### HelmholtzZentrum münchen

German Research Center for Environmental Health

## Approximate Bayesian Computation for Stochastic Single-Cell Time-Lapse Data Using Multivariate Test Statistics

Carolin Loos, PhD student ca Institute of Computational Biology Research Group Data-Driven Computational Modeling

Nantes, 16/09/15

carolin.loos@helmholtz-muenchen.de



Gene expression is known to be affected by different sources of variability<sup>1</sup>



- Gene expression is known to be affected by different sources of variability<sup>1</sup>
- Heterogeneity and stochasticity can have major consequences for the understanding of the cellular mechanisms<sup>2</sup>



- Gene expression is known to be affected by different sources of variability<sup>1</sup>
- Heterogeneity and stochasticity can have major consequences for the understanding of the cellular mechanisms<sup>2</sup>
- Single-cell data for analysis of variability required





Single-cell snapshot data

Hasenauer, PhD thesis (2013)



Single-cell snapshot data

Single-cell time-lapse data



Single-cell time-lapse data provide information about temporal cross-correlation

Derive model M with parameters  $\theta$  and fit it to data D by e.g. maximizing the likelihood function or sampling from the posterior distribution p( $\theta$ |D)

Derive model M with parameters  $\theta$  and fit it to data D by e.g. maximizing the likelihood function or sampling from the posterior distribution p( $\theta$ |D)

- D: Single-cell time-series  $\bm{x}_1,\ldots,\bm{x}_n$ 

Derive model M with parameters  $\theta$  and fit it to data D by e.g. maximizing the likelihood function or sampling from the posterior distribution p( $\theta$ |D)

- D: Single-cell time-series  $\bm{x}_1,\ldots,\bm{x}_n$
- M: Single-cell time-lapse data modeled with continuous time Markov chains (CTMCs)<sup>1</sup>

Derive model M with parameters  $\theta$  and fit it to data D by e.g. maximizing the likelihood function or sampling from the posterior distribution p( $\theta$ |D)

- D: Single-cell time-series  $\bm{x}_1,\ldots,\bm{x}_n$
- M: Single-cell time-lapse data modeled with continuous time Markov chains (CTMCs)<sup>1</sup>
- $\theta$ : Parameters of process e.g. kinetic parameters, initial conditions, ...

 Evaluation of likelihood function for CTMCs often computationally too costly

- Evaluation of likelihood function for CTMCs often computationally too costly
- Approximate Bayesian Computation (ABC) methods required

- Evaluation of likelihood function for CTMCs often computationally too costly
- Approximate Bayesian Computation (ABC) methods required
- ABC circumvents calculation of likelihood by comparing observed and simulated data set



Goal: Approximate posterior distribution of  $\theta$  given data D

Figure modified from Toni et al., Bioinformatics (2010)



Goal: Approximate posterior distribution of  $\theta$  given data D

1. Sample from prior distribution

Figure modified from Toni et al., Bioinformatics (2010)



Goal: Approximate posterior distribution of  $\theta$  given data D

- 1. Sample from prior distribution
- 2. Simulate data set

Figure modified from Toni et al., Bioinformatics (2010)





# ime series of a birth-death process



### **Approximate Bayesian Computation with Sequential Monte Carlo (ABC SMC)**

### Approxim Monte Ca



### ABC SMC (Sequential Monte Carlo)



### Approxi Monte C



ABC SMC (Sequential Monte Carlo)



7

### Approxi Monte C



### ABC SMC (Sequential Monte Carlo)



7

### Approxi Monte C



### ABC SMC (Sequential Monte Carlo)



7

ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>

- ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>
- Independent measurements for different time points

- ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>
- Independent measurements for different time points
- Maximal Kolmogorov-Smirnov (KS) distance for every time point

- ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>
- Independent measurements for different time points
- Maximal Kolmogorov-Smirnov (KS) distance for every time point



- ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>
- Independent measurements for different time points
- Maximal Kolmogorov-Smirnov (KS) distance for every time point



- ABC SMC has been used to analyse flow cytometry (snapshot) data<sup>1</sup>
- Independent measurements for different time points
- Maximal Kolmogorov-Smirnov (KS) distance for every time point



Accept parameters that produce a low KS distance
Performance and convergence of ABC highly depends on the employed distance measure

- Performance and convergence of ABC highly depends on the employed distance measure
- Project single-cell time-series into high dimensional space

- Performance and convergence of ABC highly depends on the employed distance measure
- Project single-cell time-series into high dimensional space



- Performance and convergence of ABC highly depends on the employed distance measure
- Project single-cell time-series into high dimensional space



- Performance and convergence of ABC highly depends on the employed distance measure
- Project single-cell time-series into high dimensional space





> observed samples> simulated samples



> observed samples> simulated samples

Perform a minimum weight non-bipartite matching on a complete graph



Perform a minimum weight non-bipartite matching on a complete graph



Perform a minimum weight non-bipartite matching on a complete graph



Perform a minimum weight non-bipartite matching on a complete graph



Perform a minimum weight non-bipartite matching on a complete graph

Accept parameters that produce more cross-matches

$$\mathsf{MMD} = \left(\frac{1}{n^2} \sum_{i \neq j}^n k(\mathbf{x}_i, \mathbf{x}_j) + \frac{1}{m^2} \sum_{i \neq j}^m k(\mathbf{y}_i, \mathbf{y}_j) - \frac{2}{nm} \sum_{i,j=1}^{n,m} k(\mathbf{x}_i, \mathbf{y}_j)\right)^{\frac{1}{2}}$$





> observed samples x<sub>1</sub>,...,x<sub>n</sub>
> simulated samples y<sub>1</sub>,...,y<sub>m</sub>



> observed samples x<sub>1</sub>,...,x<sub>n</sub>
> simulated samples y<sub>1</sub>,...,y<sub>m</sub>



> observed samples x<sub>1</sub>,...,x<sub>n</sub>
> simulated samples y<sub>1</sub>,...,y<sub>m</sub>

Accept parameters that produce a small MMD

### **Example: Bivariate Normal Distribution**

### **Example: Bivariate Normal Distribution**

Data: 100 samples of a bivariate normal distribution  $\bar{y} \sim N (\mu, \Sigma)$ with  $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  and  $\Sigma = \begin{pmatrix} \theta_1 & \theta_2 \\ \theta_2 & \theta_1 \end{pmatrix}$ 



Data: 100 samples of a bivariate normal distribution  $\bar{y} \sim N (\mu, \Sigma)$ with  $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  and  $\Sigma = \begin{pmatrix} \theta_1 & \theta_2 \\ \theta_2 & \theta_1 \end{pmatrix}$ 

Goal: Approximate posterior distributions of  $\theta_1$  and  $\theta_2$  using ABC SMC



Data: 100 samples of a bivariate normal distribution  $\bar{y} \sim N (\mu, \Sigma)$ with  $\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$  and  $\Sigma = \begin{pmatrix} \theta_1 & \theta_2 \\ \theta_2 & \theta_1 \end{pmatrix}$ 

Goal: Approximate posterior distributions of  $\theta_1$  and  $\theta_2$  using ABC SMC



Compare ABC SMC with cross-match test (CM), maximum mean discrepancy (MMD), and maximal Kolmogorov-Smirnov distance (KS)





• For KS only  $\theta_1$  is identifiable



- For KS only  $\theta_1$  is identifiable
- CM gives a wider posterior approximation



- For KS only  $\theta_1$  is identifiable
- CM gives a wider posterior approximation
- Only MMD gives a reasonable posterior approximation

- One stage model of gene expression with mRNA synthesis rate  $\lambda$  and degradation rate  $\gamma$ 

- One stage model of gene expression with mRNA synthesis rate  $\lambda$  and degradation rate  $\gamma$ 



- One stage model of gene expression with mRNA synthesis rate  $\lambda$  and degradation rate  $\gamma$
- Comparison of ABC SMC with multivariate statistics (MMD) and univariate statistics (KS)



### ngle-cell time-lapse data



### ngle-cell time-lapse data



### ngle-cell time-lapse data



Generate time-series using the Stochastic Simulation Algorithm<sup>1</sup>

**1** Gillespie, Annual Review of Physical Chemistry (2007)
# ngle-cell time-lapse data



- Generate time-series using the Stochastic Simulation Algorithm<sup>1</sup>
- Consider 10, 100, 1000 cells and equilibrium and non equilibrium time-series













 Comparison of posterior approximation with Finite State Projection (FSP)<sup>1</sup>



- Comparison of posterior approximation with Finite State Projection (FSP)<sup>1</sup>
- FSP gives narrower posterior approximations compared to MMD and KS

**1** Munsky et al., J. Chem. Phys. (2006)













Bigger difference between FSP and MMD for equilibrium data

## **Tree-structured data**



gray. (d) Fitted mean and variance of the autocorrelation function. (e) Comparison



gray. (d) Fitted mean and variance of the autocorrelation function. (e) Comparison



gray. (d) Fitted mean and variance of the autocorrelation function. (e) Comparison



gray. (d) Fitted mean and variance of the autocorrelation function. (e) Comparison

## **Tree-structured data**





n

# Outlook

• Evaluate information content of lineage information

- Evaluate information content of lineage information
- Incorporate efficient simulations

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>
  - Method of conditional moments<sup>2</sup>

1 Gillespie, Annual Review of Physical Chemistry (2007)2 Hasenauer et al., Journal of Mathematical Biology (2014)

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>
  - Method of conditional moments<sup>2</sup>
- Further improvements by tuning parameters of ABC

1 Gillespie, Annual Review of Physical Chemistry (2007)2 Hasenauer et al., Journal of Mathematical Biology (2014)

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>
  - Method of conditional moments<sup>2</sup>
- Further improvements by tuning parameters of ABC
  - Appropriate threshold sequence<sup>3</sup>

Gillespie, Annual Review of Physical Chemistry (2007)
 Hasenauer et al., Journal of Mathematical Biology (2014)
 Silk et al., Stat. Appl. Genet. Mol. (2013)

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>
  - Method of conditional moments<sup>2</sup>
- Further improvements by tuning parameters of ABC
  - Appropriate threshold sequence<sup>3</sup>
  - Study perturbation kernels

1 Gillespie, Annual Review of Physical Chemistry (2007)
2 Hasenauer et al., Journal of Mathematical Biology (2014)
2 Gille et al., Const. Mat. (2012)

**3** Silk et al., Stat. Appl. Genet. Mol. (2013)

- Evaluate information content of lineage information
- Incorporate efficient simulations
  - Tau-leaping<sup>1</sup>
  - Method of conditional moments<sup>2</sup>
- Further improvements by tuning parameters of ABC
  - Appropriate threshold sequence<sup>3</sup>
  - Study perturbation kernels
- Study adaptive number of simulations

**1** Gillespie, Annual Review of Physical Chemistry (2007)

**2** Hasenauer et al., Journal of Mathematical Biology (2014)

**3** Silk et al., Stat. Appl. Genet. Mol. (2013)

# Summary

• Method to analyse single-cell time-lapse data

- Method to analyse single-cell time-lapse data
- Account for temporal information of individual cells

- Method to analyse single-cell time-lapse data
- Account for temporal information of individual cells
- Identifiability in several cases only with multivariate statistics
- Method to analyse single-cell time-lapse data
- Account for temporal information of individual cells
- Identifiability in several cases only with multivariate statistics
- MMD is a suitable test statistic

- Method to analyse single-cell time-lapse data
- Account for temporal information of individual cells
- Identifiability in several cases only with multivariate statistics
- MMD is a suitable test statistic
- Enables incorporation of lineage information

- Method to analyse single-cell time-lapse data
- Account for temporal information of individual cells
- Identifiability in several cases only with multivariate statistics
- MMD is a suitable test statistic
- Enables incorporation of lineage information

Flexible framework for the analysis of single-cell time-lapse data, which might help getting deeper insight into the cellular mechanisms and therefore advance e.g. stem cell research.

## Acknowledgement



**Institute of Computational Biology** Data-Driven Computational Modelling Jan Hasenauer Carsten Marr Fabian Theis Dennis Rickert

## HelmholtzZentrum münchen

Deutsches Forschungszentrum für Gesundheit und Umwelt



## Acknowledgement



**Institute of Computational Biology** Data-Driven Computational Modelling Jan Hasenauer Carsten Marr Fabian Theis Dennis Rickert

## HelmholtzZentrum münchen

Deutsches Forschungszentrum für Gesundheit und Umwelt



carolin.loos@helmholtz-muenchen.de