

Bioinformatics 3

V 5 – Robustness and Modularity

Fri, Nov 7, 2014

Network Robustness

Network = set of connections

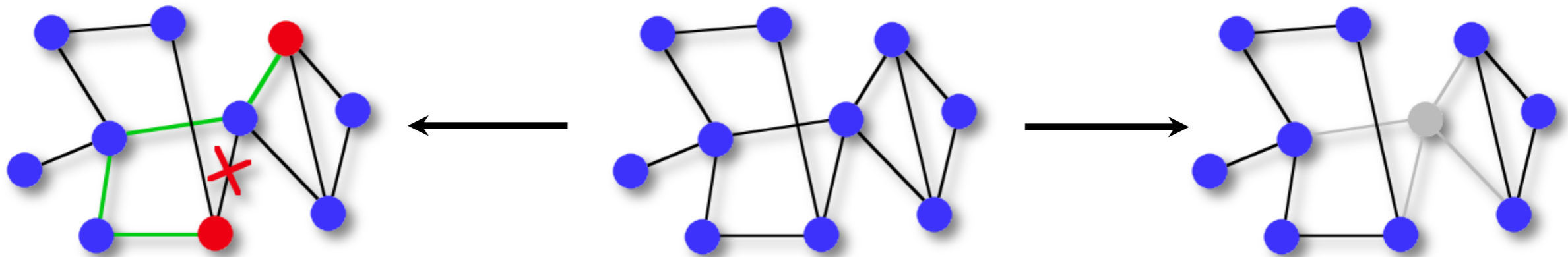
Failure events:

- loss of edges
- loss of nodes (together with their edges)

→ loss of connectivity

- paths become longer (detours required)
- connected components break apart

→ network characteristics change



→ **Robustness** = how much does the network (not) change when edges/nodes are removed

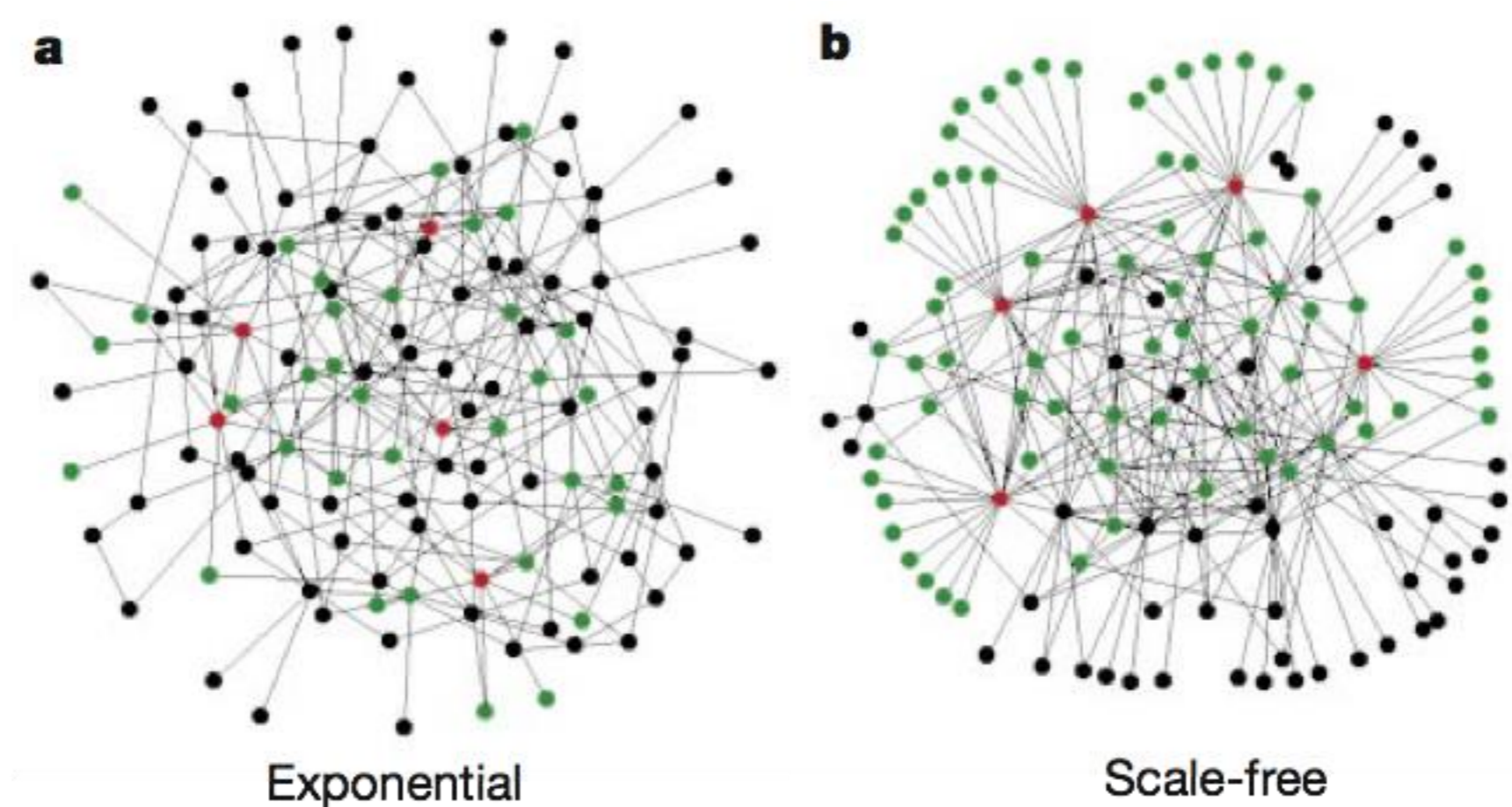
Error and attack tolerance of complex networks

Réka Albert, Hawoong Jeong & Albert-László Barabási

Department of Physics, 225 Nieuwland Science Hall, University of Notre Dame, Notre Dame, Indiana 46556, USA

Many complex systems display a surprising degree of tolerance against errors. For example, relatively simple organisms grow, persist and reproduce despite drastic pharmaceutical or environmental interventions, an error tolerance attributed to the robustness of the underlying metabolic network¹. Complex communication networks² display a surprising degree of robustness: although key components regularly malfunction, local failures rarely lead to the loss of the global information-carrying ability of the network. The stability of these and other complex systems is often attributed to the redundant wiring of the functional web defined by the systems' components. Here we demonstrate that error tolerance is not shared by all redundant systems: it is displayed only by a class of inhomogeneously wired networks,

Random vs. Scale-Free



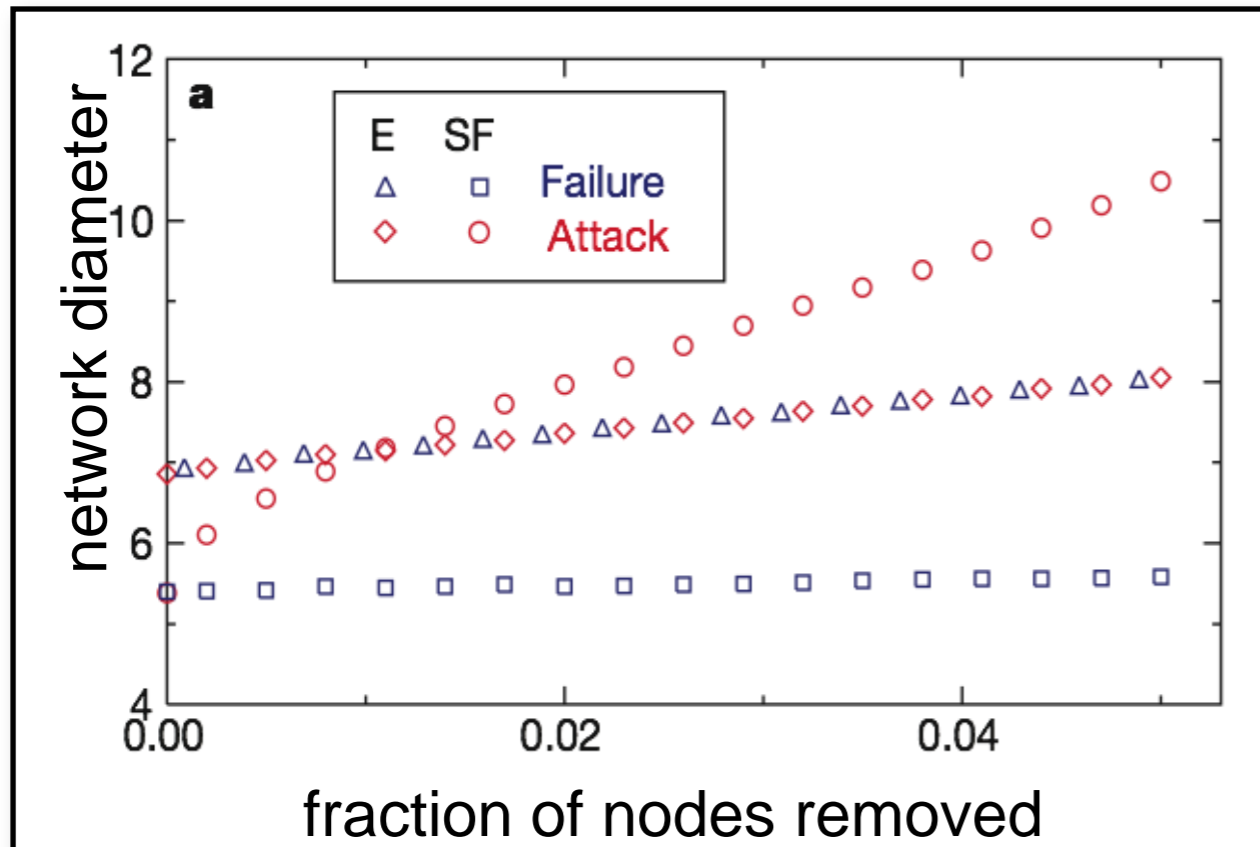
130 nodes, 215 edges

The **top 5** nodes with the highest k **connect** to...
... 27% of the network ... 60% of the network

Failure vs. Attack

Failure: remove randomly selected nodes

Attack: remove nodes with highest **degrees**



SF: scale-free network -> attack

E: exponential (random) network
-> failure / attack

SF: failure

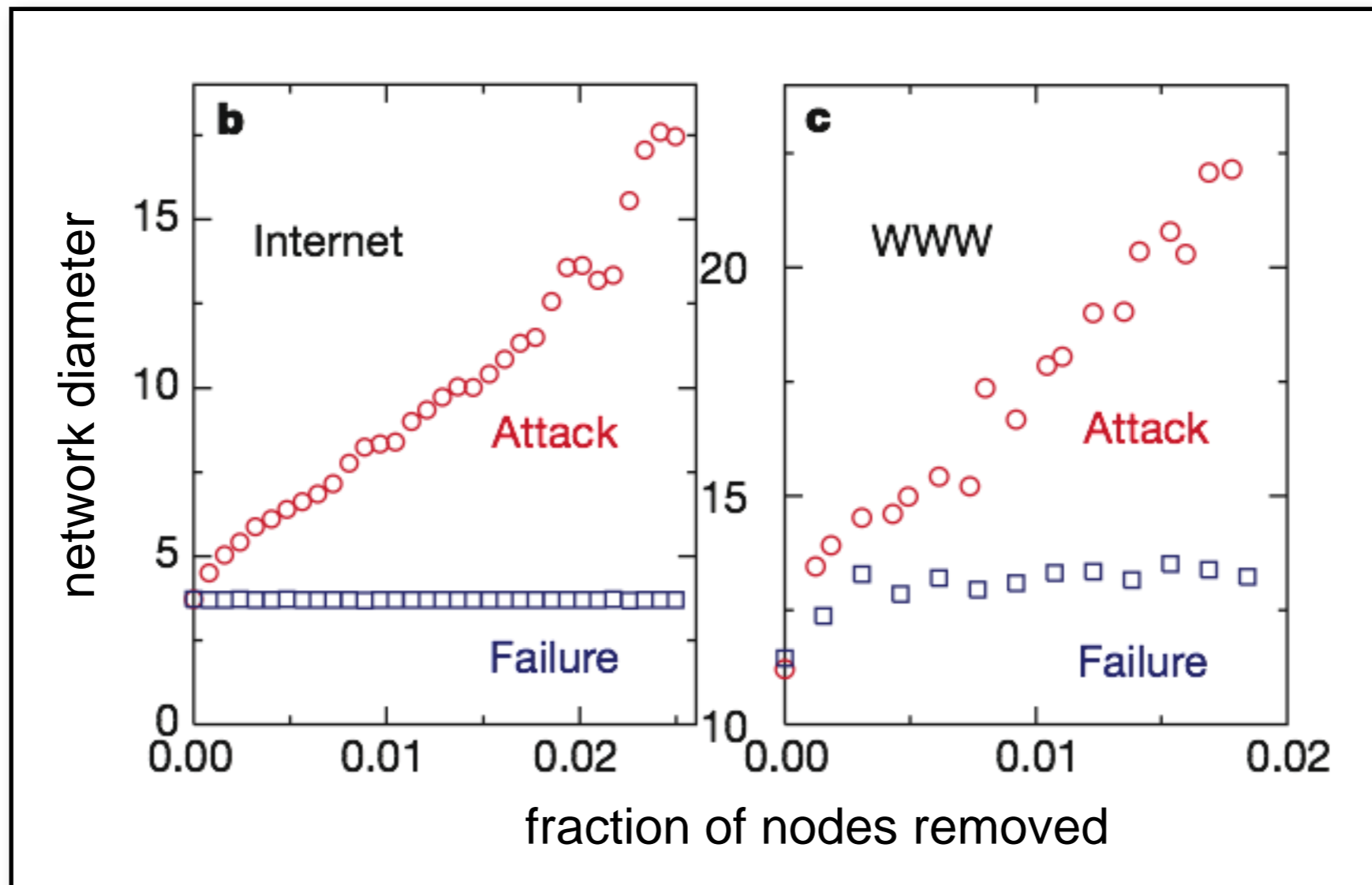
$N = 10000$, $L = 20000$, but effect is size-independent;

Interpretation:

SF network diameter increases strongly when network is attacked
but not when nodes fail randomly

Two real-world networks

- Scale-free:**
- very **stable** against random **failure** ("packet re-rooting")
 - very **vulnerable** against dedicated **attacks** ("9/11")

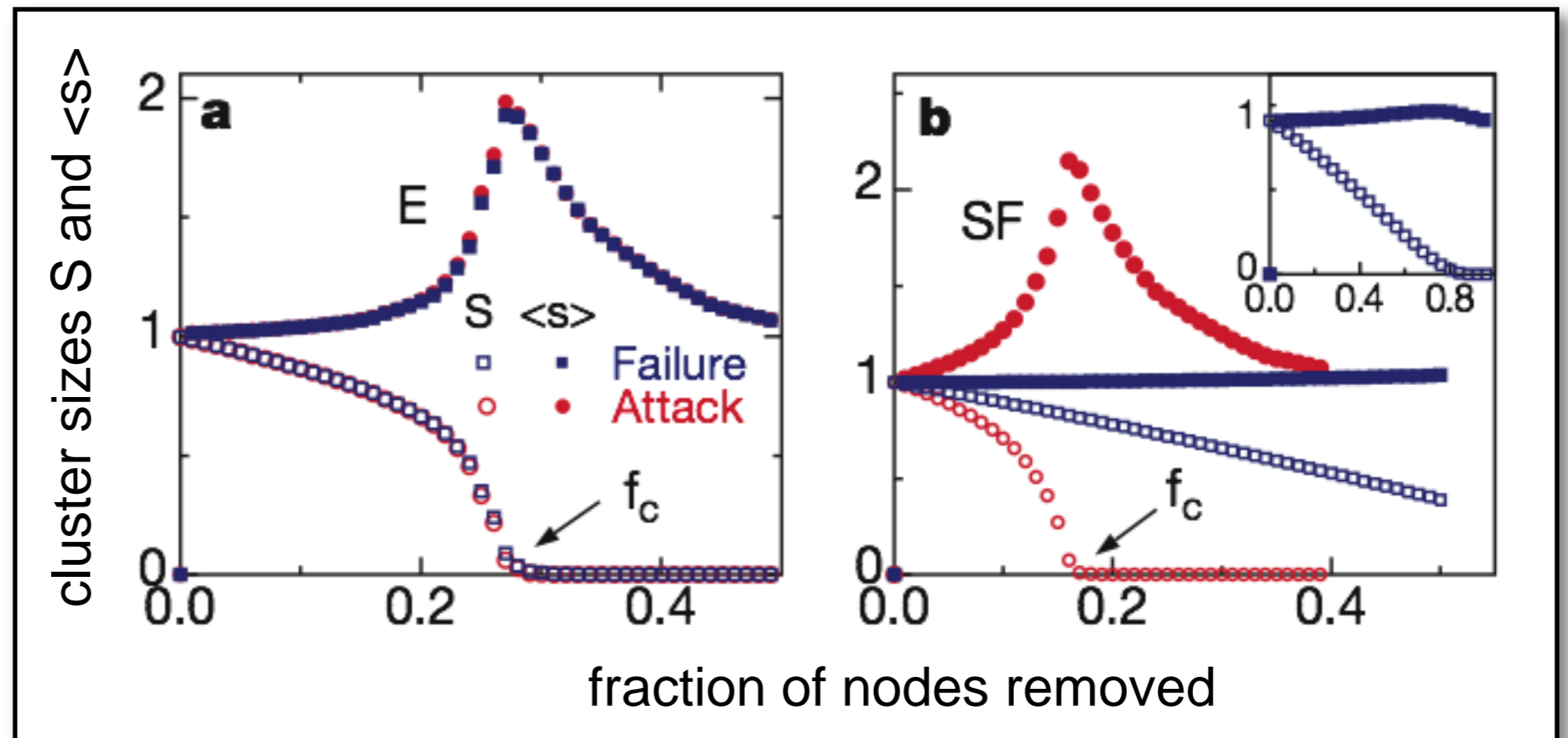


<http://moat.nlanr.net/Routing/rawdata/> :
6209 nodes and 12200 links (2000)

WWW-sample containing 325729
nodes and 1498353 links

Network Fragmentation

$\langle s \rangle$: average size of the isolated clusters (except the largest one)
 S : relative size of the largest cluster S ; this is defined as the fraction of nodes contained in the largest cluster (that is, $S = 1$ for $f = 0$)



Random network:

- **no difference** between attack and failure (homogeneity)
- fragmentation threshold at $f_c \gtrsim 0.28$ ($S \approx 0$)

Scale-free network:

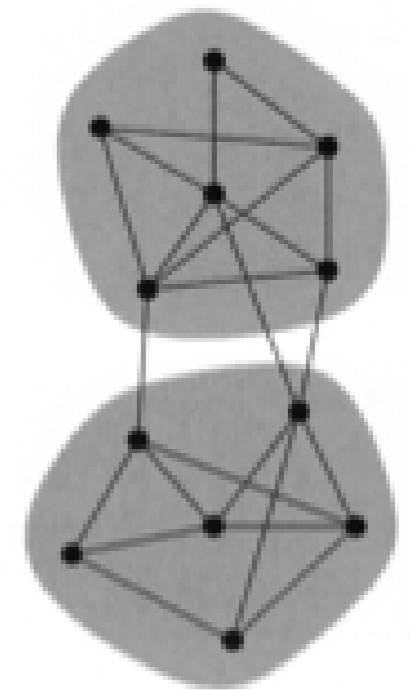
- **delayed fragmentation** and isolated nodes for failure
- critical breakdown under attack at $f_c \approx 0.18$

Modularity: an example of graph partitioning

The simplest graph partitioning problem is the division of a network into just 2 parts. This is called **graph bisection**.

If we can divide a network into 2 parts, we can also divide it further by dividing one or both of these parts ...

graph bisection problem: divide the vertices of a network into 2 non-overlapping groups of given sizes such that the **number of edges** running **between** vertices in **different groups** is **minimized**.



The number of edges between groups is called the **cut size**.

In principle, one could simply look through all possible divisions of the network into 2 parts and choose the one with smallest cut size.

Algorithms for graph partitioning

But this exhaustive search is prohibitively expensive!

Given a network of n vertices. There are $\frac{n!}{n_1!n_2!}$ different ways of dividing it into 2 groups of n_1 and n_2 vertices.

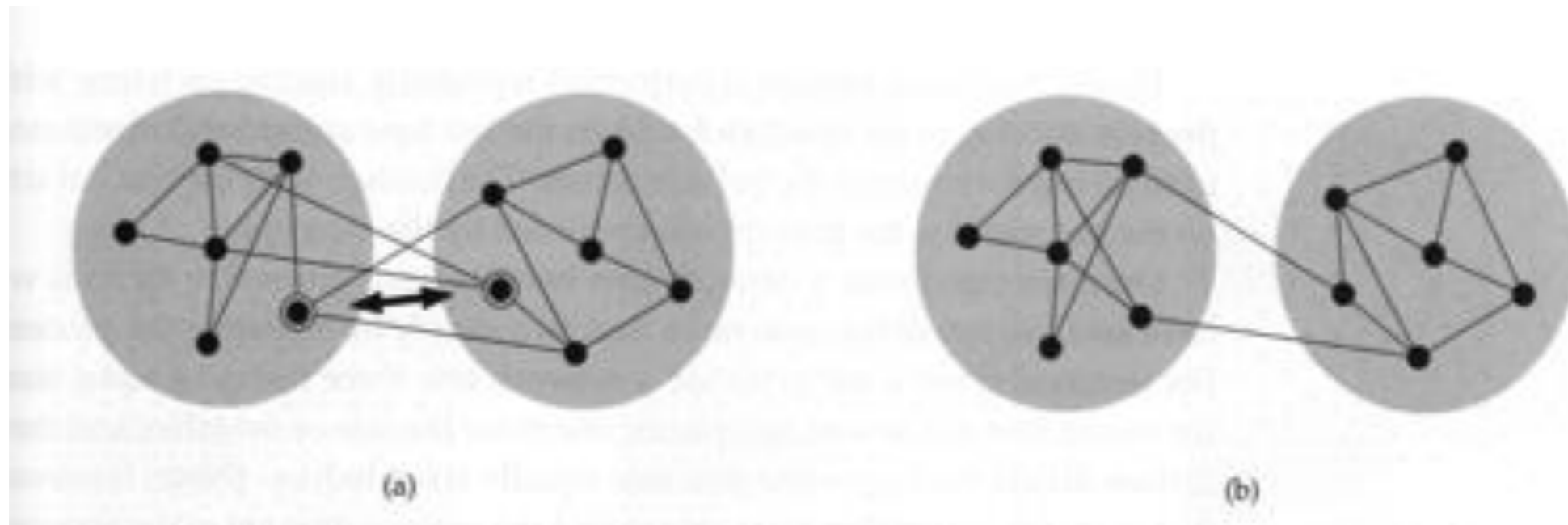
The amount of time to look through all these divisions will go up roughly exponentially with the size of the system.

Only values of up to $n = 30$ are feasible with current computers.

In computer science, either an algorithm can be clever and run quickly, but will fail to find the optimal answer in some (and perhaps most) cases, or it will always find the optimal answer, but takes an impractical length of time to do it.

The Kernighan-Lin algorithm

This algorithm proposed by Brian Kernighan and Shen Lin in 1970 is one of the simplest and best known heuristic algorithms for the graph bisection problem. (Kernighan is also one of the developers of the C language).



(a) The algorithm starts with any division of the vertices of a network into two groups (shaded) and then searches for pairs of vertices, such as the pair highlighted here, whose interchange would reduce the cut size between the groups.

(b) The same network after interchange of the 2 vertices.

The Kernighan-Lin algorithm

- (1) Divide the vertices of a given network into 2 groups (e.g. randomly)
- (2) For each pair (i,j) of vertices, where i belongs to the first group and j to the second group, calculate how much the cut size between the groups would change if i and j were interchanged between the groups.
- (3) Find the pair that reduces the cut size by the largest amount.

If no pair reduces it, find the pair that increases it by the smallest amount.

Repeat this process, but with the important restriction that each vertex in the network can only be moved once.

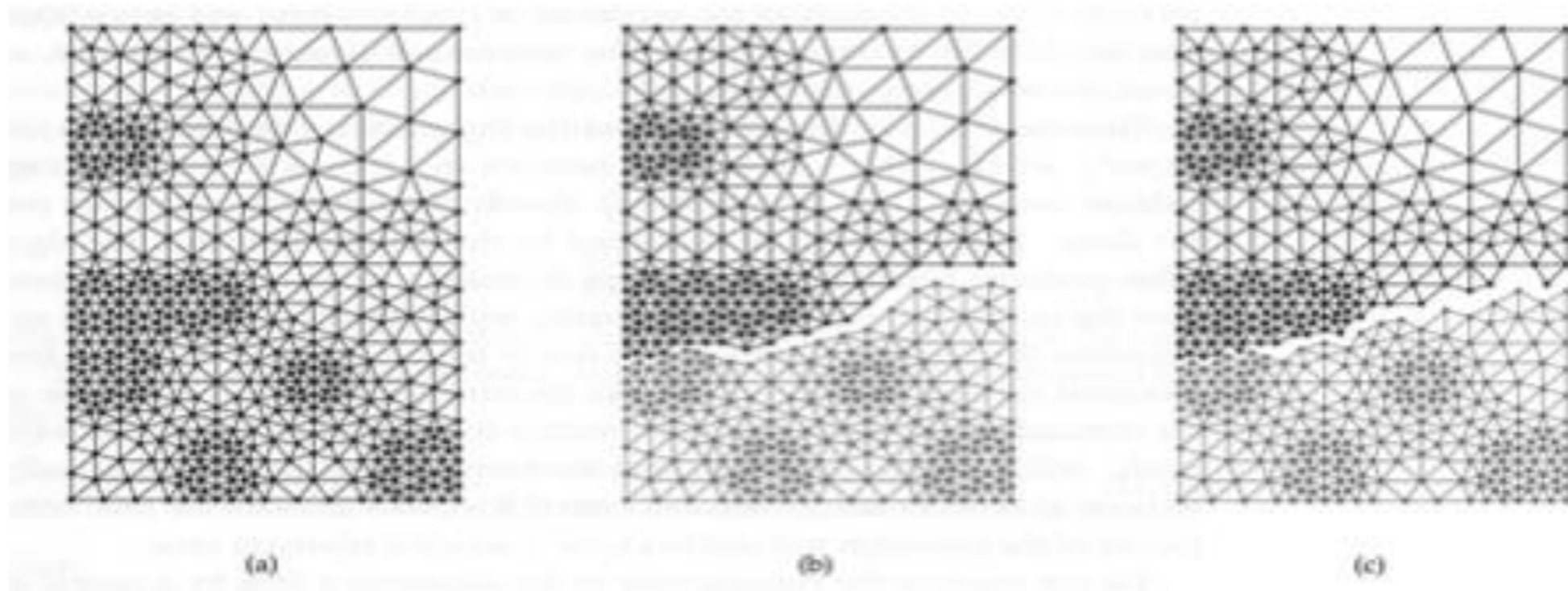
Stop when there is no pair of vertices left that can be swapped.

The Kernighan-Lin algorithm (II)

- (3) Go back through every state that the network passed through during the swapping procedure and choose among them the state in which the cut size takes its smallest value.
- (4) Perform this entire process repeatedly, starting each time with the best division of the network found in the last round.
- (5) Stop when no improvement on the cut size occurs.

Note that if the initial assignment of vertices to group is done randomly, the Kernighan-Lin algorithm may give (slightly) different answers when it is run twice on the same network.

The Kernighan-Lin algorithm (II)



(a) A mesh network of 547 vertices of the kind commonly used in finite element analysis.

(b) The best division found by the Kernighan-Lin algorithm when the task is to split the network into 2 groups of almost equal size.

This division involves cutting 40 edges in this mesh network and gives parts of 273 and 274 vertices.

(c) The best division found by spectral partitioning (alternative method).

Runtime of the Kernighan-Lin algorithm

The number of swaps performed during one round of the algorithm is equal to the smaller of the sizes of the two groups $\in [0, n / 2]$.

→ in the worst case, there are $O(n)$ swaps.

For each swap, we have to examine all pairs of vertices in different groups to determine how the cut size would be affected if the pair was swapped.

In the worst case, there are $n / 2 \times n / 2 = n^2 / 4$ such pairs, which is $O(n^2)$.

Runtime of the Kernighan-Lin algorithm (ii)

When a vertex i moves from one group to the other group, any edges connecting it to vertices in its current group become edges between groups after the swap.

Let us suppose that there are k_i^{same} such edges.

Similarly, any edges that i has to vertices in the other group, (say k_i^{other} ones) become within-group edges after the swap.

There is one exception. If i is being swapped with vertex j and they are connected by an edge, then the edge is still between the groups after the swap

→ the change in the cut size due to the movement of i is $k_i^{other} - k_i^{same} - A_{ij}$

A similar expression applies for vertex j .

→ the total change in cut size due to the swap is $k_i^{other} - k_i^{same} + k_j^{other} - k_j^{same} - 2A_{ij}$

Runtime of the Kernighan-Lin algorithm (iii)

For a network stored in adjacency list form, the evaluation of this expression involves running through all the neighbors of i and j in turn, and hence takes time on the order of the average degree in the network, or $O(m/n)$ with m edges in the network.

→ the total running time is $O(n \times n^2 \times m/n) = O(mn^2)$

On a sparse network with $m \propto n$, this is $O(n^3)$

On a dense network (with $m \rightarrow \frac{n(n-1)}{2}$), this is $O(n^4)$

This time still needs to be multiplied by the number of rounds the algorithm is run before the cut size stops decreasing.

For networks up to a few 1000 of vertices, this number may be between 5 and 10.

Mesoscale properties of networks

- identify cliques and highly connected clusters

Most relevant processes in biological networks correspond to the mesoscale (5-25 genes or proteins) not to the entire network.

However, it is computationally enormously expensive to study mesoscale properties of biological networks.

e.g. a network of 1000 nodes contains 1×10^{23} possible 10-node sets.

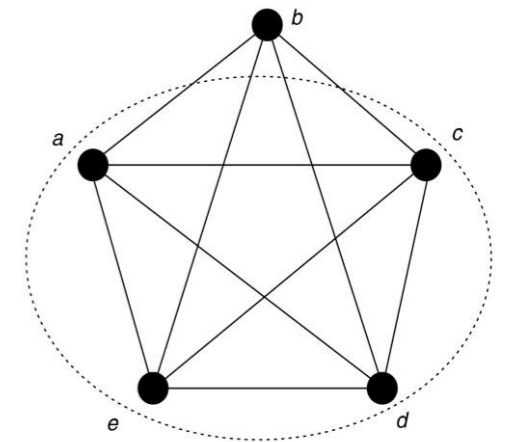
Spirin & Mirny analyzed combined network of protein interactions in *S. cerevisiae* with data from CELLZOME, MIPS, BIND: 6500 interactions.

Identify connected subgraphs

The network of protein interactions is typically presented as an undirected graph with proteins as nodes and protein interactions as undirected edges.

First aim: identify **fully connected subgraphs** (cliques)

A clique is a set of nodes that are all neighbors of each other.



The „maximum clique problem“ – finding the largest clique in a given graph is known to be NP-hard.

In this example, the whole graph is a clique and consequently any subset of it is also a clique, for example $\{a, c, d, e\}$ or $\{b, e\}$.

A **maximal clique** is a clique that is not contained in any larger clique. Here only $\{a, b, c, d, e\}$ is a maximal clique.

In general, protein complexes need not to be fully connected.

Spirin, Mirny,
PNAS 100, 12123 (2003)

Identify all fully connected subgraphs (cliques)

Although the general problem - finding all cliques of a graph - is very hard, this can be done relatively quickly for the given network because the protein interaction graph is quite sparse (the number of interactions (edges) is similar to the number of proteins (nodes)).

To find cliques of size n one needs to enumerate only the cliques of size $n-1$.

The search for cliques starts with $n = 4$, pick all (known) pairs of edges (6500 \times 6500 protein interactions) successively.

For every pair $A-B$ and $C-D$ check whether there are edges between A and C , A and D , B and C , and B and D . If these edges are present, $ABCD$ is a clique.

For every clique identified, $ABCD$, pick all known proteins successively.

For every picked protein E , if all of the interactions $E-A$, $E-B$, $E-C$, and $E-D$ exist, then $ABCDE$ is a clique with size 5.

Continue for $n = 6, 7, \dots$

The largest clique found in the protein-interaction network has size 14.

Spirin, Mirny, PNAS 100, 12123 (2003)

Identify all fully connected subgraphs (cliques)

These results include, however, many redundant cliques.

For example, the clique with size 14 contains 14 cliques with size 13.

To find all nonredundant subgraphs, mark all proteins comprising the clique of size 14, and out of all subgraphs of size 13 pick those that have at least one protein other than marked.

After all redundant cliques of size 13 are removed, proceed to remove redundant twelves etc.

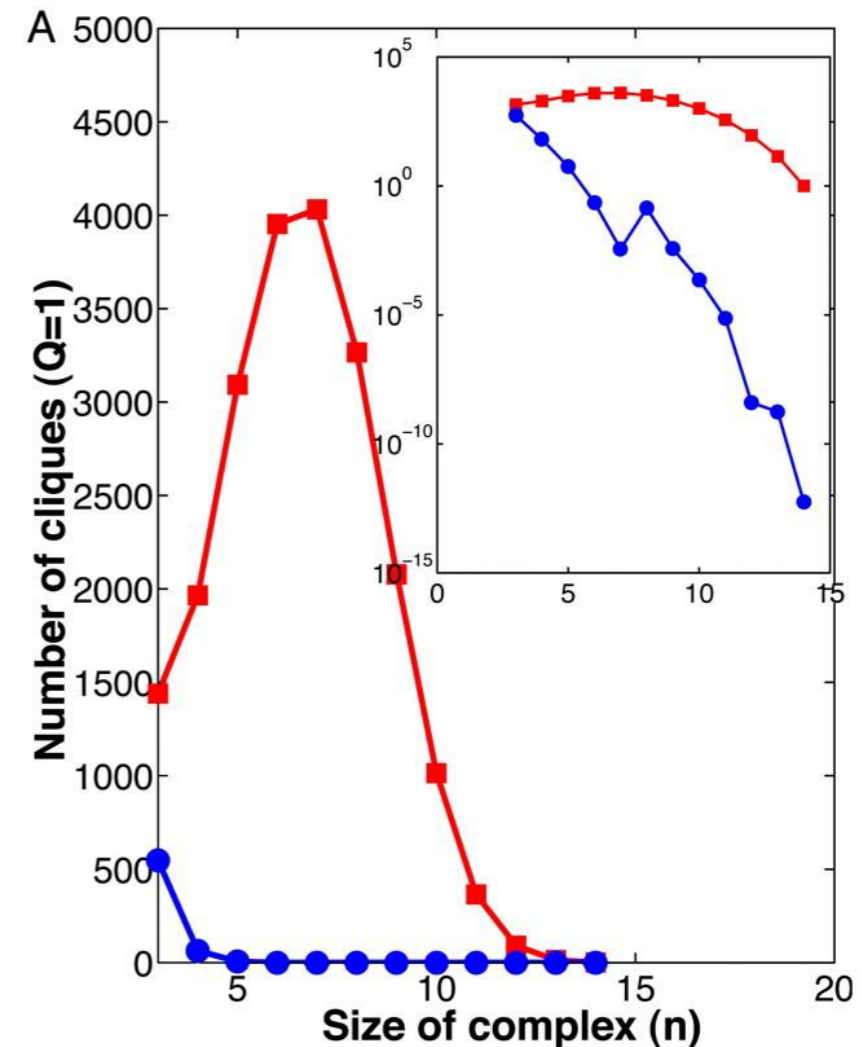
In total, only 41 nonredundant cliques with sizes 4 - 14 were found by Spirin & Mirny.

Spirin, Mirny, PNAS 100, 12123 (2003)

Statistical significance of cliques

Number of complete cliques as a function of clique size enumerated in the network of protein interactions (red) and in randomly rewired graphs (blue, averaged over >1,000 graphs where the number of interactions for each protein is preserved).

Inset shows the same plot on a log-normal scale. Note the dramatic enrichment in the number of cliques in the protein-interaction graph compared with the random graphs. Most of these cliques are parts of bigger complexes and modules.



Spirin, Mirny, PNAS 100, 12123 (2003)

Method [bioRxiv preprint doi: <https://doi.org/10.1101/000000>; this version posted May 1, 2014. The copyright holder for this preprint \(which was not certified by peer review\) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.](#)

Modular decomposition of protein-protein interaction networks

Julien Gagneur^{*†}, Roland Krause^{*}, Tewis Bouwmeester^{*} and Georg Casari^{*}

Addresses: ^{*}Cellzome AG, Meyerhofstrasse 1, 69117 Heidelberg, Germany. [†]Laboratoire de Mathématiques Appliquées aux Systèmes, Ecole Centrale Paris, Grande Voie des Vignes, 92295 Châtenay-Malabry cedex, France.



Abstract

We introduce an algorithmic method, termed modular decomposition, that defines the organization of protein-interaction networks as a hierarchy of nested modules. Modular decomposition derives the logical rules of how to combine proteins into the actual functional complexes by identifying groups of proteins acting as a single unit (sub-complexes) and those that can be alternatively exchanged in a set of similar complexes. The method is applied to experimental data on the pro-inflammatory tumor necrosis factor- α (TNF- α)/NF κ B transcription factor pathway.

Shared Components

Shared components = proteins or groups of proteins occurring in different complexes are fairly common. A shared component may be a small part of many complexes, acting as a **unit** that is constantly **reused** for its function.

Also, it may be the **main part** of the complex e.g. in a family of variant complexes that differ from each other by distinct proteins that provide functional specificity.

Aim: **identify** and properly **represent** the modularity of protein-protein interaction networks by identifying the **shared components** and the way they are arranged to generate **complexes**.

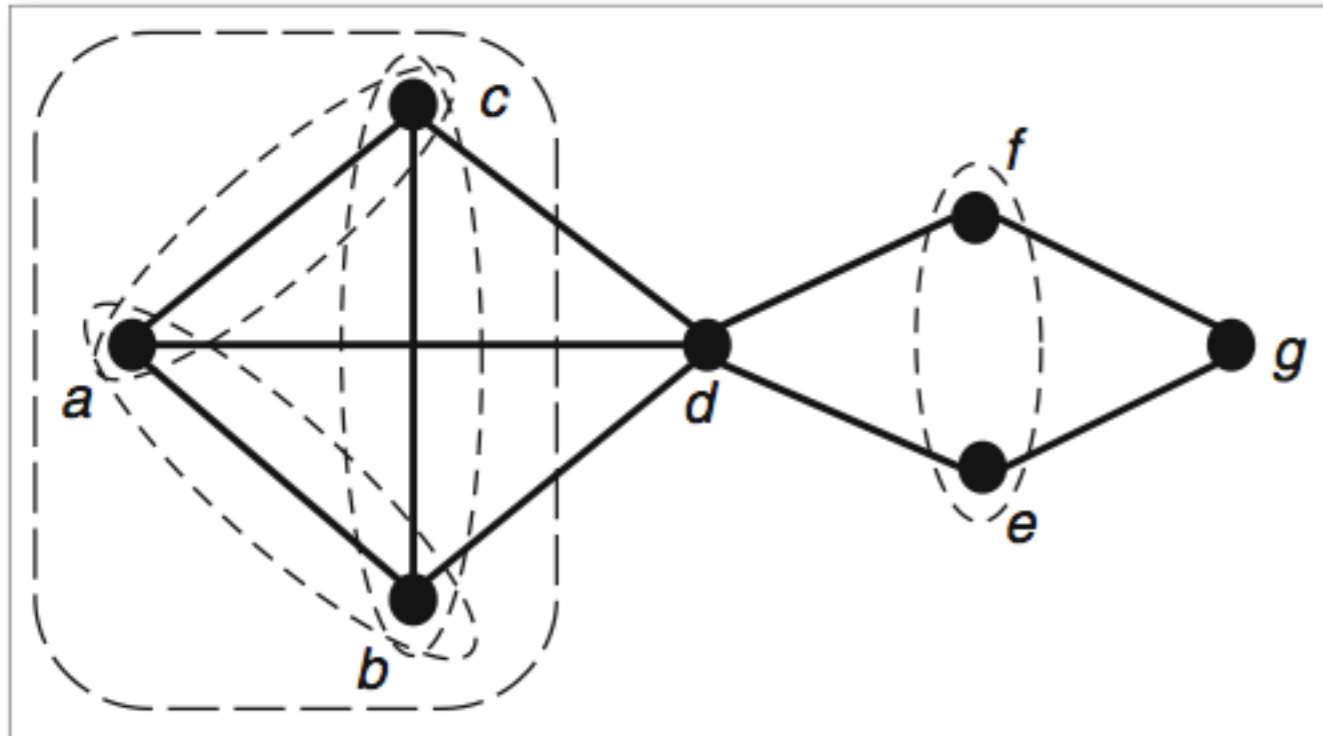


Georg Casari, Cellzome (Heidelberg)

Gagneur et al. Genome Biology 5, R57 (2004)

Modular Decomposition of a Graph

Module := set of **nodes** that have the
same neighbors outside of the module



trivial modules:

$\{a\}, \{b\}, \dots, \{g\}$

$\{a, b, \dots, g\}$

non-trivial modules:

$\{a, b\}, \{a, c\}, \{b, c\}$

$\{a, b, c\}$

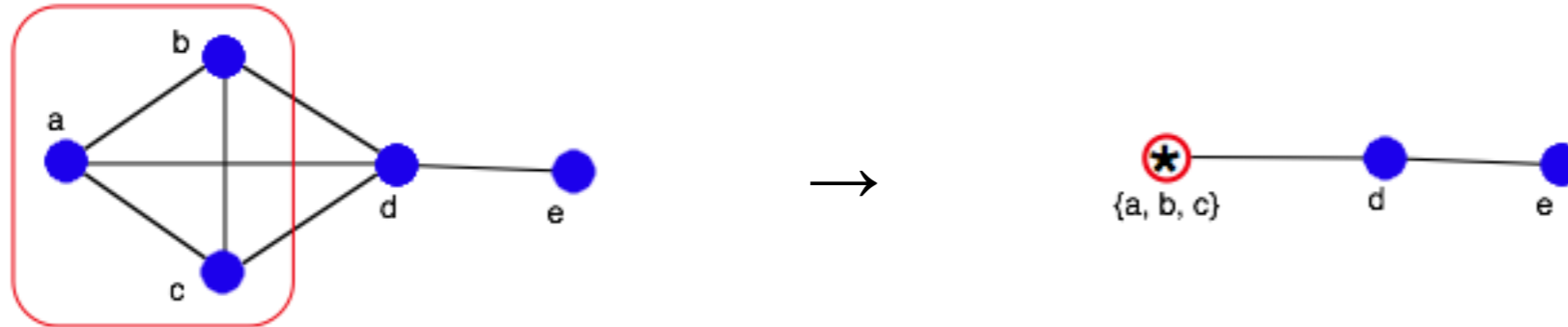
$\{e, f\}$

Quotient: representative node for a module

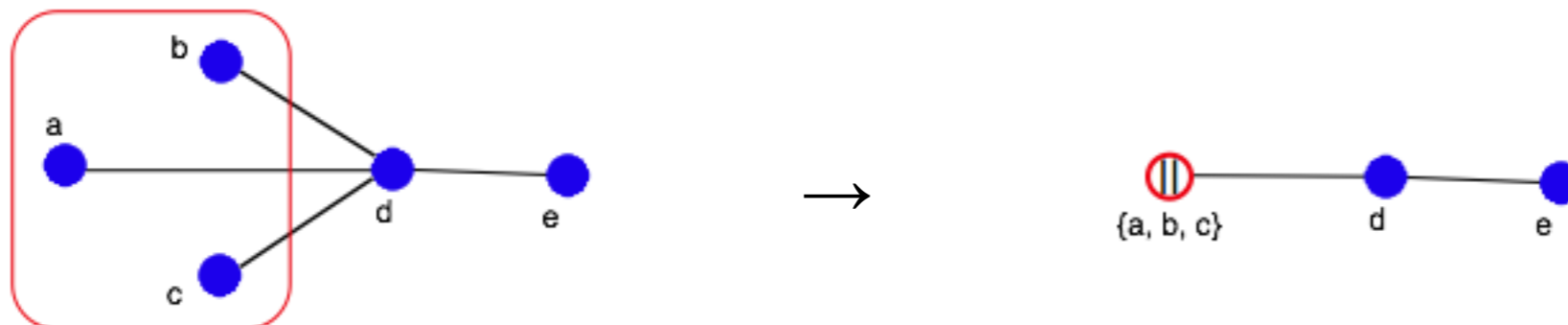
Iterated quotients → labeled tree representing the original network
→ "**modular decomposition**"

Quotients

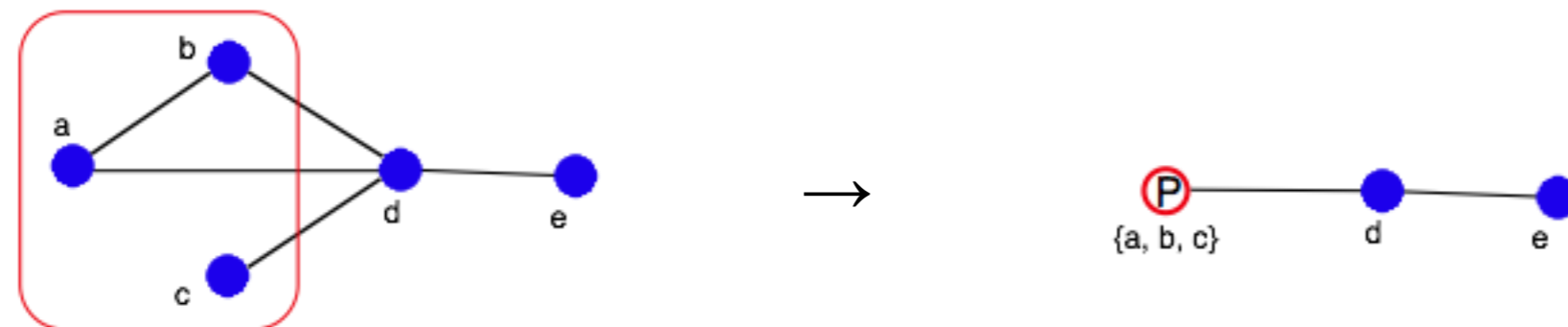
Series: all included nodes are direct **neighbors** (= **clique**)



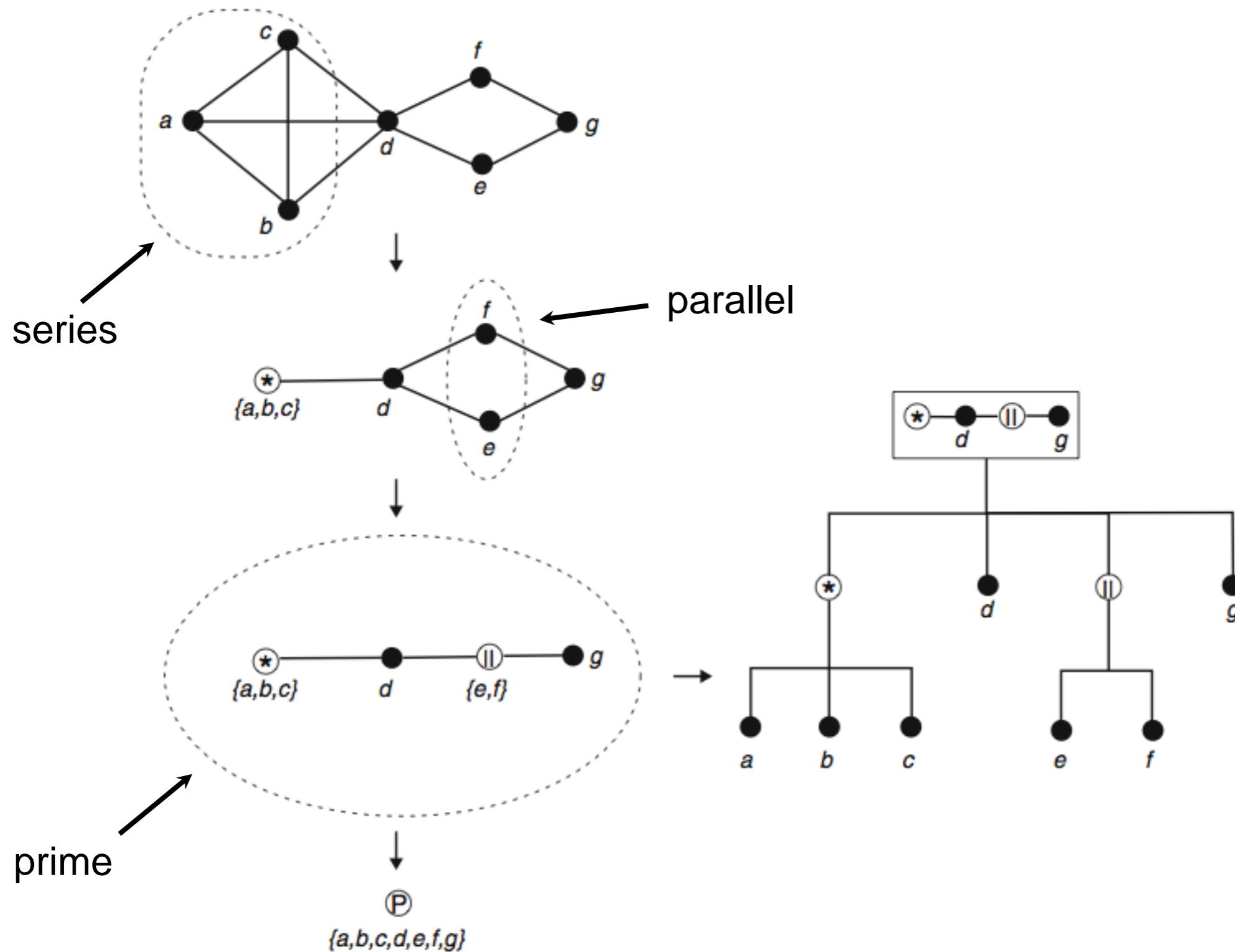
Parallel: all included nodes are **non-neighbors**



Prime: "anything else" (best labeled with the actual structure)



A Simple Recursive Example



Using data from protein complex purifications e.g. by TAP

Different types of data:

- Y2H: detects direct physical interactions between proteins
 - PCP by tandem affinity purification with mass-spectrometric identification of the protein components identifies multi-protein complexes
- Molecular decomposition will have a **different meaning** due to different **semantics** of such graphs.

Here, we focus analysis on **PCP content** from TAP-MS data.

PCP experiment: select bait protein where TAP-label is attached → Co-purify protein with those proteins that co-occur in at least one complex with the bait protein.

Gagneur et al. Genome Biology 5, R57 (2004)

Data from Protein Complex Purification

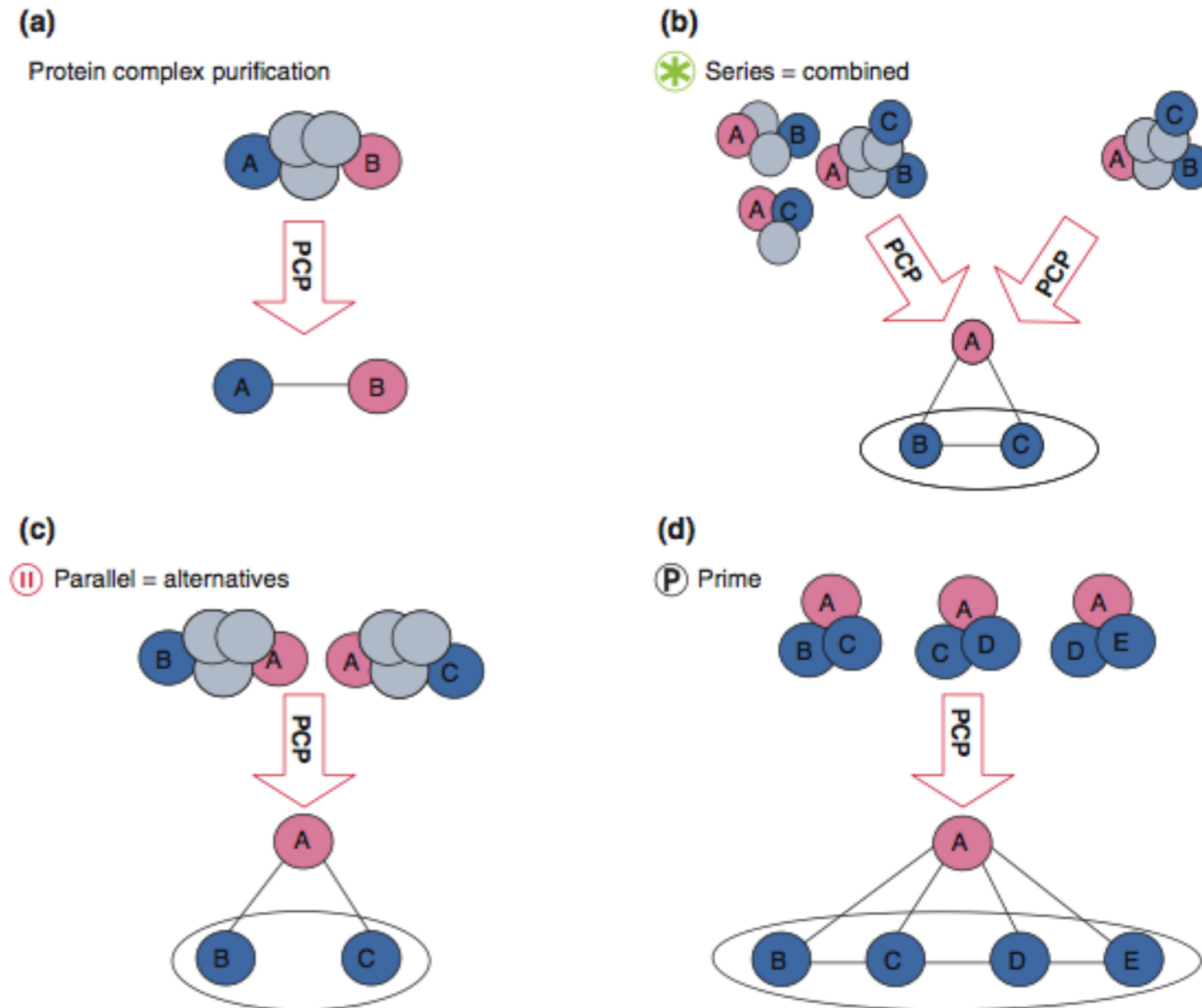
Graphs and module labels from systematic PCP experiments:

(a) Two neighbors in the network are proteins occurring in a same complex.

(b) Several potential sets of complexes can be the origin of the same observed network. Restricting interpretation to the simplest model (top right), the **series** module reads as a logical AND between its members.

(c) A module labeled **parallel** corresponds to proteins or modules working as strict alternatives with respect to their common neighbors.

(d) The **prime** case is a structure where none of the two previous cases occurs.



Gagneur et al. Genome Biology 5, R57 (2004)

Real World Examples

Two examples of modular decompositions of protein-protein interaction networks.

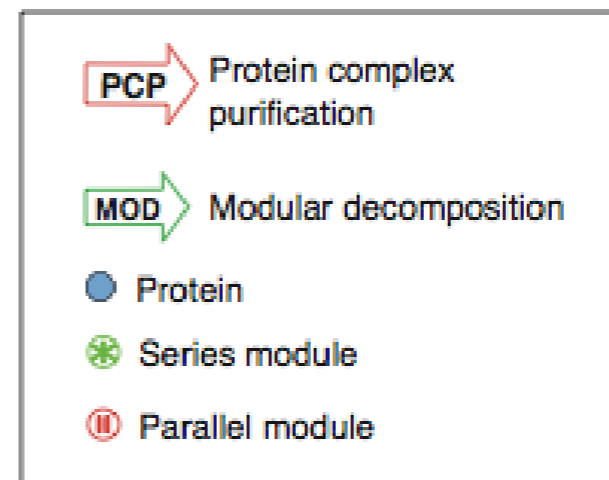
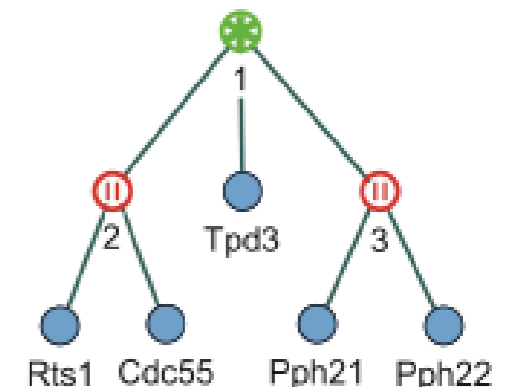
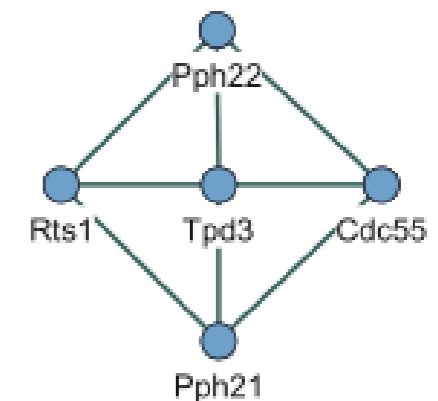
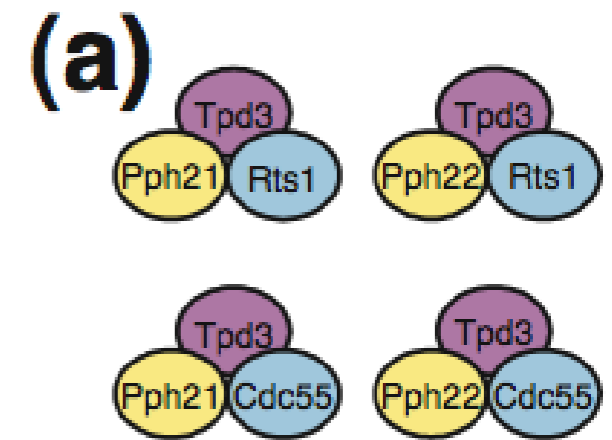
In each case from top to bottom: schemata of the complexes, the corresponding protein-protein interaction network as determined from PCP experiments, and its modular decomposition (MOD).

(a) Protein phosphatase 2A.

Parallel modules group proteins that do not interact but are functionally equivalent.

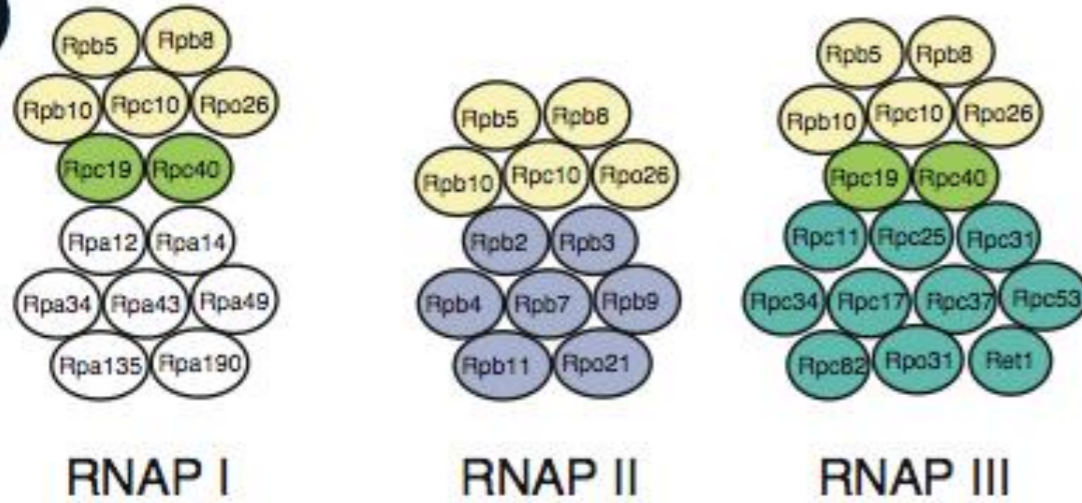
Here these are the catalytic proteins Pph21 and Pph22 (module 2) and the regulatory proteins Cdc55 and Rts1 (module 3), connected by the Tpd3 „backbone“.

Notes: • Graph does not show functional alternatives!!!
• other decompositions also possible

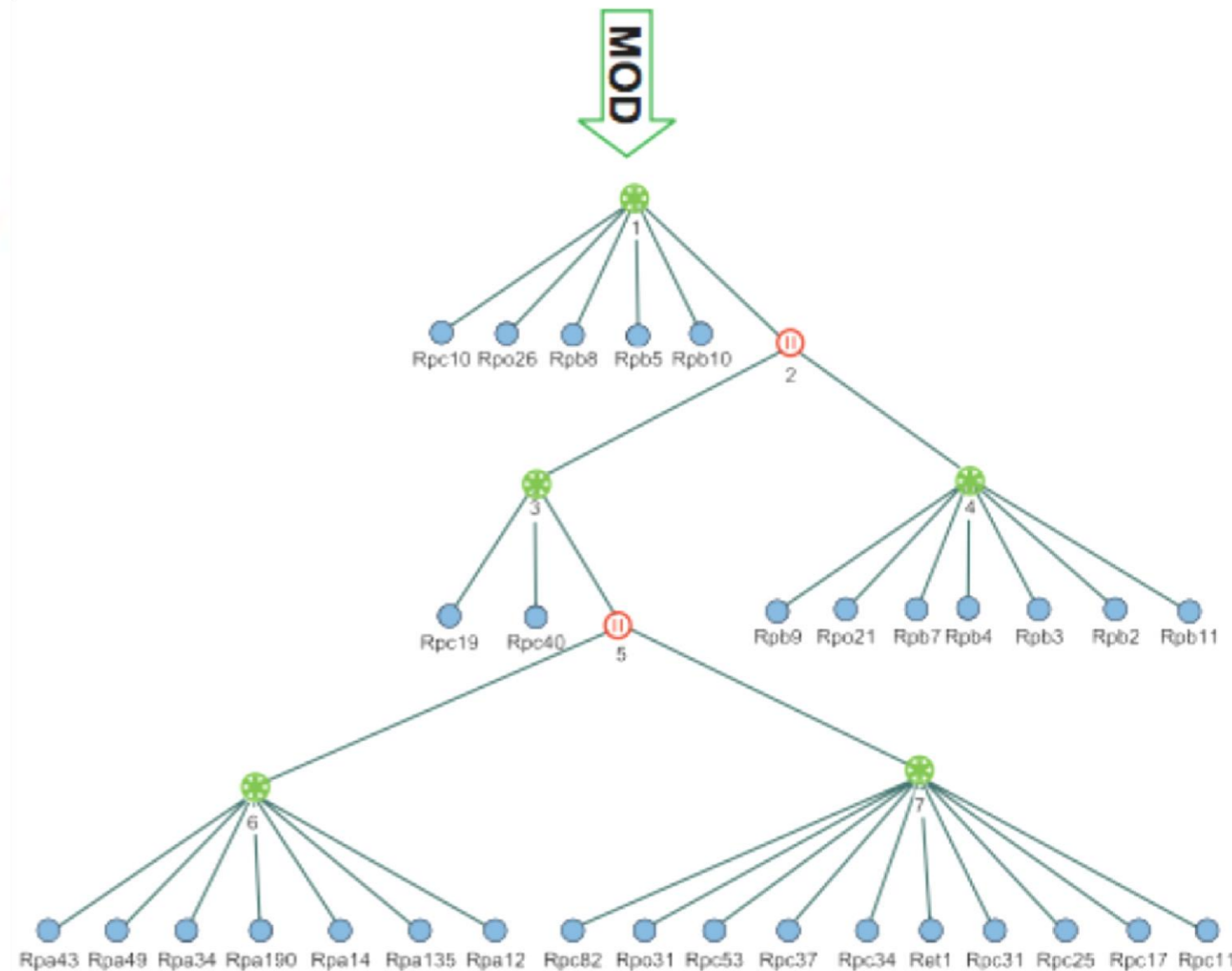
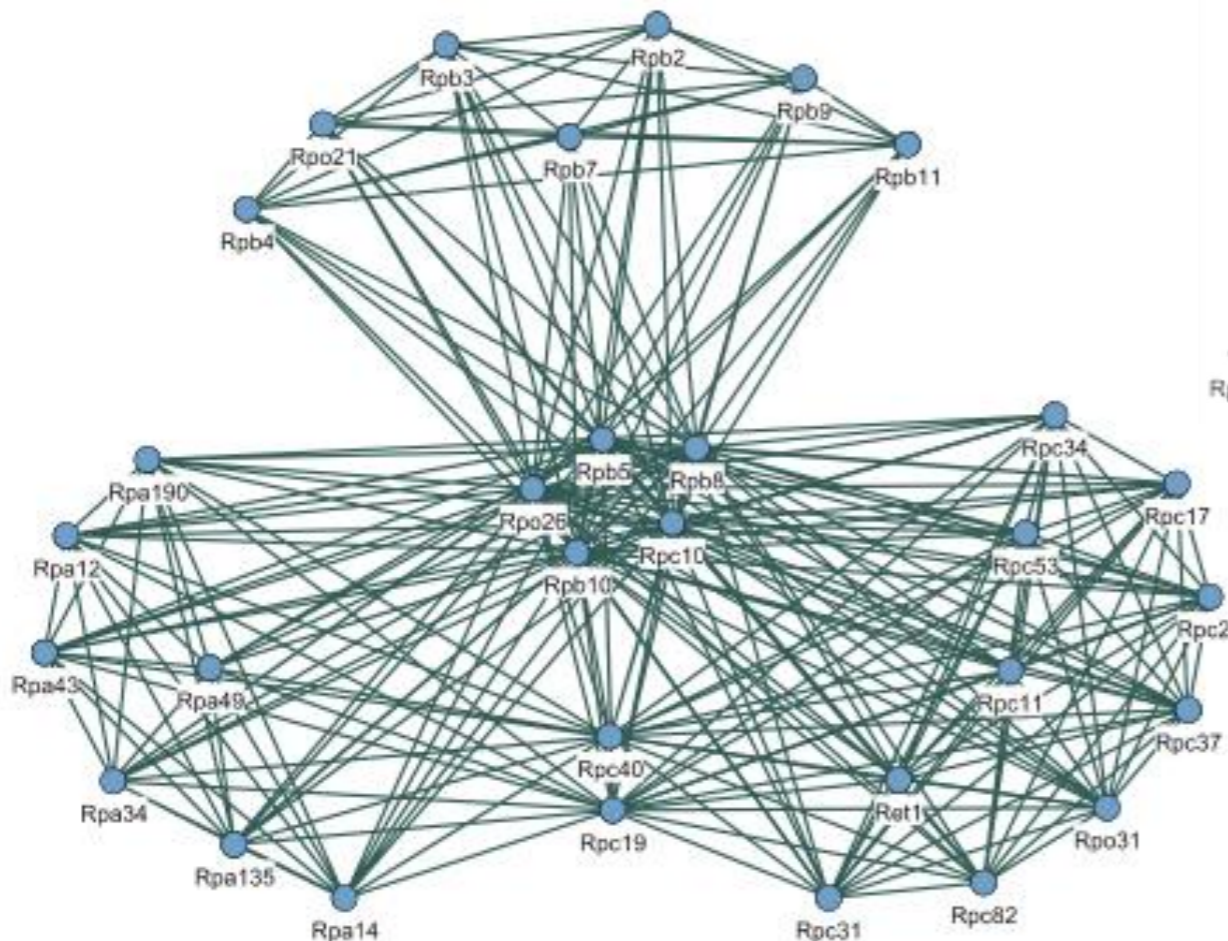


RNA polymerases I, II and III

(b)



PCP



Again: modular decomposition is much easier to understand than the connectivity graph

Gagneur et al. Genome Biology 5, R57 (2004)

Summary

Modular decomposition of graphs is a **well-defined concept**.

- One can prove thoroughly for which graphs a modular decomposition exists.
- Efficient $O(m + n)$ algorithms exist to compute the decomposition.

However, experiments have shown that **biological** complexes are **not strictly disjoint**. They often share components

→ separate complexes do not always fulfill the strict requirements of modular graph decomposition.

Also, there exists a „danger“ of false-positive or false-negative interactions.

→ **other methods**, e.g., for detecting communities (Girven & Newman) or densely connected clusters are **more suitable** for identification of **complexes** because they are more sensitive.

Power Graph Analysis

OPEN ACCESS Freely available online

PLOS COMPUTATIONAL BIOLOGY

Unraveling Protein Networks with Power Graph Analysis

Loïc Royer, Matthias Reimann, Bill Andreopoulos, Michael Schroeder*

Biotechnology Center, Technische Universität Dresden, Germany

PLoS Comp Biol 4 (2008) e1000108

Lossless compact abstract representation of graphs:

- **Power nodes** = set of nodes (criterion for grouping?)
- **Power edges** = edges between power nodes

Exploit observation that **cliques** and bi-cliques are **abundant** in real networks
→ **explicitly** represented in power graphs

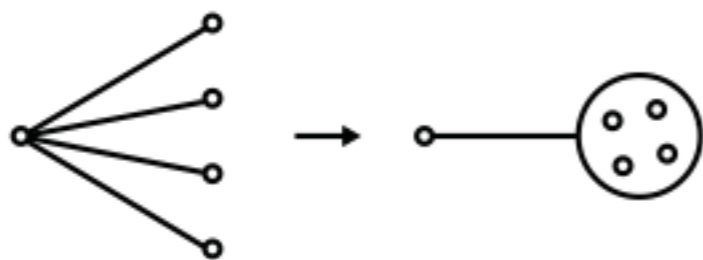
Power Nodes

In words: "... if two **power nodes** are **connected** by a **power edge** in G' , this means in G that **all nodes** of the first power node are **connected to all nodes** of the second power node.

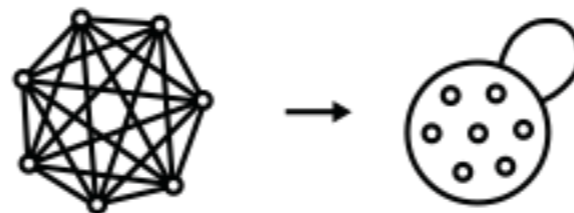
Similarly, if a power node is connected to itself by a power edge in G' , this means that all nodes in the power node are connected to each other by edges in G .

With: "real-world" graph $G = \{V, E\}$
power graph $G' = \{V', E'\}$

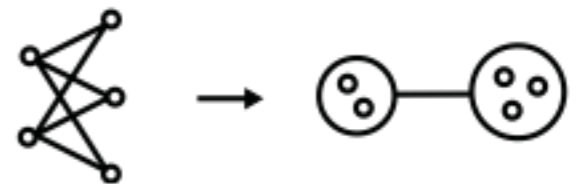
Star motif



Clique motif



Biclique motif



Power Graph Analysis Algorithm

Two **conditions**:

- power **node hierarchy** condition:

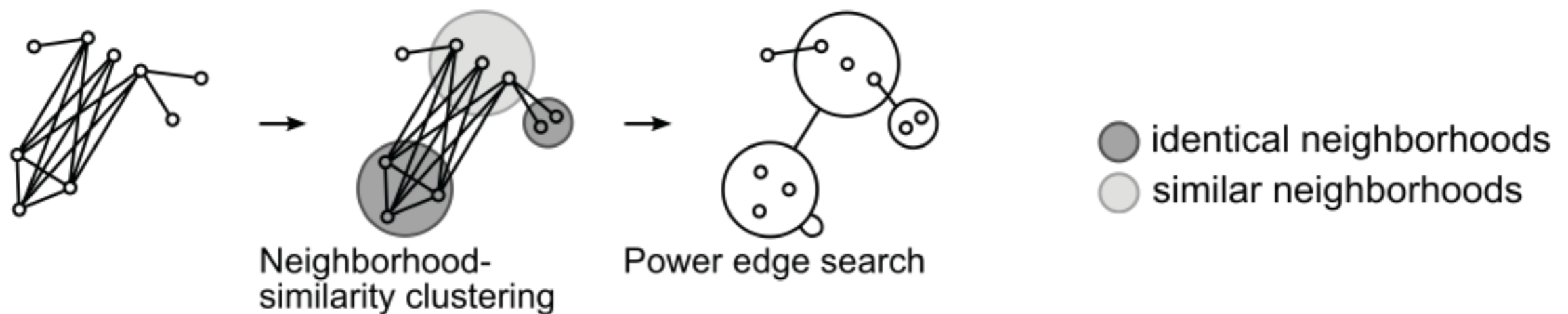
two power nodes are either disjoint, or one is included in the other one

- power **edge disjointness** condition: each edge of the original graph is represented by one and only one power edge

Algorithm:

- 1) identify potential power nodes with hierarchical clustering based on neighborhood similarity

- 2) greedy power edge search



Complex = Star or Clique?

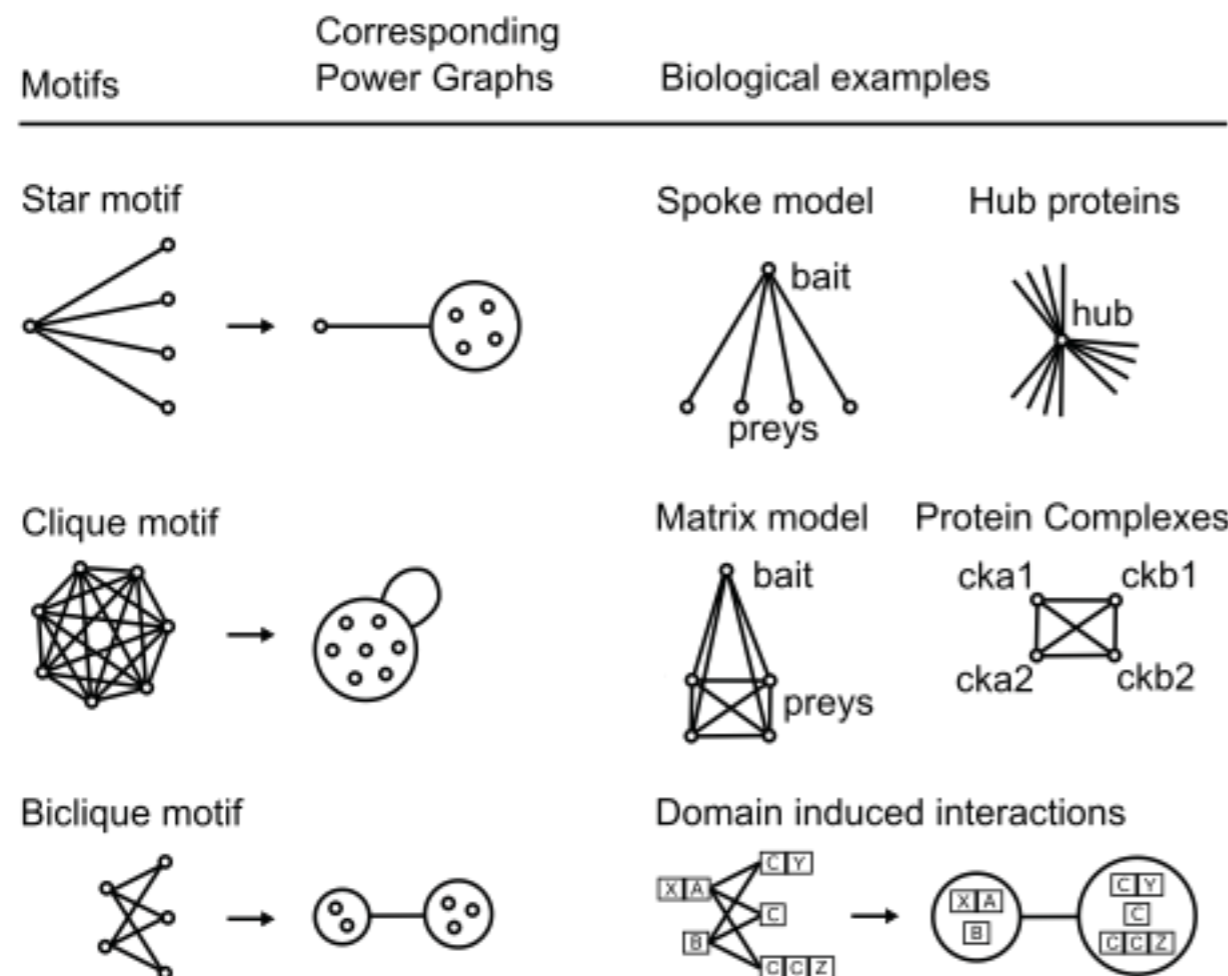


Figure 1. The Three Basic Motifs: Star, Biclique, and Clique.

Stars often occur because of hub proteins or when affinity purification complexes are interpreted using the spoke model. Bicliques often arise because of domain-domain or domain-motif interactions inducing protein interactions [25]. Power nodes are sets of nodes and power edges connect power nodes. A power edge between two power nodes signifies that all nodes of the first set are connected to all nodes of the second set. Note that nodes within a power node are not necessarily connected to each other.

doi:10.1371/journal.pcbi.1000108.g001

In **pull-down** experiments:
Bait is used to capture **complexes** of prey proteins
→ do they all just stick to the bait or to each other?

spoke model
→ **underestimates** connectivity

matrix model
→ **overestimates** connectivity

Casein Kinase II Complex

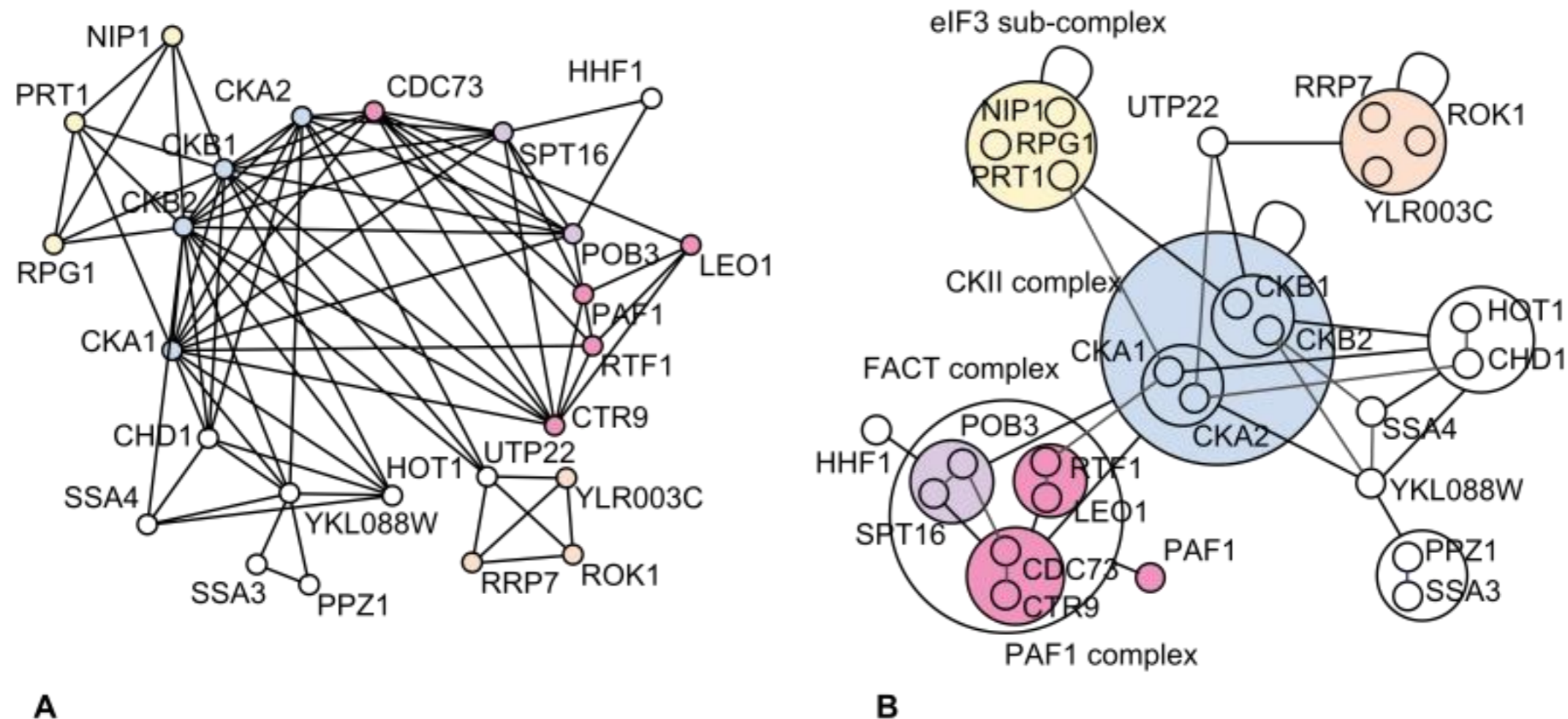


Figure 2. Casein Kinase II Complex. Two catalytic alpha subunits (CKA1, CKA2) and two regulatory beta subunits (CKB1, CKB2) interacting with the FACT complex, with sub-complex NIP1-RPG-PRT1, and with the PAF1 complex. The graph representation (A) consists of 80 edges whereas the power graph representation (B) has 30 power edges, thus an edge reduction of 62%. This simplification of the representation makes the separation of the regulatory subunits from the catalytic subunits immediately apparent without loss of information on individual interactions.
doi:10.1371/journal.pcbi.1000108.g002

→ Power graph: compressed and cleaner representation

Various Similarities

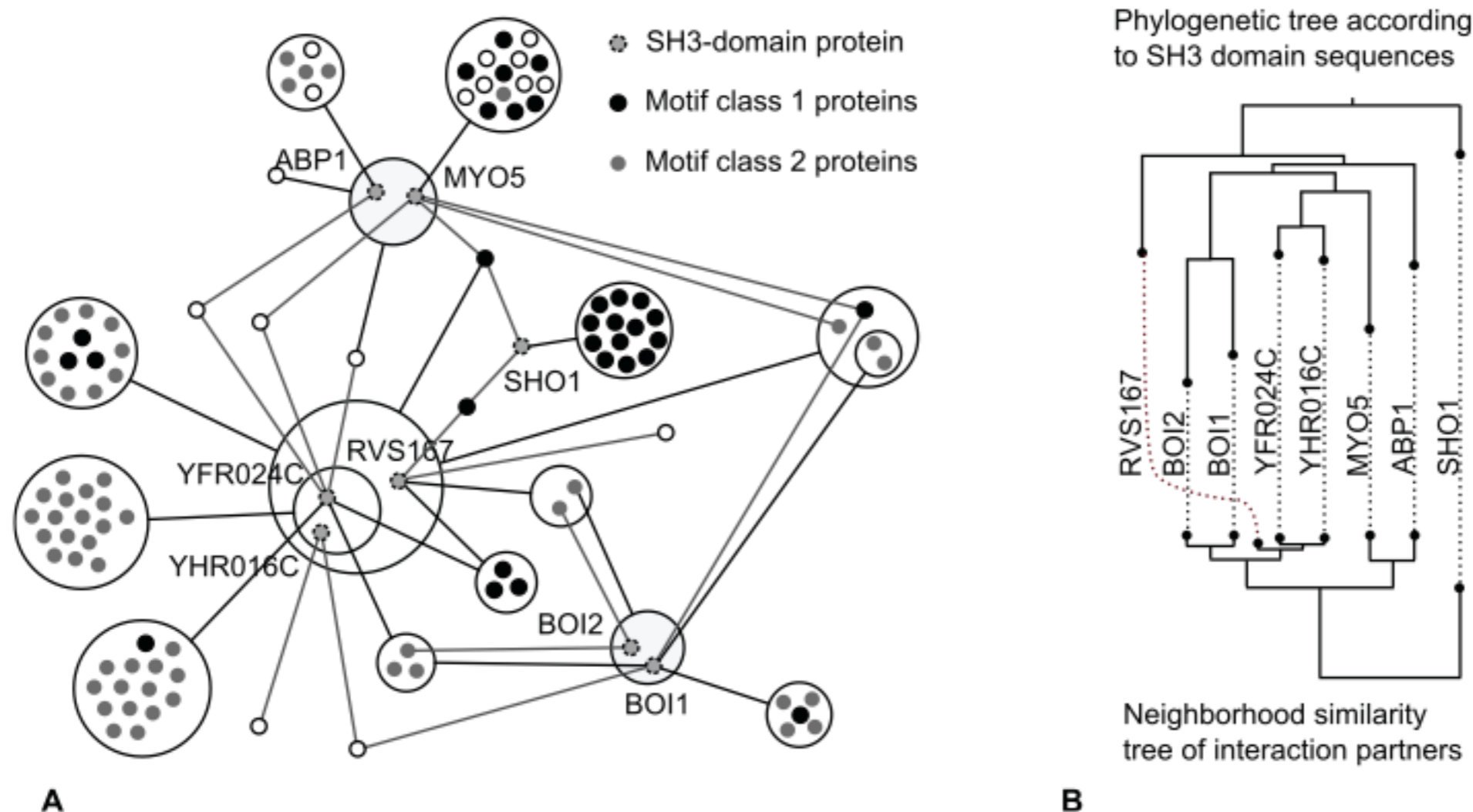


Figure 4. Interactions of SH3 Carrying Proteins. (A) Protein interaction network showing the 105 interaction partners of the SH3 domain carrying proteins: SHO1, ABP1, MYO5, BOI1, BOI2, RVS167, YHR016C and YFR024. The underlying network consists of 182 interactions represented here as 36 power edges—a reduction of 80%—leaving all but only the core information. Class 1 motif (RxxPxxP) proteins are shown in black. Class 2 motif (PxxPxR) proteins are shown in light grey [15]. Note how power graphs group proteins having similar binding motifs together. (B) Phylogeny and interaction profiles. Comparison of the phylogenetic tree of the SH3 domains sequences with the neighbourhood similarity tree of interaction partners. The neighbourhood similarity implied by the power graph reflects the sequence similarity of the SH3 domains.
doi:10.1371/journal.pcbi.1000108.g004

Network Compression

Power graph analysis: group nodes with **similar neighborhood**
 → often **functionally** related proteins end up in one power node

Lossless compression
 of graphs:
38...85% edge reduction
 for biological networks

Protein Interaction Network	# Nodes	# Edges	Avg. Degree	e.r.	c.r
Lim et al. (2006) [46]	571	701	2.45	85%	12.1
Hazbun et al. (2003) [47]	2243	3130	2.79	79%	13
Kim et al. (2006) [48]	577	1090	3.78	67%	4.1
Gunsalus et al. (2004) [49]	281	514	3.6	65%	4.6
Gavin et al. (2006) [4]	1462	6942	9.4	64%	7.2
Ewing et al. (2007) [50]	2294	6449	5.62	54%	6.6
Ito et al. (2001) [51]	3243	4367	2.69	53%	5.3
Rual et al. (2005) [12]	1527	2529	3.31	50%	4.5
Krogan et al. (2006) [6]	2708	7123	5.26	49%	4.5
Stanyon et al. (2004) [9]	478	1778	7.43	48%	5.3
Stanyon et al. (2004) [9]	478	1778	7.43	48%	5.3
Butland et al. (2005) [52]	1277	5324	8.33	43%	6.0
Arifuzzaman et al. (2006) [53]	2457	8663	7.05	39%	5.4
Lacount et al. (2005) [13]	1272	2643	4.16	38%	3.8

Average degree, edge reduction (e.r.), and edge to power node conversion rate (c.r.).
 doi:10.1371/journal.pcbi.1000108.t001

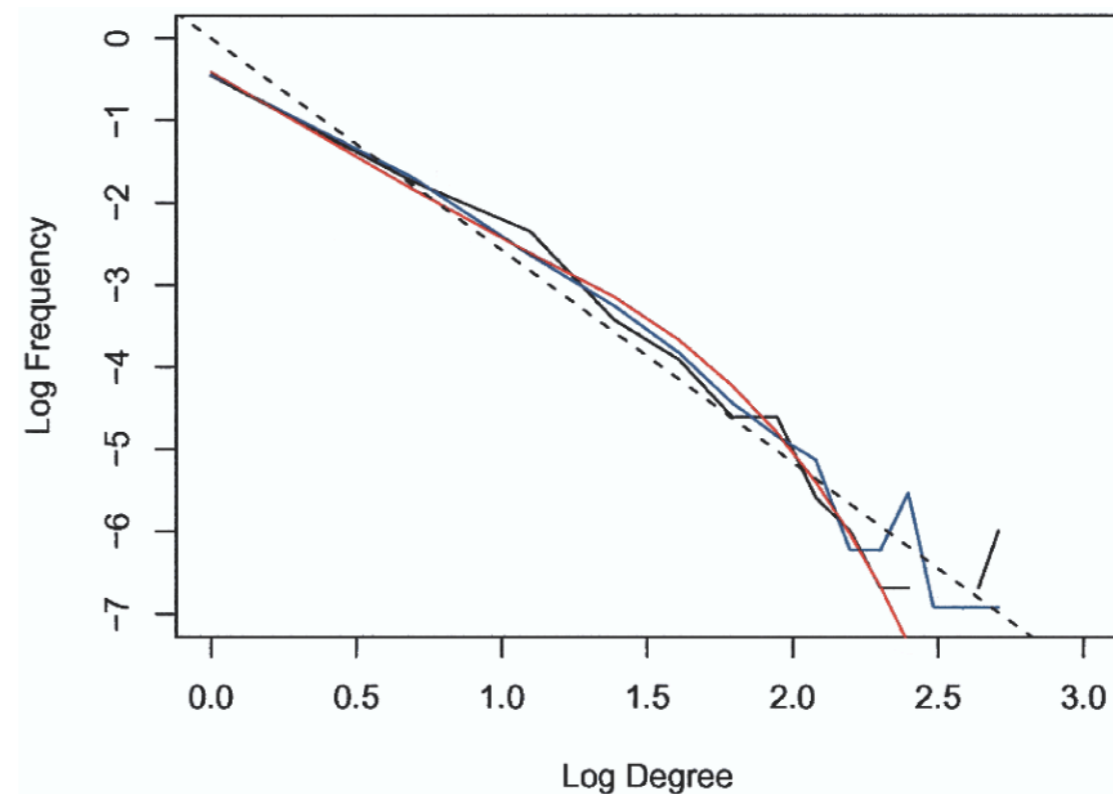
Royer et al, *PLoS Comp Biol* 4 (2008) e1000108

Some PPI Networks

For some time: "**Biological** networks are **scale-free**..."



Y2H PPI network from Uetz et al, *Nature* **403** (2003) 623



$P(k)$ compared to a power law

However, there are some doubts... → next lecture

Summary

What you learned **today**:

- Network **robustness**

- ☐ scale-free networks are failure-tolerant, but fragile to attacks
 - \Leftrightarrow the few **hubs** are important
 - \Rightarrow immunize hubs!

- **Modules** in networks

- \Rightarrow modular decomposition
- \Rightarrow power graph analysis

Next lecture:

- Are biological networks scale-free? (other models?)
- Network growth mechanisms

Short Test #1: Mon, Nov. 10

(covers lectures V1-V5)