PAPER

Hierarchical off-diagonal low-rank approximation of Hessians in inverse problems, with application to ice sheet model initialization

To cite this article: Tucker Hartland et al 2023 Inverse Problems 39 085006

View the article online for updates and enhancements.

You may also like

- Electron scattering cross sections from NHa: a comprehensive study based on Rmatrix method Yingqi Chen, Xianwu Jiang, Lufeng Yao et
- Three-dimensional structure and stability of discontinuities between unmagnetized pair plasma and magnetized electron-
- M E Dieckmann, D Folini, M Falk et al.
- Air-to-land transitions: from wingless animals and plant seeds to shuttlecocks and bio-inspired robots
 Victor M Ortega-Jimenez, Ardian Jusufi, Christian E Brown et al.

Hierarchical off-diagonal low-rank approximation of Hessians in inverse problems, with application to ice sheet model initialization

Tucker Hartland^{1,*}, Georg Stadler², Mauro Perego³, Kim Liegeois³ and Noémi Petra¹

- ¹ Department of Applied Mathematics, University of California, Merced, Merced, CA, United States of America
- ² Courant Institute of Mathematical Sciences, New York University, New York, NY, United States of America
- ³ Center for Computing Research, Sandia National Laboratories, Albuquerque, NM, United States of America

E-mail: thartland@ucmerced.edu

Received 22 September 2022; revised 7 May 2023 Accepted for publication 19 May 2023 Published 26 June 2023



Abstract

Obtaining lightweight and accurate approximations of discretized objective functional Hessians in inverse problems governed by partial differential equations (PDEs) is essential to make both deterministic and Bayesian statistical large-scale inverse problems computationally tractable. The cubic computational complexity of dense linear algebraic tasks, such as Cholesky factorization, that provide a means to sample Gaussian distributions and determine solutions of Newton linear systems is a computational bottleneck at large-scale. These tasks can be reduced to log-linear complexity by utilizing hierarchical off-diagonal low-rank (HODLR) matrix approximations. In this work, we show that a class of Hessians that arise from inverse problems governed by PDEs are well approximated by the HODLR matrix format. In particular, we study inverse problems governed by PDEs that model the instantaneous viscous flow of ice sheets. In these problems, we seek a spatially distributed basal sliding parameter field such that the flow predicted by the ice sheet model is consistent with ice sheet surface velocity observations. We demonstrate the use of HODLR Hessian approximation to efficiently sample the Laplace approximation of the posterior distribution with covariance further approximated by HODLR matrix compression. Computational studies are performed which

^{*} Author to whom any correspondence should be addressed.

illustrate ice sheet problem regimes for which the Gauss-Newton data-misfit Hessian is more efficiently approximated by the HODLR matrix format than the low-rank (LR) format. We then demonstrate that HODLR approximations can be favorable, when compared to global LR approximations, for large-scale problems by studying the data-misfit Hessian associated with inverse problems governed by the first-order Stokes flow model on the Humboldt glacier and Greenland ice sheet.

Supplementary material for this article is available online

Keywords: Hessians, inverse problems, hierarchical matrices, HODLR matrices, ice-sheet models, large-scale

(Some figures may appear in colour only in the online journal)

1. Introduction

Model-based simulations of complex physical systems play an essential role in understanding real world phenomena. These models are often characterized by partial differential equations (PDEs), and are typically subject to uncertainties stemming from unknown coefficient fields, constitutive laws, source terms, initial and/or boundary conditions, geometries, etc. When observation data exist, these parameters can be estimated by solving an inverse problem governed by the underlying model (e.g. PDE). It is well known that uncertainty is a fundamental feature of inverse problems, therefore in addition to inferring the parameters of interest, we need to quantify the uncertainty associated with this inference. This uncertainty quantification can be done via Bayesian inference. Solving Bayesian inverse problems governed by complex PDEs can be extremely challenging due to high-dimensional parameter spaces that stem from discretization of infinite-dimensional parameter fields and the need to repeatedly solve the underlying PDEs. To overcome these computational challenges, it is essential to exploit problem structure, when possible. For example, the underlying PDE solution operator is often diffusive, that observation data may be sparse or only contain limited information about the parameter field. These particularities give rise to a low-rank (LR) structure in the second derivative of the data-misfit component of the inverse problem objective (or of the negative log likelihood), hereafter referred to as the data-misfit Hessian. In previous work [1, 2] we exploited this LR structure in the context of inverse ice sheet flow problems. However, for cases when this 'low-rank' is in fact large, as is the case for many inverse problems of practical interest, where the observation data are highly informative, LR approximation is insufficient. In this article, we exploit the local sensitivity of model predictions to parameters, which gives rise to an off-diagonal LR structure. We do so by invoking hierarchical off-diagonal low-rank (HODLR) matrix approximations and detail how they can be used to reduce the computational cost to solve large-scale PDE-based inverse problems.

1.1. Related work

Global LR approximation of Hessians in inverse problems have been successfully utilized in [1, 3–6], with deterministic and randomized methods [5, 7] being available to generate said approximations. However, some problems, specifically those with highly informative observation data, are not amenable to global LR approximation, and thus other structure-exploiting strategies are needed such as those based on local translation invariance and localized

sensitivities [8–10]. Here, we focus on hierarchical LR methods for which convenient randomized methods are available [11, 12].

Hierarchical matrices have been demonstrated e.g. in [13, 14] to be an effective means to approximate covariance matrices associated to large-scale Gaussian processes. In [15], hierarchical matrix approximations with general hierarchical partitioning patterns are utilized for the construction of explicit representations of Hessian inverses. In one of the examples studied, the authors find that the diffusivity of the parameter-to-observable map and the informativeness of the observation data impact whether the data-misfit Hessian is more suited for compression with hierarchical or global LR formats. Here, we build on this study and focus on a specific inverse problem arising in land ice modeling.

1.2. Contributions

The main contributions of this work are as follows. (1) We motivate the use of HODLR compression for data-misfit Hessians in inverse problems governed by PDEs, and present a detailed study for large-scale ice sheet inverse problems, such as the Greenland ice sheet. (2) We describe a strategy that leverages the fast manipulation of HODLR matrices to efficiently generate samples from a Gaussian distribution for posterior uncertainty quantification. (3) We numerically study the influence of various problem setups on the off-diagonal LR structure of the data-misfit Hessian. The results show the effectiveness of the HODLR approximation for various problem sizes including a Greenland ice sheet inverse problem, which has a discretized parameter dimension of 3.2×10^5 .

2. Preliminaries

In section 2.1, we summarize background material regarding the solution of discretizations of infinite-dimensional inverse problems. We also briefly review HODLR matrices. Specifically, in section 2.2 we define HODLR matrices, list some of their properties and summarize the computational complexities of computing symmetric HODLR matrix approximations of symmetric matrices that are only available through a means to compute its action on vectors. We refer to [16, 17] for a more thorough discussion of hierarchical matrices and to [12] for more detail on HODLR matrices.

2.1. Bayesian inverse problems

A means to account for uncertainty in parameter inference is to employ the Bayesian approach to inverse problems [18–20], which takes as input observation data d, prior knowledge of the parameter and a model for the likelihood of observation data conditional to β . Prior knowledge of the discretized parameter β is typically determined by the expertise of domain scientists and is mathematically encoded in a probability density function $\pi_{\text{prior}}(\beta)$. The likelihood $\pi(d|\beta)$ involves the data uncertainty and the mathematical model for the parameter-to-observable process. The solution of a Bayesian inverse problem is a probability density function for the discretized parameter β , that is conditioned on the observation data according to Bayes formula

$$\pi_{\text{post}}(\boldsymbol{\beta}) = \pi(\boldsymbol{\beta}|\boldsymbol{d}) \propto \pi_{\text{prior}}(\boldsymbol{\beta}) \pi(\boldsymbol{d}|\boldsymbol{\beta}),$$

which provides a formal expression for the posterior distribution. Here, ' \propto ' means equal up to a normalization constant. For a problem with Gaussian prior $\mathcal{N}\left(\overline{\beta},\Gamma_{\text{prior}}\right)$, i.e. a Gaussian

centered at $\overline{\beta}$ with covariance matrix Γ_{prior} , and additive data error ξ described by the zero mean Gaussian $\mathcal{N}(\mathbf{0}, \Gamma_{noise})$, the Bayes formula for $\pi_{post}(\cdot)$ has the following form [19]

$$\pi_{\text{post}}(\boldsymbol{\beta}) \propto \exp\left(-\frac{1}{2}\|\boldsymbol{\mathcal{F}}(\boldsymbol{\beta}) - \boldsymbol{d}\|_{\boldsymbol{\Gamma}_{\text{noise}}^{-1}}^2 - \frac{1}{2}\|\boldsymbol{\beta} - \overline{\boldsymbol{\beta}}\|_{\boldsymbol{\Gamma}_{\text{prior}}^{-1}}^2\right),\tag{1}$$

where \mathcal{F} is the parameter-to-observable map. The notation $\|\cdot\|_A$ means that the norm is weighted by the positive-definite matrix A, i.e. $\|v\|_A := \sqrt{v^\top A v}$. The parameter-to-observable map is typically nonlinear, and consequently the posterior distribution (1) is non-Gaussian. Thus, typically no closed-form expressions for its moments are available. One characteristic of the posterior distribution is the point at which it is maximized, or equivalently the point which minimizes the negative log-posterior, the so-called maximum a posterior (MAP) point,

$$\boldsymbol{\beta}^{\star} := \arg\min_{\boldsymbol{\beta}} J(\boldsymbol{\beta}) := \frac{1}{2} \| \boldsymbol{\mathcal{F}}(\boldsymbol{\beta}) - \boldsymbol{d} \|_{\boldsymbol{\Gamma}_{\text{noise}}^{-1}}^2 + \frac{1}{2} \| \boldsymbol{\beta} - \overline{\boldsymbol{\beta}} \|_{\boldsymbol{\Gamma}_{\text{noior}}^{-1}}^2. \tag{2}$$

A means to compute the MAP point is to employ a Newton-type descent method for optimization [21], which critically relies on the availability of second derivatives. Since J is defined implicitly in terms of the parameter-to-observable map, which involves a PDE solution operator, we utilize the adjoint method [22–24] to compute its gradient and the application of its Hessian to vectors.

To explore posterior distributions beyond the MAP point, Markov chain Monte-Carlo (MCMC) techniques [25, 26] can be used. Such techniques require a proposal distribution that is easy to sample and that at least locally should reflect some of the behavior of the target posterior density. Satisfying these two requirements becomes challenging with increasing parameter dimension. One way to construct such proposals is to use Gaussian approximations of the posterior, either around a current parameter β_k or the MAP point β^* as a proposal. We denote such a proposal by $\tilde{\pi}_{post}$, given by

$$\tilde{\pi}_{\text{post}}(\boldsymbol{\beta}, \boldsymbol{\beta}_k) \propto \exp\left(-\frac{1}{2}(\boldsymbol{\beta} - \boldsymbol{\mu}_k)^{\top} \boldsymbol{H}_k(\boldsymbol{\beta} - \boldsymbol{\mu}_k)\right),$$
 (3)

where \mathbf{g}_k , \mathbf{H}_k are the gradient and Hessian of the log-posterior $J(\beta)$ at β_k , and we use the notation $\mu_k = \beta_k - \mathbf{H}_k^{-1} \mathbf{g}_k$. Note that if β_k is chosen as the MAP point β^* , then $\mu_k = \beta^*$ since the gradient vanishes at the MAP point. Thus, (3) reduces to the Laplace approximation at the MAP point. Another class of MCMC sampling approaches are generalized preconditioned Crank–Nicholson methods [27, 28]. These methods are derived through a discretization of the Langevin equation, and they require a preconditioner that is equivalent to the prior covariance matrix. An attractive choice for this required preconditioner is the Hessian at the MAP point [29].

For the above discussed (and other) MCMC samplers, one typically needs to apply the inverse Hessian H_k^{-1} or its square root $H_k^{-1/2}$ to vectors. These operations are needed repeatedly to either compute μ_k defined above, or to draw samples [2, 29]. The requirement to perform these operations efficiently motivates the study presented in this paper. In particular, in section 3.2 we discuss how HODLR approximations can be used for the fast application of the Hessian square root.

2.2. Symmetric HODLR Matrices

A HODLR matrix $A \in \mathbb{R}^{N \times N}$, is a matrix equipped with a depth $L \in \mathbb{N}$, hierarchical partitionings of the index set $\mathcal{I} = \{1, 2, \dots, N\}$ into contiguous subsets and LR off-diagonal blocks defined by the partition, which is described in greater detail in e.g. [12]. The block rank-

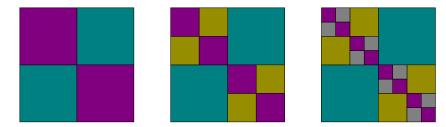


Figure 1. Rank-structure of a matrix A with hierarchical depths L=1 (left), L=2 (middle) and L=3 (right). Off-diagonal blocks are assumed to be low-rank.

structure of a HODLR matrix for various hierarchical depths is illustrated in figure 1. An HODLR matrix must satisfy two additional properties.

(i) The depth of the hierarchical partitioning scales with the logarithm of the size of the matrix, i.e.

$$L = \mathcal{O}(\log N)$$
.

(ii) The maximum rank of each hierarchical level ℓ off-diagonal block, r_{ℓ} , is bounded above by a number r that is independent of the problem size N, for each level ℓ

$$\max_{1\leqslant\ell\leqslant L}r_{\ell}\leqslant r=\mathcal{O}\left(1\right).$$

Such matrices are referred to as data-sparse since the LR blocks allow for them to be represented with less than $\mathcal{O}(N^2)$ floating point numbers. In particular, the storage of an HODLR matrix is $\mathcal{O}(N \log N)$, $\mathcal{O}(N \log N)$ flops are needed to compute a HODLR matrix-vector product [7], and $\mathcal{O}(N \log^2 N)$ flops are required for direct methods to compute an inverse HODLR matrix-vector product [30], as well as square root and inverse square root matrix-vector products [31].

2.2.1. Compression. We aim to generate HODLR approximations of data-misfit Hessians in inverse problems. For large-scale problems, the data-misfit Hessian is typically not available explicitly, but only as an operator that can be applied to vectors. Each such Hessian-vector product requires the solution of two, potentially large-scale, PDE model solves, the first being a linearized forward model and the second a linear adjoint model. Due to the complex computational procedure required of each Hessian-vector product, these products dominate the overall computation for large-scale problems and we thus use their number throughout this work as a measure of computational cost. In order to construct HODLR approximations of symmetric matrix-free operators, we employ previously developed randomized linear algebraic routines which only require the matrix-free action on a limited number of random vectors with specified null entries, referred to as structured random vectors. The Hessian action on these structured random vectors is used to sample row and column spaces of off-diagonal Hessian submatrices and allow for randomized approximate truncated singular value decompositions of the aforementioned off-diagonal submatrices. More details can be found in the appendix; see algorithm 2.

For the results that we present in section 5 a rank-adaptive symmetric matrix-free [32, 33] hierarchical compression algorithm is utilized, that is based on [12]. A similar algorithm is

presented in [34], wherein the hierarchical partitioning is more general and the LR blocks have nested bases. The rank-adaptivity provides a high probability means of resolving the off-diagonal blocks to a desired level of accuracy. By utilizing available matrix–vector product information and the Rayleigh quotient, a rank adaptive relative tolerance algorithm is made possible. More computationally efficient rank-adaptive stopping criteria are detailed in [35].

2.2.2. Computational cost of generating HODLR approximations. The number of matrix-vector products ζ , needed to compress a symmetric matrix using q oversampling vectors, into a level L HODLR matrix with off-diagonal ranks $\{r_\ell\}_{\ell=1}^L$ is given by

$$\zeta = 2(\langle r \rangle + q)L + N/2^L$$
, where $\langle r \rangle := \frac{1}{L} \sum_{\ell=1}^{L} r_{\ell}$. (4)

Equation (4) can be understood from algorithm 2 in appendix 'Randomized compression algorithms', as for each level ℓ one needs to compute $r_{\ell} + q$ Hessian-vector products, in order to compute Y (line 7 of algorithm 2) and $r_{\ell} + q$ Hessian-vector products to compute Z (line 14 of algorithm 2). The remaining $N/2^L$ Hessian-vector products arise from the need to determine the diagonal sub-blocks, which is detailed in [7]. We note that an adaptive procedure to determine an approximate basis Q, such as that in [33], for a block matrix column space, reduces the number of matrix–vector products to $\zeta_{\text{adaptive}} = 2(\langle r \rangle + q/2)L + N/2^L$ but with the additional computational burden of extra orthogonalization routine calls. We note that $\zeta = \mathcal{O}(\log N)$ matrix–vector products are needed to generate an HODLR approximation of a matrix with HODLR structure. For sufficiently large problems HODLR compression is not expected to be more computationally efficient than global LR compression, as the number of Hessian-vector products to generate a LR approximation by the randomized single-pass algorithm [12] is independent of the problems size. However, for problems of substantial size, we observe that the HODLR format does offer computational savings (see section 6).

3. HODLR matrices in inverse problems governed by PDEs

Here, we illustrate why data-misfit Hessians in inverse problems governed by PDEs may contain numerically LR off-diagonal blocks, describe how one can permute parameters to expose this HODLR structure, and show how HODLR approximations can be leveraged to draw samples from Gaussian approximations of Bayesian posterior distributions.

3.1. Motivation

Consider the following data-misfit cost functional

$$J_{ ext{misfit}}(eta) := rac{1}{2} \| \mathcal{F}(eta) - \boldsymbol{d} \|_{\Gamma_{ ext{noise}}^{-1}}^2, \quad ext{with} \quad \mathcal{F}(eta) = \boldsymbol{\mathcal{B}} u,$$

where \mathcal{B} linearly maps the PDE solution $u = u(\beta)$, for the spatially-distributed parameter field β , to the model predictions associated to the data d. Moreover, Γ_{noise} is the covariance matrix describing the Gaussian noise of the observational data. For illustration purposes, we assume that the parameter function β is defined on a region Γ_1 and the data d is observed on a region Γ_2 , which may or may not be distinct. These quantities are related through the solution of the governing PDE and the measurement operator \mathcal{B} . The characteristics of this relation depends on properties of the governing PDE. In the following, we assume that a spatially (or temporally) localized perturbation in the β field leads to a predominantly localized effect in the PDE

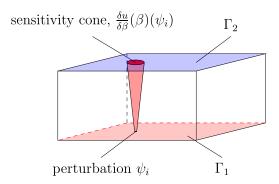


Figure 2. Sketch illustrating a case where the influence of changes in the parameter β on the PDE solution u in Γ_2 is focused in a small area. To illustrate this, we show a sensitivity cone, i.e. the PDE solution u is predominantly impacted in a cone about the support of the localized parameter perturbation.

solution u, and thus the model predictions $\mathcal{B}u$. This property is illustrated in figure 2, where a sensitivity cone depicts the influence of a local perturbation in β , defined over Γ_1 , on the PDE solution u in Γ_2 . It is well known that for an elliptic PDE, local perturbations influence the solution globally, but depending on the geometry of the domain and the equation, this global effect may rapidly decay outside a subset of Γ_2 that captures the main effects of the perturbation. For instance, in a problem as in figure 2, the influence of perturbations in β on u is likely to become more localized when the distance between Γ_1 and Γ_2 decreases.

We next discuss the relationship between properties of the PDE as discussed above and offdiagonal blocks in the Hessian matrix (or its Gauss–Newton variant). The data-misfit Hessian can be derived using the adjoint method [22–24], but we find the HODLR structure of the data-misfit Hessian to be most easily understood by studying a formal expression in terms of the first and second order sensitivities $\delta u/\delta \beta$, $\delta^2 u/\delta \beta^2$. Said expression is given by

$$\begin{split} \frac{\delta^{2}}{\delta\beta^{2}} J_{\text{misfit}}\left(\beta\right)\left(\beta_{1},\beta_{2}\right) &= \left(\boldsymbol{\mathcal{B}}\boldsymbol{u} - \boldsymbol{d}\right)^{\top} \boldsymbol{\Gamma}_{\text{noise}}^{-1} \left(\boldsymbol{\mathcal{B}} \frac{\delta^{2}\boldsymbol{u}}{\delta\beta^{2}}\left(\beta\right)\left(\beta_{1},\beta_{2}\right)\right) \\ &+ \left(\boldsymbol{\mathcal{B}} \frac{\delta\boldsymbol{u}}{\delta\beta}\left(\beta\right)\left(\beta_{1}\right)\right)^{\top} \boldsymbol{\Gamma}_{\text{noise}}^{-1} \left(\boldsymbol{\mathcal{B}} \frac{\delta\boldsymbol{u}}{\delta\beta}\left(\beta\right)\left(\beta_{2}\right)\right), \end{split}$$

where $\delta u/\delta \beta(\beta)(\beta_1)$ is the first variation [36] of u with respect to β in direction β_1 , and $\delta^2 u/\delta \beta^2(\beta)(\beta_1,\beta_2)$ is the second variation of u with respect to β in directions β_1,β_2 , that is,

$$\frac{\delta u}{\delta \beta}(\beta)(\beta_1) := \left[\frac{\mathrm{d}}{\mathrm{d}\epsilon} u(\beta + \epsilon \beta_1)\right]_{\epsilon = 0},$$

$$\frac{\delta^2 u}{\delta \beta^2}(\beta)(\beta_1, \beta_2) := \left[\frac{\mathrm{d}}{\mathrm{d}\epsilon} \frac{\delta u}{\delta \beta} (\beta + \epsilon \beta_2)(\beta_1) \right]_{\epsilon = 0}.$$

Upon discretizing β with finite elements we obtain the following formal expression for the (i,j)-entry of the data-misfit Hessian $\boldsymbol{H}_{\text{misfit}}^{\text{GN}}$ and of the Gauss-Newton data-misfit Hessian $\boldsymbol{H}_{\text{misfit}}^{\text{GN}}$

$$\left(\boldsymbol{H}_{\text{misfit}}\right)_{i,j} = \frac{\delta^2}{\delta\beta^2} \left(J_{\text{misfit}}\left(\beta\right)\right) \left(\psi_i, \psi_j\right),\tag{5}$$

$$\left(\boldsymbol{H}_{\text{misfit}}^{\text{GN}}\right)_{i,j} = \left(\boldsymbol{\mathcal{B}}\frac{\delta u}{\delta\beta}\left(\beta\right)\left(\psi_{i}\right)\right)^{\top} \boldsymbol{\Gamma}_{\text{noise}}^{-1} \left(\boldsymbol{\mathcal{B}}\frac{\delta u}{\delta\beta}\left(\beta\right)\left(\psi_{j}\right)\right),\tag{6}$$

where $\{\psi_j\}_{j=1}^N$ is a basis for the nodal finite-element space, which is used to approximate β . When sensitivities are predominantly local as discussed above and when the support of two finite element basis functions ψ_i, ψ_i are well separated, the terms

$$\left(\boldsymbol{\mathcal{B}}\frac{\delta u}{\delta \beta}\left(\beta\right)\left(\psi_{i}\right)\right)^{\top} \, \boldsymbol{\Gamma}_{\mathrm{noise}}^{-1} \! \left(\boldsymbol{\mathcal{B}}\frac{\delta u}{\delta \beta}\left(\beta\right)\left(\psi_{j}\right)\right) \quad \text{and} \quad \boldsymbol{\mathcal{B}} \! \left(\frac{\delta^{2} u}{\delta \beta^{2}}\left(\beta\right)\left(\psi_{i},\psi_{j}\right)\right),$$

are rather small (assuming diagonally dominant noise covariance matrices). This is, e.g. due to $\mathcal{B}\delta u/\delta\beta(\beta)(\psi_i)$ having small values when $\mathcal{B}\delta u/\delta\beta(\beta)(\psi_j)$ is large. Now, let \mathcal{I},\mathcal{J} be disjoint index subsets of $\{1,2,\ldots,N\}$, then the entries in the matrix block $\{(\boldsymbol{H}_{\text{misfit}})_{i\in\mathcal{I},j\in\mathcal{J}}\}$ of the datamisfit Hessian are relatively small whenever $\cup_{i\in\mathcal{I}} \text{supp}(\psi_i)$ and $\cup_{j\in\mathcal{J}} \text{supp}(\psi_j)$ are well separated. Such Hessian blocks are well suited for approximation by LR matrices. When the degrees of freedom (dofs) corresponding to the finite element basis functions ψ_i are ordered such that \mathcal{I},\mathcal{J} are contiguous, $(\boldsymbol{H}_{\text{misfit}})_{\mathcal{I},\mathcal{J}}$ is an off-diagonal sub-block of $\boldsymbol{H}_{\text{misfit}}$ and $\boldsymbol{H}_{\text{misfit}}$ tends to have HODLR structure as defined in section 2.2. The Gauss–Newton data-misfit Hessian may have HODLR structure for the same reasons. In both cases, the order of the basis functions and thus the dofs influence this structure. An ordering that maintains locality, i.e. consecutive indices correspond to basis functions with supports that are near one another, is ideal. As a consequence of such ordering, basis function supports with significantly different indices are far from each other such that the corresponding off-diagonal blocks have small entries and can be well approximated using a LR matrix approximation. We defer to section 6.2 for a discussion of methods and numerical experiments regarding the order of the dofs.

3.2. Exploiting HODLR structure for fast sampling of Gaussian approximations of the posterior distribution

In [2], the following expressions for the covariance of the Laplace approximation of the posterior distribution are provided,

$$egin{aligned} oldsymbol{\Gamma}_{
m post} &= \left(oldsymbol{H}_{
m misfit} + oldsymbol{\Gamma}_{
m prior}^{-1}
ight)^{-1} = oldsymbol{\Gamma}_{
m prior}^{1/2} \left(oldsymbol{H}_{
m misfit}' + oldsymbol{I}
ight)^{-1} oldsymbol{\Gamma}_{
m prior}^{7/2}, \ oldsymbol{H}_{
m misfit}' &= oldsymbol{\Gamma}_{
m prior}^{1/2} oldsymbol{H}_{
m misfit} oldsymbol{\Gamma}_{
m prior}^{1/2}, \ oldsymbol{\Gamma}_{
m post}^{1/2} &= oldsymbol{\Gamma}_{
m prior}^{1/2} \left(oldsymbol{H}_{
m misfit}' + oldsymbol{I}
ight)^{-1/2}, \end{aligned}$$

where the matrix square-root $A^{1/2}$ is such that $A = A^{1/2} \left(A^{1/2}\right)^{\top}$. For Bayesian inverse problems with a parameter field that is distributed spatially over a bounded subset of \mathbb{R}^m , m=2,3, a reasonable choice is to use the square of an inverse elliptic PDE operator, whose action is given by twice applying an inverse elliptic PDE operator, to define the prior covariance [20]. Furthermore, this choice permits a relatively simple means of obtaining a symmetric square root of Γ_{prior} , as an inverse elliptic PDE operator. Multigrid solvers for linear systems that arise from discretized elliptic PDEs [37], provide a scalable means to compute $\Gamma_{\text{prior}}^{1/2}x$. In previous works such as [1, 3–6], the prior-preconditioned data-misfit Hessian H'_{misfit} , was approximated by means of a global LR compression. This strategy provides an efficient means of approximating the posterior covariance matrix in inverse problems with sufficiently small amounts of

observation data. Here, we exploit HODLR problem structure and generate approximations of the posterior covariance matrix, that comes from the Laplace approximation of the posterior, by HODLR approximations of the prior-preconditioned data-misfit $\tilde{H}'_{\text{misfit}}$. We term the resulting Gaussian distribution as the HODLR Laplace approximation of the posterior. Appendix 'Error of the Laplace posterior covariance due to approximation of the prior-preconditioned data-misfit Hessian' provides an analysis on how such an approximation impacts the accuracy of the approximate posterior covariance

$$\tilde{m{\Gamma}}_{
m post} = m{\Gamma}_{
m prior}^{1/2} \left(m{ ilde{H}}_{
m misfit}^{\prime} + m{I}
ight)^{-1} m{\Gamma}_{
m prior}^{ op/2}.$$

A symmetric square-root factorization of $\tilde{\boldsymbol{H}}'_{\text{misfit}} + \boldsymbol{I}$ is then generated with $\mathcal{O}(N \log^2 N)$ flops [31]. The symmetric factorization allows for a $\mathcal{O}(N \log N)$ means of computing square root and inverse square root matrix-vector products.

4. Bayesian inverse ice sheet problems

The simulation of the dynamics of ice sheets (e.g. the Greenland or Antarctic ice sheets) is an important component of coupled climate simulations. Such simulations require estimation of a present state of the ice that is consistent with available observations, a process sometimes referred to as model initialization. This estimation problem can be formulated either as a deterministic inverse problem (i.e. as nonlinear least squares optimization governed by PDEs) or as a Bayesian inverse problem (i.e. as a statistical problem which aims to characterize a distribution of states). The latter approach, while more expensive, provides uncertainty estimates in addition to determining a best parameter fit.

Ice sheet dynamics [38] is typically governed by nonlinear Stokes equations or simplifications thereof, such as the first-order equations (see e.g. [39]). Generally, the most uncertain component in ice sheet simulations is the basal boundary condition, i.e. how the ice sheet interacts with the rock, sand, water or a mix thereof at its base. Estimating an ice sheet's effective boundary condition from velocity observations on the top surface, the ice sheet's geometry and a model for its dynamics is thus an important problem that can be mathematically formulated as an inverse problem [1, 40–43].

We summarize the formulation of this inverse problem next. As common in the literature, we use *a snapshot* optimization approach, where all the data are assumed to be collected over a short period of time during which changes in the ice geometry are negligible. We denote the bounded domain covered by ice by $\Omega \subset \mathbb{R}^m$, $m \in \{2,3\}$, and the basal, lateral and top parts of the domain boundary $\partial \Omega$ by Γ_b , Γ_l , and Γ_t , as illustrated in figure 3.

The governing equations are nonlinear incompressible Stokes equations whose solution is the ice flow velocity $u : \Omega \to \mathbb{R}^m$ and the pressure $p : \Omega \to \mathbb{R}$ given as follows:

$$-\nabla \cdot \boldsymbol{\sigma}_{\boldsymbol{u}} = \rho \boldsymbol{g} \text{ in } \Omega, \tag{7a}$$

$$\nabla \cdot \boldsymbol{u} = 0 \quad \text{in } \Omega, \tag{7b}$$

$$\sigma_{u}n = 0 \quad \text{on } \Gamma_{t}, \tag{7c}$$

$$\mathbf{u} \cdot \mathbf{n} = 0 \text{ and } \mathbf{T}(\boldsymbol{\sigma}_{\mathbf{u}}\mathbf{n} + \exp(\beta)\mathbf{u}) = \mathbf{0} \text{ on } \Gamma_b,$$
 (7*d*)

along with additional lateral boundary conditions. Here, β is a basal sliding parameter field, ρg the body force density, where ρ is the mass density of the ice and g the acceleration due to gravity. Equation (7a) describes the conservation of momentum, (7b) the conservation of

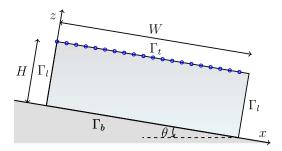


Figure 3. Schematic of two-dimensional slab of ice used for Example I in section 5. The blue circles show representative (uniformly-distributed) measurement locations. The angle θ is the slope of the ice slab.

mass, and (7c) are stress-free boundary conditions for the top surface (the ice-air interface). In normal direction, equation (7d) states a non-penetration condition, i.e. the ice cannot flow into the rock/sand layer which supports it (here n denotes the outward unit normal to the boundary $\partial\Omega$ and T the tangential operator, $Tv = v - n(n^{\top}v)$). In tangential direction, equation (7d) specifies a tangential sliding condition that relates the fraction of tangential sliding and tangential stress through the (logarithmic) basal sliding field $\beta = \beta(x)$, $x \in \Gamma_b$. We employ Glen's flow law [44], a constitutive law for ice that relates the stress tensor σ_u and the strain rate tensor $\dot{\varepsilon}_u = \frac{1}{2} \left(\nabla u + \nabla u^{\top} \right)$,

$$\sigma_{\mathbf{u}} = 2\eta(\mathbf{u})\dot{\boldsymbol{\varepsilon}}_{\mathbf{u}} - \mathbf{I}p, \text{ with } \eta(\mathbf{u}) = \frac{1}{2}A^{-1/n}\dot{\boldsymbol{\varepsilon}}_{\text{II}}^{\frac{1-n}{2n}},$$
 (8)

where η is the effective viscosity, I is the unit matrix, $\dot{\varepsilon}_{II} = \text{tr}(\dot{\varepsilon}_u^2)$ is the second invariant of the strain rate tensor, A is a flow rate factor, and n is Glen's exponent. Ice is typically modeled using $n \approx 3$, which corresponds to a shear-thinning constitutive relation, here we use n = 3.

As discussed above, the parameter containing the largest uncertainty is the (logarithmic) basal sliding field $\beta = \beta(x)$. Thus, it is usually the parameter inferred from (typically, satellite) observation data d, here in the form of surface velocity measurements. Using an appropriate point observation operator \mathcal{B} that extracts point data from the solution u of the governing equation (7), and assuming additive observation errors ξ , the relationship between model and data is now of the typical form

$$d = \mathcal{B}u + \xi. \tag{9}$$

Assuming that the observation errors ξ and the prior for the parameter field β follow Gaussian distributions, we are in the framework of Bayesian inverse problems summarized in section 2.1.

5. Example I: two-dimensional ISMIP-HOM benchmark

We first study the prospects of compressing Gauss–Newton data-misfit Hessians in a problem inspired by the ISMIP-HOM collection of ice sheet simulation benchmark problems [45]. This problem set was used to explore inverse ice sheet problems in e.g. [42, 43]. After a short description of the problem setup, we present results such as the MAP point estimate β^* and samples from the HODLR Laplace approximation of the posterior distribution. Then, we study the impact that various problem features have on the suitability of the Gauss–Newton data-misfit Hessian for compression to the HODLR and global LR formats.

5.1. Problem setup

This problem setup consists of a rectangular piece of ice on a slope, as sketched in figure 3. This simple example allows us to study the influence of the domain aspect ratio, the number of observations and the level of mesh refinement on the properties of the Gauss–Newton datamisfit Hessian matrix. The domain has a width of $W = 10^4$ (m) and a height of $H = 10^2$ (m). Periodic boundary conditions are employed along the lateral boundaries such that the setup models an infinite slab of ice on a slope. The governing equations and other boundary conditions are as discussed in equation (7).

The Stokes equations are discretized using Taylor–Hood finite elements on a mesh of 256×10 rectangles, each subdivided into two triangles, for the domain length [0, W) and height [0, H]. To compute a MAP point estimate, we generate synthetic surface velocity observation data using the 'true' logarithmic basal sliding field, $\beta_{\text{true}}(x) := \log\left(1200 + 1100\sin\left(\frac{2\pi x}{W}\right)\right)$. Given this basal sliding field, we solve equation (7), extract the tangential velocity component at 100 uniformly distributed points on the top boundary Γ_t . The synthetic observation data d is obtained by corrupting each component of the extracted tangential velocity by adding random noise ξ_i to its ith component d_i . The random noise ξ_i are independent and identically distributed according to a Gaussian with zero mean and standard deviation equal to 1% of the maximum absolute value of the extracted tangential velocity field.

It remains to define the prior distribution for the parameter field β . The average value of β_{true} is used as constant prior mean $\overline{\beta}(x)=6.73315\approx\frac{1}{W}\int_0^W\beta_{\text{true}}(s)\,\mathrm{d}s$. The prior covariance matrix Γ_{prior} is a discretization of the covariance PDE operator $\mathcal{C}:=(\delta I-\gamma\Delta)^{-1}$, with $\gamma=6\times10^2$ and $\delta=2.4\times10^{-3}$, with Robin boundary conditions [46]. These values are chosen in order to provide a relatively large prior correlation length of 10^3 (m) [47]. After discretizing $\overline{\beta}$ by finite elements, the resultant dofs define $\overline{\beta}$ as in equation (2). Next, we summarize the computation of the MAP point and the compression of the Gauss–Newton data-misfit Hessian matrix at the MAP point.

5.2. MAP point and HODLR Laplace approximation of the posterior

The nonlinear optimization problem for finding the MAP estimate is solved using an inexact Gauss–Newton minimization method with backtracking linesearch [21], where the linear systems are iteratively solved by the conjugate gradient method. The resulting MAP point is shown in figure 4. The MAP parameter field β^* , closely resembles the true parameter β_{true} , which is a consequence of the large amount of available data and relatively small noise level

Next, we sample the HODLR Laplace approximation of the posterior distribution with a HODLR compressed prior-preconditioned data-misfit Hessian H'_{misfit} (details and comparisons can be found below in section 5.3), as outlined in section 3.2. In figure 5, we compare the mean, pointwise standard deviation and samples from the prior and the posterior distributions. As expected, we find that the data updates our belief about the spatially distributed parameter field and reduces the uncertainty. In particular, the 2σ bounds on the one-dimensional point marginals $\sigma(x)$, $\sigma_i = [\Gamma_{i,i}]^{-1/2}$ of the Laplace approximation of the posterior and the prior distributions are shown, in order to verify that the samples are largely contained within two standard deviations of their respective means. The prior-preconditioned data-misfit Hessian H'_{misfit} , is compressed using a relative tolerance of 10^{-6} , so that with high probability $||H'_{\text{misfit}}||_2 \le 10^{-6}$.

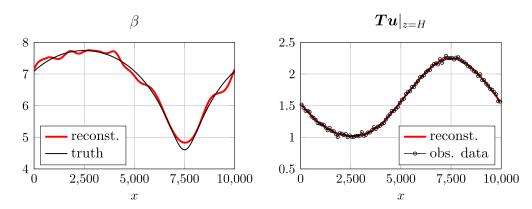


Figure 4. Shown for Example I are on the left the MAP point β^* (red) and the truth basal sliding parameter β_{true} (black) used to generate synthetic observations of the tangential velocity component on the upper surface Γ_t . Shown on the right are noisy synthetic observations (black dots) used for computing the MAP point and the associated tangential surface velocity reconstruction (red).

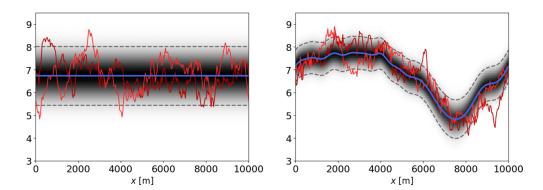


Figure 5. Results for Example I: two random samples (red), mean $\overline{\beta}$ (blue) and boundaries of the region $R = \{(x,y) \text{ such that } 0 \le x \le W \text{ and } \overline{\beta}(x) - 2\sigma(x) \le y \le \overline{\beta}(x) + 2\sigma(x)\}$ (dashed black) are shown for the prior (left) and a HODLR Laplace approximation of the posterior using the methodology described in section 3.2 (right).

5.3. Dependence of Hessian block spectra on problem setting

Next, we study how problem features impact the numerical suitability of using global LR and HODLR compressions to approximate the Gauss–Newton data-misfit Hessian. In this and subsequent sections we measure the cost to generate the matrix compression in terms of Hessian-vector products, which we also describe as Hessian applies, as each said vector product requires two linearized PDE solves and thus dominates the computational cost. We use the result of appendix 'HODLR approximation error due to the accumulation of LR approximations of off-diagonal blocks', to claim ε absolute error in a level L HODLR approximation, when there is no more than ε/L absolute error in each off-diagonal block. What is particular to this section, is that *adaptive* single-pass and HODLR algorithms are used to generate global LR and HODLR approximations, based on absolute tolerance criteria. The

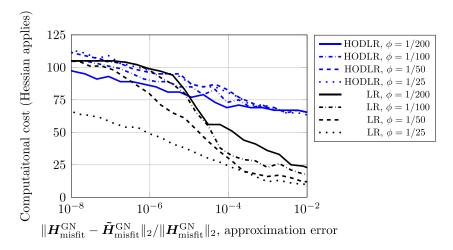


Figure 6. Comparison of HODLR and global LR compression costs of the Gauss–Newton data-misfit Hessian $H_{\text{misfit}}^{\text{GN}}$, for Example I with ice sheet aspect ratio ϕ . This figure shows that for low aspect ratios, HODLR becomes more efficient than global LR for medium levels of target accuracy.

absolute tolerance algorithmic input is scaled by the largest global LR singular value in order to report relative approximation errors. We note that additional errors are neglected in the reported approximation error such as that incurred in the peeling process [11, 12] and additional approximation assumptions in the single-pass algorithm, both of which are not expected to be significant.

5.3.1. Influence of aspect ratio. Here, we vary the aspect ratio of the domain $\phi = H/W$, where H and W are the domain height and width respectively, in order to study how it influences the block spectra of the Gauss–Newton data-misfit Hessian and ultimately the computational cost. Figure 6 shows that the global spectrum is more sensitive to changes in the relative length scale ϕ than the spectra of the off-diagonal blocks. LR approximations of the off-diagonal blocks become computationally cheaper as ϕ decreases as a result of the sensitivity cones becoming increasingly localized as the ice sheet thickness decreases. Global LR approximations become more expensive as ϕ decreases, a result of the data being more informative. We note that realistic problems, such as the Humboldt glacier and the Greenland ice sheet studied later in section 6, have small aspect ratios and are thus expected to have data-misfit Hessians that are less amenable to global LR approximation.

5.3.2. Influence of the parameter dimension. We now vary the level of mesh refinement in order to study the influence of observation data informativeness, through the discretized parameter dimension $N = \dim(\beta)$, on the computational cost to generate HODLR and global LR approximations of the Gauss–Newton data-misfit Hessian. The hierarchical depth L is incremented for every doubling of the discretized parameter dimension, in order that the hierarchical depth scales with the logarithm of the size of the Hessian matrix, a condition described in section 2.2. Figure 7 provides computational evidence of the claim made in section 2.2, that the number of applies needed to hierarchically compress an operator with HODLR structure

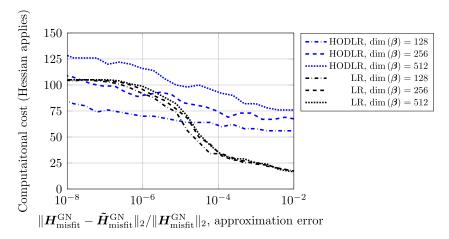


Figure 7. Dependence of HODLR and global LR compression costs of the Gauss–Newton data-misfit Hessian on $\dim(\beta)$, the dimension of the discretized logarithmic basal sliding field for Example I. The cost of global LR compression is insensitive to $\dim(\beta)$, while the cost of HODLR compression increases as the mesh is refined.

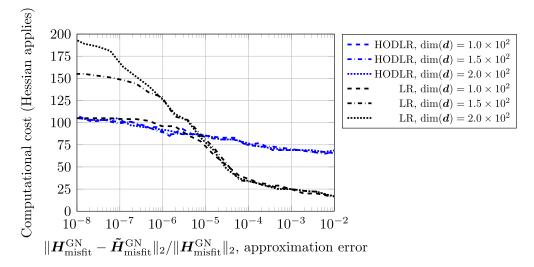


Figure 8. Dependence of HODLR and global LR compression costs of the Gauss–Newton data-misfit Hessian on $\dim(d)$, the data dimension, for Example I. The computational cost for global LR approximation increases with the amount of observation data, while the cost for HODLR compression is rather insensitive.

is $\mathcal{O}(\log N)$. On the contrary, the number of applies to generate the global LR approximation is rather insensitive to the level of mesh refinement.

5.3.3. Influence of the data dimension. Figure 8 shows that the global rank grows with the amount of observation data and thus global LR compression tends to be less efficient for problems with strongly informative observation data. The rate of spectral decay

of the (Gauss–Newton) data-misfit Hessian is related to the degree of ill-posedness of the unregularized inverse problem. As the number of observations increases, these associated model predictions are increasingly sensitive to small scale variations in the basal sliding field. Thus, more data generally makes the data set more informative about the parameter and the (Gauss–Newton) data-misfit Hessian have a weaker rate of spectral decay.

6. Example II: Humboldt glacier and Greenland ice sheet

Here, we study the scalability of the proposed methods using large-scale ice sheet problems which are typically used in climate simulations. Namely, we focus on the Humboldt glacier in North-West Greenland, and the entire Greenland ice sheet. For these simulations, we use the ice sheet model MALI [48], which relies on Albany [49], a C++ multi-physics library for the implementation of the first-order approximation of Stokes equations. This first-order approximation is based on scaling arguments motivated by the shallow nature of ice sheets and uses the incompressibility condition to reduce the unknowns to the horizontal velocities. We use PyAlbany [50] a convenient Python interface to the Albany package, which in turn builds upon Trilinos [51]. Albany is designed to support parallel and scalable finite-element discretized PDE solvers and various analysis capabilities. Details about the parameter, state, data dimensions as well as the number of cores and hierarchical levels used in the computations is provided in table 1.

The following study is partially motivated by findings made in the section 5, namely that the role of the aspect ratio between the vertical and horizontal directions (see section 5.3) influences the ability to use global LR compression and favors HODLR compression. We generate HODLR and global LR approximations and then based on the computed spectra, equation (4) and $\zeta^{LR} = r + q$, we estimate the computational cost. Additionally, we demonstrate that the ordering of the dofs impacts the spectral decay for off-diagonal blocks of the data-misfit Hessian. We present results for both, the Humboldt glacier, which expands about 4×10^2 (km) laterally, and the Greenland ice sheet, which expands about 1.8×10^3 (km). The ice is at most 3.4 (km) thick, resulting in approximate aspect ratios of 8.5×10^{-3} for Humboldt and 1.9×10^{-3} for Greenland. We use a nonuniform triangulation of the Greenland ice sheet, with mesh size ranging from 1 to 10 (km), and we then extrude it in the vertical direction, obtaining a 3D mesh having ten layers of prismatic elements. The velocity observations at the top surface of the Greenland ice sheet are obtained from satellite observations [52]. The MAP basal sliding field and the temperature fields are obtained as part of the initialization process, using a numerical optimization approach to match the ice velocity observations and constrained by the first-order flow model coupled with a temperature model [53]. Additional details about the mesh geometries and data, in particular regarding the Humboldt glacier, can be found in [54].

In figure 9, we show the surface velocity observation data d in $(m \text{ yr}^{-1})$, the MAP point estimate of the logarithmic basal sliding field β^* $(\exp(\beta^*)$ is in $(k\text{Pa} \text{ yr} \text{ m}^{-1}))$ and surface velocity in $(m \text{ yr}^{-1})$ generated by the model.

6.1. HODLR compressibility

We next generate global LR approximations of a Greenland and Humboldt data-misfit Hessian as well as LR approximations of various off-diagonal blocks. Plots of the estimated singular values are provided in figure 10. We observe that the spectrum of the Greenland ice

Table 1. Problem specifications for the Humboldt glacier and Greenland ice-sheet problems (Example II). Dimension of the discretized basal sliding field $\dim(\beta)$, dimension of the discretized velocity field $\dim(u)$, dimension of the observation data $\dim(d)$, processors employed for computations and depth of the HODLR hierarchical partitioning L.

	Humboldt	Greenland
$\dim(\boldsymbol{\beta})$	11608	320116
$\dim(\mathbf{u})$	255 376	7042552
$\dim(\boldsymbol{d})$	23 216	640232
# of cores	120	2048
L	8	10

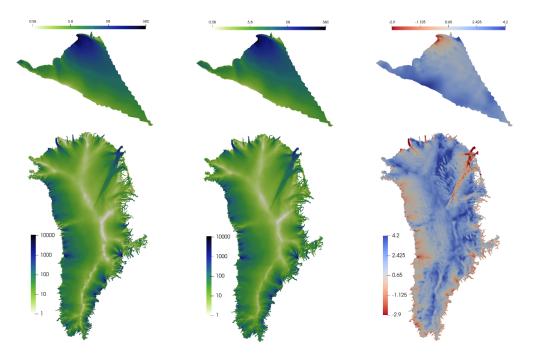


Figure 9. Data and MAP estimates for Example II. Shown are the surface velocity observation data (left), and the reconstructed surface velocity field (middle) that is based on the MAP point estimate of the logarithmic basal sliding field (right). Top row is for the Humboldt glacier and bottom row for the Greenland ice sheet.

sheet decays substantially slower than the one for the Humboldt glacier. Besides the different sizes of these two discretized problems, this is also due to the different aspect ratios. Having estimated singular values of the data-misfit Hessians and the appropriate off-diagonal blocks, one is able to estimate computational costs to compress them into the global LR and HODLR matrix formats. The computational cost as a function of Hessian approximation target accuracy is given in figure 11, wherein it is demonstrated that the HODLR compression format can offer a favorable means to approximate data-misfit Hessians for large-scale inverse problems governed by complex ice-sheet models.

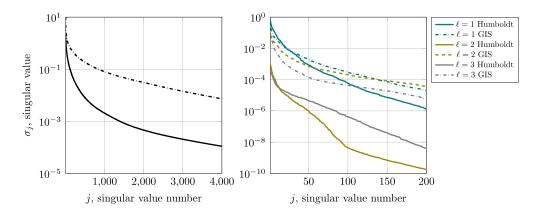


Figure 10. Singular values of the data-misfit Hessian (left figure) and various off-diagonal blocks of the data-misfit Hessian (right figure) for Example II. The color-scheme in the right most figure is consistent with figure 1. On the left, the singular values of the Humboldt and Greenland data-misfit Hessians are shown using a solid and dash-dotted line, respectively. On the right, we show the singular values of the upper most blocks, that is $A_{1,2}^{(\ell)}$ as defined in appendix 'HODLR approximation error due to the accumulation of LR approximations of off-diagonal blocks'.

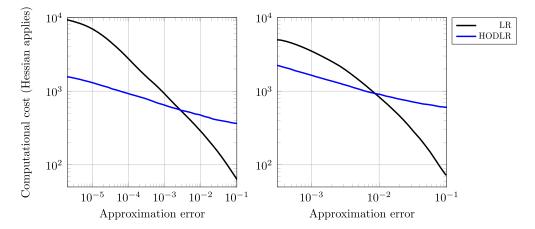


Figure 11. Estimated computational costs (measured by the number of Hessian applies) to compress the Humboldt glacier (left) and Greenland ice-sheet (right) data-misfit Hessians into the global LR and HODLR formats as a function of the approximation error $\|\boldsymbol{H}_{\text{misfit}} - \tilde{\boldsymbol{H}}_{\text{misfit}}\|_2 / \|\boldsymbol{H}_{\text{misfit}}\|_2$.

6.2. Impact of parameter degree of freedom ordering

We seek to ensure that the off-diagonal blocks, determined by the hierarchical partitioning described in section 2.2, of the data-misfit Hessian are LR. For this reason, the nodes $\{x_i\}_i$ associated to the dofs are ordered according to a kd-tree, i.e. a recursive hyperplane splitting. The ordering provided by the kd-tree is such that the (i,j)-entry of the distance matrix $D_{i,j} = \|x_i - x_j\|_2$, is typically small whenever |i - j| is small, that is the dof ordering preserves some notion of locality (see section 3.1). In particular, a sparse permutation matrix B, is determined, whose action reorders the dofs from the default ordering provided by the finite element

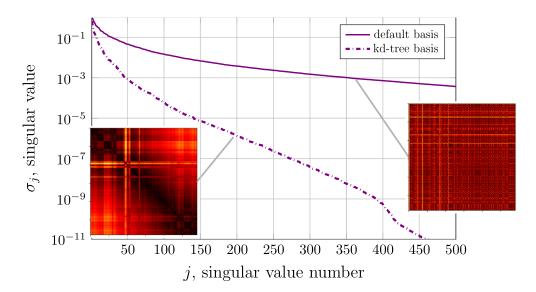


Figure 12. Singular values of the hierarchical level 1 off-diagonal block, $A_{1,2}^{(1)}$, of the Humboldt glacier data-misfit Hessian, when expressed in a kd-tree basis and the default basis. Shown also are heat maps of the distance matrices $D_{i,j} = ||x_i - x_j||_2$, wherein the nodes $\{x_i\}_i$ associated to the finite element degrees of freedom have been ordered according to a default standard and a kd-tree.

discretization to that specified by the kd-tree. The data-misfit Hessian with respect to the kdtree ordering, $H_{\text{misfit}}^{\text{kd}} := BH_{\text{misfit}} B^{\top}$, is then amenable to HODLR compression. Subsequently, $B^{\top} \tilde{H}_{\text{misfit}}^{\text{kd}} B$ is an approximation of the data-misfit Hessian with respect to the default ordering. The dof ordering has no impact on a matrix's global numerical rank but does indeed impact the numerical rank of its numerous submatrices that are defined by a fixed partitioning scheme, such as the off-diagonal blocks of an HODLR matrix (see section 2.2). Here, we study the HODLR compressibility of the Humboldt glacier data-misfit Hessian by comparing the rate of decay of an off-diagonal block's singular values using the default ordering provided by Albany and the ordering obtained by a kd-tree recursive hyperplane splitting. As observed in figure 12, the rate at which the singular values of the level-1 off-diagonal block decay, strongly depends on the dof ordering. This is because the ordering given by the kd-tree better preserves locality, and as a consequence, by the argument provided in section 3.1, the singular values decay much faster when using the kd-tree ordering. The kd-tree ordering therefore provides a substantially computationally cheaper means to generate an HODLR approximation of the data-misfit Hessian. Figure 12 also shows distance matrices for the default and kd-tree bases. These show the improved locality for the kd-orderings. Note that data-misfit Hessian matrices are expected to follow a similar structure as these distance matrices, which explains why the former's off-diagonal blocks can be compressed more effectively in the kd-order than in the default order of dofs.

7. Conclusion

In this work, we motivated why data-misfit Hessians which arise from a class of inverse problems governed by PDEs have HODLR matrix structure. HODLR matrices can efficiently be inverted and factorized, operations needed for solving inverse problems governed by PDEs by Newton's method and for MCMC sampling methods. We study inverse ice sheet problems, for which, under certain regimes, HODLR matrices provide a more computationally efficient approximation format than the global LR matrix format. These problems are those with highly informative data and small aspect ratio ice sheets. While global LR matrices are favorable for large discretized parameter dimension and small data dimension, we find that HODLR matrices can offer computational savings for large-scale inverse problems such as a Greenland ice sheet inverse problem with satellite observational data and a discretized parameter dimension that exceeds 10⁵.

The computational cost of each Hessian-vector product increases with the size of the problem and is a computational challenge that must be considered for larger-scale problems. However, HODLR remains a potential means to approximate the Hessian as the number of required Hessian-vector products has only mild logarithmic growth with respect to the problem size. For future work, we believe that the computational cost can be reduced further by utilizing hierarchical matrix partitionings that satisfy a strong admissibility condition [17], as they are better suited to exploit the data-misfit Hessian structure described in section 3.1. However, generating a hierarchical matrix approximation with such a partitioning, e.g. by the peeling method [11, 12], requires substantially more Hessian-vector products. Ultimately, to further reduce the computational cost of Hessian approximations in large-scale inverse problems governed by PDEs, exploiting further problem structure will be essential.

Data availability statement

The data cannot be made publicly available upon publication because they are not available in a format that is sufficiently accessible or reusable by other researchers. The data that support the findings of this study are available upon reasonable request from the authors.

Acknowledgments

The authors thank Trevor Hillebrand from Los Alamos National Laboratory for help with setting up the Humboldt and Greenland ice-sheet grids and datasets. Support for this work was provided by the National Science Foundation under Grant No. DMS-1840265 and CAREER-1654311 and through the SciDAC Project ProSPect, funded by the U.S. Department of Energy (DOE) Office of Science, Advanced Scientific Computing Research and Biological and Environmental Research programs. This research used resources of the National Energy Research Scientific Computing Center (NERSC), a U.S. Department of Energy Office of Science User Facility operated under Contract No. DE-AC02-05CH11231, under NERSC Award ERCAP0020130.

Disclaimer

This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government. Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of

Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc. for the U.S. Department of Energy's National Nuclear Security Administration under Contract DE-NA-0003525.

Appendix

Randomized compression algorithms

Here, for completeness we outline the matrix-free randomized double-pass global LR [32] and HODLR compression algorithms [12]. For conciseness we omit the single-pass algorithm, which exploits symmetry and that was used to compute global LR compressions for each example presented in this work. A description of the double-pass algorithm is included as it is an essential component of the HODLR compression algorithm. The essential ideas of the randomized double-pass global LR algorithm are

- (i) the application of a vector ω with random entries to a matrix A, yields a vector $y = A\omega$, which is likely aligned with the dominant left singular vectors of A;
- (ii) a matrix Q, whose columns are nearly aligned with the dominant left singular vectors of A, can be used to construct an accurate LR approximation $\tilde{A} = QQ^{\top}A$ of A.

The double-pass randomized singular value decomposition (SVD) algorithm is presented in algorithm 1 and does not significantly differ from that in [32], specifically it is lines 7,8 and 9 that are distinct. This minor modification frees us from the need to compute a (parallel) SVD of a (distributed) $N \times k$ matrix, such as \mathbf{Z} . Here, we only need to compute an SVD of the smaller $k \times k$ matrix $\mathbf{R}_{\mathbf{Z}}$. In the distributed memory parallelism setting of section 6, this algorithmic modification allows us to only require the invocation of serial SVD routines, on $\mathbf{R}_{\mathbf{Z}}$, which is typically small and available on each processor.

Algorithm 1. Double-pass randomized SVD. **Input:** $A \in \mathbb{R}^{N \times N}$, $r \in \mathbb{N}$ desired rank and oversampling parameter $q \in \mathbb{N}$. **Output:** low-rank approximation \tilde{A} of A

```
1: k = r + q
 2: \Omega = \text{randn}(N, k)
                                                                                           {Initiate random matrix}
 3: Y = A\Omega
                                                                                            {Sample column space}
 4: \mathbf{Q}_{\mathbf{Y}} = \mathtt{orthog}(\mathbf{Y})
                                                                             {Orthogonalize column samples}
 5: \mathbf{Z} = \mathbf{A}^{\top} \mathbf{Q}_{\mathbf{Y}}
                                                                                                  {Sample row space}
 6: Q_Z = \text{orthog}(Z)
                                                                                   {Orthogonalize row samples}
 7: R_Z = Q_Z^\top Z
                                                                                          {Compress row samples}
 8: \mathbf{R}_{\mathbf{Z}} = \hat{\mathbf{V}} \mathbf{\Sigma} \hat{\mathbf{U}}^{\top}
                                                       {SVD of k \times k compressed row sample matrix}
 9: V = Q_z \hat{V}
                                                                               {Project row space information}
10: U = \mathbf{Q}_{\mathbf{Y}} \hat{\mathbf{U}}
11: \tilde{\mathbf{A}} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^{\top}
                                                                         {Project column space information}
                                                                               {Form low-rank approximation}
```

The randomized HODLR algorithm proceeds by compressing off-diagonal blocks by the double-pass algorithm. The off-diagonal blocks are defined in terms of hierarchical partitions of the index set $\mathcal{I}_1^{(1)} = \{1,2,\ldots,N\}$, that is $\mathcal{I}_j^{(\ell)} = \{\frac{N}{2^\ell}(j-1)+1,\frac{N}{2^\ell}(j-1)+2,\ldots,\frac{N}{2^\ell}j\}$, for $1 \le j \le 2^\ell$ and $1 \le \ell \le L$. The larger off-diagonal blocks are compressed before the compression of the smaller off-diagonal blocks, via a peeling procedure [11]. Here, both A and A are

assumed to be symmetric as we seek compression of symmetric operators and computation of symmetric approximants. There are means to adaptively choose the desired ranks r_1, \ldots, r_L in algorithm 2 as discussed in section 2.2 and references within. The authors are not aware of an adaptive means to choose the hierarchical depth L. Such means could leverage the structure of the peeling algorithm in which the compression is done serially from the hierarchical level $\ell=1$ to $\ell=L$ and balance the HODLR compression cost as detailed in (4) and the approximation error as detailed in appendix 'HODLR approximation error due to the accumulation of LR approximations of off-diagonal blocks'.

Algorithm 2. Symmetric matrix-free randomized HODLR.

Input: symmetric $A \in \mathbb{R}^{N \times N}$, hierarchical depth $L \in \mathbb{N}$, r_1, \dots, r_L desired ranks of the off-diagonal blocks at each hierarchical depth and oversampling parameter q.

Output: symmetric HODLR approximation \hat{A} of A

```
1: for \ell = 1, 2, ..., L do
                 k_{\ell} = r_{\ell} + q
  2:
  3:
                 \Omega = \mathsf{zeros}(N, k_\ell)
                 \begin{array}{l} \mathbf{for}\, j=1,\dots,2^{\ell-1}\,\mathbf{do} \\ \mathbf{\Omega}(\mathcal{I}_{2j}^{(\ell)},:)=\mathrm{randn}\,(|\mathcal{I}_{2j}^{(\ell)}|,k_\ell) \\ \mathbf{end}\,\,\mathbf{for} \end{array}
  4:
  5:
                                                                                                                                         {Initiate structured random matrix}
  6:
                 \mathbf{Y} = \left(\mathbf{A} - \sum_{j=1}^{\ell-1} \mathbf{A}^{(j)}\right) \mathbf{\Omega}
  7:
                                                                                                                       {Sample off-diagonal block column spaces}
                 for j=1,\ldots,2^{\ell-1} do
  8:
                      \mathbf{Y}^{(j)} = \operatorname{zeros}(N, k_{\ell})
  9:
                      \mathbf{Y}^{(j)}(\mathcal{I}_{2i-1}^{(\ell)},:) = \mathbf{Y}(\mathcal{I}_{2i-1}^{(\ell)},:)
10:
                      Q_{\mathbf{Y}}^{(j)} = \operatorname{orthog}(\mathbf{Y}^{(j)})
11:
                                                                                          {Orthogonalize column samples of the level \ell off-diagonal
                                                                                                                                                                                                  blocks}
12:
                 \mathbf{Q}_{\mathbf{Y}} = \sum_{j=1}^{2^{\ell-1}} \mathbf{Q}_{\mathbf{Y}}^{(j)}
13:
                                                                                                                                                    {Row space sampling matrix}
                 \mathbf{Z} = \left(\mathbf{A} - \sum_{i=1}^{\ell-1} \mathbf{A}^{(j)}\right) \mathbf{Q}_{\mathbf{Y}}
14:
                                                                                                                               {Sample off-diagonal block row spaces}
                 for j = 1, ..., 2^{\ell-1} do
15:
                      \mathbf{Z}^{(j)} = \mathbf{Z}(\mathcal{I}_{2i}^{(\ell)},:)
16:
                      oldsymbol{Q}_{\mathbf{Z}}^{(j)} = \mathtt{orthog}\,(\mathbf{Z}^{(j)})
17:
                                                                                {Orthogonalize row samples of the level \ell off-diagonal blocks}
                     \mathbf{R}_{\mathbf{Z}}^{(j)} = \left(\mathbf{Q}_{\mathbf{Z}}^{(j)}\right)^{\top} \mathbf{Z}^{(j)}
18:
                                                                                                       {Compress level \ell off-diagonal block row samples}
                     m{R}_{m{z}}^{(j)} = m{\hat{V}}_{2i-1}^{(\ell)} m{\Sigma}_{2i-1}^{(\ell)} m{\hat{U}}_{2i-1}^{(\ell)}
19:
                                                                                                            {SVD of k_{\ell} \times k_{\ell} compressed row sample matrix}
                      V_{2i-1}^{(\ell)} = Q_Z^{(j)} \hat{V}_{2i-1}^{(\ell)}
20:
                                                                                                                                               {Project row space information}
                      U_{2i-1}^{(\ell)} = Q_{Y}^{(j)} \hat{U}_{2i-1}^{(\ell)}
21:
                                                                                                                                        {Project column space information}
                      V_{2j}^{(\ell)} = U_{2j-1}^{(\ell)}
22:
                      U_{2j}^{(\ell)} = V_{2j-1}^{(\ell)}
23:
                 egin{aligned} & \mathbf{\Sigma}_{2j}^{(\ell)} = \mathbf{\Sigma}_{2j-1}^{(\ell)} \ & \mathbf{end for} \ & \mathbf{A}^{(\ell)} = \sum_{j=1}^{2^\ell} \mathbf{U}_j^{(\ell)} \mathbf{\Sigma}_j^{(\ell)} \left(\mathbf{V}_j^{(\ell)}
ight)^{	op} \end{aligned}
24:
25:
26:
27: end for
28: obtain block diagonal D of A by sampling A - \sum_{i=1}^{L} A^{(i)}
29: \tilde{A} = D + \sum_{\ell=1}^{L} A^{(\ell)}
```

HODLR approximation error due to the accumulation of LR approximations of off-diagonal blocks

Let A be a $N \times N$ matrix and consider the following partitioning

$$\begin{split} \boldsymbol{A}^{(1)} &= \begin{pmatrix} \mathbf{0} & A_{1,2}^{(1)} \\ A_{2,1}^{(1)} & \mathbf{0} \end{pmatrix}, \\ \boldsymbol{A}^{(2)} &= \begin{pmatrix} \mathbf{0} & A_{1,2}^{(2)} & \mathbf{0} & \mathbf{0} \\ A_{2,1}^{(2)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & A_{3,4}^{(2)} \\ \mathbf{0} & \mathbf{0} & A_{4,3}^{(2)} & \mathbf{0} \end{pmatrix}, \\ \boldsymbol{D} &= \begin{pmatrix} A_{1,1}^{(2)} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & A_{2,2}^{(2)} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & A_{3,3}^{(2)} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & A_{4,4}^{(2)} \end{pmatrix}, \end{split}$$

where $A_{i,j}^{(\ell)}$ is the (i,j) block of a $2^\ell \times 2^\ell$ block partitioning of A, where $1 \leqslant \ell \leqslant L$. $A^{(\ell)}$ contains all blocks $A_{i,j}^{(\ell)}$ such that |i-j|=1 and D contains the diagonal blocks $A_{i,i}^{(L)}$. Above, we show the decomposition $A=\sum_{\ell=1}^L A^{(\ell)}+D$ for L=2 hierarchical depth but in the following analysis L is a arbitrary. Let $\mathbf{x} \in \mathbb{R}^N$, then

$$Ax = \sum_{j=1}^{L} A^{(j)}x + Dx,$$

$$A^{(1)}x = \begin{pmatrix} A_{1,2}^{(1)}x_{2}^{(1)} \\ A_{2,1}^{(1)}x_{1}^{(1)} \end{pmatrix}, x = \begin{pmatrix} x_{1}^{(1)} \\ x_{2}^{(2)} \end{pmatrix},$$

$$A^{(j)}x = \begin{pmatrix} A_{1,2}^{(j)}x_{2}^{(j)} \\ A_{2,1}^{(j)}x_{1}^{(j)} \\ \vdots \\ A_{2j-1,2j}^{(j)}x_{2j}^{(j)} \\ A_{2j}^{(j)}x_{2j-1}^{(j)} \end{pmatrix}, x = \begin{pmatrix} x_{1}^{(j)} \\ x_{2}^{(j)} \\ \vdots \\ x_{2j-1}^{(j)} \\ x_{2j}^{(j)} \end{pmatrix},$$

from which we obtain the following expression

$$\|\boldsymbol{A}^{(j)}\boldsymbol{x}\|_{2}^{2} = \sum_{k=1}^{2^{j-1}} \left(\|\boldsymbol{A}_{2k-1,2k}^{(j)}\boldsymbol{x}_{2k}^{(j)}\|_{2}^{2} + \|\boldsymbol{A}_{2k,2k-1}^{(j)}\boldsymbol{x}_{2k-1}^{(j)}\|_{2}^{2} \right).$$

Now assume that \tilde{A} is an HODLR approximation of A, whose diagonal D is equal to the diagonal of A so that

$$(\mathbf{A} - \tilde{\mathbf{A}}) = \sum_{j=1}^{L} \Delta \mathbf{A}^{(j)},$$

$$\Delta \mathbf{A}^{(j)} := (\mathbf{A}^{(j)} - \tilde{\mathbf{A}}^{(j)}).$$

Here, it is assumed that each off-diagonal block at level $j=1,2,\ldots,L$ has been approximated to some absolute tolerance $\varepsilon_j>0$, so that $\|\Delta A_{2k-1,2k}^{(j)}\|_2, \|\Delta A_{2k,2k-1}^{(j)}\| \le \varepsilon_j$ for each $k=1,2,\ldots,2^{j-1}$. For $x\in\mathbb{R}^N$ we have

$$\begin{split} \|\left(\boldsymbol{A} - \tilde{\boldsymbol{A}}\right)\boldsymbol{x}\|_{2} & \leq \sum_{j=1}^{L} \|\Delta\boldsymbol{A}^{(j)}\boldsymbol{x}\|_{2}, \\ \|\Delta\boldsymbol{A}^{(j)}\boldsymbol{x}\|_{2} & = \sqrt{\sum_{k=1}^{2^{j-1}} \left(\|\Delta\boldsymbol{A}_{2k-1,2k}^{(j)}\boldsymbol{x}_{2k}^{(j)}\|_{2}^{2} + \|\Delta\boldsymbol{A}_{2k,2k-1}^{(j)}\boldsymbol{x}_{2k-1}^{(j)}\|_{2}^{2}\right)} \\ & \leq \sqrt{\sum_{k=1}^{2^{j-1}} \left(\varepsilon_{j}^{2} \|\boldsymbol{x}_{2k}^{(j)}\|_{2}^{2} + \varepsilon_{j}^{2} \|\boldsymbol{x}_{2k-1}^{(j)}\|_{2}^{2}\right)}, \\ \|\Delta\boldsymbol{A}^{(j)}\boldsymbol{x}\|_{2} & \leq \varepsilon_{j} \sqrt{\sum_{k=1}^{2^{j-1}} \left(\|\boldsymbol{x}_{2k}^{(j)}\|_{2}^{2} + \|\boldsymbol{x}_{2k-1}^{(j)}\|_{2}^{2}\right)} = \varepsilon_{j} \|\boldsymbol{x}\|_{2}, \\ \|\left(\boldsymbol{A} - \tilde{\boldsymbol{A}}\right)\boldsymbol{x}\|_{2} & \leq \|\boldsymbol{x}\|_{2} \sum_{j=1}^{L} \varepsilon_{j}, \\ \|\boldsymbol{A} - \tilde{\boldsymbol{A}}\|_{2} &:= \sup_{\boldsymbol{x} \neq \boldsymbol{0}} \left(\frac{\|\left(\boldsymbol{A} - \tilde{\boldsymbol{A}}\right)\boldsymbol{x}\|_{2}}{\|\boldsymbol{x}\|_{2}}\right) \leq \sum_{j=1}^{L} \varepsilon_{j}. \end{split}$$

Error of the Laplace posterior covariance due to approximation of the prior-preconditioned data-misfit Hessian

Consider a symmetric matrix $A \in \mathbb{R}^{N \times N}$, whose eigenvalues are bounded below by a number greater than -1 and a symmetric approximant \tilde{A} , with discrepancy $\Delta A = A - \tilde{A}$. We signify a generic eigenvalue of S by $\lambda(S)$ so that $s_1 \leq \lambda(S) \leq s_2$ indicates that all eigenvalues of S are bounded below by s_1 and above by s_2 . Next we provide an upper bound for the error of $(I+A)^{-1} - (I+\tilde{A})^{-1}$, given that $\|\Delta A\|_2 = \varepsilon$. When, as in section 3.2, A is the prior-preconditioned data-misfit Hessian $\|(I+A)^{-1} - (I+\tilde{A})^{-1}\|_2$ quantifies the error of the covariance of the Laplace approximation of the posterior distribution that is introduced through HODLR compression

$$(I+A)^{-1} - (I+\tilde{A})^{-1} = (I+A)^{-1} - (I+A-\Delta A)^{-1}$$

$$= (I+A)^{-1} - ((I+A)(I-(I+A)^{-1}\Delta A))^{-1}$$

$$= (I+A)^{-1} - (I-(I+A)^{-1}\Delta A)^{-1}(I+A)^{-1}$$

$$= (I-(I-(I+A)^{-1}\Delta A)^{-1})(I+A)^{-1}.$$

Given that $\|\Delta A\|_2 = \varepsilon$, we have

$$\begin{split} &-\varepsilon \leqslant \lambda \left(\Delta \boldsymbol{A}\right) \leqslant \varepsilon, \\ &-\varepsilon^* \leqslant \lambda \left(\left(\boldsymbol{I} + \boldsymbol{A} \right)^{-1} \Delta \boldsymbol{A} \right) \leqslant \varepsilon^*, \\ &\varepsilon^* := \varepsilon (1 + \lambda_{\min}(\boldsymbol{A}))^{-1}, \\ &1 + \varepsilon^* \geqslant \lambda \left(\boldsymbol{I} - \left(\boldsymbol{I} + \boldsymbol{A} \right)^{-1} \Delta \boldsymbol{A} \right) \geqslant 1 - \varepsilon^*, \end{split}$$

we next assume $\varepsilon^* < 1$, so that the eigenvalues of $I - (I + A)^{-1} \Delta A$ are necessarily positive

$$(1+\varepsilon^*)^{-1} \leqslant \lambda \left(\left(\boldsymbol{I} - (\boldsymbol{I} + \boldsymbol{A})^{-1} \Delta \boldsymbol{A} \right)^{-1} \right) \leqslant (1-\varepsilon^*)^{-1}.$$

With this it follows that

$$\| (\mathbf{I} + \mathbf{A})^{-1} - (\mathbf{I} + \tilde{\mathbf{A}})^{-1} \|_{2} / \| (\mathbf{I} + \mathbf{A})^{-1} \|_{2} \leq (1 - (1 + \varepsilon^{*})^{-1})$$

$$\| (\mathbf{I} + \mathbf{A})^{-1} - (\mathbf{I} + \tilde{\mathbf{A}})^{-1} \|_{2} / \| (\mathbf{I} + \mathbf{A})^{-1} \|_{2} \leq \frac{\varepsilon^{*}}{1 + \varepsilon^{*}},$$

where, as before $\varepsilon^* = \|\Delta A\|_2/(1 + \lambda_{\min}(A))$.

ORCID iDs

Tucker Hartland https://orcid.org/0000-0002-4638-3209
Georg Stadler https://orcid.org/0000-0001-7762-6544
Mauro Perego https://orcid.org/0000-0002-2671-8032
Kim Liegeois https://orcid.org/0000-0002-1182-4078
Noémi Petra https://orcid.org/0000-0002-9491-0034

References

- [1] Isaac T, Petra N, Stadler G and Ghattas O 2015 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 J. Comput. Phys. 296 348–68
- [2] Petra N, Martin J, Stadler G and Ghattas O 2014 A computational framework for infinite-dimensional Bayesian inverse problems: part II. Stochastic Newton MCMC with application to ice sheet flow inverse problems SIAM J. Sci. Comput. 36 A1525–55
- [3] Spantini A, Solonen A, Cui T, Martin J, Tenorio L and Marzouk Y 2015 Optimal low-rank approximations of Bayesian linear inverse problems SIAM J. Sci. Comput. 37 A2451–87
- [4] Flath H P, Wilcox L C, Akçelik V, Hill J, van Bloemen Waanders B and Ghattas O 2011 Fast algorithms for Bayesian uncertainty quantification in large-scale linear inverse problems based on low-rank partial Hessian approximations SIAM J. Sci. Comput. 33 407–32
- [5] Bui-Thanh T, Ghattas O, Martin J and Stadler G 2013 A computational framework for infinitedimensional Bayesian inverse problems: part I. The linearized case, with application to global seismic inversion SIAM J. Sci. Comput. 35 A2494–523
- [6] Saibaba A K and Kitanidis P K 2015 Fast computation of uncertainty quantification measures in the geostatistical approach to solve inverse problems Adv. Water Resour. 82 124–38
- [7] Martinsson P-G 2011 A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix SIAM J. Matrix Anal. Appl. 32 1251–74
- [8] Alger N, Rao V, Meyers A, Bui-Thanh T and Ghattas O 2019 Scalable matrix-free adaptive product-convolution approximation for locally translation-invariant operators SIAM J. Sci. Comput. 41 A2296–328

- [9] Zhu H, Li S, Fomel S, Stadler G and Ghattas O 2016 A Bayesian approach to estimate uncertainty for full waveform inversion with a priori information from depth migration *Geophysics* 81 R307–23
- [10] Alger N, Hartland T, Petra N and Ghattas O 2023 Point spread function approximation of high rank Hessians with locally supported non-negative integral kernels (in preparation)
- [11] Lin L, Lu J and Ying L 2011 Fast construction of hierarchical matrix representation from matrixvector multiplication J. Comput. Phys. 230 4071–87
- [12] Martinsson P-G 2016 Compressing rank-structured matrices via randomized sampling SIAM J. Sci. Comput. 38 A1959–86
- [13] Geoga C J, Anitescu M and Stein M L 2020 Scalable Gaussian process computations using hierarchical matrices J. Comput. Graph. Stat. 29 227–37
- [14] Litvinenko A, Sun Y, Genton M G and Keyes D E 2019 Likelihood approximation with hierarchical matrices for large spatial datasets Comput. Stat. Data Anal. 137 115–32
- [15] Ambartsumyan I, Boukaram W, Bui-Thanh T, Ghattas O, Keyes D, Stadler G, Turkiyyah G and Zampini S 2020 Hierarchical matrix approximations of Hessians arising in inverse problems governed by PDEs SIAM J. Sci. Comput. 42 A3397–426
- [16] Hackbusch W 1999 A sparse matrix arithmetic based on H-matrices. Part I: introduction to H-matrices Computing 62 89–108
- [17] Hackbusch W and Börm S 2002 Data-sparse approximation by adaptive \mathcal{H}^2 -matrices Computing 69 1–35
- [18] Tarantola A 2005 Inverse Problem Theory and Methods for Model Parameter Estimation (Philadelphia, PA: SIAM)
- [19] Kaipio J and Somersalo E 2006 Statistical and Computational Inverse Problems vol 160 (Berlin: Springer)
- [20] Stuart A M 2010 Inverse problems: a Bayesian perspective Acta Numer. 19 451–559
- [21] Nocedal J and Wright S J 2006 Numerical Optimization 2nd edn (Berlin: Springer)
- [22] Borzì A and Schulz V 2011 Computational Optimization of Systems Governed by Partial Differential Equations (Philadelphia, PA: SIAM)
- [23] Gunzburger M D 2002 Perspectives in Flow Control and Optimization (Philadelphia, PA: SIAM)
- [24] Petra Nemi and Sachs E W 2021 Second order adjoints in optimization *Numerical Analysis and Optimization* ed M Al-Baali, A Purnama and L Grandinetti (Cham: Springer) pp 209–30
- [25] Keith Hastings W 1970 Monte Carlo sampling methods using Markov chains and their applications Biometrika 57 97–109
- [26] Robert C P and Casella G 1999 Monte Carlo Statistical Methods vol 2 (Berlin: Springer)
- [27] Rudolf D and Sprungk Born 2018 On a generalization of the preconditioned Crank–Nicolson Metropolis algorithm Found. Comput. Math. 18 309–43
- [28] Pinski F J, Simpson G, Stuart A M and Weber H 2015 Algorithms for Kullback-Leibler approximation of probability measures in infinite dimensions SIAM J. Sci. Comput. 37 A2733-57
- [29] Kim K-T, Villa U, Parno M, Marzouk Y, Ghattas O and Petra N 2023 hIPPYlib-MUQ: a Bayesian inference software framework for integration of data with complex predictive models under uncertainty ACM Trans. Math. Softw. 49 17
- [30] Ambikasaran S and Darve E 2013 An $\mathcal{O}(n \log n)$ fast direct solver for partial hierarchically semi-separable matrices *J. Sci. Comput.* **57** 477–501
- [31] Ambikasaran S, O'Neil M and Singh K R 2014 Fast symmetric factorization of hierarchical matrices with applications (arXiv:1405.0223)
- [32] Halko N, Martinsson P G and Tropp J A 2011 Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions SIAM Rev. 53 217–88
- [33] Xi Y, Xia J and Chan R 2014 A fast randomized eigensolver with structured LDL factorization update SIAM J. Matrix Anal. Appl. 35 974–96
- [34] Boukaram W, Turkiyyah G and Keyes D 2019 Randomized GPU algorithms for the construction of hierarchical matrices from matrix-vector operations SIAM J. Sci. Comput. 41 C339–66
- [35] Gorman C, Chávez G, Ghysels P, Mary T, Rouet Fçois-H and Li X S 2019 Robust and accurate stopping criteria for adaptive randomized sampling in matrix-free hierarchically semiseparable construction SIAM J. Sci. Comput. 41 S61–S85
- [36] Gelfand I M and Fomin S V 1963 Calculus of Variations (New York: Dover)
- [37] Bramble J H 1993 Multigrid Methods vol 294 (Boca Raton, FL: CRC Press)
- [38] Cuffey K M and Paterson W S B 2010 The Physics of Glaciers (New York: Academic)

- [39] Dukowicz J K, Price S F and Lipscomb W H 2010 Consistent approximations and boundary conditions for ice-sheet dynamics from a principle of least action *J. Glaciol.* **56** 480–96
- [40] Larour E, Seroussi H, Morlighem M and Rignot E 2012 Continental scale, high order, high spatial resolution, ice sheet modeling using the Ice Sheet System Model (ISSM) J. Geophys. Res. Earth Surf. 117 F01022
- [41] Morlighem M, Rignot E, Seroussi H, Larour E, Ben Dhia H and Aubry D 2010 Spatial patterns of basal drag inferred using control methods from a full-Stokes and simpler models for Pine Island Glacier, West Antarctica Geophys. Res. Lett. 37 L14502
- [42] Perego M, Price S and Stadler G 2014 Optimal initial conditions for coupling ice sheet models to Earth system models *J. Geophys. Res. Earth Surf.* 119 1894–917
- [43] Petra N, Zhu H, Stadler G, Hughes T J R and Ghattas O 2012 An inexact Gauss-Newton method for inversion of basal sliding and rheology parameters in a nonlinear Stokes ice sheet model J. Glaciol. 58 889–903
- [44] Glen J W 1955 The creep of polycrystalline ice Proc. R. Soc. A 228 519-38
- [45] Pattyn F et al 2008 Benchmark experiments for higher-order and full-Stokes ice sheet models (ISMIP-HOM) Cryosphere 2 95–108
- [46] Daon Y and Stadler G 2018 Mitigating the influence of boundary conditions on covariance operators derived from elliptic PDEs *Inverse Problems Imaging* 12 1083–102
- [47] Lindgren F, Rue Håvard and Lindström J 2011 An explicit link between Gaussian fields and Gaussian Markov random fields: the stochastic partial differential equation approach J. R. Stat. Soc. B 73 423–98
- [48] Hoffman M J, Perego M, Price S F, Lipscomb W H, Zhang T, Jacobsen D, Tezaur I, Salinger A G, Tuminaro R and Bertagna L 2018 MPAS-Albany land ice (MALI): a variable-resolution ice sheet model for Earth system modeling using Voronoi grids Geosci. Model Dev. 11 3747–80
- [49] Tezaur I K, Perego M, Salinger A G, Tuminaro R S and Price S F 2015 Albany/FELIX: a parallel, scalable and robust, finite element, first-order Stokes approximation ice sheet solver built for advanced analysis *Geosci. Model Dev.* 8 1197–220
- [50] Liegeois K, Perego M and Hartland T 2023 PyAlbany: a Python interface to the C++ multiphysics solver Albany J. Comput. Appl. Math. 425 115037
- [51] The Trilinos Project Team 2020 *The Trilinos Project Website* (available at: https://trilinos.github.io)
- [52] Joughin I, Smith B, Howat I and Scambos T 2015 MEaSUREs Greenland ice sheet velocity map from InSAR data, version 2
- [53] Perego M 2022 Large-scale PDE-constrained optimization for ice sheet model initialization SIAM News Online
- [54] Hillebrand T R, Hoffman M J, Perego M, Price S F and Howat I M 2022 The contribution of Humboldt Glacier, North Greenland, to sea-level rise through 2100 constrained by recent observations of speedup and retreat Cryosphere Discuss. 2022 1–33