

COMPUTATIONALLY EFFICIENT INFERENCE FOR SPARSITY-PROMOTING HIERARCHICAL BAYESIAN MODELS^a

BAYESIAN INVERSE PROBLEMS AND UQ @ ICERM

Jan Glaubitz

Joint work with Youssef Marzouk (MIT)

March 04, 2026

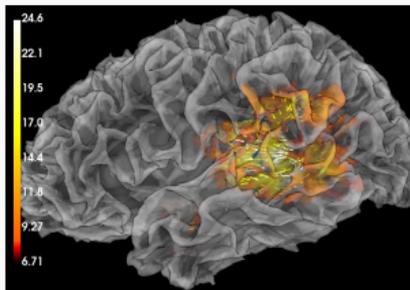
Department of Mathematics

Linköping University

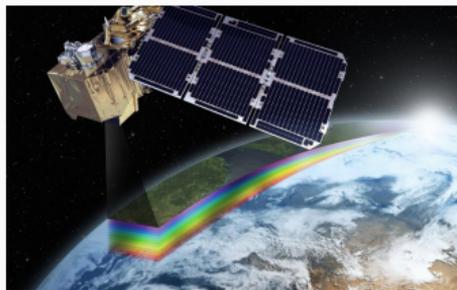
Linköping, Sweden

^aThis work was supported by the DOD (ONR MURI) grant #N00014-20-1-2595, the Swedish Research Council (VR) Starting Grant #2025-05370, the Zenith Career Development Grant #26.07, and the National Academic Infrastructure for Supercomputing in Sweden (NAISS) grants #2024/22-1207 and #2025/22-1599.

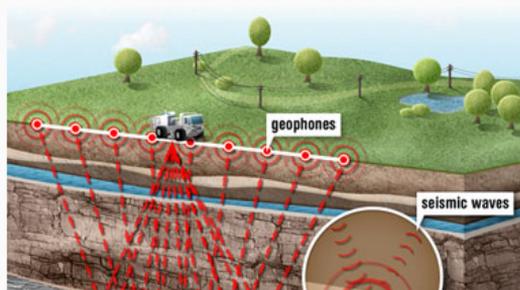
Inverse problem $y = F(x) + e$



Medical imaging (MEG, CT, MRI)

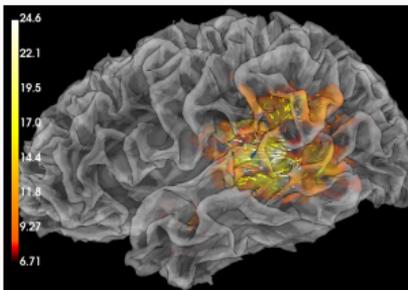


Remote sensing

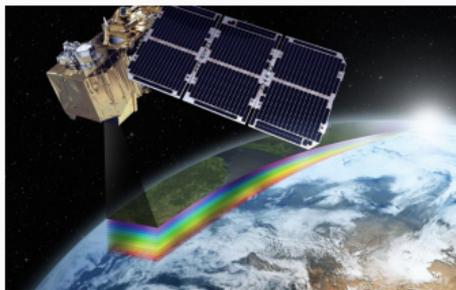


Geoscience

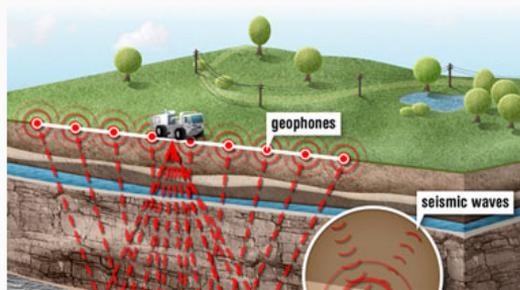
Inverse problem $y = F(x) + e$



Medical imaging (MEG, CT, MRI)



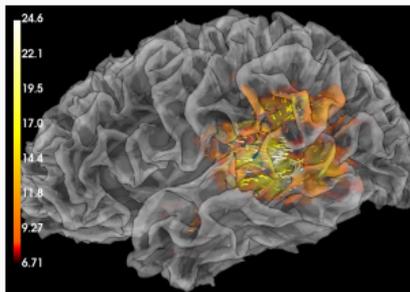
Remote sensing



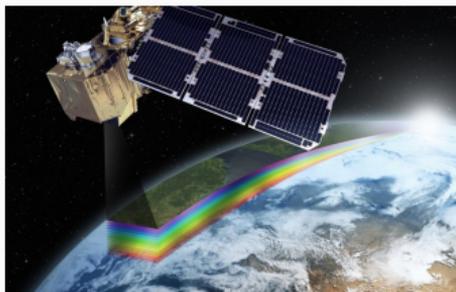
Geoscience

Belief. x or Ψx is sparse

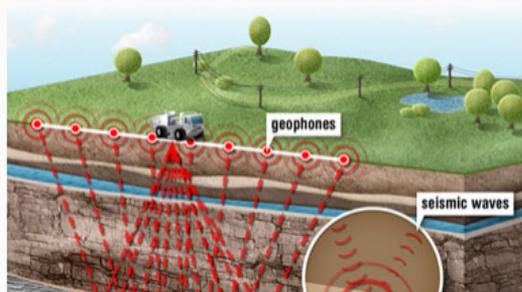
Inverse problem $y = F(x) + e$



Medical imaging (MEG, CT, MRI)



Remote sensing



Geoscience

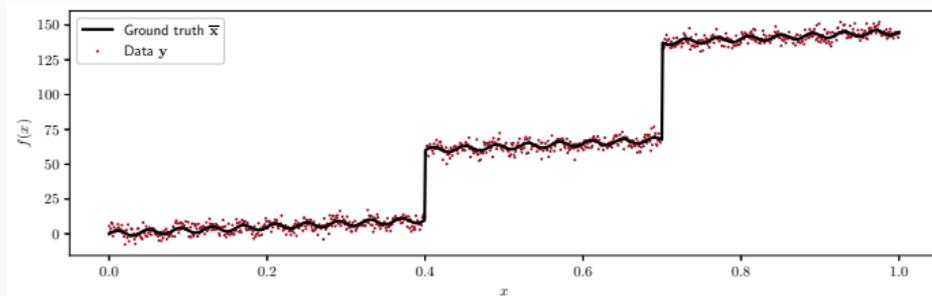
Belief. x or Ψx is sparse

Bayesian approach^{1,2}: $\pi^y(x) \propto f(x; y) \pi^0(x)$

¹Stuart, "Inverse problems: a Bayesian perspective," 2010

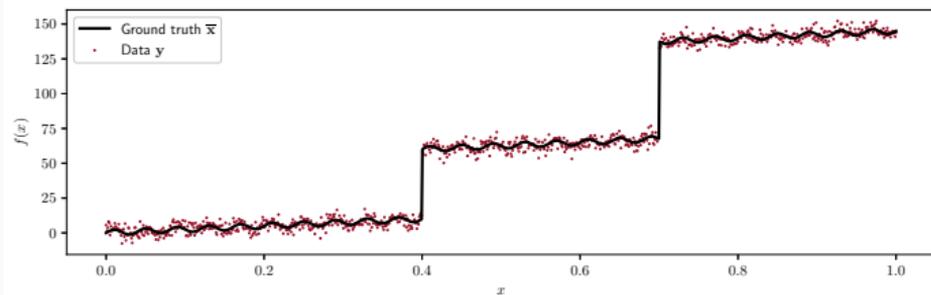
²Calvetti & Somersalo, "Bayesian Scientific Computing," 2023

HIERARCHICAL SPARSITY-PROMOTING PRIORS



Assumption: $\Psi \mathbf{x} \in \mathbb{R}^K$ is sparse—e.g., $[\Psi \mathbf{x}]_k = x_{k+1} - x_k$

HIERARCHICAL SPARSITY-PROMOTING PRIORS



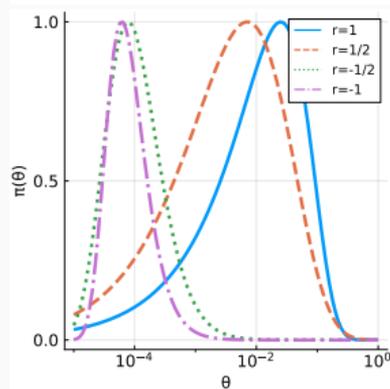
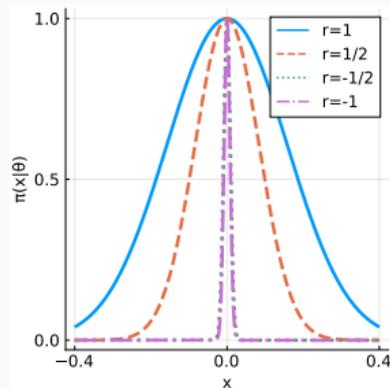
Assumption: $\Psi \mathbf{x} \in \mathbb{R}^K$ is sparse—e.g., $[\Psi \mathbf{x}]_k = x_{k+1} - x_k$

Sparse Bayesian learning (SBL) priors^{ab}

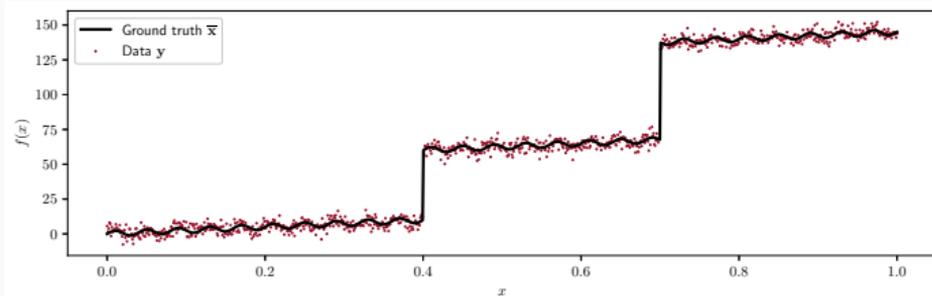
$$\begin{cases} [\Psi \mathbf{x}]_k \sim \mathcal{N}(0, \theta_k), \\ \theta_k \sim \mathcal{GG}(r, \beta, \vartheta), \end{cases} \quad k = 1, \dots, K$$

^aTipping, "Sparse Bayesian learning and the relevance vector machine," 2001

^bCalvetti, Pragliola, Somersalo, Strang, "Sparse reconstructions from few noisy data: analysis of hierarchical Bayesian models with generalized gamma hyperpriors," 2020



HIERARCHICAL SPARSITY-PROMOTING PRIORS

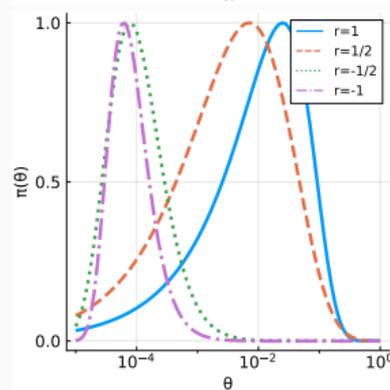
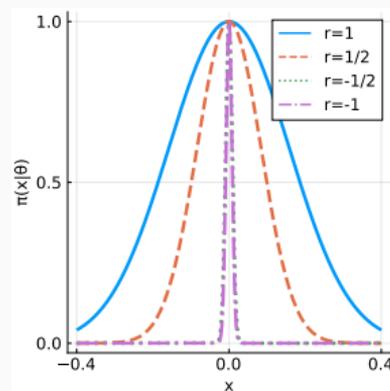


Assumption: $\Psi \mathbf{x} \in \mathbb{R}^K$ is sparse—e.g., $[\Psi \mathbf{x}]_k = x_{k+1} - x_k$

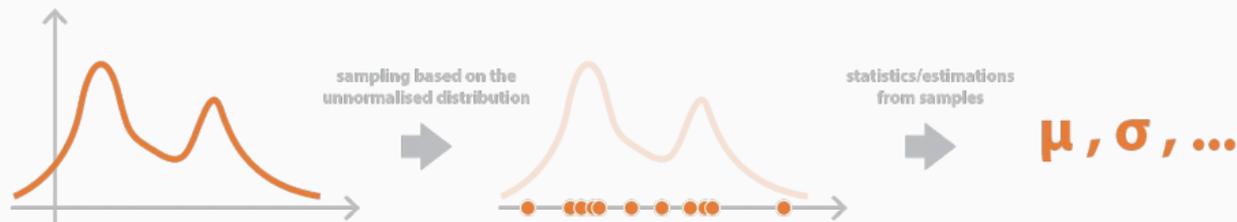
Sparse Bayesian learning (SBL) priors^{ab}

$$\begin{cases} [\Psi \mathbf{x}]_k \sim \mathcal{N}(0, \theta_k), \\ \theta_k \sim \mathcal{GG}(r, \beta, \vartheta), \end{cases} \quad k = 1, \dots, K$$

- \mathcal{GG} makes $\theta_k \approx 0$ likely, while allowing for outliers
- If $\theta_k \approx 0$, then $[\Psi \mathbf{x}]_k \approx 0$ is most likely \rightarrow no jump!
- If $\theta_k \gg 0$, then $|[\Psi \mathbf{x}]_k| \gg 0$ more likely \rightarrow jump!



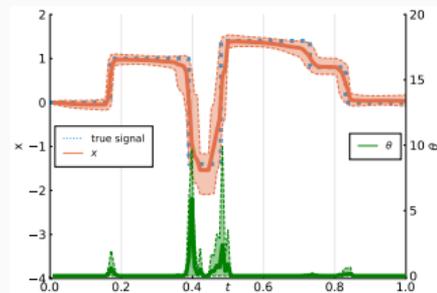
SAMPLING & UQ

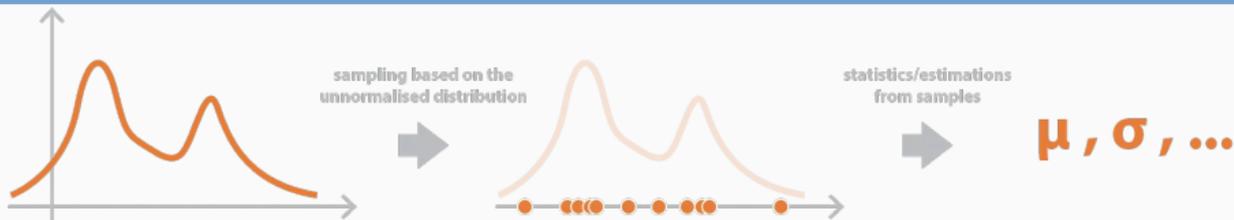


Goal. Sample from $\pi^y(\mathbf{x}, \theta) \propto f(\mathbf{x}; \mathbf{y}) \pi^0(\mathbf{x}, \theta)$

Enables approxating posterior expectations:

$$\int \mathcal{G}(\mathbf{x}, \theta) \pi^y(\mathbf{x}, \theta) d(\mathbf{x}, \theta) \approx \frac{1}{I} \sum_{i=1}^I \mathcal{G}(\mathbf{x}_i, \theta_i)$$





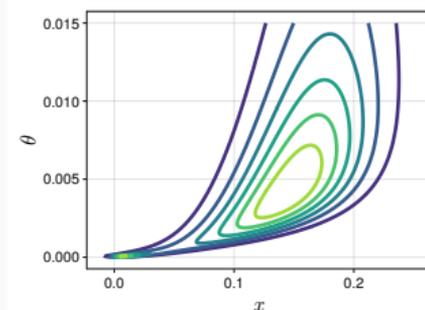
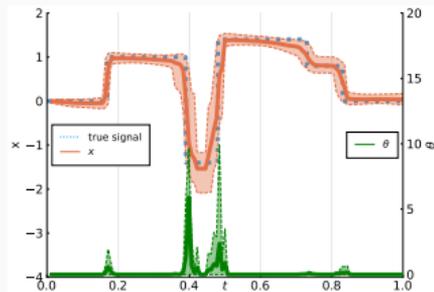
Goal. Sample from $\pi^y(x, \theta) \propto f(x; y) \pi^0(x, \theta)$

Enables approximating posterior expectations:

$$\int \mathcal{G}(x, \theta) \pi^y(x, \theta) d(x, \theta) \approx \frac{1}{I} \sum_{i=1}^I \mathcal{G}(x_i, \theta_i)$$

Challenges for sparsity-promoting priors.^a

- (i) multi-modal
- (ii) strong correlation between x and θ
- (iii) high-dimensional



^aCalvetti, Somersalo, "Computationally efficient sampling methods for sparsity promoting hierarchical Bayesian models," 2024

OUR APPROACH: HIERARCHICAL PRIOR NORMALIZATION³

Idea. Transform the SBL prior $\pi^0(\mathbf{x}, \boldsymbol{\theta})$ into a standard normal one $\phi^0(\mathbf{u}, \boldsymbol{\tau})$

Need. $S : (\mathbf{x}, \boldsymbol{\theta}) \mapsto (\mathbf{u}, \boldsymbol{\tau})$ s. t. $(\mathbf{u}, \boldsymbol{\tau})$ follows ϕ^0 whenever $(\mathbf{x}, \boldsymbol{\theta})$ follows π^0



³G. & Marzouk, "Efficient sampling for sparse Bayesian learning using hierarchical prior normalization," 2025

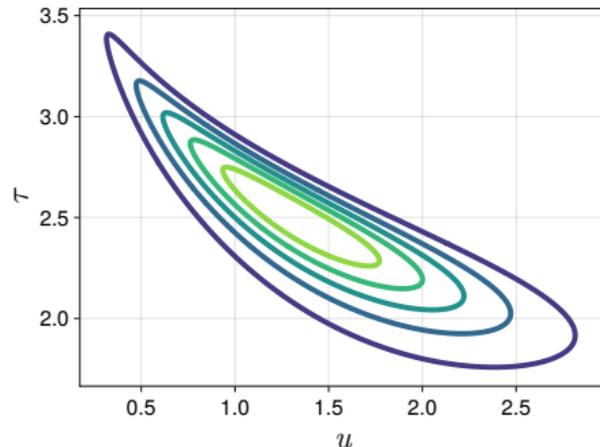
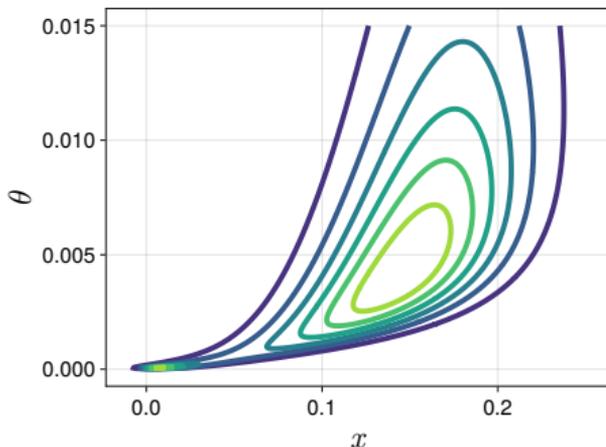
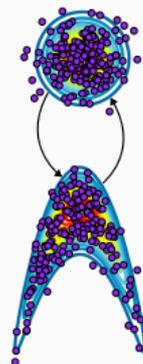
OUR APPROACH: HIERARCHICAL PRIOR NORMALIZATION³

Idea. Transform the SBL prior $\pi^0(\mathbf{x}, \theta)$ into a standard normal one $\phi^0(\mathbf{u}, \tau)$

Need. $S : (\mathbf{x}, \theta) \mapsto (\mathbf{u}, \tau)$ s. t. (\mathbf{u}, τ) follows ϕ^0 whenever (\mathbf{x}, θ) follows π^0

Get. Simpler “prior-normalized” posterior

$$\phi^y(\mathbf{u}, \tau) = (S_{\#}\pi^y)(\mathbf{u}, \tau) = \frac{1}{Z} f(S^{-1}(\mathbf{u}, \tau); \mathbf{y}) \phi^0(\mathbf{u}, \tau)$$



Product-like form. ($\Psi = I$, i.e., \mathbf{x} sparse)

$$\pi^0(\mathbf{x}, \boldsymbol{\theta}) = \prod_{i=1}^n \mathcal{N}(x_i | 0, \theta_i) \mathcal{G}\mathcal{G}(\theta_i | r, \beta, \vartheta) \implies S(\mathbf{x}, \boldsymbol{\theta}) = \begin{bmatrix} s_1(x_1, \theta_1) \\ \vdots \\ s_n(x_n, \theta_n) \end{bmatrix}$$

with $s_i : (x_i, \theta_i) \mapsto (u_i, \tau_i)$ transforming $\mathcal{N}(x_i | 0, \theta_i) \mathcal{G}\mathcal{G}(\theta_i | r, \beta, \vartheta)$ into $\mathcal{N}(\mathbf{0}_2, I_2)$

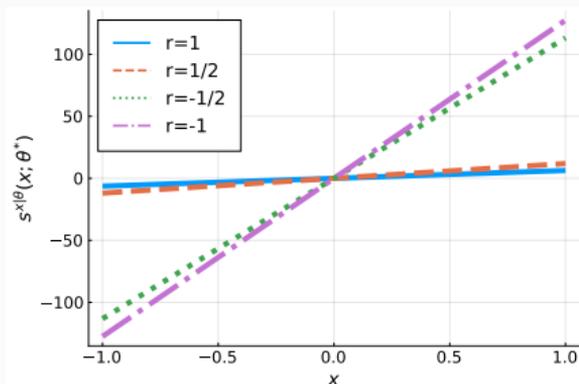
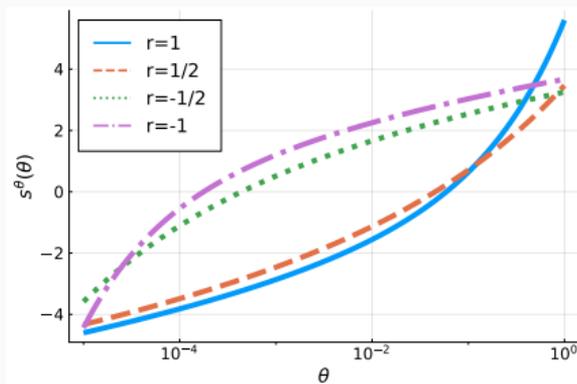
Product-like form. ($\Psi = I$, i.e., \mathbf{x} sparse)

$$\pi^0(\mathbf{x}, \theta) = \prod_{i=1}^n \mathcal{N}(x_i | 0, \theta_i) \mathcal{G}\mathcal{G}(\theta_i | r, \beta, \vartheta) \implies S(\mathbf{x}, \theta) = \begin{bmatrix} s_1(x_1, \theta_1) \\ \vdots \\ s_n(x_n, \theta_n) \end{bmatrix}$$

with $s_i : (x_i, \theta_i) \mapsto (u_i, \tau_i)$ transforming $\mathcal{N}(x_i | 0, \theta_i) \mathcal{G}\mathcal{G}(\theta_i | r, \beta, \vartheta)$ into $\mathcal{N}(\mathbf{0}_2, I_2)$

Knothe–Rosenblatt (KR) rearrangements.

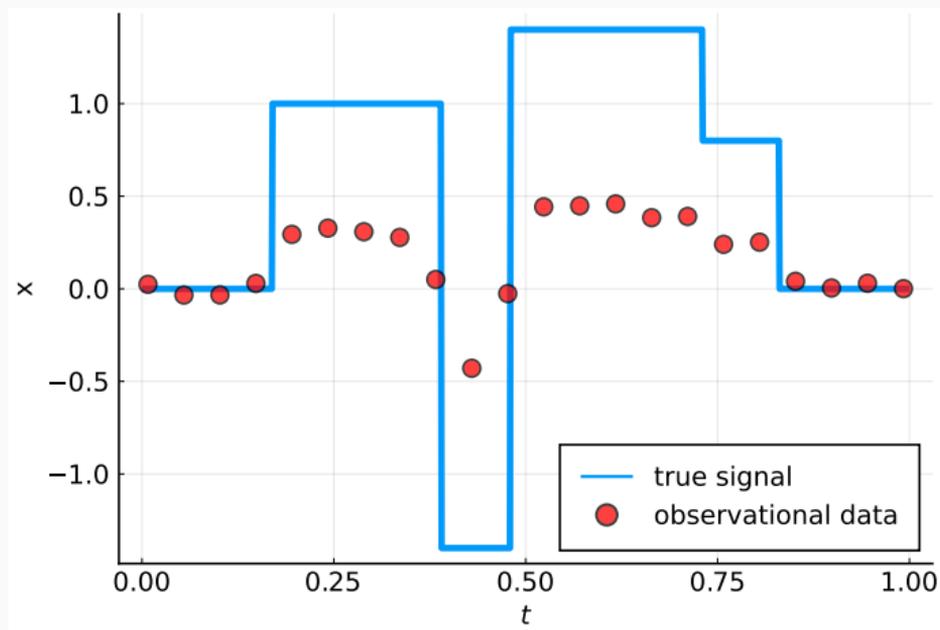
$$s(x, \theta) = \begin{bmatrix} s^\theta(\theta) \\ s^{x|\theta}(x; \theta) \end{bmatrix}, \quad \begin{cases} s^\theta = (\Psi^0)^{-1} \circ \mathcal{P}^\theta, \\ s^{x|\theta} = x / \sqrt{\theta}. \end{cases}$$



EXAMPLE 1: SIGNAL DEBLURRING

Goal. Recover $n = 128$ nodal values from $m = 22$ noisy blurry observations:

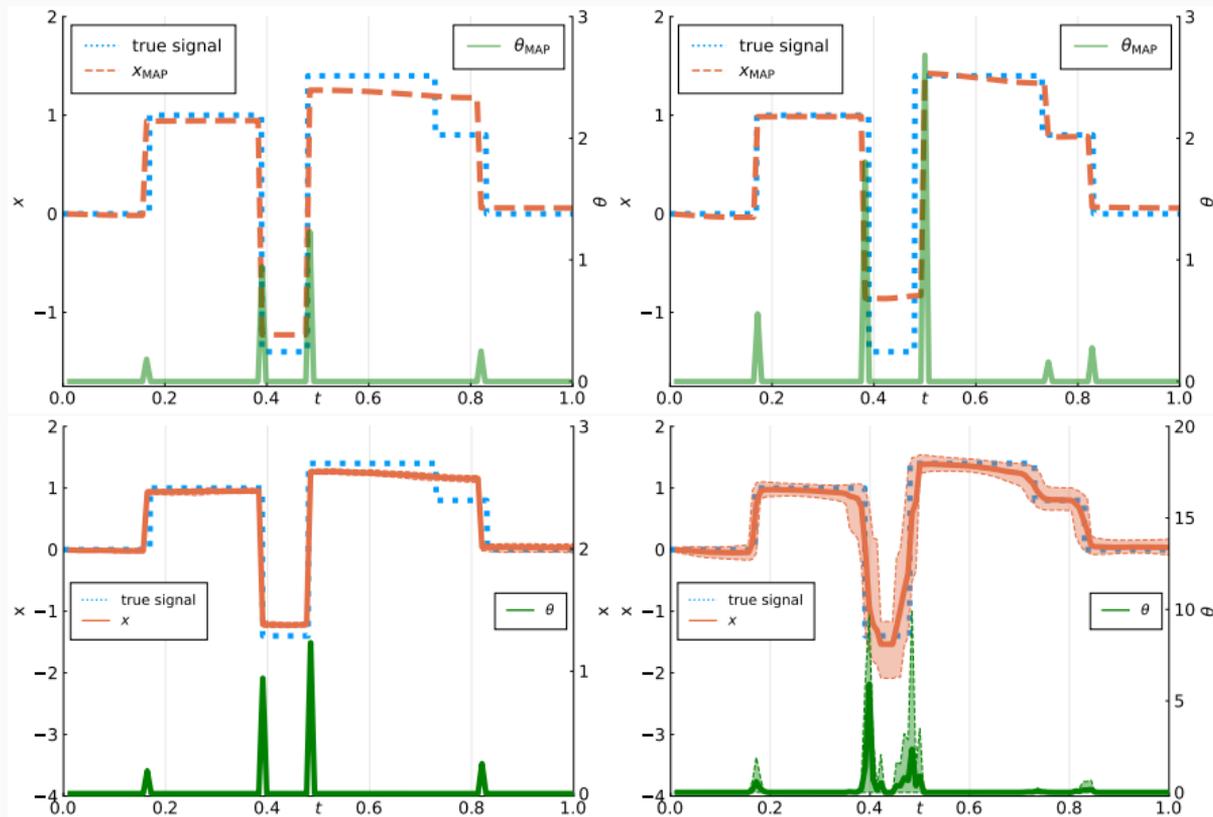
$$y = Fx + e \quad \text{where} \quad [\Psi x]_i = x_{i+1} - x_i \text{ is sparse}$$



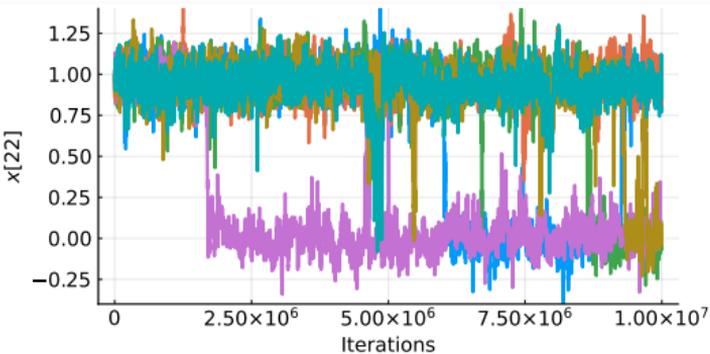
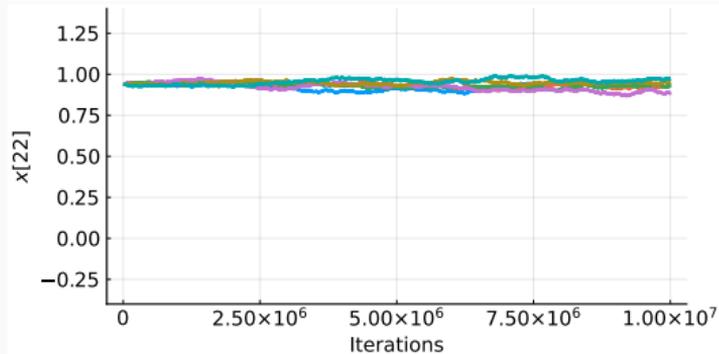
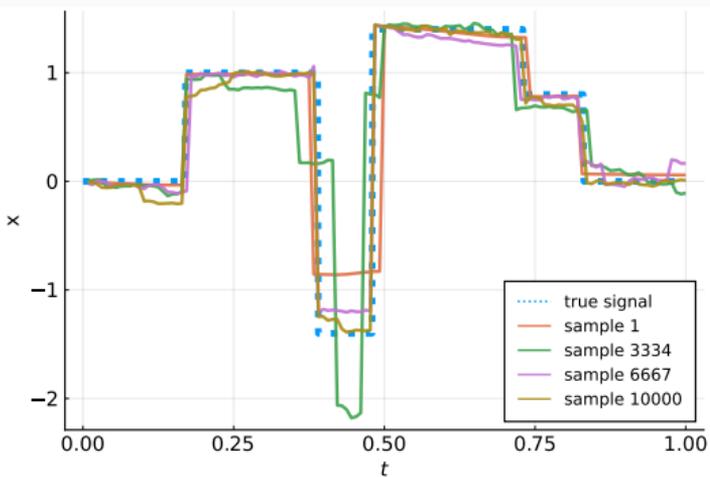
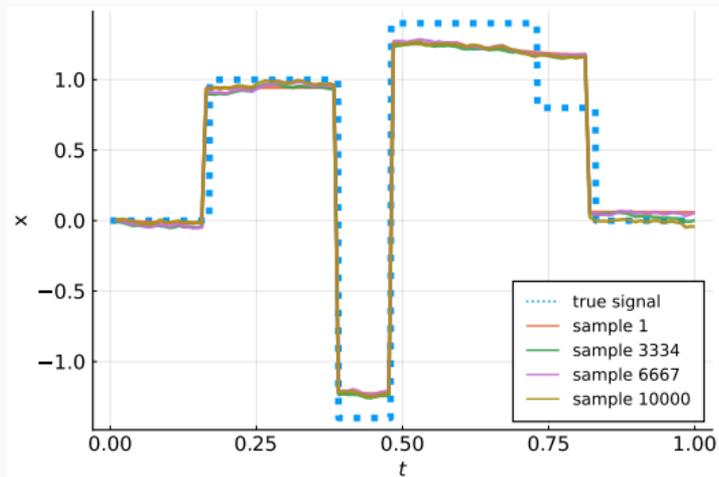
Compare. AM sampler for **original** and **prior-normalized posterior** with $r = -1$

STATISTICS (MAP INITIALIZATION)

Setup. $r = -1$, $J = 6$ chains, 10^7 samples, MAP initialization

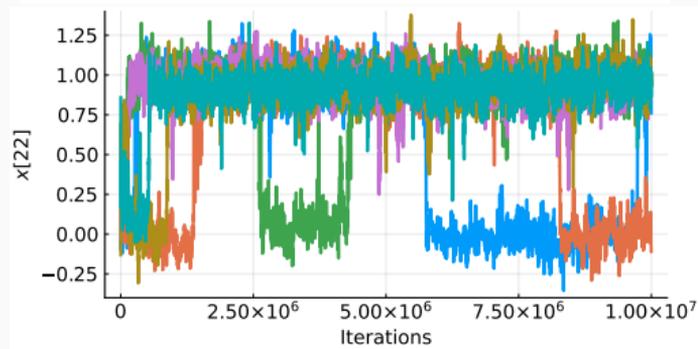
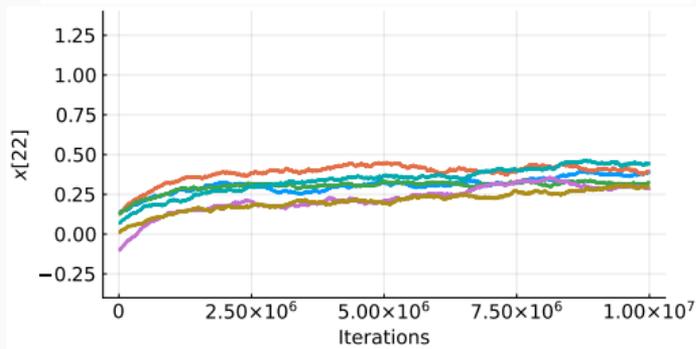
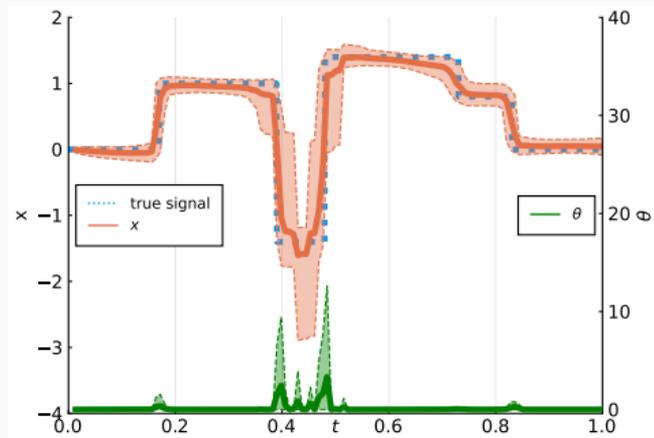
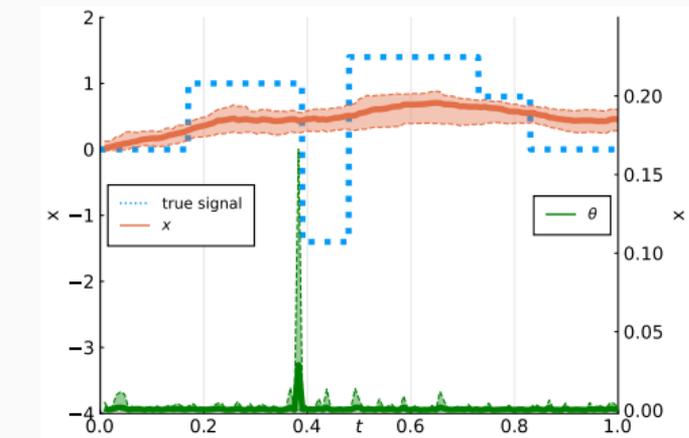


SAMPLES & TRACES (MAP INITIALIZATION)

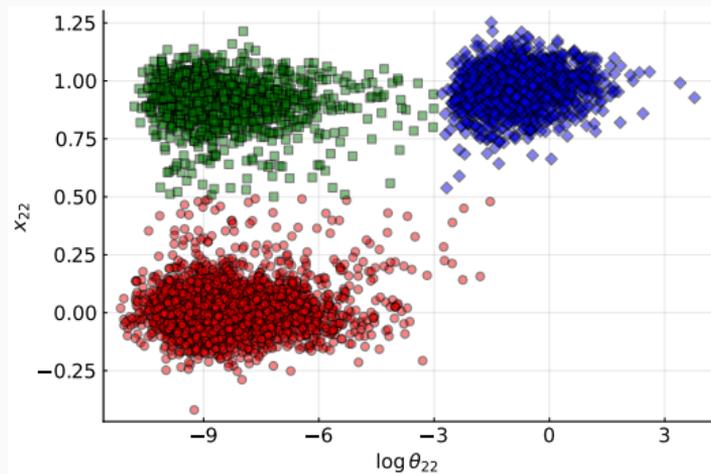


STATISTICS & SAMPLES (RANDOM INITIALIZATION)

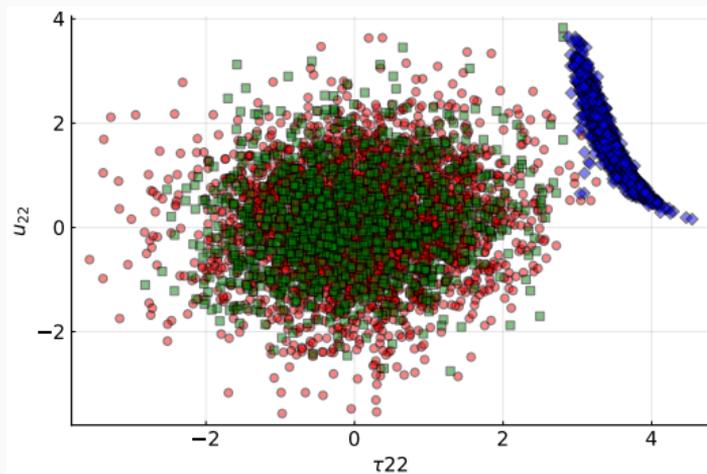
Setup. $r = -1$, $J = 6$ chains, 10^7 samples, random initialization



SCATTER PLOTS (RANDOM INITIALIZATION)



(a) Original posterior



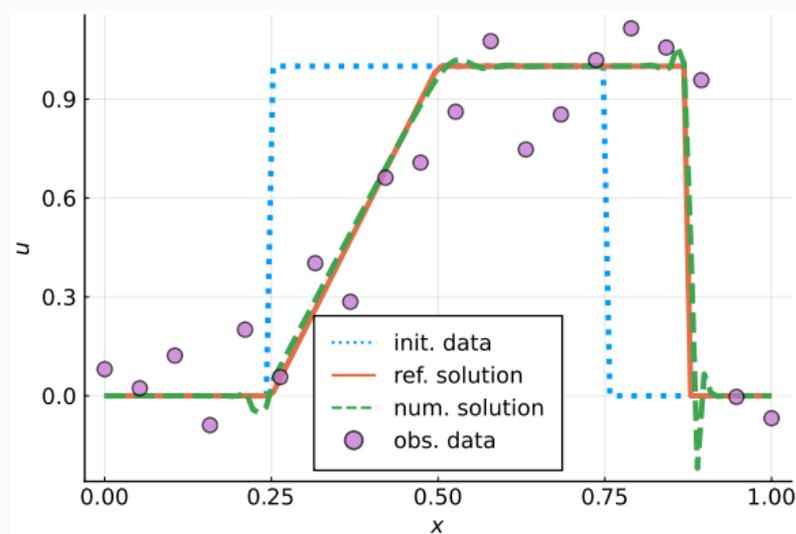
(b) Prior-normalized posterior

COMPUTATIONAL EXAMPLE II: INVERSE BURGERS' EQUATION

Burgers' equation. $\partial_t u(x, t) + \partial_x u(x, t)^2 = 0$ with $u(x, 0) = u_0(x)$

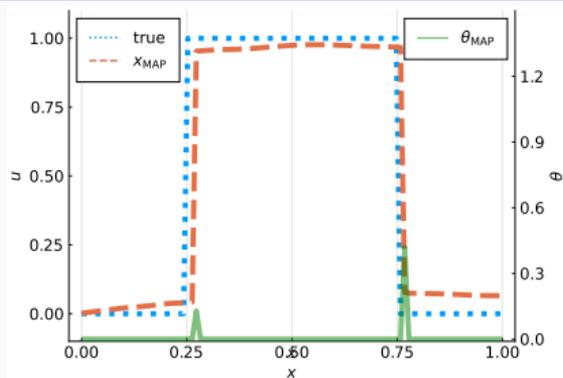
Goal. Recover u_0 at $n = 100$ points from $m = 20$ noisy solution observations:

$$\mathbf{b} = \mathcal{F}(\mathbf{u}) + \mathbf{e} \quad \text{where} \quad [\Psi \mathbf{u}]_j = u_{i+1} - u_i \text{ is sparse}$$

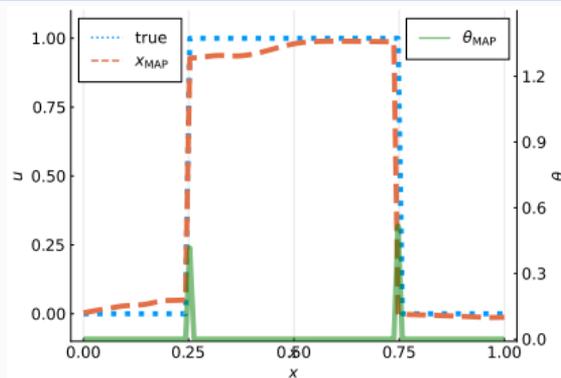


Compare. AM sampler (10^7) for original and prior-normalized posterior, $r = -1$

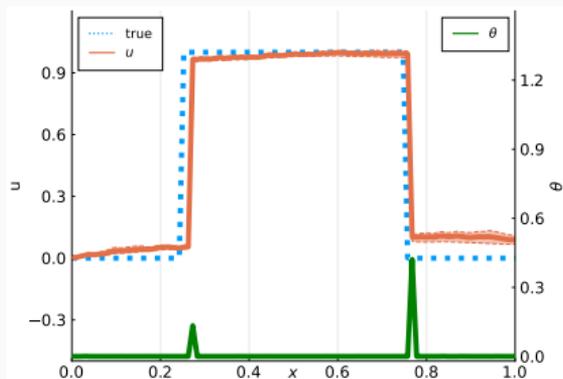
MAP ESTIMATES & STATISTICS



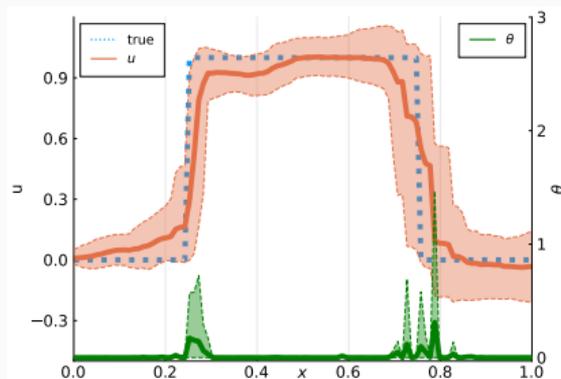
Original posterior, MAP estimate



Prior-normalized posterior, MAP estimate

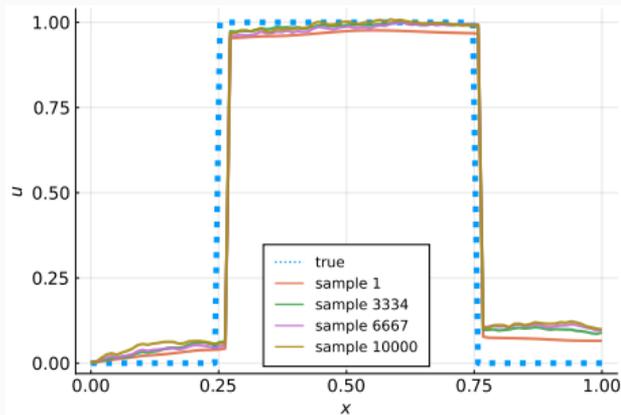


Original posterior, statistics

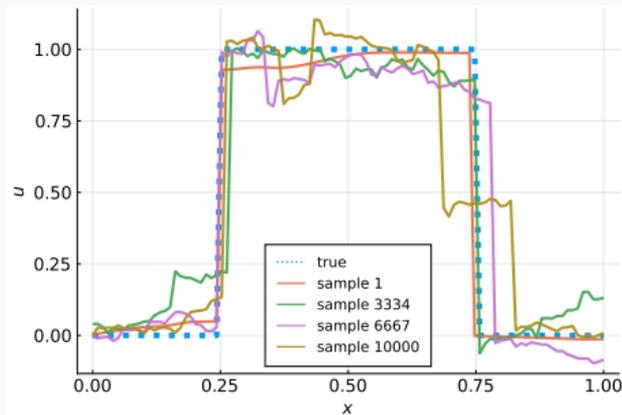


Prior-normalized posterior, statistics

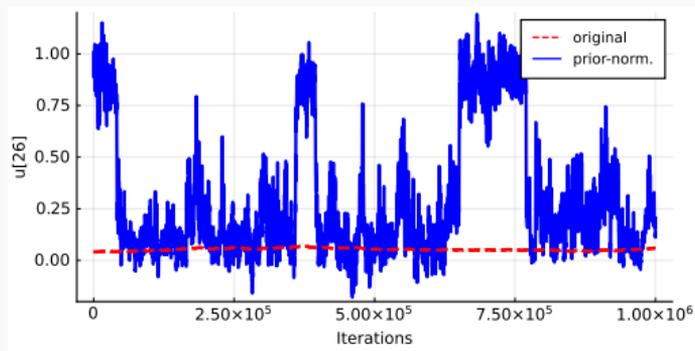
SAMPLES & TRACES: MAP INITIALIZED



Original posterior, samples

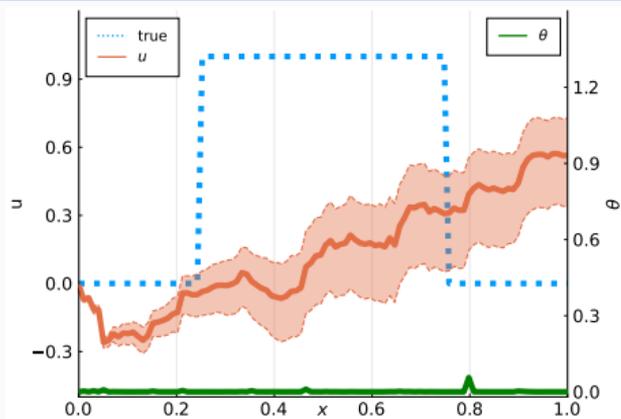


Prior-normalized posterior, samples

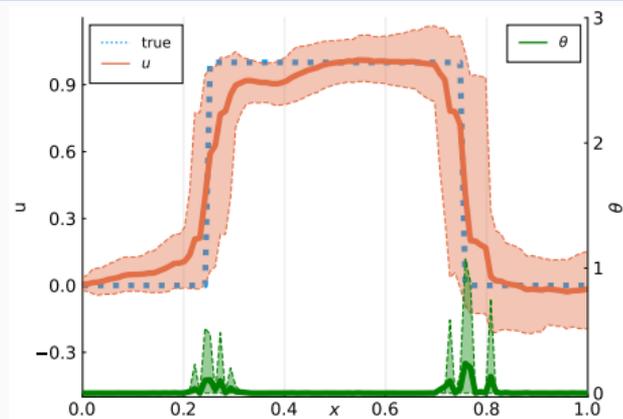


Traces

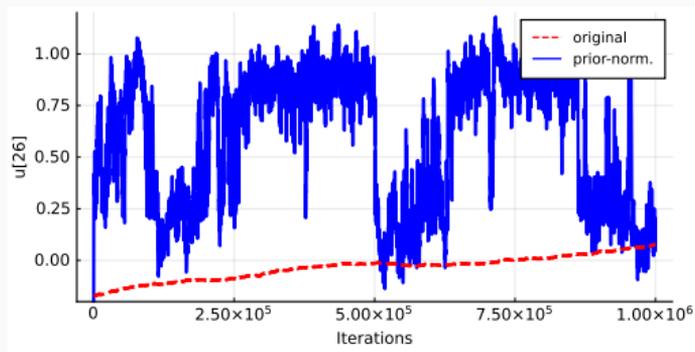
STATISTICS & SAMPLES: RANDOMLY INITIALIZED



Original posterior, statistics



Prior-normalized posterior, statistics

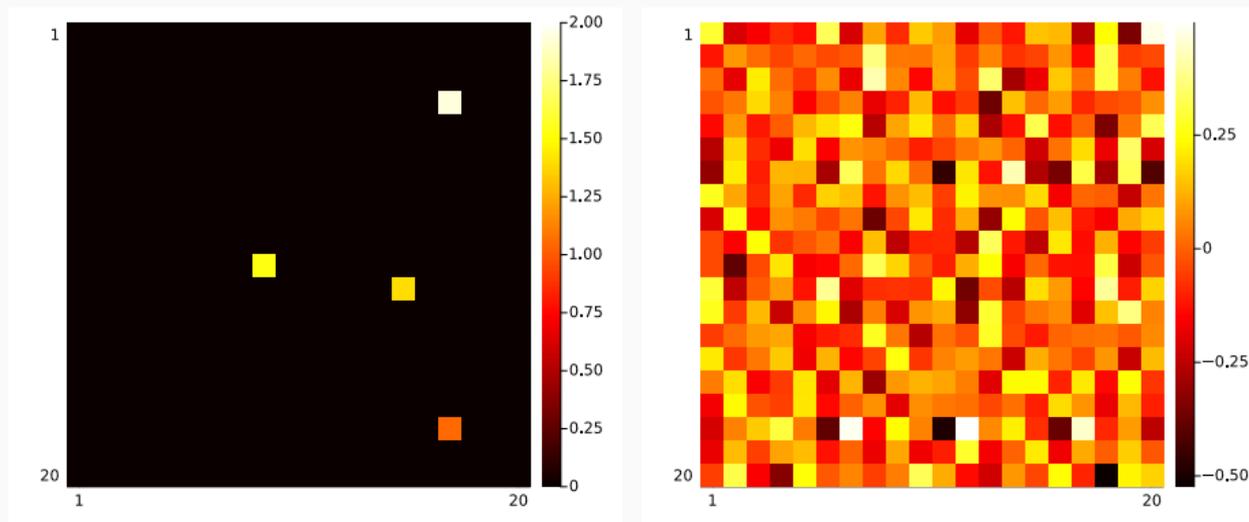


Traces

EXAMPLE 3: IMPULSE IMAGE

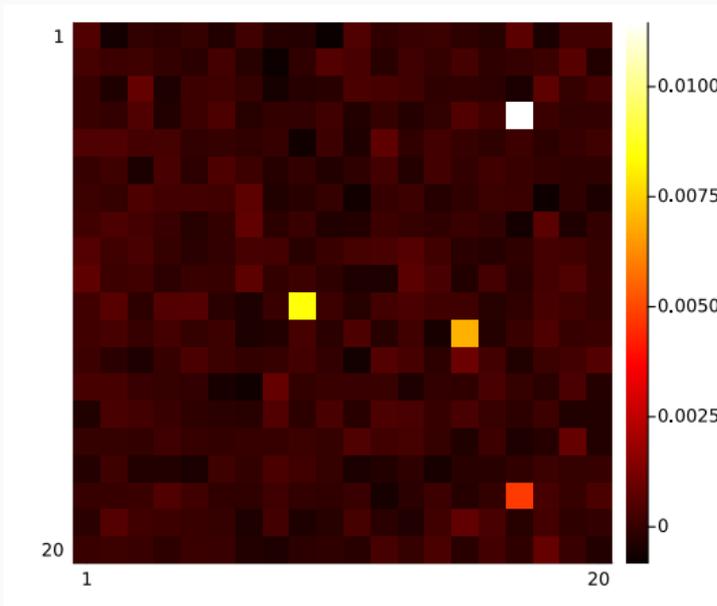
Goal. Recover 20×20 impulse image from its noisy DCT:

$$y = Fx + e \quad \text{where } x \text{ is sparse}$$

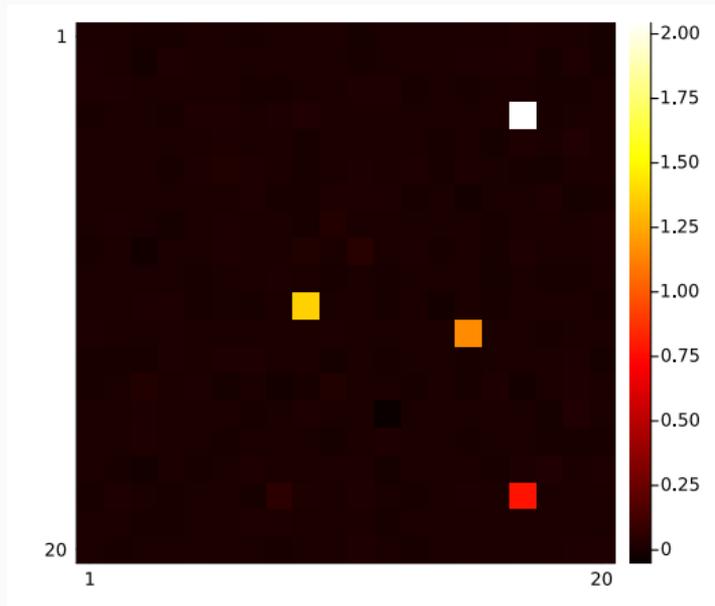


Setup. $r = -1$, $J = 4$ chains, random initialization

Compare. **Original posterior** (Gibbs sampler, 10^5 samples, ≈ 20 m runtime) vs **prior-normalized posterior** (ESS, 10^4 , ≈ 3 m runtime)



(a) Original posterior



(b) Prior-normalized posterior

- Generally applies to scale-mixtures-of-normal priors [**Jonathan Lindbloom**]
- Combine with geometry-exploiting (dimension-robust) samplers
- High-dimensional problems (dynamic inverse problems) [**Mirjeta Pasha**]
- Chemical reaction neural network [**Federica Milinanni**]
- Use as a preconditioner—combined with, e.g., flow matching or diffusion
- **What else?**

For more details:

Glaubitz and Marzouk. “Efficient sampling for sparse Bayesian learning using hierarchical prior normalization.” To appear in SIAM-ASA UQ (2026).

The End

Thank You!



VR Starting Grant #2025-05370



Zenith Career Development Grant



NAISS grant #2024/22-1207