

LINEAR-COMPLEXITY BLACK-BOX RANDOMIZED COMPRESSION OF RANK-STRUCTURED MATRICES *

JAMES LEVITT[†] AND PER-GUNNAR MARTINSSON[†]

Abstract. A randomized algorithm for computing a compressed representation of a given rank-structured matrix $\mathbf{A} \in \mathbb{R}^{N \times N}$ is presented. The algorithm interacts with \mathbf{A} only through its action on vectors. Specifically, it draws two tall thin matrices $\mathbf{\Omega}, \mathbf{\Psi} \in \mathbb{R}^{N \times s}$ from a suitable distribution, and then reconstructs \mathbf{A} from the information contained in the set $\{\mathbf{A}\mathbf{\Omega}, \mathbf{\Omega}, \mathbf{A}^*\mathbf{\Psi}, \mathbf{\Psi}\}$. For the specific case of a “Hierarchically Block Separable (HBS)” matrix (a.k.a. Hierarchically Semi-Separable matrix) of block rank k , the number of samples s required satisfies $s = O(k)$, with $s \approx 3k$ being representative. While a number of randomized algorithms for compressing rank-structured matrices have previously been published, the current algorithm appears to be the first that is both of truly linear complexity (no $N \log(N)$ factors in the complexity bound) and fully “black box” in the sense that no matrix entry evaluation is required. Further, all samples can be extracted in parallel, enabling the algorithm to work in a “streaming” or “single view” mode.

Key words. randomized approximation of matrices; rank-structured matrices; HODLR matrix; hierarchically block separable matrix; hierarchically semiseparable matrix; randomized SVD; fast direct solver.

AMS subject classifications. 65N22, 65N38, 15A23, 15A52

1. Introduction. This work describes an efficient algorithm for handling large dense matrices that have *rank structure*. To simplify slightly, this means that an $N \times N$ matrix can be tessellated into $O(N)$ blocks in such a way that each block is either small or of low numerical rank, cf. [Figure 2](#). This structure allows the matrix to be stored and applied to vectors efficiently, often with cost that scales linearly or close to linearly with N [[12](#), [3](#), [2](#), [23](#)]. Sometimes, it is also possible to compute an approximate inverse or LU factorization in linear or close to linear time [[27](#)]. Matrices of this type have turned out to be ubiquitous in both engineering and data sciences, and have been the subject of much research in recent decades, going under names such as \mathcal{H} -matrices [[2](#), [3](#), [12](#)]; HODLR matrices [[1](#), [20](#)], Hierarchically Block Separable (HBS) or Hierarchically Semi-Separable (HSS) matrices [[4](#), [5](#), [28](#)], Recursive Skeletonization [[8](#), [14](#), [19](#), [26](#)], and many more.

The specific problem we address is the following: Suppose that \mathbf{A} is an $N \times N$ matrix that we know has HBS structure, but we do not have direct access to the low-rank factors that define the compressible off-diagonal blocks. Instead, we have access to fast “black-box” algorithms that given tall thin matrices $\mathbf{\Omega}, \mathbf{\Psi} \in \mathbb{R}^{N \times s}$, evaluate the matrix-matrix products

$$\mathbf{Y} = \mathbf{A}\mathbf{\Omega} \quad \text{and} \quad \mathbf{Z} = \mathbf{A}^*\mathbf{\Psi}.$$

The problem is then to recover \mathbf{A} from the information in the set $\{\mathbf{Y}, \mathbf{\Omega}, \mathbf{Z}, \mathbf{\Psi}\}$. The algorithm described here solves the reconstruction problem using $s = O(k)$ sample vectors, where k is an upper bound on the ranks of the off-diagonal blocks that we assume is known in advance. (The scaling factor hidden in the formula $s = O(k)$ is modest, with $s \approx 3k$ being representative.)

The scheme presented has several important applications. First, it can be used to derive a rank-structured representation of any integral operator for which a fast

*Submitted to the editors DATE.

[†] Oden Institute for Computational Engineering and Sciences, University of Texas at Austin, Austin, TX (jlevitt@utexas.edu, pgm@oden.utexas.edu).

matrix-vector multiplication algorithm, such as the Fast Multipole Method [10, 11], is available. The new representation that we compute opens the door to a wider range of matrix operations such as LU factorization, matrix inversion, and sometimes even full spectral decompositions. Second, it can greatly simplify algebraic operations involving products of rank-structured matrices. For instance, the perhaps key application of rank-structured matrix algebra is the acceleration of sparse direct solvers, as the dense matrices that arise during LU factorization are often rank structured. In the course of such a solver, a typical operation would be to form a Schur complement such as $\mathbf{S}_{22} = \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$ that would arise when the top left block \mathbf{A}_{11} is eliminated from a 2×2 blocked matrix. If \mathbf{A}_{11} is rank structured, then \mathbf{A}_{11}^{-1} can easily be applied to vectors via an LU factorization. If, additionally, \mathbf{A}_{12} and \mathbf{A}_{21} are either sparse or rank structured, then \mathbf{S}_{22} can easily be applied to a vector. The technique described will then enable one to construct a data-sparse representation of \mathbf{S}_{22} . In contrast, to directly evaluate the product $\mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}$ is both onerous to code and slow to execute.

The method we describe is inspired by the randomized compression algorithm introduced in [18, 21] for compressing HBS matrices. That algorithm is similar to the one presented here in that it requires only $O(k)$ matrix-vector applications and $O(k^2N)$ floating point operations. However, the algorithm of [18, 21] is not a true black-box algorithm since, in addition to randomized samples, it also requires direct evaluation of a small number of entries of the matrix. In contrast, the method presented here is truly black box.

Remark 1.1 (Peeling algorithms). A related class of algorithms for randomized compression of rank-structured matrices is described in [17, 22], and a recent improvement is described in [16]. These techniques are “true” black-box algorithms in that they only access the matrix through the black-box matrix-vector multiplication routines, and they also apply to a broader class of rank-structured matrices than the one considered here. However, they require $O(k \log(N))$ samples and $O(k^2 N \log(N))$ floating point operations, so they do not have linear complexity.

The manuscript is structured as follows: Section 2 surveys some basic linear algebraic techniques that we rely on. Section 3 introduces our formalism for HBS matrices. Section 4 describes the new algorithm, and analyzes its asymptotic complexity. Section 5 describes numerical results.

2. Preliminaries. In introducing well-known material, we follow the presentation of [22].

2.1. Notation. Throughout the paper, we measure a vector $\mathbf{x} \in \mathbb{R}^n$ by its Euclidean norm $\|\mathbf{x}\| = (\sum_i |x_i|^2)^{\frac{1}{2}}$. We measure a matrix $\mathbf{A} \in \mathbb{R}^{m \times n}$ with the corresponding operator norm $\|\mathbf{A}\| = \sup_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|$, and in some cases with the Frobenius norm $\|\mathbf{A}\|_{\text{Fro}} = (\sum_{i,j} |\mathbf{A}(i,j)|^2)^{1/2}$. To denote submatrices, we use the notation of Golub and Van Loan [9]: If \mathbf{A} is an $m \times n$ matrix, and $I = [i_1, i_2, \dots, i_k]$ and $J = [j_1, j_2, \dots, j_l]$, then $\mathbf{A}(I, J)$ denotes the $k \times l$ matrix

$$\mathbf{A}(I, J) = \begin{bmatrix} \mathbf{A}(i_1, j_1) & \mathbf{A}(i_1, j_2) & \dots & \mathbf{A}(i_1, j_l) \\ \mathbf{A}(i_2, j_1) & \mathbf{A}(i_2, j_2) & \dots & \mathbf{A}(i_2, j_l) \\ \vdots & \vdots & & \vdots \\ \mathbf{A}(i_k, j_1) & \mathbf{A}(i_k, j_2) & \dots & \mathbf{A}(i_k, j_l) \end{bmatrix}$$

We let $\mathbf{A}(I, :)$ denote the submatrix $\mathbf{A}(I, [1, 2, \dots, n])$ and define $\mathbf{A}(:, J)$ analogously. We let \mathbf{A}^* denote the transpose of \mathbf{A} , and we say that matrix \mathbf{U} is *orthonormal* if its

columns are orthonormal, $\mathbf{U}^*\mathbf{U} = \mathbf{I}$.

2.2. The QR factorization. The full QR factorization of a matrix \mathbf{B} of size $m \times n$ takes the form

$$(2.1) \quad \begin{array}{ccc} \mathbf{B} & = & \mathbf{Q} \quad \mathbf{R}, \\ m \times n & & m \times m \quad m \times n \end{array}$$

where \mathbf{Q} is unitary and \mathbf{R} is upper-triangular.

If we further assume that \mathbf{B} has rank k and that its first k columns are linearly independent, then it has a rank- k partial QR factorization given by

$$\begin{array}{ccc} \mathbf{B}_k & = & \mathbf{Q}_k \quad \mathbf{R}_k, \\ m \times n & & m \times k \quad k \times n \end{array}$$

where \mathbf{Q}_k has orthonormal columns and \mathbf{R}_k is upper-triangular.

2.3. Randomized compression. Let \mathbf{B} be an $m \times n$ matrix that can be accurately approximated by a matrix of rank k , and suppose we seek to determine a matrix \mathbf{Q} with orthonormal columns (as few as possible) such that $\|\mathbf{B} - \mathbf{Q}\mathbf{Q}^*\mathbf{B}\|$ is small. In other words, we seek a matrix \mathbf{Q} whose columns form an approximate orthonormal basis (ON-basis) for the column space of \mathbf{B} . This task can efficiently be solved via the following randomized procedure:

1. Pick a small integer p representing how much ‘‘oversampling’’ is done ($p = 10$ is often good), and define $r = k + p$.
2. Form an $n \times r$ matrix \mathbf{G} whose entries are independent and identically distributed (i.i.d.) normalized Gaussian random numbers.
3. Form the ‘‘sample matrix’’ $\mathbf{Y} = \mathbf{B}\mathbf{G}$ of size $m \times r$.
4. Construct an $m \times r$ matrix \mathbf{Q} whose columns form an ON basis for the columns of \mathbf{Y} .

Note that each column of the sample matrix \mathbf{Y} is a random linear combination of the columns of \mathbf{B} . The probability of the algorithm producing an accurate result approaches 1 extremely rapidly as p increases, and, remarkably, this probability depends only on p (not on m or n , or any other properties of \mathbf{B}); cf. [13].

2.4. Functions for computing orthonormal bases. For matrix \mathbf{B} of rank at most k , we write

$$\mathbf{Q} = \text{col}(\mathbf{B}, k)$$

for a function call that returns a matrix \mathbf{Q} whose k columns are orthonormal and span the column space of \mathbf{B} . In practice, we implement this function using a truncated QR factorization.

For matrix \mathbf{B} with nullspace of dimension k or greater, we write

$$\mathbf{Z} = \text{null}(\mathbf{B}, k)$$

for a function call that returns matrix \mathbf{Z} whose k columns are orthonormal and belong to the nullspace of \mathbf{B} . In practice, we implement this function by taking the last k columns of the factor \mathbf{Q} of the full QR factorization of \mathbf{B}^* .

Remark 2.1. For a matrix \mathbf{B} whose first k columns have rank less than $\text{rank}(\mathbf{B})$, the first k columns of \mathbf{Q} produced by an unpivoted QR factorization algorithm might not span the column space of \mathbf{B} . A similar concern about linear dependence of rows

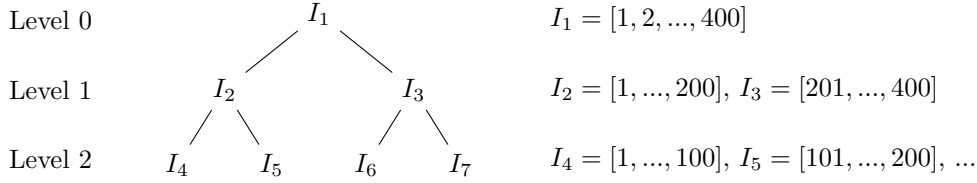


Fig. 1: A binary tree structure, where the levels of the tree represent successively refined partitions of the index vector $[1, 2, \dots, 400]$.

exists for the task of computing a basis belonging to the nullspace. However, in this work we only apply `col` and `null` to random matrices or to products involving random matrices, so any subset of k rows or columns of those matrices will with probability 1 have rank of $\text{rank}(\mathbf{B})$ as long as $\text{rank}(\mathbf{B}) \leq k$. Therefore, `col` and `null` can safely rely on unpivoted QR factorizations. (Observe that unpivoted QR is numerically stable in this context, and will guarantee that $\|\mathbf{B} - \mathbf{Q}\mathbf{Q}^*\mathbf{B}\|$ is small.)

3. Hierarchically block separable matrices. We start this section with a review of important concepts relating to HBS matrices, following the presentation of [24]. While [Subsections 3.1](#) and [3.2](#) review established material [Subsection 3.3](#) introduces a new telescoping factorization of HBS matrices that differs from what is used in prior literature. We encourage the reader to review section 3.3 carefully, as the new telescoping factorization is a key idea leading to the black-box algorithm of [section 4](#).

3.1. A tree structure. Let $I = [1, 2, \dots, N]$ be a vector of indices corresponding to the rows and columns of an $N \times N$ matrix. We define a tree \mathcal{T} , in which each node τ is associated with a contiguous subset of the indices I_τ . To the root node of the tree, we assign the full set of indices I . The two children of the root node are given index vectors $[1, \dots, \lceil N/2 \rceil]$ and $[\lceil N/2 \rceil + 1, \dots, N]$. We continue evenly splitting the indices of each node to form two child nodes until we reach a level of the tree in which the size of each node is below some given threshold m . We refer to a node with children as a parent node, and a node with no children as a leaf node. The depth of a node is defined as its distance from the root node, and level ℓ of the tree is defined as the set of nodes with depth ℓ , so that level 0 consists of only the root node, level 1 consists of the two children of the root node, and so on. The levels of the tree represent successively finer partitions of I . The depth of the tree is defined as the maximum node depth, denoted by $L \approx \log_2(N/m)$. For simplicity, we only consider fully populated binary trees. An example tree structure is depicted in [Figure 1](#).

3.2. The HBS matrix format. Let \mathcal{T} be a tree defined on index vector $I = [1, 2, \dots, N]$. Matrix \mathbf{A} of size $N \times N$ is said to be hierarchically block separable with block rank k with respect to \mathcal{T} if the following conditions are satisfied.

(1) *Assumption on the ranks of off-diagonal blocks of the finest level:* For every pair of leaf nodes τ and τ' , we define

$$\mathbf{A}_{\tau, \tau'} = \mathbf{A}(I_\tau, I_{\tau'})$$

and require that if τ and τ' are distinct, then $\mathbf{A}_{\tau, \tau'}$ must have rank of at most k . Additionally, for each leaf node τ there must exist basis matrices \mathbf{U}_τ and \mathbf{V}_τ such

that for every leaf node $\tau' \neq \tau$ we have

$$(3.1) \quad \begin{array}{ccccc} \mathbf{A}_{\tau,\tau'} & = & \mathbf{U}_\tau & \tilde{\mathbf{A}}_{\tau,\tau'} & \mathbf{V}_{\tau'}^* \\ m \times m & & m \times k & k \times k & k \times m \end{array}$$

(2) *Assumption on the ranks of off-diagonal blocks of levels $L-1, L-2, \dots, 1$:* The following conditions must hold for each level $\ell = L-1, L-2, \dots, 1$. For every pair of distinct nodes τ and τ' on level ℓ with children (α, β) and (α', β') , respectively, we define

$$\mathbf{A}_{\tau,\tau'} = \begin{bmatrix} \tilde{\mathbf{A}}_{\alpha,\alpha'} & \tilde{\mathbf{A}}_{\alpha,\beta'} \\ \tilde{\mathbf{A}}_{\beta,\alpha'} & \tilde{\mathbf{A}}_{\beta,\beta'} \end{bmatrix}$$

and require that every such matrix have rank of at most k . Additionally, for each node τ on level ℓ there must exist basis matrices \mathbf{U}_τ and \mathbf{V}_τ such that for every node $\tau' \neq \tau$ on level ℓ , we have

$$\begin{array}{ccccc} \mathbf{A}_{\tau,\tau'} & = & \mathbf{U}_\tau & \tilde{\mathbf{A}}_{\tau,\tau'} & \mathbf{V}_{\tau'}^* \\ 2k \times 2k & & 2k \times k & k \times k & k \times 2k \end{array}$$

Notably, no assumptions are made on the ranks of the on-diagonal blocks of level L , and those blocks may have full rank. An example tessellation of an HBS matrix showing compressible and incompressible blocks is given in [Figure 2](#).

3.3. Telescoping factorizations. We define the following block-diagonal basis matrices.

$$\begin{aligned} \mathbf{U}^{(\ell)} &= \text{diag}(\mathbf{U}_\tau : \tau \text{ is a node on level } \ell), & \ell = 1, 2, \dots, L \\ \mathbf{V}^{(\ell)} &= \text{diag}(\mathbf{V}_\tau : \tau \text{ is a node on level } \ell), & \ell = 1, 2, \dots, L \end{aligned}$$

Then we obtain a factorization of level L of the form

$$(3.2) \quad \mathbf{A} = \mathbf{U}^{(L)} \tilde{\mathbf{A}}^{(L)} (\mathbf{V}^{(L)})^* + \mathbf{D}^{(L)},$$

where

$$(3.3) \quad \tilde{\mathbf{A}}^{(L)} = (\mathbf{U}^{(L)})^* \mathbf{A} \mathbf{V}^{(L)}$$

$$(3.4) \quad \mathbf{D}^{(L)} = \mathbf{A} - \mathbf{U}^{(L)} \tilde{\mathbf{A}}^{(L)} (\mathbf{V}^{(L)})^*.$$

[Equation \(3.2\)](#) can be viewed as a decomposition of \mathbf{A} into a term that “fits” into the low-rank approximation using basis matrices $\mathbf{U}^{(L)}$ and $\mathbf{V}^{(L)}$ and a discrepancy term $\mathbf{D}^{(L)}$ that does not.

For successively coarser levels $\ell = L-1, L-2, \dots, 1$, we similarly have

$$(3.5) \quad \tilde{\mathbf{A}}^{(\ell+1)} = \mathbf{U}^{(\ell)} \tilde{\mathbf{A}}^{(\ell)} (\mathbf{V}^{(\ell)})^* + \mathbf{D}^{(\ell)},$$

where

$$\begin{aligned} \tilde{\mathbf{A}}^{(\ell)} &= (\mathbf{U}^{(\ell)})^* \tilde{\mathbf{A}}^{(\ell+1)} \mathbf{V}^{(\ell)} \\ \mathbf{D}^{(\ell)} &= \tilde{\mathbf{A}}^{(\ell+1)} - \mathbf{U}^{(\ell)} \tilde{\mathbf{A}}^{(\ell)} (\mathbf{V}^{(\ell)})^*. \end{aligned}$$

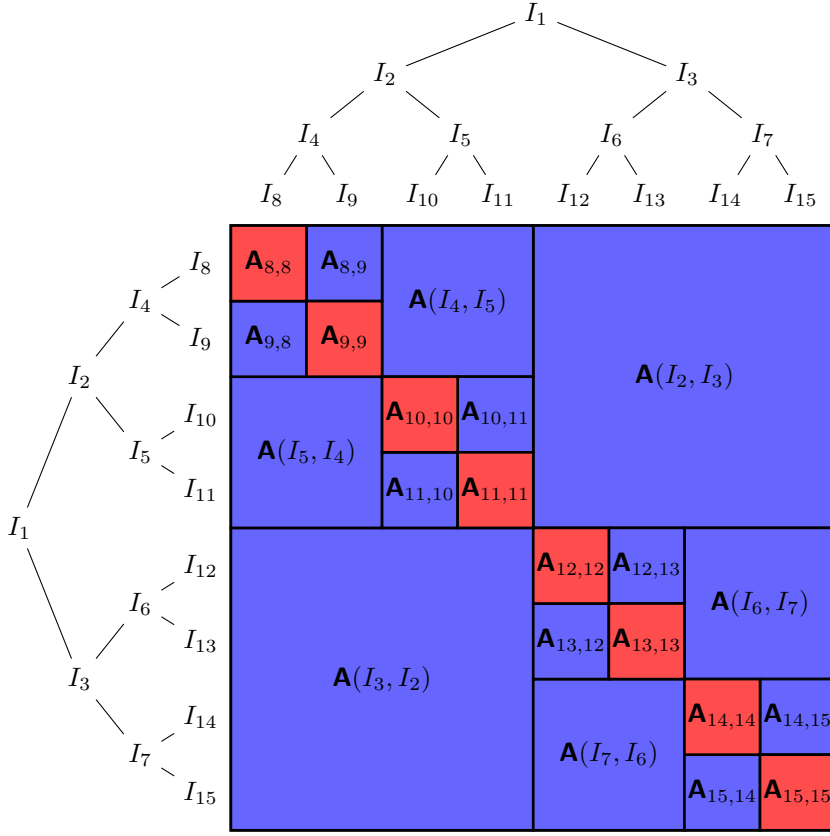


Fig. 2: Tesselation of an HBS matrix with depth 3. Low-rank blocks are shown in blue, and blocks that are not necessarily low-rank are shown in red. The blue blocks are stored in factored form, using “nested” basis matrices as described in section 3.2.

For the root level, we define

$$\mathbf{D}^{(0)} = \tilde{\mathbf{A}}^{(\ell+1)}.$$

Equations (3.2) and (3.5) define a telescoping factorization of \mathbf{A} . For example, a factorization with $L = 2$ takes the form

$$\mathbf{A} = \mathbf{U}^{(2)} \left(\mathbf{U}^{(1)} \mathbf{D}^{(0)} (\mathbf{V}^{(1)})^* + \mathbf{D}^{(1)} (\mathbf{V}^{(2)})^* + \mathbf{D}^{(2)} \right).$$

It follows from (3.1) that matrices $\mathbf{D}^{(\ell)}$, $\ell = 0, 1, \dots, L$, are also block-diagonal. Matrix $\mathbf{D}^{(\ell)}$ can be described in terms of its on-diagonal blocks as

$$\mathbf{D}^{(\ell)} = \text{diag}(\mathbf{D}_\tau : \tau \text{ is a node on level } \ell),$$

where

$$(3.6) \quad \mathbf{D}_\tau = \mathbf{A}_{\tau,\tau} - \mathbf{U}_\tau \mathbf{U}_\tau^* \mathbf{A}_{\tau,\tau} \mathbf{V}_\tau \mathbf{V}_\tau^*.$$

Algorithm 3.1 describes the process of efficiently applying the telescoping factorization to a vector.

Algorithm 3.1 Apply a compressed HBS matrix to a vector: $\mathbf{u} = \mathbf{A}\mathbf{q}$.

Upward pass

for level $\ell = L, L - 1, \dots, 1$ **do**

for node τ in level L **do**

if τ is a leaf node **then**

$$\hat{\mathbf{q}}_\tau = \mathbf{V}_\tau^* \mathbf{q}(I_\tau)$$

else

 Let α and β be the children of τ .

$$\hat{\mathbf{q}}_\tau = \mathbf{V}_\tau^* \begin{bmatrix} \hat{\mathbf{q}}_\alpha \\ \hat{\mathbf{q}}_\beta \end{bmatrix}$$

end if

end for

end for

Downward pass

for levels $\ell = 0, 1, \dots, L$ **do**

if τ is the root node **then**

 Let α and β be the children of τ .

$$\begin{bmatrix} \hat{\mathbf{u}}_\alpha \\ \hat{\mathbf{u}}_\beta \end{bmatrix} = \mathbf{D}_\tau \begin{bmatrix} \hat{\mathbf{q}}_\alpha \\ \hat{\mathbf{q}}_\beta \end{bmatrix}$$

else if τ is a parent node **then**

 Let α and β be the children of τ .

$$\begin{bmatrix} \hat{\mathbf{u}}_\alpha \\ \hat{\mathbf{u}}_\beta \end{bmatrix} = \mathbf{U}_\tau \hat{\mathbf{u}}_\tau + \mathbf{D}_\tau \begin{bmatrix} \hat{\mathbf{q}}_\alpha \\ \hat{\mathbf{q}}_\beta \end{bmatrix}$$

else

$$\mathbf{u}(I_\tau) = \mathbf{U}_\tau \hat{\mathbf{u}}_\tau + \mathbf{D}_\tau \mathbf{q}(I_\tau)$$

end if

end for

Remark 3.1. Let us stress that our definition of the HBS structure is slightly different from the more common framework of, e.g., [21, 23]. Specifically, most authors define the block diagonal matrix $\mathbf{D}^{(L)}$ as the matrix that holds the diagonal subblocks of the original matrix \mathbf{A} , which then leaves $\tilde{\mathbf{A}}^{(L)}$ with zero blocks on the diagonal. In contrast, our definition (3.4) has $\mathbf{D}^{(L)}$ holding the “remainder” of the diagonal blocks after the component that can be spanned by the basis matrices on level L has been peeled off in (3.3). This new definition of $\mathbf{D}^{(L)}$ is essential to our technique for avoiding the need to explicitly form entries of the original matrix.

4. An algorithm for compressing HBS matrices. This section describes the new compression algorithm for HBS matrices that is the main contribution of the manuscript. Let \mathbf{A} be an $N \times N$ HBS matrix with block rank k and leaf node size m , and let $r = k + p$, where p represents a small amount of oversampling ($p = 5$ or

$p = 10$ are often sufficient). Let $\mathbf{\Omega}$ and $\mathbf{\Psi}$ be $N \times s$ Gaussian test matrices, where $s \geq \max(r+m, 3r)$, and define sample matrices $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ and $\mathbf{Z} = \mathbf{A}^*\mathbf{\Psi}$. Our objective is to use the information contained in the test and sample matrices to construct a telescoping factorization of \mathbf{A} , as defined in [subsection 3.3](#).

We begin this section by describing the process of finding the level- L basis matrices $\mathbf{U}^{(L)}$ and $\mathbf{V}^{(L)}$ and the level- L discrepancy matrix $\mathbf{D}^{(L)}$. Next, we describe how to proceed to coarser levels of the tree. Finally, we analyze the asymptotic complexity of the compression algorithm.

4.1. Computing basis matrices $\mathbf{U}^{(L)}$, $\mathbf{V}^{(L)}$ at the finest level. We start by describing the process of computing the basis matrix $\mathbf{U}^{(L)}$ of the finest level, which involves finding for each τ on level L a basis matrix \mathbf{U}_τ that spans the range of $\mathbf{A}(I_\tau, I_{\tau'})$ for every node $\tau' \neq \tau$ on level L . We will compute \mathbf{U}_τ by applying the randomized algorithm described in [subsection 2.3](#) to $\mathbf{A}(I_\tau, I_\tau^c)$, where $I_\tau^c = I \setminus I_\tau$ is the set of indices that are not in I_τ . Importantly, the procedure does not require the ability to apply $\mathbf{A}(I_\tau, I_\tau^c)$ to random vectors; rather we compute the randomized samples using only information contained in $\mathbf{\Omega}$ and \mathbf{Y} .

For the purpose of illustration, suppose for now that \mathbf{A} is an HBS matrix with depth 2, and we want to find the basis matrix \mathbf{U}_4 associated with node 4. We let $\mathbf{\Omega}_4 = \mathbf{\Omega}(I_4, :)$ and $\mathbf{Y}_4 = \mathbf{Y}(I_4, :)$ denote the blocks of $\mathbf{\Omega}$ and \mathbf{Y} of size $m \times s$ associated with node 4. Then the randomized sample of \mathbf{A} takes the following form.

$$\mathbf{Y} = \mathbf{A} \mathbf{\Omega}$$

Since $\mathbf{\Omega}_4$ is of size $m \times s$, it has a nullspace of dimension at least $s - m \geq r$. Then we can find a set of r orthonormal vectors that belong to its nullspace, $\mathbf{P}_4 = \text{null}(\mathbf{\Omega}_4, r)$, so that $\mathbf{\Omega}_4 \mathbf{P}_4 = \mathbf{0}$. Then we have

$$\mathbf{Y} \mathbf{P}_4 = \mathbf{A} \mathbf{\Omega} \mathbf{P}_4$$

Notably, submatrix $(\mathbf{\Omega} \mathbf{P}_4)(I_4, :)$ is filled with zeros, and submatrix $(\mathbf{\Omega} \mathbf{P}_4)(I_4^c, :)$ is filled with values that turn out to also have a standard Gaussian distribution, cf. [Remark 4.1](#). The product $\mathbf{A} \mathbf{\Omega} \mathbf{P}_4$ can be viewed as a randomized sample of \mathbf{A} , excluding contributions from columns $\mathbf{A}(:, I_4)$, so the rows of $\mathbf{A} \mathbf{\Omega} \mathbf{P}_4$ indexed by I_4 contain a randomized sample of $\mathbf{A}(I_4, I_4^c)$. As suggested by the diagram above, we can obtain that sample inexpensively by simply multiplying $\mathbf{Y}_4 \mathbf{P}_4$. Then we orthonormalize the sample to find the basis matrix, $\mathbf{U}_4 = \text{col}(\mathbf{Y}_4 \mathbf{P}_4, r)$.

The generalization to an arbitrary leaf node τ is straightforward. We define

$$\mathbf{\Omega}_\tau = \mathbf{\Omega}(I_\tau, :) \quad \text{and} \quad \mathbf{Y}_\tau = \mathbf{Y}(I_\tau, :).$$

Then we compute the nullspace

$$\mathbf{P}_\tau = \text{null}(\boldsymbol{\Omega}_\tau, r)$$

so that $\mathbf{Y}_\tau \mathbf{P}_\tau$ contains a randomized sample of $\mathbf{A}(I_\tau, I_\tau^c)$. Finally, we orthonormalize the sample to find the basis matrix \mathbf{U}_τ ,

$$(4.1) \quad \mathbf{U}_\tau = \text{col}(\mathbf{Y}_\tau \mathbf{P}_\tau, r).$$

A similar process using $\boldsymbol{\Psi}$ and \mathbf{Z} yields basis matrices \mathbf{V}_τ ,

$$(4.2) \quad \begin{aligned} \mathbf{Q}_\tau &= \text{null}(\boldsymbol{\Psi}_\tau, r) \\ \mathbf{V}_\tau &= \text{col}(\mathbf{Z}_\tau \mathbf{Q}_\tau, r). \end{aligned}$$

Remark 4.1. We claimed that the matrix $(\boldsymbol{\Omega} \mathbf{P}_\tau)(I_\tau^c, :) = \boldsymbol{\Omega}(I_\tau^c, :)\mathbf{P}_\tau$ has a Gaussian distribution. This claim follows from (i) the fact that the distribution of Gaussian matrices is invariant under unitary transformations, and (ii) that the construction of the matrix \mathbf{P}_τ is done independently of the entries in $\boldsymbol{\Omega}(I_\tau^c, :)$.

4.2. Computing discrepancy matrix $\mathbf{D}^{(L)}$ at the finest level. Once we have computed $\mathbf{U}^{(L)}$ and $\mathbf{V}^{(L)}$, we next compute $\mathbf{D}^{(L)}$. Recall that our definition of $\mathbf{D}^{(L)}$, given in (3.6), is different from what appears in the typical telescoping factorization of HBS matrices (cf. Remark 3.1). We proceed by rewriting (3.6) as

$$(4.3) \quad \mathbf{D}_\tau = (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{A}_{\tau, \tau} + \mathbf{U}_\tau \mathbf{U}_\tau^* \mathbf{A}_{\tau, \tau} (\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*),$$

and deriving formulas for computing $(\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{A}_{\tau, \tau}$ and $\mathbf{U}_\tau \mathbf{U}_\tau^* \mathbf{A}_{\tau, \tau} (\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*)$ separately.

For $(\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{A}_{\tau, \tau}$, we first express \mathbf{Y}_τ as a blocked matrix product

$$\mathbf{Y}_\tau = \sum_{\tau' \text{ in level } \ell} \mathbf{A}_{\tau, \tau'} \boldsymbol{\Omega}_{\tau'}.$$

Multiplying $(\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*)$ and applying (3.1) gives

$$(\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{Y}_\tau = (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{A}_{\tau, \tau} \boldsymbol{\Omega}_\tau.$$

Solving a least-squares problem with $\boldsymbol{\Omega}_\tau$ gives

$$(4.4) \quad (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{A}_{\tau, \tau} = (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{Y}_\tau \boldsymbol{\Omega}_\tau^\dagger.$$

A similar derivation yields

$$(4.5) \quad \mathbf{A}_{\tau, \tau} (\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*) = ((\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*) \mathbf{Z}_\tau \boldsymbol{\Psi}_\tau^\dagger)^*.$$

Substituting (4.4) and (4.5) into (4.3) gives the formula

$$(4.6) \quad \mathbf{D}_\tau = (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{Y}_\tau \boldsymbol{\Omega}_\tau^\dagger + \mathbf{U}_\tau \mathbf{U}_\tau^* ((\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*) \mathbf{Z}_\tau \boldsymbol{\Psi}_\tau^\dagger)^*.$$

Remark 4.2. The computation of \mathbf{D}_τ involves solving least squares problems with Gaussian matrices $\boldsymbol{\Omega}_\tau$ and $\boldsymbol{\Psi}_\tau$ of size $m \times s$ for τ belonging to level L or size $2r \times s$ for τ belonging to a coarser level. Gaussian matrices with nearly square shapes have non-negligible probabilities of being ill-conditioned, but the probabilities quickly become negligible even for slightly rectangular matrices [6, 7]. Such concerns can be alleviated by choosing s to be sufficiently large. For the numerical experiments in section 5, we simply use $s = r + m = 3r$.

4.3. Compressing levels $\ell = L - 1, L - 2, \dots, 0$. After compressing level L , we proceed to the next coarser level $L - 1$. That is, we seek $\mathbf{U}^{(L-1)}, \mathbf{V}^{(L-1)}$, and $\mathbf{D}^{(L-1)}$ that satisfy (3.5). We will first obtain randomized samples of $\tilde{\mathbf{A}}^{(L)}$, and then using the same procedure as for level L we will find $\mathbf{U}^{(L-1)}, \mathbf{V}^{(L-1)}$, and $\mathbf{D}^{(L-1)}$.

To compute randomized samples of $\tilde{\mathbf{A}}^{(L)}$, we multiply (3.2) with $\mathbf{\Omega}$ to obtain

$$\mathbf{Y} = \mathbf{A}\mathbf{\Omega} = (\mathbf{U}^{(L)}\tilde{\mathbf{A}}^{(L)}(\mathbf{V}^{(L)})^* + \mathbf{D}^{(L)})\mathbf{\Omega},$$

and rearrange to obtain

$$(4.7) \quad \underbrace{(\mathbf{U}^{(L)})^*(\mathbf{Y} - \mathbf{D}^{(L)}\mathbf{\Omega})}_{\text{sample matrix}} = \tilde{\mathbf{A}}^{(L)} \underbrace{(\mathbf{V}^{(L)})^*\mathbf{\Omega}}_{\text{test matrix}}.$$

Then the columns of $(\mathbf{U}^{(L)})^*(\mathbf{Y} - \mathbf{D}^{(L)}\mathbf{\Omega})$ contain s samples of $\tilde{\mathbf{A}}^{(L)}$ taken with test matrix $(\mathbf{V}^{(L)})^*\mathbf{\Omega}$. Then for node τ on level $L - 1$ with children α and β , we define

$$(4.8) \quad \mathbf{\Omega}_\tau = \begin{bmatrix} \mathbf{V}_\alpha^*\mathbf{\Omega}_\alpha \\ \mathbf{V}_\beta^*\mathbf{\Omega}_\beta \end{bmatrix} \quad \text{and} \quad \mathbf{Y}_\tau = \begin{bmatrix} \mathbf{U}_\alpha^*(\mathbf{Y}_\alpha - \mathbf{D}_\alpha\mathbf{\Omega}_\alpha) \\ \mathbf{U}_\beta^*(\mathbf{Y}_\beta - \mathbf{D}_\beta\mathbf{\Omega}_\beta) \end{bmatrix},$$

to be the corresponding blocks of the new test and sample matrices. We define $\mathbf{\Psi}_\tau$ and \mathbf{Z}_τ analogously. Once we have $\mathbf{\Omega}_\tau, \mathbf{\Psi}_\tau, \mathbf{Y}_\tau, \mathbf{Z}_\tau$ we compute $\mathbf{U}_\tau, \mathbf{V}_\tau$, and \mathbf{D}_τ exactly as before using (4.1), (4.2), and (4.6).

This process is applied to successively coarser levels of the tree until the root node is reached. For the root node τ , we have $\mathbf{Y}_\tau = \mathbf{D}^{(0)}\mathbf{\Omega}_\tau$, so we simply solve $\mathbf{D}^{(0)} = \mathbf{Y}_\tau\mathbf{\Omega}_\tau^\dagger$. The full compression procedure is summarized in [Algorithm 4.1](#).

Remark 4.3. The sampling at the finest level that was described in [subsection 4.1](#) is directly supported by the theory for the randomized SVD [[13](#), [25](#)] since the test matrices are provably Gaussian matrices drawn independently of the matrices that are approximated, cf. [Remark 4.1](#). The sampling at the higher levels is harder to analyze. Considering equation (4.7), we see that the test matrix for sampling $\tilde{\mathbf{A}}^{(L)}$ is $(\mathbf{V}^{(L)})^*\mathbf{\Omega}$, and for level ℓ , the test matrix for sampling $\tilde{\mathbf{A}}^{(\ell)}$ is $(\mathbf{V}^{(\ell)})^*(\mathbf{V}^{(\ell+1)})^* \dots (\mathbf{V}^{(L)})^*\mathbf{\Omega}$. Since for $\ell = L - 1, \dots, 1$, $\mathbf{V}^{(\ell)}$ depends not only on \mathbf{A} and $\mathbf{\Psi}$, but also on $\mathbf{\Omega}$, the argument in [Remark 4.1](#) does not immediately apply. Extensive numerical experiments (including those reported in [section 5](#)) indicate that these test matrices work just as well as “true” Gaussian ones, but we have not yet been able to substantiate this claim through analysis.

4.4. Asymptotic complexity. Recall that $r = k + p$, where k is the block rank of \mathbf{A} and p represents the amount of oversampling. We assume for simplicity that the leaf node size is $m = 2r$ and $\mathbf{\Omega}$ and $\mathbf{\Psi}$ both have $s = 3r$ columns. [Algorithm 4.1](#) requires s matrix-vector products of \mathbf{A} and \mathbf{A}^* , and an additional $\mathcal{O}(r^3)$ operations for each node in the tree, of which there are approximately $2N/m$. Therefore, the total compression time is

$$T_{\text{compress}} = 6rN \times T_{\text{rand}} + 6r \times T_{\text{mult}} + \mathcal{O}(r^2N) \times T_{\text{flop}},$$

where T_{rand} denotes the time to sample a value from the standard Gaussian distribution, T_{mult} denotes the time to apply \mathbf{A} or \mathbf{A}^* to a vector, and T_{flop} denotes the time to carry out a floating point arithmetic operation.

Algorithm 4.1 Compressing an HBS matrix

Compute randomized samples of \mathbf{A} and \mathbf{A}^* .
 Form Gaussian random test matrices $\mathbf{\Omega}$ and $\mathbf{\Psi}$ of size $N \times s$.
 Multiply $\mathbf{Y} = \mathbf{A}\mathbf{\Omega}$ and $\mathbf{Z} = \mathbf{A}^*\mathbf{\Psi}$.

Compress level by level from finest to coarsest.

for level $\ell = L, L-1, \dots, 0$ **do**

for node τ in level ℓ **do**

 Obtain test and sample matrices associated with τ .

if τ is a leaf node **then**

$$\mathbf{\Omega}_\tau = \mathbf{\Omega}(I_\tau, :), \quad \mathbf{\Psi}_\tau = \mathbf{\Psi}(I_\tau, :)$$

$$\mathbf{Y}_\tau = \mathbf{Y}(I_\tau, :), \quad \mathbf{Z}_\tau = \mathbf{Z}(I_\tau, :)$$

else

 Let α and β denote the children of τ .

$$\mathbf{\Omega}_\tau = \begin{bmatrix} \mathbf{V}_\alpha^* \mathbf{\Omega}_\alpha \\ \mathbf{V}_\beta^* \mathbf{\Omega}_\beta \end{bmatrix}, \quad \mathbf{\Psi}_\tau = \begin{bmatrix} \mathbf{U}_\alpha^* \mathbf{\Psi}_\alpha \\ \mathbf{U}_\beta^* \mathbf{\Psi}_\beta \end{bmatrix}$$

$$\mathbf{Y}_\tau = \begin{bmatrix} \mathbf{U}_\alpha^* (\mathbf{Y}_\alpha - \mathbf{D}_\alpha \mathbf{\Omega}_\alpha) \\ \mathbf{U}_\beta^* (\mathbf{Y}_\beta - \mathbf{D}_\beta \mathbf{\Omega}_\beta) \end{bmatrix}, \quad \mathbf{Z}_\tau = \begin{bmatrix} \mathbf{V}_\alpha^* (\mathbf{Z}_\alpha - \mathbf{D}_\alpha^* \mathbf{\Psi}_\alpha) \\ \mathbf{V}_\beta^* (\mathbf{Z}_\beta - \mathbf{D}_\beta^* \mathbf{\Psi}_\beta) \end{bmatrix}$$

end if

 Compute blocks of the factorization associated with τ .

if $\ell > 0$ **then**

$$\mathbf{P}_\tau = \text{null}(\mathbf{\Omega}_\tau, r), \quad \mathbf{Q}_\tau = \text{null}(\mathbf{\Psi}_\tau, r)$$

$$\mathbf{U}_\tau = \text{col}(\mathbf{Y}_\tau \mathbf{P}_\tau, r), \quad \mathbf{V}_\tau = \text{col}(\mathbf{Z}_\tau \mathbf{Q}_\tau, r)$$

$$\mathbf{D}_\tau = (\mathbf{I} - \mathbf{U}_\tau \mathbf{U}_\tau^*) \mathbf{Y}_\tau \mathbf{\Omega}_\tau^\dagger + \mathbf{U}_\tau \mathbf{U}_\tau^* \left((\mathbf{I} - \mathbf{V}_\tau \mathbf{V}_\tau^*) \mathbf{Z}_\tau \mathbf{\Psi}_\tau^\dagger \right)^*$$

else

$$\mathbf{D}_\tau = \mathbf{Y}_\tau \mathbf{\Omega}_\tau^\dagger$$

\triangleright Here τ is the root node.

end if

end for

end for

Remark 4.4 (Comparison of information efficiency). If we view each randomized sample as carrying N values worth of information about \mathbf{A} , and still assume $m = 2r$ and $s = 3r$, we find that the compression algorithm requires a total of $6rN$ values to reconstruct the matrix. The algorithm of [21] requires only r samples of \mathbf{A} and \mathbf{A}^* , but it also requires access to $\sim mN$ matrix entries that form the block-diagonal part of \mathbf{A} as well as $\sim rN$ elements that appear in interpolative decompositions, for a total of $5rN$ values worth of information. Therefore, the algorithm in the present work requires only slightly more information to recover \mathbf{A} , while having the advantage of being truly black box.

5. Numerical experiments. In this section, we present a selection of numerical results. We report the following quantities for a number of test problems: (1) the time to compress the operator, (2) the time to apply the compressed representation to a vector, (3) the relative accuracy of the compressed representation, and (4) the storage requirements of the compressed representation measured as the number of floating point values per degree of freedom. The algorithms for compressing matrices and applying compressed representations are written in Python, and the black-box

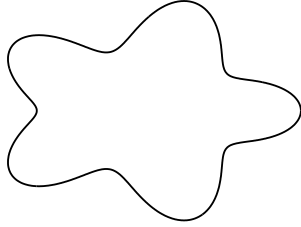


Fig. 3: Contour Γ on which the BIE (5.1) is defined.

multiplication routines are written in MATLAB. The experiments were carried out on a workstation with an Intel Core i9-10900K processor with 10 cores and 128GB of memory.

We report two measurements of compression time: (1) the time spent executing Algorithm 4.1, and (2) the execution time of Algorithm 4.1, excluding time spent within the black-box multiplication routine. In the context of black-box compression, it is assumed that the black-box multiplication routine is provided, and any steps taken to initialize the matvec routine are considered external to the compression algorithm. Therefore, the reported timings exclusively measure time spent executing Algorithm 4.1, and do not include any time initializing the matvec routines. For each of the test problems we consider in this section, the time spent initializing the matvec routines is relatively short, taking less than one tenth of the time spent for compression.

We measure the accuracy of the compressed matrices using the relative error

$$\frac{\|\tilde{\mathbf{A}} - \mathbf{A}\|}{\|\mathbf{A}\|},$$

where \mathbf{A} is the matrix defined by the matvec in each example. The operator norms are estimated by taking 20 steps of power iteration. Specifically, to estimate the norm of $\mathbf{B} \in \{\mathbf{A}, \tilde{\mathbf{A}} - \mathbf{A}\}$, we initialize a random vector \mathbf{v}_0 , compute $\mathbf{v}_i = \mathbf{B}^* \mathbf{B}(\mathbf{v}_{i-1}/\|\mathbf{v}_{i-1}\|)$, $i = 1, \dots, 20$, and take the estimate $\|\mathbf{B}\| \approx \sqrt{\|\mathbf{v}_{20}\|}$. We also report the maximum leaf node size m and the number r of random vectors per test matrix, which are inputs to the compression algorithm. We select the value of r so that the resulting approximation achieves reasonably high accuracy and then set $m = 2r$ and $s = 3r$.

5.1. Boundary integral equation. We consider a matrix arising from the discretization of the Boundary Integral Equation (BIE)

$$(5.1) \quad \frac{1}{2}q(x) + \int_{\Gamma} \frac{(x-y) \cdot n(y)}{4\pi|x-y|^2} q(y) ds(y) = f(x), \quad x \in \Gamma,$$

where Γ is the simple closed contour in the plane shown in Figure 3, and where $n(y)$ is the outwards pointing unit normal of Γ at y . The BIE (5.1) is a standard integral equation formulation of the Laplace equation with boundary condition f on the domain interior to Γ . The BIE (5.1) is discretized using the Nyström method on N equispaced points on Γ , with the Trapezoidal rule as the quadrature (since the kernel in (5.1) is smooth, the Trapezoidal rule has exponential convergence).

The fast matrix-vector multiplication is in this case furnished by the recursive skeletonization (RS) procedure of [19]. In this case, the “exact” matrix \mathbf{A} is the RS

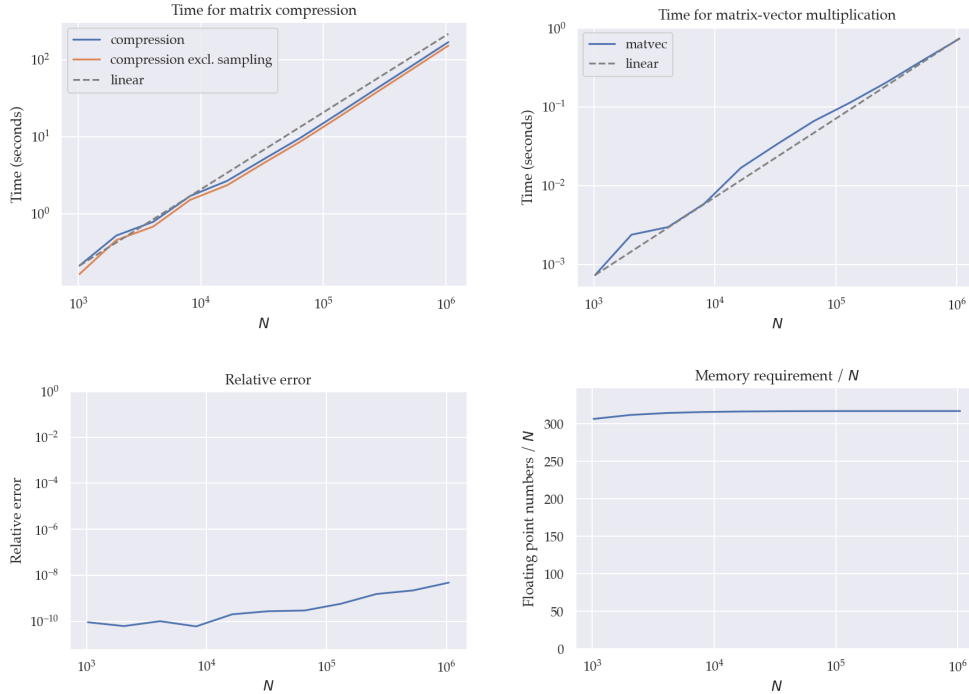


Fig. 4: Results from applying the compression algorithm to a double layer potential on a simple contour in the plane. Here $r = 30$ and $m = 60$.

approximation to the discretized integral operator. To avoid spurious effects due to the rank structure inherent in the RS representation, we set a computational tolerance in the RS approximation to $1e-15$. Additionally, we used a different tree structure for the RS compression than the one we used when running the randomized algorithm.

Results are given in [Figure 4](#).

Remark 5.1. The problem under consideration here is artificial in the sense that there is no actual need to use more than a couple of hundred points to resolve (5.1) numerically to double precision accuracy. It is included merely to illustrate the asymptotic performance of the proposed method.

5.2. Operator multiplication. We next investigate how the proposed technique performs on a matrix matrix multiplication problem. Specifically, we determine the Neumann-to-Dirichlet operator T for the contour shown in [Figure 3](#) using the well known formula

$$T = S \left(\frac{1}{2}I + D^* \right)^{-1},$$

where S is the single layer operator $[Sq](x) = \int_{\Gamma} -\frac{1}{2\pi} \log|x-y|q(y)ds(y)$, and where D^* is the adjoint of the double-layer operator $[D^*q](x) = \int_{\Gamma} \frac{n(x) \cdot (x-y)}{2\pi|x-y|^2} q(y)ds(y)$. The operators S and D are again discretized using a Nyström method on equispaced points (with sixth order Kapur-Rokhlin [15] corrections to handle the singularity in S), resulting in matrices \mathbf{S} and \mathbf{D} . The matrix $\mathbf{S}(\frac{1}{2}\mathbf{I} + \mathbf{D}^*)^{-1}$ is again applied using

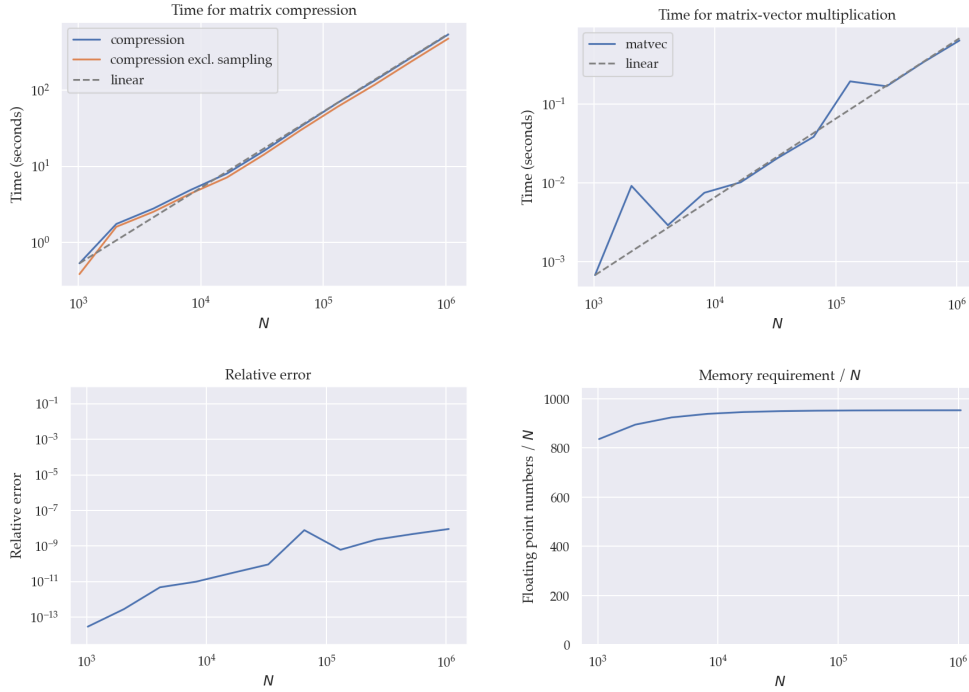


Fig. 5: Results from applying the compression algorithm to the Neumann-to-Dirichlet operator. Here $r = 100$ and $m = 200$.

the recursive skeletonization procedure of [19]. Results are given in [Figure 5](#).

5.3. Frontal matrices in nested dissection. Our next example is a simple model problem that illustrates the behavior of the proposed method in the context of sparse direct solvers. The idea here is to use rank structure to compress the increasingly large Schur complements that arise in the LU factorization of a sparse matrix arising from the finite element or finite difference discretization of an elliptic PDE, cf. [23, Ch. 21]. As a model problem, we consider a $51N \times 51N$ matrix \mathbf{C} that encodes the stiffness matrix for the standard five-point stencil finite difference approximation to the Poisson equation on a rectangle using a grid with $N \times 51$ nodes. We partition the grid into three sets $\{1, 2, 3\}$, as shown in [Figure 6](#), and then tessellate \mathbf{C} accordingly,

$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{11} & \mathbf{0} & \mathbf{C}_{13} \\ \mathbf{0} & \mathbf{C}_{22} & \mathbf{C}_{23} \\ \mathbf{C}_{31} & \mathbf{C}_{32} & \mathbf{C}_{33} \end{bmatrix},$$

where \mathbf{C}_{11} and \mathbf{C}_{22} are $25N \times 25N$, and \mathbf{C}_{33} is $N \times N$. The matrix we seek to compress is the $N \times N$ Schur complement

$$\mathbf{A} = \mathbf{C}_{33} - \mathbf{C}_{31}\mathbf{C}_{11}^{-1}\mathbf{C}_{13} - \mathbf{C}_{32}\mathbf{C}_{22}^{-1}\mathbf{C}_{23}.$$

In our example, we apply \mathbf{A} to vector by calling standard sparse direct solvers for the left and the right subdomains, respectively.

Results are given in [Figure 7](#).

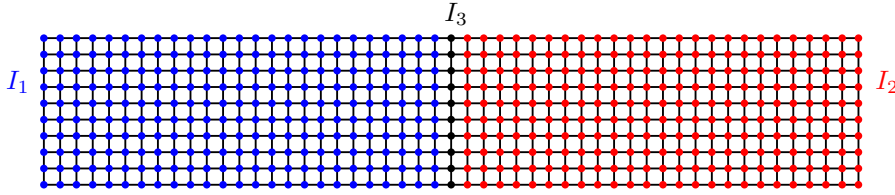


Fig. 6: An example of the grid in the sparse LU example described in Section 5.3. There are $N \times n$ points in the grid, shown for $N = 10, n = 51$.

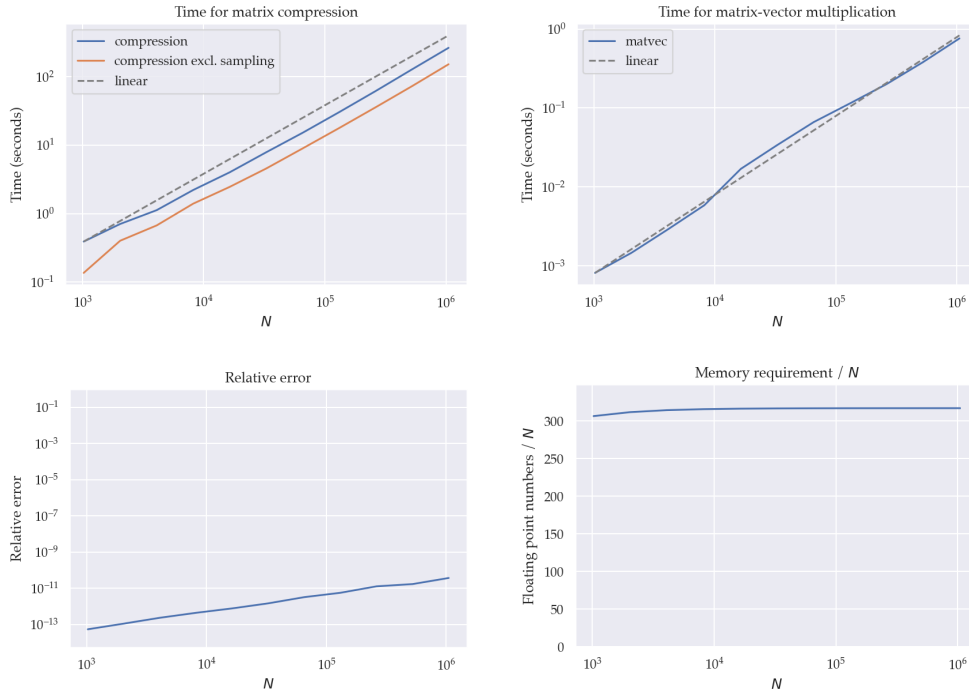


Fig. 7: Results from applying the compression algorithm to frontal matrices in the nested dissection algorithm. Here $r = 30$ and $m = 60$.

5.4. Summary of observations.

- The numerical results support our claim of linear complexity for all steps of the algorithms presented.
- The approximations achieve high accuracy in every case. In particular, the numerical results indicate that the method is numerically stable, and that there is minimal aggregation of errors as hierarchical trees deepen.
- The asymptotic storage cost, reported as number of floating point numbers per degree of freedom, appears to be independent of the problem size N . (It does, of course, depend on the rank parameter r and on the size of the diagonal blocks m .)

6. Conclusions. This paper presents an algorithm for black-box randomized compression of Hierarchically Block Separable matrices. To compress an $N \times N$ matrix \mathbf{A} , the algorithm requires only $\mathcal{O}(k)$ samples of \mathbf{A} and \mathbf{A}^* , where k is the block rank of \mathbf{A} . Numerical experiments demonstrate that the algorithm is accurate and very computationally efficient, with compression time scaling linearly in N when the cost of applying \mathbf{A} and \mathbf{A}^* to a vector is $\mathcal{O}(N)$.

Acknowledgments. The work reported was supported by the Office of Naval Research (N00014-18-1-2354), by the National Science Foundation (DMS-2313434 and DMS-1952735), and by the Department of Energy ASCR (DE-SC0022251).

REFERENCES

- [1] S. AMBIKASARAN AND E. DARVE, *An $o(n \log n)$ fast direct solver for partial hierarchically semi-separable matrices*, Journal of Scientific Computing, 57 (2013), pp. 477–501.
- [2] M. BEBENDORF, *Hierarchical matrices*, Springer, 2008.
- [3] S. BÖRM, *Efficient numerical methods for non-local operators*, vol. 14 of EMS Tracts in Mathematics, European Mathematical Society (EMS), Zürich, 2010. \mathcal{H}^2 -matrix compression, algorithms and analysis.
- [4] S. CHANDRASEKARAN, P. DEWILDE, M. GU, W. LYONS, AND T. PALS, *A fast solver for hss representations via sparse matrices*, SIAM Journal on Matrix Analysis and Applications, 29 (2007), pp. 67–81.
- [5] S. CHANDRASEKARAN, M. GU, AND T. PALS, *A fast ulv decomposition solver for hierarchically semiseparable representations*, SIAM Journal on Matrix Analysis and Applications, 28 (2006), pp. 603–622.
- [6] Z. CHEN AND J. J. DONGARRA, *Condition numbers of gaussian random matrices*, SIAM Journal on Matrix Analysis and Applications, 27 (2005), pp. 603–620.
- [7] A. EDELMAN AND B. D. SUTTON, *Tails of condition number distributions*, SIAM journal on matrix analysis and applications, 27 (2005), pp. 547–560.
- [8] A. GILLMAN, P. M. YOUNG, AND P.-G. MARTINSSON, *A direct solver with $o(n)$ complexity for integral equations on one-dimensional domains*, Frontiers of Mathematics in China, 7 (2012), pp. 217–247.
- [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix computations*, JHU press, 2013.
- [10] L. GREENGARD AND V. ROKHLIN, *A fast algorithm for particle simulations*, J. Comp. Phys, 73 (1987), pp. 325–348.
- [11] L. GREENGARD AND V. ROKHLIN, *A new version of the fast multipole method for the laplace equation in three dimensions*, Acta numerica, 6 (1997), pp. 229–269.
- [12] W. HACKBUSCH, *A sparse matrix arithmetic based on h -matrices. part i: Introduction to h -matrices*, Computing, 62 (1999), pp. 89–108.
- [13] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions*, SIAM review, 53 (2011), pp. 217–288.
- [14] K. L. HO AND L. GREENGARD, *A fast direct solver for structured linear systems by recursive skeletonization*, SIAM Journal on Scientific Computing, 34 (2012), pp. A2507–A2532.
- [15] S. KAPUR AND V. ROKHLIN, *High-order corrected trapezoidal quadrature rules for singular functions*, SIAM Journal on Numerical Analysis, 34 (1997), pp. 1331–1356.
- [16] J. LEVITT AND P.-G. MARTINSSON, *Randomized compression of rank-structured matrices accelerated with graph coloring*, arXiv preprint arXiv:2205.03406, (2022).
- [17] L. LIN, J. LU, AND L. YING, *Fast construction of hierarchical matrix representation from matrix-vector multiplication*, Journal of Computational Physics, 230 (2011), pp. 4071–4087.
- [18] P. MARTINSSON, *Rapid factorization of structured matrices via randomized sampling*, 2008. arXiv:0806.2339.
- [19] P. MARTINSSON AND V. ROKHLIN, *A fast direct solver for boundary integral equations in two dimensions*, J. Comp. Phys., 205 (2005), pp. 1–23.
- [20] P.-G. MARTINSSON, *A fast direct solver for a class of elliptic partial differential equations*, Journal of Scientific Computing, 38 (2009), pp. 316–330.
- [21] P.-G. MARTINSSON, *A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix*, SIAM Journal on Matrix Analysis and Applications, 32 (2011), pp. 1251–1274.

- [22] P.-G. MARTINSSON, *Compressing rank-structured matrices via randomized sampling*, SIAM Journal on Scientific Computing, 38 (2016), pp. A1959–A1986.
- [23] P.-G. MARTINSSON, *Fast Direct Solvers for Elliptic PDEs*, vol. CB96 of CBMS-NSF conference series, SIAM, 2019, <https://doi.org/10.1137/1.9781611976045>.
- [24] P.-G. MARTINSSON, *Fast direct solvers for elliptic PDEs*, SIAM, 2019.
- [25] P.-G. MARTINSSON AND J. A. TROPP, *Randomized numerical linear algebra: Foundations and algorithms*, Acta Numerica, 29 (2020), pp. 403–572.
- [26] V. MINDEN, K. L. HO, A. DAMLE, AND L. YING, *A recursive skeletonization factorization based on strong admissibility*, Multiscale Modeling & Simulation, 15 (2017), pp. 768–796.
- [27] J. VOGEL, J. XIA, S. CAULEY, AND V. BALAKRISHNAN, *Superfast divide-and-conquer method and perturbation analysis for structured eigenvalue solutions*, SIAM Journal on Scientific Computing, 38 (2016), pp. A1358–A1382, <https://doi.org/10.1137/15M1018812>, <https://arxiv.org/abs/https://doi.org/10.1137/15M1018812>.
- [28] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Superfast multifrontal method for large structured linear systems of equations*, SIAM Journal on Matrix Analysis and Applications, 31 (2010), pp. 1382–1411.