# Assessing the effect of unknown widespread perturbations in complex systems using the $\nu$-gap

Alberto Carignano*, Jin Junyang, Alex Webb, Jorge Gonçalves

*Abstract*— Pinpointing the exact locations of perturbations can help to detect and correct faults in large-scale systems, such as power grids (lost of transmission lines), internet (loss of servers) and biological systems (diseases). This paper outlines a mathematical framework to study the effect of perturbations on large-scale systems, with particular emphasis to biological applications. In particular, it focuses on wide-spread perturbations that target unknown components of the system. These problems are usually studied with genome-wide assays, which are becoming increasingly popular and accessible. However, analysis of the data sets produced by this technology remains challenging: genome-wide experiments are often inherently noisy, with a small number of measurements, and a low sampling rate. The paper first develops a simple yet powerful network inference tool based on LTIs. We compare this tool with the current state of the art. Then, as its major contribution, the paper develops a method for network differentiation, where it detects the effects of perturbations in large scale-systems. The method is based on the $\nu$-gap, a control engineering tool that measures the distance between linear models. A major difference between this work and others, is that we look at changes in dynamics in links, as opposed to the standard differential expression analysis that focuses on changes in node concentrations. Through an illustrative model, the paper shows how perturbations impact certain links in the network, which can then be captured by differences between LTIs with the $\nu$-gap.

## I. INTRODUCTION

One difficult, yet fundamental problem in biology is to assess the effect of unknown widespread perturbations in complex systems. For instance, a certain drug treatment might result in a particular measurable phenotype, but the exact dynamics involved are often unknown. Knowing precisely the pathways affected by the treatment would be a big step toward improving the effect of the drug, or eliminating potential side effects. However, this is not a trivial problem: a drug can have multiple targets and its mechanism of action is an inherently noisy process (the distribution of the drug, for instance, is probably not uniform across the organism). This problem becomes even more complicated if the structure of the network is unknown. Cellular components are connected through very complex regulatory mechanisms, often including mutual regulation.

A. Carignano is with Faculty of Electrical Engineering, University of Washington, Seattle, WA 98915, USA ac86@uw.edu

J. Junyang is with Faculty of Plant Sciences, University of Cambridge, CB2 3EA, UK jj415@cam.ac.uk

Prof. A.A.R. Webb is with Faculty of Plant Sciences, University of Cambridge, CB2 3EA, UK aarw2@cam.ac.uk

Prof. J. Gonçalves is with the Faculté des Sciences, de la Technologie et de la Communication, Université du Luxembourg, Belvaux, L-4366, Luxembourg jorge.goncalves@uni.lu

The topic of network inference has been central for computational system biology. Since the first application of Bayesian networks on gene expression data ([8]), several network inference methods have been suggested, along with procedures to assess their accuracy ([16], [5], [23]). A recent work has implemented a comparison of the most commonly-used network inference methods using both simulated and real data from the Arabidopsis circadian clock ([1]). Methods were assessed according to a scoring system (AUROC) that evaluates the number of false positives and correct answers as a smooth function of the acceptance threshold. The methods considered in the paper all infer causal networks, and include most of the well-known methodologies, such as Graphical Gaussian Models (GGM, [21], [7]), Sparse Regression (Lasso [22] and Elastic Net [24]), Time-varying Sparse Regression (Tesla [2]), Hierarchical Bayesian Regression models (HBR, [1]), Non-Homogeneous Hierarchical Bayesian models ([11]), Automatic Relevance Determination in the context of Sparse Bayesian Regression (ARD-SBR, [20]), Bayesian Spline Autoregression (BSA, [17]), State Space Models (SSM, [3]), Gaussian processes (GP, [4]), Mutual information methods (ARACNE, [16]), Mixture Bayesian network models (MBN, [14]), and Gaussian Bayesian networks (BGe, [9]).

Data were simulated using the computational model proposed by the Millar's group (Millar2010, [19]), in both deterministic and stochastic settings. AUROC scores were measured for each method and for each data set, and an ANOVA index measured the general performance. According to the study, MBN and ARACNE show a significantly low performance, while all the other methods perform in a similar AUROC score range. The method that performs best is HBR, followed closely by BGe.

Here we propose a new methodology for network inference that proves to be especially useful for poorly-sampled data corrupted with noise (as those produced by a microarray experiment). The method uses Linear Time Invariant (LTI) models to represent biological regulation: gene A regulates gene B if a linear model could be estimated and validated using the two genes as an input-output pair. The use of LTIs to explain biological regulations has been shown previously ([6], [12]), but the application to network inference for genome-wide data sets is novel. LTIs models have several advantages over nonlinear models. First of all, they have been extensively studied in the field of control theory and signal processing. In particular, LTIs have a frequency description with an easy visual interpretation, which can be used to infer their stability and performance. Moreover, a linear

description incorporates in $x(t)$ hidden variables, biological steps that are not part of the input/output pair but are necessary to describe the dynamics. The number of hidden variables is useful to formulate hypothesis on the nature of the biological process (for instance if the regulation is at the transcriptional or the translational level). Finally, LTIs unlike nonlinear models, have a fixed structure that prevents the user to bias the results towards specific interpretations. Under mild hypothesis on the input, there is only one optimal linear model that could explain the data, and there are qualitative and quantitative ways of assessing the goodness of the performance.

A perturbation of a biological regulation can either affect the way the regulation works, or inhibit it. Inhibition corresponds to a change in topology: the network does not preserve its structure. This should be fairly visible assuming the networks are correctly estimated. On the other hand, regulation changes are harder to capture as the network structure should not, in principle, change. We will show how LTIs are appropriate descriptions to capture the effect of perturbations. This can be achieved using a tool called $\nu$-gap ([10]). In particular, a model-example will be used to describe how our method performs in practice

## II. NETWORK INFERENCE USING LINEAR TIME INVARIANT MODELS

In this section we outline the LTI network inference method, and test it on a chosen simulated model. Focus will be on data with low sampling time and affected by noise, in order to reproduce common experimental conditions of biological system.

### A. Methodology

An LTI model has the form:

$$\frac{dx}{dt}(t) = Ax(t) + Bu(t) + Ke(t)$$
$$y(t) = Cx(t) + Du(t) + e(t) \qquad (1)$$

where $x(t)$ represents the internal dynamics, and $e(t)$ represents the inherent white noise of the system and measurements. This concept can be generalized to higher dimensions ($x \in \mathcal{R}^n$) by considering the parameters in 1 as matrices with appropriate dimensions. LTIs are often seen in the equivalent, frequency-domain version:

$$Y(s) = \mathbf{G}(\mathbf{s})U(s) \qquad (2)$$

where $\mathbf{G} = C(s\mathbb{I} - A)^{-1}B + D$ and $U(s)$ and $Y(s)$ are the Laplace transform of the original input/output pair. Estimating an LTI model consists of identifying the matrices A, B, C, D, and K, and the initial conditions $x(0)$ in 1. Several methods have been implemented that could identify the parameters of system 1, for instance the Matlab® algorithm 'pem'.

LTI models have been shown to be capable of describing biological systems ([6], [12]), and are less prone to over-fitting than most nonlinear methods. Moreover, their computational cost is known to be light. These characteristics make them a suitable choice to study large data sets. In particular, given a data set $S$ containing a set of measurements as time series, we are interested to infer the underlying system that have generated them. Our method is divided in two steps. In the first step, LTI models of different orders are estimated for all the possible input-output pairs from data set $D$: the most appropriate order is then selected, and each model is characterized by a measure of fitness between the simulated and the real data. In the second step, the noise in the system is estimated and a threshold on the fitness is set based on the noise level. All the models estimated in step 1 that are below the threshold are now discarded. The remaining models define the topology of the inferred network. Essentially, the first step estimates a fully-connected graph: every node is connected with all the others (no prior information is assumed). By applying a threshold, the second step selects the models that are most likely to correctly represent biological interactions, reducing the initial fully-connected graph to a more realistic one.

In the context of poorly-sampled noisy data set, we opted for the following specifications

- To identify the appropriate model order for the data, we used the small sample (second-order bias correction) version of the Akaike Information Criteria (AIC$_c$, [13]). Given a set of models for comparison, the one with the lowest AIC$_c$ values is the preferred choice;
- We use a performance index $f$ to compute the fitness between simulated and experimental data that is calculated according to the formula:

$$\text{f} = 100 \left( 1 - \frac{\sum\limits_{k=1}^{N} (y_k - \hat{y}_k)^2}{\sum\limits_{k=1}^{N} (y_k - \overline{y})^2} \right) \qquad (3)$$

where $\overline{y}$ is the average value of the experimental data, $y_k$ is the $k$-th data point, and $\hat{y}_k$ is the $k$-th simulated data point (in order to avoid divisions by zero, a different formula has to be used to estimate the fitness of a constant output);

- To assess which threshold is most appropriate for our particular biological system, we apply the method to a known network and we associate to each estimated LTI, the corresponding fitness value. We then consider increasing values of thresholds, and discard models with fitness below the threshold. The optimal threshold corresponds to the network with the highest number of correct answers.

### B. Comparison with the state of the art on a plant circadian clock model

To test our methodology we chose a known mathematical model of a biological system: the model proposed by the Millar's group in 2006 ([15]) of the circadian clock of the plant model organism A. *thaliana*. The model has light as its only input and five main nodes, and it takes into account three stages for each gene (mRNA level, protein

| Noise standard dev. | Millar 2006 | | Millar 2010 | | Millar 2012 | |
|---|---|---|---|---|---|---|
| | Correct '1's | Correct '0's | Correct '1's | Correct '0's | Correct '1's | Correct '0's |
| No noise | $-/75$ | $-/75$ | $-/44$ | $-/85$ | $-/33$ | $-/84$ |
| 0.01 noise | $8 \pm 0/45 \pm 7$ | $100 \pm 0/70 \pm 5$ | $50 \pm 4/38 \pm 0$ | $59 \pm 4/86 \pm 3$ | $18 \pm 1/36 \pm 2$ | $85 \pm 0/74 \pm 2$ |
| 0.03 noise | $7 \pm 4/60 \pm 6$ | $95 \pm 7/62 \pm 5$ | $49 \pm 4/38 \pm 0$ | $62 \pm 10/85 \pm 0$ | $17 \pm 3/37 \pm 0$ | $86 \pm 3/73 \pm 1$ |
| 0.05 noise | $5 \pm 5/63 \pm 0$ | $97 \pm 6/55 \pm 5$ | $56 \pm 11/58 \pm 11$ | $56 \pm 12/72 \pm 5$ | $17 \pm 7/37 \pm 0$ | $85 \pm 5/66 \pm 4$ |
| 0.07 noise | $3 \pm 5/73 \pm 6$ | $95 \pm 7/43 \pm 11$ | $53 \pm 11/63 \pm 0$ | $61 \pm 16/62 \pm 0$ | $18 \pm 3/37 \pm 0$ | $88 \pm 3/60 \pm 1$ |
| 0.1 noise | $5 \pm 5/80 \pm 7$ | $97 \pm 6/20 \pm 7$ | $50 \pm 9/63 \pm 0$ | $65 \pm 9/62 \pm 0$ | $10 \pm 7/40 \pm 2$ | $88 \pm 3/60 \pm 0$ |

level in the cytosol, and protein level in the nucleus). In addition, it has several nonlinear components (Hill's functions), which makes it the perfect test-case for our method. We add process noise to each differential equation and run a stochastic simulation of the model using the standard Euler-Maruyama method. We considered the following signal-to-noise standard deviation: $[0, 0.01, 0.03, 0.05, 0.07, 0.1]$. To obtain data for network inference, we simulated the model for 600 hours in 24 hour light/dark cycles to remove all possible transients. Then we changed the light conditions to constant light for another 96 hours. To reproduce standard whole-genome experiments (like, for instance, microarrays), characterized by low-sampled measurements, we took the simulated mRNA expression levels in the last 48 hours of constant light, and consider only the values with a 4-hour sampling frequency (12 data points for each of the five time series).

We then applied our method and estimated the optimal threshold for the deterministic case (standard deviation $= 0$). Models were estimated according to 1, imposing $K = 0$ as no process noise was present. Optimality is defined in terms of number of connections correctly identified (percentage of interactions plus percentage of lack of interactions correctly estimated). We then identify networks for the stochastic cases, applying the optimal thresholds computed for the deterministic case. Since process noise was present in these cases, we let the parameter $K$ free in the optimization. The results are shown in Table I. To assess the quality of our method, we carried out a comparison with the current state-of-the-art in network inference. We presented in the introduction the work of ([1]): their analysis highlights that hierarchical Bayesian regression models (HBR) seem to perform best compared to most of the known methods. Hence we chose HBR for benchmarking.

We implemented the algorithm in Matlab® following the outline given in the paper, and we set the hyperparameters of the gamma distribution of the 'signal-to-noise' ratio parameter $\delta_g$, to $\alpha = 2$ and $\beta = 0.2$, according to the specification suggested by the authors. The hyperparameter for the noise variance prior $\nu$ was set to $0.005$, as described in the paper. We compared the performances of the two methods using the simulated data of three known models of the circadian clock with increasing level of complexity: the previously

introduced Millar 2006 model, and the Millar 2010 ([19]), and Millar 2012 ([18]) models. Each model was simulated using Matlab® in both a deterministic and a stochastic manner: noise was added in each differential equation, with signal-to-noise standard deviation varying in the set $[0, 0.01, 0.03, 0.05, 0.07, 0.1]$. For each value of the variance, we simulated the models five times in order to prevent bias over a particular realization. Hence, we ended up with $3 + 3 \times 5 \times 5 = 78$ different data sets for our comparison. For a fair comparison, we first identified the optimal threshold for the HBR method using the datasets with 0.01 signal-to-noise. Hence, we obtained three different thresholds that allow HBR to adapt to the different structures of the networks (as we did for our method). Since HBR working assumption is that there has to be noise in the system, the deterministic data was not used. Following the protocol in the paper ([1]), we run HBR using 20000 iterations and then averaged the second half of the output sampled graphs (the first half is ignored as 'burn-in phase'): the averaged graph is the inferred network. The inferred network output is a matrix where entry $(i, j)$ contains the probability that gene i regulates gene j. The threshold is the minimum accepted probability for regulation. We performed a comparison between the two methods on the simplest model first, the Millar 2006 mode. We run HBR 10 times using the same settings, and compared the resulting graphs. The variance between each entry of the 10 averaged graphs was very small (below $0.05$): we concluded that the algorithm converged. Finally, we let the threshold vary between 0 and 1 with increments of $0.01$, compared the correspondent estimated graphs with the real one and found the optimal threshold according to the same metric used for our method. Networks were then derived for the stochastic cases using HBR and the derived optimal threshold (results shown in Table I). By adding the number of correct answers in the Millar 2006 column in Table I ('0's+'1's), one could observe that our method outperforms HBR for all but the 0.1 noise case. This is not surprising as the noisier the system, the less likely the optimal thresholds will hold.

To test whether these results were independent from the data set we used for training, we applied both methods to the remaining 52 data sets. The results are summarized in Table I, and strongly suggest that our method is better in determining the real network for this specific low-sampled and noisy data.

In particular, we noticed that while HBR optimizes towards a fully-disconnected networks for both Millar 2006 and Millar 2012, our method is better at balancing the number of correct '0's and '1's in all cases. Hence, the corresponding networks will be more likely to provide useful insights on the real structure of the network.

The results for both methods are fairly consistent over different noise variance, as it could be observed by looking at both mean and standard deviation. This suggests robustness to noise, and possible trend independence on the optimal threshold. The two methods operate on similar time scales: the analysis time is polynomial on the number of nodes in the graph $\tilde{O}(n^2)$.

One significant difference with respect to Bayesian methods is that the identified models contain dynamical information on the nature of the regulation. As such, we can distinguish between positive and negative regulation, and we can make predictions for different experimental conditions. Moreover, models can be compared to study dynamical differences: this feature will be exploited in the next section.

## III. ASSESSING BIOLOGICAL ALTERATION USING THE $\nu$-GAP

This section outlines a method to assess the effect of widespread perturbations in complex systems. The method relies on the machinery that we introduced in the previous sections.

### A. Methodology

Assuming that we have two data sets of measurements, one taken in normal and one in 'perturbed' conditions, the method follows these three steps, as shown in Fig.1:

- A network is estimated for the control data set;
- A network is estimated for the perturbed data set;
- A comparison between the two networks highlights how the perturbations affect the system.

Using the mathematical formalism introduced in the previous section, the problem of assessing regulation changes is equivalent to comparing 2 models: one that is estimated using measurements taken in normal conditions (*control data set*), and one estimated using measurements taken in the experimental conditions (*Perturbed data set*). The natural mathematical tools for this comparison are functions called 'metric' that define and measure the distance between two objects inside their space. While several metrics exist, we decided to use the $\nu$-gap ([10]), a tool developed in the control engineering community. The $\nu$-gap has been designed to address the stability properties of the closed loop system, defined using the same controller for both plants, based only on the open loop transfer function. Let $P$ be a model that explains the regulation between gene A and gene B in control conditions in a certain biological system. The genes are known to be involved in a feedback loop, represented by $C$. The system is then affected by a biological perturbation, and a novel model $P_1$ is estimated using data collected after the perturbation takes place. The $\nu$-gap measured between $P$ and $P_1$ returns a value between 0 and 1 that compare

the two systems from the controllability point of view. In particular it sets a lower bound for the difference in norm between the two closed loop transfer functions $[P, C]$ and $[P_1, C]$, where $C$ is the controller. Ideally, in a biological context this property can be used to assess how perturbing a single regulation can affect the whole feedback system. This is ideal in studying biological systems, where genes regulate each other through multiple interlocked feedback loops. In the next section we will use a toy network to define an empirical map between magnitude of perturbation and values of the $\nu$-gap that will be used as benchmark for future analysis.

### B. Effect of alterations on a toy network

To test the performance of the $\nu$-gap analysis, we used it to study the effect of known perturbations on a toy network. We opted for a toy network because parameter perturbations on the circadian clock models introduced previously, often lead to non-linearities that cannot be described using LTIs. Moreover, we wanted to investigate the potential of the $\nu$-gap approach in a more general setting, to avoid being biased towards specific network patterns.

We designed the toy network as follows:

- It follows a circular structure: each node regulates only one other, and is regulated by one other node only. We opted for this simple structure in order to guarantee stability to perturbations of the nominal parameters. Moreover, its regularity allows us to study the effect of perturbation on the neighboring nodes independently of
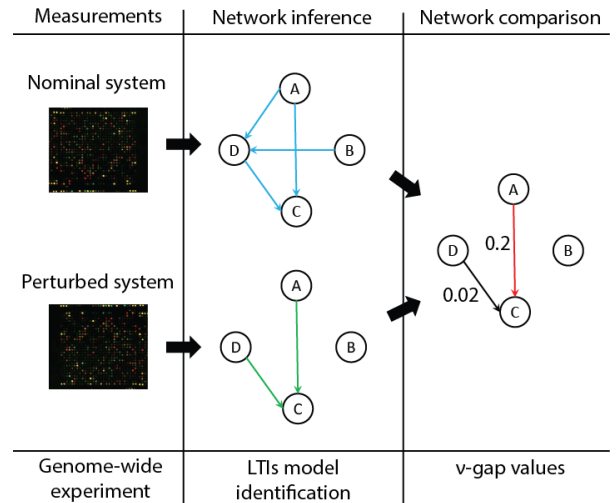


Fig. 1. **Schematic of the proposed method for measuring alterations in biological systems.** The figure summaries the strategy established in this chapter to capture alterations in biological systems: (1) Biological data are collected using a whole genome technique, like microarrays. Data are collected for the 'nominal system' and for the 'perturbed system' of interest; (2) Using LTI models, a network is inferred for both nominal and perturbed conditions data set. Causal regulations are represented by directional arrows; (3) The two estimated networks are compared and only the common nodes between the two are considered: the $\nu$-gap network measures change in regulation activity (changes are represented with red arrows, no change is represented with a black arrow. The value of the $\nu$-gap is next to each arrow). In this example, gene A is therefore a candidate for future analysis.
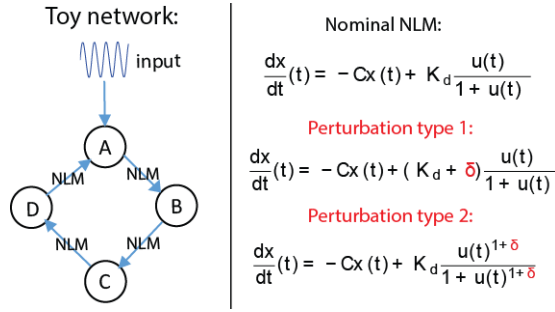
Fig. 2. **Toy network study** The linear and nonlinear networks are both defined as a 4-node circle network with one external sinusoidal input on one of the node (in the figure, the input is drawn on node A, although the node is chosen arbitrarily during the simulation). For the purpose of our analysis, we alter either the parameter $K_d$, or the exponent $n$ of one of the regulations. In the figure, the regulation that is perturbed is always B→C, although this is randomly selected in our simulations.



Fig. 3. $\nu$-gap analysis on a toy model with nonlinear dynamics: the coefficient $K_d$ was perturbed from $25\%$ to $100\%$ its nominal value.

the target regulation. For simplicity, we set the number of nodes to 4;

- Each node regulates the following one according to a nonlinear differential equation:

$$\frac{dx}{dt}(t) = -Cx(t) + K_d \frac{u(t)}{1 + u(t)}$$

For instance, if node 1 regulates node 2, then node 1 is represented by $u(t)$, and node 2 by $x(t)$. The non-linearity was chosen to be a Hill's function (with parameters $K_d$ drawn from a uniform distribution to prevent bias towards certain parameter values, and a Hill coefficient $n = 1$), since this is used to represent most biological regulations.

Being the network stable, we excited it using a sinusoid of period 24 (to simulate a circadian period as in the Millar's models) applied as external input to one randomly chosen node. The network was simulated for 100 unit times, with a sampling time of 1 time unit. The data set was then used for network inference as outlined in the previous section. The optimal thresholds were identified to maximize the previously-introduced cost function. We will call the inferred network *nominal network* hereafter.

We simulated the model twice: the first time using its nominal values, and the second time we selected one of the parameters and perturbed (increased or decreased) it randomly from $0\%$ up to $100\%$ its nominal value (granted this would lead to a stable system: perturbations that violate this conditions were discarded). We call *perturbation type 1* perturbations of the parameter $K_d$, and *perturbation type 2* those of the parameter $n$. The process is summarized in Fig.2. We run 20 simulations for each parameter perturbation ($K_d$, and $n$), and for each one estimated the nominal and the perturbed network from the simulated data using the same thresholds. Finally, we compared the two resulting networks using the $\nu$-gap. To facilitate the interpretation of the results, we split the results of the $\nu$-gap analysis in four classes: the gap measured on the target node/regulation, and the gap measured on its first, second and third neighborhood (nodes at
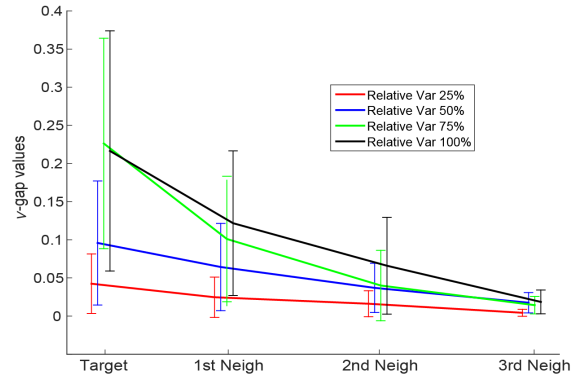
distance 1, 2 and 3 from the target node). This split analysis aims to test if the $\nu$-gap can accurately predict the regulation that was perturbed, or if it is unable to distinguish between the affected regulation and its neighboring ones. Notice that the $\nu$-gap has been measured only in the neighborhood of the frequencies excited by the input (in this case, the period of the input sinusoid): outside the frequencies of excitation, the estimated models lose reliability as they are not derived from data. A summary of the results (averaged over the 100 different runs) according to the two perturbation types is presented in Fig.3 and 4 respectively. For perturbations of $K_d$, the $\nu$-gap is in average above $0.2$ for the target node, and this holds in the case of the Hill coefficient for perturbations above $50\%$. Interestingly, the effect of perturbations on neighboring nodes is generally lower according to the $\nu$-gap, even when the original perturbation is big (for perturbation up to $100\%$ the nominal value, the maximum variation is $0.185$ for the first neighborhood, with the only exception of $0.23$ for a $100\%$ perturbation of $K_d$). Hence, this analysis suggests that values above around $0.2$ could be used to infer the main target of the perturbation. Perturbing the parameter $K_d$ as well as the exponent $n$ of the Hill's equation leads to similar results: the increase of the $\nu$-gap is proportional
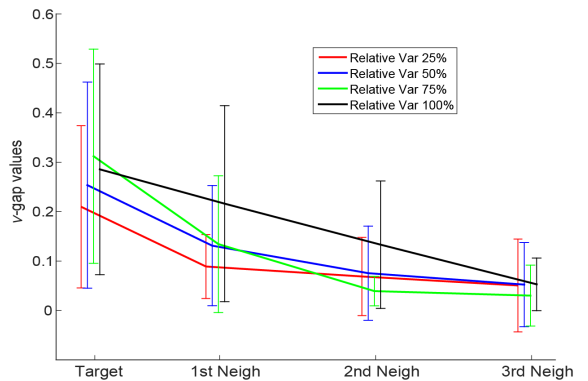


Fig. 4. $\nu$-gap analysis on a toy model with nonlinear dynamics: $n$, the exponent of the Hill function in the nonlinear network, was perturbed from $25\%$ to $100\%$ its nominal value. Perturbations at $100\%$ corresponded to a non-linearity that could not be identified using LTIs, so no values of the $\nu$-gap could be estimated for the target node.

to some degree to the amount of perturbation (although a precise mapping is not feasible). This analysis shows that the $\nu$-gap is capable of capturing the effects of network perturbation. However, it's worth noticing that our method is not capable of inferring the target regulation when the perturbed system becomes too nonlinear for LTIs to infer it. This suggests that connections that are not estimated in the perturbed case (increase in nonlinear dynamic) are potentially interesting, and should not be discarded from further considerations.

## IV. CONCLUSIONS

We outlined a methodology to analyze the effect of unknown wide-spread perturbations on unknown biological networks. The method requires as only input a data set that contains the measurements of all the signals involved in the control and in the experimental conditions.

In summary, the method follows three steps:

- The data is analyzed and all the signals that cannot be used for dynamical modeling are filtered out;
- Two networks are inferred from the filtered data: one from the control data, and one from the experimental conditions data;
- The 2 networks are compared using the $\nu$-gap, and a subset of regulations is returned as candidate target of the perturbation.

In particular, this approach was developed to analyze genome-wide experimental measurements. Consequently, it is fast, accurate, and requires little prior knowledge of the system. Moreover, it has been tested on data sets with low sampling rate and low signal-to-noise ratio in order to reproduce real experimental conditions. The method here proposed should now be tested on real data. Preliminary results on the circadian clock of A. *Thaliana* (unpublished data) and experimental validation suggest that the method can correctly identify targets of drug-induced wide-spread perturbations.

Although the methodology, as presented here, has been tuned to analyze circadian data, it has the potential to analyze any biological data set as long as the system is near linear regime. Furthermore, this technique does not impose any limitations on the nature of the experimental conditions.

## V. ACKNOWLEDGMENTS

## REFERENCES

[1] Andrej Aderhold, Dirk Husmeier, and Marco Grzegorczyk. Statistical inference of regulatory networks for circadian regulation. *Statistical Applications in Genetics and Molecular Biology*, 13(3), 2014.

[2] Amr Ahmed and Eric P Xing. Recovering time-varying networks of dependencies in social and biological studies. *Proceedings of the National Academy of Sciences*, 106(29):11878–11883, 2009.

[3] Matthew J Beal, Francesco Falciani, Zoubin Ghahramani, Claudia Rangel, and David L Wild. A bayesian approach to reconstructing genetic regulatory networks with hidden factors. *Bioinformatics*, 21(3):349–356, 2005.

[4] Matthew James Beal. *Variational algorithms for approximate Bayesian inference*. PhD thesis, University of London, 2003.

[5] Allister Bernard and Alexander J Hartemink. Evaluating algorithms for learning biological networks. 2006.

[6] Neil Dalchau, Katharine E Hubbard, Fiona C Robertson, Carlos T Hotta, Helen M Briggs, Guy-Bart Stan, Jorge M Gonçalves, and Alex AR Webb. Correct biological timing in arabidopsis requires multiple light-signaling pathways. *Proceedings of the National Academy of Sciences*, 107(29):13171–13176, 2010.

[7] Jerome Friedman, Trevor Hastie, and Robert Tibshirani. Sparse inverse covariance estimation with the graphical lasso. *Biostatistics*, 9(3):432–441, 2008.

[8] Nir Friedman, Michal Linial, Iftach Nachman, and Dana Pe'er. Using bayesian networks to analyze expression data. *Journal of computational biology*, 7(3-4):601–620, 2000.

[9] Dan Geiger and David Heckerman. Learning gaussian networks. *arXiv preprint arXiv:1302.6808*, 2013.

[10] Keith Glover, G Vinnicombe, and G Papageorgiou. Guaranteed multi-loop stability margins and the gap metric. In *Decision and Control, 2000. Proceedings of the 39th IEEE Conference on*, volume 4, pages 4084–4085. IEEE, 2000.

[11] Marco Grzegorczyk and Dirk Husmeier. A non-homogeneous dynamic bayesian network with sequentially coupled interaction parameters for applications in systems and synthetic biology. *Statistical applications in genetics and molecular biology*, 11(4), 2012.

[12] Eva Herrero, Elsebeth Kolmos, Nora Bujdoso, Ye Yuan, Mengmeng Wang, Markus C Berns, Heike Uhlworm, George Coupland, Reena Saini, Mariusz Jaskolski, et al. Early flowering4 recruitment of early flowering3 in the nucleus sustains the arabidopsis circadian clock. *The Plant Cell Online*, 24(2):428–443, 2012.

[13] Clifford M Hurvich and Chih-Ling Tsai. Bias of the corrected aic criterion for underfitted regression and time series models. *Biometrika*, 78(3):499–509, 1991.

[14] Younhee Ko, ChengXiang Zhai, and Sandra Rodriguez-Zas. Inference of gene pathways using mixture bayesian networks. *BMC systems biology*, 3(1):54, 2009.

[15] James CW Locke, László Kozma-Bognár, Peter D Gould, Balázs Fehér, Eva Kevei, Ferenc Nagy, Matthew S Turner, Anthony Hall, and Andrew J Millar. Experimental validation of a predicted feedback loop in the multi-oscillator clock of arabidopsis thaliana. *Molecular systems biology*, 2(1), 2006.

[16] Adam A Margolin, Ilya Nemenman, Katia Basso, Chris Wiggins, Gustavo Stolovitzky, Riccardo D Favera, and Andrea Califano. Aracne: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC bioinformatics*, 7(Suppl 1):S7, 2006.

[17] Edward R Morrissey, Miguel A Juárez, Katherine J Denby, and Nigel J Burroughs. Inferring the time-invariant topology of a nonlinear sparse gene regulatory network using fully bayesian spline autoregression. *Biostatistics*, 12(4):682–694, 2011.

[18] Alexandra Pokhilko, Aurora Piñas Fernández, Kieron D Edwards, Megan M Southern, Karen J Halliday, and Andrew J Millar. The clock gene circuit in arabidopsis includes a repressilator with additional feedback loops. *Molecular systems biology*, 8(1), 2012.

[19] Alexandra Pokhilko, Sarah K Hodge, Kevin Stratford, Kirsten Knox, Kieron D Edwards, Adrian W Thomson, Takeshi Mizuno, and Andrew J Millar. Data assimilation constrains new connections and components in a complex, eukaryotic circadian clock model. *Molecular systems biology*, 6(1), 2010.

[20] Simon Rogers and Mark Girolami. A bayesian regression approach to the inference of regulatory networks from gene expression data. *Bioinformatics*, 21(14):3131–3137, 2005.

[21] Juliane Schäfer and Korbinian Strimmer. An empirical bayes approach to inferring large-scale gene association networks. *Bioinformatics*, 21(6):754–764, 2005.

[22] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288, 1996.

[23] Jing Yu, V Anne Smith, Paul P Wang, Alexander J Hartemink, and Erich D Jarvis. Advances to bayesian network inference for generating causal networks from observational biological data. *Bioinformatics*, 20(18):3594–3603, 2004.

[24] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.