



HELMHOLTZ  
MUNICH

# Amortized inference for pharmacological mixed-effects models with incomplete data

JOACHIM  
HERZ  
STIFTUNG



Yannik Schälte

JHS Add-On 8th Cohort First Meeting, Hamburg 2023-02-11



@yannik\_schaelte



yannikschaelte

[yannikschaelte.github.io/pres\\_npe\\_hjs\\_2023](https://yannikschaelte.github.io/pres_npe_hjs_2023)

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

**mechanistic modeling**

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

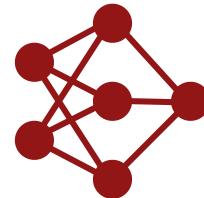
## **mechanistic modeling**

- theory-driven
- interpretation
- testability

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

## **mechanistic modeling**

- theory-driven
- interpretation
- testability

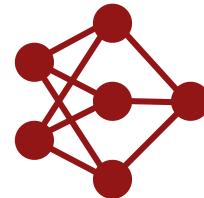


## **machine learning**

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

## **mechanistic modeling**

- theory-driven
- interpretation
- testability



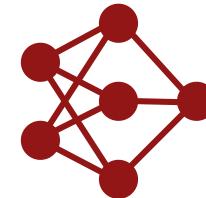
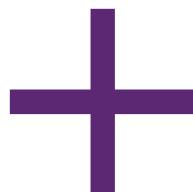
## **machine learning**

- data-driven
- large data
- automation

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

## mechanistic modeling

- theory-driven
- interpretation
- testability



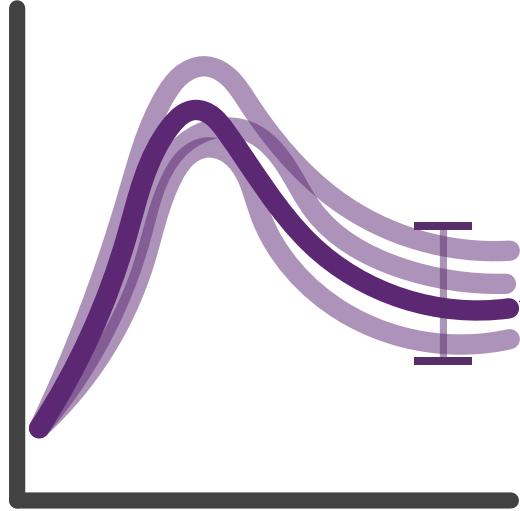
## machine learning

- data-driven
- large data
- automation

**sciML**

model-based data-efficient ML

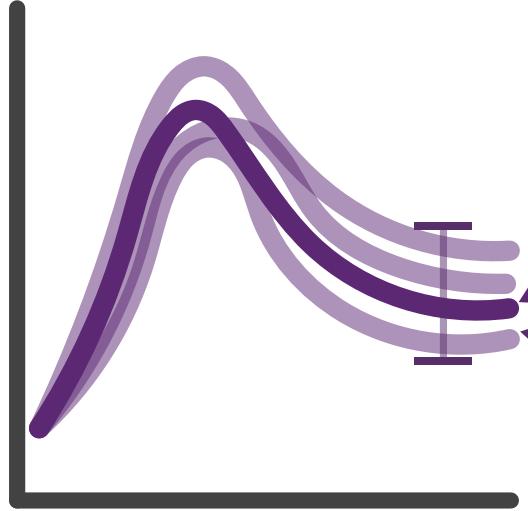
# MIXED-EFFECTS MODELING



fixed effects  $\alpha$   
random effects  $\beta$



# MIXED-EFFECTS MODELING



fixed effects  $\alpha$   
random effects  $\beta$



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

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

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

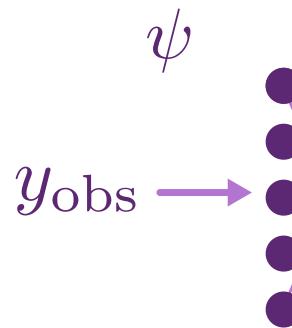
# THE PROBLEM

# THE PROBLEM

# AMORTIZED INFERENCE VIA INVERTIBLE NEURAL NETWORKS

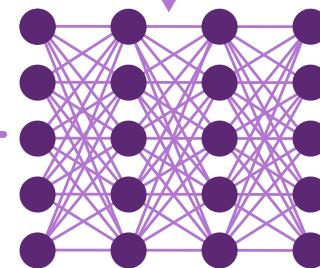
you have to solve many similar problems? amortize the solution!

Summary network



Invertible network

$\phi$

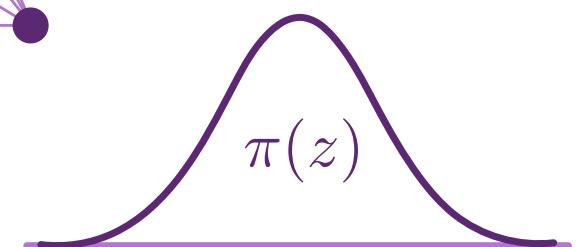


Sampling

$$z \sim \mathcal{N}_d(0, I)$$

Approximate posterior

$$\pi_\phi(\theta | y = \tilde{y}_{\text{obs}})$$





- parameter estimation requires repeatedly evaluating an integral

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

- parameter estimation requires repeatedly evaluating an integral

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

## OUR METHOD

- parameter estimation requires repeatedly evaluating an integral

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

## OUR METHOD

- parameter estimation requires repeatedly evaluating an integral

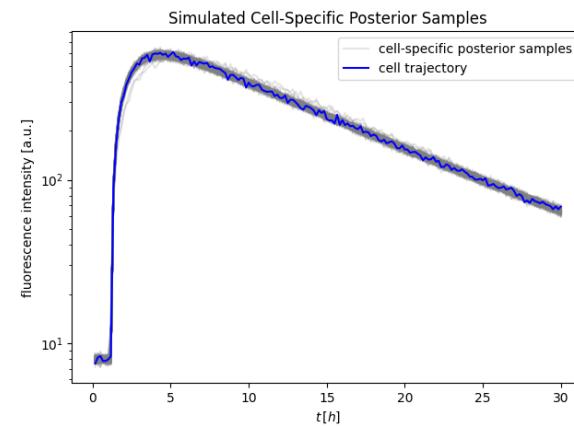
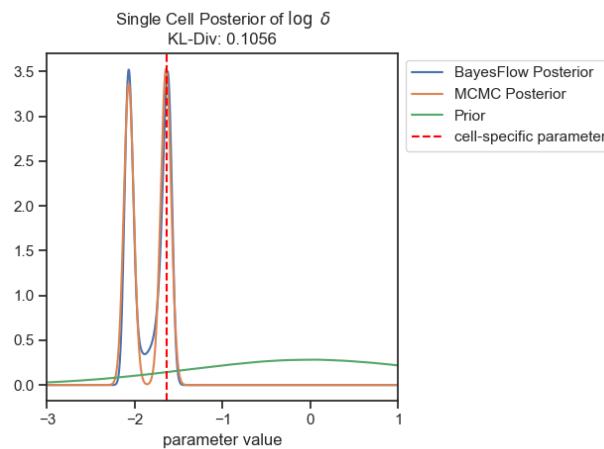
$$\pi(y_{\text{obs}} | \alpha, \Sigma) = \prod_i \int \pi(y_i | \theta) \pi(\theta | \alpha, \Sigma) d\theta$$

## OUR METHOD

# RESULTS

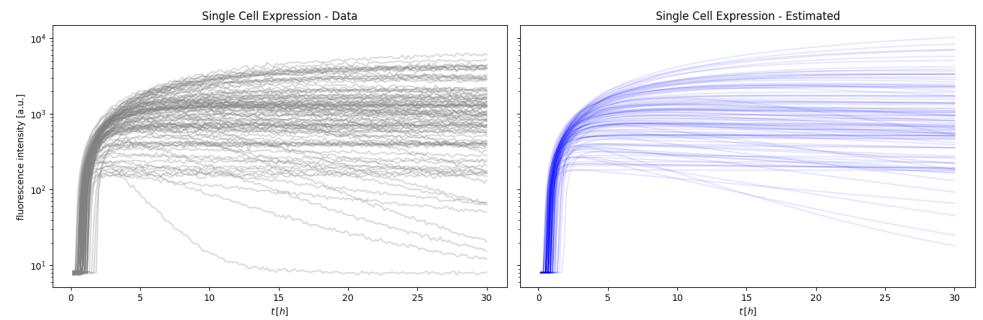
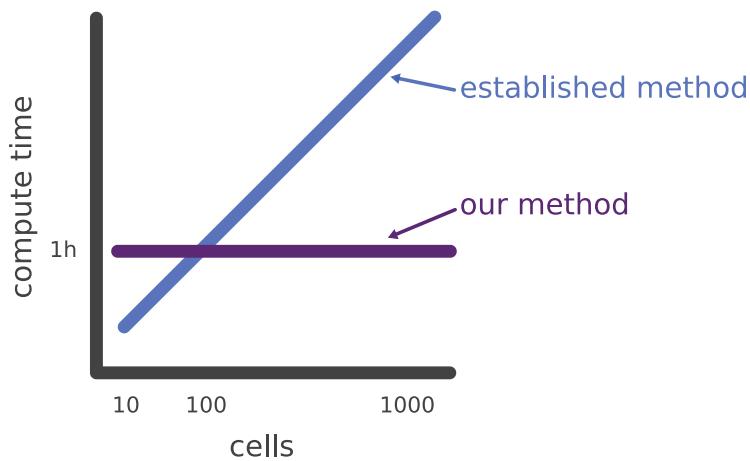
# RESULTS

- ✓ NN fits the posterior well



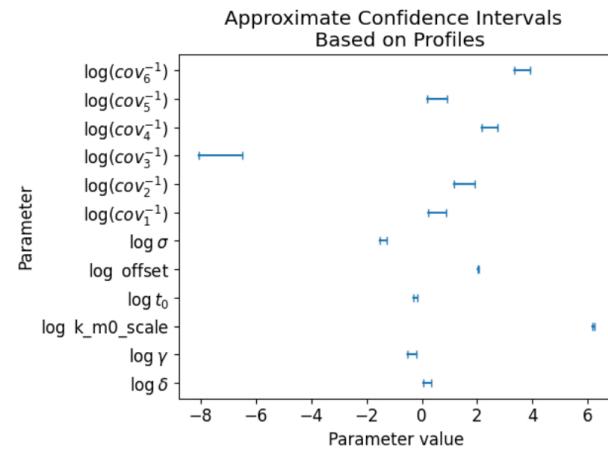
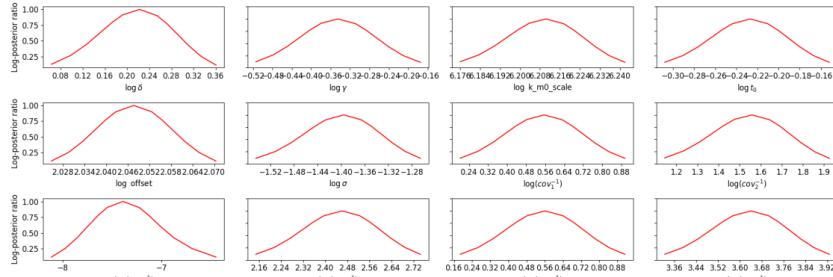
# RESULTS

- ✓ NN fits the posterior well
- ✓ after training once ( $1 - 2h$ ), the optimization is **very fast** ( $s - min$ )



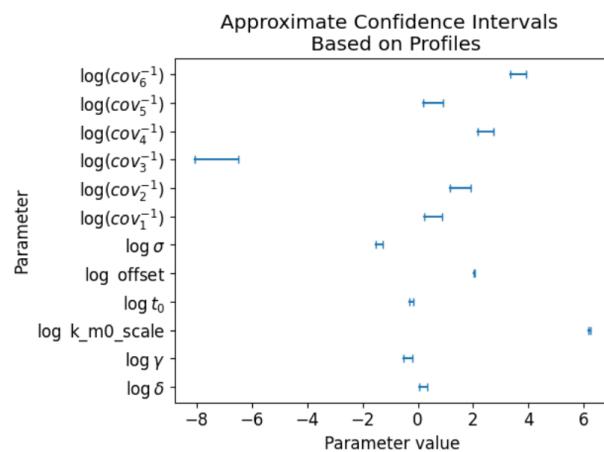
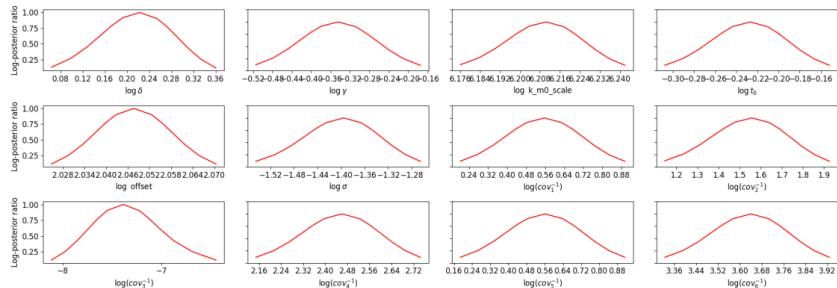
# RESULTS

- ✓ NN fits the posterior well
- ✓ after training once ( $1 - 2h$ ), the optimization is very fast ( $s - min$ )
- ✓ even uncertainty analysis easily possible



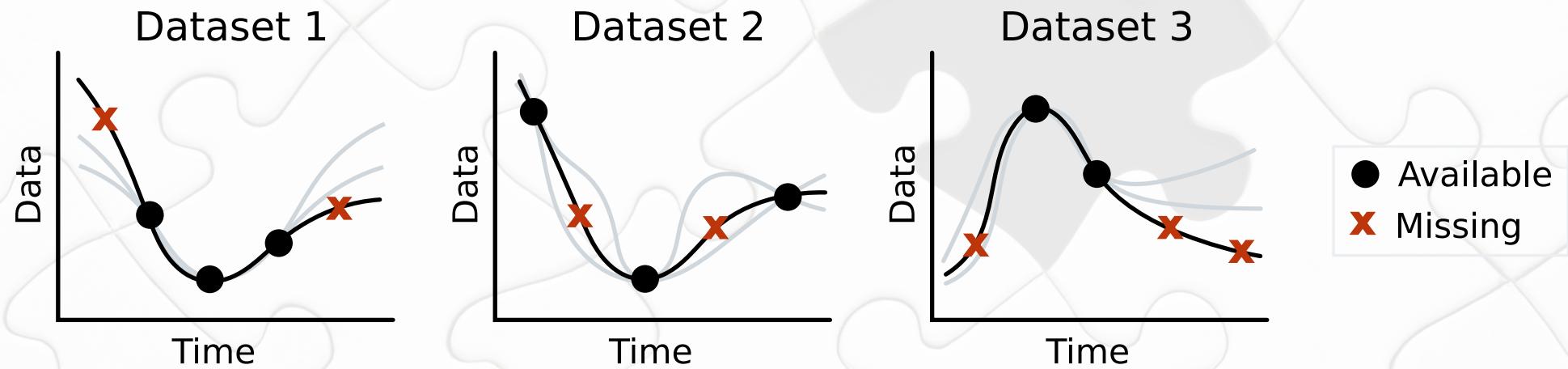
# RESULTS

- ✓ NN fits the posterior well
- ✓ after training once ( $1 - 2h$ ), the optimization is very fast ( $s - min$ )
- ✓ even uncertainty analysis easily possible
- ✓ can handle multiple experiment conditions



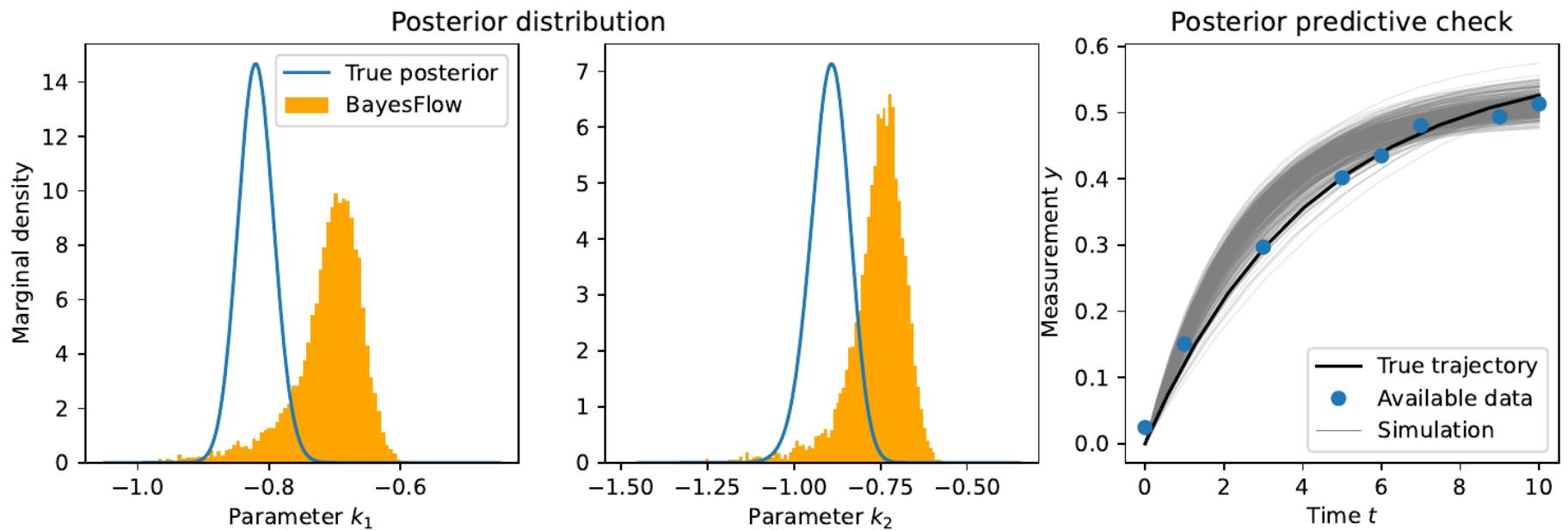


# HOW TO HANDLE MISSING DATA IN AMORTIZED INFERENCE?

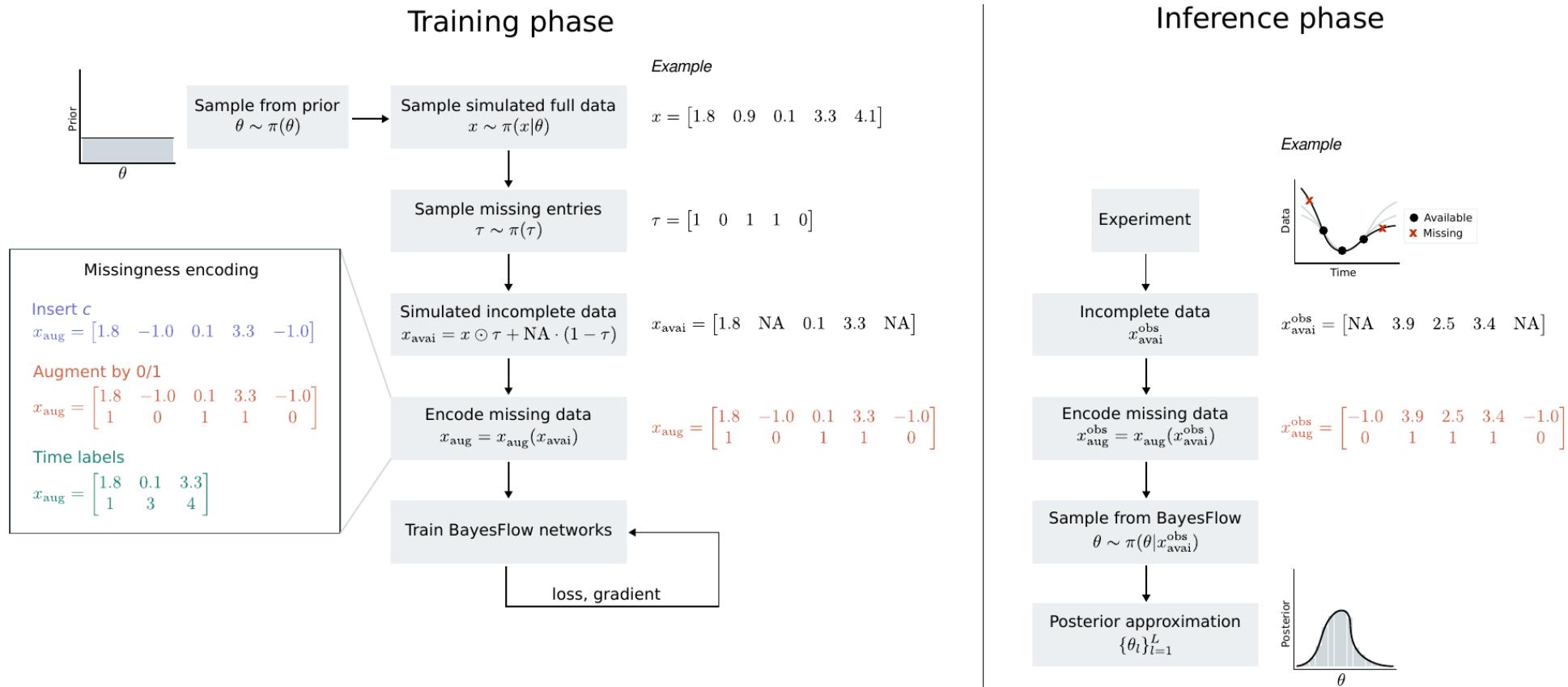


# HOW TO HANDLE MISSING DATA IN AMORTIZED INFERENCE?

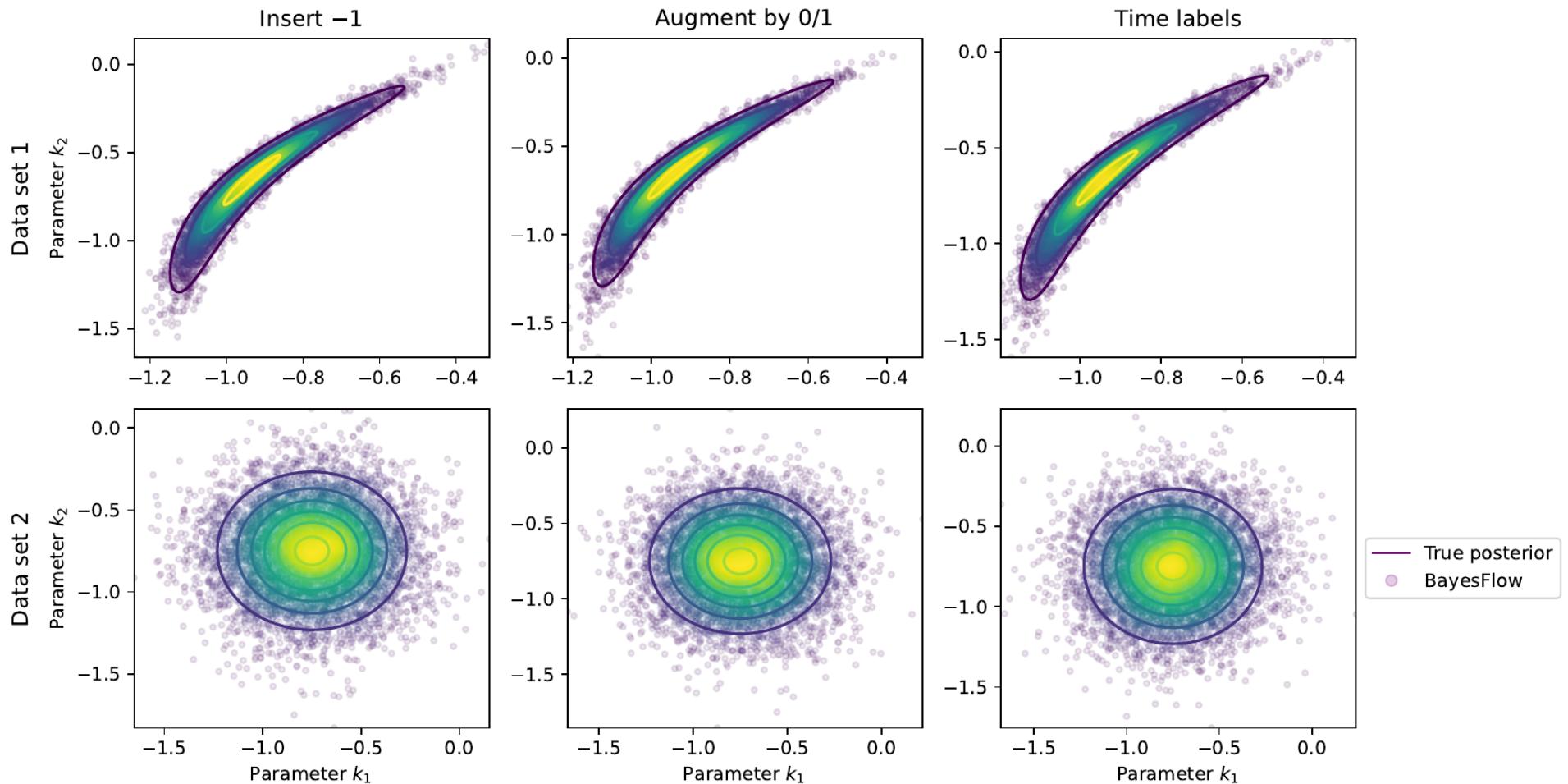
# PROBLEM: BAYESFLOW CANNOT INTERPRET THE DATA



# ENCODE MISSING DATA

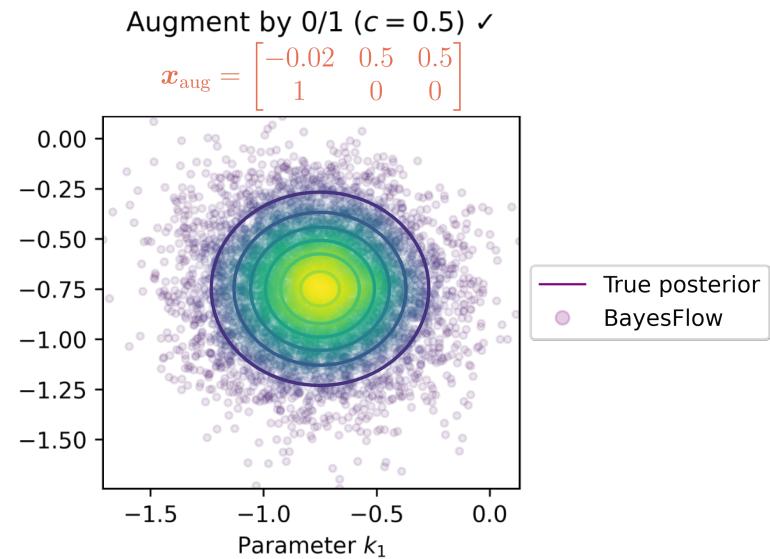
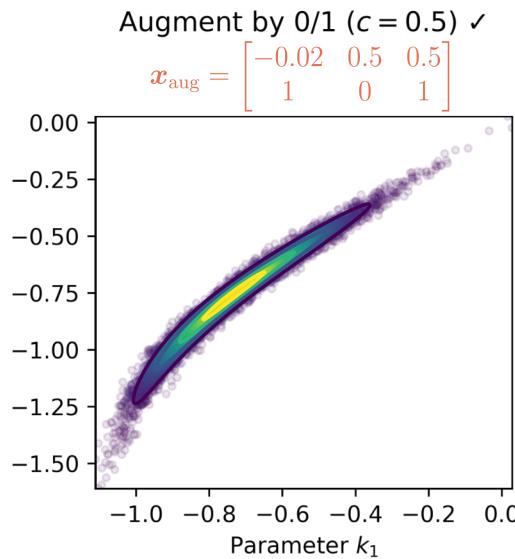
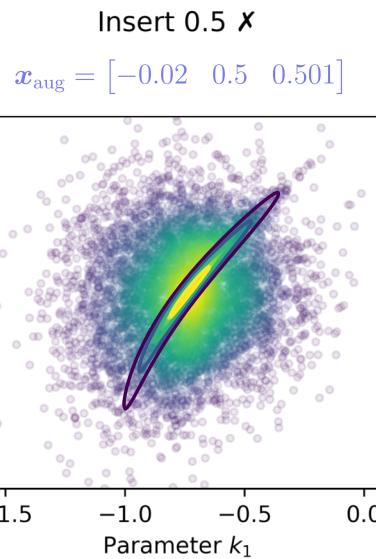


# All approaches perform well on simple test problem

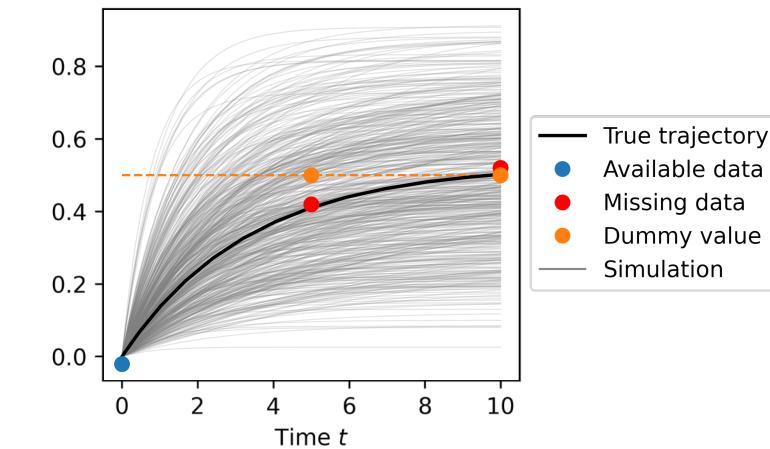
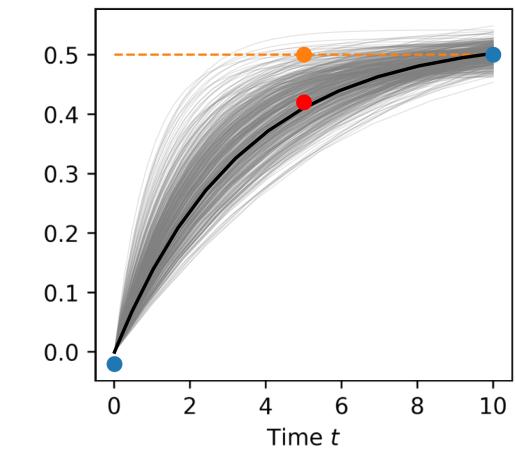
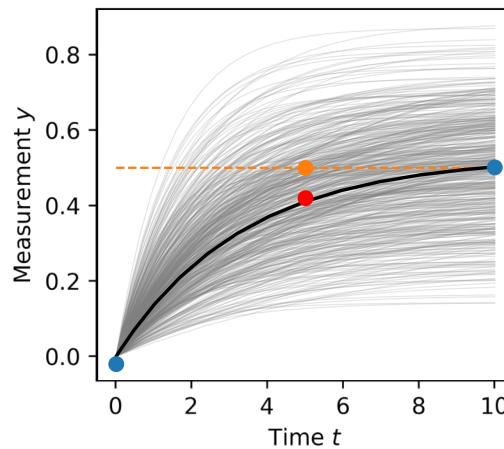


# Binary indicator augmentation more robust for ambiguous fill-in values

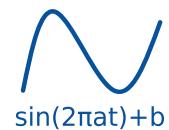
Posterior distribution

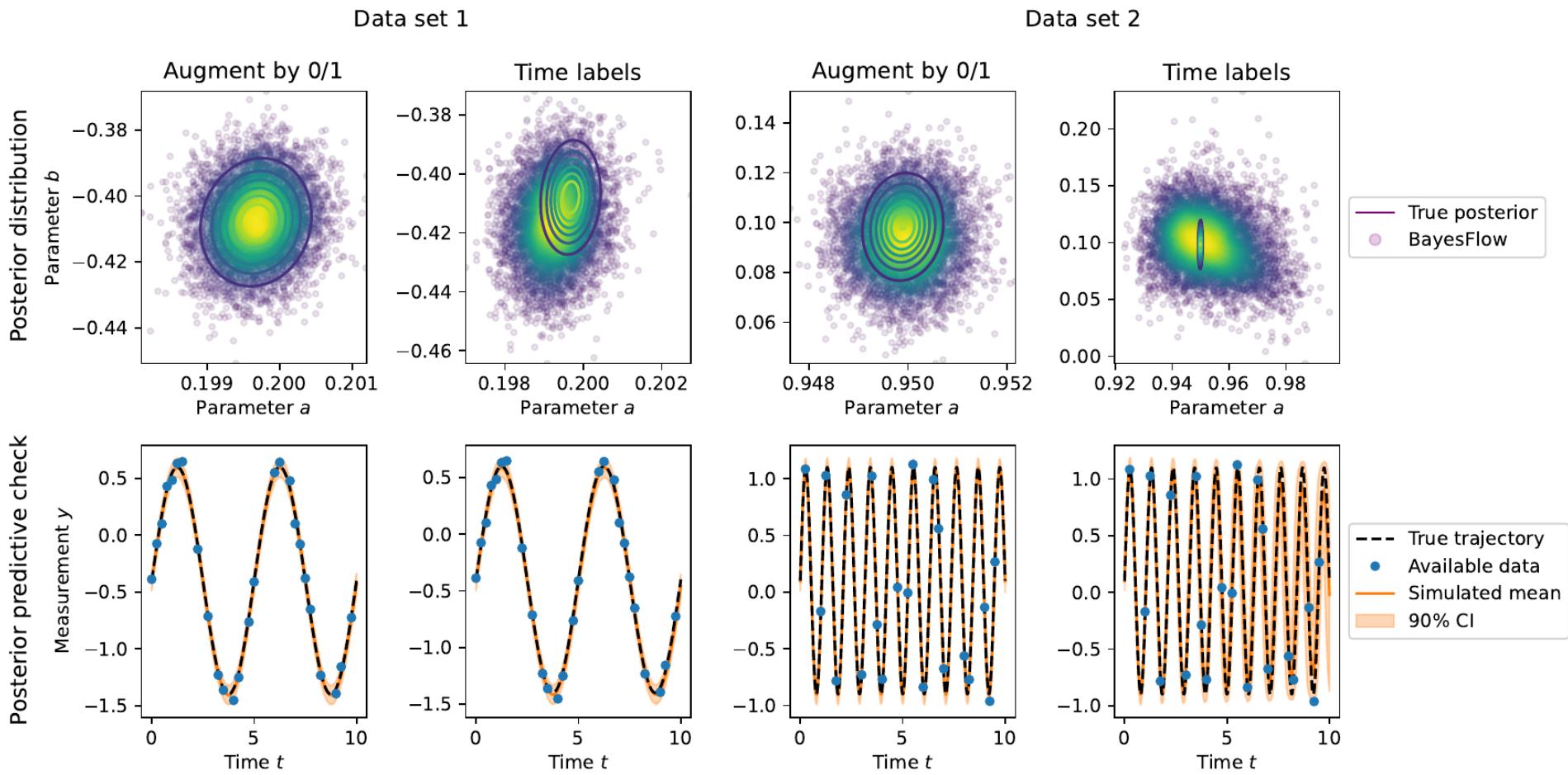


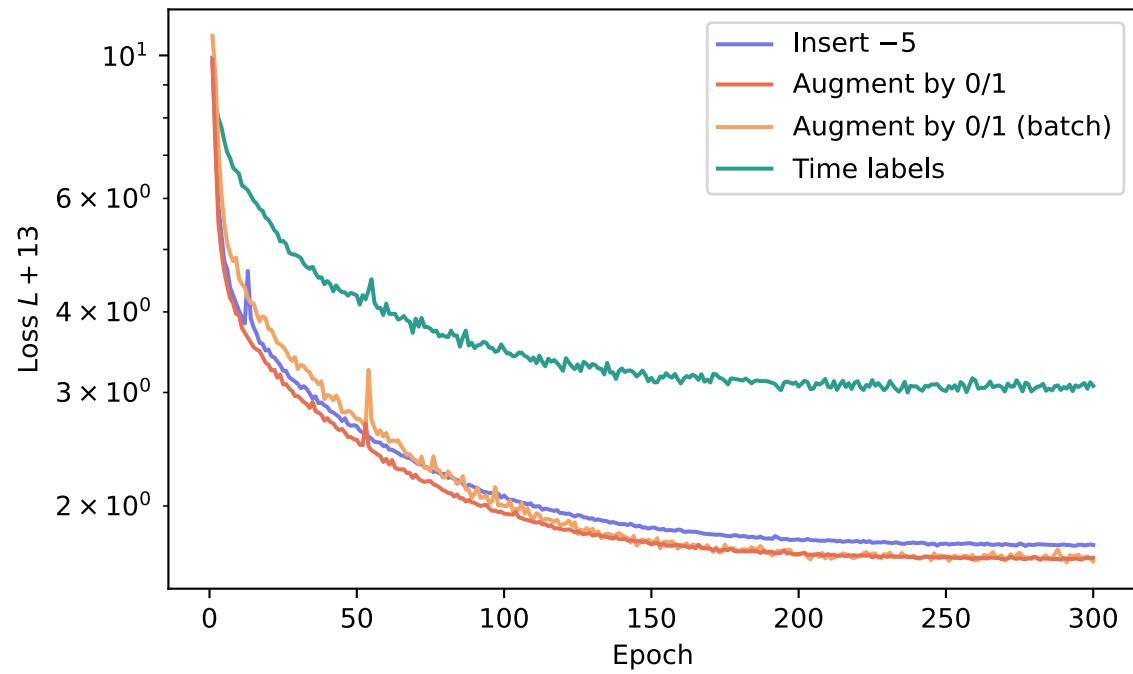
Posterior predictive check



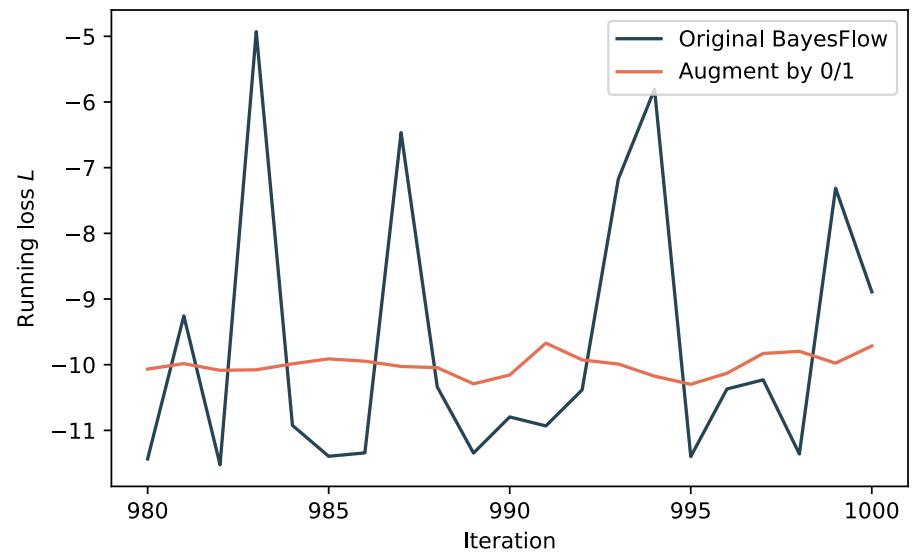
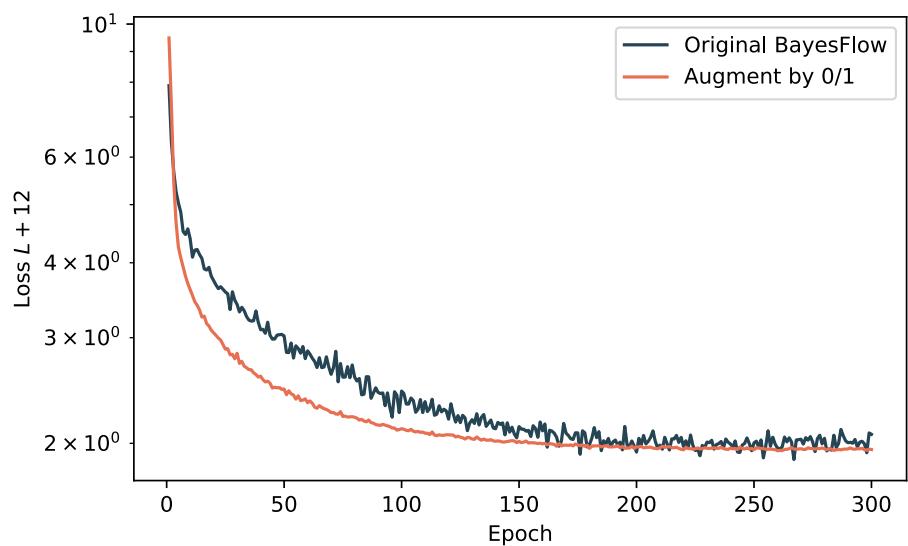
# Time labels encoding performs not robustly in case of oscillatory data


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



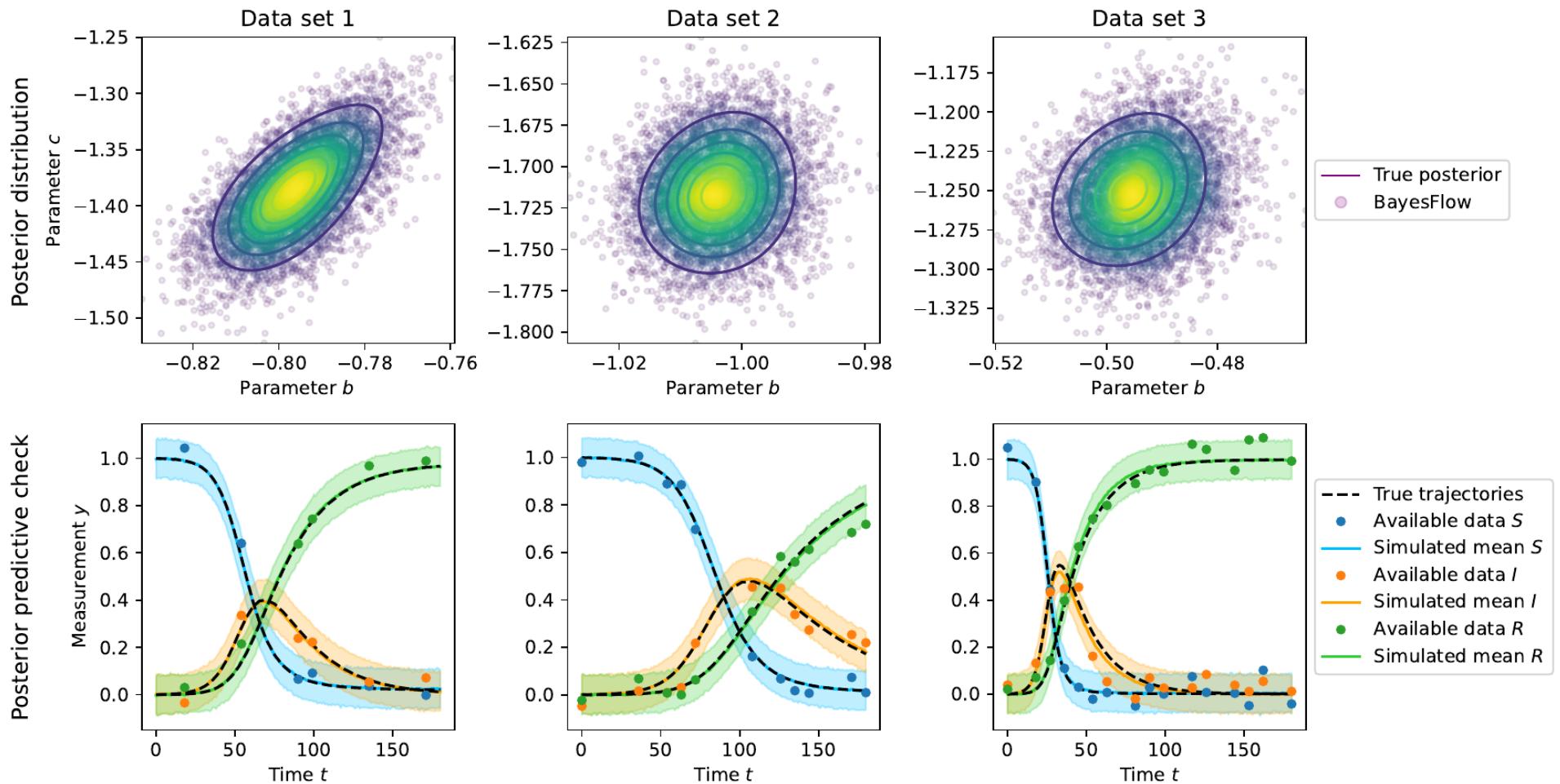


## Variable dataset size as a special case of missing data

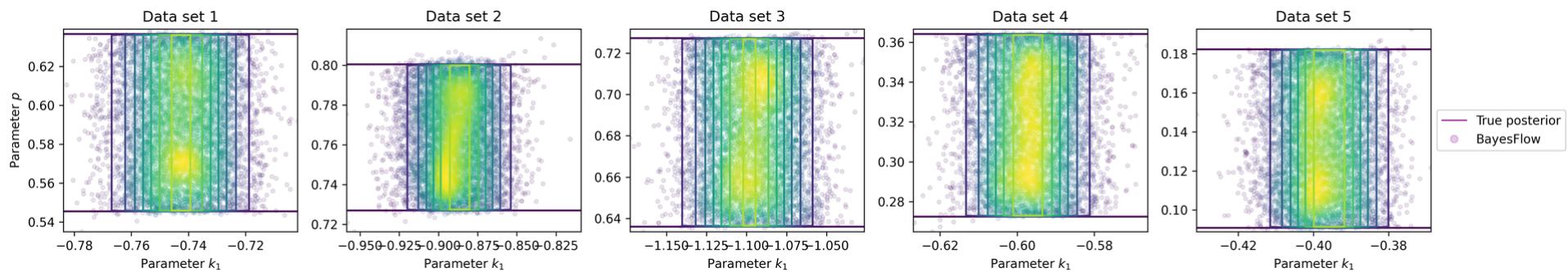


Augment by 0/1 improves performance due to better cost function approximation with individual-specific missingness

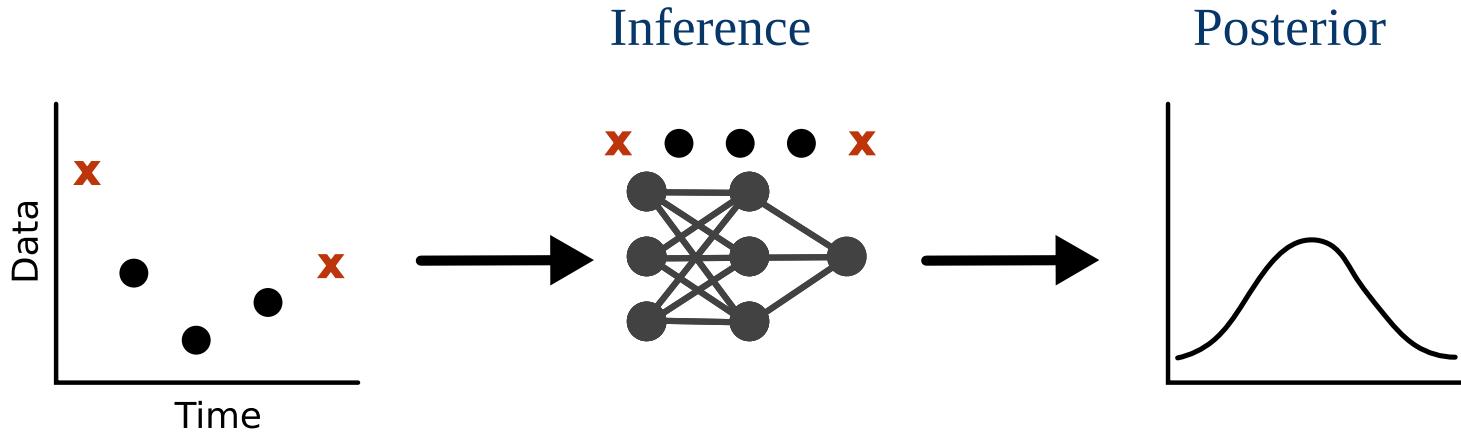
# Scales to more complex inference problems



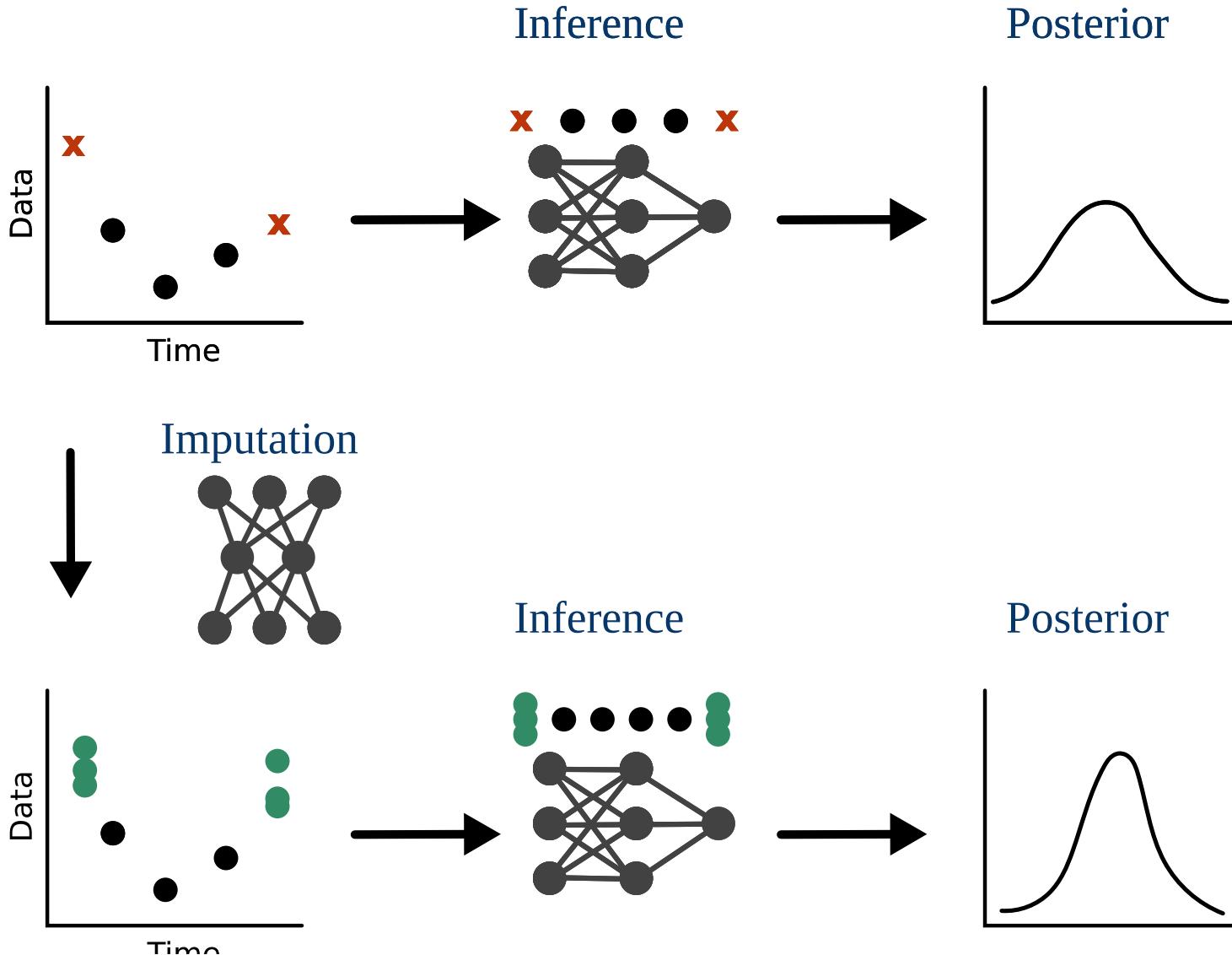
## Able to unravel parameter-dependent missingness



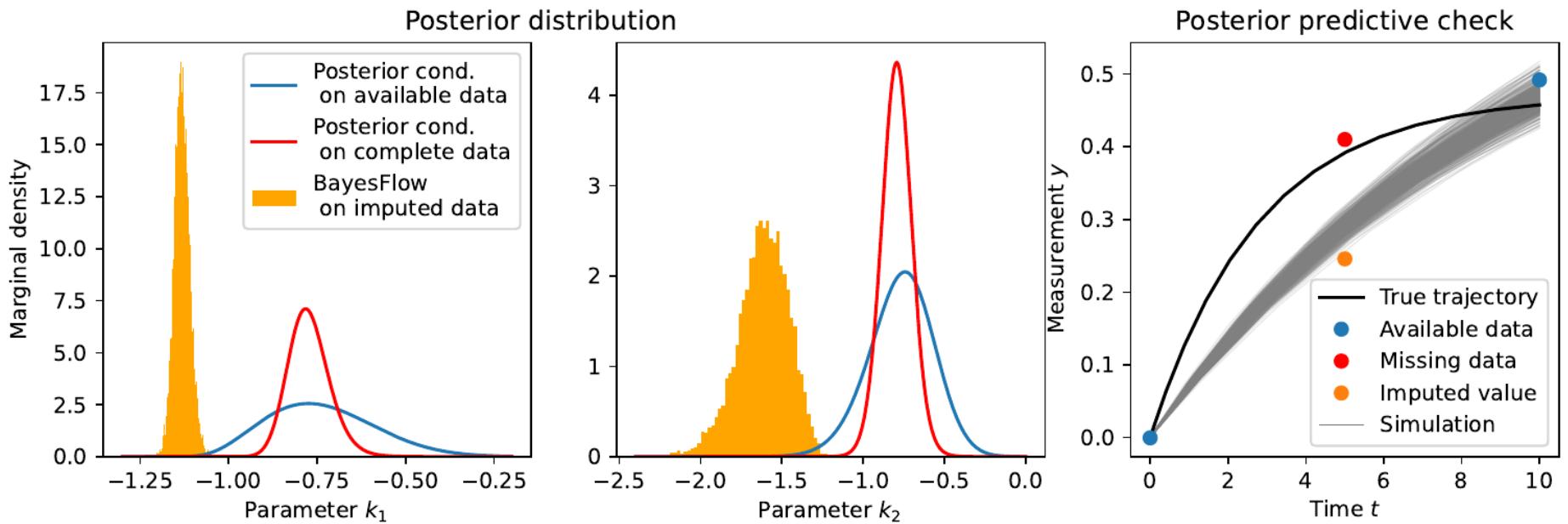
# CAN WE JUST IMPUTE MISSING VALUES?



# 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  $\pi(\bar{x}|x)\pi(x) = \pi(\bar{x}, x)$ .

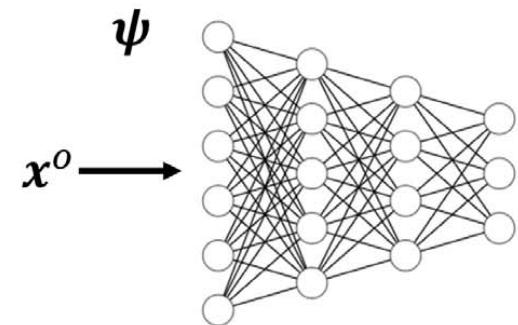
However, if we integrate out over all possible realizations of full data, we obtain  $\int \pi(\theta|\bar{x})\pi(\bar{x}|x)d\bar{x} = \pi(\theta|x)$  (or similarly  $\pi(\theta|x, \tau)$ ).

**TLDR:** When doing uncertainty quantification properly, we just recover the same posterior.

**THANKS! QUESTIONS?**

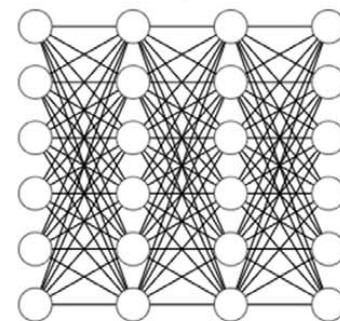
# BAYESFLOW

Summary network



Invertible network

$\phi$



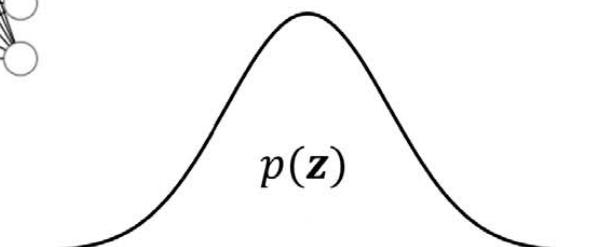
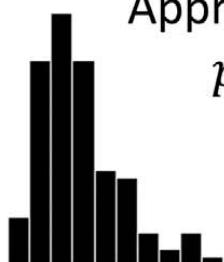
Sampling

$$\mathbf{z} \sim \mathcal{N}_d(\mathbf{0}, \mathbf{I})$$

$\theta$

Approximate posterior

$$p_\phi(\theta | x = \tilde{x}^o)$$



from: Radev et al, IEEE Transactions on Neural Networks 2020

## THE PROBLEM

- Classical simulation-based inference is case-based + slow +  $\varepsilon$ -approximate
- What if we want to fit the same model to multiple datasets?

## THE IDEA

- Learn a global estimator for the probabilistic mapping  $y \mapsto \theta$  via cINNs
- Once trained, amortize inference on arbitrarily many datasets
- Embed data via summary statistics model

# GENERATIVE MODELS

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

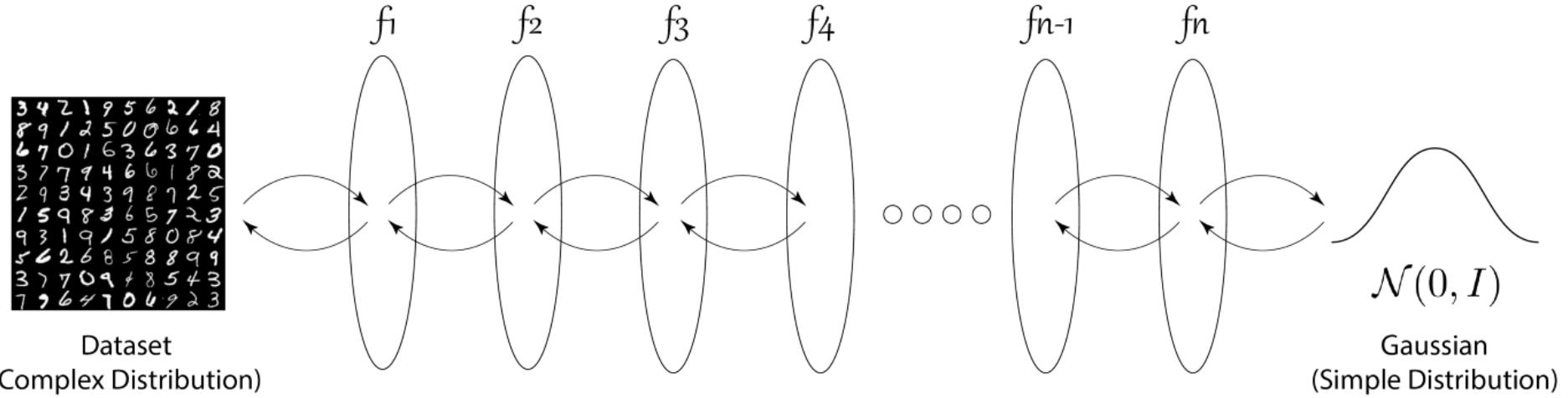


from: Kingma et al, NeurIPS 2019

e.g.: GANs, VAEs, Flows

# NORMALIZING FLOWS

generative models based on an invertible transformation



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

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

# THE PROBLEM

- forward model  $x_i \sim p(x|\theta) \Leftrightarrow x_i = g(\theta, \xi_i)$  with  $\xi_i \sim p(\xi)$
- Bayesian inference  $p(\theta|x_{1:N}) \propto p(x_{1:N}|\theta)p(\theta)$
- aim: train an invertible neural network that approximates the true posterior  $p_\phi(\theta|x) \approx p(\theta|x) \forall \theta, x$

# THE METHOD

Goal: Approximate the true posterior  $p_\phi(\theta|x) \approx p(\theta|x) \forall \theta, x.$

Parameterize  $p_\phi$  in terms of a cINN given via a bijective  $f_\phi : \mathbb{R}^D \rightarrow \mathbb{R}^D$ ,  $\theta \mapsto z$ , which implements a normalizing flow between  $\theta$  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  $\forall x$ , giving the objective ...

# 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:

# SUMMARY STATISTICS LEARNING

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

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

Learning 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 Bayesian inference

