

1 Using approximate Bayesian computation to quantify cell-cell
2 adhesion parameters in a cell migratory process

3 Robert J. H. Ross ^{*1}, R. E. Baker ^{†1}, Andrew Parker ^{‡1}, M. J. Ford ^{§2}, R. L. Mort ^{¶3},
4 and C. A. Yates ^{||4}

5 ¹Wolfson Centre for Mathematical Biology, Mathematical Institute, University of
6 Oxford, Radcliffe Observatory Quarter, Woodstock Road, Oxford, OX2 6GG

7 ²MRC Human Genetics Unit, MRC IGMM, Western General Hospital, University of
8 Edinburgh, Edinburgh, EH4 2XU

9 ³Division of Biomedical and Life Sciences, Faculty of Health and Medicine, Furness
10 Building, Lancaster University, Bailrigg, Lancaster, LA1 4YG

11 ⁴Centre for Mathematical Biology, Department of Mathematical Sciences, University of
12 Bath, Claverton Down, Bath, BA2 7AY

13 January 13, 2017

14 **Abstract**

15 In this work we implement approximate Bayesian computational methods to improve the
16 design of a wound-healing assay used to quantify cell-cell interactions. This is important
17 as cell-cell interactions, such as adhesion and repulsion, have been shown to play a role in
18 cell migration. Initially, we demonstrate with a model of an *unrealistic* experiment that we
19 are able to identify model parameters that describe agent motility and adhesion, given we
20 choose appropriate summary statistics for our model data. Following this, we replace our

*ross@maths.ox.ac.uk (corresponding author, +44(0)7766721659)

†baker@maths.ox.ac.uk

‡parker@maths.ox.ac.uk

§matthew.ford@ed.ac.uk

¶r.mort@lancaster.ac.uk

||c.yates@bath.ac.uk

21 model of an unrealistic experiment with a model representative of a practically realisable
22 experiment. We demonstrate that, given the current (and commonly used) experimental
23 set-up, our model parameters cannot be accurately identified using approximate Bayesian
24 computation methods. We compare new experimental designs through simulation, and show
25 more accurate identification of model parameters is possible by expanding the size of the
26 domain upon which the experiment is performed, as opposed to increasing the number of
27 experimental replicates. The results presented in this work therefore describe time *and*
28 cost-saving alterations for a commonly performed experiment for identifying cell motility
29 parameters. Moreover, this work will be of interest to those concerned with performing
30 experiments that allow for the accurate identification of parameters governing cell migratory
31 processes, especially cell migratory processes in which cell-cell adhesion or repulsion are
32 known to play a significant role.

33 **Keywords:** Cell migration, adhesion, wound-healing, summary statistics, parameter identifica-
34 tion, experimental design, approximate Bayesian computation, agent-based model, simulation.

35
36 **Corresponding author:** Robert J. H. Ross, ross@maths.ox.ac.uk, +44(0)7766721659, Wolf-
37 son Centre for Mathematical Biology, Mathematical Institute, University of Oxford, Radcliffe
38 Observatory Quarter, Woodstock Road, Oxford, OX2 6GG.

39 1 Introduction

40 Cell-cell interactions are known to play an important role in several cell migration processes.
41 For example, multiple different cell-cell interactions, such as cell-cell signalling and cell-cell ad-
42hesion [1], have been identified as promoting metastasis in breast cancer. Repulsive interactions
43 mediated via ephrins on the surface of neural crest stem cells are known to coordinate the early
44 stages of melanoblast migration away from the neural tube [2]. More fundamentally, it is hy-
45pothesised that the emergence of cell-cell interactions over one billion years ago helped establish
46 the necessary conditions for multicellular organisms [3].

47
48 A well-established approach for studying cell migration is to construct an agent-based model
49 (ABM) to simulate the cell migratory process of interest [4-8]. Typically, this involves using
50 a computational model to simulate a population of agents on a two-dimensional surface, or in

51 a three-dimensional volume. The agents in the ABM represent cells, and each agent is able to
52 move and interact with other agents in the ABM. In this work we use an ABM to simulate a
53 wound-healing assay¹, an experiment commonly used for studying cell motility [9–15]. Other
54 modelling approaches apart from ABMs have been employed to study wound-healing. For in-
55 stance, a huge amount of research has been completed using continuum methods to model the
56 wound-healing process (see Flegg et al. [16] for a recent review of the field). However, we employ
57 an ABM in this work because they provide an intuitive representation of cells, and allow for
58 complex behaviours representing biological processes, such as cell-cell interactions and volume
59 exclusion, to be easily assigned to agents in the ABM.

60

61 If an ABM is an *effective*² representation of a cell migration process it can be used for a
62 number of purposes. One such purpose for an ABM is to perform *in silico* experiments to
63 test scientific hypotheses. For instance, a recent study used an ABM to demonstrate that a
64 simple mechanism of undirected cell movement and proliferation could account for neural crest
65 stem cell colonisation of the developing epidermis in the embryonic mouse [4]. Other studies
66 involving ABMs have tested hypotheses concerning the influence of matrix stiffness and matrix
67 architecture on cell migration [17], and the mechanism by which cranial neural crest stem cells
68 become ‘leaders’ or ‘followers’ in the embryonic chick to facilitate their collective migration [6–8].

69

70 ABMs can also be used to *identify* parameters in experimental data (with the caveat that
71 the parameters are model-dependent). The reasoning behind using an ABM to identify pa-
72 rameters in experimental data is as follows: if an ABM is an effective representation of an
73 experiment, then the parameter values the ABM requires to reproduce the experimental data
74 may be representative of the parameter values in the biological process that is the focus of
75 the experiment. For instance, the value of a parameter that describes cell proliferation rate.
76 Even if the parameter values in the parameterised ABM are not representative of the parameter
77 values in the biological process, the parameterised ABM may still be used to make predictions
78 about the process of interest by performing *in silico* experiments, as described above. These
79 predictions can then be experimentally tested.

¹Wound-healing assays are also often referred to as scratch assays.

²By an effective representation we mean the ABM captures the salient features of the process of interest, and is therefore a viable research tool with which to study the process of interest.

80

81 Alternatively, if the ABM is an effective representation of an experiment (i.e. the experimental
82 data can be reproduced), but the parameters of the ABM are not identifiable, this may suggest
83 the experiment is not well-designed (that is, if the experiment has been designed to estimate
84 parameters). By parameters not being identifiable we mean that different parameter values in
85 the ABM can reproduce the same experimental data. If this is the case, the ABM can then be
86 used to suggest improvements to the experiment's design, namely by altering the ABM design
87 such that the ABM parameters become identifiable. These alterations can then be applied
88 to the experiment to improve parameter identifiability. For example, a recent study using an
89 ABM has examined the time-points at which data should be collected from an experiment to
90 maximise the identifiability of ABM parameters [11]. Other theoretical work has shown how to
91 maximise the information content of an experiment by choosing an appropriate experimental
92 set-up [18].

93

94 The focus of our study is to determine the experimental conditions, and experimental data,
95 required for the accurate identification of cell motility and adhesion parameters in an ABM of a
96 wound-healing assay. To do so we employ approximate Bayesian computation (ABC), a proba-
97 bilistic approach whereby a probability distribution for the parameter(s) of interest is estimated,
98 as opposed to a point estimate [10, 19, 20]. Although ABC is well-established in some fields, for
99 instance in population genetics [21], its applicability for ABMs representing cell migration is still
100 an area of active research [9–11, 22–24]. Recent studies combining ABC and ABMs have been
101 able to identify motility and proliferation rates in cell migratory processes [10], and improve the
102 experimental design of scratch assays [11]. However, as far as we are aware no study to date has
103 used ABC methods to examine the experimental conditions, and experimental data, required
104 for the accurate identification of cell motility and adhesion parameters in a wound-healing assay.

105

106 Other methods to identify parameters from experimental data using ABMs also exist. For
107 instance, a standard approach is to generate point estimates of model parameters that best re-
108 produce statistics of the experimental data in the ABM. For example, the generation of motility
109 and proliferation rates for agents in an ABM representing a biological process [4]. This approach,

110 while applicable in some circumstances, often gives little insight into how much uncertainty ex-
111 ists in the parameters chosen, a factor that can be of importance when analysing biological
112 systems. For example, relationships between parameter uncertainty and system robustness are
113 thought to be connected in biological function at a systems level [25].

114

115 The outline of this work is as follows: in Section 2 we introduce our ABM and define the
116 cell-cell interactions we implement. We also outline the method of ABC, and the summary
117 statistics we use to analyse the ABM output. In Section 3 we present results and demonstrate
118 that, given an ABM representing an unrealistic experiment, we are able to identify ABM pa-
119 rameters for agent motility and adhesion. Following this, we replace our ABM representing
120 an unrealistic experiment with an ABM that simulates a practically realisable experiment. In
121 doing so we show that agent motility and adhesion parameters cannot be successfully identi-
122 fied using ABC given the current experimental design. To improve parameter identifiability
123 we compare different experimental set-ups, and show that identification of ABM parameters
124 is made more accurate if the size of the domain upon which the experiment is performed is
125 expanded, as opposed to the number of experimental replicates increased. Experimentally, ex-
126 panding the size of the domain is equivalent to increasing the field of view of the microscope
127 used to collect the experimental data. For instance, generating five experimental replicates on
128 a larger domain enables more accurate identification of ABM parameters than generating 500
129 experimental replicates on a domain eight times smaller. In Section 4 we discuss the results
130 presented in this work.

131 **2 Methods**

132 In this section we first introduce the ABM. We then define our summary statistics and explain
133 ABC and its implementation.

134 **2.1 Agent-based model**

135 An ABM is a computational model for simulating the behaviour of autonomous agents. The
136 agents in the ABM represent cells, and each agent is able to move and interact with other
137 agents. The ABM is simulated on a two-dimensional square lattice with lattice spacing Δ [26]

138 and size L_x by L_y , where L_x is the number of lattice sites in each row, and L_y is the number of
 139 sites in each column. Each agent is initially assigned to a lattice site, from which it can move
 140 into adjacent sites. If an agent attempts to move into a site that is already occupied by another
 141 agent, the movement event is aborted. Processes such as this whereby one agent is allowed per
 142 site are often referred to as exclusion processes [26]. In the ABM time evolves continuously,
 143 and as our ABM can be modelled as a continuous-time Markov process we use the Gillespie
 144 algorithm [27] to generate sample paths. Attempted agent movement events occur with rate
 145 P_m per unit time. $P_m \delta t$, therefore, is the probability of an agent attempting to move in the
 146 next infinitesimally small time interval δt . In our ABM a lattice site is denoted by $v = (i, j)$,
 147 where i indicates the column number and j the row number. Each lattice site has four adjacent
 148 lattice sites (except for those sites situated on nonperiodic boundaries), and so the number of
 149 nearest neighbour lattice sites that are occupied by an agent, denoted by n , is $0 \leq n \leq 4$. We
 150 denote the set of unoccupied nearest neighbour lattice sites by \mathcal{U} .

151

152 The ABM domain size for simulations representing unrealistic experiments is $L_x = 100$ by
 153 $L_y = 100$, and the lattice sites indexed by $1 \leq j \leq L_y$ and $1 \leq i \leq 10$, and $1 \leq j \leq L_y$ and
 154 $91 \leq i \leq L_x$ are initially occupied by agents. In Fig. 1 the initial conditions in the ABM for the
 155 unrealistic experiment can be seen. The initial condition in Fig. 1 represents a ‘wound’, in that
 156 agents are positioned either side of a space, the ‘wound’, that they can migrate into. The agent
 157 migration into this space simulates one aspect of the wound-healing process. We refer to this
 158 simulation as unrealistic because the uniformity of the initial conditions would not be possible
 159 in a realistic experimental setting. The initial condition is also improved from our experimen-
 160 tally realisable simulation as it is ‘double-sided’, as opposed to the ‘single-sided’ experimental
 161 data that we will later simulate for our ABM of a realistic experiment. It has been shown that
 162 double-sided initial conditions can provide more information than single-sided initial conditions
 163 for some model parameters [11]. For instance, double-sided initial conditions can improve pa-
 164 rameter identifiability if increasing the number of agents in a simulation improves parameter
 165 identifiability. For the ABM of an unrealistic experiment all simulations have periodic bound-
 166 ary conditions at the top and bottom of the domain (i.e. for lattice sites indexed by $j = 1$ or
 167 $j = L_y$), and no-flux boundary conditions at the left-hand and right-hand boundaries of the

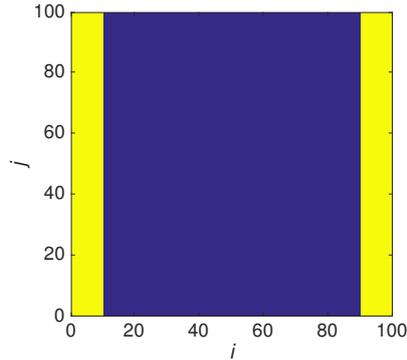


Figure 1: The initial condition in the ABM for the unrealistic experiment. Yellow indicates a site occupied by an agent and blue indicates an empty lattice site.

168 domain (i.e. for lattice sites indexed by $i = 1$ or $i = L_x$).

169

170 It is important to stress that throughout this work we assume that cellular processes such
 171 as migration have constant parameter values associated with them. Inference procedures do ex-
 172 ist in which the parameter values associated with cell processes are not assumed to be constant,
 173 but are instead treated as a random variable sampled from a distribution. These methods are
 174 often important for sensitivity analysis, or if the data is sampled from a heterogeneous pop-
 175 ulation [28–30]. However, we do not implement these methods in this work as it would serve
 176 to prematurely complicate our research question. It is also important to acknowledge that in
 177 migrating cell populations there are often many more factors at play than simply cell motility
 178 and adhesion. For instance, the cell cycle and a cell’s response to environmental cues may
 179 be important factors in a cell’s behaviour. Again, however, we have purposely simplified our
 180 model to first ascertain if we can accurately estimate parameters associated with cell motility
 181 and adhesion.

182 2.2 Cell-cell adhesion models

183 In the ABM cell-cell interactions are simulated by altering the probability of an agent attempting
 184 to move, depending on the number of nearest occupied neighbours, n , an agent has. We employ
 185 two models to simulate cell-cell interactions in the ABM, one of which has been published
 186 before [13, 31]. We define $T(v'|v)$ as the transition probability that an agent situated at site v ,
 187 having been selected to move, attempts to move to site v' , where v' indicates one of the nearest

188 neighbour sites of v . Therefore, $T(v'|v)$ is only non-zero if v and v' are nearest neighbours. The
 189 transition probability in the first model, which we refer to as model A, is defined as

$$190 \quad T_A(v'|v) = \frac{1 - n\alpha}{4}, \quad (1)$$

192 where α is the adhesion parameter. The subscript A on the transition probability in Eq. (1)
 193 indicates that this is the transition probability for model A. If $\alpha > 0$ Eq. (1) models cell-cell
 194 adhesion, and if $\alpha < 0$ Eq. (1) models cell-cell repulsion. The transition probabilities stated in
 195 Eq. (1) must satisfy

$$196 \quad 0 \leq \sum_{v' \in \mathcal{U}} T_A(v'|v) \leq 1. \quad (2)$$

198 Inequality (2) ensures the probability of an agent, if selected to move, attempting to move to
 199 any of its unoccupied nearest neighbour sites never exceeds unity, and so constrains the value α
 200 can take. The transition probability in the second model, which we refer to as model B [13, 31],
 201 is defined as

$$202 \quad T_B(v'|v) = \frac{(1 - \alpha)^n}{4}, \quad (3)$$

204 and must satisfy

$$205 \quad 0 \leq \sum_{v' \in \mathcal{U}} T_B(v'|v) \leq 1. \quad (4)$$

207 As in model A if $\alpha > 0$ Eq. (3) models cell-cell adhesion, and if $\alpha < 0$ Eq. (3) models cell-cell
 208 repulsion.

209
 210 Models A and B simulate different types of cell-cell adhesion. In model A the transition proba-
 211 bility is a linear function of n . Meanwhile, in model B the transition probability is a nonlinear
 212 function of n . Not only may these different types of cell-cell adhesion be relevant for different
 213 cell types, but implementing two models of cell-cell adhesion allows us to test the robustness of
 214 the methods we present in this work for identifying cell-cell adhesion parameters.

215 **2.3 Summary statistics**

216 Summary statistics are lower-dimensional summaries of data that provide a tractable means to
217 compare different sets of data. Summary statistics are important because experimental data is
218 often of high dimensionality, and if we want to use experimental data to efficiently guide com-
219 putational algorithms we require ways to accurately summarise it. We now define the summary
220 statistics we apply to the ABM output and experimental data. Following this we describe how
221 we utilise these summary statistics to implement ABC.

222

223 We initially use three summary statistics to evaluate the ABM output, all of which have been
224 considered previously [9, 31, 32]. Our aim is to ascertain which summary statistic (or combina-
225 tion of summary statistics) is most effective for the identification of agent motility and adhesion
226 parameters in the ABM.

227 **Average horizontal displacement of agents**

228 The average horizontal displacement of all agents, \bar{i} , in a given time interval, $[t_i, t_f]$, in the ABM
229 is calculated as

$$230 \quad \bar{i} = \frac{1}{N} \sum_{k=1}^N |i_{t_i}^k - i_{t_f}^k|, \quad (5)$$

231

232 where \bar{i} is the average horizontal displacement of agents, N is the total number of agents in the
233 simulation, $i_{t_i}^k$ is the column position of agent k at time t_i , and $i_{t_f}^k$ is the column position of
234 agent k at time t_f . We only look at the horizontal displacement of agents as this is the direction
235 in which the majority of agent displacement occurs, due to the initial conditions of the ABM
236 (Fig. 1). It has previously been shown that different cell-cell interactions have different effects
237 on the average displacement of agents in an ABM [31]. As may be expected, repulsive (adhesive)
238 interactions between agents tend to increase (decrease) the average displacement of agents, and
239 so the average displacement of agents may be a useful summary statistic for distinguishing
240 between repulsive and adhesive cell-cell interactions in the ABM.

241 **Agent density profile**

242 The agent density profile at time t in the ABM is calculated as

$$243 \quad C_t(i) = \frac{1}{L_y} \sum_{j=1}^{L_y} \mathbb{1}\{v\}. \quad (6)$$

244

245 Here $C_t(i)$ is the agent density profile and $\mathbb{1}$ is the indicator function for the occupancy of a
 246 lattice site v (i.e. 1 if an agent occupies lattice site v , and 0 if it is not occupied by an agent).
 247 We have shown previously that different cell-cell interactions have different effects on the agent
 248 density profile [31]. For instance, repulsive interactions between agents can create a concave
 249 agent density profile, whereas adhesive interactions between agents can create a convex agent
 250 density profile. Therefore, the agent density profile may be an effective summary statistic for
 251 distinguishing between repulsive and adhesive cell-cell interactions in the ABM.

252 **Pairwise-correlation function**

253 The final summary statistic we consider is the pairwise-correlation function (PCF). The PCF
 254 provides a measure of the spatial clustering between agents in an ABM, and has been used
 255 frequently in the analysis of cell migratory processes [4, 9, 33, 34]. The PCF has also been
 256 successfully used as a summary statistic for the parameterisation of ABMs of cell migration
 257 [10]. We use i_t^k to denote the column position of agent k at time t , i_t^l to denote the column
 258 position of agent l at time t , and define $c_t(m)$ to be the number of occupied pairs of lattice sites
 259 for each *nonperiodic*³ horizontal pair distance $m = 1, \dots, L_x - 1$ at time t . This means $c_t(m)$ is
 260 given by

$$261 \quad c_t(m) = \sum_{k=1}^N \sum_{l=k+1}^N \mathbb{1}\{|i_t^k - i_t^l| = m\}, \quad \forall m = 1, \dots, L_x - 1, \quad (7)$$

262

263 where $\mathbb{1}$ is the indicator function equal to unity if $|i_t^k - i_t^l| = m$, and is equal to zero otherwise.
 264 In Eq. (7) only the pair agent distances in the horizontal direction are counted. Given the
 265 translational invariance of the initial conditions in the vertical direction of the ABM, the ma-
 266 jority of important spatial information will be in the horizontal direction⁴. Binder and Simpson

³By nonperiodic it is meant the distance measured between two agents cannot cross the ABM boundary.

⁴This approach is in agreement with previous studies [34], which showed the most relevant information from the PCF summary statistic is perpendicular to the wound axis in a wound-healing assay.

267 [34] demonstrated that is necessary to normalise Eq. (7) to account for volume exclusion. The
 268 normalisation term is

$$269 \hat{c}_t(m) = L_y^2(L_x - m)\rho\hat{\rho}, \quad \forall m = 1, \dots, L_x - 1, \quad (8)$$

271 where $\rho = N/(L_x L_y)$, and $\hat{\rho} = (N - 1)/(L_x L_y - 1)$. Equation (8) describes the expected
 272 number of pairs of occupied lattice sites, for each nonperiodic horizontal pair distance, m , in a
 273 population distributed uniformly at random on the domain. Combining Eqs. (7) and (8), the
 274 PCF is

$$275 q_t(m) = \frac{c_t(m)}{\hat{c}_t(m)}, \quad (9)$$

277 where $q_t(m)$, the PCF, is a measure of how far $c_t(m)$ departs from describing the expected
 278 number of occupied lattice pairs for each horizontal distance of an agent population spatially
 279 distributed uniformly at random on the ABM domain.

280

281 It is important to briefly discuss why we chose these summary statistics and not others that
 282 have also been used to analyse cell migration [10, 22, 24]. Other summary statistics were ini-
 283 tially implemented in this study, such as the concavity of agent trajectories, the *total* distance
 284 travelled by agents, and the leading edge of the agent population. However, these summary
 285 statistics were found not to be informative for the identification of agent motility and adhesion
 286 parameters in our ABM, and so were excluded from this work. The three summary statistics we
 287 implement are encapsulated in Table 1 for the reader's convenience, in addition to the properties
 288 each summary statistic summarises in the agent population.

289

290 2.4 Approximate Bayesian computation

291 Here we introduce our ABC algorithm [19]. We define M as a stochastic model that takes
 292 parameters Θ and produces data D . This relationship can be written as $D \sim M(\Theta)$. For the
 293 ABM presented in this work $\Theta = (P_m, \alpha)$, where Θ is sampled from a prior distribution, π , and
 294 so this relationship can be written as $\Theta \sim \pi$. The relationship between π and Θ is often written

Summary statistic	Description
Average horizontal displacement of agents	<p>Summarises the displacement of agents into the ‘wound’. This displacement is affected by the adhesion of agents and their motility rate. Mathematically the average horizontal displacement of agents is defined as</p> $\bar{i} = \frac{1}{N} \sum_{k=1}^N i_{t_i}^k - i_{t_f}^k .$
Agent density profile	<p>Summarises the macroscopic shape of the population as it moves into the ‘wound’. We have previously shown this shape is partly determined by agent interactions and motility [31]. Mathematically the agent density profile is defined as</p> $C_t(i) = \frac{1}{L_y} \sum_{j=1}^{L_y} \mathbb{1}\{v\}.$
Pair-wise correlation function	<p>Summarises the spatial correlations/structure established by agent movement and interactions. Mathematically the pair-wise correlation function is defined as</p> $q_t(m) = \frac{c_t(m)}{\hat{c}_t(m)}.$

Table 1: The summary statistics we implement and the properties of the agent population they summarise.

295 as $\Theta \sim \pi(\Theta)$, which indicates that a new Θ sampled from the prior distribution may depend on
296 the previous Θ . This relationship will be relevant later on in this work, however, initially each
297 Θ sampled from π is independent of the previous Θ .

298

299 The identification of ABM parameters in this work centres around the following problem: given
300 a stochastic model, M , and data, D , what is the probability density function that describes Θ
301 being the model parameters that produced data D ? More formally, we seek to obtain a poste-
302 rior distribution, $p(\Theta|D)$, which is the conditional probability of Θ given D (and the model, M).

303

304 Typically, to compute the posterior distribution a likelihood function, $L(D|\Theta)$, is required.

305 This is because the likelihood function and posterior distribution are related in the following

306 manner by Bayes' theorem:

$$307 \quad p(\Theta|D) \propto L(D|\Theta)\pi(\Theta). \quad (10)$$

308

309 That is, the posterior distribution is proportional to the product of the likelihood function and
310 the prior distribution. Approximate Bayesian computation is a well-known method for esti-
311 mating posterior distributions of model parameters in scenarios where the likelihood function
312 is *intractable* i.e. it is impossible or computationally prohibitive to obtain [19].

313

314 In many cases for ABC, due to the high dimensionality of the data, D , it is necessary to
315 utilise a summary statistic, $S = S(D)$. The summary statistics we employ in this work are
316 of varying dimension. For instance, the agent density profile at time t has L_x data points,
317 whereas the average agent displacement at time t has one data point. Therefore we write $S(D)$
318 as $S(D)_{r,t}$, where $S(D)_{r,t}$ is the r^{th} data point in the summary statistic at the t^{th} sampling time.

319

320 The ABC method proceeds in the following manner: we wish to estimate the posterior dis-
321 tribution of Θ given D . We now simulate model M with parameters Θ , sampled from π , and
322 produce data \tilde{D} . We calculate the difference between a summary statistic applied to D and \tilde{D}
323 with

$$324 \quad d = \sum_{t=1}^T \sum_{r=1}^R |S(D)_{r,t} - S(\tilde{D})_{r,t}|, \quad (11)$$

325

326 where R is the number of data points in $S(D)$ and T is the number of sampling times. We
327 repeat the above process many times, that is, sample Θ from π , produce \tilde{D} , calculate d with
328 Eq. (11), and only accept Θ for which d is below a user defined certain threshold (alternatively,
329 a predefined number of Θ that minimise d can be accepted). This enables us to generate a
330 distribution for Θ that is an approximation of the posterior distribution, $p(\Theta|D)$, given M [35].
331 More specific details of the ABC algorithms we implement are introduced when necessary in
332 the text.

333 **3 Results**

334 We begin by demonstrating that for an ABM representing an unrealistic experiment we are able
335 to identify model parameters, given appropriate summary statistics.

336 **3.1 Unrealistic experiment**

337 To ascertain the effectiveness of the chosen summary statistics to identify model parameters, we
338 attempt to identify Θ from data generated *synthetically*. Synthetic data is ABM data generated
339 with fixed parameter values, and so can be thought of as a simulation equivalent of experimental
340 data. To generate the synthetic data using the ABM we proceed as follows:

- 341 1. We choose parameters Θ to identify. To help clarify this explanation let us make these
342 parameters $\Theta = (P_m, \alpha) = (0.5, 0.1)$ in model A⁵.
- 343 2. For model A we perform a simulation of the ABM with $\Theta = (0.5, 0.1)$, generate data,
344 D , and calculate summary statistics, $S(D)$, from the simulation at our time-points of
345 interest. These times are $t = [240, 480, 720]$. We choose these times as they are the
346 times (in minutes) we will later analyse for the simulations of the practically realisable
347 experiment, and correspond to 4 hours, 8 hours and 12 hours into an experiment.
- 348 3. We repeat step 2 ten times and calculate the ensemble average for each summary statistic
349 for each individual time-point.

350 This procedure generates synthetic data for which we will now attempt to identify the param-
351 eters. In this work we present representative results using $P_m = 0.5$ and $\alpha = 0.1$ for model A,
352 and $P_m = 0.5$ and $\alpha = 0.25$, and $P_m = 0.5$ and $\alpha = -0.1$ for model B.

353

354 Throughout this work we sample P_m and α for our model from uniform priors. In the case
355 of model A, $P_m \in [0, 1]$ and $\alpha \in [-0.2, 0.25]$, and for model B, $P_m \in [0, 1]$ and $\alpha \in [-0.2, 1.0]$.

356 We stipulate these lower and upper bounds for α for both models A and B to make sure in-
357 equalities (2) and (4) are satisfied.

⁵A value of $P_m = 0.5$, given that the simulation time will later be defined to be in minutes, and the length of a lattice site represents cell length (typically between $10\mu\text{m}$ - $100\mu\text{m}$), means that the motility of the agents is biologically realistic. The parameter α is dimensionless. The experimental realism of these parameters will be expanded on when we address the simulation of a practically realisable experiment.

359 We begin by implementing an ABC rejection algorithm that proceeds as follows:

- 360 1. Run 10^4 ABM simulations, in each case using Θ sampled uniformly at random from the
361 prior distribution.
- 362 2. Compute the distance d as defined in Eq. (11) for simulation times $t = [240, 480, 720]$.
- 363 3. Accept the 100 parameter values, Θ , that give the smallest values of d .

364 In Fig. 2 the posteriors generated using each of the three summary statistics applied to data
365 from simulations of an unrealistic experiment are displayed. The most effective summary statis-
366 tic for identifying the synthetic data parameters is the PCF. This is evident in the location of
367 the posterior distribution density relative to the red dot (the red dot represents the synthetic
368 data parameter values), and the narrow spread of the posterior distribution density as indicated
369 by the scale bar in Fig. 2 (c), (f) and (i). The agent density profile summary statistic performs
370 less well than the PCF for parameter identification, especially for model A (Fig. 2 (b)). In
371 the case of the average agent displacement summary statistic many combinations of P_m and α
372 lead to the same average agent displacement, which results in an extended region of possible
373 parameter values. To some extent this is to be expected, as increasing either P_m or α will have
374 opposing effects on the average agent displacement. This means that using agent displacement
375 as a summary statistic results in parameter identifiability issues in this example.

376

377 To quantify the difference between the performance of the different summary statistics we
378 use the Kullback-Leibler divergence (KLD), which is a measure of the information gained in
379 moving from the prior distribution to the posterior distribution [36]. The KLD for a discrete
380 probability distribution is defined as follows:

$$381 \quad D_{KL}(p|\pi) = \sum_l p(\Theta_l|D) \log \left(\frac{p(\Theta_l|D)}{\pi(\Theta_l)} \right), \quad (12)$$

382

383 where the index l accounts for all possible discretised parameter pairs (i.e. all combinations of
384 P_m and α). A larger $D_{KL}(p|\pi)$ value suggests that more information is obtained (the entropy
385 of the distribution is reduced) when moving from the prior distribution to the posterior distri-
386 bution. However, this does not necessarily mean the posterior distribution is a more accurate

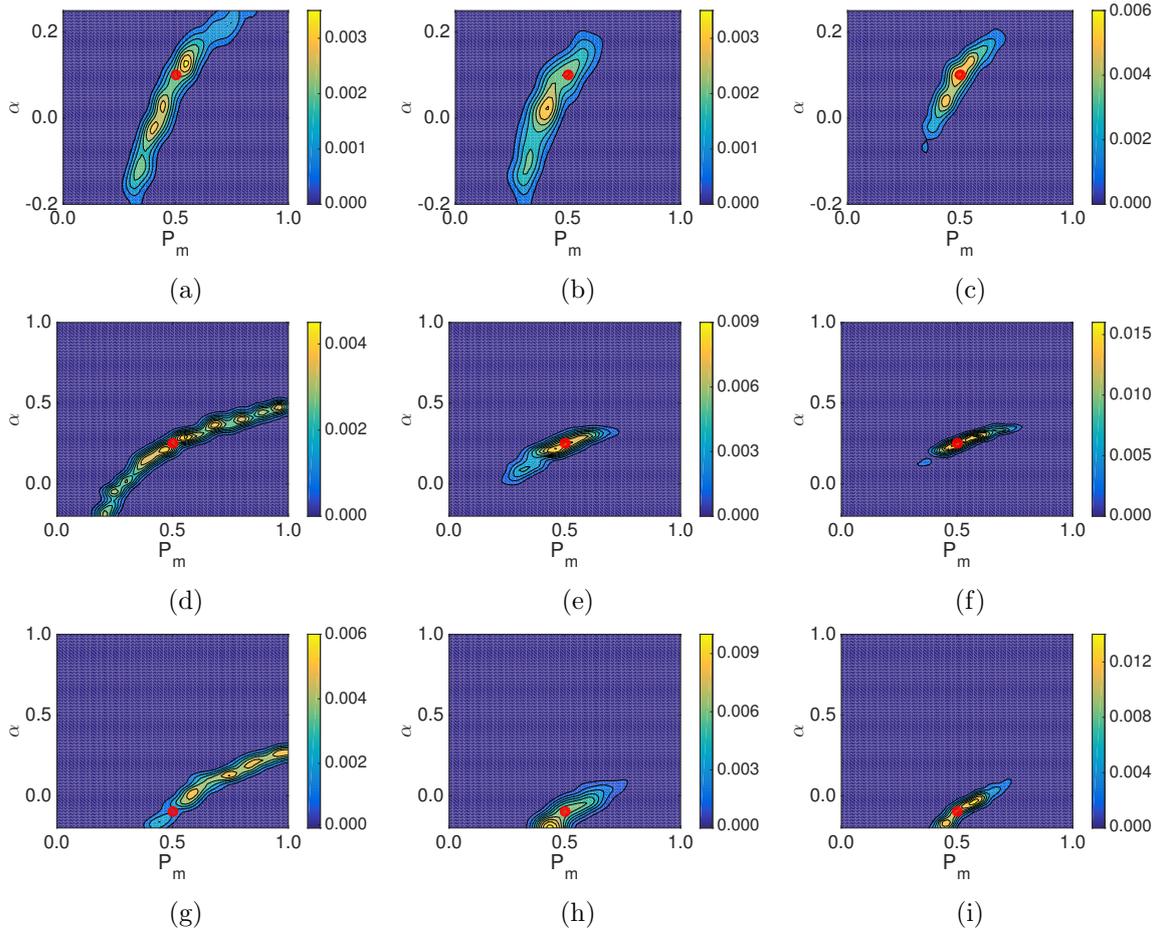


Figure 2: (a)-(c) Posterior distributions for model A for an unrealistic experiment with different summary statistics: (a) average displacement of agents in the horizontal direction; (b) agent density profile; (c) PCF. In all cases the red dot indicates the value of the parameters used to generate the synthetic data, $P_m = 0.5$, $\alpha = 0.1$. As indicated by the colour bar the yellow regions indicate areas of high relative density of the posterior distribution, while the blue regions indicate areas of low relative density of the posterior distribution. (d)-(f) Model B, $P_m = 0.5$, $\alpha = 0.25$: (d) average displacement of agents in the horizontal direction; (e) agent density profile; (f) PCF. (g)-(i) Model B, $P_m = 0.5$, $\alpha = -0.1$: (g) average displacement of agents in the horizontal direction; (h) agent density profile; (i) PCF.

387 representation of the parameter distribution. Therefore, the KLD should not be seen as ubiq-
 388 uitously applicable to inference problems similar to those described in this work. In particular,
 389 the KLD should be used with caution in scenarios in which an informative prior is used. In
 390 such scenarios, other methods to measure the improvement of an inference procedure have been
 391 examined and may be more suitable [37].

392

393 To compute the KLD we discretise our posterior distribution onto a lattice with 2^6 equally
 394 spaced values of P_m and 2^6 equally spaced values of α . Computing $D_{KL}(p|\pi)$ for all nine plots

395 in Fig. 2 gives: (a) 1.77; (b) 1.70; (c) 2.32; and (d) 2.15; (e) 2.57; (f) 3.35; and (g) 2.45; (h)
396 2.72; (i) 3.27. In tandem with the proximity of the peak of the posterior distribution densities
397 to the red dots in Fig. 2 (c), (f) and (i), compared to Fig. 2 (a)-(b), (d)-(e) and (g)-(h), this is
398 increase in the KLD suggests that the PCF summary statistic is more effective for parameter
399 identification than the average agent displacement and agent density profile summary statistics.

400 **3.2 Practically realisable experiment**

401 In the previous section we demonstrated that for unrealistic experimental conditions the PCF
402 summary statistic is best able to identify synthetic data parameters (for data generated from an
403 ABM of an unrealistic experiment), and so moving forward we will only use the PCF summary
404 statistic for parameter identification. Previous work has combined summary statistics to im-
405 prove parameter identification, and how best to combine summary statistics has been the focus
406 of a significant amount of research, with a wide range of different methods examined [10, 37–40].
407 However, in this case combining our summary statistics results in a negligible improvement to
408 the posterior distribution⁶.

409
410 We now replace our ABM that represents an unrealistic experiment with an ABM that rep-
411 resents an actual experiment, and examine if synthetic data parameters can be identified in
412 the ABM. That is, from this point on, we generate all synthetic data from an ABM based on
413 a realistic experimental set-up. We provide brief details of the experiment here, however, a
414 more detailed description can be found in the supplementary material (Section S2). In Fig. 3
415 a typical initial frame of the experimental data can be seen.

416
417 In total we have data from five replicates of the experiment. Therefore, we now generate
418 our synthetic data from five replicates of the ABM, using the same procedure as described in
419 Section 3.1. One key difference between the unrealistic and practically realisable experiments
420 is the size of the domain and, because of this, the number of agents in a simulation.

421
422 The experimental images were captured by a microscope with a field of view of 597.24 μm

⁶An example of a posterior distribution generated by combining all three summary statistics can be found in the supplementary material (Section S1).

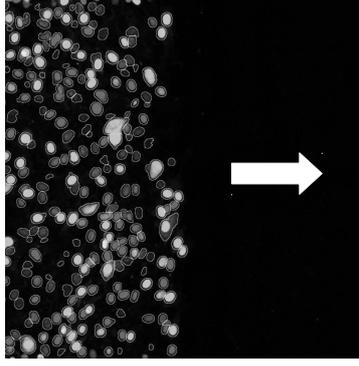


Figure 3: Typical initial frame of the experimental data. The cells are positioned such that they will migrate primarily horizontally into the space without cells, this space represents a wound (the direction of migration is indicated by the white arrow).

423 by $597.24 \mu\text{m}$. The cell size in the experimental images is consistent with each cell occupying a
 424 $26 \mu\text{m}$ by $26 \mu\text{m}$ square lattice site. Given the size of the microscope field of view this means
 425 the ABM domain size is $L_x = 23$ by $L_y = 23$. We use the average initial conditions from the
 426 experiment to generate the initial conditions in the ABM of a realistic experiment. Exact details
 427 of how the initial condition is generated in the ABM, and how experimental data is mapped to
 428 a lattice, can be found in the supplementary material (Section S3).

429

430 We also alter the ABM to have flux (nonperiodic) boundary conditions at the left-hand and
 431 right-hand boundaries of the domain (i.e. for lattice sites with $j = 1$ or $j = N_y$). The left-most
 432 column is kept at or above a constant density throughout the simulation time course. That is,
 433 after any movement event from the left-most column in the simulation the column density of
 434 the left-most column is calculated, and if found to be below a certain density agents are added
 435 to empty sites in this column chosen uniformly at random until the required density is achieved.
 436 This mechanism ensures that the agent density profile in the ABM replicates the evolution of
 437 the experimental data throughout the simulation. Further details regarding the implementation
 438 of this boundary condition are provided in the supplementary material (Section S3). The top
 439 and bottom boundaries of the ABM domain remain periodic as cells were seen to move in and
 440 out of the microscope field at these boundaries in the experimental images, at an approximately
 441 equal rate.

442

443 To reduce computational time we now implement a Markov Chain Monte Carlo variant of

444 ABC [19]. Details of the implementation of the algorithm are given in the supplementary ma-
445 terial (Section S4). As before we sample from uniform priors $P_m \in [0, 1]$ and $\alpha \in [-0.2, 0.25]$
446 for model A, and $P_m \in [0, 1]$ and $\alpha \in [-0.2, 1.0]$ for model B, and collect simulation data at
447 $t = [240, 480, 720]$. We collect simulation data at three time-points so that the computational
448 time is of practical length (our longest ABC Markov Chain Monte Carlo implementations took
449 approximately 192 hours). A value of $P_m = 0.5$, given that the simulation time is in minutes,
450 and the length of a lattice site is 26 μm , means that the motility of the agents is biologically
451 realistic. To be precise, the agents here are approximately five times faster than cell motility
452 rates previously published [4, 9]⁷. However, the cells considered in [4, 9] are not thought to
453 exhibit cell-cell adhesion, and so a higher motility rate for the agents is sensible as agent move-
454 ment is reduced by cell-cell adhesion in our ABM.

455

456 In Fig. 4 it can be seen that the synthetic data parameters cannot be accurately identified
457 using ABC, with the PCF summary statistic, given the current ABM design. This is evident
458 in the location of the red dots (indicating the parameter values used to generate the synthetic
459 data) relative to the posterior distributions, and the wide spread of the posterior distributions
460 (indicated by the scale bar in Fig. 4). We have included the ABC Markov chain Monte Carlo
461 traces corresponding to Fig. 4 in the supplementary material (Section S5).

462

463 A possible reason why the synthetic data parameters cannot be identified is that the synthetic
464 data does not accurately represent the parameter values used to generate it, making parameter
465 identification infeasible. To examine this possibility we calculated the variance in the PCF
466 synthetic data. In Fig. 5 (a)-(c) the blue line indicates the variance in the PCF synthetic data
467 for the current simulation design generated from five replicates of the ABM on a domain of
468 dimension $L_x = 23$ by $L_y = 23$.

469

470 If the variance in the summary statistics of the synthetic data precludes accurate identifica-
471 tion of model parameters using ABC, a sensible strategy may be to examine methods to reduce
472 the variance in the summary statistics of the synthetic data. Reducing the variance of the
473 summary statistics may mean the synthetic data is a more accurate reflection of the parameters

⁷Using the relationship that the diffusion coefficient is equal to $P_m \Delta^2$.

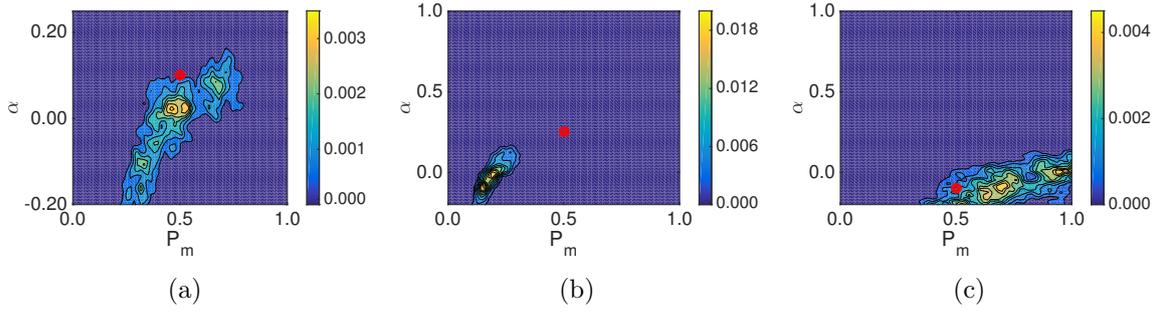


Figure 4: Posterior distributions for simulations of the realistic experiment described in Section 2.5 using the PCF as a summary statistic for an ABM of dimension $L_x = 23$ and $L_y = 23$. The synthetic data is generated from five replicates of the ABM. (a) Model A: $P_m = 0.5$, $\alpha = 0.1$, (b) model B: $P_m = 0.5$, $\alpha = 0.25$, (c) model B: $P_m = 0.5$, $\alpha = -0.1$. In all cases the red dot indicates the value of the parameters used to generate synthetic data.

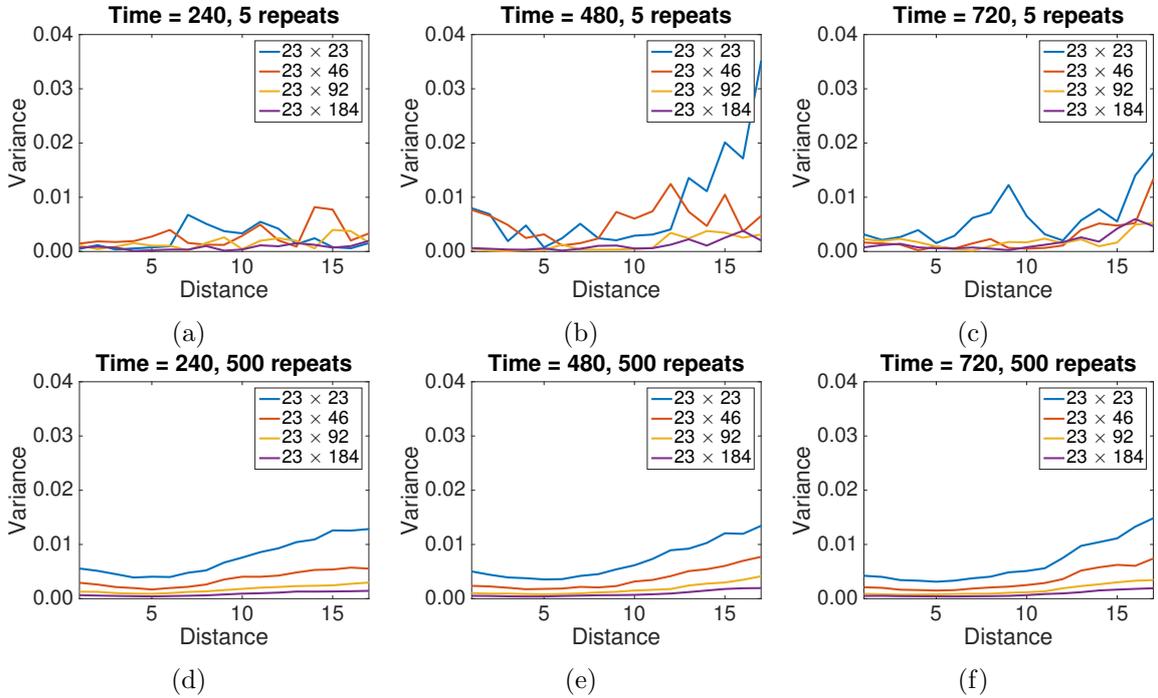


Figure 5: The variance in the PCF synthetic data for model B with $P_m = 0.5$, $\alpha = 0.25$ and different ABM domain sizes. Panels (a)-(c) display synthetic data generated from five replicates of the ABM, panels (d)-(f) display synthetic data generated from 500 replicates of the ABM. The domain size is indicated in the legend.

474 values used to generate it. This may also explain why parameter identification for the unreal-
 475 istic experiment was successful, as the variance in the summary statistics of the synthetic data
 476 was much smaller than for the practically realisable experiment (data not shown).

477

478 We conjectured that the variance in the summary statistics of the synthetic data could be

479 reduced in two ways:

- 480 1. increasing the number of ABM replicates used to generate the synthetic data;
- 481 2. increasing the size of the ABM domain while keeping the column density of the initial
482 conditions invariant. An example of this proposed initial condition is given in Fig. 6 (b),
483 in which the domain is twice the size of that in Fig. 6 (a). Importantly, increasing the
484 size of the ABM domain increases the number of agents in the simulation, and can be
485 thought of as equivalent to increasing the field of view of the microscope.

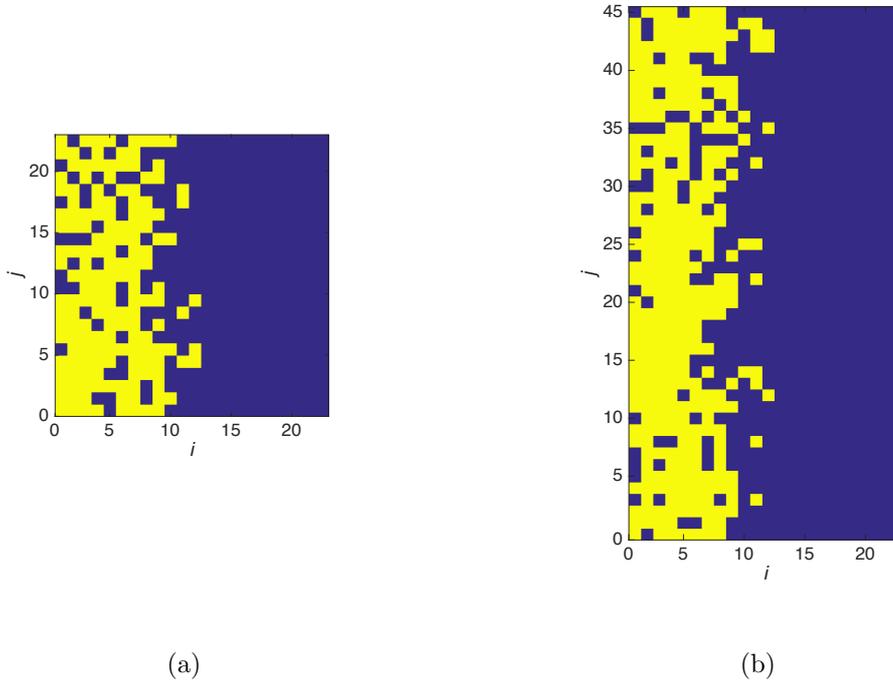


Figure 6: Increasing the size of the simulation domain while keeping the initial column densities the same. The domain in (b) is twice the size of that in (a), however, the average initial density of each column is the same in both (a) and (b).

486 In Fig. 5 the variance in the PCF synthetic data for model B with $P_m = 0.5$ and $\alpha = 0.25$
487 for different domain sizes and varying numbers of replicates can be seen. It is evident that the
488 variance in the PCF calculated from 500 replicates of our ABM on a $L_x = 23$ by $L_y = 23$ sized
489 domain (blue line in Fig. 5 (d)-(f)) is greater than the variance in the PCF calculated from five
490 replicates of our ABM on a $L_x = 23$ by $L_y = 184$ sized domain (purple line in Fig. 5 (a)-(c)).
491 This can be understood by considering Eq. (7): the number of occupied lattice pairs for each
492 horizontal pair distance used to generate the PCF does not increase linearly with the number

493 of agents. Specifically, the number of occupied lattice pairs for each horizontal pair distance
494 that generates the PCF is proportional to⁸

$$495 \frac{N(N-1)}{2}. \quad (13)$$

496

497 Therefore, the identification of parameters in experimental data using the PCF as a summary
498 statistic may be best facilitated by increasing the size of the domain upon which the experiment
499 is performed, rather than increasing the number of replicates of an experiment with a smaller
500 domain. Further variance plots for models A and B for the PCF summary statistic can be found
501 in the supplementary material (Section S6).

502

503 It is important to note that it is also the case for the *agent density profile* synthetic data,
504 that increasing the size of the domain is more effective at reducing variance in the synthetic
505 data than increasing the number of replicates. If generated from 500 replicates of our ABM on
506 an $L_x = 23$ by $L_y = 23$ sized domain, the agent density profile synthetic data will have greater
507 variance than the agent density profile synthetic data generated from five replicates of our ABM
508 on an $L_x = 23$ by $L_y = 184$ sized domain (data not shown). In this case the reduction in vari-
509 ance is an artefact of the lattice-based model. This is because the density of each column in the
510 ABM can take on a greater range of values between 0 and 1 as the column length is increased,
511 leading to a reduction in variance in the agent density profile synthetic data (especially in the
512 initial conditions of the simulations used to generate the synthetic data). However, as we do not
513 use the agent density profile summary statistic to identify parameters in the current simulation
514 design we do not pursue this matter further.

515 3.3 Improving the experimental design

516 We now confirm that more accurate identification of synthetic data parameters can be carried
517 out by expanding the domain upon which the experiment is performed, as opposed to increasing
518 the number of experimental replicates.

519

⁸This is not quite correct as a distance of ‘0’ between agents, that is they share the same column, is not accounted for in Eq. (7). To make Eq. (13) exact is not trivial as the expected number of agents each agent shares a column with depends on both the column position and simulation time.

520 In Fig. 7 (a)-(c) we plot the posterior distribution for synthetic data generated from 500
521 replicates of our ABM on a $L_x = 23$ by $L_y = 23$ sized domain, while in Fig. 7 (d)-(f) we plot
522 the posterior distribution generated from synthetic data generated from five replicates of our
523 ABM on a $L_x = 23$ by $L_y = 184$ sized domain⁹. As predicted, it is apparent that increasing the
524 domain size is more effective for parameter identification than increasing the number of repli-
525 cates used to generate the synthetic data. This is evident in the location (and narrow spread) of
526 the posterior distribution relative to the red dot, whereby the peak of the posterior distribution
527 is closer to the red dot in the case of Fig. 7 (d)-(f) compared to Fig. 7 (a)-(c). Despite this,
528 the identification of the parameters for repulsive interactions remains somewhat elusive (Fig. 7
529 (f)). A possible reason for this is that the repulsive interaction we present here is a weak one,
530 due to the constraint of Eqs. (2) and (4), and larger values of $|\alpha|$ are easier to identify as they
531 have a more profound effect on the behaviour of the agent population.

532

533 Computing $D_{KL}(p|\pi)$ for all six plots in Fig. 7 gives: (a) 2.55; (b) 2.69; (c) 1.53; and (d)
534 3.69; (e) 2.97; (f) 3.54. In tandem with the proximity of the peak of the posterior distribution
535 densities to the red dots in Fig. 7 (d)-(f) compared to Fig. 7 (a)-(c), this increase in the KLD
536 suggests that generating synthetic data on a larger domain is more effective for improving pa-
537 rameter identification than increasing the number of replicates used to generate the synthetic
538 data.

539 4 Discussion

540 In this work we have presented methods to identify motility and adhesion parameters in an
541 ABM of a wound-healing assay. Our findings suggest that for a commonly performed exper-
542 iment increasing the size of the experimental domain can be more effective in improving the
543 accuracy of parameter identification, when compared to increasing the number of replicates
544 of the experiment. This is because increasing the size of the domain, which is equivalent to
545 increasing the number of cells in the experiment, more effectively reduces the variance in the
546 summary statistics of the synthetic data from which the parameters are identified. The reason
547 for this reduction in variance is explained by Eq. (7), where the number of agent pair counts that

⁹A Markov chain Monte Carlo trace corresponding to Fig. 7 (e) can be found in the supplementary material (Section S5).

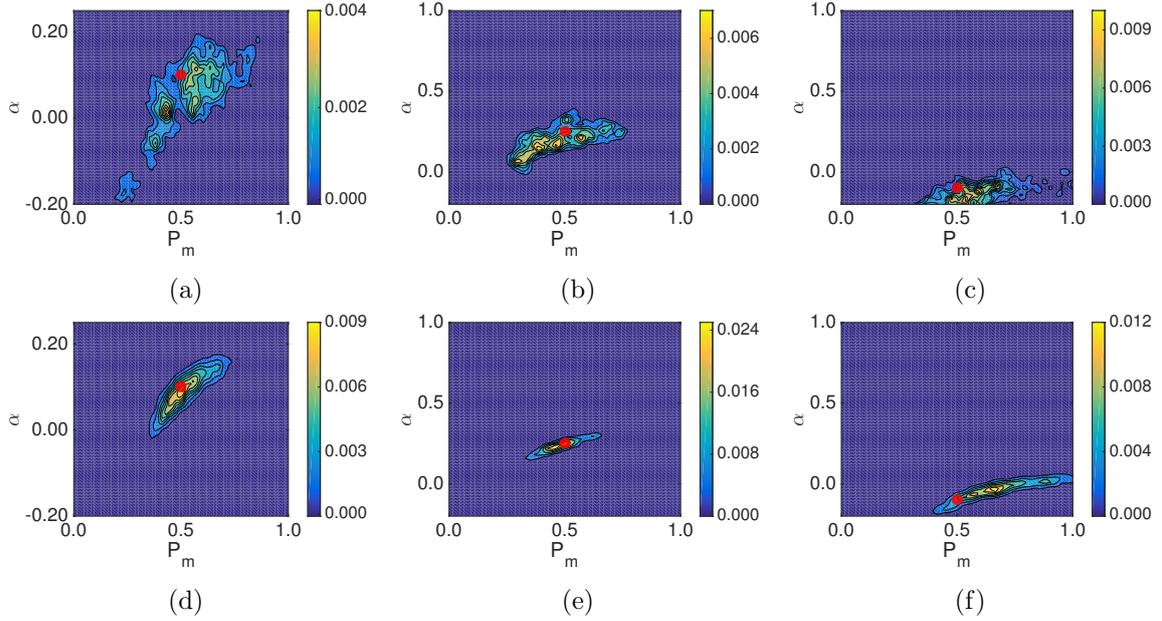


Figure 7: (a)-(c) Posterior distributions for simulations of the realistic experiment using the PCF as a summary statistic for an ABM simulated on a domain of dimension $L_x = 23$ by $L_y = 23$ with synthetic data generated from 500 replicates. (a) Model A: $P_m = 0.5$, $\alpha = 0.1$, (b) model B: $P_m = 0.5$, $\alpha = 0.25$, (c) model B: $P_m = 0.5$, $\alpha = -0.1$. (d)-(f) Posterior distribution plots for simulations of the experiment using the PCF as a summary statistic for an ABM simulated on a domain of size $L_x = 23$ by $L_y = 184$ with synthetic data generated from five replicates. (a) Model A: $P_m = 0.5$, $\alpha = 0.1$, (b) model B: $P_m = 0.5$, $\alpha = 0.25$, (c) model B: $P_m = 0.5$, $\alpha = -0.1$. Further figure information can be found in Fig. 4.

548 generate the PCF increases nonlinearly with the number of agents on the domain. In addition,
549 increasing the size of the experimental domain may make the collection of experimental data
550 less time-consuming, as potentially fewer replicates of the experiment will have to be conducted.
551 For instance, five replicates of the experiment on a larger domain provides more information
552 about parameters than 500 replicates of the experiment on a smaller domain (in the examples
553 we have presented in this work). Therefore, a comprehensive study of all summary statistics
554 commonly used for analysing cell migration, to understand how their variance scales with the
555 size of the experimental domain, is an interesting avenue for further research.

556

557 We also studied using the average horizontal displacement of agents and the agent density
558 profile as summary statistics. These were found to be less effective than the PCF in parameter
559 identification. This was especially the case for the averaged agent displacement, whereby a
560 range of adhesion and motility parameters could result in the same average agent displacement.
561 This result suggests that agent displacement may not be a suitable summary statistic for iden-

562 tifying cell motility and adhesion parameters, due to parameter identifiability issues.

563

564 The most obvious extension to the work presented here is to experimentally validate the find-
565 ings. That is, expand the wound-healing experimental domain and demonstrate: i) the cell
566 migratory process can be effectively described by the model we have presented here; and ii)
567 the experimental parameters are identifiable given a larger experimental domain. If validated,
568 evidence may be provided that demonstrates which adhesion model, A or B, is more applicable
569 to the cell type under consideration. Subsequently, we could add further agent behaviours to
570 the ABM, such as the role of the cell cycle. This may allow us to better capture the behaviour of
571 the cell populations we have studied here, and so produce more realistic models of cell migration.

572

573 To conclude, the findings presented in this work will be of particular interest to those con-
574 cerned with performing experiments that enable the effective parameterisation of cell migratory
575 processes. In particular, cell migratory processes in which cell-cell adhesion or repulsion are
576 known to play an important role. More generally, we have also suggested time and cost-saving
577 alterations to a commonly performed experiment for identifying cell motility parameters.

578 **Acknowledgements**

579 RJHR would like to thank the UK's Engineering and Physical Sciences Research Council (EP-
580 SRC, EP/G03706X/1) for funding through a studentship at the Systems Biology programme
581 of The University of Oxford's Doctoral Training Centre. RLM was supported by a Medical
582 Research Scotland Project Grant (436FRG). The authors declare no competing interests.

583 **Contributions**

584 RJHR, REB, AP and CAY conceived the work, and performed the mathematical and compu-
585 tational analysis. Data collection and analysis was performed by RLM and MJF. RJHR, REB
586 and CAY drafted the manuscript. All authors agree with manuscript results and conclusions.
587 All authors approved the final version.

588 References

- 589 [1] K. J. Cheung and A. J. Ewald. Illuminating breast cancer invasion: diverse roles for cell–cell
590 interactions. *Current Opinion in Cell Biology*, 30:99–111, 2014.
- 591 [2] A. Santiago and C. A. Erickson. Ephrin-B ligands play a dual role in the control of neural
592 crest cell migration. *Development*, 129(15):3621–3632, 2002.
- 593 [3] J. J. Fredberg. Power steering, power brakes, and jamming: Evolution of collective cell-cell
594 interactions. *Physiology*, 29(4):218–219, 2014.
- 595 [4] R. L. Mort, R. J. H. Ross, K. J. Hainey, O. Harrison, M. A. Keighren, G. Landini, R. E.
596 Baker, K. J. Painter, I. J. Jackson, and C. A. Yates. Reconciling diverse mammalian
597 pigmentation patterns with a fundamental mathematical model. *Nature Communications*,
598 7(10288), 2016.
- 599 [5] B. J. Binder, K. A. Landman, D. F. Newgreen, J. E. Simkin, Y. Takahashi, and D. Zhang.
600 Spatial analysis of multi-species exclusion processes: application to neural crest cell migra-
601 tion in the embryonic gut. *Bulletin of Mathematical Biology*, 74(2):474–90, 2012.
- 602 [6] R. McLennan, L. Dyson, K. W. Prather, J. A. Morrison, R. E. Baker, P. K. Maini, and
603 P. M. Kulesa. Multiscale mechanisms of cell migration during development: theory and
604 experiment. *Development*, 139(16):2935–2944, 2012.
- 605 [7] R. McLennan, L. J. Schumacher, J. A. Morrison, J. M. Teddy, D. A. Ridenour, A. C. Box,
606 C. L. Semerad, H. Li, W. McDowell, D. Kay, P. K. Maini, R. E. Baker, and P. M. Kulesa.
607 Neural crest migration is driven by a few trailblazer cells with a unique molecular signature
608 narrowly confined to the invasive front. *Development*, 142(11):2014–2025, 2015.
- 609 [8] R. McLennan, L. J. Schumacher, J. A. Morrison, J. M. Teddy, D. A. Ridenour, A. C.
610 Box, C. L. Semerad, H. Li, W. McDowell, D. Kay, P. K. Maini, R. E. Baker, and P. M.
611 Kulesa. VEGF signals induce trailblazer cell identity that drives neural crest migration.
612 *Developmental Biology*, 407(1):12–25, 2015.
- 613 [9] S. T. Johnston, M. J. Simpson, D. L. S. McElwain, B. J. Binder, and J. V. Ross. Interpreting

- 614 scratch assays using pair density dynamics and approximate Bayesian computation. *Open*
615 *Biology*, 4(9):140097, 2014.
- 616 [10] S. T. Johnston, M. J. Simpson, and D. L. S. McElwain. How much information can be
617 obtained from tracking the position of the leading edge in a scratch assay? *Journal of The*
618 *Royal Society Interface*, 11(97):20140325, 2014.
- 619 [11] S. T. Johnston, J. V. Ross, B. J. Binder, D. L. . McElwain, P. Haridas, and M. J. Simpson.
620 Quantifying the effect of experimental design choices for *in vitro* scratch assays. *Journal*
621 *of Theoretical Biology*, 400:19–31, 2016.
- 622 [12] T. Callaghan, E. Khain, L. M. Sander, and R. M. Ziff. A stochastic model for wound
623 healing. *Journal of Statistical Physics*, 122(5), 2006.
- 624 [13] E. Khain, L. M. Sander, and C. M. Schneider-Mizell. The role of cell-cell adhesion in wound
625 healing. *Journal of Statistical Physics*, 128(1-2):209–218, 2007.
- 626 [14] A. Q. Cai, K. A. Landman, and B. D. Hughes. Multi-scale modeling of a wound-healing
627 cell migration assay. *Journal of Theoretical Biology*, 245(3):576–594, 2007.
- 628 [15] M. Holcombe, S. Adra, M. Bicak, S. Chin, S. Coakley, A. I. Graham, J. Green, C. Gree-
629 nough, D. Jackson, M. Kiran, S. MacNeil, A. Maleki-Dizaji, P. McMinn, M. Pogson,
630 R. Poole, E. Qwarnstrom, F. Ratnieks, M. D. Rolfe, R. Smallwood, T. Sun, and D. Worth.
631 Modelling complex biological systems using an agent-based approach. *Integrative Biology*,
632 4(1):53–64, 2012.
- 633 [16] J. A. Flegg, S. N. Menon, P. K. Maini, and D. L. S. McElwain. On the mathematical
634 modeling of wound healing angiogenesis in skin as a reaction-transport process. *Frontiers*
635 *in Physiology*, 6(262), 2015.
- 636 [17] D. K. Schlüter, I. Ramis-Conde, and M. A. J. Chaplain. Computational modeling of single-
637 cell migration: the leading role of extracellular matrix fibers. *Biophysical Journal*, 103(6):
638 1141–1151, 2012.
- 639 [18] J. Liepe, S. Filippi, M. Komorowski, and M. P. H. Stumpf. Maximizing the information
640 content of experiments in systems biology. *PLoS Computational Biology*, 9(1):e1002888,
641 2013.

- 642 [19] P. Marjoram, J. Molitor, V. Plagnol, and S. Tavaré. Markov chain Monte Carlo without
643 likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328, 2003.
- 644 [20] E. van der Vaart, M. A. Beaumont, A. S. A. Johnston, and R. M. Sibly. Calibration and
645 evaluation of individual-based models using Approximate Bayesian Computation. *Ecolog-
646 ical Modelling*, 312:182–190, 2015.
- 647 [21] M. A. Beaumont, W. Zhang, and D. J. Balding. Approximate Bayesian computation in
648 population genetics. *Genetics*, 162(4):2025–2035, 2002.
- 649 [22] B. N. Vo, C. C. Drovandi, A. N. Pettitt, and G. J. Pettet. Melanoma cell colony expansion
650 parameters revealed by approximate Bayesian computation. *PLoS Computational Biology*,
651 11(12):e1004635, 2015.
- 652 [23] B. N. Vo, C. C. Drovandi, A. N. Pettitt, and M. J. Simpson. Quantifying uncertainty in
653 parameter estimates for stochastic models of collective cell spreading using approximate
654 Bayesian computation. *Mathematical Biosciences*, 263:133–142, 2015.
- 655 [24] P. J. M. Jones, A. Sim, H. B. Taylor, L. Bugeon, M. J. Dallman, B. Pereira, M. P. H.
656 Stumpf, and J. Liepe. Inference of random walk models to describe leukocyte migration.
657 *Physical Biology*, 12(6):066001, 2015.
- 658 [25] H. Kitano. Biological robustness. *Nature Reviews Genetics*, 5(11):826–837, 2004.
- 659 [26] T. M. Liggett. *Stochastic Interacting Systems: Contact, Voter, and Exclusion Processes*.
660 Springer-Verlag, Berlin, 1999.
- 661 [27] D. T. Gillespie. Exact stochastic simulation of coupled chemical reactions. *Journal of
662 Physical Chemistry*, 81(25):2340–2361, 1977.
- 663 [28] M. Fonoberova, V. A. Fonoberov, and I. Mezić. Global sensitivity/uncertainty analysis for
664 agent-based models. *Reliability Engineering & System Safety*, 118:8–17, 2013.
- 665 [29] J. Lee, T. Filatova, A. Ligmann-Zielinska, B. Hassani-Mahmooui, F. Stonedahl,
666 I. Lorscheid, A. Voinov, J. G. Polhill, Z. Sun, and D. C. Parker. The complexities of
667 agent-based modeling output analysis. *Journal of Artificial Societies and Social Simula-
668 tion*, 18(4):4, 2015.

- 669 [30] S. Marino, I. B. Hogue, C. J. Ray, and D. E. Kirschner. A methodology for perform-
670 ing global uncertainty and sensitivity analysis in systems biology. *Journal of Theoretical*
671 *Biology*, 254(1):178–196, 2008.
- 672 [31] R. J. H. Ross, C. A. Yates, and R. E. Baker. Inference of cell–cell interactions from
673 population density characteristics and cell trajectories on static and growing domains.
674 *Mathematical Biosciences*, 264:108–118, 2015.
- 675 [32] M. J. Simpson, K. K. Treloar, B. J. Binder, P. Haridas, K. J. Manton, D. I. Leavesley,
676 D. L. S. McElwain, and R. E. Baker. Quantifying the roles of cell motility and cell prolifer-
677 ation in a circular barrier assay. *Journal of The Royal Society Interface*, 10(82):20130007,
678 2013.
- 679 [33] D. J. G. Agnew, J. E. F. Green, T. M. Brown, M. J. Simpson, and B. J. Binder. Distin-
680 guishing between mechanisms of cell aggregation using pair-correlation functions. *Journal*
681 *of Theoretical Biology*, 352:16–23, 2014.
- 682 [34] B. J. Binder and M. J. Simpson. Quantifying spatial structure in experimental observations
683 and agent-based simulations using pair-correlation functions. *Physical Review E*, 88(2):
684 022705, 2013.
- 685 [35] T. Toni, D. Welch, N. Strelkowa, A. Ipsen, and M. P. H. Stumpf. Approximate Bayesian
686 computation scheme for parameter inference and model selection in dynamical systems.
687 *Journal of the Royal Society Interface*, 6(31):187–202, 2009.
- 688 [36] K. P. Burnham and D. R. Anderson. *Model Selection and Multimodel Inference: A Practical*
689 *Information-Theoretic Approach*. Berlin, Germany: Springer, 2002.
- 690 [37] M. A. Nunes and D. J. Balding. On optimal selection of summary statistics for approximate
691 Bayesian computation. *Statistical Applications in Genetics and Molecular Biology*, 9(1),
692 2010.
- 693 [38] D. Prangle, P. Fearnhead, M. P. Cox, P. J. Biggs, and N. P. French. Semi-automatic
694 selection of summary statistics for ABC model choice. *Statistical Applications in Genetics*
695 *and Molecular Biology*, 13(1):67–82, 2014.

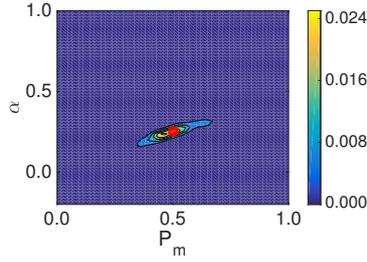
- 696 [39] C. P. Robert, J. Cornuet, J. Marin, and N. S. Pillai. Lack of confidence in approximate
697 Bayesian computation model choice. *Proceedings of the National Academy of Sciences*, 108
698 (37):15112–15117, 2011.
- 699 [40] C. P. Barnes, S. Filippi, M. P. H. Stumpf, and T. Thorne. Considerate approaches to
700 constructing summary statistics for ABC model selection. *Statistics and Computing*, 22
701 (6):1181–1197, 2012.

15 S1: Combining summary statistics

In Fig. S1 we plot the posterior distribution generated from combining all three summary statistics¹. As described in the main text, there is little difference between Fig. S1 in the supplementary material and Fig. 7 (e) of the main text. We quantify the difference between the posterior distributions in Fig. S1 and Fig. 7 (e) as follows:

$$\text{Difference} = \frac{1}{N} \sum_n^N |p^A(\Theta_n|D) - p^B(\Theta_n|D)|, \quad (\text{S1})$$

16 where the index n accounts for all possible discretised parameter pairs, $p^A(\Theta_n|D)$ is the posterior
 17 distribution in Fig. S1, and $p^B(\Theta_n|D)$ is the posterior distribution in Fig. 7 (e). The difference
 18 between the posterior distributions in Fig. S1 and Fig. 7 (e) is 0.00006, which shows that the
 19 performance of all three summary statistics is little different from the performance of the PCF
 20 summary statistic individually. By means of comparison the difference between the posterior
 21 distributions in Fig. 7 (b) and Fig. 7 (e) is 0.00031.



(a)

Figure S1: Posterior distribution plot for simulations of the experiment using all three summary statistics for an ABM simulated on a domain of dimension $L_x = 23$ by $L_y = 184$ with synthetic data generated from five replicates. Model B: $P_m = 0.5$, $\alpha = 0.25$.

¹To combine all three summary statistics we implement Eq. (11) (equivalently Eq. (S3)). If the condition stipulated in Section S4 fails for any individual summary statistic the parameter values are rejected.

22 **S2: Experimental methods**

23 The details of the experiment we aim to identify cell motility and adhesion parameters from is as
24 follows: Fucci2a 3T3 flp-In cells were maintained in dulbeccos modified eagle medium (DMEM)
25 containing 10% fetal calf serum, 1% Penicillin/Streptomycin and 100 μ g/ml Hygromycin B [1].
26 A silicon well (Ibidi) was attached to the surface of a 24-well glass-bottomed plate (Greiner
27 bio-one) by surface tension and allowed to attach overnight. Cells were plated within the in-
28 sert and allowed to attach phenol-red free DMEM (Biochrom) containing 10% fetal calf serum,
29 and 1% Penicillin/Streptomycin. Cells migrating from the leading edge of the cell mass were
30 then imaged with a 20x objective using a Nikon A1R inverted confocal microscope in a heated
31 chamber supplied with 5% CO2 in air. All image analysis tasks (required to generate the initial
32 conditions for the ABM of a practically realisable experiment) were performed using custom
33 written macros for the Fiji [2] distribution of ImageJ an open source image analysis package
34 based on NIH Image [3]. The cell nucleus of each cell was identified by merging of the green
35 and red channels containing the Fucci signal followed by segmentation. The centre of mass of
36 each object in the segmented image was then determined automatically.

37

38 In total we have data from five replicates of the experiment. Each data set contains cell
39 track data for every cell for sixty-four hours imaged at twenty minute intervals. Therefore,
40 we have the information required to apply our summary statistics to the experimental data.
41 More specifically, we have the position of all cells at each time interval so that the expected
42 horizontal displacement of cells, cell density profile, and PCF may be computed.

43 **S3: Practically realisable experiment ABM design**

44 **Initial conditions**

To map the position of cells in the experimental images where cell position is a continuous
variable, (x, y) , to a discrete lattice site, (i, j) , we use the following formulae

$$i = \left\lceil \frac{x}{\Delta} \right\rceil, \quad j = \left\lceil \frac{y}{\Delta} \right\rceil, \quad (\text{S2})$$

45 where $\lceil \cdot \rceil$ denotes the ceiling function and Δ is as defined in the main text. Given the experi-
 46 mental data and the lattice size no two cells were mapped to the same lattice site².

47

48 The application of Eq. (S2) to the initial frames of the five experiments allowed the average
 49 initial condition for the experimentally realistic ABM to be calculated. These initial condi-
 50 tions are expressed in terms of the average initial density of each column. These average initial
 51 column densities are:

Column	Initial density
1 st	0.8261
2 nd	0.7826
3 rd	0.8261
4 th	0.8261
5 th	0.8261
6 th	0.7391
7 th	0.6957
8 th	0.6087
9 th	0.5217
10 th	0.2609
11 th	0.2174
12 th	0.0870
13 th – 23 rd	0

52

53 To generate the initial conditions at the start of each ABM realisation each site in a column
 54 receives an agent uniformly at random at a probability equal to the average initial column
 55 density of the column the site is in. Therefore, the initial condition in the ABM is generated
 56 such that an ensemble average of the initial conditions of many realisations would equal the
 57 averaged initial conditions from the experiment. This initial condition is then used in the
 58 experimentally realistic ABM simulations. An example of this initial condition can be seen in
 59 the main text.

²If two cells did map to the same lattice site, one of these cells would be placed in the nearest unoccupied lattice site to the original lattice site. If there was more than one nearest unoccupied lattice site, one of these sites would be chosen uniformly at random for the cell to be mapped to.

60 **Boundary conditions**

61 Following the start of the simulation the density of the first column is checked after each agent
62 movement event out of the first column in the ABM. If the first column's density is below
63 0.6, agents are added uniformly at random to empty sites in the first column until the density
64 of the first column is greater than 0.6. This mechanism and density ensures that the agent
65 density profile in the ABM matches the experimental density profile for the entire course of the
66 experimental data throughout the simulation.

67 **S4: Markov chain Monte Carlo ABC algorithm**

68 We define a transition kernel w that proposes Θ' values as a bivariate uniform distribution. The
69 transition kernel ensures $P_m \in [0, 1]$ and $\alpha \in [-0.2, 0.25]$ for the model A, and $P_m \in [0, 1]$ and
70 $\alpha \in [-0.2, 1.0]$ for the model B. The parameter d^* is a constant selected so that approximately
71 one percent of the proposed parameter sets are accepted, the value of which is obtained through
72 trial and error.

73

74 To implement a Markov chain Monte Carlo method (Metropolis-Hastings algorithm) we proceed
75 as follows [4]:

76 **R1** If at Θ step to Θ' according to a transition kernel $w(\Theta \rightarrow \Theta')$.

R2 Simulate \tilde{D} from the model using Θ' and calculate the summary statistic $S(\tilde{D})$ at each
sampling point. That is, for each individual $t = [240, 480, 720]$ calculate d :

$$d = \sum_{r=1}^R |S(D)_{r,t} - S(\tilde{D})_{r,t}|, \quad (\text{S3})$$

77 If $d > d^*$ (at any t) reject Θ' and return to R1.

78 **R3** Calculate

$$79 \quad h = \min \left(1, \frac{\pi(\Theta')w(\Theta' \rightarrow \Theta)}{\pi(\Theta)w(\Theta \rightarrow \Theta')} \right).$$

80 **R4** Accept Θ' with probability h .

81 **R5** Return to **1** until 10^6 steps have been attempted.

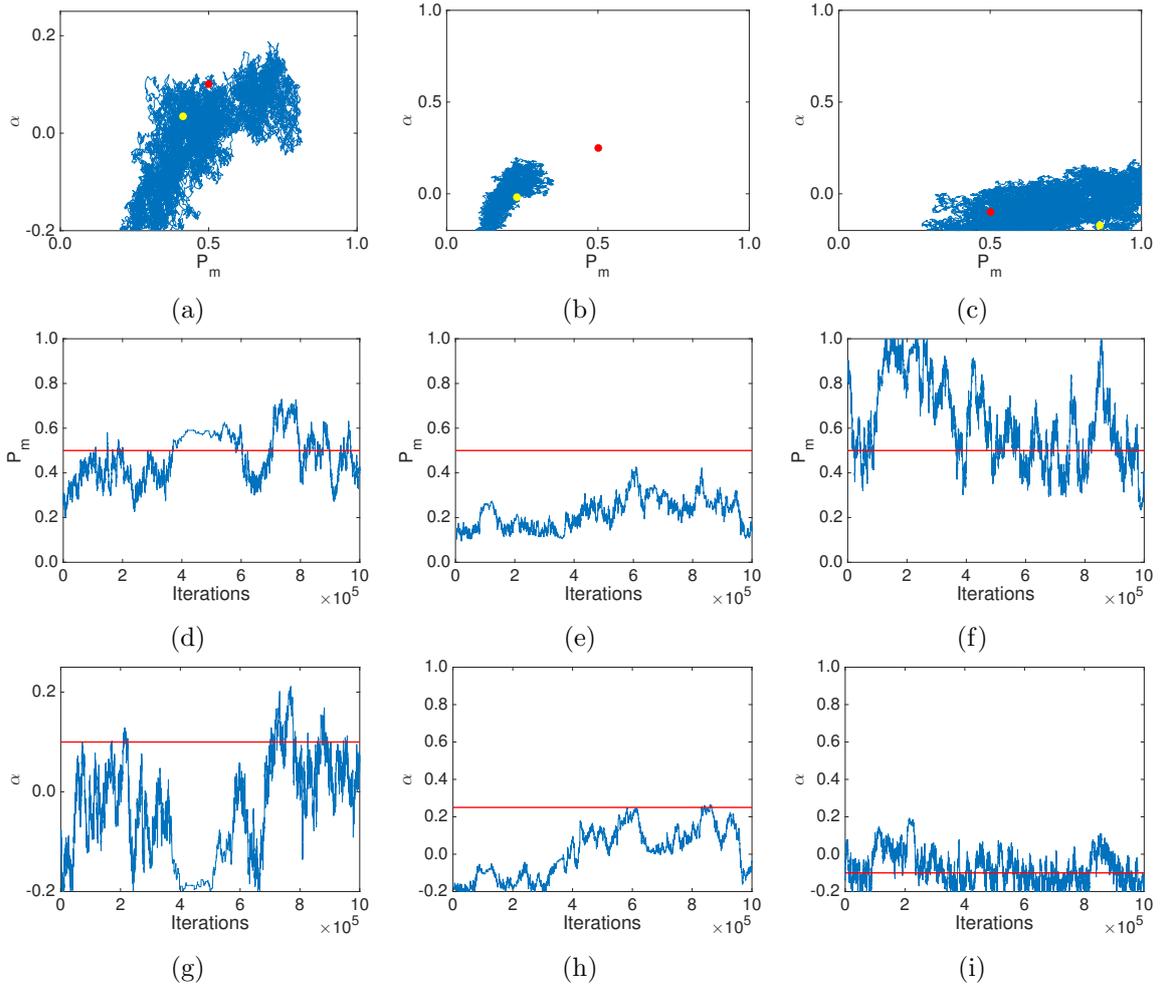
82 Initially, we sample Θ randomly from the prior distribution until a parameter set has been
83 accepted (**R4**).

84 **S5: Markov chains: trace plots**

85 In Fig. S2 (d)-(i) the Markov chain traces for the posterior distributions for Fig. 4 in the main
86 text are displayed. The mean and variance values for these chains are: (d) mean = 0.4722,
87 variance = 0.0115; (e) 0.2202, 0.0050; (f) 0.6236, 0.0356; (g) -0.0377, 0.0102; (h) 0.0087, 0.0176;
88 (i) -0.0734, 0.0074. The reason as to why the estimation of the values of P_m and α in the
89 synthetic data is inaccurate in Fig. S2 is because the synthetic data (in conjunction with the
90 PCF summary statistic) does not provide an accurate enough representation of the parameters
91 with which the synthetic data was generated i.e. the parameters are not identifiable. Therefore,
92 the Markov chain Monte Carlo ABC algorithm is not able to work effectively.

93

94 In the case of Fig. S3 (corresponds to Fig. 7 (e) in the main text) the same algorithm ac-
95 curately estimates the parameter values used to generate the synthetic data. This is because
96 the synthetic data in this case is an accurate representation of the parameters used to generate
97 it. The mean and variance values for these chains are: (b) 0.5627, 0.0086; (c) 0.2718, 0.0017.



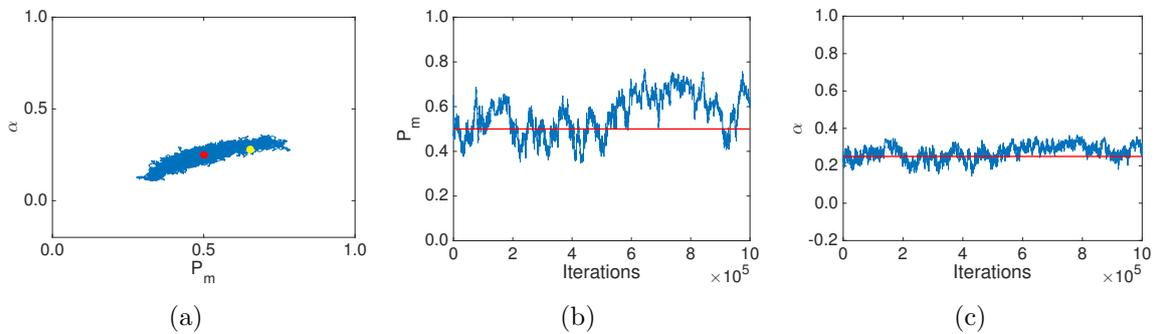


Figure S3: Markov chain Monte Carlo trace plots for Fig. 7 (e) in the main text. In panel (a) the yellow dot indicates the initial value of the chain used to generate the posterior distribution, and the red dot indicates the parameter values used to generate the synthetic data. In panels (b) and (c) the red line indicates the value of the parameter used to generate the synthetic data. Panels (b) and (c) display individual parameter trace plots. Panels (a), (b) and (c) correspond to model B, $P_m = 0.5$, $\alpha = 0.25$.

98 **S6: Further variance plots for models A and B for the PCF**
 99 **summary statistic.**

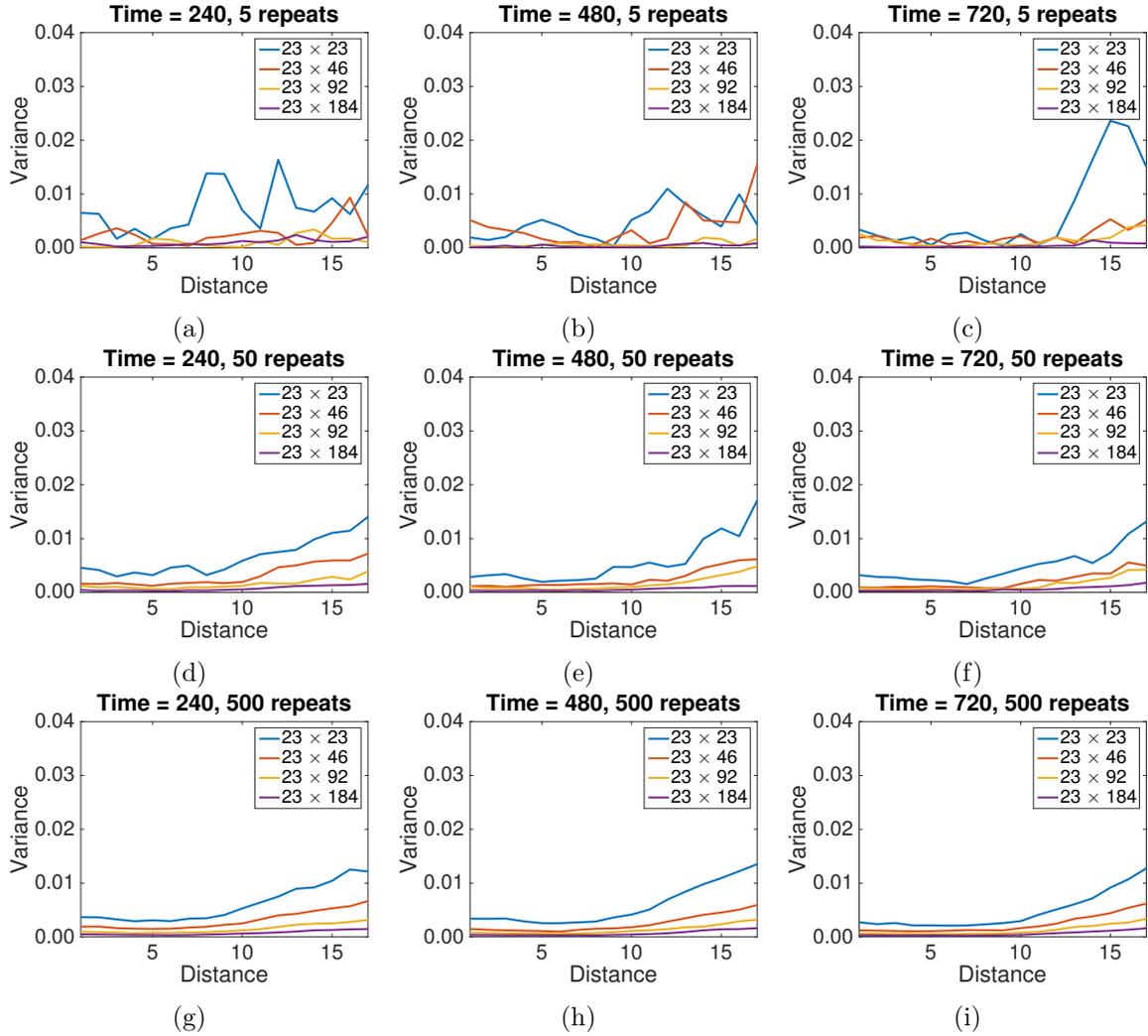


Figure S4: The variance in PCF synthetic data for model A with $P_m = 0.5$, $\alpha = 0.1$ for different ABM domain sizes. Panels (a)-(c) display synthetic data generated from five replicates of the ABM, panels (d)-(f) display synthetic data generated from 50 replicates of the ABM and panels (g)-(i) display synthetic data generated from 500 replicates of the ABM.

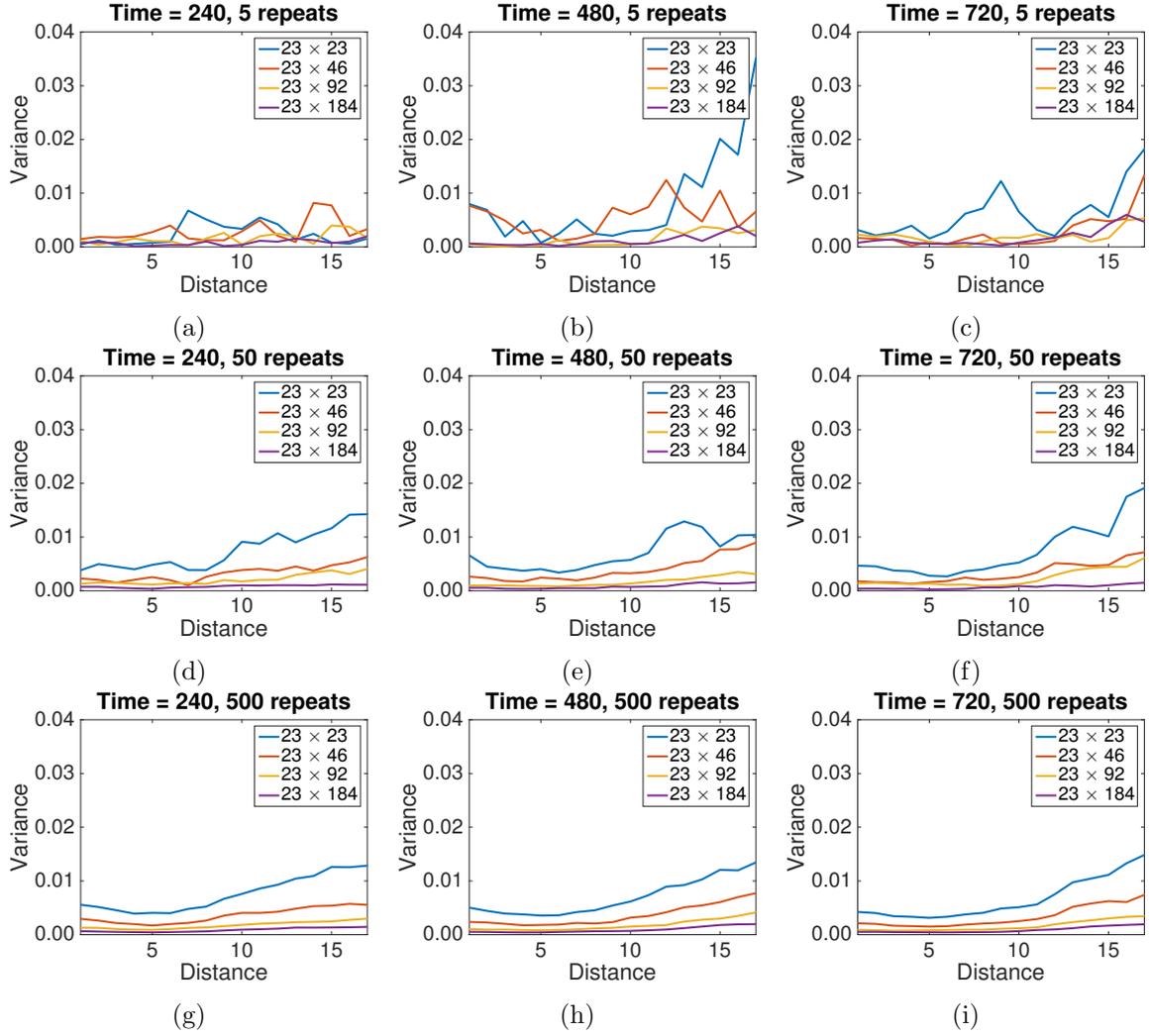


Figure S5: The variance in the synthetic data for model B with $P_m = 0.5$, $\alpha = 0.25$ for different ABM domain sizes. Panels (a)-(c) display synthetic data generated from five replicates of the ABM, panels (d)-(f) display synthetic data generated from 50 replicates of the ABM and panels (g)-(i) display synthetic data generated from 500 replicates of the ABM.

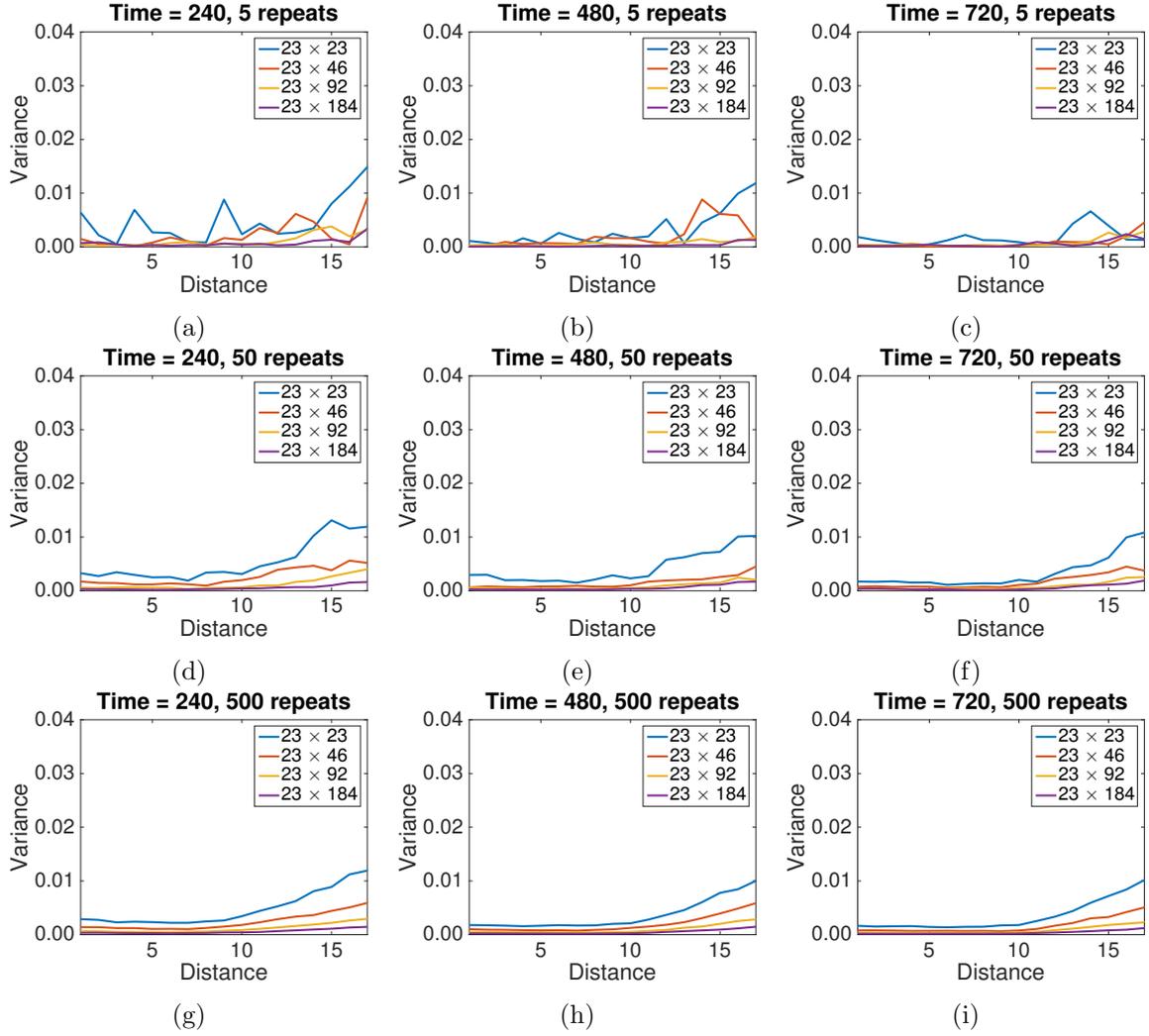


Figure S6: The variance in PCF synthetic data for model B with $P_m = 0.5$, $\alpha = -0.1$ for different ABM domain sizes. Panels (a)-(c) display synthetic data generated from five replicates of the ABM, panels (d)-(f) display synthetic data generated from 50 replicates of the ABM and panels (g)-(i) display synthetic data generated from 500 replicates of the ABM.

100 **References**

- 101 [1] R. L. Mort, M. J. Ford, A. Sakaue-Sawano, N. O. Lindstrom, A. Casadio, A. T. Douglas,
102 M. A. Keighren, P. Hohenstein, A. Miyawaki, and I. J. Jackson. Fucci2a: a bicistronic cell
103 cycle reporter that allows Cre mediated tissue specific expression in mice. *Cell Cycle*, 13
104 (17):2681–2696, 2014.
- 105 [2] J. Schindelin, I. Arganda-Carreras, E. Frise, V. Kaynig, M. Longair, T. Pietzsch,
106 S. Preibisch, C. Rueden, S. Saalfeld, B. Schmid, J. Y. Tinevez, D. J. White, V. Hartenstein,
107 K. Eliceiri, P. Tomancak, and A. Cardona. Fiji: an open-source platform for biological-image
108 analysis. *Nature Methods*, 9(7):676–82, 2012.
- 109 [3] C. A. Schneider, W. S. Rasband, and K. W. Eliceiri. NIH image to ImageJ: 25 years of
110 image analysis. *Nature Methods*, 9(7):671–5, 2012.
- 111 [4] P. Marjoram, J. Molitor, V. Plagnol, and S. Tavaré. Markov chain Monte Carlo without
112 likelihoods. *Proceedings of the National Academy of Sciences*, 100(26):15324–15328, 2003.