

# Bioinformatics I Biological Networks

Andreas Wagner

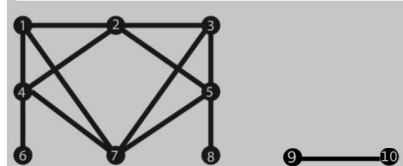
Institute of Evolutionary Biology  
and Environmental Studies, UZH  
andreas.wagner@ieu.uzh.ch

## Available on OLAT

Homework exercises for Bioinformatics I, Bio390  
Biological networks, Andreas Wagner

Note: These exercises are for you to solve on your own. You do not have to turn them in and they will not be graded. Even though solutions are provided at the end of this document, we highly recommend that you solve them and do so before looking at the solutions, because similar (not necessarily identical) problems will occur on the final exam.

Exercise 1: (Graph Representation)



# Further reading

## Complex networks in general

Newman, MEJ. The structure and function of complex networks. *SIAM Review* **45**, 167-256, 2003.

Fortunato, S., Hric, D. Community detection in networks: A user guide. *Physics Reports* **659**, 1-44, 2016.

## Protein interaction networks

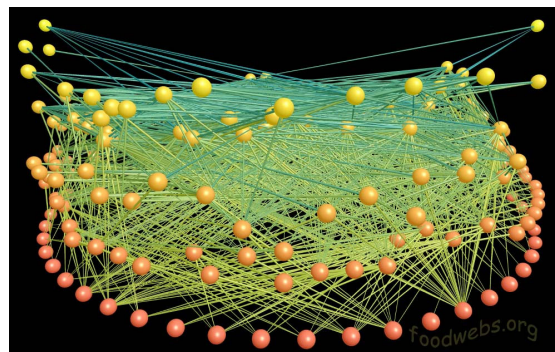
Xia et al. Analyzing cellular biochemistry in terms of molecular networks. *Annu. Rev. Biochem.* **73**:1051-87, 2004

Rajagopala et al., The binary protein-protein interaction landscape of *Escherichia coli*. *Nature Biotechnology* **32**, 285-290, 2014

## Metabolic networks

Price et al. *Nature Reviews Microbiology* **2**, 886-897, 2004

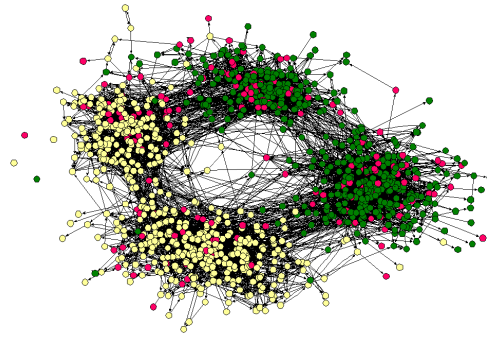
# Networks everywhere



## El Verde Rainforest trophic web, Puerto Rico

Dunne, J.A., R.J. Williams, and N.D. Martinez. 2002. Food-web structure and network theory: The role of connectance and size. *PNAS*, vol. 99, no. 20, pp. 12917-12922.

## Networks everywhere



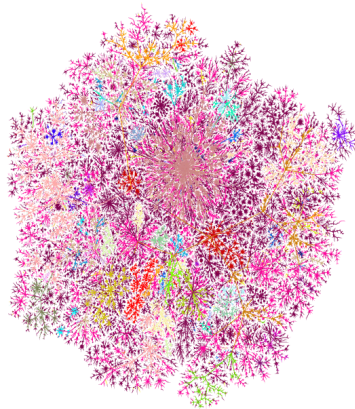
**Middle and High school friendship network in a US school**

Yellow - White Race; Green - Black Race; Pink - Other

The split from the lower left to the upper right is according to age (middle/high school)

James Moody, Race, school integration, and friendship segregation in America, *American Journal of Sociology* **107**, 679-716 (2001).

## Networks everywhere



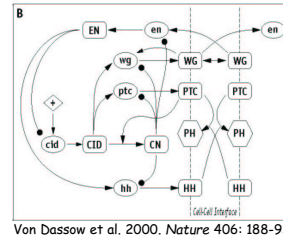
**Internet IP addresses, colored by ISP**

Bill Cheswick (<http://www.cheswick.com/ches/>)

## Cell-biological networks

### 1. Small networks dedicated to a specific task (up to dozens of gene products)

Chemotaxis  
Cell-cycle regulation  
Fruit fly segmentation  
Flower development  
...



Mathematical characterization based on detailed,  
quantitative biochemical information

## Cell-biological networks

### 2. Genome-scale networks (hundreds to thousands of gene products)

**Protein interaction networks**  
**Metabolic networks**  
Transcriptional regulation networks  
Genetic interaction networks  
...



Mathematical characterization based on qualitative  
understanding of network topology

## Protein interaction networks



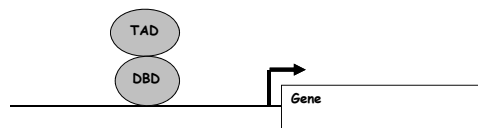
## The yeast two-hybrid assay

A technique to identify interacting proteins

Relies on the modularity of eukaryotic transcriptional regulators

DBD: DNA binding domain

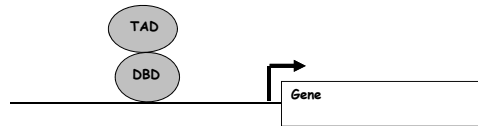
TAD: transcriptional activation domain



# The yeast two -hybrid assay

Carried out in cells of the yeast *Saccharomyces cerevisiae*

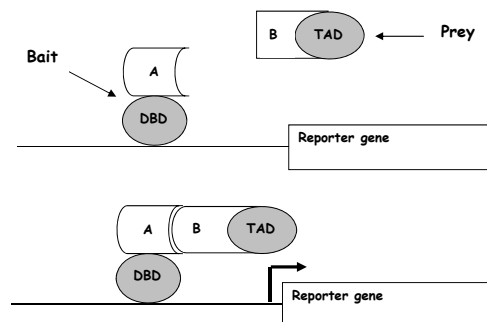
Can be applied to any two proteins (not just yeast proteins)

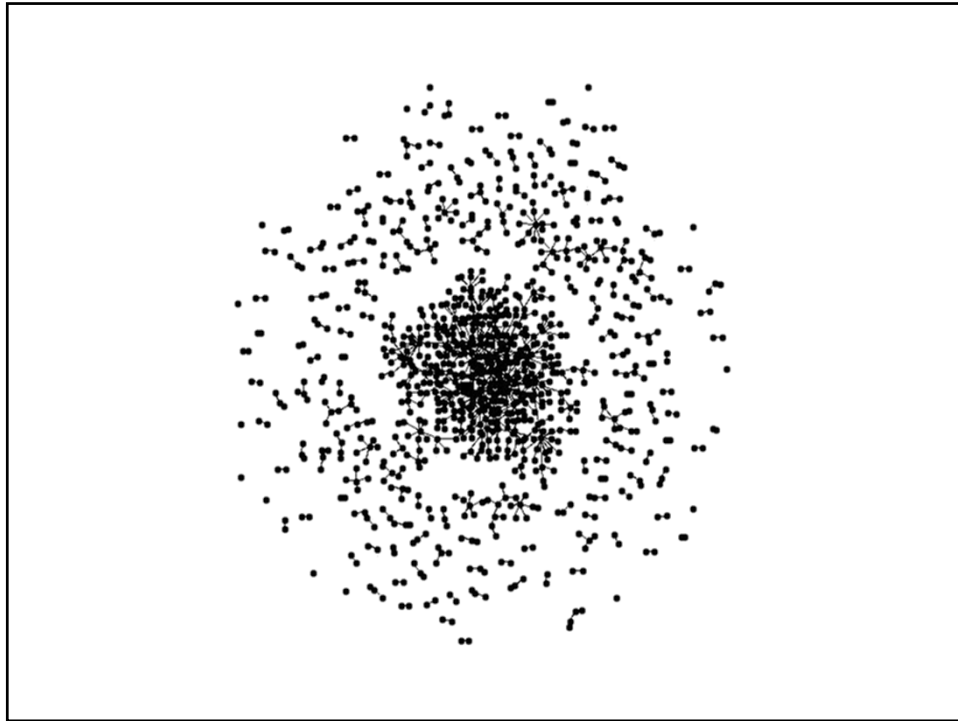


# The yeast two -hybrid assay

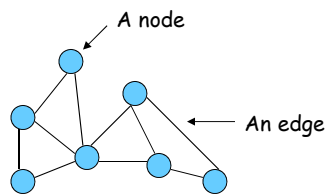
A,B: two proteins whose interaction is to be assayed

Reporter gene: a gene whose activity is easily monitored





## Graphs



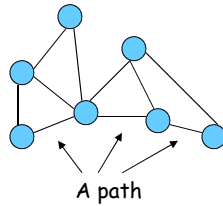
A graph  $G=(V,E)$  comprises  
a set  $V$  of nodes (vertices)  
a set  $E$  of edges

$$V = \{V_1, \dots, V_n\}$$

$$E = \{(V_i, V_j), \dots, (V_k, V_l)\}$$

Protein interaction networks are undirected graphs  
(Individual node pairs in  $E$  are unordered.)

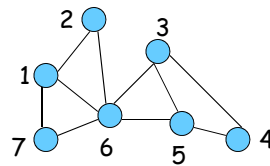
## Graphs



A path is a sequence of alternating nodes and edges in which no node is visited more than once

A geodesic is the shortest path between two nodes.

## Graphs can be represented by matrices

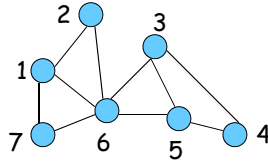


Adjacency matrix  $A=(a_{ij})$

$$\begin{aligned} a_{ij} &= 1 & (V_i, V_j) \in E \\ a_{ij} &= 0 & \text{otherwise} \end{aligned}$$

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 1 & 0 & 1 & 0 & 1 \\ 1 & 0 & 0 & 0 & 0 & 1 & 0 \end{pmatrix}$$



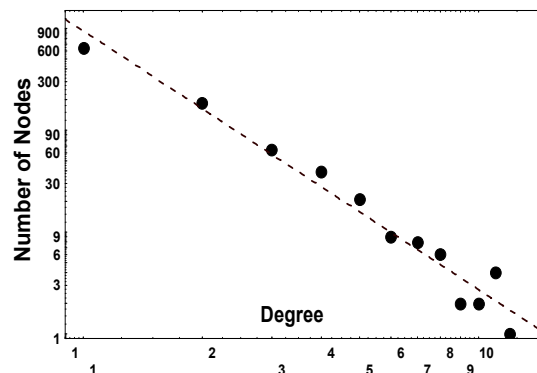


The degree (connectivity)  $k_i$  of a node  $V_i$  is the number of edges incident with the node (e.g.,  $k_1=3$ ,  $k_6=5$ ).

$$k_i = \sum_j a_{ij}$$

Graphs can be characterized according to their degree distribution  $P(k)$ , the fraction of nodes having degree  $k$ .

Protein interaction networks (and many other networks) have broad-tailed degree distributions.



Wagner A, Proc. Roy. Soc. London 2003

## The best-studied mathematical models of graphs

### k-regular graphs

N nodes,  $K=kN$  edges  
every node has degree k

### Erdős-Rényi random graphs

N nodes, K edges

edges connect pairs of randomly chosen nodes  
(avoiding multiple edges)

Degree distribution is Poisson

$$P(k) = \exp(-\langle k \rangle) \frac{\langle k \rangle^k}{k!}$$

**Biological networks are more complex and heterogeneous than predicted by these models**

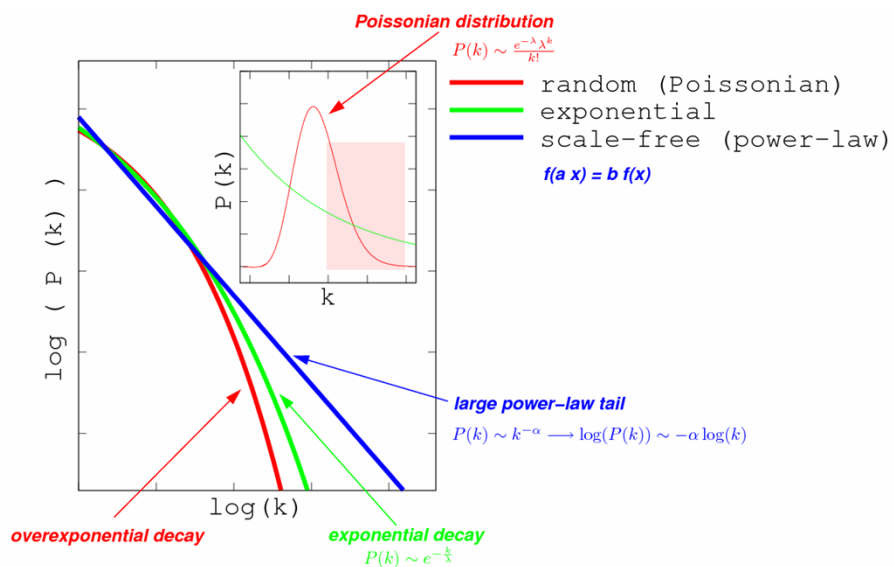
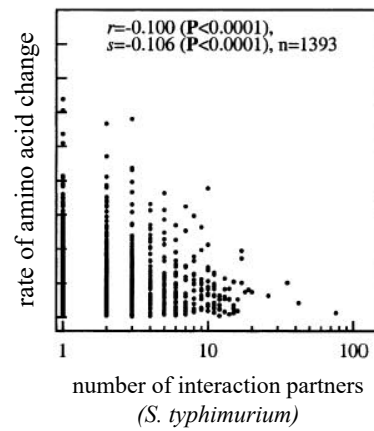


Figure courtesy of A. Caflisch

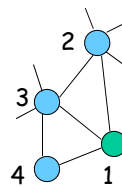
**Highly connected proteins tolerate fewer amino acid substitutions in their evolution**



Hahn et al. Journal of Molecular Evolution 2004

**The degrees of nodes in a graph may be correlated**

Average nearest neighbor degree of a node



$$k_1 = 3$$

$$k_2 = 5$$

$$k_3 = 5$$

$$k_4 = 2$$

$$k_{nn,i} = (1/3)(5 + 5 + 2) = 4$$

$$k_{nn,i} = \frac{1}{k_i} \sum_{j, \text{ nearest neighbors of } i} k_j$$

## The degrees of nodes in a graph may be correlated

Average nearest neighbor degree of all nodes with degree  $k$

$$k_{nn,i} = \frac{1}{k_i} \sum_{j, \text{ nearest neighbors of } i} k_j$$

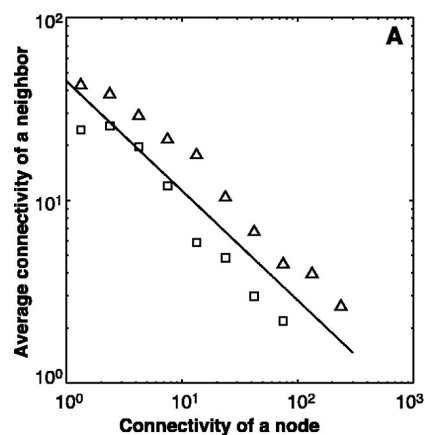
$N_k$ ...number of nodes with degree  $k$

$$k_{nn}(k) = \frac{1}{N_k} \left( \sum_{\text{nodes with degree } k} k_{nn,i} \right)$$

A graph is assortative if  $k_{nn}(k)$  increases with  $k$   
nodes connect to nodes of similar connectivity

A graph is disassortative if  $k_{nn}(k)$  decreases with  $k$

## Protein interaction networks are disassortative



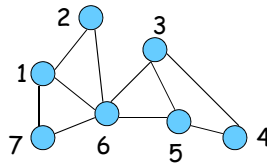
Few interactions between hubs

Many interactions between hubs and neighbors with low degree

Plot of  $P_{nn}(k)$  against  $k$  for the yeast protein interaction network (triangles) and the transcriptional regulation network (squares)

Maslov and Sneppen, Science 2002

**Path length and diameter are measures of graph compactness**



**Matrix of shortest paths  $D=(d_{ij})$**

$$D = \begin{pmatrix} 0 & 1 & 2 & 3 & 2 & 1 & 1 \\ 1 & 0 & 2 & 3 & 2 & 1 & 2 \\ 2 & 2 & 0 & 1 & 1 & 1 & 2 \\ 3 & 3 & 1 & 0 & 1 & 2 & 3 \\ 2 & 2 & 1 & 1 & 0 & 1 & 2 \\ 1 & 1 & 1 & 2 & 1 & 0 & 1 \\ 1 & 2 & 2 & 3 & 2 & 1 & 0 \end{pmatrix}$$

**Connected graph:  $d_{ij} < \infty$  for all  $i, j$**

**Path length and diameter are measures of graph compactness**

Diameter of a graph:  $\max_{i,j} d_{ij}$

Mean (arithmetic) shortest path length  
or characteristic path length

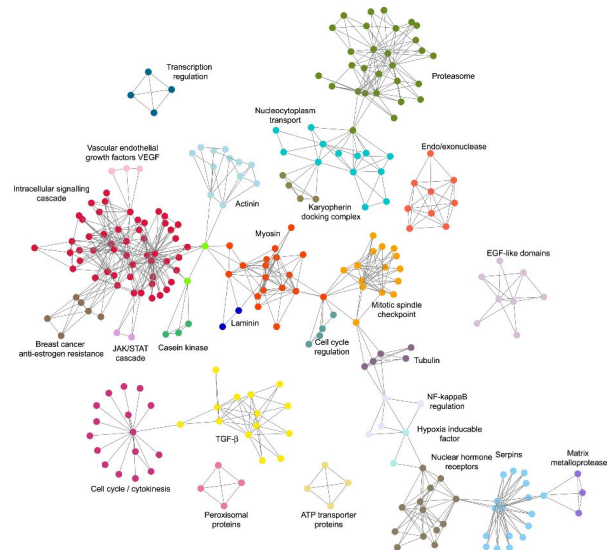
$$L = \frac{1}{N(N-1)} \sum_{i,j, i \neq j} d_{ij}$$

Mean (harmonic) shortest path length  
or "efficiency" of a graph

$$L = \frac{1}{N(N-1)} \sum_{i,j, i \neq j} \frac{1}{d_{ij}}$$

(Better suited than characteristic path length for disconnected graphs)

## Many graphs can be subdivided into “communities”

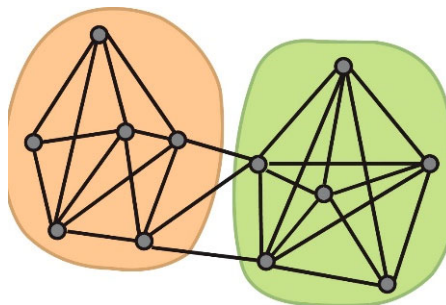


Community structure of a rat protein interaction network

Fortunato and Hric, Physics Reports 659, 1-44, 2016

## Many graphs can be subdivided into “communities”

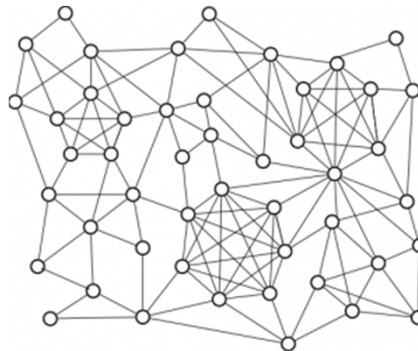
In a community (cluster, module) nodes can be subdivided into groups that share more edges with each other than with nodes outside the community



Fortunato and Hric, Physics Reports 659, 1-44, 2016

## The most densely connected communities are cliques

**clique:** a largest complete (=fully connected) subgraph



A graph with multiple cliques

<http://skipperkongen.dk/2010/11/>

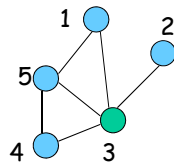
## The clustering coefficient is a measure of edge density

**Clustering coefficient  $c_i$  of a node  $i$ .**

The fraction of a node's neighbors that are neighbors of each other

$$c_i = \frac{E_i}{\frac{k_i(k_i-1)}{2}}$$

$E_i$  ... number of edges among neighbors of  $i$   
 $k_i$  ... degree of  $i$



$$c_3 = \frac{2}{\frac{4(3)}{2}} = \frac{1}{3}$$

**Clustering coefficient  $c$  of a graph**

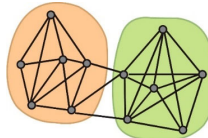
The average of the clustering coefficients of all nodes

(In a clique, all nodes have  $c_i=1$ , so  $c=1$  for a graph that is a clique.)

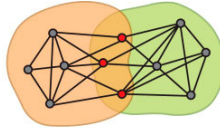
## Many computational methods aim to detect communities in networks

Some require information about the total number of communities (easier), others don't (more difficult).

**Hard-clustering** methods generate non-overlapping communities (easier)



**Soft-clustering** methods allow overlapping communities (more difficult)



Optimization methods for community detection aim to maximize a quantity that indicates to what extent a network clusters into different communities

A very popular such quantity

**Modularity Q** for a network that is subdivided into  $n$  modules

$$Q = \sum_{i=1 \dots n} (e_{ii} - a_i^2)$$

$e_{ij}$ ...fraction of edges that connect nodes in module  $i$  and module  $j$

$e_{ii}$ ...fraction of edges that connect nodes within module  $i$ .

$$a_i = \sum_{j=1 \dots n} e_{ij} \text{ fraction of edges that begin or end in module } i$$



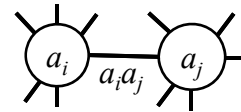
**Optimization methods for community detection aim to maximize a quantity that indicates to what extent a network clusters into different communities**

A very popular such measure

**Modularity Q** 
$$Q = \sum_{i=1 \dots n} (e_{ii} - a_i^2)$$

If you have subdivided a graph into  $n$  putative modules, but these modules do not reflect the graph's actual structure, then the fraction of edges that connect two such "spurious" modules  $i$  and  $j$  is given by the product rule of probabilities as  $e_{ij} = a_i a_j$ .

A special case is  $e_{ii} = a_i^2$



Thus, if a graph does not have a modular structure, then  $Q \approx 0$ .

**Optimization methods aim to maximize a quantity that indicates to what extent the network clusters into different communities**

A very popular such measure

**Modularity Q** 
$$Q = \sum_{i=1 \dots n} (e_{ii} - a_i^2)$$

Q is larger for graphs and communities where pairs of connected nodes tend to reside in the same module.

$Q \approx 1$  for graphs with the most pronounced modular structure

This occurs if all values of  $e_{ii}$  are large, i.e., almost all edges connect nodes within the same module, while the  $a_i^2$  is small, that is, by chance alone one would expect that very few edges connect nodes within the module

## The Girvan-Newman algorithm is a popular heuristic to cluster large graphs

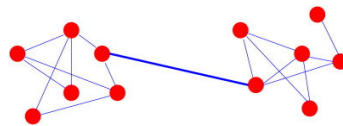
It does not guarantee to find the best possible clustering

It relies on the concept of edge betweenness

Edge betweenness (centrality, load):

the number of shortest paths passing through an edge  $i$

$$b_i = \sum_{j,k, j \neq k} \frac{n_{jk}(i)}{n_{jk}}$$



$n_{jk}(i)$  number of shortest paths connecting node  $j$  and  $k$  and passing through edge  $i$   
 $n_{jk}$  number of shortest paths connecting node  $j$  and  $k$

## The Girvan-Newman algorithm is a popular heuristic to cluster large graphs

It is an iterative divisive clustering algorithm

Idea: Edges between modules are those with the highest edge betweenness

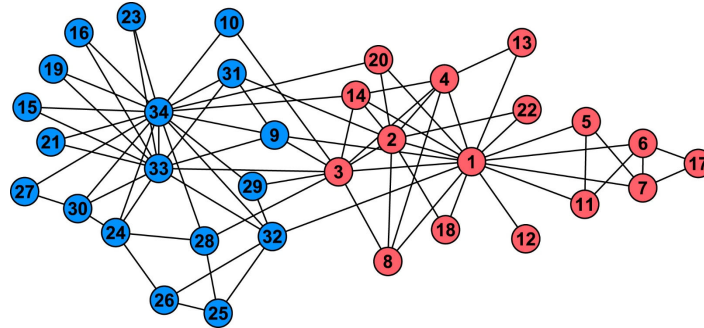
Remove those edges and you get good module separation

Procedure

1. Remove the edge with the highest betweenness
2. Recalculate edge betweenness for the now-reduced graph
- (3. Determine modularity  $Q$ )
4. Back to one until all nodes are isolated

The optimal partition is that with the highest  $Q$

The “Karate club” network of Zachary has served as a benchmark for many community detection algorithms



Fortunato and Hric, Physics Reports 659, 1-44, 2016

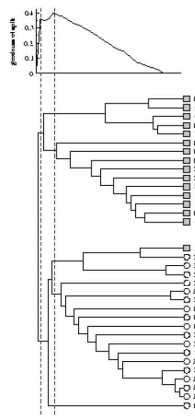
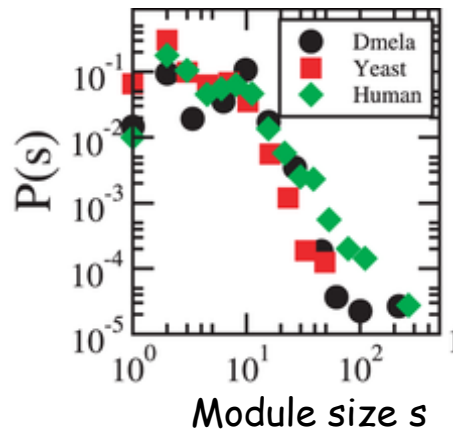


Figure 7.5: Modularity  $Q$  (top panel) and dendrogram (bottom panel) obtained by the application of the GN algorithm to the karate club network. The modularity has two maxima corresponding respectively to *i*) a split into two communities, which matches closely the real split of the club (only node 3 is incorrectly classified), and *ii*) a split into five communities. Figure taken from Ref. [51].

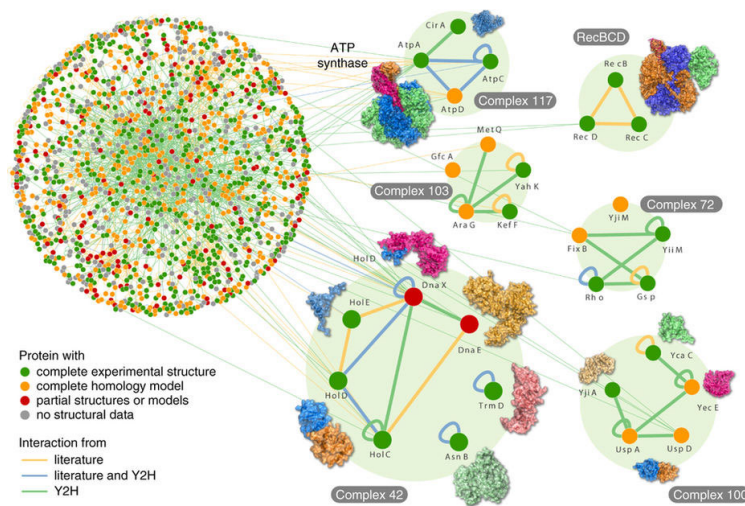
Boccaletti et al. 2006

## Module sizes in protein interaction networks have a broad-tailed distribution



Lancichinetti et al. (2010) Characterizing the Community Structure of Complex Networks. PLOS ONE 5(8): e11976.

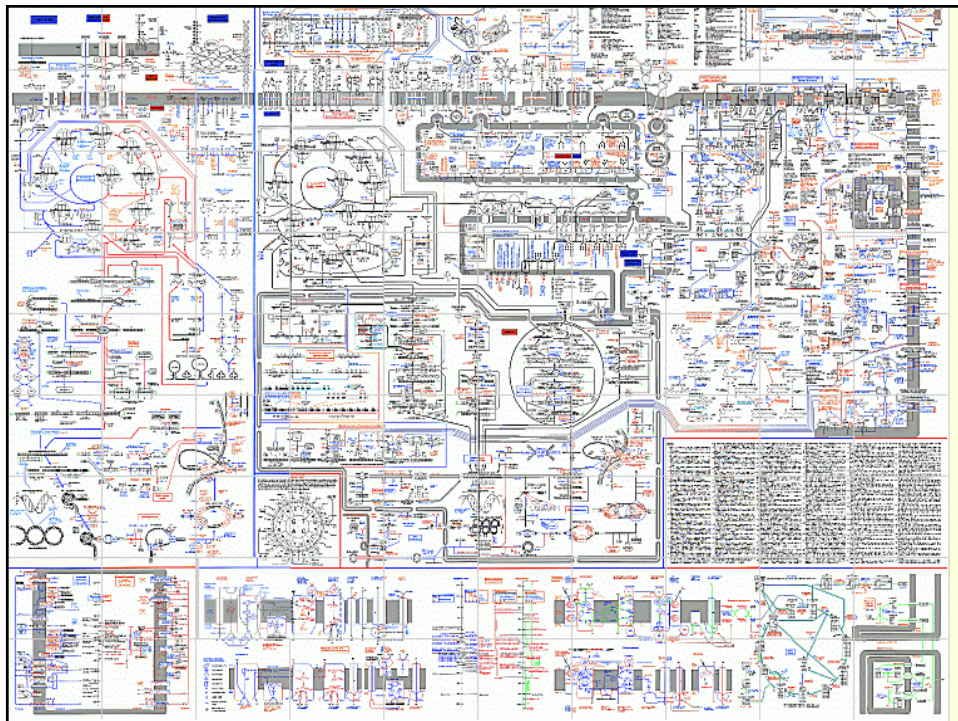
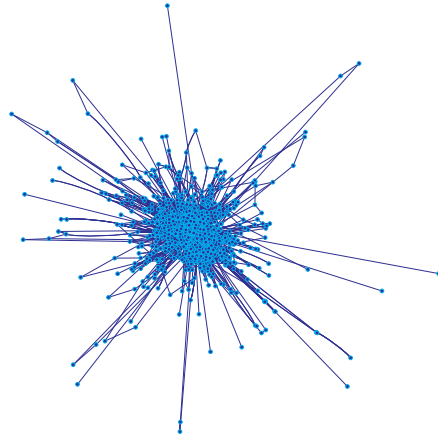
## The best maps of protein interaction networks integrate different kinds of information



An E.coli protein interaction network

Rajagopala et al., Nature Biotechnology 2014

# Metabolic networks



**A metabolic network is a set of chemical reactions that produces**

*energy*

(for maintenance of cell functions and for biosyntheses)

*molecular building blocks for biosyntheses*

**These reactions are catalyzed by enzymes that are encoded by genes.**

**In free-living heterotrophic organisms, several hundred such enzymatic reactions are necessary to fulfill these functions.**

**Graphs can (crudely) represent large chemical reaction networks**

#### Stoichiometric Equations

1 Glucose 6-phosphate (G6P) + 1 NADP<sup>+</sup>

1 6-Phosphoglucono δ-lactone + 1 H<sub>2</sub>O

1 6-Phosphogluconate + 1 NADP<sup>+</sup>

1 Ribulose 5-phosphate

$\xrightarrow{zwf}$

$\xrightarrow{pgl}$

$\xrightarrow{gnd}$

$\xrightarrow{rpe}$

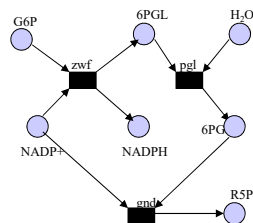
1 6-Phosphoglucono δ-lactone (6PGL) + 1 NADPH

1 6-Phosphogluconate (6PG)

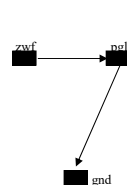
1 Ribulose 5-phosphate (R5P) + 1 NADPH

1 Xylulose 5-phosphate (X5P)

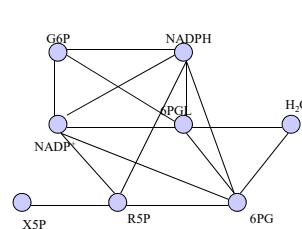
**Bipartite graph**



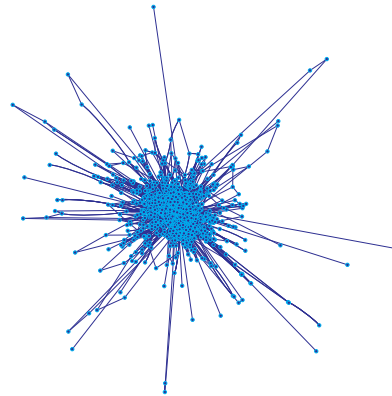
**Enzyme graph**



**Substrate graph**

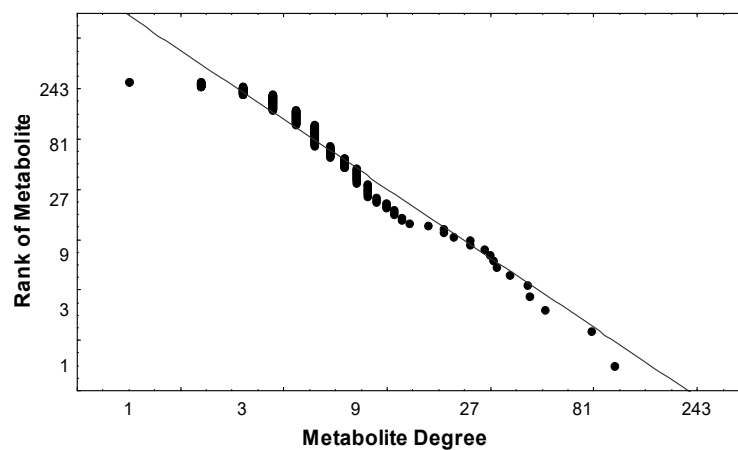


**An enzyme graph representation of the metabolic network  
of the yeast *Saccharomyces cerevisiae***



 Pajek

**Metabolic networks have a broad-tailed degree distribution**

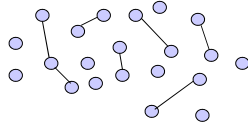


Substrate network of *E. coli*

Wagner and Fell, Proc. Roy. Soc. London B 2001

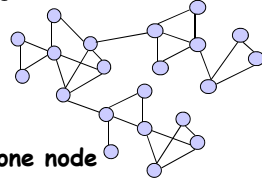
## Key features of small-world graphs

1. They are sparse



2. They are "cliquish"

as measured by a high clustering coefficient



3. Despite 1 and 2, paths from any one node to any other node are VERY short

(short mean path length, "small-worldness")

Watts and Strogatz, Nature 1998

## The *E. coli* core metabolism is a small-world network

It is sparse

It is highly clustered

It has short characteristic path length



## Many graphs have “small-world” features

<i>Graph</i>	<i>Nodes</i>	<i>Edges</i>
<b>Computer networks</b>	Computers	Data transmission lines
<b>Friendship networks</b>	People	Being acquainted
<b>The world wide web</b>	Web pages	Hyperlinks
<b>Actor collaboration graph</b>	Actors	Having acted in the same movie
<b>Power grids</b>	Transformers	Power lines
<b>Citation network</b>	Publication	Citation
<b>Nematode CNS</b>	Nerve cells	Axons

## Why are metabolic networks small-world networks?

Signals propagate VERY rapidly in small world networks.

Perhaps compact network structure allows the cell to adapt rapidly to changing conditions.

Studying only the structure of metabolic networks neglects their function

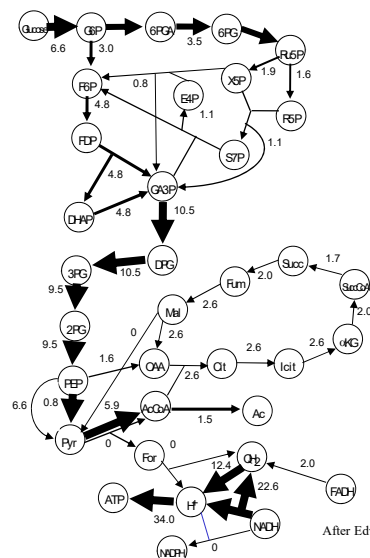
One needs to analyze the flow (flux) of matter through these networks

For optimal cell growth, metabolic networks need to produce biochemical precursors in well-balanced amounts.

This necessitates a specific distribution of metabolic fluxes through enzymatic reactions in the network.

(Metabolic flux: the rate at which an enzyme converts substrate into product per unit time.)

Metabolic flux through central carbon metabolism of *E.coli* growing at a maximally possible rate in a glucose-minimal medium



After Edwards JS, Palsson BO. 2000. *PNAS* 97: 5528-33

**Flux balance analysis** requires a list of chemical reactions known to be catalyzed by enzymes in a given organism.

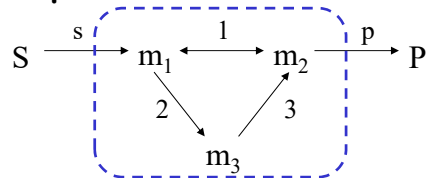
(For example, in yeast  
 >1100 reactions,  
 >500 metabolites,  
 >100 nutrients or waste products.)

**Flux balance analysis** has two tasks

Identify **allowable** metabolic fluxes through a metabolic network (fluxes that do not violate the law of mass conservation)

Within the set of allowable fluxes, identify fluxes that are associated with desirable properties (e.g., maximal rate of biomass production, maximal biomass yield per unit of carbon source.)

### A simple chemical reaction network



Metabolite concentrations  $m_i$  change according to the equations

$$\frac{dm_1}{dt} = v_s - v_1 - v_2$$

$$\frac{dm_2}{dt} = v_1 + v_3 - v_p$$

$$\frac{dm_3}{dt} = v_2 - v_3$$

$$\frac{d\vec{m}}{dt} = \mathbf{S}\vec{v}$$

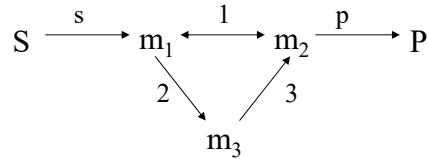
$$\mathbf{S} = \begin{pmatrix} 1 & -1 & -1 & 0 & 0 \\ 0 & 1 & 0 & 1 & -1 \\ 0 & 0 & 1 & -1 & 0 \end{pmatrix}$$

**Stoichiometry matrix**

$v_i$  metabolic flux through reaction  $i$

$$\vec{v} = (v_s, v_1, v_2, v_3, v_p)^T$$

Rows: metabolites  
 Columns: reactions



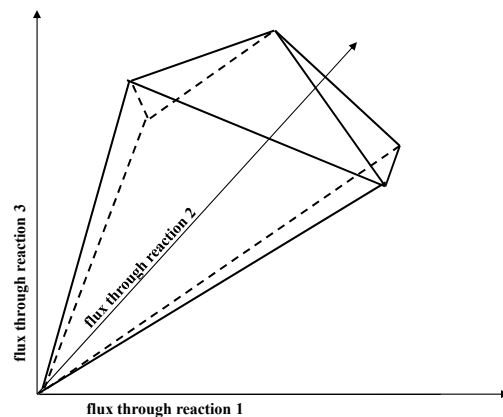
FBA assumes that metabolism is in a steady state where the concentrations of metabolites no longer change

$$\frac{d\vec{m}}{dt} = 0$$

$$S\vec{v} = 0$$

The solutions of these equations are the allowable metabolic fluxes. They form the so-called null space of  $S$

The null space of a metabolic network forms a high-dimensional "flux cone" (a convex polytope)



Several important properties of a metabolic network can be expressed as weighted sums of fluxes

$$Z(\vec{v}) = \sum_{i=1}^m c_i v_i$$

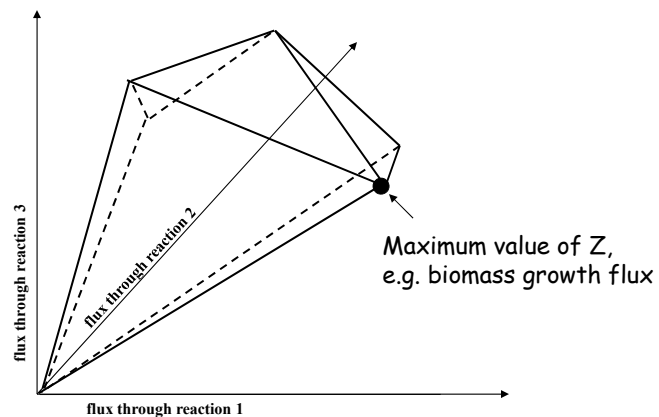
Example:

In the biomass growth flux,

$v_i$  is the rate at which essential biochemical precursor  $i$  is produced by a metabolic network.

$c_i$  is a constant that reflects the relative contribution of precursor  $i$  to biomass (can be estimated from the biomass composition of a cell.)

Linear programming can be used to determine regions within the flux cone where some linear function  $Z$  of the fluxes will be maximized.



## Example questions for flux balance analysis

Can a given organism (metabolism) survive in environment X?

How fast could it grow in this environment?

Why are many enzymatic reactions dispensable in any one environment?

Why do some metabolisms have many reactions, while others have few?

Does network function and flux influence network evolution

Is it possible to design "resistance-proof" antimetabolic drugs?

## Summary

Among the most prominent examples of genome-scale cell-biological networks are

protein interaction networks  
metabolic networks

Graph theory can be used to characterize these networks via

degree distribution and correlation  
characteristic path lengths and diameter  
clustering coefficient  
indicators of modularity  
...

# Summary

The biological significance of many aspects of network structure is still unclear

Analyses of network function need to go beyond graph theory

Flux balance analysis