

Linking systems biology models to data: a stochastic kinetic model of p53 oscillations

D. A. Henderson*, R. J. Boys* C. J. Proctor[†] & D. J. Wilkinson*

January 5, 2009

Abstract

This chapter considers the assessment and refinement of a dynamic stochastic process model of the cellular response to DNA damage. The proposed model is a complex nonlinear continuous time latent stochastic process. It is compared to time course data on the levels of two key proteins involved in this response, captured at the level of individual cells in a human cancer cell line. The primary goal of this study is to “calibrate” the model by finding parameters of the model (kinetic rate constants) that are most consistent with the experimental data. Significant amounts of prior information are available for the model parameters. It is therefore most natural to consider a Bayesian analysis of the problem, using sophisticated MCMC methods to overcome the formidable computational challenges.

1 Introduction

1.1 Overview

Systems biology is an exciting new paradigm for life science research in the post-genomic era. It is a development of molecular biology in which the focus has moved from trying to understand the function of individual biomolecules (or pairs of biomolecules) to understanding how collections of biomolecules of varying types act together to accomplish the observed dynamic biological system behaviour. Systems biology involves a combination of mathematical modelling, biological experimentation and quantitative data generation. In particular, it crucially depends on the ability to adjust models in the light of experimental data. Further, there is now overwhelming evidence that intrinsic stochasticity is an important feature of intra-cellular processes. Statistical methods are therefore likely to play an increasingly important role in systems biology as models become more realistic and quantitative dynamic data becomes more routinely available (Wilkinson, 2009).

This chapter considers the assessment and refinement of a dynamic stochastic process model of the cellular response to DNA damage. The proposed model is compared to time course data on the levels of two key proteins involved in this response, captured at the level of individual cells in a human cancer cell line. The primary goal of this study is to “calibrate” the model by finding parameters of the model (kinetic rate constants) that are most consistent with the experimental data. The model is a complex nonlinear continuous time latent stochastic process model and so Markov chain Monte Carlo (MCMC) methods are a natural way to approach the inferential analysis from a computational perspective. In addition to being computationally difficult, the problem is also conceptually hard as the data-poor scenario means that some parameters of interest are only weakly

*School of Mathematics & Statistics, Newcastle University, Newcastle upon Tyne, NE1 7RU, UK

[†]Institute for Ageing & Health, Newcastle University, Newcastle upon Tyne, NE4 6BE, UK

identifiable. Fortunately, the mechanistic nature of the model means that all of the parameters are clearly interpretable, and a significant level of prior information is available for many of these from biological experts or the biological literature. The problem is therefore ideally suited to a Bayesian analysis using MCMC for computation.

The model concerns oscillations observed in the levels of two proteins in single living cancer cells subsequent to gamma irradiation. The two proteins, p53 and Mdm2, appear to oscillate out of phase with one another — system behaviour typically associated with some kind of negative feedback loop. This is consistent with current biological knowledge as p53 is known to enhance the production of Mdm2, and Mdm2 is known to inhibit p53. However, these oscillations are only observed at the single cell level and are not present in data derived from cell populations. It is therefore of interest to develop a simple mechanistic model, consistent with current biological knowledge, which explains the oscillatory behaviour and also explains why it is observed only at the single cell level (Proctor and Gray, 2008). Stochasticity is the key feature required to reconcile the apparent discrepancy between the single cell and population level data, with noisy oscillations being “averaged out” in the population level data. The stochastic process model contains several parameters whose values are uncertain. This chapter considers the problem of using time course data on levels of p53 and Mdm2 in several individual cells to improve our knowledge regarding plausible parameter values, and also to assess the extent to which the proposed stochastic model is consistent with the available data.

1.2 Biological background

The p53 tumour suppressor protein plays a major role in cancer as evidenced by the high incidence of TP53 gene mutations in human tumours (Hainaut and Hollstein, 2000). The TP53 gene encodes a transcription factor with target genes that are involved in DNA repair, cell cycle arrest and apoptosis. It has been described as the “guardian of the genome” (Lane, 1992), blocking cell cycle progression to allow the repair of damaged DNA. Under normal homeostatic conditions, the cellular levels of p53 protein are kept at a low level. There is basal transcription of the p53 gene (TP53) even in unstressed cells but the protein product does not accumulate as it has a short half-life of about 15-30 minutes (Finlay, 1993) and is usually bound to Mdm2, an ubiquitin E3 ligase, which targets p53 to the proteasome for degradation (Haupt *et al.*, 1997; Clegg *et al.*, 2008). Mdm2-binding prevents the transcriptional activity of p53 (Thut *et al.*, 1997), a phenomenon that is dependent on the catalytic activity of Mdm2 (Christophorou *et al.*, 2005). Mdm2 also has a short half-life and is a substrate of its own E3 ligase activity *in vitro* (Fang *et al.*, 2000). The transcription of Mdm2 is regulated by p53 (Barak *et al.*, 1993) and so under normal conditions, both p53 and Mdm2 are kept at low levels.

It is well known that stress induces an increase in levels of p53 which in turn leads to an increase in the transcription of Mdm2 (Mendrysa and Perry, 2000). One pathway for stabilization of p53 is via the kinase ATM, which is activated by DNA damage and phosphorylates p53 close to its Mdm2 binding site, so blocking its interaction with Mdm2 (Vogelstein *et al.*, 2000). In addition, ATM phosphorylates Mdm2 which not only interferes with its ability to bind to p53 but also enhances the degradation of Mdm2 (Pereg *et al.*, 2005; Khosravi *et al.*, 1999), providing an additional route for p53 stabilization. Another mechanism for the increase in p53 levels is the activation of ARF (known as p14ARF in humans), a nucleolar protein that senses DNA damage (Khan *et al.*, 2004). Although ARF responds to DNA damage, it is better known for its response to aberrant growth signals which are triggered by oncogenes (mutated forms of normal cellular genes which when activated can induce cancer). ARF binding enhances the degradation of Mdm2, resulting in p53 stabilisation (Khan *et al.*, 2004; Zhang *et al.*, 1998). Since an increase in p53 leads to an increase in Mdm2 transcription, and Mdm2 targets p53 for degradation, p53 levels are again inhibited,

providing a negative feedback loop.

Negative feedback loops have been found in several systems of interacting proteins (e.g. Hes1 in Notch signalling (Hirata *et al.*, 2002), NF- κ B signalling system (Nelson *et al.*, 2004)) and have attracted the attention of mathematical modellers. In particular, models have been produced to analyse the oscillations of p53 and Mdm2 in previously published single-cell fluorescent reporter assays (Ciliberto *et al.*, 2005; Geva-Zatorsky *et al.*, 2006; Lev Bar-Or *et al.*, 2000; Ma *et al.*, 2005; Tiana *et al.*, 2007; Zhang *et al.*, 2007). The single cell assays have been very informative, revealing that increasing DNA damage results in an increased number of oscillations, but not an increased magnitude in the response (Geva-Zatorsky *et al.*, 2006; Lahav *et al.*, 2004). The data also show that there is large intercellular variation with a fraction of cells showing no response or a slowly fluctuating signal. In the cells in which oscillations were detected, there was a wide fluctuation in the amplitude (about 70%) and smaller variations in the period of the peaks (about 20%) (Geva-Zatorsky *et al.*, 2006). The oscillations in these data showed a period of about 5.5 hours with a delay of about 2 hours between p53 and Mdm2 peaks (Geva-Zatorsky *et al.*, 2006).

All previous models to date have used a deterministic approach to analyse the oscillatory behaviour. These models have used differential equations and mathematical functions requiring a fairly large number of parameters with the generation of oscillations being very dependent on the range of parameter values chosen. Geva-Zatorsky *et al.* (2006) constructed six different models and found that the simplest model, which contained one intermediary and one negative feedback loop with a delay, was unable to produce multiple oscillations and that it was necessary to either introduce a positive feedback loop or a time delay term (see Figure 6 of Geva-Zatorsky *et al.* (2006)). However, these additions were not sufficient for robustness over a wide range of parameter values. The addition of a non-linear negative feedback loop, a linear positive feedback loop or a second negative feedback loop produced models that were able to demonstrate sustained oscillations over a wide range of parameters. As the models are deterministic, the outcome only depends on the initial conditions and so they cannot easily be used to investigate inter- and intra-cell variability. Geva-Zatorsky *et al.* (2006) incorporated some random noise in protein production in their models and found that the introduction of low-frequency noise resulted in variability in the amplitude of the oscillations as observed experimentally. Ma *et al.* (2005) also incorporated a stochastic component for the DNA damage component of their model which resulted in variability in the number of oscillations. However, for a simulated dose of 2.5Gy, they found that the majority of cells had only one peak and that a step input of DNA damage was required to obtain sustained oscillations.

We built a mechanistic model (Proctor and Gray, 2008) within a discrete stochastic chemical kinetic framework (Wilkinson, 2006), so that the intercellular variability could be accounted for in a natural way. Our approach meant that we did not need to include complex rate laws — mass action stochastic kinetics were assumed throughout — or any forced time delay terms.

1.3 Construction of the stochastic kinetic model

We assume that p53 production consists of two steps: transcription to form messenger RNA (p53_mRNA in the model) and then translation to form protein (p53). Under normal conditions p53 is usually bound to Mdm2 to form a complex, Mdm2-p53 (Mdm2_p53). Mdm2 targets p53 to the proteasome for degradation. We assume that p53 is only transcriptionally active when not bound to Mdm2, and so the production of Mdm2 mRNA (Mdm2_mRNA) is dependent on the pool of unbound p53. The synthesis of Mdm2 depends on the level of Mdm2 mRNA and so is also dependent on the level of unbound p53. Thus Mdm2 mRNA provides the intermediary link between p53 and Mdm2 to provide the necessary delay in the negative feedback loop. We also include degradation of Mdm2, Mdm2 mRNA and p53 mRNA. ATM is included in the model in two states: either inactive (ATM_I) or active (ATM_A). Initially all ATM is in its inactive state. After DNA

damage, ATM is activated and is then able to phosphorylate both p53 and Mdm2. Phosphorylated p53 and Mdm2 are presented in the model by the species **p53_P** and **Mdm2_P** respectively. We assume that the phosphorylated proteins are unable to bind to one another and so phosphorylated p53 is not degraded. However, phosphorylation of Mdm2 leads to its enhanced degradation and **p53_P** is transcriptionally active. We also include steps for de-phosphorylation. Further details of the model are given in Proctor and Gray (2008). To carry out a “virtual experiment”, whereby the cell is subject to irradiation, the species that represents DNA damage (**damDNA**) is set to a large value for the initial period of the simulation. Damaged DNA is repaired at a rate determined by the parameter *krepair*.

The model was encoded using SBML-shorthand (Wilkinson, 2006) and then converted into the Systems Biology Markup Language (SBML) (Hucka *et al.*, 2003). SBML is a well-known modelling standard, allowing models to be shared in a form that other researchers can use in different hardware and software environments.

1.4 The experiment and data

The experiments were carried out in the laboratory of Uri Alon (Department of Molecular Cell Biology, Weizmann Institute of Science, Israel). Full details of the experimental procedure can be found in Geva-Zatorsky *et al.* (2006). The cell line used was MCF7, which are human breast cancer epithelial cells. They used a clone (all cells genetically identical) which was stably transfected with p53 fused to cyan fluorescent protein (CFP) and Mdm2 fused to yellow fluorescent protein (YFP). They irradiated the cells with different doses of gamma irradiation and obtained time-lapse fluorescence microscopy movies of the cells over time periods of about 30 hours. Images were captured every 10-20 minutes. Overall they collected data for 1000 individual cells in different experiments with different doses of irradiation. We were supplied with raw data for the cells which were irradiated at a dose of 2.5 Gy (141 cells) and at 5 Gy (146 cells). The units for the raw data are relative fluorescence units and the values for Mdm2 had been normalized (by multiplying the YFP fluorescence by a constant) so that the values for p53 and Mdm2 lie in the same range; see, for example, the time-course data for the seven cells in the second movie displayed in Figure 1.

2 Stochastic kinetic model

The stochastic kinetic model describes the evolution of $k = 10$ species by using 19 reactions. Each reaction occurs at a rate governed by the numbers of molecules of the reacting species and associated

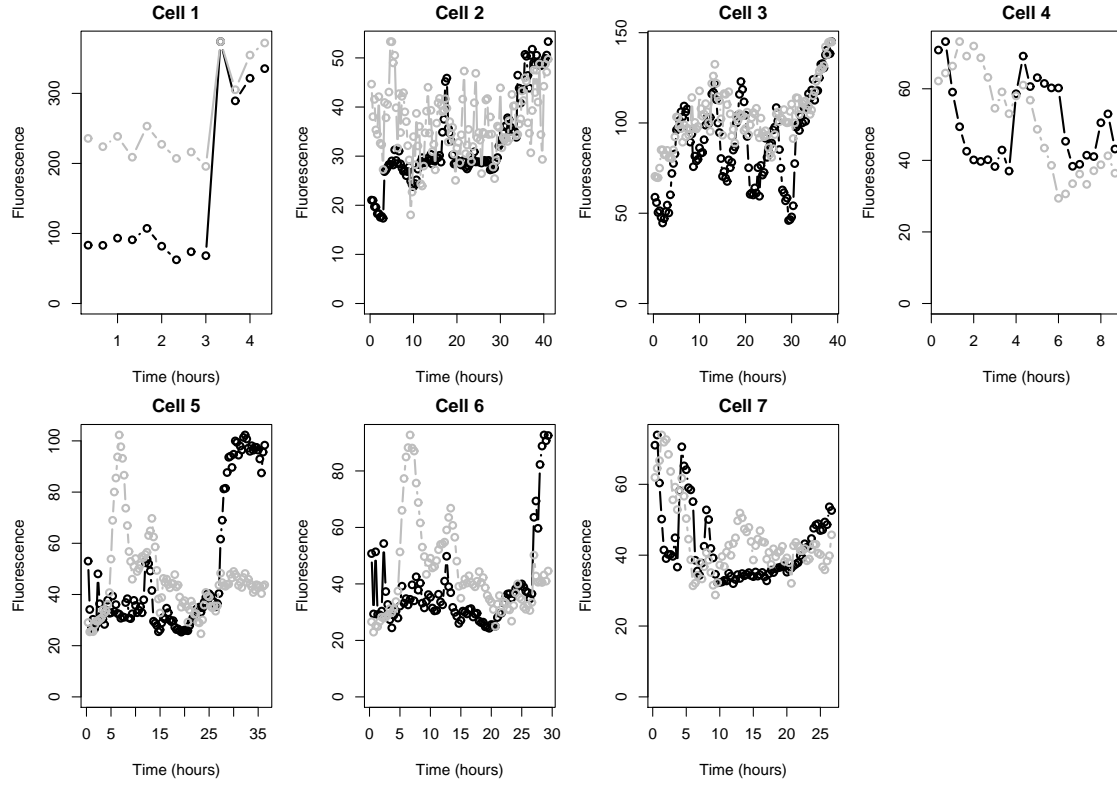


Figure 1: Measured (normalized) fluorescence levels for the seven cells in movie 2; p53 (black) and Mdm2 (grey).

rate constants. The list of reactions is



It will be convenient to work with the rate constants (*eg. ksynp53mRNA*) on a log scale and so we denote the collection of the $r = 19$ “calibration” parameters by $\theta = (\theta_1, \theta_2, \dots, \theta_r)'$, where θ_j is the log of the j th rate constant.

Let $Y_t = (Y_{t,1}, \dots, Y_{t,k})$ denote the state of the system at time t , where $Y_{t,j}$ is the number of molecules of species j at time t . The $k = 10$ species (and their corresponding index used in the notation) are listed in Table 1. Note however that these 10 species are not linearly independent due to the presence of a conservation law in the system that makes species 6 and 7 linearly related. Such conservation laws need to be preserved by inference algorithms — an example of how this is achieved is presented in Section 5. Also let $Y = (Y_{t_0}, Y_{t_1}, \dots, Y_{t_n})$ denote the state of the system at the time points (t_0, t_1, \dots, t_n) . The kinetic model is a Markov jump process and so the joint probability of Y factorises as

$$p(Y|\theta) = p(Y_{t_0}|\theta) \prod_{i=1}^n p(Y_{t_i}|Y_{t_{i-1}}, \theta),$$

where we assume that the initial state Y_{t_0} is independent of the reaction constants θ , that is, $p(Y_{t_0}|\theta) = p(Y_{t_0})$.

Given full information on the process, that is, the times and types of each reaction that take place, closed form expressions can be found for the conditional probabilities $p(Y_{t_i}|Y_{t_{i-1}}, \theta)$, and hence for the joint probability $p(Y|\theta)$. However, in the present application (as with many other practical scenarios) experimental techniques do not provide this full information, perhaps only giving the levels of some species at a limited number of time points. Here different strategies are required for analysing such partial information; see Boys *et al.* (2008) for details. The strategy we employ in this chapter is based on the fact that, for given reaction constants θ , it is possible to forward simulate realisations Y_t of the model exactly using, for example, the Gillespie algorithm (Gillespie, 1977).

3 Data

Suppose that the data available consist of time-course information on C cells. Specifically, the data on cell i are the scaled fluorescence measurements of two quantities, p53 and Mdm2, measured at n_i time points t_1, \dots, t_{n_i} . We will sometimes refer to the p53 and Mdm2 measurements by the colours of the fluorescent proteins used, namely the “cyan channel” and the “yellow channel” respectively. In these data, measurements are taken every $\tau = 1200$ seconds, that is, at times $t_j = j\tau$, $j = 0, 1, 2, \dots$ and so we simplify the notation by referring to time by its index, $j = 0, 1, 2, \dots$. Also

Species index	Species name
1	Mdm2
2	p53
3	Mdm2_p53
4	Mdm2_mRNA
5	p53_mRNA
6	ATMA
7	ATMI
8	p53_P
9	Mdm2_P
10	damDNA

Table 1: Species names and their corresponding index.

the measurement pair at the t th time point in the i th cell is denoted by $z_t^i = (z_{t,c}^i, z_{t,y}^i)'$, where the subscripts c and y refer to the channel colours.

The data on cell i is denoted by $z^i \equiv z_{1:n_i}^i = \{z_1^i, \dots, z_{n_i}^i\}$, and the full dataset for all C cells is denoted by $z = \{z^1, \dots, z^C\}$. Finally, an additional complication in these data is that the observed measurements of Mdm2 (yellow) have been scaled so that the maximum value is the same as the maximum observed value of p53 (cyan). This post-processing step adds an extra layer of complexity into the modelling task.

4 Linking the model to the data

4.1 Modelling the raw measurements

Let $Y_{t,c}^i$ and $Y_{t,y}^i$ denote the total amount of p53 and Mdm2, respectively, in the i th cell at the t th time point. Each of these amounts is the sum of three species counts (at time point t), namely

$$Y_{t,c}^i = Y_{t,2}^i + Y_{t,3}^i + Y_{t,8}^i \quad \text{and} \quad Y_{t,y}^i = Y_{t,1}^i + Y_{t,3}^i + Y_{t,9}^i.$$

Note that $Y_{t,3}^i$ is common to both $Y_{t,c}^i$ and $Y_{t,y}^i$. Let $Y_c^i = (Y_{1,c}^i, \dots, Y_{n_i,c}^i)$ and $Y_y^i = (Y_{1,y}^i, \dots, Y_{n_i,y}^i)$ denote the amounts (by channel) for cell i , and let the amounts over all cells be denoted by $Y_c = (Y_c^1, \dots, Y_c^C)$, $Y_y = (Y_y^1, \dots, Y_y^C)$ and $Y = \{Y_c, Y_y\}$.

The true fluorescence levels in the cyan and yellow channels are assumed to be directly proportional to the numbers of molecules of p53 and Mdm2, respectively, with unknown proportionality constants β_j , $j \in \{c, y\}$, that is

$$\gamma_{t,c}^i = \beta_c Y_{t,c}^i \quad \text{and} \quad \gamma_{t,y}^i = \beta_y Y_{t,y}^i.$$

The raw measurements of these true fluorescence levels (after adjusting for background noise), denoted by $\zeta_{t,c}^i$ and $\zeta_{t,y}^i$, are assumed to be independent and normally distributed with means $\gamma_{t,c}^i$ and $\gamma_{t,y}^i$ respectively. The measurement processes in the cyan and yellow channels are sufficiently similar that we will assume a common precision ϕ to these processes. Thus

$$\zeta_{t,c}^i | Y_{t,c}^i, \beta_c, \phi \sim N(\beta_c Y_{t,c}^i, \phi^{-1}) \quad \text{and} \quad \zeta_{t,y}^i | Y_{t,y}^i, \beta_y, \phi \sim N(\beta_y Y_{t,y}^i, \phi^{-1}).$$

Prior beliefs about ϕ are modelled through a gamma $Ga(a_\phi, b_\phi)$, distribution, with density

$$p(\phi) = \frac{b_\phi^{a_\phi} \phi^{a_\phi-1} \exp(-b_\phi \phi)}{\Gamma(a_\phi)},$$

where $\Gamma(\cdot)$ denotes the gamma function. We take $a_\phi = 2$ and $b_\phi = 50$ and this reflects fairly strong beliefs that the prior precision is close to its mean of $E(\phi) = 1/25$.

4.2 Modelling the scaling process

Beliefs about the scaling constants β_c and β_y are modelled via independent normal distributions

$$\beta_c \sim N(a_{\beta_c}, b_{\beta_c}^2) \quad \text{and} \quad \beta_y \sim N(a_{\beta_y}, b_{\beta_y}^2).$$

Prior means of $a_{\beta_c} = a_{\beta_y} = 1$, representing no scaling, were adopted. The prior standard deviations b_{β_c} and b_{β_y} were chosen to be $1/3$, so that the central 95% prior probability interval for β_c is approximately (0.347, 1.653). Strictly speaking, these prior distributions are not consistent with the requirement that the scaling constants should take positive values. However, there is a significant benefit in choosing this form for the prior distributions, as we shall see shortly. Also, for these prior distributions, the probability of the scaling constants taking negative values is negligibly small.

4.3 Marginal model for the data

The scaling constants β_c and β_y are essentially nuisance parameters and their values are not of direct interest. Our choice of normal distributions for β_c and β_y allows us to marginalise them analytically from the normal models for the raw measurements, and this may result in computational benefits. This gives the marginal distributions of the raw measurements as

$$\zeta_{t,c}^i | Y_{t,c}^i, \phi \sim N(a_{\beta_c} Y_{t,c}^i, b_{\beta_c}^2 (Y_{t,c}^i)^2 + \phi^{-1}) \quad \text{and} \quad \zeta_{t,y}^i | Y_{t,y}^i, \phi \sim N(a_{\beta_y} Y_{t,y}^i, b_{\beta_y}^2 (Y_{t,y}^i)^2 + \phi^{-1}).$$

Note that the variance of the marginal measurement error distribution now depends on the signal, that is, on the true numbers of molecules.

Unfortunately the available data are not simply a collection of raw measurements: the data have been normalized so that the values for the two channels lie on the same scale. However, only data recorded for the yellow channel (Mdm2) is affected. The data recorded for the cyan channel (p53), $z_{t,c}^i$, are the raw measurements for that channel and so $z_{t,c}^i = \zeta_{t,c}^i$. However, the data recorded for the yellow channel (Mdm2), $z_{t,y}^i$, are normalized versions of the raw measurements. The normalized measurement for Mdm2 in the i th cell at time t , $z_{t,y}^i$, is obtained from the raw measurement $\zeta_{t,y}^i$ by dividing by the maximum raw Mdm2 measurement in the i th cell, $(\zeta_y^i)^{\max} = \max(\zeta_{1,y}^i, \dots, \zeta_{n_i,y}^i)$, and then multiplying by the maximum measurement in the cyan channel, $(z_c^i)^{\max} = \max(z_{1,c}^i, \dots, z_{n_i,c}^i)$. Thus

$$z_{t,y}^i = s(\zeta_y^i, z_c^i) \equiv \zeta_{t,y}^i \frac{(z_c^i)^{\max}}{(\zeta_y^i)^{\max}}.$$

The probability of the observed scaled measurements z_y^i is

$$p(z_y^i | \zeta_y^i, z_c^i) = \begin{cases} 1, & \text{if } s(\zeta_y^i, z_c^i) = z_y^i, \\ 0, & \text{otherwise,} \end{cases}$$

and this depends on z_c^i only through $(z_c^i)^{\max}$. The joint density of the observed data (the likelihood function) is therefore

$$\begin{aligned} p(z_c, z_y | \zeta_y, Y_c, \phi) &= p(z_c | Y_c, \phi) p(z_y | \zeta_y, z_c) \\ &= p(z_c | Y_c, \phi) \prod_{i=1}^C p(z_y^i | \zeta_y^i, z_c^i), \end{aligned}$$

where $p(z_c | Y_c, \phi)$, the joint density of the raw p53 measurements, is

$$\begin{aligned} p(z_c | Y_c, \phi) &= \prod_{i=1}^C \prod_{t=1}^{n_i} p(z_{t,c}^i | Y_{t,c}^i, \phi) \\ &= \prod_{i=1}^C \prod_{t=1}^{n_i} (2\pi)^{-1/2} \{b_{\beta_c}^2 (Y_{t,c}^i)^2 + \phi^{-1}\}^{-1/2} \exp \left\{ -\frac{(z_{t,c}^i - a_{\beta_c} Y_{t,c}^i)^2}{2\{b_{\beta_c}^2 (Y_{t,c}^i)^2 + \phi^{-1}\}} \right\}. \end{aligned}$$

When constructing the posterior distribution (given by Equation (2) in Section 5) we will also need the joint density of the raw Mdm2 measurements:

$$\begin{aligned} p(\zeta_y | Y_y, \phi) &= \prod_{i=1}^C \prod_{t=1}^{n_i} p(\zeta_{t,y}^i | Y_{t,y}^i, \phi) \\ &= \prod_{i=1}^C \prod_{t=1}^{n_i} (2\pi)^{-1/2} \{b_{\beta_y}^2 (Y_{t,y}^i)^2 + \phi^{-1}\}^{-1/2} \exp \left\{ -\frac{(\zeta_{t,y}^i - a_{\beta_y} Y_{t,y}^i)^2}{2\{b_{\beta_y}^2 (Y_{t,y}^i)^2 + \phi^{-1}\}} \right\}. \end{aligned}$$

Parameter name	Value
<i>ksynMdm2</i>	-7.611
<i>kdegMdm2</i>	-7.745
<i>ksynp53</i>	-5.116
<i>kdegp53</i>	-7.100
<i>kbinMdm2p53</i>	-6.764
<i>krelMdm2p53</i>	-11.369
<i>ksynMdm2mRNA</i>	-9.210
<i>kdegMdm2mRNA</i>	-9.210
<i>kdegATMMdm2</i>	-7.824
<i>kproteff</i>	0.000
<i>ksynp53mRNA</i>	-6.908
<i>kdegp53mRNA</i>	-9.210

Table 2: Values (on a log scale) for known model calibration parameters.

Index	Parameter name	Lower limit	Upper limit
		a_{θ_i}	b_{θ_i}
1	<i>kactATM</i>	-18.210	-0.210
2	<i>kinactATM</i>	-16.601	1.399
3	<i>kphosp53</i>	-13.601	-1.601
4	<i>kdephosp53</i>	-8.996	0.702
5	<i>kphosMdm2</i>	-7.609	2.088
6	<i>kdephosMdm2</i>	-8.996	0.702
7	<i>krepair</i>	-13.820	-7.820

Table 3: The lower and upper limits of the uniform prior distributions for the unknown model calibration parameters (on a log scale).

4.4 Prior specification for the model calibration parameters

Information on likely values for these parameters can be found in Proctor and Gray (2008). We found that the parameters naturally grouped into four classes, ranging from those that were known fairly accurately to those with a fair amount of uncertainty. We have fixed the parameters that were known fairly accurately to their suggested values and these are given (on a log scale) in Table 2. This reduces the complexity of the analysis to finding plausible values for the remaining $r^* = 7$ parameters. We have taken independent uniform prior distributions for these calibration parameters (the logged kinetic rate constants) θ , with

$$\theta_i | a_{\theta_i}, b_{\theta_i} \sim U(a_{\theta_i}, b_{\theta_i}), \quad i = 1, \dots, r^*$$

and used the information in Proctor and Gray (2008) to determine reasonable values for the upper and lower limits of these distributions; see Table 3.

Therefore, suppressing dependence on the fixed hyperparameters a_{θ_i} and b_{θ_i} , the joint prior density is

$$p(\theta) = \prod_{i=1}^{r^*} p(\theta_i) = \prod_{i=1}^{r^*} (b_{\theta_i} - a_{\theta_i})^{-1}.$$

4.5 Prior specification for the initial species counts

The joint probability of the states $Y|\theta$ is given by the product

$$p(Y|\theta) = \prod_{i=1}^C \left\{ p(Y_0^i) \prod_{t=1}^{n_i} p(Y_t^i | Y_{t-1}^i, \theta) \right\}, \quad (1)$$

in which the conditional probabilities are determined by the dynamics of the stochastic kinetic model. We model the unobserved initial state of the system in cell i using independent Poisson distributions

$$Y_{0,j}^i | \lambda_j \sim Po(\lambda_j), \quad j = 1, \dots, 10.$$

Note that this imposes the same prior distribution for the initial state in each of the C cells. Therefore the marginal probability of the initial state of the system for cell i is

$$p(Y_0^i) = \prod_{j=1}^{10} p(Y_{0,j}^i) = \prod_{j=1}^{10} \lambda_j^{Y_{0,j}^i} e^{-\lambda_j} / \Gamma(Y_{0,j}^i + 1).$$

The means of these Poisson priors (λ_j) have been chosen using information in the literature (*e.g.* Proctor and Gray (2008)) on the most likely initial state of the system: $(\lambda_1, \dots, \lambda_{10}) = (6, 6, 96, 11, 11, 1, 201, 1, 1, 76)$.

5 Posterior computation

Inferences about the values of the unknown quantities in the model are based on their joint posterior distribution, which has density proportional to the product of the joint prior density and the likelihood, that is

$$p(\theta, \phi, Y, \zeta_y | z_c, z_y) \propto p(\theta) p(\phi) p(Y|\theta) p(\zeta_y | Y, \phi) p(z_c | Y, \phi) p(z_y | \zeta_y, z_c), \quad (2)$$

where all the terms on the right-hand side of (2) have been described previously.

Posterior computation is made difficult by the intractability of the conditional probabilities $p(Y_t^i | Y_{t-1}^i, \theta)$ from the stochastic kinetic model, which enter the Bayesian model through Equation (1). It is possible to avoid computation of the $p(Y_t^i | Y_{t-1}^i, \theta)$ by constructing an algorithm which uses realisations from the stochastic kinetic model; see, for example, Henderson *et al.* (2009). However, such algorithms require the generation of many model realisations and so, for this approach to work well, each model realisation must be quick to simulate. Unfortunately, this is not the case for the stochastic kinetic model considered here. For example, simulating $Y_t^i | Y_{t-1}^i, \theta$ for some values of Y_{t-1}^i and θ takes a matter of milliseconds, yet for other values it can take several seconds, even on a reasonably powerful computer (2.2GHz, 8GB RAM) using an efficient C implementation of Gillespie's exact discrete event simulation algorithm.

One way of dealing with the slow simulation speed of the Gillespie algorithm is to use a fast approximate simulation algorithm, such as the τ -leap method; see Wilkinson (2006). Here there is a trade-off between simulation speed (and therefore speed of the inference procedure) and the exactness of the inferences made. Initial investigations with some fast approximate algorithms revealed that none were able to provide the sort of improvements in simulation speed needed for this analysis.

An alternative to using an approximate simulation algorithm is to use an exact simulation algorithm for an approximation to the stochastic kinetic model. Various authors have sought solutions along

these lines. For example, Golightly and Wilkinson (2005) use an approximation based on the chemical Langevin equation (CLE; Gillespie, 2000), namely, the diffusion process that approximates most closely the Markov jump process defined by the stochastic kinetic model. These authors build on this work (in Golightly and Wilkinson (2006b)) and demonstrate how a combination of particle filtering and MCMC methods can be used to sample from the posterior distribution of the rate constants given the CLE approximation model and data observed partially, discretely and with error.

A different strategy, and one we adopt in this chapter, is to *emulate* the stochastic kinetic model. We do this by constructing a tractable approximation to the conditional probability distribution of the stochastic kinetic model, $p(Y(t + \tau)|Y(t), \theta)$, where $Y(t)$ denotes the state of the system at any particular time t (seconds), and $Y(t + \tau)$ denotes the state of the system τ seconds in the future. For the data in this application, the time step is $\tau = 1200$ seconds. The approximate probability distribution, which we refer to as the emulator, is denoted by $p^*(\cdot|\cdot, \theta)$, and the objective is to use it in place of the probability distribution $p(\cdot|\cdot, \theta)$ in Equation (1). The construction of emulators is commonplace in the computer models literature where complex deterministic functions are modelled via tractable stochastic processes; see Kennedy and O’Hagan (2001), O’Hagan (2006), Santner *et al.* (2003), and references therein. In this application, the function to be emulated is discrete, multivariate and stochastic. The emulator is constructed by fitting simple statistical models to output obtained by simulating the stochastic kinetic model for τ seconds from a designed collection of values of the model inputs $\{Y(t), \theta\}$. The approach to fitting the emulator that we follow is described in more detail in Appendix B.

Based on the emulator, the joint probability of the states is

$$p^*(Y|\theta) = \prod_{i=1}^C \left\{ p(Y_0^i) \prod_{t=1}^{n_i} p^*(Y_t^i | Y_{t-1}^i, \theta) \right\},$$

and this replaces the exact probability $p(Y|\theta)$ in (2) to give an expression for the posterior density based on this approximation, namely

$$p^*(\theta, \phi, Y, \zeta_y | z_c, z_y) \propto p(\theta)p(\phi)p^*(Y|\theta)p(\zeta_y|Y, \phi)p(z_c|Y, \phi)p(z_y|\zeta_y, z_c). \quad (3)$$

Here the superscript $*$ distinguishes probabilities (or densities) based on the emulator approximation rather than the true stochastic kinetic model. Also note that all terms in the right-hand side of (3) are tractable and so this formulation leads to a workable solution.

Sampling from the posterior distribution (3) is possible using a Metropolis-Hastings (MH) within Gibbs MCMC scheme. In the scheme, we update each set of unknown quantities in turn from their full conditional distribution by using a MH step if the full conditional distribution cannot be sampled from directly. The MCMC algorithm is constructed so that the distribution of sampled values tends to the posterior distribution as the number of iterations increases. The sampled values are then used to approximate features of the posterior distribution. For computational reasons, we find it beneficial to work with transformed values of the calibration parameters θ . Specifically we work with transformed parameters $\Lambda = (\Lambda_1, \dots, \Lambda_{r^*})$, where $\Lambda_i = \log(\theta_i - a_{\theta_i}) - \log(b_{\theta_i} - \theta_i)$. This corresponds to a logit transformation of θ_i after it has been re-scaled to lie on the unit interval. It follows that the Λ_i have independent (standard) logistic distributions, and therefore that the joint density of Λ is

$$p(\Lambda) = \prod_{i=1}^{r^*} p(\Lambda_i) = \prod_{i=1}^{r^*} \frac{\exp(\Lambda_i)}{\{1 + \exp(\Lambda_i)\}^2}.$$

The sampled values of Λ_i are simply back transformed to give sampled values of θ_i .

In outline, one iteration of the MCMC scheme entails the following steps:

- Update $\Lambda|\dots$ by using a MH step with a symmetric multivariate normal random walk proposal, centred at the current sampled value. The acceptance ratio for the proposed move from Λ to $\tilde{\Lambda}$ is

$$A_{\Lambda} = \frac{p(\tilde{\Lambda}) p^*(Y|\tilde{\theta})}{p(\Lambda) p^*(Y|\theta)},$$

where, for instance, $\theta_i = \{a_{\theta_i} + b_{\theta_i} \exp(\Lambda_i)\} / \{1 + \exp(\Lambda_i)\}$.

- Update $\phi|\dots$ by proposing a new value $\tilde{\phi}$ from a proposal distribution with density $q(\tilde{\phi}|\phi)$. The acceptance ratio for the proposed move from ϕ to $\tilde{\phi}$ is

$$A_{\phi} = \frac{p(\tilde{\phi}) p(\zeta_Y|Y, \tilde{\phi}) p(z_c|Y, \tilde{\phi}) q(\phi|\tilde{\phi})}{p(\phi) p(\zeta_Y|Y, \phi) p(z_c|Y, \phi) q(\tilde{\phi}|\phi)}.$$

- Update $Y_{\mathcal{S} \setminus \{6,7\}}|\dots$, where $\mathcal{S} = \{1, 2, \dots, 10\}$ as follows. Here $Y_{\mathcal{S} \setminus \{6,7\}}$ denotes the values of the states (in each cell) excluding those for species 6 (ATMA) and 7 (ATMI). Species 6 and 7 are treated separately as their sum is fixed throughout the time course.

For cells $i = 1, \dots, C$:

- for $t = 0$, propose independent Poisson candidate values $\tilde{Y}_{t,j}^i | Y_{t,j}^i \sim Po(Y_{t,j}^i + a_Y)$ for $j \in \mathcal{S} \setminus \{6, 7\}$, where $a_Y > 0$ is a positive tuning constant which is chosen to be small. Denote the proposal probability by $q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)$. The acceptance ratio for the proposed move from $Y_{t,\mathcal{S} \setminus \{6,7\}}^i$ to $\tilde{Y}_{t,\mathcal{S} \setminus \{6,7\}}^i$ is

$$A_{Y_{0,\mathcal{S} \setminus \{6,7\}}^i} = \frac{\prod_{j \in \mathcal{S} \setminus \{6,7\}} p(\tilde{Y}_{t,j}^i) p^*(Y_{t+1}^i | \tilde{Y}_t^i, \theta) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(Y_{t,j}^i | \tilde{Y}_{t,j}^i)}{\prod_{j \in \mathcal{S} \setminus \{6,7\}} p(Y_{t,j}^i) p^*(Y_{t+1}^i | Y_t^i, \theta) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)},$$

where \tilde{Y}_t^i denotes the vector of proposed candidate values together with the current values for species 6 and 7.

- For $t = 1, \dots, n-1$, propose independent Poisson candidate values $\tilde{Y}_{t,j}^i | Y_{t,j}^i \sim Po(Y_{t,j}^i + a_Y)$ for $j \in \mathcal{S} \setminus \{6, 7\}$. Denote the proposal probability by $q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)$. The acceptance ratio for the proposed move from $Y_{t,\mathcal{S} \setminus \{6,7\}}^i$ to $\tilde{Y}_{t,\mathcal{S} \setminus \{6,7\}}^i$ is

$$A_{Y_{t,\mathcal{S} \setminus \{6,7\}}^i} = \frac{p^*(\tilde{Y}_t^i | Y_{t-1}^i, \theta) p^*(Y_{t+1}^i | \tilde{Y}_t^i, \theta) p(z_{t,c}^i | \tilde{Y}_t^i, \phi) p(\zeta_{t,y}^i | \tilde{Y}_t^i, \phi) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(Y_{t,j}^i | \tilde{Y}_{t,j}^i)}{p^*(Y_t^i | Y_{t-1}^i, \theta) p^*(Y_{t+1}^i | Y_t^i, \theta) p(z_{t,c}^i | Y_t^i, \phi) p(\zeta_{t,y}^i | Y_t^i, \phi) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)}.$$

- For $t = n$, propose independent Poisson candidate values $\tilde{Y}_{t,j}^i | Y_{t,j}^i \sim Po(Y_{t,j}^i + a_Y)$ for $j \in \mathcal{S} \setminus \{6, 7\}$. Denote the proposal probability by $q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)$. The acceptance ratio for the proposed move from $Y_{t,\mathcal{S} \setminus \{6,7\}}^i$ to $\tilde{Y}_{t,\mathcal{S} \setminus \{6,7\}}^i$ is

$$A_{Y_{n,\mathcal{S} \setminus \{6,7\}}^i} = \frac{p^*(\tilde{Y}_t^i | Y_{t-1}^i, \theta) p(z_{t,c}^i | \tilde{Y}_t^i, \phi) p(\zeta_{t,y}^i | \tilde{Y}_t^i, \phi) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(Y_{n,j}^i | \tilde{Y}_{t,j}^i)}{p^*(Y_t^i | Y_{t-1}^i, \theta) p(z_{t,c}^i | Y_t^i, \phi) p(\zeta_{t,y}^i | Y_t^i, \phi) \prod_{j \in \mathcal{S} \setminus \{6,7\}} q(\tilde{Y}_{t,j}^i | Y_{t,j}^i)}.$$

- For cells $i = 1, \dots, C$, update $Y_{t,\{6,7\}}^i|\dots$ for $t = 0, 1, \dots, n_i$ as follows. First, for $t = 0$ and $j = 6, 7$, propose independent Poisson candidate values $\tilde{Y}_{t,j}^i | Y_{t,j}^i \sim Po(Y_{t,j}^i + a_Y)$. Denote these proposal probabilities by $q(\tilde{Y}_{0,j}^i | Y_{0,j}^i)$. This gives a proposal for the new sum of the two species in the i th cell, $\tilde{N}^i = \tilde{Y}_{t,6}^i + \tilde{Y}_{t,7}^i$. Then for $t = 1, \dots, n_i$ propose

$$\tilde{Y}_{t,6}^i | Y_{t,6}^i, Y_{t,7}^i, \tilde{N}^i \sim Bin(\tilde{N}^i, (Y_{t,6}^i + a_{Y_6}) / (Y_{t,6}^i + Y_{t,7}^i + a_{Y_6}))$$

where $a_{Y_6} > 0$ is a positive tuning constant which is chosen to be small. Then set $\tilde{Y}_{t,7}^i = \tilde{N}^i - \tilde{Y}_{t,6}^i$. Denote these proposal probabilities by $q(\tilde{Y}_{t,6}^i | Y_{t,6}^i)$. The acceptance ratio for the proposed move from $Y_{\{6,7\}}^i$ to $\tilde{Y}_{\{6,7\}}^i$ is

$$A_{Y_{\{6,7\}}^i} = \prod_{j \in \{6,7\}} \left\{ \frac{p(\tilde{Y}_{0,j}^i) q(Y_{0,j}^i | \tilde{Y}_{0,j}^i)}{p(Y_{0,j}^i) q(\tilde{Y}_{0,6}^i | Y_{0,6}^i)} \right\} \prod_{t=1}^{n_i} \left\{ \frac{p^*(Y_t^i | \tilde{Y}_{t-1}^i, \theta) q(Y_{t,6}^i | \tilde{Y}_{t,6}^i)}{p^*(Y_t^i | Y_{t-1}^i, \theta) q(\tilde{Y}_{t,6}^i | Y_{t,6}^i)} \right\}.$$

Note that there is no contribution from the data in the above acceptance ratio since species 6 and 7 do not contribute to the total amounts of p53 or Mdm2.

- Update $\zeta_y | \dots$ by proposing a candidate $\tilde{\zeta}_y^{\max}$ using a symmetric multivariate normal random walk centred at the current value $\zeta_y^{\max} = \{(\zeta_y^1)^{\max}, \dots, (\zeta_y^C)^{\max}\}$, the density of which is denoted $q(\tilde{\zeta}_y^{\max} | \zeta_y^{\max})$. Then, for each cell i , set

$$\tilde{\zeta}_{t,y}^i = z_{t,y}^i \frac{(\tilde{\zeta}_y^i)^{\max}}{(z_c^i)^{\max}}, \quad t = 1, \dots, n_i.$$

The acceptance probability for the proposed move from ζ_y to $\tilde{\zeta}_y$ is

$$A_{\zeta_y} = \frac{p(\tilde{\zeta}_y | Y, \phi) p(z_y | \tilde{\zeta}_y, z_c) q(\zeta_y | \tilde{\zeta}_y)}{p(\zeta_y | Y, \phi) p(z_y | \zeta_y, z_c) q(\tilde{\zeta}_y | \zeta_y)}.$$

The proposal ratio $q(\zeta_y | \tilde{\zeta}_y) / q(\tilde{\zeta}_y | \zeta_y) = 1$ as the proposal is symmetric. The term $p(z_y | \tilde{\zeta}_y, z_c)$ checks the validity of the proposal for ζ_y , in the sense that it is compatible with the observed z_y . Because of the form of the proposal, $\tilde{\zeta}_y$ is always compatible with z_y as we use z_y as part of the proposal. Therefore the acceptance ratio reduces to

$$A_{\zeta_y} = \frac{p(\tilde{\zeta}_y | Y, \phi)}{p(\zeta_y | Y, \phi)}.$$

6 Inference based on single cell data

We begin our analysis by studying the time-course information in a single cell. This cell has been chosen at random and is the third cell from the second movie. The time-course covers 38 hours and 40 minutes with data sampled every 1200 seconds, giving 116 time points; see Figure 1.

Several independent Markov chains were simulated from different starting points using the MCMC algorithm outlined in the previous section (with $C = 1$). We report here the results of one of these chains. The MCMC algorithm was run for 250,000 iterations after discarding an initial 250,000 iterates as burn-in. The output was then thinned to remove some of the high autocorrelation by taking every 25th iterate, leaving 10,000 sampled values from the posterior distribution on which to base inferences.

Figure 2 shows density histograms of the MCMC output for the model calibration parameters. Recall that each of these parameters has been given a constant (uniform) prior density. Estimates of posterior means and standard deviations are given in Table 4. Clearly, uncertainty regarding all seven model calibration parameters has reduced and the ranges of plausible values have narrowed significantly. An image plot representing the posterior correlations between pairs of calibration parameters is displayed in Figure 3. It shows that several pairs of parameters are highly correlated and, in particular, there is a strong positive correlation between *kactATM* and *kinactATM* and a strong negative correlation between *kdephosp53* and *kphosMdm2*.

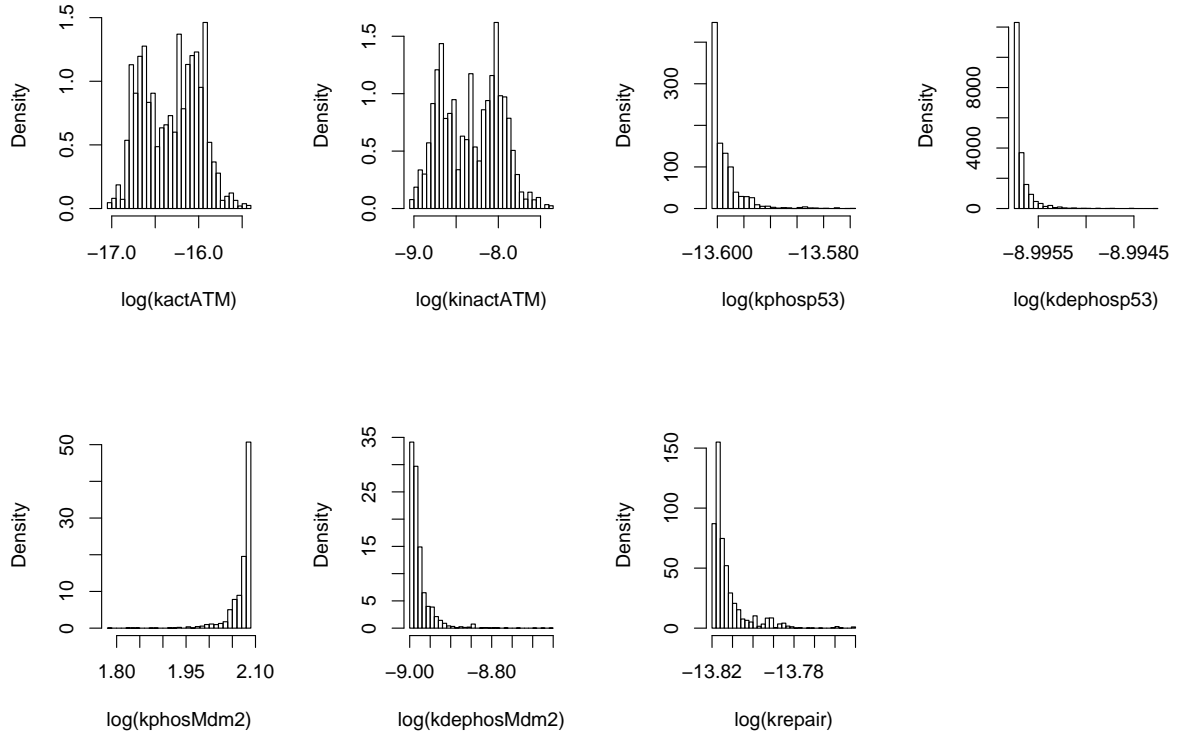


Figure 2: Model calibration parameters: marginal posterior densities.

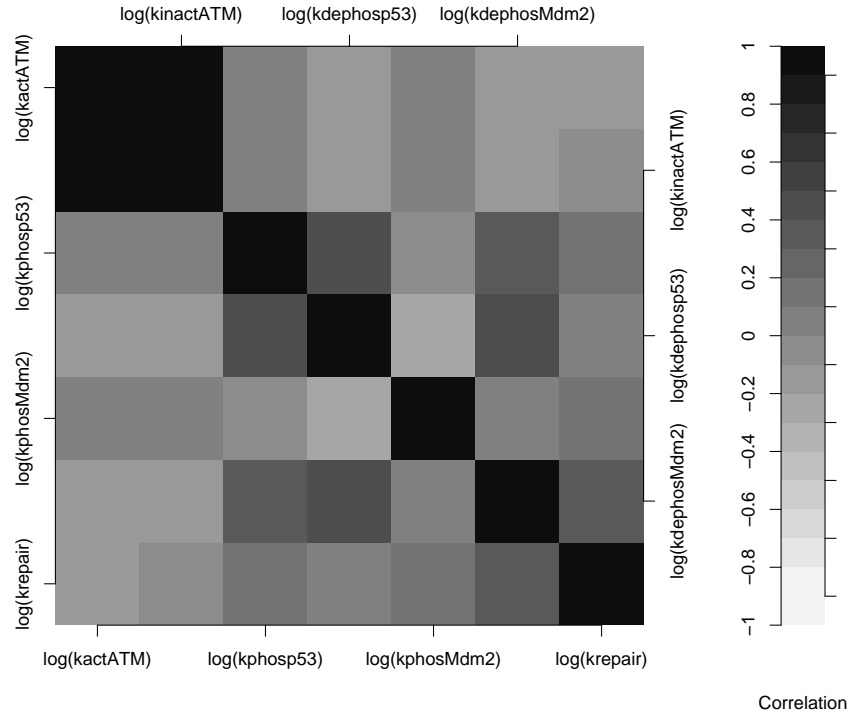


Figure 3: Model calibration parameters: posterior correlations.

Index	Parameter name	Mean	Standard deviation
1	$kactATM$	-16.314	0.3269
2	$kinactATM$	-8.320	0.3496
3	$kphosp53$	-13.599	0.0030
4	$kdephosp53$	-8.996	0.0001
5	$kphosMdm2$	2.071	0.0255
6	$kdephosMdm2$	-8.977	0.0256
7	$krepair$	-13.813	0.0087

Table 4: Model calibration parameters: posterior means and standard deviations.

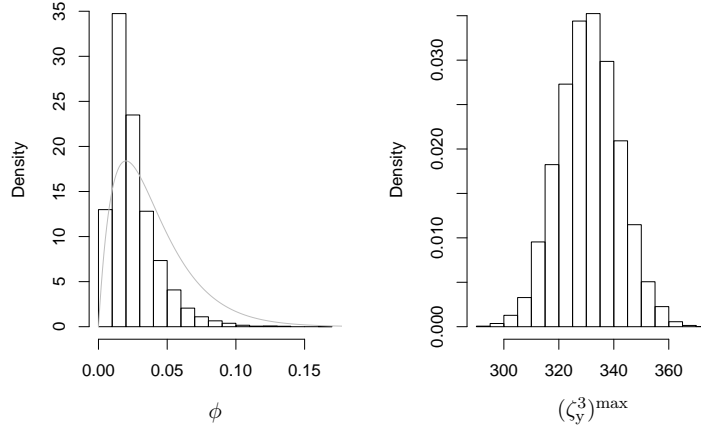


Figure 4: Marginal posterior densities of the measurement error precision ϕ (left panel) and the true maximum value in the yellow (Mdm2) channel for cell 3, $(\zeta_y^3)^{\max}$ (right panel). The marginal prior density for ϕ is shown by the grey curve.

Figure 4 displays a histogram of the sampled values from the marginal posterior distributions of the measurement error precision ϕ and the true maximum value in the yellow (Mdm2) channel, $(\zeta_y^3)^{\max}$. The posterior density for ϕ is fairly similar to its prior density, indicating that the data have not been particularly informative about likely values of ϕ . The true maximum value for Mdm2, $(\zeta_y^3)^{\max}$, has posterior mean 330.8 and equal-tailed 95% credible interval (310.1, 352.7). Note that these values are considerably larger than the observed scaled maximum value of 145.0667.

We can gain some confidence in the validity of our fitted model by comparing predictive simulations from the model with the observed data z . This model validation by predictive simulations is advocated by Gelman *et al.* (2004). Figure 5 shows a plot of the time-course data for cell 3 (circles), together with shading representing point-wise equal-tailed 95% posterior predictive probability intervals, and a line representing the estimated posterior predictive mean. From Figure 5 we see that the model fits the data reasonably well. There is room for improvement in the fit of the p53 data, in terms of both the predictive mean, and the predictive variance. Despite the fact that the Mdm2 data have been normalised, we still achieve an acceptable fit to the Mdm2 data, although the predictive mean lies above the majority of the datapoints.

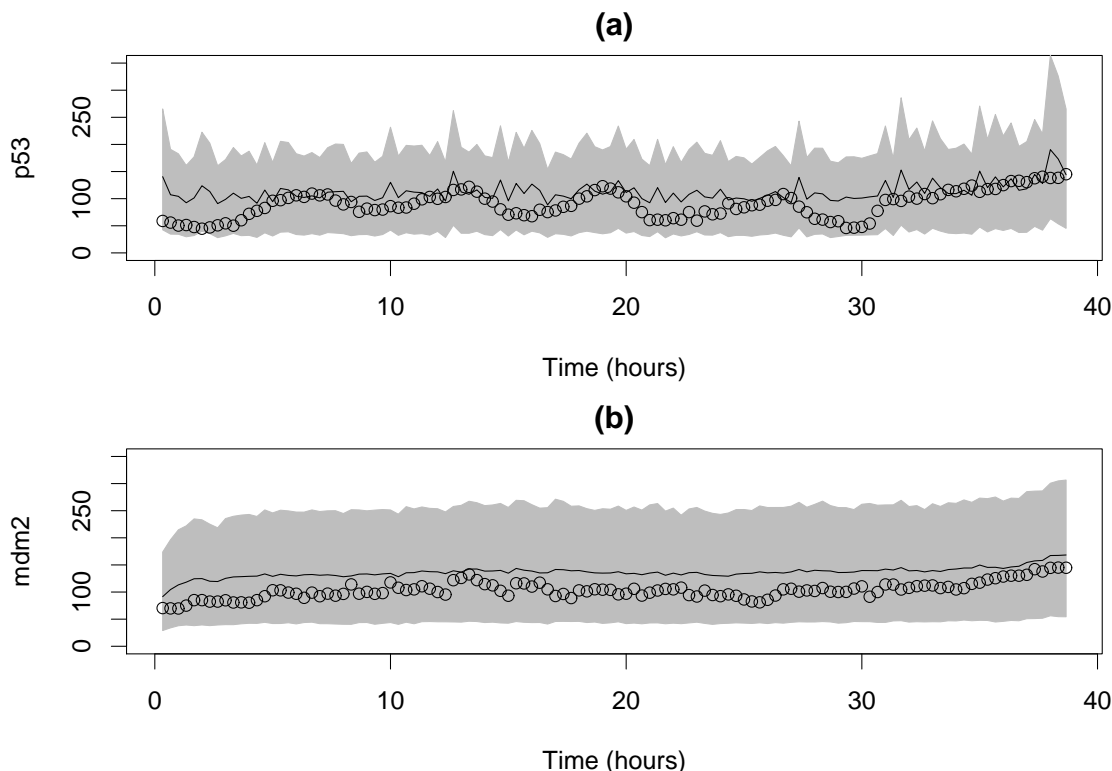


Figure 5: Time course data (circles) together with (posterior) predictive means (lines) and equal-tailed point-wise 95% predictive intervals (shading); (a) p53 fluorescence and (b) Mdm2 fluorescence.

7 Inference based on multiple cells

Although it is useful to look at what is learned about the model calibration parameters and the other unknown quantities in the model from the data on a single cell, it is natural to try and use all the available experimental data (or as much of it as is feasible) in order to make inferences.

A full Bayesian calibration of the model based on all 141 available cells is not computationally feasible at this time so here we look at all the cells in one particular movie. We chose the second movie, which consists of seven cells. The seven time-courses are all of different lengths, ranging from 4 hours 20 minutes (for cell 1) to 41 hours (for cell 2), and are shown in Figure 1.

Several independent Markov chains were simulated from different starting points, using the MCMC algorithm outlined in Section 5, with $C = 7$. We report here the results of one of these chains. The MCMC algorithm was run for 100,000 iterations after discarding an initial 150,000 iterates as burn-in. The output was thinned by taking every 25th iterate, leaving a sample of 4,000 iterates on which to base inferences.

Figure 6 displays histograms of the sampled values of the model calibration parameters. A comparison of this figure with Figure 2 shows that the marginal inferences for these parameters based on all seven cells are not too dissimilar to those based only on the third cell. However these densities are less butted up against the boundaries imposed by the prior distributions than those obtained in the single cell analysis. Also the data from the additional six cells have helped to reduce further the uncertainty about these parameters. These comments are reinforced by a comparison of the values for the posterior means and standard deviations in the seven cell analysis, given in Table 5, with those for the single cell analysis in Table 4. Overall, the data have dramatically reduced

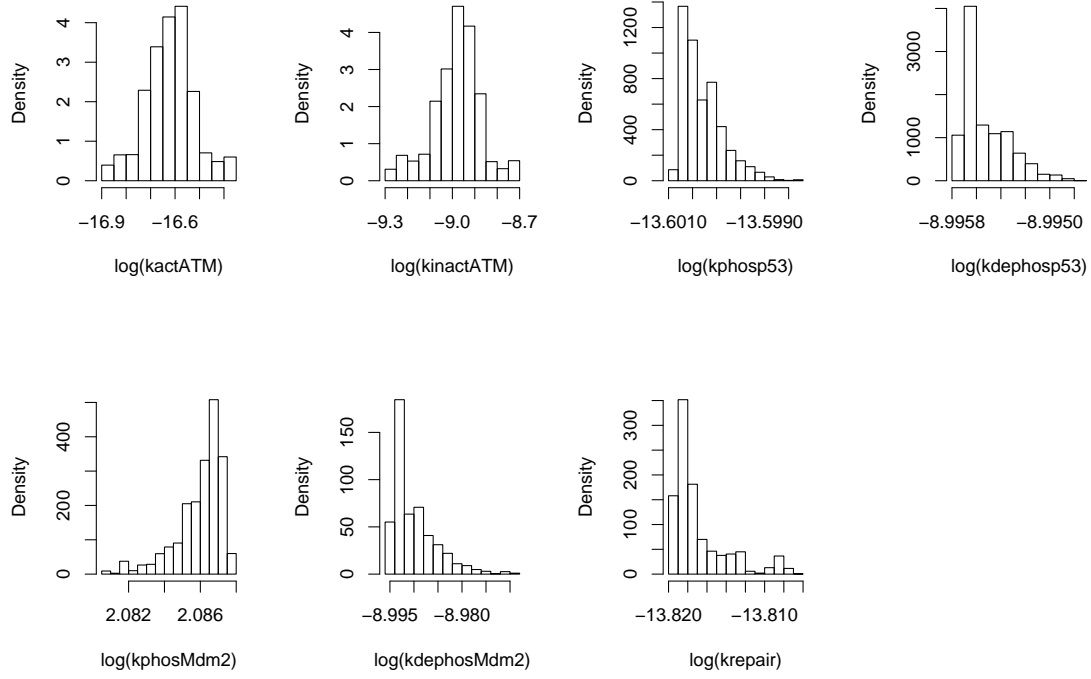


Figure 6: Model calibration parameters: marginal posterior densities based on all seven cells from movie 2.

Index	Parameter name	Posterior (prior)	
		Mean	Standard deviation
1	<i>kactATM</i>	-16.621 (-9.210)	0.1015 (5.1962)
2	<i>kinactATM</i>	-8.982 (-7.601)	0.1046 (5.1962)
3	<i>kphosp53</i>	-13.600 (-7.601)	0.0004 (3.4641)
4	<i>kdephosp53</i>	-8.996 (-4.147)	0.0002 (2.8000)
5	<i>kphosMdm2</i>	2.086 (-2.761)	0.0013 (2.8000)
6	<i>kdephosMdm2</i>	-8.990 (-4.147)	0.0045 (2.8000)
7	<i>krepair</i>	-13.817 (-10.820)	0.0028 (1.7321)

Table 5: Model calibration parameters: prior and posterior means and standard deviations based on all seven cells in movie 2.

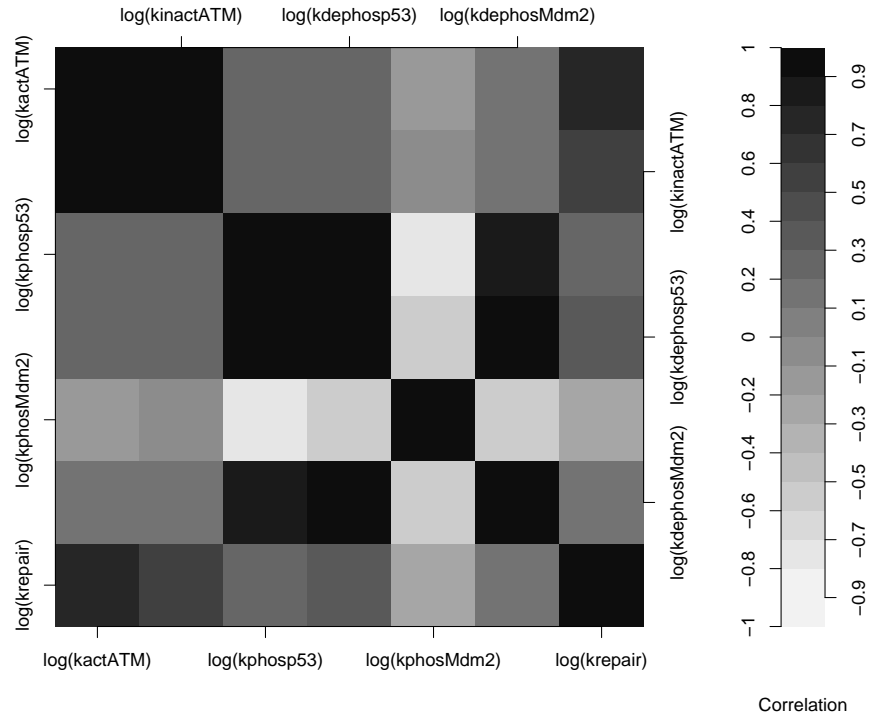


Figure 7: Model calibration parameters: posterior correlations based on all seven cells from movie 2.

uncertainty about the model calibration parameters. In comparison to their prior means, the data have been strongly suggestive that the kinetic rate constants are around two (or more) orders of magnitude smaller (on the original non-logged scale), with the exception of *kphosMdm2* which is around two orders of magnitude larger.

An image plot representing posterior correlations between pairs of calibration parameters is shown in Figure 7. This plot reveals that incorporating the data on all seven cells has increased the already high posterior correlation between *kphosp53* and *kdephosp53*. In addition, it highlights the strong negative posterior correlation between *kdephosp53* and *kphosMdm2* that was found in the single cell analysis. Further, it now shows that there is a strong negative posterior correlation between *kphosp53* and *kphosMdm2* that was not evident in the single cell analysis.

Figure 8 displays histograms of sampled values from the marginal posterior distribution of the measurement error precision ϕ and the true maximum values of fluorescence in the yellow channels for each of the seven cells. Our inferences about the value of the measurement error precision ϕ have changed considerably after incorporating data from multiple cells into the analysis. In particular, the data from all seven cells have been much more informative about ϕ than that from the third cell alone, suggesting that the measurement process is much less precise than expected *a priori*, though this could simply be a comment on an overall lack of fit of the model. Inferences on the maximum fluorescence values in the yellow (Mdm2) channels (Figure 8, panels (b)-(h)) confirm that they take posterior mean values, ranging from 294.1 for cell 4 to 495.3 for cell 5, which are much larger than their observed scaled maximum values.

How well does the model fit the data? Figures 9 and 10 show the data for the seven cells together with summaries of the posterior predictive distributions obtained by sampling replicate data from the calibrated model (in a similar fashion to Figure 5). The fit of the model to the data on cells 1 and 3 seems satisfactory, but the fit to data on the other cells shows considerable room

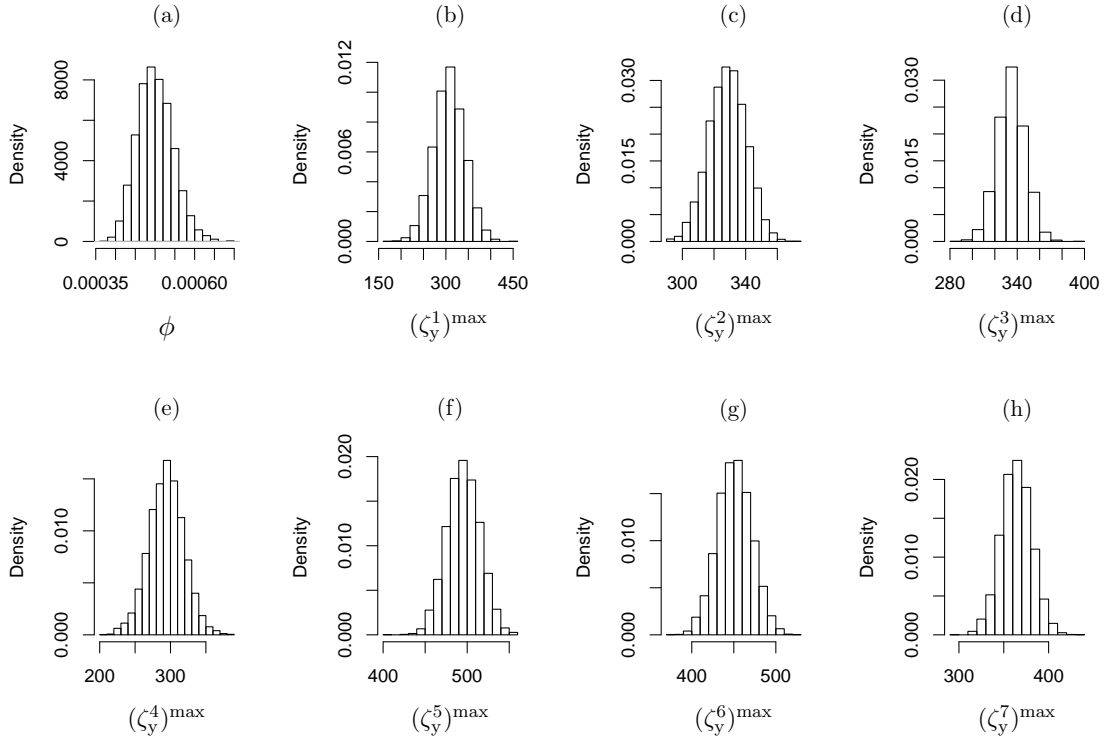


Figure 8: Marginal posterior densities of (a) the measurement error precision ϕ and (b)-(h) the true maximum values in the yellow (Mdm2) channel for cells 1 to 7.

for improvement. The predictive means are almost always larger than the observed datapoints, and there is considerable predictive uncertainty. The lack of fit in these predictive plots goes some way to explain the preference for small values of the measurement error precision ϕ . However, the lack of fit may not be too surprising given the inferential challenges posed by the scaling and normalisation of the data.

Another issue that may go some way to explaining the smaller than anticipated measurement error precision (suggestive of some lack of fit) is the distinction between endogenous p53 and Mdm2 and the exogenous fluorescent fusion proteins that are actually measured. The distinction between these is ignored in the model as it is argued in Geva-Zatorsky *et al.* (2006) that they should behave similarly *in vivo*. However, ideally these different species would be modelled separately in the stochastic kinetic model and the data linked only to the fusion proteins. Explicit modelling of fusion proteins separately from the targets being reported on is currently in its infancy, and is the subject of on-going modelling work.

8 Further Discussion

In this chapter we have demonstrated how it is possible to develop computationally intensive MCMC-based procedures for conducting a Bayesian analysis of an intra-cellular stochastic systems biology model using single-cell time course data. The information provided by this analysis is very rich. However, there are clearly several extensions of this work that merit further study. First, the model for multiple cells would benefit from the introduction of a random-effects layer in order to allow for the separation of inter- and intra-cell variation. Second, an integrated analysis allowing the comparison of competing model structures would be extremely valuable. For example, it is possible

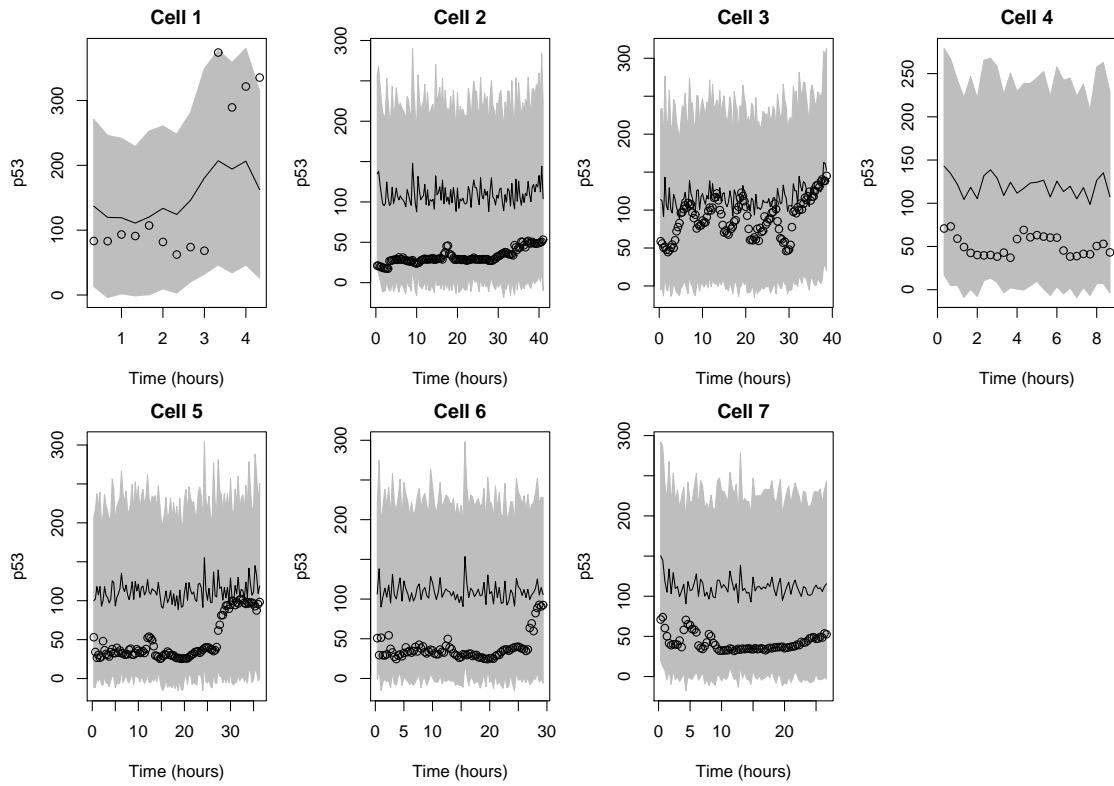


Figure 9: Time course data for p53 (circles) together with posterior predictive means (lines) and equal-tailed point-wise 95% predictive intervals (shading).

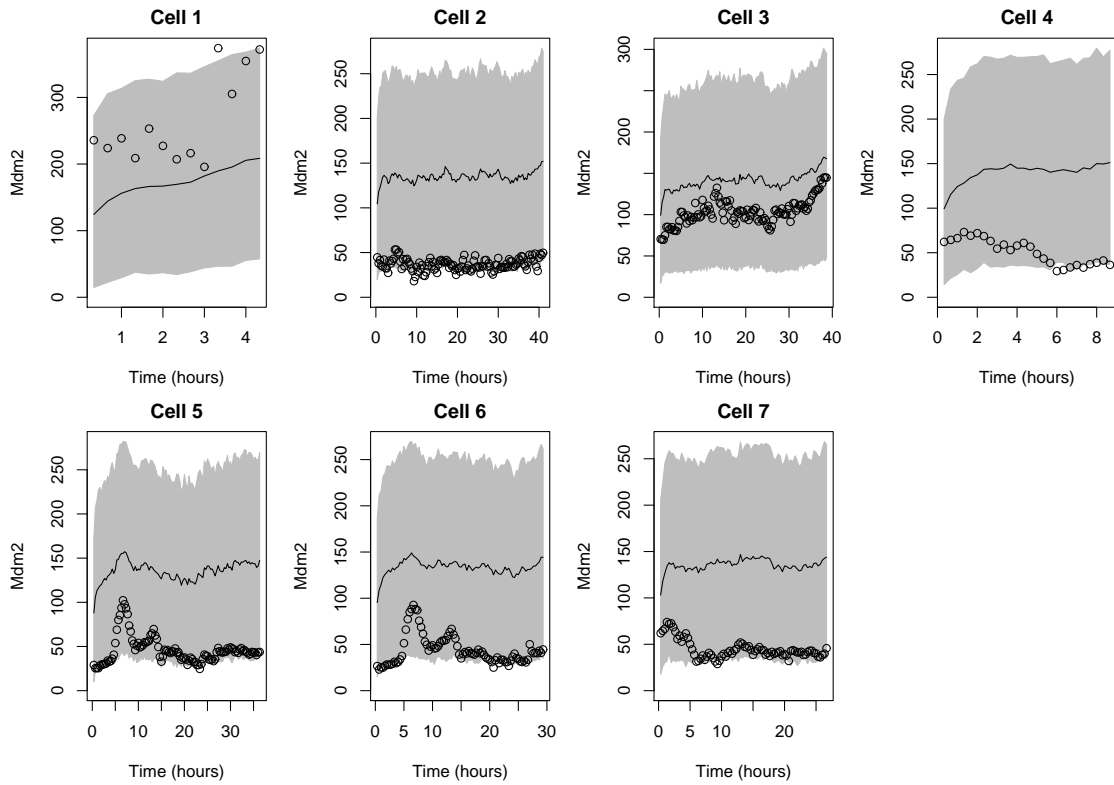


Figure 10: Time course data for Mdm2 (circles) together with posterior predictive means (lines) and equal-tailed point-wise 95% predictive intervals (shading).

to develop an alternative (competing) model by replacing the role of ATM with ARF. In this alternative model, the species ATMA and ATMI are removed, along with the reactions involving these species. The species p53-P and Mdm2-P are also removed as ARF works by a different mechanism to phosphorylation. ARF is initially set to zero but its level increases in the presence of damaged DNA. ARF binds to Mdm2 with a higher affinity than p53 and so levels of unbound p53 increase. This results in an increase of p53 transcriptional activity and so it is reasonable to predict an increase in levels of Mdm2 mRNA, followed by an increase in Mdm2. Since it is known that ARF increases the degradation rate of Mdm2, this model assumes that Mdm2 which is bound to ARF is degraded at a higher rate than normal. ARF is also degraded which allows the damage signal to decline as the damaged DNA is repaired. However, this mechanism seems to play a minor role in response to irradiation compared to ATM and so it is generally believed that this model is less appropriate for the experimental data. Nevertheless, there is genuine interest in knowing whether the data provide support for the ATM-based model considered here in favour of an ARF-based model considered less plausible by some biological experts. Extension of the algorithm to allow computation of the relevant Bayes factor ought to be straightforward in principle, but is likely to be quite difficult computationally.

More generally, the problem of constructing MCMC algorithms for models of this type is currently expensive in terms of both development and computation time. It is therefore natural to seek more straightforward and more automated approaches. Sequential Monte Carlo approaches may well offer considerable advantages in this area, given the Markovian nature of the underlying processes. Such sequential algorithms have already been considered for models of this type (Golightly and Wilkinson, 2006b), but considerable work remains to be done before they can be applied routinely to this general class of problems. Moreover, the normalisation that has been applied to the data described in this chapter means that these data do not lend themselves well to analysis via sequential methods. These and other related problems are the subject of current study within the CaliBayes project (<http://www.calibayes.ncl.ac.uk>).

Acknowledgements

We thank Uri Alon and Naama Geva-Zatorsky (Weizmann Institute of Science, Israel) for providing us with their raw experimental data, and Douglas Gray (Ottawa Health Research Institute) for advice on the biological aspects of the model construction. This work was funded by the UK Research Councils (BBSRC and EPSRC). In particular, most of the work for this chapter was directly funded by the BBSRC Bioinformatics and e-Science Initiative (BBSB16550) and the BBSRC Centre for Integrated Systems Biology of Ageing and Nutrition (BBC0082001).

Appendices

A: Broader Context and Background

The problem considered in this chapter is a special case of the general problem of conducting inference for the parameters of a Markovian stochastic process model using time course data. Although it is possible to consider stochastic processes that are intrinsically non-Markovian, it turns out that the Markovian class is very large, covering the vast majority of models that are derived from physical considerations. The class of Markov process models that may be considered is itself very large, but can be further categorised into sub-classes in various ways. It turns out that the most important attribute for classification purposes is whether the underlying stochastic process model is naturally formulated in discrete or continuous time. Time course data is typically discrete

in time, but as is the case in this chapter, it is nevertheless often natural to formulate the model in continuous time. Discrete time models have been studied more widely as they are technically simpler to work with, and often quite adequate if prediction is more important than parameter inference. The class of discrete time models can be further divided depending on whether the state space of the discrete time stochastic process is discrete or continuous. If it is discrete, then the model most likely falls into the general class of *hidden Markov models* (HMMs). Bayesian inference for HMMs is a well-studied problem, with Scott (2002) providing a comprehensive review and details of computation; also see Boys and Henderson (2004) and Fearnhead (2006). Alternatively, if the state space is continuous, then the model is often referred to as a linear or non-linear *state space* or *dynamic model*. The linear case is referred to as the *dynamic linear (state space) model* (DLM), and is studied in detail in West and Harrison (1997).

Clearly the problem considered in this chapter falls into the class of problems concerned with parameter inference for continuous time Markov process models. Here again it is helpful to subdivide this class depending on the state space of the model. If the state space is discrete and finite, then the model is a *continuous time hidden Markov model* (CTHMM), and can be tackled using techniques very similar to those used for discrete time HMMs (Fearnhead, 2008). If the state space is infinite, then inference is less straightforward (Fearnhead and Meligkotsidou, 2004). A very large class of such models (including stochastic kinetic models, and non-spatial Markovian stochastic epidemic models) is covered (in principle) by the algorithms developed in Boys *et al.* (2008), though the techniques described there do not scale well to problems of realistic size and complexity. For large and complex problems, several alternative possibilities exist. One approach is to exploit *stochastic emulators* of the process of interest, as advocated in Henderson *et al.* (2009) and this chapter. Another possibility is to exploit a combination of *sequential Monte Carlo* (Doucet *et al.*, 2001) methods and *likelihood-free MCMC* (Marjoram *et al.*, 2003). A rather different solution to the problem is to approximate the discrete state continuous time model with a continuous state model (the diffusion approximation), and then use methods for models described by stochastic differential equations (Golightly and Wilkinson, 2005).

Markovian models continuous in both time and state are typically described by stochastic differential equations. Inference for stochastic differential equation models is a rather technical topic, and problematic due to the fact that “obvious” MCMC algorithms are subject to pathological mixing problems. An excellent introduction to the topic, describing the essential structure of basic algorithms, the inherent problems with the obvious approach, and an elegant solution for the univariate case, is given in Roberts and Stramer (2001). An effective sequential Monte Carlo algorithm for the multivariate case is described in Golightly and Wilkinson (2006a) and applied to realistic systems biology scenarios in Golightly and Wilkinson (2006b). An effective global MCMC algorithm is described in Golightly and Wilkinson (2008) and more generally in Golightly and Wilkinson (2009). See Wilkinson (2006) and Wilkinson (2009) for further details of stochastic process models in the context of systems biology.

B. Construction of an emulator

Our emulator is a tractable approximation to the conditional probability distribution given by $p(Y(t + \tau)|Y(t), \theta)$, where $\tau = 1200$ seconds is the time difference between each of the datapoints. It is constructed based on methodology developed in the deterministic computer models literature and successfully applied to a stochastic computer model in Henderson *et al.* (2009). In essence, we model the joint distribution by carefully chosen univariate marginal and conditional distributions, each having the form of a standard probability distribution, but whose parameters are smooth functions of the model calibration parameters and other additional inputs.

In order to construct the emulator, the stochastic kinetic model was forward simulated for $\tau = 1200$

seconds from 2000 randomly chosen state vectors $Y(t)$ and parameter values θ , and the output in the form of the values of the 10 species was recorded. The 2000 design points were generated by using an approximate maximin Latin hypercube sample. Latin hypercube sampling was popularized as a strategy for generating input points for computer experiments by McKay *et al.* (1979). A maximin Latin hypercube sample is a Latin hypercube sample which maximises the minimal distance between pairs of design points. Maximin Latin hypercube designs are described in more detail in Santner *et al.* (2003).

As the model we are approximating is stochastic rather than deterministic we ran the model independently 100 times at each of the 2000 design points. The total simulation took around 3 days of CPU time. This was split over 50 2.2GHz processors, and so took less than two hours of real time.

The main complicating issue is that the conditional distribution $Y(t + \tau)|Y(t), \theta$ is 10-dimensional. A common practice when emulating multivariate outputs is to build an emulator for each output independently of all the others (although there have been recent developments on the construction of dynamic, multivariate emulators (Conti and O’Hagan, 2007; Conti *et al.*, 2007)). Independent emulators naturally ignore any correlations between the outputs and so are likely to be poorer approximations to the underlying stochastic kinetic model when such correlations exist. The approach we have used in this chapter is to construct an emulator for each univariate component of the factorisation of the joint conditional probability of the form

$$p(Y(t + \tau)|Y(t), \theta) = p(Y_1(t + \tau)|Y(t), \theta) \times p(Y_2(t + \tau)|Y_1(t + \tau), Y(t), \theta) \\ \times \cdots \times p(Y_{10}(t + \tau)|Y_1(t + \tau), \dots, Y_9(t + \tau), Y(t), \theta).$$

This reduces the task to that of fitting 10 univariate emulators to the output of the stochastic kinetic model. There is no natural ordering of the 10 species in the factorisation of the joint distribution, and so we focus on a particular ordering based on computational considerations. Because of the discrete nature of the output from the stochastic kinetic model we model the univariate component probabilities using Poisson distributions, with species 6 (ATMA) and 7 (ATMI) being modelled using a binomial distribution because their sum is constrained. By looking at the means and variances of the simulator output over the 100 replications we found that the Poisson assumption was not totally appropriate. However, since the output showed both over- and under-dispersion relative to the Poisson we decided to stick with the Poisson as a highly tractable compromise. The parameters of the Poisson or binomial distributions were assumed to be functions of the covariates (which are the species and model calibration parameters that they are conditioned on). For example, we model

$$Y_{10}(t + \tau)|Y_1(t + \tau), \dots, Y_9(t + \tau), Y(t), \theta \sim Po(\exp\{f(Y_1(t + \tau), \dots, Y_9(t + \tau), Y(t), \theta)\}),$$

where the Poisson mean is a function of the covariates. We have found that taking the function $f(\cdot)$ to be a low-order polynomial (quadratic or cubic) in the covariates to be adequate. In particular, we have found for this particular example that no improvement in fit (to some independent validation data) was obtained by additionally allowing for code uncertainty through inclusion of a Gaussian process term in $f(\cdot)$. It would appear that allowing for stochastic variation in the output through our construction of simple probability models can account for several of the standard sources of uncertainty encountered in the analysis of complex computer models and outlined in Section 2.1 of Kennedy and O’Hagan (2001). Each component distribution in the resulting emulator is either Poisson or binomial with parameters that are deterministic functions of the appropriate set of covariates. Therefore, we have a tractable approximation to the joint distribution of interest which we denote $p^*(\cdot|\cdot, \theta)$.

References

- Barak, Y., Juven, T., Haffner, R. and Oren, M. (1993) Mdm2 expression is induced by wild type-p53 activity. *Embo Journal*, **12**(2), 461–468.
- Boys, R. J. and Henderson, D. A. (2004) A Bayesian approach to DNA sequence segmentation (with discussion). *Biometrics*, **60**, 573–588.
- Boys, R. J., Wilkinson, D. J. and Kirkwood, T. B. L. (2008) Bayesian inference for a discretely observed stochastic kinetic model. *Statistics and Computing*, **18**, 125–135.
- Christophorou, M. A. Martin-Zanca, D., Soucek, L., Lawlor, E., Brown-Swigart, L., Verschuren, E. W. and Evan, G. I. (2005) Temporal dissection of p53 function in vitro and in vivo. *Nature Genetics*, **37**(7), 718–726.
- Ciliberto, A., Novak, B. and Tyson, J. J. (2005) Steady states and oscillations in the p53/Mdm2 network. *Cell Cycle*, **4**(3), 488–493.
- Clegg, H. V., Itahana, K. and Zhang, Y. (2008) Unlocking the Mmd2-p53 loop — ubiquitin is the key. *Cell Cycle*, **7**(3), 1–6.
- Conti, S., Gosling, J. P., Oakley, J. E. and O’Hagan, A. (2007) Gaussian process emulation of dynamic computer codes. Research Report No. 571/07, Department of Probability and Statistics, University of Sheffield.
- Conti, S. and O’Hagan, A. (2007) Bayesian emulation of complex multi-output and dynamic computer models. Research Report No. 569/07, Department of Probability and Statistics, University of Sheffield.
- Doucet, A., de Freitas, N. and Gordon, N. (eds) (2001) *Sequential Monte Carlo Methods in Practice*. New York: Springer.
- Fang, S. Y., Jensen, J. P., Ludwig, R. L., Vousden, K. H. and Weissman, A. M. (2000) Mdm2 is a RING finger-dependent ubiquitin protein ligase for itself and p53. *Journal of Biological Chemistry*, **275**(12), 8945–8951.
- Fearnhead, P. (2006) Exact and efficient Bayesian inference for multiple changepoint problems. *Statistics and Computing*, **16**(2), 203–213.
- (2008) Computational methods for complex stochastic systems: a review of some alternatives to MCMC. *Statistics and Computing*, **18**(2), 151–171.
- Fearnhead, P. and Meligkotsidou, L. (2004) Exact filtering for partially observed continuous time models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, **66**, 771–789.
- Finlay, C. A. (1993) The Mdm-2 oncogene can overcome wild-type p53 suppression of transformed cell growth. *Molecular Cell Biology*, **13**(1), 301–306.
- Gelman, A., Carlin, J. B., Stern, H. S. and Rubin, D. B. (2004) *Bayesian Data Analysis*. Boca Raton, Florida: Chapman and Hall/CRC, 2nd edn.
- Geva-Zatorsky, N., Rosenfeld, N., Itzkovitz, S., Milo, R., Sigal, A., Dekel, E., Yarnitzky, T., Liron, Y., Polak, P., Lahav, G. and *et al* (2006) Oscillations and variability in the p53 system. *Molecular Systems Biology*, **2**, 2006.0033.

- Gillespie, D. T. (1977) Exact stochastic simulation of coupled chemical reactions. *The Journal of Physical Chemistry*, **81**, 2340–2361.
- (2000) The chemical Langevin equation. *Journal of Chemical Physics*, **113**, 297–306.
- Golightly, A. and Wilkinson, D. J. (2005) Bayesian inference for stochastic kinetic models using a diffusion approximation. *Biometrics*, **61**, 781–788.
- (2006a) Bayesian sequential inference for nonlinear multivariate diffusions. *Statistics and Computing*, **16**, 323–338.
- (2006b) Bayesian sequential inference for stochastic kinetic biochemical network models. *Journal of Computational Biology*, **13**, 838–851.
- (2008) Bayesian inference for nonlinear multivariate diffusion models observed with error. *Computational Statistics and Data Analysis*, **52**(3), 1674–1693.
- (2009) Markov chain Monte Carlo algorithms for SDE parameter estimation. In *Learning and Inference for Computational Systems Biology*, MIT Press. In press.
- Hainaut, P. and Hollstein, M. (2000) p53 and human cancer: The first ten thousand mutations. *Advances in Cancer Research*, **77**, 81–137.
- Haupt, Y., Maya, R., Kazaz, A. and Oren, M. (1997) Mdm2 promotes the rapid degradation of p53. *Nature*, **387**(6630), 296–299.
- Henderson, D. A., Boys, R. J., Krishnan, K. J., Lawless, C. and Wilkinson, D. J. (2009) Bayesian emulation and calibration of a stochastic computer model of mitochondrial DNA deletions in substantia nigra neurons. *Journal of the American Statistical Association*. In press.
- Hirata, H., Yoshiura, S., Ohtsuka, T., Bessho, Y., Harada, T., Yoshikawa, K. and Kageyama, R. (2002) Oscillatory expression of the bHLH factor Hes1 regulated by a negative feedback loop. *Science*, **298**(5594), 840–843.
- Hucka, M., Finney, A., Sauro, H. M., Bolouri, H., Doyle, J. C., Kitano, H., Arkin, A. P., Bornstein, B. J., Bray, D., Cornish-Bowden, A., Cuellar, A. A., Dronov, S., Gilles, E. D., Ginkel, M., Gor, V., Goryanin, I. I., Hedley, W. J., Hodgman, T. C., Hofmeyr, J.-H., Hunter, P. J., Juty, N. S., Kasberger, J. L., Kremling, A., Kummer, U., Novere, N. L., Loew, L. M., Lucio, D., Mendes, P., Minch, E., Mjolsness, E. D., Nakayama, Y., Nelson, M. R., Nielsen, P. F., Sakurada, T., Schaff, J. C., Shapiro, B. E., Shimizu, T. S., Spence, H. D., Stelling, J., Takahashi, K., Tomita, M., Wagner, J. and Wang, J. (2003) The Systems Biology Markup Language (SBML): a medium for representation and exchange of biochemical network models. *Bioinformatics*, **19**(4), 524–531.
- Kennedy, M. C. and O’Hagan, A. (2001) Bayesian calibration of computer models (with discussion). *Journal of the Royal Statistical Society, Series B*, **63**, 425–464.
- Khan, S., Guevara, C., Fujii, G. and Parry, D. (2004) P14ARF is a component of the p53 response following ionizing irradiation of normal human fibroblasts. *Oncogene*, **23**(36), 6040–6046.
- Khosravi, R., Maya, R., Gottlieb, T., Oren, M., Shiloh, Y. and Shkedy, D. (1999) Rapid ATM-dependent phosphorylation of MDM2 precedes p53 accumulation in response to DNA damage. *Proceedings of the National Academy of Sciences*, **96**(26), 14,973–14,977.
- Lahav, G., Rosenfeld, N., Sigal, A., Geva-Zatorsky, N., Levine, A. J., Elowitz, M. B. and Alon, U. (2004) Dynamics of the p53-Mdm2 feedback loop in individual cells. *Nature Genetics*, **36**(2), 147–150.

- Lane, D. P. (1992) p53, guardian of the genome. *Nature*, **358**(6381), 15–16.
- Lev Bar-Or, R., Maya, R., Segel, L. A., Alon, U., Levine, A. J. and Oren, M. (2000) Generation of oscillations by the p53-Mdm2 feedback loop: A theoretical and experimental study. *Proceedings of the National Academy of Sciences*, **97**(21), 11,250–11,255.
- Ma, L., Wagner, J., Rice, J. J., Hu, W. W., Levine, A. J. and Stolovitzky, G. A. (2005) A plausible model for the digital response of p53 to DNA damage. *Proceedings of the National Academy of Sciences*, **102**(40), 14,266–14,271.
- Marjoram, P., Molitor, J., Plagnol, V. and Tavaré, S. (2003) Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences, USA*, **100**, 15,324–15,328.
- McKay, M. D., Beckman, R. J. and Conover, W. J. (1979) A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, **21**, 239–245.
- Mendrysa, S. M. and Perry, M. E. (2000) The p53 tumor suppressor protein does not regulate expression of its own inhibitor, MDM2, except under conditions of stress. *Molecular Cell Biology*, **20**(6), 2023–2030.
- Nelson, D. E., Ihekwebi, A. E. C., Elliott, M., Johnson, J. R., Gibney, C. A., Foreman, B. E., Nelson, G., See, V., Horton, C. A., Spiller, D. G. and *et al* (2004) Oscillations in NF-kappaB signaling control the dynamics of gene expression. *Science*, **306**(5696), 704–708.
- O’Hagan, A. (2006) Bayesian analysis of computer code outputs: a tutorial. *Reliability Engineering and System Safety*, **91**, 1290–1300.
- Pereg, Y., Shkedy, D., de Graaf, P., Meulmeester, E., Edelson-Averbukh, M., Salek, M., Biton, S., Teunisse, A. F. A. S., Lehmann, W. D., Jochemsen, A. G. and *et al* (2005) Phosphorylation of Hdmx mediates its Hdm2- and ATM-dependent degradation in response to DNA damage. *Proceedings of the National Academy of Sciences*, **102**(14), 5056–5061.
- Proctor, C. J. and Gray, D. A. (2008) Explaining oscillations and variability in the p53–Mdm2 system. *BMC Systems Biology*, **2**, 75.
- Roberts, G. O. and Stramer, O. (2001) On inference for non-linear diffusion models using Metropolis-Hastings algorithms. *Biometrika*, **88**(3), 603–621.
- Santner, T. J., Williams, B. J. and Notz, W. I. (2003) *The Design and Analysis of Computer Experiments*. New York: Springer.
- Scott, S. L. (2002) Bayesian methods for hidden Markov models: recursive computing in the 21st Century. *Journal of the American Statistical Association*, **97**, 337–351.
- Thut, C. J., Goodrich, J. A. and Tjian, R. (1997) Repression of p53-mediated transcription by MDM2: a dual mechanism. *Genes & Development*, **11**(15), 1974–1986.
- Tiana, G., Krishna, S., Pigolotti, S., Jensen, M. H. and Sneppen, K. (2007) Oscillations and temporal signalling in cells. *Physical Biology*, **4**(2), R1–R17.
- Vogelstein, B., Lane, D. and Levine, A. J. (2000) Surfing the p53 network. *Nature*, **408**(6810), 307–310.
- West, M. and Harrison, J. (1997) *Bayesian Forecasting and Dynamic Models*. New York: Springer, 2nd edn.

- Wilkinson, D. J. (2006) *Stochastic Modelling for Systems Biology*. Boca Raton, Florida: Chapman & Hall/CRC.
- (2009) Stochastic modelling for quantitative description of heterogeneous biological systems. *Nature Reviews Genetics*. In press.
- Zhang, L. J., Yan, S. W. and Zhuo, Y. Z. (2007) A dynamical model of DNA-damage derived p53-Mdm2 interaction. *Acta Physica Sinica*, **56**(4), 2442–2447.
- Zhang, Y., Xiong, Y. and Yarbrough, W. G. (1998) ARF promotes MDM2 degradation and stabilizes p53: ARF-INK4a locus deletion impairs both the Rb and p53 tumor suppression pathways. *Cell*, **92**(6), 725–734.