# Biological Network Analysis

Mark B Gerstein
Yale

slides at

**Lectures.GersteinLab.org**

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

# GersteinLab.org Research Overview: Bioinformatics



- **Genome Annotation**
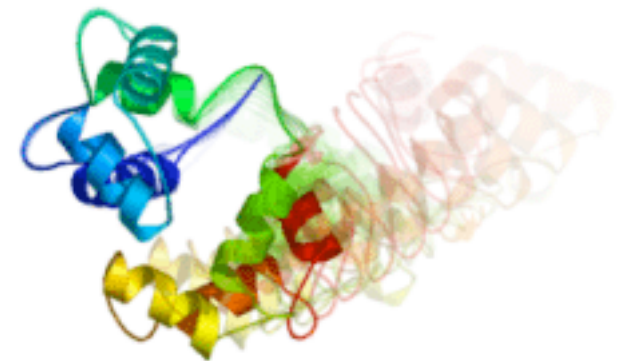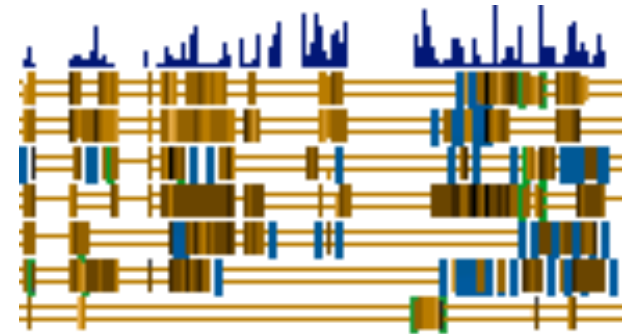  - ◊ <u>Characterizing non-coding regions</u> of the genome, focusing on protein fossils and novel RNAs
    (Pseudogene.org + GenomeTech.GersteinLab.org)
  - ◊ <u>Personal Genomics</u> – esp. related to SVs

- **Molecular Networks**
  - ◊ Using molecular networks to integrate & mine functional genomics information and describe genefunction on a large-scale
    (Networks.GersteinLab.org)



- **Macromolecular Motions**
  - ◊ Analyzing select populations of 3D-structures in detail, trying to understand their flexibility in terms of packing
    (MolMovDB.org)

# 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.

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

- Fundamental complexities
  - ◊ Role Conflation:
    molecular, cellular, phenotypic

  - ◊ Often >2 proteins/function

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

    - 

[HSP from Craig et al, Rev Physiol Biochem Pharmacol (2006) 156:1 ; Terms from Seringhaus et al. GenomeBiology (2008)]
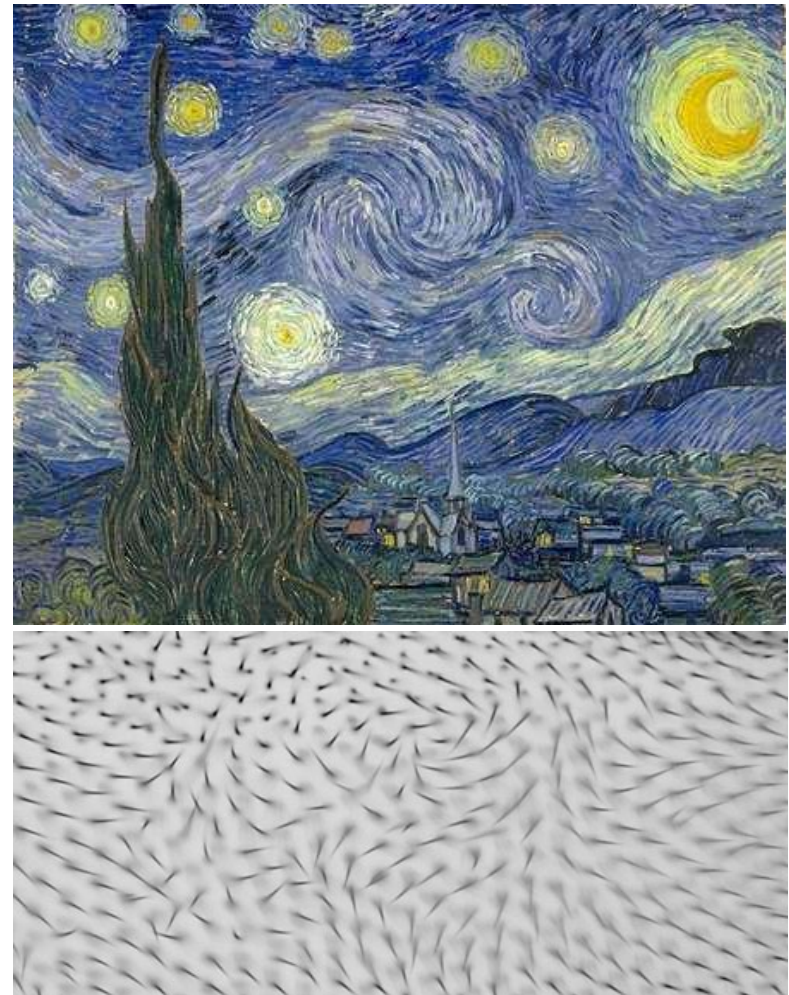
# Some obvious issues in scaling single molecule definition to a genomic scale

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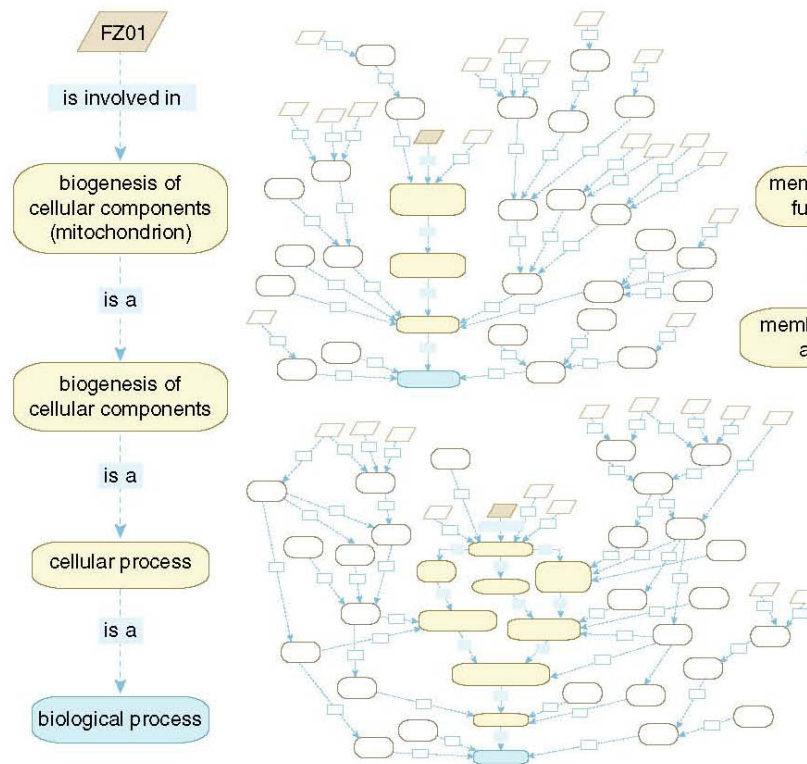
- Fun terms… but do they scale?....
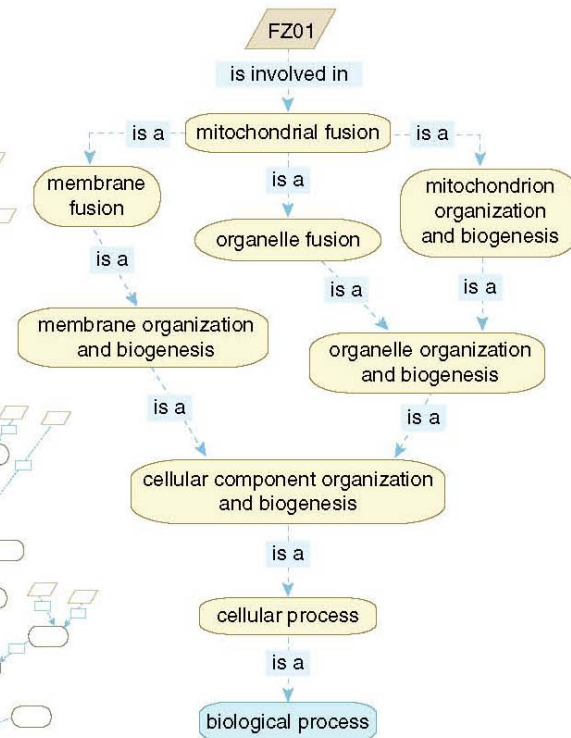  - ◊ **Starry night** (P Adler, '94)



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

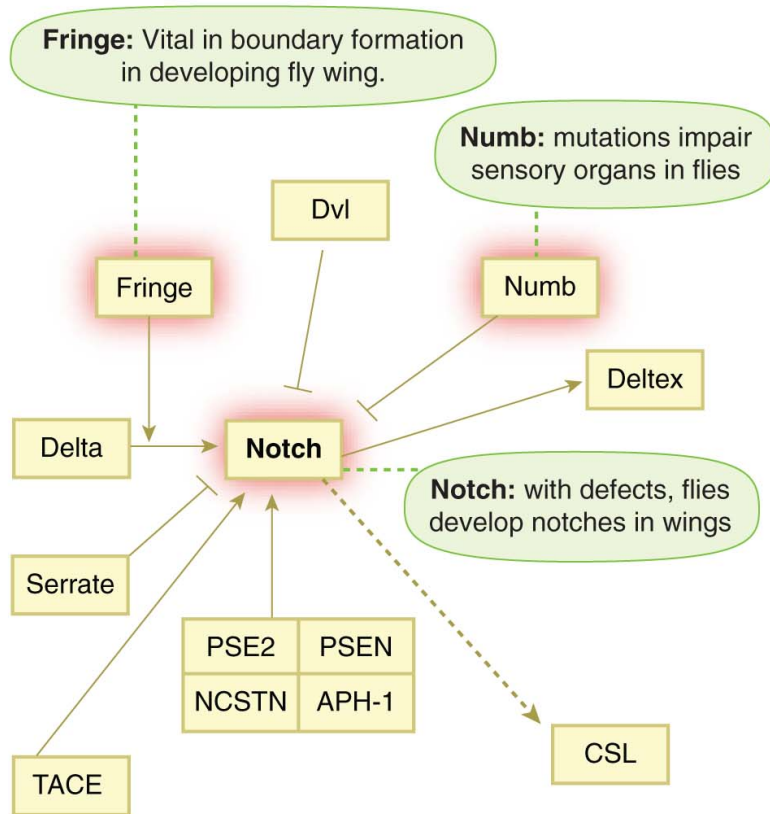# Hierarchies & DAGs of controlled-vocab terms but still have issues...
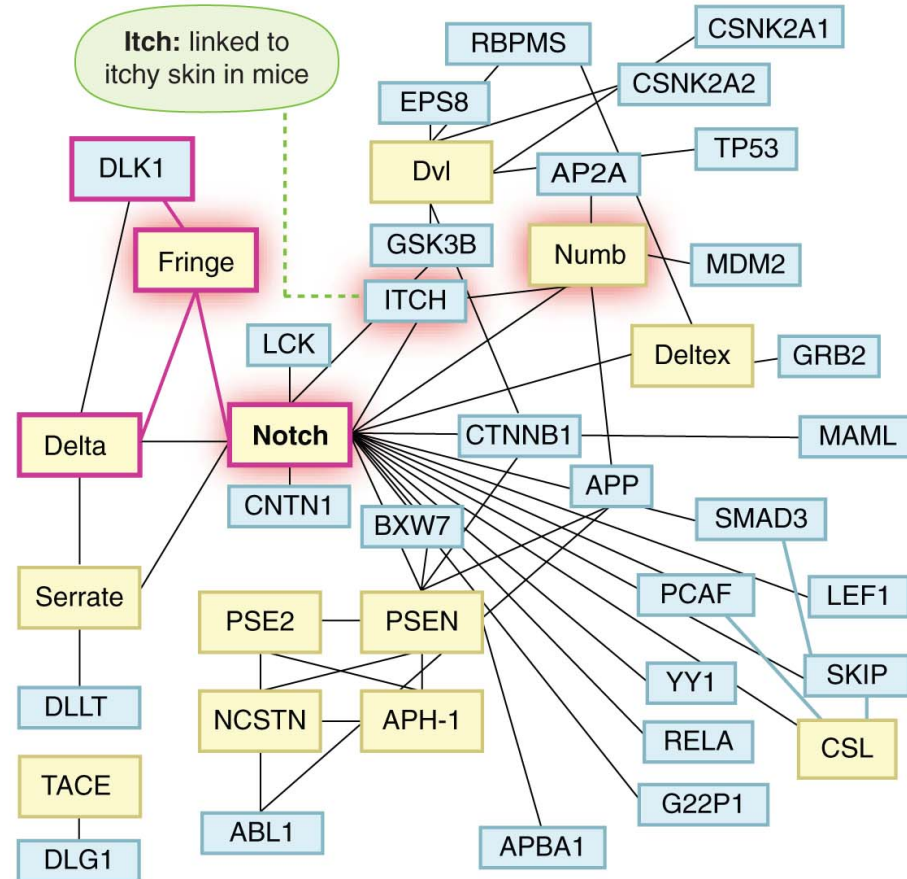


**MIPS (Mewes et al.)**
          **GO (Ashburner et al.)**

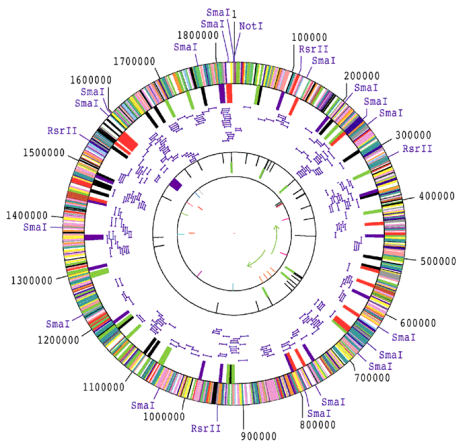[Seringhaus & Gerstein, Am. Sci. '08]

# Networks (Old & New)



Classical KEGG pathway

Same Genes in High-throughput Network

[Seringhaus & Gerstein, Am. Sci. '08]
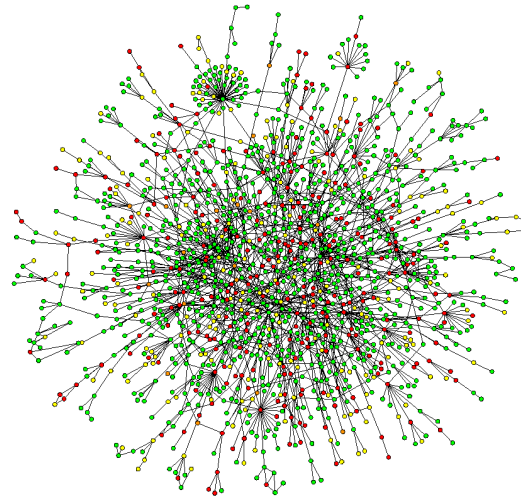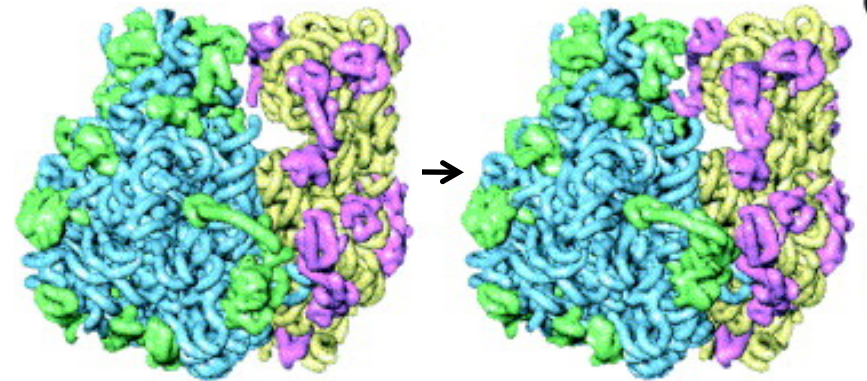
# Networks occupy a midway point in terms of level of understanding
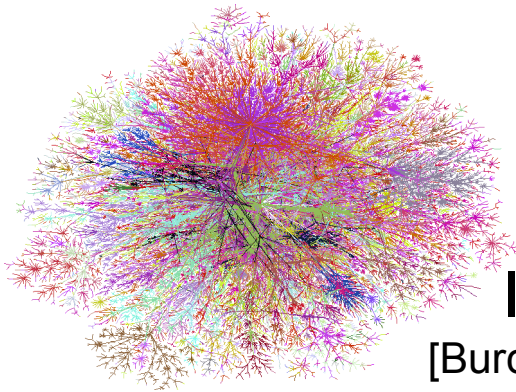


1D: Complete
Genetic Partslist

~2D: Bio-molecular
Network
Wiring Diagram

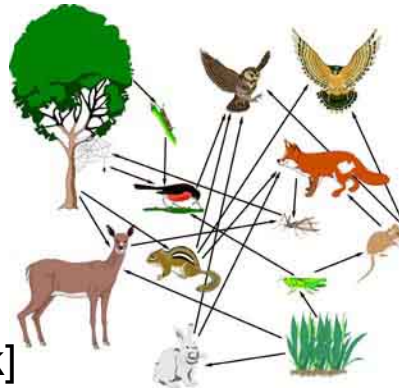3D and 4D:
Detailed structural understanding
of cellular machinery
(e.g. ribosome in different
functional states)

**[Fleischmann et al., Science, 269 :496]**        **[Jeong et al. Nature, 41:411]**        **[Chiu et al. Trends in Cell Biol, 16:144]**
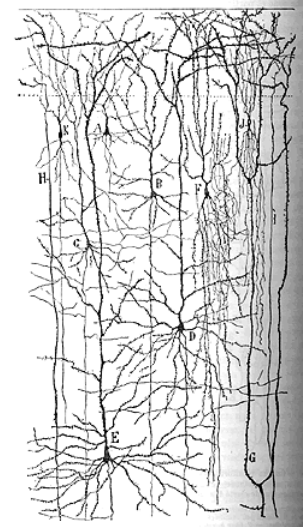
# Networks as a universal language

Internet
[Burch & Cheswick]

Food Web

Electronic
Circuit

Neural Network
[Cajal]

Disease
Spread
[Krebs]

Albert-László
Barabási

L I N K E D

The New Science
of Networks

Protein
Interactions
[Barabasi]

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

Social Network

# Using the position in networks to describe function

## 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



Part of the TCA cycle

→ **Metabolic pathway**

┄┄► **Transcriptional regulatory network**

── Physical protein-protein Interaction

┄┄┄ Co-expression Relationship

Genetic interaction (synthetic lethal) Signaling pathways

- Why Networks?
- Generating Networks
  - ◊ Processing Protein Chips
    `(yeast & human nets)`
  - ◊ Propagating Known Information
    `(yeast ppi)`
- Central Points in Networks
  - ◊ Hubs & Bottlenecks
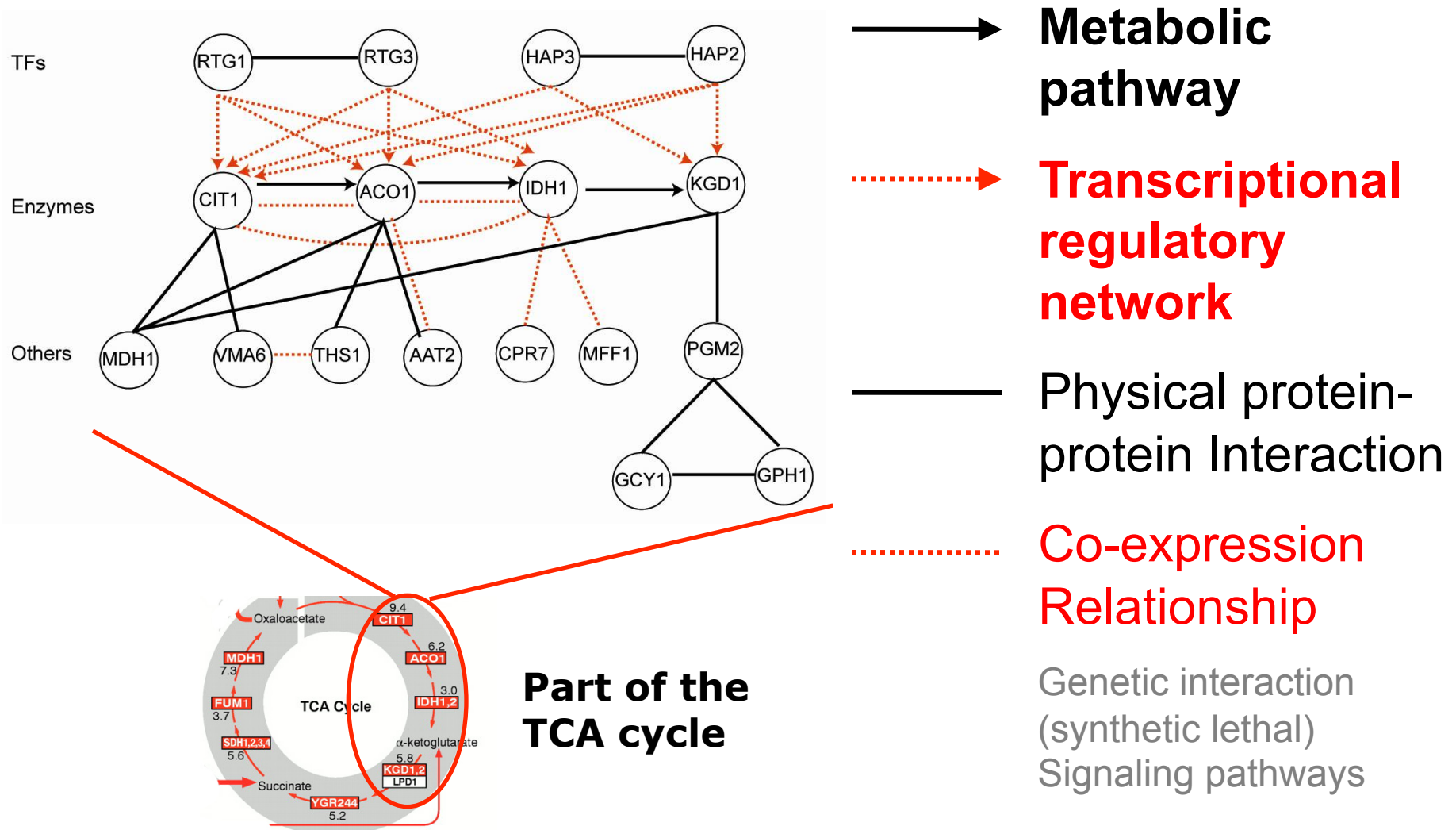    `(yeast ppi & reg. net)`
  - ◊ Tops of Heirarchies
    `(yeast reg. net)`
  - ◊ Identified by score
    `(human miRNA-targ. net)`
- Dynamics of Networks
  - ◊ Across environments
    `(prokaryote metab. pathways)`
- Protein Networks & Variation
  `(human ppi & miRNA-targ. net)`
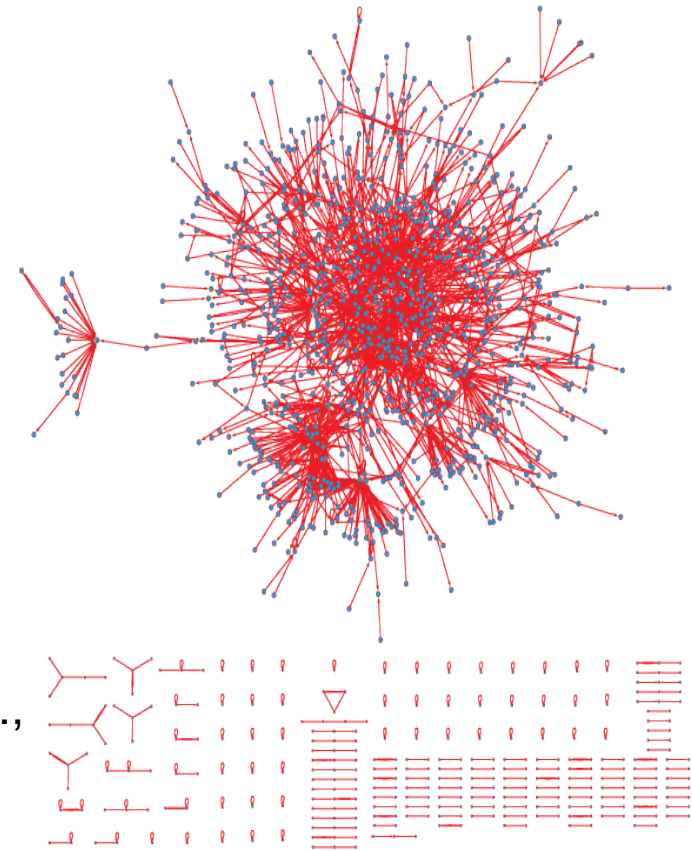
# Outline: Molecular Networks

# Example: yeast PPI network



Actual size:

◊ ~6,000 nodes
  → Computational cost: ~18M pairs

◊ Estimated ~15,000 edges
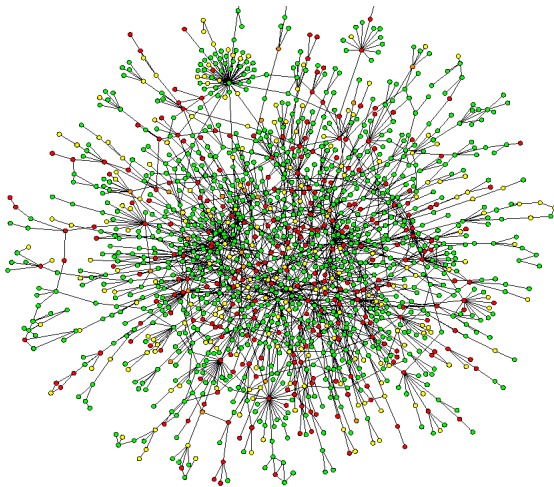  → Sparseness: 0.08% of all pairs (Yu et al., 2008)

Known interactions:

◊ Small-scale experiments: accurate but few
  → Overfitting: ~5,000 in BioGRID, involving ~2,300 proteins

◊ Large-scale experiments: abundant but noisy
  → Noise: false +ve/-ve for yeast two-hybrid data up to 45% and 90% (Huang et al., 2007)
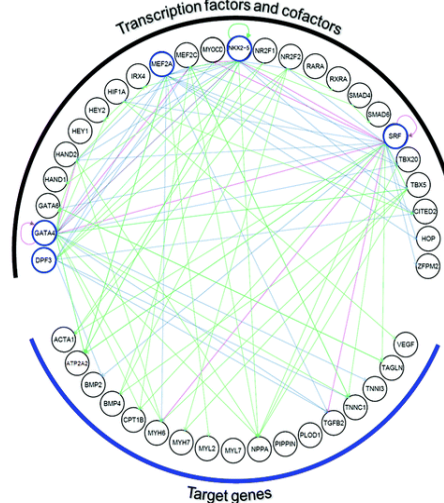
# Different Types of Molecular Networks



**Protein-protein Interaction networks**



**TF-target-gene Regulatory networks**



**Undirected**



**Metabolic pathway networks**



**miRNA-target networks**



**Directed**

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

# Generating Networks

**How do we construct large molecular networks.**
**From processing high-throughput protein array data?**

# Protein Networks from Processing Protein Chip Data



~6000 yeast proteins on a chip, Zhu et al. Science ('01)

- Array functional proteins on a chip

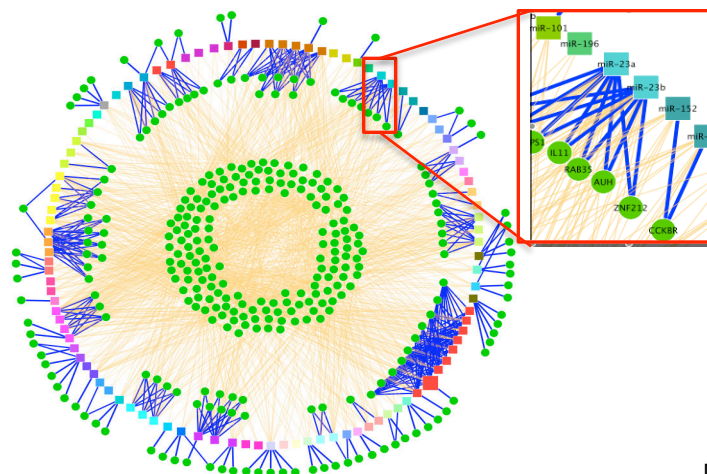- Readout can show presence of proteins in sera (via autoantibodies), small mol. interactions, enzymatic activity, & **protein interactions**

- Technical issues in processing protein chips similar but not identical to those for DNA chips

  ◊ Hybridization v protein binding

  ◊ Background correction & denoising, Normalizing across chips & replicates, Calling "hits"

  ◊ ProCAT (Zhu et al., GenomeBiology, '06) & **RLM (Sboner et al., J Proteome Sci. '09)**



4200 phosphorylations involving 1325 proteins, Ptacek et al. Nature ('05)

# Signal Distribution & Metrics

[Sboner et al., J Proteome Sci. '09]



Protein Chip Sig. Intensity Distribution from different applied sera
(**NEG**ative & with **POS**itive sera)



Representative DNA Chip Sig. Dist.

**Goal:** Decreasing variation betw. replicates (both inter- & intra- array), measured by **CV**, & increasing **separability (S)** betw. known positive & negative samples



intra-array variability
cv

Control Proteins
constant across array

$$CV = \frac{\sigma}{\mu}$$

"Positive spots"

"Positive"

"Normal"

inter-array variability
cv

inter-array variability
cv

$$S = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$$

sample separability
Fisher's signal-to-noise ratio S

# RLM Normalization, how it compares?

NORMALIZATION

– **Global**

• A single scaling factors

– **Quantile**

• Signals are normalized robustly according to the quantiles of a reference distribution

– **Robust Linear Model** **RLM**

$$y_{ijkr} = \alpha_i + \beta_j + \tau_k + \varepsilon_{ijkr}$$

$\alpha_i$    **Slide-effect (inter slide)**

$\beta_j$    **Sub-array effect (intra slide)**
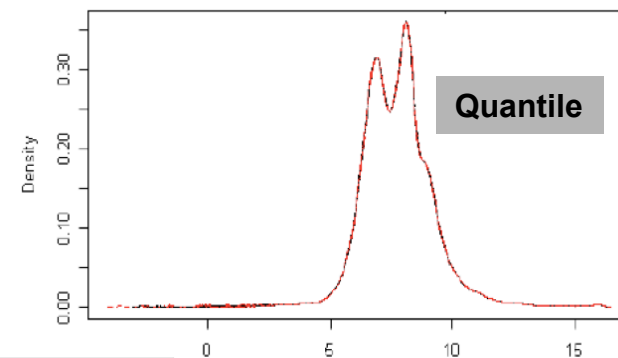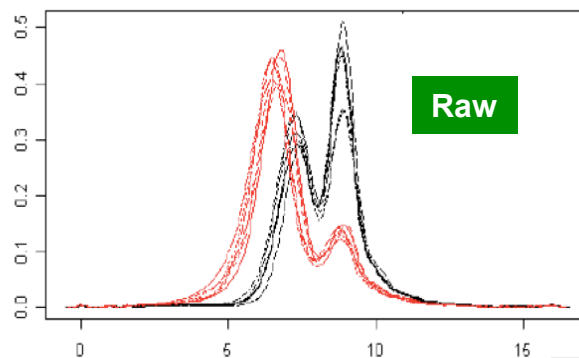
$\tau_k$    **Signal**
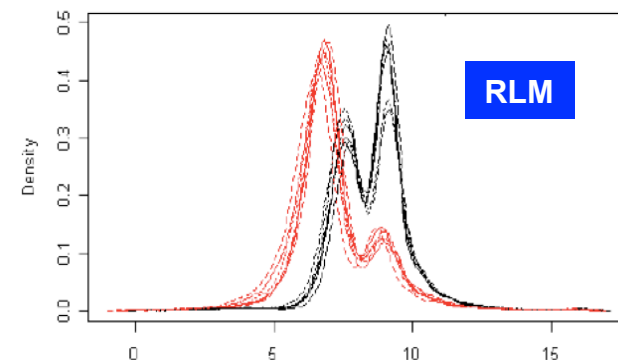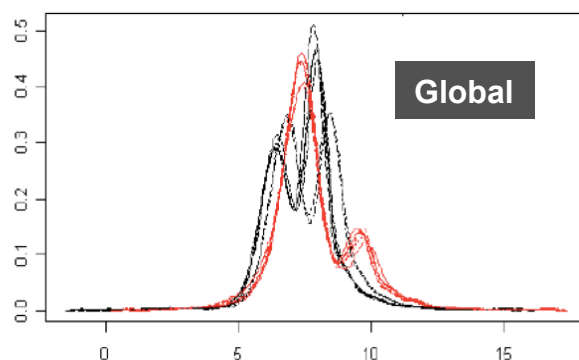
$\varepsilon_{ijkr}$    **Random error**



Raw

Quantile

[Sboner et al., J Proteome Sci. '09]

■ "normal" serum
■ "positive" serum

Global

RLM

Inter-array CV: "positive" serum

Inter-array CV: "normal" serum

**SAMPLE SEPARABILITY**

# Check #2: How Signal Intensity Correlates over a Titration

Protein signal vs. serum concentration

— strong correlation
- - - no correlation

density / correlation coefficient



"Positive" 100%    "Normal" 100%

"Positive" 75% "Normal" 25%   "Positive" 50% "Normal" 50%   "Positive" 25% "Normal" 75%

## Expectation

**"Positive" protein signal should positively correlate with "Positive" serum dilution**

**Higher number of "hits" for the "Positive" serum**

[Sboner et al., J Proteome Sci. '09]

## Correlation of signal intensity with "positive" serum concentration

orig
quantile
global
rlm.IgG.r
rlm.IgG+V5.r
rlm.V5.r

RLM

Global

Raw

Quantile

Density / Pearson's correlation

# Generating Networks #2

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

# Training sets



Known interactions

Known non-interactions

Unknown

# Network prediction: features

- Example 1: gene expression



Gasch et al., 2000

$$x_1 = (0.2, 2.4, 1.5, \ldots)$$
$$x_2 = (0.8, 2.2, 1.5, \ldots)$$
$$x_3 = (4.3, 0.1, 7.5, \ldots)$$
$$\ldots$$
$$\text{sim}(x_1, x_2) = 0.62$$
$$\text{sim}(x_1, x_3) = -0.58$$
$$\ldots$$

**Similarity scale:**

1 ▬▬▬▬▬▬▬▬▬▬ -1

# Network prediction: features

- Example 2: sub-cellular localization



http://www.scq.ubc.ca/wp-content/yeasttwohybridtranscript.gif

$$x_1 = (1, 1, 0, 0, \ldots)$$
$$x_2 = (1, 1, 1, 0, \ldots)$$
$$x_3 = (1, 0, 1, 0, \ldots)$$
$$\ldots$$
$$sim(x_1, x_2) = 0.81$$
$$sim(x_1, x_3) = 0.12$$
$$\ldots$$

**Similarity scale:**

1 ■■■■■■■■ -1

# Data integration & Similarity Matrix



|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 1.00 | 0.57 | 0.55 | 0.40 |
| 2 | 0.57 | 1.00 | 0.66 | 0.89 |
| 3 | 0.55 | 0.66 | 1.00 | 0.79 |
| 4 | 0.40 | 0.89 | 0.79 | 1.00 |

# Learning methods

## An endless list:

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

**Let's compare in a public challenge!**
**(DREAM: Dialogue for Reverse Engineering Assessment and Methods)**

# Our work: efficiently propagating known information

Training set expansion

- Motivation: lack of training examples
- Expand training sets horizontally

Multi-level learning

- Motivation: hierarchical nature of interaction
- Expand training sets vertically

DREAM3 *in silico* regulatory network reconstruction challenge

| Local model 1 | → | Local model 2 |

PPI predictions

↕

DDI predictions

↕

RRI predictions

# Protein interaction



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

**Protein-level features for interaction prediction: functional genomic information**

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

# Domain interaction



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

**Domain-level features for interaction prediction: evolutionary information**

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

# Residue interaction



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

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

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

# Combining the three problems

Protein interactions

Domain interactions

Residue interactions

i. Independent levels    ii. Unidirectional flow    iii. Bidirectional flow

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

# Empirical results (AUCs)

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

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

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

# Finding Central Points in Networks: Hubs & Bottlenecks

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

# Global topological measures

Indicate the gross topological structure of the network



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



log(Frequency)

$\log P(k)$

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

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"

# 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

**Bottlenecks & Hubs**

Bottleneck

Hub-bottleneck **node**

Non-hub-bottleneck **node**

Hub-non-bottleneck **node**

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

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

# Bottlenecks are what matters in regulatory networks



P < 10<sup>-20</sup>

Hub-non-bottlenecks
Bottleneck-non-hubs

P < 10<sup>-4</sup>

Fraction of essential genes

60%
50%
40%
30%
20%
10%
0%

Interaction Network          Regulatory Network

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

# Finding Central Points in Networks #2:
# Tops of the Hierarchy

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

# Social Hierarchy

THE GOVE[RNMENT OF THE] UNITED STATES

**LEGISLATIVE BRANCH**

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[ECUTIVE BRANCH]**

EX[ECUTIVE OFFICE OF THE] PRESIDENT

WHITE HOUSE OFFICE
OFFICE OF THE VICE PRESIDE[NT]
COUNCIL OF ECONOMIC ADVIS[ORS]
COUNCIL ON ENVIRONMENTAL [QUALITY]
NATIONAL SECURITY COUNCIL
OFFICE OF ADMINISTRATION

[MA]NAGEMENT AND BUDGET
[NA]TIONAL DRUG CONTROL POLICY
[PO]LICY DEVELOPMENT
[SC]IENCE AND TECHNOLOGY POLICY
[TH]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
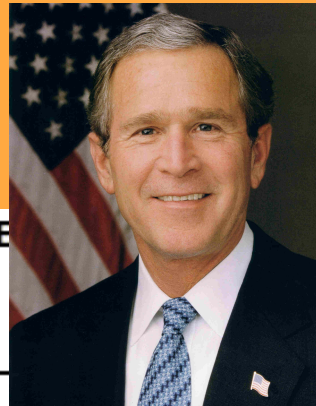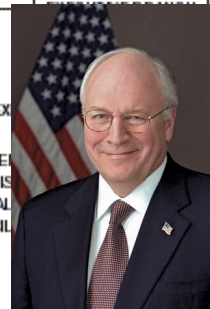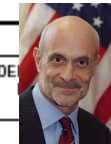FEDERAL RESERVE SYSTEM
FEDERAL RETIREMENT THRIFT INVESTMENT BOARD
FEDERAL TRADE COMMISSION
GENERAL SERVICES ADMINISTRATION
INTER-AMERICAN FOUNDATION
MERIT SYSTEMS PROTECTION BOARD
NATIONAL AERONAUTICS AND SPACE ADMINISTRATION
NATIONAL ARCHIVES AND RECORDS ADMINISTRATION
NATIONAL CAPITAL PLANNING COMMISSION
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



1

2

3

4

*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*

Level in hierarchy

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

P < 0.01

P < 6 X 10⁻⁴

P < 10⁻¹⁵

# of genes

[Yu et al., PNAS (2006)]

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

B. Governmental hierarchy of a represen-tive city (Macao)



Average # of regulated people (out-degree)
# of managers at each level

# 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)]

# Finding Central Points in Networks #3: Points of Maximal Regulatory Effect

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

## RE-score: Another way to identify "important" network nodes



Cheng et al., Genome Biology, 2009

RE score = $\bar{R}_n - \bar{R}_t$

Calculating RE scores of a miRNA in each sample

Comparing the RE scores between ER+ and ER-

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

*Cheng et al., Genome Biology, 2009*

# RE-score can be used to classify cancers

(3) Clustering based on RE score divides samples into 2 main types of cancer

(4) Clustering better than based on indiv. gene expression levels

ER+
ER-

hsa-miR-342
hsa-miR-193a
hsa-miR-145
hsa-miR-127
hsa-miR-122a
hsa-miR-588
hsa-miR-517a
hsa-miR-769-5p

(1) RE-score profile for diff. miRNA in 1 cancer sample.
(2) Tabulate over many different breast cancer samples

*Cheng et al., Genome Biology, 2009*

# Network Dynamics #2: Environments

**How do molecular networks change across environments?**
**What pathways are used more ?**
**Used as a biosensor ?**

# What is metagenomics?

## Genomics Approach

**Culture Microbes** → **Extract DNA** → **Sequence**

ATCGTATA
CGCGAAG
ACGTCTGA
AGTGCTGCT

→ **Assemble and Annotate**

Contig 1

**PROBLEM:** Estimated that less than 1% can be cultured in the lab

## Metagenomics Approach

**Collect Sample** → **Extract DNA** → **Sequence**

ATCGTGATAGATGATAGTAGA
ATGCTGCATGCATCTAGCACT
ACAGTAGCTAGCTACGTACTA
CAGCTGACTAGCTAGCTAGCT
ACGTAGCATGCTAGCTAGCAG
ACGTACGTAGCTAGCTAGCTAG
ACGTACGTACGTAGCTAGCATC
AGTCGACTGAGCCAGTGATGAT
ACGATGCATGAGCAGATGCTAC
AGATCGTAGCATGCTAGCATGCT
ACGTACGTAGCTAGCTAGCTAAG
AGCTAGCATGCTAGTAGCATGAG
ACGATGCTAGCTAGCTAGCTGATA
TCGATCAGCATGCTACGATGCAAG
ACGATCGATGCTAGCTAGCTAGCAT
AGCTAGCTAGTCAGCTAGCTAGATG

→ **Partially Assemble and Annotate**

**PROBLEM:** Lose information about which gene belongs to which microbe.

# 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----------
CCGTGACGCGATGC------
CCGTGAGCACGATGCGC ATGCTCATGCT--------
    ATCGTGACGCGGATGC------
        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**

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

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

## 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) ]

Environmental Features

Chlorophyll    Temp

P
a
t
h
w
a
y
s

Cobalamin Biosynthesis

Photosystem II

Photosystem I

Carbon Fixation (Dark rx)

Glutamine Degradation

$r^2 = .68$

Predicted Temperature

Actual Temperature

# Canonical Correlation Analysis: Simultaneous weighting



UPI = a GRE + b (books) + c GPA

GPI = a' (journals) + b' PowerPoint + c' (money)

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

# Canonical Correlation Analysis: Simultaneous weighting



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

$$C = \frac{\Sigma_X \quad \Sigma_{X,Y}}{\Sigma_Y \quad \Sigma_{Y,X}}$$

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

[ 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

[ 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) ]

# Biosensors:
# Beyond Canaries in a Coal Mine



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

# 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

**Hapmap/Perlegen**

SNPs

**Database of Genomic Variants**

CNVs + SDs

Map to ENSEMBL genes

**Ensembl Genes**

ENSG000XXXX:
rsSNP00XXX
CNV_XXX
DN/DS XXXX
Recombination rate

Map to proteins in the interaction network

**Interactome**

~30000 interactions from HPRD and Y2H screens

**Result**

- Dataset of network position / parameters (e.g. degree centrality or betweenness centrality) in relationship to SNPs, CNV's, recombination rates and positive selection tests

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

Source: PMK

# ADAPTIVE EVOLUTION CAN BE SEEN ON TWO DIFFERENT LEVELS

**Intra-species variation**

**Fixed mutations (differences to other species)**

**Single-basepair**

Positive Selection

**Single-Nucleotide Polymorphisms**

**Fixed Differences**

**Structural variation**

Positive Selection

**Copy Number Variants**

**Segmental Duplications**

Source: PMK

# POSITIVE SELECTION LARGELY TAKES PLACE AT THE NETWORK PERIPHERY

**Positive selection in the human interactome**



- High likelihood of positive selection
- Lower likelihood of positive selection
- Not under positive selection
- No data about positive selection

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

# CENTRAL PROTEINS ARE LESS LIKELY TO BE UNDER POSITIVE SELECTION

**Degree vs. Positive Selection**



Spearman Rank P−value: 1.2e−06

x-axis: Betweenness Centrality (× 10^6)
y-axis: Positive Selection Test Likelihood Ratio

Network periphery → Network center

**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)

# CENTRAL NODES ARE LESS LIKELY TO LIE INSIDE OF SDs

## Centrality vs. SD occurrence



## Reasoning

- This result also confirms our initial hypothesis – peripheral nodes tend to lie in regions rich in SDs.

- Since segmental duplications are a different mechanism of ongoing evolution, the less constrained peripheral proteins are enriched in them.

- Note that despite the small size of our dataset for known SD's we get significant correlations. It is to be expected that the correlations will get clearer as more data emerges*

* Specifically, a number of the SDs are likely not fixed, but rather common CNVs in the reference genome

Source: Database of genetic variation, HPRD, Rual et al. *Nature* (2005), and Kim et al. PNAS (2007)

# IS RELAXED CONSTRAINT OR ADAPTIVE EVOLUTION THE REASON FOR THE PREVALENCE OF BOTH SELECTED GENES AND SDs AT THE NETWORK PERIPHERY?

ILLUSTRATIVE

|  | **Relaxed Constraint** | **Adaptive Evolution** |
|---|---|---|
| **Inter-Species Variation (Fixed differences)** | • Increases inter-species variation – more variable loci are under less negative selection<br><br>• Can be seen in higher Ka/Ks ratio or SD occurrence | • Increases inter-species variation – more variable loci are under less negative selection<br><br>• Can be seen in higher Ka/Ks ratio or SD occurrence |
| **Intra-Species Variation (Polymorphisms)** | • Increases intra-species variation – for the very same reason<br><br>• Can be seen in both SNPs or CNVs | • Should not have effects on intra-species variation |

Source: Kim et al. PNAS (2007)

# SOME, BUT NOT ALL OF THE SINGLE-BASEPAIR SELECTION AT THE PERIPHERY IS DUE TO RELAXED CONSTRAINT

## Inter vs. Intra-Species Variation in Networks

**Inter-Species (Fixed differences)**

4.37

2.71

Betweenness Centrality (x 10⁴)

p<<0.01

Genes with dN/dS>1     Genes with dN/dS<=1

**Intra-Species (SNPs) [ Variability ]**

4.08     4.35

Betweenness Centrality (x 10⁴)

p<0.05

Genes with pN/pS>1     Genes with pN/pS<=1

## Reasoning

- There is a difference in **variability** (in terms of SNPs) between the network periphery and the center

- However, this difference is much smaller than the difference in **selection**

- This most likely means, that part of the effect we're seeing is due to relaxed constraint (and higher variability)

- But, not the entire effect*

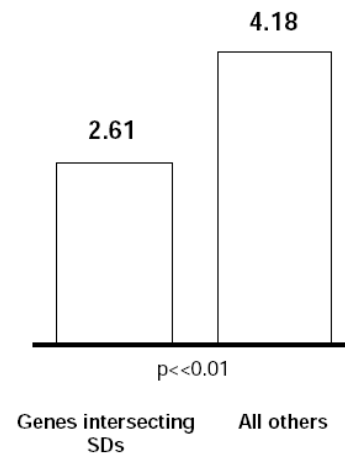  \* But it's hard to quantify

Source: Kim et al. (2007) PNAS

# Similar Results for Large-scale Genomic Changes (CNVs and SDs)
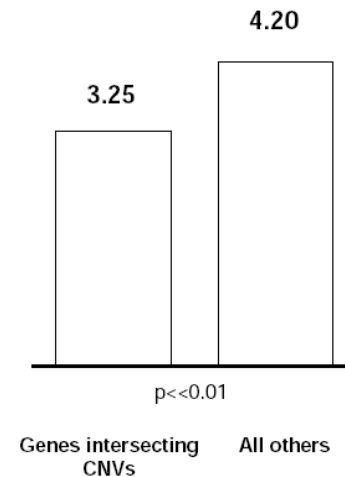
## Inter vs. Intra-Species Variation in Networks

**Inter-Species (SDs)**

Betweenness Centrality (x $10^4$)

2.61

4.18

p<<0.01

Genes intersecting SDs     All others

**Intra-Species (CNVs) [ Variability ]**

Betweenness Centrality (x $10^4$)

3.25

4.20

p<<0.01

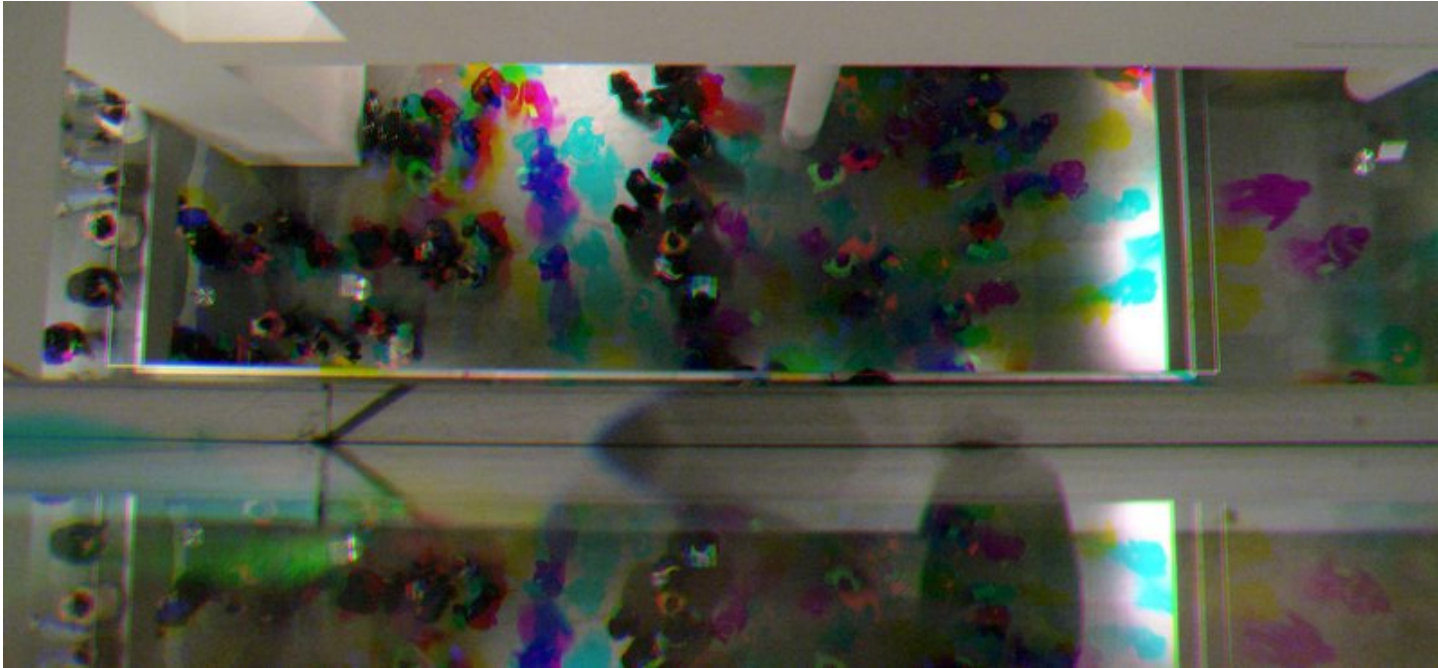Genes intersecting CNVs     All others

## Reasoning

- There a small difference in **variability** (in terms of CNVs) between the network periphery and the center

- But, there is a (as shown before) marked difference in fixed (and hence, presumably, **selected**) SDs at the network periphery and center

Source: Kim et al. (2007) PNAS

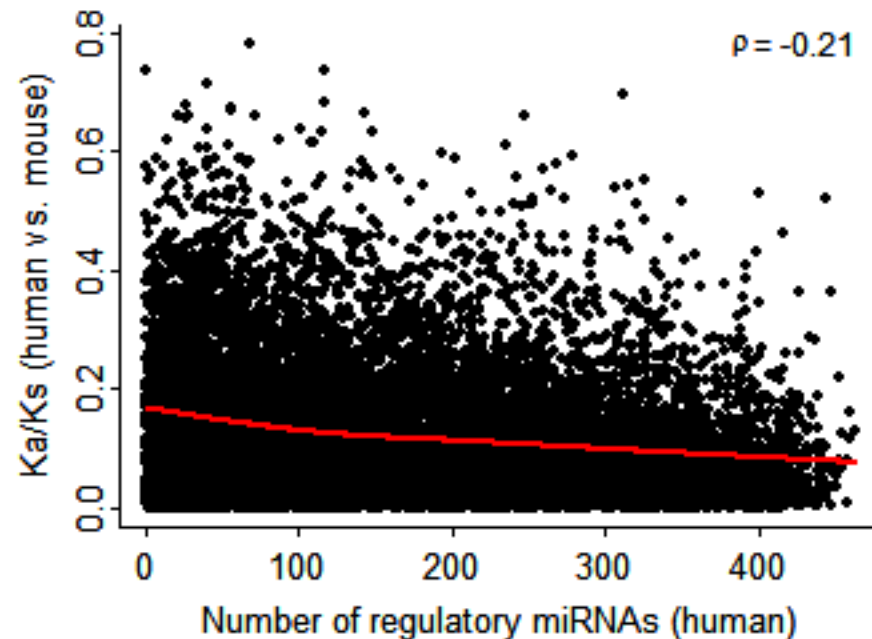# Networks & Variation 2
**Variation in the miRNA network**

# Analyze Regulation in microRNA-target Network

- Relationship between target in degree
  (number of micro-RNAs that regulate gene)
  & evolutionary rate of gene?

  ◊ In deg. related 3' UTR size


- Expectation: more regulation, more constraint

# Relationship between microRNA regulation and protein evolution



$\rho = -0.21$

Ka/Ks (human vs. mouse) vs. Number of regulatory miRNAs (human)

**Important genes are regulated more intensively regulated by the microRNAs**

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

[Cheng et al., BMC Genomics, 2009 (in press)]
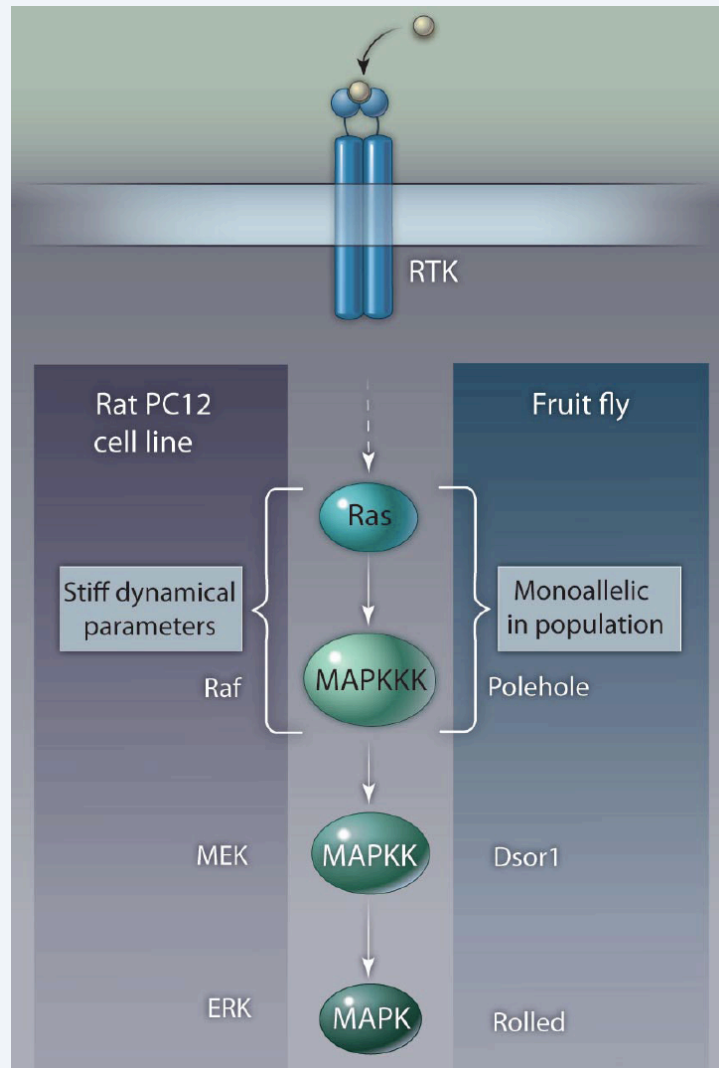
# MicroRNA regulation:
# a two-way strategy

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

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

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

# Network dynamics constrain evolution



**Hypothesis: Nodes in a molecular network with the strongest impact on dynamic behavior should be under strong purifying selection and thus exhibit the least genetic variation.**

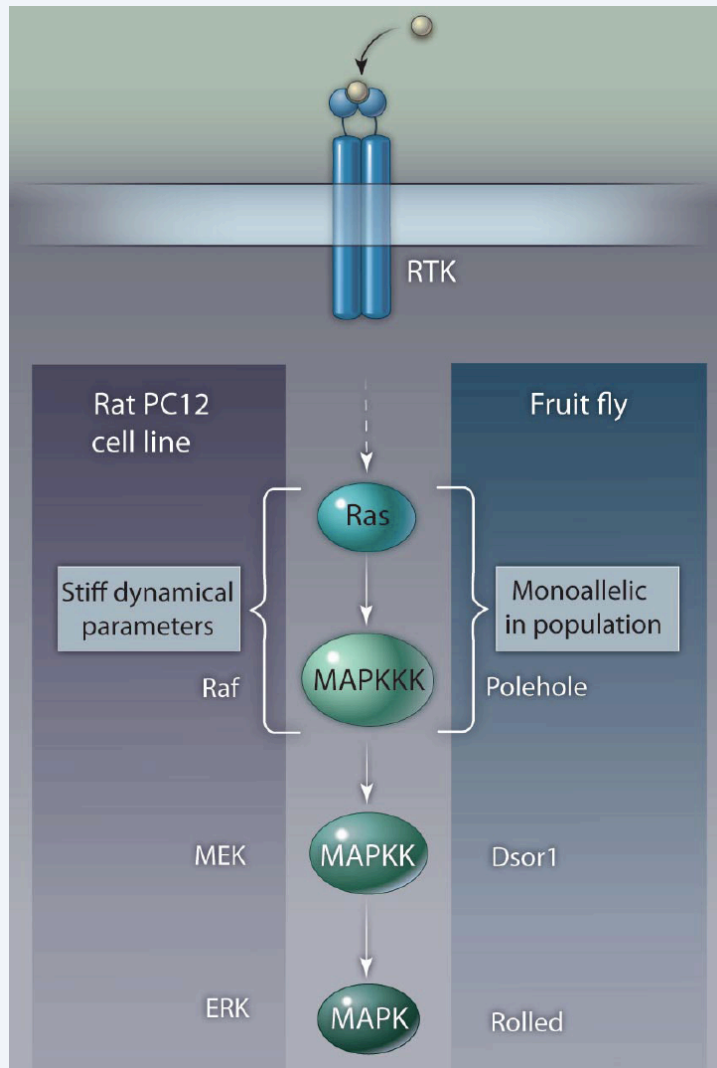**Alexander et al.** *Sci. Signal.* **(2009) 2: pe44**
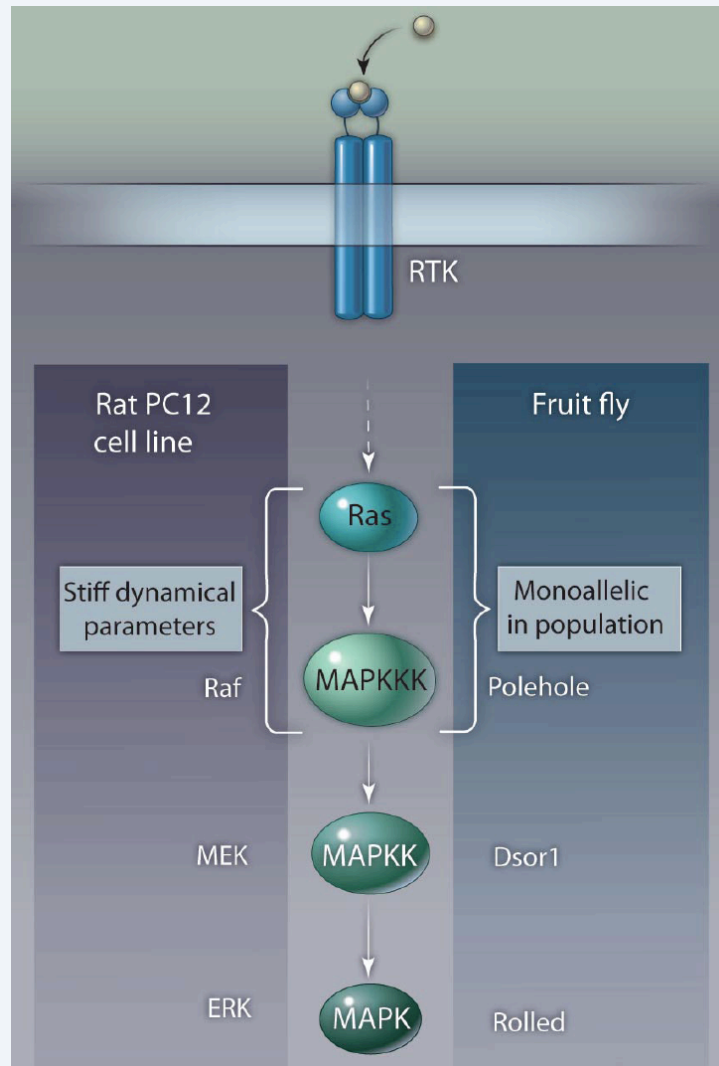
# Network dynamics constrain evolution



**Hypothesis: Nodes in a molecular network with the strongest impact on dynamic behavior should be under strong purifying selection and thus exhibit the least genetic variation.**

**Algorithm:**
**1) Reconstruct families of molecular networks from genomic data.**
**2) Map some kind of genetic variation onto the networks.**
**3) Analyze sensitivity of dynamical model of the generic network.**

**Alexander et al.** *Sci. Signal.* **(2009) 2: pe44**

# Speculation: Why more tightly regulated gene might have less variation



Example: MAP Kinase singaling pathway

**Dynamic model:**
- ODE model with Michaelis-Menten kinetics
- parameters fit
   to time series data of protein activities
   in response to EGF and NGF
   from rat PC12 cell line

In sensitivity analysis,
   stiff parameters cluster around Ras and Raf.

**Population study in fruit flies:**
- allele variation based on
   PCR of pathway genes

Ras and Raf have less allele variation
   than other proteins in the network.

**Alexander et al. *Sci. Signal.* (2009) 2: pe44**

**Brown et al. *Phys. Biol.* (2004) 1: 184**
**Riley et al. *Molec. Ecol.* (2003) 12: 1315**

- Why Networks?
- Generating Networks
  - ◊ Processing Protein Chips
    (yeast & human nets)
  - ◊ Propagating Known Information
    (yeast ppi)
- Central Points in Networks
  - ◊ Hubs & Bottlenecks
    (yeast ppi & reg. net)
  - ◊ Tops of Heirarchies
    (yeast reg. net)
  - ◊ Identified by score
    (human miRNA-targ. net)
- Dynamics of Networks
  - ◊ Across environments
    (prokaryote metab. pathways)
- Protein Networks & Variation
  (human ppi & miRNA-targ. net)

# Conclusions on Networks: Generation



- Networks from processing protein chip data
  - ◊ RLM normalization surpresses quantile

- Predicting Networks
  - ◊ Extrapolating from the Training Set
  - ◊ Principled ways of using known information in the fullest possible fashion
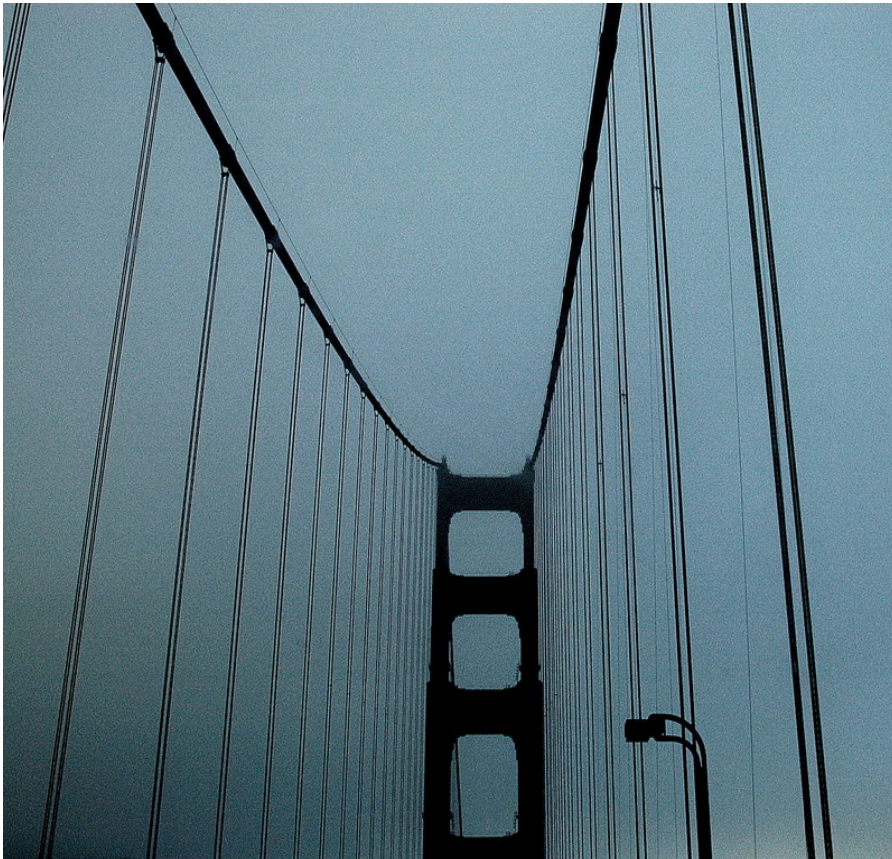    - Multi-level learning

# Conclusions:
## Analysis of Network Structure



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

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

# Conclusions:
# Points of Network Centrality



- RE-score measures degree of (down) regulation of targets v. non-targets

- Application to miRNA network

- Use in cancer classification

# Conclusions: Networks Dynamics across Environments



- Developed and adapted techniques to connect quantitative features of environment to metabolism.

- Applied to available aquatic datasets, we identified footprints that were predictive of their environment (potentially could be used as biosensor).

- Strong correlation exists between a community's energy conversion strategies and its environmental parameters (e.g. temperature and chlorophyll).

- Suggest that limiting amounts of cofactor can (partially) explain increased import of amino acids in nutrient-limited conditions.

# Conclusions: Connecting Networks & Variation



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

# TopNet
## Topology of Networks

# tYNA

# - an automated web tool
(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

**H Yu**
**P Kim**
**K Yip**
**T Gianoulis**
**C Cheng**
**A Sboner**

G Chen
M Smith
D Mattoon
L Freeman-Cook
P Patel
A Karpikov
A Paccanaro
P Alves
N Bhardwaj
R Alexander
P Cayting          J Raes
M Seringhaus       **T Emonet**
Y Xia              **P Bork**
J Korbel           **B Schweitzer**
E Franzosa         **M Snyder**



## Networks.GersteinLab.org

Job opportunities currently for postdocs & students

# More Information on this Talk

**SUBJECT:** Networks

**DESCRIPTION**:
CSHL, Cold Spring Harbor, NY; 2010.01.06, 12:00-13:00; [I:CSHL2]
(Long networks talk, derived from [I:MBINETS], including rlm* & new
intro. for 1st time)

(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 )