

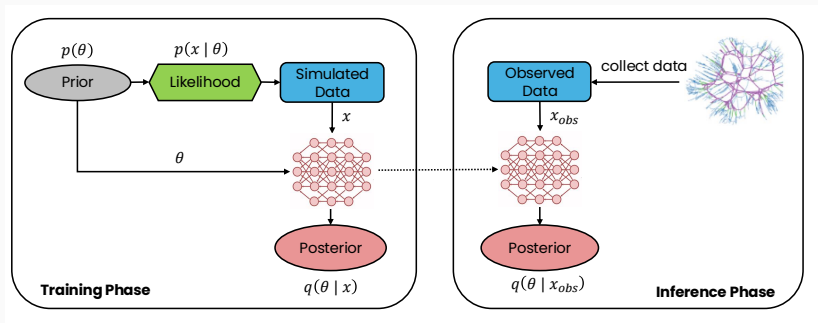
# Robust Amortized Bayesian Inference with Self-Consistency Losses on Unlabeled Data

<https://arxiv.org/abs/2501.13483>

---

Aayush Mishra, Daniel Habermann, Marvin Schmitt, Stefan T. Radev,  
Paul-Christian Bürkner

# Amortized Bayesian Inference (ABI)



Parameters  $\theta$ , data  $x$ , neural approximator  $q(\theta | x)$  for the posterior  $p(\theta | x)$

I care about *the* (analytic) posterior  $p(\theta | x)$  corresponding to my specified probabilistic model  $p(x, \theta) = p(x | \theta) p(\theta)$

# Standard neural posterior estimation (NPE)

General form of (standard) NPE losses in SBI:

$$\text{NPELoss}(q) = \mathbb{E}_{(\theta, x) \sim p(\theta, x)} [S(q(\cdot \mid x), \theta)]$$

For normalizing flows with invertible neural networks:

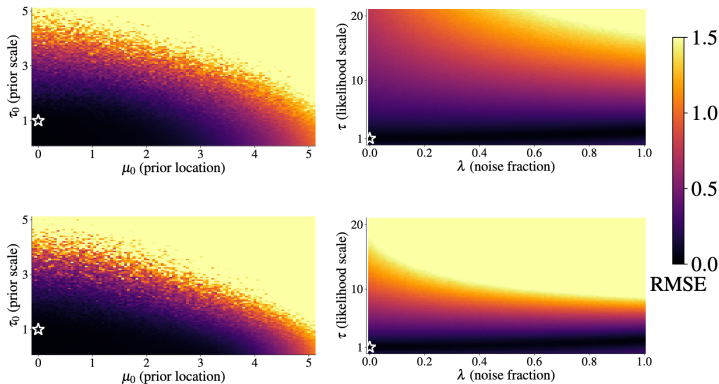
$$\text{NPELoss}(q) = \mathbb{E}_{(\theta, x) \sim p(\theta, x)} [-\log q(\theta \mid x)]$$

# Standard NPE on misspecified models fails

Summary Network  
minimal  
overcomplete

## Model Misspecification

Prior Simulator & noise



Source: <https://arxiv.org/abs/2406.03154>

# Bayesian Self-consistency

For any set of parameter values  $\theta^{(1)}, \dots, \theta^{(L)}$ , the following holds:

$$p(x) = \frac{p(x \mid \theta^{(1)}) p(\theta^{(1)})}{p(\theta^{(1)} \mid x)} = \dots = \frac{p(x \mid \theta^{(L)}) p(\theta^{(L)})}{p(\theta^{(L)} \mid x)}.$$

This implies that the variance of the log-ratios must be zero:

$$\text{Var}_{l=1}^L \left[ \log \left( \frac{p(x \mid \theta^{(l)}) p(\theta^{(l)})}{p(\theta^{(l)} \mid x)} \right) \right] = 0$$

Our initial paper on Bayesian self-consistency:

<https://arxiv.org/abs/2310.04395>

# Bayesian self-consistency loss

Replace the true posterior  $p(\theta \mid x)$  with the neural approximate posterior  $q(\theta \mid x)$ .

For any (**unlabeled**) dataset  $x^*$  and any parameter generating distribution  $\tilde{p}(\theta)$ , we define:

$$\text{SCLoss}(\mathbf{q}) = \text{Var}_{\theta \sim \tilde{p}(\theta)} [\log p(x^* \mid \theta) + \log p(\theta) - \log q(\theta \mid x^*)]$$

$\Rightarrow$  We can use **real-world data** as  $x^*$  to train our SC loss!

The SC-Loss alone doesn't work well most of the time so we combine it with the standard NPE loss:

$$\text{SemiSupervisedLoss}(\mathbf{q}) = \text{NPELoss}(\mathbf{q}) + \lambda \cdot \text{SCLoss}(\mathbf{q}).$$

# Bayesian Self-Consistency losses a strictly proper

Let  $C$  be a score that is globally minimized if and only if its functional argument is constant across the support of the posterior  $p(\theta | x)$  almost everywhere. Then,  $C$  applied to the Bayesian self-consistency ratio with known likelihood

$$C \left( \frac{p(x | \theta) p(\theta)}{q(\theta | x)} \right)$$

is a strictly proper loss: It is globally minimized if and only if  $q(\theta | x) = p(\theta | x)$  almost everywhere.

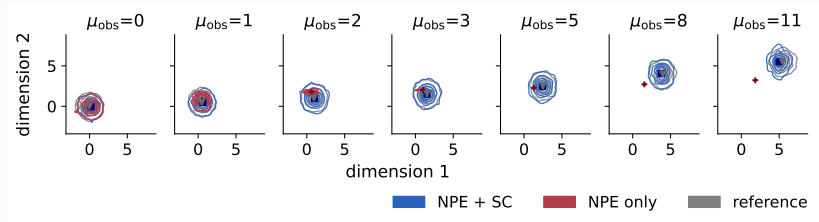
This implies that also the semi-supervised loss is strictly proper.

# Case Study 1: Multivariate normal model

$$\theta \sim \text{Normal}(\mu_{\text{prior}}, I_D), \quad x \sim \text{Normal}(\theta, I_D)$$

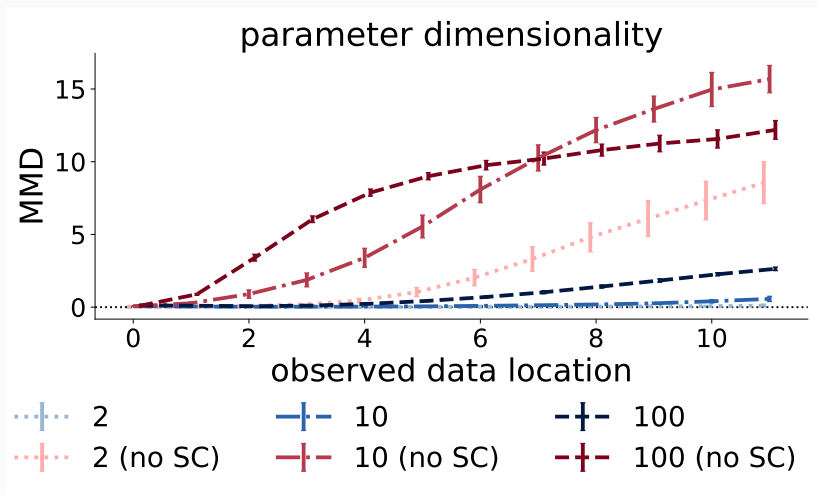
- For the NPE loss, we simulate from the model with  $\mu_{\text{prior}} = 0$
- For the SC loss, we simulate **few unlabeled datasets** from the model with  $\mu_{\text{prior}} = 2$

Illustrative results:

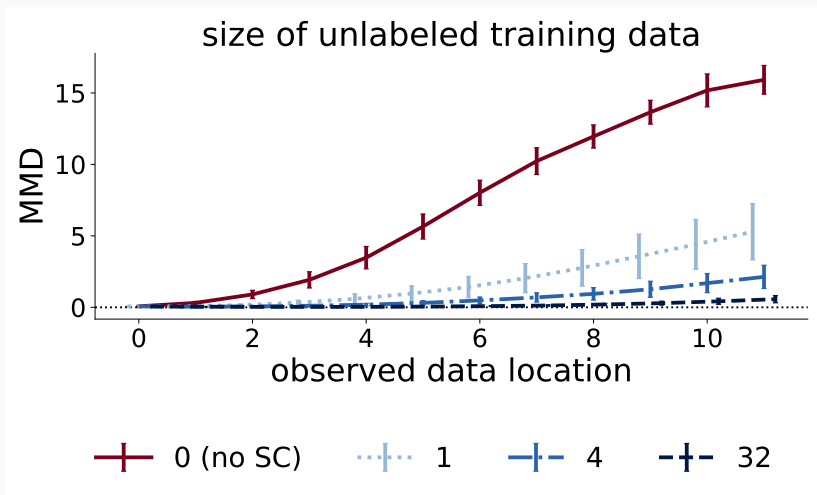




## Case Study 1: More Results



## Case Study 1: More Results



## Case Study 2: Time Series of Air Traffic data

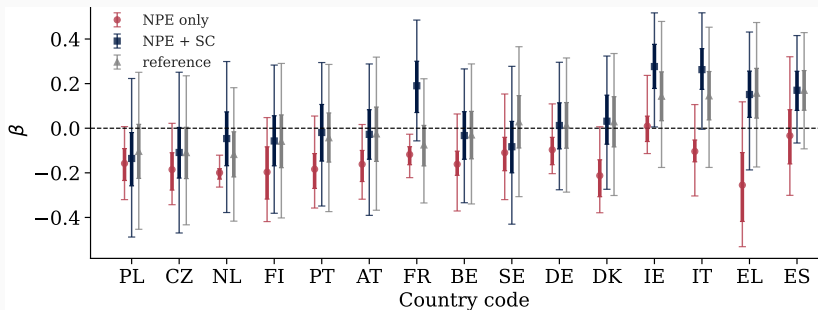
Predicting the change in air traffic for different European countries

$$y_{j,t+1} \sim \text{Normal}(\alpha_j + \beta_j y_{j,t} + \dots, \sigma_j)$$

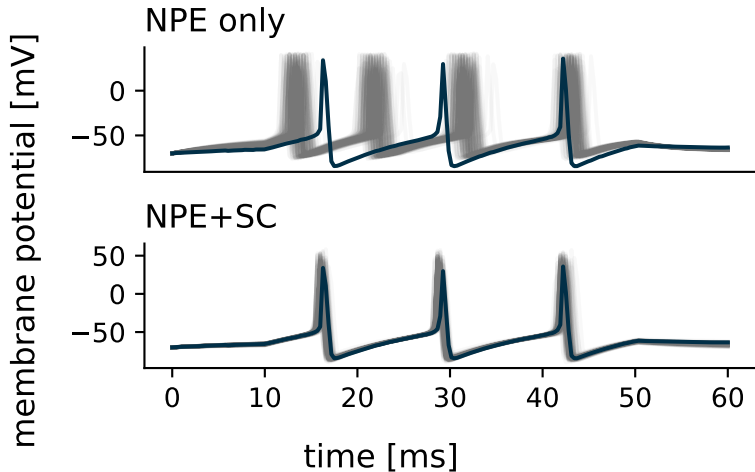
- $y_{j,t}$  number of passengers for country  $j$  at year  $t$
- $\alpha_j$  intercept parameter
- $\beta_j$  auto-correlation parameter
- $\sigma_j$  residual standard deviation

We have data of 15 countries, 4 of which are used as training data in our SC loss.

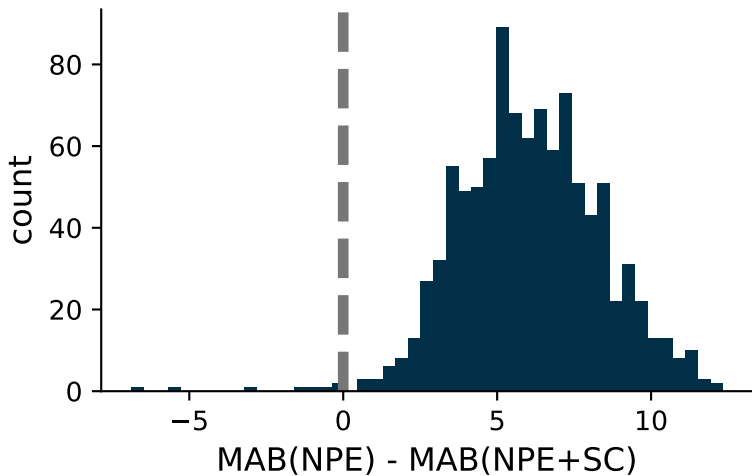
## Case Study 2: Results



## Case Study 3: Hodgkin-Huxley model of neuron activation



### Case Study 3: More results



# Intermediate Summary

- The SC loss can strongly improve robustness to model misspecification
- The SC loss requires no data labels so we may even use **real-world data** for training
- The SC loss is strictly proper so it has the same target (the true posterior) as the NPE loss
- Challenge 1: The SC loss requires a known or estimated likelihood density: stronger robustness in the known case
- Challenge 2: We need neural approximators that have fast density evaluation

# Our current reserach on SC losses

- Model comparison: SC works great when done correctly (<https://arxiv.org/abs/2508.20614>)
- Continual learning: SC losses may lead to catatrophic forgetting if applied alone but we can mitigate that
- Unknown likelihood densities: We have some promising ideas how to adjust SC losses and training



# What is the target of inference?

- **Target 1:** The analytic posterior  $p(\theta \mid x_{\text{obs}}) \propto p(x_{\text{obs}} \mid \theta) p(\theta)$  of the assumed probabilistic model given the observed data  $x_{\text{obs}}$ .
- **Target 2:** A posterior  $p(\theta \mid \tilde{x}_{\text{obs}}) \propto p(\tilde{x}_{\text{obs}} \mid \theta) p(\theta)$  of the assumed probabilistic model given *adjusted data*  $\tilde{x}_{\text{obs}}$ .
  - Equivalent: A posterior given an *adjusted likelihood*
- **Target 3:** A posterior  $\tilde{p}(\theta \mid x_{\text{obs}}) \propto p(x_{\text{obs}} \mid \theta) \tilde{p}(\theta)$  from an *adjusted prior*  $\tilde{p}(\theta)$  given the observed data  $x_{\text{obs}}$ .

Source: <https://arxiv.org/abs/2502.04949>

# Ways to achieve robust/trustworthy inference

- Unsupervised likelihood-based losses
  - Example: SC losses
  - Aims at Target 1
  - Drawback: Requires the likelihood density
- Unsupervised domain adaptation
  - Example: Minimize the distance of simulated and real-world data in the summary space
  - Aims at Target 2
  - Drawback: Adjusts the target posterior implicitly
- Supervised real-world data calibration
  - Calibrate posterior on **labeled** real-world data
  - Aims at Target 2 (I think)
  - Drawback: Requires real-world data with labels

- Mathematical theory of SC losses
  - In the style of <https://arxiv.org/abs/2411.12068>
- ABI robustness methods
- Architecture improvements
- Software: BayesFlow
- Applications of ABI