



HAL
open science

Bayesian filtering for model predictive control of stochastic gene expression in single cells

Zachary Fox, Gregory Batt, Jakob Ruess

► **To cite this version:**

Zachary Fox, Gregory Batt, Jakob Ruess. Bayesian filtering for model predictive control of stochastic gene expression in single cells. *Physical Biology*, 2023, 10.1088/1478-3975/ace094 . hal-04148503

HAL Id: hal-04148503

<https://inria.hal.science/hal-04148503>

Submitted on 3 Jul 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

Bayesian filtering for model predictive control of stochastic gene expression in single cells

Zachary R Fox*

Gregory Batt**

Jakob Ruess**

Abstract—This study describes a method for controlling the production of protein in individual cells using stochastic models of gene expression. By combining modern microscopy platforms with optogenetic gene expression, experimentalists are able to accurately apply light to individual cells, which can induce protein production. Here we use a finite state projection based stochastic model of gene expression, along with Bayesian state estimation to control protein copy numbers within individual cells. We compare this method to previous methods that use population based approaches. We also demonstrate the ability of this control strategy to ameliorate discrepancies between the predictions of a deterministic model and stochastic switching system.

I. INTRODUCTION

Modern microscopy platforms are revolutionizing the quality and quantity of biological data. Synthetic biology, and in particular optogenetics provide novel ways to quantify and perturb biological systems at the single-cell level. These platforms allow for the first time online control of protein production in single cells [1], [2], [3], [4], [5]. However, protein production in single cells is stochastic [6], and novel methods must be developed to use models to control single cell gene expression. When a model is available, model predictive control (MPC) is known to have many benefits. From a control perspective, controlling the protein expression in individual cells should reduce the error of each cell with respect to the target. Until now, most attempts to control gene expression have been at the population level [3], [7], [8]. The distinction is shown in Fig. 1.

To date, few studies have attempted to perform model predictive control on multiple individual cells simultaneously using modern platforms [2], [5]. Of these two studies, only [5] uses a stochastic model of the process to control gene expression in individual bacteria. However, their model uses moment-based approximations of the stochastic dynamics, and it is known that such approximations are not always appropriate [9].

In this study, we introduce a new method to perform single-cell control of optogenetic production of protein. This method uses a simple, but stochastic, reaction-based model to control the stochastic gene expression within each cell in the population. For each cell, we estimate the intracellular protein levels based on fluorescence measurements and the

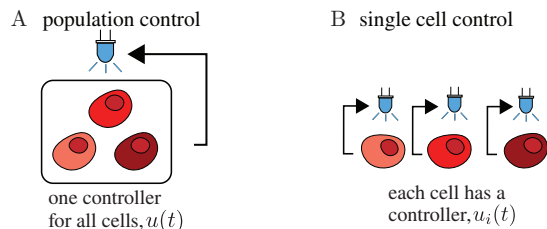


Fig. 1. Comparison of population control and single-cell control (A) In population control, there is a single controller for the entire ensemble of cells, and measurements are of the average cellular fluorescence. (B) Each cell is controlled and measured independently.

stochastic model of the process using a Bayesian filter. Once this intracellular probability distribution has been estimated, we use receding horizon model predictive control, along with the stochastic model, to drive intracellular protein values to pre-determined target values.

We compare the single-cell approach to the population based approach. In the single-cell approach, each cell has its own controller, while in the population approach, one controller is used for all cells. We evaluate these approaches computationally for two different models of optogenetically controlled gene expression. The first is a model of a light-driven gene expression system using the EL222 transcription factor [10], which has been used experimentally in [11], [2], [12]. The second is a simplified single-species nonlinear model of a self-regulated gene.

II. STOCHASTIC MODELING AND CONTROL OF GENE EXPRESSION

A. Modeling gene expression

Gene expression is inherently stochastic [6], due to the random diffusion of the various molecules involved in the process of transcription and translation. In general, the chemical master equation governs the time evolution of an infinite state continuous time Markov chain with generator matrix $\mathbf{A}(u(t), \theta)$. We consider the matrix \mathbf{A} as a function of control parameter $u(t)$, and of kinetic rate parameters θ . The i^{th} discrete state in the Markov chain is a node on the lattice, $\mathbf{x}_i = \{\zeta_1, \zeta_2, \dots, \zeta_n\}$ for the n species in the chemical system, where each ζ_i is an integer molecule count. The vector \mathbf{p} collects the probabilities of each state or configuration as $\mathbf{p} = [p(\mathbf{x}_1), p(\mathbf{x}_2), \dots]^T$. We denote the total set of states as \mathcal{X} . Note that this vector is often infinite in its dimension. The CME is the set of ODEs which describes the time evolution

* Corresponding author. ZRF is currently with the Computational Science and Engineering Division at Oak Ridge National Lab in Oak Ridge, TN, USA, and with Inria Paris and Institut Pasteur, Paris, France. Email: foxzr@ornl.gov

**GB and JR are with Inria, Institut Pasteur, Université Paris Cité, Paris, France. Email: {gregory.batt, jakob.ruess}@inria.fr

of these probabilities, i.e.

$$\frac{d}{dt}\mathbf{p} = \mathbf{A}(u(t), \theta)\mathbf{p}. \quad (1)$$

We use the finite state projection (FSP) approach to solve the chemical master equation (CME)[13], [14], [15]. The FSP approach splits this infinite state Markov chain into two distinct sets, \mathcal{J} and \mathcal{J}' ,

$$\frac{d}{dt} \begin{bmatrix} \mathbf{p}_{\mathcal{J}} \\ \mathbf{p}_{\mathcal{J}'} \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{\mathcal{J}\mathcal{J}} & \mathbf{A}_{\mathcal{J}\mathcal{J}'} \\ \mathbf{A}_{\mathcal{J}'\mathcal{J}} & \mathbf{A}_{\mathcal{J}'\mathcal{J}'} \end{bmatrix} \begin{bmatrix} \mathbf{p}_{\mathcal{J}} \\ \mathbf{p}_{\mathcal{J}'} \end{bmatrix}, \quad (2)$$

and collects the (infinite) set of states \mathcal{J}' into a single absorbing sink state $g(t)$,

$$\frac{d}{dt} \begin{bmatrix} \mathbf{p}_{\mathcal{J}} \\ g(t) \end{bmatrix} = \begin{bmatrix} \mathbf{A}_{\mathcal{J}\mathcal{J}} & \mathbf{0} \\ -\mathbf{1}^T \mathbf{A}_{\mathcal{J}\mathcal{J}} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{p}_{\mathcal{J}} \\ g(t) \end{bmatrix}. \quad (3)$$

The sink state $g(t)$ is the error in the FSP solution, which monotonically increases with t [15]. As such, the FSP dimension (i.e. the cardinality of \mathcal{J}) can be modified to decrease the error to a prescribed value for any finite time. Solving the FSP can be extremely computationally expensive in general, as the number of elements in the matrix \mathbf{A} grows exponentially with the number of species considered in the model, and therefore one typically does not include more than three or four species in a given model [16]. For convenience, for the remainder of the work we will refer to $\mathbf{A}_{\mathcal{J}\mathcal{J}}$ as \mathbf{A} , and $\mathbf{p}_{\mathcal{J}}$ as \mathbf{p} .

In this study, we are considering controllers that are piecewise-constant, and therefore \mathbf{A} is constant over fixed time intervals. Furthermore, we are considering a control which has either an on or off state, i.e. there is either light shining on a given cell or not. When \mathbf{A} is constant, the solution at t_f given $\mathbf{p}(t_0)$ is $\mathbf{p}(t_f) = e^{\mathbf{A}(t_f-t_0)}\mathbf{p}(t_0)$. We can pre-compute the infinitesimal generator for both the on and off state of the controller, where we define the ‘‘on’’ generator matrix as \mathbf{A}_{on} and the ‘‘off’’ generator matrix as \mathbf{A}_{off} . Given a sequence of input $u(t)$ the FSP solution can be iteratively evaluated on the constant time intervals,

$$\mathbf{p}(t + \Delta t) = \begin{cases} e^{\mathbf{A}_{\text{on}}\Delta t}\mathbf{p}(t) & \text{if } u(s) = 1 \text{ for } s \in (t, t + \Delta t) \\ e^{\mathbf{A}_{\text{off}}\Delta t}\mathbf{p}(t) & \text{if } u(s) = 0 \text{ for } s \in (t, t + \Delta t). \end{cases} \quad (4)$$

Equation 4 gives the forward map of the probability distribution $\mathbf{p}(t)$ to $\mathbf{p}(t + \Delta t)$, which is key for the Bayesian filter discussed next. For the remainder of the manuscript, we consider a system which is measured every Δt time units, for a total of N_t measurements.

B. Bayesian filtering for single cells using the FSP

In general, the probability distribution of molecule counts modeled by the chemical master equation is non-Gaussian; the domain is discrete and low-copy effects abound [9]. Therefore, state estimation techniques that assume symmetric, Gaussian variability such as the Kalman filter are not appropriate, and instead we turn to the more general Bayesian filtering [17].

The Bayesian filter in this work serves as a way to filter out observation noise when we have arbitrary state

distributions which come from the solution of the chemical master equation solved with the FSP. A Bayesian filtering approach considers a prediction step and an estimation step, assuming an initial distribution $\mathbf{p}(t_0)$. We start with the distribution of the estimated state of the cell at the previous time $\mathbf{p}((k-1)\Delta t) = [p_0, p_1, \dots, p_N]$. This distribution can be propagated from $(k-1)\Delta t$ to $k\Delta t$ according to $u(t)$ and Eq. 4.

To obtain an update rule for the distribution of the estimated state at measurement time points, we need to specify a model that describes how the state is mapped to the observable outputs in experiments. Because measurements of fluorescently labeled biomolecules are intrinsically noisy due to the photophysics of fluorescent molecules and detection of photons, the fluorescent measurement of a single cell is described probabilistically. Each molecule emits a random number of photons, some of which are detected by the camera. The distribution of the number of photons emitted by a single fluorescent protein is typically taken to be Poisson, but in the limit of many emissions is approximately Gaussian, i.e. $v_k \sim \mathcal{N}(\mu_{\text{FP}}, \sigma_{\text{FP}}^2)$. For n_i such molecules in a cell, the probability of measuring fluorescence z_k is given by

$$f_{\text{FP}}(z) = \frac{1}{\sqrt{2n_i\pi\sigma_{\text{FP}}^2}} \exp\left(-\frac{(z - n_i\mu_{\text{FP}})^2}{2n_i\sigma_{\text{FP}}^2}\right), \quad (5)$$

as the sum of Gaussian random variables is a Gaussian random variable. Together, Eqns. 4 and 5 can be used to iteratively estimate $\mathbf{p}(k\Delta t)$ based on the measurement z_k . The FSP-based Bayesian filtering is detailed in Algorithm 1. The algorithm is initialized by defining an initial distribution of molecules $\mathbf{p}(t_0)$. For each measurement period Δt , the posterior estimate of the probability distribution $\tilde{\mathbf{p}}((k-1)\Delta t)$ from the previous measurement is propagated to the current time $k\Delta t$. The appropriate generator \mathbf{A}_{on} or \mathbf{A}_{off} is used depending on the prescribed actuation during the time interval. Finally, the measurement at time $k\Delta t$, z_k is incorporated to the model according to Eq. 5 and Bayes’ rule. The computation of the normalizing factor \mathcal{Z} in line 10 of Algorithm 1 is not feasible in general. This is because computation of the marginal likelihood requires the integration over continuous variables, except for cases such as the Gaussian and other simple parametric likelihoods and conjugate priors. However, for the FSP-based approach, the marginalization is simply a sum over the discrete and finite state space,

$$\mathcal{Z} = \sum_{x \in \mathcal{X}} f_{\text{FP}}(z_k|x; \mu_{\text{FP}}, \sigma_{\text{FP}}^2)p(x). \quad (6)$$

This fact also enables easy computation when dealing with hidden species; the marginalization can be taken over the unobserved variables.

To evaluate our approach, we will compare the FSP-based Bayesian filtering algorithms described above to the more commonly used population based Kalman filter.

C. Kalman Filtering for population state estimation

To contrast single-cell control, we introduce here the notion of a population model. In the population model, indi-

Algorithm 1 FSP-based Bayesian filtering

```

1:  $k = 1$ 
2:  $\tilde{\mathbf{p}}(k\Delta t) = \mathbf{p}(t_0)$ 
3: while  $k < \Delta t N_t$  do
4:   if  $u(s) = 1$  for  $s \in ((k-1)\Delta t, k\Delta t)$  then
5:      $\mathbf{p}(k\Delta t) = e^{\mathbf{A}_{\text{on}}\Delta t} \tilde{\mathbf{p}}((k-1)\Delta t)$ 
6:   else
7:      $\mathbf{p}(k\Delta t) = e^{\mathbf{A}_{\text{off}}\Delta t} \tilde{\mathbf{p}}((k-1)\Delta t)$ 
8:   end if
9:    $z_k = \text{measurement}$ 
10:   $\tilde{\mathbf{p}}(k\Delta t) = \frac{1}{Z} f_{\text{FP}}(z_k | \mathbf{n}; \mu_{\text{FP}}, \sigma_{\text{FP}}^2) \mathbf{p}(k\Delta t)$ 
11:   $k = k + 1$ 
12: end while

```

vidual cells are not measured, and the fluorescence measurements come from averaging the fluorescence across the entire population of cells. We use the Linear Noise Approximation for stochastic chemical kinetics to approximate the dynamics of the mean and variance in the CME in Eq. 1 [13]. The LNA is valid when the number of interacting molecules and system volume are sufficiently large. These dynamics are given by

$$\frac{d}{dt} \bar{\mathbf{x}} = \phi(\bar{\mathbf{x}}, \theta, u(t)) \quad (7)$$

$$\frac{d}{dt} \Sigma = \rho(\bar{\mathbf{x}}, \Sigma, \theta, u(t)). \quad (8)$$

The equations for ϕ and ρ can be found in [18]. The mean molecule number for the population is approximated by $\bar{\mathbf{x}}(t)$. The variance about $\bar{\mathbf{x}}(t)$ is approximated by $\Sigma(t)$ for each single cell, and therefore the variance of the population mean is simply given by $\mathbf{P} = \Sigma(t)/N_c$, where N_c is the number of cells that are averaged. Under the LNA, we consider the state and process noise to follow a normal distribution $\mathcal{N}(\bar{\mathbf{x}}, \mathbf{P})$. Therefore, we can use the LNA to construct a Kalman filter to iteratively estimate the population's mean molecule number using population-averaged fluorescence measurements, y_k . The previous state and process noise estimate at the $(k-1)$ th time, $\tilde{\mathbf{x}}((k-1)\Delta t) \equiv \tilde{\mathbf{x}}_{k-1|k-1}$ and $\tilde{\mathbf{P}}((k-1)\Delta t) \equiv \tilde{\mathbf{P}}_{k-1|k-1}$, are used as the initial conditions in Eq. 7 and propagated forward Δt to find the $\tilde{\mathbf{x}}_{k|k-1}$ and $\tilde{\mathbf{P}}_{k|k-1}$.

Finally, because not all species are observed, we define an observability matrix, \mathbf{C} which linearly combines the state variables to an observable state. Therefore, we can write the filtering equations

$$\tilde{\mathbf{x}}_{k|k} = \tilde{\mathbf{x}}_{k|k-1} + \mathbf{K}_k (y_k - \mathbf{C} \tilde{\mathbf{x}}_{k|k-1}) \quad (9)$$

$$\tilde{\mathbf{P}}_{k|k} = (\mathbf{I} - \mathbf{K}_k \mathbf{C}) \tilde{\mathbf{P}}_{k|k-1} \quad (10)$$

where \mathbf{K}_k is the Kalman gain,

$$\mathbf{K}_k = \tilde{\mathbf{P}}_{k|k-1} \mathbf{C}^T \left(\mathbf{C} \tilde{\mathbf{P}}_{k|k-1} \mathbf{C}^T + R \right)^{-1}, \quad (11)$$

where R is the additive observation noise which comes from the measurement process.

D. Receding Horizon Model Predictive Control

For both the population model and the single-cell stochastic model we implemented a receding horizon model predictive control. Thus far, we have left the Bayesian and Kalman filtering methods to arbitrary dimension. For the remainder of the manuscript we focus on controlling a particular molecular species abundance within the system, and we denote the molecular abundance of this controlled species as x^c and the mean molecular abundance of this species as \bar{x}^c . For a prescribed set of times $i/\Delta t \in [1, 2, \dots, N_t]$, we aim to control the molecular abundance of either single-cells or the population to follow a target value, T_i . In the examples below, the target is in terms of molecule number of fluorescent protein. This cost function is evaluated after each measurement update $k\Delta t$ for a finite horizon H . The horizon H is defined in units of Δt . For the standard population control, we defined a sum of squares cost function of the prediction horizon H ,

$$J_{\text{pop}} = \frac{1}{H} \sum_{i=k}^{k+H} (T_i - \bar{x}_i^c)^2. \quad (12)$$

In the L_2 cost in Eq. 12, the error is computed with respect to the mean \bar{x}^c to be consistent with the assumption that population model makes; namely that the population mean \bar{x} is Gaussian distributed. However, it is well known that individual cells' molecule counts often follow more complex distribution and approaches that rely on averaging and the central limit theorem are not valid [9]. Our goal is to find the control that moves this complex distribution towards a target molecule count. We chose the expected absolute deviation, as it allows the cost function to account for intrinsic asymmetries and non-Gaussian features in $\tilde{\mathbf{p}}$:

$$J_{\text{sc}} = \frac{1}{H} \sum_{i=k}^{k+H} \sum_{j=0}^N |T_i - x_j^c| \tilde{p}(x_j^c, i\Delta t). \quad (13)$$

The summation over j is taken sum over all values of the observable molecular abundance in the FSP analysis. We perform receding horizon MPC, in which we aim to optimize the cost function J over the finite time horizon H using the current state estimate, i.e. to find the light input sequence $u(t)$ which minimizes J over $[k\Delta t, k\Delta t + H]$, i.e.

$$u(t) = \arg \min_{u(t)} J(u(t)). \quad (14)$$

For these systems, the photostimulation is taken either to be on or off between $k\Delta t$ and $(k+1)\Delta t$, therefore there are $2^{H/\Delta t}$ possible light combinations over the horizon H . We optimized by exhaustively evaluating Eq. 14 for each possible light sequence. Other tree-based searches or "greedy algorithms" may be useful in optimizing finding controllers without exhaustively searching through all possible combinations.

III. COMPUTATIONAL RESULTS

A. Control of a light-driven expression system

Optogenetically controlled protein production is now possible in bacteria, yeast, and higher eukaryotes [19]. The

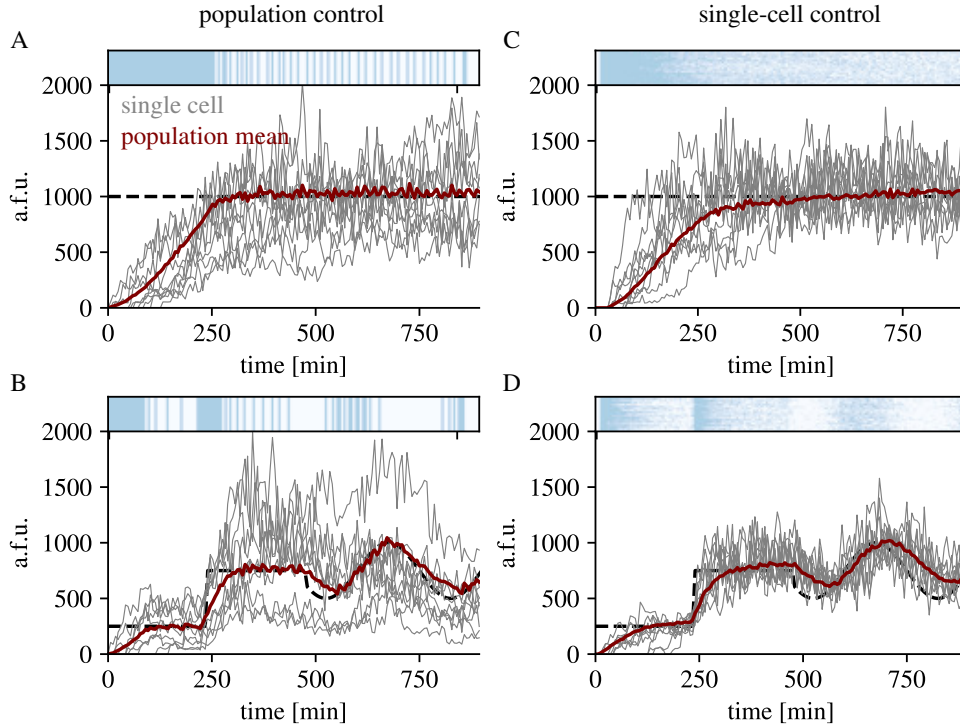


Fig. 2. (A-B) Controlling the population mean (red), (A-B) and individual cells (gray), (B-D) for a static target (A,C) and target which changes over time (B,D). We simulated a population of 100 cells, 10 of which are shown as individual gray lines. The light inputs are shown in blue at the top of each plot, in which dark blue indicates blue light shown to the entire population (A-B) or individual cells (B-C).

general idea is that light at specific wavelengths can be used to activate transcription factors that “turn on” protein production within cells [20], [10]. Optogenetic systems often react faster than traditional chemically induced protein production, and downstream responses can be manipulated using frequency and/or amplitude modulation of the input light stimulation [11].

Here, we work specifically with the EL222 optogenetic system. In this system, EL222 molecules are created and degraded at all times within the cell. Under blue light photostimulation, EL222 molecules dimerize and then become activated transcription factors, which bind to the pEL222 promoter, activating the production a protein. We model the system of light-induced protein production using the following set of biochemical reactions, describing the constitutive transcription and translation of EL222 molecules, which we denote X and the production of fluorescent protein, Y . The control parameter, $u(t)$, represents a step function that follows from the light either being on or off. We assume a delay τ between the time the light is applied and fluorescent

protein is produced according to the model.



B is a random variable which describes a “burst” of EL222 expression, and is geometrically distributed [21]. This model is valid when mRNA are short-lived compared to protein.

Under this assumption, each individual translation event can be ignored, and replaced with a burst of B protein. The mean burst size b is given by the number of translation events per mRNA, which may be derived from the ratio of the degradation rate of mRNA to the translation rate. B is a geometrically distributed random variable, corresponding to the number of protein made by an mRNA before it was degraded. Using this bursting model of gene expression, we

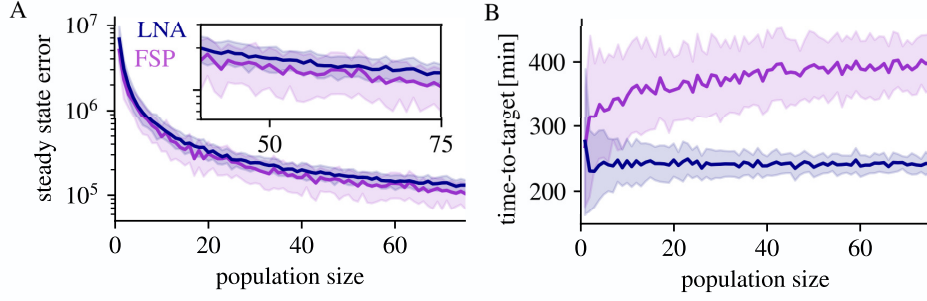


Fig. 3. Population level (A) steady state errors and (B) timing errors of single-cell controls and population control as a function of population size N_{ens} . The blue curves show the error of the population mean at steady state (A) and time-to-target (B) for cells under population control. The purple curves show these same errors for cells under single-cell control. Shaded regions indicate \pm one standard deviation of the errors for different populations of size N_{ens} . The inset in (A) shows steady state errors for population sizes from 40 to 75.

have the following set of propensities for the system,

$$\begin{aligned} w_1 &= k_1 & w_2(x_i) &= \gamma_1 x_i \\ w_3(x_i) &= k_2 x_i & w_4(y_i) &= \gamma_2 y_i, \end{aligned}$$

where x_i and y_i indicate the integer counts of each molecular species. The infinitesimal generator can therefore be written

$$\mathbf{A}_{i,j} = \begin{cases} -\sum_{k=1}^4 w_k(x_i, y_i) & \text{for } i = j \\ w_1 b(1-b)^{k-1} & \text{for } (i, j) \text{ such that } \mathbf{x}_j = \mathbf{x}_i + [k, 0] \\ & \text{for } k \in \mathbb{Z} : k \in \{1, N_x - x_i\} \\ w_2(x_i) & \text{for } (i, j) \text{ such that } \mathbf{x}_j = \mathbf{x}_i + [-1, 0] \\ w_3(x_i, y_i) & \text{for } (i, j) \text{ such that } \mathbf{x}_j = \mathbf{x}_i + [0, 1] \\ w_4(y_i) & \text{for } (i, j) \text{ such that } \mathbf{x}_j = \mathbf{x}_i + [0, -1] \end{cases} \quad (16)$$

We considered both a constant target and a dynamic target. Individual cells were simulated using a modified form of the stochastic simulation algorithm [22] which allowed us to dynamically update the simulation by applying the control $u(t)$. For the population measurement we computed the average fluorescence by taking the molecular abundance of the fluorescent protein s_i from the Gillespie simulation for the i^{th} cell at time k ,

$$y_k = \frac{1}{N_{\text{ens}}} \sum_{i=1}^{N_{\text{ens}}} f_k^i \quad (17)$$

where f_k^i is sampled for each cell from $\mathcal{N}(s_i \mu_{\text{FP}}, s_i \sigma_{\text{FP}}^2)$, and N_{ens} is the number of cells in the simulation. To numerically compare these approaches, we picked a constant target T_1 , Fig. 2A and C, and a dynamic target T_2 , Fig. 2(B)

and (D). We define the population mean error,

$$SSE_{\text{pop}} = \sum_{k=1}^{N_t} (\mu_{\text{FP}} T_k - y_k)^2 \quad (18)$$

and the single-cell error

$$SSE_{\text{sc}} = \frac{1}{N_{\text{ens}}} \sum_{i=1}^{N_{\text{ens}}} \sum_{k=1}^{N_t} (\mu_{\text{FP}} T_k - f_k^i)^2, \quad (19)$$

where the total number of time points is indicated by N_t . In contrast to Eqns. 12-13, the purpose of these errors is to evaluate the quality of the fluorescent measurements realized by each control strategy. We compared the model-based controllers at the population and single-cell level to naive bang-bang controllers. For all controllers, we set $\Delta t = 6$ minutes, and used $H=4$ for the model-based controllers. The naive controller stimulates the cells at $k(\Delta t + 1)$ if the current measurement is below the target fluorescence ($y_k < \mu_{\text{FP}} T_k$ for population, $f_k^i < \mu_{\text{FP}} T_k$ for single-cells), and does not stimulate cells if the current measurement is above the target fluorescence. Population and single-cell errors for the time-varying target in Fig. 2 are reported in Table 1 for 100 cells.

As is apparent in Fig. 2A and C, the population mean under population control reaches the target faster than under single-cell control. On the other hand, after arriving to the target the single-cell control achieves a lower error. We analyzed this trend as a function of population size in Fig. 3. Time-to-target is the first time at which the population mean reaches the target. Steady state error is the mean squared error of the population mean after the target has first been reached. For the LNA population control, 50 independent cell populations were simulated at each population size. For

TABLE I
SUMMARY OF NUMERICAL RESULTS FOR A TIME-VARYING TARGET

	single-cell FSP-MPC	population LNA-MPC	single-cell naive	population naive
single-cell error	5.80×10^6	1.38×10^7	8.24×10^7	7.62×10^7
population error	2.37×10^7	1.11×10^6	7.64×10^7	6.15×10^7

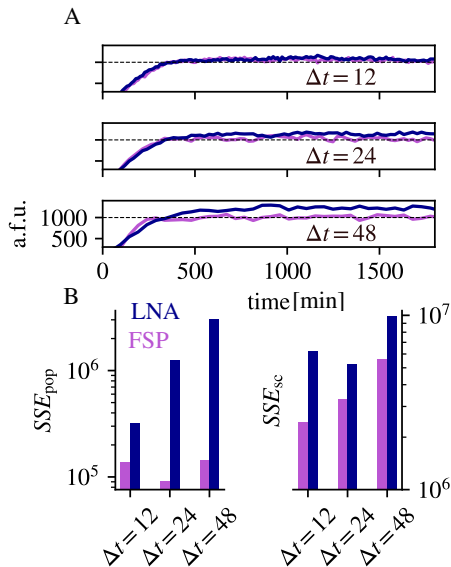


Fig. 4. (A) Population means of single cells controlled with different measurement periods $\Delta t = [12, 24, 48]$ minutes using either the LNA applied to individual cells (blue) or the FSP control (purple). (B) The population errors SSE_{pop} and single-cell errors SSE_{sc} for the LNA control (blue) and FSP control (purple).

each FSP population control, populations at each size were sampled 50 times from a set of 200 independently-controlled cells without replacement to construct standard deviations of the errors for different population sizes. For the steady state error, we show that single-cell control results in marginally lower population errors on average. However, this trend is reversed for time-to-target; across all population sizes, the single-cell control leads to longer time-to-target.

Next, single-cell controller performance as a function of measurement time Δt was quantified. Only FSP-based and simple bang-bang single-cell controllers have been considered thus far; now we evaluate the LNA-based population control (with population size 1) as another single-cell controller against the FSP-based controller. Figure 4 compares the LNA-based single-cell control with the FSP-based one. We considered three measurement times; $\Delta t = [12, 24, 48]$

min. When Δt is small, the LNA and FSP give similar solutions, and therefore a Gaussian approximation leads to accurate control. However, when Δt is larger, mismatch between the Gaussian approximation and the true probability distribution (predicted by the FSP) leads to inaccurate control. Single cell and population errors are shown in Fig. 4B.

Computationally, the single-cell based control is significantly more expensive than the population approach. It requires solving Eq. 4 $2^{H/\Delta t}$ times for each cell within Δt (i.e. between measurement updates), and therefore the main bottleneck is the computation of Eq. 4, which can be efficiently evaluated using Krylov subspace methods [23]. Algorithm 1 scales with $t_{FSP} \times 2^{H/\Delta t} \times N_t$, in which the time to solve the FSP from $k\Delta t$ to $(k+1)\Delta t$ is denoted t_{FSP} . The FSP solution cost mostly depends on the computational expense of an orthogonalization step in the Krylov subspace methods (see [24] for a detailed discussion), and is therefore problem dependent. As long as the time to evaluate over the specified number of horizon steps, $t_{FSP} \times 2^{H/\Delta t}$ is less than the measurement times Δt , single-cell control using the FSP is feasible for at least one cell. Yet, for each cell the computation is “embarrassingly parallel”, and one can scale up the number of cells that are controlled arbitrarily with the number of processors that are available. In the example presented here, the FSP with 2,400 states can be solved in ≈ 0.05 seconds on a single CPU. This amounts to about 3 seconds for each cell with $H = 6\Delta t$. Using realistic experimental time scales of $\Delta t \approx 6$ minutes [12], [25], [5], one could sequentially evaluate control for 120 cells on a single processor.

B. Restoring deterministic behavior for a self-regulated gene

In this example, we start with a simple self regulated gene

$$\emptyset \xrightarrow{f(x,u)} \mathcal{X} \quad (20)$$

$$\mathcal{X} \xrightarrow{\gamma} \emptyset, \quad (21)$$

where the function $f(x)$ is a Hill function which we assume can be affected by the control by modulating the basal

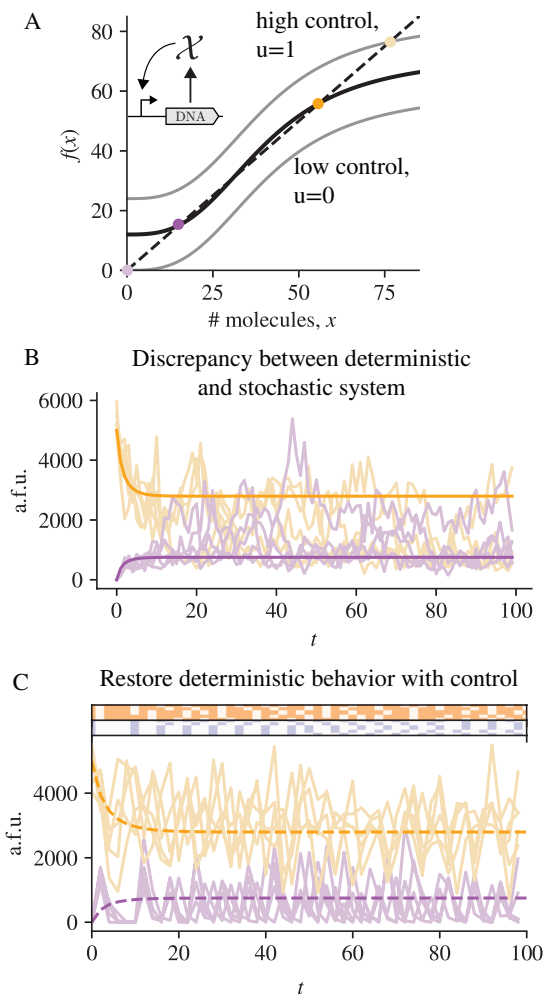


Fig. 5. (A) Hill function (solid line) and degradation (dashed line) for self-regulated gene. The gray lines represent two modes of the controlled system, in which k_b is multiplied by $u(t)$. (B) In the deterministic system (dark orange and purple lines), the system is bistable, and the final state could be either a low or high value depending on the initial condition. However, when simulated with the SSA, state switching is observed and fluorescence switches between the high state and low state (light purple and orange lines). (C) Using the single-cell control, the state switching can be eliminated, restoring the deterministic behavior of the system.

production rate k_b ,

$$f(x, u) = k_b u(t) + \frac{k_f x^n}{K_m^n + x^n}. \quad (22)$$

As shown in Fig. 5, under certain parameter settings, the system is bistable. When the same system is simulated stochastically the trajectories bounce between the two stable equilibria (Fig. 5B), a phenomenon which is not seen in the deterministic setting. We then used the single-cell control

strategy to prevent state switching by supplying the deterministic model behavior (dashed lines, Fig. 5C) as the target for cells that either start with 0 protein molecules (shown in purple) or 100 protein molecules (shown in orange). The controller is effectively able to prevent stochastic switching and stabilize the equilibria of the system, shown in Fig. 5A, by anticipating when cells are likely to pass the unstable region and switch from low to high expression levels.

IV. DISCUSSION

The advent of optogenetic control of single gene expression systems promises precise manipulation of key biological processes, yet many challenges remain [5], [11], [12], [25]. We have compared single-cell and population control (Figs. 2-4). We show that population control tends to have a faster time to target, but has a higher errors for single-cells and the population mean (Fig. 2-3, Table 1). While the improvement of single-cell control is marginal for $\Delta t = 6$ minutes, further numerical experiments between FSP and LNA single-cell controllers in Fig. 4 suggest that model mismatch could lead to larger discrepancies at longer Δt . This indicates that non-Gaussianity plays a key role in model predictive control of cell populations at both the single-cell and population level.

Using an accurate single-cell model can be important for performing single-cell control, particularly if the time between measurements is large and Gaussian approximations are not valid, as shown in Fig. 4B, right panel. This model mismatch can even affect the ability to accurately control the population mean (Fig. 4). While we have limited our work to analyzing how these controllers perform at different Δt , one may also consider changing the prediction horizon, H . The horizon will be target dependent; for example, the optimal H would likely be different for the time-varying target compared to the constant target.

While we have presented FSP-based approaches and demonstrated their feasibility and practical importance, other approaches may also yield precise control. One example is particle filtering, which typically use a sequential importance sampling scheme to draw samples from the posterior state estimate $\tilde{\mathbf{p}}$. For single-cell control, this amounts to managing a set of particles for each cell independently, and could become computationally burdensome to achieve similar accuracy as the FSP approach presented here. However, if accuracy is not a concern, such a sampling scheme may allow one to consider more species than is possible using the approach in this work. It would also be interesting to investigate the use of Gillette's algorithm [22] as a way to generate/propagate the particles through time. Recently, reinforcement learning has showed promise for controlling biological systems for which we do not have accurate models [25].

Future model predictive control approaches could stream real-time measurements to HPC resources, enabling more complex analyses such as model ensembles, online learning, and online experiment design. This could allow one to consider mismatch between the model performing the control and the system under study. Along these same lines, one could devise automated schemes to switch between computationally cheap approaches, such as the LNA, and computationally intense approaches such as the FSP.

V. ACKNOWLEDGMENTS

This study was supported in part by the Center for Nonlinear Studies at Los Alamos National Laboratory, which is operated by Triad National Security, LLC under the auspices of the National Nuclear Security Administration of U.S. Department of Energy under Contract No. 89233218CNA000001. This work was also supported in part by ANR grants CyberCircuits (ANR-18-CE91-0002), MEMIP (ANR-16-CE33-0018), and Cogex (ANR-16-CE12-0025), by the H2020 Fet-Open COSY-BIO grant (grant agreement no. 766840) and by the Inria IPL grant COSY. ZRF thanks Huy Vo and Anatoly Zlotnik for helpful discussions regarding the presentation of this work. This manuscript has been authored by UT-Battelle, LLC under Contract No. DE-AC05-00OR22725 with the U.S. Department of Energy. The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a non-exclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. The Department of Energy will provide public access to these results of federally sponsored research in accordance with the DOE Public Access Plan (<http://energy.gov/downloads/doe-public-access-plan>).

REFERENCES

- [1] J.-B. Lugagne and M. J. Dunlop, "Cell-machine interfaces for characterizing gene regulatory network dynamics," *Current Opinion in Systems Biology*, vol. 14, pp. 1 – 8, 2019, synthetic biology.
- [2] M. Rullan, D. Benzinger, G. W. Schmidt, A. Miliias-Argeitis, and M. Khammash, "An Optogenetic Platform for Real-Time, Single-Cell Interrogation of Stochastic Transcriptional Regulation." *Molecular cell*, vol. 70, no. 4, pp. 745–756.e6, May 2018.
- [3] J. Uhlenendorf, A. Miermont, T. Delaveau, G. Charvin, F. Fages, S. Bottani, G. Batt, and P. Hersen, "Long-term model predictive control of gene expression at the population and single-cell levels," *Proceedings of the National Academy of Sciences*, vol. 109, no. 35, pp. 14 271–14 276, Aug. 2012.
- [4] A. Miliias-Argeitis, M. Rullan, S. K. Aoki, P. Buchmann, and M. Khammash, "Automated optogenetic feedback control for precise and robust regulation of gene expression and cell growth." *Nature communications*, vol. 7, p. 12546, Aug. 2016.
- [5] R. Chait, J. Ruess, T. Bergmiller, G. Tkačik, and C. C. Guet, "Shaping bacterial population behavior through computer-interfaced control of individual cells," *Nature communications*, vol. 8, no. 1, p. 2557, Nov. 2017.
- [6] M. B. Elowitz, A. J. Levine, E. D. Siggia, and P. S. Swain, "Stochastic gene expression in a single cell," *Science*, vol. 297, no. 5584, pp. 1183–1186, 2002.
- [7] J. Melendez, M. Patel, B. L. Oakes, P. Xu, P. Morton, and M. N. McClean, "Real-time optogenetic control of intracellular protein concentration in microbial cell cultures." *Integrative biology : quantitative biosciences from nano to macro*, vol. 6, no. 3, pp. 366–372, Mar. 2014.
- [8] G. Fiore, G. Perrino, M. di Bernardo, and D. di Bernardo, "In Vivo Real-Time Control of Gene Expression: A Comparative Analysis of Feedback Control Strategies in Yeast." *ACS synthetic biology*, vol. 5, no. 2, pp. 154–162, Feb. 2016.
- [9] B. Munsky, G. Li, Z. R. Fox, D. P. Shepherd, and G. Neuert, "Distribution shapes govern the discovery of predictive models for gene regulation." *Proceedings of the National Academy of Sciences*, Jun. 2018.
- [10] L. B. Motta-Mena, A. Reade, M. J. Mallory, S. Glantz, O. D. Weiner, K. W. Lynch, and K. H. Gardner, "An optogenetic gene expression system with rapid activation and deactivation kinetics." *Nature chemical biology*, vol. 10, no. 3, pp. 196–202, Mar. 2014.
- [11] D. Benzinger and M. Khammash, "Pulsatile inputs achieve tunable attenuation of gene expression variability and graded multi-gene regulation," *Nature communications*, vol. 9, no. 1, p. 3521, Aug. 2018.
- [12] Z. R. Fox, S. Fletcher, A. Fraisse, C. Aditya, S. Sosa-Carrillo, J. Petit, S. Gilles, F. Bertaux, J. Ruess, and G. Batt, "Enabling reactive microscopy with micromator," *Nature Communications*, vol. 13, no. 1, p. 2199, 2022.
- [13] N. G. Van Kampen and N. Godfried, *Stochastic processes in physics and chemistry*. Elsevier, 1992.
- [14] D. A. McQuarrie, "Stochastic Approach to Chemical Kinetics," *Journal of Applied Probability*, vol. 4, no. 3, p. 413, Dec. 1967.
- [15] B. Munsky and M. Khammash, "The finite state projection algorithm for the solution of the chemical master equation." *The Journal of Chemical Physics*, vol. 124, no. 4, p. 044104, Jan. 2006.
- [16] H. Vo and B. Munsky, "A parallel implementation of the finite state projection algorithm for the solution of the chemical master equation," *bioRxiv*, 2020.
- [17] S. Säräkk, *Bayesian Filtering and Smoothing*, ser. Institute of Mathematical Statistics Textbooks. Cambridge University Press, 2013.
- [18] M. Komorowski, M. J. Costa, D. A. Rand, and M. P. H. Stumpf, "Sensitivity, robustness, and identifiability in stochastic chemical kinetics models." *Proceedings of the National Academy of Sciences of the United States of America*, vol. 108, no. 21, pp. 8645–8650, May 2011.
- [19] O. Poleskaya, A. Baranova, S. Bui, N. Kondratev, E. Kananykhina, O. Nazarenko, T. Shapiro, F. B. Nardia, V. Kornienko, V. Chandhoke, I. Stadler, R. Lanzafame, and M. Myakishev-Rempel, "Optogenetic regulation of transcription," *BMC Neuroscience*, vol. 19, no. Suppl 1, Apr. 2018. [Online]. Available: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5998900/>
- [20] Castillo-Hair, Sebastian M, Baerman, Elliot A, Fujita, Masaya, Igoshin, Oleg A, and Tabor, Jeffrey J, "Optogenetic control of Bacillus subtilis gene expression." *Nature communications*, vol. 10, no. 1, pp. 3099–11, Jul. 2019.
- [21] V. Shahrezaei and P. S. Swain, "Analytical distributions for stochastic gene expression," *Proceedings of the National Academy of Sciences*, vol. 105, no. 45, pp. 17 256–17 261, 2008.
- [22] D. T. Gillespie, "Exact stochastic simulation of coupled chemical reactions," *The Journal of Physical Chemistry*, vol. 81, no. 25, pp. 2340–2361, Dec. 1977.
- [23] R. B. Sidje, "Expokit: a software package for computing matrix exponentials," *ACM Transactions on Mathematical Software (TOMS)*, vol. 24, no. 1, pp. 130–156, 1998.

- [24] H. D. Vo and R. B. Sidje, "Approximating the large sparse matrix exponential using incomplete orthogonalization and krylov subspaces of variable dimension," *Numerical linear algebra with applications*, vol. 24, no. 3, p. e2090, 2017.
- [25] J.-B. Lugagne, C. M. Blassick, and M. Dunlop, "Deep model predictive control of gene expression in thousands of single cells," *bioRxiv*, pp. 2022–10, 2022.