

IMPERIAL

Inference methods for studying single cells a tutorial lecture

Philipp Thomas
June 2024

Resources



Tutorial 1 Machine learning-based inference

Tutorial 2 Approximate Bayesian Computation

Tutorial 3 Model selection

Can be found at:

<https://shorturl.at/hDps9>

Python requirements:

numpy, scipy, matplotlib, pandas, scikit-learn

Single-cell data

Timelapse data of bacteria expressing green fluorescent protein

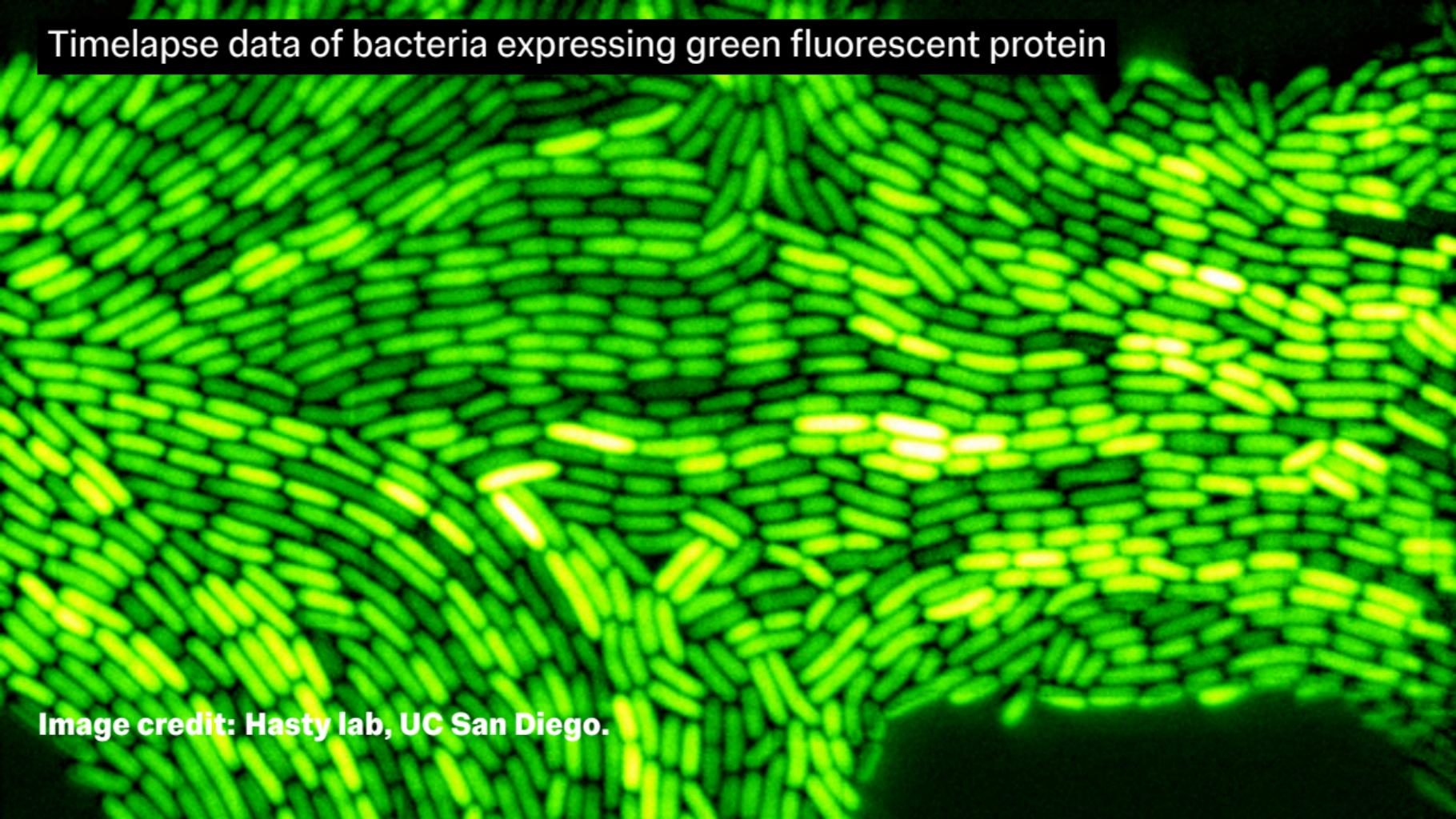


Image credit: Hasty lab, UC San Diego.

Snapshot of individual transcripts in bacteria (smFISH)

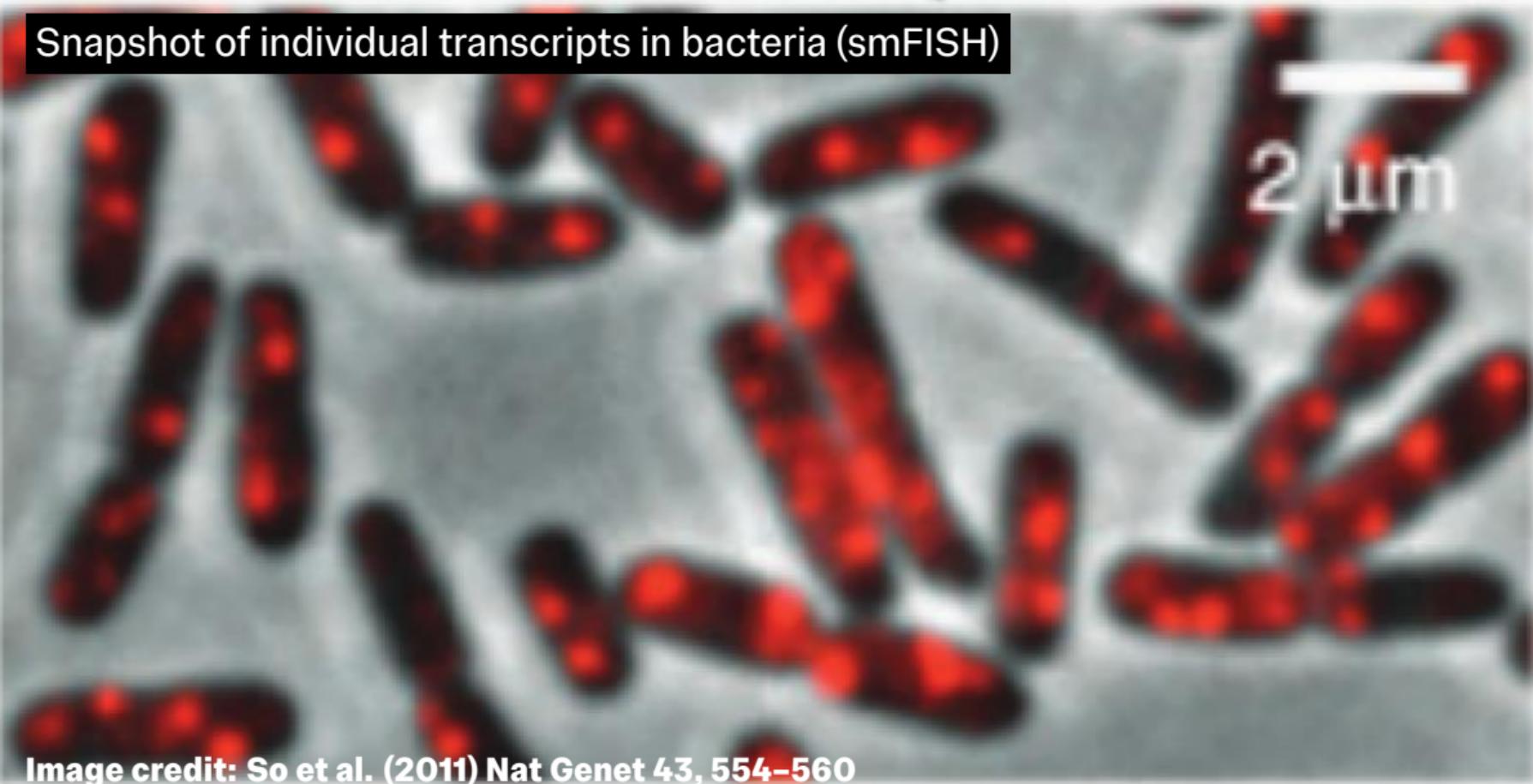
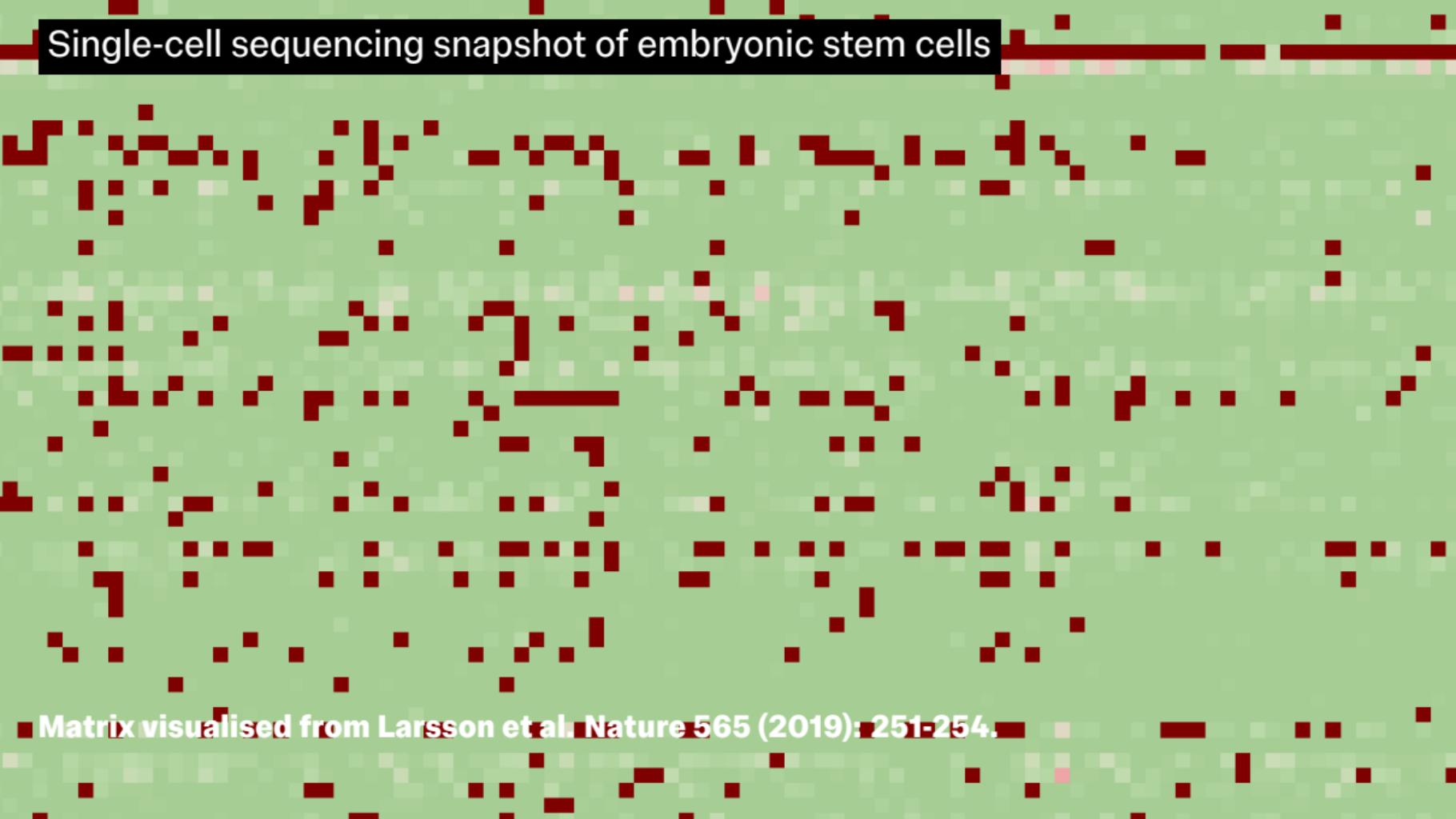


Image credit: So et al. (2011) Nat Genet 43, 554–560

Single-cell sequencing snapshot of embryonic stem cells



Matrix visualised from Larsson et al. *Nature* 565 (2019): 251-254.

Single-cell data

Focus: analysis of snapshot data

high-throughput

transcript numbers for thousands of genes in thousands of cells

beyond average

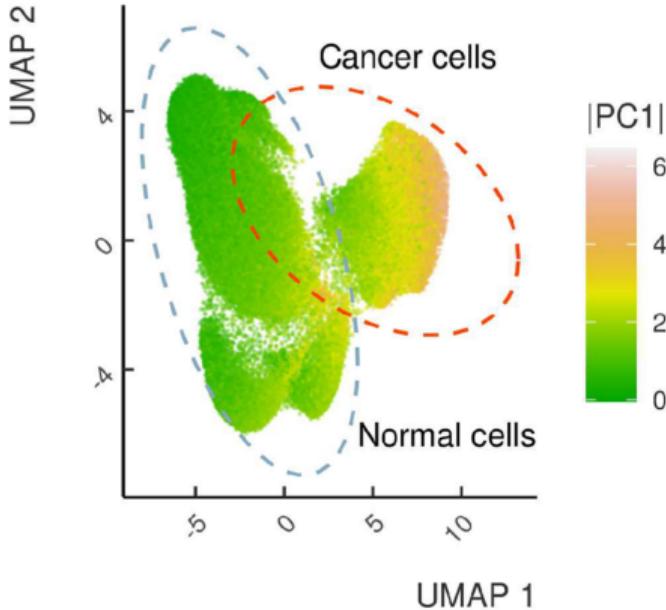
harness cell-to-cell variation (distributions) to understand processes inside individual cells

mathematical models

capture distributions to make sense of data and biological mechanisms

inference

parameterise models, which allows testing hypotheses with data



Agenda

- 01** Data, modelling and likelihood
- 02** Maximum likelihood estimation
- 03** Moment-based inference
- 04** Machine learning-assisted inference
- 05** Bayesian inference
- 06** Approximate Bayesian Computation
- 07** Model selection
- 08** Applications to real data

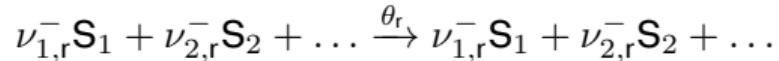
Galileo Galilei: “Measure what can be measured, and make measurable what cannot be measured.”

Modelling gene expression

Modelling gene expression

Chemical Master Equation

Biochemical kinetics of a general reaction:



Master equation:

$$\frac{dp}{dt}(m|\theta) = \sum_r [w_r(m - \nu_r, \theta_r)p(m - \nu_r|\theta) - w_r(m, \theta_r)p(m|\theta)],$$

where $\nu_r = \nu_r^+ - \nu_r^-$ is the net-change of an reaction.

Infinite system of coupled ordinary differential equations:
cannot generally be solved.

Modelling gene expression

Poisson model

Biochemical kinetics of transcription and degradation:



Master equation:

$$\dot{p}(m) = k_1 p(m-1) - k_1 p(m) + k_2(m+1)p(m+1) - k_2 m p(m),$$

In steady state:

$$\dot{p}(m) = 0 \implies 0 = k_1 p(m-1) - k_2 m p(m)$$

Poisson distribution:

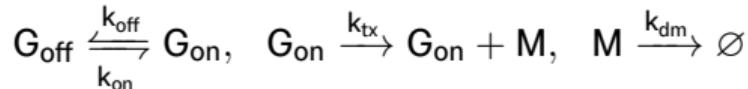
$$P(m) = p_0 \times \mu \times \mu/2 \times \mu/3 \times \dots = p_0 \frac{\mu^m}{m!},$$

where $\mu = \frac{k_1}{k_2}$ and $p_0 = e^{-\mu}$.

Modelling gene expression

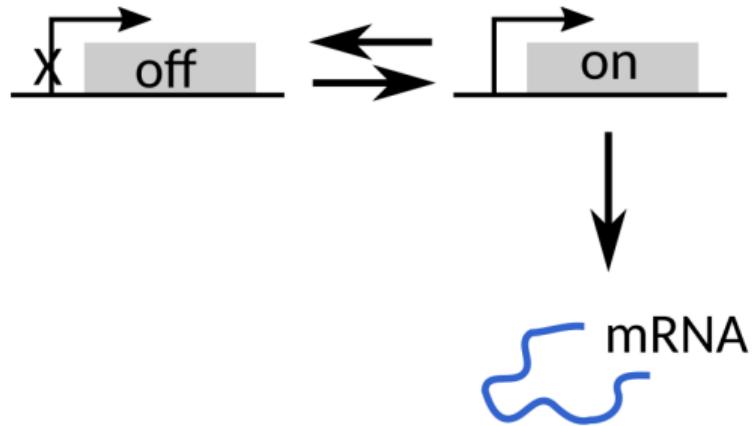
Telegraph model

Biochemical kinetics of transcription and degradation:



Solved by Beta-Poisson distribution ($k_2 = 1$):

$$m|p \sim \text{Poisson}(pk_{\text{tx}}), \quad p \sim \text{Beta}(k_{\text{on}}, k_{\text{off}})$$



Determinants of gene expression noise

- **burst frequency** $a = k_{\text{on}}/k_{\text{dm}}$
- **burst size** $b = k_{\text{tx}}/k_{\text{off}}$
- **switching frequency** $s = k_{\text{on}}/k_{\text{off}}$

Modelling gene expression

What's the likelihood?

Definition: Assuming independent data points, the probability that the data X was generated under the model θ is:

$$\mathcal{L}(\theta|X) = p(x_1|\theta) \times p(x_2|\theta) \times \dots \times p(x_N|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

For example, we have two parameters θ_1 and θ_2 then if

$$\mathcal{L}(\theta_1|X) > \mathcal{L}(\theta_2|X)$$

then θ_1 fits the data "better".

- Likelihood includes effects of
 - sampling noise (N), i.e., measurement could also arise by chance
 - biological noise (p), i.e., variability across
- Likelihood measures **goodness of fit** rather than the probability that θ is the truth.

Maximum likelihood principle

The maximum likelihood principle is a method for estimating parameters in statistical models. It involves finding the parameter values that make the observed data most probable. This is typically done by maximizing the likelihood function, which is the probability of the observed data given the parameters.

The likelihood function is defined as:

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

where θ is the vector of parameters, x_i are the observed data points, and $f(x_i; \theta)$ is the probability density or mass function of the data given the parameters.

To find the maximum likelihood estimates, we take the natural logarithm of the likelihood function and then differentiate it with respect to the parameters. This results in the log-likelihood function:

$$\ell(\theta) = \ln L(\theta) = \sum_{i=1}^n \ln f(x_i; \theta)$$

We then find the values of θ that maximize this function. This can be done using various optimization techniques, such as gradient descent or Newton-Raphson methods.

The maximum likelihood principle has many applications in statistics and machine learning, including parameter estimation, model selection, and hypothesis testing.

Maximum likelihood principle

Negative log-likelihood acts as an error function:

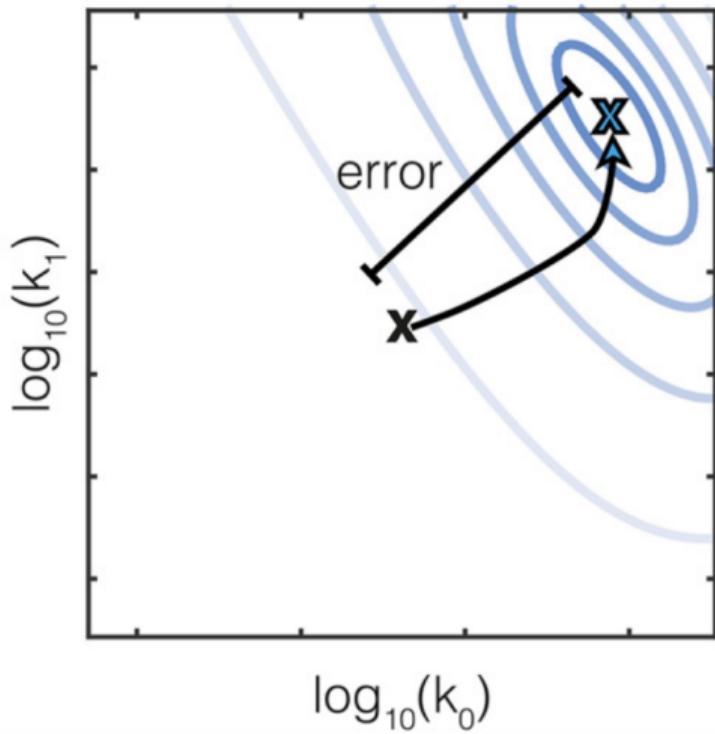
$$\widehat{\text{err}}(\theta; \mathbf{x}) = -\ln \mathcal{L}(\theta | \mathbf{X}) = -\sum_{i=1}^N \ln p(x_i | \theta)$$

For example, assume normal distribution $p(x_i | \theta) \propto e^{-(x_i - \theta)^2}$ with mean θ , then

$$\widehat{\text{err}}(\theta; \mathbf{x}) = \sum_{i=1}^N (x_i - \theta)^2 + \text{const}$$

Idea: Maximising the likelihood is minimising the error between the model and data.

$$\theta^* = \operatorname{argmax}_{\theta} \mathcal{L}(\theta | \mathbf{X})$$



log-likelihood surface: blue cross is maximum of likelihood

Maximum likelihood principle

Example: Birth-death process

Poisson distribution: $P(m|\theta) = e^{-\mu} \frac{\mu^m}{m!}$

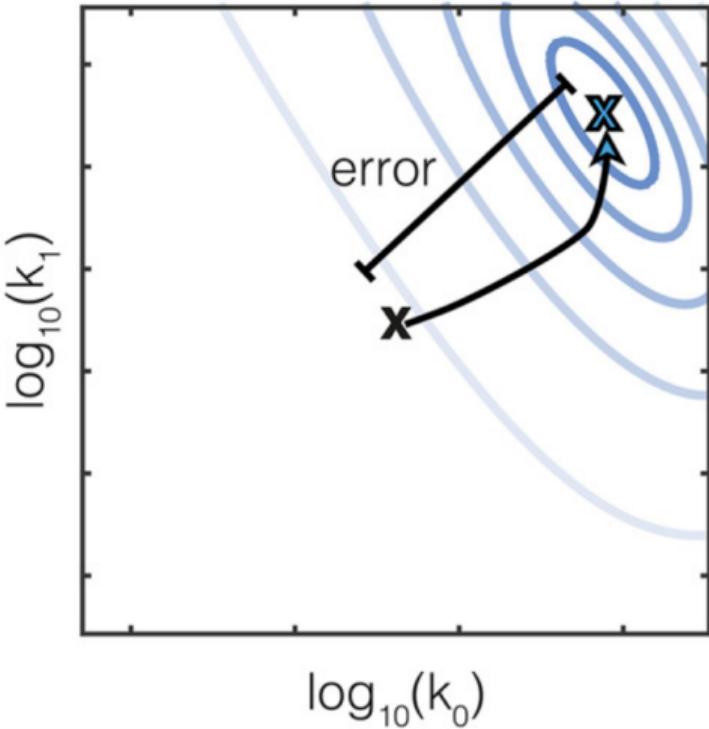
The error is

$$\widehat{\text{err}}(\theta; x) = \sum_{i=1}^N (x_i \ln \mu - \mu) + \text{const}$$

Minimising this gives

$$\nabla_\mu \widehat{\text{err}}(\theta; x) = \sum_{i=1}^N (x_i/\mu - 1) = 0$$

and hence $\mu^* = \frac{1}{N} \sum_{i=1}^N x_i$



log-likelihood surface: blue cross is maximum of likelihood

Maximum likelihood principle

Example: Telegraph process

Solve the Chemical Master Equation:

Beta-Poisson distribution^a

$$P(m|\theta) = \frac{k_1^m \Gamma(m+k_{on}) \Gamma(k_{off}+k_{on}) {}_1F_1(m+k_{on}; m+k_{off}+k_{on}; -k_{tx})}{\Gamma(m+1) \Gamma(k_{on}) \Gamma(m+k_{off}+k_{on})}$$

Likelihood can be obtained in terms of special functions.

Maximise likelihood:

$$\nabla \widehat{\text{err}}(\theta; x) = - \sum_i \ln p(x_i | \theta)$$

has no closed-form solution.

^aRaj et al. PLoS Biology 4 (2006): e309.

Maximum likelihood principle

Error guarantees

Asymptotic error guarantees (large sample size):

- Unbiased
- Smallest variances

Limitations:

- Implicitly assumes large sample size (how much is sufficient?)
- Computationally intractable optimisation; non-convex problem can suffer local likelihood maxima
- Computationally intractable likelihoods

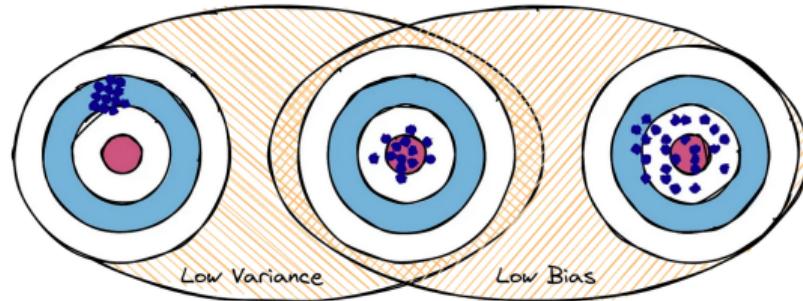


Illustration of bias and variance
(Credit: Ivan Reznikov on medium.com)

Moment-based inference

Moment-based inference

Approximation methods

Problem: Likelihood is not available

Idea: Distribution of $\hat{\mu} = \frac{1}{N} \sum_i x_i$ and $\hat{\Sigma} = \frac{1}{N} \sum_i (x_i - \mu)^2$ is asymptotically normal

Likelihood of moment data:

$$(\hat{\mu}, \hat{\Sigma}) \sim \text{Normal}((\mu(\theta), \Sigma(\theta)), s(N))$$

where s is the matrix standard errors of mean and variance:

$$s(\theta, N) = \frac{1}{N} \begin{bmatrix} \Sigma & E[(x - \mu)^3] \\ E[(x - \mu)^3] & E[(x - \mu)^4] - \frac{n-3}{n-1} \Sigma^2 \end{bmatrix}$$

This means

$$\begin{aligned} -2 \ln \mathcal{L}(\theta) = & \ln(|s(\theta, N)|) \\ & + (\hat{\mu} - \mu(\theta)) s_{11}^{-1}(\theta) (\hat{\mu} - \mu(\theta))^T \\ & + (\hat{\Sigma} - \Sigma(\theta)) s_{22}^{-1}(\theta) (\hat{\Sigma} - \Sigma(\theta))^T + \dots \end{aligned}$$

Observations:

- Readily generalised to higher order moments
- Roughly minimised by matching mean and variance of model and data
- Penalty for high variance of summary statistics (N small)

Moment-based inference

Moments of telegraph process

- Obtain equations for the moments:

$$m_i(\theta) = E_p[x^i | \theta]$$

- Solve them in steady state:

$$m_1 = \frac{k_1 k_{on}}{k_{off} + k_{on}}$$

$$m_2 = m_1 \frac{k_1 (k_{on} + 1)}{(k_{off} + k_{on} + 1)}$$

$$m_3 = m_2 \frac{k_1 (k_{on} + 2)}{(k_{off} + k_{on} + 2)}$$

- Express parameters in terms of moments

Result:

burst frequency:

$$a = \frac{(2m_1^2 - m_2) m_3 - m_1 m_2^2}{m_2 (m_1^2 - 2m_2) + m_1 m_3},$$

burst size:

$$b = m_1 \left(\frac{m_1}{m_1^2 - m_2} + \frac{2m_2}{m_3 - m_1 m_2} \right)$$

These expressions match mean, variance and skewness of model and data ($N \gg 1$).

Moment-based inference

Approximation methods

Need for further approximation methods when moments cannot be obtained in closed-form.

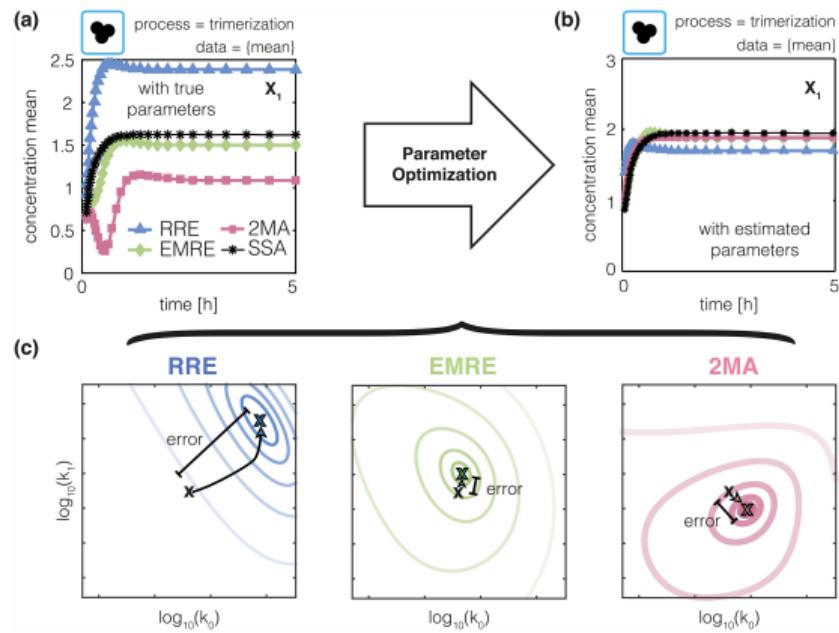
Linear noise approximation: Mean given by deterministic rate equations:

$$\frac{d}{dt}\mu = \sum_r \nu_r f_r(\mu)$$

Covariance matrix:

$$\frac{d}{dt}\Sigma = J\Sigma + \Sigma J^T + D$$

where J is the Jacobian matrix, $D = \nu \text{diag}(f) \nu^T$ is noise matrix.

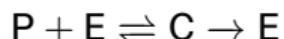


Approximations can fit data but have different likelihood profiles

Moment-based inference

Example: enzyme degradation process

Model:

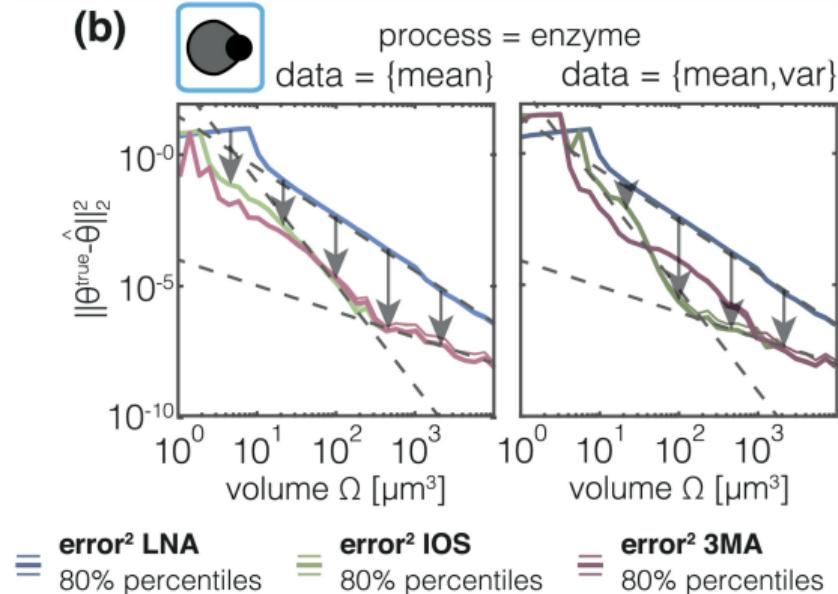


Solve using linear noise approximation or higher-order approximations.

Accuracy increases with the system size Ω that determines the total number of molecules.

Fröhlich, Thomas et al. PLoS Computational Biology (2016) 12, e1005030.

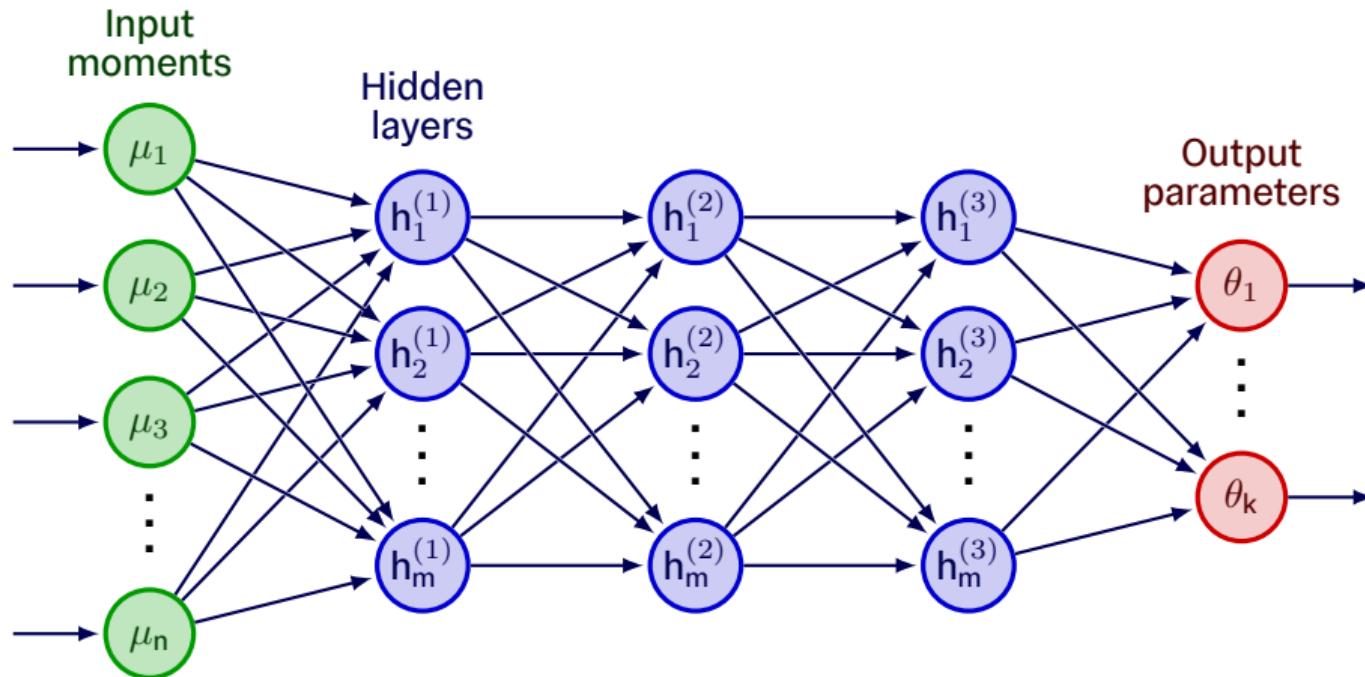
Unsurprisingly, more accurate approximations yield more accurate parameter estimates.



Machine learning-assisted inference

Machine learning-assisted inference

Neural network method



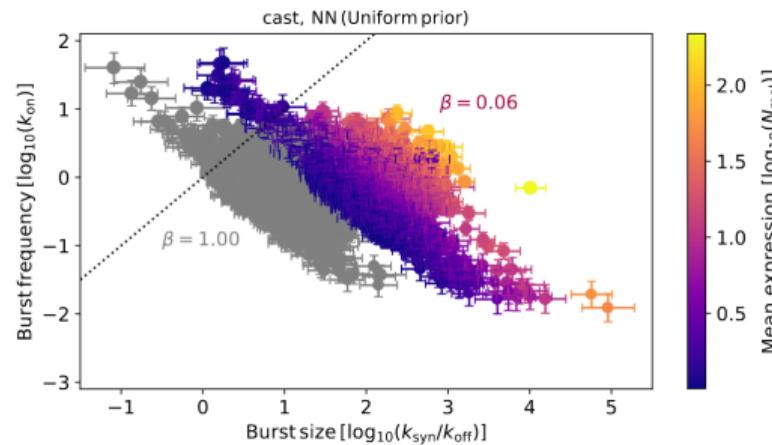
Machine learning-assisted inference

Neural network method

Regression-based algorithm

1. Simulate from the proposal (prior) $\theta' \sim \pi(\theta)$
2. Use θ' and simulate your model $x'_i \sim p(\cdot | \theta')$ for $i = 1, \dots, N$.
3. Compute a number of summary statistics from your data $S_1(x'), \dots, S_m(x')$
4. Pair the prior samples with the summary statistics to train a neural network.

Justification to follow... now onto the coding example



Bayesian inference

Bayesian inference

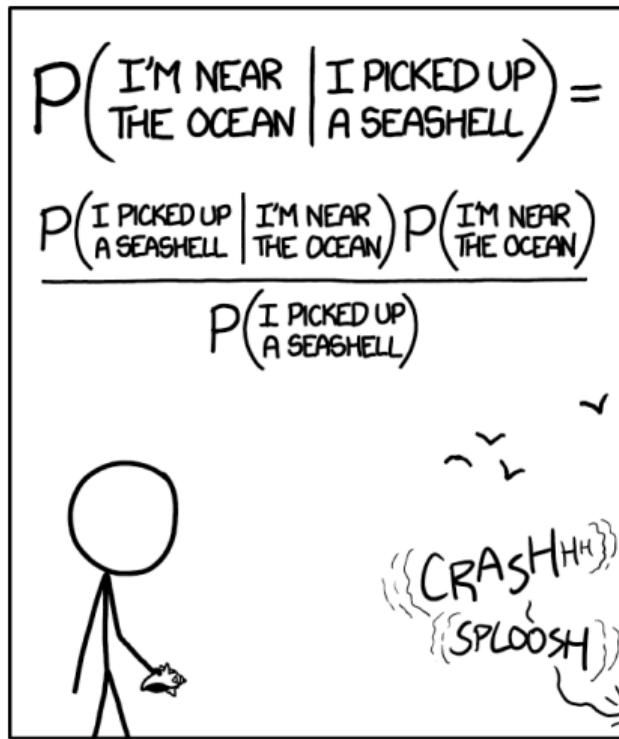
Bayes' theorem

Views both the data D and the parameters as random variables with joint distribution:

$$P(D, \theta) = P(D|\theta)P(\theta)$$

using conditional probability:

$$P(\theta|D = X) = \underbrace{P(D = X|\theta)}_{=\mathcal{L}(\theta|X) \text{ likelihood}} \underbrace{P(\theta)}_{\pi(\theta) \text{ prior}} / \underbrace{P(D = X)}_{\text{evidence}}$$



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND DON'T HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.

Image credit: XKCD needs a mechanistic model

Bayesian inference

Example: Poisson distribution

$$\mathcal{L}(X|\mu) = \prod e^{-\mu} \frac{\mu^{x_i}}{x_i!} \propto e^{-\mu N} \mu^{\sum_i x_i} = e^{-\mu N} \mu^{N\bar{\mu}}$$

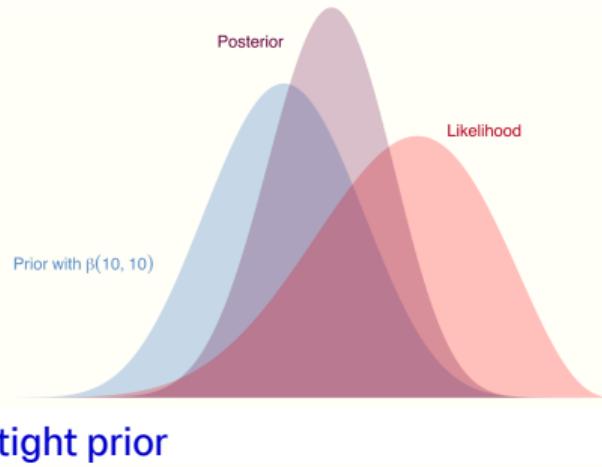
Maximum likelihood occurs when $\mu = \bar{\mu}$

Prior: $\mu \sim \Gamma(\alpha, \beta) : \pi(\mu) \propto \mu^{\alpha-1} e^{-\beta\mu}$

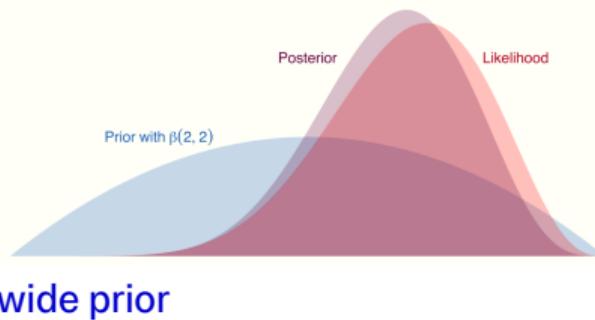
Posterior:

$$\begin{aligned} p(\mu|X) &= \mathcal{L}(\mu|X)\pi(\mu) \\ &\propto e^{-\mu N} \mu^{N\bar{\mu}} \mu^{\alpha-1} e^{-\beta\mu} \\ &= e^{-\mu N - \beta\mu} \mu^{N\bar{\mu} + \alpha - 1} \end{aligned}$$

$$\mu|X \sim \Gamma(\alpha + N\bar{\mu}, N + \beta)$$



tight prior



wide prior

Bayesian inference

MCMC

Markov Chain Monte Carlo (MCMC)

- Idea: construct a stochastic process that has $p(\theta|X)$ as a stationary distribution
- **Advantages:** exact, many excellent libraries exist
- **Drawback:** requires analytically form of the likelihood
- MCMC will not be covered here

MODIFIED BAYES' THEOREM:

$$P(H|X) = P(H) \times \left(1 + P(C) \times \left(\frac{P(x|H)}{P(x)} - 1 \right) \right)$$

H: HYPOTHESIS

X: OBSERVATION

P(H): PRIOR PROBABILITY THAT H IS TRUE

P(x): PRIOR PROBABILITY OF OBSERVING X

P(C): PROBABILITY THAT YOU'RE USING
BAYESIAN STATISTICS CORRECTLY

Image credit: XKCD on model uncertainty

Approximate Bayesian Computation

Simulation-based Bayesian Inference

Exact rejection sampling

Exact algorithm (Rubin 1984^a)

1. Simulate from the prior $\theta' \sim \pi(\theta)$
2. Use θ' and simulate your model $x'_i \sim p(\cdot|\theta)$ for $i = 1, \dots, N$.
3. If $x' = X$ then accept θ'

Each accepted θ' is such that $\theta' \sim \pi(\theta|X)$ exactly.

^aRubin, DB (1984). "Bayesianly Justifiable and Relevant Frequency Calculations for the Applied Statistician". The Annals of Statistics. 12 (4): 1151–1172.
doi:10.1214/aos/1176346785.

Justification

- Joint distribution of accepted samples and parameters is $f(\theta', x') = \underbrace{\pi(\theta)}_{1.} \prod_{i=1}^N \underbrace{p(x'_i|\theta)}_{2.} \underbrace{\mathbb{I}_X(x')}_{3.}$
- Then the marginal distribution is $\int dx' f(\theta', x') = \pi(\theta') \mathcal{L}(\theta'|X) \propto \pi(\cdot|X)$

Problems

- Since x' is discrete, this is manageable if N is small.
- If N is large, the chance of generating the exact same sequence is essentially zero.

Approximate Bayesian Computation

ABC rejection sampling

Exact algorithm (Rubin 1984)

1. Simulate from the prior $\theta' \sim \pi(\theta)$
2. Use θ' and simulate your model $x'_i \sim p(\cdot|\theta')$ for $i = 1, \dots, N$.
3. If $x' = X$ then accept θ'

Each accepted θ' is such that $\theta' \sim \pi(\theta|X)$ when S is a sufficient set of summary statistics and $\epsilon = 0$.

ABC algorithm

1. Simulate from the prior $\theta' \sim \pi(\theta)$
2. Use θ' and simulate your model $x'_i \sim p(\cdot|\theta')$ for $i = 1, \dots, N$.
3. Compute a number of summary statistics from your data $S_1(x'), \dots, S_m(x')$
4. If $d(S(x'), S(X)) < \epsilon$ then accept θ' for a given distance function d .

Justification

- The algorithm is exact when the summary statistics are sufficient and $\epsilon = 0$.
- When $\epsilon > 0$ and the summary statistics are arbitrary, this is a computationally tractable approximation.

Model selection

Frequentist vs. Bayesian inference

Akaike Information Criterion (AIC)

$$AIC = 2k - \max \ln \mathcal{L}(\theta|X)$$

Occam's razor: compromise between model fit ($\max L$) and model complexity (number of parameters k).

ABC model selection algorithm

1. Simulate model m and parameter θ from the prior $m, \theta' \sim \pi$
2. Use θ' and simulate model m according to $x'_i \sim p_m(\cdot|\theta')$ for $i = 1, \dots, N$.
3. Compute a number of summary statistics from your data $S_1(x'), \dots, S_m(x')$
4. If $d(S(x'), S(X)) < \epsilon$ then accept θ' and m for a given distance function d .

The frequency of accepted models determines the model probability under the prior π .

Occam's razor in-built: simpler models are easier to fit.

Comparison of inference methods

Likelihood-based vs. Likelihood-free approaches

Maximum likelihood

Provides point estimate.

Bayesian inference

Provides distribution over parameters based on prior.

Moment-based inference

Uses moments instead of likelihoods.

Approximates the maximum likelihood estimate.

ABC

Simulates moments from prior.

Provides approximate posterior distributions.

Maximum posterior agrees with maximum likelihood if prior is flat.

Neural network simulation-based inference

Simulates moments from prior.

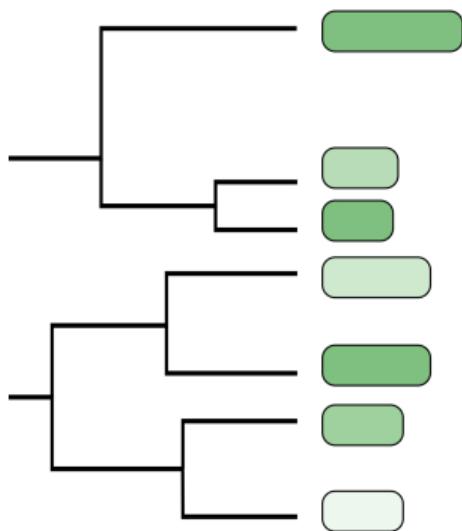
Approximates the maximum of the posterior distribution.

Application to single cell data

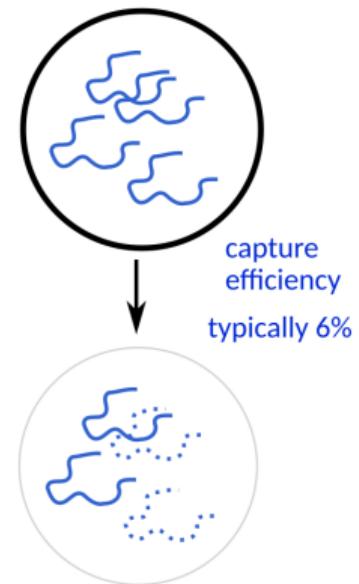
Extrinsic noise models

Stochastic concentration homeostasis: Thomas, Philipp, and Vahid Shahrezaei (2021). Royal Society Interface 18:20210274.

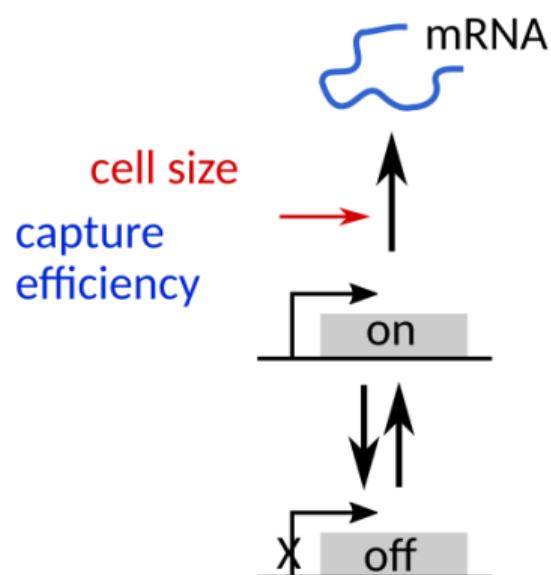
Agent-based modelling



Sequencing (technical) noise



Extrinsic noise models

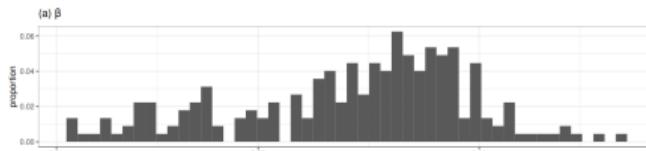


Application to single cell data

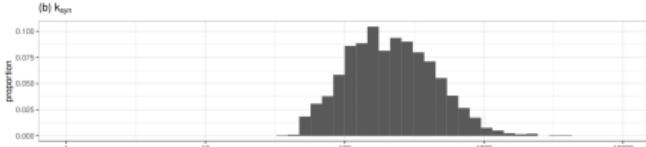
Modelling capture efficiency improves estimates of transcription rate

Mouse fibroblast data from Larsson et al. Nature 565 (2019): 251-254.

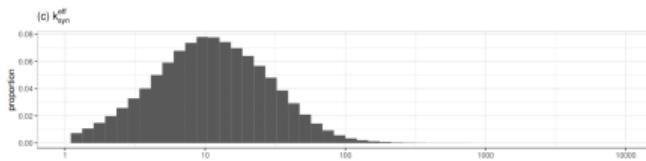
cell-specific capture efficiency



gene-specific transcription rate



effective transcription rate $\beta_{ij} = k_i c_j s_j$



Gene expression matrix of telegraph model:

$$m_{ij}|(p_i, \beta_{i,j}) \sim \text{Poisson}(p_i \beta_{i,j}), \\ p_i \sim \text{Beta}(k_{on,i}, k_{off,i})$$

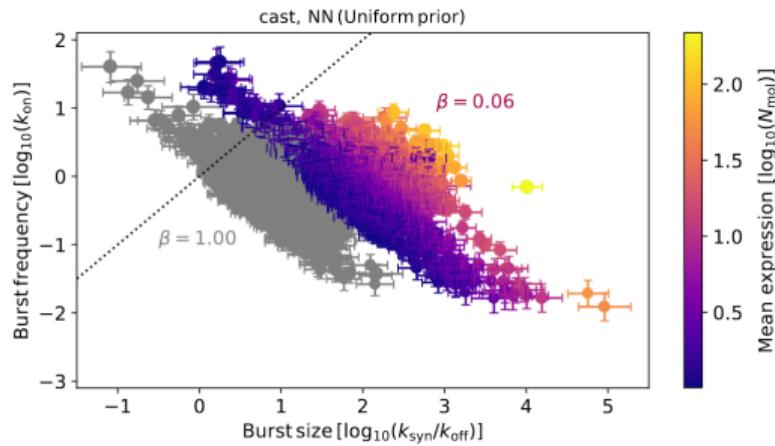
i is gene- and j is cell index.

Decomposing transcription rate into cell- and gene-specific components.

- k's gene-specific rate constants
- $\beta_{i,j}$ effective transcription rate comprising
 - transcription rate
 - cell size
 - cell-specific technical noise

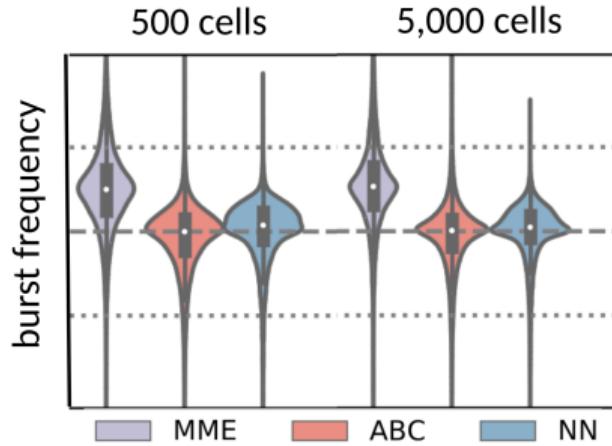
Application to single cell data

Inference of gene expression parameters



Extrinsic noise models more accurately predict burst size and frequency.

Tang, Wenhao, et al. "Modelling capture efficiency of single-cell RNA-sequencing data improves inference of transcriptome-wide burst kinetics." Bioinformatics 39.7 (2023): btad395.

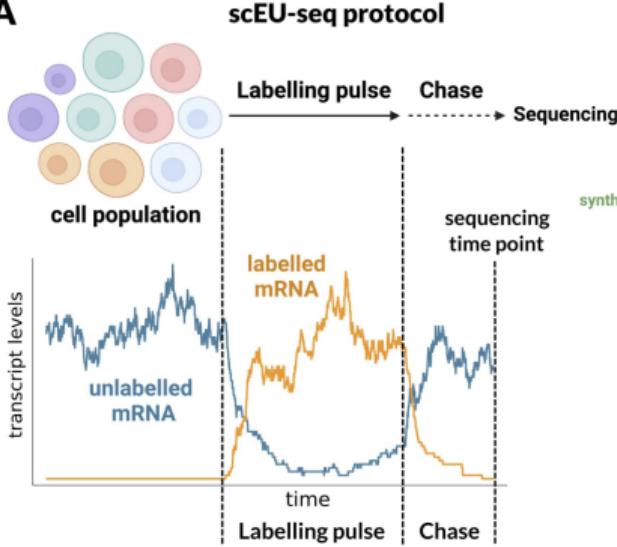


ABC and neural network approaches have similar accuracy.

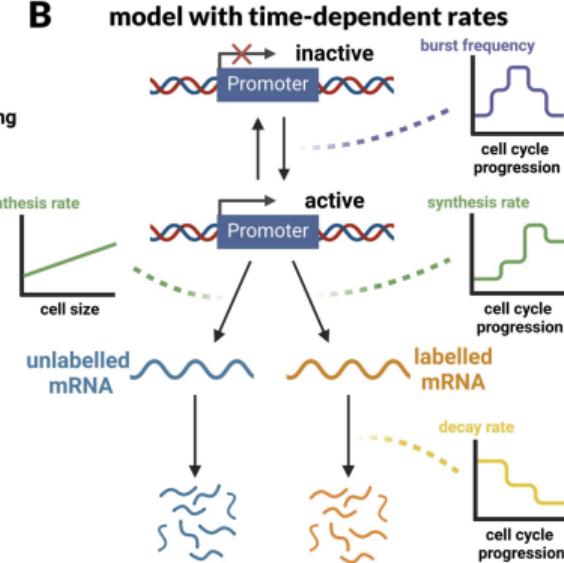
Application to single cell data

Agent-based modelling of time-resolved scRNA

A



B

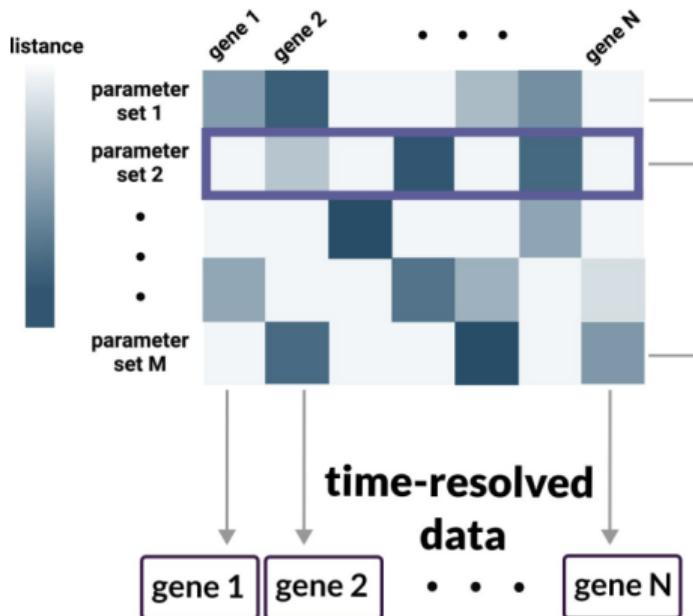


We fitted **every possible model to more than 2,000 genes** under a range of conditions.
Sounds alright?

Application to single cell data

Amortised inference framework of time-resolved scRNA

ABC rejection sampling

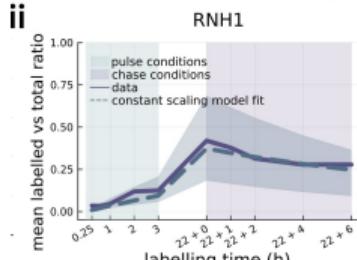


ABC rejection sampling scales easily to genome-scale datasets.

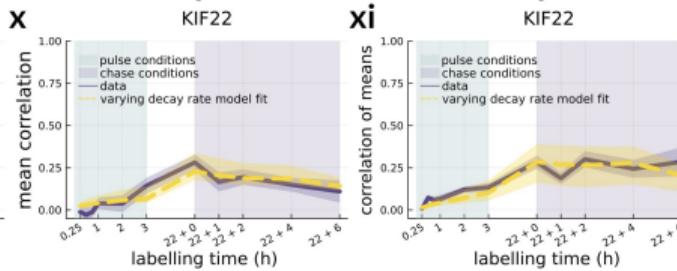
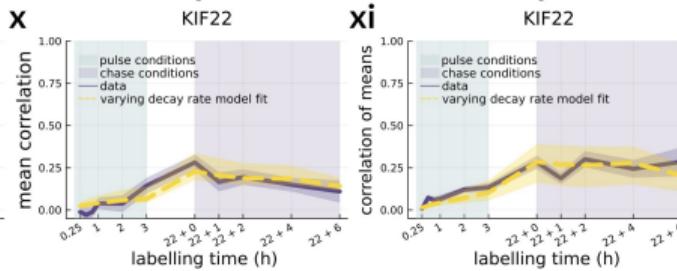
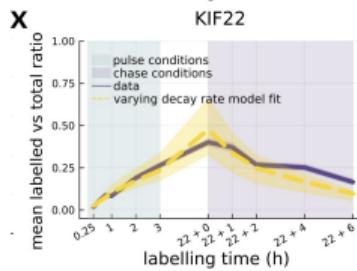
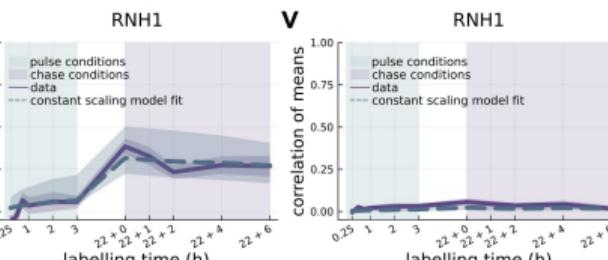
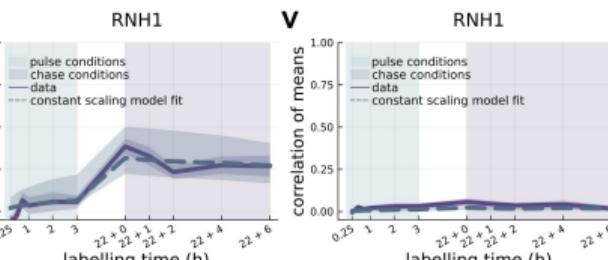
Application to single cell data

ABC fits time-resolved scRNA data

mean labelled vs total expression over labelling time



correlation coefficients of labelled-unlabelled expression over labelling time



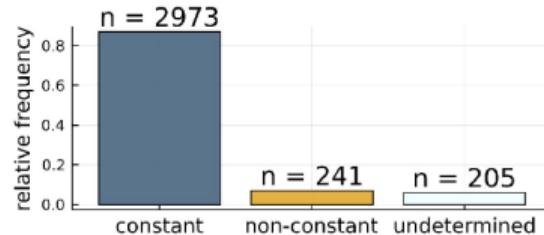
Different models are required to fit individual genes.

Human cell line data (RPE1) from Battich et al. Science 367 (2020): 1151-1156.

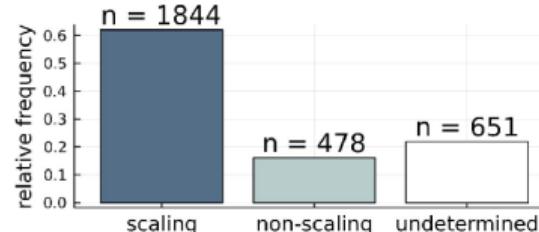
Application to single cell data

Genome-wide regulation underlies distinct regulatory mechanisms

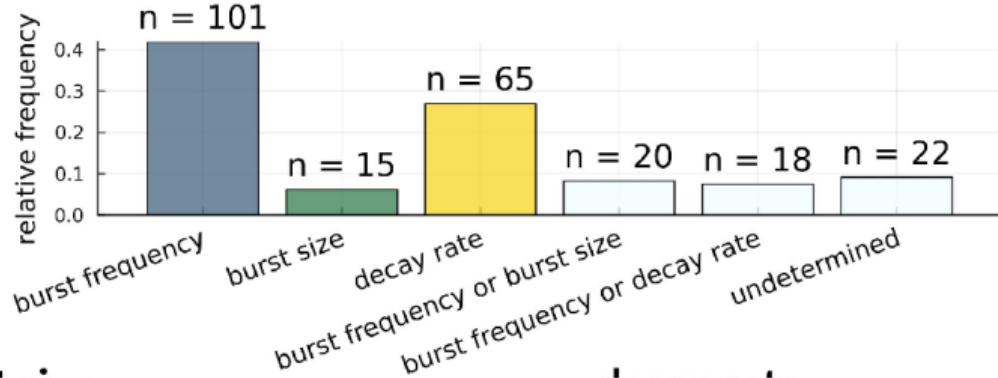
i. overall gene classification



ii. constant gene classification

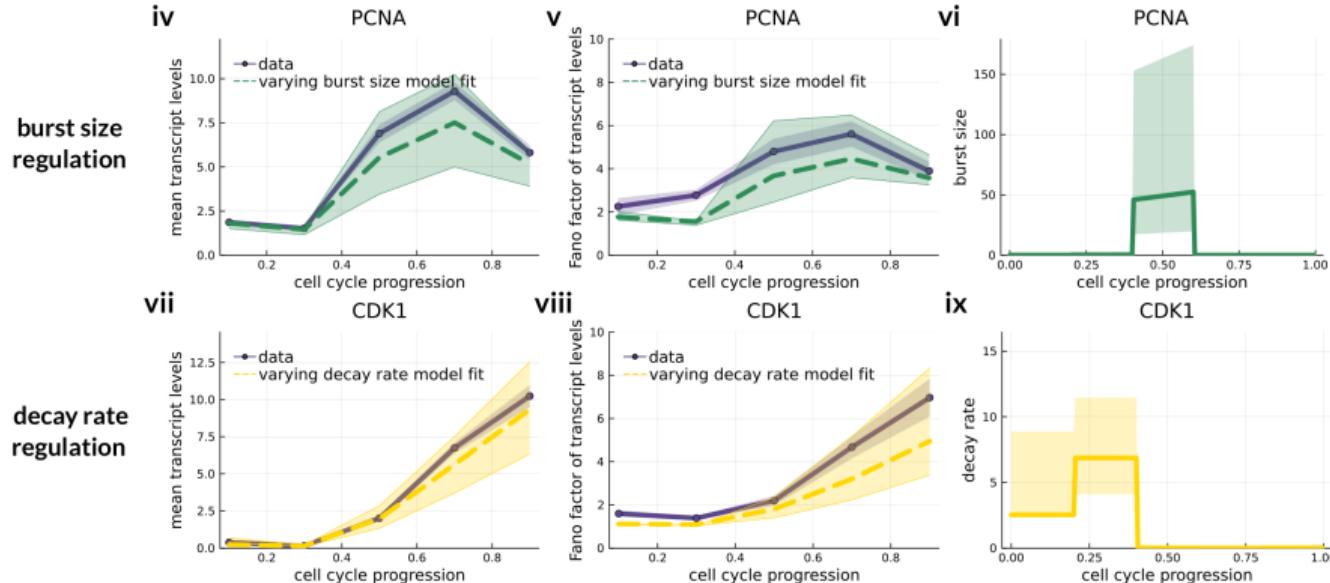


iii. non-constant gene classification



Application to single cell data

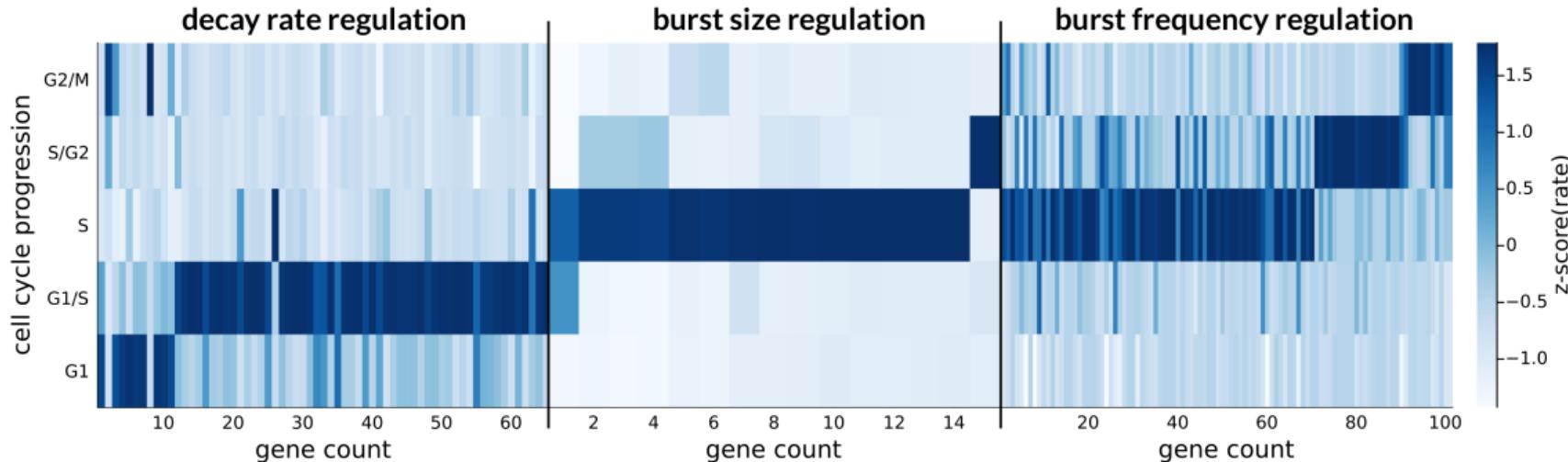
Cell cycle dependence of transcription regulation



Transcription and degradation rates of cell cycle-dependent peak at distinct cell cycle stages.

Application to single cell data

Waves of transcriptional regulation



Global regulation of cell-cycle-dependent transcripts occurs in waves.

Summary

likelihood-based inference

can be computationally challenging
or infeasible depending on model and
data complexity

simulation-based methods

ABC and ML are flexible tools for
mechanistic modelling of data
grounded in Bayesian inference and
provide uncertainty quantification

Outlook

ABC

ML emulators
→ for more efficient
sampling in ABC

Tankhilevich et al.
Bioinformatics 36.10
(2020): 3286-3287.

ML-assisted inference

Bayesian neural
networks
→ posterior distributions

Tang et al.
Bioinformatics 39.7
(2023): btad395.

Moment-based inference

Error guarantees and
bounds on parameters

Li et al. arxiv (2024):
0.48550/arXiv.2406.17434.

IMPERIAL

References:

Volteras, Shahrezaei, Thomas

"Global transcription regulation revealed from dynamical correlations in time-resolved single-cell RNA-sequencing."

Cell Systems (2004, in press)

Tang, et al. "Modelling capture efficiency of single-cell RNA-sequencing data improves inference of transcriptome-wide burst kinetics."

Bioinformatics (2023) 39:btad395.

Fröhlich, Thomas et al. "Inference for stochastic chemical kinetics using moment equations and system size expansion." **PLoS Computational Biology** (2016) 12, e1005030.

Collaborators:

inference of scRNA-seq data:

Dimitris Volteras, Wenhao Tang,
Andreas Joergensen, Vahid
Shahrezaei

inference of tree-structured data:

Fern Hughes, Alasdair Daniels, Alexis
Barr

optimisation-based inference:

Zekai Li, Mauricio Barahona

Thank you