

# A rate-distortion theory for gene regulatory networks and its application to logic gate consistency

Giuseppe Facchetti, Giovanni Iacono, Giovanna De Palo and Claudio Altafini\*

SISSA, Int. School for Advanced Studies, via Bonomea 265, 34136 Trieste, Italy

Associate Editor: Alfonso Valencia

## ABSTRACT

**Motivation:** A gene regulatory network in which the modes (activation/inhibition) of the transcriptional regulations are known and in which gene expression assumes boolean values can be treated as a system of linear equations over a binary field, i.e. as a constraint satisfaction problem for an information code.

**Results:** For currently available gene networks, we show in this article that the distortion associated with the corresponding information code is much lower than expected from null models, and that it is close to (when not lower than) the Shannon bound determined by the rate-distortion theorem. This corresponds to saying that the distribution of regulatory modes is highly atypical in the networks, and that this atypicality greatly helps in avoiding contradictory transcriptional actions.

Choosing a boolean formalism to represent the gene networks, we also show how to formulate criteria for the selection of gates that maximize the compatibility with the empirical information available on the transcriptional regulatory modes. Proceeding in this way, we obtain in particular that non-canalizing gates are upper-bounded by the distortion, and hence that the boolean gene networks are more canalizing than expected from null models.

**Contact:** altafini@sissa.it.

**Supplementary information:** Supplementary data are available at *Bioinformatics* online.

Received on October 11, 2012; revised on January 16, 2013; accepted on March 1, 2013

## 1 INTRODUCTION

A gene regulatory network consists of a set of transcription factors (TF) regulating the expression of the genes of a given genome. As our knowledge of these regulatory mechanisms grows steadily, the problem of understanding the principles behind their organization and their functioning is becoming a crucial question of Systems Biology, see Bonneau, 2008; Cosentino Lagomarsino *et al.*, 2007; Mangan and Alon, 2003; Materna and Davidson, 2007; Shmulevich *et al.*, 2002. While the current literature is focused mainly on the topological aspects of a gene regulatory network, the perspective that we take in this article is to study whether the *transcriptional regulatory modes* (i.e. the *activation/inhibition role* of a TF on its target genes) of the gene networks currently available are distributed randomly on the given network or are organized according to some

criterion, and how this organization can affect the dynamics of the gene network.

To formalize this question and to try to understand the rules of this design, we use tools from information theory (Cover and Thomas, 2006; Mezard and Montanari, 2009) and in particular we treat a gene network as a ‘code’ for which the signs of the regulatory actions constitute a particular ‘source word’, which can be compared with the typical words generated by a corresponding probabilistic model (further details on the terminology in section 2). As a comparison criterion, we use the ‘level of coherence’ of the regulatory actions along the network. Two regulatory orders emanating from a TF and acting on the same target gene (possibly through intermediate genes) are considered coherent when they induce the same behavior on the target gene (i.e. they both induce activation or repression); they are considered incoherent when they induce conflicting behaviors. As can be easily deduced from simple examples like Feed-Forward Loops (FFL) (Mangan and Alon, 2003), in a gene regulatory network, incoherence is associated with negative (undirected) cycles on the signed graph having the genes as nodes, the regulations as edges and the modes of the regulations as signs of the edges, see Iacono *et al.*, 2010. An undirected cycle is negative when it contains an odd number of inhibitions. Provided we associate binary values to the expression of the genes, the main feature of a negative undirected cycle is that no choice of expression can satisfy all constraints imposed by the regulations. This satisfiability can be tested by formulating the ‘compatibility’ as a system of linear algebraic equations over a binary field. In combinatorial optimization, such problems are well known to be equivalent to constraint satisfaction problems of exclusive OR (XOR) type, see Mezard and Montanari, 2009.

If the gene regulatory network is a code, and the signs of the regulations specify a codeword, then the problem can be studied as a (*lossy*) *source compression problem*, see Ciliberti and Mézard, 2005; Wainwright *et al.*, 2010. In this framework, in particular, the level of coherence of the transcriptional regulations can be rigorously computed as the *distortion* introduced by the source compression problem. Computing this distortion is a hard problem. In the constraint satisfaction literature, it is sometimes referred to as Maximum XOR Satisfiability (MAX-XORSAT) problem, and consists in computing the binary gene expression assignment that maximizes the number of satisfied (SAT) linear equations at steady state, see Correale *et al.*, 2006b; Cosentino Lagomarsino *et al.*, 2005. For the gene networks currently available, the distortion can be quantified with sufficient precision and can be compared with two important quantities: (i) the average distortion of a typical word associated

\*To whom correspondence should be addressed.

to the same ‘code’ (i.e. a gene network with the same topology but signs reshuffled); (ii) the best average distortion achievable by any gene network with the same ratio of genes/regulations. The latter value corresponds to the Shannon bound provided by the rate-distortion theorem (Cover and Thomas, 2006).

We show in the article that the distortion of the currently available gene networks is much lower than the one of a typical sequence of the same code, and that it is comparable with (when not better than) the Shannon bound. This atypicality implies that in our gene networks the signs of the transcriptional regulations are highly organized and far from random. In particular, the origin of the low distortion can be traced in the scarcity of dual-mode TF, i.e. of TF acting both as activators and as repressors. Our calculation suggests a practical reason for such an organization: single-mode TF lower the distortion and hence help in avoiding conflictual transcriptional orders in which different TF induce contradictory actions on a downstream gene.

In a gene network represented as a signed graph, the regulatory actions represent two-body interaction terms (i.e. the value of a gene acting as a TF, multiplied by that of its target gene, with the sign of the corresponding edge) or, in system-theoretic language, Single-Input Single-Output (SISO) regulations. When multiple genes act simultaneously as TF on a target gene, the overall regulation can be described as a superposition of these SISO terms. An alternative assumption, common in the gene network literature (Balleza *et al.*, 2008; Buchler *et al.*, 2003; Correale *et al.*, 2006b; Kauffman, 1993; Kauffman *et al.*, 2004; Shmulevich *et al.*, 2005; Silva-Rocha and de Lorenzo, 2008; van Hijum *et al.*, 2009), is that the transcription of a gene having multiple regulators is decided by some logical combination of the inputs. This corresponds to replacing superpositions of two-body terms with a single multibody term (MISO: Multi-Input Single-Output action). One of the major drawbacks of the boolean networks obtained in this way is the arbitrariness in the choice of the gates, due to the lack of systematic methods for gates disambiguation based on experimental evidence (see e.g. Bonneau, 2008; Chowdhury *et al.*, 2010; Kim *et al.*, 2007; Lau *et al.*, 2007; Shmulevich *et al.*, 2002; Silva-Rocha and de Lorenzo, 2008; Zou, 2010 for a few attempts in this direction, mostly based on inference from microarrays).

The approach that we take in this article is to try to identify classes of boolean gates such that the steady-state MISO predictions for the regulation overlap as much as possible with the corresponding SISO predictions. In terms of perturbative expansions, this corresponds to asking that the two-body projections of the complex multibody terms be ‘coherent’ with the corresponding two-body terms, whenever the latter are not ambiguous. The *rationale* behind this choice is that the only type of information currently available in large-scale gene networks is precisely compendia of SISO regulatory signs. We will show in the article that under this assumption, a natural choice is to associate canalizing gates (i.e. AND-OR type of gates, see Kauffman, 1993) to positive undirected cycles and (possibly) non-canalizing gates to negative undirected cycles. The property of low distortion of the gene networks then reflects in a low amount of non-canalizing gates in the boolean formulations of the gene networks with respect to null models.

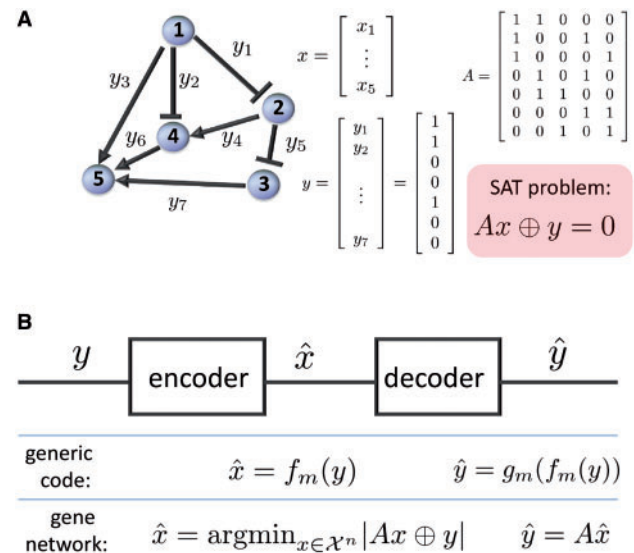
Looking at the corresponding dynamics, we observe that, in spite of the higher canalization, our gene networks seem to be

more sensitive to perturbations than the corresponding null models. This appears to be due to the near-acyclic and massively parallel feed-forward architecture of the networks, in which non-canalizing gates have the effect of breaking the steady-state global symmetry of the basic input-output motifs.

## 2 METHODS

**Gene regulatory networks and SISO constraint satisfaction problems.** Consider a gene regulatory network composed of  $n$  nodes  $\mathbf{x} = [x_1 \dots x_n]^T$  representing the genes and  $m$  directed edges  $\mathbf{y} = [y_1 \dots y_m]^T$  representing regulatory actions of activation/inhibition of a gene on another gene. Assume both  $\mathbf{x}$  and  $\mathbf{y}$  are represented in boolean terms,  $x_i \in \mathcal{X} = \{0, 1\} = \{\text{‘low’}, \text{‘high’}\}$ , and  $y_i \in \mathcal{Y} = \{0, 1\}$ , where we use the convention that 0 stands for activation (i.e. ‘+’) and 1 for inhibition (i.e. ‘−’), see example in Figure 1A. Then both  $\mathcal{Y}$  and  $\mathcal{X}$  can be identified with  $\mathbb{Z}_2$ , the Galois field with two elements endowed with the operation  $\oplus$  (addition mod2).

In this section, we are interested in formulating a description of the fixed points of the signed graph (of nodes  $\mathbf{x}$  and edges  $\mathbf{y}$ ) representing the gene regulatory network. We shall for now assume that each regulatory action is SISO and that multiple regulations acting on the same gene happen simultaneously, independently and in parallel. Consider for example a single regulatory action  $y_1$  between the two genes  $x_1$  and  $x_2$  (i.e.  $m = 1$ ,  $n = 2$ ). For any value  $y_1 \in \mathcal{Y}$ , it is always possible to find (at least) a combination of  $x_1, x_2 \in \mathcal{X}$ , which is ‘compatible’ with the expected action described by the sign  $y_1$ . If, for example,  $y_1 = 0$  (i.e. activation), then the regulation is SAT when  $x_1, x_2$



**Fig. 1.** Gene networks and rate-distortion theory. (A) A toy example of signed gene regulatory network, and its formulation as a SAT problem. (B) Rate-distortion scheme. In the regime  $n < m$ , the encoding/decoding scheme is normally referred to as a *lossy source compression problem* (Cover and Thomas, 2006), as a length- $n$  sequence  $\hat{\mathbf{x}}$  is used to represent a length- $m$  word  $\mathbf{y}$ . The distortion corresponds to the relative Hamming distance  $d(\mathbf{y}, \hat{\mathbf{y}})/m$ . In our gene networks the distortion of a network is a measure of potential ‘conflicts’ (or contradictory orders) in the gene regulatory program of an organism

assume the same value, opposite values when  $y_1 = 1$ , see Supplementary Figure S1. Using the formalism of constraint satisfaction problems (Mezard *et al.*, 2003; Mezard and Montanari, 2009), we can rewrite this ‘compatibility’ condition as a linear equation over a binary field:  $x_1 \oplus x_2 = y_1$ . This exclusive-OR satisfiability (XORSAT) problem can also be written as  $x_1 \oplus x_2 \oplus y_1 = 0$ .

Let us make use of an  $m \times n$  matrix  $A$  to describe the topology of the gene network. Each row of  $A$  identifies an edge, and has two non-zero entries (equal to 1) in correspondence of the two genes involved in the regulatory action (the order is irrelevant for our purposes), see Figure 1A for an example. Finding gene expression assignments  $\mathbf{x} \in \mathcal{X}^n$  compatible with all the regulatory signs of a given  $\mathbf{y} \in \mathcal{Y}^m$  means solving a linear system over the  $\mathbb{Z}_2$  field:

$$A\mathbf{x} \oplus \mathbf{y} = 0. \quad (1)$$

The XORSAT problem (1) may or may not have a solution depending on  $\mathbf{y}$ . When (1) has no solution (it is said UNSAT in this case), then one can look for the  $\mathbf{x} \in \mathcal{X}^n$  that solves most of the  $m$  constraints of (1). This problem is called MAX-XORSAT, see Mezard and Montanari, 2009. Denoting  $d(\mathbf{y}, A\mathbf{x})$  the Hamming distance between  $\mathbf{y}$  and  $A\mathbf{x}$ , then solving the MAX-XORSAT for a given  $\mathbf{y}$  means finding  $\hat{\mathbf{x}}$  that minimizes  $d(\mathbf{y}, A\mathbf{x})$ :

$$\hat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}^n} d(\mathbf{y}, A\mathbf{x}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}^n} |A\mathbf{x} \oplus \mathbf{y}|. \quad (2)$$

The relative Hamming distance

$$D_{A,\mathbf{y}} = \frac{1}{m} \min_{\mathbf{x} \in \mathcal{X}^n} d(\mathbf{y}, A\mathbf{x}) \quad (3)$$

is called the distortion of the ‘word’  $\mathbf{y}$  associated to the ‘code’ (1)-(2). We shall call average distortion of the gene regulatory network associated to (1)-(2) the expectation of (3) over the entire alphabet  $\mathcal{Y}^m$ :

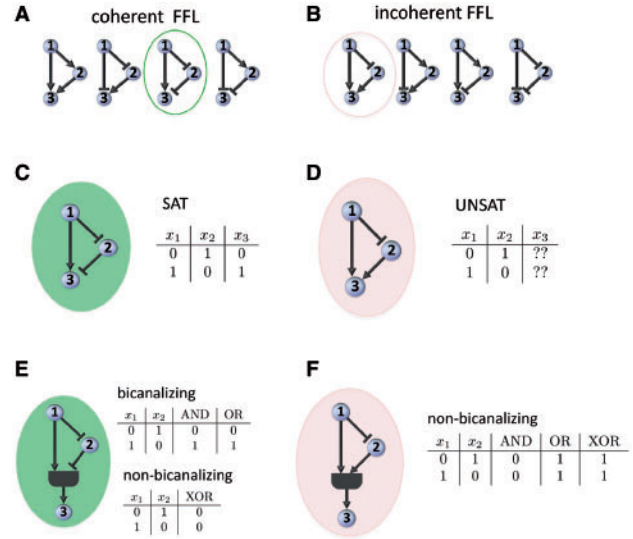
$$D_A = \frac{1}{m} \sum_{\mathbf{y} \in \mathcal{Y}^m} p(\mathbf{y}) d(\mathbf{y}, A\hat{\mathbf{x}}) \quad (4)$$

where, for each  $\mathbf{y}$ ,  $\hat{\mathbf{x}}$  solves (2). Denoting  $\hat{\mathbf{y}} = A\hat{\mathbf{x}}$  the estimate of  $\mathbf{y}$  obtained through (2) (see Fig. 1B), then  $D_A = 0$  if and only if  $\hat{\mathbf{y}} = \mathbf{y} \forall \mathbf{y} \in \mathcal{Y}^m$ .

**Example: FFL.** Let us consider as an example the FFL of Figure 2. In a FFL,  $n = m = 3$ , and hence both alphabets  $\mathcal{X}^n$  and  $\mathcal{Y}^m$  consist of eight words, and we may ask if for each choice of  $\mathbf{y} \in \mathcal{Y}^3$  there is always an  $\mathbf{x} \in \mathcal{X}^3$  such that (1) is SAT. In this case, the connectivity matrix  $A$  is the following

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 1 \end{bmatrix} \quad (5)$$

and it is straightforward to verify that only in four of the eight choices of  $\mathbf{y}$  (1) is SAT. These cases (Fig. 2A) are known in the literature as *coherent* FFL, whereas the 4 UNSAT cases (Fig. 2B) are called *incoherent* in Mangan and Alon, 2003. In a coherent FFL  $d(\mathbf{y}, A\hat{\mathbf{x}}) = 0$ , whereas in an incoherent FFL  $d(\mathbf{y}, A\hat{\mathbf{x}}) = 1$ , i.e.  $D_{A,\mathbf{y}} = 1/3$ . See Supplementary Information and Supplementary Table S1 for a count of the FFL motifs in the gene networks of Table 1.



**Fig. 2.** FFL motifs and their steady-state values. The FFL is the simplest motif forming an undirected cycle in a graph. Of the eight possible FFL, four are coherent (A) and four incoherent (B); see Mangan and Alon, 2003. As SISO systems, at steady state the coherent FFL are SAT (C), while the incoherent FFL are UNSAT (D). Therefore, coherent FFL admit bicanalizing steady-state behavior if the boolean logic gate is itself canalizing (E). No bicanalizing steady state is possible for incoherent FFL (F)

**Table 1.** Gene regulatory networks and their distortion

| Network             | $n$  | $m$  | $R$   | $q$   | $D_{A,\mathbf{y}}^{\text{emp}}$ | $D_A$            |
|---------------------|------|------|-------|-------|---------------------------------|------------------|
| <i>E.coli</i>       | 1461 | 3220 | 0.454 | 0.416 | [0.1134, 0.1152]                | [0.1767, 0.2043] |
| <i>S.cerevisiae</i> | 690  | 1082 | 0.638 | 0.204 | [0.0379, 0.0379]                | [0.1077, 0.1091] |
| <i>B.subtilis</i>   | 918  | 1324 | 0.693 | 0.256 | [0.0536, 0.0536]                | [0.1040, 0.1043] |

$n$  and  $m$  are the number of nodes and edges of the directed graph representing the gene network;  $R$  is the ratio nodes/edges;  $q$  is the fraction of negative edges;  $D_{A,\mathbf{y}}^{\text{emp}}$  and  $D_A$  are respectively the distortion of the true edge sign assignment and of the null models, see (3) and (4) (lower and an upper bound are provided; see Iacono and Altafini, 2010 for the details).

**Computing distortion for a gene regulatory network.** In information theory, the process of obtaining  $\hat{\mathbf{y}}$  from  $\mathbf{y}$  through (1)-(2) in the regime  $n < m$  is called a *lossy source compression problem*, see Figure 1 and the Supplementary Information. In particular, solving the combinatorial optimization problem (2) corresponds to encoding the length- $m$  binary sequence  $\mathbf{y}$  into the length- $n$  binary sequence  $\hat{\mathbf{x}}$  through the ‘channel’ given by the connectivity matrix  $A$ . Decoding then corresponds to constructing an estimate  $\hat{\mathbf{y}}$  of  $\mathbf{y}$  out of  $\hat{\mathbf{x}}$  (Fig. 1).

**Boolean networks and steady-state consistency.** A Boolean network or, more generally, a discrete dynamical system on  $\mathbb{Z}_2$  can be written as

$$\mathbf{x}(t+1) = \phi(\mathbf{x}(t)), \quad \mathbf{x}(t+1), \mathbf{x}(t) \in \mathcal{X}^n \quad (6)$$

where  $\phi = (\phi_1, \dots, \phi_n) : \mathcal{X}^n \rightarrow \mathcal{X}^n$  is the boolean state update map. If  $k_i$  is the in-degree of a fan-in node  $x_i$ , then  $k_i > 1$ , and the



value  $x_i(t+1)$  is decided univocally from the values of its precursory inputs, denoted  $x_{i_1}, \dots, x_{i_{k_i}}$ , according to a logical function:  $x_i(t+1) = \phi_i(x_{i_1}(t), \dots, x_{i_{k_i}}(t))$ .

As in the previous section, we aim at studying fixed points of the dynamics, where a fixed point for the system (6) is  $\mathbf{x} \in \mathcal{X}^n$  such that  $\mathbf{x} = \phi(\mathbf{x})$ . Setting the problem as in (1) corresponds to saying that all transcriptional actions happen independently and in parallel at each fan-in node. The concept is natural in continuous-time models, where superposition of the effects can take place, but it is ambiguous for discrete dynamics in  $\mathbb{Z}_2$ . Placing a function  $\phi_i$  on each fan-in node resolves all ambiguities and guarantees existence and uniqueness of the solution of (6).

In particular, we aim at formulating a constraint satisfaction problem containing all our information about transcriptional regulation and at exploring what functions  $\phi$  are compatible with it. In the boolean networks literature, this approach is referred to as *consistency problem*, see Shmulevich *et al.*, 2002.

A  $k$ -input boolean logic gives rise to  $2^k$  possible boolean gates. Little is known empirically on the actual update rules at fan-in nodes in real genome-wide transcriptional networks, except for some case-by-case analysis (Albert and Othmer, 2003; Harris *et al.*, 2002; Materna and Davidson, 2007) and some in-depth studies of specific transcriptional units (Buchler *et al.*, 2003; Setty *et al.*, 2003). Especially when  $k_i > 2$ , it is not even an obvious assumption that the  $k_i$  inputs necessarily have to form a single coordinated transcriptional logic unit, rather than multiple concomitant, independent subunits, perhaps active in different environmental conditions.

In random boolean networks (Kauffman, 1993),  $k_i$  and  $\phi_i(\cdot)$  of each node are chosen randomly. In our gene networks, such random choice would disregard the information we have available about the topology and about the regulatory action (activation/inhibition) from each precursor  $x_j$  to  $x_i$ . This information can be used to prune a large part of the possible gates. If adopting the topology of the true networks means fixing the  $k_i$  precursors of each node (Correale *et al.*, 2006a; Kauffman *et al.*, 2003), imposing in the boolean logic the compatibility with the direct SISO regulations requires assuming that all the functions  $\phi_i(\cdot)$  are unate in each of the arguments (Sontag *et al.*, 2008). This corresponds to saying that multiple appearances of the same  $x_j$  in one of the  $\phi_i(\cdot)$  must be characterized by the same sign. In turn,  $\phi$  unate implies that it is possible to fix the sign of each input of the gate equal to the value of the direct SISO transcriptional regulation. Once this compatibility condition is taken into account, further negations on the gates should be avoided, thereby reducing drastically the number of possible choices.

Formally, we have to solve the following steady-state consistency problem, see Correale *et al.*, 2006b; Cosentino Lagomarsino *et al.*, 2005

Find  $\phi_i : \mathcal{X}^{k_i} \rightarrow \mathcal{X}$ ,  $i = 1, \dots, n$ , s.t.

$$x_i \oplus \phi_i(*, \dots, *, x_j, *, \dots, *) = y_\ell \quad (7)$$

$$\forall j, \ell \text{ s.t. } \exists A_{\ell,i} = A_{\ell,j} = 1.$$

Because  $\phi_i$  is unate, the edge sign  $y_\ell$  can be explicitly attributed to  $x_j$  by replacing (7) with

$$x_i \oplus \phi_i(*, \dots, *, y_\ell \oplus x_j, *, \dots, *) = 0 \quad (8)$$

$$\forall j, \ell \text{ s.t. } \exists A_{\ell,i} = A_{\ell,j} = 1.$$

The argument above allows to restrict the search to  $\phi_i$ , which do not contain any negation neither on the output nor in the inputs (except for the  $y_\ell$ ).

Let us consider for example the 16 gates of a two-input  $\phi_i$ , see Table S3. These gates can be classified into four groups (Correale *et al.*, 2006b):

1. Constant functions (2 gates);
2. Projections (4 gates);
3. XOR class (2 gates);
4. AND-OR class (8 gates).

Excluding the trivial functions in class 1 and 2, and excluding the gates containing negations [in (8) the  $y_\ell$  are not considered part of the function  $\phi_i$ ] leaves us with only three gates: AND, OR and XOR. Qualitatively, the AND-OR gates differ from the XOR gate in the sense that the former are canalizing in one of the inputs, while XOR is not. A gate is said canalizing if at least one of its inputs has a value that alone decides the output of the gate, regardless of the values of the other variables. The canalizing value is 0 for AND and 1 for OR. For XOR instead no input is canalizing. For random boolean networks, it is well known that canalizing functions represent mechanisms associated with ordered dynamical behavior (as opposed to chaotic behavior), see Kauffman, 1993; Kauffman *et al.*, 2004. In the compendium of transcriptional mechanisms analyzed in Harris *et al.*, 2002, for example, canalizing functions are neatly overrepresented.

In the following we show that further constraints on the gates, and in particular a criterion to discern canalizing from (potentially) non-canalizing gates, can be deduced by imposing that on certain positive undirected cycles, the ‘multibody’ steady-state behavior must be compatible with the direct SISO regulations. Let us consider Feed-Forward Circuits (FFC), i.e. FFL-like subgraphs with a single root (i.e. node with zero in-degree) and a single sink (i.e. node with zero out-degree) but branches of any length. For these FFC, positivity of the undirected cycle amounts to having the same signature along the two-directed paths connecting the single root to the single sink. In terms of constraint satisfaction conditions (1)-(2), the corresponding subproblems must be SAT, and this means that a steady state is uniquely determined along the cycle by the value taken by the root. This means also that at steady state the root-sink boolean relationship relies on at least two independent transcriptional routes. For positive FFC, in particular, the two relationships are coherent in the sense that they agree on the root-sink steady-state value. Root-sink coherent redundancies like this can be taken as criteria for discerning the type of gate to be placed at multiinput nodes involved in FFC motifs.

*Example: consistency for boolean FFL.* An isolated FFL with its single root, single sink and two-input logic gate is an example of undirected cycle splittable into two-directed paths. If we consider a coherent FFL as in Figure 2C, then at a fixed point the direct SISO constraints are  $x_2 = \bar{x}_1$  and  $x_3 = x_1$  (after the double negation). Replacing the direct SISO actions on  $x_3$  with a two-input gate  $\phi_3(x_1, x_2)$  then (8) becomes:

$$x_3 \oplus \phi_3(x_1, \bar{x}_2) = x_3 \oplus \phi_3(x_1, x_1) = 0.$$

It is straightforward to verify that  $\phi_3 = \{\text{AND}, \text{OR}\}$  are both compatible with the root-sink actions of each directed path, while XOR is not, see Figure 2E.

If instead we consider the incoherent FFL in Figure 2D, the SISO steady-state compatibility conditions reduce to  $x_2 = \bar{x}_1$  because the value of  $x_3$  is ambiguous and (8) is

$$x_3 \oplus \phi_3(x_1, x_2) = x_3 \oplus \phi_3(x_1, \bar{x}_1) = 0.$$

Depending on the choice of  $\phi_3$  (see Fig. 2F), at steady state we have  $\forall x_1$

$$x_3 = \begin{cases} 0 & \text{if } \phi_3 = \text{AND} \\ 1 & \text{if } \phi_3 = \text{OR, XOR} \end{cases}.$$

As expected, for incoherent FFL, no gate can be compatible with the root-sink action on both branches simultaneously. Each of the basic gates opts for one of the choices. It is consequently impossible to discriminate between AND, OR or XOR with the SISO information available. Notice how in correspondence of each choice of  $\phi_3$  the steady-state value of  $x_3$  is blocked for all possible values of  $x_1$ . No steady-state block appears instead in the coherent FFL with canalizing gates, where in addition at steady state AND and OR are indistinguishable and the output is completely determined by  $x_1$ . This can be rephrased in terms of global symmetry of the corresponding SAT subproblem. Steady-state blocks, like in the incoherent FFL, break the global symmetry on the corresponding motif: flipping the value of the root does not result in a steady-state flip at the sink. We shall denote FFL motifs (and more generally FFC motifs) respecting the global input-output steady-state symmetry *bicanalizing*.

### 3 RESULTS/DISCUSSION

*The distortion of the gene networks of Table 1.* Consider the three gene regulatory networks of Table 1. In a source compression scheme, the difficult task is the encoding part [i.e. solving the MAX-XORSAT problem (2)]. For two-body constraints like (2), it is well known to be equivalent to solving a maximum cut problem or to computing the ground state of an Ising spin glass with bimodal bonds, see Mezard and Montanari, 2009.

For the three gene-regulatory networks of Table 1, efficient heuristics are discussed in Iacono *et al.*, 2010; Iacono and Altafini, 2010, see also Supplementary Information for a recap. In the following, we denote  $D_{A,y}^{\text{emp}}$  the distortion in correspondence of the true ('empirical') edge sign assignments  $y$  of a given network. The values reported in Table 1 are computed from the 'distance to monotonicity' of Iacono and Altafini, 2010, using (3). The upper and lower bounds computed for  $D_{A,y}^{\text{emp}}$  are fairly tight, see Table 1.

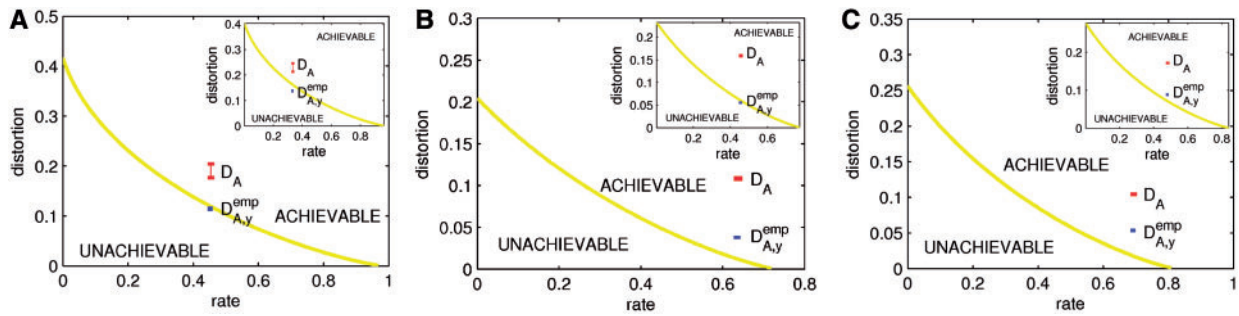
Needless to say, it is impossible to evaluate exactly quantities such as (4), which are computed exhaustively over  $\mathcal{Y}^m$ . To estimate  $D_A$ , we can repeat the same optimization as in (2) on a sufficiently large number of null models, obtained drawing length- $m$  independent identically distributed edge sign words from a Bernoulli distribution  $\mathcal{B}(q)$ , where  $q = m^-/m$  is the fraction of negative edge signs of the original  $y$ ; see Table 1 and the Supplementary Information for more details. As can be seen on Table 1,  $D_{A,y}^{\text{emp}} < D_A$  for all three gene networks, meaning that the

true gene networks have less distortion than the corresponding null models (coherently with the results reported in Iacono *et al.*, 2010; Iacono and Altafini, 2010). Notice that only nodes involved in undirected cycles contribute to the distortion. In particular, then, intending the networks as undirected graphs and restricting to the bicomponents (i.e. dropping the nodes/edges not involved in undirected cycles) means changing the values of  $n$  and  $m$ , and hence of  $D_{A,y}^{\text{emp}}$  and of  $D_A$ , see Supplementary Table S2.

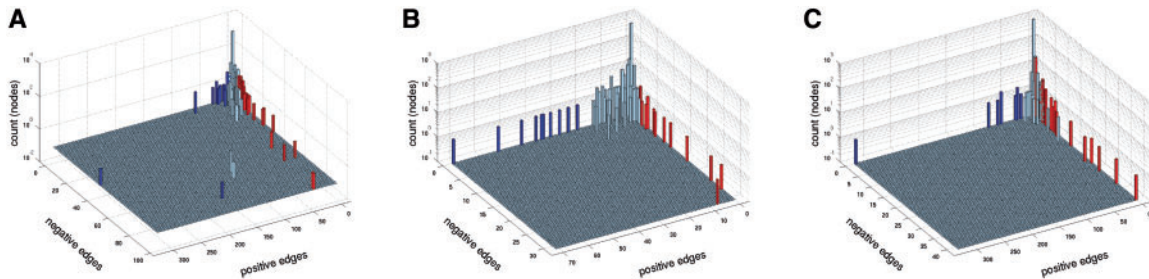
Both  $D_{A,y}^{\text{emp}}$  and  $D_A$  are obtained in correspondence of the given topology, described by the connectivity matrix  $A$ . If we allow also the topology to vary, then we can use the rate-distortion theorem of information theory (Cover and Thomas, 2006) to determine the admissible region for the distortion  $D$  for all possible compression rates  $R = n/m$  in correspondence of a  $\mathcal{B}(q)$  source of edge sign words. As explained in the Supplementary Information, the boundary of such an admissible region represents a Shannon-type bound, and it is achieved in correspondence of a 'best' network topology. For Bernoulli sources, the distortion  $D$  on such bound can be computed explicitly; see (S2) on the Supplementary Information. In Figure 3, the Shannon bound  $D$  is compared with  $D_{A,y}^{\text{emp}}$  and  $D_A$  for the three gene networks of Table 1. As can be observed, the values of  $D_{A,y}^{\text{emp}}$  are close to the corresponding  $D$ . In particular, once we restrict to the bicomponents (insets in Fig. 3A–C),  $D_{A,y}^{\text{emp}} < D$  in two of the three networks. On the contrary,  $D_A \gg D$  in all three cases. The meaning is that in spite of non-optimal topologies, the distortions of our gene regulatory networks (which, from  $D_{A,y}^{\text{emp}} < D_A$ , we know to be much lower than expected) are also at the level expected for a 'best' network topology.

A statistical physics analog of an XORSAT problem is an Ising spin glass (Mezard and Montanari, 2009); see Supplementary Information. In this context, the distortion  $D_{A,y}^{\text{emp}}$  has the interpretation of 'frustration' encoded in the undirected cycles, i.e. of linearly independent undirected cycles having negative sign (meaning an odd number of inhibitions). Our result, therefore, implies that frustration is largely absent in these signed graphs. Notice that direct counts of the basic frustrated/non-frustrated motifs such as the incoherent/coherent FFL are largely inconclusive; see Supplementary Information and Table S1. The true distortion can be computed only genome-wide, and its calculation confirms that indeed conflictual orders are largely avoided.

*Low distortion and single-mode TF.* It is worth mentioning that the origin of the low-distortion of the gene networks lies in the highly skewed distribution of the signs of the actions of the TF. As can be seen in Figure 4, the vast majority of the TF tends to operate in a single-mode fashion on all their target genes. Dual-mode TF are statistically rare with respect to the null models (cumulative binomial test,  $P$ -value  $10^{-2}$ ). See Supplementary Tables S4 and S5–S7 for more details. While this skewness is expected, as the physical interaction mechanisms of an activator and of an inhibitor are normally different, its consequences for the regulation on a genome-wide scale have rarely been assessed, except on small motifs like FFL. Following the arguments of Facchetti *et al.*, 2011, Iacono and Altafini, 2010, and in particular the notion of gauge equivalence discussed therein, it can be shown that such a pattern is responsible for the limited amount of distortion of these networks. On the contrary, the signs on the



**Fig. 3.** Distortion and Shannon bounds. For the three gene networks of Table 1 [(A) *E.coli*; (B) *S.cerevisiae*; (C) *B.subtilis*], the distortion in correspondence of the ‘true’ edge signs  $\mathbf{y}$  ( $D_{A,y}^{\text{emp}}$ ) and in correspondence of the null models ( $D_A$ ) is compared with the rate-distortion theorem (yellow curve); see Supplementary Information. Upper and lower bounds on  $D_{A,y}^{\text{emp}}$  and on  $D_A$  are normally close; see Table 1 (one exception is  $D_A$  for *E.coli*). In the insets, the same quantities are shown for the bicomponents of the networks (where nodes and edges not involved in undirected cycles are dropped; see Supplementary Table S2). In this last case,  $D_{A,y}^{\text{emp}}$  is below the Shannon bound in two of the three networks (and close in the third one). In all cases,  $D_{A,y}^{\text{emp}} < D_A$ , i.e. the distortion is atypical for a  $B(q)$  source of words (i.e. of edge sign assignments)



**Fig. 4.** Single-mode action of the TF. The histograms show the number of positive and negative edges in the out-degree of the TF for the three networks [(A) *E.coli*; (B) *S.cerevisiae*; (C) *B.subtilis*]. The histograms are highly skewed, meaning that the majority of TF have a single mode of action. In particular, the TF significantly single mode (with respect to a cumulative binomial test,  $P$ -value  $10^{-2}$ ) are highlighted in color: blue for activators, red for repressors. See also Supplementary Tables S5–S7 for a list of the corresponding genes. For *E.coli*, the few TF having both positive and negative edges are well-known dual-mode regulators, such as *crp*, *fnr*, *ihf*, *fis*, *arcA* and *narL*. See Supplementary Table S4

incoming edges at the target genes do not show any deviation from the null models. See Supplementary Information and Supplementary Figure S2 for more details.

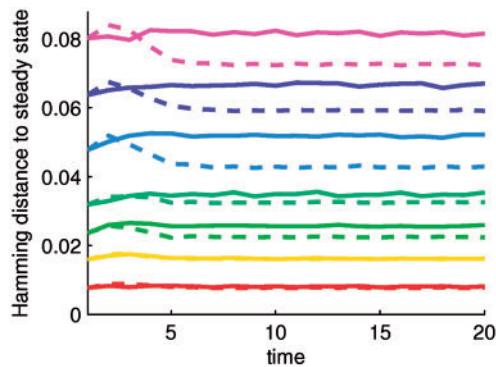
*Distortion as an upper bound on non-canonicalizing functions of boolean networks.* In non-isolated FFL, the bicanalizing properties of boolean logic mentioned in section 2 extend to output nodes that are fan-in of multiple coherent FFL, see Supplementary Figure S3, while in presence of one or more incoherent FFL the situation is more complex, and more difficult to analyze in detail; see e.g. Supplementary Figure S4. Not only the type of gates chosen matters in this case, but also the order in which logical operations are carried out; see table in Supplementary Figure S4. The case in which the node  $x_2$  of an FFL has in-degree greater than 1 is equally complex because the motifs have more than one root. See Supplementary Figures S5–S7 for a few examples. While the nested FFL case of Supplementary Figure S5 is always bicanalizing, in Supplementary Figure S6 even though all undirected cycles are positive, the multiinput motif is bicanalizing only for some combinations of AND, OR. In Supplementary Figure S7, instead, no logic gate can achieve bicanalization. Even on these simple examples, it is possible to observe how the choice of non-canonicalizing

gates along positive undirected cycles can lead to paradoxical input-output steady-state behaviors. Notice for example in Supplementary Figure S6 how XOR can induce ‘wrong’ bicanalizations, i.e. bicanalization with respect to the wrong input.

For these reasons, and to maximize the ‘overlap’ between SISO transcriptional actions and MISO actions induced by the boolean logic, it is reasonable to assume that positive undirected cycles are preferentially endowed with canalizing gates, while in negative undirected cycles no priority can be imposed on the type of gates. This choice corresponds to solving the steady-state consistency problem while maximizing the satisfiability to both (1) and (8). Following Correale *et al.*, 2006a, the problem can be stated rigorously as a MAX-SAT problem subject jointly to (1) and (8); see Supplementary Information. Proceeding in this way, we obtain that in our boolean gene networks the number of non-canonicalizing gates is upper bounded by  $D_{A,y}^{\text{emp}}$  (and for null models by  $D_A$ ). Because  $D_{A,y}^{\text{emp}} < D_A$  and  $D_{A,y}^{\text{emp}}$  approaches the Shannon bound, this criterion implies that non-canonicalizing gates for the true edge sign words  $\mathbf{y}$  are remarkably underrepresented with respect to the null models on the same topology.

*A perturbation analysis of the boolean dynamics.* In the rest of the article we consider unate boolean networks with the topology





**Fig. 5.** Steady-state perturbation of the boolean network of *E.coli*. Fifty different boolean networks are constructed on the topology of the *E.coli* gene network, for both the true sign assignment  $\gamma$  and for 50 random sign words drawn from  $B(q)$  (null models). The networks are unary, and non-canalizing gates (XOR) are allowed only in correspondence of negative undirected cycles (AND, OR and XOR are chosen at random in this case). Each curve of the plot represents an average over 500 trajectories obtained perturbing a steady state for both the true (solid) and the null (dashed) models. The size of the perturbation at  $t = 0$  varies between 1 and 8% of the number of genes. In all curves, perturbations on the (more canalized) true boolean network tend to settle at a larger distance from the steady state than in the (less canalized) null models

of our gene networks, in which the gates are drawn at random, but respecting the rule above that non-canalizing gates are admissible only in presence of negative undirected cycles. We focus in particular on the more complex *Escherichia coli* gene network (the other two have a lower size and a much simpler topology; see Supplementary Figure S2). It is interesting to compare the dynamics obtained in this way with those of the null models, in which the number of non-canalizing gates is higher. For random boolean networks, non-canalizing is taken as a proxy for complex, chaotic-like behavior; see Kauffman, 1993; Kauffman *et al.*, 2004. Topologically, our gene networks differ from random networks in many aspects, most importantly in the limited number of feedback loops and in the highly parallel feed-forward architecture (Cosentino Lagomarsino *et al.*, 2007); see Supplementary Table S8. For a boolean network with these characteristics, rather than rendering the dynamics chaotic, the presence of XOR gates has the effect of constraining the possible steady states admissible by the system. In this respect, the whole network replicates the blocking phenomenon we have observed on small motifs like FFL. In particular, in Figure 5, a perturbation analysis is shown. In the true network, perturbations of a steady state tend to persist more than in a null model, and the new steady state tends to remain more distant from the original one than in the null models. This behavior is consistent across all perturbation sizes and is not explainable in terms of the sign of the feedback loops: we recover in fact the same behavior when we consider the maximal directed acyclic graph extracted from the gene network, see Supplementary Figure S8.

In Figure 5, the fact that the perturbed initial conditions in average do not increase their distance to the steady state as time grows (apart from a transient of limited amplitude for the null models) is an indication of a non-chaotic regime. The fact that

the distance does not decrease much in time is less predictable because in an ordered regime one would expect perturbations to die out consistently. Comparing Figure 5 and Supplementary Figure S8, we can observe how the presence of a few directed cycles can impact significantly the response to perturbations, rendering the gene networks much less stable than one would expect from their nearly acyclic, massively feed-forward topology and their high level of canalization. This confirms, if necessary, that the predictions one obtains through a random boolean network (even those obtained choosing a flat distribution of boolean rules on a given topology) are of limited significance for models closer to the reality of currently available gene networks.

**Funding:** C.A. acknowledges financial support from Ministero dell'Istruzione, dell'Università e della Ricerca. The EU-IndiaGRID2 project (European FP7 e-Infrastructure Grant 246698) is acknowledged for the use of its grid infrastructure.

**Conflict of Interest:** none declared.

## REFERENCES

- Albert, R. and Othmer, H.G. (2003) The topology of the regulatory interactions predicts the expression pattern of the segment polarity genes in *Drosophila melanogaster*. *J. Theor. Biol.*, **223**, 1–18.
- Balleza, E. *et al.* (2008) Critical dynamics in genetic regulatory networks: examples from four kingdoms. *PLoS One*, **3**, e2456.
- Bonneau, R. (2008) Learning biological networks: from modules to dynamics. *Nat. Chem. Biol.*, **4**, 658–664.
- Buchler, N.E. *et al.* (2003) On schemes of combinatorial transcription logic. *Proc. Natl Acad. Sci. USA*, **100**, 5136–5141.
- Chowdhury, S. *et al.* (2010) Information propagation within the genetic network of *Saccharomyces cerevisiae*. *BMC Syst. Biol.*, **4**, 143.
- Ciliberti, S. and Mézard, M. (2005) The theoretical capacity of the parity source coder. *J. Stat. Mech.*, P10003.
- Correale, L. *et al.* (2006a) The computational core and fixed point organization in boolean networks. *J. Stat. Mech.*, P03002.
- Correale, L. *et al.* (2006b) Core percolation and onset of complexity in Boolean networks. *Phys. Rev. Lett.*, **96**, 018101.
- Cosentino Lagomarsino, M. *et al.* (2005) Logic backbone of a transcription network. *Phys. Rev. Lett.*, **95**, 158701.
- Cosentino Lagomarsino, M. *et al.* (2007) Hierarchy and feedback in the evolution of the *Escherichia coli* transcription network. *Proc. Natl Acad. Sci. USA*, **104**, 5516–5520.
- Cover, T.M. and Thomas, J.A. (2006) *Elements of Information Theory* (Wiley Series in Telecommunications and Signal Processing). Wiley-Interscience, Hoboken, NJ.
- Facchetti, G. *et al.* (2011) Computing global structural balance in large-scale signed social networks. *Proc. Natl Acad. Sci. USA*, **108**, 20953–20958.
- Harris, S.E. *et al.* (2002) A model of transcriptional regulatory networks based on biases in the observed regulation rules. *Complexity*, **7**, 23–40.
- Iacono, G. and Altafini, C. (2010) Monotonicity, frustration, and ordered response: an analysis of the energy landscape of perturbed large-scale biological networks. *BMC Syst. Biol.*, **4**, 83.
- Iacono, G. *et al.* (2010) Determining the distance to monotonicity of a biological network: a graph-theoretical approach. *IET Syst. Biol.*, **4**, 223–235.
- Kauffman, S. (1993) *The Origins of Order: Self-Organization and Selection in Evolution*. Oxford University Press, NY.
- Kauffman, S. *et al.* (2003) Random boolean network models and the yeast transcriptional network. *Proc. Natl Acad. Sci. USA*, **100**, 14796–14799.
- Kauffman, S. *et al.* (2004) Genetic networks with canalizing boolean rules are always stable. *Proc. Natl Acad. Sci. USA*, **101**, 17102–17107.
- Kim, H. *et al.* (2007) Boolean networks using the chi-square test for inferring large-scale gene regulatory networks. *BMC Bioinformatics*, **8**, 37.
- Lau, K.Y. *et al.* (2007) Function constrains network architecture and dynamics: a case study on the yeast cell cycle Boolean network. *Phys. Rev. E Stat. Nonlin. Soft Matter Phys.*, **75**, 051907.

- Mangan, S. and Alon, U. (2003) Structure and function of the feed-forward loop network motif. *Proc. Natl Acad. Sci. USA*, **100**, 11980–11985.
- Materna, S.C. and Davidson, E.H. (2007) Logic of gene regulatory networks. *Curr. Opin. Biotechnol.*, **18**, 351–354.
- Mezard, M. and Montanari, A. (2009) *Information, Physics, and Computation*. Oxford University Press, New York, NY, USA.
- Mezard, M. *et al.* (2003) Two solutions to diluted p-spin models and XORSAT problems. *J. Stat. Phys.*, **111**, 505–533.
- Setty, Y. *et al.* (2003) Detailed map of a cis-regulatory input function. *Proc. Natl Acad. Sci. USA*, **100**, 7702–7707.
- Shmulevich, I. *et al.* (2002) From Boolean to probabilistic Boolean networks as models of genetic regulatory networks. *Proc. IEEE*, **90**, 1778–1792.
- Shmulevich, I. *et al.* (2005) Eukaryotic cells are dynamically ordered or critical but not chaotic. *Proc. Natl Acad. Sci. USA*, **102**, 13439–13444.
- Silva-Rocha, R. and de Lorenzo, V. (2008) Mining logic gates in prokaryotic transcriptional regulation networks. *FEBS Lett.*, **582**, 1237–1244.
- Sontag, E. *et al.* (2008) The effect of negative feedback loops on the dynamics of Boolean networks. *Biophys. J.*, **95**, 518–526.
- van Hijum, S.A. *et al.* (2009) Mechanisms and evolution of control logic in prokaryotic transcriptional regulation. *Microbiol. Mol. Biol. Rev.*, **73**, 481–509.
- Wainwright, M.J. *et al.* (2010) Lossy source compression using low-density generator matrix codes: analysis and algorithms. *IEEE Trans. Inf. Theory*, **56**, 1351–1368.
- Zou, Y.M. (2010) Modeling and analyzing complex biological networks incorporating experimental information on both network topology and stable states. *Bioinformatics*, **26**, 2037–2041.