

Extracting metabolic pathways activity from gene expression data

Jean-Philippe Vert
Ecole des Mines de Paris, France

Minoru Kanehisa
Kyoto University, Japan

ECCB 2003, Paris, September 27, 2003.

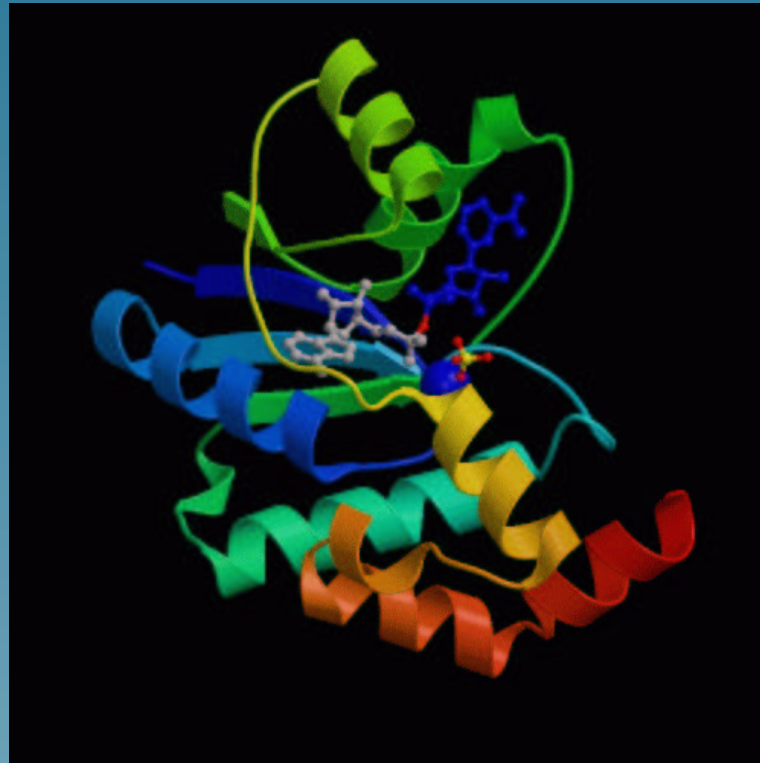
Overview

1. The problem
2. Using expression data only
3. Using a pathway database
4. Combining expression and pathways
5. Experiments

Part 1

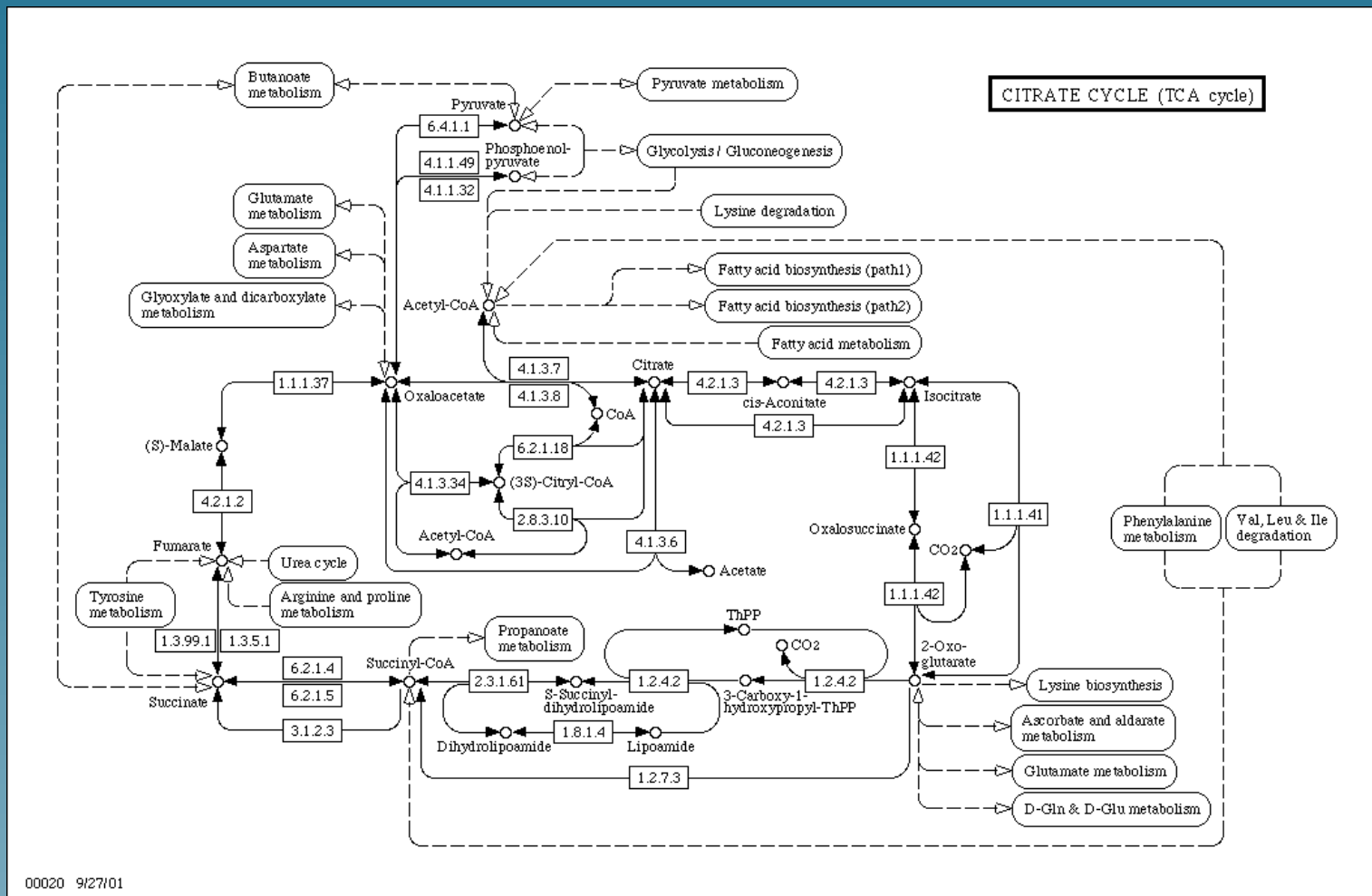
The problem

Genes encode proteins which can catalyse chemical reactions

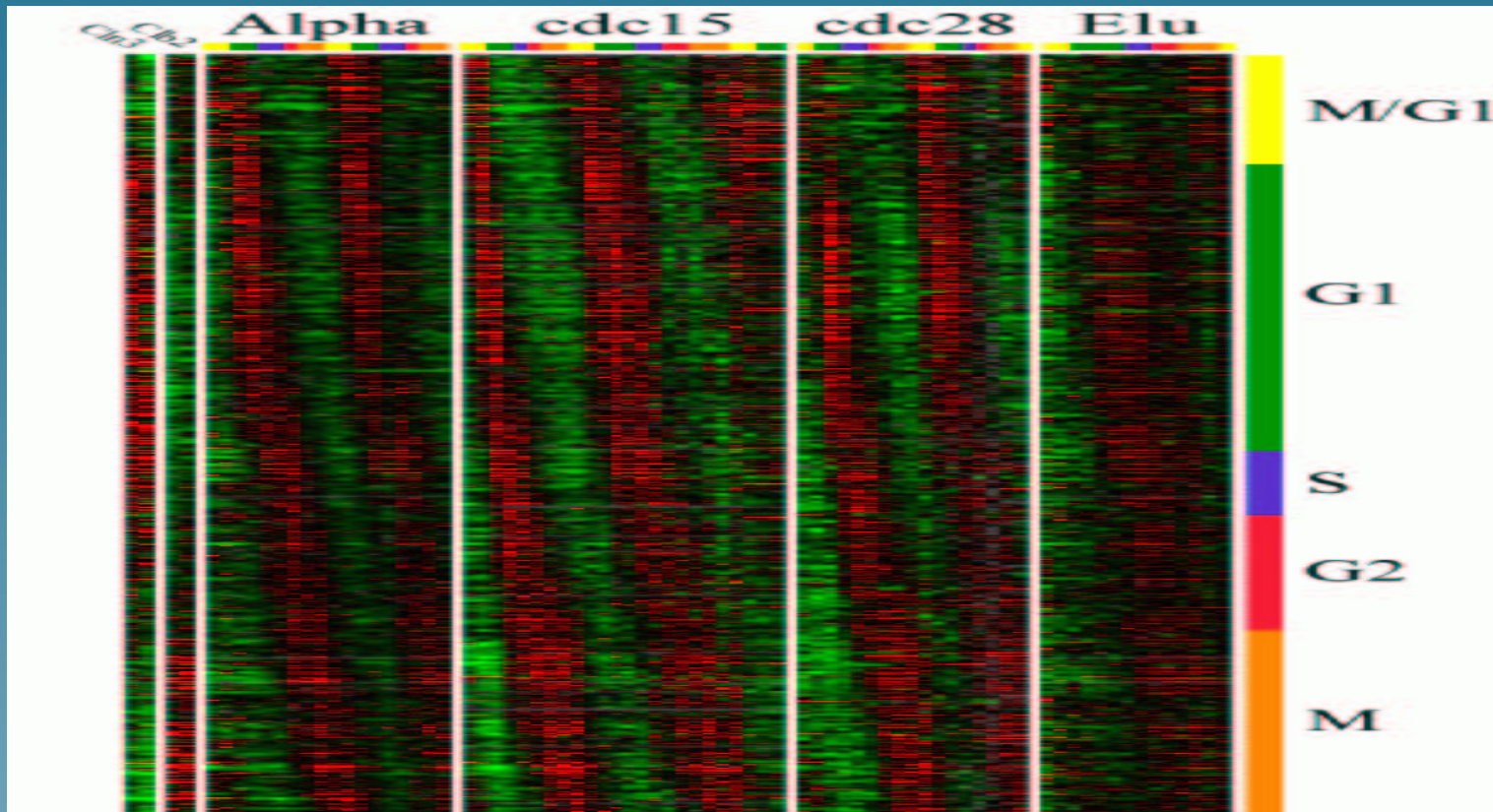


Nicotinamide Mononucleotide Adenylyltransferase With Bound Nad⁺

Chemical reactions are often parts of pathways

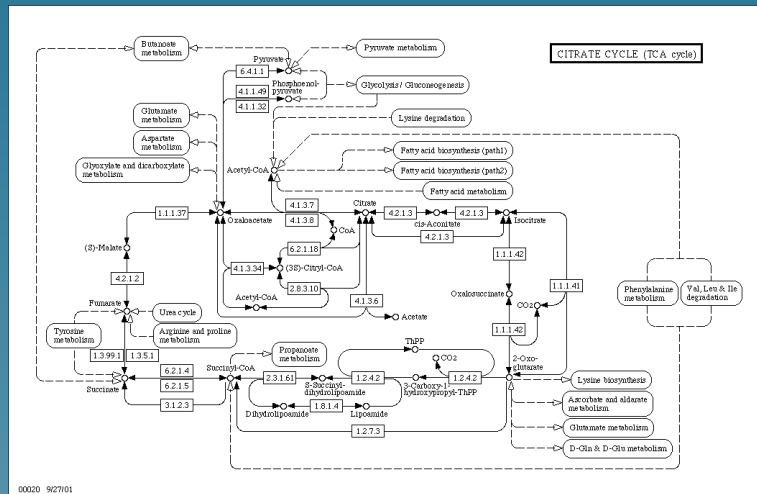


Microarray technology monitors RNA quantity

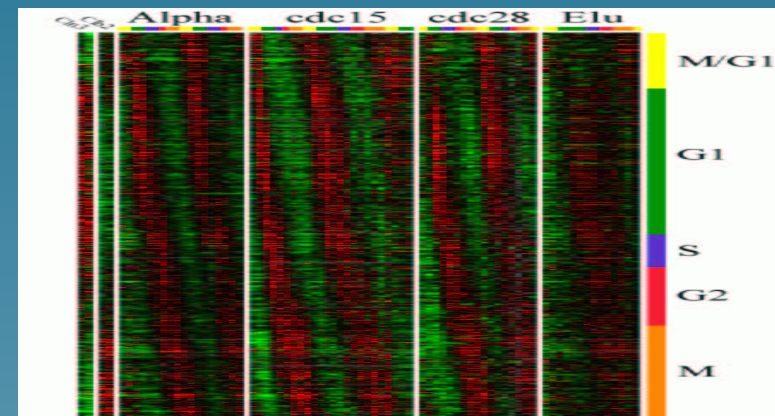


(From Spellman et al., 1998)

Comparing gene expression and pathway databases

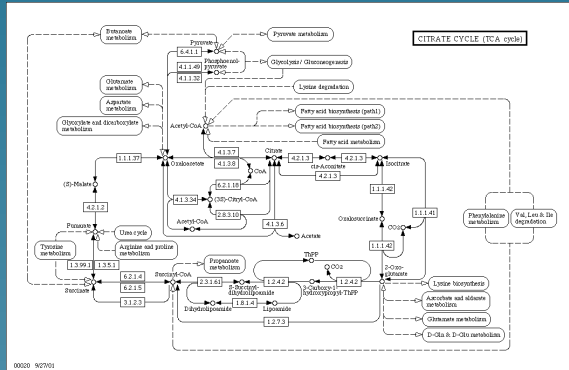


VS

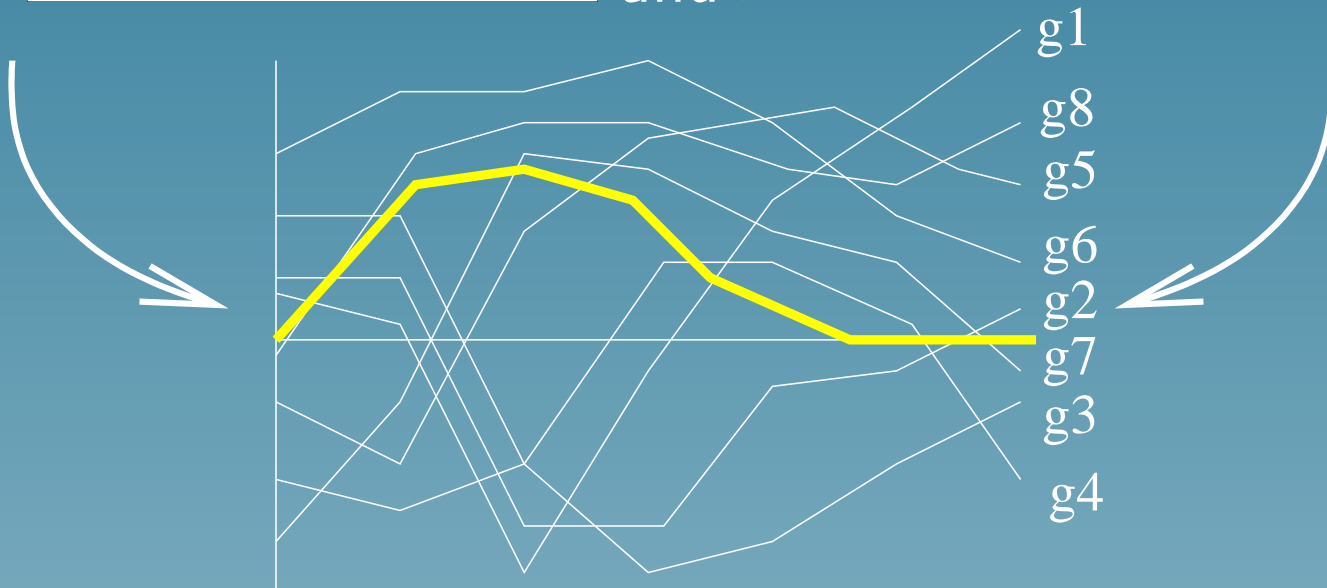


Detect active pathways? Denoise expression data?
 Denoise pathway database? Find new pathways?
 Are there “correlations”?

A useful first step



and



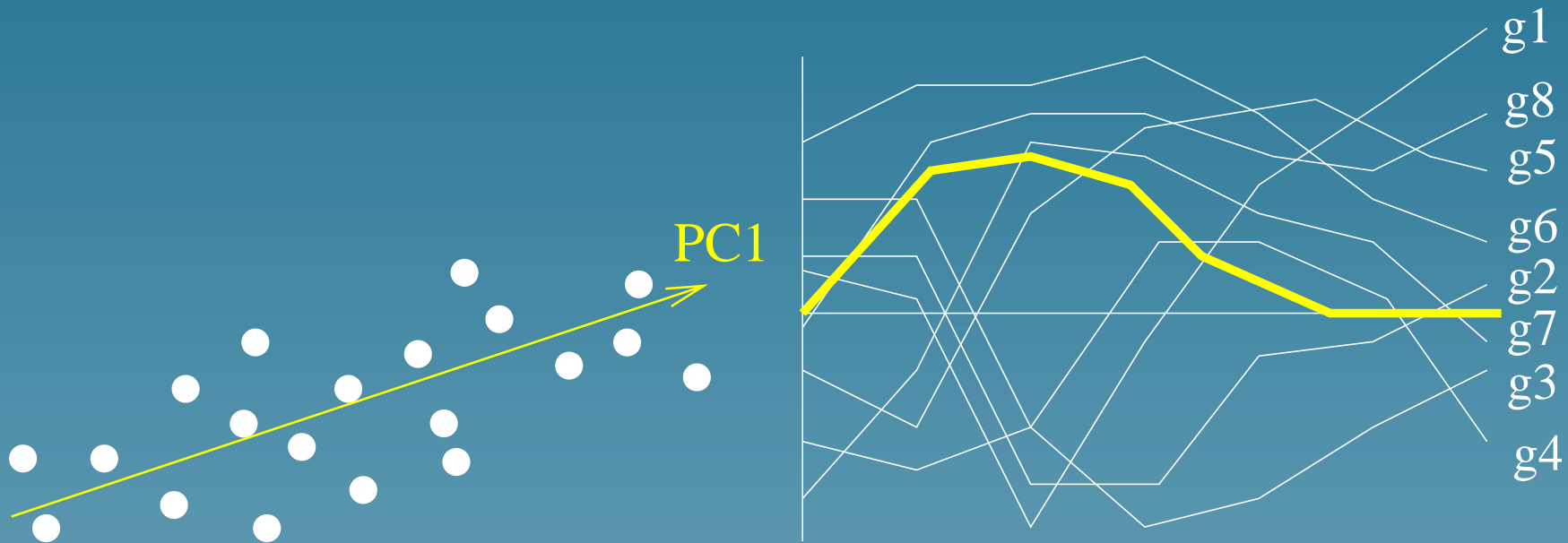
Part 1

Using expression data only

Motivation

- Pathways and biological events involve the coordinated action of several genes
- Co-regulation is an important way to coordinate the action of several genes
- Systematic variations in the set of gene expression profiles might be an indicator of an underlying biological phenomenon

Principal component analysis (PCA)



PCA finds the directions (*profiles*) explaining the **largest amount of variations** among expression profiles.

PCA formulation

- Let $f_v(i)$ be the **projection** of the i -th profile onto v .
- The **amount of variation** captured by f_v is:

$$h_1(v) = \sum_{i=1}^N f_v(i)^2$$

- PCA finds an orthonormal basis by solving successively:

$$\max_v h_1(v)$$

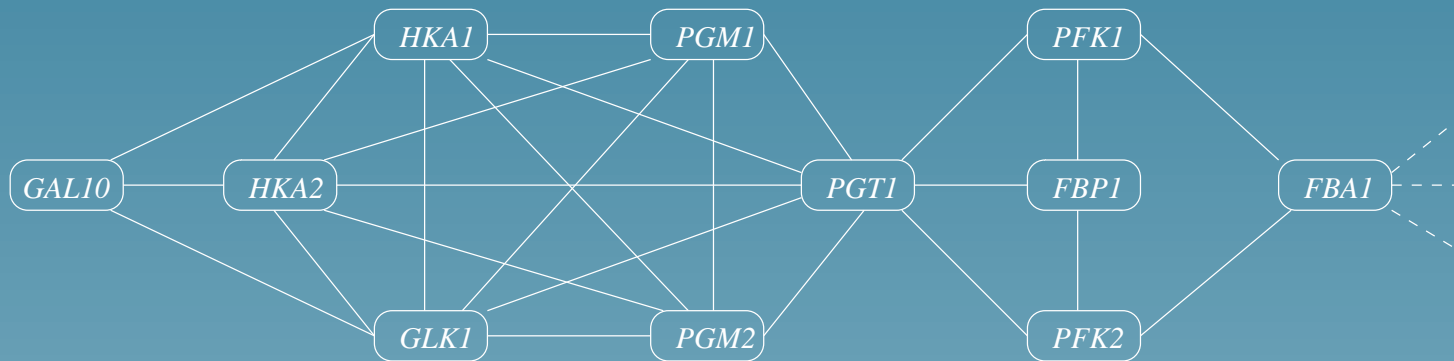
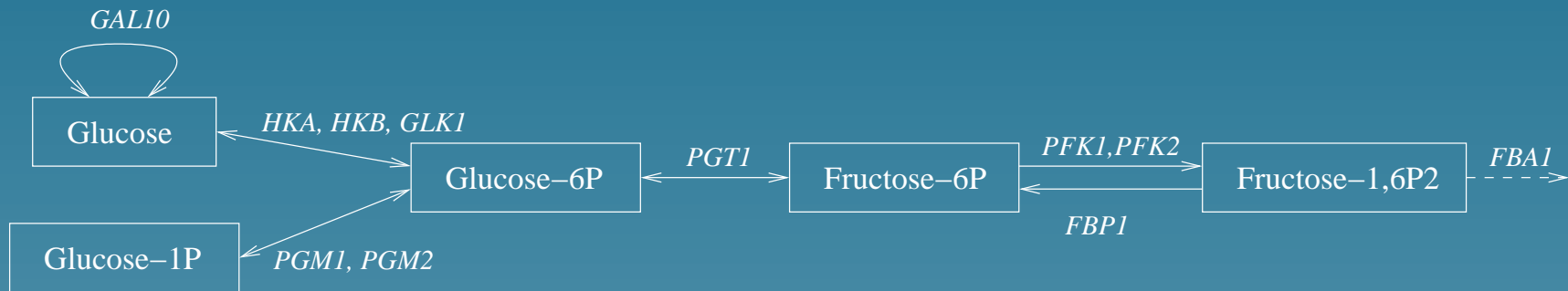
Part 3

Using the metabolic database

Motivation

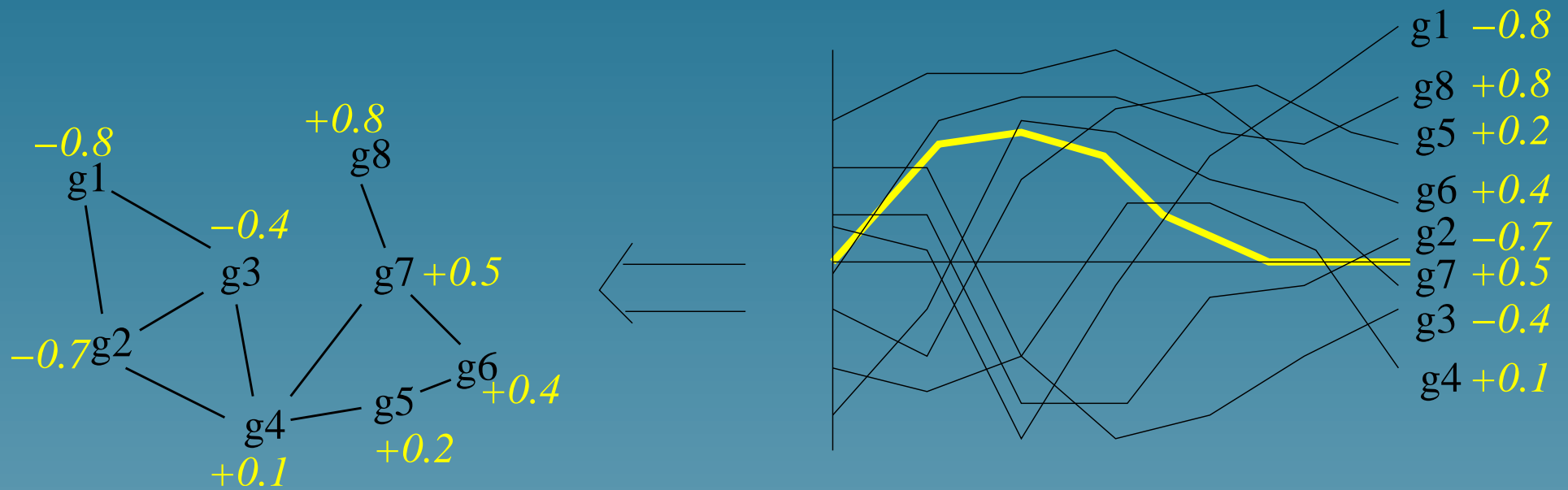
- PCA is useful if there is a small number of strong signal
- In concrete applications, we observe a **noisy superposition** of many events
- Using a prior knowledge of metabolic networks can help denoising the information detected by PCA

The metabolic gene network



Link two genes when they can **catalyze two successive reactions**

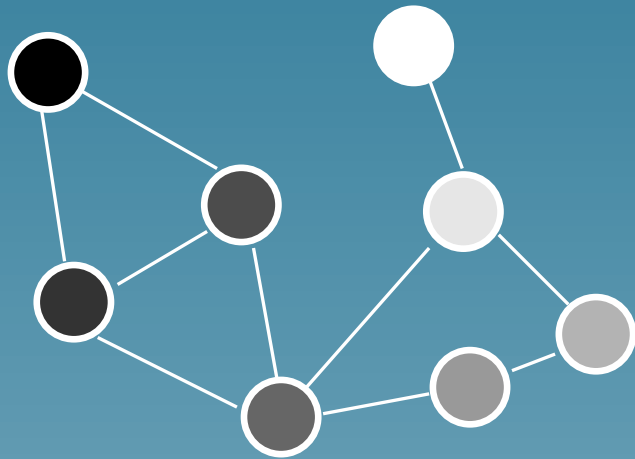
Mapping f_v to the metabolic gene network



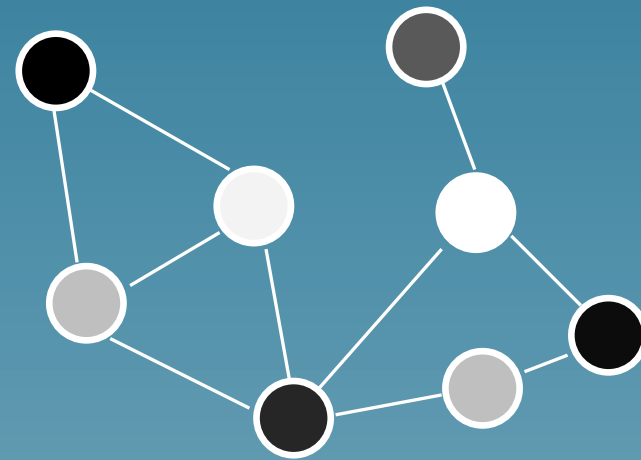
Does it look interesting or not?

Important hypothesis

If v is related to a metabolic activity, then f_v should **vary** "smoothly" on the graph

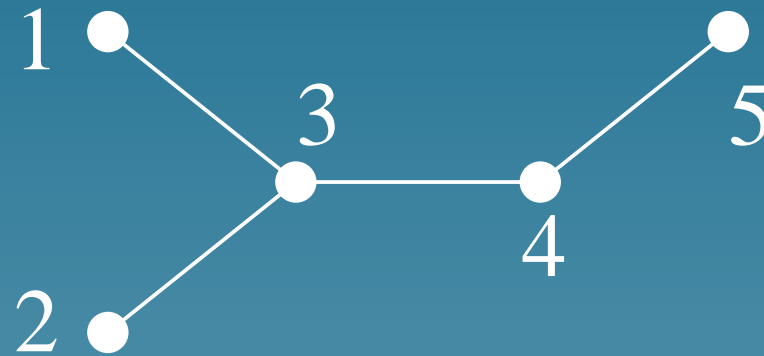


Smooth



Rugged

Graph Laplacian $L = D - A$

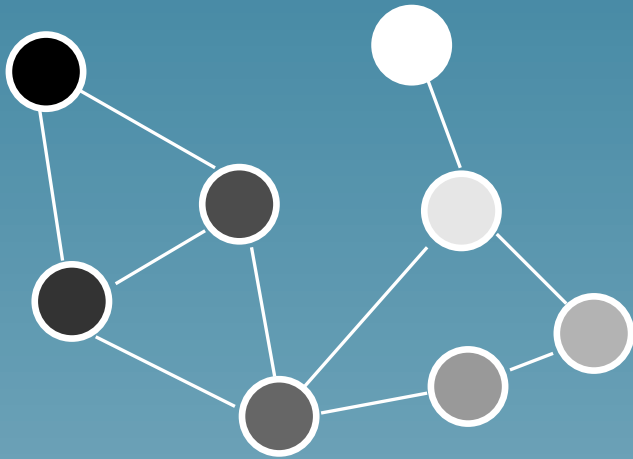


$$L = \begin{pmatrix} -1 & 0 & 1 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 \\ 1 & 1 & -3 & 1 & 0 \\ 0 & 0 & 1 & -2 & 1 \\ 0 & 0 & 0 & 1 & -1 \end{pmatrix}$$

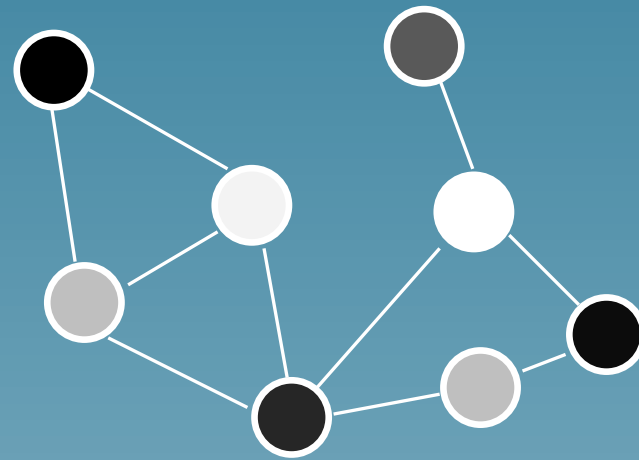
Smoothness quantification

$$h_2(f) = \frac{f^\top \exp(-\beta L) f}{f^\top f}$$

is large when f is smooth



$$h(f) = 2.5$$



$$h(f) = 34.2$$

Part 3

Combining expression and metabolic pathways

Motivation

For a candidate profile v ,

- $h_1(f_v)$ is large when v captures a lot of natural variation among profiles
- $h_2(f_v)$ is large when f_v is smooth on the graph

Try to maximize both terms in the same time

Problem reformulation

Find a function f_v and a function f_2 such that:

- $h_1(f_v)$ be large
- $h_2(f_2)$ be large
- $corr(f_1, f_2)$ be large

by solving:

$$\max_{(f_1, v)} corr(f_1, f_2) \times \frac{h_1(f_v)}{h_1(f_v) + \delta} \times \frac{h_2(f_2)}{h_2(f_2) + \delta}$$

Solving the problem

This formulation is equivalent to a generalized form of CCA (**Kernel-CCA**, Bach and Jordan, 2002), which is solved by the following generalized eigenvector problem

$$\begin{pmatrix} 0 & K_1 K_2 \\ K_2 K_1 & 0 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \rho \begin{pmatrix} K_1^2 + \delta K_1 & 0 \\ 0 & K_2^2 + \delta K_2 \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

where $[K_1]_{i,j} = e_i^\top e_j$ and $K_2 = \exp(-L)$.
Then, $f_v = K_1 \alpha$ and $f_2 = K_2 \beta$.

Part 4

Experimental results

Data

- **Gene network:** two genes are linked if they catalyze successive reactions in the KEGG database (669 yeast genes)
- **Expression profiles:** 18 time series measures for the 6,000 genes of yeast, during two cell cycles

First pattern of expression

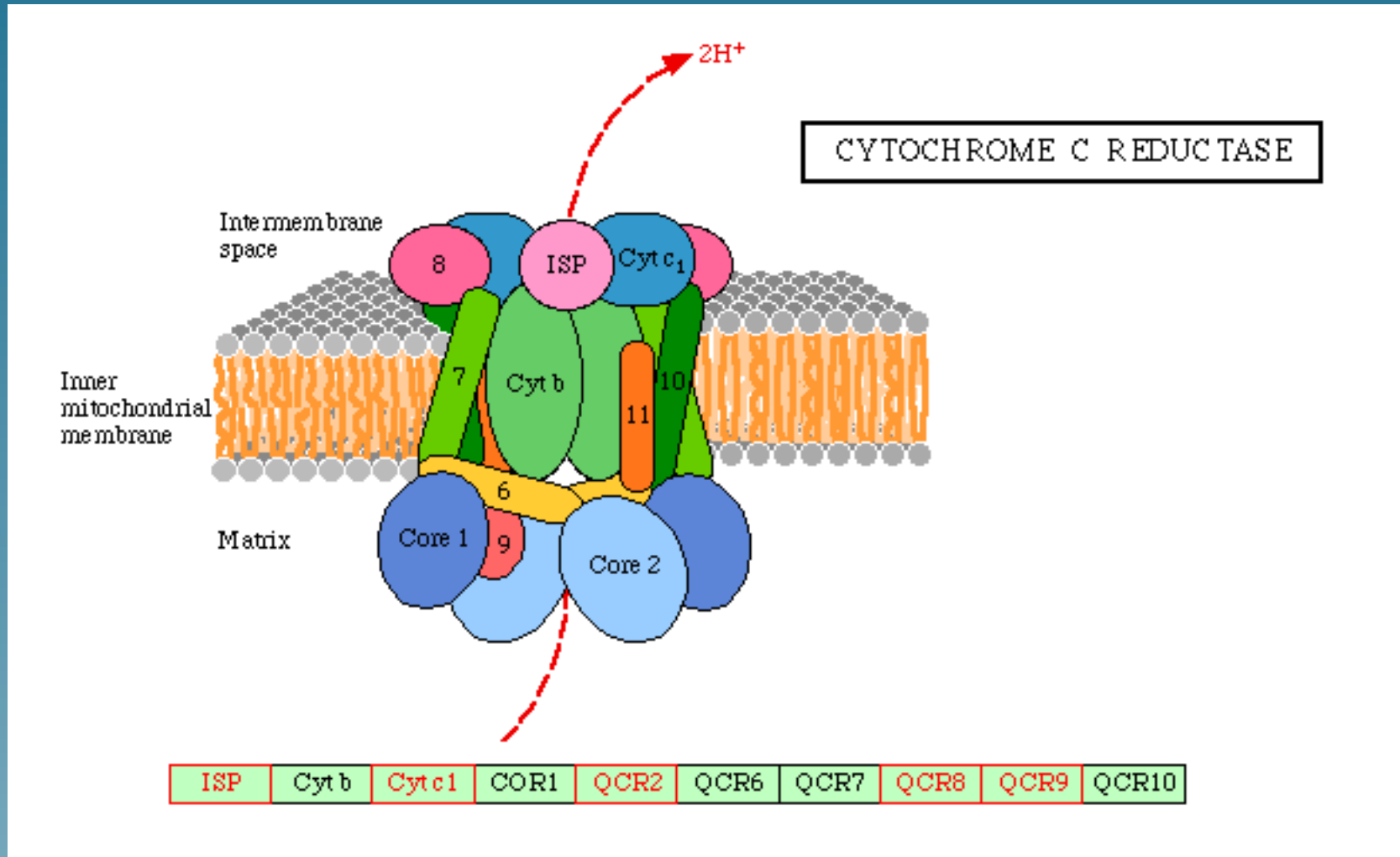


Related metabolic pathways

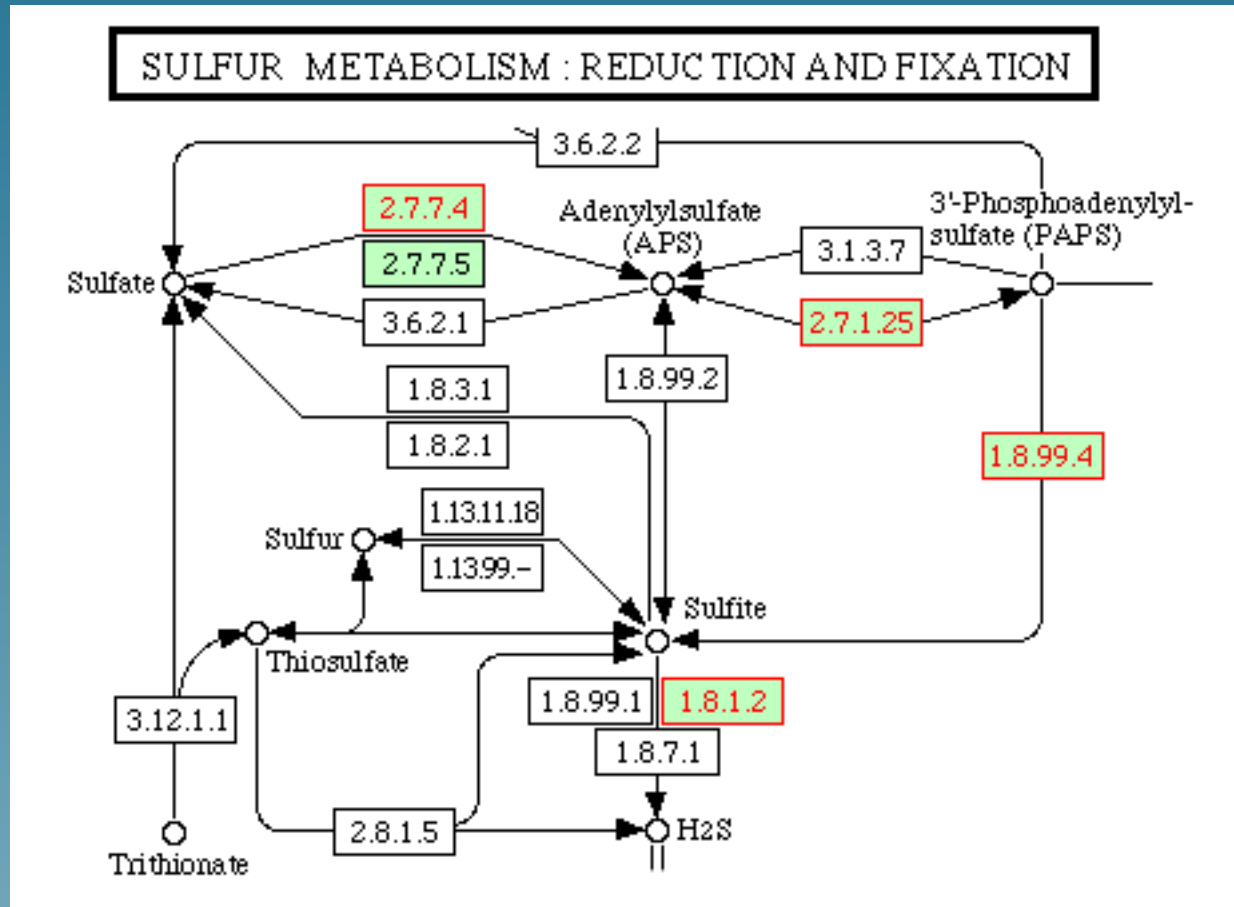
50 genes with highest $s_2 - s_1$ belong to:

- Oxidative phosphorylation (10 genes)
- Citrate cycle (7)
- Purine metabolism (6)
- Glycerolipid metabolism (6)
- Sulfur metabolism (5)
- Selenoaminoacid metabolism (4) , etc...

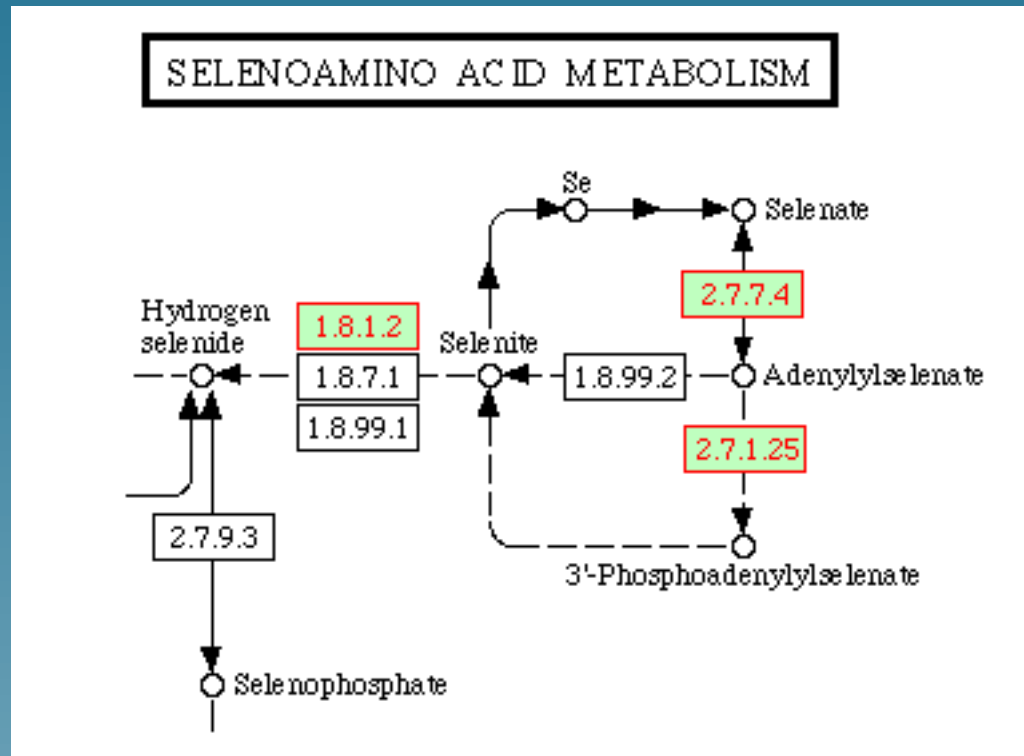
Related genes



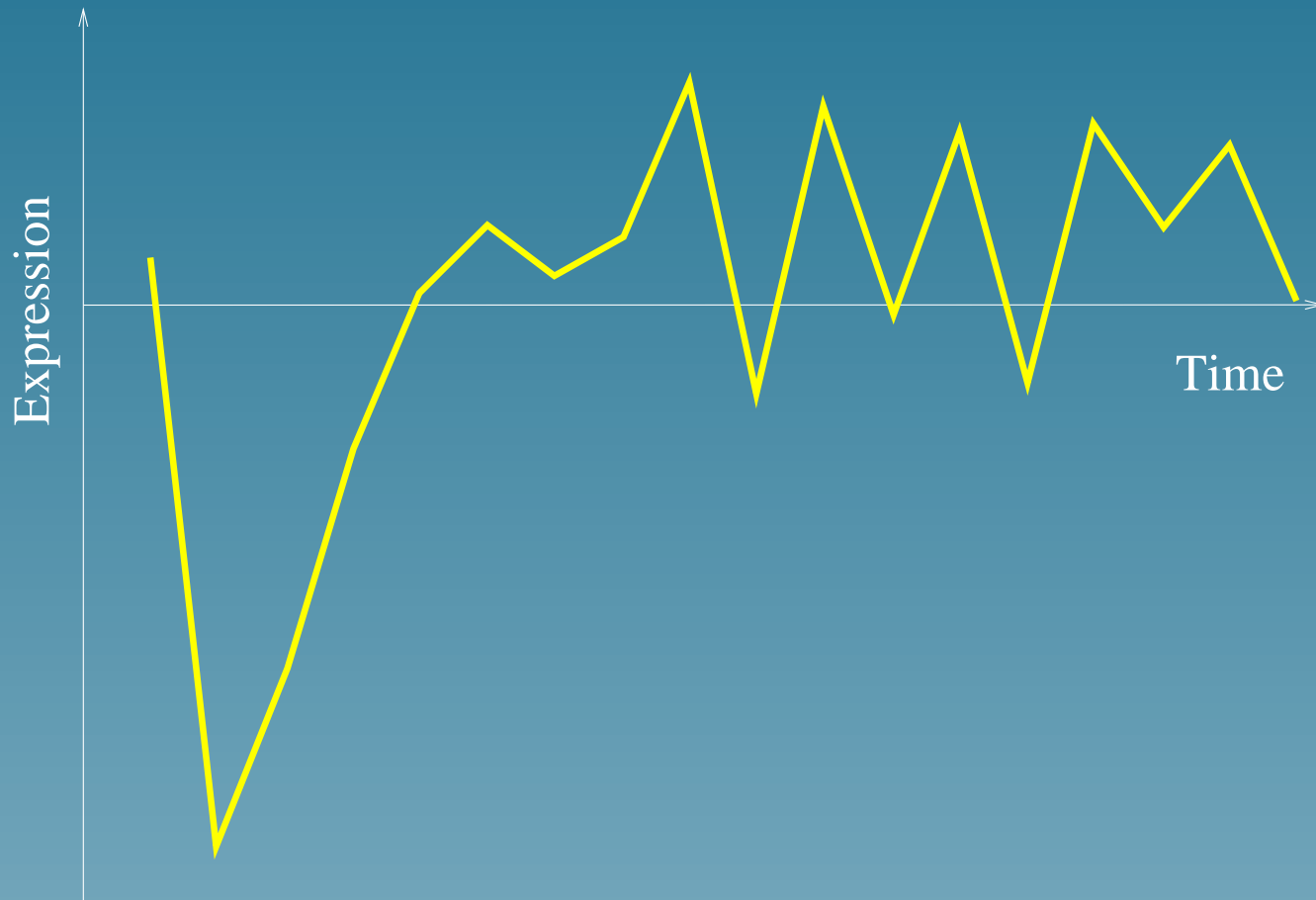
Related genes



Related genes



Opposite pattern



Related genes

- RNA polymerase (11 genes)
- Pyrimidine metabolism (10)
- Aminoacyl-tRNA biosynthesis (7)
- Urea cycle and metabolism of amino groups (3)
- Oxidative phosphorylation (3)
- ATP synthesis(3) , etc...

Related genes

RNA POLYMERASE

RNA polymerase II (*Saccharomyces cerevisiae*)

Eukaryotic Pol II

B2	B3	B4	B5	B6	B7
B1	B8	B9	B10	B11	B12

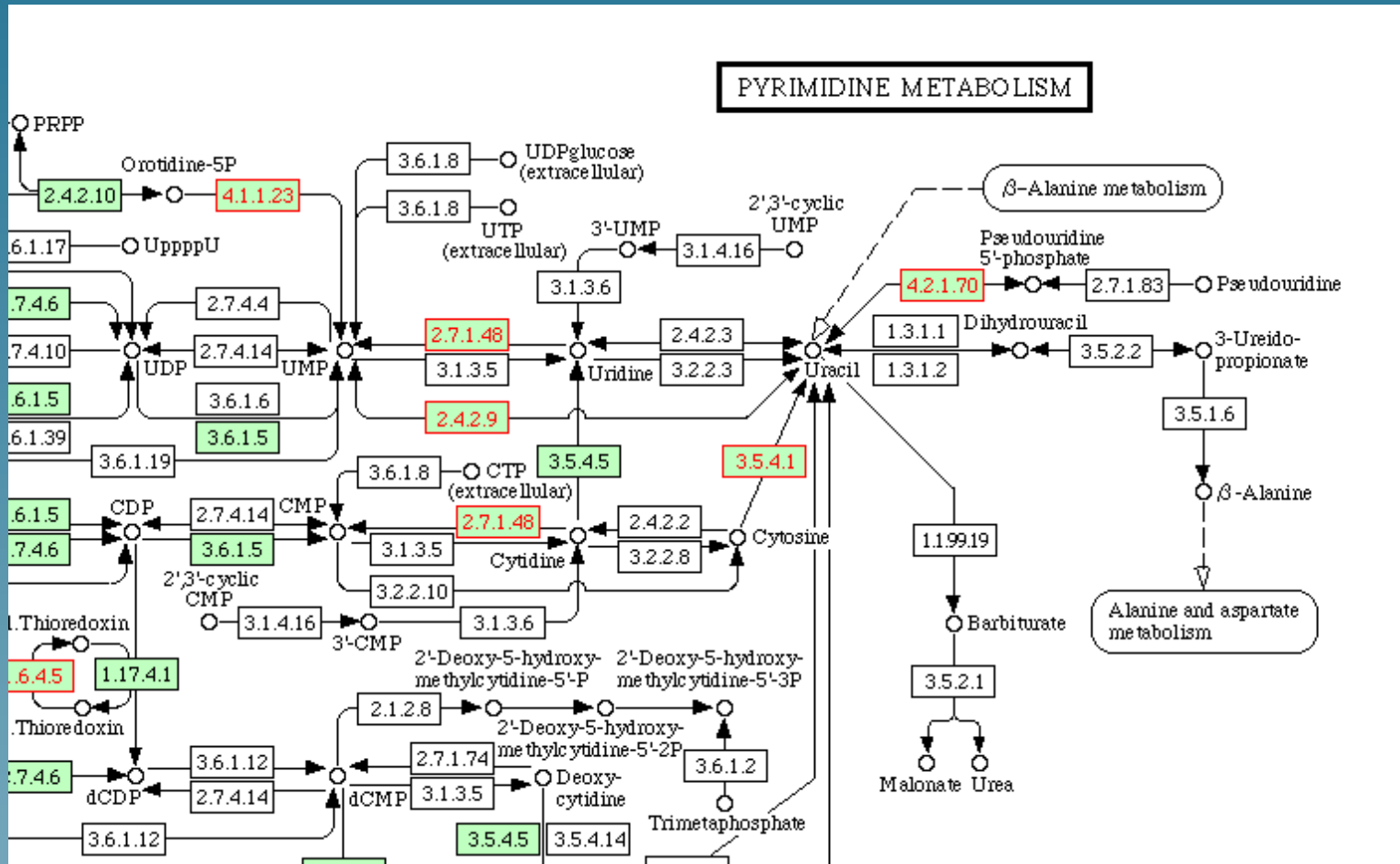
Eukaryotic Pol III

C2	C3	C4	C5	C11
C1	C19	C25	C31	C34

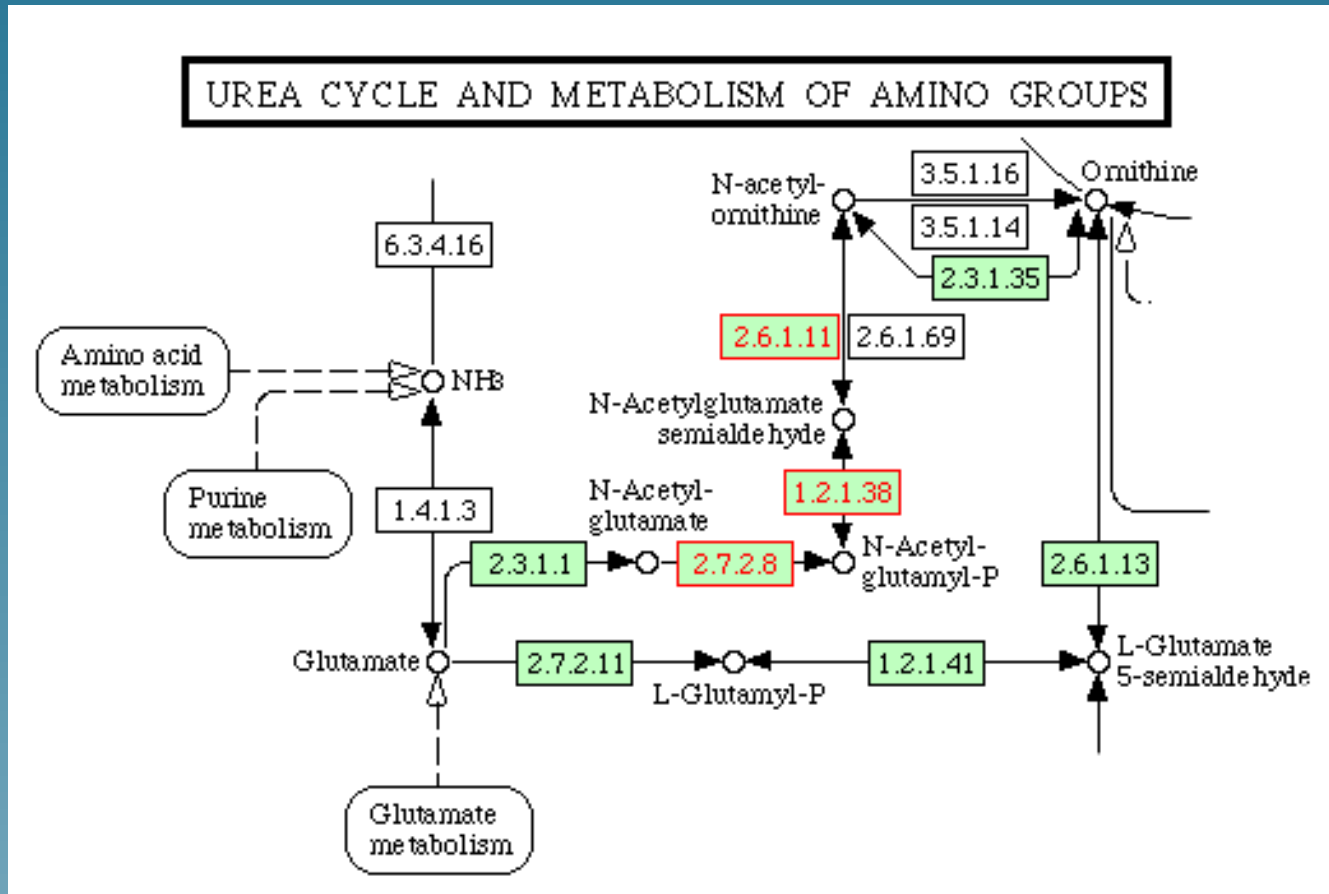
Eukaryotic Pol I

A2	A12	A14	A34	A43	A49
A1					

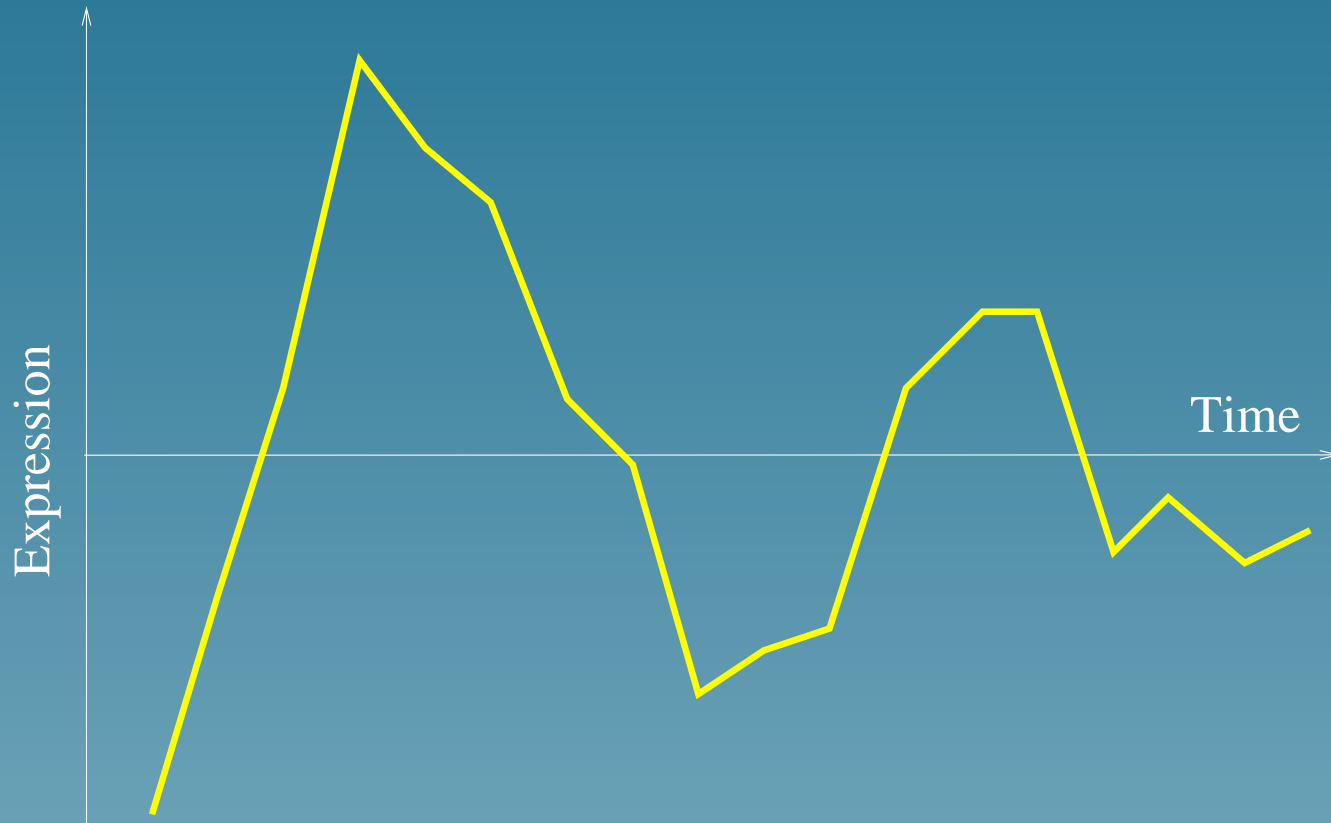
Related genes



Related genes



Second pattern



Conclusion

Conclusion

- An approach to **robustify PCA** using side information
- An approach to **integrate heterogeneous data**
- A particular case of more generic methods (**kernel methods**)
- Generalization to **other types of data** and **more than two datasets** is possible (see Yamanishi et al., ISMB 2003)