

# Hierarchical Off-Diagonal Low-Rank Approximation for Hessians in Bayesian Inference

Tucker Hartland<sup>1</sup>

joint work with:  
Noémi Petra<sup>1</sup>, Georg Stadler<sup>2</sup>,  
Ruanui Nicholson<sup>3</sup>,

<sup>1</sup>School of Natural Sciences, University of California, Merced

<sup>2</sup>Courant Institute of Mathematical Sciences, New York University

<sup>3</sup>Department of Engineering Science, University of Auckland

SIAM Conference on Computational  
Science and Engineering  
February 27, 2019

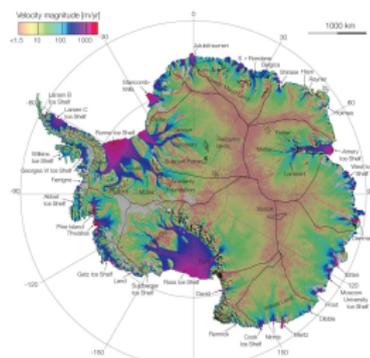


Research Supported By NSF-Grant CAREER-1654311

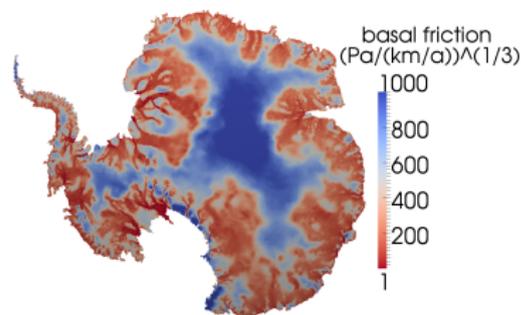


# Outline

- Motivation: infer basal friction field in an ice sheet model and characterize the associated uncertainties.



Observed surface flow velocity from InSAR (Rignot et. al, 2011)



Antarctic ice sheet inversion for the basal friction parameter field from InSAR surface velocities

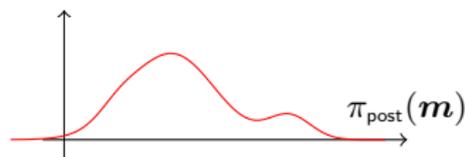
- Bayesian formulation of an inverse problem.
- The role and importance of the Hessian (of the likelihood) in Bayesian inference.
- Hierarchical off-diagonal low-rank (HODLR) approximation of the Hessian.
- Numerical example: inference of membrane log stiffness field in a Poisson model problem.

# Bayesian approach to inverse problems

Inverse problem: given (possibly noisy) data  $\mathbf{d}$  and (a possibly uncertain) model  $\mathbf{f}$ , infer parameters  $\mathbf{m}$  that characterize the model, i.e.,

$$\mathbf{f}(\mathbf{m}) + \mathbf{e} = \mathbf{d}$$

Interpret  $\mathbf{m}$ ,  $\mathbf{d}$  as random variables; solution of inverse problem is the “posterior” probability density function  $\pi_{\text{post}}(\mathbf{m})$  for  $\mathbf{m}$ :



## Remarks:

- Bayesian framework quantifies uncertainty in the inverse solution, given uncertainty in the prior, the data, and the model.
- Prior incorporates known information and in infinite dimensions chosen to act as a regularization.
- Bayesian solution is a probability density function in as many dimensions as the number of parameters.

# Exploring the posterior

If the prior is Gaussian with mean  $\mathbf{m}_{\text{pr}}$  and covariance  $\mathbf{\Gamma}_{\text{prior}}$ , and we assume additive Gaussian noise in the measurements, i.e.,  $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_{\text{noise}})$ , then using Bayes theorem we obtain the posterior pdf:

$$\pi_{\text{post}}(\mathbf{m}) \propto \exp\left(-\frac{1}{2} \|\mathbf{f}(\mathbf{m}) - \mathbf{d}\|_{\mathbf{\Gamma}_{\text{noise}}^{-1}}^2 - \frac{1}{2} \|\mathbf{m} - \mathbf{m}_{\text{pr}}\|_{\mathbf{\Gamma}_{\text{prior}}^{-1}}^2\right)$$

The “maximum a posteriori” point is

$$\begin{aligned}\mathbf{m}_{\text{MAP}} &\stackrel{\text{def}}{=} \arg \max_{\mathbf{m}} \pi_{\text{post}}(\mathbf{m}) \\ &= \arg \min_{\mathbf{m}} \frac{1}{2} \|\mathbf{f}(\mathbf{m}) - \mathbf{d}\|_{\mathbf{\Gamma}_{\text{noise}}^{-1}}^2 + \frac{1}{2} \|\mathbf{m} - \mathbf{m}_{\text{pr}}\|_{\mathbf{\Gamma}_{\text{prior}}^{-1}}^2\end{aligned}$$

⇒ deterministic inverse problem with appropriate weighted norms!

Details in: J. Kaipio and E. Somersalo, Statistical and Computational Inverse Problems, 2005

# The Hessian (of the negative log posterior) plays a critical role in inverse problems

- Its spectral properties characterize the degree of ill-posedness.
- The Hessian drives Newton-type optimization algorithms for solving the inverse problem.
- The inverse of the Hessian locally characterizes the uncertainty in the solution of the inverse problem (under the Gaussian assumption, it is precisely the posterior covariance matrix).
- **Goal:** rapidly perform linear algebraic operations, i.e., manipulation of the Hessian (and its inverse and inverse square root) actions needed by sampling or CG solvers, hence **seek approaches to approximate the Hessian(-applies)**.
- These approximations will then be used as **pre-conditioners**, and to **build MCMC proposals** based on local Gaussian approximations.

# Scalability of the inverse solver

Inexact Newton-CG method applied to an ice sheet inverse problem

#sdof	#pdof	#N	#CG	avgCG	#Stokes
95,796	10,371	42	2718	65	7031
233,834	25,295	39	2342	60	6440
848,850	91,787	39	2577	66	6856
3,372,707	364,649	39	2211	57	6193
22,570,303	1,456,225	40	1923	48	5376

- **#sdof**: number of degrees of freedom for the state variables;
- **#pdof**: number of degrees of freedom for the inversion parameter field;
- **#N**: number of Newton iterations;
- **#CG, avgCG**: total and average (per Newton iteration) number of CG iterations;
- **#Stokes**: total number of linear(ized) Stokes solves (from forward, adjoint, and incremental forward and adjoint problems)

Details in the SIAG CSE19 Best Paper: T. Isaac, N. Petra, G. Stadler, and O. Ghattas. *Scalable and efficient algorithms for the propagation of uncertainty from data through inference to prediction for large-scale problems, with application to flow of the Antarctic ice sheet*, Journal of Computational Physics, 296, 348-368 (2015).

# MCMC sampling: stochastic Newton

Performance results / Convergence diagnostics

	<b>MPSRF</b>	<b>IAT</b>	<b>ESS</b>	<b>MSJ</b>	<b>ARR</b>	<b>#Stokes</b>	<b>time (s)</b>
SN	1.348	600	875	64	2	8400	420

- **MPSRF**: multivariate potential scale reduction factor
- **IAT**: integrated autocorrelation time
- **ESS**: effective sample size
- **MSJ**: mean squared jump distance
- **Statistics**: 21 parallel chains (each 25k); # samples: 525k; dof: 139; rank Hessian: 15
- **ARR**: average rejection rate
- **#Stokes**: # of Stokes solves per independent sample
- **time**: time per independent sample

Proposal density:

$$q(\mathbf{m}_k, \mathbf{m}) \propto \exp\left(-\frac{1}{2} \langle \mathbf{m} - (\mathbf{m}_k - \mathbf{H}^{-1} \mathbf{g}_k), \mathbf{H} (\mathbf{m} - (\mathbf{m}_k - \mathbf{H}^{-1} \mathbf{g}_k)) \rangle_M\right)$$

Details in: N. Petra, J. Martin, G. Stadler, O. Ghattas. *A computational framework for infinite-dimensional Bayesian inverse problems: Part II. Stochastic Newton MCMC with application to ice sheet inverse problems*, SIAM Journal on Scientific Computing, 2014

# The role of the Hessian in Stochastic Newton proposal sampling and deterministic inversion

To efficiently sample from the Stochastic Newton proposal density we need computationally effective means to,

- apply  $\mathbf{H}^{-1}$  to vectors in order to compute the mean,  $\mathbf{m}_k - \mathbf{H}^{-1}\mathbf{g}_k$ ;
- apply  $\mathbf{H}^{-1/2}$  the square root of the covariance operator  $\mathbf{H}^{-1}$ , to vectors.

To efficiently solve the deterministic inverse problem by means of a Newton-CG method we need an estimate of the Hessian that will serve as a preconditioner.

---

## Algorithm 1 Inexact Newton-CG

---

**while** Not sufficiently near a minimizer **do**

    Solve( $\mathbf{H}\mathbf{p} = -\mathbf{g}$ )

$\vdots$

$\mathbf{m}^{(k+1)} = \mathbf{m}^{(k)} + \alpha\mathbf{p}$

**end while**

---

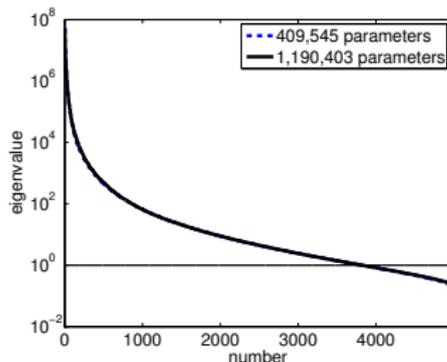
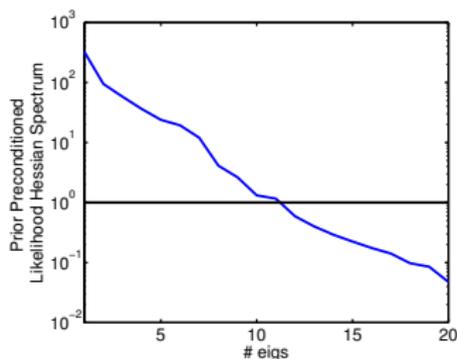
# Low rank approximation of the Hessian of the likelihood

Under the assumption of Gaussian noise and linearized parameter-to-observable map (i.e., Gaussian posterior)

- Invoke low-rank approximation given by the inverse of the Hessian:

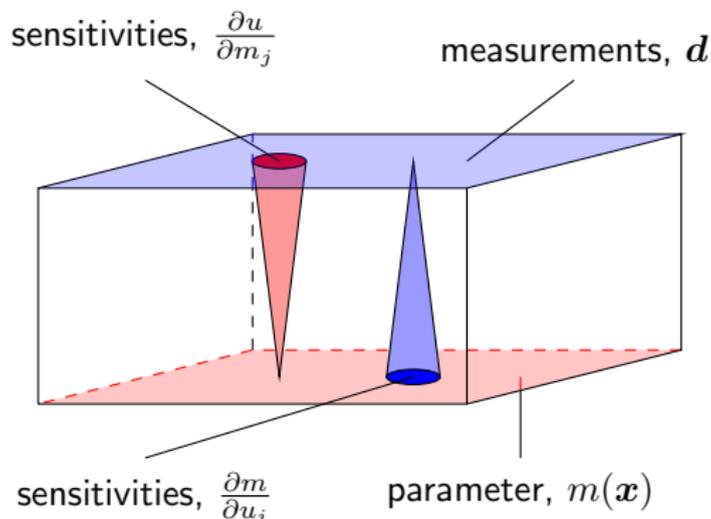
$$\mathbf{\Gamma}_{\text{post}} = \mathbf{H}^{-1} = \left( \mathbf{F}^T \mathbf{\Gamma}_{\text{noise}}^{-1} \mathbf{F} + \mathbf{\Gamma}_{\text{prior}}^{-1} \right)^{-1} \approx \mathbf{\Gamma}_{\text{prior}}^{1/2} (\mathbf{V}_r \mathbf{\Lambda}_r \mathbf{V}_r^T + \mathbf{I})^{-1} \mathbf{\Gamma}_{\text{prior}}^{1/2}$$

- $\mathbf{V}_r$  and  $\mathbf{\Lambda}_r$  are the eigenvectors/values of  $\mathbf{F}^T \mathbf{\Gamma}_{\text{noise}}^{-1} \mathbf{F} \mathbf{v}_i = \lambda_i \mathbf{\Gamma}_{\text{prior}}^{-1} \mathbf{v}_i$ .
- $\mathbf{F}$  is the Jacobian matrix of the parameter to observable map  $\mathbf{f}(\mathbf{m})$ .
- Spectrum of the prior-preconditioned likelihood Hessian for the Arolla glacier (139 parameters, left) and Antarctica (1.19M parameters, right).



# Exploiting structure

Taking advantage of the parameter to observable mapping

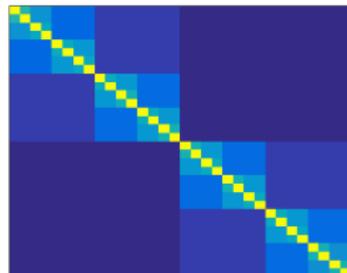
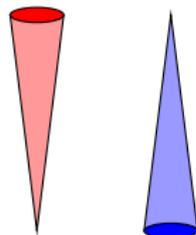


- We expect (hope) for narrow sensitivities.
- $\Rightarrow$  We hope (for a suitable ordering of the discrete degrees of freedom) that the Hessian  $\mathbf{H}$ , will be concentrated on the diagonal.

# Exploiting structure

Taking advantage of the parameter to observable mapping

- We expect (hope) for narrow sensitivities.
- $\Rightarrow$  We hope that  $\mathbf{H}$  is concentrated on the diagonal.
- We wish to maintain only *useful* information in  $\mathbf{H}$ .
- $\Rightarrow$  We expect *useful* information concentrated in certain blocks of  $\mathbf{H}$ .
- Apply techniques of hierarchical ( $\mathcal{H}$ ) matrices: Specifically hierarchical off-diagonal low-rank (HODLR) matrices.





## Desired characteristics of the preconditioner for the Newton-CG system:

- 1 accurately approximates the Hessian,
- 2 the application of the inverse to vectors is computationally inexpensive.

## Compression:

- Apply  $\mathbf{H}$  to random vectors with specified blocks of zeros, in order to sample the column spaces of the off-diagonal blocks.

The samples are likely aligned with dominant modes, thus allowing the generation of low rank approximations of the off-diagonal blocks.

Details in: P.G. Martinsson, *Compressing Rank-Structured Matrices Via Randomized Sampling*, SIAM Journal on Scientific Computing, 2016

## Direct Solve:

- At  $\mathcal{O}(N \log^2 N)$  cost one obtains a factorization of a HODLR matrix.
- Employing the factorization one has a  $\mathcal{O}(N \log N)$  means of applying  $(\mathbf{H}_{\text{HODLR}})^{-1}$  to vectors.

Details in: S. Ambikasaran, E. Darve, *An  $\mathcal{O}(N \log N)$  Fast Direct Solver for Partial Hierarchically Semi-Separable Matrices*, Springer Journal of Scientific Computing, 2013

# Example problem

Log stiffness coefficient field inversion in an elliptic PDE

**Goal:** determine the field  $m(\mathbf{x})$  such that the solution  $u(\mathbf{x})$ , of the elliptic PDE is closest to the observation data,  $u_{\text{obs}}(\mathbf{x})$ .

$$\min_m J[m] := \underbrace{\frac{1}{2} \int_{\Omega} (u(\mathbf{x}) - u_{\text{obs}}(\mathbf{x}))^2 \, d\mathbf{x}}_{\text{Misfit}} + \underbrace{\frac{\gamma}{2} \int_{\Omega} (\nabla m \cdot \nabla m) \, d\mathbf{x}}_{\text{Regularization}}$$

Where  $u(\mathbf{x})$  is uniquely defined by  $m(\mathbf{x})$ , as the solution of the following PDE:

$$\begin{aligned} -\nabla \cdot (\exp(m) \nabla u) &= f(\mathbf{x}), & \mathbf{x} \text{ in } \Omega \\ u(\mathbf{x}) &= 0, & \mathbf{x} \text{ on } \Gamma \end{aligned}$$

- $u(\mathbf{x})$ , state, solution of PDE for a given  $m(\mathbf{x})$ ,
- $u_{\text{obs}}(\mathbf{x})$ , observation data,
- $m(\mathbf{x})$ , log stiffness parameter field,
- $\gamma \geq 0$ , regularization parameter,
- $f(\mathbf{x})$ , source term,
- $\Gamma$ , boundary of spatial domain  $\Omega$ .

# Estimating the minimizer of a deterministic PDE constrained optimization problem

## Method of Solution

- 1 Derive 1st and 2nd order derivative information via adjoints.
- 2 Setup an inexact Newton-CG system.

---

## Algorithm 2 Inexact Newton-CG

---

**while** Not sufficiently near a minimizer **do**

Solve( $\mathbf{H}\mathbf{p} = -\mathbf{g}$ )

$\vdots$

$\mathbf{m}^{(k+1)} = \mathbf{m}^{(k)} + \alpha\mathbf{p}$

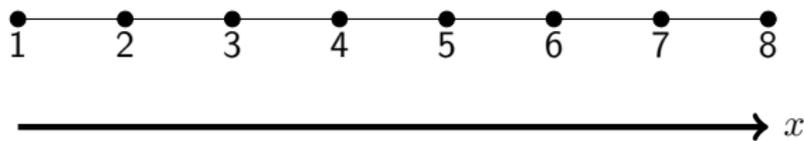
**end while**

---

**Computational Cost:** an application of the Hessian,  $\mathbf{H}$  to a vector requires two PDE solves.

# Example problem in one dimension

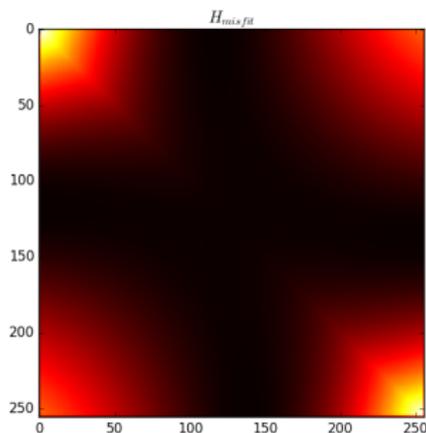
Ordering the degrees of freedom



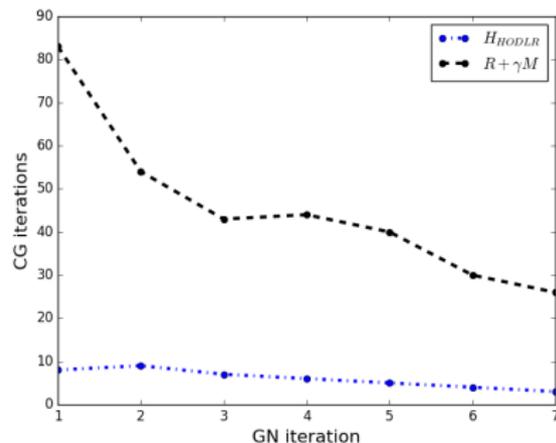
Ordering in 1d naturally respects locality :

- If points  $\mathbf{x}_i, \mathbf{x}_j \in \Omega$  are close to each other then  $|i - j|$  is relatively small.  
If this was *not true*, then the Hessian would have large magnitude elements far from the diagonal.
- If  $i, j$  are such that  $|i - j|$  is relatively small then  $\|\mathbf{x}_i - \mathbf{x}_j\|$  is small.  
If this was *not true*, then the Hessian would have small magnitude elements near the diagonal.

# HODLR preconditioning in one dimension



Matrix plot of the Hessian



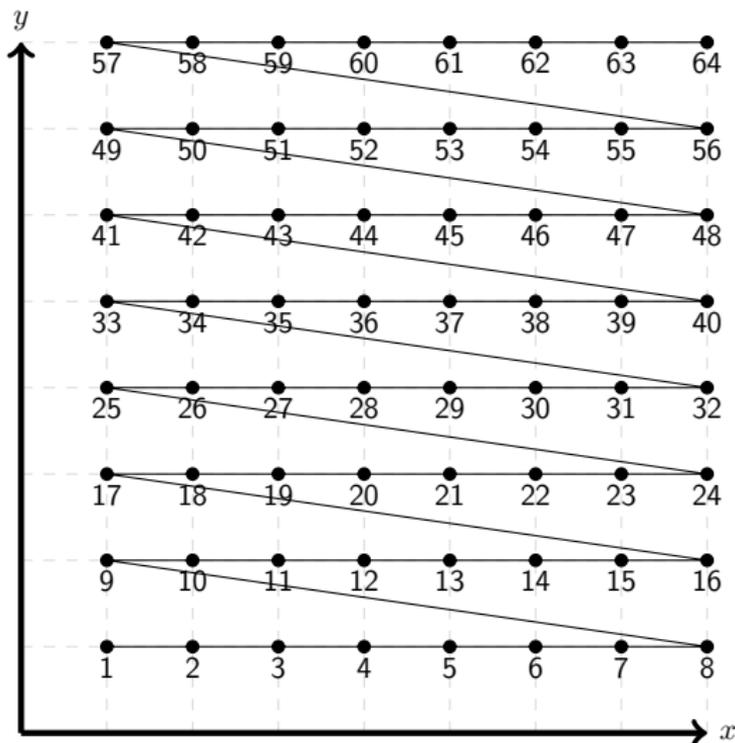
CG iterations vs Newton iteration  
256 dof, 3 levels of H-refinement, OD rank  
 $k = 20$

- 1 The structure of the Hessian indicates that it can be well approximated by a HODLR compression.
- 2 The HODLR compressed Hessian serves as an effective preconditioner for the

# Example problem in two dimensions

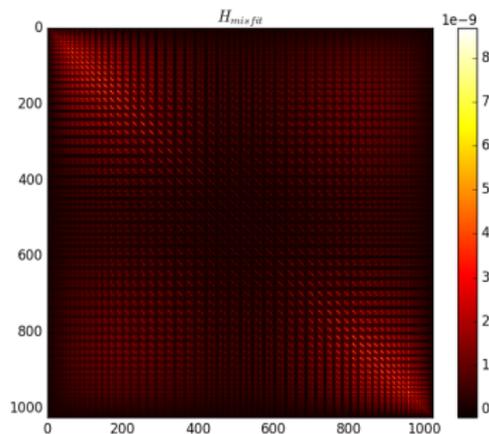
Ordering the degrees of freedom by row

Ordering in 2d does *not* naturally respect locality.



# Hessian with a row ordered basis

Naive row ordering basis leads to Hessians unsuitable for hierarchical compression



Matrix plot of the Hessian

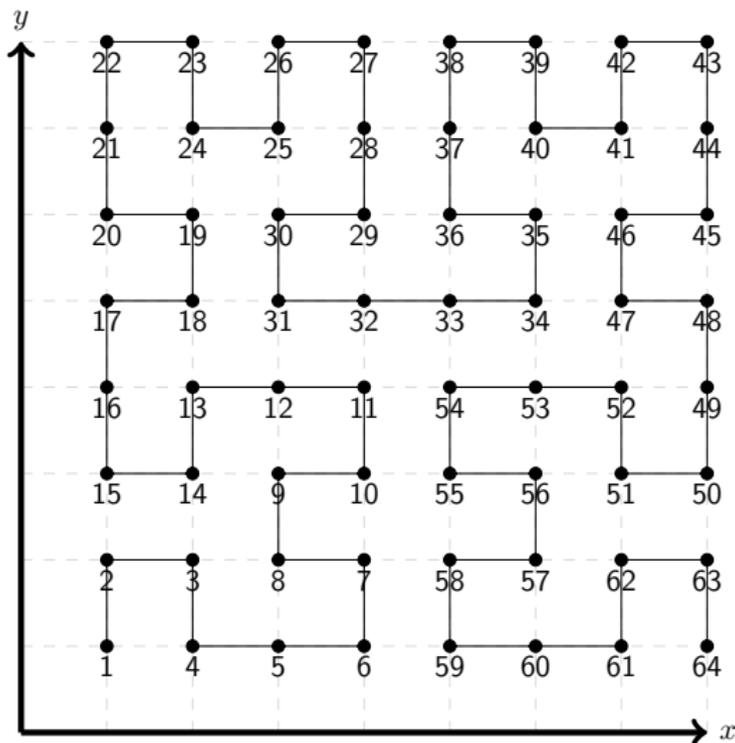
# of Conjugate-Gradient iterations consistently equal to the problem size

When using a default ordering of the degrees of freedom in 2d:

- 1 the Hessian does *not* possess HODLR structure,
- 2 the HODLR compressed Hessian *fails* to effectively precondition the Newton Conjugate-Gradient system.

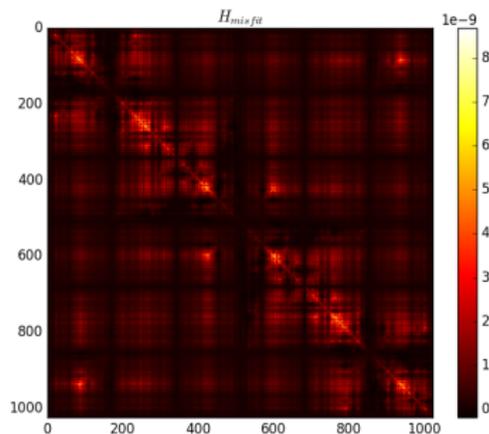
# Example problem in two dimensions

Ordering the degrees of freedom according to a Hilbert curve

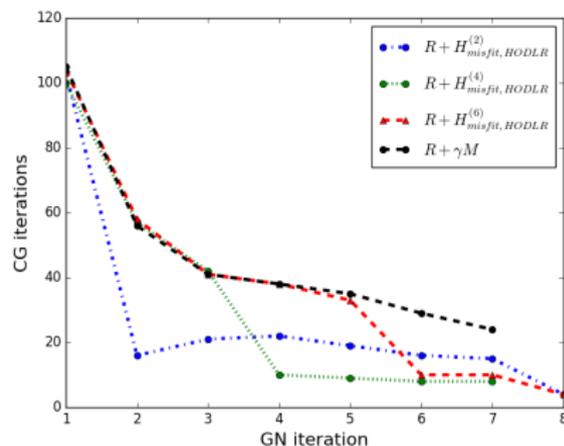


# Hessian with a Hilbert curve ordered basis

Hilbert curve ordering leads to Hessians more suitable for hierarchical compression



Matrix plot of the Hessian



CG iterations vs Newton iteration

When using a Hilbert curve ordering of the degrees of freedom in 2d:

- 1 the Hessian possesses HODLR structure,
- 2 the HODLR compressed Hessian effectively preconditions the Newton Conjugate-Gradient system.

# HODLR compressed Hessian in two dimensions

Positive definiteness is not preserved!

- For moderate levels of hierarchical refinement and/or moderately low off-diagonal rank, the HODLR compressed Hessian **fails to maintain positive definiteness**, thereby failing to adequately precondition the Newton-CG system.
- We hope to moderately extend recent work on generating symmetric positive definite hierarchical semi-separable (HSS) approximations to HODLR matrices.

Details in: X. Xing, E. Chow, *Preserving Positive Definiteness in Hierarchically Semiseparable Matrix Approximations*, SIAM Journal on Matrix Analysis and Applications, 2018

## Conclusions:

- Numerical evidence demonstrates that Hessian manipulations can be made tractable by off-diagonal low-rank approximations.
- So that the localized PDE sensitivities manifest as a diagonally concentrated Hessian, it is necessary to reorder the degrees of freedom in a way that respects locality.

## Ongoing and future work:

- Generate hierarchical approximations that are guaranteed to preserve positive definiteness.
- Apply all of the above to the full Stokes ice sheet problem.
- Apply HODLR approximations of the Hessian to efficiently generate Gaussian proposal distributions for MCMC sampling.
- Look to apply general  $\mathcal{H}$ -matrices to approximate  $\mathbf{H}$ .

# Diagonal concentration of the Hessian

Manifestation of localized sensitivities

Let  $\mathbf{u}(\mathbf{m})$  be the solution of the PDE for a given parameter field  $\mathbf{m}$ .  $\mathbf{M}$ ,  $\mathbf{K}$  the mass and stiffness matrices arising from the discretization of the PDE.

$$J(\mathbf{u}(\mathbf{m}), \mathbf{m}) = \frac{1}{2} (\mathbf{u} - \mathbf{u}_{\text{obs}})^\top (\mathbf{u} - \mathbf{u}_{\text{obs}}) + \frac{\gamma}{2} \mathbf{m}^\top \mathbf{K} \mathbf{m}$$

We take a discretize and then optimize (DTO) approach, to get a sense of the structure of the Hessian in terms of the first and second order sensitivities.

$$\mathbf{H}_{i,j} = \frac{\partial^2}{\partial m_i \partial m_j} J(\mathbf{u}(\mathbf{m}), \mathbf{m})$$

# Diagonal concentration of the Hessian

Manifestation of localized sensitivities

$$\mathbf{H}_{i,j} = (\mathbf{u} - \mathbf{u}_{\text{obs}})^\top \mathbf{M} \frac{\partial^2 \mathbf{u}}{\partial m_i \partial m_j} + \frac{\partial \mathbf{u}}{\partial m_i}^\top \mathbf{M} \frac{\partial \mathbf{u}}{\partial m_j} + \gamma \mathbf{K}_{i,j}$$

- The only components of  $\mathbf{u}$  which will be sensitive to variations of the parameter  $m$  at node  $i$  are those components of  $\mathbf{u}$  that are near node  $i$ .
- If parameter dof nodes  $i$  and  $j$  are far from one another, then no component of  $\mathbf{u}$  will be sensitive to both variations.