



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Optimally Designed vs Intuition-Driven Inputs: The Study Case of Promoter Activity Modelling (I)

### Citation for published version:

Bandiera, L, Kothamachu, VB, Balsa-Canto, E, Swain, P & Menolascina, F 2018, Optimally Designed vs Intuition-Driven Inputs: The Study Case of Promoter Activity Modelling (I). in Proceedings of the 59th Conference on Decision and Control. 59th IEEE Conference on Decision and Control, Seogwipo, Korea, Republic of, 8/12/20.

### Link:

[Link to publication record in Edinburgh Research Explorer](#)

### Published In:

Proceedings of the 59th Conference on Decision and Control

### General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

### Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Optimally designed vs intuition-driven inputs: the study case of promoter activity modelling

L. Bandiera<sup>1,2</sup>, V. Kothamachu<sup>1</sup>, E. Balsa-Canto<sup>3</sup>, P. S. Swain<sup>2</sup> and F. Menolascina<sup>1,2</sup>

**Abstract**—Synthetic biology is an emerging engineering discipline that aims at synthesising logical circuits into cells to accomplish new functions. Despite a thriving community and some notable successes, the basic task of assembling predictable gene circuits is still a key challenge. Mathematical models are uniquely suited to help solve this issue. Yet in biology they are perceived as expensive and laborious to obtain because *low-information* experiments have often been used to infer model parameters. How much additional information can be gained using optimally designed experiments? To tackle this question we consider a building block in Synthetic Biology, an inducible promoter in yeast *S. cerevisiae*. Using *in vivo* data we re-fit a mathematical model for such a system; we then compare *in silico* the quality of the parameter estimates when model calibration is done using typical (e.g. step inputs) and optimally designed experiments. We find that Optimal Experimental Design leads to  $\sim 70\%$  improvement in the predictive ability of the inferred models. We conclude providing suggestions on how optimally designed experiments can be implemented *in vivo*.

## I. INTRODUCTION

Synthetic Biology is an emerging discipline that seeks to implement *de novo* tasks in cells. Despite a booming community and the opportunities that Synthetic Biology offers [1], the assembly of synthetic circuits with predictable functions remains a challenge. If Synthetic Biology is to advance towards application, it is necessary to increase the predictability of gene network dynamics. Mathematical models offer a means to achieve this goal, yet their use in Synthetic Biology has so far only been limited [2]. The reason for the low adoption of models largely lies in the limitations of traditional experimental platforms in biology (e.g. microplate readers). “Pulses” or “steps” of chemicals are the *de facto* standard stimuli to probe cell behaviour. Such designs, however, often allow for poorly-informative experiments. Indeed, chemical “steps” and “pulses” are low-pass filtered by molecular diffusion, which limits their frequency content. Furthermore, the inherent non-linearity of biological networks prevents the adoption of results on *persistent excitation* developed for the identification of linear models [3].

\*This project is partially supported by EC funding 766840-COSY-BIO and a Royal Society of Edinburgh-MoST grant to F.M.; L.B. is supported by EPSRC funding EP/P017134/1-CONDSYC; EBC acknowledges funding from Spanish MINECO, grant ref. AGL2015-67504-C3-2-R.

<sup>1</sup>School of Engineering, Institute for Bioengineering, The University of Edinburgh, Edinburgh, EH9 3DW, UK. [lucia.bandiera@ed.ac.uk](mailto:lucia.bandiera@ed.ac.uk), [filippo.menolascina@ed.ac.uk](mailto:filippo.menolascina@ed.ac.uk)

<sup>2</sup>SynthSys - Centre for Synthetic and Systems Biology, The University of Edinburgh, Edinburgh, EH9 3BF.

<sup>3</sup>(Bio)Process Engineering Group, IIM-CSIC Spanish National Research Council, Vigo, Spain.

These facts raise the questions of whether it is possible to design informative experiments for the identification of gene networks and how to do so.

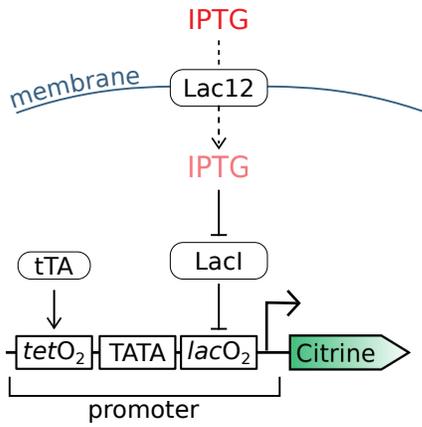
Model-based Optimal Experimental Design (MBOED) allows the design of maximally informative experiments and has recently been adopted in the identification of biological systems. For example, Bandara *et al.* showed that optimally designed (yet technologically constrained) experiments lead to a 60-fold reduction in the mean variance of parameter estimates over experience-based schemes [4]. Similarly, Ruess *et al.* emphasised the improvement of optimised dynamic inputs over random stimulation patterns in the characterisation of a light-inducible promoter [5].

The adoption of MBOED in (Synthetic) Biology generally faces a large amount of inertia: optimally designed experiments are difficult to implement with traditional experimental platforms and the skills to design them are not widespread in wet laboratories. Technological developments (e.g. microfluidics) and computational tools (e.g. AMIGO2 [6]) allow this limitation to be overcome but they have steep learning curves. The question is then: does the gain in information OED offers justify the efforts of adopting it?

To address this question, here we consider the identification of a mathematical model of a building block in Synthetic Biology: an inducible promoter. Many synthetic promoters are available, but since they generally use DNA sequences from the same organisms they are engineered for, they suffer from unwanted regulation from other genes in the genome. This makes disentangling and modelling promoter activity a non-trivial task. To overcome this issue, we focus on an orthogonal promoter [7], i.e. a promoter built in a species (*S. cerevisiae*) using DNA sequences from a different one (*E. coli*). This promoter, designed by Gnügge *et al.* [7] (Fig. 1), drives the expression of a fluorescent reporter, Citrine, when cells are exposed to the chemical IPTG. IPTG enters the cell through the permease Lac12 and binds the LacI protein, thereby relieving its repression on the promoter activity. Binding of the constitutively expressed tTA to the tetO<sub>2</sub> site results in expression of Citrine.

Based on published characterisation data [7], we first refine a mathematical model of the inducible promoter ( $M_{PLac}$ ), obtaining  $M_{PLac,r}$ . We then define a reduced model structure,  $M_{3D}$ , able to mimic the dynamics of  $M_{PLac,r}$  ( $M_{IP,r}$ ). We hence simulate the response of  $M_{IP,r}$  to both optimal and intuition-driven inputs and compare the amount of information provided by each input class using the posterior distributions of the inferred parameters.

Not only do our results suggest that MBOED allows the



**Fig. 1:** Schematic of the inducible promoter implemented in [7]. A native *S. cerevisiae* promoter was engineered by cloning the  $(tetO)_2$  and  $(lacO)_2$  operator sequences upstream and downstream of the TATA box respectively. The construct was integrated into the genome of a budding yeast strain constitutively expressing the heterologous transcription factors - tetracycline responsive transactivator (*tTA*) and *LacI* repressor, and the lactose permease (*Lac12*). The activity of the resulting promoter, regulated by the exogenous, non-metabolizable inducer  $\beta$ -D-1 thiogalactopyranoside (IPTG), is reported by the expression of the Citrine fluorescent reporter.

design of more informative experiments for the characterisation of synthetic promoters, they also provide a conservative estimate of the improvement in parameter accuracy that can be achieved via MBOED.

The manuscript is organised as follows: Section II discusses how we recalibrate a model of the inducible promoter, define a lower-order model and use it to compare the informativeness of different input classes. Section III elaborates on the importance of OED for the design of more informative experiments in Synthetic Biology. Section IV details our *in silico* experiments, the comparison of informativeness of different input classes and the design of optimal experiments. Finally, Section V presents our conclusions and future directions.

## II. RESULTS

### A. Refitting Gnügge *et al.*'s Model

As starting point of our analysis we consider  $M_{PLac}$ , the model proposed by Gnügge and colleagues [7]. We first seek to independently assess the ability of this model to capture the experimental data reported in the original paper [7], comprising several IPTG dose-response curves sampled at five equidistant time points after induction. We note that at intermediate concentrations model predictions appear to systematically underestimate the measured steady states by 20-30% (Fig. 2, grey line).

Reasoning that this discrepancy would offer an opportunity to refine  $M_{PLac}$ , we re-calibrate the model using enhanced Scatter Search (eSS) and obtain a new model,  $M_{PLac,r}$  (Fig. 2a, cyan), that generally better fits the available experimental data (Fig. 2b).  $M_{PLac,r}$  offers a 56% improvement in fit, as quantified by the sum of squared errors of predictions (SSE) (Fig. 2c).

### B. A reduced-order model captures the dynamics of the inducible promoter

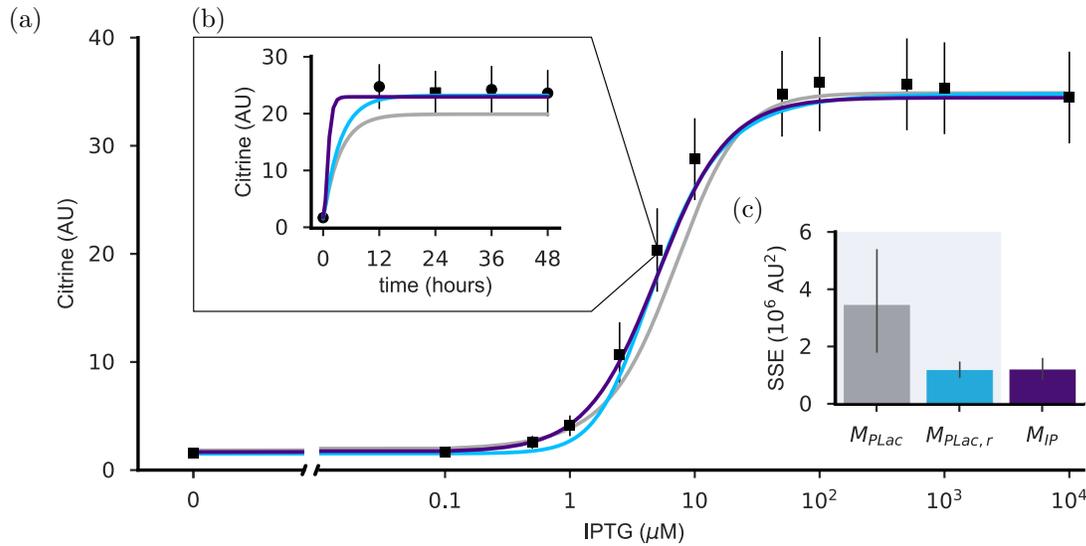
To constrain the number of parameters to be identified and the computational cost associated to optimal experimental design, we develop a lower-order model structure ( $\mathcal{M}_{3D}$ ). The model structure reads as follows:

$$\mathcal{M}_{3D} = \begin{cases} \frac{dR}{dt} &= \alpha + v \frac{IPTG^h}{K_r^h + IPTG^h} - \gamma R \\ \frac{dP_f}{dt} &= k_p R - (\gamma_f + k_f) P_f \\ \frac{dP_m}{dt} &= k_f P_f - \gamma_f P_m, \end{cases}$$

where  $R$ ,  $P_f$  and  $P_m$  are the concentrations of Citrine mRNA, immature folded protein and matured (fluorescent) protein, respectively. The model features 8 parameters:  $\alpha$  and  $v$  are the basal and maximal transcriptional rate respectively;  $h$ , the Hill coefficient;  $K_r$ , the Michaelis-Menten coefficient;  $k_p$ , the translation rate and the rate of maturation of the folded protein,  $k_f$ . All biochemical species are subject to linear degradation, occurring at rates  $\gamma$  for mRNA and  $\gamma_f$  for protein. This model structure builds on the assumption that the expression of *LacI* and *Lac12*, as well as the binding of *LacI*-dimer to the operator sites and to IPTG, occurs on faster time scales than Citrine expression. Fitting all parameters of  $\mathcal{M}_{3D}$  to the time-series data in [7], we obtain  $M_{IP}$ .

When compared with  $M_{PLac}$  and  $M_{PLac,r}$ , we find that  $M_{IP}$  best fits the measured steady-states, i.e. the dose-response curve (Fig. 2a), as well as the experimental data acquired at time-points different from 24 hours (see Fig. 2b for an example). Despite its lower order,  $M_{IP}$  achieves predictive capabilities comparable to  $M_{PLac,r}$ . To show this, we calculate the SSE over the whole set of experimental data (Fig. 2c). It is interesting to note that  $M_{IP}$  is characterised by a smaller rise time (1.8 hours) than both  $M_{PLac}$  and  $M_{PLac,r}$  (7.9 hours) (Fig. 2b). Here, the rise time is defined as the time required for the output to rise from 10% to 90% of the steady-state. We also note that the long sampling intervals used in the original study [7] does not allow further constraining the characteristic time-scale of the system.

As we aim to compare the informative content of different input classes (II, section C), we need our reduced model to mimic as closely as possible the dynamics of the genetic system of interest. We therefore consider  $M_{PLac,r}$  as our nominal model and generate a set of pseudo-experimental data to re-calibrate  $M_{IP}$ ; in so doing we obtain  $M_{IP,r}$ . It is worth noting that we could use  $M_{IP}$  for MBOED, however the ability to generate additional datasets and further constrain parameter calibration made us prefer  $M_{Lac,r}$  as our reference. Considering the limited complexity of the underlying biological system, we decided to use stepwise, pulses, ramp wise and stepwise random inputs in the pseudo-data generation. The results of  $M_{IP,r}$  identification show that this model recapitulates the dynamics of  $M_{PLac,r}$  (Fig. 3).



**Fig. 2:** Comparison between  $M_{PLac}$ ,  $M_{PLac,r}$  and  $M_{IP}$  model structures. (a) Dose-response curve after 24 hours of incubation with the specified IPTG concentrations. Experimental data, median (filled squares) and inter-quartile range (errorbars) of Citrine distributions, were retrieved from [7]. Solid lines show the in-silico dose-response curve for  $M_{PLac}$  (grey),  $M_{PLac,r}$  (cyan) and  $M_{IP}$  (purple). (b) The full data includes the dynamics of Citrine. An example is shown for the induction with 5  $\mu\text{M}$  of IPTG. (c) Barplot of the sum of squared errors of predictions (SSE), quantifying the predicted deviations from empirical data.

### C. Intuition-driven inputs are equally (poorly) informative

We first seek to compare the informative yield of intuition-driven stimuli (step, pulse and random) for the calibration of  $\mathcal{M}_{3D}$ . With this aim, we generated  $N_j = 100$  input profiles for each of the three classes (Methods, section IV-B). By simulating the output of  $M_{IP,r}$  for each input, we obtain pseudo-experimental data we use for the calibration of  $\mathcal{M}_{3D}$ . We formulate parameter estimation as a non-linear optimisation problem and use eSS to solve it. As the posterior distributions of parameter estimates are not Gaussian, we cannot use standard metrics (e.g. z-score) to assess the

statistical significance of the distance between nominal and estimated parameter value. To overcome this limitation, we compute the *relative error* ( $\varepsilon_i^{(j)}$ ) between each parameter estimate ( $p_i^{(j)}$ ) and its nominal value ( $p_i^*$ ):

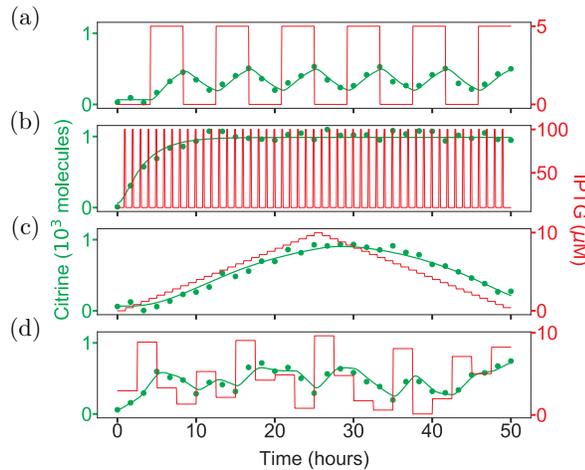
$$\varepsilon_i^{(j)} = \left| \log_2 \left( \frac{p_i^{(j)}}{p_i^*} \right) \right|. \quad (1)$$

where  $i$  identifies the  $i^{\text{th}}$  entry in the parameter vector and  $j$  is the index of the input profile yielding the parameter estimate  $p_i^{(j)}$ . Notably,  $\varepsilon_i^{(j)} = 0$  when the parameter estimate equals its nominal value, while the absolute value ensures that under and over estimates are treated equally.

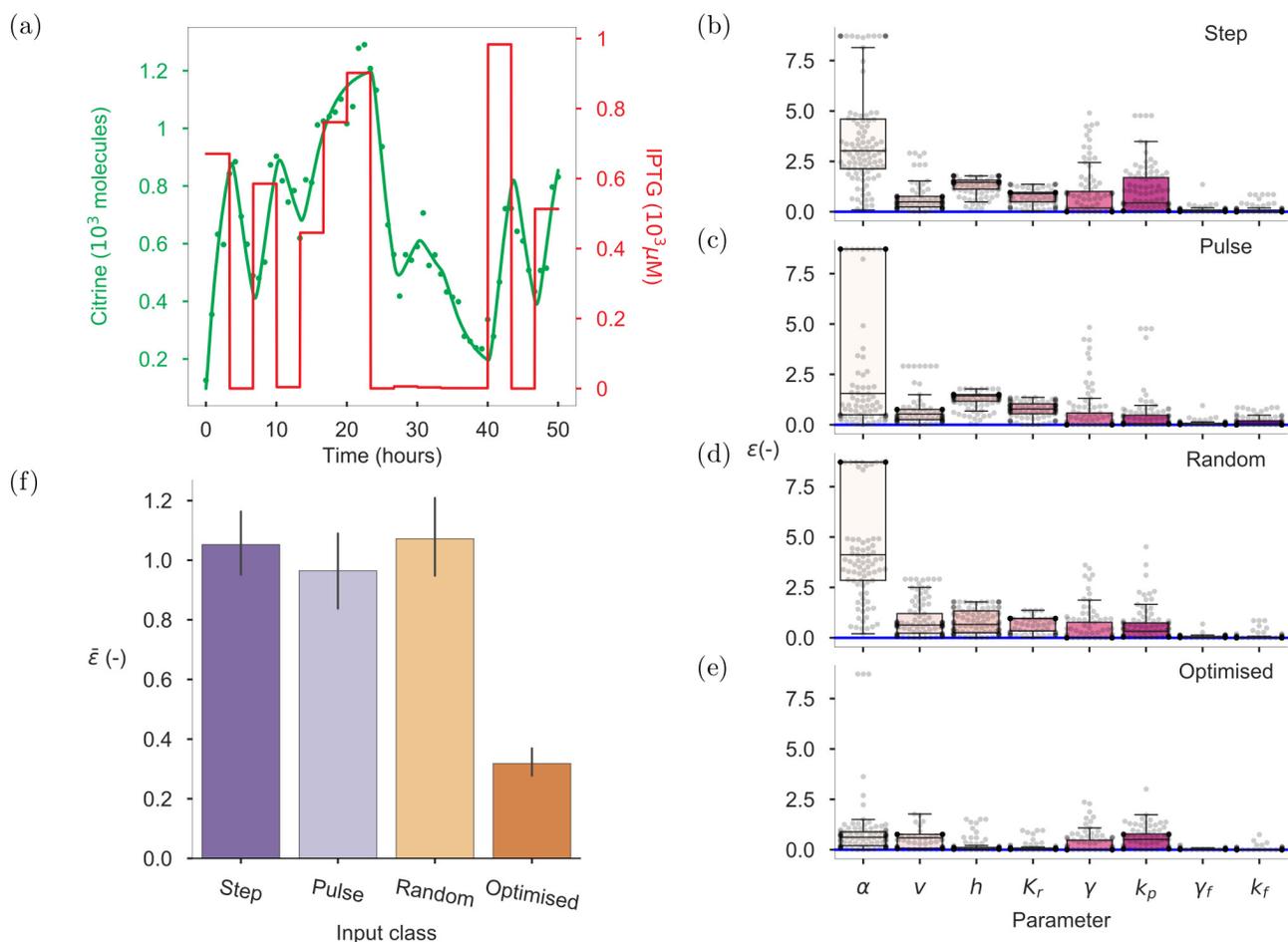
The distributions of relative error ( $\varepsilon_i$ ) for the 100 input profiles highlight a differential sensitivity of the output to the parameters (Fig. 4b-d). It is worth noting that the high variability in the estimates of  $\alpha$ ,  $v$  and  $\gamma$  agrees with a preliminary identifiability analysis (results not shown), suggesting high correlation between these parameters. Practical identifiability issues have indeed the potential to hinder our ability to identify the affected parameters with high confidence. Overall, the  $\varepsilon_i$  distributions suggest that the intuition-driven inputs convey a similar amount of information (Fig. 4b-d). This is further confirmed by the absence of a statistically significant difference in the *average relative error* ( $\bar{\varepsilon}$ ) metric (Fig. 4f), defined as:

$$\bar{\varepsilon} = \frac{1}{N_p N_j} \sum_{i=1}^{N_p} \sum_{j=1}^{N_j} \varepsilon_i^{(j)} \quad (2)$$

where  $N_p$  is the number of parameters in the model structure.



**Fig. 3:** Pseudo-experiments for the identification of  $M_{IP,r}$ . Step (a), pulse (b), ramp (c) and random (d) inputs (red line) were applied to  $M_{PLac,r}$  to simulate Citrine dynamics and to obtain pseudo-data (green circles). The response of the calibrated  $M_{IP,r}$  is shown as a green, solid line.



**Fig. 4:** Comparison of the informative content of different input classes for model identification. (a) Example of an optimally designed input (red line) applied to  $M_{IP,r}$  to simulate Citrine dynamics (green circles). The output of the system, upon calibration to the Citrine dynamic data, is shown as a green solid line. Box plots, overlaid with swarmplots, of the relative error ( $\epsilon$ ) of parameter estimates for step (b), pulse (c), random (d) and optimised (e) inputs. (f) Bar plot of the average relative error ( $\bar{\epsilon}$ ) provided by each input class.

#### D. Optimal Input Design (OID) enhances model calibration

We next test the improvement in the accuracy of parameter inference enabled by optimally designed experimental schemes. Having fixed the duration of the experiment, the sampling frequency and the switching times in the stepwise optimised input (Methods, section IV-D), we cast OID as a constrained optimisation problem that searches for the IPTG concentrations, (i.e. steps amplitude) that maximise the experimental information. We quantify information as the determinant of the Fisher Information Matrix ( $\mathcal{F}$ ) [3], [8]. This corresponds to the adoption of the highly popular D-optimality criterion [9]. To compare with the intuition-driven classes of input, we design  $N_j = 100$  optimised stimulation profiles (see Fig. 4a for an example), apply them to  $M_{IP,r}$  to obtain pseudo-data and solve the parameter estimation problem. The results show that the use of optimised inputs leads to a marked reduction in  $\epsilon$  when compared to experience-based stimulation patterns (Fig. 4a-e). The improvement in the accuracy of parameter estimates, noticeable for the poorly identifiable parameters  $\alpha$ ,  $v$  and  $\gamma$ , translates in a 69% reduction in  $\bar{\epsilon}$  for the optimally designed input over the

intuition-driven counterparts.

### III. DISCUSSIONS

Streamlining the inference of predictive mathematical models would foster their systematic use in Synthetic Biology. Here, by comparing the informative content of different input classes, we highlight optimal experimental design as a key strategy towards accurate and efficient model calibration. This conclusion was drawn considering the calibration of a deterministic model for the orthogonal, inducible promoter designed by Gnügge *et al.* [7]. We choose to focus on an inducible promoter for the key role these parts play in Synthetic Biology. Furthermore, it is commonly believed that the low complexity of synthetic promoters helps the experimentalist with the definition of informative experimental schemes based on intuition only. Our analysis clearly shows that this is not the case (Fig. 4f).

To compare the informativeness of the different input classes, we first retrieve the model structure ( $M_{PLac}$ ) proposed by the authors in [7]. The observed gap between the numerical and empirical transition-region of the dose-response curve en-

couraged us to attempt a model refinement by re-calibrating  $M_{PLac}$ . We frame model inference as a multi-experimental fitting problem and, unlike Gnügge and colleagues, address it using cross validation. While  $M_{PLac,r}$  yields a 56% improvement in the fitting over  $M_{PLac}$  (Fig. 2c), only speculations can be made on the cause of this difference.

Coherent with the *tenet* that the informative content of stimulation patterns depends on the (*a priori* unknown) dynamic properties of the system under investigation [5], we included step, pulse, ramp and random inputs in the pseudo-experiments. The match between model predictions from  $M_{IP,r}$  and pseudo-experimental data from  $M_{PLac,r}$  suggests that the dynamics of the latter can be fully recapitulated by the former (Fig. 3). This further supports using  $M_{IP,r}$  as a representative model of the true biological system when comparing the informativeness of different classes of inputs for model calibration.

We find that experiments with optimised inputs provide more accurate parameter estimates than intuition-driven inputs (Fig. 4f). However, it is important to note that the lower average error provided by the optimised input does not imply that all parameter estimates improve. This is evident in our results; for example, pulse inputs allow attaining a narrower  $\varepsilon$  distribution for  $k_p$  (Fig. 4b-d). Nevertheless, optimally designed inputs help tackling practically identifiability issues affecting some of the parameters (Fig. 4b-d).

We remark that the *a posteriori* analysis of the convergence curves of the input optimisation (results not shown) suggests that the  $\varepsilon$  we report should be considered an upper bound for the attainable improvement due to OED, rather than a precise estimate.

Taken together, these results suggest that a combination of *in silico* and experimental tools has the potential to significantly improve our ability to identify reliable and predictive models of biological systems and eventually enable the development of a Model-Based Biosystems Engineering framework in Synthetic Biology.

## IV. METHODS

### A. Generating Pseudo Experimental Data for the identification of $M_{IP,r}$

To re-calibrate parameter values in  $M_{IP}$ , and obtain  $M_{IP,r}$ , we choose to simulate the response of  $M_{PLac,r}$  to step, pulse, ramp and random inputs over 3000-minute long experiments. For each of these 4 input classes we define a generating function; we then design 3 inputs for each class. Step inputs are obtained using:

$$u_{\text{step}}(t) = \begin{cases} a, & \text{if } c \leq (t \bmod 2c) < 2c \\ b, & \text{if } 0 \leq (t \bmod 2c) < c \end{cases}$$

where  $a$ ,  $b$  and  $c$  are set to  $[5 \mu\text{M}, 0 \mu\text{M}, 250 \text{ min}]$  respectively for the first of the three time-profiles (Fig. 3A),  $[10 \mu\text{M}, 0 \mu\text{M}, 500 \text{ min}]$  for the second and  $[1000 \mu\text{M}, 10 \mu\text{M}, 500 \text{ min}]$  for the third.

To obtain pulse inputs we use the following definition:

$$u_{\text{pulse}}(t) = \begin{cases} a, & \text{if } 50 \text{ min} \leq (t \bmod 60 \text{ min}) < 60 \text{ min} \\ b, & \text{if } 0 \text{ min} \leq (t \bmod 60 \text{ min}) < 50 \text{ min} \end{cases}$$

where  $a$ ,  $b$  are set to  $[10 \mu\text{M}, 5 \mu\text{M}]$  for the first time-profile,  $[100 \mu\text{M}, 10 \mu\text{M}]$  for the second input (Fig. 3B) and  $[1000 \mu\text{M}, 600 \mu\text{M}]$  for the third.

As generating function of the ramp input we use:

$$u_{\text{ramp}}(t) = \begin{cases} \frac{a t}{1500}, & \text{if } 0 \text{ min} \leq t < 1500 \text{ min} \\ a - \frac{a t}{1500}, & \text{otherwise} \end{cases}$$

where  $a$  is set to  $10 \mu\text{M}$ ,  $100 \mu\text{M}$  (Fig. 3C) and  $1000 \mu\text{M}$  for each of the three inputs generated for this class. It should also be noted that a Zero Order Holder filter with a window of 60, 150 and 250 min was applied to the first, second and third input respectively.

Finally, the pseudo-random inputs are defined as:

$$u_{\text{random}}(t) = \begin{cases} a, & \text{if } 0 \text{ min} \leq (t \bmod c) < c \end{cases}$$

where  $a$ ,  $c$  are set to  $[\mathcal{U}(0 \mu\text{M}, 10 \mu\text{M}), 60 \text{ min}]$  for the first time-profile (Fig. 3D),  $[\mathcal{U}(0 \mu\text{M}, 90 \mu\text{M}), 150 \text{ min}]$  for the second and  $[\mathcal{U}(0 \mu\text{M}, 900 \mu\text{M}), 250 \text{ min}]$  for the third.

In all simulations, we add a 5% Gaussian noise and assign the initial conditions of the system to the steady state values derived from a 24 hour simulation of  $M_{PLac,r}$  with  $0 \mu\text{M}$  IPTG as the input. All experiments are simulated in AMIGO2 [6] and Citrine is sampled every 5 minutes. For more details on these procedures we refer the reader to our GitHub repository [10].

### B. Generating Pseudo Experimental Data for the comparison of input classes

The inputs we used to compare the informative content of different stimuli were defined as follows:

$$u_{\text{step}}(t) = \begin{cases} a, & \text{if } 0 \text{ min} \leq (t \bmod 200) < 100 \text{ min} \\ b, & \text{if } 100 \text{ min} \leq (t \bmod 200) < 200 \text{ min} \end{cases}$$

where, for each of the  $N_j$  inputs,  $a$  and  $b$  are two random values extracted from  $\mathcal{U}(0 \mu\text{M}, 1000 \mu\text{M})$ .

$$u_{\text{pulse}}(t) = \begin{cases} 0, & \text{if } 10 \text{ min} \leq (t \bmod 60 \text{ min}) < 60 \text{ min} \\ a, & \text{if } 0 \text{ min} \leq (t \bmod 60 \text{ min}) < 10 \text{ min} \end{cases}$$

where  $a$  is drawn from  $\mathcal{U}(0 \mu\text{M}, 1000 \mu\text{M})$ .

$$u_{\text{random}}(t) = \begin{cases} a, & \text{if } 0 \text{ min} \leq (t \bmod 80 \text{ min}) < 80 \text{ min} \end{cases}$$

where  $a$  is drawn from  $\mathcal{U}(0 \mu\text{M}, 1000 \mu\text{M})$ .

In all simulations, we add a 5% Gaussian noise and set the initial conditions of the system to the analytical steady-state of  $M_{IP,r}$  with IPTG equal to  $0 \mu\text{M}$ ; all experiments are simulated in AMIGO2 [6] and Citrine is sampled every 5 minutes. For more details on these procedures we refer the reader to our GitHub repository here [10].

### C. Parameter Estimation

Parameter estimation was formulated as a non-linear optimisation problem, whose objective is to identify the parameter values that minimise a scalar measure of the distance between model predictions and (pseudo) experimental data. We use the weighted least squares as a cost function, with weights set to the inverse of the experimental noise. To solve the optimisation problem, we rely on eSS [11]: a hybrid method that combines a global and a local search to speed up convergence to optimal solutions. In the initial phase, eSS explores the space of solutions, then, as local search, the algorithm employs the nonlinear least squares solver. To strengthen the predictive capabilities of the calibrated models, we use cross validation in the identification of  $M_{PLac,r}$  and  $M_{IP,r}$ . In both cases, the available experimental datasets are randomised and split into training (66%) and test (33%) sets. Parameter estimation is run on the training set starting from 100 initial guesses for the parameter vector. The latter are obtained as latin hypercube samples within the allowed boundaries for the parameters. Among the optimal solutions, the one that minimises the SSE on the test set is selected as the vector of parameter estimates. It is worth noting that, when comparing the informative content of different input classes, parameter estimation was not performed using cross validation. Details on the allowed bounds for the parameters and the scripts used for parameter estimation are provided in the GitHub repository [10].

### D. Optimal Experimental Design

To reflect wet-lab experimental constraints, we fix the sampling times (1 every 5 minutes) and the experiment duration (3000 minutes). We further set the initial condition to the steady-state in absence of induction. As a result, we restrict the optimisation to identifying the input (IPTG) time profile that maximises the information yield of the experiment. Here, information is quantified as a metric defined on the Fisher Information Matrix ( $\mathcal{F}$ ) [3], [8]:

$$\mathcal{F} = \sum_{i=1}^N \frac{1}{\sigma_i^2} [\nabla_{\theta} y]^T [\nabla_{\theta} y] \quad (3)$$

where  $y$  is the observable (Citrine) and  $\sigma_i^2$  represents the variance of the signal at the  $i^{\text{th}}$  sampling instant. The  $\mathcal{F}$  sets a lower bound on the variance of the parameter estimates through the Cramér-Rao inequality:

$$\mathcal{C} \geq \mathcal{F}^{-1} \quad (4)$$

where  $\mathcal{C}$  is the covariance matrix. Intuitively, as the eigenvalues of the  $\mathcal{F}$  are related to the inverse of parametric variances, attempting to maximise the determinant of  $\mathcal{F}$  (D-optimality) corresponds to minimising the product of the parametric variances.

In order to find the most informative input ( $u^*$ ), we formulate MBOED as an optimal control problem and search for:

$$u^* = \arg \max_u |\mathcal{F}(M_{IP,r}(p, u))| \quad (5)$$

where  $p$  is the parameter vector. We use Differential Evolution (DE) [12], a global optimisation method featuring good convergence properties and suitable for parallelisation, to solve the optimisation problem. We empirically [13] set the population size, crossover threshold and differential weight to 150, 0.3 and 0.5, respectively and adopt the strategy rand-to-best/1/exp.

## V. CONCLUSIONS AND FUTURE WORK

In this study we highlight MBOED as a key strategy for the accurate calibration of mathematical models of biological parts in Synthetic Biology. Our *in-silico* results suggest that optimally designed input profiles substantially improve the predictive ability of the inferred models, outperforming intuition-driven stimuli. While further studies are needed to explore the scalability of the computational cost for systems of higher complexity, we propose that combining flexible experimental platforms (e.g. microfluidics) and MBOED will enable the widespread adoption of mathematical models in Synthetic Biology. Beyond the required *in vivo* validation, our results encourage efforts towards the implementation of platforms to automate model calibration, in which MBOED and *in vivo* experiments are combined in an identification loop.

## ACKNOWLEDGMENTS

We would like to thank Prof. Jörg Stelling and Ms. Dharmarajan Lekshmi for clarifications on the analysis of experimental data and on the model implementation in [7]. We thank Mr. Alastair Hume for his contribution to the simulation code used in this study. We are further grateful to Prof. Diego di Bernardo and his lab for insightful discussions.

## REFERENCES

- [1] D. E. Cameron, C. J. Bashor, and J. J. Collins, "A brief history of synthetic biology," *Nature Reviews Microbiology*, vol. 12, no. 5, pp. 381–390, 2014.
- [2] T. Ellis, X. Wang, and J. J. Collins, "Diversity-based, model-guided construction of synthetic gene networks with predicted functions," *Nature biotechnology*, vol. 27, no. 5, pp. 465–471, 2009.
- [3] L. Ljung, "System identification: Theory for the user, ptr prentice hall information and system sciences series," ed: *Prentice Hall, New Jersey*, 1999.
- [4] S. Bandara, J. P. Schlöder, R. Eils, H. G. Bock, and T. Meyer, "Optimal experimental design for parameter estimation of a cell signaling model," *PLoS computational biology*, vol. 5, no. 11, p. e1000558, 2009.
- [5] J. Ruess, F. Parise, A. Miliadis-Argeitis, M. Khammash, and J. Lygeros, "Iterative experiment design guides the characterization of a light-inducible gene expression circuit," *Proceedings of the National Academy of Sciences*, vol. 112, no. 26, pp. 8148–8153, 2015.
- [6] E. Balsa-Canto, D. Henriques, A. Gábor, and J. R. Banga, "Amigo2, a toolbox for dynamic modeling, optimization and control in systems biology," *Bioinformatics*, vol. 32, no. 21, pp. 3357–3359, 2016.
- [7] R. Gnugge, L. Dharmarajan, M. Lang, and J. Stelling, "An orthogonal permease-inducer-repressor feedback loop shows bistability," *ACS synthetic biology*, vol. 5, no. 10, pp. 1098–1107, 2016.
- [8] E. Walter and L. Pronzato, *Identification of parametric models from experimental data*. Springer Verlag, 1997.
- [9] E. Balsa-Canto, A. A. Alonso, and J. R. Banga, "Computational procedures for optimal experimental design in biological systems," *IET systems biology*, vol. 2, no. 4, pp. 163–172, 2008.
- [10] <https://github.com/csynbiosysIBioEUoE>.

- [11] J. A. Egea, E. Balsa-Canto, M.-S. G. Garca, and J. R. Banga, "Dynamic optimization of nonlinear processes with an enhanced scatter search method," vol. 48, no. 9, pp. 4388–4401. [Online]. Available: <https://doi.org/10.1021/ie801717t>
- [12] R. Storn and K. Price, "Differential evolution—a simple and efficient heuristic for global optimization over continuous spaces," *Journal of global optimization*, vol. 11, no. 4, pp. 341–359, 1997.
- [13] D. Zaharie, "Critical values for the control parameters of differential evolution algorithms," in *Proceedings of MENDEL*, vol. 2, 2002, p. 6267.