

A Monte Carlo method to estimate cell population heterogeneity

Ben Lambert^{1,2*}, David J. Gavaghan³, Simon Tavener⁴.

1 Department of Zoology, University of Oxford, Oxford, Oxfordshire, U.K.

2 MRC Centre for Global Infectious Disease Analysis, School of Public Health, Imperial College London, London W2 1PG, UK.

3 Department of Computer Science, University of Oxford, Oxford, U.K.

4 Department of Statistics, Colorado State University, Fort Collins, Colorado, U.S.A.

*ben.c.lambert@gmail.com.

Revision date & time: 2019-05-21 18:37

1 Abstract

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Variation is characteristic of all living systems. Laboratory techniques such as flow cytometry can probe individual cells and, after decades of experimentation, it is clear that even members of seemingly homogeneous cell populations can exhibit differences. To understand whether this variation is biologically meaningful, it is essential to discern its source. Mathematical models of biological systems are tools that can be used to investigate causes of cell-to-cell variation. From mathematical analysis and simulation of these models, biological hypotheses can be posed and investigated, then parameter inference can determine which of these is most compatible with experimental data. Data from laboratory experiments often takes the form of “snapshots” representing distributions of cellular properties at different points in times, rather than individual cell trajectories. This data is not straightforward to fit using hierarchical Bayesian methods since it requires inferring the identities of the groups to which individual cells belong. Here, we introduce a computational sampling method we call “Contour Monte Carlo” for estimating mathematical model parameters from snapshot distributions which is straightforward to implement and does not require explicitly assigning cells to categories. Our method is most applicable to systems where the dominant source of uncertainty is heterogeneity in cellular processes rather than experimental measurement error which, due to the increasingly finescale resolution of laboratory techniques, may be the case for a wide class of systems. In this paper, we illustrate the use of our method by quantifying cellular variation for three biological systems of interest and provide code in the form of a Julia notebook which allows others to apply this approach to their problem.

2 Introduction

27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50

Variation rather than homogeneity is the rule rather than exception in biology. Indeed, without variation, biology as a discipline would not exist, since as evolutionary biologist JBS Haldane wrote, variation is the “raw material” of evolution. The Red Queen Hypothesis asserts that organisms must continually evolve in order to survive when pitted against other - also evolving - organisms [1]. A corollary of this hypothesis is that multicellular organisms may evolve cellular phenotypic heterogeneity to allow faster adaptation to changing environments, which may explain the observed variation in a range of biological systems [2]. Whilst cell population variation can confer evolutionary advantages, it can also be costly in other circumstances. In biotechnological processes, heterogeneity in cellular function can lead to reduced yields of biochemical products [3]. In human biology, variation across cells can enable pathologies to develop and also prevents effective medical treatment, since medical interventions typically aim to steer modal cellular properties and hence fail to influence key subpopulations. For example, cellular heterogeneity likely contributes to the persistence of some cancerous tumours [4] and may also allow them to evolve resistance to chemotherapies over time [5]. Identifying and quantifying sources of variation in populations of cells is important for a wide range of applications because it allows us to determine whether this variability is benign or alternatively requires remedy.

Mathematical models are essential tools for understanding cellular systems, whose emergent properties are the result of complex interactions

between various actors. Perhaps the simplest flavour of mathematical model used in biological systems are ordinary differential equations (ODEs) that lump individual actors into partitions according to structure or function, and seek to model the mean behaviour of each partition. Data from population-averaged experimental assays can be a powerful resource to understand whether such models faithfully reproduce system behaviours and can allow quantification of the interactions of various cellular components of complex metabolic, signalling and transcriptional networks. The worth of such models however is determined by whether averages mask differences in behaviour of individual cells that result in functional consequences [6]. In some cases, differences in cellular protein abundances due to biochemical “noise” may not be meaningful biologically [7] and so mean cell behaviour suffices as a description of the system, whereas in others there are functional consequences. For example, a recent study demonstrated that subpopulations of clonally derived hematopoietic progenitor cells with low or high expression of a particular stem cell marker produced different blood lineages [8].

To accommodate cell population heterogeneity in mathematical models, a variety of modelling choices are available, each posing different challenges for parameter inference, and are described in a recent review [9]. These include modelling biochemical processes stochastically, with properties of ensembles of cells represented by probability distributions evolving according to chemical master equations (see [10] for a tutorial on stochastic reaction-diffusion processes; RDEs). Alternatively, population balance equations (PBEs) can be used to dictate the evolution of the “number density” of differing cell types, whose properties are represented as points in \mathbb{R}^n which, in turn, affect their function, including their rate of death and cell division (see [11] for an introduction to PBEs). In a PBE approach, variation in measured quantities results primarily due to differing functional properties of heterogeneous cell types and variable initial densities of each type.

The approach we follow here is similar to that of [12], wherein dynamic cellular variation is generated by describing the evolution of each cell’s state using an ODE, but with individual cell differences in the rate parameters of the process. To our knowledge, this flavour of model is unnamed and so, for sake of reference, we term them “heterogenous ODE” models (HODEs). In HODEs, the aim of inference is to estimate the distributions of parameter values across cells consistent with observed distributions of measurements at various timepoints. A benefit of using HODEs to model cell heterogeneity is that these models are computationally straightforward to simulate and, arguably, simpler to parameterise than PBEs. In these models the predominant source of variation is due to differences in biological processes across cells not inherent stochasticity in biochemical reactions within cells, as in stochastic RDEs.

The difficulty of parameter inference for HODEs is partly due to experimental hurdles in generating data of sufficient quality to allow identification. Unlike models which represent a population by a single scalar ODE, since HODEs are individual-based they ideally require individual cell data for estimation. A widely-used method for generating data for individual cells is flow cytometry, where a large number of cells are streamed individually through a laser beam and, for example, abundance measurements are made of proteins labelled with fluorescent markers [13]. Alternatively, experimental techniques such as Western blotting and cytometric fluorescence microscopy can generate single cell measurements [14, 15]. A property of

these experimental methods is that they are destructive, meaning that individual cells are sacrificed as part of the measurement process. This means that the measurements of cell properties obtained at a certain point in time represent what are termed “snapshots” of the underlying population [15]. These snapshots are often described by histograms [12] or density functions [9] fit to the underlying data at different points in time. Since HODEs represent the underlying state of individual cells as evolving continuously through time, corresponding data showing individual cell trajectories constitutes a richer data resource. The demands of obtaining this data are higher however and typically involve either tracking individual cells through imaging methods [16] or trapping cells in a spatial position where their individual dynamics can be readily monitored [17]. These techniques impose restrictions on experimental practices meaning they cannot be realised in all circumstances, including for online monitoring of biotechnological processes or analysis of *in vivo* studies. For this reason, snapshot data continues to play an important role for determining cell level variability in many applications.

A variety of approaches have been proposed to estimate cell-to-cell variability by fitting HODE models to snapshot data. In HODEs, parameter values vary across cells according to a to-be-determined probability distribution meaning that in order to solve the exact inverse problem, the underlying ODE system needs to be simulated for each individual. Since the numbers of cells in these experiments are typically $>\sim 10^4$ [15], this usually precludes exact inference due to its computational burden and instead the raw snapshot data is approximated by probability densities [12, 15, 18, 19]. Hasenauer et al. (2011) presents a Bayesian approach to inference for HODEs, which models the input parameter space using mixtures of ansatz densities. The authors then use their method to reproduce population substructure on synthetic data generated from a model of tumour necrosis factor stimulus. Hasenauer et al. (2014) uses mixture models to model subpopulation structure in snapshot data with multiple-start local optimisation employed to maximise the non-convex likelihood, which they then apply to synthetic and real data from signalling pathway models. Loos et al. (2018) also uses mixture models to represent subpopulation structure and a maximum likelihood approach that allows for estimation of within- and between-subpopulation variability which permits fitting to multivariate output distributions with complex correlation structures. Dixit et al. (2018) discretises cell abundances into bins, then uses a maximum entropy approach as part of a Bayesian framework to fit the distribution representing cell-to-cell variability.

The framework we present here is Bayesian although is distinct from the traditional Bayesian inferential paradigm used to fit many dynamic models since the source of stochasticity arises solely due to cell-to-cell parameter variation not measurement noise. Our approach is hence most suitable when measurement error is a minor contributor to observed experimental variability. Our computational method is a two-step Monte Carlo approach which, for reasons described in §3, we term “Contour Monte Carlo” (CMC). Unlike many of the existing methods however CMC is relatively computationally straightforward to implement and does not require extensive computation time. CMC uses MCMC in its second step to sample from the posterior distribution over parameter values and hence does not require specification of ansatz densities. It also does not require *a priori* representation of subpopulation structure using mixture components, rather,

subpopulations emerge as modes in the posterior parameter distributions.
 Like [19] CMC can fit multivariate snapshot data and unlike [12], does
 not require this data to be discretised into bins. As more experimental
 techniques are developed which elucidate single cell behaviour, there is
 likely to be more interest in methods which can be used to recapitulate
 the observed snapshots. We argue that due to its simplicity and generality,
 CMC is a useful addition to the modeller’s toolkit and can be used to
 perform inference on the proliferation of rich single cell data.

Outline of the paper: In §3, we present the details of our methodological
 framework and detail the CMC algorithm used to generate samples from
 the posterior parameter distribution. In §4, we use CMC to estimate cell
 population heterogeneity in three systems of biological interest.

3 Method

In this section, we first describe the probabilistic framework
 that underlies the CMC algorithm, before introducing CMC in pseudocode
 (Algorithm 1). We also detail the workflow we have found useful in applying
 this approach to analyse cell snapshot data and suggest practical remedies
 to issues we have encountered in using CMC (Figure 4). A glossary of all
 the variables used in this paper is included as Table 1.

Experimental methods such as flow cytometry can measure single cell
 characteristics at a given point in time. Cells are typically destroyed by
 the measurement process and so rather than providing time series for each
 individual cell, the data consists of cross-sections or “snapshots” of sampled
 individuals from the population (Figure 1).

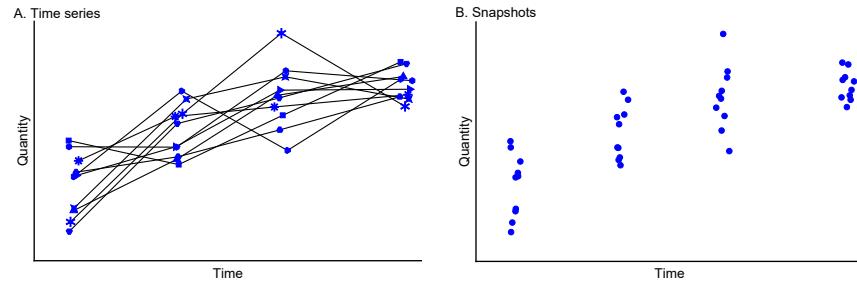


Figure 1: Time series data (A) versus snapshot data (B) typical of single cell experiments. In A that the cell identities are retained at each measurement point (indicated by given plot markers) whereas in the snapshot data in B, either this information is lost or, more often, cells are destroyed by the measurement process and so each observation corresponds to a distinct cell.

We model the processes of an individual cell using a system of ordinary differential equations (ODEs), where each element of the system describes the governing dynamics of a particular quantity of interest (for example, protein levels, RNA concentrations and so on),

$$\dot{\mathbf{x}}(t) = f(\mathbf{x}(t); \boldsymbol{\theta}). \quad (1)$$

Here $\mathbf{x}(t) = (x_1(t), x_2(t), \dots, x_k(t))$ is a vector of states for each compartment in the system at time t and $f(\cdot)$ is a function of these states and parameters

$\theta \in \mathbb{R}^p$. Note that in most circumstances, the initial state of the system, $\mathbf{x}(0)$, is unknown and it is convenient to include these as elements of θ to be estimated. The solution of eq. (1) is given by $\mathbf{x}(t) = g(t; \theta)$, where $\mathbf{x}(t) \in \mathbb{R}^k$ is a vector of outputs at time t and $g(\cdot)$ is a function that typically won't be analytically-determined; instead approximated via a numerical integration scheme.

In this paper, we assume variation characterised by snapshot data arises due to between-cell heterogeneity in the underlying parameters θ . Therefore, the evolution of the underlying state of cell i is described by an idiosyncratic ODE,

$$\dot{\mathbf{x}}^i(t) = f(\mathbf{x}^i(t); \theta^i), \quad (2)$$

with solution $\mathbf{x}^i(t) = g(t; \theta^i)$. The traditional (non-hierarchical) state-space approach to modelling dynamic systems supposes that measurement randomness generates output variation (Figure 2A). Our approach, by contrast, relies on the assumption that stochasticity in outputs is solely the result of variability in parameter values between cells (Figure 2B). Whether the assumption of “perfect” measurements is reasonable depends on the experimental details of the system under investigation but we argue that our method nevertheless provides a useful approximation in many cases where the signal to noise ratio is high.

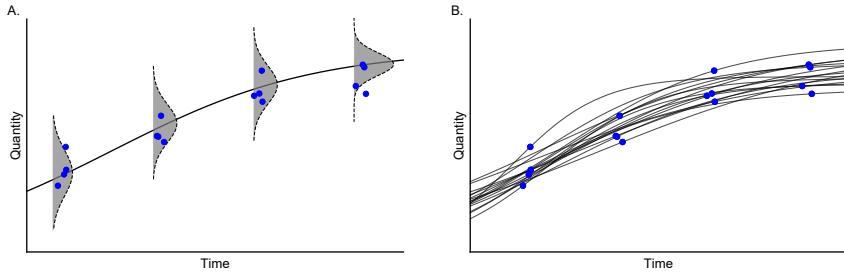


Figure 2: Two ways to generate randomness in measured outputs: the state-space model (A) versus the parameter heterogeneity model (B). For non-hierarchical state-space models (A), there is assumed to be a single “true” latent state where observations result from a noisy measurement process (grey histograms). For models with parameter heterogeneity (B), the uncertainty is generated by differences in cellular processes (black lines) between cells. Note that in both cases, individual cells are measured only once in their lifetime.

Our generative model used to produce output observations for a single cell consists of two stages: (i) sample $\theta^i \sim p(\theta)$, where $p(\theta)$ is a probability distribution characterising heterogeneity in cellular processes and, (ii) calculate output values using $\mathbf{x}_j^i(t') = g(t'; \theta)$ where, for each cell, a measurement of a subset j of output states $\mathbf{x}_j^i(t') \in \mathbf{x}^i$ is made at a single point in time $t' \in (t_1, t_2, \dots, t_o)$. The experimental observations for a single timepoint then consists of the collection of individual cell measurements $\mathbf{X}_j(t') = (\mathbf{x}_j^1(t'), \mathbf{x}_j^2(t'), \dots, \mathbf{x}_j^n(t'))$, where n is the number of cells measured at time t' . The entire observation dataset is the union of such sets across all measured timepoints $\mathbf{X}(\mathbf{t}) = (\mathbf{X}_{j_1}(t_1), \mathbf{X}_{j_2}(t_2), \dots, \mathbf{X}_{j_o}(t_o))$, where $\mathbf{t} = (t_1, t_2, \dots, t_o)$ is the vector of observation times and the subscript a on each j_a allows for measurement of different elements of the system at distinct timepoints.

Variable	Definition	Dimension
$\boldsymbol{x}(t)$	state of modelled system at time t	\mathbb{R}^k
$x_j(t) \in \boldsymbol{x}(t)$	individual state j of modelled system at time t	\mathbb{R}
$\boldsymbol{\theta}$	parameters of ODE system	\mathbb{R}^p
$f(\boldsymbol{x}(t); \boldsymbol{\theta})$	function specifying RHS of ODE system	\mathbb{R}^k
$g(t; \boldsymbol{\theta})$	solution of ODE at time t	\mathbb{R}^k
$g_j(t; \boldsymbol{\theta})$	solution of ODE for state j at time t	\mathbb{R}
$\boldsymbol{x}^i(t)$	state of modelled system of cell i at time t	\mathbb{R}^k
$\boldsymbol{x}_j^i(t) \in \boldsymbol{x}^i(t)$	state of subset j of modelled system of cell i at time t	$\mathbb{R}^{[j] \times 1}$
$\mathbf{X}_j(t) = (\boldsymbol{x}_j^1(t), \dots, \boldsymbol{x}_j^n(t))$	collection of n individual cell measurements at time t	$\mathbb{R}^{[j] \times n}$
$\mathbf{t} = (t_1, t_2, \dots, t_o)$	unique observation times	\mathbb{R}^o
$\mathbf{X}(\mathbf{t}) = (\mathbf{X}_{j_1}(t_1), \dots, \mathbf{X}_{j_o}(t_o))$	all observations collected at times \mathbf{t}	$\dim(\mathbf{X}_{j_1}) \times \dots \times \dim(\mathbf{X}_{j_o})$
Φ	parameters characterising output target distribution $p(\boldsymbol{x} \Phi)$	-
\hat{a}	estimates of any quantity a	-
$\tilde{\mathbf{t}} = (t_1, \dots, t_m)$	times when each observable is recorded	\mathbb{R}^m
$\tilde{\boldsymbol{x}}(\tilde{\mathbf{t}}) = (x_{j_1}(t_1), \dots, x_{j_m}(t_m))$	system observables	\mathbb{R}^m
$\mathcal{V}(\tilde{\boldsymbol{x}})$	the “volume” of parameter space mapping to an output of value $\tilde{\boldsymbol{x}}$	\mathbb{R}^+
$\mathbf{g}(\boldsymbol{\theta}) = \mathbf{g}(\tilde{\mathbf{t}}; \boldsymbol{\theta}) = (g_{j_1}(t_1; \boldsymbol{\theta}), \dots, g_{j_m}(t_m; \boldsymbol{\theta}))$	solution of ODE for each observable at respective times \mathbf{t}	\mathbb{R}^m
V	total volume of parameter space with uniform priors used for all parameters	\mathbb{R}^+
$\Omega(\tilde{\boldsymbol{x}}) = \{\boldsymbol{\theta} : \mathbf{g}(\boldsymbol{\theta}) = \tilde{\boldsymbol{x}}\}$	region of parameter space mapping to output $\tilde{\boldsymbol{x}}$	\mathbb{R}^p
Ψ	parameters characterising output prior distribution $p(\tilde{\boldsymbol{x}} \Psi)$	-
Ξ	parameters characterising parameter prior distribution $p(\boldsymbol{\theta} \Xi)$	-

Table 1: **Glossary of variable names used in this paper.** The dimensions of Φ , Ψ and Ξ are listed as “-” since they depend on the form of the density used to represent the process and can be anywhere from \mathbb{R}^1 to \mathbb{R}^∞ . The variables are listed in the approximate order in which they appear in the text.

Raw snapshot data consists of measurements of individual cells with exact inference requiring simulating the underlying ODE system for each individual. This is cumbersome and impractical for the numbers of cells sampled in typical experimental setups and so, instead, we follow previous work and instead represent snapshot data using probability distributions [12, 15, 18, 19]. The snapshots themselves can either be distributions of a single species or multiple species, which can be approximated by univariate and multivariate probability distributions respectively. These probability distributions are characterised by parameter estimates $\hat{\Phi}$ determined by the output observations $\mathbf{X}(\mathbf{t})$. The dimensionality of these probability distributions depends on the set of m distinct observables $\tilde{\boldsymbol{x}}(\tilde{\mathbf{t}}) = (x_{j_1}(t_1), x_{j_2}(t_2), \dots, x_{j_m}(t_m))$ recorded by experimental measurements. Note that, $\tilde{\boldsymbol{x}}(\tilde{\mathbf{t}})$ corresponds to a particular set of measurements from a hypothetical cell and is distinct from $\mathbf{X}(\mathbf{t})$, which represents the full set of experimental outputs. The vector $\tilde{\boldsymbol{x}}(\tilde{\mathbf{t}})$ is hypothetical because in reality each cell is measured at a single timepoint (although we suppose measurements of different cellular attributes are possible contemporaneously).

The goal of our inference process is to characterise the probability distribution $p(\boldsymbol{\theta}|\mathbf{X}(\mathbf{t}))$ representing heterogeneity in cellular processes. The first step in our inference workflow is to fit the output distributions using probability distributions (Figure 4(i)). We assume that the volume of observational data means the estimated probability distributions are approximate sufficient statistics of the outputs, meaning $p(\boldsymbol{\theta}|\hat{\Phi}) \approx p(\boldsymbol{\theta}|\mathbf{X}(\mathbf{t}))$.

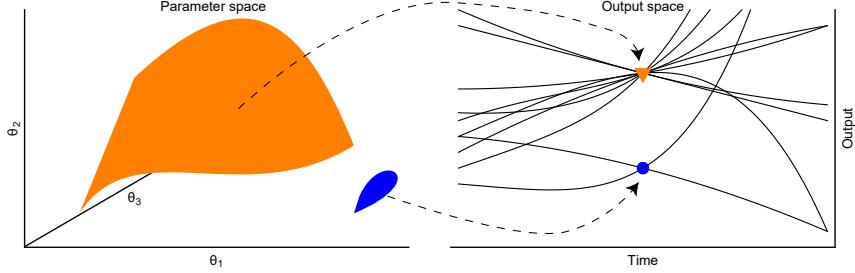


Figure 3: The non-linear mapping from parameter values (left panel) to outputs (black lines; right panel) means different sized regions of parameter space (orange and blue surfaces; left panel) correspond to distinct output values (orange triangle and blue square; right panel). In the right panel, each black line represents a distinct model simulation $g(t; \theta_1, \theta_2, \theta_3)$ and the triangle and square indicate outputs at a given point in time.

The models we seek to fit to snapshot data mostly cannot be identified from the observations. This is often because the number of model parameters exceeds the dimensionality of the output distribution (that is, $m > k$) meaning there typically exist non-singular sets of parameter values mapping to a single set of output values. That is, each vector of observed outputs $\tilde{\mathbf{x}} \in \mathbb{R}^m$, can often be caused by many combinations of parameters although, due to the non-linearity of the map from parameters to outputs, the “volume” of these regions of parameter space, $\mathcal{V}(\tilde{\mathbf{x}})$, is a function of output (Figure 3). In what follows, we make clear the distinction between observables $\tilde{\mathbf{x}}(\tilde{\mathbf{t}})$ and the vector-valued function representing modelled outputs $\mathbf{g}(\tilde{\mathbf{t}}; \boldsymbol{\theta}) = (g_{j_1}(t_1; \boldsymbol{\theta}), g_{j_2}(t_2; \boldsymbol{\theta}), \dots, g_{j_m}(t_m; \boldsymbol{\theta})) \in \mathbb{R}^m$ since the latter is a function whereas that latter is a numeric value; we also drop the $\tilde{\mathbf{t}}$ notation from following expressions to minimise clutter.

A consequence of this non-linear parameter to output geometry is that any target output distribution $p(\tilde{\mathbf{x}}|\hat{\Phi})$ does not correspond to a unique parameter distribution. For example, suppose $g(\theta_1, \theta_2) = \theta_1 + \theta_2$: the target distribution $\tilde{\mathbf{x}} \sim \mathcal{N}(0, 1)$ can be generated by any member of the set of parameter distributions $\sqrt{\eta}\theta_1 + \sqrt{1 - \eta}\theta_2$, where $\eta \in [0, 1]$ and $\theta_1, \theta_2 \sim \mathcal{N}(0, 1)$. This means that in order to ensure uniqueness of the “posterior” parameter distributions, we are required to specify “prior” distributions for the parameters, as in more traditional Bayesian inference. An additional consequence of the degeneracy of the mapping from parameters to outputs is that any sampling algorithm aimed at exploring posterior parameter space must account for the differential volumes of iso-output contours. Whilst we refer the interested reader to our companion paper on this subject [citation for tutorial paper published in Open Science], we provide a quick derivation of the posterior parameter distribution which accounts for the non-linear mapping.

To derive the posterior distribution of parameter values $p(\boldsymbol{\theta}|\hat{\Phi})$, we consider the joint density of parameters and outputs $p(\boldsymbol{\theta}, \tilde{\mathbf{x}}|\hat{\Phi})$. This can be decomposed in two ways,

$$p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \hat{\Phi}) \times p(\tilde{\mathbf{x}}|\hat{\Phi}) = p(\boldsymbol{\theta}, \tilde{\mathbf{x}}|\hat{\Phi}) = p(\tilde{\mathbf{x}}|\boldsymbol{\theta}, \hat{\Phi}) \times p(\boldsymbol{\theta}|\hat{\Phi}). \quad (3)$$

The left and right hand sides of eq. (3) can be equated and rearranged to

obtain the posterior parameter distribution,

274

$$p(\boldsymbol{\theta}|\hat{\Phi}) = \frac{p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \hat{\Phi}) \times p(\tilde{\mathbf{x}}|\hat{\Phi})}{p(\tilde{\mathbf{x}}|\boldsymbol{\theta}, \hat{\Phi})}. \quad (4)$$

Given parameters $\boldsymbol{\theta}$, the mapping from parameters to outputs is deterministic meaning $p(\tilde{\mathbf{x}}|\boldsymbol{\theta}, \hat{\Phi}) = \delta(\mathbf{g}(\boldsymbol{\theta}))$ is the Dirac delta function centred at $\tilde{\mathbf{x}} = \mathbf{g}(\boldsymbol{\theta})$. In what follows, we assume that the conditional distribution $p(\boldsymbol{\theta}|\tilde{\mathbf{x}}, \hat{\Phi})$ is independent of the data, meaning it represents a conditional “prior”, which can be manipulated by Bayes’ rule,

275
276
277
278
279

$$p(\boldsymbol{\theta}|\mathbf{g}(\boldsymbol{\theta})) = \frac{p(\boldsymbol{\theta})}{p(\mathbf{g}(\boldsymbol{\theta}))}, \quad (5)$$

where we have used the Dirac delta function for $p(\tilde{\mathbf{x}}|\boldsymbol{\theta})$. This results in the form of the posterior parameter distribution targeted by our sampling algorithm,

280
281
282

$$p(\boldsymbol{\theta}|\hat{\Phi}) = \frac{p(\boldsymbol{\theta})}{p(\mathbf{g}(\boldsymbol{\theta}))} p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Phi}). \quad (6)$$

Again, we refer to our companion piece [citation] for detailed explanation of eqs. (5) & (6) and instead here provide brief interpretation when considering a uniform prior on parameter space. In this case, $p(\boldsymbol{\theta}) = \frac{1}{V}$, where V is the total volume of parameter space. The denominator term of eq. (5) is the prior induced on output space by the prior over parameter space. For a uniform prior on parameter values, this is just proportion of parameter space where $\mathbf{g}(\boldsymbol{\theta}) = \tilde{\mathbf{x}}$, meaning,

283
284
285
286
287
288
289

$$p(\boldsymbol{\theta}|\mathbf{g}(\boldsymbol{\theta})) = \frac{1}{\mathcal{V}(\mathbf{g}(\boldsymbol{\theta}))}, \quad (7)$$

where $\mathcal{V}(\mathbf{g}(\boldsymbol{\theta}))$ is the volume of parameter space occupied by the iso-output region $\Omega(\tilde{\mathbf{x}}) = \{\boldsymbol{\theta} : \mathbf{g}(\boldsymbol{\theta}) = \tilde{\mathbf{x}}\}$. Therefore a uniform prior over parameter space implies a prior structure where all parameter values resulting in the same output $\tilde{\mathbf{x}}$ are given equal weighting.

290
291
292
293

The denominator term of eq. (5) cannot be calculated apart from for some toy examples, meaning that exact sampling from the posterior parameter distribution of eq. (6) is not, in general, possible. We propose instead a computationally efficient sampling method to estimate $p(\mathbf{g}(\boldsymbol{\theta}))$, which forms the first step of our so-called “Contour Monte Carlo” (CMC) algorithm (Algorithm 1; Figure 4(ii)), where we estimate the volume of iso-output contours with output value $\mathbf{g}(\boldsymbol{\theta})$. This step involves repeated independent sampling from the prior distribution of parameters $\boldsymbol{\theta}^i \sim p(\boldsymbol{\theta}|\Xi)$, where, for completeness, we have conditioned on Ξ parameterising our probability density. Each parameter sample is then converted into an output value $\tilde{\mathbf{x}}^i = \mathbf{g}(\boldsymbol{\theta}^i)$. The collection of output samples is then fitted using a vine copula kernel density estimator (KDE) [20], $(\tilde{\mathbf{x}}^1, \tilde{\mathbf{x}}^2, \dots, \tilde{\mathbf{x}}^{N_1}) \sim p(\tilde{\mathbf{x}}|\hat{\Psi})$. Throughout the course of development of CMC, we have tested many forms of KDE and have found vine copula KDE is best suited to approximating the higher dimensional probability distributions required in practice.

294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312

The second step in our algorithm then uses Markov chain Monte Carlo (MCMC) to sample from an approximate version of eq. (6) with the estimated density $p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Psi})$ replacing its corresponding estimand (Algorithm 1; Figure 4(iii)),

$$p(\boldsymbol{\theta}|\hat{\Phi}, \Xi, \hat{\Psi}) = \frac{p(\boldsymbol{\theta}|\Xi)}{p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Psi})} p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Phi}). \quad (8)$$

Algorithm 1 Pseudocode for the Contour Monte Carlo algorithm for sampling from the posterior parameter distribution of eq. (8). Here we provide code for the Random Walk Metropolis algorithm for the MCMC sampling but for the examples in §4, we use an adaptive MCMC algorithm [21]. A definition of all variables is provided in Table 1.

```

procedure CMC( $\mathbf{X}(t), \Xi, N_1, N_2$ )  $\triangleright$  Sample from posterior parameter distribution
     $\hat{\Phi} = \text{SNAPSHOTESTIMATOR}(\mathbf{X}(t))$ 
     $\hat{\Psi} = \text{CONTOURVOLUMEESTIMATOR}(\Xi)$ 
     $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}) = \text{MCMC}(\hat{\Phi}, \Xi, \hat{\Psi})$ 
    converged = COMPAREOUTPUTTOTARGET(( $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}$ ),  $\hat{\Psi}$ )
    if converged  $\neq 1$  then
         $(\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}) = \text{CONTOURVOLUMEESTIMATOR}(\hat{\Phi}, \Xi, N'_1, N'_2)$ 
        where,  $N'_1 > N_1$  and/or  $N'_2 > N_2$   $\triangleright$  Rerun contour volume estimation and/or
        MCMC with larger sample sizes
    end if
    return ( $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}$ )
end procedure

procedure SNAPSHOTESTIMATOR( $\mathbf{X}(t)$ )  $\triangleright$  Fit density to snapshot observations
     $\mathbf{X}(t) \sim p(\tilde{\mathbf{x}}|\hat{\Phi})$ 
    return  $\hat{\Phi}$ 
end procedure

procedure CONTOURVOLUMEESTIMATOR( $\Xi$ )  $\triangleright$  Estimate volume of contours
    for  $i$  in  $1 : N_1$  do
         $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta}|\Xi)$   $\triangleright$  Sample from prior density
         $\tilde{\mathbf{x}}_i = \mathbf{g}(\boldsymbol{\theta}_i)$   $\triangleright$  Calculate corresponding output value
    end for
     $(\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_{N_1}) \sim p(\tilde{\mathbf{x}}|\hat{\Psi})$   $\triangleright$  Fit vine copula kernel density estimator to output
    values.
    return  $\hat{\Psi}$ 
end procedure

procedure MCMC( $\hat{\Phi}, \Xi, \hat{\Psi}$ )  $\triangleright$  Random Walk Metropolis algorithm targeting posterior
    parameter distribution.
     $\boldsymbol{\theta}_0 \sim \pi(\cdot)$   $\triangleright$  Sample from arbitrary initialisation distribution
    for  $i$  in  $1 : N_2$  do
         $\boldsymbol{\theta}'_i \sim \mathcal{N}(\boldsymbol{\theta}_{i-1}, \Sigma)$   $\triangleright$  Propose new parameter values for parameters
         $r = \left[ p(\boldsymbol{\theta}'|\Xi) p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Psi}) p(\mathbf{g}(\boldsymbol{\theta}')|\hat{\Phi}) \right] / \left[ p(\boldsymbol{\theta}|\Xi) p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Psi}) p(\mathbf{g}(\boldsymbol{\theta})|\hat{\Phi}) \right]$   $\triangleright$  Metropolis
        acceptance ratio.
         $u \sim U(0, 1)$   $\triangleright$  Sample from uniform distribution
        if  $r > u$  then
             $\boldsymbol{\theta}_i = \boldsymbol{\theta}'_i$   $\triangleright$  Accept proposal
        else
             $\boldsymbol{\theta}_i = \boldsymbol{\theta}_{i-1}$   $\triangleright$  Reject proposal
        end if
    end for
    return ( $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}$ )
end procedure

procedure COMPAREOUTPUTTOTARGET(( $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_{N_2}$ ),  $\hat{\Psi}$ )  $\triangleright$  Check output distribution
    close to target
    for  $i$  in  $1 : N_2$  do
         $\tilde{\mathbf{x}}_i = \mathbf{g}(\boldsymbol{\theta}_i)$   $\triangleright$  Compute output for each parameter sample
    end for
    if  $(\tilde{\mathbf{x}}_1, \tilde{\mathbf{x}}_2, \dots, \tilde{\mathbf{x}}_{N_2}) \sim p(\tilde{\mathbf{x}}|\hat{\Psi})$ ? then
        return 1  $\triangleright$  Compare outputs with target
         $\triangleright$  If outputs sufficiently close then converged
    else
        return 0
    end if
end procedure

```

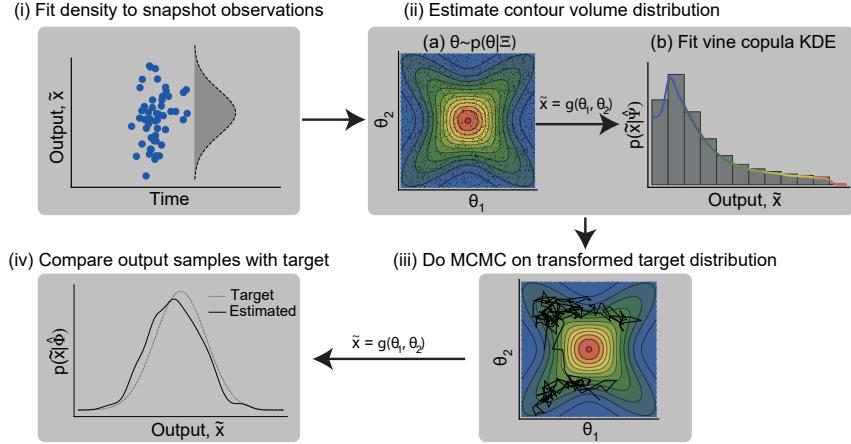


Figure 4: **The workflow for using Contour Monte Carlo to estimate cell population heterogeneity.** In (iii), the distribution targeted is given by eq. (8). The variables used in this figure are defined in the text and Table 1.

The final step in CMC is to transform parameter samples from the MCMC into outputs, then compare the sampled outputs with the target distribution (Figure 4(iv)). Asymptotically (in terms of the sample size of both sampling steps), CMC produces a sample of parameter values ($\theta^1, \theta^2, \dots$) which, when transformed to outputs, corresponds to the target distribution $p(\mathbf{x}|\hat{\Psi})$. In developing CMC, we have found that a finite sample of modest size for both steps of CMC results in parameter samples that, when transformed, often represent reasonable approximations of the target. There are however occasions when this is not the case and we have found this final confirmatory step indispensable since it frequently highlights inadequacies in the contour volume estimation or the MCMC, meaning more samples from either or both of these steps are required. It may also be necessary to tweak hyperparameters of the KDE to ensure reasonable approximation in the contour volume estimation step. If the target distribution is sensitive to the contour volume estimates, this may also indicate that the target snapshot distribution is incompatible with the model: here, we make no claims on existence of a solution to the inverse problem, only that, if one should exist, Contour Monte Carlo is a pragmatic approach to approximate it by sampling. A useful way to diagnose whether the target distribution can be produced from the model and specified priors is to examine the output values from the contour volume estimation step of CMC. If the majority of probability mass of the target lies outside the bounds of the bulk simulated output values obtained by independent sampling from the prior, then the model and/or chosen prior is unlikely to be invertible to this particular target.

In generating our results in §4, for the contour volume estimation step, we assumed sample sizes were sufficient if the output samples from the MCMC provided a reasonable approximation to the target, although we recognise that future work should refine this process further. For the MCMC step, we use adaptive covariance MCMC (see SOM of [21]) to sample from the target distribution, as we have found that it provides a considerable speed-up over Random Walk Metropolis [22, 23]. We also use the Gelman-Rubin convergence statistic \hat{R} which provides a heuristic

313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345

measurement of convergence [23, 24], and use a threshold of $\hat{R} \leq \sim 1.1$ to 346
diagnose convergence. 347

4 Results

348

In this section, we use CMC to estimate the posterior parameter distribution 349
for three biological systems of interest targeting synthetic parametric densi- 350
ties. That is, we assume that the first step of CMC (“SnapshotEstimator” 351
within Algorithm 1) has already been undertaken and we are faced with 352
inferring a parameter distribution which, when transformed to outputs, 353
recapitulates the target density. Alongside the text, we also provide a Julia 354
notebook used to generate the results, which we hope will be of use to 355
others wanting to apply CMC to estimate cell population heterogeneity. 356

Model	Target density	Parameter	Prior density	Prior parameter 1	Prior parameter 2
Growth factor	two-dimensional normal	R_T	uniform	2.5×10^5	8×10^5
		k_1	uniform	0.25	3.0
		k_{-1}	uniform	2.0	20.0
		k_{deg}	uniform	0.005	0.03
		k_{deg}^*	uniform	0.1	0.5
Growth factor	two-dimensional normal	R_T	normal	5×10^5	1×10^5
		k_1	normal	0.5	0.1
		k_{-1}	normal	3.0	1.0
		k_{deg}	normal	0.02	0.005
		k_{deg}^*	normal	0.3	0.1
Michaelis-Menten kinetics	bimodal normal	k_f	uniform	0.2	15
		k_r	uniform	0.2	2.0
		k_{cat}	uniform	0.5	3.0
Michaelis-Menten kinetics	four-dimensional normal	k_f	uniform	0.2	15
		k_r	uniform	0.2	2.0
		k_{cat}	uniform	0.2	3.0
		E_0	uniform	3.0	5.0
		S_0	uniform	5.0	10.0
		ES_0	uniform	0.0	0.2
		P_0	uniform	0.0	0.2
TNF signalling	bivariate normal	a_1	uniform	0.4	0.8
		a_2	uniform	0.1	0.7
		a_3	uniform	0.3	0.7
		a_4	uniform	0.1	0.3
		b_1	uniform	0.5	0.7
		b_2	uniform	0.4	0.6
		b_3	uniform	0.4	0.6
		b_4	uniform	0.2	0.4
		b_5	uniform	0.2	0.4
		a_1	uniform	0.5	0.7
TNF signalling	bimodal normal	a_2	uniform	0.1	0.3
		a_3	uniform	0.1	0.3
		a_4	uniform	0.4	0.6
		b_1	uniform	0.3	0.5
		b_2	uniform	0.6	0.8
		b_3	uniform	0.2	0.4
		b_4	uniform	0.4	0.6
		b_5	uniform	0.3	0.5

Table 2: The priors used for each problem in §4.

4.1 Growth factor model

357

Here we consider the “growth factor model” introduced by [12], which concerns the dynamics of inactive ligand-free cell surface receptors R and active ligand-bound cell surface receptors P , modulated by an exogenous

ligand L . The governing dynamics are determined by the following system,

$$\dot{R}(t) = R_T k_{deg} + k_1 L R(t) + k_{-1} P(t) - k_{deg} R(t) \quad (9)$$

$$\dot{P}(t) = k_1 L R(t) - k_{-1} P(t) - k_{deg}^* P(t), \quad (10)$$

where $\theta = (R_T, k_1, k_{-1}, k_{deg}, k_{deg}^*)$ are parameters to be determined. In this example, we use measurements of the active ligand-bound receptors P to estimate cellular heterogeneity in processes. We denote the solution of eq. (10) as $P(t; \theta, L)$ and seek to determine the parameter distribution consistent with an output distribution,

$$\begin{pmatrix} P(10; \theta, 2) \\ P(10; \theta, 10) \end{pmatrix} \sim \mathcal{N} \left[\begin{pmatrix} 2 \times 10^4 \\ 3 \times 10^4 \end{pmatrix}, \begin{pmatrix} 1 \times 10^5 & 0 \\ 0 & 1 \times 10^5 \end{pmatrix} \right]. \quad (11)$$

To start, we specify a uniform prior for each of the five parameters, with bounds given in Table 2. To estimate the posterior parameter distribution, we use CMC, with adaptive covariance MCMC [21] for the second step.

In Figure 5A, we show the sampled outputs (blue points) versus the contours of the target distribution (black solid closed curves), illustrating a good correspondence between the sampled and target densities. Above and to the right of the main panel, we also display the marginal target densities (solid black lines) versus kernel density estimator reconstructions of the output marginals from the CMC samples (dashed blue lines), which again highlights the fidelity of the CMC sampled density to the target.

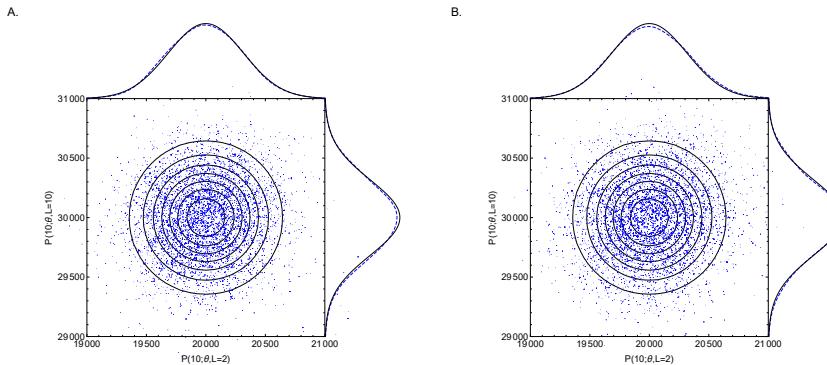


Figure 5: The target joint output distribution (solid contour lines) and target marginal distributions (solid lines; above and to side of figure) versus outputs sampled by CMC (blue points and dashed lines) for (A) uniform and (B) normal parameter priors. In CMC, 100,000 independent samples were used in the “ContourVolumeEstimator” step and 10,000 MCMC samples across each of 4 Markov chains were used in the second step, with the first half of the chains discarded as “warm-up” [23]. For the reconstructed marginal densities in the plots, we use Mathematica’s “SmoothKernelDistribution” function specifying bandwidths of 100 with Gaussian kernels [25].

In Figure 6A, we plot the joint posterior parameter distribution for k_1 , the rate of ligand binding to inactive receptors, and k_{-1} , which dictates the rate of the reverse reaction, where the ligands unbind. The output

measurements we used to fit the model correspond to levels of the bound ligands, which can be generated whenever the ratio of k_1 to k_{-1} is approximately given by the corresponding steady state ratio. Because of this, the distribution representing cell process heterogeneity contains linear positive correlations between these parameters. In Figure 6B, we show the posterior parameter distribution for k_{deg} , the rate of degradation of ligand-free cell surface receptors and R_T , which dictates the rate of introduction of ligand-free cell surface receptors, which shows a concentrated region of posterior probability mass. Why is it that we are better able to resolve (k_{deg}, R_T) compared to (k_1, k_{-1}) from our measurements? To answer this, it is useful to calculate the sensitivity of $P(t; \theta, L)$ to changes in each of the parameters. To account for the differing magnitudes of each parameter, we calculate elasticities, the proportional changes in measured output for a proportional change in parameter values, using the forward sensitivities method described in [26], which are shown in Figure 7. When the exogenous ligand is set $L = 2$, these indicate the active ligand-bound receptor concentration is most elastic to changes in R_T and k_{deg} , meaning that their range is more restricted by the output measurement than for k_1 and k_{-1} , which have elasticities at $t = 10$ closer to 0. In Table 3, we show the posterior quantiles for the estimated parameters and, in the last column, indicate the ratio of the 25%-75% posterior interval widths to the uniform prior range for each parameter. These were strongly negatively correlated with the magnitude of the elasticities for each parameter ($\rho = 0.95$, $t = -5.22$, $df = 3$, $p = 0.01$ Pearson's product-moment correlation), indicating the utility of sensitivity analyses for optimal experimental design. We would suggest however that CMC can also be used for this purpose, using synthetic data in place of real measurements.

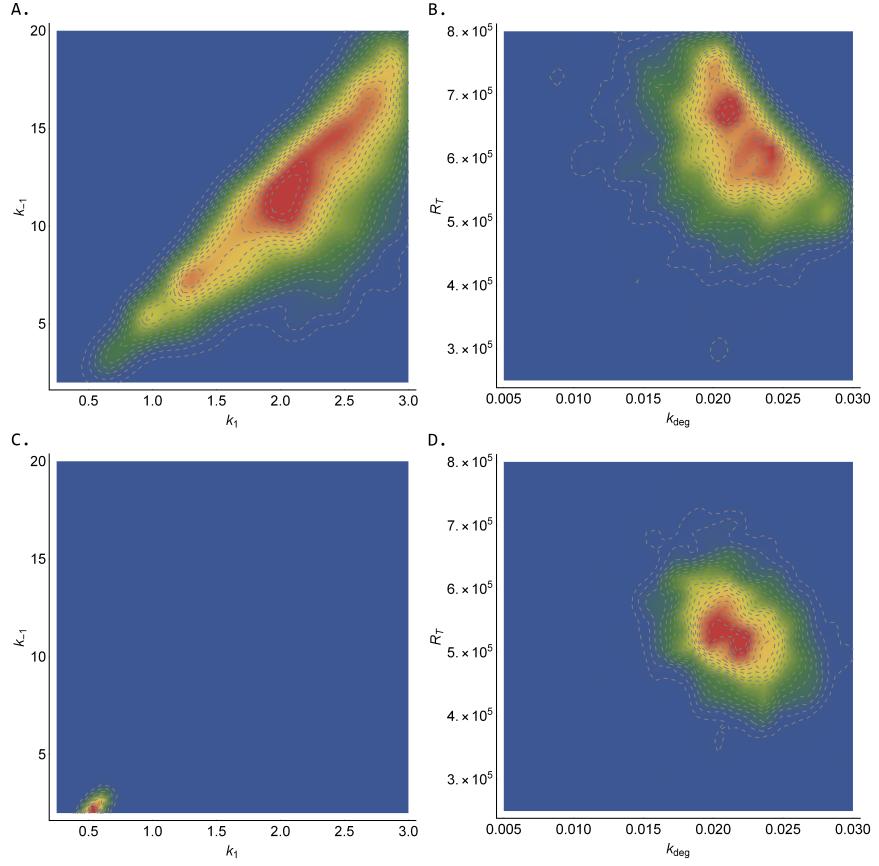


Figure 6: The joint distributions of (k_1, k_{-1}) (left-column) and (k_{deg}, R_T) for the growth factor model using uniform priors (top row) and normal priors (bottom row). See Figure 5 caption for CMC details and Table 2 for the priors used.

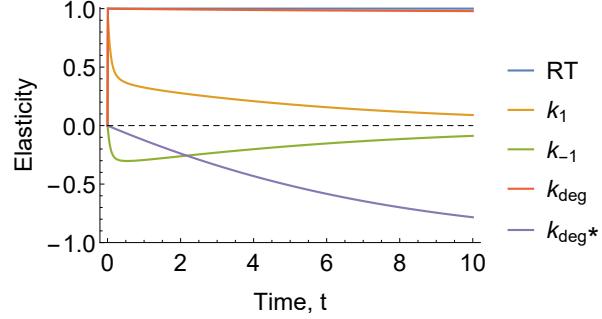


Figure 7: The elasticities of the measured concentration of active ligand-bound receptors P versus time when $L = 2$. When calculating the elasticities of each parameter, the other parameters were set to their posterior medians given in Table 3.

For an unidentified model, there is typically a multitude of possible probability distributions over parameters which map to the same target output distribution. To reduce the space of posterior parameter distributions to one,

403

404

405

Prior	Parameter	2.5%	25%	Quantiles 50%	75%	97.5%	Posterior 25%-75% concentration
Uniform	R_T	441,006	548,275	606,439	677,055	772,484	23%
	k_1	0.90	1.69	2.17	2.56	2.95	32%
	k_{-1}	4.35	8.35	11.23	14.23	18.71	33%
	k_{deg}	0.013	0.019	0.021	0.024	0.029	20%
	k_{deg}^*	0.20	0.34	0.40	0.44	0.49	27%
Normal	R_T	408,396	487,372	529,558	577,970	678,632	16%
	k_1	0.39	0.49	0.54	0.60	0.70	4%
	k_{-1}	1.39	1.92	2.26	2.63	3.35	4%
	k_{deg}	0.016	0.020	0.022	0.024	0.027	16%
	k_{deg}^*	0.22	0.29	0.33	0.38	0.46	21%

Table 3: **Estimated quantiles from CMC samples for the growth factor model with uniform and normal priors.** The last column indicates the proportion of the uniform prior bounds occupied by the 25%-75% posterior interval in each case. The particular priors used in each case are given in Table 2.

it is therefore necessary to specify a prior parameter distribution. It is also preferable to allow priors to influence estimates in studies of cellular heterogeneity, since this allows incorporation of pre-existing biological knowledge with compensatory reductions in estimator variance. CMC accommodates different prior choices, with both the “ContourVolumeEstimation” step and the acceptance ratio in the “MCMC” step (Algorithm 1) being affected in such a way that posterior parameter distribution maps to the same output target. We now use CMC to estimate the posterior parameter distribution when changing from uniform priors to more concentrated normal priors (prior hyperparameters shown in Table 2). As desired, the target output distribution appears invariant (Figure 5B) although with substantial changes in the posterior parameter distributions (Figure 6C&D). In particular, the posterior distributions obtained from shifting to the normal prior are more concentrated in parameter space compared to the uniform case (rightmost column of Table 3). The differences in posterior distribution resultant from changes to priors are likely to be more marked the less definitive a guide the data provides on the underlying process and, hence, can be used to stabilise the resultant inferences according to external knowledge about the system.

4.2 Michaelis-Menten kinetics

In this section, we use CMC to invert output measurements from the Michaelis-Menten model of enzyme kinetics (see, for example, [27]); illustrating the capability of CMC to resolve population substructure from multimodality of the output distribution. The Michaelis-Menten model of enzyme kinetics describes the dynamics of concentrations of an enzyme (E), a substrate (S), an enzyme-substrate complex (ES), and a product (P). Specifically,

$$\begin{aligned}\dot{E}(t) &= -k_f E(t)S(t) + k_r ES(t) + k_{cat}ES(t), \\ \dot{S}(t) &= -k_f E(t)S(t) + k_r ES(t), \\ \dot{ES}(t) &= k_f E(t)S(t) - k_r ES(t) - k_{cat}ES(t), \\ \dot{P}(t) &= k_{cat}ES(t),\end{aligned}\tag{12}$$

with initial conditions,

433

$$E(0) = E_0, S(0) = S_0, ES(0) = ES_0, P(0) = P_0, \quad (13)$$

where k_f is the rate constant for the forward reaction $E + S \rightarrow ES$, k_r is the rate of the reverse reaction $ES \rightarrow E + S$, and k_{cat} is the catalytic rate at which the product is formed by the reaction $ES \rightarrow E + P$.

434

435

436

When subpopulations of cells, each with distinct dynamics, are thought to exist, determining their characteristics - proportions of overall cell number, likely parameter values, and so on - is often of key interest [15, 19]. Before formal inference occurs, multi-modality of the output distribution may signal the existence of fragmented subpopulations of cells. Here we target the following bimodal bivariate normal distribution,

$$f(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma_1, \boldsymbol{\mu}_2, \Sigma_2) = \frac{1}{2} (\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \Sigma_1) + \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \Sigma_2)), \quad (14)$$

where $\mathbf{x} = (E(2; \boldsymbol{\theta}), S(1; \boldsymbol{\theta}))$ with each element corresponding to the solutions of eq. (12) for the enzyme and substrate at times $t = 2$ and $t = 1$, respectively, and $\boldsymbol{\theta} = (k_f, k_r, k_{cat})$. The parameters of the mixture normal output distribution we target are $\boldsymbol{\mu}_1 = [2.2, 1.6]'$, $\Sigma_1 = \begin{pmatrix} 0.018 & -0.013 \\ -0.013 & 0.010 \end{pmatrix}$, $\boldsymbol{\mu}_2 = [2.8, 1.0]'$ and $\Sigma_2 = \begin{pmatrix} 0.020 & -0.010 \\ -0.010 & 0.020 \end{pmatrix}$. In what follows, we specify uniform priors on each of the elements of $\boldsymbol{\theta}$ (see Table 2).

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

Using a modest number of samples in each step, CMC was able to recapitulate the output target distribution (Figure 8A). Without specifying *a priori* information on the subpopulations of cells, two distinct clusters of cells emerged from application of CMC (orange and blue points in Figure 8B), each corresponding to distinct modes of the output distribution (corresponding coloured points in Figure 8A). It is worth noting however that the issues inherent with MCMC sampling of multimodal distributions similarly apply here and so, whilst here adaptive MCMC [21] sufficed to explore the posterior surface, it may be necessary to use MCMC methods known to be robust to such geometries (for example, population MCMC [28]).

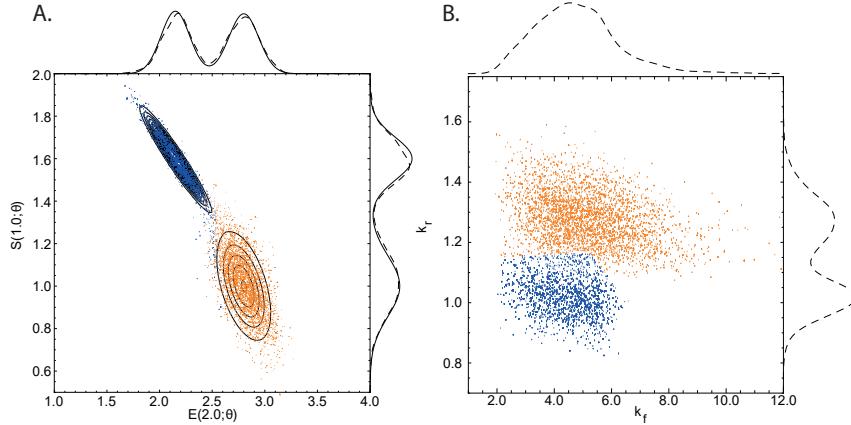


Figure 8: **(A)** bimodal target distribution (solid contour lines) versus output samples (points) for the Michaelis-Menten model and **(B)** posterior parameter samples (points). The solid and dashed lines above and to the side of panel A indicate the target and estimated marginal output distributions, respectively. The orange (blue) points in A were generated by the orange (blue) parameter samples in B. See Figure 5 caption for CMC details. For the reconstructed marginal densities in the plots, we use Mathematica’s “SmoothKernelDistribution” function with Gaussian kernels with (A) default bandwidths and (B) bandwidths of 0.3 (horizontal axis) and 0.03 (vertical axis) [25]. The clusters were identified by applying Mathematica’s “ClusteringComponents” function to the pairs of parameter samples displayed in B [25].

Loos et al. (2018) consider a multidimensional output distribution, with correlations between system characteristics that evolve over time. Our approach allows arbitrary covariance structure between measurements, and to exemplify this, we now target a four-dimensional output distribution, with paired measurements of enzyme and substrate at $t = 1$ and $t = 2$,

$$\begin{pmatrix} E(1.0; \theta) \\ S(1.0; \theta) \\ E(2.0; \theta) \\ S(2.0; \theta) \end{pmatrix} \sim \mathcal{N} \left[\begin{pmatrix} 0.5 \\ 2.8 \\ 0.9 \\ 1.4 \end{pmatrix}, \begin{pmatrix} 0.02 & -0.05 & 0.04 & -0.05 \\ -0.05 & 0.30 & -0.15 & 0.20 \\ 0.04 & -0.15 & 0.12 & -0.17 \\ -0.05 & 0.20 & -0.17 & 0.30 \end{pmatrix} \right]. \quad (15)$$

Since this system has four output measurements, and the Michaelis-Menten model has three rate parameters (k_f, k_r, k_{cat}), the system is over-identified and so CMC cannot be straightforwardly applied. Instead, we allow the four initial states (E_0, S_0, ES_0, P_0) to be uncertain quantities, bringing the total number of parameters to 7, and ensuring that the system is in the unidentified regime where CMC applies. We set uniform priors on all parameters (see Table 2) and to check that the model and priors were consistent with the output distribution given by eq. (15), we plotted the output measurements used to estimate contour volumes (in the first step of the “ContourVolumeEstimator” method in Algorithm 1) against the target (Figure 9). Since the main support of the densities (black contours) lies within a region of output space reached by independent sampling of the priors (blue points), this indicated that the distribution given by eq. (15) could feasibly be generated from this model and priors, and we proceeded to estimation by CMC.

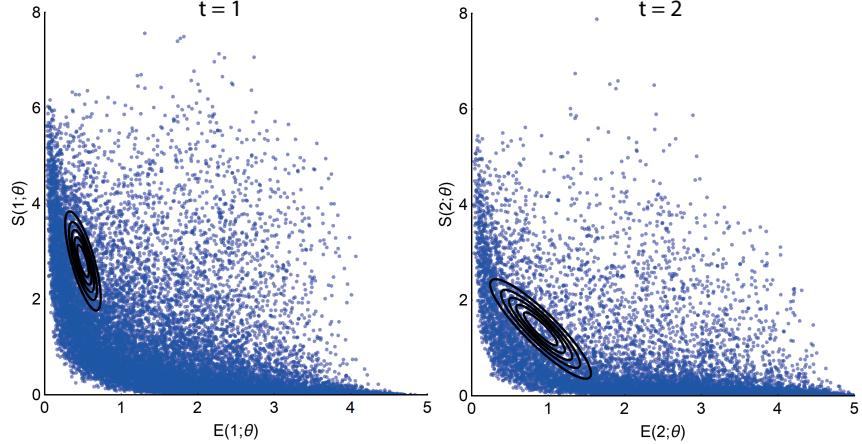


Figure 9: **The output measurements (blue points)** of enzyme and substrate at times $t = 1$ (left panel) and $t = 2$ (right panel) obtained by independently sampling the priors of the 7 parameter Michaelis-Menten model versus the target distribution (black solid contours). In this plot, we show 20,000 output samples, where each set of four measurements was obtained from a single sample of the 7 parameters. The output target distribution shown by the contours corresponds to the marginal densities of each pair of enzyme-substrate measurements given by eq. (15).

Figure 10 plots the output samples of enzyme and substrate from the last step of CMC for $t = 1$ (blue points) and $t = 2$ (orange points) versus the contours (black lines) of the joint marginal distributions of eq. (15). The distribution of paired enzyme-substrate samples illustrates that the CMC output samples approximated the target density, itself representing dynamic evolution of the covariance between enzyme and substrate measurements. The target marginal distributions (solid lines) along with their approximations from kernel density estimation (dashed lines) are also shown above and beside the main panel of Figure 10, and largely indicate correspondence.

474
475
476
477
478
479
480
481
482
483

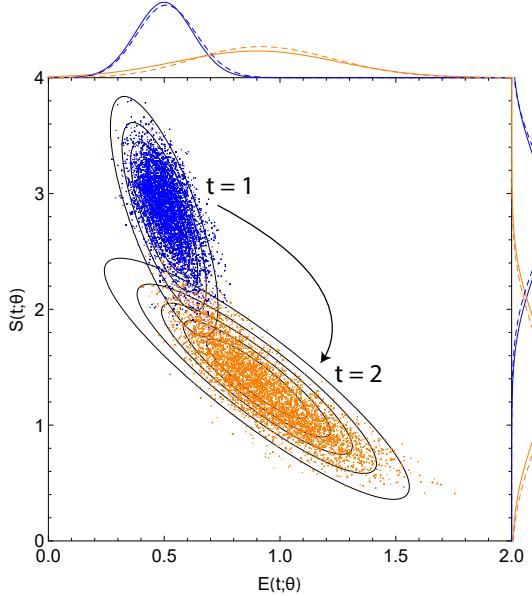


Figure 10: **Posterior output samples from CMC (coloured points) versus the contour plots of the joint marginal distributions of eq. (15) (black solid lines).** The blue and orange points indicate output observations of enzyme (horizontal axis) and substrate (vertical axis) at times $t = 1$ and $t = 2$ respectively. The solid and dashed coloured lines outside of the panels indicate the true target marginals of eq. (15) and those estimated from CMC, respectively. The line colours correspond with the point colours and indicate the marginals for $t = 1$ (blue) and $t = 2$ (orange). In CMC, 200,000 independent samples were used in the “ContourVolumeEstimator” step and 10,000 MCMC samples across each of 4 Markov chains were used in the second step, with the first half of the chains discarded as “warm-up” [23]. For the reconstructed marginal densities in the plots, we use Mathematica’s “SmoothKernelDistribution” function with Gaussian kernels [25] of bandwidths varying from 0.1 to 0.4.

4.3 TNF signalling pathway

We now illustrate how CMC can be applied to an ODE system of larger size, the tumour necrosis factor (TNF) signalling pathway model introduced in [29] and used by [15] to illustrate a Bayesian approach to cell population variability estimation. The model incorporates known activating and inhibitory interactions between four key species within the TNF pathway: active caspase 8 (X_1) and active caspase 3 (X_2), a nuclear transcription factor (X_3) and its inhibitor (X_4),

$$\begin{aligned} \dot{X}_1(t) &= -X_1(t) + \frac{1}{2} (\text{inh}_4(X_3(t))\text{act}_1(u(t)) + \text{act}_3(X_2(t))) \\ \dot{X}_2(t) &= -X_2(t) + \text{act}_2(X_1(t))\text{inh}_3(X_3(t)) \\ \dot{X}_3(t) &= -X_3(t) + \text{inh}_2(X_2(t))\text{inh}_5(X_4(t)) \\ \dot{X}_4(t) &= -X_4(t) + \frac{1}{2} (\text{inh}_1(u(t)) + \text{act}_4(X_3(t))), \end{aligned} \quad (16)$$

484
485
486
487
488
489
490
491

where the functions act_i and inh_j represent activating and inhibitory interactions respectively,

$$\begin{aligned}\text{act}_i(X) &= \frac{X^2}{a_i^2 + X^2} \\ \text{inh}_j(X) &= \frac{b_j^2}{b_j^2 + X^2},\end{aligned}\tag{17}$$

and the parameters a_i for $i \in (1, 2, 3, 4)$ and b_j for $j \in (1, 2, 3, 4, 5)$ represent activation and inhibition thresholds. We assume that the initial states of the system are $(X_1(0), X_2(0), X_3(0), X_4(0)) = (0.0, 0.0, 0.29, 0.625)$, which correspond to the steady state of the system when $X_2 = 0$. The function $u(t)$ represents a TNF stimulus which is given by a top hat function,

$$u(t) = \begin{cases} 1, & \text{if } t \in [0, 2], \\ 0, & \text{otherwise.} \end{cases}\tag{18}$$

When there are fewer output measurements than parameters, models tend to be underdetermined meaning that many combinations of parameters can lead to the same combination of output values. A consequence of this unidentifiability is that we cannot perform “full circle” inference: that is, using a known parameter distribution to generate an output distribution does not result in that parameter distribution being recapitulated through inference. We illustrate this idea by generating an output distribution by varying a single parameter value between runs of the forward model corresponding to the solution of eq. (16) and performing inference on all nine system parameters, whilst collecting only two output measurements. Specifically, we vary $a_1 \sim \mathcal{N}(0.6, 0.05)$, whilst holding the other parameters constant $(a_2, a_3, a_4, b_1, b_2, b_3, b_4, b_5) = (0.2, 0.2, 0.5, 0.4, 0.7, 0.3, 0.5, 0.4)$ and measure $X_1(2.0)$ and $X_2(1.0)$. By running the forward model, we obtain an output distribution that is well approximated by the bivariate normal distribution,

$$\begin{pmatrix} X_1(2.0) \\ X_2(1.0) \end{pmatrix} \sim \mathcal{N} \left[\begin{pmatrix} 0.26 \\ 0.07 \end{pmatrix}, \begin{pmatrix} 2.1 \times 10^{-4} & 5.9 \times 10^{-5} \\ 5.9 \times 10^{-5} & 1.8 \times 10^{-5} \end{pmatrix} \right].\tag{19}$$

We now apply CMC to the target output distribution given by eq. (19) to estimate a posterior distribution over all nine parameters of eq. (16). Apart from a few cases, the priors for each parameter were chosen to exclude the values that were used to generate the output distribution (see Table 2), to illustrate the non-equivalence between the recovered posterior distribution and the data generating process. In Figure 11A, we plot the actual parameter values (horizontal axis) used in the true data generating process versus the inferred values (vertical axis). This illustrates that apart from a_1 , where the estimated parameter values correspond well with the range of values used to generate the data, due to the choice of priors there is a disjunction between the actual and estimated values. Despite these differences, due to the model being underdetermined, it is nonetheless possible to use CMC to sample from an output distribution that well approximates the target (Figure 11B).

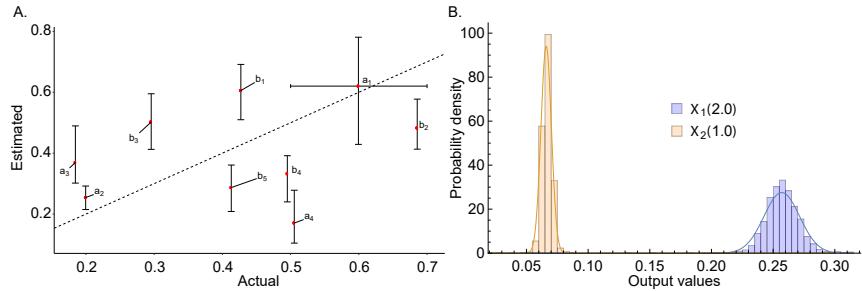


Figure 11: (A) actual parameter values versus estimated quantiles for the output distribution for the TNF signalling pathway model with output distribution given and (B) the marginal output target (solid lines) given by eq. (19) and sampled output distribution (histograms). In A, in the vertical direction, the red points indicate the 50% posterior quantiles and the upper and lower whiskers indicate the 97.5% and 2.5% quantiles, respectively; in the horizontal direction, with the exception of a_1 , the red points indicate the parameter values used to generate the data; for a_1 the red point indicates the mean of the normal distribution used to generate the data and the whiskers indicate the 95% quantiles of this distribution. In CMC, 10,000 independent samples were used in the ‘ContourVolumeEstimator’ step and 5,000 MCMC samples across each of 4 Markov chains were used in the second step, with the first half of the chains discarded as ‘warm-up’ [23].

Cell populations may be well described by subpopulations which each evolve along characteristic trajectories over time. We now apply CMC to investigate a bimodal output distribution for the TNF signalling pathway model similar to that investigated by [15]. In particular, we aim to find a distribution over parameter values which, when used as inputs to the solution to the ODE system, results in the following output distribution,

$$\begin{aligned} \mathbf{X}_2(1.0) &\sim \mathcal{N}(0.06, 0.01) \\ \mathbf{X}_2(2.0) &\sim \frac{1}{2} (\mathcal{N}(0.1, 0.01) + \mathcal{N}(0.14, 0.01)) \\ \mathbf{X}_2(4.0) &\sim \frac{1}{2} (\mathcal{N}(0.1, 0.01) + \mathcal{N}(0.20, 0.01)), \end{aligned} \quad (20)$$

where via a slight abuse of notation, the target distributions for $\mathbf{X}_2(2.0)$ and $\mathbf{X}_2(4.0)$ indicate mixtures of univariate normals, and the priors used are shown in Table 2. This target distribution, along with the unique trajectories obtained by applying the CMC algorithm for 5,000 MCMC steps, are shown in Figure 12. This figure illustrates that given by bimodality of the output distribution, CMC estimates a corresponding subpopulation structure in the parameter distribution without *a priori* specification of the number of clusters.

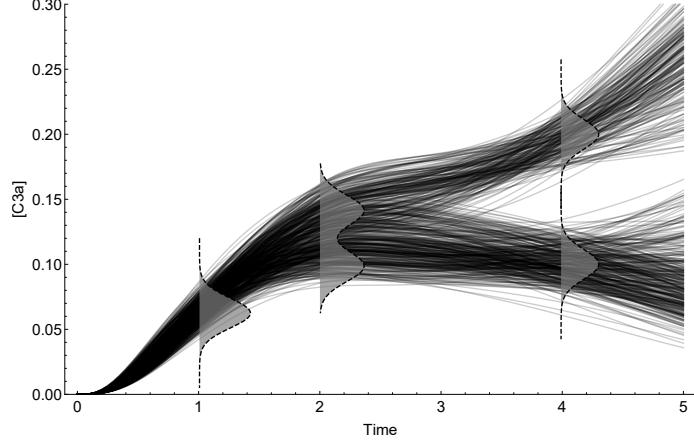


Figure 12: **The target output distribution (dashed plots with grey filling) and unique trajectories (black solid lines) obtained from the posterior parameter distribution.** In CMC, 10,000 independent samples were used in the “ContourVolumeEstimator” step and 5,000 MCMC samples across each of 4 Markov chains were used in the second step, with the first half of the chains discarded as “warm-up” [23].

5 Discussion

Determining the cause of variability in cellular processes is crucial in many applications, ranging from bioengineering to drug development. In this paper, we introduce a Bayesian method for estimating cellular heterogeneity from “snapshot” measurements of cellular properties, taken at discrete intervals throughout the experimental course. Our approach assumes what we term a “heterogeneous ordinary differential equation” (HODE) framework, in which biochemical processes in individual cells are assumed to follow dynamics governed by a common ODE, although with idiosyncratic differences in parameter values. In this framework, estimating heterogeneity in cellular processes amounts to determining the probability distributions over parameter values of the governing ODE. Our method of estimation is a two-step Monte Carlo sampling process we term “Contour Monte Carlo” (CMC) which does not require *a priori* specification of cell population substructure unlike other approaches. CMC can be used to process high volumes of individual cellular measurements since the framework involves fitting a kernel density estimator to raw experimental data and using these distributions rather than data as the target outcome. CMC also allows for arbitrary multivariate structure in the measurement space, meaning it can capture correlations that occur between the same cellular species at different timepoints or, for example, contemporaneous correlations between different cellular compartments. Being a Bayesian approach, CMC uses prior distributions over parameter values to ensure uniqueness of the posterior distribution, allowing pre-experimental knowledge to be used to improve estimation robustness. The flexible and robust framework that CMC provides means it can be used to do automatic inference for wide-ranging systems of practical interest.

As well as providing a framework for estimating cellular variation, our approach also provides a natural way to test that it is working as

desired. By feeding the posterior parameter samples obtained by CMC into the forward model, this results in a distribution over output values that can be compared to the target. Indeed, we have found this comparison indispensable in applying CMC in practice and include it as the last step in the CMC algorithm (Algorithm 1). Discrepancies between the target output distribution and samples from it by CMC can occur either as a result of poor estimates of the “contour volume distribution” in the first stage of the algorithm or due to insufficient MCMC samples in the second. Either of these issues can often be easily addressed and although kernel density estimation in high dimensional spaces remains an open research problem, we have found vine copula kernel density estimation works well for the dimensionality of output measurements we investigate here [20].

Failure to reproduce a given output distribution can also indicate that the generating model (the priors and the forward model) are incongruent with experimental results. This may either be due to misspecification of the ODE system or the inadequacy of the assumed deterministic framework. Our approach currently assumes that output stochasticity is dominated by cellular variation in the parameter values of the underlying ODE, with measurement noise making a minor contribution. Whether this is a reasonable assumption depends on the system under investigation and, more importantly, on experimental details. We recognise that neglecting measurement noise when it is an important determinant of the observed data is likely to mean we overstate the degree of cellular variation. It may also mean that some output distributions cannot be obtained through our assumed model system. Future work allowing inclusion of a stochastic noise process or, more generally, including stochastic cellular mechanisms is thus likely to be worthwhile.

Whilst we have labelled the approach we follow here as Bayesian, since it involves explicit estimation of probability distributions and involves priors over parameter values, we recognise that it is not in the form typically utilised by exponents of this framework. This is because rather than aiming to formulate a model that describes output observations, instead, our approach aims to recapitulate output distributions. Others [30], (including us [31]), have considered this problem before; perhaps most notably by Albert Tarantola in his landmark work on inverse problem theory (see, for example, [32]). In Tarantola’s framework, a joint input parameter & output space is considered, where prior knowledge and experimental theory combine elegantly to produce a posterior distribution whose marginal output distribution matches the experimentally obtained one. This work has seen considerable interest in areas such as the geosciences [33,34], and we propose that these methods may prove useful for the biosciences. In particular, we posit that this framework may permit a generalised Bayesian inference that can more parsimoniously encompass output data and output distributions as their outcome measures.

The natural world is rife with variation. Mathematical models represent frameworks for understanding the causes of such variation. Typically, the state of biological knowledge is such that one effect, a given pattern of variation, has many possible causes, and observational or experimental data are necessary to apportion weight to each of them, in a process which amounts to solving inverse an problem. The approach we describe here follows the Bayesian paradigm of inverse problem solving whereby uncertainty in potential causes is reflected by probability distributions. Here, we illustrate the utility of our method by applying it to estimate cellular

heterogeneity in biochemical processes however, it could equally well be used
 to understand the inversion of deterministic systems more generally. Whilst
 describing the inversion process of deterministic models using probability
 distributions may sound contradictory, it is worth acknowledging that many
 ODE systems are structurally unidentified meaning there is incompressible
 uncertainty over some regions of parameter space. Contour Monte Carlo
 provides an automatic framework for doing inference on these deterministic
 systems and the use of priors allows for robust and precise parameter
 estimation unattainable through the data alone. 624
625
626
627
628
629
630
631
632

6 Author contributions

633

BL, DJG and SJT conceived the study. BL carried out the analysis. All
 authors helped to write and edit the manuscript. 634
635

References

- [1] M Ridley. *The red queen: Sex and the evolution of human nature*. Penguin UK, 1994.
- [2] D Fraser and M Kaern. A chance at survival: gene expression noise and phenotypic diversification strategies. *Molecular Microbiology*, 71(6):1333–1340, 2009.
- [3] F Delvigne, Q Zune, AR Lara, W Al-Soud, and SJ Sørensen. Metabolic variability in bioprocessing: implications of microbial phenotypic heterogeneity. *Trends in Biotechnology*, 32(12):608–616, 2014.
- [4] RA Gatenby, K Smallbone, PK Maini, F Rose, J Averill, Raymond B Nagle, L Worrall, and RJ Gillies. Cellular adaptations to hypoxia and acidosis during somatic evolution of breast cancer. *British Journal of Cancer*, 97(5):646, 2007.
- [5] PM Altrock, LL Liu, and F Michor. The mathematics of cancer: integrating quantitative models. *Nature Reviews Cancer*, 15(12):730, 2015.
- [6] SJ Altschuler and LF Wu. Cellular heterogeneity: do differences make a difference? *Cell*, 141(4):559–563, 2010.
- [7] MB Elowitz, AJ Levine, ED Siggia, and PS Swain. Stochastic gene expression in a single cell. *Science*, 297(5584):1183–1186, 2002.
- [8] HH Chang, M Hemberg, M Barahona, DE Ingber, and S Huang. Transcriptome-wide noise controls lineage choice in mammalian progenitor cells. *Nature*, 453(7194):544, 2008.
- [9] S Waldherr. Estimation methods for heterogeneous cell population models in systems biology. *Journal of The Royal Society Interface*, 15(147):20180530, 2018.
- [10] R Erban, J Chapman, and P Maini. A practical guide to stochastic simulations of reaction-diffusion processes. *arXiv preprint arXiv:0704.1908*, 2007.

- [11] D Ramkrishna and MR Singh. Population balance modeling: current status and future prospects. *Annual Review of Chemical and Biomolecular Engineering*, 5:123–146, 2014.
- [12] P Dixit, E Lyashenko, M Niepel, and D Vitkup. Maximum entropy framework for inference of cell population heterogeneity in signaling network dynamics. *bioRxiv*, page 137513, 2018.
- [13] WG Telford, T Hawley, F Subach, V Verkhusha, and RG Hawley. Flow cytometry of fluorescent proteins. *Methods*, 57(3):318–330, 2012.
- [14] AJ Hughes, DP Spelke, Z Xu, CC Kang, DV Schaffer, and AE Herr. Single-cell western blotting. *Nature methods*, 11(7):749, 2014.
- [15] J Hasenauer, S Waldherr, M Doszczak, N Radde, P Scheurich, and F Allgöwer. Identification of models of heterogeneous cell populations from population snapshot data. *BMC Bioinformatics*, 12(1):125, 2011.
- [16] O Hilsenbeck, M Schwarzfischer, S Skylaki, B Schauberger, PS Hoppe, D Loeffler, KD Kokkaliaris, S Hastreiter, E Skylaki, A Filipczyk, et al. Software tools for single-cell tracking and quantification of cellular and molecular properties. *Nature Biotechnology*, 34(7):703, 2016.
- [17] FSO Fritzsch, C Dusny, O Frick, and A Schmid. Single-cell analysis in biotechnology, systems biology, and biocatalysis. *Annual Review of Chemical and Biomolecular Engineering*, 3:129–155, 2012.
- [18] J Hasenauer, C Hasenauer, T Hucho, and FJ Theis. Ode constrained mixture modelling: a method for unraveling subpopulation structures and dynamics. *PLOS Computational Biology*, 10(7):e1003686, 2014.
- [19] C Loos, K Moeller, F Fröhlich, T Hucho, and J Hasenauer. A hierarchical, data-driven approach to modeling single-cell populations predicts latent causes of cell-to-cell variability. *Cell Systems*, 6(5):593–603, 2018.
- [20] T Nagler and C Czado. Evading the curse of dimensionality in non-parametric density estimation with simplified vine copulas. *Journal of Multivariate Analysis*, 151:69–89, 2016.
- [21] RH Johnstone, ETY Chang, R Bardenet, TP De Boer, DJ Gavaghan, P Pathmanathan, RH Clayton, and GR Mirams. Uncertainty and variability in models of the cardiac action potential: Can we build trustworthy models? *Journal of Molecular and Cellular Cardiology*, 96:49–62, 2016.
- [22] N Metropolis, AW Rosenbluth, MN Rosenbluth, AH Teller, and E Teller. Equation of state calculations by fast computing machines. *The Journal of Chemical Physics*, 21(6):1087–1092, 1953.
- [23] B Lambert. *A Student’s Guide to Bayesian Statistics*. Sage Publications Ltd., 2018.
- [24] A Gelman and DB Rubin. Inference from iterative simulation using multiple sequences. *Statistical Science*, pages 457–472, 1992.
- [25] Inc. Wolfram Research. Mathematica 8.0. <https://www.wolfram.com>.

- [26] AC Daly, DJ Gavaghan, J Cooper, and SJ Tavener. Inference-based assessment of parameter identifiability in nonlinear biological models. *Journal of The Royal Society Interface*, 15, 2018.
- [27] JD Murray. *Mathematical biology: I. An Introduction (interdisciplinary applied mathematics)(Pt. 1)*. New York, Springer, 2007.
- [28] A Jasra, DA Stephens, and CC Holmes. On population-based simulation for static inference. *Statistics and Computing*, 17(3):263–279, 2007.
- [29] M Chaves, T Eissing, and F Allgower. Bistable biological systems: A characterization through local compact input-to-state stability. *IEEE Transactions on Automatic Control*, 53(Special Issue):87–100, 2008.
- [30] T Butler, D Estep, SJ Tavener, C Dawson, and JJ Westerink. A measure-theoretic computational method for inverse sensitivity problems iii: Multiple quantities of interest. *SIAM/ASA Journal on Uncertainty Quantification*, 2(1):174–202, 2014.
- [31] B Lambert, D Gavaghan, and SJ Tavener. Inverse sensitivity analysis of mathematical models avoiding the curse of dimensionality. *BioRxiv*, page 432393, 2018.
- [32] A Tarantola. *Inverse problem theory and methods for model parameter estimation*, volume 89. SIAM, 2005.
- [33] K Mosegaard and A Tarantola. Monte carlo sampling of solutions to inverse problems. *Journal of Geophysical Research: Solid Earth*, 100(B7):12431–12447, 1995.
- [34] T Vukicevic and D Posselt. Analysis of the impact of model nonlinearities in inverse problem solving. *Journal of the Atmospheric Sciences*, 65(9):2803–2823, 2008.