

# Dynamic Algorithm for Inferring Qualitative Models of Gene Regulatory Networks

Zheng Yun and Kwoh Chee Keong

Bioinformatics Research Center  
School of Computer Engineering  
Nanyang Technological University  
Singapore 639798

Tel: +65-67906613, +65-67906057

Fax: +65-63162780, +65-67906559

Email: [pg04325488@ntu.edu.sg](mailto:pg04325488@ntu.edu.sg), [asckkwoh@ntu.edu.sg](mailto:asckkwoh@ntu.edu.sg)

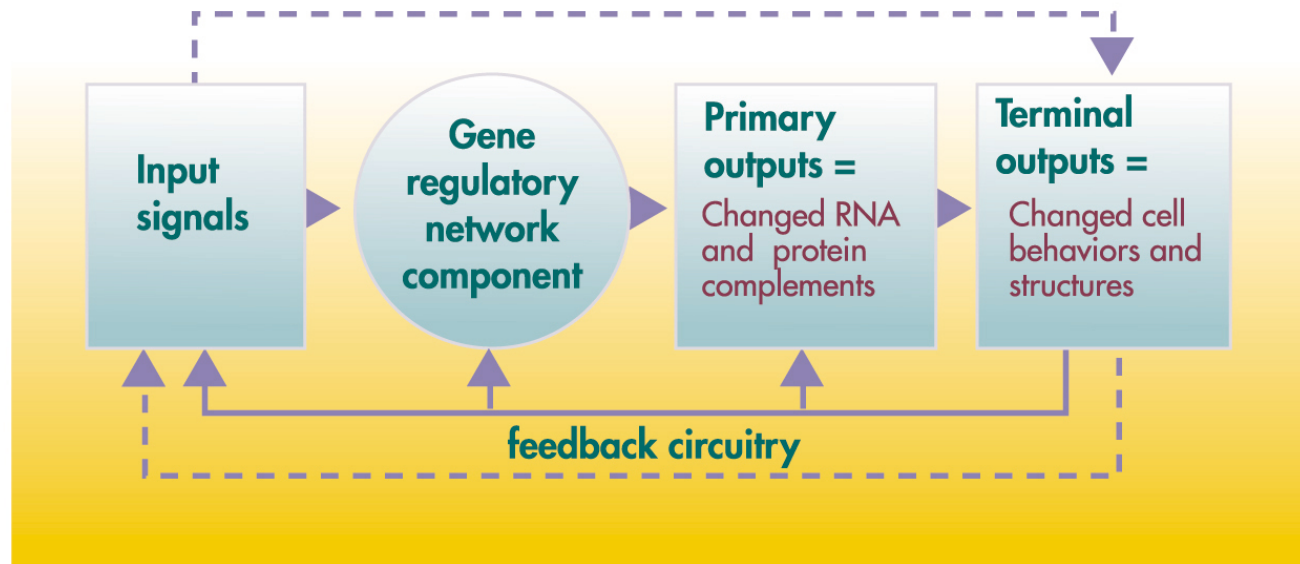


# Topics

---

- Gene Regulatory Networks (GRNs)
- Reverse engineering methods for reconstructing GRNs
- Information theory foundation for our method
- The DFL algorithm
- Experimental results
- Conclusion
- Summary

# GRNs macro-view

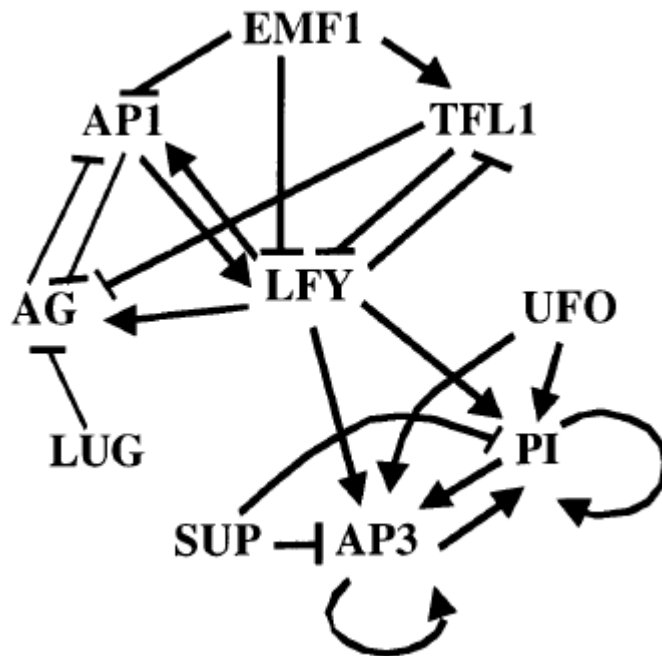


YGG 01-0086

Gene regulatory networks, Courtesy of Genomes to Life Program of U.S. Department of Energy, <http://www.doegenomestolife.org>.

Input signals are from both intra-cellular and inter-cellular sources. The upper dashed arrow is the signaling responses, which may act directly on cell behaviors and structures. The solid and the dashed arrows at the bottom are direct and indirect feedbacks respectively.

# GRN for *Arabidopsis thaliana* flower morphogenesis



Generalized logical  
model of flower  
morphogenesis,  
Courtesy of Mendoza  
*et al.* 1999



# Experimental data supporting the GRN of *Arabidopsis thaliana* flower morphogenesis

Interactions	Main evidence	Main references
<i>AG</i> —  <i>API</i>	<i>API</i> mRNA accumulates uniformly in <i>ag-1</i> mutant flowers	Gustafson-Brown <i>et al.</i> (1994)
<i>API</i> —  <i>AG</i>	Sepals are replaced by carpels, and petals by stamens in <i>ap1</i> mutants. <i>AG</i> mRNA found in all the flower primordium of <i>ap1-1</i> plants	Bowman <i>et al.</i> (1993) Weigel and Meyerowitz (1993)
<i>API</i> → <i>LFY</i>	Reduction of <i>LFY</i> mRNA in <i>ap1 cal</i> double mutants (independent pathways)	Weigel and Nilsson (1995) Bowman <i>et al.</i> (1993) Kempin <i>et al.</i> (1995)
<i>AP3/PI</i> → <i>AP3/PI</i>	<i>AP3</i> and <i>PI</i> mRNA levels are not maintained in <i>ap3-3, pi-1</i> or double mutants. Co-immunoprecipitation of AP3 and PI proteins	Goto and Meyerowitz (1994) Jack <i>et al.</i> (1992)
<i>EMF1</i> —  <i>API, LFY</i>	Inferred by morphological evidence that <i>EMF1</i> inhibits the flowering promoting genes	Mendoza and Alvarez-Buylla (1998)
<i>EMF1</i> → <i>TFL1</i>	Inferred by morphological evidence that <i>EMF1</i> activates the late late-flowering genes	Mendoza and Alvarez-Buylla (1998)
<i>LFY</i> → <i>AG</i>	Early expression of <i>AG</i> is abnormally low in <i>lfy-6</i> flowers	Weigel and Meyerowitz (1993)
<i>LFY</i> → <i>API</i>	<i>API</i> mRNA delayed in <i>lfy</i> mutants. Earlier <i>API</i> promoter induction in plants overexpressing <i>LFY</i>	Weigel and Nilsson (1995) Parcy <i>et al.</i> (1998)
<i>LFY</i> → <i>AP3</i>	Amount and domain of <i>AP3</i> expression reduced in <i>lfy-6</i> mutants	Weigel and Meyerowitz (1993)
<i>LFY</i> → <i>PI</i>	Amount and domain of <i>PI</i> expression reduced in <i>lfy-6</i> mutants	Weigel and Meyerowitz (1993)
<i>LFY</i> —  <i>TFL1</i>	Plants overexpressing <i>LFY</i> are very similar to <i>tfl1</i> mutants	Weigel and Nilsson (1995)
<i>LUG</i> —  <i>AG</i>	Ectopic expression of <i>AG</i> in <i>lug-1</i> mutants	Liu and Meyerowitz (1995)
<i>SUP</i> —  <i>AP3</i>	Ectopic expression of <i>AP3</i> in <i>sup-1</i> mutants	Sakai <i>et al.</i> (1995)
<i>SUP</i> —  <i>PI</i>	Contrary to wild type, <i>PI</i> expression is not reduced in the center of <i>sup-1</i> flowers	Goto and Meyerowitz (1994)
<i>TFL1</i> —  <i>AG</i>	Inferred from morphological evidence. Double mutants <i>ap1-1 ap2-2</i> have a disrupted C activity, which is rescued with the addition of <i>tfl1</i> mutation	Mendoza and Alvarez-Buylla (1998)
<i>TFL1</i> —  <i>LFY</i>	Precocious appearance of floral buds expressing <i>LFY</i> in <i>tfl1-2</i> plants	Weigel <i>et al.</i> (1992)
<i>UFO</i> → <i>AP3</i>	<i>AP3</i> protein and messenger levels reduced in <i>ufo-2</i> plants	Levin and Meyerowitz (1995)
<i>UFO</i> → <i>PI</i>	<i>PI</i> mRNA reduced in early stages of flower development in <i>ufo-2</i> plants	Levin and Meyerowitz (1995)

Table Courtesy of Mendoza *et al.* 1999



# Reverse engineering?

---

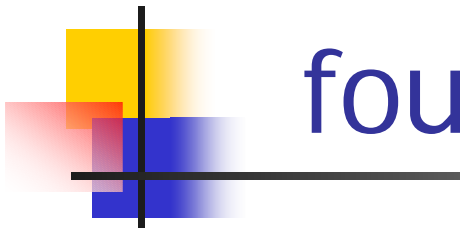
- By mapping the output of each gene to the inputs of other genes, it is possible to reverse engineer developmental circuits and even whole networks.

Meredith L. Howard and Eric H. Davidson. *Development* 271:109–118.  
2004

- If the expression of gene  $A$  is regulated by proteins  $B$  and  $C$ , then  $A$ 's expression level is a function of the joint activity levels of  $B$  and  $C$ .

Nir Friedman. *SCIENCE* 303:799-805.2004

# Information theory foundation of our approach



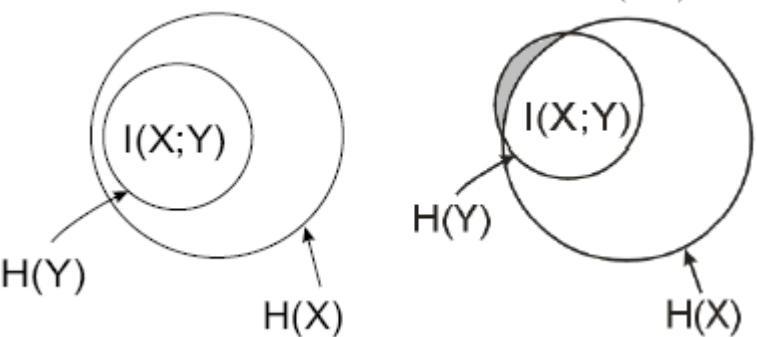
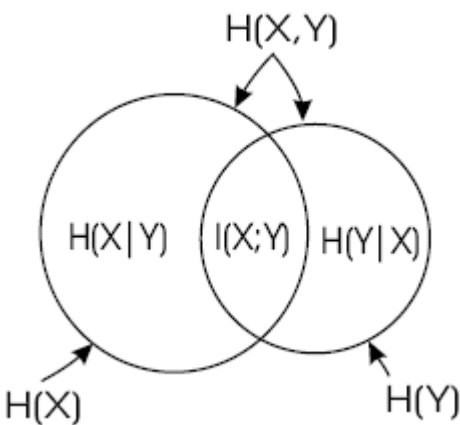
$$H(X) = - \sum_x P(X = x) \log P(X = x)$$

$$H(Y|\mathbf{X}) = - \sum_{\mathbf{x}} \sum_y p(\mathbf{x}, y) \log p(y|\mathbf{x})$$

$$I(\mathbf{X}; Y) = H(Y) - H(Y|\mathbf{X}) = H(\mathbf{X}) - H(\mathbf{X}|Y)$$

**Theorem 2.1** *If the mutual information between  $\mathbf{X}$  and  $Y$  is equal to the entropy of  $Y$ , i.e.,  $I(\mathbf{X}; Y) = H(Y)$ , then  $Y$  is a function of  $\mathbf{X}$ .*

**Definition 5.1** *If  $H(Y) - I(\mathbf{X}; Y) \leq \epsilon \times H(Y)$ , then  $Y = f_\epsilon(\mathbf{X})$  where  $\epsilon$  is a significant factor.*





# Models under consideration

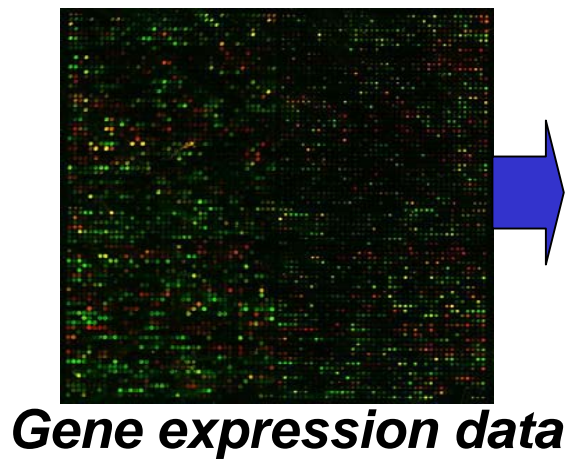
---

$$X_i(t + 1) = f_i(X_{i1}(t), \dots, X_{ik}(t)),$$

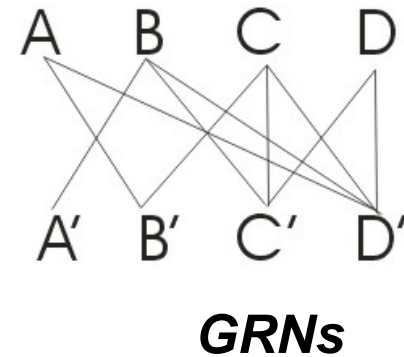
- Boolean networks: (Liang, Fuhrman. & Somogyi 1998), (Akutsu, Miyano & Kuhara 1999), (Wuensche 1998), etc.
- Generalized Logical Formalism (GLF): (Sanchez & Thieffy 2001), (Thomas & d'Ari 1990), (Thomas, Thieffy & Kaufman 1995), etc.
- Partial Linear Differential Equations (PLDE): (Mendoza, et al. 1999.), (Mestl et al. 1995), (de Jong et al. 2002), etc.



# Reverse engineering approach



A	B	C	D	→	A'	B'	C'	D'
0	0	0	0		0	0	0	0
0	0	0	1		0	0	0	0
0	0	1	0		0	1	0	0
0	0	1	1		0	1	1	1
0	1	0	0	→	1	0	0	0
0	1	0	1		1	0	1	0
0	1	1	0		1	1	1	0
0	1	1	1		1	1	1	1
1	0	0	1		0	1	0	0
1	0	1	0		0	1	0	0
1	0	1	1		0	1	1	1
1	1	0	0		1	1	0	1



**State-transition pairs**



# How much data

---

**Theorem 3.1 (Akutsu 1998)**  $\Omega(2^k + k \log_2 n)$  transition pairs are necessary in the worst case to identify the Boolean network of maximum indegree  $\leq k$ .

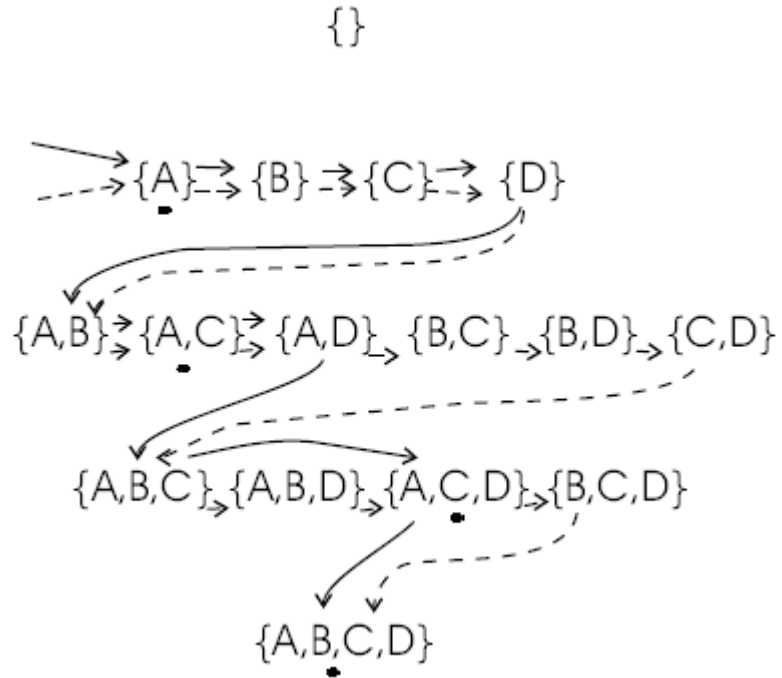
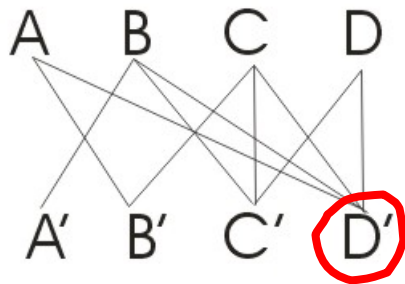
We generalize Theorem by Akutsu 1998 to meet the multilevel datasets.

**Theorem 3.2**  $\Omega(b^k + k \log_b n)$  transition pairs are necessary in the worst case to identify the qualitative GRN models of maximum indegree  $\leq k$  and the maximum number of discrete level for variables  $\leq b$ .

$$N = c \times (b^k + k \log_b n).$$



# The DFL algorithm



Search methods:

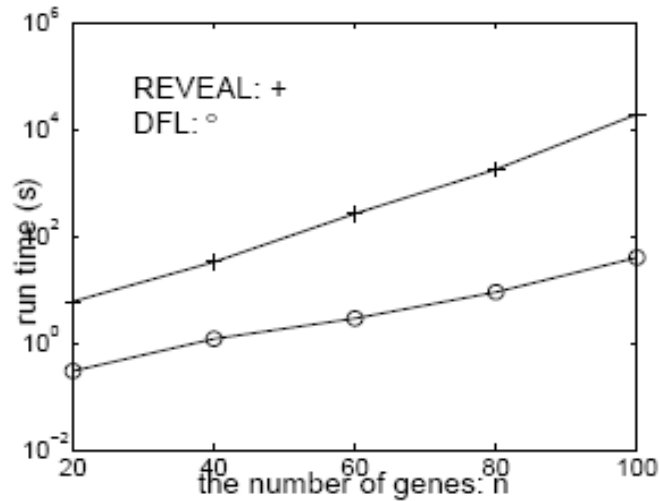
REVEAL(Liang *et al.*  
1998): dashed line

$$O((b^k + k \log_b n) n^{k+1})$$

DFL: solid line

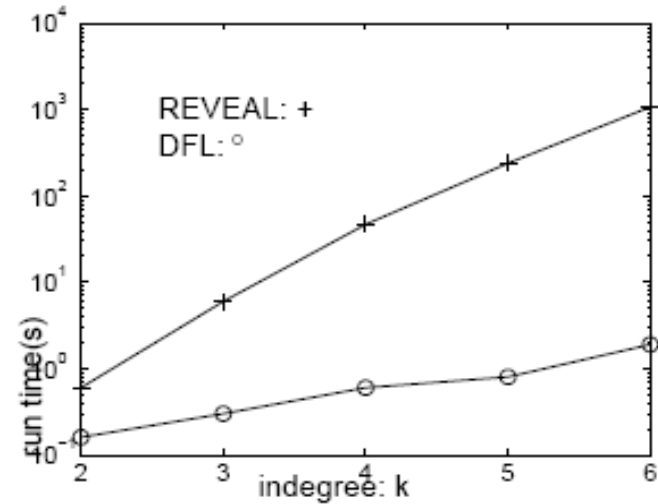
$$O((k b^k + k^2 \log_b n) n^2)$$

# Experiments for time complexity



(a)

$k = 3, c = 3$



(b)

$n = 20, c = 3$



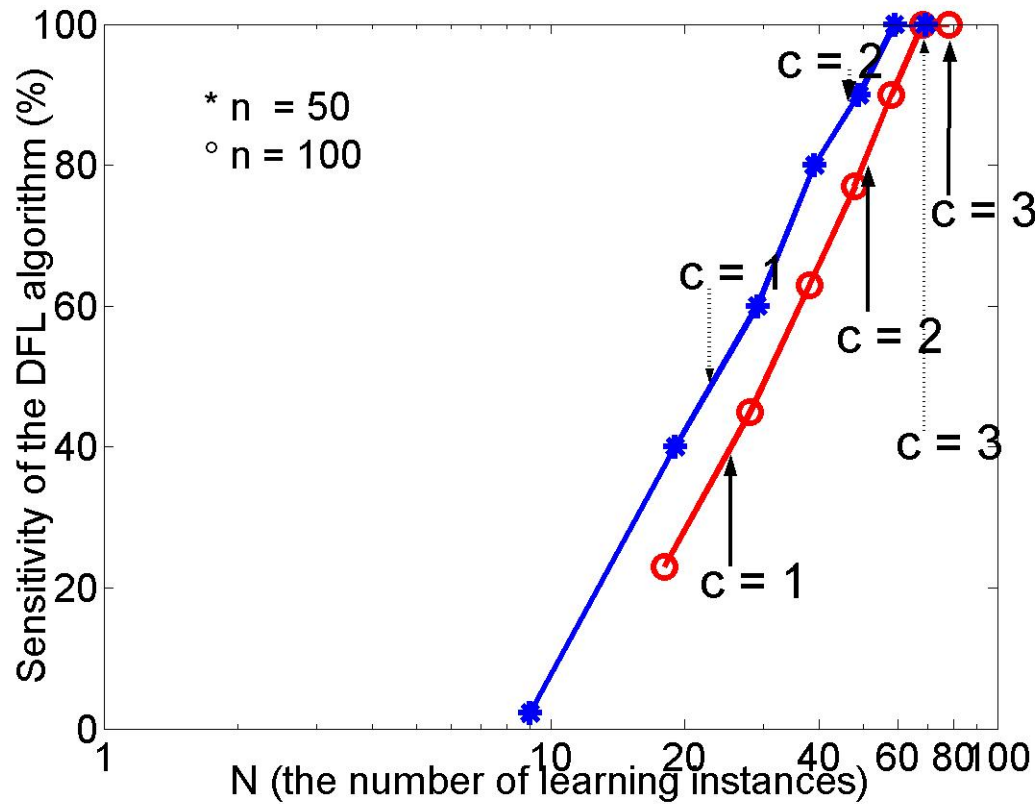
# Sensitivity

	predicted as positive	predicted as negative
positive	TP	FN
negative	FP	TN

$$\begin{aligned}\text{Sensitivity} &= \frac{\text{No. of correct positive predictions}}{\text{No. of positives}} \\ \text{wrt positives} & \\ &= \frac{TP}{TP + FN}\end{aligned}$$



# Experiments for sensitivity

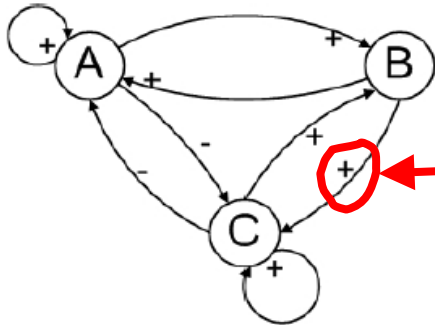


# Reconstruct GLF Model

abc	ABC
0 0 0	0 0 0
0 0 1	0 0 1
0 0 2	0 1 1
0 1 0	1 0 1
0 1 1	0 0 2
0 1 2	0 1 2

A simple example of GLF from (Thieffry and Thomas 1998).

↓ DFL

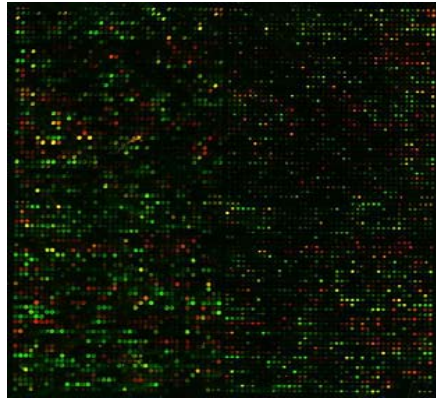


	$A'$	$B'$	$C'$
$A$	0.2	0.6	-0.7
$B$	0.3	0	0.3
$C$	-0.7	0.6	0.2

**Table 4. The correlation coefficient matrix of the GLF example in Figure 6.**



# Experiments on yeast cycle-cycle gene expression profiles



Cell-cycle expression profiles, from Cho et al. 1998, cover approximately two full cell cycles.

DFL

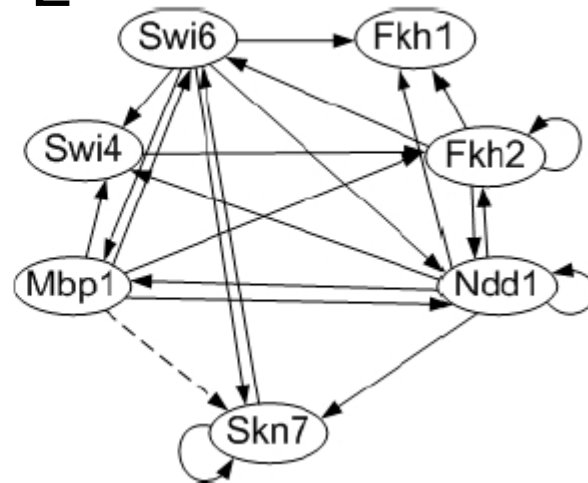
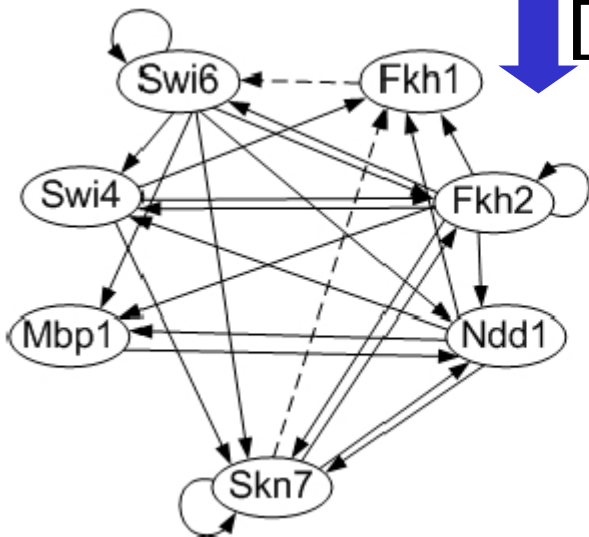


Figure 7. The learned GRN model. (a) The number of discrete levels for gene expression value is 3 and the indegree of the GRN is set to 5. (b) Idem, where the base for gene expression value is 4. The regulators are represented by ovals. The directed edge from Gene A to Gene B means that Gene A is a regulator of Gene B. The solid edges represent regulatory relations that have been verified by other approaches. The dashed edges represent regulatory relations that have not been verified.





# Literature Evidences

Gene	Regulator (Protein)							Evidence
	M1	S4	S6	F1	F2	N1	S7	
MBP1	*3	*	*34		*34	*34		[19], [26]
SWI4	*34	*3	*34		*3	*34	*	[19], [26]
SWI6	* 4	*	*34	3	*34	*	* 4	[19], [26]
FKH1	* 4	*3	*4	*	*34	*34	3	[19], [26]
FKH2	* 4	*34	*3	*3	*34	* 4	*3	[19], [26]
NDD1	*34	*	*34	*	*34	* 4	*3	[19], [26]
SKN7	34	*3	*34		*3	*34	*34	[19]

“\*” means regulatory relations. For example, “\*” in the first cell of first line means that Mbp1 gives MBP1 gene autoregulation [19]. “3” and “4” represent the regulatory relations found with the DFL algorithm when the bases for expression values are 3 and 4 respectively. M1, S4, S6, F1, F2, N1 and S7 are Mbp1, Swi4, Swi6, Fkh1, Fkh2, Ndd1 and Skn7 respectively.

**Table 5. The literature evidences for the GRN model in Figure 7 and Figure 8.**

# The $\epsilon$ -Function method

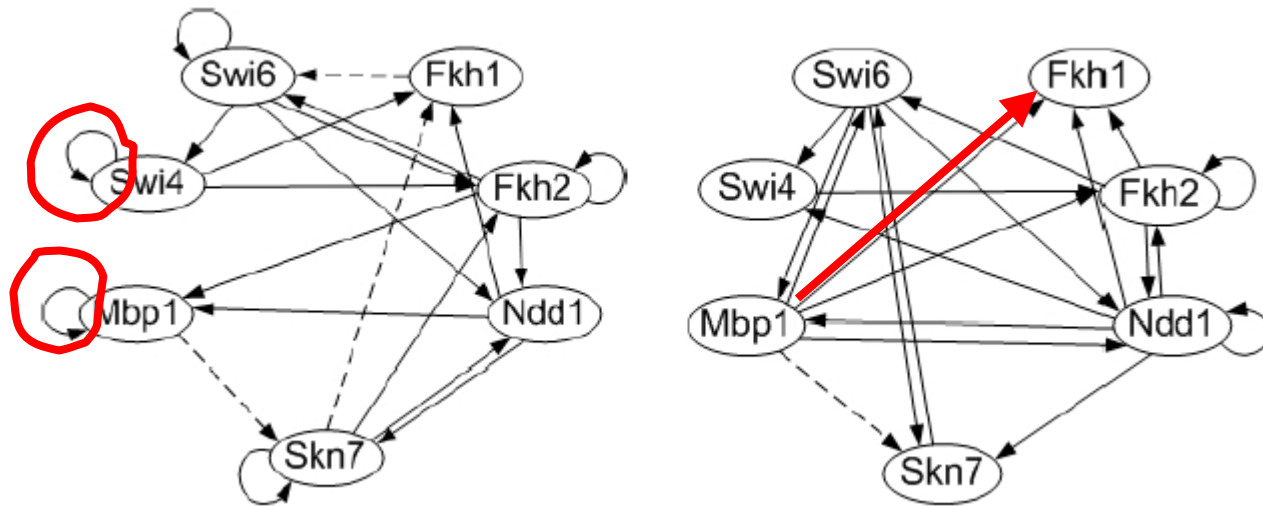
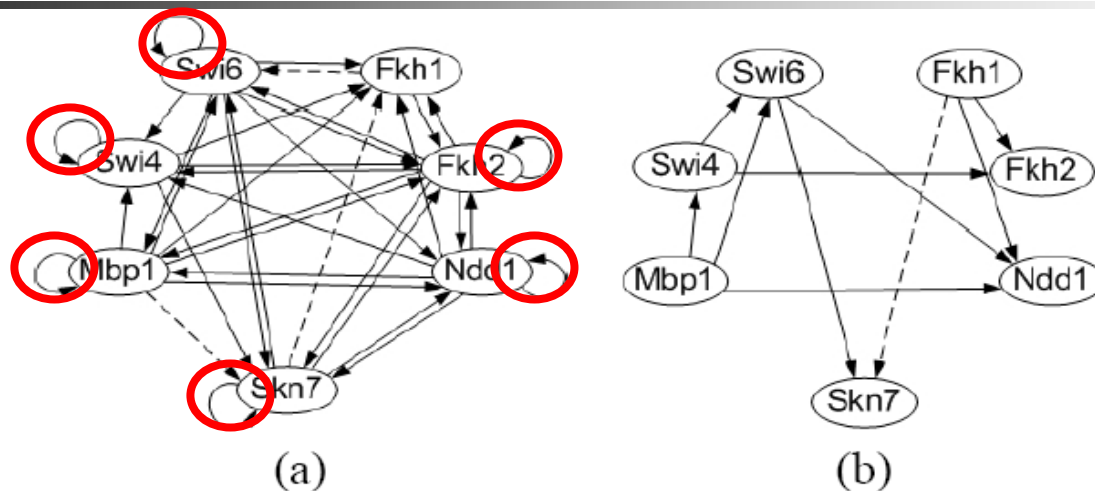


Figure 8. The learned GRN model for yeast cell cycle with the  $\epsilon$  function method. (a) The base for gene expression value is 3, the indegree of the GRN is 5, and the  $\epsilon$  is 0.2. (b) The base for gene expression value is 4, the indegree of the GRN is 5, and the  $\epsilon$  is 0.15. The legends are the same as those of Figure 7.

# A comparison with K2 for learning Bayesian networks



**Figure 9. The combined GRN models. (a) Combined model of Figure 7 and Figure 8. (b) Combined Bayesian network structure learned with the K2 algorithm where the base for expression value is set to 3 and 4 respectively. The legends are the same as those of those of Figure 7.**



# Accuracy and precision

---

$$\text{Accuracy} = \frac{\text{No. of correct predictions}}{\text{No. of predictions}}$$

$$= \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}}$$

$$\text{Precision} = \frac{\text{No. of correct positive predictions}}{\text{No. of positives predictions}}$$

wrt positives

$$= \frac{\text{TP}}{\text{TP} + \text{FP}}$$



# The comparison of prediction performances

	Accur.	Sensi.	Preci.
DFL ( $b = 3$ )	65	67	90
DFL ( $b = 4$ )	63	60	96
DFL (Combined)	80	83	92
K2 ( $b = 3$ )	27	17	88
K2 ( $b = 4$ )	22	12	83
K2 (Combined)	33	24	91

**Table 6. The accuracy, sensitivity and precision of the DFL algorithm and the K2 algorithm.**





# Conclusion

---

- The DFL algorithm is more efficient than current algorithms for reconstructing qualitative models of GRNs without loss of prediction performances.
- The  $\varepsilon$ -Function method is a good supplement to the DFL algorithm.
- The DFL algorithm identifies biologically meaningful GRN models from a limit gene expression profile.



# Questions and suggestions

---

Thanks for your interests!



# Acknowledgements

---

We appreciate Prasanna R Kolatkar and Ng See-Kiong for their reviews on an early version of this paper.

We also thank the two anonymous reviews of the paper.