

1   **POINT SPREAD FUNCTION APPROXIMATION OF HIGH RANK**  
2   **HESSIANS WITH LOCALLY SUPPORTED NON-NEGATIVE**  
3   **INTEGRAL KERNELS\***

4   NICK ALGER<sup>†</sup>, TUCKER HARTLAND<sup>‡</sup>, NOEMI PETRA<sup>‡</sup>, AND OMAR GHATTAS<sup>§</sup>

5   **Abstract.** We present an efficient matrix-free point spread function (PSF) method for approximating operators that have locally supported non-negative integral kernels. The **PSF-based** method  
6   computes impulse responses of the operator at scattered points, and interpolates these impulse re-  
7   sponses to approximate entries of the integral kernel. To compute impulse responses efficiently, we  
8   apply the operator to Dirac combs associated with batches of point sources, which are chosen by  
9   solving an ellipsoid packing problem. The ability to rapidly evaluate kernel entries allows us to con-  
10   struct a hierarchical matrix (H-matrix) approximation of the operator. Further matrix computations  
11   are then performed with fast H-matrix methods. **This end-to-end procedure is illustrated on a blur**  
12   **problem.** We demonstrate the **effectiveness of this PSF-based method's effectiveness** by using it to  
13   build preconditioners for the Hessian operator arising in two inverse problems governed by partial  
14   differential equations (PDEs): inversion for the basal friction coefficient in an ice sheet flow problem  
15   and for the initial condition in an advective-diffusive transport problem. While for many ill-posed  
16   inverse problems the Hessian of the data misfit term exhibits a low rank structure, and hence a low  
17   rank approximation is suitable, for many problems of practical interest the numerical rank of the  
18   Hessian is still large. The Hessian impulse responses on the other hand typically become more local  
19   as the numerical rank increases, which benefits the PSF-based method. Numerical results reveal  
20   that the **PSF-based** preconditioner clusters the spectrum of the preconditioned Hessian near one,  
21   yielding roughly  $5 \times 10 \times$  reductions in the required number of PDE solves, as compared to classical  
22   regularization-based preconditioning and no preconditioning. We also present a comprehensive nu-  
23   merical study for the influence of various parameters (that control the shape of the impulse responses  
24   and the rank of the Hessian) on the effectiveness of the advection-diffusion Hessian approximation.  
25   The results show that the PSF-based **preconditioners are method is** able to form good approximations  
26   of high-rank Hessians using only a small number of operator applications.  
27

28   **Key words.** data scalability, Hessian, hierarchical matrix, high-rank, impulse response, local  
29   translation invariance, matrix-free, moment methods, operator approximation, PDE constrained  
30   inverse problems, point spread function, preconditioning, product convolution

31   **AMS subject classifications.** 35R30, 41A35, 47A52, 47J06, 65D12, 65F08, 65F10, 65K10,  
32   65N21, 86A22, 86A40

33   **1. Introduction.** We present an efficient *matrix-free* point spread function (PSF)  
34   method for approximating operators  $\mathcal{A} : L^2(\Omega) \rightarrow L^2(\Omega)'$  that have locally supported  
35   non-negative integral kernels. Here,  $\Omega \subset \mathbb{R}^d$  is a bounded domain, and  $L^2(\Omega)'$  is the  
36   space of real-valued continuous linear functionals on  $L^2(\Omega)$ . By “**non-negative integral**  
37   **kernel,** we mean that entries of  $\mathcal{A}$ 's integral kernel are non-negative numbers; this is  
38   not the same as positive semi-definiteness of  $\mathcal{A}$ . Such operators appear, for instance,  
39   as Hessians in optimization and inverse problems governed by partial differential equa-  
40   tions (PDEs) [14, 20, 47], Schur complements in Schur complement methods for solv-  
41   ing partial differential equations (**PDEs**) and Poincare-Steklov operators in domain

---

\*Submitted to the editors DATE.

**Funding:** This research was supported by the National Science Foundation under Grant No. DMS-1840265 and CAREER-1654311, DOD grant FA9550-21-1-0084, and DOE grants DE-SC0021239 and DE-SC0019303.

<sup>†</sup>Oden Institute for Computational Engineering and Sciences, The University of Texas at Austin, Austin, TX ([nalger@oden.utexas.edu](mailto:nalger@oden.utexas.edu)).

<sup>‡</sup>Department of Applied Mathematics, University of California, Merced, Merced, CA. ([tchartland@ucmerced.edu](mailto:tchartland@ucmerced.edu), [npetra@ucmerced.edu](mailto:npetra@ucmerced.edu)).

<sup>§</sup>Oden Institute for Computational Engineering and Sciences and Walker Department of Mechanical Engineering, The University of Texas at Austin, Austin, TX ([omar@oden.utexas.edu](mailto:omar@oden.utexas.edu)).

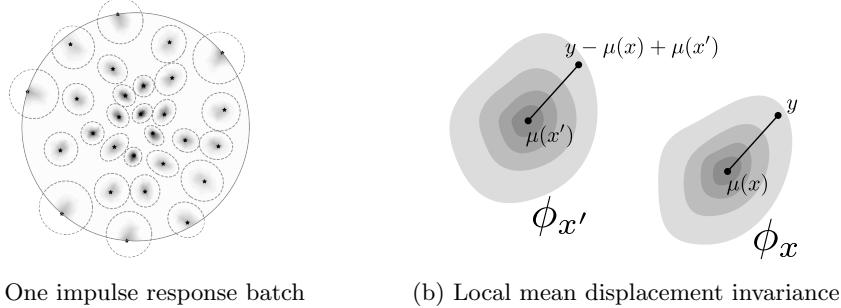


Fig. 1: (1a) One batch,  $\eta_b$ , of normalized impulse responses,  $\phi_x$ , that arise from applying  $\mathcal{A}$  to a weighted sum of scattered point sources (see Section 5.2). Here,  $\mathcal{A}$  is the ice sheet inverse problem data misfit Gauss-Newton Hessian described in Section 7. Black stars are point source locations. Shading shows the magnitude of the normalized impulse responses (darker means larger function values). Dashed gray ellipses are estimated impulse response support ellipsoids based on the moment method in Section 4.1. The large circle is  $\partial\Omega$ . (1b) Illustration of impulse responses,  $\phi_x$  and  $\phi_{x'}$ , corresponding to points  $x$  and  $x'$ . The operator  $\mathcal{A}$  is locally mean displacement invariant (Section 4.2) if  $\phi_x(y) \approx \phi_{x'}(y - \mu(x) + \mu(x'))$  when  $x$  is close to  $x'$ . Here,  $\mu(z)$  denotes the mean (center of mass) of  $\phi_z$ .

42 decomposition methods (e.g., Dirichlet-to-Neumann maps) [16, 67, 73], covariance op-  
43 erators in spatial statistics [17, 36, 37, 56], and blurring operators in imaging [22, 60].  
44 Here, “matrix-free” means that we may apply  $\mathcal{A}$  and its transpose<sup>1</sup>,  $\mathcal{A}^T$ , to functions,

45 (1.1) 
$$u \mapsto \mathcal{A}u \quad \text{and} \quad w \mapsto \mathcal{A}^T w,$$

46 via a black box computational procedure, but cannot easily access entries of  $\mathcal{A}$ ’s  
47 integral kernel. Evaluating the maps in (1.1) may require solving a subproblem that  
48 involves PDEs, or performing other costly computations. By “non-negative integral  
49 kernel,” we mean that entries of the integral kernel are non-negative numbers; this is  
50 not the same as positive semi-definiteness of  $\mathcal{A}$ .

51 We

52 The idea of the proposed method, which we refer to throughout the paper as  
53 the “PSF-based method,” is to use *impulse response interpolation* to form a high  
54 rank approximation of  $\mathcal{A}$  using a small number of operator applications. The impulse  
55 response,  $\phi_x$ , associated with a point,  $x$ , is the Riesz representation<sup>2</sup> of the linear  
56 functional that results from applying  $\mathcal{A}$  to a delta distribution (i.e., point source,  
57 impulse) centered at  $x$ . We compute batches of impulse responses by applying  $\mathcal{A}$   
58 to weighted sums of delta distributions associated with batches of points scattered  
59 throughout the domain (see Figure 1a). Then we interpolate translated and scaled  
60 versions of these impulse responses to approximate entries of the operator’s integral

<sup>1</sup>Recall that  $\mathcal{A}^T : L^2(\Omega) \rightarrow L^2(\Omega)'$  is the unique operator satisfying  $(\mathcal{A}u)(w) = (\mathcal{A}^T w)(u)$  for all  $u, w \in L^2(\Omega)$ , where  $\mathcal{A}u \in L^2(\Omega)'$  is the result of applying  $\mathcal{A}$  to  $u \in L^2(\Omega)$ , and  $(\mathcal{A}u)(w)$  is the result of applying that linear functional to  $w \in L^2(\Omega)$ , and similar for operations with  $\mathcal{A}^T$ .

<sup>2</sup>Recall that the Riesz representative of a functional  $\rho \in L^2(\Omega)'$  with respect to the  $L^2$  inner product is the unique function  $\rho^* \in L^2(\Omega)$  such that  $\rho(w) = (\rho^*, w)_{L^2(\Omega)}$  for all  $w \in L^2(\Omega)$ .

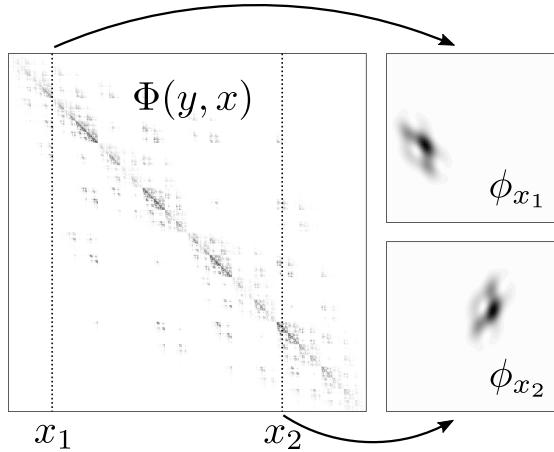


Fig. 2: Left: Matrix created by evaluating the integral kernel  $\Phi$  for  $\mathcal{A}$  (Equation 3.1) at all pairs of mesh vertices. This illustration is for the integral kernel in Equation 7.4. Dark colors indicate large entries and light colors indicate small entries. Rows and columns are ordered according to a kd-tree hierarchical clustering. Right: Impulse responses associated with points  $x_1, x_2 \in \Omega$ , shown by the two dotted vertical lines. Intuitively, one may think of impulse responses as ‘‘columns’’ of the integral kernel.

61 kernel. Picking the batches of points requires us to estimate the supports of the  
 62 impulse responses  $\phi_x$  before they are computed. Ellipsoid Batches of impulse responses  
 63 may be thought of intuitively as sets of ‘‘columns’’ of the kernel (Figure 2). To choose  
 64 the batches, we form ellipsoid estimates for the supports of all  $\phi_x$  are determined  
 65 a-priori via a moment method (Figure 3) that involves applying  $\mathcal{A}^T$  to a small number  
 66 of polynomial functions polynomials (see Section 4.1). We are inspired by resolution  
 67 analysis in seismic imaging, in which  $\mathcal{A}^T$  is applied to a random noise function, and  
 68 the width of the support  $\phi_x$  is estimated to be the autocorrelation length of the  
 69 resultant function near  $x$  [31, 74]. The moment method that we use estimates the  
 70 support of  $\phi_x$  more accurately than random noise probing in resolution analysis, at  
 71 the cost of the additional constraint that  $\mathcal{A}$  has a non-negative integral kernel then use  
 72 a greedy ellipsoid packing algorithm (Figure 4) to maximize the number of impulse  
 73 responses per batch. Then we interpolate translated and scaled versions of these  
 74 impulse responses to approximate entries of the operator’s integral kernel (Figure 5).  
 75 Adding more batches yields impulse responses at more points, increasing the accuracy  
 76 of the approximation approximation accuracy at the cost of one operator application  
 77 per batch (Figure 6).

78 The proposed operator approximation method may be categorized as a PSF  
 79 method that

80 The PSF-based method we propose is loosely based on ‘‘product convolution’’  
 81 (PC) approximations, which are approximations of an operator by weighted sums of  
 82 convolution operators with spatially varying weights. PC and PSF methods have a  
 83 long history dating back several decades. We note the following papers (among many  
 84 others) in which the convolution kernels are constructed from sampling impulse re-  
 85 sponds of the operator to scattered point sources: [1, 5, 12, 27, 29, 30, 32, 60, 78]. For  
 86 background on PC and PSF methods, we recommend the following papers: [23, 28, 35].

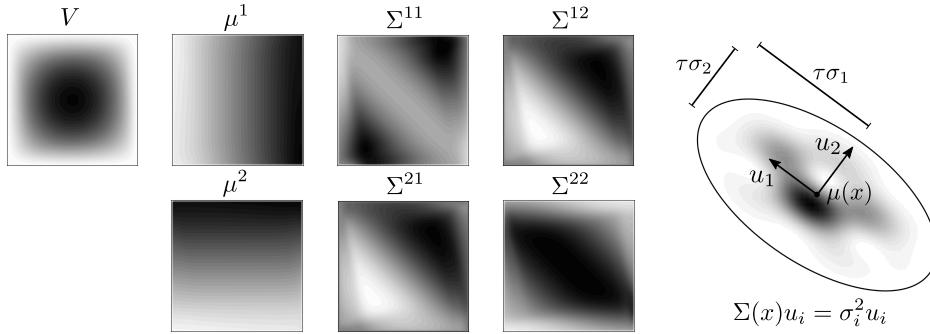


Fig. 3: Left: Impulse response moments. Scaling factor ( $V$ ), mean ( $\mu$ ), and covariance ( $\Sigma$ ). For each point  $x \in \Omega$ , the quantity  $V(x)$  is the integral of  $\phi_x$  over  $\Omega$ ,  $\mu(x)$  is the location that  $\phi_x$  is centered at, and  $\Sigma(x)$  is a matrix with eigenvectors and eigenvalues that characterize the width of the support of  $\phi_x$  about  $\mu(x)$  (see Section 4.1). Right: Ellipsoid support for an impulse response. This ellipsoid is the set of points within  $\tau$  standard deviations of the mean of the Gaussian distribution with mean  $\mu(x)$  and covariance  $\Sigma(x)$ . The scaling factor  $V(x)$  characterizes the magnitude of  $\phi_x$ .

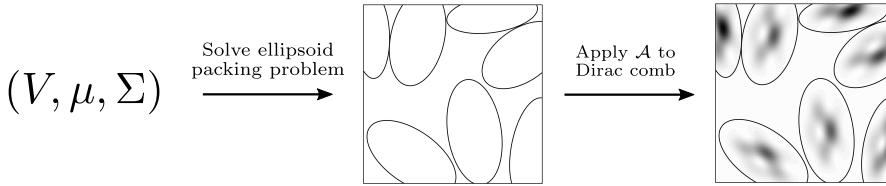


Fig. 4: Illustration of the process to compute one impulse response batch. Impulse response moments are first used to form ellipsoid shaped estimates of the supports of impulse responses (Equation 4.6). Then, an ellipsoid packing problem is solved to choose batches of non-overlapping support ellipsoids (Section 5.1). Finally,  $\mathcal{A}$  is applied to a Dirac comb associated with the points  $x_i$ , which correspond to the ellipsoids (Section 5.2). The process is repeated to form more batches.

87 The proposed PSF-based method improves upon existing PC and PSF methods in  
 88 the following ways: (1) While PC approximations are and PSF approximations are  
 89 typically based on an assumption of local translation invariance, the method we pro-  
 90 pose is based on a more general assumption we call “local mean displacement in-  
 91 variance” (Section 4.2 and Figure 1b), which improves the interpolation of the im-  
 92 pulse responses. (2) In our previous work [5], we chose point sources optimally in  
 93 an adaptive grid via a sequential adaptive procedureprocedure; the refinements to  
 94 the adaptive grid were chosen to maximally reduce the error at each step. However,  
 95 in that work each point source required a separate operator application, making the  
 96 previous method expensive when a large number of impulse responses is desired. In  
 97 this paper, we use a new moment method (Section 4.1) which permits computation  
 98 of many impulse responses (e.g., 50) per operator application. We are inspired by  
 99 resolution analysis in seismic imaging, in which  $\mathcal{A}^T$  is applied to a random noise  
 100 function, and the width of the support of  $\phi_x$  is estimated to be the autocorrelation

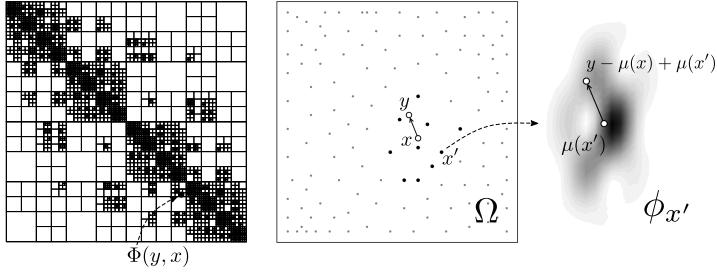


Fig. 5: Left: H-matrix structure for  $\Phi$ . Computing an entry of this matrix requires evaluating the integral kernel,  $\Phi(y, x)$ , at a pair of points  $(y, x) \in \Omega \times \Omega$ . Center: Kernel evaluation points  $x$  and  $y$  (black circles), sample points for the approximation (light gray and black dots), and the  $k_n$  sample points,  $x'$ , that are nearest to  $x$  (black dots). Right: Known impulse response at  $x'$ . Using radial basis function interpolation, the desired kernel entry is approximated as a weighted linear combination of translated and scaled versions of impulse responses at the points  $x'$  (Section 5.3).

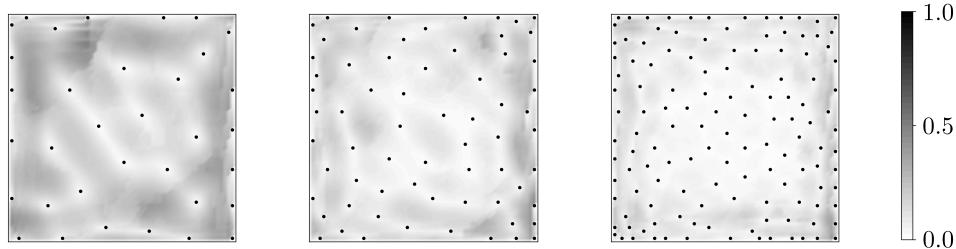


Fig. 6: Relative error,  $\|\Phi(\cdot, x) - \tilde{\Phi}(\cdot, x)\|/\|\Phi(\cdot, x)\|$ , in the approximation of the “column” of the integral kernel associated with  $x$ , using 5 (left), 10 (center) and 20 (right) impulse response batches. Sample points are indicated by black dots. The error associated with the point  $x$  is the shade of the image at location  $x$ , with white indicating zero error and black indicating 100% error. At the sample points, the error is zero. The further the point  $x$  is from the sample points, the larger the error. Adding more batches yields a more accurate approximation.

length of the resultant function near  $x$  [31, 74]. The moment method that we use estimates the support of  $\phi_x$  more accurately than random noise probing in resolution analysis, at the cost of the additional constraint that  $\mathcal{A}$  has a non-negative integral kernel. (3) The **PSF-based** method we propose never evaluates computed impulse responses outside of their domain of definition. This eliminates **boundary artifacts** which “boundary-artifact” errors (see [5, Section 1.1]) that plague conventional PC and **PSF** methods.

The ability to rapidly approximate entries of  $\mathcal{A}$ 's integral kernel allows one to approximate discretized versions of  $\mathcal{A}$  using the full arsenal of tools for matrix approximation that rely on fast access to matrix entries. In this work, we form a hierarchical matrix(**H-matrix**) [13, 42] approximation of a discretized version of  $\mathcal{A}$ . H-matrices are a **compressed** matrix format in which the rows and columns of the matrix are re-ordered, then the matrix is recursively subdivided into blocks in such a way that

many off-diagonal blocks are low rank, even though the matrix as a whole may be high rank. H-matrix methods permit us to perform matrix-vector products cheaply, and perform other useful linear algebra operations that cannot be done easily using the original operator. These operations include matrix-matrix addition, matrix-matrix multiplication, matrix factorization, and matrix inversion. The work and memory required to perform these operations for an  $N \times N$  H-matrix scales nearly linearly in  $N$  (i.e.,  $O(N^{1+\epsilon})$  for any  $\epsilon > 0$ ). The exact cost depends with rank  $k_h$  blocks scales as  $O(k_h^a N \log(N)^b)$  where  $a, b \in \{0, 1, 2, 3\}$  are constants which depend on the type of H-matrix used, and the operation being performed, and the rank of the off-diagonal blocks [40] [40][52, Section 2.1].

## 2. Why we need more efficient approximations of high rank Hessians.

While the PSF-based method proposed in this paper may be used to approximate any operator that has a locally supported non-negative integral kernel, we are primarily motivated by approximation of high-rank Hessians in distributed parameter inverse problems governed by PDEs. In this section, we provide a brief background on this topic, and explain why existing Hessian approximation methods are not satisfactory.

In distributed parameter inverse problems governed by PDEs, one seeks to infer an unknown spatially varying parameter field from limited observations of a state variable that depends on the parameter implicitly through the solution of a PDE. Conventionally, the inverse problem is formulated using either a deterministic framework [9, 76], or a Bayesian probabilistic framework [49, 70, 72]. In the deterministic framework, one solves an optimization problem to find the parameter that best fits the observations, subject to appropriate regularization [25, 76]. In the probabilistic framework, Bayes' theorem combines the observations with prior information to form a posterior distribution over the space of all possible parameter fields, and computations are performed to extract statistical information about the parameter from this posterior. The Hessian of the objective function with respect to the parameter in the deterministic optimization problem and the Hessian of the negative log posterior in the Bayesian setting are equal or approximately equal under typical noise, regularization, and prior models, so we refer to both of these Hessians as “the Hessian.” The Hessian consists of a data misfit term (the *data misfit Hessian*), which depends on a discrepancy between the observations and the associated model predictions, and a regularization or prior term (the *regularization Hessian*) which does not depend on the observations. For more details on the Hessian, see [4, 38, 75].

Hessian approximations and preconditioners are highly desirable because the Hessian is central to efficient solution of inverse problems in both the deterministic and Bayesian settings. When solving the deterministic optimization problem with Newton-type methods, the Hessian is the coefficient operator for the linear system that must be solved or approximately solved at every Newton iteration. Good Hessian preconditioners reduce the number of iterations required to solve these Newton linear systems with the conjugate gradient method [66]. In the Bayesian setting, the inverse of the Hessian is the covariance of a local Gaussian approximation of the posterior. This Gaussian distribution can be used directly as an approximation of the posterior, or it can be used as a proposal for Markov chain Monte-Carlo methods for drawing samples from the posterior. For instance, see [50, 62] and the references therein.

Owing to the implicit dependence of predicted observations on the parameter, entries of the Hessian are not easily accessible. Rather, the Hessian may be applied to a vector via a computational process that involves solving a pair of forward and adjoint

163 PDEs which are linearizations of the original PDE [38, 63]. The most popular matrix-  
 164 free Hessian approximation methods are based on low rank approximation of either the  
 165 data misfit Hessian, or the data misfit Hessian preconditioned by the regularization  
 166 Hessian, e.g., [15, 19, 33, 62, 68]. Krylov methods such as Lanczos or randomized  
 167 methods [18, 44] are typically used to construct these low rank approximations by  
 168 applying the Hessian to vectors. Using these methods, the required number of Hessian  
 169 applications (and hence the required number of PDE solves) is proportional to the  
 170 rank of the low rank approximation. Low rank approximation methods are justified  
 171 by arguing that the numerical rank of the data misfit Hessian is insensitive to the  
 172 dimension of the discretized parameter. This means that the required number of  
 173 PDE solves remains the same as the mesh used to discretize the parameter is refined.  
 174 However, in many inverse problems of practical interest the numerical rank of the  
 175 data misfit Hessian, while mesh independent, is still large, which makes it costly to  
 176 approximate the Hessian using low rank approximation methods [7, 15, 48].

177 Examples of inverse problems with high rank data misfit Hessians include large-  
 178 scale ice sheet inverse problems [45, 48], advection dominated advection-diffusion in-  
 179 verse problems [2][34, Chapter 5], high frequency wave propagation inverse prob-  
 180 lems [15], inverse problems governed by high Reynolds number flows, and more gen-  
 181 erally, all inverse problems in which the observations highly inform the parameter.  
 182 The eigenvalues of the data misfit Hessian characterize how informative the data are  
 183 about components of the parameter in the corresponding eigenvector directions, hence  
 184 more informative data leads to larger eigenvalues and a larger numerical rank [3][4,  
 185 Section 1.4 and Chapter 4]. Roughly speaking, the numerical rank of the data mis-  
 186 fit Hessian is the dimension of the subspace of parameter space that is informed by  
 187 the data. The numerical rank of the regularization preconditioned data misfit Hes-  
 188 sian may be reduced by increasing the strength of the regularization, but this throws  
 189 away useful information: components of the parameter that could be learned from  
 190 the observations would instead be reconstructed based on the regularization [6, Sec-  
 191 tion 4][76, Chapters 1 and 7]. Hence, low rank approximation methods suffer from a  
 192 predicament: if the data highly inform the parameter and the regularization is chosen  
 193 appropriately, then a large number of operator applications are required to form an  
 194 accurate approximation of the Hessian using low rank approximation methods. High  
 195 rank Hessian approximation methods are thus needed.

196 Recently there have been improvements in matrix-free H-matrix construction  
 197 methods in which an operator is applied to structured random vectors, and the re-  
 198 sponse of the operator to those random vectors is processed to construct an H-matrix  
 199 approximation [54, 55, 57, 58, 59]. These methods (which we do not use here) have  
 200 been used to approximate Hessians in PDE constrained inverse problems [7, 45]. Al-  
 201 though these methods are promising, the required number of operator applications is  
 202 still large (e.g., hundreds to thousands). For example, using the method in [55], the  
required number of operator applies to construct an  $H^1$  matrix with hierarchical rank  
 $r$  for problems in a 2D domain discretized with a regular grid is  $\#levels \cdot 64 \cdot (r + c)$ ,  
where  $\#levels$  is the depth of the hierarchical partitioning,  $r$  is the rank of the  
blocks (hierarchical rank), and  $c$  is an oversampling parameter (see [55, Section 2.4]  
). On a  $64 \times 64$  grid with depth 4, hierarchical rank 10, and oversampling parameter  
 $c = 5$ , this works out to  $4 \cdot 64 \cdot (10 + 5) = 3840$  operator applies. In Section 7.3, we  
see numerically that the randomized hierarchical off-diagonal low rank (HODLR)  
method in [58] requires hundreds to thousands of matrix-vector products to construct  
approximations of the integral kernel for a blur problem example with modest (e.g.,  
10%) relative error. Matrix-free H-matrix construction is currently an active area of

213 research, hence these costs may decrease as new algorithms are developed. In this  
 214 paper, we also form an H-matrix Hessian-approximation. However, to reduce the  
 215 required number of operator applications, we first form a PSF approximation of the  
 216 data misfit Hessian by exploiting locality and non-negative integral kernel properties,  
 217 then form the H-matrix using classical H-matrix techniques. By using techniques.  
 218 Using this two stage approach, we reduce the number of operator applications to a  
 219 few dozen at most. One limitation of our method is that not

220 Not all data misfit Hessians satisfy the local non-negative integral kernel properties  
 221 (for example, We note, in particular, that the wave inverse problem data misfit  
 222 Hessian has a substantial amount and Gauss-Newton Hessian have a substantial  
 223 proportion of negative entries in its integral kernel) their integral kernels. In this  
 224 case more specialized techniques have been developed using, eg., pseudodifferential  
 225 operator theory [21, 71], and sparsity in the wavelet domain [46]. However, many  
 226 data misfit Hessians of practical interest do satisfy these the local non-negative integral  
 227 kernel properties (either exactly or approximately), and the proposed method PSF-based  
 228 method we propose is targeted at approximating these Hessians.

229 **3. Preliminaries.** Let  $\Omega \subset \mathbb{R}^d$  be a bounded domain (typically  $d = 1, 2$ , or  $3$ ).  
 230 We seek to approximate integral operators  $\mathcal{A} : L^2(\Omega) \rightarrow L^2(\Omega)'$  of the form

$$231 \quad (3.1) \quad (\mathcal{A}u)(w) := \int_{\Omega} \int_{\Omega} w(y) \Phi(y, x) u(x) dx dy.$$

232 The linear functional  $\mathcal{A}u \in L^2(\Omega)'$  is the result of applying  $\mathcal{A}$  to  $u \in L^2(\Omega)$ , and the  
 233 scalar  $(\mathcal{A}u)(w)$  is the result of applying that linear functional to  $w \in L^2(\Omega)$ . The  
 234 integral kernel,  $\Phi : \Omega \times \Omega \rightarrow \mathbb{R}$ , exists but is not easily accessible. In this section we  
 235 describe how to extend the domain of  $\mathcal{A}$  to distributions, which allows us to define  
 236 impulse responses (Section 3.1), we then state the conditions on  $\mathcal{A}$  that the proposed  
 237 PSF-based method requires (Section 3.2), and detail finite element discretization (Sec-  
 238 tion 3.3).

239 **3.1. Distributions and impulse responses.** The operator  $\mathcal{A}$  may be applied  
 240 to distributions<sup>3</sup> if  $\Phi$  is sufficiently regular. Given  $\rho \in L^2(\Omega)'$ , let  $\rho^* \in L^2(\Omega)$  denote  
 241 the Riesz representative of  $\rho$  with respect to the  $L^2(\Omega)$  inner product. We have

$$242 \quad (3.2a) \quad (\mathcal{A}\rho^*)(w) = \int_{\Omega} \int_{\Omega} w(y) \Phi(y, x) \rho^*(x) dx dy$$

$$243 \quad (3.2b) \quad = \int_{\Omega} w(y) \int_{\Omega} \Phi(y, x) \rho^*(x) dx dy = \int_{\Omega} w(y) \rho(\Phi(y, \cdot)) dy,$$

245 where  $\Phi(y, \cdot)$  denotes the function  $x \mapsto \Phi(y, x)$ . Now let  $\mathcal{D}(\Omega) \subset L^2(\Omega)$  be a suitable  
 246 space of test functions and let  $\rho : \mathcal{D}(\Omega) \rightarrow \mathbb{R}$  be a distribution. In this case,  $\rho^*$  may  
 247 not exist, so the derivation in (3.2) is not valid. However, if  $\Phi$  is sufficiently regular  
 248 such that the function  $y \mapsto \rho(\Phi(y, \cdot))$  is well-defined for almost all  $y \in \Omega$ , and if  
 249 this function is in  $L^2(\Omega)$ , then the right hand side of (3.2b) is well-defined. Hence, we  
 250 *define* the application of  $\mathcal{A}$  to the distribution  $\rho$  to be the right hand side of (3.2b).  
 251 We denote this operator application by “ $\mathcal{A}\rho^*$ ,” even if  $\rho^*$  does not exist.

252 Let  $\delta_x$  denote the delta distribution<sup>4</sup> (i.e., point source, impulse) centered at the  
 253 point  $x \in \Omega$ . The *impulse response* of  $\mathcal{A}$  associated with  $x$  is the function  $\phi_x : \Omega \rightarrow \mathbb{R}$ ,

$$254 \quad (3.3) \quad \phi_x := (\mathcal{A}\delta_x^*)^*,$$

<sup>3</sup>I.e., generalized functions such as the Dirac delta distribution. See, for example, [8, Chapter 5].

<sup>4</sup>Recall that the delta distribution  $\delta_x : \mathcal{D}(\Omega) \rightarrow \mathbb{R}$  is defined by  $\delta_x(w) = w(x)$  for all  $w \in \mathcal{D}(\Omega)$ .

255 that is formed by applying  $\mathcal{A}$  to  $\delta_x$  (per the generalized notion of operator “application”<sup>5</sup>  
 256 defined above), then taking the Riesz representation of the resulting linear functional.  
 257 Using (3.2b) and the definition of the delta distribution, we see that  $\phi_x$  may also be  
 258 written as the function  $\phi_x(y) = \Phi(y, x)$ .

259 **3.2. Required conditions.** We focus on approximating operators that satisfy  
 260 the following conditions:

- 261 1. The kernel  $\Phi$  is sufficiently regular so that  $\phi_x$  is well-defined for all  $x \in \Omega$ .  
 262 2. The supports of the impulse responses  $\phi_x$  are contained in localized regions.  
 263 3. The integral kernel is non-negative<sup>5</sup> in the sense that

$$\Phi(y, x) \geq 0 \quad \text{for all } (y, x) \in \Omega \times \Omega.$$

263 The proposed PSF-based method may still perform well if these conditions are relaxed  
 264 slightly.

265 It is acceptable if the support of  $\phi_x$  is not perfectly contained in a localized region  
 266 (violating Assumption 2), so long as the bulk of the “mass” of  $\phi_x$  is contained in a  
 267 localized region. The integral kernel does not need to be non-negative for all pairs  
 268 of points  $(y, x) \in \Omega \times \Omega$  so long as it is non-negative for the vast majority of pairs  
 269 of points  $(y, x)$ , and so long as In principle, the PSF-based method can be applied  
 270 even if the impulse responses are widely dispersed. However, in this case only a small  
 271 number of impulse responses can be computed per batch, which means more batches,  
 272 and hence more operator applies, are needed to form an accurate approximation.

273 If there are negative numbers in the integral kernel (violating Assumption 3), the  
 274 ellipsoid estimation procedure may incur errors or fail, leading to poor performance  
 275 or failure of the PSF-based method. In Figure 7 we investigate the robustness of the  
 276 ellipsoid support estimation procedure to violations of Assumption 3. We study two  
 277 integral kernel examples, both of which are parameterized by a quantity that controls  
 278 how negative the kernels are. We make the following observations:

- 279 • The larger and more numerous the negative numbers are comparatively small.  
   If these conditions are violated, the method will incur additional error, and  
   depending on the severity of , the more inaccurate the ellipsoid support  
   estimate is.
- 283 • The further away from the center of the ellipsoid the negative numbers are,  
   the violation, the proposed method may fail more influence they have on the  
   quality of the ellipsoid support estimate. This is because moment formulas  
   (Equations 4.2 and 4.3) assign more weight to entries in the kernel that are  
   further from the center.
- 288 • Negative numbers affect the ellipsoid estimation method more if they are  
   isolated, and less if they are balanced by nearby positive numbers.
- 290 • As kernels become more negative, the ellipsoid estimation performs well up  
   to a certain threshold that depends on the spatial distribution of negative  
   and positive entries. After that threshold is crossed, the estimation rapidly  
   transitions to performing poorly and ultimately failing.

294 For the kernel in the left two columns of Figure 7, negative numbers are interspersed  
 295 with positive numbers, allowing us to include a large amount of negative numbers  
 296 before the ellipsoid estimation fails. For the kernel in the right two columns, the

---

5 Note that having a non-negative integral kernel is different from positive semi-definiteness. The operator  $\mathcal{A}$  need not be positive semi-definite to use the proposed PSF-based method, and positive semi-definite operators need not have a non-negative integral kernel.

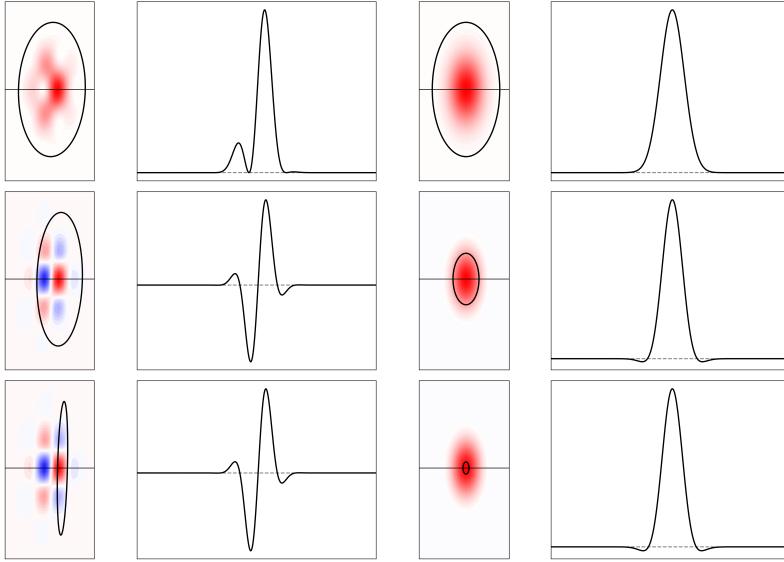


Fig. 7: Illustration of the influence of negative numbers in the integral kernel on the robustness of the ellipsoid estimates for the supports of impulse responses. Left two columns: Blur kernel given in Equation 7.4. Right two columns: Ricker wavelet-type kernel given by  $\Phi(y, x) = (1 - a\gamma) \exp(-\gamma/2)$ , where  $\gamma = (y - x)^T \Sigma^{-1} (y - x)$ , and  $\Sigma = \text{diag}(0.0025, 0.01)$ . Ordered from top to bottom, the results are obtained with  $a \in \{1.0, 20.0, 27.0\}$  for the left two columns, and  $a \in \{0.0, 0.23, 0.249\}$  for the right two columns. Columns 1 and 3: impulse responses with estimated support ellipsoids indicated by the black ellipses. Red and blue represent positive and negative numbers in the integral kernel, respectively. Columns 2 and 4: one-dimensional slice along the horizontal line indicated in the two-dimensional plots. The dashed gray line is at zero.

297 ellipsoid estimate fails with tiny amounts of negative numbers because the negative  
 298 numbers are far away from the mean and not balanced by positive numbers at similar  
 299 distances and angles. In the bottom two rows, we see the aforementioned threshold  
 300 effect, in which the ellipsoid estimation method rapidly transitions from performing  
 301 reasonably well to performing poorly with only a small change to the integral kernel.

302 **3.3. Finite element discretization.** In computations, functions are discretized  
 303 and replaced by finite-dimensional vectors, and operators mapping between infinite-  
 304 dimensional spaces are replaced by operators mapping between finite-dimensional spa-  
 305 ces. In this paper we discretize using continuous Galerkin—the functions that  $\mathcal{A}$  and  
 306  $\mathcal{A}^T$  are applied to using continuous finite elements satisfying the Kronecker property  
 307 (defined below). With minor modifications, the proposed PSF-based method could be  
 308 used with more general finite element methods, or other discretization schemes such  
 309 as finite differences or finite volumes. These restrictions on discretization only apply  
 310 to functions  $u$  that  $\mathcal{A}$  and  $\mathcal{A}^T$  are applied to. Other functions that arise internally  
 311 during the process of computing actions of  $\mathcal{A}$ , such as state variables in a PDE that  
 312 is solved in a subproblem, may be discretized with any method.

313 Let  $\psi_1, \psi_2, \dots, \psi_N$  be a set of continuous Galerkin finite element basis functions

314 used to discretize the problem on a mesh with mesh size parameter  $h$ , let  $V_h :=$   
 315  $\text{span}(\psi_1, \psi_2, \dots, \psi_N)$  be the corresponding finite element space under the  $L^2$  inner  
 316 product, and let  $p_i \in \mathbb{R}^d$ ,  $i = 1, \dots, N$  be the Lagrange nodes associated with the  
 317 functions  $\psi_i$ . We assume that the finite element basis satisfies the Kronecker property,  
 318 i.e.,  $\psi_i(p_i) = 1$  and  $\psi_i(p_j) = 0$  if  $i \neq j$ . For  $u_h \in V_h$  we write  $\mathbf{u} \in \mathbb{R}_{\mathbf{M}}^m$  to denote  
 319 the coefficient vector for  $u_h$  with respect to the finite element basis, i.e.,  $u_h(x) =$   
 320  $\sum_{i=1}^N \mathbf{u}_i \psi_i(x)$ . Linear functionals  $\rho_h \in V'_h$  have coefficient dual vectors  $\boldsymbol{\rho} \in \mathbb{R}_{\mathbf{M}^{-1}}^m$ ,  
 321 with entries  $\boldsymbol{\rho}_i = \rho_h(\psi_i)$  for  $i = 1, \dots, m$ . Here,  $\mathbf{M} \in \mathbb{R}^{N \times N}$  denotes the sparse finite  
 322 element mass matrix which has entries  $\mathbf{M}_{ij} = \int_{\Omega} \psi_i(x) \psi_j(x) dx$  for  $i, j = 1, \dots, N$ .  
 323 The space  $\mathbb{R}_{\mathbf{M}}^N$  is  $\mathbb{R}^N$  with the inner product  $(\mathbf{u}, \mathbf{w})_{\mathbf{M}} := \mathbf{u}^T \mathbf{M} \mathbf{w}$ , and  $\mathbb{R}_{\mathbf{M}^{-1}}^N$  is the  
 324 analogous space with  $\mathbf{M}^{-1}$  replacing  $\mathbf{M}$ . Direct calculation shows that  $\mathbb{R}_{\mathbf{M}}^N$  and  $\mathbb{R}_{\mathbf{M}^{-1}}^N$   
 325 are isomorphic to  $V_h$  and  $V'_h$  as Hilbert spaces, respectively.

326 After discretization, the operator  $\mathcal{A} : L^2(\Omega) \rightarrow L^2(\Omega)'$  is replaced by an operator  
 327  $A_h : V_h \rightarrow V'_h$ , which becomes an operator  $\mathbf{A} : \mathbb{R}_{\mathbf{M}}^N \rightarrow \mathbb{R}_{\mathbf{M}^{-1}}^N$  under the isomorphism  
 328 discussed above. The proposed-PSF-based method is agnostic to the computational  
 329 procedure for approximating  $\mathcal{A}$  with  $\mathbf{A}$ . What is important is that we do not have  
 330 direct access to matrix entries  $\mathbf{A}_{ij}$ . Rather, we have a computational procedure  
 331 that allows us to compute matrix-vector products  $\mathbf{u} \mapsto \mathbf{A}\mathbf{u}$  and  $\mathbf{w} \mapsto \mathbf{A}^T\mathbf{w}$ , and  
 332 computing these matrix-vector products is costly. The proposed-PSF-based method  
 333 mitigates this cost by performing as few matrix-vector products as possible. Of course,  
 334 matrix entries can be computed via matrix-vector products as  $\mathbf{A}_{ij} = (\mathbf{A}\mathbf{e}_j)_i$ , where  
 335  $\mathbf{e}_j = (0, \dots, 0, 1, 0, \dots, 0)^T$  is the length  $N$  unit vector with one in the  $j$ th coordinate  
 336 and zeros elsewhere. But computing the matrix-vector product  $\mathbf{e}_j \mapsto \mathbf{A}\mathbf{e}_j$  is costly,  
 337 and therefore wasteful if we do not use other matrix entries in the  $j$ th column of  $\mathbf{A}$ .  
 338 Hence, methods for approximating  $\mathbf{A}$  are computationally intractable if they require  
 339 accessing scattered matrix entries from many different rows and columns of  $\mathbf{A}$ .

340 The operator  $A_h : V_h \rightarrow V'_h$  can be written in integral kernel form, (3.1), but  
 341 with  $\Phi$  replaced by a slightly different integral kernel,  $\Phi_h$ , which we do not know,  
 342 and which differs from  $\Phi$  due to discretization error. Since the functions in  $V_h$  are  
 343 continuous at  $x$ , the delta distribution  $\delta_x$  is a continuous linear functional on  $V_h$ , which  
 344 has a discrete dual vector  $\boldsymbol{\delta}_x \in \mathbb{R}_{\mathbf{M}^{-1}}^N$  with entries  $(\boldsymbol{\delta}_x)_i = \psi_i(x)$  for  $i = 1, \dots, N$ .  
 345 Additionally, it is straightforward to verify that the Riesz representation,  $\rho_h^* \in V_h$ ,  
 346 of a functional  $\rho \in V'_h$  has coefficient vector  $\boldsymbol{\rho}^* = \mathbf{M}^{-1}\boldsymbol{\rho}$ . Therefore, the formula for  
 347 the impulse response from (3.3) becomes  $\phi_x = (A_h \delta_x^*)^* = \mathbf{M}^{-1} \mathbf{A} \mathbf{M}^{-1} \boldsymbol{\delta}_x$ , and the  
 348  $(y, x)$  kernel entry of  $\Phi_h$  may be written as  $\Phi_h(y, x) = \boldsymbol{\delta}_y^T \phi_x = \boldsymbol{\delta}_y^T \mathbf{M}^{-1} \mathbf{A} \mathbf{M}^{-1} \boldsymbol{\delta}_x$ .  
 349 Now define  $\boldsymbol{\Phi} \in \mathbb{R}^{N \times N}$  to be the following dense matrix of kernel entries evaluated  
 350 at all pairs of Lagrange nodes:

351 (3.4) 
$$\boldsymbol{\Phi}_{ij} := \Phi_h(p_i, p_j).$$

352 Because of the Kronecker property of the finite element basis, we have  $\boldsymbol{\delta}_{p_i} = \mathbf{e}_i$ . Thus,  
 353 we have  $\Phi_h(p_i, p_j) = (\mathbf{M}^{-1} \mathbf{A} \mathbf{M}^{-1})_{ij}$ , which implies

354 (3.5) 
$$\mathbf{A} = \mathbf{M} \boldsymbol{\Phi} \mathbf{M}.$$

355 Broadly, we will construct an H-matrix approximation of  $\mathbf{A}$  by forming a-an H-matrix  
 356 approximation of  $\boldsymbol{\Phi}$ , then multiplying  $\boldsymbol{\Phi}$  by  $\mathbf{M}$  (or a lumped mass version of  $\mathbf{M}$ ) on  
 357 the left and right using H-matrix methods. Classical H-matrix construction methods  
 358 require access to arbitrary matrix entries  $\boldsymbol{\Phi}_{ij}$ , but these matrix entries are not easily  
 359 accessible. The bulk of the proposed-PSF-based method is therefore dedicated to  
 360 forming approximations of these matrix entries that can be evaluated rapidly.

361     *Lumped mass matrix.* At the continuum level,  $\Phi$  is assumed to be non-negative.  
 362     However, entries of  $\Phi$  involve inverse mass matrices, which typically contain negative  
 363     numbers. We therefore recommend replacing the mass matrix,  $\mathbf{M}$ , with a positive  
 364     diagonal *lumped mass* approximation. Here, we use the lumped mass matrix in which  
 365     the  $i$ th diagonal entry of the lumped mass matrix is the sum of all entries in the  $i$ th  
 366     row of the mass matrix. Other mass lumping techniques may be used.

367     **4. Key innovations.** In this section we present two key innovations that the  
 368     proposed-PSF-based method is based on. First, we define moments of the impulse  
 369     responses,  $\phi_x$ , show how these moments can be computed efficiently, and use these  
 370     moments to form ellipsoid shaped a-priori estimates for the supports of the impulse  
 371     responses (Section 4.1). Second, we describe an improved method to approximate  
 372     impulse responses from other nearby impulse responses, which we call “normalized  
 373     local mean displacement invariance” (Section 4.2).

374     **4.1. Impulse response moments and ellipsoid support estimate.** The  
 375     impulse response  $\phi_x$  may be interpreted as a scaled probability distribution because  
 376     of the non-negative integral kernel property. Let  $V : \Omega \rightarrow \mathbb{R}$ ,

$$377 \quad (4.1) \quad V(x) := \int_{\Omega} \phi_x(y) dy,$$

378     denote the spatially varying scaling factor, and for  $i, j = 1, \dots, d$  define  $\mu : \Omega \rightarrow \mathbb{R}^d$   
 379     and  $\Sigma : \Omega \rightarrow \mathbb{R}^{d \times d}$  as follows:

$$380 \quad (4.2) \quad \mu^i(x) := \frac{1}{V(x)} \int_{\Omega} (\phi_x(y) / V(x)) y^i dy$$

$$381 \quad (4.3) \quad \Sigma^{ij}(x) := \frac{1}{V(x)} \int_{\Omega} (\phi_x(y) / V(x)) (y^i - \mu^i(x)) (y^j - \mu^j(x)) dy,$$

383     where  $\mu^i(x)$  denotes and  $y^i$  denote the  $i$ th component of the vector components of  
 384     the vectors  $\mu(x)$  and  $y$ , respectively, and  $\Sigma^{ij}(x)$  denotes the  $(i, j)$  entry of the matrix  
 385      $\Sigma(x)$ . The vector quantities  $\mu(x) \in \mathbb{R}^d$  and the matrix  $\Sigma(x) \in \mathbb{R}^{d \times d}$  are the mean  
 386     and covariance of the normalized version of  $\phi_x$ , respectively.

387     The direct approach to compute  $V(x)$ ,  $\mu(x)$ , and  $\Sigma(x)$  is to apply  $\mathcal{A}$  to a point  
 388     source centered at  $x$  to obtain  $\phi_x$ , per (3.3). Then one can post process  $\phi_x$  to de-  
 389     termine  $V(x)$ ,  $\mu(x)$ , and  $\Sigma(x)$ . However, this direct approach is not feasible because  
 390     we need our algorithm for picking sample points (Section 5.1 and Figure 4) needs  
 391     to know  $V(x)$ ,  $\mu(x)$ , and  $\Sigma(x)$  before we compute  $\phi_x$  in order choose the point  $x$ .  
 392     Computing  $\phi_x$  in order to determine  $V(x)$ ,  $\mu(x)$ , and  $\Sigma(x)$  would be extremely com-  
 393     putationally expensive, and defeat the purpose of the proposed-PSF-based method,  
 394     which is to reduce the computational cost by computing impulse responses in batches.  
 395     Fortunately, it is possible to compute  $V(x)$ ,  $\mu(x)$ , and  $\Sigma(x)$  indirectly, for all points  
 396      $x \in \Omega$  simultaneously, by applying  $\mathcal{A}^T$  to one constant function,  $d$  linear functions,  
 397     and  $d(d+1)/2$  quadratic functions (e.g., 6 total operator applications in two spatial  
 398     dimensions and 10 in three spatial dimensions). This may be motivated by analogy  
 399     to matrices. If  $\mathbf{A} \in \mathbb{R}^{N \times N}$  is a matrix with  $i$ th column  $\mathbf{a}_i$  and  $\mathbf{w} \in \mathbb{R}^N$ , then

$$400 \quad \mathbf{A}^T \mathbf{w} = \begin{bmatrix} \mathbf{a}_1^T & \mathbf{a}_2^T & \cdots \\ \vdots & \vdots & \ddots \\ \mathbf{a}_N^T & \mathbf{a}_N^T & \cdots \end{bmatrix} \mathbf{w} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{w} \\ \vdots \\ \mathbf{a}_N^T \mathbf{w} \end{bmatrix}.$$

401 By computing one matrix-vector product of  $\mathbf{A}^T$  with  $\mathbf{w}$ , we compute the inner product  
 402 of each column of  $\mathbf{A}$  with  $\mathbf{w}$  simultaneously. The operator case is analogous, with  $\phi_x$   
 403 taking the place of a matrix column. We have

404 (4.4) 
$$(\mathcal{A}^T w)^*(x) = \int_{\Omega} \Phi(y, x) w(y) dy = (\phi_x, w)_{L^2(\Omega)}.$$

405 By computing one operator application of  $\mathcal{A}^T$  to  $w$ , we compute the inner product of  
 406 each  $\phi_x$  with  $w$ , for all points  $x$  simultaneously.

407 Let  $C$ ,  $L^i$ , and  $Q^{ij}$  be the following constant, linear, and quadratic functions:

408 
$$C(x) := 1, \quad L^i(x) := x^i, \quad Q^{ij}(x) := x^i x^j$$

409 for  $i, j = 1, \dots, d$ . Using the definition of  $V$  in (4.1) and using (4.4), we have

410 
$$V(x) = \int_{\Omega} \phi_x(y) C(y) dy = (\phi_x, C)_{L^2(\Omega)} = (\mathcal{A}^T C)^*(x).$$

411 Hence, we compute  $V(x)$  for all  $x$  simultaneously by applying  $\mathcal{A}^T$  to  $C$ . Analogous  
 412 manipulations show that  $\mu(x)$  and  $\Sigma(x)$  may be computed for all points  $x$  simultane-  
 413 ously by applying  $\mathcal{A}^T$  to the functions  $L^i$  and  $Q^{ij}$ , respectively. We have

414 (4.5a) 
$$V = (\mathcal{A}^T C)^*$$

415 (4.5b) 
$$\mu^i = (\mathcal{A}^T L^i)^*/V$$

416 (4.5c) 
$$\Sigma^{ij} = (\mathcal{A}^T Q^{ij})^*/V - \mu^i \cdot \mu^j$$

418 for  $i, j = 1, \dots, d$ . Here  $u/w$  denotes pointwise division,  $(u/w)(x) = u(x)/w(x)$ , and  
 419  $u \cdot w$  denotes pointwise multiplication,  $(u \cdot w)(x) = u(x)w(x)$ .

420 We approximate the support of  $\phi_x$  with the ellipsoid

421 (4.6) 
$$E_x := \{x' \in \Omega : (x' - \mu(x))^T \Sigma(x)^{-1} (x' - \mu(x)) \leq \tau^2\},$$

422 where  $\tau$  is a fixed constant (see Figure 3). The ellipsoid  $E_x$  is the set of points  
 423 within  $\tau$  standard deviations of the mean of the Gaussian distribution with mean  
 424  $\mu(x)$  and covariance  $\Sigma(x)$ , i.e., the Gaussian distribution which has the same mean  
 425 and covariance as the normalized version of  $\phi_x$ . The quantity  $\tau$  is a parameter that  
 426 must be chosen appropriately. The larger  $\tau$  is, the larger the ellipsoid  $E_x$  is, and  
 427 the more conservative the estimate is for the support of  $\phi_x$ . However, in Section 5.1  
 428 we will see that the cost of the proposed PSF-based method depends on how many  
 429 non-overlapping ellipsoids  $E_x$  we can “pack” in the domain  $\Omega$  (more ellipsoids is  
 430 better), and choosing a larger value of  $\tau$  means that fewer ellipsoids will fit in  $\Omega$ . In  
 431 practice, we find that  $\tau = 3.0$  yields a reasonable balance between these competing  
 432 interests, and use  $\tau = 3.0$  in the numerical results all numerical results, except for  
 433 Figure 14, where we study the effects of varying  $\tau$ . The fraction of the “mass” of  $\phi_x$   
 434 residing outside of  $E_x$  is less than  $1/\tau^2$  by Chebyshev’s inequality, though this bound  
 435 is typically conservative.

436 **4.2. Local mean displacement invariance.** Let  $x$  and  $x'$  be points in  $\Omega$  that  
 437 are close to each other, and consider the following approximations:

438 (4.7) 
$$\phi_x(y) \approx \phi_{x'}(y)$$

439 (4.8) 
$$\phi_x(y) \approx \phi_{x'}(y - x + x')$$

440 (4.9) 
$$\phi_x(y) \approx \phi_{x'}(y - \mu(x) + \mu(x'))$$

441 (4.10) 
$$\phi_x(y) \approx \phi_{x'}(y - \mu(x) + \mu(x')) V(x)/V(x').$$

443 These are four different ways to approximate an impulse response by a nearby impulse response, with each successive approximation building upon the previous ones.  
 444 The ~~proposed~~-PSF-based method uses (4.10), which is the most sophisticated. Approximation (4.7) says that  $\phi_x$  can be approximated by  $\phi_{x'}$  when  $x$  and  $x'$  are close.  
 445 Operators satisfying (4.7) can be well approximated via low rank CUR approximation.  
 446 However, the required rank in the low rank approximation can be large, which makes  
 447 algorithms based on (4.7) expensive. Operators that satisfy (4.8) are called “locally  
 448 translation invariant” because integral kernel entries  $\Phi(y, x)$  for such operators are  
 449 approximately invariant under translation of  $x$  and  $y$  by the same displacement, i.e.,  
 450  $x \rightarrow x + h$  and  $y \rightarrow y + h$ . It is straightforward to show that if equality holds in (4.8),  
 451 then  $\mathcal{A}$  is a convolution operator. Locally translation invariant operators act like  
 452 convolutions locally, and can therefore be well approximated by PC approximations.  
 453

454 Approximation (4.9) improves upon (4.7) and (4.8), and generalizes both. On one  
 455 hand, if (4.7) holds, then  $\mu(x) \approx \mu(x')$ , and so (4.9) holds. On the other hand, translating a distribution translates its mean, so if (4.8) holds, then  $\mu(x') - \mu(x) \approx x' - x$ , so  
 456 again (4.9) holds. But approximation (4.9) can hold in situations where neither (4.7)  
 457 nor (4.8) holds. For example, because the expected value commutes with affine trans-  
 458 formations, (4.9) will hold when  $\mathcal{A}$  is locally translation invariant with respect to a  
 459 translated and rotated frame of reference, while (4.8) will not. Additionally, (4.9)  
 460 generalizes to operators  $\mathcal{A} : L^2(\Omega_1) \rightarrow L^2(\Omega_2)'$  that map between function spaces on  
 461 different domains  $\Omega_1$  and  $\Omega_2$ , and even operators that map between domains with  
 462 different spatial dimensions. In contrast, (4.8) does not naturally generalize to op-  
 463 erators that map between function spaces on different domains, because the formula  
 464  $y - x + x'$  requires vectors in  $\Omega_2$  and  $\Omega_1$  to be added together. We call (4.9) “local  
 465 mean displacement invariance,” and illustrate (4.9) in Figure 1b.

466 We use approximation (4.10), which is the same as (4.9), except for the factor  
 467  $V(x)/V(x')$ . This factor makes the approximation more accurate if  $V(x)$  varies widely.  
 468 Approximation (4.10) is equivalent to (4.9), but with  $\phi_x$  replaced by its normalized  
 469 version,  $\phi_x/V(x)$ . We call (4.10) *normalized local mean displacement invariance*.

## 472 5. Operator approximation algorithm. Before presenting the technical details 473 of the algorithm in Sections 5.1–5.5, we first provide an overview.

474 We use (4.5) to compute  $V$ ,  $\mu$ , and  $\Sigma$  by applying  $\mathcal{A}^T$  to polynomial functions.  
 475 Then we use (4.6) to form ellipsoid shaped estimates for the support of the  $\phi_x$ 's,  
 476 without computing them (see Figure 3). This allows us to compute large numbers of  
 477  $\phi_{x_i}$  in “batches,”  $\eta_b$  (see Figure 1a<sup>1</sup> and Figure 4). We compute one batch, denoted  
 478  $\eta_b$ , by applying  $\mathcal{A}$  to a weighted sum of point sources (Dirac comb) associated with  
 479 a batch,  $S_b$ , of points  $x_i$  scattered throughout  $\Omega$  (Section 5.2). The batch of points,  
 480  $S_b$ , is chosen via a greedy ellipsoid packing algorithm so that, for  $x_i, x_j \in S_b$ , the  
 481 support ellipsoid for  $\phi_{x_i}$  and the support ellipsoid for  $\phi_{x_j}$  do not overlap if  $i \neq j$   
 482 (Section 5.1). Because these supports do not overlap (or do not overlap much), we  
 483 can post process  $\eta_b$  to recover the functions  $\phi_{x_i}$  associated with all points  $x_i \in S_b$ .  
 484 With one application of  $\mathcal{A}$ , we recover many  $\phi_{x_i}$  (Section 5.2). The process is repeated  
 485 until a desired number of batches is reached.

486 Once the batches  $\eta_b$  are computed, we approximate the integral kernel  $\Phi(y, x)$   
 487 at arbitrary points  $(y, x)$  by interpolation of translated and scaled versions of the  
 488 computed  $\phi_{x_i}$  (Section 5.3 and Figure 5). The key idea behind the interpolation is the  
 489 normalized local mean displacement invariance assumption discussed in Section 4.2.  
 490 Specifically, we approximate  $\Phi(y, x) = \phi_x(y)$  by a weighted linear combination of the  
 491 values  $\frac{V(x)}{V(x_i)}\phi_{x_i}(y - \mu(x) + \mu(x_i))$  for a small number of sample points  $x_i$  near  $x$ . The

492 weights are determined by radial basis function (RBF) interpolation.

493 The ability to rapidly evaluate approximate kernel entries  $\Phi(y, x)$  allows us to  
 494 construct an H-matrix approximation,  $\Phi_H \approx \Phi$ , using the conventional adaptive  
 495 cross H-matrix construction method (Section 5.4). In this method, one forms low rank  
 496 approximations of off-diagonal blocks of the matrix by sampling rows and columns of  
 497 those blocks. We then convert  $\Phi_H$  into an H-matrix approximation  $\mathbf{A}_H \approx \mathbf{A}$ .

498 When  $\mathcal{A}$  is symmetric positive semi-definite,  $\mathbf{A}_H$  may be non-symmetric and  
 499 indefinite due to errors in the approximation. In this case, one may optionally sym-  
 500 metrize  $\mathbf{A}_H$ , then modify it via low rank updates to remove erroneous negative ei-  
 501 genvalues (Section 5.5). The complete algorithm for constructing  $\mathbf{A}_H$  is shown in  
 502 Algorithm 1. The computational cost is discussed in Section 6.

---

**Algorithm 1:** Construct PSF H-matrix approximation

---

**Input :** Linear operator  $\mathcal{A}$ , parameter  $n_b$

**Output:** H-matrix  $\mathbf{A}_H$

- 1 Compute  $V, \mu$ , and  $\Sigma$  (Equations (4.5) in Section 4.1)
  - 2 **for**  $k = 1, 2, \dots, n_b$  **do**
  - 3     Choose a batch of sample points,  $S_k$  (Section 5.1)
  - 4     Compute impulse response batch  $\eta_k$  by applying  $\mathcal{A}$  to the Dirac comb for  
 $S_k$  (Section 5.2)
  - 5     Form H-matrix approximation  $\Phi_H$  of integral kernel (Sections 5.3 and 5.4)
  - 6     Form H-matrix approximation  $\mathbf{A}_H$  of  $\mathcal{A}$  (Section 5.4)
  - 7     (optional) Modify  $\mathbf{A}_H$  to make it symmetric and remove negative  
eigenvalues (Section 5.5)
- 

503 **5.1. Sample point selection via greedy ellipsoid packing.** We choose sam-  
 504 ple points,  $x_i$ , in batches  $S_k$ . We use a greedy ellipsoid packing algorithm to choose  
 505 as many points as possible per batch, while ensuring that there is no overlap between  
 506 the support ellipsoids,  $E_{x_i}$ , associated with the sample points within a batch.

507 We start with a finite set of candidate points  $X$  and build  $S_k$  incrementally with  
 508 points selected from  $X$ . For simplicity of explanation, here  $S_k$  and  $X$  are mutable  
 509 sets that we add points to and remove points from. First we initialize  $S_k$  as an empty  
 510 set. Then we select the candidate point  $x_i \in X$  that is the farthest away from all  
 511 points in previous sample point batches  $S_1 \cup S_2 \cup \dots \cup S_{k-1}$ . Candidate points for the  
 512 first batch  $S_1$  are chosen randomly from  $X$ . Once  $x_i$  is selected, we remove  $x_i$  from  
 513  $X$ . Then we perform the following checks:

- 514 1. We check whether  $x_i$  is sufficiently far from all of the previously chosen points  
in the current batch, in the sense that  $E_{x_i} \cap E_{x_j} = \emptyset$  for all  $x_j \in S_k$ .
- 515 2. We make sure that  $V(x_i)$  is not too small, by checking whether  $V(x_i) >$   
 $\epsilon_V V_{\max}$ . Here,  $V_{\max}$  is the largest value of  $V(x_j)$  over all points  $q$  in the  
initial set of candidate points, and  $\epsilon_V$  is a small threshold parameter (we use  
 $\epsilon_V = 10^{-5}$ ).
- 516 3. We make sure that all eigenvalues of  $\Sigma(x_i)$  are positive, and the aspect ratio of  
 $E_{x_i}$  (square root of the ratio of the largest eigenvalue of  $\Sigma(x_i)$  to the smallest)  
is bounded by a constant  $1/\epsilon_\Sigma$  (we use  $1/\epsilon_\Sigma = 20$ ). Negative integral kernel  
entries due to discretization error can cause  $\Sigma(x_i)$  to be indefinite or highly  
ill-conditioned.
- 517 4. If  $x_i$  passes these checks (i.e., if  $x_i$  is sufficiently far from other points in the batch, and  
 $V(x_i)$  and  $\Sigma(x_i)$  are acceptable) then we add  $x_i$  to  $S_k$ . Otherwise we discard  $x_i$ . This

527 process repeats until there are no more points in  $X$ . We repeat the point selection  
 528 process to construct several batches of points  $S_1, S_2, \dots, S_{n_b}$ . For each batch,  $X$  is  
 529 initialized as the set of all Lagrange nodes for the finite element basis functions used  
 530 to discretize the problem, except for points in previous batches.

531 We check whether  $E_{x_i} \cap E_{x_j} = \{\}$  in a two stage process. First, we check whether  
 532 the axis aligned bounding boxes for the ellipsoids intersect. This quickly rules out  
 533 intersections of ellipsoids that are far apart. Second, if the bounding boxes intersect,  
 534 we check if the ellipsoids intersect using the ellipsoid intersection test in [39].

535 **5.2. Impulse response batches.** We compute impulse responses,  $\phi_{x_i}$ , in batches  $\blacksquare$   
 536 by applying  $\mathcal{A}$  to Dirac combs. The Dirac comb,  $\xi_k$ , associated with a batch of sam-  
 537 ple points,  $S_k$ , is the following weighted sum of Dirac distributions (point sources)  
 538 centered at the points  $x_i \in S_k$ :

$$539 \quad \xi_k := \sum_{x_i \in S_k} \delta_{x_i} / V(x_i).$$

540 We compute the *impulse response batch*,  $\eta_k$ , by applying  $\mathcal{A}$  to the Dirac comb:

$$541 \quad (5.1) \quad \eta_k := (\mathcal{A}\xi_k^*)^* = \sum_{x_i \in S_k} \phi_{x_i} / V(x_i).$$

542 The last equality in (5.1) follows from linearity and the definition of  $\phi_{x_i}$  in (3.3). Since  
 543 the points  $x_i$  are chosen so that the ellipsoid  $E_{x_i}$  that (approximately) supports  $\phi_i$ ,  
 544 and the ellipsoid  $E_{x_j}$  that (approximately) supports  $\phi_j$  do not overlap when  $i \neq j$ ,  
 545 we have (approximately)

$$546 \quad (5.2) \quad \phi_{x_i}(z) = \begin{cases} \eta_k(z)V(x_i), & z \in E_{x_i} \\ 0, & \text{otherwise} \end{cases}$$

547 for all  $x_i \in S_k$ . By applying the operator once,  $\xi_k \mapsto (\mathcal{A}\xi_k^*)^*$ , we recover  $\phi_{x_i}$  for every  
 548 point  $x_i \in S_k$ .

549 Each point source,  $\delta_{x_i}$ , is scaled by  $1/V(x_i)$  so that the resulting scaled impulse  
 550 responses within  $\eta_k$  are comparable in magnitude. Without this scaling, the portion  
 551 of  $\phi_{x_i}$  outside of  $E_{x_i}$ , which we neglect, may overwhelm  $\phi_{x_j}$  for a nearby point  $x_j$   
 552 if  $V(x_i)$  is much larger than  $V(x_j)$ . Note that we are not in danger of dividing  
 553 by zero, because the ellipsoid packing procedure from Section 5.1 excludes  $x_i$  from  
 554 consideration as a sample point if  $V(x_i)$  is smaller than a predetermined threshold.

555 **5.3. Approximate integral kernel entries.** Here, we describe how to rapidly  
 556 evaluate arbitrary entries of an approximation to the integral kernel by performing  
 557 radial basis function interpolation of translated and scaled versions of nearby known  
 558 impulse responses. In Section 5.4 we use this procedure for rapidly evaluating kernel  
 559 entries to construct the H-matrix approximation of  $\mathbf{A}$ .

560 Given  $(y, x) \in \Omega \times \Omega$ , let  $z_i := y - \mu(x) + \mu(x_i)$  and define

$$561 \quad (5.3) \quad f_i := \frac{V(x)}{V(x_i)} \phi_{x_i}(z_i)$$

562 for  $i = 1, \dots, k_n$ , where  $\{x_i\}_{i=1}^{k_n}$  are the  $k_n$  nearest sample points to  $x$ , excluding  
 563 points  $x_i$  for which  $z_i \notin \Omega$ . Here,  $k_n$  is a small user-defined parameter, e.g.,  $k_n = 10$ .  
 564 We find the  $k_n$  nearest sample points to  $x$  by querying a precomputed kd-tree [11] of

565 all sample points. We check whether  $z_i \in \Omega$  by querying a precomputed axis aligned  
 566 bounding box tree (AABB tree) [26] of the mesh cells used to discretize the problem.  
 567 Note that  $\phi_{x_i}(z_i)$  is well-defined because  $z_i \in \Omega$ , and  $\frac{V(x)}{V(z_i)}$  is well-defined because  
 568 the sample point choosing procedure in Section 5.1 ensures that  $V(x_i) > 0$ . Per the  
 569 discussion in Section 4.2, we expect  $\Phi(y, x) \approx f_i$  for  $i = 1, \dots, k_n$ . The closer  $x_i$  is to  
 570  $x$ , the better we expect the approximation to be. We therefore approximate  $\Phi(y, x)$   
 571 by interpolating the (point,value) pairs  $\{(x_i, f_i)\}_{i=1}^{k_n}$  at the point  $x$ . Interpolation is  
 572 performed using the following radial basis function [77] scheme:

$$573 \quad (5.4) \quad \Phi(y, x) \approx \tilde{\Phi}(y, x) := \sum_{i=1}^{k_n} c_i \varphi(\|x - x_i\|),$$

574 where  $c_i$  are weights, and  $\varphi(r) := \exp\left(-C_{\text{RBF}}^2 \frac{r^2}{2L^2}\right)$ ,  $\varphi(r) := \exp\left(-\frac{1}{2}\left(C_{\text{RBF}} \frac{r}{r_0}\right)^2\right)$  is  
 575 a Gaussian kernel radial basis function. Here  $L := \text{diam}\left(\{x_i\}_{i=1}^{k_n}\right)$ ,  $r_0 := \text{diam}\left(\{x_i\}_{i=1}^{k_n}\right)$  is  
 576 the diameter of the set of sample points used in the interpolation, and  $C_{\text{RBF}}$  is a  
 577 user-defined shape parameter that controls the width of the kernel function (we use  
 578  $C_{\text{RBF}} = 3.0$ ). The vector of weights,  $c = (c_1, c_2, \dots, c_{k_n})^T$ , is found as the solution to  
 579 the  $k_n \times k_n$  linear system

$$580 \quad (5.5) \quad Bc = f,$$

581 where  $B \in \mathbb{R}^{k_n \times k_n}$ ,  $B_{ij} := \varphi(\|x_i - x_j\|)$ , and  $f \in \mathbb{R}^{k_n}$  has entries  $f_i$  from (5.3).

582 To evaluate  $f_i$ , we check whether  $z_i \in E_{x_i}$  using (4.6). If  $z_i \notin E_{x_i}$ , then  $z_i$  is  
 583 outside the estimated support of  $\phi_{x_i}$ , so we set  $f_i = 0$ . If  $z_i \in E_{x_i}$ , we look up the  
 584 batch index  $b$  such that  $x_i \in S_b$ , and evaluate  $f_i$  via the formula  $f_i = V(x)\eta_b(z_i)$ ,  
 585 per (5.2). Note that  $z_i$  is typically not a gridpoint of the mesh used to discretize  
 586 the problem, even if  $y$ ,  $x$ , and  $x_i$  are gridpoints. Hence, evaluating  $\eta_b(z_i)$  requires  
 587 determining which mesh cell contains  $z_i$ , then evaluating finite element basis functions  
 588 on that mesh cell. Fortunately, the mesh cell containing  $z_i$  was determined as a side  
 589 effect of querying the AABB tree of mesh cells when we checked whether  $z_i \in \Omega$ .

590 The shape parameter,  $C_{\text{RBF}}$ , mediates a tradeoff between accuracy and stability.  
 591 Small  $C_{\text{RBF}}$  is required for RBF interpolation with Gaussian kernels to achieve high  
 592 accuracy, but small  $C_{\text{RBF}}$  also makes RBF interpolation less robust to errors or  
 593 nonsmoothness in the function being interpolated. For our numerical results involving  
 594 Hessians in inverse problems governed by PDEs (Sections 7.1 and 7.2), high accuracy  
 595 is not required because the PSF-based method is used to construct a preconditioner.  
 596 Hence, for these Hessian approximations we use a conservative choice of  $C_{\text{RBF}} = 3.0$  to  
 597 ensure robustness. For our numerical results for the blur problem example (Section 7.3),  
 598 we use a smaller value of  $C_{\text{RBF}} = 0.5$  so that the RBF interpolation accuracy is not  
 599 a limiting factor as we study convergence of the PSF-based method.

600 **5.4. Hierarchical matrix construction.** We form an H-matrix approximation  
 601  $\mathbf{A}_H \approx \mathbf{A}$  by forming an H-matrix representation  $\Phi_H$  of  $\Phi$  then multiplying  $\Phi$  with  
 602 mass matrices  $\mathbf{M}$  per (3.5) to form  $\mathbf{A}_H = \mathbf{M}\Phi_H\mathbf{M}$ . Here, we use a diagonal lumped  
 603 mass matrix, so these matrix-matrix multiplications are trivial. If a non-diagonal  
 604 mass matrix is used, one may form a H-matrix representation of the mass matrix,  
 605 then perform the matrix-matrix multiplications in (3.5) using H-matrix methods. We  
 606 use H1 matrices in the numerical results, but any other H-matrix format could be  
 607 used instead. For more details on H-matrices, see [43].

We form  $\Phi_H$  using the standard geometrical clustering/adaptive cross method implemented within the HLIBpro software package [51]. For details about the algorithms used for geometrical clustering, H-matrix construction, and H-matrix operations in HLIBpro, we refer the reader to [13, 41, 52]. Although  $\Phi$  is a dense  $N \times N$  matrix, constructing  $\Phi_H$  only requires evaluation of  $O(k_h N \log N)$  kernel entries  $\Phi_{ij} = \tilde{\Phi}(p_i, p_j)$  (see [10]), and these entries are computed via the radial basis function interpolation method described in Section 5.3. Here,  $k_h$  is the rank of the highest rank block in the H-matrix. We emphasize that the dense matrix  $\Phi$  is never formed.

**5.5. Symmetrizing and flipping negative eigenvalues (optional).** In many applications, one seeks to approximate an operator  $\mathcal{H} = \mathcal{A} + \mathcal{R}$ , where  $\mathcal{A}$  is a symmetric positive semi-definite operator that we approximate with the PSF-PSF-based method to form an H-matrix  $\mathbf{A}_H$ , and  $\mathcal{R}$  is a symmetric positive definite operator that may be easily converted to an H-matrix  $\mathbf{R}_H$  without using the PSF-PSF-based method. For example, in inverse problems  $\mathcal{H}$  is the Hessian,  $\mathcal{A}$  is the data misfit term in the Hessian which is dense and available only matrix-free, and  $\mathcal{R}$  is the regularization term, which is typically an elliptic differential operator that becomes a sparse matrix after discretization.

The PSF-PSF-based approximation  $\mathbf{A}_H$ , and therefore  $\mathbf{A}_H + \mathbf{R}_H$ , may be non-symmetric and indefinite because of approximation error. This is undesirable because symmetry and positive semi-definiteness are important properties which should be preserved if possible. Also, lacking these properties may prevent one from using highly effective algorithms to perform further operations involving  $\mathbf{A}_H + \mathbf{R}_H$ , such as using  $\mathbf{A}_H + \mathbf{R}_H$  as a preconditioner in the conjugate gradient method.

We modify  $\mathbf{A}_H$  to make it symmetric and remove negative eigenvalues via the following procedure. First, we symmetrize  $\mathbf{A}_H$  via  $\mathbf{A}_H^{\text{sym}} := \frac{1}{2}(\mathbf{A}_H + \mathbf{A}_H^T)$ . Next, we find negative eigenvalues and their corresponding eigenvectors for the generalized eigenvalue problem  $\mathbf{A}_H^{\text{sym}}\mathbf{u} = \lambda\mathbf{R}_H$  using a Cayley shift-and-invert Krylov scheme [53]. We flip the signs of these eigenvalues to be positive instead of negative (i.e.,  $\lambda \rightarrow |\lambda|$ ) by performing a low rank update to  $\mathbf{A}_H^{\text{sym}}$ . We observe that the eigenvectors associated with large erroneous negative eigenvalues tend to be directions that are nevertheless “important” to  $\mathcal{A}$ , so flipping the eigenvalues instead of setting them to zero tends to yield better approximations. The primary computational task in the Cayley shift-and-invert scheme is the solution of shifted linear systems of the form  $(\mathbf{A}_H^{\text{sym}} + \mu_i \mathbf{R}_H) \mathbf{x} = \mathbf{b}$ , for a small number of positive shifts  $\mu_i$ . We solve these linear systems by factorizing the matrices  $\mathbf{A}_H^{\text{sym}} + \mu_i \mathbf{R}_H$  using fast H-matrix methods. We compute and flip all eigenvalues  $\lambda < \epsilon_{\text{flip}}$  which are less than some threshold  $\epsilon_{\text{flip}} \in (-1, 0]$ . By choosing  $\epsilon_{\text{flip}} > -1$ , we ensure that the modified version of  $\mathbf{A}_H^{\text{sym}} + \mathbf{R}_H$  is positive definite. Choosing  $\epsilon_{\text{flip}} = 0$  would remove all erroneous negative eigenvalues. However, this is computationally infeasible if  $\mathcal{A}$  has a large or infinite cluster of eigenvalues near zero, a common situation for Hessians in ill-posed inverse problems. We therefore recommend choosing  $\epsilon_{\text{flip}} < 0$ . In our numerical results, we use  $\epsilon_{\text{flip}} = -0.1$ .

**6. Computational cost.** The computational cost of the proposed-PSF-based method may be divided into the costs to perform the following tasks: (1) Computing impulse response moments and batches (Lines 1 and 4 in Algorithm 1); (2) Building the H-matrix (Lines 5 and 6 in Algorithm 1); (3) Performing linear algebra operations with the H-matrix. This may optionally include the symmetric positive semi-definite modifications described in Section 5.5. In target applications, (1) is the dominant cost because applying  $\mathcal{A}$  to a vector requires an expensive computational procedure such as solving a PDE, and (1) is the only step that requires applying  $\mathcal{A}$  to vectors. All

Symbol	Typical size	Variable name
$N$	$10^3\text{--}10^9$	Number of finite element degrees of freedom
$n_b$	1–25	Number of batches
$k_h$	5–50	H-matrix rank
$k_n$	5–15	Number of nearest neighbors for RBF interpolation
$d$	1–3	Spatial dimension
$m$	$10^1\text{--}10^4$	Total number of sample points (all batches)
$ S_i $	1–500	Number of sample points in the $i$ th batch

Table 1: Symbols used for variables in computational cost estimates, and approximate ranges for their sizes in practice.

657 operations that do not require applications of  $\mathcal{A}$  to vectors are nearly linear polylog  
 658 linear (i.e.,  $O(N \log(N)^b)$  for some  $b$ ), and therefore scalable, in the size of the problem,  
 659  $N$ . We now describe these costs in detail. For convenience, Table 1 lists variable  
 660 symbols and their approximate sizes.

661 (1) *Computing impulse response moments and batches.* Computing  $V$ ,  $\mu$ , and  $\Sigma$   
 662 requires applying  $\mathcal{A} - \mathcal{A}^T$  to 1,  $d$ , and  $d(d+1)/2$  vectors, respectively. This works  
 663 out to 3 applications of  $\mathcal{A}^T$  in one spatial dimension, 6 in two dimensions, and 10 in  
 664 three dimensions. Computing each  $\eta_i$  requires applying  $\mathcal{A}$  to one vector, so computing  
 665  $\{\eta_i\}_{i=1}^{n_b}$  requires  $n_b$  operator applications. In total, computing all impulse response  
 666 moments and batches therefore requires

$$667 \quad 1 + d + d(d+1)/2 + n_b \quad \text{operator applications.}$$

668 In a typical application one might have  $d = 2$  and  $n_b = 5$ , in which case a modest 11  
 669 operator applications are required.

670 Computing the impulse response batches also requires choosing sample point  
 671 batches via the greedy ellipsoid packing algorithm described in Section 5.1. Choosing  
 672 the  $i$ th batch of sample points may require performing up to  $N|S_i|$  ellipsoid intersec-  
 673 tion tests, where  $|S_i|$  is the number of sample points in the  $i$ th batch. Choosing all  
 674 of the sample points therefore requires performing at most

$$675 \quad Nm \quad \text{ellipsoid intersection tests,}$$

676 where  $m$  is the total number of sample points in all batches. The multiplicative  
 677 dependence of  $N$  with  $m$  is undesirable since  $m$  may be large, and reducing this cost  
 678 is possible with more involved computational geometry methods. However, from a  
 679 practical perspective, the cost of choosing sample points is small compared to other  
 680 parts of the algorithm, and hence such improvements are not pursued here.

681 (2) *Building the H-matrix.* Classical H-matrix construction techniques require  
 682 evaluating  $O(k_h N \log N)$  matrix entries of the approximation [10], where  $k_h$  is the H-  
 683 matrix rank, i.e, the maximum rank among the blocks of the H-matrix. To evaluate  
 684 one matrix entry, first one must find the  $k_n$  nearest sample points to a given point,  
 685 where  $k_n$  is the number of impulse responses used in the RBF interpolation. This is  
 686 done using a precomputed kd-tree of sample points, and requires  $O(k_n \log m)$  floating  
 687 point and logical elementary operations. Second, one must find the mesh cells that  
 688 the points  $\{z_i\}_{i=1}^{k_n}$  reside in. This is done using an AABB tree of mesh cells, and  
 689 requires  $O(k_n \log N)$  elementary operations. Third, one must evaluate finite element

690 basis functions on those cells, which requires  $O(k_n)$  elementary operations. Finally,  
 691 the radial basis function interpolation requires solving a  $k_n \times k_n$  linear system, which  
 692 requires  $O(k_n^3)$  elementary operations. Therefore, building the H-matrix requires

$$693 \quad O((k_h N \log N)(k_n \log N + k_n^3)) \quad \text{elementary operations.}$$

694 (3) *Performing linear algebra operations with the H-matrix.* It is well known that  
 695 H-matrix methods for matrix-vector products, matrix-matrix addition, matrix-matrix  
 696 multiplication, matrix factorization, matrix inversion, and low rank updates require  
 697 performing  $O(k_h^a N \log(N)^b)$  elementary operations, where  $a, b \in \{0, 1, 2, 3\}$  are con-  
 698 stants which depend on the type of H-matrix used and the operation being per-  
 699 formed [40][52, Section 2.1]. ~~In the numerical results (Section 7)~~ For our numerical  
 700 results involving Hessians (Sections 7.1 and 7.2), we use one matrix-matrix addition  
 701 to add the H-matrix approximation of the data misfit term in the Hessian to the  
 702 regularization term in the Hessian. Symmetrizing  $\mathbf{A}_H$  requires one matrix-matrix  
 703 addition. Flipping negative eigenvalues to be positive requires a handful (typically  
 704 around 5) of matrix-matrix additions and matrix factorizations to factor the required  
 705 shifted linear systems, and a number of factorized solves that is proportional to the  
 706 number of erroneous negative eigenvalues.

707 In summary, computing all the necessary ingredients to evaluate kernel entries of  
 708 the ~~PSF-PSF-based~~ approximation requires a handful of operator applications (e.g.,  
 709  $6+n_b$  operator applications in two dimensions, or  $10+n_b$  operator applications in three  
 710 dimensions, with  $n_b$  typically in the range 1–25), plus comparatively cheap additional  
 711 overhead costs, most notably performing ellipsoid intersection tests while choosing  
 712 sample point batches. Once these ingredients are computed, no more operator appli-  
 713 cations (and thus PDE solves) are required, and approximate kernel entries can be  
 714 evaluated rapidly. Constructing the H-matrix from kernel entries requires a number  
 715 of elementary operations that scales polylog linearly in  $N$ . Using the ~~resulting~~ H-  
 716 matrix to perform ~~further~~ linear algebra operations also scales ~~nearly~~-polylog linearly  
 717 in  $N$ , though the details of these costs depend heavily on the type of H-matrix and  
 718 operation being performed.

719 **7. Numerical results.** We use the ~~proposed~~ PSF-based method to approximate  
 720 the Newton (or Gauss-Newton) Hessians in inverse problems governed by PDEs which  
 721 model steady state ice sheet flow [64] (Section 7.1) and advective-diffusive transport  
 722 of a contaminant [63] (Section 7.2). ~~These inverse, and to approximate the integral~~  
 723 ~~kernel in a blur problem that is not based on PDEs (Section 7.3).~~ These problems  
 724 are described in detail in their respective sections. ~~In both cases~~

725 ~~In both PDE-based inverse problems (Sections 7.1 and 7.2),~~ to reconstruct the  
 726 unknown parameter fields, denoted  $q$ , the inverse problems are formulated as nonlinear  
 727 least squares optimization problems, whose objective functions consist of a data misfit  
 728 term (between the observations and model output) and a bi-Laplacian regularization  
 729 term following [75]. The regularization is centered at a constant function  $q_0(x)$ . To  
 730 mitigate boundary effects we use a constant coefficient Robin boundary condition as  
 731 in [65]. The parameters for the bi-Laplacian operator are chosen so that the Green's  
 732 function of the Hessian of the regularization has a characteristic length of 0.25 of the  
 733 domain radius. For the specific setup, we refer the reader to [75, Section 2.2]. In all  
 734 numerical results we choose the regularization parameter (which controls the overall  
 735 strength of the regularization) using the Morozov discrepancy principle [76].

736 We solve the ice sheet inverse problem with an inexact Newton preconditioned  
 737 conjugate gradient (PCG) scheme and a globalizing Armijo line search [61]. The

738 Newton search directions,  $\hat{\mathbf{q}}$ , are obtained by solving

739 (7.1) 
$$\mathbf{H}\hat{\mathbf{q}} = -\mathbf{g} \quad \text{or} \quad \mathbf{H}_{\text{gn}}\hat{\mathbf{q}} = -\mathbf{g},$$

740 wherein we choose the initial guess as the discretization of the constant function  $q_0$ .  
 741 Here,  $\mathbf{g}$ ,  $\mathbf{H}$  and  $\mathbf{H}_{\text{gn}}$  are the discretized gradient, Hessian, and Gauss-Newton Hessian  
 742 of the inverse problem objective function, respectively, evaluated at the current  
 743 Newton iterate. To ensure positive definiteness of the Hessian we use  $\mathbf{H}_{\text{gn}}$  for the  
 744 first five iterations, and  $\mathbf{H}$  for all subsequent iterations. The Newton iterations are  
 745 terminated when  $\|\mathbf{g}\| < 10^{-6}\|\mathbf{g}_0\|$ , where  $\mathbf{g}_0$  is the gradient evaluated at the initial  
 746 guess. Systems (7.1) are solved inexactly using an inner PCG iteration, which is ter-  
 747 minated early based on the Eisenstat-Walker [24] and Steihaug [69] conditions. The  
 748 inverse problem governed by the advection-diffusion PDE is linear, hence Newton's  
 749 method converges in one iteration. In this case the Newton linear system, (7.1), is  
 750 solved using PCG, using termination tolerances described in Section 7.2.

751 We use the framework described in this paper to generate Hessian preconditioners.  
 752 We build H-matrix approximations,  $\mathbf{A}_H$ , of the data misfit Gauss-Newton Hessian  
 753 (the term in  $\mathbf{H}_{\text{gn}}$  that arises from the data misfit). The approximations are indicated  
 754 by “PSF ( $n_b$ )”, where  $n_b$  is the number of impulse response batches used to build the  
 755 approximation. The Hessian of the regularization term is a combination of stiffness  
 756 and mass matrices, which are sparse. Therefore, we form H-matrix representations of  
 757 these matrices and combine them into ~~a-an~~ H-matrix approximation of the regular-  
 758 ization term in the Hessian,  $\mathbf{R}_H$ , using standard sparse H-matrix techniques. Then,  
 759 H-matrix approximations of the Gauss-Newton Hessian,  $\mathbf{H}_{\text{gn}} \approx \tilde{\mathbf{H}} := \mathbf{A}_H + \mathbf{R}_H$ ,  
 760 are formed by adding  $\mathbf{A}_H$  to  $\mathbf{R}_H$  using fast H-matrix arithmetic. We modify  $\tilde{\mathbf{H}}$  to  
 761 be (approximately) symmetric positive semi-definite via the procedure described in  
 762 Section 5.5. We factor  $\tilde{\mathbf{H}}$  using fast H-matrix methods, then use the factorization  
 763 as a preconditioner. We approximate  $\mathbf{H}_{\text{gn}}$  rather than  $\mathbf{H}$  because  $\mathbf{H}$  more often has  
 764 negative values in its integral kernel. The numerical results show that  $\tilde{\mathbf{H}}$  is a good  
 765 preconditioner for both  $\mathbf{H}_{\text{gn}}$  and  $\mathbf{H}$ .

766 **7.1. Example 1: Inversion for the basal friction coefficient in an ice**  
 767 **sheet flow problem.** For this example, we consider a sheet of ice flowing down a  
 768 mountain (see Figure 8a). Given observations of the tangential component of the ice  
 769 velocity on the top surface of the ice, we invert for the logarithm of the unknown  
 770 spatially varying basal friction Robin coefficient field, which governs the resistance  
 771 to sliding along the base of the ice sheet. The setup, which we briefly summarize,  
 772 follows [48, 64]. The region of ice is denoted by  $\mathcal{D} \subset \mathbb{R}^3$ . The basal, lateral and top  
 773 parts of the boundary  $\partial\mathcal{D}$  are denoted by  $\Gamma_b$ ,  $\Gamma_l$ , and  $\Gamma_t$ , respectively. The governing  
 774 equations are the linear incompressible Stokes equations,

775 (7.2a) 
$$-\nabla \cdot \sigma(v, p) = f \text{ and } \nabla \cdot v = 0 \quad \text{in } \mathcal{D},$$

776 (7.2b) 
$$\sigma(v, p)\nu = 0 \quad \text{on } \Gamma_t,$$

777 (7.2c) 
$$v \cdot \nu = 0 \text{ and } T(\sigma(v, p)\nu + \exp(q)v) = 0 \quad \text{on } \Gamma_b,$$

778 (7.2d) 
$$\sigma(v, p)\nu + sv = 0 \quad \text{on } \Gamma_l.$$

780 The solution to these equations is the pair  $(v, p)$ , where  $v$  is the ice flow velocity field<sup>6</sup>  
 781 and  $p$  is the pressure field. Here,  $q$  is the unknown logarithmic basal friction field

---

<sup>6</sup>We do not use bold to denote vector or tensor fields to avoid confusion with vectors that arise from finite element discretizations, which are already denoted with bold.

(large  $q$  corresponds to large resistance to sliding) defined on the surface  $\Gamma_b$ . The quantity  $f$  is the body force density due to gravity,  $s = 10^6$  is a Robin boundary condition constant,  $\nu$  is the outward unit normal and  $T$  is the tangential projection operator that restricts a vector field to its tangential component along the boundary. We employ a Newtonian constitutive law,  $\sigma(v, p) = 2\eta\dot{\varepsilon}(v) - Ip$ , where  $\sigma$  is the stress tensor and  $\dot{\varepsilon}(v) = \frac{1}{2}(\nabla v + \nabla v^\top)$  is the strain rate tensor [48]. Here,  $\eta$  is the viscosity and  $I$  is the identity operator. Note that while the PDE is linear, the parameter-to-solution map,  $q \mapsto (v, p)$ , is nonlinear.

The pressure,  $p$ , is discretized with first order scalar continuous Galerkin finite elements defined on a mesh of tetrahedra. The velocity,  $v$ , is discretized with second order continuous Galerkin finite elements on the same mesh. The parameter  $q$  is discretized with first order scalar continuous Galerkin finite elements on the mesh of triangles that results from restricting the tetrahedral mesh to the basal boundary,  $\Gamma_b$ . Note that  $\Gamma_b$  is a two-dimensional surface embedded in three dimensions due to the mountain topography. The proposed PSF-based method involves translating impulse responses. Hence it requires either a flat domain, or a notion of local parallel transport. We therefore generate a flattened version of  $\Gamma_b$ , denoted by  $\Omega \subset \mathbb{R}^2$ , by ignoring the height coordinate. The parameter  $q$  is viewed as a function on  $\Gamma_b$  for the purpose of solving the Stokes equations, and as a function on  $\Omega$  for the purpose of building Hessian approximations and defining the regularization. The observations are generated by adding multiplicative Gaussian noise to the tangential component of the velocity field restricted to the top surface of the geometry. We use 5% noise in all cases, except for Figure 9 and Table 3 where the noise is varied from 1% to 25% and the regularization is determined by the Morozov discrepancy principle for each noise level. The true basal friction coefficient and resulting velocity fields, which are obtained by solving (7.2), are shown in Figure 8.

Table 2 shows the performance of the preconditioner for accelerating the solution of the optimization problem to reconstruct  $\mathbf{q}$  from observations with 5% noise. We build the PSF (5) preconditioner in the third Gauss-Newton iteration, and reuse it for all subsequent Gauss-Newton and Newton iterations. No preconditioning is used in the iterations before the PSF (5) preconditioner is built. We compare the proposed PSF-based method with the most commonly used existing preconditioners: no preconditioning (NONE), and preconditioning by the regularization term in the Hessian (REG). The results show that using PSF (5) reduces the total number of Stokes PDE solves to 70, as compared to 908 for regularization preconditioning and 308 for no preconditioning, a reduction in cost of roughly  $5\times$ – $10\times$ . For problems with a larger physical domain and correspondingly more observations, such as continental scale ice sheet inversion, the speedup will be even greater. This is because the rank of the data misfit Hessian will increase, while the locality of the impulse responses will remain the same. In Figure 9 we show reconstructions for 1%, 5%, and 25% noise.

Next, we build PSF (1), PSF (5), and PSF (25) preconditioners based on the Gauss-Newton Hessian evaluated at the converged solution  $\mathbf{q}$  (note:  $k$  in PSF ( $k$ ) refers to the number of batches; this is not to be confused with the noise levels which range over the same numerical values). We use PCG to solve a linear system with the Hessian as the coefficient operator and a right hand side vector with random independent and identically distributed (i.i.d.) entries drawn from the standard Gaussian distribution. In Figure 10 (left) we compare the convergence of PCG for solving this linear system using the PSF (1), PSF (5), PSF (25), REG, and NONE preconditioners. PCG converges fastest with the PSF-based preconditioners, with PSF (25) converging fastest, followed by PSF (5), followed by PSF (1), as expected. In Figure 10 (right) we

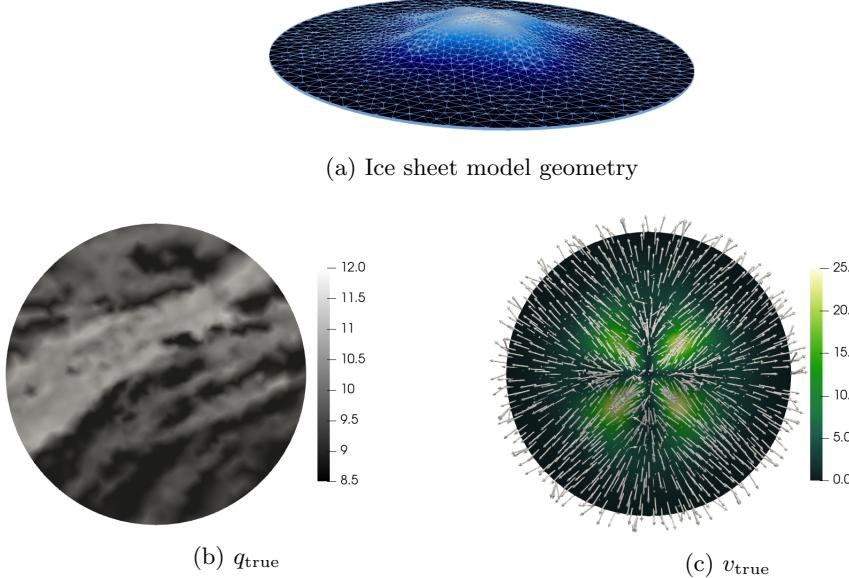


Fig. 8: (Ice sheet) (8a) Bird's eye view of the ice sheet discretized by a mesh of tetrahedra. Color indicates the height of the base of the ice sheet (i.e., the mountain topography). The radius of the domain is  $10^4$  meters, the maximum height of the mountain is  $2.1 \times 10^3$  meters, and the average thickness of the ice sheet is 250 meters. (8b) True parameter,  $q_{\text{true}}$ . (8c) True velocity,  $v_{\text{true}}$ . Arrows indicate the direction of  $v_{\text{true}}$  and color indicates the magnitude of  $v_{\text{true}}$ .



Fig. 9: (Ice sheet) The log basal friction parameter, with color scale as in Figure 8b, computed from the PDE constrained optimization problem with noise levels: 25% (left), 5.0% (middle), and 1.0% (right).

832 show the generalized eigenvalues for the generalized eigenvalue problem  $\mathbf{H}\mathbf{u} = \lambda \tilde{\mathbf{H}}\mathbf{u}$ .  
 833 The matrix  $\tilde{\mathbf{H}}$  is one of the PSF (1), PSF (5), or PSF (25) Gauss-Newton Hessian  
 834 approximations, the regularization Hessian (REG), or the identity matrix (NONE).  
 835 With the PSF-based preconditioners, the generalized eigenvalues cluster near one,  
 836 with more batches yielding better clustering.

837 In Table 3, we show the condition number of the preconditioned Hessian for noise  
 838 levels ranging from 1% to 25%. Note that the condition number using PSF-based

Iter	PSF (5)			REG			NONE		
	#CG	#Stokes	$\ \mathbf{g}\ $	#CG	#Stokes	$\ \mathbf{g}\ $	#CG	#Stokes	$\ \mathbf{g}\ $
0	1	4	1.9e+7	3	8	1.9e+7	1	4	1.9e+7
1	2	6	6.1e+6	8	18	8.4e+6	2	6	6.1e+6
2	4	10	2.6e+6	16	34	4.1e+6	4	10	2.6e+6
3	2	6+22	6.9e+5	34	70	1.8e+6	14	30	6.9e+5
4	3	8	4.4e+4	52	106	5.6e+5	29	60	1.3e+5
5	5	12	2.2e+3	79	160	9.4e+4	38	78	1.0e+4
6	0	2	1.1e+1	102	206	6.5e+3	58	118	1.8e+2
7	—	—	—	151	304	1.2e+2	0	2	5.5e-1
8	—	—	—	0	2	2.9e-1	—	—	—
Total	17	70	—	445	908	—	146	308	—

Table 2: (Ice sheet) Convergence history for solving the Stokes inverse problem using inexact Newton PCG to tolerance  $10^{-6}$ . Preconditioners shown are the [proposed PSF](#) [PSF-based](#) method with five batches (PSF (5)) constructed at the third iteration, regularization preconditioning (REG), and no preconditioning (NONE). Columns #CG show the number of PCG iterations used to solve the Newton system for  $\hat{\mathbf{q}}$ . Columns  $\|\mathbf{g}\|$  show the  $l^2$  norm of the gradient at  $\mathbf{q}$ . Columns #Stokes show the total number of Stokes PDE solves performed in each Newton iteration. Under PSF (5) and in row Iter 3, we write  $6 + 22$  to indicate that 6 Stokes solves were used during the standard course of the iteration, and 22 Stokes solves were used to build the PSF (5) preconditioner.

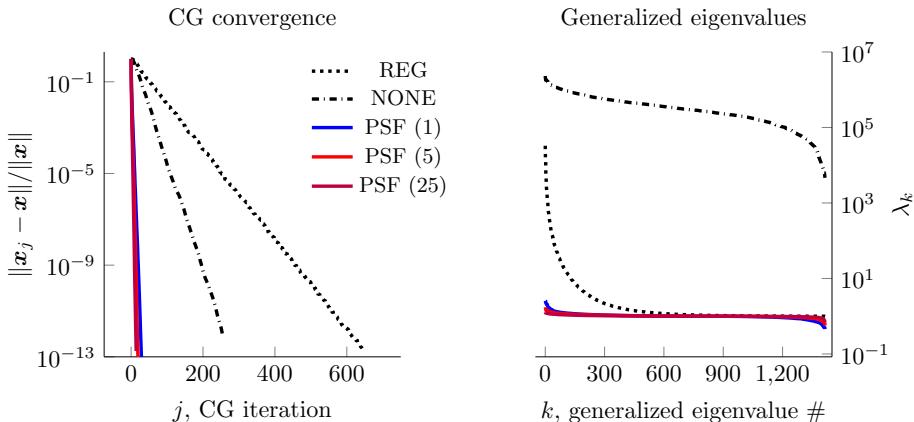


Fig. 10: (Ice sheet) Left: [convergence](#) history for solving  $\mathbf{Hx} = \mathbf{b}$  using PCG, where  $\mathbf{b}$  has i.i.d. random entries drawn from the standard Gaussian distribution and  $\mathbf{H}$  is evaluated at the solution of the inverse problem. Results in these figures are shown for the PSF-based preconditioners with 1, 5, and 25 batches (PSF (1), PSF (5), and PSF (25), respectively), regularization preconditioning (REG), and no preconditioning (NONE). The preconditioner is constructed using  $\mathbf{H}_{\text{gn}}$ . Right: [generalized](#) [Generalized](#) eigenvalues for generalized eigenvalue problem  $\mathbf{Hu}_k = \lambda_k \tilde{\mathbf{H}} \mathbf{u}_k$ . Here  $\tilde{\mathbf{H}}$  is the Hessian and the matrices  $\tilde{\mathbf{H}}$  are the same Hessian approximations used in the left sub-figure, with NONE corresponding to the identity matrix.

noise level	COND( $\tilde{\mathbf{H}}^{-1}\mathbf{H}$ )				
	REG	NONE	PSF (1)	PSF (5)	PSF (25)
25%	1.01e+3	2.96e+3	1.34e+0	1.30e+0	1.18e+0
11%	7.40e+3	1.05e+3	2.27e+0	1.55e+0	1.31e+0
5.0%	3.29e+4	4.96e+2	5.61e+0	3.06e+0	1.92e+0
2.2%	1.66e+5	8.89e+2	1.58e+1	8.07e+0	4.03e+0
1.0%	5.36e+5	1.61e+3	7.17e+1	1.93e+1	9.19e+0

Table 3: (Ice sheet) Condition number for  $\tilde{\mathbf{H}}^{-1}\mathbf{H}$  for the PSF-based preconditioners with 1, 5, and 25 batches (PSF (1), PSF (5), and PSF (25), respectively), no preconditioner (NONE) and regularization preconditioning (REG). All operators are evaluated at the soutions of the inverse problems for their respective noise levels.

839 preconditioners is extremely small (ranging between 1 and 10) and relatively stable  
 840 over this range of noise levels. As expected, PSF (25) outperforms PSF (5), which  
 841 outperforms PSF (1). All PSF-based preconditioners outperform regularization and  
 842 no preconditioning by several orders of magnitude for all noise levels.

843 **7.2. Example 2: Inversion for the initial condition in an advective-  
 844 diffusive transport problem.** Here, we consider a time-dependent advection-diffusion  
 845 equation in which we seek to infer the unknown spatially varying initial condition,  
 846  $q$ , from noisy observation of the full state at a final time,  $T$ . This PDE models  
 847 advective-diffusive transport in a domain  $\Omega \subset \mathbb{R}^d$ , which is depicted in Figure 11. In  
 848 this case, the state,  $c(x, t)$ , could be interpreted as the concentration of a contami-  
 849 nant. The problem description below closely follows [63, 75]. The domain boundaries  
 850  $\partial\Omega$  include the outer boundaries as well as the internal boundaries of the rectangles,  
 851 which represent buildings. The parameter-to-observable map  $\mathcal{F}$  in this case maps an  
 852 initial condition  $q \in L^2(\Omega)$  to the concentration field at a final time,  $c(x, T)$ , through  
 853 solution of the advection-diffusion equation given by

$$\begin{aligned} c_t - \kappa \Delta c + v \cdot \nabla c &= 0 && \text{in } \Omega \times (0, T), \\ 854 \quad (7.3) \quad c(\cdot, 0) &= q && \text{in } \Omega, \\ \kappa \nabla c \cdot \nu &= 0 && \text{on } \partial\Omega \times (0, T). \end{aligned}$$

855 Here,  $\kappa > 0$  is a diffusivity coefficient,  $\nu$  is the boundary unit normal vector, and  $T > 0$   
 856 is the final time. The velocity field,  $v : \Omega \rightarrow \mathbb{R}^d$ , is computed by solving the steady-  
 857 state Navier-Stokes equations for a two dimensional flow with Reynolds number 50,  
 858 with boundary conditions  $v(x) = (0, 1)$  on the left boundary,  $v(x) = (0, -1)$  on the  
 859 right boundary, and  $v(x) = (0, 0)$  on the top and bottom boundaries, as in [63, Section  
 860 3]. We use a checkerboard image for the initial condition (Figure 11a) and add 5%  
 861 multiplicative noise to generate a synthetic observation at the final time,  $T$ . The  
 862 initial condition, velocity field, noisy observations, and reconstructed initial condition  
 863 are shown in Figure 11. We use  $\kappa = 3.2e-1$  and  $T = 1.0$  for all results, except for  
 864 Table 4 and Figure 12 where we vary  $\kappa$  and  $T$ .

865 In Table 4 we show the number of PCG iterations,  $j$ , required to solve the Newton  
 866 linear system to a relative error tolerance of  $\|\hat{\mathbf{q}} - \hat{\mathbf{q}}_j\| < 10^{-6}\|\hat{\mathbf{q}}\|$ . The solution of  
 867 the Newton system to which we compare,  $\hat{\mathbf{q}}$ , is found via another PCG iteration with  
 868 a relative residual tolerance of  $10^{-11}$ . We show results for  $T$  ranging from 0.5 to 2.0

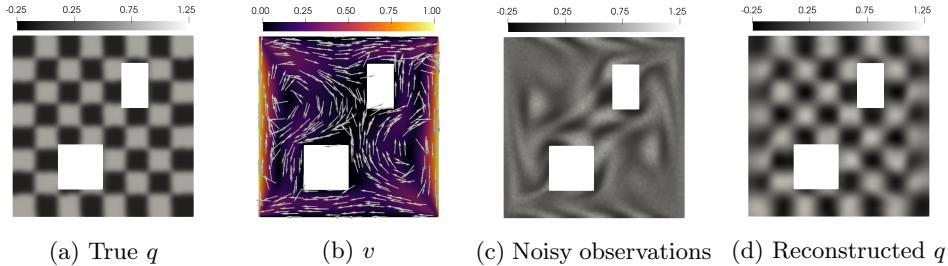


Fig. 11: (Advection-diffusive transport) (11a) True initial condition. (11b) Velocity field. Color indicates magnitude of velocity vector. (11c) Noisy observations of concentration at the final time. (11d) Reconstructed initial condition.

and  $\kappa$  ranging from  $10^{-4}$  to  $10^{-3}$  using the PSF-based preconditioners with 1, 5, and 25 batches, regularization preconditioning, and no preconditioning. The results show that PSF-based preconditioning outperforms regularization preconditioning and no preconditioning in all cases except one. The exception is  $T = 2.0$  and  $\kappa = 1.0e-3$ , in which PSF (1) performs slightly worse than regularization preconditioning but better than no preconditioning. Adding more batches yields better results, and the impact of adding more batches is more pronounced here than in the ice sheet example. For example, in the mid-range values  $T = 1.0$  and  $\kappa = 3.2e-4$ , PSF (1), PSF (5), and PSF (25) require  $1.3\times$ ,  $2.3\times$ , and  $5.1\times$  fewer PCG iterations, respectively, as compared to no preconditioning, and exhibit greater improvements as compared to regularization preconditioning. The PSF preconditioners perform best in the high rank regime where  $T$  is small, which makes sense given that short simulation times yield more localized impulse responses (see Figure 4). For example, for  $\kappa = 3.2e-4$  using the PSF (5) preconditioner yields 140, 202, and 379 iterations for  $T = 0.5$ , 1.0, and 2.0, respectively. The performance of the PSF preconditioners as a function of  $\kappa$  does not have as clear of a trend. Reducing  $\kappa$  makes the impulse responses thinner and hence easier to fit in batches, but also increases the complexity of the impulse response shapes, which may reduce the accuracy of the RBF interpolation. The greatest improvements are seen for  $T = 0.5$  and  $\kappa = 1e-3$ , for which PSF (25) requires roughly  $10\times$  and  $20\times$  fewer PCG iterations than no preconditioning and regularization preconditioning, respectively.

In Figure 13 (left) we show the convergence of PCG for solving  $\mathbf{H}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{b}$  has i.i.d. random entries drawn from a standard Gaussian distribution. The preconditioners used,  $\tilde{\mathbf{H}}$ , are the PSF-based preconditioners with 1, 5, or 25 batches, the regularization Hessian, and the identity matrix (i.e., no preconditioning). The results show that PCG converges fastest with the PSF-based preconditioners, with more batches yielding faster convergence. In Figure 13 (right), we show the eigenvalues for the generalized eigenvalue problem  $\mathbf{H}\mathbf{u} = \lambda\tilde{\mathbf{H}}\mathbf{u}$ , where the  $\tilde{\mathbf{H}}$  are the preconditioners stated above. With the PSF-based preconditioners the eigenvalues cluster near one, and more batches yields better clustering. With the regularization preconditioner the trailing eigenvalues cluster near one, while the leading eigenvalues are amplified.

**7.3. Example 3: Spatially varying blurring problem.** Here we define a PDE-free spatially varying blur problem, in which the impulse response,  $\phi_x$ , is a bumpy blob that is centered near  $x$ , and is rotated and scaled in a manner that

	$\kappa$	REG	NONE	PSF (1)	PSF (5)	PSF (25)
$T = 0.5$	1.0e-4	584	317	311	151	56
	3.2e-4	685	311	233	140	44
	1.0e-3	702	324	122	71	33
$T = 1.0$	1.0e-4	634	449	539	288	100
	3.2e-4	681	459	350	202	90
	1.0e-3	574	520	266	260	208
$T = 2.0$	1.0e-4	609	591	548	520	165
	3.2e-4	524	645	318	379	170
	1.0e-3	349	786	381	262	158

Table 4: (Advection-diffusive transport) Number of PCG iterations required to solve the Newton linear system to tolerance  $\|\hat{\mathbf{q}}_j - \hat{\mathbf{q}}\| < 10^{-6}\|\hat{\mathbf{q}}\|$ , where  $\hat{\mathbf{q}}_j$  is the  $j$ th iterate, and  $\hat{\mathbf{q}}$  is the solution of the Newton linear system. Iteration counts are shown for a variety of different diffusion parameters  $\kappa$ , simulation times  $T$ , and preconditioners.

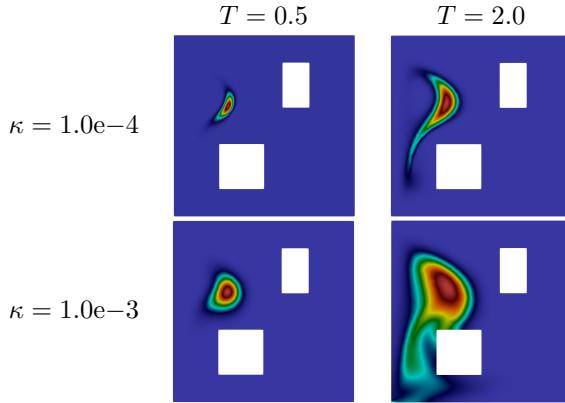


Fig. 12: (Advection-diffusive transport) Impulse responses for small and large diffusion parameters  $\kappa$  and simulation times  $T$ .

depends on  $x$  (see Figure 2). This blur problem is used in Sections 1 and 3 to visually illustrate various stages and aspects of the PSF-based method, and robustness (or lack thereof) of the PSF-based method to violations of the non-negative kernel assumption (Section 3.2). In this section, the blur problem is used to compare the PSF-based method to hierarchical off-diagonal low rank (HODLR) and global low rank (GLR) methods, to study convergence of the PSF-based method, and to investigate the effect of the ellipsoid size parameter  $\tau$ . The closed form expression for the integral kernel is given by

$$(7.4) \quad \Phi(y, x) = (1 - af(y, x))g(x) \exp\left(-\frac{1}{2}(h(y, x)^T C^{-1} h(y, x)\right),$$

where  $f(y, x) = \cos\left(h^1(y, x)/\sqrt{c_1/2}\right) \sin\left(h^2(y, x)/\sqrt{c_2/2}\right)$ , with  $h^i(y, x)$  the  $i^{th}$  component of  $R(\theta(x))(y - x)$ , with  $R(\theta(x))$  a two-dimensional rotation matrix by angle  $\theta(x) = (x^1 + x^2)\pi/2$ ,  $g(x) = x^1(1 - x^1)x^2(1 - x^2)$ , and with  $C = L^2 \text{diag}(c_1, c_2)$ . The constant  $L$  controls

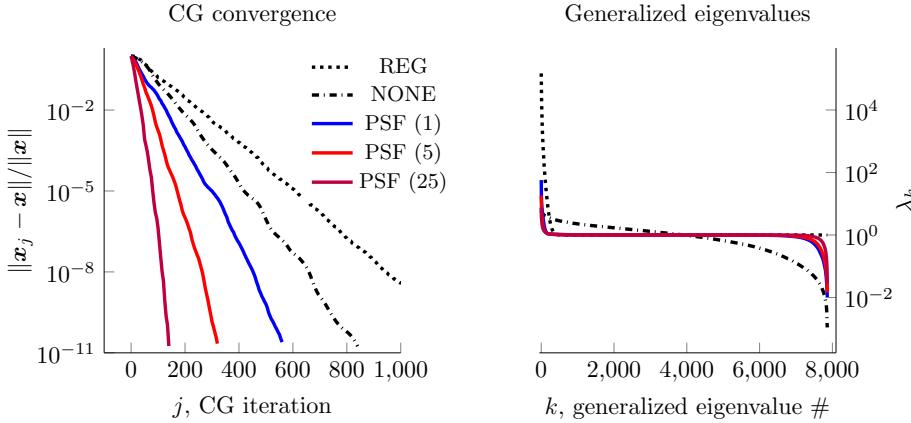


Fig. 13: (Advection-diffusive transport) Left: [convergence](#)—Convergence history for solving  $\mathbf{H}\mathbf{x} = \mathbf{b}$  using PCG, where  $\mathbf{b}$  has i.i.d. random entries drawn from the standard Gaussian distribution. Right: [generalized](#)—Generalized eigenvalues for generalized eigenvalue problem  $\mathbf{H}\mathbf{u}_k = \lambda_k \tilde{\mathbf{H}}\mathbf{u}_k$ . Here,  $\mathbf{H}$  is the Hessian and the preconditioner,  $\tilde{\mathbf{H}}$ , is the [PSF](#)—[PSF-based](#) approximation for 1, 5, or 25 batches (PSF (1), PSF (5), and PSF (25), respectively), the regularization Hessian (REG), or the identity matrix (NONE).

916 the width of the blob, and  $c_1/c_2$  controls its aspect ratio. The constant  $a$  represents  
 917 deviation from a Gaussian. When  $a = 0$ ,  $\phi_x$  is a Gaussian, and as  $a$  increases,  $\phi_x$   
 918 becomes non-Gaussian. When  $a > 1$ , the integral kernel contains negative values,  
 919 which allows us to study the robustness of the PSF-based method to violations of  
 920 Assumption 3 (Section 3.2).

921 In Table 5 we compare the cost to approximate the blur kernel from Equation (7.4)  
 922 using the PSF-based method, the randomized HODLR (hierarchical off diagonal low  
 923 rank) method [58, 45] with 8 levels, and GLR (global low rank) approximation using  
 924 double-pass randomized SVD [44]. For these results we vary the quantity  $L$  to scale  
 925 the width of the impulse responses and hence the rank of the operator. For each case,  
 926 we calculate the relative error in the approximation of the kernel measured in the  
 927 Frobenius norm,  $\|\Phi - \tilde{\Phi}\|_{\text{Fro}}/\|\Phi\|_{\text{Fro}}$  and show the number of operator applications  
 928 required to achieve 20%, 10%, and 5% relative error by each method. The results  
 929 reveal superior performance of the PSF method as compared to the HODLR and  
 930 GLR methods for all cases. We note that as we increase the rank and decrease the  
 931 error tolerance, the performance of the HODLR and GLR methods deteriorates.

932 In Figure 14 we show the convergence of the PSF-based method on the blur  
 933 kernel as a function of the total number of impulse responses (left), and the number  
 934 of batches (right). We show convergence for several ellipsoid size parameters  $\tau$ ,  
 935 ranging from 2.0–4.0. The results in Figure 14 (left) show that the relative error  
 936 decreases as  $\text{constant} \times (\#\text{impulse responses})^{-1}$  suggesting linear convergence. The  
 937 linear convergence stalls at a limit that depends on  $\tau$ . Increasing  $\tau$  lowers this limit,  
 938 allowing the PSF-based method to achieve higher accuracy. In Figure 14 (right), the  
 939 results show that before this limit is reached, the convergence is faster for smaller  $\tau$   
 940 in terms of the number of batches. This is expected because smaller  $\tau$  results in more

	Error tol.	#applies PSF	#applies HODLR	#applies RSVD
$L = 1$	20%	11	592	354
	10%	16	772	520
	5%	22	924	674
$L = 1/2$	20%	8	852	1316
	10%	9	1144	1916
	5%	12	1404	2456
$L = 1/3$	20%	7	932	2624
	10%	8	1264	3734
	5%	8	1520	4660

Table 5: (Blur) Comparison of cost to approximate the blur kernel from Equation (7.4) using the PSF-based method, the randomized HODLR (hierarchical off diagonal low rank) method, and GLR (global low rank) approximation using randomized SVD. The quantity  $L$  scales the width of the impulse responses, hence it influences the rank of the operator. Large  $L$  means low rank, and small  $L$  means high rank. The second column (“Error tol”) is the relative error in the approximation of the kernel measured in the Frobenius norm,  $\|\Phi - \tilde{\Phi}\|_{\text{Fro}}/\|\Phi\|_{\text{Fro}}$ . The remaining three columns show the number of operator applies required to achieve the given error tolerances, using the PSF, HODLR, and GLR methods.

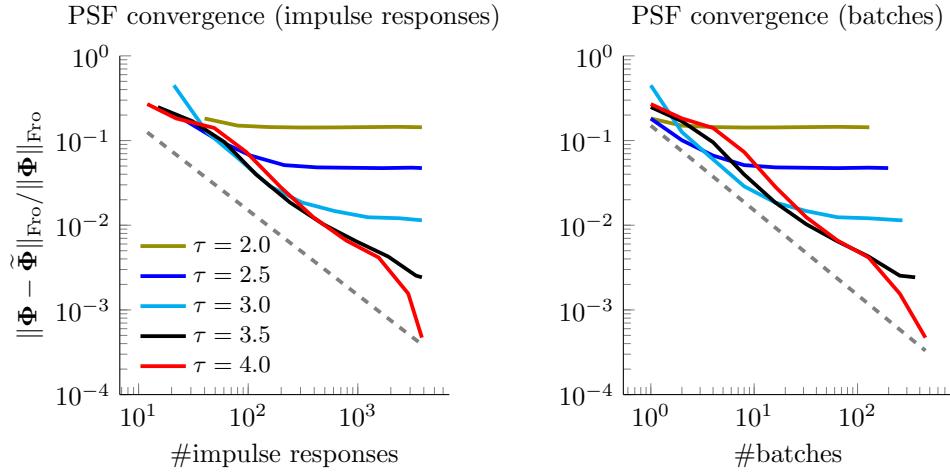


Fig. 14: (Blur) Relative error for different ellipsoid size parameters,  $\tau$ , vs. the total number of impulse responses (left), and the number of batches (right). The dashed gray lines show linear convergence rates, i.e., constant  $\times (\# \text{impulse responses})^{-1}$  on the left, and constant  $\times \# \text{batches}^{-1}$  on the right.

941 impulse responses per batch. Larger  $\tau$  causes the PSF-based method to converge  
 942 more slowly than smaller  $\tau$ , but with larger  $\tau$  the PSF-based method stalls at a lower  
 943 level of error than it does with smaller  $\tau$  (see, e.g.,  $\tau = 4.0$  vs.  $\tau = 2.5$ ).

944 **8. Conclusions.** We presented an efficient matrix-free **PSF-PSF-based** method  
 945 for approximating operators with locally supported non-negative integral kernels. The

946 PSF-based method requires access to the operator only via application of the operator  
 947 to a small number of vectors. The idea of the PSF-based method is to compute  
 948 batches of impulse responses by applying the operator to Dirac combs of scattered  
 949 point sources, then interpolate these impulse responses to approximate entries of the  
 950 operator's integral kernel. The interpolation is based on a new principle we call "local  
 951 mean displacement invariance," which generalizes classical local translation invari-  
 952 ance. The ability to quickly approximate arbitrary integral kernel entries permits us  
 953 to form an H-matrix approximation of the operator. Fast H-matrix arithmetic is then  
 954 used to perform further linear algebra operations that cannot be performed easily with  
 955 the original operator, such as matrix factorization and inversion. The supports of the  
 956 impulse responses are estimated to be contained in ellipsoids, which are determined  
 957 a-priori via a moment method that involves applying the operator to a small number  
 958 of polynomial functions. Point source locations for the impulse response batches are  
 959 chosen using a greedy ellipsoid packing procedure, in which we choose as many im-  
 960 pulse responses per batch as possible, while ensuring that the corresponding ellipsoids  
 961 do not overlap. We applied the PSF-based method to approximate the Gauss-Newton  
 962 Hessians in an ice sheet flow inverse problem governed by a linear Stokes PDE, and  
 963 an advective-diffusive transport inverse problem governed by an advection-diffusion  
 964 PDE. We saw that preconditioners based on the proposed-PSF-based approximation  
 965 cluster the eigenvalues of the preconditioned Hessian near one, and allow us to solve  
 966 the inverse problems using roughly  $5 \times 10 \times$  fewer PDE solves. For larger domains  
 967 with more observations, the rank of the data misfit Hessian will increase, while the  
 968 locality of impulse responses will remain the same. Hence, we expect the speedup will  
 969 be even greater for such problems. Although the method we presented-PSF-based  
 970 method is not applicable to all Hessians, it is applicable to many Hessians of practical  
 971 interest. For these Hessians, the proposed-PSF-based method offers a *data scalable*  
 972 alternative to conventional low rank approximation, due to the ability to form high  
 973 rank approximations of an operator using a small number of operator applications,  
 974 and thus PDE solves.

975 **Acknowledgments.** We thank J.J. Alger, Longfei Gao, Mathew Hu, and Rami  
 976 Nammour for helpful discussions. We thank Trevor Heise, Nicole Aretz, Trevor Heise,  
 977 and the anonymous reviewers for editing suggestions. We thank Georg Stadler for  
 978 help with the domain setup for the ice sheet problem.

- 980 [1] H.-M. ADOLF, *Towards HST restoration with a space-variant PSF, cosmic rays and other*  
 981 *missing data*, in The Restoration of HST Images and Spectra-II, 1994, p. 72.
- 982 [2] V. AKÇELIK, G. BIROS, A. DRĂGĂNESCU, O. GHATTAS, J. HILL, AND B. VAN BLOEMAN WAAN-  
 983 DERS, *Dynamic data-driven inversion for terascale simulations: Real-time identification*  
 984 *of airborne contaminants*, in Proceedings of SC2005, Seattle, 2005.
- 985 [3] A. ALEXANDERIAN, P. J. GLOOR, AND O. GHATTAS, *On Bayesian A-and D-optimal experimen-*  
 986 *tal designs in infinite dimensions*, Bayesian Analysis, 11 (2016), pp. 671–695.
- 987 [4] N. ALGER, *Data-scalable Hessian preconditioning for distributed parameter PDE-constrained*  
 988 *inverse problems*, PhD thesis, The University of Texas at Austin, 2019.
- 989 [5] N. ALGER, V. RAO, A. MYERS, T. BUI-THANH, AND O. GHATTAS, *Scalable matrix-free adap-*  
 990 *tive product-convolution approximation for locally translation-invariant operators*, SIAM  
 991 *Journal on Scientific Computing*, 41 (2019), pp. A2296–A2328.
- 992 [6] N. ALGER, U. VILLA, T. BUI-THANH, AND O. GHATTAS, *A data scalable augmented Lagrangian*  
 993 *KKT preconditioner for large-scale inverse problems*, SIAM Journal on Scientific Comput-  
 994 *ing*, 39 (2017), pp. A2365–A2393.
- 995 [7] I. AMBARTSUMYAN, W. BOUKARAM, T. BUI-THANH, O. GHATTAS, D. KEYES, G. STADLER,

- 996 G. TURKIYYAH, AND S. ZAMPINI, *Hierarchical matrix approximations of Hessians arising*  
 997 *in inverse problems governed by PDEs*, SIAM Journal on Scientific Computing, 42 (2020),  
 998 pp. A3397–A3426.
- 999 [8] T. ARBOGAST AND J. L. BONA, *Methods of Applied Mathematics*, University of Texas at Austin,  
 1000 2008. Lecture notes in applied mathematics.
- 1001 [9] H. T. BANKS AND K. KUNISCH, *Estimation techniques for distributed parameter systems*,  
 1002 Springer, 1989.
- 1003 [10] M. BEBENDORF AND S. RIASANOW, *Adaptive low-rank approximation of collocation matrices*,  
 1004 Computing, 70 (2003), pp. 1–24.
- 1005 [11] J. L. BENTLEY, *Multidimensional binary search trees used for associative searching*, Communications  
 1006 of the ACM, 18 (1975), pp. 509–517.
- 1007 [12] J. BIGOT, P. ESCANDE, AND P. WEISS, *Estimation of linear operators from scattered impulse  
 1008 responses*, Applied and Computational Harmonic Analysis, 47 (2019), pp. 730–758.
- 1009 [13] S. BÖRM, L. GRASEDYCK, AND W. HACKBUSCH, *Introduction to hierarchical matrices with  
 1010 applications*, Engineering analysis with boundary elements, 27 (2003), pp. 405–422.
- 1011 [14] A. BORZÌ AND V. SCHULZ, *Computational optimization of systems governed by partial differ-  
 1012 ential equations*, SIAM, 2011.
- 1013 [15] T. BUI-THANH, O. GHATTAS, J. MARTIN, AND G. STADLER, *A computational framework for  
 1014 infinite-dimensional Bayesian inverse problems Part I: The linearized case, with appli-  
 1015 cation to global seismic inversion*, SIAM Journal on Scientific Computing, 35 (2013),  
 1016 pp. A2494–A2523.
- 1017 [16] T. F. CHAN AND T. P. MATHEW, *Domain decomposition algorithms*, Acta Numerica, 3 (1994),  
 1018 pp. 61–143.
- 1019 [17] J. CHEN AND M. L. STEIN, *Linear-cost covariance functions for Gaussian random fields*, Journal  
 1020 of the American Statistical Association, (2021), pp. 1–18.
- 1021 [18] H. CHENG, Z. GIMBUTAS, P.-G. MARTINSSON, AND V. ROKHLIN, *On the compression of low  
 1022 rank matrices*, SIAM Journal on Scientific Computing, 26 (2005), pp. 1389–1404.
- 1023 [19] T. CUI, J. MARTIN, Y. M. MARZOUK, A. SOLONEN, AND A. SPANTINI, *Likelihood-informed di-  
 1024 mension reduction for nonlinear inverse problems*, Inverse Problems, 30 (2014), p. 114015.
- 1025 [20] J. C. DE LOS REYES, *Numerical PDE-constrained optimization*, Springer, 2015.
- 1026 [21] L. DEMANET, P.-D. LÉTOURNEAU, N. BOUMAL, H. CALANDRA, J. CHIU, AND S. SNELSON,  
 1027 *Matrix probing: a randomized preconditioner for the wave-equation Hessian*, Applied and  
 1028 Computational Harmonic Analysis, 32 (2012), pp. 155–168.
- 1029 [22] L. DENIS, E. THIÉBAUT, AND F. SOULEZ, *Fast model of space-variant blurring and its ap-  
 1030 plication to deconvolution in astronomy*, in Image Processing (ICIP), 2011 18th IEEE  
 1031 International Conference on, IEEE, 2011, pp. 2817–2820.
- 1032 [23] L. DENIS, E. THIÉBAUT, F. SOULEZ, J.-M. BECKER, AND R. MOURYA, *Fast approximations of  
 1033 shift-variant blur*, International Journal of Computer Vision, 115 (2015), pp. 253–278.
- 1034 [24] S. C. EISENSTAT AND H. F. WALKER, *Choosing the forcing terms in an inexact newton method*,  
 1035 SIAM Journal on Scientific Computing, 17 (1996), pp. 16–32.
- 1036 [25] H. W. ENGL, M. HANKE, AND A. NEUBAUER, *Regularization of Inverse Problems*, vol. 375 of  
 1037 Mathematics and Its Applications, Springer Netherlands, 1996.
- 1038 [26] C. ERICSON, *Real-time collision detection*, Crc Press, 2004.
- 1039 [27] P. ESCANDE AND P. WEISS, *Sparse wavelet representations of spatially varying blurring op-  
 1040 erators*, SIAM Journal on Imaging Sciences, 8 (2015), pp. 2976–3014.
- 1041 [28] P. ESCANDE AND P. WEISS, *Approximation of integral operators using product-convolution  
 1042 expansions*, Journal of Mathematical Imaging and Vision, 58 (2017), pp. 333–348.
- 1043 [29] P. ESCANDE AND P. WEISS, *Fast wavelet decomposition of linear operators through product-  
 1044 convolution expansions*, IMA Journal of Numerical Analysis, 42 (2022), pp. 569–596.
- 1045 [30] P. ESCANDE, P. WEISS, AND F. MALGOYRES, *Spatially varying blur recovery. Diagonal ap-  
 1046 proximations in the wavelet domain*, in Proceedings of ICPRAM, 2013.
- 1047 [31] A. FICHTNER AND T. v. LEEUWEN, *Resolution analysis by random probing*, Journal of Geo-  
 1048 physical Research: Solid Earth, 120 (2015), pp. 5549–5573.
- 1049 [32] D. FISH, J. GROCHMALICKI, AND E. PIKE, *Scanning singular-value-decomposition method for  
 1050 restoration of images with space-variant blur*, JOSA A, 13 (1996), pp. 464–469.
- 1051 [33] H. P. FLATH, L. C. WILCOX, V. AKÇELIK, J. HILL, B. VAN BLOEMEN WAANDERS, AND O. GHAT-  
 1052 TAS, *Fast algorithms for Bayesian uncertainty quantification in large-scale linear inverse  
 1053 problems based on low-rank partial Hessian approximations*, SIAM Journal on Scientific  
 1054 Computing, 33 (2011), pp. 407–432.
- 1055 [34] P. H. FLATH, *Hessian-based response surface approximations for uncertainty quantification in  
 1056 large-scale statistical inverse problems, with applications to groundwater flow*, PhD thesis,  
 1057 The University of Texas at Austin, 2013.

- 1058 [35] M. GENTILE, F. COURBIN, AND G. MEYLAN, *Interpolating point spread function anisotropy*,  
1059 *Astronomy & Astrophysics*, 549 (2013).
- 1060 [36] M. G. GENTON, D. E. KEYES, AND G. TURKIYYAH, *Hierarchical decompositions for the com-*  
1061 *putation of high-dimensional multivariate normal probabilities*, *Journal of Computational*  
1062 *and Graphical Statistics*, 27 (2018), pp. 268–277.
- 1063 [37] C. J. GEOGA, M. ANITESCU, AND M. L. STEIN, *Scalable Gaussian process computations us-*  
1064 *ing hierarchical matrices*, *Journal of Computational and Graphical Statistics*, 29 (2020),  
1065 pp. 227–237.
- 1066 [38] O. GHATTAS AND K. WILLCOX, *Learning physics-based models from data: perspectives from*  
1067 *inverse problems and model reduction*, *Acta Numerica*, 30 (2021), pp. 445–554.
- 1068 [39] I. GILITSCHENSKI AND U. D. HANEBECK, *A robust computational test for overlap of two*  
1069 *arbitrary-dimensional ellipsoids in fault-detection of Kalman filters*, in 2012 15th Inter-  
1070 *national Conference on Information Fusion*, IEEE, 2012, pp. 396–401.
- 1071 [40] L. GRASEDYCK AND W. HACKBUSCH, *Construction and arithmetics of H-matrices*, *Computing*,  
1072 70 (2003), pp. 295–334.
- 1073 [41] L. GRASEDYCK, R. KRIEMANN, AND S. LE BORNE, *Parallel black box  $\mathcal{H}$ -LU preconditioning*  
1074 *for elliptic boundary value problems*, *Computing and visualization in science*, 11 (2008),  
1075 pp. 273–291.
- 1076 [42] W. HACKBUSCH, *A sparse matrix arithmetic based on H-matrices. Part I: Introduction to H-*  
1077 *matrices*, *Computing*, 62 (1999), pp. 89–108.
- 1078 [43] W. HACKBUSCH, *Hierarchical matrices: algorithms and analysis*, vol. 49, Springer, 2015.
- 1079 [44] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probab-*  
1080 *ilistic algorithms for constructing approximate matrix decompositions*, *SIAM review*, 53  
1081 (2011), pp. 217–288.
- 1082 [45] T. HARTLAND, G. STADLER, M. PEREGO, K. LIEGEOIS, AND N. PETRA, *Hierarchical off-diagonal*  
1083 *low-rank approximation of hessians in inverse problems, with application to ice sheet model*  
1084 *initialization*, *Inverse Problems*, 39 (2023), p. 085006.
- 1085 [46] F. J. HERRMANN, P. MOGHADDAM, AND C. C. STOLK, *Sparsity-and continuity-promoting seis-*  
1086 *mic image recovery with curvelet frames*, *Applied and Computational Harmonic Analysis*,  
1087 24 (2008), pp. 150–173.
- 1088 [47] M. HINZE, R. PINNAU, M. ULRICH, AND S. ULRICH, *Optimization with PDE constraints*,  
1089 vol. 23, Springer Science & Business Media, 2008.
- 1090 [48] T. ISAAC, N. PETRA, G. STADLER, AND O. GHATTAS, *Scalable and efficient algorithms for*  
1091 *the propagation of uncertainty from data through inference to prediction for large-scale*  
1092 *problems, with application to flow of the Antarctic ice sheet*, *Journal of Computational*  
1093 *Physics*, 296 (2015), pp. 348–368.
- 1094 [49] J. KAPIO AND E. SOMERSALO, *Statistical and Computational Inverse Problems*, vol. 160 of  
1095 *Applied Mathematical Sciences*, Springer-Verlag New York, 2005.
- 1096 [50] K.-T. KIM, U. VILLA, M. PARNO, Y. MARZOUK, O. GHATTAS, AND N. PETRA, *hIPPYlib-MUQ:*  
1097 *A Bayesian inference software framework for integration of data with complex predictive*  
1098 *models under uncertainty*, *ACM Transactions on Mathematical Software*, (2023).
- 1099 [51] R. KRIEMANN, *HLIBpro user manual*, Max-Planck-Institute for Mathematics in the Sciences,  
1100 Leipzig, 48 (2008).
- 1101 [52] R. KRIEMANN,  *$\mathcal{H}$ -LU factorization on many-core systems*, *Computing and Visualization in*  
1102 *Science*, 16 (2013), pp. 105–117.
- 1103 [53] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK users' guide: solution of large-scale*  
1104 *eigenvalue problems with implicitly restarted Arnoldi methods*, vol. 6, SIAM, 1998.
- 1105 [54] J. LEVITT AND P.-G. MARTINSSON, *Linear-complexity black-box randomized compression of*  
1106 *hierarchically block separable matrices*, arXiv preprint arXiv:2205.02990, (2022).
- 1107 [55] L. LIN, J. LU, AND L. YING, *Fast construction of hierarchical matrix representation from*  
1108 *matrix–vector multiplication*, *Journal of Computational Physics*, 230 (2011), pp. 4071–  
1109 4087.
- 1110 [56] F. LINDGREN, H. RUE, AND J. LINDSTRÖM, *An explicit link between Gaussian fields and Gauss-*  
1111 *ian Markov random fields: the stochastic partial differential equation approach*, *Journal of*  
1112 *the Royal Statistical Society: Series B (Statistical Methodology)*, 73 (2011), pp. 423–498.
- 1113 [57] P.-G. MARTINSSON, *A fast randomized algorithm for computing a hierarchically semiseparable*  
1114 *representation of a matrix*, *SIAM Journal on Matrix Analysis and Applications*, 32 (2011),  
1115 pp. 1251–1274.
- 1116 [58] P.-G. MARTINSSON, *Compressing rank-structured matrices via randomized sampling*, *SIAM*  
1117 *Journal on Scientific Computing*, 38 (2016), pp. A1959–A1986.
- 1118 [59] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Foundations and*  
1119 *algorithms*, *Acta Numerica*, 29 (2020), pp. 403–572. Section 20.

- 1120 [60] J. G. NAGY AND D. P. O'LEARY, *Restoring images degraded by spatially variant blur*, SIAM  
1121 Journal on Scientific Computing, 19 (1998), pp. 1063–1082.
- 1122 [61] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, 1999.
- 1123 [62] N. PETRA, J. MARTIN, G. STADLER, AND O. GHATTAS, *A computational framework for infinite-  
1124 dimensional Bayesian inverse problems, Part II: Stochastic Newton MCMC with applica-  
1125 tion to ice sheet flow inverse problems*, SIAM Journal on Scientific Computing, 36 (2014),  
1126 pp. A1525–A1555.
- 1127 [63] N. PETRA AND G. STADLER, *Model variational inverse problems governed by partial differential  
1128 equations*, Tech. Report 11-05, The Institute for Computational Engineering and Sciences,  
1129 The University of Texas at Austin, 2011.
- 1130 [64] N. PETRA, H. ZHU, G. STADLER, T. HUGHES, AND O. GHATTAS, *An inexact Gauss-Newton  
1131 method for inversion of basal sliding and rheology parameters in a nonlinear Stokes ice  
1132 sheet model*, Journal of Glaciology, 58 (2012), pp. 889–903, doi:10.3189/2012JoG11J182.
- 1133 [65] L. ROININEN, J. M. HUTTUNEN, AND S. LASANEN, *Whittle-Matérn priors for Bayesian statis-  
1134 tical inversion with applications in electrical impedance tomography*, Inverse Problems &  
1135 Imaging, 8 (2014), p. 561.
- 1136 [66] Y. SAAD, *Iterative methods for sparse linear systems*, Society for Industrial and Applied Math-  
1137 ematics, Philadelphia, PA, second ed., 2003.
- 1138 [67] B. SMITH, P. BJORSTAD, AND W. GROPP, *Domain Decomposition Parallel Multilevel Methods  
1139 for Elliptic Partial Differential Equations*, Cambridge University Press, 1996.
- 1140 [68] A. SPANTINI, A. SOLONEN, T. CUI, J. MARTIN, L. TENORIO, AND Y. MARZOUK, *Optimal low-  
1141 rank approximations of Bayesian linear inverse problems*, SIAM Journal on Scientific Com-  
1142 puting, 37 (2015), pp. A2451–A2487.
- 1143 [69] T. STEIHAUG, *Local and superlinear convergence for truncated iterated projections methods*,  
1144 Mathematical Programming, 27 (1983), pp. 176–190.
- 1145 [70] A. STUART, *Inverse problems: a Bayesian perspective*, Acta Numerica, 19 (2010), pp. 451–559.
- 1146 [71] W. W. SYMES, *Approximate linearized inversion by optimal scaling of prestack depth migration*,  
1147 Geophysics, 73 (2008), pp. R23–R35.
- 1148 [72] A. TARANTOLA, *Inverse Problem Theory and Methods for Model Parameter Estimation*, SIAM,  
1149 Philadelphia, PA, 2005.
- 1150 [73] A. TOSELLI AND O. WIDLUND, *Domain decomposition methods-algorithms and theory*, vol. 34,  
1151 Springer Science & Business Media, 2004.
- 1152 [74] J. TRAMPERT, A. FICHTNER, AND J. RITSEMA, *Resolution tests revisited: the power of random  
1153 numbers*, Geophysical Journal International, 192 (2013), pp. 676–680.
- 1154 [75] U. VILLA, N. PETRA, AND O. GHATTAS, *hIPPYlib: an extensible software framework for large-  
1155 scale inverse problems governed by PDEs: Part I: deterministic inversion and linearized  
1156 Bayesian inference*, ACM Transactions on Mathematical Software, 47 (2021), pp. 1–34.
- 1157 [76] C. R. VOGEL, *Computational Methods for Inverse Problems*, Frontiers in Applied Mathematics,  
1158 Society for Industrial and Applied Mathematics (SIAM), Philadelphia, PA, 2002.
- 1159 [77] H. WENDLAND, *Scattered data approximation*, vol. 17, Cambridge University Press, 2004.
- 1160 [78] H. ZHU, S. LI, S. FOMEL, G. STADLER, AND O. GHATTAS, *A Bayesian approach to estimate  
1161 uncertainty for full waveform inversion with a priori information from depth migration*,  
1162 Geophysics, 81 (2016), pp. R307–R323.