



40 identified, which is then compressed and approximated by a low rank representation.  
 41 Traditional truncated singular value decomposition (SVD) and rank-revealing QR fac-  
 42 torization are popular ways to construct such low rank approximations. But the cost  
 43 of implementing SVD or QR themselves are much higher than a single matrix-vector  
 44 multiplication, which is the typical core operation required for an iterative solution of  
 45 the linear systems related to the kernel matrices. The fast multipole method (FMM)  
 46 [12] is one of the most important fast algorithms for kernel summation, in which tar-  
 47 get and source points are hierarchically divided as well-separated sets, and on each  
 48 set, the kernel function is low-rank approximated by using multipole expansions. In  
 49 the original FMM, kernel function is approximated by analytical tools (either with  
 50 addition theorems of special functions or Taylor expansions) [12, 4, 6, 10, 5]. To over-  
 51 come the difficulties when analytic formulation of kernel functions is not available,  
 52 various semi-analytic[1, 11, 18] and algebraic FMMs [21, 22, 23] were developed in  
 53 recent decades. In some other approaches [16, 17], the whole kernel matrix is split  
 54 into block matrices with various ranks, on each of which the SVD was implemented  
 55 and then a truncated summation was used. More recently, a random interpolative  
 56 decomposition method was developed in the framework of FMM [20, 19, 24].

57 To take advantage of modern computational architecture and especially for high  
 58 dimensional data, randomized methods are found to provide efficient approach of low-  
 59 rank approximations to handle transient data sets, where access to the full matrices  
 60 may not be possible or too expensive as only single pass or constant number of passes  
 61 of the matrices are preferred or realistic. Some fast Monte-Carlo algorithms for large  
 62 matrix-matrix multiplications and low-rank approximations have been developed in [7,  
 63 8]. Still, those algorithms are generally not more efficient than *a single* implementation  
 64 of Eq. (1). For example, it was proposed [9][15] that by randomly sampling the  
 65 column and rows of a matrix, based on a the magnitude of the row/column vector, will  
 66 give a good low rank approximation in high probability. However, such an approach  
 67 requires the calculation of the  $L_2$  norm of the column/row vectors, which is already of  
 68 complexity  $O(MN)$  as the direct kernel summation cost. Furthermore, error analysis  
 69 in [7, 8] is for general matrices and relies on large number of samples. Number of  
 70 samples will be greatly limited in practice if high efficiency is pursued and the low-  
 71 rank property of the matrix should give some additional benefits in error analysis.

72 In this paper, we develop a novel hierarchical random compression algorithm for  
 73 kernel matrices in order to enhance the efficiency of kernel summation and reduce  
 74 data storage at a very large scale. This hierarchical approach is different from that of  
 75 FMM methods; it requires only a one-way top-down pass, and is simple to implement  
 76 in a recursive manner. More importantly, this algorithm does not use any analytic  
 77 expansion of the kernel function. Thus it is especially useful when the analytic formula  
 78 of kernel does not exist, such as Green's functions in complicated inhomogeneous  
 79 geometries. To achieve this goal, we first apply the randomized sampling of the column  
 80 and row space technique for the kernel matrix resulting from far-field interactions,  
 81 i.e. the target and source points set are well-separated. We show that for such a  
 82 scenery, the uniform sampling distribution will be sufficient to give a good low-rank  
 83 approximation, thus removing the need and associated cost of computing sampling  
 84 distributions. The expectation of error depends on the number of sampled column  
 85 (row) and the diameter-distance ratio of the two sets. Next, for general source-  
 86 field configurations, we will employ the  $\mathcal{H}$ -matrix framework [13, 14] to construct  
 87 a hierarchical multilevel tree structure, and apply the far-field randomized low-rank  
 88 approximation for admissible interactions at each level, recursively.

89 The rest of the paper is organized as follows. Section 2 describes the randomized

90 compression method for far field interaction matrices. Analysis of the algorithm for  
 91 the far field case is given in Section 3. Section 4 introduces the hierarchical random  
 92 compression method (HRCM) where  $\mathcal{H}$ -matrix framework is used with the random-  
 93 ized compression method applied to far field on different hierarchical level of the data  
 94 set. Numerical tests for the proposed HRCM are provided for kernels from electro-  
 95 static and Helmholtz Green's functions in Section 5. Finally, conclusion and discussion  
 96 are given in Section 6.

97 **2. Basic random algorithms for well-separated sources and targets.** For  
 98 convenience, we list some basic notations in linear algebra. For a vector  $\mathbf{x} \in \mathbb{K}^N$ ,  
 99 where  $\mathbb{K}$  represents either real  $\mathbb{R}$  or complex  $\mathbb{C}$ , denote Euclidean norm by  $|\mathbf{x}| =$   
 100  $\left(\sum_{i=1}^N |x_i|^2\right)^{1/2}$ . For a matrix  $\mathbf{A} \in \mathbb{K}^{M \times N}$ , let  $A^{(j)}$  and  $A_{(i)}$  denote its  $j$ -th column

101 and  $i$ -th row, respectively. The Frobenius norm is  $\|\mathbf{A}\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n A_{ij}^2}$ , and the

102 product of  $\mathbf{AB}$  can be written as  $\mathbf{AB} = \sum_{j=1}^N A^{(j)} B_{(j)}$ .

103 In this section and Section 3 we restrict ourselves to low-rank matrices corre-  
 104 sponding to the approximation of kernel function for well-separated target and source  
 105 points.

106 Define the diameter of the target set  $\{\mathbf{r}_i\}$  as

107 (2) 
$$\text{diam}(T) := \max_{i,i'} |\mathbf{r}_i - \mathbf{r}_{i'}|,$$

108 using the Euclidean norm in  $\mathbb{R}^d$ . Diameter of the source set  $\text{diam}(S)$  is defined  
 109 similarly. Additionally, we will need the distance of the two groups as

110 (3) 
$$\text{dist}(T, S) := |\mathbf{r}_T^* - \mathbf{r}_S^*|,$$

111 where  $\mathbf{r}_T^*$  and  $\mathbf{r}_S^*$  are the Chebyshev centers of sets  $T$  and  $S$ , respectively. Then we say  
 112 the source and target charges are well-separated if  $\text{dist}(T, S) \geq \frac{1}{2}(\text{diam}(T) + \text{diam}(S))$ .

113 **2.1. Low-rank characteristics of kernel functions.** We say the kernel func-  
 114 tion  $\mathcal{K}(\mathbf{r}, \mathbf{r}')$  is a generalized asymptotically smooth function [3, 2] if

115 (4) 
$$|\partial_{\mathbf{r}}^\alpha \partial_{\mathbf{r}'}^\beta \mathcal{K}(\mathbf{r}, \mathbf{r}')| \leq c(|\alpha|, |\beta|)(1 + k|\mathbf{r} - \mathbf{r}'|)^{|\alpha| + |\beta|} |\mathbf{r} - \mathbf{r}'|^{-|\alpha| - |\beta| - \tau},$$

116 where parameters  $k, \tau \geq 0$  and  $\alpha, \beta$  are  $d$ -dimensional multi-indices. The constant  
 117  $C(|\alpha|, |\beta|)$  only depends on  $|\alpha|$  and  $|\beta|$ . Condition (4) covers a wide-range of kernel  
 118 functions. For example, for Green's functions of Laplace equation,  $k = 0$ ,  $\tau = 0$   
 119 for the 2D case and  $k = 0$ ,  $\tau = 1$  for the 3D case, respectively. For 3D Green's functions  
 120 of Helmholtz equation, one has  $k$  as wave number and  $\tau = 1$ .

121 Consider the Taylor expansion of  $\mathcal{K}(\mathbf{r}, \mathbf{r}')$  around  $\mathbf{r}^*$ , which is the Chebyshev  
 122 center of  $T$ :  $\mathcal{K}(\mathbf{r}, \mathbf{r}') = \tilde{\mathcal{K}}(\mathbf{r}, \mathbf{r}') + R$  with the polynomial

123 (5) 
$$\tilde{\mathcal{K}}(\mathbf{r}, \mathbf{r}') = \sum_{|v|=0}^{m-1} \frac{1}{v!} (\mathbf{r} - \mathbf{r}^*)^v \frac{\partial^v \mathcal{K}(\mathbf{r}^*, \mathbf{r}')}{\partial \mathbf{r}^v},$$

124 and the remainder  $R$  satisfies

125 (6) 
$$|R| = |\mathcal{K}(\mathbf{r}, \mathbf{r}') - \tilde{\mathcal{K}}(\mathbf{r}, \mathbf{r}')| \leq \frac{1}{m!} |\mathbf{r} - \mathbf{r}^*|^m \max_{\zeta \in T, |\gamma|=m} \left| \frac{\partial^\gamma \mathcal{K}(\zeta, \mathbf{r}')}{\partial \zeta^\gamma} \right|.$$

126 If the well-separated target (T) and source (S) charges satisfies

$$127 \quad (7) \quad \frac{\max\{\text{diam}(T), \text{diam}(S)\}}{\text{dist}(T, S)} < \eta,$$

128 for a parameter  $0 < \eta < 1$ , then Eq. (6) has the following estimate,

$$129 \quad (8) \quad |\mathcal{K}(\mathbf{r}, \mathbf{r}') - \tilde{\mathcal{K}}(\mathbf{r}, \mathbf{r}')| \leq c(m)\eta^m(1+k|\mathbf{r}-\mathbf{r}'|)^m|\mathbf{r}-\mathbf{r}'|^{-\tau}, \quad \mathbf{r} \in T, \mathbf{r}' \in S.$$

130 Estimate (8) implies that for small wavenumber parameter  $k$ , if one replaces  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j)$   
131 in Eq. (1) by  $\tilde{\mathcal{K}}(\mathbf{r}_i, \mathbf{r}_j)$  as defined in Eq. (5) with error  $O(\eta^m)$ , the corresponding  
132 interaction matrix  $\tilde{\mathcal{K}}$  is of low rank, i.e.

$$133 \quad (9) \quad \text{rank}(\tilde{\mathcal{K}}) = K(m) \ll \min\{M, N\} = \text{rank}(\mathcal{K}),$$

134 where

$$135 \quad (10) \quad K(m) = \sum_{p=0}^{m-1} \frac{(d-1+p)!}{(d-1)!p!}.$$

136 In the case of 2D, we have  $d = 2$  and  $K(m) = m(m+1)/2$ .

137 This property indicates that we can replace matrix  $\mathcal{K}$  by the low-rank matrix  $\tilde{\mathcal{K}}$   
138 with enough accuracy for well-separated target and source points. Thus, the efficiency  
139 of computation could be greatly improved. However, in many situations, there is  
140 not easily available explicit formula for  $\tilde{\mathcal{K}}$ , as for layered Green's function, thus the  
141 matrix  $\tilde{\mathcal{K}}$  cannot not be explicitly computed by Eq. (5). In our approach, we will  
142 take advantage of the fact that  $\mathcal{K}$  has redundant information (essential low rank  
143 characteristics) and will use randomized sampling methods to select a small amount  
144 of its rows or columns and approximate the full matrix with the sampled sub-matrices.  
145 It should be noted that the low-rank characteristics depend on parameter  $k$  and  $\tau$  in  
146 the decaying condition of the kernel (4).

147 **2.2. Random kernel compression algorithms.** For matrix  $\mathbf{A} \in \mathbb{K}^{M \times N}$  with  
148 a rank  $K \ll \min\{M, N\}$ , the matrix-vector multiplication with vector  $\mathbf{x} \in \mathbb{K}^N$  can  
149 be represented as

$$150 \quad (11) \quad \mathbf{Ax} = \sum_{i=1}^K \sigma_A^i \mathbf{U}_A^{(i)} \mathbf{V}_A^{(i)*} \mathbf{x},$$

151 where  $\mathbf{A} = \mathbf{U}_A \Sigma_A \mathbf{V}_A^*$ ,  $\mathbf{U}_A \in \mathbb{K}^{M \times M}$ ,  $\mathbf{V}_A \in \mathbb{K}^{N \times N}$  is the SVD decomposition of  
152 matrix  $\mathbf{A}$ . Although Eq. (11) includes small amount of summation, performing SVD  
153 to obtain  $\sigma_A^i$ ,  $\mathbf{U}_A^i$  and  $\mathbf{V}_A^i$  is much more expensive than the direct multiplication.  
154 Instead, we will approximate the singular values and unitary matrices by fast Monte  
155 Carlo methods. This process is outlined as follows:

156 Let  $\mathbf{C} \in \mathbb{K}^{M \times c}$  be the matrix made of  $c$  columns sampled from matrix  $\mathbf{A}$  and  
157 denote  $\mathbf{C} = \mathbf{U}_c \Sigma_c \mathbf{V}_c^*$ , where  $\mathbf{U}_c \in \mathbb{K}^{M \times M}$  and  $\mathbf{V}_c \in \mathbb{K}^{c \times c}$ . Further, let  $\mathbf{C}_r \in \mathbb{K}^{r \times c}$   
158 be the matrix made of  $r$  rows sampled from  $\mathbf{C}$  and denote  $\mathbf{C}_r = \mathbf{U}_r \Sigma_r \mathbf{V}_r^*$ , where  
159  $\mathbf{U}_r \in \mathbb{K}^{r \times r}$  and  $\mathbf{V}_r \in \mathbb{K}^{c \times c}$ . Then

- 160 • SVD is performed for the much smaller matrix  $\mathbf{C}_r$ , in which  $r, c \leq K \ll$   
161  $\min\{M, N\}$  and independent of  $M$  and  $N$ . Due to rapid decay of singular  
162 values, only first  $t_0$  columns of  $\mathbf{V}_r$  are needed, where  $t_0$  is determined by  
163 checking  $|\sigma_r^{t_0+1} - \sigma_r^{t_0}| < \epsilon$ , where  $\epsilon$  is a preset accuracy criteria.

164

- 165 • Since  $\mathbf{V}_r$  is now available, we can implement  $\mathbf{C}\mathbf{V}_r$ , which will be used to  
 166 approximate  $\mathbf{U}_c \Sigma_c = \mathbf{C}\mathbf{V}_c$ . Then a QR factorization of  $\mathbf{C}\mathbf{V}_r$  is computed,  
 167 i.e.,  $\mathbf{C}\mathbf{V}_r = \tilde{\mathbf{U}}_c \mathbf{R}$ , where  $\mathbf{R}$  is an upper triangle matrix. The resulting  $\tilde{\mathbf{U}}_c$  is  
 168 considered as the approximation of  $\mathbf{U}_c$ .  
 169
- 170 • Finally we take the approximation  $\mathbf{U}_A \approx \tilde{\mathbf{U}}_A = \tilde{\mathbf{U}}_c$  and  $\mathbf{A}\mathbf{x} \approx \tilde{\mathbf{U}}_A (\tilde{\mathbf{U}}_A^* \mathbf{A}) \mathbf{x}$ .  
 171 Note that exact computation of  $\tilde{\mathbf{U}}_A^* \mathbf{A}$  still requires  $O(t_0 MN)$  operations, so  
 172 the Monte Carlo Basic matrix multiplication algorithm in [7] is adopted, i.e.,  
 173 only  $c$  columns from  $\tilde{\mathbf{U}}_A^*$  and  $c$  rows from  $\mathbf{A}$  are sampled and computed,  
 174 this approximation reduces the complexity of matrix-matrix multiplication  
 175 to  $O(ct_0 N)$ .

176 Detailed implementations are given in Algorithm 1, where  $\Pi \mathbf{A}$  is represented  
 177 by orthonormal vectors. Note that overall the random kernel compression will have  
 178  $O(\max\{M, N\})$  operations.

---

**Algorithm 1** Random kernel compression
 

---

**Input:** matrix  $\mathbf{A} \in \mathbb{K}^{M \times N}$ ,  $c, r, \epsilon \in \mathbb{N}$ , such that  $1 < c \ll N$  and  $1 < r \ll M$ . A  
 constant  $\epsilon > 0$ .

**Output:**  $\sigma_t \in \mathbb{R}^+$ , orthonormal vectors  $\mathbf{U}_t \in \mathbb{K}^M$ , and  $\mathbf{V}_t \in \mathbb{K}^N$ , for  $t = 1, 2, \dots, t_0 \ll$   
 $\min(M, N)$ , such that

$$(12) \quad \mathbf{A} \approx \Pi \mathbf{A} = \sum_{t=1}^{t_0} \sigma_t \mathbf{U}_t \mathbf{V}_t^*.$$

1. Construct matrix  $\mathbf{C} \in \mathbb{K}^{M \times c}$ ;  
**for**  $t = 1$  to  $c$  **do**  
 (a) pick  $i_t \in 1, 2, \dots, N$  randomly using uniform distribution;  
 (b) set  $\mathbf{C}^{(t)} = \sqrt{N/c} \mathbf{A}^{(i_t)}$ ;  
**end for**
  2. Construct matrix  $\mathbf{C}_r \in \mathbb{R}^{r \times c}$ ;  
**for**  $t = 1$  to  $r$  **do**  
 (a) pick  $j_t \in 1, 2, \dots, M$  randomly using uniform distribution;  
 (b) set  $\mathbf{C}_{r(t)} = \sqrt{M/r} \mathbf{C}_{(j_t)}$ ;  
**end for**
  3. Perform SVD on matrix  $\mathbf{C}_r$ , i.e.,  $\mathbf{C}_r = \mathbf{U}_r \Sigma_r \mathbf{V}_r^*$  and denote  $\sigma(\mathbf{C}_r)$  as the  
 singular values.
  4. Let  $t_0 = \min\{r, c, \max\{j, \sigma_j(\mathbf{C}_r)\} > \epsilon\}$ ;  
**for**  $t = 1$  to  $t_0$  **do**  
 (a)  $\sigma_t = \sigma_t(\mathbf{C}_r)$ ;  
 (b)  $\mathbf{U}_t^c = \mathbf{C}\mathbf{V}_r$ ;  
**end for**
  5. Perform QR decomposition on  $\{\mathbf{U}_t^c\}$  to obtain the output orthonormal vectors  
 $\{\mathbf{U}_t\}_{t=1}^{t_0}$ .
  6. Approximate  $\mathbf{V} = \mathbf{A}^* \mathbf{U}_t$  by the Monte Carlo Basic matrix multiplication algo-  
 rithm in [7], with  $c$  columns from  $\mathbf{A}^*$  and  $c$  rows from  $\mathbf{U}_t$ ;
  7. Perform QR decomposition on  $\mathbf{V} \in \mathbb{K}^{N \times t_0}$  to obtain the output orthonormal  
 vectors  $\{\mathbf{V}_t\}_{t=1}^{t_0}$ .
- 

179 Note that in this algorithm, only the small matrix  $\mathbf{C}_r$  is stored in memory and

180 entries of other matrices are calculated on the fly. The idea of Algorithm 1 is similar to  
 181 the ‘‘ConstantTimeSVD algorithm’’ in [8]. However, in that work, the columns  $\mathbf{U}_t$  are  
 182 directly calculated as  $\mathbf{C}\mathbf{V}_r/\sigma_t(\mathbf{C}_r)$  and not orthonormal. More importantly, when  $\mathbf{V}_r$   
 183 is not close enough to  $\mathbf{V}_c$  (otherwise efficiency will be compromised), dividing rapidly  
 184 decaying singular values is not numerically stable. There will be no such issues in  
 185 Algorithm 1.

186 **3. Analysis of the compression algorithm for well-separated sets.** In  
 187 this section we investigate the error of the random compression operator  $\Pi$  in (12) for  
 188  $\mathbf{A}$  as generated by Algorithm 1,

$$189 \quad (13) \quad \|(\mathbf{A} - \Pi\mathbf{A})\mathbf{x}\| \leq \|(\mathbf{A} - \Pi\mathbf{A})\| \|\mathbf{x}\|.$$

190 Unfortunately, a full analysis on  $\Pi\mathbf{A}$  with general sampling probability is technically  
 191 complicated. We would rather consider a theoretically simpler case about

$$192 \quad (14) \quad \mathbf{A} - \tilde{\Pi}\mathbf{A}, \quad \text{with } \tilde{\Pi} = \mathbf{U}_K\mathbf{U}_K^*$$

193 where  $\mathbf{U}_K$  is the matrix containing only the first  $K$  columns of  $\mathbf{U}_c$  in  $\mathbf{C} = \mathbf{U}_c\Sigma_c\mathbf{V}_c^*$ .  
 194 Matrix  $\mathbf{C}$  here is the same as in Algorithm 1, but could be associated with arbitrary  
 195 sampling probability. Construction of the projector  $\tilde{\Pi}$  is equivalent to the ‘‘Lin-  
 196 earTimeSVD’’ algorithm in [8]. It is computationally inefficient but theoretically  
 197 simple. We cite the result from [8]:

198 **THEOREM 3.1** (Theorem 2 and 3 in [8]). *Suppose  $\mathbf{A} \in \mathbb{K}^{M \times N}$  and  $\mathbf{C} \in \mathbb{K}^{M \times c}$*   
 199 *being the column sampled matrix from  $\mathbf{A}$  as in Algorithm 1. Let  $\tilde{\Pi}$  be the projector*  
 200 *defined in Eq. (14), then*

$$201 \quad (15) \quad \|\mathbf{A} - \tilde{\Pi}\mathbf{A}\|_F^2 \leq \|\mathbf{A} - \mathbf{A}_K\|_F^2 + 2\sqrt{K}\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F;$$

202 and

$$203 \quad (16) \quad \|\mathbf{A} - \tilde{\Pi}\mathbf{A}\|_2^2 \leq \|\mathbf{A} - \mathbf{A}_K\|_2^2 + 2\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_2;$$

204 where  $\mathbf{A}_K$  the best  $K$ -rank approximation of  $\mathbf{A}$ .

205 Although the theoretical projector  $\tilde{\Pi}$  is different from projector  $\Pi$  of (12) in  
 206 Algorithm 1, we still can obtain meaningful insights about sampling strategies and  
 207 accuracy of Algorithm 1 by examining Theorem 3.1.

208 Since  $\mathbf{A}$  is already assumed as low-rank, we can only focus on estimating  $\|\mathbf{A}\mathbf{A}^* -$   
 209  $\mathbf{C}\mathbf{C}^*\|_\xi, \xi = 2, F$ . We look for a practical sampling probability and derive an error  
 210 estimate for well-separated source and target points. For simplicity, we assume  $a =$   
 211  $\text{diam}(T) = \text{diam}(S)$ ,  $\delta = \text{dist}(T, S)$ , and  $a/\delta \leq \eta$  for some  $\eta \in (0, 1)$ .

212 **3.1. Nearly optimal uniform sampling for far field kernel matrices.** The  
 213 low rank property of a matrix  $\mathbf{A}$  means most of its columns/rows are linearly depen-  
 214 dent, while each column/rows may contribute much differently to the overall matrix  
 215 property. For example the matrix  $A = vv^T$  with  $v^T = (1, 2, 3, \dots, N)$  is only of rank  
 216 one, but the last column/row is the most important. Therefore, the column/row  
 217 sampling probability is a critical factor in minimizing the error  $\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_\xi$ .

218 The optimal probability of (column) sampling to perform a Monte Carlo matrix-  
 219 matrix multiplication  $\mathbf{A}\mathbf{B}$  is proposed in [7].

220 **DEFINITION 3.2** (Optimal probability). *For  $\mathbf{A} \in \mathbb{K}^{M \times N}, \mathbf{B} \in \mathbb{K}^{N \times P}, 1 \leq c \leq$*   
 221  *$N$ , and  $\{p_j\}_{j=1}^N$  such that*

$$222 \quad (17) \quad p_j = \frac{|A^{(j)}||B_{(j)}|}{\sum_{j'=1}^N |A^{(j')}||B_{(j')}|}, \quad j = 1, 2, \dots, N,$$

223 then for  $t = 1$  to  $c$ , pick  $i_t \in \{1, 2, \dots, N\}$  with  $\Pr[i_t = j] = p_j, j = 1, 2, \dots, N$  inde-  
 224 pendently and with replacement. Set  $C^{(t)} = A^{(i_t)}/\sqrt{cp_{i_t}}$  and  $R_{(t)} = B_{(i_t)}/\sqrt{cp_{i_t}}$ , the  
 225 expectation value  $\mathbf{E} [\|\mathbf{A}\mathbf{B} - \mathbf{C}\mathbf{R}\|_F]$  is minimized.

226 Applying this definition to  $\mathbf{A}\mathbf{A}^*$ , and using Lemma 4 in [7], it is easy to conclude  
 227 that if  $c$  columns are sampled, the expectation of error  $\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F^2$  is minimized  
 228 as

$$229 \quad (18) \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F^2] = \frac{1}{c} (\|\mathbf{A}\|_F^4 - \|\mathbf{A}\mathbf{A}^*\|_F^2).$$

230 However, unless known in advance, utilizing the optimal probability is not practical  
 231 because its computation is already more expensive than the actual matrix-vector  
 232 multiplication.

233 In practice, a nearly optimal probability introduced in [7] will be used for our  
 234 approach.

235 DEFINITION 3.3 (Nearly optimal probability). For the same conditions in Defi-  
 236 nition 3.2, the probability  $\{p_j\}$  is called a nearly optimal probability if

$$237 \quad (19) \quad p_j \geq \frac{\beta |A^{(j)}| |B_{(j)}|}{\sum_{j'=1}^N |A^{(j')}| |B_{(j')}|}, \quad j = 1, 2, \dots, N,$$

238 for some  $0 < \beta \leq 1$ .

239 With the nearly optimal probability, one has

$$240 \quad (20) \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F^2] \leq \frac{1}{\beta c} \|\mathbf{A}\|_F^4 - \frac{1}{c} \|\mathbf{A}\mathbf{A}^*\|_F^2.$$

241 Now we will present the following result.

242 THEOREM 3.4. For well-separated source and target points, uniform sampling pro-  
 243 vides a nearly optimal probability.

244 Proof. For Eq. (18), the optimal probability is

$$245 \quad (21) \quad p_j = \frac{|A^{(j)}|^2}{\|\mathbf{A}\|_F^2}.$$

246 Recall entries of matrix  $\mathbf{A}$  are  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j)q_i$  and denote the distance  $r_{ij} = |\mathbf{r}_i - \mathbf{r}_j|$  for  
 247 simplicity. Then, for the well-separated source and target points, we have the bound

$$248 \quad \delta - a \leq r_{ij} \leq \delta + a,$$

249 since  $|\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j)|^2$  is monotonically decreasing with  $r_{ij}$ , we have

$$250 \quad (22) \quad p_j = \frac{|A^{(j)}|^2}{\|\mathbf{A}\|_F^2} \leq \frac{|\mathcal{K}(\delta - a)|^2 \sum_{i=1}^M |q_i|^2}{N |\mathcal{K}(\delta + a)|^2 \sum_{i=1}^M |q_i|^2}.$$

251 Thus, the uniform sampling probability

$$252 \quad (23) \quad \hat{p}_j = \frac{1}{N} > \frac{|\mathcal{K}(\delta + a)|^2 |A^{(j)}|^2}{|\mathcal{K}(\delta - a)|^2 \|\mathbf{A}\|_F^2}, \quad \square$$

253 is the the nearly optimal probability with  $\beta = \frac{|\mathcal{K}(\delta + a)|^2}{|\mathcal{K}(\delta - a)|^2} < 1$ .

254 Theorem 3.4 indicates that for a fixed number of samples and without additional  
 255 computational effort (uniform sampling), we can achieve the nearly optimal accuracy  
 256 as in estimate (20) in the kernel compression for well-separated source and target  
 257 points. However, the bound in Eq. (20) depends on parameter  $\beta$ , which is an indicator  
 258 of how “well” the two sets are separated. In case the two sets “touching” each other,  
 259 one has  $\eta \rightarrow 1$  or  $\delta - a \rightarrow 0$  and hence  $\beta \rightarrow 0$ . As a result, the error bound (20) fails.

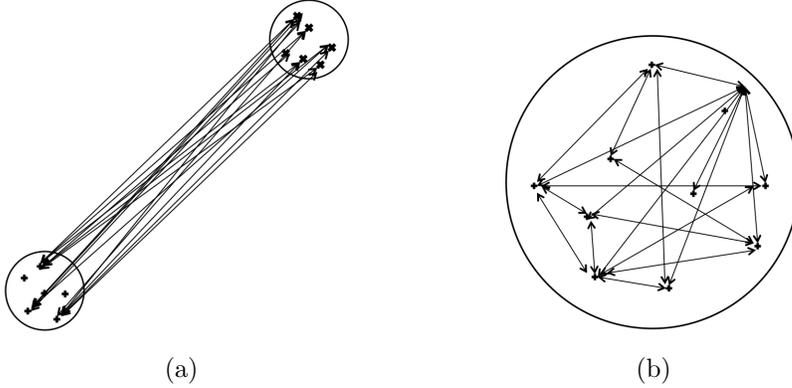


FIG. 1. (a) well-separated; (b) not well-separated.

260 The reason why uniform sampling works can also be illustrated heuristically by  
 261 Fig.1: When target and source sets are well-separated and have a small diameter-  
 262 distance ratio  $\eta$  as in Fig. 1(a), if we associate the interaction between target/source  
 263 points as lines connecting the sources and targets, it can be seen that all interactions  
 264 are “similar”, in terms of direction and magnitude, to each other. Then in this case,  
 265 columns or row vectors in the matrix have fairly the same contribution, so uniform  
 266 sampling will perform well with a small error. On the other hand, if target and source  
 267 points are mixed and belong to the same set as in Fig. 1(b), the interaction lines are  
 268 rather different from each other. Even the matrix maybe low rank, uniform sampling  
 269 will yield in uncontrollable error due to the complexity of the “interaction” lines.

270 In the next subsection we provide a further quantitative analysis of Eqs. (18) or  
 271 (20) for the separation. We show that even the error (18) or (20) depends on not only  
 272 sample number  $c$ , but also the diameter-distance ratio  $\eta$ .

273 **3.2. Error analysis on target-source separation.** Error estimates (18) and  
 274 (20) provided in [8] are for a general matrix  $\mathbf{A}$ , stating that the error is small if the  
 275 samples are large enough. But in practice, the sample number  $c$  can not be too large  
 276 due to efficiency requirement. Here, we give a finer estimate for  $\mathbf{A}$  corresponding to  
 277 well-separated target and source points. We show that for some type of kernels, if the  
 278 target and source sets are far way enough, the error is small enough regardless of the  
 279 sample numbers. Since we are more interested in the dependence of error bound on  
 280 diameter-distance ration, we will consider Eq. (18) for simplicity.

281 **THEOREM 3.5.** Suppose  $\mathbf{A} \in \mathbb{K}^{M \times N}$  is the matrix generated from the kernel func-  
 282 tion  $\mathcal{K}(\mathbf{r}, \mathbf{r}')$  satisfying condition (4), and  $\mathbf{C} \in \mathbb{K}^{M \times c}$  being the column sampled matrix  
 283 from  $\mathbf{A}$ . Let  $\tilde{\Pi}$  be the orthogonal projector defined in Eq. (14), then

$$284 \quad (24) \quad \|\mathbf{A} - \tilde{\Pi}\mathbf{A}\|_F^2 \leq \|\mathbf{A} - \mathbf{A}_K\|_F^2 + 2\sqrt{K}\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F;$$

285 *or*

$$286 \quad (25) \quad \|\mathbf{A} - \tilde{\Pi}\mathbf{A}\|_2^2 \leq \|\mathbf{A} - \mathbf{A}_K\|_2^2 + 2\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_2;$$

287 *and*

$$288 \quad (26) \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F] \leq C(\delta, a, \tau, k, M, N) \frac{1}{\sqrt{c}} \frac{2\eta}{2 - \eta} \|\mathbf{A}\|_F^2,$$

289 *where*  $\delta$  *and*  $a$  *are the distance between and diameter for the target and source sets,*  
 290 *respectively,  $0 < \eta < 1$  is the diameter-over-distance ratio of target and source sets,*  
 291 *and*

$$292 \quad (27) \quad C(\delta, a, \tau, k, M, N) = \left(\frac{\delta}{2}\right)^{-\tau} \frac{(1 + 2k\delta) \sqrt{M}}{|\mathcal{K}(\delta + a)| \sqrt{N}}.$$

293 *Proof.* Equations (24) and (25) have been stated in Theorem 3.1. To prove (26),  
 294 write the  $i$ -th row of  $\mathbf{A}$  as

$$295 \quad (28) \quad A_{(i)} = q_i(\mathcal{K}(\mathbf{r}_i, \mathbf{r}_1), \mathcal{K}(\mathbf{r}_i, \mathbf{r}_2), \dots, \mathcal{K}(\mathbf{r}_i, \mathbf{r}_N)), \quad i = 1, 2, \dots, M.$$

296 Then, the difference between  $q_{i'}A_{(i)}$  and  $q_iA_{(i')}$  is, approximated to the first order,

$$297 \quad (29) \quad \Delta_{ii'} = q_i q_{i'} (\Delta_{ii'}^{(1)}, \Delta_{ii'}^{(2)}, \dots, \Delta_{ii'}^{(N)}),$$

298 where

$$299 \quad (30) \quad \Delta_{ii'}^{(j)} = \mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) - \mathcal{K}(\mathbf{r}_{i'}, \mathbf{r}_j) = (\mathbf{r}_i - \mathbf{r}_{i'}) \cdot \frac{\partial}{\partial \mathbf{r}} \mathcal{K}(\mathbf{r}, \mathbf{r}_j)|_{\mathbf{r}=\mathbf{r}_i'}.$$

300 Recall assumption in Eq. (4) and condition (7), we have estimate

$$301 \quad (31) \quad |\Delta_{ii'}^{(j)}| \leq \frac{a}{\delta - \frac{a}{2}} \left(1 + k \left(\delta + \frac{a}{2}\right)\right) \left(\delta - \frac{a}{2}\right)^{-\tau} \leq \frac{2\eta}{2 - \eta} (1 + 2k\delta) \left(\frac{\delta}{2}\right)^{-\tau}.$$

302 Rewrite Eq. (18) as

$$\begin{aligned} 303 \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F^2] &= \frac{1}{c} (\|\mathbf{A}\|_F^4 - \|\mathbf{A}\mathbf{A}^*\|_F^2) \\ 304 \quad &= \frac{1}{c} \left[ \left( \sum_{i=1}^M |A_{(i)}|^2 \right)^2 - \sum_{i=1}^M \sum_{i'=1}^M \langle A_{(i)}, A_{(i')} \rangle^2 \right] \\ 305 \quad (32) \quad &= \frac{1}{c} \sum_{i=1}^M \sum_{i'=1}^M \frac{1}{q_i^2 q_{i'}^2} [\langle q_{i'} A_{(i)}, q_{i'} A_{(i)} \rangle \langle q_i A_{(i')}, q_i A_{(i')} \rangle - \langle q_{i'} A_{(i)}, q_i A_{(i')} \rangle^2], \\ 306 \end{aligned}$$

307 where  $\langle \cdot, \cdot \rangle$  represents inner product.

308 Since  $q_i A_{(i')} = q_{i'} A_{(i)} + \Delta_{ii'}$ , using the linearity of inner product, we have

$$\begin{aligned} 309 \quad &\langle q_{i'} A_{(i)}, q_{i'} A_{(i)} \rangle \langle q_i A_{(i')}, q_i A_{(i')} \rangle - \langle q_{i'} A_{(i)}, q_i A_{(i')} \rangle^2 \\ 310 \quad &= \langle q_{i'} A_{(i)}, q_{i'} A_{(i)} \rangle \langle \Delta_{ii'}, \Delta_{ii'} \rangle - \langle q_{i'} A_{(i)}, \Delta_{ii'} \rangle^2 \\ 311 \quad (33) \quad &\leq \langle q_{i'} A_{(i)}, q_{i'} A_{(i)} \rangle \langle \Delta_{ii'}, \Delta_{ii'} \rangle. \end{aligned}$$

313 Plugging (33) in (32) and using estimate (31), we arrive at

$$\begin{aligned}
314 \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F^2] &\leq \frac{1}{c} \sum_{i=1}^M \sum_{i'=1}^M \langle A_{(i)}, A_{(i')} \rangle \langle \Delta_{ii'} \Delta_{ii'} \rangle \\
315 \quad (34) \quad &\leq \frac{M\|Q\|_2^2}{c} \left(\frac{2\eta}{2-\eta}\right)^2 (1+2k\delta)^2 \left(\frac{\delta}{2}\right)^{-2\tau} \|\mathbf{A}\|_F^2, \\
316
\end{aligned}$$

317 where  $Q = (q_1, q_2, \dots, q_M)$  is the density of target points. By Jensen's inequality

$$\begin{aligned}
318 \quad \mathbf{E} [\|\mathbf{A}\mathbf{A}^* - \mathbf{C}\mathbf{C}^*\|_F] &\leq \frac{\sqrt{M}\|Q\|_2}{\sqrt{c}\|\mathbf{A}\|_F} (1+2k\delta) \left(\frac{\delta}{2}\right)^{-\tau} \frac{2\eta}{2-\eta} \|\mathbf{A}\|_F^2. \\
319 \quad &\leq \frac{\sqrt{M}\|Q\|_2(1+2k\delta) \left(\frac{\delta}{2}\right)^{-\tau}}{\sqrt{N}\|Q\|_2 |\mathcal{K}(\delta+a)|} \frac{1}{\sqrt{c}} \frac{2\eta}{2-\eta} \|\mathbf{A}\|_F^2 \\
320 \quad (35) \quad &= C(\delta, a, \tau, k, M, N) \frac{1}{\sqrt{c}} \frac{2\eta}{2-\eta} \|\mathbf{A}\|_F^2,
\end{aligned}$$

321 where  $C(\delta, a, \tau, k, M, N) = \left(\frac{\delta}{2}\right)^{-\tau} \frac{(1+2k\delta) \sqrt{M}}{|\mathcal{K}(\delta+a)| \sqrt{N}}$ .  $\square$

322 This result indicates that the error of kernel compression also depends on the  
323 diameter-distance ratio of the well-separated target and source points. It is impor-  
324 tant to point out that in the constant  $C(\delta, a, \tau, k, M, N)$ , wavenumber parameter  $k$   
325 is critical to the compression error. A typical example is the Green's function for  
326 Helmholtz equation, for which the high frequency problem is still a challenge for any  
327 kernel compression algorithm.

328 **4. Hierarchical matrix ( $\mathcal{H}$ -matrix) structure for general data sets.** For  
329 general cases when target and source are not well-separated, in fact typically they  
330 belong to the same set, we will partition the whole matrix into blocks, each of which  
331 corresponds to a far-field interaction sub-matrix and will have a low-rank approxi-  
332 mation, thus the kernel compression algorithm can be applied hierarchically at dif-  
333 ferent scales. The resulting method will be termed "hierarchical random compression  
334 method (HRCM)".

#### 335 4.1. Review of $\mathcal{H}$ -matrix.

336 **DEFINITION 4.1** (Hierarchical matrices ( $\mathcal{H}$ -matrix)). *Let  $I$  be a finite index set*  
337 *and  $P_2$  be a (disjoint) block partitioning (tensor or non-tensor) of  $I \times I$  and  $K \in \mathbb{N}$ .*  
338 *The underlying field of the vector space of matrices is  $\mathbb{K} \in \{\mathbb{R}, \mathbb{Z}\}$ . The set of  $\mathcal{H}$ -*  
339 *matrix induced by  $P_2$  is*

$$340 \quad (36) \quad \mathcal{M}_{\mathcal{H}, K} := \{\mathbf{M} \in \mathbb{K}^{I \times I} : \text{each block } \mathbf{M}^b, b \in P_2, \text{ satisfies } \text{rank}(\mathbf{M}^b) \leq K\}.$$

341 **Remarks:** (1) The index set  $I$  can be the physical coordinates of target/source  
342 points; we denote a matrix  $\mathbf{A}$  as R- $K$  matrix if  $\text{rank}(\mathbf{A}) \leq K$ ; (2) A specific  $\mathcal{H}$ -matrix  
343 is defined through 4.1 recursively. Full definition, description, and construction of  $\mathcal{H}$ -  
344 matrices are given in details in [13, 14]. Two simple examples are given as follows;  
345 (3) Since  $\mathcal{H}$ -matrix is recursive, we always assume  $\mathbf{A} \in \mathbb{K}^{N \times N}$  and  $N = 2^p$  for the  
346 following context.

347 **Example One:**  $\mathbf{A} \in \mathcal{M}_{\mathcal{H}, K}$  if either  $2^p = K$  or it has the structure

$$348 \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \text{ with } \mathbf{A}_{11}, \mathbf{A}_{22} \in \mathcal{M}_{\mathcal{H}, K} \text{ and R-}K \text{ matrices } \mathbf{A}_{12}, \mathbf{A}_{21}.$$

349 This is the simplest  $\mathcal{H}$ -matrix. Such a matrix with three levels of division is visualized  
 350 in the left of Fig. 2. All blocks are R- $K$  matrices, except the smallest ones.

351 The next example is more complicated and it includes another two recursive  
 352 concepts of neighborhood matrix ( $\mathcal{N}$ -type and  $\mathcal{N}^*$ -type). Construction of this  $\mathcal{H}$ -  
 353 matrix includes three steps.

354 **Example Two:** (i) Neighborhood matrix of  $\mathcal{N}$ -type, or  $\mathcal{M}_{\mathcal{N},K}$ : if either  $2^p = K$  or  
 355 it has the structure

$$356 \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \text{ with } \mathbf{A}_{21} \in \mathcal{M}_{\mathcal{N},k} \text{ and R-}K \text{ matrices } \mathbf{A}_{11}, \mathbf{A}_{12}, \mathbf{A}_{22}.$$

357 (ii) Neighborhood matrix of  $\mathcal{N}^*$ -type, or  $\mathcal{M}_{\mathcal{N}^*,K}$ :  $\mathbf{A} \in \mathcal{M}_{\mathcal{N}^*,k}$  if  $\mathbf{A}^* \in \mathcal{M}_{\mathcal{N},K}$ .

358 (iii)  $\mathbf{A} \in \mathcal{M}_{\mathcal{H},K}$  if either  $2^p = K$  or it has the structure

$$359 \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \text{ with } \mathbf{A}_{11}, \mathbf{A}_{22} \in \mathcal{M}_{\mathcal{H},K}, \mathbf{A}_{12} \in \mathcal{M}_{\mathcal{N},K}, \text{ and } \mathbf{A}_{21} \in \mathcal{M}_{\mathcal{N}^*,K}.$$

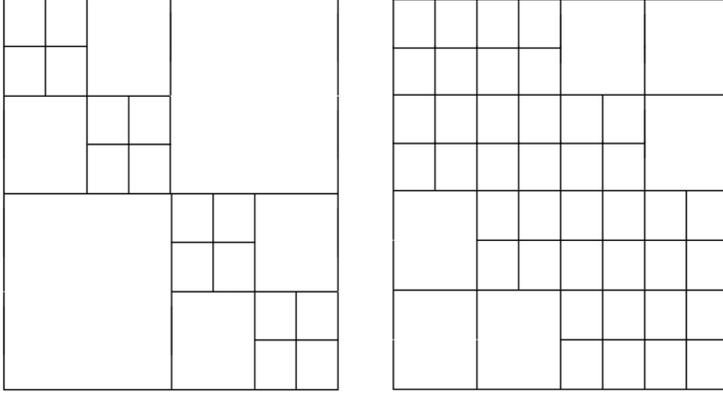


FIG. 2. Two examples of  $\mathcal{H}$ -matrix. Left: example one; right: example two.

360

361 Visualization of the second example with three levels of division is given in the  
 362 right subfigure of Fig. 2. Similarly, only the larger blocks are R- $K$  matrices.

363 With such a decomposition, matrix-vector product will be performed as

$$364 \quad (37) \quad \mathbf{Ax} = \mathbf{A}_{11}\mathbf{x}_1 + \mathbf{A}_{12}\mathbf{x}_2 + \mathbf{A}_{21}\mathbf{x}_1 + \mathbf{A}_{22}\mathbf{x}_2.$$

365 where  $\mathbf{x}^T = (\mathbf{x}_1^T, \mathbf{x}_2^T)$ . For  $\mathbf{A}$  in Example One, random compression can be immedi-  
 366 ately applied to  $\mathbf{A}_{12}\mathbf{x}_2$  and  $\mathbf{A}_{21}\mathbf{x}_1$ , while recursive division and random compression  
 367 need to be implemented on  $\mathbf{A}_{11}\mathbf{x}_1$  and  $\mathbf{A}_{22}\mathbf{x}_2$ , until a preset minimum block is reached  
 368 where direct matrix-vector multiplication is used.

369 But in Example Two, none of the four terms in Eq. (37) is R- $K$  matrix ready for  
 370 compression. Each  $\mathbf{A}_{ij}$  needs to be further divided and investigated. For instance,  
 371  $\mathbf{A}_{12} \in \mathcal{M}_{\mathcal{N},K}$  by definition (iii), then by (i),  $\mathbf{A}_{12_{11}}$ ,  $\mathbf{A}_{12_{12}}$ , and  $\mathbf{A}_{12_{22}}$  are R- $K$  matrices  
 372 and the compression algorithm can be applied. In contrast,  $\mathbf{A}_{12_{21}}$  needs to be further  
 373 divided.

374 Matrices in Example One and Two can be understood as interactions of target/-  
 375 source points along a 1D geometry. Example One indicates that interactions of all

subset of  $I$  can be approximated as low-rank compressions except self-interactions of the subsets, which requires further division. While the blocks in Example Two require further division for both self-interacting and immediate neighboring subsets of  $I$ .

Structures of  $\mathcal{H}$ -matrices for target/source points in high-dimensional geometry are much more complicated. It is difficult to partition them into blocks as shown in Fig 2. Instead, we construct the  $\mathcal{H}$ -matrix logically through a partition tree of  $I$  and the concept of admissible clusters.

**4.2. Tree structure of  $\mathcal{H}$ -matrix for two dimensional data.** We illustrate algorithms in two-dimensional (2D) case. Here the dimension refers to the geometry where the target and source points are located instead of the dimension of the kernel function. For simplicity, let  $\Omega = [0, L] \times [0, L]$  and consider a regular grid index

$$(38) \quad I = \{(i, j) : 1 \leq i, j \leq N_1\}, \quad N_1 = 2^p.$$

Each index  $(i, j) \in I$  is associated with the square

$$(39) \quad X_{ij} : \{(x, y) : (i-1)h \leq x \leq ih, (j-1)h \leq y \leq jh\}, \quad h = L/N_1,$$

where a certain amount of target/source points are located. The partitioning  $T(I)$  uses a quadtree, with children (or leaves):

$$(40) \quad t_{\alpha, \beta}^l := \{(i, j) : 2^{p-l}\alpha + 1 \leq i \leq 2^{p-l}(\alpha + 1), 2^{p-l}\beta + 1 \leq j \leq 2^{p-l}(\beta + 1)\},$$

with  $0 \leq \alpha, \beta \leq 2^l - 1$  belong to level  $0 \leq l \leq p$ .

We consider a target quadtree  $I_t$  and a source quadtree  $I_s$ , which could be same or different. Then blocks of interaction matrix are defined as a block  $b = (t_1, t_2) \in T(I_t \times I_s)$ , where  $t_1, t_2 \in I_t \times I_s$  belong to the the same level  $l$ . Then, following Eq. (2)-(3), we define the diameters and distances of  $t_1$  and  $t_2$ , and the admissibility condition

$$(41) \quad \max\{\text{diam}(t_1), \text{diam}(t_2)\} \leq \eta \text{dist}(t_1, t_2), \quad 0 < \eta < 1,$$

for the block  $b = (t_1, t_2)$ . A block  $b$  is called admissible or an admissible cluster if either  $b$  is a leaf or the admissibility condition holds. If  $b$  is admissible, no matter how many points are in  $t_1$  and  $t_2$ , the block matrix  $\mathbf{M}^b$ , consisting of entries from the original matrix  $\mathbf{M}$  with row indices  $t_1$  and column indices  $t_2$ , has rank up to  $K$ , thus low rank approximation algorithms are used. Otherwise, both  $t_1$  and  $t_2$  will be further partitioned into children to be further investigated. And the above process is implemented, recursively.

Figure 3(a) shows the index set  $I$ , where black solid, gray solid and gray lines are for partition at level  $l = 1, 2, 3$ , respectively. For different values of  $\eta$ , admissible clusters are different for a given child. For the square marked with star in Fig 3(a), if  $\eta = \sqrt{2}/2$ , the non-admissible clusters are itself and the eight immediate surrounding squares. While for  $\eta = 1/2$ , any squares within the the red lines are non-admissible.

Figure 3(b) displays the quadtree that divides each square. Four children of each branch are labeled as 0, 1, 2, 3 and ordered counter-clock wisely. If there are  $N = 4^p$  target (source) points, the depth of the target (source) tree is  $p - p_0$ , where  $p_0$  is the number of points in each leaf, or direct multiplication is performed when matrix size is down to  $4^{p_0} \times 4^{p_0}$ .

Figure 4 illustrates the admissible and non-admissible clusters. Initially the 2D set is partitioned into four subdomain  $A, B, C$  and  $D$ . Any two of the subdomains

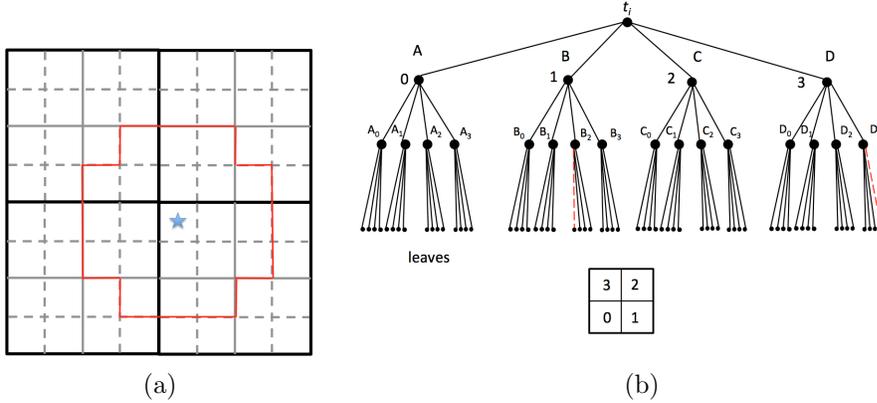


FIG. 3. Illustration of the index set  $I$ : (a) admissible cluster for the starred square; (b) quadtree structure.

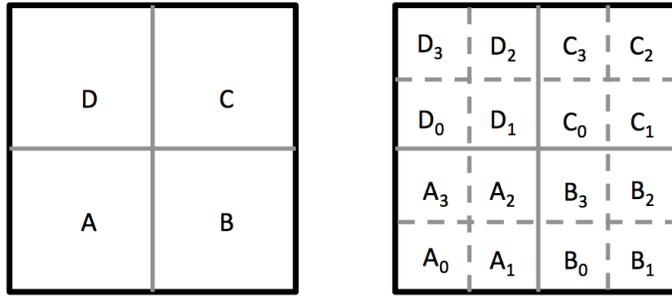


FIG. 4. Illustration of the partitioning of admissible clusters

419 are non-admissible for  $\eta = \sqrt{2}/2$ . Then each of them are further divided into four  
 420 children domain, as label on the right of Fig 4, among which the interactions are  
 421 examined. For example, the interactions of  $A$  and  $B$  can be viewed as the sum of  
 422 interactions of  $A_i$  and  $B_j$ ,  $i, j = 0, 1, 2, 3$ . If we take  $\eta = 1/2$ , only  $A_3$  and  $B_1$ , and  
 423  $A_0$  and  $B_2$  are admissible clusters. But if  $\eta = \sqrt{2}/2$ , only  $A_2$  and  $B_0$ , and  $A_1$  and  
 424  $B_3$  are non-admissible pairs. For both values of  $\eta$ , only  $A_2$  and  $C_0$  are non-admissible  
 425 clusters in the interactions of  $A$  and  $C$ . Self-interactions, such as interaction between  
 426  $A$  and  $A$ , can be viewed as the same process of interactions among  $A, B, C$  and  $D$ ,  
 427 but for  $A_0, A_1, A_2$ , and  $A_3$ .

428 Direct and low rank approximation of matrix-vector multiplications are imple-  
 429 mented on the quadtree. To perform the algorithms, all the index in  $I$  is ordered  
 430 as  $i = \{0, 1, \dots, N - 1\}$  with  $N = 4^p$ . Note that there are totally  $4^{p-l}$  points in each  
 431 child/leaf at level  $l$ . Matrix column (row) sampling is achieved through sampling of  
 432 children from level  $l + 1$  to level  $p$ . For example, cluster  $b = (B_2, D_3)$  in Fig. 3 (b)  
 433 is admissible and assume it is at level  $l$ . If we generate  $s_i \in \{0, 1, 2, 3\}$ ,  $i = l + 1, \dots, p$   
 434 randomly and choose only the  $s_i$ -th child of  $B_2$  (red dash line) at  $i$ -th level, then one  
 435 row sampling of the corresponding (block) kernel matrix is accomplished assuming  
 436  $B_2 \in I_t$ . Similarly column sampling is the child-picking process on  $D_3 \in I_s$ .

The full HRCM algorithms are summarized as in the following:

---

**Algorithm 2** HRCM: Direct product on source and target quadtrees

---

subroutine name: DirectProduct(\*target, \*source, int level)

**Input:** Root pointers of source and target quadtree information for  $\mathbf{A}$  and  $\mathbf{x}$ . Current level  $l$  and maximum level  $p$ .

**Output:** Product  $\mathbf{y} = \mathbf{Ax}$ , where  $\mathbf{y}$  is stored in the target quadtree.

```

if level == maxlevel then
    perform and scalar product, and return
else
    for j = 0; j < 4; j++ do
        for i = 0; i < 4; i++ do
            DirectProduct(target->child[j], source->child[i], level + 1)
        end for
    end for
end if

```

---

437

---

**Algorithm 3** HRCM: Low-rank product on source and target quadtrees

---

subroutine name: LowRankProduct(\*target, \*source, int level)

**Input:** Root pointers of source and target quadtree information for  $\mathbf{A}$  and  $\mathbf{x}$ . Current level  $l$  and maximum level  $p$ .

**Output:**  $\mathbf{y} \approx \sum_{t=1}^{t_0} \sigma_t \mathbf{U}_t \mathbf{V}_t^* \mathbf{x}$ , where  $\mathbf{y}$  is stored in the target quadtree.

1. Column sampling: On the source tree, pick the “random path” from level  $l$  to maxlevel  $p$  by only randomly choosing one child from each level;
  2. Row sampling: On the target tree, pick the “random path” from level  $l$  to maxlevel  $p$  by only randomly choosing one child from each level;
  3. Extract matrix entries from the source and target tree by the column/row sampling; perform Algorithm 1
  4. Instore  $\mathbf{y}$  into the target tree with root \*target.
- 

438 **4.3. Efficiency analysis.** It is easy to perform efficiency analysis of the hierar-  
439 chical kernel compression method by constructing an interaction pattern tree. With  
440  $\eta = \sqrt{2}/2$ , all the non-admissible clusters can be classified into three interaction pat-  
441 terns: the self-, edge-contact, and vertex-contact interactions. And these interactions  
442 are labeled as S, E and V as shown in Fig. 5. Assume it is currently level  $l$  and those  
443 target and source boxes need to be further divided into four children in level  $l + 1$ .  
444 In level  $l + 1$ , those children form 16 interactions that fall into the S, E, V patterns.  
445 It is easy to check that from level  $l$  to level  $l + 1$ , as displayed in Fig 5, S-interaction  
446 forms 4 S-, 8 E- and 4 V-interactions at level  $l + 1$ . On the other hand, E-interaction  
447 forms 2 E-interaction, 2 V-interactions and 12 admissible clusters for which low-rank  
448 approximation (LR) applies. Additionally, V-interaction forms 1 V-interactions and  
449 15 LR approximations.

450 Complexity of all direct calculations. We assume the direct computation is im-  
451 plemented when the matrix scale is down to  $4^{p_0} \times 4^{p_0}$ . So we start from S-interaction  
452 as the root at level  $l$  and just need to count how many E- and V-interactions at level  
453  $p - p_0 - 1$  are generated. From Fig. 5 we can see that at level  $l + 1$ , S generates eight Es  
454 and four Vs. And afterwards, from level  $l + 1$  to level  $p - p_0 - 1$ , each E generates two  
455 Es and two Vs, and each V generates just one E. Thus, the total number of generated

---

**Algorithm 4** HRCM:  $\mathcal{H}$ -matrix product on source and target quadrees

subroutine name: HmatrixProduct(\*target, \*source, int level)

**Input:** Root pointers of source and target quadtree for  $\mathbf{A}$  and  $\mathbf{x}$ . Current level  $l$  and maximum level  $p$ .

**Output:**  $\tilde{\mathbf{y}} \approx \mathbf{y} = \mathbf{A}\mathbf{x}$ , where  $\tilde{\mathbf{y}}$  is stored in the target quadtree.

**if** matrix small enough **then**

DirectProduct(\*target, \*source, level)

**else**
**if** clusters rooted from \*target, \*source are admissible **then**

LowRankProduct(\*target, \*source, level)

**else**
**for**  $j = 0; j < 4; j++$  **do**
**for**  $i = 0; i < 4; i++$  **do**

HmatrixProduct(target-&gt;child[j], source-&gt;child[i], level + 1)

**end for**
**end for**
**end if**
**end if**


---

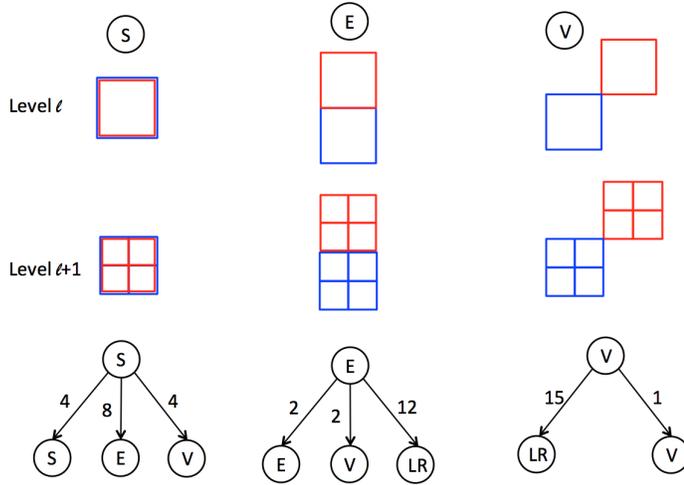


FIG. 5. Evolution of non-admissible clusters.  $S$ : self-interaction clusters;  $E$ : clusters touch by edge;  $V$ : clusters touch by vertex;  $LR$ : admissible clusters with low-rank approximation.

456 E and V and level  $p - p_0 - 1$  from the S at level  $l$  is

457 (42) 
$$\underbrace{8 \cdot 2^{p-p_0-1-l-1}}_{\text{generated Es}} + \underbrace{8 \cdot 2 \cdot 1^{p-p_0-1-l-2} + 4 \cdot 1^{p-p_0-1-l-1}}_{\text{generated Vs}}.$$

458 Since there are  $4^l$  S-interactions at level  $l$ , the total complexity for performing direct com-  
459 putation is

460 (43) 
$$16 \cdot O((4^{p_0})^2) \left[ \sum_{l=0}^{p-p_0-3} 4^l (8 \cdot 2^{p-p_0-l-2} + 8 \cdot 2 + 4) + 4^{p-p_0-2} \cdot 16 \right] = O(4^p) = O(N).$$

461

462 Complexity of all low-rank compressions. We follow the similar strategy and count  
 463 the low rank interactions (LR) generated from S at level  $l$  to level  $p - p_0 - 1$ . Again  
 464 we start from S as the root at level  $l$  and according to Fig. 5, the resulting E and  
 465 V start to generate LR at the  $(l + 2)$ -th level and continue to the last level. At level  
 466  $l + 1$ , eight Es are generated, each of which generates 12 LRs at level  $l + 2$ ; at the  
 467 same time, four Vs are generated and each of them generates 15 LRs at level  $l + 2$ .  
 468 Then at level  $l + 2$ , totally  $8 \times 12 + 4 \times 15$  LRs are generated. Recall at level  $l$  the  
 469 complexity of performing LR is  $O(4^{p-l})$ , so the complexity of performing LR starting  
 470 from S at level  $l$  is

$$\begin{aligned}
 471 \quad & \underbrace{(8 \cdot 12 + 4 \cdot 15)O(4^{p-l-2})}_{\text{from E and V at } l+2 \text{ level}} + \underbrace{\sum_{l'=l+3}^{p-p_0-1} 8 \cdot 2^{l'-l-2} \cdot 12 \cdot O(4^{p-l'})}_{\text{from all the E to the bottom}} \\
 472 \quad (44) \quad & + \underbrace{\sum_{l'=l+3}^{p-p_0-1} 8 \cdot 2 \cdot 1^{l'-l-3} \cdot 15 \cdot O(4^{p-l'}) + \sum_{l'=l+3}^{p-p_0-1} 4 \cdot 1^{l'-l-2} \cdot 15 \cdot O(4^{p-l'})}_{\text{from all the V to the bottom}},
 \end{aligned}$$

473 where the latter three terms in (44) account for LR generated by E and V from level  
 474  $l + 3$  all down to level  $p - p_0 - 1$ . Note that the first item in (44) dominates and there  
 475 are  $4^l$  S-interactions in level  $l$ . The total complexity of low-rank approximation in the  
 476 HRCM is in the order of

$$477 \quad (45) \quad \sum_{l=0}^{p-p_0-1} 4^l O(4^{p-l-2}) = O(p \times 4^p) = O(N \log N).$$

478 Combining Eqs. (43) and (45), the complexity of the HRCM is  $O(N \log N)$ .

479 **5. Numerical results.** In this section, we present the statistical analysis, accu-  
 480 racy and efficiency of the proposed HRCM in 2D computations. For all the following  
 481 simulations, we take  $N = 4^p$  target/source points uniformly distributed in square do-  
 482 mains of length  $L$ . In the random kernel compression Algorithm 1, the total numbers  
 483 of sampled columns and rows are denoted as  $c = r = K = 4^{p_c}$ , respectively.

484 **5.1. Uncertainty quantification.** In this section, we study the statistical prop-  
 485 erties of the single-realization and multiple-realization of HRCM in kernel summa-  
 486 tions. In order to perform the algorithm with large number of replications, we take  
 487 a matrix with moderate size  $N = 4^7$  ( $p = 7$ ), and pick  $K = 16, 64$ , and  $256$ , or  
 488  $p_c = 2, 3, 4$ . The 16384 target and source points are distributed in two squares  
 489 with length  $L = 8$  that are separated 16 apart. The kernel function is taken as  
 490  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = e^{-ikR}/R$  with  $k = 0.5$ .

491 Single-realization of HRCM: Given sample size  $K$ , one can apply the HRCM al-  
 492 gorithm by one realization to obtain the compressed matrix and continue to next  
 493 implementations such as kernel summation. This treatment gives the best algo-  
 494 rithm efficiency but we concern the reliability. The algorithm has been replicated  
 495 by  $N_s = 20,000$  times and the corresponding relative errors of single realization are  
 496 displayed as histograms in Fig. 6 (a) for  $K = 16$ , Fig. 6 (b) for  $K = 64$ , and Fig. 6  
 497 (c) for  $K = 256$ .

498 Define the sample mean error of single realization (sMESR) as

$$499 \quad (46) \quad \text{sMESR} = \frac{1}{N_s} \sum_{s=1}^{N_s} \frac{\|(\mathbf{A} - \Pi_s(K)\mathbf{A})\mathbf{x}\|_2}{\|\mathbf{A}\mathbf{x}\|_2}.$$

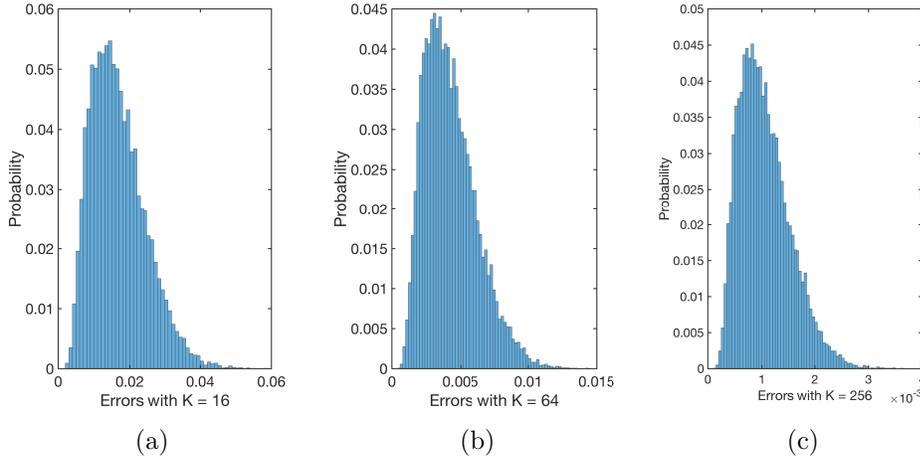


FIG. 6. Histograms of relative errors of 20,000 single-realizations of matrix compression. (a) With  $K = 16$ , sample mean  $1.699 \times 10^{-2}$ , variance  $6.002 \times 10^{-5}$ , 95<sup>th</sup> percentile  $3.118 \times 10^{-2}$ ; (b) With  $K = 64$ , sample mean  $4.213 \times 10^{-3}$ , variance  $3.694 \times 10^{-6}$ , 95<sup>th</sup> percentile  $7.84 \times 10^{-3}$ ; (c) With  $K = 256$ , sample mean  $1.056 \times 10^{-3}$ , variance  $2.250 \times 10^{-7}$ , 95<sup>th</sup> percentile  $1.928 \times 10^{-3}$ .

500 where  $\Pi_s(K)$  is the  $s$ -th realization of the compression projector. The quantitative  
 501 results corresponding to Fig. 6 are summarized in Table 1. We conclude that a single

TABLE 1  
 Statistics of single realization of HRCM with 20,000 replications.

	sMESR	Variance	95 <sup>th</sup> percentile
$K = 16$	1.699E-2	6.002E-5	3.118E-2
$K = 64$	4.213E-3	3.694E-6	7.840E-3
$K = 256$	1.056E-3	2.250E-7	1.928E-3

502 realization of HRCM can provide reliable results: for very small sample size  $K$ , it  
 503 offers two-digits relative errors with high confidence. And when  $K$  is increased, the  
 504 HRCM shows convergence with respect to  $K$  at the same percentile of confidence.

505 **Random versus Deterministic sampling:** Secondly, we illustrate the importance of  
 506 *random sampling*, although a single-realization of HRCM can be trusted. In Algorithm  
 507 4.2, sampling is achieved by randomly choosing one child from the target/source tree.  
 508 For comparison, we also perform the HRCM by deterministically picking a specific  
 509 child each time. This is equivalent to sample columns/rows from the matrix evenly  
 510 according to indices, e.g., the  $iN/K$ -th columns/rows ( $i = 0, 1, 2, \dots, K - 1$ ) are picked.  
 511 The corresponding relative errors of single-realization of HRCM with such sampling  
 512 strategies are listed in Table 2. Patterns I, II, III, and IV correspond to the situations  
 513 that each time the  $j$ -th ( $j = 0, 1, 2, 3$ ) child is picked. It clearly indicates that the  
 514 errors are larger than the sMESR and even above the 95-th percentile of the HRCM  
 515 with random samples. So we can conclude that it is statistically significant that  
 516 deterministic sampling will yield larger errors in the HRCM.

517 **Multiple-realization of HRCM:** To improve the accuracy of HRCM with a given  
 518  $K$ , one can also apply it multiple times, take average of the compressed matrix and  
 519 then implement the kernel summation. At last, we present the efficiency of the

TABLE 2  
*Errors of the HRCM with some deterministic samples*

Patterns	I	II	III	IV
$K = 16$	3.080E-2	9.515E-2	4.320E-2	9.883E-2
$K = 64$	1.600E-2	4.218E-2	1.905E-2	4.140E-2
$K = 256$	7.323E-3	1.904E-2	8.362E-3	1.956E-2

520 multiple-realization of HRCM. Define the error of multiple realization (EMR):

$$521 \quad (47) \quad \text{EMR} = \frac{\|(\mathbf{A} - \frac{1}{N_s} \sum_{s=1}^{N_s} \Pi_s(K)\mathbf{A})\mathbf{x}\|_2}{\|\mathbf{A}\mathbf{x}\|_2},$$

522 we calculate the EMR with different number of  $N_s$  (from  $N_s = 5$  to 10,000) and  
 523 display the results in Fig. 7. It can be concluded that the EMR converges in the scale  
 of  $1/\sqrt{N_s}$  and reaches an equilibrium state (expectation value) depending on  $K$ .

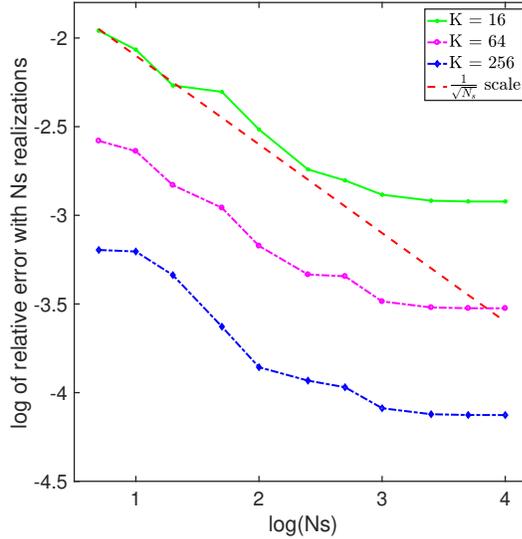


FIG. 7. *Convergence of EMR with respect to  $N_s$ .*

524

525 By comparing statistics of the sMESR and EMR we can see that one can trust the  
 526 one-time realization of HRCM, which is especially useful for some situations such as  
 527 the positions of target/source points are dynamic so everything needs to be computed  
 528 on the fly. On the other hand, averaged matrix compressions with multiple realizations  
 529 of HRCM will improve accuracy with convergence order roughly as  $1/\sqrt{N_s}$ . Since the  
 530 HRCM is applied multiple times, extra computational efforts are committed. But this  
 531 treatment is useful if the compressed matrix can be stored and used repeatedly, such  
 532 as iterative method for solving linear systems.

533 In the following numerical results regarding accuracy and efficiency of HRCM, we  
 534 will only consider the single realization.

535 **5.2. Accuracy and efficiency for well-separated sets.** First, we investigate  
 536 the decay of singular values for the kernel matrix formed by the well-separated target

537 and source points. The kernel function is taken as  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = e^{-0.01R}/R$  with  $L =$   
 538 8, and the SVD of relatively small matrices with  $N = 1024$  are calculated, with  
 539 diameter/distance ratios being  $\eta = 0.5, 0.36,$  and  $0.25$ . Logarithmic values of the first  
 540 18 singular values for each case are displayed in Fig. 8 (a). It clearly shows that  
 541 singular values of the matrix decay faster as the corresponding target and source sets  
 542 are further away. For the fixed  $\eta = 0.5$ , approximated singular values from randomly  
 543 sampled matrix  $\mathbf{C}_r \in \mathbb{R}^{K \times K}$ , with  $K = 4^2, 4^3$  and  $4^4$  are presented in Fig. 8 (b). The  
 544 relative error with respect to the largest singular value is small enough after several  
 singular values even for a very small amount of samples.

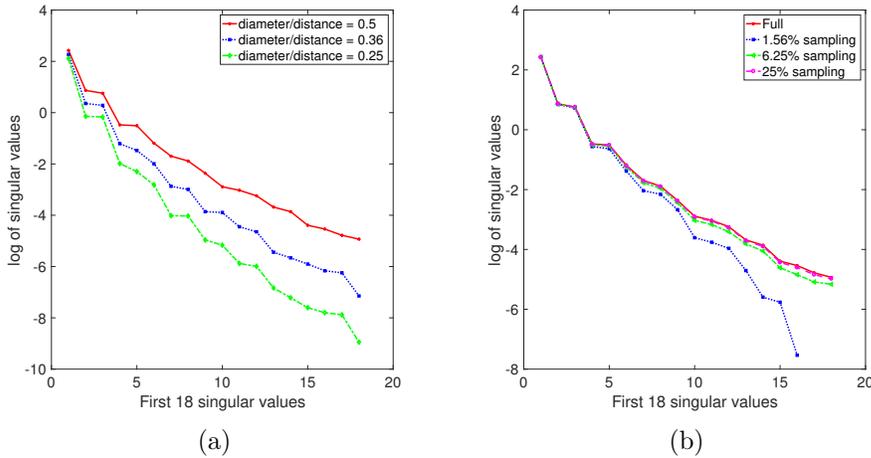


FIG. 8. (a) Singular values for matrices from well-separated points (1024 targets and 1024 source points) with various cluster diameter/distance ratios; (b) singular values comparison between different amounts of matrix samplings:  $K = 16, 64, 256$  against  $N = 1024$ .

545

546 Next, we check the the algorithm accuracy. A total of  $N$  target points and source  
 547 points are uniformly assigned in two  $8 \times 8$  boxes with centers 16 units apart. Then, the  
 548 direct multiplication (1) and Algorithm 1 are performed with parameter  $c = r = K$   
 549 and  $\epsilon = 1.0 \times 10^{-8}$ . For each comparison, sMESR and variance of errors are calculated  
 550  $N_s = 20$ . Errors and variances for kernels  $\mathcal{K}(\mathbf{r}, \mathbf{r}') = \log(\sqrt{(x-x')^2 + (y+y')^2}) -$   
 551  $\log(\sqrt{(x-x')^2 + (y-y')^2})$  and  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \exp(-0.01R)/R$  are displayed in Tables  
 552 3-4, respectively, with various  $N$  and  $K$  and  $\eta = 0.5$ . We can clearly observe the  
 553 convergence of the mean errors against  $K$  in these tables.

554 Based on the numerical results from Tables 3-4, we conclude that that, given  
 555 the fixed diameter/distance ratio of the target/source point sets, the accuracy of the  
 556 low-rank compression algorithm does not depend significantly on the total number  
 557  $N$  but on the sample number  $K$ . It suggests that in computational applications, as  
 558 long as two boxes are admissible clusters, it does not matter how many target/source  
 559 points are in them, the algorithm accuracy is mainly controlled by the diameter/ratio  
 560 distance and number of row/column samples.

561 Table 5 summarizes the corresponding computational time in seconds for the  
 562 matrix-vector product, for direct computation and the low-rank compression method.  
 563 If the target and source points are well-separated, the algorithm is very efficient and  
 564 the computational time is linear both in sample size  $K$  and matrix size  $N$ . CPU times  
 565 for the two kernel functions are similar so only one of them is presented.

TABLE 3

Errors and variances for a pair of well-separated target/source point sets. Kernel function  $\mathcal{K}(\mathbf{r}, \mathbf{r}') = \log(\sqrt{(x-x')^2 + (y+y')^2}) - \log(\sqrt{(x-x')^2 + (y-y')^2})$

	$N = 1,024$	$N = 4096$	$N = 16,384$	$N = 65,536$	$N = 262,144$
$K = 16$					
sMESR	2.79E-2	3.07E-2	3.51E-2	3.78E-2	4.01E-2
Variance	3.58E-4	4.13E-4	5.38E-4	6.32E-4	6.95E-4
$K = 64$					
sMESR	8.06E-3	8.54E-3	9.70E-3	9.84E-3	1.01E-2
Variance	5.46E-6	5.89E-6	4.92E-6	6.27E-6	6.18E-6
$K = 256$					
sMESR	2.25E-3	2.39E-3	2.52E-3	2.90E-3	2.75E-3
Variance	4.76E-6	4.71E-6	5.37E-6	5.63E-6	5.86E-6

TABLE 4

Errors and variances for a pair of well-separated target/source point sets. Kernel function  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \exp(-0.01R)/R$

	$N = 1,024$	$N = 4096$	$N = 16,384$	$N = 65,536$	$N = 262,144$
$K = 16$					
sMESR	2.67E-2	3.39E-2	3.07E-2	3.02E-2	3.51E-2
Variance	7.51E-4	4.44E-4	6.28E-4	6.62E-4	9.95E-4
$K = 64$					
sMESR	7.46E-3	7.58E-3	6.70E-3	8.51E-3	8.40E-3
Variance	1.41E-5	2.78E-5	3.89E-5	4.17E-5	4.86E-5
$K = 256$					
sMESR	1.62E-3	1.85E-3	1.92E-3	2.10E-3	2.30E-3
Variance	1.76E-6	1.47E-6	1.37E-6	3.53E-6	3.53E-6

566 Figure 9 shows the algorithm error against the diameter-distance ratios with  
 567  $N = 262,144$  and different values of  $K$ . As expected, the relative error decays as  
 568  $\eta$  increases. This graph is for kernel  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \exp(-0.01R)/R$ , the one for kernel  
 569  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \log(R)$  is similar.

TABLE 5

CPU time (second) comparison for well-separated target and source points.

	$N = 1,024$	$N = 4096$	$N = 16,384$	$N = 65,536$	$N = 262,144$
Direct	0.047	0.75	12	204	3,264
$K = 16$	0.01	0.039	0.16	0.625	2.5
$K = 64$	0.04	0.16	0.66	2.6	12.0
$K = 256$	0.18	0.8	2.7	12.5	47.0

570 **5.3. Accuracy and efficiency for a common source and target set .** Next,  
 571 we test the accuracy and efficiency of the HRCM for target and source points in the  
 572 same set. In total  $N = 4^p$  points with  $p = 6, 7, 8, 9, 10, 11$  are uniformly distributed in  
 573 the domain  $[0, 8] \times [0, 8]$ . For best computation efficiency, we only present the results  
 574 with  $K = 16$  and  $64$ . Note that it has been found that once sample size  $K$  is fixed,  
 575 the accuracy of the low-rank compression algorithm for a pair of admissible cluster

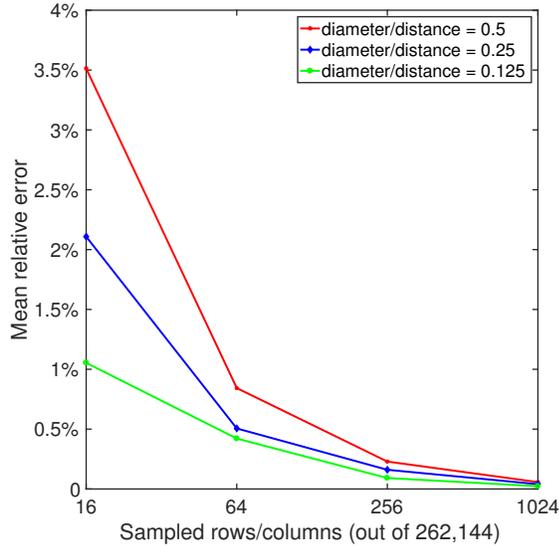


FIG. 9. Error of low-rank compression algorithm against diameter-distance ratio  $\eta$ .

576 does not change too much regardless of number of points in them.

577 The error for kernel  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \exp(-0.01R)/R$  with different  $K$  and  $N$  are  
 578 summarized in Table 6. Note these values are generally smaller than those in Table 4.  
 579 Because Table 4 is for a single pair of well-separated target/source sets with diameter-  
 580 distance  $\eta = 0.5$ . But in the HRCM there exist a mixture of  $\eta$  with  $\eta = 0.5$  as the  
 largest value. Similar convergence of error with respect to  $K$  is shown in the table.

TABLE 6  
 Accuracy of the HRCM for kernel summation with  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \exp(-0.01R)/R$

Matrix size	$N = 16, 384$	$N = 65, 536$	$N = 262, 144$	$N = 1, 048, 576$
$k = 16$				
sMESR	2.87E-3	3.32E-3	3.46E-3	3.53E-3
Variance	6.82E-7	7.32E-7	7.65E-7	8.30E-7
$k = 64$				
sMESR	6.09E-4	7.43E-4	6.26E-4	7.32E-4
Variance	7.03E-8	6.49E-8	6.63E-8	7.56E-8

581

582 The efficiency of the HRCM with such a kernel function is presented in Fig. 10 as  
 583 log-log CPU time and matrix size  $N$ . For better comparison, the curves of CPU time  
 584 for the direct method and an ideal  $O(N \log N)$  scale are also displayed. For HRCM  
 585 with  $K = 16$  and  $K = 64$ , the curves are almost parallel to the ideal  $O(N \log N)$   
 586 scale. Combining Fig. 10 and Table 6, we can conclude that the break-even point of  
 587 the HRCM comparing to the direct method with three or four digits in relative error  
 588 is  $N$  slightly larger than  $10^4$  for this screened Coulomb potential. If higher accuracy  
 589 is desired, one needs to increase number of  $K$  and the break-even point will be larger.

590 **5.4. Wave number  $k$  dependence.** In this section, we will study the perfor-  
 591 mance of the HRCM for compressing the kernel matrices for wave interactions. We

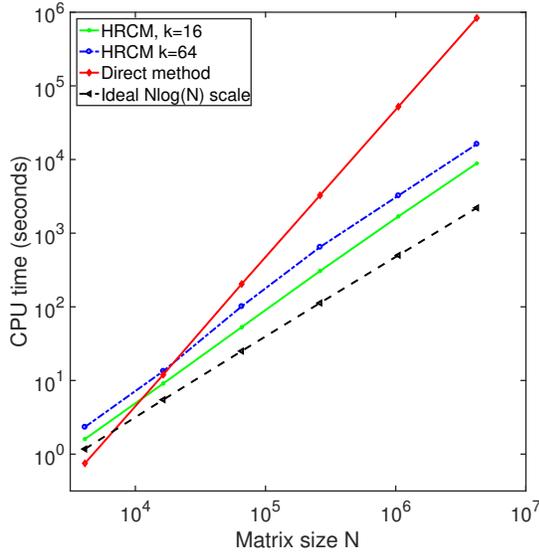


FIG. 10. CPU costs of HRCM with  $K = 16$  and  $K = 64$ , for various matrix sizes. For comparison, the CPU time for direct method is shown in red and the ideal  $N \log(N)$  curve is in black.

592 consider the Green's function for the 3-D Helmholtz equation,  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_i) = \exp(-ikR)/R$ ,  
 593 when the wavenumber  $k$  appears in the error estimate constant  $C$  in (27). In the fol-  
 594 lowing results, we take  $N$  points in the domain  $[0, 2\pi]^2$ , so the wavenumber  $k$  is  
 595 the same as the number of wavelengths along each direction of the domain. For  
 596 small value of  $k$ , the performance of HRCM is similar to the cases of non-oscillatory  
 597 screened Coulomb kernels. Errors and variances of the HRCM for this kernel with  
 598  $k = 0.25, k = 0.5$  and  $k = 1$  are presented in Tables 7, where the relative errors are  
 599 all below 0.2% for small sample size  $K = 16$  or 64.

600 For wave problems with larger wavenumber  $k$ , one needs  $N$  large enough and  
 601 proportional to  $k$  in order to resolve the wave structure of the numerical solution of  
 602 the Helmholtz equation. In the following simulations, we use 8 points per wavelength  
 603 in each directions of the domain and implement the HRCM for wave number  $k$  up  
 604 to 64. The relation of errors and wavenumber  $k$  is revealed in Fig. 11 (a). For each  
 605 fixed sample size  $K$ , the relation is roughly linear as indicated by the error analysis  
 606 in Eq. (35). As a consequences, large sample size  $K$  is required to maintain a desired  
 607 accuracy. For example, with relative error 1%, the sample size  $K = 16$  can only  
 608 handle the wavenumber  $k$  up to 24, and if the relative error 0.3% is desired, one has  
 609 to take the sample size  $K = 256$  for wavenumber  $k = 64$ . On the other hand, large  
 610 sample size will undermine the efficiency of the HRCM. For sample size  $K = 16$  and  
 611 64, the HRCM is more efficient than the direct method, while for  $K = 256$ , the HRCM  
 612 is only superior when  $N$  exceeds  $10^6$ , as shown in Fig. 11 (b).

613 **6. Conclusion and discussion.** Kernel summation in large scale is a common  
 614 challenge in a wide range of applications, from problems in computational sciences  
 615 and engineering to statistical learning. In this work, we have developed a novel hier-  
 616 archical random compression method (HRCM) to tackle this difficulty. The HRCM  
 617 is a fast Monte-Carlo method that can reduce computation complexity from  $O(N^2)$

TABLE 7  
*Accuracy of the HRCM for  $\mathcal{K}(\mathbf{r}_i, \mathbf{r}_j) = \frac{\exp(-ikR)}{R}$  with small  $k = 0.25, 0.5, 1$ .*

Matrix size	$N = 16,384$	$N = 65,536$	$N = 262,144$	$N = 1,048,576$
$K = 16, k = 0.25$				
sMESR	2.56E-3	2.68E-3	2.71E-3	2.89E-3
Variance	6.11E-7	3.56E-7	1.35E-7	1.01E-7
$K = 64, k = 0.25$				
sMESR	5.42E-4	5.51E-4	5.57E-4	5.89E-4
Variance	1.43E-8	7.69E-8	2.18E-9	1.32E-9
$K = 16, k = 0.5$				
sMESR	2.86E-3	2.98E-3	3.17E-3	3.65E-3
Variance	5.97E-7	6.26E-7	6.54E-7	7.83E-7
$K = 64, k = 0.5$				
sMESR	6.91E-4	6.99E-4	7.95E-4	9.01E-4
Variance	2.55E-8	3.61E-8	6.21E-8	7.53E-8
$K = 16, k = 1$				
sMESR	5.36E-3	5.64E-3	6.23E-3	7.47E-3
Variance	2.12E-6	3.26E-6	4.25E-6	6.37E-6
$K = 64, k = 1$				
sMESR	1.12E-3	1.37E-3	1.51E-3	1.86E-3
Variance	3.48E-7	4.16E-7	6.09E-7	1.11E-6

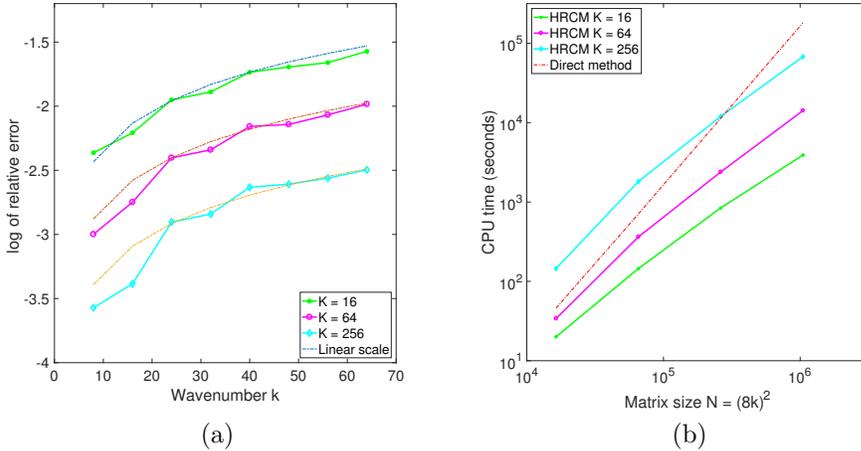


FIG. 11. (a) Relative errors and (b) Efficiency of HRCM against wavenumber  $k$ .

618 to  $O(N \log N)$  for a given accuracy. The method can be readily applied to iterative  
619 solver of linear systems resulting from discretizing surface/volume integral equations  
620 of Poisson equation, Helmholtz equation or Maxwell equations, as well as fractional  
621 differential equations. It also applies to machine learning methods such as regression  
622 or classification for massive volume and high dimensional data.

623 In designing HRCM, we first developed a random compression algorithm for kernel  
624 matrices resulting from far-field interactions, based on the fact that the interaction

625 matrix from well-separated target and source points is of low-rank. Therefore, we  
 626 could sample a small number of columns and rows from the large-scale matrix, in-  
 627 dependent of matrix sizes and only dependent on the separation distance between  
 628 source and target locations, and then perform SVD on the small matrix, resulting in  
 629 a low-rank approximation to the original matrix. A key factor in the HRCM is that a  
 630 uniform sampling, implemented without cost of computing the usual sampling distri-  
 631 bution based on the magnitude of sampled columns/rows, can yield a nearly optimal  
 632 error in the low-rank approximation algorithm. The HRCM is kernel-independent  
 633 without the need for analytic expansions of the kernels. Furthermore, an error bound  
 634 of the algorithm with some assumption on kernel function was also provided in terms  
 635 of the smoothness of the kernel, the number of samples and diameter-distance ratio  
 636 of the well-separated sets.

637 For general source and target configurations, we utilized the concept of  $\mathcal{H}$ -matrix  
 638 to hierarchically divide the whole matrix into logical block matrices, for which the  
 639 developed low-rank compression algorithm can be applied if blocks correspond to low-  
 640 rank far field interactions at an appropriate scale, or they are divided further until  
 641 direct summation is needed. Different from analytic or algebraic FMMs, the recursive  
 642 structure nature of HRCM only executes an one-time, one way top-to-down path along  
 643 the hierarchical tree structure: once a low-rank matrix is compressed, the whole block  
 644 is removed from further consideration, and has no communications with the remaining  
 645 entries of the whole kernel matrix. As the HRCM combines the  $\mathcal{H}$ -matrix structure  
 646 and low-rank compression algorithms, it has an  $O(N \log N)$  computational complexity.

647 Numerical simulations are provided for source and targets in two-dimensional  
 648 (2D) geometry for several kernel functions, including 2D and 3D Green's function  
 649 for Laplace equation, Poisson-Boltzmann equation and Helmholtz equation. As a  
 650 Monte Carlo method, its reliability was analyzed in terms of single-realization and  
 651 multiple-realizations. We concluded that, a single-realization of HRCM is statistically  
 652 reliable. Multi-realization of the algorithm is more accurate but efficiency needs to be  
 653 balanced. In various cases, the mean relative errors of the HRCM against direct kernel  
 654 summation show convergence in terms of number of samples and diameter-distance  
 655 ratios. The computational cost was shown numerically as  $O(N \log N)$ . Additionally,  
 656 the break-even point with direct method is in the order of thousands, with three  
 657 or four digit relative error. At last, the HRCM is implemented for high frequency  
 658 wave problems for Helmholtz equations. As shown by our simulations, the mean  
 659 error is linearly proportional wave number  $k$ , thus it could be significantly large in  
 660 high frequency region. Increasing numbers of sampled columns or rows in the low-  
 661 compression algorithm is one of the ways to reduce the error, but may not be the best  
 662 way. The HRCM algorithm needs to be improved to handle high frequency problem.

663 One direction in the future work is to extend the HRCM for higher dimensional  
 664 data. The next work could be solving volume integral equation (VIE) of Helmholtz  
 665 or Maxwell's equations in 3-D using the HRCM. In this case, target and source points  
 666 are distributed in 3-D geometries so an octree will be used to construct the  $\mathcal{H}$ -matrix,  
 667 but the random compression algorithm for low rank matrix can be applied exactly the  
 668 same way. Efficiency and accuracy of HRCM of the VIE method will be analyzed. For  
 669 even higher dimensional data with dimension  $d > 3$ , the  $d$ -d tree is a more suitable  
 670 space partitioning data structure. However, it is one type of binary trees different  
 671 from quadtrees or octrees. The overall efficiency of the hierarchical structure and  
 672 possible changes in the low rank compression algorithm for ultra-high dimensional  
 673 data require careful investigations.

674

## REFERENCES

- 675 [1] C. R. ANDERSON, *An implementation of the fast multipole method without multipoles*, SIAM  
676 Journal on Scientific and Statistical Computing, 13 (1992), pp. 923–947.
- 677 [2] L. BANJAI AND W. HACKBUSCH, *Hierarchical matrix techniques for low-and high-frequency*  
678 *helmholtz problems*, IMA journal of numerical analysis, 28 (2008), pp. 46–79.
- 679 [3] M. BEBENDORF, *Approximation of boundary element matrices*, Numerische Mathematik, 86  
680 (2000), pp. 565–589.
- 681 [4] W. C. CHEW, E. MICHELSEN, J. SONG, AND J.-M. JIN, *Fast and efficient algorithms in*  
682 *computational electromagnetics*, Artech House, Inc., 2001.
- 683 [5] M. H. CHO, J. HUANG, D. CHEN, AND W. CAI, *A Heterogeneous FMM for Layered Media*  
684 *Helmholtz Equation I: Two Layers in  $\mathbb{R}^2$* , Journal of Computational Physics, in press  
685 (2018).
- 686 [6] E. DARVE, *The fast multipole method i: error analysis and asymptotic complexity*, SIAM  
687 Journal on Numerical Analysis, 38 (2000), pp. 98–128.
- 688 [7] P. DRINEAS, R. KANNAN, AND M. W. MAHONEY, *Fast monte carlo algorithms for matrices i:*  
689 *Approximating matrix multiplication*, SIAM Journal on Computing, 36 (2006), pp. 132–  
690 157.
- 691 [8] P. DRINEAS, R. KANNAN, AND M. W. MAHONEY, *Fast monte carlo algorithms for matrices ii:*  
692 *Computing a low-rank approximation to a matrix*, SIAM Journal on computing, 36 (2006),  
693 pp. 158–183.
- 694 [9] A. FRIEZE, R. KANNAN, AND S. VEMPALA, *Fast monte carlo algorithms for finding low-rank*  
695 *approximations*, J. ACM, 51 (2004), pp. 1025–1041.
- 696 [10] Z. GIMBUTAS AND V. ROKHLIN, *A generalized fast multipole method for nonoscillatory kernels*,  
697 SIAM Journal on Scientific Computing, 24 (2003), pp. 796–817.
- 698 [11] A. G. GRAY AND A. W. MOORE, *N-body problems in statistical learning*, in Advances in neural  
699 information processing systems, 2001, pp. 521–527.
- 700 [12] L. GREENGARD AND V. ROKHLIN, *A fast algorithm for particle simulations*, Journal of Com-  
701 putational Physics, 135 (1997), pp. 280–292.
- 702 [13] W. HACKBUSCH, *A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices. Part I: Introduction to*  
703  *$\mathcal{H}$ -matrices*, Computing, 62 (1999), pp. 89–108.
- 704 [14] W. HACKBUSCH AND B. KHOROMSKIJ, *A sparse  $\mathcal{H}$ -matrix Arithmetic. Part II: Application to*  
705 *Multi-Dimensional Problems*, Computing, 64 (2000), pp. 21–47.
- 706 [15] N. HALKO, P.-G. MARTINSSON, AND J. A. TROPP, *Finding structure with randomness: Probabilistic*  
707 *algorithms for constructing approximate matrix decompositions*, SIAM Review, 53  
708 (2011), pp. 217–288.
- 709 [16] S. KAPUR AND D. E. LONG, *N-body problems: Ies 3: Efficient electrostatic and electromagnetic*  
710 *simulation*, IEEE Computational Science and Engineering, 5 (1998), pp. 60–67.
- 711 [17] S. KAPUR AND J. ZHAO, *A fast method of moments solver for efficient parameter extraction*  
712 *of mcms*, in Proceedings of the 34th annual Design Automation Conference, ACM, 1997,  
713 pp. 141–146.
- 714 [18] D. LEE, P. SAO, R. VUDUC, AND A. G. GRAY, *A distributed kernel summation framework*  
715 *for general-dimension machine learning*, Statistical Analysis and Data Mining: The ASA  
716 Data Science Journal, 7 (2014), pp. 1–13.
- 717 [19] W. B. MARCH AND G. BIROS, *Far-field compression for fast kernel summation methods in high*  
718 *dimensions*, Applied and Computational Harmonic Analysis, 43 (2017), pp. 39–75.
- 719 [20] W. B. MARCH, B. XIAO, AND G. BIROS, *Askit: Approximate skeletonization kernel-independent*  
720 *treecode in high dimensions*, SIAM Journal on Scientific Computing, 37 (2015), pp. A1089–  
721 A1110.
- 722 [21] P.-G. MARTINSSON, V. ROKHLIN, AND M. TYGERT, *A randomized algorithm for the decompo-*  
723 *sition of matrices*, Applied and Computational Harmonic Analysis, 30 (2011), pp. 47–68.
- 724 [22] N. YARVIN AND V. ROKHLIN, *Generalized gaussian quadratures and singular value decomposi-*  
725 *tions of integral operators*, SIAM Journal on Scientific Computing, 20 (1998), pp. 699–718.
- 726 [23] L. YING, G. BIROS, AND D. ZORIN, *A kernel-independent adaptive fast multipole algorithm in*  
727 *two and three dimensions*, Journal of Computational Physics, 196 (2004), pp. 591–626.
- 728 [24] C. D. YU, J. LEVITT, S. REIZ, AND G. BIROS, *Geometry-oblivious fmm for compressing dense*  
729 *spd matrices*, in Proceedings of the International Conference for High Performance Com-  
730 puting, Networking, Storage and Analysis, ACM, 2017, p. 53.