

A divide-and-conquer method for analyzing high-dimensional noisy gene expression networks

Zhou Fang^{a,1}, Ankit Gupta^{a,1}, Sant Kumar^a, and Mustafa Khammash^{a,*}

^aDepartment of Biosystems Science and Engineering, ETH Zurich, CH-4058 Basel, Switzerland

This manuscript was compiled on October 6, 2023

Intracellular gene expression systems are inevitably random due to low molecular counts. Consequently, mechanistic models for gene expression should be stochastic, and central to the analysis and inference of such models is solving the Chemical Master Equation (CME), which characterizes the probability evolution of the randomly evolving copy-numbers of the reacting species. While conventional methods such as Monte-Carlo simulations and finite state projections exist for estimating CME solutions, they suffer from the curse of dimensionality, significantly decreasing their efficacy for high-dimensional systems. Here, we propose a new computational method that resolves this issue through a novel divide-and-conquer approach. Our method divides the system into a *leader* system and *several* conditionally independent *follower* subsystems. The solution of the CME is then constructed by combining Monte Carlo estimation for the leader system with stochastic filtering procedures for the follower subsystems. We develop an optimized system decomposition, which ensures the low-dimensionality of the sub-problems, thereby allowing for improved scalability with increasing system dimension. The efficiency and accuracy of the method are demonstrated through several biologically relevant examples in high-dimensional estimation and inference problems. We demonstrate that our method can successfully identify a yeast transcription system at the single-cell resolution, leveraging mRNA time-course microscopy data, allowing us to rigorously examine the heterogeneity in rate parameters among isogenic cells cultured under identical conditions. Furthermore, we validate this finding using a novel noise decomposition technique introduced in this study. This technique exploits experimental time-course data to quantify intrinsic and extrinsic noise components, without requiring supplementary components, such as dual-reporter systems.

stochastic reaction systems | chemical master equation | stochastic filtering | modularization | Rao-Blackwellization

Advances in modern biological technology (e.g., flow cytometry (1), time-lapse microscopy (2, 3), and fluorescence proteins (4)) have substantially improved scientists' ability to investigate living biological cells. It is now well-established that the dynamics within cells are intrinsically noisy, with significant cell-to-cell heterogeneity. For example, it is known that gene expression systems typically possess a high degree of randomness due to the low molecular counts of the involved biomolecular species (5–7). To appropriately describe such systems, stochastic continuous-time Markov chain models (8) are often employed, where each node in the chain corresponds to a specific cellular state, and the transitions between these nodes correspond to the reactions happening in the cell.

To calibrate these stochastic models, single-cell measurement data is required. One such experimental technique is Flow Cytometry which records the dynamical evolution of a heterogeneous cell population over time, providing valuable information for inferring intracellular dynamics (9). Time-lapse microscopy is another technology that can enable such an infer-

ence task by providing time-course data of each individual cell (10–12). Such data is richer in information than population-level data, as the same cells are tracked over time, preserving temporal correlations between their dynamic states. This facilitates inference of parameters, *localized* to each tracked cell (12), thereby affording a more nuanced understanding of cellular dynamics. A typical issue that arises is that only a few state variables can be dynamically tracked. However, given these measurement trajectories and a stochastic model for the dynamics, the conditional probability distribution of the unobserved state-variables can be estimated by solving what is known as the *stochastic filtering problem* (13). Estimation of this time-varying conditional probability distribution provides valuable insight into the dynamical behaviors of the intracellular systems and enables better feedback control design for these systems (14).

Central to stochastic modeling and analysis of biological systems is the mathematical problem of solving chemical master equations (CMEs), which consist of a collection of linear ordinary differential equations (ODEs) that determine the time-evolution of the probability distribution of the random state-vector comprising of molecular counts of all the species. The CME plays a key role in investigating the effects of noise on biological processes, such as chemotaxis (15) and cell cycle (16, 17). Solutions of the CME can also be applied to the rational design of systems, such as the genetic repressilator (18) and the antithetic integral feedback controller (19). In the problem of stochastic model inference from single-cell measurement data, several instances of the CME need to be solved to identify the model that best fits the data (20). The filtering problem associated with single-cell time-course data is often solved recursively by prediction steps and correction steps (21). The prediction steps rely heavily on solving the CME for computing predicted probability distributions (21).

All these applications, and many others, require efficient and accurate methods for solving CMEs. However, estimating CME solutions is extremely challenging, **as the number of ODEs that form the CME is usually very large, and in fact in most cases of interest, it is infinite.** Exact solutions of CMEs exist in very special cases (22–25), and in general one has to resort to numerical approaches to estimate CME solutions. **Even though a significant amount of research effort has been devoted towards developing efficient CME solvers, the current methods are ill-equipped for solving CMEs for high-dimensional systems.** This has been a major bottleneck

Author contributions: Z.F., A.G., and M.K. designed research; Z.F. and A.G. performed research; S.K. conducted the biological experiment; Z.F., A.G., and M.K. wrote the paper.

The authors declare no competing interest.

¹These authors contributed equally: Zhou Fang and Ankit Gupta.

*To whom correspondence should be addressed. E-mail: mustafa.khammash@bsse.ethz.ch

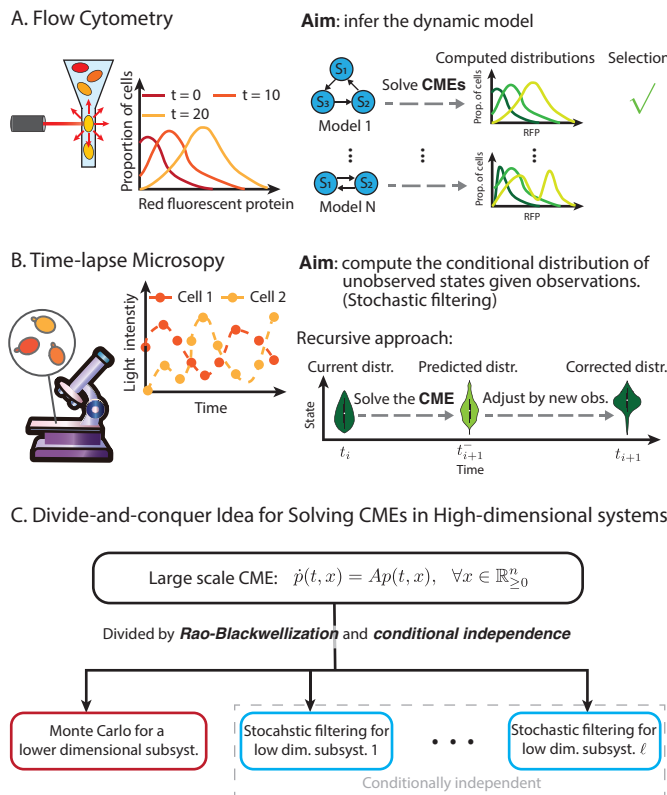


Fig. 1. At the heart of stochastic modeling and analysis of biological systems is the mathematical problem of solving chemical master equations (CMEs). (A) Inference problem associated with single-cell measurement data. To solve this inference problem, researchers need to solve CMEs to obtain single-cell distributions for many candidate models and then select the model based on the computed distributions. (B) Stochastic filtering based on single-cell time-course data. The aim is to compute the conditional probability distribution of unobserved species within a cell using the cell's time-course data. The problem is often solved in a recursive manner involving prediction steps and correction steps. The prediction steps need to solve CMEs to obtain the predicted probability distributions. (C) Our divide-and-conquer idea for solving large-scale CMEs. Our method utilizes Rao-Blackwellization and stochastic filtering (for conditional independence) to divide the original CME into several manageable sub-problems for low dimensional subsystems.

in the widespread adoption of stochastic models in biology.

The most commonly used methods for solving CMEs are the Monte-Carlo methods and the direct approach involving state-space truncations that reduce the CME system to a finite tractable system of ODEs. Monte-Carlo methods simulate multiple trajectories using algorithms such as the Gillespie algorithm (26, 27), tau-leaping (28–30), or hybrid models (31, 32). The empirical distribution obtained from these simulations is used to approximate the exact probability distribution. This method can be computationally efficient, but its accuracy decreases as the dimension of the system increases, making it unsuitable for high-dimensional settings. In contrast, direct approaches involving truncations, such as the finite state projection (FSP) approach (33–35), are very accurate, as they solve a finite-dimensional approximation of the CME directly and they provide a computable error-bound. However, since the size of the required truncated state-space scales exponentially with the system dimension, the FSP is computationally infeasible in large dimensional settings.

In some situations, CMEs for high-dimensional systems can be effectively approximated by some parametric methods. For

example, when the system size is large, the linear noise approximation (36) approximates the exact probability distribution by a Gaussian distribution whose mean and variance are the parameters that need to be computed based on the underlying system. Similarly, the moment closure method closes the moment equations based on a chosen family of distributions and then tracks the first few moments (37–44). More recently, some deep learning approaches (45–47) were established for solving CMEs thanks to the universal approximation properties of neural networks. These parametric methods are quite successful in many applications. However, their validity largely depends on the suitable choice of the family of distributions, which is not trivial for generic systems.

When the system is multiscale, i.e. it has reactions firing at different time-scales this curse of dimensionality can be somewhat mitigated by the method of timescale separation (e.g., (32, 48–50)), where a reduced order CME system can be derived by applying the quasi-stationary assumption. Though this method is not always applicable, it is effective for analyzing many biological processes, provided they are multiscale in nature.

Dimension reduction ideas have also been applied to some hybrid methods for solving CMEs. The method of conditional moments (51) derives a reduced order CME for selected species by integrating out the moments of other species conditioned on the state of these selected ones using moment closure methods. Similarly, the uncoupled simulation method (52, 53) generates Monte-Carlo samples only for several selected species by integrating out the effect of other species using moment closure methods. These hybrid approaches (consisting of the moment closure method and some others) are computationally efficient and accurate in estimating the selected species (51–53). However, since these approaches use the moment closure approach to approximate the conditional distributions of non-selected species, their accuracy and reliability are not always guaranteed, especially when Hill-type kinetics are involved. In addition, converting approximated moments into distributions to estimate non-selected species incurs another source of error. Finally, a systematic decomposition of the reacting system for optimal performance of a hybrid approach is yet to be determined.

Overall, there is still a paucity of methods in the current literature to effectively deal with high-dimensional stochastic reaction systems. Would not it be advantageous if we could divide a high dimensional system into smaller pieces, thereby mitigating the curse of dimensionality? Motivated by this idea, we propose in this paper a modularization-based method for solving CMEs.

Our approach is inspired by a divide-and-conquer strategy enabled by probabilistic independence. Specifically, when all the species are probabilistically independent, the joint probability distribution can be derived by forming a product of all the marginal distributions. In this case, instead of directly solving the CME for the large dimensional system, we can first compute the marginal distributions by solving the reduced order CME for each species and then combine them to obtain the solution of the original CME. Suppose such a system has n -species with each species up to $T - 1$ copies, and all the CMEs are solved by the FSP. Then, this divide-and-conquer strategy reduces the computational complexity from $O(T^{3n})$ to $O(nT^3)$, which can result in significant improvement when the

system dimension (n) is large. Moreover, this strategies enables parallel computation for the marginal distributions, which could further accelerate the computational process. Similar levels of complexity reduction can also be achieved with other computational methods, e.g., Monte Carlo.

Of course, the independence of species would generally not hold as the species are interacting. However, we shall adopt a new modularization strategy that exploits conditional independence between certain parts of the system, given the trajectories of intermediate species acting like conduits for inter-modular interactions. This new approach, combined with filtering (for conditional independence), makes modularization applicable to general high-dimensional networks and can significantly reduce the computational effort for solving the CME. Following this idea, our method divides the whole chemical reacting system into a leader system and several conditionally independent follower subsystems in a principled manner (see Figure 2 for a graphic illustration). Given the decomposition, we solve the CME by employing the Rao-Blackwellization technique (54, 55), which applies the Monte-Carlo method to the leader system and a proper filtering approach to computing the conditional distributions of the follower subsystems (see Figure 2.C). More importantly, the system decomposition involved in our method is well-designed so that the leader system and all the follower subsystems are low dimensional (Figure 2.B). Consequently, our approach breaks down the original large scale problem into several manageable sub-problems for low-dimensional subsystems, resulting in reduced overall computational cost.

Our approach combines aspects of both the Monte-Carlo method and the stochastic filtering approach, making it a hybrid method. In this sense, our optimized system-decomposition algorithm offers a way to balance the strengths of both methods for improved results. At one extreme, when all species are treated as leader-level species, our approach reduces to the conventional Monte-Carlo method. At the other extreme, when all species are treated as follower-level species, it becomes equivalent to applying the chosen filtering approach for solving the CME. In this paper, we specifically employ the filtered FSP method (56) for solving the filtering problems associated with follower subsystems. However, it is not the only choice; one can also choose other suitable filtering algorithms. Particularly, when the Moment closure approach is applied for the filtering sub-problems, our method boils down to the method in (51–53) but with a principled and far more optimal system decomposition.

In this paper, we demonstrate the efficacy of our approach both computationally and experimentally. First, we consider several biologically relevant *in silico* models and demonstrate the superior performance of our method in solving both the CMEs and the associated stochastic filtering problems. We then further develop our method and show how it can successfully leverage experimental time-course data for identifying a yeast transcription model at the single-cell resolution.

This analysis illuminated the significant heterogeneities in rate parameters, even among identically cultured isogenic cells. These parameter heterogeneities can be viewed as the “extrinsic” component of the overall cell-to-cell heterogeneity (57), while the “intrinsic” component is generated by the randomness in the firing of intracellular reactions. Decomposition of total cellular homogeneity into extrinsic and intrinsic

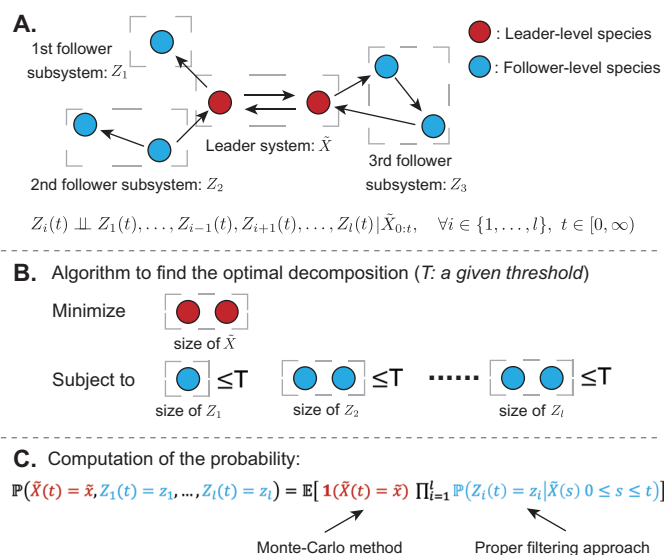


Fig. 2. Illustration of our approach. (A) System decomposition. Our approach works by a system decomposition, which first divides a chemical reaction system into a leader system (red species) and a follower system (blue species). Then, the follower system is further decomposed into several subsystems according to some topological conditions such that the follower subsystems are conditionally independent given the trajectory of the leader system. This conditional independence is denoted by $\perp\!\!\!\perp$. (B) Algorithm to find the optimal decomposition. Our algorithm chooses the optimal decomposition that minimizes the size of the leader system (or maximizes the size of the whole follower system) while keeping the size of each follower subsystem below a given threshold. By this algorithm, all the subsystems are low dimensional. (C) The computation of the probability distribution. Our approach solves the CME by applying the Monte-Carlo method to the leader system and a filtering approach (e.g., filtered FSP (56)) to each follower subsystem. Since each subsystem is low dimensional by our system decomposition algorithm, our computational approach scales more favorably with the system dimension.

components is an active research problem (58–61), and it is fundamentally important not just for deciphering the source of variability, but also for understanding how noise affects intracellular signal processing and control (62, 63). Traditionally, the decomposition of noise in cellular heterogeneity has been conducted using experimental Flow Cytometry data in conjunction with dual reporter assays (58, 60, 61). However, the synthetic implementation of these assays presents numerous challenges. These include the need to ensure the statistical equivalence of the dual reporters and to confirm the conditional independence of the generated estimates, given the input variables.

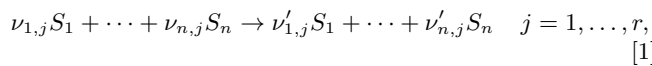
To validate the estimates of the intrinsic and extrinsic noise from the inferred transcription models, this paper introduces a novel method for noise decomposition in gene expression from experimental data that does not rely on dual-reporter systems. This method leverages the information-richness of time-course data and capitalizes on the inherent stability of the stochastic gene-expression network. We employed this method to analyze single-cell time-course data for the yeast transcription model and compared it with the noise decomposition result generated by the inferred transcription models. The results exhibited a high degree of concordance between the two sets of estimates, thereby confirming the efficacy of our identification method.

Table 1. Terminology Table. The terms highlighted in bold are new concepts developed in this paper.

Terminology	Meaning
CME:	the chemical master equation Eq. [3].
Filtered CME:	the filtering equation Eq. [5] characterizing the conditional distributions of follower subsystems.
FSP:	the finite state projection method (33).
Filtered FSP:	the FSP method for the filtered CME (56).
MC:	the Monte-Carlo method
PF:	the particle filter, applying MC to filtering problems.
RB-CME solver:	the divide-and-conquer method we proposed for solving CMEs by exploiting Rao-Blackwellization and conditional independence. The method combines MC and a filtering approach, determined by the user's preference. In the following examples, we selected the filtered FSP as our choice of the filtering approach.
RB-PF:	the Rao-Blackwellized particle filter we developed, which applies the RB-CME solver to the prediction step. Again, in the following examples, we used the filtered FSP as the chosen filtering approach within the RB-CME solver.

1. Method and theoretical results

A. Chemical master equations and the curse of dimensionality. We consider an intracellular reaction system that has r reactions,



where S_1, \dots, S_n are n different chemical species, and $\nu_{i,j}$ and $\nu'_{i,j}$ are the stoichiometric coefficients. Due to the low molecular counts, this chemical reaction system is often modeled by a continuous time Markov chain (8)

$$X(t) = X(0) + \sum_{j=1}^r \zeta_j R_j \left(\int_0^t \lambda_j(X(s)) ds \right) \quad [2]$$

where $X(t)$ is an n -dimensional vector representing the molecular count of each species at time t , the vector ζ_j equals to $(\nu'_{1,j} - \nu_{1,j}, \dots, \nu'_{n,j} - \nu_{n,j})^\top$ indicating the state change after a firing of the j -th reaction, $R_j(t)$ are independent unit rate Poisson processes, and $\lambda_j(\cdot)$ are the propensities indicating the rates of these reactions. In this paper, we consider several mild technical assumptions for this dynamical system (see [SI Appendix, section S1](#)) so that the process $X(t)$ is regular enough. Under these assumptions, the probability distribution of the system Eq. [2] is characterized by the chemical master equation (CME) (8)

$$\frac{dp(t, x)}{dt} = \sum_{j=1}^r \lambda_j(x - \zeta_j) p(t, x - \zeta_j) - \sum_{j=1}^r \lambda_j(x) p(t, x), \quad [3]$$

where $x \in \mathbb{Z}_{\geq 0}^n$ is the value of the state, and $p(t, x) \triangleq \mathbb{P}(X(t) = x)$ is the probability at x . By solving the CME, scientists can gain many insights into the considered biological processes.

Usually, CMEs are difficult to solve explicitly. Conventional methods to numerically solve the CME include the simulation-based Monte-Carlo methods and the finite state projection (FSP) method. However, both of them scale poorly with the size of the state space and the system dimension.

A simulation-based Monte Carlo method first simulates N trajectories of the system Eq. [2], denoted by $x_1(t), \dots, x_N(t)$, and then uses the empirical distribution $p_{MC}(t, x) \triangleq \frac{1}{N} \sum_{j=1}^N \mathbb{1}(x_1(t) = x)^*$ to approximate the exact probability distribution. The error of this method can be evaluated by the L_1 distance between $p_{MC}(t, \cdot)$ and $p(t, \cdot)$, which upper bounds the largest possible error of the Monte Carlo in estimating any particular probability. Mathematically, this error converges at the rate of $1/\sqrt{N}$, and its pre-convergence rate factor (defined by $\lim_{N \rightarrow \infty} \sqrt{N} \mathbb{E} [\| \hat{p}_{MC}(t, \cdot) - p(t, \cdot) \|_1]$) lies in a particular range shown as follows ([SI Appendix, section S2.A](#)):

$$\begin{aligned} & \lim_{N \rightarrow \infty} \sqrt{N} \mathbb{E} [\| \hat{p}_{MC}(t, \cdot) - p(t, \cdot) \|_1] \\ &= \sqrt{\frac{2}{\pi}} \left[\left(\sum_{x \in \mathbb{Z}_{\geq 0}^n} \sqrt{p(t, x)} \right) \pm 1 \right]. \end{aligned} \quad [4]$$

This suggests that given a large sample size N , the error of the Monte Carlo largely depends on the value of $\sum_{x \in \mathbb{Z}_{\geq 0}^n} \sqrt{p(t, x)}$. This value tends to scale poorly with the size of the state space containing most of the probability mass. Particularly, when the probability mass is uniformly distributed on T states, this quantity equals \sqrt{T} , which can be very large when T is big. Also, since this state size often grows exponentially with the number of species, this quantity $\sum_{x \in \mathbb{Z}_{\geq 0}^n} \sqrt{p(t, x)}$ tends to scale poorly with the system dimension n . When the molecular counts of different chemical species are independent, this sum does grow exponentially with n , as its value equals to $\prod_{i=1}^n \left(\sum_{x_i \in \mathbb{Z}_{\geq 0}} \sqrt{p_i(t, x_i)} \right)$ with $p_i(t, \cdot)$ the marginal distribution for the i -th species. In summary, the error of a Monte-Carlo method tends to scale unfavorably with the size of state space and system dimension, and, thus, this method usually performs poorly in high-dimensional problems.

In contrast, the FSP approach (33), which directly solves a truncated CME on a large but finite state space, is very accurate for arbitrary reaction systems. However, the computational complexity of this method scales cubically with the size of the truncated state space. Moreover, since the size of the finite state space scales exponentially with n , the FSP is computationally demanding for high-dimensional problems.

In conclusion, both of the methods above suffer the curse of dimensionality.

B. Motivating example: beat the curse of dimensionality through probabilistic independence. Despite the challenges mentioned above, the CME can be efficiently solved for large dimensional systems when all the states in $X(t)$ are probabilistically independent. In such cases, the probability distribution can be obtained by combining all the marginal distributions. Therefore, the large-scale CME can be divided into several manageable sub-problems, each of which only solves a CME for obtaining one marginal distribution.

To further elaborate this divide-and-conquer approach, we consider a toy example where all the species are independent, each with no more than $T - 1$ copies, and their marginal probability mass is uniformly distributed in the state space. When straightforwardly applying the FSP to such a system, we encounter a computational complexity of $O(T^{3n})$, as discussed

* $\mathbb{1}(\cdot)$ is the indicator function whose value equals to 1 if its argument is true, otherwise 0.

in Section A. In contrast, the divide-and-conquer strategy reduces the complexity to $O(nT^3)$, which is much more favorable when both n and T are large. In addition to this complexity reduction, this divide-and-conquer strategy also enables parallel computation for the marginal distributions, which could further reduce the computational time. Such an improvement also occurs when the Monte-Carlo method is employed. According to Eq. [4], the classical Monte-Carlo method needs to generate $\frac{T^n}{\epsilon^2}$ samples to reach an accuracy of ϵ . In contrast, each marginal distribution only needs $\frac{T}{\epsilon^2}$ samples to reach this accuracy. Since the errors in estimating these marginal distributions additively contribute to the error in the joint distribution estimate (SI Appendix, section S2.A), the divide-and-conquer strategy only needs $\frac{n^3 T}{\epsilon^2}$ samples to attain the accuracy of ϵ . This figure is significantly less than the sample size ($\frac{T^n}{\epsilon^2}$) required by the conventional method. Overall, this modularization method powered by probabilistic independence can significantly reduce the required computational complexity and beat the curse of dimensionality.

Nevertheless, species are generally not probabilistically independent as they are interacting in most cases. In what follows, we propose a new modularization approach exploiting conditional independence together with Rao-Blackwellization. This new strategy, combined with filtering (for conditional independence), makes modularization applicable to general high-dimensional networks and can also significantly reduce the computational effort for solving CMEs.

C. Modularization-based Rao-Blackwell method for solving CMEs. We propose a new modularization method for solving CMEs by exploiting Rao-Blackwellization and conditional independence. Specifically, our method first uses an automated algorithm to decompose the system into two parts: $\tilde{X}(t)$ (denoted as the leader system) and $Z(t)$ (termed as the follower system) (see Figure 2). The detailed decomposition strategy will be discussed later in this subsection. Based on this leader-follower decomposition, we also divide reaction vectors ζ_j into $\zeta_j^{\tilde{X}}$ and ζ_j^Z , where $\zeta_j^{\tilde{X}}$ indicates the state change of the leader system, and ζ_j^Z indicates the state change of the follower system. Moreover, we further decompose the follower system $Z(t)$ into several lower-dimensional subsystems, $Z_1(t), \dots, Z_l(t)$, such that the following topological conditions are satisfied.

- C1 Each reaction involves a maximum of one follower subsystem (meaning that at most one follower subsystem can influence the reaction's propensity or have its state altered by the reaction).
- C2 The reactions with the same non-zero $\zeta_j^{\tilde{X}}$ involve a maximum of one follower subsystem (meaning that at most one follower subsystem can influence the propensities of these reactions or have its state altered by these reactions).

For the ease of notations, we rearrange the order of species such that $X(t) = (\tilde{X}(t), Z_1(t), \dots, Z_l(t))$; also, for every state x , we write $x = (\tilde{x}, z_1, \dots, z_l)$, where \tilde{x} is the state of the leader system, and z_i ($i = 1, \dots, l$) is the state of the follower system. Under the conditions above, the follower subsystems are conditionally independent given the trajectory of the leader system (see the proof in SI Appendix, Section S3.A). Moreover, the conditional probability distribution $\pi_{Z_i|\tilde{X}}(t, z_i) \triangleq \mathbb{P}(Z_i(t) = z_i | \tilde{X}(s), 0 \leq s \leq t)$ is characterized

by a set of differential equations with jumps (heuristically derived in (53, 64) and rigorously verified in (56)):

$$\begin{aligned} d\pi_{Z_i|\tilde{X}}(t, z_i) &= f_1(\tilde{X}(t), z_i, \pi_{Z_i|\tilde{X}}(t, \cdot)) dt + f_2(\tilde{X}(t), z_i, \pi_{Z_i|\tilde{X}}(t, \cdot)) dt \\ &+ \sum_{j=1}^n \mathbb{1}(\tilde{X}(t) - \tilde{X}(t^-) = \zeta_j^{\tilde{X}}) g_j(\tilde{X}(t^-), z_i, \pi_{Z_i|\tilde{X}}(t^-, \cdot)) \end{aligned} \quad [5]$$

whose detailed expression is given in SI Appendix, Section S3.A. Here, f_1 represents the prediction of the conditional distribution based on the dynamical model, and the remaining terms correspond to the corrections to the estimates in accordance with the dynamics of the observable species. In this paper, we call Eq. [5] the filtered CME because it computes the conditional distribution, which is a solution of a filtering problem. A schematic illustration of the decomposition introduced above is presented in Figure 2.A. Here, we should note that this conditional independence holds only when the entire trajectory of the leader species is given. If only the current state of the leader species is provided, this conditional independence may not necessarily hold.

By this decomposition and the law of total expectation, we can rewrite the probability distribution by $p(t, \tilde{x}, z_1, \dots, z_l) = \mathbb{E}[\mathbb{1}(\tilde{X}(t) = \tilde{x}) \prod_{i=1}^l \pi_{Z_i|\tilde{X}}(z_i)]$. Based on it, we design an Rao-Blackwell method for CMEs, which applies Monte Carlo to the leader system and a filtering approach to each follower subsystems. Concretely, we first generate N simulations of the system Eq. [2] and denote their leader-system parts by $\tilde{x}_1(t), \dots, \tilde{x}_N(t)$, respectively. Then, for each simulated trajectory and each follower subsystem, we calculate the conditional probability distribution $q_j^i(t, z_i) \triangleq \mathbb{P}(Z_i(t) = z_i | \tilde{X}(s) = \tilde{x}_j(s), 0 \leq s \leq t)$ using a stochastic filtering approach. Users have the flexibility to choose suitable filtering approaches (e.g., Monte-Carlo method (64), filtered FSP (56) etc.) for computing these conditional distributions. In later examples, we specifically employ the filtered FSP (56), which solves Eq. [5] directly on a large but finite state space. Finally, the exact probability distribution is approximated by the quantity

$$\hat{p}_{\text{RB}}(t, x) = \frac{1}{N} \sum_{j=1}^N \left[\mathbb{1}(\tilde{x}_j(t) = \tilde{x}) \prod_{i=1}^l q_j^i(t, z_i) \right]. \quad [6]$$

We name this algorithm the Rao-Blackwellized CME solver (RB-CME solver). Mainly, this Rao-Blackwellized method transforms the original problem of solving the CME into several potentially low dimensional subproblems. Therefore, our method tends to scale more favorably with the system dimension.

The RB-CME solver can be seen as a principled way to combine the Monte-Carlo method and the chosen filtering approach. Particularly, when all the species are classified as leader-level species, this method becomes a Monte-Carlo method. When all the species are classified as follower-level species, the filtered CME becomes the CME, and the RB-CME solver is equivalent to applying the chosen filtering approach to the CME. By classifying chemical species differently, one can choose whether this method is similar to the Monte-Carlo method or to the chosen filtering approach. When the filtered FSP is applied to the follower subsystems, our method is more

computationally practical than the original FSP approach in high dimensional cases, as our method applies the filtered FSP to each follower subsystem separately.

Given the same samples size N , the RB-CME solver is no less accurate than the conventional Monte-Carlo method, if all the conditional distributions $q_j^i(t, z_i)$ are computed precisely. Basically, the first layer of the RB-CME solver is a Monte Carlo approach[†]; therefore, its L_1 error has a convergence rate \sqrt{N} and its pre-convergence rate factor depends on the variance of the random variable it generates. The expression of this pre-convergence factor is given in [SI Appendix, Section S3.C](#). Note that the variance of $\mathbb{1}(\tilde{X}(t) = \tilde{x}) \prod_{i=1}^l \pi_{Z_i|\tilde{X}}(z_i)$ (used for the RB-CME solver) is no greater than that of $\mathbb{1}(X(t) = x)$ (used for the conventional Monte-Carlo method) due to the law of total variance. So, we can conclude the superior performance of our method in the sense of the pre-convergence rate factor, i.e., (see [SI Appendix, Section S3.C](#))

$$\begin{aligned} & \lim_{N \rightarrow \infty} \sqrt{N} \mathbb{E} [\|\hat{p}_{\text{RB}}(t, \cdot) - p(t, \cdot)\|_1] \\ & \leq \lim_{N \rightarrow \infty} \sqrt{N} \mathbb{E} [\|\hat{p}_{\text{MC}}(t, \cdot) - p(t, \cdot)\|_1] \end{aligned}$$

Particularly, if the conditional probability distribution $\pi_{Z_i|\tilde{X}}(z_i)$ only depends on the final state of the leader system $\tilde{X}(t)$ and is irrelevant to any of its historical information, then the RB-CME solver is equivalent to the time scale separation method, which eliminates all dynamics of the follower-level species using quasi-stationary assumption. In this case, our method's error equals to the Monte-Carlo method's error for estimating $\tilde{X}(t)$ solely [SI Appendix, Section S3.C](#), and its pre-convergence rate factor is far less than that in Eq. [4] when the dimension of $\tilde{X}(t)$ is much less than that of $X(t)$. This indicates that the RB-CME solvers can be far more accurate than the Monte-Carlo method in general high-dimensional problems.

Despite this error reduction, given the same sample size N , the RB-CME solver is more time-consuming and memory-demanding than the Monte-Carlo method, as our method needs to additionally compute and store the marginal distributions $q_j^i(t, z_i)$. Then, the question is when our method has superior performance under the same time cost or at the same accuracy level. This proposes a challenging mathematical problem. Specifically, this requires an exact computation of the variance of $\mathbb{1}(\tilde{X}(t) = \tilde{x}) \prod_{i=1}^l \pi_{Z_i|\tilde{X}}(z_i)$, which is case-dependent and hard (if not impossible) to obtain in a general explicit form. Instead of pursuing a universal solution to this problem, we aim to obtain insights through a case study of the system considered in Section 1.B. Recall that for this n -dimensional system, the Monte-Carlo method needs $\frac{T^n}{\epsilon^2}$ samples (equivalent to a complexity of $O(\frac{T^n}{\epsilon^2})$) to reach an accuracy of ϵ , where T is the state space size for each species. Now, we apply our RB-CME solver to this problem by classifying n_1 species as leader species and the remaining as follower-level species. We also assume that these follower-level species can be divided into $\frac{n-n_1}{n_2}$ follower subsystems, each containing n_2 species, and the error of the RB-CME solver equals the error in estimating the leader species. In this scenario, the RB-CME solver needs $\frac{T^{n_1}}{\epsilon^2}$ samples to reach an accuracy of ϵ . Moreover, when we apply the filtered FSP to compute the filtered CME for each follower subsystem, the RB-CME solver has a complexity of

$O(\frac{T^{n_1}}{\epsilon^2} \times \frac{n-n_1}{n_2} T^{3n_2}) = O(\frac{n-n_1}{n_2 \epsilon^2} T^{n_1+3n_2})$. We can observe that this complexity is much reduced compared with that of the Monte-Carlo method $O(\frac{T^n}{\epsilon^2})$, when the system dimension n is large, and n_1 and n_2 are small. In conclusion, the RB-CME solver has better performance if the leader system and all the follower subsystems have small state spaces.

As demonstrated in the above example, system decomposition has an important role in determining the performance of the RB-CME solver. Here, we present a principled method for system decomposition aimed at maximizing the efficiency of our approach. For each follower subsystem, we define SS_i as the size of its truncated state space containing the most probability mass. Intuitively, the RB-CME solver is more accurate and efficient when the leader system and all the follower subsystems have small sizes in (truncated) state spaces. Mathematically, this suggests that $\prod_{i=1}^l \text{SS}_i$ (the size of the whole follower system) should be large, while $\max_i \text{SS}_i$ (the size of the largest follower subsystem) should be small. In most cases, these two optimization objectives cannot not be achieved simultaneously. Consequently, we choose a leader-follower decomposition (among all) that maximizes the size of the whole follower system ($\prod_{i=1}^l \text{SS}_i$) while keeping the size of each individual follower subsystem (SS_i) below a given threshold. The detailed algorithm following this strategy is provided in [SI Appendix, Section S3.D](#). With this decomposition algorithm, we effectively divide the problem of solving CMEs into several lower-scale computational sub-problems, and, consequently, our approach scales more favorably with the system dimension in terms of accuracy and efficiency.

Finally, we summarize the procedure of the RB-CME solver as follows (also see a schematic illustration in Figure 2).

1. Decompose the system into a leader system and several follower subsystems using the algorithm in [SI Appendix, Section S3.D](#).
2. Generate N simulations of the leader system, and denote them by $\tilde{x}_1(t), \dots, \tilde{x}_N(t)$.
3. For each trajectory $\tilde{x}_j(\cdot)$ and each subsystem $Z_i(t)$, solve $q_j^i(t, z_i) \triangleq \mathbb{P}(Z_i(t) = z_i | \tilde{X}(s) = \tilde{x}_j(s), 0 \leq s \leq t)$ by applying the filtered FSP to the filtered CME Eq. [5].
4. Compute the result by Eq. [6].

D. Rao-Blackwell method for solving the filtering problem.

With modern microscope, scientists can measure some signals (e.g., fluorescent reporters) and use these measurements to infer the dynamical states of unobserved species (e.g., the gene state). This process, known as stochastic filtering for intracellular reaction systems, allows scientists to gain insights into unobserved chemical species. Also, this in turn can lead to the development of improved control strategies for the reacting process. Here, we apply the proposed RB-CME solver to this filtering problem.

Mathematically, we can model the observation channels by

$$Y(t_i) = h(X(t_i)) + \sigma W_i \quad [7]$$

where t_i are the observation time points, $Y(t_i)$ is a vector of observations with each element corresponding to a particular light frequency, $h(\cdot)$ is a vector-valued function indicating the ideal relation between the measurement and the state, σ is a diagonal matrix indicating observation noise intensities, and W_i are

[†] In the rest of this paper, we refer the Monte-Carlo method to the conventional one we introduced in the previous subsection.

vectors of independent standard Gaussian noise. Stochastic filtering aims to compute the conditional probability distribution of the state $X(t_i)$ (for every $i = 1, 2, \dots$) given the observations up to time t_i , i.e., $\pi_{t_i}(x) \triangleq \mathbb{P}(X(t_i) = x | Y(t_s), 1 \leq s \leq i)$. Let us denote $\rho_{t_{i+1}}(x) \triangleq \mathbb{P}(X(t_{i+1}) = x | Y(t_s), 1 \leq s \leq i)$. By Bayes' rule, $\pi_{t_i}(x)$ satisfies the following recursive formulas (65)

$$\rho_{t_{i+1}}(x) = \sum_{x' \in \mathbb{Z}_{\geq 0}^n} \mathbb{P}(X(t_{i+1}) = x | X(t_i) = x') \pi_{t_i}(x') \quad [8]$$

$$\pi_{t_{i+1}}(x) \propto L(Y(t_{i+1}) | x) \rho_{t_{i+1}}(x) \quad [9]$$

where $L(y|x)$ is the density function of the distribution $\mathbb{P}(Y(t_{i+1}) \in dy | X(t_{i+1}) = x)$ (usually called the likelihood function). We can interpret Eq. [8] as the prediction of the state at the next time point t_{i+1} using the observation up to the current time t_i . We interpret Eq. [9] as the adjustment of the prediction according to the new observation. Note that the prediction step Eq. [8] is actually solving a CME with $\pi_{t_i}(\cdot)$ being the initial probability distribution and $\rho_{t_{i+1}}(\cdot)$ being the final solution. Consequently, the filtering problem can be seen as a combination of the usual CME and an adjustment step.

Following the idea above, one can solve the filtering problem by applying various CME solvers to the prediction step. Conventional methods include the particle filter that uses the Monte-Carlo method for Eq. [8]] (e.g., (21, 66, 67), to cite a few) and the direct approach that uses the FSP for Eq. [8]. Similar to the situation in solving CMEs, these two approaches scale poorly with the system dimension when solving the filtering problem. Specifically, the particle filter has a similar L_1 error to the Monte-Carlo method (for solving CMEs) when they are applied to the same system (SI Appendix, Section S2.B), and, therefore, the particle filter can be inaccurate when the system dimension is large. Besides, applying the FSP to the filtering problem can be time-consuming because the size of the state space scales exponentially with n .

Here, we solve the filtering problem by applying the RB-CME solver to Eq. [8]; we call this approach the Rao-Blackwellized particle filter (RB-PF). In this case, the follower subsystems need to be conditionally independent given the trajectories of both the leader system and the observation. To this end, we require the following condition for the system decomposition in addition to C1 and C2. (The proof is given in SI Appendix, Section S4.A.)

- C3 Each observation channel cannot be affected by two follower subsystems. In other words, each entry of $h(\cdot)$ can depend on one follower subsystem at most (besides the leader system).

This requirement automatically holds in many practical applications where fluorescent reporters are used as the probe. In these cases, experimentalists usually give different colors to different genes, thereby establishing a one-by-one correspondence between observed signals and actual gene products. We provide the modified leader-follower decomposition algorithm for the filtering problem in SI Appendix, Section S4.A, and the detailed algorithm of the RB-PF in SI Appendix, Section S4.B.

We found that when applied to the same chemical reacting system, the RB-PF (for the filtering problem) usually has similar accuracy to the RB-CME solver (for solving CMEs) SI

Appendix, Section S4.C, which suggests that the RB-PF also scales favorably with the system dimension. Also, this means that for a given reaction system, if the RB-CME solver can accurately solve its CME, then the RB-PF can also accurately solve its filtering problem and vice versa. To conclude, this consistency result guarantees good performance of the RB-PF.

E. Rao-Blackwell method for cell-specific model identification.

The biological dynamics occurring within cells are not exactly known to scientists. This fact gives rise to another important topic in biology, i.e., model identification, which aims to develop robust and accurate math models for biological processes from given datasets. Securing a good math model can provide a deep understanding of the underlying biological mechanisms and allow for more accurate predictions of system behaviors. In this subsection, we present a novel method for the identification of cell-specific models by exploiting the RB-CME solver.

In the identification problem, we consider that each cell can undergo r reactions as in Eq. [1], consisting of all possible biological mechanisms within the cell. Also, the system follows a dynamical equation $X(t) = X(0) + \sum_{j=1}^r \zeta_j R_j \left(\int_0^t \lambda_j(\Theta, X(s)) ds \right)$, which is similar to Eq. [2] with the exception that the propensity functions are also dependent on unknown parameters $\Theta \in \mathbb{R}^{\tilde{r}}$ (e.g., reaction constants and hill coefficients). Usually, these parameters can be any value within certain ranges; however, for simplicity, we consider that Θ only takes values in a discrete set Θ that provides a fine-grained representation of these continuous regions. As in the filtering problem, we consider that a cell is tracked and measured under a microscope at different time points (t_1, \dots, t_{n_f}) , with the observations being represented by $Y(t_1), \dots, Y(t_{n_f})$. Ultimately, this cell-specific identification problem aims to calculate the conditional probability of the parameters Θ given the measurements, i.e., $P(\Theta = \theta | Y(t_s), 1 \leq s \leq n_f)$.

Similar to stochastic filtering, this identification problem can also be solved in a recursively manner. Essentially, the parameters Θ can be viewed as the state of some additional special chemical species in the system, which can take non-integer values and remain constant over time. From this perspective, this cell-specific identification problem aims to infer the hidden states of these special species, a task closely aligned with the filtering problem introduced in the preceding subsection. Therefore, by denoting the condition distributions $\pi_{t_i}(\theta, x) \triangleq \mathbb{P}(\Theta = \theta, X(t_i) = x | Y(t_s), 1 \leq s \leq i)$ and $\rho_{t_{i+1}}(\theta, x) \triangleq \mathbb{P}(\Theta = \theta, X(t_{i+1}) = x | Y(t_s), 1 \leq s \leq i)$, we can solve this identification problem using the following recursive formulas:

$$\rho_{t_{i+1}}(\theta, x) = \sum_{x' \in \mathbb{Z}_{\geq 0}^n} \mathbb{P}(X(t_{i+1}) = x | \Theta = \theta, X(t_i) = x') \pi_{t_i}(\theta, x') \quad [10]$$

$$\pi_{t_{i+1}}(\theta, x) \propto L(Y(t_{i+1}) | x) \rho_{t_{i+1}}(\theta, x) \quad [11]$$

$$\mathbb{P}(\Theta = \theta | Y(t_s), 1 \leq s \leq t_{n_f}) = \sum_{x \in \mathbb{Z}_{\geq 0}^n} \pi_{t_{n_f}}(\theta, x). \quad [12]$$

Here, Eq. [10] represents the prediction of the entire state $(\Theta, X(t))$ at the subsequent time point t_{i+1} using the observation up to the current time t_i . Meanwhile, Eq. [11] adjusts this prediction using the new measurement at time t_{i+1} . Finally,

Eq. [12] obtains the result of model identification by marginalization. Since Θ can be viewed as the state of additional chemical species, the probability distribution of $(\Theta, X(t))$ evolves according to an augmented CME (SI Appendix, Section S5.B). Consequently, the prediction step Eq. [10] corresponds to solving this augmented CME with $\pi_{t_i}(\cdot)$ as the initial probability distribution and $\rho_{t_{i+1}}(\cdot)$ as the solution at the final time. Thus, this model identification problem can also be interpreted as a combination of solving the augmented CMEs and making subsequent adjustments.

From these viewpoints, we propose using the RB-PF to solve this identification problem, i.e., applying the RB-CME solver to the prediction step. To achieve this, we need to classify all the chemical species and parameters into leader and follower systems, and the follower subsystems must satisfy C1–C3 so that they are conditionally independent given the trajectories of both the leader system and the measurement. In addition, it is noteworthy that classical particle filtering is generally ineffective for the inference of static hidden state variables (e.g., model parameters) due to sample degeneracy (68). For this reason, we require our identification algorithm to classify all the parameters as follower components, which allows their inference to be aided by a filtering approach (e.g., filtered FSP) rather than being identified purely by classical particle filtering.

C4 All the model parameters Θ are classified as follower components of the system.

The detailed algorithm for the leader-follower decomposition adhering to C1–C4 is provided in SI Appendix, Section S5.C.1, and the Rao-Blackwell algorithm for model identification is presented in SI Appendix, Section S5.C.2.

It is also worth noting that the idea of adapting the Rao-Blackwell method for model identification has also been explored in the literature (11). The method described in that paper explicitly marginalizes out the uncertainty of parameters given the dynamics of the chemical species, and it demonstrated strong performance with both numerical and experimental data. Several significant differences exist between our approach and that one. First, the method in (11) is tailored specifically for systems with mass-action kinetics, which limits its applicability to many real-world biological scenarios that feature non-mass-action kinetics (e.g., Michaelis-Menten Kinetics and Hill-type dynamics). In contrast, our approach is not constrained by the type of kinetics and therefore has much broader applicability. Second, (11) relies on classical particle filtering for the inference of all chemical species, whereas our approach can classify some species as follower components, which allows us to infer them with the assistance of a suitable filtering approach (e.g., the filtered FSP). To conclude, though both approaches draw inspiration from the Rao-Blackwellization technique, our approach is more applicable and scalable.

2. Numerical case studies

We first illustrate our approach for solving CMEs and stochastic filtering problems through several biologically relevant numerical examples. Unless stated otherwise, all experiments were performed on the Euler computing cluster at ETH Zurich, utilizing computational nodes with 2.25-GHz, 12-core CPUs. The code is available on GitHub: “<https://github.com/ZhouFang92/Rao-Blackwellized-CME-solver>”.

A. Application of the RB-CME solver to a class of linear reaction systems. To demonstrate the scalability of the RB-CME solver, we applied our approach to a class of expandable linear networks shown in Figure 3A and compared its performance with the traditional approaches. Specifically, the linear network consists of n chemical species and three types of reactions: the production ($\emptyset \rightarrow S_i$), the degradation ($S_i \rightarrow \emptyset$), and the conversion of S_i into S_{i+1} ($i \leq n-1$). We modeled the reactions to follow mass-action kinetics with the rate constants presented in the caption of Figure 3. At the initial time, the molecular counts of different species are independent and have a Poisson probability with mean 0.5. In this setting, the associated CME can be solved by a multivariate Poisson distribution whose mean evolves according to the deterministic dynamics of the system (22).

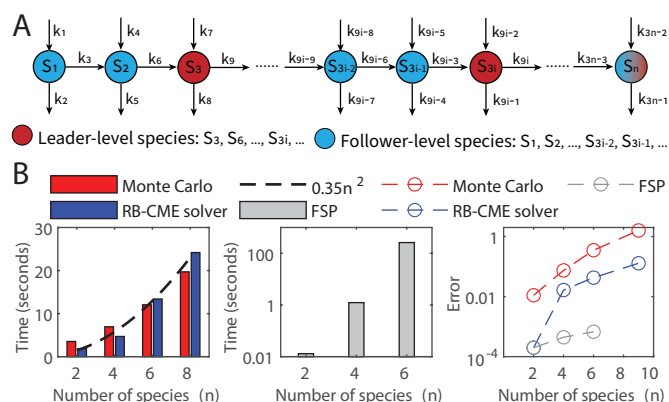


Fig. 3. Scalability of the RB-CME solver in a class of linear networks. (A) A class of linear networks that consists of three types of reactions: the production, degradation, and conversion of S_i into S_{i+1} ($i = 1, \dots, n-1$). All the reactions follow mass-action kinetics, and their reaction constants are $k_1 = 2.4$, $k_{3n-1} = 1.6$, $k_{3i+1} = 0.9$, $k_{3i-1} = 0.6$, and $k_{3i} = 1$ ($i = 1, \dots, n-1$). At the initial time, each species has a Poisson probability with mean 0.5, and all of them are independent. (B) The scalability of the Monte-Carlo method (with 10^5 samples), RB-CME solver (with 10^4 samples), and finite state projection approach (with the truncated space $\bigotimes_{i=1}^n \{0, 1, \dots, 9\}$) when solving the CME at time 10. We used the filtered FSP as our chosen filtering approach in the RB-CME solver. In the third block, the error is evaluated by the L_1 distance between the numerical solution and the exact solution of the CME. This panel tells that the computational time of the Monte-Carlo method and that of the RB-CME solver both grow quadratically with the system dimension (n), whereas the computational time of the FSP method grows exponentially with n . Moreover, the error of the Monte-Carlo method and that of the RB-CME solver both grow exponentially with n , but the latter grows much slower than the former (see the slopes of the linear-like curves in the log-domain). Notably, the RB-CME solver is as accurate as the FSP method when $n = 2$. It is because, in this case, no leader-level species exist, and, therefore, the RB-CME solver with the filtered FSP as the chosen filtering approach is equivalent to the FSP method.

We first compared the scalability of different approaches with respect to accuracy and efficiency. For the FSP method, we truncated the state space by $\bigotimes_{i=1}^n \{0, 1, \dots, 9\}$, which contains most of the probability. Similarly, with the RB-CME solver, we truncated the state space for each follower-level species by $\{0, 1, \dots, 9\}$ and required each follower subsystem to contain no more than 100 states. In this setting, our leader-follower decomposition algorithm consistently classifies species S_{3i} ($i = 1, 2, \dots$) as leader-level species and the rest as follower-level species. Each pair of S_{3i-2} and S_{3i-1} ($i = 1, 2, \dots$) form a follower subsystem. In this example, we selected the filtered FSP as the filtering approach to be used in the RB-CME solver. Also, we set the RB-CME solver and the Monte-Carlo method to have respectively 10^4 samples and 10^5 samples so that their

computational time is relatively the same. The experimental results are shown in Figure 3B.

Figure 3B indicates that the computational time of the RB-CME solver scales well with the system dimension (n), and its accuracy is much better than that of the Monte-Carlo method at the same time cost. Specifically, the first block in Figure 3B shows that the computational time of the Monte-Carlo method grows quadratically with n , which is because the Gillespie method has the computational complexity $\mathcal{O}(\text{\#reaction channels} \times \text{\#reaction firing events})$ (see (69, Section III)), and both of these quantities grow linearly with n in this example. For the FSP method, the algorithm's time-complexity is linear with the truncated space size, which scales exponentially with the dimension n in this example, so its cost also grows exponentially with n (see the second block of Figure 3B). In contrast, though the RB-CME solver utilizes the filtered FSP to the follower system whose size also grows exponentially with n , its computational time still scales quadratically with n (see the first block of Figure 3B). This reduced computational complexity is because we apply the filtering approach to each follower subsystems separately rather than the whole system. Therefore, in our algorithm, the computational cost of the FSP part becomes $\mathcal{O}(\text{size of the largest follower subsystem} \times \text{\#follower subsystems})$ where the first term is fixed ($= 100$), and the second term scales linearly with n . Additionally, parallel computing further aids in reducing computational time. As for the accuracy, the third block of Figure 3B tells that the error of the Monte-Carlo method and that of the RB-CME solver both scale exponentially with n , but the latter grows much slower than the former. Notably, when the system dimension is two, i.e., all the species are follower-level species, the RB-CME solver (with the filtered FSP as the chosen filtering approach) is equivalent to the FSP method, and both approaches are equivalently accurate. All these results indicate that the RB-CME solver is a good compromise between the Monte-Carlo method and the FSP approach, and it is more favorable for high-dimensional problems.

To further understand the benefit of the RB-CME solver, we investigated its use in more detail on a linear network with six species (see Figure 4). From the results, we can observe that both the Monte-Carlo method and the RB-CME solver converge to the exact probability distribution at the rate of $1/\sqrt{N}$, which agrees with the law of large numbers (see Figure 4B). Moreover, given the same sample size (resp., the same computational time), the RB-CME solver is significantly more accurate than the Monte-Carlo method with an improvement of 20 times (resp., 8 times), respectively. We also studied the performance of both approaches in estimating the marginal distributions when the time costs are relatively the same (see Figure 4C). The result shows that both methods accurately approximate the marginal distribution for individual species, but their performance in estimating joint probabilities is very different. Concretely, the RB-CME solver is more accurate in estimating the follower system, especially the first follower subsystem consisting of S_1 and S_2 , but it is less accurate in estimating the leader systems (see the last block in Figure 4B). The relative inaccuracy of the RB-CME solver for the leader system is attributed to the fact that both the RB-CME solver and the Monte-Carlo method use the same protocol to estimate the leader system, and for the same time cost, the RB-CME

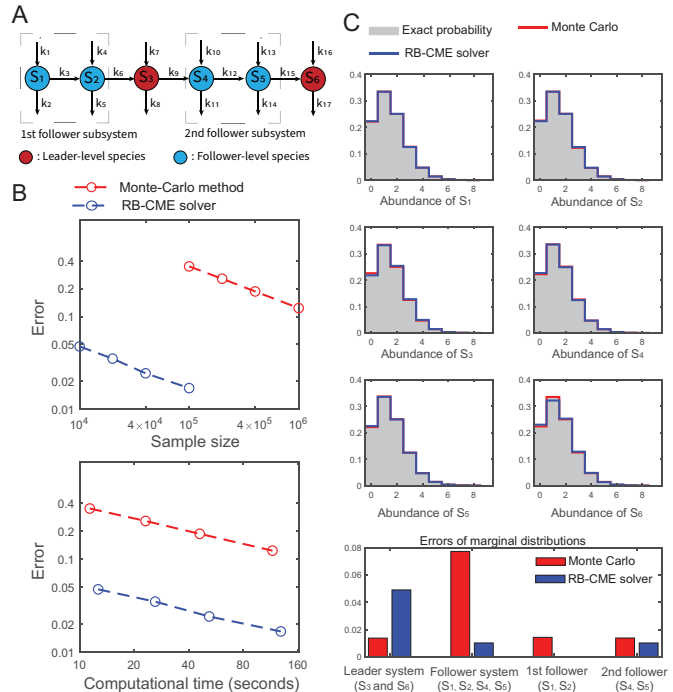


Fig. 4. Comparison of the RB-CME solver and Monte-Carlo method in the linear network with six species. (A) The diagram of the linear network with six species: all the settings are the same as that in Figure 3. (B) Convergence of both approaches in terms of the L_1 error. Both methods converge at the rate of $1/\sqrt{N}$ or $1/\sqrt{T}$, where N and T are the sample size and computational time, respectively. This also implies that for both approaches, the computational time is proportional to the sample size (see the last block). With the same sample size, the RB-CME solver is 20 times more accurate than the Monte-Carlo method, and at the same time cost, the RB-CME solver is 8 times more accurate. (C) Performance of the RB-CME solver (with 10^4 samples) and Monte-Carlo method (with 10^5 samples) in estimating marginal distributions. Both methods accurately approximate the marginal distributions for individual species, but they perform quite differently in estimating joint probabilities. Specifically, the RB-CME solver is more accurate in estimating the follower system, particularly the first follower subsystem consisting of S_1 and S_2 , but it is less accurate in estimating the leader system.

solver has fewer samples than the Monte-Carlo method. Despite this issue, the RB-CME solver still has a much better performance in estimating the whole system because of the greater benefit gained from the estimation of the follower part.

B. Application of the RB-CME solver to the repressilator. To demonstrate the performance of the RB-CME solver for non-linear systems, we consider a well-known genetic circuit called the repressilator (see Figure 5), where three gene expression systems cyclically repress each other. The repressilator was first engineered by Elowitz and Leibler in 2000 (18), and its name comes from the cyclical repression topology and its oscillatory dynamic behavior. In this system, the nonlinearity comes from the mRNA production processes, whose propensities are hill functions. More details about the modeling are presented in *SI Appendix, Section S6.A*.

In this model, we intended to compare the RB-CME solver with other approaches. With the parameters set in *SI Appendix, Section S6.A*, most of the probability is contained in the state space where each mRNA has no more than 20 copies, and each protein has no more than 200 copies. Notice that storing a probability distribution on this state space requires about 500 GB ($8 \text{ bytes/state} \times (20^3 \times 200^3) \text{ states}$),

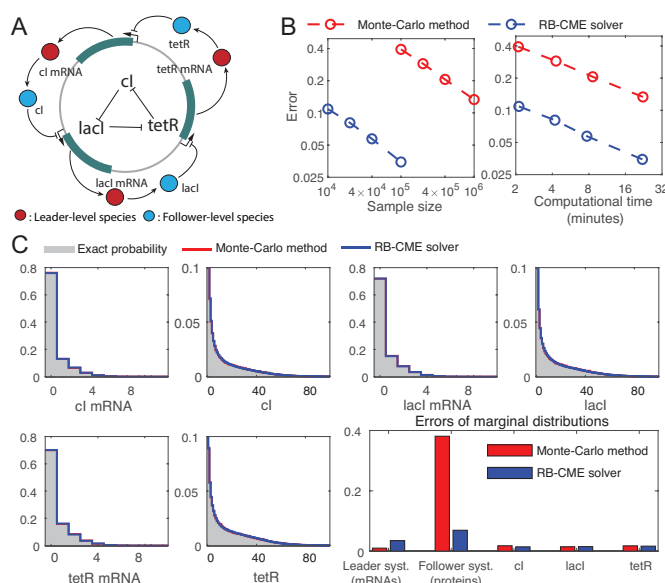


Fig. 5. Performance of the RB-CME solver in the repressilator. (A) Diagram of the repressilator with three gene expression systems whose proteins cyclically repress each other. Our method classifies all the mRNAs as leader-level species and all the proteins as follower-level species. (B) Convergence of the Monte-Carlo method and the RB-CME solver (with the filtered FSP as the chosen filtering approach). We depict the error by the sum of the L_1 errors in estimating the leader system and the follower system. The exact probability distribution is approximated by the Monte-Carlo method with 3×10^9 samples. Given the same sample size or the same time cost, the RB-CME solver is much more accurate than the Monte-Carlo method. (C) Performance of the RB-CME solver (with 10^4 samples) and the Monte-Carlo method (with 10^5 samples) in estimating the marginal distributions. Both approaches have relatively the same computational time, and they both accurately estimate the marginal distributions of individual species. The superior performance of the RB-CME solver is attributed to the estimation of the follower system, which dominates in the whole estimation problem.

so applying the FSP approach to this example is impractical. Alternatively, to get an accurate solution of the CME, we simulated 3×10^9 trajectories of this dynamical system and used the empirical distribution to approximate the exact probability distribution. This whole procedure took 10 graphics processing units (GPUs) about 24 hours! For the RB-CME solver, we set the truncated state space for each mRNA to be $\{0, 1, \dots, 19\}$ (if the mRNA is classified as a follower-level species), the truncated state space for each protein to be $\{0, 1, \dots, 199\}$ (if the protein is classified as a follower-level species), and each follower subsystem to contain no more than 200 states. In this setting, the leader-follower decomposition algorithm classifies all the mRNAs as the leader-level species and all the proteins as the follower-level species. Finally, we applied both the RB-CME solver (with the filtered FSP as the chosed filtering approach) and Monte-Carlo method to the repressilator; the results are shown in Figure 5.

Again, the results show that the RB-CME solver is efficient and accurate in this example. Specifically, Figure 5B tells that given the same sample size, the RB-CME solver is about 15 times more accurate than the Monte-Carlo method, and given the same time cost, the RB-CME solver is 4 times more accurate. Notably, by only taking one CPU and about half an hour, our method can provide a very accurate solution with an L_1 error of 3% in estimating marginal distributions of the mRNAs and the proteins. Furthermore, we compared the performance of both approaches in estimating the marginal

distributions when their computational time is relatively the same (see Figure 5C). From this panel, we can observe that both methods accurately estimate the marginal distribution of individual species, and the superior performance of the RB-CME solver is attributed to the estimation of the follower system. More specifically, the proteins have much more molecular counts than the mRNAs, and, therefore, the marginal probability distribution of the proteins tends to be more dispersed than that of the mRNAs. According to our analysis in Section 1.A, this explains why the estimation error of the follower system (the proteins) dominates in the Monte-Carlo method (see Figure 5C). In contrast, the RB-CME solver has a much lower error in estimating the proteins because it applies a FSP method to the follower system, and this advantage greatly dominates the disadvantage of the RB-CME solver in estimating the mRNAs (see Figure 5C). More interestingly, the improvement of the RB-CME solver in estimating the whole follower system is much greater than that in estimating each follower subsystem (see Figure 5C) thanks to our modularization strategy, which further demonstrates the scalability of our approach with increasing system dimensions. These observations also imply that the RB-CME solver can accurately estimate other genetic circuits if all the species with large molecular counts are classified as follower-level species.

C. Applying the RB-PF to stochastic filtering for the genetic toggle switch. Now, we consider the filtering problem for another well-known genetic circuit, called the genetic toggle switch, which was first engineered by Gardner, Cantor, and Collins in 2000 (70). This circuit consists of two gene expression systems, whose protein products repress each other's expression (see Figure 6A), and their trajectories exhibit switching behaviors (see Figure 6B). We assume that the first protein is fluorescent and measured by a microscope at several time points, and our goal is to infer the hidden dynamical state given the observations. More details about the modeling are presented in *SI Appendix, Section S6.B*.

In this example, we intended to compare the RB-PF with other filtering approaches. First, we simulated a trajectory of the system and generated observations at 10 different time points. Our task was to infer the hidden states based on these generated observations. To get an accurate approximation of the exact filter, we applied the FSP to the problem with a state space where each protein has fewer than 200 copies. This procedure took 6.5 hours! For the RB-PF, we set the truncated state space for each protein to be $\{0, 1, \dots, 199\}$ (if it is classified as a follower-level species) and the truncated state space for genes to be $\{0, 1\}$ (if it is classified as a follower-level species). By letting each follower subsystem have no more than 200 states, our approach classifies all the proteins as the follower-level species and the rest as leader-level species. Finally, we applied both the RB-PF and conventional particle filter to the filtering problem. In the prediction step, this RB-PF uses an RB-CME solver that adopts the filtered FSP as the selected filtering algorithm in its framework. The results are shown in Figure 6 and Figure 7.

In Figure 6, we compare the performance of the RB-PF (10^4 samples) and PF (10^5 samples), which have the similar time cost (see Figure 7B). The numerical results show that though the RB-PF and the PF have similar performance in estimating the conditional mean and variance of the second protein (Figure 6B), the RB-PF is far more accurate in estimating the

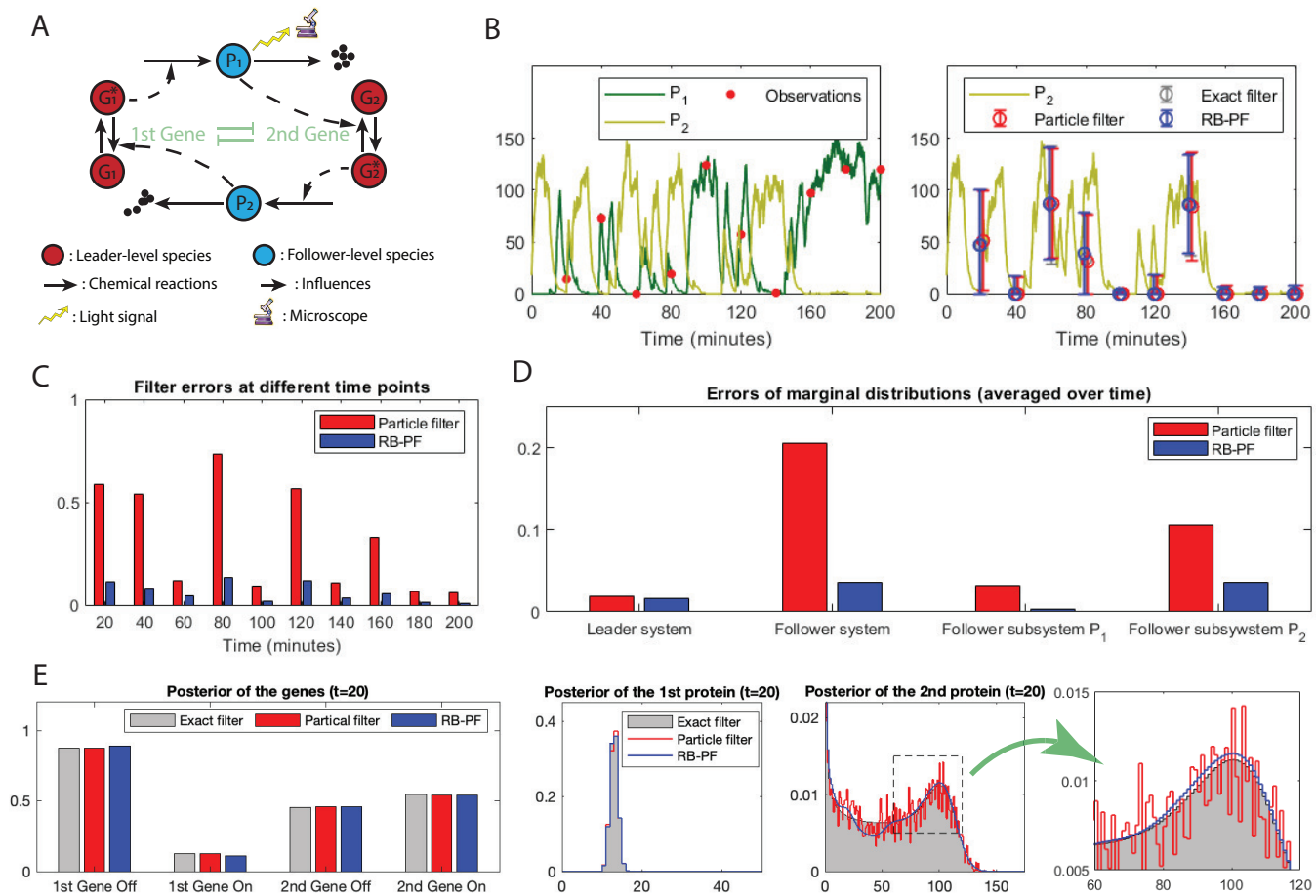


Fig. 6. Filtering results for the genetic toggle switch. (A) Description of the genetic toggle switch with two gene expression systems whose proteins repress each others' genes. The protein of the first gene is fluorescent and can be measured by a microscope. The observation noise intensity is set to be 1. In this setting, our approach classifies all the genes as the leader-level species and all the proteins as the follower-level species. (B) Performance of different filters in estimating the conditional mean of the second protein given a discrete-time trajectory of the first protein. In this example, the exact filter is approximated by the FSP method. The particle filter and the Rao-Blackwellized particle filter (RB-PF) have, respectively, 10^5 samples and 10^4 samples, and, in this setting, they consume relatively the same computational time. The RB-PF uses an RB-CME solver with the filtered FSP as the chosen filtering algorithm in its framework. This panel tells that all the filters have similar performance in estimating the conditional mean of the second protein. (C) L_1 errors of the PF and the RB-PF at different time points. It shows that the RB-PF outperforms the PF at every observation time point, and the improvement ratio is quite significant. (D) Average L_1 errors of the PF and the RB-PF in estimating the marginal distributions. (E) Evaluation of marginal posteriors at the first observation time ($t = 20$). The last two panels show that the major difference between the PF and the RB-PF lies in the estimates of the follower system.

whole conditional probability distribution. Specifically, the RB-PF consistently performs better than the PF at different observation time points, and the improvement ratio is always significant (Figure 6C). When looking at the marginal distributions, we can observe that the superior performance of the RB-PF also comes from the estimation of the follower system (Figure 6D). Particularly, in this example, the estimation error of the follower system (the proteins) dominates in the PF (see Figure 6D), as its conditional probability distribution is more dispersed than that of the leader system. Thanks to the modularization strategy and the filtered FSP that our method applies to the follower system, the RB-PF is more accurate than the PF in estimating the follower system (see Figure 6D), and, therefore, the RB-PF is significantly more accurate in estimating the whole conditional probability distribution. From Figure 6E, we can observe that for the follower system, the solution of the RB-PF is more smooth and accurate, whereas the result of the PF is fuzzier and more inaccurate, which again shows the advantage of RB-PF. More interestingly, in

this filtering problem, the RB-PF and the PF have similar performance in estimating the leader system (Figure 6D and Figure 6E), quite different from the situation in solving the CME where the Rao-Blackwell method is less accurate than the conventional Monte-Carlo method in estimating the leader (see Figure 4 and Figure 5 in previous examples). The reason for this phenomenon is that in the adjustment step of the filtering algorithm, the conditional probability distribution of the leader system is influenced by the conditional probability distribution of the follower system via the likelihood function, and, therefore, the leader system can also benefit from Rao-Blackwellization.

We also compared the performance of the RB-PF and the PF under various different conditions. First, the performance of both filters is quite robust to the variation of the observation noise intensity Figure 7A. Moreover, the RB-PF and the RB-CME solver perform similarly in their associated problems (see Figure 7B and Figure 7C); the same is true for the PF and the Monte-Carlo method. These results agree well with our

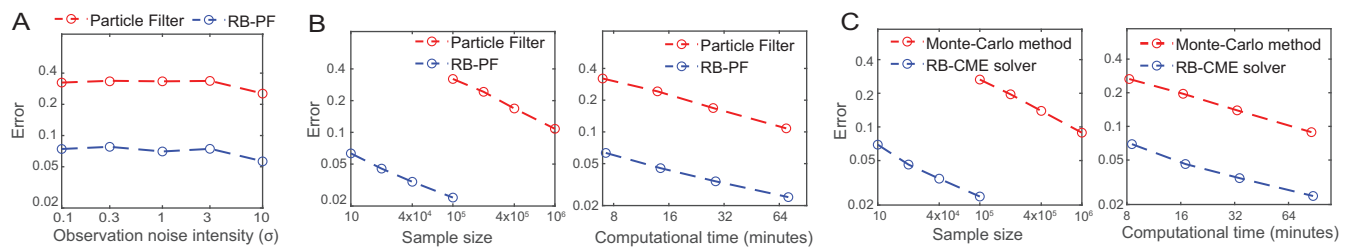


Fig. 7. Performance of the Rao-Blackwellized particle filter (RB-PF) in different settings. In this figure, we still consider the filtering problem for the genetic toggle switch (see Figure 6) and evaluate the performance of the RB-PF under different problem settings. Here, the RB-PF uses an RB-CME solver that adopts the filtered FSP as the selected filtering algorithm in its framework. (A) Performance of the particle filter (PF) and the RB-PF under different observation noise intensity. The PF and RB-PF have 10^5 samples and 10^4 samples, respectively, and the performance is evaluated by the average L_1 error over 10 observation time points. The panel shows that the performance of both filters are robust to the intensity of the observation noise. (B) Convergence of the PF and the RB-PF under a fixed observation noise intensity ($\sigma = 1$). Again, the performance is evaluated by the average L_1 error over 10 observation time points. (C) Convergence of the Monte-Carlo method and the RB-CME solver in computing the CME of the genetic toggle switch at time 200. In this panel, the performance is evaluated by the L_1 error between the numerical result and the exact solution, where the exact solution is approximated by the FSP method. The last two panels show that the RB-CME solver and the RB-PF have similar performance in their associated problems given the same sample size or the same computational time. The same is also true for the Monte-Carlo method and PF.

analysis in the theoretical part (also see *SI Appendix, Section S2.B and Section S4.C*). Also, they suggest a consistency result that for one system, if the RB-CME solver can accurately solve its CME, then the RB-PF should also perform well in its filtering problem and vice versa. In addition, similar to the situation in solving the CME, the RB-PF is orders of magnitude more accurate than the PF given the same sample size or the same time cost Figure 7B. All these conclusions indicate the reliability of our approach.

3. Identifying transcription dynamics in Yeast cells from single-cell time-course trajectories

In this section, we applied our identification algorithm to an experimental dataset of yeast cells and used the result to analyze the sources of cell-to-cell variability.

Cell-to-cell heterogeneity is influenced by two main factors: the random firing of reactions within individual cells (known as intrinsic noise) and the variability in the parameters that determine the reaction propensities across the population (known as extrinsic noise). Understanding the role of intrinsic and extrinsic noise in shaping the cell-to-cell variability is an important topic in biology. This section investigates this issue in the context of transcription dynamics in yeast cells, using our parameter identification method introduced in Section 1.E.

We focused on genetically identical cells built in (71) (see Figure 8A for an illustration of the synthetic circuit). Within this gene circuit, VP-EL222 homodimerizes in the presence of light, which then binds to its cognate promoter (a fusion of several EL222-binding sites) to stimulate the expression of a downstream gene. The produced RNAs contain stem-loops that can be recognized and bound by a fluorescent reporter (tdPCP-tdmRuby3); this structure allows for the measurement of RNA dynamics under a microscope.

A reaction network model for this cell system is presented in Figure 8B. Since the DNA contains several sites, a simple telegraph model containing only two gene states (ON and OFF) may not adequately represent this practical system. Consequently, we considered a three-gene-state model, where the gene has one inactive state and two active states. In these different active states, the RNAs are transcribed at different rates. Particularly, when the reaction rate k_3 equals zero, this model degenerates to the conventional telegraph model. Also, we need to point out that the method in (11) is not

applicable to this system. In that method, the system needs to have distinct reaction vectors when expressed by mass-action kinetics; however, this requirement does not hold here due to the mRNA transcription dynamics.

In this section, we aimed to identify the dynamical parameters of each cell from the time-trajectory of its RNA dynamics and then investigate the contribution of intrinsic and extrinsic noise to the overall cell-to-cell variability. In particular, we compare the intrinsic and extrinsic noise estimates obtained from our cell-specific inference results with the estimates obtained directly from the time-course experimental data. To perform this direct estimation, we propose a novel decomposition technique which does not require dual-reporters (58, 60, 61), but produces equivalent results under the assumption of stability of the underlying stochastic reaction network.

The code for the analysis in this section is available on GitHub: "https://github.com/ZhouFang92/Rao-Blackwell-method-for-cell-specific-model-identification".

A. In-silico verification of the proposed identification method.

First, we examined the accuracy of our Rao-Blackwell method in identifying model parameters through numerical simulation. In this simulation study, we assumed that k_1, \dots, k_4 were drawn from the set $\{0, 0.05, \dots, 1\}$, and k_{p_1}, k_{p_2} were within the set $\{1, 2, \dots, 10\}$. All these parameters were in units of minute^{-1} and had uniform prior distributions, with the exception of k_3 which represents the rate of switching from the first active state G_1 to the second G_2 . For k_3 , half of its initial probability mass was allocated at zero to indicate the uncertainty of whether the number of active gene states is one or two, and the rest of the probability was uniformly distributed over states $0.05, 0.1, \dots, 1$. At the initial time, we set the gene state to G_0 and the mRNA count to zero. Moreover, the mRNA count was measured every minute with observation noise intensity of 0.1.

We first generated a simulated trajectory of the given model with randomly selected parameters. Then, we applied our Rao-Blackwell method to inferring these model parameters using the simulated time-course measurements of this system. We truncated the space for the mRNA count to be $\{0, 1, \dots, 20\}$ and that for each gene to be $\{0, 1\}$; also, we set the sample size in our algorithm to 10,000. By requiring the size of the maximum follower subsystem to be less than 30,000, our algorithm classified G_1 and G_2 as leader species and assigned the remain-

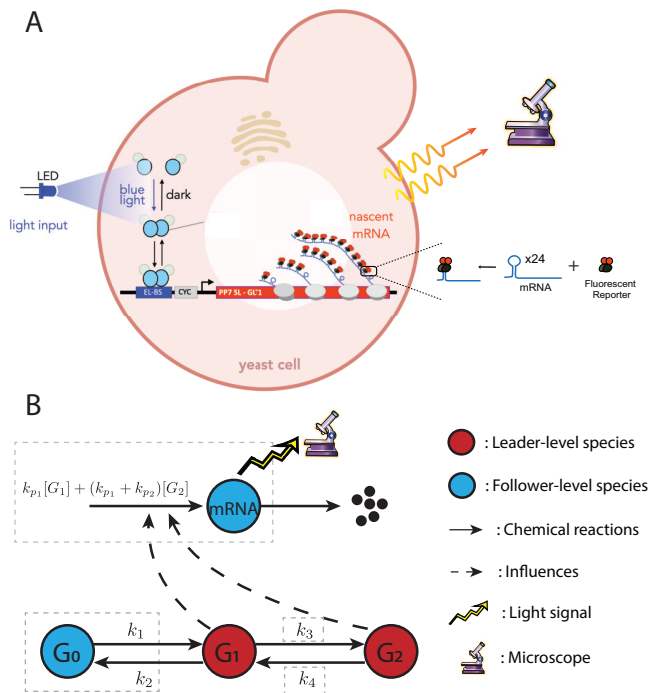


Fig. 8. Transcription system in Yeast cells. (A) Synthetic circuit in yeast cells. In the presence of light, EL222 can dimerize, bind to EL222-binding sites (EL-BS), and activate the expression of the downstream gene. The transcribed RNAs contain stem-loops to which fluorescent reporters can attach, allowing visualization of RNA dynamics. (B) Reaction network model for the gene circuit. In this model, the gene has three states: one inactive state G_0 and two active states G_1 and G_2 . In the first active state G_1 , the mRNA are transcribed at a rate k_{p1} . In the second active state G_2 , the mRNA is transcribed at a higher rate $k_{p1} + k_{p2}$, where k_{p2} represents an additional rate of transcription beyond k_{p1} . The parameters k_1, \dots, k_4, k_{p1} , and k_{p2} are unknown but fixed over time, and the mRNA degrades at a linear rate with the rate constant 1. Under the setting in Section 3.A and Section 3.B, our algorithm consistently classifies G_1 and G_2 as leader species, with the remaining elements categorized as follower components. The follower part has four subsystems: the first subsystem consists of G_0 , k_1 , and k_2 , the second subsystem consists of k_3 , the third subsystem consists of k_4 , and the last subsystem consists of k_{p1} , k_{p2} and mRNA.

ing components into four follower subsystems (see Figure 8B for the decomposition). Moreover, in prediction steps Eq. [10], our method uses an RB-CME solver that adopts the filtered FSP as the selected filtering algorithm in its framework. The numerical result of our algorithm is presented in Figure 9A.

The grey box in Figure 9A illustrates that in this inference problem, our algorithm accurately identifies the model parameters, with the maximum a posteriori (MAP) estimates being exactly or very close to the real values. Moreover, the provided conditional distributions of the parameters are quite narrow, indicating that these estimates have high confidence.

In practical applications, the actual parameter values for a specific cell are usually unknown, which makes it impractical to validate the maximum a posteriori (MAP) estimate by directly comparing it with the true values. Consequently, there is a need for a method to validate the identification results using experimental data, without the necessity for precise parameter information. To this end, we propose a new verification method which compares the stationary probability distributions of the target cell and the inferred model. Note that the stationary probability distribution of the target cell is still not directly available; by definition, it requires the

measurement of a cell population having identical parameters to the target cell. However, thanks to the stability[‡] is that the occupation time distribution $\mathbb{P}_{oc}(T, x) \triangleq \frac{1}{T} \int_0^T \mathbb{1}(X(t) = x) dt$ converges to the stationary distribution as $T \rightarrow \infty$. of the system (see SI Appendix, Section S7.B for the proof), this stationary distribution can be approximated by the occupation time distribution of the target cell, which is available from the mRNA measurements. Moreover, to account for the measurement noise, we rounded each measurement to the nearest integer and used the rounded numbers to compute the occupation time distribution. Meanwhile, the stationary distribution of the inferred model is approximated by a distribution computed by the FSP at a large time point. The time point is chosen so that the stationary distribution is approximately reached. In conclusion, this new verification method, through comparing the occupation time distribution (of the target cell) and the stationary distribution (of the inferred model), is feasible in the experimental setting.

The comparison of these two distributions is presented in the bottom-left panel of Figure 9A. We can observe that the two distributions almost overlap perfectly, with a Kullback–Leibler (KL) divergence of 0.02. The slight difference might be attributed to the small difference between the MAP estimates and the true parameter values (see the grey box in Figure 9A) and the imperfect approximation of the stationary distribution via the occupation time distribution. In general, the closely matched distributions suggest that our identification approach is accurate in this example.

Next, we tested whether our approach could correctly identify the model when the target system had only two gene states, i.e., $k_3 = 0$. Similar to the previous case, we first simulated a trajectory of the system (with parameters $k_1 = 0.3$, $k_2 = 0.4$, $k_3 = 0$, $k_4 = 0$, $k_{p1} = 3$, and $k_{p2} = 0$), and then we used our Rao-Blackwell method to identify the model via the simulated measurements. We maintained the same settings for the identification algorithm as in the previous numerical example in this subsection. The result, presented in Figure 9B, shows that our algorithm accurately infers the parameters k_1 , k_2 , k_3 , and k_{p1} , with the MAP estimates either matching or closely resembling the true value. Notably, the algorithm provided a high-confidence estimate of $k_3 = 0$ with a conditional probability of 0.9 and therefore successfully recognized the correct two-gene-state model of the system. Since k_4 and k_{p2} have no effect on the dynamics of the two-gene-state model, the inference of these parameters is unimportant in this context, and quite expectedly, our algorithm yields conditional distributions close to uniform distributions (the prior distribution). Moreover, the stationary distributions of the target system and the inferred model also match perfectly (see the bottom-left panel in Figure 9B), suggesting the accuracy of the inferred model. To conclude, our approach can accurately identify the dynamical model when the target system has only two gene states.

B. Analysis of experimental data. We investigated the performance of our identification method in experimental data and applied this to understanding the contribution of intrinsic and extrinsic noise to the cell-to-cell dynamical variability.

[‡] Here by stability we mean that the underlying continuous-time Markov chain is ergodic (72). One of consequences of this ergodicity

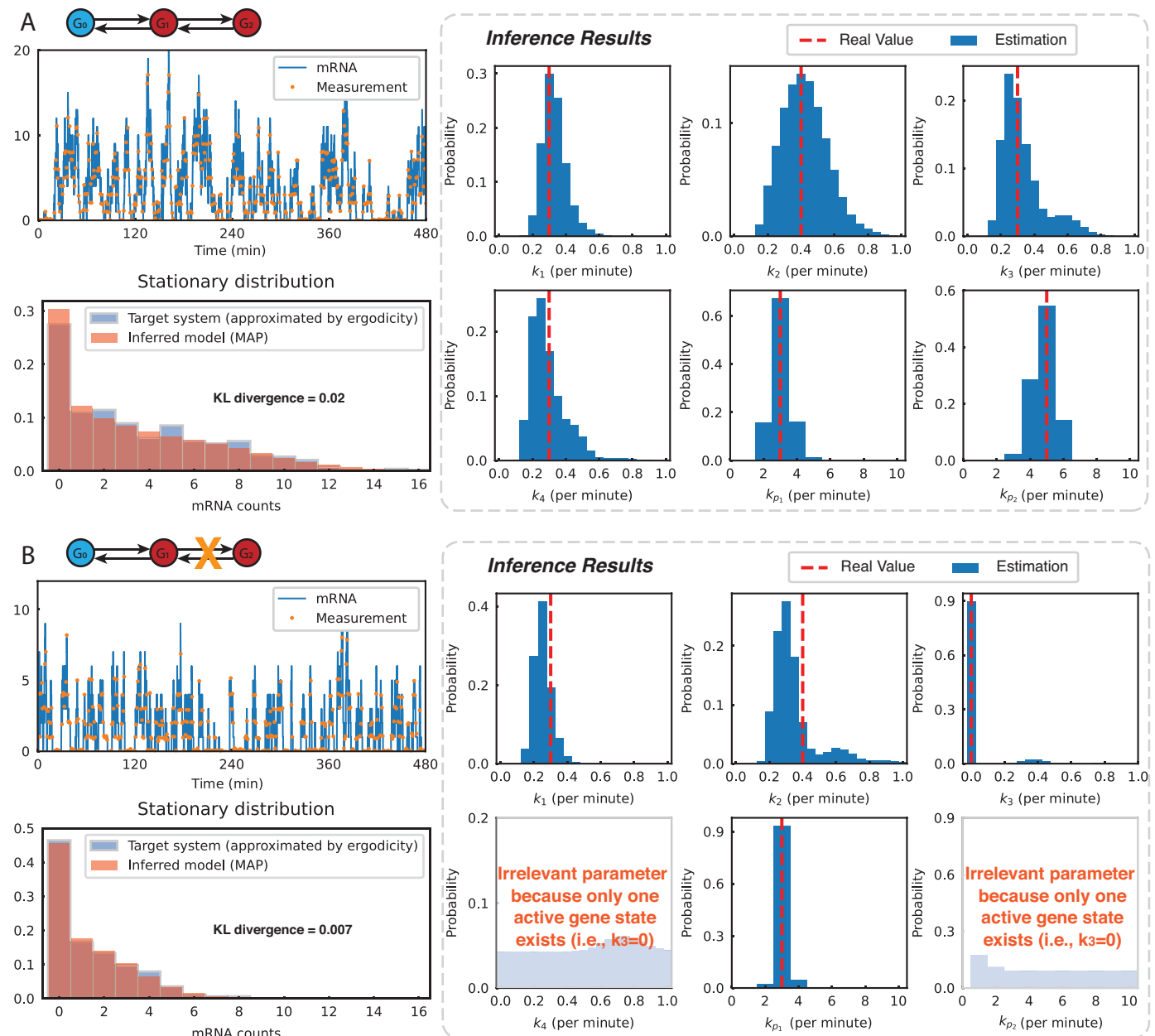


Fig. 9. Performance of the Rao-Blackwell identification algorithm in numerical examples.

(A) Inference of a 3-gene-state system. We first simulated a 3-gene-state system (see Figure 8B) with parameters $k_1 = 0.3$, $k_2 = 0.4$, $k_3 = 0.3$, $k_4 = 0.3$, $k_{p1} = 3$, and $k_{p2} = 5$ (all in units per minute); its mRNA dynamics and time-course measurements are depicted in the top-left panel. Next, we utilized our algorithm to identify the model parameters using the simulated measurements. The sample size of our algorithm was set to 10,000, and the result is presented in the box surrounded by the dash lines. The results illustrate that our algorithm accurately infers the hidden model parameters. The bottom-left panel compares the stationary distributions of the target system and the inferred model (with parameters being the maximum a posteriori estimates). Due to ergodicity, the stationary distribution of the target system was approximated by the occupation time distribution of the mRNA measurements. This bottom-left panel shows a close match between the two stationary distributions, with a KL divergence of 0.02, suggesting the accuracy of the inferred model.

(B) Inference of a 2-gene-state system. We tested whether our approach could accurately identify the gene circuit when the system had only two gene states. We first simulated the same system with parameters $k_1 = 0.3$, $k_2 = 0.4$, $k_3 = 0$, $k_4 = 0$, $k_{p1} = 3$, and $k_{p2} = 0$ (all in units per minute), i.e., the system had only one active gene state. The mRNA dynamics and its time-course measurements are presented in the top-left panel. Then, we used our algorithm to identify the model parameters with the sample size set to 10,000; the result is presented in the box surrounded by the dash lines. The results illustrate that our algorithm accurately infers the parameters k_1 , k_2 , k_3 , and k_{p1} ; as the maximum a posteriori estimate of k_3 is zero, our algorithm correctly identifies the two-gene-state model. Since k_4 and k_{p2} do not affect on the dynamics of the two-gene-state model, the inference of these parameters is unimportant, and our algorithm gives conditional distributions close to uniform distributions (the prior distribution). The bottom-left panel shows a close match between the stationary distributions of the target system and the inferred model, with a KL divergence of 0.007, suggesting the accuracy of the identification result.

B.1. Model identification for yeast cells. The experiment used the yeast cell constructed in (71). As mentioned in (71), all strains were derived from BY4741 and BY4742 (Euroscarf, Germany). Before the experiments, all the cells were kept in a dark environment, ensuring the gene started in the inactive state and the mRNA count to be initially zero. Then, the cells were placed under a microscope platform developed in (71, 73), and they were stimulated by being exposed to constant light. Subsequently, their mRNA fluorescence was measured every 2 minutes for a total period of 4 hours. Finally, this experimental process resulted in the collection of single-cell time-course data from 130 cells. Due to background noise, the platform can only provide the readout when the mRNA count is greater than 7, i.e., $h(x) = x\mathbb{1}(x > 7)$ (see Eq. [7] for the meaning of $h(\cdot)$). Moreover, the mRNA measurements provided by this platform are considered fairly accurate, though the specific magnitude of the measurement noise has not been exactly quantified in the literature (71, 73). Therefore, we treated the measurement noise intensity σ as 1.

To account for our limited knowledge about the parameter values, we considered large ranges for the values of the parameters. Specifically, we assumed k_1 , k_2 , and k_3 to be within the set $\{0, 0.05, \dots, 1\}$, k_4 within the set $\{0, 0.1, \dots, 2\}$, k_{p_1} within $\{0, 8, 16, \dots, 80\}$, and k_{p_2} within $\{20, 30, 40, \dots, 120\}$. All these parameters were assumed to have uniform prior distributions over their specified range, with the exception of k_3 . For k_3 , we assumed that half of its prior probability mass is concentrated at state zero, reflecting the uncertainty regarding the actual number of gene states; the rest of the mass was uniformly distributed on the remaining values in its range. Also, we truncated the state space for the mRNA count to be $\{0, 1, \dots, 150\}$ and that for each gene (G_0 , G_1 , and G_2) to be $\{0, 1\}$. Moreover, the sample size of our algorithm was set to 3,000, and the size of each follower subsystem was required to be less than 30,000. Under this setting, our algorithm classifies G_1 and G_2 as leader species and assigns the remaining components into four follower subsystems (see Figure 8B for the decomposition). Some of the identification results are presented in Figure 10.

Figure 10A presents the inference results for cell #78. Our algorithm provides sharp estimates for each parameter of the cell system, as evidenced by the relatively narrow conditional distributions. To verify the estimation result, we also compared the stationary distributions of the inferred model and the actual cell system, using the approach introduced in Section 3.A. The result shows that these two distributions both have a bimodal shape and align well with a KL divergence of 0.031 (see the top-middle panel of Figure 10A). We noticed that the stationary distribution is sensitive to the parameters, with the elasticity (defined by $\frac{\partial \mathbb{E}_{\text{st}}[X_{\text{mRNA}}|\theta^{\text{inf}}]}{\partial \theta_i^{\text{inf}}} / \frac{\mathbb{E}_{\text{st}}[X_{\text{mRNA}}|\theta^{\text{inf}}]}{\theta_i^{\text{inf}}}$) for each inferred parameter θ_i^{inf} being 35%, -36%, 32%, -30%, 69%, and 31% for k_1 , k_2 , k_3 , k_4 , k_{p_1} , and k_{p_2} , respectively. Therefore, this consistency between the stationary distributions indicates that the inferred parameters are consistent with the true system parameters. The slight difference between these distributions can be attributed to the imperfect approximation of the stationary distribution (of the actual system) by the occupation time distribution and the imperfect inference result due to the limited length of the mRNA trajectory. Overall, the results indicate the accuracy of the inference result.

The identification result indicates that the parameter k_3

was positive with a high probability and, therefore, suggests that the system had three gene states. To further validate this result, we also inferred the system assuming only two gene states existed. In this identification problem, we set $k_3 = k_{p_2} = 0$, considered k_{p_1} within the set $\{0, 1, \dots, 121\}$, and kept the other settings as before. We also compared the stationary distributions of the inferred model and the actual cell system and found that the two distributions had a big mismatch (see the grey box in Figure 10A). Specifically, in this case, the inferred model fails to capture the bimodal distribution of the actual cell system, and the KL divergence between these two distributions is relatively large, with a value of 0.145 (about five times the value in the previous case). All these results support the validity of the three-gene-state model for the real system.

Figure 10B shows the inference results for some typical cells in the population. These cells displayed a variety of behaviors: some cells spent the majority ($\geq 90\%$) of the time in the inactive gene state, while others spent only half of their time in this state. Some cells were highly active, exhibiting a unimodal stationary distribution with a peak away from the origin, while others showed switch-like dynamics, exhibiting a bimodal stationary distribution. Despite this variability, our algorithm consistently provides accurate inference results for all the cells, as evident by the close match between the stationary distributions of the inferred models and the actual cells. These results demonstrate the effectiveness of our method in identifying dynamical models of real biological cells from their individually measured time-trajectories.

B.2. Analysis of noise decomposition. The inference results show that the cells displayed significant heterogeneity in their system parameters (see Figure 11A). We then employed the inference result to investigate how this heterogeneity affected the variability of mRNA counts in the cell population and whether this is a dominant noise source.

The total variability of mRNA counts (at the stationary probability distribution) across the cell population is $\text{Var}(X_{\text{mRNA}}^*)$. Here, Var is the notation of variance, X_{mRNA}^* represents the mRNA count (at the stationary probability distribution) of a randomly selected cell, and θ is a vector of system parameters of this selected cell. As discussed in (58, 60) we can apply the law of total variance to decompose this variability or *noise* into extrinsic and intrinsic components as

$$\underbrace{\text{Var}(X_{\text{mRNA}}^*)}_{\text{total noise}} = \underbrace{\mathbb{E}[\text{Var}(X_{\text{mRNA}}^*|\theta)]}_{\text{intrinsic noise}} + \underbrace{\text{Var}(\mathbb{E}[X_{\text{mRNA}}^*|\theta])}_{\text{extrinsic noise}}.$$

The first term on the right is the intrinsic noise, because in the conditional variance $\text{Var}(X_{\text{mRNA}}^*|\theta)$, the parameter θ is fixed and so only the noise due to the random firing of reactions contributes to this term. On the other hand, the second term quantifies the extrinsic noise because the conditional expectation $\mathbb{E}[X_{\text{mRNA}}^*|\theta]$ filters out the noise from the random firing of reactions, and hence its variance quantifies the noise-contribution due to the variability in the extrinsic parameter θ .

Direct estimation of these two terms from experimental single-cell data is difficult as computing the conditional expectation and variance, given a fixed θ -value, would require us to have measurements from a population of cells that share

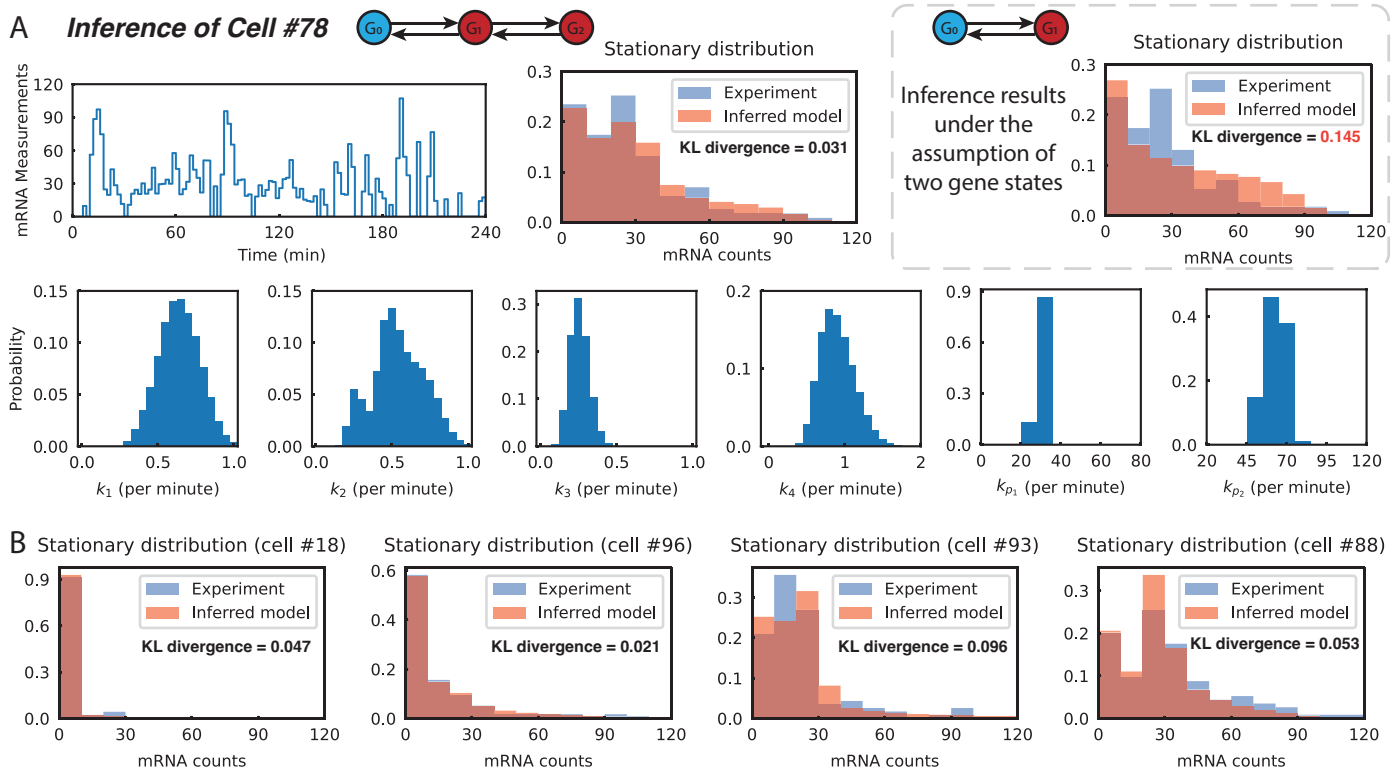


Fig. 10. Performance of the Rao-Blackwell identification algorithm in experimental data from yeast cells.

(A) Inference result of cell #78. The top-left panel shows mRNA dynamics measured every two minutes over a 4-hour period. The lower row presents the parameter estimates, with the narrow conditional probability distributions illustrating the high confidence of these estimates. The top-middle panel compares the stationary distributions of the inferred model (using the MAP estimates) and the actual cell (approximated by the occupation time distribution of mRNA measurements). In these distributions, states have been organized in groups of ten. The good agreement of these two distributions (evident from the bimodal structure and low KL divergence) underscores the accuracy of the identification result.

Our results indicate that k_3 is positive, implying that the system comprises three gene states. We also inferred the system assuming only two gene states; these results are shown in the grey box. In this case, the stationary distribution of the inferred model does not align with the distribution of the real cell. Specifically, the inferred model fails to capture the bimodality of the actual stationary distribution, leading to a relatively substantial KL divergence of 0.145 between the two distributions. This observation further supports the validity of the three-gene-state model for the real cell.

(B) Inference results of some typical cells. We present the inference results of several cells exhibiting different behaviors. Cell #18 took about 90% of the time staying in the inactive gene state; cell #96 had a shorter duration in the inactive gene state; cell #93 took even less time in the inactive state, and its stationary distribution has a peak different from the origin; cell #88 exhibited a bimodal stationary distribution. In all these cases, the stationary distribution of the inferred model closely matches that measured in the experiment in terms of the shape and KL divergence, indicating the validity of our algorithm.

this θ -value. As we cannot have such measurements, the strategy proposed in the current literature to measure the two noise components relies on having dual-reporters within each (58, 60, 61). These two reporters should not only provide conditionally independent measurements given θ , but also these measurements must have the same first two moments[§]. These two conditionals make dual-reporter systems hard to realise, and therefore it is of interest to find other ways to directly estimate noise components from single-cell experimental data without needing such dual-reporter systems.

We now propose a novel approach to measure the two noise components, in the situation where the underlying stochastic dynamics is stable (i.e. ergodic) and we have time-course experimental data for each measured cell. In this case the required conditional expectation and variance (given θ) can be approximated by time-averages over a large interval $[0, T]$, i.e. $\mathbb{E}[X_{\text{mRNA}}^*|\theta] \approx \frac{1}{T} \int_0^T X_{\text{mRNA}}^\theta(t) dt$ and $\text{Var}(X_{\text{mRNA}}^*|\theta) \approx \frac{1}{T} \int_0^T (X_{\text{mRNA}}^\theta(t))^2 dt - (\mathbb{E}[X_{\text{mRNA}}^*|\theta])^2$, where $X_{\text{mRNA}}^\theta(t)$ is the mRNA dynamics of a cell with parameters θ (see SI

Appendix, Section S8 for more details). Therefore extrinsic and intrinsic noise can be estimated as

$$\text{Extrinsic noise} \approx \text{Var} \left(\frac{1}{T} \int_0^T X_{\text{mRNA}}^\theta(t) dt \right) \quad [13]$$

$$\text{Intrinsic noise} \approx \mathbb{E} \left[\frac{1}{T} \int_0^T (X_{\text{mRNA}}^\theta(t))^2 dt - \left(\frac{1}{T} \int_0^T X_{\text{mRNA}}^\theta(t) dt \right)^2 \right]. \quad [14]$$

We employed this method to evaluate the intrinsic and extrinsic noise in the yeast cell population, and the result is presented in Figure 11.B and Figure 11.C (labeled as “experiments”). Our cell-specific identification results provide an indirect inference-based approach for the estimation of the two noise components. Specifically, our inference result provided cell-specific estimates for the parameters, thereby allowing for the computation of $\mathbb{E}[X_{\text{mRNA}}^*|\theta]$ and $\text{Var}(X_{\text{mRNA}}^*|\theta)$ for each cell. Then, by combining these conditional mean and variance, we obtained the estimates to intrinsic and extrinsic noise, as illustrated in Figure 11.C (labeled as “inferred models”).

[§]Under these conditions the covariance between the measurements of the two reporters provides an estimate of the extrinsic noise

The noise decomposition results obtained from both methods are quite consistent (see Figure 11), indicating the accuracy of our identification method and the reliability of the noise decomposition result. The result in Figure 11.C suggests that the extrinsic noise only accounted for a small fraction (about 18%) of the total variation, though the heterogeneity in dynamic parameters was significant. This result is mainly attributed to the fact that the intrinsic noise in this transcription dynamics was extremely strong. Figure 10.A tells that the mRNA dynamics of the 78th cell fluctuated greatly from 0 to 100, with the coefficient of variation (defined as $\frac{\text{standard deviation}}{\text{mean}}$) being almost one. This fluctuation of mRNA dynamics was also observed in other cells, as a result of which the intrinsic noise was around 800 (Figure 11.C). On contrast, the extrinsic noise was only about 170 (Figure 11.C). Overall, in this transcription dynamics, the intrinsic noise dominated the total variation.

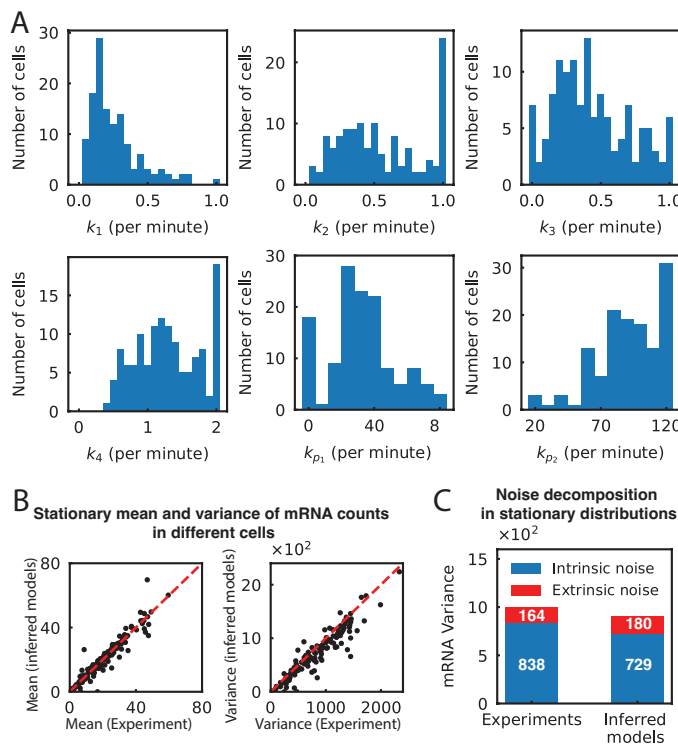


Fig. 11. Noise analysis in gene expression of Yeast cells. (A) Distribution of inferred parameters across the cell population. The plots indicate significant variability in system parameters within the yeast population. (B) Stationary mean and variance of mRNA counts in different cells. Here, we contrast the mean and variance estimates from the inferred models with those estimated from experimental data. Recall that the stationary distribution of each real cell was estimated using the occupation time distribution of mRNA measurements. The results show that our estimates are consistent with the experimental data, with all the dots distributed around the diagonal lines. Also, the results demonstrate that the yeast cells displayed considerable heterogeneity in terms of stationary mean and variance. (C) Noise decomposition in stationary distributions. By the law of total variance, the variance of mRNA counts across the population can be expressed as $\text{Var}(X_{\text{mRNA}}) = \mathbb{E}[\text{Var}(X_{\text{mRNA}}|\theta)] + \text{Var}(\mathbb{E}[X_{\text{mRNA}}|\theta])$, where \mathbb{E} and Var are the notations of mean and variance in the stationary distribution, X_{mRNA} is the mRNA count of a cell randomly selected from the population, and θ is a vector of system parameters of that cell. The quantity $\mathbb{E}[\text{Var}(X_{\text{mRNA}}|\theta)]$ is termed as the intrinsic noise, and $\text{Var}(\mathbb{E}[X_{\text{mRNA}}|\theta])$ is termed as the extrinsic noise. Here, we compare the intrinsic noise and extrinsic noise obtained from the experimental data sets and the inferred models. The plot demonstrates a substantial consistency between the noise decomposition results derived in both ways.

4. Discussion

The modeling and analysis of intracellular reactions, which are subject to inherent randomness in living cells, is typically achieved through formulating the dynamics as a continuous-time Markov chain and finding solutions to the associated Chemical Master Equation (CME). This equation describes the probability evolution of the underlying chemical species abundances and is crucial for studying noisy intracellular reaction systems. However, solving the CME for high-dimensional systems remains a significant challenge. To address this issue, we have developed a new divide-and-conquer approach called the Rao-Blackwellized CME Solver (RB-CME Solver) that combines modularization and stochastic filtering.

Specifically, the RB-CME solver works by an optimized decomposition of the system, which transforms the large-scale problem into several lower scale ones and then solves each of them using either the Monte-Carlo method or a filtering approach (e.g., the filtered FSP). We showed that our method can be remarkably efficient and accurate in analyzing high-dimensional systems. Compared with the traditional Monte-Carlo method, our approach is far more accurate in estimating the follower system and less accurate in estimating the leader system. However, the well-designed system decomposition compensates for the additional costs of the filtering algorithm, making our method favorable for estimating CME solution when compared to the classical Monte Carlo method, given the same computational time-constraints.

We also developed the method for the stochastic filtering problem (named the Rao-Blackwellized particle filter or the RB-PF) and showed its superior performance over conventional approaches. Specifically, the RB-PF utilizes the RB-CME solver to compute the CME in the prediction step of the filtering problem. In RB-PFs, we found that the leader system can also benefit from Rao-Blackwellization, as the conditional probability distribution of the follower interacts with the probability distribution of the follower in the correction step. More interestingly, the RB-CME solver and the RB-PF have relatively the same performance when dealing with the same system. Overall, our method is a powerful tool to solve CMEs and the associated stochastic filtering problems for high-dimensional chemical reaction systems.

Furthermore we extended our RB-PF method for cell-specific model identification by including the cell-specific rate parameters as unobserved species to be inferred, in the sense of conditional probability distribution, from the measurement of the cell's time-trajectory. We successfully applied this technique to an experimental time-course data-set of transcription dynamics in yeast cells, and our results revealed significant variation in rate parameters among isogenic and identically cultured yeast cells. To validate our model inference we showed a close match between the cell-specific stationary distributions computed with our inferred model and estimated through time-averages of the measured trajectory. The parameter-heterogeneity among cells can be viewed as extrinsic noise (57) whose contribution to the overall cell-to-cell variability can be easily estimated from our cell-specific inferred models, along with the corresponding intrinsic noise that is due to the random firing of reactions. To estimate these two noise components directly from time-course experimental data, we proposed a novel noise decomposition approach that relies on the stability of the underlying stochastic dynamics, and does not require

complex dual-reporter systems (58, 60, 61) that are traditionally employed to perform this noise separation. We illustrate that both the model-based and the direct approach provide comparable estimates for the two noise components. This concordance further emphasizes the validity and real-world applicability of our cell-specific inference method.

Since the CME plays the central role in the stochastic analysis of biochemical reaction systems, our method can also be generalized to consider other computational problems in biology. For instance, it is worth extending our method for the parameter sensitivity analysis and power spectrum analysis of stochastic reaction systems, which could also be computationally demanding in high-dimensional cases (12, 74–76). Also, our method might be helpful in estimating the probability of rare events, which can be difficult to achieve using conventional importance sampling approaches (77). Moreover, the divide-and-conquer idea used in our method can also be applied to computing other equations that capture the probability evolution of stochastic systems, e.g., the Fokker-Planck equation.

Finally, there is room for improvement with our method in terms of generality and computational efficiency. To tackle more complex systems, our method needs further development to incorporate additional biological mechanisms such as time delays, multiscale phenomena, and cell-to-cell interactions. Additionally, it is worth exploring ways to enhance the performance of the RB-CME solver, such as incorporating deep learning methods (45–47) to our divide-and-conquer approach.

Materials and Methods

Mathematical derivations and additional information are provided in [SI Appendix](#).

ACKNOWLEDGMENTS. We acknowledge funding from the Swiss National Science Foundation under grant 182653.

1. A Adan, G Alizada, Y Kiraz, Y Baran, A Nalbant, Flow cytometry: basic principles and applications. *Critical reviews biotechnology* **37**, 163–176 (2017).
2. DJ Stephens, VJ Allan, Light microscopy techniques for live cell imaging. *science* **300**, 82–86 (2003).
3. C Vonesch, F Aguet, JL Vonesch, M Unser, The colored revolution of bioimaging. *IEEE signal processing magazine* **23**, 20–31 (2006).
4. J Zhang, RE Campbell, AY Ting, RY Tsien, Creating new fluorescent probes for cell biology. *Nat. reviews Mol. cell biology* **3**, 906–918 (2002).
5. HH McAdams, A Arkin, Stochastic mechanisms in gene expression. *Proc. Natl. Acad. Sci.* **94**, 814–819 (1997).
6. A Arkin, J Ross, HH McAdams, Stochastic kinetic analysis of developmental pathway bifurcation in phase λ -infected escherichia coli cells. *Genetics* **149**, 1633–1648 (1998).
7. N Fedoroff, W Fontana, Small numbers of big molecules. *Science* **297**, 1129–1131 (2002).
8. DF Anderson, TG Kurtz, *Stochastic analysis of biochemical systems*. (Springer) Vol. 674, (2015).
9. B Munsky, B Trinh, M Khammash, Listening to the noise: random fluctuations reveal gene network parameters. *Mol. systems biology* **5**, 318 (2009).
10. I Golding, J Paulsson, SM Zawilski, EC Cox, Real-time kinetics of gene activity in individual bacteria. *Cell* **123**, 1025–1036 (2005).
11. C Zechner, M Unger, S Pelet, M Peter, H Koepl, Scalable inference of heterogeneous reaction kinetics from pooled single-cell recordings. *Nat. methods* **11**, 197–202 (2014).
12. A Gupta, M Khammash, Frequency spectra and the color of cellular noise. *Nat. Commun.* **13**, 1–18 (2022).
13. A Bain, D Crisan, *Fundamentals of stochastic filtering*. (Springer Science & Business Media) Vol. 60, (2008).
14. C Briat, M Khammash, Noise in biomolecular systems: Modeling, analysis, and control implications. *Annu. Rev. Control. Robotics, Auton. Syst.* **6**, 283–311 (2023).
15. E Korobkova, T Emonet, JM Vilar, TS Shimizu, P Cluzel, From molecular noise to behavioural variability in a single bacterium. *Nature* **428**, 574–578 (2004).
16. S Kar, WT Baumann, MR Paul, JJ Tyson, Exploring the roles of noise in the eukaryotic cell cycle. *Proc. Natl. Acad. Sci.* **106**, 6471–6476 (2009).
17. R Perez-Carrasco, C Beentjes, R Grima, Effects of cell cycle variability on lineage and population measurements of messenger rna abundance. *J. Royal Soc. Interface* **17**, 20200360 (2020).
18. MB Elowitz, S Leibler, A synthetic oscillatory network of transcriptional regulators. *Nature* **403**, 335–338 (2000).
19. SK Aoki, et al., A universal biomolecular integral feedback controller for robust perfect adaptation. *Nature* **570**, 533–537 (2019).
20. G Neuert, et al., Systematic identification of signal-activated stochastic gene regulation. *Science* **339**, 584–587 (2013).
21. Z Fang, A Gupta, M Khammash, Stochastic filtering for multiscale stochastic reaction networks based on hybrid approximations. *J. Comput. Phys.* **467**, 111441 (2022).
22. T Jahnke, W Huisinga, Solving the chemical master equation for monomolecular reaction systems analytically. *J. mathematical biology* **54**, 1–26 (2007).
23. JJ Vastola, Solving the chemical master equation for monomolecular reaction systems and beyond: a doi-peliti path integral view. *J. Math. Biol.* **83**, 1–82 (2021).
24. Z Cao, R Grima, Analytical distributions for detailed models of stochastic gene expression in eukaryotic cells. *Proc. Natl. Acad. Sci.* **117**, 4682–4692 (2020).
25. V Shahrezaei, PS Swain, Analytical distributions for stochastic gene expression. *Proc. Natl. Acad. Sci.* **105**, 17256–17261 (2008).
26. DT Gillespie, A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *J. computational physics* **22**, 403–434 (1976).
27. DT Gillespie, Exact stochastic simulation of coupled chemical reactions. *The journal physical chemistry* **81**, 2340–2361 (1977).
28. DT Gillespie, Approximate accelerated stochastic simulation of chemically reacting systems. *The J. chemical physics* **115**, 1716–1733 (2001).
29. Y Cao, DT Gillespie, LR Petzold, Efficient step size selection for the tau-leaping simulation method. *The J. chemical physics* **124**, 044109 (2006).
30. M Rathinam, LR Petzold, Y Cao, DT Gillespie, Stiffness in stochastic chemically reacting systems: The implicit tau-leaping method. *The J. Chem. Phys.* **119**, 12784–12794 (2003).
31. EL Haseltine, JB Rawlings, Approximate simulation of coupled fast and slow reactions for stochastic chemical kinetics. *The J. chemical physics* **117**, 6959–6969 (2002).
32. HW Kang, TG Kurtz, Separation of time-scales and model reduction for stochastic reaction networks. *The Annals Appl. Probab.* **23**, 529–583 (2013).
33. B Munsky, M Khammash, The finite state projection algorithm for the solution of the chemical master equation. *The J. chemical physics* **124**, 044104 (2006).
34. V Kazeze, M Khammash, M Nip, C Schwab, Direct solution of the chemical master equation using quantized tensor trains. *PLoS computational biology* **10**, e1003359 (2014).
35. IG Ion, C Wildner, D Loukrezis, H Koepl, H De Gersem, Tensor-train approximation of the chemical master equation and its application for parameter inference. *The J. Chem. Phys.* **155**, 034102 (2021).
36. NG Van Kampen, *Stochastic processes in physics and chemistry*. (Elsevier) Vol. 1, (1992).
37. L Bronstein, H Koepl, A variational approach to moment-closure approximations for the kinetics of biomolecular reaction networks. *The J. chemical physics* **148**, 014105 (2018).
38. CA Gomez-Urbe, GC Verghese, Mass fluctuation kinetics: Capturing stochastic effects in systems of chemical reactions through coupled mean-variance computations. *The J. chemical physics* **126**, 024109 (2007).
39. I Näsell, An extension of the moment closure method. *Theor. population biology* **64**, 233–239 (2003).
40. J Hespanha, Moment closure for biochemical networks in 2008 3rd International Symposium on Communications, Control and Signal Processing. (IEEE), pp. 142–147 (2008).
41. MJ Keeling, Multiplicative moments and measures of persistence in ecology. *J. Theor. Biol.* **205**, 269–281 (2000).
42. A Singh, JP Hespanha, Lognormal moment closures for biochemical reactions in *Proceedings of the 45th IEEE Conference on Decision and Control*. (IEEE), pp. 2063–2068 (2006).
43. P Smadbeck, YN Kaznessis, A closure scheme for chemical master equations. *Proc. Natl. Acad. Sci.* **110**, 14261–14265 (2013).
44. J Ruess, A Milias-Areitis, S Summers, J Lygeros, Moment estimation for chemically reacting systems by extended kalman filtering. *The J. chemical physics* **135**, 10B621 (2011).
45. A Gupta, C Schwab, M Khammash, Deepcme: A deep learning framework for computing solution statistics of the chemical master equation. *PLoS computational biology* **17**, e1009623 (2021).
46. Q Jiang, et al., Neural network aided approximation and parameter inference of non-markovian models of gene expression. *Nat. communications* **12**, 1–12 (2021).
47. Y Tang, J Weng, P Zhang, Neural-network solutions to stochastic reaction networks. *Nat. Mach. Intell.* **5**, 376–385 (2023).
48. J Goutsias, Quasiequilibrium approximation of fast reaction kinetics in stochastic biochemical systems. *The J. chemical physics* **122**, 184102 (2005).
49. Y Cao, DT Gillespie, LR Petzold, The slow-scale stochastic simulation algorithm. *The J. chemical physics* **122**, 014116 (2005).
50. LD EW, E Vanden-Eijnden, Nested stochastic simulation algorithm for chemical kinetic systems with disparate rates. *J Chem Phys* **123**, 194107 (2005).
51. J Hasenauer, V Wolf, A Kazerooni, FJ Theis, Method of conditional moments (mcm) for the chemical master equation. *J. mathematical biology* **69**, 687–735 (2014).
52. C Zechner, H Koepl, Uncoupled analysis of stochastic reaction networks in fluctuating environments. *PLoS computational biology* **10**, e1003942 (2014).
53. L Duso, C Zechner, Selected-node stochastic simulation algorithm. *The J. chemical physics* **148**, 164108 (2018).
54. CR Rao, Information and the accuracy attainable in the estimation of statistical parameters in *Breakthroughs in statistics*. (Springer), pp. 235–247 (1992).
55. D Blackwell, Conditional expectation and unbiased sequential estimation. *The Annals Math. Stat.* pp. 105–110 (1947).
56. ES D’Ambrosio, Z Fang, A Gupta, M Khammash, Filtered finite state projection method for the analysis and estimation of stochastic biochemical reaction networks. *bioRxiv* (2022).
57. C Zechner, et al., Moment-based inference predicts bimodality in transient gene expression. *Proc. Natl. Acad. Sci.* **109**, 8340–8345 (2012).
58. PS Swain, MB Elowitz, ED Siggia, Intrinsic and extrinsic contributions to stochasticity in gene expression. *Proc. Natl. Acad. Sci.* **99**, 12795–12800 (2002).

59. MB Elowitz, AJ Levine, ED Siggia, PS Swain, Stochastic gene expression in a single cell. *Science* **297**, 1183–1186 (2002).
60. A Hilfinger, J Paulsson, Separating intrinsic from extrinsic fluctuations in dynamic biological systems. *Proc. Natl. Acad. Sci.* **108**, 12167–12172 (2011).
61. CG Bowsher, PS Swain, Identifying sources of variation and the flow of information in biochemical networks. *Proc. Natl. Acad. Sci.* **109**, E1320–E1328 (2012).
62. JM Raser, EK O'Shea, Control of stochasticity in eukaryotic gene expression. *science* **304**, 1811–1814 (2004).
63. C Zechner, G Seelig, M Rullan, M Khammash, Molecular circuits for dynamic noise filtering. *Proc. Natl. Acad. Sci.* **113**, 4729–4734 (2016).
64. M Rathinam, M Yu, State and parameter estimation from exact partial state observation in stochastic reaction networks. *The J. Chem. Phys.* **154**, 034103 (2021).
65. D Crisan, Particle filters—a theoretical perspective in *Sequential Monte Carlo methods in practice*. (Springer), pp. 17–41 (2001).
66. Z Fang, A Gupta, M Khammash, Stochastic filters based on hybrid approximations of multiscale stochastic reaction networks in *2020 59th IEEE Conference on Decision and Control (CDC)*. (IEEE), pp. 4616–4621 (2020).
67. Z Fang, A Gupta, M Khammash, Convergence of regularized particle filters for stochastic reaction networks. *SIAM J. on Numer. Analysis* **61**, 399–430 (2023).
68. J Liu, M West, Combined parameter and state estimation in simulation-based filtering in *Sequential Monte Carlo methods in practice*. (Springer), pp. 197–223 (2001).
69. KR Sanft, HG Othmer, Constant-complexity stochastic simulation algorithm with optimal binning. *The J. chemical physics* **143**, 074108 (2015).
70. TS Gardner, CR Cantor, JJ Collins, Construction of a genetic toggle switch in *escherichia coli*. *Nature* **403**, 339–342 (2000).
71. M Rullan, D Benzinger, GW Schmidt, A Milias-Areitis, M Khammash, An optogenetic platform for real-time, single-cell interrogation of stochastic transcriptional regulation. *Mol. cell* **70**, 745–756 (2018).
72. A Gupta, C Briat, M Khammash, A scalable computational framework for establishing long-term behavior of stochastic reaction networks. *PLoS computational biology* **10**, e1003669 (2014).
73. S Kumar, M Rullan, M Khammash, Rapid prototyping and design of cybergenetic single-cell controllers. *Nat. communications* **12**, 5651 (2021).
74. DF Anderson, An efficient finite difference method for parameter sensitivities of continuous time markov chains. *SIAM J. on Numer. Analysis* **50**, 2237–2258 (2012).
75. P Dürrenberger, A Gupta, M Khammash, A finite state projection method for steady-state sensitivity analysis of stochastic reaction networks. *The J. Chem. Phys.* **150**, 134101 (2019).
76. A Gupta, M Khammash, Unbiased estimation of parameter sensitivities for stochastic chemical reaction networks. *SIAM J. on Sci. Comput.* **35**, A2598–A2620 (2013).
77. M Ahmadi, et al., A comparison of weighted stochastic simulation methods for the analysis of genetic circuits. *ACS Synth. Biol.* (2022).