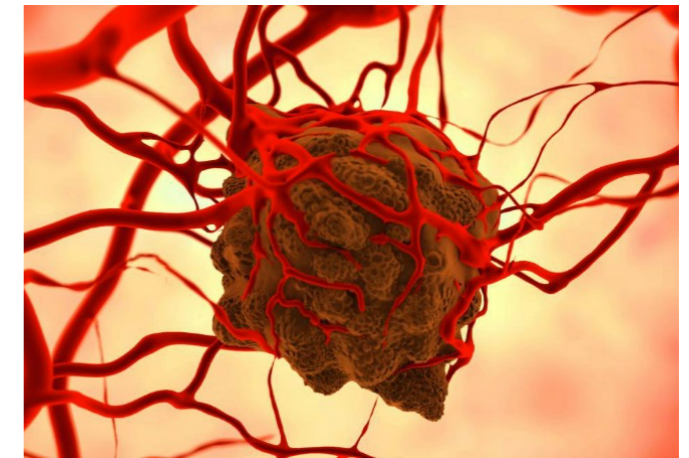# Leveraging concepts from stochastic simulation and machine learning for efficient Bayesian inference

Ruth Baker

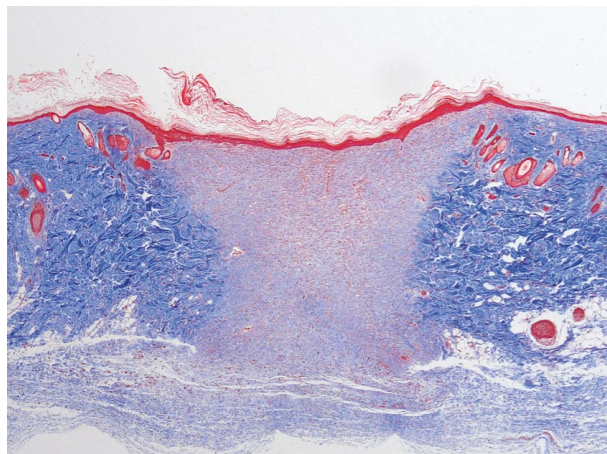@ruth_baker

# Applied research focus

To understand the mechanisms driving collective cell motility, proliferation and death and their contributions to complex biological processes, such as those associated with development, disease and repair.
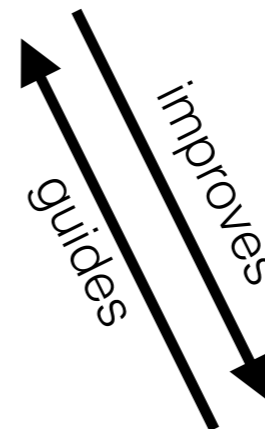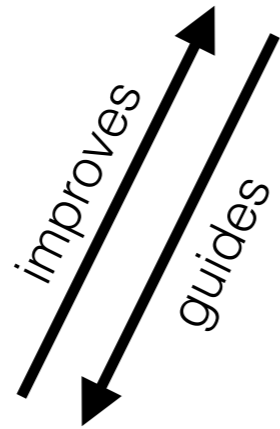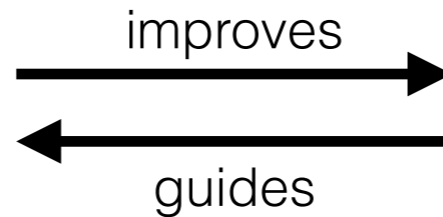


tumour growth



embryo development



wound healing

Combine mechanistic mathematical models and statistical / machine learning tools to provide new biological insights.

# The inverse problem

- Given quantitative data, can we estimate model parameters?



- Parameter inference using a Bayesian framework:

$$\mathbb{P}(\theta \,|\, \mathcal{D}) \propto \mathcal{L}(\mathcal{D} \,|\, \theta)\mathbb{P}(\theta)$$

posterior   likelihood   prior

# Approximate Bayesian computation

- Likelihood is generally intractable, so we use ABC:

  - A **very large number** of times:

    - sample parameter from some distribution;

> **When a model takes minutes (or longer) to simulate, and / or the model has many parameters, ABC methods can be infeasible.**

    - evaluate how close model output is to data using a summary statistic and distance function;

    - assign a weight to the parameter - depends on this distance.
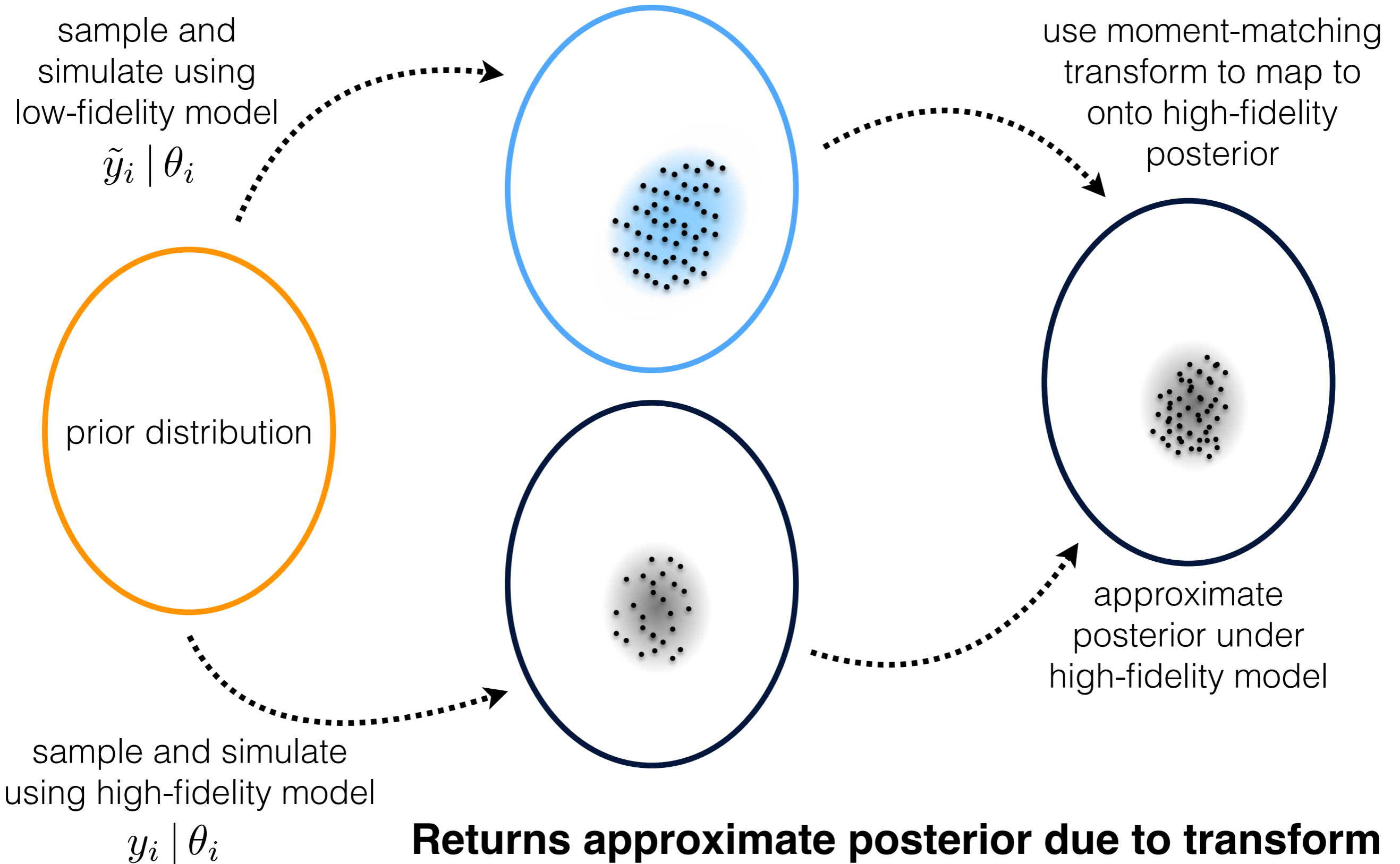
# Increasing the efficiency of ABC

- Use a more intelligent exploration of parameter space.

  - Importance sampling, sequential Monte Carlo *etc.*

  - **Pre-conditioned ABC** - Use a simple model to help transition from prior to posterior.

- Make the weights less expensive to calculate.

  - **Moment-matching ABC** - Use a simple model to help estimate posterior.

  - **Multifidelity ABC** - Use hierarchies of models in the weight calculations.

  - **Minibatch ABC** - Use subsets of the data in evaluating the distance function.

sample and simulate using low-fidelity model

$\tilde{y}_i \mid \theta_i$

sample, perturb and simulate using high-fidelity model

$y_i \mid \theta_i$

prior distribution

posterior under low-fidelity model

posterior under high-fidelity model

**Weight accepted particles so that ABC posterior is returned.**

Warne, Baker and Simpson, *J. Comp. Graphical Stat.* (2022).

# Moment-matching ABC
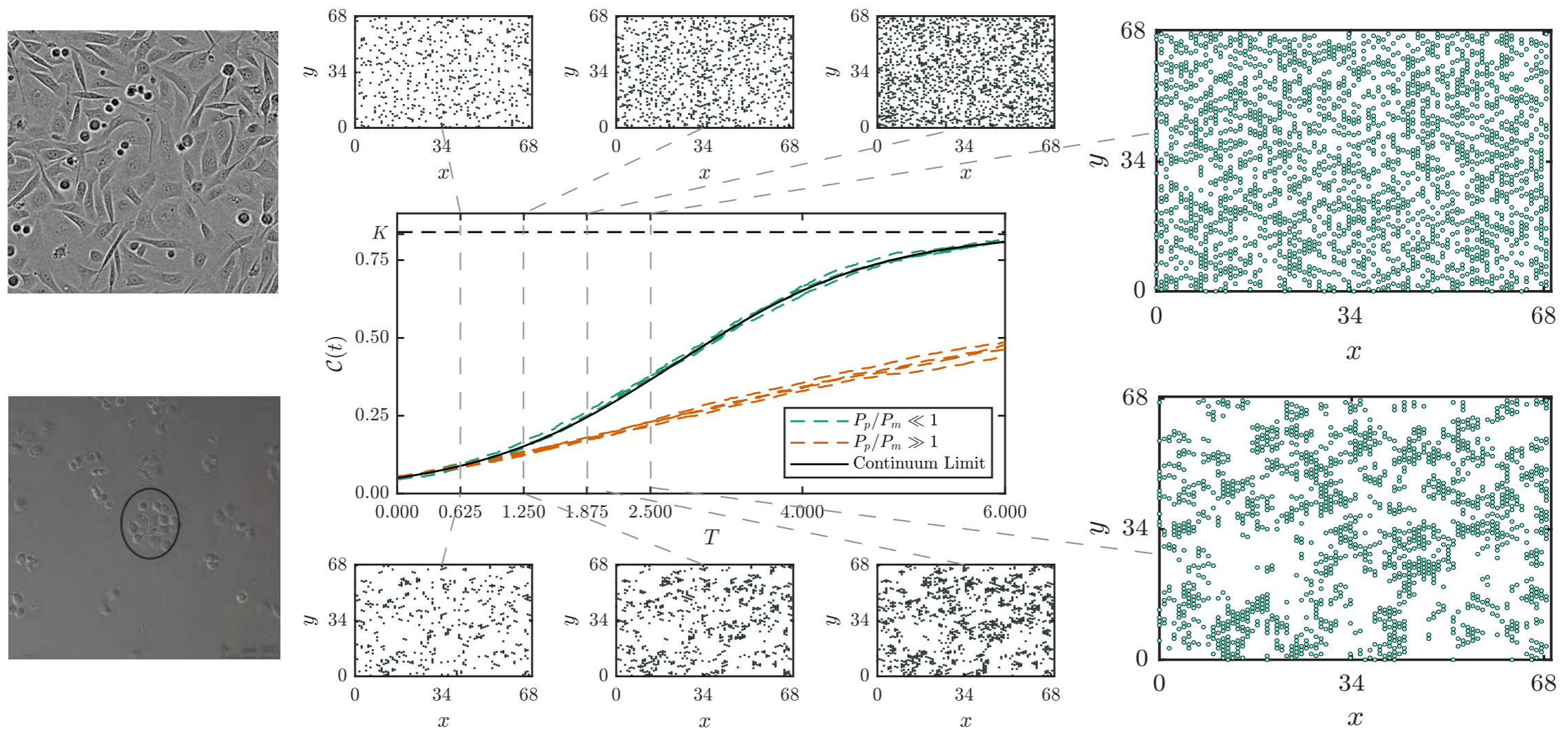


sample and
simulate using
low-fidelity model
$\tilde{y}_i \mid \theta_i$

use moment-matching
transform to map to
onto high-fidelity
posterior

prior distribution

sample and simulate
using high-fidelity model
$y_i \mid \theta_i$

approximate
posterior under
high-fidelity model

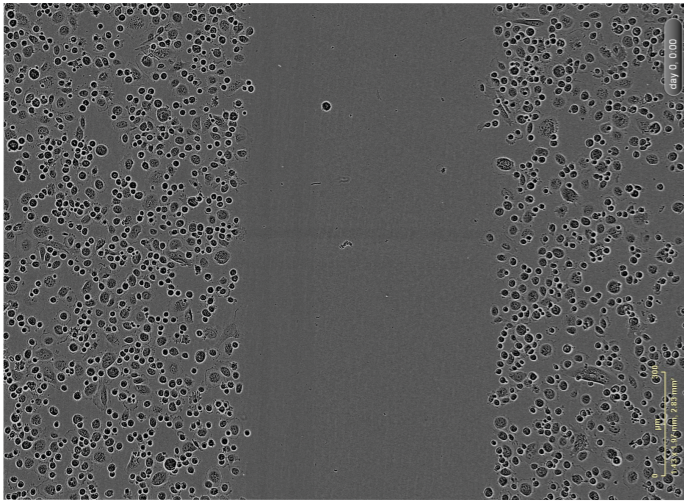**Returns approximate posterior due to transform**

# Pre-conditioned / MM ABC

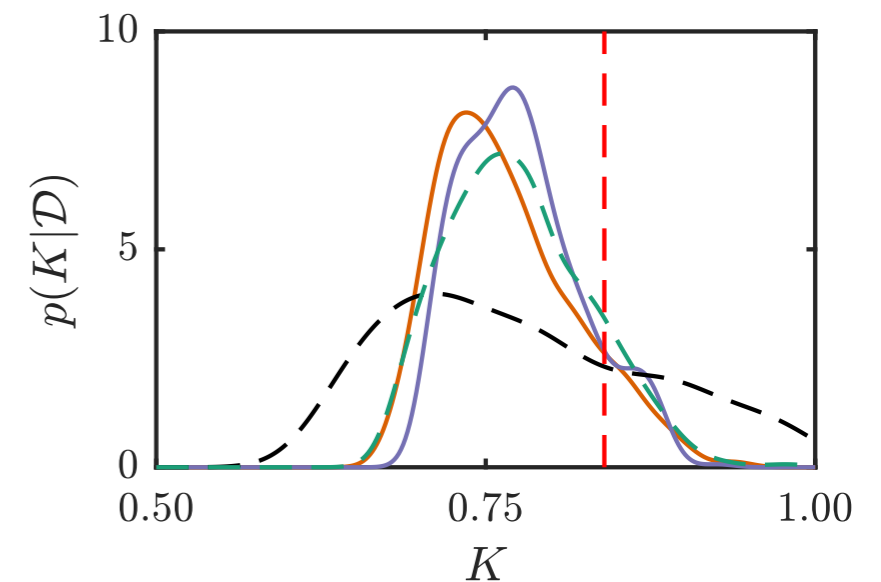- Discrete, on-lattice random walk model of cell migration and proliferation.
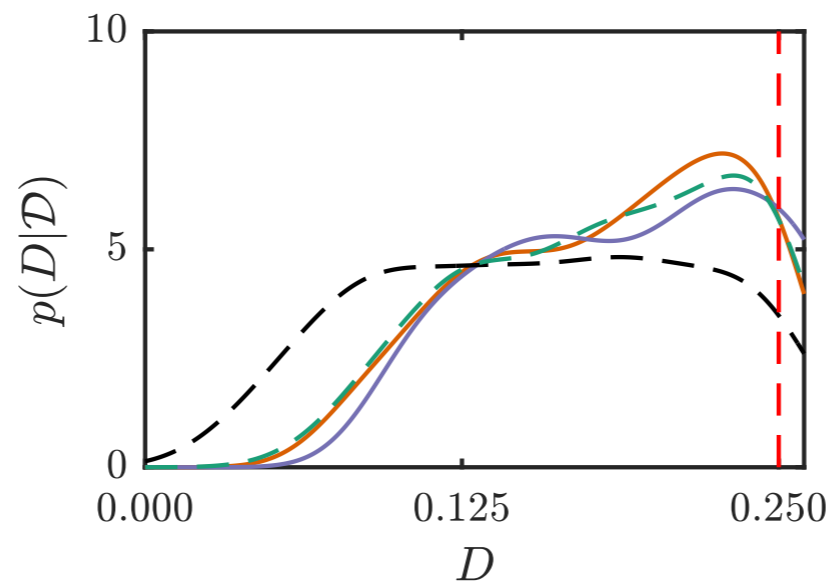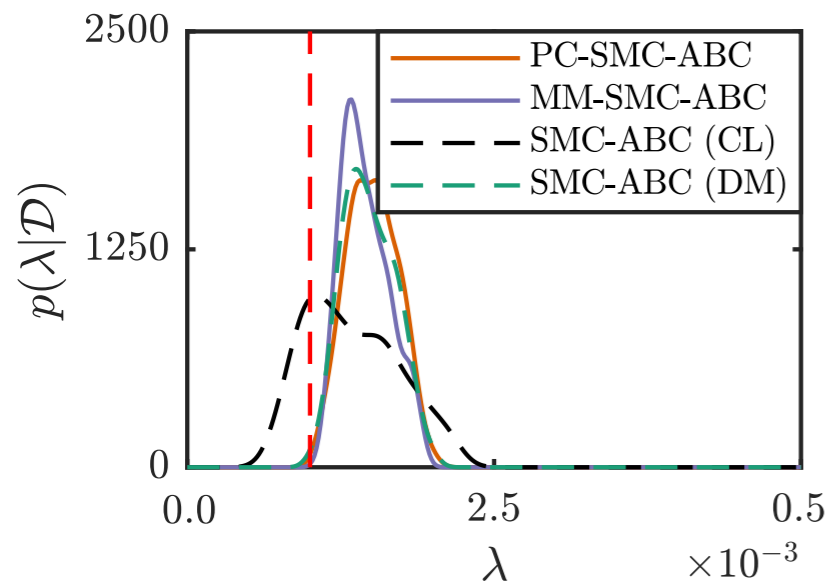


Warne, Baker and Simpson, *J. Comp. Graphical Stat.* (2022).

UNIVERSITY OF OXFORD

Continuum limit - Fisher-KPP equation:

$$\frac{\partial C}{\partial t} = D\frac{\partial^2 C}{\partial x^2} + \lambda C\left(1 - \frac{C}{K}\right)$$
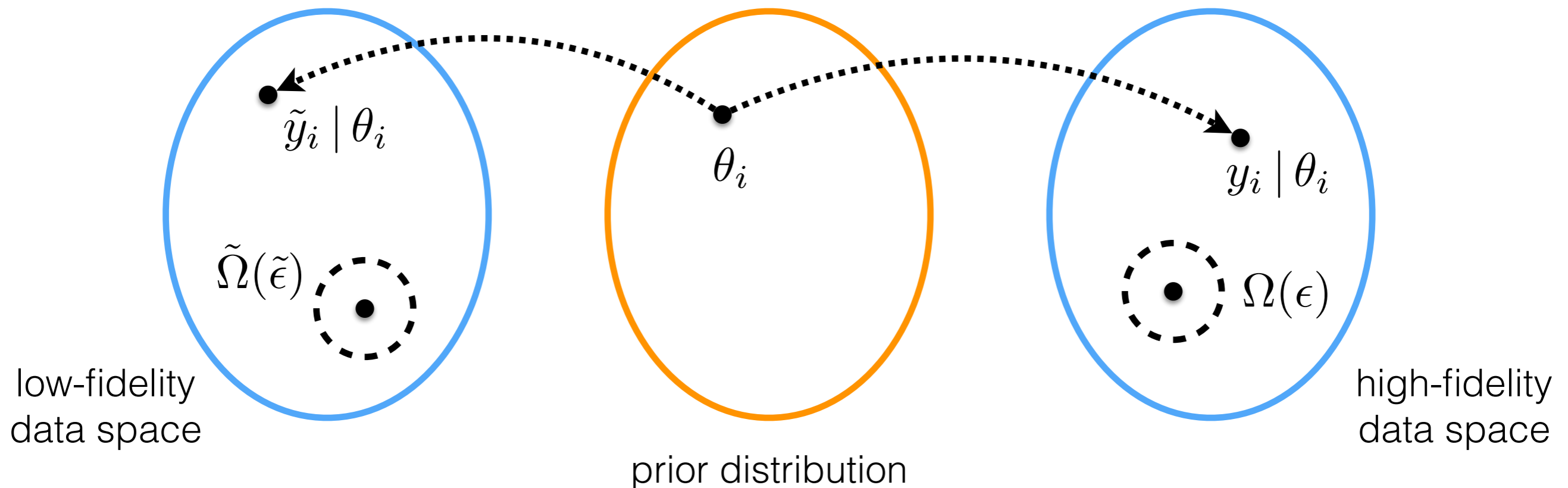
Marginal posterior parameter distributions



Warne, Baker and Simpson, *J. Comp. Graphical Stat.* (2022).

- Pre-conditioned and moment-matching ABC can provide significant time savings, through the combined use of high-fidelity and low-fidelity models.

- Need to explore trade-off between fraction of high-fidelity and low-fidelity samples.

- Doesn't require the low-fidelity model to be particularly accurate, just that the model outputs depend in a qualitatively similar way on the input parameters.
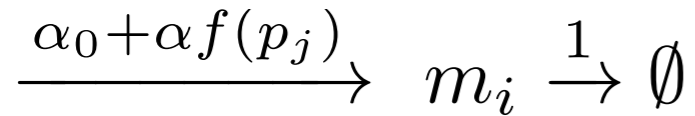
Warne, Baker and Simpson, *J. Comp. Graphical Stat.* (2022).

- Can we use the low-fidelity model to allocate some weights (accept or reject parameters) and not bias the result?
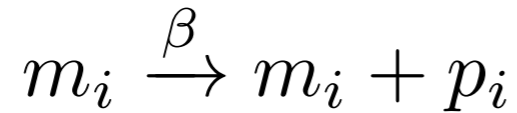


low-fidelity
data space

prior distribution

high-fidelity
data space

Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

$$\xrightarrow{\alpha_0 + \alpha f(p_j)} \quad m_i \xrightarrow{1} \emptyset \qquad \text{for } (i,j) = (1,3), (2,1), \text{ and } (3,2)$$

$$m_i \xrightarrow{\beta} m_i + p_i \qquad \text{for } i = 1, 2, 3$$

$$p_i \xrightarrow{\beta} \emptyset \qquad \text{for } i = 1, 2, 3 \qquad f(p) = \frac{K_h^n}{(K_h^n + p^n)}$$

Distance from data: multifidelity



- **Matching estimator values**
- **False positive**
- **False negative**

High-fidelity model - data distance

Low-fidelity model - data distance

- How to combine the outputs of the two models, so that the result is unbiased weights?

- How can we make this process efficient?

Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

- Attempt to make an "early decision" using the low-fidelity model, and "sometimes" check that decision using the high-fidelity model.

- Decision to check is made uniformly at random, with probability $\alpha(\tilde{y}, \theta_n)$

- Here, w            ⋯ming

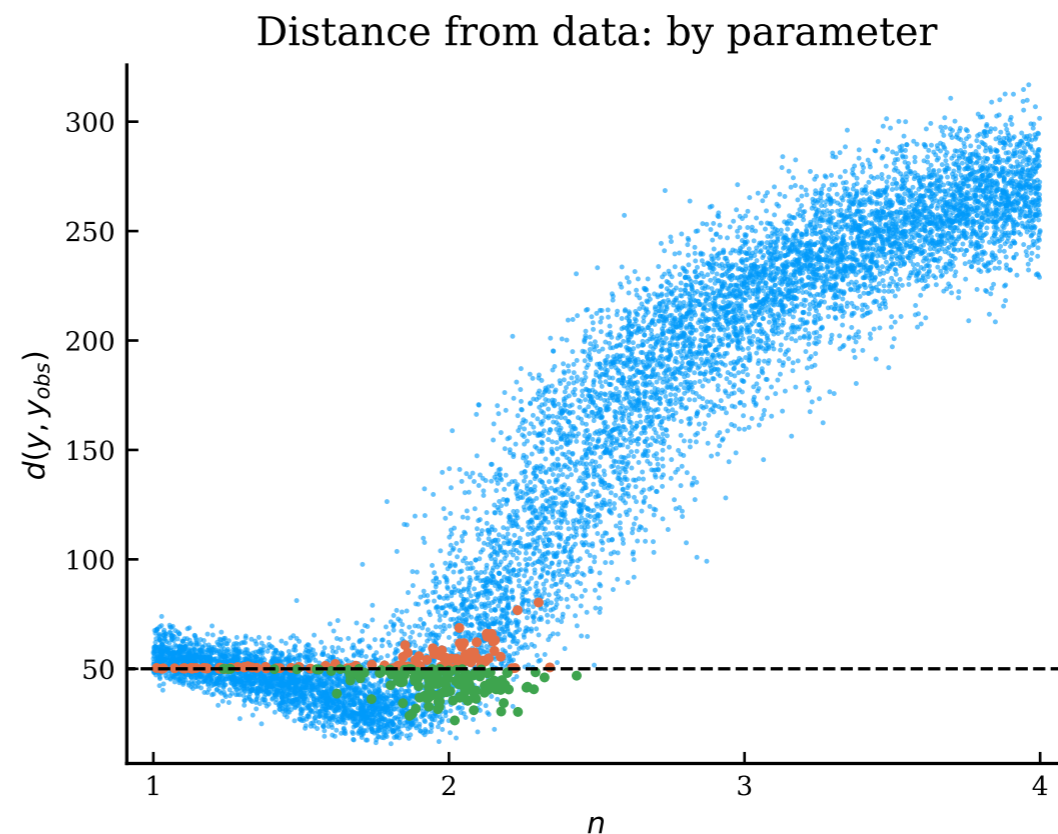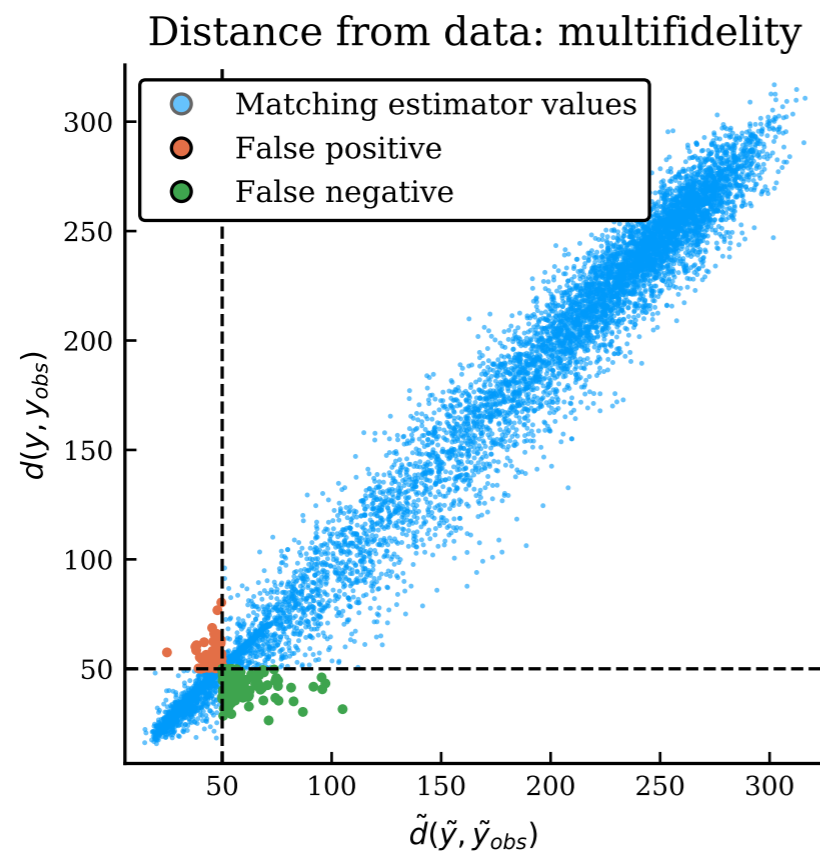**Using a common noise input makes this a very efficient approach.**

$$\alpha(y, \theta_n) = \eta_1 \mathbb{1}(y \in \Omega(\epsilon)) + \eta_2 \mathbb{1}(y \notin \Omega(\epsilon))$$

if low-fidelity model
is **close to** data

if low-fidelity model
is **far from** data

Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

- Derive analytical expressions for the optimal continuation probabilities, given estimates of various quantities.

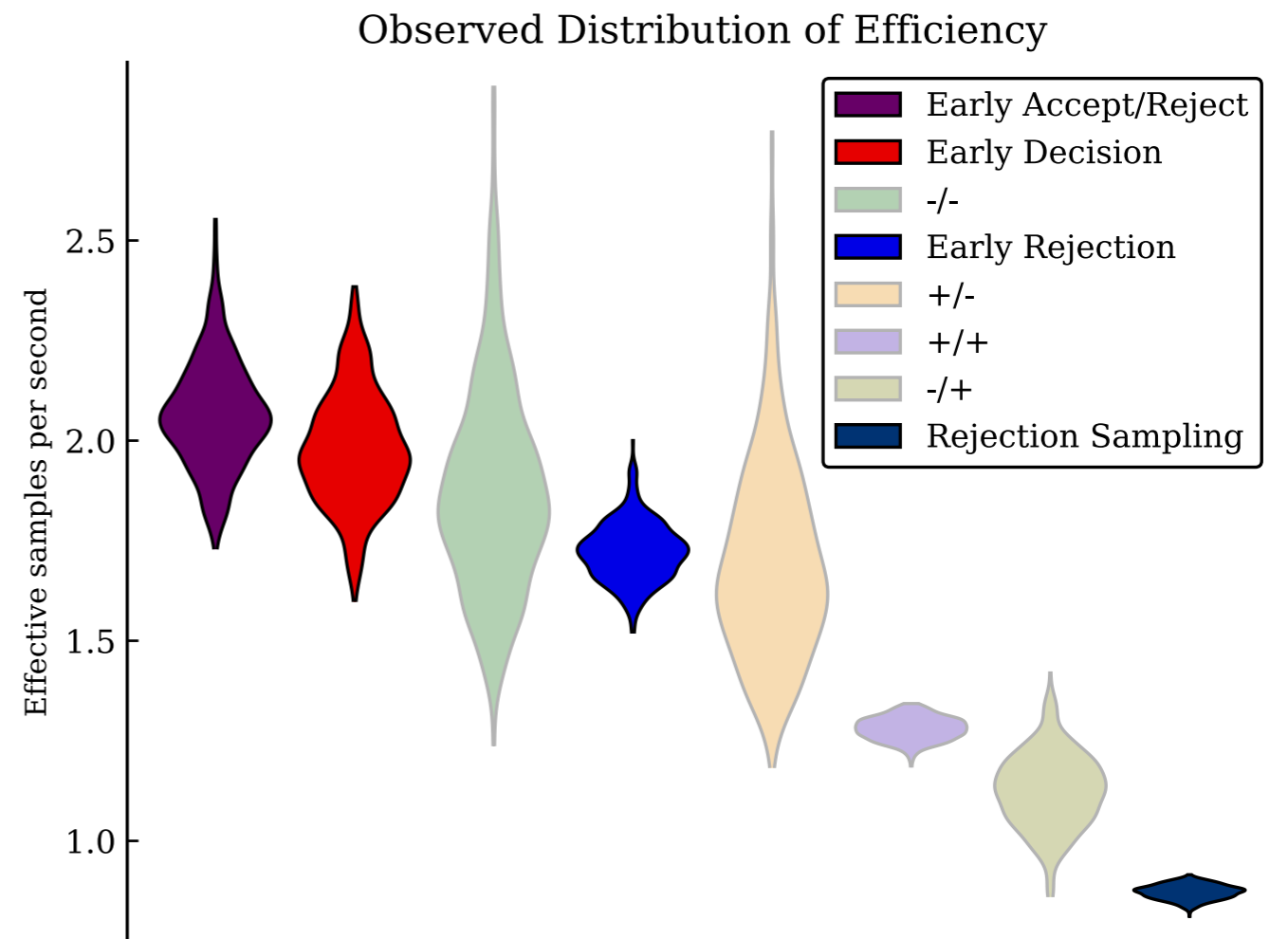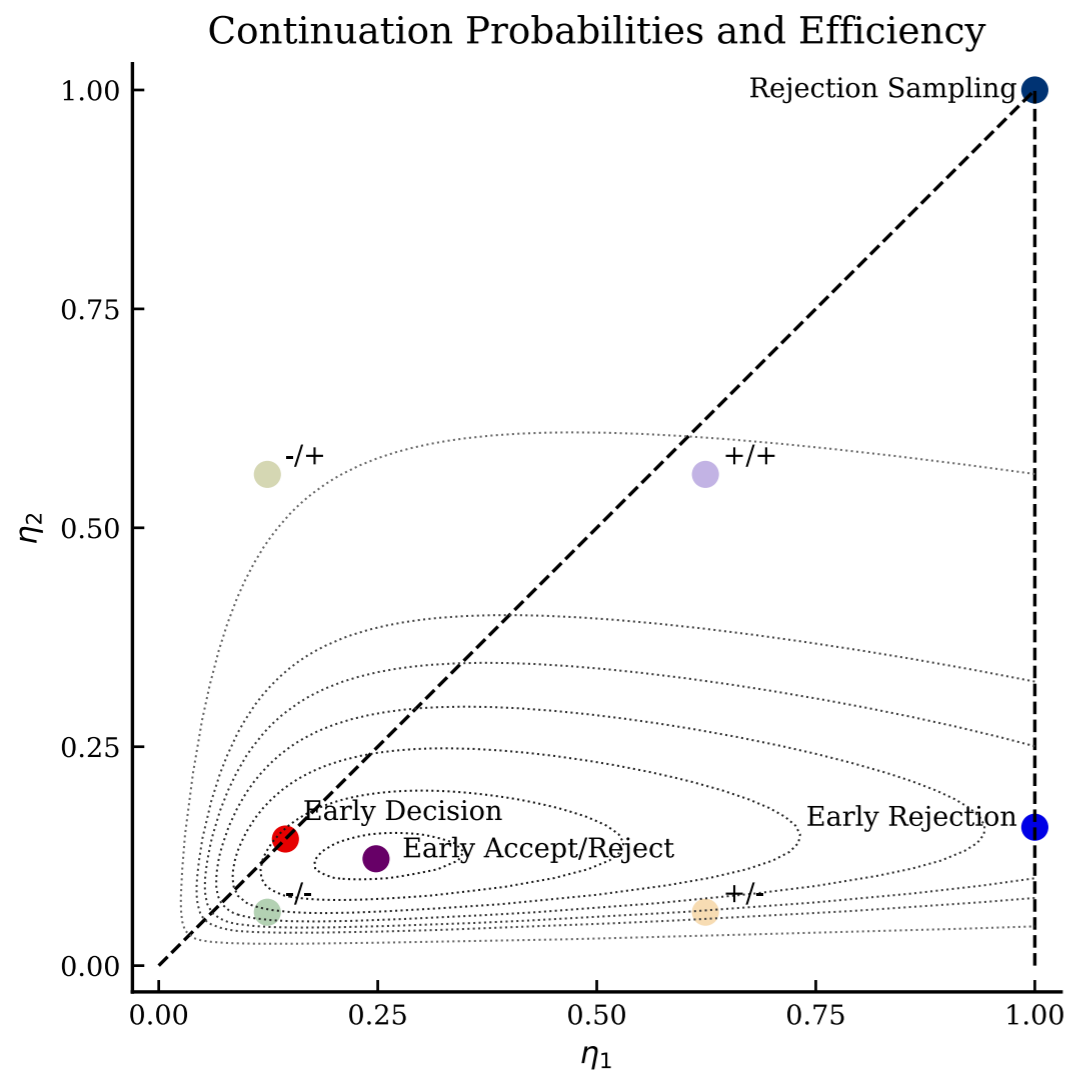  - In practice: adapt the continuation probabilities "on the fly", as samples are generated…



Distance from data: multifidelity

Distance from data: by parameter

Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

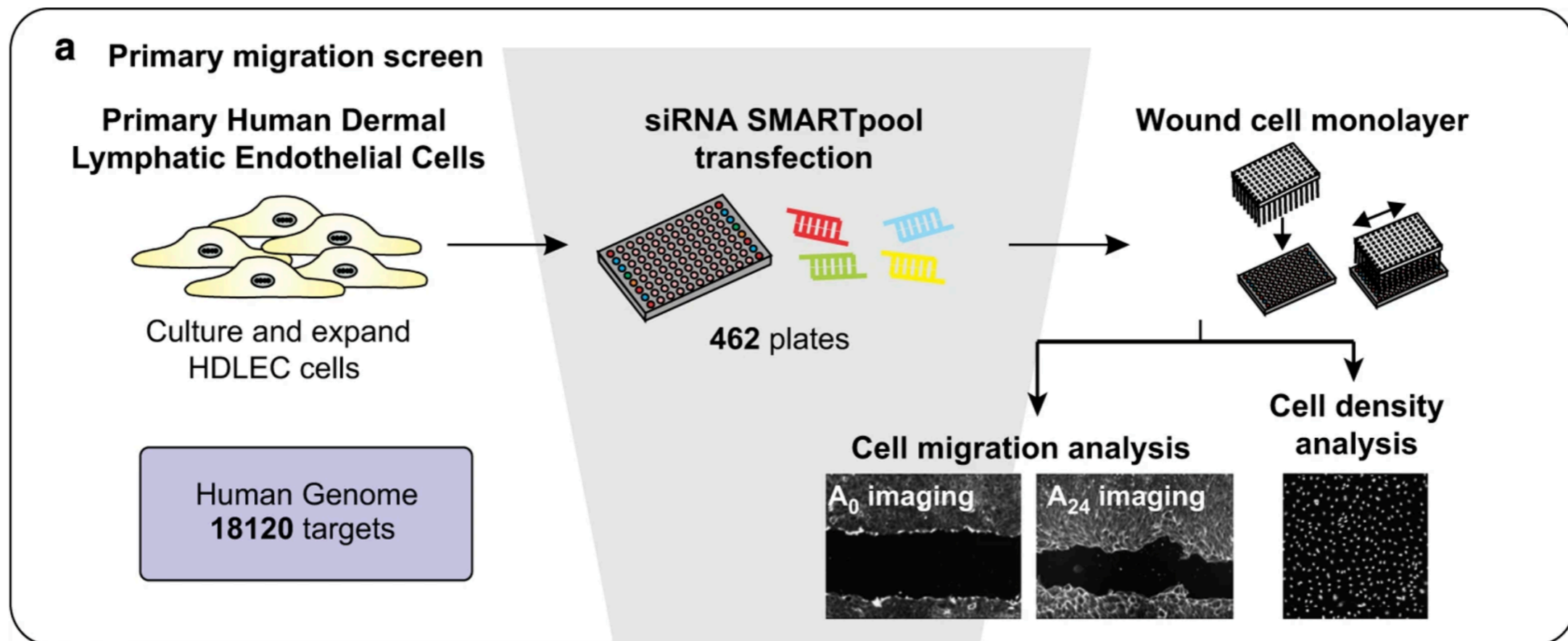- Comparing results for a range of continuation probabilities:



Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

- Multifidelity ABC can provide time savings, through the combined use of high- and low-fidelity models.

- Can "learn" optimal continuation probabilities as the algorithm proceeds, separately controlling rates of checking early acceptance and early rejection.

- Rates of false positives and negatives can be reduced by generating the high-fidelity model output conditional on the low-fidelity model output.

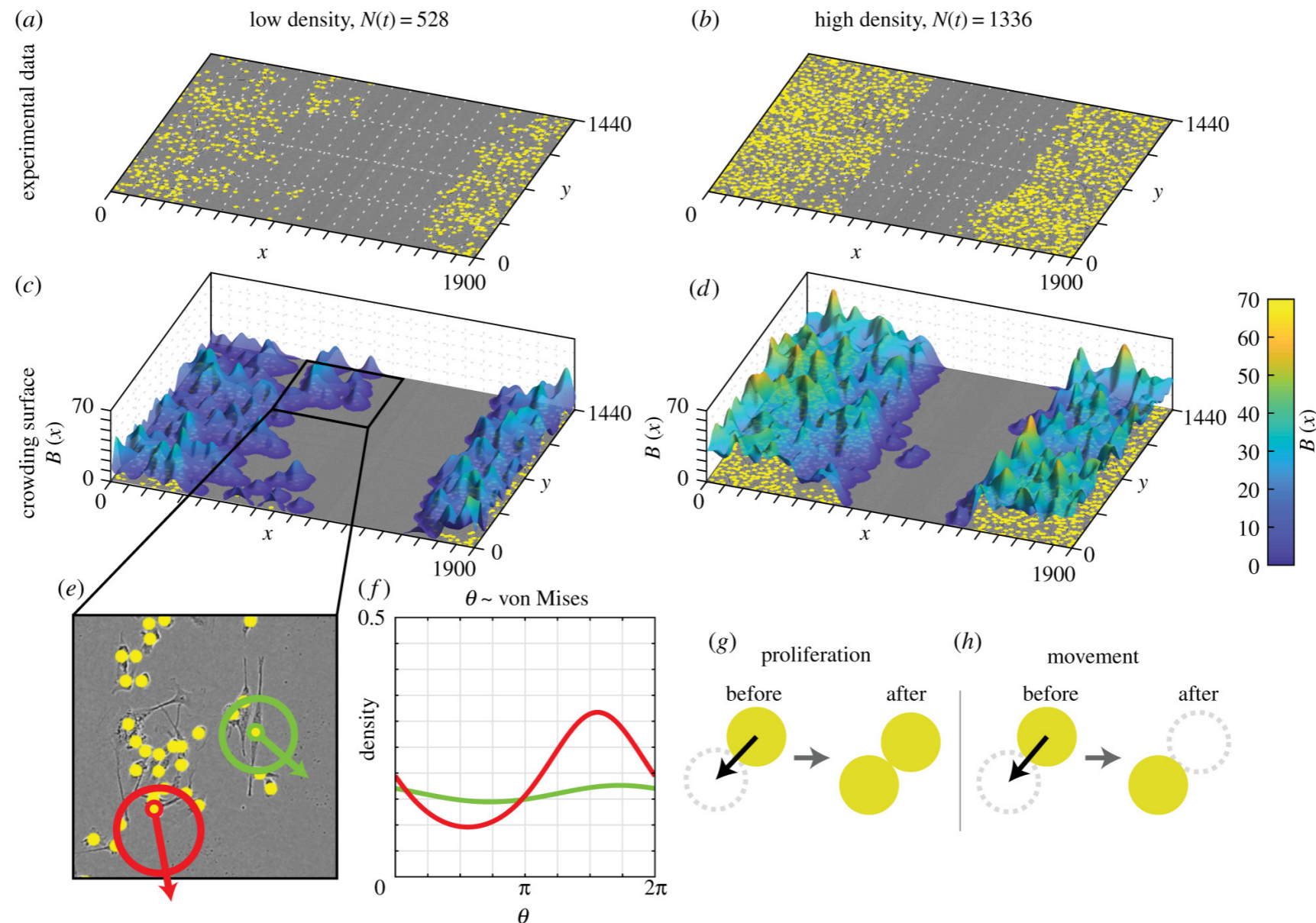  - Enables smaller continuation probabilities and hence simulation cost.

Prescott and Baker, *SIAM / ASA J. UQ* (2020) and (2021).

- Designed for use with high-throughput data i.e. large quantities of data.

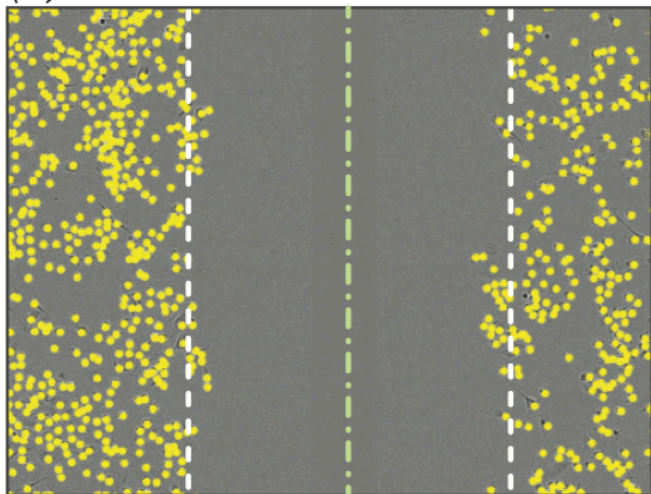- Motivated by genome-wide RNAi screen of endothelial cells.



Williams *et al. Sci. Data* 4 (2017).

- Stochastic, off-lattice individual-based model of cell motility and proliferation - including density-dependent effects.
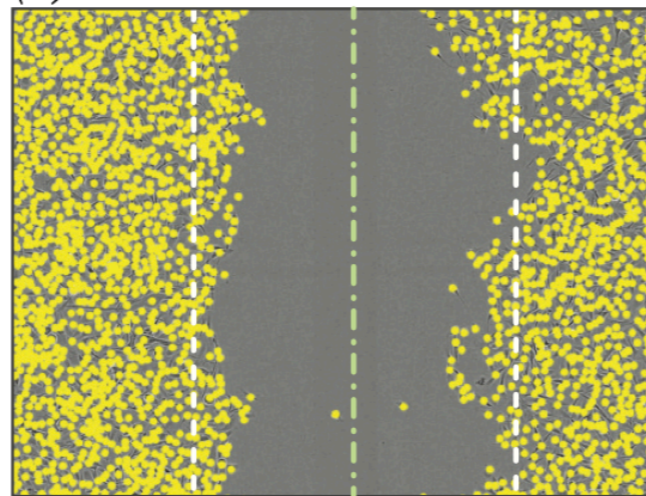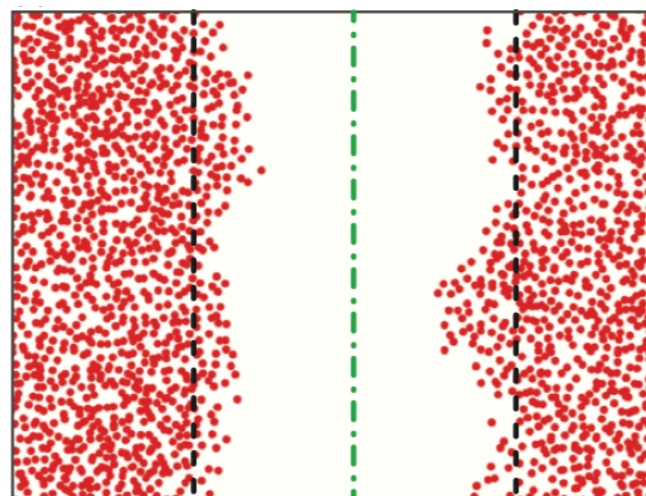


Browning *et al. J. Roy. Soc. Interface* 17 (2020).

# Model - data comparison

Data at t=24 hrs



Data at t=0 hrs



Model at t=24 hrs



Parameter inference using a Bayesian framework:

$$\mathbb{P}(\theta \,|\, \mathcal{D}) \propto \mathcal{L}(\mathcal{D} \,|\, \theta)\mathbb{P}(\theta)$$
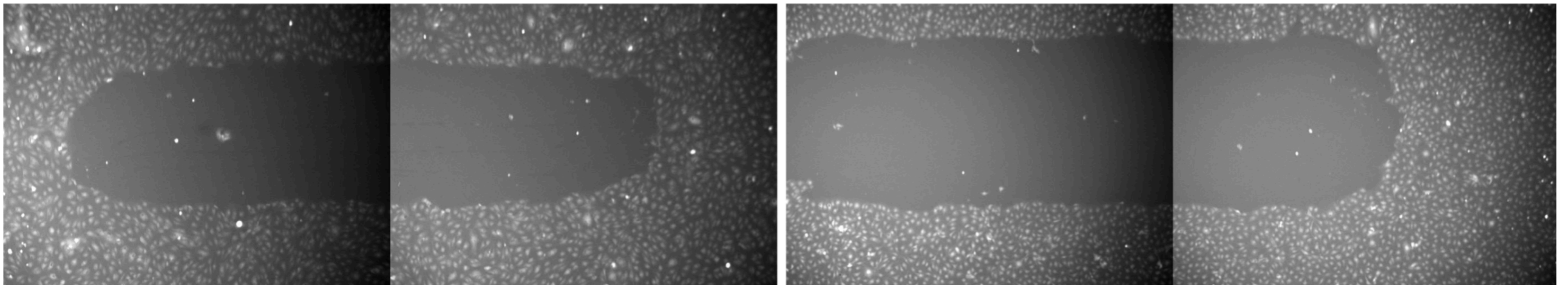
with

$$\theta = (m, p, d, \gamma_b, \gamma_m, \gamma_p)$$

Browning *et al. J. Roy. Soc. Interface* 17 (2020).

- Large numbers of replicates (~100) for some knockdowns.

  - Huge variability in the initial wound size / shape (initial condition).
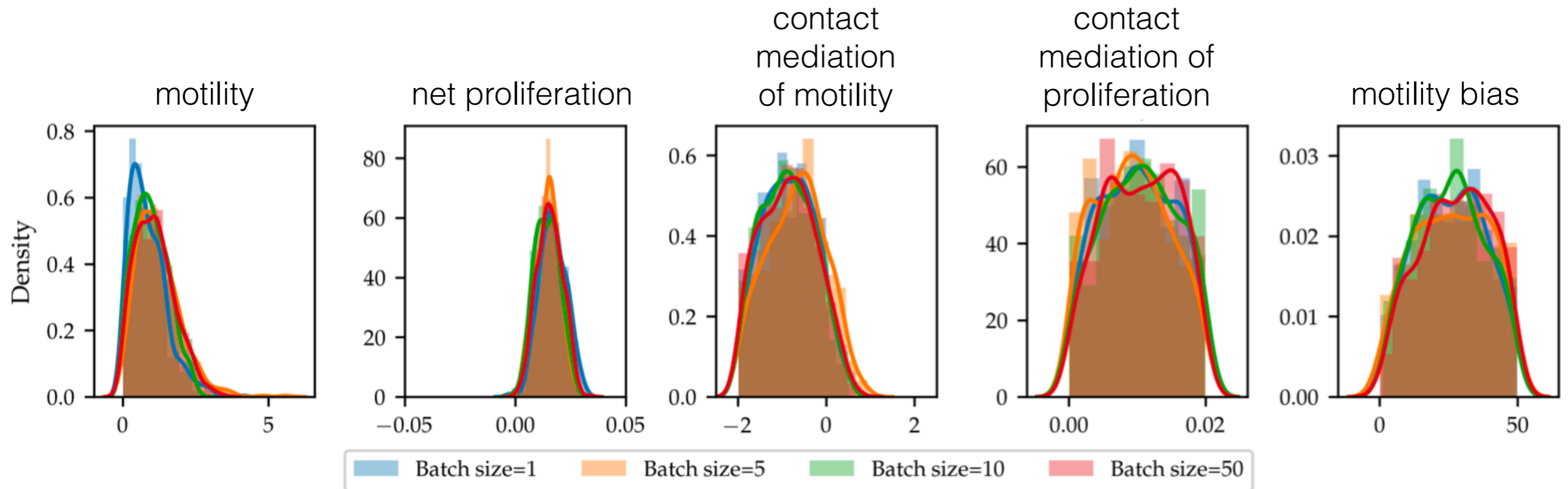


  - Need to generate individual simulations / make individual comparisons for each replicate.

- A very large number of times:

  - sample parameter from prior distribution;

  - sample a **minibatch** of the data;

  - simulate model using this parameter - **do this for each sample from the minibatch individually**;

  - evaluate how close model output is to data using summary statistics and distance function - **do this for each sample from the minibatch individually**;

  - assign a weight to the parameter - depends on this distance.

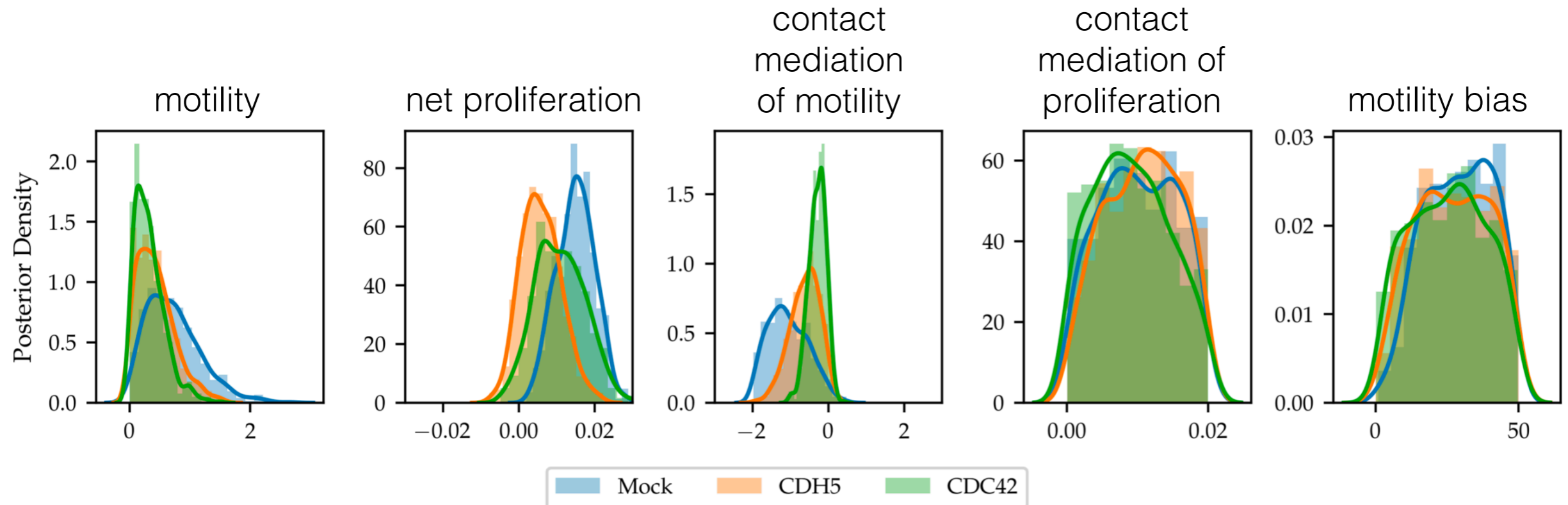Martina Perez, Sailem and Baker, *PLOS Comp. Biol.* (2022).

- Posterior distribution very similar over the different batch sizes.

- Confidence in estimates of motility, net proliferation, contact-mediation of motility parameters.

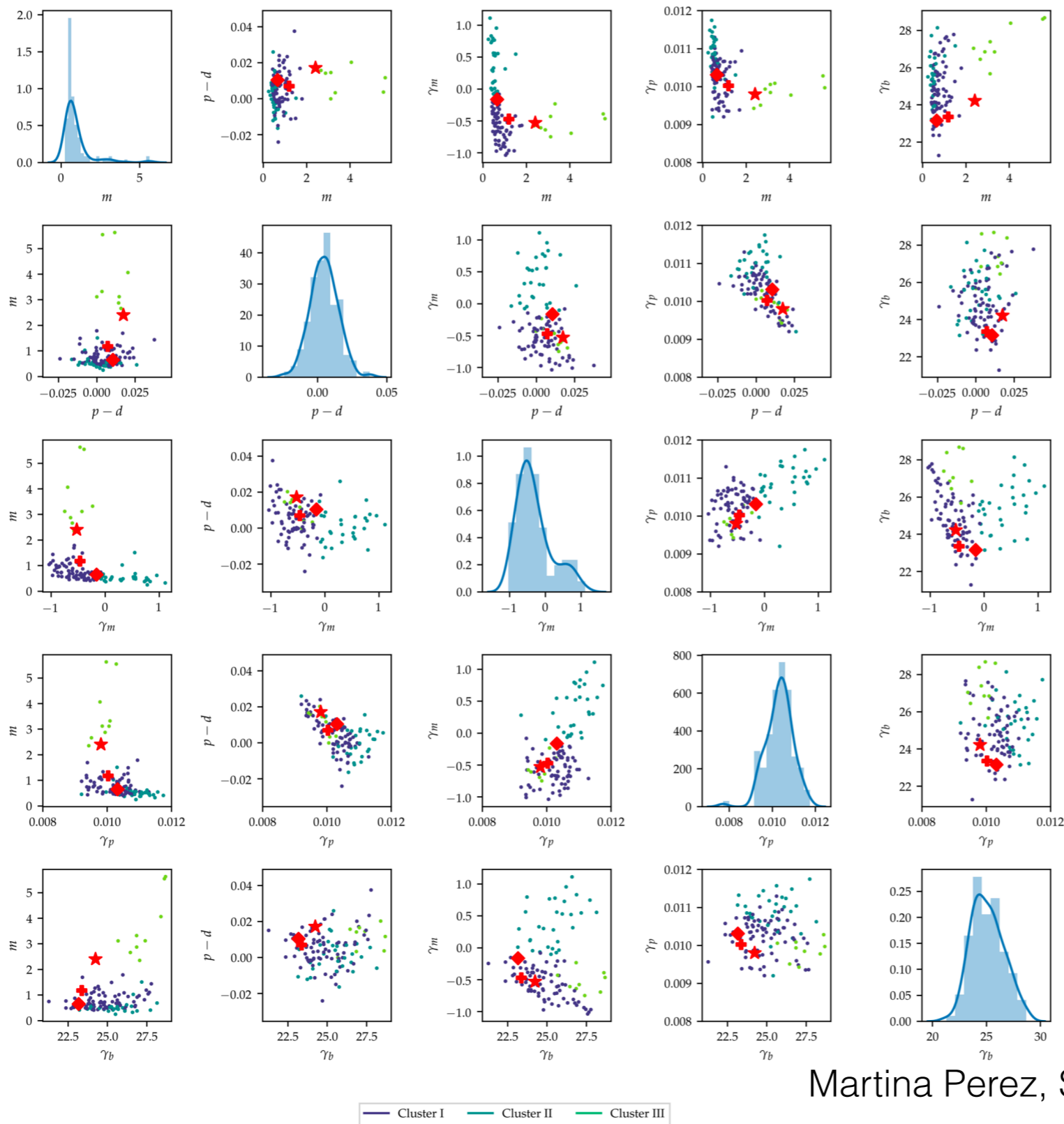Martina Perez, Sailem and Baker, *PLOS Comp. Biol.* (2022).

- Marginal posterior distributions



- Wildtype - motility strongly up-regulated in regions of high density compared to CDC42 and CDH5.

  - Consistent with current understanding of the roles of CDC42 and CDH5.

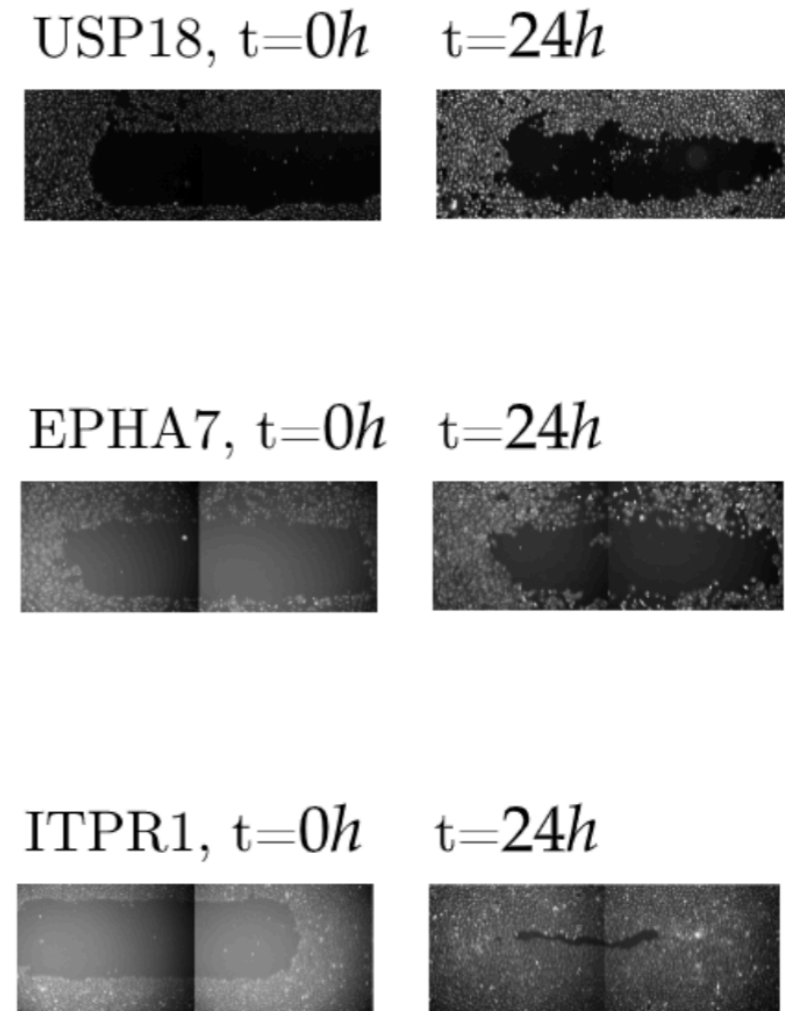Martina Perez, Sailem and Baker, *PLOS Comp. Biol.* (2022).

- Repeat for the other genes in the RNAi screen…

- Plots of the posterior means for each knockdown.

- Used K-means clustering to establish three main clusters in phenotype space.

Martina Perez, Sailem and Baker, *PLOS Comp. Biol.* (2022).

Martina Perez, Sailem and Baker, *PLOS Comp. Biol.* (2022).

# Summary

- Possible to calibrate more complicated models to data, and use them to infer greater detail on the mechanisms driving wound closure in a range of gene knockdowns.

- Complex relationships between cell motility and proliferation drive wound closure.

- Mini-batch ABC provides a means to apply ABC approaches to high-throughput datasets, without huge computational costs.

# Summary

- ABC is a fantastic tool for calibrating models to data since it relies only on forward simulation of the model.

- However, for modern mathematical biology studies the computational costs of naive forms of the method are prohibitive.

- Proposed three novel approaches to ABC - pre-conditioned, multi-fidelity and mini-batch.

- Importantly, ALL of these approaches can be combined with existing approaches e.g. ABC-SMC to provide further improvements.

- There's still much to do to optimise each approach!

# Acknowledgements

**Oxford**
Tom Prescott
Simon Martina-Perez

**Further afield**
David Warne (QUT)
Mat Simpson (QUT)

- First, need to integrate importance sampling into the multifidelity ABC framework:

---

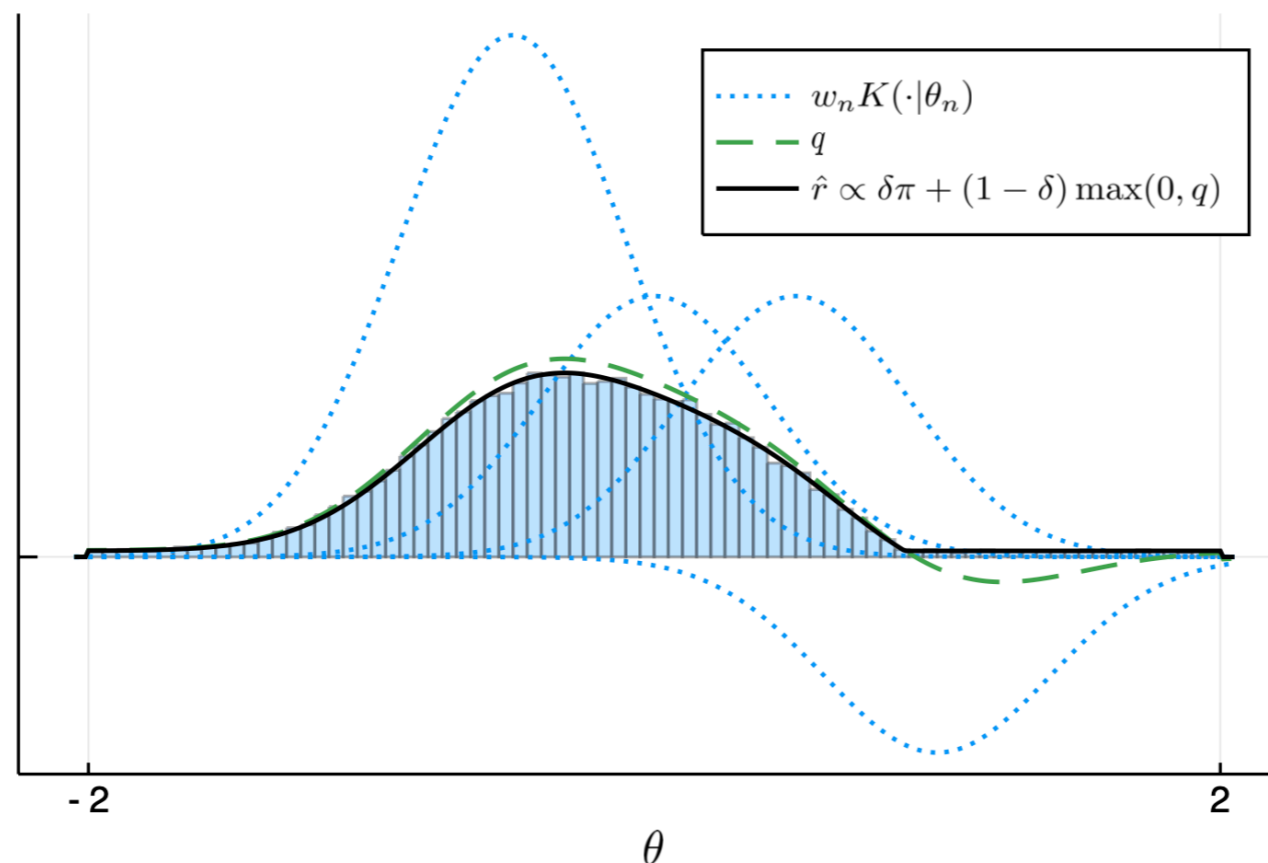**Algorithm 5** Multifidelity ABC importance sampling (MF-ABC-IS)

---

**Input:** Data $y_{\text{obs}}$ and neighbourhood $\Omega_\epsilon$; prior $\pi$; models $\tilde{f}(\cdot \mid \theta)$, $f(\cdot \mid \tilde{y}, \theta)$; continuation probability function $\alpha = \alpha(\tilde{y}, \theta)$; sample index $n = 0$; importance distribution $\hat{q}$ pr...

**Output:**

1: **rep...**
2:    In...
3:    Ge...
4:    Si...
5:    Se...
6:    Generate $u_n \sim \text{Uniform}(0,1)$.
7:   **if** $u_n < \alpha(\tilde{y}_n, \theta_n)$ **then**
8:     Simulate $y_n \sim f(\cdot \mid \tilde{y}_n, \theta_n)$.
9:     Update $w_n \leftarrow w_n + \left[\mathbb{I}(y_n \in \Omega_\epsilon) - w_n\right]/\alpha(\tilde{y}_n, \theta_n)$.
10:   **end if**
11:   Update $w_n \leftarrow \left[\pi(\theta_n)/q(\theta_n)\right] w_n$.
12: **until** $S = \texttt{true}$.

---

**For SMC: how do we sample from the importance distribution, given the weights that result from multifidelity ABC can be negative?**

- Use defensive importance sampling, first defining a new (non-negative) importance distribution.



- Estimate continuation probabilities for each generation "on the fly", using information from the previous generations.

UNIVERSITY OF OXFORD

- Kuramoto oscillator network:

angular velocities drawn
from Cauchy distribution
median - $\omega_0$
dispersion - $\gamma$

$$\dot{\phi}_i = \omega_i + \frac{K}{M} \sum_{j=1}^{M} \sin\left(\phi_j - \phi_i\right)$$

- Low-fidelity model - based on tracking Daido order parameters:

$$Z_n(t) = \frac{1}{M} \sum_{j=1}^{M} \exp(in\phi_j)$$
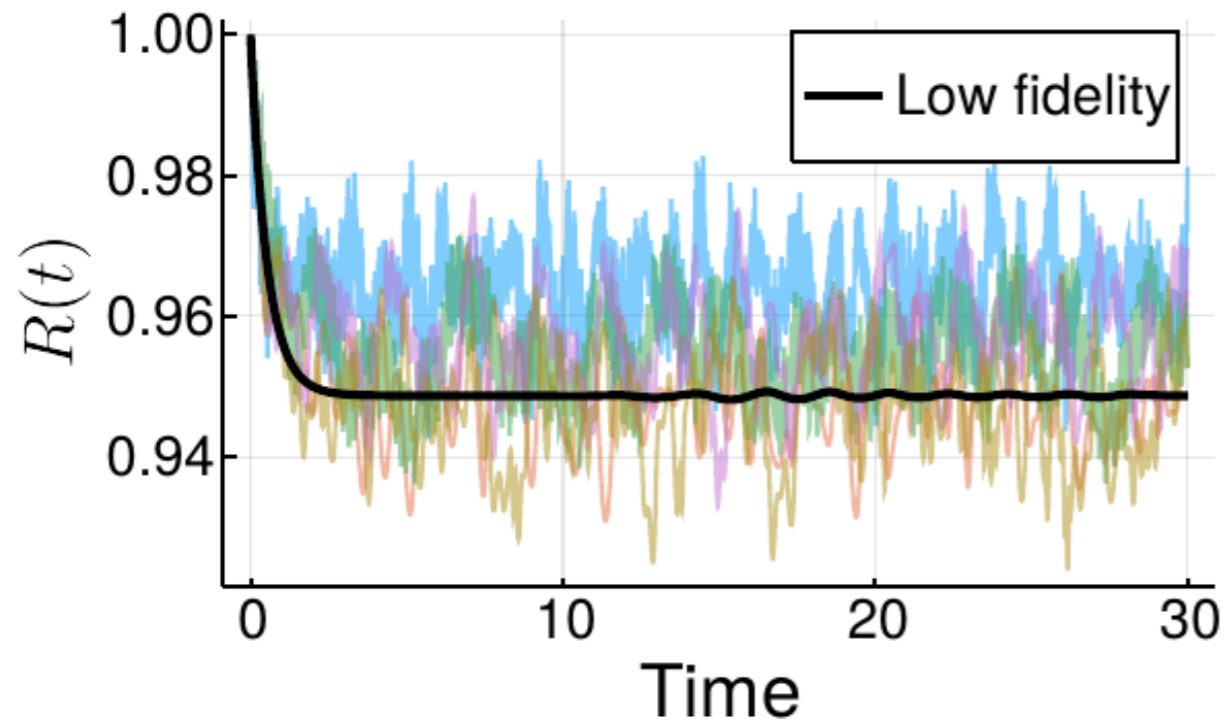
assume $\quad Z_n(t) = Z_1(t)^n$

to get
$$\dot{\tilde{R}} = \left(\frac{K}{2} - \gamma\right)\tilde{R} - \frac{K}{2}\tilde{R}^3 \qquad \text{(magnitude)}$$
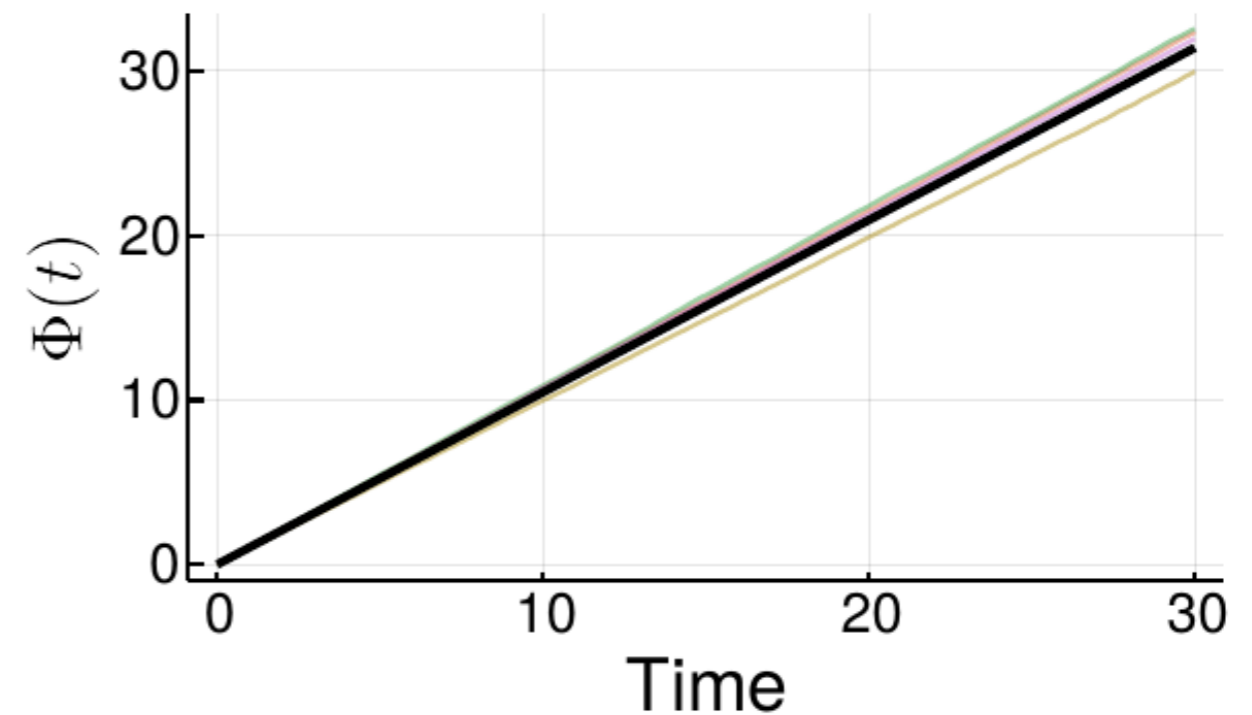$$\dot{\tilde{\Phi}} = \omega_0 \qquad \text{(phase)}$$
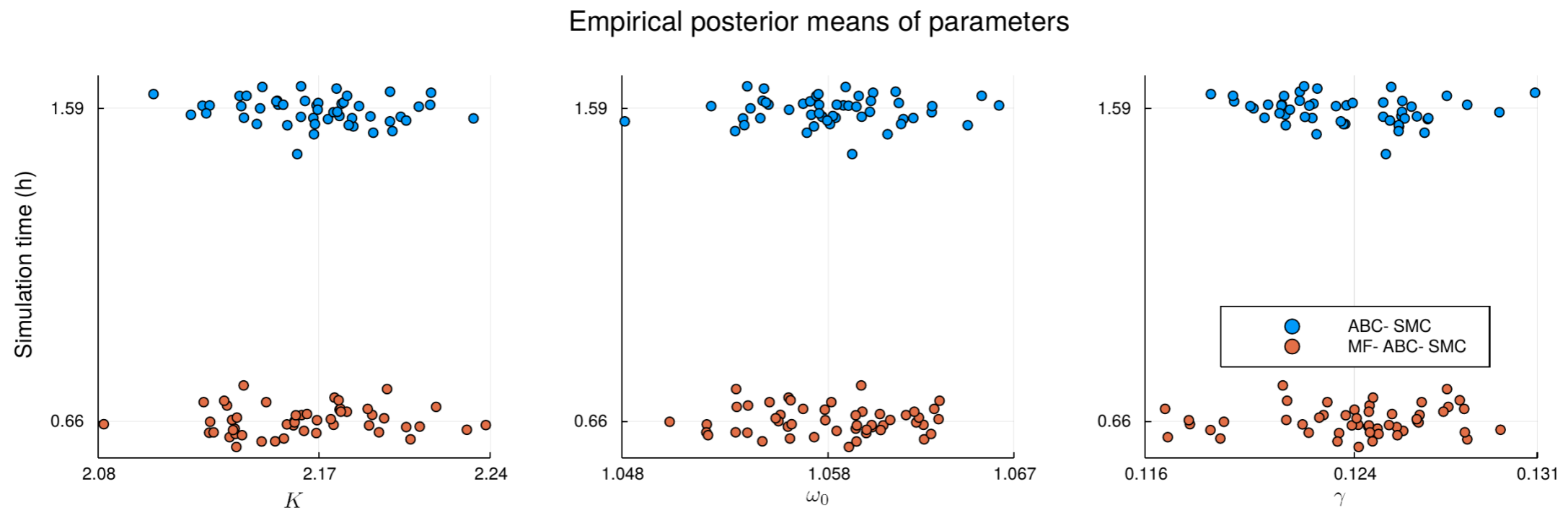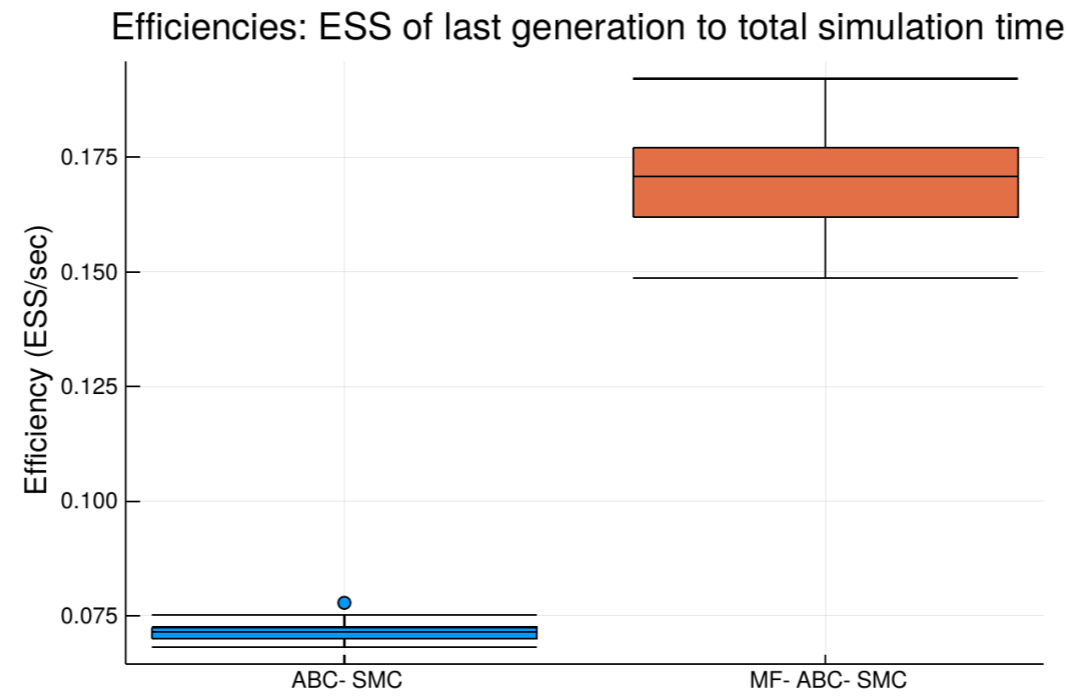
- Typical simulation output:



Kuramoto parameter: magnitude

Kuramoto parameter: phase

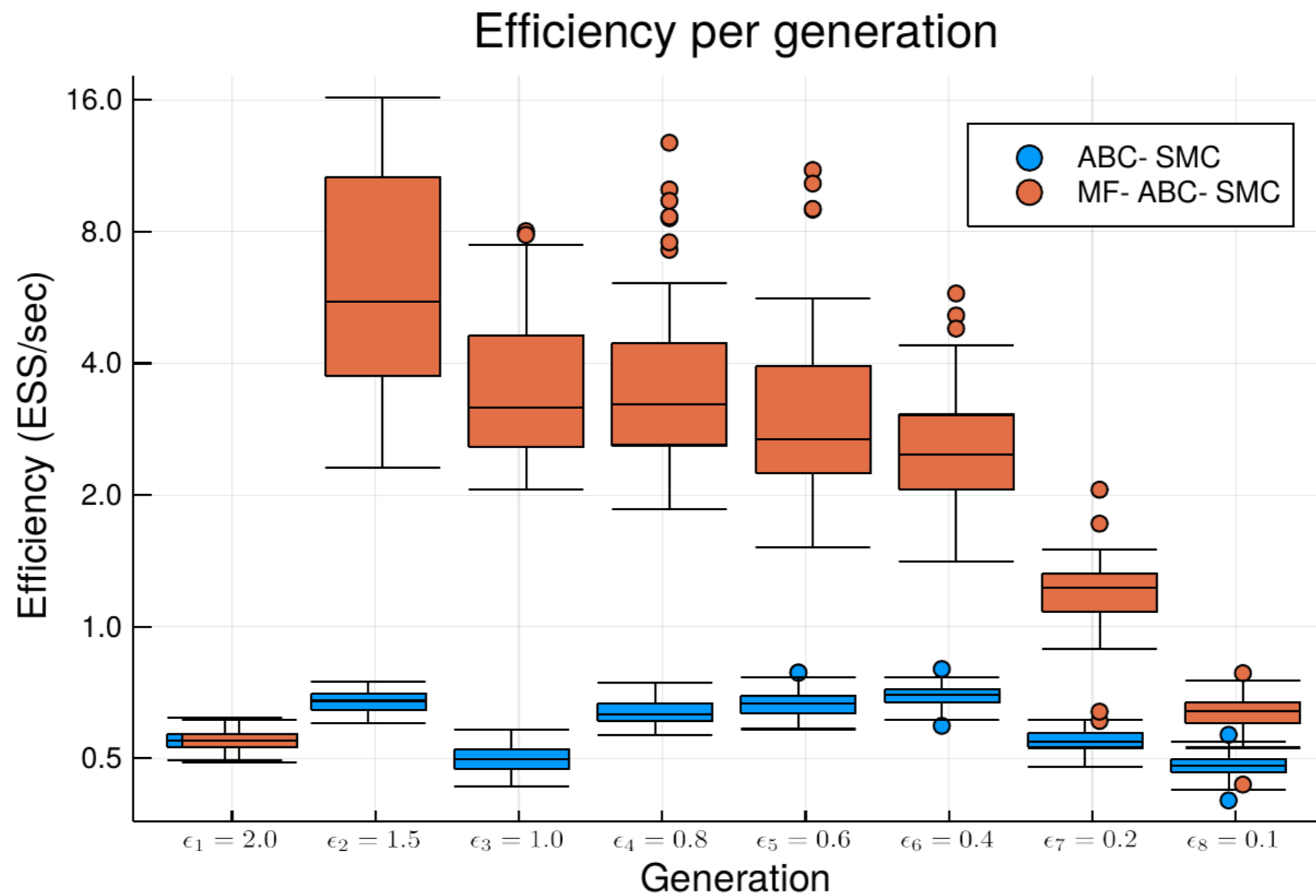Efficiencies: ESS of last generation to total simulation time

Empirical posterior means of parameters

- Stopping criterion at each generation: $\text{ESS} \geq 400$.



Efficiency per generation

Continuation probabilities by generation