

# Singular Likelihoods to Prevent Particle Filter Collapse

Gregor Robinson, Ian Grooms, William Kleiber  
*Applied Mathematics, CU Boulder*

2018 April 16



Particle filtering is an importance sampling approach to the Bayesian estimation problem for dynamical systems.

Main points:

- ▶ Particle filtering basics
- ▶ Why particle filters are cool in theory
- ▶ Mode of failure in high-dimensional practice
- ▶ Our approach to reduce the curse of dimensionality
- ▶ Example on a toy model

# IMPORTANCE SAMPLING

To approximate 'target' distribution  $\mu$ :

1. Draw ensemble of  $N_e$  'particles' from some 'proposal' distribution (that's easy to simulate)
2. assign each particle a positive weight  $0 < w^{(i)} < 1$  such that  $\sum_i w^{(i)} = 1$  and specially rigged so that

$$\mu \approx \mu_{N_e} := \sum_{i=1}^{N_e} w^{(i)} \delta(\mathbf{x} - \mathbf{x}^{(i)})$$

Simple computations can give us the right weights so that for any  $\varphi \in \mathcal{C}_b$ ,  $\langle \mu_{N_e}, \varphi \rangle \rightarrow \langle \mu, \varphi \rangle$  as  $N_e \rightarrow \infty$ .  
(Weak convergence.)

# SEQUENTIAL IMPORTANCE SAMPLING (SIS)

1. Start with posterior (analysis) importance sample (ensemble) at assimilation timestep  $k - 1$
2. Propagate by proposal kernel  $\mathbb{P} \left( \mathbf{x}_k^{(i)} | \mathbf{x}_{0:k-1}^{(i)}, \mathbf{y}_k \right)$
3. Reweight so ensemble  $\left\{ \left( \mathbf{x}_k^{(i)}, w_k^{(i)} \right) \right\}$  becomes valid importance approximation of posterior at assimilation timestep  $k$ :

$$w_k^{(i)} \propto w_{k-1}^{(i)} \frac{\mathbb{P} \left( \mathbf{y}_k | \mathbf{x}_k^{(i)} \right) \mathbb{P} \left( \mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)} \right)}{\mathbb{P} \left( \mathbf{x}_k^{(i)} | \mathbf{x}_{0:k-1}^{(i)}, \mathbf{y}_k \right)}.$$

with likelihood  $\mathbb{P} \left( \mathbf{y}_k | \mathbf{x}_k \right)$  and transition prior  $\mathbb{P} \left( \mathbf{x}_k^{(i)} | \mathbf{x}_{k-1}^{(i)} \right)$

## HERDING CATS (PARTICLES)

Particles can wander from observations.

Closest particle then has much higher likelihood, picking up almost all the weight.

Specifically:

$$w_k^{(i)} \propto \mathbb{P}(\mathbf{y}_k | \mathbf{x}_k^{(i)}) w_{k-1}^{(i)}$$

If one particle has  $w^{(i)} \approx 1$  then the UQ is bad. Define

$$\text{Effective Sample Size} = \text{ESS} = \left( \sum_i (w^{(i)})^2 \right)^{-1}$$

$$1 \leq \text{ESS} \leq N_e.$$

# COLLAPSE



The PF has 'collapsed' when  $ESS \ll N$ .

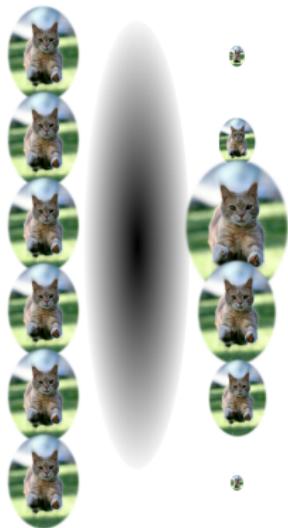
To fix, resample: eliminate particles with small weights, replicate ones with large weights, resetting all weights to  $1/N$  (gross simplification).

# Prior



Prior

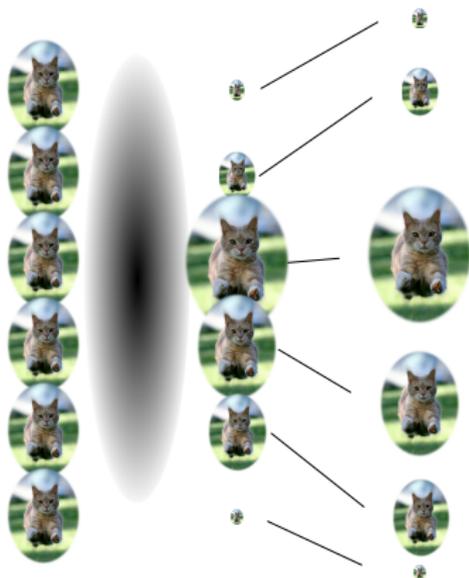
Observation



Prior

Observation

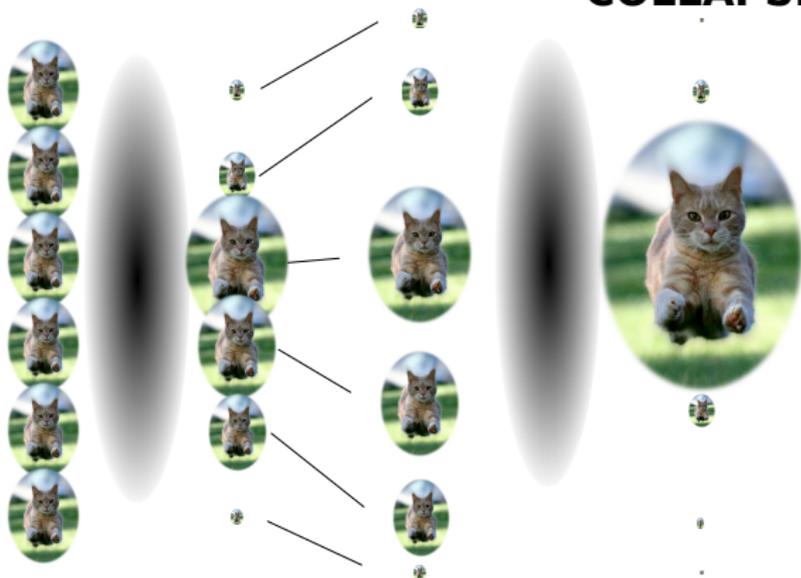
Dynamic



Prior

Observation

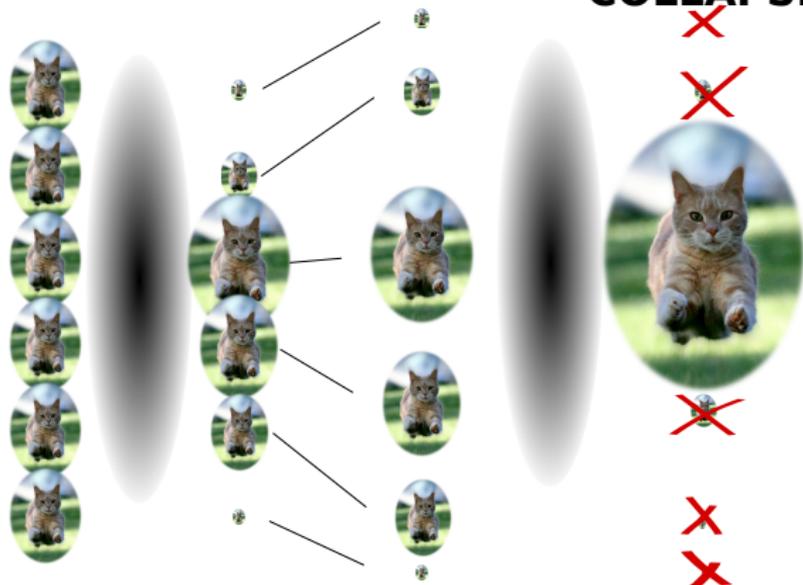
Dynamic

**COLLAPSE!**

Prior

Observation

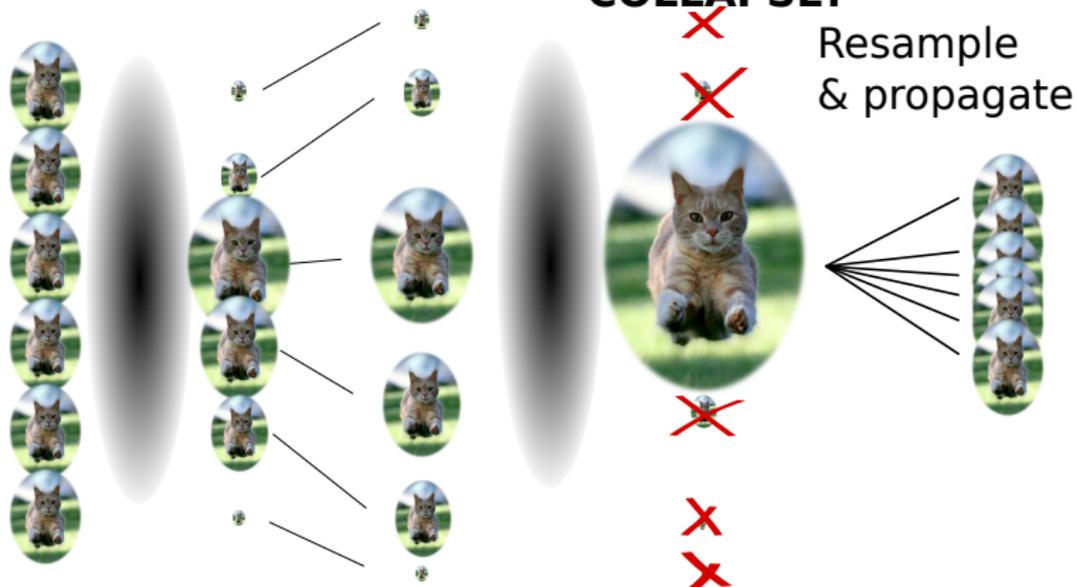
Dynamic

**COLLAPSE!**

Prior

Observation

Dynamic



# SEQUENTIAL IMPORTANCE SAMPLING WITH RESAMPLING (SIR)

SIR weakly converges<sup>1</sup> to the Bayesian posterior as  $N_e \rightarrow \infty$  with extremely permissive constraints on prior, transition kernel, and likelihood.

That's useful for UQ of non-Gaussian and nonlinear problems!

The SIR PF does not work well for high-dimensional problems.

---

<sup>1</sup>D. Crisan

Snyder et al.<sup>2</sup> consider SIR with linear Gaussian obs:

$$\mathbf{y} = \mathbf{H}\mathbf{x} + \boldsymbol{\xi}, \quad \boldsymbol{\xi} \sim \mathcal{N}(0, \mathbf{R})$$

To avoid collapse, they show you need  $N_e \sim \exp\{\tau^2/2\}$ , where

$$\tau^2 = \sum_k \lambda_k^2 \left( \frac{3}{2} \lambda_k^2 + 1 \right)$$

Where, for the ‘standard proposal’,

$$\mathbf{P} = \text{Cov}[\mathbf{R}^{-1/2}\mathbf{H}\mathbf{x}]$$

and  $\lambda_k^2$  are the eigenvalues of  $\mathbf{P}$ .

---

<sup>2</sup>Bengtsson, Bickel, & Li 2008; Snyder, Bengtsson, Bickel, & Anderson 2008

∃ optimal proposal and related ways to improve required number of particles by incorporating observations into proposal:

- ▶ Snyder, Bengtsson, & Morzfeld (2015) showed that the number of required particles is exponential in the ‘effective dimension’ even in best case scenario of optimal proposal<sup>3</sup>
- ▶ Chorin, Morzfeld, & Tu (2009–present) develop an ‘implicit’ PF that approximates the optimal proposal
- ▶ Ades & van Leeuwen (2009–present) develop an ‘equivalent weights’ PF that is related to (but not equivalent to) the optimal proposal

---

<sup>3</sup>See also Agapiou, Papaspiliopoulos, Sanz-Alonso, and Stuart (2017) for precise non-asymptotic results

# FIDDLING WITH THE LIKELIHOOD

In practice (e.g. for satellite obs) the likelihood is rarely known precisely. Common to assume spatially-uncorrelated observation errors for serial processing in EnKF.

If we can choose  $\mathbf{R}$  to increase EnKF computational efficiency, why not choose it to avoid PF collapse?

# FIDDLING WITH THE LIKELIHOOD

Revisit Snyder estimate: need  $N_e \sim \exp\{\tau^2/2\}$ , where

$$\tau^2 = \sum_k \lambda_k^2 \left( \frac{3}{2} \lambda_k^2 + 1 \right)$$

$$\mathbf{P} = \text{Cov}[\mathbf{R}^{-1/2} \mathbf{H}\mathbf{x}]$$

Idea: increase variance at scales that don't matter that much.

For geophysical forecast, viscous damping means that small scales don't matter much.

Consider the case where  $\mathbf{H} = \mathbf{I}$ , and  $\mathbf{R}$  and  $\text{Cov}[\mathbf{x}]$  are simultaneously diagonalizable. Then

$$\lambda_k^2 = \frac{\sigma_k^2}{\gamma_k^2}$$

where  $\sigma_k^2$  are eigvals of  $\text{Cov}[\mathbf{x}]$ , and  $\gamma_k^2$  are eigvals of  $\mathbf{R}$ .

$\mathbf{x}$  is an ordinary random field – realizations are, e.g., continuously differentiable – so the spectrum must decay.

$$\lim_{k \rightarrow \infty} \sigma_k^2 = 0$$

$\sigma_k \rightarrow 0$  is good news: even if the grid is refined ad infinitum,  $\tau^2$  might converge to a finite value

$$\tau^2 = \sum_k \lambda_k^2 \left( \frac{3}{2} \lambda_k^2 + 1 \right), \quad \lambda_k^2 = \frac{\sigma_k^2}{\gamma_k^2}$$

But presumably the observation error is also a continuous field, so  $\gamma_k \rightarrow 0$  also. This is BAD for  $\tau^2$ .

But it's widely accepted to use a spatially-uncorrelated obs error model, which has constant  $\gamma_k$ . This is good for  $\tau^2$ .

Why not use  $\gamma_k^2 \rightarrow \infty$ , which is even better for  $\tau^2$ ?

$\sigma_k \rightarrow 0$  is good news: even if the grid is refined ad infinitum,  $\tau^2$  might converge to a finite value

$$\tau^2 = \sum_k \lambda_k^2 \left( \frac{3}{2} \lambda_k^2 + 1 \right), \quad \lambda_k^2 = \frac{\sigma_k^2}{\gamma_k^2}$$

But presumably the observation error is also a continuous field, so  $\gamma_k \rightarrow 0$  also. This is BAD for  $\tau^2$ .

But it's widely accepted to use a spatially-uncorrelated obs error model, which has constant  $\gamma_k$ . This is good for  $\tau^2$ .

Why not use  $\gamma_k^2 \rightarrow \infty$ , which is even better for  $\tau^2$ ?

**Changing the likelihood changes the posterior.**

**Changing the likelihood changes the posterior.** Yes, but...

- ▶ a white likelihood already does that (sometimes)
- ▶ and you needn't change the posterior *much*:

Consider again a fully-Gaussian case with  $\mathbf{H} = \mathbf{I}$  and simultaneously-diagonalizable  $\mathbf{R}$  and  $\text{Cov}[\mathbf{x}]$ .

The spectrum of the posterior is

$$\frac{\sigma_k^2 \gamma_k^2}{\sigma_k^2 + \gamma_k^2}$$

At small scales (large  $k$ ), the posterior variance is small, regardless of how you choose  $\gamma_k$  (because  $\sigma_k$  is small).

As long as  $\gamma_k$  is correct at large scales, the posterior will be correct at large scales.

## CONNECTION TO SMOOTHING

Our approach is like truncation/projection onto a large-scale subspace, but with a gradual cutoff.

Assuming jagged observation errors is equivalent to smoothing the observations by applying  $\mathbf{R}^{-1/2}$ , and then assuming uncorrelated errors

$$\hat{\mathbf{y}} = \mathbf{R}^{-1/2}\mathbf{y} = \mathbf{R}^{-1/2}\mathbf{H}\mathbf{x} + \hat{\boldsymbol{\xi}}, \quad \hat{\boldsymbol{\xi}} \sim \mathcal{N}(0, \mathbf{I})$$

Note that this is *not* equivalent to assuming uncorrelated obs error and then smoothing.

## EXAMPLE

As a the linear SPDE

$$\frac{du}{dt} = \left( -b - c \frac{d}{dx} + \nu \frac{d^2}{dx^2} \right) u + F_t, \quad (1)$$

in a  $2\pi$ -periodic domain where the forcing is Gaussian, white in time, with spatial spectrum  $(1 + |k|)^{-1}$ .

2048 Fourier modes,  $b = 1$ ,  $c = 2\pi$ , and  $\nu = 1/9$ .

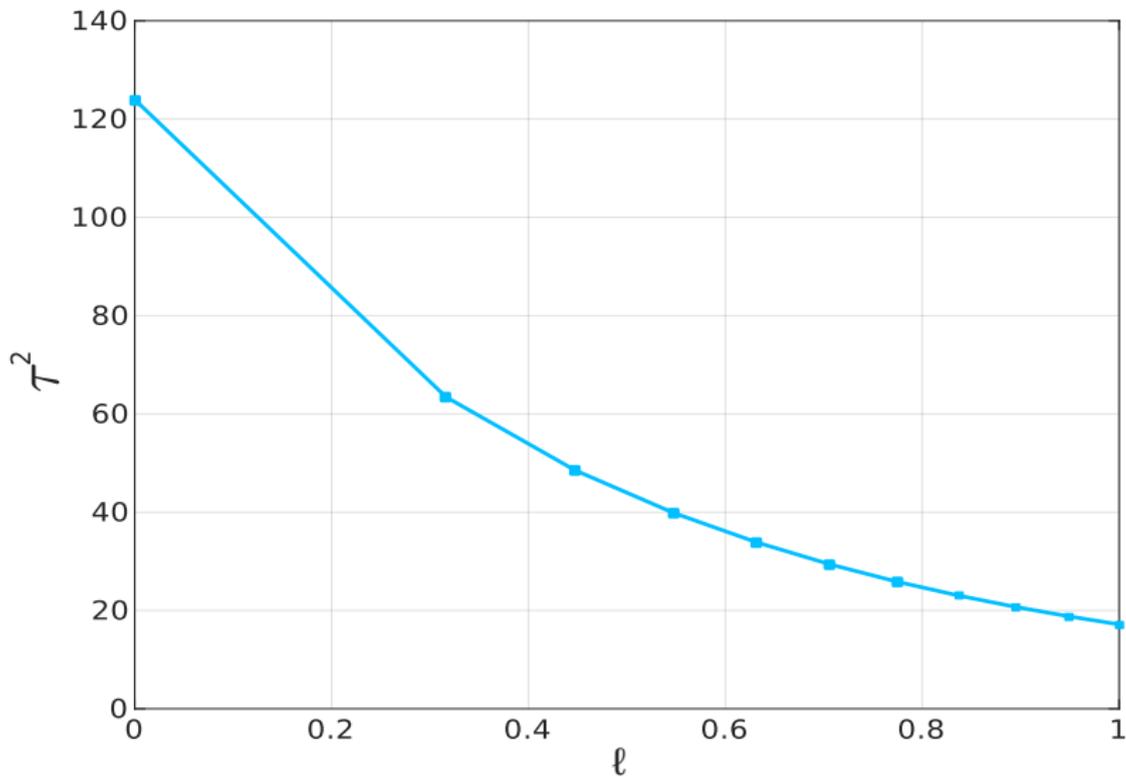
Observations are taken on a regular spatial grid, at discrete times; the true obs errors are smooth (Gaussian, zero-mean).

We compare posteriors using true  $\mathbf{R}$ ,  $\mathbf{R} = \gamma^2 \mathbf{I}$ , or a second-order finite-difference discretization of  $\gamma^2(1 - \ell^2 \partial_x^2)$ .

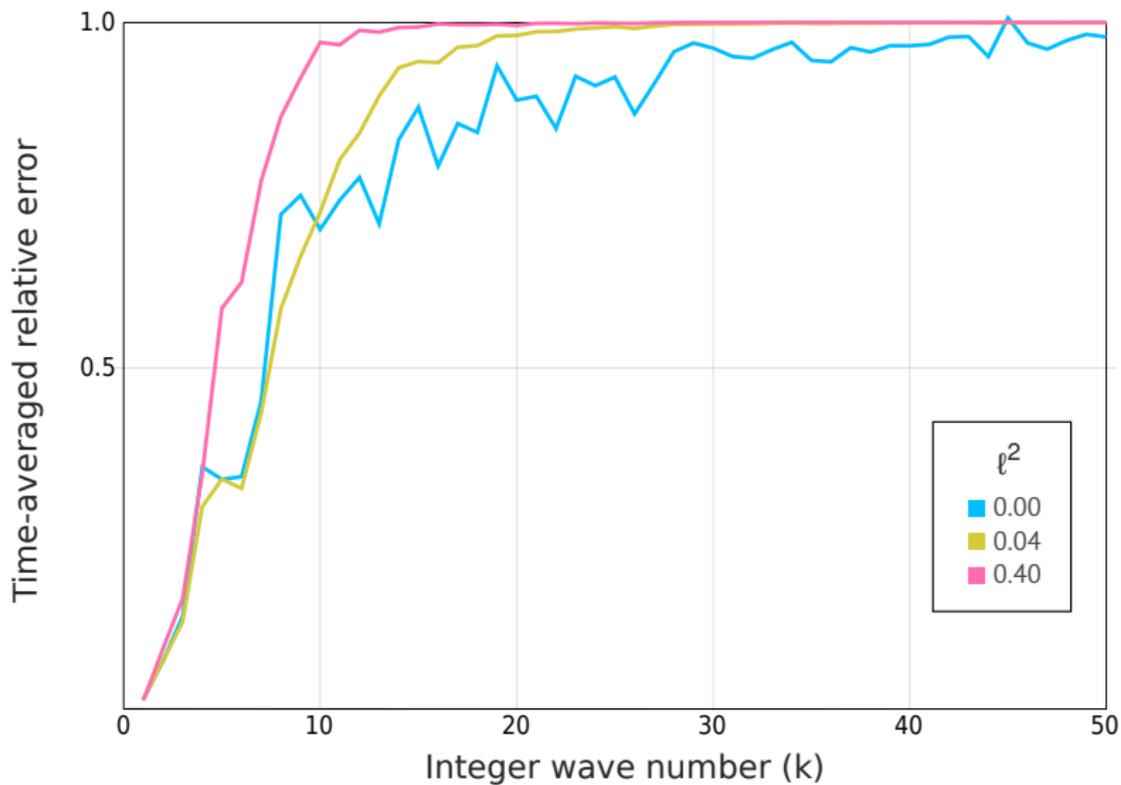
First we compute the exact solution of the filtering problem using the Kalman filter. The filter covariance converges exponentially to a steady state; we use this information to compute  $\tau^2$ :

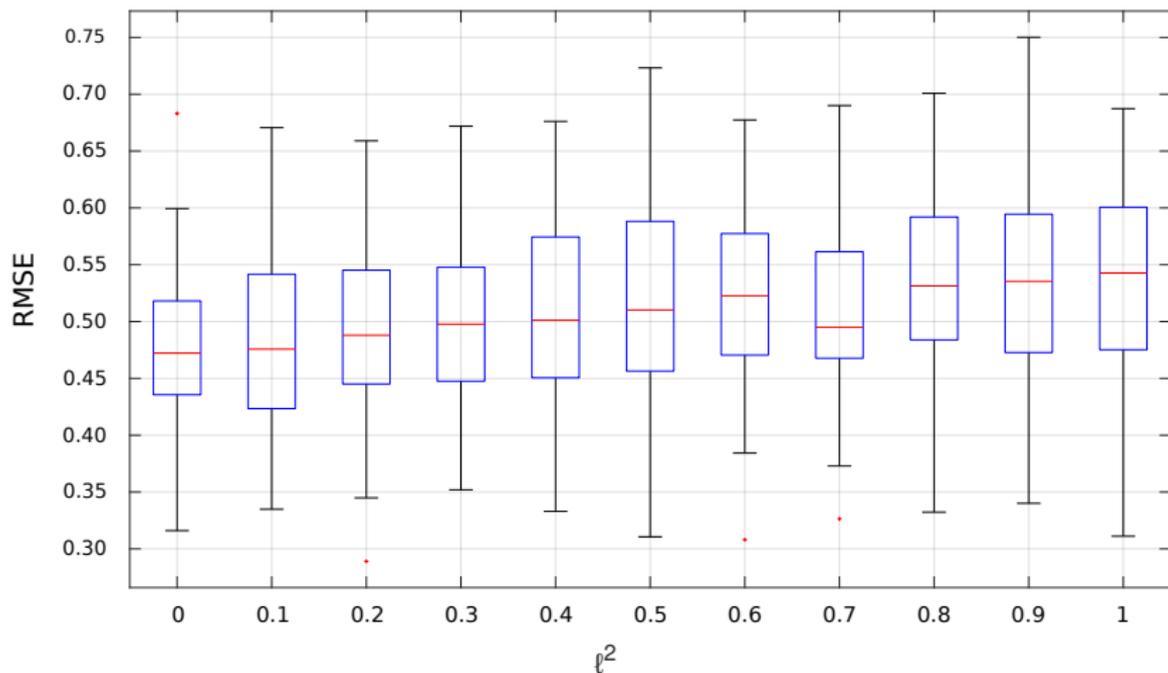
**64 obs:**

- ▶ True  $\mathbf{R}$ :  $\tau^2 = 87, N \approx 10^{19}$
- ▶  $\mathbf{R} = \gamma^2 \mathbf{I}$ :  $\tau^2 = 162, N \approx 10^{35}$
- ▶ Generalized,  $\ell^2 = 2$ :  $\tau^2 = 16, N \approx 3000$

**$\tau^2$  versus GRF length scale**

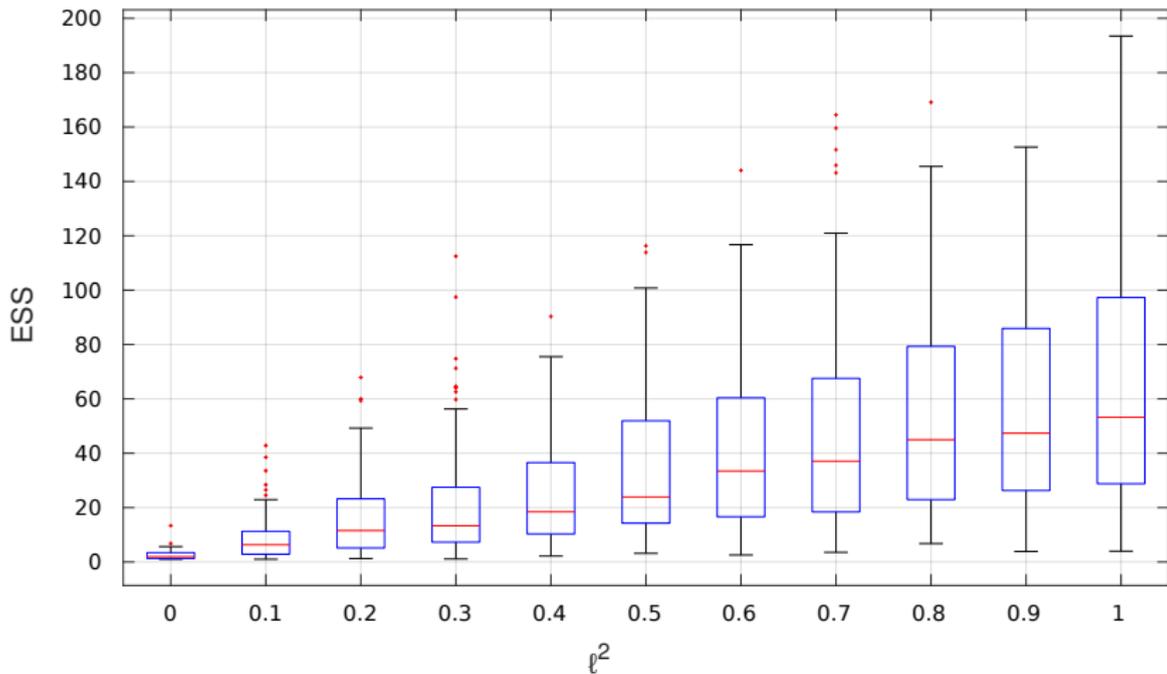
### Relative error of KF mean in Fourier space



**RMSE vs GRF length scale squared**

RMSE minimally suffers.

### Effective sample size vs GRF length scale squared

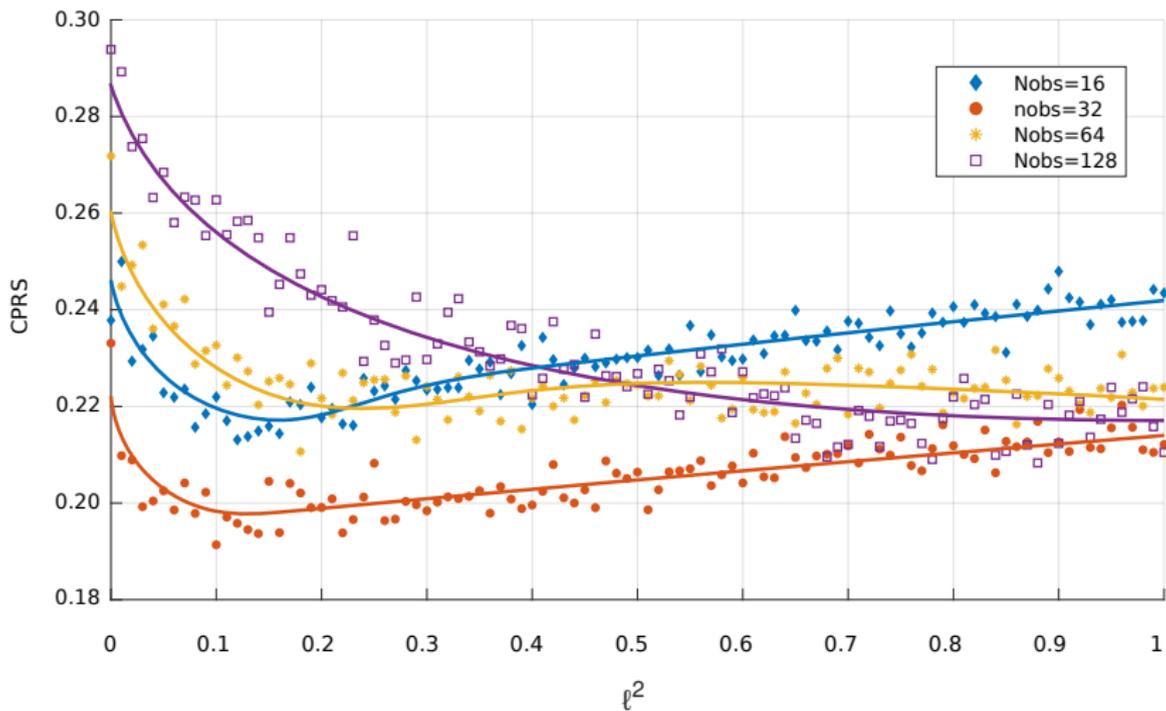


Ensemble size = 400

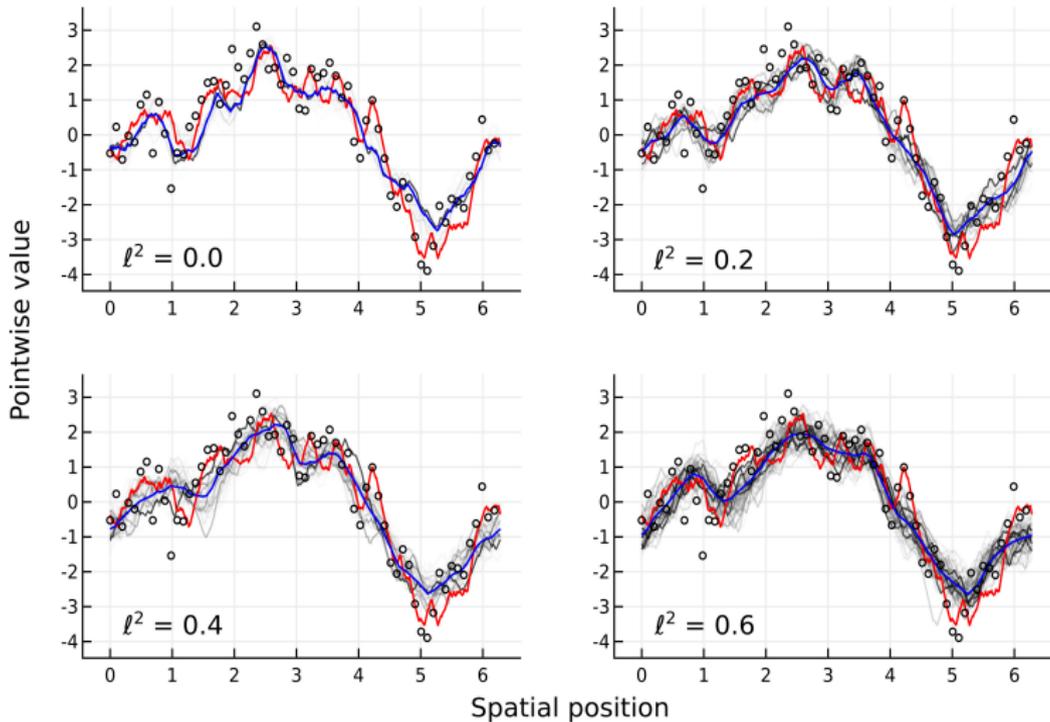
# QUANTIFYING UNCERTAINTY QUANTIFICATION

We use the Continuous Ranked Probability Score (CRPS) to quantify the skill of our full PF estimate.

## CRPS vs GRF length scale squared



### Weighted samples for different $\ell^2$



## WHY USE $\gamma^2(1 - \ell^2\partial_x^2)$ ?

There's a link between PDEs and Gaussian random fields: the discretization of an elliptic, self-adjoint PDE approximates the precision matrix of a random field. <sup>4</sup>

We reverse this:

a discretization of a self-adjoint elliptic PDE approximates the covariance matrix of a jagged random field.

Admits covariance structure that is exploitable for computational efficiency (including multiresolution) even on nonuniform grids.

---

<sup>4</sup>E.g. Lindgren, Rue, & Lindström, J R Stat Soc 2011

## Conclusions & Future Directions

- ▶ Using a generalized random field for the obs error model can reduce incidence of PF collapse.
- ▶ The price to pay is that the posterior is only accurate on large scales; in practice that might be OK.
- ▶ We plan to continue development of approaches to discretizing  $\mathbf{R}$ , esp. for scattered obs, and to apply to real meteorological data.
- ▶ Our approach probably won't be a silver bullet, but can be combined with implicit sampling/optimal-proposal and with localization.

Thanks to Jeff Anderson, Greg Beylkin & Chris Snyder for helpful discussions.

# THE END!

Questions?

`gregor.robinson@colorado.edu`

`@precompact`