

Bayesian Optimization for Design of Multiscale Biological Circuits

Published as part of the ACS Synthetic Biology virtual special issue "AI for Synthetic Biology".

Charlotte Merzbacher, Oisín Mac Aodha, and Diego A. Oyarzún*



Cite This: *ACS Synth. Biol.* 2023, 12, 2073–2082



Read Online

ACCESS |



Metrics & More



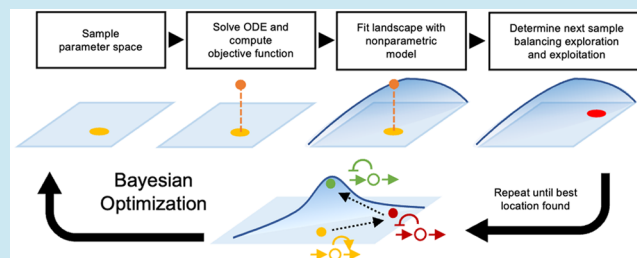
Article Recommendations



Supporting Information

ABSTRACT: Recent advances in synthetic biology have enabled the construction of molecular circuits that operate across multiple scales of cellular organization, such as gene regulation, signaling pathways, and cellular metabolism. Computational optimization can effectively aid the design process, but current methods are generally unsuited for systems with multiple temporal or concentration scales, as these are slow to simulate due to their numerical stiffness. Here, we present a machine learning method for the efficient optimization of biological circuits across scales. The method relies on Bayesian optimization, a technique commonly used to fine-tune deep neural networks, to learn the shape of a performance landscape and iteratively navigate the design space toward an optimal circuit. This strategy allows the joint optimization of both circuit architecture and parameters, and provides a feasible approach to solve a highly nonconvex optimization problem in a mixed-integer input space. We illustrate the applicability of the method on several gene circuits for controlling biosynthetic pathways with strong nonlinearities, multiple interacting scales, and using various performance objectives. The method efficiently handles large multiscale problems and enables parametric sweeps to assess circuit robustness to perturbations, serving as an efficient *in silico* screening method prior to experimental implementation.

KEYWORDS: Bayesian optimization, machine learning, dynamic pathway control, genetic circuit design, multiscale biological systems, metabolic engineering



1. INTRODUCTION

The design of molecular circuits with prescribed functions is a core task in synthetic biology.¹ These circuits can include components that operate across various scales of cellular organization, such as gene expression, signaling pathways,² or metabolic processes.³ Computational methods are widely employed to discover circuits with specific dynamics,^{4–6} and, in particular, optimization-based strategies can be employed to search over design space and single out circuits predicted to fulfill a desired function.^{7–10} However, circuit design requires the specification of circuit architecture, i.e., the circuit “wiring diagram”, as well as the strength of interactions among molecular components. Since circuit architectures are discrete choices and molecular interactions depend on continuous parameters such as binding rate constants, circuit design leads to mixed-integer optimization problems that can be notoriously difficult to solve.¹¹ Moreover, when circuits operate across multiple scales, their computational models become numerically stiff,¹² resulting in extremely slow simulations that make their mixed-integer optimization challenging or even impossible to solve.

Previous work on computational circuit design has largely focused on genetic circuits that operate in isolation from other layers of the cellular machinery (Figure 1A). A range of

techniques have been employed to identify functional circuits, including exhaustive search,^{4–6,13} computational optimization,^{7,8} systems theoretic approaches,^{14–18} Bayesian design,^{19,20} and machine learning.^{9,21} While these methods differ on their specific modeling strategies and assumptions, they all require computational simulations at many locations, typically thousands to millions, in the design space. But since multiscale systems often cannot be simulated at such scale, the computational costs limit the applicability of current optimization methods.

A notable example of this challenge appears in genetic circuits for dynamic control of metabolic pathways.^{22–26} These systems are receiving substantial attention thanks to several successful implementations that improved yields as compared to classic techniques in metabolic engineering.^{27,28} The key principle is to put enzymatic genes under the control of metabolite-responsive mechanisms that couple heterologous

Received: February 23, 2023

Published: June 20, 2023



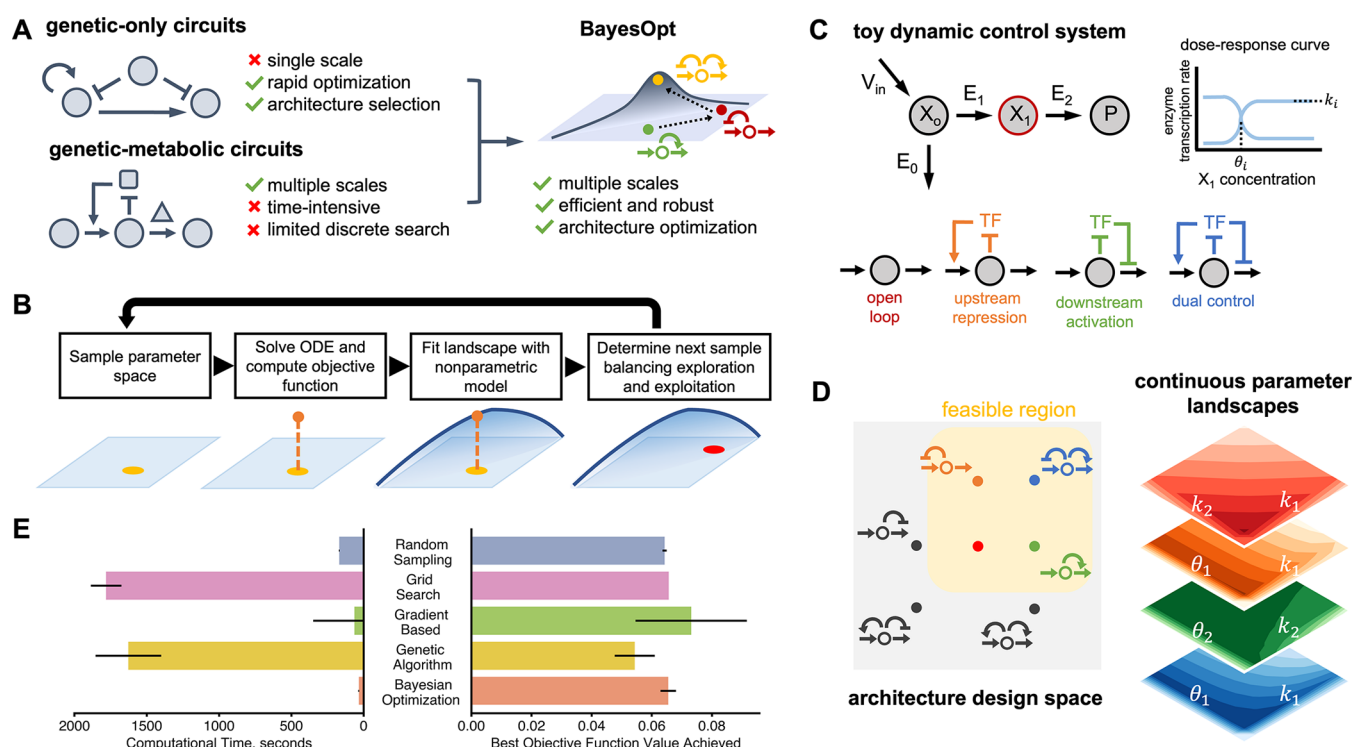


Figure 1. Bayesian optimization for the design of circuit architectures and parameters. (A) Previous optimization methods have focused on genetic circuits in isolation from other cellular processes. For multiscale circuits, optimization approaches become infeasible due to the difficulty of simulating stiff dynamical systems in many locations of the design space; a common example of such multiscale systems are gene circuits that control metabolic production.³ We propose the use of Bayesian optimization (BayesOpt) for efficient optimization of architectures and parameters in multiscale circuits. (B) Schematic of a mixed-integer Bayesian optimization loop; the objective function is regarded as a random variable to be optimized over an input space comprised of continuous parameters and a set of discrete circuit architectures. At each iteration, the algorithm computes the value of the objective function from the solution of an ordinary differential equation (ODE) model at a single location in the input space. The algorithm learns the shape of the objective landscape using a nonparametric statistical model,³⁷ which is employed to propose a new location in the input space through an acquisition function designed to balance exploration and exploitation of the input space; more details in the [Methods](#). The algorithm iteratively learns the shape of the performance landscape until convergence to a global optimum. (C) Example metabolic pathway under gene regulation. We consider three negative feedback architectures plus open loop control; the architectures are named based on the net effect of the metabolite on gene expression. The intermediate X_1 binds a transcription factor (TF) that controls the expression of pathway enzymes, either as an activator or repressor. V_{in} is the constant influx to the engineered pathway from native metabolism. The TF dose-response curve (at right) is described by three parameters, k_i , θ_i , and n_i where $i = 1, 2$. The aim is to find designs with optimal architecture and dose-response parameters (k_i , θ_i); for simplicity the Hill coefficient was fixed to $n_i = 2$. (D) Performance landscapes of the four feasible circuit architectures. We exclude architectures with positive feedback loops as these are prone to multistability.⁴⁷ The shape of the performance landscape defined in [eq 3](#) shows substantial variation across the four architectures. This leads to a highly nonconvex mixed-integer optimization problem. Heatmaps show the value of the objective J computed on a regular grid of the indicated parameters. (E) Comparison of BayesOpt against other strategies using the toy model as a benchmark; lower objective function values are better. Shown are the results for random sampling ($N = 1,000$ samples), grid search ($N = 40,000$), a genetic algorithm⁵⁵ ($N = 100$ individuals, $N = 1,000$ generations), and a gradient-based optimizer to find optimal continuous parameter values for each architecture.⁴⁸

expression to the concentration of a pathway intermediate.³ This creates feedback loops between enzyme expression and pathway intermediates that allow the control of pathway activity in response to upstream changes in growth conditions or precursor availability. Such dual genetic-metabolic systems are particularly challenging to simulate efficiently because metabolites and enzymes vary in different time scales, from milliseconds (enzyme kinetics) to minutes (enzyme expression), and they also appear in vastly different concentrations; in bacteria enzymes are typically expressed in nanomolar concentrations, while metabolites are found typically above the millimolar range.²⁹ Moreover, the implementation of these systems is costly and requires substantial experimental fine-tuning. As a result, a central task prior to implementation is the choice of a suitable feedback control loop between metabolites and enzymatic genes, and the strength of interactions between metabolites and actuators of gene expression such as

transcription factors³⁰ or riboregulators.³¹ The design of control architectures is particularly important, because there are many ways of building similar control loops,³² for example by employing combinations of transcriptional activators and repressors,^{33,34} that may differ in their performance and cost of implementation.

Here, we present a fast and scalable machine learning approach for optimization of multiscale circuit architectures and parameters ([Figure 1A](#)). The method is based on Bayesian optimization coupled with differential equation models, and we highlight its utility in various models of metabolic pathways under genetic feedback control.³⁵ Using a toy example for a simple pathway, we first show that the method converges rapidly and outperforms other optimizers by a substantial margin. We then consider real world models of metabolic pathways in *Escherichia coli* for the production of several relevant precursors: glucaric acid,³⁶ fatty acids,³³ and *p*-

aminostyrene.³⁴ We use these pathways to illustrate how the speed of our method enables screening optimal designs in realistic design tasks that would otherwise be infeasible to compute, including the impact of uncertain enzyme kinetic parameters, the use of layered architectures that combine metabolic and genetic control, and the optimization of a complex model with 23 differential equations, 27 candidate control architectures, and 16 parameters to be optimized. The method can help speed up the design of synthetic biological circuits and presents a novel approach to explore the design space ahead of implementation.

2. RESULTS

2.1. Bayesian Optimization for Joint Optimization of Circuit Architecture and Parameters. In general, a circuit design task can be stated as the following mixed-integer optimization problem:

$$\begin{aligned} \min_{p_d, p_c} J(x, p_c, p_d) \\ \text{subject to:} \\ dx/dt = h(x) \\ p_c \in C, \quad p_d \in D \end{aligned} \quad (1)$$

where $J(x, p_c, p_d)$ is a performance objective to be optimized over a space of continuous parameters p_c and a discrete set of circuit architectures p_d . The ordinary differential equation (ODE) in eq 1 describes the temporal dynamics of circuit components and are typically built from mass balance relations comprised in the nonlinear function $h(x)$. Common examples of continuous parameters in applications are binding affinities between DNA and regulatory proteins, or the strength of protein-protein interactions. Conversely, circuit architecture would typically involve various combinations of positive and negative feedback loops among molecular species. We have stated the problem as minimization of J , but similar formulations can be posed as a maximization problem.

In this paper, we propose to solve the design problem in eq 1 with Bayesian optimization (BayesOpt), a class of algorithms designed for problems with objective functions that are expensive to compute. BayesOpt is a global optimization technique that treats the objective function as a random variable with a prior distribution on it. The algorithm creates a statistical model of the objective through subsequent evaluations, which are employed to build a posterior distribution and determine the next set of inputs to evaluate (Figure 1B). A typical application of BayesOpt is in design of experiments³⁵ where the objective function requires measuring data with costly and/or slow experimental work. In deep learning, BayesOpt is widely employed for model selection, as traditional grid search approaches require large compute resources to train many architectures with combinations of various layer sizes and other hyperparameters.^{37,38}

For the circuit design task in eq 1, when the biological system has multiple scales the computation of objective J requires solving a stiff ODE at many locations of the mixed-integer search space, which can rapidly become infeasible. To illustrate the utility of BayesOpt in a range of design problems, we focus on genetic control circuits for metabolic pathways that synthesize high-value products. In this case, the ODE in eq 1 contains two sets of equations:

$$\begin{aligned} ds/dt &= f(s, e) - \lambda s \\ de/dt &= u(s, p_c, p_d) - \lambda e \end{aligned} \quad (2)$$

where s and e are vectors of metabolite and enzyme concentrations, respectively. The term $f(s, e)$ describes the biochemical reactions between pathway intermediates, while the parameter λ models the dilution effect by cell growth. The vector $u(s, p_c, p_d)$ describes the enzyme expression rates controlled by some pathway intermediates, and typically take the form of sigmoidal dose-response curves that lump together processes such as metabolite-TF or metabolite-riboregulator interactions.³⁰ The continuous parameters p_c model the dose-response curves of the feedback mechanisms, whereas the discrete parameters p_d specify the gene control architecture. The number of heterologous enzymes determines the number of genes in the control circuit. In pathways under dynamic control as in eq 2, both sets of species change in different time scales; metabolic reactions operate in the millisecond range or faster,³⁹ while enzyme expression changes in the scale of minutes or longer. Moreover, metabolites and enzymes are also present in different ranges of concentrations, from nM for enzymes to mM and higher for metabolites.²⁹ As a result, simulation of the ODE in eq 2 is computationally expensive, particularly when this has to be done many times as part of an optimization-based search.

The performance objective J can be flexibly used to model common design goals such as production flux, yield or titer, as well as cost-benefit tasks that balance production with the deleterious impact of the pathway on the physiology of the host. To first establish a baseline for the performance of our method, we employed a simple toy pathway model that displays common features found in real metabolic pathway (Figure 1C). The model includes a metabolic branch point through a heterologous pathway with two enzymatic steps. As a performance objective we considered the minimization of

$$J = \alpha_1 J_{\text{prod}} + \alpha_2 J_{\text{cost}} \quad (3)$$

where J_{prod} was designed so that its minimization is equivalent to maximization of the production flux, and J_{cost} penalizes total amount of enzyme expressed during the culture. The parameters α_1 and α_2 are positive weights used to control the balance between the costs and benefits of expressing the heterologous pathway. Details on the objective function can be found in the Supporting Information.

We considered the four control architectures shown in Figure 1C, which include open loop control as well as three different implementations of negative feedback control using a metabolite-responsive transcription factor. Negative feedback is widely employed in gene circuits as it has substantial benefits in terms of robustness and performance, and their properties have been extensively studied in the literature.^{40–42} To illustrate the challenge of jointly optimizing circuit architecture and parameters, in Figure 1D we show a schematic of the design space. The four control architectures under consideration reside at different discrete points in the architecture space. Within each architecture, we observe substantial variations in the shape of the performance landscape J as a function of the dose-response parameters p_c . There are cases with convex landscapes with a clear optimum (e.g. dual control) and landscapes with flat basins where most optimization algorithms would struggle to find the optimum (e.g. downstream activation). When searching over the space

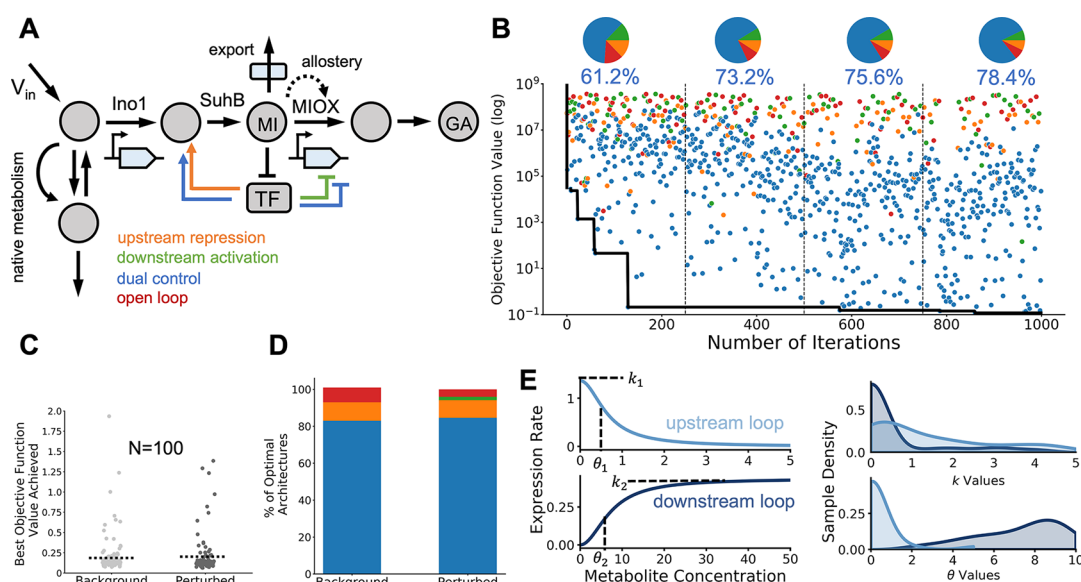


Figure 2. Robustness of optimal circuits to parameter uncertainty. (A) Schematic of a dynamic pathway for production of glucaric acid in *Escherichia coli*.²⁶ The pathway includes allosteric inhibition and export of an intermediate to the extracellular space. The core pathway components myoinositol (MI) and glucaric acid (GA) are modeled explicitly, as are the enzymes Ino1 and MIOX. The enzyme SuhB is not rate-limiting and is not modeled explicitly. V_{in} is the constant influx to the engineered pathway from native metabolism. As in Figure 1C, the architectures are named based on the net effect of the metabolite on gene expression. (B) Sample run of the BayesOpt algorithm for 1,000 iterations of the loop in Figure 1B. Black line shows the descent on the value of the objective function. Dots show all samples explored by the algorithm; pie charts show the fraction of architectures explored by the algorithm, and the fraction of samples taken from the majority architecture (dual control). The first quarter of the run had the most exploration of architectures other than dual control, with 38.6% of samples coming from nonmajority architectures. This percentage steadily decreased over the iterations but did not drop below 20%, illustrating the global nature of the optimization routine. (C) To examine the robustness of the optimal solutions to parameter uncertainty, we computed optimal solutions for many perturbed parameters of the allosteric activation of MIOX by its substrate myoinositol (MI). Strip plot shows the best objective function values achieved for background and perturbed kinetic parameters ($V_{m,MIOX}$, a_{MIOX} , $k_{d,MIOX}$) in eq 4. Kinetic parameters were perturbed using Latin Hypercube sampling⁵⁶ in the range (−100%, +100%) of the nominal values (Supporting Information). We observed little difference between background and perturbed values; dashed line denotes the mean value of the objective function. Only one of the $N = 100$ runs for perturbed parameters failed to converge the optimum. (D) Optimal architectures across runs with background and perturbed parameter values. Both background and perturbed systems resulted in over 80% of runs selecting dual control as the optimal architecture. (E) Average dose-response curves and distribution of optimal parameters for the dual control architecture with perturbed allosteric parameters. The repressive and activatory loops have substantially different dose-response curves on average. The distributions of the dose-response parameters (right) show important variations in their mean and dispersion. The parameter k_i and θ_i determine the maximal enzyme expression rate and regulatory threshold, respectively.

of architectures and parameters simultaneously, the problem becomes a mixed-integer, nonconvex optimization that is extremely challenging to solve with traditional approaches.

We implemented a BayesOpt routine to jointly compute the architecture (p_d) and dose-response parameters (p_c) that minimize the performance objective in eq 3. We benchmarked its performance against several other methods, including a random search, an exhaustive grid search, a gradient based method, and a genetic algorithm (Figure 1E). The algorithm was able to compute optimal solutions rapidly (average 27 seconds per run across 100 runs) and robustly (standard deviation less than 2.5% of the mean optimal objective function value). BayesOpt runs significantly faster than the other methods, and provides a 30-fold improvement over a genetic algorithm. The accuracy of the optimum, quantified by the minimal value of the objective function, is on average 11.4% worse than the genetic algorithm, but this falls within the variation of the latter across several runs. We also note that the traditional gradient-based optimizer proved unreliable and failed to converge on 14.5% of runs.

A key advantage of Bayesian methods is that they are not gradient-based, and therefore are not constrained to navigate the space smoothly in the direction of steepest descent. Gradient-based methods can get trapped in local minima and

struggle to find the global optimum, especially in highly nonconvex landscapes like the ones presented here. In contrast, BayesOpt does not converge by chasing minima directly but rather by modeling the entire objective function landscape, which results in rapid and reliable results. The method can perform multiple “jumps” between distant locations in the discrete-continuous search space, where each subsequent sample is selected to maximize the expected improvement on the best sample found so far.

The speed of our approach enables the computation of large solution ensembles under model perturbations such as sweeps of key model parameters. In addition, our method can search high-dimensional mixed-integer design spaces. We next illustrate the versatility of the approach in a range of relevant real world pathways that require solving the optimization problem for large samples of parameter values.

2.2. Robustness of Control Circuits to Uncertainty in Enzyme Kinetic Parameters. A challenge in building pathway models is the substantial uncertainty on the enzyme kinetic parameters; this is particularly critical for pathways that include regulatory mechanisms such as allostery or product inhibition, which are often poorly characterized. Databases such as BRENDA⁴³ often have insufficient data on enzyme kinetics for a particular host strain or substrate of interest.

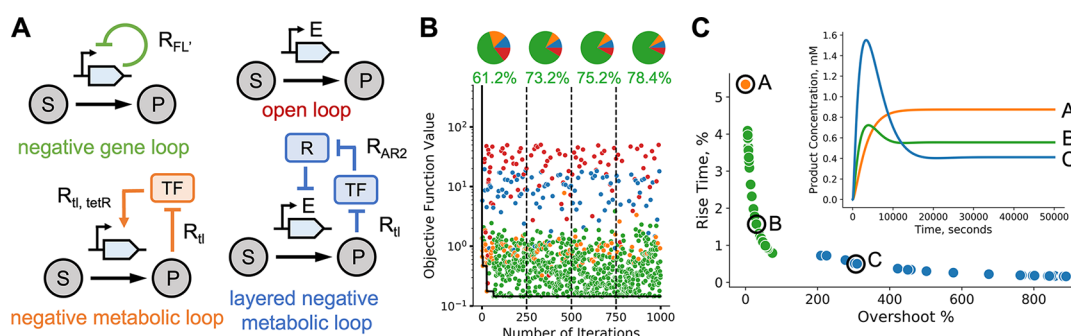


Figure 3. Optimization of metabolite dynamics in a fatty acid synthesis pathway. (A) Pathway diagrams with various control architectures implemented in *Escherichia coli*.³³ The metabolic loop employs a metabolite-responsive transcription factor, whereas the gene loop includes only a repressor expressed on the same promoter as the enzyme. (B) Representative run of BayesOpt with cost-benefit objective showing the best objective function value (black line). All samples are colored by their architecture. Pie charts of each quarter of the run show continued exploration of all architectures despite clear stratification in losses. (C) Optimal trade-off curve between overshoot and rise time. The objective weight α was swept from $\alpha = 0.01$ to $\alpha = 10,000$ and BayesOpt was run for 100 iterations at each α value. The optimal parameter values were used to compute the rise time and overshoot for visualization. The inset shows three sample trajectories illustrating how different optimal architectures navigate the trade-off between overshoot and rise time.

Since pathway dynamics can strongly depend on enzyme kinetics, the parametric uncertainty requires extensive sweeps of kinetic parameters to determine the robustness of a specific control architecture deemed to be optimal.

We focused on a pathway for synthesis of glucaric acid in *E. coli* (Figure 2A), a key precursor for many downstream products.³⁶ The pathway branches from glucose-6-phosphate (g6p) in upper glycolysis and contains three enzymatic steps (Ino1, SuhB, and MIOX). Doong and colleagues implemented a dynamic control circuit using the dual transcriptional regulator IpsA which responds to the intermediate myoinositol (MI).²⁶ The pathway enzyme MIOX is allosterically activated by its own precursor, and one intermediate (MI) can be exported to the extracellular space. We employed a previously developed ODE model¹⁰ that was parametrized using a combination of enzyme kinetic data and omics measurements, and considered the same four control architectures as in the previous example, including various alternative implementations of negative feedback control.

The results in Figure 2B show a typical run of the optimizer when using the cost-benefit objective in eq 3 (details in Supporting Information), together with the fraction of samples in which the algorithm explored each control architecture across the successive iterations. The optimal architecture (dual control in this case) was found quickly and the algorithm was able to further decrease the value of the objective function by exploring the space of dose-response parameters of IpsA. We observe that as the iterations progress, the algorithm shows a remarkable ability to explore other architectures despite their larger objective function values, thus highlighting the global nature of the algorithm.

To explore the impact of uncertain enzyme kinetics, we perturbed the parameters of the rate-limiting MIOX allosteric reaction:

$$V_{\text{MIOX}} = \frac{V_{\text{m,eff}} \text{MI}}{k_{\text{m,MIOX}} + \text{MI}}$$

$$\text{given } V_{\text{m,eff}} = V_{\text{m,MIOX}} \frac{1 + a_{\text{MIOX}} \text{MI}}{k_{\text{a,MIOX}} + \text{MI}} \quad (4)$$

where $V_{\text{m,MIOX}}$ is the maximum rate of reaction, $k_{\text{m,MIOX}}$ is the Michaelis-Menten constant, and $k_{\text{a,MIOX}}$ and a_{MIOX} are

allosteric activation constants. We solved the optimization problem for 1,000 combinations of these three parameters, which took under 16 hours on a Macbook Air with Apple M1 processor and 8 GB of RAM running MacOS Monterey. Perturbing the kinetic parameters of the glucaric acid pathway did not significantly affect the minimum objective function value achieved, indicating that the optimum is robust to uncertainty in the kinetic parameters (Figure 2C). However, the mean optimal objective function value was not significantly higher among the perturbed samples. We found that the dual control architecture was chosen as optimal in more than 85% of samples (Figure 2D). We thus sought to examine the optimal dose-response parameters of this architecture in more detail.

The maximal enzyme expression rates (k) and regulatory thresholds (θ) control the shape of the dose-response curves. As shown in Figure 2E, we found that the upstream repressive loop and downstream activatory loop had different optimal dose-response curves, corresponding to different optimal values of the continuous parameters. Optimal values of the upstream repression threshold θ_1 are low (mean value 0.64) and compressed into a narrow range as compared to the larger standard deviation of the downstream repression threshold θ_2 (mean value 7.24). This is reflected on a larger variation in the shape of the dose response curve for the downstream loop. Experimental fine-tuning of a dual control circuit might target parameters with optimal values with a wide range, such as k_1 , as varying these parameters is less likely to impair circuit function. Overall, these results show the robustness of the glucaric acid dual control system to kinetic parameter uncertainty and demonstrate the possibilities enabled by the speed of BayesOpt.

2.3. Exploration of Alternative Objective Functions.

In the previous case studies we employed a cost-benefit objective designed to account for the trade-off between heterologous production and the cost of expressing pathway enzymes, as in eq 3. To demonstrate the flexibility of the method with other objective functions, here we consider the optimization of the temporal trajectories of pathway metabolites.

We focused on the joint optimization of the rise time and overshoot in a model of a fatty acid production pathway considered previously in the literature.³³ Fatty acids are an

essential energy source and cellular membrane component. In addition, hydrocarbons derived from fatty acids have attracted attention as a potential biofuel source.^{24,44} Recent work engineering metabolic and genetic control loops showed that negative feedback control could speed up the rise to steady-state conditions.³³ The pathway built in literature expressed a thioesterase under transcriptional control, shown as the negative metabolic loop (NML) architecture in Figure 3A. In addition to transcription-factor mediated negative feedback loops, this model also includes individually implemented direct genetic loops where a repressor is expressed on the same promoter as the enzyme. These two different scales of loops interface with different levels of cellular organization. We explore several control architectures previously proposed in the literature³³ (Figure 3A).

We first considered a similar objective function as in eq 3 so as to compare convergence against the previous case studies. We implemented a modified production flux objective which takes the reciprocal of the product flux to convert the optimization to a minimization problem. The pathway cost J_{cost} is measured by summing the expression of all heterologous enzymes and varies across the different architectures. Details on the objective function can be found in the Supporting Information. A representative optimization run for this objective (Figure 3B) shows that the negative gene loop (NGL, green) and negative metabolic loop (NML, orange) architectures perform, on average, better than the other two architectures. BayesOpt samples taken from the open-loop architecture were, on average, 2 orders of magnitude worse than samples taken from NML and NGL architectures. Despite such hierarchy of loss values across the four architectures, the method effectively explores all architectures throughout the optimization run.

We next considered the optimization of percent overshoot and rise time presented in the literature.³³ The percent overshoot objective, J_{os} , measures the maximal deviation of product from its steady state concentration and is defined as the percent difference between the maximum fatty acid concentration and the steady state fatty acid concentration. The rise time, J_{rt} , is a measure of how fast fatty acid production rises to steady state and is defined as the first time point where the fatty acid concentration reaches 50% of the steady state value, normalized by the total integration time. We minimized the sum of the overshoot and rise time with a scaling weight α :

$$J = \alpha J_{\text{rt}} + J_{\text{os}} \quad (5)$$

Adjusting α balances the relative importance of the two optimization criteria; details on the calculation of rise time and overshoot are in the Supporting Information. Higher values of α correspond to optimal circuits with low rise times, while lower values of α prioritize circuits with low overshoot. Rise time is a measure of circuit speed, while overshoot is a measure of circuit accuracy. We found that when α is varied across several orders of magnitude, the optimal circuits form a optimal trade-off curve (Figure 3C). Different architectures occupy different parts of the optimal trade-off curve and display markedly different dynamics. The NML optima occupies a single point in the loss space, indicating that multiple continuous parameter values give the same loss function value for multiple values of α . The NML also has the lowest absolute loss function value of all the architectures considered. The NGL and layered negative metabolic loop (LNML) architectures occupy larger ranges on the curve, with

NGL giving a low overshoot and LNML a low rise time. The optimal NML circuit has no overshoot but the slowest rise time, while the LNML has a rapid rise time but overshoots the steady-state value by more than double. These opposing trade-offs demonstrate the importance of balancing multiple circuit design objectives.

2.4. Scalability to Large Pathway Models. Our previous case studies have been limited to circuits with a single metabolite controlling gene expression and a relatively small number of control architectures. We now study a large model for the synthesis of *p*-aminostyrene (p-AS), an industrially relevant vinyl aromatic monomer, in *E. coli* (Figure 4A)⁴⁵ using a cost-benefit objective similar to eq 3 tailored to the specific pathway (details in Supporting Information). This model has two possible metabolites that can regulate gene expression, namely *p*-aminocinnamic acid (p-ACA) and *p*-aminophenylalanine (p-AF), both of which can act as ligands for aptazyme-regulated expression device (aRED) transcription factors,⁴⁶ and three genes to be controlled. The aRED transcription factors can also act as dual regulators (activators or repressors) on any of the three promoters involved in the pathway. For simplicity, we limit the design space to control architectures without positive feedback loops, as these are prone to bistability.⁴⁷ This results in 27 possible control architectures and 16 continuous parameters to be optimized. The model also has a number of additional complexities. It contains operon-based gene expression commonly found in bacterial systems (genes *papA*, *papB*, and *papC* are expressed on the *papABC* operon), it includes a detailed description of mRNA dynamics and protein folding, which results in a large model with 23 differential equations, and it can also display oscillatory dynamics.

In addition to expression of heterologous enzymes, the accumulation of toxic intermediates is another major source of genetic burden to host organisms. The p-AS model has several sources of toxicity present in the pathway.^{34,45} The intermediate p-ACA and the efflux pump used to remove p-ACA from cells are both cytotoxic, while another intermediate, p-AF, leaks from cells.³⁴ The pathway enzyme L-Amino Acid Oxidase (LAAO) depletes key aromatic amino acid metabolites and creates toxic hydrogen peroxide as a byproduct. The model incorporates these various types of toxicity in the form of a toxicity factor τ . This toxicity factor is of the form

$$\tau = \frac{k_i}{k_i + \frac{p_{\text{ACA}}}{t_a} + \frac{p_{\text{efflux}}}{t_p} + \frac{\text{LAAO}}{t_l}} \quad (6)$$

where t_b , t_a , and t_p are chemical-specific toxicity factors. Enzyme-induced toxicity t_l scales the key metabolite depletion rate driven by the enzyme LAAO. Metabolite-induced toxicity t_a scales the impact of toxic intermediate p-ACA concentration. Finally, protein-induced toxicity t_p reflects the toxicity caused by efflux pump expression. The toxicity factor acts as a scaling coefficient on the pathway synthesis, degradation, and folding reaction rates.

Despite the complexity and size of the p-AS model, we observe that BayesOpt explores many of the 27 possible architectures and converges to a low value of the objective function (Figure 4B); this was also achieved at a reasonable computational cost (mean run time under 2 min). The best architecture selected in the sample run was a double upstream repression, single downstream activation loop controlled by p-AF (Figure 4B, inset), but there is no clear best architecture

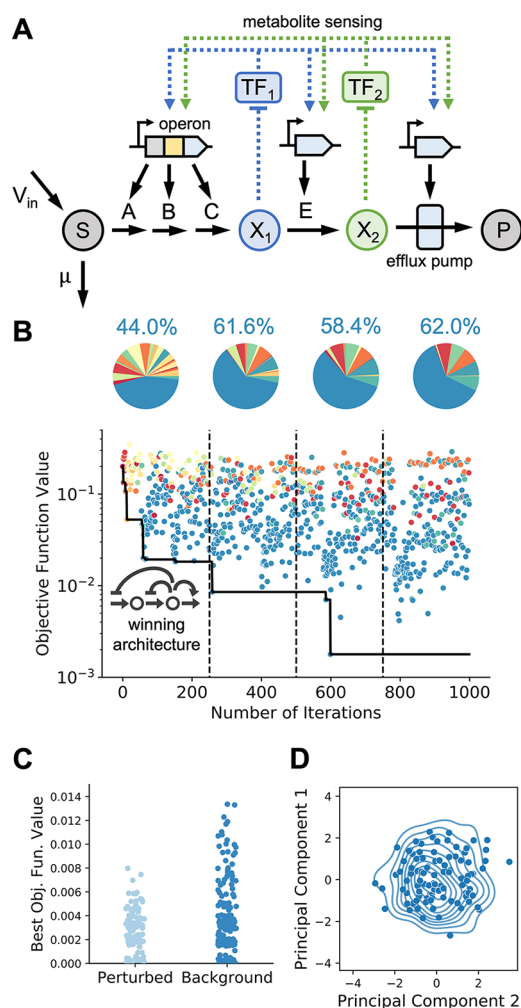


Figure 4. Bayesian optimization in a complex pathway. (A) Schematic of pathway for production of *p*-aminostyrene.³⁴ Two intermediates can act as ligands for metabolite-dependent riboregulators, and three promoter sites of control. The optimization problem has 16 continuous decision variables and 27 circuit architectures. The substrate *S* is converted by enzymes *A*, *B*, and *C* to X_1 , which is then converted by *E* to X_2 . The toxic substrate X_2 is then pumped out of the cell via an efflux pump to form the product *P*. Both X_1 and X_2 can act on the transcription factors TF_1 and TF_2 . V_{in} is the constant influx to the engineered pathway from native metabolism. (B) Representative run of the BayesOpt algorithm; the method samples many architectures before settling on the optimal one. Pie charts show continued exploration of a large number of architectures. The winning architecture is shown in the inset. (C) The *p*-aminostyrene pathway has several forms of substrate, protein, and enzyme toxicity expressed via a toxicity factor τ (see eq 6). To explore the effects of protein and metabolite toxicity, we perturbed the toxicity factor. Metabolite-induced toxicity was perturbed on the nominal range (1×10^{-3} , 1×10^{-4}) and protein-induced toxicity on the range (1×10^{-4} , 1×10^1) respectively. Both ranges were selected to match the ranges provided in the literature.³⁴ Latin Hypercube sampling was used to generate $N = 100$ perturbed parameter values, and the optimal solutions were compared to an equal number of background solutions using the nominal parameter values. (D) Visualization of the optimal solutions; scatter plot of principal components of the optimal parameter values for the model with perturbed toxicity parameters ($N = 100$). Contour plots show the background distribution of parameter values.

when the optimization is run many times. No architecture is optimal for more than 15% of test runs, demonstrating that

there are combinations of architectures and parameter values that achieve a similar optimal loss. We also found that several architectures can display oscillatory solutions, which we chose to exclude from the search by applying a peak detection algorithm⁴⁸ and adding a large regularization term to the loss.

To investigate the robustness to chemical toxicity, we perturbed the metabolite-induced toxicity t_a and protein-induced toxicity t_p in eq 6. The optimal loss values were found to be comparable between perturbed and background systems (Figure 4C). Additionally, when projected onto a 2-dimensional space using principal component analysis, the distribution of background parameter values was similar to the distribution of perturbed solutions, indicating that the perturbation did not significantly affect the optimal parameters selected (Figure 4D).⁴⁹ The *p*-AS pathway lies at the far end of what is currently possible to build experimentally and thus illustrates the broad applicability of BayesOpt to realistic design tasks in metabolic engineering.

3. DISCUSSION

Progress in synthetic biology allows the construction of circuits of increased complexity across various levels of biological organization. However, large design spaces and multiple scales can become substantial challenges for the design of functional systems. In this paper, we presented the use of Bayesian optimization for the design of biological circuits. The method can rapidly find circuit architectures and parameters that optimize a performance objective that captures the target circuit functionality.

The method is particularly well suited for cases in which the multiple scales prevent efficient simulation of ODE models. Gene circuits designed to control metabolic pathways are an excellent example of such multiscale systems, as they combine fast metabolic time scales with the much slower dynamics of gene expression. Moreover, the choice of regulators, control points, and control architectures adds multiple degrees of freedom that are infeasible to explore experimentally. Previously implemented metabolic control systems have been built primarily based on application-specific knowledge of pathway features.^{27,50} We have shown that Bayesian optimization can aid the design of such systems prior to implementation and serve as tools for *in silico* screening of competing designs that may have similar performance but entail different cost of wetlab implementation. We showed the efficiency and scalability of the method in several real world case studies from metabolic engineering. In particular, the *p*-aminostyrene pathway is more complex than systems typically implemented in literature so far, which suggests that the method is applicable across a range of relevant design tasks.

We anticipate several novel applications of this work to other problem areas where discovery or tuning of multiscale circuits has been previously infeasible. For instance, this method could be employed to fit temporal circuit dynamics to data or discern which of several discrete circuit mechanisms most closely matches observed behavior. As with other design strategies based on ODE models, a challenge of our approach is the significant domain knowledge required to construct models for a target pathway, both in terms of the enzyme kinetics and the downstream metabolic processes that affect pathway activity. Machine learning has already proved useful in a range of metabolic engineering tasks⁵¹ and is gaining substantial interest in other areas of synthetic biology.^{52,53} In this paper we have shown how such methods can also benefit dynamic pathway

Table 1. Summary of Pathway Models Studied in This Paper. The ODEs in the *p*-Aminostyrene Pathway Also Include mRNA and Folding Dynamics

product	parameters (p_c)	architectures (p_d)	metabolites	enzymes	ODEs
toy pathway	4	4	2	2	4
glucaric acid ^{10,26}	4	4	3	2	5
fatty acid ³³	2	4	1	2	3
<i>p</i> -aminostyrene ³⁴	16	27	7	6	23

engineering by using optimization as a means to navigate the design space prior to system prototyping.

4. METHODS

4.1. Bayesian Optimization. We employed the Bayesian optimization routine implemented in the Python HyperOpt package.³⁷ Bayesian optimization is commonly employed for hyperparameter tuning in deep neural networks. We employed Expected Improvement as an acquisition function and a tree-structured Parzen estimator (TPE) as a nonparametric statistical model for the loss landscape. We performed a grid search over the TPE hyperparameter γ which controls the balance between exploration and exploitation but found little impact on the algorithm performance; we thus used the default value of $\gamma = 15$ (Supporting Figure S1).

Constraints on the continuous and discrete decision variables were incorporated directly into the HyperOpt search space. At each run of the Bayesian optimization routine, the initial guess for the continuous decision variables were sampled from uniform distributions, with upper and lower bounds were taken from literature.^{10,34,44} Architectures were chosen uniformly from the set of architectures without positive feedback loops.

4.2. Model Pathways. We considered four exemplar pathways modeled via ordinary differential equations (ODEs): the toy system in Figure 1C, the glucaric acid pathway in Figure 2A, the fatty acid pathway in Figure 3A, and the *p*-aminostyrene pathway in Figure 4A. Table 1 contains a summary of the four considered models. In all cases, pathway models include ODEs for both metabolites and pathway enzymes. In each case, we define the various control architectures and incorporate them as discrete decision variables in the optimization problem, i.e., p_d in eq 1; the continuous decision variables, i.e., p_c in eq 1, appear in the expression rates of the pathway enzymes. For the toy model and the glucaric acid pathway, enzyme expression was parametrized using a lumped Hill equation model to describe the interaction between a regulatory metabolite and a transcription factor. For the fatty acid and *p*-aminostyrene pathways, expression rates were parametrized with bespoke nonlinear functions describing specific biochemical processes. The discrete control architectures were defined in two different ways. For the toy, glucaric acid, and *p*-aminostyrene models, the architectures were defined using a binary matrix to encode the mode of transcriptional control. For the fatty acid model we instead defined each architecture as a categorical choice and switched between model functions correspondingly. We note that the *p*-aminostyrene pathway also contains ODEs for mRNA abundance and folded/unfolded proteins. All models and their parameters are described in the Supporting Information.

The ODE models were solved with scikit-odes, a Python wrapper for the SUNDIALS suite of solvers.⁵⁴ In all cases, the initial concentrations of heterologous pathway enzymes were

assumed to be zero. Initial concentrations for native metabolites were determined by first solving a model without the heterologous enzymes up to steady state. Simulation times and initial conditions are detailed in the Supporting Information for each model.

4.3. Loss Function. In all cases the loss function J in eq 3 was instantiated to each pathway. Generally, the loss is defined as a linear combination of costs and benefits of pathway activity so as to balance opposing design goals commonly found in applications. Since both components of the loss function have different magnitudes, for each model we first swept the weights α_1 and α_2 across many model simulations, and chose values that led to similar values for both components; this prevents the optimizer from biasing the search towards low loss values caused by the scaling effects alone. For the fatty acid model in Figure 3 we also optimized the circuit rise time and overshoot % defined in eq 5. Details on all objective functions can be found in the Supporting Information.

■ ASSOCIATED CONTENT

Data Availability Statement

The Python code for this paper is available on Zenodo at <https://doi.org/10.5281/zenodo.7926205>.

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acssynbio.3c00120>.

Details of model construction for toy, glucaric acid, fatty acid, and *p*-aminostyrene models, including parameter values and model equations; Additional supporting figures related to hyperparameter tuning (PDF)

■ AUTHOR INFORMATION

Corresponding Author

Diego A. Oyarzún — School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, U.K.; The Alan Turing Institute, London NW1 2DB, U.K.; School of Biological Sciences, University of Edinburgh, Edinburgh EH9 3JH, U.K.; orcid.org/0000-0002-0381-5278; Email: d.oyarzun@ed.ac.uk

Authors

Charlotte Merzbacher — School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, U.K.; orcid.org/0009-0000-2853-1864

Oisín Mac Aodha — School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, U.K.; The Alan Turing Institute, London NW1 2DB, U.K.

Complete contact information is available at: <https://pubs.acs.org/doi/10.1021/acssynbio.3c00120>

Author Contributions

CM built the optimization pipeline, ran simulations, and produced figures. CM and DAO analyzed results. DAO and OMA supervised the research.

Funding

CM and DAO were supported by the United Kingdom Research and Innovation (grant EP/S02431X/1, UKRI Centre for Doctoral Training in Biomedical AI).

Notes

The authors declare no competing financial interest.

REFERENCES

- (1) Brophy, J. A. N.; Voigt, C. A. Principles of genetic circuit design. *Nat. Methods* **2014**, *11*, 508–520.
- (2) Shaw, W. M.; Yamauchi, H.; Mead, J.; Gowers, G. O. F.; Bell, D. J.; Öling, D.; Larsson, N.; Wigglesworth, M.; Ladds, G.; Ellis, T. Engineering a Model Cell for Rational Tuning of GPCR Signaling. *Cell* **2019**, *177*, 782–796.
- (3) Zhang, F.; Carothers, J. M.; Keasling, J. D. Design of a dynamic sensor-regulator system for production of chemicals and fuels derived from fatty acids. *Nat. Biotechnol.* **2012**, *30*, 354–9.
- (4) Ma, W.; Trusina, A.; El-Samad, H.; Lim, W. A.; Tang, C. Defining network topologies that can achieve biochemical adaptation. *Cell* **2009**, *138*, 760–773.
- (5) Li, Z.; Liu, S.; Yang, Q. Incoherent inputs enhance the robustness of biological oscillators. *Cell Systems* **2017**, *5*, 72–81.
- (6) Qiao, L.; Zhao, W.; Tang, C.; Nie, Q.; Zhang, L. Network topologies that can achieve dual function of adaptation and noise attenuation. *Cell Systems* **2019**, *9*, 271–285.
- (7) Dasika, M. S.; Maranas, C. D. OptCircuit: An optimization based method for computational design of genetic circuits. *BMC Syst. Biol.* **2008**, *2*, 1–19.
- (8) Otero-Muras, I.; Banga, J. R. Automated Design Framework for Synthetic Biology Exploiting Pareto Optimality. *ACS Synth. Biol.* **2017**, *6*, 1180–1193.
- (9) Hiscock, T. W. Adapting machine-learning algorithms to design gene circuits. *BMC Bioinformatics* **2019**, *20*, 1–13.
- (10) Verma, B. K.; Mannan, A. A.; Zhang, F.; Oyarzún, D. A. Trade-offs in biosensor optimization for dynamic pathway engineering. *ACS Synth. Biol.* **2022**, *11*, 228–240.
- (11) Banga, J. Optimization in computational systems biology. *BMC Syst. Biol.* **2008**, *2*, 47.
- (12) Hairer, E.; Wanner, G. *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*; Springer-Verlag, 1996.
- (13) Blanchini, F.; Franco, E.; Giordano, G. A structural classification of candidate oscillatory and multistationary biochemical systems. *Bulletin of Mathematical Biology* **2014**, *76*, 2542–69.
- (14) Briat, C.; Gupta, A.; Khammash, M. Antithetic proportional-integral feedback for reduced variance and improved control performance of stochastic reaction networks. *J. R. Soc., Interface* **2018**, *15*, 20180079.
- (15) Briat, C.; Gupta, A.; Khammash, M. Antithetic integral feedback ensures robust perfect adaptation in noisy biomolecular networks. *Cell Systems* **2016**, *2*, 15–26.
- (16) Bhattacharya, P.; Raman, K.; Tangirala, A. K. Discovering adaptation-capable biological network structures using control-theoretic approaches. *PLOS Computational Biology* **2022**, *18*, No. e1009769.
- (17) Araujo, R. P.; Liotta, L. A. The topological requirements for robust perfect adaptation in networks of any size. *Nat. Commun.* **2018**, *9*, 1757.
- (18) Drengstig, T.; Ueda, H. R.; Ruoff, P. Predicting perfect adaptation motifs in reaction kinetic networks. *J. Phys. Chem. B* **2008**, *112*, 16752–16758.
- (19) Woods, M. L.; Leon, M.; Perez-Carrasco, R.; Barnes, C. P. A statistical approach reveals designs for the most robust stochastic gene oscillators. *ACS Synth. Biol.* **2016**, *5*, 459–470.
- (20) Gonzalez, J.; Longworth, J.; James, D. C.; Lawrence, N. D. Bayesian optimization for synthetic gene design. *arXiv*, May 7, 2015. DOI: 10.48550/arXiv.1505.01627.
- (21) Shen, J.; Liu, F.; Tu, Y.; Tang, C. Finding gene network topologies for given biological function with recurrent neural network. *Nat. Commun.* **2021**, *12*, 1–10.
- (22) Zhang, F.; Carothers, J. M.; Keasling, J. D. Design of a dynamic sensor-regulator system for production of chemicals and fuels derived from fatty acids. *Nat. Biotechnol.* **2012**, *30*, 354–359.
- (23) Oyarzún, D. A.; Stan, G.-B. V. Synthetic gene circuits for metabolic control: design trade-offs and constraints. *J. R. Soc., Interface* **2013**, *10*, 20120671.
- (24) Xu, P.; Li, L.; Zhang, F.; Stephanopoulos, G.; Koffas, M. Improving fatty acids production by engineering dynamic pathway regulation and metabolic control. *Proc. Natl. Acad. Sci. U. S. A.* **2014**, *111*, 11299–11304.
- (25) Dunlop, M. J.; Keasling, J. D.; Mukhopadhyay, A. A model for improving microbial biofuel production using a synthetic feedback loop. *Systems and Synthetic Biology* **2010**, *4*, 95–104.
- (26) Doong, S. J.; Gupta, A.; Prather, K. L. Layered dynamic regulation for improving metabolic pathway productivity in *Escherichia coli*. *Proc. Natl. Acad. Sci. U. S. A.* **2018**, *115*, 2964–2969.
- (27) Ni, C.; Dinh, C. V.; Prather, K. L. Dynamic Control of Metabolism. *Annu. Rev. Chem. Biomol. Eng.* **2021**, *12*, 519.
- (28) Hartline, C. J.; Schmitz, A. C.; Han, Y.; Zhang, F. Dynamic control in metabolic engineering: Theories, tools, and applications. *Metabolic Engineering* **2021**, *63*, 126–140.
- (29) Tonn, M. K.; Thomas, P.; Barahona, M.; Oyarzún, D. A. Stochastic modelling reveals mechanisms of metabolic heterogeneity. *Commun. Biol.* **2019**, *2*, 108.
- (30) Mannan, A. A.; Liu, D.; Zhang, F.; Oyarzún, D. A. Fundamental Design Principles for Transcription-Factor-Based Metabolite Biosensors. *ACS Synth. Biol.* **2017**, *6*, 1851–1859.
- (31) Zhou, L.-B.; Zeng, A.-P. Exploring Lysine Riboswitch for Metabolic Flux Control and Improvement of L-Lysine Synthesis in *Corynebacterium glutamicum*. *ACS Synth. Biol.* **2015**, *4*, 729–734.
- (32) Chaves, M.; Oyarzún, D. A. Dynamics of complex feedback architectures in metabolic pathways. *Automatica* **2019**, *99*, 323–332.
- (33) Liu, D.; Zhang, F. Metabolic feedback circuits provide rapid control of metabolite dynamics. *ACS Synth. Biol.* **2018**, *7*, 347–356.
- (34) Stevens, J. T.; Carothers, J. M. Designing RNA-based genetic control systems for efficient production of engineered metabolic pathways. *ACS Synth. Biol.* **2015**, *4*, 107–115.
- (35) Frazier, P. I. A tutorial on Bayesian optimization. *arXiv*, July 8, 2018. DOI: 10.48550/arXiv.1807.02811.
- (36) Moon, T. S.; Yoon, S.-H.; Lanza, A. M.; Roy-Mayhew, J. D.; Prather, K. L. J. Production of glucaric acid from a synthetic pathway in recombinant *Escherichia coli*. *Applied and Environmental Microbiology* **2009**, *75*, 589–595.
- (37) Bergstra, J.; Yamins, D.; Cox, D. Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In *International Conference on Machine Learning*; ICML, 2013; pp 115–123.
- (38) Snoek, J.; Larochelle, H.; Adams, R. P. Practical bayesian optimization of machine learning algorithms. In *Advances in Neural Information Processing Systems*; NeurIPS, 2012; Vol. 25.
- (39) Bar-Even, A.; Noor, E.; Savir, Y.; Liebermeister, W.; Davidi, D.; Tawfik, D. S.; Milo, R. The moderately efficient enzyme: evolutionary and physicochemical trends shaping enzyme parameters. *Biochemistry* **2011**, *50*, 4402–4410.
- (40) Venayak, N.; Anesiadis, N.; Cluett, W. R.; Mahadevan, R. Engineering metabolism through dynamic control. *Curr. Opin. Biotechnol.* **2015**, *34*, 142–152.
- (41) Hartline, C.; Mannan, A.; Liu, D.; Zhang, F.; Oyarzún, D. Metabolite sequestration enables rapid recovery from fatty acid depletion in *Escherichia coli*. *mBio* **2020**, *11*, S90943.
- (42) Zhu, Y.; Li, Y.; Xu, Y.; Zhang, J.; Ma, L.; Qi, Q.; Wang, Q. Development of bifunctional biosensors for sensing and dynamic

control of glycolysis flux in metabolic engineering. *Metabolic Engineering* **2021**, 68, 142–151.

(43) Schomburg, I.; Jeske, L.; Ulbrich, M.; Placzek, S.; Chang, A.; Schomburg, D. The BRENDA enzyme information system—From a database to an expert system. *Journal of biotechnology* **2017**, 261, 194–206.

(44) Zhang, Y.; Nielsen, J.; Liu, Z. Metabolic engineering of *Saccharomyces cerevisiae* for production of fatty acid-derived hydrocarbons. *Biotechnology and Bioengineering* **2018**, 115, 2139–2147.

(45) Goikhman, M. Y.; Yevlampieva, N.; Kamanina, N.; Podeshvo, I.; Gofman, I.; Mil'tsov, S.; Khurchak, A.; Yakimanskii, A. New polyamides with main-chain cyanine chromophores. *Polymer Science Series A* **2011**, 53, 457–468.

(46) Ellington, A. D.; Szostak, J. W. In vitro selection of RNA molecules that bind specific ligands. *Nature* **1990**, 346, 818–822.

(47) Oyarzún, D. A.; Chaves, M. Design of a bistable switch to control cellular uptake. *Journal of The Royal Society Interface* **2015**, 12, 20150618.

(48) Virtanen, P.; Gommers, R.; Oliphant, T. E.; Haberland, M.; Reddy, T.; Cournapeau, D.; Burovski, E.; Peterson, P.; Weckesser, W.; Bright, J.; et al. SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **2020**, 17, 261–272.

(49) Abdi, H.; Williams, L. J. Principal component analysis. *Wiley Interdisciplinary Reviews: Computational Statistics* **2010**, 2, 433–459.

(50) Liu, D.; Mannan, A. A.; Han, Y.; Oyarzún, D. A.; Zhang, F. Dynamic metabolic control: towards precision engineering of metabolism. *Journal of Industrial Microbiology and Biotechnology* **2018**, 45, 535–543.

(51) Radivojević, T.; Costello, Z.; Workman, K.; García Martín, H. A machine learning Automated Recommendation Tool for synthetic biology. *Nat. Commun.* **2020**, 11, 1–14.

(52) Carbonell, P.; Radivojevic, T.; García Martín, H. Opportunities at the Intersection of Synthetic Biology, Machine Learning, and Automation. *ACS Synth. Biol.* **2019**, 8, 1474–1477.

(53) Nikolados, E.-M.; Wongprommoon, A.; Aodha, O. M.; Cambray, G.; Oyarzún, D. A. Accuracy and data efficiency in deep learning models of protein expression. *Nat. Commun.* **2022**, 13, 7755.

(54) Gardner, D. J.; Reynolds, D. R.; Woodward, C. S.; Balos, C. J. Enabling new flexibility in the SUNDIALS suite of nonlinear and differential/algebraic equation solvers. *ACM Trans. Math. Softw.* **2022**, 48, 1.

(55) Solgi, R. M. *Geneticalgorithm Package*; 2020. <https://pypi.org/project/geneticalgorithm/>.

(56) Loh, W.-L. On Latin hypercube sampling. *Ann. Statist.* **1996**, 24, 2058–2080.