



# The effective graph reveals redundancy, canalization, and control pathways in biochemical regulation and signaling

Alexander J. Gates<sup>a,1</sup> , Rion Brattig Correia<sup>b,c</sup> , Xuan Wang<sup>d</sup> , and Luis M. Rocha<sup>b,d,e,1</sup>

<sup>a</sup>Network Science Institute, Northeastern University, Boston, MA 02115; <sup>b</sup>Instituto Gulbenkian de Ciéncia, 2780-156 Oeiras, Portugal; <sup>c</sup>Coordenação de Aperfeiçoamento de Pessoal de Nível Superior, Ministry of Education of Brazil, 70040-020 Brasília, DF, Brazil; <sup>d</sup>Center for Social and Biomedical Complexity, Luddy School of Informatics, Computing & Engineering, Indiana University, Bloomington, IN 47408; and <sup>e</sup>Department of Systems Science and Industrial Engineering, Binghamton University, Binghamton, NY 13902

Edited by Herbert Levine, Northeastern University, Boston, MA, and approved February 17, 2021 (received for review November 25, 2020)

**The ability to map causal interactions underlying genetic control and cellular signaling has led to increasingly accurate models of the complex biochemical networks that regulate cellular function. These network models provide deep insights into the organization, dynamics, and function of biochemical systems: for example, by revealing genetic control pathways involved in disease. However, the traditional representation of biochemical networks as binary interaction graphs fails to accurately represent an important dynamical feature of these multivariate systems: some pathways propagate control signals much more effectively than do others. Such heterogeneity of interactions reflects canalization—the system is robust to dynamical interventions in redundant pathways but responsive to interventions in effective pathways. Here, we introduce the effective graph, a weighted graph that captures the nonlinear logical redundancy present in biochemical network regulation, signaling, and control. Using 78 experimentally validated models derived from systems biology, we demonstrate that 1) redundant pathways are prevalent in biological models of biochemical regulation, 2) the effective graph provides a probabilistic but precise characterization of multivariate dynamics in a causal graph form, and 3) the effective graph provides an accurate explanation of how dynamical perturbation and control signals, such as those induced by cancer drug therapies, propagate in biochemical pathways. Overall, our results indicate that the effective graph provides an enriched description of the structure and dynamics of networked multivariate causal interactions. We demonstrate that it improves explainability, prediction, and control of complex dynamical systems in general and biochemical regulation in particular.**

biochemical regulation | Boolean network | canalization | complex networks | complex networks

Increasing evidence indicates that nonlinear interactions between biochemical variables—such as cell signaling, protein interactions, and genetic regulation and suppression—are pervasive (1–5), yet linear models of biochemical regulation fail to capture these key features of network causality (6). The simplest way to model such causal interdependent nonlinear dynamics is with multivariate discrete dynamical systems, also known as automata networks. Boolean networks (BNs), for instance, are canonical models of complex systems that exhibit a wide range of dynamical behaviors (3, 7). They have been successfully used to reveal insights into the dynamics of biochemical regulation (8), cell signaling (9), metabolism (10), anticancer drug response (11), and neuronal action potentials (12), among other things (13). In addition, BNs provide a convenient modeling framework to explore general properties of complex systems, such as self-organization, criticality, causality, canalization, robustness, and evolvability (3, 14–19).

The success of BNs can be attributed largely to three features of these models (7, 13, 20, 21): 1) qualitative thresholds

to measure transitions in concentration/expression of biochemical molecules in experimental data without the need for precise parameter estimation; 2) interaction graphs that synthesize complex multivariate dynamics to reveal the topology of the causal organization of biological systems; and 3) discrete dynamics that facilitate the prediction of critical behavior, self-organization, robustness, evolvability, and controllability. The first feature makes BNs very useful for estimating predictive systems biology models from data, especially because many processes in biology—such as gene expression and immune or neuron activation—are characterized by switch-like transitions between the presence or absence of a biochemical molecule or signal (13, 21). The second and third features of BNs make them ideal models to explore the interplay between the organization and the dynamics of complex systems (22, 23). Traditionally, the organization and dynamics of BNs are captured by general probabilistic parameters of the system variables (e.g., the mean number of interactions between variables or mean node bias) that are used to predict features of system-wide behavior, such as the transition from order to chaos (14, 18). Interactions between BN variables are usually represented as directed graphs, with arrows indicating when one node variable is an input to the logical rules governing another node variable. Thus,

## Significance

Many biological networks are modeled with multivariate discrete dynamical systems. Current theory suggests that the network of interactions captures salient features of system dynamics, but it misses a key aspect of these networks: some interactions are more important than others due to dynamical redundancy and nonlinearity. This unequivalence leads to a canalized dynamics that differs from constraints inferred from network structure alone. To capture the redundancy present in biochemical regulatory and signaling interactions, we present the effective graph, an experimentally validated mathematical framework that synthesizes both structure and dynamics in a weighted graph representation of discrete multivariate systems. Our results demonstrate the ubiquity of redundancy in biology and provide a tool to increase causal explainability and control of biochemical systems.

Author contributions: A.J.G. and L.M.R. designed research; A.J.G., R.B.C., X.W., and L.M.R. analyzed data; and A.J.G., R.B.C., X.W., and L.M.R. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Published under the [PNAS license](#).

See [online](#) for related content such as Commentaries.

<sup>1</sup>To whom correspondence may be addressed. Email: a.gates@northeastern.edu or rocha@binghamton.edu.

This article contains supporting information online at <https://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2022598118/-DCSupplemental>.

Published March 18, 2021.

interaction graphs treat all inputs to a variable with equal importance, even though each input may have a weaker or stronger role in determining state transitions.

In reality, the states of almost all biochemical node variables are robust to dynamic perturbations from many of their input variables but highly responsive to just a few (3, 19). Such dynamical redundancy is a ubiquitous hallmark of the third feature of BNs that has been used to study canalization in biological complexity (17, 24)—a concept Waddington (25) introduced to characterize the mechanisms organisms use to buffer development, regulation, and evolution against perturbations. Indeed, the presence of canalization can drastically alter the functional interaction topology of BNs, with profound consequences to the stability and controllability of biological systems (17, 20, 24).

To better capture the functionally relevant pathways of BN models of biochemical regulation and signaling, we introduce the effective graph. It uses a measure of collective canalization to take into account nonlinear effects present when several inputs are needed to regulate a variable. This way, the effective graph integrates all of the dynamical redundancy present in BN dynamics, thus revealing the most important interactions and pathways in determining state transitions. Through an analysis of 78 experimentally validated biological models across a wide range of different biochemical systems and cell types (*SI Appendix, section 2*), we show that interactions in biological networks are on average much less effective at generating state transitions than interactions in random Boolean automata. For instance, in gene regulation, this means that a gene on its own is less likely to regulate the expression of another gene it interacts with than what would be expected from the set of possible gene–gene interactions.

The effective graph provides a probabilistic characterization of multivariate interactions and dynamics in a causal graph form. It also captures how conditioning the system on known input states, such as when administering a drug intervention, can modify the remaining biochemical interactions. The conditional effective graph thus provides a mechanistic explanation for how control propagates through biochemical models and how causal, nonlinear, microlevel interactions integrate to define macrolevel biological function. We leverage this analytical tool to study a model of signal transduction in ER+ breast cancer (26) to reveal why and how certain drugs drive cancer cells to proliferate or die and identify the modular pathway dynamics that facilitate or hinder this control.

Finally, the redundancy observed in BN models from systems biology also reveals that only a fraction of causal interactions is typically needed to determine convergence to dynamical attractors, which represent biological function in these models. This suggests that the regulatory dynamics of biological networks are robust to random dynamical perturbations yet controllable via the most effective pathways revealed by the effective graph. To demonstrate this observation, we show that the effective graph is consistently better than the original interaction graph at predicting the impact of dynamical perturbations across random networks, a model of floral organ specification in the flowering plant *Arabidopsis thaliana* (27, 28), and the ER+ breast cancer model (26). Given the widespread applicability of BNs, our framework opens a promising research direction in the control of complex dynamical systems and can facilitate the design of interventions in systems biology models, especially those for development and disease. By synthesizing structure and dynamics into a single-graph formalism, the effective graph increases the predictability and explainability of actionable models of biochemical regulation and signaling and causal automata models in general.

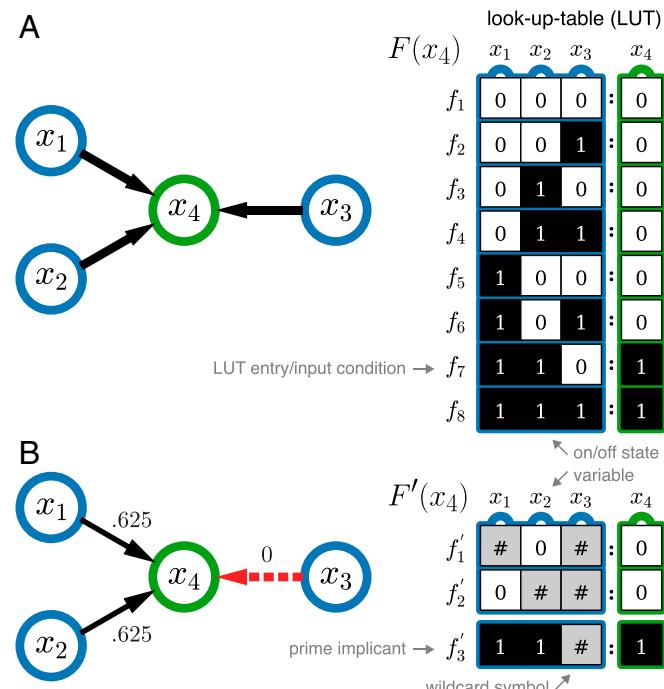
### Canalization of Boolean Automata

A Boolean automaton is a binary variable,  $x \in \{0, 1\}$ , whose state is updated in discrete time steps,  $t$ , according to a determinis-

tic state-transition function relating the states of  $k$  inputs to its own state at the next time step:  $x^{t+1} = f(x_1^t, \dots, x_k^t)$ . This logical function,  $f : \{0, 1\}^k \rightarrow \{0, 1\}$ , is defined by a look-up (truth) table (LUT),  $F \equiv \{f_\alpha : \alpha = 1, \dots, 2^k\}$ , with one entry for each of the  $2^k$  combinations of input states and a mapping to the automaton's next state (transition or output),  $x^{t+1}$ . The bias,  $\rho$ , of the automata is the fraction of transitions to state 1 in the output column of the LUT. An exemplar Boolean automaton with its LUT is shown in Fig. 1A.

A BN is a graph  $\mathcal{B} \equiv (X, C)$ , where  $X$  is a set of  $N$  Boolean automata nodes  $x_i \in X, i = 1, \dots, N$  and  $C$  is a set of directed edges,  $c_{ji} \in C : x_i, x_j \in X$ , that represent the interaction network, denoting that automaton  $x_j$  is an input to automaton  $x_i$ , as computed by  $f_i(x_1, \dots, x_j, \dots, x_k)$  with LUT  $F_i$ : for example, the interaction graph for the BN model of the floral organ development in the *A. thaliana* plant (28) (see Fig. 4A). The set of inputs into automaton  $x_i$  is denoted by  $X_i = \{x_j \in X : c_{ji} \in C\}$ , and its cardinality,  $k_i = |X_i|$ , is the in-degree of node  $x_i$ . At any given time  $t$ ,  $\mathcal{B}$  is in a specific configuration of automata states,  $\mathbf{x}^t = \langle x_1^t, x_2^t, \dots, x_N^t \rangle$ —we use the terms state for individual automata ( $x_i^t$ ) and configuration ( $\mathbf{x}^t$ ) for the collective network state (i.e., the vector of states of all automata of the BN at time  $t$ ). The set of all possible network configurations is denoted by  $\mathcal{X} \equiv \{0, 1\}^N$ , where  $|\mathcal{X}| = 2^N$ . BNs update synchronously (all automata simultaneously at time  $t$ ) or asynchronously (some automata are selected randomly or via a schedule at time  $t$ ). However, the effective graph does not depend on the chosen update policy since it is constructed from the redundancy parameters of each automaton considered separately.

The canalization of an automaton reflects the fact that not all input states are equally important for determining its state



**Fig. 1.** Constructing the effective graph. (A, Left) The interaction graph of automaton  $x_4$  (green node), with  $k = 3$  input variables (blue nodes,  $x_1, x_2, x_3$ ) and (A, Right) its corresponding Boolean logic given by the LUT, with bias  $\rho(x_4) = 1/4$ . (B, Left) The effective graph of automata  $x_4$  is built from the wild card redescription of the LUT (B, Right),  $F'$ , which shows that input  $x_3$  is always redundant (only wild cards in its column) and that  $x_4 = x_1 \wedge x_2$ . Edge thickness denotes edge effectiveness,  $e_{ji}$ , with the fully redundant edge shown in dashed red. The total input redundancy of automaton  $x_4$  is  $k_r(x_4) = 1.75$ , and therefore, its effective connectivity is  $k_e(x_4) = 1.25$ .

transition (3). We follow Marques-Pita and Rocha (17) by quantifying canalization through the amount of logical redundancy present in the automata. Specifically, we use the first step of the Quine–McCluskey Boolean minimization algorithm (29) to identify inputs of an automaton  $x$ , which are redundant given the state of its other inputs. This procedure compresses an LUT into the set of all distinct prime implicants of  $f$ , represented as a set of wild card schemata,  $F' \equiv \{f'_v\}$ , in which the wild card or “don’t care” symbol,  $\#$ , denotes an input whose state is redundant for determining the automaton transition given the states of other necessary inputs. For instance, schema  $f'_1$  in the Fig. 1B example specifies that when input  $x_2 = 0$ , the states of inputs  $x_1$  and  $x_3$  are redundant to determine the next state of  $x_4$ , which is guaranteed to be 0. In this process, the original LUT  $F$  is redescribed into a complete set of schemata  $F'$  (Fig. 1).

Every wild card schema  $f'_v \in F'$  redescribes a subset of entries in the original LUT, denoted by  $\Upsilon_v \equiv \{f_\alpha : f_\alpha \rightarrowtail f'_v\} \subseteq F$ , where  $\rightarrowtail$  means “is redescribed by.” For example, schema  $f'_1$  in the Fig. 1 example redescribes the set of LUT entries  $\Upsilon_1 \equiv \{f_1, f_2, f_5, f_6\} \subseteq F$ . The set of (overlapping) schemata  $F'$  is complete since it contains all unique prime implicants that redescription all entries of the original LUT, here described as wild card schemata. In Boolean minimization, the set of prime implicants can be further reduced (via the additional steps of the Quine–McCluskey algorithm or equivalent methods), but because our goal is to tally all possible minimal transition conditions of an automaton, we preserve all prime implicants; ref. 17 and *SI Appendix* have details. Notice that all measures that ensue are computed from the entire population of prime implicants and are thus parameters, not sampled statistics, of logical functions  $f_i$ .

The amount of canalization present in the logic of an automaton can be quantified by probabilistic parameters derived from the schema redescription of its LUT. Input redundancy,  $k_r(x)$ , measures the number of inputs that, on average, are not needed to determine the state of automaton  $x$ , assuming that all input combinations are equally likely. It is quantified by tallying the mean number of wild card symbols present in schemata set  $F'(x)$  that redescribes LUT  $F(x)$ :

$$k_r(x) = \frac{\sum_{f_\alpha \in F} \text{avg}_{v:f_\alpha \in \Upsilon_v} (n_v^\#)}{|F|}, \quad [1]$$

where  $n_v^\#$  is the number of  $\#$  symbols in schema  $f'_v$ . In computing  $k_r(x)$ , we assume that each entry  $f_\alpha$  of LUT  $F$  can be redescribed with equal likelihood by any of the schemata  $f'_v$  in  $F'(x)$  that includes it ( $f_\alpha \in \Upsilon_v$ ). Thus, we use the average operator ( $\text{avg}$ ) in Eq. 1. This is the same as assuming that any schema (or prime implicant) is a viable intervention possibility to change the state of automaton  $x$ . Other redundancy aggregations are possible (17), but averaging over all possible schemata allows the per-edge separation of redundancy we pursue below (*SI Appendix*, section 1 has additional discussion).

A complementary parameter of the redundancy of automaton  $x$  is its effective connectivity:

$$k_e(x) = k(x) - k_r(x), \quad [2]$$

which yields the number of inputs that are on average necessary to determine the automaton’s state. Whereas  $k(x)$  is the number of inputs to automaton  $x$  present in the interaction graph of the BN (in-degree),  $k_e(x)$  measures the number of such inputs that are actually (on average) necessary to determine the state of  $x$ —the effective connectivity of  $x$ . In the Fig. 1B example, because six entries of LUT  $F$  ( $f_1 \dots f_6$ ) are redescribed by schemata with two wild cards ( $f'_1, f'_2$ ) and two entries ( $f_7, f_8$ ) are

redescribed by schemata with one wild card ( $f'_3$ ), via Eqs. 1 and 2 we obtain  $k_r(x_4) = (6 \times 2 + 2 \times 1)/8 = 1.75$  and  $k_e(x_4) = 3 - 1.75 = 1.25$ . In other words, on average, 1.75 inputs to  $x_4$  are redundant, and thus, its effective connectivity is 1.25—in contrast to its in-degree of three.

Other automata parameters, distinct from Eqs. 1 and 2, can be used to measure canalization. For instance, sensitivity (30) also aims to measure the effective dynamics of a Boolean automaton, but as we discuss below, it does not capture the nonlinear effects of collective canalization. We can also extract additional redundancy from the symmetries that exist in the schemata set  $F'$ , thus providing a further compression of this set (17, 31), but we do not consider symmetry redundancy in the present analysis. Additional algorithmic details as well as relationships between canalization, control, robustness, and modularity of BN models are presented in refs. 17 and 20 and *SI Appendix*.

Most automata contain some amount of input redundancy; only the two parity functions for any  $k$  have  $k_r = 0$  (e.g., the exclusive OR, XOR function and its negation for  $k = 2$ ). Therefore, the original interaction graph of a BN misses the high amount of redundancy present in most BNs and does not capture how automata truly influence one another in a network.

## The Effective Graph and Redundancy in Models of Biochemical Regulation and Signaling

The input redundancy and effective connectivity of Eqs. 1 and 2 reveal that, on average, the interaction graph overestimates the number of inputs needed to determine transitions. However, these parameters do not specify which of the interactions are actually more effective and how they combine to form pathways that transmit signals through the network. To measure how input redundancy is distributed over the individual inputs to an automaton, we introduce the per-input parameters of redundancy and effectiveness. The latter is then used to compute the edge weights of the effective graph, which provides a (probabilistic) synthesis of the canalizing dynamics of a BN.

Edge redundancy,  $r_{ji} \in [0, 1]$ , tells us, on average, how redundant an incoming edge from automaton  $x_j$  is in determining the state of automaton  $x_i$ . This is computed by counting the average number of schema in  $F'_i$  in which input  $x_j$  is specified by a wild card symbol:

$$r_{ji} = \frac{\sum_{f_\alpha \in F_i} \text{avg}_{v:f_\alpha \in \Upsilon_v} (j \rightarrowtail \#)_v}{|F_i|}, \quad [3]$$

where  $(j \rightarrowtail \#)_v$  is a logical condition that assumes the truth value one if input  $x_j$  is a wild card in schema  $f'_v$  and zero otherwise;  $\text{avg}$  is the average operator. Similarly, edge effectiveness,  $e_{ji} \in [0, 1]$ , captures the extent to which an incoming edge from automaton  $x_j$  is on average necessary to determine the value of automaton  $x_i$ :

$$e_{ji} = 1 - r_{ji}. \quad [4]$$

Naturally,  $k_r(x_i) = \sum_j r_{ji}$  and  $k_e(x_i) = \sum_j e_{ji}$ , meaning that the canalization of an automaton is additive over its incoming edges.

We can now define the effective graph of a BN to capture the varying influence of each input edge on the dynamics of automata nodes. Specifically,  $\mathcal{E} \equiv (X, E)$ , where  $X$  is the set of automata and  $E$  is the set of directed edges, weighted by their effectiveness  $e_{ji}$  as defined by Eq. 4. Note that an edge (interaction) can be fully redundant if its effectiveness is null,  $e_{ji} = 0$ . This is the case of input  $x_3$  in the Fig. 1B example, which is always redescribed by a wild card in  $F'(x_4)$  and thus,  $e_{34} = 0$ . In practice, fully redundant edges should be completely removed, but in this article, to catalog their existence, we emphasize them as red dashed edges (e.g., edge  $e_{34}$  in Fig. 1B).

One may think that fully redundant edges should not occur in well-constructed networks; however, they are fairly common in systems biology models. Indeed, we analyzed 78 Boolean models stored on the Cell Collective (9) and found that 17 of them (22%) contained at least one fully redundant edge, with 87 fully redundant edges in total. The inclusion of fully redundant edges in these models may result from inference methods based on information theory that can fail to capture polyadic relationships (32), but the most likely reason is an incomplete record of experimental observation. Typically, systems biology models integrate many experimental studies conducted by many different teams in different scenarios, which are available in interaction databases and the published literature (33). Modelers who integrate such scientific evidence have to make decisions about conflicting or weak evidence (13). For instance, the *A. thaliana* model studied below (see Fig. 4) contains three fully redundant edges, which ultimately result from “subjective decisions given alternatives with equivalent results” (27). *SI Appendix, section 2B* has a more detailed discussion of how these issues can lead to fully redundant edges in systems biology models. Certainly, our methodology can serve as a logical check on these models to remove completely redundant interactions.

The effective graph is a probabilistic synthesis of the dynamical redundancy of a BN model given all its possible initial conditions. However, in systems biology we often want to study a model under specific initial conditions: for example, cells in a cancer state or under the influence of a particular drug, as pursued below in the analysis of the estrogen receptor positive (ER+) breast cancer model. Since the set of possible initial conditions in such cases is reduced, the interaction topology of the effective graph changes. This is easily captured in our methodology by

conditioning the effective graph  $\mathcal{E}$  on a set of variables,  $K \subseteq X$ , that are fixed to specific constant states. The resulting conditional effective graph,  $\mathcal{E}|K$ , can have a drastically altered effective topology, for instance, with many more interactions revealed to be fully redundant.

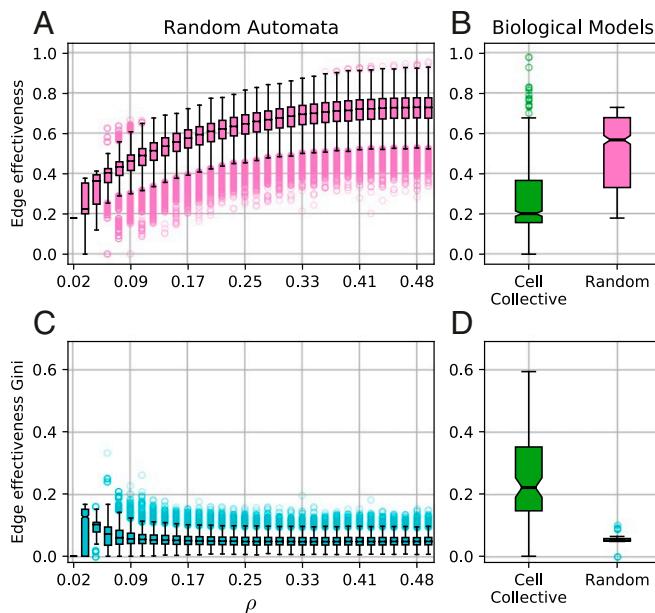
The computational complexity of our canalization parameters and the effective graph scale linearly with the number of nodes  $N$  and can thus be computed for large BNs (17), unlike most methodologies used to analyze the dynamics of BNs. Instead, the computational complexity bottleneck to derive the effective graph is bounded by the Quine–McCluskey algorithm (29) on the largest degree node in the BN: that is, the automaton with the largest  $k_i$ . When this value is very large, one can sample the prime implicant population, but none of the analysis here pursued required such estimation. We provide a full implementation of all canalization parameters (Eqs. 1–4) and the effective graph in the open-source CANA python package (31).

### Effectiveness of Biochemical Interactions

Input redundancy (Eq. 1) is prevalent in random Boolean automata. In BNs, this leads to a lower effective connectivity (Eq. 2) for automata nodes than the interaction graph (in-degree) specifies, with varying edge effectiveness (Eq. 4) distributed across inputs. The prevalence and variation of edge redundancy are shown in Fig. 2A for random Boolean automata of degree  $k = 6$ . For all values of bias ( $\rho$ ), we observe much variation in edge effectiveness, although its median value goes from  $\bar{e}_{ji} \approx 0.18$  at the lowest bias ( $\rho = \frac{1}{64}$ ) to  $\bar{e}_{ji} \approx 0.75$  at the highest bias ( $\rho = \frac{1}{2}$ ). The upward shift of the distribution of edge effectiveness indicates that inputs tend to become more important for determining the state transition of the automata as bias increases. The behavior for automata with other  $k$  is similar (*SI Appendix, Fig. S3*).

The observed distribution of edge effectiveness for random Boolean automata provides context for next question: how much redundancy is present in experimentally validated biochemical interactions? To answer this question, we calculated the edge effectiveness of all 8,220 interaction edges from the 78 BN models in the Cell Collective (*SI Appendix, section 2*). We compare this distribution with an ensemble of random automata matching the degree ( $k$ ) and bias ( $\rho$ ) of the Cell Collective automata. Specifically, for each automaton from the systems biology models, we sample  $10^3$  random automata with exactly the same degree and bias. We observe that the mean edge effectiveness of interactions in the biochemical networks is much smaller than that of interactions in the random ensemble (Fig. 2B). For simplicity but without loss of generality, Fig. 2B depicts the distribution of effectiveness for 630 incoming edges to 105 automata of degree  $k = 6$  in the systems biology models, as compared with that of the bias-matched random ensemble of same  $k$ ; distribution comparisons for other values of  $k$  are shown in *SI Appendix, Fig. S3*. A two-sample independent  $t$  test for the difference in the means between the experimentally validated biochemical (0.27) and the random interactions (0.51) confirms the statistical differences between these distributions for  $k = 6$  automata, with a  $P$  value  $< 10^{-100}$ .

Edge effectiveness allows us to differentiate network interactions based on how much they contribute to determining automata transitions and to identify the inputs that most control a given automaton. In contrast, the original interaction graph of a BN does not differentiate the inputs to an automaton. Let us look at how edge effectiveness differentiates input importance with an example. Consider two automata,  $b$  (balanced) and  $u$  (unbalanced), each with  $k = 4$  inputs:  $x_1, x_2, x_3, x_4$ . In the first case, the transition function is specified by  $f_b = x_1 \wedge x_2 \wedge x_3 \wedge x_4$ , a symmetric logic since all inputs are interchangeable. In this balanced case, the edge effectiveness is the same for all incoming edges  $e_{j,b} \approx 0.3$ . In the second case, the transition function is



**Fig. 2.** Central tendency, variation, and heterogeneity of edge effectiveness of Boolean automata in biochemical regulation and random ensembles. (A) The distributions of edge effectiveness for ensembles of  $10^4$  automata with degree  $k = 6$  at each bias  $\rho$ . (B) The distribution of edge effectiveness of the 630 incoming interactions to 105 automata with degree  $k = 6$  in Cell Collective models (green) compared with a bias-matched sample of random Boolean automata (pink). (C) The distributions of edge effectiveness Gini coefficients for inputs to automata in each of the random ensembles from A. (D) The distribution of edge effectiveness Gini coefficients for inputs to the 105 automata with degree  $k = 6$  in the Cell Collective models (green) compared with the bias-matched ensemble of random Boolean automata (cyan).

specified by  $f_u = x_1 \vee (x_2 \wedge x_3 \wedge x_4)$ , a logic where input  $x_1$  primarily determines the transition. Accordingly, in this unbalanced case the edge effectiveness varies by incoming edge:  $e_{1,u} \approx 0.91$  and  $e_{2,b} = e_{3,b} = e_{4,b} \approx 0.24$ , which reflect the importance of input  $x_1$  in determining the state of the automaton.

To ascertain the heterogeneity of effectiveness in biochemical interactions, we again compare the systems biology models in the Cell Collective with similar random ensembles. For each automaton, we compute the Gini inequality coefficient to obtain a real number between zero and one, where zero denotes that all inputs have exactly the same effectiveness and one denotes maximum inequality among the inputs (i.e., one input completely dominates over the others) (*SI Appendix, section 3*). When applied to the balanced  $b$  and unbalanced  $u$  example automata, we obtain Gini coefficients of 0 and 0.31, respectively. This accurately reflects the effectiveness equality of the inputs to  $b$  and the effectiveness inequality of the inputs to  $u$ .

The Gini coefficient calculated for the automata in the Cell Collective database reveals greater effectiveness inequality in their interaction than in comparable random automata. In other words, a smaller subset of inputs plays a more important role in controlling the biochemical variables in these models than is expected by chance. Consider first Fig. 2C, where edge effectiveness in random automata with degree  $k=6$  is characterized by a relatively small Gini coefficient for all biases. This indicates that in random automata, while redundancy is pervasive (Fig. 2A), it is distributed similarly over the inputs—incoming edges are roughly equally effective in determining the state of the random automata. In contrast, as shown in Fig. 2D, the Gini coefficient of the edge effectiveness of automata with  $k=6$  in the 78 biological models varies much more but is consistently higher than the bias-matched random automata. This is further supported by a two-sample independent  $t$  test for the difference in the means of the automata distributions in biochemical models (0.22) and in the random ensemble (0.05), which confirms the statistical differences between these distributions with a  $P$  value  $< 10^{-100}$ .

In summary, our analysis of edge effectiveness demonstrates not only that automata used to model biochemical regulation in the Cell Collective contain more redundancy than expected in random automata but also, that this redundancy is unevenly distributed over their inputs. In other words, in the models of biochemical regulation, only a few interactions are effective in controlling variable transitions, while most interactions are redundant and not very dynamically effective.

### Collective Canalization in Dynamical Regulation

A crucial feature of Boolean automata is the potential for highly nonlinear integration over their inputs (14, 15). Canalization is one such nonlinear phenomenon whereby a subset of inputs jointly determines the state of an automaton while rendering redundant the complement subset of inputs (3, 17). However, existing measures of canalization do not consider the full range of nonlinear joint interaction.

We can measure the extent to which nonlinear collective canalization is present in an automaton by comparing our per-input canalization and redundancy parameters—that capture joint dependencies—with parameters that consider each input independently. One such measure assuming independence is the activity of an input  $x_j$  to a Boolean automaton  $x_i$ :  $a_j(x_i)$ . It is the probability,  $P(\neg x_i^{t+1} | \neg x_j^t)$ , that automaton  $x_i$  flips its state at  $t+1$  when its input  $x_j$  flips its state at  $t$ , given a uniform distribution of input states at  $t$  (30). In turn, the sensitivity of automaton  $x_i$  is the sum of all its input activities  $s(x_i) = \sum_j a_j(x_i)$ .

Interestingly, our formulation of canalization via schema redescription also yields the activity of an input with a simple modification to formula (Eqs. 3 and 4), by substituting the maximum operator (max) for the average operator (avg):

$$a_j(x_i) = 1 - \frac{\sum_{f_\alpha \in F_i} \max_{v: f_\alpha \in \Upsilon_v^i} (j \rightarrow \#)_v}{|F_i|}. \quad [5]$$

*SI Appendix, section 1.C* has a proof of this formulation of activity. From here, it follows that  $e_{ji} \geq a_j(x_i)$  and  $k_e(x_i) \geq s(x_i)$ . This fact allows us to measure how much of the effective connectivity of an automaton  $x_i$  and the effectiveness of its inputs  $x_j$  derives from joint interactions among the inputs:

$$k_c(x_i) = k_e(x_i) - s(x_i), \quad c_{ji} = e_{ji} - a_j(x_i). \quad [6]$$

In other words,  $k_c(x_i)$  and  $c_{ji}$  measure the portion of canalization that derives from collective canalization at the node and input levels of a BN—in excess of sensitivity and activity, respectively.

Because collective canalization is very common, especially as the number of inputs ( $k$ ) increases (3), the distinction between effective connectivity and sensitivity is quite relevant for understanding the true regulatory dynamics in BNs, especially in systems biology models. Indeed, even for Boolean automata of  $k=2$ , the sensitivity parameter does not discriminate between such common Boolean functions as conjunction/disjunction and proposition/negation:  $s(x_1 \wedge x_2) = s(x_1 \vee x_2) = s(x_1) = s(\neg x_1) = 1$ . In contrast, effective connectivity correctly accounts for the additional collective canalization that is present in the conjunction/disjunction (and other) functions:  $k_e(x_1 \wedge x_2) = k_e(x_1 \vee x_2) = 5/4 = 1.25$ , while  $k_e(x_1) = k_e(\neg x_1) = 1$ .

Collective canalization is at play even in the small BN shown in Fig. 1. The edges of its effective graph are  $e_{14} = e_{24} = 0.625$ ,  $e_{34} = 0$ , whereas the activity measured for the same interactions is  $a_1(x_4) = a_2(x_4) = 0.5$ ,  $a_3(x_2) = 0$ . The discrepancy occurs because  $x_1$  and  $x_2$  jointly determine  $x_4$  with a collective canalization of  $c_{14} = c_{24} = 0.125$ . Indeed, on average one input is not sufficient to determine the state of  $x_4$ , as sensitivity  $s(x_4) = 1$  implies. For one-quarter of the input configurations (two of eight entries in the LUT redescribed by schema  $f'_3$ ), both inputs  $x_1$  and  $x_2$  are needed to jointly determine the state of  $x_4$ , and thus, its collective canalization is  $k_c(x_4) = 0.25$ . Clearly, the effective connectivity value of  $k_e(x_4) = 1.25$  is a more accurate characteristic of how inputs jointly determine the state of  $x_4$ , by aggregating both their individual and collective contributions. On average, 1.25 inputs are needed to specify the transition of  $x_4$  (conversely,  $k_r(x_4) = 1.75$  inputs are on average redundant); *SI Appendix, section 1* has an additional example. While the concept of “ $c$  sensitivity” (34) extends sensitivity to subsets of  $c$  inputs, it results in a vector of values, which is less intuitive in a network context than the scalar parameter  $k_e$ .

Collective canalization is measured at node and edge levels unequivocally via Eq. 6 and characterizes differences in the canalization of Boolean functions that sensitivity and activity do not measure. The question of how much the collective canalization captured by our parameters affects the dynamics of BNs is beyond the scope of this paper. However, collective canalization has already been shown to lead to more accurate predictions of critical behavior across a wide range of BN connectivity and dynamical behavior (19).

Together, our results show that our canalization parameters (Eqs. 1–4) capture redundancy and dynamical effectiveness in BNs at the automaton node and edge levels. They encompass parameters such as sensitivity and activity and importantly, also account for the nonlinear effects of collective canalization. One goal of the effective graph is to precisely quantify the true impact of interactions in spreading perturbations and control signals in BN models of biochemical regulation. The edges, therefore, are weighted according to their dynamical effectiveness (Eq. 4) to capture both their activity and (nonlinear) collective canalization contributions (Eq. 6). Next, we study the utility of the effective graph in predicting the spread of dynamical perturbations and identifying control pathways.

## Effective Graph Predicts the Spread of Perturbations

An important goal of systems biology is to quantify and predict the spread of perturbations and control signals across networked regulatory pathways (35, 36). While the interaction graph of a BN is useful for a theory of perturbations because it specifies which automata are topologically reachable in a given number of time steps, it fails to capture the varying effectiveness of each interaction in propagating signals through network pathways. The effective graph, on the other hand, provides enriched information about dynamical redundancy and canalization from each constituent automaton. In this section, we demonstrate that the effective graph's enriched portrait better captures the spread of perturbations in both random BNs and experimentally validated systems biology models.

Many types of perturbations can be considered, including those that change the structure or logic of the original model such as edge shuffling or deletion (22). Here, we focus on how specific, fixed models of biochemical regulation respond to different dynamical conditions—such as cellular response to drug regimens in ER+ breast cancer. Thus, unless otherwise noted, by perturbation we mean negating the logical state of an automaton at time  $t$  (also known as bit-flip perturbation). The impact of such a perturbation to an automaton in a BN is quantified by the Boolean analogue of the partial derivative (37):

$$\partial_t^{(i)} x_j(\mathbf{x}_\alpha) = |x_j^t(\mathbf{x}_\alpha) - x_j^t(\mathbf{x}_\alpha^{-i})|, \quad [7]$$

where  $x_j^t(\mathbf{x}_\alpha)$  denotes the state (truth value) of automaton node  $x_j$  at time  $t$  when the BN is initiated with configuration  $\mathbf{x}^0 = \mathbf{x}_\alpha$  at time  $t=0$  and  $\mathbf{x}_\alpha^{-i}$  denotes configuration  $\mathbf{x}_\alpha$  with the state of automaton  $x_i$  negated. The partial derivative yields one if flipping the state of  $x_i$  in initial configuration  $\mathbf{x}^0$  leads to  $x_j$  flipping its state at time  $t$  and zero otherwise. The total impact on automaton  $x_j$  of perturbations to automaton  $x_i$  after  $t$  steps is the average over all initial configurations:

$$\iota_{ij}(t) = 2^{-N} \sum_{\alpha=1}^{2^N} \partial_t^{(i)} x_j(\mathbf{x}_\alpha). \quad [8]$$

For large BNs,  $\iota_{ij}(t)$  must be estimated by averaging over a random sample of initial network configurations.

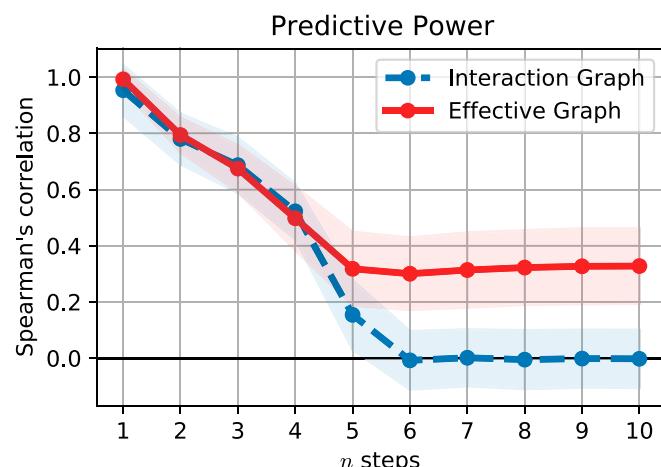
We now study how well the interaction and effective graphs predict the total impact of perturbations, using a different spreading model for each:  $\mathcal{M}_{IG}$  and  $\mathcal{M}_{EG}$ , respectively. To set up  $\mathcal{M}_{IG}$ , we consider that all nodes  $x_j$  connected via a path of at most  $t$  edges starting from node  $x_i$  are equally impacted by a perturbation to node  $x_i$ , where  $t$  is the number of time steps since the perturbation—the “light cone” of  $x_i$  as signals to any  $x_j$  cannot travel faster than the minimum number of edges (*SI Appendix*, Fig. S1). The second model  $\mathcal{M}_{EG}$  is similar except that the (weighted) effectiveness edges in the effective graph are assumed to proportionally constrain the spread of a perturbation (*SI Appendix*, section 4). This constraint is given by the product of edge weights in the strongest path between  $x_i$  and  $x_j$  (*SI Appendix*, Fig. S1), limited by the light cone such that the number of edges in the path is smaller than the number of elapsed time steps. By choosing the path with maximum product of edge weights as a surrogate measure for the total impact of perturbations, we assume [as in linear control (38)] that a signal can propagate without restriction via a connected path in the interaction graph model ( $\mathcal{M}_{IG}$ ), but edge effectiveness deferentially restricts propagation in the effective graph model ( $\mathcal{M}_{EG}$ ). To measure how well each model predicts which nodes  $x_j$  are most affected by perturbations to node  $x_i$ , we compute the average Spearman's rank correlation between each model and the true  $\iota_{ij}(t)$  at each time step.

We illustrate the superior predictive power of  $\mathcal{M}_{EG}$  first with an experiment using random BNs of  $N = 100$  nodes, fixed degree  $k = 3$ , and average bias  $\bar{\rho} = 0.4$  (*SI Appendix*, section 4). For each BN, we select 10 nodes  $x_i$  at random to perturb, approximating the total impact  $\iota_{ij}(t)$  on the other nodes  $x_j$  with a sample of  $10^4$  random initial configurations. As shown in Fig. 3, the rank correlation of  $\iota_{ij}(t)$  with the effective graph model,  $\mathcal{M}_{EG}$  (red), is consistently better than with the interaction graph model,  $\mathcal{M}_{IG}$  (blue). Indeed, after the full network is encompassed in the light cone,  $\mathcal{M}_{IG}$  is unable to differentiate which nodes might have been impacted by the perturbation, and Spearman's correlation is zero. In contrast, the effective graph retains predictive information about which variables are impacted by perturbations as measured by a significant positive Spearman correlation with the true dynamical impact. As shown below, the ability of the effective graph to predict perturbation spread in experimentally validated models of biochemical regulation is even more striking.

## Effective Graph Reveals How Control Pathways Function in Models of Biochemical Regulation

The characterization of control strategies in biomedicine can help focus experiments, aid the design of advanced disease therapeutics (1, 39), and even suggest intervention strategies to reprogram cells (40) (e.g., to revert a mutant cell to a wild-type state). It is well known that when the set of automata nodes  $X$  of a BN is large, enumeration of all configurations  $\mathbf{x} \in \mathcal{X}$  of its state-transition graph (STG) becomes difficult, making the controllability of BNs a nondeterministic-polynomial-time hard problem (41). Therefore, control methodologies that leverage the interaction graph or otherwise approximate the dynamics are highly desirable since they can greatly simplify the complexity of BN control (1, 17).

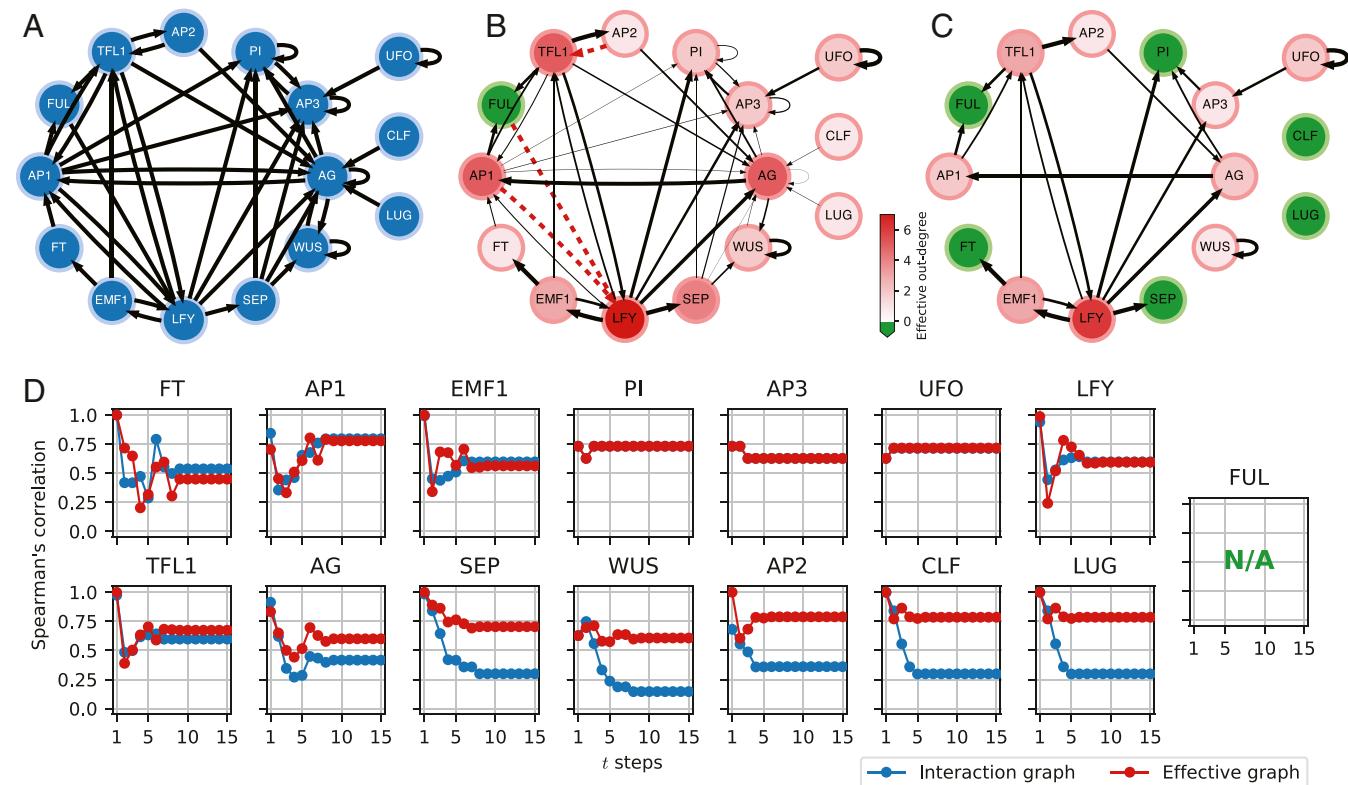
**Effective Graph Enhances Structure-Only Control Inference.** Several recent methodologies aim to determine the controllability of complex dynamical systems based solely on the graph of interactions between variables: structural controllability (SC) (38), minimum dominating set (MDS) (42), and feedback-vertex set control (FVC) (1, 43). By using only the interaction graph to predict minimum sets of variables (driver nodes) that are needed to control a network, these methods make predictions about



**Fig. 3.** The effective graph captures the spread of perturbations. The predictive power of the edge-product approximation using the two models,  $\mathcal{M}_{IG}$  (blue) and  $\mathcal{M}_{EG}$  (red), measured by the Spearman rank correlation (vertical axis) with the total impact,  $\iota_{ij}(t)$ , sampled from  $10^4$  trajectories, after  $t$  steps (horizontal axis). The shaded region denotes one SD for a sample of 100 random networks and 10 perturbed nodes per network.

the entire ensemble of dynamical systems that fit the same interaction graph (20). In contrast, since the effective graph is obtained by removing dynamical redundancy in a specific BN, the ensemble of possible dynamical systems that can fit is much smaller. Therefore, the effective graph is likely to lead to more precise inferences of control pathways in specific systems biology (BN) models than those derived from structure-only methods, such as SC, MDS, and FVC.

The removal of fully redundant edges from the interaction graph can reduce the number of feedback loops, revealing a smaller or distinct set of driver nodes than those predicted by FVC (*SI Appendix, section 5*). Consider the interaction and effective graphs of the *A. thaliana* flower development BN (27, 28) (TBN) in Fig. 4A and B. This gene regulatory model integrates experimental evidence of causal relationships among 15 genes (and the proteins they encode) that regulate cell-fate determination during floral organ specification in this plant. The loop between *Terminal Flower 1* (*TFL1*) and *Floral homeotic Apetala* 2 proteins disappears because the edge from the latter to the former is completely redundant. While FVC predicts that *TFL1* is required to control the network (to control this nonexistent loop), analysis of its STG (*SI Appendix, section 6*) reveals that *TFL1* can be replaced by the *AP1* (*Floral homeotic Apetala* 1) protein, which is not in this loop, to control the network. Interestingly, *AP1* is not in the set of driver nodes FVC predicts are needed for control. Similarly, the completely redundant interaction between *AP1* and *LFY* (*Leafy*) removes the loop between these two proteins, allowing *LFY* to be replaced by the *EMF1* (*Embryonic flower* 1) protein to control the network under



**Fig. 4.** Study of the *A. thaliana* BN model. (A) The interaction graph for the *A. thaliana* BN. (B) The effective graph. Edge thickness denotes effectiveness,  $e_{ij}$ ; dashed red indicates fully redundant edges (Table 1 shows parameter values); node color intensity denotes effective out-degree; and green nodes denote cases of null effective out-degree ( $k_e^{out} = 0$ ). (C) A threshold effective graph showing only edges with  $e_{ij} \geq 0.4$  to enhance visibility of the largest connected component that allows *LFY* to function as a master regulator and reveals that *WUS* functions simply as an autoregulator; green nodes denote cases of null effective out-degree at this threshold level. (D) Spearman's rank correlation (vertical axis) between the true impact of perturbing each node [ $\nu_{ij}(t)$ ] and respective path-length approximation predictions using the interaction (blue) and effective (red) graphs after  $t$  steps (horizontal axis); *FUL* cannot be computed (N/A, not available) because it has null impact on other variables (validating our observation of a fully redundant output).

50% effective connectivity, as seen in columns  $k_r^*$  and  $k_e^*$  in Table 1.

#### Effective Graph Aids Explanation of Biological Mechanism.

Structure-only control theories yield a set of driver nodes that are needed for control, but they do not provide a mechanistic explanation of how those nodes control the network or which nodes are more effective at control and signal propagation. The examples above demonstrate that the effective graph includes important dynamical redundancy information pertaining to the specific BN being analyzed. It reveals a more accurate portrait of how control operates, including alternative, actionable intervention strategies—such as the possibility of using *API* or *EMF1* instead of *TFL1* or *LFY*, respectively, in the set of driver nodes that control the TBN (by pinning).

Beyond identification of accurate driver variables, an analysis of the strongest paths of the effective graph reveals a more precise mechanistic understanding of how control propagates in biochemical regulation models. Consider the control roles of the *LFY* and *WUS* (*Wuschel*) transcription factor proteins in the Thaliana model. The most general form of BN control allows perturbations at any stage of the dynamics (to any configuration of the STG)—a more general form of control than the FVC pinning control assumptions (*SI Appendix, section 5*). In this case, via an STG enumeration method (20), we observe that the TBN is fully controllable by interventions to the trivial inputs {*UFO*, *LUG*, *CLF*} and additional driver set {*LFY*, *WUS*} alone. This makes sense because in the interaction graph in Fig. 4A, there is a path from *WUS* or *LFY* to any other node (except the three input nodes); so, in principle, signals from these nodes could reach any other node. However, in the effective graph in Fig. 4B, *WUS* is connected to the remainder of the network via a single very low-effectiveness edge with *AG* (*AG* transcription factor):  $e_{WUS,AG} = 0.1$ . Therefore, *WUS* is, in effect, dynamically decoupled from the remainder of the network. In contrast, *LFY* preserves paths with high edge effectiveness to all other nodes in the effective graph. The threshold effective graphs in Fig. 4C and *SI Appendix, Fig. S4* clarify the very distinct functional roles of these two proteins in the dynamics of this development model.

We validate these inferences with the analysis of perturbation spread on the TBN effective graph, as shown in Fig. 4D. The predictive power of the path-length approximation is very similar for both the interaction and effective graphs in the case of *LFY*, but it is completely different for *WUS* where the effective graph leads to a much higher correlation with the true impact of perturbing the latter variable. In other words, *WUS* does not behave at all like the original interaction graph would suggest. The effective graph reveals that these two transcription factors function very differently in how they control the TBN dynamics

**Table 1. Canalization parameters for variables with  $k \geq 2$  in the *A. thaliana* model (*SI Appendix, Table S4*)**

$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k^{out}$	$k_e^{out}$	$k_e^{out}/k^{out}$
AG	9	6.9	2.1	0.77	0.23	5	1.9	0.38
AP3	7	4.7	2.3	0.68	0.32	2	0.8	0.4
PI	6	3.8	2.2	0.64	0.36	2	0.47	0.24
AP1	4	2.4	1.6	0.59	0.41	6	1.4	0.23
LFY	4	2.8	1.2	0.69	0.31	7	4.8	0.69
TFL1	4	2.8	1.2	0.69	0.31	5	2.8	0.57
WUS	3	1.4	1.6	0.48	0.52	2	0.91	0.46
FUL	2	0.75	1.2	0.38	0.62	1	0	0

$k$ ,  $k_r$ , and  $k_e$  denote in-degree, input redundancy, and effective connectivity, respectively;  $k_r^*$  and  $k_e^*$  denote versions of  $k_r$  and  $k_e$  normalized by  $k$ ; and  $k^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively.

of this model. While *WUS* is only an autoregulator, *LFY* is a master regulator mechanism (44). Thus, even though the driver set for this network is {*LFY*, *WUS*}, except to control *WUS* itself, *LFY* is sufficient and a much more effective candidate for experimental intervention.

Also striking in the TBN effective graph is the case of the DNA-binding *FUL* protein. The interaction graph, built from published pairwise experiments, depicts that it causally affects *LFY*. However, this interaction is completely redundant in the model's logic for *LFY*. The *FUL* protein, therefore, has no impact in this model, as shown in the effective graph in Fig. 4B and confirmed by our perturbation analysis. Notice that because perturbations to *FUL* lead to null impact on other variables, we cannot compute Spearman's correlation to its predictive power for the interaction and effective graphs (hence, the N/A in Fig. 4D). Although the interaction graph implies that signals from *FUL* can reach almost all other variables in the model, the effective graph clearly reveals it reaches none.

We note that the effective graph is a probabilistic representation of the underlying dynamics, so even an edge with very low effectiveness may on rare occasions play a key role in determining dynamics. Still, statistically, edges with very low effectiveness are likely to play a reduced role in propagating control signals. This is demonstrated by the fact that the effectiveness-weighted paths of the effective graph are much more predictive of (correlated with) spreading dynamics after variable perturbation than paths of the original interaction graph, for both random graphs in Fig. 3 and the Thaliana network in Fig. 4D. Thus, strong paths in the effective graph are likely good control channels in systems biology models because they are better at propagating signals than other paths in the original interaction graph.

**Effective Graph Enhances Understanding of Signaling in Large Network Cancer Models.** The effective graph reveals multivariate canalizing dynamics by removing redundancy from automata networks and allows for a more precise characterization of perturbation and control signals. Since the effective graph is computed from the scalable schema redescription methodology, we can apply it to large networks for which full enumeration of the configuration space, and thus, computation of true control behavior or identification of all attractors, is not possible (17, 31). To demonstrate how it allows us to understand canalizing dynamics and identify effective control pathways, we study two large signal transduction networks involved in leukemia (45) and ER+ breast cancer (26), whose interaction and effective graphs are shown in Fig. 5 and *SI Appendix, Figs. S7–S13*.

The ER+ breast cancer network is a multistate automata network that has been converted to a fully equivalent, 80-variable BN (26). The goal of this model, built from experimental evidence, is to study resistance mechanisms to *PI3K* (phosphatidylinositol 3-kinase) inhibitors in *ER+*, *HER2+*, and *PIK3CA*-mutant breast cancer cells. Seven drugs that inhibit specific targets of interest are included in the model. For instance, alpelisib is a *PI3K* inhibitor (a drug that inhibits phosphoinositide 3-kinase enzymes involved in cell growth signaling pathways). The model is used to study known and novel combinatorial interventions that combine *PI3K* inhibition with other strategies (26). The objective is not so much to find the attractors of the entire multivariate dynamical system but simply to identify the final state of specific outcome nodes that model cancer cell death (apoptosis) or proliferation. In other words, the model is constructed to study which dynamical interventions control cancer cells to their programmed death or at least inhibit their proliferation.

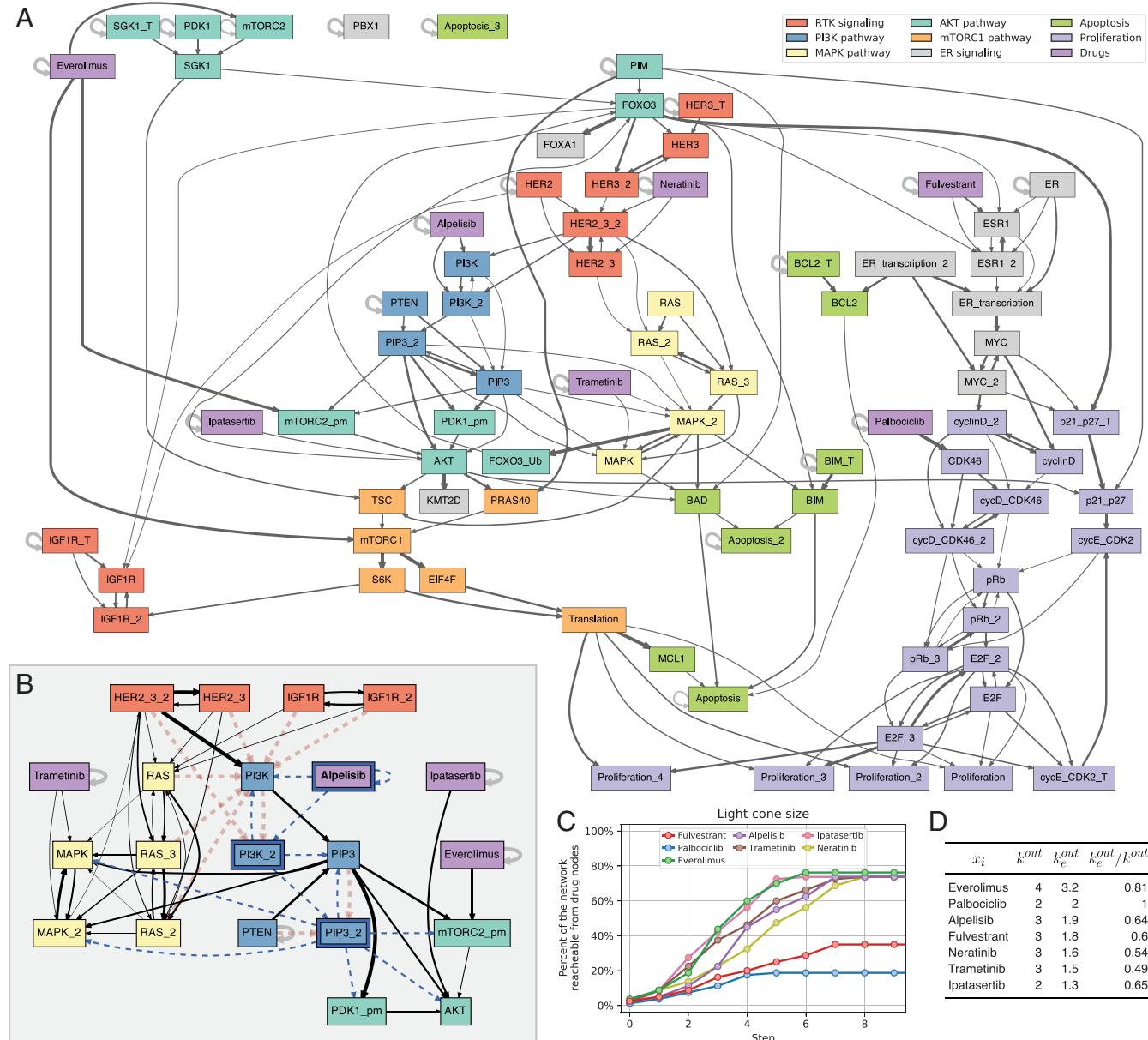
The effective graph of this model (*SI Appendix, Fig. S11*) reveals that much redundancy is present in its dynamics, and edge effectiveness is highly variable. Some interactions are

almost completely redundant with effectiveness as small as 0.065. The maximum edge effectiveness, 1.0, is only observed for automata with a single input, where redundancy cannot exist by definition. *SI Appendix, Table S7* shows canalization parameters for all of the variables in this model; Fig. 5D shows key parameters for the seven drugs included in this model.

The interaction graph (80 nodes) has 23 (29%) autoregulator (self-loop) nodes, of which 18 (23%) are input nodes (*SI Appendix, Fig. S10*). This means that a large proportion of nodes cannot be controlled via other nodes. Still, reachability is ultimately formed by a single weakly connected component and 45 strongly connected components, the largest of which has 24

nodes (*SI Appendix, Tables S2 and S3*). This implies that signals from input nodes could in principle reach the entire network via the weakly connected component and 30% of the nodes could regulate each other via the largest strongly connected component.

The effective graph, however, reveals a different, clearer understanding. The network dynamics is effectively separated into various modules, perhaps because the model is a synthesis of six pathways implementing distinct resistance mechanisms to PI3K inhibition that affect the apoptosis and proliferation pathways (26). Indeed, the most effective edges form very few connections among subsystems that can effectively propagate



**Fig. 5.** Study of the ER+ breast cancer BN model. (A) Hierarchical rendering of the effective graph for the BN model of ER+ breast cancer. Edge thickness denotes its effectiveness, thresholded to  $e_{ij} > 0.2$ ; node color denotes constituent pathways (legend is in the top right corner). (B) Conditional effective graph with Alpelisib=ON (pinned state denoted with bold text and blue border), revealing how it renders much of the influence from RTK (receptor tyrosine kinases) pathway redundant (red dashed edges) while fixing the state of several variables in the PI3K pathway, such as the phospholipid PIP3 (phospholipid); variables whose state becomes fixed (constants) are denoted by a blue border, and edges that transmit a constant input state are denoted by a dashed blue color. (C) Spreading dynamics of perturbations to each of the seven drugs in the model and the proportion of network effectively reachable. (D) Effectiveness of outgoing edges of drug variables;  $k_e^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively.

signals. This is shown by the existence of many strongly connected components (and input nodes) when only the reasonably effective edges are considered—especially in comparison with other models studied (*SI Appendix*, Tables S2 and S3). Consider the threshold effective graph with edge effectiveness greater than or equal to 0.2 shown in Fig. 5A. The largest strongly connected component is composed of only 17 (21%) nodes, and 45 (56%) nodes form strongly connected components of a single node. This means that there is little effective cross-regulation dynamics or long-range signaling in the model. Most of the effective dynamics can only be driven by direct intervention to many individual nodes or short pathways involving few nodes. The network becomes even more splintered at a threshold of 0.4 (*SI Appendix*, Figs. S12 and S13), resulting in the largest strongly connected component of only 3 (4%) nodes, with 61 (76%) nodes forming strongly connected components of a single node. Reachability is also quite diminished; for an edge effectiveness threshold of 0.4, there are 12 weakly connected components, the largest of which is composed of only 52 nodes (65% of the network). Indeed, two of the apoptosis nodes and one of the proliferation nodes, all key targets of the model, become isolated at effectiveness threshold 0.4, and one of the apoptosis nodes becomes isolated at effectiveness threshold 0.2.

These results are consistent with the known behavior of the model, whereby control of cancer apoptosis or inhibition of proliferation requires interventions to many nodes, including the *PI3K* inhibitor, other drugs, and every input node (26). The effective graph, however, reveals that the dynamics of this network is very robust to perturbation and hard to control because its subsystems are effectively decoupled. That is, canalization works by preventing propagation of signals and cross-regulation. Indeed, most of the (nondrug) variables that have an impact on cancer apoptosis or proliferation, when working in tandem with the *PI3K* inhibitor and baseline (table 3 in ref. 26), have short paths to those target variables (at most three edges) in the effective graph.

The connectivity of the effective graph thus reveals that the overall dynamics of the ER+ breast cancer network is very modular with many effectively decoupled subsystems—substantially more than the other experimentally validated biochemical models considered, as seen in *SI Appendix*, Tables S2 and S3. In contrast, in the TBN discussed above, canalization enables *LFY* to function as a single master regulator gene that can effectively propagate signals (effectiveness at or above 0.4). It reaches all other nodes in a large, weakly connected component of 12 nodes (80% of network), except for *WUS* (which remains in a decoupled component) and the input nodes. Similarly, in the case of the T cell survival in leukemia network (45), the effective graph maintains a single weakly connected component of 58 nodes (97% of the network), even for a high 0.4 effectiveness, which reveals a greater ability to propagate effective control signals through this network.

Let us now use the effective graph to study how differently the seven drugs are capable of controlling cancer cells to apoptosis or proliferation in this model. The goal of the original model is to find interventions—especially single-node interventions—that synergize with the *PI3K* inhibitor Alpelisib (26). Focusing on the remaining six drugs, the model reveals that Fulvestrant and Palbociclib best synergize with Alpelisib to increase apoptosis or decrease proliferation of cancer cells. Everolimus also modestly increases apoptosis, although not as much as the other two drugs. In contrast, Neratinib, Trametinib, and *Ipatercept* were shown to not synergize with Alpelisib (table 3 in ref. 26).

An initial observation of the effective graph, summarized in Fig. 5D, is consistent with those results: Alpelisib and the three drugs that best synergize with it are the top four with the largest effective out-degrees ( $k_e^{out}$ ). Thus, the most outwardly effective drugs are also those previously shown to lead to greatest control

of cancer apoptosis or proliferation. More importantly, the effective graph reveals why the seven drugs affect the cancer dynamics the way they do. The hierarchical rendering of the (0.2) threshold effective graph shown in Fig. 5A clearly reveals why Fulvestrant and Palbociclib synergize so well with the *PI3K* inhibitor Alpelisib: they act on the estrogen (*ER*) signaling and cell proliferation pathways that Alpelisib cannot effectively reach [except by indirectly reaching the terminal proliferation nodes via the *mTORC1* (mechanistic target of rapamycin complex 1) pathway]. This is demonstrated by comparing the conditional effective graph for Alpelisib= ON with those for a combined intervention Alpelisib=Fulvestrant= ON or Alpelisib=Palbociclib= ON (*SI Appendix*, Figs. S16, S18, and S21): only the combination interventions are capable of fully resolving the state of the proliferation variables. This explains why these drugs in combination with Alpelisib can drive cancer proliferation to zero in this model, while Alpelisib on its own cannot (table 3 in ref. 26). Moreover, Fulvestrant can also effectively reach some of the apoptosis pathway, which explains why, in combination with Alpelisib, it can increase apoptosis of cancer cells in this model but Palbociclib does not—the latter is only effective on the proliferation pathway and is not effective on apoptosis. These observations are also corroborated by the study of spreading perturbations. In Fig. 5C, we can see that Fulvestrant and Palbociclib reach a distinct, smaller part of the network, with Fulvestrant reaching more of the network (*ER* signaling, proliferation, and apoptosis pathways) than palbociclib (only the proliferation pathway).

The drugs that were shown not to synergize with Alpelisib (*Ipasertib*, *Neratinib*, and *Trametinib*) not only have the lowest values of  $k_e^{out}$  in Fig. 5D but are also shown in Fig. 5A and the respective conditional effective graphs (*SI Appendix*, Figs. S19, S20, and S22) to only contribute to the same pathways that Alpelisib already acts on. Fig. 5C also shows perturbing these three drugs ultimately spreads only to the same subgraph of the network that Alpelisib already acts upon. Indeed, the conditional effective graph for an intervention to Alpelisib alone (plus the baseline cancer-state input variables) shown in *SI Appendix*, Fig. S16 reveals that the drugs *Ipasertib*, *Neratinib*, and *Trametinib* are rendered completely redundant—interestingly, *Neratinib* is actually redundant even without Alpelisib but just with the *ER* + /*Her2*– cancer cell-state baseline (*SI Appendix*, Fig. S15). This highlights how the effective graph methodology provides an analytical explanation of the causal relationships in the model; one does not need to run ensemble Monte Carlo simulations of the BN model to know that *Ipasertib*, *Neratinib*, and *Trametinib* have no effect on apoptosis and proliferation of ER+ cancer cells in this model when Alpelisib is present.

Finally, the case of *Everolimus* is also well explained by the effective graph. While it is very outwardly effective (largest  $k_e^{out}$  in Fig. 5D), it also acts mostly on pathways already under downstream control by Alpelisib, as can be seen in Fig. 5A and in comparisons between the respective conditional effective graphs in *SI Appendix*, Figs. S16 and S17. Thus, while the simulations in ref. 26 report a very modest effect on apoptosis in synergy with Alpelisib ( $\approx 4\%$  increase), our results predict that, in this model, the effect on apoptosis of a combined *Alpelisib* + *Everolimus* intervention (Alpelisib alone) is causally negligible. It is noteworthy that *Everolimus* retains an effective edge to the *AKT* (protein kinase B) pathway (via *mTORC2*), providing some control of a subset of this pathway not under Alpelisib control (the top left in Fig. 5A and *SI Appendix*, Fig. 19). Indeed, the spreading dynamics experiments summarized in Fig. 5C show that perturbations to *Everolimus* spread just a little farther than perturbations to Alpelisib. *Everolimus* is therefore not as redundant to the overall dynamics as are *Ipasertib*, *Neratinib*, and *Trametinib*. Moreover, it preserves very effective pathways to both the *AKT* and *mTORC1*

pathways even at a high effectiveness threshold of 0.4 (*SI Appendix*, Fig. S13), which can play a part if Alpelisib becomes inactive.

It should be noted, similarly to the TBN model (Fig. 4D), that perturbation analysis of the *ER+* breast cancer network for the seven drugs studied shows that the effective graph is always more correlated with impact on dynamics than is the interaction graph (*SI Appendix*, Fig. S14). Therefore, the inferences derived above from the effective graph are grounded on a more realistic description of the model's true dynamics than inferences made directly from the interaction graph.

In summary, analysis of the effective graph and its dynamics provides a more complete understanding of why the seven drugs behave as reported in previous experiments with this model (26)—including why some are redundant. Removal of redundancy, furthermore, reveals analytically how canalization affects the mechanisms of apoptosis and proliferation of *ER+* cancer cells in this model. In particular, some drugs are more effective than others due to how decoupled from overall dynamics their pathways become. Indeed, the *ER+* breast cancer network is one of the most “fractured” of all of the experimentally validated biochemical models we studied—an issue we discuss in detail in *SI Appendix*, section 7 by studying their dynamical modularity via the analysis of strongly and weakly connected network components for all effectiveness threshold levels.

## Discussion and Conclusion

The effective graph we introduce synthesizes both the causal interaction structure and the nonlinear dynamics of BNs into a single scalable graph formalism. We use 78 experimentally validated BN models from systems biology to demonstrate that biochemical interactions contain significantly more redundancy than expected by chance, and this leads to very canalized nonlinear dynamics. This observation is consistent with Waddington's idea that canalization is pervasive in biological systems (25), whereby most random dynamical perturbations are not effective and only a few interactions control changes in network dynamics. This suggests that evolution in biological regulation has selected for redundancy, which has long been hypothesized as a requirement for the robustness to random perturbations that is necessary for evolvability (46, 47).

In addition to systems biology models, we use artificial models to show that effective graphs provide a more precise characterization of the (nonlinear) causal interaction logic of automata networks than do interaction graphs. These examples demonstrate that the effective graph is a better predictor of how perturbation signals propagate than is the original interaction graph, and thus, it is a useful construct to predict how control signals propagate. The effective graph can greatly aid the construction, refinement, and analysis of systems biology models by revealing how evidence from pairwise biochemical regulation experiments is integrated. Indeed, 22% of the biological models from the Cell Collective contain at least one fully redundant edge, and all contain much redundancy (19). Thus, the effective graph can aid in the simplification of biochemical network models to reveal their most essential regulatory pathways.

In comparison with the original automata networks, edge effectiveness reflects a loss of causal detail about which specific input combinations result in downstream variable-state changes. However, the effective graph is not proposed as a substitute for the causal interaction details that the original automata network contains. It is rather a revision of the original interaction graph that provides a much more precise, probabilistic accounting of causal dynamics and can be conditioned on different input assumptions with the conditional effective graph. Therefore, the loss of specific causal detail yields a powerful approach

for analysts who want to identify the most effective intervention strategies, those most likely to steer dynamics to desirable behavior.

Other methods have been proposed to integrate structure and dynamics into enhanced network representations. The general idea is to capture all of the possible roles that variables and interactions play in the logic of automata networks with additional formalism such as hyperedges (48) or distinct node types for variable states (49). The removal of redundancy via Boolean minimization can also be used to obtain parsimonious enhanced network representations, as shown in previous work (17). While these methods can preserve all possible causal interactions, even the rarest ones, they increase the complexity of the network representation. In contrast, the effective graph is a directed, weighted graph with a single node type, which is simpler and more amenable to the well-known graph-theoretical analysis and methods of network science (50). Moreover, the node- and edge-level effectiveness parameters are directly interpretable and provide an aggregate but accurate quantification of the causal pathways that are of greater interest for analysis and intervention in biochemical networks.

Without additional knowledge, our probabilistic characterization first assumes a uniform distribution over the likelihood of all input-state combinations to a given automaton. While this assumption is valid for automata in isolation, the presence of a biologically relevant subset of states or the convergence of dynamics onto attractors can alter the distribution of input states. The conditional effective graph allows us to explore such distinct input assumptions—as we do to study the causal roles of specific drugs in the cellular processes involved in *ER+* breast cancer. It also provides a promising direction for future work toward integrating the dynamically evolving likelihood of input states into a temporal effective graph.

Because we are interested in studying the (ontogenetic) dynamics of specific biochemical regulation systems, we focus on dynamical perturbations that change the state of biochemical variables. In future work, the effective graph is likely to be very useful to study the impact of structural perturbations, such as edge deletions or changes in logical transition rules (14). Indeed, one would expect greater dynamical disruptions from structural perturbations to effective pathways than to redundant pathways (20). Thus, our methodology can also be a tool to study the robustness and evolvability of function in biochemical networks—including developmental and disease control—especially in synergy with methods that hitherto have used only the original interaction graph (1, 38, 42, 43).

To demonstrate that the effective graph is useful in designing interventions in a specific systems biology model of development, disease, and biochemical regulation, we focus on the analysis of a small BN model of flower development, *A. thaliana*, as well as a large BN model of signal transduction in a model of *ER+* breast cancer. In these models, the effective graph allows us to demonstrate how different biochemical molecules or signals control dynamics. Whereas existing methods can identify driver variables that control dynamics, we show that by removing dynamical redundancy, the effective graph not only can help identify a more precise (smaller) set of driver variables but can also show how these variables function. For instance, our method distinguishes between an autoregulator gene (*WUS*) that is only needed to control itself and a master regulator gene (*LFY*) that controls most of the *A. thaliana* network. This enhanced explainability can also be used to reveal alternative, actionable control strategies, such as using *AP1* or *EMF1* instead of *TFL1* or *LFY* to control the Thaliana model.

Similarly, the effective graph of the *ER+* breast cancer model allows us to understand why and how some *PI3K* inhibitor drugs are more effective than others at controlling apoptosis or growth. Specifically, the methodology provides an analytical explanation

of the causal relationships that arise in the macrolevel network dynamics rather than observations from Monte Carlo simulations. This allows us to show analytically how Fulvestrant synergizes with Alpelisib to best control the *ER+* cancer cell line model, as well as why several drugs in the model are completely redundant. Such accurate explanations of how control interventions propagate throughout a biochemical system are important for the design of advanced disease therapeutics (39). Indeed, explainability is an important feature to derive actionable complex systems models in biomedicine and elsewhere. It can lead not only to model refinement (for example, by testing and potentially removing interactions predicted to be redundant) but also, to a deeper understanding of how causal, nonlinear, microlevel interactions integrate to define macrolevel biological functions. Our approach thus enhances understanding of multilevel complexity in biochemical regulation and multivariate dynamical systems at large.

1. J. G. T. Zanudo, G. Yang, R. Albert, Structure-based control of complex networks with nonlinear dynamics. *Proc. Natl. Acad. Sci. U.S.A.* **114**, 7234–7239 (2017).
2. S. Klamt, U. U. Haus, F. Theis, Hypergraphs and cellular networks. *PLoS Comput. Biol.* **5**, e1000385 (2009).
3. C. J. O. Reichhardt, K. E. Bassler, Canalization and symmetry in Boolean models for genetic regulatory networks. *J. Phys. Math. Theor.* **40**, 4339–4350 (2007).
4. A. J. Gates, D. M. Gygi, M. Kellis, A. L. Barabási, A wealth of discovery built on the human genome project—by the numbers. *Nature* **590**, 212–215 (2021).
5. J. M. Perez-Perez, H. Candela, J. L. Micó, Understanding synergy in genetic interactions. *Trends Genet.* **25**, 368–376 (2009).
6. E. Davidson, M. Levin, Gene regulatory networks. *Proc. Natl. Acad. Sci. U.S.A.* **102**, 4935 (2005).
7. S. Bornholdt, Boolean network models of cellular regulation: Prospects and limitations. *J. R. Soc. Interface* **5**, S85–S94 (2008).
8. F. Li, T. Long, Y. Lu, Q. Ouyang, C. Tang, The yeast cell-cycle network is robustly designed. *Proc. Natl. Acad. Sci. U.S.A.* **101**, 4781–4786 (2004).
9. T. Helikar et al., The cell collective: Toward an open and collaborative approach to systems biology. *BMC Syst. Biol.* **6**, 96 (2012).
10. G. Chechik et al., Activity motifs reveal principles of timing in transcriptional control of the yeast metabolic network. *Nat. Biotechnol.* **26**, 1251–1259 (2008).
11. M. Choi, J. Shi, Y. Zhu, R. Yang, K. H. Cho, Network dynamics-based cancer panel stratification for systemic prediction of anticancer drug response. *Nat. Commun.* **8**, 1940 (2017).
12. K. E. Kurten, Correspondence between neural threshold networks and Kauffman Boolean cellular automata. *J. Phys. Math. Gen.* **21**, L615 (1988).
13. R. Albert, J. Thakar, Boolean modeling: A logic-based dynamic approach for understanding signaling and regulatory networks and for making useful predictions. *Wiley Interdiscip. Rev. Syst. Biol. Med.* **6**, 353–369 (2014).
14. S. A. Kauffman, *The Origins of Order: Self-Organization and Selection in Evolution* (OUP USA, 1993).
15. C. Gershenson, Guiding the self-organization of random Boolean networks. *Theor. Biosci.* **131**, 181–191 (2012).
16. W. Marshall, H. Kim, S. I. Walker, G. Tononi, L. Albantakis, How causal analysis can reveal autonomy in models of biological systems. *Phil. Trans. Math. Phys. Eng. Sci.* **375**, 20160358 (2017).
17. M. Marques-Pita, L. M. Rocha, Canalization and control in automata networks: Body segmentation in *Drosophila melanogaster*. *PLoS One* **8**, e55946 (2013).
18. M. Aldana, Boolean dynamics of networks with scale-free topology. *Phys. Nonlinear Phenom.* **185**, 45–66 (2003).
19. S. Manicka, M. Marques-Pita, L. M. Rocha, Effective connectivity determines the critical dynamics of biochemical networks. arXiv[Preprint] (2021). <https://arxiv.org/abs/2101.08111> (Accessed 22 January 2021).
20. A. J. Gates, L. M. Rocha, Control of complex networks requires both structure and dynamics. *Sci. Rep.* **6**, 24456 (2016).
21. D. Bérenguier et al., Dynamical modeling and analysis of large cellular regulatory networks. *Chaos* **23**, 025114 (2013).
22. S. A. Kauffman, Metabolic stability and epigenesis in randomly constructed genetic nets. *J. Theor. Biol.* **22**, 437–467 (1969).
23. R. Thomas, Boolean formalization of genetic control circuits. *J. Theor. Biol.* **42**, 563–585 (1973).
24. S. Kauffman, C. Peterson, B. Samuelsson, C. Troein, Genetic networks with canalizing Boolean rules are always stable. *Proc. Natl. Acad. Sci. U.S.A.* **101**, 17102–17107 (2004).
25. C. H. Waddington, Canalization of development and the inheritance of acquired characters. *Nature* **150**, 563–565 (1942).
26. J. Zanudo, M. Scaltriti, R. Albert, A network modeling approach to elucidate drug resistance mechanisms and predict combinatorial drug treatments in breast cancer. *Canc. Converg.* **1**, 5 (2017).
27. C. Espinosa-Soto, P. Padilla-Longoria, E. R. Alvarez-Buylla, A gene regulatory network model for cell-fate determination during *Arabidopsis thaliana* flower development that is robust and recovers experimental gene expression profiles. *Plant Cell* **16**, 2923–2939 (2004).
28. Á. Chaos et al., From genes to flower patterns and evolution: Dynamic models of gene regulatory networks. *J. Plant Growth Regul.* **25**, 278–289 (2006).
29. E. J. McCluskey, Minimization of Boolean functions. *Bell. Syst. Tech. J.* **35**, 1417–1444 (1956).
30. I. Shmulevich, S. A. Kauffman, Activities and sensitivities in Boolean network models. *PRL* **93**, 048701 (2004).
31. R. B. Correia, A. J. Gates, X. Wang, L. M. Rocha, Cana: A python package for quantifying control and canalization in Boolean networks. *Front. Physiol.* **9** (2018).
32. R. James, J. Crutchfield, Multivariate dependence beyond Shannon information. *Entropy* **19**, 531 (2017).
33. A. Kolchinsky, A. Lourenco, H. Y. Wu, L. Li, L. M. Rocha, Extraction of pharmacokinetic evidence of drug-drug interactions from the literature. *PLoS One* **10**, e0122199 (2015).
34. C. Kadelka, J. Kuipers, R. Laubenbacher, The influence of canalization on the robustness of Boolean networks. *Phys. Nonlinear Phenom.* **353**, 39–47 (2017).
35. A. Kolchinsky, A. J. Gates, L. M. Rocha, Modularity and the spread of perturbations in complex dynamical systems. *Phys. Rev. E* **92**, 060801 (2015).
36. M. Santolini, A. L. Barabási, Predicting perturbation patterns from the topology of biological networks. *Proc. Natl. Acad. Sci. U.S.A.* **115**, E6375–E6383 (2018).
37. B. Luque, R. V. Solé, Lyapunov exponents in random Boolean networks. *Phys. Stat. Mech. Appl.* **284**, 33–45 (2000).
38. Y. Y. Liu, J. J. Slotine, A. L. Barabási, Controllability of complex networks. *Nature* **473**, 167–173 (2011).
39. R. Zhang et al., Network model of survival signaling in large granular lymphocyte leukemia. *Proc. Natl. Acad. Sci. U.S.A.* **105**, 16308–16313 (2008).
40. R. S. Wang, R. Albert, Elementary signaling modes predict the essentiality of signal transduction network components. *BMC Syst. Biol.* **5**, 44 (2011).
41. T. Akutsu, M. Hayashida, W. K. Ching, M. K. Ng, Control of Boolean networks: Hardness results and algorithms for tree structured networks. *J. Theor. Biol.* **244**, 670–679 (2007).
42. J. C. Nacher, T. Akutsu, Structural controllability of unidirectional bipartite networks. *Sci. Rep.* **3**, 1647 (2013).
43. B. Fiedler, A. Mochizuki, G. Kurosawa, D. Saito, Dynamics and control at feedback vertex sets. I. Informative and determining nodes in regulatory networks. *J. Dynam. Differ. Eq.* **25**, 563–604 (2013).
44. S. S. K. Chan, M. Kyba, What is a master regulator? *J. Stem Cell Res. Ther.* **3**, 1000e114 (2013).
45. A. Saadatpour et al., Dynamical and structural analysis of a T cell survival network identifies novel candidate therapeutic targets for large granular lymphocyte leukemia. *PLoS Comput. Biol.* **7**, e1002267 (2011).
46. M. Conrad, The geometry of evolution. *Biosystems* **24**, 61–81 (1990).
47. M. Pigliucci, Is evolvability evolvable? *Nat. Rev. Genet.* **9**, 75–82 (2008).
48. S. Klamt, J. Saez-Rodríguez, J. A. Lindquist, L. Simeoni, E. D. Gilles, A methodology for the structural and functional analysis of signaling and regulatory networks. *BMC Bioinf.* **7**, 56 (2006).
49. G. Yang, J. Gómez Tejeda Zanudo, R. Albert, Target control in logical models using the domain of influence of nodes. *Front. Physiol.* **9**, 454 (2018).
50. A. L. Barabási et al., *Network Science* (Cambridge University Press, 2016).

## Data and Code Availability

All simulations and data used to support the findings of this study are freely available in the CANA package (31) or the Cell Collective (9).

**Data Availability.** All study data are included in the article and/or *SI Appendix*.

**ACKNOWLEDGMENTS.** We thank Santosh Manicka and Manuel Marques-Pita for helpful discussions, Deborah Rocha for editing the manuscript, and Alice Grishchenko for graphics consultation. R.B.C. was partially funded by Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) Grant 18668127 and Fundação para a Ciência e a Tecnologia (FCT) Grant PTDC/MEC-AND/30221/2017. L.M.R. was partially funded by NIH, National Library of Medicine Grant 1R01LM012832; by a Fulbright Commission fellowship; and by National Science Foundation Research Traineeship “Interdisciplinary Training in Complex Networks and Systems” Grant 1735095. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.



1

<sup>2</sup> **Supplementary Information for**

<sup>3</sup> **The effective graph reveals redundancy, canalization, and control pathways in biochemical  
4 regulation and signaling**

<sup>5</sup> **Alexander J. Gates, Rion Brattig Correia, Xuan Wang, and Luis M. Rocha**

<sup>6</sup> **Alexander J. Gates and Luis M. Rocha.**

<sup>7</sup> **E-mails: [a.agates@northeastern.edu](mailto:a.agates@northeastern.edu) and [rocha@indiana.edu](mailto:rocha@indiana.edu)**

<sup>8</sup> **This PDF file includes:**

<sup>9</sup> **Supplementary text**

<sup>10</sup> **Figs. S1 to S24**

<sup>11</sup> **Tables S1 to S7**

<sup>12</sup> **References for SI reference citations**

13 **Supporting Information Text**

14 **1. Schema Redescription Theory**

15 **A. Boolean minimization and prime implicants.** The minimization of Boolean functions is a well known problem in electrical  
16 engineering and computer science. In both fields, the goal is typically to reduce the components and complexity needed to  
17 implement Boolean logic in electronic circuitry by removing redundancy from the computation of multivariate logical functions.  
18 The Quine-McCluskey algorithm (Q-M) is one of the best-known methods to remove the redundancy from logical functions  
19 (1, 2). Boolean minimization in the Q-M algorithm is achieved by uncovering the *prime implicants* (PI) of a Boolean function,  
20 i.e. an implication that resolves the logical value of the function (output) utilizing the least number of input variable states (or  
21 literals) possible. PI are thus minimal conditions (a conjunction of literals) to achieve automata state transition—minimal in  
22 the sense that the removal of any literal from the implication results in not knowing the output state (1). The first step of  
23 Q-M yields the set of all PI, or the *Blake Canonical Form* of a Boolean function—a disjunctive normal form of all PI. In our  
24 method, we represent PI by wildcard schemata, and the set of schemata that redescribe the automaton is a disjunction of all  
25 PI, equivalent to its Blake Canonical Form(3).

26 In the standard application of Q-M for Boolean minimization, the set of all PI is further reduced to uncover the *essential*  
27 *prime implicants*, which are PI that cover at least an entry of the LUT that cannot be covered by a combination of other PI.  
28 These essential PI are always needed to minimize a Boolean function. In contrast, redundant PI cover LUT entries that are  
29 already covered by a combination of essential PI. Finally, selective PI are neither essential nor redundant, and are selected at  
30 the end of Q-M to ensure coverage of all entries of the LUT. In the typical application of Boolean minimization, all redundant  
31 PI are removed and the final form of the Boolean function is a disjunction of all essential PI, plus selected PI needed to cover  
32 the entire LUT—selective PI are chosen to minimize the number of literals (maximizing wildcards).

33 In contrast to the engineering view of Boolean minimization, we are interested in characterizing the redundancy present in  
34 all possible interventions, even if some of those interventions can be built by combinations of other interventions. Therefore, we  
35 keep the entire set of PI, or the full Blake canonical form in schemata form. This allows us to preserve all possible *mechanisms*  
36 that change dynamical state in automata and are thus, in principle, amenable to control or perturbation interventions. In this  
37 sense, our goal is distinct from Boolean minimization of circuit design where the focus is on guaranteeing correct function  
38 computation with the least amount of wiring: to remove all redundancy in the computation of the full function. See (3) for  
39 additional details, including why the set of all PI is also necessary to infer symmetry constraints in schemata.

To appreciate why it is important to use all PI in the computation our proposed measures, consider the following example  
LUT of a Boolean function of  $k = 3$  inputs  $f(x_1, x_2, x_3)$ :

$$\begin{array}{l} 000 \rightarrow 0 \\ 001 \rightarrow 0 \\ \mathbf{010 \rightarrow 1} \\ 011 \rightarrow 0 \\ \mathbf{100 \rightarrow 0} \\ \mathbf{101 \rightarrow 1} \\ 110 \rightarrow 1 \\ 111 \rightarrow 1 \end{array}$$

Its schema redescription, or Blake canonical form with all PI is:

$$\begin{array}{l} 00\# \rightarrow 0 \\ \mathbf{0\#1 \rightarrow 0} \\ \#00 \rightarrow 0 \\ \#10 \rightarrow 1 \\ \mathbf{1\#1 \rightarrow 1} \\ 11\# \rightarrow 1 \end{array}$$

40 where the essential PI are shown in bold, with corresponding unique LUT entries also shown in bold in LUT. Note that a PI is  
41 essential if it redescribes at least one LUT entry that is not redescribed by any other PI (e.g.  $010 \rightarrow 1$  is only redescribed by  
42  $\#10 \rightarrow 1$ ). From this redescription, the measures from main text eqs. 3-5 yield:  $r_1 = r_2 = 3/8, r_3 = 1/4, e_1 = e_2 = 5/8, e_3 =$   
43  $3/4, a_1 = a_2 = a_3 = 1/2$ . If we considered the full Boolean minimization process, we would use (in this case) only the essential  
44 PI shown in bold, yielding:  $r_1 = r_2 = 1/2, r_3 = 0, e_1 = e_2 = 1/2, e_3 = 1, a_1 = a_2 = 1/2, a_3 = 1$ . Notably, using only the  
45 essential PI would lead to incorrect inferences about the result of possible interventions. First, the effectiveness of input  $x_3$   
46 would be inferred as maximal (1), when the non-essential PI  $00\# \rightarrow 0$  and  $11\# \rightarrow 1$  show that in reality it is possible to  
47 determine the state of  $f$  not knowing the state of  $x_3$ , or conversely, to control  $f$  with  $x_1$  and  $x_2$  alone (collective canalization).  
48 This can also be seen in traditional logical form; the minimized expression is  $f = (x_2 \wedge \neg x_3) \vee (x_1 \wedge x_3)$ , whereas adding the  
49 non-essential PI yields  $f = (x_2 \wedge \neg x_3) \vee (x_1 \wedge x_3) \vee (x_1 \wedge x_2)$  (showing an extra term without  $x_3$ .)

50 Our analysis with all PI reveals the fact that while  $x_3$  is more effective in determining the state of  $f$  ( $e_3 > e_1 = e_2$ ), it  
 51 has some redundancy ( $r_3 = 1/4$ ) and is therefore not fully effective ( $e_3 = 3/4$ ). Moreover, our formula (eq. 5 in main text)  
 52 for deriving activity is only accurate when using all PI (see proof below). Indeed, the correct activity for input  $x_3$  in this  
 53 automaton is  $a_3 = 1/2$  (only half the time this input changes its state leads  $f$  to change its output), not  $a_3 = 1$  as using only  
 54 the essential PI would have us infer. This example demonstrates that if our goal is to study which minimal interventions lead  
 55 to state changes (rather than minimizing Boolean functions), we must include all PI in schema redescription.

56 Finally, this example also serves to highlight the difference between edge effectiveness and activity. While the latter does  
 57 not distinguish the role of the three inputs ( $a_1 = a_2 = a_3 = 1/2$ ), the former characterizes input  $x_3$  as a little more effective at  
 58 changing the state of  $f$  than inputs  $x_1$  and  $x_2$  ( $e_1 = e_2 = 5/8, e_3 = 3/4$ ). Indeed, this is the only input that does not appear as  
 59 wildcard in any of the essential PI (appears always as literal), thus it is a little more effective (less redundant) than the other  
 60 too. Our measures distinguish this behavior because the LUT entries (000, 001, 110, and 111) redescribed by the non-essential  
 61 schemata/PI with wildcards in the  $x_3$  position (00# and 11#), can also be redescribed by essential PI which have no wildcards  
 62 in the  $x_3$  position (#00, 0#1, #10, and 1#1), and all other LUT entries are redescribed by schemata that also do not have  
 63 wildcards in the  $x_3$  position. In contrast, in the case of input variables  $x_1$  and  $x_2$ , there are two (out of eight) LUT entries  
 64 for each that can only be redescribed by essential PI/schemata that have wildcards in the  $x_1$  or  $x_2$  positions: 010 and 100  
 65 redescribed by #10 and #00 for  $x_1$ , and 011 and 101 redescribed by 0#1 and 1#1 for  $x_2$ , respectively. In summary, there are  
 66 more possibilities (PI) to intervene on the state of function  $f$  that must include  $x_3$  than is the case for the other two inputs,  
 67 which are thus a little more redundant as measured by our collective canalization measures, but not by activity. Naturally, the  
 68 node-level measures of collective canalization (eqs. 1-2 in main text) also paint a more accurate description of redundancy than  
 69 sensitivity:  $k_r(f) = 1, k_e(f) = 2, s(f) = 3/2$ . As can be seen by the PI/schemata, one always needs two inputs to determine  
 70 the state of  $f$ , so only one is redundant on average and the effective connectivity of the automaton is two inputs on average. In  
 71 contrast, because it does not account for collective canalization ( $k_c(f) = 1/2$ ), sensitivity posits that  $f$  is sensitive to only 1.5  
 72 inputs.

73 **B. Aggregation of prime implicant influence.** The per-node and per-edge measures of redundancy given by eqs 1 and 3 in main  
 74 paper, aggregate the redundancy of PI/schemata that redescribe each entry of a LUT for the entire function or per input,  
 75 respectively. The idea is to tally the redundancy of inputs of an automaton as conveyed by every PI/schema (as a possible  
 76 intervention strategy). In the analysis pursued here, the aggregation is computed via the average operator (avg). This is the  
 77 same as assuming that any PI/schema is a viable intervention possibility to change the state of automaton  $x$ . In (3) it was  
 78 shown that the aggregation for input redundancy is bound by using minimum and maximum operators. The lower (upper)  
 79 bound, obtained by substituting min (max) for avg in eq. 1, assumes that the schema that bests redescribes a given LUT entry  
 80 is the one with least (most) number of wildcards.

81 Because we have no reason to prioritize a PI/schema over another one that redescribes the same LUT entry, we use the  
 82 average operator, thus assuming that all PI are equally likely and important—as shown above, it is important to consider  
 83 the influence of all PI as possible interventions. Moreover, averaging over all possible schemata, allows the additive per-edge  
 84 separation of redundancy necessary for the effective graph (eqs. 3-4), whereby  $k_r(x_i) = \text{SUM}_j r_{ji}$  and  $k_e(x_i) = \text{SUM}_j e_{ji}$ , as well  
 85 as the clear relationship with activity and sensitivity (see proof below). One could consider a general case where a probability or  
 86 weighting of each PI is considered, e.g. assigning greater importance to essential PI, but at this point no advantage to pursuing  
 87 such a route was identified. The CANA package (4) defaults to calculating input redundancy and effective connectivity using  
 88 the average operator, but has the available functionality to use their lower or upper bounds as well.

89 **C. Activity and sensitivity via prime implicants.** Here we show that the activity of input  $x_j$  to automaton  $x_i$  can be defined  
 90 in terms of the schema redescription (the set of all PI) of  $x_i$  by eq. 5 in main text with the max operator in place of the  
 91 average operator (as used for edge effectiveness in eq. 3 in main text). To see this, let us review the definition of activity (5):  
 92  $a_j(x_i) = P(\overline{x_i^{t+1}} | \overline{x_j^t})$ , which is the probability that automaton  $x_i$  flips its state at  $t + 1$  when its input  $x_j$  flips its state at  $t$ ,  
 93 given a uniform distribution of input states at  $t$ . To compute the activity from the LUT of  $x_i$  we consider that for every entry  
 94  $f_\alpha$  in the LUT  $F_i$ , if flipping the state of input  $x_j$  (flipping the  $j$ th bit) leads  $x_i$  to change its state output, then that entry will  
 95 add  $1/|F_i| = 1/2^k$  to the activity. Summing over all LUT entries we have (5):

$$a_j(x_i) = \frac{1}{2^k} \sum_{f_\alpha \in F_i} \frac{\partial f_\alpha}{\partial x_j} \quad [1]$$

97 where  $\partial f_\alpha / \partial x_j = 1$  if flipping the  $j$ th bit of  $f_\alpha$  changes  $x_i$ , and 0 otherwise.

98 **Theorem 1.**

$$a_j(x_i) = \frac{1}{|F_i|} \sum_{f_\alpha \in F_i} \frac{\partial f_\alpha}{\partial x_j} = \frac{1}{|F_i|} \sum_{f_\alpha \in F_i} \left( 1 - \max_{v: f_\alpha \in \Upsilon_v^i} (j \rightarrow \#)_v \right) = 1 - \frac{1}{|F_i|} \sum_{f_\alpha \in F_i} \max_{v: f_\alpha \in \Upsilon_v^i} (j \rightarrow \#)_v \quad [2]$$

100 where  $(j \rightarrow \#)_v$  is a logical condition that is 1 / True if input  $x_j$  is a wildcard in schema  $f'_v$ , and 0 / False otherwise.

101 **Proof.** To prove the equality of equation 2, we have to show that:

$$102 \quad \frac{\partial f_\alpha}{\partial x_j} = 1 - \max_{v: f_\alpha \in \Upsilon_v^i} (j \rightarrow \#)_v \quad [3]$$

103 which can rephrased as an existence statement:

$$104 \quad \frac{\partial f_\alpha}{\partial x_j} = 0 \iff \exists v : f_\alpha \in \Upsilon_v^i \wedge (j \rightarrow \#)_v, \quad f_\alpha \in F_i. \quad [4]$$

105 In other words, LUT entry  $f_\alpha \in F_i$  contributes 0 to  $a_{ji}$  if and only if it is redescribed by at least one schema  $f'_v$  with a wildcard  
106 in its  $x_j$  position. Rephrased in terms of the PI, we have that LUT entry  $f_\alpha \in F_i$  contributes 0 to  $a_{ji}$  if and only if it is  
107 covered by at least one PI with no literals for the  $x_j$  logical variable.

108 Proving this statement, requires us to prove 3 simple lemmas.

109 **Lemma 1:** If LUT entry  $f_\alpha \in F_i$  is redescribed by at least one schema/PI  $f'_v$  with a wildcard in its  $x_j$  position, it will  
110 contribute 0 to activity ( $\partial f_\alpha / \partial x_j = 0$ ).

111 The proof of the lemma follows from the definition of partial derivative of an automaton (eq. 7 in main text). Assume  
112 that schema/PI  $f'_v$  redescribes  $f_\alpha \in F_i$  and has a wildcard in its  $x_j$  position. Then the state of  $x_j$  does not influence the  
113 configuration of literals (other inputs)  $f'_v$  specifies and cannot contribute to changing the state of  $x_i$ .

114 **Lemma 2,** the converse proposition of Lemma 1: if LUT entry  $f_\alpha \in F_i$  contributes 0 to the activity of input  $x_j$  to automaton  
115  $x_i$  ( $\partial f_\alpha / \partial x_j = 0$ ), then there must exist at least one schema/PI of  $x_i$ ,  $f'_v \in F'_i$ , which redescribes  $f_\alpha$  and has a wildcard in its  
116  $x_j$  position, that is,  $f_\alpha \in \Upsilon_v^i$ .

117 For an LUT entry  $f_\alpha = s_1 \dots s_j \dots s_k$  to have  $\partial f_\alpha / \partial x_j = 0$ , then we must have  $(s_1 \dots 0 \dots s_k \rightarrow s) \wedge (s_1 \dots 1 \dots s_k \rightarrow s)$ , which  
118 is equivalent to  $(s_1 \dots \# \dots s_k \rightarrow s)$ , where  $s$  is any truth value (or literal for the input variables). Though the implicant  
119  $(s_1 \dots \# \dots s_k \rightarrow s)$  might not be a PI itself, Lemma 3 below says there must be a schema/PI with a wildcard in position  $x_j$ .

120 **Lemma 3:** if LUT entry  $f_\alpha \in F_i$  can be covered by an implicant with a wildcard in its  $x_j$  position, then there must be at  
121 least one schema/PI having a wildcard in its  $x_j$  position.

122 This is straightforward considering the definition of prime implicant. Merging this implicant with other implicants to create  
123 a prime implicant will only add new wildcards without removing any existing wildcard. If the implicant cannot be covered by  
124 any prime implicant, then it is itself a prime implicant.

125 Combining all three lemmas, we prove our theorem.

126 The reader can notice that lemma 3 is contingent on using all possible PI. If only essential PIs were used, lemma 3 will not  
127 hold true as the example LUT in A shows.

## 128 2. Systems Biology Models

129 **A. Cell Collective Data set.** All experimentally-validated biochemical regulation and signalling BN models were retrieved from  
130 the Cell Collective (6) as of Aug 5<sup>th</sup> 2020 or were retrieved from literature and implemented in CANA(4). Table S1 shows the  
131 complete list of BN, including their respective Cell Collective and PubMed identifiers.

132 Note that two of the BN studied here contain a Boolean automata that is a full contradiction. That is, the Boolean automata  
133 is constant and does not depend on any of its inputs, despite having 4 or 8 specified inputs in its logical transition function.  
134 The presence of such logic irregularities further emphasizes the need for the effective connectivity methods presented here.

135 **B. Existence of Fully Redundant Interactions.** We found that 17 of the Cell Collective models (22%) contained at least one  
136 fully redundant edge, with 87 fully redundant edges in total. One possible explanation for the prevalence of fully redundant  
137 edges in gene regulatory or protein-interaction networks is that these models are often inferred from experimental data via  
138 information-theoretic measures, e.g. mutual information or transfer entropy (69), that can fail to discriminate between dyadic  
139 and polyadic relationships (70), and can thus miss the true multivariate dependency structure. But perhaps the main reasons  
140 are: 1) an incomplete record of experimental observation, 2) integration of experimental studies conducted by many different  
141 teams in different scenarios, and 3) modeling decisions about conflicting or weak evidence. This is especially problematic when  
142 the number of possible input combinations ( $k$ ) is large and experimentally testing all possible control conditions becomes  
143 unfeasible. Thus, fully redundant edges may be included because they refer to interactions that were not fully observed.  
144 Moreover, interactions observed in a given experimental setup may be subsequently rendered redundant when considering  
145 additional experimental controls or different thresholds for interaction significance; the reverse is also possible, whereby a  
146 redundant interaction is subsequently considered necessary with additional experimental evidence or change in criterion for  
147 strength of interaction.

148 This can be appreciated using the example in Figure 1 in main text. Let us imagine that it is assembled by integrating  
149 distinct interaction inference studies. Suppose that it describes how genes  $x_1, x_2$ , and  $x_3$  regulate the expression of gene  
150  $x_4$ . Imagine that one interaction study does not include the effect of  $x_3$ , revealing a rather strong interaction relationship:  
151  $x_4 = x_1 \wedge x_2$ . A separate study, on the other hand, does not control for the effect of  $x_1$ , observing not as clear an interaction  
152 effect of  $x_2$  and  $x_3$  on  $x_4$ . The study may reveal unequivocally that  $x_2 = 0 \Rightarrow x_4 = 0$  (per LUT entries  $f_1, f_2, f_5$  and  $f_6$  in  
153 Figure 1 in main text). But because  $x_1$  is not controlled in the study, correlational inference (e.g. via information-thererical  
154 measures (69)) is more uncertain as to how the expression of  $x_2$  and  $x_3$  affect the expression of  $x_4$ . It could be that  $x_1$  was

**Table S1. The 78 Boolean network models used in the analysis. Network names appear *ipsis litteris* for BN obtained from the Cell Collective.**

	Network	Cell col. ID	PMID	Ref
1	Thaliana flower development	-	15486106	(7, 8)
2	budding yeast cell cycle	-	15037758	(9)
3	ER+ breast cancer signal transduction	-	29623959	(10)
4	Signal Transduction in Fibroblasts	1557	18250321	(11)
5	Signaling in Macrophage Activation	1582	18433497	(12)
6	Mammalian Cell Cycle	1607	19118495	(13)
7	FA BRCA pathway	1778	22267503	(14)
8	HGF Signaling in Keratinocytes	1969	22962472	(15)
9	Cortical Area Development	2035	20862356	(16)
10	Death Receptor Signaling	2084	20221256	(17)
11	Yeast Apoptosis	2135	23233838	(18)
12	Cardiac development	2136	23056457	(19)
13	Guard Cell Abscisic Acid Signaling	2161	16968132	(20)
14	T Cell Receptor Signaling	2171	17722974	(21)
15	Cholesterol Regulatory Pathway	2172	19025648	(22)
16	T-LGL Survival Network 2008	2176	18852469	(23)
17	Neurotransmitter Signaling Pathway	2202	17010384	(24)
18	IL-1 Signaling	2214	21968890	(25)
19	Differentiation of T lymphocytes	2215	23743337	(26)
20	EGFR & ErbB Signaling	2309	19662154	(27)
21	IL-6 Signalling	2314	21968890	(25)
22	Apoptosis Network	2329	19422837	(28)
23	Body Segmentation in Drosophila 2013 (We uploaded this)	2341	23520449	(3)
24	B cell differentiation	2394	26751566	(29)
25	Mammalian Cell Cycle 2006	2396	16873462	(30)
26	Budding Yeast Cell Cycle	2404	23049686	(31)
27	T-LGL Survival Network 2011	2407	22102804	(32)
28	Budding Yeast Cell Cycle 2009	2423	19185585	(33)
29	Wg Pathway of Drosophila Signalling Pathways	2663	23868318	(34)
30	VEGF Pathway of Drosophila Signaling Pathway	2667	23868318	(34)
31	Toll Pathway of Drosophila Signaling Pathway	2668	23868318	(34)
32	Processing of Spz Network from the Drosophila Signaling Pathway	2669	23868318	(34)
33	Cell Cycle Transcription by Coupled CDK and Network Oscillators	2681	18463633	(35)
34	T-Cell Signaling 2006	2691	16464248	(36)
35	BT474 Breast Cell Line Long-term ErbB Network	2697	24970389	(37)
36	HCC1954 Breast Cell Line Long-term ErbB Network	2698	24970389	(37)
37	BT474 Breast Cell Line Short-term ErbB Network	2699	24970389	(37)
38	HCC1954 Breast Cell Line Short-term ErbB Network	2700	24970389	(37)
39	SKBR3 Breast Cell Line Short-term ErbB Network	2701	24970389	(37)
40	SKBR3 Breast Cell Line Long-term ErbB Network	2703	24970389	(37)
41	HIV-1 interactions with T Cell Signalling Pathway	2738	25431332	(38)
42	T cell differentiation	2901	16542429	(39)
43	Influenza A Virus Replication Cycle	3481	23081726	(40)
44	TOL Regulatory Network	3491	23171249	(41)

*Continues on the next page*

155 (unknowingly) more frequently expressed when testing condition  $x_2 = x_3 = 1$  (matching  $f_8$  more often than  $f_4$ ), but was more  
156 frequently inhibited when testing condition  $x_2 = 1 \wedge x_3 = 0$  (matching  $f_3$  more often than  $f_7$ ). This would lead the study  
157 to conclude that the most likely interaction is represented by  $x_4 = x_2 \wedge x_3$ , though with an observed weak effect due to how  
158 frequently  $x_1$  was expressed or not in the study for each condition. The weak effect would likely be reported and attributed to  
159 unknown causes of the expression of gene  $x_4$ —since  $x_1$  was not controlled in this second hypothetical study.

160 When systems biologists synthesize both studies into a network model, such as that of in Figure 1 in main text, decisions  
161 must be made. They may consider that the correct relationship is  $x_4 = (x_1 \wedge x_2) \vee (x_2 \wedge x_3)$ . This would result in the same  
162 LUT as Figure 1 in main text, except  $f_4 \equiv 011 : 1$ . But because the reported effect of the second experiment (second term) is  
163 much weaker than the observed effect in the first experiment (first term), the modelers may opt to characterize the synthesis of  
164 both experiments as  $x_4 = (x_1 \wedge x_2) \vee (x_1 \wedge x_2 \wedge x_3)$ , assuming that the observed unknown factor in the second experiment is  $x_1$   
165 which was tested in the first experiment—because they have no evidence of any other genes being involved in the expression  
166 of  $x_4$ . This decision would lead to the example in Figure 1 in main text, which by associativity and absorption is simply  
167  $x_4 = x_1 \wedge x_2$ , rendering  $x_3$  a redundant input.

168 A scenario that leads to a similar result is imagining a experimental study that does control for all three input genes, but

Table S1 - *Continued from previous page*

	Network	Cell col. ID	PMID	Ref
45	Bordetella bronchiseptica	3492	22253585	(42)
46	Trichostrongylus retortaeformis	3493	22253585	(42)
47	HH Pathway of Drosophila Signaling Pathways	3506	23868318	(34)
48	B bronchiseptica and T retortaeformis coinfection	3509	22253585	(42)
49	FGF pathway of Drosophila Signalling Pathways	3510	23868318	(34)
50	Glucose Repression Signaling 2009	3511	19144179	(43)
51	Oxidative Stress Pathway	3512	23134720	(44)
52	CD4 T cell signaling	3521	25538703	(45)
53	Colitis-associated colon cancer	4601	26446703	(46)
54	Septation Initiation Network	4705	26244885	(47)
55	Predicting Variabilities in Cardiac Gene	4706	26207376	(48)
56	PC12 Cell Differentiation	4775	27148350	(49)
57	Human Gonadal Sex Determination	4779	26573569	(50)
58	IGVH mutations in chronic lymphocytic leukemia.	4783	26088082	(51)
59	Fanconi anemia and checkpoint recovery	4790	26385365	(52)
60	Arabidopsis thaliana Cell Cycle	4837	26340681	(53)
61	Bortezomib Responses in U266 Human Myeloma Cells	4850	26163548	(54)
62	Stomatal Opening Model	4932	27542373	(55)
63	Pro-inflammatory Tumor Microenvironment in Acute Lymphoblastic Leukemia	4942	27594840	(56)
64	CD4+ T Cell Differentiation and Plasticity	5025	26090929	(26)
65	Lac Operon	5128	21563979	(57)
66	Metabolic Interactions in the Gut Microbiome	5731	26102287	(58)
67	Tumour Cell Invasion and Migration	5884	26528548	(59)
68	CD4+ T cell Differentiation	6678	22871178	(6)
69	Regulation of the L-arabinose operon of Escherichia coli.	6885	28639170	(60)
70	Aurora Kinase A in Neuroblastoma	7916	26616283	(61)
71	Iron acquisition and oxidative stress response in aspergillus fumigatus.	7926	25908096	(62)
72	MAPK Cancer Cell Fate Network	7984	24250280	(63)
73	Treatment of Castration-Resistant Prostate Cancer	8048	28361666	(64)
74	Lymphopoiesis Regulatory Network	8080	26408858	(65)
75	Lymphoid and myeloid cell specification and transdifferentiation	8186	28584084	(66)
76	T-LGL Survival Network 2011 Reduced Network	8227	22102804	(32)
77	Senescence Associated Secretory Phenotype	11863	29206223	(67)
78	Signaling Pathway for Butanol Production in Clostridium beijerinckii NRRL B-598	36604	30718562	(68)

where we observe very strong effects for some conditions and not others. For instance, imagine condition  $x_1 = 0 \wedge x_2 = 1 \wedge x_3 = 1$  leads to uncertain results about the expression of  $x_4$ ; say,  $x_4$  is expressed in only 60% of the experiments for that condition. In contrast, all other conditions lead to very certain observations of expression or inhibition of  $x_4$ . With this information, modelers, who have to decide on a acceptable threshold for evidence, may chose the relationship as  $x_4 = (x_1 \wedge x_2) \vee (x_1 \wedge x_2 \wedge x_3)$  (LUT in Figure 1 in main text with  $f_4 \equiv 011 : 0$ ) rather than  $x_4 = (x_1 \wedge x_2) \vee (x_2 \wedge x_3)$  (LUT with  $f_4 \equiv 011 : 1$ ).

Notice that in network inference, modelers often consider each node's LUT condition independently, as they can result from different experimental evidence. This means that the resulting LUTs are not necessarily further checked for logical redundancy or even incoherence (tautologies and contradictions). For instance, the *Thaliana* model (Figure 4, main text) contains three fully redundant edges, which ultimately result from “subjective decisions given alternatives with equivalent results” (7) regarding the expression of the *LFY* (*Leafy*) and *TFL1* (*Terminal Flower 1*) proteins. Certainly, our methodology to quantify redundancy in automata networks can also serve as an additional logical check on these models to avoid and understand the existence of completely redundant interactions derived from incomplete experimental evidence or modeling decisions.

### 3. Effectiveness Gini

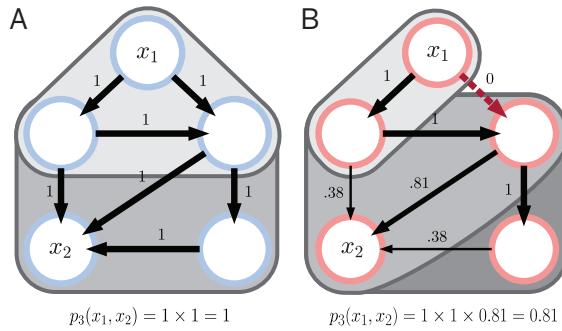
Given a vector of  $n$  values  $[x_1, x_2, \dots, x_n]$ , the Gini coefficient is defined as:

$$G = \frac{\sum_{i,j} |x_i - x_j|}{2n^2\bar{x}} \quad [5]$$

where  $\bar{x}$  is the mean of the values. Note that, due to finite-size effects, the maximal Gini coefficient for automata of degree  $k = 6$  is only 0.833, found when the vector is  $[0, 0, 0, 0, 0, 1]$ . As the number of samples increases, the maximal Gini coefficient approaches 1.

### 4. The effective graph and the spread of perturbations

Here we consider bit-flip perturbations to individual variables, defined as altering the logical state of the variable. In principle, the same framework can be used to study more elaborate classes of dynamic perturbations, including multi-variable bit flips or



**Fig. S1. The effective graph captures the spread of perturbations.** **A** In the  $\mathcal{M}_{\text{IG}}$  model, perturbations from automata  $x_1$  spread equally to all connected variables at  $t$  steps away in the interaction graph of an example BN. **B** In the  $\mathcal{M}_{\text{EG}}$  model, the effective graph of a BN constrains the spread of perturbations.

190 pinning perturbations. Changes to the system structure by changing variable transition functions through the addition or  
191 removal of an input are not considered.

192 The impact of a perturbation on an automaton in a BN is quantified by the Boolean analogue of the partial derivative (71):

$$193 \quad \partial_t^{(i)} x_j(\mathbf{x}_\alpha) = |x_j^t(\mathbf{x}_\alpha) - x_j^t(\mathbf{x}_\alpha^{\neg i})| \quad , \quad [6]$$

194 where  $x_j^t(\mathbf{x}_\alpha)$  denotes the state (truth value) of automaton node  $x_j$  at time  $t$  when the BN is initiated with configuration  
195  $\mathbf{x}^0 = \mathbf{x}_\alpha$  at time  $t = 0$ , and  $\mathbf{x}_\alpha^{\neg i}$  denotes configuration  $\mathbf{x}_\alpha$  with the state (truth value) of automaton  $x_i$  negated. In other words,  
196 the partial derivative yields 1 if flipping the state of  $x_i$  in initial configuration  $\mathbf{x}^0$  leads to  $x_j$  flipping its state at time  $t$ , and 0  
197 otherwise. The total impact on automaton  $x_j$  of perturbations to automaton  $x_i$  after  $t$  steps is found by averaging over all  
198 initial configurations:

$$199 \quad \iota_{ij}(t) = 2^{-N} \sum_{\alpha=1}^{2^N} \partial_t^{(i)} x_j(\mathbf{x}_\alpha). \quad [7]$$

200 For large BNs, the exact calculation of  $\iota_{ij}(t)$  becomes computationally infeasible and is approximated by averaging over a  
201 random sample of initial network configurations.

202 The interaction graph defines the light-cone of perturbation spreading, but it cannot differentiate the potential impact  
203 to nodes within the cone. Specifically, since signals between two nodes cannot travel faster than the minimum number of  
204 edges (thus, time steps  $t$ ) between them, the interaction graph provides an upper bound on the number of variables potentially  
205 affected by a perturbation after  $t$  time steps. We capture this upper bound in model  $\mathcal{M}_{\text{IG}}$ . The model  $\mathcal{M}_{\text{IG}}$  partitions the  
206 nodes into two groups: all nodes connected via a path of at most  $t$  edges starting from node  $x_j$  are equally impacted by a  
207 perturbation to node  $x_j$ , where  $t$  is the number of time steps since the perturbation, and all other nodes are not impacted by  
208 the perturbation. For simplicity, this model considers only the minimum path length from the perturbed node  $x_j$ , and does not  
209 account for loops or multiple paths.

210 The effective graph has at least two advantages for capturing the spread of perturbations: 1) it more accurately defines  
211 the light-cone of perturbation spreading since it removes interactions that are fully redundant, and 2) it provides edge weights  
212 that can differentiate the potential impact to nodes within the cone. The second model  $\mathcal{M}_{\text{EG}}$  thus ranks the node variables  
213 within the perturbation light-cone based on the weights along the most effective path from  $x_j$  to the variable. Specifically, the  
214 propensity for a perturbation to be transmitted along a path in  $\mathcal{M}_{\text{EG}}$  is given by the maximum product of edge strengths,  
215 constrained by the perturbation light-cone such that the number of edges in the path is less than the number of elapsed time  
216 steps. The maximum product path is calculated by finding the minimum additive path of negative log edge effectiveness.

217 Implementations of both models for perturbation spreading provided in the CANA python package(4).

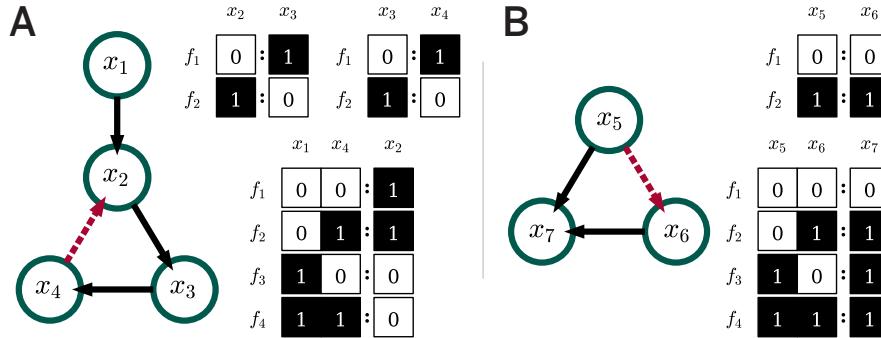
218 For the random BN experiment we generate 100 sample networks with  $N = 100$  nodes, and in which all nodes have an  
219 in-degree of 3, and average bias  $\bar{\rho} = 0.4$  (SI, S3). The network topology is randomized using a configuration model such that  
220 multi-edges are not allowed (each multi-edge is randomly swapped with another edge until no multi-edges remain). Self-loops  
221 are allowed. For each of the 100 randomly constructed networks, we sample 10 nodes at random to perturb. For each of the  
222 focus nodes, the dynamical impact on all other nodes is approximated using trajectories starting from  $10^4$  random initial  
223 configurations.

## 224 5. The effective graph improves structure-based control

225 The discovery of control strategies in BN models is a central problem in systems biology and biomedicine because predictions  
226 about controllability can help focus experimental interventions on genes, proteins and even medications more likely to result in  
227 the desired phenotype or medical outcome. Accurate control predictions would facilitate, for example, the design of advanced  
228 disease therapeutics(72, 73) or intervention strategies to reprogram cells(74), e.g. to revert a mutant cell to a wild-type state.

It is well known that when the set of automata nodes  $X$  of a BN is large, enumeration of all configurations  $\mathbf{x} \in \mathcal{X}$  of its STG becomes difficult, making the controllability of BN an NP-hard problem (75). Therefore, control methodologies which leverage the interaction graph or otherwise simplify the dynamics are highly desirable since they can greatly simplify the complexity of BN control(3, 73).

Two recent methodologies aim to determine the controllability of complex dynamical systems based solely from the graph of interactions between variables: structural controllability (SC)(76), and feedback-vertex set control (FVC)(73, 77). By using only the structural graph to predict minimum sets of variables that are needed to control a network (a.k.a. driver nodes), both of these methods make predictions about the entire ensemble of dynamical systems that fit the same interaction graph—rather than a specific multivariate dynamical system, such as a BN(78).



**Fig. S2. Fully redundant edges in the effective graph and structural control.** Two examples of BN in which canalization alters the predictions of structure-based control methods. **(A)**, A small BN with a fully canalized interaction between nodes  $x_4$  and  $x_2$  (dashed red). The predictions of SC and FVC on the structural graph suggest two nodes ( $x_1$  and one of  $x_2$ ,  $x_3$ , or  $x_4$ ) are required to control the network dynamics, while the same methods applied to the effective graph accurately identify that only node  $x_1$  is required to fully control the network. **(B)**, A small BN with a fully canalized interaction between nodes  $x_5$  and  $x_6$  (red). The predictions of SC, MDS and FVC on the structural graph suggest only node  $x_5$  is required to control the network dynamics, while the same methods applied to the effective graph accurately identify that both nodes  $x_5$  and  $x_6$  are required to fully control the network.

Here, we demonstrate that the effective graph is a more accurate representation of interactions between variables with important consequences for structural approximations of control, such as SC and FVC. Specifically, the redundancy of some logical functions means that the effective structure of interactions is reduced: fully-canalized edges of the structural graph play no role in determining the transitions between configurations. The presence of such canalization can both decrease or increase the estimated driver variable sets using structure-based control methods.

We illustrate how canalization can reduce the number of predicted driver variables using the BN shown in Fig. S2A. Both SC and FVC would predict that interventions on two nodes are required to control the system dynamics ( $x_1$  and one of  $x_2$ ,  $x_3$ , or  $x_4$ ). However, the interaction between nodes  $x_4$  and  $x_2$  is fully canalized (red edge), meaning that it should be disregarded by the structure-based methods. Applying both SC and FVC to the effective graph correctly reveals that only node  $x_1$  is required to control the system.

On the other hand, canalization can also increase the number of predicted driver variables as illustrated by the BN shown in Fig. S2B. In this case, both SC and FVC would predict that interventions on node  $x_5$  are sufficient to control the system dynamics. The effective graph reveals that the interaction between nodes  $x_5$  and  $x_6$  is full canalized (red edge); using the effective graph, SC and FVC correctly predict that both nodes  $x_5$  and  $x_6$  are required to control the system.

## 6. Effective graph and control

Similar insights about the control patterns from individual driver variables can be seen when comparing the minimum driver variable set for pinning controllability predicted by FVC (73, 77), with the real one obtained by full enumeration of all possible node sets in the STG (78). Recall that pinning controllability specifies a system can be controlled from any initial configuration to any of its attractors via “pinning” the driver variables to their state(s) in the target attractor (77). For the *Arabidopsis thaliana* BN model, in addition to the input nodes, FVC predicts the network to be pinning controllable with interventions to the 6 additional variables  $D_{FVC} = \{WUS, AP3, AG, TFL1, LFY, PI\}$ . However, enumeration of the STG (78) reveals that there are actually 3 equivalent minimum driver variable sets required for pinning control, each with only 5 additional variables:  $D_{pin} = \{WUS, AP3, AG\} \cup \{AP1, LFY\} \cup \{TFL1, EMF1\} \cup \{TFL1, LFY\}$ . The effective graph reveals why FVC overestimates this minimum set, and why a multiplicity of sets are equally effective in pinning control. First, note that all of the self-loops to  $PI$ ,  $AP3$ , and  $AG$  have negligible edge effectiveness, and thus do not need to be controlled because its self-loop is negligible and edges with stronger effectiveness exist from  $LFY$ ,  $AP3$  and  $AG$  that can control it (see Fig. 4C).

The effective graph reveals that some loops are removed entirely with fully redundant edges (i.e.  $TFL1 \leftrightarrow AP2$  and  $AP1 \leftrightarrow LFY$ ), or almost removed with very low effectiveness edges (e.g. the  $PI$  and  $AG$  self-loops), which help explain the differences between the FVC-predicted and real pinning control driver sets. Several interesting observations derive by looking at the effective graph and comparing  $D_{FVC}$  with  $D_{Pin}$ . Recall that FVC theory requires that all loops, including

variable self-loops, need to be controlled by at least one variable that interrupts the loop. The fully redundant edge between *AP2* and *TFL1* removes the loop between these variables which means that *TFL1* is not always necessary to control the network and *AP1* with *LFY* can control *TFL1* and *AP2*. Similarly, in the real pinning control driver set, *EMF1* can take the place of *LFY*, which is not allowed by the FVC prediction. The effective graph shows that this happens because the fully redundant edge between *AP1* and *LFY* removes the loop between these variables, and so *LFY* is not always needed to control the network—since all the loops with effective edges ( $e_{ij} > 0$ ) in which *LFY* participates can be interrupted by pinning *TFL1* and *EMF1*.

These observations demonstrate that while FVC predicts the necessary driver set to control the entire ensemble of BN that fit the same interaction graph (Fig. 4A), the effective graph of a specific BN can reduce the size of the necessary driver set and help identify alternative control strategies and most important variables to control, as is the case of the effective graph of the TBN model (Fig. 4B). These features make the effective graph useful to analyze control propagation in BN systems biology models, by providing more specific understanding of the effective pathways to control dynamics. This suggests that applying FVC to the effective graph can lead to more accurate (pinning) control predictions. The development of such a method is beyond the scope of this article, but we can at least demonstrate that in the case of the TBN model the effective graph with a threshold of  $e_{ij} \geq 0.2$  (Fig. S4) best explains the real pinning control driver set. If we apply FVC to this graph the result is  $D_{Pin} - AG$ . That is, it only misses the *AG* variable, which is needed for pinning control. A closer inspection of the full effective graph (Fig. 4B) reveals that *AG* participates in three low-effectiveness loops\* that are not interrupted by other variables in  $D_{Pin}$ . In contrast, *PI* which is not needed for real control (not in  $D_{Pin}$  though in  $D_{FVC}$ ), participates in a single low-effectiveness loop (its own self-loop) that is not interrupted by other variables in  $D_{Pin}$ . This strongly suggests that the number of low-effectiveness loops may be cumulative and needs to be accounted for, a hypothesis we will address in future work.

Similarly, in the real pinning control driver set, *EMF1* can take the place of *LFY*, which is not allowed by the FVC prediction. The effective graph shows that this happens because the fully redundant edge between *AP1* and *LFY* removes the loop between these variables, and so *LFY* is not always needed to control the network—since all the loops with effective edges ( $e_{ij} > 0$ ) in which *LFY* participates can be interrupted by pinning *TFL1* and *EMF1*.

In summary, we need to control *LFY* because it can control every other node, but we need to control *WUS* only for its own sake. This type of information is very useful for considering intervention strategies in Biology. If *WUS*'s final state is not of high importance, we can control most of the network by intervening only on *LFY*: its impact/power on the rest of the network is much higher.

## 7. Effective graph reveals dynamically-decoupled modules

The analysis of perturbation spread in complex dynamical systems can reveal dynamical modules that constrain this spread (79). Such modules are often related to well-known pathways with important roles in biochemical regulation and signalling. As shown for the ER+ breast cancer and TBN models, by thresholding the effective graph and eliminating edges with small effectiveness, we reveal subgraphs with greater dynamical influence, as well as those that are more or less decoupled from the rest of the network. One way to characterize such dynamical modularity is to compute the strongly and weakly connected components of the effective graph for various effectiveness thresholds, and compare them to the interaction graph; Tables S2 & S3 in SI show such an analyses for the four networks studied in detail above. Strongly connected components reveal modules where every node can in principle perturb every other node in same module, and weakly connected components those where some of the (driver) nodes can perturb all nodes in module. The size of such components can also reveal how fractured the dynamics of a network is.

As shown in Fig. S23A, B, the 78 biochemical BN models from the Cell Collective vary in how dynamically-decoupled modules arise as we change the effectiveness threshold. For example, the TBN model (green) splits into several weakly connected components at a relatively low effectiveness value ( $e_{ij} = 0.1$ ), but the largest component contains most of the network for a wide range of effectiveness values, at least 80% for  $e_{ij} \leq 0.5$ —which demonstrates the existence of the single primary dynamical module driven by the *LFY* protein described above. In contrast, the *ER+* Breast Cancer model (orange) fractures into many small components at  $e_{ij} \approx 0.4$ , with the largest weakly connected component comprising less than 20% of the network—which is coherent with the patchwork composition of several dynamically distinct modules discussed above.

Overall, analysis of the 78 network models in the Cell Collective reveals that for edge effectiveness  $e_{ij} \leq 0.2$  or even  $e_{ij} \leq 0.4$ , the majority of networks remain connected in a single or largest weakly connected component comprised of most nodes. As shown in Fig. S24, for  $e_{ij} \leq 0.2$ , about 90% of the networks have a largest weakly connected component comprised of at least 80% of the network. For  $e_{ij} \geq 0.4$ , on the other hand, most networks quickly lose a substantial largest weakly connected component, and for  $e_{ij} \geq 0.65$ , no networks have a largest weakly connected component comprised of even 70% of the network. That is, most networks break into many small components when  $e_{ij} \in [0.4, 0.6]$ . In this sense, connectivity, signal transmission, and dynamical control in these networks tends to be robust to removal of edges with  $e_{ij} \leq 0.2$  even up to  $e_{ij} \leq 0.4$ . This suggests  $e_{ij} \in [0.2, 0.4]$  is an optimum range for effectiveness, whereby redundant edges are removed but effective edges remain to reliably send signals through mostly connected networks.

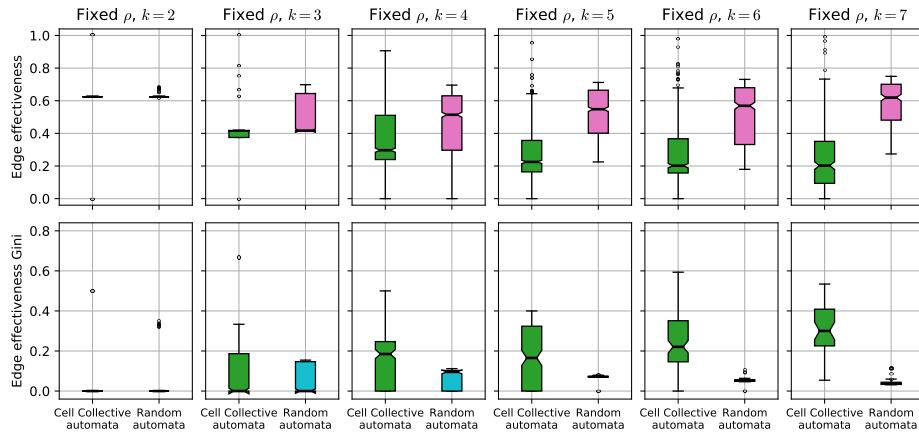
\*  $AG \longleftrightarrow WUS$ , and  $AG \longleftrightarrow AP1$ ,  $AG \longleftrightarrow AG$ .

324 **References**

- 325 1. Quine WV (1955) A Way to Simplify Truth Functions. *American Mathematical Monthly* 62:627–631.
- 326 2. McCluskey EJ (1956) Minimization of boolean functions. *The Bell System Technical Journal* 35(6):1417–1444.
- 327 3. Marques-Pita M, Rocha LM (2013) Canalization and control in automata networks: body segmentation in Drosophila melanogaster. *PloS ONE* 8(3):e55946.
- 328 4. Correia RB, Gates AJ, Wang X, Rocha LM (2018) Cana: A python package for quantifying control and canalization in boolean networks. *Frontiers in Physiology* 9.
- 329 5. Shmulevich I, Kauffman SA (2004) Activities and sensitivities in boolean network models. *PRL* 93(4):048701.
- 330 6. Helikar T, et al. (2012) The cell collective: Toward an open and collaborative approach to systems biology. *BMC Systems Biology* 6:96.
- 331 7. Espinosa-Soto C, Padilla-Longoria P, Alvarez-Buylla ER (2004) A gene regulatory network model for cell-fate determination during Arabidopsis thaliana flower development that is robust and recovers experimental gene expression profiles. *The Plant Cell Online* 16(11):2923–2939.
- 332 8. Chaos Á, et al. (2006) From Genes to Flower Patterns and Evolution: Dynamic Models of Gene Regulatory Networks. *Journal of Plant Growth Regulation* 25(4):278–289.
- 333 9. Li F, Long T, Lu Y, Ouyang Q, Tang C (2004) The yeast cell-cycle network is robustly designed. *PNAS* 101:4781–4786.
- 334 10. Zañudo J, Scaltriti M, Albert R (2017) A network modeling approach to elucidate drug resistance mechanisms and predict combinatorial drug treatments in breast cancer. *Cancer Convergence* 1(1).
- 335 11. Helikar T, Konvalina J, Heidel J, Rogers JA (2008) Emergent decision-making in biological signal transduction networks. *Proceedings of the National Academy of Sciences* 105(6):1913–1918.
- 336 12. Raza S, et al. (2008) A logic-based diagram of signalling pathways central to macrophage activation. *BMC Systems Biology* p. 15.
- 337 13. Sahin Ö, et al. (2009) Modeling ERBB receptor-regulated G1/S transition to find novel targets for de novo trastuzumab resistance. *BMC Systems Biology* p. 20.
- 338 14. Rodríguez A, et al. (2012) A boolean network model of the FA/BRCA pathway. *Bioinformatics* 28(6):858–866.
- 339 15. Singh A, Nascimento JM, Kowar S, Busch H, Boerries M (2012) Boolean approach to signalling pathway modelling in hgf-induced keratinocyte migration. *Bioinformatics* 28(18):i495–i501.
- 340 16. Giacomantonio CE, Goodhill GJ (2010) A boolean model of the gene regulatory network underlying mammalian cortical area development. *PLoS Computational Biology* 6(9):13.
- 341 17. Calzone L, et al. (2010) Mathematical modelling of cell-fate decision in response to death receptor engagement. *PLoS Computational Biology* 6(3):15.
- 342 18. Kazemzadeh L, Cvijovic M, Petranovic D (2012) Boolean model of yeast apoptosis as a tool to study yeast and human apoptotic regulations. *Frontiers in physiology* 3:446.
- 343 19. Herrmann F, Groß A, Zhou D, Kestler HA (2012) A boolean model of the cardiac gene regulatory network determining first and second heart field identity. *PLOS ONE* 7(10):10.
- 344 20. Li S, Assmann SM (2006) Predicting essential components of signal transduction networks: A dynamic model of guard cell abscisic acid signaling. *PLoS Biology* 4(10):17.
- 345 21. Saez-Rodriguez J, et al. (2007) A logical model provides insights into t cell receptor signaling. *PLoS Computational Biology* 3(8):11.
- 346 22. Kervizic G, Corcos L (2008) Dynamical modeling of the cholesterol regulatory pathway with boolean networks. *BMC Systems Biology* p. 14.
- 347 23. Zhang R, et al. (2008) Network model of survival signaling in large granular lymphocyte leukemia. *Proceedings of the National Academy of Sciences* 105(42):16308–16313.
- 348 24. Gupta S, Bisht SS, Kukreti R, Jain S, Brahmachari SK (2007) Boolean network analysis of a neurotransmitter signaling pathway. *Journal of Theoretical Biology* p. 7.
- 349 25. Ryll A, Samaga R, Schaper F, Alexopoulos LG, Klamt S (2011) Large-scale network models of IL-1 and IL-6 signalling and their hepatocellular specification. *Molecular BioSystems* p. 18.
- 350 26. Martinez-Sanchez ME, Mendoza L, Villarreal C, Alvarez-Buylla ER (2015) A minimal regulatory network of extrinsic and intrinsic factors recovers observed patterns of CD4+ t cell differentiation and plasticity. *PLoS Computational Biology* p. 23.
- 351 27. Samaga R, Saez-Rodriguez J, Alexopoulos LG, Sorger PK (2009) The logic of EGFR/ErbB signaling: Theoretical properties and analysis of high-throughput data. *PLoS Computational Biology* 5(8):19.
- 352 28. Mai Z, Liu H (2009) Boolean network-based analysis of the apoptosis network: Irreversible apoptosis and stable surviving. *Journal of Theoretical Biology* p. 10.
- 353 29. Méndez A, Mendoza L (2016) A network model to describe the terminal differentiation of b cells. *PLoS Computational Biology* p. 26.
- 354 30. Fauré A, Naldi A, Chaouiya C, Thieffry D (2006) Dynamical analysis of a generic boolean model for the control of the mammalian cell cycle. *Bioinformatics* 22(14):e124–e131.
- 355 31. Todd RG (2012) Ergodic sets as cell phenotype of budding yeast cell cycle. *PLOS ONE* 7(10):10.
- 356 32. Saadatpour A, et al. (2011) Dynamical and structural analysis of a t cell survival network identifies novel candidate therapeutic targets for large granular lymphocyte leukemia. *PLoS Comput Biol* 7(11):e1002267.

- 385 33. Irons DJ (2009) Logical analysis of the budding yeast cell cycle. *Journal of Theoretical Biology* p. 17.
- 386 34. Mbodj A, Junion G, Brun C, Furlong EEM, Thieffry D (2013) Logical modelling of drosophila signalling pathways.  
387 *Molecular BioSystems* 9(9):2248.
- 388 35. Orlando DA, et al. (2008) Global control of cell-cycle transcription by coupled CDK and network oscillators. *Nature* 453:5.
- 389 36. Klamt S, Saez-Rodriguez J, Lindquist JA, Simeoni L, Gilles ED (2006) A methodology for the structural and functional  
390 analysis of signaling and regulatory networks. *BMC Bioinformatics* p. 26.
- 391 37. der Heyde S, et al. (2014) Boolean ErbB network reconstructions and perturbation simulations reveal individual drug  
392 response in different breast cancer cell lines. *BMC Systems Biology* 8(1):75.
- 393 38. Oyeyemi OJ, Davies O, Robertson DL, Schwartz JM (2015) A logical model of hiv-1 interactions with the t-cell activation  
394 signalling pathway. *Bioinformatics* 31(7):1075–1083.
- 395 39. Mendoza L, Xenarios I (2006) A method for the generation of standardized qualitative dynamical systems of regulatory  
396 networks. *Theoretical Biology and Medical Modelling* p. 18.
- 397 40. Madrahimov A, Helikar T, Kowal B, Lu G, Rogers J (2013) Dynamics of influenza virus and human host interactions  
398 during infection and replication cycle. *Bulletin of Mathematical Biology* 75(6):988–1011.
- 399 41. Silva-Rocha R, de Lorenzo V (2013) The tol network of *p seudomonas putida* mt-2 processes multiple environmental  
400 inputs into a narrow response space. *Environmental microbiology* 15(1):271–286.
- 401 42. Thakar J, Pathak AK, Murphy L (2012) Network model of immune responses reveals key effectors to single and co-infection  
402 dynamics by a respiratory bacterium and a gastrointestinal helminth. *PLoS Computational Biology* 8(1):19.
- 403 43. Christensen TS, Oliveira AP, Nielsen J (2009) Reconstruction and logical modeling of glucose repression signaling pathways  
404 in *saccharomyces cerevisiae*. *BMC Systems Biology* p. 15.
- 405 44. Sridharan S, Layek R, Datta A, Venkatraj J (2012) Boolean modeling and fault diagnosis in oxidative stress response.  
406 *BMC Genomics* 13(Suppl 6):S4.
- 407 45. Conroy BD, et al. (2014) Design, assessment, and in vivo evaluation of a computational model illustrating the role of cavl  
408 in cd4+ t-lymphocytes. *Frontiers in immunology* 5:599.
- 409 46. Lu J, et al. (2015) Network modelling reveals the mechanism underlying colitis-associated colon cancer and identifies novel  
410 combinatorial anti-cancer targets. *Scientific reports* 5:14739.
- 411 47. Chasapi A, et al. (2015) An extended, boolean model of the septation initiation network in *s.pombe* provides insights into  
412 its regulation. *PLOS ONE* p. 22.
- 413 48. Grieb M, et al. (2015) Predicting variabilities in cardiac gene expression with a boolean network incorporating uncertainty.  
414 *PLOS ONE* p. 15.
- 415 49. Offermann B, et al. (2016) Boolean modeling reveals the necessity of transcriptional regulation for bistability in pc12 cell  
416 differentiation. *Frontiers in genetics* 7:44.
- 417 50. Ríos O, et al. (2015) A boolean network model of human gonadal sex determination. *Theoretical Biology and Medical  
418 Modelling* 12(1):1–18.
- 419 51. Álvarez Silva MC (2015) Proteins interaction network and modeling of IgVH mutational status in chronic lymphocytic  
420 leukemia. *Theoretical Biology and Medical Modelling* p. 15.
- 421 52. Rodríguez A (2015) Fanconi anemia cells with unrepaired DNA damage activate components of the checkpoint recovery  
422 process. *Theoretical Biology and Medical Modelling* p. 22.
- 423 53. Ortiz-Gutiérrez E, García-Cruz K, Azpeitia E, Castillo A (2015) A dynamic gene regulatory network model that recovers  
424 the cyclic behavior of *arabidopsis thaliana* cell cycle. *PLoS Computational Biology* p. 28.
- 425 54. Chudasama VL, Ovacik MA, Abernethy DR, Mager DE (2015) Logic-based and cellular pharmacodynamic modeling of  
426 bortezomib responses in u266 human myeloma cells. *Journal of Pharmacology and Experimental Therapeutics* 354(3):448–  
427 458.
- 428 55. Gan X (2016) Analysis of a dynamic model of guard cell signaling reveals the stability of signal propagation. *BMC Systems  
429 Biology* p. 14.
- 430 56. Enciso J, Mayani H, Mendoza L, Pelayo R (2016) Modeling the pro-inflammatory tumor microenvironment in acute  
431 lymphoblastic leukemia predicts a breakdown of hematopoietic-mesenchymal communication networks. *Frontiers in  
432 physiology* 7:349.
- 433 57. Veliz-Cuba A, Stigler B (2011) Boolean models can explain bistability in the lac operon. *Journal of computational biology*  
434 18(6):783–794.
- 435 58. Steinway SN, Biggs MB, Loughran TP, Papin JA, Albert R (2015) Inference of network dynamics and metabolic interactions  
436 in the gut microbiome. *PLoS Computational Biology* p. 25.
- 437 59. Cohen DPA, et al. (2015) Mathematical modelling of molecular pathways enabling tumour cell invasion and migration.  
438 *PLoS Computational Biology* p. 29.
- 439 60. Jenkins A, Macauley M (2017) Bistability and asynchrony in a boolean model of the l-arabinose operon in *escherichia coli*.  
440 *Bulletin of mathematical biology* 79(8):1778–1795.
- 441 61. Dahlhaus M (2016) Boolean modeling identifies greatwall/MASTL as an important regulator in the AURKA network of  
442 neuroblastoma. *Cancer Letters* p. 11.
- 443 62. Brandon M, Howard B, Lawrence C, Laubenbacher R (2015) Iron acquisition and oxidative stress response in *aspergillus  
444 fumigatus*. *BMC Systems Biology* p. 18.
- 445 63. Grieco L, Calzone L, Bernard-Pierrot I (2013) Integrative modelling of the influence of MAPK network on cancer cell fate

- 446 decision. *PLOS Computational Biology* 9(10):15.
- 447 64. Arshad OA (2017) Towards targeted combinatorial therapy design for the treatment of castration-resistant prostate cancer.  
*BMC Bioinformatics* p. 11.
- 448 65. Mendoza L (2015) A dynamical model of the regulatory network controlling lymphopoiesis. *Biosystems* p. 8.
- 449 66. Collombet S, et al. (2017) Logical modeling of lymphoid and myeloid cell specification and transdifferentiation. *Proceedings  
450 of the National Academy of Sciences* 114(23):5792–5799.
- 451 67. Meyer P, et al. (2017) A model of the onset of the senescence associated secretory phenotype after dna damage induced  
452 senescence. *PLoS computational biology* 13(12):e1005741.
- 453 68. Patakova P, et al. (2019) Acidogenesis, solventogenesis, metabolic stress response and life cycle changes in clostridium  
454 beijerinckii nrrl b-598 at the transcriptomic level. *Scientific reports* 9(1):1–21.
- 455 69. Bitbol AF (2018) Inferring interaction partners from protein sequences using mutual information. *PLoS Computational  
456 Biology* 14(11):e1006401.
- 457 70. James R, Crutchfield J (2017) Multivariate dependence beyond shannon information. *Entropy* 19(10):531.
- 458 71. Luque B, Solé RV (2000) Lyapunov exponents in random Boolean networks. *Physica A: Statistical Mechanics and its  
459 Applications* 284(1-4):33–45.
- 460 72. Zhang R, et al. (2008) Network model of survival signaling in large granular lymphocyte leukemia. *PNAS* 105:16308–16313.
- 461 73. Zanudo JGT, Yang G, Albert R (2017) Structure-based control of complex networks with nonlinear dynamics. *PNAS*  
462 114(28):7234–7239.
- 463 74. Wang RS, Albert R (2011) Elementary signaling modes predict the essentiality of signal transduction network components.  
*BMC Systems Biology* 5.
- 464 75. Akutsu T, Hayashida M, Ching WK, Ng MK (2007) Control of Boolean networks: hardness results and algorithms for  
465 tree structured networks. *Journal of Theoretical Biology* 244(4):670–679.
- 466 76. Liu YY, Slotine JJ, Barabási AL (2011) Controllability of complex networks. *Nature* 473:167–173.
- 467 77. Fiedler B, Mochizuki A, Kurosawa G, Saito D (2013) Dynamics and control at feedback vertex sets. i: Informative and  
468 determining nodes in regulatory networks. *Journal of Dynamics and Differential Equations* 25(3):563–604.
- 469 78. Gates AJ, Rocha LM (2016) Control of complex networks requires both structure and dynamics. *Scientific Reports* 6:24456.
- 470 79. Kolchinsky A, Gates AJ, Rocha LM (2015) Modularity and the spread of perturbations in complex dynamical systems.  
*Phy. Rev. E* 92(6):060801.
- 471
- 472
- 473



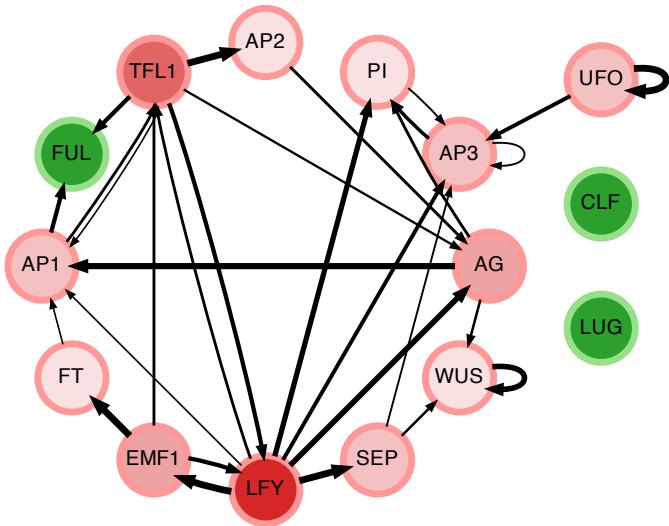
**Fig. S3. Edge effectiveness of Boolean automata in biochemical regulation and random ensembles of varying degree.** Top) Edge effectiveness of the 240 incoming edges (interactions) to 40 automata with varying degree in Cell Collective models (green) compared to a bias-matched sample of random Boolean automata (pink). Bottom) The Gini coefficient of the effectiveness of the incoming edges to the 40 automata with varying degree (green) in the Cell Collective models compared to the bias-matched ensemble of random Boolean automata (blue).

		Arabidopsis thaliana			
		IG	EG	EG 0.2	EG 0.4
nodes		15			
self-loops		5 (33%)			
input nodes		3 (20%)	3 (20%)	3 (20%)	4 (27%)
weakly conn. comp.		1 (15, 100%)	1 (15, 100%)	3 (13, 87%)	4 (12, 80%)
' strongly conn. comp.		5 (10, 67%)	6 (9, 60%)	8 (7, 47%)	10 (6, 40%)
		Saccharomyces cerevisiae (yeast)			
		IG	EG	EG 0.2	EG 0.4
nodes		12			
self-loops		8 (67%)			
input nodes		1 (8%)	1 (8%)	1 (8%)	1 (8%)
weakly conn. comp.		1 (12, 100%)	1 (12, 100%)	1 (12, 100%)	1 (12, 100%)
strongly conn. comp.		3 (10, 83%)	3 (10, 83%)	3 (10, 83%)	3 (10, 83%)
		Leukemia			
		IG	EG	EG 0.2	EG 0.4
nodes		60			
self-loops		11 (18%)			
input nodes		6 (10%)	6 (10%)	6 (10%)	10 (17%)
weakly conn. comp.		1 (60, 100%)	1 (60, 100%)	1 (60, 100%)	2 (58, 97%)
strongly conn. comp.		12 (48, 80%)	12 (48, 80%)	27 (29, 48%)	47 (9, 15%)
		Breast Cancer			
		IG	EG	EG 0.2	EG 0.4
nodes		80			
self-loops		23 (29%)			
input nodes		18 (23%)	18 (23%)	21 (26%)	29 (36%)
weakly conn. comp.		1 (80, 100%)	1 (80, 100%)	3 (78, 98%)	12 (52, 65%)
strongly conn. comp.		45 (24, 30%)	45 (24, 30%)	52 (17, 21%)	70 (3, 4%)

**Table S2. Structural characteristics of four biochemical regulation networks. Number (and proportion) of nodes that are self-loops and inputs. Number of weakly and strongly connected components that exist for each graph; for each case, also shown in brackets is the number of nodes in largest component, followed by the proportion of network in largest component.**

{size: number of comp.}		Thaliana			
		IG	EG	EG 0.2	EG 0.4
weakly conn. comp.		{15: 1}	{15: 1}	{13: 1, 1: 2}	{12: 1, 1: 3}
strongly conn. comp.		{10: 1, 2: 1, 1: 3}	{9: 1, 2: 1, 1: 4}	{7: 1, 2: 1, 1: 6}	{6: 1, 1: 9}
		Saccharomyces cerevisiae (yeast)			
		IG	EG	EG 0.2	EG 0.4
weakly conn. comp.		{12: 1}	{12: 1}	{12: 1}	{12: 1}
strongly conn. comp.		{10: 1, 1: 2}	{10: 1, 1: 2}	{10: 1, 1: 2}	{10: 1, 1: 2}
		Leukemia			
		IG	EG	EG 0.2	EG 0.4
weakly conn. comp.		{60: 1}	{60: 1}	{60: 1}	{58: 1, 2: 1}
strongly conn. comp.		{48: 1, 2: 1, 1: 10}	{48: 1, 2: 1, 1: 10}	{29: 1, 4: 1, 2: 2, 1: 23}	{9: 1, 4: 1, 2: 2, 1: 43}
		Breast Cancer			
		IG	EG	EG 0.2	EG 0.4
weakly conn. comp.		{80: 1}	{80: 1}	{78: 1, 1: 2}	{52: 1, 13: 1, 4: 1, 3: 1, 1: 8}
strongly conn. comp.		{24: 1, 8: 1, 2: 5, 1: 38}	{24: 1, 8: 1, 2: 5, 1: 38}	{17: 1, 8: 1, 2: 5, 1: 45}	{3: 1, 2: 8, 1: 61}

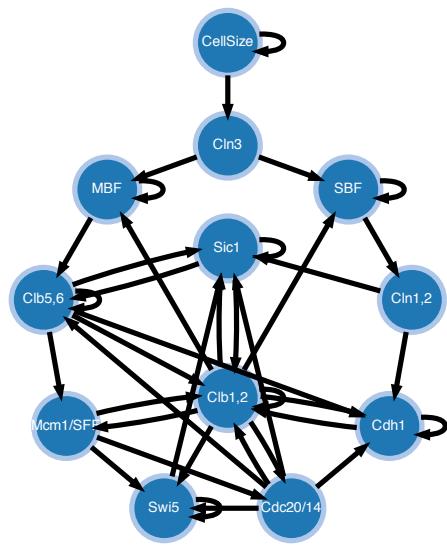
Table S3. Size and Number of structural components of four biochemical regulation models. For each type, components are listed in decreasing order of size (number of nodes), with the number of components of that size shown after each ':'.



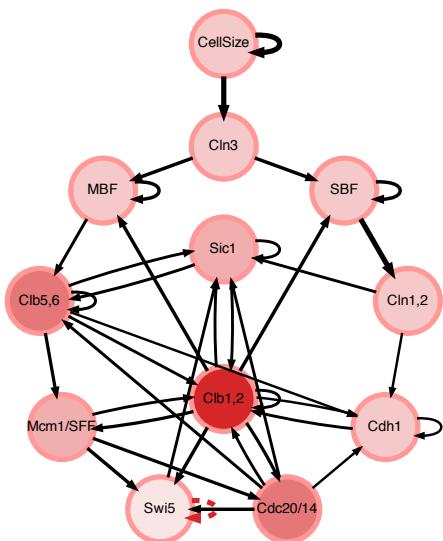
**Fig. S4. The effective graph in the *Arabidopsis thaliana* BN model.** The effective graph for the BN model of the *Arabidopsis thaliana*, in which edge thickness denotes its effectiveness, thresholded to  $e_{ji} > 0.2$ ; node color intensity denotes the node effective out-degree; green nodes denote variables with no effective out-degree.

$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k^{out}$	$k_e^{out}$	$k_e^{out}/k^{out}$
<i>AG</i>	9	6.9	2.1	0.77	0.23	5	1.9	0.38
<i>AP3</i>	7	4.7	2.3	0.68	0.32	2	0.8	0.4
<i>PI</i>	6	3.8	2.2	0.64	0.36	2	0.47	0.24
<i>AP1</i>	4	2.4	1.6	0.59	0.41	6	1.4	0.23
<i>LFY</i>	4	2.8	1.2	0.69	0.31	7	4.8	0.69
<i>TFL1</i>	4	2.8	1.2	0.69	0.31	5	2.8	0.57
<i>WUS</i>	3	1.4	1.6	0.48	0.52	2	0.91	0.46
<i>FUL</i>	2	0.75	1.2	0.38	0.62	1	0	0
<i>UFO</i>	1	0	1	0	1	2	1.6	0.79
<i>FT</i>	1	0	1	0	1	1	0.24	0.24
<i>EMF1</i>	1	0	1	0	1	3	2	0.68
<i>AP2</i>	1	0	1	0	1	2	0.43	0.22
<i>SEP</i>	1	0	1	0	1	4	0.9	0.22
<i>LUG</i>	0	0	1	0	0	1	0.1	0.1
<i>CLF</i>	0	0	1	0	0	1	0.1	0.1

**Table S4.** Canalization measures for variables in the *Arabidopsis thaliana* model.  $k$ ,  $k_r$ , and  $k_e$  denote in-degree, input redundancy and effective connectivity, respectively;  $k_r^*$  and  $k_e^*$  denote versions of  $k_r$ , and  $k_e$  normalized by  $k$ ;  $k^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively. Nodes with  $k = 1$  have no redundancy ( $k_r = 0$ ,  $k_e = 1$ ), and input nodes have no incoming edges.



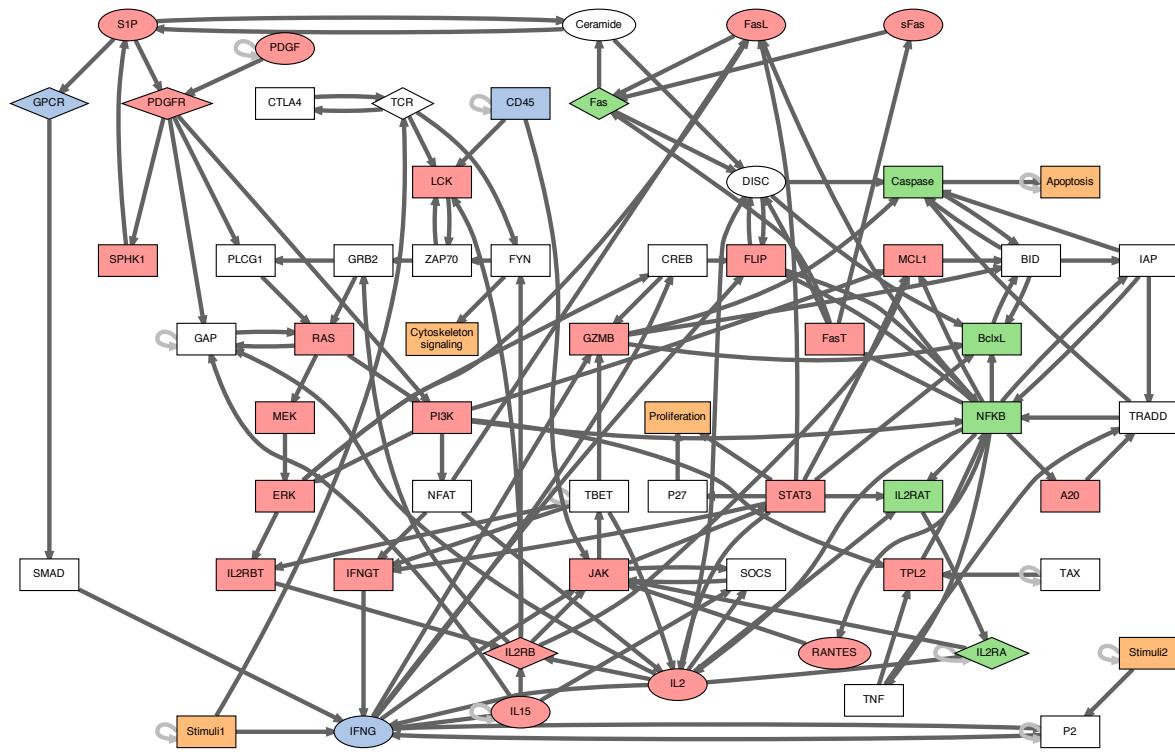
**Fig. S5. The interaction graph in the *Saccharomyces cerevisiae* (yeast) BN model.** The interaction graph for the BN model of *Saccharomyces cerevisiae*.



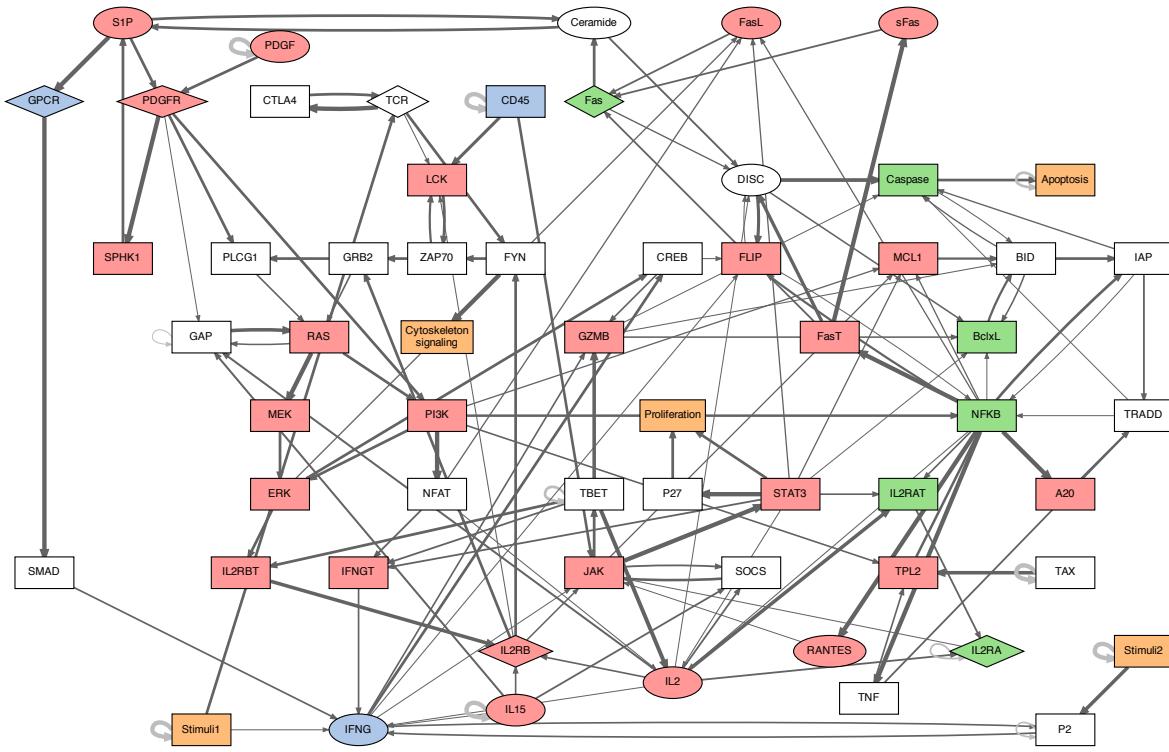
**Fig. S6. The effective graph in the *Saccharomyces cerevisiae* (yeast) BN model.** The effective graph for the BN model of *Saccharomyces cerevisiae*, in which edge thickness denotes its effectiveness,  $e_{ji}$ . Notice the fully redundant self-loop (dashed red) edge on the *Swi5* transcription factor node.

**Table S5. Canalization Measures for Variables in the *Saccharomyces cerevisiae* (yeast) model.**  $k$ ,  $k_r$ , and  $k_e$  denote in-degree, input redundancy and effective connectivity, respectively;  $k_r^*$  and  $k_e^*$  denote versions of  $k_r$ , and  $k_e$  normalized by  $k$ ;  $k_e^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively. Nodes with  $k = 1$  have no redundancy ( $k_r = 0$ ,  $k_e = 1$ ), and input nodes have no incoming edges.

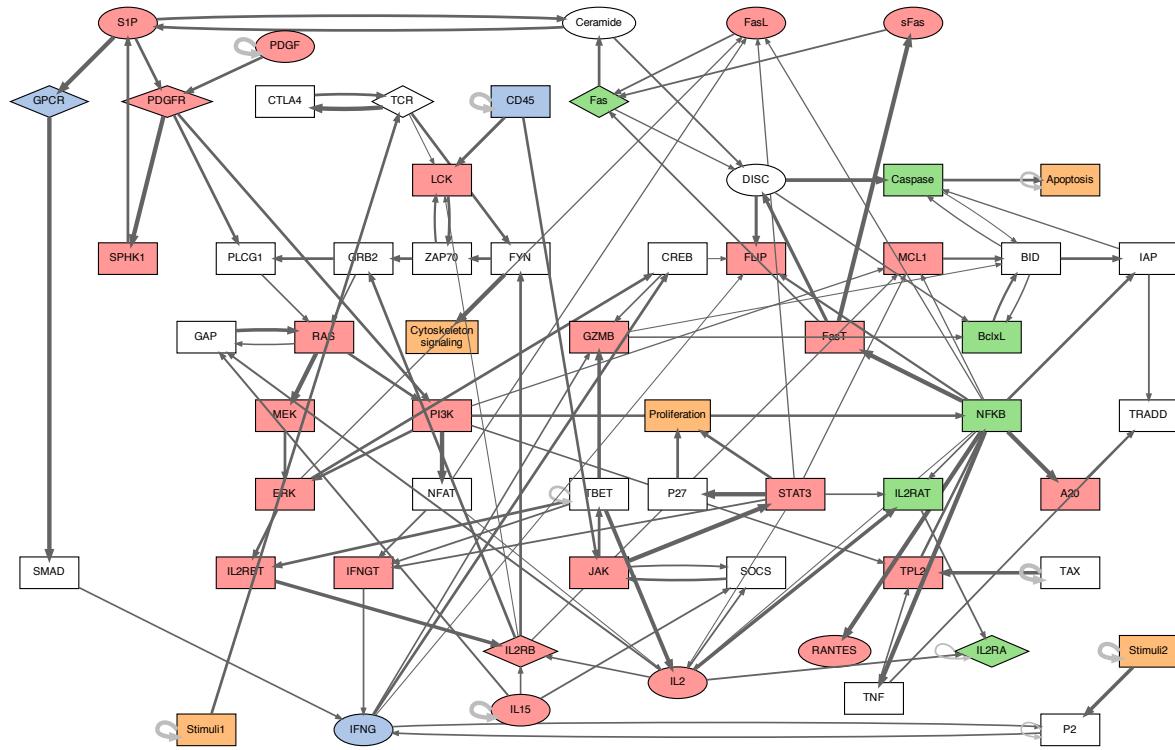
$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k_e^{out}$	$k_e^{out}/k_e^{out}$	
Sic1	6	2.7	3.3	0.44	0.56	3	1.7	0.56
Clb1,2	6	2.7	3.3	0.44	0.56	8	4.8	0.6
Cdh1	5	2.6	2.4	0.53	0.47	2	1	0.52
Clb5,6	4	1.7	2.3	0.42	0.58	5	2.8	0.56
Swi5	4	2	2	0.5	0.5	2	0.56	0.28
SBF	3	1	2	0.33	0.67	2	1.7	0.83
MBF	3	1	2	0.33	0.67	2	1.2	0.62
Mcm1/SFF	2	0.75	1.2	0.38	0.62	3	1.8	0.62
Cdc20/14	2	0.75	1.2	0.38	0.62	5	2.8	0.57
CellSize	1	0	1	0	1	2	2	1
Cln3	1	0	1	0	1	2	1.3	0.67
Cln1,2	1	0	1	0	1	2	1	0.52



**Fig. S7. The interaction graph in the Leukemia BN model.** The interaction graph for the BN model of leukemia.



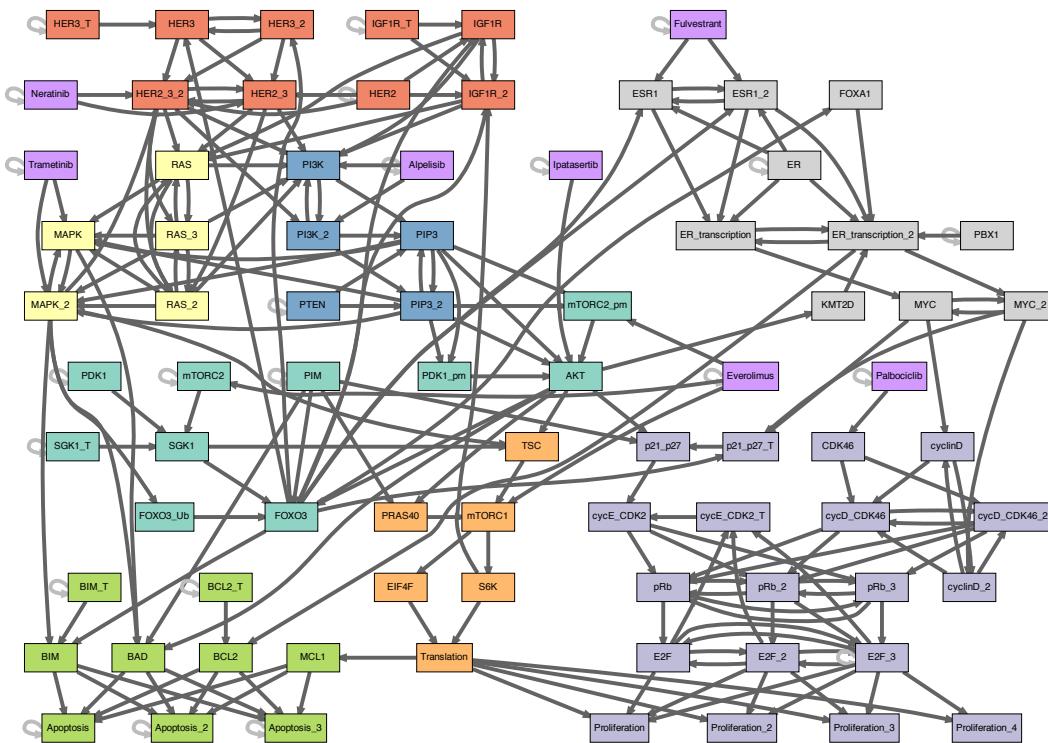
**Fig. S8. The effective graph in the Leukemia BN model.** The effective graph for the BN model of leukemia, in which edge thickness denotes its effectiveness,  $e_{ji}$ .



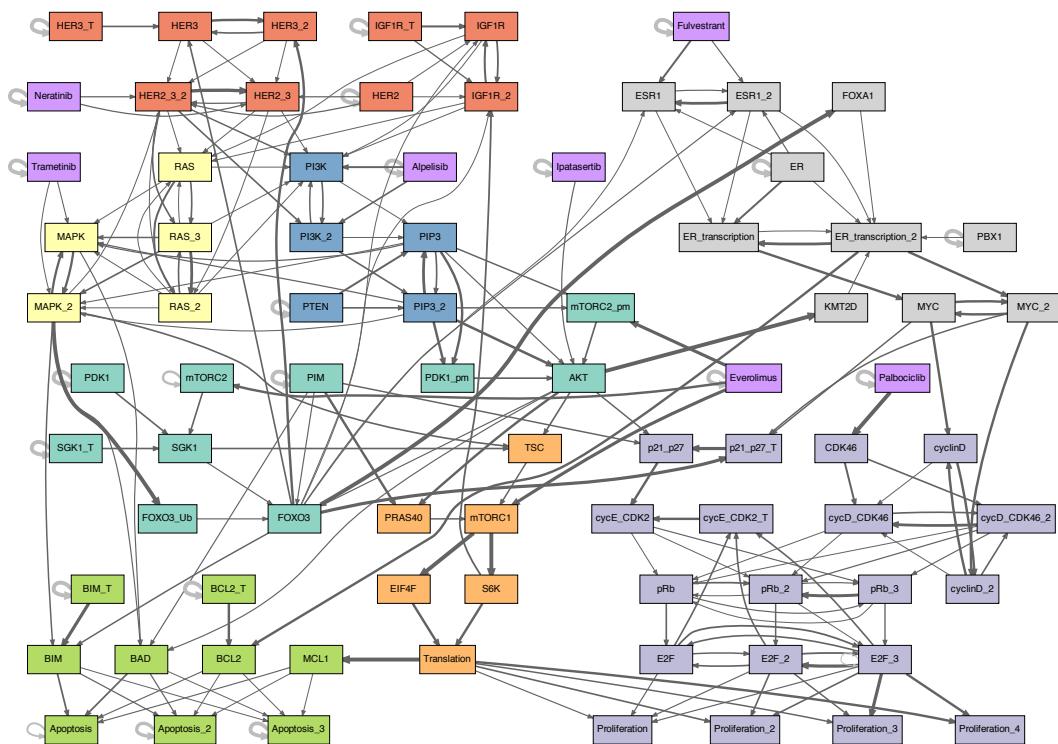
**Fig. S9. The effective graph (thresholded to 0.2 in the Leukemia BN model).** The effective graph for the BN model of leukemia, in which edge thickness denotes its effectiveness, thresholded to  $e_{ji} > 0.2$ .

**Table S6. Canalization measures for variables in the *Leukemia* model.**  $k$ ,  $k_r$ , and  $k_e$  denote in-degree, input redundancy and effective connectivity, respectively;  $k_r^*$  and  $k_e^*$  denote versions of  $k_r$ , and  $k_e$  normalized by  $k$ ;  $k_e^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively. Nodes with  $k = 1$  have no redundancy ( $k_r = 0, k_e = 1$ ), and input nodes have no incoming edges.

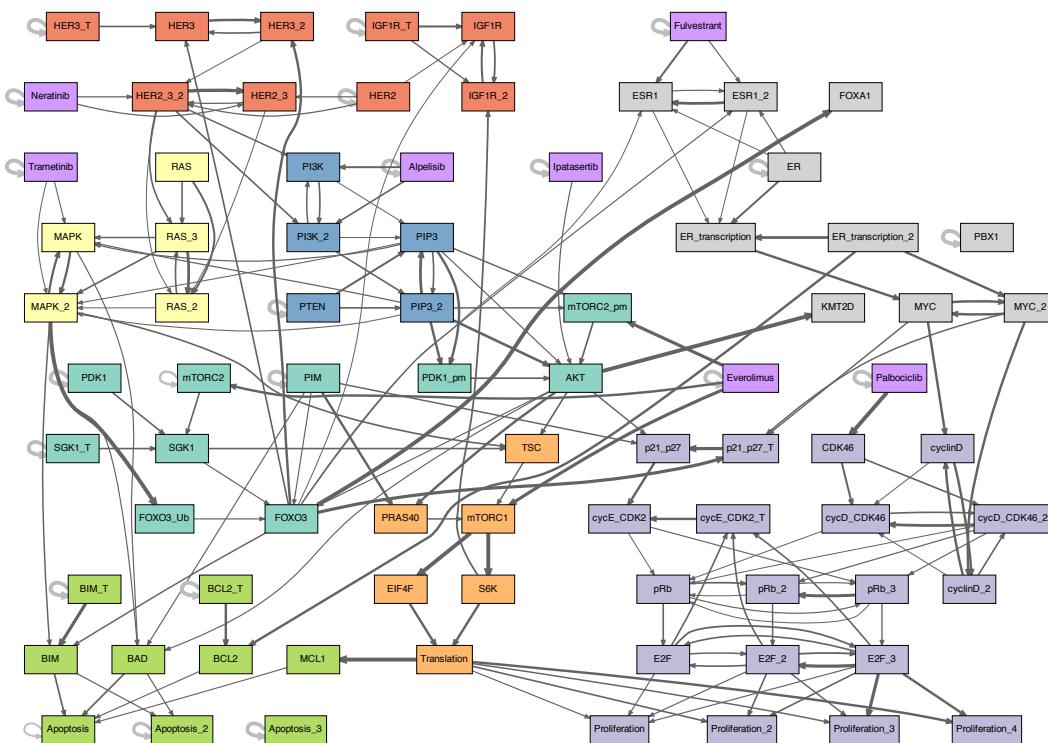
$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k_e^{out}$	$k_e^{out}$	$k_e^{out}/k_e^{out}$
IFNG	6	4.6	1.4	0.76	0.24	5	1.7	0.34
JAK	6	4.4	1.6	0.74	0.26	3	2	0.68
GAP	5	3.4	1.6	0.68	0.32	2	0.98	0.49
Caspase	5	3.2	1.8	0.65	0.35	2	0.86	0.43
DISC	5	3.2	1.8	0.64	0.36	3	1.9	0.64
NFKB	5	3.4	1.6	0.68	0.32	11	6.5	0.59
BclxL	5	3.6	1.4	0.72	0.28	1	0.52	0.52
BID	4	2.5	1.5	0.62	0.38	3	1.3	0.43
FLIP	4	2.3	1.7	0.58	0.42	2	0.29	0.15
LCK	4	2.3	1.7	0.58	0.42	1	0.62	0.62
MCL1	4	2.8	1.2	0.7	0.3	1	0.52	0.52
FasL	4	2.8	1.2	0.7	0.3	1	0.42	0.42
IL2	4	2.4	1.6	0.59	0.41	7	2.7	0.39
IFNGT	3	1.8	1.2	0.58	0.42	1	0.39	0.39
IL2RB	3	1.4	1.6	0.48	0.52	5	1.9	0.38
Fas	3	1.7	1.3	0.58	0.42	2	0.94	0.47
IL2RA	3	1.7	1.3	0.58	0.42	2	0.51	0.26
IL2RAT	3	1.4	1.6	0.48	0.52	1	0.42	0.42
TRADD	3	1.7	1.3	0.58	0.42	2	0.29	0.15
P2	3	1.4	1.6	0.48	0.52	2	0.76	0.38
TPL2	3	1.4	1.6	0.48	0.52	1	0.57	0.57
GZMB	3	1.4	1.6	0.48	0.52	3	0.74	0.25
RAS	3	1.4	1.6	0.48	0.52	3	2	0.66
SOCS	3	1.7	1.2	0.58	0.42	1	0.6	0.6
GRB2	2	0.75	1.2	0.38	0.62	2	1	0.5
FYN	2	0.75	1.2	0.38	0.62	2	1.6	0.81
PDGFR	2	0.75	1.2	0.38	0.62	4	2.4	0.6
Apoptosis	2	0.75	1.2	0.38	0.62	1	0.62	0.62
IAP	2	0.75	1.2	0.38	0.62	3	0.88	0.29
Ceramide	2	0.75	1.2	0.38	0.62	2	1.1	0.53
TBET	2	0.75	1.2	0.38	0.62	5	3.4	0.68
TCR	2	0.75	1.2	0.38	0.62	3	1.9	0.62
PLCG1	2	0.75	1.2	0.38	0.62	1	0.38	0.38
ZAP70	2	0.75	1.2	0.38	0.62	2	1.1	0.57
Proliferation	2	0.75	1.2	0.38	0.62	0	0	-
S1P	2	0.75	1.2	0.38	0.62	3	2.2	0.75
ERK	2	0.75	1.2	0.38	0.62	3	1.5	0.52
CREB	2	0.75	1.2	0.38	0.62	2	0.61	0.31
IL2RBT	2	0.75	1.2	0.38	0.62	1	0.81	0.81
PI3K	2	0.75	1.2	0.38	0.62	5	2.9	0.57
Stimuli	1	0	1	0	1	3	1.7	0.57
Stimuli2	1	0	1	0	1	2	1.8	0.91
GPCR	1	0	1	0	1	1	1	1
IL15	1	0	1	0	1	5	2.3	0.47
CD45	1	0	1	0	1	3	2.3	0.77
SMAD	1	0	1	0	1	1	0.39	0.39
SPHK1	1	0	1	0	1	1	0.62	0.62
PDGF	1	0	1	0	1	2	1.6	0.81
CTLA4	1	0	1	0	1	1	0.62	0.62
A20	1	0	1	0	1	1	0.42	0.42
sFas	1	0	1	0	1	1	0.42	0.42
FasT	1	0	1	0	1	3	2.2	0.73
TNF	1	0	1	0	1	2	0.79	0.4
P27	1	0	1	0	1	1	0.62	0.62
STAT3	1	0	1	0	1	8	3.4	0.43
RANTES	1	0	1	0	1	1	0.095	0.095
NFAT	1	0	1	0	1	3	0.95	0.32
MEK	1	0	1	0	1	1	0.62	0.62
Cytoskeleton_signaling	1	0	1	0	1	0	0	-
TAX	1	0	1	0	1	2	1.8	0.91



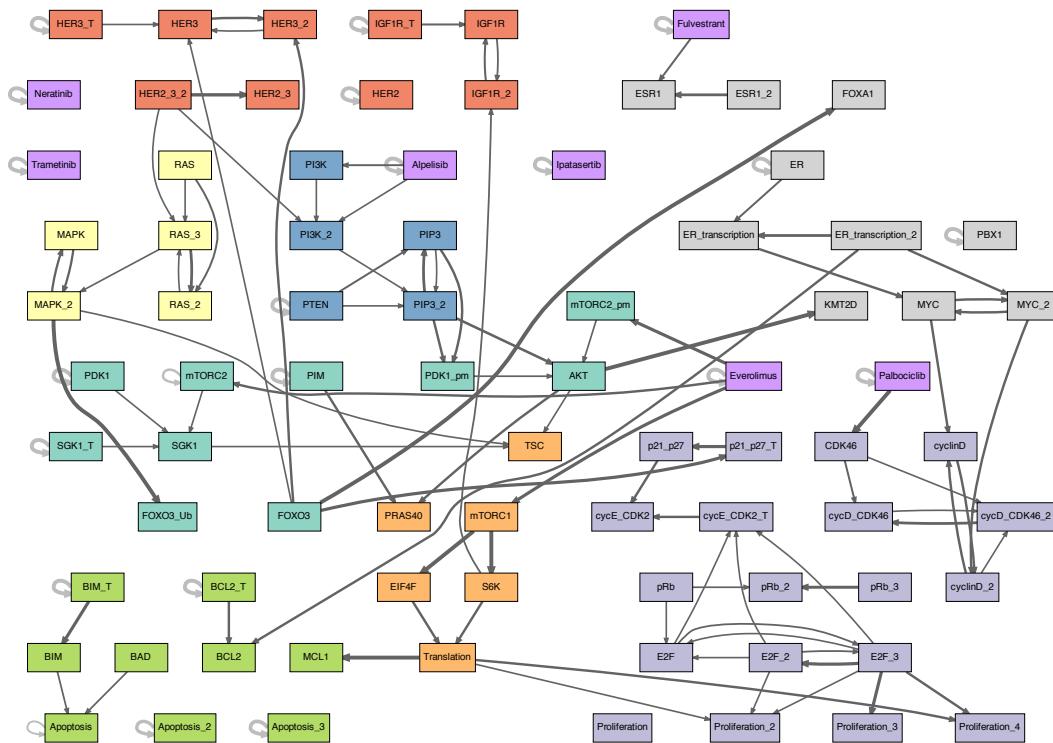
**Fig. S10. The interaction graph in the ER+ breast cancer BN model.** The goal of this model is to find interventions—especially single-node interventions—that synergize with the PI3K inhibitor *Alpelisib*, with particular interest on six other drugs used in cancer treatment: *Fulvestrant*, *Palbociclib*, *Everolimus*, *Neratinib*, *Trametinib*, and *Ipatasertib* (all drug nodes in Purple). Specifically, the goal is to study how well these drugs control cancer cells to apoptosis or proliferation, which in this model are specific variables (10). This is done by running Monte-Carlo simulations of the BN while setting *Alpelisib* to ON, in addition to setting baseline nodes to the cancerous state (Figs. S15–6), followed by setting the interventions to be tested to the appropriate state, for example, another Drug set to ON. See details in (10).



**Fig. S11. The effective graph in the ER+ breast cancer BN model.** Edge thickness denotes its effectiveness,  $e_{ji}$ .



**Fig. S12. The effective graph (thresholded to 0.2) in the ER+ breast cancer BN model.** The effective graph for the BN model of ER+ breast cancer, in which edge thickness denotes its effectiveness, thresholded to  $e_{ji} > 0.2$ .



**Fig. S13. The effective graph (thresholded to 0.4) in the ER+ breast cancer BN model.** The effective graph for the BN model of ER+ breast cancer, in which edge thickness denotes its effectiveness, thresholded to  $e_{ji} > 0.4$ .

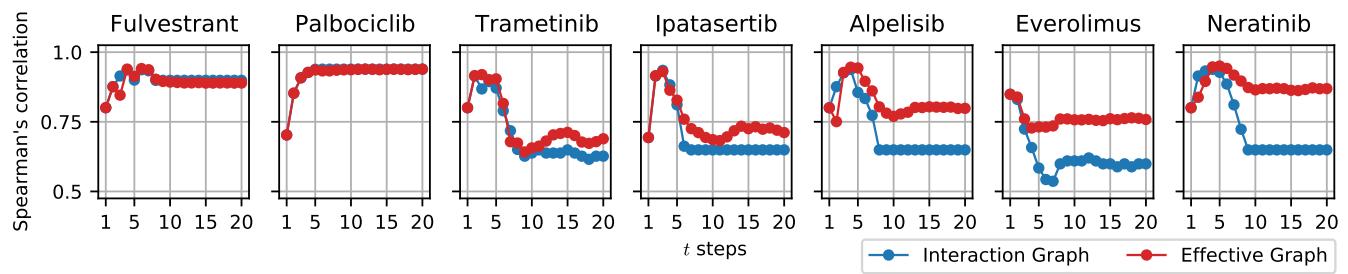
**Table S7. Canalization measures for variables in the *Breast Cancer* model.**  $k$ ,  $k_r$ , and  $k_e$  denote in-degree, input redundancy and effective connectivity, respectively;  $k_r^*$  and  $k_e^*$  denote versions of  $k_r$ , and  $k_e$  normalized by  $k$ ;  $k_e^{out}$  and  $k_e^{out}$  denote out-degree and effective out-degree, respectively. Nodes with  $k = 1$  have no redundancy ( $k_r = 0$ ,  $k_e = 1$ ), and input nodes have no incoming edges.

$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k_e^{out}$	$k_e^{out}$	$k_e^{out}/k_e^{out}$
PI3K	9	7.4	1.6	0.82	0.18	2	0.66	0.33
MAPK	7	4.9	2.1	0.7	0.3	2	0.88	0.44
MAPK_2	6	4	2	0.67	0.33	6	2.8	0.47
ER_transcription_2	6	4.9	1.1	0.82	0.18	3	2	0.65
RAS	6	4.9	1.1	0.82	0.18	4	1.1	0.28
HER2_3_2	6	4.5	1.5	0.76	0.24	6	2.5	0.41
pRb_2	5	3.2	1.8	0.64	0.36	4	1.1	0.27
AKT	5	2.9	2.1	0.57	0.43	6	3	0.5
HER2_3	5	3.2	1.8	0.65	0.35	4	0.81	0.2
pRb	5	3.9	1.1	0.78	0.22	3	1.2	0.38
Apoptosis_3	5	3.4	1.6	0.68	0.32	1	0.95	0.95
IGF1R_2	5	3.4	1.6	0.68	0.32	3	0.77	0.26
Apoptosis	5	3	2	0.61	0.39	1	0.49	0.49
E2F_3	5	3.4	1.6	0.68	0.32	8	4	0.5
Apoptosis_2	5	3.2	1.8	0.65	0.35	1	0.85	0.85
ER_transcription	4	2.3	1.7	0.58	0.42	2	0.8	0.4
PIP3	4	2.3	1.7	0.58	0.42	6	2.3	0.38
IGF1R	4	2.5	1.5	0.62	0.38	3	0.71	0.24
FOXO3	4	2.8	1.2	0.7	0.3	9	4.2	0.46
RAS_2	4	2.3	1.7	0.58	0.42	5	1	0.2
ESR1	4	2.3	1.7	0.58	0.42	2	0.54	0.27
Proliferation	4	2.8	1.2	0.7	0.3	0	0	-
cycD_CDK46	4	2.3	1.7	0.58	0.42	3	0.79	0.26
BAD	4	2.8	1.2	0.7	0.3	3	0.92	0.31
ESR1_2	4	2.8	1.2	0.7	0.3	3	1.1	0.38
pRb_3	4	2.8	1.2	0.7	0.3	3	1.3	0.44
BIM	3	1.4	1.6	0.48	0.52	3	0.92	0.31
mTORC1	3	1.4	1.6	0.48	0.52	2	2	1
TSC	3	1.8	1.2	0.58	0.42	1	0.38	0.38
cycE_CDK2_T	3	1.8	1.2	0.58	0.42	1	0.62	0.62
p21_p27	3	1.4	1.6	0.48	0.52	1	0.62	0.62
p21_p27_T	3	1.4	1.6	0.48	0.52	1	0.81	0.81
E2F	3	1.8	1.2	0.58	0.42	4	1.5	0.39
mTORC2_pm	3	1.4	1.6	0.48	0.52	1	0.43	0.43
RAS_3	3	1.7	1.3	0.58	0.42	5	1.8	0.35
E2F_2	3	1.4	1.6	0.48	0.52	6	2.4	0.4
cycD_CDK46_2	3	1.7	1.3	0.58	0.42	4	1.5	0.39
SGK1	3	1.7	1.3	0.58	0.42	2	0.71	0.36
Proliferation_2	3	1.8	1.2	0.58	0.42	0	0	-
Proliferation_3	3	1.4	1.6	0.48	0.52	0	0	-
PIP3_2	3	1.7	1.3	0.58	0.42	6	3	0.49
PI3K_2	3	1.7	1.3	0.58	0.42	3	1	0.34
HER3	3	1.8	1.2	0.58	0.42	3	0.89	0.3
cyclinD_2	2	0.75	1.2	0.38	0.62	3	1.3	0.43
Translation	2	0.75	1.2	0.38	0.62	5	2.7	0.54
MYC	2	0.75	1.2	0.38	0.62	3	1.6	0.54
MYC_2	2	0.75	1.2	0.38	0.62	3	1.6	0.54
cyclinD	2	0.75	1.2	0.38	0.62	2	0.86	0.43
BCL2	2	0.75	1.2	0.38	0.62	3	0.6	0.2
Proliferation_4	2	0.75	1.2	0.38	0.62	0	0	-
PRAS40	2	0.75	1.2	0.38	0.62	1	0.38	0.38
mTORC2	2	0.75	1.2	0.38	0.62	2	1	0.52
cycE_CDK2	2	0.75	1.2	0.38	0.62	3	0.67	0.22
PDK1_pm	2	0.75	1.2	0.38	0.62	1	0.43	0.43
HER3_2	2	0.75	1.2	0.38	0.62	3	0.81	0.27

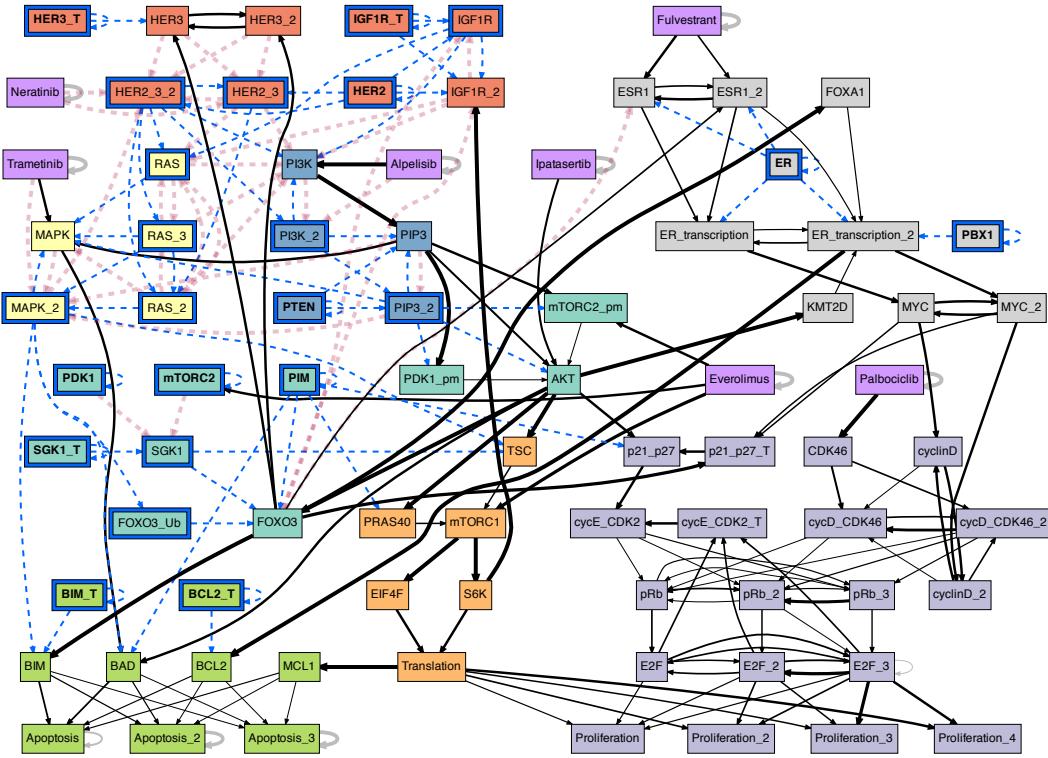
*Continues on the next page*

Table S7 - *Continued from previous page*

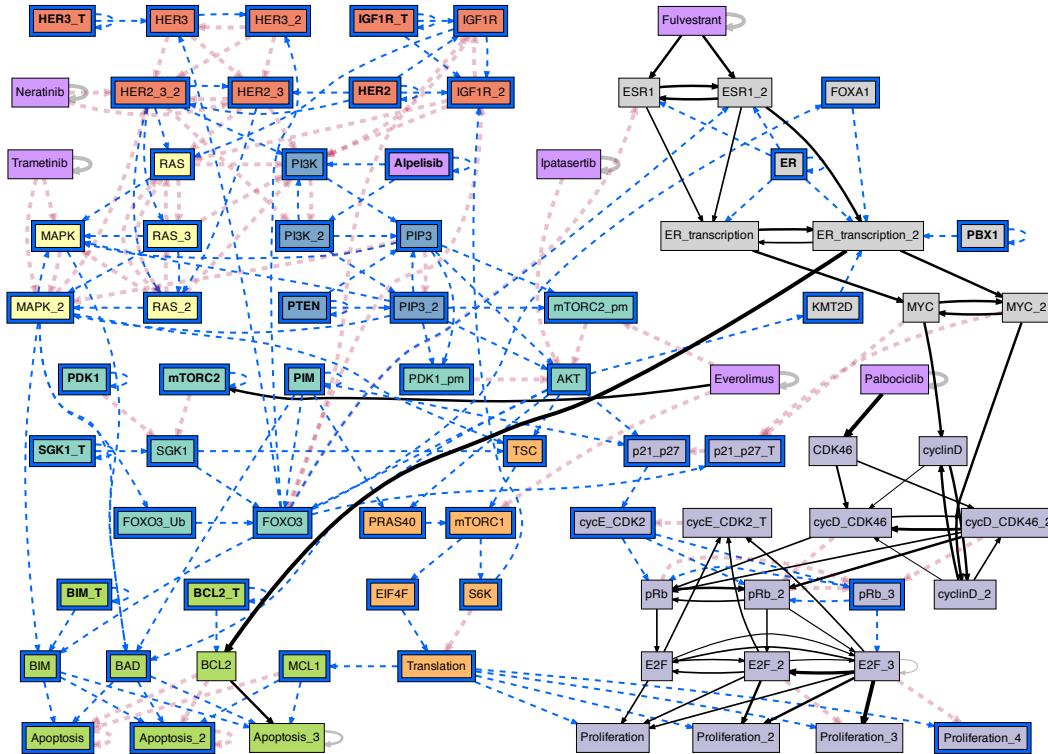
$x_i$	$k$	$k_r$	$k_e$	$k_r^*$	$k_e^*$	$k^{out}$	$k_e^{out}$	$k_e^{out}/k^{out}$
S6K	1	0	1	0	1	2	1.1	0.54
Fulvestrant	1	0	1	0	1	3	1.8	0.6
Alpelisib	1	0	1	0	1	3	1.9	0.64
Everolimus	1	0	1	0	1	4	3.2	0.81
Trametinib	1	0	1	0	1	3	1.5	0.49
Ipatasertib	1	0	1	0	1	2	1.3	0.65
Palbociclib	1	0	1	0	1	2	2	1
Neratinib	1	0	1	0	1	3	1.6	0.54
HER2	1	0	1	0	1	5	2	0.41
HER3_T	1	0	1	0	1	2	1.4	0.71
PDK1	1	0	1	0	1	2	1.4	0.71
PIM	1	0	1	0	1	5	2.6	0.52
SGK1_T	1	0	1	0	1	2	1.4	0.71
ER	1	0	1	0	1	5	2.2	0.45
CDK46	1	0	1	0	1	2	0.93	0.46
PTEN	1	0	1	0	1	3	1.9	0.64
KMT2D	1	0	1	0	1	1	0.18	0.18
FOXO3_Ub	1	0	1	0	1	1	0.3	0.3
BIM_T	1	0	1	0	1	2	1.8	0.91
BCL2_T	1	0	1	0	1	2	1.6	0.81
PBX1	1	0	1	0	1	2	1.2	0.59
FOXA1	1	0	1	0	1	1	0.18	0.18
MCL1	1	0	1	0	1	3	0.6	0.2
EIF4F	1	0	1	0	1	1	0.62	0.62
IGF1R_T	1	0	1	0	1	3	1.9	0.62



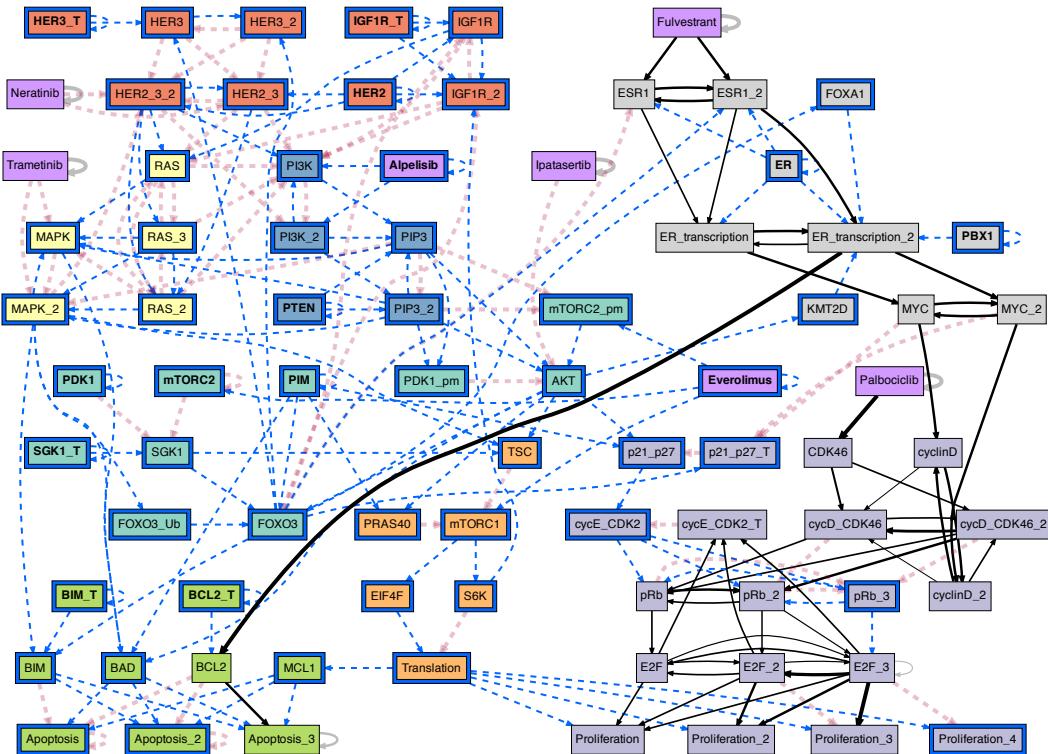
**Fig. S14. The spread of perturbations in the ER+ breast cancer BN model.** For each of the drug variable nodes, the predictive power of the path-length approximation using the interaction graph (blue), and the effective graph (red); measured by the Spearman's rank correlation (vertical axis) to the total impact of respective variable, after 20 steps (horizontal axis).



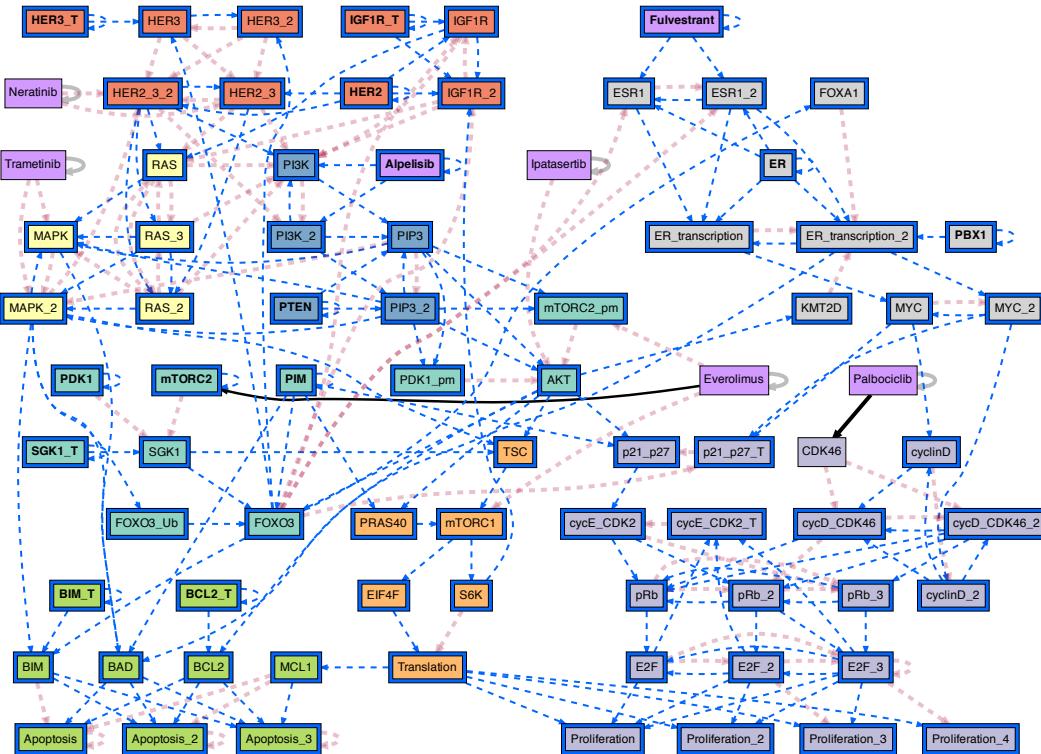
**Fig. S15. The conditional effective graph of the ER+ breast cancer BN model, conditioned on the ER+/Her2- cancer cell state baseline.** The baseline is defined by  $K = \{ER = ON, HER2 = OFF, HER3\_T = OFF, IGF1R\_T = ON, PBX1 = ON, PTEN = OFF, SGK1\_T = OFF, PIM1 = OFF, PDK1\_T = OFF, mTORC2 = OFF, BIM\_T = OFF, and BCL2\_T = OFF\}$ . Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that the  $ER + /Her2-$  cancer cell state baseline alone resolves a substantial portion of the possible dynamics in comparison to the non-conditioned effective graph (Fig. S10). Strikingly, *Neratinib* becomes redundant under this baseline initial condition, along with much of the *HER* pathway. Thus, *Neratinib* has no effect on this model under  $ER + /Her2-$  cancer cell state. Indeed, *Neratinib* is one of the drugs that were shown not to synergize with *Alpelisib* in (10); the others are *Ipatasertib* and *Trametinib*. We can see that these three drugs only contribute to the same pathways that *Alpelisib* already acts on and become redundant when *Alpelisib* is present (see Fig. S16, as well as Figs. S19, S20, S22).



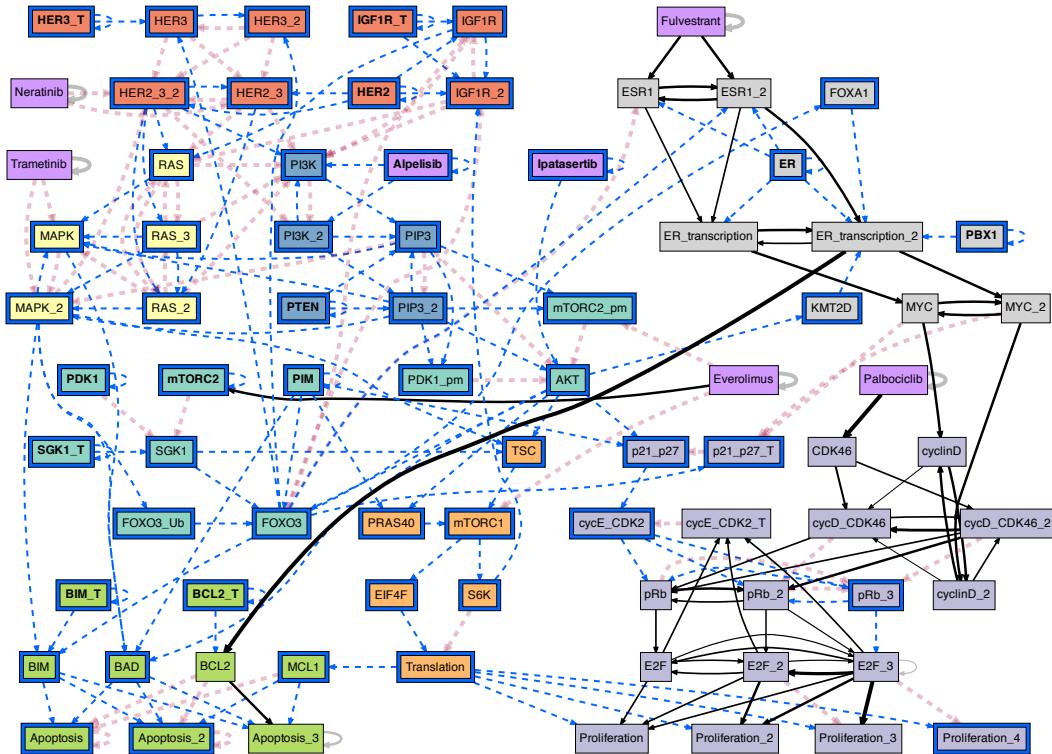
**Fig. S16. The conditional effective graph of the ER+ breast cancer BN model, conditioned on the ER+/Her2- cancer cell state baseline + Alpelisib=ON.**  $K$  is comprised of baseline nodes (see S15 caption) + {Alpelisib = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that Alpelisib with the  $ER + /Her2-$  cancer cell state baseline resolves a majority of the possible dynamics in comparison to the non-conditioned effective graph (Fig. S10); unresolved dynamics is circumscribed almost entirely to  $ER$  and proliferation pathways. Interestingly, Ipatasertib, Neratinib, and Trametinib become redundant under this initial condition: they have no effect on model dynamics under Alpelisib+  $ER + /Her2-$  cancer cell state. The effective graph, however, reveals that the dynamics of this network is very robust to perturbation and hard to control because its subsystems are effectively decoupled. In particular, the baseline + Alpelisib=ON condition reveals that canalization works by preventing propagation of signals and cross-regulation. Indeed, most of the (non-drug) variables that have an impact on cancer apoptosis or proliferation under this condition (see Table 3 in (10)) have short paths to those target variables (at most 3 edges) in the effective graph. Exceptions in Table 3 in (10) are only a few nodes involved in the estrogen receptor ( $ER$ ) transcription and signaling pathway such as the  $MYC$  oncogene and  $KMT2D$  epigenetic transcription activator. Interestingly, the conditional effective graph reveals that  $KMT2D$  becomes fixed under this condition, so any impact on proliferation can only occur by perturbing it out of its fixed state.



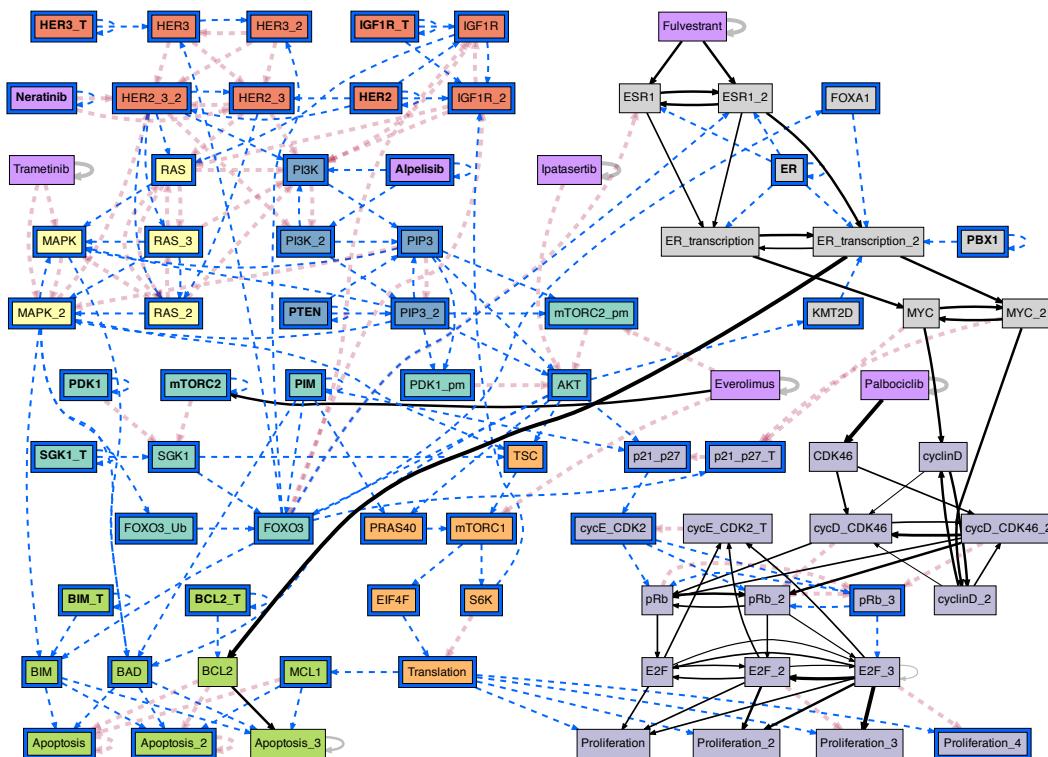
**Fig. S17. The conditional effective graph of the ER+ breast cancer BN model, conditioned on the *ER+/Her2-* cancer cell state baseline + *Alpelisib=Everolimus=ON*.**  $K$  is comprised of baseline nodes (see S15 caption) + {*Alpelisib* = ON, *Everolimus* = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Everolimus* hardly resolves any additional dynamics to what *Alpelisib* with the *ER+ / Her2-* cancer cell state baseline already do (Fig. S16); only a few connections to *AKT* pathway.



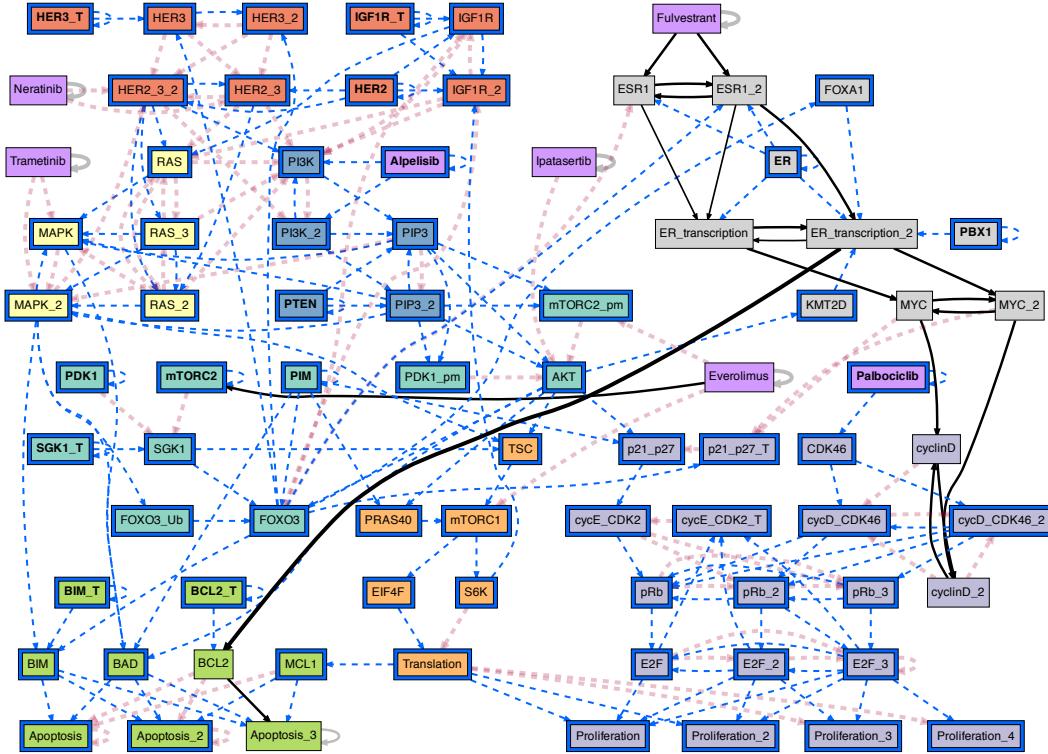
**Fig. S18. The conditional effective graph of the ER+ breast cancer BN model, conditioned on the ER+/Her2- cancer cell state baseline + Alpelisib=Fulvestrant=ON.**  $K$  is comprised of baseline nodes (see S15 caption) + { $Alpelisib = \text{ON}$ ,  $Fulvestrant = \text{ON}$ }. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Fulvestrant* resolves almost the remaining dynamics that was not yet resolved by *Alpelisib* with the *ER + /Her2-* cancer cell state baseline (Fig. S16). In particular, every apoptosis and proliferation node gets resolved. This combination strategy is the most powerful, rendering all other drugs redundant, except for the influence *Everolimus* has in *AKT* pathway, which does not propagate anyway.



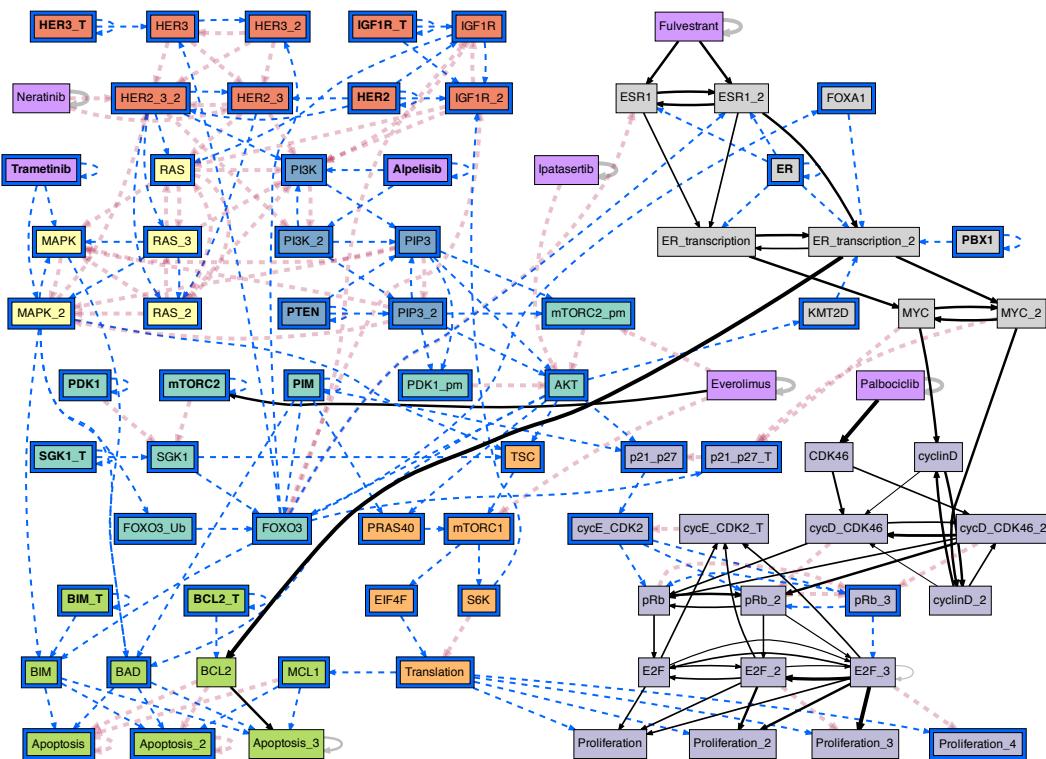
**Fig. S19.** The conditional effective graph of the ER+ breast cancer BN model, conditioned on the *ER+/Her2-* cancer cell state baseline + *Alpelisib=Ipatasertib=ON*.  $K$  is comprised of baseline nodes (see S15 caption) + {*Alpelisib* = ON, *Ipatasertib* = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Ipatasertib* does resolve any additional dynamics to what *Alpelisib* with the *ER+ / Her2-* cancer cell state baseline already do (Fig. S16).



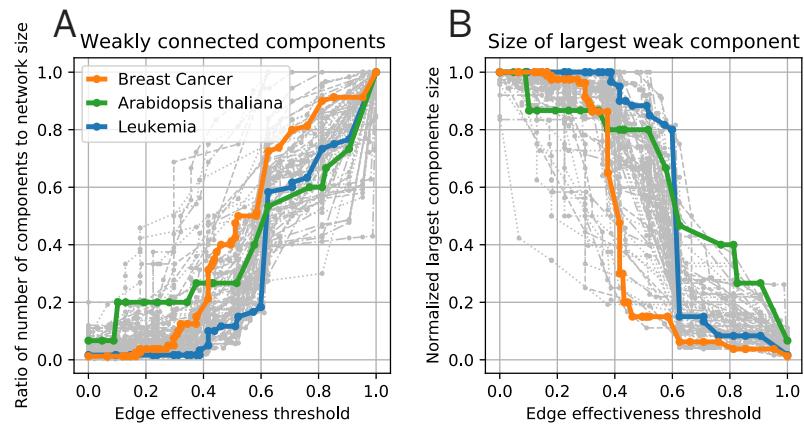
**Fig. S20.** The conditional effective graph of the ER+ breast cancer BN model, conditioned on the *ER+/Her2-* cancer cell state baseline + *Alpelisib=Neratinib=ON*.  $K$  is comprised of baseline nodes (see S15 caption) + {*Alpelisib* = ON, *Neratinib* = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Neratinib* does resolve any additional dynamics to what *Alpelisib* with the *ER+ / Her2-* cancer cell state baseline already do (Fig. S16).



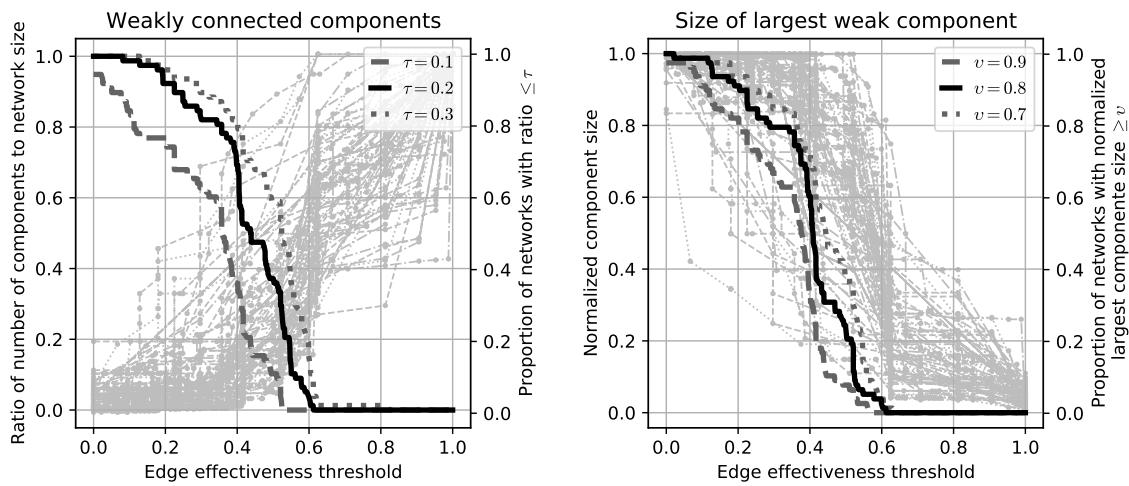
**Fig. S21. The conditional effective graph of the ER+ breast cancer BN model, conditioned on the *ER+/Her2-* cancer cell state baseline + *Alpelisib=Palbociclib=ON*.**  $K$  is comprised of baseline nodes (see S15 caption) + {*Alpelisib* = ON, *Palbociclib* = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Palbociclib* resolves the proliferation pathway in addition to the dynamics that *Alpelisib* with the *ER+ / Her2-* cancer cell state baseline resolved on their own (Fig. S16). In particular, every proliferation node gets resolved, which means that *Palbociclib* displays useful synergy with *Alpelisib*. However, unlike the case of *Fulvestrant*, *Palbociclib* has no effect on apoptosis in this model for the *ER+ / Her2-* cancer cell state baseline.



**Fig. S22.** The conditional effective graph of the ER+ breast cancer BN model, conditioned on the *ER+ / Her2-* cancer cell state baseline + *Alpelisib = Trametinib = ON*.  $K$  is comprised of baseline nodes (see S15 caption) + {*Alpelisib* = ON, *Trametinib* = ON}. Variables in  $K$  (those initially pinned) are shown with a blue border and bold text; variables whose state becomes fixed (become constants), are shown with a blue border only. Edges that transmit a constant input state are denoted by a dashed blue color, while unresolved edges are denoted by black color with thickness proportional to their effectiveness,  $e_{ji}$ , with the fully redundant edges shown in dashed red. We can see that *Trametinib* does resolve any additional dynamics to what *Alpelisib* with the *ER+ / Her2-* cancer cell state baseline already do (Fig. S16).



**Fig. S23. Weakly connected components of threshold effective graphs reveal dynamical modules.** **A** Ratio of the number of weakly connected components to network size. **B** Size of the largest weakly connected component, normalized by network size. In both graphs, the *ER+* breast cancer (orange), leukemia (blue), and *Arabidopsis thaliana* (blue) networks are highlighted.



**Fig. S24. Analysis of weakly connected components per edge effectiveness threshold for all Cell Collective models.** Gray thin lines denote each network in dataset, thick lines denote overall statistics per legend. **Left.** Ratio of number of weakly connected components to network size (left vertical axis) for a given edge effectiveness threshold (horizontal axis); ratio is 1 when every node in network is a separate component, and very small when there is a single weakly connected component. Also shown are three statistics for the proportion of networks with ratio  $\leq \tau$  (right vertical axis); e.g. for  $\tau = 0.2$ , black thick line denotes the proportion of networks whose ratio of number of weakly connected components to network size is smaller than 0.2 at a given edge effectiveness threshold. We can see, for instance, that for an edge effectiveness of 0.2, more than 90% of the networks have a small number of weakly connected components, specifically, less than  $\tau = 20\%$  of the network size; for edge effectiveness larger than 0.4, on the other hand, most networks quickly break into many components, and for edge effectiveness larger than 0.6, no networks have a ratio of number of components to network size smaller than 20% (or even 30%, per dotted thick line). **Right.** Size of largest weakly connected component relative to network size (left vertical axis) for a given edge effectiveness threshold (horizontal axis); ratio is 1 when there is a single weakly connected component, and very small when every node is its separate component. Also shown are three statistics for the proportion of networks with largest normalized component size  $\geq v$  (right vertical axis); e.g. for  $v = 0.8$ , black thick line denotes the proportion of networks whose largest normalized component size is larger than 0.8 at a given edge effectiveness threshold. We can see, for instance, that for an edge effectiveness of 0.2, about 90% of the networks have a largest weakly connected component comprised of at least  $v = 80\%$  of the network; for edge effectiveness larger than 0.4, on the other hand, most networks quickly lose a substantial largest weakly connected component, and for edge effectiveness larger than 0.6, no networks have a largest weakly connected component comprised of at least 80% of the network (or even 70%, per dotted thick line)