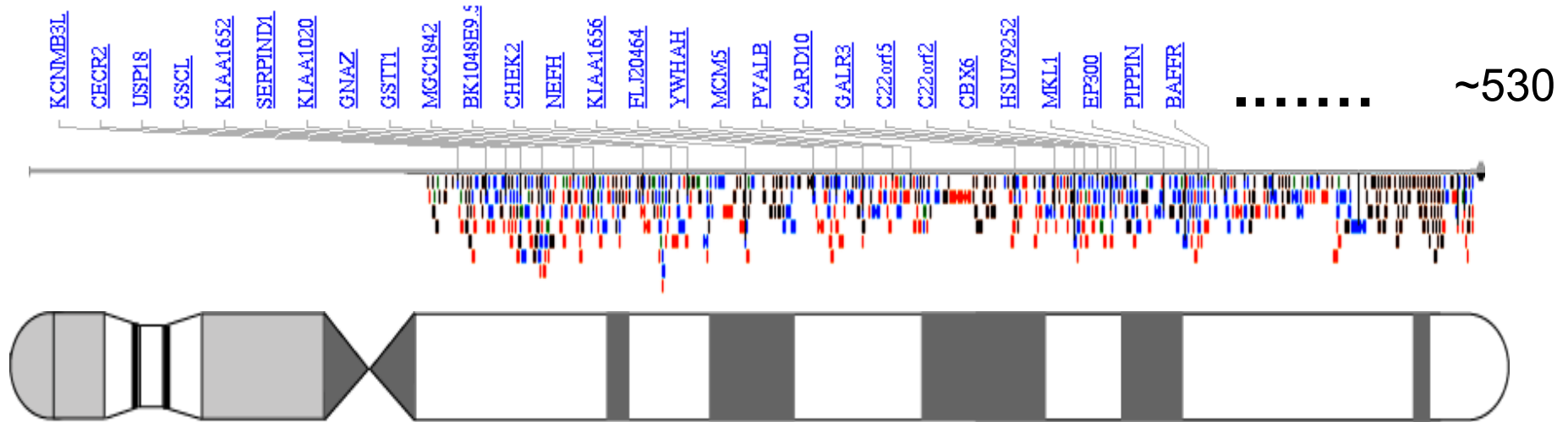# Understanding Protein Function on a Genome-scale through the Analysis of Molecular Networks

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?



~530

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

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

# EF2_YEAST

## Traditional single molecule way to integrate evidence & describe function

**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
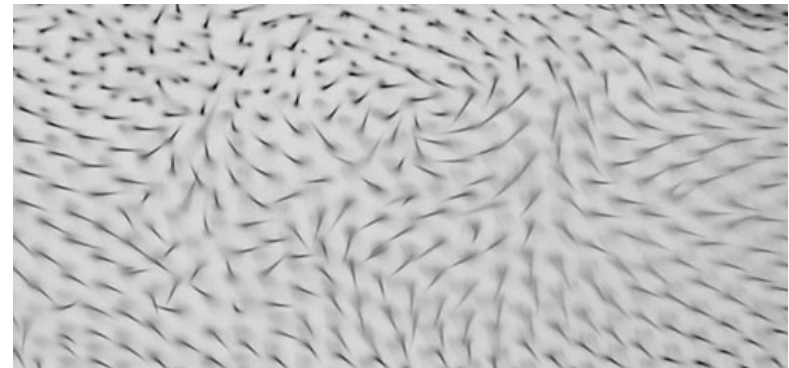    - ◊ Role Conflation:
      molecular, cellular, phenotypic

# 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?....
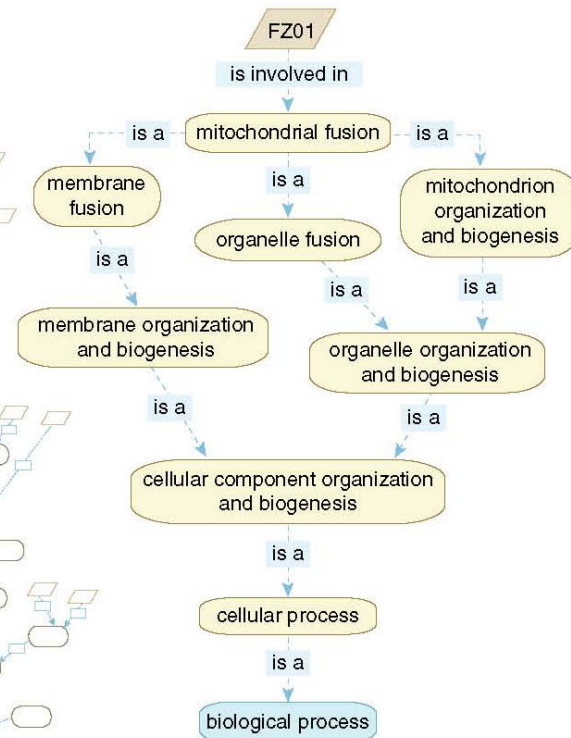    - ◊ **Starry night** (P Adler, '94)



[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)
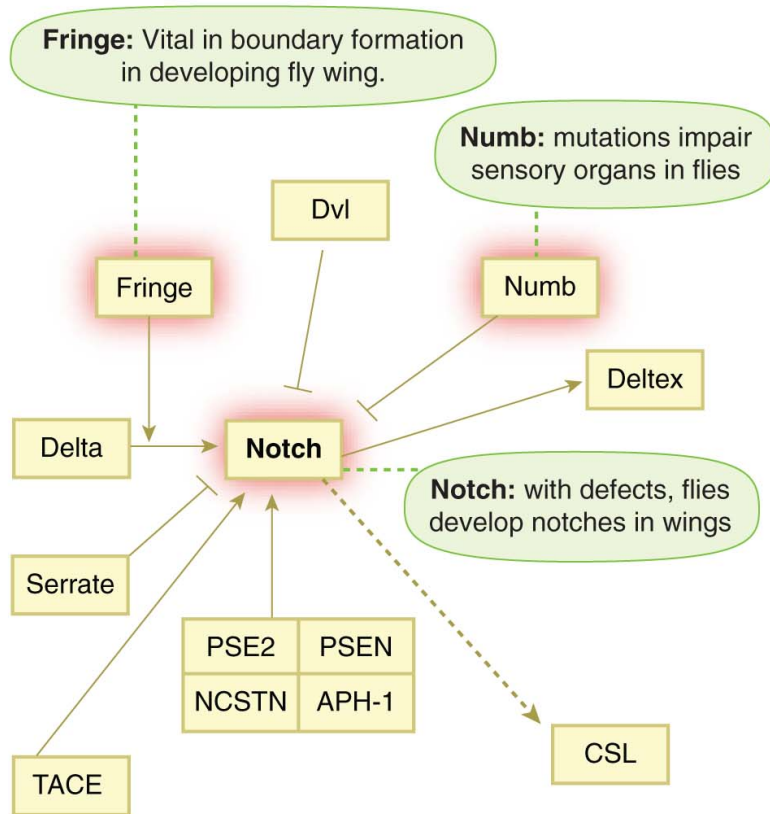
- Function as a vector

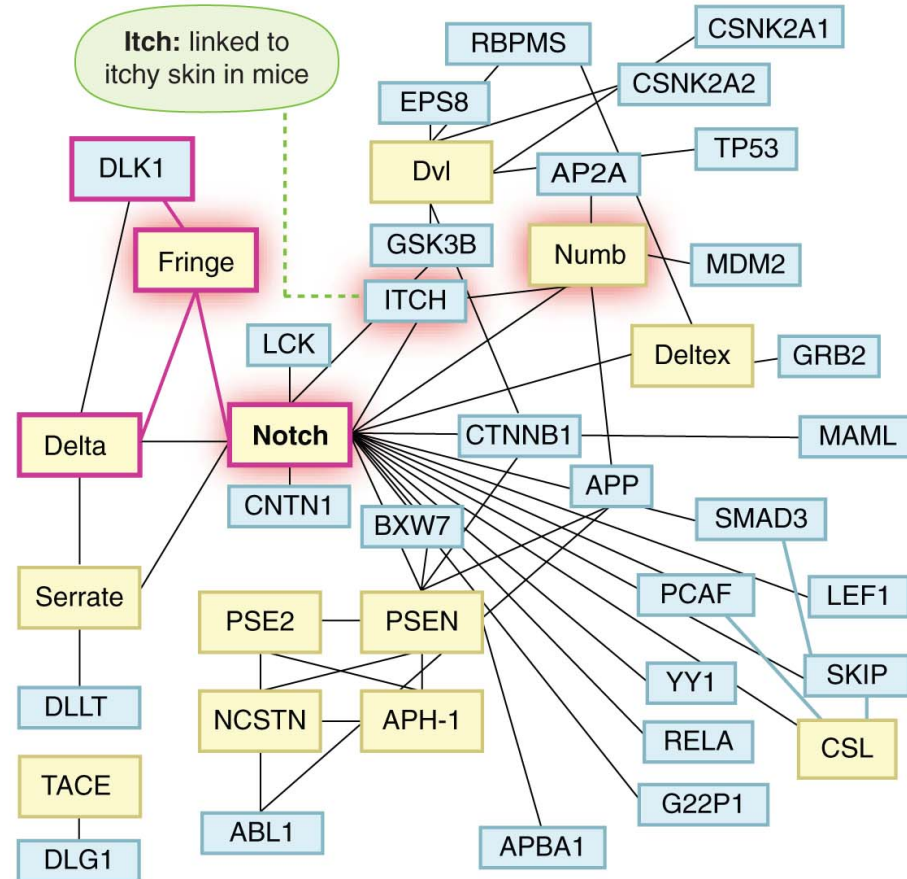|  | nucleic acids | | small molecules | | | | | proteins | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
|  | DNA | RNA | ATP | Metal | CoA | NAD | ...... | G protein | CDC28 | Calmodulin | ...... |
| protein 1 | 1.0 | 0 | 0 | 0 | 0 | 0 | ...... | 0 | 0 | 0 | ...... |
| protein 2 | 0 | 0.9 | 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)



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]

# 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]
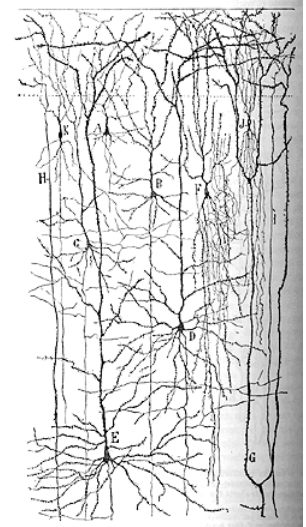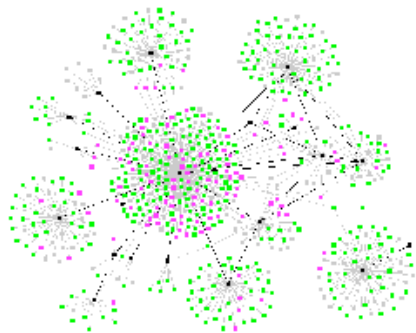
# Networks as a universal language



Internet
[Burch & Cheswick]

Food Web

Electronic
Circuit

Neural Network
[Cajal]

Disease
Spread
[Krebs]

LINKED

Albert-László
Barabási

The New Science
of Networks

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

Protein
Interactions
[Barabasi]

Social Network

# Networks as a Central Theme in Systems Biology



**Reductionist Approach**

**Integrative Approach**

**Systems Biology**

**[Adapted from H Yu]**

# Network pathology & pharmacology



**Interactome networks**

[Adapted from H Yu]

# Using the position in networks to describe function



Blame Game

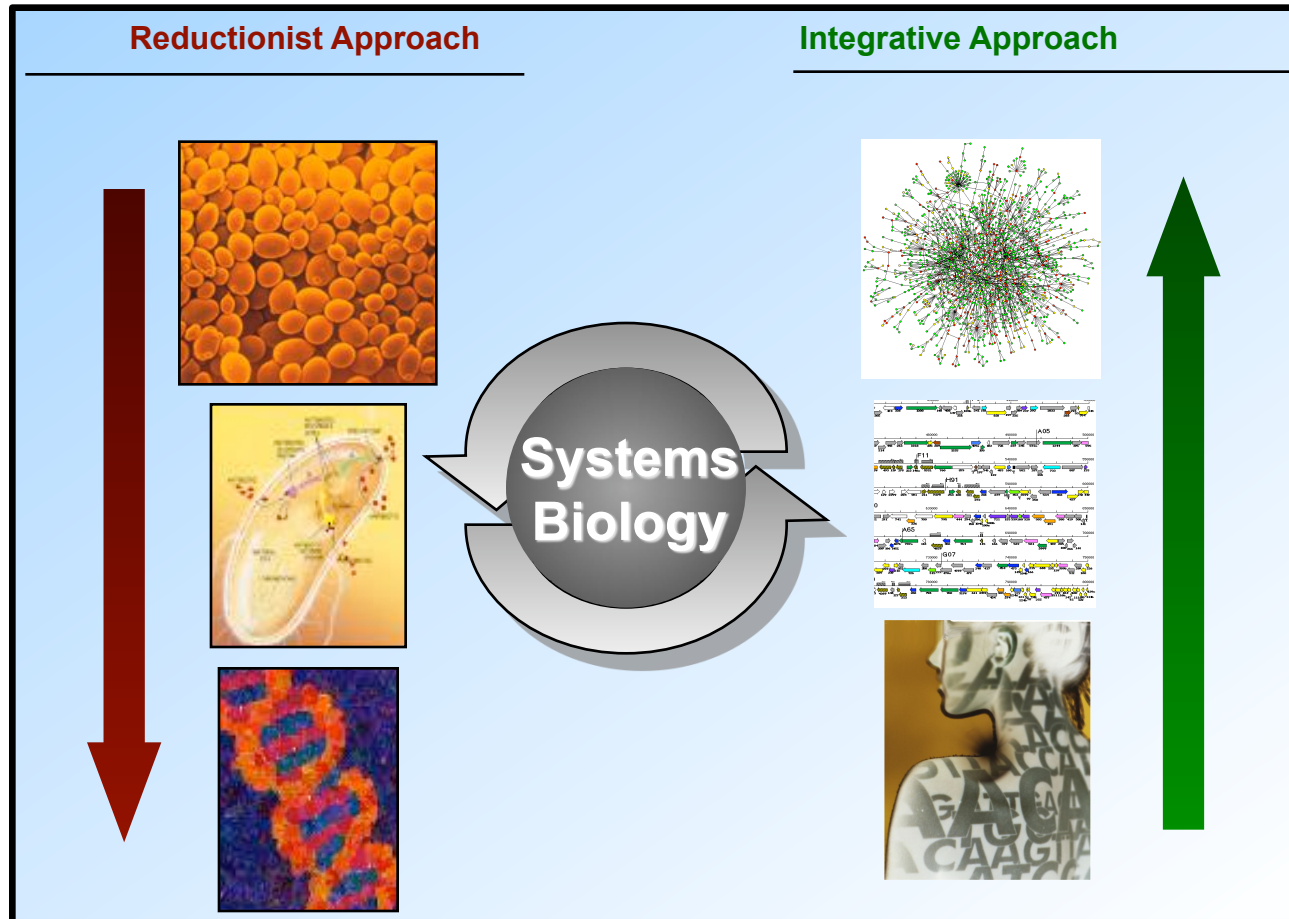When failures made the disaster in New Orleans worse, everybody found fault with somebody. And the fingers haven't stopped pointing. Last week, Michael Brown, the former head of FEMA, had his turn in a Congressional hearing. Here is a sampling of notables and their views.

Bill Marsh/The New York Times

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

# **Types of Networks**



Interaction networks



Regulatory networks



Metabolic 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]

# 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

- Predicting Networks
  - ◊ Training set expansion
- Properties of Protein Networks
  - ◊ Hubs
- Dynamics of Networks
  - ◊ Dynamics across cellular states
  - ◊ Dynamics across environments
- Protein Networks and Human Variation

# 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.**

# Network Prediction

- Only small portions are already known
- Many other kinds of data available
- → Use them to learn models for predicting the unknown portions

**Known**

**New**



**Ex. of Predicted Membrane Protein Interactome in Xia et al. JMB (2006)**

**Figure 6:** A map of known and a subset of predicted interactions among helical membrane proteins. Nodes represent helical

# 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)

# Learning

Concepts in machine learning:

- Training sets:
  - ◊ Positive set: known interactions
  - ◊ Negative set: known non-interactions

- Features:
  - ◊ Data describing the objects

- Model:
  - ◊ A function that takes two objects as input and predicts whether they interact

# Training sets



Known interactions

Known non-interactions

Unknown

# Features

- Example 1: gene expression



Gasch et al., 2000

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

**Similarity scale:**

**1** ——————————— **-1**

# **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

# **Evaluation**

- Computational:
    - ◊ Cross-validation
    - ◊ Indirect evidence (e.g. same GO category)
- Experimental:
    - ◊ Validation of de novo predictions

# **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 fairly in a public challenge! (DREAM)

# Kernels

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



**Objects in an feature space**

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 |

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

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



?

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?)

# 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 |

# 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



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

# 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

# **Experiments**

- Gold-standard interactions: BioGRID, from studies that report less than 10 interactions
- Features:

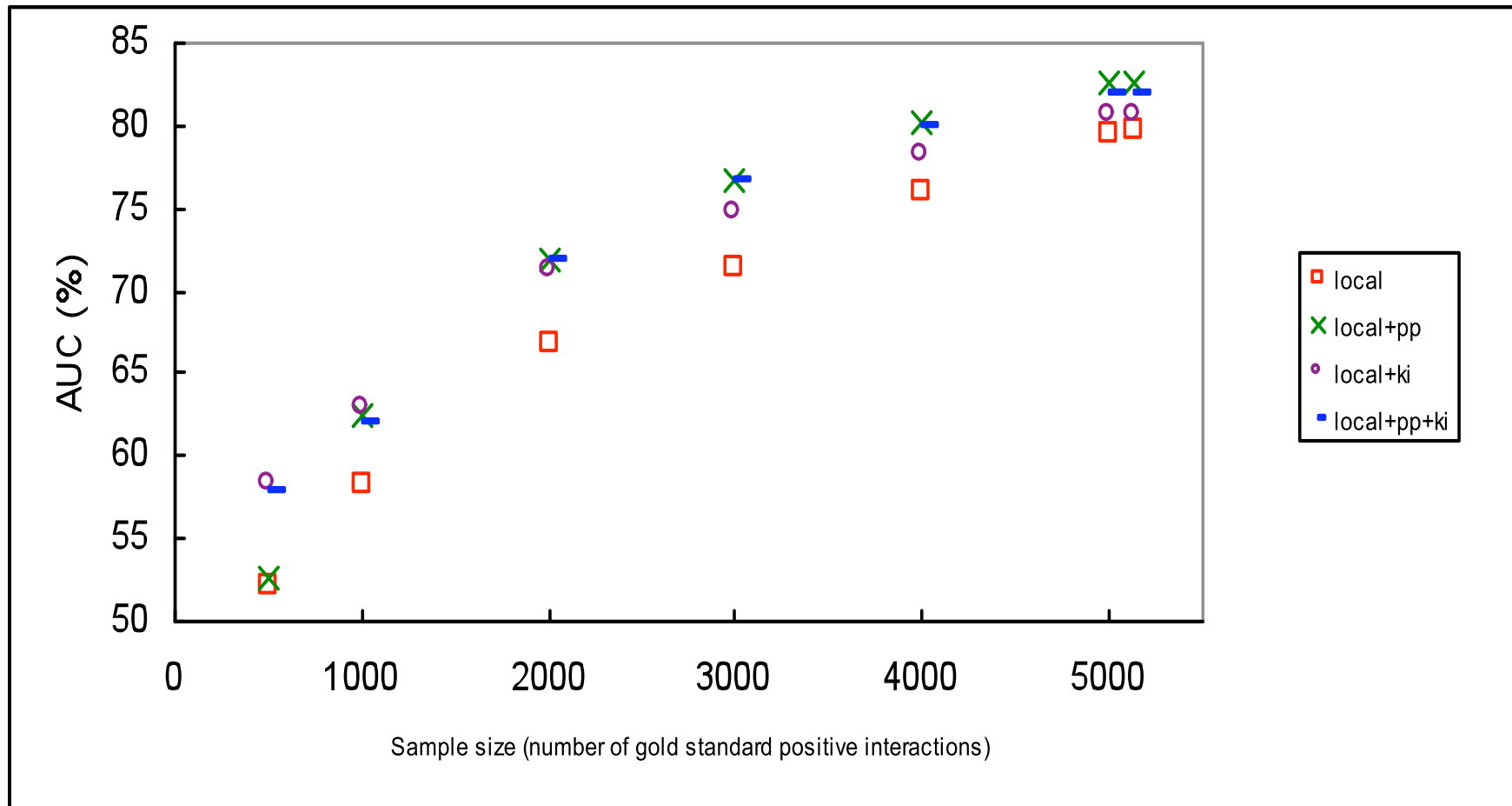| Code | Data type | Source | Kernel |
|------|-----------|--------|--------|
| phy | Phylogenetic profiles | COG v7 (Tatusov et al., 1997) | RBF ($\sigma$=3,8) |
| loc | Sub-cellular localization | (Huh et al., 2003) | Linear |
| exp-gasch | Gene expression (environmental response) | (Gasch et al., 2000) | RBF ($\sigma$=3,8) |
| exp-spellman | Gene expression (cell-cycle) | (Spellman et al., 1998) | RBF ($\sigma$=3,8) |
| y2h-ito | Yeast two-hybrid | (Ito et al., 2000) | Diffusion ($\beta$=0.01) |
| y2h-uetz | Yeast two-hybrid | (Uetz et al., 2000) | Diffusion ($\beta$=0.01) |
| tap-gavin | Tandem affinity purification | (Gavin et al., 2006) | Diffusion ($\beta$=0.01) |
| tap-krogan | Tandem affinity purification | (Krogan et al., 2006) | Diffusion ($\beta$=0.01) |
| int | Integration | | Summation |

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

# Prediction accuracy

| | 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
- Overfitting of local modeling without training set expansion
- Comparing prediction propagation and kernel initialization

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

# Complementarity of the two methods



[**Yip and Gerstein, Bioinformatics ('09, in press)**]
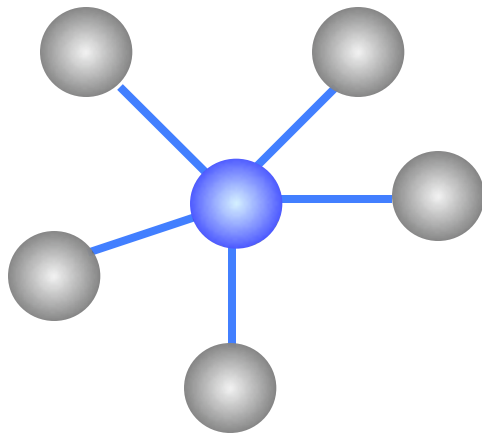
# Network Dynamics #1: Cellular States

**How do networks change across different cellular states?**
**How can this be used to assign function to a protein?**

# 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 $P(k)$

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

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)**

**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"

# *Dynamic* Yeast TF network

Transcription Factors



Target Genes

- Analyzed network as a static entity

- But network is *dynamic*
  ◊ Different sections of the network are active under different cellular conditions

- Integrate gene expression data

Luscombe et al. Nature 431: 308

# Gene expression data for five cellular conditions in yeast

| Cellular condition |
| --- |
| Cell cycle |
| Sporulation |
| Diauxic shift |
| DNA damage |
| Stress response |

Multi-stage { Cell cycle, Sporulation

Binary { Diauxic shift, DNA damage, Stress response

[Brown, Botstein, Davis….]

# Backtracking to find active sub-network



- Define differentially expressed genes

- Identify TFs that regulate these genes

- Identify further TFs that regulate these TFs

Active regulatory sub-network

# Network usage under different conditions

## static



Luscombe et al. Nature 431: 308

# Network usage under different conditions

## cell cycle



(c) 2008

# Network usage under different conditions

## sporulation

# Network usage under different conditions

## diauxic shift

# Network usage under different conditions

## DNA damage

# Network usage under different conditions

## stress response

# Network usage under different conditions

**Cell cycle**   **Sporulation**   **Diauxic shift**   **DNA damage**   **Stress**



## SANDY:

**1. Standard graph-theoretic statistics:**
- Global topological measures
- Local network motifs

**2. Newly derived follow-on statistics:**
- Hub usage
- Interaction rewiring

**3. Statistical validation of results**

Luscombe et al. Nature 431: 308

# Network usage under different conditions

**Cell cycle**  **Sporulation**  **Diauxic shift**  **DNA damage**  **Stress**



### SANDY:
**1. Standard graph-theoretic statistics:**
- Global topological measures
- Local network motifs

**2. Newly derived follow-on statistics:**
- Hub usage
- Interaction rewiring

**3. Statistical validation of results**

**Analysis of condition-specific subnetworks in terms of global topological statistics**

Outdegree
- Cell cycle: 7.9
- Sporulation: 6.5
- Diauxic shift: 17.1
- DNA damage: 15.0
- Stress response: 9.0

Indegree
- Cell cycle: 2.0
- Sporulation: 1.9
- Diauxic shift: 1.6
- DNA damage: 1.6
- Stress response: 1.6

Pathlength
- Cell cycle: 4.5
- Sporulation: 3.4
- Diauxic shift: 2.1
- DNA damage: 2.0
- Stress response: 2.2

Clustering coefficient
- Cell cycle: 0.15
- Sporulation: 0.14
- Diauxic shift: 0.09
- DNA damage: 0.09
- Stress response: 0.08

Cell cycle | Sporulation | Diauxic shift | DNA damage | Stress response
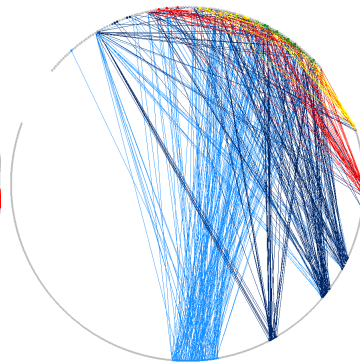
Multi-stage
Controlled, ticking
over of genes
at different stages

Binary
Quick, large-scale
turnover of genes

Luscombe et al. Nature 431: 308

Single-input module

Multi-input module

Feed-forward loop

32%  39%  57%  56%  59%

24%  17%  24%  27%  20%

44%  45%  19%  17%  21%

Cell cycle

Sporulation

Diauxic shift

DNA damage

Stress response

Multi-stage
Controlled, ticking
over of genes
at different stages

Binary
Quick, large-scale
turnover of genes

**Analysis of condition-specific subnetworks in terms of occurrence of local motifs**

Luscombe et al. Nature 431: 308

**Cell cycle** **Sporulation** | **Diauxic shift** **DNA damage** **Stress**

multi-stage conditions | binary conditions

## Summary

| | | |
|---|---|---|
| less pronounced | Hubs | more pronounced |
| longer | Path Lengths | shorter |
| more | TF inter-regulation | less |
| complex loops (FFLs) | Motifs | simpler (SIMs) |

# Transient Hubs



Regulatory hubs

- Questions:
  - ◊ Do hubs stay the same or do they change over between conditions?
  - ◊ Do different TFs become important?

- Our Expectations
  - ◊ Literature:
    - • Hubs are permanent features of the network regardless of condition
  - ◊ Random networks (sampled from complete regulatory network)
    - • Random networks converge on same TFs
    - • 76-97% overlap in TFs classified as hubs (*ie* hubs are permanent)

Luscombe et al. Nature 431: 308

cell cycle, sporulation, diauxic shift, DNA damage, stress response

cell cycle
- YMR016C
- YLR183C
- YIL131C
- SWI4
- YDR451C
- SWI6
- STE12
- MBP1
- MCM1
- YDR146C
- YLR131C

sporulation
- UME6
- IME1
- YNL216W
- SIN3
- YIR023W
- YPL038W
- YNL103W
- YMR021C
- CBF1
- YBL021C
- YIL122W

diauxic shift
- HAP4
- HAP2

DNA damage
- YHR206W
- YAP1
- HSF1
- YPL089C
- YCR065W
- CIN5
- YDR310C

stress response
- YDR259C
- MSN2
- YDR501W
- MSN4
- YGL096W
- PDR1
- YLR403W
- YGL071W
- YIR018W

all conditions
- YKL043W
- YLR013W
- YGL209W
- YML027W
- YFR034C
- YEL009C
- YBR049C
- YGL035C
- YKL112W
- YDR043C
- YPR065W

transient hubs

permanent hubs

- Some permanent hubs
  ◊ house-keeping functions

- Most are transient hubs
  ◊ Different TFs become key regulators in the network

- Implications for condition-dependent vulnerability of network

Luscombe et al. Nature 431: 308

cell cycle
sporulation
diauxic shift
DNA damage
stress response

**cell cycle**

YMR016C
YLR183C
YIL131C
SWI4
YDR451C
SWI6
STE12
MBP1
MCM1
YDR146C
YLR131C

Swi4, Mbp1

**sporulation**

UME6
IME1
YNL216W
SIN3
YIR023W
YPL038W
YNL103W
YMR021C
CBF1
YBL021C
YIL122W

Ime1, Ume6

**diauxic shift**

HAP4
HAP2

**DNA damage**

YHR206W
YAP1
HSF1
YPL089C
YCR065W
CIN5
YDR310C

**stress response**

YDR259C
MSN2
YDR501W
MSN4
YGL096W
PDR1
YLR403W
YGL071W
YIR018W

Msn2, Msn4

**all conditions**

YKL043W
YLR013W
YGL209W
YML027W
YFR034C
YEL009C
YBR049C
YGL035C
YKL112W
YDR043C
YPR065W

Luscombe et al. Nature 431: 308

cell cycle
sporulation
diauxic shift
DNA damage
stress response

**cell cycle**

YMR016C
YLR183C
YIL131C
SWI4
YDR451C
SWI6
STE12
MBP1
MCM1
YDR146C
YLR131C

**sporulation**

UME6
IME1
YNL216W
SIN3
YIR023W
YPL038W
YNL103W
YMR021C
CBF1
YBL021C
YIL122W

**diauxic shift**

HAP4
HAP2

**DNA damage**

YHR206W
YAP1
HSF1
YPL089C
YCR065W
CIN5
YDR310C

**stress response**

YDR259C
MSN2
YDR501W
MSN4
YGL096W
PDR1
YLR403W
YGL071W
YIR018W

**all conditions**

YKL043W
YLR013W
YGL209W
YML027W
YFR034C
YEL009C
YBR049C
YGL035C
YKL112W
YDR043C
YPR065W

Unknown functions

Luscombe et al. Nature 431: 308

# 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

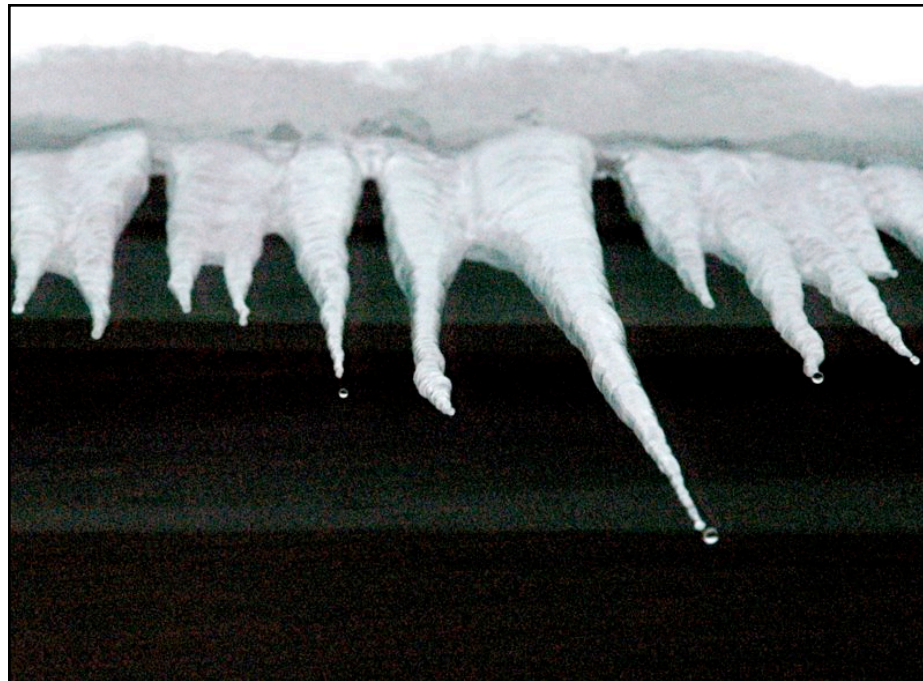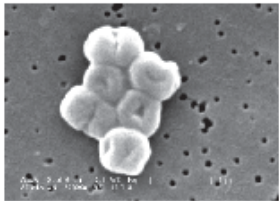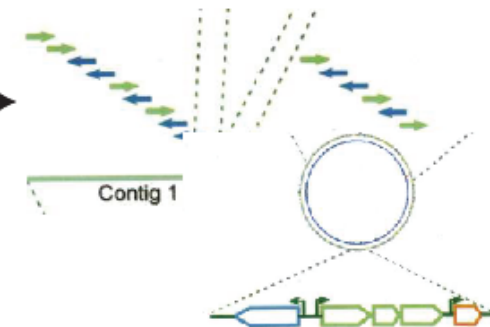**Culture Microbes** → **Extract DNA** → **Sequence** → **Assemble and Annotate**

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.

# Comparative Metagenomics



Water

Soil

Do the proportions of pathways represented in these two samples differ?

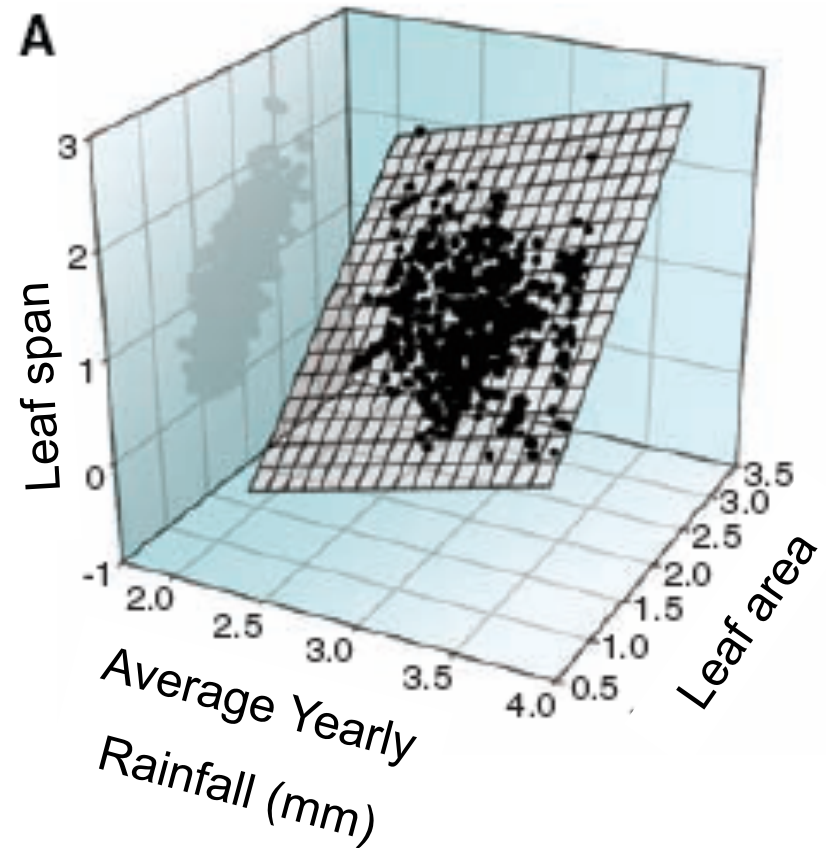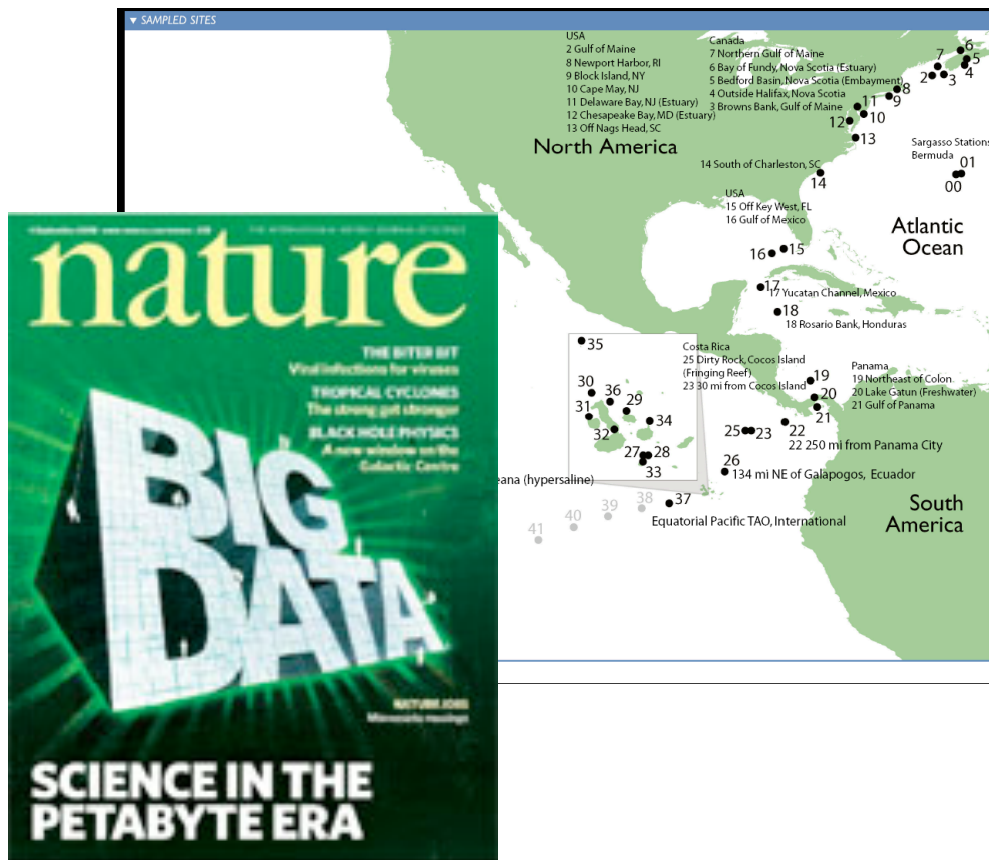Dinsdale et. al., Nature 2008

# Trait-based Biogeography

Charles River, MA

Long Island Sound, CT

Do the proportions of pathways represented in these two samples CHANGE as a function of their environments?



Green et. al., Science 2008

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

Sites →

## Environmental Features

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

Sites →

---

**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



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

# Canonical Correlation Analysis: Simultaneous weighting



UPI = a GRE + b 📚 + c GPA

GPI = a' 📚 + b' PowerPoint + 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



[ 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) ]

# Conclusion #4: Cofactor (Metal) Optimization

**Methionine salvage, synthesis, and uptake, transport**

**IS DEPENDENT-ON**

**Methionine synthesis**

Cobalamin biosynthesis

Cobalt transporters



Methionine

**IS NEEDED FOR**

**Methionine degradation**

S-adenosyl Methionine Biosynthesis (synthesize SAM one of the most important methyl donors)

Polyamine biosynthesis

**RELIES ON**

**Methionine Salvage**

Spermidine/Putrescine transporters

Arg/His/Ornithine transporters

[ 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**

**Database of Genomic Variants**

SNPs

CNVs + SDs

Map to ENSEMBL genes

**Ensembl Genes**

ENSG000XXXX:
rsSNP00XXX
CNV_XXX
DN/DS XXXX
Recombination rate

**Interactome**

~30000 interactions from HPRD and Y2H screens

Map to proteins in the interaction network

**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**

# CENTRAL PROTEINS ARE LESS LIKELY TO BE UNDER POSITIVE SELECTION

Hubs

## Degree vs. Positive Selection

Spearman Rank P-value: 1.2e-06

(Y-axis: Positive Selection Test Likelihood Ratio, 0 to 5)
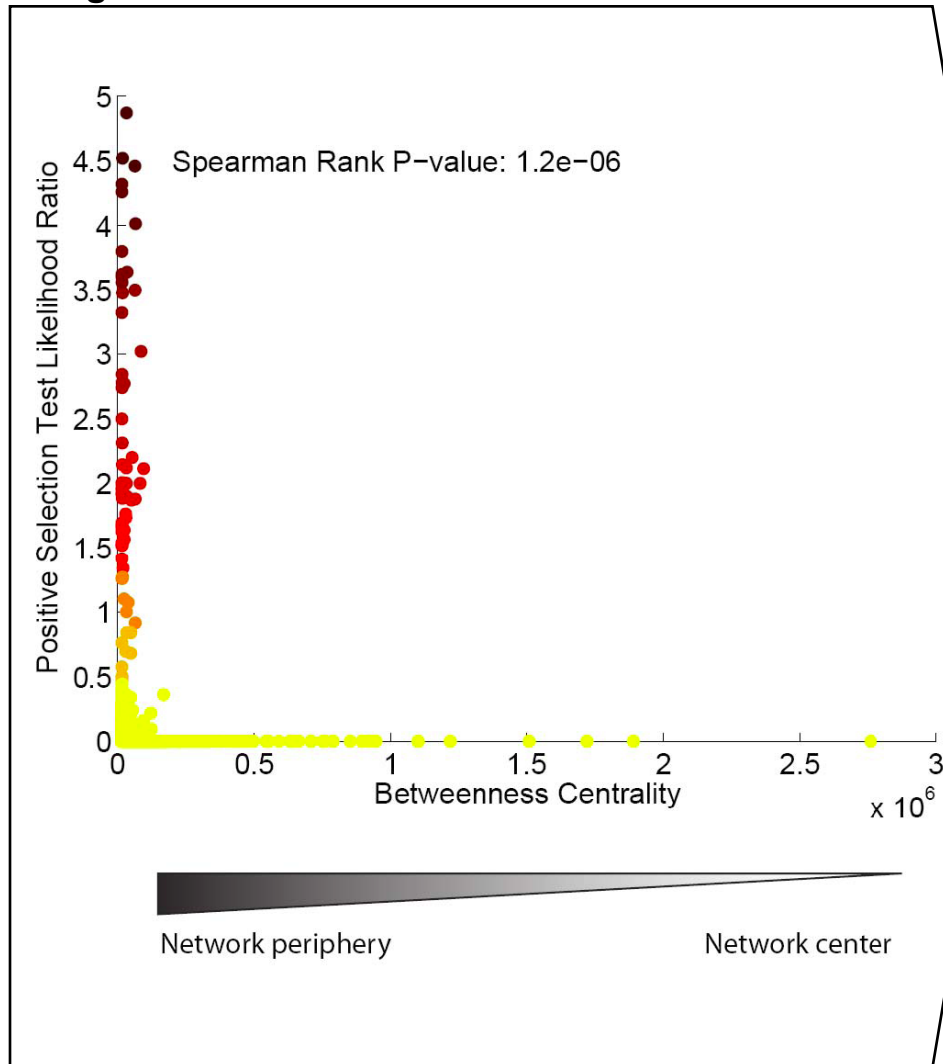(X-axis: Betweenness Centrality, 0 to 3 x 10^6)

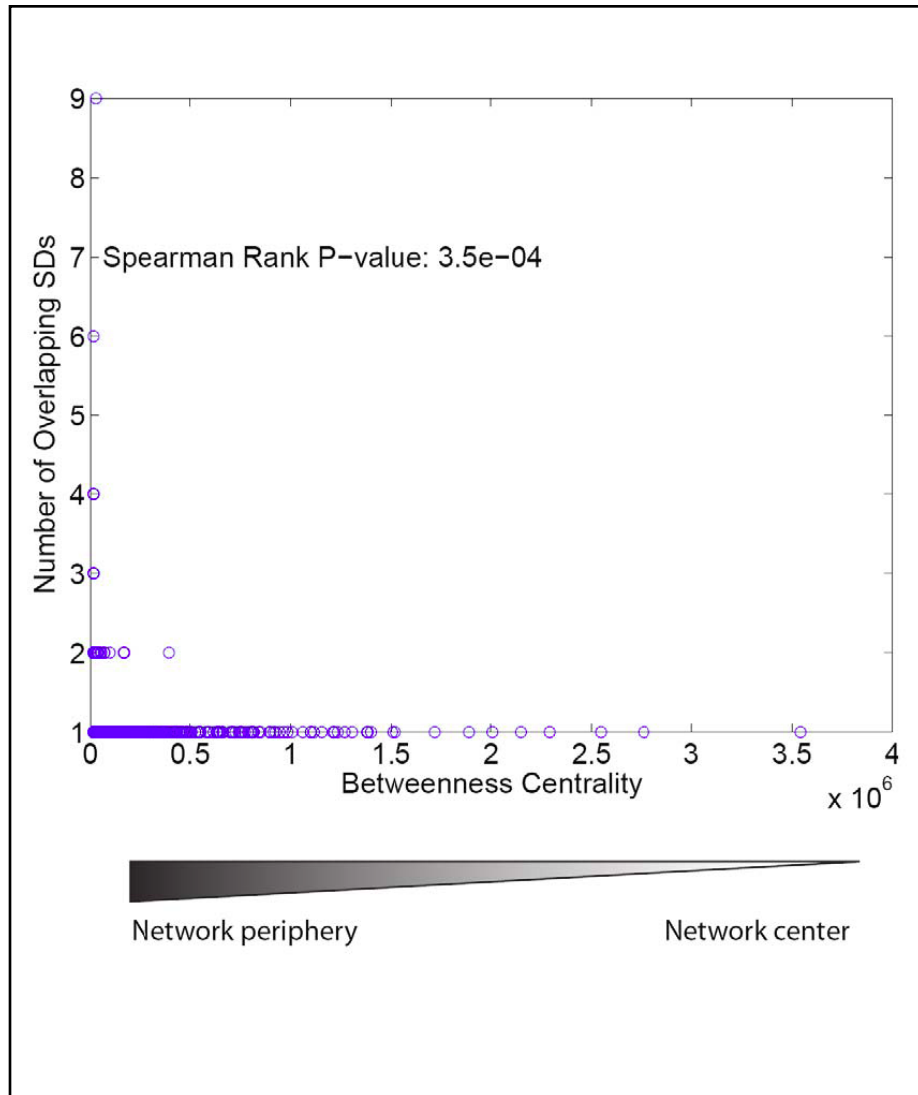Network periphery — Network center

## 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
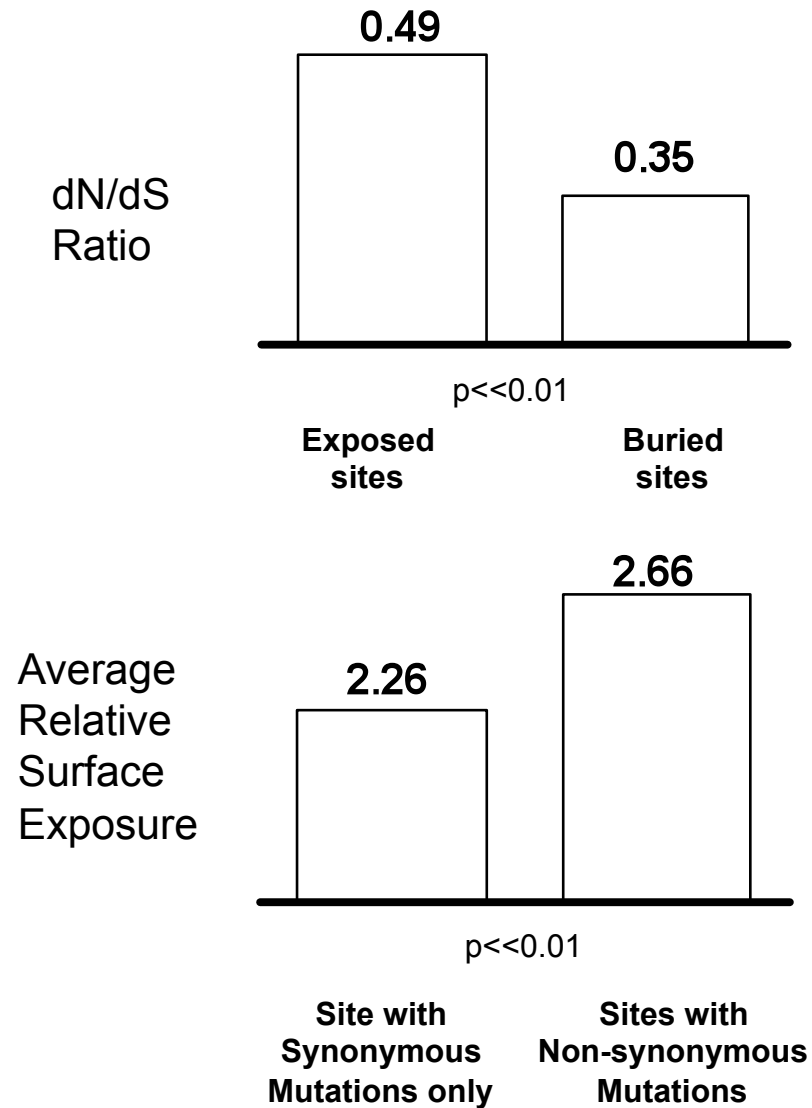
**Centrality vs. SD occurrence**



**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*

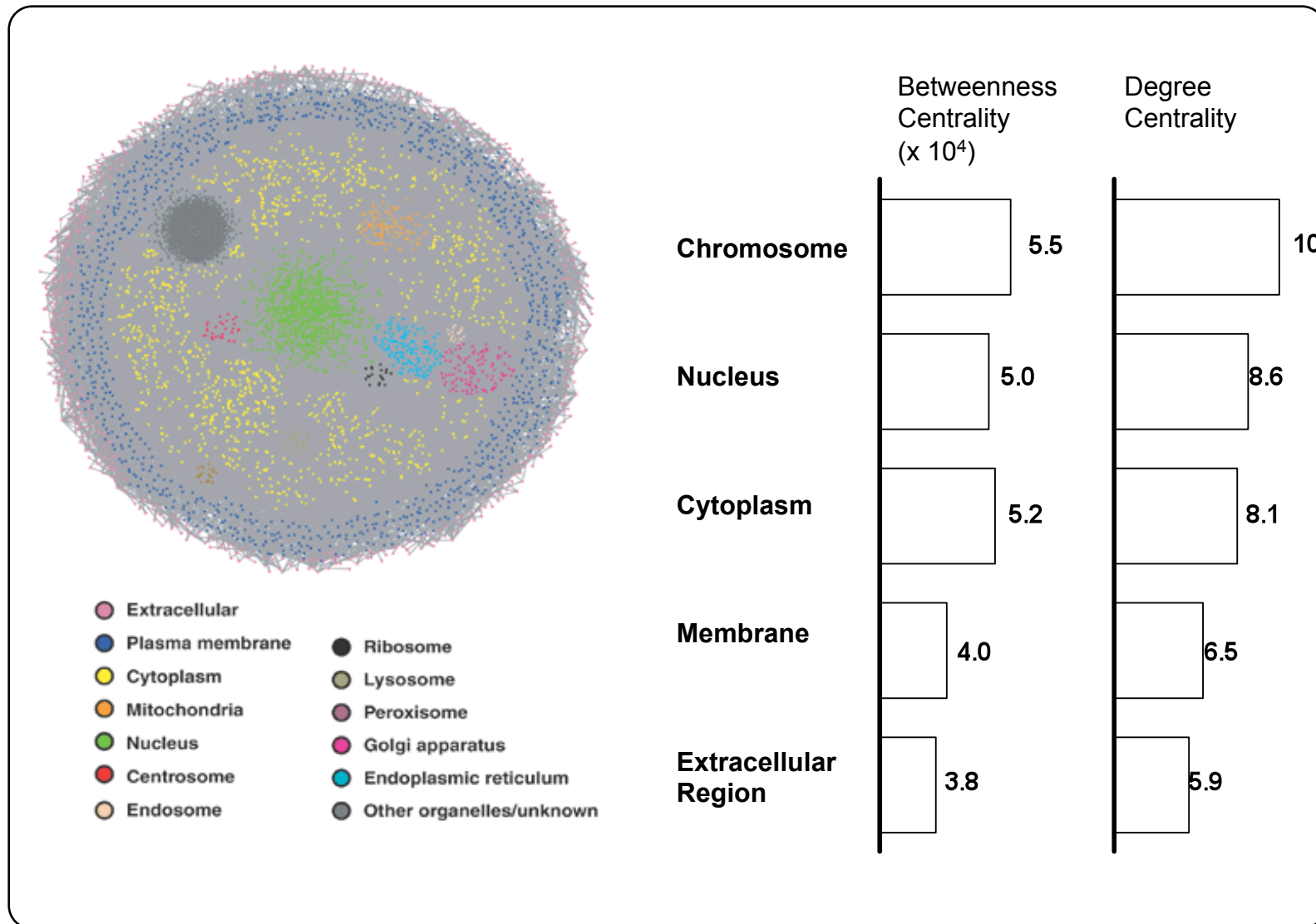*Specifically, a number of the SDs are likely not fixed, but rather common CNVs in the reference genome

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<<0.01

**Exposed sites**

**Buried sites**

Average Relative Surface Exposure

2.66

2.26

p<<0.01

**Site with Synonymous Mutations only**

**Sites with Non-synonymous Mutations**

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



| | Betweenness Centrality ($\times 10^4$) | 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

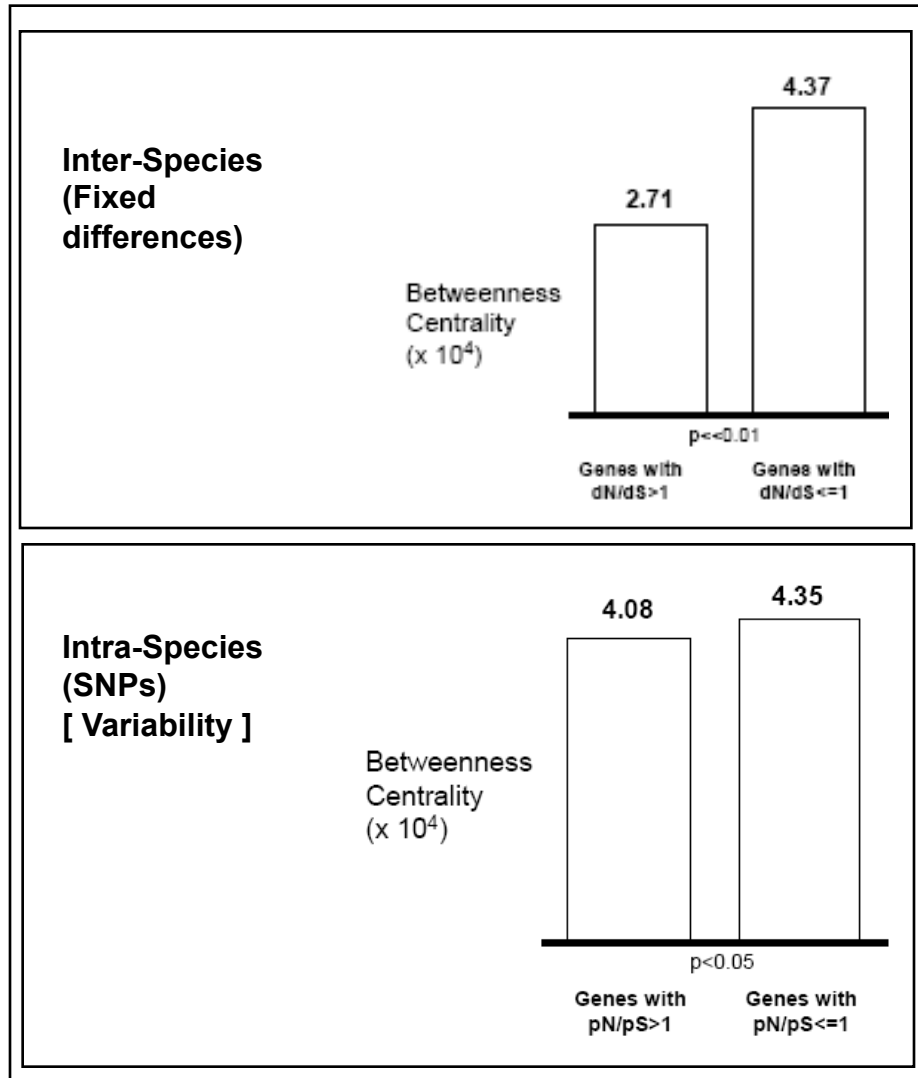# IS RELAXED CONSTRAINT OR ADAPTIVE EVOLUTION THE REASON FOR THE PREVALENCE OF BOTH SELECTED GENES AND SDs AT THE NETWORK PERIPHERY?

| | 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 ($\times 10^4$)

p<<0.01

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

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

4.08    4.35

Betweenness Centrality ($\times 10^4$)

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*

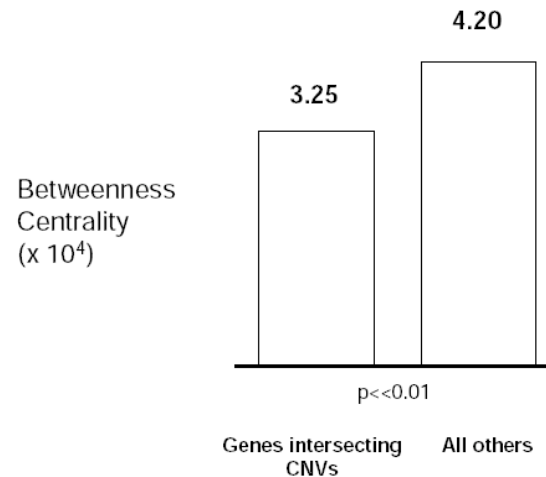*But it's hard to quantify

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

# Conclusions:
# Net Intro. + Predicting Networks



- Developing Standardized Descriptions of Protein Function
  ◊ Gene Naming

- Predicting Networks
  ◊ Extrapolating from the Training Set
  ◊ Principled ways of using the training set data in the fullest possible fashion
    - Prediction Propagation
    - Kernel Initialization

# Conclusions: Network Dynamics across Cellular States

- Merge expression data with Networks

- Active network markedly different in different conditions

- Identify transient hubs associated with particular conditions

- Use these to annotate genes of unknown function

# 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 & Human 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,,,

# TopNet – an automated web tool (vers. 2 : "TopNet-like Yale Network Analyzer")
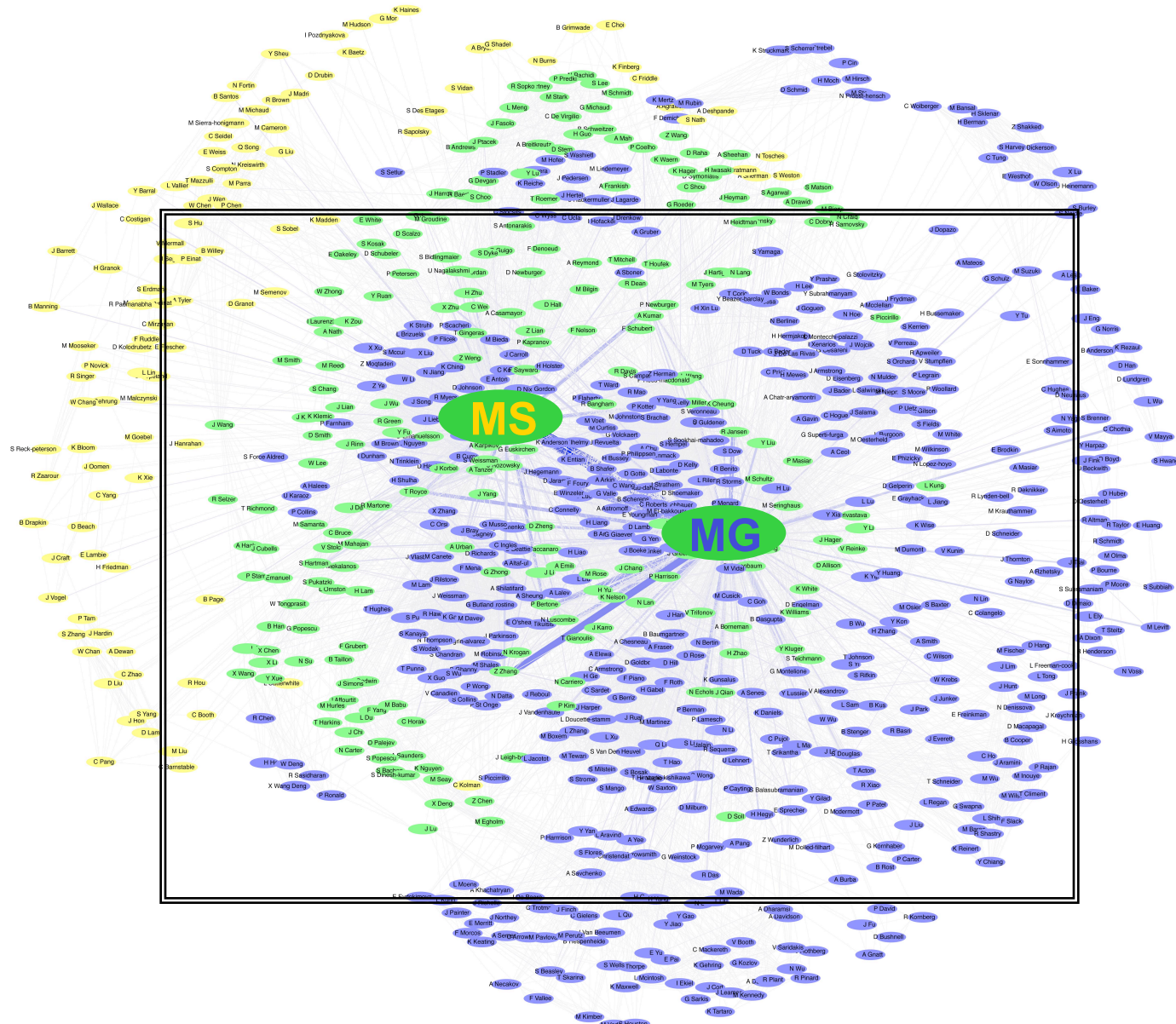


Normal website + Downloaded code (JAVA)
+ Web service (SOAP) with Cytoscape plugin

[Yu et al., NAR (2004); Yip et al. Bioinfo. (2006);
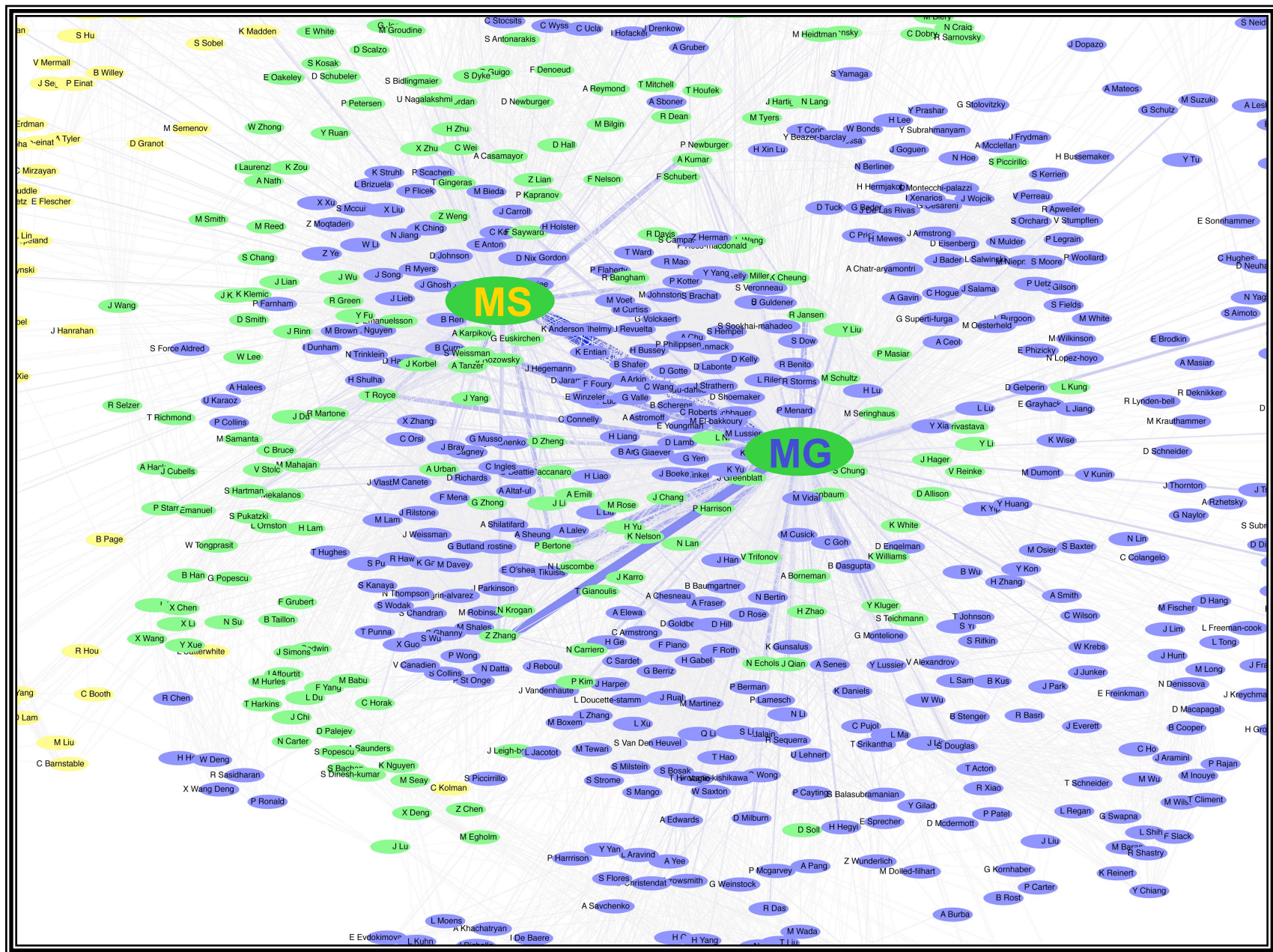Similar tools include Cytoscape.org, Idekar, Sander et al]
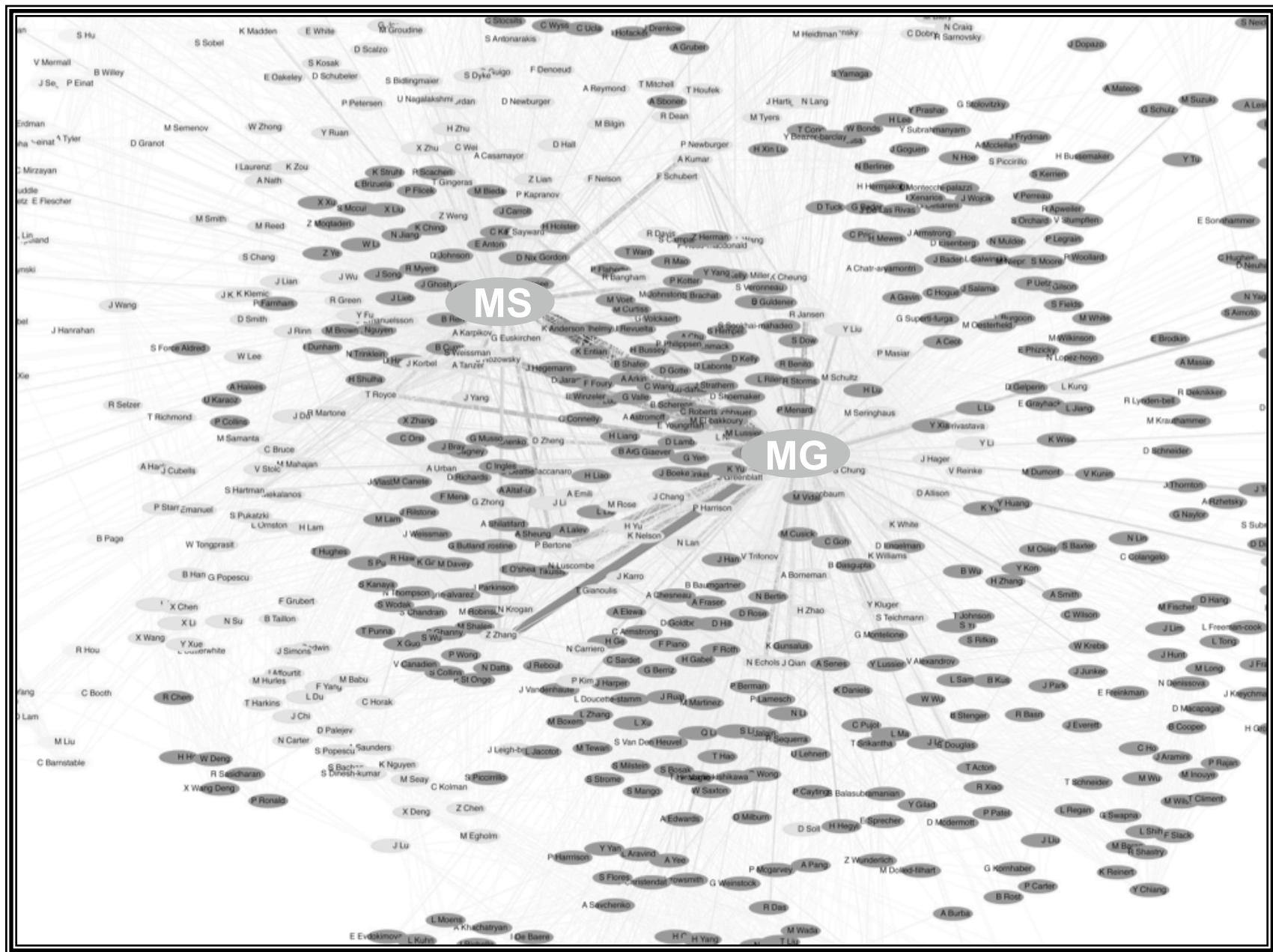
# Acknowledgements

## TopNet.GersteinLab.org
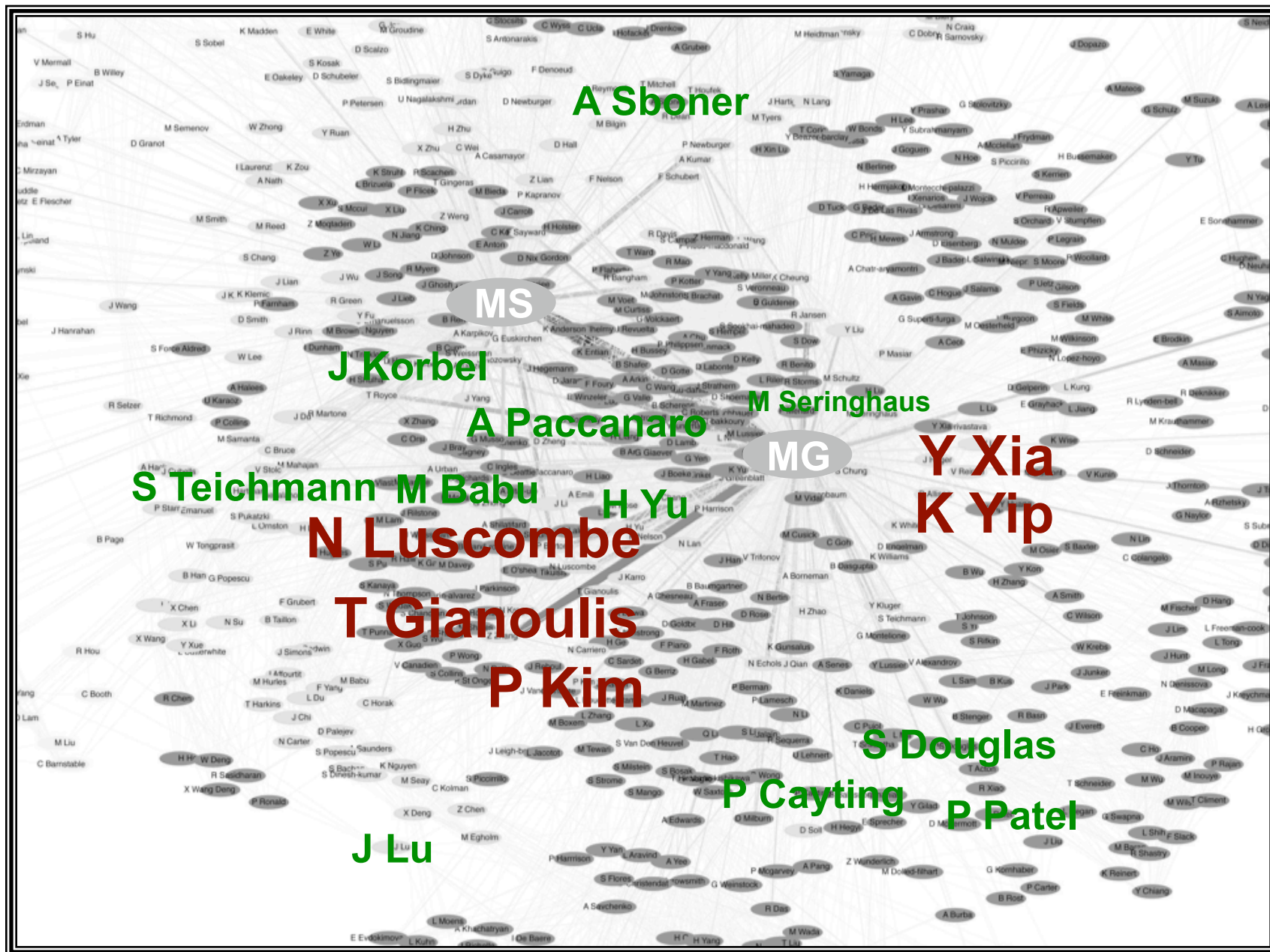
# Acknowledgements

## TopNet.GersteinLab.org

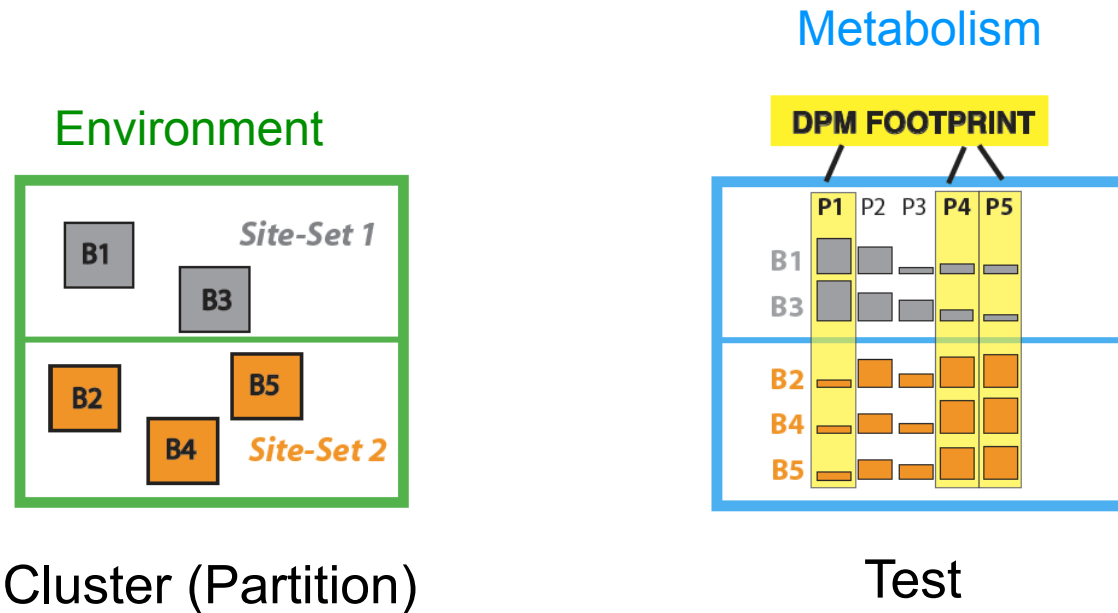# Acknowledgements

## TopNet.GersteinLab.org

**P Bork, J Raes**

**Job opportunities currently for postdocs & students**



97 Gerstein.Info Talks (c) 2006

# Extra

# DPM: Discriminative Partition Matching
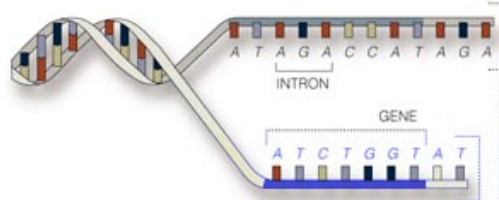
Metabolism

Environment
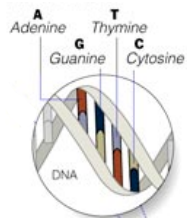


Cluster (Partition)

Test

Taurine biosynthesis
Heme biosynthesis
Asparagine degradation
Nitrogen fixation
Acylglycerol degradation
Asparagine biosynthesis
Cysteine Metabolism

| Functional class | pval |
|---|---|
| InfoStorage & Processing | .07 |
| Cellular Process | .08 |
| Metabolism | $4 \times 10^{-14}$ |

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

# Interactome

**GENOME**

protein-DNA interactions

**PROTEOME**

**protein-protein interactions**

**METABOLISM**

**Protein-small molecule interactions**

[From H Yu]

Citrate Cycle

100 Gerstein.info/talks (c) 2008

**Networks help us understand biological processes**

[From H Yu]

# More Information on this Talk

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

SUBJECT: Networks

DESCRIPTION:
Cornell Medical School, Physiology, Biophysics and Systems Biology (PBSB) graduate program, 2009.01.26, 16:00-17:00; [I:CORNELL-PBSB] (Long networks talk, incl. the following topics:
why networks w. amsci*, funnygene*, net. prediction intro, memint*, tse*, essen*, sandy*, metagenomics*, netpossel*, tyna*+ topnet*, & pubnet* . Fits easily into 60' w. 10' questions. PPT works on mac & PC and has many photos w. EXIF tag kwcornellpbsb .)

(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 )