

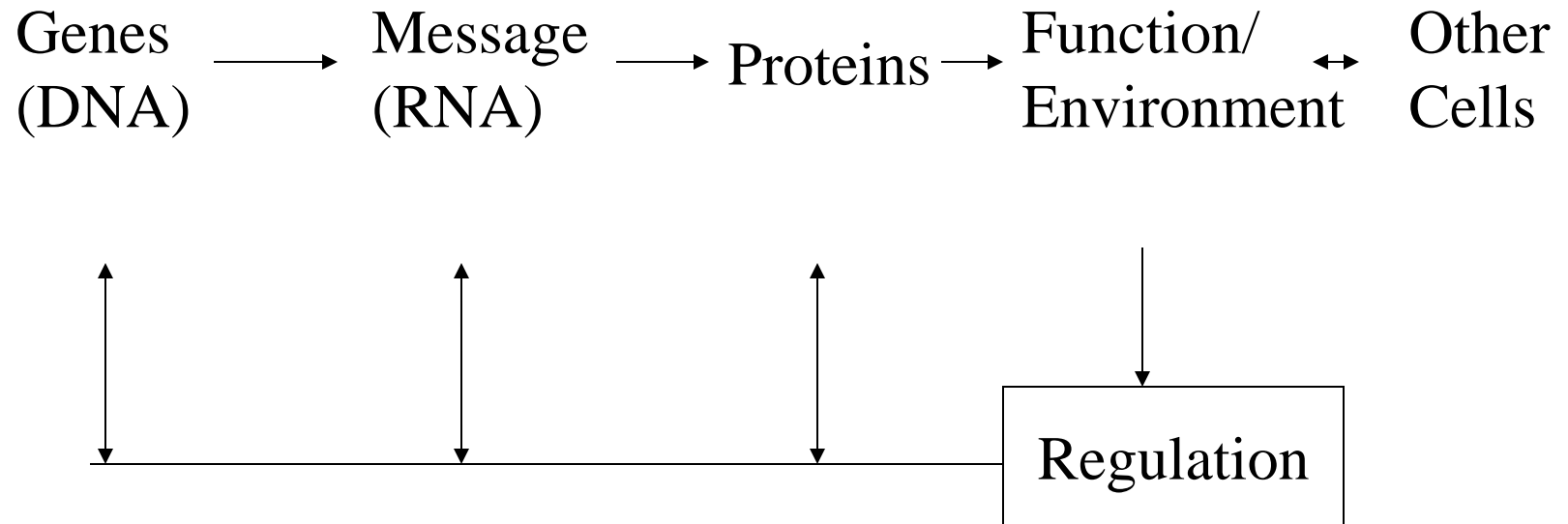
CISC 636 Computational Biology & Bioinformatics (Fall 2016)

Genetic networks and gene
expression data

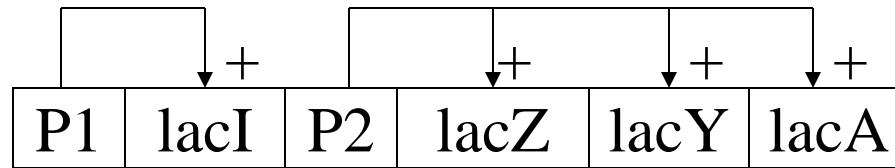
Gene Networks

- **Definition:** A gene network is a set of molecular components, such as genes and proteins, and interactions between them that collectively carry out some cellular function. A genetic regulatory network refers to the network of controls that turn on/off gene transcription.
- **Motivation:** Using a known structure of such networks, it is sometimes possible to describe behavior of cellular processes, reveal their function and the role of specific genes and proteins
- **Experiments**
 - DNA microarray : observe the expression of many genes simultaneously and monitor gene expression at the level of mRNA abundance.
 - Protein chips: the rapid identification of proteins and their abundance is becoming possible through methods such as 2D polyacrylamide gel electrophoresis.
 - 2-hybrid systems: identify protein-protein interactions
 - (Stan Fields' lab <http://depts.washington.edu/sfields/>)

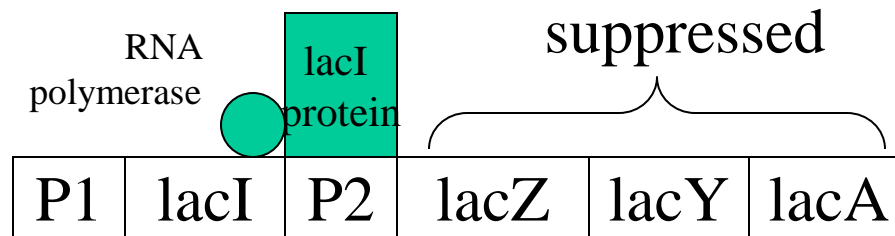
Regulation



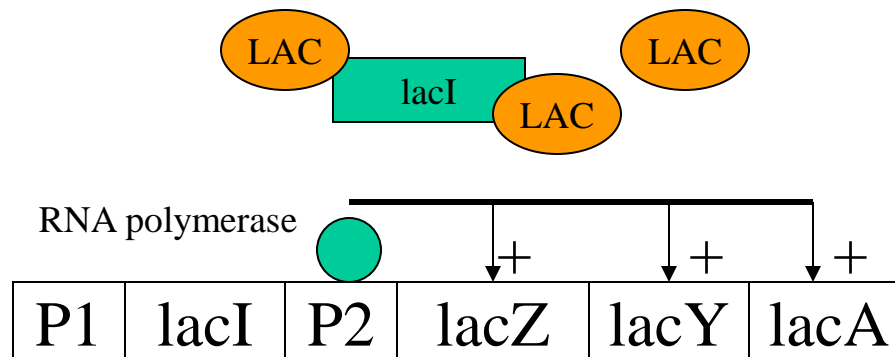
Operon



lac operon on *E. coli*



Repressor protein coded by lacI, bind to P2 preventing transcription of lacZ, lacY and lacA



Lactose binds with lacI, allowing RNA polymerase to bind to P2 and transcribe the structural genes

Genetic Network Models

- **Linear Model:** expression level of a node in a network depends on linear combination of the expression levels of its neighbors.
- **Boolean Model:** The most promising technique to date is based on the view of gene systems as a logical network of nodes that influence each other's expression levels. It assumes only two distinct levels of expression: 0 and 1. According to this model a value of a node at the next step is boolean function of the values of its neighbors.
- **Bayesian Model:** attempts to give a more accurate model of network behavior, based on Bayesian probabilities for expression levels.

Evaluation of Models

- Inferential power
- Predictive power
- Robustness
- Consistency
- Stability
- Computational cost

Boolean Networks: An example

X0	X1	X2	X3	
1	1	1	0	P0
-	1	0	1	P1
1	-	0	0	P2
1	1	-	1	P3
1	1	1	+	P4

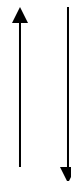
1: induced

0: suppressed

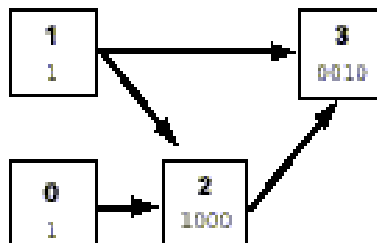
-: forced low

+: forced high

Interpreting data



Reverse Engineering



A A directed graph structure with numbered nodes connected by edges

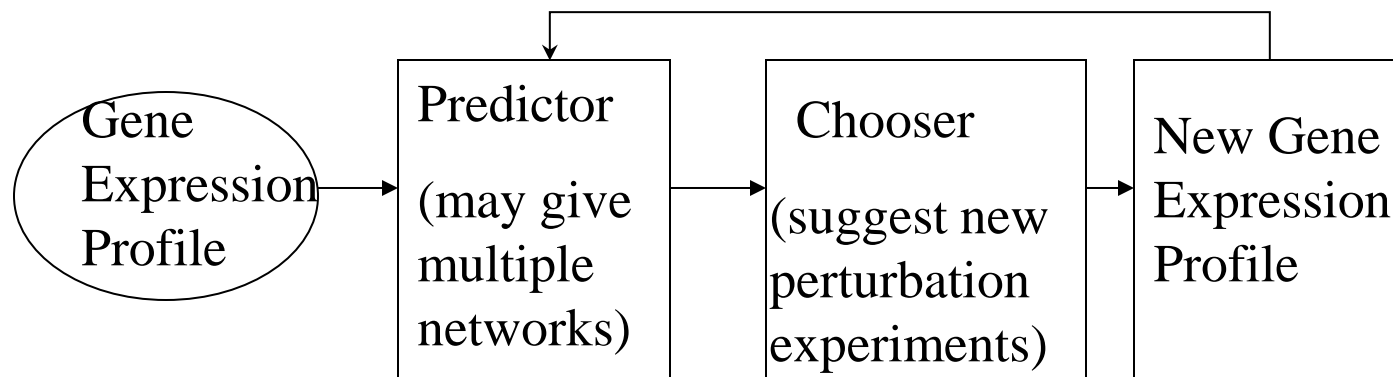
x1	1	0	1	0
x2	1	1	0	0
x3	0	0	1	0

B The truth table (shown for node 3 only)

```
x0 := 1
x1 := 1
x2 := x0 and x1
x3 := x1 and not x2
```

C The logic equations for each node

Boolean networks: A Predictor/Chooser scheme



Predictor

- A population of cells containing a target genetic network T is monitored in the steady state over a series of M experimental perturbations.
- In each perturbation p_m ($0 \leq m < M$) any number of nodes may be forced to a low or high level.

	x_0	x_1	x_2	x_3		
$E =$	1	1	1	0	p_0	← Wild-type state
	-	1	0	1	p_1	
	1	-	0	0	p_2	:- forced low
	1	1	-	1	p_3	
	1	1	1	+	p_4	+: forced high

Figure 2: Example expression matrix generated from the genetic network in fig. 1.

Step 1. For each gene x_n , find all pairs of rows (i, j) in E in which the expression level of x_n differs, excluding rows in which x_n was forced to a high or low value.

$$E = \begin{array}{c|cccc|c} & x_0 & x_1 & x_2 & x_3 & \\ \hline & 1 & 1 & 1 & 0 & p_0 \\ & - & 1 & 0 & 1 & p_1 \\ & 1 & - & 0 & 0 & p_2 \\ & 1 & 1 & - & 1 & p_3 \\ & 1 & 1 & 1 & + & p_4 \end{array}$$

Figure 2: Example expression matrix generated from the genetic network in fig. 1.

For x_3 , we find:

(p0, p1),

(p0, p3),

(p1, p2),

(p2, p3)

Step 2. For each pair (i,j), S_{ij} contains all other genes whose expression levels also differ between experiments i and j. Find the *minimum cover set* S_{min} , which contains at least one node from each set S_{ij}

$$E = \begin{array}{c|cccc|c} & x_0 & x_1 & x_2 & x_3 & \\ \hline & 1 & 1 & 1 & 0 & p_0 \\ & - & 1 & 0 & 1 & p_1 \\ & 1 & - & 0 & 0 & p_2 \\ & 1 & 1 & - & 1 & p_3 \\ & 1 & 1 & 1 & + & p_4 \end{array}$$

Figure 2: Example expression matrix generated from the genetic network in fig. 1.

Step 1:

(p0,p1),
(p0, p3),
(p1,p2),
(p2,p3)



Step 2:

(p0, p1)-> $S_{01}=\{x_0, x_2\}$
(p0, p3)-> $S_{03}=\{x_2\}$
(p1, p2)-> $S_{12}=\{x_0, x_1\}$
(p2, p3)-> $S_{23}=\{x_1\}$

So, now the S_{min} is $\{x_1, x_2\}$

Step 3. use the nodes in S_{\min} as input, x_n as output, build truth table to find out f_n (In this example, $n=3$)

Now the S_{\min} is $\{x_1, x_2\}$

$$E = \begin{array}{c|cccc|c} & x_0 & x_1 & x_2 & x_3 & \\ \hline & 1 & 1 & 1 & 0 & p_0 \\ & - & 1 & 0 & 1 & p_1 \\ & 1 & - & 0 & 0 & p_2 \\ & 1 & 1 & - & 1 & p_3 \\ & 1 & 1 & 1 & + & p_4 \end{array}$$

x_1 1 0 1 0

x_2 1 1 0 0

x_3 0 * 1 0

So $f_3 = 0 * 1 0$

* cannot be determined

Figure 2: Example expression matrix generated from the genetic network in fig. 1.

Chooser

For L hypothetical equiprobable networks generated by the predictor, choose perturbation p that would best discriminate between the L networks, by maximizing entropy H_p as defined below.

$$H_p = - \sum_{s=1}^S (l_s/L) \log_2 (l_s/L)$$

where l_s is the number of networks giving the state s

Note: ($1 \leq s \leq S$), and ($1 \leq S \leq L$)

Result and Evaluation

- Evaluation of Predictor
- construct a target network T : size = N , and maximum in-degree = k (where the in-degree of a node is its number of incoming edges)
- *sensitivity* is defined as the percentage of edges in the target network that were also present in the inferred network, and *specificity* is defined as the percentage of edges in the inferred network that were also present in the target network.

<i>A</i>	<i>B</i>	<i>C</i>	<i>D</i>	<i>E</i>	<i>F</i>	<i>G</i>	<i>H</i>	<i>I</i>	<i>J</i>
N	k	Total Sim. Edges	Num. Inferred Networks	Total Inferred Edges	Num. Shared Edges	Sensitivity	Specificity	Num. Nodes w/ 1 Soln.	CPU Time (sec)
5	2	4 (0.1)	1 (.02)	3 (0.1)	3 (0.1)	77%	99%	5 (0.0)	0.1 (0.0)
10	2	12 (0.1)	60 (50)	9 (0.1)	9 (0.1)	71%	95%	9 (0.1)	0.1 (0.0)
20	2	27 (0.2)	3×10^7 (10^7)	21 (0.2)	19 (0.1)	71%	92%	18 (0.1)	0.2 (0.0)
50	2	72 (0.2)	1×10^{12} (10^{12})	57 (0.3)	51 (0.3)	71%	90%	45 (0.2)	0.8 (0.0)
100	2	146 (0.7)	3×10^{26} (10^{26})	119 (0.9)	104 (0.7)	70%	88%	89 (0.5)	6.6 (0.3)
20	4	44 (0.3)	2×10^6 (10^6)	28 (0.3)	23 (0.2)	51%	84%	16 (0.1)	0.2 (0.0)
20	6	57 (0.5)	2×10^7 (10^7)	33 (0.3)	27 (0.2)	42%	82%	14 (0.2)	0.2 (0.0)
20	8	69 (0.7)	9×10^7 (10^8)	38 (0.4)	31 (0.3)	35%	82%	13 (0.2)	0.2 (0.0)

Discussions

- Incorporate pre-existing information
- Boolean to multi-levels
- Cyclic networks
- Noise tolerance

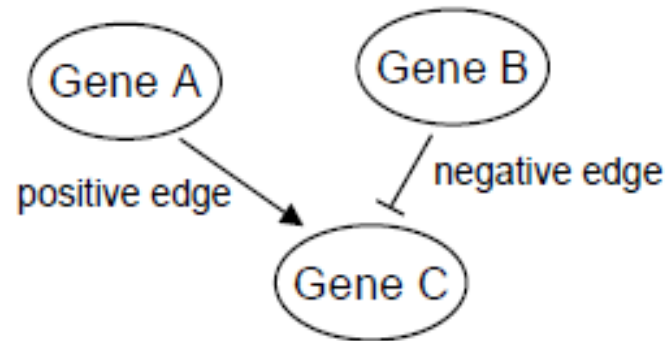
References

- Ideker, Thorsson, and Karp, PSB 2000, 5: 302-313.

Bayesian Networks

Biological processes are stochastic

- Data can be noisy as well.

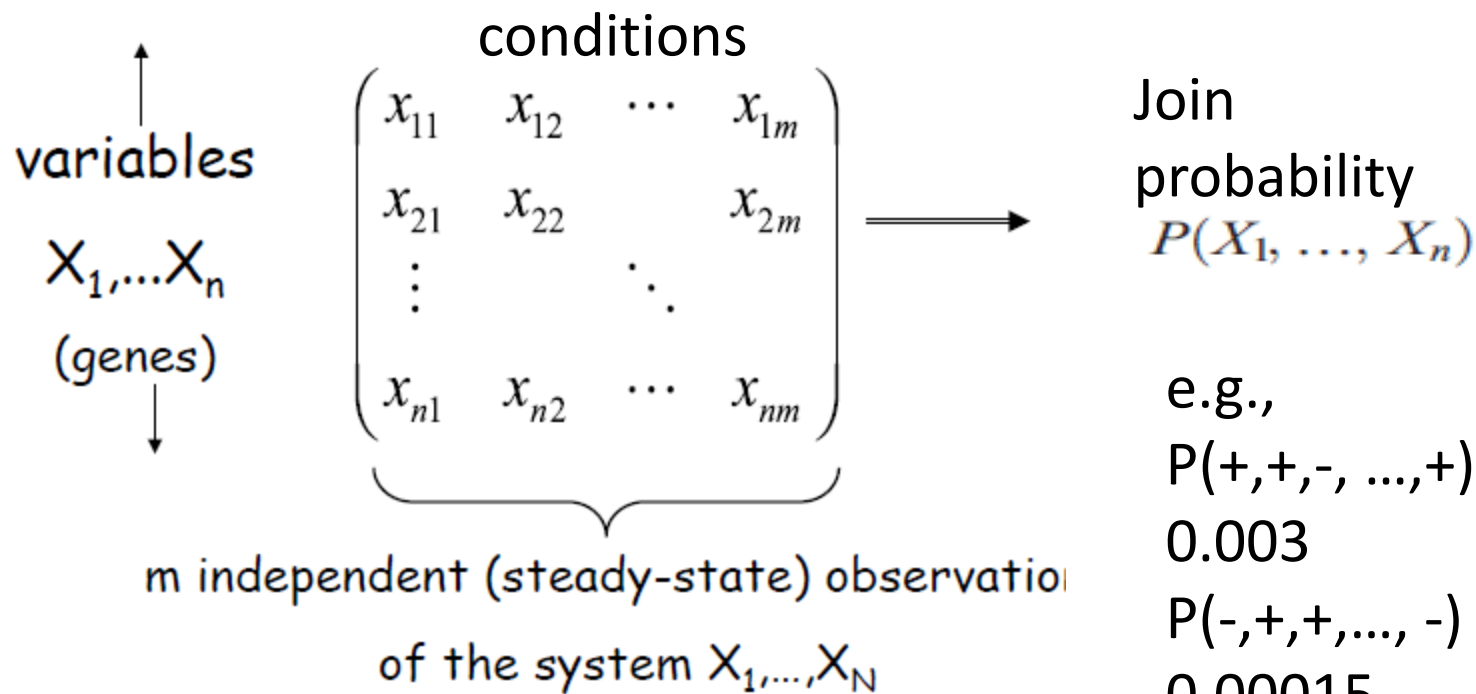


§ Quantitative part:

Gene A	Gene B	$P(C+ AB)$	$P(C- AB)$
+	+	0.6	0.4
-	+	0.01	0.99
+	-	0.99	0.01
-	-	0.4	0.6



This row indicates that when Gene A and Gene B are up-regulated, then Gene C has a 60% probability to be up-regulated and a 40% probability to be down-regulated.



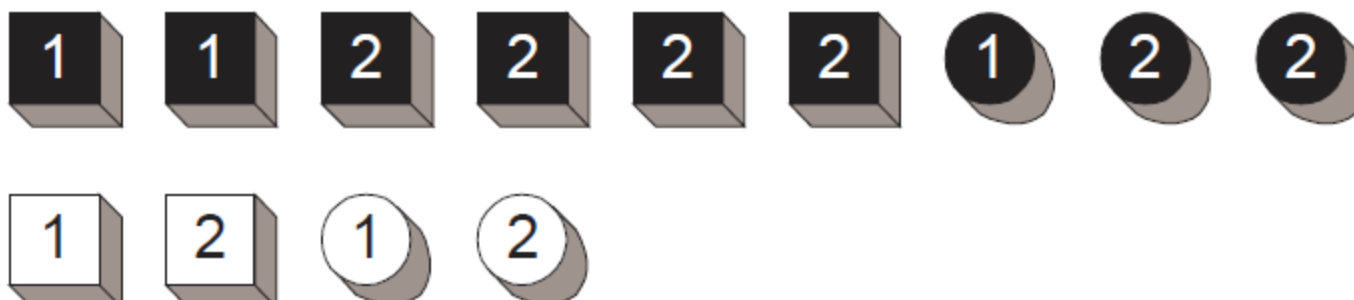
e.g.,
 $P(+, +, -, \dots, +) = 0.003$
 $P(-, +, +, \dots, -) = 0.00015$

\dots
 2^N

Query/Inference: $P(X_1 \mid X_6, X_7)$?

How many
 combinations?

Conditional Probability and Conditional Independence



$$P(\text{One}) = \frac{5}{13}$$

$$P(\text{One}|\text{Square}) = \frac{3}{8}$$

$$P(\text{One}|\text{Black}) = \frac{3}{9} = \frac{1}{3}$$

$$P(\text{One}|\text{Square} \cap \text{Black}) = \frac{2}{6} = \frac{1}{3}$$

$$P(\text{One}|\text{White}) = \frac{2}{4} = \frac{1}{2}$$

$$P(\text{One}|\text{Square} \cap \text{White}) = \frac{1}{2}$$

So One and Square are not independent, but they are conditionally independent given Black and given White.

Bayesian Network as an efficient way to factorize the Joint Probability

Factorization of joint probability

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | X_1, \dots, X_{i-1})$$

of parameters = $2^N - 1$

Conditional independence

$$P(X_1, \dots, X_n) = \prod_{i=1}^n P(X_i | Pa_i)$$

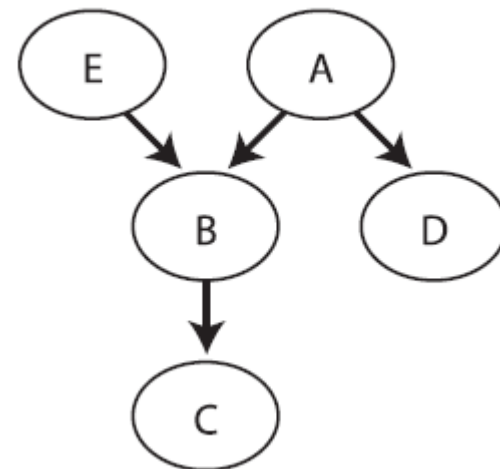
Assuming max in-degree k , the number of parameters is reduced to $2^k N$

Example:

$$P(A, E, B, C, D) = P(A)P(E|A)P(B|A, E)P(C|A, E, B)P(D|A, E, B, C)$$

of parameters = $1 + 2 + 4 + 8 + 16 = 31$

$$P(A, E, B, C, D) = P(A)P(E)P(B|A, E)P(C|B)P(D|A)$$



of parameters = $1 + 1 + 4 + 1 + 1 = 10$

A greedy Algorithm to Learn Bayesian Network from the data

Input

D // a data set

G_o // initial network structure

Output

G // final network structure

Greedy-structure-search

$G_{\text{best}} = G_o$

repeat // apply best possible operator to G in each iteration

$G = G_{\text{best}}$

foreach operator o // (each edge addition, deletion, or reversal on G)

$G^o = o(G)$ // apply to G

if G^o is cyclic **continue**

if $\text{scoreBDe}(G^o : D) > \text{scoreBDe}(G_{\text{best}} : D)$

$G_{\text{best}} = G^o$

until $G == G_{\text{best}}$ // no change in structure improves score

Parameter Estimation

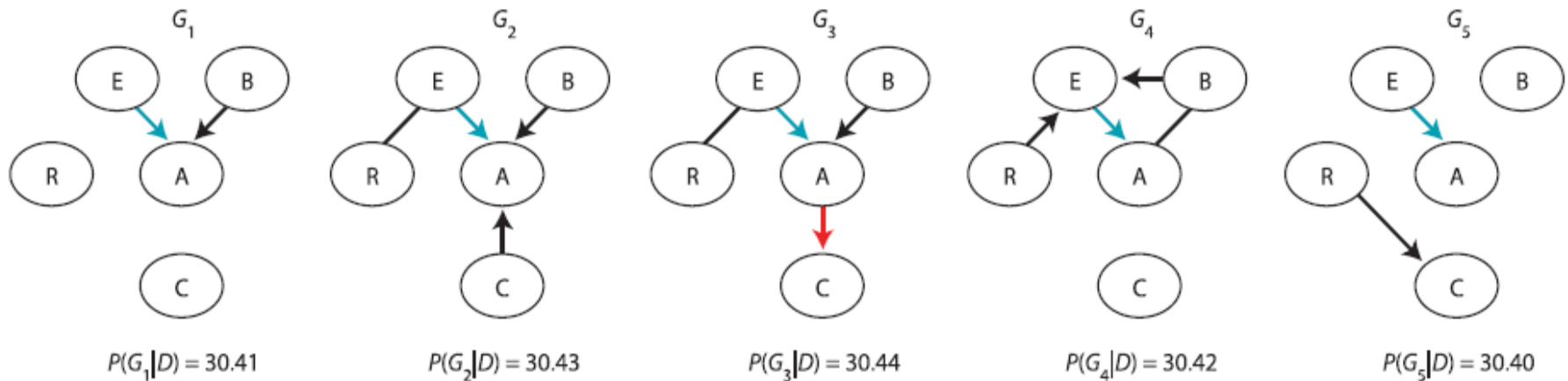
- Maximum Likelihood
- Bayesian approach
 - Dirichlet priors are used for model parameters.

Structure evaluation

$$\begin{aligned}\text{BayesianScore}(M) &= \log[P(M \mid D)] \\ &= \log[P(M)] + \log[P(D \mid M)] + c\end{aligned}$$

- Where M = model, D = microarray data, c = constant

Model Averaging



$$P[f(G)|D] = \sum_G f(G)P(G|D)$$

Feature f : edge $X \rightarrow Y$ is in the network.

$f(G) = 1$, if G has the feature
 $= 0$, otherwise.

How to compute $P[f(G)|D]$?

- Enumerate all high scored networks
- Sampling (MCMC)
- Bootstrap

Bootstrap

- For $i = 1, \dots, m$, construct a data set D_i by sampling, with replacement, M instances from D . Then, apply the learning procedure on D_i to induce a network structure G_i .
- For each feature f of interest, calculate

$$\text{conf}(f) = \frac{1}{m} \sum_{i=1}^m f(G_i)$$

Research article

Open Access

From gene expression to gene regulatory networks in *Arabidopsis thaliana*

Chris J Needham^{*1}, Iain W Manfield², Andrew J Bulpitt¹,
Philip M Gilmartin^{2,4} and David R Westhead³

Address: ¹School of Computing, University of Leeds, Leeds, LS2 9JT, UK, ²Institute of Integrative and Comparative Biology, University of Leeds, Leeds, LS2 9JT, UK, ³Institute of Molecular and Cellular Biology, University of Leeds, Leeds, LS2 9JT, UK and ⁴Current address : School of Biological and Biomedical Sciences, Durham University, Durham, UK

Email: Chris J Needham^{*} - C.Needham@leeds.ac.uk; Iain W Manfield - I.Manfield@leeds.ac.uk; Andrew J Bulpitt - A.J.Bulpitt@leeds.ac.uk; Philip M Gilmartin - Philip.Gilmartin@durham.ac.uk; David R Westhead - D.R.Westhead@leeds.ac.uk

^{*} Corresponding author

Published: 3 September 2009

Received: 7 April 2009

BMC Systems Biology 2009, 3:85 doi:10.1186/1752-0509-3-85

Accepted: 3 September 2009

This article is available from: <http://www.biomedcentral.com/1752-0509/3/85>

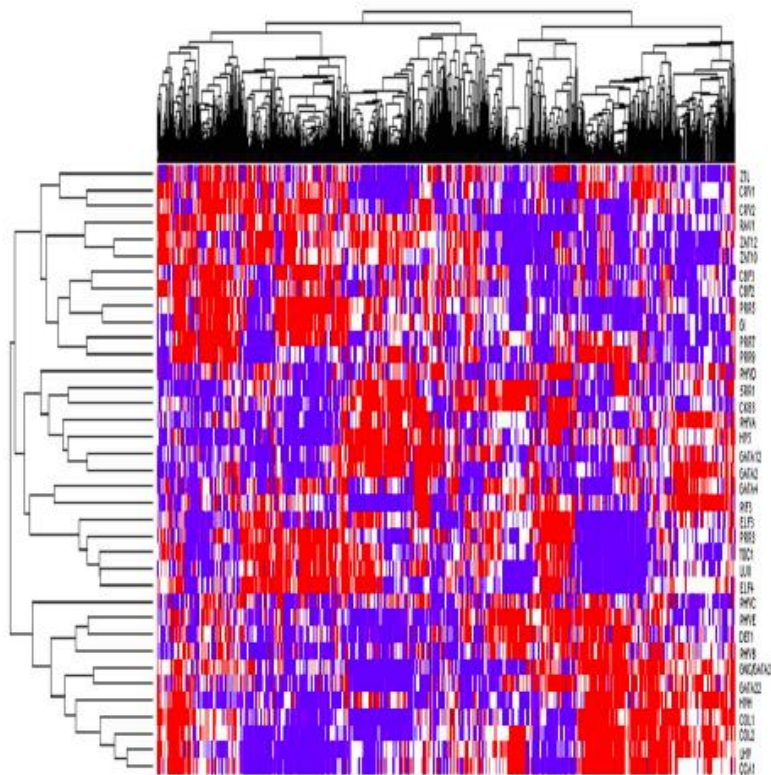


Figure 4
Clustergram of quantized gene expression profiles for 37 genes of interest, over 2904 microarrays. Both genes and experiments have been clustered. The three classes representing the low, medium and high classes are coloured blue, white and red respectively.

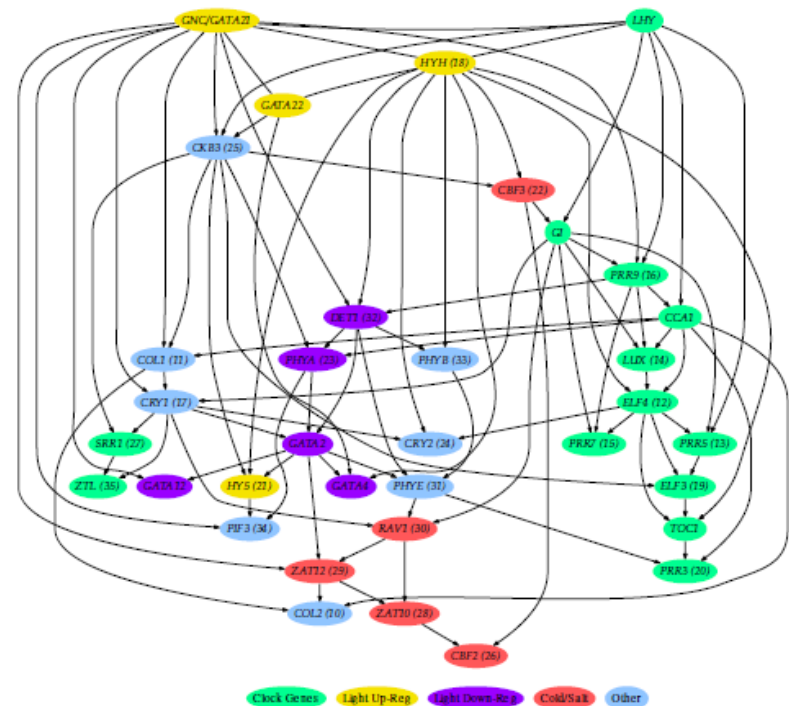


Figure 5
Learned regulatory network for other networks and poorly-characterized genes. The learned network structure starting from a set of nine genes (four clock and five GATA genes of interest), with additional genes added to the network from a selection of 37 genes. The number in parentheses next to the gene name denotes the order it was added to the network. Most of these genes were added to the network in early iterations, however, genes such as *SRR1* and *ZTL* with *bona fide* roles in the clock were added late and only indirectly linked to other clock components. All these interactions are very similar throughout the later iterations, once most of these components have been added to the network.