

Transport methods for Bayesian computation: Part 2

Youssef Marzouk

joint work with Daniele Bigoni, Matthew Parno,
Alessio Spantini, & Olivier Zahm

Department of Aeronautics and Astronautics
Center for Computational Engineering
Statistics and Data Science Center

Massachusetts Institute of Technology
<http://uqgroup.mit.edu>

Support from AFOSR, DARPA, DOE

25–26 September 2019

What to do when $T_{\#}\eta \neq \pi$?

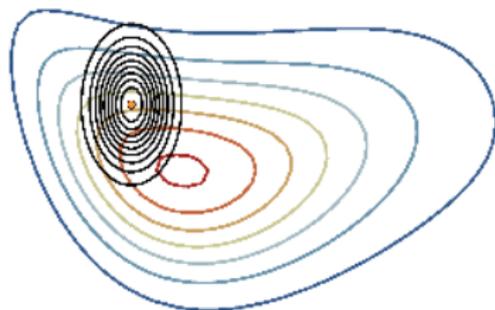
- ▶ Maybe close enough? Can evaluate variance diagnostic $\mathbb{V}\text{ar}_{\eta}[\log(\eta/T_{\#}^{-1}\bar{\pi})]$, bound $\text{Tr}(H_{T_{\#}^{-1}\pi})$, etc.
- ▶ Enrich T , e.g., add a layer or expand \mathcal{T}_{Δ}^h in the given layer
- ▶ **Sample the pullback:** treat $T_{\#}^{-1}\pi$ with an asymptotically exact scheme, e.g., Markov chain Monte Carlo

What to do when $T_{\#}\eta \neq \pi$?

- ▶ Maybe close enough? Can evaluate variance diagnostic $\text{Var}_{\eta}[\log(\eta/T_{\#}^{-1}\bar{\pi})]$, bound $\text{Tr}(H_{T_{\#}^{-1}\pi})$, etc.
- ▶ Enrich T , e.g., add a layer or expand \mathcal{T}_{Δ}^h in the given layer
- ▶ **Sample the pullback:** treat $T_{\#}^{-1}\pi$ with an asymptotically exact scheme, e.g., Markov chain Monte Carlo

One possible construction: **transport-accelerated MCMC**

- ▶ Transport map “preconditions” MCMC target; use MCMC iterates in maps-from-samples construction
- ▶ Can be understood in the framework of *adaptive MCMC*



- ▶ Effective MCMC proposal = **adapted to the target**
 - ▶ Can we transform *proposals* or, equivalently, *targets* for better sampling?

$$\min_{S \in \mathcal{S}_{\Delta}^h} \mathcal{D}_{KL}(S_{\#}\pi \parallel \eta) = \min_{S \in \mathcal{S}_{\Delta}^h} \mathcal{D}_{KL}(\pi \parallel S_{\#}^{-1}\eta)$$

- ▶ Suppose we have Monte Carlo samples $\{x_i\}_{i=1}^M \sim \pi$
- ▶ For standard Gaussian η , this problem is **convex** and **separable** for any π
- ▶ This is *density estimation via transport!* (cf. Tabak & Turner 2013)

Recall maps-from-samples construction

$$\min_{S \in \mathcal{S}_{\Delta}^h} \mathcal{D}_{KL}(S_{\#}\pi \parallel \eta) = \min_{S \in \mathcal{S}_{\Delta}^h} \mathcal{D}_{KL}(\pi \parallel S_{\#}^{-1}\eta)$$

- ▶ Suppose we have Monte Carlo samples $\{x_i\}_{i=1}^M \sim \pi$
- ▶ For standard Gaussian η , this problem is **convex** and **separable** for any π
- ▶ This is *density estimation via transport!* (cf. Tabak & Turner 2013)
- ▶ Equivalent to maximum likelihood estimation of S

$$\hat{S} \in \arg \max_{S \in \mathcal{S}_{\Delta}^h} \frac{1}{M} \sum_{i=1}^M \log \underbrace{S_{\#}^{-1} \eta(x_i)}_{\text{pullback}}, \quad \eta = \mathcal{N}(0, \mathbf{I}_n),$$

- ▶ Each component \hat{S}^k of \hat{S} can be computed *separately*, via smooth convex optimization

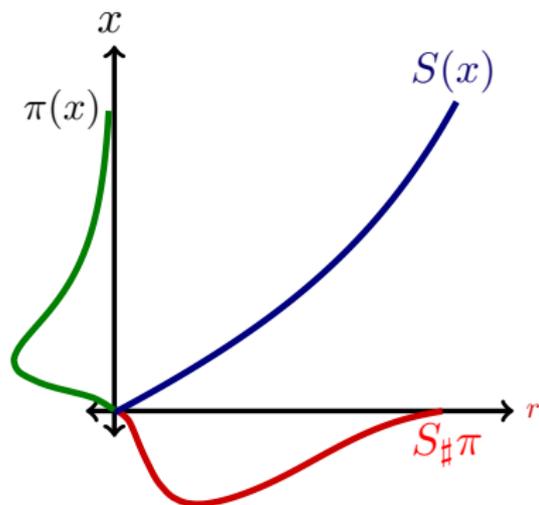
$$\hat{S}^k \in \arg \min_{S^k \in \mathcal{S}_{\Delta,k}^h} \frac{1}{M} \sum_{i=1}^M \left(\frac{1}{2} S^k(x_i)^2 - \log \partial_k S^k(x_i) \right)$$

- ▶ View $\hat{S}_{\#}\pi$ as the “preconditioned” target
 - ▶ In the MCMC setting, $\{x_i\}_{i=1}^M$ comprises *dependent* MCMC samples
 - ▶ $\hat{S}_{\#}\pi$ may be far from standard Gaussian for small M and/or crude \mathcal{S}_{Δ}^h

Map-accelerated MCMC

- **Ingredient #1: static map**

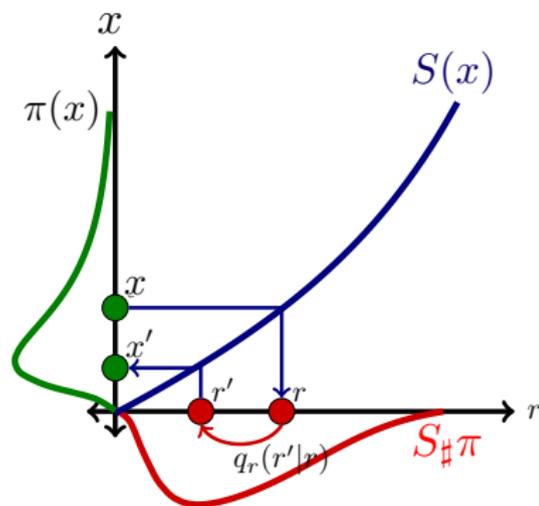
- Perform MCMC in the reference space, on the “preconditioned” density
- Simple proposal in reference space (e.g., random walk) corresponds to a more complex/tailored proposal on target



Map-accelerated MCMC

- **Ingredient #1: static map**

- Perform MCMC in the reference space, on the “preconditioned” density
- Simple proposal in reference space (e.g., random walk) corresponds to a more complex/tailored proposal on target



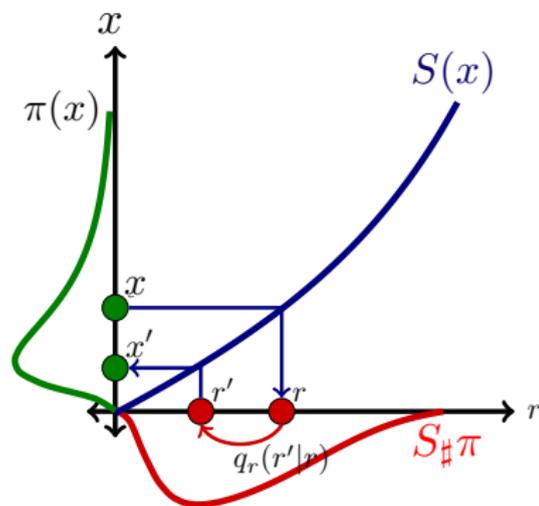
$$\alpha = \frac{\pi(S^{-1}(r')) \left| \nabla S^{-1} \right|_{r'} q_r(r | r')}{\pi(S^{-1}(r)) \left| \nabla S^{-1} \right|_r q_r(r' | r)}$$

simple proposal q_r on **pushforward of target through map**

Map-accelerated MCMC

- **Ingredient #1: static map**

- Perform MCMC in the reference space, on the “preconditioned” density
- Simple proposal in reference space (e.g., random walk) corresponds to a more complex/tailored proposal on target



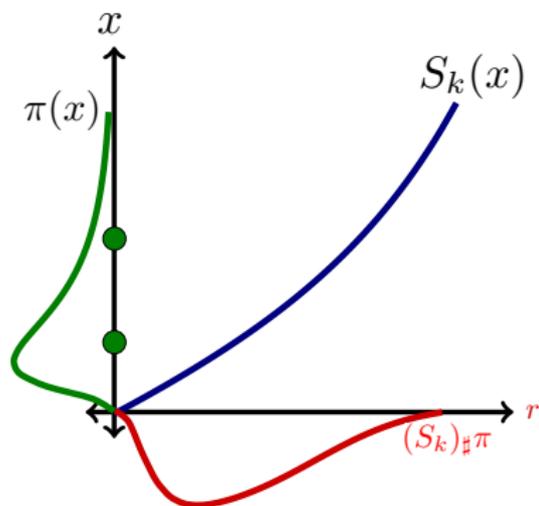
$$\alpha = \frac{\pi(S^{-1}(r')) \left| \nabla S^{-1} \right|_{r'} q_r(r | r')}{\pi(S^{-1}(r)) \left| \nabla S^{-1} \right|_r q_r(r' | r)}$$

more complex proposal, directly on
target distribution

Map-accelerated MCMC

- **Ingredient #2: adaptive map**

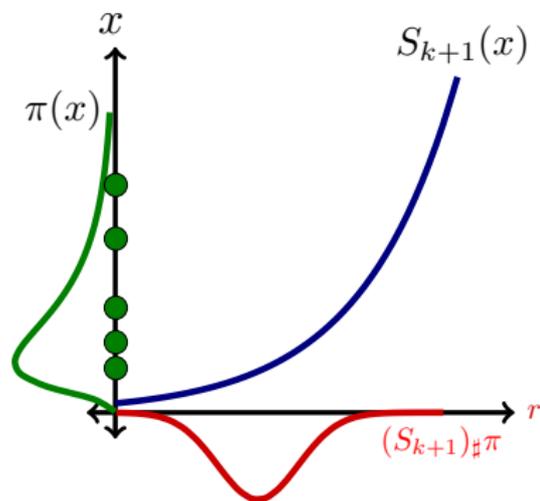
- Update the map with each MCMC iteration:
more samples, more accurate \mathbb{E}_π , better S
- Adaptive MCMC [Haario 2001, Andrieu 2006], but with nonlinear transformation to capture non-Gaussian structure



Map-accelerated MCMC

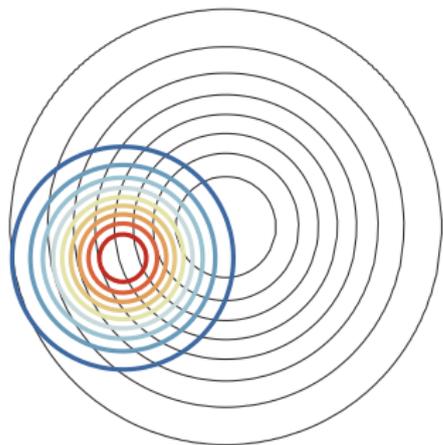
- **Ingredient #2: adaptive map**

- Update the map with each MCMC iteration:
more samples, more accurate \mathbb{E}_π , better S
- Adaptive MCMC [Haario 2001, Andrieu 2006], but with nonlinear transformation to capture non-Gaussian structure

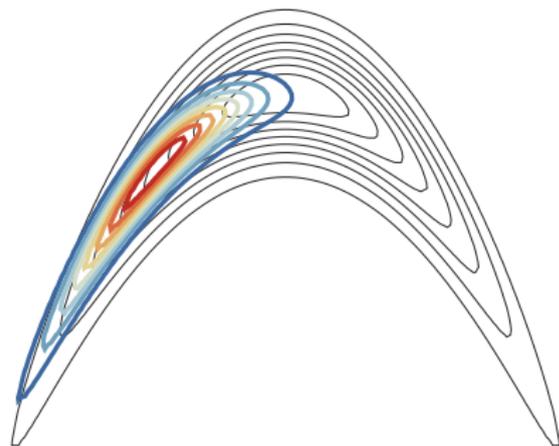


- **Ingredient #3: global proposals**

- If the map becomes sufficiently accurate, would like to avoid random-walk behavior



reference RW proposal

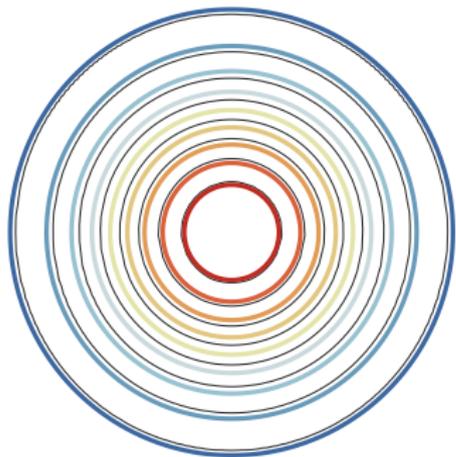


mapped RW proposal

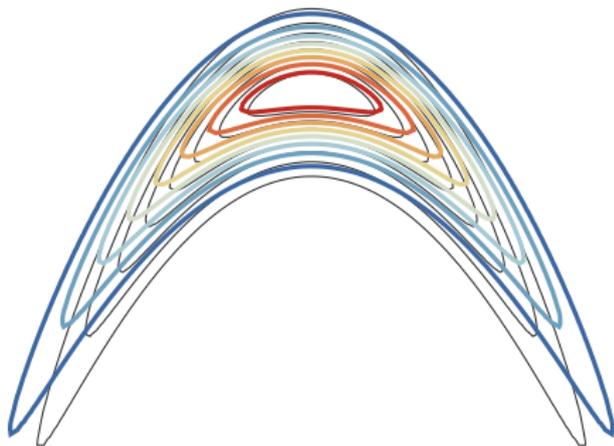
Map-accelerated MCMC

- **Ingredient #3: global proposals**

- If the map becomes sufficiently accurate, would like to avoid random-walk behavior



reference independence proposal



mapped independence proposal

- **Ingredient #3: global proposals**
 - If the map becomes sufficiently accurate, would like to avoid random-walk behavior
 - Solution: **delayed rejection** MCMC [Mira 2001]
 - First proposal = independent sample from η (global, more efficient); second proposal = random walk (local, more robust)
- Entire scheme is provably **ergodic** with respect to the exact posterior measure [Parno & M, *SIAM JUQ* 2018]
 - Requires enforcing some regularity conditions on maps, to preserve tail behavior of transformed target

Example: biological oxygen demand model

- ▶ Likelihood model:

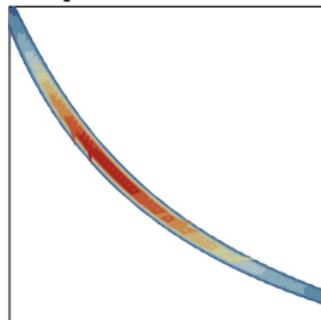
$$d = \theta_1(1 - \exp(-\theta_2 x)) + \epsilon$$
$$\epsilon \sim N(0, 2 \times 10^{-4})$$

- ▶ 20 noisy observations at

$$x = \left\{ \frac{5}{5}, \frac{6}{5}, \dots, \frac{25}{5} \right\}$$

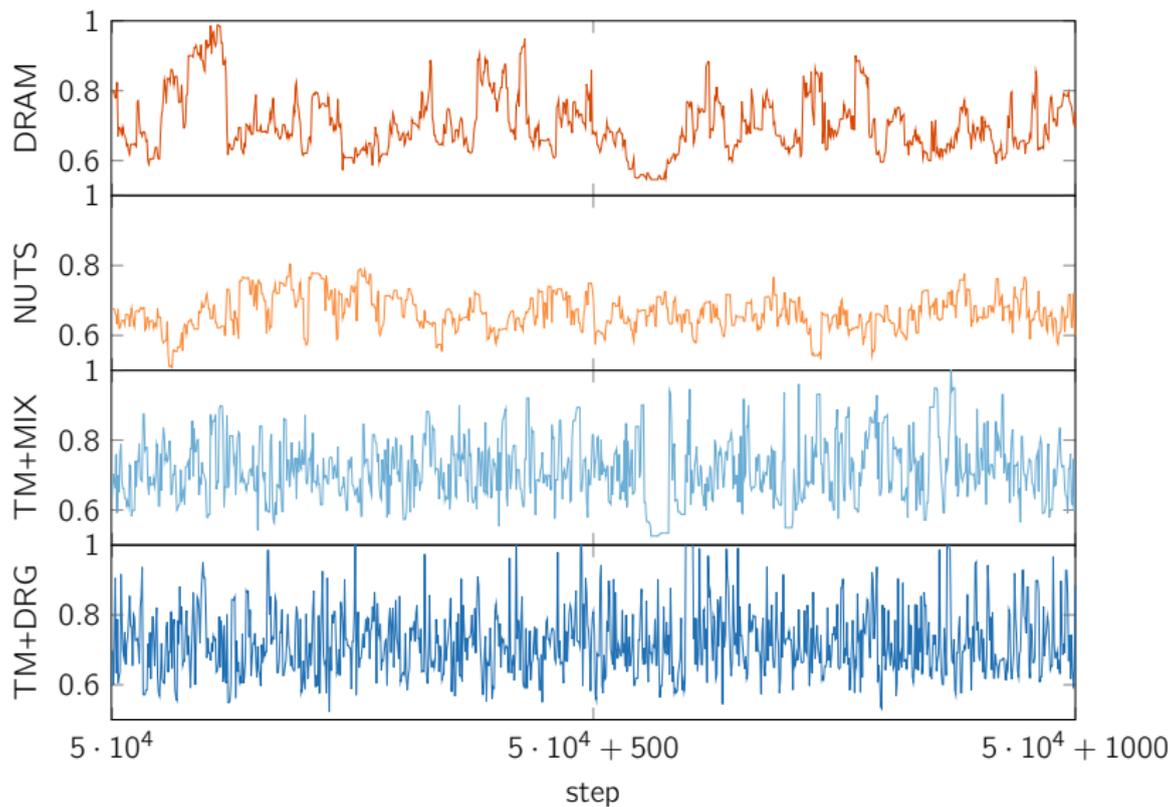
- ▶ Degree-three polynomial map

True posterior density

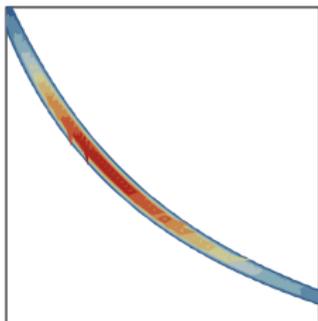


Results: MCMC chain

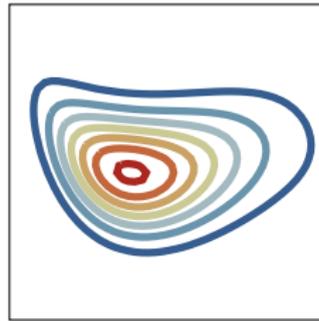
θ_1 component of MCMC chain



Transformed distribution



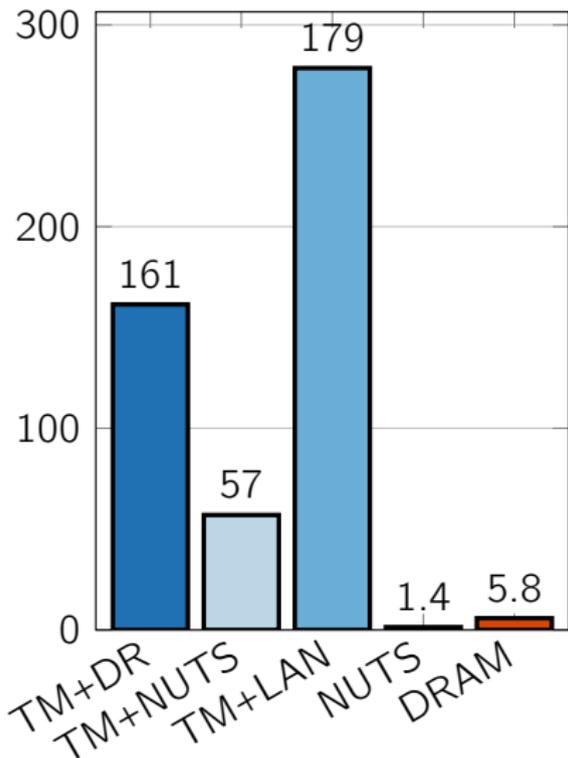
Original posterior π



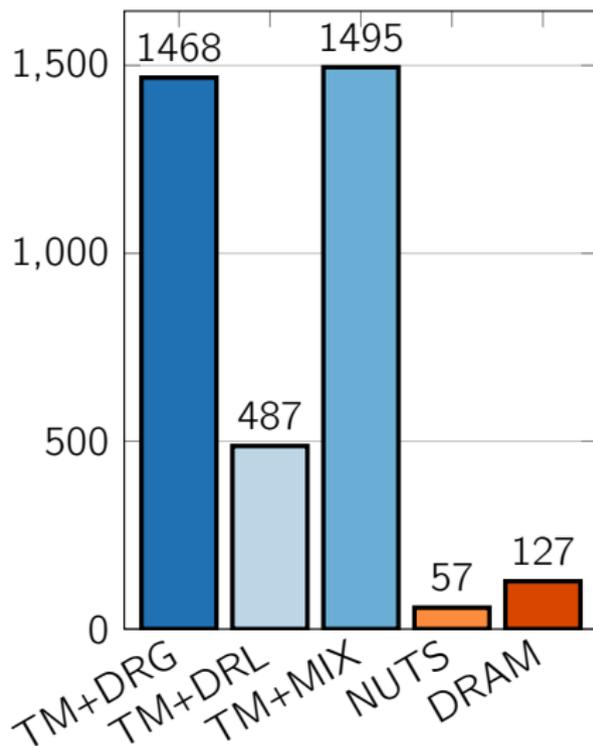
Pushforward of posterior through
learned map, $S_{\#}\pi$

Results: ESS per computational effort

ESS/(1,000 Evaluations) – θ_1



ESS/second – θ_1



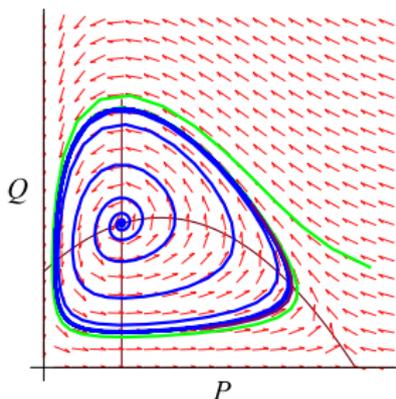
Example #2: predator-prey model

- ▶ Six parameter ODE population model

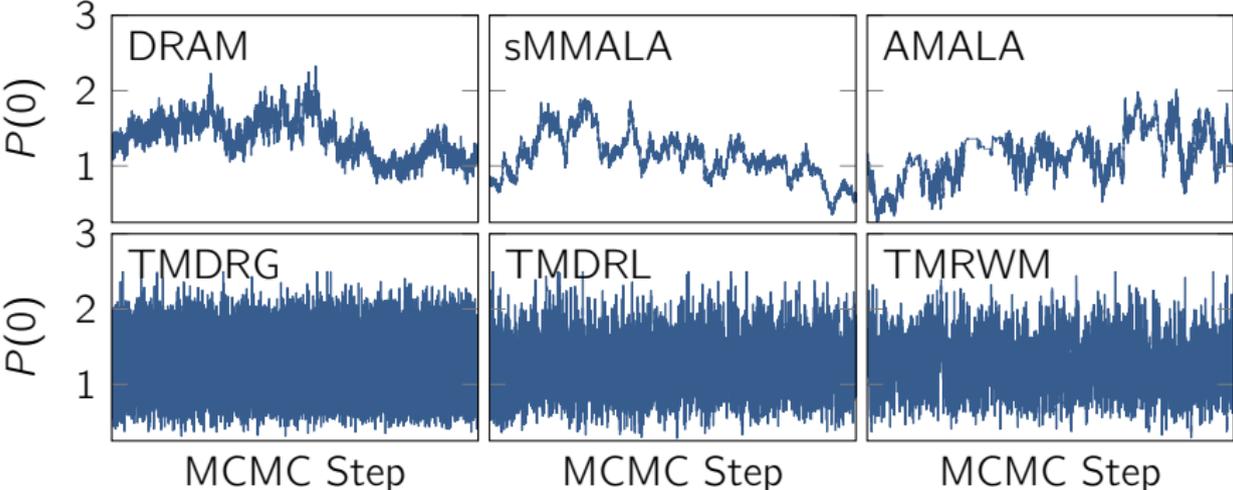
$$\frac{dP}{dt} = rP \left(1 - \frac{P}{K} \right) - s \frac{PQ}{a + P}$$

$$\frac{dQ}{dt} = u \frac{PQ}{a + P} - vQ$$

- ▶ Five noisy observations of both populations
- ▶ Infer 6 parameters + 2 initial values; uniform priors



Predator-prey model: chains



Example: maple sap dynamics model

- ▶ Coupled PDE system for ice, water, and gas locations [Ceseri & Stockie 2013]
- ▶ Measure gas pressure in vessel
- ▶ Infer 10 physical model parameters
- ▶ Very challenging posterior!

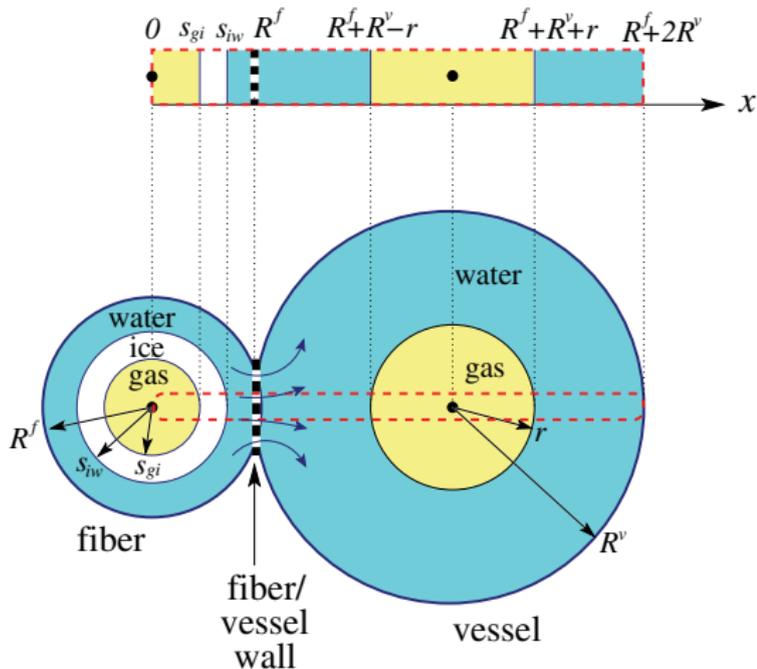
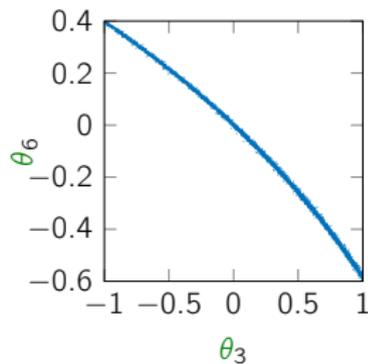
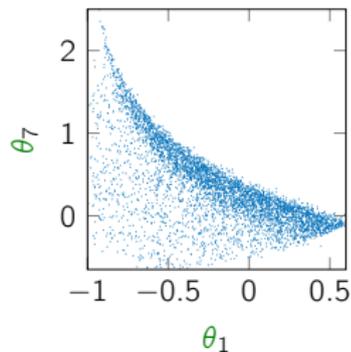
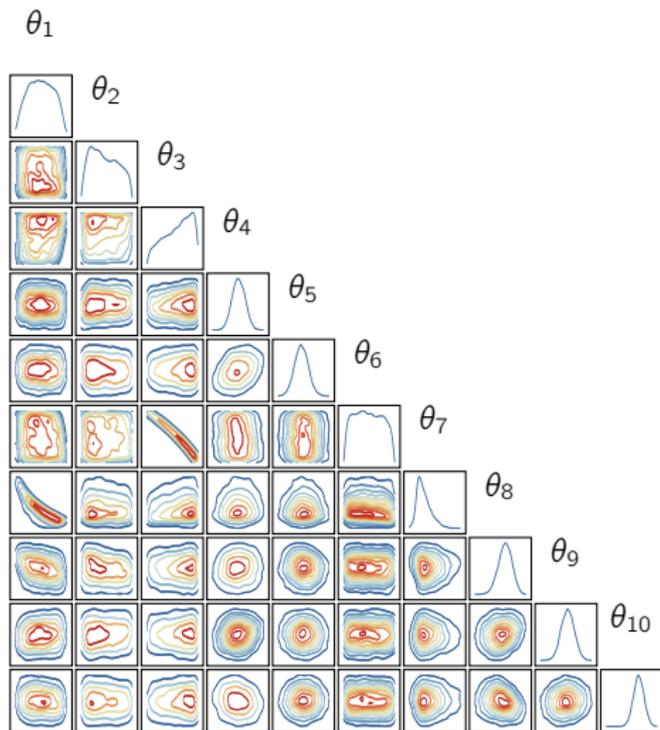


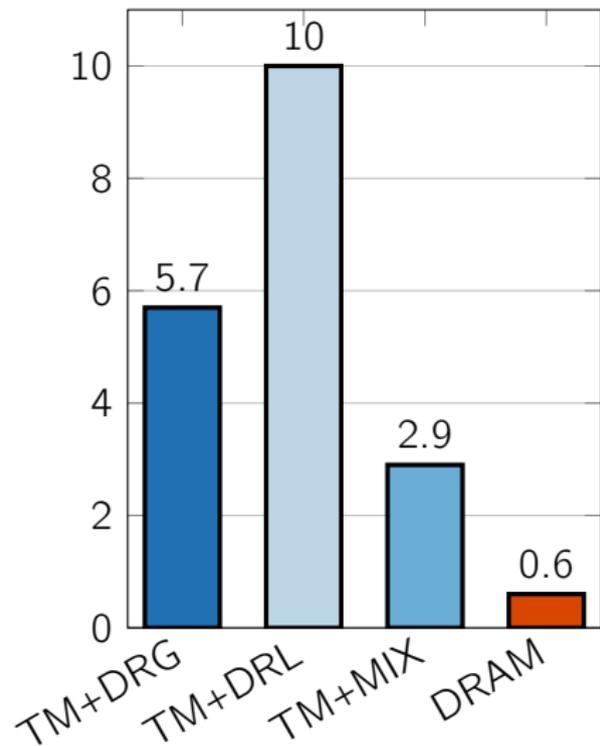
Image from *Ceseri and Stockie, 2013*

Maple posterior distribution

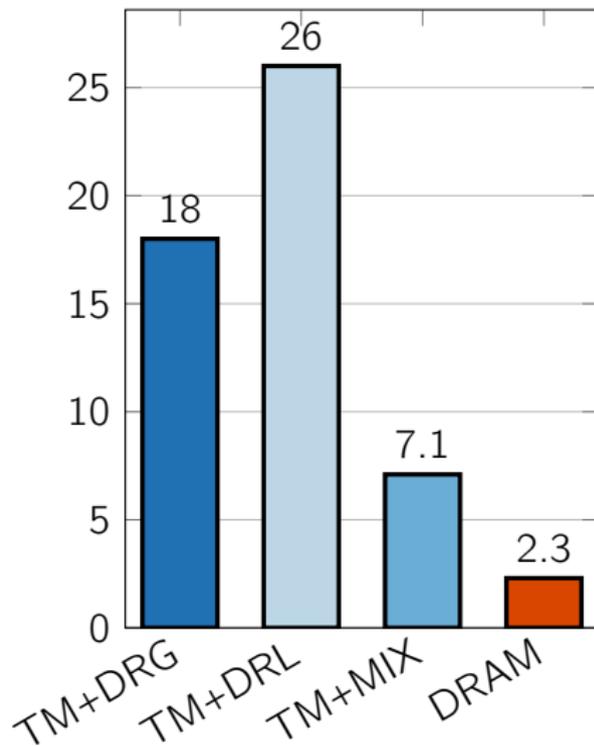


Results: ESS per computational effort

ESS/(10,000 Evaluations)



ESS/(1000 seconds)



Useful characteristics of the algorithm:

- ▶ Map construction is easily parallelizable
- ▶ Requires no gradients from posterior density

Generalizes many current MCMC techniques:

- ▶ Adaptive Metropolis: map enables **non-Gaussian proposals** *and* a natural mixing between local and global moves
- ▶ Manifold MCMC [Girolami & Calderhead 2011]: map also defines a Riemannian metric

Looking to higher dimensions: regularized estimation of S

For simplicity, consider map components $S^k(\mathbf{x}) = \sum_j \beta_j \psi_j(x_{1:k-1}) + \alpha_k x_k$

$$\hat{S}^k \in \arg \min_{S^k \in \mathcal{S}_{\Delta, k}^h} \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2} S^k(\mathbf{x}_i)^2 - \log \partial_k S^k(\mathbf{x}_i) \right) + \lambda_N \|\boldsymbol{\beta}\|_1$$

Assume sub-Gaussian π and basis functions $\psi_j(\mathbf{x})$

Theorem [BZM]

For polynomial maps of degree m with sparsity s , with high probability

$$\mathbb{E}_{\pi} \left[D_{KL} \left(\pi(\mathbf{x}_k | \mathbf{x}_{1:k-1}) \parallel \hat{S}_k^{\#} \eta \right) \right] \lesssim \sqrt{\frac{s^2 m \log k}{N}}$$

Takeaways

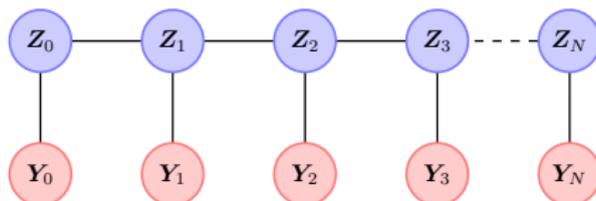
- ▶ Accurate estimation is feasible in high dimensions with $N \ll k$
- ▶ From factorization property of density, error in conditionals ensures

$$D_{KL}(\pi \parallel \hat{S}^{\#} \eta) \lesssim d \sqrt{\frac{s^2 m \log d}{N}}$$

Next topic: ensemble filtering via transport

► **Nonlinear/non-Gaussian** state-space model:

- Transition density $\pi_{\mathbf{z}_k|\mathbf{z}_{k-1}}$
- Observation density (likelihood) $\pi_{\mathbf{y}_k|\mathbf{z}_k}$

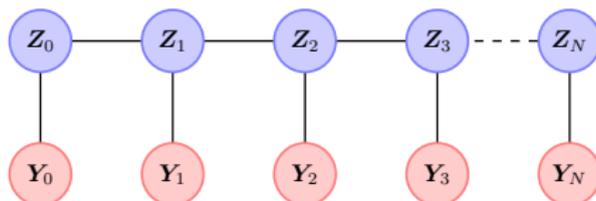


► Interested in **recursively** updating the **full Bayesian solution**:

$$\pi_{\mathbf{z}_{0:k} | \mathbf{y}_{0:k}} \rightarrow \pi_{\mathbf{z}_{0:k+1} | \mathbf{y}_{0:k+1}} \text{ (smoothing)}$$

► **Nonlinear/non-Gaussian** state-space model:

- Transition density $\pi_{\mathbf{z}_k | \mathbf{z}_{k-1}}$
- Observation density (likelihood) $\pi_{\mathbf{y}_k | \mathbf{z}_k}$



► Interested in **recursively** updating the **full Bayesian solution**:

$$\pi_{\mathbf{z}_{0:k} | \mathbf{y}_{0:k}} \rightarrow \pi_{\mathbf{z}_{0:k+1} | \mathbf{y}_{0:k+1}} \text{ (smoothing)}$$

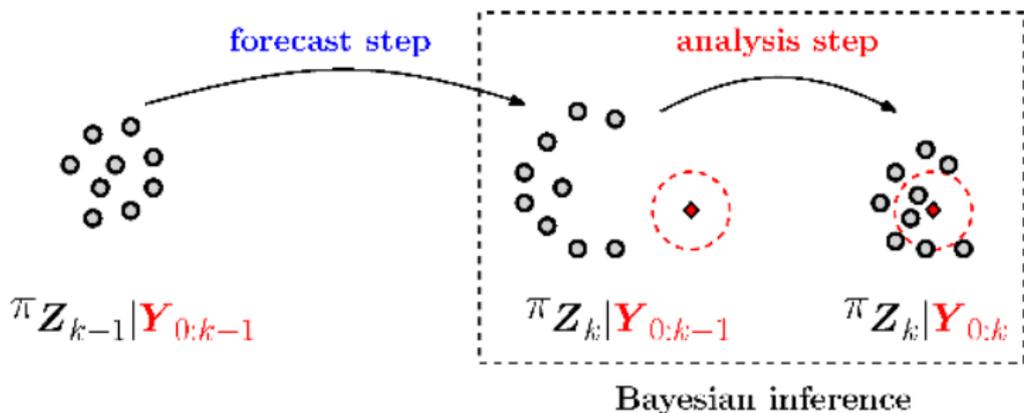
► Or focus on approximating the **filtering distribution**:

$$\pi_{\mathbf{z}_k | \mathbf{y}_{0:k}} \rightarrow \pi_{\mathbf{z}_{k+1} | \mathbf{y}_{0:k+1}} \text{ (marginals of the full Bayesian/smoothing solution)}$$

- ▶ Consider the filtering of state-space models with:
 - ① High-dimensional states
 - ② Challenging nonlinear dynamics (e.g., chaotic systems)
 - ③ Intractable transition kernels: can only obtain *forecast* samples, i.e., draws from $\pi_{\mathbf{z}_{k+1} | \mathbf{z}_k}$
 - ④ Limited model evaluations, e.g., small ensemble sizes
 - ⑤ Sparse and local observations in space/time
- ▶ These constraints reflect typical challenges faced in numerical weather prediction, geophysical data assimilation

Ensemble Kalman filter

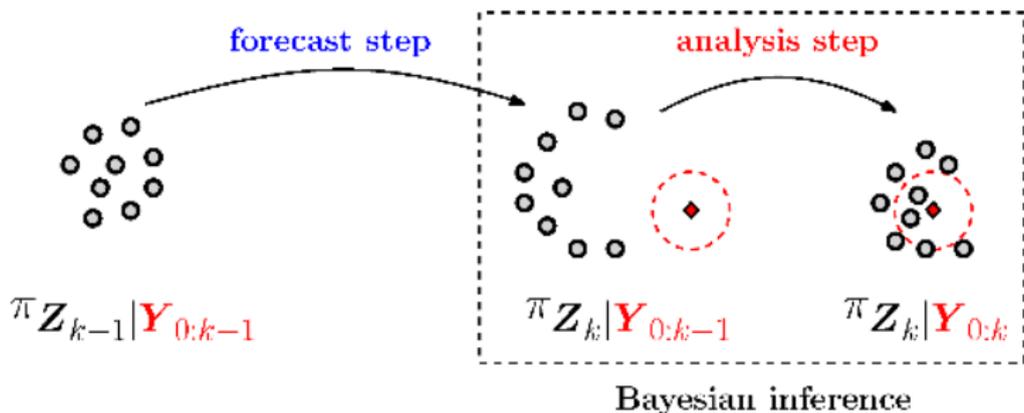
- ▶ State-of-the-art results (in terms of tracking) are typically obtained with the ensemble Kalman filter (EnKF)



- ▶ Move samples via an **affine** transformation; no weights or resampling!
- ▶ Yet ultimately **inconsistent**: does not converge to the true posterior

Ensemble Kalman filter

- ▶ State-of-the-art results (in terms of tracking) are typically obtained with the ensemble Kalman filter (EnKF)

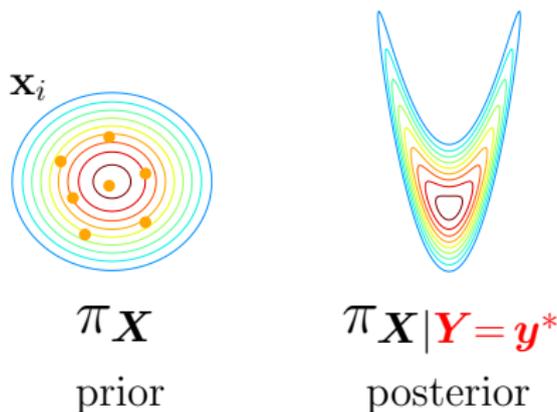


- ▶ Move samples via an **affine** transformation; no weights or resampling!
- ▶ Yet ultimately **inconsistent**: does not converge to the true posterior

Can we improve and *generalize* the EnKF, preserving scalability, via **nonlinear** transformations?

Assimilation step

At any assimilation time k , we have a Bayesian inference problem:



- ▶ $\pi_{\mathbf{X}}$ is the forecast distribution on \mathbb{R}^n
- ▶ $\pi_{\mathbf{Y}|\mathbf{X}}$ is the likelihood of the observations $\mathbf{Y} \in \mathbb{R}^d$
- ▶ $\pi_{\mathbf{X}|\mathbf{Y}=\mathbf{y}^*}$ is the filtering distribution for a realization \mathbf{y}^* of the data

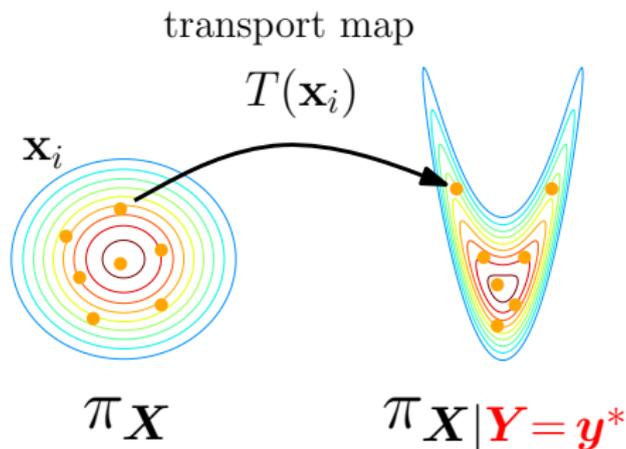
Goal: sample the posterior given only M prior samples $\mathbf{x}_1, \dots, \mathbf{x}_M$

Inference as a transportation of measures

- ▶ Seek a map T that pushes forward prior to posterior

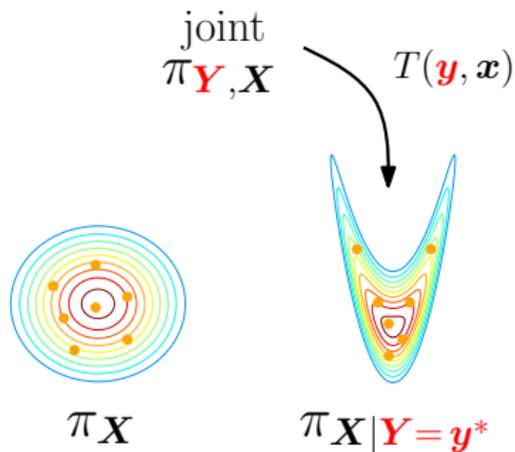
$$(\mathbf{x}_1, \dots, \mathbf{x}_M) \sim \pi_{\mathbf{X}} \implies (T(\mathbf{x}_1), \dots, T(\mathbf{x}_M)) \sim \pi_{\mathbf{X}|\mathbf{Y}=\mathbf{y}^*}$$

- ▶ The map induces a coupling between prior and posterior measures



How to construct a “good” coupling from very few prior samples?

Consider the joint distribution of state and observations

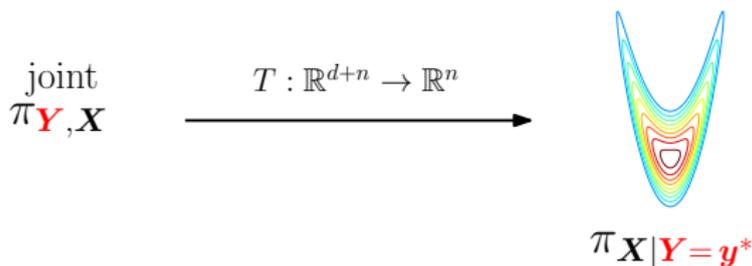


- ▶ Construct a map T from the joint distribution $\pi_{\mathbf{Y}, \mathbf{X}}$ to the posterior
- ▶ T can be computed via [convex optimization](#) given samples from $\pi_{\mathbf{Y}, \mathbf{X}}$
- ▶ Sample $\pi_{\mathbf{Y}, \mathbf{X}}$ using the forecast ensemble and the likelihood

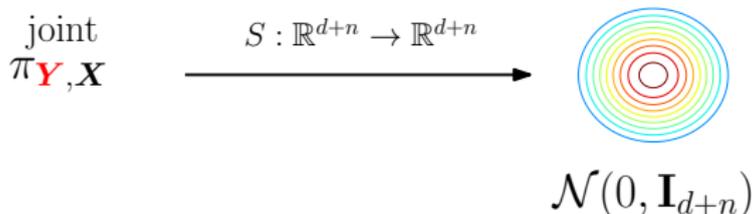
$$(\mathbf{y}_i, \mathbf{x}_i) \quad \mathbf{y}_i \sim \pi_{\mathbf{Y} | \mathbf{X} = \mathbf{x}_i}$$

- ▶ **Intuition:** a generalization of the “perturbed observation” EnKF

Couple the joint distribution with a standard normal



We can find T by computing a Knothe–Rosenblatt (KR) rearrangement S between $\pi_{\mathbf{Y}, \mathbf{X}}$ and $\mathcal{N}(0, \mathbf{I}_{d+n})$



► We will show how to derive T from S ...

- ▶ **Definition:** for any pair of absolutely continuous densities π, η on \mathbb{R}^m , there exists a unique **triangular** and **monotone** map $S : \mathbb{R}^m \rightarrow \mathbb{R}^m$ such that

$$S_{\#}\pi = \eta$$

- ▶ Triangular function (nonlinear generalization of a triangular matrix):

$$S(x_1, \dots, x_m) = \begin{bmatrix} S^1(x_1) \\ S^2(x_1, x_2) \\ \vdots \\ S^m(x_1, x_2, \dots, x_m) \end{bmatrix}$$

- ▶ Existence stems from general factorization properties of a density,

$$\pi = \pi_{\mathbf{X}_1} \pi_{\mathbf{X}_2|\mathbf{X}_1} \cdots \pi_{\mathbf{X}_m|\mathbf{X}_1, \dots, \mathbf{X}_{m-1}}$$

$$S(x_1, \dots, x_m) = \begin{bmatrix} S^1(x_1) \\ S^2(x_1, x_2) \\ \vdots \\ S^m(x_1, x_2, \dots, x_m) \end{bmatrix}$$

- ▶ Each component S^k links marginal conditionals of π and η
- ▶ For instance, if $\eta = \mathcal{N}(0, \mathbf{I})$, then for all $x_1, \dots, x_{k-1} \in \mathbb{R}^{k-1}$

$\xi \mapsto S^k(x_1, \dots, x_{k-1}, \xi)$ pushes $\pi_{\mathbf{x}_k | \mathbf{x}_{1:k-1}}(\xi | \mathbf{x}_{1:k-1})$ to $\mathcal{N}(0, 1)$

- ▶ **Simulate the conditional** $\pi_{\mathbf{x}_k | \mathbf{x}_{1:k-1}}$ by inverting a 1-D map $\xi \mapsto S^k(\mathbf{x}_{1:k-1}, \xi)$ at Gaussian samples (*need triangular structure*)

Filtering: the analysis map

- ▶ We are interested in the KR map S that pushes $\pi_{\mathbf{Y}, \mathbf{X}}$ to $\mathcal{N}(0, \mathbf{I}_{d+n})$
- ▶ The KR map immediately has a block structure

$$S(\mathbf{y}, \mathbf{x}) = \begin{bmatrix} S^{\mathbf{Y}}(\mathbf{y}) \\ S^{\mathbf{X}}(\mathbf{y}, \mathbf{x}) \end{bmatrix},$$

which suggests **two properties**:

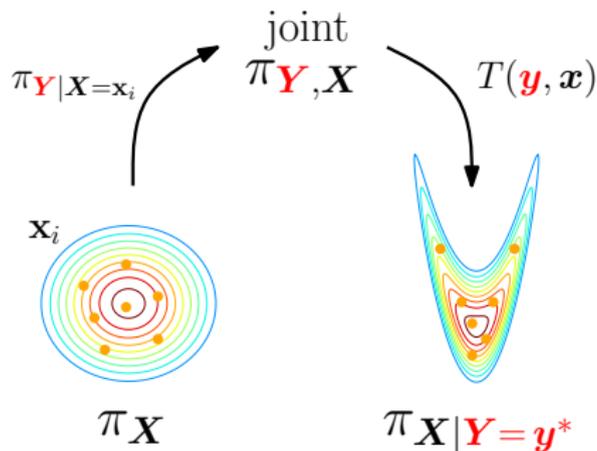
$S^{\mathbf{X}}$ pushes $\pi_{\mathbf{Y}, \mathbf{X}}$ to $\mathcal{N}(0, \mathbf{I}_n)$

$\xi \mapsto S^{\mathbf{X}}(\mathbf{y}^*, \xi)$ pushes $\pi_{\mathbf{X}|\mathbf{Y}=\mathbf{y}^*}$ to $\mathcal{N}(0, \mathbf{I}_n)$

- ▶ The **analysis map** that pushes $\pi_{\mathbf{Y}, \mathbf{X}}$ to $\pi_{\mathbf{X}|\mathbf{Y}=\mathbf{y}^*}$ is then given by

$$T(\mathbf{y}, \mathbf{x}) = S^{\mathbf{X}}(\mathbf{y}^*, \cdot)^{-1} \circ S^{\mathbf{X}}(\mathbf{y}, \mathbf{x})$$

A novel filtering algorithm with maps



Transport map ensemble filter

- 1 Compute forecast ensemble $\mathbf{x}_1, \dots, \mathbf{x}_M$
- 2 Generate samples $(\mathbf{y}_i, \mathbf{x}_i)$ from $\pi_{\mathbf{Y}, \mathbf{X}}$ with $\mathbf{y}_i \sim \pi_{\mathbf{Y}|\mathbf{X}=\mathbf{x}_i}$
- 3 Build an estimator \hat{T} of T
- 4 Compute analysis ensemble as $\mathbf{x}_i^a = \hat{T}(\mathbf{y}_i, \mathbf{x}_i)$ for $i = 1, \dots, M$

- ▶ Recall the form of S :

$$S(\mathbf{y}, \mathbf{x}) = \begin{bmatrix} S^Y(\mathbf{y}) \\ S^X(\mathbf{y}, \mathbf{x}) \end{bmatrix}, \quad S_{\#} \pi_{Y, X} = \mathcal{N}(0, \mathbf{I}_{d+n}).$$

- ▶ We propose the following estimator \hat{T} of T :

$$\hat{T}(\mathbf{y}, \mathbf{x}) = \hat{S}^X(\mathbf{y}^*, \cdot)^{-1} \circ \hat{S}^X(\mathbf{y}, \mathbf{x}),$$

where \hat{S} is a **maximum likelihood estimator** of S

Estimating the KR rearrangement from samples

Given samples $\mathbf{x}_1, \dots, \mathbf{x}_M$ from a distribution π on \mathbb{R}^{d+n} , estimate the KR rearrangement S that pushes forward π to $\mathcal{N}(0, \mathbf{I}_{d+n})$

- ▶ Constrained MLE for S

$$\hat{S} \in \arg \max_{S \in \mathcal{S}_{\Delta}^h} \frac{1}{M} \sum_{i=1}^M \log \underbrace{S_{\#}^{-1} \eta(\mathbf{x}_i)}_{\text{pullback}}, \quad \eta = \mathcal{N}(0, \mathbf{I}_n),$$

where \mathcal{S}_{Δ}^h is an approximation space for the rearrangement

- ▶ Each component \hat{S}^k of \hat{S} can be computed **separately**, via smooth **convex optimization**

$$\hat{S}^k \in \arg \min_{S^k \in \mathcal{S}_{\Delta,k}^h} \frac{1}{M} \sum_{i=1}^M \left(\frac{1}{2} S^k(\mathbf{x}_i)^2 - \log \partial_k S^k(\mathbf{x}_i) \right)$$

$$\hat{S}^k \in \arg \min_{S^k \in \mathcal{S}_{\Delta,k}^h} \frac{1}{M} \sum_{i=1}^M \left(\frac{1}{2} S^k(\mathbf{x}_i)^2 - \log \partial_k S^k(\mathbf{x}_i) \right)$$

- ▶ In general, convex optimization
- ▶ Optimization is not needed for nonlinear separable parameterizations of the form $\hat{S}^k(x_{1:k}) = \alpha x_k + g(x_{1:k-1})$ (just *linear regression*)
- ▶ **Connection to EnKF:** a linear parameterization of \hat{S}^k yields a particular form of EnKF with “perturbed observations”
- ▶ Choice of approximation space allows **control of the bias and variance** of \hat{S}
- ▶ Richer parameterizations yield less bias, but potentially higher variance

Strategy: depart gradually from the linear ansatz by introducing local nonlinearities + regularization

Example: Lorenz-63

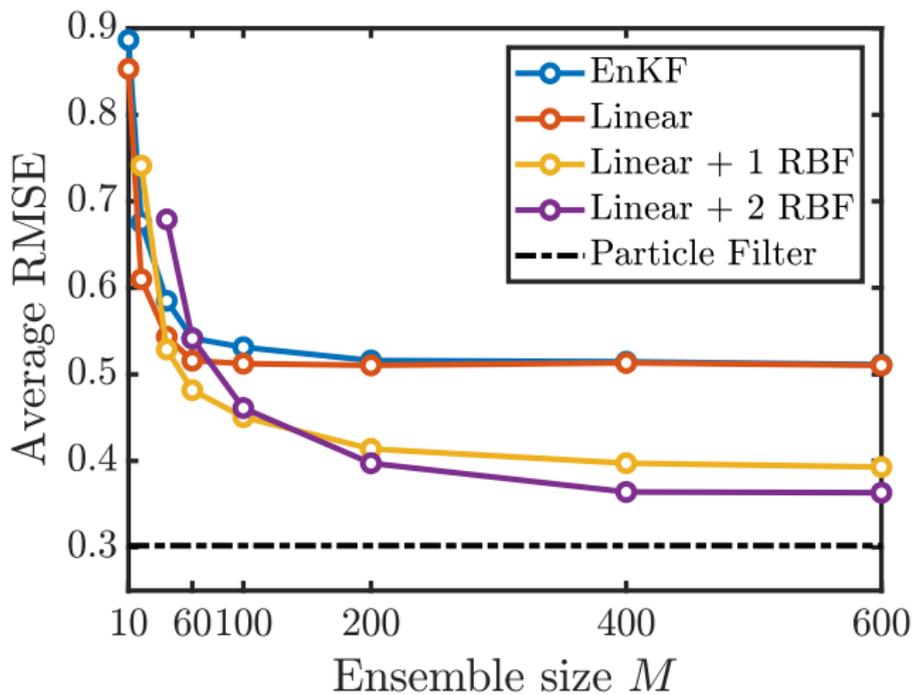
Simple example: three-dimensional Lorenz-63 system

$$\begin{aligned}\frac{dX_1}{dt} &= \sigma(X_2 - X_1), \\ \frac{dX_2}{dt} &= X_1(\rho - X_3) - X_2 \\ \frac{dX_3}{dt} &= X_1X_2 - \beta X_3\end{aligned}$$

- ▶ Chaotic setting: $\rho = 28$, $\sigma = 10$, $\beta = 8/3$
- ▶ Fully observed, with additive Gaussian observation noise $\mathcal{E}_j \sim \mathcal{N}(0, 2^2)$
- ▶ Assimilation interval $\Delta t = 0.1$
- ▶ Results computed over 2000 assimilation cycles, following spin-up
- ▶ **Map parameterizations:** $S^k(x_{1:k}) = \sum_{i \leq k} \Psi_i(x_i)$, with $\Psi_i = \text{linear} + \{\text{RBFs or sigmoids}\}$

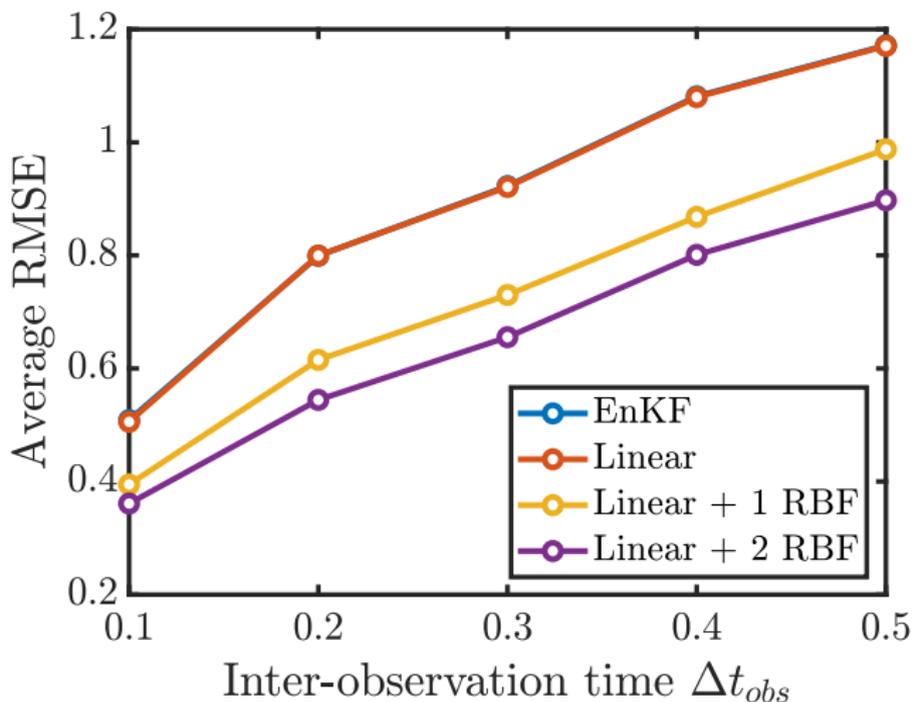
Example: Lorenz-63

Mean “tracking” error vs. ensemble size and choice of map



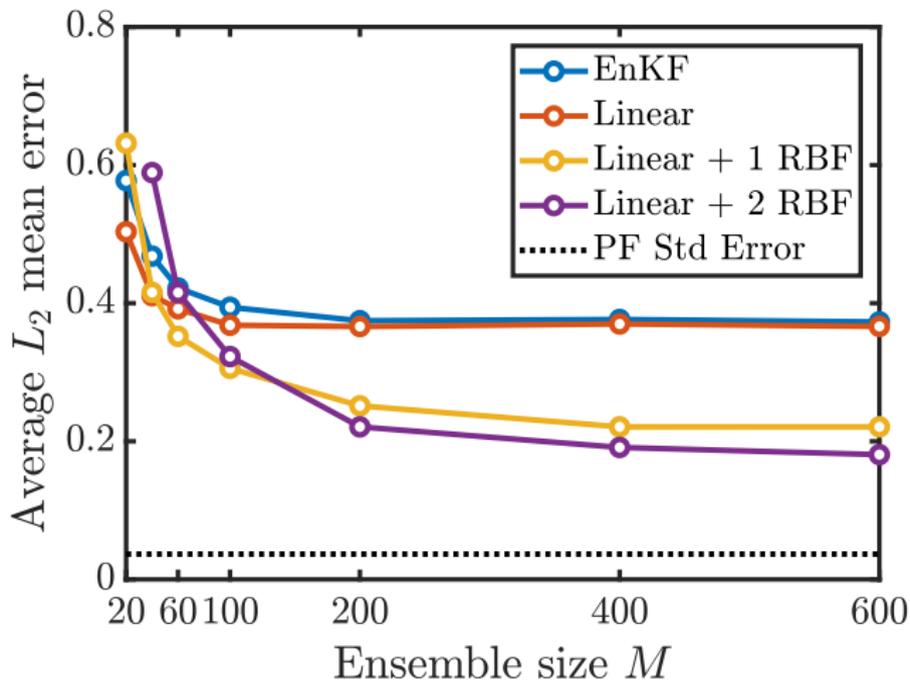
Example: Lorenz-63

How do $M \rightarrow \infty$ “plateaus” depend on assimilation interval?



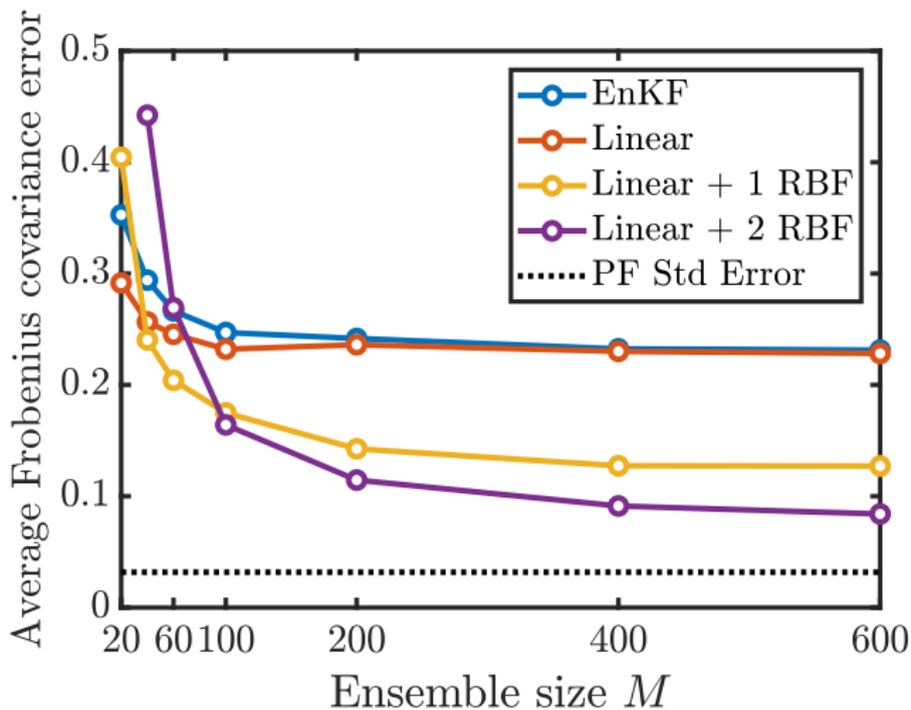
Example: Lorenz-63

What about comparison to the *true Bayesian solution*?



Example: Lorenz-63

What about comparison to the *true Bayesian solution*?



“Localize” the map in high dimensions

- ▶ Regularize the estimator \hat{S} of S by imposing **sparsity**, e.g.,

$$\hat{S}(x_1, \dots, x_4) = \begin{bmatrix} \hat{S}^1(x_1) \\ \hat{S}^2(x_1, x_2) \\ \hat{S}^3(x_2, x_3) \\ \hat{S}^4(x_3, x_4) \end{bmatrix}$$

- ▶ The sparsity of the k th component of S depends on the **sparsity of the marginal conditional** function $\pi_{\mathbf{x}_k | \mathbf{x}_{1:k-1}}(x_k | \mathbf{x}_{1:k-1})$
- ▶ **Localization heuristic:** let each \hat{S}^k depend on variables $(x_j)_{j < k}$ that are within a distance ℓ from x_k in state space. Estimate optimal ℓ offline
- ▶ Explicit link between sparsity of S and conditional independence in non-Gaussian graphical models

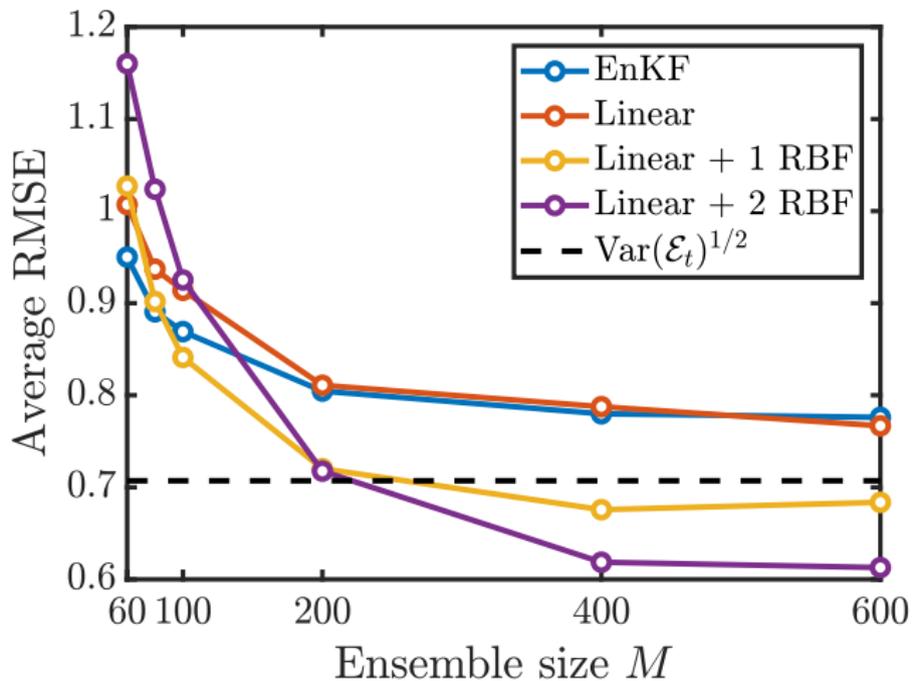
Lorenz-96 in chaotic regime (40-dimensional state)

- ▶ A **hard** test-case configuration [Bengtsson et al. 2003]:

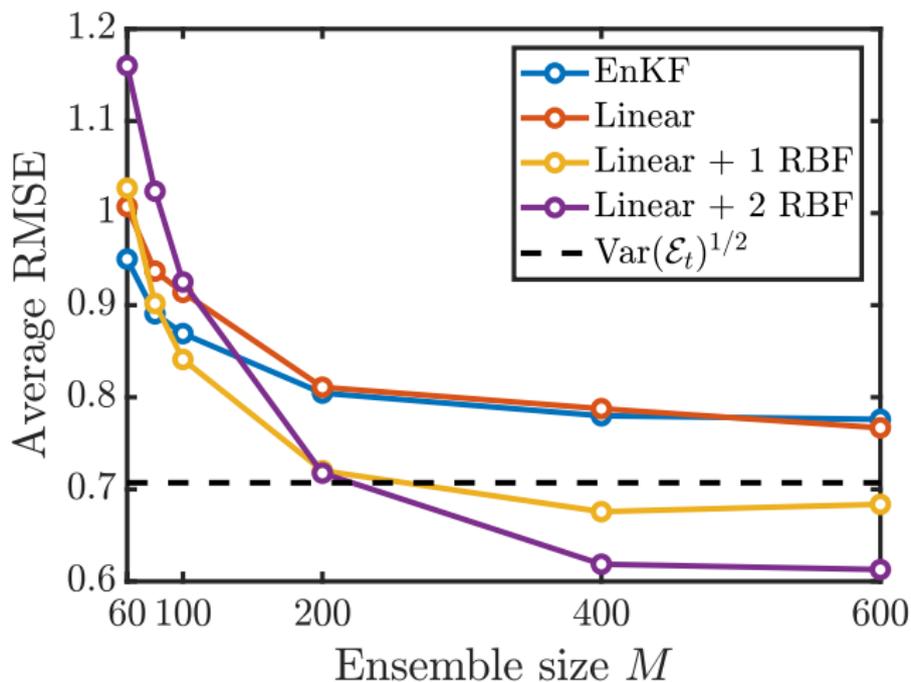
$$\begin{aligned}\frac{d\mathbf{X}_j}{dt} &= (\mathbf{X}_{j+1} - \mathbf{X}_{j-2})\mathbf{X}_{j-1} - \mathbf{X}_j + F, & j = 1, \dots, 40 \\ \mathbf{Y}_j &= \mathbf{X}_j + \mathcal{E}_j, & j = 1, 3, 5 \dots, 39\end{aligned}$$

- ▶ $F = 8$ (chaotic) and $\mathcal{E}_j \sim \mathcal{N}(0, 0.5)$ (**small noise for PF**)
- ▶ Time between observations: $\Delta_{\text{obs}} = 0.4$ (**large**)
- ▶ Results computed over 2000 assimilation cycles, following spin-up

Lorenz-96: "hard" case

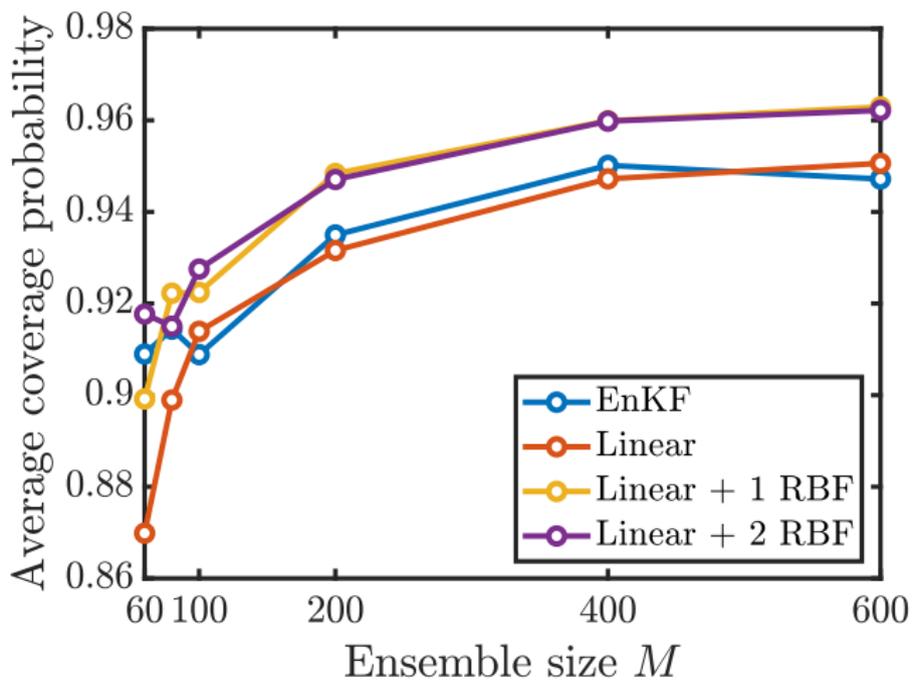


Lorenz-96: "hard" case



- ▶ The nonlinear filter is $\approx 25\%$ more accurate in RMSE than EnKF

Lorenz-96: "hard" case

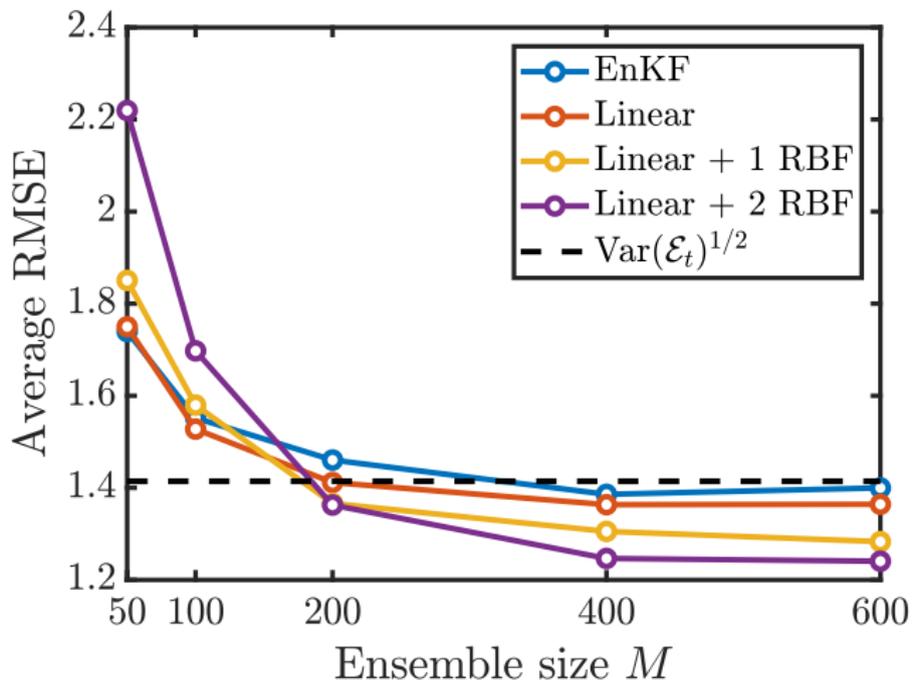


- ▶ A heavy-tailed noise configuration:

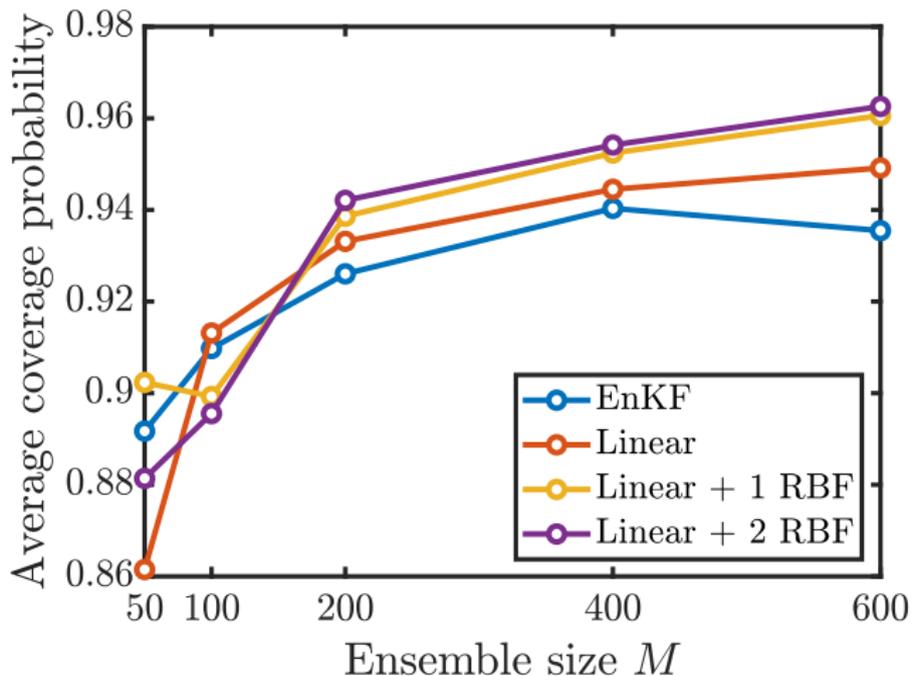
$$\begin{aligned}\frac{d\mathbf{X}_j}{dt} &= (\mathbf{X}_{j+1} - \mathbf{X}_{j-2})\mathbf{X}_{j-1} - \mathbf{X}_j + F, & j = 1, \dots, 40 \\ \mathbf{Y}_j &= \mathbf{X}_j + \mathcal{E}_j, & j = 1, 5, 9, 13, \dots, 37\end{aligned}$$

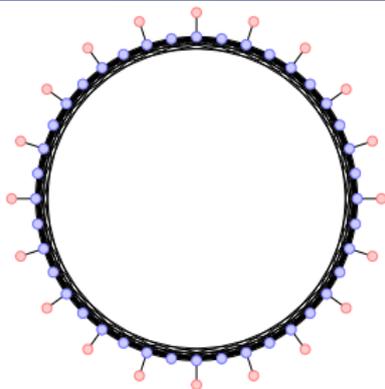
- ▶ $F = 8$ (chaotic) and $\mathcal{E}_j \sim \text{Laplace}(\lambda = 1)$
- ▶ Time between observations: $\Delta_{\text{obs}} = 0.1$
- ▶ Results computed over 2000 assimilation cycles, following spin-up

Lorenz-96: non-Gaussian noise



Lorenz-96: non-Gaussian noise





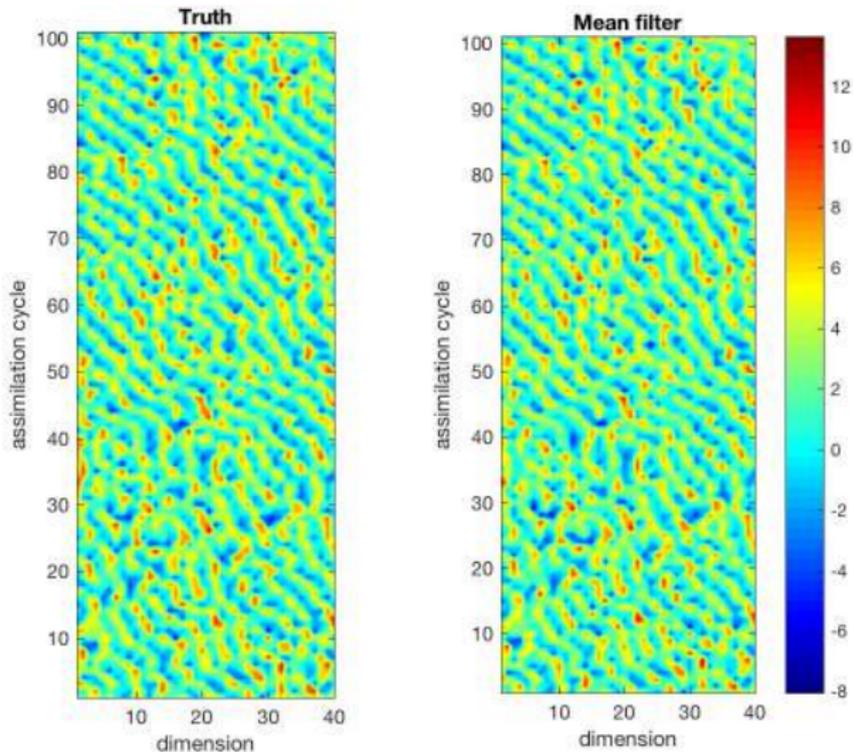
- ▶ Observations were assimilated one at a time
- ▶ Impose sparsity of the map with a 5-way interaction model (*above*)
- ▶ Separable and nonlinear parameterization of each component

$$\widehat{S}^k(x_{j_1}, \dots, x_{j_p}, x_k) = \psi(x_{j_1}) + \dots + \psi(x_{j_p}) + \widetilde{\psi}(x_k),$$

where $\psi(x) = a_0 + a_1 \cdot x + \sum_{i>1} a_i \exp(-(x - c_i)^2/\sigma)$.

- ▶ Much **more general** parameterizations are of course possible

Lorenz-96: tracking performance of the filter



- ▶ Simple and localized nonlinearities have significant impact!

- ▶ Nonlinear generalization of the EnKF: move the ensemble members via local nonlinear transport maps, *no weights or degeneracy*
- ▶ Learn non-Gaussian features via nonlinear continuous transport and *convex optimization*
- ▶ Choice of map basis and **sparsity** provide regularization (e.g., *localization*)

- ▶ Nonlinear generalization of the EnKF: move the ensemble members via local nonlinear transport maps, *no weights or degeneracy*
- ▶ Learn non-Gaussian features via nonlinear continuous transport and *convex optimization*
- ▶ Choice of map basis and **sparsity** provide regularization (e.g., *localization*)
- ▶ In principle, filter is consistent as \mathcal{S}_{Δ}^h is enriched and $M \rightarrow \infty$, but a careful *error analysis* is needed!
- ▶ What is a good or even optimal choice of \mathcal{S}_{Δ}^h for any fixed ensemble size M ?
- ▶ Can regularization penalties (e.g., ℓ_1) help identify sparse structure, and/or learn sparse maps from few samples?

Regularized estimation of S

For simplicity, consider map components $S^k(\mathbf{x}) = \sum_j \beta_j \psi_j(x_{1:k-1}) + \alpha_k x_k$

$$\hat{S}^k \in \arg \min_{S^k \in \mathcal{S}_{\Delta, k}^h} \frac{1}{N} \sum_{i=1}^N \left(\frac{1}{2} S^k(\mathbf{x}_i)^2 - \log \partial_k S^k(\mathbf{x}_i) \right) + \lambda_N \|\boldsymbol{\beta}\|_1$$

Assume sub-Gaussian π and basis functions $\psi_j(\mathbf{x})$

Theorem [BZM]

For polynomial maps of degree m with sparsity s , with high probability

$$\mathbb{E}_{\pi} \left[D_{KL} \left(\pi(\mathbf{x}_k | \mathbf{x}_{1:k-1}) \parallel \hat{S}_k^{\#} \eta \right) \right] \lesssim \sqrt{\frac{s^2 m \log k}{N}}$$

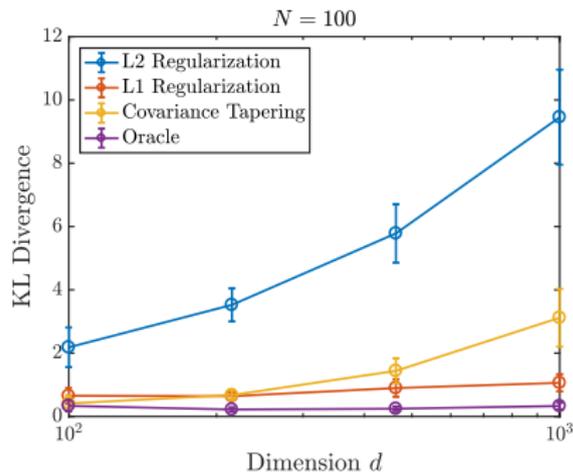
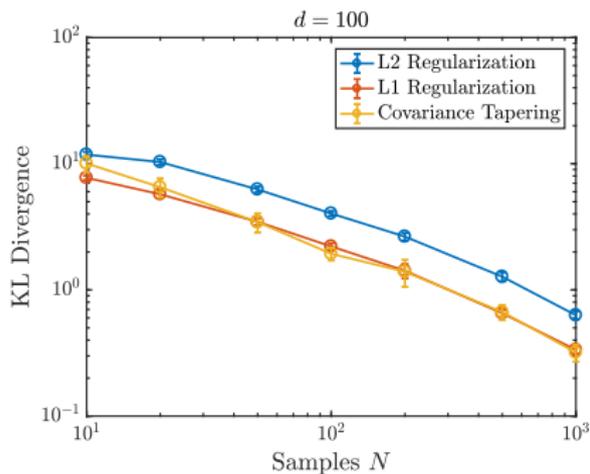
Takeaways

- ▶ Accurate estimation is feasible in high dimensions with $N \ll k$
- ▶ From factorization property of density, error in conditionals ensures

$$D_{KL}(\pi \parallel \hat{S}^{\#} \eta) \lesssim d \sqrt{\frac{s^2 m \log d}{N}}$$

Linear–Gaussian problem

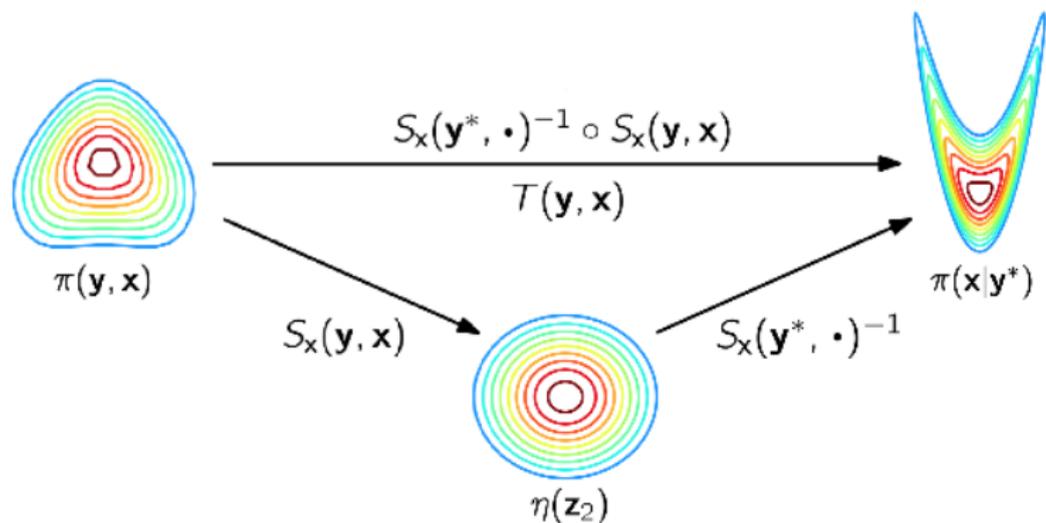
- ▶ Prior: $\mathbf{X} \sim \mathcal{N}(\mu, \Sigma_{pr})$ with exponential covariance
- ▶ Likelihood: Local observations $\mathbf{Y} = \mathbf{H}\mathbf{X} + \mathcal{E}$ with $\mathcal{E} \sim \mathcal{N}(0, I)$



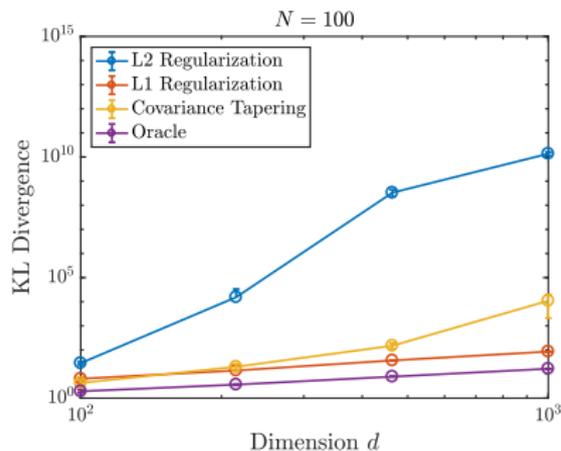
Takeaway

- ▶ Learning sparse prior-to-posterior map **matches oracle scaling**

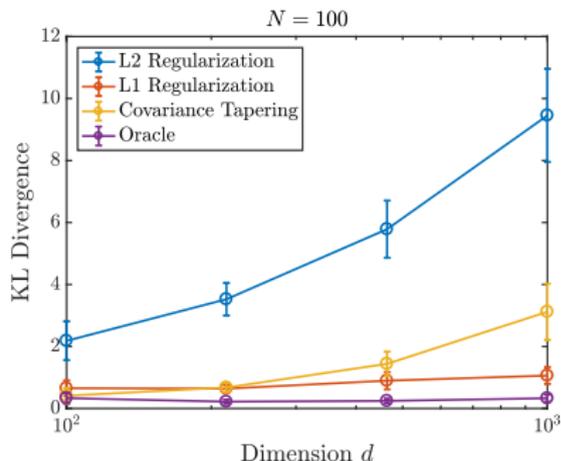
Compare two approaches for posterior sampling



Compare two approaches for posterior sampling



$$\mathbf{X}|\mathbf{y}^* \sim \widehat{\mathcal{S}}_{\mathbf{X}}(\mathbf{y}^*, \cdot)_{\#}^{-1} \eta$$



$$\mathbf{X}|\mathbf{y}^* \sim \widehat{\mathcal{T}}_{\#} \pi_{\mathbf{y}, \mathbf{x}} \text{ for } \widehat{\mathcal{T}} = \widehat{\mathcal{S}}_{\mathbf{X}}(\mathbf{y}^*, \cdot)^{-1} \circ \widehat{\mathcal{S}}_{\mathbf{X}}(\cdot, \cdot)$$

- ▶ Propagating forecast through **composed maps** has lower error
- ▶ This is in fact a **general** approach to likelihood-free inference/ABC

- ▶ There is also a “square root” version of the nonlinear ensemble filter
- ▶ Continuous-time formulations?
- ▶ Nonlinear ensemble smoothers
- ▶ Open questions about estimation and regularization of continuous transport maps:
 - ▶ How to choose and *adapt* approximation space/basis to the forecast ensemble?
 - ▶ Properties of the estimator \hat{T} , e.g., consistency, sample size requirements and scaling
 - ▶ Other forms of low-dimensional structure and regularization
- ▶ Applications to inference in “likelihood-free” settings

- ▶ A. Spantini, R. Baptista, Y. Marzouk. “Coupling techniques for nonlinear ensemble filtering.” arXiv:1907.00389.
- ▶ D. Bigoni, O. Zahm, A. Spantini, Y. Marzouk. “Greedy inference with layers of lazy maps.” arXiv:1906.00031.
- ▶ O. Zahm, T. Cui, K. Law, A. Spantini, Y. Marzouk. “Certified dimension reduction in nonlinear Bayesian inverse problems.” arXiv:1807.03712.
- ▶ A. Spantini, D. Bigoni, Y. Marzouk. “Inference via low-dimensional couplings.” *JMLR* 19(66): 1–71, 2018.
- ▶ M. Parno, Y. Marzouk, “Transport map accelerated Markov chain Monte Carlo.” *SIAM JUQ* 6: 645–682, 2018.
- ▶ G. Detomasso, T. Cui, A. Spantini, Y. Marzouk, R. Scheichl, “A Stein variational Newton method.” NeurIPS 2018.
- ▶ R. Morrison, R. Baptista, Y. Marzouk. “Beyond normality: learning sparse probabilistic graphical models in the non-Gaussian setting.” NeurIPS 2017.
- ▶ Y. Marzouk, T. Moselhy, M. Parno, A. Spantini, “An introduction to sampling via measure transport.” *Handbook of Uncertainty Quantification*, R. Ghanem, D. Higdon, H. Owhadi, eds. Springer (2016). arXiv:1602.05023.
- ▶ **General python code at <http://transportmaps.mit.edu>**