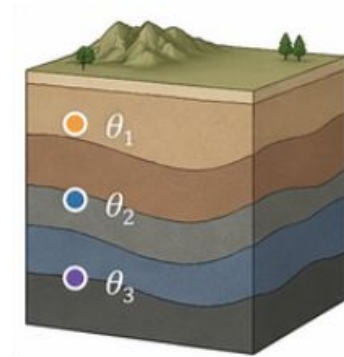
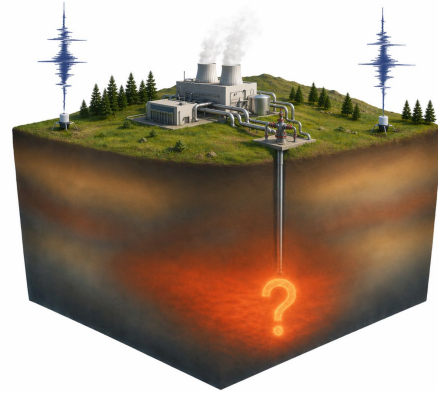


Learning Earth Structure

Generative Posterior Inference for Seismic Inversion under Model
Misspecification

Context and motivation

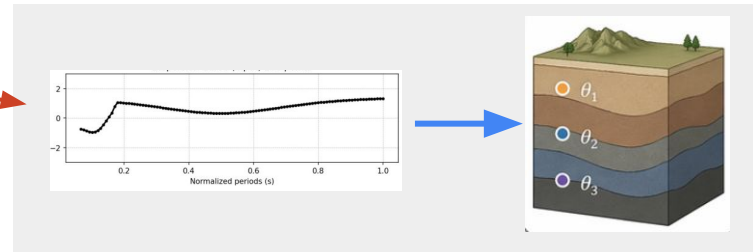
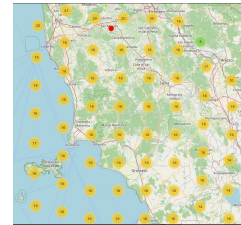
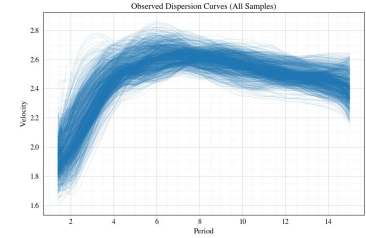
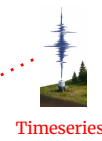
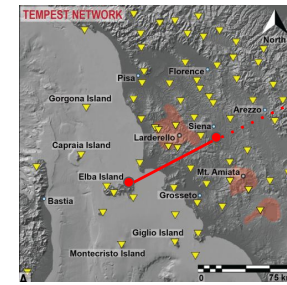
- **Geothermal exploration** requires imaging the subsurface.
- But exploiting it safely requires **knowing the underground structure** (subsurface).
- Classical subsurface imaging remains expensive, especially at the inversion stage.



The 2D-1D ANT Pipeline

2D-1D Inversion

1. **Correlation of ambient noise** : cross-correlation between stations ;
2. **Extraction of dispersion curves** : how groups of frequencies propagate between a pair of stations ;
3. **Reconstruct a 2-D map** : extract local dispersion curves at each point (x,y) ;
4. **1D-inversion** : local curves can be inverted independently to recover 1-D velocity profiles ;



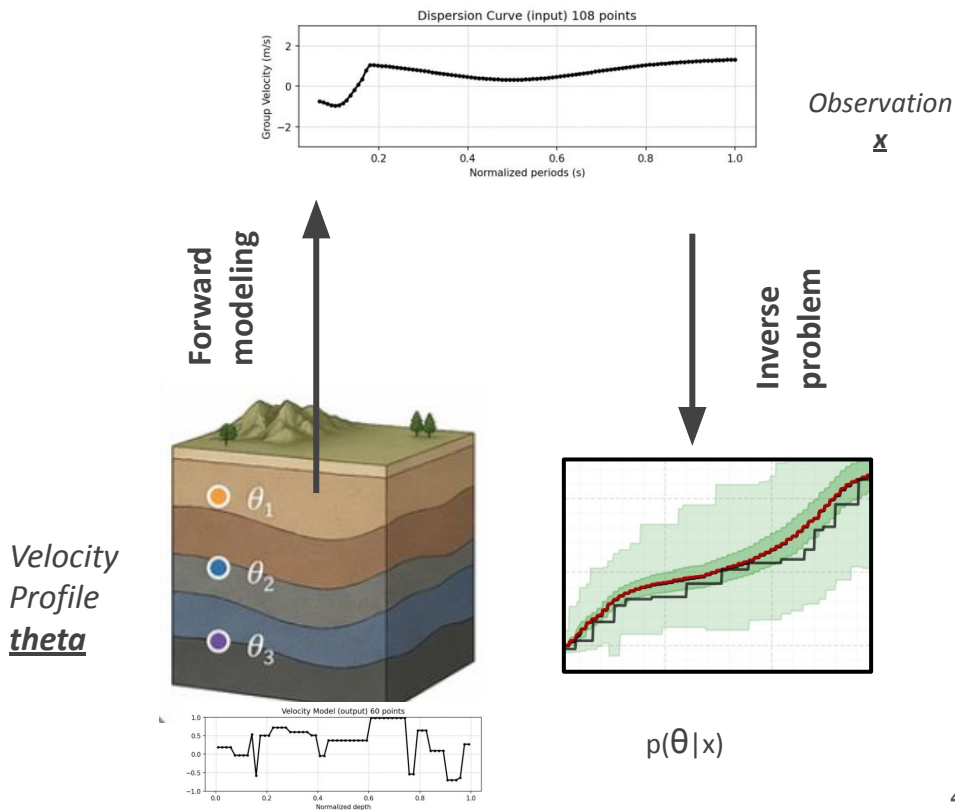
1D inversion of dispersion curves

Forward Modeling

- **Input:** 1D layered model of the subsurface
- **Solves:** dispersion equation derived from elastic wave equations
- **Output:** group velocity curve $v_g(p)$ as a function of *period* p

Inverse Problem

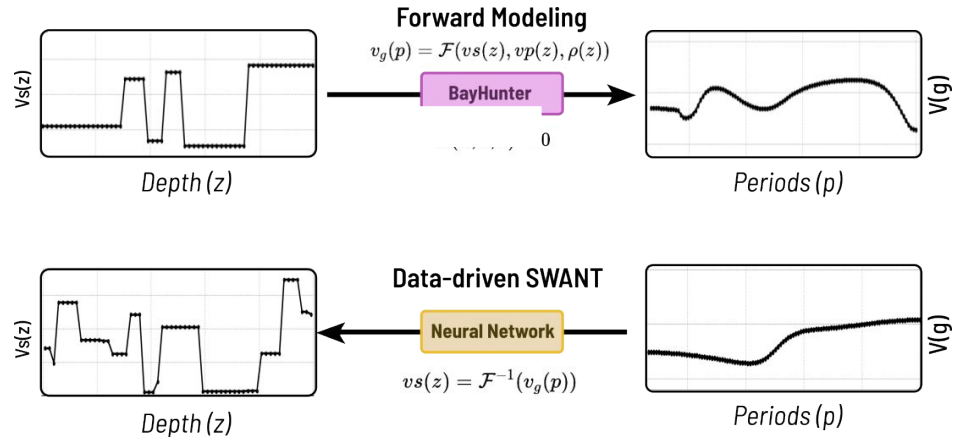
- **Ill-posed** (non-unique solution) and **trans-dimensional**
- **Learn the posterior distribution** $p(\theta|x)$
- **Validation:** do the generated model explain the observation ?



1D inversion of dispersion curves

Forward Modeling

- **Input:** 1D layered model of the subsurface
- **Solves:** dispersion equation derived from elastic wave equations
- **Output:** group velocity curve $v_g(p)$ as a function of *period* p



Inverse Problem

- **Ill-posed** (non-unique solution) and **trans-dimensional**
- **Learn the posterior distribution** $p(\theta|x)$
- **Validation:** do the generated model explain the observation ?

Inversion problem: $v_g(p) \rightarrow v_s(z)$

SBI Formulation

Main points

- We want the **conditional posterior** :

$$p(\theta | x)$$

- Likelihood-based inference is difficult because:
 - forward model is nonlinear / ill-defined
 - simulator is costly
 - number of layers is unknown
- Instead, we sample training pairs:

$$(\theta_i, x_i) \sim p(\theta) p(x | \theta)$$

- Then train a neural posterior estimator (NPE)

$$\hat{p}_\phi(\theta | x) \approx p(\theta | x)$$

- The cost is shifted from per-observation sampling to one offline training stage, enabling fast amortized posterior inference.

$$p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\int p(x|\theta)p(\theta) d\theta}$$

RJ-MCMC (Bayhunter)

- **Goal** : infer subsurface structure from observed data (e.g., dispersion curves, travel times, receiver functions).
- **RJ-MCMC** samples both parameter values (V_s) and model dimension (*number of layers*).
- **Within-model move**: perturb parameters with fixed (k) (e.g., velocity, depth, thickness).
- **Birth move**: add one component (k to $k+1$).
- **Death move**: remove one component (k to $k-1$).
- Move probabilities strongly affect efficiency and mixing.
- Build a Markov chain whose stationary distribution is the posterior.
- **At each iteration** :
 - a. propose a candidate model,
 - b. accept or reject with a probability rule,
 - c. keep the accepted state.
- After burn-in, samples approximate the full posterior uncertainty.
- **Computationally expensive** (many forward solves).
- Start from a prior :
 - a. V_s : [0.5, 4.5]
 - b. Layers : [2, 20]
 - c. V_p/V_s : [1.65, 1.85]

Synthetic Dataset & Prior

Main points

- We generate paired samples (θ, x) with the forward simulator **surf96**
- Observation:

$$x \in \mathbb{R}^{108}$$

- fundamental-mode Rayleigh-wave **group-velocity curve** sampled from **1 to 30 s**.

$$\theta \in \mathbb{R}^{480}$$

- Fixed-depth **velocity map** $V_s(z)$ on $[0, 15]$ km
 - Total dataset : 10m synthetic pairs

Layered prior

- N_l in $[2, 20]$ layers
- Depth range : 0-15 km
- V_s in $[0.5, 4.0]$ km/s
- Increasing V_s with 10% probability of LVZ (low velocity zone)
- Minimum layer thickness
- V_p/V_s in $[1.65, 1.85]$
- Key equation : $x = \mathcal{F}(\theta)$

$$x \in \mathbb{R}^{108} \quad \theta_{\text{layered}} = \{(V_{S,i}, h_i)\}_{i=1}^{N_l} \quad \theta_{\text{velmap}} \in \mathbb{R}^{480}$$

From sampling to learning

Normalizing Flow (MAF, NSF)

- Learn an invertible conditional transport from Gaussian latent space to parameters.
- Conditional mapping:

$$z \sim \mathcal{N}(0, I) \quad \theta = g_\phi(z, x)$$

- Exact density via change of variables:

$$p_\phi(\theta | x) = p_Z(g_\phi^{-1}(\theta, x)) \left| \det \frac{\partial g_\phi^{-1}(\theta, x)}{\partial \theta} \right|$$

- Training

$$\mathcal{L}(\phi) = -\frac{1}{N} \sum_{i=1}^N \log p_\phi(\theta_i | x_i)$$

Flow Matching (UNet-1D, FUNet-1D)

- Start from Gaussian noise and gradually move toward a target parameter sample.

$$\theta_0 \sim \mathcal{N}(0, I)$$

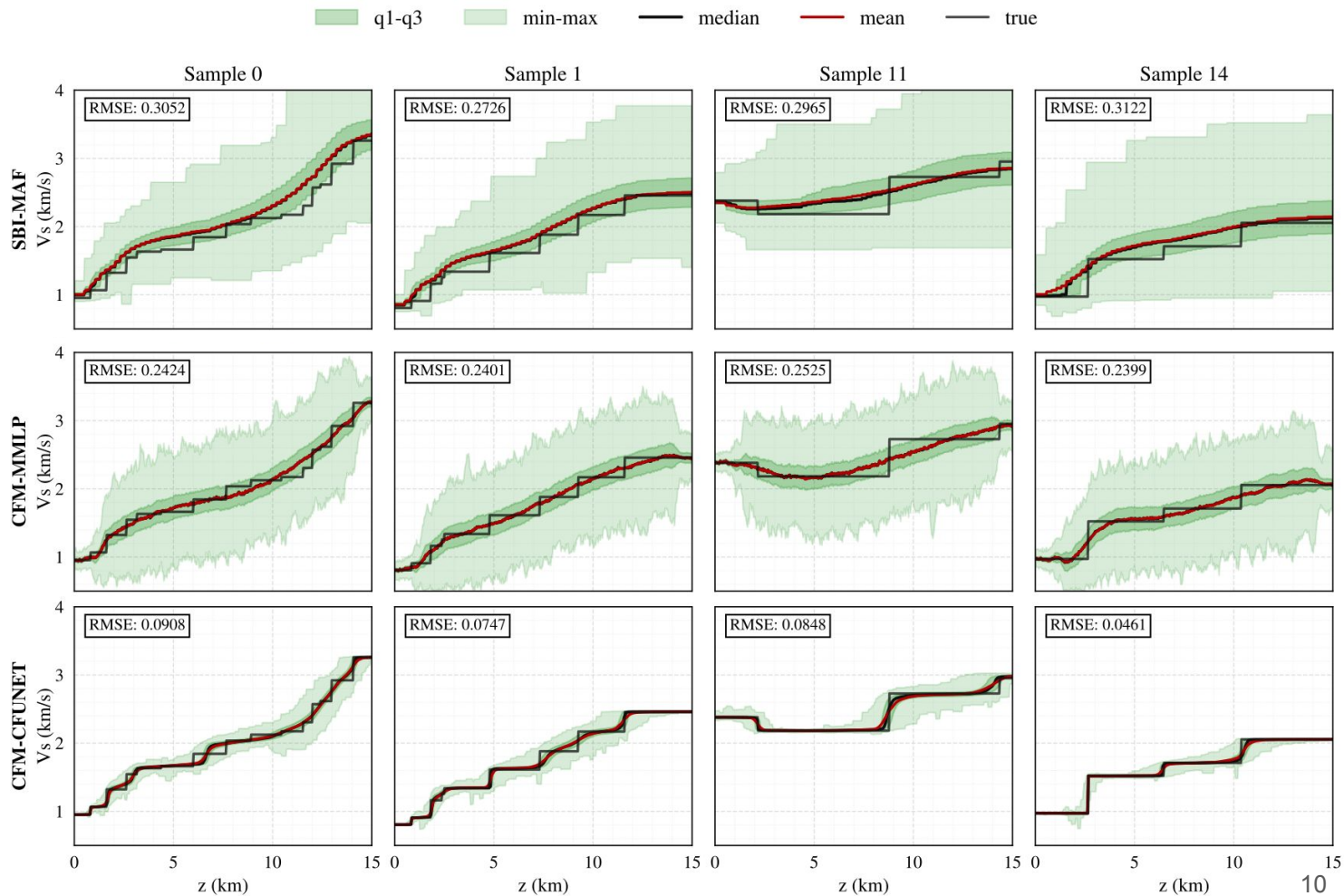
- Learn the vector field

$$\frac{d\theta_t}{dt} = v_\phi(\theta_t, x, t)$$

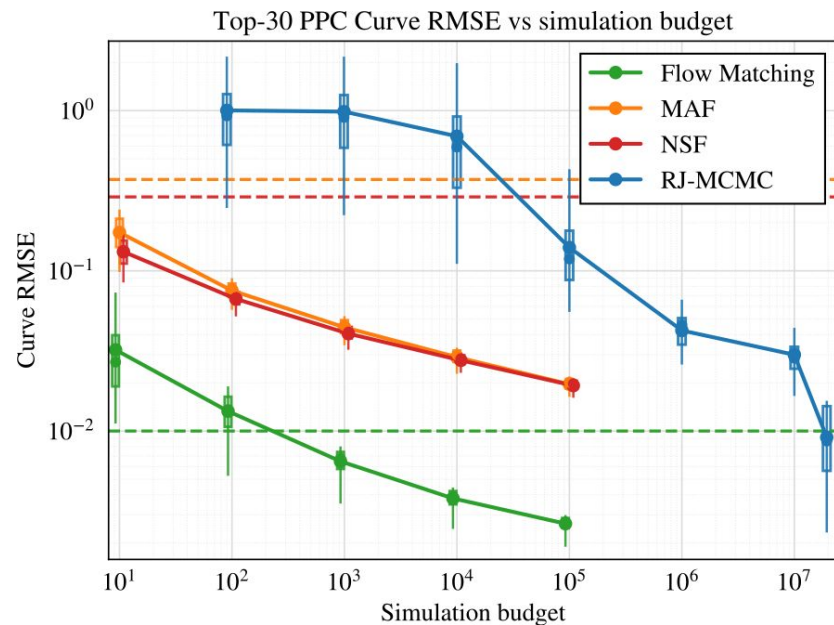
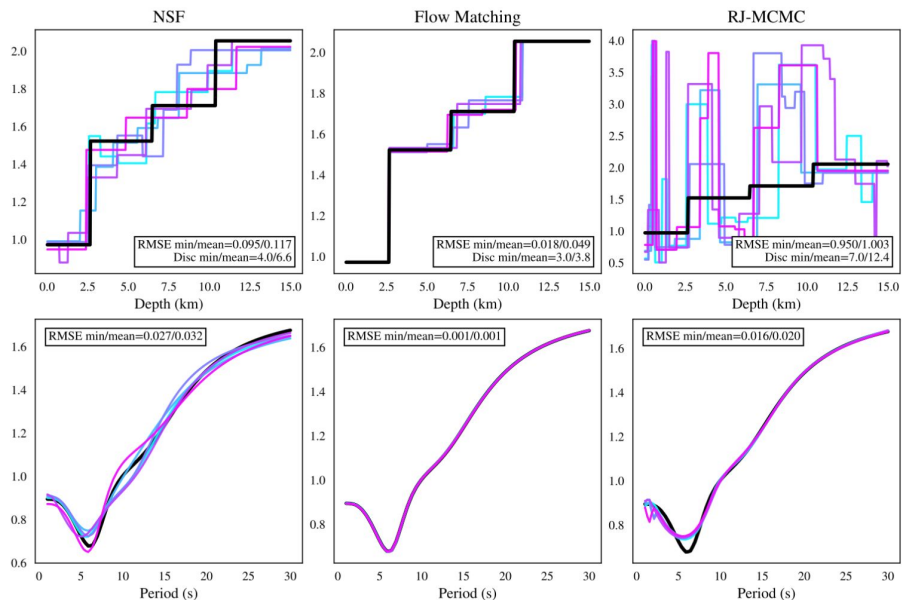
- Train by matching this predicted motion to the true transport direction.

$$\theta_t = (1 - t)\theta_0 + t\theta \quad \mathcal{L}(\phi) = \mathbb{E} \left[\|v_\phi(\theta_t, x, t) - (\theta - \theta_0)\|^2 \right]$$

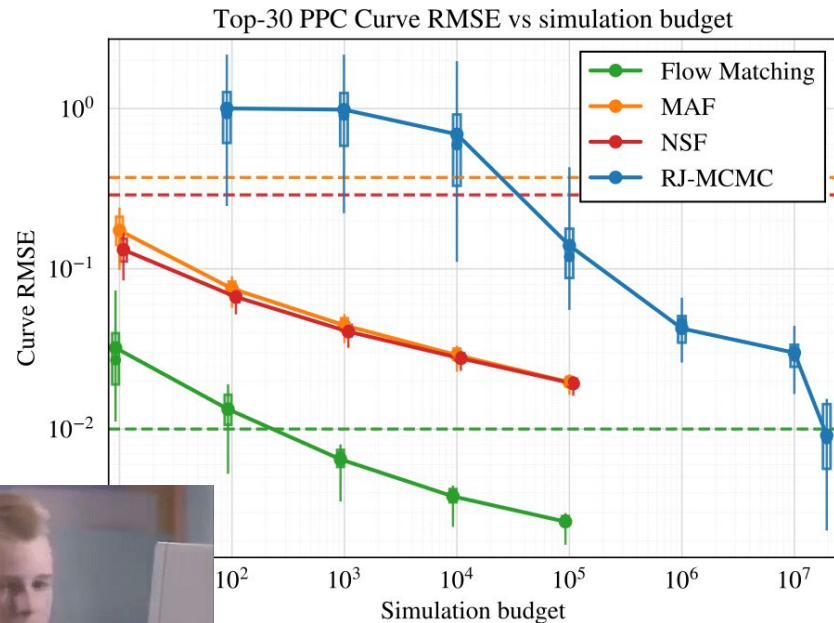
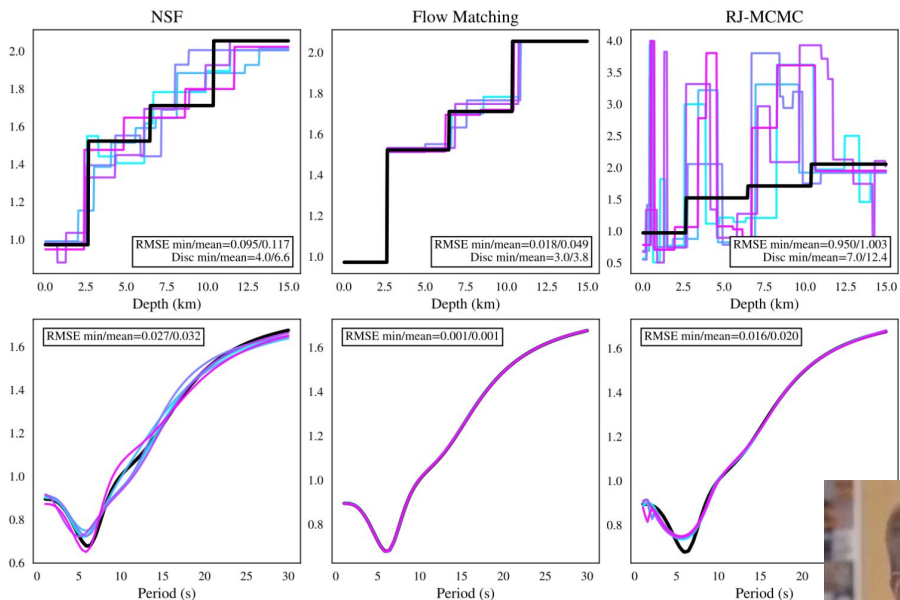
Visualization of the posterior



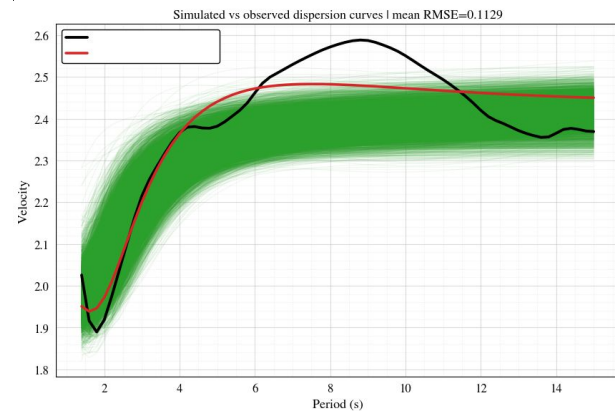
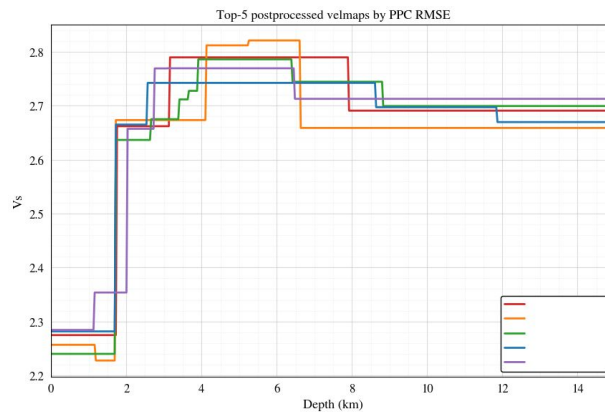
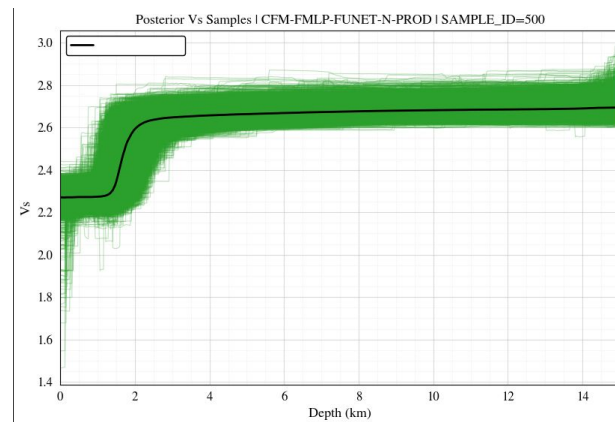
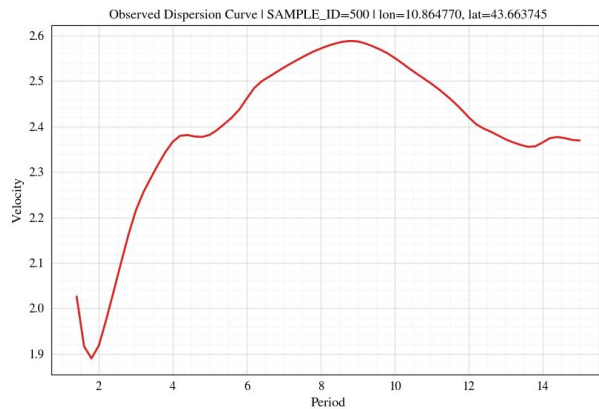
Efficiency and Computational Budget



Efficiency and Computational Budget

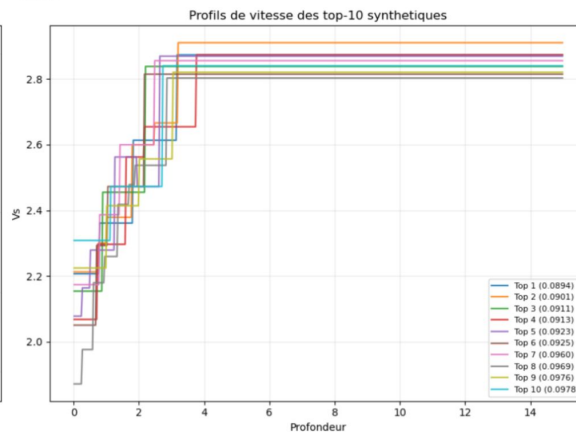
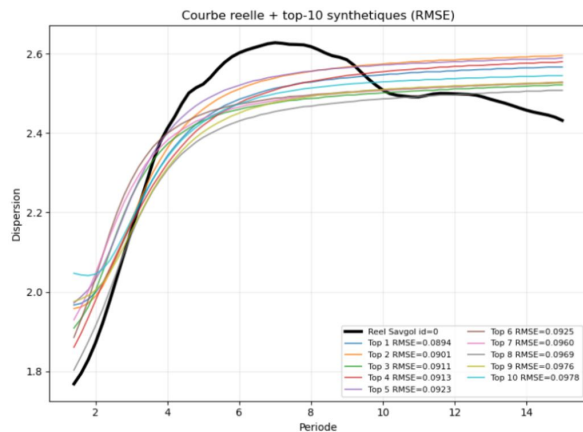
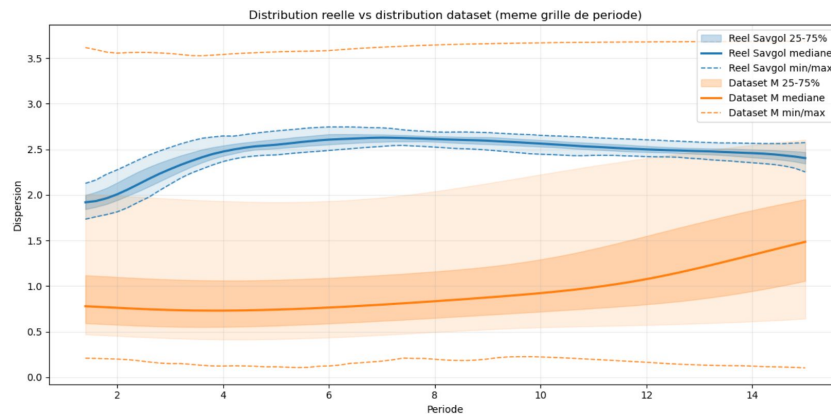


But...



Model misspecification

Unrealistic prior



Model Misspecification

A **computer simulator** with parameter $\theta \in \Theta \subset \mathbb{R}^d$ defines implicitly a parametric family of distributions

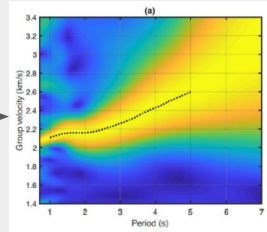
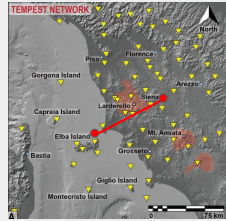
$$\mathbb{P}_\theta := \{p(\cdot | \theta); \theta \in \Theta\}.$$

Running the simulator on θ by sampling from the likelihood $x \sim p(x|\theta)$. We have access to one or more observations y_i drawn from an unknown data-generating process p^* .

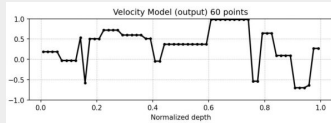
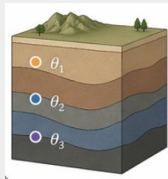
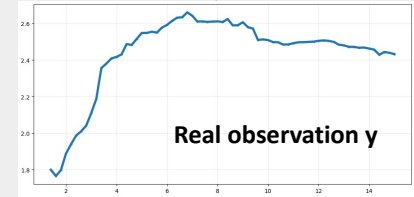
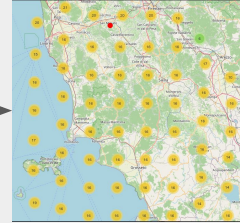
- Well-specified scenario: $p^* \in \mathbb{P}_\theta$ with $p^* = p(\cdot | \theta^*)$ for some θ^*
- Misspecified scenario: $p^* \notin \mathbb{P}_\theta$ then \mathbb{P}_θ is misspecified with y

$\mathcal{X} := \left\{ x : \int p(x | \theta)\pi(\theta) d\theta > 0 \right\}$ is the **observational support** of the model

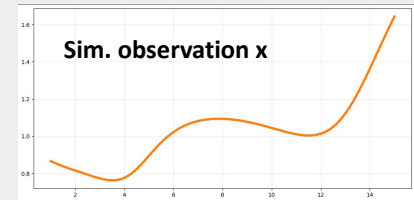
Model Misspecification



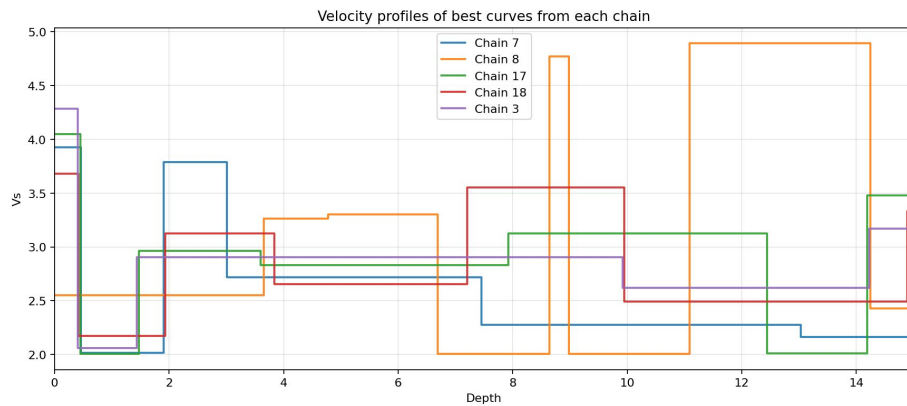
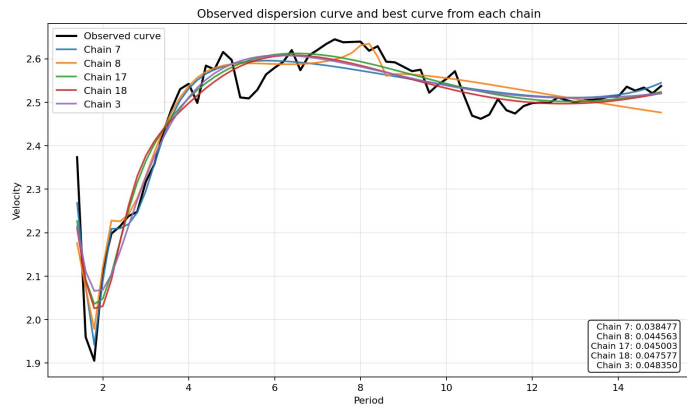
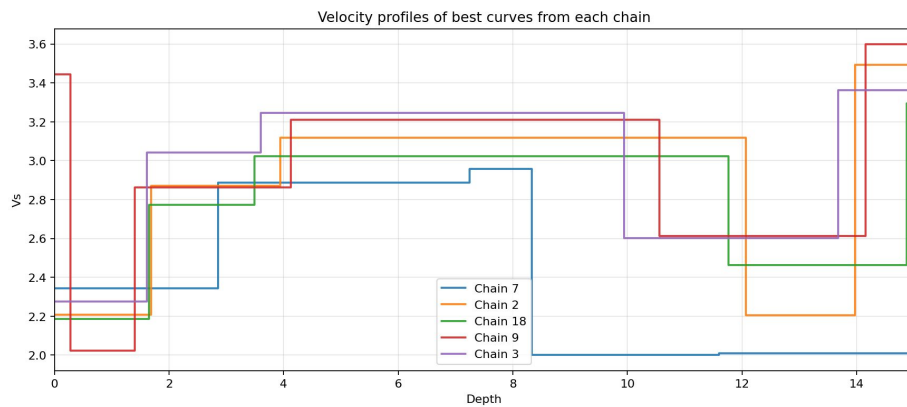
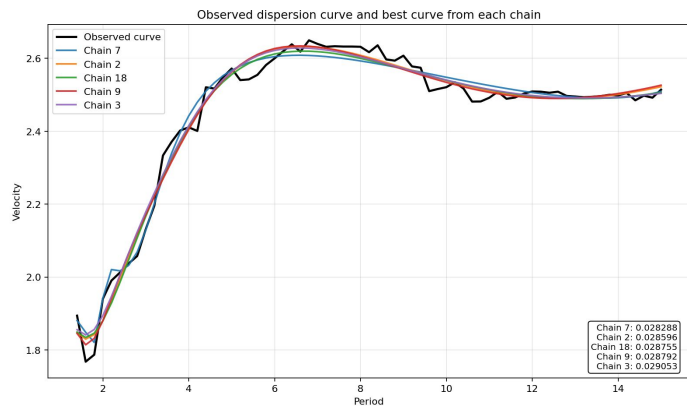
Hand picking



Surfdisp96



Observational support of Surfdisp96



Robust Neural Posterior Estimation

We can explicitly model discrepancies between observed data and simulations with an **error model** :

$$p(\theta, x, y) = p(y|x)p(x|\theta)p(\theta)$$

θ and \mathbf{x} are respectively the parameters and output of the simulator and \mathbf{y} is the observation. \mathbf{x} is treated as an unobserved latent variable.

$$p(\theta|\mathbf{y}) \approx \frac{1}{M} \sum_{m=1}^M p(\theta|\tilde{x}_m), \quad \tilde{x}_1, \dots, \tilde{x}_M, \stackrel{\text{i.i.d.}}{\sim} p(\mathbf{x}|\mathbf{y})$$

If we have access to $p(\mathbf{x}|\mathbf{y})$ and $p(\theta|\mathbf{x})$, we could sample posterior parameters by first sampling from the posterior over \mathbf{x} and then from $p(\theta|\mathbf{x})$.

Example error model

In practice, the true error model is unknown and should describe our prior belief over the discrepancy between the data generating process and simulations. **Spike and slab** error model,

$$\begin{aligned}x &\sim q(x), \\z_j &\sim \text{Bernoulli}(\rho), \\y_j \mid x_j, z_j &\sim \begin{cases} N(x_j, \sigma^2) & \text{if } z_j = 0 \\ \text{Cauchy}(x_j, \tau) & \text{if } z_j = 1 \end{cases}\end{aligned}$$

A key advantage of this error model : given by the latent variable z . The probability of being in the slab $\Pr(z_j = 1 \mid y)$ can be used as **posterior misspecification probability**.

RNPE Algorithm

Algorithm 1: Robust neural posterior estimation (RNPE)

Generate dataset

1 **for** i in $1 : N$ **do**

2 1 Sample $\boldsymbol{\theta}_i \sim p(\boldsymbol{\theta})$

2 2 Simulate $\mathbf{x}_i \sim p(\mathbf{x} \mid \boldsymbol{\theta})$

end

3 Train NPE $q(\boldsymbol{\theta} \mid \mathbf{x})$ on $\{(\boldsymbol{\theta}_i, \mathbf{x}_i)\}_{i=1}^N$

4 Train $q(\mathbf{x})$ on $\{\mathbf{x}_i\}_{i=1}^N$

5 Sample $\tilde{\mathbf{x}}_m \sim \hat{p}(\mathbf{x} \mid \mathbf{y}_o) \propto p(\mathbf{y}_o \mid \mathbf{x})q(\mathbf{x})$, $m = 1, \dots, M$ using MCMC

6 Sample $\tilde{\boldsymbol{\theta}}_m \sim q(\boldsymbol{\theta} \mid \tilde{\mathbf{x}}_m)$, $m = 1, \dots, M$

return $\{(\tilde{\boldsymbol{\theta}}_m, \tilde{\mathbf{x}}_m)\}_{m=1}^M$, samples drawn approximately from $p(\boldsymbol{\theta}, \mathbf{x} \mid \mathbf{y}_o)$

Sample denoised \mathbf{x}_m with MCMC

Learning Robust Statistics

- A common practice is to project both the simulated and the observed data onto a low-dimensional space of summary statistics \mathcal{S} via the mapping $\eta : \mathcal{X}^n \rightarrow \mathcal{S}$ such that $s_{\text{obs}} = \eta(\mathbf{y}_{1:n})$
- Here η is not known **a priori** but is learned using a NN, the **summary networks** η_ψ
- Map each statistic **vector** s to the posterior $p(\theta|s)$ via **inference network** h_φ trained on simulated data.
- Perform poorly when the model is misspecified because the observed statistics s_{obs} becomes OOD under misspecification.
- $\eta\#\mathbb{P}_\theta^n$ is the **distribution of simulated** statistics, $\eta\#\mathbb{Q}^n$ is the **distribution of real statistics**. We define :

$$\varepsilon_\eta = \inf_{\theta \in \Theta} \mathcal{D}(\eta\#\mathbb{P}_\theta^n, \eta\#\mathbb{Q}^n)$$

- ε_η is the minimum non-reducible mismatch. \mathcal{D} is distribution distance (MMD).

Learning Robust Statistics

- ε_n is very slow to compute, but we can compute an upper bound by averaging the distances over some θ sampled from $p(\theta)$.
- To avoid collapse, the loss must be a tradeoff between being informative about θ reducing mismatch,

$$\mathcal{L}_{\text{RS}}(\omega, \psi) = \mathcal{L}(\omega, \psi) + \lambda \mathbb{E}_{p(\theta)} [\mathcal{D}(\eta \# \mathbb{P}_\theta^n, \eta \# \mathbb{Q}^n)]$$

- ω are the **inference model** parameters, ψ the summary network parameters and λ the regularization parameter.
- We don't have access to the complete distribution $\eta \# \mathbb{P}_\theta^n$ and $\eta \# \mathbb{Q}^n$ only to some samples of these ($x_{1:n,i} \sim \mathbb{P}_{\theta_i}$ and $y_{1:n,i} \sim \mathbb{Q}$).
- $\lambda \rightarrow 0$ converge to classical NPE methods. $\lambda \rightarrow \infty$ converges to the prior distribution.
- λ encodes the trade-off between efficiency and robustness.