Determination, optimization and taxonomy of regulatory networks The example of *Arabidopsis thaliana* flower morphogenesis

Nicolas Glade

TIMC-IMAG Laboratory, CRNS UMR5525, UJF & AGIM Laboratory, FRE UJF Email: nicolas.glade@imag.fr Telephone.: +33-4-56-52-00-26 Fax: +33-4-76-76-88-44

Fabien Corblin

TIMC-IMAG Laboratory, CRNS UMR5525 & UJF Email: fabien.corblin@imag.fr Telephone.: +33-4-56-52-00-27

Jacques Demongeot

TIMC-IMAG Laboratory, CRNS UMR5525, UJF & AGIM Laboratory, FRE UJF Email: jacques.demongeot@imag.fr Telephone.: +33-4-56-52-01-08

Abstract—This paper aims at warning modellers in systems biology against several traps encountered in the modelling of Boolean thresholded automata networks, *i.e.* the Hopfield-like networks that are often used in the context of neural and genetic networks. It introduces a new manner based on inverse methods to conceive such models. Using these techniques, we re-visit the model of regulatory network of *Arabidopsis thaliana* morphogenetic network. In this context, we discuss about the non-uniqueness of models, on a possible taxonomy of the set of valid models and on the sense of the relative size of the basin of attractions within or between these models.

I. INTRODUCTION

In the present paper, in the first part, we discuss principally of the methodological paradigms that usually drive the understanding of biological regulatory networks. Nowadays, biological complexity (in the sense of modelling and computational complexity) and the necessity to ask correctly Adrien Elena

AGIM Laboratory, FRE UJF Email: adrien.elena@imag.fr

Eric Fanchon

TIMC-IMAG Laboratory, CRNS UMR5525 & UJF Email: fabien.corblin@imag.fr Telephone.: +33-4-56-52-00-27

Hedi Ben Amor

TIMC-IMAG Laboratory, CRNS UMR5525, UJF & AGIM Laboratory, FRE UJF Email: hedi.ben-amor@imag.fr Telephone.: +33-4-56-52-00-26

the biological question by using the best theoretical representations, is an evidence. The democratization in biology of modelling techniques and tools coming from mathematics or computer sciences is a consequence of this quest for the understanding of biological complex systems and it led to the development of numerous models and many breakthroughs in biology. However, their pragmatical and often naive use is dangerous and many traps hide from the researchers. Here, we raise the problem of the good use of formalisms and that of the non-uniqueness of consistent models in biological representations of regulatory networks in the case of Hopfieldlike networks. The so-common trial-error approach (a trialerror looping process) that reaches only one solution of model lets indeed the biologist believing in the uniqueness of this solution in the context of the available knowledge. We will show that a constraint-based approach (an inverse method) is infinitely more adapted to the search for valid models of

978-0-7695-4338-3/11 \$26.00 © 2011 IEEE DOI 10.1109/WAINA.2011.155 Boolean thresholded regulatory networks.

A second part is devoted to the taxonomic classification of the set of valid models into homogeneous groups of models based on robustness criteria. We discuss the importance of the size of basins of attraction and their relation to the robustness of biological systems. The regulatory network of *Arabidopsis thaliana* flower morphogenesis, as modelled by Mendoza et al [1], is chosen to illustrate these points. This very interesting work constitued a real advancement in the understanding of *A. thaliana* flower morphogenesis. However what we discuss here about their paper (their modelling approach and the formalism they used) is very emblematic of biological modelling nowadays.

In addition to introducing new features in the search for valid models in the context of Hopfield-like networks, our paper is a warning to the modellers, especially biologists, against the apparent simplicity and universality of application of this formalism.

II. FORMALISM AND METHODOLOGICAL CONSIDERATIONS

A. Trial-error vs constraint-based approaches

Unless having a detailed knowledge on the kinetics of the studied system, the modeller often starts by the elaboration of a model of network based on the formalization of a qualitative experimental knowledge. This knowledge is composed of two aspects : the structure of the network (the nodes and their interactions) and its dynamics (stationary or transitional states). These aspects are dual since the one can be obtained from the other on condition that the knowledge is complete. The most common approach is to see the dynamics of the network as the result of its structure. In an experimental context the knowledge is always incomplete making its modelling difficult and subject to discrepancy. A classical method (Fig. 1 Top) is to formalize the a priori knowledge concerning the structure of the network, and then, by using a trial-error looping process, to try to converge towards an optimal model that minimizes, by comparison, the error with the experimental observations on the dynamics. The learning process of artificial neural networks based on a cost function corresponds to this approach. All the set of parameters can be explored so as to check the behaviour of all possible models. Monte carlo methods, sequential simulated annealing or genetic algorithms can be also used to explore the set of models and find one of them that fits with the expected - experimentally observed - behaviour, 'The Optimal One'. Rapidly, one will realize the expensive computational cost of this approach. Moreover it does not guarantee the complete adequacy of the 'solution' to the experimental knowledge : one can find a model that minimizes the error between simulated and experimental dynamics, but nothing guarantees that it is absolutely conform to the observations.

A declarative approach (like in Fig. 1 Bottom) was proposed by [2], [3] as an alternative to this trial-error process in the context of Thomas' networks [4]. These networks take into account the cellular context in which the interaction occurs. This implies the existence of multivalued arcs that represent the different possible kinetics between two entities (chemical species) as a function of their concentration levels. Its application to the re-examination of an existing model of the nutritional stress in *Escherishia* coli allowed to show an incoherence of the obtained model. Experimental data are indeed not exhaustive and, in consequence, the intuitions of the modeller constitute, in case of critical uncompleteness, a considerable contribution for determining solutions. The reexamination of the model by this declarative approach could reject an intuition and propose automatically an alternative.

Such an approach is a potential link between structure and dynamics of the network. The works realized by [5]-[8] focus particularly on the relations between regulatory network structure and dynamics, in the case of Hopfield-like networks, but also on the robustness of the dynamics obtained depending on the update modes of the network (parallel, sequential, bloc-sequential) [5]-[8]. In the same context (Hopfield-like networks), we developed a similar constraint-based approach in order to guarantee the flexibility of the modelling process and the adequacy of the obtained solutions [9]-[11]. The knowledge concerning the structure and the dynamics is formalized in the form of a series of constraints. Then, a query in the normal conjonctive form is defined [12] and submitted to a satisfiability solver [13]-[16]. We obtain, without defining any cost function, a set of valid instances of models (that can be empty in case of thin consistency), *i.e.* they are all in adequacy with experimental data (structure and dynamics).



Fig. 1. **The modelling approaches** (Top) The classical trial-error approach starts with data on interactions between the elements (*e.g.* genes) of a supposed network, and via a trial–error and validation looping process one progressively gets closer to the fitness corresponding to the experimental knowledge on the dynamics. At the end, only one model is selected by using a certain criterion of optimality. (Bottom) On the contrary, the constraint-based approach, starting from the global knowledge, *i.e.* both structural and dynamical knowledge, does not need any trial-error and validation looping process. From that global knowledge converted into a series of constraints, a logical formula is written that defines the space of consistent models. All models in this space are equally valid.

B. Hopfield-like networks

Biological regulatory networks are often abstracted as interaction graphs. They are modelled in the form of Boolean - automata - networks composed of nodes representing the components of the system (genes, proteins, cells) linked together by oriented arrows indicating the relationships between them. Boolean automata networks (introduced by Kauffman to study global properties of genetic nets [17]) are among the most used models in biological modelling of regulatory networks. Hopfield-like network formalism (Fig. 2) is based on a Boolean thresholded automaton. This model is similar to Hopfield's model [18] but it is more flexible since there are no conditions of symmetry imposed on the weights and self-interaction loops are authorized. A collection of two-state (but it can be extended to n-valued) and thresholded nodes is in interaction. Interactions are encoded in the form of arcs linking two nodes in a directional way. Each arc has a certain strength called weight that will determine the nature of the interaction. The update of the node states is determined by a test called transition function implying a target node and a set of efferent nodes acting on it : the sum of the products of efferent state values and the weights of the connecting arcs is compared to the threshold of the target node ; when superior, the state of the node changes.



Fig. 2. **Hopfield-like networks.** Hopfield-like networks, inspired from Hopfield's formalism for neural networks but more flexible, are composed of n nodes of state S_i having a threshold of activation θ_i . The arcs are the relations of inhibition or activation and their strengths are noted w_{ij} . When the sum of efferent states S_j modulo the weights w_{ij} of the arcs is superior to its threshold θ_i , the target node i is activated. The Hopfield's transition function is given at the bottom of the figure. The example given here corresponds to the case $W_{ij} \ge 0$ and $\theta_j \ge 0$. One of its behaviours is a cycle of length 2 (100 \leftrightarrow 001).

C. A formalism is an abstraction that cannot represent the overall knowledge

When working on the Arabidopsis thaliana flower morphogenesis regulatory network (see below), we discovered that some relations in the Hopfield-like model of this network by Mendoza et al [1] were not well-founded. In fact, the authors tried to inject all the available biological knowledge in their model, including the relative levels of expression of the different genes. We point out that this is not possible in a Boolean thresholded automaton. Comparing the weights of different lines in the interaction matrix is absolutely nonsense in this formalism [19] and generates some unsatisfiable

situations. For example, if one consider some of the constraints given by Mendoza et al [1] (see Fig. 4) like a > b > l > 0, then a > 2; however, the weight a plays a role in a regulation composed by 3 elements (on gene 4 = AP1), then after [19], $-2 \le a \le 2$ which is unsatisfiable with the constraint a > 2. Mendoza et al, as probably many biologists, will be tempted to describe in the form of weights the relative strengths of interactions between genes acting on different targets, because it has sense in biological terms (a gene q_1 can be much more sensible to the action of the product of another one g_2 , than another third gene g_3 would be for the products of a fourth one g_4 , and this because of different levels of expression of genes g_2 and g_3 , and because of the different efficiencies of the promoters of g_1 and g_3). This is allowed in Hopfield's formalism only for interactions that concern the same transition function; i.e. the same target gene. In other situations, another, more adapted, formalism must be used.

D. Minimal models

Another point that must be mentioned is the notion of minimal model. It is an evidence that weight or thresholds of a Hopfield-like network can take all the possible signed integer values. All the works written in articles or presented in conferences on such networks do not take an interest in this question of the values of interaction matrix or of the thresholds. Because biologists try to inject in the most realistic form (the closer form to biology) their knowledge in their model, they use values of different importance that can be compared between each others. For the reasons evoked in the previous section, we know that it does not make sense most of the time (every time it concerns comparisons between lines in the interaction matrix). Moreover, even if the values are carefully chosen to make sense (local comparisons only), they can be reduced in such a way the transition functions has the same functioning [19].

Two models are equivalent when they have the same transition function (behaviour). We call the parameters interval the smallest interval $I \subset \mathbb{Z}$ such that all the parameters of the model, *i.e.* weights and thresholds, belong to I. A model O, having I_O as the parameters interval, is called minimal when there is no equivalent model M having I_M as a parameters interval and such that $I_M \subset I_O$ and $I_M \neq I_O$. An example is given in figure 3 to explain the notion of equivalence and minimality. A non-minimal model M_1

$$W = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 1 & 2 & 0 \end{bmatrix} \qquad \qquad \theta = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$$
$$I_{M_1} = [0, 2]$$

The corresponding minimal model M_2



Fig. 3. Minimal AND function in Hopfield's formalism. An AND function can be represented by 3 neurons : the state of the nodes N_1 and N_2 (resp. N_3) represent the inputs (resp. the output). Models M_1 and M_2 are equivalent. Their respective matrices of interaction and thresholds are given. Models M_1 and M_2 have the same transition function but M_2 is minimal because there is no other equivalent model having its parameters interval smaller than I_{M_2} .

A transition function can be described by a unique model, the minimal model. Without loss of generality one can say that a minimal model defines a unique behaviour and reciprocally, a behaviour corresponds to only one minimal model. Then, when one finds a solution to a regulatory network, he must ensure that it is a minimal model or reduce it so as to compare it easily to other minimal models. Moreover it allows, as we will see, to reduce drastically the size of the set of valid instances of models.

III. DETERMINING THE REGULATORY NETWORK OF Arabidopsis thaliana FLOWER MORPHOGENESIS

Using the constraint-based approach described above, we looked for all the valid models of the *Arabidopsis thaliana* flower morphogenesis regulatory network that are in adequacy with the experimental knowledge on the structure and the dynamics given in [1] (Fig. 4). We translated the structural knowledge (relations between weights and values of thresholds) and the dynamical knowledge (set of attractors describing the differentiation behaviours of the flower) into a constraint normal form. We also removed or modified (See Fig. 4) the constraints that did not make sense (like the relation a > b > l > 0 that causes unsatisfiable situations when confronted to the intrinsic constraints of Hopfield's formalism). Dynamical constraints are such that the network, depending on initial conditions converge all to at least 4 physiological tissues

(stationary points), sepals, petals, carpel and stamens, plus a stationary point attractor called 'no flower'. More than 39 millions of models (not yet reduced) were compatible with structural knowledge. After addition of the dynamical knowledge, there were only 3360 non-minimal models. Finally, among them, we found 532 minimal models (these results will be presented in another article [11]). All converge to the experimentally observed stationary points but some showed other stationary points, some showed different cycles ... As expected, Mendoza's do not belong not in them because of the presence in this model of unsatisfiable constraints.



Fig. 4. **Hopfield-like regulatory network of the** *Arabidopsis thaliana* **flower morphogenesis.** (Top left) The scheme of the regulatory networks with its 12 nodes, as modelled by Mendoza et al [1]. (Bottom left) A photograph of some *A. thaliana* flowers. (Right) The list of structural and dynamical constraints as given by [1]. We indicate highlighted in grey the structural constraints we conserved in the minimal model. All dynamical constraints were preserved. All the other relations bring into play weights of arcs targeting different nodes, *i.e.* they are non-sense relations.

Beyond the absolute result concerning the models themselves, this clearly shows the non-uniqueness of models. Biologist must not be too much confident in their reasoning based on an accumulated knowledge and their intuitions. They will be able to find only few models (perhaps one or two different ones) without imagining the huge number of other models that exist. The constraint-based approach allows to reveals the range of the set of different models that do exactly the same thing. The role of the biologist will then be to find new constraints that will divide this set and reject progressively the models until only one or few ones remain.

IV. LETS GO FARTHER : TAXONOMY, OPTIMALITY AND ROBUSTNESS

Many different minimal models of *A. thaliana* flower morphogenesis regulatory network exist (532). What can we do with that ? Which model can we choose and on which criterion ? Are there groups of homogeneous models ? Do the size of the basins of attraction makes sense in biology ?

A first manner to apprehend the set of models is to look at the correlation between the size of the basins of attractions within the set of valid models as shown in Fig. 5. Here, and in the following, we will only concentrate on the stationary points corresponding to the 4 physiological tissues of the flower and the 'no flower' stationary point. One observes first that there exists a strong correlation between these sizes, e.g. petals and sepals are positively correlated (i.e. when the size of the basins of attraction of sepal increases, that of petals increases too), idem for carpel and stamen, but sepals or petals are negatively correlated with carpel or stamen. If one consider these models not anymore as different independent models but as variants of an ideal model, this result makes sense. Stamen and carpel are co-localized and have a common function (reproduction), and that the same for petals and sepals that play a protective role together (lets call them 'co-tissues'). If a variation (a mutation) occurs in the network that changes the size of the basin of attraction of one tissue of the co-tissues, then it makes sense that the variation will have an effect on both tissues of a cotissue.

The other observation made on this graphic is the presence of homogeneous groups of models.



Fig. 5. Correlations between the basins of attraction sizes.

The size of the basins of attraction of the different tissues gives an idea of their robustness against fluctuations of the network state. One could calculate the probabilities of transition from a basin of attraction to another in the presence of noise. The bigger the basins are, the more robust they are because it is probable that the network state formed by change in the state of one gene only gives another network state in the same basin [5]–[8]. Then, the size of an basin of attraction is a good criterion of robustness. We used a score based on the sizes of the basins of attraction of the 5 tissues to classify the models (or morphogenetic landscapes). In figure 6, two examples of classification based on on a different calculus of the scores are given. A work concerning a good manner to calculate such a score would be welcome but the essential is to make appearing the homogeneous groups the figure 5 showed. At the right of figure 6, one can see a coarse grain classification of models into 4 groups.



Fig. 6. Classification of valid models by scoring their robustness. (Bottom) The score is displayed in a growing order. It is calculated (Left) as the sum of the sizes of the 4 basin of attractions of petals, sepals, carpel and stamen, minus that of the no flower attractor (score optimal value = 0), or (Right) as the product of Gaussian functions centered on the 'optimal size' of each of the basin of attractions of the tissues, that of the no flower being centered on 0 (score optimal value = 1). The optimal size is fixed here as the average size of the 4 tissues taken together. (Top) Using the scores, one can classify the models in a barplot graphic showing the cumulated size of the basin of attractions. (Top right) Roughly, 4 distinct groups appear : (A) big sepal basins of attraction, (B) big carpel basins of attraction, (C) all equally-distributed.

Using a principal component analysis in the space of models (or morphogenetic landscapes), we confirm the observation of figure 6. The principal components are linear combinations of the 532 models. By projecting the models into the new spaces formed by the principal components, we are able to separate them into 4 homogeneous groups.

Many other methods of taxonomic classification are possible like those using the Hamming distances between the states belonging to the basins of attraction of each model. Far for being a simple curiosity, providing a taxonomy of models is a powerful tool for dividing the set of valid models in a Constraint-based approach. If it exists, and if it has any sense to compare reality to formal models, *the model* that represents the real system belongs to one of these groups. Finding criteria to separate these groups, criteria that can be translated in the form of constraints, is then a very efficient manner to find useful constraints.



Fig. 7. **Principal component analysis in the model space.** A principal component analysis is realized on the model space (covariance matrix is 532x532 size). Here, the principal components represent new models (or new morphogenetic landscapes) that are combinations of 532 different models. (Top left) Four principal components only (mostly 3) are significantly different to 0. This means that reasonably 4 to 5 homogeneous groups of models can be identified. (Others) Projections of morphogenetic landscapes (models) on principal component (new morphogenetic landscapes or models) spaces.

The importance of 'less important' attractors such as limit cycles has not been discussed yet. We think that one cannot ignore them anymore. They are network states that the system can reach depending on initial conditions. Real systems are always fluctuating and sometimes the state of a system can jump from one attractor to another one. Of course it will be easier if the source basin is smaller than the target basin. Such source basin could be limit cycles or other fixed points having small basins of attraction. They could correspond to stem cells of different differentiation degrees, the differenciation process occurring due to fluctuations of the network states of these stem basins. A future study should integrate them as a criterion for classifying the models and provide an analysis of the differentiation landscape from the size of the basins. The optimal size of the basins is to be determined from biological experiments on the robustness of the development of the different tissues under exogenous perturbations : if the differentiation or the homeostasis of a tissue is robust to perturbation, then its basin of attraction must be large.

V. SUMMARY AND PERSPECTIVES

Along the first part of this paper we illustrate the different traps encountered when conceiving models of biological regulatory networks. We introduced a new approach (the constraint approach) very few considered in systems biology but that having a great potential in the – logical – inference of models from experimental knowledge. The entire knowledge is considered from the beginning and have the same importance. The notion of minimal model is also introduced to complete our inference approach in Hopfield-like networks. These points will be developed in two separate articles.

We are now interested in using this approach to detect eventual incoherences in the experimental knowledge and to infer new knowledge concerning the relation structuredynamics such as the association of recurrent motifs such as positive or negative circuits having particular dynamics. We indeed think that it is possible to identify limit cycles being able to synchronize easily through the structural motifs they produce. A limit cycle that easily synchronize is a limit cycle that, after a perturbation, relaxes following a particular phase whatever is its initial phase. The perturbation is elementary, *i.e.* it only affects one component of the network. At the scale of a population of networks, this results in a global synchrony favouring a collective behaviour. We are also currently integrating in the form of new constraints the possible update modes of the network (parallel, sequential, bloc-parallel).

All these new developments will allow us imagining 'intelligent' laboratory books : the biologist will feed his laboratory book in experimental knowledge such as a series of experiments of activation or inhibition of genes and the observed behaviours obtained for example from micro-array records. As the records go along, the data will be automatically translated in the form of constraints injected in a satisfiability solver that will infer the set of valid models in agreement with the experimental knowledge.

ACKNOWLEDGMENT

This work was supported by the Virtual Physiological Human Network of Excellence of the European Community (VPH-NoE).

REFERENCES

- L. Mendoza and E. Alvarez-Buylla, "Dynamics of the genetic regulatory network: Arabidopsis thaliana flower morphogenesis," *Journal of Theoretical Biology*, vol. 193, no. 2, pp. 307–319, 1998.
- [2] F. Corblin, S. Tripodi, E. Fanchon, D. Ropers, and L. Trilling, "A declarative constraint-based method for analysing discrete genetic regulatory networks," *Biosystems*, vol. 98, pp. 91–104, 2009.
- [3] F. Corblin, E. Fanchon, and L. Trilling, "Applications of a formal approach to decipher discrete genetic networks," *BMC Bioinformatics*, vol. 11:385, 2010.
- [4] R. Thomas, "On the relation between the logical structure of systems and their ability to generate multiple steady states or sustained oscillations," *Springer series in synergetics*, vol. 9, pp. 180–193, 1980.
- [5] J. Demongeot, A. Elena, and S. Sené, "Robustness in neural and genetic networks," Acta Biotheor., vol. 56, pp. 27–49, 2008.
- [6] A. Elena, H. Ben-Amor, N. Glade, and J. Demongeot, "Motifs in regulatory networks and their structural robustness," in *IEEE Transactions on Information Technology in Biomedicine*, 8th IEEE International Conference on BioInformatics and BioEngineering, IEEE Proceedings, Piscataway, Ed., 2008, pp. 234–242.
- [7] H. Ben-Amor, J. Demongeot, and S. Sené, "Structural sensitivity of neural and genetic networks," in LNCS 5317 Proceedings of 7th Mexican International Conference on Artificial Intelligence, 2008 (MICAI'08), Springer, Ed., 2008, pp. 973–986.

- [8] H. Ben-Amor, S. Cadau, A. Elena, D. Dhouailly, and J. Demongeot, "Regulatory networks analysis: Robustness in biological regulatory networks," in *IEEE Proceedings of International Conference on Advanced Information Networking and Applications Workshops*, 2009 (AINA'09), IEEE, Ed., 2009, pp. 924–929.
- [9] H. Ben-Amor, "Applications of a constraint-based approach to build and re-examine biological regulatory and formal neural networks," in 3rd Franco-Japanese Symposium on Knowledge Discovery in Systems Biology, Corsica, France, 2009.
- [10] —, "Formal methods for biological regulatory networks," in 3rd international conference of the SFBT, Institut Pasteur of Tunis, Tunisia, 2010.
- [11] H. Ben-Amor, F. Corblin, E. Fanchon, A. Elena, L. Trilling, J. Demongeot, and N. Glade, "Formal methods for hopfield-like networks," 2011, submitted.
- [12] K. Apt, Principles of Constraint Programming. Cambridge University Press, 2003.
- [13] M. Carlsson, G. Ottosson, and B. Carlson, "An open-ended finite domain constraint solver," in *Proc. Programming Languages: Implementations, Logics, and Programs*, 1997.
- [14] N. Eén and N. Sörensson, "An extensible SAT-solver," in SAT'2003 Theory and Applications of Satisfiability Testing, LNCS 2919, 2004.
- [15] N. Eén and A. Biere, "Effective preprocessing in SAT through variable and clause elimination," in SAT'2005 – Theory and Applications of Satisfiability Testing, LNCS 3569, 2005.
- [16] F. Corblin, L. Bordeaux, E. Fanchon, Y. Hamadi, and L. Trilling, "Connections and integration with SAT solvers: A survey and a case study in computational biology," in *Hybrid Optimization: the 10 YEARs* of CPAIOR. Springer, 2009, accepted.
- [17] S. Kauffman, "Metabolic stability and epigenesis in randomly constructed genetic nets," J. Theor. Biol., vol. 22, pp. 437–467, 1969.
- [18] J. J. Hopfield, "Neural networks and physical systems with emergent collective computational abilities," *Proc. Natl. Acad. Sci. USA*, vol. 79, pp. 2554–2558, 1982.
- [19] A. Elena, "Robustesse des réseaux d'automates booléens à seuil aux modes d'itération. application à la modélisation des réseaux de régulation génétique," Ph.D. dissertation, 2009.