

Gene Regulatory Networks: Introduction of multiplexes into R. Thomas' modelling

Z. Khalis^{1,2}, G. Bernot^{1,2}, J.-P. Comet¹ and Observability Group²

¹ I3S, CNRS & Univ. of Nice-Sophia Antipolis, F-06903 Sophia Antipolis, France

² Epigenomics Project, Genopole, F-91034 Évry, France

email: {khalis,bernot,comet}@i3s.unice.fr

Abstract

When modelling gene regulatory networks, the cornerstone of the modelling process is the search of parameter values which are consistent with the known properties of the system. These parameters drive the dynamics of the system. In this article, we give a formal definition of a slight extension of the R. Thomas' modelling framework, with explicit information about cooperative, concurrent or more complex molecular interactions. It considerably decreases the number of parameters and determining parameter values becomes less time consuming, making possible the study of larger systems.

1 Introduction

To study complex biological systems, formal modelling is often mandatory since the complexity of the interleaved interactions between constituents makes intuitive reasoning error prone. Numerous mathematical modelling frameworks have been proposed to model gene regulatory networks, see for example [7, 13, 20, 8]. Common approaches are quantitative, based on differential or stochastic equations, providing numerical simulations of the system. Nevertheless actual predictions often remain only qualitative because the parameter values of these systems are not precisely known. Several other modelling frameworks are based on a qualitative view, see for example boolean networks and their generalizations [16, 19], Petri nets [4, 6], hybrid modellings [12, 1], and stochastic π -calculus [5]. Each modelling framework highlights some views of models and allows one to detail or to abstract different biological aspects.

We focus here on Thomas' modelling, in which the gene regulatory system is represented by an interaction graph and a set of parameters. The interaction graph is composed, on the one hand, of nodes which abstract genes and their proteins, and on the other hand, of edges which represent the interactions between the genes. The values assigned to the parameters permit one to deduce the dynamics of the system from the interaction graph. Even in a qualitative perspective, the lack of reliable data about the system leads to a typical difficulty of the modelling approach : How to select the parameter values of the model?

For determining values of parameters, we proposed in [3] to test the set of all possible parameterizations against temporal properties. It is finite in the

case of Thomas' modelling. This approach can be computer aided [3] using formal temporal logics and systematic model checking. Even if the set of possible parameterizations is finite, it exponentially grows with the size of the interaction graph. Several theorems established in the Thomas' framework considerably reduce the number of generated parameter sets, nevertheless, an entire exploration is not conceivable for large networks.

In order to reduce the time required by this exploration step, it becomes crucial to introduce in the modelling framework more biological information (when available). In this chapter, we propose to take into account information about how constituents of the system act on their targets. For example (Figure 1), if two genes act positively on a common target *via* the formation of a complex (e.g. the transcription factor of the common target contains the complex), then it is obvious that the common target has in fact a unique predecessor (the complex instead of two genes separately) and only two possibilities (instead of four) can occur: The complex is present and the transcription can take place or the complex is not present. Indeed this idea is far from being new but it has never been formalized up to now. R. Thomas remarked that this kind of information can be taken into account in its modelling framework through the valuation of parameters, but he did not explicitly include such information in the interaction graph [18].

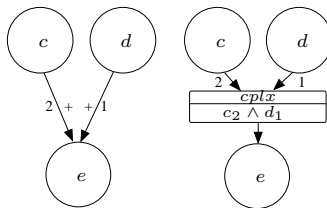


Figure 1: Example of cooperative action

Here, we propose a modelling framework in which the interaction graph makes such cooperative or concurrent biological phenomenon explicit. The decreasing of the number of parameters coupled with the methodology developed in [3], will make possible the study of larger systems.

The chapter is organized as follows. We firstly define our new interaction graph: *Multiplexes* are formally defined to take into account available biological information describing the cooperation or concurrency between constituents acting on a common target. Then we define when a multiplex has an effective action on its targets, and we construct the associated dynamics. We show that Thomas' and multiplex frameworks have the same power of expression but we illustrate, through the classical example of the lac operon, how multiplexes allow us to be more legible and terse. Lastly we present the benefits of this multiplex modelling.

2 GENE REGULATORY GRAPHS WITH MULTIPLEXES

Formal modelling frameworks for gene regulatory networks represent interactions between entities (genes, proteins, *etc.*) *via* a *static graph*. Then, dynamics focus on the evolution of entity expression levels and ask for more elaborated mathematical stuff with many parameters.

In our framework, we represent the static part by a directed graph composed of two types of vertices: Variables which correspond to genes and their products, and multiplexes which correspond to interactions between variables. Multiplexes abstract biological phenomena like complex forming or more elaborated phenomena. The predecessors of a multiplex are either variables or other multiplexes brought into play in the interaction; the successors are called the targets of the interaction.

2.1 Formal Definition

The following notation will be useful.

Notation 1 Given a directed graph G and a node v of G , $G^{-1}(v)$ is the set of all nodes v' of G such that (v', v) is an edge of G (set of predecessors of v).

A multiplex is provided with a formula in a propositional logic which encodes the situations in which the interaction occurs. For example, if a complex formed with proteins a and b is required in cooperative action and if the complex $(a-b)$ is inactive in the presence of a protein c , then the corresponding formula looks like “ $a \wedge b \wedge \neg c$,” where the symbols “ \wedge ” and “ \neg ” stand for “and” and “not” respectively.

Definition 1 A gene regulatory graph with multiplexes, RG *for short*, is a tuple $G = (V, M, E_V, E_M)$ such that:

1. $(V \cup M, E_V \cup E_M)$ constitutes a (labelled) directed graph whose set of nodes is $V \cup M$ and set of edges is $E_V \cup E_M$, with $E_V \subset V \times \mathbb{N} \times M$ and $E_M \subset M \times (V \cup M)$.
2. V and M are disjoint finite sets. Nodes of V are called variables and nodes of M are called multiplexes. An edge (v, s, m) of E_V is denoted $(v \xrightarrow{s} m)$ where s is called the threshold.
3. Each variable v of V is labelled with a positive integer b_v called the bound of v .
4. Each multiplex m of M is labelled with a formula belonging to the language L_m inductively defined by:
 - If $(v \xrightarrow{s} m) \in E_V$, then v_s is an atom of L_m , and if $(m' \rightarrow m) \in E_M$ then m' is an atom of L_m .

- If ϕ and ψ belong to L_m then $\neg\phi$, $(\phi \wedge \psi)$, $(\phi \vee \psi)$ and $(\phi \Rightarrow \Psi)$ belong to L_m .

5. All cycles of the underlying graph $(V \cup M, E_V \cup E_M)$ contain at least one node belonging to V .

Note: Condition 5 is necessary for the definition of dynamics (see Definition 3).

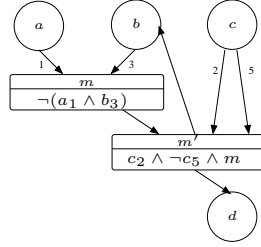


Figure 2: Graphical conventions

Figure 2 provides graphical conventions. In this figure, a, b, c, d are variables; m, m' are multiplexes; m and c are the inputs of m' and b and d are its outputs; the cycle b, m, m' contains the variable b .

In addition to these standard graphical conventions, we allow “light” additional graphical notation abuse:

- If a variable is an input of a multiplex with only one threshold, we then allow to omit the threshold in the formula. For example, in Figure 2, the formula of multiplex m can be simply written as “ $\neg(a \wedge b)$.” Of course, this light form is not possible for m' .
- Multiplexes with a formula reduced to a unique atom can be removed from the diagram. In figure 3a, removing the multiplex m allows us to retrieve the usual diagrammatic convention of R. Thomas for activations.
- Similarly, in figure 3b, we retrieve usual inhibitions, either by adding the minus sign, or by using the “inhibition arrow” usual in biology.

In figure 2, we also see that in multiplex formulas the variables are indexed by their thresholds. This is useful when a given variable acts on a multiplex at several thresholds. The multiplex m' means that the expression level of c must be both greater than 2 and lower than 5 in order to participate to the induction of d .

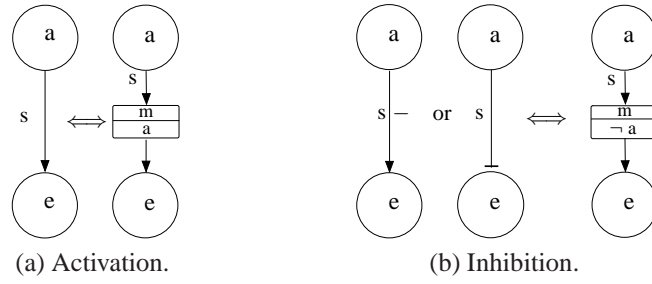


Figure 3: Light graphical convention for activation and inhibition.

2.2 States and resources

A gene regulatory graph with multiplexes constitutes the static representation of the system. We have now to focus on the dynamics of the system, abstracted by the evolutions of expression levels of variables. Let us first define the states of a system.

Definition 2 A state of a RG $G = (V, M, E_V, E_M)$ is a map $\eta : V \rightarrow \mathbb{N}$ such that for each variable v belonging to V , $\eta(v) \leq b_v$. $\eta(v)$ is called the expression level of v .

A multiplex does not have any expression level because it is a logical composition of variables at a given time. So, we consider only the expression level of all the variables at that time and from this current state it is possible to deduce if the multiplex is active or not *via* the interpretation of its propositional formula.

According to a current state, the set of resources of a variable a is the set of multiplexes which can help a to express its product. More precisely a resource r of a variable a is a multiplex belonging to $G^{-1}(a)$ whose formula is satisfied. So graphically, edges of interaction graphs have no sign but negative actions are taken into account through multiplexes with the operator \neg . For example, in Figure 2 the multiplex m represents an inhibition (the complex a - b inhibits b and d via m').

Definition 3 Given a RG $G = (V, M, E_V, E_M)$ and a state η of G , the set of resources of a variable $v \in V$ for the state η is the set of multiplexes m of $G^{-1}(v)$ such that the formula φ_m of the multiplex m is satisfied. The interpretation of φ_m in m is inductively defined by:

- If φ_m is reduced to an atom v_s of $G^{-1}(m)$ then φ_m is satisfied iff $\eta(v) \geq s$.
- If φ_m is reduced to an atom $m' \in M$ of $G^{-1}(m)$ then φ_m is satisfied iff $\varphi_{m'}$ of m' is satisfied.

- If $\varphi_m \equiv \psi_1 \wedge \psi_2$ then φ_m is satisfied if ψ_1 and ψ_2 are satisfied; and we proceed similarly for all other connectives.

We note $\rho(v, \eta)$ the set of resources of v for the state η .

This definition is actually inductive because RG never contain a cycle of multiplex (item 5 of Definition 1). If cycle of multiplexes were allowed then indeterminations or contradictions would be possible. For instance, consider the graph in figure 4. Suppose that the expression level of a is greater or equal to the threshold s :

- If the formula of m' is assumed to be satisfied, then the formula of m must be satisfied and so the formula of m' cannot be satisfied. So, we get an inconsistency.
- If the formula of m' is assumed to be unsatisfied, then the formula of m must be unsatisfied and so the formula of m' must be satisfied. So, whatever we assume, we always get an inconsistency.

Let us consider now, the graph in figure 4 where the formula associated with m' is m instead of $\neg m$. Suppose again that the expression level of a is greater or equal to the threshold s . Then, the two interpretations of m' are consistent and compatible with the current state. There is an indetermination which is similar to the notion of schizophrenic cycles of [15].

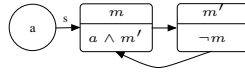


Figure 4: Cycle of multiplexes

To avoid these inconsistencies and indeterminations, cycles of multiplexes are not allowed. This motivates the item 5 of Definition 1.

3 GENE NETWORKS WITH MULTIPLEXES

We call network a graph associated with the parameters which determine the dynamics.

Definition 4 A gene regulatory network with multiplexes (RN) is a couple (G, \mathcal{K}) where

- $G = (V, M, E_V, E_M)$ is a RG.
- $\mathcal{K} = \{k_{v,\omega}\}$ is a family of parameters indexed by $v \in V$ and $\omega \subset G^{-1}(v)$ such that all $k_{v,\omega}$ are integers and $0 \leq k_{v,\omega} \leq b_v$.

Notice that each variable v admits 2^n parameters of the form $k_{v,\omega}$ where n is the in-degree of v in G .

Additional restrictions for the choice of parameters can be considered. The Snoussi's hypotheses [14] which ensure the consistency of qualitative behaviours with some underlying differential equation system, are well-known: If $\omega \subset \omega'$ then $k_{v,\omega} \leq k_{v,\omega'}$. These hypotheses signify that an effective resource cannot induce the decrease of the expression level of v . Moreover, we can always ignore the parameters $k_{v,\omega}$ such that the conjunction of the formulas associated with the multiplexes in ω is unsatisfied for all states.

3.1 Dynamics

The value of the parameter $k_{v,\rho(v,\eta)}$ (where ρ is defined in definition 3 above), indicates how the expression level of v can evolve from the state η . It can increase (respectively decrease) if the parameter value is greater (respectively less) than $\eta(v)$. The expression level must stay constant if both values are equal. The tendency (increasing, decreasing, unchanging) of variables are given by the directional map associated with each state:

Notation 2 Given a RN $N = (G, \mathcal{K})$ and a state η of $G = (V, M, E_V, E_M)$, the directional map $d : V \rightarrow \{-1, 0, 1\}$ is defined by:

$$\forall v \in V, d(v) = \begin{cases} -1 & \text{if } \eta(v) > k_{v,\rho(v,\eta)} \\ 0 & \text{if } \eta(v) = k_{v,\rho(v,\eta)} \\ 1 & \text{if } \eta(v) < k_{v,\rho(v,\eta)} \end{cases}$$

The probability that two variables change their expression level at the same time is negligible *in vivo*; following the Thomas' approach a state transition of the model modifies only one of the involved variables at a time.

Definition 5 Let $N = (G, \mathcal{K})$ be a RN, and let η be a state of G . A state η' of G is a successor of the state η if and only if:

- There exists a variable u such that $\eta'(u) = \eta(u) + d(u)$ and $d(u) \neq 0$
- For any other variable $v \neq u$ we have $\eta'(v) = \eta(v)$

In each state transition, at most one variable is modified; this procedure is called *asynchronous update* in Thomas' framework.

Definition 6 The asynchronous state graph of a RN $N = (G, \mathcal{K})$ is the graph S defined by:

- The set of vertices of S is the set of possible states of G (isomorphic to the Cartesian product $\prod_{v \in V} [0, b_v]$).
- The set of edges of S is the set of couples (η, η') such that η' is a successor of η .

4 RELATIVE TERSENESS WITH RESPECT TO THE CLASSICAL FRAMEWORK

Obviously our framework with multiplexes embeds the classical Thomas' framework [17] as it is sufficient to translate an activation (resp. an inhibition) with a multiplex whose formula is reduced to the input variable (resp. its negation), see Figure 3. Conversely, a non atomic formula in a multiplex obviously corresponds to a constraint on the parameters [18] following an induction similar to the one of Definition 3.

Our conviction is that this kind of knowledge is a static knowledge and consequently it should be present in the interaction graph (formulas in multiplexes). When we know, for biological reasons, the nature of combined influences, this information should be included in the model as soon as possible because it considerably reduces the number of possible parameters, as shown in the example below. Of course, the nature of combined influences is not always *a priori* known and, in this case, according to our formalism, variables have then several inputs in the regulatory graph.

4.1 Example of lactose operon.

The cell needs carbon. Carbon is preferably obtained from glucose *via* a given catalytic pathway. When glucose is absent, lactose is used *via* an alternative catalytic pathway.

Lactose operon in E.coli is the first genetic regulatory system elucidated, by François Jacob and Jacques Monod [9]. The induction of this system requires two conditions: Absence of glucose and presence of lactose.

An operon is a set of contiguous genes whose transcription is controlled simultaneously by a unique transcription factor. This transcription factor has an affinity with a DNA area at the beginning of the operon, called operator and denoted O.

The lactose operon is formed by three genes denoted by Z, Y and A. The genes Z, Y and A produce respectively the enzymes β -galactosidase, permease and thiogalactoside transacetylase.

When glucose is absent, the alternative pathway is controlled as follows:

- CAP (Catabolite gene Activator Protein) forms a complex with cAMP (cyclic Adenosine MonoPhosphate), and binds to DNA to increase the transcription of the operon. This is a positive regulation.
- The transcription of the operon is possible only if the DNA area O is free. The regulatory protein lacI binds to O, this is a negative regulation. However, when lactose is present, a lactose isomer binds to lacI and lacI loses its affinity for O. So the operator O becomes free.

When glucose is present, the alternative pathway is inhibited as follows: Glucose inhibits indirectly cAMP and leads to the absence of complex CAP-cAMP. Consequently, there is no transcription even if lacI is present.

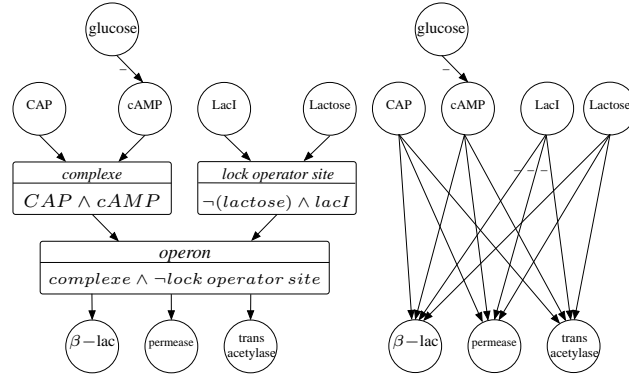


Figure 5: Lactose operon metabolism graph with and without multiplexes.

In Figure 5, the interaction graph of the lactose operon is represented in the multiplex framework (the left part of the figure) and in the classical Thomas' framework (the right part of the figure). The first advantage of the multiplex framework is its legibility: The left hand side of the figure is, to some extent, more legible than the textual description given before. On the contrary, the right hand side of the figure cannot be understood without the textual description.

The second advantage of the multiplex framework is methodological. When we try to elucidate a biological system using Thomas' framework, we do not know the values of the parameters: The $k_{v,\omega}$ have to be inferred from *in vivo* behaviours. Consequently, models with a small number of parameters allow us to rapidly converge towards the elucidation of the studied biological system. On the contrary, models with large numbers of parameters can be so heavy to manipulate that they obstruct the discovery process. On this small lactose operon example, the total number of parameters according to the multiplex approach is 12, while the total number of parameters according to the classical approach is 54. Putting as much static information as possible explicitly in the graph (instead of putting it later manually in the dynamics) considerably reduces the complexity of the modelling methodology. Indeed, formalizing cooperative actions of several variables on the same target *via* multiplexes enables one to merge into a single multiplex the different acting resources.

The knowledge formalised into multiplexes can lead to reduce even more the number of useful parameters. In figure 6, multiplexes m_1 and m_2 cannot be satisfied for the same state: m_1 is active only if expression level of a is strictly less than 2 whereas m_2 is active when expression level of a is greater or equal to 2. Among the set of formal parameters $\mathcal{K} = \{k_{c,\{ \}}, k_{c,\{m_1\}}, k_{c,\{m_2\}}, k_{c,\{m_1,m_2\}}\}$, $k_{c,\{m_1,m_2\}}$ is never used. More generally, when two multiplexes having the same target v have two *mutually exclusive formulas* ϕ_1 and ϕ_2 , all parameters of the form $K_{v,\omega \cup \{m_1,m_2\}}$ can be ignored and the number of relevant parameters is reduced.

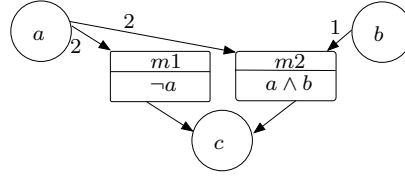


Figure 6: Example of RG which contains mutually exclusive formulas

5 Application

The software SMBionet-3.0 [11] has been designed to facilitate the modelling process of genetic regulatory systems. It allows one to select models of given RG according to their temporal properties. It takes as input a RG and a formula in temporal logic expressing the known or hypothetical temporal properties of the system. It gives as output all the models satisfying the formula.

In both modelling frameworks (with or without multiplexes), we have to give a value to each parameter in order to deduce the dynamics of the system. Because parameter values are not *a priori* known this leads us to consider an enormous number of parameterizations. Indeed, each variable v admits 2^n parameters of the form $k_{v,\omega}$ where n is the in-degree of v in G ($\omega \subset G^{-1}(v)$). Each of these parameters can take $b_v + 1$ different values where b_v is the bound of v . The number of parameterizations is then given by $\prod_{v \in V} (b_v + 1)^{2^n}$ where n is the in-degree of v . For the example of lactose operon in Thomas' framework, the number of parameterizations is on the order of 2.27×10^8 whereas in our multiplex framework, the number of parameterizations is 1296. For instance, in Thomas' framework, the variable *permease* has 2^4 parameters, generating 2^{2^4} (65536) different parameter settings while in our framework, *permease* has 2 parameters, generating 2^2 (4) different parameter settings. The difference resides in the addition of the multiplexes, which reduces the number of inward edges to *permease* and so the number of possible parameter settings. Consequently, taking into account information about cooperation between variables (through multiplexes) leads to a significant decreasing of the number of possible models: Here, the set of possible models is cut down by a factor of 175000.

We used SMBionet-3.0 to exhibit models which present characteristic alternative catalytic pathway when glucose is absent. Under the Snoussi's hypotheses (see section 3 Biological Regulatory Networks with multiplexes) and for a given logical formula, all possible parameter settings in our framework have been explored in 27 seconds whereas all possible parameter settings in Thomas' framework have been explored in approximately 1000 hours. Notice that the ratio between both time is less than 175000 because SMBionet-3.0 optimizes the exploration of the model set.

6 CONCLUSION

We rigorously introduced propositional logic elements in the R. Thomas' framework in order to take into account available information concerning the cooperation or concurrency between genes or genes products acting on the same targets.

This idea is rather natural: R. Thomas introduced in [17] a notation that allows the representation of *several* actions of a *unique* gene on another one. Moreover, dozens of articles can be cited which use similar ideas in different frameworks:[2, 10], *etc.* Up to our knowledge, our contribution is the first one which rigorously *formalizes* this more elaborated framework.

The introduction of multiplexes makes models terser because this framework allows the gathering of edges into a single multiplex.

The major advantage of multiplex modelling is methodological: It reduces the number of parameters by formalizing additional biological information. So, the step which searches parameter values consistent with known or hypothetical properties of the system is significantly improved. These advantages open perspectives to study larger gene regulatory networks.

Another advantage of multiplexes is to facilitate manipulations of networks. For example, we may develop graph folding methods in order to reduce the number of variables, at the price of possibly long formulas in multiplexes. However the role of some variables in a path is essentially to delay the global process. Consequently to improve the biological usefulness of such abstractions, it seems necessary to take delays into account. One of our future works will be to introduce delays in multiplexes.

References

- [1] J. Ahmad, G. Bernot, J.-P. Comet, D. Lime, and O. Roux. Hybrid modelling and dynamical analysis of gene regulatory networks with delays. *ComPlexUs*, 3(4):231–251, 2006 Cover Date: November 2007).
- [2] R. Albert and H. Othmer. The topology of the regulatory interactions predicts the expression pattern of the segment polarity genes in *drosophila melanogaster*. *J. Theor. Biol.*, 223:1–18, 2003.
- [3] G. Bernot, J.-P. Comet, A. Richard, and J. Guespin. Application of formal methods to biological regulatory networks: Extending Thomas' asynchronous logical approach with temporal logic. *J. Theor. Biol.*, 229(3):339–347, 2004.
- [4] C. Chaouiya, E. Remy, P. Ruet, and D. Thieffry. Qualitative modelling of genetic networks: From logical regulatory graphs to standard Petri nets. In *ICATPN 2004*, LNCS 3099, pages 137–156. Springer-Verlag, 2004.

- [5] F. Ciocchetta and C. Priami. Biological transactions for quantitative models. *ENTCS*, 171(2):55–67, 2007.
- [6] J.-P. Comet, H. Klaudel, and S. Liauzu. Modeling multi-valued genetic regulatory networks using high-level Petri nets. In *ICATPN 2005*, volume 3536 of *LNCS*, pages 208–227, 2005.
- [7] H. de Jong. Modeling and simulation of genetic regulatory systems: a literature review. *J. Comput. Biol.*, 9(1):67–103., 2002.
- [8] S. Fujita, M. Matsui, H. Matsuno, and S. Miyano. Modeling and simulation of fission yeast cell cycle on hybrid functional Petri net. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences*, E87-A(11):2919–2928, 2004.
- [9] F. Jacob and J. Monod. Genetic regulatory mechanisms in the synthesis of protein. *Journal of molecular biology*, 3:318–356, 1961.
- [10] S. Klamt, J. Saez-Rodriguez, J. A. Lindquist, L. Simeoni, and E. D. Gilles. A methodology for the structural and functional analysis of signaling and regulatory networks. *BMC Bioinformatics*, 7:56, 2006.
- [11] A. Richard. *Modèle formel pour les réseaux de régulation génétique et influence des circuits de rétroaction*. PhD thesis, Univ. Evry Val d’Essonne, 2006.
- [12] H. Siebert and A. Bockmayr. Context sensitivity in logical modelling with time delays. In LNBI, editor, *Computational Methods in Systems Biology, CMSB 2007*, volume 4695, pages 64–79. Springer, 2007.
- [13] P. Smolen, D.A. Baxter, and J.H. Byrne. Modeling transcriptional control in gene networks: Methods, recent results, and future directions. *Bulletin of Mathematical Biology*, 62(2):247–292, 2000.
- [14] E.H. Snoussi. Necessary conditions for multistationarity and stable periodicity. *Journal of Biological Systems*, 6:3–9, 1998.
- [15] O. Tardieu and R. de Simone. Curing schizophrenia by program rewriting in Esterel. *Formal Methods and Models for Co-Design*, pages 39–48, 2004.
- [16] R. Thomas. Boolean formalization of genetic control circuits. *J. Theor. Biol.*, 42:563–585, 1973.
- [17] R. Thomas. Regulatory networks as seen asynchronous automata: A logical description. *J. theor. Biol.*, 153:1–23, 1991.
- [18] R. Thomas and R. d’Ari. *Biological Feedback*. CRC Press, 1990.

- [19] R. Thomas, D. Thieffry, and M. Kaufman. Dynamical behaviour of biological regulatory networks. *Bull. Math. Biol.*, 57(2):247–276, 1995.
- [20] D.J. Wilkinson. *Stochastic Modelling for Systems Biology*. Chapman & Hall/CRC, 2006.