

Optimal Experimental Design for Gene Regulatory Networks in the Presence of Uncertainty

Roozbeh Dehghannasiri, *Student Member, IEEE*, Byung-Jun Yoon, *Senior Member, IEEE*, and Edward R. Dougherty, *Fellow, IEEE*

Abstract

Of major interest to translational genomics is the intervention in gene regulatory networks (GRNs) to affect cell behavior; in particular, to alter pathological phenotypes. Owing to the complexity of GRNs, accurate network inference is practically challenging and GRN models often contain considerable amounts of uncertainty. Considering the cost and time required for conducting biological experiments, it is desirable to have a systematic method for prioritizing potential experiments so that an experiment can be chosen to optimally reduce network uncertainty. Moreover, from a translational perspective it is crucial that GRN uncertainty be quantified and reduced in a manner that pertains to the operational cost that it induces, such as the cost of network intervention. In this work, we utilize the concept of mean objective cost of uncertainty (MOCU) to propose a novel framework for optimal experimental design. In the proposed framework, potential experiments are prioritized based on the MOCU expected to remain after conducting the experiment. Based on this prioritization, one can select an optimal experiment with the largest potential to reduce the pertinent uncertainty present in the current network model. We demonstrate the effectiveness of the proposed method via extensive simulations based on synthetic and real regulatory networks.

Index Terms

Mean objective cost of uncertainty (MOCU), experimental design, gene regulatory network (GRN), network intervention.



1 INTRODUCTION

Since the earliest days of modern science, it has been recognized that experimental design is critical for the efficient observation of nature. Today, a salient objective of translational systems biology is to determine beneficial interventions in gene regulatory networks (GRNs) for the purpose of identifying potential drug targets. A precondition for using GRNs to design intervention strategies is network identification. Hence, given a model possessing uncertainty, the aim of an experiment is to reduce that uncertainty as it pertains to the intervention objective. Optimal experimental design will depend on the model, the uncertainty, and the objective. In developing an experimental design methodology, it is insufficient to depend merely on the uncertainty, without taking into account the translational objective. Thus, entropy alone is inadequate. One needs a measure that incorporates both the uncertainty and the objective. To that end, in this paper we propose a new experimental design procedure for GRN identification based on the previously introduced mean objective cost of uncertainty (MOCU) [1].

From the earliest days of high-throughput gene-expression measurements, the intervention problem has been addressed from two perspectives: (1) dynamical intervention by altering one or more regulatory outputs (expressions) over time [2], and (2) structural intervention via a one-time change of one or more regulatory functions constituting the network [3]. Dynamical intervention interferes with signaling and does not alter network wiring, whereas structural intervention constitutes a one-time alteration of the physical

-
- *R. Dehghannasiri is with the Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX, 77843.
E-mail: roozbehdn@email.tamu.edu*
 - *B.J. Yoon is with the College of Science, Engineering, and Technology, Hamad bin Khalifa University (HBKU), Doha, Qatar, and also with the Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX, 77843, Center for Bioinformatics and Genomic Systems Engineering, Texas A&M University, College Station, TX, USA.*
 - *E.R. Dougherty is with the Department of Electrical and Computer Engineering, Texas A&M University, College Station, TX, 77843, Center for Bioinformatics and Genomic Systems Engineering, Texas A&M University, College Station, TX, USA.
E-mail: byoon@qa.orf.qa, edward@ece.tamu.edu.*

network. Both approaches have mainly developed in the context of probabilistic Boolean networks (PBNs) [4]. Structural intervention, which concerns us here, has been studied from a logical perspective to achieve a desired alteration of the attractor structure of a PBN [5] and in the framework of Markov Chain perturbation theory to derive an altered transition probability matrix that optimally reduces undesirable (pathological) steady-state probability mass [6].

In the vast majority of methods considered for both dynamical and structural intervention, the GRN is assumed to be known, which in the case of Markovian networks means that the transition probability matrix is known. However, given the complex regulatory machinery of the cell and the lack of sufficient data for accurate inference, there is typically significant uncertainty in GRN models. Hence, rather than assume that the model is fully known, it can be beneficial to assume that the true GRN belongs to an *uncertainty class* of networks and the problem is to find a *robust* intervention strategy that is optimal across the uncertainty class. In the case of dynamical intervention in PBNs, robust control policies have been found under two scenarios: (1) no knowledge is assumed concerning the distribution of the networks in the uncertainty class and optimality is defined via a *minimax criterion* – find the control policy that has the best worst-case performance across the uncertainty class [7]; and (2) there is a *prior distribution* governing the networks in the uncertainty class and optimality is defined via a *Bayesian criterion* – find the control policy that has the best expected performance across the uncertainty class [8]. Robust design has also been addressed in structural intervention, where one searches for the optimal regulatory function alteration relative to the uncertainty class [1].

It should be recognized that the uncertainty problem is inherent to computational biology owing to the complexity of biological systems and the ubiquity of samples that are small relative to the number of variables. This is why many works analyze gene regulatory networks with uncertainty [9]–[13]. In particular, not only is model uncertainty an issue that must be addressed for network intervention (optimal therapy), it is also an issue for biomarker design (optimal diagnosis) and in this context has been treated in the context of uncertainty classes of feature-label distributions [14], [15].

From an experimental perspective, one would like to reduce model uncertainty and

thereby improve intervention performance. For smaller uncertainty classes, it is more likely that the performance of a designed robust intervention strategy is close to the performance of the optimal intervention for the actual network. This brings up the issue of experimental design. In the abstract, one would like to obtain accurate estimates of all uncertain parameters; that is, fill the uncertain gaps within the network model. Unfortunately, a GRN involves a large number of variables and therefore many uncertain parameters might exist in the network model, which might require an inordinate number of experiments, where experiments can be expensive, time consuming, and dependent on unavailable specimens. The issue, then, is to evaluate potential experiments to find out which ones are most informative relative to the problem at hand. Experimental design has roots in statistics and machine learning [16]–[18]. In [16], the information gain of each experiment is measured in terms of reduction in the entropy of the model. Later, experimental design is utilized in the inference of gene regulatory networks to reduce the entropy of the network model [19]–[22]. Here we take the viewpoint that, when designing an intervention strategy we are not so much concerned with reducing model uncertainty from a general perspective, say, entropy; rather, our goal is to reduce uncertainty that will retard the effectiveness of our designed strategy.

In this paper, we present an experimental design method based on the concept of *mean objective cost of uncertainty* (MOCU), introduced in [1]. MOCU is an uncertainty quantification for dynamical models that quantifies the increased cost due to uncertainty, where the cost function depends on one's objective. In the context of controlling GRNs, MOCU measures uncertainty in terms of the differential cost between applying the robust and true-model optimal interventions. According to our proposed method, we conduct experiments to estimate unknown parameters in such a way as to maximize the expected reduction of MOCU. By computing the expected remaining MOCU after conducting each experiment, we select the experiment that results in the minimum expected remaining MOCU. We desire experiments that estimate these uncertain parameters and would like to know which experiment should be conducted first. To evaluate our proposed experimental design method, we perform simulations on both synthetic and real networks. The simulation results demonstrate the effectiveness of the proposed method.

The remainder of the paper is organized as follows. Section 2 provides an overview of Boolean networks and mean objective cost of uncertainty, and presents our proposed experimental design method. In Section 3, a comprehensive performance analysis of the proposed experimental design method for both synthetic and real networks is given. Finally, we conclude the paper in Section 4.

2 METHODS

Our interest is in designing optimal experiments for improving gene regulatory network (GRN) models possessing uncertainty. Because the optimization is computationally intensive, we restrict our examples to Boolean networks so that we can perform extensive simulations; however, the proposed experimental design method is fairly general and can be applied to different models and applications in a straightforward manner.

2.1 Boolean Networks: A Brief Overview

Boolean networks (BNs) [23] and their probabilistic extension [24] are widely accepted models for studying GRNs and have been shown to be effective in capturing much of the complex dynamics of gene regulatory networks [25]–[28]. An n -gene Boolean network (BN) is a pair (\mathbf{V}, \mathbf{F}) , where $\mathbf{V} = \{X_1, X_2, \dots, X_n\}$ is a set of nodes representing the binary expression states of genes and $\mathbf{F} = \{f_1, f_2, \dots, f_n\}$ is a set of Boolean functions such that $f_i : \{0, 1\}^{k_i} \rightarrow \{0, 1\}$ is the Boolean function that determines the expression state of X_i . It is commonplace to refer to gene i as X_i . The binary values $X_i = 0$ and $X_i = 1$ correspond to the gene being turned “off” or “on” respectively. The vector $X(t) = (X_1(t), \dots, X_n(t))$ of gene values at time t is called the *gene activity profile* (GAP). It reflects the “state” of the network at time t . The value of gene i at the next time point, $X_i(t+1) = f_i(X_{i_1}(t), X_{i_2}(t), \dots, X_{i_{k_i}}(t))$, is determined by the values of k_i predictor genes at time t . In a Boolean network with perturbation (BNp), each gene may randomly flip its value at a given time with a perturbation probability p , independently from other genes. Hence, for a BNp, $X(t+1) = \mathbf{F}(X(t))$ with probability $(1-p)^n$ when there is no perturbation, but $X(t+1)$ may take a different value with probability $1 - (1-p)^n$, when there exists one or more random perturbations. In a BNp, the sequence of states over time can be regarded as a Markov chain where transitions are made according to

a fixed transition probability matrix (TPM) \mathbf{P} . Therefore, classical Markov chain theory can be applied for analyzing network dynamics. The general formula of a TPM using Boolean functions and perturbation probability has been derived in [7]. When $p > 0$, the resulting Markov chain is ergodic, irreducible, and possesses a steady-state distribution (SSD) $\pi^T = \pi^T \mathbf{P}$, where the k th element, π_k , of the column vector π corresponds to the steady-state probability of state k and T denotes the transpose operator.

The long-run behavior of a GRN is characterized by its steady-state distribution. In the context of translational genomics, the state space of a network can typically be partitioned into undesirable states (U), corresponding to abnormal (disease) phenotypes, and desirable states (D), corresponding to normal (healthy) phenotypes. The goal in controlling GRNs via interventions is to decrease the probability that the network will enter the undesirable set of states. In other words, intervention aims at minimizing the overall steady-state probability mass $\pi_U = \sum_{i \in U} \pi_i$ in undesirable states. Structural interventions [1], [3], [5], [6] alter the long-run behavior of a network via a one-time change of the underlying network structure (wiring); dynamical interventions [2], [7], [8], [29] utilize Markov decision theory to flip (or not flip) the value of a control gene at each time instant.

In this paper, we focus on the structural intervention method proposed in [6]. In [6], intervention is performed via a rank-1 function perturbation such that the relation between the transition probability matrices of the original and perturbed networks is $\tilde{\mathbf{P}} = \mathbf{P} + \mathbf{a}\mathbf{b}^T$, where $\tilde{\mathbf{P}}$ is the transition probability matrix after perturbation and $\mathbf{a}\mathbf{b}^T$ is the rank-1 perturbation matrix, \mathbf{a} and \mathbf{b} being two arbitrary vectors, and $\mathbf{b}^T \mathbf{e} = 0$ for \mathbf{e} (all unity column vector). We use a single-gene perturbation, which is a special case of a rank-1 function perturbation in which the output state for only one input state changes and the output states of other states remain unchanged.

2.2 MOCU-based Experimental Design

Let $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ be a vector of parameters that characterizes the gene regulatory network. We assume that θ is uncertain and belongs to an uncertainty class Θ of possible networks. We refer to θ as "uncertainty vector". For any $\theta \in \Theta$, let $\xi_\theta(\psi)$ be the cost of applying the intervention $\psi \in \Psi$, a class of potential interventions, to the

network defined by the uncertainty vector θ . For instance, $\xi_\theta(\psi)$ might be the steady-state probability mass in undesirable states after applying the intervention. Let $\psi(\theta) \in \Psi$ denote an optimal intervention relative to ξ_θ , meaning that $\xi_\theta(\psi(\theta)) \leq \xi_\theta(\psi)$ for any $\psi \in \Psi$. $\psi(\theta)$ is an optimal intervention for the network with uncertainty vector θ .

An *intrinsically Bayesian robust* (IBR) intervention is defined as

$$\psi_{IBR}(\Theta) = \arg \min_{\psi \in \Psi} E_\theta \left[\xi_\theta(\psi) \right] \quad (1)$$

[1]. The expectation E_θ is taken over the probability distribution $f(\theta)$ of θ .

The *mean objective cost of uncertainty* (MOCU) relative to an uncertainty class Θ of networks and a class Ψ of interventions is defined as

$$M_\Psi(\Theta) = E_\theta \left[\xi_\theta(\psi_{IBR}(\Theta)) - \xi_\theta(\psi(\theta)) \right] \quad (2)$$

[1]. MOCU is the expected cost increase that results from applying a robust intervention over all networks in Θ instead of the optimal intervention for the true network, which is unknown.

When it is computationally infeasible to search through the class Ψ to identify an optimal IBR intervention, we can confine the search to the set of model-specific optimal interventions for networks within Θ . We define a *model-constrained Bayesian robust* (MCBR) intervention by

$$\psi_{MCBR}(\Theta) = \arg \min_{\psi(\phi): \phi \in \Theta} E_\theta \left[\xi_\theta(\psi(\phi)) \right]. \quad (3)$$

An MCBR intervention is suboptimal relative to an IBR intervention. Empirical results in [1] indicate that, at least for binary PBNs with up to ten genes, the MCBR structural intervention provides an extremely accurate approximation of the IBR structural intervention. Since the large number of MOCU computations required for the simulations performed in the current study would be computationally prohibitive using IBR intervention, we employ MCBR intervention. Historically, MCBR filtering goes back to binary filtering [30] and a general theory of IBR filtering has recently been established [31]. MCBR dynamical intervention was presented in [8]. Using the MCBR intervention, rather than the IBR intervention, we can obtain an approximation of the true MOCU in (2) by replacing the optimal IBR operator, $\psi_{IBR}(\Theta)$, by the optimal MCBR operator,

$\psi_{MCBR}(\Theta)$. In what follows, we will refer to the approximate MOCU computed based on MCBR intervention as MOCU. We will denote an optimal MCBR intervention $\psi_{MCBR}(\Theta)$ as $\psi^*(\Theta)$.

Consider a GRN possessing k uncertain parameters $\theta_1, \theta_2, \dots, \theta_k$. Suppose there exists a corresponding set of k experiments T_1, T_2, \dots, T_k , where performing experiment T_i would completely determine θ_i such that we would be sufficiently confident about the value of θ_i that we would no longer consider it to be uncertain. In practice, more than one actual experiment might be needed to be conducted for the true estimation of an uncertain parameter but we can consider these experiments collectively as one experiment for our analysis. For simplicity, let us assume that θ_i is a binary variable and that experiment T_i can determine whether $\theta_i = 0$ or $\theta_i = 1$. Our aim is to decide which experiment T_i among the k potential experiments should be conducted first in order to optimally reduce the uncertainty based on a single experiment. Let $\theta_\phi^{(i)} = \theta | (\theta_i = \phi)$ be the *conditional uncertainty vector* composed of all uncertain parameters other than θ_i , with $\theta_i = \phi$, and let $\Theta_{i,\phi} = \{\theta | \theta \in \Theta, \theta_i = \phi\}$ be the reduced uncertainty class of networks obtained by assuming that $\theta_i = \phi$. Let $M_\Psi(\Theta_{i,\phi})$ be the remaining MOCU given $\theta_i = \phi$:

$$M_\Psi(\Theta_{i,\phi}) = E_{\theta_\phi^{(i)}} \left[\xi_{\theta_\phi^{(i)}}(\psi^*(\Theta_{i,\phi})) - \xi_{\theta_\phi^{(i)}}(\psi(\theta_\phi^{(i)})) \right], \quad (4)$$

where the expectation is taken over the conditional probability distribution $f(\theta_\phi^{(i)}) = f(\theta | \theta_i = \phi)$ of the remaining uncertain parameters given $\theta_i = \phi$ and $\psi^*(\Theta_{i,\phi})$ is the optimal MCBR intervention for the reduced uncertainty class $\Theta_{i,\phi}$:

$$\psi^*(\Theta_{i,\phi}) = \arg \min_{\psi(\tau): \tau \in \Theta_{i,\phi}} E_{\theta_\phi^{(i)}} \left[\xi_{\theta_\phi^{(i)}}(\psi(\tau)) \right]. \quad (5)$$

We define the cost function by

$$\xi_{\theta_\phi^{(i)}}(\psi(\tau)) = \tilde{\pi}_{U, \theta_\phi^{(i)}}(\psi(\tau)), \quad (6)$$

where $\tilde{\pi}_{U, \theta_\phi^{(i)}}(\psi(\tau))$ is the steady-state probability mass in undesirable states after applying intervention $\psi(\tau)$ to the network defined by the uncertainty vector $\theta_\phi^{(i)}$ in the reduced uncertainty class $\Theta_{i,\phi}$. We define the expected remaining MOCU after determining the value of θ_i via experiment T_i by

$$M_\Psi(\Theta, i) = E_{\theta_i} \left[M_\Psi(\Theta_{i,\theta_i}) \right], \quad (7)$$

where the expectation is taken over the marginal probability density function, $f(\theta_i)$, for the uncertain parameter θ_i . In order to optimally reduce the uncertainty in the current uncertainty class Θ , we should select the experiment T_{i^*} such that

$$i^* = \arg \min_{i \in \{1, 2, \dots, k\}} M_{\Psi}(\Theta, i), \quad (8)$$

since T_{i^*} is expected to minimize the remaining MOCU by determining the value of the parameter θ_{i^*} .

To calculate $M_{\Psi}(\Theta, i)$, we need to define the class of interventions. Consider single-gene perturbations [6] for the class Ψ of structural interventions. Let $\tilde{\mathbf{F}} = \{\tilde{f}_1, \tilde{f}_2, \dots, \tilde{f}_n\}$ be the list of Boolean functions for the perturbed BNp. The structural intervention for input state j solely changes the output state for input state j and leaves the rest unaltered: $s = \tilde{\mathbf{F}}(j) \neq \mathbf{F}(j) = r$ and $\tilde{\mathbf{F}}(i) = \mathbf{F}(i)$ for $i \neq j$. The transition probability matrix $\tilde{\mathbf{P}}$ of the perturbed network will be identical to the transitional probability matrix \mathbf{P} of the original network, except for $\tilde{p}_{jr} = p_{jr} - (1-p)^n$ and $\tilde{p}_{js} = p_{js} + (1-p)^n$. The SSD of the perturbed BNp can be obtained by

$$\tilde{\pi}_i(j, s) = \pi_i + \frac{(1-p)^n \pi_j (z_{si} - z_{ri})}{1 - (1-p)^n (z_{sj} - z_{rj})}, \quad (9)$$

where π_i is the steady-state probability for the state i , $z_{si}, z_{ri}, z_{sj}, z_{rj}$ are elements of the fundamental matrix of the BNp, and $\tilde{\pi}_i(j, s)$ is the perturbed steady-state probability for state i after applying the aforementioned intervention [6]. The fundamental matrix of a BNp can be computed as $\mathbf{Z} = [\mathbf{I} - \mathbf{P} + \mathbf{e}\boldsymbol{\pi}^T]^{-1}$, where \mathbf{I} is the $n \times n$ identity matrix and \mathbf{e} is the all unity column vector. Let $\tilde{\pi}_{i,\theta}(j, s)$ be the steady-state probability of state i in the network with uncertainty vector θ after intervention (j, s) . Then $\tilde{\pi}_{U,\theta}(j, s) = \sum_{i \in U} \tilde{\pi}_{i,\theta}(j, s)$ is the steady-state probability mass in undesirable states after applying the single-gene perturbation structural intervention. For a BNp defined by a given uncertainty vector θ , the optimal single-gene perturbation structural intervention $(j(\theta), s(\theta))$ is the one that minimizes $\tilde{\pi}_{U,\theta}(j, s)$:

$$(j(\theta), s(\theta)) = \arg \min_{j, s \in \{1, 2, 3, \dots, 2^n\}} \tilde{\pi}_{U,\theta}(j, s) \quad (10)$$

For each network $\theta \in \Theta$, we find the optimal intervention $\psi(\theta) = (j(\theta), s(\theta))$. The MCBR intervention $\psi^*(\Theta) = (j^*(\Theta), s^*(\Theta))$ is chosen from the set $\{\psi(\theta), \theta \in \Theta\}$ such that it can minimize the expected error over the uncertainty class as shown in (3).

2.3 Simulation Set-Up

The simulations involve GRNs with genes regulated according to the commonly used majority vote rule [32]. Regulations in the network are governed by a regulatory matrix \mathbf{R} , where R_{ij} represents the regulatory relation from gene j to gene i as follows:

$$R_{ij} = \begin{cases} 1 & \text{the relation from } j \text{ to } i \text{ is activating} \\ -1 & \text{the relation from } j \text{ to } i \text{ is suppressive} \\ 0 & \text{there is no relation from } j \text{ to } i \end{cases} \quad (11)$$

A given gene takes the value 1 if the majority of its regulator genes up-regulate it and the value 0 if the majority of the predictor genes down-regulate it; otherwise, it remains unchanged. Under this rule,

$$X_i(t+1) = f_i(X(t)) = \begin{cases} 1 & \text{if } \sum_j R_{ij}X_j(t) > 0 \\ 0 & \text{if } \sum_j R_{ij}X_j(t) < 0 \\ X_i(t) & \text{if } \sum_j R_{ij}X_j(t) = 0 \end{cases} \quad (12)$$

We assume that for certain gene pairs, we are aware of the existence of regulatory relations based on prior biological knowledge; however, the precise type of regulation (i.e., activating or suppressive) may not be known. Therefore, the uncertain parameters in our simulations would be these regulatory relations. Each uncertain parameter θ_i , corresponding to an uncertain regulatory relation of an unknown type, can take on two different values: 1 for activating regulation and -1 for suppressive regulation. For a network with k uncertain regulations, the uncertainty class Θ contains 2^k potential networks that differ in one or more of these uncertain regulations. The proposed experimental design method is used to decide which uncertain parameter would be better to determine first, or equivalently, which experiment should be conducted first, in order to maximally reduce the uncertainty in the current network model and thereby optimally improve the performance of structural intervention.

After performing the optimal experiment, we are left with a smaller number of uncertain parameters that lead to a reduced uncertainty class of networks. Suppose we have performed an experiment to estimate the parameter θ_i and that the experiment has identified the true value to be $\theta_i = \mu_i$. We denote the reduced uncertainty class as Θ_{i,μ_i} and the robust intervention for this reduced uncertainty class as $\psi^*(\Theta_{i,\mu_i})$. An effective

experiment selection strategy should allow us to find out the best parameter θ_{i^*} to be determined first, such that on average the optimal robust intervention $\psi^*(\Theta_{i^*,\mu_{i^*}})$ for the reduced uncertainty class $\Theta_{i^*,\mu_{i^*}}$ would outperform other robust interventions on the true (unknown) network after identifying θ_j ($j \neq i^*$).

To illustrate the proposed experimental design strategy, consider $k = 5$ uncertain parameters in the GRN. Suppose the five potential experiments, each identifying one of the five parameters, $\theta_1, \theta_2, \dots, \theta_5$, have been ranked to obtain an ordered $\theta_{1'}, \theta_{2'}, \dots, \theta_{5'}$. Performing the experiment $T_{i'}$ leads to the identification of the unknown parameter $\theta_{i'}$ and results in the expected remaining MOCU $M_\Psi(\Theta, i')$, such that

$$M_\Psi(\Theta, 1') < M_\Psi(\Theta, 2') < \dots < M_\Psi(\Theta, 5'). \quad (13)$$

To measure the overall gain for performing the optimal experiment $T_{1'}$ relative to other suboptimal experiments, we define

$$\eta_i = \xi_\mu \left(\psi^*(\Theta_{(i+1)',\mu_{(i+1)'}}) \right) - \xi_\mu \left(\psi^*(\Theta_{1',\mu_{1'}}) \right), \quad (14)$$

where μ is the vector of true parameter values corresponding to θ . For example, η_1 denotes the difference between the cost $\xi_\mu \left(\psi^*(\Theta_{2',\mu_{2'}}) \right)$ of applying the robust intervention, derived for the reduced uncertainty class that results from conducting the second best experiment $T_{2'}$, to the true network and the cost $\xi_\mu \left(\psi^*(\Theta_{1',\mu_{1'}}) \right)$ of applying the robust intervention obtained from conducting the optimal experiment $T_{1'}$. η_i ($i = 1, 2, \dots, k-1$) quantifies the expected benefit of performing the best experiment predicted by the proposed strategy compared to other experiments, in terms of the operational cost that could be further reduced by performing the selected experiment.

3 RESULTS AND DISCUSSION

3.1 Performance Evaluation Based on Synthetic BNp

To evaluate the performance of the proposed experimental design strategy, we have performed simulations based on synthetic BNps. In our simulations, $k = 2, 3, 4, 5$ uncertain parameters are considered, assuming a uniform distribution $f(\theta)$ for all potential networks $\theta \in \Theta$. The analysis can be easily extended to other distributions. To make the simulations computationally tractable, we consider networks with six genes, X_1, \dots, X_6 ,

TABLE 1: The average gain of conducting the optimal experiment predicted by the proposed experimental design strategy in comparison to other suboptimal experiments.

	Average η_1	Average η_2	Average η_3	Average η_4
$k = 2$	0.0584	N/A	N/A	N/A
$k = 3$	0.0544	0.0718	N/A	N/A
$k = 4$	0.0545	0.0750	0.0855	N/A
$k = 5$	0.0474	0.0696	0.0803	0.0863

where each gene has three predictor genes. To generate a random BNp, we randomly select three predictor genes for each gene with uniform probability and randomly assign 1 (up-regulation) or -1 (down-regulation) to the corresponding entries in the regulatory matrix \mathbf{R} . The perturbation probability is set to $p = 0.001$. States for which $X_1 = 1$ are assumed to be undesirable, so that the set of undesirable states is $U = \{32, \dots, 63\}$. For a given k , we generate 1,000 synthetic BNps and randomly select 50 different sets of k edges (i.e., regulations) for each network. In each case, the regulatory information of other edges is retained while that of the k selected edges is assumed to be unknown.

From a translational perspective, the salient issue in evaluating an experimental design scheme using synthetic networks is controllability. Unlike real biological networks, which are controllable to a certain extent, many randomly generated networks may not be controllable. In other words, regardless of the intervention applied to the network, the SSD shift that results from the intervention may be negligible. For such networks, the difference between optimal and suboptimal experiments may be insignificant. For this reason, to examine the practical impact of experimental design, we must take controllability into account. In this work, the percentage decrease of total steady-state mass in undesirable states after intervention is used as a measure of controllability:

$$\Delta = \frac{\pi_U - \tilde{\pi}_U(j^*, s^*)}{\pi_U} \times 100\%,$$

where controllable networks have a larger Δ .

Table 1 summarizes the average gain of performing the optimal experiment predicted by the proposed strategy over other suboptimal experiments. The average is taken over

TABLE 2: The average gain of conducting the optimal experiment predicted by the proposed experimental design strategy in comparison to a randomly selected experiment.

	Average Gain
$k = 2$	0.0291
$k = 3$	0.0430
$k = 4$	0.0533
$k = 5$	0.0571

different sets of uncertain regulations and different networks with $\Delta \geq 40\%$. For $k = 2$, we calculate η_1 ; for $k = 3$, we calculate η_1 and η_2 ; and so on. As we can see in Table 1, the average gain is always positive. The results in Table 1 clearly show that the robust intervention derived from the uncertainty class reduced by conducting the optimal experiment outperforms the robust intervention that results from any other suboptimal experiment on average. We can also see that the average η_i gets larger for a larger i . For example, for $k = 4$, average $\eta_1 = 0.0545 < \text{average } \eta_3 = 0.0855$, which shows that, on average, the gain of determining $\theta_{1'}$ over $\theta_{4'}$ is larger than that of determining $\theta_{1'}$ over $\theta_{2'}$. This demonstrates that $M_\Psi(\Theta, i)$ can serve as an effective measure for prioritizing potential experiments. Furthermore, this suggests that we could expect larger gains when we compare the optimal experiment with an experiment that has a larger $M_\Psi(\Theta, i)$.

A salient question is how much we can gain by conducting an optimal experiment predicted by the proposed method over a randomly selected experiment. Since we would normally have to randomly pick an experiment unless there are reasons to prefer a specific experiment over the rest, such comparison would be useful in demonstrating the efficacy of the proposed method in a practical setting. We calculate the average gain of applying the robust intervention derived from the reduced uncertainty class obtained by conducting the optimal experiment instead of the intervention that results from a randomly chosen experiment, for all networks with $\Delta \geq 40\%$. The simulation results are shown in Table 2. It should be noted that the randomly chosen experiment may be identical to the optimal experiment (in fact, they are identical with probability $1/k$),

which is the main reason that the performance gain shown Table 2 is typically smaller than the gain shown in Table 1. For example, the average η_1 in Table 1 for $k = 2$ is almost two times the average gain for $k = 2$ in Table 2, which is due to the fact that the randomly picked experiment will be identical to the optimal experiment predicted by our method about 50% of the time. We can also see in Table 2 that the average gain increases for a larger k . For example, while the average gain for $k = 2$ is 0.0291, it is 0.0571 when $k = 5$. This implies that the performance gap between optimal and random selection is expected to increase as the uncertainty of the network increases.

As mentioned earlier, previous works for experimental design in gene regulatory networks are based on entropy reduction of the model. In [19], the information gain for each experiment is defined as the difference between the model entropy before experiment and the conditional entropy of conducting the experiment:

$$\begin{aligned} I(\Theta; T_i) &= H(\Theta) - H(\Theta|T_i) \\ &= H(\Theta) + \sum_{\theta, \phi} p(\theta, \theta_i = \phi) \log_2 p(\theta|\theta_i = \phi), \quad i = 1, 2, \dots, k, \end{aligned} \quad (15)$$

where $H(\Theta)$ is the model entropy and $p(\cdot)$ is the probability operator. The chosen experiment according to [19] is the one that maximizes (15). In our setting (uniform distribution and independent uncertain parameters), with k uncertain parameters, $H(\Theta)$ would be k and $I(\Theta; T_i)$ would be $k - 1$ for each potential experiment. Therefore, this experimental design scheme does not discriminate between potential experiments and as a result it would perform like a random selection approach. This makes sense because (15) only takes into account the stochastic properties of the model without considering the objective. Throughout this section, whenever we compare our method with the random experiment strategy, in fact, we are also comparing our method with experimental design methods based on entropy, such as [19].

We have compared $\xi_\mu(\psi^*(\Theta_{1', \mu_{1'}}))$ and $\xi_\mu(\psi^*(\Theta_{i', \mu_{i'}}))$ ($i \neq 1$) and measured the proportion of “success” (predicted optimal experiment $T_{1'}$ outperforms the suboptimal experiment $T_{i'}$), “failure” ($T_{i'}$ outperforms $T_{1'}$), and “tie” ($T_{1'}$ and $T_{i'}$ provide identical intervention performance). These results are summarized in Table 3. In this table, $\theta_{1'} \sim \theta_{i'}$ denotes the comparison between $\xi_\mu(\psi^*(\Theta_{1', \mu_{1'}}))$ and $\xi_\mu(\psi^*(\Theta_{i', \mu_{i'}}))$. When comparing $\theta_{1'} \sim \theta_{i'}$, a “tie” means that conducting either of the two experiments

TABLE 3: The proportion of success, failure, and tie of the optimal experiment predicted by the proposed strategy in comparison to other suboptimal experiments.

	$\theta_{1'} \sim \theta_{2'}$			$\theta_{1'} \sim \theta_{3'}$			$\theta_{1'} \sim \theta_{4'}$			$\theta_{1'} \sim \theta_{5'}$		
	Success	Failure	Tie									
$k = 2$	38.07%	15.29%	46.64%	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
$k = 3$	40.76%	22.32%	36.92%	42.84%	15.30%	41.86%	N/A	N/A	N/A	N/A	N/A	N/A
$k = 4$	40.97%	25.82%	33.21%	42.98%	19.21%	37.82%	43.75%	15.95%	40.30%	N/A	N/A	N/A
$k = 5$	43.00%	28.76%	28.24%	45.02%	22.62%	32.36%	45.63%	18.32%	36.05%	46.17%	15.96%	37.87%

TABLE 4: The proportion of overall success, overall failure, and overall tie of the optimal experiment predicted by the proposed strategy in comparison to all other suboptimal experiments.

	$\theta_{1'} \sim \theta_{i'} (i \neq 1)$		
	Success	Failure	Tie
$k = 2$	38.07%	15.29%	46.64%
$k = 3$	44.86%	28.33%	26.81%
$k = 4$	44.90%	37.10%	18.00%
$k = 5$	44.71%	43.13%	12.16%

results in the same intervention performance after the uncertainty reduction, a “success” means that $\xi_{\mu}(\psi^*(\Theta_{1',\mu_{1'}})) < \xi_{\mu}(\psi^*(\Theta_{i',\mu_{i'}}))$, and a “failure” means that $\xi_{\mu}(\psi^*(\Theta_{1',\mu_{1'}})) > \xi_{\mu}(\psi^*(\Theta_{i',\mu_{i'}}))$. We can see that the “success” proportion is consistently larger than the “failure” proportion, which explains why the gain in Table 1 is always positive. For $k > 2$, the proportion of “failure” decreases and the proportion of “success” increases as we compare $\psi^*(\Theta_{1',\mu_{1'}})$ with $\psi^*(\Theta_{i',\mu_{i'}})$, $i \neq 1$, for a larger i . Moreover, for $k = 2$, the proportion of “tie” is larger than that for $k > 2$. This is because the size of the uncertainty class of networks is small for $k = 2$ and therefore it is more likely that conducting either experiment yields the same robust intervention.

Table 4 shows the proportions of “overall success”, “overall failure”, and “overall tie” for the proposed experimental design strategy. Here, an “overall success” means that

$\xi_\mu(\psi^*(\Theta_{1',\mu_{1'}})) \leq \xi_\mu(\psi^*(\Theta_{i',\mu_{i'}}))$ for all $i \neq 1$ (except in the case that $\xi_\mu(\psi^*(\Theta_{1',\mu_{1'}})) = \xi_\mu(\psi^*(\Theta_{i',\mu_{i'}}))$ for all $i \neq 1$). An “overall tie” means that $\xi_\mu(\psi^*(\Theta_{1',\mu_{1'}})) = \xi_\mu(\psi^*(\Theta_{i',\mu_{i'}}))$ for all $i \neq 1$. Finally, an “overall failure” means that $\xi_\mu(\psi^*(\Theta_{1',\mu_{1'}})) > \xi_\mu(\psi^*(\Theta_{i',\mu_{i'}}))$ for at least one $i \neq 1$. As this table shows, the proportion of “overall success” is larger than that of “overall failure” for all k . The proportion of “tie” decreases with increasing k , as the size of the uncertainty class of networks increases. While the proportion of “overall tie” decreases with increasing k , the proportion of “overall failure” increases. This is intuitive, since by increasing the number of uncertain regulations k , it becomes more difficult to have an “overall success”, since $\psi^*(\Theta_{1',\mu_{1'}})$ has to outperform all other robust interventions, whose number increases with k .

Now, let us consider the difference between the expected remaining MOCU of the optimal experiment and that of a suboptimal experiment:

$$\Delta MOCU = M_\Psi(\Theta, i') - M_\Psi(\Theta, 1')$$

for $i' \neq 1'$. Figure 1 shows the empirical conditional expectation, $E[\eta_i | \Delta MOCU]$, of η_i ($i = 1, 2, 3, 4$) given $\Delta MOCU$ estimated based on all random networks with $\Delta \geq 40\%$. The average gain is positive for all $\Delta MOCU$. This shows that, on average, the robust interventions obtained by conducting the optimal experiments predicted by our proposed method outperform the robust interventions obtained from other suboptimal experiments when applied to the true network. Moreover, as $\Delta MOCU$ increases, the average gain increases in a more or less linearly proportional manner. Another interesting observation is that $E[\eta_i | \Delta MOCU]$ does not significantly differ for different i . This result is intuitive, since we expect the gain to depend on the estimated $\Delta MOCU$, and not the predicted rank of the suboptimal experiment.

To see how the controllability Δ measured in terms of the SSD shift that can be achieved by optimal intervention, affects the average gain of the proposed experimental design strategy, we compute the average η_i ($i = 1, 2, 3, 4$) for random networks whose controllability (i.e., Δ) exceeds a certain minimum value, where we consider minimum Δ ranging between 0% and 90%. According to Fig. 2, the average gain of η_i increases as the minimum Δ increases, regardless of i and the number of uncertain regulations k . For example, for $k = 5$ uncertain regulations in the network, the average gain based

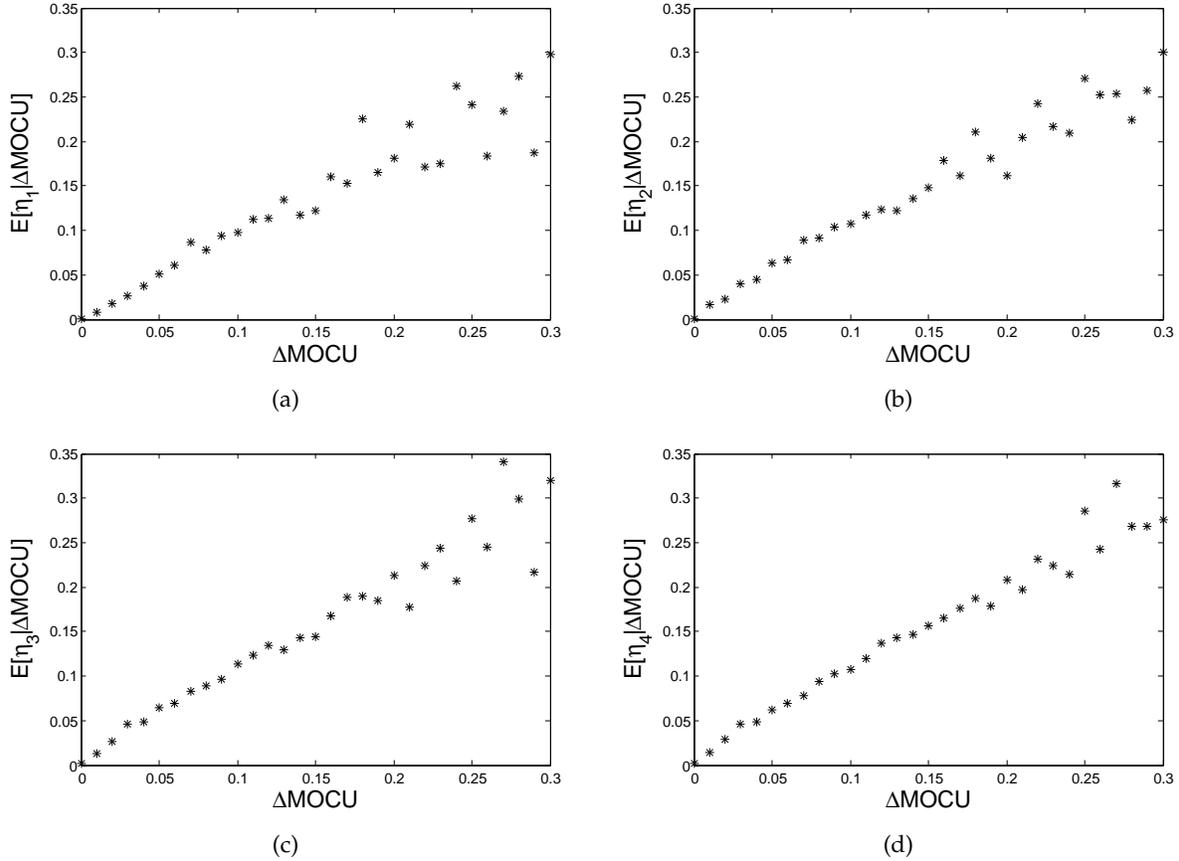


Fig. 1: The empirical conditional expectation of the gain $E[\eta_i | \Delta MOCU]$ given the difference in MOCU between the optimal and suboptimal experiments. Synthetic BNps with five uncertain regulations are considered. (a) $\theta_{1'} \sim \theta_{2'}$. (b) $\theta_{1'} \sim \theta_{3'}$. (c) $\theta_{1'} \sim \theta_{4'}$. (d) $\theta_{1'} \sim \theta_{5'}$.

on comparing $\theta_{1'}$ to $\theta_{5'}$ is slightly below 0.07 for all networks, but it increases to almost 0.1 when we consider only highly controllable networks with $\Delta \geq 90\%$.

Figure 3 compares the performance of the proposed experimental design method and that of the random selection approach based on a sequence of experiments. Assuming $k = 5$ uncertain regulations in each network, we perform 5 consecutive experiments until the network does not contain any uncertainty. First, we consider adopting the proposed experimental design strategy, where at each step, we select the optimal experiment predicted by our method, conduct the experiment to reduce the uncertainty class, and repeat this process until the network is fully identified. For comparison, we

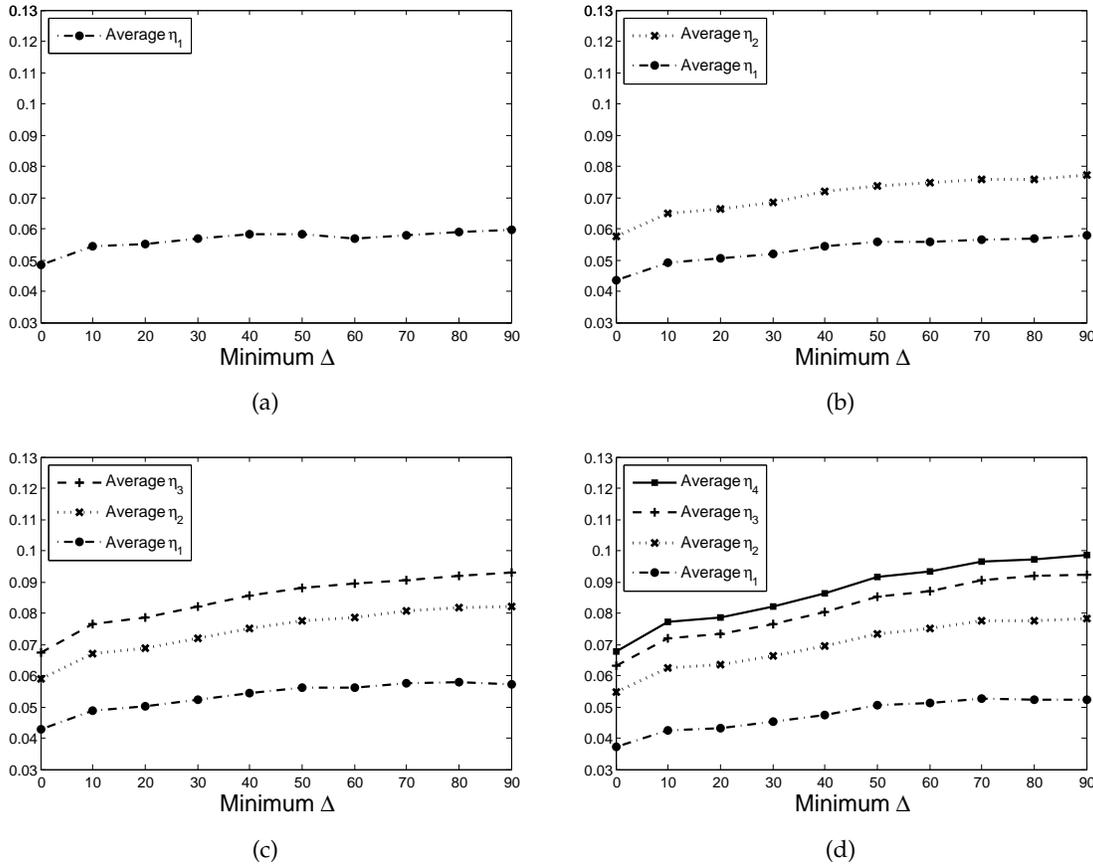


Fig. 2: Effect of the controllability of the synthetic BNp on the average performance gain of the proposed experimental design method. (a) Networks with 2 uncertain regulations. (b) Networks with 3 uncertain regulations. (c) Networks with 4 uncertain regulations. (d) Networks with 5 uncertain regulations.

perform similar simulations by conducting a randomly selected experiment at each step until there is no uncertainty about the network. In both cases, the network will be fully identified after conducting 5 experiments. To compare the performance of the two approaches, after conducting each experiment, we derive the robust intervention based on the reduced network class, apply it to the true (unknown) network, and measure the cost of intervention (i.e., total steady-state mass in undesirable states). The average performance is estimated based on 1,000 synthetic BNps and 50 different sets of uncertain regulations Δ for each of these networks. Let ψ_{opt}^* denote the robust intervention obtained by taking the proposed strategy and let ψ_{rnd}^* denote the robust

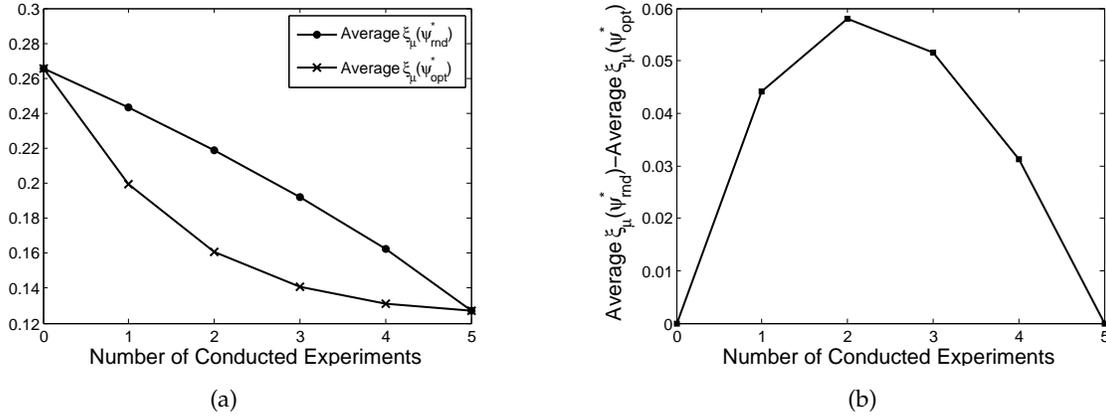


Fig. 3: Performance comparison based on a sequence of experiments. (a) The average cost of robust intervention after performing the sequence of experiments predicted by the proposed strategy and the average cost after performing randomly selected experiments. (b) The performance difference between the proposed approach and the random selection approach.

intervention obtained by performing randomly selected experiments. As seen in Fig. 3a, the curves corresponding to these two methods begin and end at the same average cost, but the curve that corresponds to the proposed experimental design strategy drops much more sharply at the beginning compared to the random selection approach. This clearly demonstrates the effectiveness of the proposed method in reducing the network uncertainty. Figure 3b plots the difference between the average $\xi_{\mu}(\psi_{rnd}^*)$ and the average $\xi_{\mu}(\psi_{opt}^*)$. In both figures, the performance difference is especially large for the first few experiments and $\xi_{\mu}(\psi_{opt}^*)$ quickly approaches the minimum cost attained by the optimal intervention. This fast convergence is important, considering the difficulty of performing a large number of experiments in real applications.

3.2 Computational Complexity Analysis

The proposed experimental design method is objective-based and in this paper the objective is network intervention. Therefore, the computational burden of the design method is mainly based on the associated network intervention strategy. A salient issue for network intervention methods is their inherent computational complexity [33]–

TABLE 5: The approximate processing time (seconds) needed for the proposed experimental design for different number of uncertain regulations k and different network size n .

	$k = 2$	$k = 3$	$k = 4$	$k = 5$	$k = 6$
$n = 5$	0.9	1.7	3.4	6.8	13.7
$n = 6$	5	10	20	39	78
$n = 7$	34	67	133	268	534
$n = 8$	260	500	1000	2000	4000

[35]. The complexity of network intervention grows exponentially with network size. Computational complexity for experimental design is much greater because we need to find the optimal intervention for every potential network inside the uncertainty class. Owing to this high computational complexity, our future work will focus on finding reliable approximations for the expected remaining MOCU and reducing network size in a way that has minor effect on the decision making process for the optimal experiment.

Suppose a BNp has n genes and k uncertain parameters, where each uncertain parameter can take m possible values and states from 2^{n-1} to $2^n - 1$ are assumed to be undesirable. Then there are m^k networks inside the uncertainty class and for each network we need to compute equation (9) $2^n \times 2^n \times 2^{n-1}$ times to find the optimal structural intervention. Therefore, the complexity of the proposed experimental design method is $O(m^k \times 2^{3n-1})$. The complexity grows exponentially with n and k , and grows polynomially with m . The complexity is highly dependent on the network size n . By adding one gene to the network, the complexity is multiplied by 8. Table 5 compares the approximate processing time needed to run the proposed experimental design method under the same setting as Section 3.1 for each set of k uncertain regulations in an n -gene BNp on a 2.9 GHz Intel Xeon, 4 GB RAM machine with MATLAB implementation.

3.3 Performance Evaluation Based on the Mammalian Cell Cycle Network

In this section, we evaluate the performance of the proposed experimental design strategy based on the mammalian cell cycle network. The cell cycle involves a sequence

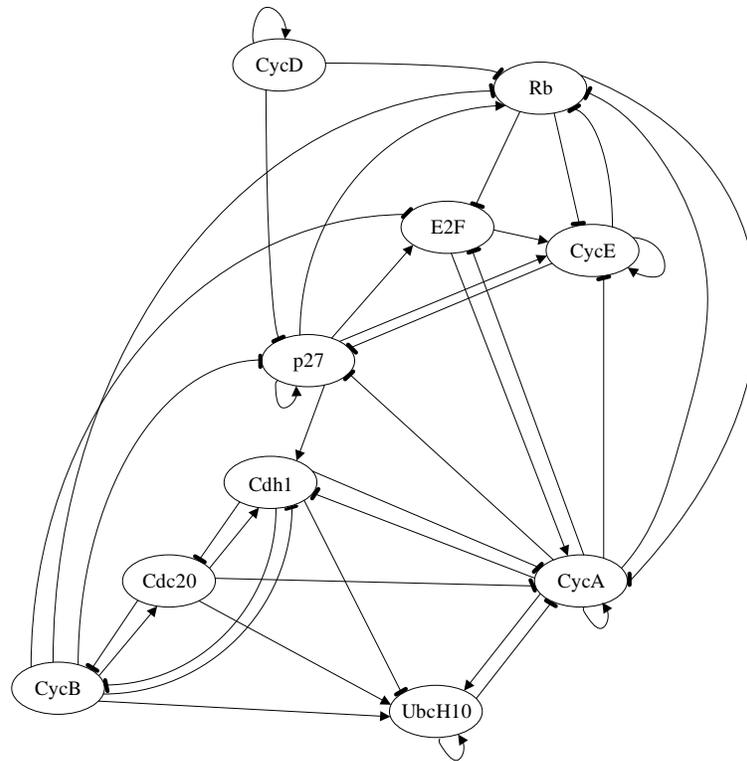


Fig. 4: A gene regulatory network model of the mammalian cell cycle. Normal arrows represent activating regulations and blunt arrows represent suppressive regulations.

of events resulting in the duplication and division of the cell. It occurs in response to growth factors and, under normal conditions, it is a tightly controlled process. A regulatory model for the mammalian cell cycle is proposed in [36]. This model contains 10 genes: CycD, Rb, p27, E2F, CycE, CycA, Cdc20, Cdh1, UbcH10, and CycB. We represent this gene regulatory network by a BN p , where the perturbation probability is set to $p = 0.001$ and genes are numbered in the previous order. The regulatory model for this network is shown in Fig. 4. The blunt arrows represent suppressive regulations and the normal arrows represent activating regulations. The cell cycle in mammals is controlled via extra-cellular stimuli. Positive stimuli activate Cyclin D (CycD) in the cell leading to cell division. CycD inactivates Rb protein, a tumor suppressor, via phosphorylation. When gene p27 and either CycE or CycA are active, the cell cycle stops, because Rb can be expressed even in the presence of cyclins. States in which the cell cycle continues even in the absence of stimuli are associated with cancerous phenotypes. For this reason,

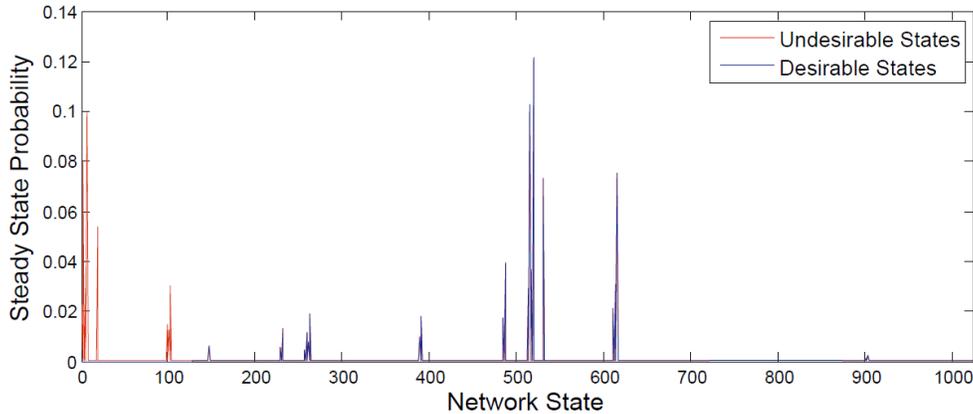


Fig. 5: The steady state distribution of the mammalian cell cycle network modeled by a BNp with perturbation probability $p = 0.001$.

we regard states with down-regulated CycD, Rb, and p27 ($X_1 = X_2 = X_3 = 0$) as undesirable states, representing cancerous phenotypes. The network SSD is shown in Fig. 5, where the total undesirable steady-state mass is $\pi_U = 0.3461$ without intervention. Suppose we want to reduce the steady-state probability mass of the set of undesirable states, $U = \{0, \dots, 127\}$, via structural intervention. The optimal intervention is to change the transition from the input state 0000000111 to the output state 1110001011 by perturbing the regulatory function such that $\tilde{\mathbf{F}}(0000000111) = 1110001011$. For all other states, their output states remain unchanged after the intervention.

To evaluate the proposed experimental design method based on the given network, we again assume that k ($= 2, 3, 4, 5$) regulations are unknown. For each k , we randomly select 50 different sets of k regulations from the network, for which we assume their regulatory information is not known, and apply the experimental design strategy to predict the optimal experiment to be performed. Table 6 summarizes the average gain of the predicted optimal experiment over other suboptimal experiments for different values of k . The average gain is positive in all cases, as in our simulations based on synthetic BNps. Furthermore, the average gain η_i increases with i . For example, when $k = 5$,

$$\text{average } \eta_4 > \text{average } \eta_3 > \text{average } \eta_2 > \text{average } \eta_1.$$

Table 7 shows the proportion of “success”, “failure”, and “tie” for applying the

TABLE 6: The average gain of conducting the optimal experiment predicted by the proposed experimental design strategy in comparison to other suboptimal experiments. The 10-gene mammalian cell cycle network with k unknown regulations are considered.

	Average η_1	Average η_2	Average η_3	Average η_4
$k = 2$	0.0208	N/A	N/A	N/A
$k = 3$	0.0207	0.0261	N/A	N/A
$k = 4$	0.0217	0.0337	0.0379	N/A
$k = 5$	0.0365	0.0389	0.0395	0.0425

TABLE 7: The proportion of success, failure, and tie of the optimal experiment predicted by the proposed strategy in comparison to other suboptimal experiments. The 10-gene mammalian cell cycle network with k unknown regulations are considered.

	$\theta_{1'} \sim \theta_{2'}$			$\theta_{1'} \sim \theta_{3'}$			$\theta_{1'} \sim \theta_{4'}$			$\theta_{1'} \sim \theta_{5'}$		
	Success	Failure	Tie									
$k = 2$	40.00%	24.00%	36.00%	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
$k = 3$	52.00%	30.00%	18.00%	54.00%	26.00%	20.00%	N/A	N/A	N/A	N/A	N/A	N/A
$k = 4$	48.00%	24.00%	28.00%	56.00%	16.00%	28.00%	60.00%	12.00%	28.00%	N/A	N/A	N/A
$k = 5$	56.00%	26.00%	18.00%	60.00%	20.00%	20.00%	60.00%	14.00%	26.00%	68.00%	10.00%	22.00%

proposed experimental design strategy. The results based on the mammalian cell cycle network are consistent with the results obtained from the synthetic networks. The “success” rate is consistently and significantly higher than the “failure” rate in all cases, thereby demonstrating the effectiveness of the proposed method. The proportion of “success” increases when we compare the optimal experiment with an experiment with larger $M_{\Psi}(\Theta, i')$ (i.e., for larger i'), which shows that the MOCU provides a sound mathematical basis for predicting the effectiveness of potential experiments.

3.4 Performance Evaluation Based on a p53 Network

We now investigate performance of the proposed experimental design method on a p53 network [37]. p53 is a tumor suppressor gene which plays a major role in DNA

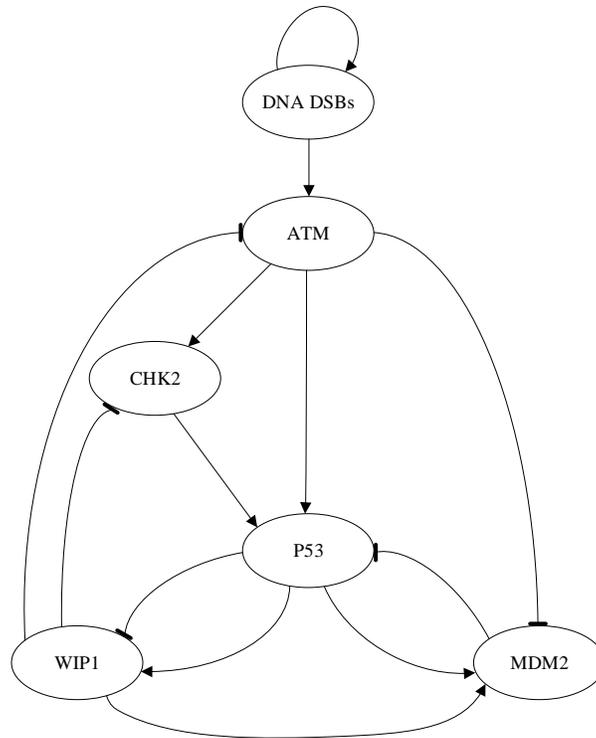


Fig. 6: A gene regulatory model for the p53 network. Normal arrows represent activating regulations and blunt arrows represent suppressive regulations.

damage regulation and programmed cell death (apoptosis). It has been observed that p53 is mutated in 30-50% of commonly occurring human cancers [38]. Under normal conditions, the expression level of p53 remains low via the control of MDM2, an oncogene that is often highly expressed in tumor cells. When DNA damage occurs, p53 is up-regulated and either activates other genes involved in DNA repair or it initiates apoptosis. Figure 6 shows key pathways that involve the regulation of p53 (see [39] for a detailed dynamical analysis of a very similar network). In the model of Fig. 6, when a DNA double strand break occurs, DNA DSBs becomes 1. The model contains five genes: MDM2, p53, WIP1, CHK2, and ATM.

Like the mammalian cell cycle network, we model the p53 network as a BNp with perturbation probability $p = 0.001$ and 6 nodes: X_1 (DNA DSBs), X_2 (MDM2), X_3 (p53), X_4 (WIP1), X_5 (CHK2), and X_6 (ATM). The presence of DNA damage ($X_1 = 1$), permanently up-regulated MDM2 ($X_2 = 1$) and permanently down-regulated p53 ($X_3 = 0$)

TABLE 8: The average gain of conducting the optimal experiment predicted by the proposed experimental design strategy in comparison to other suboptimal experiments. The 6-gene P53 network with k unknown regulations is considered.

	Average η_1	Average η_2	Average η_3	Average η_4
$k = 2$	0.0386	N/A	N/A	N/A
$k = 3$	0.0466	0.0434	N/A	N/A
$k = 4$	0.0343	0.0489	0.0657	N/A
$k = 5$	0.0387	0.0597	0.0622	0.0632

would result in an abundance of cancerous cells. For example, TCGA studies on 138 patients with glioblastoma (a kind of brain tumor) have shown that 32% and 12% of them had mutated p53 and MDM2 genes, respectively. Therefore, states with $X_1 = 1$, $X_2 = 1$, and $X_3 = 0$ are considered as the undesirable states; i.e., $U = \{48, \dots, 55\}$. The steady-state probability mass of undesirable states before and after optimal structural intervention is 0.3478 and 0.0289, respectively. Our simulations use the same settings as for the mammalian cell cycle network analysis.

Table 8 shows the average gain of conducting optimal experiments instead of other suboptimal experiments. The average gain is always positive and in most cases average η_i increases with i . However, there is an anomaly for $k = 3$, where average η_1 is larger than average η_2 . Because the average results are obtained based on a single network and 50 different selections of uncertain parameters, we should expect such occurrences since we are not averaging over a large set of simulations as with the synthetic networks. Table 9 evaluates the performance of the predicted optimal experiments in terms of percentages of “success”, “failure”, and “tie”. The “success” percentage is always larger than the “failure” percentage and it becomes larger when we compare the optimal experiment with an experiment corresponding to a larger i' . Again, there are a few anomalies, such as the decrease of “success” percentage for $k = 3$ and $k = 4$ when the optimal experiment is compared against the second and third optimal experiments – again not surprising given the small number of observations.

TABLE 9: The proportion of success, failure, and tie of the optimal experiment predicted by the proposed strategy in comparison to other suboptimal experiments. The 6-gene P53 network with k unknown regulations is considered.

	$\theta_{1'} \sim \theta_{2'}$			$\theta_{1'} \sim \theta_{3'}$			$\theta_{1'} \sim \theta_{4'}$			$\theta_{1'} \sim \theta_{5'}$		
	Success	Failure	Tie									
$k = 2$	26.00%	8.00%	66.00%	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A	N/A
$k = 3$	34.00%	6.00%	60.00%	30.00%	4.00%	66.00%	N/A	N/A	N/A	N/A	N/A	N/A
$k = 4$	46.00%	24.00%	30.00%	44.00%	16.00%	40.00%	52.00%	14.00%	34.00%	N/A	N/A	N/A
$k = 5$	62.00%	18.00%	20.00%	62.00%	8.00%	30.00%	64.00%	8.00%	28.00%	66.00%	4.00%	30.00%

4 CONCLUSION

Prioritization of potential experiments is of great practical import in systems biology and translational medicine. In this work, we have proposed a novel framework for evaluating the expected impact of a potential experiment in reducing the amount of uncertainty present in a dynamic network model. We estimate the mean objective cost of uncertainty expected to remain after conducting a specific experiment and select the one expected to optimally reduce network uncertainty. Extensive simulations based on both synthetic and actual networks show that the proposed experimental design strategy significantly outperforms random selection. Since computational complexity is an impediment for large uncertainty classes, we are currently investigating two approaches to complexity reduction. One is to discover heuristics that can be used to efficiently compute an approximate MOCU that preserves the ranking of potential experiments. A second approach is network reduction. Here the situation is analogous to the reduction of GRNs to facilitate design of optimal controllers [33]–[35], except that reduction must be accomplished in such a way as to preserve (to the extent possible) the MOCU calculations.

Finally, it is worth noting that the problem considered in this work bears conceptual similarity to the online learning problems that have been gaining broad interest in recent years. In online learning, sequential measurements are made, one at a time, to improve an uncertain model. The online knowledge gradient (KG) algorithm is an

interesting example that deals with a general class of such online learning problems [40]. It is assumed that one of M alternatives can be measured at each time step, which yields a random reward with an unknown mean and known variance (corresponding to measurement error). The main goal is to make sequential measurements that will maximize the expected total reward to be collected over a time period. To achieve this goal, in every time step, one tries to identify the optimal KG policy that will allow one to choose the single best measurement (among the M available alternatives) that is expected to bring forth the largest improvement. The alternative measurements (or rewards) are typically assumed to be independent Gaussian random variables, but one can incorporate prior beliefs about the measurements and their correlations into the problem via their joint distribution. Although the online learning problem and the aforementioned KG algorithm bear some conceptual similarities to the sequential experimental design problem considered in this paper and our MOCU-based strategy, there are critical differences. For example, our approach does not require direct modeling of the distribution of the reward (or cost). Instead, we focus on the uncertainty regarding the underlying network as it pertains to the cost of the operation of interest. Even though our ultimate goal is minimizing the cost, it is indirectly attained by optimally improving our knowledge regarding the network in a way that is pertinent to the operation (and its cost) to be performed based on the network.

REFERENCES

- [1] B. J. Yoon, X. Qian, and E. R. Dougherty, "Quantifying the objective cost of uncertainty in complex dynamical systems," *Signal Processing, IEEE Transactions on*, vol. 61, no. 9, pp. 2256–2266, 2013.
- [2] I. Shmulevich, E. R. Dougherty, and W. Zhang, "Gene perturbation and intervention in probabilistic boolean networks," *Bioinformatics*, vol. 18, no. 10, pp. 1319–1331, 2002.
- [3] —, "Control of stationary behavior in probabilistic boolean networks by means of structural intervention," *Biological Systems*, vol. 10, no. 4, pp. 431–446, 2002.
- [4] —, "From boolean to probabilistic boolean networks as models of genetic regulatory networks," *Proceedings of the IEEE*, vol. 90, no. 11, pp. 1778–1792, 2002.
- [5] Y. Xiao and E. R. Dougherty, "The impact of function perturbations in boolean networks," *Bioinformatics*, vol. 23, no. 10, pp. 1265–1273, 2007.
- [6] X. Qian and E. R. Dougherty, "Effect of function perturbation on the steady-state distribution of genetic regulatory networks: Optimal structural intervention," *Signal Processing, IEEE Transactions on*, vol. 56, no. 10, pp. 4966–4976, 2008.

- [7] R. Pal, A. Datta, and E. R. Dougherty, "Robust intervention in probabilistic boolean networks," *Signal Processing, IEEE Transactions on*, vol. 56, no. 3, pp. 1280–1294, 2008.
- [8] —, "Bayesian robustness in the control of gene regulatory networks," *Signal Processing, IEEE Transactions on*, vol. 57, no. 9, pp. 3667–3678, 2009.
- [9] G. Chesi, "On the steady states of uncertain genetic regulatory networks," *Systems, Man and Cybernetics, Part A: Systems and Humans, IEEE Transactions on*, vol. 42, no. 4, pp. 1020–1024, July 2012.
- [10] Z. Wang, H. Wu, J. Liang, J. Cao, and X. Liu, "On modeling and state estimation for genetic regulatory networks with polytopic uncertainties," *NanoBioscience, IEEE Transactions on*, vol. 12, no. 1, pp. 13–20, March 2013.
- [11] F.-X. Wu, "Global and robust stability analysis of genetic regulatory networks with time-varying delays and parameter uncertainties," *Biomedical Circuits and Systems, IEEE Transactions on*, vol. 5, no. 4, pp. 391–398, Aug 2011.
- [12] S. Z. Denic, B. Vasic, C. D. Charalambous, and R. Palanivelu, "Robust control of uncertain context-sensitive probabilistic boolean networks," *Systems Biology, IET*, vol. 3, no. 4, pp. 279–295, July 2009.
- [13] M. S. Esfahani, B.-J. Yoon, and E. R. Dougherty, "Probabilistic reconstruction of the tumor progression process in gene regulatory networks in the presence of uncertainty." *BMC Bioinformatics*, vol. 12, no. S-10, p. S9, 2011.
- [14] L. A. Dalton and E. R. Dougherty, "Bayesian minimum mean-square error estimation for classification error - part i: Definition and the bayesian mmse error estimator for discrete classification," *Signal Processing, IEEE Transactions on*, vol. 59, no. 1, pp. 115–129, 2011.
- [15] —, "Optimal classifiers with minimum expected error within a bayesian framework part i: Discrete and gaussian models," *Pattern Recognition*, vol. 46, no. 5, pp. 1288–1300, 2013.
- [16] A. C. Atkinson and A. N. Donev, *Optimum Experimental Designs*. Oxford: Oxford University Press, 1992.
- [17] H. Raiffa and R. Schlaifer, *Applied statistical decision theory*. Division of Research, Graduate School of Business Administration, Harvard University, 1961.
- [18] V. V. Fedorov, *Theory Of Optimal Experiments*. Elsevier Science, 1972.
- [19] C.-H. Yeang, H. C. Mak, S. McCuine, C. Workman, T. Jaakkola, and T. Ideker, "Validation and refinement of gene-regulatory pathways on a network of physical interactions," *Genome Biology*, vol. 6, no. 7, p. R62, 2005.
- [20] T. E. Ideker, V. Thorsson, and R. M. Karp, "Discovery of regulatory interactions through perturbation: Inference and experimental design," in *Proceedings of the Pacific Symposium on Biocomputing, 2000*, pp. 302–313.
- [21] R. D. King, K. E. Whelan, F. M. Jones, P. G. K. Reiser, C. H. Bryant, S. H. Muggleton, D. B. Kell, and S. G. Oliver, "Functional genomic hypothesis generation and experimentation by a robot scientist," *Nature*, no. 6971, p. 247252, 2004.
- [22] A. Almudevar and P. Salzman, "Using a bayesian posterior density in the design of perturbation experiments for network reconstruction," in *Proceedings of the IEEE Symposium on Computational Intelligence in Bioinformatics and Computational Biology*, Nov. 2005, pp. 1–7.
- [23] S. A. Kauffman, *The Origins of Order*. Oxford: Oxford University Press, 1993.
- [24] I. Shmulevich, E. R. Dougherty, S. Kim, and W. Zhang, "Probabilistic boolean networks: a rule-based uncertainty model for gene regulatory networks," *Bioinformatics*, vol. 18, no. 2, pp. 261–274, 2002.
- [25] R. Albert and H. G. Othmer, "The topology of the regulatory interactions predicts the expression pattern of the segment polarity genes in drosophila melanogaster," *Journal of Theoretical Biology*, vol. 223, no. 1, pp. 1 – 18, 2003.

- [26] S. Kauffman, C. Peterson, B. Samuelsson, and C. Troein, "Random boolean network models and the yeast transcriptional network," *Proceedings of the National Academy of Sciences*, vol. 100, no. 25, pp. 14 796–14 799, 2003.
- [27] M. I. Davidich and S. Bornholdt, "Boolean network model predicts cell cycle sequence of fission yeast," *PLoS ONE*, vol. 3, no. 2, 2008.
- [28] K.-Y. Lau, S. Ganguli, and C. Tang, "Function constrains network architecture and dynamics: A case study on the yeast cell cycle boolean network," *Phys. Rev. E*, vol. 75, p. 051907, May 2007.
- [29] G. Vahedi, B. Faryabi, J. Chamberland, A. Datta, and E. R. Dougherty, "Intervention in gene regulatory networks via a stationary mean-first-passage-time control policy," *Biomedical Engineering, IEEE Transactions on*, vol. 55, no. 10, pp. 2319–2331, Oct 2008.
- [30] A. M. Grygorian and E. R. Dougherty, "Design and analysis of robust optimal binary filters in the context of a prior distribution for the states of nature," *Mathematical Imaging and Vision*, vol. 11, no. 3, pp. 239–254, 1999.
- [31] L. A. Dalton and E. R. Dougherty, "Intrinsically optimal bayesian robust filtering," *Signal Processing, IEEE Transactions on*, vol. 62, no. 3, pp. 657–670, 2014.
- [32] K. Lau, S. Ganguli, and C. Tang, "Function constrains network architecture and dynamics: A case study on the yeast cell cycle boolean network," *Phys Rev E*, vol. 75, p. 051907, 2007.
- [33] X. Qian, N. Ghaffari, I. Ivanov, and E. R. Dougherty, "State reduction for network intervention in probabilistic Boolean networks," *Bioinformatics*, vol. 26, no. 24, pp. 3098–3104, Dec 2010.
- [34] I. Ivanov, P. Simeonov, N. Ghaffari, X. Qian, and E. R. Dougherty, "Selection policy-induced reduction mappings for boolean networks," *Signal Processing, IEEE Transactions on*, vol. 58, no. 9, pp. 4871–4882, Sept 2010.
- [35] I. Ivanov, R. Pal, and E. R. Dougherty, "Dynamics preserving size reduction mappings for probabilistic boolean networks," *Signal Processing, IEEE Transactions on*, vol. 55, no. 5, pp. 2310–2322, May 2007.
- [36] A. Faure, A. Naldi, C. Chaouiya, and D. Thieffry, "Dynamical analysis of a generic Boolean model for the control of the mammalian cell cycle," *Bioinformatics*, vol. 22, no. 14, pp. 124–131, 2006.
- [37] E. Batchelor, A. Loewer, and G. Lahav, "The ups and downs of p53: understanding protein dynamics in single cells," *Nature Reviews Cancer*, vol. 9, no. 5, pp. 371–377, Apr. 2009.
- [38] R. A. Weinberg, *The Biology of Cancer*. New York: Garland Science, 2007.
- [39] R. K. Layek, A. Datta, and E. R. Dougherty, "From biological pathways to regulatory networks," *Mol. BioSyst.*, vol. 7, pp. 843–851, 2011.
- [40] I. O. Ryzhov, W. B. Powell, and P. I. Frazier, "The knowledge gradient algorithm for a general class of online learning problems," *Operations Research*, vol. 60, no. 1, pp. 180–195, 2012.