

# *Design of Probabilistic Genetic Networks*

*Applications to Malaria and Cell Cycle.*

Junior Barrera

**DCC – IME/USP**

**Universidade de São Paulo**

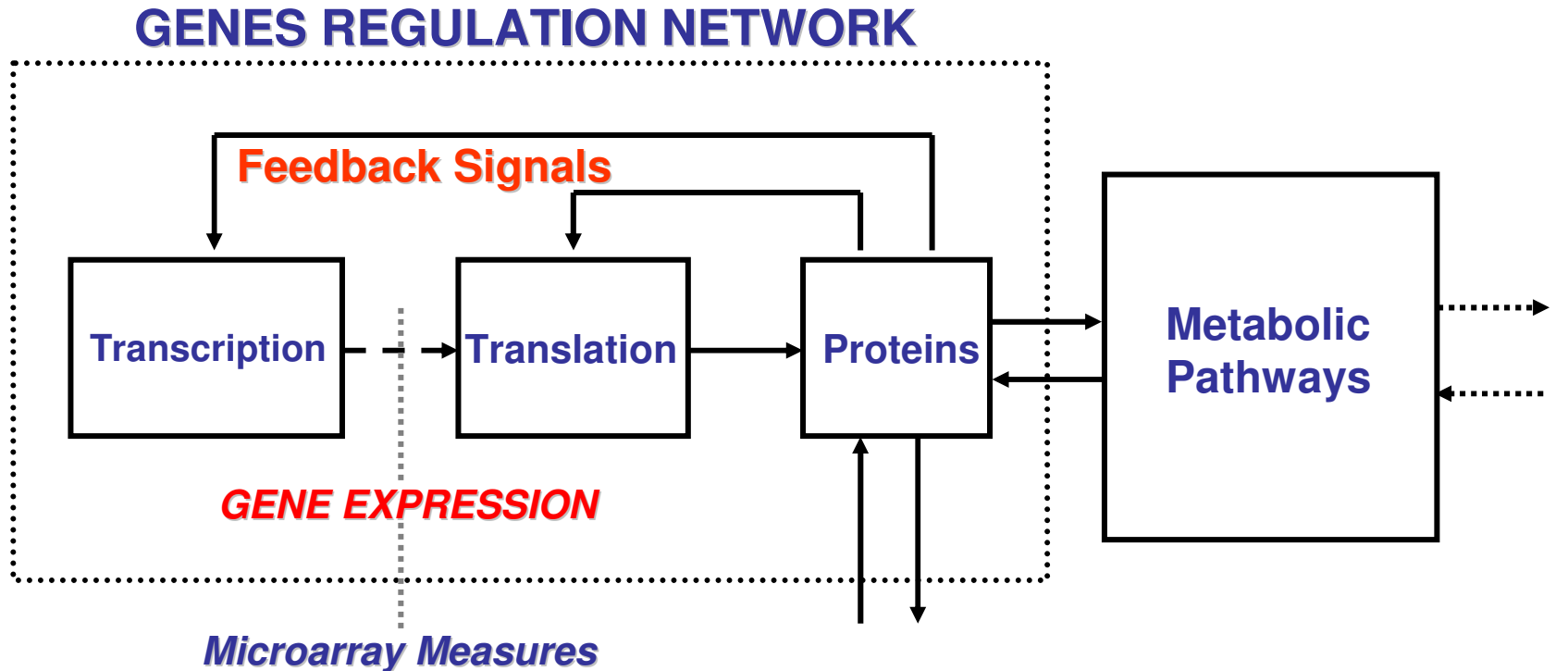
# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- Malaria
- Cell Cycle
- Future works

# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- Malaria
- Cell Cycle
- Future works

# Genetic Control System



- - - → mRNA
- → peptide
- ..... → other signals

# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- Malaria
- Cell Cycle
- Future works

# Model characteristics and Biological motivations

- Describes genes expression dynamics

Gene expression is the only signal measured

- Discrete time and range

At the time-space resolution level considered, molecular synthesis and interaction are sequential discrete phenomena

- Genes expression dynamics are stochastic signals

There is noise in molecular synthesis and interaction

- Genes are non linear causal stochastic gates

Transcription is regulated by the interaction between proteins and DNA

# Model characteristics and Biological motivations

- PGN is built by the connection of non linear causal stochastic gates

Genes network is the interaction of genes

- PGN dynamics is a vector of discrete stochastic signals

Genes network dynamics is a vector of stochastic signals

## Model formalization

Number of genes in the PGN:  $N$

Expression of gene  $i$  at time  $t$  :  $x_i(t) \in R \subset Z, |R| \text{ finite}$

State of the PGN at time  $t \in \{0,1,2,\dots\}$  :  $x = \begin{bmatrix} x_1(t) \\ x_2(t) \\ \vdots \\ x_N(t) \end{bmatrix}$



## Model formalization

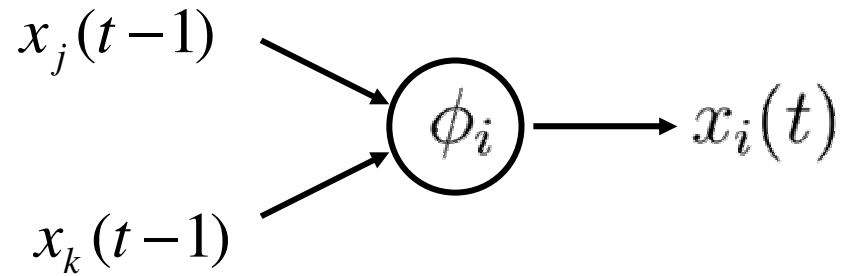
Dynamics:

$$x_i(t) = \phi_i(x(t-1)),$$
$$t > 1 \quad \text{and} \quad i \in \{1, 2, \dots, N\}$$

$\phi_i$  is a stochastic function :

$$x_i(t) = \begin{cases} r_1 & p(r_1 | (x(t-1))) \\ r_2 & p(r_2 | (x(t-1))) \\ \vdots & \vdots \\ r_{|R|} & p(r_{|R|} | (x(t-1))) \end{cases}$$

# Nomenclature



Predictors

Target

## Model formalization

The sequence of random vectors  $X_0, X_1, \dots, X_t$ , with observations in  $R^N$  is a **stochastic process**

**Markov chain** is a stochastic process such that

$$P(X_t = x(t) \mid P(X_0 = x(0), X_1 = x(1), \dots, X_{t-1} = x(t-1))) = P(X_t = x(t) \mid X_{t-1} = x(t-1))$$

A Markov chain is **homogenous** when

$$P(X_t = x(t) \mid X_{t-1} = x(t-1)) \text{ is the same for every } t$$

An **homogeneous Markov** chain is **characterized** by the Distribution of  $X_0$  and the state transition probability matrix:

$$(\pi_0, \pi_{Y/X})$$

## Model formalization

A PGN is an homogenous Markov chain  $(\pi_0, \pi_{Y|X})$  which is

- conditionally independent,

$$P(x[t+1] | x[t]) = \prod_{i=1}^n p(x_i[t+1] | x[t])$$

- almost deterministic and

$$\forall t, \forall x(t) \in R^N, \forall i \in \{1, \dots, N\}, \exists r \in R : p(r | x(t)) \approx 1$$

- has limited dependence.

$$\forall i \in \{1, \dots, N\}, \exists j, j \ll N : \forall t, \forall x(t) \in R^N, \forall r \in R, p(r | x(t)) = p(r | x_{1:j}^r(t))$$

# Characterization

- Markov chain

$$P(x(t) | x(t-1))$$

Matrix dimension:  $|R|^N \times |R|^N$

- PGN,

$$p(x_i(t) | x(t-1))$$

Matrix dimension:  $|R|^j \times |R| \times N$

# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- **Estimation of PGNs**
- Architecture estimation
- Malaria
- Cell Cycle
- Future works

## Distribution estimation

Training data for the target gene  $i$

$$(x_1^j, y_1), (x_2^j, y_2), \dots, (x_m^j, y_m)$$

Estimator

$$n = \left| \left\{ (x_k^j, y_k) : x_k^j = x^j \right\} \right|$$

$$p = \left| \left\{ (x_k^j, y_k) : x_k^j = x^j \wedge y_k = r \right\} \right|$$

$$\hat{P}_i(X^j = x^j) = \frac{n}{m} \qquad \hat{P}_k(Y = r | X^j = x^j) = \frac{p}{n}$$

# Distribution estimation

## Problem

Lack of data, usually there are non observed  $x^j$

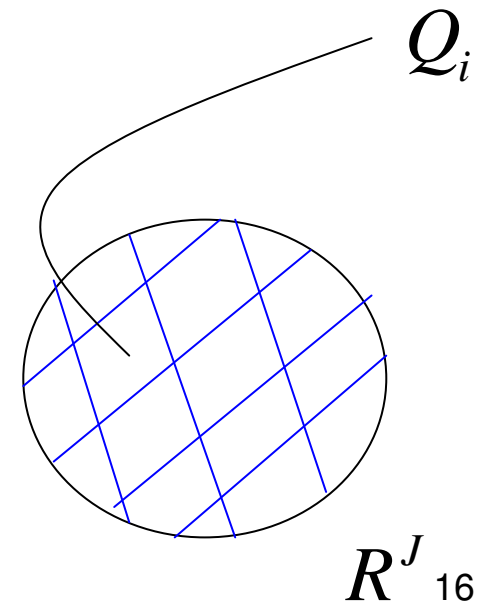
## Solution

Constraint the probability families considered  
partitioning  $R^J$

$Q = \{Q_1, Q_2, \dots, Q_n\}$  is a **partition** of  $R^J$

$$\forall x_k^j \in Q_i : P(y | x_k^j) = P(y | Q_i)$$

$$P(Q_i) = \sum_{x^j \in Q_i} P(x^j)$$



$R^J$  16



## Distribution estimation

Training data for the target gene  $i$

$$(x_1^j, y_1), (x_2^j, y_2), \dots, (x_m^j, y_m)$$

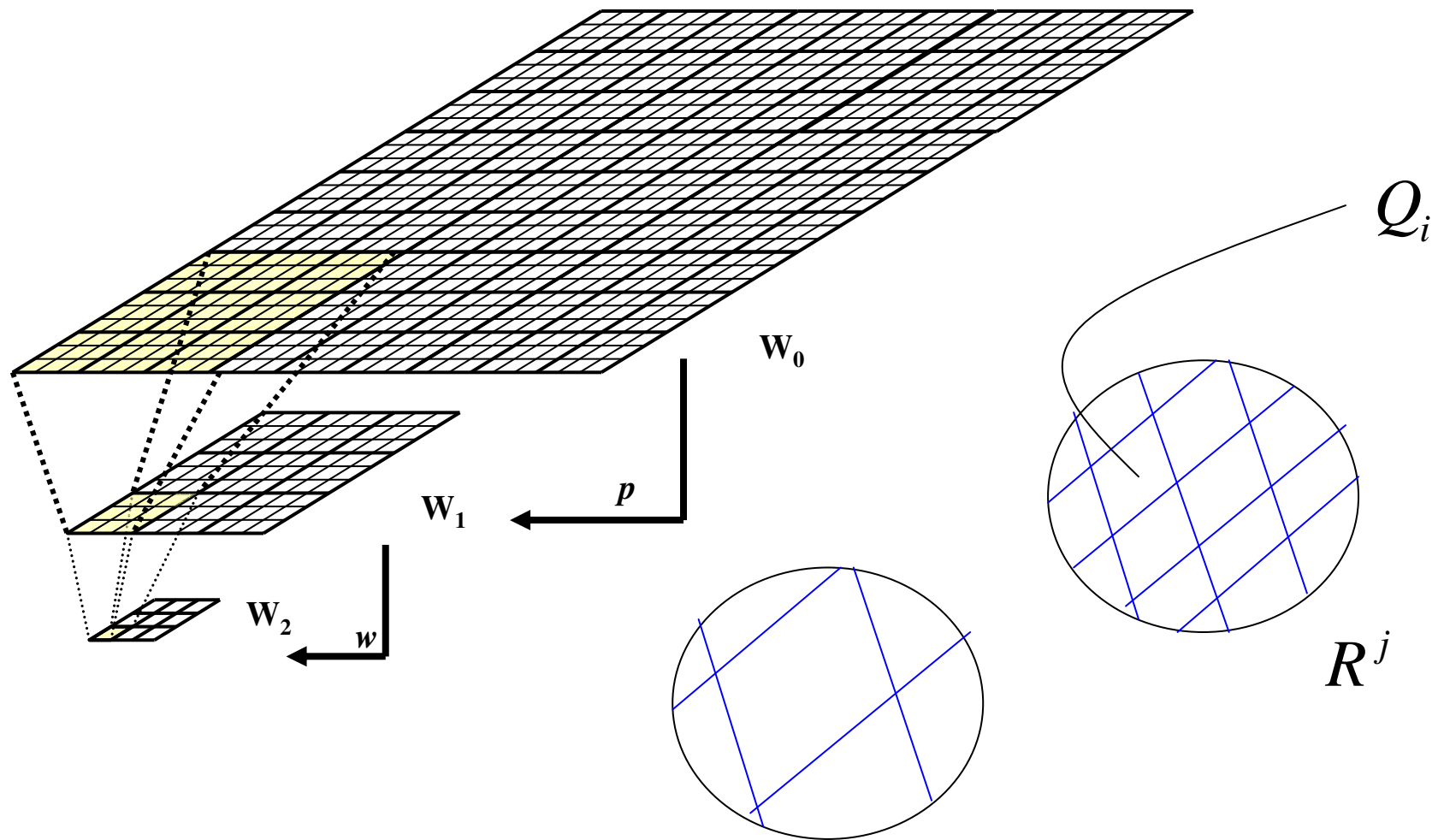
Distributions constraint estimator

$$N_{Q_i} = \sum_{k=1}^m c_{Q_i}(x_k^j) \quad c_{Q_i}(x^j) = 1 \Leftrightarrow x^j \in Q_i$$

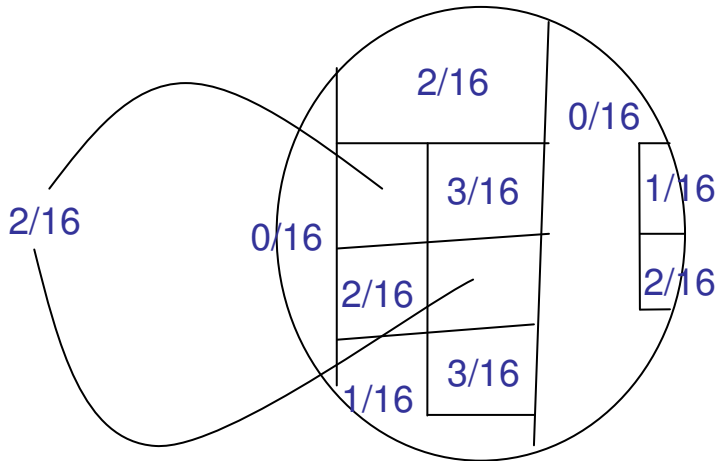
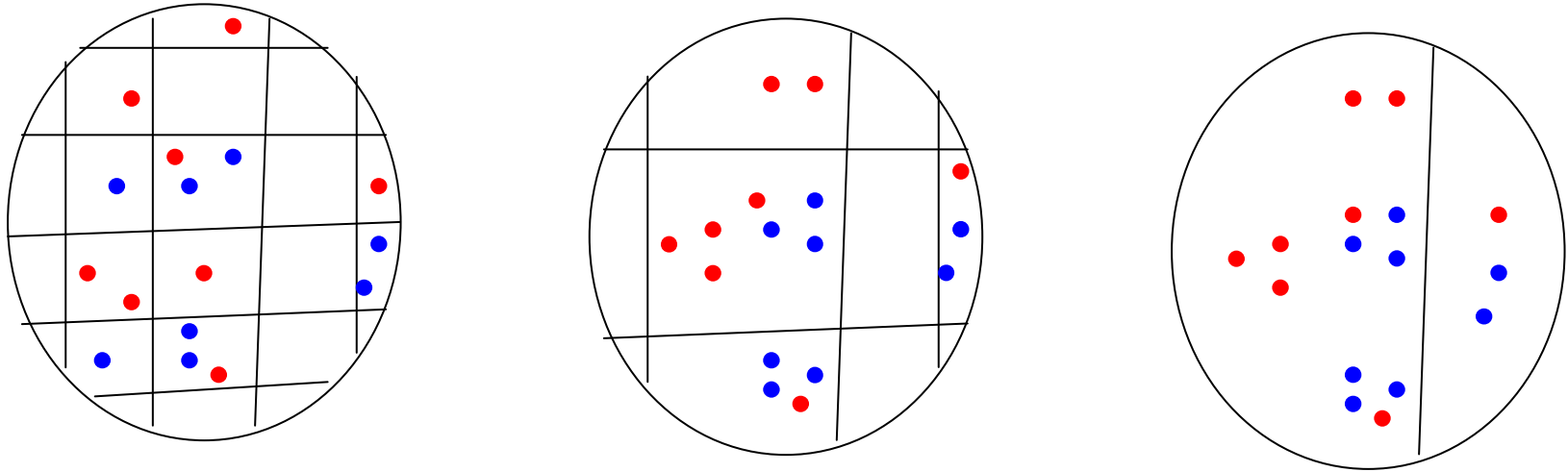
$$L_{Q_i,r} = \sum_{k=1}^m l_{Q_i,r}(x_k^j, y_k) \quad l_{Q_i,r}(x^j, y) = 1 \Leftrightarrow x^j \in Q_i \wedge y = r$$

$$\hat{P}(Q_i) = \frac{N_{Q_i}}{m} \quad \hat{P}(r | Q_i) = \frac{L_{Q_i,r}}{N_{Q_i}}$$

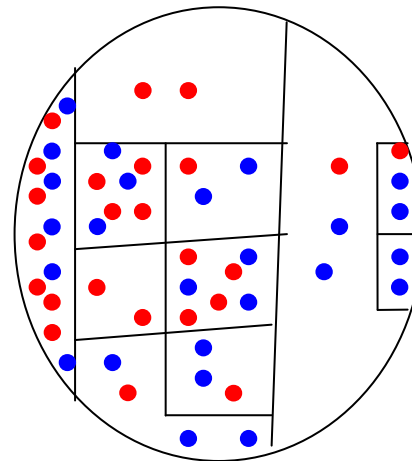
# Multi-resolution



# Example: multi-resolution estimation



$P(X)$



$P(Y|X)$

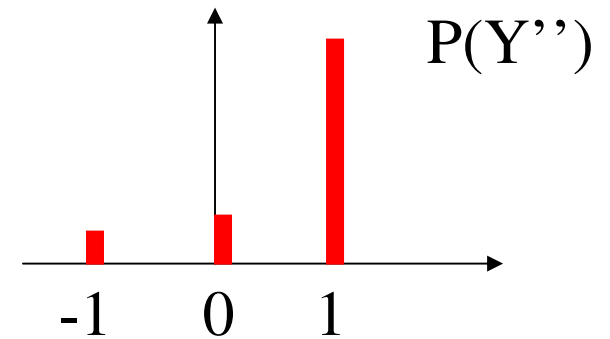
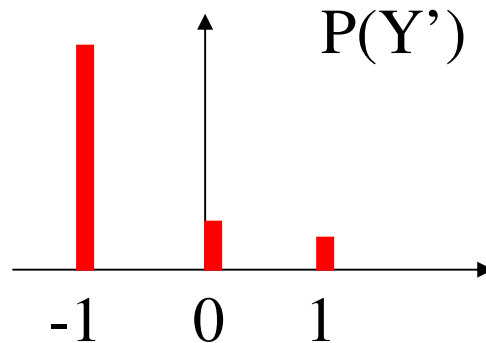
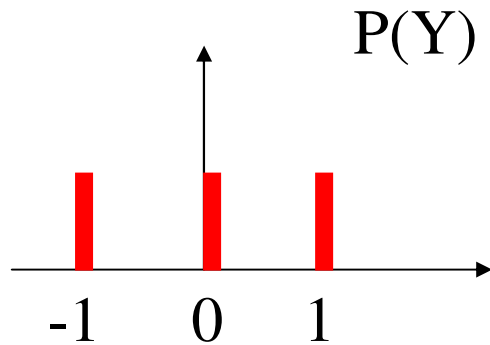
# Measuring estimator quality

## Entropy

$$H(Y) = - \sum_{y \in \{-1,0,1\}} P(y) \log P(y)$$

## Mutual information

$$I(X, Y) = H(Y) - H(Y | X) \geq 0$$



## Distributions of Y

$$H(Y) > H(Y') = H(Y'')$$

# Measuring estimator quality

Mean conditional entropy

$$E[H(Y | X)] = \sum P(X)H(Y | X)$$

Mean mutual information

$$E[I(X, Y)] = H(Y) - E[H(Y | X)]$$

Mean mutual information estimation

$$\hat{E}[H(Y | X)] = -\sum \hat{P}(X) \sum \hat{P}(Y | X) \log(\hat{P}(Y | X)).$$

$$\hat{E}[I(X, Y)] = H(\hat{Y}) - \hat{E}[H(Y | X)]$$

# Optimal estimator

Mean conditional entropy for a constraint distribution

$$E[H(Y | X)] = \sum_{i=1}^n P(Q_i) H(Y | Q_i)$$

Choose a space of partitions and choose the ones that minimize

$$\hat{E}[H(Y | X)]$$

Example: all possible partitions generated by a family of classifiers; partitions generated by a family of pyramid's

# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- **Architecture estimation**
- Malaria
- Cell Cycle
- Future works

$(x_1, y_1), (x_2, y_2), \dots, (x_m, y_m)$



**Feature selection**

$x \equiv x^j$

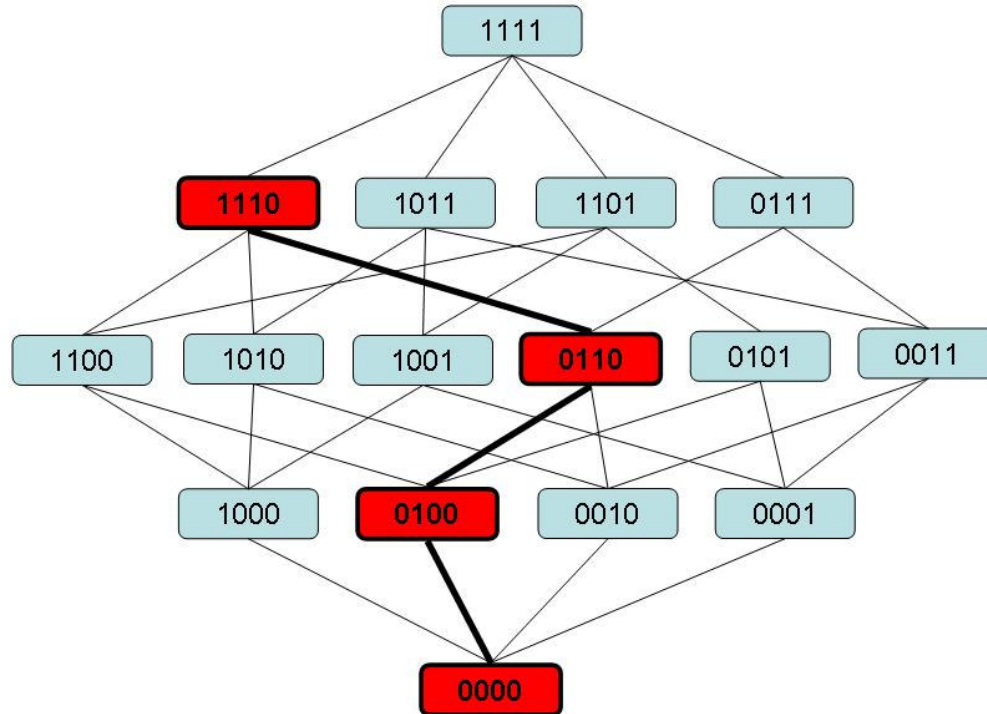


$A \subset \{1, 2, \dots, N\}$

$j = |A|$



# Feature selection



**Boolean Lattice**

# Feature selection

- A particular case of distribution estimation  
Equivalence classes built eliminating features
- Processing a node  
For each feature set several partitions may be compared
- Linear partitions  
May be used for generating the search space for a feature set
- Walking through the search space  
Some heuristics: SFS, SFFS, U-curve

## Inhibitory and excitatory interactions

Interaction parameter:  $a_i \in K = \{-k, \dots, 0, \dots, k\}$

$$\sum_{i=1}^N a_i x_i(t) = \sum_{i=1}^N a_i x_i(t + \Delta) \Rightarrow$$

$$P(r | x(t)) = P(r | x(t + \Delta))$$

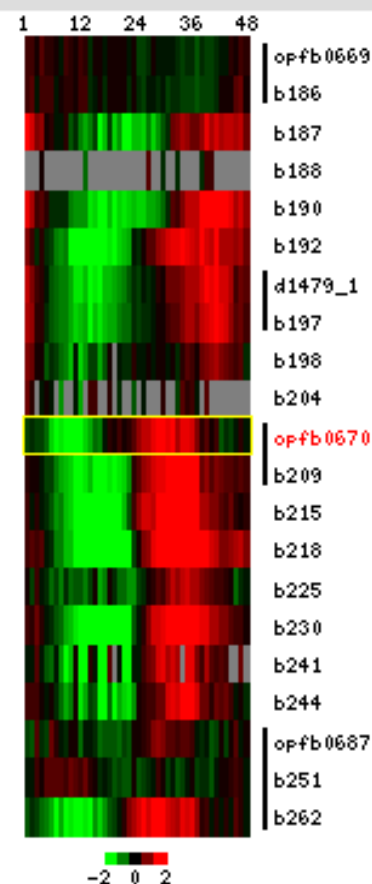
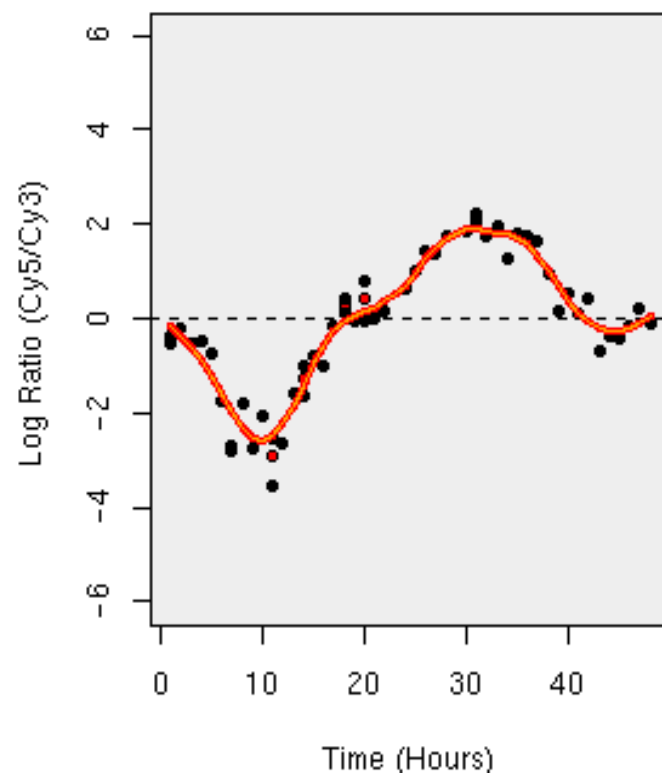
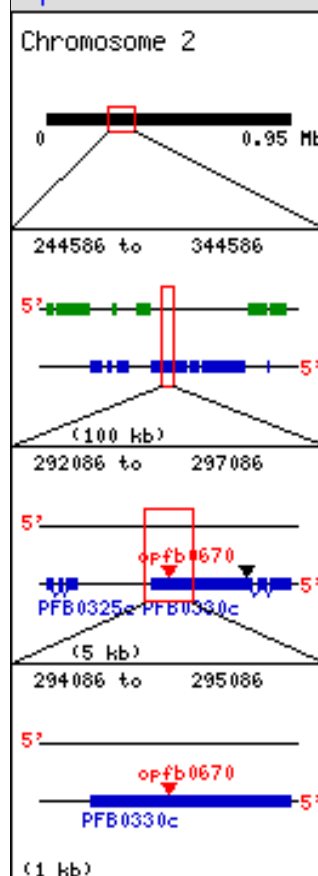
- A vector of parameters  $a \in K^N$  defines a partition
- The search space is generated by  $a \in K^N$
- The best vector gives interaction weight and signal

# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- **Malaria**
- Cell Cycle
- Future works



OligoID	Status	Maximum Hour	Minimum Hour	Amplitude Score (log2)	Score (%)	Phase (-Pi to +Pi)	CGH %3D7	Avg. Med. Intensity
<a href="#">opfb0670</a>	UNIQUE	30	10	4.5	87	0.06	89	3211.57



[← OLIGO →](#)

PlasmoDB ID

Description

[PFB0330c](#)

cysteine protease, putative

Oligo Sequence

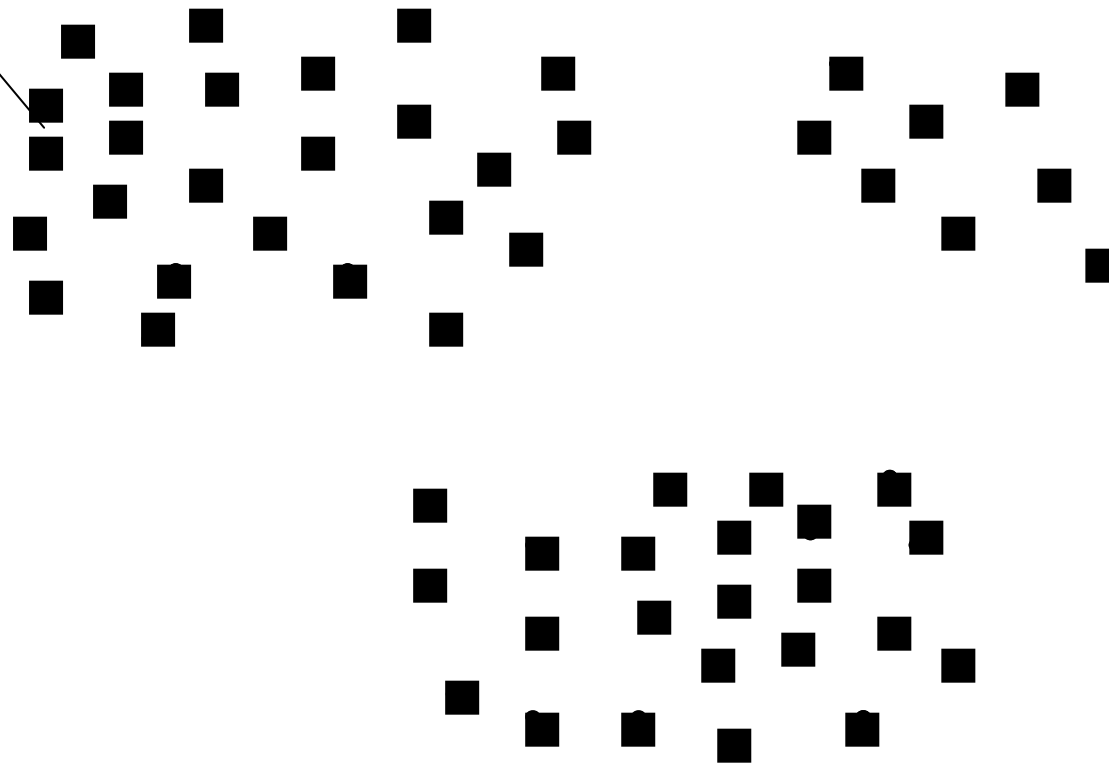
[BLAST @ PlasmoDB](#)

5'

CTGCCCAAGATGAGCCACCTACTGATAATGTAGAATCACAAGCAGAAAATAACAAAAAACAGAAATTTA

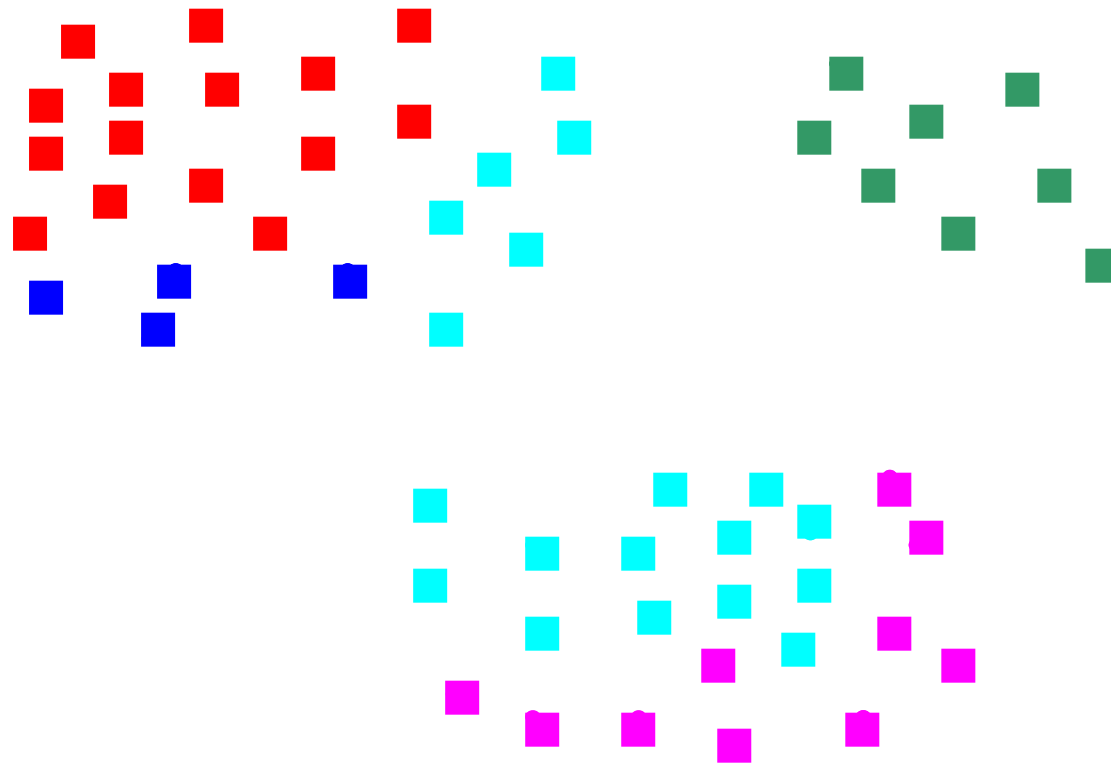
# Malaria parasite genes with almost sinusoidal signals

Sinuousoidal signals



DeRisi, 2003.

# Functional Classification



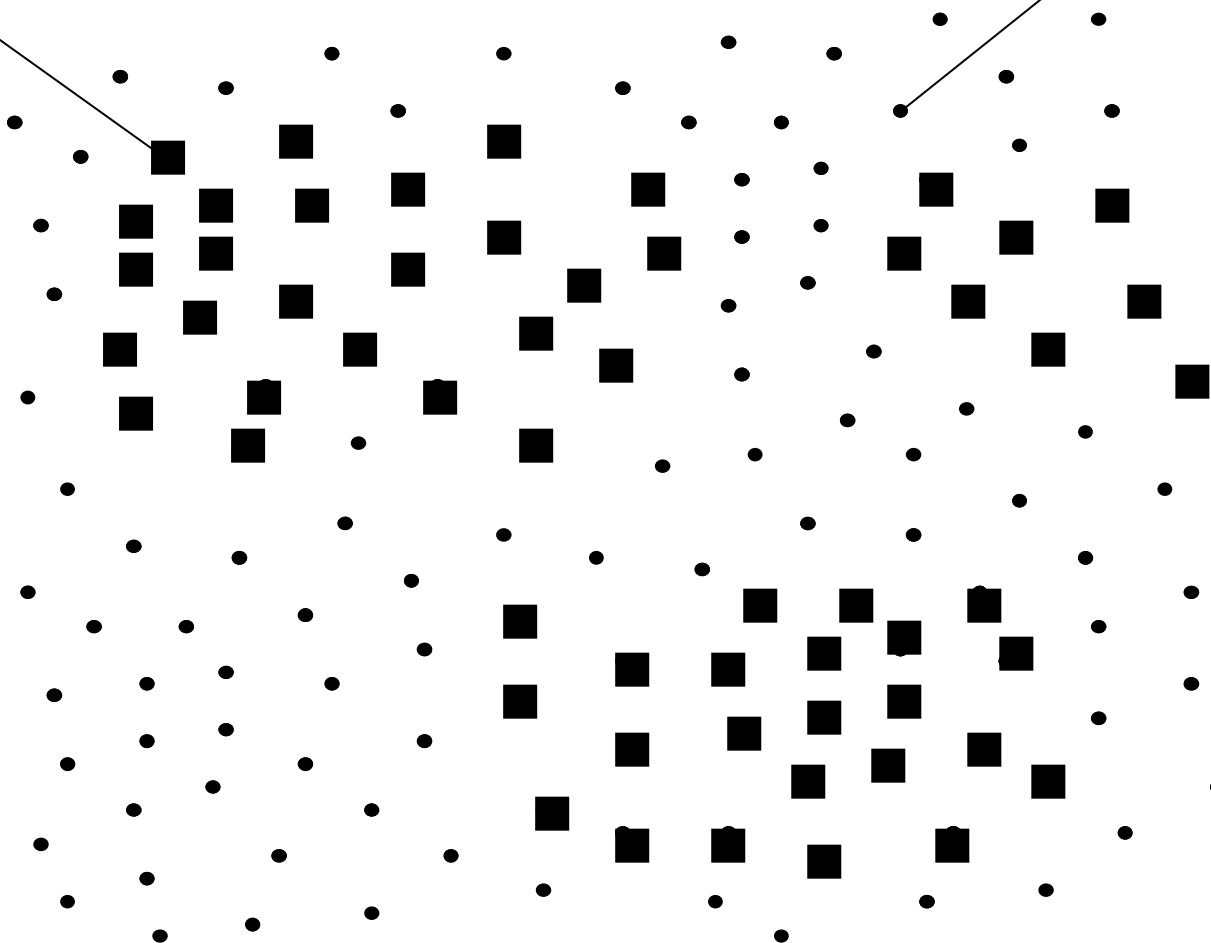
DeRisi, 2003.



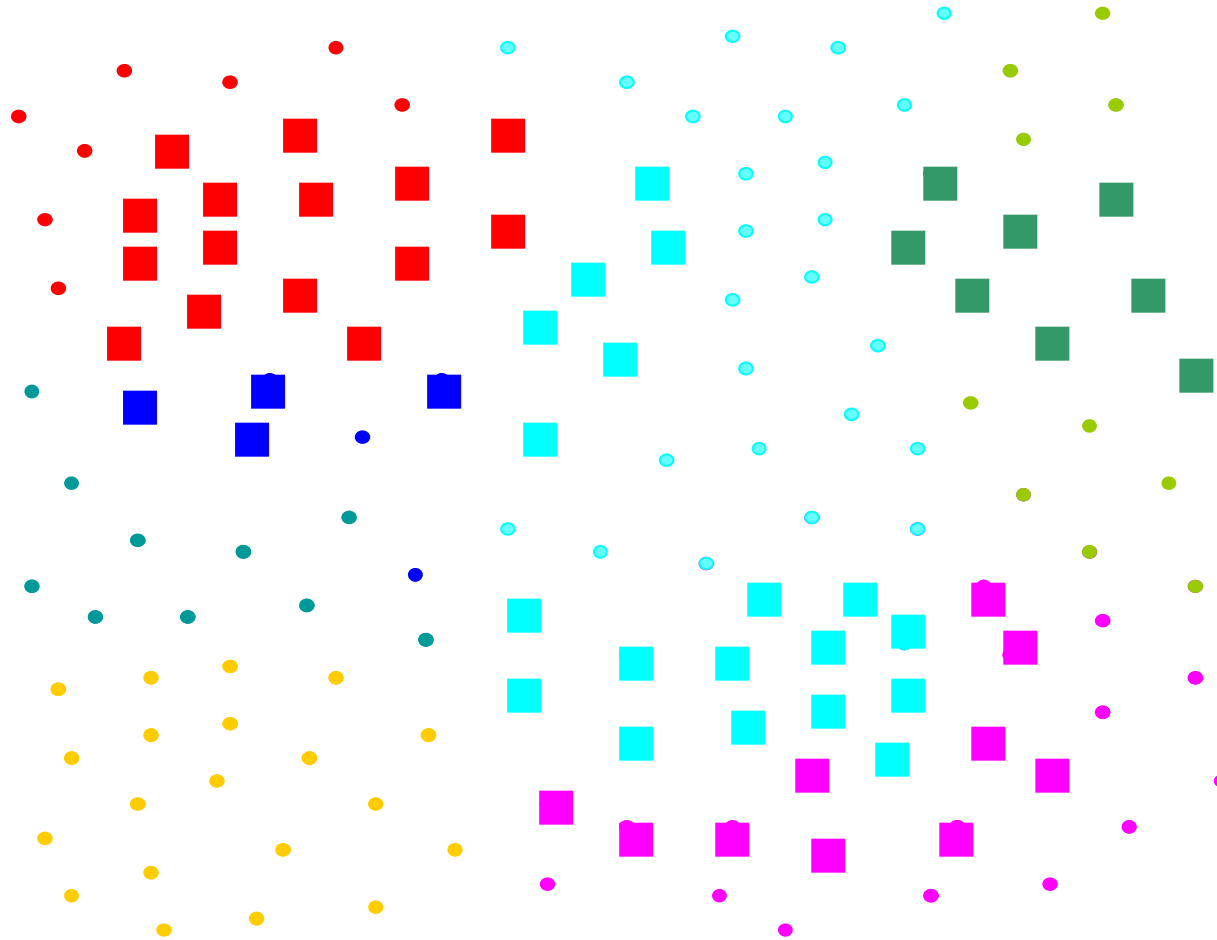
# Malaria parasite genes

Sinuousoidal signals

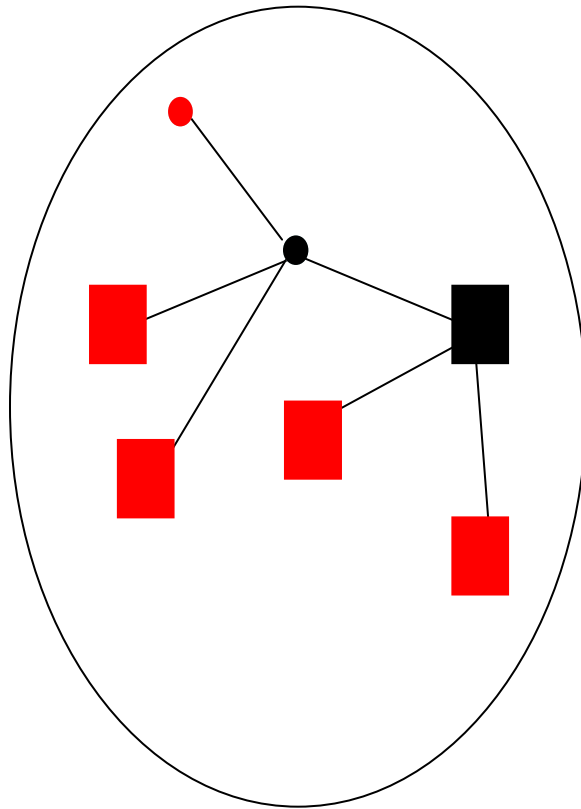
Non sinuousoidal signals



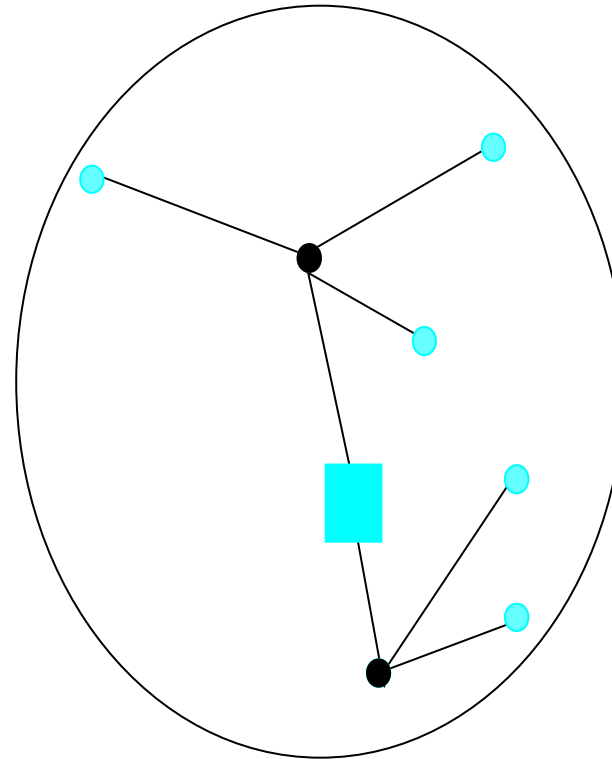
# Functional Classification



# Interaction Graph

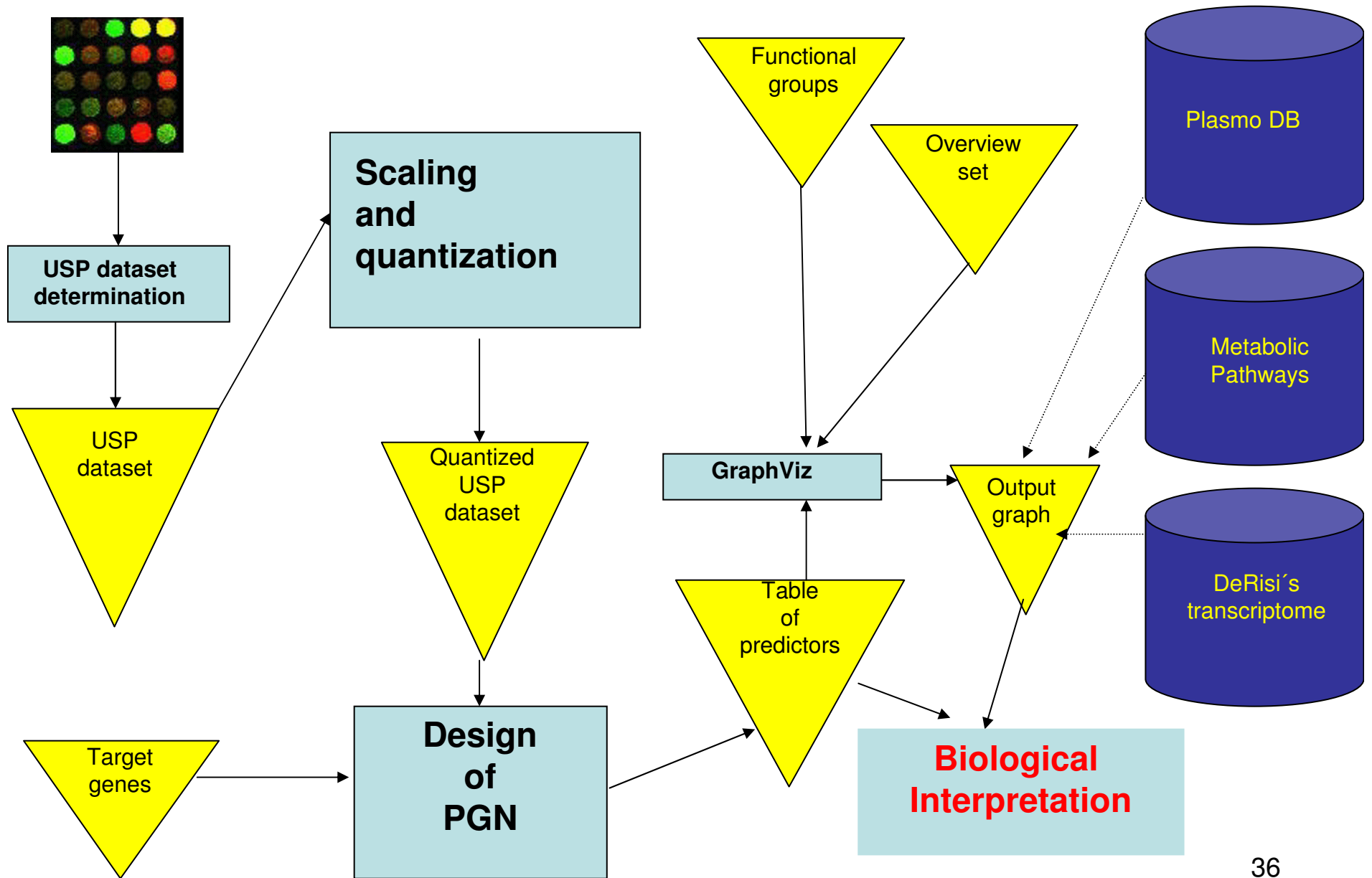


Glycolysis

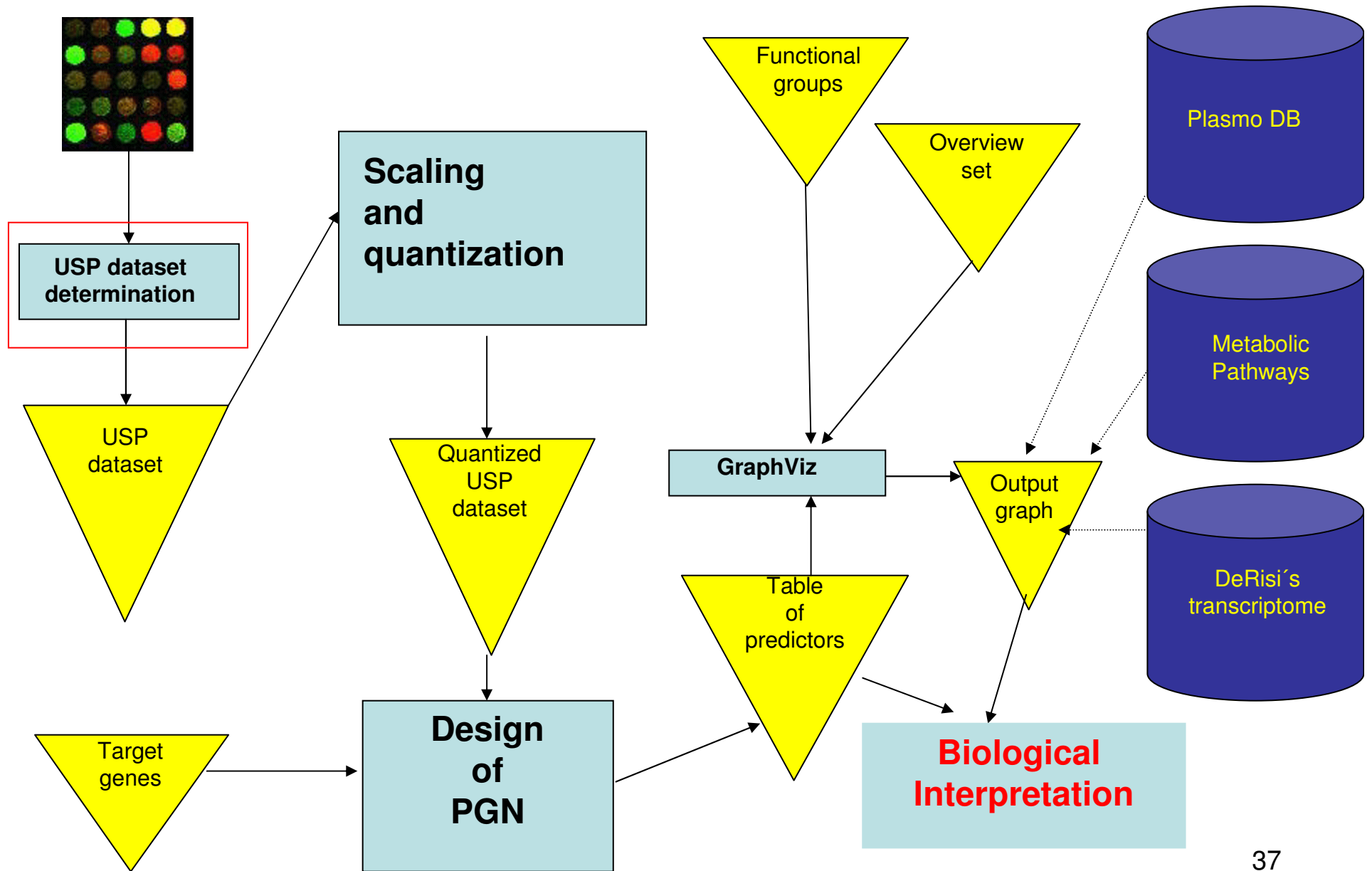


Apicoplast

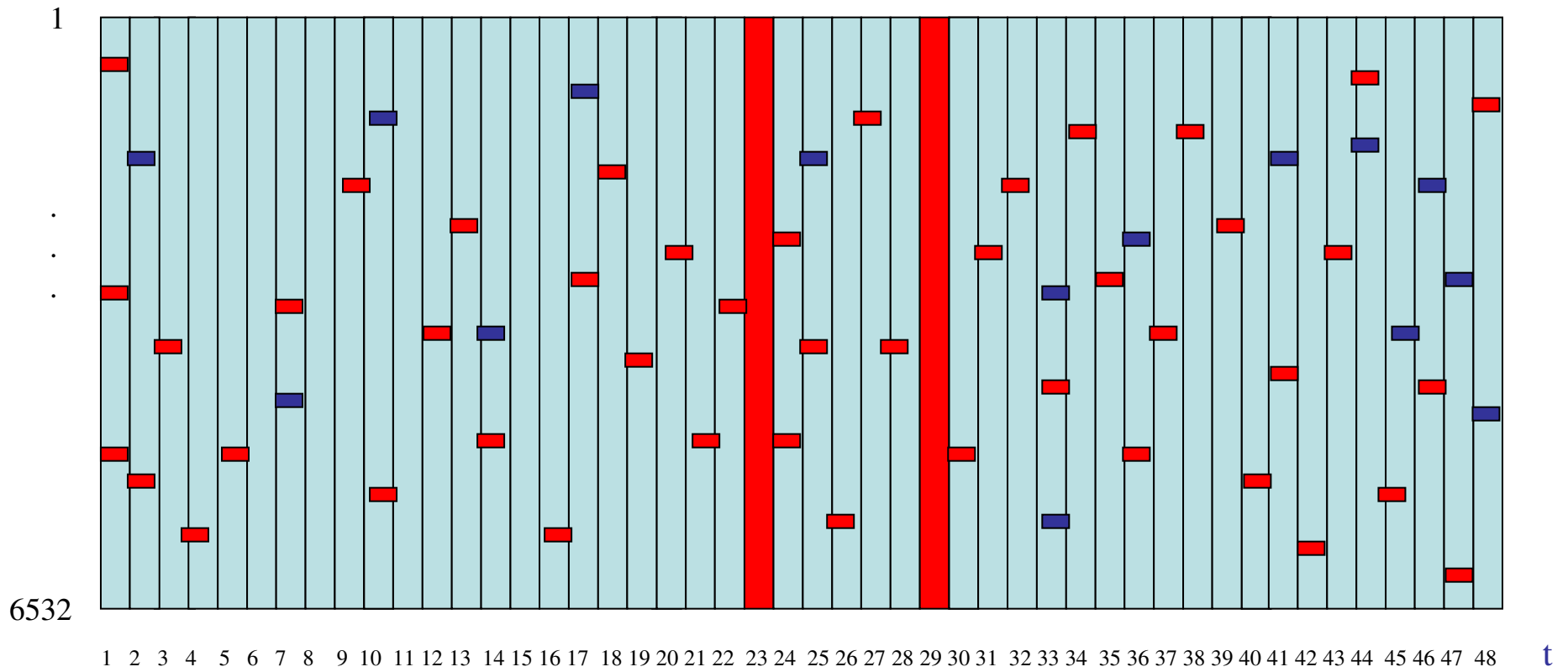
# System architecture



# System architecture



# Genes



Good spots



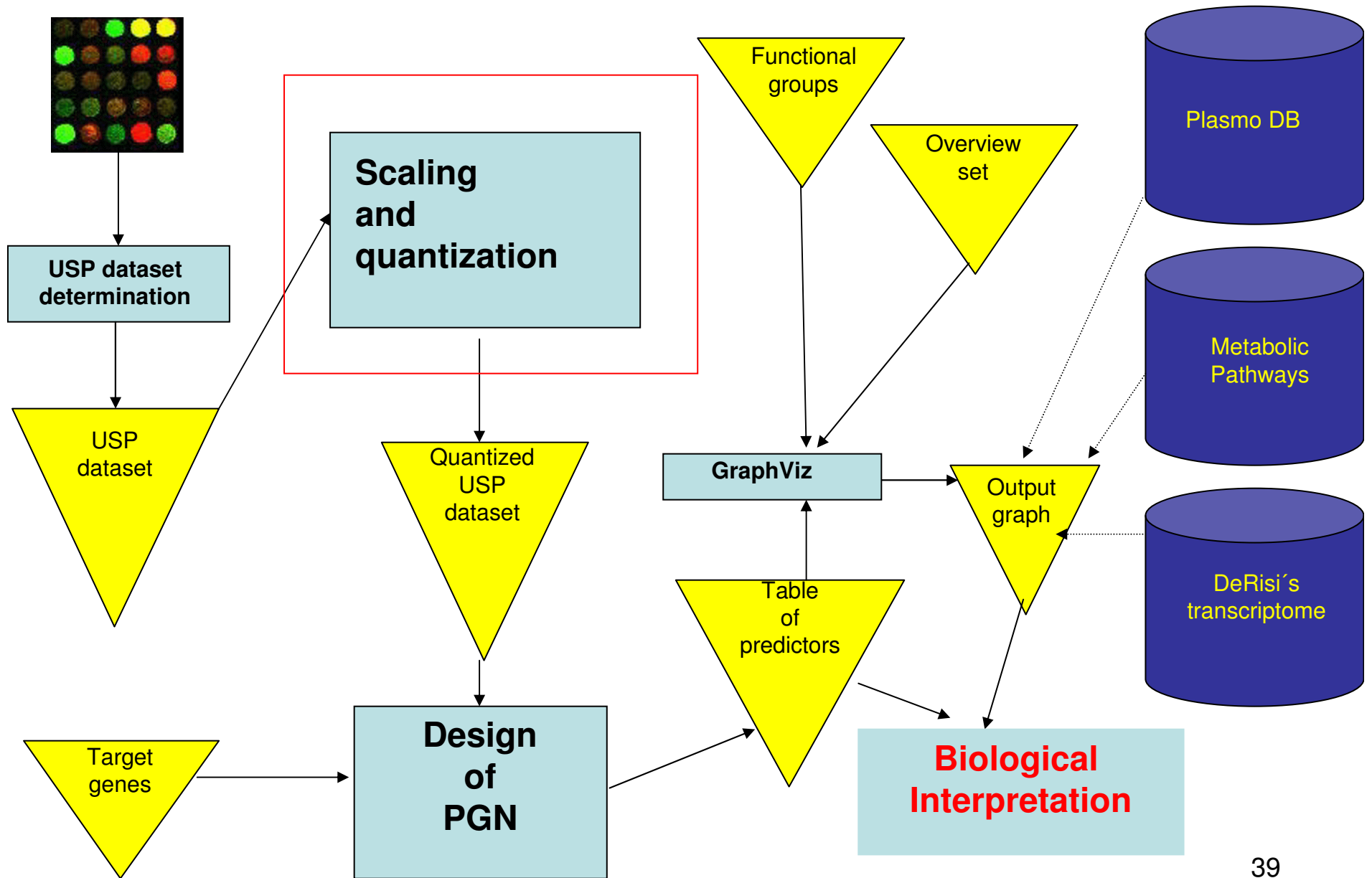
Weak spots



Bad spots

**NO INTERPOLATION**

# System architecture



# Scaling

For each  $i$ , estimate the mean  $\hat{E}[x_i[t]]$   
and standard deviation  $\hat{\sigma}[x_i[t]]$   
of the spots

Scale normalization of the spots

$$n_i[t] = \frac{x_i[t] - \hat{E}[x_i[t]]}{\hat{\sigma}[x_i[t]}}$$



# Quantization

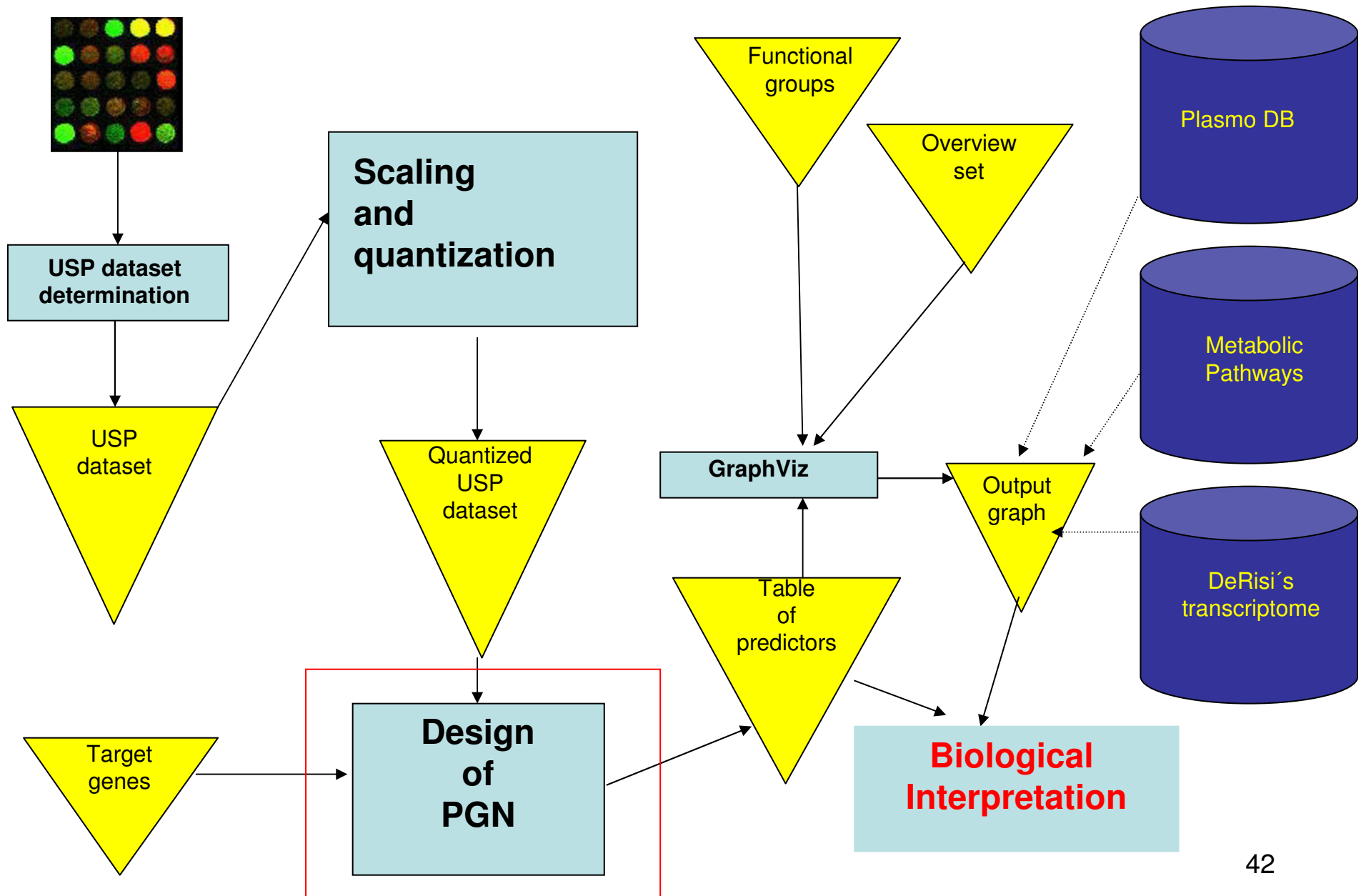
Let  $n_i^+[t]$  and  $n_i^-[t]$  denote, respectively, the normalized signals greater and lower than zero at  $t$ .

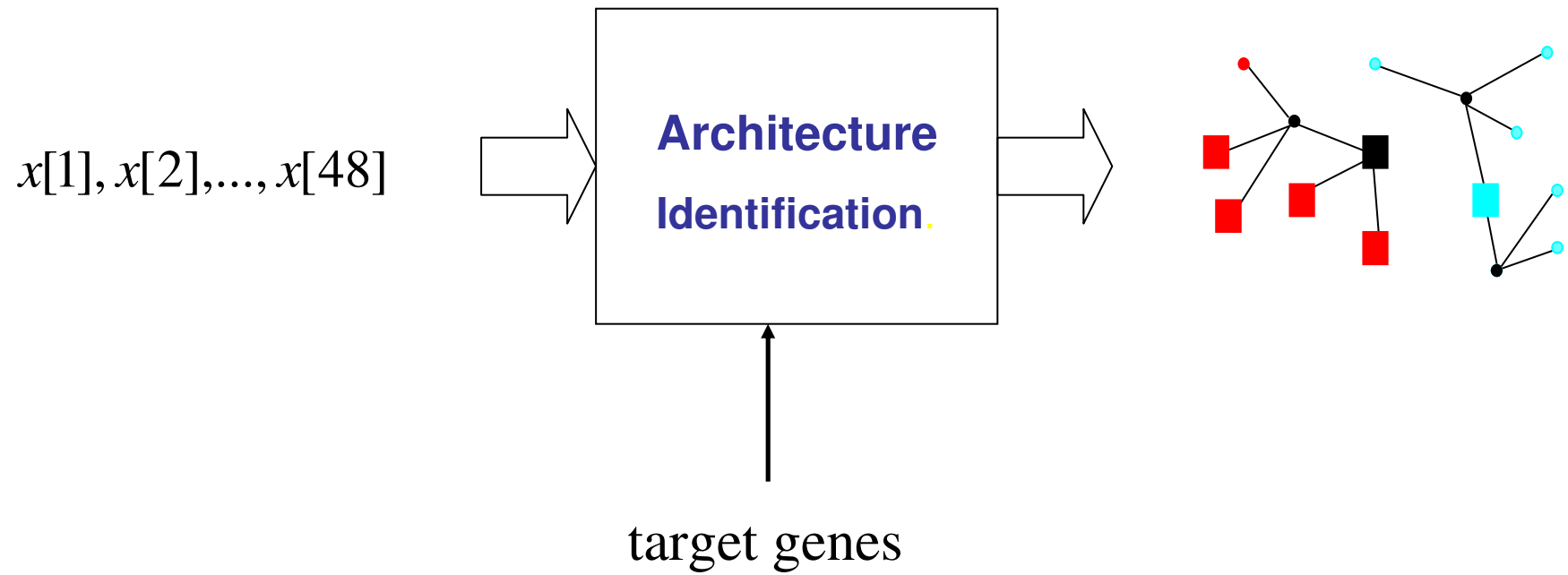
If  $n_i^+[t] > \hat{E}[n_i^+[t]]$ , then  $x_i[t] = +1$

If  $n_i^-[t] \geq \hat{E}[n_i^-[t]]$  and  $n_i^+[t] \leq \hat{E}[n_i^+[t]]$ , then  $x_i[t] = 0$

If  $n_i^-[t] < \hat{E}[n_i^-[t]]$ , then  $x_i[t] = -1$

# System architecture





## Estimation of $P(Y|X)$

Y: the target gene at  $t+1$ , that is,  $Y = x_i[t+1]$

X: the predictors at  $t$ , that is,  $X = (x_j[t], x_k[t])$

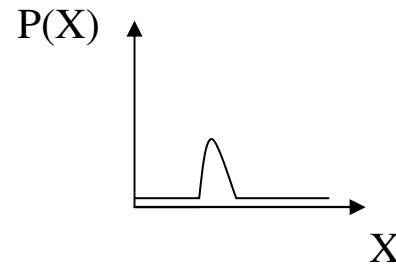
For a fixed parameter  $n$

If  $\#(X=(a,b)) \geq n$ , then  $\hat{P}(Y=c | X=(a,b)) = \frac{\#((Y=c) \wedge X=(a,b))}{\#(X=(a,b))}$

If  $\#(X=(a,b)) < n$ , then  $\hat{P}(Y | X=(a,b))$  is uniform

## Estimation of $P(X)$ for a fixed parameter $n$

$$X = (x_j[t], x_k[t])$$



$$N^+ = \sum_{\#(X=(a,b)) \geq n, \forall (a,b)} \#(X=(a,b))$$

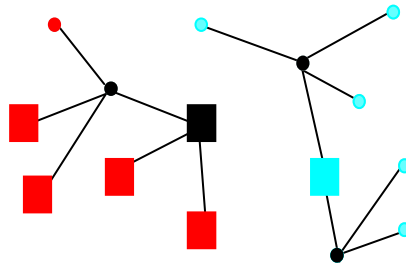
$$N^- = \sum_{\#(X=(a,b)) < n, \forall (a,b)} \#(X=(a,b))$$

$$\text{If } \#(X=(a,b)) \geq n, \text{ then } \hat{P}(X=(a,b)) = \frac{N^+}{N^- + N^+} \times \frac{\#(X=(a,b))}{N^+}$$

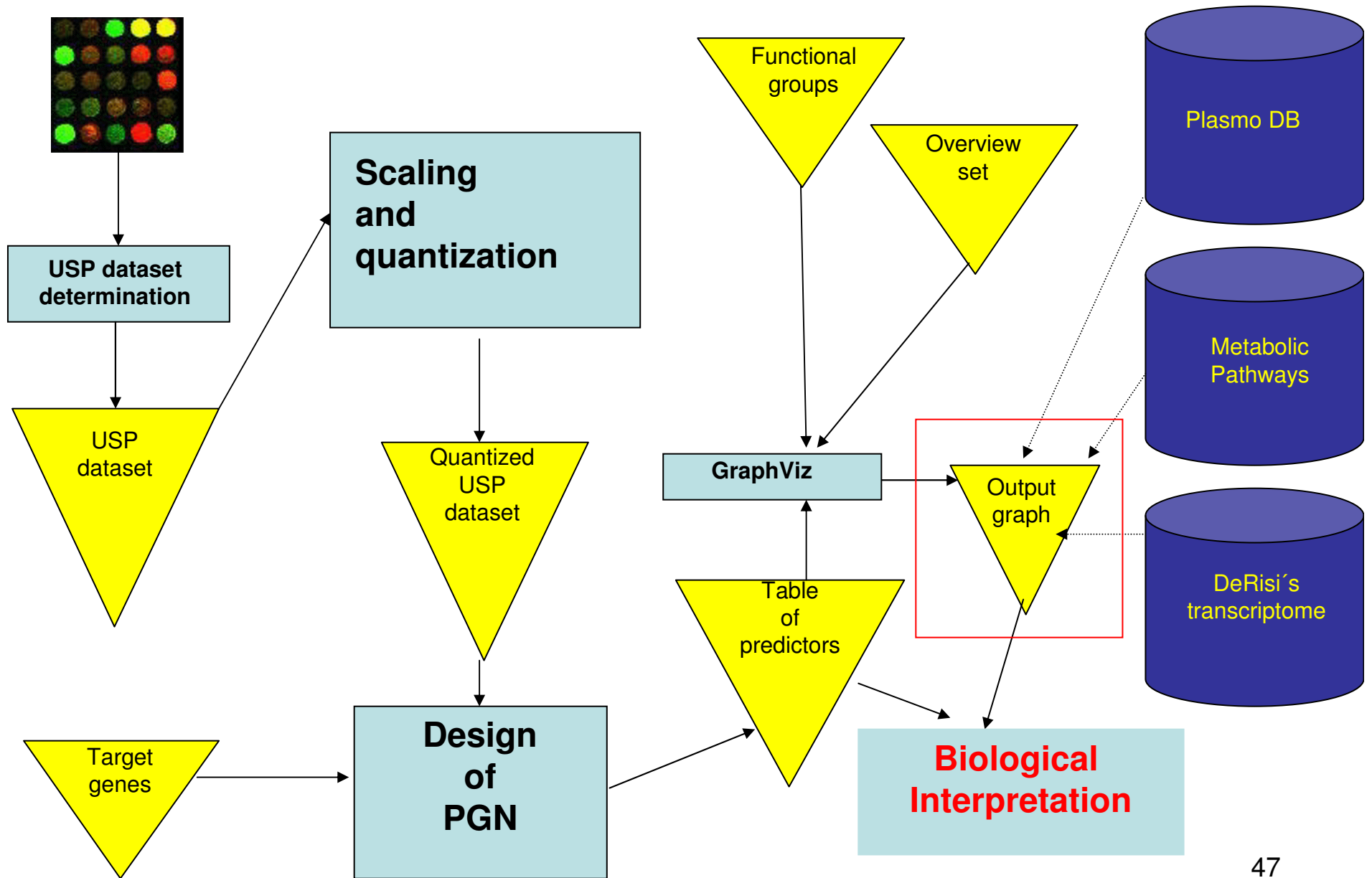
$$\text{If } \#(X=(a,b)) < n, \text{ then } \hat{P}(X=(a,b)) = \frac{N^-}{N^- + N^+} \times \frac{1}{3^2 - |\{(a,b) : \#(X=(a,b)) \geq n\}|}$$

## Building Interaction Graphs

- For each target gene, rank the couples of all genes by their estimated mutual information and sample size;
- When two mutual information are equal, the one estimated from a larger sample comes first;
- Choose the best couples;
- Design the interaction graph



# System architecture

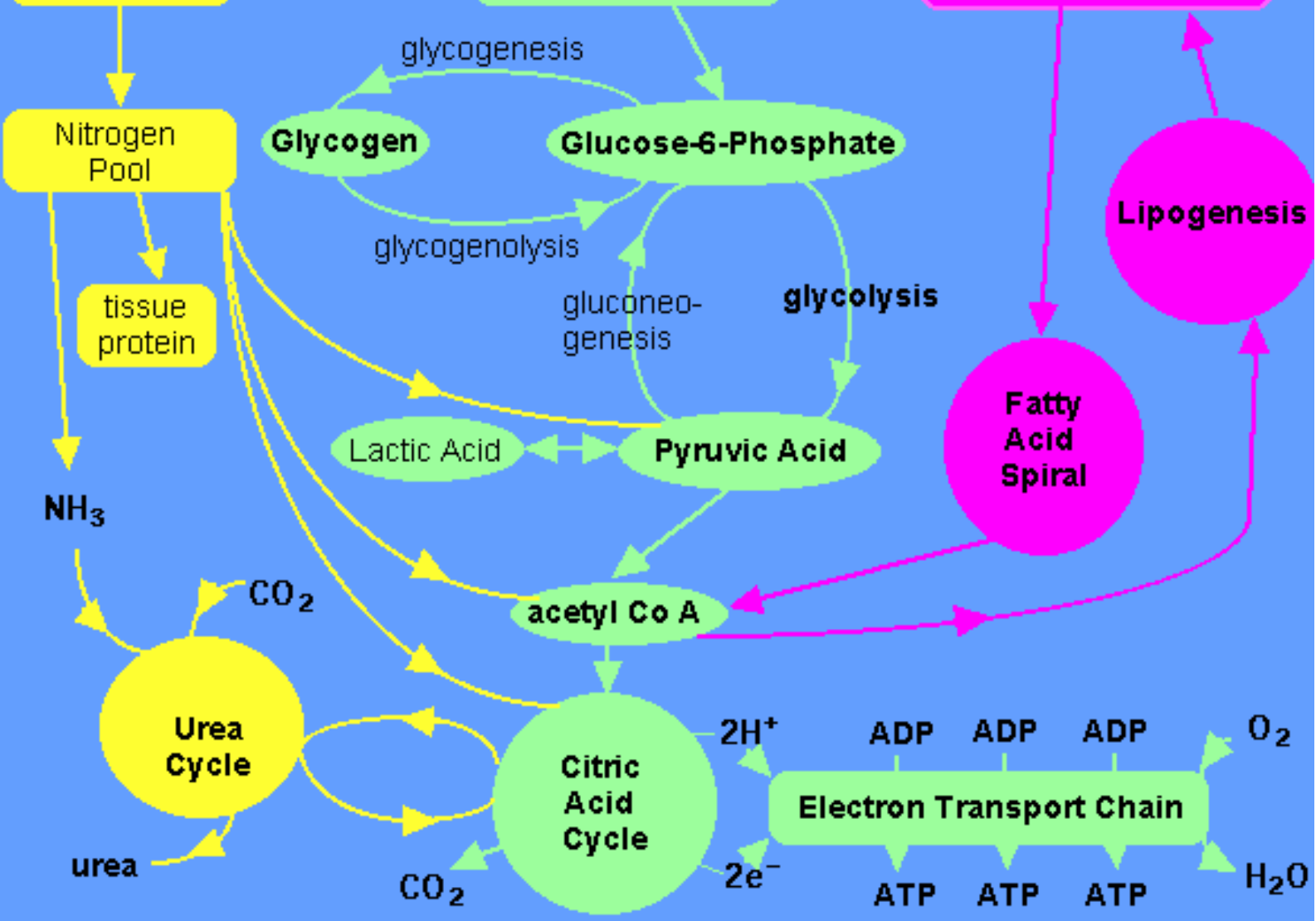


# Metabolism Summary

**Proteins**  
amino acids

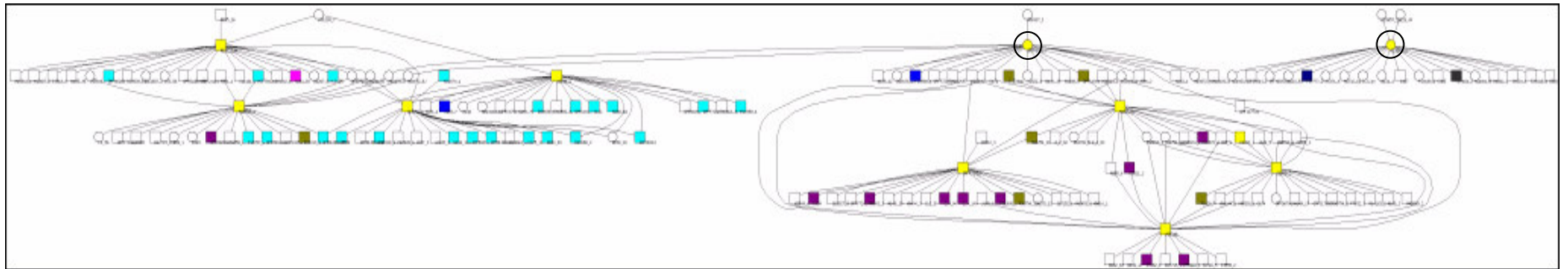
**Carbohydrates**  
glucose, fructose, galactose




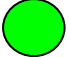






**Fats and Lipids**  
fatty acid, glycerol



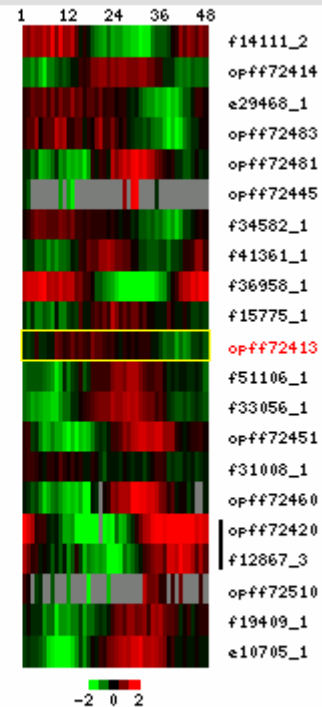
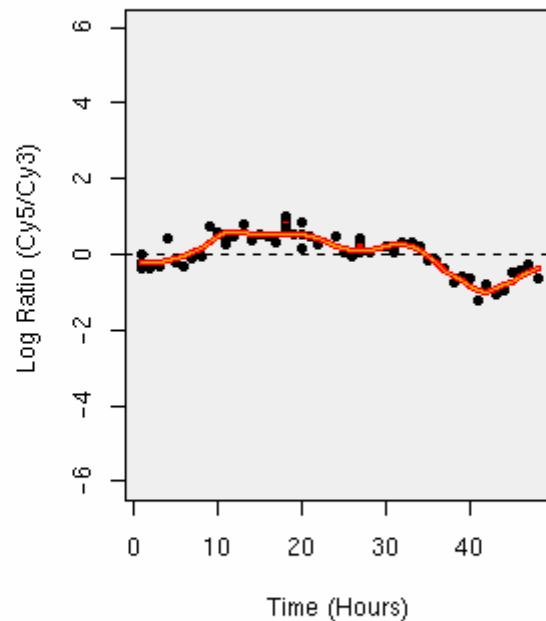
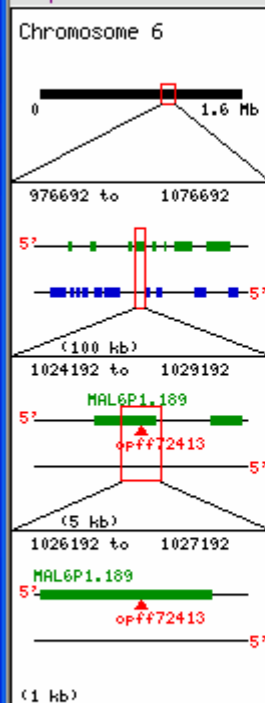


## Glycolytic PGN network (single genes)



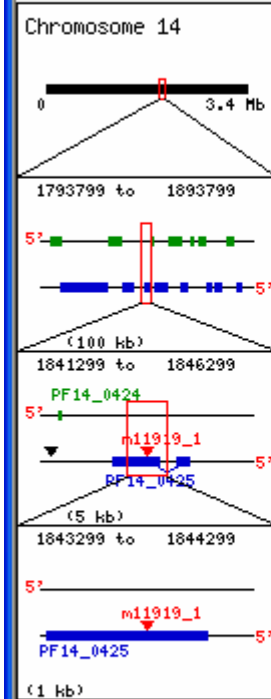
- |   |                                  |   |                                     |
|---|----------------------------------|---|-------------------------------------|
|    | <b>glycolysis</b>                |    | <b>proteasome</b>                   |
|   | <b>transcription machinery</b>   |   | <b>plastid genome</b>               |
|  | <b>cytoplasmic translation</b>   |  | <b>merozoite invasion (kinases)</b> |
|  | <b>ribonucleotide synthesis</b>  |  | <b>actin myosin motors</b>          |
|  | <b>deoxynucleotide synthesis</b> |  | <b>early ring transcripts</b>       |
|  | <b>DNA replication</b>           |   |                                     |

OligoID	Status	Maximum Hour	Minimum Hour	Amplitude (log2)	Score (%)	Phase (-Pi to +Pi)	CGH %3D7	Avg. Med. Intensity
opff72413	UNIQUE	12	42	1.6	69	2.05	82	14088.2

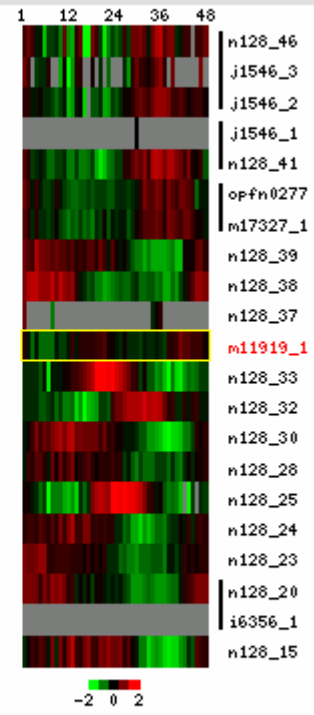
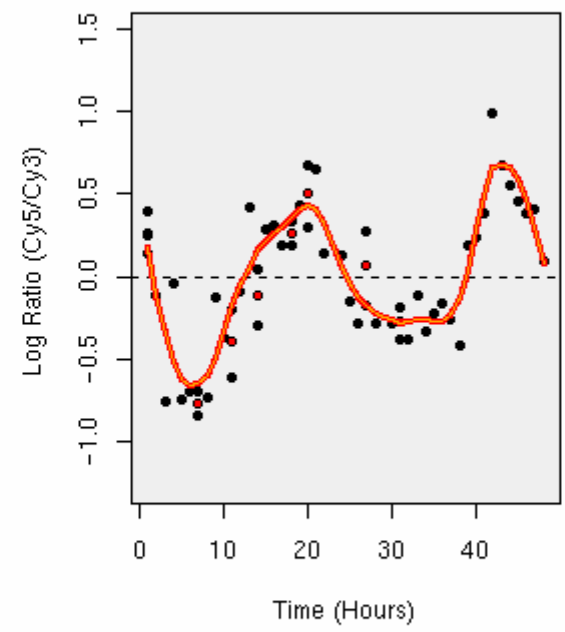


PlasmoDB ID	Description
MAL6P1.189	hexokinase

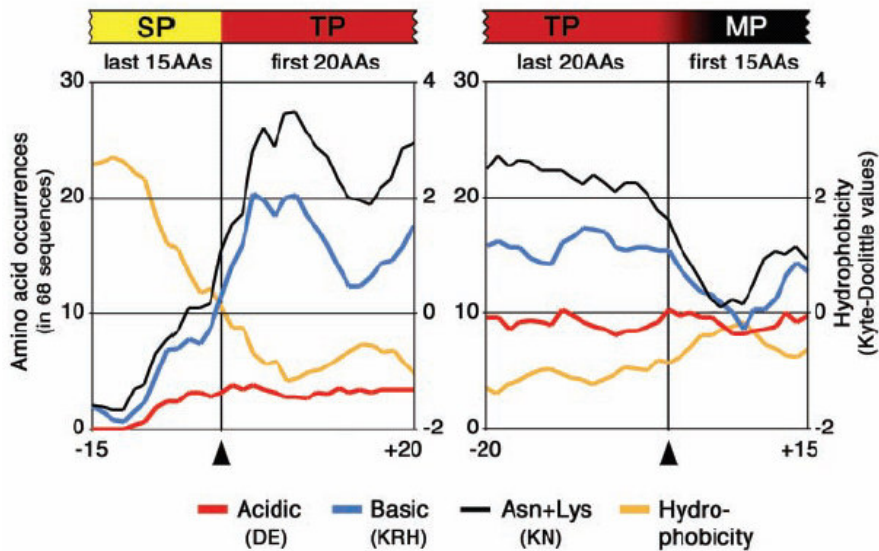
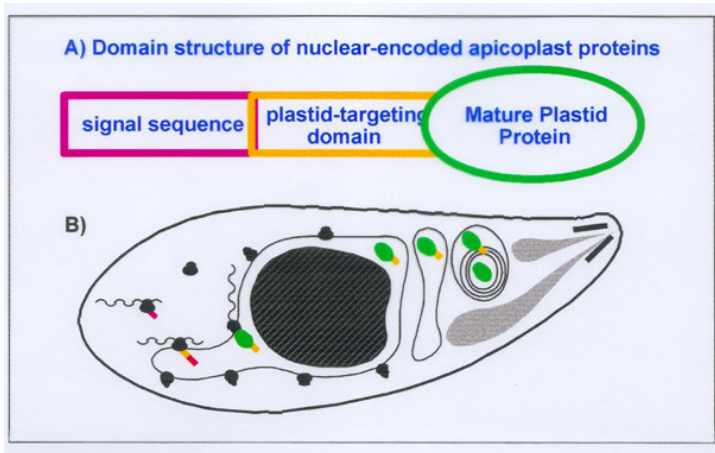
OligoID	Status	Maximum Hour	Minimum Hour	Amplitude (log2)	Score (%)	Phase (-Pi to +Pi)	CGH %3D7	Avg. Med. Intensity
m11919_1	UNIQUE	43	6	1.3	72	-1.53	80	50799.8



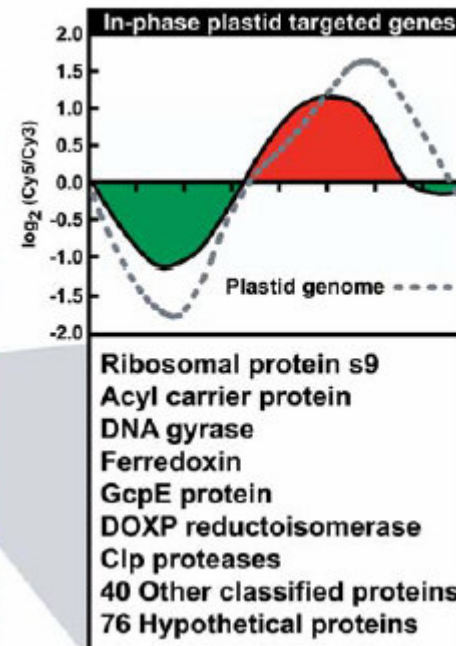
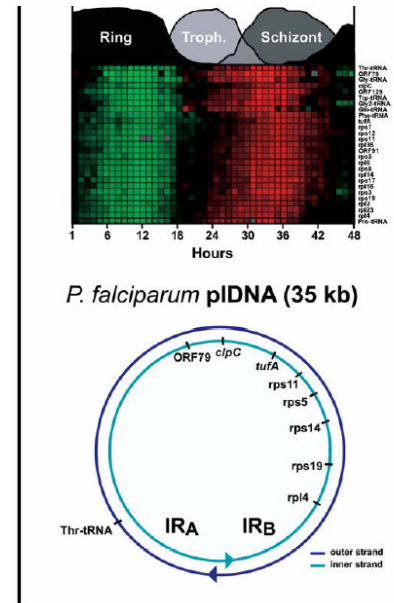
← OLIGO →



PlasmoDB ID	Description
PF14_0425	fructose-bisphosphate aldolase

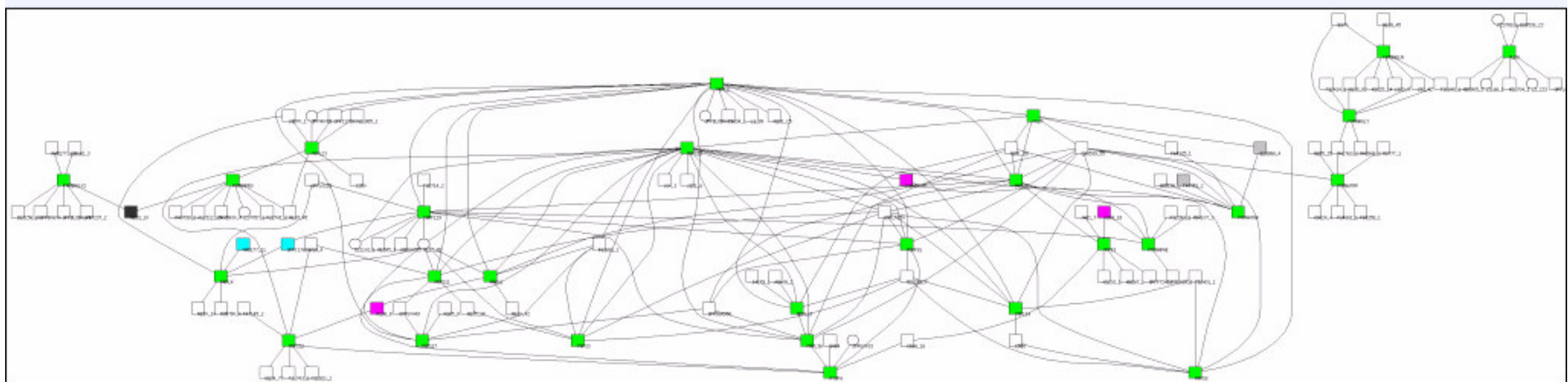


**550 apicoplast proteins**

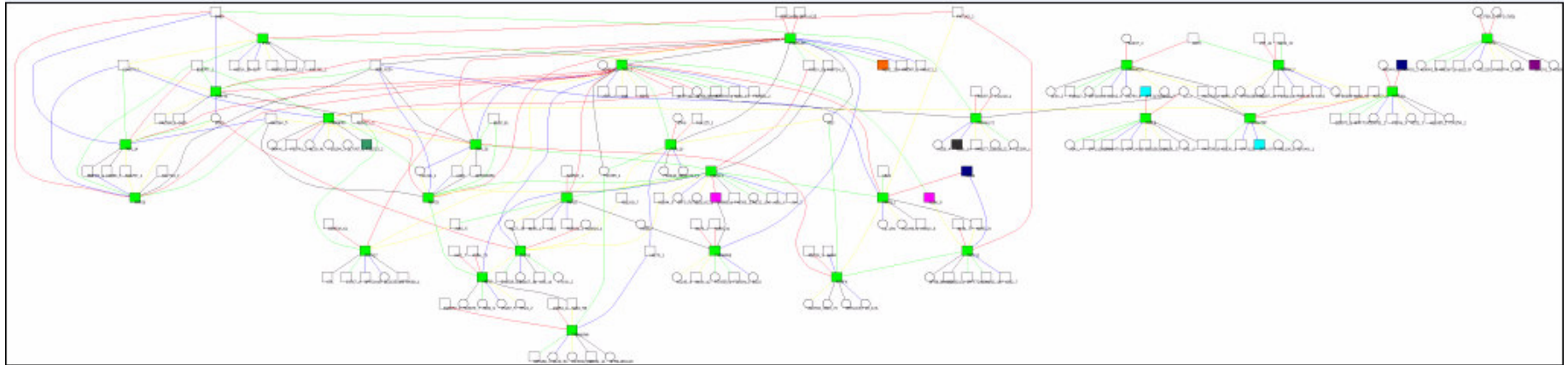


**124 apicoplast proteins**

## Apicomplast PGN network (single genes)

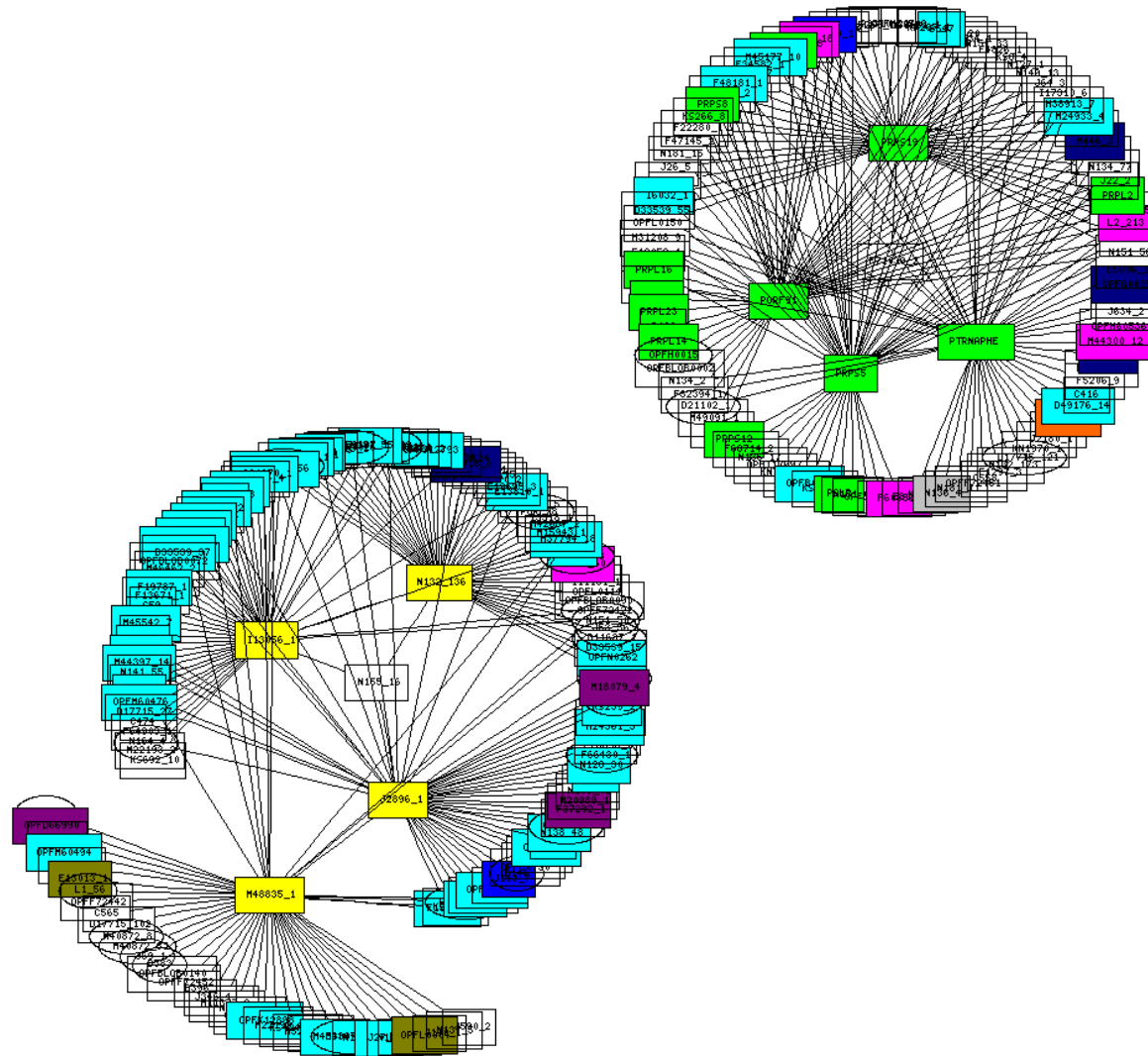


## Apicoplast PGN network (double genes)



- |   |                                  |   |                               |
|---|----------------------------------|---|-------------------------------|
|    | <b>glycolysis</b>                |    | <b>proteasome</b>             |
|  | <b>transcription machinery</b>   |  | <b>plastid genome</b>         |
|  | <b>cytoplasmic translation</b>   |  | <b>merozoite invasion</b>     |
|  | <b>ribonucleotide synthesis</b>  |  | <b>actin myosin motors</b>    |
|  | <b>deoxynucleotide synthesis</b> |  | <b>early ring transcripts</b> |
|  | <b>DNA replication</b>           |   |                               |

# Subsystems identification





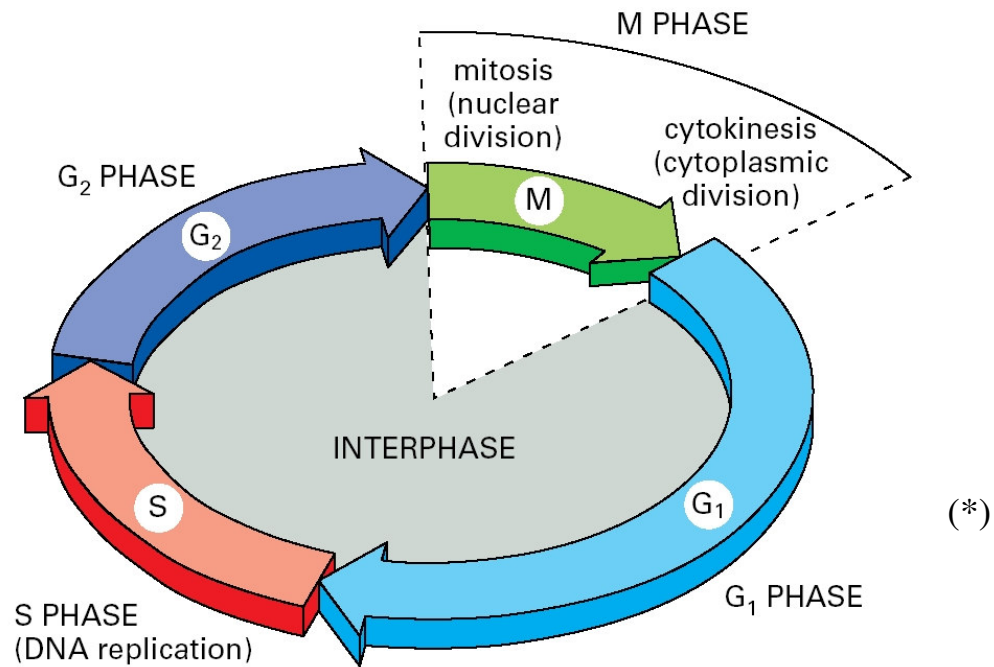




# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- Malaria
- **Cell Cycle**
- Future works

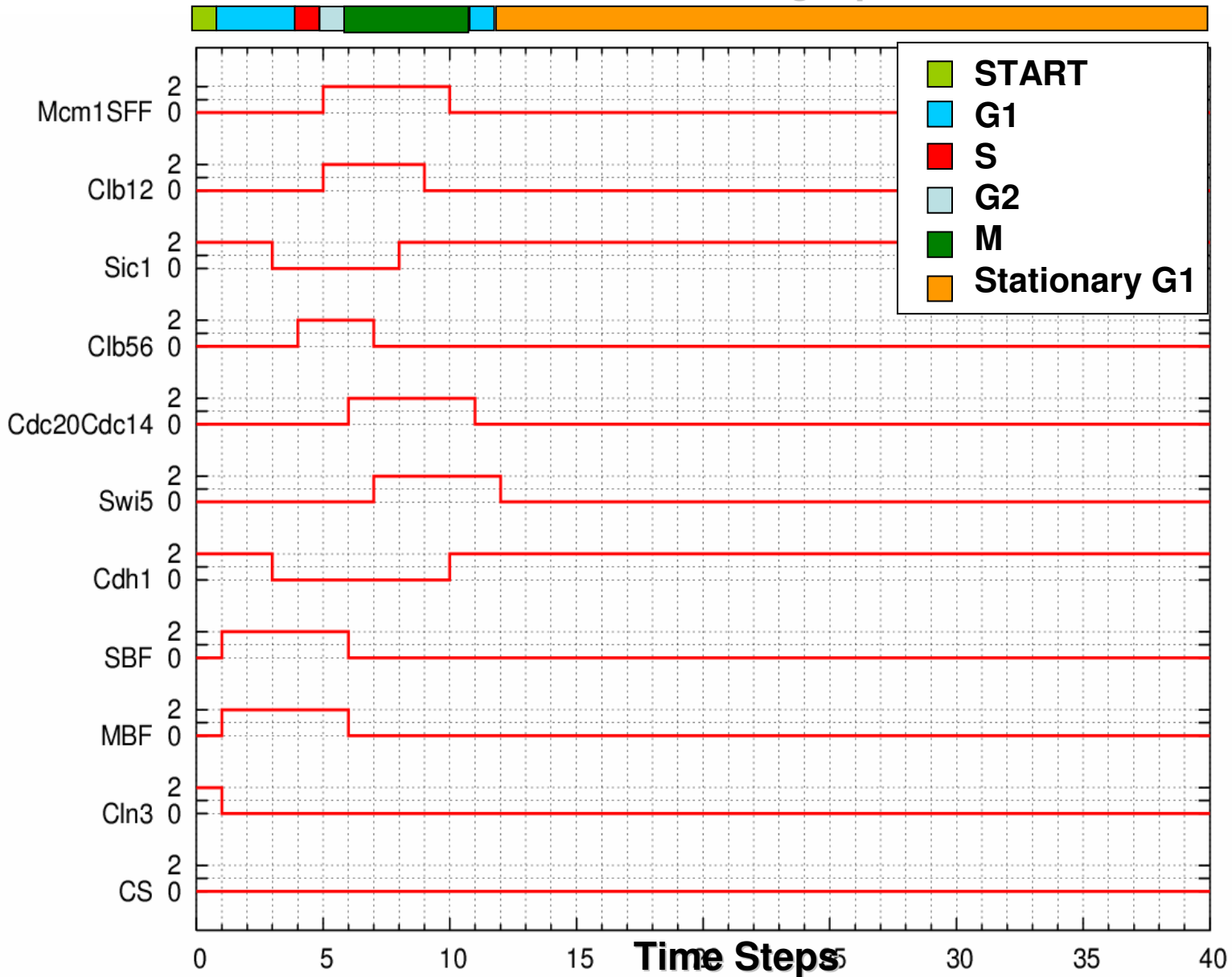
# Phases of the Cell Cycle





# Deterministic

One single pulse of  $CS = 2$  at  $t = -1$

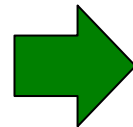


**Binary**

**1** → **2**  
**0** → **0**

**3 Levels (0, 1, 2)**

$\sum_j a_{ij} S_j(0)$	$S_j(t+1)$	
	$S_j(t)=0$	$S_j(t)=1$
⋮	⋮	⋮
3	1	1
2	1	1
1	1	1
0	0	1
-1	0	0
-2	0	0
-3	0	0
⋮	⋮	⋮



$\sum_j a_{ij} S_j(0)$	$y_j(t+1)$		
	$S_j(t)=0$	$S_j(t)=1$	$S_j(t)=2$
⋮	⋮	⋮	⋮
3	2	2	2
2	2	2	2
1	1	2	2
0	0	1	2
-1	0	0	1
-2	0	0	0
-3	0	0	0
⋮	⋮	⋮	⋮

# Stochastic Transition Function

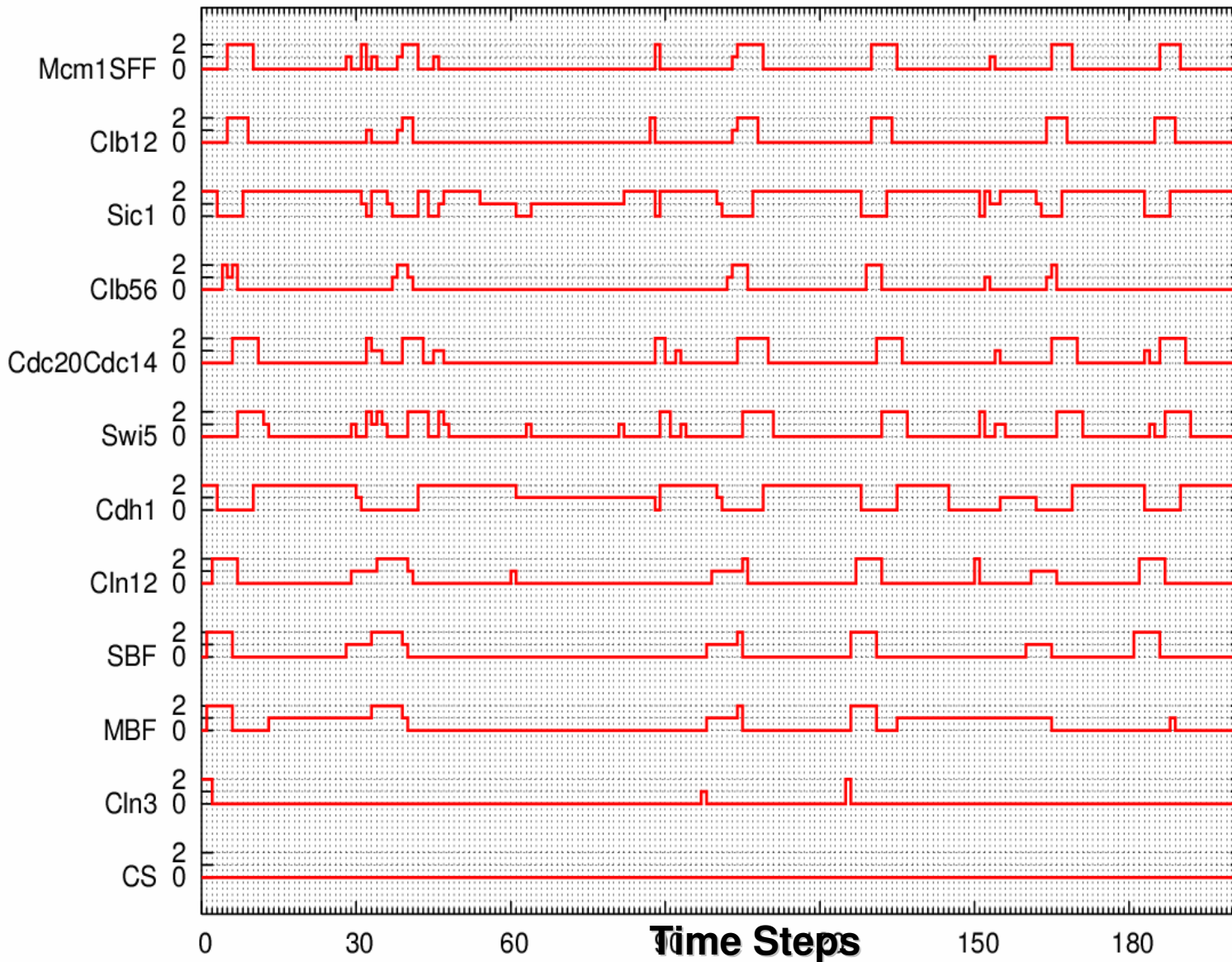
$\sum_j a_{ij} S_j(0)$	$y_j(t+1)$		
	$S_j(t) = 0$	$S_j(t) = 1$	$S_j(t) = 2$
$\vdots$	$\vdots$	$\vdots$	$\vdots$
<b>3</b>	<b>2</b>	<b>2</b>	<b>2</b>
<b>2</b>	<b>2</b>	<b>2</b>	<b>2</b>
<b>1</b>	<b>1</b>	<b>2</b>	<b>2</b>
<b>0</b>	<b>0</b>	<b>1</b>	<b>2</b>
<b>-1</b>	<b>0</b>	<b>0</b>	<b>1</b>
<b>-2</b>	<b>0</b>	<b>0</b>	<b>0</b>
<b>-3</b>	<b>0</b>	<b>0</b>	<b>0</b>
$\vdots$	$\vdots$	$\vdots$	$\vdots$

$$x_i(t+1) = \begin{cases} y_i(t+1) & \text{with } P = 0.99 \\ a & \text{with } P = 0.005 \\ b & \text{with } P = 0.005 \end{cases}$$

where  $a, b \in \{0, 1, 2\} - \{y_i(t+1)\}$

# PGN with $P = 0.99$

One single pulse of  $CS = 2$  at  $t = -1$



# Our gene model

Total input signal driving a generic variable

$$x_i(t) \in \{0, 1, 2\} \quad (1 \leq i \leq N)$$

$$d_i(t) = \sum_{j=1}^N \sum_{k=1}^m a_{ji}^k x_j(t-k)$$

**Driving function** for  $x_i$

$m$ : memory of the system

$a_{ji}^k$ : weight for variable  $x_j$  at time  $t - k$



# Our gene model

$$y_i(t+1) = \begin{cases} 2 & \text{if } d_i(t) \geq th_{i2} \\ 1 & \text{if } th_{i1} \leq d_i(t) \leq th_{i2} \\ 0 & \text{if } d_i(t) < th_{i1} \end{cases}$$

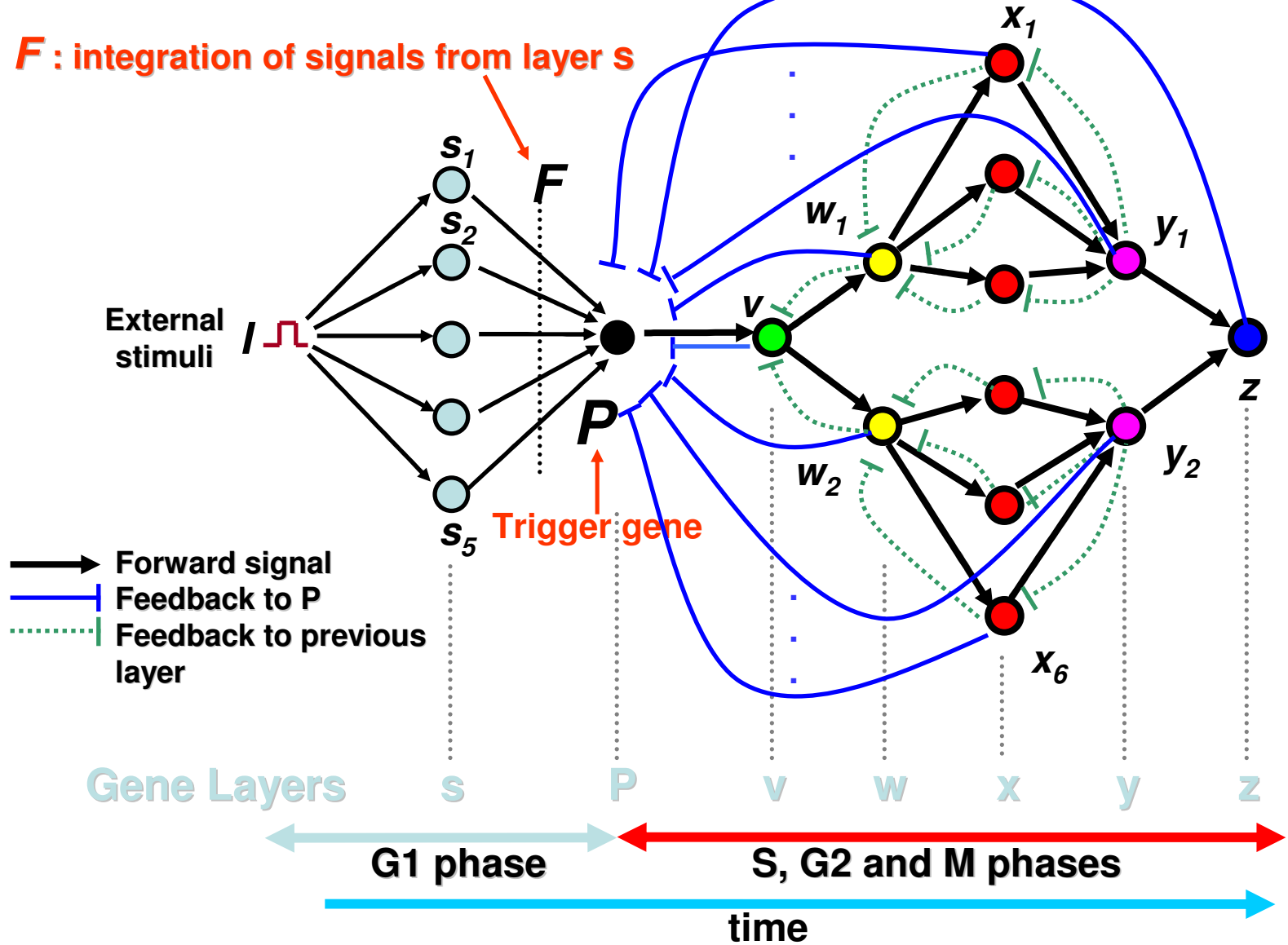
**Stochastic  
Transition  
Function**

$$x_i(t) = \begin{cases} y_i(t) & \text{with } P \approx 1 \\ a & \text{with } (1 - P)/2 \\ b & \text{with } (1 - P)/2. \end{cases}$$

$$a, b \in \{0, 1, 2\} - \{y_i(t)\}$$

# Our network architecture

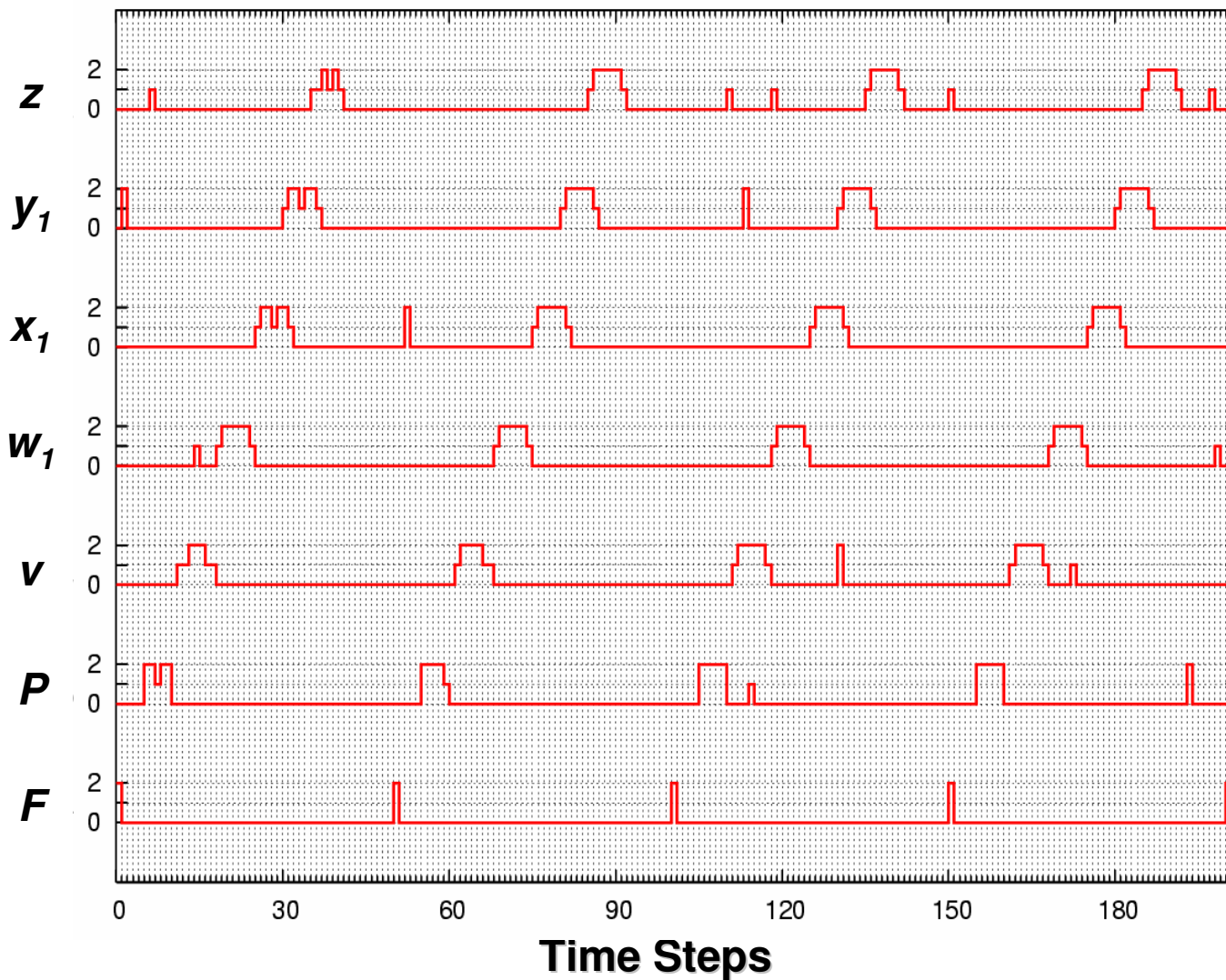
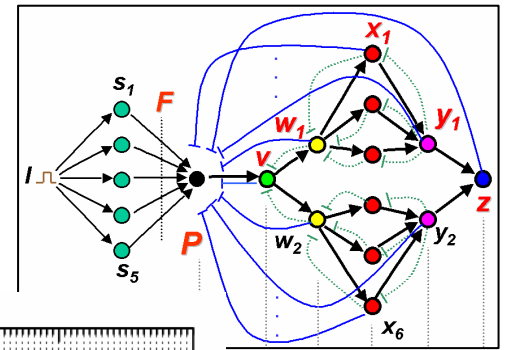
**F** : integration of signals from layer **s**





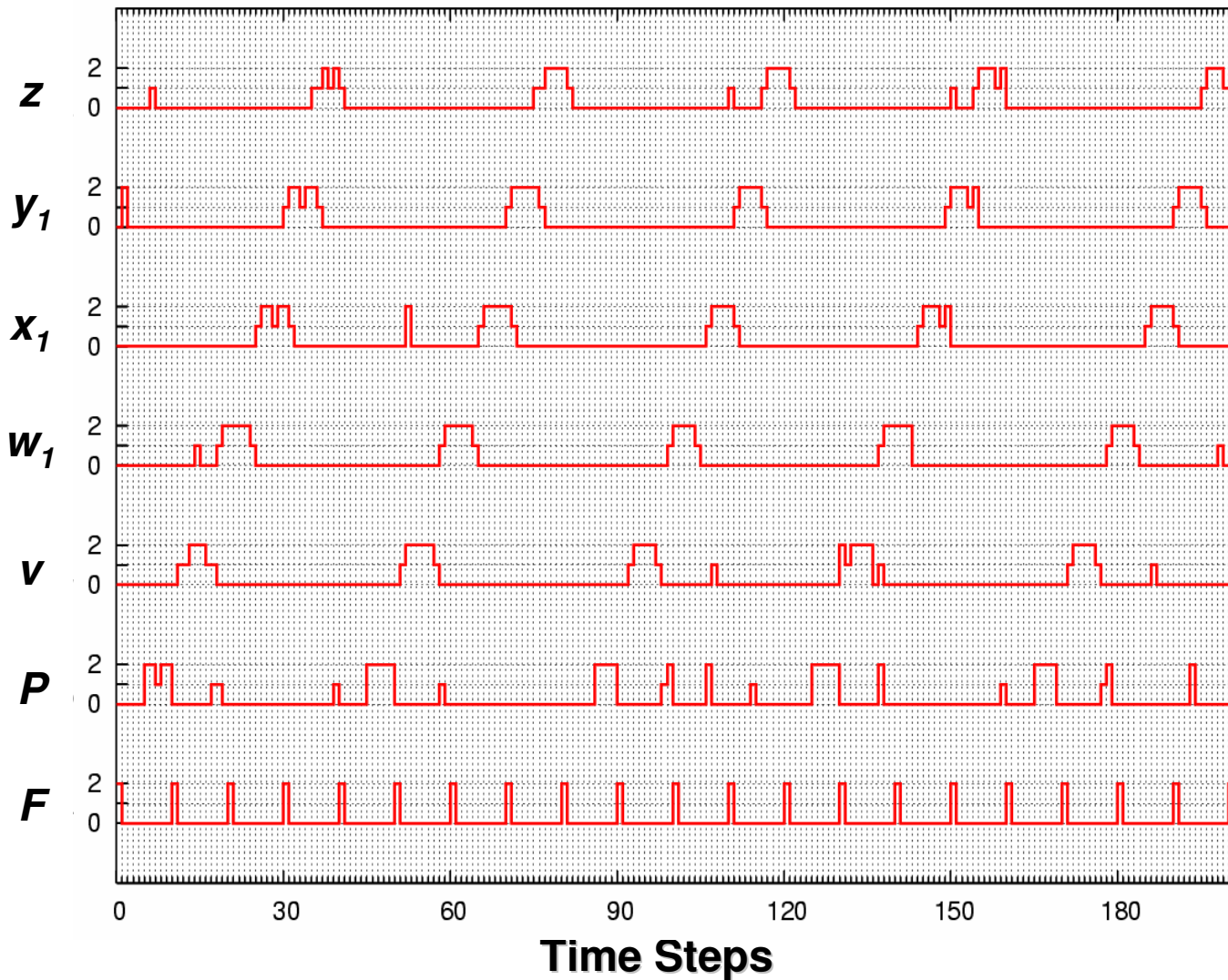
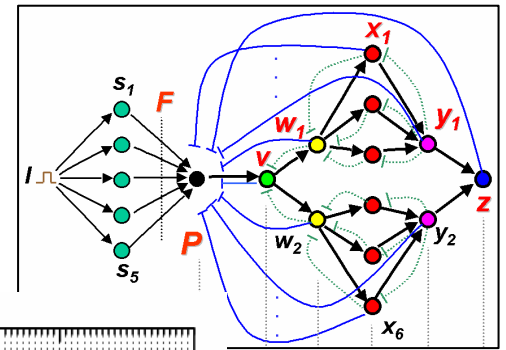
# PGN with $P = 0.99$

Signal  $F$  = period 50 oscillator



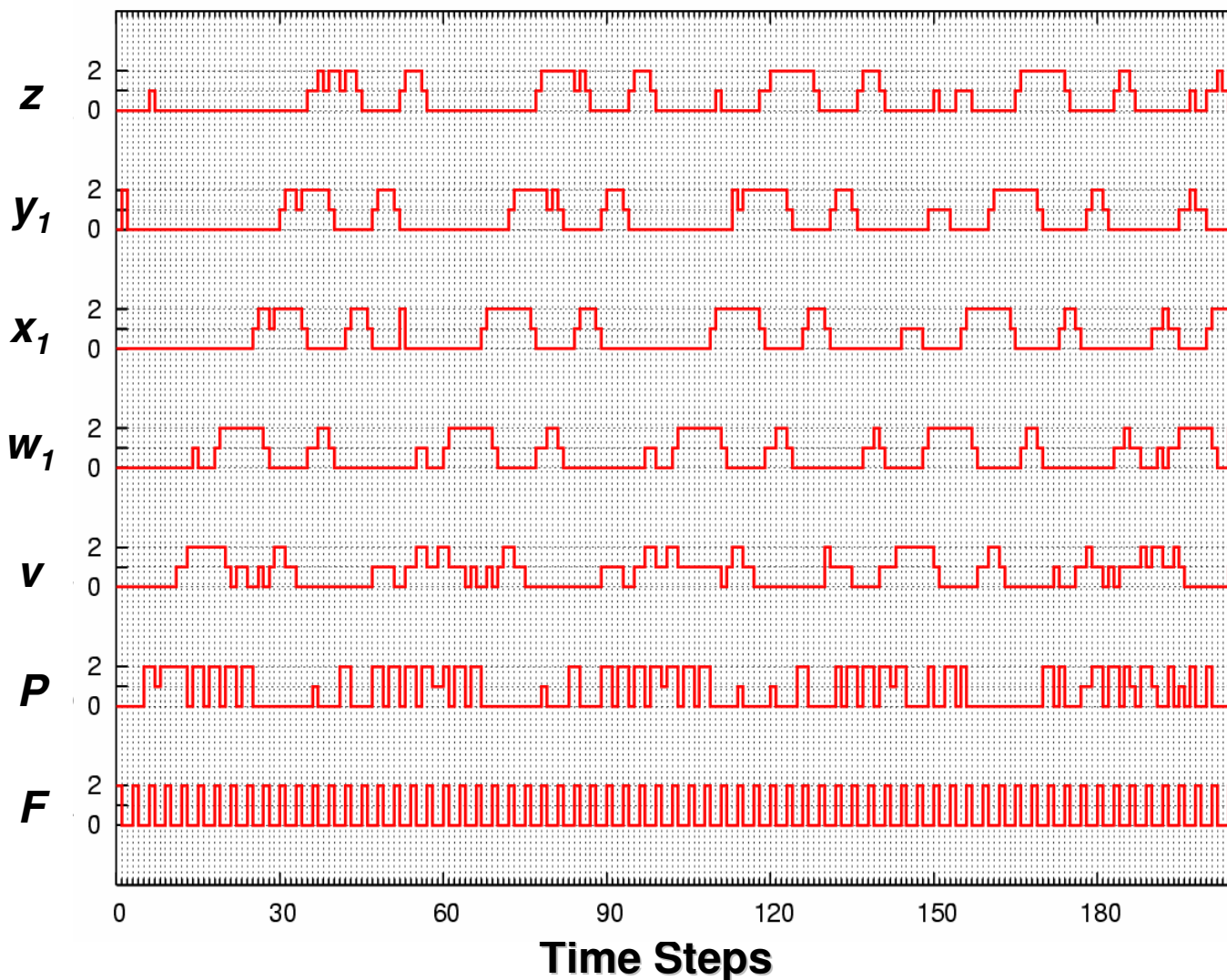
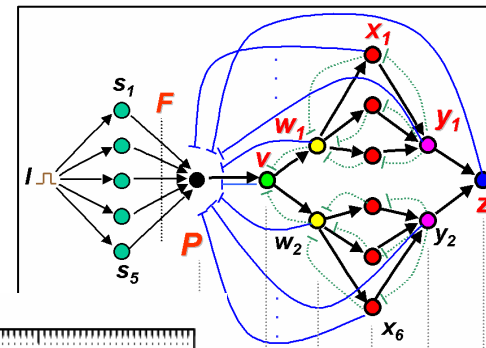
# PGN with $P = 0.99$

Signal  $F =$  period 10 oscillator



# PGN with $P = 0.99$

Signal  $F =$  period 3 oscillator



# Layout

- Introduction
- Probabilistic Genetic Networks (PGN)
- Estimation of PGNs
- Architecture estimation
- Malaria
- Cell Cycle
- Future works

- **Estimating cell cycle network from data**  
Use the developed model as constraint
- **Add new constraints**  
Representing phenomena like stability, robustness, protein interactions
- **Create dynamical criteria for network estimation**  
A kind of measure on the time sequence distribution
- **Design dedicated search algorithms for particular partitions spaces**  
Different partitions may have common parts that does not need to be calculated again.
- **Design smart search algorithms for feature selection**  
These algorithms should learn while walking through the search space