

Data mining in genetics

Junior Barrera

BIOINFO-USP

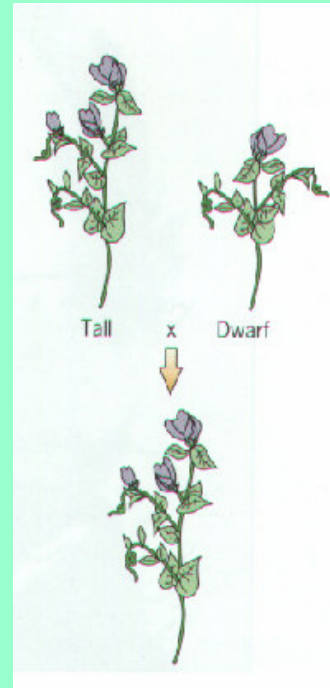
Layout

- Introduction
- Data mining
- Mapping of rare genes
- **Expression analysis:** measure, genes differentially expressed; clustering expression signals; identification of gene regulation networks

Introduction

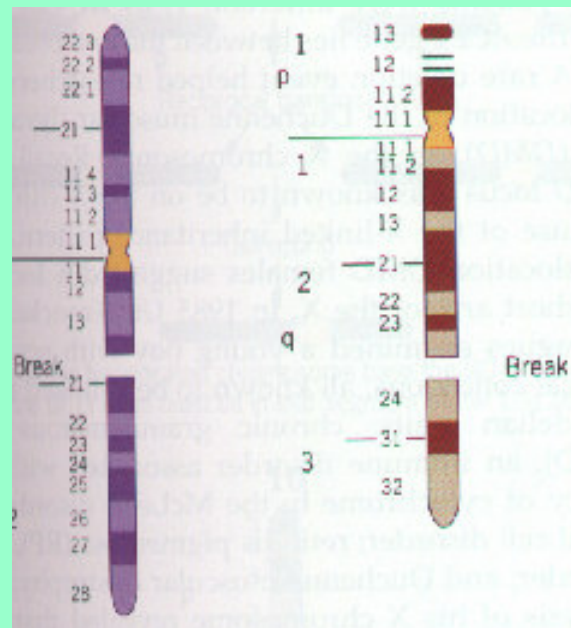
Knowledge evolution in genetics

- Heredity - Mendel (1866)
- The phenotypes of an individual depends on genes of his parents.



Knowledge evolution in genetics

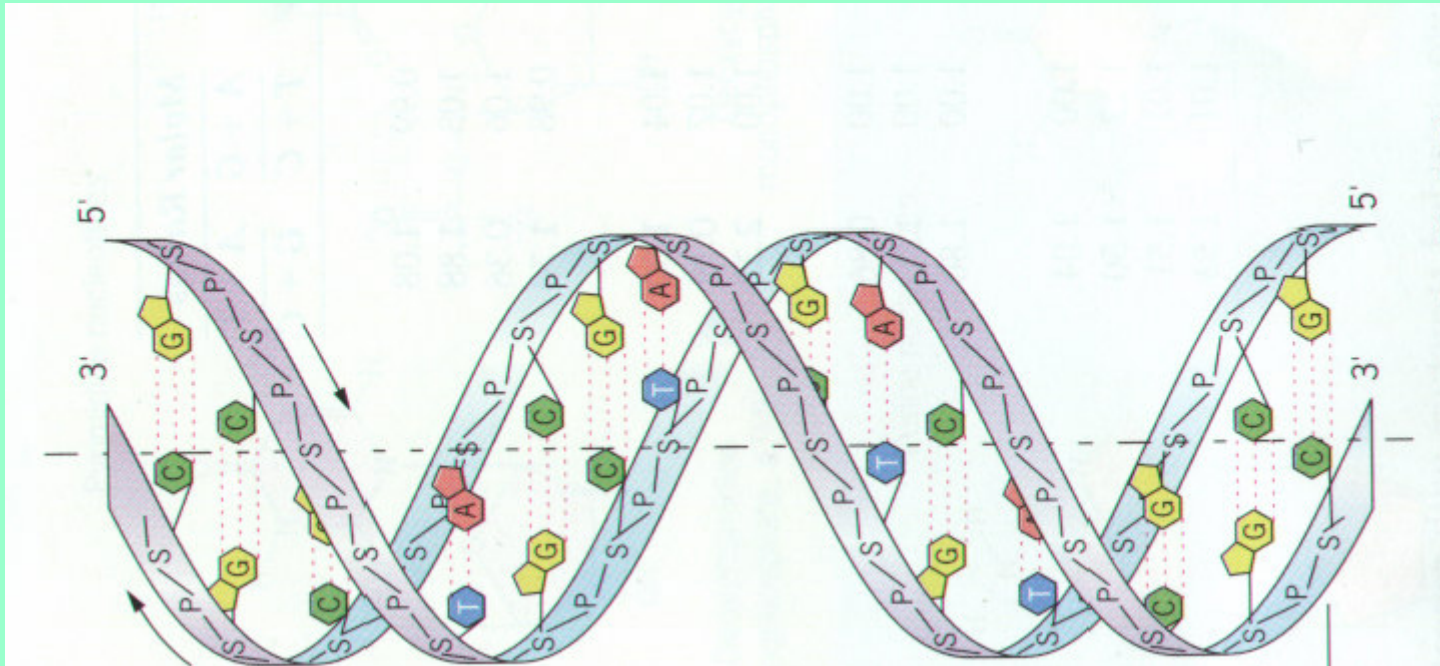
- Chromosome Theory - Morgan (1910)
- Genes were situated in chromosomes



Knowledge evolution in genetics

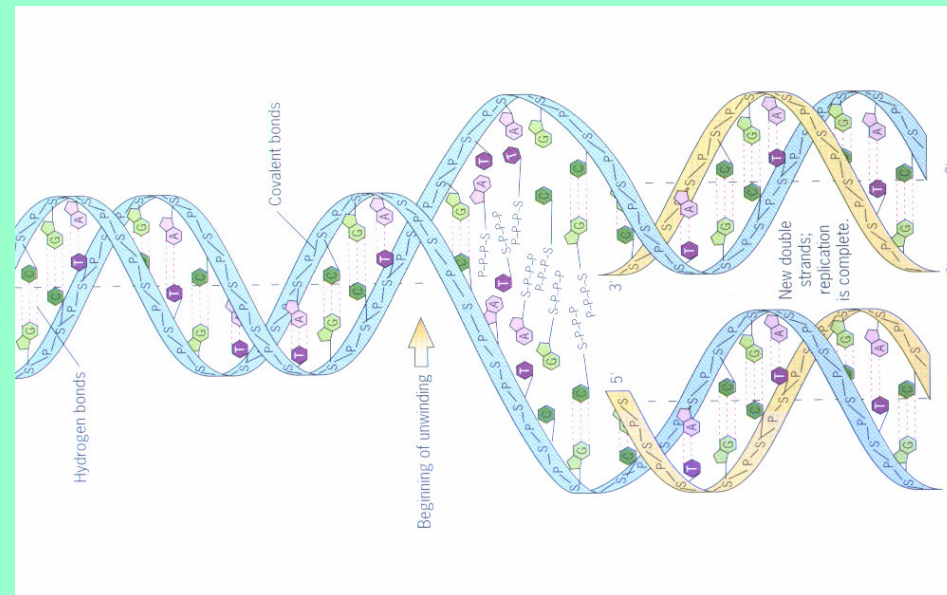
- The molecular structure of chromosomes
(Watson and Crick - 1953)
- DNA structure: the double helix
- Four basis: adenine(A), guanine(G),
thymine(T), cytosine(C)
- genes are sequences of nucleotides

Knowledge evolution in genetics



Knowledge evolution in genetics

- DNA manipulation
- cut, replication and decoding



Knowledge evolution in genetics

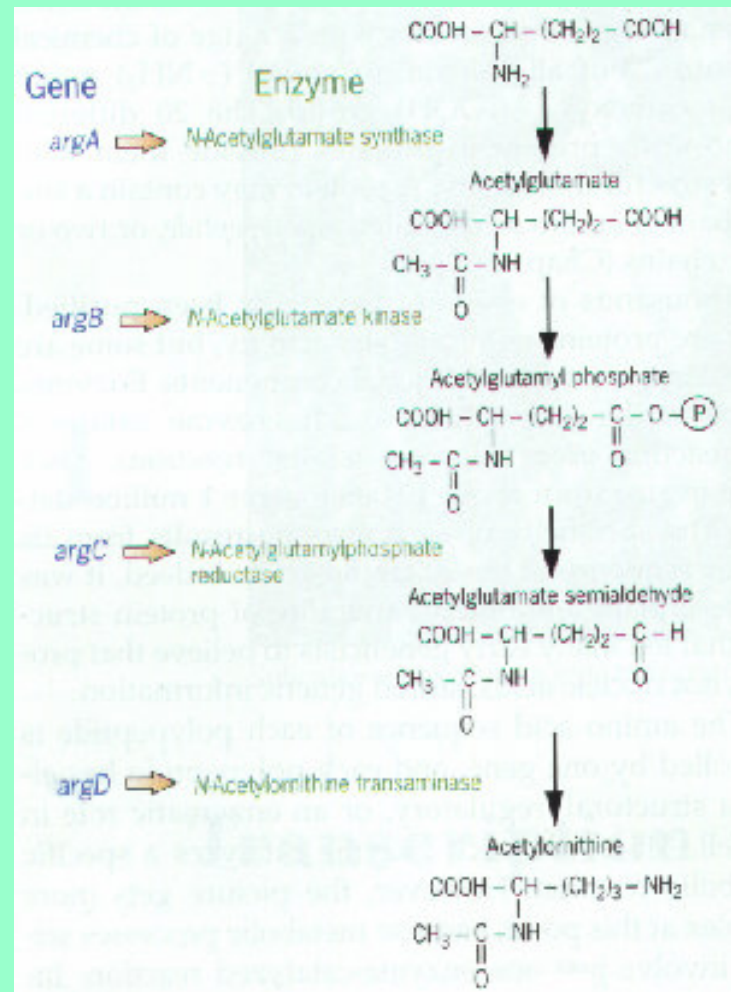
- Genetic engineering
- species modification, drug production



Knowledge evolution in genetics

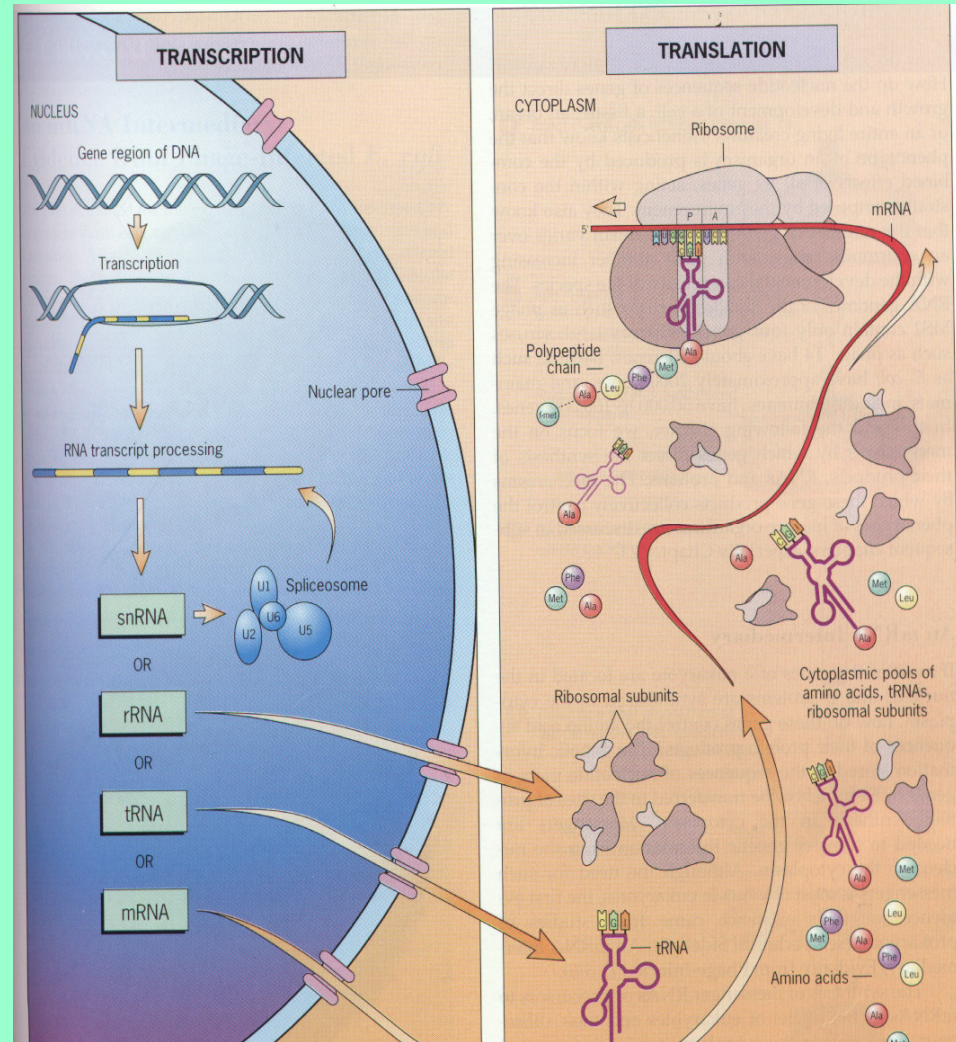
- Genes control the metabolism
- Metabolism occurs by sequences of enzyme-catalyzed reactions.
- Enzymes are specified by one or more genes

Knowledge evolution in genetics

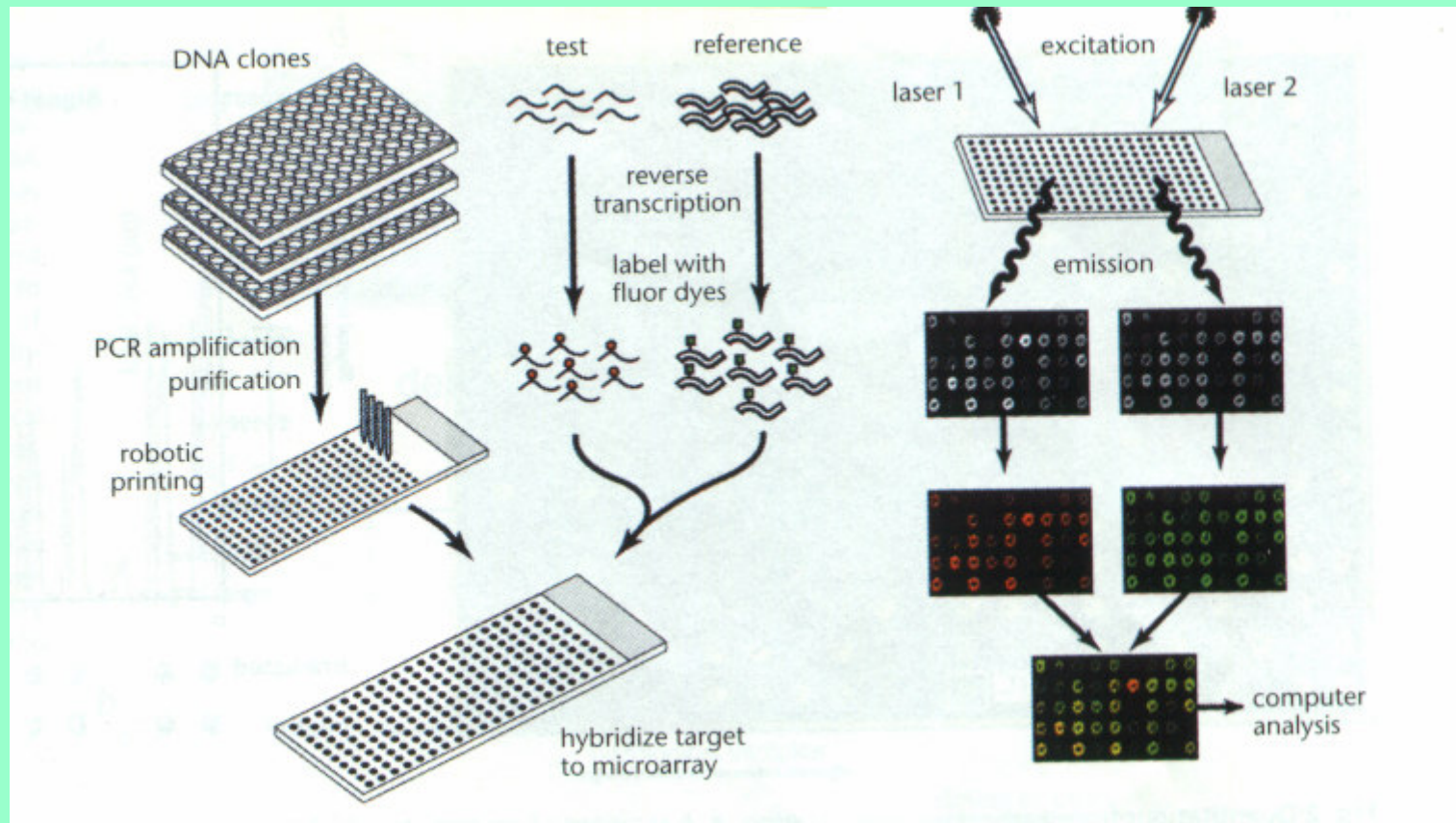


Knowledge evolution in genetics

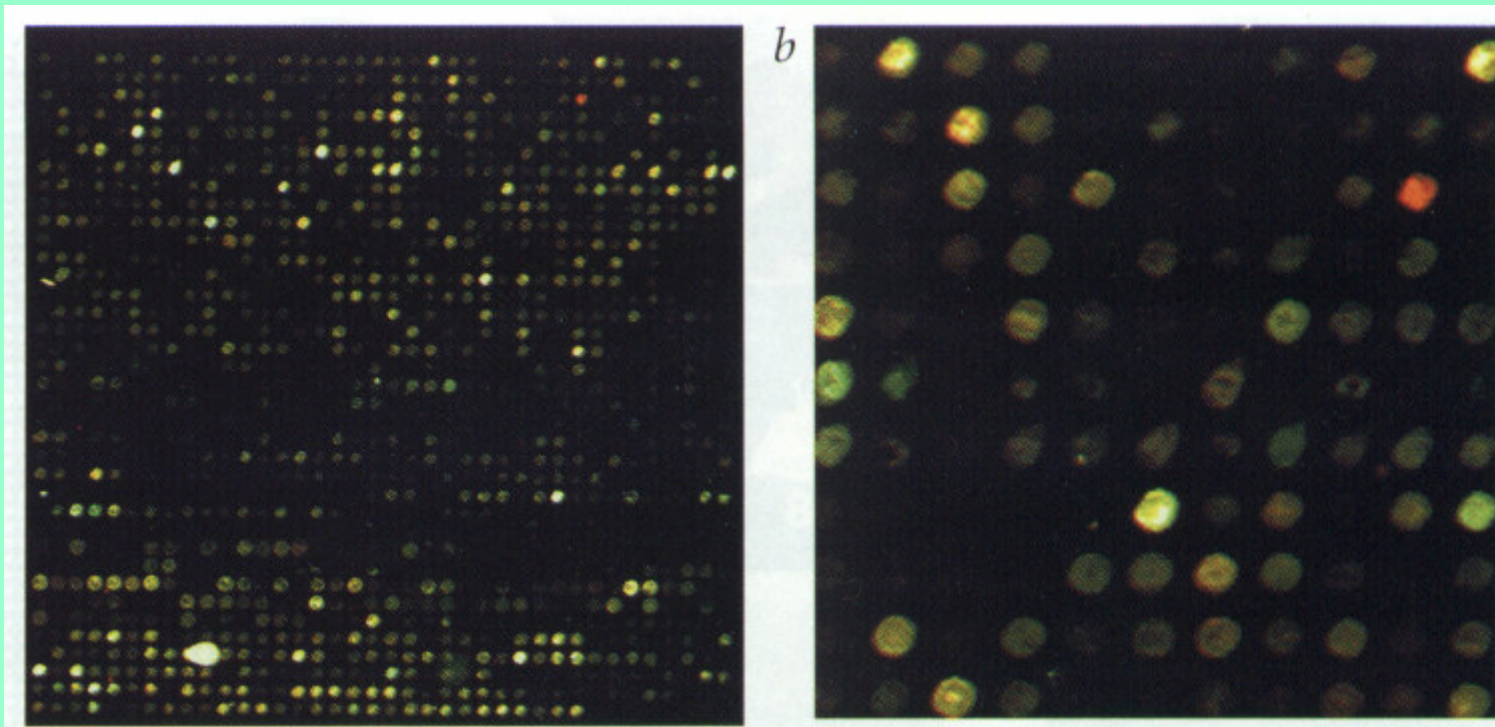
- Gene expression



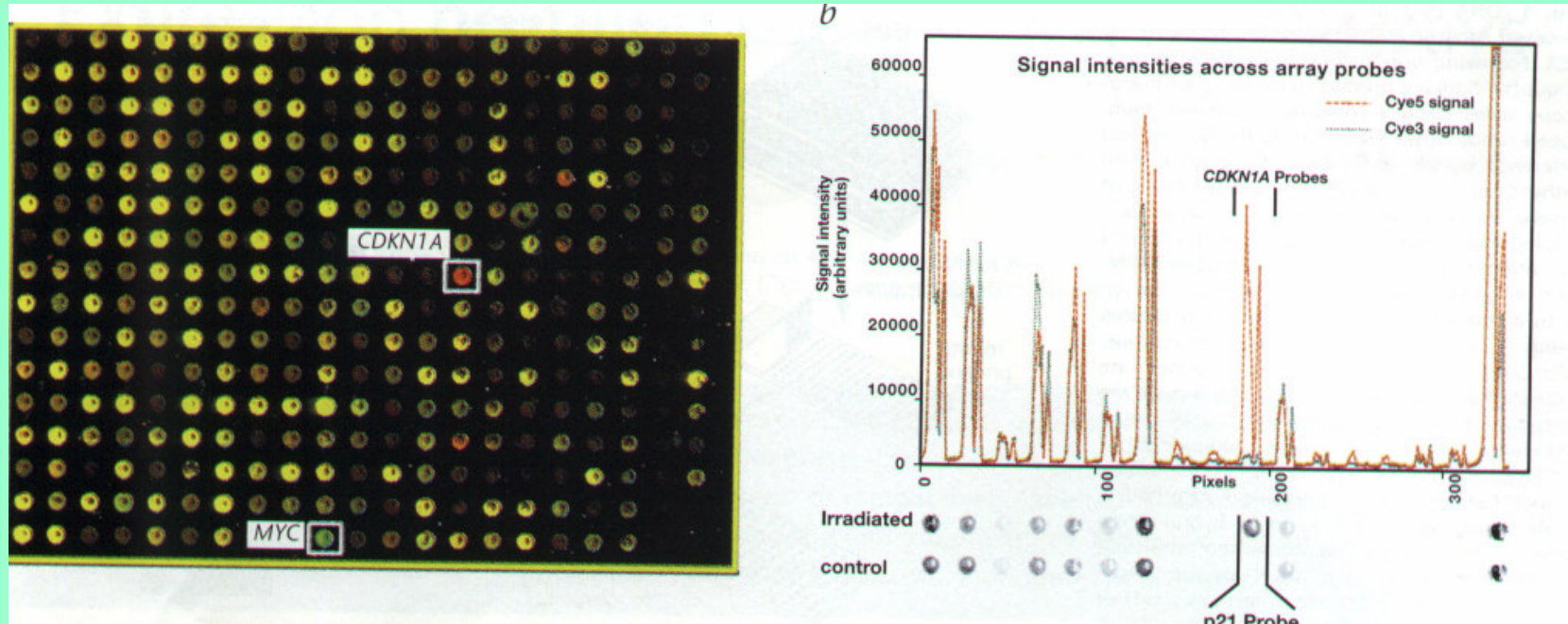
Data acquisition



Data acquisition



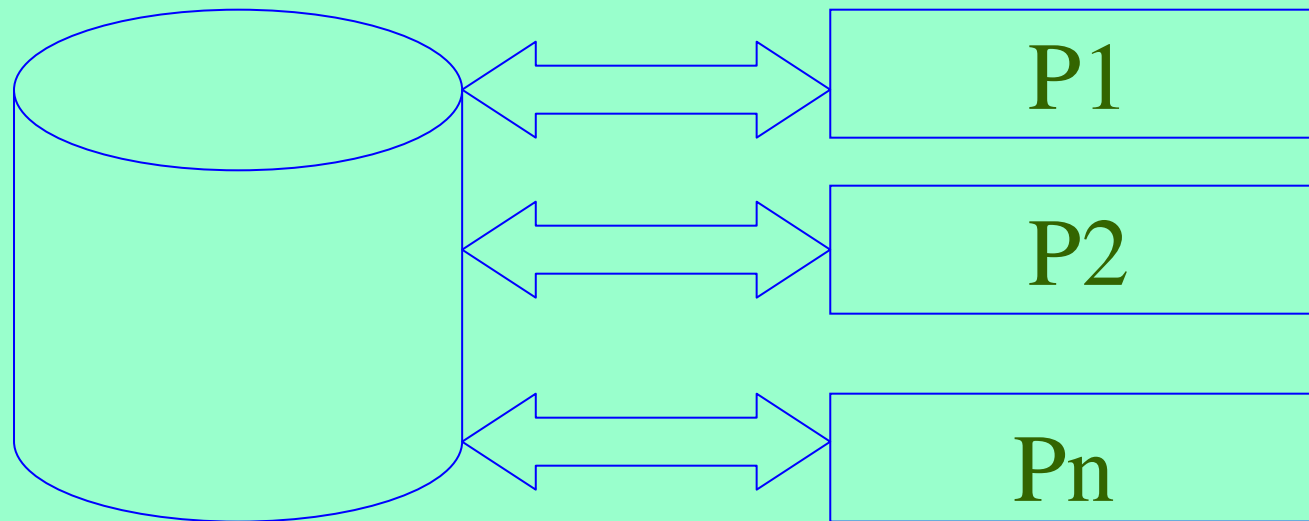
Data acquisition



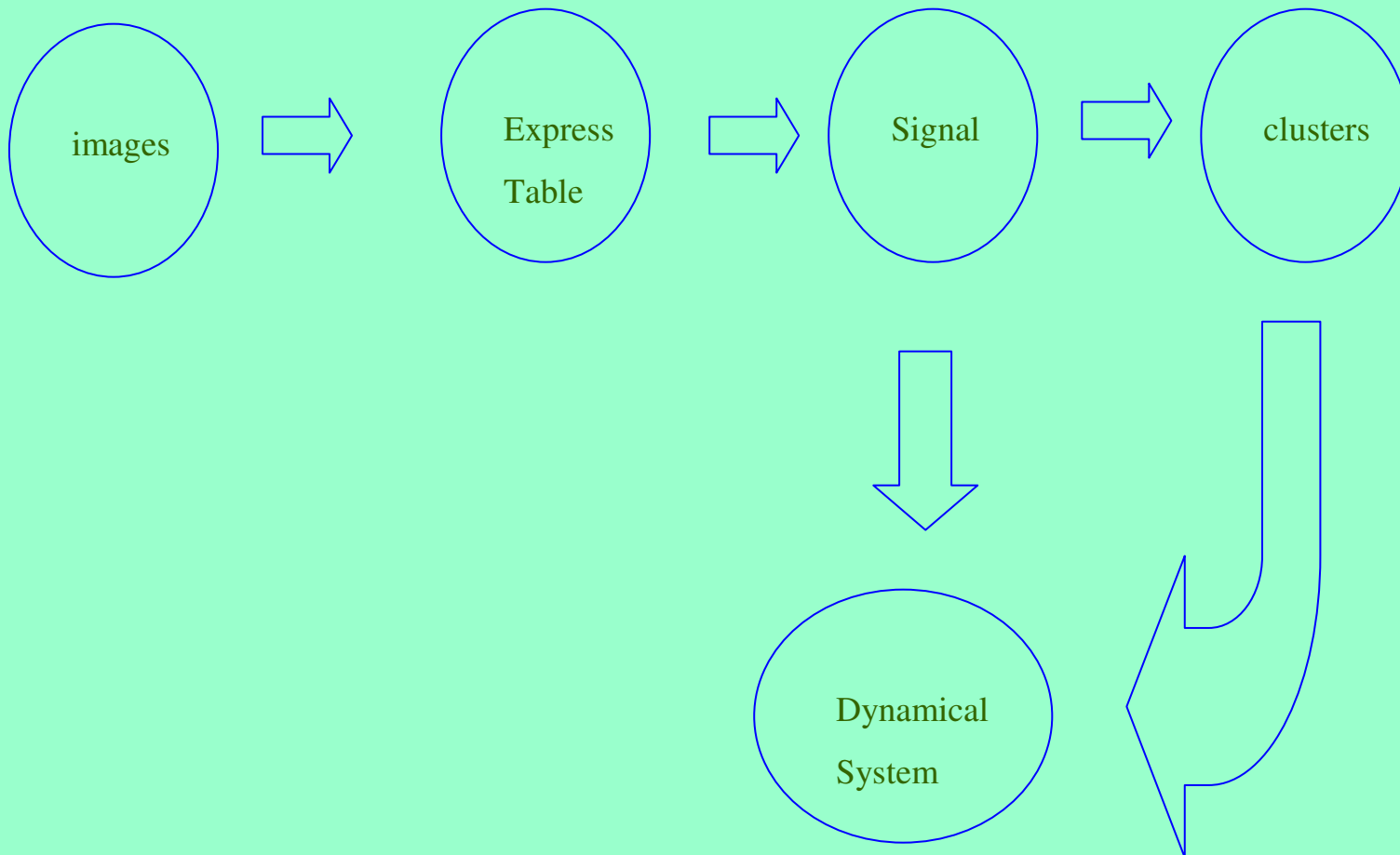
Quantization - $\{-1,0,1\}$

Data mining

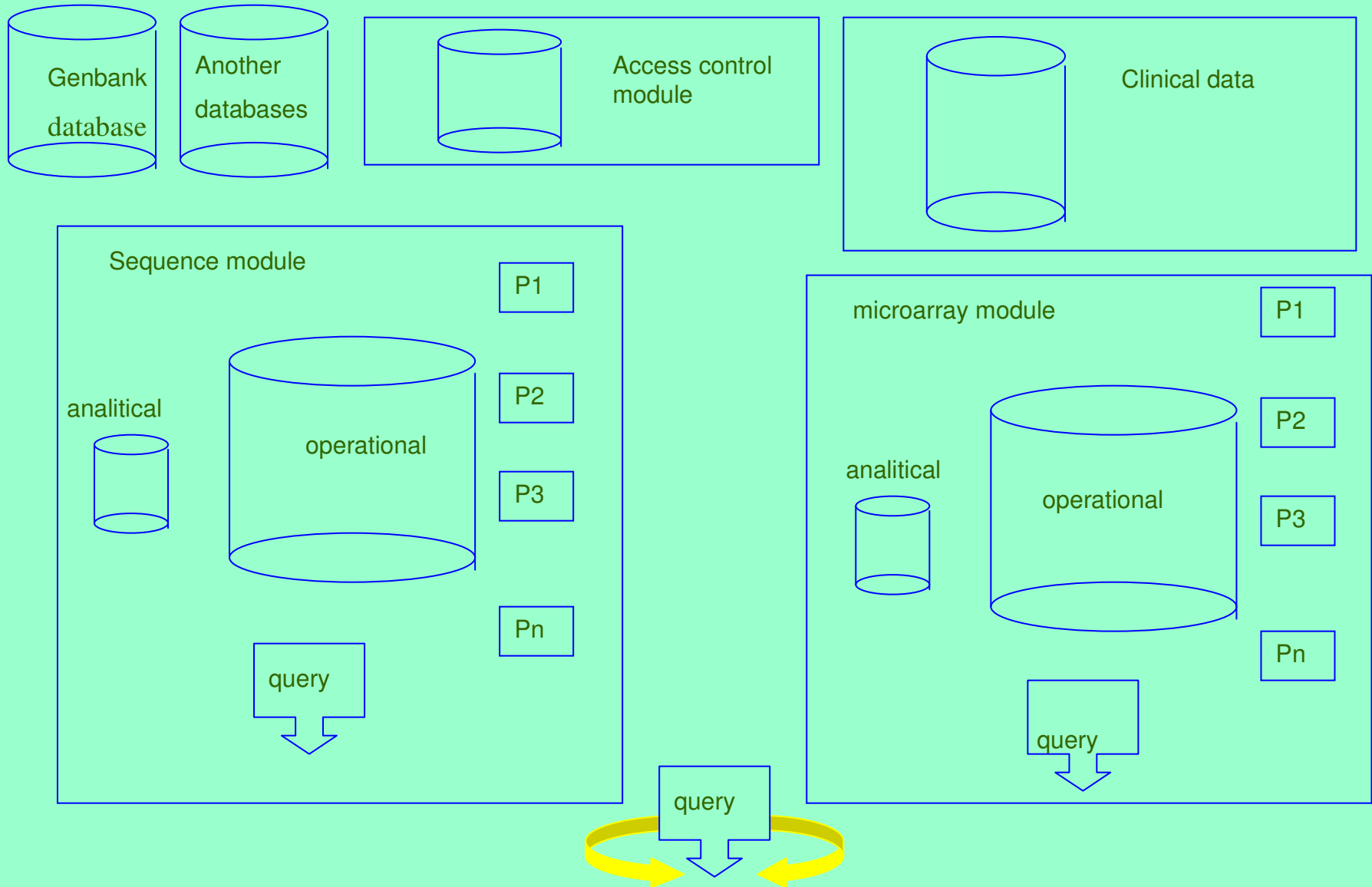
Object oriented database

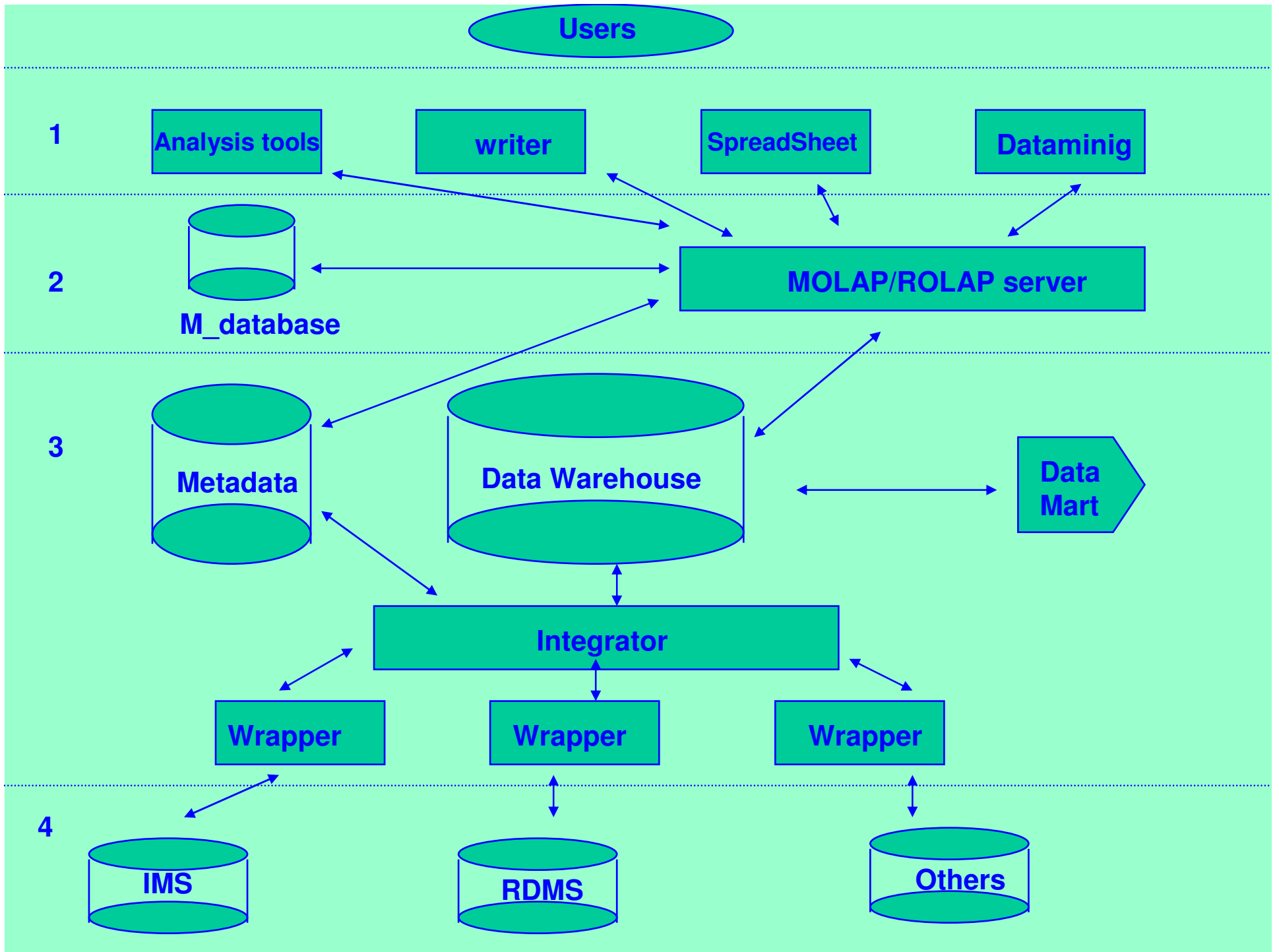


P_i : analytical and mining procedures (kernel parallel)

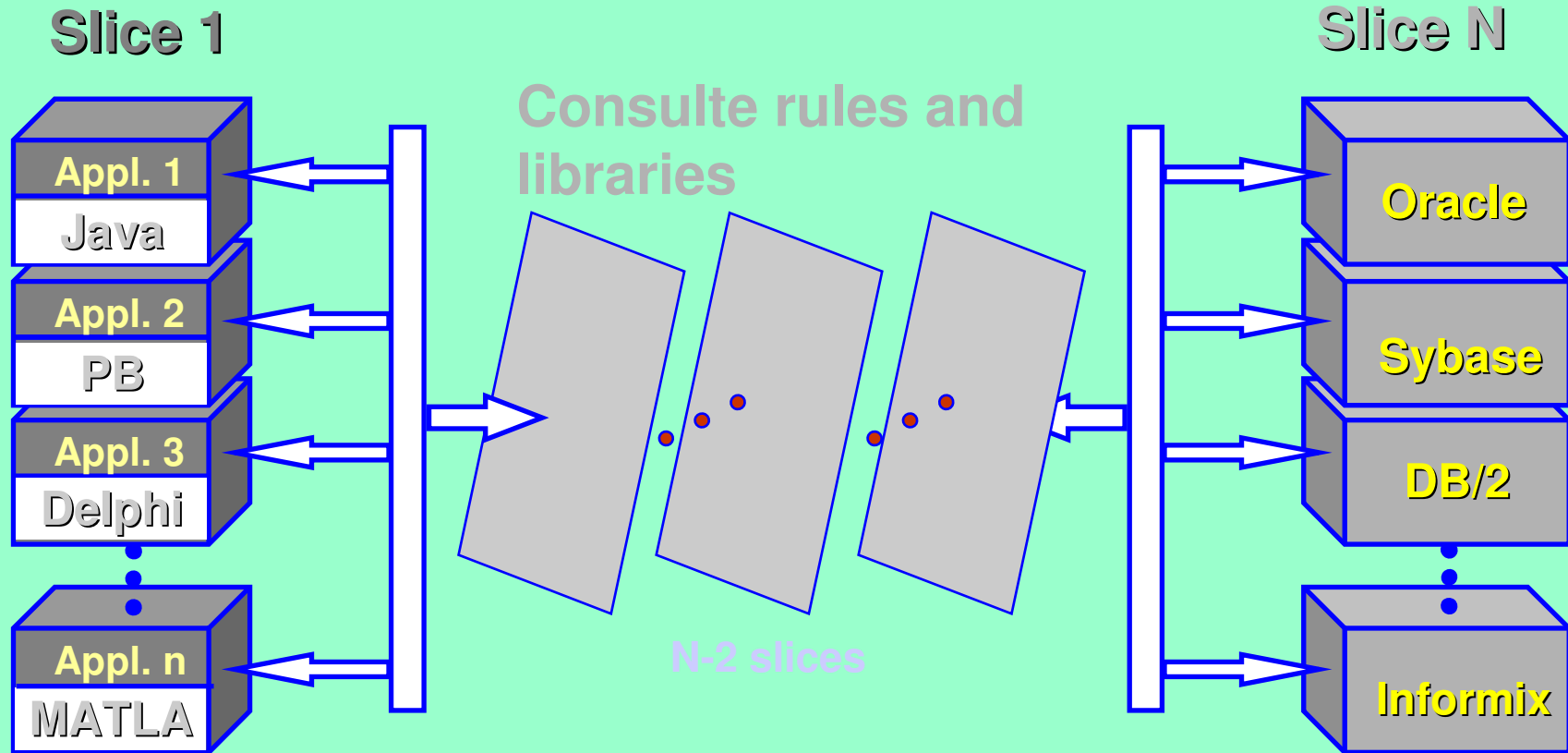


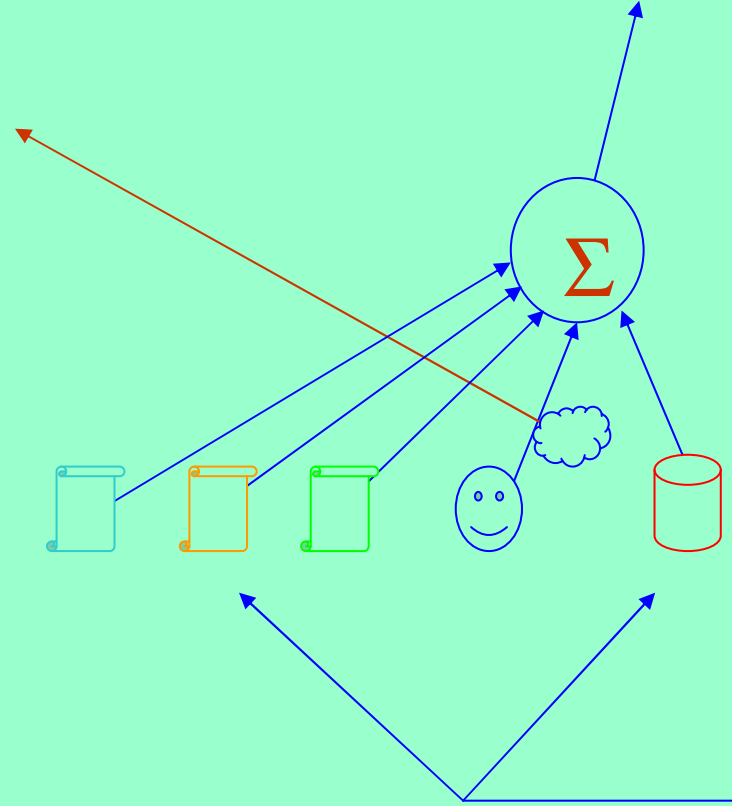
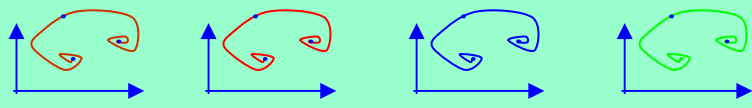
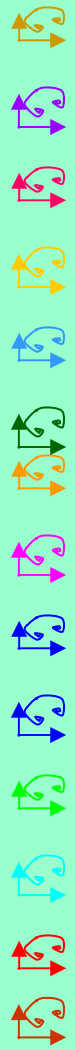
Integrated Environment





System Architecture

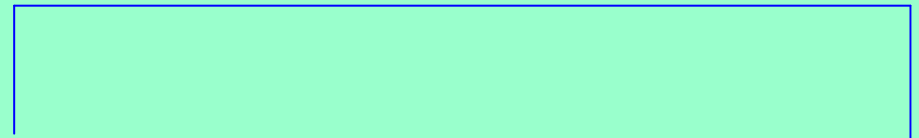
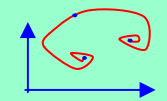




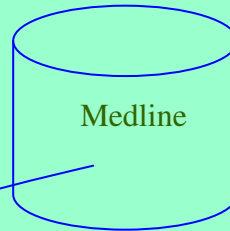
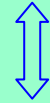
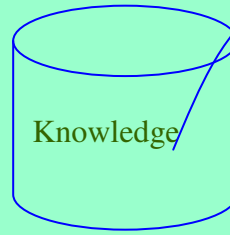
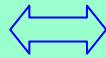
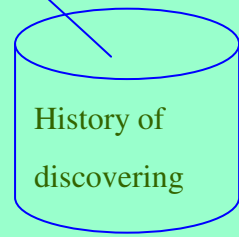
What genes regulate the pathway A->B->C->D ?

- Proteome
- Transcriptome
- Genome
- Pathways

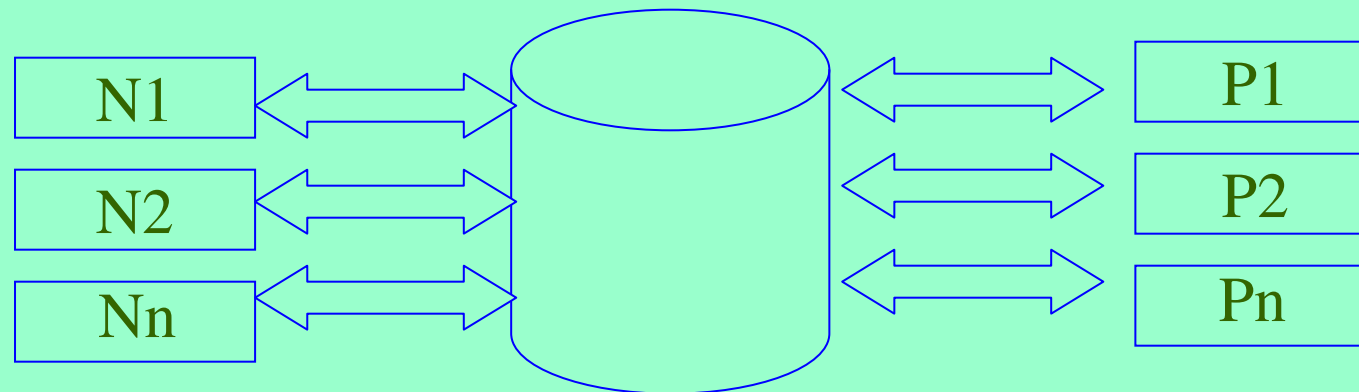
Wet Lab



Pi, Ni

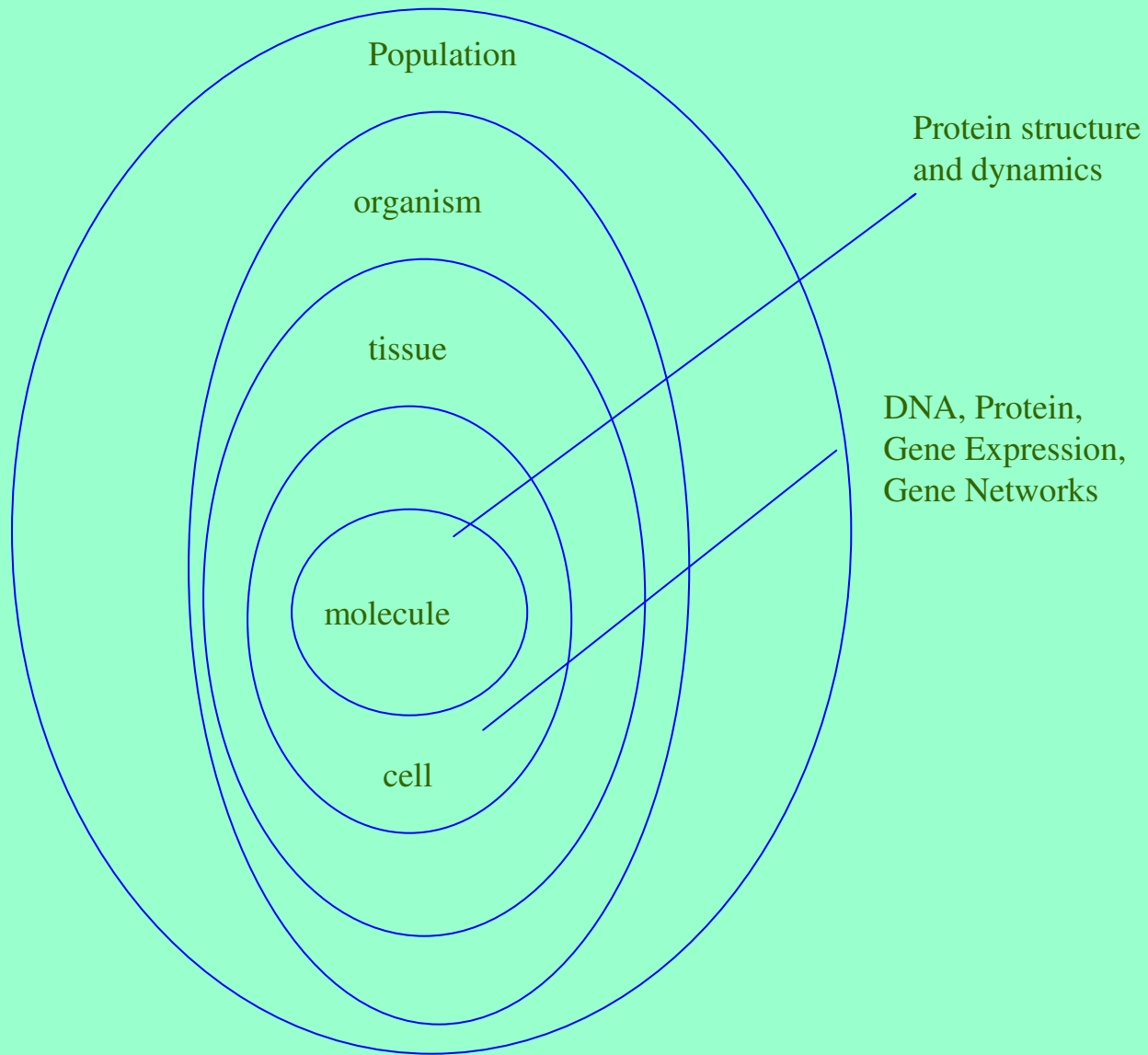


Object oriented database

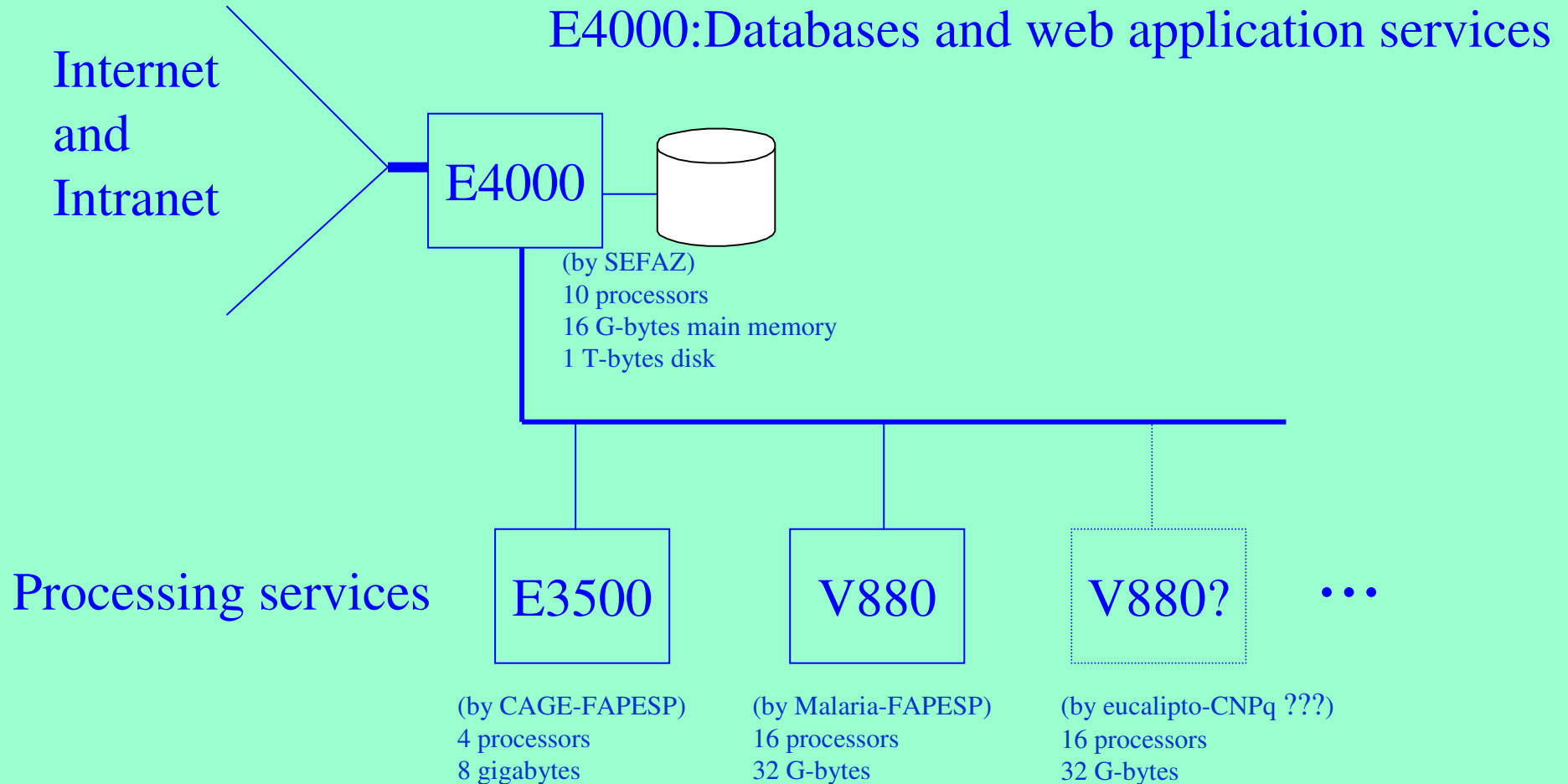


P_i : analytical and mining procedures (kernel parallel)

N_i : knowledge discovering procedures (kernel parallel)



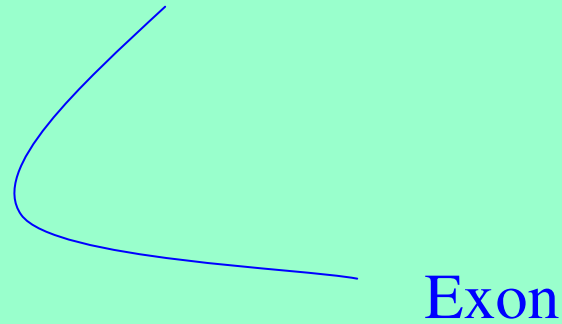
GRID Computer - DCC-IME-USP



Mapping of rare genes

ACGAATCTAGAGAATTAATTAACCGAGTTAAGA

ACGAATCTAGAGAATTAATTAACCGAGTTAAGA



Training from known genes

Expression Analysis

Analysis Phases

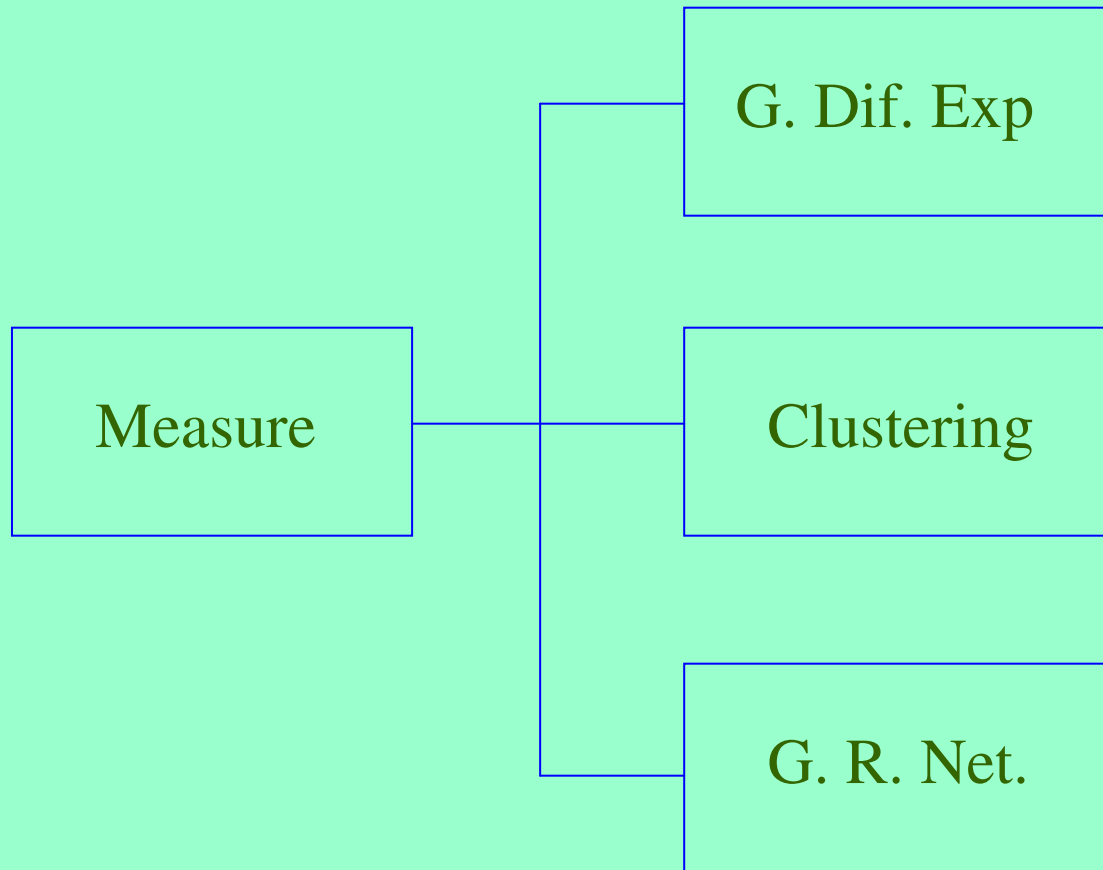


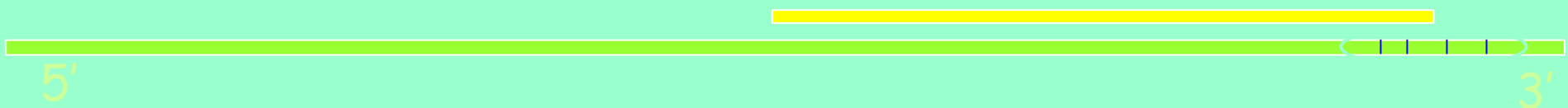
Image Analysis

Selection of clones

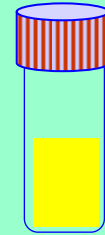
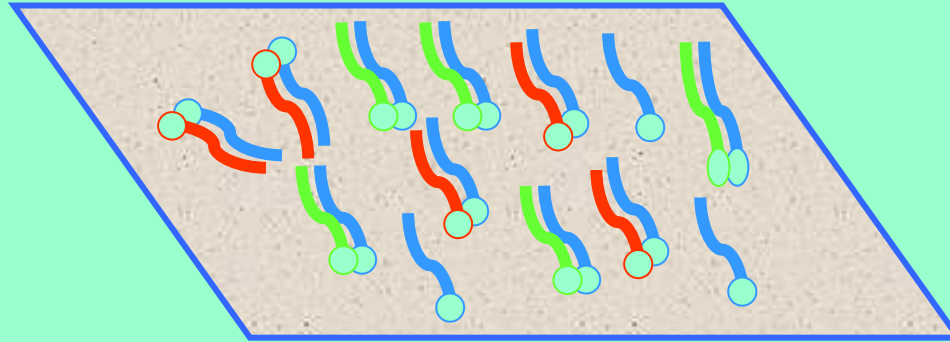
1. Clusters of the same gene



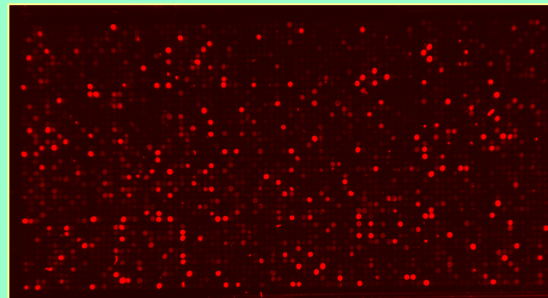
2. Choice of a representation for the gene



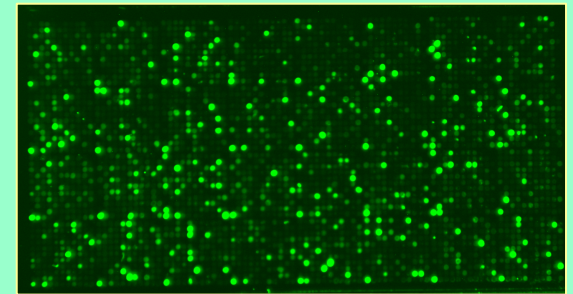
Hibridization



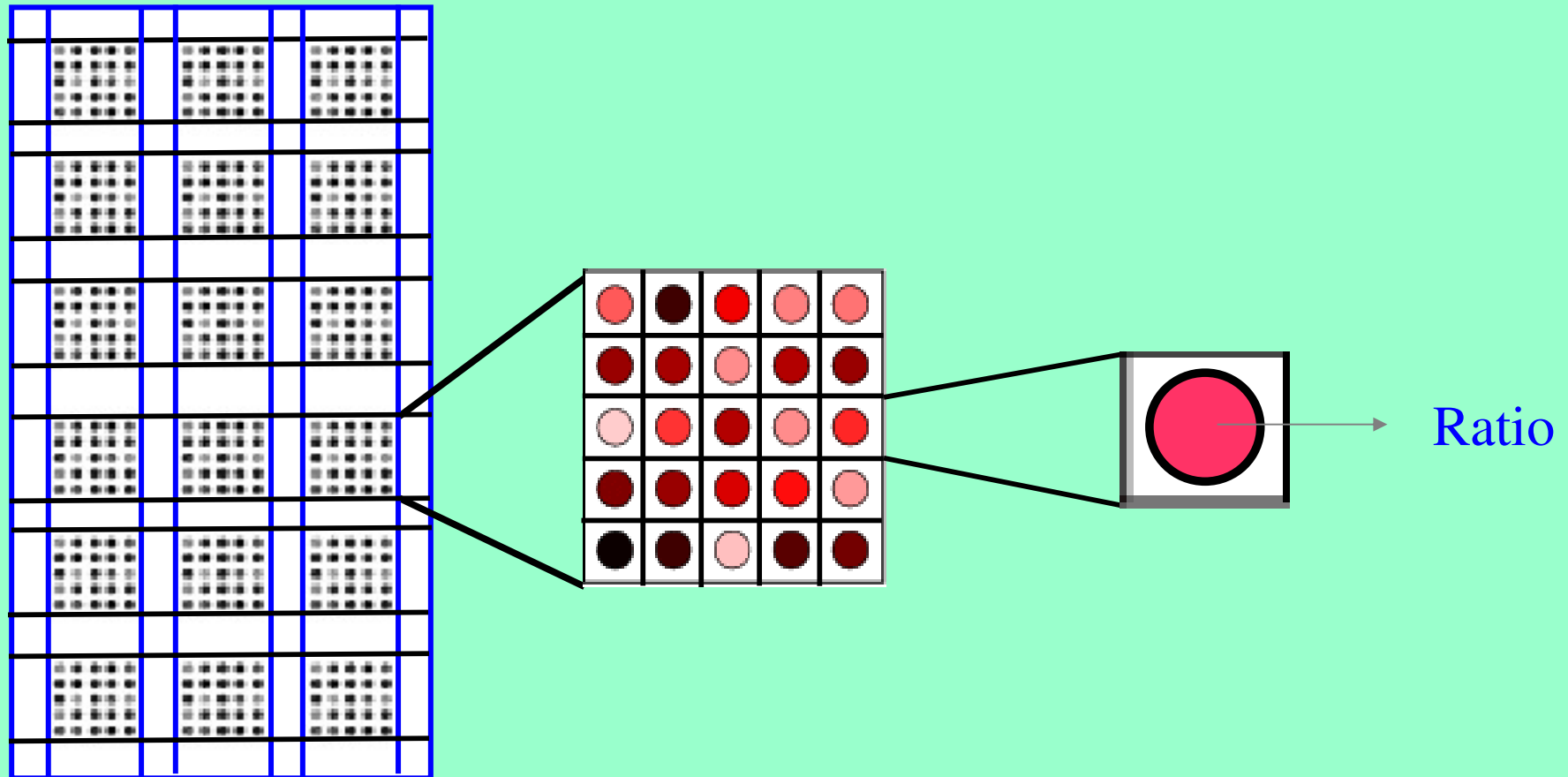
Cy5



Cy3



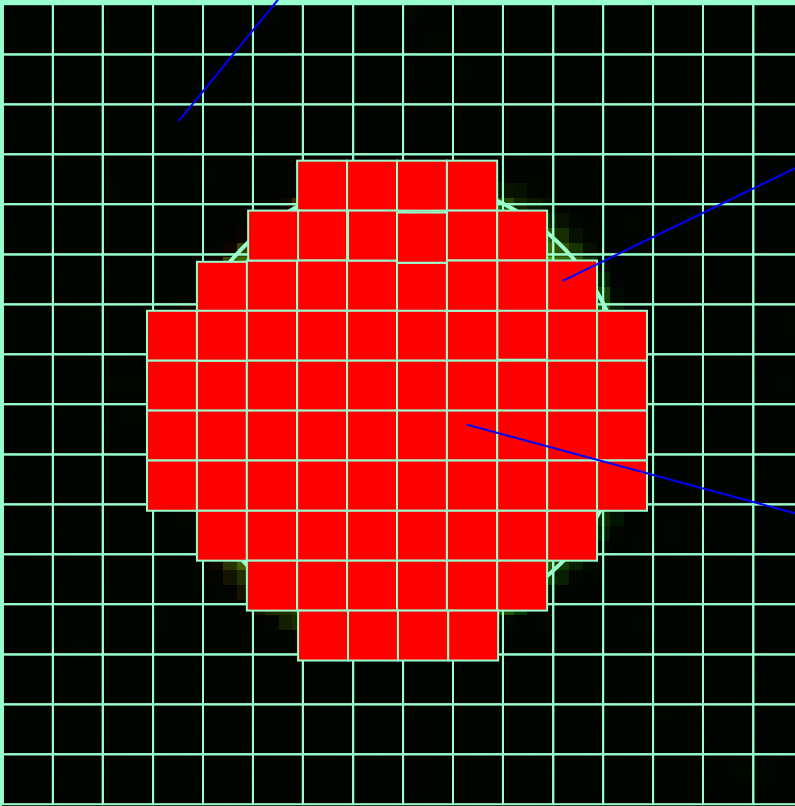
Expression Calculus



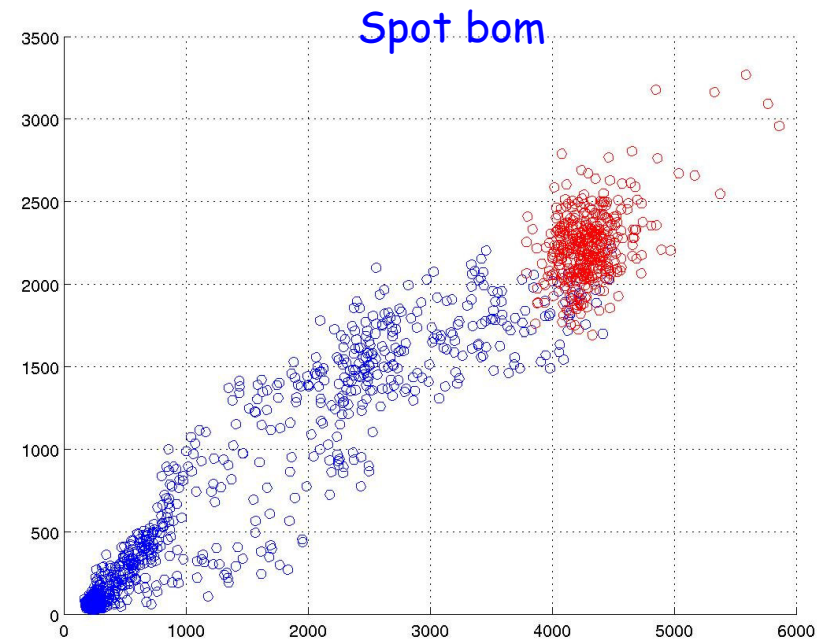
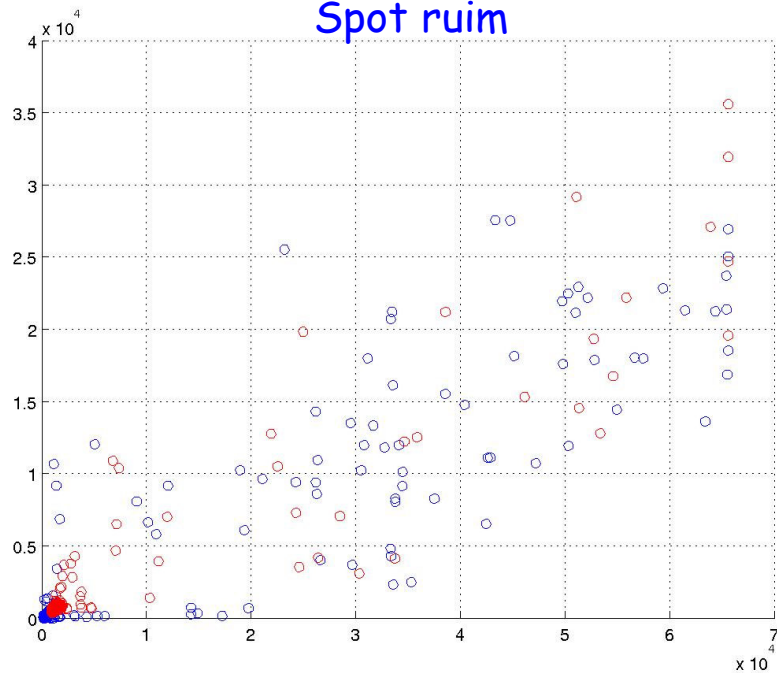
Background

signal + noise

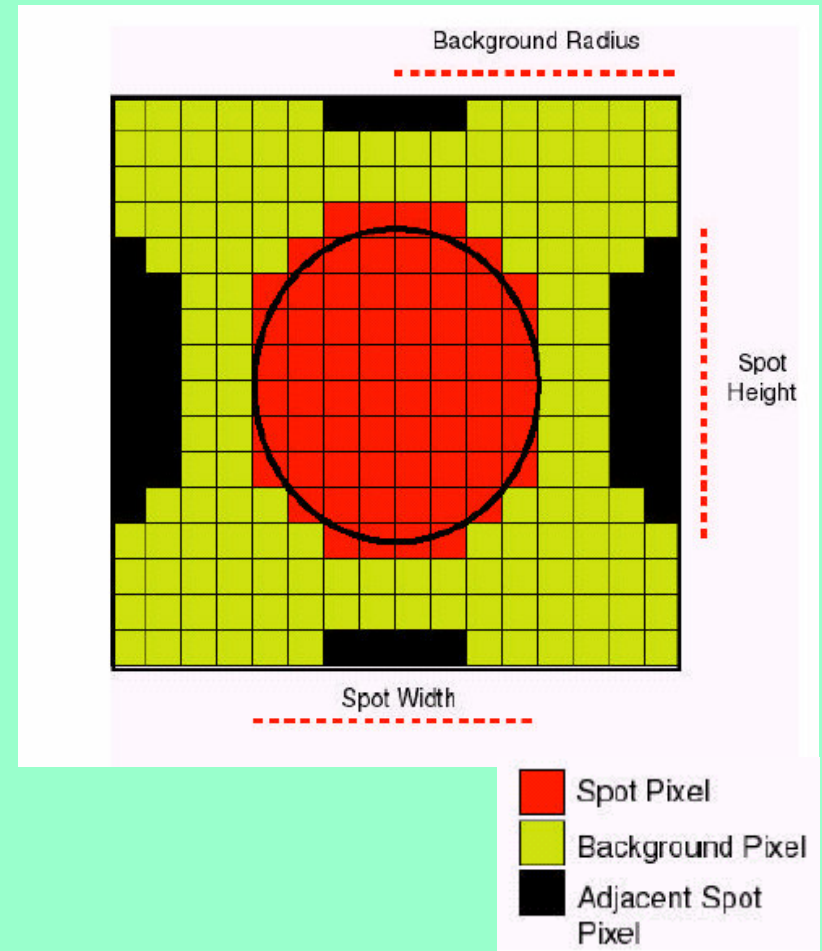
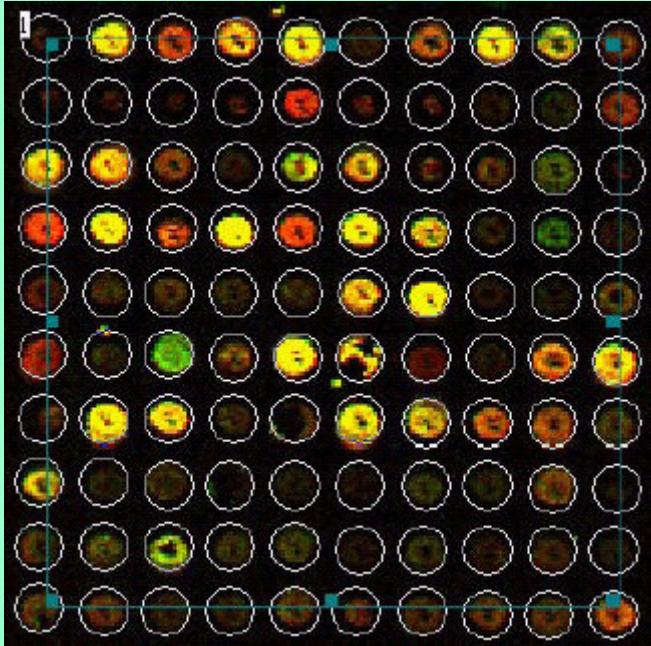
Foreground



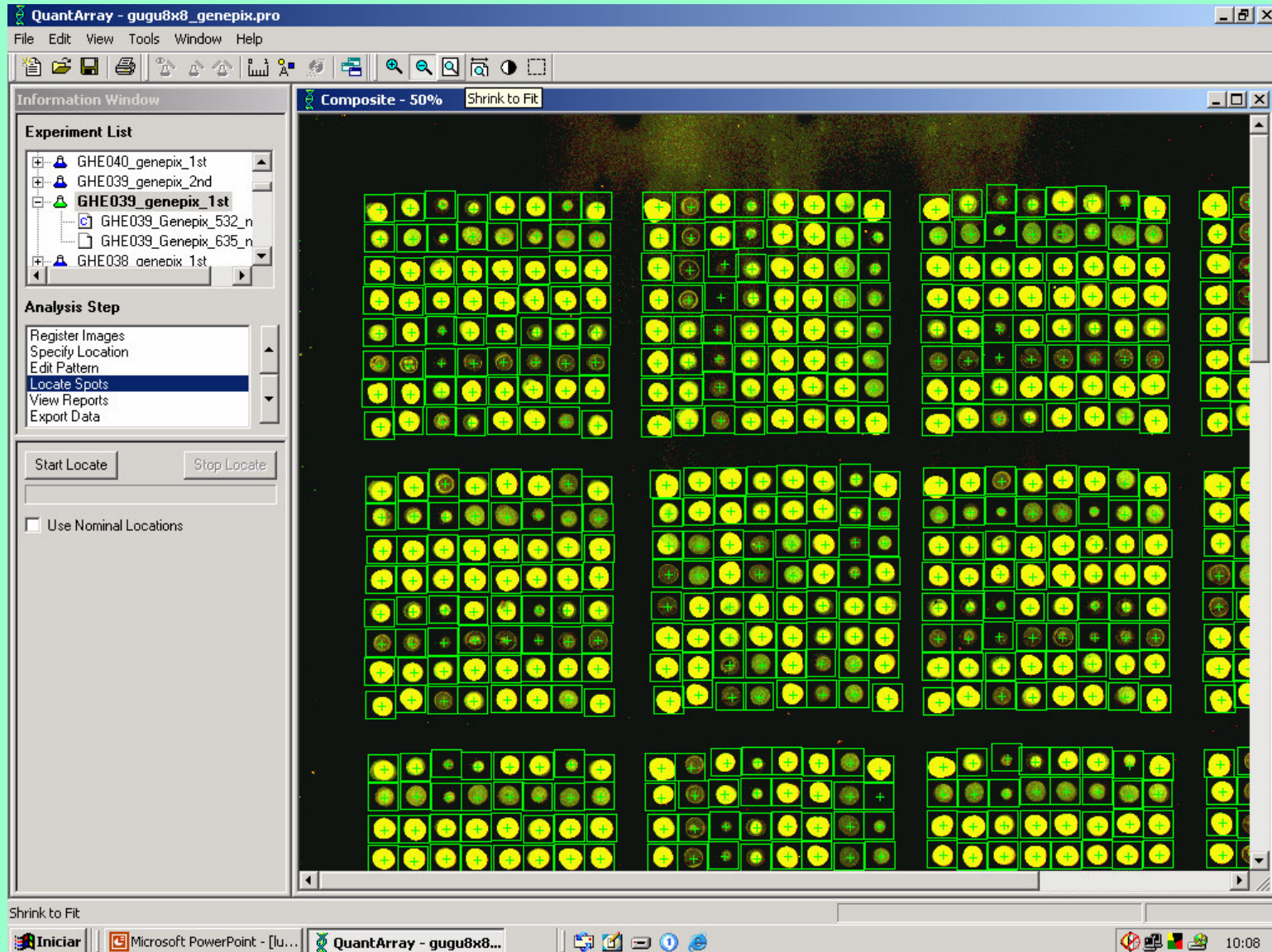
Dispersion of cy3 and cy5 give an idea of the spot quality



ScanAlyze



QuantArray®



Spot

Spot Control Panel

Spot Operation:
Place grids de novo

Output file prefix (suffixes will be appended):

Input file (for Hint):

De Novo + Hint

Dapi Image:

Test Image:

Reference Image:

Options
Array slop:

Layout adjustment options
Arr X (col): Arr Y (row): XTweak: YTweak:

Multi-Tweak

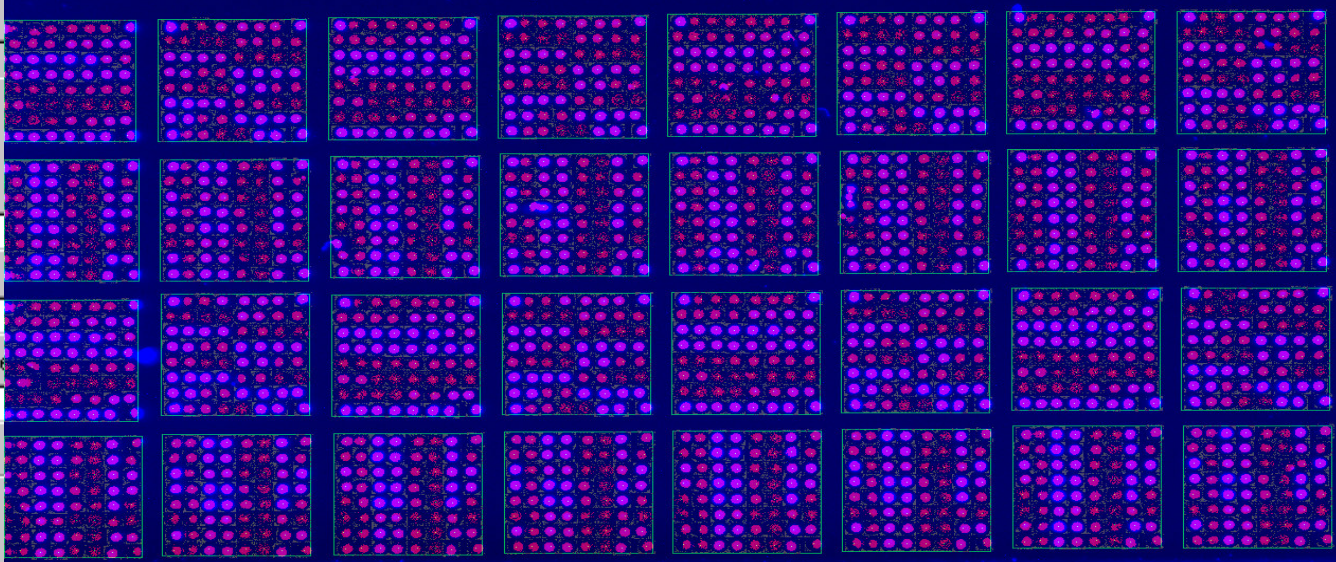
spacing hint: spacing hint:

Skip Hard Optimization

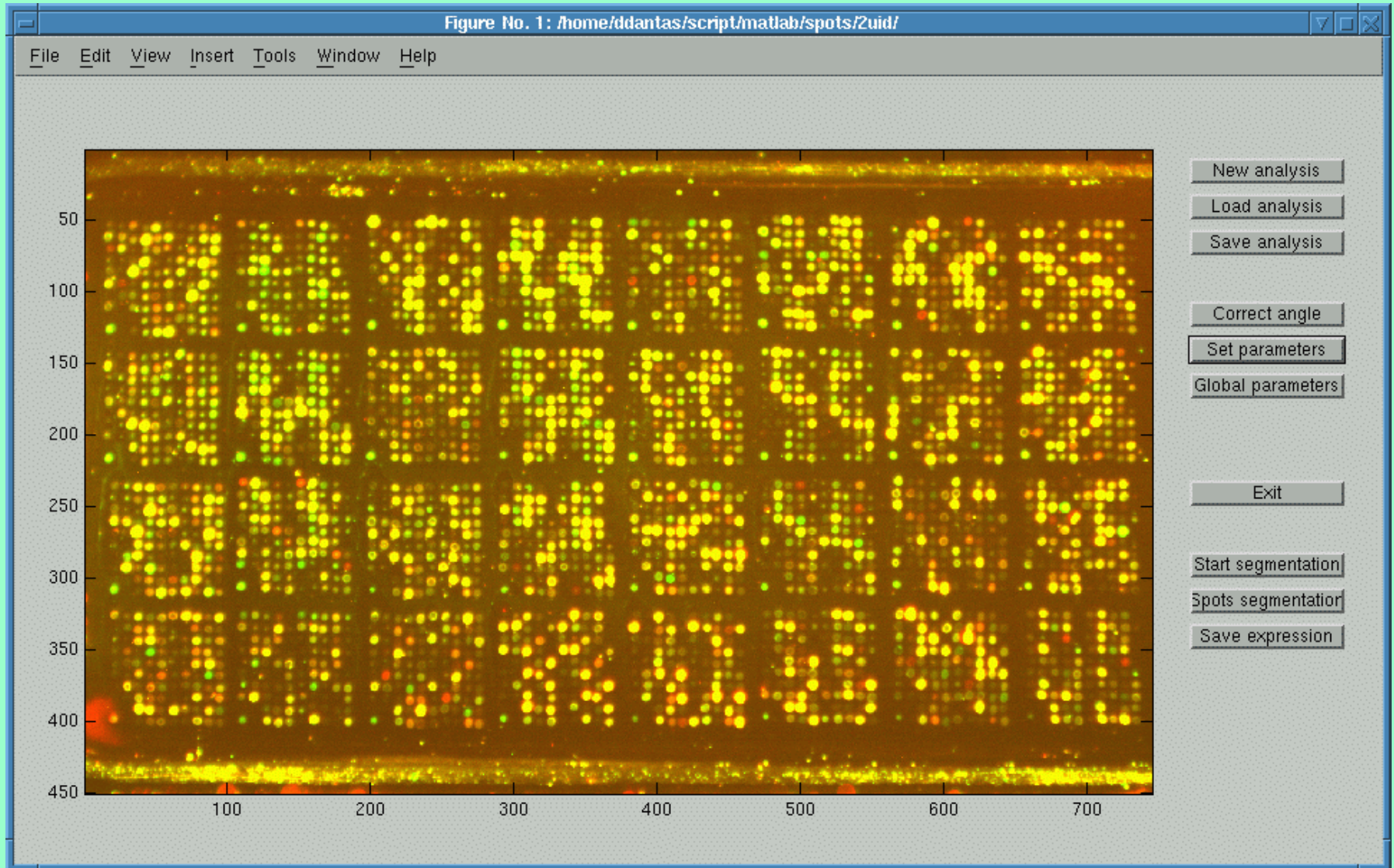
Fore Pct: Back Pct: Spot scale: Enh. Radius

Skip median bk filter
No Spot Enhance

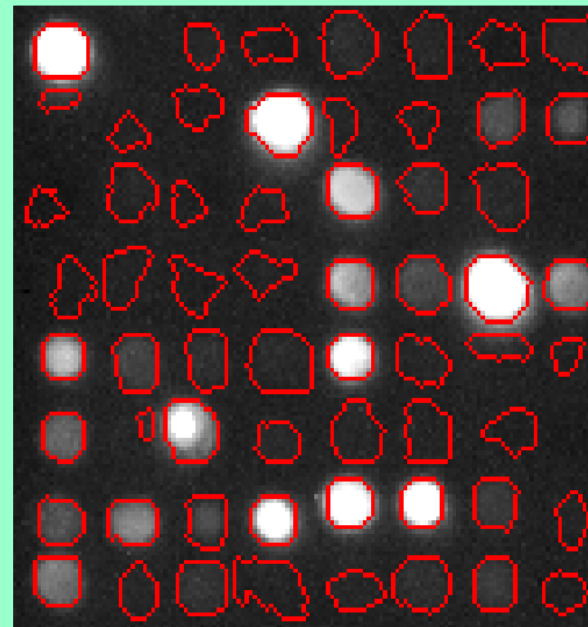
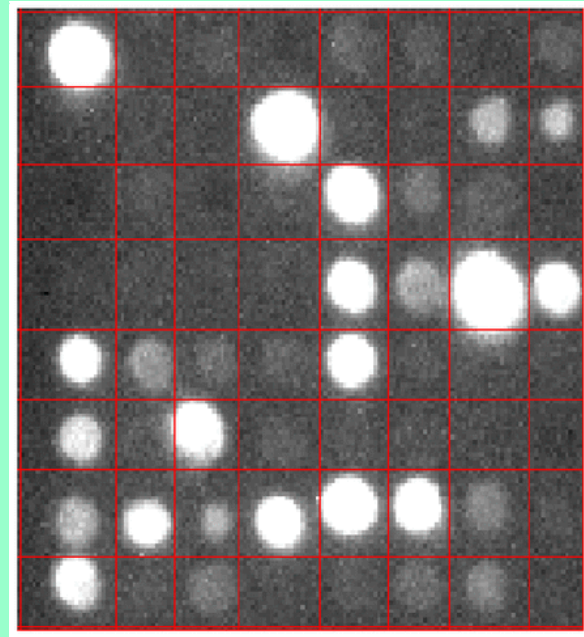
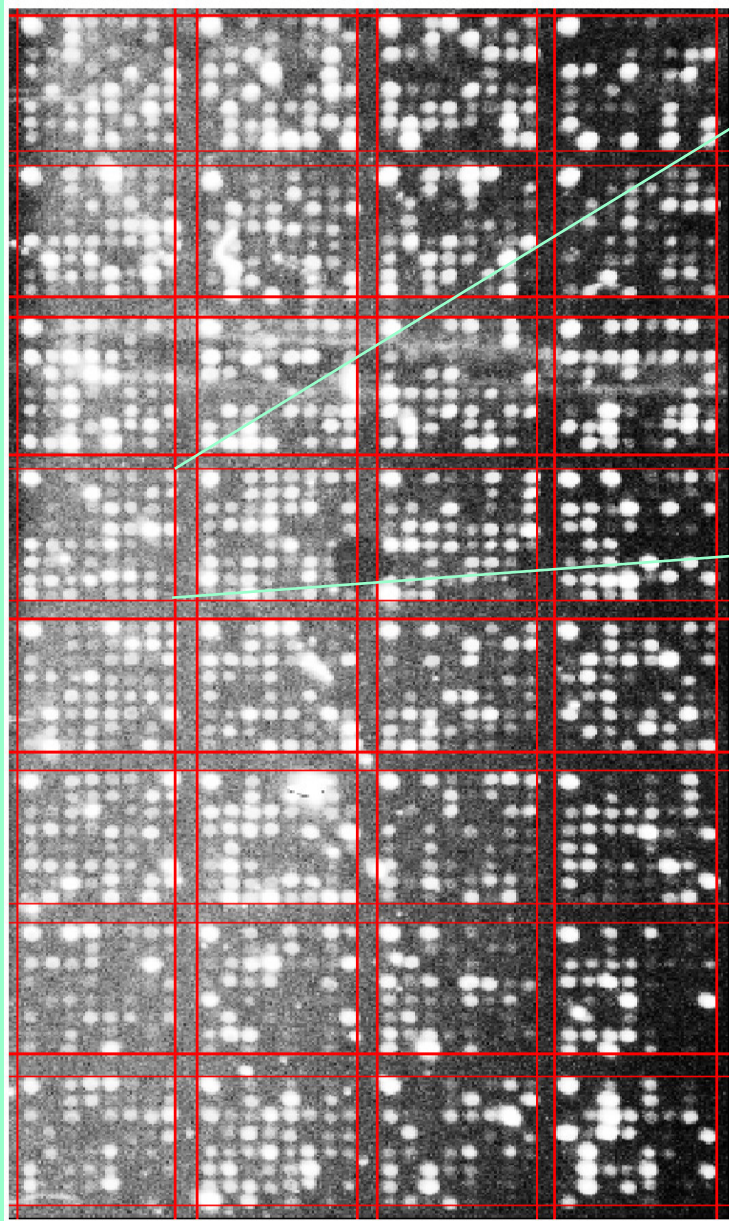
Run
Exit



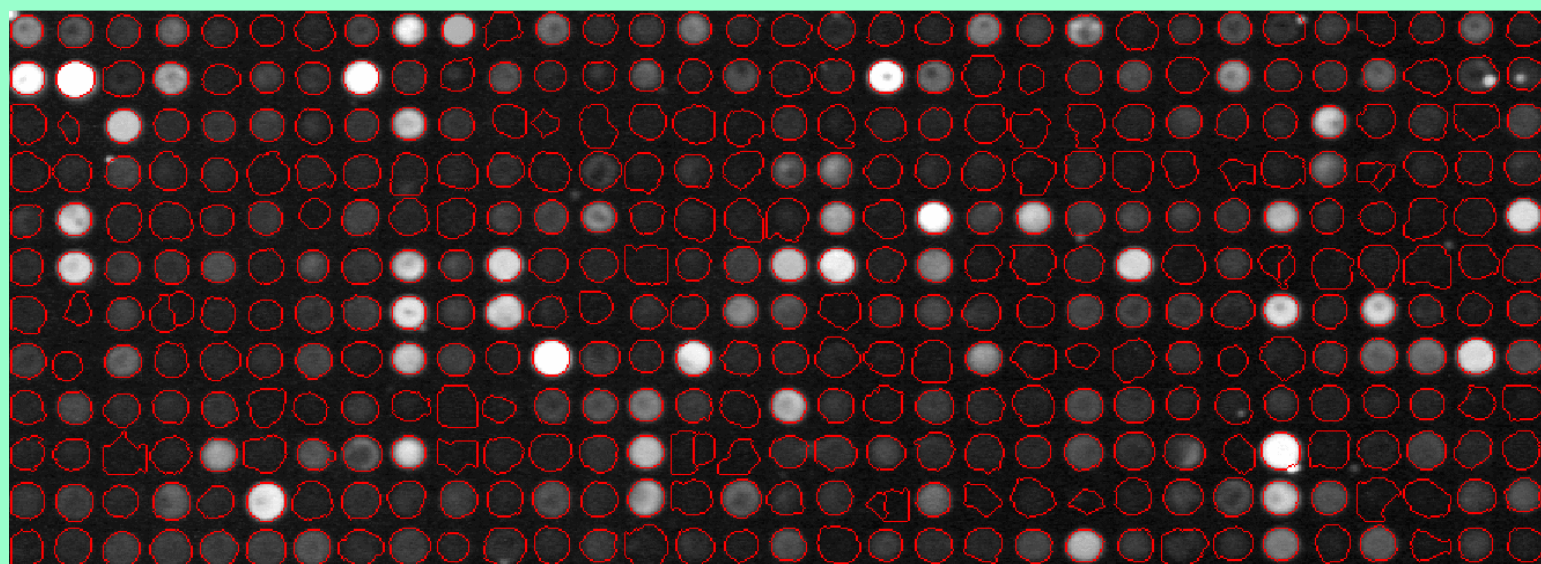
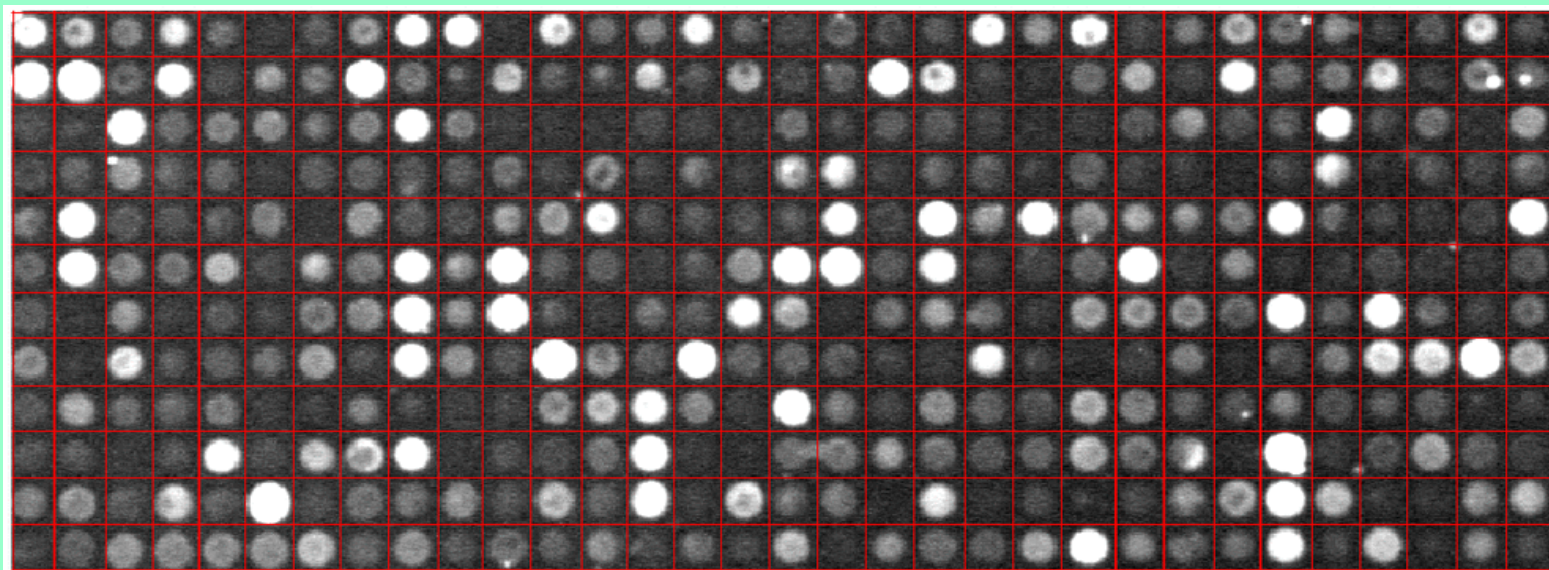
BIOINFO - USP



Example - Segmentation



Exemplo - Segmentação

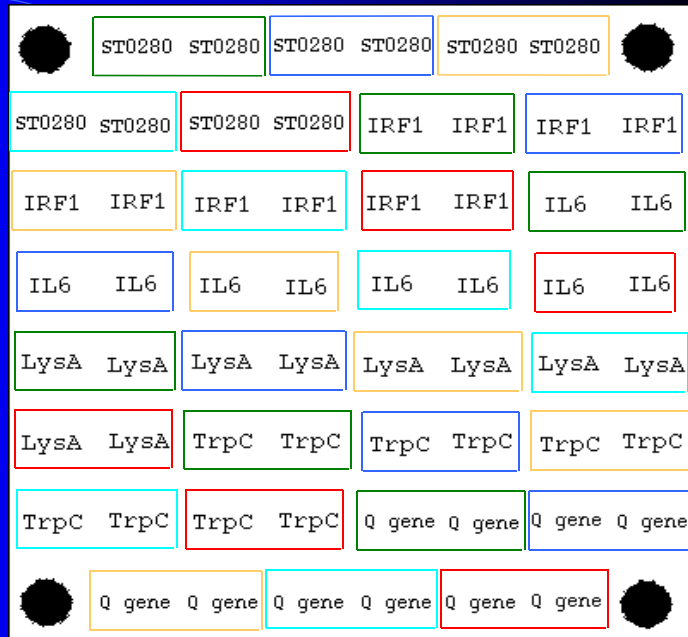


cDNAs used

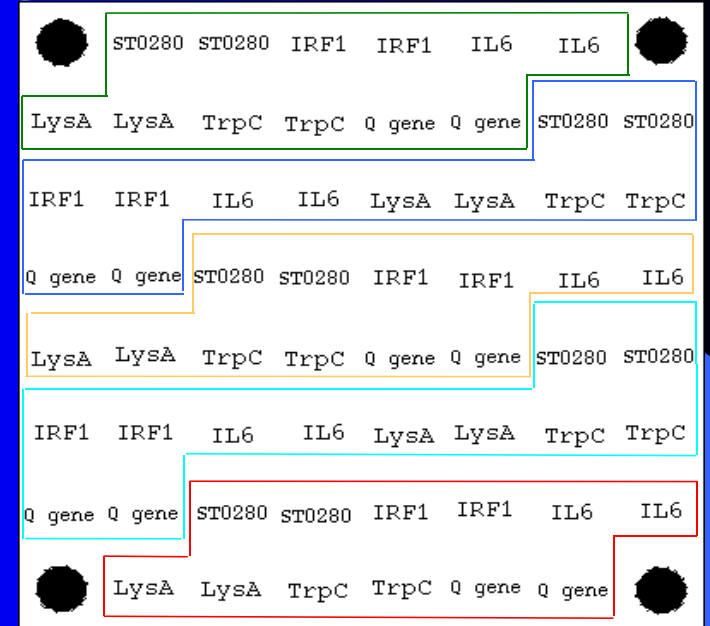
- ✓ *lysA* - 303pb, 47,2% de C e G
- ✓ *trpC* - 338pb, 45,3% de C e G
- ✓ Q gene - 637pb, 52,3% de C e G
- ✓ ST0280 (ORESTES) - 659pb, 34,6% de C e G
- ✓ IL6 - 948pb, 37,7% de C e G
- ✓ IRF1 - 2069pb, 52,5% de C e G

Dilution of fixed material

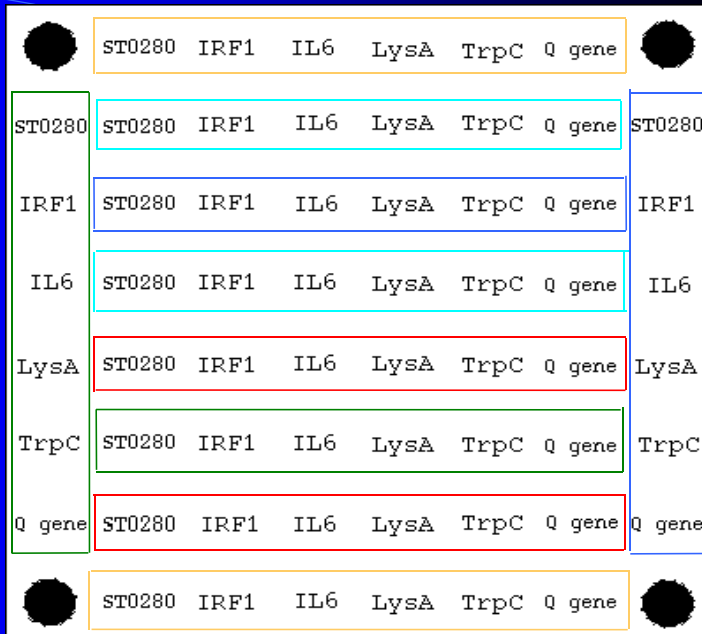
Each cDNA will be fixed in the dilutions
1/1, 1/2, 1/4, 1/8, 1/16

B1**Legenda**

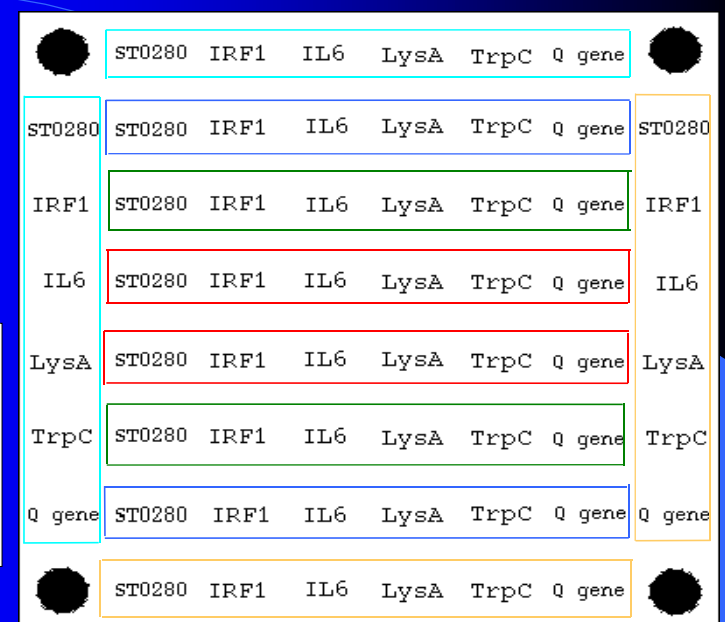
- Diluição 1
- Diluição 2
- Diluição 3
- Diluição 4
- Diluição 5

B2**Legenda**

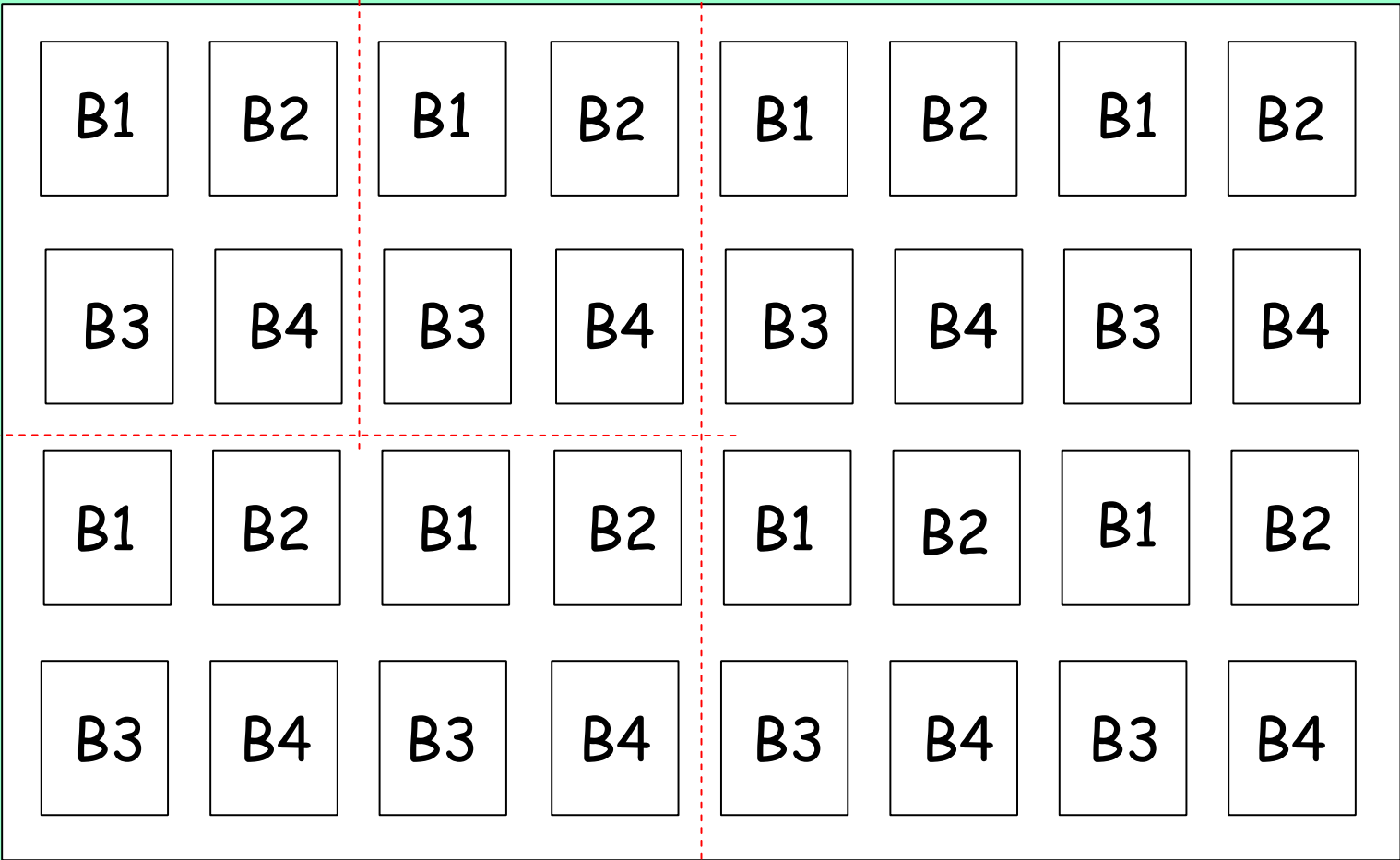
- Diluição 1
- Diluição 2
- Diluição 3
- Diluição 4
- Diluição 5

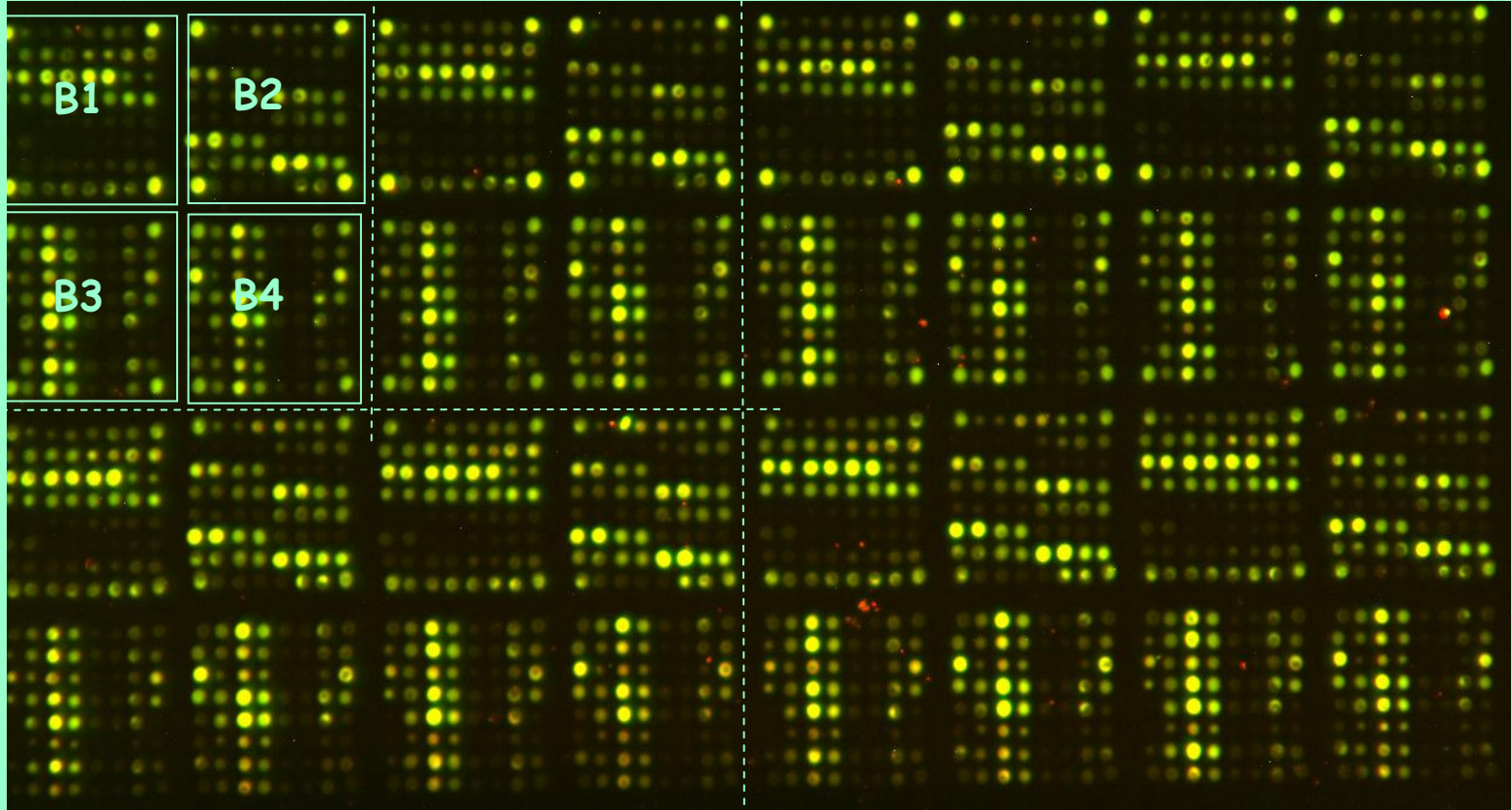
B3**Legenda**

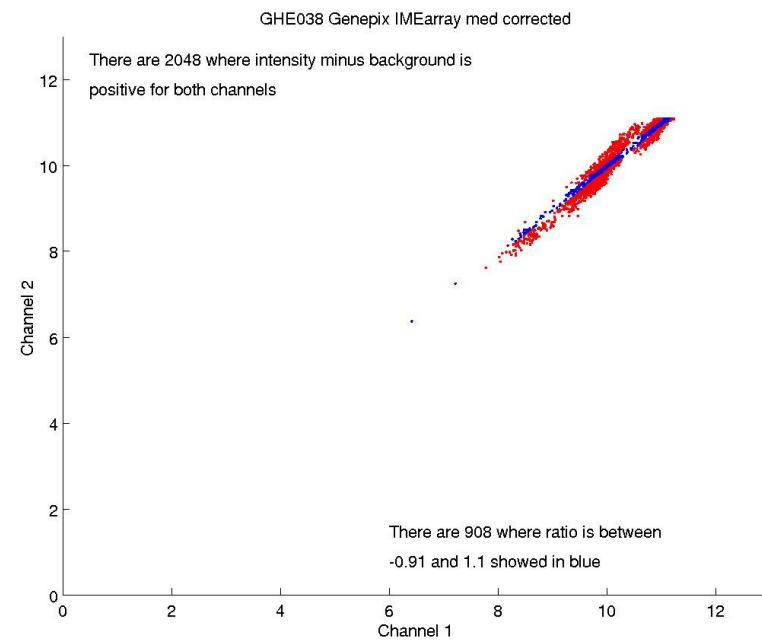
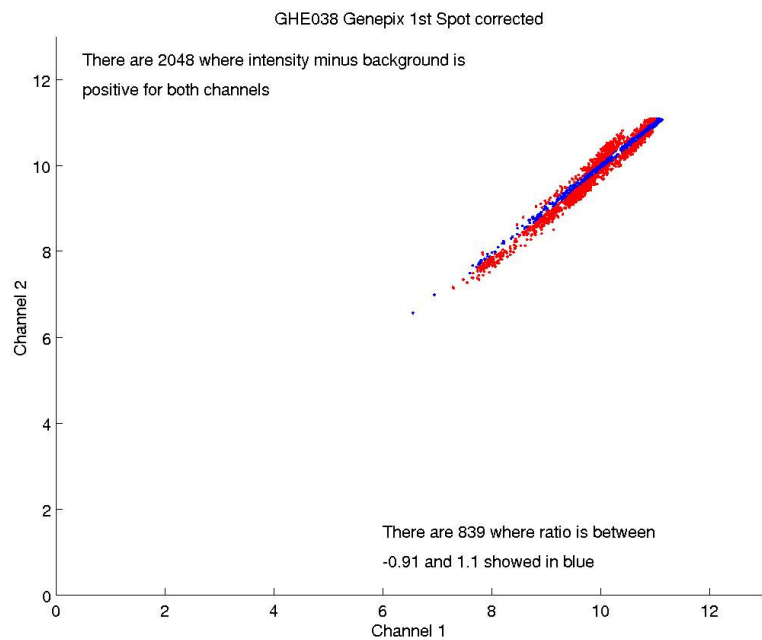
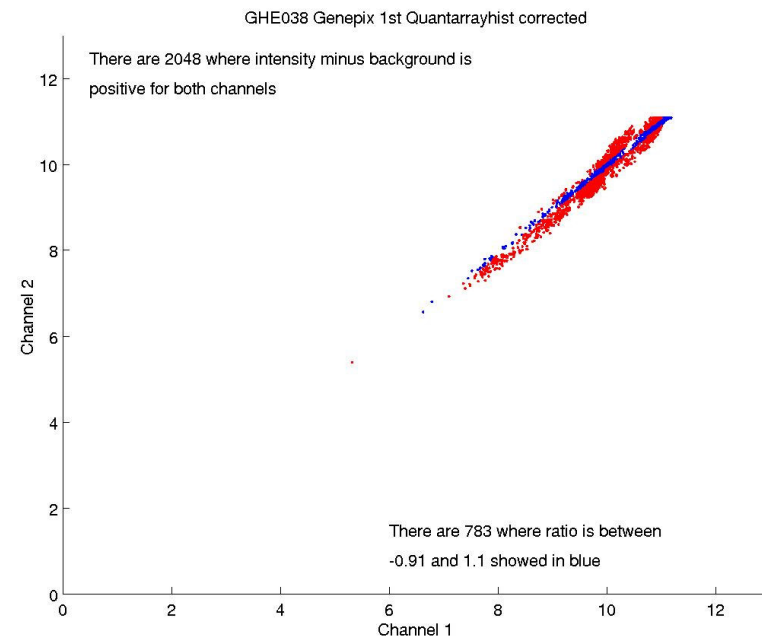
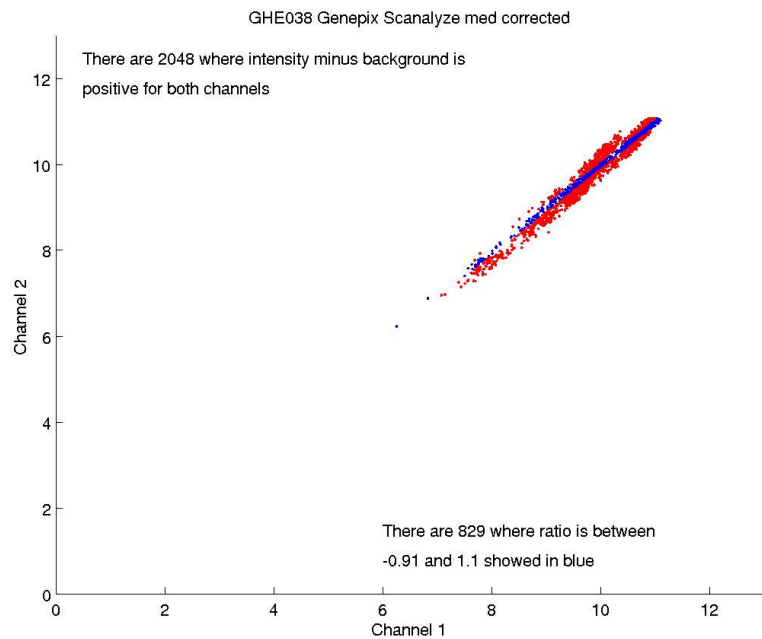
- Diluição 1
- Diluição 2
- Diluição 3
- Diluição 4
- Diluição 5

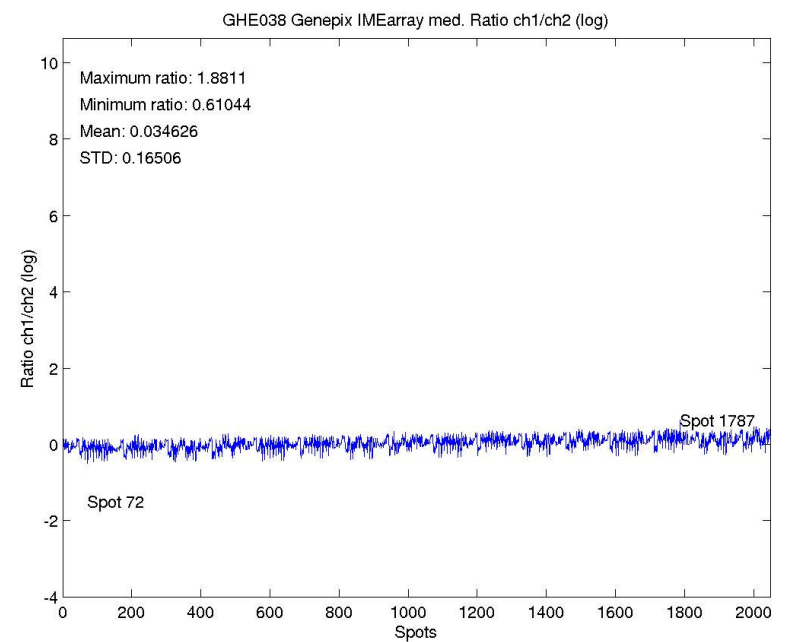
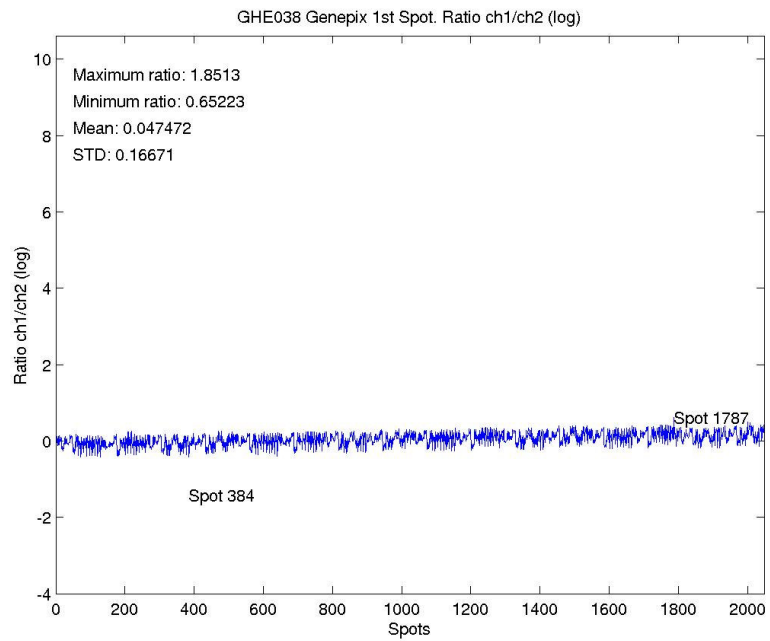
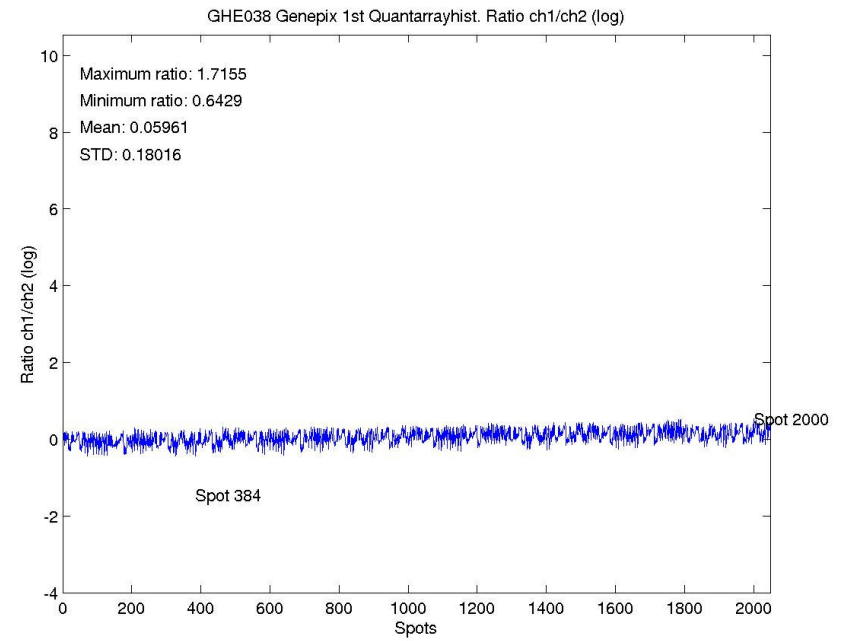
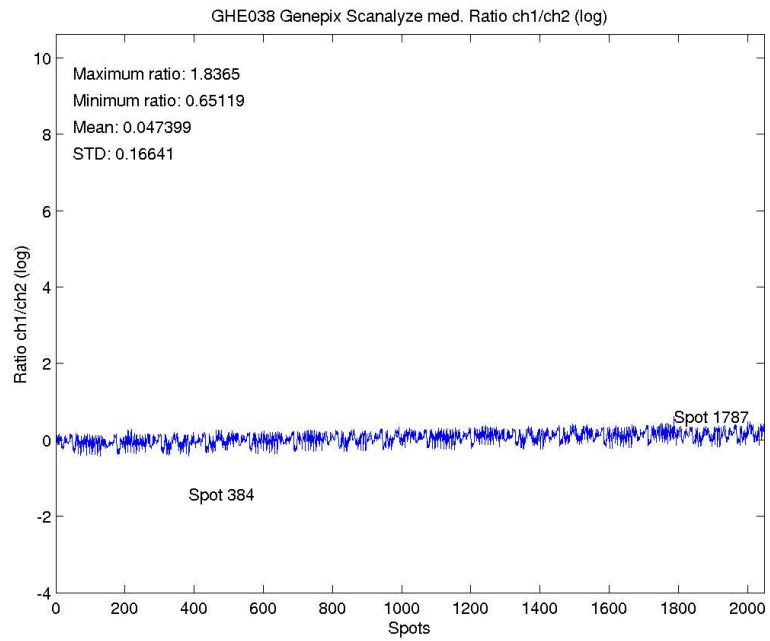
B4**Legenda**

- Diluição 1
- Diluição 2
- Diluição 3
- Diluição 4
- Diluição 5









For a good signal:

- the linear regression is a good estimator for ratio (background estimation is avoided)
- Swap permits to normalize cy3 and cy5
- Confidence intervals increase inversely with the signal intensity

Genes differentially expressed

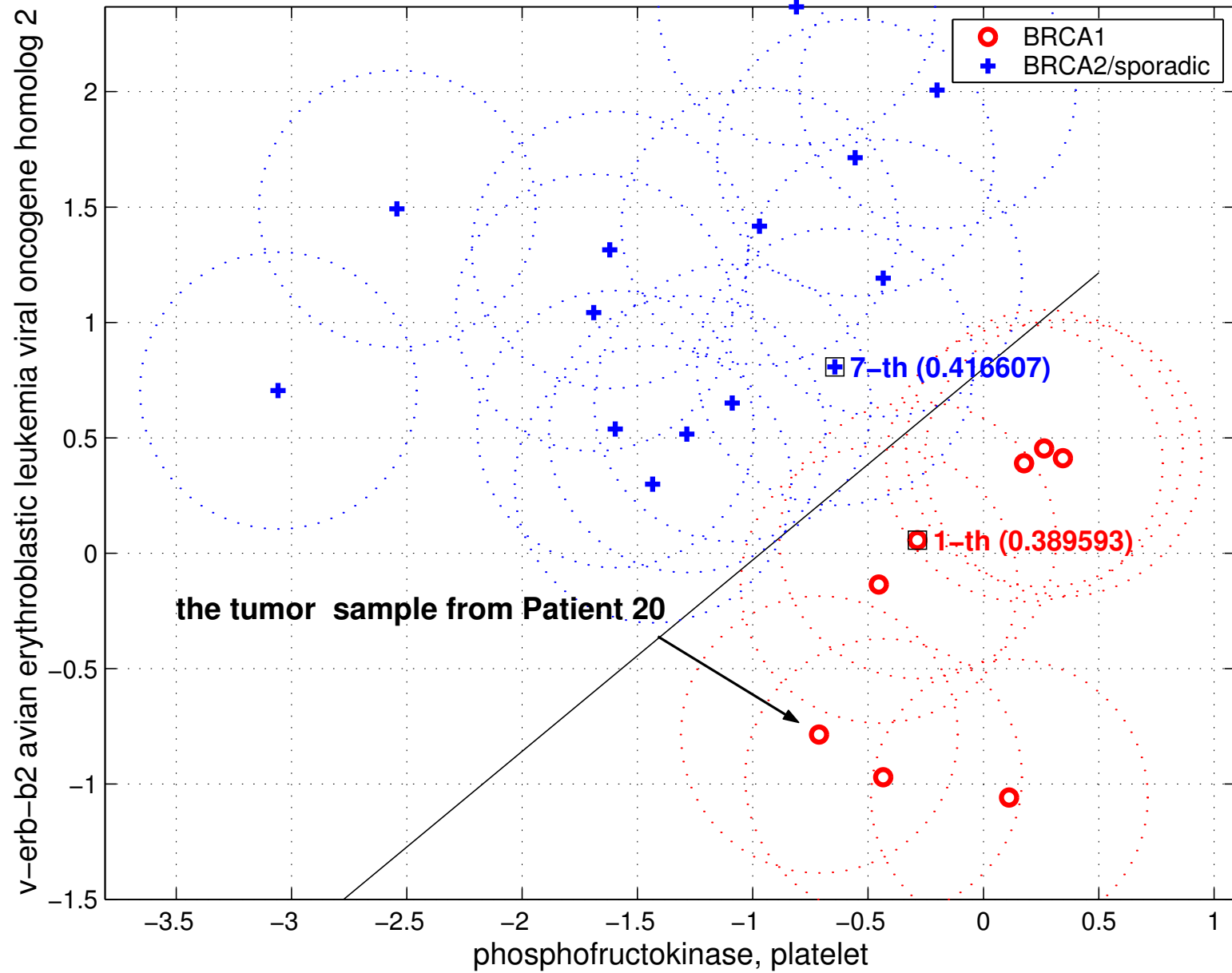
Experiment

- Choose a population of men and measure their physical characteristics
- Ask men to move a 150 kg object
- Separate men that succeed and the ones that do not
- Find common characteristics between men that succeed and between men that do not

Genes differentially expressed

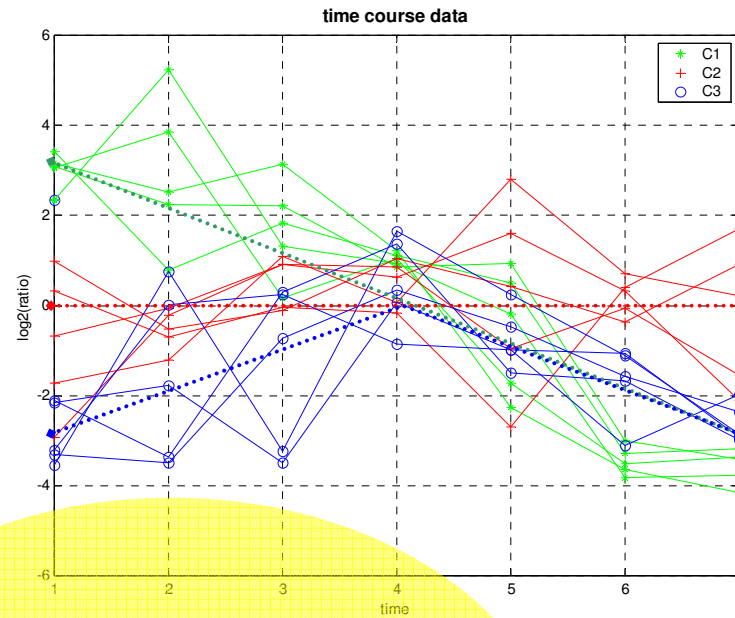
- Choose a set of genes
- Measure the expression of these genes on two different Biological states
- Choose subsets of genes that are enough to characterize each Biological state

LINEAR CLASSIFIER (DISPERSED-GAUSSIAN) w/ $\sigma = 0.600$

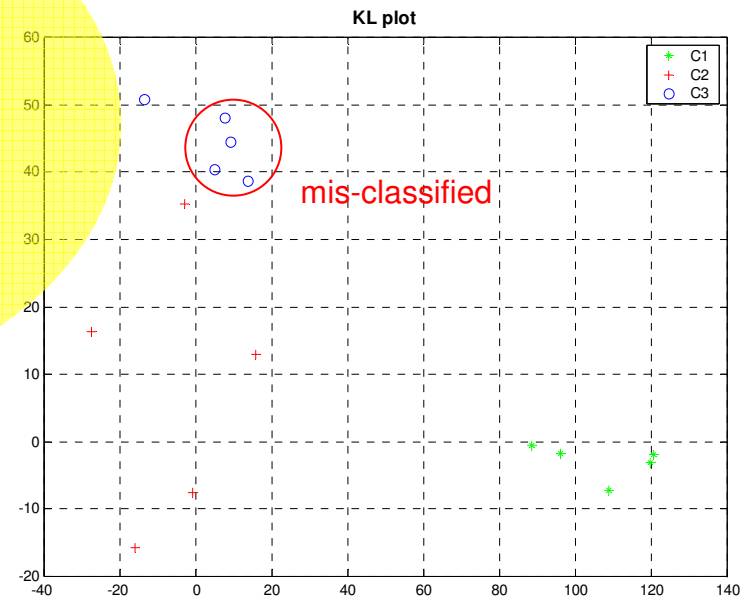


Clustering

Time course data



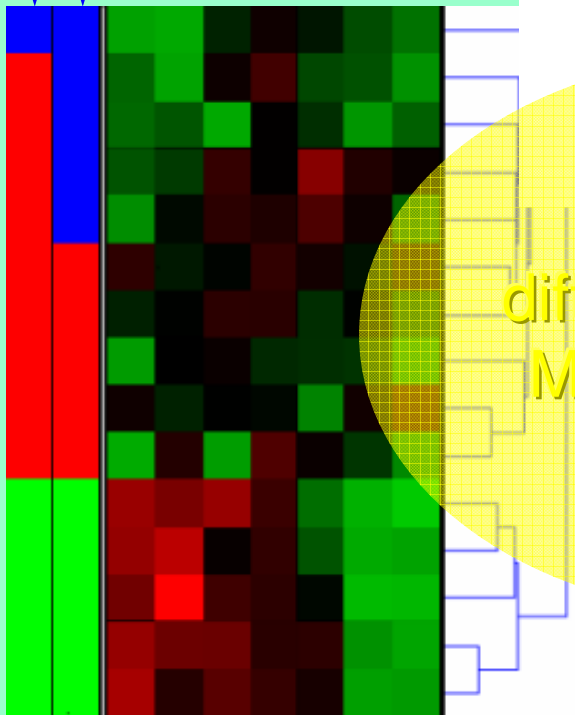
KL plot
multidimensional space



Clustered by dendrogram

Original clusters

Dendrogram

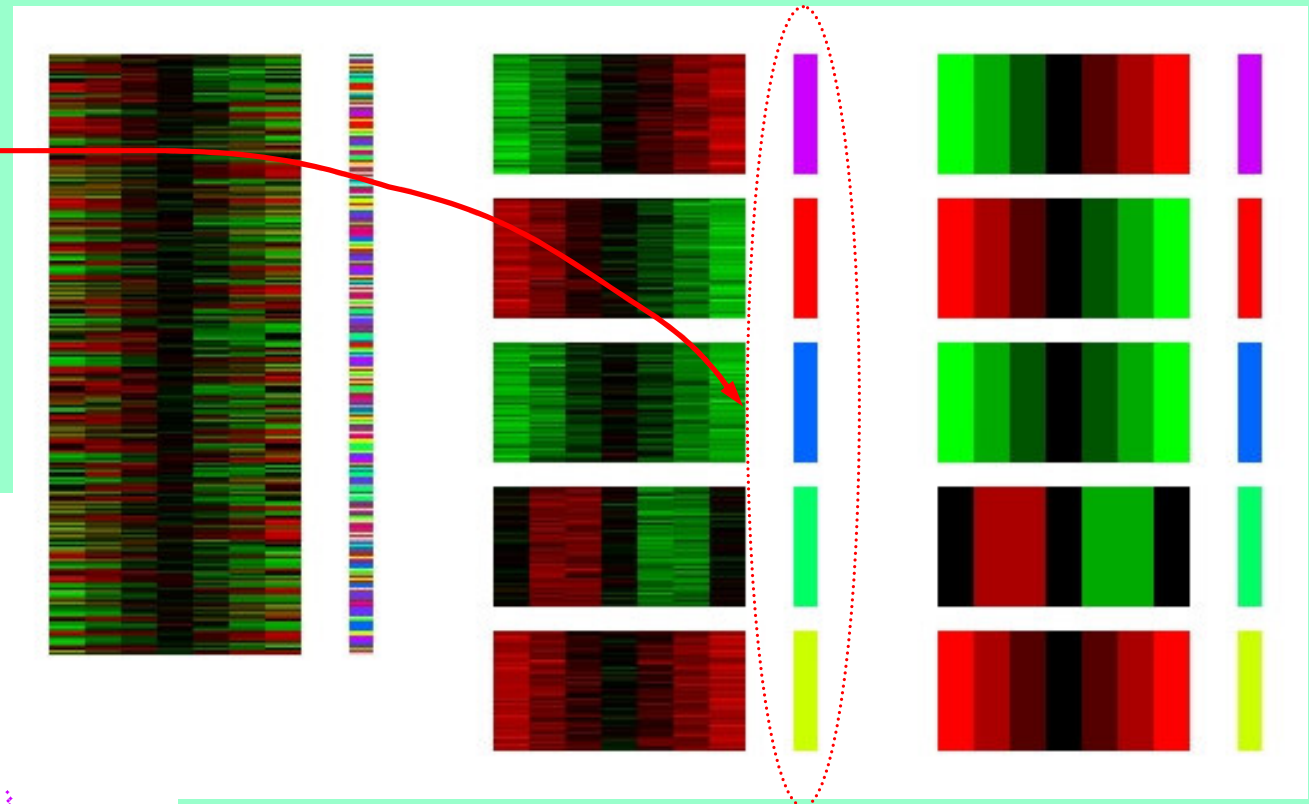
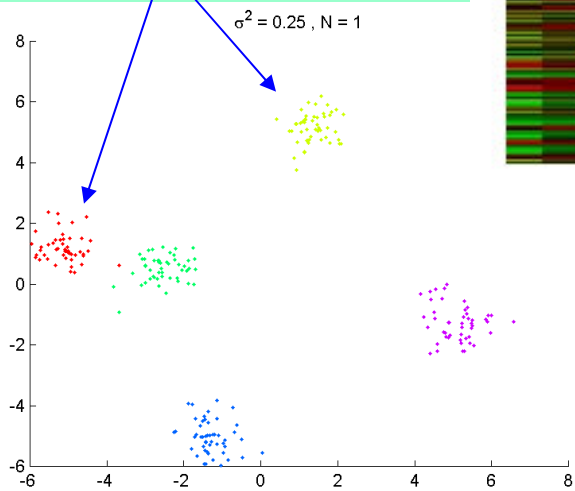


different views of
Microarray data

Example

No error!

Tighter clusters due to small variance



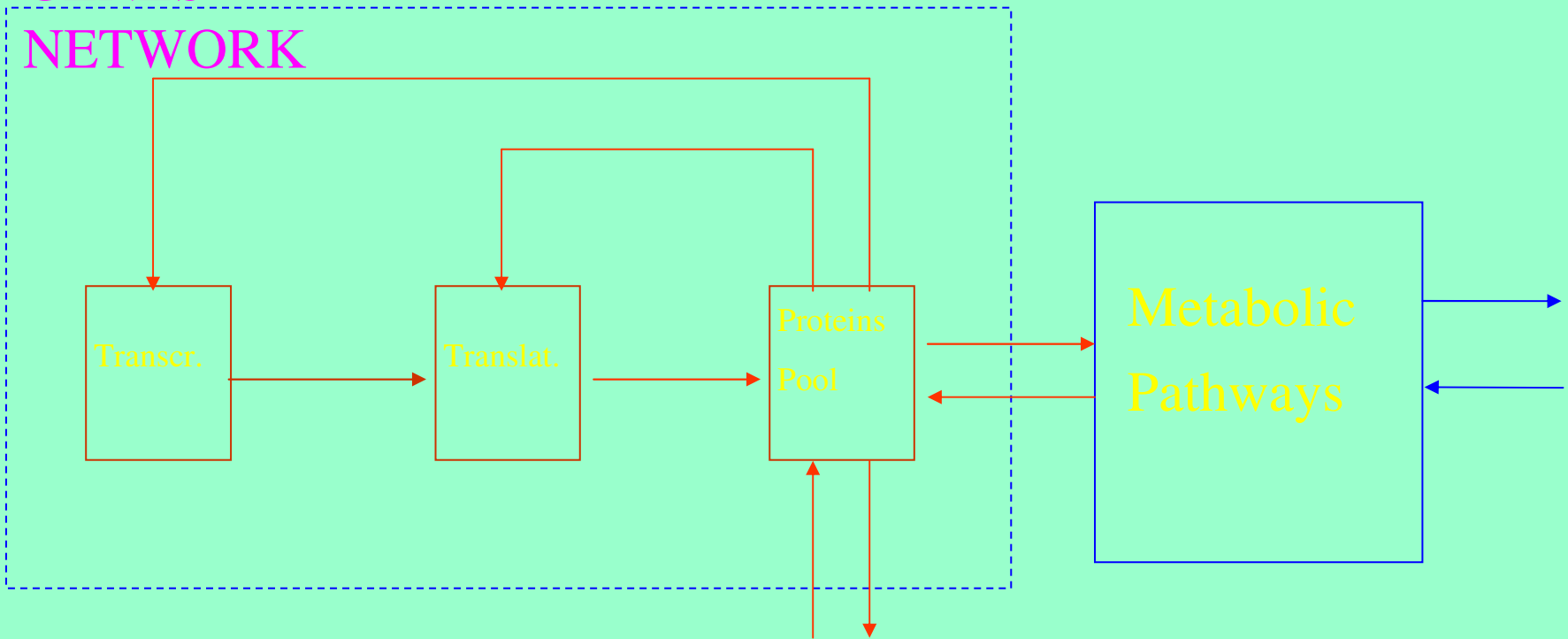
Results from Fuzzy c-means

Gene Regulation Networks

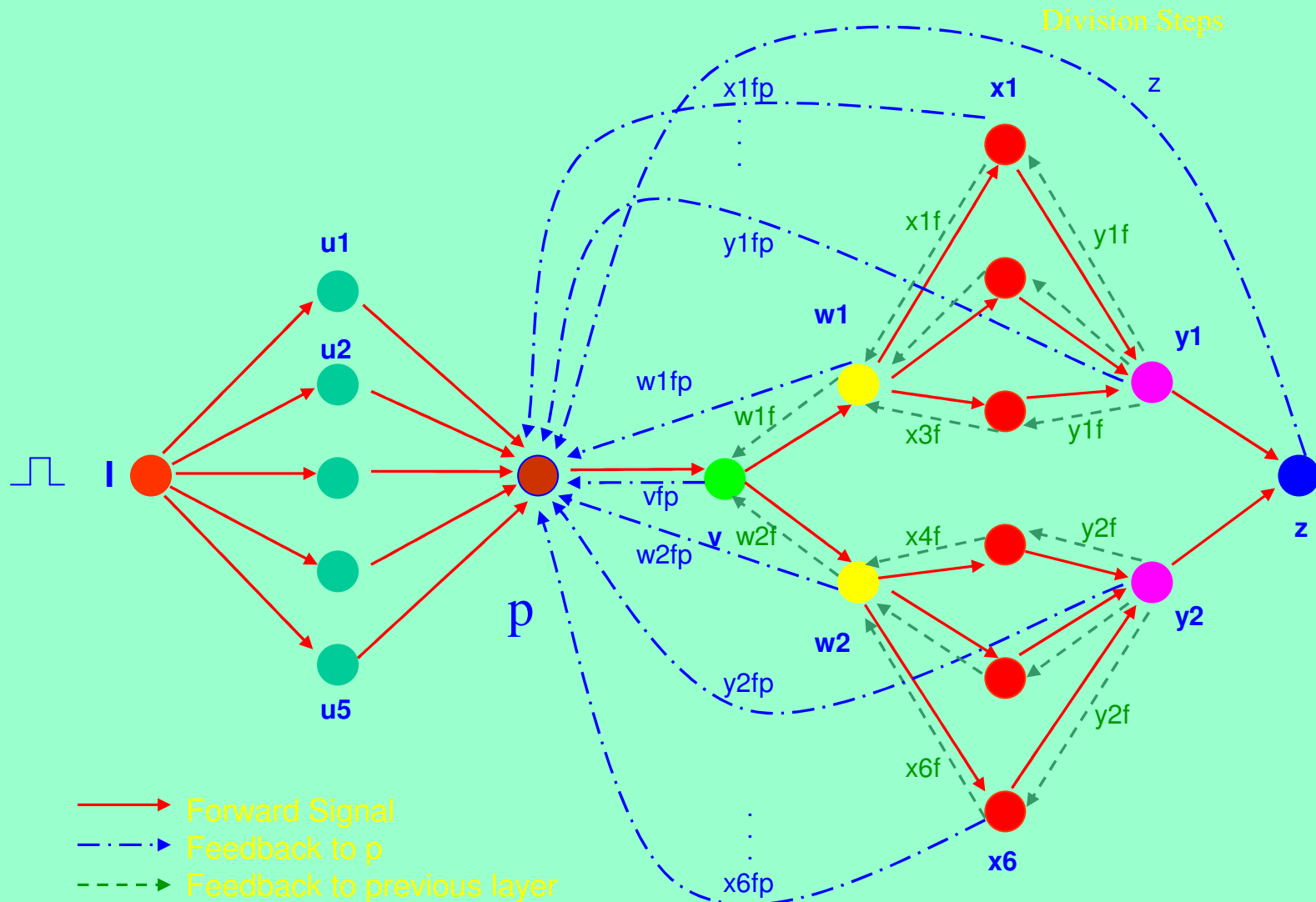
Cell

- peptide
- other signals
- mRNA

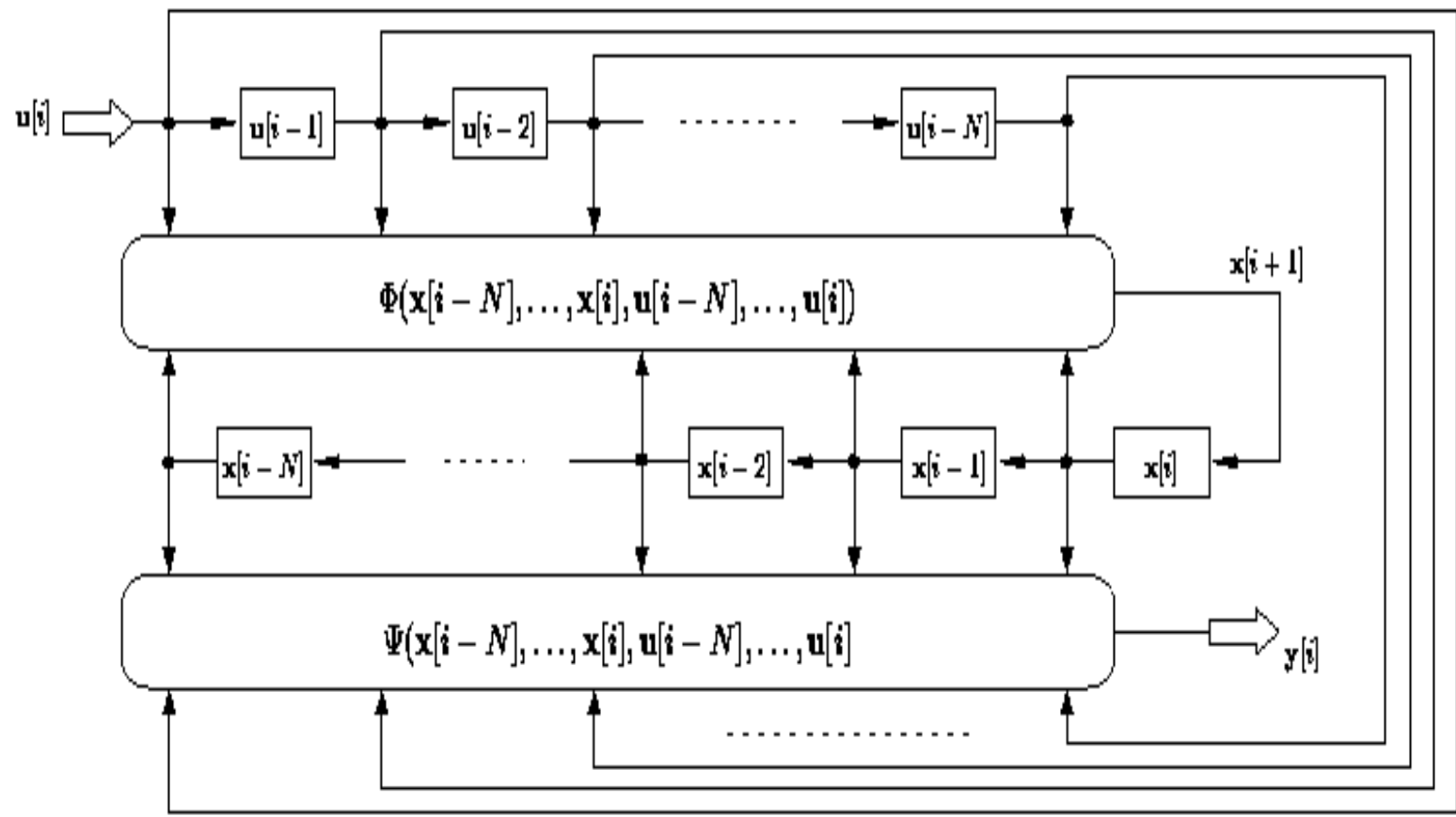
GENES NETWORK



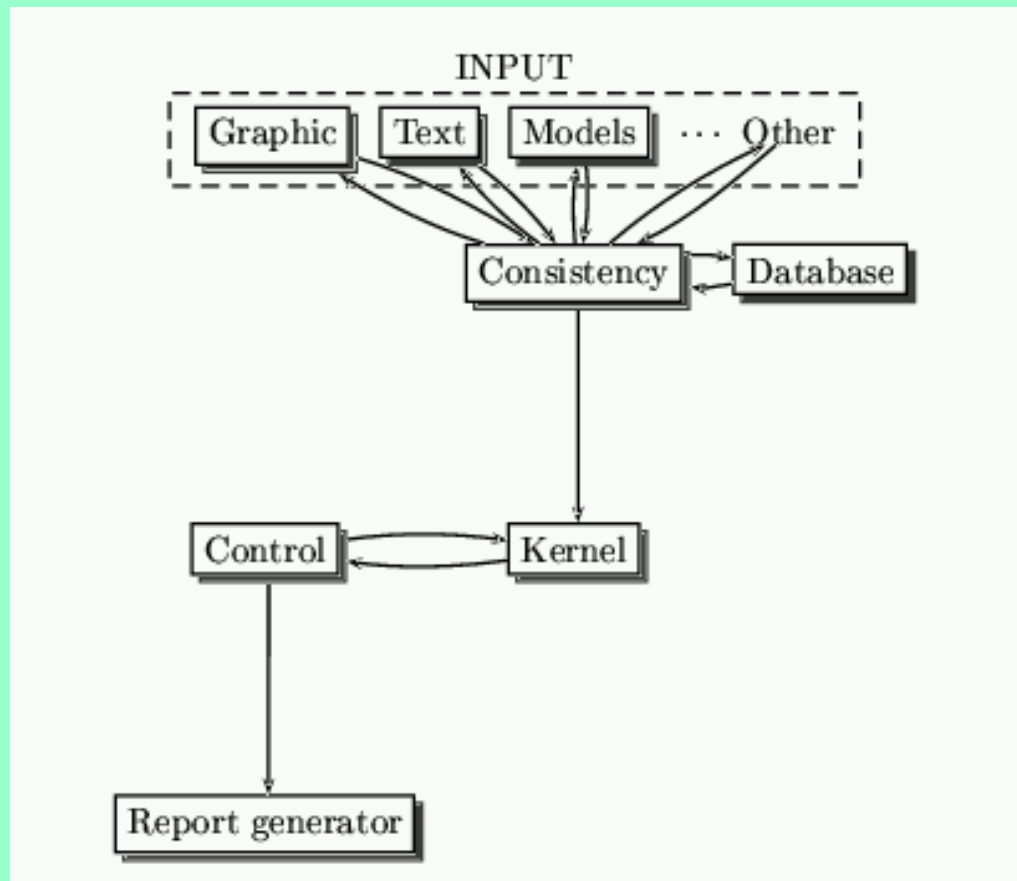
Cell Cycle Modeling



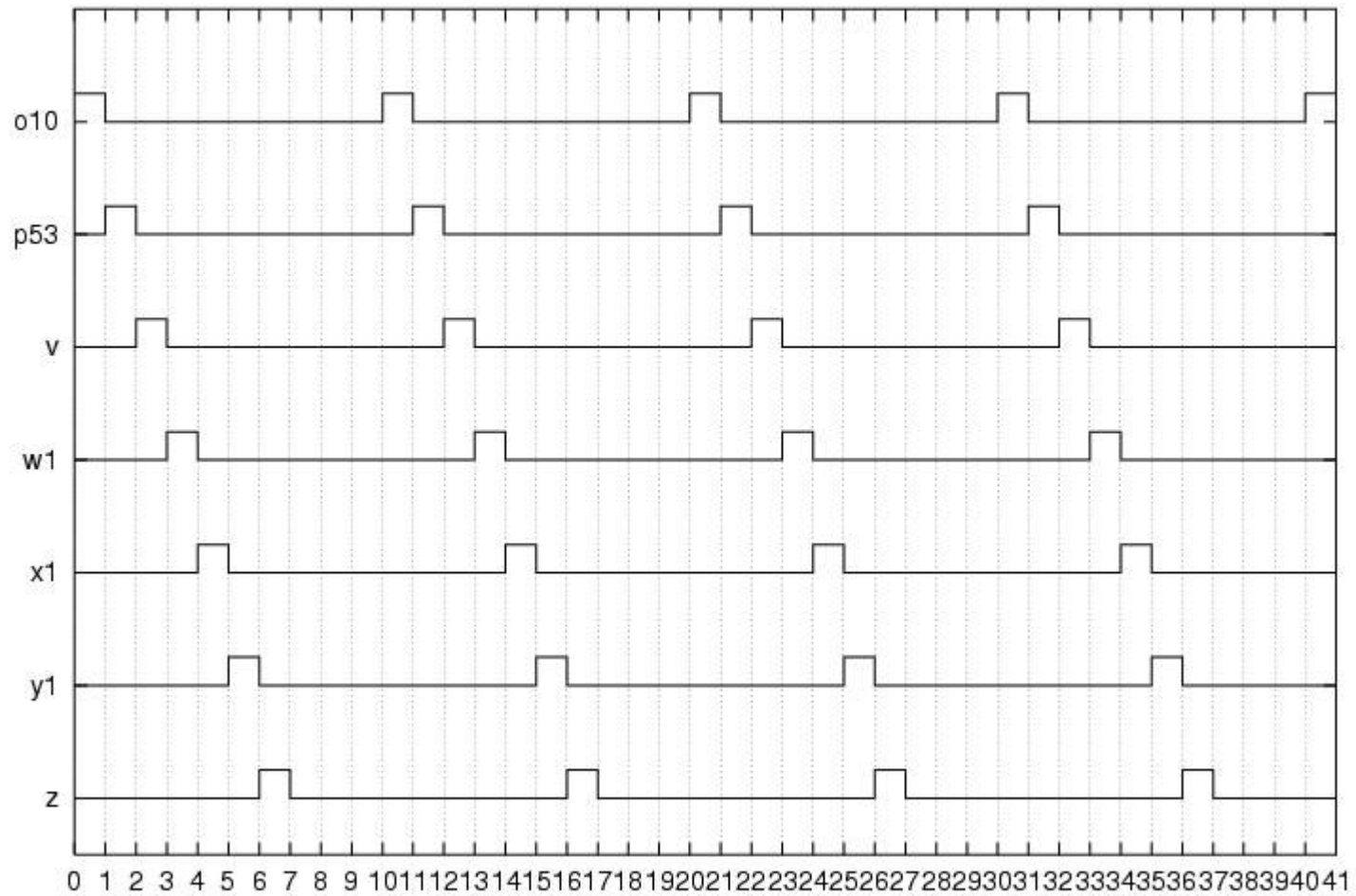
Modeling Dynamical Systems



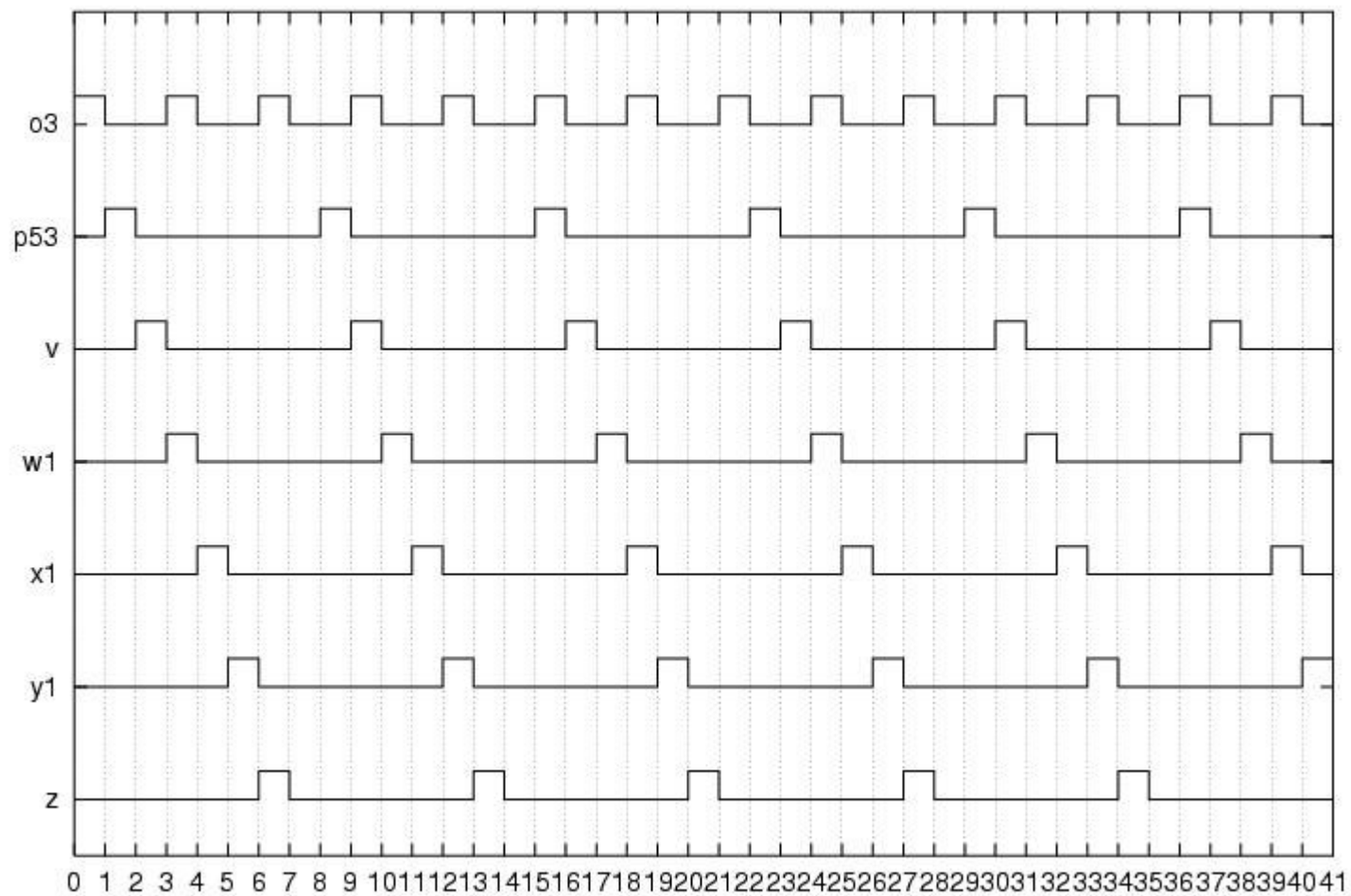
Simulator Architecture



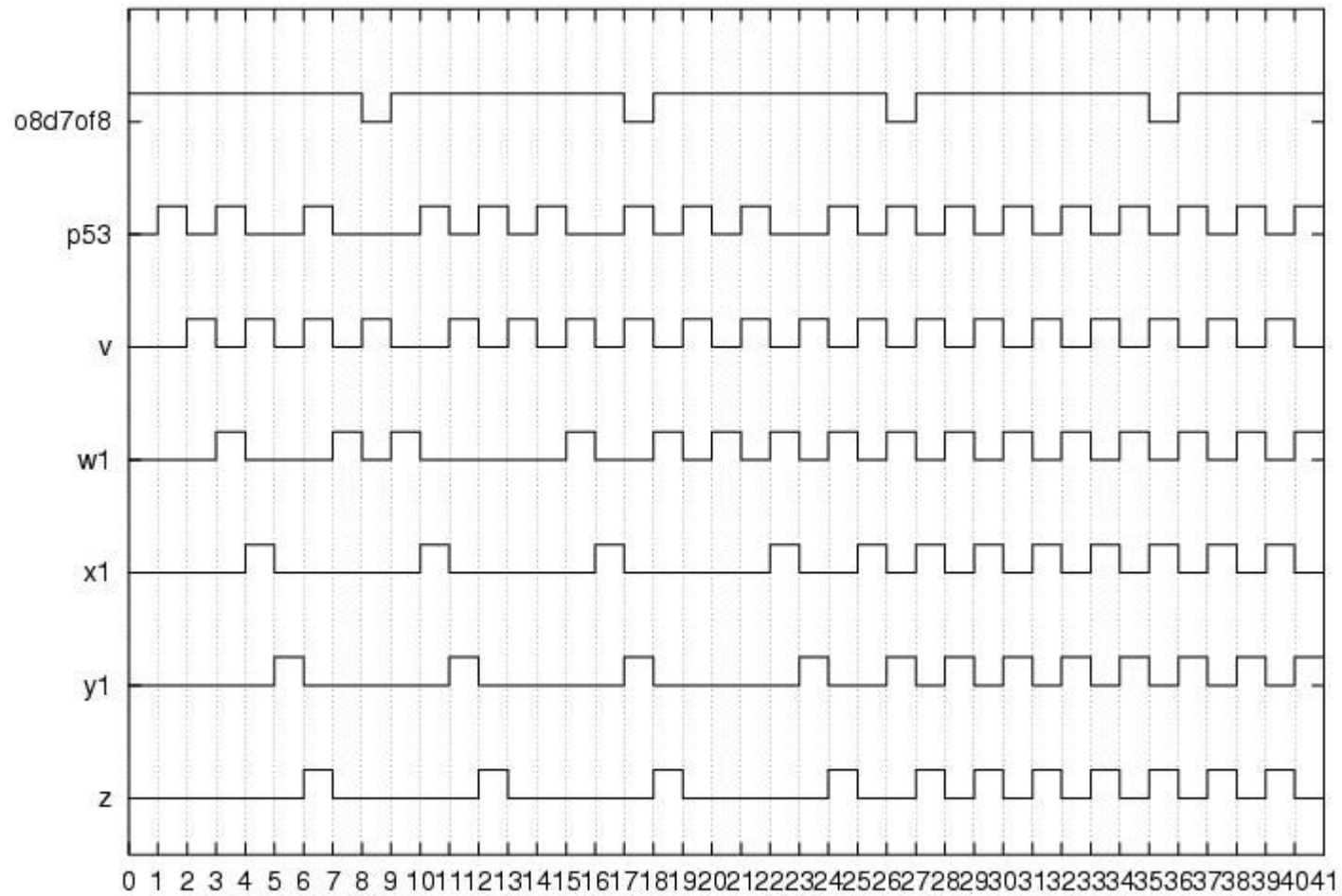
Oscilador de Período 10: FUNCIONAMENTO GERAL (parte_B-t4A-o10.sim)



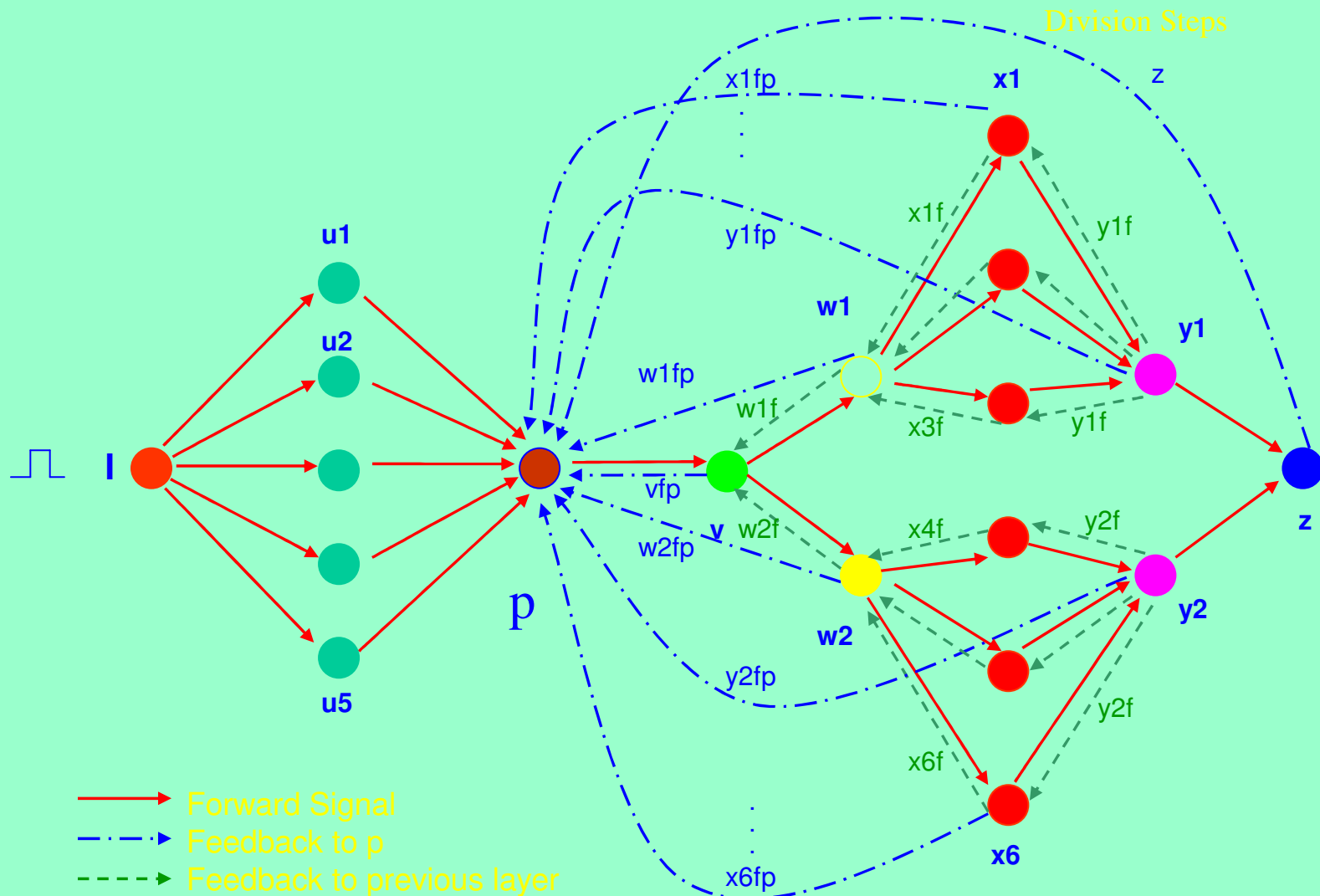
Oscilador de Período 3: FUNCIONAMENTO GERAL (parte_B-t4A-o3.sim)



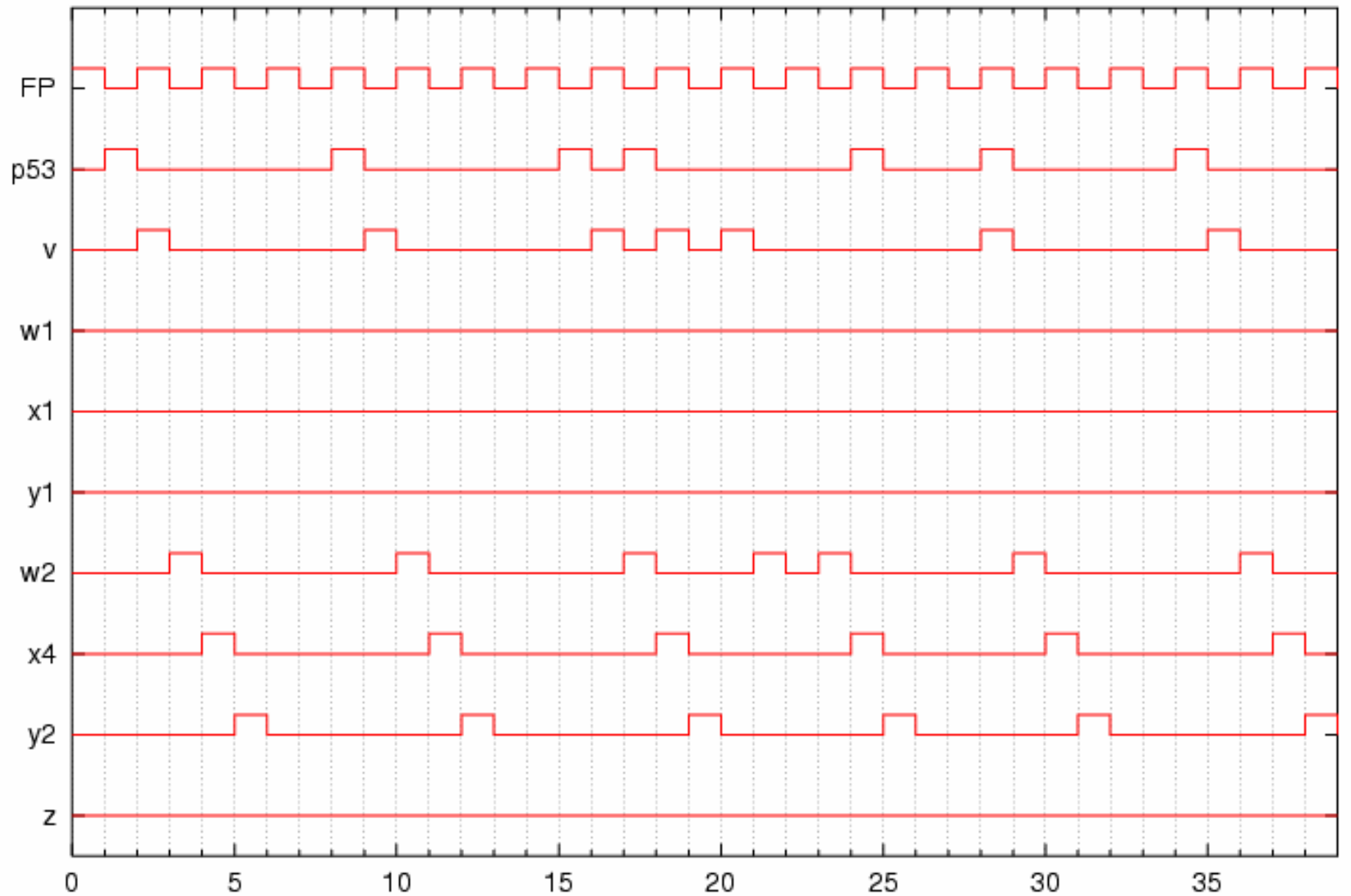
Sinal Periodico 7 ligados 1desligado: FUNCIONAMENTO GERAL (parte_B-t4-o8-7of8.sim)



Knockout



SYSTEM BEHAVIOUR WITH FP = Period 2 Oscillator AND w1 KNOCK OUT



Challenges

- Architecture identification
- Dynamics transition function identification