

RESEARCH ARTICLE

Systematic comparison of modeling fidelity levels and parameter inference settings applied to negative feedback gene regulation

Adrien Coulier¹, Prashant Singh², Marc Sturrock^{3†}, Andreas Hellander^{1†*}

1 Department of Information Technology, Uppsala University, Uppsala, Sweden, **2** Science for Life Laboratory, Department of Information Technology, Uppsala University, Uppsala, Sweden, **3** Department of Physiology, Royal College of Surgeons in Ireland, Dublin, Ireland

†These authors are joint senior authors on this work.

* andreas.hellander@it.uu.se



Abstract

Quantitative stochastic models of gene regulatory networks are important tools for studying cellular regulation. Such models can be formulated at many different levels of fidelity. A practical challenge is to determine what model fidelity to use in order to get accurate and representative results. The choice is important, because models of successively higher fidelity come at a rapidly increasing computational cost. In some situations, the level of detail is clearly motivated by the question under study. In many situations however, many model options could qualitatively agree with available data, depending on the amount of data and the nature of the observations. Here, an important distinction is whether we are interested in inferring the true (but unknown) physical parameters of the model or if it is sufficient to be able to capture and explain available data. The situation becomes complicated from a computational perspective because inference needs to be approximate. Most often it is based on likelihood-free Approximate Bayesian Computation (ABC) and here determining which summary statistics to use, as well as how much data is needed to reach the desired level of accuracy, are difficult tasks. Ultimately, all of these aspects—the model fidelity, the available data, and the numerical choices for inference—interplay in a complex manner. In this paper we develop a computational pipeline designed to systematically evaluate inference accuracy for a wide range of true known parameters. We then use it to explore inference settings for negative feedback gene regulation. In particular, we compare a detailed spatial stochastic model, a coarse-grained compartment-based multiscale model, and the standard well-mixed model, across several data-scenarios and for multiple numerical options for parameter inference. Practically speaking, this pipeline can be used as a preliminary step to guide modelers prior to gathering experimental data. By training Gaussian processes to approximate the distance function values, we are able to substantially reduce the computational cost of running the pipeline.

OPEN ACCESS

Citation: Coulier A, Singh P, Sturrock M, Hellander A (2022) Systematic comparison of modeling fidelity levels and parameter inference settings applied to negative feedback gene regulation. *PLoS Comput Biol* 18(12): e1010683. <https://doi.org/10.1371/journal.pcbi.1010683>

Editor: Pedro Mendes, University of Connecticut School of Medicine, UNITED STATES

Received: January 17, 2022

Accepted: October 25, 2022

Published: December 15, 2022

Copyright: © 2022 Coulier et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All code and generated data underlying the study is available at <https://github.com/prasi372/PipelineforParameterInference>. Data used for some experiments are taken from Hofmann H, Kafadar K, Wickham H. (2011), and is publicly available as a json file at <https://github.com/Aratz/MultiscaleCompartmentBasedModel/tree/master/data>.

Funding: This work has been funded by the Swedish research council (2015-03964 to AH), by

the eSENCE strategic collaboration of eScience, and by the NIH NIH/2R01EB014877-04A1 (to AH). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Author summary

Computational models play a vital role in modern biology and are commonly used to compare theory with data. These models can take different forms, and there is often a trade-off between model detail and the computational resources needed to simulate them. Furthermore a choice must be made regarding how to compare the model output with the available data with several different distance metrics available. The choice of model and distance metric may also be impacted by the amount of available data. Therefore deciding how best to infer model parameters from available experimental data is a challenging problem. In this paper we have developed a computational pipeline designed to systematically evaluate inference accuracy for a wide range of true known parameters. To demonstrate its use, we applied it to a well studied gene regulation model. In particular, we compared a simple model, mid-complexity model and a complex model for several data-scenarios and for multiple numerical options for parameter inference. We believe this pipeline can be used as a preliminary step to guide scientists prior to gathering experimental data. This could prevent experimentalists from gathering unnecessary expensive experimental data or modelers from expending huge computational resources on simulating superfluously complex models.

1 Introduction

Mathematical modeling is an important tool to study gene regulatory networks (GRNs) in single cells. These models exist on many levels of fidelity, ranging from deterministic Ordinary Differential Equations (ODE) to discrete stochastic well-mixed models, to detailed spatial stochastic models. In practice, the choice of a mathematical model has a subjective element to it and the choice depends both on the question under study, the available data, and the computational budget.

There have been many studies where modelers have simulated models at different levels of fidelity to study the same biological system. Often these studies have had the same aim of capturing either qualitative properties or relatively coarse-grained summary statistics, both of which could in principle be captured by any of the models. One commonly studied biological system is the Hes1 negative feedback system. This system produces oscillatory dynamics which are thought to provide a clock for the segmentation of the somites during embryogenesis [1]. This system has spawned models of different levels of biological fidelity, including models comprised of ordinary differential equations [2], delay differential equations [3, 4], partial differential equations (PDE) [5], stochastic differential equations [6] and spatial stochastic simulations [7]. These models were used to simulate the Hes1 system with the aim of producing oscillatory dynamics of the Hes1 protein and messenger RNA (mRNA) with a period of approximately 2 hours. All models were able to produce the desired oscillatory behaviour, but in the ODE case an extra intermediate reaction had to be introduced. This deployment of a variety of modelling approaches is not unique to the Hes1 intracellular pathway, indeed, other pathways which exhibit oscillatory dynamics including the p53 signalling pathway [8, 9] and NF- κ B pathway [10] have similar variety in the fidelity of models produced to capture simple summary statistics. It is not clear in these studies whether the data alone necessitated a higher fidelity model. Outside of the intracellular modelling space, there have been many similar studies in the cancer modelling space that have used models of differing fidelity to capture tumour growth time series data. These models have often been either ODE or PDEs or even agent based models with huge differences in computational expense and model complexity [11].

Again though, it is not clear if and when the available data clearly motivated one approach over another. A central issue is whether the models are developed with the goal to infer parameter values from experimental data or if they are developed to predict some future state of the system. In the latter case, some have argued that the precise parameter values may be of less concern and may even take biologically infeasible values so long as the model is a good predictor of the system [12].

In some situations the question at hand clearly motivates a spatial model over a well-mixed one, such as if we were to study the effect of location and numbers of membrane receptors on downstream signaling in a signaling cascade. However, in situations where the question could in principle be addressed using models of various different levels of fidelity, and when the choice is driven by observed data, the situation is less clear. In particular, it is common that quantitative experimental observations are “well-mixed” in nature, such as time series data of total protein or mRNA counts in cells, or distribution data from large cell populations from e.g., fluorescent activated cell sorting (FACS). An interesting question then is in which situations it is motivated to use a high-fidelity spatial model even though the observations are more coarse-grained in nature. This question can be expected to depend critically on what the goal of the modeling is, for example, is the goal simply to capture the qualitative trends in data, or is the goal to identify model parameters that are in good quantitative agreement with the true biochemical and physical parameters, such as diffusion constants and kinetics rate constants? This is fundamentally a hard problem to address due to the lack of ground truth (both for the model and for the parameters). But if enough computational power is available, we can scan through the space of possible parameters and evaluate how inference would perform were the postulated parameters the ground truth (synthetic data). While the true posterior distribution is out of reach, we can still analyze how the estimated posterior relates to the true parameters for various choices of models and for various types of observations. In this paper, we develop such a computational pipeline in which we generate synthetic ground truth data using a high-fidelity spatial model (simulated using Smoldyn [13]) for a wide range of possible true parameters (controlling the degree of diffusion limited dynamics), then systematically compare inference tasks for the spatial model, the well-mixed model and a coarse-grained multiscale model.

However, since it is necessary to use likelihood-free, or simulation-based approximate inference, there are several numerical considerations for accurate inference, apart from the question of the simulator and how much data is available. In particular, ABC, the most widely used method, relies critically on the chosen distance metric and summary statistics. In the end, both the model fidelity, the data, and the numerical choices for parameter inference need to be studied simultaneously. The fact that ABC requires a large number of potentially expensive simulations becomes a practical hindrance to such large-scale studies. Here we train Gaussian Processes (GPs) to approximate the distance metric values, and are in this way able to substantially reduce the computational cost of running the pipeline.

While ideally models would be constrained using sufficiently rich experimental data, in practice there are often limitations to the amount and kind of data available. The modeling of cellular functions requires sensitive measurement of various molecular species, such as mRNA and proteins. For data captured at the intracellular level, there are often trade offs between the richness of the time series, the number of replicates and the level of spatial information captured. Traditionally used population-averaged techniques like Western blots, Northern blots and enzyme-linked immunosorbent assay (ELISA) do not capture the important details at the single-cell level. More modern techniques such as mass spectrometry (MS) generally lack the sensitivity to detect the small amounts of proteins present in individual cells [14, 15]; however, recent developments in MS have made progress towards uncovering single-cell proteomes [16]. Flow cytometry and mass cytometry (e.g., CyTOF) can detect proteins in single cells, but

developing sample standards for quantification has proved challenging [17]. In very recent years, the digital proximity ligation assay (dPLA) was developed. dPLA provides the ability to take direct digital measurements of protein and mRNA copy numbers in single mammalian cells [18]. In dPLA, digital PCR (dPCR) was used to quantify proteins detected with a pair of oligonucleotide-tagged antibodies called PLA probes. Previously published PLA methods enabled multiplexed simultaneous protein and mRNA measurements from single cells. It was noted though that their quantitative polymerase chain reaction (qPCR) readout limits the sensitivity of the measurements [19]. The use of the dPCR readout provides significantly improved resolution and limits of detection [20], which allows direct quantification of protein copy numbers in individual mammalian cells. This advance, along with the use of dPCR for mRNA quantification allows for simultaneous measurement of both mRNA and protein albeit at low time resolution (with readings captured every 10 minutes) [21]. All of these experimental technologies come with different costs and can require different levels of experience to gather, hence it is important to know when enough data is captured to warrant the use of a more sophisticated model. To that end, in this study we use a synthetic data set that mimics the cutting edge of what is possible experimentally, i.e. simultaneously capturing mRNA and protein copy number data at the single cell level for various cells and time points, to address the question of how much data is needed to warrant a higher fidelity model.

There are a growing number of studies investigating parameter inference in the presence of different kinds of data. In [6], it was demonstrated that MCMC methods for stochastic differential equations provide practical algorithms for estimating the parameters of simple dynamic regulatory and signaling systems even when the time series data are coarse. Furthermore, it was reported that if one had access to good quality temporally resolved data, one could also obtain information about stochastic modeling parameters and population sizes. In [22], Kur-sawe *et al.* investigated the performance of parameter inference using a vertex model for cell mechanics and image data. They showed that estimating the noise by having several realizations of the observed data was critical for reliable inference. Harrison *et al.* [23] quantified the effect of noise and data density on the posterior estimate and compared ABC to the particle Markov Chain Monte Carlo method (pMCMC). They showed that, when applicable, pMCMC performed better, although ABC was more general and more easily parallelizable.

Some doubts have been raised regarding the validity of the posterior distributions generated with ABC [24, 25]. Specifically, and although they are fundamental requirements for the well-posedness of the inverse problem, identifiability [26] and sufficiency [27] may not be attainable in practice. Yet, this need not be the end of the story. Firstly because identifiability *can* be demonstrated for some simpler models [28], and secondly because insights can still be gained from models where the true solution is known [29].

In this paper, we show that such an approach is indeed feasible for ABC. We propose a computational pipeline to evaluate the accuracy of ABC in different scenarios. Specifically we evaluate the performance of ABC throughout the parameter space while keeping the cost down using Gaussian processes to approximate the distance metric values between simulated and observed data. We then analyze the accuracy of the resulting posterior distribution with respect to the true parameters. We can then repeat this procedure for different models, summary statistics or even data sets (e.g., when measuring proteins or when measuring mRNA) and determine which setup gives the best performance. This preliminary analysis can then be used to guide practitioners to choose models and design experiments. Contrary to other analytical approaches to guide modelers and optimize experiments [30, 31], our approach is purely computational and does not require the ability to derive analytical formulae from the model formulation. Hence, models of arbitrarily high complexity can be used. With this approach we can answer questions such as:

- How much data is needed to reach a given level of inference accuracy?
- What features are worth measuring in an experiment if the goal is to identify parameters?
- Which modeling fidelity is appropriate?
- Which summary statistics or distance metric should be used?

We exemplify this procedure through different scenarios based on a canonical negative feedback gene regulation network motif for three models of various fidelities.

2 Results

In what follows we are concerned with likelihood-free inference, and in particular how different modeling fidelity levels, the amount of data, and inference settings impact our ability to accurately infer parameters when the observed data comes from a high-fidelity spatial stochastic model of negative feedback gene regulation. In mathematical terms, the setting is as follows: let $y(d)$ be a random variable/process representing the observed data from an experimental protocol d , and let θ be a vector of biophysical parameters that one wishes to estimate based on y . A model, as we define it in this manuscript, is a stochastic mapping f (taken from a family of models F) between θ and y , that is, $y = f(\theta)$. ABC allows one to approximately solve the inverse problem to estimate θ by sampling an approximate posterior distribution $\pi_{ABC}(\theta | f, g(y), d)$ where $g(y)$ is a vector of summary statistics derived from y . In the Methods section we detail the the models f considered.

2.1 A computational pipeline for systematic parameter inference evaluation

Given that we want to be able to use models of high complexity, it is usually not possible to rely on *a priori* mathematical analysis to determine the validity of ABC for a given setup. For example, information on system identifiability is often impractical to approximate numerically unless an exact solution to the stochastic model is available, and numerical identifiability needs to be studied empirically. We can, however, systematically evaluate the performance of ABC using synthetic observed data sampled throughout the parameter space. Given the high computational cost, we use an approximation scheme based on Gaussian processes. This makes it possible to only generate the data once, prior to inferring parameters with ABC in various configurations.

We have developed a pipeline made up of two main parts, a data generation step and a parameter inference step. [Fig 1](#) illustrates how these parts are combined. Starting with a prior distribution, we first simulate ground truth data for one parameter point using the highest model fidelity ([Fig 1A](#)) and then generate simulated data across the entire parameters range of the prior with the model meant for Bayesian inference ([Fig 1B](#)). We then use this data to approximate the distance between observed (synthetic data generated with the highest fidelity) and simulated data. This approximate distance map ([Fig 1C](#)) is then used to perform parameter inference without the need to simulate more data during this process ([Fig 1D](#)).

This entire pipeline can then be executed multiple times using synthetic observed data sets from different regions of the parameter space. The core idea is to generate the data once and then reuse it first as synthetic observed data, and then as training data to approximate the distance metric used in ABC. Thus, the same data set can be reused in various configurations, e.g. with different summary statistics or by subsampling the data in terms of number of trajectories or time samples. This in turn generates many posterior distributions which can be compared to the true parameters. By systematically measuring the discrepancy between posterior

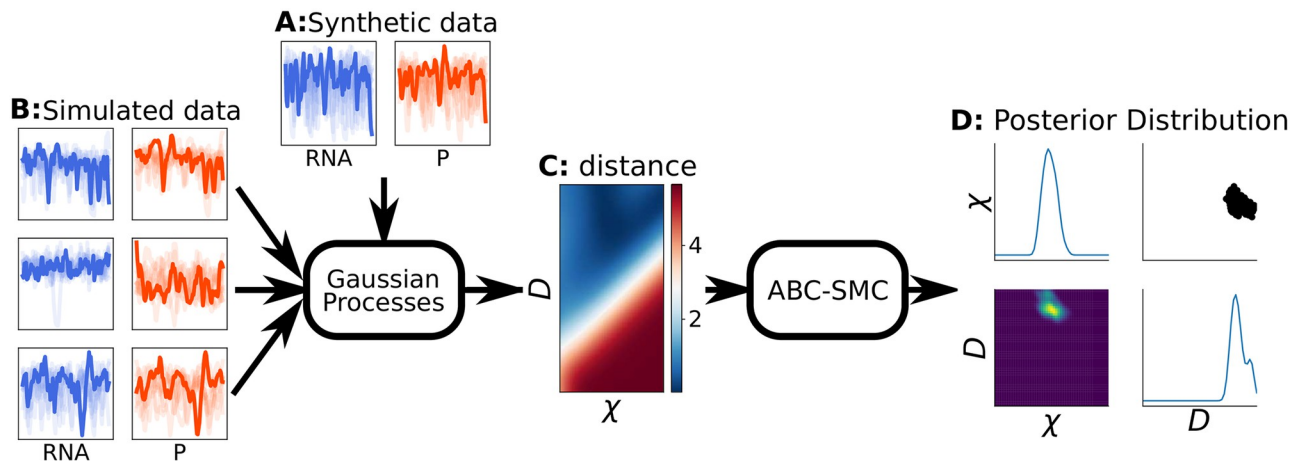


Fig 1. Parameter inference pipeline based on Gaussian processes. A map of distance metric values based on one synthetic data set (A) is trained using all other simulated data sets (B). This distance metric map (C) is then used in pyABC to infer the parameters of the synthetic data set using ABC-SMC (D), i.e. in this way we avoid sampling additional data points during ABC inference points using expensive simulations.

<https://doi.org/10.1371/journal.pcbi.1010683.g001>

distribution and true known parameters, we can then build up error maps showing regions in parameter space where inference performs better or worse.

Analyzing the thus obtained error maps from various combinations of models, summary statistics and amounts of data can offer insights into which combination would be more or less likely to perform well when used against experimental data. We believe that this pipeline can be used as a pre-processing step to calibrate inference pipelines when the true parameters and model is unknown, assuming we believe that the model of the highest fidelity is fundamentally the best representation of reality of the simulators we consider. As we will show, the pipeline aims to identify the appropriate model fidelity to use for a given data set in systems biology projects.

In what follows we use the above described pipeline to investigate different scenarios for inferring parameters. We first detail the model used to generate the data in Section 2.2 and then elaborate on how we measure the posterior error in Section 2.3. We then investigate how the amount of data in terms of time sampling density, throughput and observed species (protein, mRNA or both) can influence accuracy and in what situations we clearly benefit from using the higher spatial model fidelity, as motivated by data. We then investigate where in the parameter space each model performs best, and which combination of model and distance metrics gives the best results overall.

2.2 Design of the computational experiments—GRN models, synthetic data generation and distance metrics

In [32], we studied a negative feedback motif motivated by the Hes1 GRN, in which a gene represses its own expression through a negative feedback loop: mRNA is transcribed in the nucleus and diffuses out into the cytoplasm where it is translated to proteins. These proteins then diffuse back to the nucleus and repress the gene. This process is illustrated in Fig 2. The delay between the moment mRNA is produced and the moment proteins diffuse back into the nucleus and bind to the gene tends to generate oscillations in gene expression level. These chemical reactions are described in Eqs 1–5 while their parameters are summarized in Table 1.

Three approaches were used to model this network: the first approach consists of a detailed particle model based on the Smoluchowski diffusion limited equations. We use the widely-

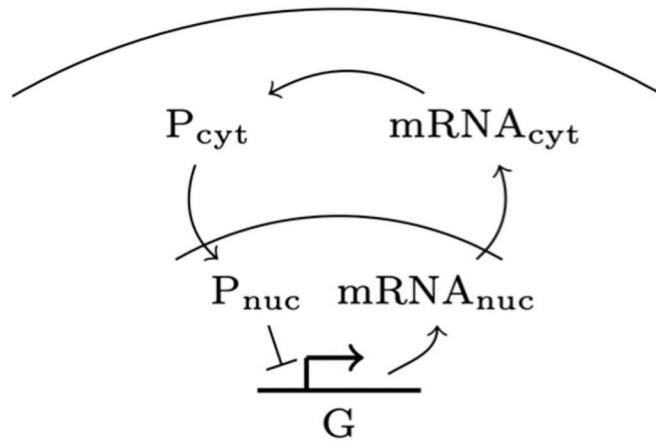


Fig 2. Sketch of the genetic motif studied in this article. A gene, placed at the center of the nucleus of the cell, transcribes mRNA. mRNA then diffuses out of the nucleus and into the cytoplasm, where it is translated to proteins. These proteins then diffuse back into the nucleus, where they repress the expression of the gene.

<https://doi.org/10.1371/journal.pcbi.1010683.g002>

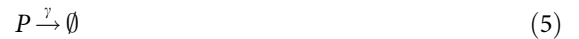
used software Smoldyn [13] to simulate the model, thus we refer to it simply as *Smoldyn* in the remainder of this paper. The second approach is a cheap multiscale approximation from [32]. Here the cell geometry is divided into two compartments (cytosol and nucleus) which are themselves considered to be well-mixed. Transition rates between these two compartments are then derived using hitting-time analysis on the Smoluchowski model, thus capturing some of the spatial effects. The model is simulated with the standard SSA and is referred to as the Compartment-Based Model (CBM). Finally, the entire cell is considered to be well-mixed and the model (WMM) is simulated using SSA. Note that we use diffusion-limited reaction rates in the association reactions between the gene and protein, thus all three models explicitly involve all physical constants, enabling direct comparison. For all SSA simulations we use the software Gillespy2 in the StochSS suite of tools [33, 34].



Table 1. Base parameters as presented in [32]. The parameters to be estimated with Bayesian inference are highlighted in grey and are varied over several orders of magnitude in the synthetic data sets. Parameters μ , κ and γ are varied simultaneously by multiplying them by a common variable, noted χ . This variable and the diffusion constant are the targets to be inferred.

Parameter	Description	Localization	Base Value	Unit
R	cell radius		6.0	μm
r	nucleus radius		2.5	μm
D	diffusion constant		0.6	$\mu\text{m}^2 \text{min}^{-1}$
k_a	binding rate	nucleus	1.00×10^9	$\text{M}^{-1} \text{min}^{-1}$
k_d	unbinding rate	nucleus	0.1	min^{-1}
μ	transcription rate	nucleus	3.0	min^{-1}
κ	translation rate	cytoplasm	1.0	min^{-1}
γ	degradation rate	entire cell	0.04	min^{-1}

<https://doi.org/10.1371/journal.pcbi.1010683.t001>



Our goal is to systematically evaluate the quality of parameter estimates inferred with ABC, depending on which of the three models are used to simulate the data, and on the amount of data available. We use the well-established Sequential Monte-Carlo (SMC) variant of ABC as provided in the pyABC Python package [35]. The prior is set to a log-uniform distribution between 0.25 and 16 for χ and between 0.0039 and 16 for the diffusion constant. For each setup, we infer parameters from 256 different synthetic data sets generated with Smoldyn (i.e., we sample from a 16×16 grid of the diffusion constant, χ parameter space), and compare the posterior distributions to the true parameters. This gives us a map of inference performance for the systems ranging from the strongly diffusion-limited regime to the well-mixed regime.

We consider three distance metrics to measure the discrepancy between candidate particles and the observed data:

1. First we consider four common summary statistics, namely the mean value, the minimum value, the maximum value and the standard deviation of a trajectory. We then take the expected value over all trajectories and for both species, and compute the distance between simulated data and synthetic data using the L_2 norm. This setting represents a likely first setting formulated manually by a modeler. We refer to this setting as *naïve statistics*.
2. Second, we select optimal summary statistics using the AS algorithm described in Subsection 4.2. The selected statistics are: the longest strike below the mean, the longest strike above the mean, the mean absolute change, the maximum, the minimum and the variance. The distance is computed as in 1 above. We refer to this setting as *optimized statistics*.
3. Third, following [36], we observe the distribution of molecular counts for each species and at every time point, i.e. we compute a histogram density approximation of the cumulative density function (CDF) based on the observed trajectories for each time point, and then compute the average Kolmogorov distance between the two data sets over these distributions. We refer to this setting as *Kolmogorov distance*.

The data used here are taken from a previous study [32] and is publicly available as a .json file at GitHub, <https://github.com/Aratz/MultiscaleCompartmentBasedModel/blob/master/data/data.zip>. All the code used is available on GitHub at <https://github.com/prasi372/PipelineforParameterInference>. For each pair of diffusion and reactivity coefficient, the data set contains 64 trajectories over 100 time samples for the two species of interest in the system (namely mRNA and proteins), and for each of three models. A burn in period was used at the beginning of each simulation to make sure all trajectories are uncorrelated from the initial condition.

In the data-scenarios detailed in this study, the pipeline is run separately for each model (WMM, CBM and Smoldyn) and then for each of our 256 synthetic data sets. Fig 3 illustrates this process. In total, 768 inferences are performed for each setup. In each scenario we vary the amount of data and the distance metric used each time. All the computations were run on Rackham, a high performance computing cluster provided by the Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX). Each nodes consists of two 10-core

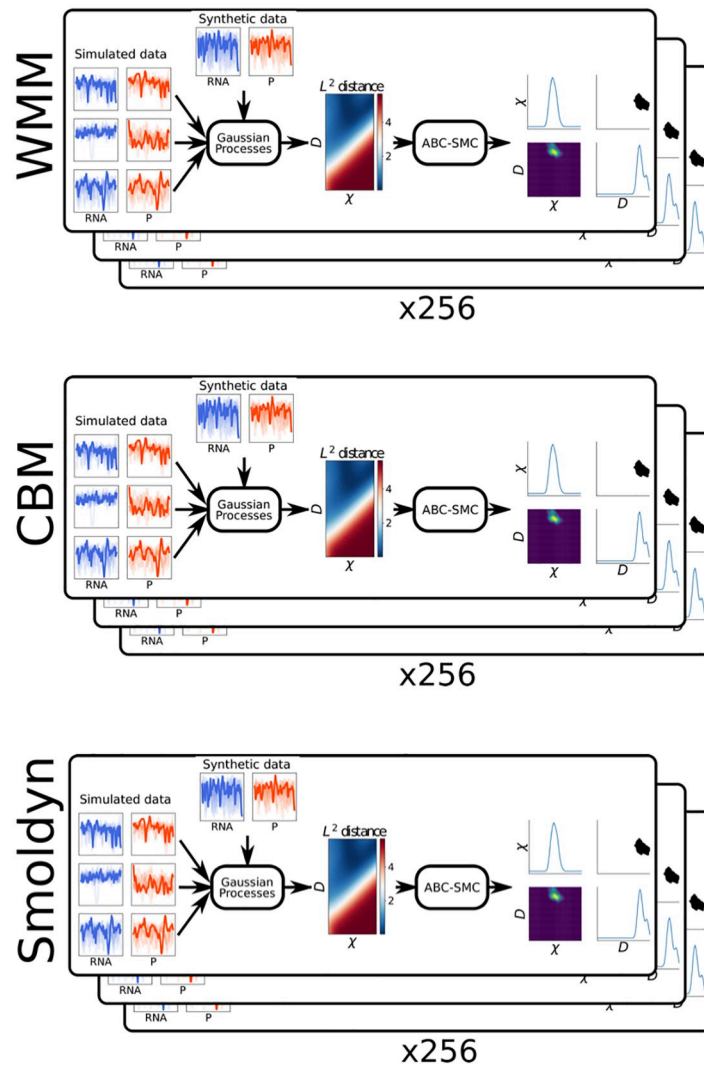


Fig 3. Illustration of the computational experiments performed. For the different data-scenarios and for each combination of distance metric and amount of data, we execute the pipeline using all 256 synthetic data sets as observed data, and for each of the three models.

<https://doi.org/10.1371/journal.pcbi.1010683.g003>

Xeon E5–2630 V4 processors at 2.2 GHz and 128 GB of memory. Running the pipeline for one setup (i.e. 256×3 inferences) took approximately 200 core-hours. Specifically, this is the cost of generating the results as presented e.g., in Fig 4. We emphasize that executing our pipeline is only made possible by the use of Gaussian Processes to approximate the distance metric values between simulated and observed data.

2.3 Depending on distance metric, a detailed spatial model can be motivated also by non-spatial observations

In this section, we compare the performance of the three model fidelities for parameter inference using the complete observed dataset (64 trajectories sampled at 100 time points) for each of the three distance metrics.

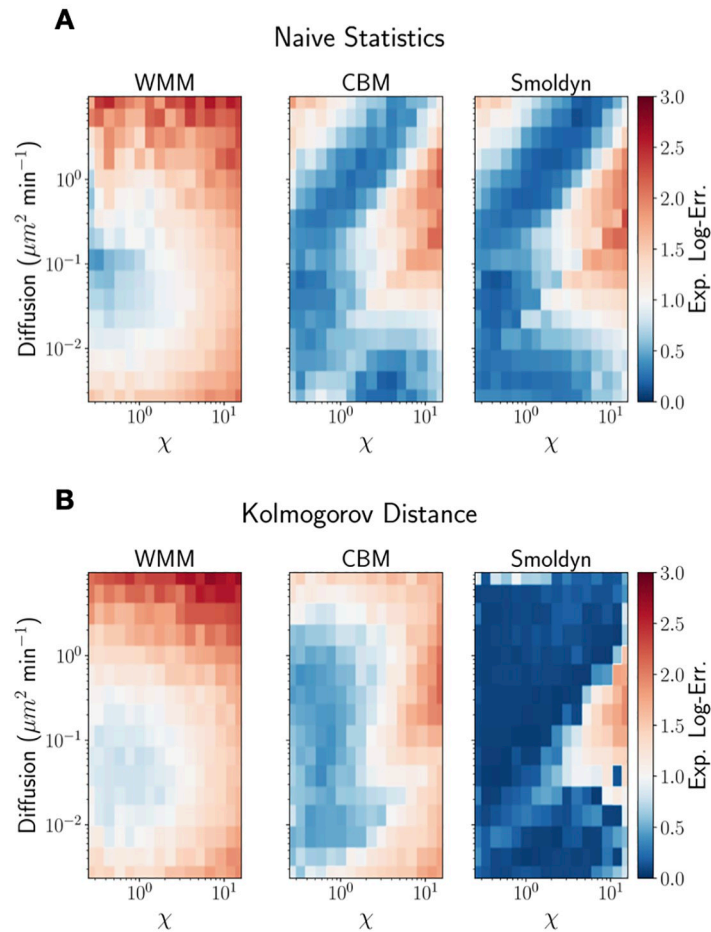


Fig 4. Heatmaps of expected log-error for different combinations of the diffusion coefficient (D, y-axes) and reactivity coefficient (χ , x-axes). The results presented in (A) are based on summary statistics while those presented in (B) are based on the Kolmogorov distance. The title of each heatmap corresponds to which model fidelity was used. The CBM performs almost as well as Smoldyn when using summary statistics. When using the Kolmogorov distance, Smoldyn outperforms the two other models.

<https://doi.org/10.1371/journal.pcbi.1010683.g004>

In cases where the true parameters are known (e.g., as in the considered test problem), it is possible to compute the error between the estimated parameters and the true parameters, or report if they fall within given confidence intervals.

Since the considered parameter space spans several orders of magnitude, computing the relative error in the mean posterior parameters is not very meaningful, i.e., when the true parameters are close in magnitude to the lower bound describing the parameter space, the relative error tends to be large, while if the true parameters are close to the upper bound describing the parameter space, the relative error will be close to 1. Instead, we report the *expected log-error* with respect to the posterior with the following formula:

$$\varepsilon = \mathbb{E}_{\hat{\theta}} \left[\|\log_{10} \hat{\theta} - \log_{10} \theta_o\|_2 \right].$$

Here θ_o are the true parameters and $\hat{\theta}$ is the set of samples comprising the estimated posterior. We therefore compare the estimates in terms of order of magnitudes from the true parameters. This measure also penalizes wide posteriors, which will have a larger expected error than

tighter posteriors. This combined metric for accuracy and uncertainty allows us to compare the 256 inferences performed in each experiment and for each model. The preference towards tighter posterior distributions enforced via the expected log error metric also allows identification of regions in the parameter space where parameter inference is unlikely to accurately yield the true parameters. We note that there is no ideal error metric, and the benefits and limitations of the considered metrics should dictate usage on a case by case basis. In summary, the chosen error measure takes smaller values for inferences with tight posteriors around the correct expected parameter point and increases both with bias and spread of the posterior (lower inference quality). Fig 4 shows these expected log-error maps for all three models, the naïve statistics, and the Kolmogorov distance. For the sake of comparison, we also investigated using the frequently employed root-mean-square-error (RMSE) in place of our expected log-error and present the RMSE maps in S6 Fig. While there are quantitative differences, the same general trends persist. We emphasise that users of our pipeline should confirm their results are not artefacts of one particular error metric and should check consistency of conclusions with multiple error metrics.

As seen from Fig 4, we obtain substantially better inference performance when using Smoldyn as the simulator compared to when using the WMM both when using the naïve set of statistics and the Kolmogorov distance, even in the well-mixed regions of parameter space. We also see that using the Kolmogorov distance together with the detailed spatial simulator gives the best inference quality overall. This answers one of our initial questions—it is best practice while performing parameter inference to use models with explicit spatial detail even though experimental observations are more coarse-grained. The nature of observation and the distance metric employed have a large influence here: when using summary statistics rather than Kolmogorov distance, the CBM leads to approximately the same inference quality as Smoldyn, suggesting that the model is able to capture the critical spatial effects on the statistics. We also clearly see the variation of inference quality throughout parameter space—in simple words, some regions are easier to infer than others, and for some regions the expected log error is unacceptably high (2 orders of magnitude or more) even for the best configurations. In particular, with this dataset size we are not able to accurately identify parameters even if using the ground truth model in those regions. We suggest that computing this type of error map will also aid the development of models using real experimental data—if the inferred parameter falls in a region of the map where the error is large, it is a good indication that care needs to be taken in interpreting the results. In one selected case (well-mixed model with naïve statistics), we estimated the inference error due the GP surrogates by running ABC with the true model on 32 randomly selected synthetic datasets. The results presented in Fig 4 varied by $17\% \pm 10$.

Although the error map in Fig 4 provides a detailed view of the regions where the inference has the lowest error, the level of detail makes it hard to compare two models quantitatively. For instance, in Fig 4, it is unclear how much more accurate Smoldyn is compared to the CBM. Thus, when comparing results for two different settings, we build an enhanced box plot, referred to as a ‘Boxen plot’. Contrary to the regular box plot, where only the quartiles are shown, the extended box plot also displays the next 2^n -quantiles above the upper quartile and below the lower quartile, thus giving a better view on the distribution of tail values [37]. Fig 5 illustrates such visualization technique. By representing the error in this way, we trade local information in parameter space for easier, global comparison.

Here we also show results using the optimized statistics. Although they do improve inference for the synthetic data sets where it was already accurate with the naïve summary statistics (Fig 4), overall, they do not bring significant improvements in inference quality, and will in some cases also lead to worse performance than the naïve set. This illustrates the challenge in choosing good statistics.

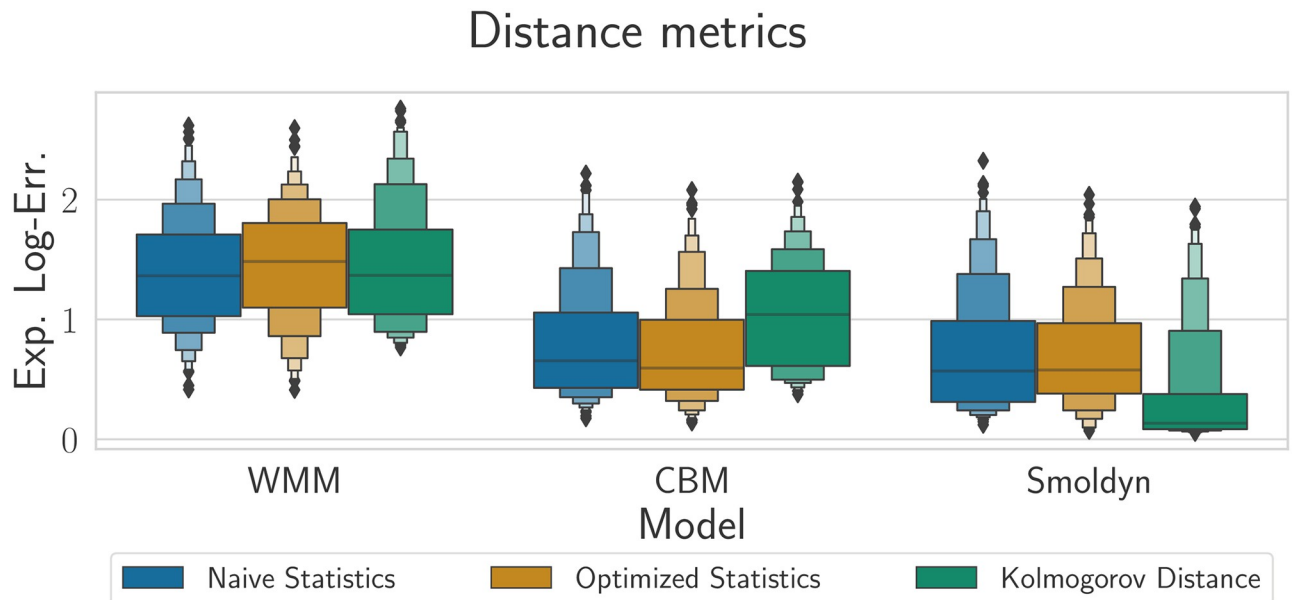


Fig 5. Boxen plots of expected log-error for WMM (blue), CBM (yellow) and Smoldyn (green) models respectively when used with Naive statistics, optimized statistics and Kolmogorov distance metrics. There is general convergence when adding more details to the model, although for a given model, no metric is consistently more accurate than the others.

<https://doi.org/10.1371/journal.pcbi.1010683.g005>

For the WMM, inference quality does not depend on the distance metric used, suggesting that the model error is dominating. When it comes to the Kolmogorov distance, we see minor improvements when going from the WMM to the CBM, and indeed distance metrics based on summary statistics outperform the Kolmogorov distance for the CBM. This can be understood in light of the approximation properties of the CBM: in [32], we showed that, when using summary statistics, the CBM could better approximate Smoldyn than when using Kolmogorov distance. There is a significant improvement in using the CBM versus the WMM. When it comes to Smoldyn, however, the Kolmogorov distance largely outperforms both naïve and optimized summary statistics. In fact, this is the only case where the parameters were inferred with acceptable accuracy over the majority of the parameter space. This suggests that this distance is more robust to outliers than summary statistics based metrics. Thus, when combined with an accurate model, this distance is capable of producing more accurate results when inferring the parameters.

We conclude with a comment on the accuracy of inference versus the local approximation quality of the coarse grained model alternatives. In [32], we showed that, when using summary statistics, the CBM could accurately approximate Smoldyn throughout the entire parameter space considered in this study. When using the Kolmogorov distance, we showed that it was only highly accurate in the upper left half of the parameter space, namely when diffusion is high and chemical reactions are slow. Regardless which distance metric was used, the WMM was only accurate in the upper left half of the parameter space. Inspecting Fig 4 and comparing it to the Fig 3 in [32], we can see that the accuracy of inference does not directly depend on the accuracy of the coarse-grained model at location of the true parameters, i.e. inference can be relatively accurate when the coarse-grained model is not a good approximation, and it can also be very inaccurate even when the coarse-grained model is a good approximation. That being said, we can still see that parameter inference becomes more accurate when a more detailed model is used, regardless of the distance metric in use.

All in all, this shows that accurate inference does not depend only on how well the coarse-grained model fits the detailed spatial model at the true parameters, but rather how well it globally fits the fully detailed model over the parameter space.

2.4 FACS-like data with Kolmogorov distance is able to discriminate between low- and high model fidelities even for coarse time samples

In this section, we compare the performance of our three models in terms of parameter inference in a fluorescence activated cell sorting (FACS) like setting. Flow cytometry is a particularly powerful tool because it allows a researcher to rapidly and accurately acquire population data related to many parameters from a heterogeneous fluid mixture containing live cells. Flow cytometry is used extensively throughout the life and biomedical sciences, and can be applied in any scenario where a researcher needs to rapidly profile a large population of loose cells in a liquid media. FACS differs from conventional flow cytometry in that it allows for the physical separation, and subsequent collection, of single cells or cell populations [38]. FACS is useful for applications such as establishing cell lines carrying a transgene, enriching for cells in a specific cell cycle phase, or studying the transcriptome, or genome, or proteome, of a whole population on a single cell level.

In order to mimic a FACS experimental setup, we followed the method of [36]. In [36] measurements are taken at regular intervals. For a given interval, the Kolmogorov distance between the observed distribution and the simulated distribution is computed and the average distance across all measurements is reported. Our data set contains a total of 100 time measurements (one every ten minutes). We reduce this data set to contain only 12, 6 and 3 measurements. We then run our inference pipeline on each coarsened data set with each model. The expected log-error is reported in Fig 6.

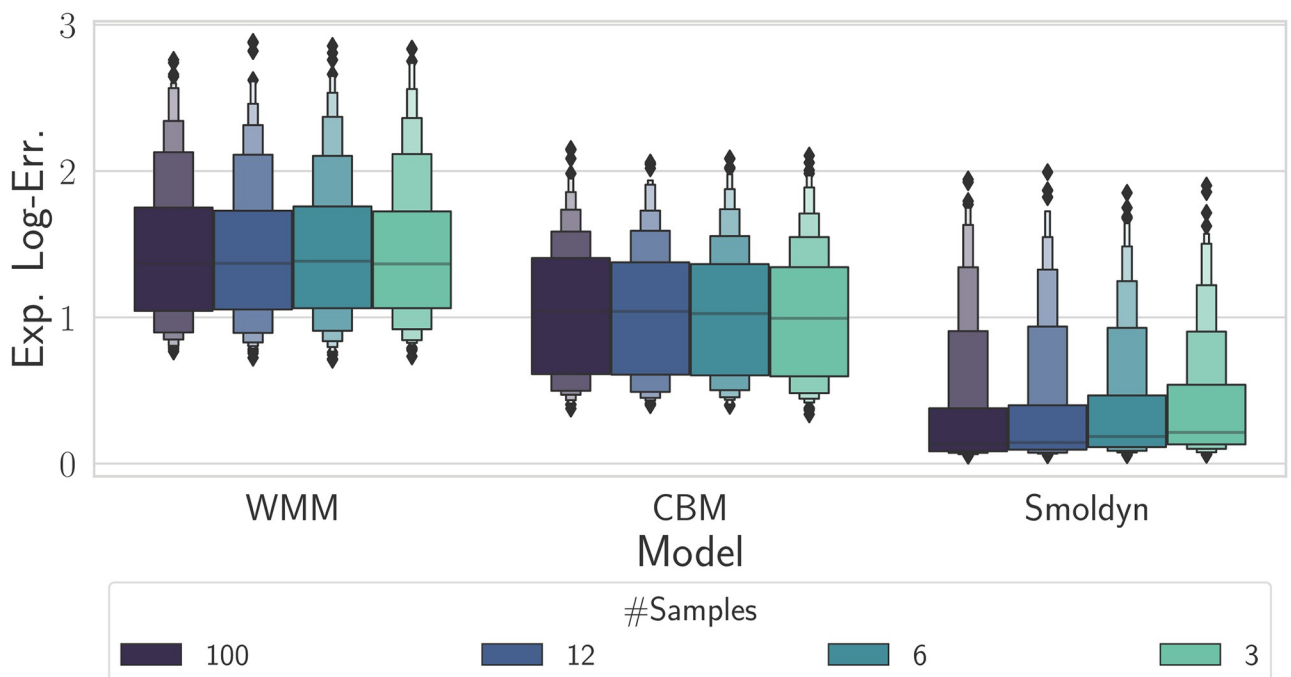


Fig 6. Boxen plots of expected log-error based on the Kolmogorov distance when increasing the number of time samples from 3 to 100. The colour of the boxen plot corresponds to the number of time samples used, with the lighter the colour corresponding to fewer time points. The error only decreases substantially when Smoldyn is used, suggesting model error is the limiting factor in the case of the WMM and the CBM.

<https://doi.org/10.1371/journal.pcbi.1010683.g006>

Strikingly, we found that, while Smoldyn had the lowest error, reducing the amount of data available had only little effect when using the WMM or the CBM, suggesting that model error was the limiting factor. Although the CBM achieved a lower error level than the WMM, it did not improve with larger numbers of samples, suggesting again that the model is not able to accurately capture the data. Put another way, there was no statistically significant improvement in the error distribution across the considered parameter space for 3 samples versus 100 samples. Increasing the number of time samples only had some effect in the case of Smoldyn. Overall the results of this section suggest that using more time samples is only beneficial if enough computational power is available to use the detailed model.

We also investigated the impact of reducing the number of samples per time point as in [39]. We present the results of this in S5 Fig. We found that decreasing the number of samples per time point had a greater impact when using the Kolmogorov distance—most notably in conjunction with the CBM and Smoldyn models where the accuracy diminished substantially.

2.5 Protein measurements are more important for inference accuracy than mRNA measurements

Depending on time and/or budgetary constraints, researchers may have access to mRNA data, protein data or both mRNA and protein data. While some dynamical models have been used to infer networks using only mRNA data [40], others have been constrained using both mRNA and protein data [41]. However, to the best of our knowledge the relative importance of mRNA and/or protein data for model inference is not well studied.

In this section, we compare the performance of our three models in three different data scenarios using three different distance measures. In terms of data, we compare using only mRNA data, using only protein data or using both mRNA and protein data. In terms of distance measures, we compare naïve statistics, optimized statistics and the Kolmogorov distance metric. We present our findings for this section in Fig 7.

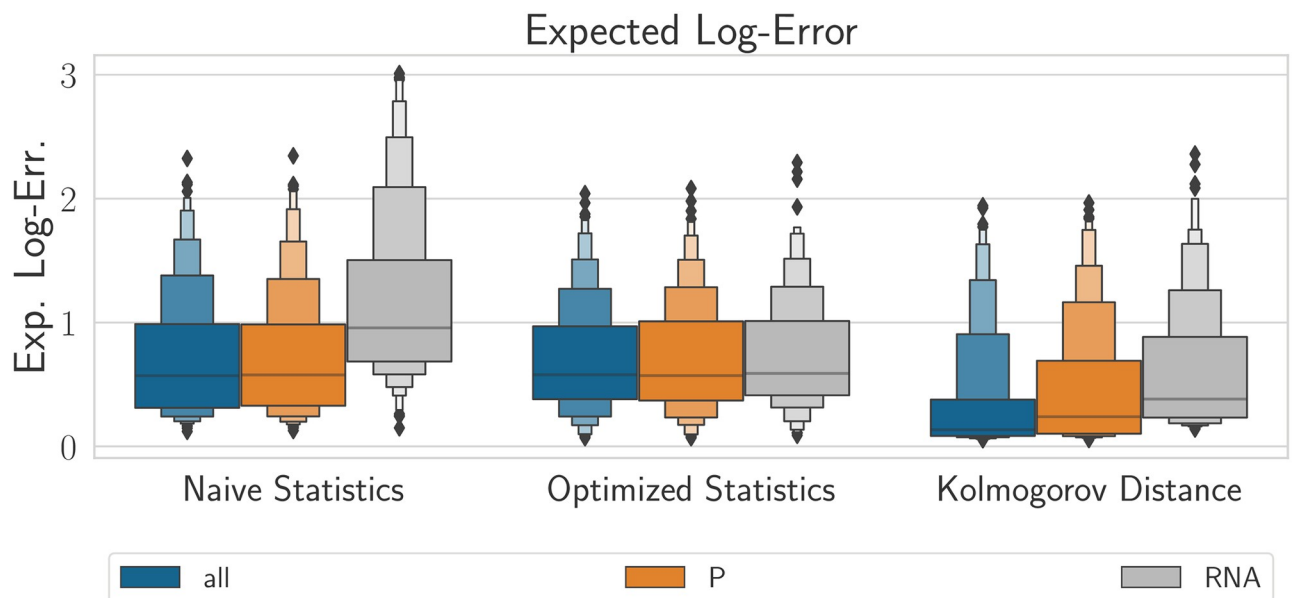


Fig 7. Boxen plots of expected log-error when using all species (blue), only proteins (orange) or only mRNA (grey) based on all three distance metrics (as shown on x-axis). The Smoldyn model was used for this comparison. Using only RNA tends to decrease the accuracy of the inference. S2 Fig shows the same plots with the CBM and the WMM.

<https://doi.org/10.1371/journal.pcbi.1010683.g007>

Analyzing data simulated using Smoldyn, we found that only measuring mRNA levels had worse accuracy than measuring only protein levels or measuring both mRNA and protein levels. We also found that collection of mRNA and protein data was only beneficial when used in conjunction with the Kolmogorov distance. In contrast, we found that data simulated using the WMM or CBM showed little difference between all three distance metrics for all three data scenarios (S2 Fig). The results here suggest that model error is the limiting factor with respect to inference accuracy and that only measuring one species is enough if only the WMM or CBM models can be used. These conclusions are valid for our particular choice of relative true parameters—if another baseline set are of interest, the inference pipeline would need to be re-run to confirm the relative importance of mRNA versus protein measurements. We would like to point out that the particular choice of reference value here is motivated by a previous study of the Hes1 GRN where the parameters were observed to result in spatial simulations agreeing qualitatively and quantitatively with experimental data [7].

3 Discussion

In this manuscript, we built a computational pipeline to systematically investigate the performance of ABC with respect to the choice of model, the nature and amount of observed data and the way we compare true and simulated data in likelihood-free inference with ABC (i.e., the distance metric). We applied this pipeline to a canonical model of negative feedback regulation and several experimentally motivated data-scenarios. We showed that the pipeline can be used to reveal insights into which combination of model, amount of data and distance metric can be expected to lead to the most accurate inference result.

When analyzing inference error over the parameter space, one can then identify areas where inference is most likely to be accurate when executed against experimental data. Using the complete set of observations (64 trajectories each with 100 equidistant temporal observations) we found that parameter inference was most accurate overall using the most detailed model (as expected) and the Kolmogorov distance (as opposed to summary statistics).

One key question we sought to answer using our simulation-driven inference pipeline was under which inference conditions we were able to see a clear benefit from using the full spatial model (Smoldyn). Since the true synthetic data was generated with Smoldyn we expected that, if inference conditions allowed, using that simulator should lead to superior inference accuracy. In our experiments, we observed a clear benefit from using Smoldyn over the WMM model when using both naïve and optimized summary statistics (Fig 5). Surprisingly, when using these summary statistics, there was no significant benefit in using Smoldyn over a simpler multiscale compartment model that incorporates some spatial features of the simulation. However, when distribution data was used with the Kolmogorov distance metric instead of the mean values, only the fully detailed spatial model managed to leverage the information in the data (Fig 5 and S4 Fig).

From an experimental point of view, another interesting question we sought to answer using our pipeline is whether it is desirable to observe mRNA, protein, or both, in order to minimize parameter inference error. In order to answer this question, we conducted experiments for the specific case of the negative feedback model and observed that solely measuring mRNA levels result in highest parameter inference error. Measuring only proteins was typically as accurate as measuring both species except when Smoldyn and the Kolmogorov distance were used, in which case a relatively small gain was seen from observing both species.

We considered two sets of summary statistics, a naïve pool consisting of typical statistics a modeler might choose for a first attempt at inference, and a set of statistics obtained by state-

of-the-art summary statistic selection from a larger pool of time series features. These optimized summary statistics did not show major improvements compared to the naively chosen statistics, highlighting the fact that statistics selection is a challenging problem. It is possible, of course, that there exists a possible set of statistics that performs better. Nevertheless, in our experiments the Kolmogorov distance approach gave more robust results. Here, using the Kolmogorov distance to compare summary statistics (as opposed to taking expectations) shows a gain for the spatial model, with the standard Kolmogorov distance measure led to more accurate results still, see [S4 Fig](#). An interesting future direction would be to also compare to an emerging class of methods that automatically learn good or optimal statistics by training a regression model [42] to predict the posterior mean parameters [43–45]. This approach is particularly useful in cases where optimal summary statistics may not even exist in the candidate pool to select from, but comes with an overhead of requiring a sizeable amount of training data (pairs of the form y, θ) to obtain the regression models.

Executing the pipeline over such a large section of the parameter space we used was only made computationally tractable by the use of Gaussian processes to approximate the distance metric values when assessing the discrepancy of candidate particles. Indeed, as shown in [32], inferring the parameters for a single synthetic data set with the WMM or the CBM took between 10 and 100 core-hours, and a lower bound estimate for running ABC using Smoldyn ranged from 780 to 4635 core-hours depending on parameter values. Clearly, running this process on a large amount of synthetic data sets and in various configurations in terms of distance metrics and amount of data is not computationally tractable. In comparison, inferring the parameters of one data set with the Gaussian process approximation took less than one core-hour including the time to train the Gaussian processes. There is of course some error associated with using such an approximation, and this error should be monitored to differentiate it from other sources. In [S1 Fig](#), we looked at the utility as defined by Järvenpää *et al.* [46] and showed that it was relatively constant, suggesting our results are consistent over the parameter space. In some selected cases, it is even possible to estimate this error by running ABC with the full model on a few selected cases. Once confidence has been gained that the GPs are reasonably accurate, the pipeline can be used to scan other configurations of models, summary statistics and observed datasets.

Executing this pipeline should only be seen as a preliminary step to select the components involved in parameter inference. Once the setup has been calibrated and experimental data has been collected, regular ABC can be used, provided simulating the model is not too computationally expensive.

In this study, we used the model of highest fidelity when it comes to biophysical realism as a proxy for the ground truth. This enabled us to compare different model fidelities to each other. When attempting to make extrapolations of the pipeline results to real experimental conditions, it is important to recall that in practice, real experimental data will be noisier than simulated data, first because every model is wrong (although some of them are useful, as the quote says), and then because even if the model is an accurate representation of reality there will always be measurement noise due to limitations in the experimental protocol. Here we opted for “perfect synthetic data” to make interpretation of numerical settings and model comparison easier, though we note that it would be possible and interesting to repeat the numerical experiments using a measurement error model. In this way it would be possible to also study the different settings and scenarios (e.g. summary statistics vs. Kolmogorov distance) with respect to robustness to noise.

In an inference scenario using real experimental data the “true” parameters are unknown and it will be difficult to validate parameter inference, given that the model error is unknown. As discussed in the introduction, most often a modeler will favor one model type. As our study

revealed, it can be important to compare different models to each other and in particular compare different inference strategies.

Bayesian parameter inference is an efficient technique to find areas of the parameter space that resembles observed data. If the parameters are identifiable, the estimated parameters will correspond to the “true” parameters. Showing identifiability, however, is only possible for a limited class of models [47]. Our approach provides an alternative where it is possible to quantify potential error due to model choice, summary statistics or lack of data in the context of Bayesian inference.

Finally, we note that in this work we were interested in using the model to accurately learn the underlying model parameters. For a given set of observed data and for given numerical inference settings, it is then clear that a modeler should favor the computationally cheapest model that results in good parameter inference. Note here that, due to practical aspects of likelihood-free inference, it is not enough if a coarse-grained approximation is accurate only at the true parameter point, it needs to be accurate throughout the support of the prior distribution. Given that true parameters are unknown in practice, this means that we either need some prior knowledge about in which regimes the true parameter will fall (in which case we can use a pipeline like ours to suggest the stability of inference to different model choices), or we need to seek a globally accurate approximation. We emphasize that this problem is different from a typical ABC-based Bayesian model selection problem, in which we seek to use simulators for different models to compute the probability of the models generating the observed data. In that setting we allow models to take “wrong” physical parameter values as long as that model configuration is capable of generating trajectories close to the observed data. We plan to in future work apply our developed pipelines to investigate this aspect in more detail.

4 Methods

4.1 Stochastic models for chemical kinetics

Stochastic chemical kinetics in single cells can be modeled at various fidelities, from stochastic differential equation to detailed particle models [48]. Yet, the choice of modeling fidelity level is not always easy to do *a priori*. In particular, models including spatial details about the distribution of molecules throughout the cells can reveal new insights but come with a significant increase in computational cost [49].

A popular modeling framework is the Chemical Master Equation (CME) [50]. In the CME formalism, the system is represented by a state vector \mathbf{x} where each row represent the molecular count of a given species. The probability distribution for a system of n species and m reactions is given by the solution of the master equation:

$$\partial_t p(\mathbf{x}, t | \mathbf{x}_0, t_0) = \sum_{i=1}^m a_i(\mathbf{x} - \mathbf{v}_i) p(\mathbf{x} - \mathbf{v}_i | \mathbf{x}_0, t_0) - a_i(\mathbf{x}) p(\mathbf{x}, t | \mathbf{x}_0, t_0), \quad (6)$$

where \mathbf{x} is the state vector of the system, $a_i(\mathbf{x})$ and \mathbf{v}_i are the propensity and the stoichiometric vector of reaction i , respectively, and \mathbf{x}_0 is the state of the system at time t_0 .

Unfortunately, solving the CME numerically is in most practical cases intractable. It is however possible to generate realizations of the CME using Gillespie’s Stochastic Simulation Algorithm (SSA) [51]. One fundamental assumption of the CME is that molecules inside the cell are *well-mixed*, i.e. there is enough time between reactions for the molecules to diffuse uniformly across the cell. In other words, the CME does not include spatial details about the location of each molecules.

In [32], we presented a technique to include some degrees of spatial details into the CME framework. By dividing the cell into compartments and computing transition rates between these compartments, we are able to include some spatial information and approximate more detailed models for only a marginal increase in computational cost.

Another, more standard generalization of the CME to include spatial details is the Reaction-Diffusion Master Equation (RDME) [52, 53]. In the RDME framework, space is discretized into small voxels. Every voxel is assumed to be well mixed and reaction can only occur between molecules belonging to the same voxel. Additionally, molecules can diffuse to neighboring voxels, depending on the geometry of the discretization and of the diffusion rate.

Other, even more detailed methods track the position of each molecule in continuous space. For instance, in the Smoluchowski diffusion limited model, molecules diffuse in space following Fick's second law:

$$\frac{dp}{dt}(\mathbf{r}, t) = D\Delta p(\mathbf{r}, t) \quad (7)$$

where \mathbf{r} is the position of a molecule, D its diffusion constant and p the probability distribution of its position. Molecules are then modeled as hard spheres that react upon collision with a given probability. Solving the Smoluchowski equation in the general case is an open problem. One approach to circumvent this issue is to discretize time and rely on approximations to determine when two molecules collide and potentially trigger a chemical reaction. This approach is used in e.g. Smoldyn [13]. Another approach consists in isolating pairs of molecules or single molecules in protective domains where the Smoluchowski equations can be solved analytically using Green's Function Reaction Dynamics. This is the approach used in e.g. eGFRD [54].

There is a well defined mathematical hierarchy between these approaches [48, 55]. As a matter of fact, it is known that in the limit of infinite diffusion, spatial approaches will converge towards the CME. However, in practice, it is difficult to determine if diffusion is "fast enough" for the CME to be a valid approximation of the chemical system under scrutiny. Using a more detailed model will be more accurate, but at the price of a higher computational cost. Balancing these two aspects is a critical question in modeling.

In a previous study [32], we considered three stochastic models including various level of details and compared how they related to each other in terms of accuracy in the context of the Hes1 system presented in [7]. In Section 2 we will use the same models to illustrate how our pipeline can be used to compare these models in a Bayesian parameter inference setting.

4.2 Likelihood-free Bayesian inference

Given a mathematical model $\mathbf{y} = f(\theta)$, the goal of parameter inference is to fit f to observed data \mathbf{y}_o , i.e., estimate the parameters θ_o that give rise to simulated data $\mathbf{y}_{sim} = f(\theta_o)$ such that $\mathbf{y}_{sim} = \mathbf{y}_o$. As the model f is stochastic, in reality the equality condition is too strict and can never be fulfilled exactly. Typically the equality condition is replaced by a relaxed form involving a threshold. Also, we note that for all but very simple models of academic interest, the likelihood function corresponding to f is either unavailable or computationally impractical to approximate. Therefore, parameter estimation must proceed in a likelihood-free manner making use of the observed data, and access to the simulation model f .

The most popular family of parameter estimation methods in the likelihood-free setting is approximate Bayesian computation (ABC) [56]. ABC parameter inference begins with the specification of a *prior* distribution $p(\theta)$ over the parameters θ , representing the parameter search space. The ABC rejection sampling algorithm then samples $\theta_{sim} \sim p(\theta)$, and simulates

$y_{\text{sim}} = f(\theta_{\text{sim}})$. The simulated time series y_{sim} must now be compared to y_o to validate whether the two time series were sufficiently close, i.e. if the distance $d(y_o, y_{\text{sim}}) \leq \epsilon$, where ϵ is a user specified acceptance threshold, and d is a distance metric, often chosen to be the Euclidean distance in practice. If so, then θ_{sim} is deemed to be *accepted*, else *rejected*. This sample-simulate-compare rejection sampling cycle is repeated until enough amount of accepted samples are obtained. The set of accepted samples then form the estimated posterior distribution $p(\theta|y_o)$, solving the parameter inference problem.

The comparison between time series' y_o and y_{sim} is typically performed in terms of k low-dimensional *summary statistics* $\mathbf{S} = S_1(\mathbf{y}), \dots, S_k(\mathbf{y})$ or features of the time series (e.g., statistical moments). This is due to the *curse of dimensionality* when comparing rich high-dimensional time series (detailed discussion in Chapter 5 of [56]).

Summary statistic selection is a well-studied problem, and there exist methods to select k informative statistics out of a candidate pool of m total statistics. A thorough treatment of the topic can be found in [27, 56]. A well motivated method of selecting summary statistics is based on the notion of approximate sufficiency (AS) [57]. Summary statistics \mathbf{S} are sufficient if adding a statistic S_{new} to \mathbf{S} does not change the approximated posterior $p(\theta|y_o)$. The AS algorithm initiates tests for different statistics in random order [58], therefore in this work we will repeat summary statistic selection several times to compute the frequency of selection of each statistic. The most frequently selected statistics will be used in the parameter inference process. The candidate pool of statistics to choose from include the following statistics/features for each species.

- sum of values
- absolute energy
- mean absolute change
- mean change
- median
- mean
- length
- standard deviation
- skewness
- kurtosis
- longest strike below mean
- longest strike above mean
- last location of maximum
- first location of maximum
- last location of minimum
- first location of minimum
- maximum
- minimum

Therefore, in total there are 36 candidate statistics to choose from—18 for each of the species mRNA and protein (see model definition in Section 2.2).

4.3 Gaussian processes for likelihood-free parameter inference

ABC typically entails slow convergence towards the estimated posterior, and may require several thousands of simulations to provide a reliable estimate. When simulations are expensive, and when we need to perform many such inference computations, this can represent a serious computational bottleneck. In this study, we set to make this cost as small as possible, so that we can repeat parameter inference experiments over wide prior ranges and many inference setups.

A Gaussian process is a generalization of the Gaussian probability distribution and as such can be thought of as a distribution over functions $f(\mathbf{x})$. The distribution is specified using a mean function $\mu(\mathbf{x})$, and a (positive semidefinite) covariance function $k(\mathbf{x}, \mathbf{x}')$, where the pairing $(\mathbf{x}, \mathbf{x}')$ covers all possible data point pairs in the training set. The mean function μ is often set to a constant, while the kernel function k is used to enforce certain prior beliefs (e.g., smoothness via the squared exponential kernel function).

In [46], Järvenpää *et al.* describe how Gaussian processes can be used to approximate the distance metric values, or discrepancy, between simulated data and observed data and demonstrate how this technique can be used to efficiently infer parameters when simulations are too costly to be used directly in the inference algorithm. Here we use a similar approach and train a Gaussian process model as a surrogate to approximate and replace the distance metric values between simulated and observed data (one for each simulation model we consider). The surrogate model therefore indirectly also approximates the simulation model.

We set up the Gaussian processes with Scikit-learn [59] and set the kernel to the sum between a rational quadratic kernel and a white kernel:

$$\begin{aligned} k_1(\mathbf{x}, \mathbf{x}') &= \left(1 + \frac{d(\mathbf{x}, \mathbf{x}')^2}{2\alpha l^2}\right)^{-\alpha}, \\ k_2(\mathbf{x}, \mathbf{x}') &= \gamma_{noise} \text{ if } \mathbf{x} == \mathbf{x}' \text{ else } 0, \\ k(\mathbf{x}, \mathbf{x}') &= k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}'), \end{aligned}$$

where l is the length-scale hyperparameter controlling the correlation strength, α is the scale mixture hyperparameter while γ_{noise} signifies the amount of white noise. The training process maximizes the log marginal likelihood using the L-BFGS-B algorithm in order to optimize the hyperparameters. The training data comes from a previous study [32], and is composed of 512 samples.

The pipeline can be summarized as follows:

1. We train Gaussian processes to approximate the distance metric values between synthetic, observed data and simulated data.
2. This surrogate distance is used by pyABC to evaluate candidate particles. No extra simulations are performed during this process.

Gaussian processes not only estimate the mean value of the distance metric at a given point in parameter space, they also estimate the uncertainty of this value. In our case, this is important because of the stochastic aspect of particle acceptance in ABC. Specifically, the same particle may be accepted or rejected depending on the simulated data from this particle, especially if it comes from a stochastic model. Thus, by modeling the stochastic variations around the measured distance with Gaussian processes, we can reproduce this aspect.

Naturally, the GP surrogates come at the cost of an approximation error. Since we do not know the true posterior distribution, and because executing ABC with the full models is computationally expensive (or even intractable for the most detailed models), it is difficult to quantify what is the effect of this approximation on our results. The reader is referred to [60] for a discussion on the effect of GP approximation error within the context of the likelihood-free parameter inference problem. Järvenpää *et al.* [46] introduced a measure of utility to quantify the goodness of fit of the Gaussian processes. Although in absolute terms, this raw number is not very enlightening, it can be used to compare how the approximation performs in different configurations. In all our experiments, the utility was never correlated with the patterns exhibited on the error maps (see S1 Fig), suggesting it had only a low impact on the inference error.

In conclusion, by using this approximate distance, we can greatly reduce the computational cost of running ABC, regardless of the model used to simulate the data. In particular, even detailed models which would be far too computationally expensive to be used as is in ABC can be plugged in into our pipeline. This makes it possible to set a baseline in terms of what could be achieved in terms of accuracy when using simpler models.

Supporting information

S1 Fig. Heatmaps of computed model utilities when using the Kolmogorov distance for WMM (left), CBM (middle) and Smoldyn (right). The y-axis shows how the model utility varies with the diffusion constant and x-axis shows how it varies with the reactivity constant. All axes are displayed in log scale and so too is the colour bar. Overall, the model utilities do not correlate with the error estimates presented in Fig 4.

(PDF)

S2 Fig. Boxen plots showing a comparison of expected log-error for the WMM (top row) and the CBM (bottom row), for all three different distance metrics and when measuring only protein levels (orange), only mRNA levels (grey), or both (blue). Overall no big difference can be observed, contrary to the case where the Smoldyn was used (see Fig 7).

(PDF)

S3 Fig. Heatmaps showing expected log-error for WMM (left), CBM (middle) and Smoldyn (right) based on optimized statistics. The y-axis shows how the expected log-error varies with the diffusion constant and x-axis shows how it varies with the reactivity constant. The error is slightly lower than when only basic summary statistics are used.

(PDF)

S4 Fig. Comparison of distance measure based on expected value or based on distribution. When comparing simulated and true data via summary statistics for datasets with multiple trajectories the most straightforward and the most common way is to compute the statistic for each trajectory and then compare the expected values for true and simulated data. This is what is done in the main manuscript when using summary statistics. As a more elaborate alternative we can also compare the summary statistics using the Kolmogorov distance (Kolmogorov statistics). This entails computing the histogram CDF for the statistic at each timepoint, then taking the Kolmogorov distance between true and observed data. As can be seen, this approach is advantageous when using the Smoldyn simulator, however it does not lead to as low errors as directly comparing the copy numbers as done in the main manuscript (Kolmogorov distance). For the CBM model, however, taking a distribution measure leads to higher error.

(PDF)

S5 Fig. Boxen plots showing a comparison of expected log-error for all three models using summary statistics (top row) or the Kolmogorov distance (bottom row) with either 64 simulated trajectories (dark blue) or 4 simulated trajectories (light blue). In general, increasing the granularity of the data has a greater effect when using the Kolmogorov distance, where the accuracy improves for both the CBM and Smoldyn.

(PDF)

S6 Fig. The data from Fig 4 in the main manuscript plotted with Root Mean Squared Error (RMSE) instead of the expected log-error metric. While there are quantitative differences, the same general trend persists.

(PDF)

S1 Text. Summary statistic definitions. Explanation of each of the summary statistics used.

(PDF)

Acknowledgments

We are grateful for computational resources provided by the Swedish National Infrastructure for Computing (SNIC) through Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX) under Project 2019/8-227.

Author Contributions

Conceptualization: Marc Sturrock, Andreas Hellander.

Formal analysis: Adrien Coulier.

Funding acquisition: Andreas Hellander.

Investigation: Adrien Coulier.

Methodology: Adrien Coulier, Prashant Singh, Marc Sturrock, Andreas Hellander.

Project administration: Adrien Coulier, Andreas Hellander.

Software: Adrien Coulier.

Supervision: Marc Sturrock, Andreas Hellander.

Validation: Adrien Coulier.

Visualization: Adrien Coulier.

Writing – original draft: Adrien Coulier, Prashant Singh, Marc Sturrock, Andreas Hellander.

Writing – review & editing: Adrien Coulier, Prashant Singh, Marc Sturrock, Andreas Hellander.

References

1. Hirata H, Yoshiura S, Ohtsuka T, Bessho Y, Harada T, Yoshikawa K, et al. Oscillatory expression of the bHLH factor Hes1 regulated by a negative feedback loop. *Science*. 2002; 298(5594):840–843. <https://doi.org/10.1126/science.1074560> PMID: 12399594
2. Bernard S, Čajavec B, Pujol-Menjouet L, Mackey MC, Herzog H. Modelling transcriptional feedback loops: the role of Gro/TLE1 in Hes1 oscillations. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*. 2006; 364(1842):1155–1170. <https://doi.org/10.1098/rsta.2006.1761> PMID: 16608701
3. Monk NA. Oscillatory expression of Hes1, p53, and NF- κ B driven by transcriptional time delays. *Current Biology*. 2003; 13(16):1409–1413. [https://doi.org/10.1016/S0960-9822\(03\)00494-9](https://doi.org/10.1016/S0960-9822(03)00494-9) PMID: 12932324

4. Jensen M, Sneppen K, Tiana G. Sustained oscillations and time delays in gene expression of protein Hes1. *Febs Letters*. 2003; 541(1-3):176–177. [https://doi.org/10.1016/S0014-5793\(03\)00279-5](https://doi.org/10.1016/S0014-5793(03)00279-5) PMID: 12706840
5. Chaplain M, Ptashnyk M, Sturrock M. Hopf bifurcation in a gene regulatory network model: Molecular movement causes oscillations. *Mathematical Models and Methods in Applied Sciences*. 2015; 25(06):1179–1215. <https://doi.org/10.1142/S021820251550030X>
6. Heron EA, Finkenstädt B, Rand DA. Bayesian inference for dynamic transcriptional regulation; the Hes1 system as a case study. *Bioinformatics*. 2007; 23(19):2596–2603. <https://doi.org/10.1093/bioinformatics/btm367> PMID: 17660527
7. Sturrock M, Hellander A, Matzavinos A, Chaplain MA. Spatial stochastic modelling of the Hes1 gene regulatory network: intrinsic noise can explain heterogeneity in embryonic stem cell differentiation. *Journal of The Royal Society Interface*. 2013; 10(80):20120988. <https://doi.org/10.1098/rsif.2012.0988> PMID: 23325756
8. Eliaš J, Dimitrio L, Clairambault J, Natalini R. The dynamics of p53 in single cells: physiologically based ODE and reaction–diffusion PDE models. *Physical biology*. 2014; 11(4):045001. <https://doi.org/10.1088/1478-3975/11/4/045001> PMID: 25075792
9. Sturrock M, Terry AJ, Xirodimas DP, Thompson AM, Chaplain MA. Spatio-temporal modelling of the Hes1 and p53-Mdm2 intracellular signalling pathways. *Journal of theoretical biology*. 2011; 273(1):15–31. <https://doi.org/10.1016/j.jtbi.2010.12.016> PMID: 21184761
10. Williams RA, Timmis J, Qwarnstrom EE. Computational models of the NF-KB signalling pathway. *Computation*. 2014; 2(4):131–158. <https://doi.org/10.3390/computation2040131>
11. Enderling H, Chaplain M AJ. Mathematical modeling of tumor growth and treatment. *Current pharmaceutical design*. 2014; 20(30):4934–4940. <https://doi.org/10.2174/1381612819666131125150434> PMID: 24283955
12. Gutenkunst RN, Waterfall JJ, Casey FP, Brown KS, Myers CR, Sethna JP. Universally sloppy parameter sensitivities in systems biology models. *PLoS Comput Biol*. 2007; 3(10):e189. <https://doi.org/10.1371/journal.pcbi.0030189> PMID: 17922568
13. Andrews SS, Addy NJ, Brent R, Arkin AP. Detailed simulations of cell biology with Smoldyn 2.1. *PLoS Comput Biol*. 2010; 6(3):e1000705. <https://doi.org/10.1371/journal.pcbi.1000705> PMID: 20300644
14. Aebersold R, Mann M. Mass spectrometry-based proteomics. *Nature*. 2003; 422(6928):198–207. <https://doi.org/10.1038/nature01511> PMID: 12634793
15. Schirle M, Heurtier MA, Kuster B. Profiling core proteomes of human cell lines by one-dimensional PAGE and liquid chromatography-tandem mass spectrometry. *Molecular & Cellular Proteomics*. 2003; 2(12):1297–1305. <https://doi.org/10.1074/mcp.M300087-MCP200> PMID: 14532353
16. Budnik B, Levy E, Harmange G, Slavov N. SCoPE-MS: mass spectrometry of single mammalian cells quantifies proteome heterogeneity during cell differentiation. *Genome biology*. 2018; 19(1):1–12. <https://doi.org/10.1186/s13059-018-1547-5> PMID: 30343672
17. Bendall SC, Simonds EF, Qiu P, El-ad DA, Krutzik PO, Finck R, et al. Single-cell mass cytometry of differential immune and drug responses across a human hematopoietic continuum. *Science*. 2011; 332(6030):687–696. <https://doi.org/10.1126/science.1198704> PMID: 21551058
18. Albayrak C, Jordi CA, Zechner C, Lin J, Bichsel CA, Khammash M, et al. Digital quantification of proteins and mRNA in single mammalian cells. *Molecular cell*. 2016; 61(6):914–924. <https://doi.org/10.1016/j.molcel.2016.02.030> PMID: 26990994
19. Darmanis S, Gallant CJ, Marinescu VD, Niklasson M, Segerman A, Flamourakis G, et al. Simultaneous multiplexed measurement of RNA and proteins in single cells. *Cell reports*. 2016; 14(2):380–389. <https://doi.org/10.1016/j.celrep.2015.12.021> PMID: 26748716
20. Whale AS, Huggett JF, Cowen S, Speirs V, Shaw J, Ellison S, et al. Comparison of microfluidic digital PCR and conventional quantitative PCR for measuring copy number variation. *Nucleic acids research*. 2012; 40(11):e82–e82. <https://doi.org/10.1093/nar/gks203> PMID: 22373922
21. Lin J, Jordi C, Son M, Van Phan H, Drayman N, Abasiyanik MF, et al. Ultra-sensitive digital quantification of proteins and mRNA in single cells. *Nature communications*. 2019; 10(1):1–10. <https://doi.org/10.1038/s41467-019-11531-z> PMID: 31391463
22. Kursawe J, Baker RE, Fletcher AG. Approximate Bayesian computation reveals the importance of repeated measurements for parameterising cell-based models of growing tissues. *Journal of theoretical biology*. 2018; 443:66–81. <https://doi.org/10.1016/j.jtbi.2018.01.020> PMID: 29391171
23. Harrison JU, Baker RE. The impact of temporal sampling resolution on parameter inference for biological transport models. *PLoS computational biology*. 2018; 14(6):e1006235. <https://doi.org/10.1371/journal.pcbi.1006235> PMID: 29939995

24. Robert CP, Cornuet JM, Marin JM, Pillai NS. Lack of confidence in approximate Bayesian computation model choice. *Proceedings of the National Academy of Sciences*. 2011; 108(37):15112–15117. <https://doi.org/10.1073/pnas.1102900108> PMID: 21876135
25. Sunnåker M, Busetto AG, Numminen E, Corander J, Foll M, Dessimoz C. Approximate bayesian computation. *PLoS Comput Biol*. 2013; 9(1):e1002803. <https://doi.org/10.1371/journal.pcbi.1002803> PMID: 23341757
26. Maclaren OJ, Nicholson R. What can be estimated? Identifiability, estimability, causal inference and ill-posed inverse problems. *arXiv preprint arXiv:190402826*. 2019;.
27. Prangle D. Summary statistics in approximate Bayesian computation. *arXiv preprint arXiv:151205633*. 2015;.
28. Browning AP, Warne DJ, Burrage K, Baker RE, Simpson MJ. Identifiability analysis for stochastic differential equation models in systems biology. *Journal of the Royal Society Interface*. 2020; 17(173):20200652. <https://doi.org/10.1098/rsif.2020.0652> PMID: 33323054
29. Macklin P. When seeing isn't believing: How math can guide our interpretation of measurements and experiments. *Cell Systems*. 2017; 5(2):92–94. <https://doi.org/10.1016/j.cels.2017.08.005> PMID: 28837815
30. Warne DJ, Baker RE, Simpson MJ. Using experimental data and information criteria to guide model selection for reaction–diffusion problems in mathematical biology. *Bulletin of Mathematical Biology*. 2019; 81(6):1760–1804. <https://doi.org/10.1007/s11538-019-00589-x> PMID: 30815837
31. Fox ZR, Munsky B. The finite state projection based Fisher information matrix approach to estimate information and optimize single-cell experiments. *PLoS computational biology*. 2019; 15(1):e1006365. <https://doi.org/10.1371/journal.pcbi.1006365> PMID: 30645589
32. Coulier A, Hellander S, Hellander A. A multiscale compartment-based model of stochastic gene regulatory networks using hitting-time analysis. *The Journal of Chemical Physics*. 2021; 154(18):184105. <https://doi.org/10.1063/5.0010764> PMID: 34241042
33. Jiang R, Jacob B, Geiger M, Matthew S, Rumsey B, Singh P, et al. Epidemiological modeling in StochSS Live! *Bioinformatics*. 2021;.
34. Drawert B, Hellander A, Bales B, Banerjee D, Bellesia G, Daigle BJ Jr, et al. Stochastic simulation service: bridging the gap between the computational expert and the biologist. *PLoS computational biology*. 2016; 12(12):e1005220. <https://doi.org/10.1371/journal.pcbi.1005220> PMID: 27930676
35. Klinger E, Rickert D, Hasenauer J. pyABC: distributed, likelihood-free inference. *Bioinformatics*. 2018; 34(20):3591–3593. <https://doi.org/10.1093/bioinformatics/bty361> PMID: 29762723
36. Lillacci G, Khammash M. The signal within the noise: efficient inference of stochastic gene regulation models using fluorescence histograms and stochastic simulations. *Bioinformatics*. 2013; 29(18):2311–2319. <https://doi.org/10.1093/bioinformatics/btt380> PMID: 23821649
37. Hofmann H, Kafadar K, Wickham H. Letter-value plots: Boxplots for large data. *had.co.nz*; 2011.
38. Julius M, Masuda T, Herzenberg L. Demonstration that antigen-binding cells are precursors of antibody-producing cells after purification with a fluorescence-activated cell sorter. *Proceedings of the National Academy of Sciences*. 1972; 69(7):1934–1938. <https://doi.org/10.1073/pnas.69.7.1934> PMID: 4114858
39. Vo HD, Fox Z, Baetica A, Munsky B. Bayesian estimation for stochastic gene expression using multifidelity models. *The Journal of Physical Chemistry B*. 2019; 123(10):2217–2234. <https://doi.org/10.1021/acs.jpcc.8b10946> PMID: 30777763
40. Matsumoto H, Kiryu H, Furusawa C, Ko MS, Ko SB, Gouda N, et al. SCODE: an efficient regulatory network inference algorithm from single-cell RNA-Seq during differentiation. *Bioinformatics*. 2017; 33(15):2314–2321. <https://doi.org/10.1093/bioinformatics/btx194> PMID: 28379368
41. Schwanhäusser B, Busse D, Li N, Dittmar G, Schuchhardt J, Wolf J, et al. Global quantification of mammalian gene expression control. *Nature*. 2011; 473(7347):337–342. <https://doi.org/10.1038/nature10098> PMID: 21593866
42. Fearnhead P, Prangle D. Constructing summary statistics for approximate Bayesian computation: semi-automatic approximate Bayesian computation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*. 2012; 74(3):419–474. <https://doi.org/10.1111/j.1467-9868.2011.01010.x>
43. Åkesson M, Singh P, Wrede F, Hellander A. Convolutional Neural Networks as Summary Statistics for Approximate Bayesian Computation. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*. 2021;.
44. Jiang B, Wu TY, Zheng C, Wong WH. Learning Summary Statistic for Approximate Bayesian Computation via Deep Neural Network. *Statistica Sinica*. 2017; 27(4):1595–1618.

45. Wiqvist S, Mattei PA, Picchini U, Frellsen J. Partially Exchangeable Networks and Architectures for Learning Summary Statistics in Approximate Bayesian Computation. In: International Conference on Machine Learning; 2019. p. 6798–6807.
46. Järvenpää M, Gutmann MU, Vehtari A, Marttinen P, et al. Gaussian process modelling in approximate Bayesian computation to estimate horizontal gene transfer in bacteria. *Annals of Applied Statistics*. 2018; 12(4):2228–2251.
47. Simpson MJ, Baker RE, Vittadello ST, Maclaren OJ. Practical parameter identifiability for spatio-temporal models of cell invasion. *Journal of the Royal Society Interface*. 2020; 17(164):20200055. <https://doi.org/10.1098/rsif.2020.0055> PMID: 32126193
48. Gillespie DT, Hellander A, Petzold LR. Perspective: Stochastic algorithms for chemical kinetics. *The Journal of chemical physics*. 2013; 138(17):05B201_1. <https://doi.org/10.1063/1.4801941> PMID: 23656106
49. Burrage K, Burrage PM, Leier A, Marquez-Lago T, Nicolau DV. Stochastic simulation for spatial modeling of dynamic processes in a living cell. In: *Design and Analysis of Biomolecular Circuits*. Springer; 2011. p. 43–62.
50. Gillespie DT. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of computational physics*. 1976; 22(4):403–434. [https://doi.org/10.1016/0021-9991\(76\)90041-3](https://doi.org/10.1016/0021-9991(76)90041-3)
51. Gillespie DT. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry*. 1977; 81(25):2340–2361. <https://doi.org/10.1021/j100540a008>
52. Elf J, Doncic A, Ehrenberg M. Mesoscopic reaction-diffusion in intracellular signaling. In: *Fluctuations and noise in biological, biophysical, and biomedical systems*. vol. 5110. International Society for Optics and Photonics; 2003. p. 114–124.
53. Stundzia AB, Lumsden CJ. Stochastic simulation of coupled reaction–diffusion processes. *Journal of computational physics*. 1996; 127(1):196–207. <https://doi.org/10.1006/jcph.1996.0168>
54. Sokolowski TR, Paijmans J, Bossen L, Miedema T, Wehrens M, Becker NB, et al. eGFRD in all dimensions. *The Journal of chemical physics*. 2019; 150(5):054108. <https://doi.org/10.1063/1.5064867> PMID: 30736681
55. Smith S, Grima R. Spatial stochastic intracellular kinetics: A review of modelling approaches. *Bulletin of mathematical biology*. 2019; 81(8):2960–3009. <https://doi.org/10.1007/s11538-018-0443-1> PMID: 29785521
56. Sisson SA, Fan Y, Beaumont M. *Handbook of approximate Bayesian computation*. CRC Press; 2018.
57. Joyce P, Marjoram P. Approximately sufficient statistics and Bayesian computation. *Statistical applications in genetics and molecular biology*. 2008; 7(1).
58. Nunes MA, Balding DJ. On optimal selection of summary statistics for approximate Bayesian computation. *Statistical Applications in Genetics & Molecular Biology*. 2010; 9(1). PMID: 20887273
59. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in Python. *the Journal of machine Learning research*. 2011; 12:2825–2830.
60. Jarvenpää M. Gaussian Process Surrogate Methods for Sample-Efficient Approximate Bayesian Computation. Aalto University publication series, DISSERTATIONS 121/2020. 2020;.