Models and algorithms for deciphering transcriptional regulations coordinating the response to environmental stresses

Vicente Acuña Aguayo

Scientist Centro de Modelamiento Matemático (CMM), Universidad de Chile

Workshop : Desafíos matemáticos e informáticos para la construcción y análisis de redes de regulación biológica

U. de Concepción 22-23 Abril 2016



イロト イポト イヨト イヨト

э



◆□▶ ◆□▶ ◆□▶ ◆□▶

590





イロト イポト イヨト イヨト

э

### Levels of regulation in bacterial gene expression



FIGURE 17.1 Gene Expression in Bacteria Can Be Regulated at Three Levels.

EXERCISE Label the mode of regulation that is the slowest in response time, and that which is fastest. Label the most efficient and least efficient in resource use.

# **Overview of transcriptional regulation**



.0

# Positive and negative regulation





V. Acuña (CMM, University of Chile) Enumerating substructures in bio networks. Abril 2015, U. de Concepción 6 / 22

### Input Data

Let  ${\mathbb G}$  be the set of genes in the studied organism. The method requires:

Affinity pairs: A set A ⊆ G × G of gene pairs obtained by TF/BS sequence affinity and the *p*-values associated. This forms a directed graph G.

### Input Data

Let  ${\mathbb G}$  be the set of genes in the studied organism. The method requires:

- Affinity pairs: A set A ⊆ G × G of gene pairs obtained by TF/BS sequence affinity and the *p*-values associated. This forms a directed graph G.
- Co-expressed pairs: A set  $C \subseteq \mathbb{G} \times \mathbb{G}$  of pairs of co-expressed genes.

## Input Data

Let  ${\mathbb G}$  be the set of genes in the studied organism. The method requires:

- Affinity pairs: A set A ⊆ G × G of gene pairs obtained by TF/BS sequence affinity and the *p*-values associated. This forms a directed graph G.
- Co-expressed pairs: A set  $C \subseteq \mathbb{G} \times \mathbb{G}$  of pairs of co-expressed genes.
- Validated pairs: Optionally, G can also include a set V ⊆ G × G of gene pairs corresponding to independent experimentally validated regulations, if available.

- 4 同 1 - 4 回 1 - 4 回 1

- 3

#### Definition

Given a pair  $(A, B) \in C$  of co-expressed genes, an *explanation for* (A, B) in G is a set of arcs  $\mathcal{E}$  that satisfy any of the following conditions:

- $\mathcal{E}$  is a directed path from A to B;
- $\mathcal{E}$  is a directed path from B to A;
- $\mathcal{E}$  is the union of two divergent directed paths starting from some gene *C* and arriving respectively at *A* and *B*, having only vertex *C* in common.

- But indeed, we don't want to recover any explanation for each pair, but just to foster simple (short) and confident explanations.
- For each TF/BS affinity, we define a cost on each arc depending on the *p*-value associated to the affinity.
- The cost of an explanation is the sum of the cost of its arcs.



All explanations for the co-expression  $\{F, I\} \in \mathcal{C}$  and the cost associated.

We say that a subgraph  $\mathcal{G}' \subseteq \mathcal{G}$  explains  $\mathcal{C}$  if for every pair  $(A, B) \in \mathcal{C}$  the subgraph  $\mathcal{G}'$  contains an explanation for (A, B).



In the figure: three subgraphs explaining  $C = \{\{A, B\}, \{F, I\}, \{G, H\}\}$ . Only the first and the last are minimal ones. **Modelling objective**: to define a *good* subgraph explaining all co-regulations in  $\mathcal{C}$ .

Algorithmic objective: to enumerate all those subgraphs.

How difficult is enumerate all minimal subgraphs (without considering costs)?

#### Problem

 $\begin{array}{l} {\rm ENUMCOHE}(\mathcal{G},\mathcal{C}): \mbox{ Given an oriented graph } \mathcal{G} \mbox{ and a set of pairs of } \\ {\rm vertices } \ \mathcal{C}, \mbox{ enumerate all minimal subgraphs of } \mathcal{G} \mbox{ that explain } \mathcal{C}. \end{array}$ 

How difficult is enumerate all minimal subgraphs (without considering costs)?

#### Problem

 $\begin{array}{l} {\rm ENUMCOHE}(\mathcal{G},\mathcal{C}): \mbox{ Given an oriented graph } \mathcal{G} \mbox{ and a set of pairs of } \\ {\rm vertices } \ \mathcal{C}, \mbox{ enumerate all minimal subgraphs of } \mathcal{G} \mbox{ that explain } \mathcal{C}. \end{array}$ 

- ENUMCOHE is hard: enumerate all minimal subgraphs that explain C cannot be done in polynomial total time unless P = NP (we reduce from *path conjunction problem*).
- We can still try to develop some heuristics...
- but indeed this model is not very interesting: it generates too many solutions.

(4回) (4回) (4回)

#### More interesting:

#### Problem

 $\begin{array}{l} {\rm ENUMMINCOHE}(\mathcal{G},\mathcal{C}): \mbox{ Given an oriented graph } \mathcal{G} \mbox{ and a set of pairs of } \\ {\rm vertices } \ \mathcal{C}, \mbox{ enumerate all subgraph explaining } \ \mathcal{C} \mbox{ of minimum total cost.} \end{array}$ 

#### More interesting:

#### Problem

 $EnumMinCohe(\mathcal{G}, \mathcal{C}): \ \text{Given an oriented graph } \mathcal{G} \text{ and a set of pairs of vertices } \mathcal{C}, \ enumerate \ all \ subgraph \ explaining \ \mathcal{C} \ of \ minimum \ total \ cost. }$ 

- Unfortunately, even finding **one** subgraph of minimum cost is *NP*-hard (reduction from Steiner Weighted Directed Tree problem).
- Again, we can use some heuristics...
- but indeed this model is not very robust: adding a new pair in C can (in theory) change dramatically the set of solutions.

A E F A E F

Maybe the cell (or evolution) uses a "local parsimony":

3. 3

Maybe the cell (or evolution) uses a "local parsimony":

### Definition

We say that an explanation  $\mathcal{E}$  for the pair (A, B) in  $\mathcal{C}$  is *optimal* if it is of minimum cost among all the explanations for the pair.

Maybe the cell (or evolution) uses a "local parsimony":

### Definition

We say that an explanation  $\mathcal{E}$  for the pair (A, B) in  $\mathcal{C}$  is *optimal* if it is of minimum cost among all the explanations for the pair.

Now, we can define the solutions that contain an optimal explanation for every pair in  $\ensuremath{\mathcal{C}}.$ 

#### Definition

We say that a subgraph is an *optimal* subgraph explaining C if it is the union of |C| optimal explanations, one for each pair  $(A, B) \in C$ .

▲ 臣 ▶ | ▲ 臣 ▶

#### Definition

We say that a subgraph is an *optimal* subgraph explaining C if it is the union of |C| optimal explanations, one for each pair  $(A, B) \in C$ .

- Instead of enumerating all *optimal* subgraphs, we define G<sub>L</sub> the union of all *optimal* subgraph explaining C.
- To compute  $\mathcal{G}_L$ , we need simply to compute all optimal explanations for each pair in  $\mathcal{C}$ .
- The software LOMBARDE uses an algorithm of minimum paths in modified graph to compute  $\mathcal{G}_L$ .



990

E

◆□▶ ◆□▶ ◆□▶ ◆□▶

Network	Explained co-expressions		No . of vertices	No. of arcs	No. of arcs in ${\cal V}$
TRN $\mathcal{G}_{\mathcal{V}}$ built from $\mathcal{V}$	3,990	(6.5 %)	823	1,652	1,652
E. coli ab initio $\mathcal{G}_{\mathcal{A}}$	56,044	(91.1 %)	2,390	25,604	444
Lombarde output for $\mathcal{G}_\mathcal{A}$	56,044	(91.1 %)	2,336	4,922	295
E. coli extended $\mathcal{G}_{\mathcal{AV}}$	56,789	(92.3 %)	2,434	26,812	1,652
Lombarde output for $\mathcal{G}_{\mathcal{AV}}$	56,789	(92.3 %)	2,370	4,374	1,520

## **Table 1** Characteristics of the *a priori* graphs and LOMBARDE output networks

The a priori graphs explained most of the co-expressions. The LOMBARDE results kept most of the vertices, significantly reduced the number of arcs, and kept most of the validated arcs

イロト イポト イヨト イヨト

#### p-value distribution on putative TRN











p-value

イロト イポト イヨト イヨト

E



< 冊

590

э

-



<ロ> (日) (日) (日) (日) (日)

E



< A

-