

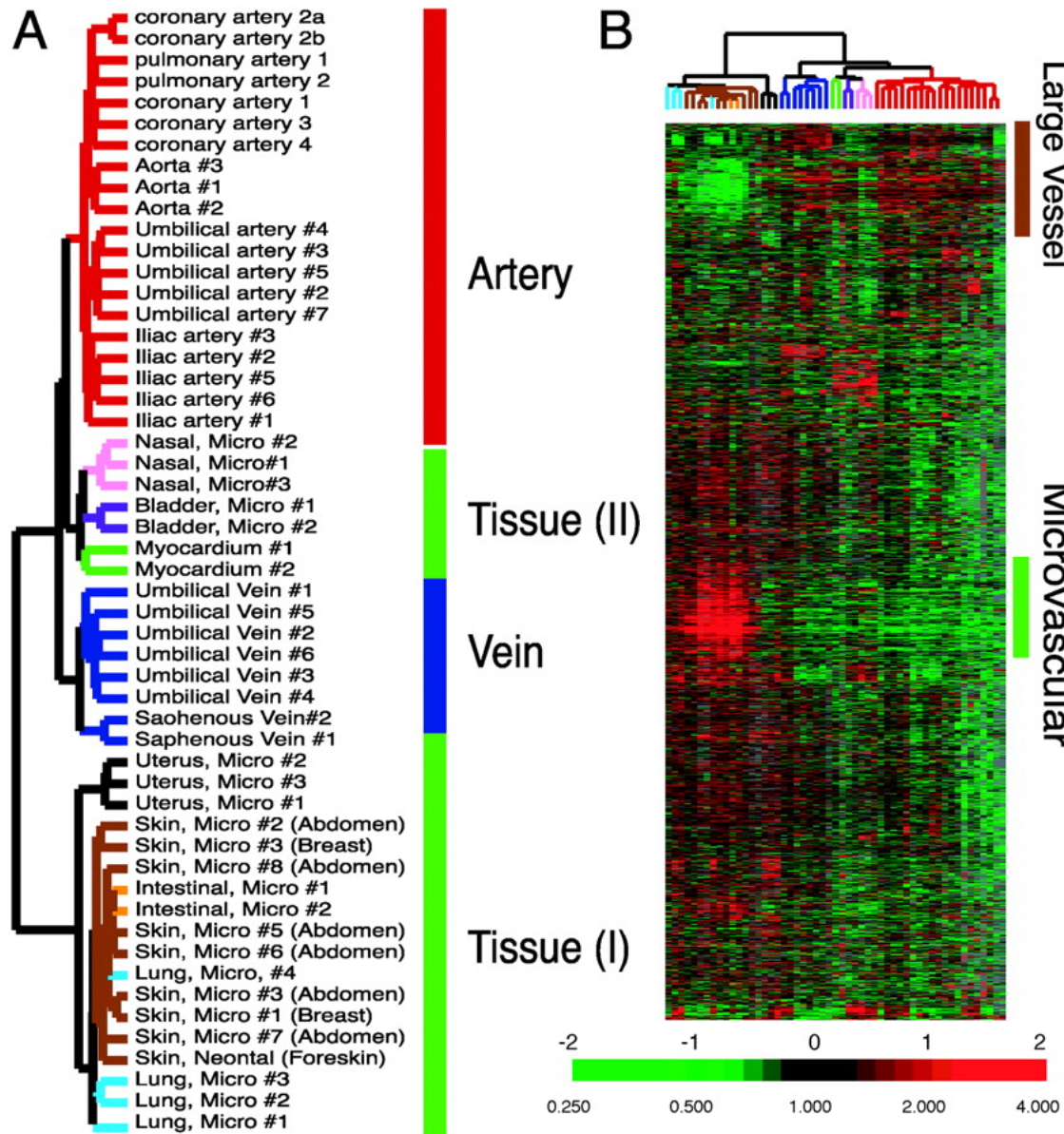
# Gene expression 3

ROC curves

Multiple testing

Gene expression networks

# Hierarchical clustering results



Chi et al., PNAS | **September 16, 2003** | vol. 100 | no. 19 | 10623-10628

“Endothelial cell diversity revealed by global expression profiling”

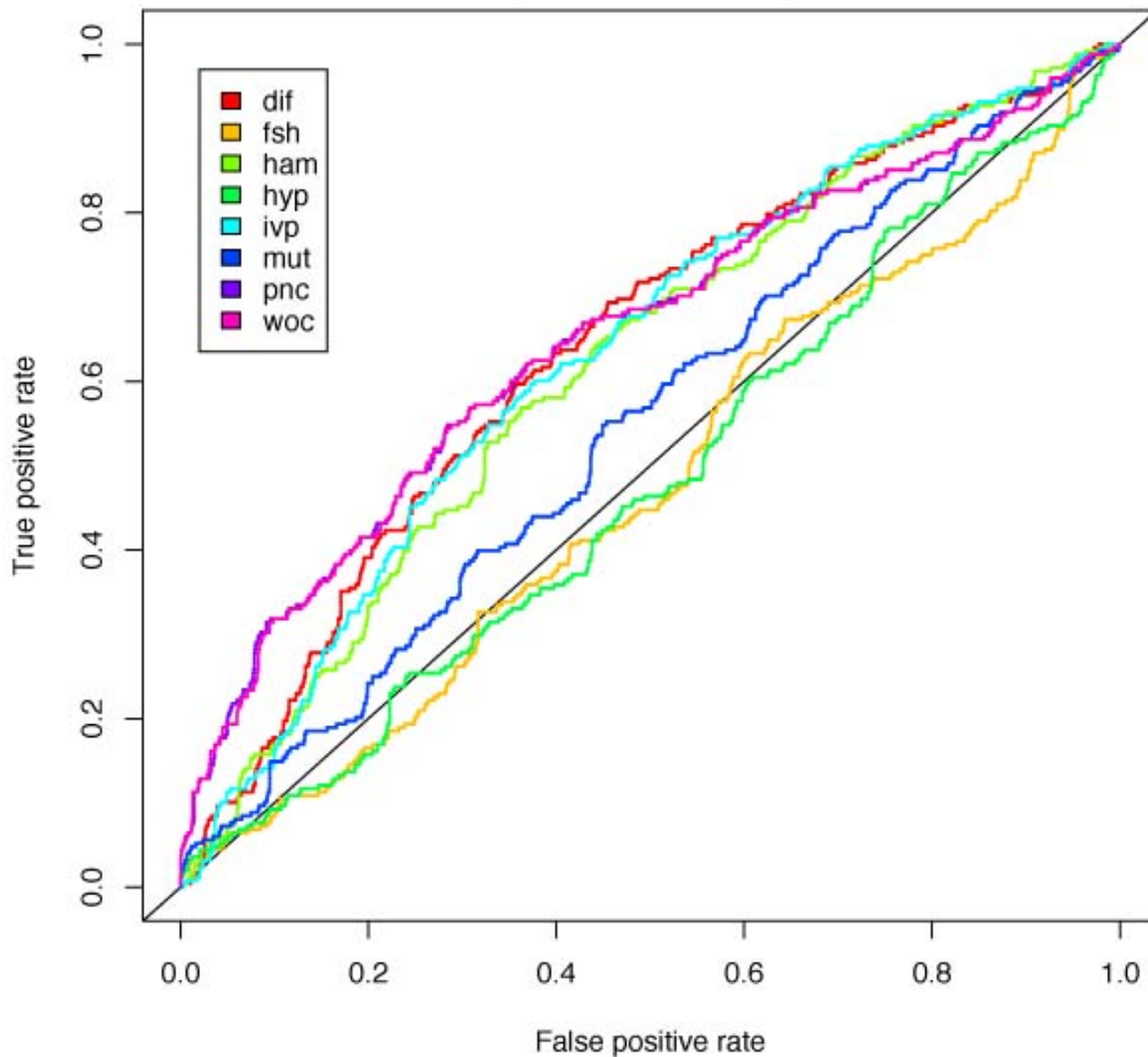
# Receiver operating characteristic

- A framework to compare the performance of binary classifiers
- Plot of *false positive rate (sensitivity)* vs *true positive rate (1-specificity)*

# Gütemaße

- Sensitivität/*Recall*  
 $TPR = TP/P = TP/(TP+FN)$
- Spezifität  
 $FPR = FP/N = FP/(FP+TN)$
- *Precision (positive predictive value)*  
 $PPV = TP/(TP+FP)$
- *False discovery rate*  
 $FDR = FP/(TP+FP)$

ROC curves with data set 'falMP'  
AUC is 0.653 for 'pnc' with rank p-value 3.40e-016  
Fisher's exact test p-value for 1st 250 genes is 4.41e-016





# Hypothesis Testing

- $H_0$  : Null hypothesis vs.  $H_1$  : Alternative Hypothesis
- $T$  : test statistics       $C$  : critical value
- If  $|T| > C$ ,  $H_0$  is rejected. Otherwise  $H_0$  is retained
- Example  
 $H_0 : \mu_1 = \mu_2$  vs.  $H_1 : \mu_1 \neq \mu_2$   
 $T = (\bar{x}_1 - \bar{x}_2) / \text{pooled standard error (se)}$
- If  $|T| > z(1 - \alpha/2)$ ,  $H_0$  is rejected at the significance level  $\alpha$
- $C_\alpha$

# Hypothesis Testing

		Hypothesis Result	
		Retained	Rejected
Truth	H0		Type I error
	H1	Type II error	

- Type I error rate = false positives ( $\alpha$  : significance level )
- Type II error rate = false negatives
- Power :  $1 - \text{Type II error rate}$
  
- P-values :  $p = \inf\{\alpha \mid H_0 \text{ is rejected at the significance level } \alpha\}$



# Issues in Multiple Comparison

- Given n treatments, which two treatments are significantly different ? (simultaneous testing)
- Is treatment A different from treatment B ?
- m treatment means :  $\mu_1, \dots, \mu_n$   
 $H_j : \mu_i = \mu_j$  where  $i \neq j$      $T_j = (x_i - x_j) / \text{pooled SE}$
- Type I error when testing each at 0.05 significance level one by one :  $1 - (0.95)^n$
- Inflated Type I error, ex)  $\alpha = 1 - (0.95)^{10} = 0.401263$
- Remedies : Bonferroni Method  
Type I error rate =  $\alpha / \#$  of comparison

# Type I Error Rates

		Hypothesis Result		Total
		#retained	#rejected	
Truth	H0	U	V	m0
	H1	T	S	m1
Total		m-R	R	m

- Per-comparison error rate ( PCER ) =  $E(V) / m$
- Per-family error rate ( PFER ) =  $E(V)$
- Family-wise error rate =  $\text{pr} ( V \geq 1 )$
- False discovery rate ( FDR ) =  $E(Q)$ ,  $Q \begin{cases} V/R, & \text{if } R > 0 \\ 0, & \text{if } R = 0 \end{cases}$

# Type I Error Rates

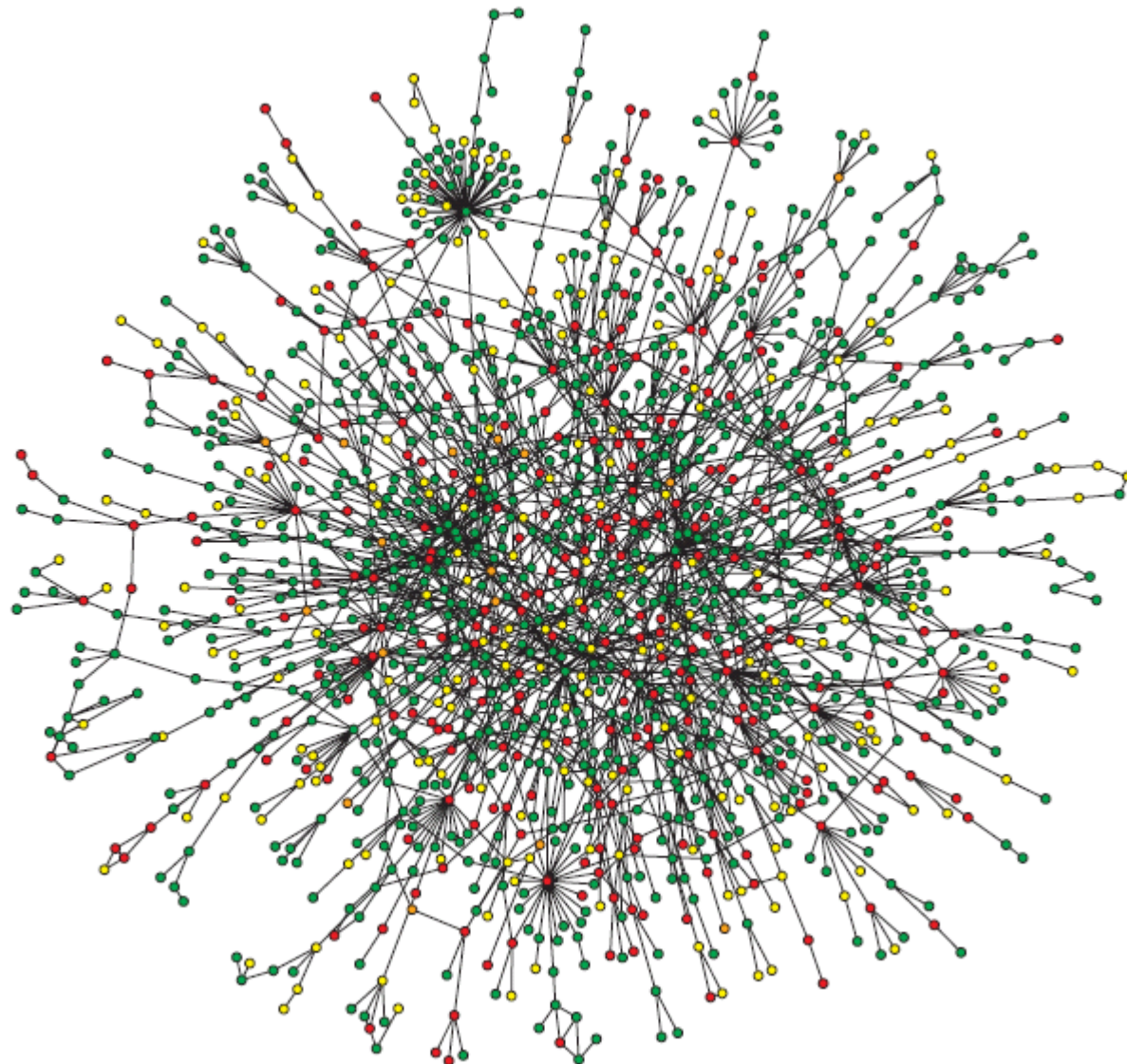
Under the complete null hypothesis, each  $H_j$  has Type I error rate  $\alpha_j$ .

- $PCER = E(V) / m = (\alpha_1 + \dots + \alpha_m) / m$
- $PFER = E(V) = \alpha_1 + \dots + \alpha_m$
- $FWER = \Pr ( V \geq 1 ) = 1 - \Pr ( H_j, j=1, \dots, m, \text{ not rejected } )$
- $FDR = E(V / R) = FWER$

$$PCER = (\alpha_1 + \dots + \alpha_m) / m \leq \max (\alpha_1 + \dots + \alpha_m) \\ \leq PWER = FDR \leq PFER = \\ \alpha_1 + \dots + \alpha_m$$

# Types of comparisons

- Assume  $H_j$ ,  $j=1, \dots, m$ , with their test statistics  $T_j$ ,  $j=1, \dots, m$ , which has a MN with mean  $\mu=(\mu_1, \dots, \mu_m)$  and identity covariance vector
- Let  $R_j = I$  ( $H_j$  is rejected) and  $r_j$  is observed value of  $R_j$
- Let  $\gamma_j = \Pr ( H_j \text{ rejected under } H_j )$ .
- $\text{PFER} = \sum_{j=1}^m \gamma_j$  (Per family error rate)
- $\text{PCER} = \sum_{j=1}^m \gamma_j / m$  (Per comparison error rate)
- $\text{FWER} = 1 - \prod_{j=1}^m (1 - \gamma_j)$  (Family wise error rate)
- $\text{FDR} = \sum_{r_1=0}^1 \dots \sum_{r_m=0}^1 (\sum_{j=1}^m r_j / \sum_{j=1}^m r_j) \prod_{j=1}^m \gamma_j r_j (1 - \gamma_j)^{1-r_j}$  (False discovery rate)



# Networks

# Considerations for the analysis

- Directed vs undirected graphs
- Analysis of confounding factors
- How to assign weights?
  - Repetitions in screen
  - Outgoing and incoming edges
  - External data
- Hubs

# Topological analysis

- **Small worlds**

- Shortest path lengths are small
- Degrees of separation

- **Modular**

- Clustering co-efficient

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

- Degree distribution

- Random model

$$P(\text{deg}(v) = k) = \binom{n-1}{k} p^k (1-p)^{n-1-k},$$

- Poisson with  $\text{max} = P(\langle k \rangle)$

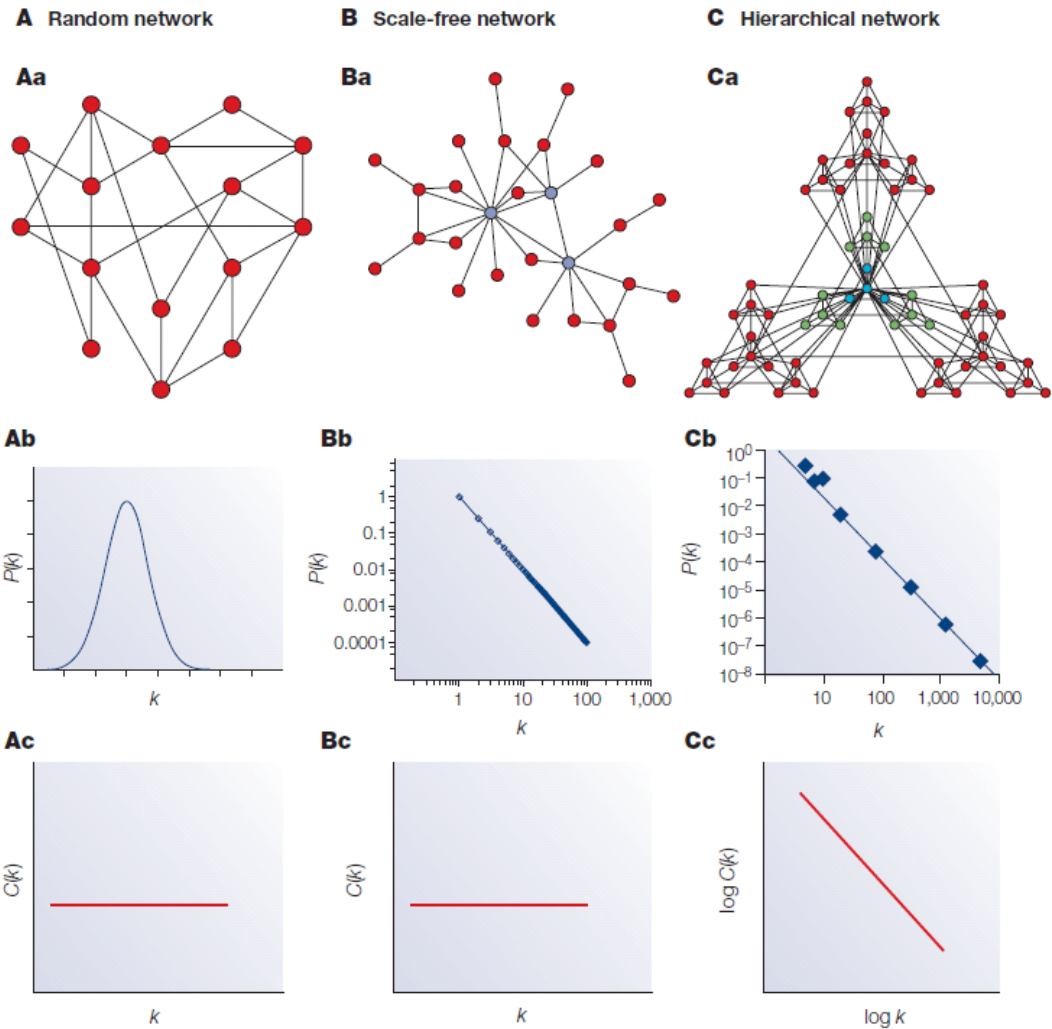
- **Scale free**

- $P(k) \sim k^{-r}$
- $-1.5 > r > -3$



# Different networks

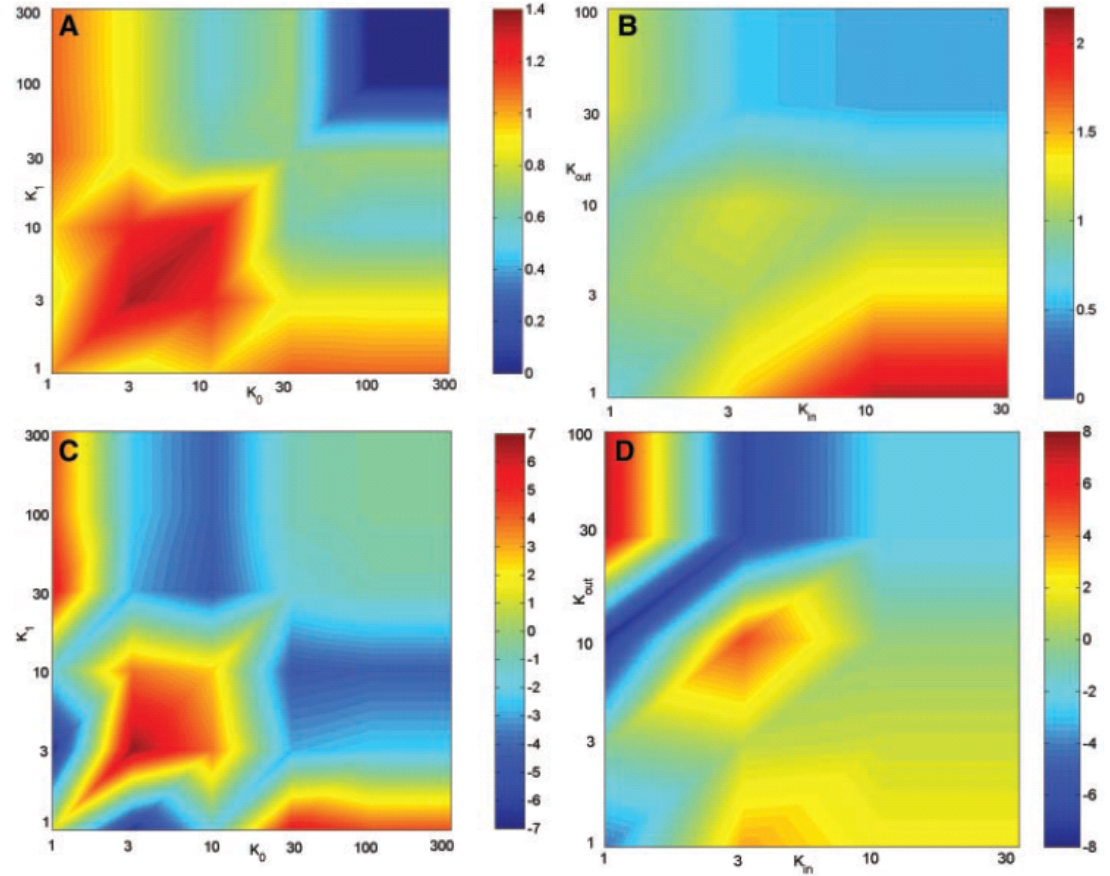
From Barabási (2004), Nature Reviews Genetics



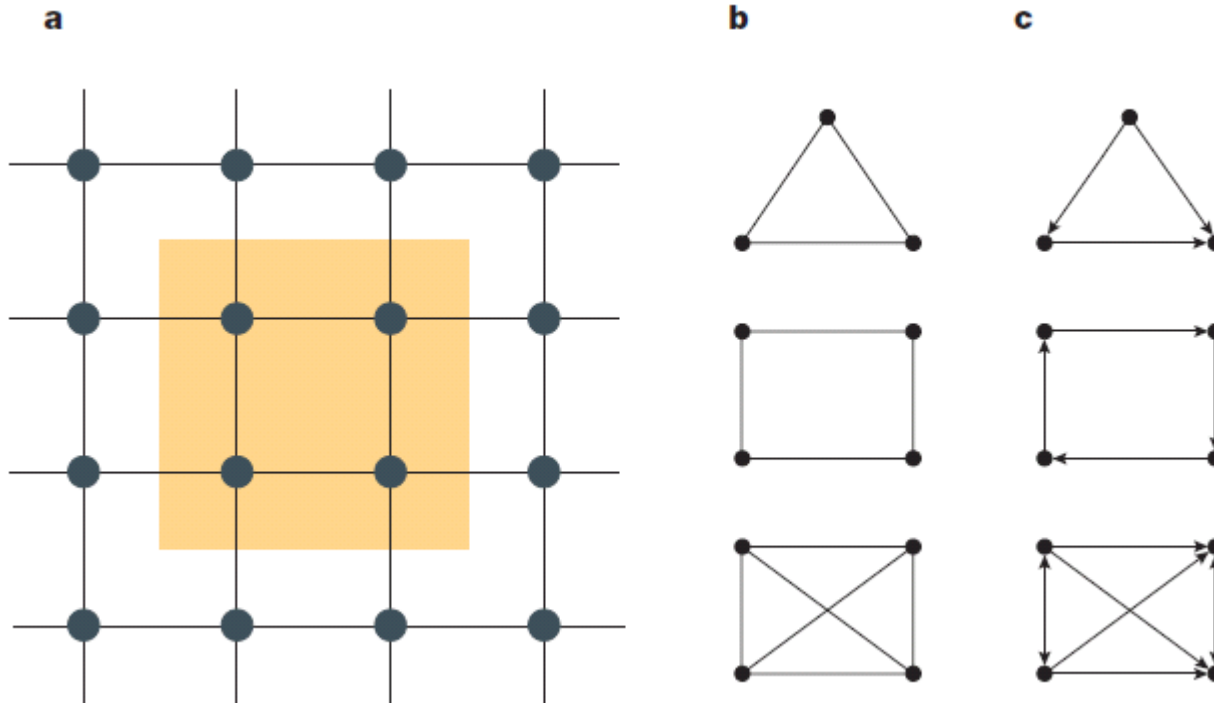
# Connections between hubs

Maslov and Sneppen (2002)  
Science

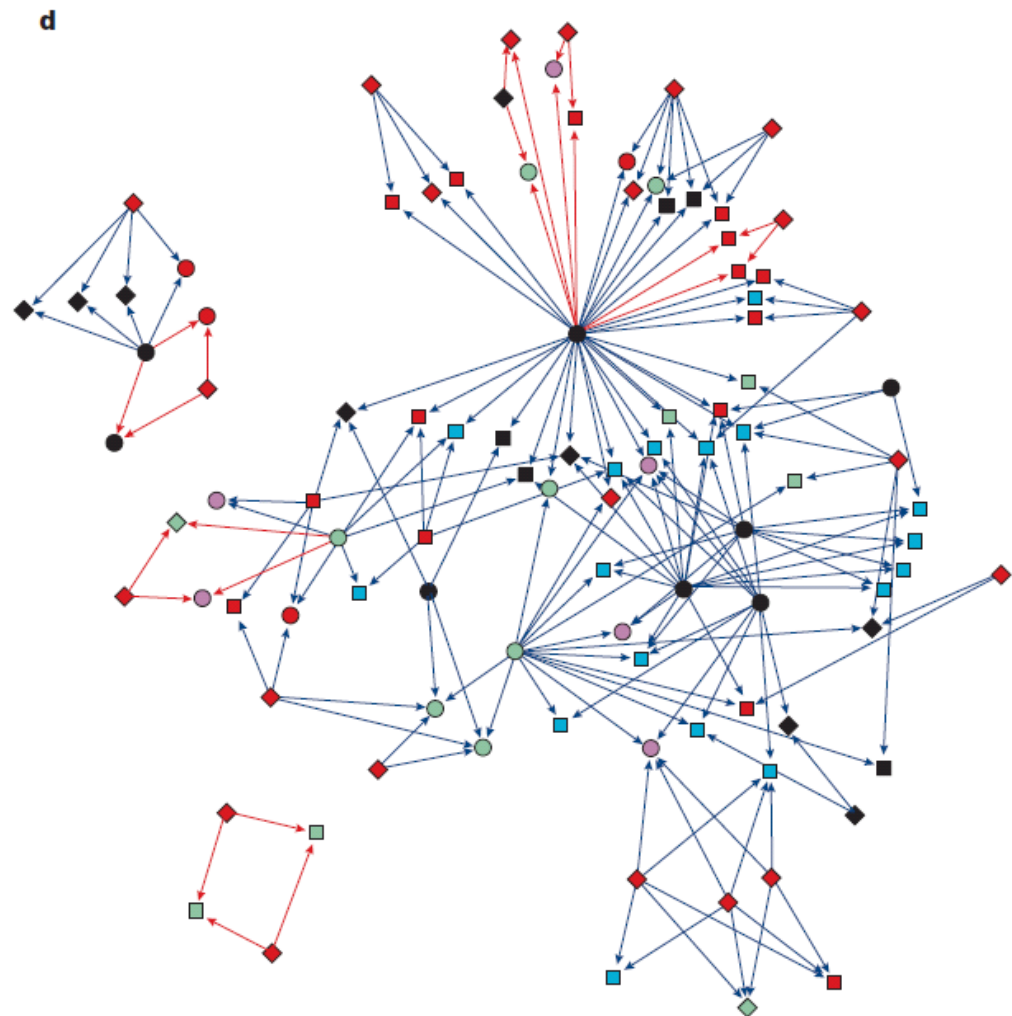
Hubs are connected to proteins of low degree, not between each other



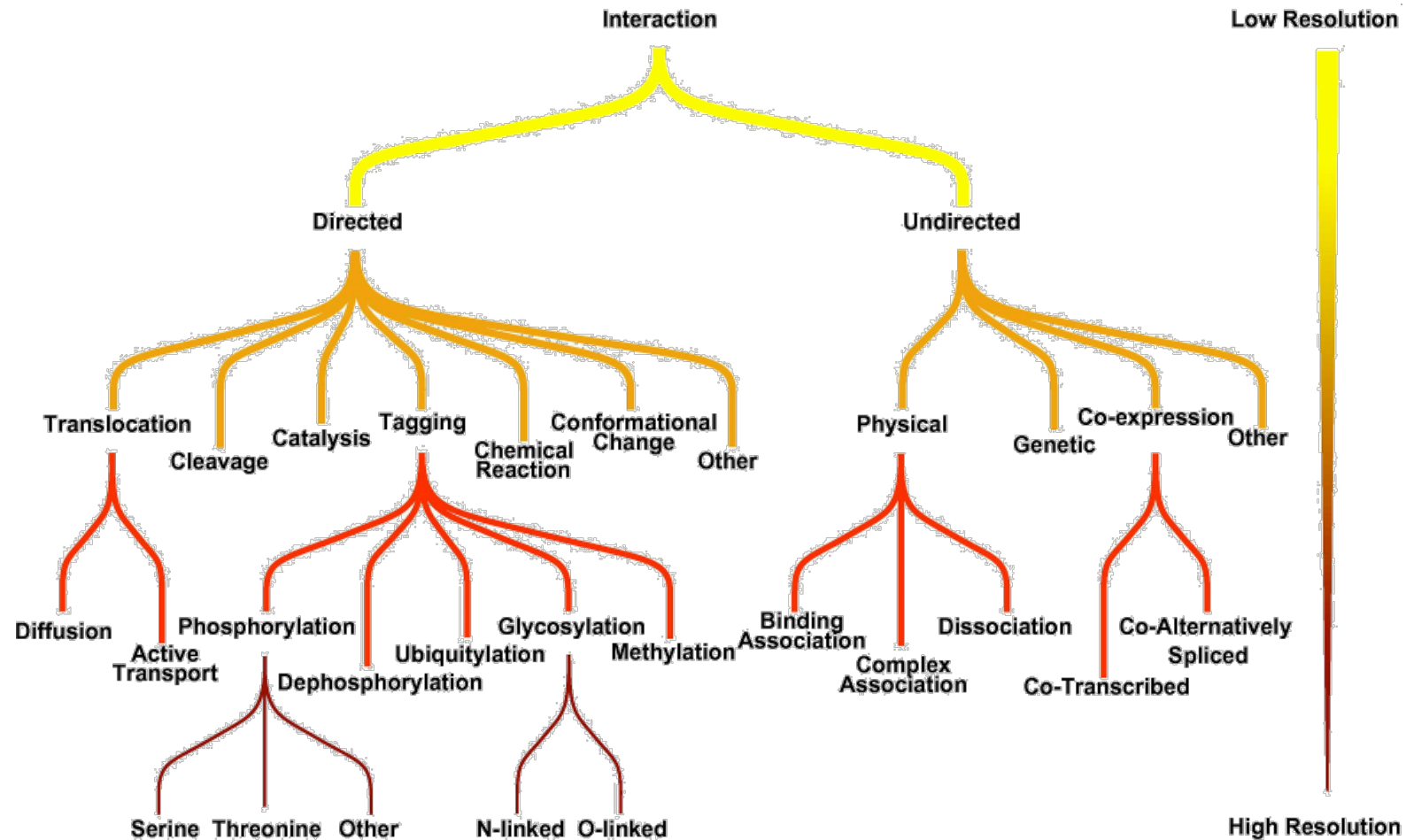
# Motifs and subgraphs



# Motifs in real networks



# Biological types of interactions



A proposed ontology for interactions (Lu et al.)

- Unweighted graphs
  - Hamming distance
- Weighted graphs
  - Euclidean distance
  - Correlation
    - Pearson
    - Spearman
- Boolean networks
- Probabilistic networks
  - Markov Random Fields
  - Bayesian networks