# Reconstructing networks of pathways via significance analysis of their intersections

Embedding biological knowledge in genomic statistical analysis

Mirko Francesconi, Daniel Remondini, Nicola Neretti, John Sedivy, Ettore

Verondini, Luciano Milanesi, Leon N Cooper, Gastone Castellani
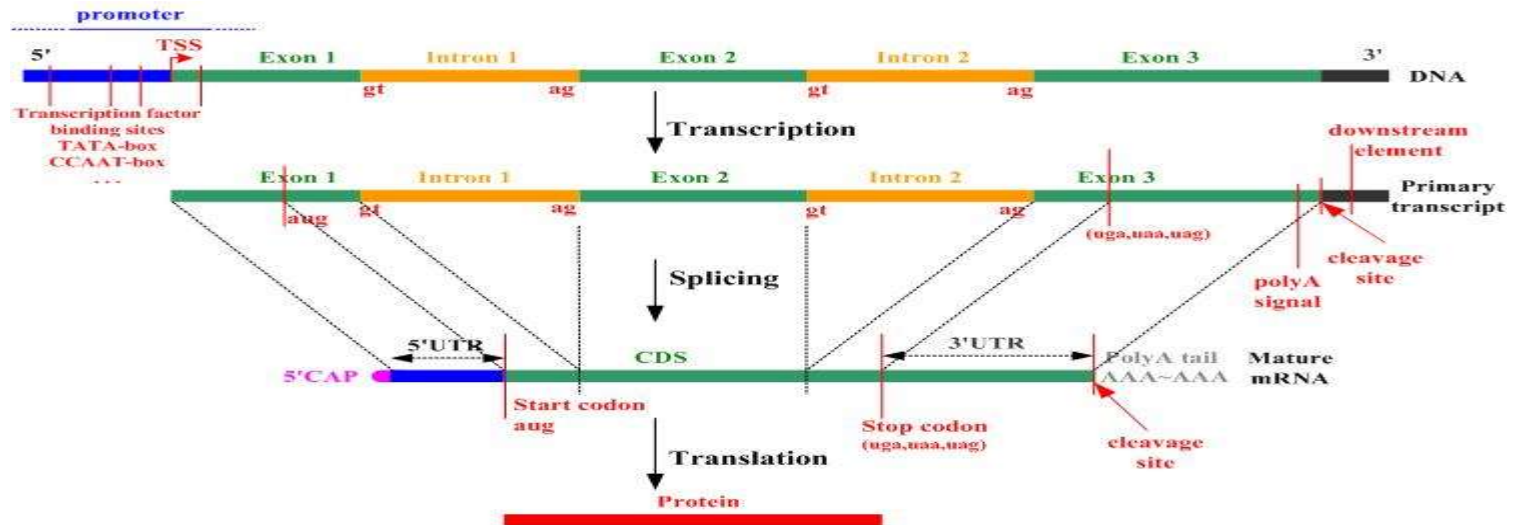
# Collaboration Bologna-Brown

| Brown University | Bologna University |
|---|---|
| **Brain Research Center** | **CIG-BBB Biophysics** |
| **Genomic Protemic Center** | **BioComplexity Bioinformatics** |
| Theoretical Physics | Systems Biophysics |
| Molecular Biology | **Unilever Research Center** |
| | ITB CNR Milano |

# Gene expression



- Regulation of transcription

We have generated and analyzed/ing several datasets

1) c-myc dataset (enginered rat fibroblasts)
2) TAC dataset (mouse)
3) Ewing sarcoma dataset (human)
4) Aging dataset (human time series & monozygotic twins)
5) c-myc exon array dataset (enginered rat fibroblasts)

# Probe selection

- Time series (myc on and myc off data sets, cardiac hypertrophy dataset)

- Linear model with empirical bayes shrinkage of variance (limma, Bioconductor).

- Contrasts of any time point with respect to zero time point

# Significance analysis:

## ANOVA-MULTIPLE TEST COMPARISON

- Preprocessing for **"dimensionality reduction"** of the probeset number

- Identify genes with **significative expression levels difference** between the two conditions (perturbed and unperturbed)

- Differences are analyzed over all times

- Significance analysis applied to all probesets and **eventual correction with FDR**

# c-Myc-triggered gene expression

- C-Myc encode for **transcriptional regulators** whose inappropriate expression is correlated with a wide array of human malignancies.

- Up-regulation of Myc enforces growth, antagonizes cell cycle withdrawal and differentiation, and in some situations promotes apoptosis.

- c-*myc*-/- cells reconstituted with the conditionally active, tamoxifen-specific c-Myc-estrogen receptor fusion protein (MycER) allows the fine and selective change of of c-Myc activity  by Tamoxifen .

**Time series experiment** with **5 time points** in triplicate and 9000 probes

From the J.M. Sedivy lab          O'Connel et al JBC 2004

# Evaluation of global gene expression of left ventricular tissue in animal model of left ventricular hypertrophy (LVH) induced by transverse aortic constriction (TAC).
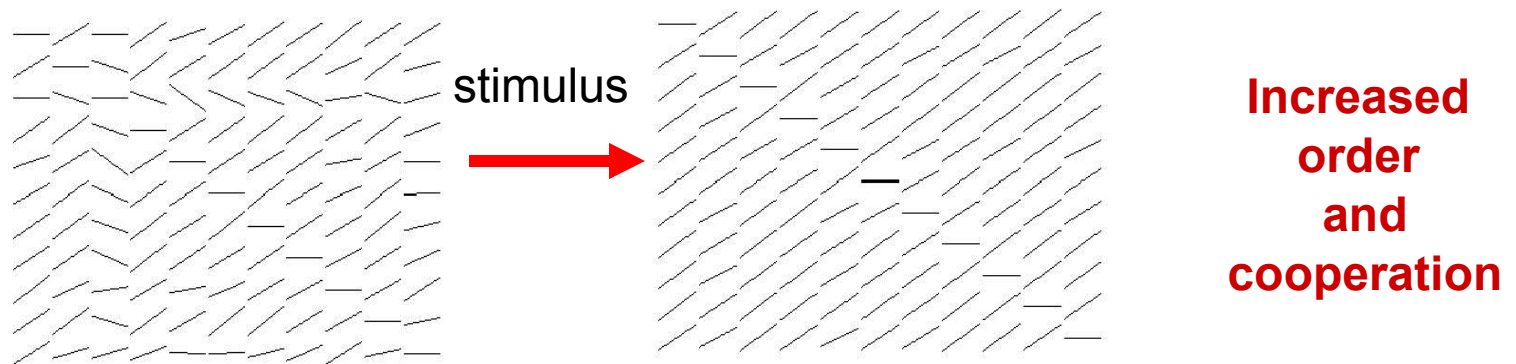
- Time series experimental design

- Measurements were done by 15 Affymetric chips at T1=0, T2=2,or T3=4 weeks after TAC.

- Each time point have been repeated with 5 replica

# Genomic analysis drawbacks

- single gene analysis is not sufficient to understand cell mechanisms undergoing experimental conditions

- cell behaviour is a complex phenomenon: several elements (e.g. genes) act together in order to generate it

# Perturbation approach

- These experiments can be conceptualized as **"perturbation"** of a "basal state" (cell growth, metabolism, young phenotype, cancer phenotype etc)

- "External perturbations" like temperature in physical systems are realized by gene activation via transcription factor triggering (c-myc, dfoxo-nutrition, aging)

- Emergent properties arising in the context of perturbation theory are the so called **"phase transitions"** (superconductivity, superfluidity,etc) and **"condensation"** phenomena.
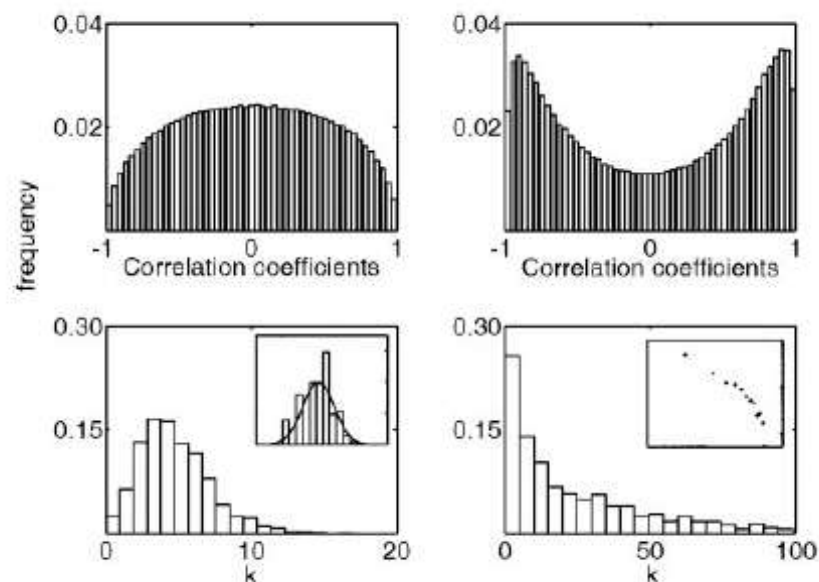
stimulus

**Increased order and cooperation**

# Targeting c-Myc-activated genes with a correlation method: Detection of global changes in large gene expression network dynamics

D. Remondini*[†‡], B. O'Connell[§], N. Intrator[¶‖], J. M. Sedivy[§], N. Neretti[†¶], G. C. Castellani*[†‡¶]**, and L. N. Cooper[¶]**[††‡‡]

*Dipartimento di Fisica and [†]Galvani Center for Biocomplexity, Università di Bologna, Bologna 40127, Italy; Departments of [§]Molecular Biology, Cell Biology, and Biochemistry, [††]Physics, and [‡‡]Neuroscience and [¶]Institute for Brain and Neural Systems, Brown University, Providence, RI 02912; [‖]School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel; and [‡]Dipartimento di Morfofisiologia Veterinaria e Produzioni Animali , Università di Bologna, Ozzano Emilia 40064, Italy

This work studies the dynamics of a gene expression time series network. The network, which is obtained from the correlation of gene expressions, exhibits global dynamic properties that emerge after a cell state perturbation. The main features of this network appear to be more robust when compared with those obtained with a network obtained from a linear Markov model. In particular, the network properties strongly depend on the exact time sequence relationships between genes and are destroyed by random temporal data shuffling. We discuss in detail the problem of finding targets of the c-myc protooncogene, which encodes a transcriptional regulator whose inappropriate expression has been correlated with a wide array of malignancies. The data used for network construction are a time series of gene expression, collected by microarray analysis of a rat fibroblast cell line expressing a conditional Myc-estrogen receptor oncoprotein. We show that the correlation-based model can establish a clear relationship between network structure and the cascade of c-myc-activated genes.

# Targeting c-Myc-activated genes with a correlation method: Detection of global changes in large gene expression network dynamics

D. Remondini*[†‡], B. O'Connell[§], N. Intrator[¶‖], J. M. Sedivy[§], N. Neretti[†¶], G. C. Castellani*[†‡¶**], and L. N. Cooper[¶**††‡‡]

*Dipartimento di Fisica and [†]Galvani Center for Biocomplexity, Università di Bologna, Bologna 40127, Italy; Departments of [§]Molecular Biology, Cell Biology, and Biochemistry, [††]Physics, and [‡‡]Neuroscience and [¶]Institute for Brain and Neural Systems, Brown University, Providence, RI 02912; [‖]School of Computer Science, Tel Aviv University, Tel Aviv 69978, Israel; and [‡]Dipartimento di Morfofisiologia Veterinaria e Produzioni Animali, Università di Bologna, Ozzano Emilia 40064, Italy

**Fig. 3.** Network of selected Myc-influenced pathways showing positive and negative correlations. The red and blue arrows denote positive and negative co-regulation, respectively. The thickness of the arrows is proportional to the magnitude, or absolute value, of the co-regulation. A network with these properties is called a weighted directed graph.

# Multiscale correlation for co-regulation detection

•Capture **correlation profile  changes  at several scales** (whole array, gene family and pathways) and is informative of significative activity

•pathways synthesis into single functional forms (**Fluxes**) or index such as Subgraph Conductance.

•assessment of co-regulation between and within several  pathways

•When the perturbation is conditionally switched on, the correlation between genes with a significant change in their expression level is altered  on a genomic scale

We have strong indications that **a similar transition is conserved  on different scales** and is indicative of **co-regulation changes**

To reduce the **dimensionality** of the problem and introduce "**a-priori biological knowledge**", we will extend this method  by mapping the gene arrays data onto gene pathways and ontologies.

*Castellani et al, PNAS 2001*

*Castellani et al, Learning and Memory 2005, BMC Bioinformatics 2007, IJCB 2007*

# A biophysical model of bidirectional synaptic plasticity: Dependence on AMPA and NMDA receptors

Gastone C. Castellani*, Elizabeth M. Quinlan[†], Leon N Cooper[‡§¶], and Harel Z. Shouval[‡‖]

*Physics Department, CIG and Dimorfipa Bologna University, Bologna 40121, Italy; [†]Department of Biology, University of Maryland, College Park, MD 20742; and [‡]Institute for Brain and Neural Systems, [§]Department of Neuroscience, and [¶]Department of Physics, Brown University, Providence, RI 02912
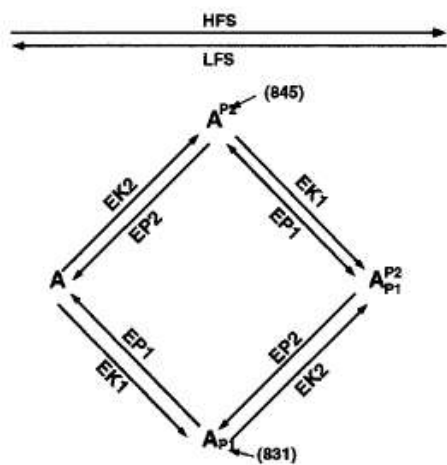


Fig. 1. An idealized model for the cycle of GluR1 phosphorylation/dephosphorylation at two sites. The model assumes two specific kinases (EK1, EK2) and two opposing specific phosphatases (EP1, EP2). It is assumed that high-frequency stimulation preferentially stimulates the activity of protein kinases, resulting in GluR1 phosphorylation, whereas low-frequency stimulation preferentially stimulates the activity of protein phosphatases, resulting in GluR1 dephosphorylation.



Fig. 3. Synaptic strength, measured as AMPAR conductance depicted as a function of presynaptic stimulation frequency ($f$) and postsynaptic membrane voltage ($V$). (a) A two-dimensional plot depicting postsynaptic membrane potential as a function of presynaptic stimulation frequency. The grey scale indicates the conductance level of the AMPAR. At low stimulation frequencies and postsynaptic voltages, the conductance is below baseline, defined as $f = 0, V = -100$. The diagonal line indicates a linear $f - V$ relation, which we assume to extract the results in $b$. (b) AMPAR conductance as a function of presynaptic stimulation frequency, where a linear dependence of $V$ on $f$ is assumed (as shown in $a$). Low-frequency stimulation induces LTD, whereas high-frequency stimulation induces LTP.

# A model of bidirectional synaptic plasticity: From signaling network to channel conductance

Gastone C. Castellani,[1,2,6] Elizabeth M. Quinlan,[4] Ferdinando Bersani,[1]
Leon N. Cooper,[2,3] and Harel Z. Shouval[2,5]

[1]Physics Department, DIMORFIPA, CIG, Bologna University, Bologna 40137, Italy; [2]Institute for Brain and Neural Systems and [3]Physics and Neuroscience Department, Brown University, Providence, Rhode Island 02912, USA; [4]Neuroscience and Cognitive Sciences Program, University of Maryland, College Park, Maryland 20742, USA; [5]Department of Neurobiology and Anatomy, University of Texas Medical School at Houston, Houston, Texas 77030, USA

In many regions of the brain, including the mammalian cortex, the strength of synaptic transmission can be bidirectionally regulated by cortical activity (synaptic plasticity). One line of evidence indicates that long-term synaptic potentiation (LTP) and long-term synaptic depression (LTD), correlate with the phosphorylation/dephosphorylation of sites on the $\alpha$-Amino-3-hydroxy-5-methyl-4-isoxazolepropionic acid (AMPA) receptor subunit protein GluRI. Bidirectional synaptic plasticity can be induced by different frequencies of presynaptic stimulation, but there is considerable evidence indicating that the key variable is calcium influx through postsynaptic N-methyl-D-aspartate (NMDA) receptors. Here, we present a biophysical model of bidirectional synaptic plasticity based on $[Ca^{2+}]$-dependent phospho/dephosphorylation of the GluRI subunit of the AMPA receptor. The primary assumption of the model, for which there is wide experimental support, is that the postsynaptic calcium concentration, and consequent activation of calcium-dependent protein kinases and phosphatases, is the trigger for phosphorylation/dephosphorylation at GluRI and consequent induction of LTP/LTD. We explore several different mathematical approaches, all of them based on mass-action assumptions. First, we use a first order approach, in which transition rates are functions of an activator, in this case calcium. Second, we adopt the Michaelis-Menten approach with different assumptions about the signal transduction cascades, ranging from abstract to more detailed and biologically plausible models. Despite the different assumptions made in each model, in each case, LTD is induced by a moderate increase in postsynaptic calcium and LTP is induced by high $Ca^{2+}$ concentration.
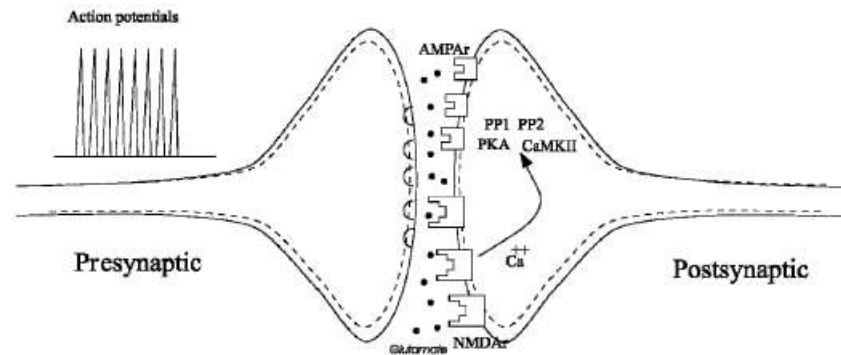
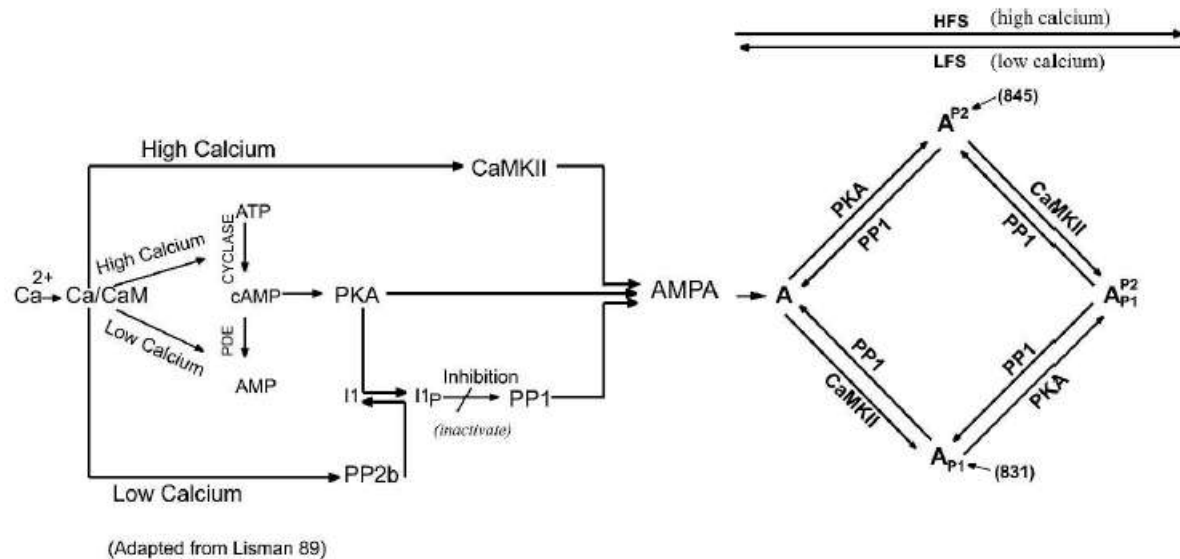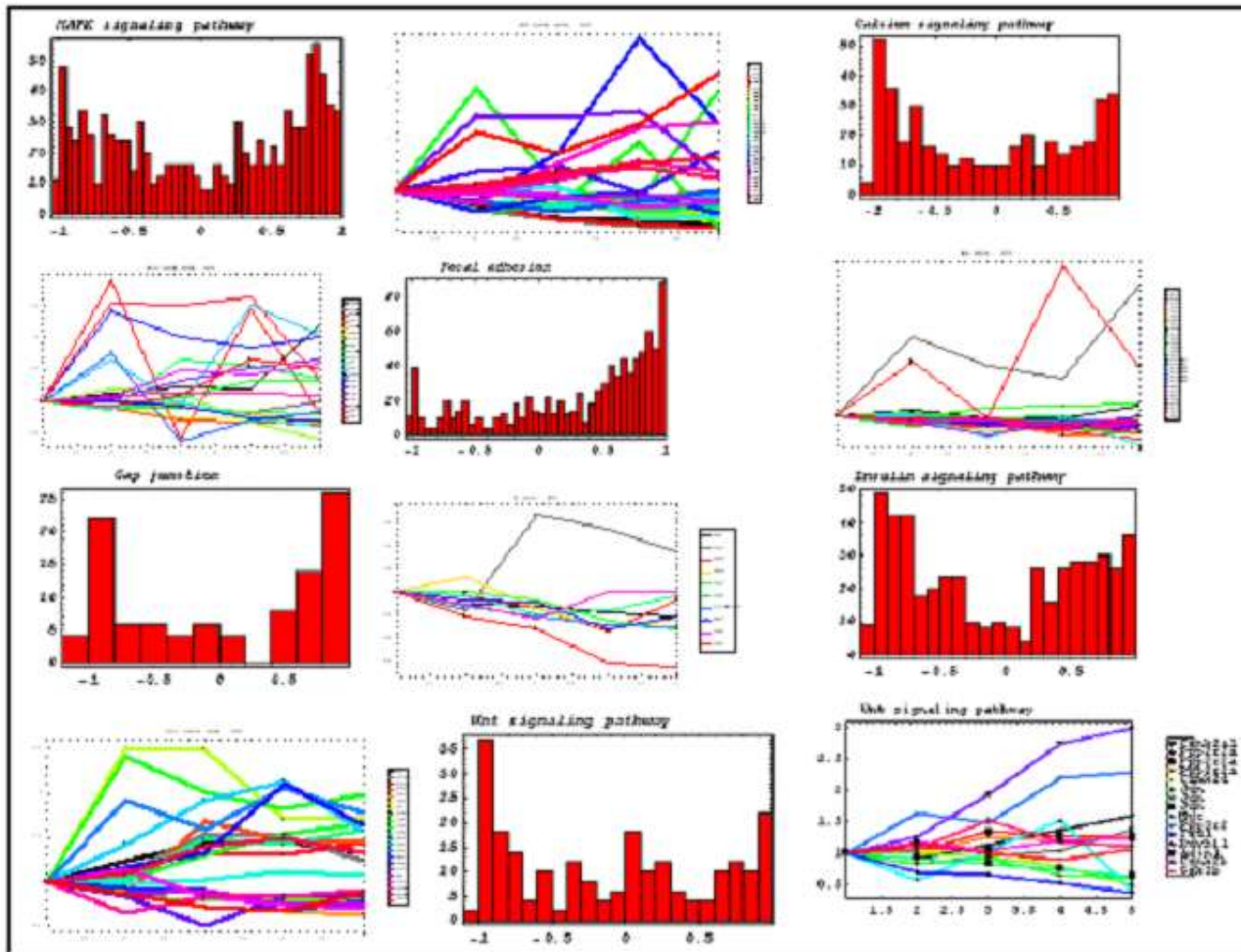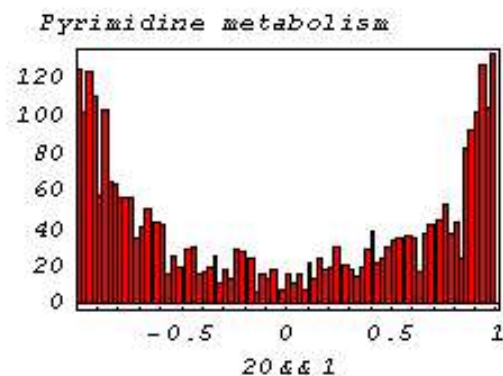**Figure 1.** A schematic of an excitatory glutamatergic synapse. Action

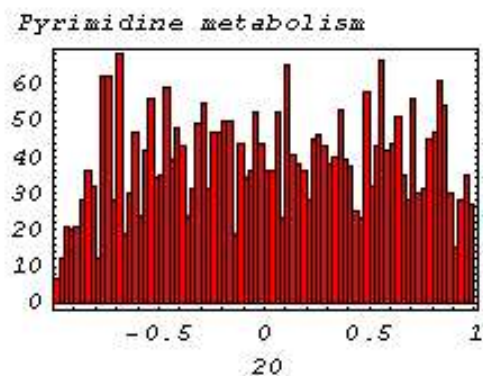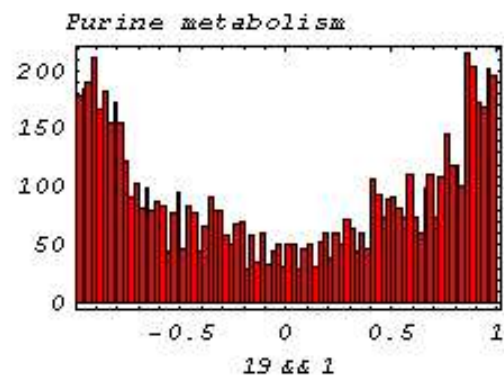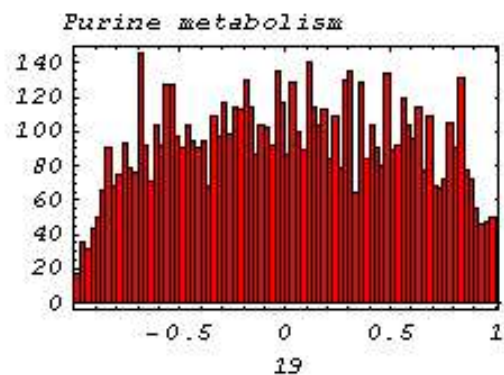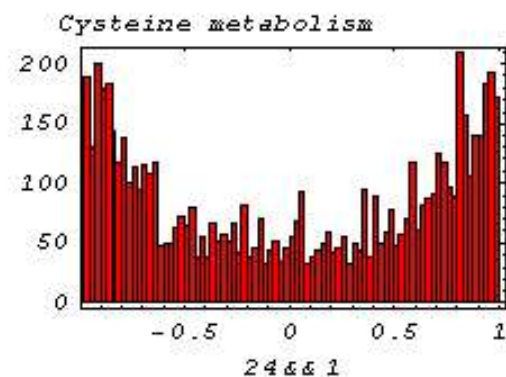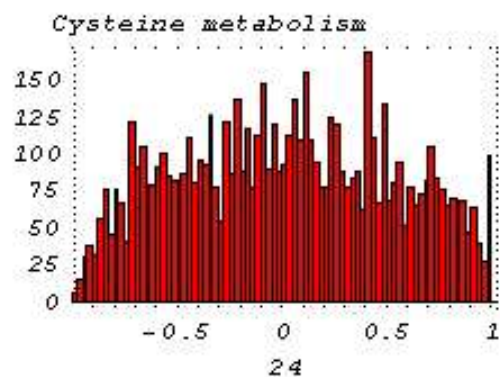(Adapted from Lisman 89)

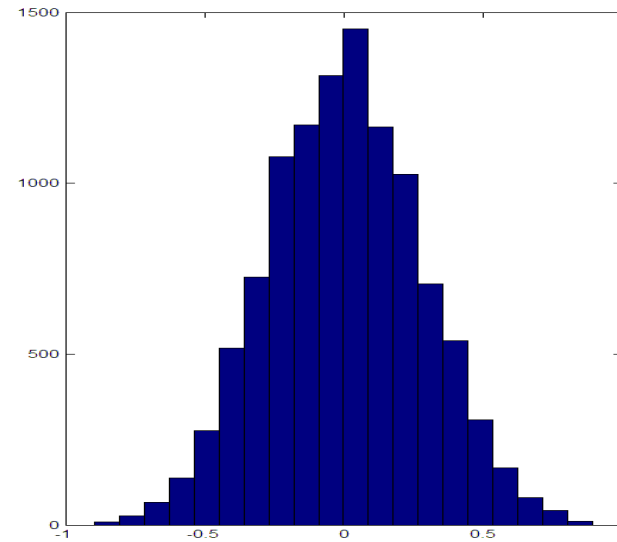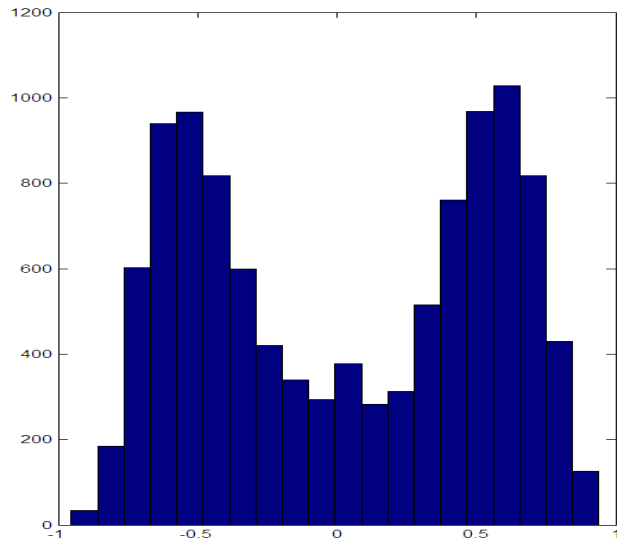**Figure 6.** Calcium-dependent activity of the kinase/phosphatase network. (*Left*) Introduction of
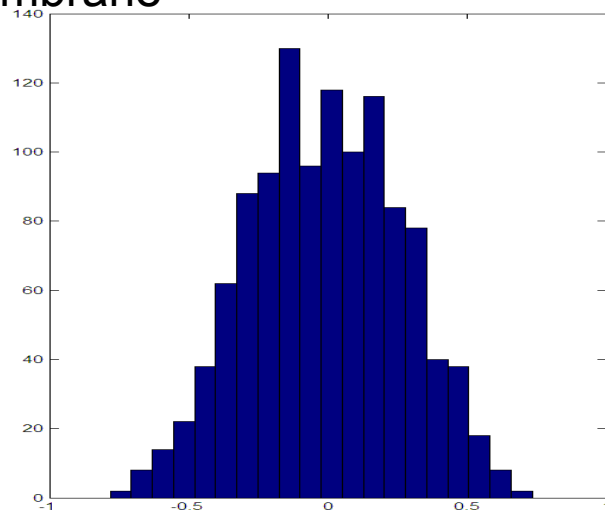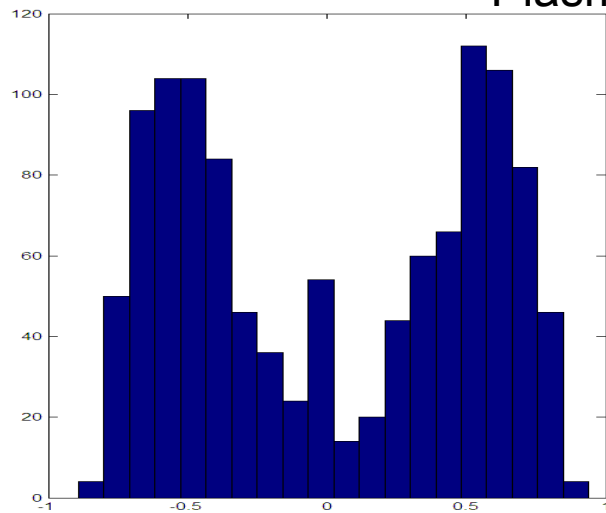
# Multiscale Correlation Model:
## c-Myc results

Cysteine metabolism
24

Cysteine metabolism
24 && 1

Purine metabolism
19

Purine metabolism
19 && 1

Pyrimidine metabolism
20

Pyrimidine metabolism
20 && 1

# Multiscale Correlation Model: human aging results



Protein Binding

Plasma Membrane

Castellani et al International Journal of Chaos and Bifurcation 2007

**HUMAN AGING**

| | | | |
|---|---|---|---|
| 1 | PPAR SigPath | 26 | Apoptosis |
| 2 | Adipocytokine SigPath | 27 | Carbon fixation |
| 3 | Inositol phosphate Met | 28 | Colorectal cancer |
| 4 | Jak-STAT SigPath | 29 | Glutathione metabolism |
| 5 | Phosphatidylinositol SigSyst | 30 | $\gamma$-ExaCloCE Degr |
| 6 | Purine metabolism | 31 | Antigen ProcAndPres |
| 7 | Glyo and Dicarbo xylate Met | 32 | Cyanoamino Ac Met |
| 8 | Cysteine metabolism | 33 | Gap junction |
| 9 | B cell receptor SigPath | 34 | Taur HypoTaur Met |
| 10 | Glycolysis-Gluconeogenesis | 35 | ALA-ASP Met |
| 11 | Styrene degradation | 36 | Leuk tr-e migration |
| 12 | Long-term depression | 37 | Atrazine Deg |

| | | | |
|---|---|---|---|
| 13 | Alkaloid Bios I | 38 | Nitrogen metabolism |
| 14 | Tyrosine Met | 39 | Hematopoietic cell lineage |
| 15 | mTOR SigPath | 40 | Glycan STR-Bios 1 |
| 16 | Fc ε RI SigPath | 41 | VEGF SigPath |
| 17 | Bisphenol A Degr | 42 | Focal adhesion |
| 18 | Val Leu ILeu Bios | 43 | Nicotinate and nicotinamide metabolism |
| 19 | Complement and Coag | 44 | Ribosome |
| 20 | Pyrimidine metabolism | 45 | Insulin SigPath |
| 21 | Pyruvate metabolism | 46 | Cell cycle |
| 22 | Benzoate degradation | 47 | Cytk-Citk RecInt |
| 13 | Type II Diab Mell | 48 | Glutamate Met |
| 14 | PhenylAla Met | 49 | Propanoate Met |
| 15 | T cell Reec SigPath | 50 | Toll-like Rec SigPath |

"Databases" like KEGG have also an interesting **network structure**, it is possible that biologically relevant informations can be retrieved from the **topological structure** of **nodes** (pathways) and **edges** (common genes between two pathways)

The most relevant edges can be **focal areas** from which biological messages are spread throughout the network (like the **hubs** for the nodes)
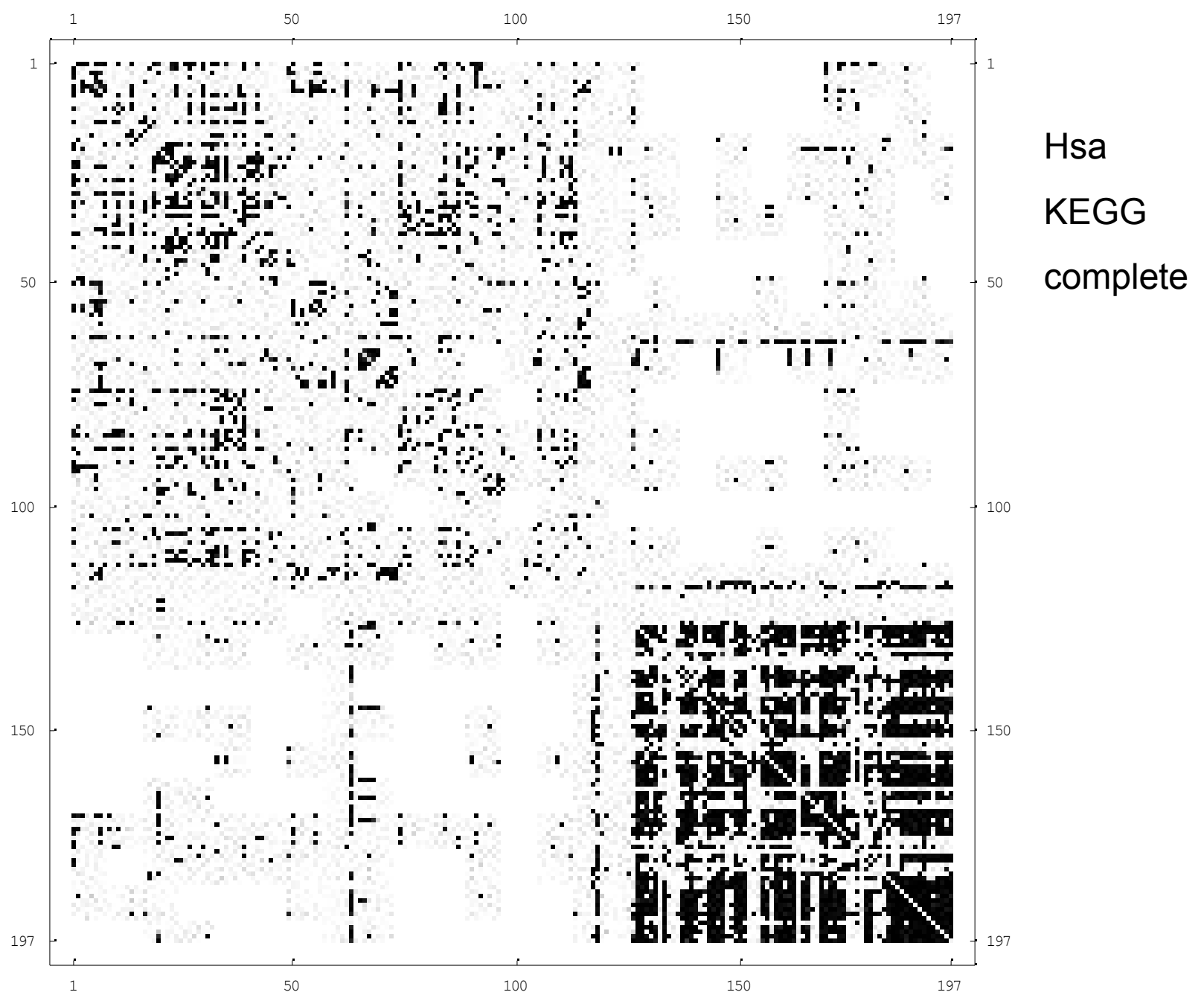
# Pathway network analysis

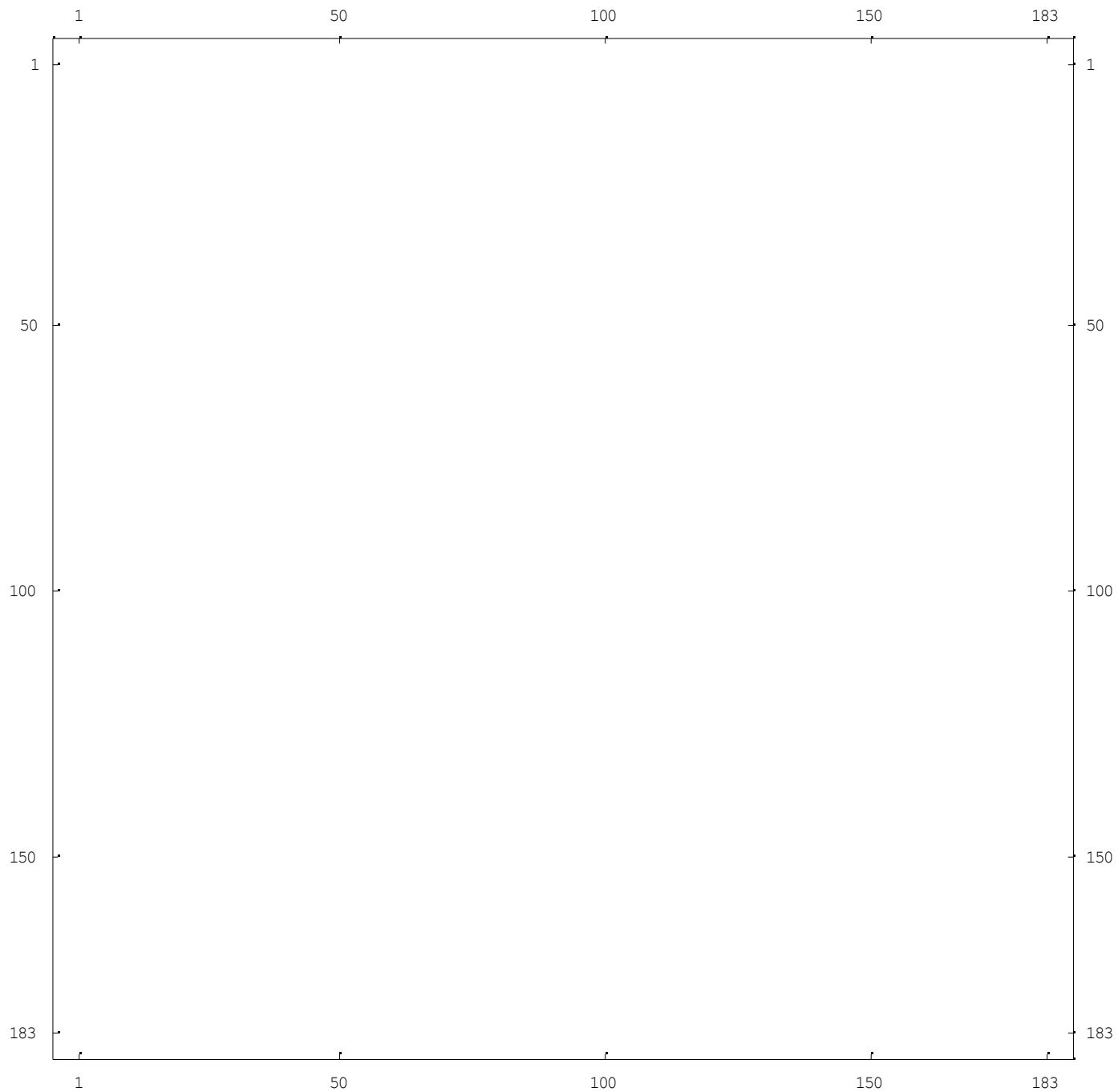Given significant nodes and edges, the **pathway network** can be **reconstructed**.

Edges and nodes can be **ranked** based on their **centrality** in the network (connectivity degree and betweenness)

# Betweenness centrality

Betweenness centrality is a very interesting parameter because:

- it can be calculated both for nodes and edges

- it is a measure of the possible **information flow** through that element, thus if it is affected by experimental conditions it is very likely that such **perturbation** can **spread** to the whole system more easily

Hsa

KEGG

complete

rno

KEGG

complete

Histogram of betweenness centrality
of pathways extracted from KEGG hsa

Plot of betweenness centrality
of pathways extracted from KEGG hsa

| | | |
|---|---:|---|
| 0.06869 | 7 | Galactose metabolism |
| 0.053446 | 169 | Insulin signaling pathway |
| 0.049498 | 20 | Purine metabolism |
| 0.043014 | 39 | Tryptophan metabolism |
| 0.039993 | 33 | Tyrosine metabolism |
| 0.039176 | 62 | Glycerolipid metabolism |
| 0.032689 | 176 | Alzheimer's disease |
| 0.031585 | 17 | Androgen and estrogen metabolis |
| 0.031433 | 173 | Type II diabetes mellitus |
| 0.02946 | 1 | Glycolysis / Gluconeogenesis |
| 0.029339 | 191 | Prostate cancer |
| 0.022151 | 24 | Glycine, serine and threonine me |
| 0.021969 | 172 | Adipocytokine signaling pathway |
| 0.020961 | 126 | PPAR signaling pathway |
| 0.020138 | 22 | Glutamate metabolism |
| 0.019782 | 30 | Lysine degradation |
| 0.018842 | 87 | Butanoate metabolism |
| 0.01853 | 96 | Nicotinate and nicotinamide meta |
| 0.018316 | 50 | Starch and sucrose metabolism |
| 0.018112 | 115 | Glycan structures - biosynthesis |

Top 20 pathways
 extracted
from KEGG Database
ranked for their
betwennes centrality

# Pathway significance analysis

Node (pathway) or edge (intersection) significance analysis can be performed by considering the total number of genes represented in KEGG and the total number of statistically significant genes, compared with the significant genes found in a node or edge and their total number of elements (e.g. by a test based on the **hypergeometric distribution**)

## 2.4 Fisher's exact test (Draghici *et al.*, 2003)

We consider that there are $N$ single-symbol-annotated genes on the microarray (replicates were averaged by calculating the mean), which are either significantly differentially expressed ($S$) or not ($F$), and either belong to a pre-defined pathway list ($P$) or not ($NP$), see Table 2. If we pick randomly $P$ genes, we would like to estimate the probability of having exactly $\alpha$ genes in $S$. The $p$-value of having $\alpha$ genes or fewer in $S$ can be calculated by summing the probabilities of a random list of $K$ genes having $1, 2, \ldots, \alpha$ genes in $S$:

$$p = 1 - \sum_{i=0}^{\alpha} \frac{\binom{S}{i}\binom{F}{P-i}}{\binom{N}{P}} \tag{1}$$

This is a one-sided test in which the $P$ values correspond to over-represented lists of genes.

A review about similar current tools used for group testing on the level of Gene Ontology (GO) terms was given by Khatri and Draghici (2005).

|        | 0   | 1   | Totals |
|--------|-----|-----|--------|
| **1**  | a   | b   | a+b    |
| **0**  | c   | d   | c+d    |
| **Totals** | a+c | b+d | n  |

$$\mu_{ij} = \frac{T_{R_i} \times T_{C_j}}{T_G}$$

Null table is constructed

by the multinomial

distribution and then

tested by a $\chi^2$ test

Fisher exact test for a 2x2 contingency table

|  | 0 | 1 | Totals |
|---|---|---|---|
| **1** | a | b | a+b |
| **0** | c | d | c+d |
| **Totals** | a+c | b+d | n |

The probability
Is due by the
Hypergeometric
distribution

$$\frac{\dfrac{(a+c)!}{a!c!} \times \dfrac{(b+d)!}{b!d!}}{\dfrac{n!}{(a+b)!(c+d)!}}$$

# Pathways and their intersections significance analysis

- calculated considering the hypergeometric distribution:

  $p(x) = choose(m, x)\ choose(n, k-x)\ /\ choose(m+n, k)$

- where
  - p= probability.
  - x = number of significant probes in a pathway (or intersection)
  - m = total number of significant probes.
  - n = total number of non significant probes.
  - k = number of probes in a pathway.
- P <0.05 was considered as significant

# Network representation

- Significantly underrepresented: (-1)
- Significantly overrepresented: 1
- Not significant: 0
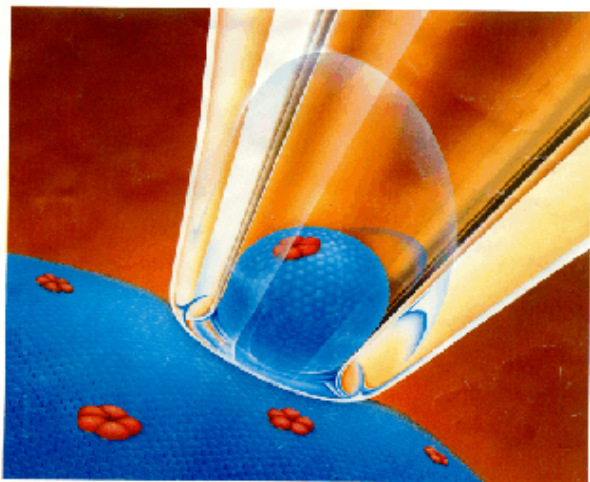
c-Myc off

# c-Myc on

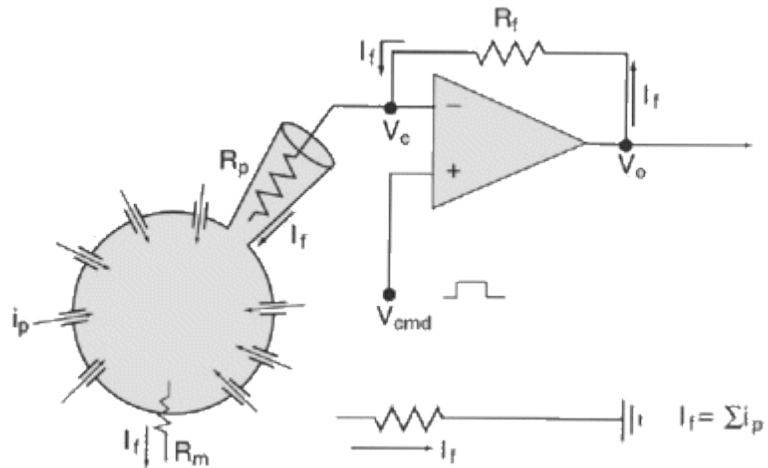# cardiac hypertrophy
# 2 weeks
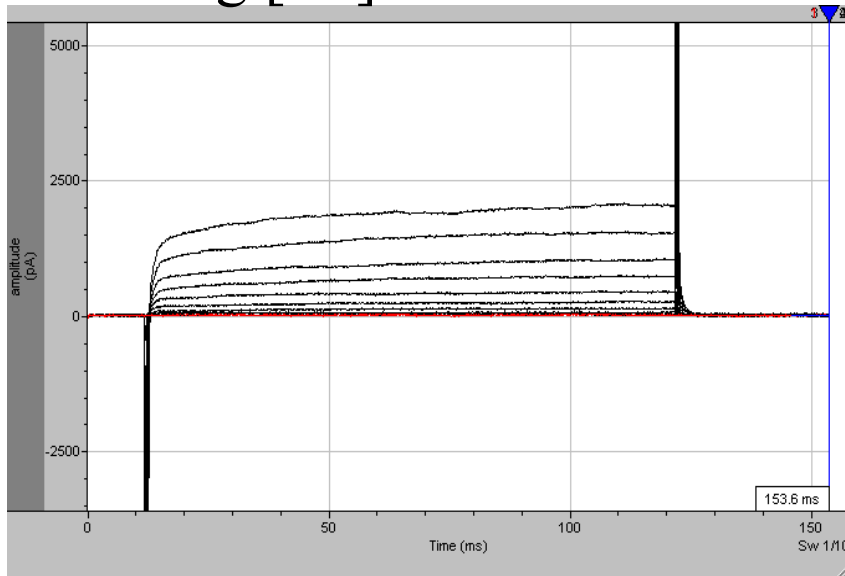
cardiac hypertrophy
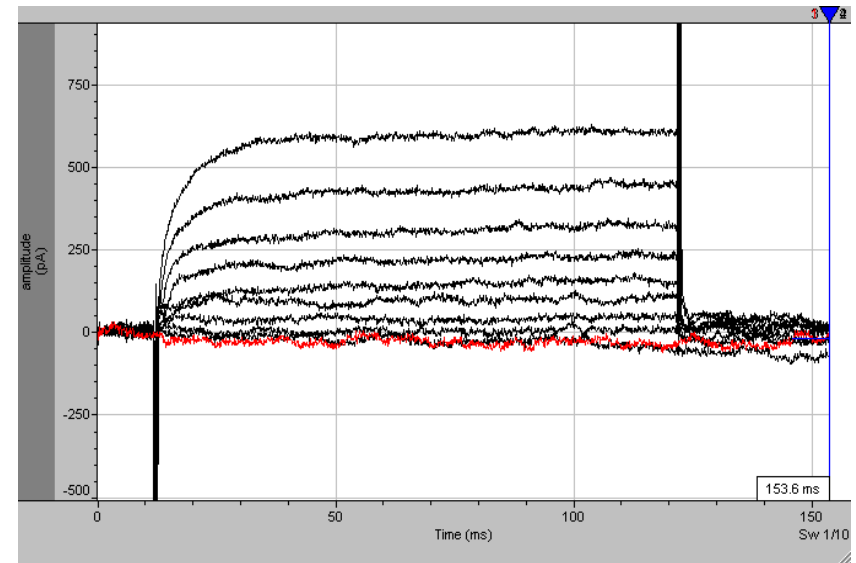4 weeks
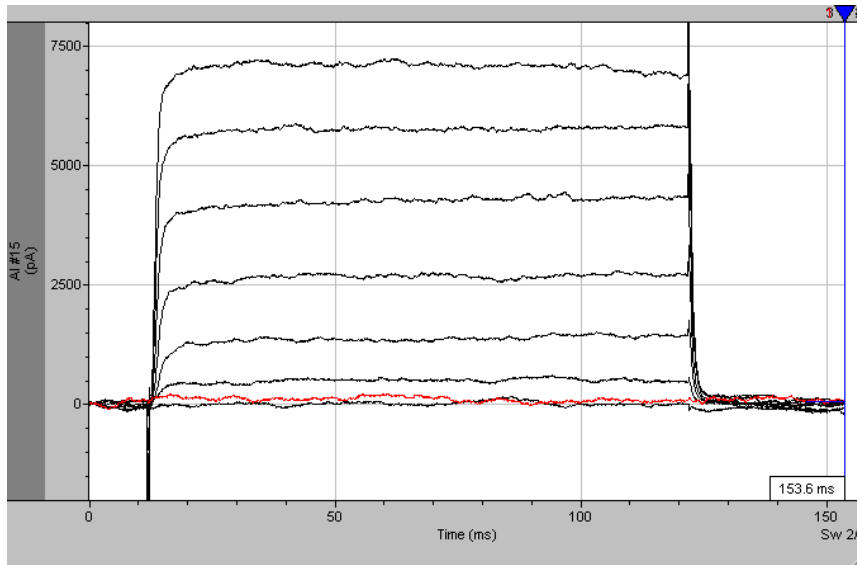
# CONCLUSIONS AND PERSPECTIVES

Young [Ca] 1 mm

Centenarian [Ca] 1 mm

Young [Ca] 10 mm

Centenarians [Ca] 10 mm