



Accelerating mechanistic inference via neural density estimation

Yannik Schaelte

AstraZeneca, Computational Pathology, 2023-05-16

Journal Club



yannik-schaelte



@yannik_schaelte

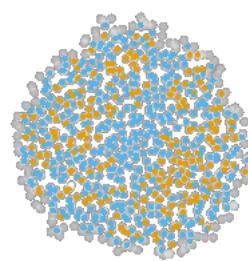
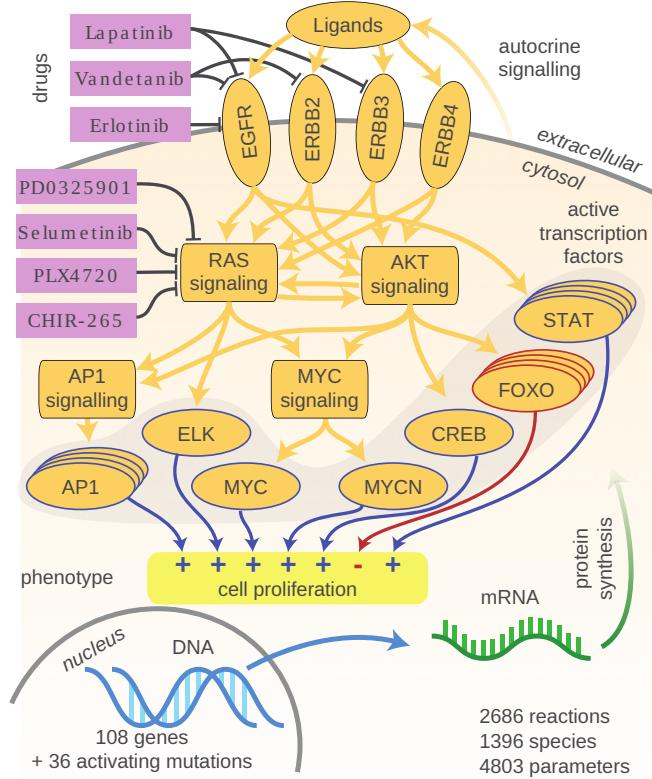


yannikschaelte

OVERVIEW

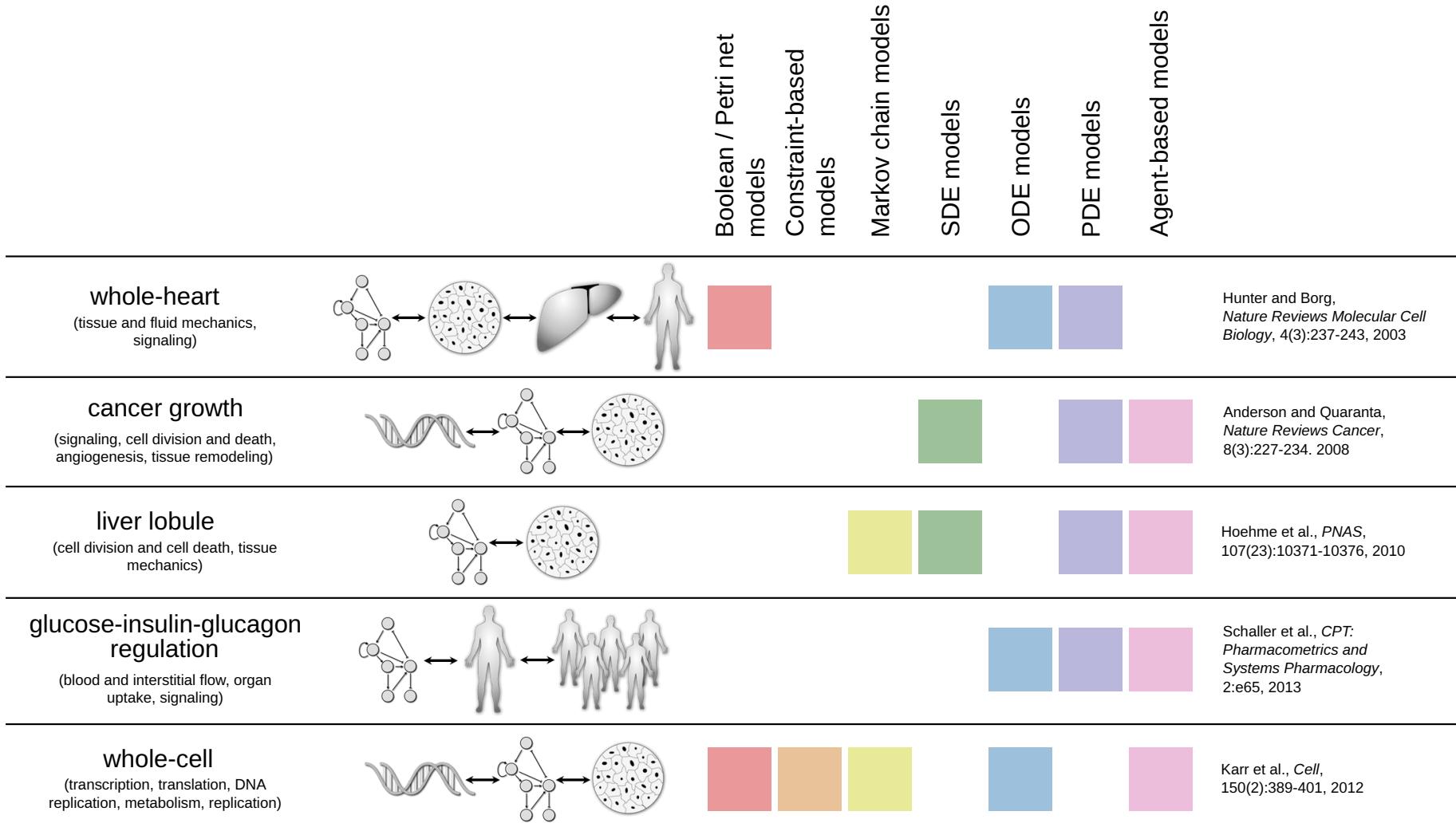
- Quantitative systems biology
- Amortized inference
- Missing data
- Mixed-effects modeling

BIOLOGICAL PROCESSES ARE COMPLICATED



Fröhlich et al., Cell Systems, 2018, and Jagiella et al., Cell Systems, 2017

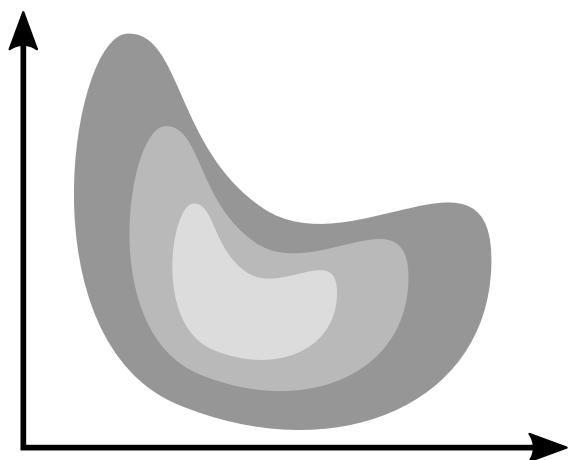
MODEL TYPES



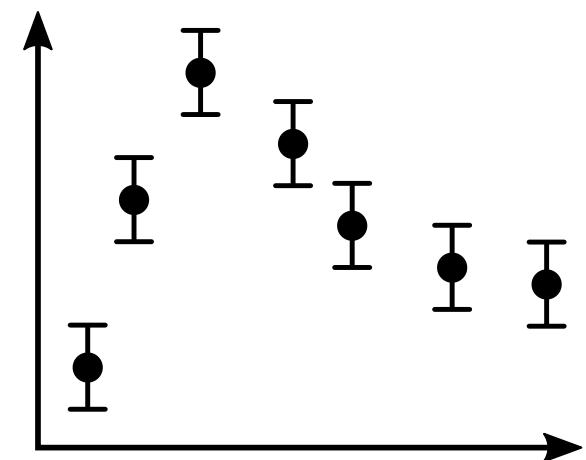
based on Hasenauer et al., Coupl. Sys. 2015

THE INVERSE PROBLEM

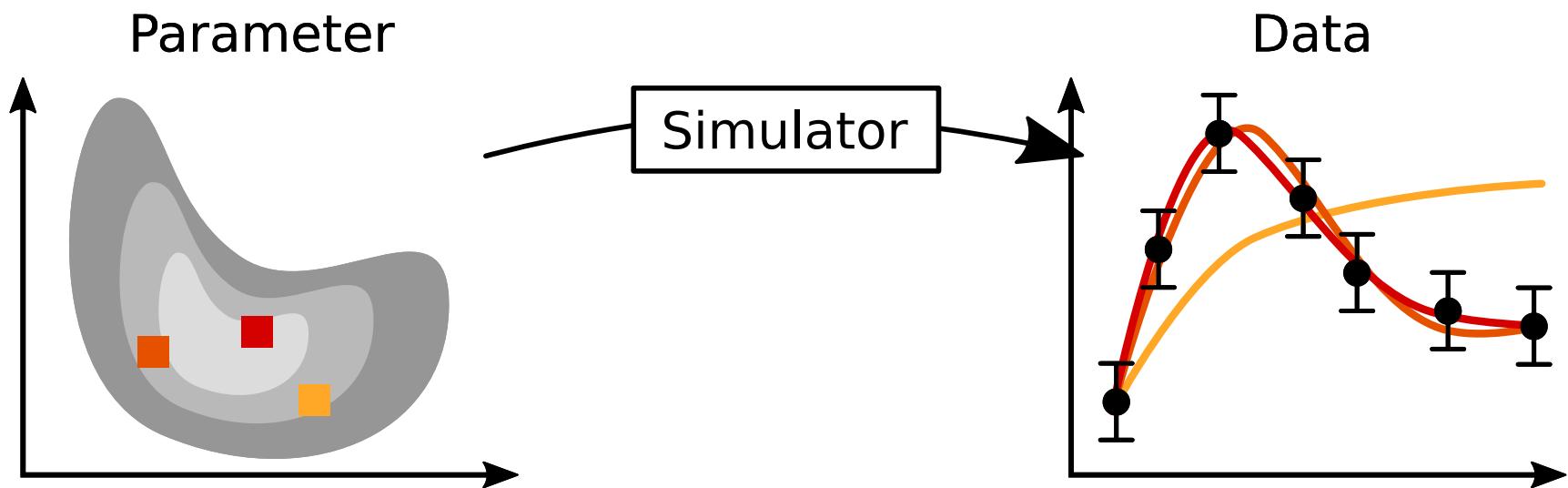
Parameter



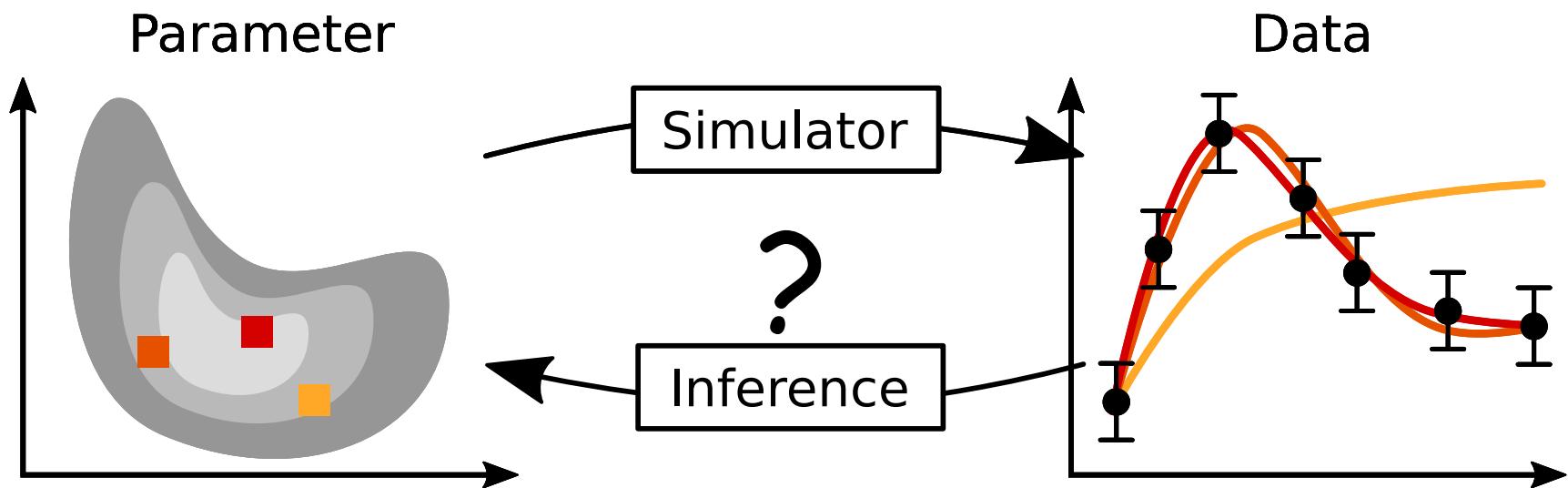
Data



THE INVERSE PROBLEM

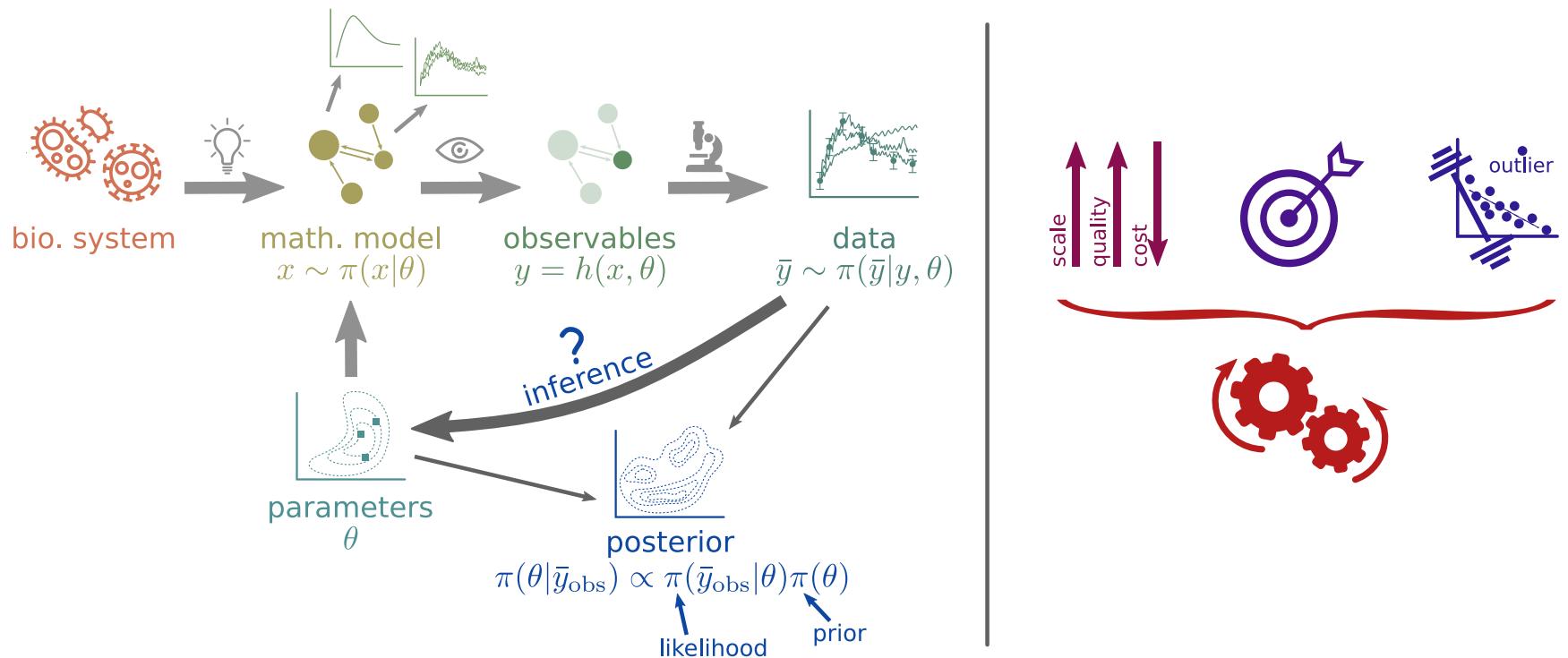


THE INVERSE PROBLEM



Efficient accurate and robust statistical inference for
deterministic and stochastic models of biochemical systems

Efficient accurate and robust statistical inference for deterministic and stochastic models of biochemical systems



SOFTWARE



ODE simulation and sensitivity
analysis



analysis of multi-scale multi-cellular
models



standard to specify estimation
problems



likelihood-free inference
for stochastic models



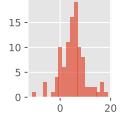
parameter inference: optimization,
sampling, profiles, ...

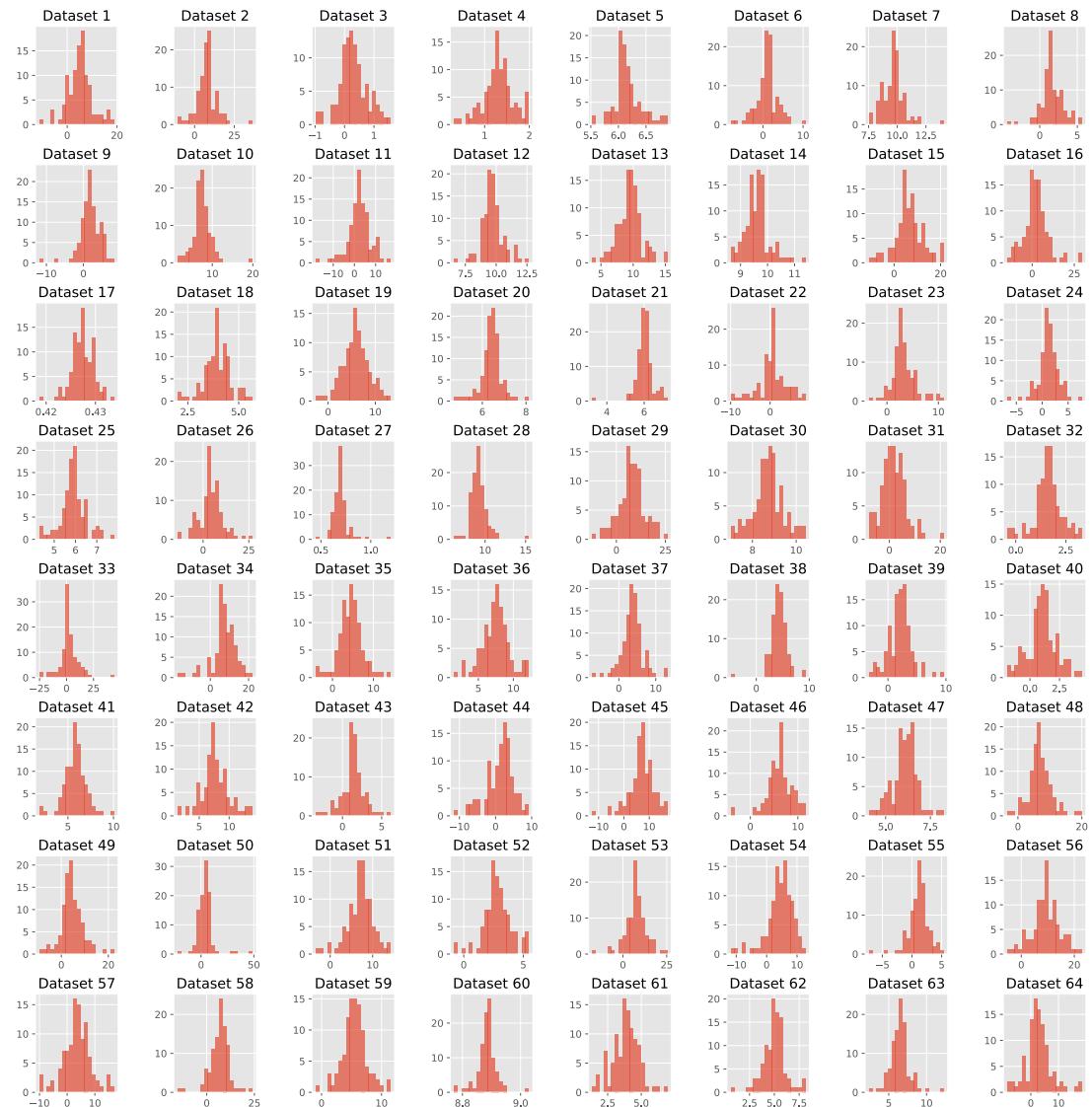


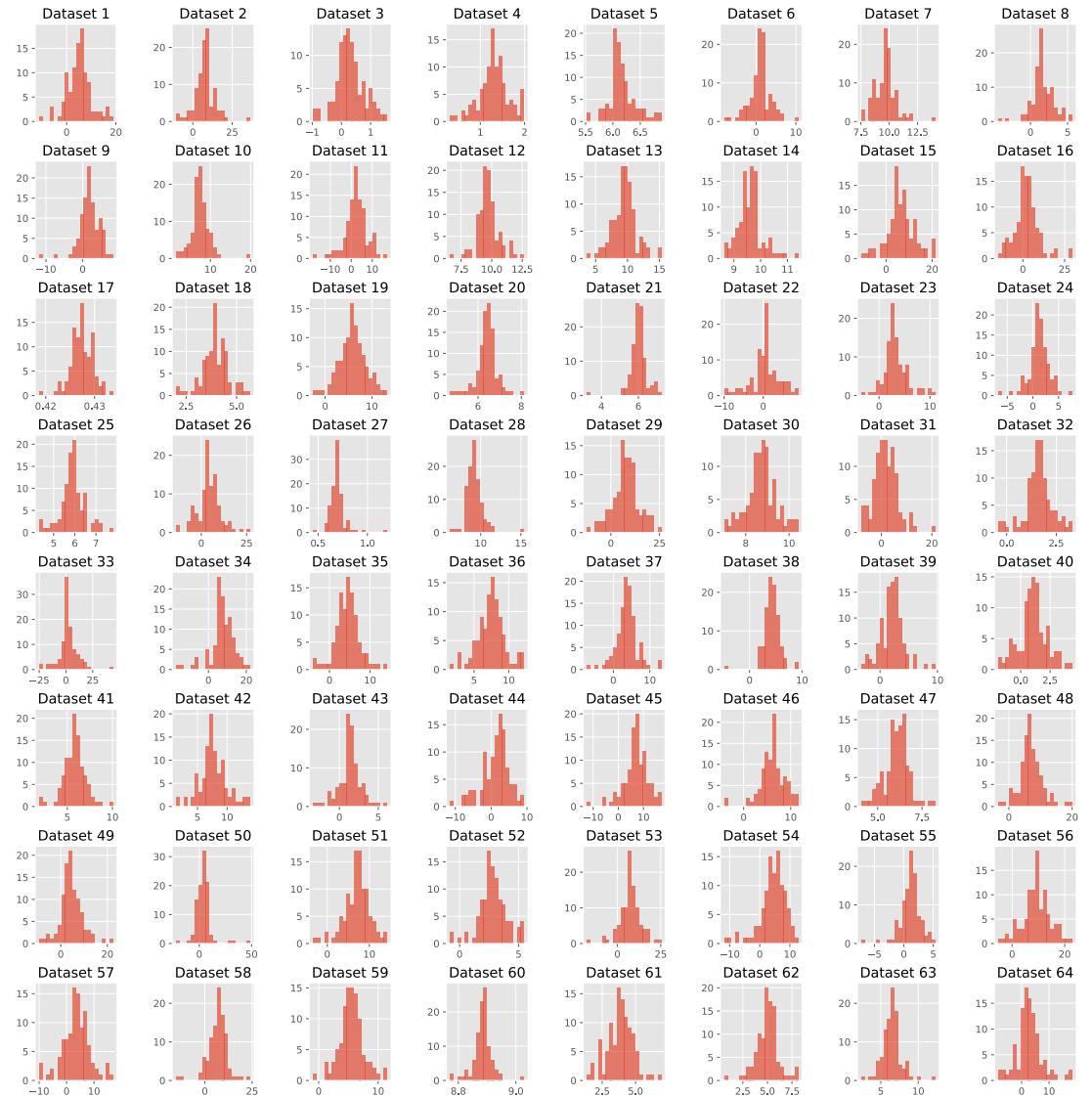
definition of ODE problems in YAML
and conversion to SBML

AMORTIZED INFERENCE

Dataset 1







FITTING A MODEL TO MANY DATASETS?

IN BRIEF

IN BRIEF

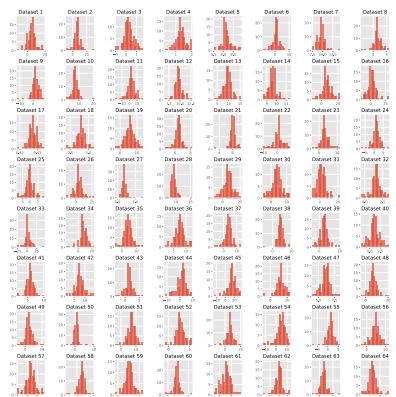
- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate

IN BRIEF

- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate
- ✗ What if we want to fit the same model to multiple datasets?

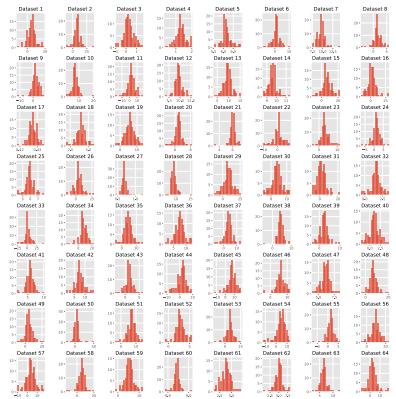
IN BRIEF

- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate
- ✗ What if we want to fit the same model to multiple datasets?



IN BRIEF

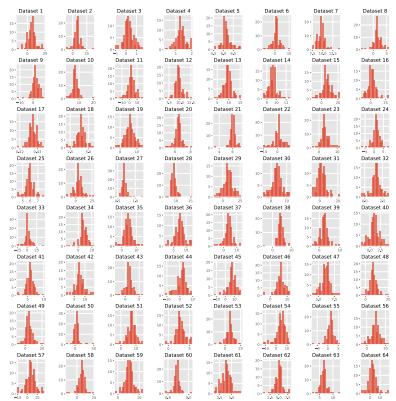
- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate
- ✗ What if we want to fit the same model to multiple datasets?



- ✓ Learn a global estimator for the probabilistic mapping from data to parameters

IN BRIEF

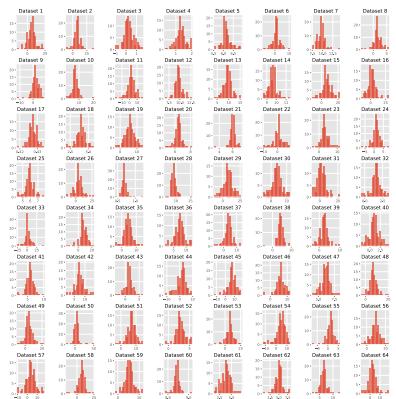
- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate
- ✗ What if we want to fit the same model to multiple datasets?



- ✓ Learn a global estimator for the probabilistic mapping from data to parameters
- ✓ Once trained, amortize inference on arbitrarily many datasets

IN BRIEF

- ✗ Classical (simulation-based) parameter estimation is case-based + slow + approximate
- ✗ What if we want to fit the same model to multiple datasets?



- ✓ Learn a global estimator for the probabilistic mapping from data to parameters
- ✓ Once trained, amortize inference on arbitrarily many datasets
- ✓ Embed data via summary statistics model

GENERATIVE MODELS

generate new data instances, $x \sim \pi(X|Y = y)$

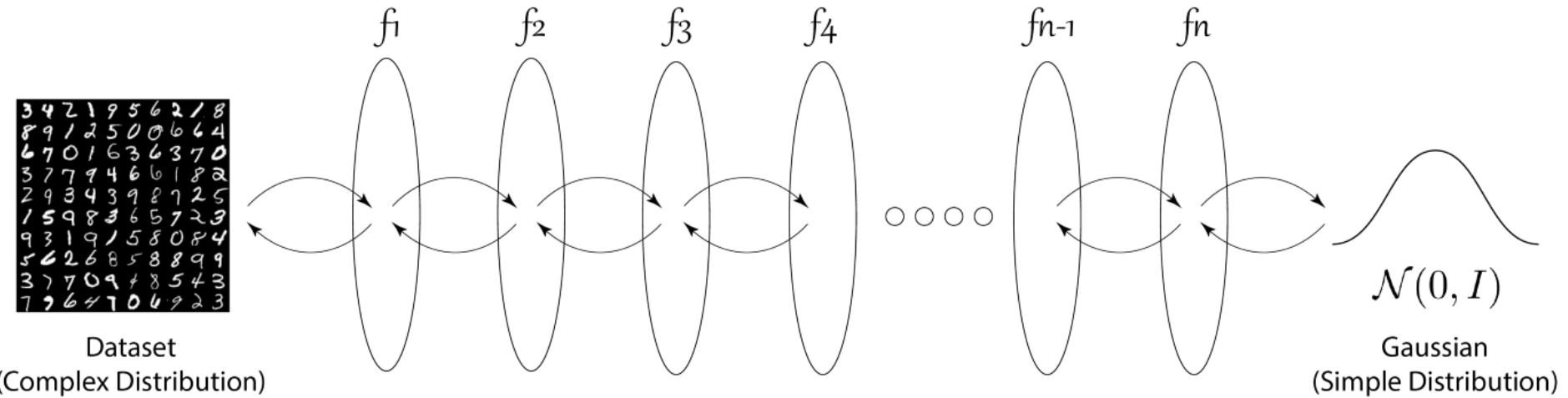


from: Kingma et al., NeurIPS 2019

e.g.: GANs, VAEs, Flows

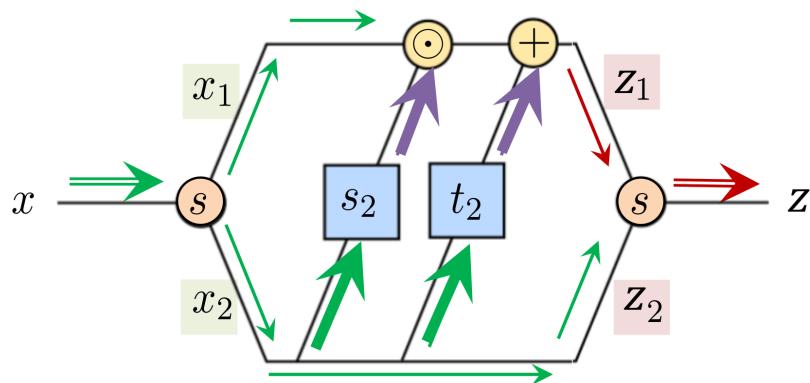
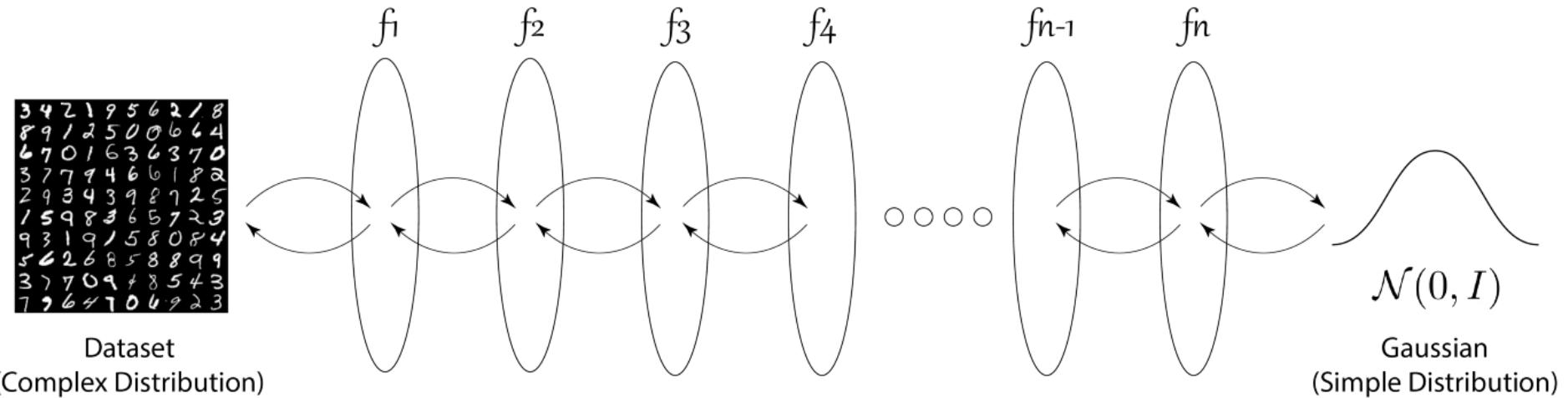
NORMALIZING FLOWS

generative models based on an invertible transformation



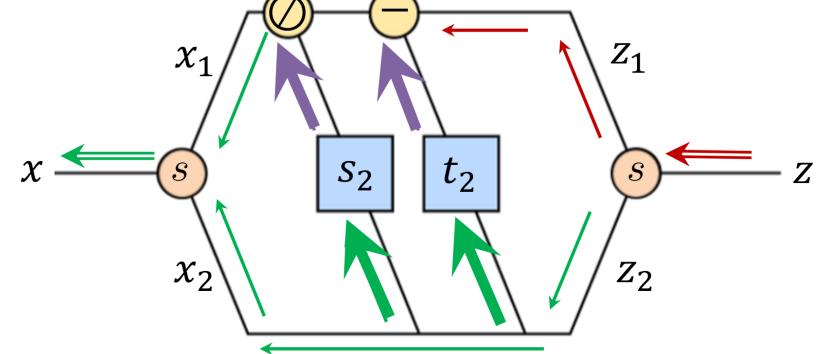
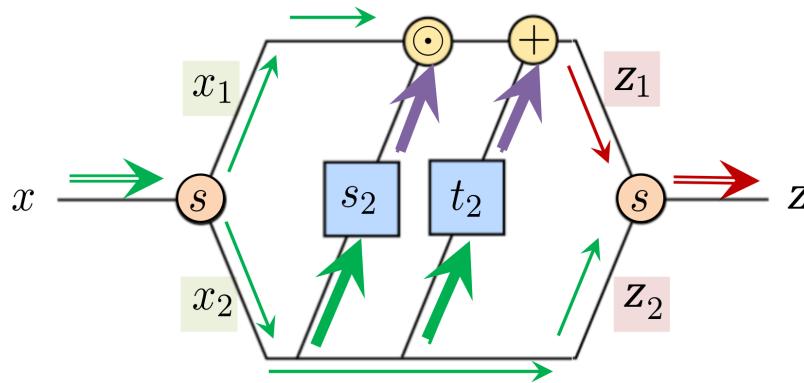
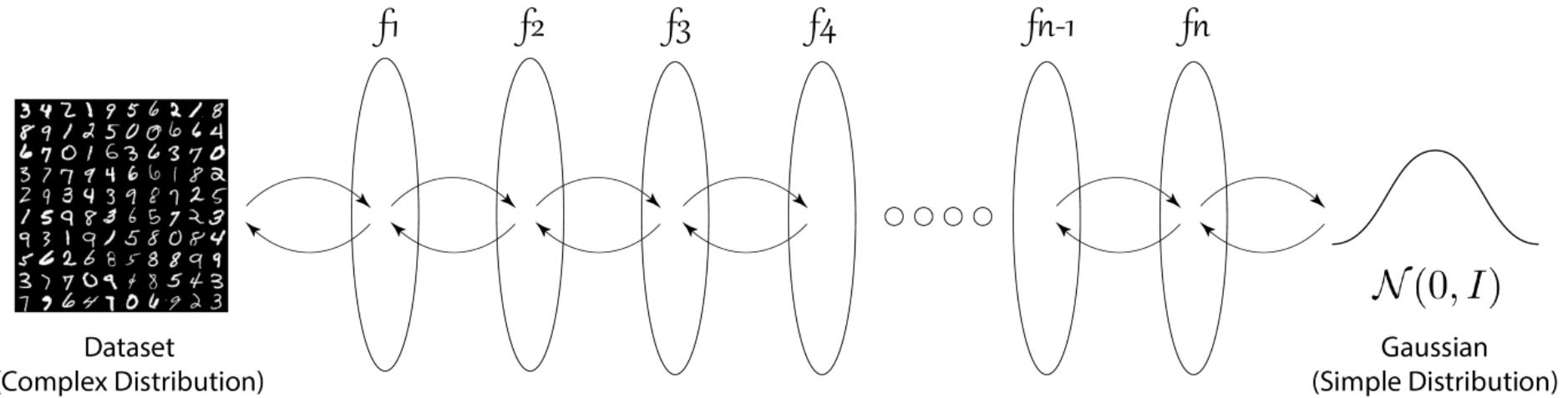
NORMALIZING FLOWS

generative models based on an invertible transformation



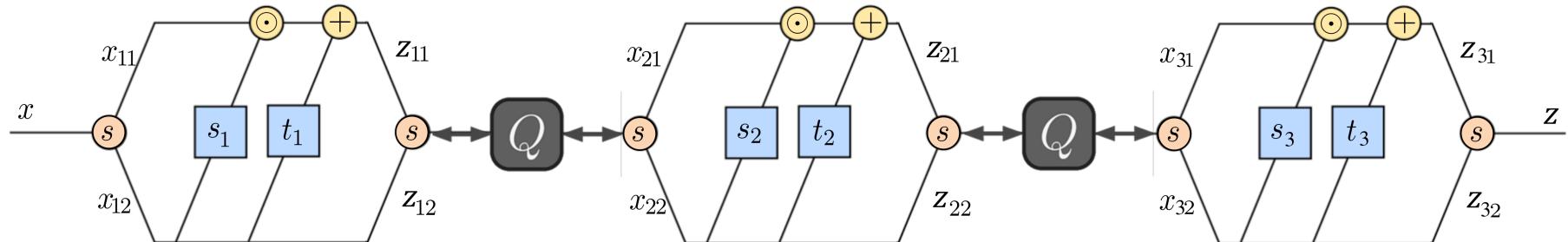
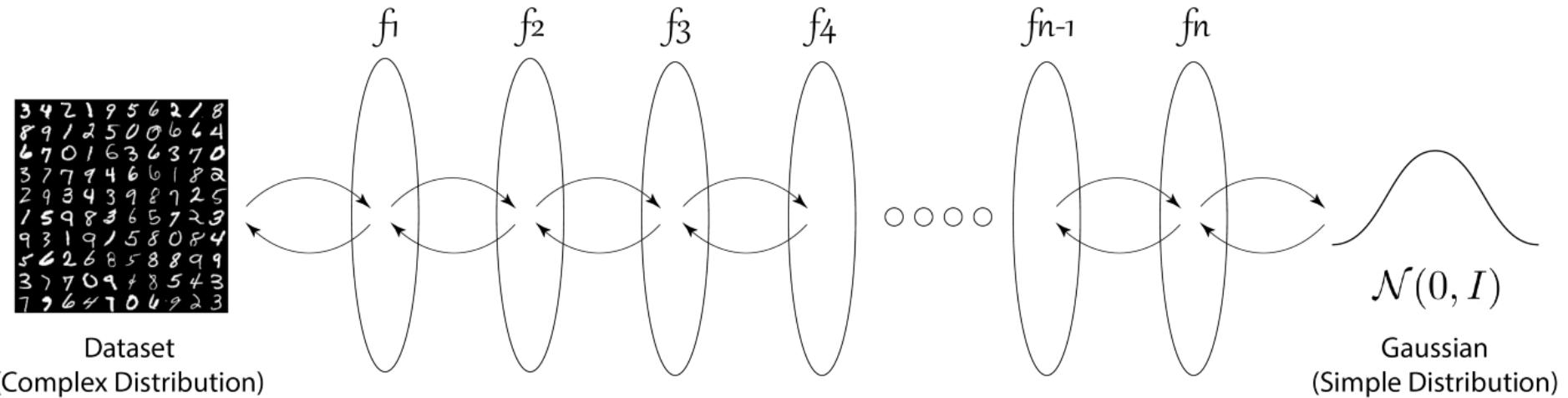
NORMALIZING FLOWS

generative models based on an invertible transformation



NORMALIZING FLOWS

generative models based on an invertible transformation

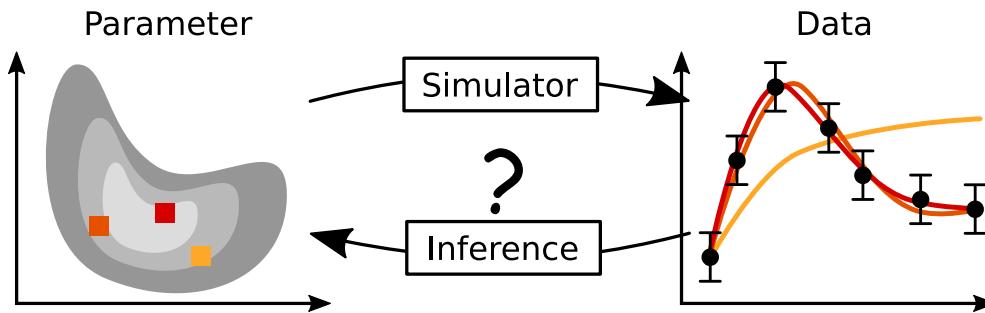


Let $z \sim \mathcal{N}(0, I)$ and $f : z \mapsto x$ bijective. Then via change of variable,s the pdf of $x = f(z)$ is given as

$$p_x(x) = p_z(f^{-1}(x)) \cdot \left| \det\left(\frac{df^{-1}}{dx}(x)\right) \right|.$$

- in training, transform data points to simple distribution
- trained via negative log-likelihood
- afterwards, generate samples via $f^{-1}(z)$ with $z \sim \mathcal{N}(0, I)$

POSTERIOR LEARNING



- parameters θ
 - observations x
 - forward model $x \sim p(x|\theta) \Leftrightarrow x = g(\theta, \xi)$ with $\xi \sim p(\xi)$
 - Bayesian posterior $p(\theta|x) \propto p(x|\theta)p(\theta)$
-
- aim: find a NF $p_\phi(\theta|x) \approx p(\theta|x) \forall \theta, x$

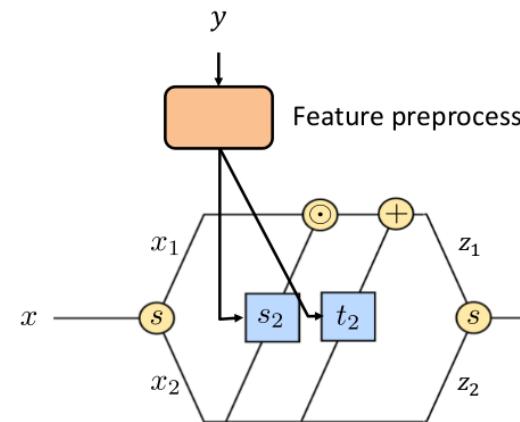
THE METHOD

Parameterize p_ϕ in terms of a cINN given via a bijective

$$f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D, \quad \theta \mapsto z,$$

which implements a normalizing flow between θ and a Gaussian latent variable z ,

$$\theta \sim p_\phi(\theta|x) \Leftrightarrow \phi = f_\phi^{-1}(z; x) \quad \text{with} \quad z \sim \mathcal{N}_D(z|0, I).$$



Seek neural network parameters $\hat{\phi}$ that minimize the KL divergence between true and approximate posterior $\forall x$

THE METHOD

$$\begin{aligned}\hat{\phi} &= \arg \min_{\phi} \mathbb{E}_{p(x)}[\text{KL}(p(\theta|x) || p_{\phi}(\theta|x))] \\ &= \arg \max_{\phi} \iint p(x, \theta) \log p_{\phi}(\theta|x) dx d\theta \\ &= \arg \max_{\phi} \iint p(x, \theta) (\log p(f_{\phi}(\theta; x)) + \log |\det J_{f_{\phi}}|) dx d\theta\end{aligned}$$

Approximate via Monte-Carlo sample:

$$\begin{aligned}\hat{\phi} &= \arg \min_{\phi} \frac{1}{M} \sum_{m=1}^M (-\log p(f_{\phi}(\theta^{(m)}; x^{(m)})) - \log |\det J_{f_{\phi}}^{(m)}|) \\ &= \arg \min_{\phi} \frac{1}{M} \sum_{m=1}^M \left(\frac{|f_{\phi}(\theta^{(m)}; x^{(m)})|_2^2}{2} - \log |\det J_{f_{\phi}}^{(m)}| \right)\end{aligned}$$

LEARN SUMMARY STATISTICS

If data $x_{1:N}$ are high-dimensional: Jointly learn a summary network
 $\tilde{x} = h_\psi(x_{1:N})$, giving the objective

$$\hat{\phi}, \hat{\psi} = \arg \max_{\phi, \psi} \mathbb{E}_{p(x, \theta, N)} [\log p_\phi(\theta | h_\psi(x_{1:N}))]$$

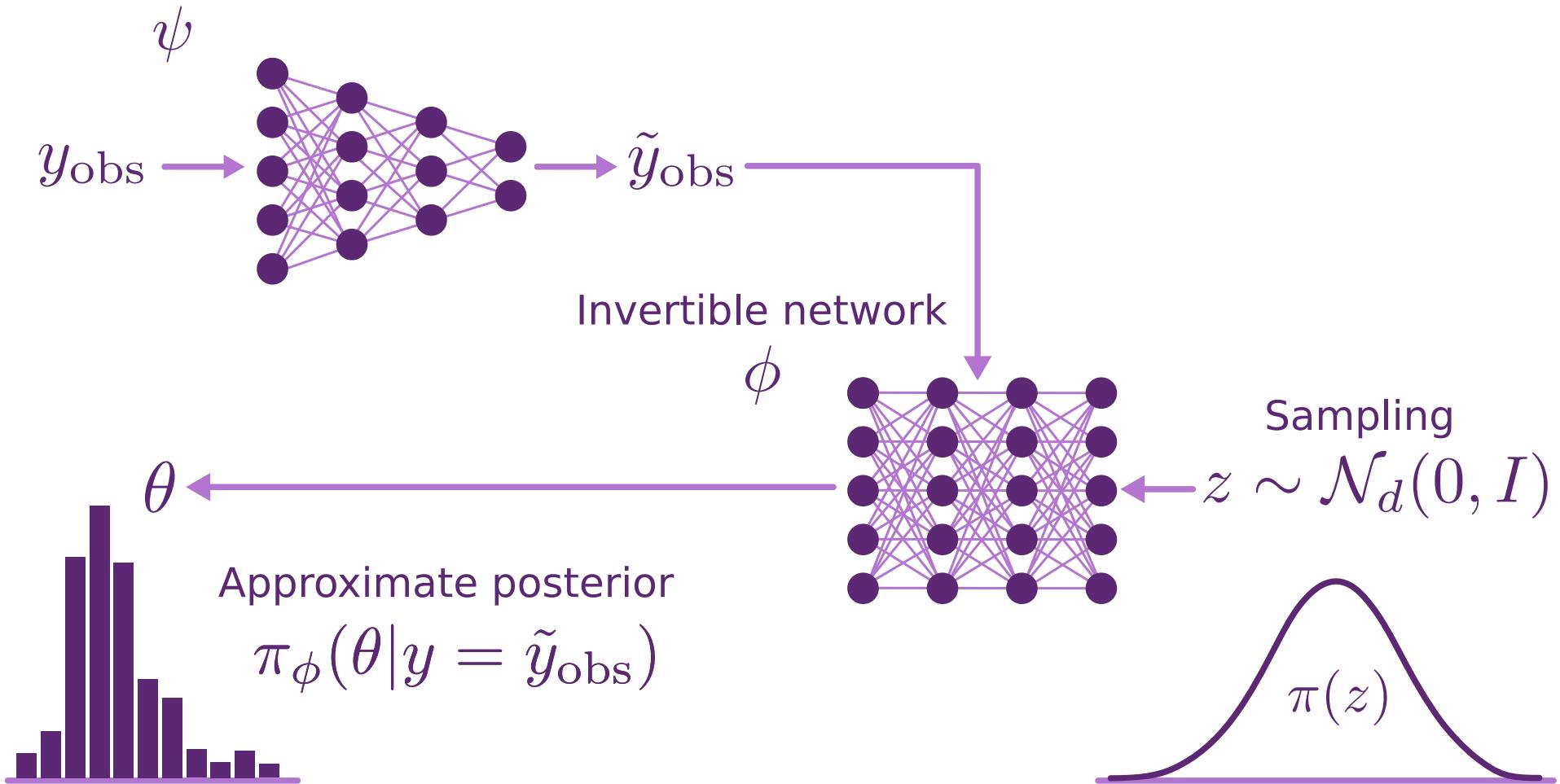
with Monte-Carlo estimate

$$\hat{\phi}, \hat{\psi} = \arg \min_{\phi, \psi} \frac{1}{M} \sum_{m=1}^M \left(\frac{|f_\phi(\theta^{(m)}; h_\psi(x_{1:N}^{(m)})|_2^2}{2} - \log |\det(J_{f_\phi}^{(m)})| \right)$$

AMORTIZED INFERENCE VIA INNS

you have to solve many similar problems? amortize the solution!

Summary network



TRAINING AND INFERENCE

Training phase:

- create plenty of synthetic data $(y_i, \theta_i) \sim \pi(y, \theta)$
- train a cINN in forward mode

Inference phase:

- sample many latent $z_i \sim \pi(z)$
- run cINN backwards, $\theta_i = g(z_i; y_{\text{obs}}) \sim \pi(\theta | y_{\text{obs}})$

✓ fast + accurate amortized simulation-based Bayesian inference

$$\frac{\partial y}{\partial x}$$

mechanistic modeling

$$\frac{\partial y}{\partial x}$$

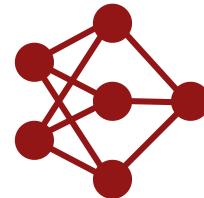
mechanistic modeling

- theory-driven
- interpretation
- testability

$$\frac{\partial y}{\partial x}$$

mechanistic modeling

- theory-driven
- interpretation
- testability

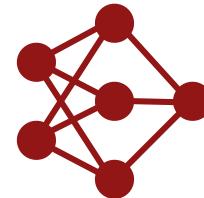


machine learning

$$\frac{\partial y}{\partial x}$$

mechanistic modeling

- theory-driven
- interpretation
- testability



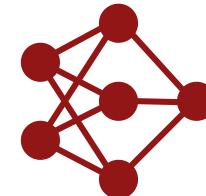
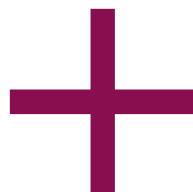
machine learning

- data-driven
- large data
- automation

$$\frac{\partial y}{\partial x}$$

mechanistic modeling

- theory-driven
- interpretation
- testability



machine learning

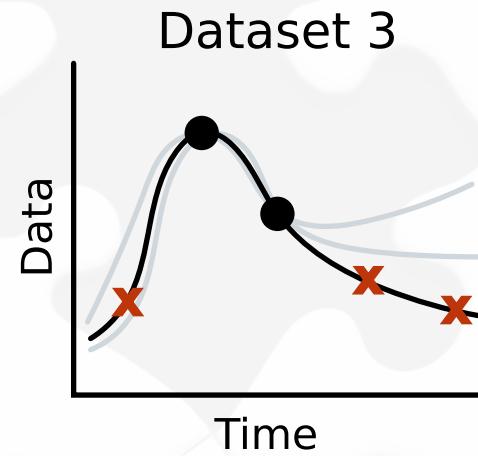
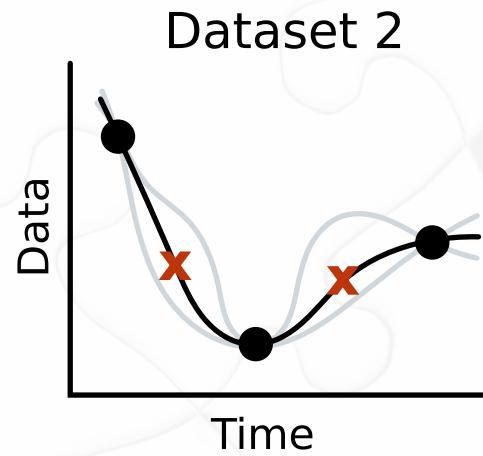
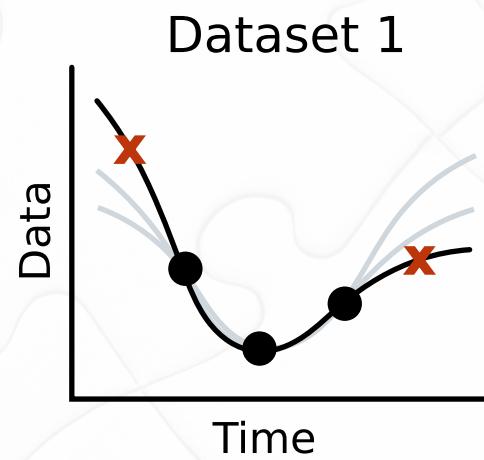
- data-driven
- large data
- automation

sciML

model-based data-efficient ML

MISSING DATA

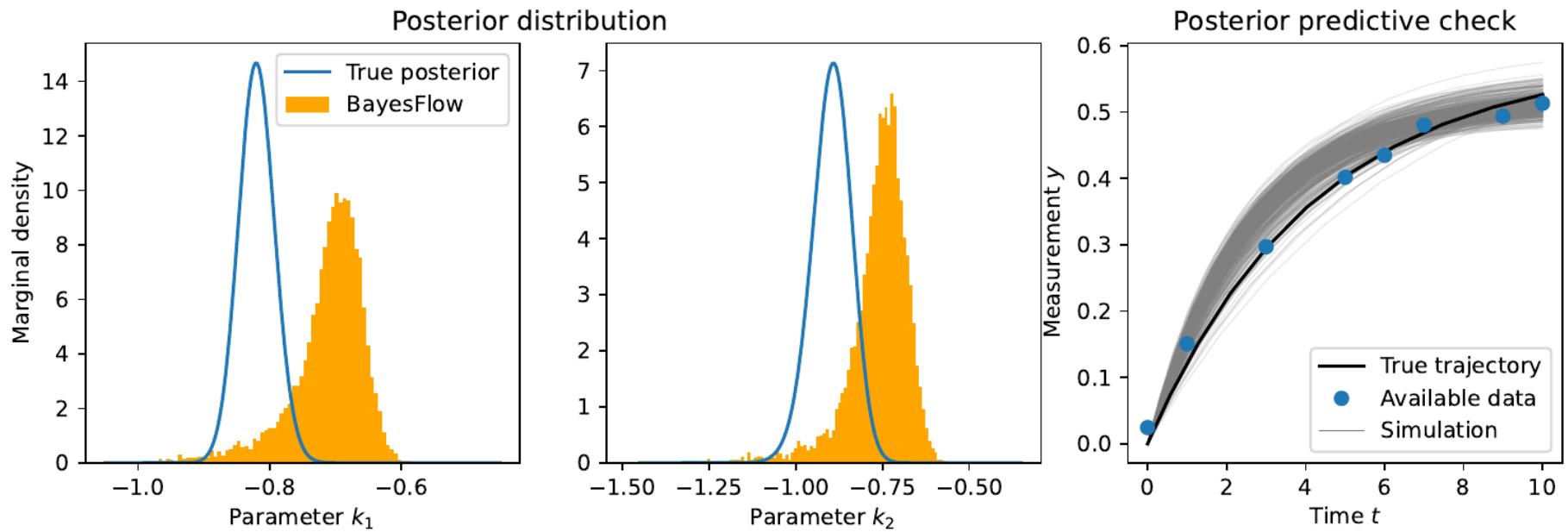
HOW TO HANDLE MISSING DATA IN AMORTIZED INFERENCE?



Available
Missing

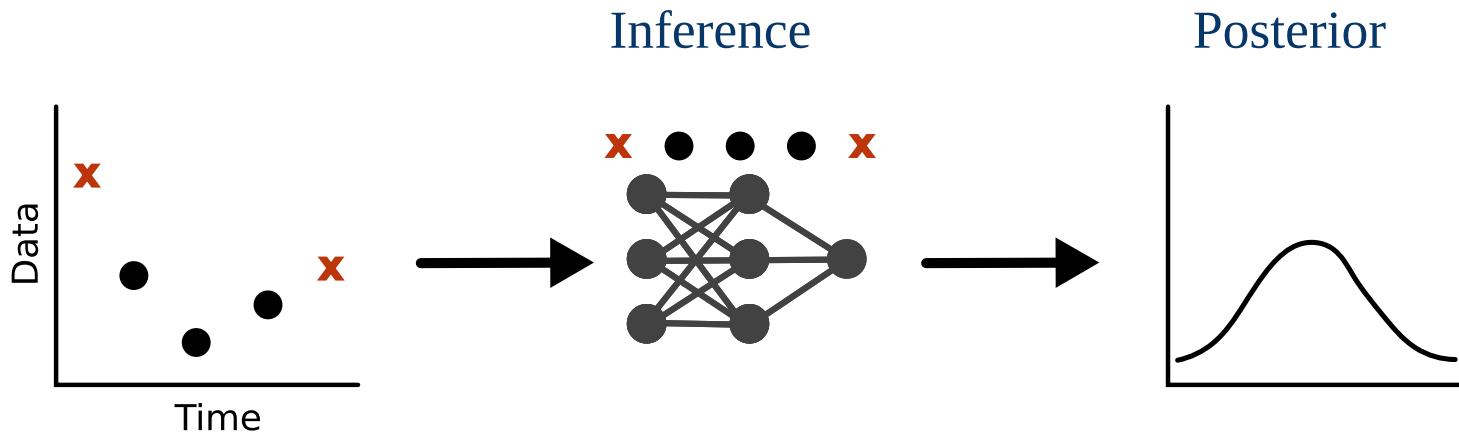
HOW TO HANDLE MISSING DATA IN AMORTIZED INFERENCE?

PROBLEM: INN CANNOT INTERPRET THE DATA

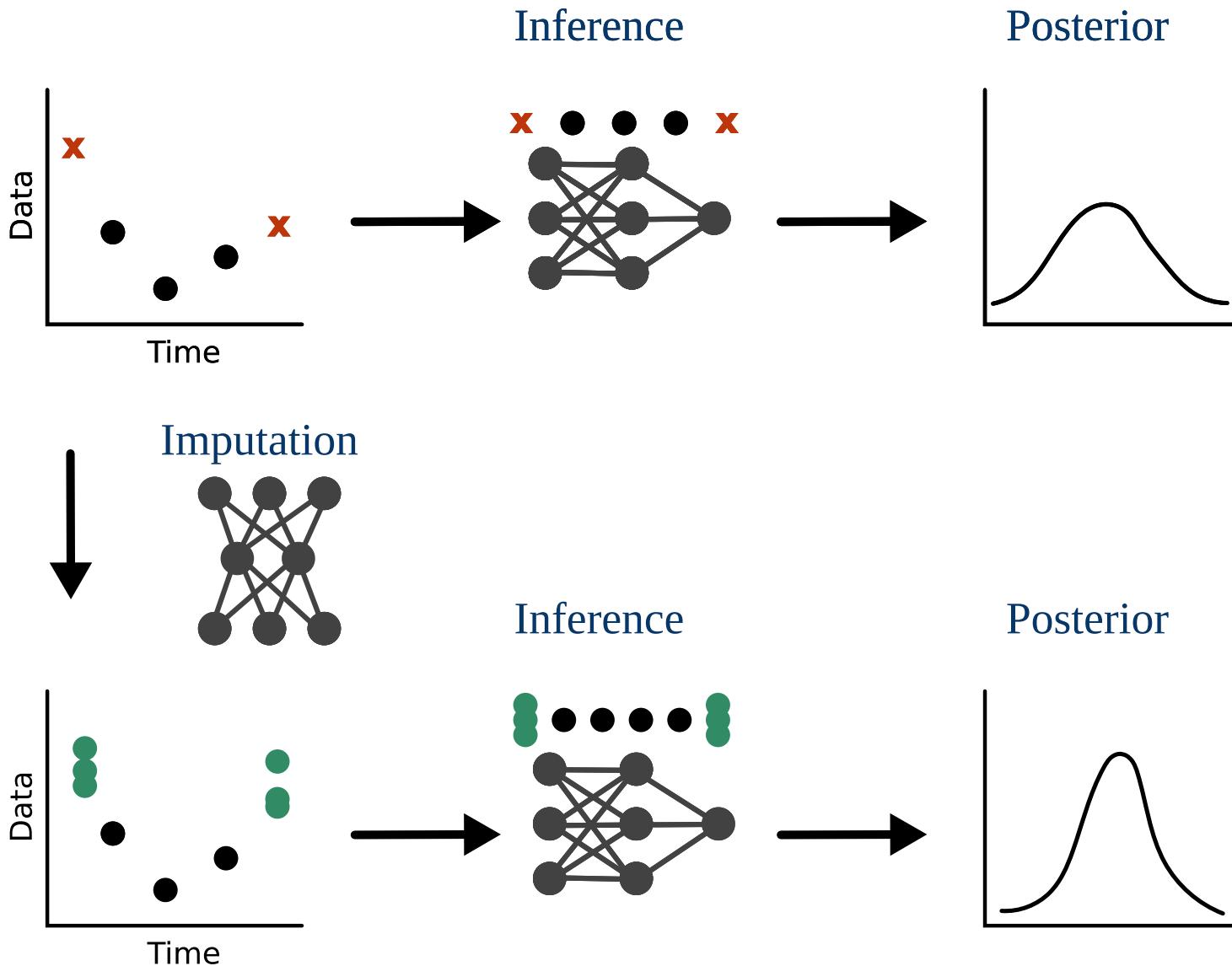


(besides iid + time-series data of heterogeneous length)

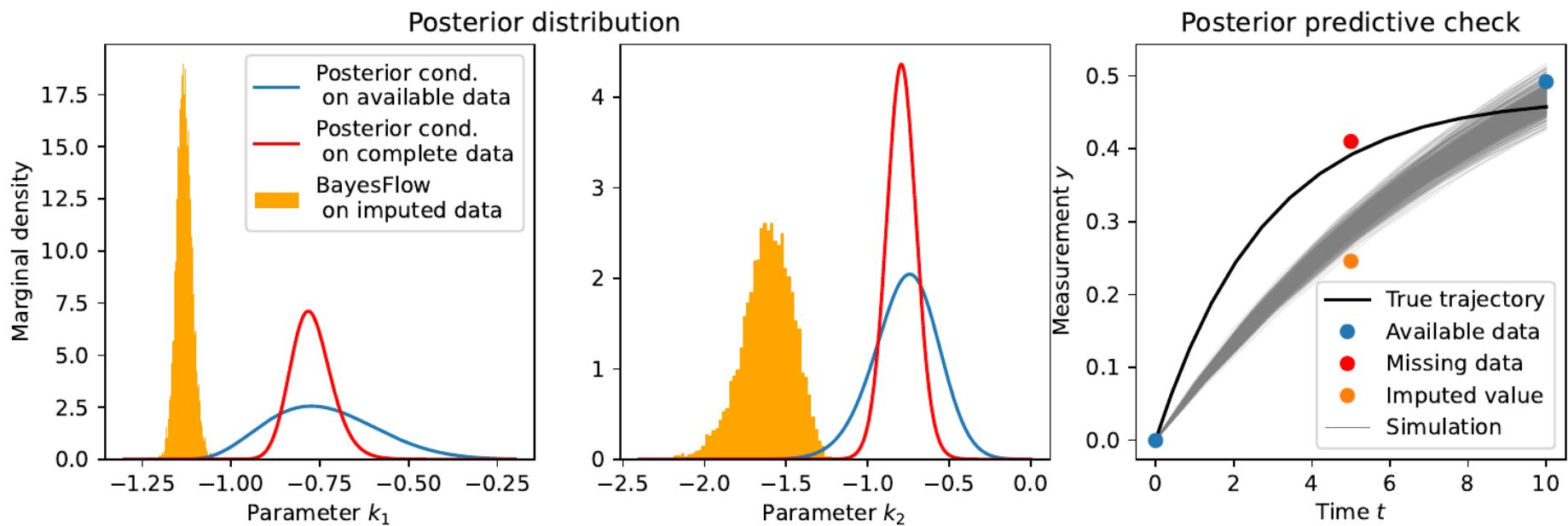
CAN WE JUST IMPUTE MISSING VALUES?



CAN WE JUST IMPUTE MISSING VALUES?



INAPPROPRIATE IMPUTATION CAN LEAD TO BIASED RESULTS



THERE ARE NO FREE DATA

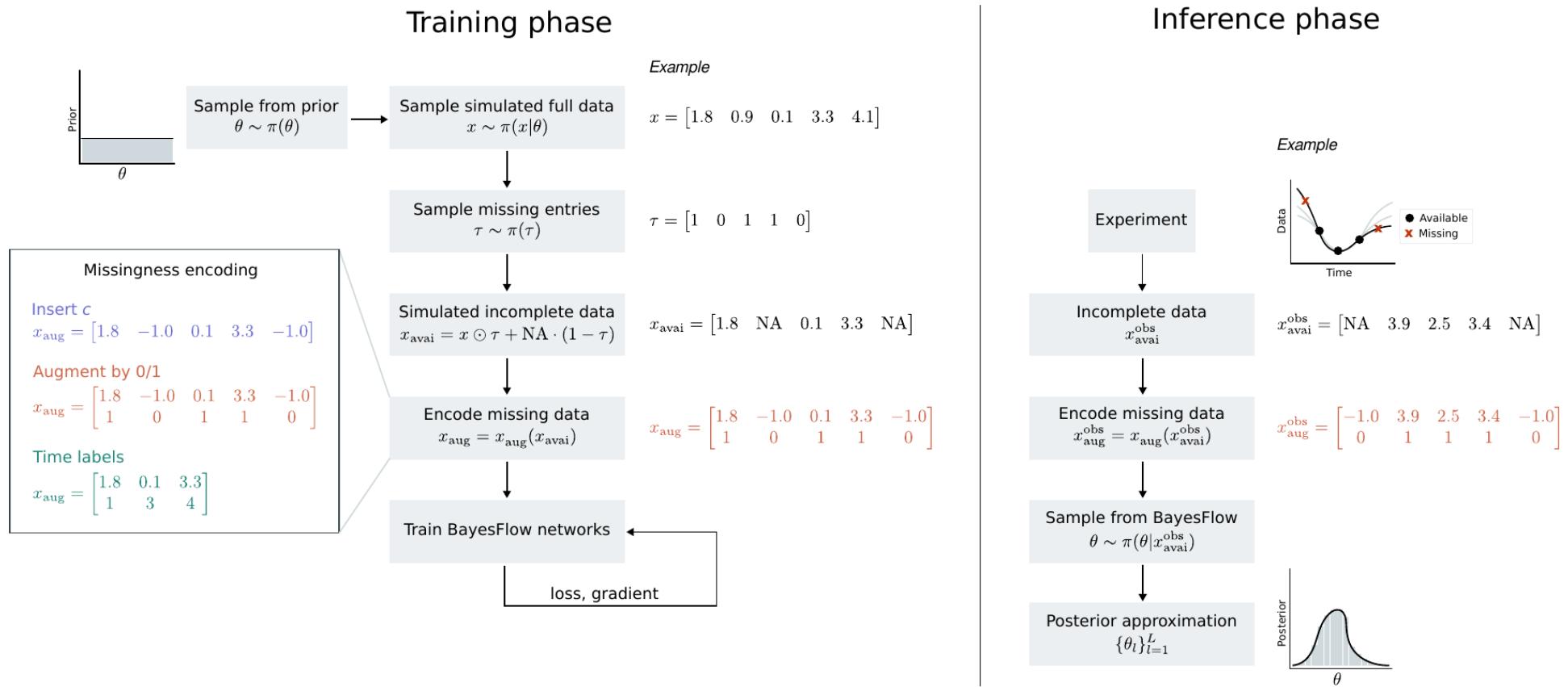
Imputation means that instead of working with available data x , we try to reconstruct the full data \bar{x} , and estimate parameter probabilities $\pi(\theta|\bar{x})$ instead of $\pi(\theta|x)$. However, the true full data are unknown, therefore we need to take uncertainty in \bar{x} into account, considering a full distribution of values $\pi(\bar{x}|x)$.

We must either make up a distribution (introducing a bias), or use a faithful approximation $p(\bar{x}|x) = \pi(\bar{x}|x)$ where $\pi(\bar{x}|x)\pi(x) = \pi(\bar{x}, x)$.

However, if we integrate out over all possible realizations of full data, we obtain
 $\int \pi(\theta|\bar{x})\pi(\bar{x}|x)d\bar{x} = \pi(\theta|x)$ (or similarly $\pi(\theta|x, \tau)$).

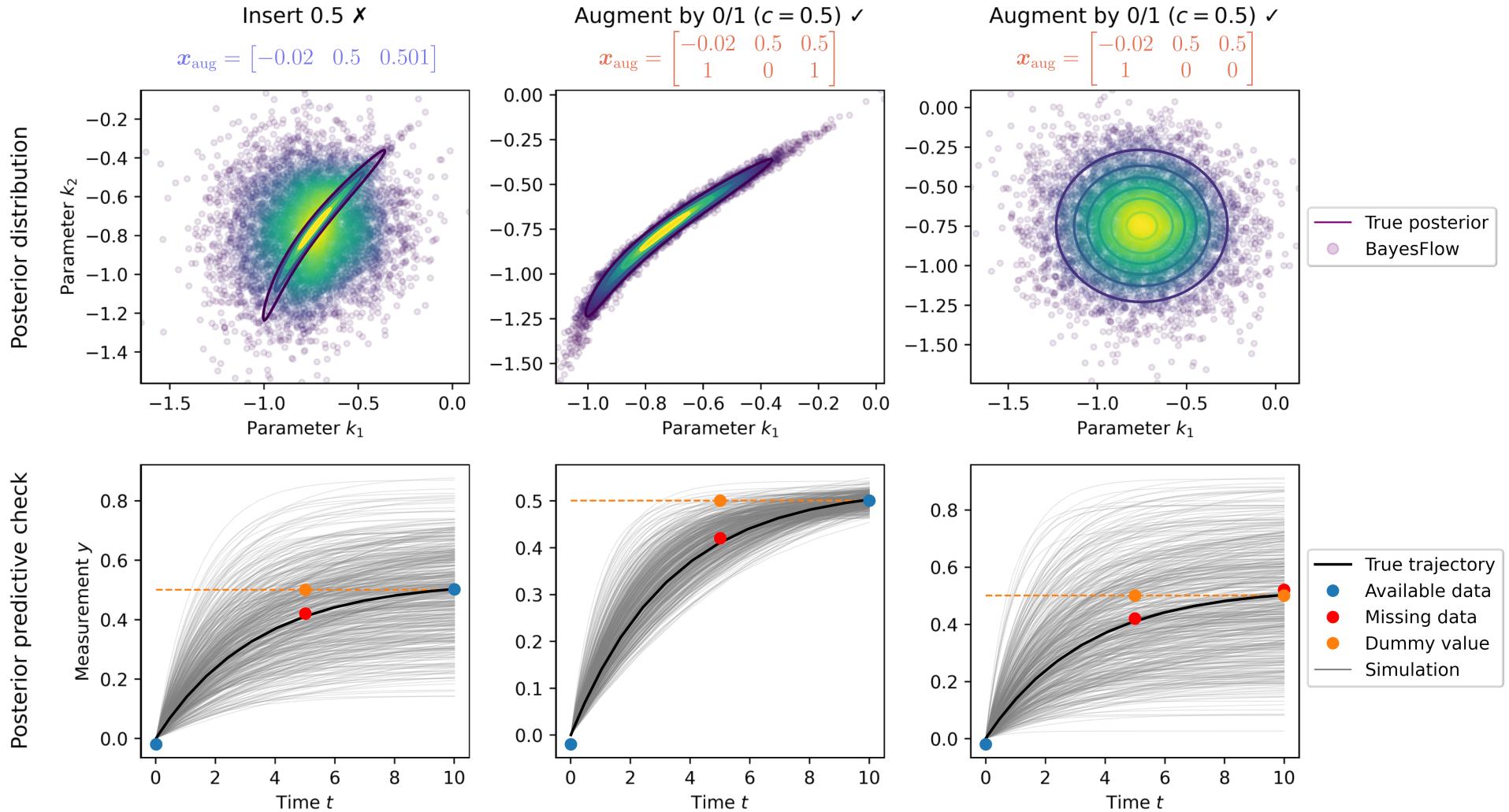
TLDR: With correct uncertainty quantification (which is hard), we just recover the same posterior.

ENCODE MISSING DATA

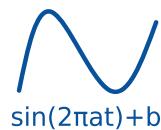


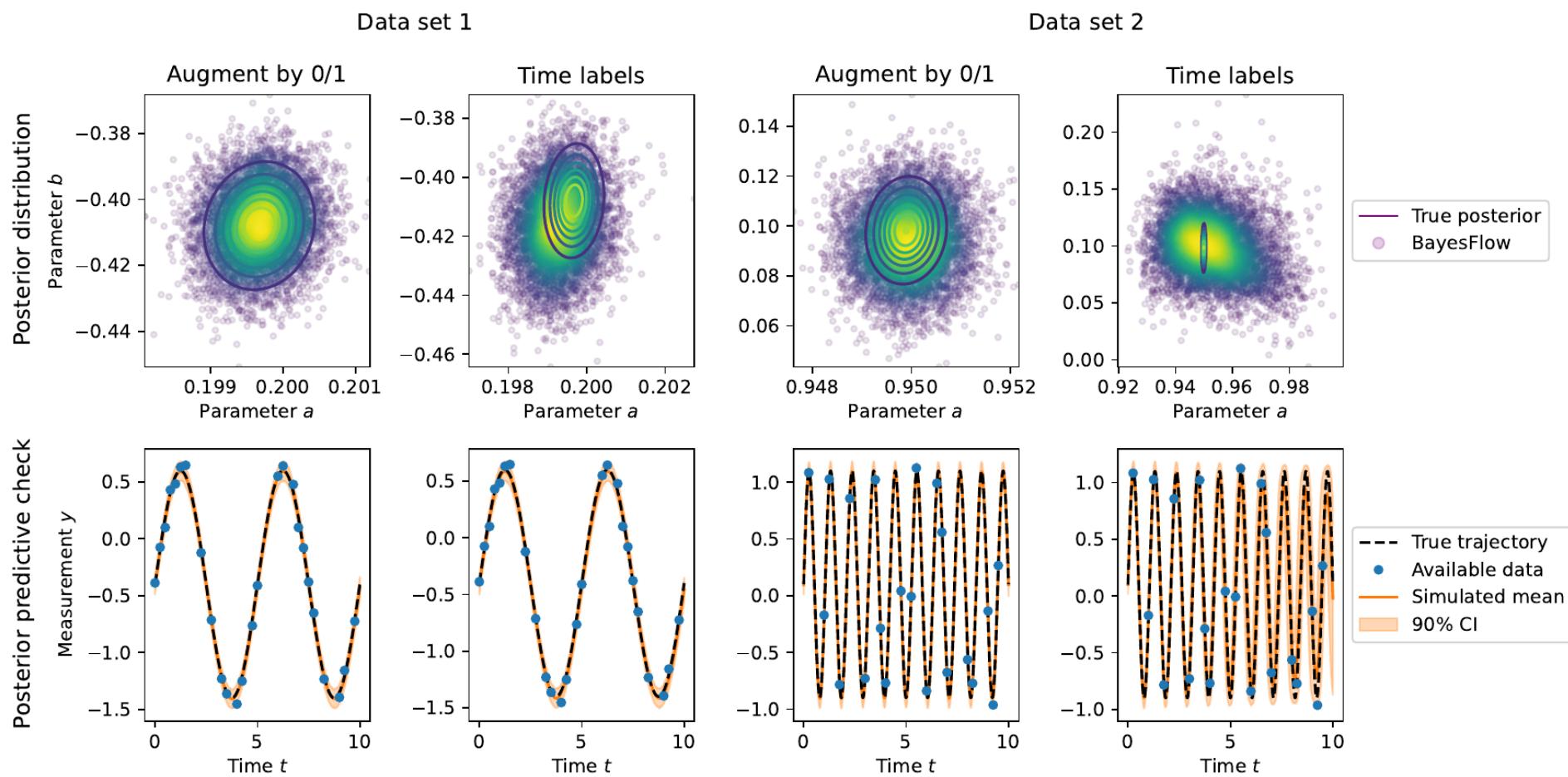
ALL APPROACHES PERFORM WELL ON SIMPLE TEST PROBLEM

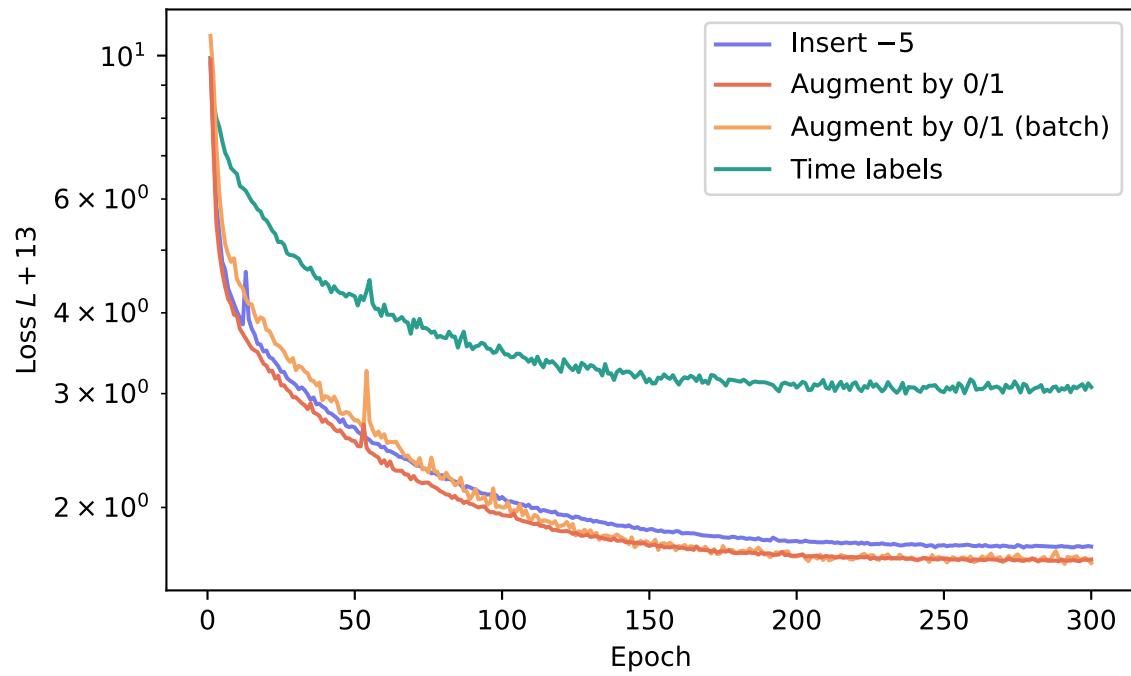
BINARY INDICATOR AUGMENTATION MORE ROBUST FOR AMBIGUOUS FILL-IN VALUES



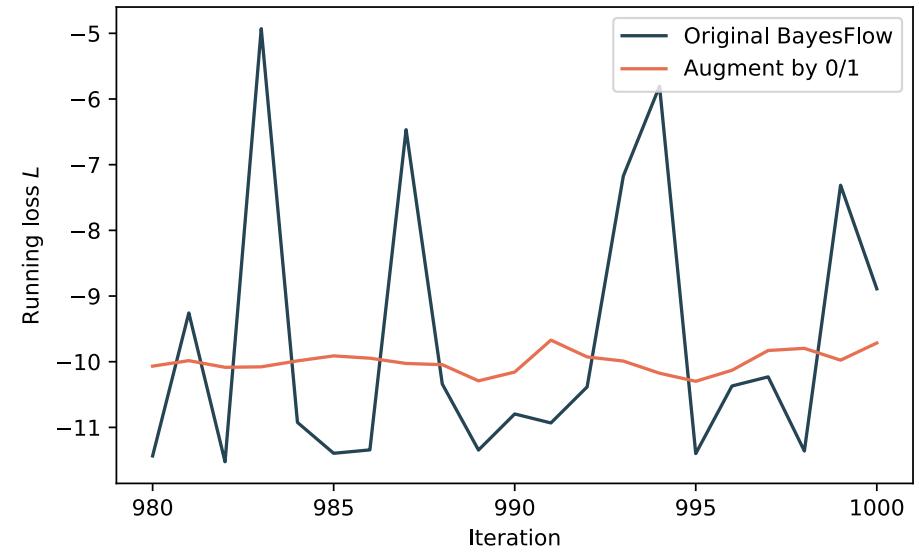
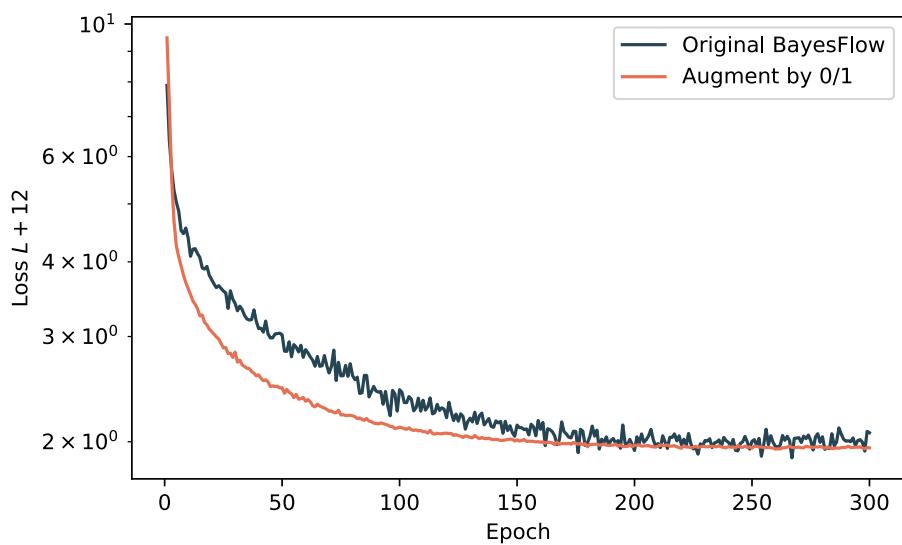
POSITIONAL ENCODING NOT ROBUST ON OSCILLATORY DATA


$$\sin(2\pi at) + b$$



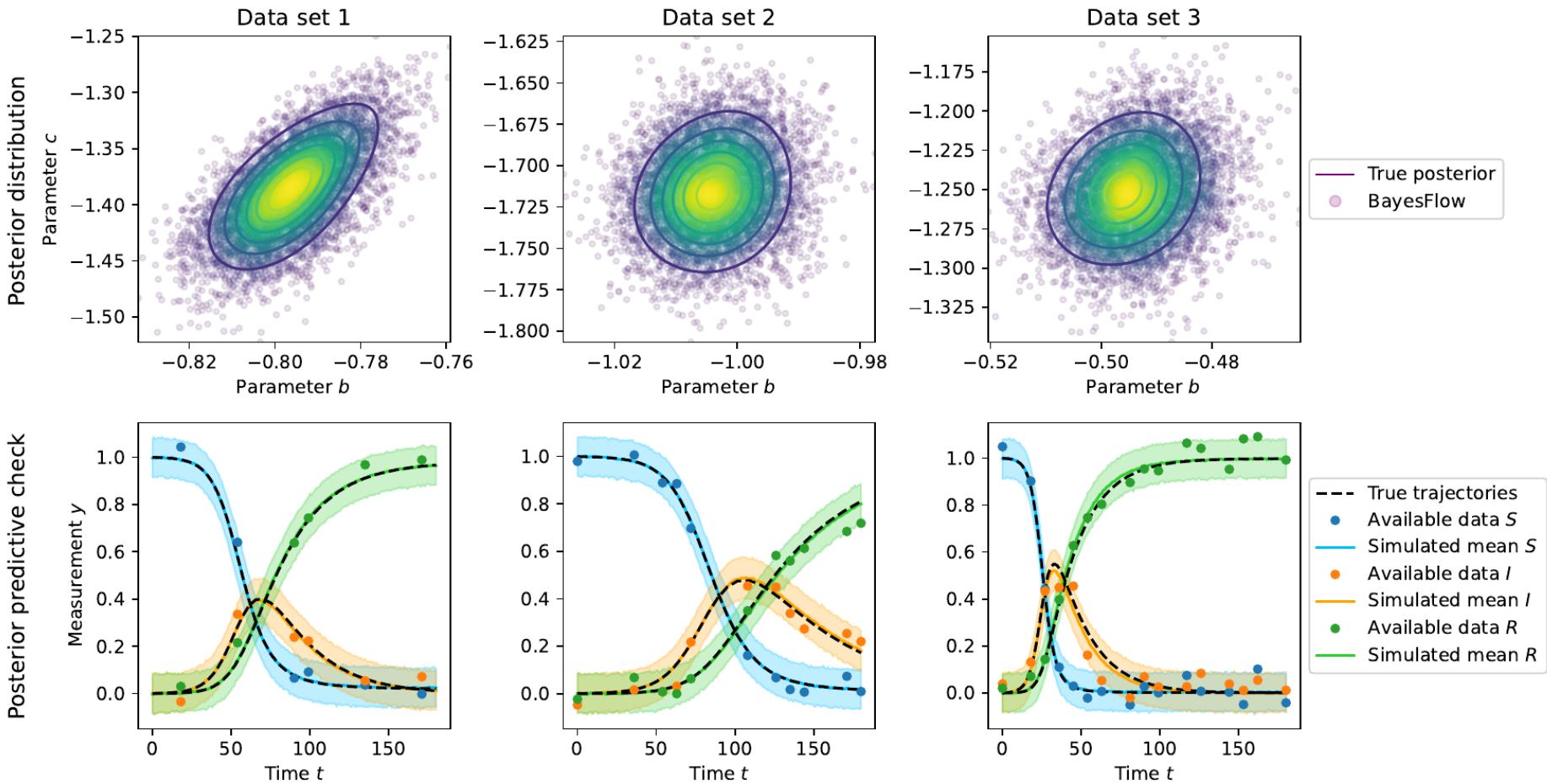


VARIABLE DATASET SIZE AS A SPECIAL CASE OF MISSING DATA

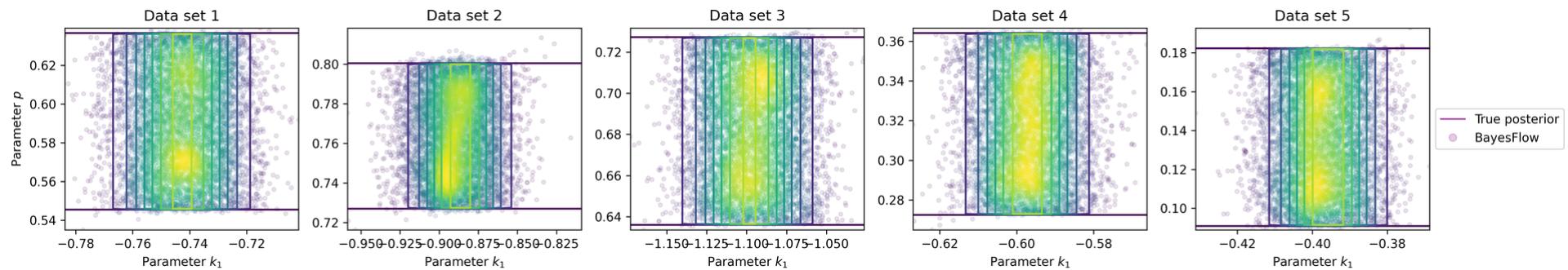


Augment by 0/1 improves performance due to better cost function approximation with individual-specific missingness

SCALES TO COMPLEX DYNAMICS

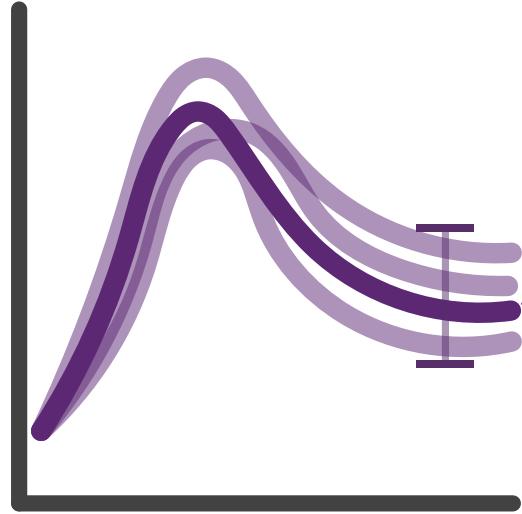


ABLE TO UNRAVEL PARAMETER-DEPENDENT MISSINGNESS



MIXED-EFFECTS MODELS

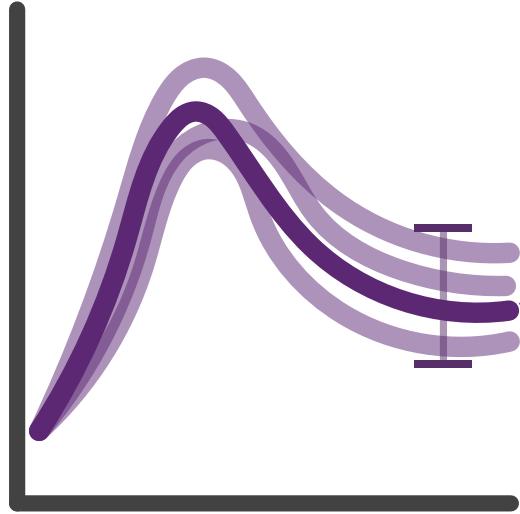
MIXED-EFFECTS MODELING



fixed effects α
random effects β



MIXED-EFFECTS MODELING



fixed effects α
random effects β



dynamical model: $\dot{x} = f(x, \theta)$

observables: $y = h(x, \theta) + \varepsilon$

parameters: $\theta = A\alpha + B\beta, \quad \beta \sim N(0, \Sigma)$

PROBLEM

- **estimate parameters:** maximize over α and Σ the likelihood of data y_{obs} , marginalized over random effects β ,

$$\pi(y_{\text{obs}} | \alpha, \Sigma) = \prod_i \int \color{red} \pi(y_i | \theta) \pi(\theta | \alpha, \Sigma) d\theta$$

PROBLEM

- **estimate parameters:** maximize over α and Σ the likelihood of data y_{obs} , marginalized over random effects β ,

$$\pi(y_{\text{obs}} | \alpha, \Sigma) = \prod_i \int \pi(y_i | \theta) \pi(\theta | \alpha, \Sigma) d\theta$$

- **problem: evaluating these (high-dim) integrals is challenging,** especially with many individuals



AN AMORTIZED APPROACH

AN AMORTIZED APPROACH

- idea: rewrite in terms of an **individual-specific posterior**:

$$\begin{aligned}\pi(y_{\text{obs}} | \alpha, \Sigma) &= \prod_i \pi(y_i) \int \pi(\theta | y_i) \frac{\pi(\theta | \alpha, \Sigma)}{\pi(\theta)} d\theta \\ &= \prod_i \pi(y_i) \mathbb{E}_{\theta \sim \pi(\theta | y_i)} \left[\frac{\pi(\theta | \alpha, \Sigma)}{\pi(\theta)} \right]\end{aligned}$$

AN AMORTIZED APPROACH

- idea: rewrite in terms of an **individual-specific posterior**:

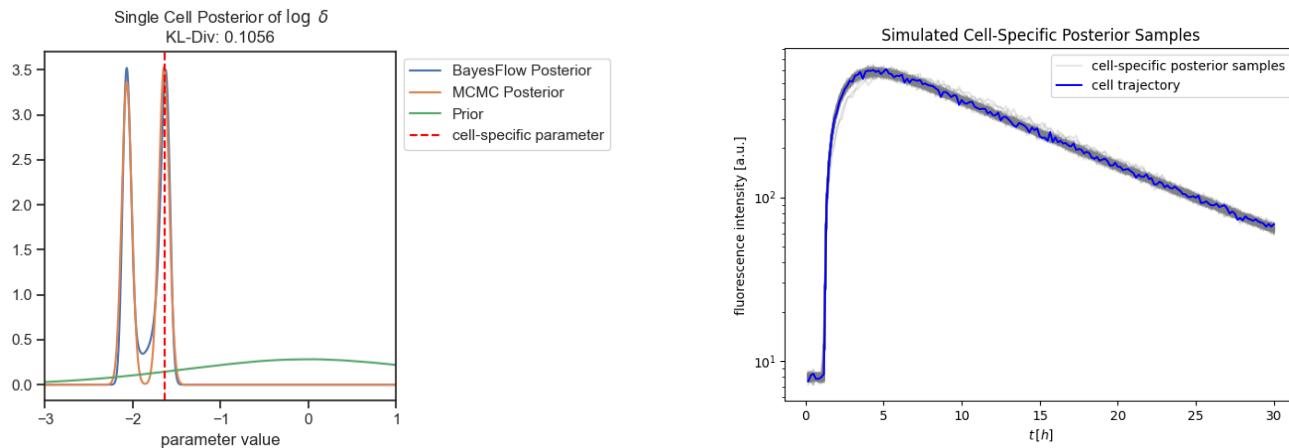
$$\begin{aligned}\pi(y_{\text{obs}} | \alpha, \Sigma) &= \prod_i \pi(y_i) \int \pi(\theta | y_i) \frac{\pi(\theta | \alpha, \Sigma)}{\pi(\theta)} d\theta \\ &= \prod_i \pi(y_i) \mathbb{E}_{\theta \sim \pi(\theta | y_i)} \left[\frac{\pi(\theta | \alpha, \Sigma)}{\pi(\theta)} \right]\end{aligned}$$

- ... and approximate the posterior using a **neural density estimator** trained on synthetic data!

EVALUATION

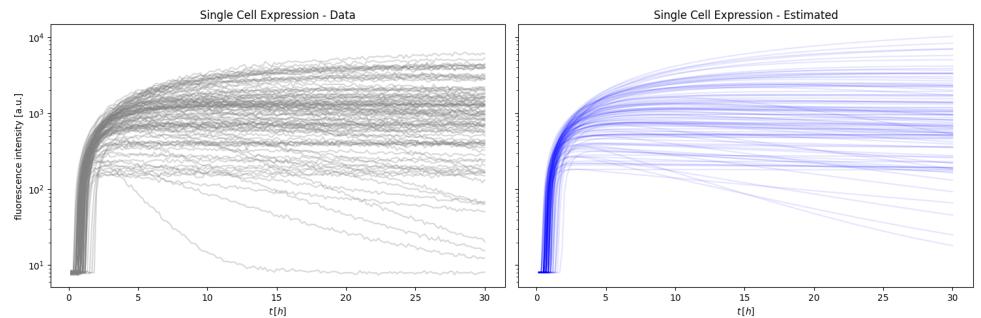
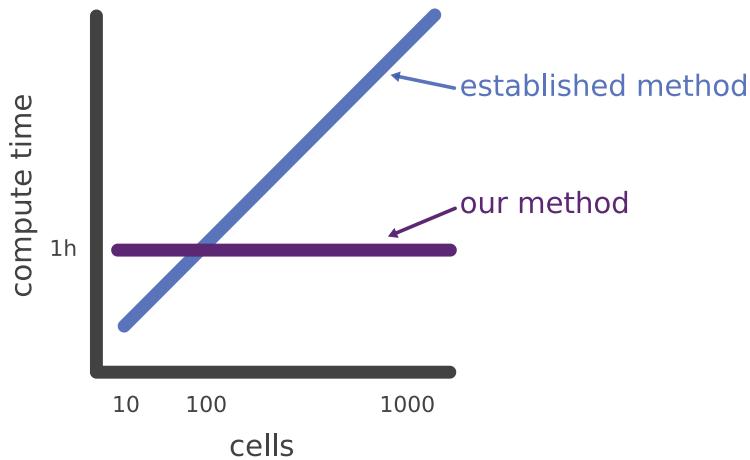
- ✓ INN **fits the posterior** well
- ✓ after training once (h), the optimization is **very fast** ($s - \min$)
- ✓ allows to **easily test hypotheses**
- ✓ allows considering **stochastic models**
- ✓ even **uncertainty analysis** easily possible

EVALUATION



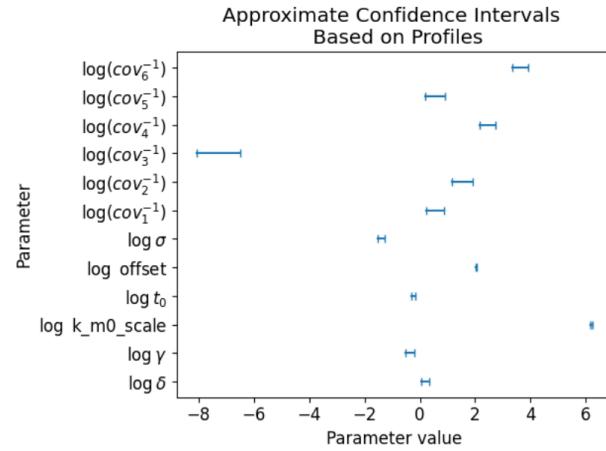
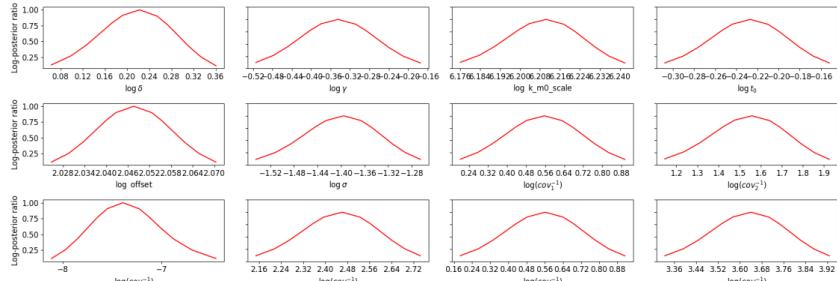
- ✓ INN fits the posterior well
- ✓ after training once (h), the optimization is very fast ($s - min$)
- ✓ allows to easily test hypotheses
- ✓ allows considering stochastic models
- ✓ even uncertainty analysis easily possible

EVALUATION



- ✓ INN fits the posterior well
- ✓ after training once (h), the optimization is very fast ($s - min$)
- ✓ allows to easily test hypotheses
- ✓ allows considering stochastic models
- ✓ even uncertainty analysis easily possible

EVALUATION



- ✓ INN fits the posterior well
- ✓ after training once (h), the optimization is very fast ($s - min$)
- ✓ allows to easily test hypotheses
- ✓ allows considering stochastic models
- ✓ even uncertainty analysis easily possible

OUTLOOK

OUTLOOK

- applications in pharmacology and single-cell biology
- mechanistic insights to inform drug therapy
- combine with equation learning
- integrate image and tabular data in one framework
- large-scale modeling
- federated learning
- software

THANKS! QUESTIONS?



yannik-schaelte



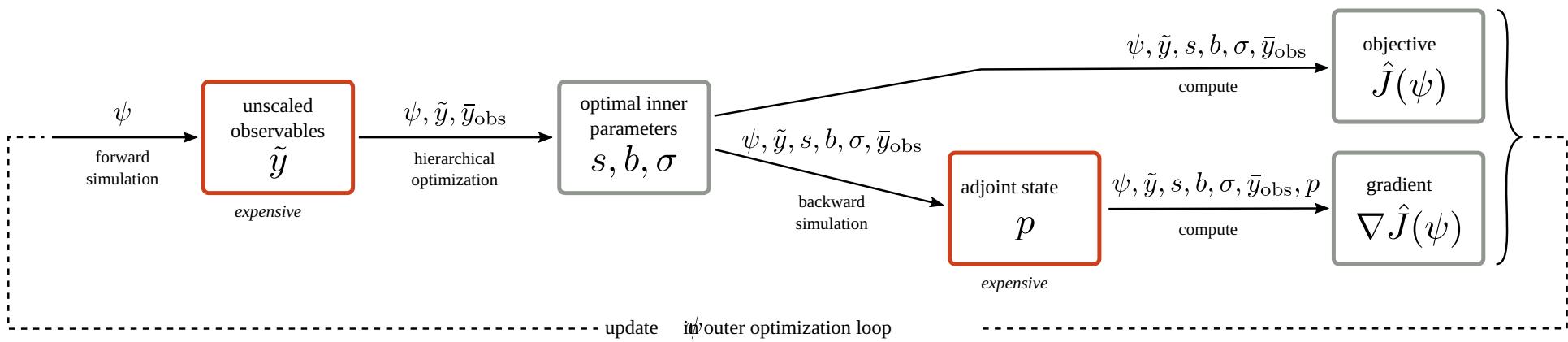
@yannik_schaelte



yannikschaelte

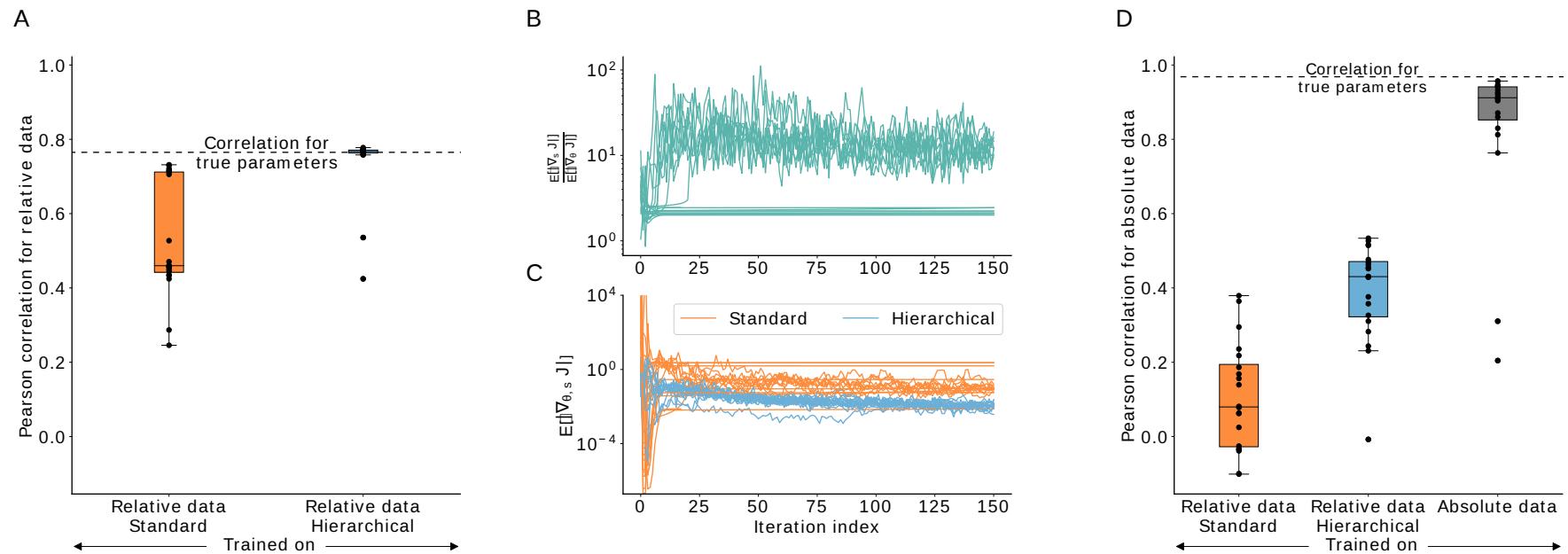
BACKUP

ILLUSTRATION OF ADJOINT-HIERARCHICAL EVALUATION SCHEME



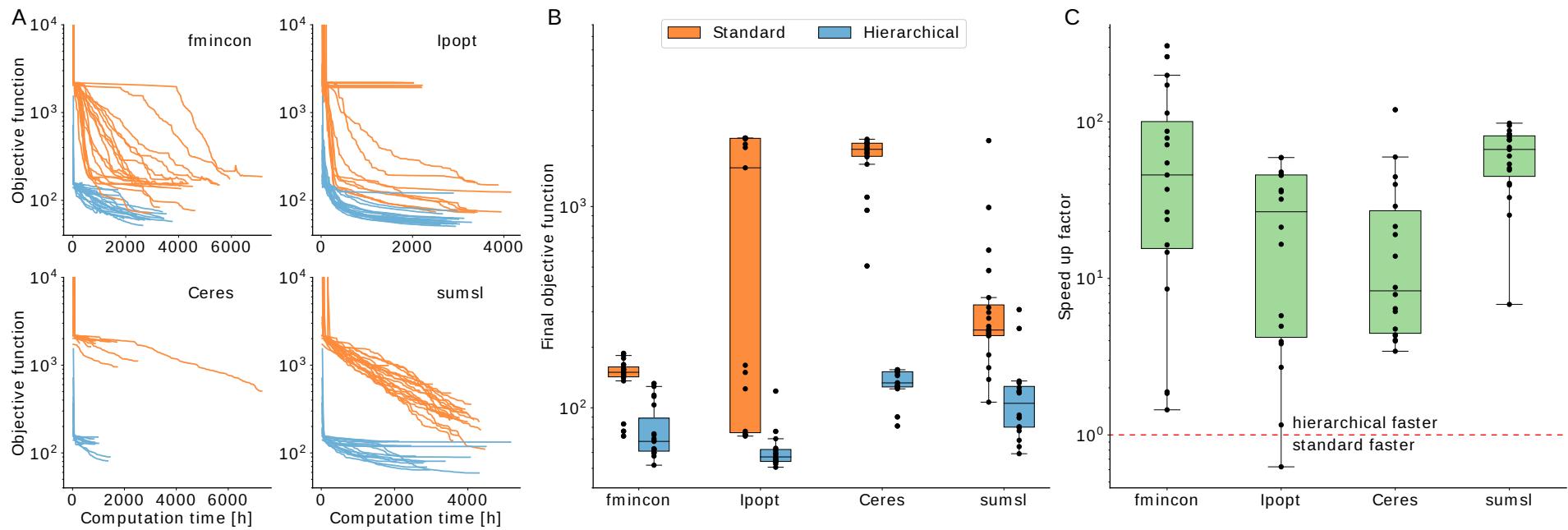
from: Schmiester, Schälte et al., Bioinformatics 2020

CONVERGENCE OF STANDARD AND HIERARCHICAL OPTIMIZATION



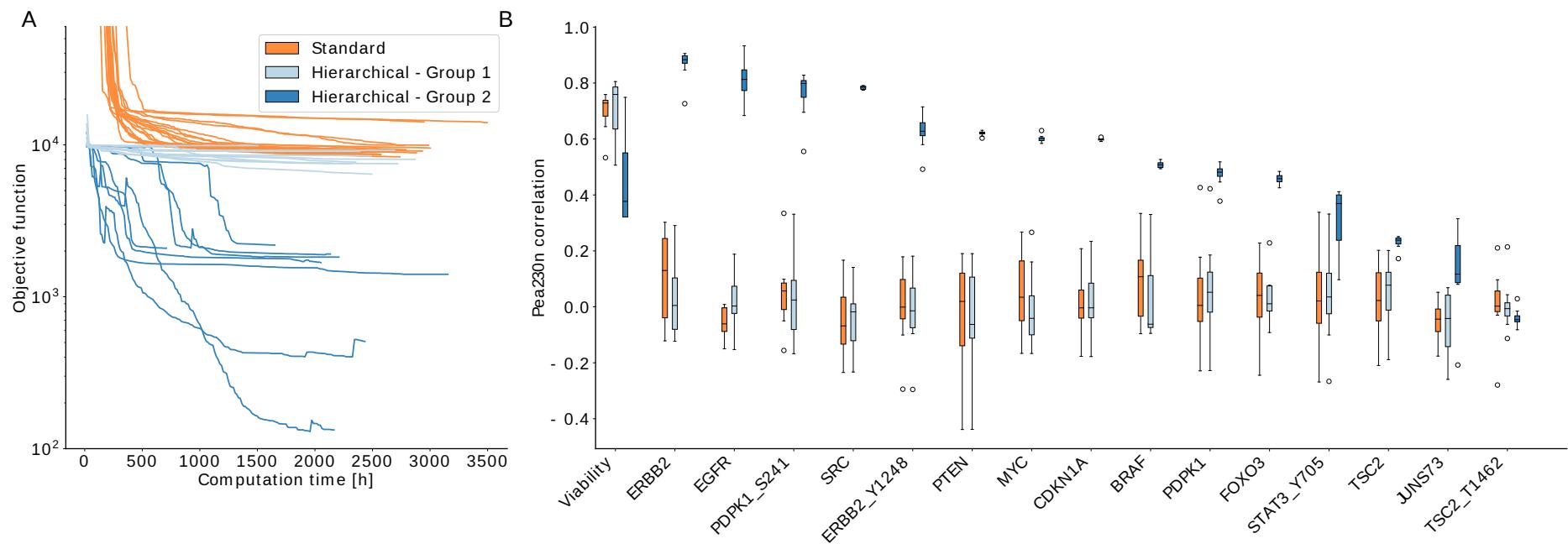
from: Schmiester, Schälte et al., Bioinformatics 2020

COMPUTATIONAL EFFICIENCY OF STANDARD AND HIERARCHICAL OPTIMIZATION



from: Schmiester, Schälte et al., Bioinformatics 2020

INTEGRATION OF HETEROGENEOUS DATA USING HIERARCHICAL OPTIMIZATION



from: Schmiester, Schälte et al., Bioinformatics 2020



github.com/icb-dcm/pyabc

Klinger et al., Bioinformatics 2018 and Schälte et al., JOSS 2022

```
# specify problem and parallelization
abc = ABCSMC(model, prior, distance, sampler)
# pass data
abc.new(db, data)
# run it
abc.run()
```



flexible

modular
configure
extend



user-friendly

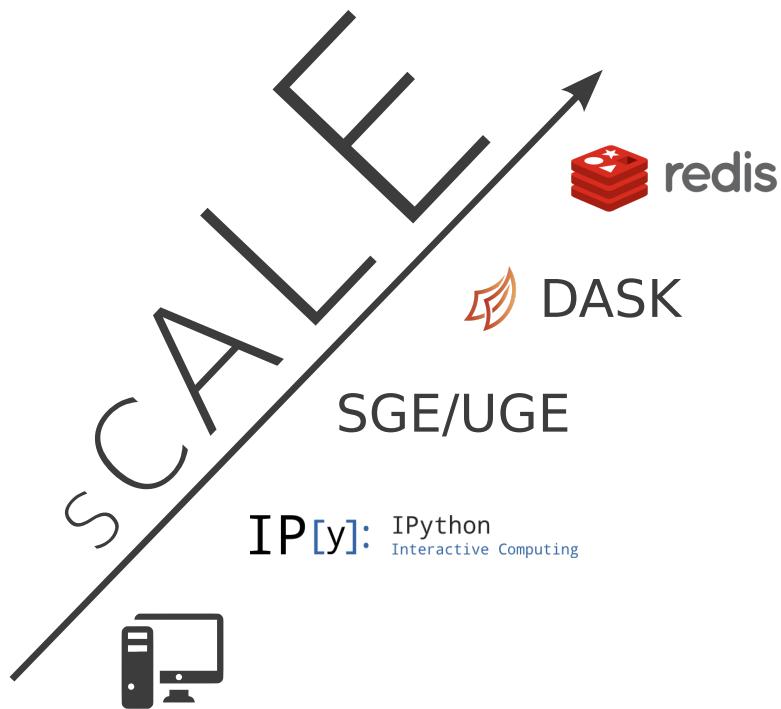
documented
self-tuning
robust



scalable

1 to 1,000s cores
efficient
tested

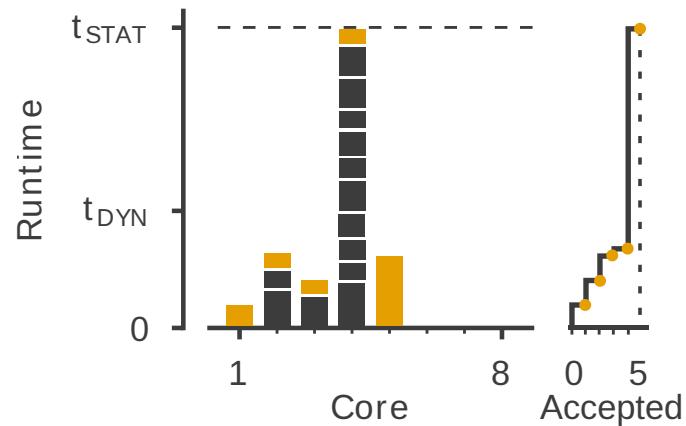
PARALLEL BACKENDS: 1 TO 1,000S CORES



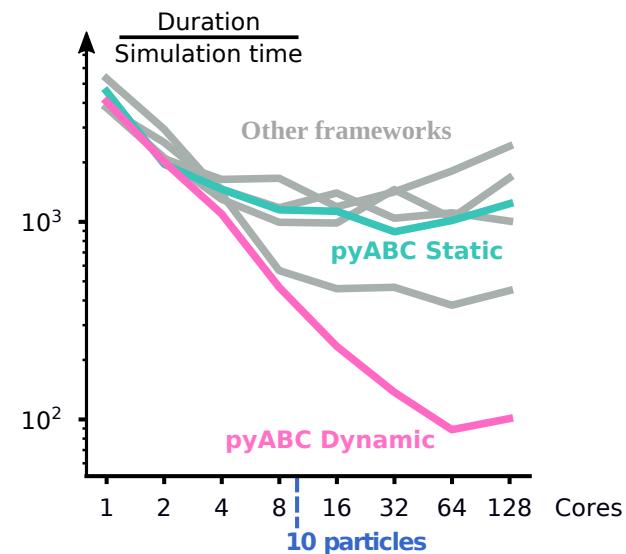
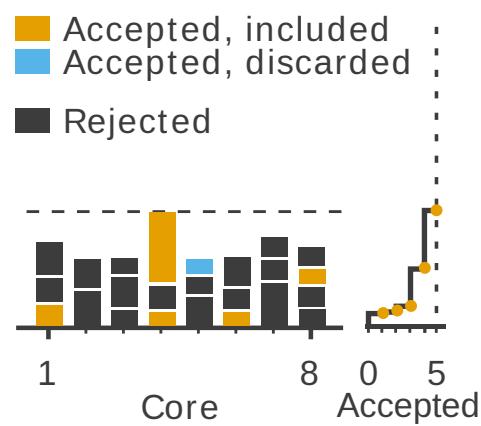
PARALLELIZATION STRATEGIES

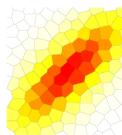
Klinger et al., CMSB Proceedings 2017

Static Scheduling



Dynamic Scheduling

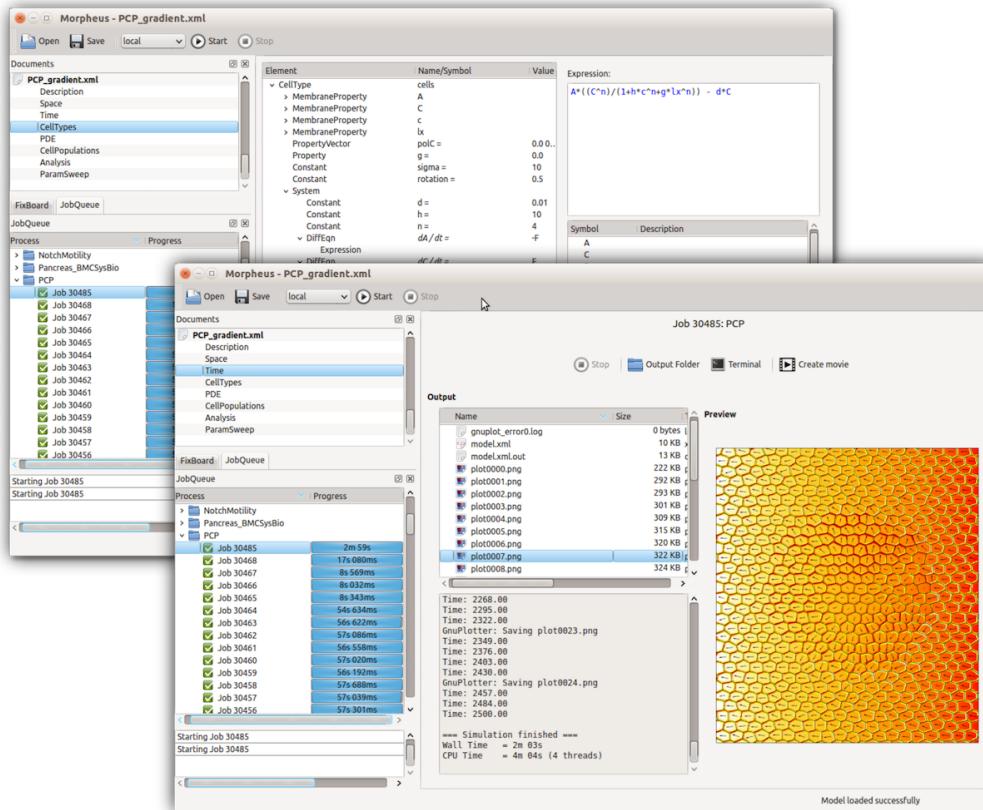




FitMultiCell

fitmulticell.gitlab.io

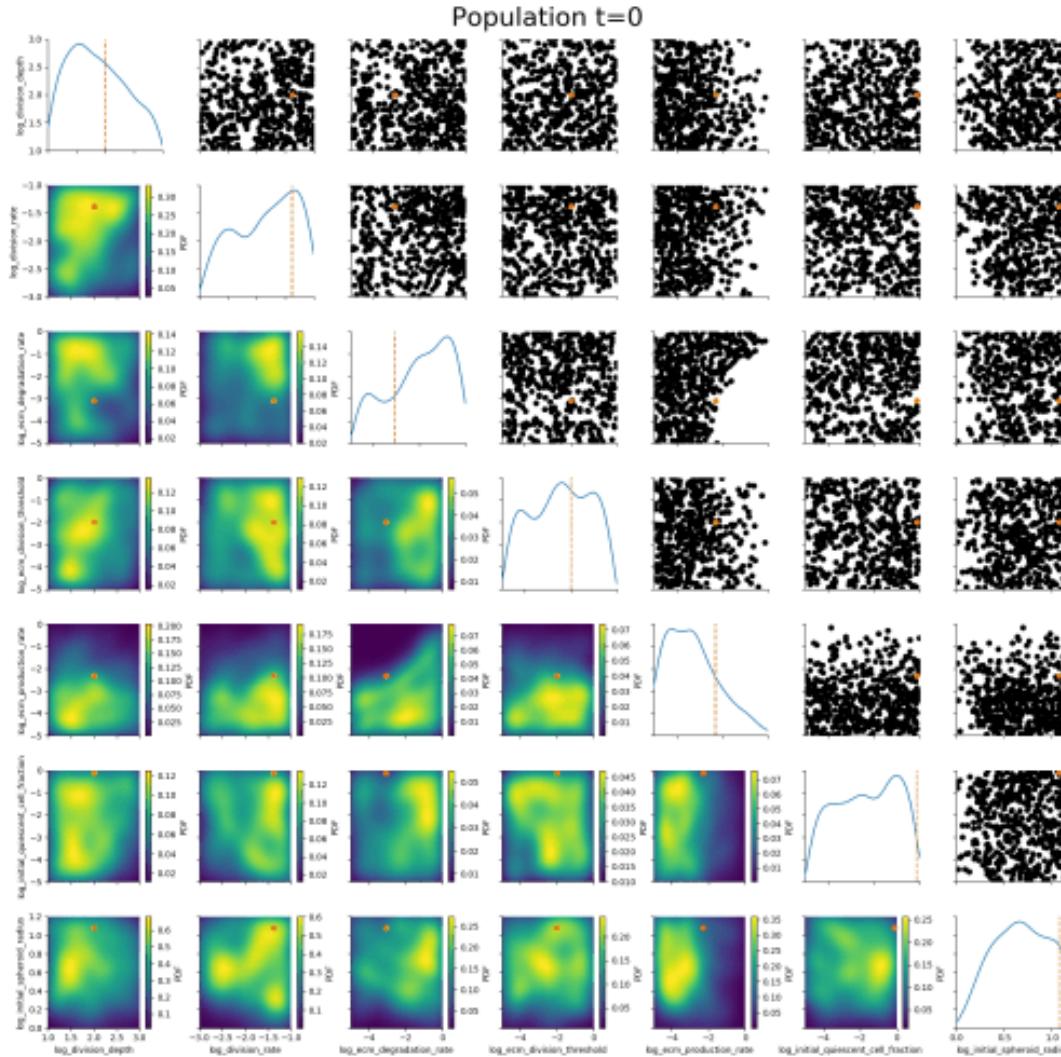
a platform for modeling, simulation and inference for multi-scale multi-cellular



models

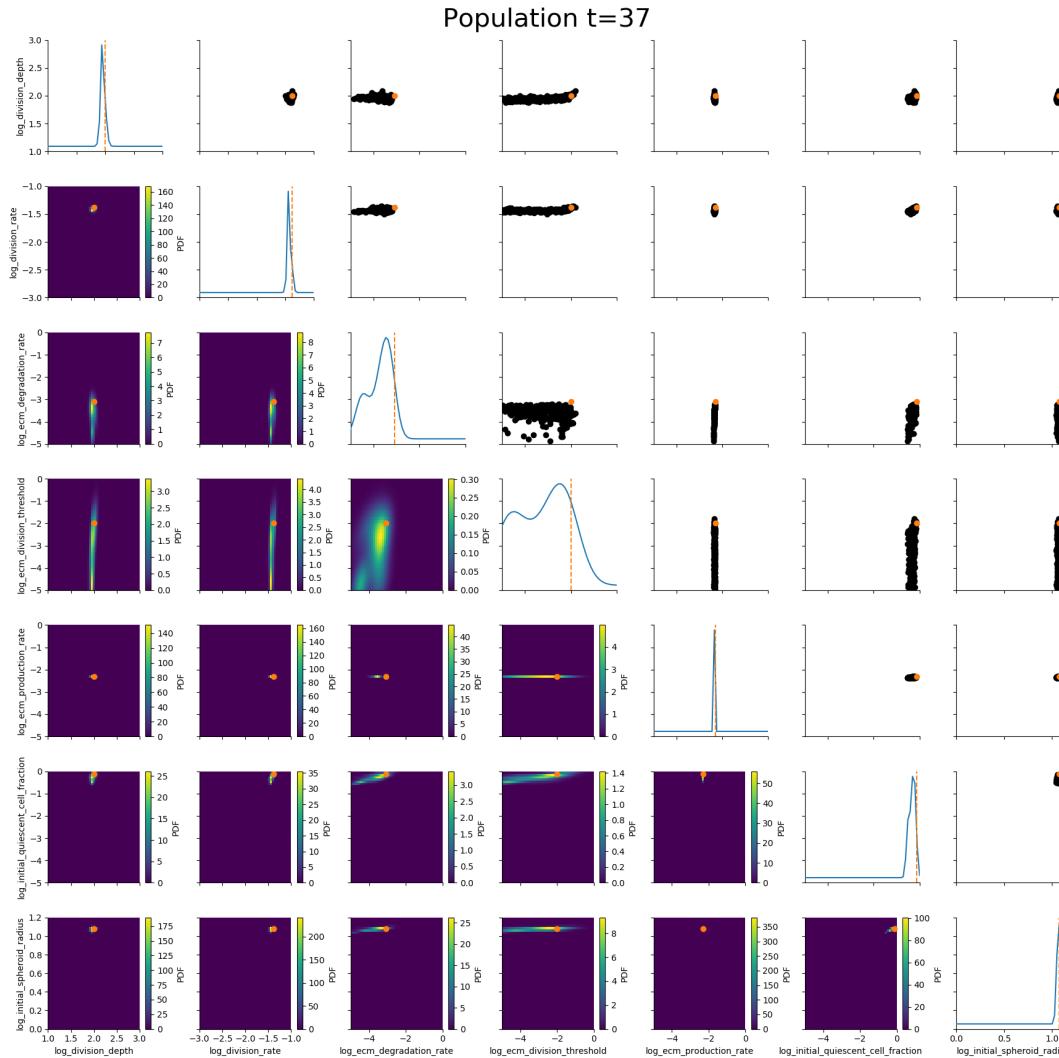
Starruß et al., Bioinformatics 2014; Alamoudi et al., NIC Proc. 2022; Alamoudi et al., bioRxiv 2023

APPLICATION EXAMPLE

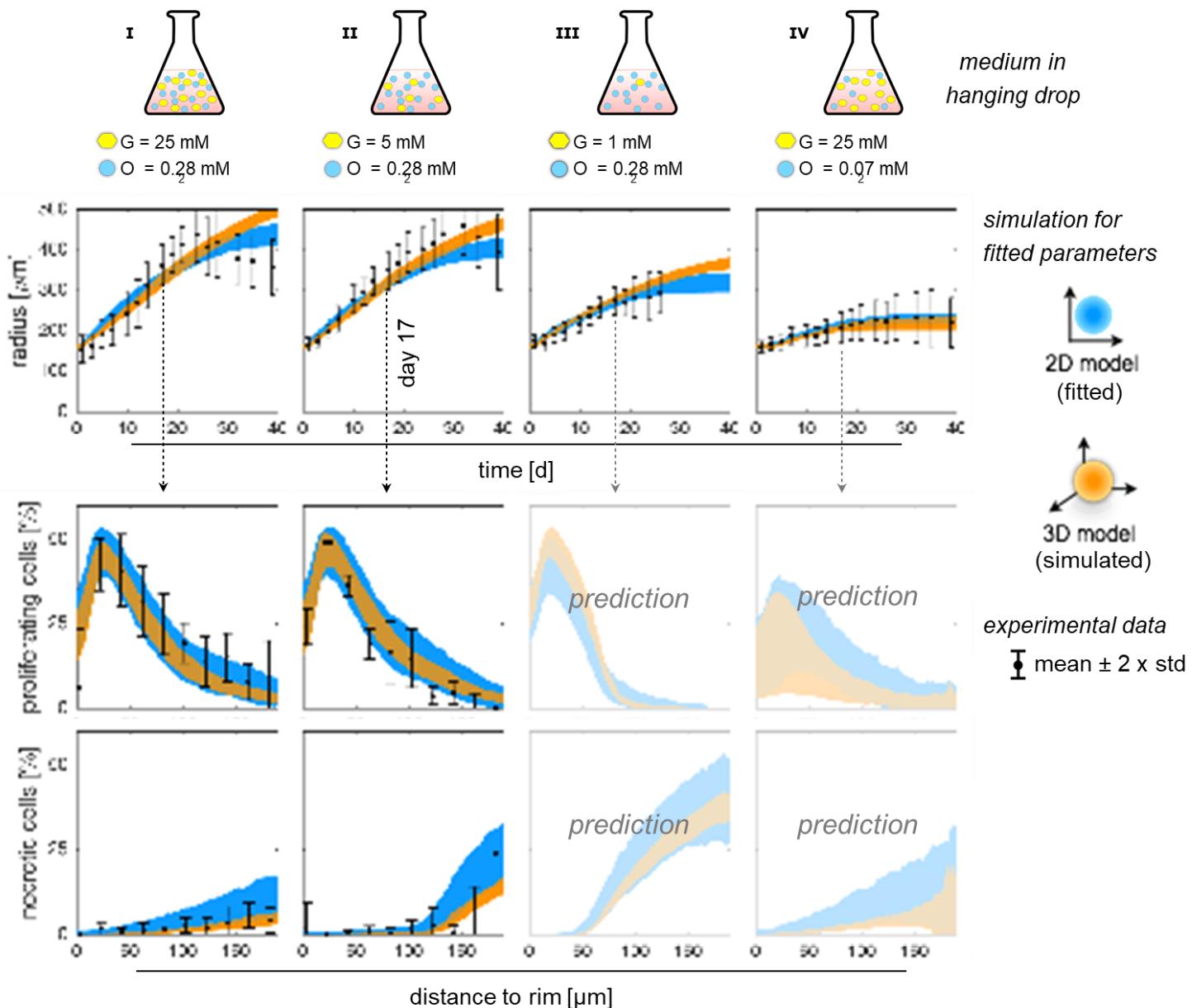


ABC worked where many other methods had failed.

APPLICATION EXAMPLE

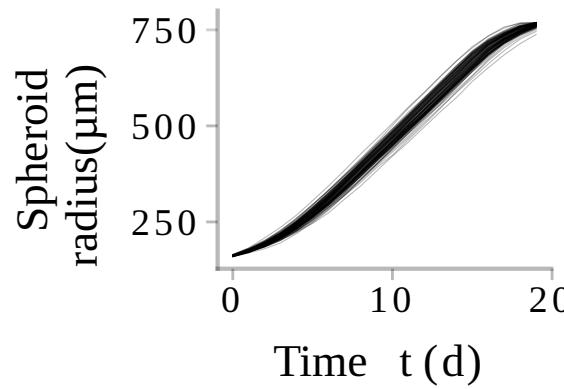


ABC worked where many other methods had failed.

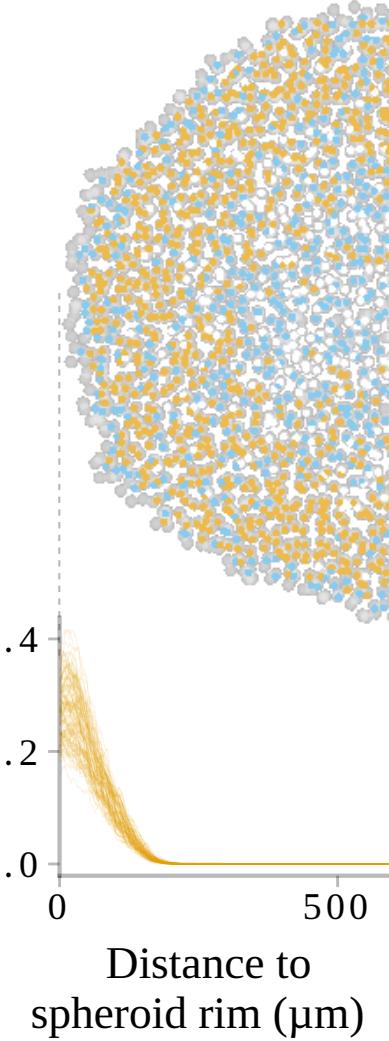


Uncertainty-aware predictions, easy data integration.

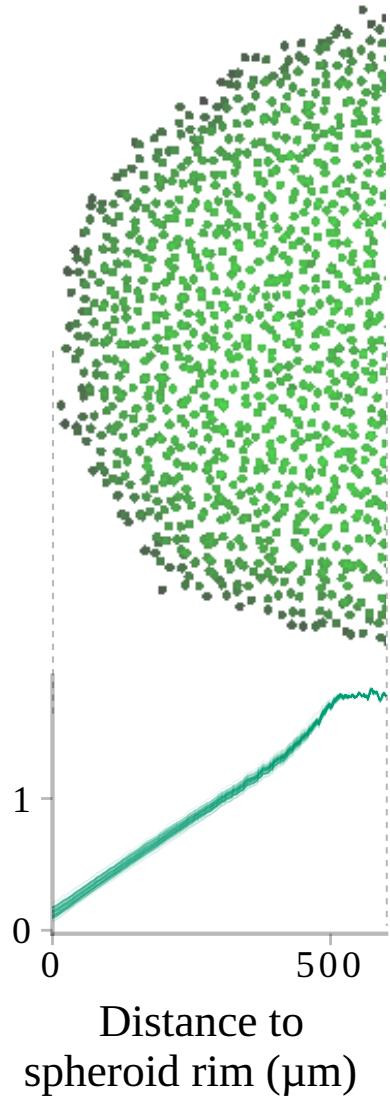
DEFINE SUMMARY STATISTICS

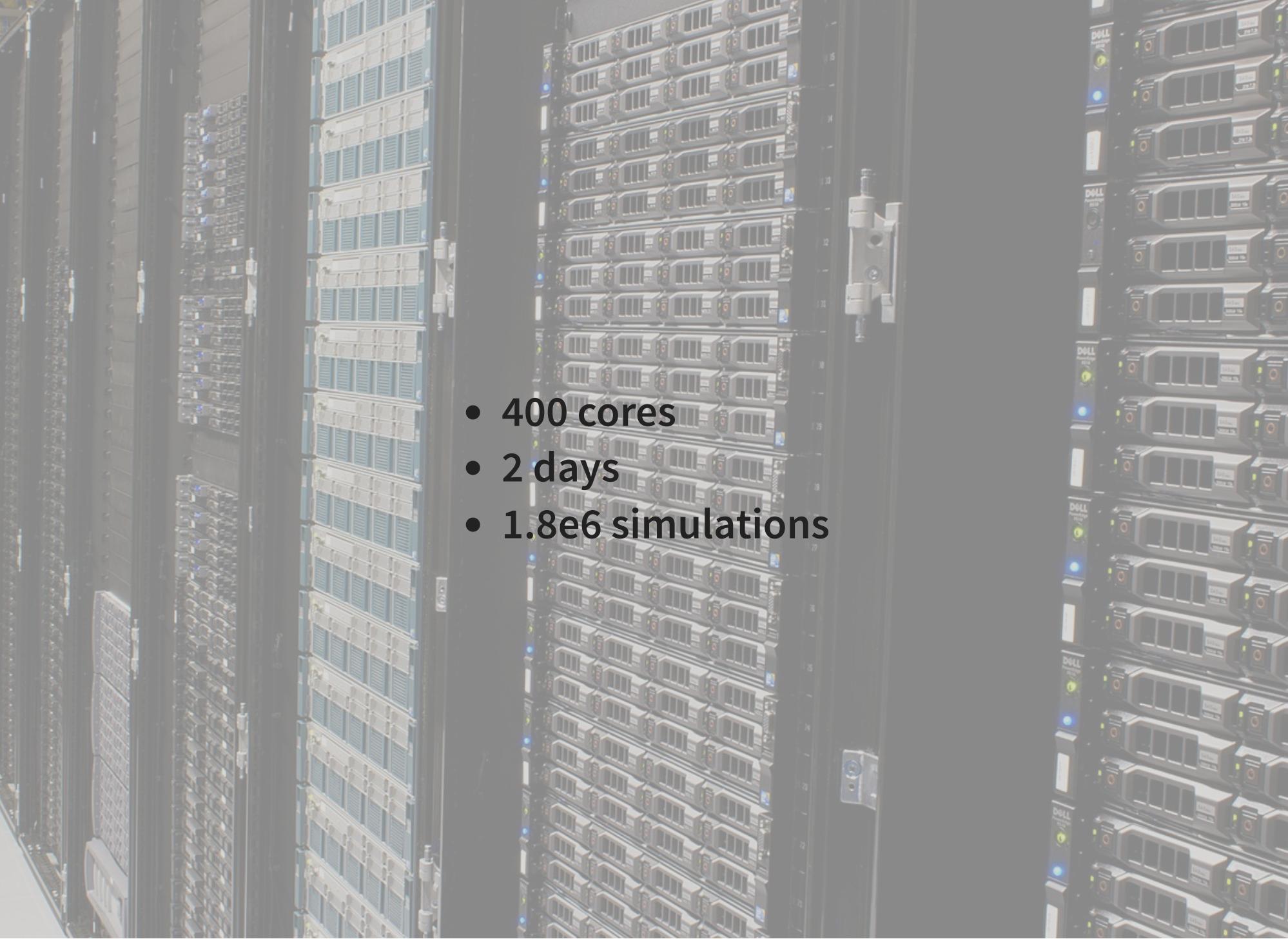


Fraction of
proliferating
cells



Extracellular
matrix
density



- 
- A photograph of a server rack filled with Dell PowerEdge server units. The servers are arranged in a grid pattern, with many blue indicator lights visible on their front panels. The rack has a dark frame and is set against a white background.
- 400 cores
 - 2 days
 - 1.8e6 simulations

THEOREM (EXACT INFERENCE)

Using the **modified kernel** with $c > \pi(\bar{y}_{\text{obs}}|y, \theta) \forall y, \theta$, we sample from the **true posterior**

$$\pi_{\text{ABC}}(\theta|\bar{y}_{\text{obs}}) = \pi(\theta|\bar{y}_{\text{obs}}) \propto \int \pi(\bar{y}_{\text{obs}}|y, \theta)p(y|\theta) dy \cdot \pi(\theta)$$

assuming **noisy data** $\bar{y}_{\text{obs}} \sim \pi(\bar{y}|y, \theta)$.

- non-trivial noise allows to do **exact likelihood-free inference**
- applicable to **stochastic models**
- **parameterized noise model**

THEOREM (OPTIMAL SUMMARY STATISTICS)

[...] Given $\lambda : \mathbb{R}^{n_\theta} \rightarrow \mathbb{R}^{n_\lambda}$ such that $\mathbb{E}_{\pi(\theta)}[|\lambda(\theta)|] < \infty$, define **summary statistics** as the conditional expectation

$$s(y) := \mathbb{E}[\lambda(\Theta)|Y = y] = \int \lambda(\theta)\pi(\theta|y)d\theta.$$

Then, it holds $\|\mathbb{E}_{\pi_{ABC,\varepsilon}}[\lambda(\Theta)|s(y_{obs})] - s(y_{obs})\| \leq \varepsilon$, and therefore

$$\lim_{\varepsilon \rightarrow 0} \mathbb{E}_{\pi_{ABC,\varepsilon}}[\lambda(\Theta)|s(y_{obs})] = \mathbb{E}[\lambda(\Theta)|Y = y_{obs}].$$

In practice: Train regression model $s : y \mapsto \lambda(\theta) = (\theta^1, \dots, \theta^k)$.

SCALE-NORMALIZING AND OUTLIER-ROBUST ADAPTIVE DISTANCES

