

# Inferring parameters for a lattice-free model of cell migration and proliferation using experimental data.

Alexander P. Browning<sup>a</sup>, Scott W. McCue<sup>a</sup>, Rachelle N. Binny<sup>b,c,d</sup>, Michael J. Plank<sup>b,d</sup>, Esha T. Shah<sup>e</sup>, Matthew J. Simpson<sup>a,\*</sup>

<sup>a</sup>*School of Mathematical Sciences, Queensland University of Technology (QUT), Brisbane, Australia.*

<sup>b</sup>*Landcare Research, Lincoln, Canterbury, New Zealand.*

<sup>c</sup>*Biomathematics Research Centre, University of Canterbury, Christchurch, New Zealand.*

<sup>d</sup>*Te Pūnaha Matatini, a New Zealand Centre of Research Excellence, New Zealand.*

<sup>e</sup>*Ghrelin Research Group, Translational Research Institute, QUT, 37 Kent St, Woolloongabba, Queensland, Australia.*

---

## Abstract

Collective cell spreading takes place in spatially continuous environments, yet it is often modelled using discrete lattice-based approaches. Here, we use data from a series of cell proliferation assays, with a prostate cancer cell line, to calibrate a spatially continuous individual based model (IBM) of collective cell migration and proliferation. The IBM explicitly accounts for crowding effects by modifying the rate of movement, direction of movement, and the rate of proliferation by accounting for pair-wise interactions. Taking a Bayesian approach we estimate the free parameters in the IBM using rejection sampling on three separate, independent experimental data sets. Since the posterior distributions for each experiment are similar, we perform simulations with parameters sampled from a new posterior distribution generated by combining the three data sets. To explore the predictive power of the calibrated IBM, we forecast the evolution of a fourth experimental data set. Overall, we show how to calibrate a lattice-free IBM to experimental data, and our work highlights the importance of interactions between individuals. Despite great care taken to distribute cells as uniformly as possible experimentally, we find evidence of significant spatial clustering over short distances, suggesting that standard mean-field models could be inappropriate.

**Keywords:** individual based model, cell migration, model calibration, cell proliferation assay, approximate Bayesian computation

---



---

\*Corresponding author at: Mathematical Sciences, QUT, Brisbane, Australia. Tel.: +617 3138 5241; fax: +617 3138 2310

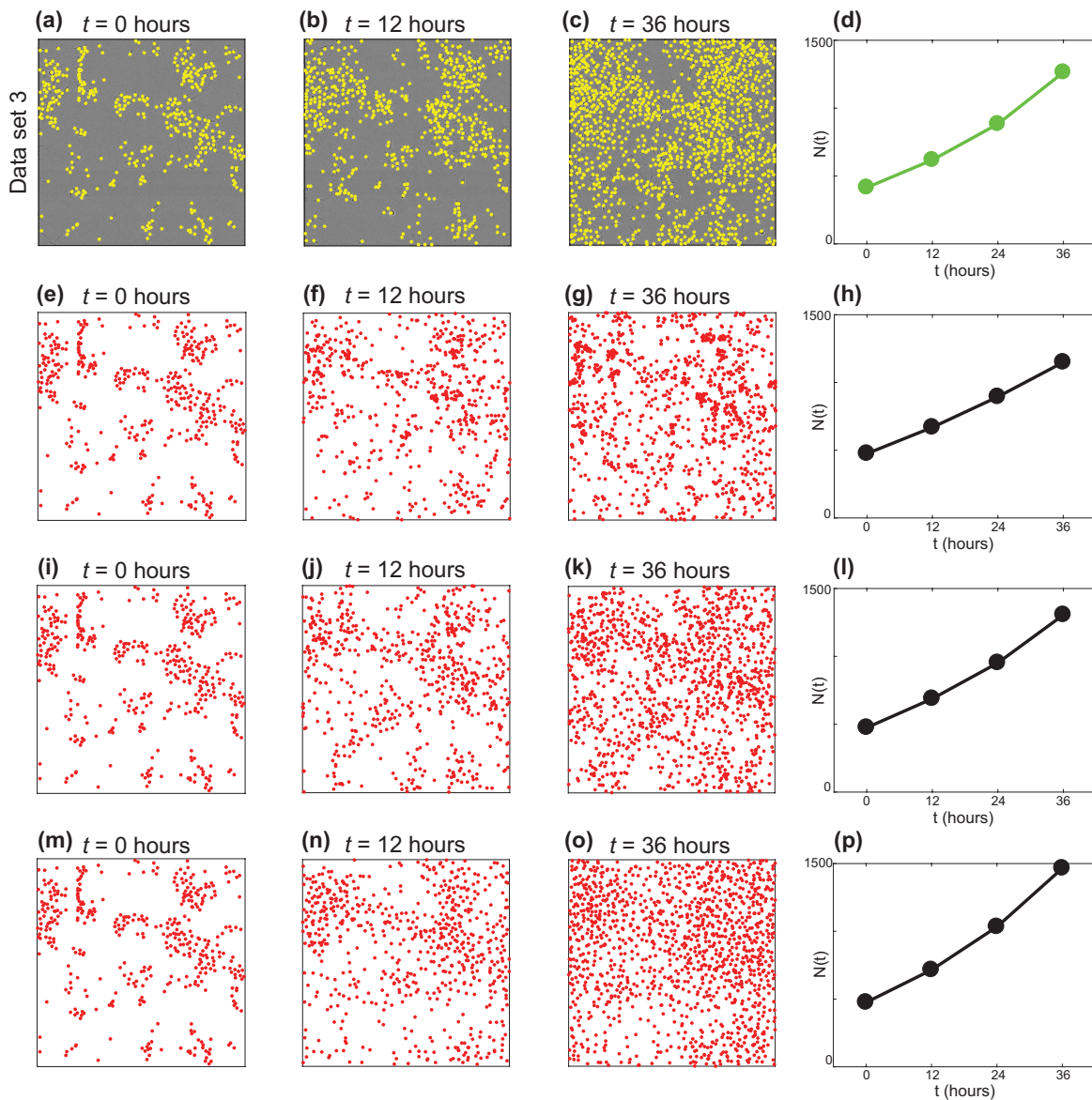
Email address: [matthew.simpson@qut.edu.au](mailto:matthew.simpson@qut.edu.au) (Matthew J. Simpson)

# 1. Introduction

One of the most common *in vitro* cell biology experiments is called a *cell proliferation assay* (Bosco et al., 2015; Bourseguin et al., 2016; Browning et al., 2017). These assays are conducted by placing a monolayer of cells, at low density, on a two-dimensional substrate. Individual cells undergo proliferation and movement events, and the assay is monitored over time as the density of cells in the monolayer increases (Tremel et al., 2009). One approach to interpret a cell proliferation assay is to use a mathematical model (Warne et al. 2017). Calibrating the solution of a mathematical model to data from a cell proliferation assay can provide quantitative insight into the underlying mechanisms, by, for example, estimating the cell proliferation rate (Tremel et al., 2009; Sengers et al., 2007). A standard approach to modelling a cell proliferation assay is to use a mean-field model, which is equivalent to assuming that individuals within the population interact in proportion to the average population density and that there is no spatial structure, such as clustering (Tremel et al., 2009; Sengers et al., 2007; Maini et al., 2004b; Sarapata and de Pillis, 2014; Sherratt and Murray, 1990). More recently, increased computational power has meant that individual based models (IBMs) have been used to directly model the cell-level behaviour (Binny et al., 2016a; Frascoli et al., 2013; Johnston et al., 2014). IBMs are attractive for modelling biological phenomena because they can be used to represent properties of individual agents, such as cells, in the system of interest (Binny et al., 2016a,b; Frascoli et al., 2013; Peirce et al., 2004; Read et al., 2012; Treloar et al., 2013). Typical IBMs use a lattice, meaning that both the position of agents, and the direction of movement, are restricted (Codling et al., 2008). In contrast, lattice-free IBMs are more realistic because they enable agents to move in continuous space, in any direction. However, this extra freedom comes at the cost of higher computational requirements (Plank and Simpson, 2012).

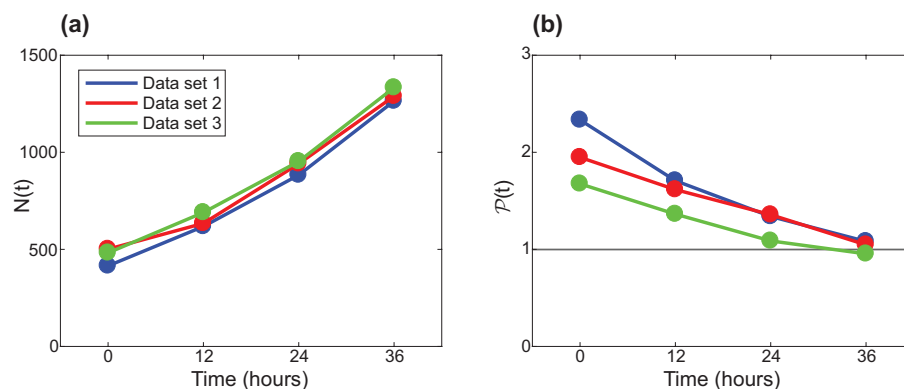
In this work we consider a continuous-space, continuous-time IBM (Binny et al., 2016b). This IBM is well-suited to studying experimental data from a cell proliferation assay with PC-3 prostate cancer cells (Kaighn et al., 1979), as shown in Figure 1(a)-(d). The key mechanisms in the experiments include cell migration and cell proliferation, and we note that there is no cell death in the experiments on the time scale that we consider. Therefore, agents in the IBM are allowed to undergo both proliferation and movement events. Crowding effects that are often observed in two-dimensional

cell biology experiments (Cai et al., 2007) are explicitly incorporated into the IBM as the rates of proliferation and movement in the model are inhibited in regions of high agent density. In this study we specifically choose to work with the PC-3 cell line because these cells are known to be highly migratory, mesenchymal cells (Kaighn et al., 1979). This means that cell-to-cell adhesion is minimal for this cell line, and cells tend to migrate as individuals. We prefer to work with a continuous-space, lattice-free IBM as this framework gives us the freedom to identically replicate the initial location of all cells in the experimental data when we specify the initial condition in the IBM. In addition, lattice-free IBMs do not restrict the direction of movement like a lattice-based approach.



**Fig. 1:** (a)-(c) Experimental data set 3 at  $t = 0, 12$  and  $36$  hours. The position of each cell is identified with a yellow marker. The field of view is a square of length  $1440 \mu\text{m}$ . (d) Population size,  $N(t)$  for experimental data set 3. (e)-(h) One realisation of the IBM with  $\gamma_b = 0 \mu\text{m}$ , leading to an overly clustered distribution of agents. (i)-(l) One realisation of the IBM with  $\gamma_b = 6.0 \mu\text{m}$ , leading to a distribution of agents with similar clustering to the experimental data. (m)-(p) One realisation of the IBM with  $\gamma_b = 20 \mu\text{m}$ , leading to an overly segregated distribution of agents. All IBM simulations are initiated using the same distribution of agents as in (a), with  $m = 1.0$  /hour,  $p = 0.040$  /hour, and  $\sigma = 24 \mu\text{m}$ .

37 A key contribution of this study is to demonstrate how the IBM can be calibrated to experimen-  
 38 tal data. In particular, we use approximate Bayesian computation (ABC) to infer the parameters  
 39 in the IBM. Four sets of experimental images (Supplementary Material 1), each corresponding to  
 40 an identically-prepared proliferation assay, are considered. The experiments are conducted over a  
 41 duration of 36 hours, which is unusual because proliferation assays are typically conducted for no  
 42 more than 24 hours (Browning et al., 2017). Data from the first three sets of experiments (Figure  
 43 2) are used to calibrate the IBM and data from the fourth set of images is used to examine the  
 44 predictive capability of the calibrated IBM. The IBM that we work with was presented very recently  
 45 (Binny et al., 2016b). The description of the IBM by Binny et al. (2016b) involves a discussion of  
 46 the mechanisms in the model and the derivation of a spatial moment continuum description (Binny  
 47 et al., 2016b). Our current work is the first time that experimental data has been used to provide  
 48 parameter estimates for this new IBM.



**Fig. 2:** Summary statistics for experimental data sets 1, 2 and 3, shown in blue, red and green, respectively. (a) Population size,  $N(t)$ . (b) Local measure of spatial structure,  $\mathcal{P}(t)$ , given by Equation 10. Unprocessed experimental data are given in Supplementary Material documents 1 and 2.

Taking a Bayesian approach, we assume that cell proliferation assays are stochastic processes, and model parameters are random variables, allowing us to update information about the model parameters using ABC (Collis et al., 2017; Tanaka et al., 2006). For this purpose we perform a large number of IBM simulations using parameters sampled from a prior distribution. Previous work, based on mean-field models, suggests that the proliferation rate and cell diffusivity for PC-3 cells is  $\lambda \approx 0.05$  /hour and  $D \approx 175 \mu\text{m}^2/\text{hour}$ , respectively (Johnston et al., 2015). The prior distribution for the IBM parameters are taken to be uniform and to encompass these previous estimates. We generate  $10^6$  realisations of the IBM using parameters sampled from the prior distribution, and accept the top 1% of simulations that provide the best match to the experimental data. Our approach to connect the experimental data and the IBM is novel, we are unaware of any previous work that has used ABC to parameterise a lattice-free IBM of a cell proliferation assay. One possible reason why ABC methods are not routinely used to calibrate lattice-free models of cell migration and cell proliferation with crowding effects is because of high computational requirements (Fröhlich et al., 2016). For example, we find that the typical run time to simulate our experiments is approximately 2 seconds on a standard desktop machine using C++. This means that simulating  $10^6$  realisations for inference with three unknown parameters becomes challenging. All work presented here is simulated on a High Performance Computing cluster to manage these computational limitations (QUT High Performance Computing, 2017).

Applying the ABC algorithm to data from three sets of identically prepared experiments leads to three similar posterior distributions. This result provides confidence that the IBM is a realistic representation of the cell proliferation assays and leads us to produce a combined posterior distribution from which we use the mean to give point estimates of the model parameters. To provide further validation of the IBM, we use the combined posterior distribution and the IBM to make a prediction of the fourth experimental data set. Simulating the IBM with parameters sampled from the combined posterior distribution allows us to predict both the time evolution of the population size,  $N(t)$ , and a measure of the density of pairs of cells,  $\mathcal{P}(t)$ , which provides a measure of spatial structure. These results indicate that the *in silico* predictions are consistent with the experimental observations.

This manuscript is organised as follows. Sections 2.1-2.2 describe the experiments and the IBM, respectively. In Section 2.3 we explain how to apply the ABC algorithm to estimate the IBM parameters. In Section 3 we present the marginal posterior distributions of the IBM parameters using data from the first three sets of experiments. The predictive power of the calibrated IBM is demonstrated by using the combined posterior distributions to predict the fourth experimental data set. The predictive power of the calibrated IBM is compared with a stochastic analogue of the standard mean-field logistic equation (Murray, 2002). While both models can accurately predict  $N(t)$ , the logistic equation provides no information about the spatial structure in the experimental data. Finally, in Section 4, we conclude and summarise opportunities for further research.

## 2. Material and methods

### 2.1. Experimental methods

We perform a series of proliferation assays using the IncuCyte ZOOM™ live cell imaging system (Essen BioScience, MI USA) (Jin et al., 2017). All experiments are performed using the PC-3 prostate cancer cell line (Kaighn et al., 1979). These cells, originally purchased from American Type Culture Collection (Manassas, VA, USA), are a gift from Lisa Chopin (April, 2016). Cells are propagated in RPMI 1640 medium (Life Technologies, Australia) with 10% foetal calf serum (Sigma-Aldrich, Australia), 100 U/mL penicillin, and 100  $\mu$ g/mL streptomycin (Life Technologies), in plastic tissue culture flasks (Corning Life Sciences, Asia Pacific). Cells are cultured in 5% CO<sub>2</sub> and 95% air in a Panasonic incubator (VWR International) at 37 °C. Cells are regularly screened for *Mycoplasma*.

Approximately 8,000 cells are distributed in the wells of the tissue culture plate as uniformly as possible. After seeding, cells are grown overnight to allow for attachment and some subsequent growth. The plate is placed into the IncuCyte ZOOM™ apparatus, and images showing a field of view of size  $1440 \times 1440 \mu\text{m}$  are recorded every 12 hours for a total duration of 36 hours. Experimental images for experimental data set three is shown in Figure 1(a)-(c). Images from the other three data sets are provided in Supplementary Material 1. ImageJ is used to determine the approximate locations of individual cells in all images, this data is given in Supplementary Material 2. Summary statistics,  $N(t)$  and  $\mathcal{P}(t)$ , for the first three experimental data sets are given in 2(a)–(b).



## 2.2. Mathematical model

### 2.2.1. Individual based model

We consider an IBM describing the proliferation and movement of individual cells (Binny et al., 2016a,b). Since cell death is not observed in the experiments, the IBM does not include agent death. The IBM allows the net proliferation rate and the net movement rate of agents to depend on the spatial arrangement of other agents. To be consistent with previous experimental observations, the IBM incorporates a biased movement mechanism so that agents tend to move away from nearby crowded regions (Cai et al., 2007). We use the IBM to describe the dynamics of a population of agents on a square domain of length  $L = 1440 \mu\text{m}$  to match the field-of-view of the experimental data (Figure 1(a)-(c)). Agents in the model are treated as a series of points which we may interpret as a population of uniformly-sized discs with diameter  $\sigma = 24 \mu\text{m}$  (Supplementary Material 1). Each agent has location  $\mathbf{x}_n = (x_1, x_2)$ , for  $n = 1, \dots, N(t)$ . Since the field-of-view of each image is much smaller than the size of the well in the tissue culture plate, we apply periodic boundary conditions (Jin et. al., 2017).

Proliferation and movement events occur according to a Poisson process over time (Binny et al., 2016b). The  $n$ th agent is associated with neighbourhood-dependent rates,  $P_n \geq 0$  and  $M_n \geq 0$ , of proliferation and movement, respectively. These rates consist of intrinsic components,  $p > 0$  and  $m > 0$ , respectively. Crowding effects are introduced by reducing the intrinsic rates by a contribution from other neighbouring agents. These crowding effects are calculated using a kernel,  $w^{(\cdot)}(r)$ , that depends on the separation distance,  $r \geq 0$ , so that

$$P_n = \max \left( 0, p - \sum_{i \neq n}^{N(t)} w^{(p)}(r) \right), \quad (1)$$

$$M_n = \max \left( 0, m - \sum_{i \neq n}^{N(t)} w^{(m)}(r) \right). \quad (2)$$

Following Binny et al.,(2016), we specify the kernels to be Gaussian with width corresponding to the

125 cell diameter,  $\sigma$ , giving

$$w^{(p)}(r) = \gamma_p \exp\left(-\frac{r^2}{2\sigma^2}\right), \quad (3)$$

$$w^{(m)}(r) = \gamma_m \exp\left(-\frac{r^2}{2\sigma^2}\right). \quad (4)$$

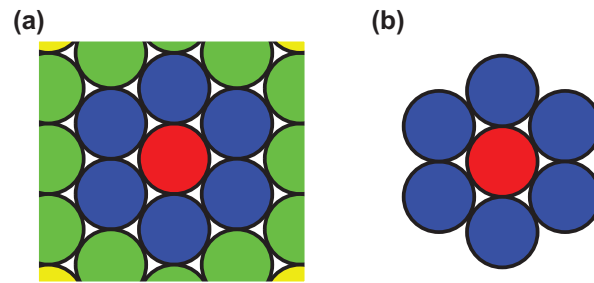
126 Here,  $\gamma_p$  is the value of  $w^{(p)}(0)$  and  $\gamma_m$  is the value of  $w^{(m)}(0)$ . These parameters provide a measure  
 127 of the strength of crowding effects on agent proliferation and movement, respectively. The kernels,  
 128  $w^{(p)}(r)$  and  $w^{(m)}(r)$ , ensure that the interactions between pairs of agents separated by more than  
 129 roughly 2-3 cell diameters lead to a negligible contribution. For computational efficiency, we truncate  
 130 the Gaussian kernels so that  $w^{(p)}(r) = w^{(m)}(r) = 0$ , for  $r \geq 3\sigma$  (Law et al., 2003).

131 To reduce the number of unknown parameters in the IBM, we specify  $\gamma_p$  and  $\gamma_m$  by invoking an  
 132 assumption about the maximum packing density of the population. Here we suppose that the net  
 133 proliferation and net movement rates reduce to zero when the agents are packed at the maximum  
 134 possible density, which is a hexagonal packing (Figure 3(a)). For interactions felt between the nearest  
 135 neighbours only (Figure 3(b)), we obtain

$$\gamma_p = \frac{p}{6} \exp\left(\frac{1}{2}\right), \quad (5)$$

$$\gamma_m = \frac{m}{6} \exp\left(\frac{1}{2}\right), \quad (6)$$

136 which effectively specifies a relationship between  $\gamma_p$  and  $p$ , and between  $\gamma_m$  and  $m$ . Note that this  
 137 assumption does not preclude a formation of agents in which some pairs have a separation of less  
 138 than  $\sigma$  and densities greater than hexagonal packing, which can occur by chance.



**Fig. 3:** (a) Hexagonal packing of uniformly sized discs. The focal agent (red) is surrounded by six nearest neighbouring agents (blue), and twelve next nearest neighbouring agents (green). (b) Hexagonal packing around a focal agent (red) showing the six nearest neighbours only.

When an agent at  $\mathbf{x}_n$  proliferates, the location of the daughter agent is selected by sampling from a bivariate normal distribution with mean  $\mathbf{x}_n$  and variance  $\sigma^2$  (Binny et al., 2016b). Since mesenchymal cells in two-dimensional cell culture are known to move with a directional movement bias away from regions of high density (Cai et al., 2007), we allow the model to incorporate a bias so that the preferred direction of movement is in the direction of decreasing agent density. For simplicity, the distance that each agent steps is taken to be a constant, equal to the cell diameter,  $\sigma$  (Plank and Simpson, 2012).

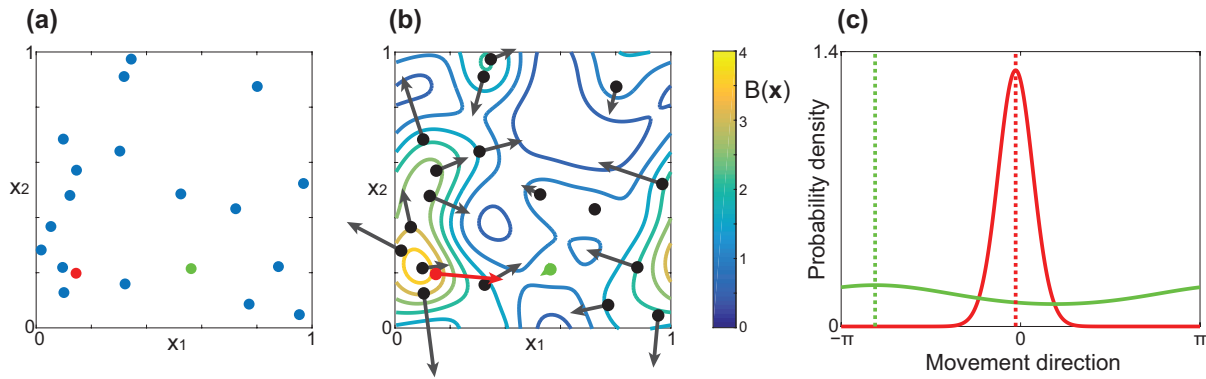
To choose the movement direction, we use a crowding surface,  $B(\mathbf{x})$ , to measure the local crowdedness at location  $\mathbf{x}$ , given by

$$B(\mathbf{x}) = \sum_{i=1}^{N(t)} w^{(b)}(\|\mathbf{x} - \mathbf{x}_i\|). \quad (7)$$

The crowding surface is the sum of contributions from every agent, given by a bias kernel,  $w^{(b)}(r)$ . The contributions depend on the distance between  $\mathbf{x}$  and the location of the  $i$ th agent,  $\mathbf{x}_i$ , given by  $r = \|\mathbf{x} - \mathbf{x}_i\|$ . Again, we choose  $w^{(b)}$  to be Gaussian, with width equal to the cell diameter, and repulsive strength,  $\gamma_b \geq 0$ , so that

$$w^{(b)}(r) = \gamma_b \exp\left(-\frac{r^2}{2\sigma^2}\right), \quad (8)$$

where  $\gamma_b$  is value of  $w^{(b)}(0)$ , and has dimensions of length. Note that  $B(\mathbf{x})$  is an increasing function of local density, and approaches zero as the local density decreases. A typical crowding surface is shown in Figure 4(b) for the arrangement of agents in Figure 4(a).



**Fig. 4:** (a) Example distribution of agents on a  $1 \times 1$  periodic domain. (b) Level curves of the corresponding crowding surface,  $B(\mathbf{x})$ , for this arrangement of agents. The arrows show the preferred direction of movement,  $\mathbf{B}_n$ . To illustrate how the direction of movement is chosen, (c) shows the probability density of the von Mises distribution for the red and green agents highlighted in (a) and (b). The preferred direction,  $\arg(\mathbf{B}_n)$ , is shown as dotted vertical lines for both agents. The red agent is in a crowded region so  $\|\mathbf{B}_n\|$  is large, meaning that the agent is likely to move in the preferred direction  $\arg(\mathbf{B}_n)$ . The green agent is in a low density region and  $\|\mathbf{B}_n\|$  is small, meaning that the bias is very weak and the agent's direction of movement is almost uniformly distributed. To illustrate the effects of the crowding surface as clearly as possible, we set  $\gamma_b = 1$ ,  $\sigma = 0.1$ ,  $L = 1$  in this schematic figure to draw attention to the gradient of the crowding surface.

To determine the direction of movement we use the shape of  $B(\mathbf{x})$  to specify the bias, or preferred direction, of agent  $n$ ,  $\mathbf{B}_n$ , given by

$$\mathbf{B}_n = -\nabla B(\mathbf{x}_n), \quad (9)$$

which gives the magnitude and direction of steepest descent. Results in Figure 4(b) show  $\mathbf{B}_n$  for the arrangement of agents in Figure 4(a). To determine the direction of movement, we consider the magnitude and direction of  $\mathbf{B}_n$ , and sample the actual movement direction from a von Mises distribution,  $\text{von Mises}(\arg(\mathbf{B}_n), \|\mathbf{B}_n\|)$  (Binny et al., 2016b; Forbes et al., 2011). Therefore, agents are always most likely to move in the direction of  $\mathbf{B}_n$ , however as  $\|\mathbf{B}_n\| \rightarrow 0$ , the preferred direction becomes uniformly distributed.

To illustrate how the direction of movement is chosen, we show, in Figure 4(b), the bias vector for each agent,  $\mathbf{B}_n$ . Note that  $\mathbf{B}_n$  does not specify the movement step length, and the direction of  $\mathbf{B}_n$  does not necessarily specify the actual direction. Rather,  $\arg(\mathbf{B}_n)$  specifies the preferred direction. To illustrate this property, we highlight two agents in Figure 4(a). The red agent is located on a relatively steep part of the crowding surface, so  $\|\mathbf{B}_n\|$  is large. The green agent is located on a relatively flat part of the crowding surface, so  $\|\mathbf{B}_n\|$  is close to zero. Figure 4(c) shows the von Mises distributions for the red and green agent. Comparing these movement distributions confirms that the crowded red agent is more likely to move in the direction of  $\mathbf{B}_n$ . The bias is weak for the green agent, so the direction of movement is almost uniformly distributed since  $\|\mathbf{B}_n\|$  is smaller.

IBM simulations are performed using the Gillespie algorithm (Gillespie, 1977). To initialise each simulation we specify the initial number and initial location of agents to match to the experimental images at  $t = 0$  hours (Supplementary Material 1) for experimental data sets 1, 2, 3 and 4. In all simulations we set  $\sigma = 24 \mu\text{m}$  and  $L = 1440 \mu\text{m}$ . The remaining three parameters,  $m$ ,  $p$  and  $\gamma_b$ , are varied with the aim of producing posterior distributions using a Bayesian framework.

If  $\gamma_m = \gamma_b = 0$ , and the variance of the dispersal distribution is large, the IBM corresponds to logistic growth (Binny et al., 2016b, Browning et al. 2017). Under these simplified conditions, a uniformly distributed initial population of agents will grow, at rate  $p$ , to eventually reach a uniformly distributed maximum average density of  $p/(2\pi\gamma_p\sigma_p^2)$ . We do not consider this case here as our initial distribution of cells in the experiments is clustered, and so the logistic growth model is, strictly

speaking, not valid (Binny et al., 2016b).

### 2.2.2. Summary statistics

To match the IBM simulations with the experimental data we use properties that are related to the first two spatial moments (Law et al., 2003). The first spatial moment, the average density, is characterised by the number of agents in the population,  $N(t)$ . The second spatial moment characterises how agents are spatially distributed, and is often reported in terms of a pair correlation function (Binny et al., 2016a,b; Law et al., 2003). In the Supplementary Material 1 document we present the pair correlation function for all four experimental data sets. These results show that we have a fairly typical pair correlation function that contains, at most, one maximum (Binder and Simpson, 2013). Therefore, instead of using all details contained in the pair correlation function, we use a simplified measure of spatial structure. We consider a local measure of pair density within a distance of  $R$   $\mu\text{m}$ , given by

$$\mathcal{P}(t) = \frac{L^2 \sum_{i=1}^{N(t)} \sum_{\substack{j=1 \\ j \neq i}}^{N(t)} \mathbb{I}_{\|\mathbf{x}_i - \mathbf{x}_j\| \leq R}}{N(t)^2 \pi R^2}, \quad (10)$$

where  $\mathbb{I}$  is an indicator function so that the double sum in Equation 10 gives twice the number of distinct pairs within a distance  $R$ . For all results presented in the main document we set  $R = 50$   $\mu\text{m}$ . Therefore,  $\mathcal{P}(t)$  is the ratio of the number of pairs of agents, separated by a distance of less than 50  $\mu\text{m}$ , to the expected number of pairs of agents separated by a distance of less than 50  $\mu\text{m}$ , if the agents were randomly distributed. This means that,  $\mathcal{P}(t) = 1$  corresponds to randomly placed agents;  $\mathcal{P}(t) > 1$  corresponds to a locally clustered distribution; and,  $\mathcal{P}(t) < 1$  corresponds to a locally segregated distribution.

To ensure that our choice of setting  $R = 50$   $\mu\text{m}$  is adequate, we also repeat some results with  $R = 100$   $\mu\text{m}$ . This exercise leads to very similar posterior distributions, confirming that working with  $R = 50$   $\mu\text{m}$  is sufficient (Supplementary Material 1).

### 2.3. Approximate Bayesian computation

We consider  $m, p$  and  $\gamma_b$  as random variables, and the uncertainty in these parameters is updated using observed data (Collis et al., 2017; Tanaka et al., 2006). To keep the description of the inference

algorithm succinct, we refer to the unknown parameters as  $\Theta = \langle m, p, \gamma_b \rangle$ .

In the absence of any experimental observations, information about  $\Theta$  is characterised by specified prior distributions. The prior distributions are chosen to be uniform on an interval that is wide enough to encompass previous estimates of  $m$  and  $p$  (Johnston et al., 2015). To characterise the prior for  $\gamma_b$ , we note that this parameter is related to a length scale over which bias interactions are felt. Preliminary results (not shown) use a prior in the interval  $0 \leq \gamma_b \leq 20 \mu\text{m}$  and suggest that a narrow prior in the interval  $0 \leq \gamma_b \leq 10 \mu\text{m}$  is appropriate. In summary, our prior distributions are uniform and independent, given by

$$\pi(m) = \text{U}(0, 10) / \text{hour}, \quad (11)$$

$$\pi(p) = \text{U}(0, 0.1) / \text{hour}, \quad (12)$$

$$\pi(\gamma_b) = \text{U}(0, 10) \mu\text{m}. \quad (13)$$

We always summarise data,  $\mathbf{X}$ , with a lower-dimensional summary statistic,  $S$ . Data and summary statistics from the experimental images are denoted  $\mathbf{X}_{\text{obs}}$  and  $S_{\text{obs}}$ , respectively. Similarly, data and summary statistics from IBM simulations are denoted  $\mathbf{X}_{\text{sim}}$  and  $S_{\text{sim}}$ , respectively. Information from the prior is updated by the likelihood of the observations,  $\pi(S_{\text{obs}}|\Theta)$ , to produce posterior distributions,  $\pi(\Theta|S_{\text{obs}})$ . We employ the most fundamental ABC algorithm, known as ABC rejection (Liepe et al., 2014; Tanaka et al., 2006), to sample from the approximate posterior distribution. The approximate posterior distributions are denoted  $\pi_u(\Theta|S_{\text{obs}})$ .

In this work we use a summary statistic that is a combination of  $N(t)$  and  $\mathcal{P}(t)$  at equally spaced time intervals. A discrepancy measure,  $\rho(S_{\text{obs}}, S_{\text{sim}})$ , is used to assess the closeness of  $S_{\text{obs}}$  and  $S_{\text{sim}}$ ,

$$\rho(S_{\text{obs}}, S_{\text{sim}}) = \sum_t \left( \frac{[N_{\text{sim}}(t) - N_{\text{obs}}(t)]^2}{N_{\text{obs}}(t)^2} + \frac{[\mathcal{P}_{\text{sim}}(t) - \mathcal{P}_{\text{obs}}(t)]^2}{\mathcal{P}_{\text{obs}}(t)^2} \right), \quad (14)$$

Algorithm 1 is used to obtain  $10^6 u$  samples,  $\{\Theta_i\}_{i=1}^{10^6 u}$ , from the approximate joint posterior distribution,  $\pi_u(\Theta|S_{\text{obs}})$ , for each data set. Here,  $u \ll 1$  is the accepted proportion of samples.

To present marginal posterior samples, we use a kernel density estimate to form smooth, approximate marginal posterior distributions, for each parameter, and each data set using the `ksdensity`



---

**Algorithm 1** ABC rejection sampling algorithm to obtain  $10^6 u$  samples from the approximate posterior distribution,  $\pi_u(\Theta|S_{\text{obs}})$ .

---

- 1: Set  $\sigma = 24 \mu\text{m}$ ,  $L = 1440 \mu\text{m}$ , and set  $\mathbf{x}_n$  to match experimental data  $\mathbf{X}_{\text{obs}}$  at  $t = 0$ .
  - 2: Draw parameter samples from the prior  $\Theta_i \sim \pi(\Theta)$ .
  - 3: Simulate cell proliferation assay with  $\Theta_i$  and  $t \leq 36$  hours.
  - 4: Record summary statistic  $S_{\text{sim}_i} = \{N_{\text{sim}}(t), \mathcal{P}(t)\}_t$ , where  $t = 12, 24$  and  $36$  hours.
  - 5: Compute the discrepancy measure  $\epsilon_i = \rho(S_{\text{obs}}, S_{\text{sim}_i})$ , given in Equation 14.
  - 6: Repeat steps 2-5 until  $10^6$  samples  $\{\Theta_i, \epsilon_i\}_{i=1}^{10^6}$  are simulated.
  - 7: Order  $\{\Theta_i, \epsilon_i\}_{i=1}^{10^6}$  by  $\epsilon_i$  such that  $\epsilon_1 < \epsilon_2 < \dots$ .
  - 8: Retain the first 1% ( $u = 0.01$ ) of prior samples  $\Theta_i$ , as posterior samples,  $\{\Theta_i\}_{i=1}^{10^6 u}$ .
- 

function in MATLAB (Mathworks, 2017). Point estimates of parameters are always given as the mean of the posterior distribution, and always presented to two significant figures.

### 2.3.1. Sampling from the combined posterior distribution

Samples from the posterior distributions for each experimental data set are given in Supplementary Material 1. Kernel density estimates for the marginal posterior distributions for each experimental data set are given in Figure 5. Visually, the posterior distributions for each experimental data set appear to be similar, therefore we are motivated to form a combined posterior distribution,  $\pi_u(\Theta|\{S_{\text{obs}}^{(k)}\}_{k=1}^3)$ , where  $S_{\text{obs}}^{(k)}$  is the summary statistic from the  $k$ th experimental data set. We use ABC rejection to sample from the combined posterior distribution according to Algorithm 2. That is, Algorithm 2 is designed to sample the combined posterior distribution by retaining the top 1% of parameter combinations that provide the best fit to all three experimental data sets.

---

**Algorithm 2** ABC rejection sampling algorithm to obtain  $10^6 u$  samples from the approximate combined posterior distribution,  $\pi_u(\Theta|\{S_{\text{obs}}^{(k)}\}_{k=1}^3)$ .

---

- 1: Set  $\sigma = 24 \mu\text{m}$ ,  $L = 1440 \mu\text{m}$ .
  - 2: Draw parameter samples from the prior  $\Theta_i \sim \pi(\Theta)$ .
  - 3: For experimental data sets  $k = 1, 2$  and  $3$ :
    - 3.1: Set  $\mathbf{x}_n$  to match experimental data set  $k$ ,  $\mathbf{X}_{\text{obs}}^{(k)}$  at  $t = 0$ .
    - 3.2: Simulate cell proliferation assay with  $\Theta_i$  and  $t \leq 36$  hours.
    - 3.3: Record summary statistic  $S_{\text{sim}_i}^{(k)} = \{N_{\text{sim}}(t), \mathcal{P}(t)\}_t$ , where  $t = 12, 24$  and  $36$  hours.
  - 4: Compute the discrepancy measure  $\epsilon_i = \sum_{k=1}^3 \rho(S_{\text{obs}}^{(k)}, S_{\text{sim}_i}^{(k)})$ , where  $\rho$  is given in Equation 14.
  - 5: Repeat steps 2-4 until  $10^6$  samples  $\{\Theta_i, \epsilon_i\}_{i=1}^{10^6}$  are simulated.
  - 6: Order  $\{\Theta_i, \epsilon_i\}_{i=1}^{10^6}$  by  $\epsilon_i$  such that  $\epsilon_1 < \epsilon_2 < \dots$ .
  - 7: Retain the first 1% ( $u = 0.01$ ) of prior samples  $\Theta_i$ , as posterior samples,  $\{\Theta_i\}_{i=1}^{10^6 u}$ .
-

### 2.3.2. Predicting experimental data set 4 using the combined posterior distribution

To test the predictive power of the calibrated IBM, we use  $10^4$  parameter samples from the combined posterior distribution and simulate the IBM initialised with the actual initial arrangement of cells in data set 4 at  $t = 0$ . For each parameter combination  $S_{\text{sim}}$  is recorded at 12 hour intervals, and used to construct distributions of  $N(t)$  and  $\mathcal{P}(t)$ . These distributions are represented as box plots and compared with summary statistics from experimental data set 4.

### 2.3.3. Calibrating the standard mean-field logistic model to the experimental data

To illustrate the importance of considering individual level details in the IBM, we also calibrate the logistic growth model to experimental data sets 1, 2 and 3. The logistic growth model describes the IBM when spatial structure is neglected (Law et al., 2003; Binny et al., 2016b). The logistic growth model is given by

$$\frac{dN(t)}{dt} = \lambda N(t) \left( 1 - \frac{N(t)}{N_{\text{max}}} \right), \quad (15)$$

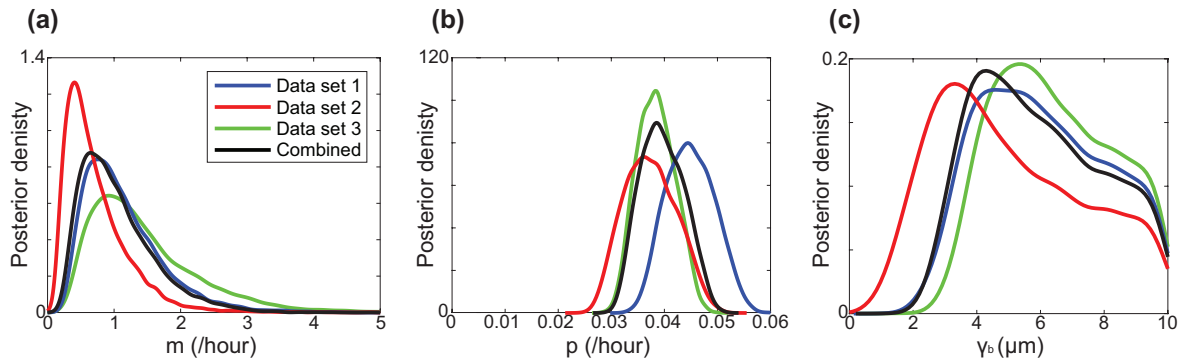
where  $\lambda$  is the cell proliferation rate and  $N_{\text{max}}$  is the maximum number of agents. To find estimates of  $\lambda$  and  $N_{\text{max}}$  to best match our experimental data we simulate the stochastic logistic model using the Gillespie algorithm (Gillespie, 1977; Fröhlich et. al., 2016). Proliferation events are treated as a Poisson process, with the rate given by the right hand side of Equation 15. Details of the ABC rejection algorithm used to estimate  $\lambda$  and  $N_{\text{max}}$  are given in Supplementary Material 1.

## 3. Results and discussion

To qualitatively illustrate the importance of spatial structure we show, in rows 2-4 of Figure 1, snapshots from the IBM with different choices of parameters. In each case the IBM simulations evolve from the initial condition specified in Figure 1(a). Results in the right-most column of Figure 1 compare the evolution of  $N(t)$  and we see that the parameter combination in the second row underestimates  $N(t)$ , the parameter combination in the fourth row overestimates  $N(t)$ , and the parameter combination in the third row produces a reasonable match to the experimental data. A visual comparison of the spatial arrangement of agents in rows 2-4 of Figure 1 suggests that these different parameter combinations may lead to different spatial structures. This illustration of how the IBM results vary with the choice of parameters motivates us to use ABC rejection to estimate the

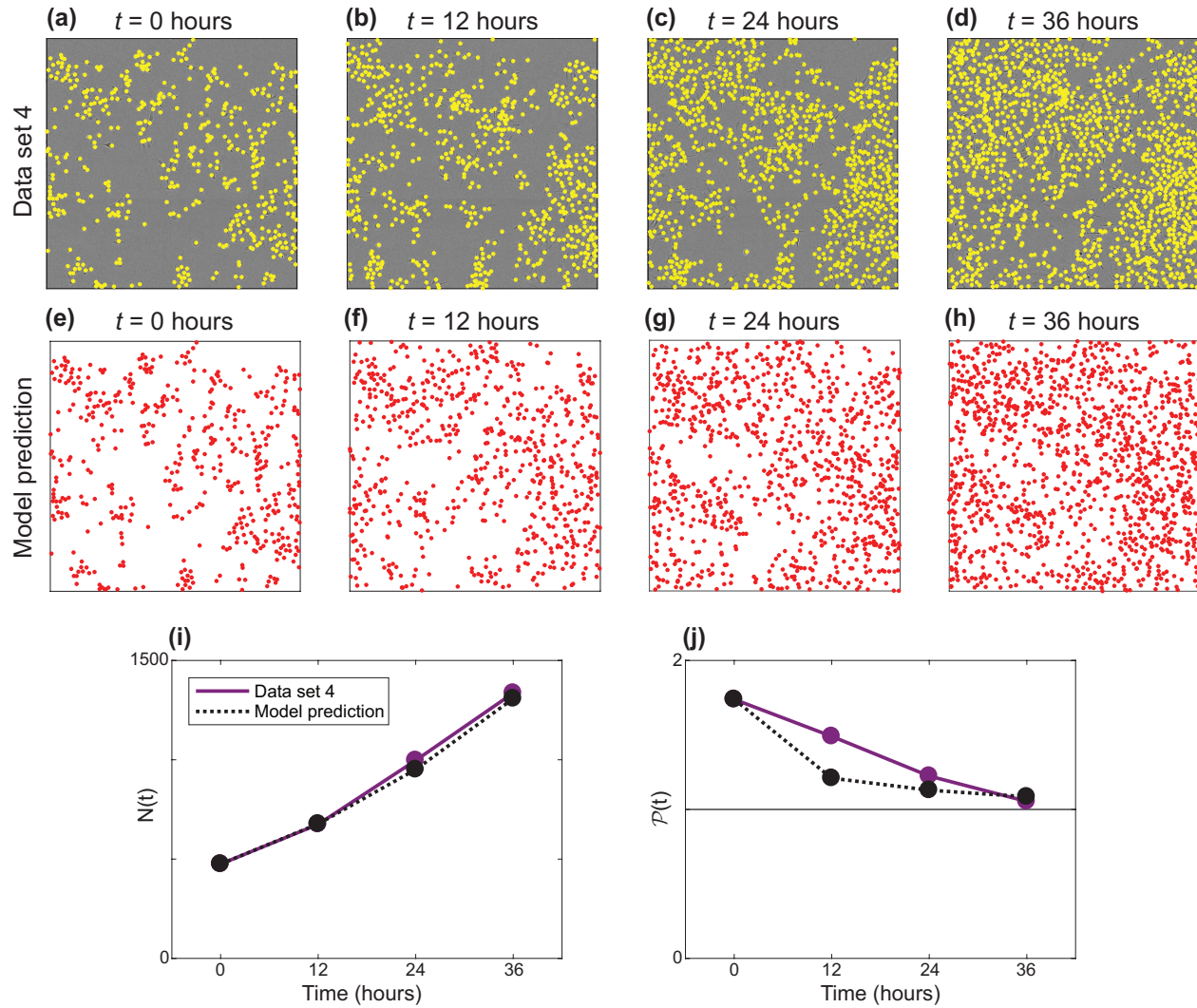
joint distribution of the parameters. To do this we will use summary statistics from three identically prepared, independent sets of experiments. The summary statistics for these experiments,  $N(t)$  and  $\mathcal{P}(t)$ , are summarised in Figure 2, and tabulated in Supplementary Material 1.

The approximate marginal posterior distributions for  $m$ ,  $p$  and  $\gamma_b$  are shown in Figure 5(a)-(c), respectively, for experimental data sets 1, 2 and 3. There are several points of interest to note. In each case, the posterior support is well within the interior of the prior support, suggesting that our choice of priors is appropriate. An interesting feature of the marginal posterior distributions for all parameters is that there is significant overlap for each independent experimental data set. There is some variation in the mean between experimental data sets, for each parameter, which could arise as a consequence of some other variation among experiments, or under the assumption that cell proliferation assays are stochastic processes. The combined marginal posterior distributions are superimposed, and the mean is given by 1.0 /hour, 0.040 /hour and 6.0  $\mu\text{m}$  for  $m$ ,  $p$  and  $\gamma_b$ , respectively. These point estimates of  $p$  and  $m$  give a cell doubling time of  $\ln(2)/p \approx 17$  hours, and a cell diffusivity of approximately 150  $\mu\text{m}^2/\text{hour}$ , which are typical values for PC-3 cells at low density (Johnston et al., 2015). All results in the main document correspond to retaining the top 1% of samples ( $u = 0.01$ ) and additional results (Supplementary Material 1) confirm that the results are relatively insensitive to this choice.



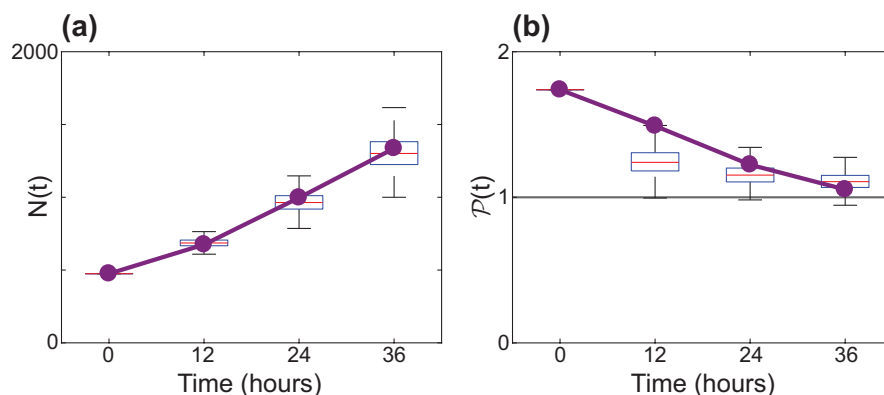
**Fig. 5:** (a)-(c) Kernel-density estimates of the approximate marginal posterior distributions for each data set, for parameters  $m$ ,  $p$  and  $\gamma_b$ , respectively, with  $u = 0.01$ . The combined posterior distribution (black) is superimposed. The point estimates from the combined posterior distribution are  $m = 1.0$  /hour,  $p = 0.040$  /hour and  $\gamma_b = 6.0$   $\mu\text{m}$ . All distributions are scaled so that the area under the curve is unity.

282 To assess the predictive power of the calibrated IBM, we attempt to predict the time evolution  
 283 of a separate, independently collected data set, experimental data set 4, as shown in Figure 6(a)-(d).  
 284 We use the mean of the combined posterior distribution and the initial arrangement of agents in  
 285 experimental data set 4 to produce a typical prediction in Figure 6(e)-(h). Visual comparison of the  
 286 experimental data and the IBM prediction suggests that the IBM predicts a similar number of agents,  
 287 and a similar spatial structure, with some short range clustering present. To quantify our results,  
 288 we compare the evolution of  $N(t)$  in Figure 6(i) which reveals an excellent match. Furthermore, we  
 289 predict the evolution of  $\mathcal{P}(t)$  in Figure 6(j) confirming similar trends.



**Fig. 6:** (a)-(d) Experimental images for data set 4. The position of each cell is identified with a yellow marker. The field of view is a square of length  $1440 \mu\text{m}$ . (e)-(h) One realisation of the IBM with parameters corresponding to the posterior mean:  $m = 1.0$  /hour,  $p = 0.040$  /hour and  $\gamma_b = 6.0 \mu\text{m}$ , with the same initial arrangement of agents as in (a). (i)  $N(t)$  for the experimental data (purple) and the IBM prediction (dashed black). (j)  $\mathcal{P}(t)$  for the experimental data (purple) and the IBM prediction (dashed black). The discrepancy measure (Supplementary Material 1) is given by  $\text{MSE} = 968$  for (i), and  $\text{MSE} = 0.0296$  for (j).

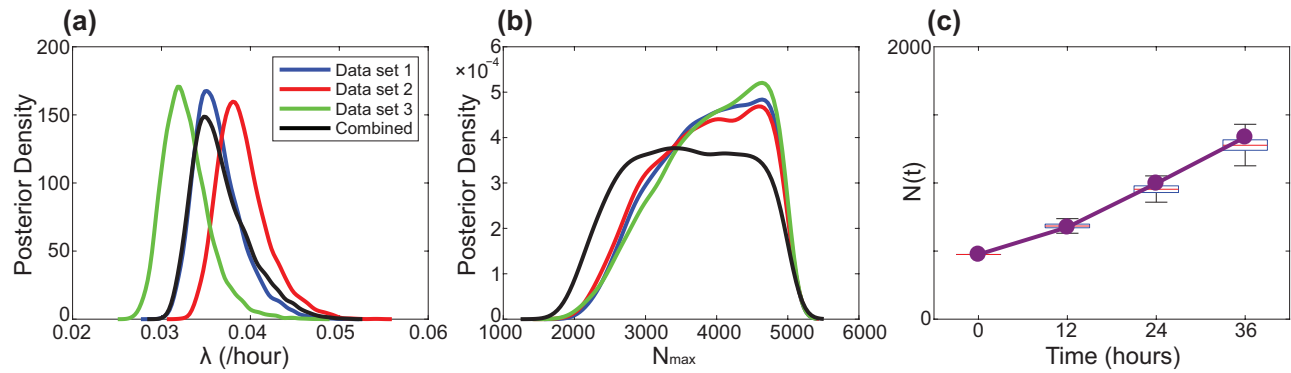
290 In addition to examining a single, typical realisation of the calibrated model, we now examine a  
 291 suite of realisations of the calibrated IBM, and compare results with experimental data set 4. The  
 292 suite of IBM realisations is obtained by sampling from the joint posterior distribution. Results in  
 293 Figure 7(a) compare  $N(t)$  from experimental data set 4 with distributions of  $N(t)$  from the suite of  
 294 IBM simulations, showing an excellent match. The spread of the distributions of  $N(t)$  increases with  
 295 time, which is expected. Results in Figure 7(b) compare the evolution of  $\mathcal{P}(t)$  from experimental  
 296 data set 4 with distributions of  $\mathcal{P}(t)$  from the suite of IBM simulations, showing the predicted  
 297 distributions of  $\mathcal{P}(t)$  overlap with the experimental data. Overall, the quality of the match between  
 298 the prediction and the experimental data is high, as the prediction captures both qualitative and  
 299 quantitative features of the data.



**Fig. 7:** Predictive distributions for  $N(t)$  and  $P(t)$ , respectively, generated using the IBM.  $10^4$  parameter samples are taken from the combined posterior distribution, and a model realisation produced for each sample, initiated as in Figure 6(a). Box plots show the distribution of  $N(t)$  and  $P(t)$  across these realisations in (a) and (b), respectively. The discrepancy measure (Supplementary Material 1) is taken to be between the mean of each boxplot and the observed data, given by  $MSE = 682$  for (a), and  $MSE = 0.0225$  for (b).



300 We now use ABC rejection to form combined posterior distributions of the parameters in the  
 301 standard logistic growth model, Equation 15,  $\lambda$  and  $N_{\max}$ . Results are shown in Figure 8(a)–(b).  
 302 The point estimates of the combined posterior distributions are  $\lambda = 0.037$  /hour and  $N_{\max} = 3600$ .  
 303 This estimate leads to a doubling time of approximately 19 hours, which is slightly longer than the  
 304 doubling time predicted using the calibrated IBM. We then examine a suite of solutions of Equation  
 305 15, where we sample from the combined posterior distribution for  $\lambda$  and  $N_{\max}$ . The predicted  
 306 distribution of  $N(t)$  is compared with experimental data set 4 in Figure 8(c), revealing an excellent  
 307 match.



**Fig. 8:** (a)-(b) Kernel-density estimates of the marginal posterior distributions are shown for each data set, for parameters in the stochastic logistic model, Equation 15,  $\lambda$  and  $N_{\max}$  in (a) and (b), respectively. The combined posterior distribution (black) is superimposed. The point estimates are  $\lambda = 0.037$  /hour and  $N_{\max} = 3600$ . All marginal distributions are scaled to an area of unity. (c) A predictive distribution for  $N(t)$ , generated from the stochastic logistic model, Equation 15.  $10^4$  parameter samples are taken from the combined posterior distribution, and a model realisation produced for each sample. Boxplots show the distribution of  $N(t)$  across these realisations. The procedure for sampling from the combined posterior distribution for the stochastic logistic model, and the procedure for solving the stochastic logistic model, are outlined in Supplementary Material 1. The discrepancy measure (Supplementary Material 1) is taken to be between the mean of each boxplot and the observed data, given by  $\text{MSE} = 1668$ .

Therefore, while both calibrated models provide good predictions for the observed evolution of  $N(t)$ , the IBM offers additional insights relating to spatial structure in the cell population, while the logistic model does not provide this level of information. The differences in the way that the logistic model and the IBM treat interactions between individuals could explain why the calibration process leads to different estimates of the proliferation rate. These differences suggest that the interactions between individuals appear to be relevant for our experimental data.

## 4. Conclusions

In this work we explore how to connect a spatially continuous IBM of cell migration and cell proliferation to novel data from a cell proliferation assay. Previous work parameterising IBM models of cell migration and cell proliferation to experimental data using ABC have been restricted to lattice-based IBMs (Johnston et al., 2014). This is partly because ABC methods require large numbers of IBM simulations, and lattice-based IBMs are far less computationally expensive than lattice-free IBMs (Plank and Simpson, 2012). We find it is preferable to work with a lattice-free IBM when dealing with experimental data as a lattice-based IBM requires approximations when mapping the distribution of cells from experimental images to a lattice (Johnston et al., 2014; Johnston et al., 2016). This mapping can be problematic. For example, if multiple cells in an experimental image are equally close to one lattice site, *ad hoc* assumptions have to be introduced about how to arrange those cells on the lattice without any overlap. These issues are circumvented using a lattice-free method.

To help overcome the computational cost of using ABC with a lattice-free IBM, we introduce several realistic simplifying assumptions. The IBM originally presented by Binny et al. (2016b) involves 12 free parameters, which is a relatively large number for standard inference techniques (Schnoerr et al. 2016). The model is simplified by noting that our experiments do not involve cell death, and specifying the width of the interaction kernels to be constant, given by the cell diameter. Another simplification is given by assuming that crowding effects reduce the proliferation and movement rates to zero when the agents are packed at the maximum hexagonal packing density. This leads to a simplified model with three free parameters:  $m$ ,  $p$  and  $\gamma_b$ . Using ABC rejection, we arrive at posterior distributions for these parameters for three independent experimental data sets. The marginal posterior distributions for the three parameters are similar, leading us to form a

combined posterior distribution. The point estimates from the combined posterior distributions for  $m$  and  $p$  are consistent with previous parameter estimates (Johnston et al., 2015) and the point estimate for  $\gamma_b$  is consistent with previous observations that mesenchymal cells in this kind of two-dimensional experiment tend to move away from regions of high cell density (Cai et al., 2007).

In the field of mathematical biology, questions about how much detail to include in a mathematical model, and what kind of mathematical model is preferable for understanding a particular biological process are often settled in an *ad hoc* manner, as discussed by Maclaren et al. (2015). Our approach in this work is to use a mathematical model that incorporates just the key mechanisms, with an appropriate number of unknown parameters. Other approaches are possible, such as using much more complicated mathematical models that describe additional mechanisms such as: (i) detailed information about the cell cycle in individual cells (Fletcher et al., 2012); (ii) concepts of leader and follower cells (Kabla, 2012); (iii) explicitly coupling cell migration and cell proliferation to the availability of nutrients and growth factors (Tang et al., 2014); or (iv) including mechanical forces between cells (Stichel et al., 2017). However, we do not include these kinds of detailed mechanisms because our experimental data does not suggest that these mechanisms are relevant to our situation. Furthermore, it is not always clear that using a more complicated mathematical model, with additional mechanisms and additional unknown parameters, necessarily leads to improved biological insight. In fact, simply incorporating additional mechanisms and parameters into the mathematical model often leads to a situation where multiple parameter combinations lead to equivalent predictions which limits the usefulness of the mathematical model (Simpson et al., 2006). In this study, our approach is to be guided by experimental data and our ability to infer the parameters in a mathematical model based on realistic amounts of experimental data (Maclaren et al. 2015). In particular we use three experimental data sets to calibrate the IBM, and an additional data set to separately examine the predictive capability of the calibrated IBM. We find that the process of calibrating the IBM leads to well defined posterior distributions of the model parameters, and that the calibrated IBM produces a reasonable match to the experimental data. The process of calibrating the IBM, and then separately testing the predictive capability of the calibrated IBM, provides some confidence that the level of model complexity is appropriate for our purposes.

An interesting feature of all experimental data at early time, when the cell density is relatively low, is that  $\mathcal{P}(t)$  suggests that the cells are clustered at short intervals, and that this clustering becomes less pronounced with time. This observation is very different to the way that previous theoretical studies have viewed the role of spatial structure. For example, previous simulation-based studies assume that some initial random spatial arrangement of cells can lead to clustering at later times (Baker and Simpson, 2010). In contrast, our experimental data suggests it could be more realistic to consider that the spatial structure is imposed by the initial arrangement of cells. Moreover, since all of our experimental data involves some degree of spatial clustering, our work highlights the importance of using appropriate models to provide a realistic representation of key phenomena. Almost all continuum models of collective behaviour in cell populations take the form of ordinary differential equations and partial differential equations that implicitly invoke a mean-field assumption (Tremel et al., 2009; Sengers et al., 2007; Maini et al., 2004b; Sarapata and de Pillis, 2014; Sherratt and Murray, 1990). Such assumptions ignore the role of spatial structure. While pairwise models that avoid mean-field assumptions are routine in some fields, such as disease spreading (Sharkey et al., 2006; Sharkey, 2008) and ecology (Law et al., 2003), models that explicitly account for spatial structure are far less common for collective cell behaviour.

Using our parameter estimates, the continuum spatial moment description could be used to interpret experimental data sets with larger numbers of cells (Binny et al., 2016b), such as experimental images showing a wider field-of-view, or experiments initiated with a higher density of cells. Our approach to estimate the parameters in the model is to work with the IBM since this allows us more flexibility in connecting with the experimental data, such as choosing the initial locations of the agents in the IBM to precisely match the initial locations of cells in the experimental images.

There are many ways that our study could be extended. For example, here we choose a summary statistic encoding information about the first two spatial moments. However, other summary statistics may provide different insight, and it could be of interest to explore the effect of this choice. For example, here we describe the spatial structure over a relatively short interval, approximately  $2\sigma$ . It could be of interest to repeat our analysis using the entire pair correlation function, accounting for spatial structure at all distances. However, here we take a simpler approach and we provide evidence

that our results are insensitive to our measure of spatial structure as we obtain similar results when we consider spatial structure over larger distances. Another limitation of our work is that the IBM, which explicitly accounts for interactions between agents, can be computationally expensive to simulate. This limitation can be particularly problematic for computational inference and severely limits the number of parameters that can be dealt with by taking a purely individual-based approach. One promising way of overcoming this difficulty is to make use of more theoretical ways to treat interactions between individuals in an IBM, and to perform inference using a stochastic continuum description, such as the recent work by Schnoerr et al. (2016; 2017).

## 5. Acknowledgements

This work is supported by the Australian Research Council (DP140100249, DP170100474) and the Royal Society of New Zealand Marsden Fund (11-UOC-005). Computational resources provided by the High Performance Computing and Research Support Group at QUT are appreciated. We thank David Warne for technical advice. We also thank the two anonymous referees for their helpful comments.

## 6. References

- [1] Baker RE, Simpson MJ. 2010. Correcting mean-field approximations for birth-death-movement processes. *Phys Rev E* **82**, 041905.
- [2] Binder, BJ, Simpson, MJ. 2013. Quantifying spatial structure in experimental observations and agent-based simulations using pair-correlation functions. *Phys Rev E* **88**, 022705.
- [3] Binny RN, Haridas P, James A, Law R, Simpson MJ, Plank MJ. 2016a. Spatial structure arising from neighbour-dependent bias in collective cell movement. *PeerJ* **4**, e1689.
- [4] Binny RN, James A, Plank MJ. 2016b. Collective cell behaviour with neighbour-dependent proliferation, death and directional bias. *Bull Math Biol* **78**, 2277–2301.
- [5] Bosco DB, Kenworthy R, Zorio DAR, Sang QXA. 2015. Human mesenchymal stem cells are resistant to paclitaxel by adopting a non-proliferative fibroblastic state. *PLoS One* **10**, e0128511.

- 417 [6] Bourseguin J, Bonet C, Renaud E, Pandiani C, Boncompagni M, Giuliano S, Pawlikowska P,  
418 Karmous-Benailly H, Ballotti R, Rosselli F, Bertolotto C. 2016. FANCD2 functions as a critical  
419 factor downstream of MiTF to maintain the proliferation and survival of melanoma cells. *Sci Rep*  
420 **6**, 36539.
- 421 [7] Browning AP, McCue SW, Simpson MJ. 2017. A Bayesian computational approach to explore  
422 the optimal duration of a cell proliferation assay. *Bull Math Biol* **10**, 1888-1906.
- 423 [8] Cai AQ, Landman KA, Hughes BD. 2007. Multi-scale modeling of a wound-healing cell migration  
424 assay. *J Theor Biol* **245**, 576–594.
- 425 [9] Codling EA, Plank MJ, Benhamou S. 2008. Random walk models in biology. *J R Soc Interface*  
426 **5**, 813–834.
- 427 [10] Collis J, Connor AJ, Paczkowski M, Kannan P, Pitt-Francis J, Byrne HM, Hubbard ME. 2017.  
428 Bayesian calibration, validation and uncertainty quantification for predictive modelling of tumour  
429 growth: a tutorial. *Bull Math Biol* **79**, 939–974.
- 430 [11] Fletcher AG, Breward CJW, Chapman SJ. 2012. Mathematical modelling of monoclonal con-  
431 version in the colonic crypt. *J Theor Biol* **300**, 118–133.
- 432 [12] Forbes C, Evans M, Hastings N, Peacock B. 2011. *Statistical distributions*. 4th ed. John Wiley  
433 & Sons, New Jersey.
- 434 [13] Frascoli F, Hughes BD, Zaman MH, Landman KA. 2013. A computational model for collective  
435 cellular motion in three dimensions: general framework and case study for cell pair dynamics. *PLoS*  
436 *ONE* **8**, e59249.
- 437 [14] Fröhlich, F, Thomas, P, Kazeroonian, A, Theis, FJ, Grima, R, Hasenauer, J, 2016. Inference  
438 for stochastic chemical kinetics using moment equations and system size expansion. *PLoS Comput*  
439 *Biol* **12**, e1005030.
- 440 [15] Gillespie DT. 1977. Exact stochastic simulation of coupled chemical reactions. *J Phys Chem* **81**,  
441 2340–2361.

- 442 [16] QUT High Performance Computing. 2017. [https://www.student.qut.edu.au/technology/research-](https://www.student.qut.edu.au/technology/research-computing/high-performance-computing)  
443 [computing/high-performance-computing](https://www.student.qut.edu.au/technology/research-computing/high-performance-computing). Accessed: October 2017.
- 444 [17] Kabla AJ. 2012. Collective cell migration: leadership, invasion and segregation. *J R Soc Interface*  
445 **9**, 20120448.
- 446 [18] Jin W, Shah ET, Penington CJ, McCue SW, Maini PK, Simpson MJ. 2017. Logistic proliferation  
447 of cells in scratch assays is delayed. *Bull Math Biol* **79**, 1028–1050.
- 448 [19] Johnston ST, Simpson MJ, McElwain DLS, Binder BJ, Ross JV. 2014. Interpreting scratch  
449 assays using pair density dynamics and approximate Bayesian computation. *Open Biol* **4**, 140097.
- 450 [20] Johnston ST, Shah ET, Chopin LK, McElwain DLS, Simpson MJ. 2015. Estimating cell diffu-  
451 sivity and cell proliferation rate by interpreting IncuCyte ZOOM™ assay data using the Fisher-  
452 Kolmogorov model. *BMC Syst Biol* **9**, 38.
- 453 [21] Johnston ST, Ross JV, Binder BJ, McElwain DLS, Haridas P, Simpson MJ. 2016. Quantifying  
454 the effect of experimental design choices for in vitro scratch assays. *J Theor Biol* **400**, 19–31.
- 455 [22] Kaighn ME, Narayan KS, Ohnuki Y, Lechner JF, Jones LW. 1979. Establishment and charac-  
456 terization of a human prostatic carcinoma cell line (PC-3). *Invest Urol* **17**, 16–23.
- 457 [23] Law R, Murrell DJ, Dieckmann U. 2003. Population growth in space and time: Spatial logistic  
458 equations. *Ecology* **84**, 252–262.
- 459 [24] Liepe J, Kirk P, Filippi S, Toni T, Barnes CP, Stumpf MPH. 2014. A framework for parameter  
460 estimation and model selection from experimental data in systems biology using approximate  
461 Bayesian computation. *Nat Protoc* **9**, 439–456.
- 462 [25] Maclaren OJ, Byrne HM, Fletcher AG, Maini PK. 2015. Models, measurement and inference in  
463 epithelial tissue dynamics. *Math Biosci Eng* **12**, 1321.
- 464 [26] Maini PK, McElwain DLS, Leavesley DI. 2004. Traveling wave model to interpret a wound-  
465 healing cell migration assay for human peritoneal mesothelial cells. *Tissue Eng* **10**, 475–482.



- [27] Mathworks. 2017. Kernel smoothing function estimate for univariate and bivariate data.  
http://www.mathworks.com/help/stats/ksdensity.html. Accessed: October 2017.
- [28] Murray JD. 2002. *Mathematical Biology*. Springer, Berlin.
- [29] Peirce SM, Van Gieson EJ, Skalak TC. 2004. Multicellular simulation predicts microvascular  
patterning and in silico tissue assembly. *FASEB J* **18**, 731–733.
- [30] Plank MJ, Simpson MJ. 2012. Models of collective cell behaviour with crowding effects: com-  
paring lattice-based and lattice-free approaches. *J R Soc Interface* **9**, 2983–2996.
- [31] Read M, Andrews PS, Timmis J, Kumar V. 2012. Techniques for grounding agent-based sim-  
ulations in the real domain: a case study in experimental autoimmune encephalomyelitis. *Math  
Comp Model Dyn* **18**, 67–86.
- [32] Sarapata EA, de Pillis LG. 2014. A comparison and catalog of intrinsic tumor growth models.  
*Bull Math Biol* **76**, 2010–2024.
- [33] Sengers BG, Please CP, Oreffo ROC. 2007. Experimental characterization and computational  
modelling of two-dimensional cell spreading for skeletal regeneration. *J R Soc Interface* **4**, 1107.
- [34] Sharkey KJ, Fernandez C, Morgan KL, Peeler E, Thrush M, Turnbull JF, Bowers RG. 2006.  
Pair-level approximations to the spatio-temporal dynamics of epidemics on asymmetric contact  
networks. *J Math Biol* **53**, 61–85.
- [35] Sharkey KJ. 2008. Deterministic epidemiological models at the individual level. *J Math Biol* **57**,  
311–331.
- [36] Sherratt JA, Murray JD. 1990. Models of epidermal wound healing. *P Roy Soc Lond B* **241**, 29.
- [37] Schnoerr, D, Grima, R, Sanguinetti, G. 2016. Cox process representation and inference for  
stochastic reaction-diffusion processes. *Nat Commun* **7**, 11729.
- [38] Schnoerr, D, Sanguinetti, G, Grima, R. 2017. Approximation and inference methods for stochas-  
tic biochemical kineticsa tutorial review. *J Phys A* **50**, 093001.

- 490 [39] Simpson MJ, Landman KA, Hughes BD, Newgreen DF. 2006. Looking inside an invasion wave  
491 of cells using continuum models: Proliferation is the key. *J Theor Biol* **243**, 343–360.
- 492 [40] Stichel D, Middleton AM, Müller BF, Depner S, Klingmüller U, Breuhahn K, Matthäus F.  
493 2017. An individual-based model for collective cancer cell migration explains speed dynamics and  
494 phenotype variability in response to growth factors. *NPJ Syst Biol Appl* **3**, 5.
- 495 [41] Tanaka MM, Francis AR, Luciani F, Sisson SA. 2006. Using approximate Bayesian computation  
496 to estimate tuberculosis transmission parameters from genotype data. *Genetics* **173**, 1511–1520.
- 497 [42] Tang L, van de Ven AL, Guo D, Andasari V, Cristini V, Li KC, Zhou X. 2014. Computational  
498 modeling of 3D tumor growth and angiogenesis for chemotherapy evaluation. *PLoS One* **9**, e83962.
- 499 [43] Treloar KK, Simpson MJ, Haridas P, Manton KJ, Leavesley DI, McElwain DLS, Baker RE. 2013.  
500 Multiple types of data are required to identify the mechanisms influencing the spatial expansion  
501 of melanoma cell colonies. *BMC Syst Biol* **7**, 137.
- 502 [44] Tremel A, Cai A, Tirtaatmadja N, Hughes BD, Stevens GW, Landman KA, O’Connor AJ. 2009.  
503 Cell migration and proliferation during monolayer formation and wound healing. *Chem Eng Sci*  
504 **64**, 247–253.
- 505 [45] Warne DJ, Baker RE, Simpson MJ. 2017. Optimal quantification of contact inhibition in cell  
506 populations. In press *Biophysical Journal* <https://doi.org/10.1016/j.bpj.2017.09.016>