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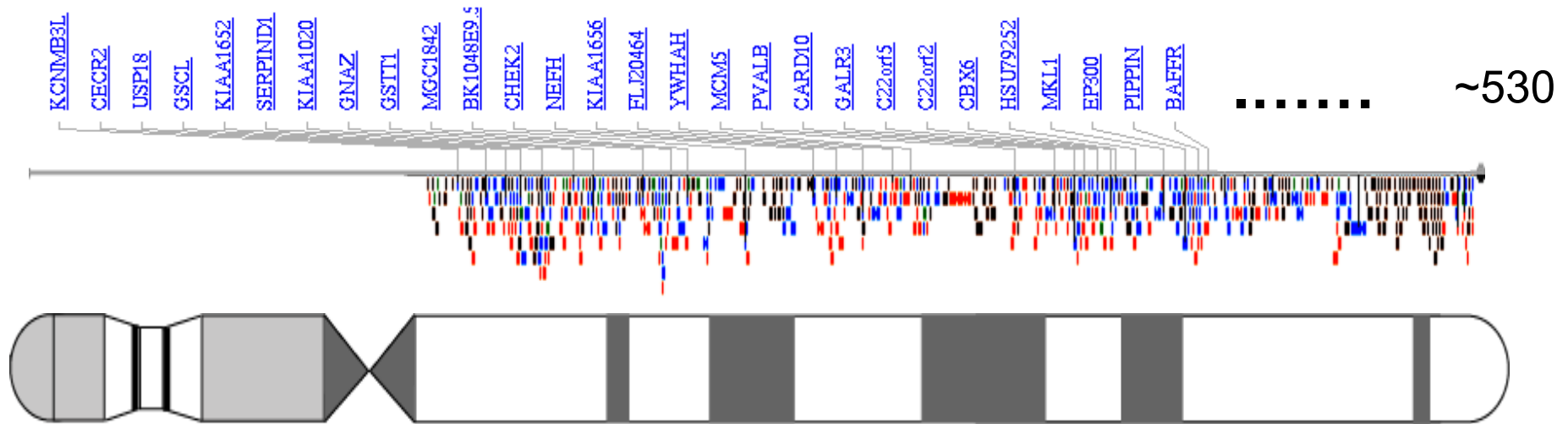
Mark B Gerstein
Yale

slides at
**Lectures.GersteinLab.org**

# 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)

2 - Lectures.GersteinLab.org (c) '09

# 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
  - ◊ Often >2 proteins/function
  - ◊ Multi-functionality:
    2 functions/protein
  - ◊ Role Conflation:
    molecular, cellular, phenotypic
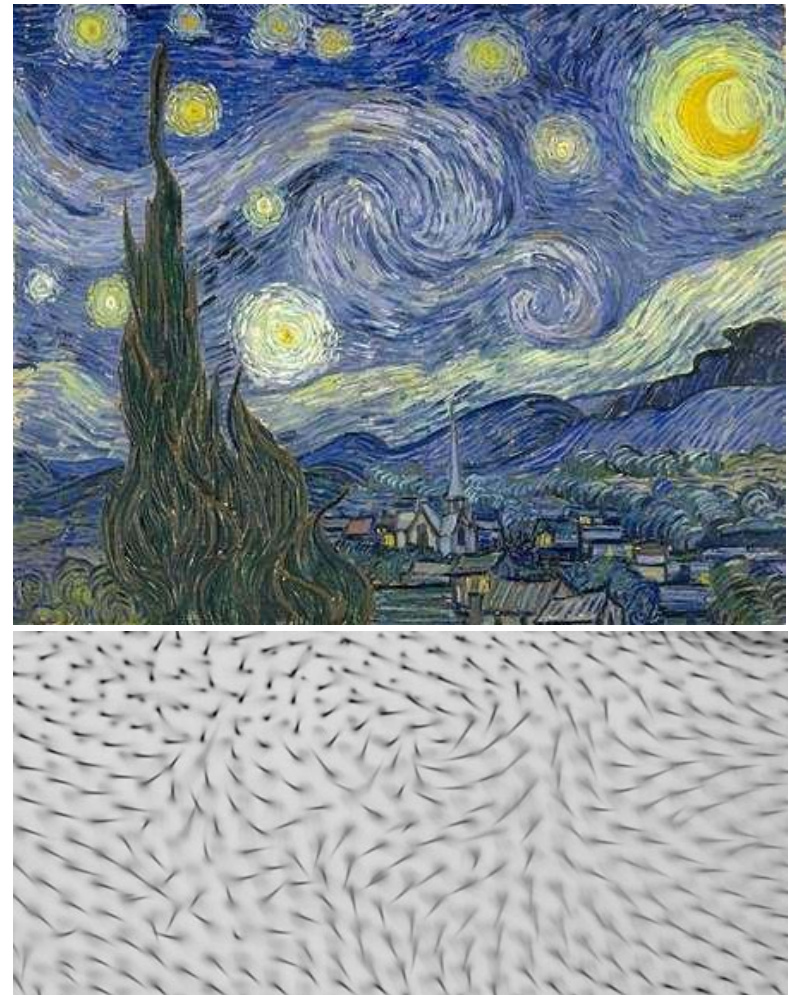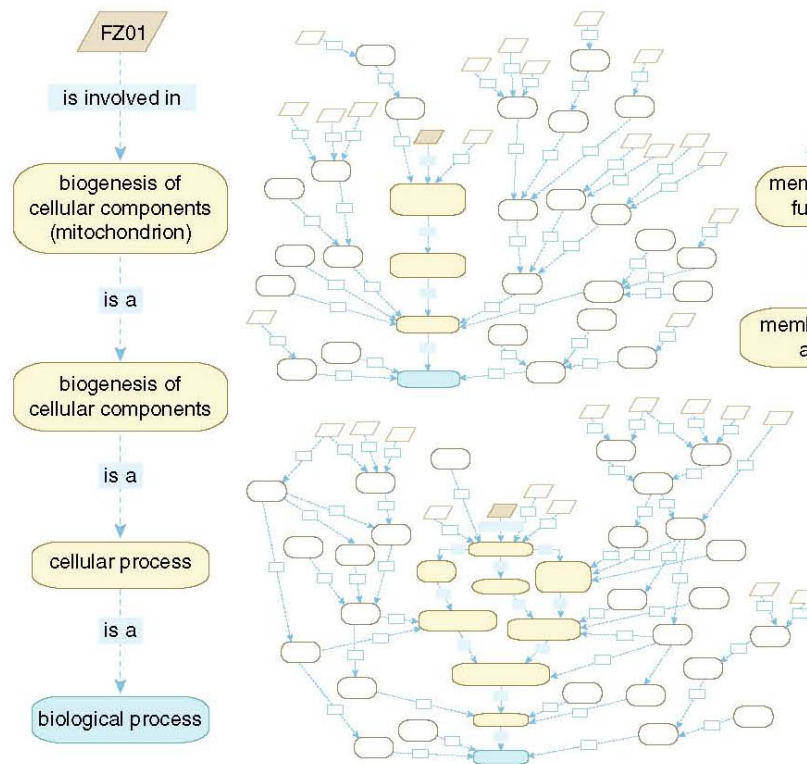
# Some obvious issues in scaling single molecule definition to a genomic scale

- 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)
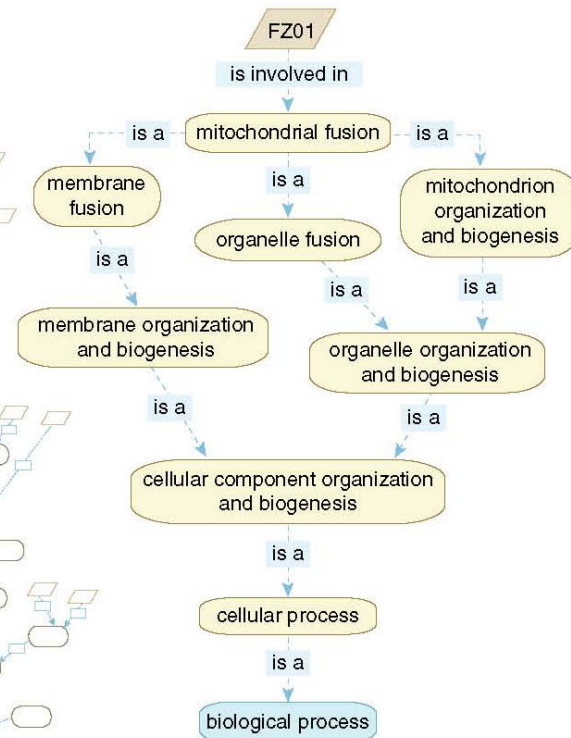


[Seringhaus et al. GenomeBiology (2008)]

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



**MIPS (Mewes et al.)**

**GO (Ashburner et al.)**

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

# Towards Developing Standardized Descriptions of Function

- Subjecting each gene to standardized expt. and cataloging effect
  - ◊ KOs of each gene in a variety of std. conditions => phenotypes
  - ◊ Std. binding expts for each gene (e.g. prot. chip)
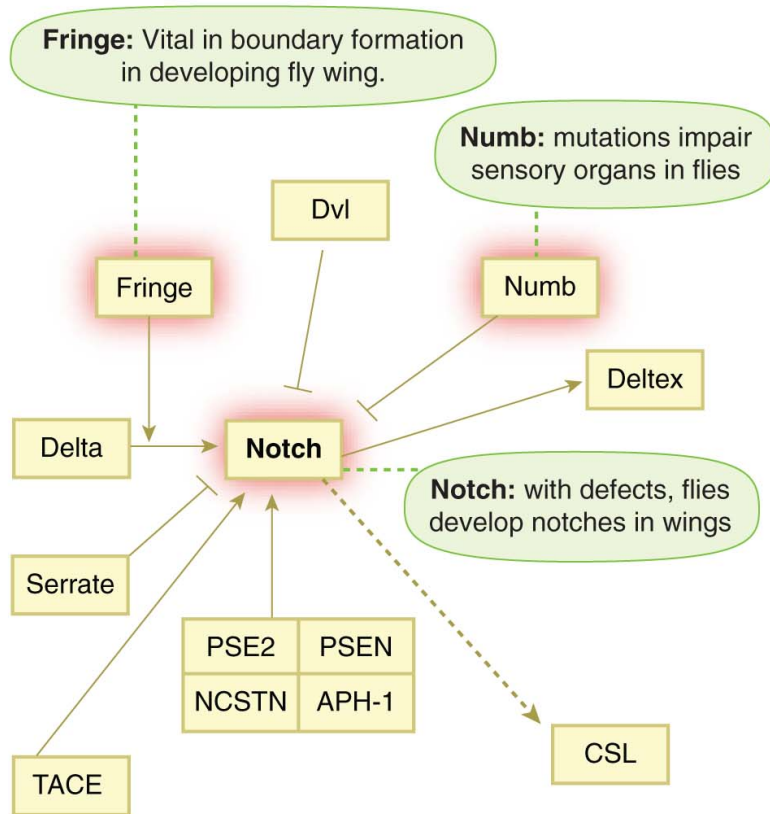
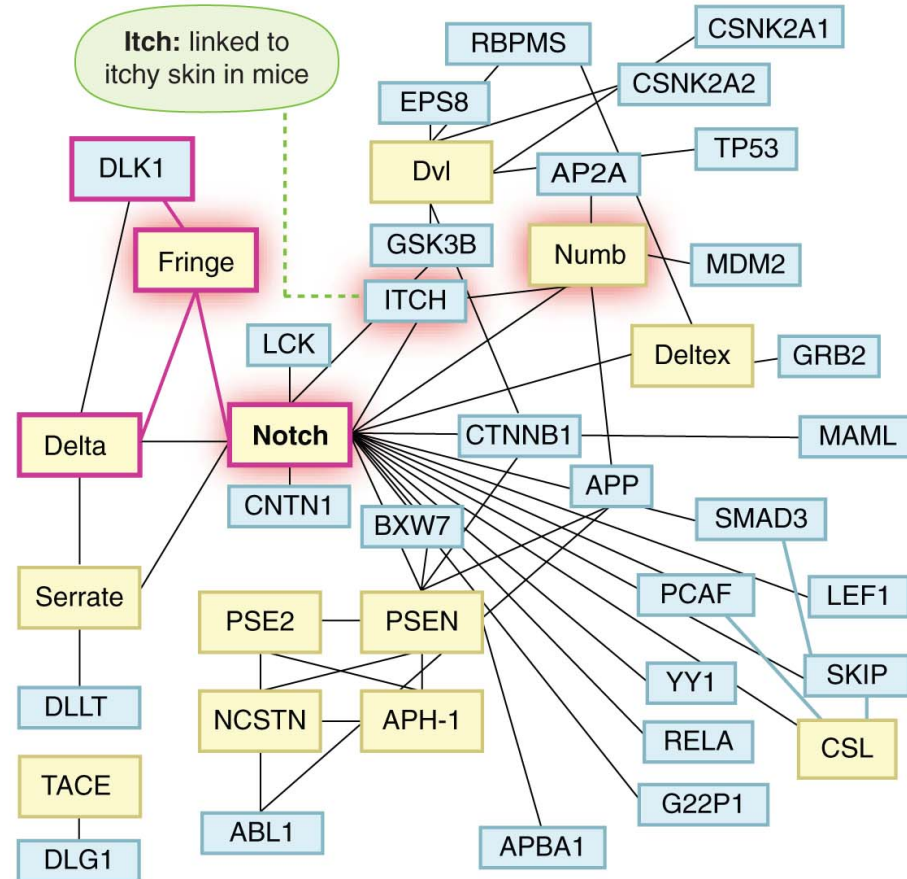- Function as a vector

| | nucleic acids | | small molecules | | | | | proteins | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | DNA | RNA | ATP | Metal | CoA | NAD | ...... | G protein | CDC28 | Calmodulin | ...... |
| protein 1 | 1.0 | 0 | 0 | 0 | 0 | 0 | ...... | 0 | 0 | 0 | ...... |
| protein 2 | 0 | 0.9 | 0 | 0 | 0 | 0 | ...... | 0 | 0 | 0 | ...... |
| protein 3 | 1.0 | 0 | 1.0 | 0 | 0 | 0 | ...... | 0 | 0 | 0 | ...... |
| protein 4 | 0 | 0 | 0 | 0 | 0.8 | 0 | ...... | 0 | 0 | 1.0 | ...... |
| protein 5 | 1.0 | 0 | 0 | 0 | 0 | 0 | ...... | 0 | 0.9 | 0 | ...... |
| protein 6 | 0.9 | 0 | | | | | ...... | | | | ...... |
| protein 7 | 0 | 0.8 | | | | | ...... | | | | ...... |
| ...... | ...... | ...... | ...... | ...... | ...... | ...... | ...... | ...... | ...... | ...... | ...... |

**Interaction Vectors** [Lan et al, IEEE 90:1848]

# 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: Detailed
structural
understanding of
cellular machinery

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

[Jeong et al. Nature, 41:411]

# **Networks as a universal language**



Internet
[Burch & Cheswick]

Food Web

Electronic
Circuit

Neural Network
[Cajal]

Disease
Spread
[Krebs]

Albert-László
Barabási

LINKED

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

# Outline: Molecular Networks

- Why Networks?

- Predicting Networks (yeast ppi)

   ◊ Propagating known information

- 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
      (in prokaryote metab. pathways)

- Protein Networks & Variation
   (human ppi & miRNA-targ. net)

# Example: yeast PPI network



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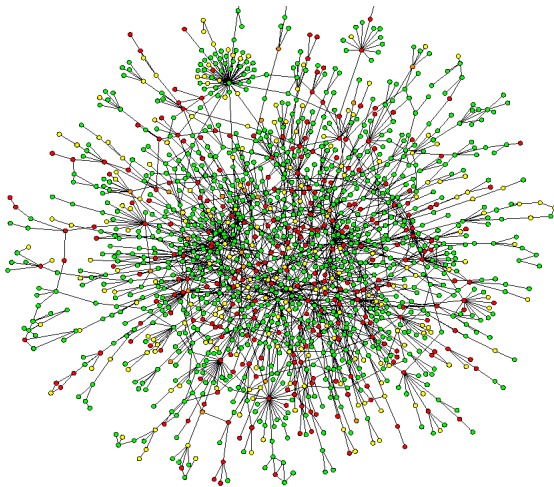
◊ ~6,000 nodes
  → Computational cost: ~18M pairs

◊ 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]

# Predicting Networks

**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

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

**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)**

# DREAM3: *in silico* regulatory network reconstruction

**Actual network**          **Expression data**          **Modeling**          **Predictions**



**Deletion strains**

Prob(signal|point)
= 2Φ((point − ref) / s) − 1

**Noise models**

**Time series after initial perturbation**

$$\frac{dy_j}{dt} = a_{j0} - a_{jj}y_j + \sum_{k \in S} a_{jk}y_k$$

$$\frac{dy_j}{dt} = \frac{b_{j1}}{1 + \exp\left(a_{j0} + \sum_{k \in S} a_{jk}y_k\right)} - b_{j2}y_j$$

$$\frac{dy_j}{dt} = a_{j0} \prod_{k_1 \in S_1}\left(\frac{b_{jk_1}}{y_{k_1}^{c_{jk_1}} + b_{jk_1}}\right)\prod_{k_2 \in S_2}\left(\frac{y_{k_2}^{c_{jk_2}}}{y_{k_2}^{c_{jk_2}} + b_{jk_2}}\right) - a_{j1}y_j$$

**Expression rate models**

| Accuracy (AUC) | E. Coli 1 | E. Coli 2 | Yeast 1 | Yeast 2 | Yeast 3 |
|---|---|---|---|---|---|
| Size-10 | 0.928 | 0.912 | 0.949 | 0.747 | 0.714 |
| Size-50 | 0.930 | 0.924 | 0.917 | 0.792 | 0.805 |
| Size-100 | 0.948 | 0.960 | 0.915 | 0.856 | 0.783 |

# 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

# Kernels

Kernel: a similarity matrix that is positive semi-definite (p.s.d.)



Compute
inner products

→

←

p.s.d. implies

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 1.00 | 0.72 | 0.45 | -0.56 |
| 2 | 0.72 | 1.00 | -0.30 | -0.98 |
| 3 | 0.45 | -0.30 | 1.00 | 0.49 |
| 4 | -0.56 | -0.98 | 0.49 | 1.00 |

**Objects in an feature space**

**Similarity matrix**

Good for integrating heterogeneous datasets (protein sequences, PSSM, gene expression, …)

– no need to explicitly place them in a common feature space

# Kernel methods

Use the kernel as proxy to work in the feature space

Example: SVM (finding the best separating hyperplane)



Equivalent to

Maximize $\sum_i \lambda_i - \dfrac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle$

Subject to $\lambda \geq 0$

$\sum_i \lambda_i y_i = 0$

The only thing that we need to know about the objects: their similarity values (inner products)

# Kernel methods for predicting networks: local vs. global modeling



Global modeling: build one model for the whole network

Pairwise kernel: consider object pairs instead of individual objects

Problem: $O(n^2)$ instances, $O(n^4)$ kernel elements



Direct methods: threshold the kernel to make predictions

Problem: One single global model, may not be able to handle subclasses

|   | 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 |

Threshold: 0.7 →

|   | 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 |

# Kernel methods for predicting networks: local vs. global modeling



Local modeling: build one model for each node

**Model for node 3:**



Problem: insufficient and unevenly distributed training data (what if node 3 has no known interactions at all?)

# Our work: training set expansion

- Goal:
    ◊ Utilize the flexibility of local modeling
    ◊ Tackle the problem of insufficient training data
- Idea: generate auxiliary training data
    ◊ Prediction propagation
    ◊ Kernel initialization

**[Yip and Gerstein, Bioinformatics ('09, in press)]**

# Prediction propagation

- Motivation: some objects have more examples than others

- Our approach:
  - ◊ Learn models for objects with more examples first
  - ◊ Propagate the most confident predictions as auxiliary examples of other objects

# Kernel initialization

- Motivation: what if most objects have very few examples?

- Our approach (inspired by the direct method):

  ◊ Add the most similar pairs in the kernel as positive examples

  ◊ Add the most dissimilar pairs in the kernel as negative examples

|   | 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 |

[Yip and Gerstein, Bioinformatics ('09, in press)]

# Remarks

- Can be used in combination
- Prediction propagation theoretically related to co-training (Blum and Mitchell, 1998)
  - ◊ Semi-supervised
    - Similarity with PSI-BLAST
- Algorithm complexity $O(nf(n))$ of local modeling vs. $O(f(n^2))$ of global modeling

**[Yip and Gerstein, Bioinformatics ('09, in press)]**

# Prediction accuracy (AUC)

| | phy | loc | exp-gasch | exp-spellman | y2h-ito | y2h-uetz | tap-gavin | tap-krogan | int |
|---|---|---|---|---|---|---|---|---|---|
| Mode 1 | | | | | | | | | |
| direct | 58.04 | 66.55 | 64.61 | 57.41 | 51.52 | 52.13 | 59.37 | 61.62 | 70.91 |
| kCCA | 65.80 | 63.86 | 68.98 | 65.10 | 50.89 | 50.48 | 57.56 | 51.85 | 80.98 |
| kML | 63.87 | 68.10 | 69.67 | 68.99 | 52.76 | 53.85 | 60.86 | 57.69 | 73.47 |
| em | 71.22 | 75.14 | 67.53 | 64.96 | 55.90 | 53.13 | 63.74 | 68.20 | 81.65 |
| local | 71.67 | 71.41 | 72.66 | 70.63 | 67.27 | 67.27 | 64.60 | 67.48 | 75.65 |
| local+pp | 73.89 | 75.25 | 77.43 | 75.35 | 71.60 | 71.51 | 74.62 | 71.39 | 83.63 |
| local+ki | 71.68 | 71.42 | 75.89 | 70.96 | 69.40 | 69.05 | 70.53 | 72.03 | 81.74 |
| local+pp+ki | 72.40 | 75.19 | 77.41 | 73.81 | 70.44 | 70.57 | 73.59 | 72.64 | 83.59 |

## Observations:

- Highest accuracy by training set expansion
- Over fitting of local modeling without training set expansion
- Prediction propagation theoretically related to co-training (Blum and Mitchell, 1998)
  - ◊ Semi-supervised (Similarity with PSI-BLAST)

[Yip and Gerstein, Bioinformatics ('09)]

# Complementarity of the two methods

[**Yip and Gerstein, Bioinformatics ('09, in press)**]

# From horizontal to vertical

## Training set expansion

Local model 1 → Local model 2

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

## Multi-level learning

PPI predictions

↕

DDI predictions

↕

RRI predictions

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

# 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(Degree)

$\log P(k)$

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

$\log k$

***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"

[Yu et al., 2003, TIG]

# 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**

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



[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... UNITED STATES

| LEGISLATIVE BRANCH | EXECUTIVE BRANCH | JUDICIAL BRANCH |
|---|---|---|
| THE CONGRESS | | THE SUPREME COURT OF THE UNITED STATES |
| SENATE   HOUSE | VICE PRESIDENT | UNITED STATES COURTS OF APPEALS |
| | | UNITED STATES DISTRICT COURTS |
| ARCHITECT OF THE CAPITOL | WHITE HOUSE OFFICE ... MANAGEMENT AND BUDGET | TERRITORIAL COURTS |
| UNITED STATES BOTANIC GARDEN | OFFICE OF THE VICE PRESIDE... ...TIONAL DRUG CONTROL POLICY | UNITED STATES COURT OF INTERNATIONAL TRADE |
| GENERAL ACCOUNTING OFFICE | COUNCIL OF ECONOMIC ADVIS... ...OLICY DEVELOPMENT | UNITED STATES COURT OF FEDERAL CLAIMS |
| GOVERNMENT PRINTING OFFICE | COUNCIL ON ENVIRONMENTAL... ...IENCE AND TECHNOLOGY POLICY | UNITED STATES COURT OF APPEALS FOR THE ARMED FORCES |
| LIBRARY OF CONGRESS | NATIONAL SECURITY COUNCIL... ...E U.S. TRADE REPRESENTATIVE | UNITED STATES TAX COURT |
| CONGRESSIONAL BUDGET OFFICE | OFFICE OF ADMINISTRATION | 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
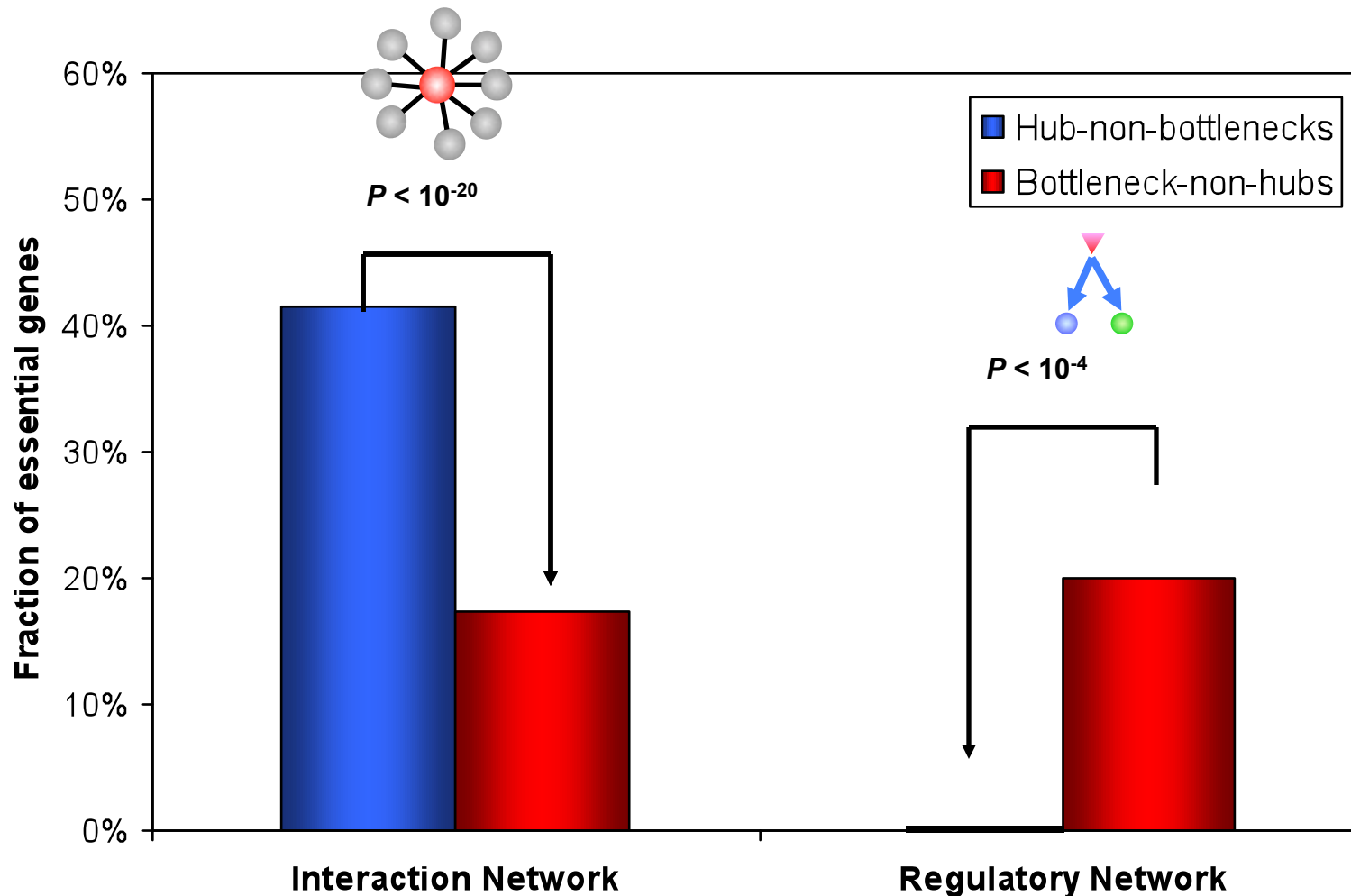NATIONAL LABOR RELATIONS BOARD
NATIONAL MEDIATION BOARD
NATIONAL RAILROAD PASSENGER CORPORATION (AMTRAK)
NATIONAL SCIENCE FOUNDATION
NATIONAL TRANSPORTATION SAFETY BOARD
NUCLEAR REGULATORY COMMISSION
OCCUPATIONAL SAFETY AND HEALTH REVIEW COMMISSION
OFFICE OF GOVERNMENT ETHICS
OFFICE OF PERSONNEL MANAGEMENT
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)

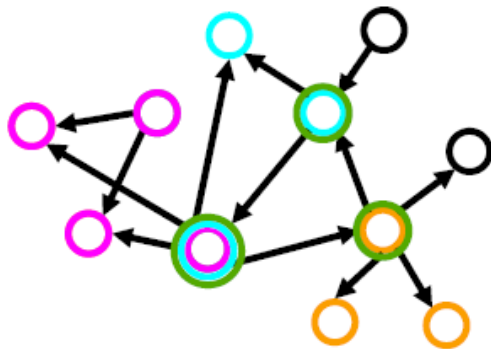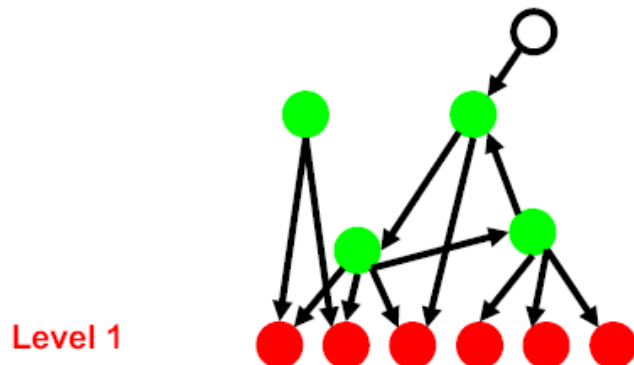III. Finding mid-level nodes (Green)

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*

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

P < 0.01

P < 6 X 10⁻⁴

P < 10⁻¹⁵

[Yu et al., PNAS (2006)]

# 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

[Yu et al., PNAS (2006)]

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



**A.** Level in Hierarchy vs. # of affected genes in knock-out experiments

Level 4: Data unavailable

P < 10^{-44}

**C.** Level in Hierarchy vs. Fraction of essential genes in *S. cerevisiae*

Level 4: P < 10^{-10}
Level 3: P < 0.02
Level 2: P = 0.42

[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 micro-RNA-target networks easy to calculate, as all influence is down-regulation
  - ◊ target prediction methods: 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 measure "importance" in networks



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

Cheng et al., Genome Biology, 2009

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

*Cheng et al., Genome Biology, 2009*

# miRNA RE-scores can be used to classify cancers



ER+
ER-

RE-score profiles for 8 miRNAs

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

*Cheng et al., Genome Biology, 2009*

# Differential expression of miRNA processing genes

**Distribution of ER-/ER+ T-scores for all miRNAs**



**The majority of miRNAs have higher RE-score in ER- than in ER+**

# 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** → **Assemble and Annotate**

Sequence:
ATCGTATA
CGCGAAG
ACGTCTGA
AGTGCTGCT

Contig 1

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

## Metagenomics Approach

**Collect Sample** → **Extract DNA** → **Sequence** → **Partially Assemble and Annotate**

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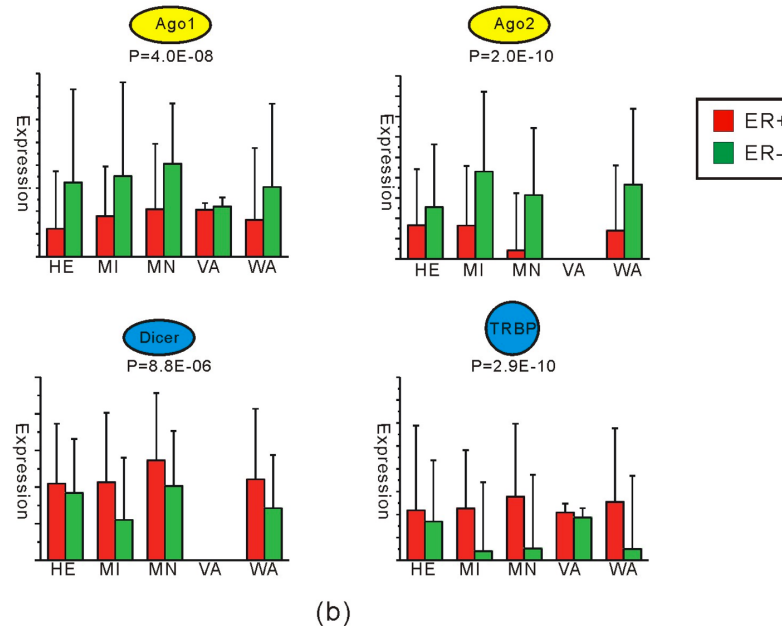ATGCTGCATGCATCTAGCACT
ACAGTAGCTAGCTACGTACTA
CAGCTGACTAGCTAGCTAGCT
ACGTAGCATGCTAGCTAGCAG
ACGTACGTAGCTAGCTAGCTAG
ACGTACGTACGTAGCTAGCATC
AGTCGACTGAGCCAGTGATGAT
ACGATGCATGAGCAGATGCTAC
AGATCGTAGCATGCTAGCATGCT
ACGTACGTAGCTAGCTAGCTAAG
AGCTAGCATGCTAGTAGCATGAG
ACGATGCTAGCTAGCTAGCTGATA
TCGATCAGCATGCTACGATGCAAG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

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)

Environmental
Features

READS → PROTEIN FAMILIES → PATHWAYS

$P_1 = f_1 + f_2 + f_3$

$P_2 = f_4 + f_5 + f_6$

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



**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 | | |
| ↓ | ---- | --- | ----- | | |

**Environmental Features**

Chlorophyll    Temp

Pathways

Cobalamin Biosynthesis

Photosystem II

Photosystem I

Carbon Fixation (Dark rx)

Glutamine Degradation

$r^2 = .68$

Predicted Temperature

Actual Temperature

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

# **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



| 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

Environmentally invariant

Environmentally variant

CCA structural correlation

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

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

CCA structural correlation
0    0.3    1

KEGG    Module

Photosynthesis

ATPase complex

Oxidative Phosphorylation

ATPase complex

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

# Conclusion #2: Outer Membrane components vary the environment



[ 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) ]

# 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**

International HapMap Project

SNPs

Map to ENSEMBL genes

**Database of Genomic Variants**

CNVs + SDs

**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**

# 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

Positive Selection Test Likelihood Ratio (y-axis): 0, 0.5, 1, 1.5, 2, 2.5, 3, 3.5, 4, 4.5, 5

Betweenness Centrality (x-axis): 0, 0.5, 1, 1.5, 2, 2.5, 3 × 10^6

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



Spearman Rank P-value: 3.5e−04

Number of Overlapping SDs (y-axis)
Betweenness Centrality (x-axis), x 10^6

Network periphery → Network center

## 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)

Why do we observer this? Perhaps central hub proteins are involved in more interactions & have more surface buried.

**BURIED SITES ARE CONSERVED AND MUCH LESS LIKELY TO HARBOR NON-SYNONYMOUS MUTATIONS**

dN/dS Ratio

0.49

0.35

$p \ll 0.01$

**Exposed sites**

**Buried sites**

Average Relative Surface Exposure

2.66

2.26

$p \ll 0.01$

**Site with Synonymous Mutations only**

**Sites with Non-synonymous Mutations**

# Another explanation: THE NETWORK PERIPHERY CORRESPONDS TO THE CELLULAR PERIPHERY



| | Betweenness Centrality (x $10^4$) | Degree Centrality |
|---|---|---|
| **Chromosome** | 5.5 | 10 |
| **Nucleus** | 5.0 | 8.6 |
| **Cytoplasm** | 5.2 | 8.1 |
| **Membrane** | 4.0 | 6.5 |
| **Extracellular Region** | 3.8 | 5.9 |

Legend:
- Extracellular
- Plasma membrane
- Cytoplasm
- Mitochondria
- Nucleus
- Centrosome
- Endosome
- Ribosome
- Lysosome
- Peroxisome
- Golgi apparatus
- Endoplasmic reticulum
- Other organelles/unknown

Source: Gandhi et al. (*Nature Genetics* 2006), 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?
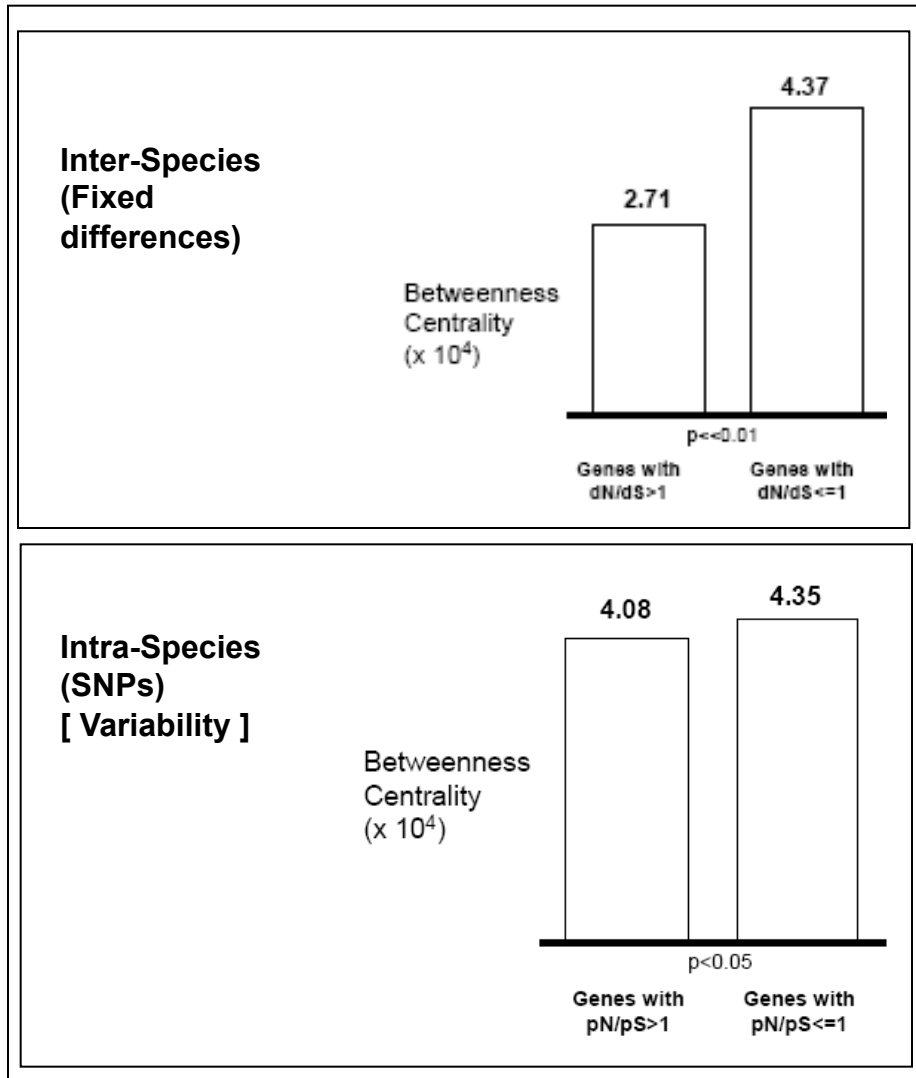
|  | **Relaxed Constraint** | **Adaptive Evolution** |
|---|---|---|
| **Inter-Species Variation (Fixed differences)** | • Increases inter-species variation – more variable loci are under less negative selection<br><br>• Can be seen in higher Ka/Ks ratio or SD occurrence | • Increases inter-species variation – 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 $10^4$)

p<<0.01

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

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



4.08    4.35

Betweenness Centrality (x $10^4$)

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*
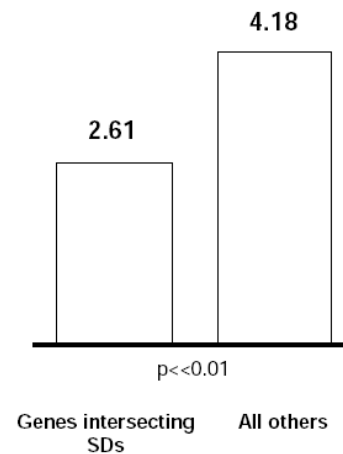
\* But it's hard to quantify

Source: Kim et al. (2007) PNAS

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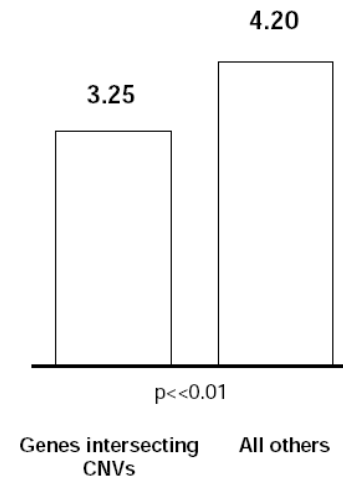
## Inter vs. Intra-Species Variation in Networks

**Inter-Species (SDs)**

4.18

2.61

Betweenness Centrality (x 10$^4$)

p<<0.01

Genes intersecting SDs

All others

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

4.20

3.25

Betweenness Centrality (x 10$^4$)
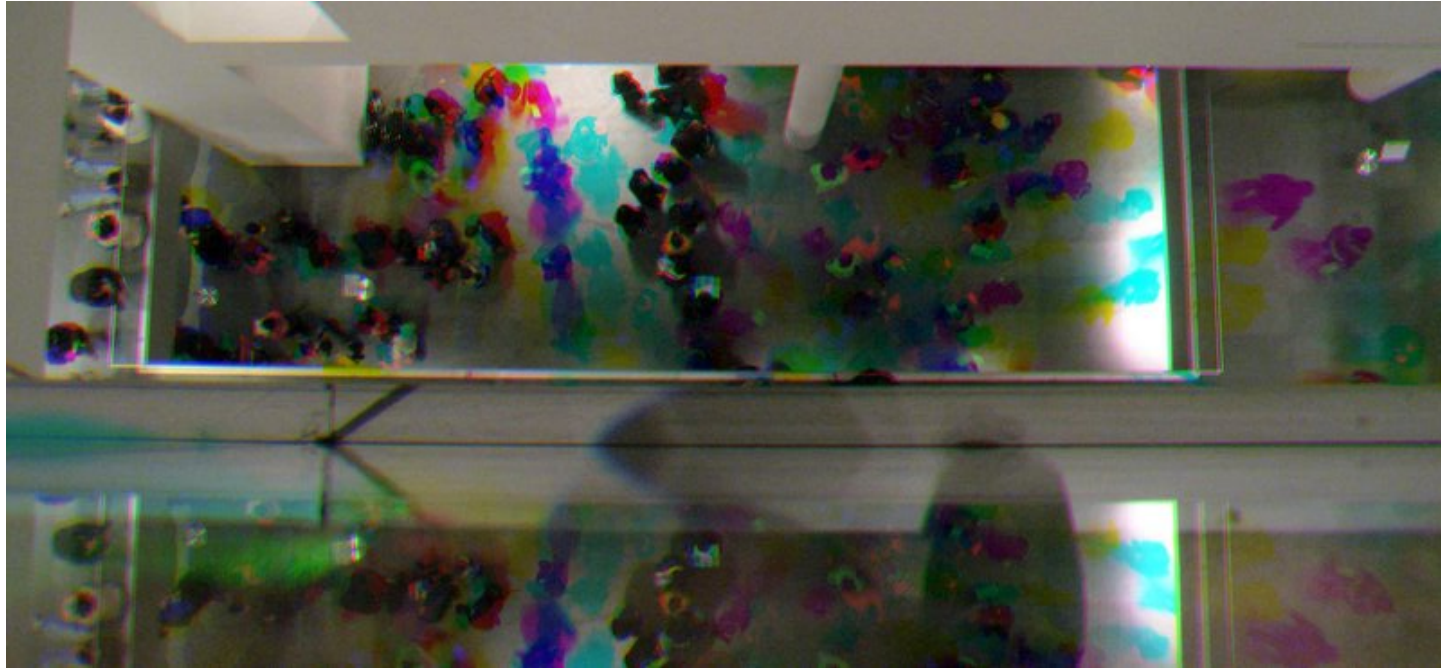
p<<0.01

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
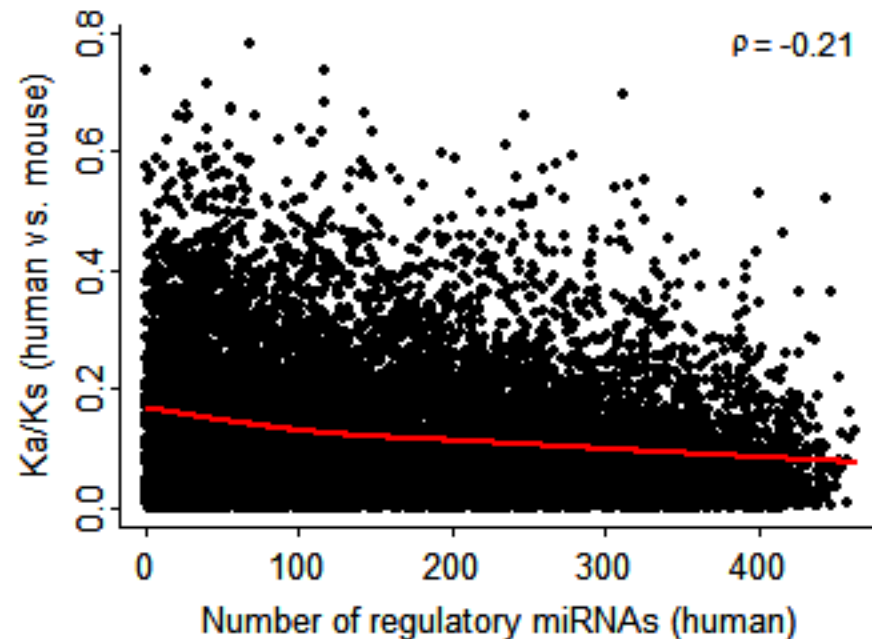
# Networks & Variation 2

**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



$\rho = -0.21$

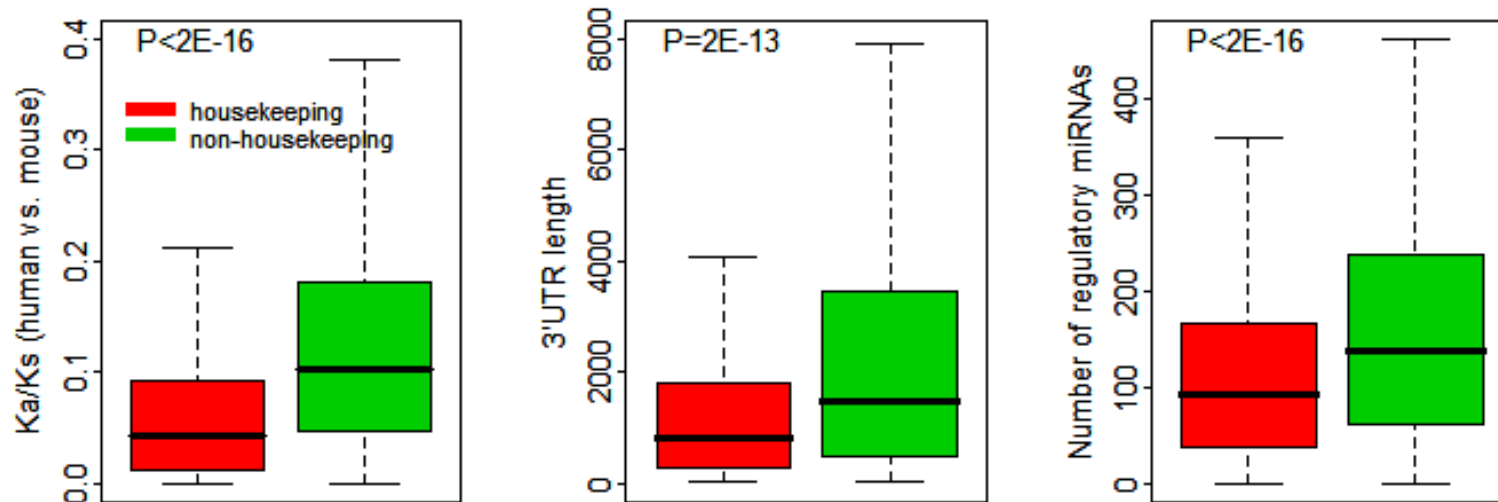**Important genes are regulated more intensively regulated by the microRNAs**

| 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)]

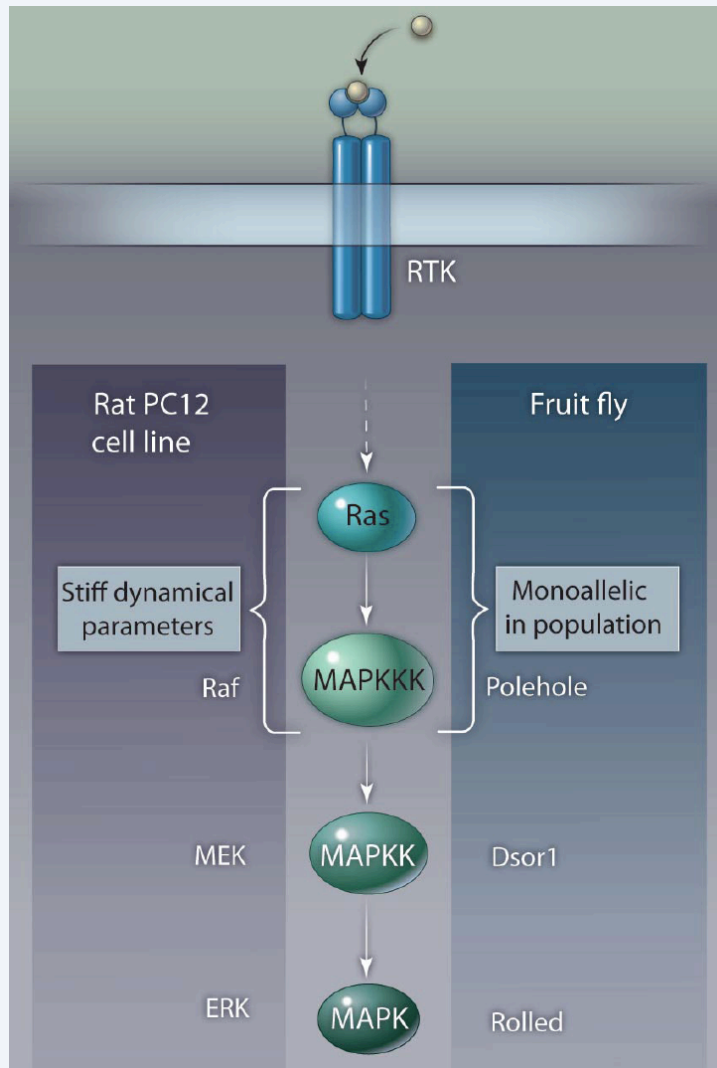# MicroRNA regulation: a two-way strategy

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

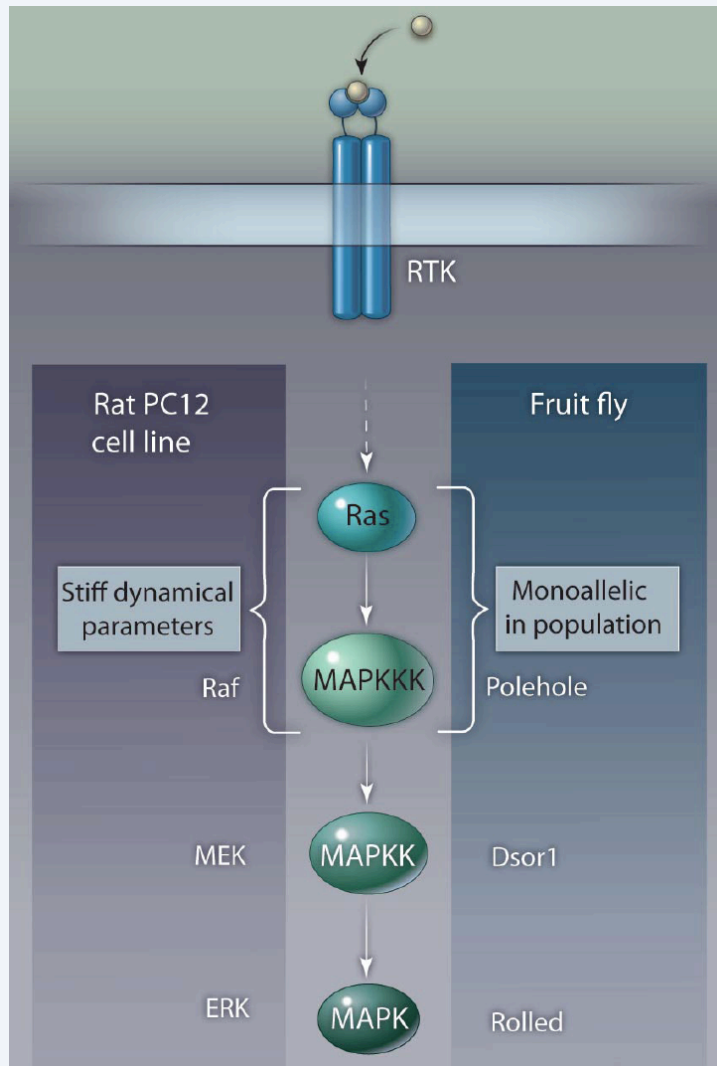housekeeping genes, however conserved, are selected to have shorter 3'UTRs to avoid miRNA regulation.



[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**
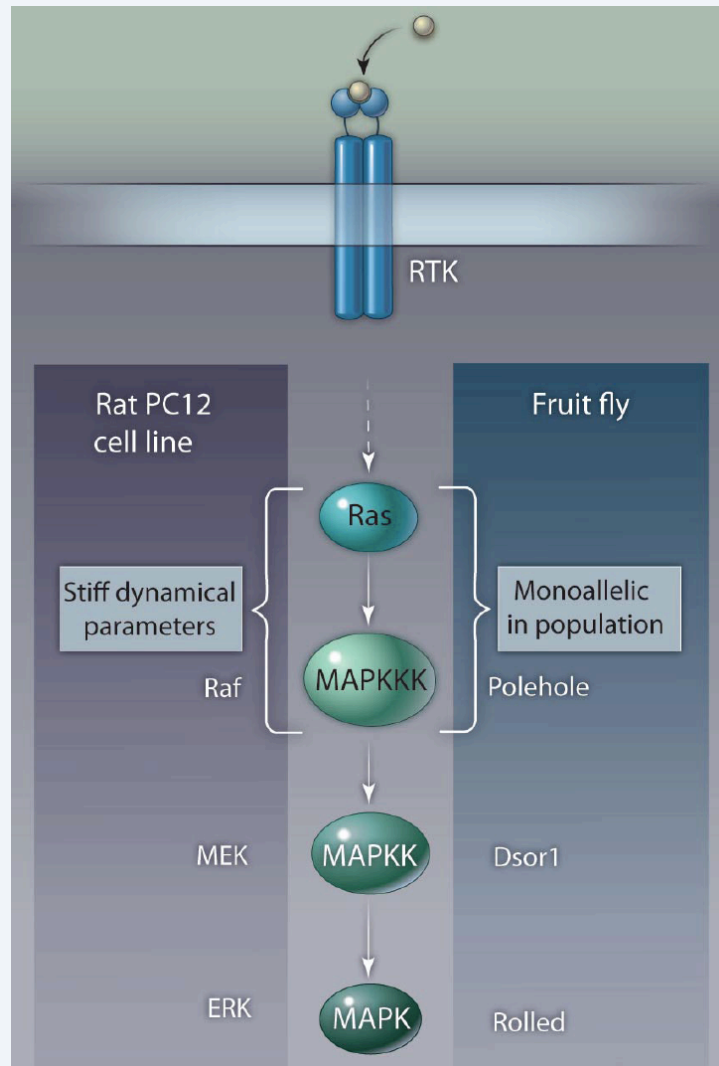
# Network dynamics constrain evolution



**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**

# Analogies show it reasonable for more variable part of network to be periphery

- Computer Networks
  - Servers in center have much depending on them; thus, can't be frequently updated & patched
  - Servers on periphery often attacked and so need frequent patches
- Social Networks
  - Individuals at center under more constraint (to conform), whereas those at periphery have more freedom to experiment

# Outline: Molecular Networks

- ## Why Networks?

- ## Predicting Networks (yeast ppi)
  - ◊ Propagating known information

- ## 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
    (in prokaryote metab. pathways)

- ## Protein Networks & Variation
  (human ppi & miRNA-targ. net)

# Conclusions on Networks: Predictions

- Predicting Networks
  - ◊ Extrapolating from the Training Set
  - ◊ Principled ways of using known information in the fullest possible fashion
    - Prediction Propagation
    - 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 and essentiality

  ◊ Topmost proteins sit at center of interaction network

# Conclusions:
# Points of Network Centrality



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

- Application to miRNA network

- Different RE-score of miRNAs can be used 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



- We find ongoing evolution (positive selection) at the network periphery.
  - ◊ This trend is present on two levels:
    - 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
  - ◊ 2 possible mechanisms for this : adaptive evolution at cellular periphery & relaxation of structural constraints at the network periphery
    - We show that the latter can only explain part of the increased variability

# Conclusions: Connecting Networks & Variation 2



- More highly regulated genes are under more constraint in miRNA-target networks

- Exception for housekeeping genes

- Speculation as to why variation at periphery is quite reasonable

# 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

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

A Paccanaro
P Alves
T Emonet
P Cayting
M Seringhaus
Y Xia
J Korbel
A Sboner
P Patel
P Bork
J Raes
E Franzosa
M Snyder
N Bhardwaj
R Alexander



**Networks.GersteinLab.org**

Job opportunities currently for postdocs & students

# More Information on this Talk

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

**SUBJECT:** Networks

**DESCRIPTION**:
Network Biology: Understanding metabolic and protein interactions, Mathematical Biosciences Institute, Columbus, OH; 2009.09.14, 13:30-14:30; [I:**MBINETS**] (Long networks talk, adding in for the first time: **rescore***, **mirnatargevolrate*** & **netdynamicsrev***. Fits easily into 55' w. 5' questions. PPT works on mac & PC and has many photos.)

(Paper references in the talk were mostly from Papers.GersteinLab.org. The above topic list can be easily cross-referenced against this website. Each topic abbrev. which is starred is actually a papers "ID" on the site. For instance,
the topic **pubnet*** can be looked up at
**http://papers.gersteinlab.org/papers/pubnet** )