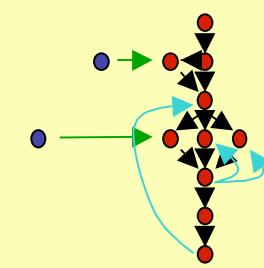


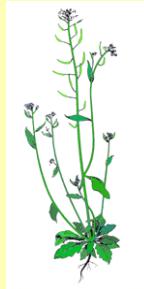
A modified graphical gaussian model approach for genetic regulatory networks

Anja Wille, SfS
Reverse Engineering Project, ETHZ

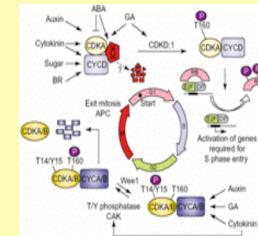
Introduction



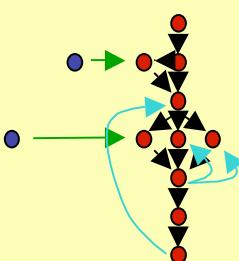
modified GGMs



biological scenarios



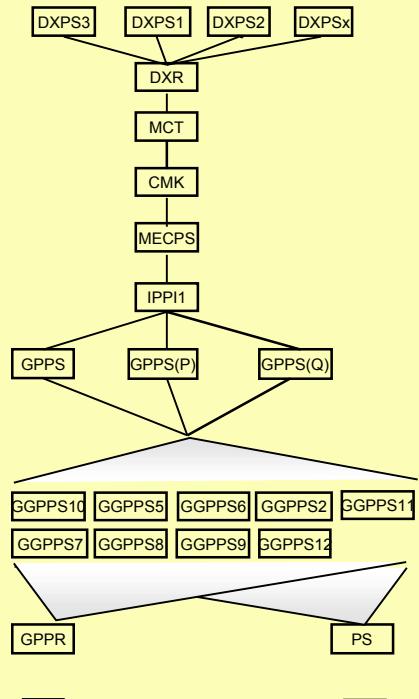
graphical gaussian models
(GGMs)



Biological scenarios

Isoprenoid pathways

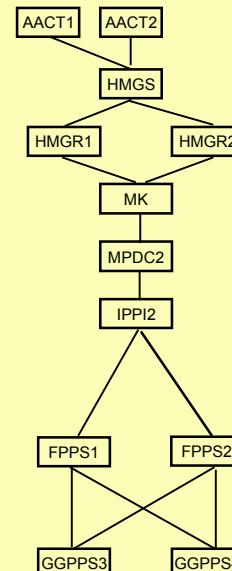
Chloroplast



Chlorophylls

Carotenoids

Cytosol

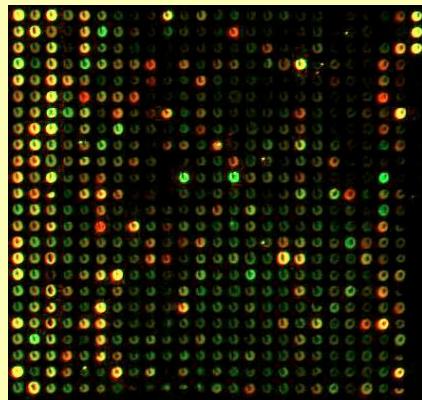


Phytosterols
Sesquiterpenes

Mitochondria



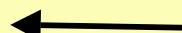
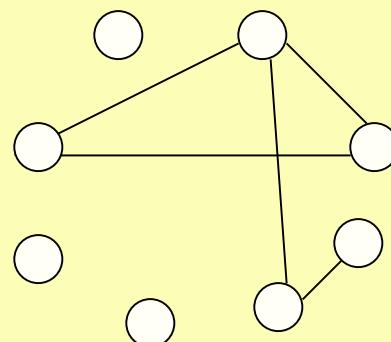
Genetic regulation



118 observations

39 genes

... signals ...



graphical model to estimate
conditional dependence
between genes

Graphical Gaussian Models

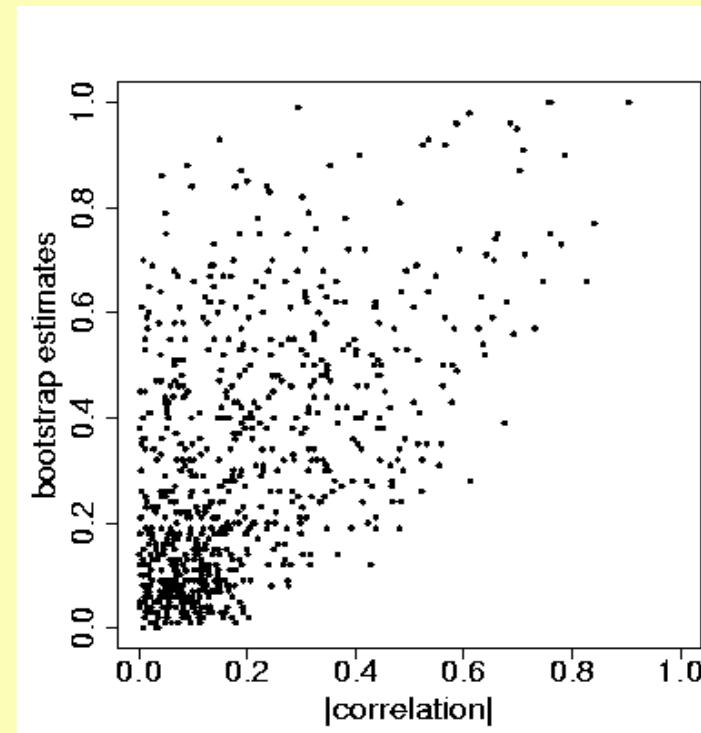
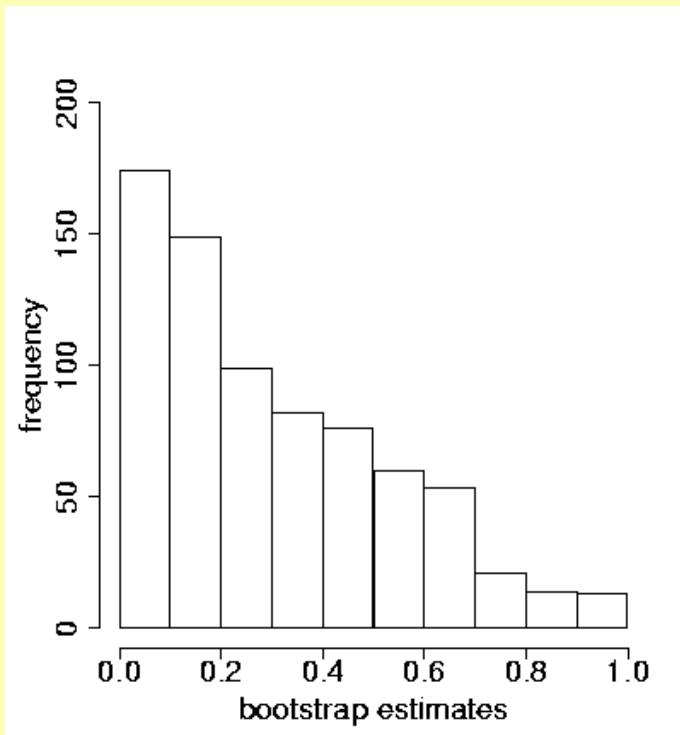
- undirected
- random variables follow a multivariate normal distribution
- log likelihood:

$$l(\boldsymbol{\Omega}, \boldsymbol{\mu}) = -\frac{N}{2} \left(q \ln(2\pi) + \ln |\boldsymbol{\Omega}| + \text{tr}(\boldsymbol{\Omega} \boldsymbol{S}) + (\bar{\mathbf{y}} \boldsymbol{\Omega} \boldsymbol{\Omega}') \boldsymbol{\Omega} (\bar{\mathbf{y}} \boldsymbol{\Omega} \boldsymbol{\Omega}) \right)$$

for a sample of N observations $\mathbf{y}^1, \dots, \mathbf{y}^N$ with
sample mean $\bar{\mathbf{y}}$, sample covariance matrix \boldsymbol{S}
and precision matrix $\boldsymbol{\Omega} = \boldsymbol{S}^{-1}$

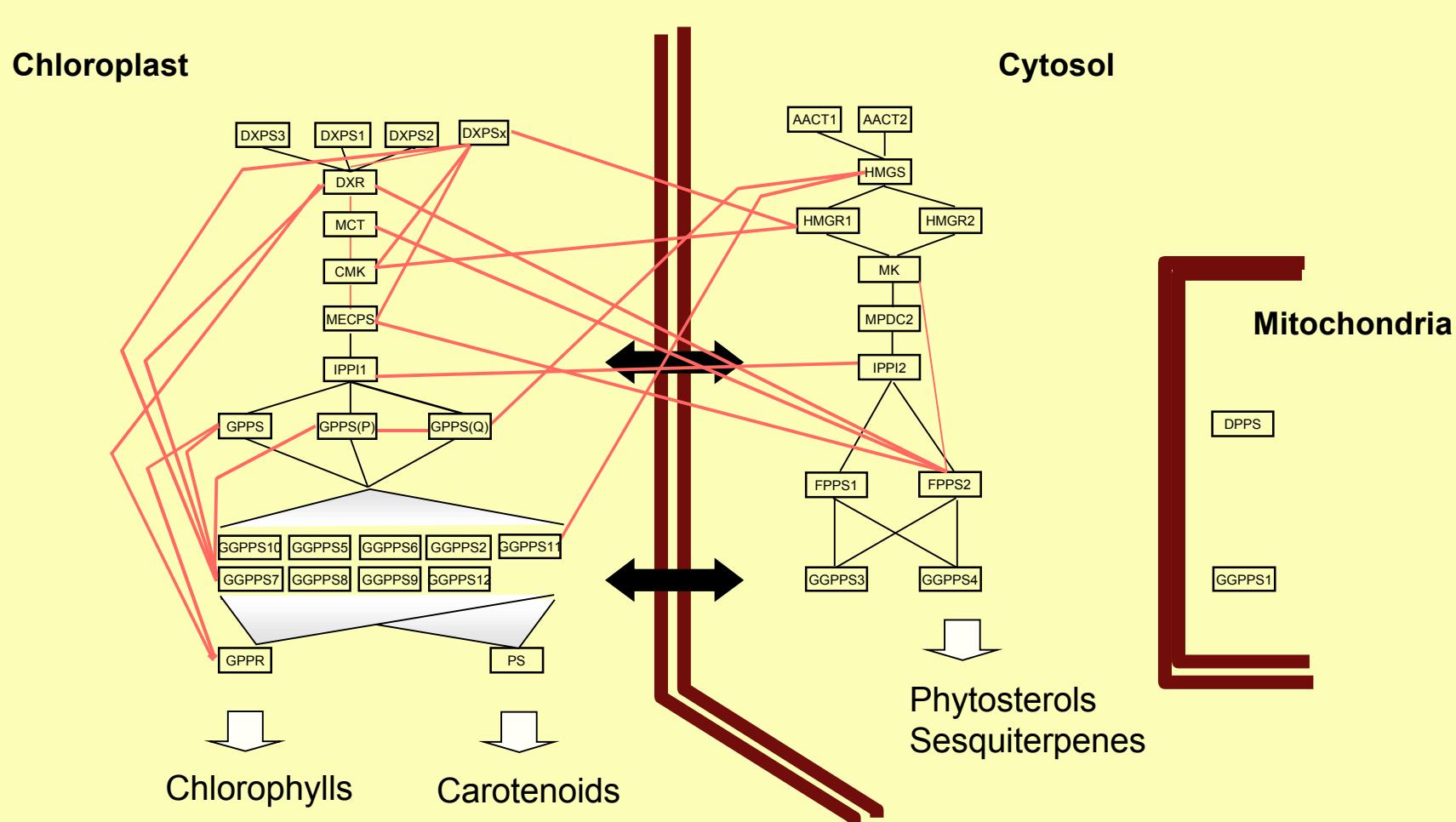
- partial correlation coefficient $\boldsymbol{\Omega}_{ij|\text{rest}}$ for gene pair ij

Bootstrapping and pairwise correlation



Application to isoprenoid pathways

Isoprenoid pathways



Problems

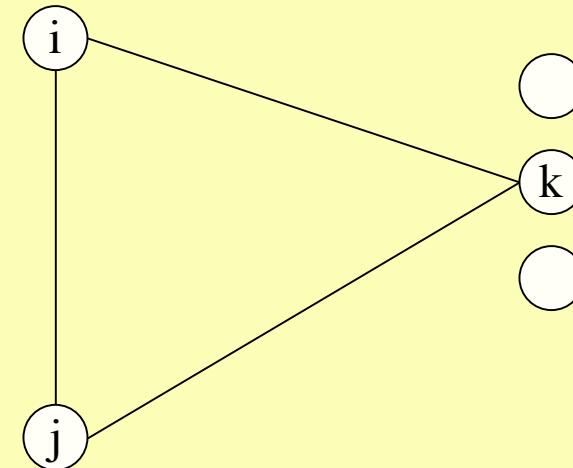
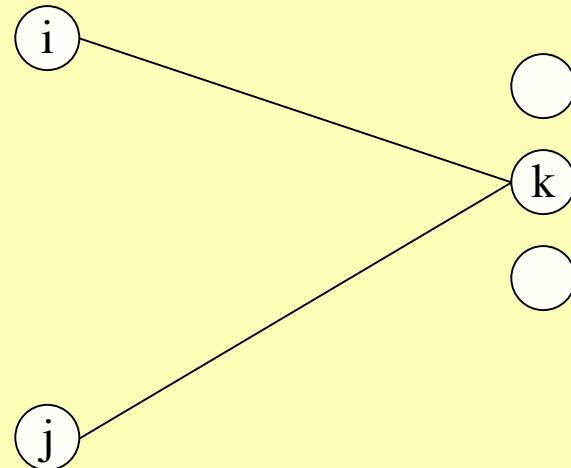
- matrix inversion rank-sensitive
- genes ↑ □ exponentially increasing number of models
 - difficult to interpret
- how to interpret high partial correlation accompanied with low pairwise correlation
- efficient procedure for attaching new genes needed

Outline for modified GGM approach

For each pair of genes i,j :

- take pairwise correlation into account
- fit GGMs with gene triples i,j,k for all remaining genes k to study the partial correlation
- combine GGMs and pairwise correlation for inference on edge ij
- 3 methods that differ with respect to statistical framework and computational costs

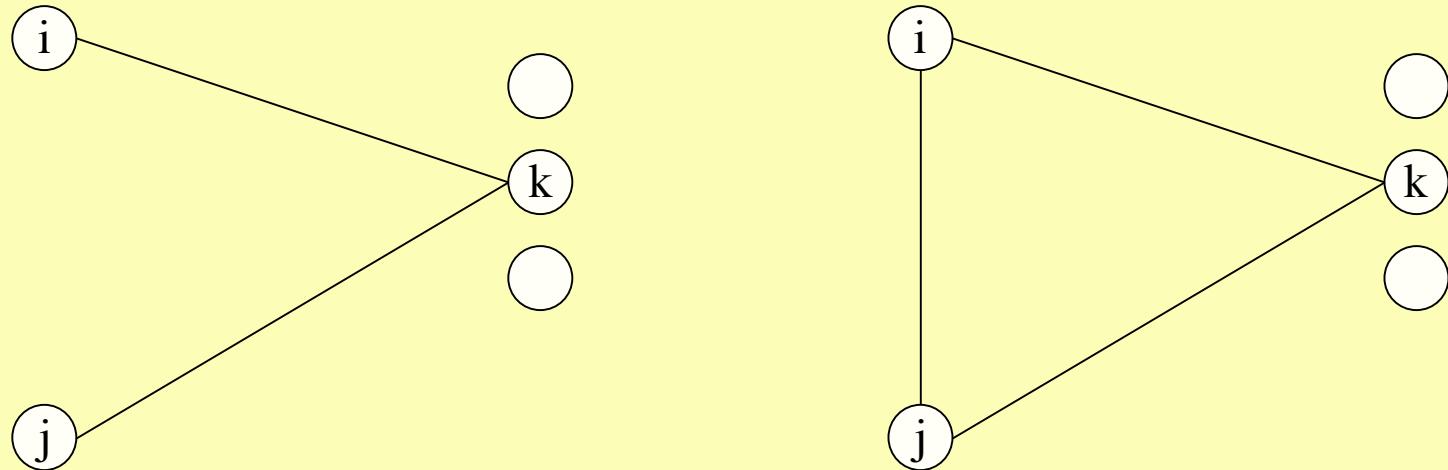
Frequentist approach



Focus on genepair ij

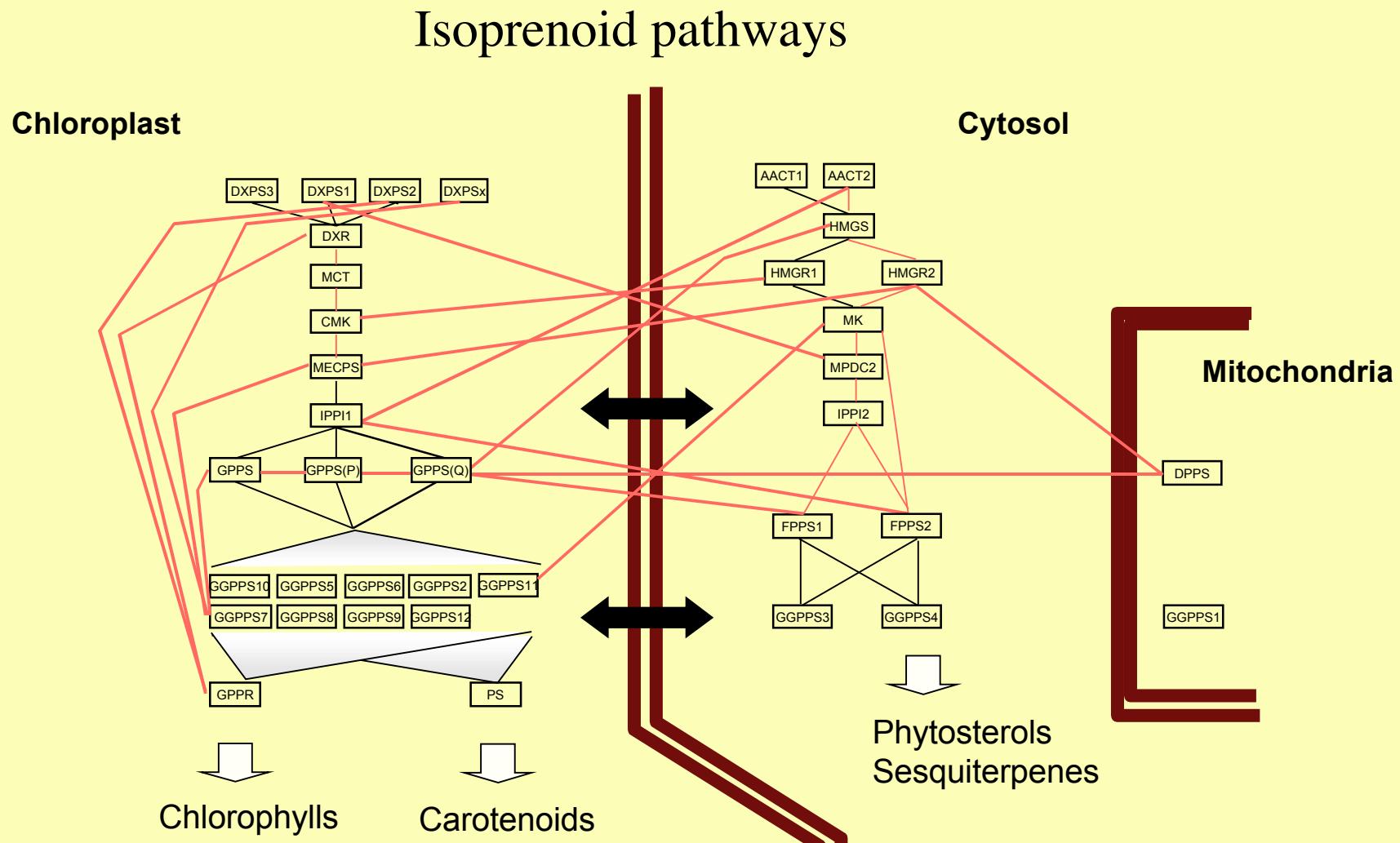
- $p_{ij|lk}$ is p-value from deviance test $\square_{ij|lk} \neq 0$ versus $\square_{ij|lk} = 0$

Frequentist approach (cont'd)

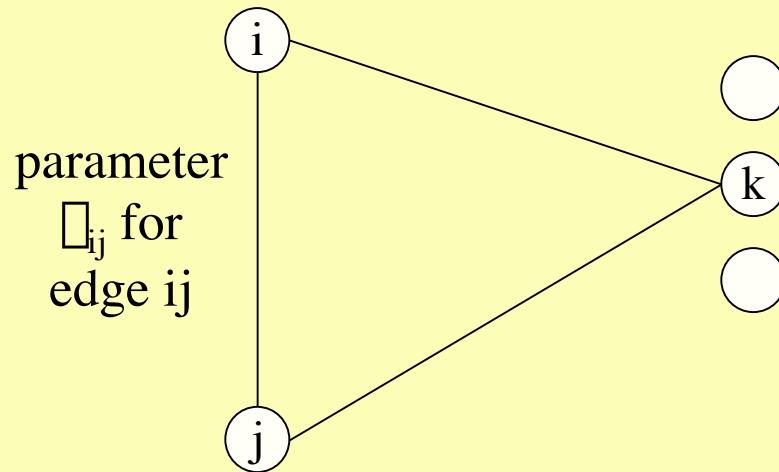


- 1) Form $p_{ij,\max} = \max(p_{ij|k} \text{ for all genes } k \neq i,j)$
- 2) Adjust $p_{ij,\max}$ according to Bonferroni-Holm or FDR
- 3) If the adjusted value $p_{ij,\max} < 0.05$, draw edge between i and j

Application to isoprenoid pathway



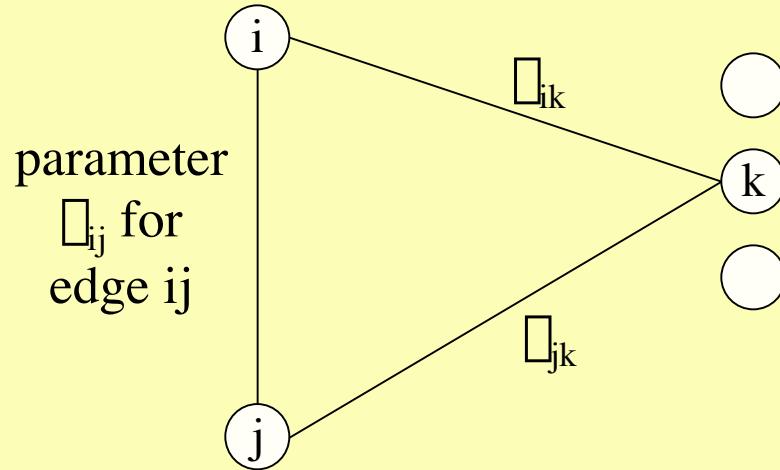
Likelihood approach with parameters θ



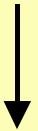
Estimate $\theta = \{\theta_{ij} \text{ for all } i,j\}$ in a maximum likelihood approach

$$L(\theta) = \prod_{g \in G} L(\theta | g) P(g)$$

EM algorithm



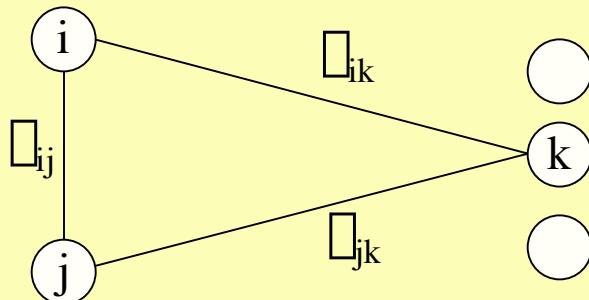
$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}') = \prod_{g \in G} L(\boldsymbol{\theta} | g) P(g | \boldsymbol{\theta}', y) \text{ to be maximized}$$



$$\alpha_{ij}^{t+1} = \frac{\prod_{g | g_{ij}=1} \alpha_{ik}^{t g_{ik}} (1 - \alpha_{ik}^t)^{1-g_{ik}} \alpha_{jk}^{t g_{jk}} (1 - \alpha_{jk}^t)^{1-g_{jk}} \cdot L(\boldsymbol{\theta}_g, \boldsymbol{\theta}_g)}{\prod_{g | g_{ij}=0} \alpha_{ik}^{t g_{ik}} (1 - \alpha_{ik}^t)^{1-g_{ik}} \alpha_{jk}^{t g_{jk}} (1 - \alpha_{jk}^t)^{1-g_{jk}} \cdot L(\boldsymbol{\theta}_g, \boldsymbol{\theta}_g)}$$

Simplification

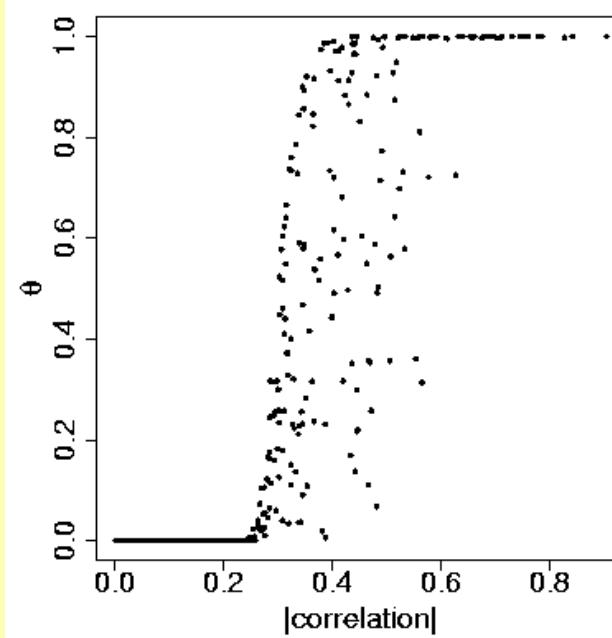
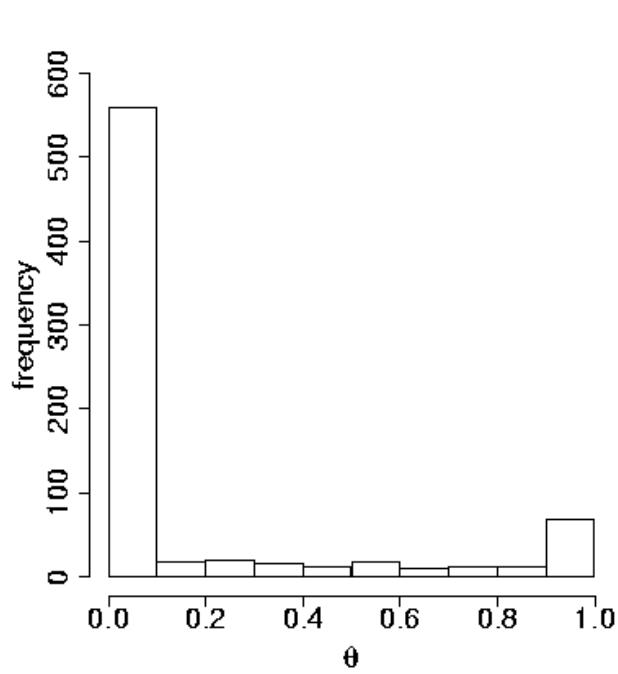
$$\square_{ij}^{t+1} = \prod_{k \neq i,j} \frac{\prod_{g|g_{ij}=1} \square_{ik}^t g_{ik} (1 - \square_{ik}^t)^{1-g_{ik}} \square_{jk}^t g_{jk} (1 - \square_{jk}^t)^{1-g_{jk}} \cdot L(\square_g, \square_g)}{\prod_g \square_{ik}^t g_{ik} (1 - \square_{ik}^t)^{1-g_{ik}} \square_{jk}^t g_{jk} (1 - \square_{jk}^t)^{1-g_{jk}} \cdot L(\square_g, \square_g)}$$



not all GGMs
with i, j, k
considered

$$\square_{ij}^{t+1} = \prod_{k \neq i,j} \square_{ik}^t \square_{jk}^t \cdot P(\square_{ijk} > 0 | y) + (1 - \square_{ik}^t \square_{jk}^t) \cdot P(\square_{ij} > 0 | y)$$

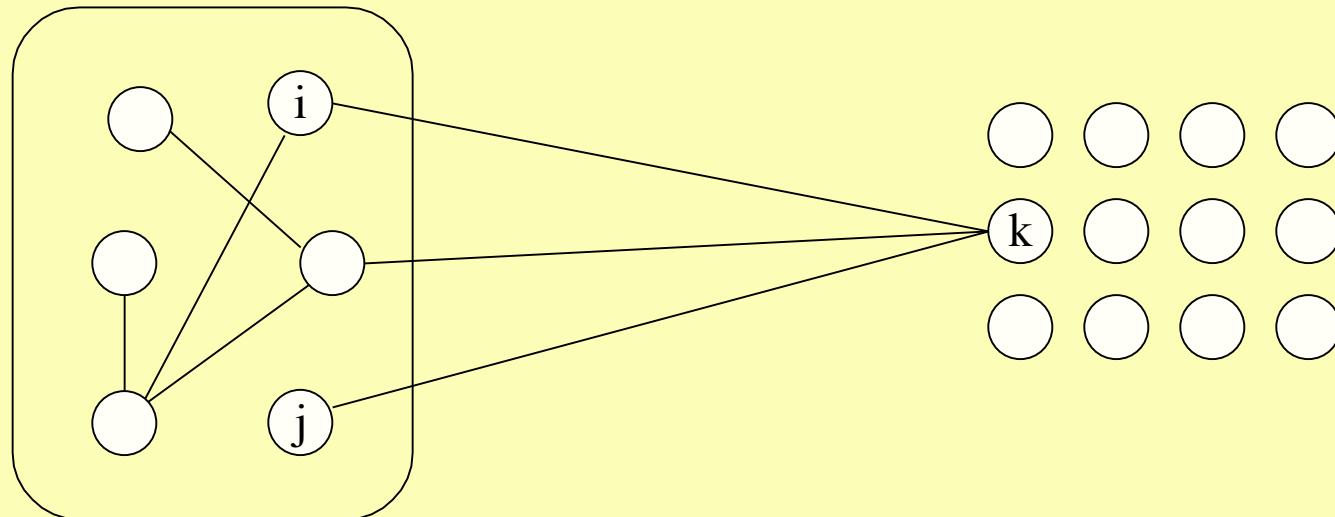
Distribution of \square_{ij}



Conclusions

- GGM of gene triples used to look whether correlation between two genes can be “explained” by a third one
- frequentist approach simple, can be applied to many genes
- approach with \Box -parameters requires iteration, tested for up to 70 genes
- a large set of additional genes can be attached to constructed network

Modeling at two levels



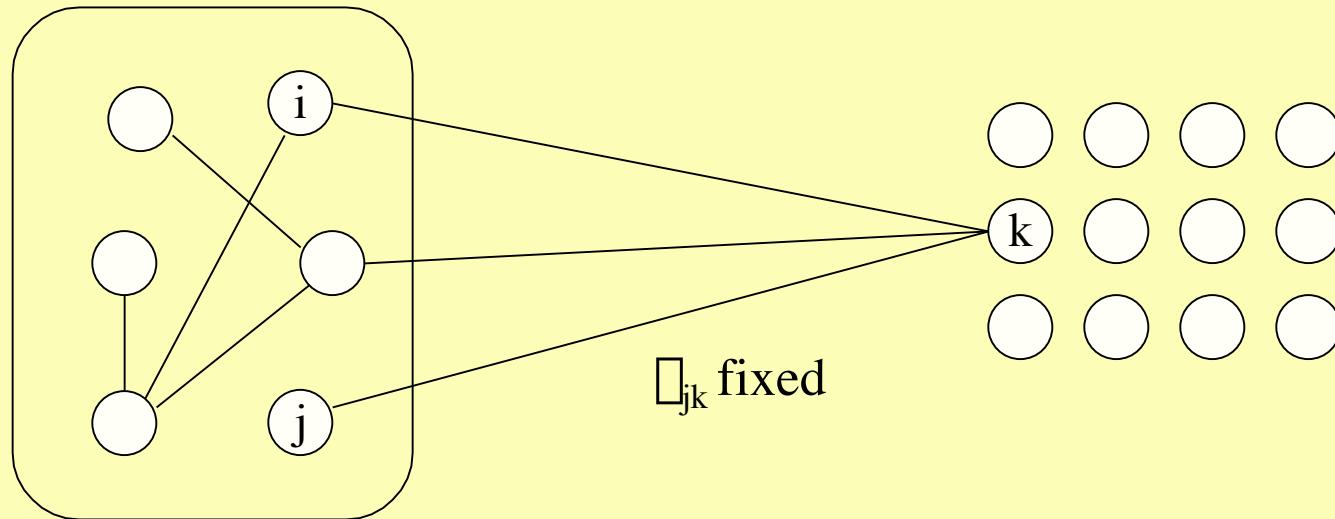
genetic network

- small number <100
- model edges

attach additional genes

- possibly 1000s
- which one “explain” edges?

Attaching additional genes

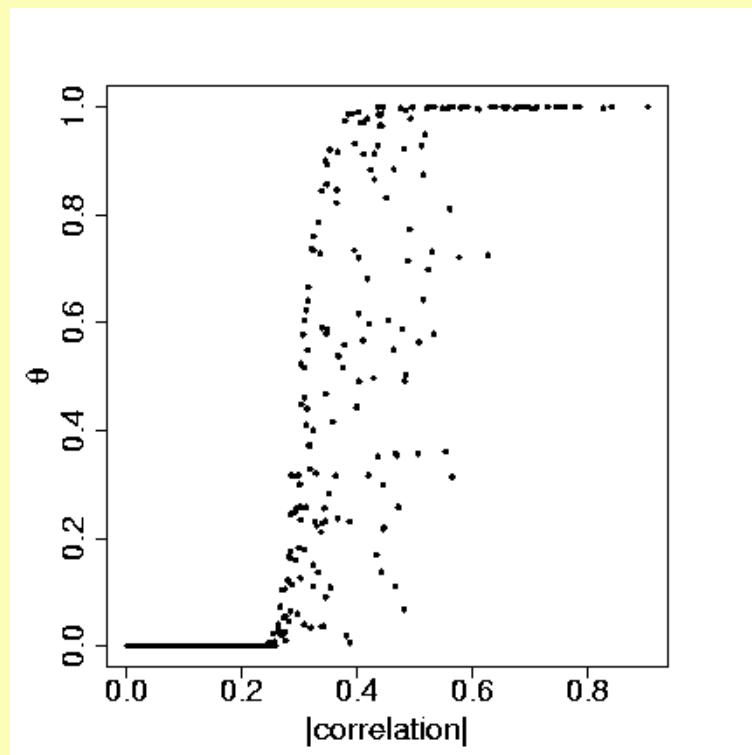


For additional genes k

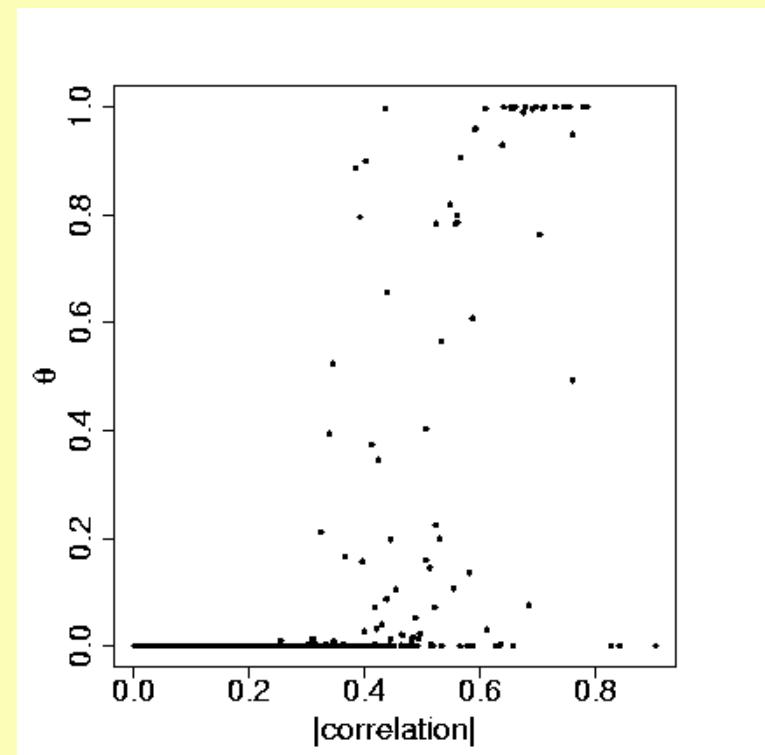
- include in computation of \square_{ij} but keep \square_{ik} and \square_{jk} fixed
- \square_{ij} decreases
- count how often k decreases \square_{ij} , validate

Attaching genes from other pathways

without additional genes



with additional genes

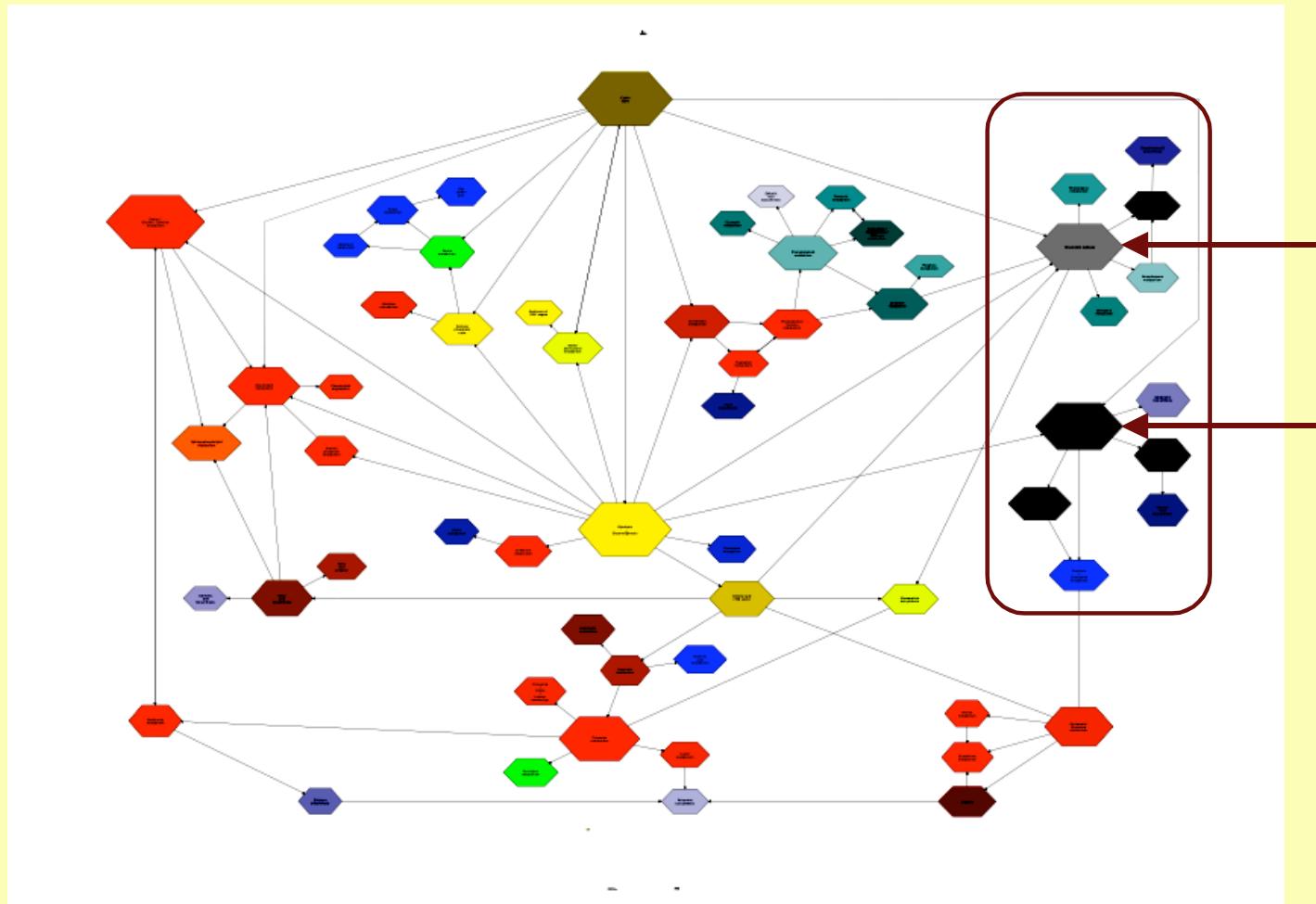


Attaching genes from other pathways

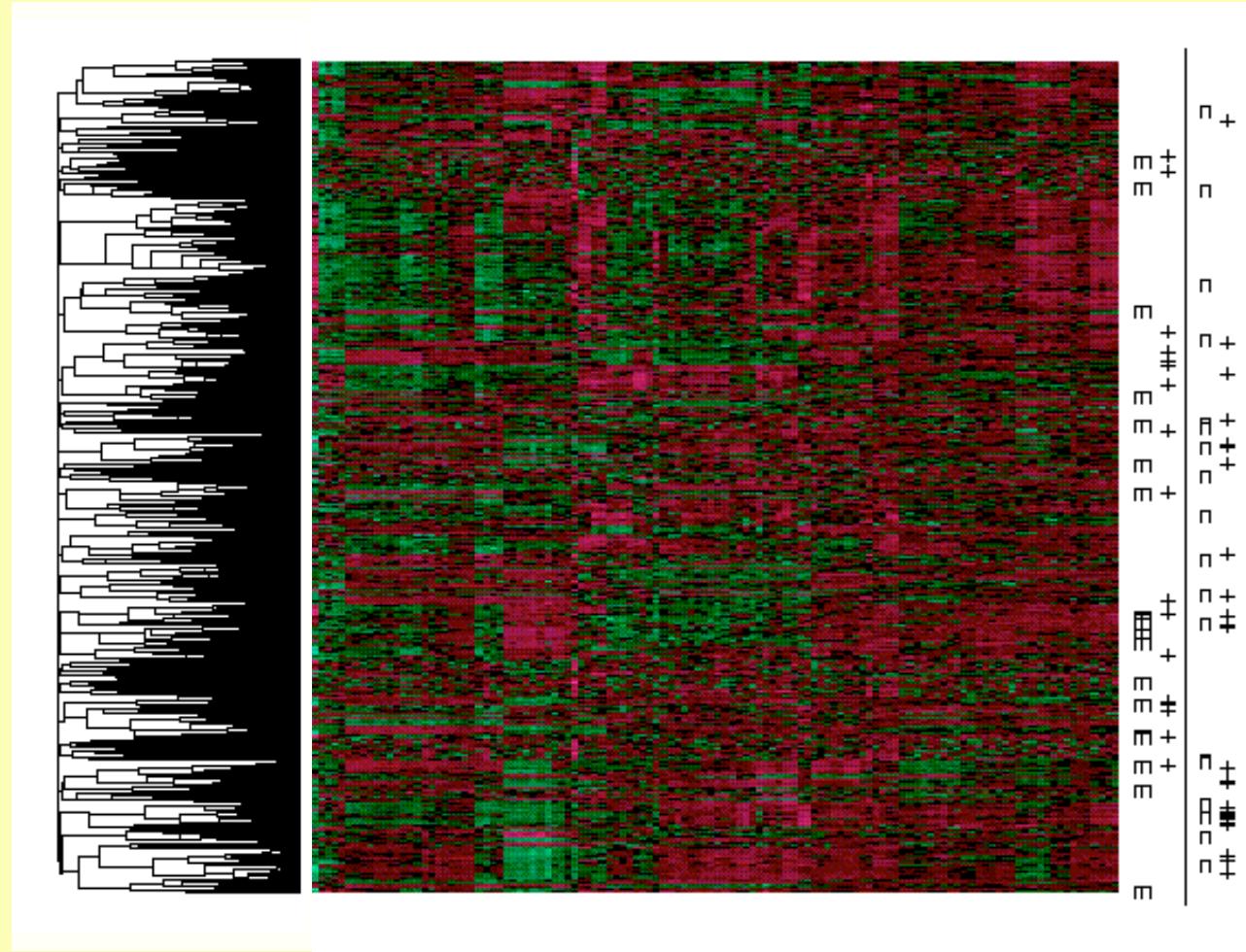
both pathways	chloroplast	cytoplasm
carotenoid* tocopherol* calvin cycle chlorophyll* abscisic acid* onecarbonpool phytosterol*	carotenoid* onecarbonpool chlorophyll* calvin cycle	tocopherol* phytosterol*

downstream pathways marked by *

Attaching genes from other pathways



Hierarchical clustering

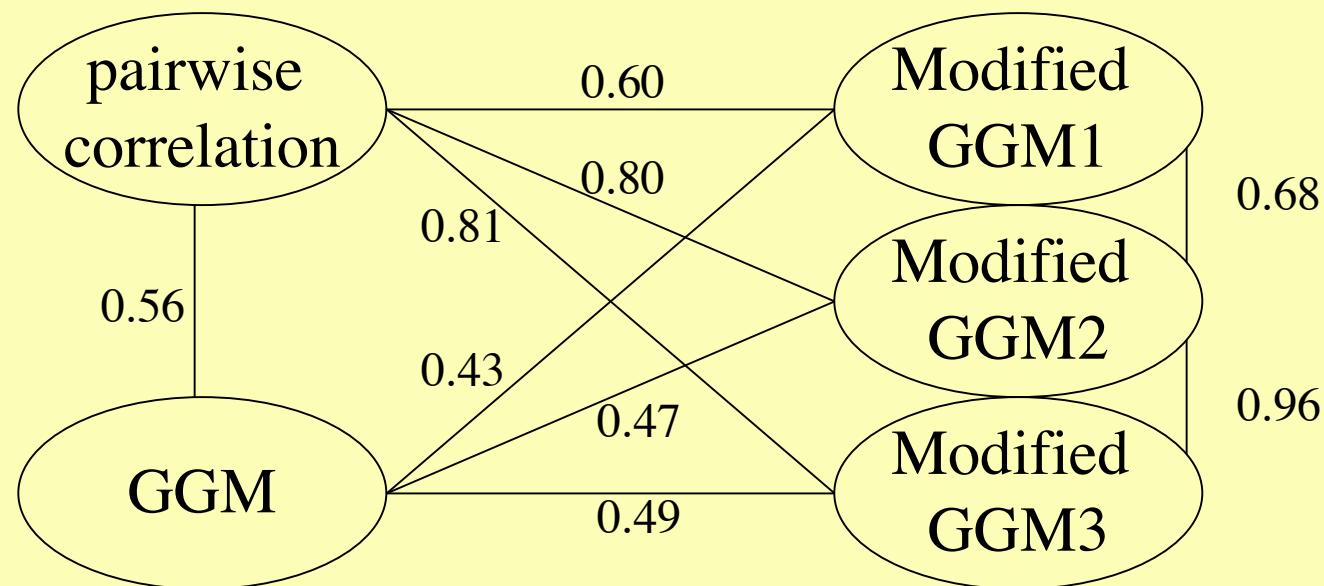


Conclusions

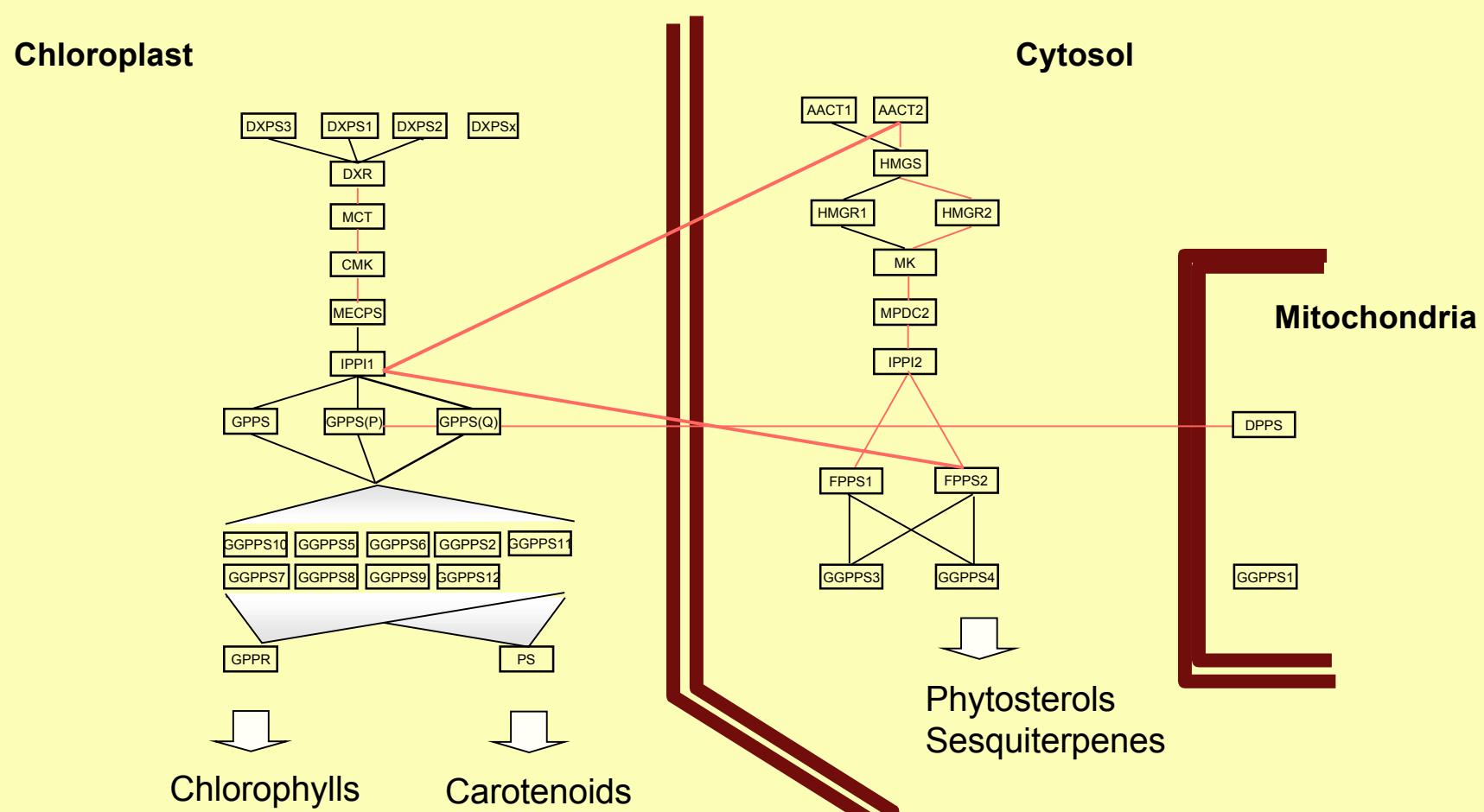
Modified GGMs

- model dependence between genes
- combine GGMs and pairwise correlation for inference on edge ij
- different statistical design, computational cost
- additional genes can be fitted into the model
- similarities in expression patterns between groups of genes can be identified (also verified in yeast data)

Comparison of different methods



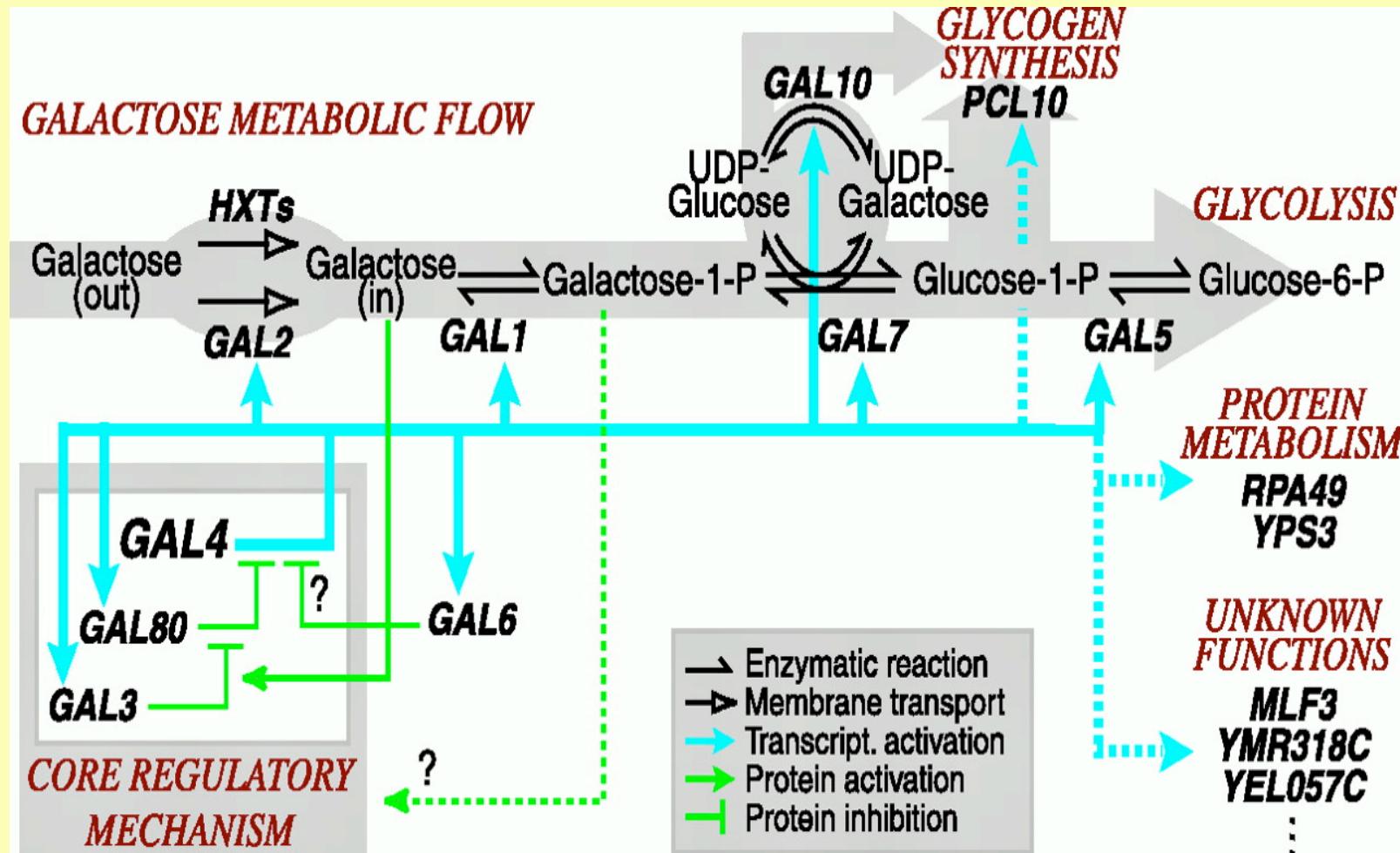
Consistent results



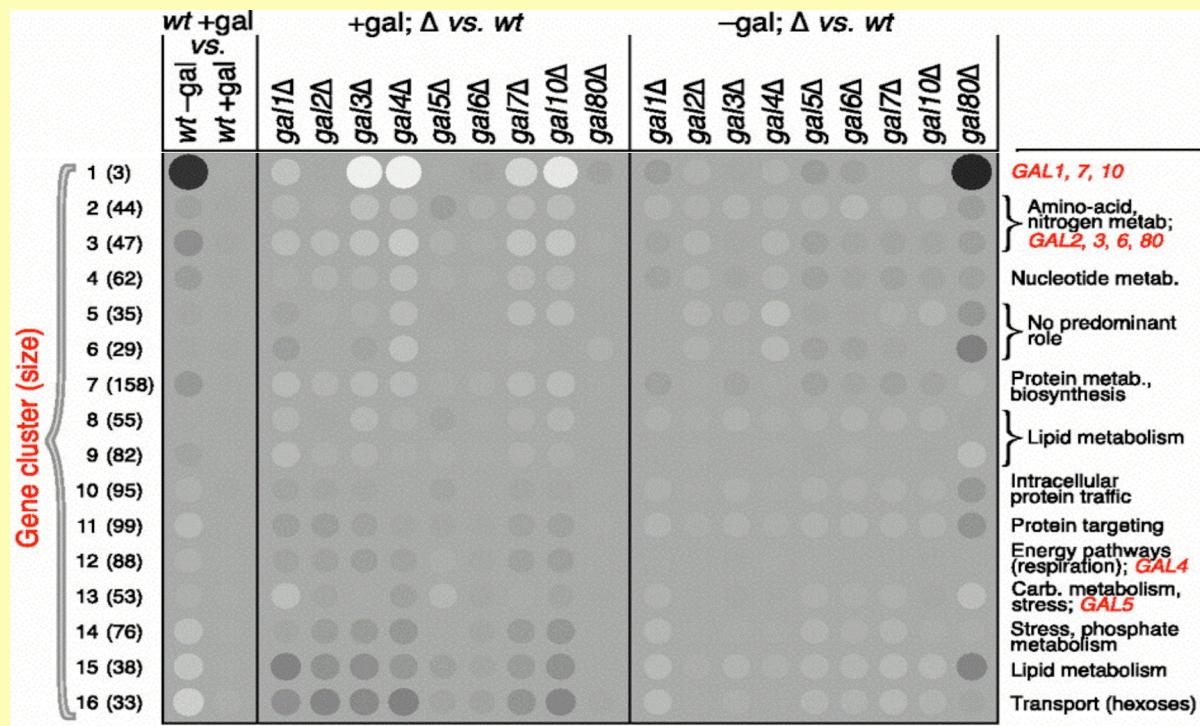
Acknowledgments

- Peter Buehlmann
- Stefan Bleuler, Amela Prelic, Eckhard Zitzler
- Philip Zimmermann, Lars Hennig, Willi Gruissem

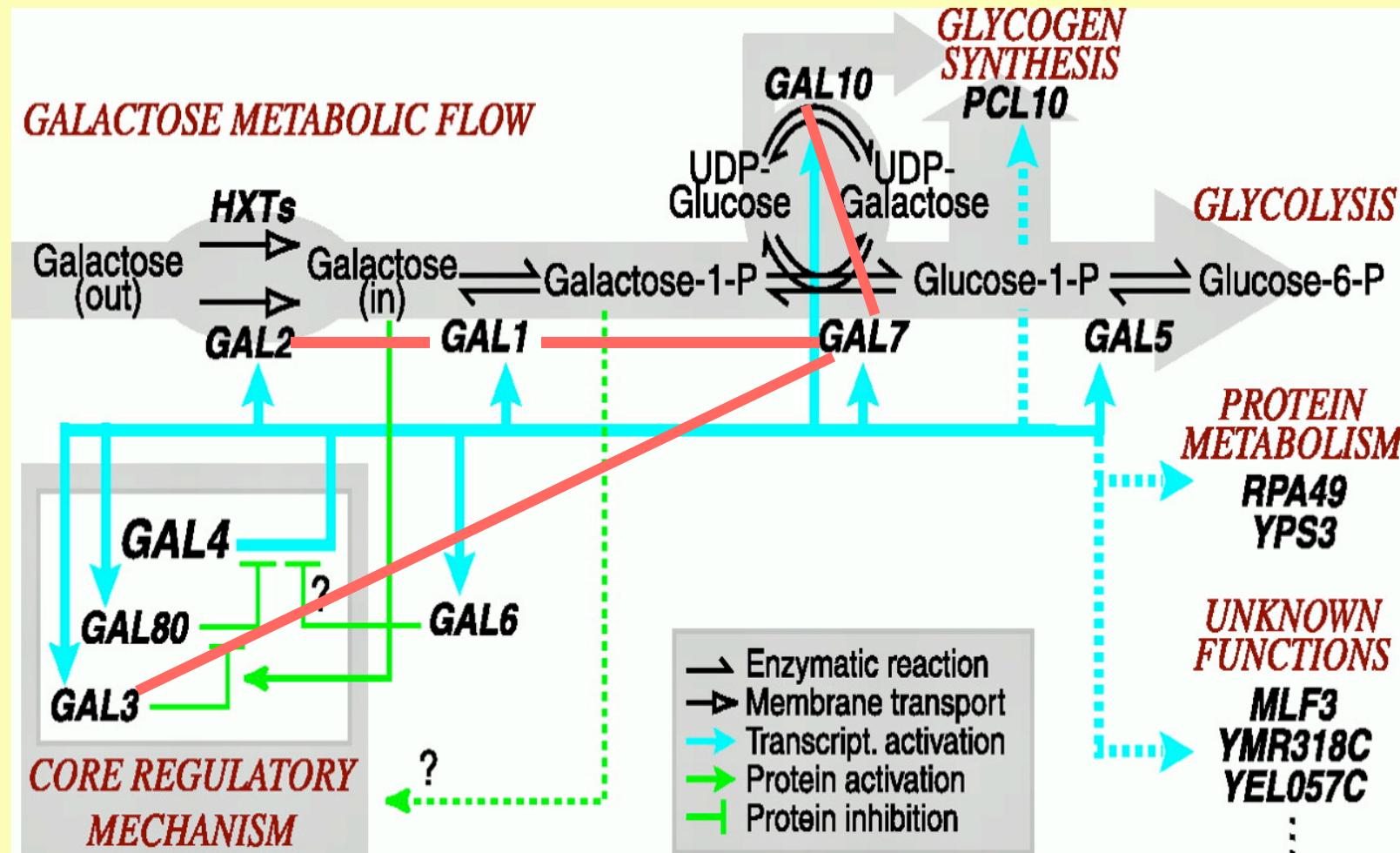
Galactose pathway in yeast



Galactose pathway in yeast



Network for galactose pathway



Network for galactose pathway

Genes that attached to network

GCY1 ←
FAR1
YDR010C
YEL057C ←
YPL066W
YOR121C
MLF3 ←

YLL058W
YJL212C
PCL10 ←
LYS1
YBR139W
YCR059C

$$l(\theta, \theta) = \frac{N}{2} (q \ln(2\pi) + \ln |\theta| + tr(\theta S) + (\bar{y} \theta \theta)' \theta (\bar{y} \theta \theta))$$

$$L(\theta) = \prod_{g \in G} L(\theta | g) P(g)$$

$$Q(\theta | \theta') = \prod_{g \in G} L(\theta | g) P(g | \theta', y)$$

$$\theta_{ij}^{t+1} = \prod_{k \neq i,j} \theta_{ik}^t \theta_{jk}^t \cdot P(\theta_{ijk} > 0 | y) + (1 - \theta_{ik}^t \theta_{jk}^t) \cdot P(\theta_{ij} > 0 | y)$$

$$\theta_{ij}^{t+1} = \prod_{k \neq i,j} \frac{\prod_{g | g_{ij}=1} \theta_{ik}^t \theta_{jk}^t (1 - \theta_{ik}^t)^{1-g_{ik}} \theta_{jk}^t (1 - \theta_{jk}^t)^{1-g_{jk}} \cdot L(\theta_g, \theta_g)}{\prod_g \theta_{ik}^t (1 - \theta_{ik}^t)^{1-g_{ik}} \theta_{jk}^t (1 - \theta_{jk}^t)^{1-g_{jk}} \cdot L(\theta_g, \theta_g)}$$