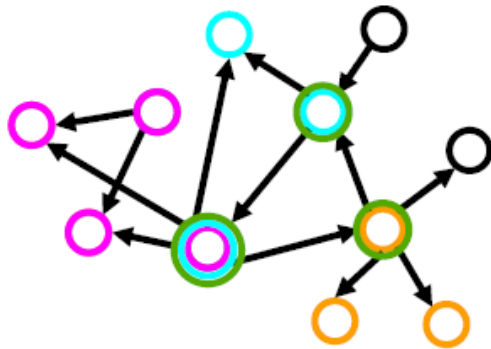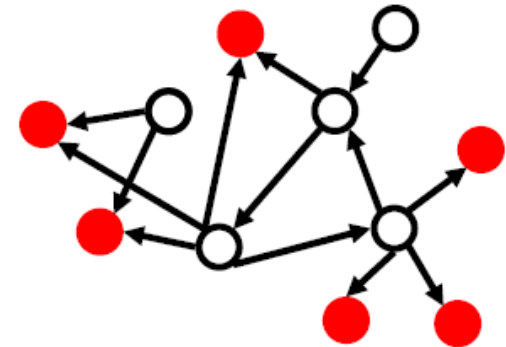
# Determination of "Level" in Regulatory Network Hierarchy with Breadth-first Search
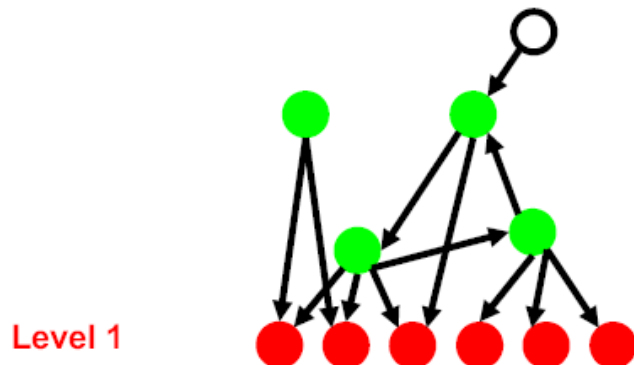


I. Example network with all 4 motifs
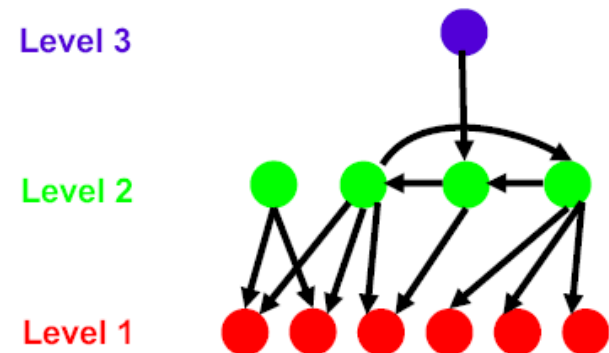
II. Finding terminal nodes (Red)

III. Finding mid-level nodes (Green)

Level 1

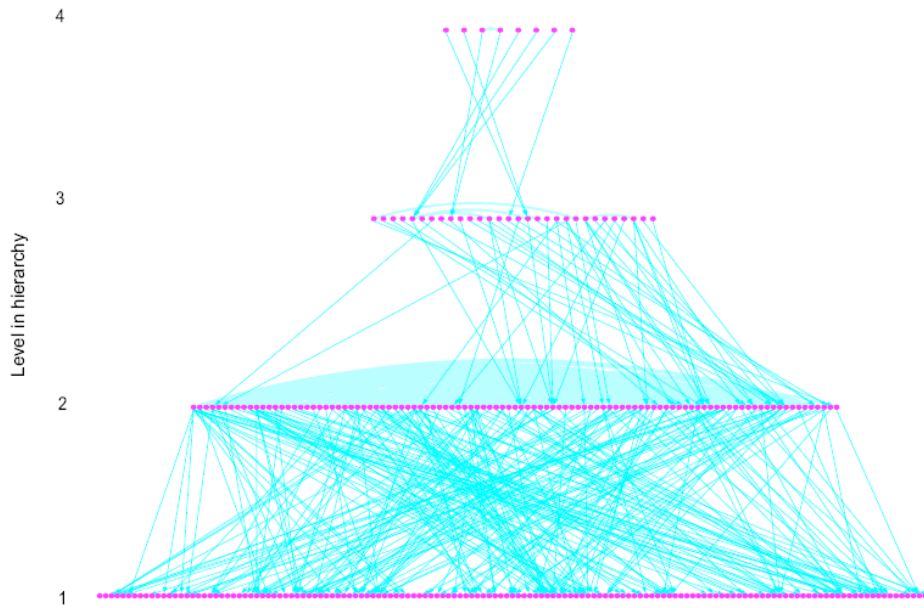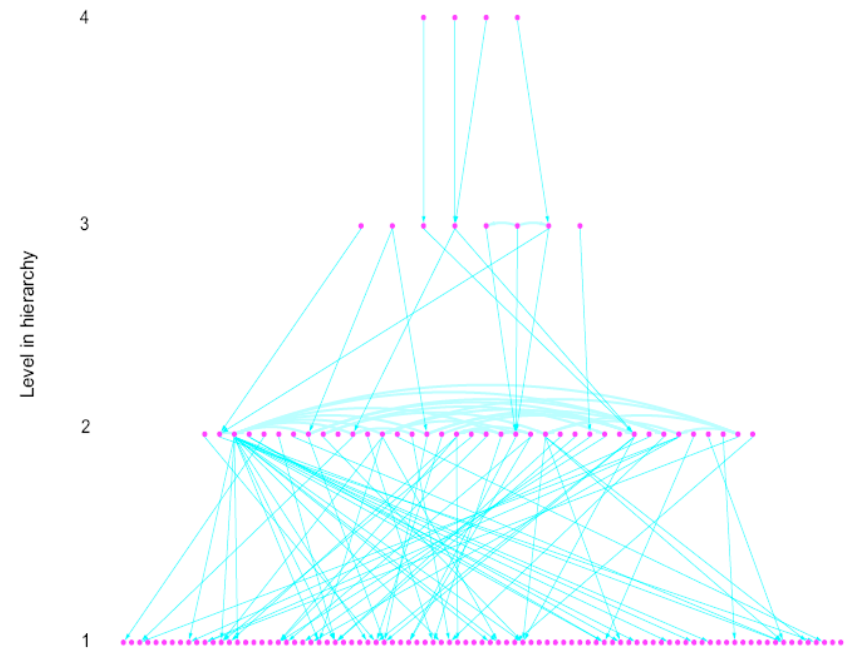IV. Finding top-most nodes (Blue)

Level 3

Level 2

Level 1

[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)]**
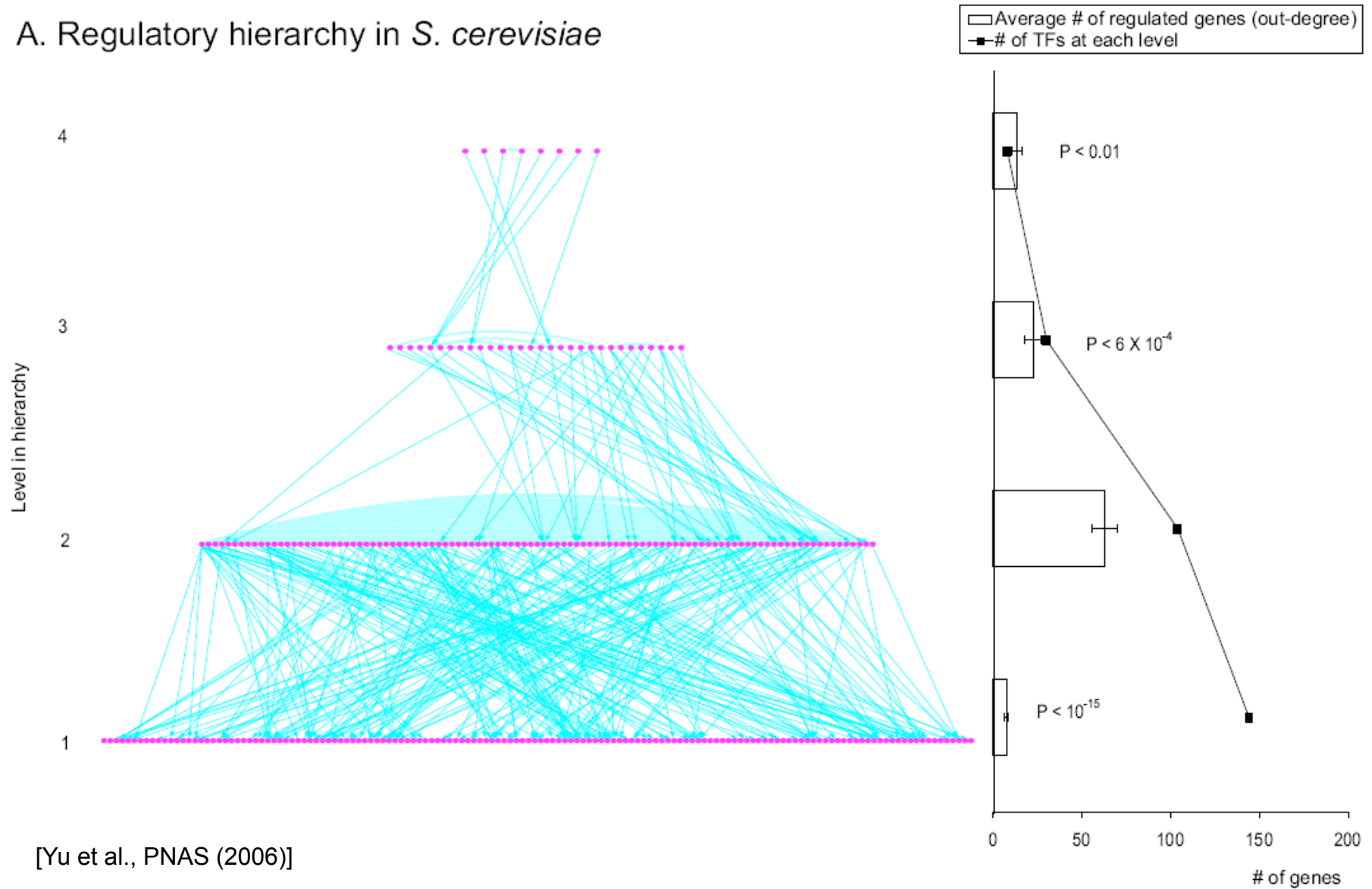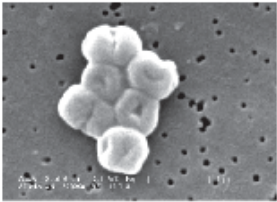
# Example of Path Through Regulatory Network



Expression of MOT3 is activated by heme and oxygen. Mot3 in turn activates the expression of NOT5 and GCN4, mid-level hubs. GCN4 activates two specific bottom-level TFs, Put3 and Uga3, which trigger the expression of enzymes in proline and nitrogen utilization.

[Yu et al., PNAS (2006)]

# Yeast Regulatory Hierarchy: the Middle-managers Rule



A. Regulatory hierarchy in *S. cerevisiae*

Average # of regulated genes (out-degree)
# of TFs at each level

P < 0.01

P < 6 X 10⁻⁴

P < 10⁻¹⁵

# of genes

Level in hierarchy

[Yu et al., PNAS (2006)]

33 Lectures.GersteinLab.org (c) 2009

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



B. Governmental hierarchy of a representive city (Macao)

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



Average betweenness at each level

P < 10⁻⁴

P < 10⁻¹¹

Level in Hierarchy

Average betweenness (x1000)

[Yu et al., PNAS (2006)]

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



[Yu et al., PNAS (2006)]

# Outline: Molecular Networks



- Why Networks?
- Network Structure: Key Positions
  - ◊ Hubs & Bottlenecks
  - ◊ Tops of a Hierachy
- Networks, Variation & the Environment
  - ◊ Which pathways change most with the environment

# What is metagenomics?

# Global Ocean Survey Statistics (GOS)



6.25 GB of data
7.7M Reads
 1 million CPU hours
to process

Rusch, et al., PLOS Biology 2007

**Pathway Sequences (Community Function)**

Metabolic Pathways

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

**Environmental Features**

Environmental Metadata

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

**READS → PROTEIN FAMILIES → PATHWAYS**

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

$P_1 = f_1 + f_2 + f_3$

$P_2 = f_4 + f_5 + f_6$

**PATHWAYS**

⬤ ▪

SITES

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

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

# 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



[ Gianoulis et al., PNAS (in press, 2009) ]

# Canonical Correlation Analysis: Simultaneous weighting



UPI = a **GRE** + b + c **GPA**

GPI = a' + b' + c'

[ Gianoulis et al., PNAS (in press, 2009) ]

# Canonical Correlation Analysis: Simultaneous weighting



| Score | # of papers published |
|-------|----------------------|
| GRE | |

| Undergraduate Performance Index (UPI) | Graduate School Performance Index (GPI) |
|---|---|
| GRE GPA | |

| Environmental Features | Metabolic Pathways |
|---|---|
| Temp      etc | Photosynthesis      etc |
| Chlorophyll | Lipid Metabolism |

[ Gianoulis et al., PNAS (in press, 2009) ]

# 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 $Y = \{y_1, y_2, ..., y_m\}$

$$C = \begin{matrix} \Sigma_X & \Sigma_{X,Y} \\ \Sigma_Y & \Sigma_{Y,X} \end{matrix}$$

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

[ Gianoulis et al., PNAS (in press, 2009) ]

Strength of Pathway co-variation with environment

CCA structural correlation

0    0.3    1

Environmentally invariant          Environmentally variant

CCA structural correlation

0    0.3

[ Gianoulis et al., PNAS (in press, 2009) ]

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



[ Gianoulis et al., PNAS (in press, 2009) ]

# Conclusion #2: Outer Membrane components vary the environment



CCA structural correlation
0    0.3    1

[ Gianoulis et al., PNAS (in press, 2009) ]

# 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?

[ Gianoulis et al., PNAS (in press, 2009) ]

# Biosensors:
# Beyond Canaries in a Coal Mine



[ Gianoulis et al., PNAS (in press, 2009) ]

# Conclusions

- Developing Standardized Descriptions of Protein Function
- Gene Naming
- Betweenness is an important global network statistic
  - ◊ Bottlenecks are more correlated with essentiality than hubs in regulatory networks
- Regulatory Network Hierarchies
  - ◊ Middle managers dominate, sitting at info. flow bottlenecks
  - ◊ Paradox of influence and essentiality
  - ◊ Topmost proteins sit at center of interaction network

# 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.

**TopNet** — *Topology of Networks*

# – 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

## TopNet.GersteinLab.org
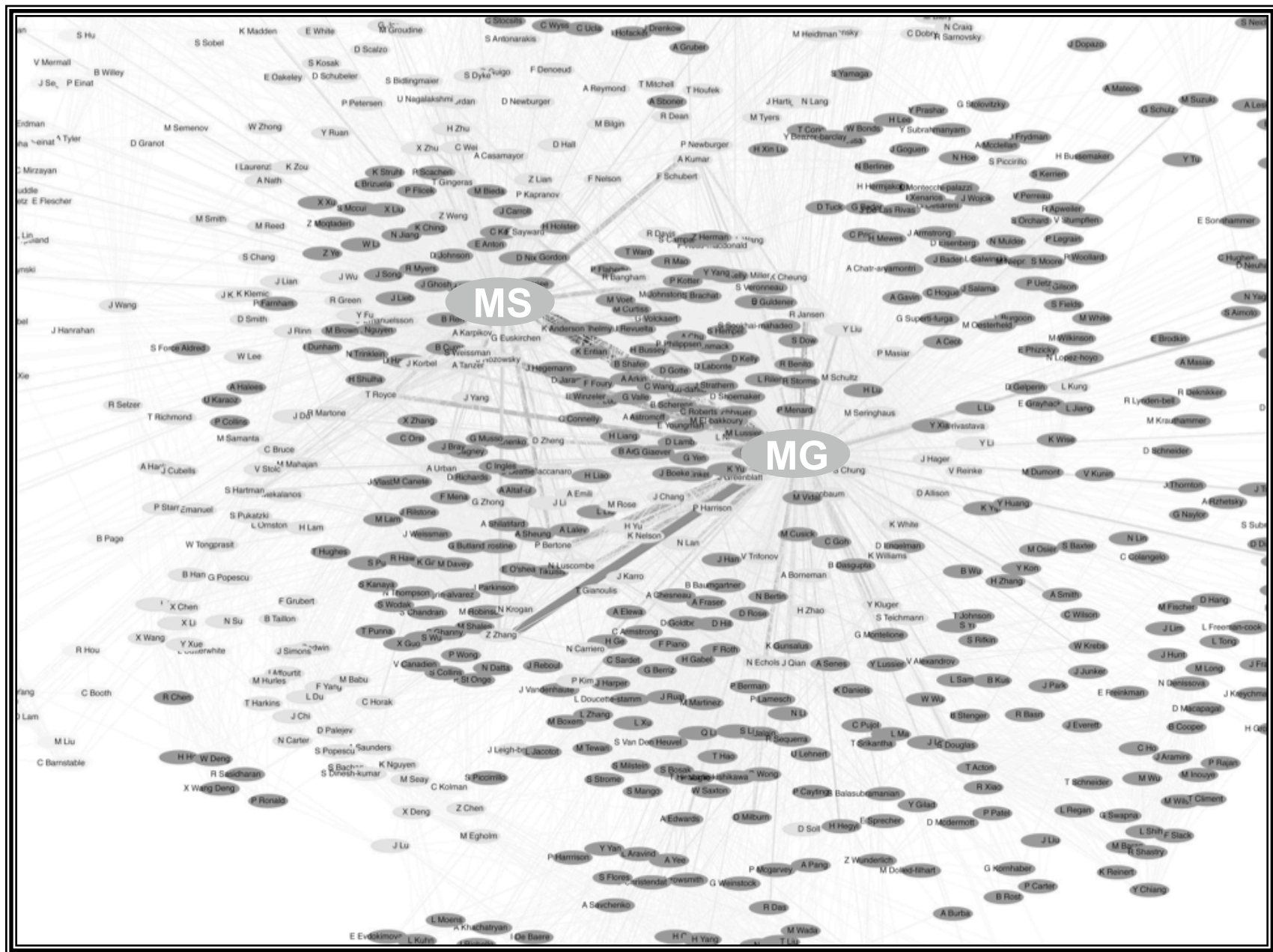
# Acknowledgements

## TopNet.GersteinLab.org
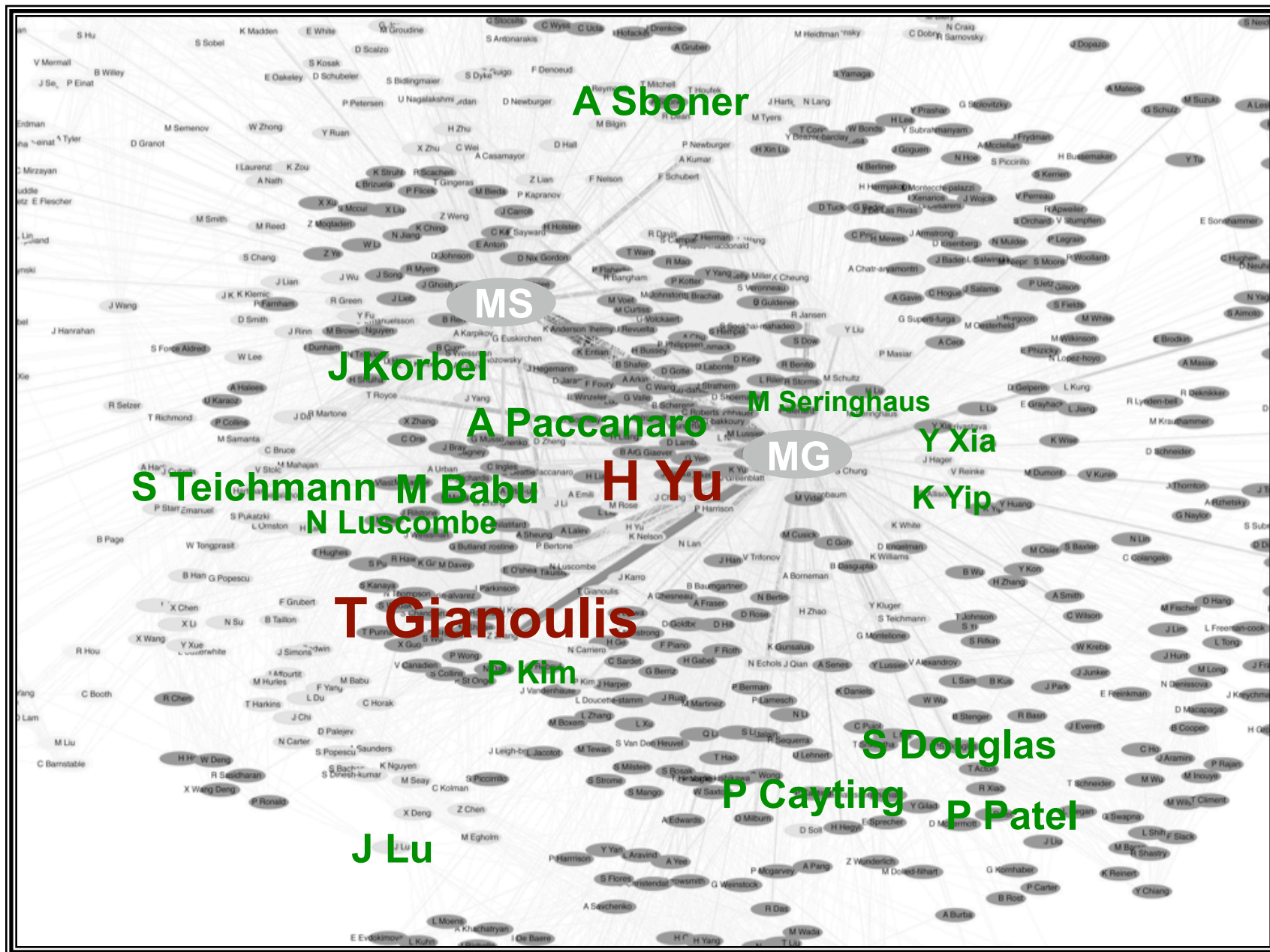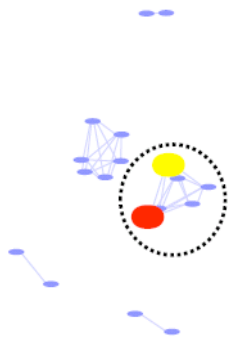
# Acknowledgements

## TopNet.GersteinLab.org

**P Bork, J Raes**

Job opportunities currently

for postdocs & students

**RNAi:
Birth of a
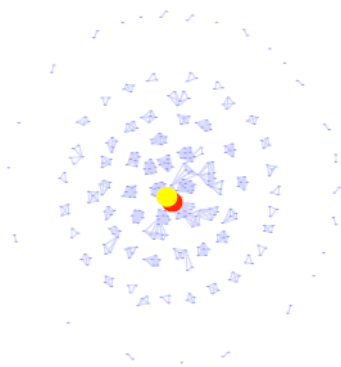Field in
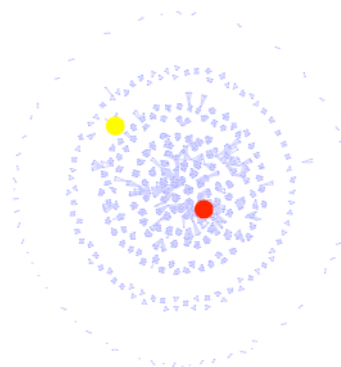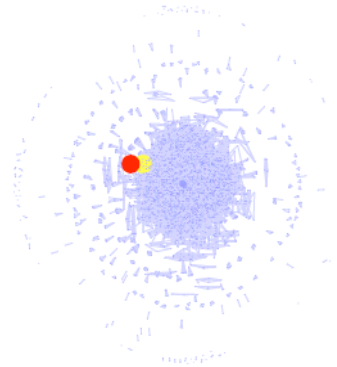the
Literature
Culmin
-ating in
the 2006
Nobel**

1998

1999

2000

2001

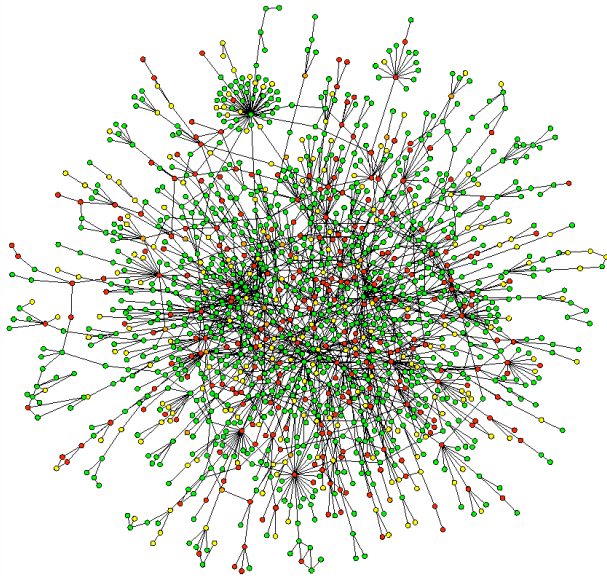2002

2003

● Andrew Fire    ● Craig Mello

# **Extra**

# Types of Networks



Regulatory networks
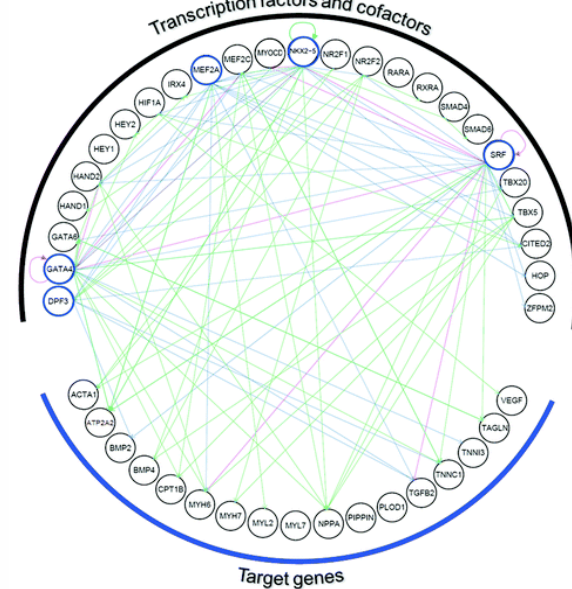


Interaction networks
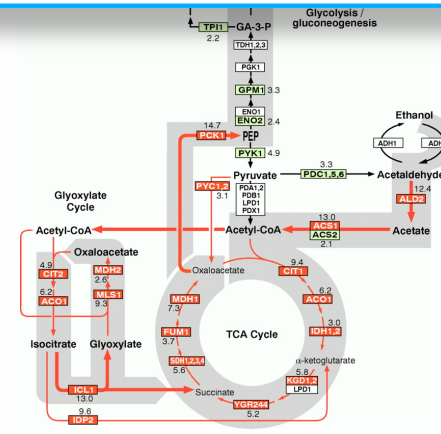
**Nodes: proteins or genes**
**Edges: interactions**

[Horak, et al, Genes & Development, 16:3017-3033]

[DeRisi, Iyer, and Brown, Science, 278:680-686]

[Jeong et al, Nature, 41:411]



Metabolic networks

# More Information on this Talk

**TITLE**: Understanding Protein Function on a Genome-scale through the Analysis of Molecular Networks

**SUBJECT**: Networks

**DESCRIPTION**:
National Academy of Engineering, Meeting at Columbia U, 2009.04.14, 14:00-14:30; [I:NAECU] (Short networks talk, incl. the following topics:
why networks w. amsci*, funnygene*, bottleneck*, nethierarchy*, metagenomics*, tyna* + topnet*, & pubnet* . Fits into 30' w. 5' questions. PPT works on mac & PC and has many photos w. EXIF tag kwtimewemet .)

(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 )