



HAL
open science

Parameter inference for stochastic biochemical models from perturbation experiments parallelised at the single cell level

Andela Davidović, Remy Chait, Gregory Batt, Jakob Ruess

► **To cite this version:**

Andela Davidović, Remy Chait, Gregory Batt, Jakob Ruess. Parameter inference for stochastic biochemical models from perturbation experiments parallelised at the single cell level. 2021. hal-03544059

HAL Id: hal-03544059

<https://inria.hal.science/hal-03544059>

Preprint submitted on 28 Jan 2022

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

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



Distributed under a Creative Commons Attribution 4.0 International License

1 Parameter inference for stochastic biochemical models
2 from perturbation experiments parallelised at the single
3 cell level

4 Andela Davidović¹, Remy Chait^{2,3}, Gregory Batt^{1,4} and Jakob Ruess^{1,4*}

¹ Department of Computational Biology, Institut Pasteur, Paris, France

² Biosciences, Living Systems Institute, University of Exeter, Exeter, The United Kingdom

³ Institute of Science and Technology Austria, Klosterneuburg, Austria

⁴ Inria Paris, Paris, France

5 **Abstract**

6 Understanding and characterising biochemical processes inside single cells requires
7 experimental platforms that allow one to perturb and observe the dynamics of such
8 processes as well as computational methods to build and parameterise models from
9 the collected data. Recent progress with experimental platforms and optogenetics has
10 made it possible to expose each cell in an experiment to an individualised input and
11 automatically record cellular responses over days with fine time resolution. However,
12 methods to infer parameters of stochastic kinetic models from single-cell longitudinal
13 data have generally been developed under the assumption that experimental data is
14 sparse and that responses of cells to at most a few different input perturbations can
15 be observed. Here, we investigate and compare different approaches for calculating
16 parameter likelihoods of single-cell longitudinal data based on approximations of the
17 chemical master equation (CME) with a particular focus on coupling the linear noise
18 approximation (LNA) or moment closure methods to a Kalman filter. We show that,

*corresponding author: jakob.ruess@inria.fr

19 as long as cells are measured sufficiently frequently, coupling the LNA to a Kalman
20 filter allows one to accurately approximate likelihoods and to infer model parameters
21 from data even in cases where the LNA provides poor approximations of the CME.
22 Furthermore, the computational cost of filtering-based iterative likelihood evaluation
23 scales advantageously in the number of measurement times and different input
24 perturbations and is thus ideally suited for data obtained from modern experimental
25 platforms. To demonstrate the practical usefulness of these results, we perform
26 an experiment in which single cells, equipped with an optogenetic gene expression
27 system, are exposed to various different light-input sequences and measured at
28 several hundred time points and use parameter inference based on iterative likelihood
29 evaluation to parameterise a stochastic model of the system.

30 **Author summary**

31 A common result for the modelling of cellular processes is that available data is not
32 sufficiently rich to uniquely determine the biological mechanism or even just to ensure
33 identifiability of parameters of a given model. Perturbing cellular processes with informative
34 input stimuli and measuring dynamical responses may alleviate this problem. With the
35 development of novel experimental platforms, we are now in a position to parallelise such
36 perturbation experiments at the single cell level. This raises a plethora of new questions. Is
37 it more informative to diversify input perturbations but to observe only few cells for each
38 input or should we rather ensure that many cells are observed for only few inputs? How
39 can we calculate likelihoods and infer parameters of stochastic kinetic models from data
40 sets in which each cell receives a different input perturbation? How does the computational
41 efficiency of parameter inference methods scale with the number of inputs and the number
42 of measurement times? Are there approaches that are particularly well-suited for such
43 data sets? In this paper, we investigate these questions using the CcaS/CcaR optogenetic
44 system driving the expression of a fluorescent reporter protein as primary case study.

45 **1 Introduction**

46 Finding appropriate mathematical models to represent biological processes inside cells is
47 one of the core challenges of computational biology. In particular, the last decade has seen
48 much work focused on the development of methods to infer parameters of mechanistic
49 models of biochemical reaction networks from experimental data [12, 16, 24, 15]. A common
50 premise in most of these works is that models are complex and unknown parameters
51 numerous, with available data being sparse and noisy [11]. While the complexity of
52 biological systems has not changed, the availability and reliability of data certainly has.
53 For instance, experimental techniques and platforms nowadays allow us to observe the
54 onset of transcription at single molecule precision, to fully automatically measure the
55 expression levels of genes every couple of minutes, to perturb and drive gene expression
56 in populations or individual cells [20, 17, 29, 6, 1], and even to let computational models
57 interact with single cell gene expression processes in real time [2]. To which degree these
58 new capacities will eventually allow us to resolve the ill-posedness of reverse engineering
59 models of biological systems from experimental data remains to be clarified, but in any case
60 the availability of new types of data calls for new methods that are capable of exploiting
61 the data in its entirety.

62 In this paper, we focus on data obtained from microscopy platforms such as the ones
63 presented in [2, 29, 6]. The key feature of these platforms is that microscopy-based
64 observation of gene expression dynamics in single cells is coupled to an optical system
65 (such as a digital micromirror device) that allows the user to target light signals at
66 individual cells. These signals are freely designable and can be time-varying. Cells are in
67 turn equipped with an optogenetic system driving the biological process of interest. In
68 summary, the platforms allow one to perturb the biological process dynamically and in a
69 different way in every cell and to observe possibly very different behaviours of the system
70 of interest at the same time. In other words, these new platforms allow us to parallelise a
71 large number of classical perturbation experiments within a single experiment.

72 It has been demonstrated in the past that biochemical reactions inside single cells are
73 inherently stochastic [18, 19]. Faithfully capturing single cell microscopy data therefore

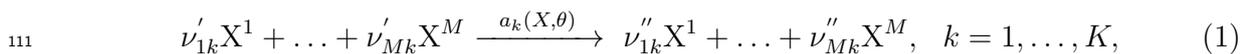
74 requires stochastic kinetic models governed by the chemical master equation (CME) [9].
75 However, inferring parameters of such models from available data is a formidable task
76 [35, 36]. Solving the CME analytically is rarely possible and numerically approximating its
77 solution based on approaches such as finite state projection [23] or stochastic simulation
78 [8] is computationally prohibitively expensive in most cases. This problem is rendered all
79 the worse if each cell in the experiment is perturbed with a different input stimulation,
80 such that the CME needs to be re-solved for every cell in the experiment. Efficient
81 approximations of the CME, such as the linear noise approximation (LNA) [14], may
82 alleviate the computational burden. However, they may become imprecise if cells are
83 observed over longer time horizons as well as impractical to use when the number of
84 measurement times becomes large and the dimensionality of the data increases. A more
85 feasible alternative may then be to evaluate likelihoods iteratively, for instance by coupling
86 the LNA to a Kalman filter or similar [3, 21, 38, 37, 5]. (See Figure 1b and Section 2.3)
87 In this paper, we couple general closure methods for approximately calculating lower order
88 moments of the CME [32] to a Kalman filter in order to approximate likelihoods of single
89 cell data. We show that the frequency at which measurements are taken determines the
90 quality of the likelihood approximation. We argue that our approach is ideally suited
91 for data obtained from parallelised single cell experiments in which the collection of
92 measurements is automatically performed every couple of minutes. We find that, if the
93 system is observed with sufficient time resolution, iterative evaluation of likelihoods allows
94 for quite precise inference of parameters even if the system is highly non-linear and the used
95 approximation of the CME is very imprecise over longer time-horizons (for instance if the
96 LNA is applied to a model of a genetic toggle switch that displays bimodal distributions).
97 Finally, to highlight the practical usefulness of the approach, we perform single-cell
98 experiments to characterise a light-inducible gene expression system. Using simulated
99 data, we show that coupling the LNA to a Kalman filter to iteratively evaluate likelihoods
100 leads to accurate parameter estimation and orders of magnitude less computation time
101 compared to a standard use of the LNA [14]. This efficiency allows us to perform Bayesian
102 parameter inference for various data sets pooled together from single cells perturbed with

103 different light stimuli and to study how much information about model parameters can be
 104 gained by parallelising perturbations in single cells. Using experimental data, we show
 105 that parameter estimation can also be effective but might require a proper treatment of
 106 non-identifiability.

107 2 Methodology

108 2.1 Stochastic biochemical reaction networks

109 Consider a reaction network of M chemical species X^1, \dots, X^M that interact stochastically
 110 according to K reactions



112 parameterised by $\theta = [\theta_1 \dots \theta_c]^\top$ and where the coefficients ν'_{ik} and ν''_{ik} determine how
 113 many molecules of the i -th species are consumed and produced in the k -th reaction,
 114 respectively. If system dynamics are influenced by an input perturbation, $u(t)$, the
 115 propensity functions $a_k(X, \theta) = a_k(X, \theta, u(t))$ will additionally depend on this input.
 116 Under the assumption that the system is well-stirred and in thermal equilibrium, the
 117 probability distribution describing the time evolution of the number of molecules of the
 118 different species is governed by the chemical master equation (CME) [9]:

$$119 \quad \dot{p}(x, t) = -p(x, t) \sum_{k=1}^K a_k(x, \theta) + \sum_{k=1}^K a_k(x - \nu_k, \theta) p(x - \nu_k, t), \quad (2)$$

120 where $x = [x^1 \dots x^M]^\top$ in \mathbb{N}^M is a possible state of the stochastic process $X(t) =$
 121 $[X^1(t) \dots X^M(t)]^\top$ that counts the numbers of molecules of all species, $p(x, t) := \mathbb{P}(X(t) = x)$,
 122 and $\nu_k = \nu''_k - \nu'_k$, where $\nu'_k = [\nu'_{1k} \dots \nu'_{Mk}]^\top$, and $\nu''_k = [\nu''_{1k} \dots \nu''_{Mk}]^\top$. Since the CME
 123 is difficult to solve in most cases, a widely used approach is to derive moment equations
 124 from it. However, except if the propensity functions $a_k(x, \theta)$ of all reactions are linear in
 125 x , the time evolution of moments up to any order depends on moments of higher order
 126 and the moment equations cannot be solved exactly and need to be approximated using
 127 moment closure methods (see Supporting Information Section A).

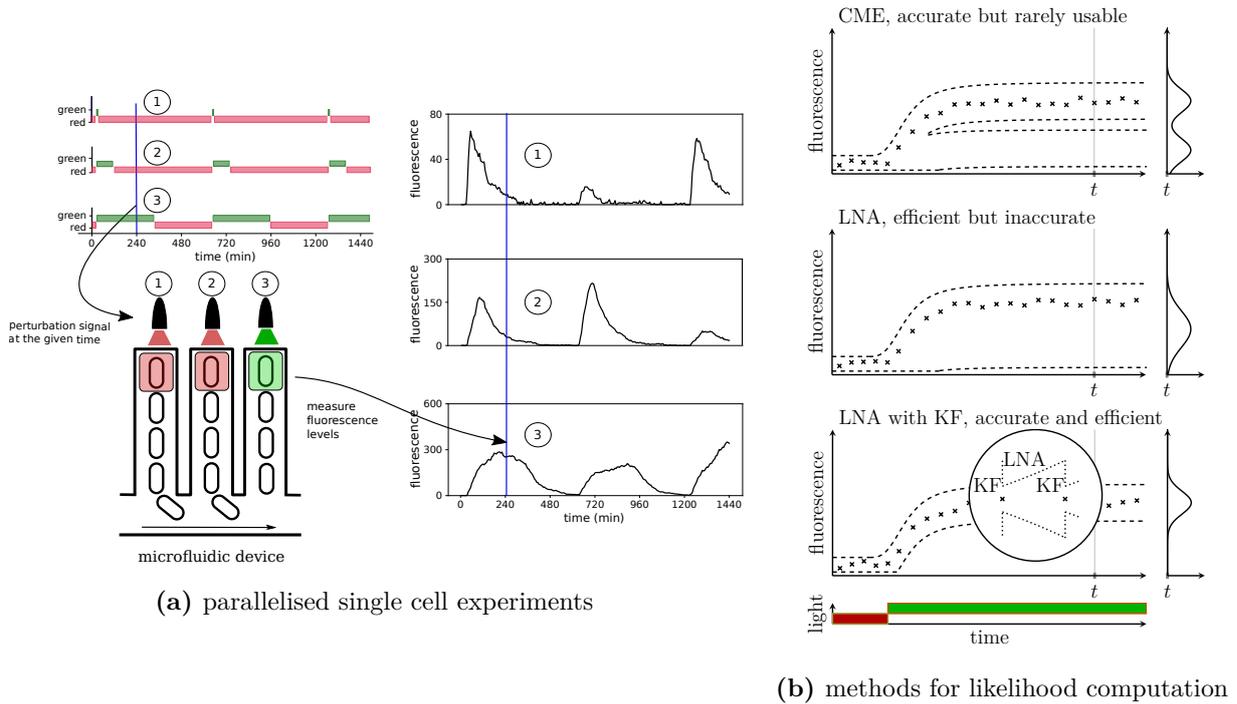


Figure 1: (a) Schematic representation of experiments parallelised at the cell level.

Top left: an example of different light patterns that drive protein production in respective cells. Green light is used to excite the light system in a single cell and to start the production of a fluorescent protein, and red light is used to stop production. Bottom left: representation of a microfluidic device called a mother machine: each microchannel is of the width of a single cell. Cells are kept in a constant environment and are growing and dividing. One cell stays at the bottom of a channel, while the others are being washed away. Each channel is used to obtain longitudinal data at the single cell level. Each channel is then given a light pattern that can be pre-defined (as in top left) or adjusted in real-time in response to incoming data. For more details see [2]. Right: fluorescence data corresponding to the light stimulation profiles at top left. **(b) Benefit of coupling the LNA to Kalman filtering for parameter inference of biochemical reactions.**

Top: The chemical master equation can be used to calculate exact likelihoods but is rarely solvable in practice. An example situation for a system that exhibits bimodality is shown (crosses: data, dashed lines: the CME solution is visualised via confidence regions). Middle: The LNA provides efficient approximations but is based on often very inaccurate Gaussian approximations of the full CME solution. Bottom: Coupling the LNA with Kalman filtering (Algorithm 1) requires only transition probabilities over short time intervals to be close to Gaussian.

128 2.2 Data likelihoods and their classical approximation

129 Here, we consider microscopy data, y^{all} , that contains information from n cells. The data
130 for a single cell is given by $y = \{y_1, \dots, y_S\}$ and described by the measurement model

$$131 \quad y_s = C \cdot X(t_s) + \xi_s, \quad s = 1, \dots, S, \quad (3)$$

132 where the ξ_s are independent technical measurement errors assumed to be Gaussian
133 $\xi_s \sim \mathcal{N}(0, \sigma^2)$, the matrix C maps the full system state to the measured output species,
134 and distances between measurement times, $t_{s+1} - t_s = t_m$, $s = 1, \dots, S - 1$, are assumed
135 to be all equal for the sake of simplicity.

136 If cells do not physically interact, the full data likelihood of the data y^{all} is then the
137 product of single cell likelihoods, $p(y^{\text{all}} | u, \theta) = \prod_{y \in y^{\text{all}}} p(y | u, \theta)$, where

$$138 \quad p(y | u, \theta) = p(y_1, \dots, y_S | u, \theta), \quad (4)$$

139 and $u = u(t)$ is the time-varying input perturbation to which the cell has been exposed.
140 The S -dimensional single cell likelihoods are determined by the distribution of the technical
141 errors ξ_s , as well as the joint distribution of the system state over all measurement time
142 points,

$$143 \quad p(x | u, \theta) = p(x_1, \dots, x_S | u, \theta) := \mathbb{P}(X(t_1) = x_1, \dots, X(t_S) = x_S | u, \theta), \quad (5)$$

144 where $x = \{x_1, \dots, x_S\}$ and $x_s = [x_s^1 \ \dots \ x_s^M]^\top$ are possible full states of the reaction
145 network at times t_s , $s = 1, \dots, S$. It is important to point out that for microscopy
146 data where individual cells are followed in time, the likelihood does not factorise over
147 time points as would be the case for other data types such as flow cytometry data [35].
148 The fundamental difficulty of any likelihood-based parameter inference scheme is that
149 evaluating (4) for given parameters, θ , and given input, u , requires one to calculate
150 $p(x | u, \theta)$ in (5), which is only possible when the CME is tractable.

151 An approach, termed *classical LNA-based inference* in the following, that has experienced
152 quite some popularity in recent years has therefore been to replace the exact likelihood
153 in (4) by an approximation derived from the linear noise approximation [13, 14, 33, 7].

154 The linear noise approximation, contrary to typical moment closure methods, does not
155 only provide approximations of moments at given time points t , but also automatically
156 approximates any inter-time distributions by Gaussian distributions. These distributions
157 can readily be calculated from the model. In sum, this approach enables evaluation of (5)
158 (and (4)) through an approximation (see [14] for details)

$$159 \quad p(x \mid u, \theta) \approx \mathcal{N}(\mu(u, \theta), \Sigma(u, \theta)). \quad (6)$$

160 While the LNA as such is computationally very efficient, it is important to point out that
161 likelihood evaluations may still become very costly for parallelised single cell experiments:
162 if cells can be tracked and measured automatically over very long time horizons (as
163 in [2]), the number of measurement time points, S , can range in the order of several
164 hundreds, implying that the distributions in (6) are extremely high dimensional and that
165 the calculation of the covariance matrix, $\Sigma(u, \theta)$, requires a significant computational effort
166 despite the otherwise efficient LNA. This problem becomes ever more severe considering
167 that $\Sigma(u, \theta)$ needs to be recalculated for every different input perturbation, u , that is
168 used in an experiment and for every point, θ , in parameter space that is explored during
169 parameter inference. In addition, it is long known that the approximation of $p(x \mid u, \theta)$
170 by a Gaussian distribution is very inaccurate or even completely useless in many cases
171 because the system dynamics are highly non-linear and the probability distribution of the
172 system is far from Gaussian. Typically, the LNA should not be used for systems such as
173 genetic toggle switches since the approximation of the real bimodal system distributions
174 with Gaussian distributions is unacceptably imprecise. One aim of this paper is to show
175 that the LNA, and other moment approximation methods, can nevertheless sometimes be
176 used for parameter inference of systems like genetic toggle switches.

177 **2.3 Approximating likelihoods using the LNA with filtering**

178 An alternative way of calculating the likelihood can be obtained if we rewrite the joint
179 probability distribution in (4) in terms of transition probabilities between measurement

180 time points as

$$181 \quad p(y_1, \dots, y_S) = p(y_1) \cdot p(y_2 | y_1) \cdots p(y_S | y_{S-1}, \dots, y_1), \quad (7)$$

182 where we have omitted the dependence on u and θ for the sake of readability. The transition
 183 probabilities between data points appearing on the right hand side can in turn be rewritten
 184 in terms of the technical noise distribution, transition probabilities between full states of
 185 the system, and posterior distributions over the system state given past observations:

$$186 \quad p(y_s | y_{s-1}, \dots, y_1) = \int \int \underbrace{p(y_s | x_s)}_{\text{tech. noise}} \cdot \underbrace{p(x_s | x_{s-1})}_{\text{trans. prob.}} \cdot \underbrace{p(x_{s-1} | y_{s-1}, \dots, y_1)}_{\text{state posterior}} dx_s dx_{s-1}. \quad (8)$$

187 While exactly evaluating (8) is clearly difficult and requires the solution of the CME due
 188 to the appearance of the transition probabilities $p(x_s | x_{s-1})$, it is important to note that
 189 reasonable approximations can be obtained under assumptions that are much weaker and
 190 much more often fulfilled than the classically used assumption in (6). In particular, if we
 191 borrow from the LNA nothing but the assumption that transition probabilities $p(x_s | x_{s-1})$
 192 between subsequent measurement times are Gaussian and use some moment approximation
 193 method to calculate their means and covariance matrices, (8) can readily be evaluated
 194 without requiring the full $p(x | u, \theta)$ in (6) to be anything close to Gaussian.

195 Concretely, (8) can be evaluated as follows using an iterative scheme (which is more or
 196 less the same as classical Kalman filtering), as graphically presented in Figure 1b.

197

198 **Algorithm 1.**

199

- 200 1. Calculate approximate moments up to order two, $\eta_{x_1}^1, \eta_{x_1}^2$, of $p(x_1)$ by moment closure
 201 (see Supporting Information Section A).
- 202 2. Approximate the true $p(x_1)$ by a Gaussian distribution that has $\eta_{x_1}^1, \eta_{x_1}^2$ as moments.
- 203 3. Given the Gaussian model of the technical measurement errors, this implies that also
 204 the distribution $p(y_1) = \int p(y_1 | x_1)p(x_1)dx_1$ is Gaussian and it can be calculated
 205 from $\eta_{x_1}^1, \eta_{x_1}^2$ and σ . Evaluate $p(y_1)$ and store it for the likelihood calculation in (7).

- 206 4. Since the state prior $p(x_1)$ and likelihood $p(y_1 | x_1)$ are Gaussian, the state posterior
207 $p(x_1 | y_1)$ is also a Gaussian distribution that can be calculated from $p(x_1)$ and
208 $p(y_1 | x_1)$ thanks to Bayes' theorem, as is classically done in Kalman filtering.
- 209 5. Extract the moments up to order two, $\eta_{x_1|y_1}^1, \eta_{x_1|y_1}^2$, of $p(x_1 | y_1)$.
- 210 6. Solve moment equations (see Supporting Information Section A) over t_m time units
211 (i.e. over $[t_1, t_2]$) using $\eta_{x_1|y_1}^1, \eta_{x_1|y_1}^2$ as initial conditions in order to obtain moments
212 $\eta_{x_2|y_1}^1, \eta_{x_2|y_1}^2$ that approximate the moments of the distribution $p(x_2 | y_1)$.
- 213 7. Iterate: $p(x_2 | y_1)$ is approximated by a Gaussian equivalently to $p(x_1)$ in the step 2.
214 and so forth.

215 At the end of computation, we obtain an approximate scheme for calculating the likelihood
216 that only requires transitions distributions between subsequent measurement times to be
217 sufficiently close to Gaussian, and that “corrects” the classical *open loop* approximation in
218 moment equations by re-conditioning on the data at every measurement time.

219 **Remark 1.** *It is important to understand that conditioning on the data to evaluate*
220 *likelihoods leads to an approximation whose accuracy depends crucially on the data and not*
221 *only on how the stochastic model is approximated. This means that we need to carefully*
222 *investigate this dependence, but overall it is an important strength of the approach as it*
223 *allows us to deal with models for which all existing approximation methods fail if applied*
224 *in open loop.*

225 3 Results

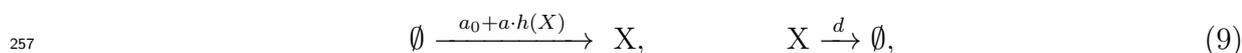
226 3.1 Accurate likelihood approximations using inaccurate approx- 227 imations of the chemical master equation

228 The core advantage of the approach in Algorithm 1 compared to standard likelihood
229 approximations based directly on equations (4) and (5) is that the approximation of

230 the CME only needs to be accurate over the time that passes in between subsequent
231 measurement steps, which, in many cases, is a much less restrictive requirement. For
232 instance, moment closure methods have been found to diverge or return negative variances
233 when solved over longer time horizons. As we mentioned earlier, Gaussian approximations,
234 such as used in the LNA, are necessarily bad for bimodal distributions displayed by some
235 systems. It is important to understand, however, that what is typically meant by “a
236 system displaying bimodal distributions” is that the solution of the CME becomes bimodal
237 when sufficient time passes but this does not necessarily imply that transition distributions
238 between subsequent measurement time points are also bimodal. For instance, when a
239 genetic toggle switch [17] is started in one of its stable equilibrium points, then after
240 a short amount of time the system will still most likely be close by and the solution
241 of the CME may still be unimodal around the equilibrium point. If measurements are
242 taken frequently, Algorithm 1 only requires approximations of such short term transition
243 distributions from given “initial states” in Step 6 and may therefore lead to accurate
244 approximations of the likelihood even if the deployed CME approximation deteriorates
245 over longer time horizons.

246 **3.1.1 Studying accuracy of likelihood approximations with a simple positive** 247 **feedback loop system**

248 To demonstrate accuracy of likelihood approximations, we start by investigating a simple
249 reaction network containing only a single chemical species, X. In this study case accurate
250 solution of the full chemical master equation using finite state projection (FSP) [23] is
251 straightforward and approximated likelihoods can readily be compared to the true likeli-
252 hoods. To nevertheless obtain an example that is sufficiently interesting and challenging,
253 we assume that the production of X is a non-linear function of the abundance of X such
254 that a positive feedback loop is formed and the system displays bimodal distributions (see
255 Figure 2a). Concretely, we assume that production and degradation of X occur according
256 to



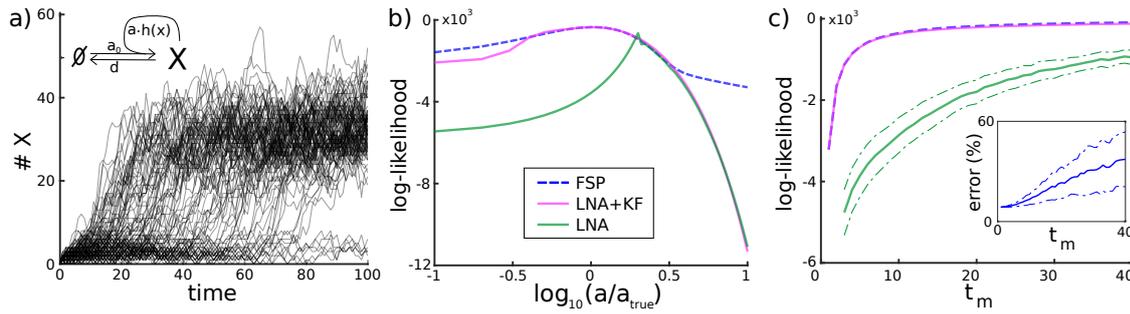


Figure 2: Accuracy of the likelihood approximation. A toy example of a reaction network (9) with a single chemical species and a positive feedback loop. **a)** $n = 100$ system trajectories obtained using Gillespie’s stochastic simulation algorithm. Molecule numbers of X are either very low, because only basal production of X is active at rate a_0 , or they switch to higher levels where the positive feedback is active. **b)** Data likelihood $p(y | \theta)$ for a single cell for the distance between measurement times $t_m = 10$ as a function of $\theta = a$ calculated using different approaches. Blue dashed: true data likelihood up to very small errors obtained by FSP with a large truncation set. Purple: approximation using the filtering approach. Green: approximation based on “open loop” use of the LNA. **c)** Mean \pm one standard deviation of data likelihood $p(y_{i,t_m} | \theta)$ averaged over $n = 100$ cells as a function of the distance between measurement times t_m . The colour coding is the same as in panel b. *Inset:* Relative error of the filtering based likelihood approximation as a function of t_m . Parameter values and initial condition for this system have been chosen as $a_0 = 2$, $a = 8$, $k = 20$, $n = 5$, $d = 0.25$, $\sigma = 4$, $X(0) = 0$, $t_S = 100$, all in arbitrary or no units.

258 where $h(x) = \frac{x^n}{x^n + k^n}$ and any possible dependence on input perturbations, $u(t)$, is, for
 259 now, omitted. We use the linear noise approximation not only because it is bound to be
 260 inaccurate for this system, but also because most moment closure methods are easy to use
 261 only for mass-action kinetics and require further adjustments when the propensities are
 262 given as Hill functions.

263 To investigate the accuracy of our likelihoods approximations, we assume that the system
 264 can be observed every t_m time units for various values of t_m and that data is collected
 265 according to the measurement model given in (3). For the sake of an easily understandable
 266 presentation, we decided to study the likelihood, $p(y | \theta)$, only for a single cell and
 267 only as a function of a single parameter, a , while the remaining parameters are fixed.

268 Figure 2b shows that, for $t_m = 10$, the true likelihood and its approximation according
269 to Algorithm 1 agree in the relevant parameter regime and have almost exactly the same
270 maximum, implying that the correct maximum-likelihood estimator would be obtained
271 with the approximated likelihood.

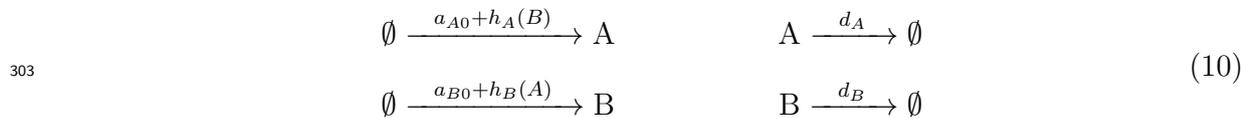
272 To compare this result to classical approaches, we also evaluated the likelihood using
273 the LNA and (6). Despite the fact that the same approximation of the CME is used
274 for the same system, the approximated likelihoods are very different and the classical
275 approach is very inaccurate almost everywhere in parameter space. This is not surprising
276 since the bimodal distributions displayed by the system are ill approximated by Gaussian
277 distributions.

278 On the other hand, since measurements are taken more frequently than the time that the
279 system needs to transition to larger molecule counts, the transition distributions between
280 subsequent measurements are sufficiently close to Gaussian to render the LNA useful when
281 it is deployed in Algorithm 1. The exception to this is when the value of the parameter a
282 is very large such that proteins are produced very quickly once the positive feedback loop
283 is triggered. In this case, transition distributions between measurement time points are
284 also bimodal and the likelihood approximation according to Algorithm 1 yields the same
285 inaccurate result as the classical approach. However, since the true value of the parameter
286 considered in this case study is not that large, this region of parameter space is not of
287 particular importance for the purpose of parameter inference.

288 To further benchmark Algorithm 1 we calculated the likelihood at the true parameters for
289 $n = 100$ cells as function of t_m . Figure 2c shows the average likelihood at the true value of
290 the parameter a as a function of t_m as well as standard deviations around it. It can be
291 seen that classical application of the LNA leads to extremely imprecise approximations for
292 almost all t_m , while Algorithm 1 is very accurate except when t_m is large and measurements
293 are too distant from each other for transition distributions to be close to Gaussian (see
294 the relative error in the inset).

295 3.1.2 Using the LNA with filtering to infer parameters of a toggle switch

296 Genetic toggle switches are systems consisting of two proteins that repress each others
297 production. Since such systems are typically constructed to implement bistability and
298 switching behaviour, they are typically considered as the role model of systems for which
299 one should not use the LNA. However, our results in the previous section suggest that it
300 might in fact be possible to accurately infer parameters of a toggle switch using the LNA
301 with filtering, as in Algorithm 1. To test this, we consider the following model of a genetic
302 toggle switch:



304 where $h_A(x) = \frac{a_A}{x^{n_A}+k_A}$ and $h_B(x) = \frac{a_B}{x^{n_B}+k_B}$ are Hill functions modelling the repression
305 of the production of protein A by protein B and of B by A, respectively. The resulting
306 system trajectories switch stochastically between two regimes where either A is present at
307 high copy numbers and the gene that produces protein B is repressed or B is present at
308 high copy numbers and the gene that produces protein A is repressed (Figure 3a). We
309 assume that only one of the proteins, B, can be measured every t_m time units in $n = 10$
310 cells with measurement noise as in (3).

311 To test if it is in principle possible to use the LNA for parameter inference for this
312 system without confounding the results by identifiability problems, we assumed that only
313 a single parameter, namely d_B , is unknown and used a Metropolis Hastings Markov chain
314 Monte Carlo (MCMC) method based on log-normal proposal distributions with fixed
315 variance for Bayesian inference, assuming a flat prior distribution for d_B . We find that the
316 MCMC algorithm generally converges very quickly to the vicinity of the true value of d_B
317 (Figure 3b). However, when the time between measurements is comparably large ($t_m = 20$)
318 the algorithm fluctuates around a value that is slightly smaller than the true value. Since
319 the data is sufficiently rich to, in principle, allow for a very precise estimation of only a
320 single parameter, we can attribute this error to approximation errors in the likelihood.

321 To investigate in more detail how this error depends on the distance between measurement

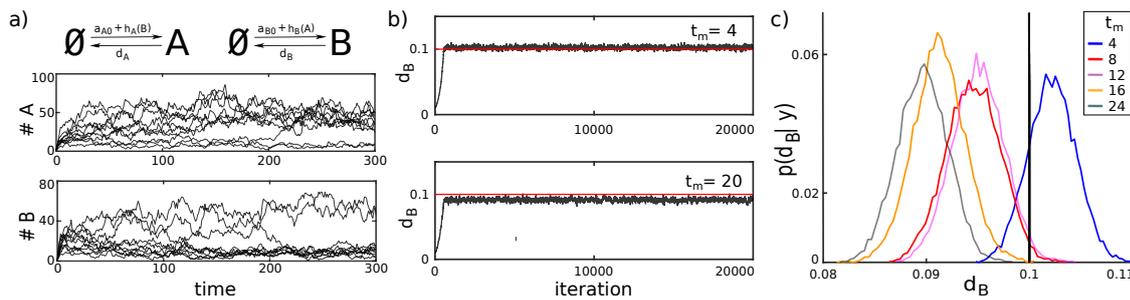


Figure 3: Parameter inference for a genetic toggle switch using the LNA with filtering. **a)** $n = 10$ system trajectories obtained using Gillespie’s stochastic simulation algorithm. **b)** The MCMC search (20000 steps) with data from $n = 10$ cells measured every $t_m = 4$ (top) and $t_m = 20$ (bottom) time units shows fast convergence from an initial guess of $d_B = 0.01$ to an approximated posterior distribution whose maximum is close to, but not exactly at, the true parameter value. The error is larger for the case where $t_m = 20$. **c)** Posterior distributions for different values of t_m : $t_m = 4$ (blue), $t_m = 8$ (red), $t_m = 12$ (purple), $t_m = 16$ (orange), $t_m = 20$ (grey). All distributions were obtained as histograms of the values visited by the corresponding Markov chain with a burn in period of 1000 steps. The magnitude of the error of the MAP estimator tends to become larger with increasing t_m as expected from the increasing error of the employed likelihood approximation. Parameters for this case study have been chosen as $a_{A0} = 0.5$, $a_{B0} = 0.5$, $a_A = a_B = 700$, $k_A = k_B = 100$, $n_A = n_B = 2$, and $d_A = d_B = 0.1$, $A(0) = B(0) = 1$, $\sigma = 4$, $t_S = 300$ all in arbitrary or no units. The model is given in (10).

322 times, t_m , we repeated the parameter search for $t_m \in \{4, 8, 12, 16, 20\}$ and found that the
 323 error of the maximum a posteriori (MAP) estimator increases with t_m but remains around
 324 10% even for the largest considered $t_m = 20$. This implies that despite the fact that the
 325 likelihood approximation becomes increasingly imprecise for large t_m , the maximum of the
 326 corresponding posterior distribution still remains at values that are reasonably close to
 327 the true parameter value. We therefore conclude that the errors in the estimation of d_B
 328 obtained here are overall very small, in particular for $t_m = 4$. This is quite remarkable
 329 considering that, to the best of our knowledge, no other successful use of any moment-based
 330 parameter estimation scheme has been reported up to date for systems like a genetic toggle
 331 switch. That said, it should be noted that the model considered here switches relatively
 332 slowly between the two regions. For faster switching systems more frequent measurements

333 would be required for a LNA-based likelihood approximation to remain accurate.

334 **3.2 Bayesian inference for experiments parallelised at the single** 335 **cell level**

336 Having demonstrated that the likelihood approximation can be accurate and used to infer
337 parameters even when the system dynamics are complex, we return to the main objective
338 of this paper: to investigate scalability of the approach and to test if it is practically usable
339 for experiments parallelised at the cell scale, where each cell is perturbed with a different
340 input.

341 To this end, we focus on the stochastic model that was used in [2] to describe single cell
342 responses to light of the CcaS/CcaR optogenetic system. Upon exposure to green light,
343 CcaS flips to an active state and phosphorylates the response regulator CcaR, which then
344 activates the production of a fluorescent protein [25]. Red light reverts this process and
345 stops gene expression very quickly.

346 In [2] the following single cell gene expression model for the CcaS/CcaR optogenetic system
347 was used:



349 where E is a generic variable that was named "cell responsiveness" and which pools together
350 sources of variability that are extrinsic to the studied gene (e.g. variations in the number
351 of ribosomes or plasmid copy number fluctuations). The second random variable, F ,
352 quantifies the amount of fluorescent protein present in a cell and is directly measured in
353 experiments (up to a scaling factor s , see Supporting Information Section B for details).
354 The single cell protein production rate is determined by the cell responsiveness, $E(t)$,
355 together with the non-linear time-dependent variable, $L(t)$, which is indirectly controlled
356 by the external light signal, $u(t)$. More precisely,

$$357 \quad L(t) = \frac{(c \cdot l(t))^n}{(c \cdot l(t))^n + k^n}, \quad (12)$$

358 where $l(t)$ evolves according to

$$359 \quad \frac{dl(t)}{dt} = u(t) - c \cdot l(t), \quad (13)$$

360 where the light signal $u(t) = 1$ if the cell is exposed to green light at time t , and $u(t) = 0$
361 if the cell is exposed to red light. Overall, the model contains 8 parameters that need
362 to be inferred from the data, $\theta = \{a, b, s, m, h, c, n, k\}$, where $m := \frac{h_0}{h}$ is the mean of the
363 stationary Poisson distribution determined by the dynamics of E .

364 **3.2.1 Inference of parameters from simulated data**

365 To test if parameters can be better learned from parallelised experiments at the single
366 cell level, we fixed the parameters of the model as in Table B.1 and used Gillespie's
367 stochastic simulation algorithm to simulate data sets in which different numbers of cells are
368 exposed to various light patterns. We then used a Metropolis-Hastings MCMC algorithm
369 to perform Bayesian inference of the model parameters for each of the data sets, in order
370 to test how well the true parameter values can be learned in these cases.

371 Concretely, we consider experiments in which each cell is exposed to one of 6 different
372 light patterns (see Figure 4). We categorise these 6 light patterns into those with short,
373 intermediate and long green signals. Light pattern 1 has a multiple short green signals
374 that last only 12 min, light patterns 2, 3, 4 and 5 fall into the intermediate category with
375 green signals lasting between 30 min and 300 min, and the light pattern 6 has a long green
376 signal that lasts until the end of experiment. We simulated 10 cells for each of the six light
377 profiles so as to obtain a data set with 60 cells in total. This set of cells forms a mixed
378 group called Group 0 (G0). Then, we simulated an additional 60 cells, all exposed to the
379 short light pattern 1, and named it Group 1 (G1). We repeated the same process for two
380 light patterns from the intermediate category, light patterns 2 and 5, and created Group 2
381 (G2) and Group 5 (G5), and again for the long green signal and created dataset Group 6
382 (G6). We decided to omit results for light patterns 3 and 4 as these results do not change
383 our conclusions and for the sake of readability of the figures that follow.

384 As in the actual experimental setup, measurements are taken every $t_m = 6$ min up to a final

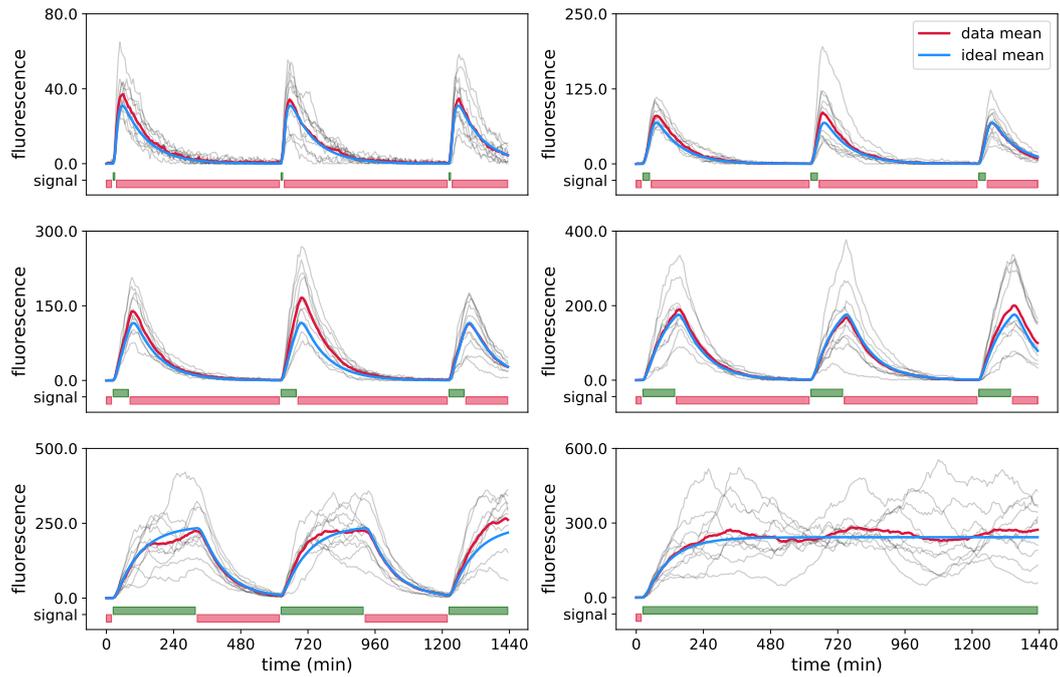


Figure 4: Simulated data for the CcaS/CcaR optogenetic system. One data set used for parameter inference for G0 is shown in the six panels. The different panels display cells exposed to different light inputs, $u(t)$. The pulses Red and green light inputs are shown in red and green bars, respectively, at the bottom of each panel. Note that each light pattern starts with red light. Simulated data for 60 cells in total, 10 cells per panel, is shown in grey, the mean of the displayed cells in red, and the true mean (i.e. for infinitely many cells) is shown in blue. Parameter values are provided in the Supporting Information Table B.1. Simulated data sets for other groups and different numbers of cells are shown in the Supporting Information Figures B.1 - B.3.

385 time of $t_S = 1440$ min, implying that there are 240 measurements per cell and that classical
 386 calculation of likelihoods based on approaches such as in (6) would be computationally
 387 difficult or unfeasible. On the other hand, likelihood evaluation using LNA with Kalman
 388 filter (Algorithm 1) takes only a fraction of a second per cell. Nevertheless, we need to
 389 explore an 8-dimensional parameter space, which remains challenging since some of the
 390 parameters affect the measured output similarly and cause identifiability problems.

391

392 Running the MCMC algorithm for 10 000 iterations, we find that the quality of parameter
 393 estimates depends strongly on the group of cells that is used for inference, meaning

394 that the identifiability of parameters depends on which light input is applied to cells.
395 Diversified light signals in cells (as in G0) lead to relatively tight posterior distributions
396 for all parameters and MAP estimates that are close to the true values of the parameters
397 (see Figure 5).

398 When all cells receive the same light input (G1, G2, G5, G6) some of the parameters are
399 not identifiable in most cases. In particular, posterior distributions for G1 and G6 are very
400 broad for some parameters. Cells in G6 are exposed permanently to green light and the
401 full dynamics of the system are never exposed and cannot be characterised (parameters
402 a, b, s and c , see Figure 5a). Cells in G1 are exposed to only very short pulses of green
403 light which makes it difficult to reliably learn parameters that characterise fluctuations in
404 the cells' responsiveness E (parameters a, s, m and h , see Figure 5a). This is because E
405 influences the measured output only in the presence of green light, that is when protein
406 production is active. Data from G2 have better posterior distributions than G6 or G1, but
407 still lead to similar identifiability problems as in G1 for certain parameters (parameters a
408 and s , see Figure 5a).

409 However, data from G5 leads to posterior distributions that are of similar quality to the
410 mixed group G0 (see Figure 5a and Supporting Information Figure B.6a.) implying that the
411 light inputs of this group excite the system dynamics in a sufficiently rich way, at least if 60
412 cells are used for parameter inference. Using fewer cells would, however, be advantageous.
413 On the one hand because evaluating likelihoods of many cells is computationally costly, and
414 on the other hand because cells in the experiment that are not needed for the calibration
415 of models can be used for other purposes, for instance for validating model predictions.

416 We therefore investigated how the quality of inferred parameter estimates for the different
417 groups changes when the total number of cells is reduced. We additionally simulated
418 smaller sets of data for the same groups G1, G2, G5 and G6, with the total number of 12
419 and 30 cells. Similarly, we simulated the corresponding smaller data sets for G0 with 2
420 and 5 cells per light pattern, respectively (see Figures B.1 - B.3). We find that when cells
421 receive diversified light signals (G0), using only 12 cells (2 each for any of the six light
422 inputs) leads to quite broad posterior distributions but using 30 cells (5 cells per light

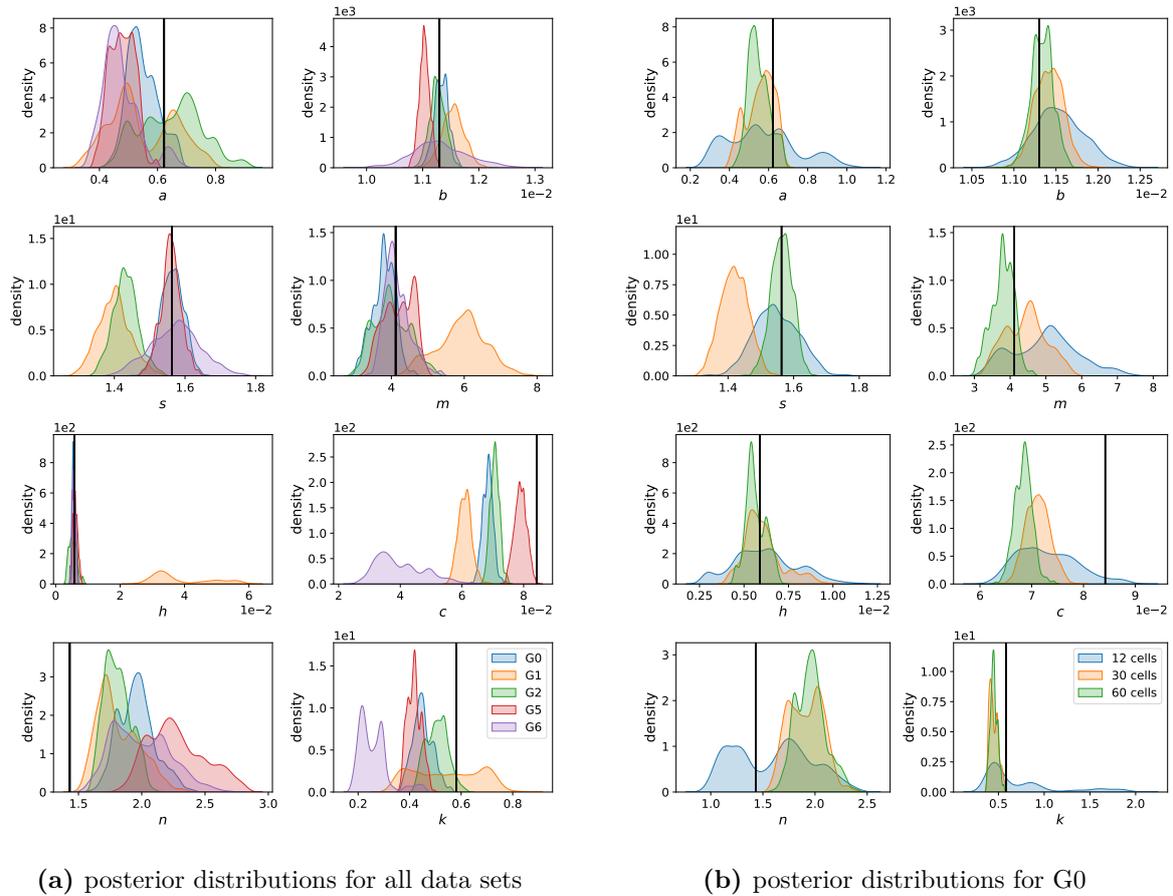


Figure 5: Inference results for the CcaS/CcaR optogenetic system. (a) One-dimensional marginals of posterior distributions obtained using 60 cells, either all exposed to the same light input (G1 - orange, G2 - green, G5 - red, G6 - purple) or to diversified light inputs containing 10 cells for each of the 6 light patterns (G0 - blue). The vertical black lines in each panel shows the true value of the corresponding parameter that was used to simulate the data sets. Results for data sets with 12 and 30 cells are provided in the Supporting Information Figure B.4. Two-dimensional marginals showing correlations between model parameters are provided in the Supporting Information Figures B.7 - B.11. (b) Posterior distributions for the mixed group, G0, for varying numbers of total cells that were used to infer parameters. Results for other groups are provided in the Supporting Information Figures B.5 - B.6. In all cases, posterior distributions have been obtained as histograms of values visited by the MCMC chain of 10 000 iterations, without a burn-in period of 4 000 iterations.

423 input) leads to posterior distributions and parameter estimates that are almost as good
424 as when 60 cells are used (Figure 5b). Using 30 cells but applying the same light input
425 to all of them, however, leads to broader posterior distributions (Supporting Information
426 Figures B.4 - B.11).

427 We can therefore conclude that exposing different cells to different perturbations within
428 the same experiment provides means to increase the information content in the data and
429 thereby allows us to learn model parameters either more precisely or with fewer measured
430 cells. In the very least, a variation of light inputs applied to cells within a single experiment
431 (like in G0) ensures that some of these inputs excite the system dynamics in informative
432 ways and that not all cells are stimulated with light inputs that alone do not provide
433 information on all model parameters (such as in G1, G2 and G6).

434 **3.2.2 Iterative likelihood evaluation reduces the computational cost by orders** 435 **of magnitude**

436 To establish that experiments parallelised at the single cell level are useful for calibrating
437 models, we deployed our likelihood approximation and showed that it is practically
438 applicable for inferring model parameters from data of such experiments. With that
439 established, we focus in more detail on the computational cost of the likelihood evaluation
440 and investigate how it scales with the number of cells, the number of different light inputs
441 that are used, and the number of measurement times for each cell. We stress that it
442 is crucial that the likelihood evaluation is very fast since this calculation needs to be
443 performed at every iteration of the MCMC algorithm, that is at least 10 000 times, to
444 infer parameters in the case studies shown in the previous section.

445 To be able to compare the computational cost of our approach to classical likelihood
446 calculations, we additionally implemented a method that uses the LNA to directly ap-
447 proximate the full likelihood according to (6). In the classical approach, once the full
448 probability distribution $p(x | u, \theta)$ is calculated for given parameters θ , data from all
449 cells that have received the input perturbation $u(t)$ simply need to be plugged into this
450 distribution to calculate the likelihood of the full data set. On the other hand, Algorithm 1

451 is always operating on the specific data of single cells and needs to be re-run for every
 452 cell irrespective of whether or not they all received the same light input. This advantage
 453 of the classical approach naturally becomes less relevant the more diversified the light
 454 inputs are that are sent to the cells in the experiment. Furthermore, when the number of
 455 measurement times per cell increases, the dimensionality of the data increases and the
 456 calculation of $p(x | u, \theta)$ becomes prohibitively costly, even if the computationally cheap
 457 Gaussian approximations of the LNA are used. In contrast to that, if Algorithm 1 is used
 458 for the likelihood calculation, additional measurements only require additional iterations
 459 of the algorithm's steps, which increases its computational cost only linearly.

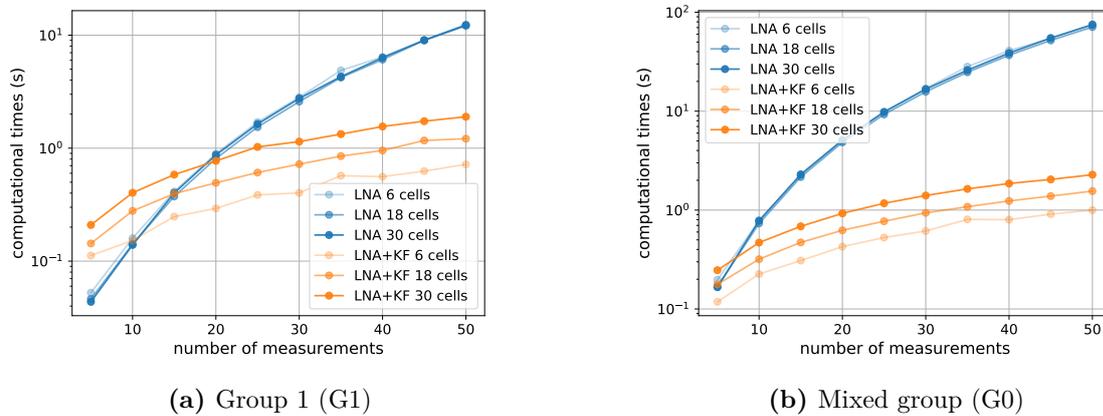


Figure 6: Evaluation of the computational efficiency of Algorithm 1 (a) The computational cost of likelihood evaluations using Algorithm 1 (orange) for different numbers of cells (differently shaded lines) is compared to classical likelihood evaluations (blue) when all cells have received the same light input (G1). (b) Same as panel (a) except that cells have been exposed to all 6 different light inputs (G0).

460 In summary, setting aside differences in the accuracy of the likelihood approximations
 461 obtained by the different approaches, we expect Algorithm 1 to computationally outperform
 462 classical approaches except if the number of measurements per cell in the experiment is
 463 very small. Concretely, we find that if all cells receive the same light input, as in the
 464 data considered for G1, Algorithm 1 is computationally superior when the number of
 465 measurements is larger than 10 if 6 cells are observed, when the number of measurements
 466 is larger than 15 if 18 cells are observed, and when the number of measurements is larger

467 than 20 if 30 cells are observed (Figure 6a). For the data set with 240 measurements per
468 cell that we considered in this paper, the classical approach would be way too costly to
469 be used at all. When cells are exposed to different light inputs, as in the data considered
470 for G0, Algorithm 1 becomes even more advantageous and computationally outperforms
471 the classical approach even for 30 cells when only a handful of measurements are taken
472 (Figure 6b).

473 3.2.3 Inference of parameters from experimental data

474 We established using simulated data that the presented likelihood approximation is com-
475 putationally efficient and can be used for inference of model parameters when experiments
476 are parallelised at the single cell level. Can we use it in practice?

477 We have recently started to use the CcaS/CcaR optogenetic system to drive parts of a
478 repressilator circuit with light. A first practical task is to quantitatively characterise the
479 optogenetic system when plasmids carrying the repressilator circuit are present in cells. To
480 test this, we constructed a strain in which the CcaS/CcaR system is driving a fluorescent
481 reporter protein, the repressilator circuit is present but the CcaS/CcaR system is not
482 coupled to it (see Supporting Information Section C). We performed an experiment with
483 the new strain in which cells are grouped and exposed to light stimulations as in the *in*
484 *silico* case study in Section 3.2.1. We then manually curated the data and extracted a
485 data set of 30 viable cells (five cells for each of the six groups, see Figure 7) in order to
486 parameterise the model.

487 Using the Metropolis-Hastings MCMC scheme for parameter inference, we observe that
488 some of the parameters are difficult to identify (Figure C.15). The reason for this might
489 be the mismatch between the model and real system coupled to unknown, and probably
490 non-Gaussian, experimental measurement noise that makes parameter inference from real
491 data difficult compared to an idealised *in silico* case study. In particular, since the light
492 activation function $L(t)$ is only indirectly connected to the observable fluorescence levels
493 in cells, the multiple parameters that were used in (12) to define its shape are difficult to
494 determine jointly.

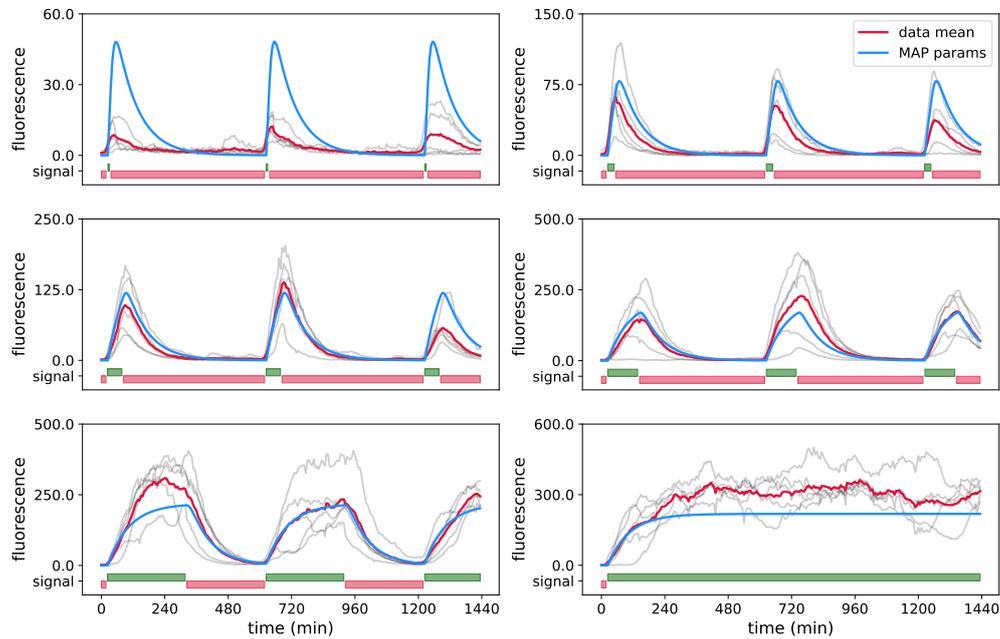


Figure 7: Experimental data for the CcaS/CcaR optogenetic system. The curated data set of 30 cells in total, 5 cells per panel, is shown in grey, where each trajectory shows the fluorescence of one cell over a period of 24 hours with the time between subsequent measurement points being 6 min. Each panel displays cells exposed to the same light inputs, $u(t)$. Red and green light inputs are shown in red and green bars, respectively, at the bottom of each panel. Note that each light pattern starts with red light. The mean of the displayed cells is shown in red, and the blue lines are ideal means predicted by the model using MAP estimator obtained from this same data set (see Supporting Information B.5 and Table C.3, MAP (4)).

495 We therefore decided to fix 4 of the 8 unknown model parameters to the maximum
496 likelihood estimates extracted from the first run of the Metropolis-Hastings MCMC
497 algorithm (see Table C.3) and to re-run the search with only b, m, c and h as unknown
498 model parameters. The results (Figure 8) show that b, m, c and h are now well identifiable.
499 Of particular interest is that the posterior distribution obtained for h is very narrow. In
500 the model, h quantifies the time scale of fluctuations of the individual cell responsiveness
501 to light. This implies that learning h requires that single cells are tracked in time and

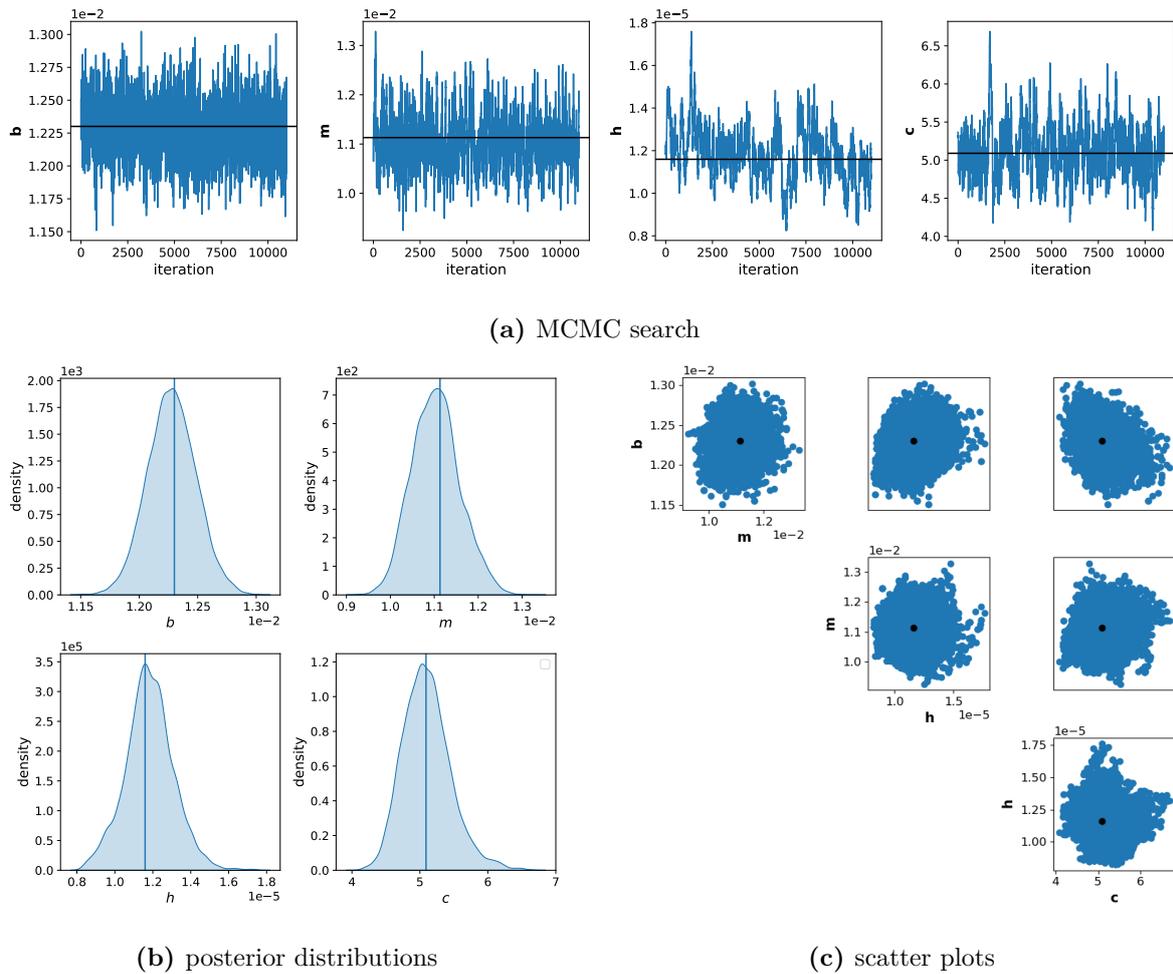


Figure 8: Parameter inference using Algorithm 1 on experimental data. (a) MCMC chains, after a burn in period, for parameters b, m, h and c . Black lines show values and the position in the chain of the MAP estimator. (b) One-dimensional marginals of posterior distributions. (c) Two-dimensional posterior marginals of the MCMC parameter search of 14 000 iterations in total and the burn-in period of 3 000 iterations. Black dots show the values of the parameters of the MAP estimator. Data used for inference: 30 cells in total, 5 cells per light group, as in G0, shown in Figure 7.

502 that resulting time-correlation information is exploited for parameter inference. The
 503 likelihood approximation presented in this work exploits this information by construction.
 504 Furthermore, h is a parameter of particular biological interest since it quantifies the
 505 fluctuations of a particular trait of cells (responsiveness to light), and hence how long
 506 the trait is inherited over cell generations. A priori, one would expect that a cell's

507 responsiveness to light is determined by quantities such as plasmid copy numbers or the
508 number of ribosomes and that these fluctuate on a time scale that is determined by the
509 cells' growth rate. However, we find that the value for h that best explains our data
510 is much smaller than what would be expected from the growth rate and that cells that
511 responded more (or less) than average to light early in the experiment still tend to be
512 more (or less) responsive several tens of cell generations later. The small value of h may
513 suggest that variability in cell responsiveness to light is driven by a biological mechanism
514 that generates atypically long memories of the state of the cell. For instance, plasmid
515 copy numbers might not fluctuate at the time scale of cellular growth due to feedback
516 regulation that maintains high or low numbers.

517 4 Discussion

518 Despite immense advances in measurement technologies in the past decade, developing
519 predictive models of natural or synthetically constructed gene networks remains a major
520 challenge. Cell populations are heterogeneous while gene expression dynamics are stochastic
521 and often regulated by intricate and poorly understood mechanisms. Understanding such
522 processes in detail requires observations of sufficiently many single cells. Stochastic kinetic
523 models are helpful to extract relevant information from the data and to test if a hypothesised
524 biological mechanism can explain the observations. However, if years of modelling in
525 systems biology have shown one thing, then it is that (some) parameters of mechanistic
526 models of biochemical processes are almost always practically unidentifiable from the
527 insufficiently informative data that is available [11]. This has led to increased efforts to
528 quantify the information content of data and to develop methods to plan perturbations
529 such that the studied system is excited in ways that reveal its dynamics [27, 28]. While
530 well chosen perturbations may resolve some identifiability problems, it is to be expected
531 that a single perturbation experiment will rarely be sufficient to ensure that the studied
532 system displays all the dynamics that need to be observed to properly understand it. Novel
533 experimental platforms that allow the user to parallelise different perturbations within a
534 single experiment may resolve this problem but raise new questions and challenges.

535 We have focused here on microscopy platforms in which perturbations can be parallelised
536 by targeting light signals at single cells. In principle, this allows one to test as many
537 perturbations in parallel as there are cells in the experiment. Yet, in the face of stochasticity
538 in the system and noise in the measurements, observing a few cells per perturbation implies
539 that we will not be able to reliably observe the system's response to any of the employed
540 perturbations. One might then wonder whether it is more informative to diversify input
541 perturbations but to observe only few cells for each input or should we rather ensure that
542 many cells are observed for only few inputs? By studying parameter inference on simulated
543 data for the model of the CcaS/CcaR optogenetic system, we show that exposing different
544 cells to different perturbations within the same experiment (G0) provides means to increase
545 the information content in the data and allows us to learn model parameters either more
546 precisely or with fewer observed cells (Figures 5, B.4, B.5, B.6).

547 To reach these results, we had to deploy a method for the approximation of likelihoods
548 that couples moment approximations, Gaussian assumptions, and a Kalman filter, since
549 classical likelihood evaluation would have been computationally infeasible for the considered
550 data. Similar methods have been proposed in the past (in particular in reference [3]
551 for an epidemiological model) but it has not been recognised that iterative likelihood
552 evaluations may lead to vastly improved precision if the chemical master equation needs
553 to be approximated. Relevant differences in approximation qualities appear when the
554 measurements are taken frequently compared to the time scale of the system. Our results
555 indicate that iterative and classical likelihood calculation lead to the same approximation
556 when the protein production rate is large in the case study of the positive feedback loop
557 (Figure 2b), and that the quality of obtained parameter estimates becomes worse when
558 measurement times are more distant for the studied model of a genetic toggle switch
559 (Figure 3b,c).

560 Additionally, this likelihood approximation method exploits the time-correlation infor-
561 mation of every tracked cell, and enables us to learn some parameters of the model
562 of the CcaS/CcaR optogenetic system from the experimental data, that are otherwise
563 unidentifiable (Figure 8).

564 In this work we demonstrate that frequent measurements and the use of Algorithm 1 lead
565 to very accurate likelihood approximations in all our case studies, even in cases where the
566 used approximation of the CME is obviously bad. Importantly, increasing numbers of
567 measurements increases the computational time of the approach only linearly (Figure 6)
568 whereas classical approaches that aim to calculate the joint distribution of all data points
569 become essentially impossible to use when the number of measurement times ranges in the
570 order of hundreds. With such data becoming more and more available, it can therefore be
571 foreseen that parameter inference for stochastic kinetic models of biochemical reaction
572 networks will have to be performed using approaches such as the one presented in this
573 paper in the future.

574 Experimental methods

575 **Bacterial strains and plasmids.** All experiments were performed with *Escherichia*
576 *coli* strain HR14. The strain is derived by transformation of *Escherichia coli* MC4100 with
577 three plasmids. Plasmids pSR43.6 and pSR58.6_3spng carry a CcaS/CcaR-based optoge-
578 netic module that drives or reduces expression from P_{cpcG2-172} promoters in green/red
579 light respectively, and synthetic pathway for chromophore phycocyanobilin [30]. Plasmid
580 pHyRep-Prg-0 carries a P_{cpcG2-172} regulated venus-YFP reporter, constitutive CFP
581 reporter, and a stabilized repressilator without *ssrA* degradation tags[30]. pSR58.6_3spng
582 also includes a repressor “sponge” region (*tetR*, *lacI*, and *cI* binding sites) from plasmid
583 pLPT145 that reduces background repressor levels and reduces variability in repressila-
584 tor period [26]. See table StrainsPlasmids.xlsx and plasmid maps PlasmidMaps.zip for
585 composition. Cells were constructed and maintained in LB broth supplemented with
586 Spectinomycin (100 µg/ml), Ampicillin (100 µg/ml), Chloramphenicol (20 µg/ml) as
587 appropriate.

588 **Microscopic culture and optogenetic stimulation setup.** Experiments were per-
589 formed as previously described, using an automated microscope platform for closed loop
590 imaging and data processing, optogenetic stimulation, and environmental regulation [2].

591 Briefly, custom software operates an Olympus IX83 fluorescence microscope fitted with
592 a 100x objective contained in an opaque, temperature-controlled incubator. The setup
593 obtains CFP (x438/29,m483/22) and YFP (x513/22,m543/22) fluorescence images, derives
594 cell size and expression reporter fluorescence levels, and delivers patterned optogenetic
595 light stimuli to cells via a modified LCD projector with custom 530nm and 660nm light
596 sources.

597 Bacterial cells are grown in microfluidic mother machines with $23\mu\text{m} \times 1.3\mu\text{m} \times 1.3\mu\text{m}$ (l,w,h)
598 growth channels at $5\mu\text{m}$ spacing along a split media trench. The microfluidic devices are
599 fabricated from degassed polydimethylsiloxane (Dow Sylgard 184, 1:10 catalyst:resin),
600 cured against epoxy replicate master molds, ports punched, and plasma-bonded (Harrick
601 PDC-002, medium power, 1 minute) to clean glass cover slips. See [2] for detailed protocol.
602 Polyethylene tubing (Instech, BTPE-50) press-fitted to 22ga luer stubs and cannulae
603 (Instech) is used to connect media and waste flows to the device. Media flow is regulated
604 by two syringe pumps (WPI, Alladin-1000).

605 **Image acquisition, processing and cell stimulation loop.** The setup cycles through
606 ten stage locations every six minutes, focusing and determining xy offsets, obtaining
607 fluorescence images and delivering optogenetic stimuli to the cells [2]. To reduce optoge-
608 netic response to fluorescence imaging, exposures are minimized and maintained across
609 experiments, and fluorescence image acquisition is immediately followed by application of
610 light stimuli to set the cells' optogenetic state.

611 Briefly, at each timepoint and stage location, image-based autofocus and xy-stage jitter and
612 drift corrections are first performed. Fluorescence images are then acquired and corrected
613 for small, slowly varying, additive camera signal offsets (by subtracting median dark images
614 acquired alongside each), and shading corrected using previously obtained, normalized
615 calibration images of a uniform fluorescent field (10% Fluorescein 0.1%NaHCO₃)[22].
616 The images are spatially registered and fluorescence-based expression estimates for the
617 constitutive and light-controlled reporters are extracted for individual mother cells as the
618 97th percentile pixel intensity within a pre-specified bounding box at each cell's image

619 location. Constitutive reporter fluorescence images are segmented to derive cell sizes, and
620 growth rates estimated via a moving average of differences in $\log_2(\text{cell length})$, excluding
621 outliers due to cell division and segmentation errors.

622 The programmed CcaSR activation ($\sim 535\text{nm}$) or deactivation ($\sim 670\text{nm}$) light stimuli for
623 each cell are then mapped to green or red boxes overlying the positions of the cells in an
624 RGB image. The image is transformed to register projector to camera image planes, and
625 projector shading corrections (using low-pass filtered reflected uniform field projections,
626 obtained at experiment initialization) applied to each color channel. To stimulate the
627 cells, this image is projected onto the field of view for ten seconds (670nm : $\sim 10.5\text{ mW/cm}^2$,
628 535nm : 7.6 mW/cm^2 ; Contrast relative to dark LCD panels, of 252 and 361, respectively;
629 crosstalk between channels $<1\%$).

630 **Experiment setup, media, and conditions** Experiments were initialized and run
631 as previously described [2]. Bacterial cells are diluted 1:100 from -80C glycerol stocks
632 into 5ml fresh LB media containing 0.01% Tween20, $20\mu\text{g/ml}$ Chloramphenicol, $100\mu\text{g/ml}$
633 Spectinomycin, and $100\mu\text{g/ml}$ Ampicillin to maintain plasmids, and incubated for 6-7
634 hours at 37C . The experimental apparatus is equilibrated, and the microfluidic device
635 prepared by filling with 0.01% Tween20 for 1 minute, then purging with air. The device is
636 prewarmed, and cell culture concentrated (centrifuge $4000 \times g$ for 4 minutes, and resuspend
637 pellet in $4\ \mu\text{l}$ supernatant) and injected into the device. Media supply and waste tubes are
638 fitted to the device and LB containing 0.4% glucose and 0.01% Tween20 delivered at 4
639 ml/hour for 1 hour, and 1.5 ml/hour - 2.0 ml/hour thereafter.

640 With cells in the device, the experiment control software is initialized and calibrations
641 (camera and projector offsets from the PDMS-glass interface, projector-camera transforms,
642 and projector shading correction) performed. Fields of view are set, measurement regions
643 for individual mother cells are specified and light sequences to be delivered are linked to
644 each cell. ~ 50 cells per field of view are distributed between seven repeating 100-pulse light
645 sequences (consisting of initial green pulses of duration 0, 2, 5, 10, 20, 50, or 100 intervals
646 followed by red pulses for the remaining intervals. All cells receive a common 100-pulse

647 red “zeroing” sequence during the first 10 hours of the experiment. The system then
648 continues to acquire data and stimulates cells according to their assigned light sequences
649 for the remainder of the experiment. For the model calibration the full red sequence is
650 not relevant and is left out and the common initial 100-pulse red “zeroing” sequence is
651 shortened to the duration of 4 intervals. Data with no pulse (0 pulse duration) is not used
652 for the model calibration as it is a trivial case and no parameter could be learnt with it.

653 **Cell classification and invalidation.** Cells that stop growing or filament in the
654 microfluidic device, or that lose plasmids or optogenetic system function, or transiently
655 shift from the detection/excitation regions can give unreliable data. Our automated setup
656 therefore continuously tests and permanently flags cells that fail presence, growth, and
657 optogenetic system function tests [2]. Briefly, the constitutively-expressed reporter (here,
658 CFP) is used to test control location presence and measurement quality (invalidating cells
659 for signal loss or noise above threshold). The constitutive reporter images are also used for
660 cell shape and growth rate extraction (invalidation for growth below threshold rate), and
661 pSR43.6 loss detection (invalidation for resulting growth rate increase and reduced CFP
662 concentration). Optogenetic response is also measured via the linked reporter (YFP) and
663 cells can be flagged for response below minimum threshold (indicating pSR58.6_3spng or
664 optogenetic system loss). Cells failing individual tests are typically classified as invalid
665 for the remainder of the experiment. As the automated classifier may fail to invalidate
666 pathological cells, we followed the standard invalidation protocol above with a manual
667 cull of remaining cells that presented improbable and outlier trajectories to our eye, for
668 example, cells that never react to the light signal or those that after some time show a
669 non-zero flat behaviour, that we assume to be dead.

670 **Data and code availability.** Data and code are publicly available as a git repository
671 on gitlab: <https://gitlab.pasteur.fr/adavidov/inferencelnakf> and on zenodo: [10.5281/zenodo.5229416](https://zenodo.org/record/5229416).
672

673 Contributions

674 Conceptualization, methodology: JR

675 Formal analysis, investigation: AD, JR

676 Resources: RC

677 Software, visualisation: AD

678 Supervision: GB, JR

679 Writing: AD, RC, GB, JR

680 Acknowledgements

681 We thank Virgile Andreani for useful discussions about the model and parameter inference.

682 We thank Johan Paulsson and Jeffrey J Tabor for kind gifts of plasmids.

683 References

684 [1] F. Bertaux, S. Sosa-Carrillo, A. Fraisse, C. Aditya, M. Furstenheim, and G. Batt.
685 Enhancing bioreactor arrays for automated measurements and reactive control with
686 reacsight. *bioRxiv*, pages 2020–12, 2021.

687 [2] R. Chait, J. Ruess, T. Bergmiller, G. Tkačik, and C. Guet. Shaping bacterial
688 population behavior through computer-interfaced control of individual cells. *Nature*
689 *Communications*, 8:1535, 2017.

690 [3] P. Fearnhead, V. Giagos, and C. Sherlock. Inference for reaction networks using the
691 linear noise approximation. *Biometrics*, 70:457–466, 2014.

692 [4] B. Finkenstädt, DJ Woodcock, M. Komorowski, CV Harper, JRE Davis, MRH White,
693 DA Rand, et al. Quantifying intrinsic and extrinsic noise in gene transcription using
694 the linear noise approximation: An application to single cell data. *The Annals of*
695 *Applied Statistics*, 7(4):1960–1982, 2013.

- 696 [5] M. Folia and M. Rattay. Trajectory inference and parameter estimation in stochastic
697 models with temporally aggregated data. *Statistics and Computing*, 28(5):1053–1072,
698 2018.
- 699 [6] Z. Fox, S. Fletcher, A. Fraisse, C. Aditya, S. Sosa-Carrillo, S. Gilles, F. Bertaux,
700 J. Ruess, and G. Batt. Micromator: Open and flexible software for reactive microscopy.
701 *bioRxiv*, 2021.
- 702 [7] F. Fröhlich, P. Thomas, A. Kazeroonian, F. Theis, F. Grima, and J. Hasenauer.
703 Inference for stochastic chemical kinetics using moment equations and system size
704 expansion. *PLOS Computational Biology*, 12(7):e1005030, 2016.
- 705 [8] D. Gillespie. A general method for numerically simulating the stochastic time evolution
706 of coupled chemical reactions. *Journal of Computational Physics*, 22(4):403–434, 1976.
707 ISSN 0021-9991.
- 708 [9] D. Gillespie. A rigorous derivation of the chemical master equation. *Physica A*, 188
709 (1-3):404–425, 1992.
- 710 [10] R. Grima. Linear-noise approximation and the chemical master equation agree up to
711 second-order moments for a class of chemical systems. *Physical Review E*, 92:042124,
712 2015.
- 713 [11] R. Gutenkunst, J. Waterfall, F. Casey, K. Brown, C. Myers, and J. Sethna. Universally
714 sloppy parameter sensitivities in systems biology models. *PLoS Computational Biology*,
715 3(10):e189, 2007.
- 716 [12] J. Hasenauer, S. Waldherr, M. Doszczak, N. Radde, P. Scheurich, and F. Allgöwer.
717 Identification of models of heterogeneous cell populations from population snapshot
718 data. *BMC Bioinformatics*, 12(1):125, 2011.
- 719 [13] M. Komorowski, B. Finkenstädt, C. Harper, and D. Rand. Bayesian inference of
720 biochemical kinetic parameters using the linear noise approximation. *BMC Bioinfor-*
721 *matics*, 10(1):343, 2009.

- 722 [14] M. Komorowski, M. Costa, D. Rand, and M. Stumpf. Sensitivity, robustness, and
723 identifiability in stochastic chemical kinetics models. *Proceedings of the National*
724 *Academy of Sciences of the USA*, 108(21):8645–8650, 2011. doi: 10.1073/pnas.
725 1015814108.
- 726 [15] J. Liepe, P. Kirk, S. Filippi, T. Toni, CP Barnes, and MPH Stumpf. A framework for
727 parameter estimation and model selection from experimental data in systems biology
728 using approximate bayesian computation. *Nature protocols*, 9(2):439, 2014.
- 729 [16] G. Lillacci and M. Khammash. Parameter estimation and model selection in
730 computational biology. *PLoS Computational Biology*, 6(3):e1000696, 2002. doi:
731 10.1371/journal.pcbi.1000696.
- 732 [17] J-B. Lugagne, S. Sosa Carrillo, M. Kirch, A. Köhler, G. Batt, and P. Hersen. Balancing
733 a genetic toggle switch by real-time feedback control and periodic forcing. *Nature*
734 *Communications*, 8:1671, 2017.
- 735 [18] H. McAdams and A. Arkin. Stochastic mechanisms in gene expression. *Proceedings*
736 *of the National Academy of Sciences of the USA*, 94(3):814–819, 1997.
- 737 [19] H. McAdams and A. Arkin. It’s a noisy business! genetic regulation at the nanomolec-
738 ular scale. *Trends in Genetics*, 15(2):65–69, 1999.
- 739 [20] A. Miliadis-Argeitis, M. Rullan, SK Aoki, P. Buchmann, and M. Khammash. Automated
740 optogenetic feedback control for precise and robust regulation of gene expression and
741 cell growth. *Nature communications*, 7:12546, 2016.
- 742 [21] P. Milner, C. Gillespie, and D. Wilkinson. Moment closure based parameter inference
743 of stochastic kinetic models. *Statistics and Computing*, 23(2):287–295, 2013.
- 744 [22] M. Model and J. Burkhardt. A standard for calibration and shading correction of
745 a fluorescence microscope. *Cytometry: The Journal of the International Society for*
746 *Analytical Cytology*, 44(4):309–316, 2001.

- 747 [23] B. Munsky and M. Khammash. The finite state projection algorithm for the solution
748 of the chemical master equation. *The Journal of chemical physics*, 124(4):044104,
749 2006.
- 750 [24] F. Neuert, B. Munsky, RZ. Tan, L. Teytelman, M Khammash, and A. van Oudenaar-
751 den. Systematic identification of signal-activated stochastic gene regulation. *Science*,
752 339:584–587, 2013.
- 753 [25] Evan J Olson, Lucas A Hartsough, Brian P Landry, Raghav Shroff, and Jeffrey J
754 Tabor. Characterizing bacterial gene circuit dynamics with optically programmed
755 gene expression signals. *Nature methods*, 11(4):449–455, 2014.
- 756 [26] L. Potvin-Trottier, N. Lord, G. Vinnicombe, and J. Paulsson. Synchronous long-term
757 oscillations in a synthetic gene circuit. *Nature*, 538(7626):514–517, 2016.
- 758 [27] J. Ruess, A. Miliias-Argeitis, and J. Lygeros. Designing experiments to understand the
759 variability in biochemical reaction networks. *Journal of The Royal Society Interface*,
760 10(88):20130588, 2013.
- 761 [28] J. Ruess, F. Parise, A. Miliias-Argeitis, M. Khammash, and J. Lygeros. Iterative
762 experiment design guides the characterization of a light-inducible gene expression
763 circuit. *Proceedings of the National Academy of Sciences*, 112(26):8148–8153, 2015.
- 764 [29] M Rullan, D. Benzinger, G. Schmidt, A. Miliias-Argeitis, and M. Khammash. An
765 optogenetic platform for real-time, single-cell interrogation of stochastic transcriptional
766 regulation. *Molecular Cell*, 70(4):745–756.e6, 2018.
- 767 [30] S. Schmidl, R. Sheth, A. Wu, and J. Tabor. Refactoring and optimization of light-
768 switchable escherichia coli two-component systems. *ACS synthetic biology*, 3(11):
769 820–831, 2014.
- 770 [31] A. Singh and J. Hespanha. Lognormal moment closures for biochemical reactions.
771 *IEEE 45th Annual Conference on Decision and Control (CDC)*. San Diego, CA, USA.,
772 pages 2063–2068, 2006. doi: 10.1109/CDC.2006.376994.

- 773 [32] Abhyudai Singh and Joao P Hespanha. Approximate moment dynamics for chemically
774 reacting systems. *IEEE Transactions on Automatic Control*, 56(2):414–418, 2010.
- 775 [33] V. Stathopoulos and M. Girolami. Markov chain monte carlo inference for markov
776 jump processes via the linear noise approximation. *Philosophical Transactions of the*
777 *Royal Society A*, 371:20110541, 2013.
- 778 [34] P. Whittle. On the use of the normal approximation in the treatment of stochastic
779 processes. *Journal of the Royal Statistical Society Series B Statistical Methodology*,
780 19:268–281, 1957.
- 781 [35] C. Zechner, J. Ruess, P. Krenn, S. Pelet, M. Peter, J. Lygeros, and H. Koepl.
782 Moment-based inference predicts bimodality in transient gene expression. *Proceedings*
783 *of the National Academy of Sciences of the USA*, 109(21):8340–8345, 2012.
- 784 [36] C. Zechner, M. Unger, S. Pelet, M. Peter, and H. Koepl. Scalable inference of
785 heterogeneous reaction kinetics from pooled single-cell recordings. *Nature methods*,
786 11(2):197, 2014.
- 787 [37] C. Zimmer. Reconstructing the hidden states in time course data of stochastic models.
788 *Mathematical Biosciences*, 269:117–129, 2015.
- 789 [38] C. Zimmer and S. Sahle. Deterministic inference for stochastic systems using multiple
790 shooting and a linear noise approximation for the transition probabilities. *IET Systems*
791 *Biology*, 9(5):181–192, 2014.