

1 **SUPERDC: SUPERFAST DIVIDE-AND-CONQUER EIGENVALUE**  
2 **DECOMPOSITION WITH IMPROVED STABILITY FOR**  
3 **RANK-STRUCTURED MATRICES\***

4 XIAOFENG OU<sup>†</sup> AND JIANLIN XIA<sup>†</sup>

5 **Abstract.** For dense symmetric matrices with small off-diagonal (numerical) ranks and in a  
6 hierarchically semiseparable form, we give a divide-and-conquer eigendecomposition method with  
7 nearly linear complexity (called SuperDC) that significantly improves an earlier basic algorithm in  
8 [Vogel, Xia, et al., SIAM J. Sci. Comput., 38 (2016)]. Some stability risks in the original algorithm are  
9 analyzed, including potential exponential norm growth, cancellations, loss of accuracy with clustered  
10 eigenvalues or intermediate eigenvalues, etc. In the dividing stage, we give a new structured low-rank  
11 updating strategy with balancing that eliminates the exponential norm growth and also minimizes  
12 the ranks of low-rank updates. In the conquering stage with low-rank updated eigenvalue solution,  
13 the original algorithm directly uses the standard fast multipole method (FMM) to accelerate function  
14 evaluations, which has the risks of cancellation, division by zero, and slow convergence. Here, we  
15 design a triangular FMM to avoid cancellation. Furthermore, when there are clustered intermediate  
16 eigenvalues, we design a novel local shifting strategy to integrate FMM accelerations into the solution  
17 of shifted secular equations. This helps achieve both the efficiency and the reliability. We also provide  
18 a deflation strategy with a user-supplied tolerance and give a precise description of the structure of  
19 the resulting eigenvector matrix. The SuperDC eigensolver has significantly improved stability while  
20 keeping the nearly linear complexity for finding the entire eigenvalue decomposition. Extensive  
21 numerical tests are used to show the efficiency and accuracy of SuperDC.

22 **Key words.** superfast eigenvalue decomposition, divide-and-conquer method, rank-structured  
23 matrix, triangular fast multipole method, shifted secular equation, local shifting

24 **AMS subject classifications.** 65F15, 65F55, 15A18, 15A23

25 **1. Introduction.** In this paper, we consider the full eigenvalue decomposition of  
26  $n \times n$  real symmetric matrices  $A$  with small off-diagonal ranks (and also  $A$  that can be  
27 approximated well by matrices with small off-diagonal ranks). Such matrices belong to  
28 the class of rank-structured matrices. Examples include banded matrices with finite  
29 bandwidth, Toeplitz matrices in Fourier space [33, 47, 55], some matrices arising  
30 from discretized PDEs and integral equations [32, 36, 51, 53], some kernel matrices  
31 [10, 56], etc. The eigenvalue decompositions of relevant matrices are very useful  
32 in computations such as matrix function evaluations [4], discretized linear system  
33 solutions [45], matrix equation solutions [35], and quadrature approximations [40].  
34 They are also very useful in fields such as optimization, imaging, Gaussian processes,  
35 and machine learning [35].

36 There are several types of rank-structured forms, such as  $\mathcal{H}/\mathcal{H}^2$  matrices [26, 27],  
37 hierarchical semiseparable (HSS) matrices [12, 54], quasiseparable/semiseparable ma-  
38 trices [11, 21, 42], BLR matrices [2], and HODLR matrices [1]. Examples of eigen-  
39 solvers for these rank-structured matrices include divide-and-conquer methods [3, 13,  
40 20, 29, 38, 44], QR iterations [7, 15, 19, 21, 41], bisection, [6, 48], and methods using  
41 accelerated characteristic polynomial evaluations [8].

42 Our work here focuses on the divide-and-conquer method for HSS matrices (that  
43 may be dense or sparse). The divide-and-conquer method has previously been well  
44 studied for tridiagonal matrices (which may be considered as special HSS forms). See,  
45 e.g., [5, 9, 16, 18, 24, 34]. In particular, a stable version is given in [24]. The algorithms

---

\*The research of Jianlin Xia was supported in part by an NSF grant DMS-1819166.

<sup>†</sup>Department of Mathematics, Purdue University, West Lafayette, IN 47907 (ou17@purdue.edu, xiaj@purdue.edu).

46 can compute all the eigenvalues in  $O(n^2)$  flops and can compute the eigenvectors in  
 47  $O(n^3)$  flops. It is also mentioned in [24] that it is possible to accelerate the operations  
 48 in the divide-and-conquer process via the fast multipole method (FMM) [23] to reach  
 49 nearly linear complexity. However, this was not actually done in [24] or later relevant  
 50 work [13, 29]. Only recently, the feasibility of the FMM acceleration of the divide-  
 51 and-conquer process was verified in a structured eigensolver in [44], which works for  
 52 HSS matrices without the need of tridiagonal reductions. For an HSS matrix with  
 53 off-diagonal ranks bounded by  $r$  (which may be a constant or a power of  $\log n$ ), the  
 54 method in [44] computes a structured eigendecomposition in  $O(r^2 n \log^2 n)$  flops with  
 55 storage  $O(rn \log n)$ . The method is then said to be *superfast*.

56 The work in [44] presents the basic framework of a rank-structured divide-and-  
 57 conquer eigensolver. It gives a proof-of-concept algorithm and verifies the feasibility  
 58 of such superfast eigenvalue solution for HSS matrices. Due to the complex nature of  
 59 the entire framework with many components, that preliminary work has some limita-  
 60 tions. It does not consider some crucial stability issues in the HSS divide-and-conquer  
 61 process, such as the risks of exponential norm growth and potential cancellations in  
 62 some function evaluations. Moreover, it does not incorporate several key stability  
 63 strategies that are used in practical tridiagonal divide-and-conquer algorithms. These  
 64 limitations are due to some major challenges in combining FMM accelerations with  
 65 those stability strategies, especially for problems with clustered eigenvalues.

66 Specifically, in the dividing stage, upper-level off-diagonal block information is  
 67 used to update lower-level diagonal blocks (also as HSS forms) of  $A$  in a hierarchical  
 68 process. The norms of the updated lower-level blocks may grow quickly during the  
 69 process, which brings stability risks and may even cause overflow. In the conquering  
 70 stage, multiple types of function evaluations are need in eigenvalue solutions (via a  
 71 modified Newton's method applied to some secular equations). The application of  
 72 FMM accelerations needs to assemble these function evaluations into matrix-vector  
 73 multiplications. However, classical stabilization techniques involve strategies such as  
 74 splitting function evaluations to avoid cancellation (see, e.g., [5, 9, 18, 28, 24]) and  
 75 solving certain *shifted* secular equations to guarantee accuracy for clustered eigen-  
 76 values [5, 9, 18, 28, 24]. Such splitting and shifting strategies depend on the each  
 77 individual eigenvalue to be sought so that it is difficult to find all the eigenvalues to-  
 78 gether with the usual FMM acceleration. (Sections 4.2.1 and 4.3.1 show the details.)  
 79 The algorithm in [44] directly applies usual FMM accelerations to standard secular  
 80 equations. This may lose accuracy or even encounter cancellations.

81 Thus, the main purpose of this paper is to overcome these limitations. We fol-  
 82 low the basic framework in [44] but provide some important stability, accuracy, and  
 83 efficiency improvements. We show how to integrate structured accelerations with sev-  
 84 eral stabilization strategies. A more reliable superfast divide-and-conquer eigensolver  
 85 (called SuperDC) is then designed to find an approximate eigenvalue decomposition  
 86 of  $A$ :

$$87 \quad (1.1) \quad A \approx Q \Lambda Q^T,$$

88 where  $\Lambda$  is a diagonal matrix for the eigenvalues and  $Q$  is for the orthogonal eigen-  
 89 vectors. For convenience, we call the matrix  $Q$  an *eigenmatrix*. (Our presentation  
 90 focuses on real symmetric  $A$ , and the ideas can be immediately extended to complex  
 91 Hermitian matrices). The main significance of the work includes the following.

- 92 1. We analyze why the original hierarchical dividing strategy in [44] can lead to  
 93 exponential norm growth. A more stable dividing strategy is designed, where  
 94 a balancing technique guarantees the norm growth is well under control. We

95 also give a strategy to choose appropriate low-rank updates in the dividing  
 96 stage so as to reduce the rank of low-rank updates and save eigenvalue solution  
 97 costs.

- 98 2. In the solution of the secular equations for the eigenvalues, we design a tri-  
 99 angular FMM to accommodate the eigenvalue-dependent splitting (for the  
 100 stability purpose as mentioned above). This enables us to quickly and stably  
 101 evaluate the functions (after splitting) through matrix-vector multiplications  
 102 that are not suitable for the standard FMM.
- 103 3. When shifted secular equations are used to handle clustered intermediate ei-  
 104 genvalues, we design a *local shifting* strategy that integrates shifts into FMM  
 105 matrices without destroying the FMM structure. This enhances the stabil-  
 106 ity of eigenvalue solution, leading to improved eigenvalue accuracy and also  
 107 better convergence of iterative secular equation solution.
- 108 4. We also provide clarifications and improvements on several aspects such as  
 109 the precise structure of the resulting eigenmatrix, the eigenvalue deflation  
 110 criterion with a user-supplied tolerance, and the stopping criterion in iterative  
 111 eigenvalue solution.
- 112 5. With all the stabilization strategies, SuperDC still nicely preserve the nearly  
 113 linear complexity. The eigendecomposition complexity is still  $O(r^2n \log^2 n)$ ,  
 114 with  $O(rn \log n)$  storage. No extra tridiagonal reduction is needed for dense  
 115 HSS matrices. We provide extensive numerical tests a SuperDC package in  
 116 Matlab. For modest matrix sizes  $n$ , SuperDC already has significantly lower  
 117 runtime and storage than some other eigensolvers while producing satisfactory  
 118 accuracies. Benefits of our stabilization strategies are also demonstrated.

119 In the remaining sections, we begin in Section 2 with a quick review of the basic  
 120 HSS divide-and-conquer eigensolver in [44]. Then the improved structured dividing  
 121 strategy is discussed in Section 3, followed by the efficient structured conquering  
 122 scheme in Section 4. Section 5 gives some numerical experiments to demonstrate the  
 123 efficiency and accuracy. Then Section 6 concludes the paper. A list of the major  
 124 algorithms is given in the supplementary materials.

125 Throughout this paper, the following notation is used.

- 126 • Lower-case letters in bold fonts like  $\mathbf{u}$  are used to denote vectors.
- 127 •  $(A_{ij})_{n \times n}$  means an  $n \times n$  matrix with the  $(i, j)$ -entry  $A_{ij}$ .
- 128 • Sometimes, a vector  $\mathbf{s}$  may be viewed as an ordered set formed by its compo-  
 129 nents  $s_i$ . Then  $s_i \in \mathbf{s}$  means  $s_i$  is a component of  $\mathbf{s}$ . Accordingly, for vectors  
 130  $\mathbf{s}$  and  $\mathbf{t}$ , a matrix  $(\kappa(s_i, t_j))_{s_i \in \mathbf{s}, t_j \in \mathbf{t}}$  may be defined by the evaluation of a  
 131 function  $\kappa(s, t)$  at the components of  $\mathbf{s}$  and  $\mathbf{t}$ .
- 132 •  $\text{diag}(\dots)$  denotes a (block) diagonal matrix.
- 133 •  $\text{rowsize}(A)$  and  $\text{colsize}(A)$  mean the row and column sizes of  $A$ , respectively.
- 134 •  $\mathbf{u} \odot \mathbf{v}$  denotes the entrywise (Hadamard) product of two vectors  $\mathbf{u}$  and  $\mathbf{v}$ .
- 135 • For a binary tree  $\mathcal{T}$ , we suppose it is in postordering so that it has nodes  
 136  $i = 1, 2, \dots, \text{root}(\mathcal{T})$ , where  $\text{root}(\mathcal{T})$  is the root.
- 137 •  $\text{fl}(x)$  denotes the floating point result of  $x$ .
- 138 •  $\epsilon_{\text{mach}}$  represents the machine precision.

139 **2. Review of the basic superfast divide-and-conquer eigensolver.** We  
 140 first briefly summarize the basic superfast divide-and-conquer eigensolver in [44],  
 141 which generalizes the classical divide-and-conquer method for tridiagonal matrices  
 142 to HSS matrices.

143 A symmetric HSS matrix  $A$  [54] may be defined with the aid of a postordered full

144 binary tree  $\mathcal{T}$  called *HSS tree*, and has a nested structure that looks like

$$145 \quad (2.1) \quad D_p = \begin{pmatrix} D_i & U_i B_i U_j^T \\ U_j B_i^T U_i^T & D_j \end{pmatrix},$$

146 where  $p \in \mathcal{T}$  has child nodes  $i$  and  $j$ , so that  $D_p$  with  $p = \text{root}(\mathcal{T})$  is the entire HSS  
147 matrix  $A$ . Here, the  $U$  matrices are off-diagonal basis matrices and also satisfy a  
148 nested relationship  $U_p = \begin{pmatrix} U_i & \\ & U_j \end{pmatrix} \begin{pmatrix} R_i \\ R_j \end{pmatrix}$ . The  $D_i, U_i, B_i$  matrices are called *HSS*  
149 *generators* associated with node  $i$ . The maximum size of the  $B$  generators is usually  
150 referred as the *HSS rank* of  $A$ . We suppose  $\text{root}(\mathcal{T})$  is at level 0, and the children of  
151 a node  $i$  at level  $l$  are at level  $l + 1$ .

152 The superfast divide-and-conquer eigensolver in [44] finds the eigendecomposition  
153 (1.1) of  $A$  through a dividing stage and a conquering stage as follows.

154 **2.1. Dividing stage.** In the dividing stage in [44],  $A$  and its submatrices are  
155 recursively divided into block-diagonal HSS forms plus low-rank updates. Start with  
156  $p = \text{root}(\mathcal{T})$  and its two children  $i$  and  $j$ .  $A = D_p$  in (2.1) can be written as

$$157 \quad (2.2) \quad D_p = \begin{pmatrix} D_i - U_i B_i B_i^T U_i^T & \\ & D_j - U_j U_j^T \end{pmatrix} + \begin{pmatrix} U_i B_i \\ U_j \end{pmatrix} \begin{pmatrix} B_i^T U_i^T & U_j^T \end{pmatrix}.$$

158 For notational convenience, we suppose the HSS rank of  $A$  is  $r$  and each  $B$  generator  
159 has column size  $r$ . Let

$$160 \quad (2.3) \quad \hat{D}_i = D_i - U_i B_i B_i^T U_i^T, \quad \hat{D}_j = D_j - U_j U_j^T, \quad Z_p = \begin{pmatrix} U_i B_i \\ U_j \end{pmatrix},$$

161 and we arrive at

$$162 \quad (2.4) \quad D_p = \text{diag}(\hat{D}_i, \hat{D}_j) + Z_p Z_p^T.$$

163 Here, the diagonal blocks  $D_i$  and  $D_j$  are modified so that a rank- $r$  update  $Z_p Z_p^T$   
164 can be used instead of a rank- $2r$  update. The column size of  $Z_p$  is referred as the *rank* of  
165 *the low-rank update* and here we have  $\text{colsize}(Z_p) = \text{colsize}(B_i)$ .

166 During this process, the blocks  $\hat{D}_i$  and  $\hat{D}_j$  remain to be HSS forms. In fact,  
167 it is shown in [44, 54] that any matrix of the form  $D_i - U_i H U_i^T$  can preserve the  
168 off-diagonal basis matrices of  $D_i$ . Specifically, the following lemma can be used for  
169 generator updates.

170 **LEMMA 2.1.** [44] *Let  $\mathcal{T}_i$  be the subtree of the HSS tree  $\mathcal{T}$  that has the node  $i$  as*  
171 *the root. Then  $D_i - U_i H U_i^T$  has HSS generators  $\tilde{D}_k, \tilde{U}_k, \tilde{R}_k, \tilde{B}_k$  for each node  $k \in \mathcal{T}_i$*   
172 *as follows:*

$$173 \quad \tilde{U}_k = U_k, \quad \tilde{R}_k = R_k,$$

$$174 \quad (2.5) \quad \tilde{B}_k = B_k - (R_k R_{k_l} \cdots R_{k_1}) H (R_{k_1}^T \cdots R_{k_l}^T R_k^T),$$

$$175 \quad \tilde{D}_k = D_k - U_k (R_k R_{k_l} \cdots R_{k_1}) H (R_{k_1}^T \cdots R_{k_l}^T R_k^T) U_k^T \quad \text{for a leaf } k,$$

177 where  $\tilde{k}$  is the sibling node of  $k$  and  $k \rightarrow k_l \rightarrow \cdots \rightarrow k_1 \rightarrow i$  is the path connecting  $k$   
178 to  $i$ . Accordingly,  $D_i - U_i H U_i^T$  and  $D_i$  have the same off-diagonal basis matrices.

179 Thus, the HSS generators of  $\hat{D}_i$  and  $\hat{D}_j$  can be conveniently obtained via the  
180 generator update procedure (2.5). Then the dividing process can continue on  $\hat{D}_i$  and  
181  $\hat{D}_j$  like above with  $p$  in (2.2) replaced by  $i$  and  $j$ , respectively.

182 **2.2. Conquering stage.** Suppose eigenvalue decompositions of the subprob-  
 183 lems  $\hat{D}_i$  and  $\hat{D}_j$  in (2.3) have been computed respectively as

$$184 \quad (2.6) \quad \hat{D}_i = Q_i \Lambda_i Q_i^T, \quad \hat{D}_j = Q_j \Lambda_j Q_j^T.$$

185 Then from (2.4), we have

$$186 \quad (2.7) \quad D_p = \text{diag}(Q_i, Q_j) \left( \text{diag}(\Lambda_i, \Lambda_j) + \hat{Z}_p \hat{Z}_p^T \right) \text{diag}(Q_i^T, Q_j^T), \quad \text{with}$$

$$187 \quad (2.8) \quad \hat{Z}_p = \text{diag}(Q_i^T, Q_j^T) Z_p.$$

189 Consequently, if we can solve the rank- $r$  updating problem

$$190 \quad (2.9) \quad \text{diag}(\Lambda_i, \Lambda_j) + \hat{Z}_p \hat{Z}_p^T = \hat{Q}_p \Lambda_p \hat{Q}_p^T,$$

191 then the eigendecomposition of  $D_p$  can be simply retrieved as

$$192 \quad (2.10) \quad D_p = Q_p \Lambda_p Q_p^T, \quad \text{with} \quad Q_p = \text{diag}(Q_i, Q_j) \hat{Q}_p.$$

193 Therefore, the main task is to compute the eigendecomposition of the low-rank  
 194 update problem (2.9). To this end, suppose  $\hat{Z}_p = (\mathbf{z}_1 \cdots \mathbf{z}_r)$ , where  $\mathbf{z}_k$ 's are the  
 195 columns of  $\hat{Z}_p$ . Then (2.9) can be treated as  $r$  rank-1 updating problems  $\text{diag}(\Lambda_i, \Lambda_j) +$   
 196  $\sum_{k=1}^r \mathbf{z}_k \mathbf{z}_k^T$ . As a result, a basic component is to quickly find the eigenvalue decom-  
 197 position of a diagonal plus rank-1 updating problem in the following form:

$$198 \quad (2.11) \quad \tilde{\Lambda} + \mathbf{v} \mathbf{v}^T = \tilde{Q} \Lambda \tilde{Q}^T,$$

199 where  $\tilde{\Lambda} = \text{diag}(d_1, \dots, d_n)$  with  $d_1 \leq \dots \leq d_n$ ,  $\mathbf{v} = (v_1, \dots, v_n)^T$ ,  $\tilde{Q} = (\tilde{\mathbf{q}}_1, \dots, \tilde{\mathbf{q}}_n)$ ,  
 200 and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_n)$ .

201 As in standard divide-and-conquer eigensolvers (see, e.g., [5, 16, 24]), the eigen-  
 202 values  $\lambda_k$  are found through the solution of the following secular equation [22]:

$$203 \quad (2.12) \quad f(x) = 1 + \sum_{k=1}^n \frac{v_k^2}{d_k - x} = 0.$$

204 Newton iterations with rational interpolations may be used and the cost for finding  
 205 all the  $n$  roots is  $O(n^2)$ . Once  $\lambda_k$  is computed, a corresponding eigenvector looks like  
 206  $\tilde{\mathbf{q}}_k = (\tilde{\Lambda} - \lambda_k I)^{-1} \mathbf{v}$ . Such an analytical form is not directly used in general for the  
 207 stability reason, since any loss of precision in the computed  $\lambda_k$  can be significantly  
 208 amplified in  $(\tilde{\Lambda} - \lambda_k I)^{-1} \mathbf{v}$ , which will result in the loss of eigenvector orthogonality  
 209 [18, 24]. A stable way to obtain  $\tilde{\mathbf{q}}_k$  is given in [24] based on Löwner's formula.

210 It is also mentioned in [24] that nearly  $O(n)$  complexity may be achieved by  
 211 assembling multiple operations into matrix-vector multiplications that can be accel-  
 212 erated by the FMM. This is first verified in [44], where the complexity of the algorithm  
 213 for finding the entire eigendecomposition is  $O(r^2 n \log^2 n)$  instead of  $O(n^3)$ , with the  
 214 eigenmatrix  $Q$  in (1.1) given in a structured form that needs  $O(rn \log n)$  storage in-  
 215 stead of  $O(n^2)$ . In the following sections, we give a series of stability enhancements  
 216 to get an improved superfast divide-and-conquer eigensolver.

217 **3. Improved structured dividing strategy.** In this section, we point out a  
 218 stability risk in the original dividing method as given in (2.2)–(2.3) and propose a  
 219 more stable dividing strategy. We also design a way to minimize  $\text{colsize}(Z_p)$ .

220 The stability risk can be illustrated as follows. Consider  $\hat{D}_i$  in (2.2) which is  
 221 the result of updating  $D_i$  in the dividing process associated with the parent  $p$  of  $i$ .  
 222 Suppose  $i$  has children  $c_1$  and  $c_2$  such that

$$223 \quad (3.1) \quad D_i = \begin{pmatrix} D_{c_1} & U_{c_1} B_{c_1} U_{c_2}^T \\ U_{c_2} B_{c_2} U_{c_1}^T & D_{c_2} \end{pmatrix}, \quad U_i = \begin{pmatrix} U_{c_1} & \\ & U_{c_2} \end{pmatrix} \begin{pmatrix} R_{c_1} \\ R_{c_2} \end{pmatrix}.$$

224 Then

$$225 \quad \hat{D}_i = D_i - U_i B_i B_i^T U_i^T = \begin{pmatrix} \tilde{D}_{c_1} & U_{c_1} \tilde{B}_{c_1} U_{c_2}^T \\ U_{c_2} \tilde{B}_{c_1}^T U_{c_1}^T & \tilde{D}_{c_2} \end{pmatrix},$$

226 where

$$227 \quad \tilde{D}_{c_1} = D_{c_1} - U_{c_1} R_{c_1} B_i B_i^T R_{c_1}^T U_{c_1}^T, \quad \tilde{D}_{c_2} = D_{c_2} - U_{c_2} R_{c_2} B_i B_i^T R_{c_2}^T U_{c_2}^T, \\ 228 \quad (3.2) \quad \tilde{B}_{c_1} = B_{c_1} - R_{c_1} B_i B_i^T R_{c_2}^T.$$

230 In HSS constructions [54], to ensure stability of HSS algorithms, the  $U$  basis gen-  
 231 erators often have orthonormal columns [46, 47]. Then due to (3.1), the  $R$  generators  
 232 also satisfy that  $\begin{pmatrix} R_{c_1} \\ R_{c_2} \end{pmatrix}$  has orthonormal columns. Then each  $B$  generator has 2-  
 233 norm equal to its associated off-diagonal block. For example,  $\|B_i\|_2 = \|U_i B_i U_j^T\|_2$ .  
 234 Furthermore,  $\|R_{c_1}\|_2 \leq 1$ ,  $\|R_{c_2}\|_2 \leq 1$ , and (3.2) means

$$235 \quad (3.3) \quad \|\tilde{B}_{c_1}\|_2 \leq \|B_{c_1}\|_2 + \|B_i\|_2^2.$$

236 If the off-diagonal block  $U_i B_i U_j^T$  has a large norm,  $\|\tilde{B}_{c_1}\|_2$  can potentially be much  
 237 larger than  $\|B_{c_1}\|_2$ . We can similarly observe the norm growth with the updated  $D$   
 238 generators. Moreover, when the dividing process proceeds on  $\tilde{D}_{c_1}$ , the norms of the  
 239 updated  $B, D$  generators at lower levels can grow exponentially.

240 **PROPOSITION 3.1.** *Suppose the  $U_k$  generator of  $A$  associated with each node  $k$   
 241 of  $\mathcal{T}$  with  $k \neq \text{root}(\mathcal{T})$  has orthonormal columns and all the original  $B_k$  generators  
 242 satisfy  $\|B_k\|_2 \leq \beta$  with  $\beta \gg 1$ . Also suppose the leaves of  $\mathcal{T}$  are at level  $l_{\max} \leq \log_2 n$ .  
 243 When the original dividing process in Section 2.1 proceeds from  $\text{root}(\mathcal{T})$  to a nonleaf  
 244 node  $i$ , immediately after finishing the dividing process associated with node  $i$ ,*

- 245 • *with  $i$  at level  $l \leq l_{\max} - 2$ , the updated  $B_k$  generator (denoted  $\tilde{B}_k$ ) associated*  
 246 *with any descendant  $k$  of  $i$  satisfies*

$$247 \quad (3.4) \quad \|\tilde{B}_k\|_2 = O(\beta^{2^l}) \leq O(\beta^{n/4}),$$

248 *where  $O(\cdot)$  denotes the asymptotic upper bound and is given in terms of the*  
 249 *highest order term in  $\beta$ ;*

- 250 • *with  $i$  at level  $l \leq l_{\max} - 1$ , the updated  $D_k$  generator (denoted  $\tilde{D}_k$ ) associated*  
 251 *with any leaf descendant  $k$  of  $i$  satisfies*

$$252 \quad (3.5) \quad \|\tilde{D}_k\|_2 = \|D_k\|_2 + O(\beta^{2^l}) \leq \|D_k\|_2 + O(\beta^{n/2}).$$

253 *Proof.* Following the update formulas in Lemma 2.1, we just need to show the  
 254 norm bound for  $\|\tilde{B}_k\|_2$ . The bound for  $\|\tilde{D}_k\|_2$  can be shown similarly.

255 After the dividing process associated with  $\text{root}(\mathcal{T})$  is finished, according to (2.5),  
 256  $\tilde{B}_k$  associated with any descendant  $k$  of a child  $i$  of  $\text{root}(\mathcal{T})$  looks like

$$257 \quad (3.6) \quad \tilde{B}_k = B_k - (R_k R_{k_{m-1}} \cdots R_{k_1}) H_i (R_{k_1}^T \cdots R_{k_{m-1}}^T R_k^T),$$

258 where  $H_i = B_i B_i^T$  if  $i$  is the left child of  $\text{root}(\mathcal{T})$  or  $H_i = I$  otherwise,  $k$  is supposed  
 259 to be at level  $m$  with sibling  $\tilde{k}$ , and  $k \rightarrow k_{m-1} \rightarrow \dots \rightarrow k_1 \rightarrow i$  is the path connecting  
 260  $k$  to  $i$  in the HSS tree  $\mathcal{T}$ . Clearly,  $\|H_i\|_2 \leq \beta^2$ . With the orthogonality condition of  
 261 the  $U$  basis generators,  $\begin{pmatrix} R_{c_1} \\ R_{c_2} \end{pmatrix}$  also has orthogonal columns. Then we get

$$262 \quad (3.7) \quad \|\tilde{B}_k\|_2 \leq \|B_k\|_2 + \|H_i\|_2 \leq \beta + \beta^2 = O(\beta^2).$$

263 Then in the dividing process associated with node  $i$  at level 1, for a child  $c$  of  $i$   
 264 (see Figure 3.1 for an illustration), the generator  $\tilde{D}_c$  is further updated to

$$265 \quad (3.8) \quad \hat{D}_c = \tilde{D}_c - U_c H_c U_c^T,$$

266 where  $H_c = \tilde{B}_c \tilde{B}_c^T$  if  $c$  is the left child of  $i$  or  $H_c = I$  otherwise. We have  $\|H_c\|_2 \leq$   
 267  $\|\tilde{B}_c\|_2^2$  for the first case and  $\|H_c\|_2 = 1$  for the second case. From (3.7), we have  
 268  $\|H_c\|_2 \leq (\beta^2 + \beta)^2$ . For any descendant  $k$  of  $c$  with sibling  $\tilde{k}$ , (3.8) needs to update  
 269 the generator  $B_k$  to

$$270 \quad (3.9) \quad \tilde{B}_k = B_k - (R_k R_{k_{m-1}} \dots R_{k_2} R_c) H_i (R_c^T R_{k_2}^T \dots R_{k_{m-1}}^T R_{\tilde{k}}^T) \\ 271 \quad - (R_k R_{k_{m-1}} \dots R_{k_2}) H_c (R_{k_2}^T \dots R_{k_{m-1}}^T) R_{\tilde{k}}^T,$$

273 where the last term on the right-hand side is because of the update associated with  
 274 the dividing of  $D_i$  like in (3.6). Then

$$275 \quad (3.10) \quad \|\tilde{B}_k\|_2 \leq \|B_k\|_2 + \|H_i\|_2 + \|H_c\|_2 \leq \beta + \beta^2 + (\beta^2 + \beta)^2 = O(\beta^4).$$

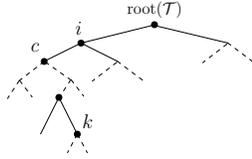


FIG. 3.1. Nodes involved in the dividing process.

276

277 If the dividing process continues to  $c$ , it is similar to obtain  $\|\tilde{B}_k\|_2 = O(\beta^8)$  for  
 278 any descendant  $k$  of a child of  $c$ . We can then similarly reach the conclusion on the  
 279 general pattern of the norm growth as in (3.4). Also, if  $i$  is at level  $l_{\max} - 1$ , then  $B_k$   
 280 associated with a child  $k$  of  $i$  is not updated, which is why only  $i$  at level  $l \leq l_{\max} - 2$   
 281 contributes to the norm growth of lower level  $B$  generators.  $\square$

282 The bound  $O(\beta^{2^l})$  in (3.4) for  $\|\tilde{B}_k\|_2$  and the bound  $\|D_k\|_2 + O(\beta^{2^l})$  in (3.5) for  
 283  $\|\tilde{D}_k\|_2$  are attainable. To see this, suppose  $i$  is a child of  $\text{root}(\mathcal{T})$  and  $\|B_i\|_2 = \beta$ .  
 284 Note the multiplicative forms like  $H_i$  below (3.6) and  $H_c$  below (3.8). Following the  
 285 proof, we can see that the asymptotic upper bounds (3.4) and (3.5) can be attained  
 286 at some leaf level node  $k$  after the dividing process associated with the parent of  $k$   
 287 is completed.

288 This proposition indicates that, during the original hierarchical dividing process,  
 289 the updated  $B, D$  generators associated with a lower-level node may potentially have  
 290 exponential norm accumulation, as long as one of its ancestors is associated with a  
 291  $B$  generator with a large norm. This can cause stability issues or even overflow, as  
 292 confirmed in the numerical tests later.

293 To resolve this, we introduce *balancing*/scaling into the updates and propose a  
 294 new dividing strategy. That is, we replace the original dividing method (2.2) by

$$295 \quad (3.11) \quad D_p = \begin{pmatrix} D_i - \frac{1}{\|B_i\|_2} U_i B_i B_i^T U_i^T & \\ & D_j - \|B_i\|_2 U_j U_j^T \end{pmatrix} \\
 296 \quad + \begin{pmatrix} \frac{1}{\sqrt{\|B_i\|_2}} U_i B_i \\ \sqrt{\|B_i\|_2} U_j \end{pmatrix} \begin{pmatrix} \frac{1}{\sqrt{\|B_i\|_2}} B_i^T U_i^T & \sqrt{\|B_i\|_2} U_j^T \end{pmatrix}. \\
 297$$

298 Then we still have (2.4), but with  
 (3.12)

$$299 \quad \hat{D}_i = D_i - \frac{1}{\|B_i\|_2} U_i B_i B_i^T U_i^T, \quad \hat{D}_j = D_j - \|B_i\|_2 U_j U_j^T, \quad Z_p = \begin{pmatrix} \frac{1}{\sqrt{\|B_i\|_2}} U_i B_i \\ \sqrt{\|B_i\|_2} U_j \end{pmatrix}.$$

300 We show how this strategy controls the norms of the updated  $B, D$  generators.

301 **PROPOSITION 3.2.** *Suppose the same conditions as in Proposition 3.1 hold, except*  
 302 *that (2.2) is replaced by (3.11) so that (2.3) is replaced by (3.12). Then (3.4) and*  
 303 *(3.5) become, respectively,*

$$304 \quad (3.13) \quad \|\tilde{B}_k\|_2 \leq 2^l \beta \leq \frac{n}{4} \beta, \quad \|\tilde{D}_k\|_2 \leq \|D_k\|_2 + 2^l \beta \leq \|D_k\|_2 + \frac{n}{2} \beta.$$

305 *Proof.* The proof follows a procedure similar to the proof for Proposition 3.1.  
 306 Again, we just show the result for  $\|\tilde{B}_k\|_2$ . After the dividing process associated with  
 307  $\text{root}(\mathcal{T})$  is finished, we still have (3.6) for any descendant  $k$  of a child  $i$  of  $\text{root}(\mathcal{T})$ ,  
 308 except that  $H_i = \frac{B_i B_i^T}{\|B_i\|_2}$  if  $i$  is the left child of  $\text{root}(\mathcal{T})$  or  $H_i = \|B_i\|_2 I$  otherwise. In  
 309 either case, we have  $\|H_i\|_2 \leq \beta$ . Then (3.7) becomes

$$310 \quad (3.14) \quad \|\tilde{B}_k\|_2 \leq 2\beta.$$

311 Then in the dividing process associated with node  $i$  at level 1, for a child  $c$  of  
 312  $i$ , the generator  $\tilde{D}_c$  is updated like in (3.8), except that  $H_c = \frac{\tilde{B}_c \tilde{B}_c^T}{\|\tilde{B}_c\|_2}$  if  $c$  is the left  
 313 child of  $i$  or  $H_c = \|\tilde{B}_c\|_2 I$  otherwise. We have  $\|H_c\|_2 \leq \|\tilde{B}_c\|_2$  for both cases. From  
 314 (3.14),  $\|H_c\|_2 \leq 2\beta$ . For any descendant  $k$  of  $c$ , (3.8) still requires the update of the  
 315 generator  $B_k$  to  $\tilde{B}_k$  like in (3.9), except that (3.10) now becomes

$$316 \quad \|\tilde{B}_k\|_2 \leq \|B_k\|_2 + \|H_i\|_2 + \|H_c\|_2 \leq \beta + \beta + 2\beta = 4\beta.$$

317 If the dividing process continues to  $c$ , it is similar to obtain  $\|\tilde{B}_k\|_2 \leq 8\beta$  for any  
 318 descendant  $k$  of the left child of  $c$ . We can similarly get the norm growth as in (3.13)  
 319 in general.  $\square$

320 Therefore, the norm growth now becomes at most linear in  $n$  and is well controlled,  
 321 in contrast to the exponential growth in Proposition 3.1. Here again, the upper bounds  
 322  $2^l \beta$  for  $\|\tilde{B}_k\|_2$  and  $\|D_k\|_2 + 2^l \beta$  for  $\|\tilde{D}_k\|_2$  are attainable.

323 Next, we can also minimize  $\text{colsize}(Z_p)$ , the rank of the low-rank update. Note  
 324 that in the original dividing method (2.2) in [44], the updates to the two diago-  
 325 nal blocks involve the  $B_i$  generator in different ways. That is,  $D_i$  is updated by  
 326  $-U_i B_i B_i^T U_i^T$  while  $D_j$  is updated by  $-U_j U_j^T$ . In fact, (2.2) may be reformulated so  
 327 that  $D_i$  is updated by  $-U_i U_i^T$  while  $D_j$  is updated by  $-U_j B_i^T B_i U_j^T$ . It is not clear  
 328 from [44] which way is better.

329 In fact, in (2.2) and also (3.11)–(3.12), the rank of the low-rank update is equal  
 330 to  $\text{colsize}(B_i)$ . In practice,  $B_i$  may not be a square matrix. Thus, (3.12) shall be used  
 331 only if  $\text{colsize}(B_i) \leq \text{rowsize}(B_i)$ . Otherwise, we replace (3.12) by the following:  
 (3.15)

$$332 \quad \hat{D}_i = D_i - \|B_i\|_2 U_i U_i^T, \quad \hat{D}_j = D_j - \frac{1}{\|B_i\|_2} U_j B_i^T B_i U_j^T, \quad Z_p = \begin{pmatrix} \sqrt{\|B_i\|_2} U_i \\ \frac{1}{\sqrt{\|B_i\|_2}} U_j B_i^T \end{pmatrix},$$

333 so that (2.4) still holds. In (3.15), the low-rank update size is now  $\text{rowsize}(B_i)$ . With  
 334 such a choice between (3.12) and (3.15), we ensure that the size of the low-rank  
 335 update  $\text{colsize}(Z_p)$  is always  $\min(\text{rowsize}(B_i), \text{colsize}(B_i))$ . This strategy benefits the  
 336 efficiency in the conquering stage since it reduces the number of rank-1 updates. With  
 337 these new ideas, we have a more stable and efficient dividing stage.

338 **4. Improved structured conquering stage.** In this section, we discuss the  
 339 solution of the eigenvalues and eigenvectors in the conquering stage via the integra-  
 340 tion of various stability strategies into FMM accelerations. We first show a flexible  
 341 deflation strategy. Then we give a triangular FMM idea for accelerating the secu-  
 342 lar equation solution and a local shifting strategy for solving shifted secular equations  
 343 and constructing structured eigenvectors. We also discuss the framework of the overall  
 344 eigendecomposition and the precise structure of the overall eigenmatrix.

345 **4.1. User-controlled deflation.** As reviewed in Section 2.2, the key problem  
 346 in the conquering stage is to quickly find the eigendecomposition of the rank-one  
 347 updating problem (2.11). Like in earlier studies in [9, 18], an eigenvalue deflation step  
 348 may be first applied to reduce the size of (2.11) if  $|v_j|$  or the difference  $|d_j - d_{j+1}|$  is  
 349 small. In the implementations of the tridiagonal divide-and-conquer eigensolver (see,  
 350 e.g., [5]), the deflation is performed in a two-step procedure with a tolerance related to  
 351  $\epsilon_{\text{mach}}$ . Here, we follow the same steps, but replace  $\epsilon_{\text{mach}}$  with a user-supplied deflation  
 352 tolerance  $\tau$  to get a more flexible deflation procedure.

353 (i) If  $|v_j| < \tau$ , without loss of generality, we assume  $j = n$ ,  $\mathbf{v} = \begin{pmatrix} \mathbf{v}_1 \\ v_n \end{pmatrix}$  and get

$$354 \quad \tilde{\Lambda} + \mathbf{v}\mathbf{v}^T = \begin{pmatrix} \tilde{\Lambda}_1 & \\ & d_n \end{pmatrix} + \begin{pmatrix} \mathbf{v}_1 \\ v_n \end{pmatrix} \begin{pmatrix} \mathbf{v}_1^T & v_n \end{pmatrix} \approx \begin{pmatrix} \tilde{\Lambda}_1 + \mathbf{v}_1 \mathbf{v}_1^T & \\ & d_n \end{pmatrix},$$

355 where the approximation has an error proportional to  $\tau$ . Then we only need to find  
 356 the eigendecomposition of the smaller problem  $\tilde{\Lambda}_1 + \mathbf{v}_1 \mathbf{v}_1^T$ .

357 (ii) If  $|(d_j - d_{j+1})v_j v_{j+1}| < (v_j^2 + v_{j+1}^2)\tau$ , we can find a Givens rotation matrix  $G$   
 358 such that  $G \begin{pmatrix} v_j \\ v_{j+1} \end{pmatrix} = \begin{pmatrix} 0 \\ w \end{pmatrix}$  with  $w = \sqrt{v_j^2 + v_{j+1}^2}$ . Then

$$359 \quad G \left( \begin{pmatrix} d_j & \\ & d_{j+1} \end{pmatrix} + \begin{pmatrix} v_j \\ v_{j+1} \end{pmatrix} \begin{pmatrix} v_j & v_{j+1} \end{pmatrix}^T \right) G^T = \begin{pmatrix} d_j & \mu \\ \mu & d_{j+1} \end{pmatrix} + \begin{pmatrix} 0 \\ w \end{pmatrix} \begin{pmatrix} 0 & w \end{pmatrix}^T \\ 360 \quad \approx \begin{pmatrix} d_j & \\ & d_{j+1} \end{pmatrix} + \begin{pmatrix} 0 \\ w \end{pmatrix} \begin{pmatrix} 0 & w \end{pmatrix}^T = \begin{pmatrix} d_j & \\ & d_{j+1} + w^2 \end{pmatrix},$$

362 where  $\mu = \frac{(d_j - d_{j+1})v_j v_{j+1}}{v_j^2 + v_{j+1}^2}$  and the approximation has a 2-norm error  $|\mu| < \tau$ . This  
 363 then leads to a diagonal subproblem.

364 After the above deflation steps, the problem size of (2.11) is reduced and the  
 365 simplified problem satisfies

$$366 \quad (4.1) \quad |v_j| \geq \tau \quad \text{and} \quad |d_j - d_{j+1}| \geq \frac{(v_j^2 + v_{j+1}^2)\tau}{|v_j v_{j+1}|}.$$

367 The parameter  $\tau$  offers the flexibility to control the accuracy of the eigenvalues.  
 368 When only moderate accuracy is needed, a larger  $\tau$  can be used for a more significant  
 369 reduction in the problem size. Moreover, this can sometimes avoid the need to deal  
 370 with situations where  $|\lambda_j - d_j|$  or  $|\lambda_j - d_{j+1}|$  is too small.

371 **4.2. Fast secular equation solution.** Assume (4.1) holds for (2.11) so that  
 372 no deflation is needed. We consider the solution of the secular equation (2.12) for  
 373 its eigenvalues  $\lambda_k, k = 1, 2, \dots, n$ . Without loss of generality, suppose the diagonal  
 374 entries  $d_k$  of  $\tilde{\Lambda}$  are ordered from the smallest to the largest.

375 **4.2.1. Standard FMM accelerations and the limitation.** When the modi-  
 376 fied Newton's method is used to solve for  $\lambda_k$ , it needs to evaluate  $f$  (referred to as the  
 377 secular function) in (2.12) and its derivative  $f'$  at certain  $x_k \in (d_k, d_{k+1})$ . The idea in  
 378 [13, 24, 44] is to assemble the function evaluations for all  $k$  together as matrix-vector  
 379 multiplications that can be accelerated by the standard FMM. That is, let

$$380 \quad \mathbf{f} = (f(x_1) \ \cdots \ f(x_n))^T, \quad \mathbf{f}' = (f'(x_1) \ \cdots \ f'(x_n))^T,$$

$$381 \quad (4.2) \quad \mathbf{v} = (v_1 \ \cdots \ v_n)^T, \quad \mathbf{w} = \mathbf{v} \odot \mathbf{v}, \quad \mathbf{e} = (1 \ \cdots \ 1)^T,$$

$$382 \quad (4.3) \quad C = \left( \frac{1}{d_j - x_i} \right)_{n \times n}, \quad S = \left( \frac{1}{(d_j - x_i)^2} \right)_{n \times n},$$

$$383 \quad (4.4) \quad \mathbf{f} = \mathbf{e} + C\mathbf{w}, \quad \mathbf{f}' = S\mathbf{w}.$$

385 The vectors  $\mathbf{f}$  and  $\mathbf{f}'$  can be quickly evaluated by the FMM with the kernel functions  
 386  $\kappa(s, t) = \frac{1}{s-t}$  and  $\kappa(s, t) = \frac{1}{(s-t)^2}$ , respectively.

387 A basic idea of the FMM for computing, say,  $C\mathbf{w}$  is as follows. Note that  $C$  is  
 388 the evaluation of  $\kappa(s, t) = \frac{1}{s-t}$  at interlaced points  $s \in \{d_j\}_{1 \leq j \leq n}$  and  $t \in \{x_i\}_{1 \leq i \leq n}$ :

$$389 \quad (4.5) \quad d_i < x_i < d_{i+1} < x_{i+1}, \quad 1 \leq i \leq n-1.$$

390 The sets  $\{x_i\}_{1 \leq i \leq n}$  and  $\{d_j\}_{1 \leq j \leq n}$  together are treated as one set and then hier-  
 391 archically partitioned. This also naturally leads to a hierarchical partition of both  
 392  $\{x_i\}_{1 \leq i \leq n}$  and  $\{d_j\}_{1 \leq j \leq n}$ . Consider two subsets produced in this partitioning:

$$393 \quad (4.6) \quad \mathbf{s}_x \subset \{x_i\}_{1 \leq i \leq n}, \quad \mathbf{s}_d \subset \{d_j\}_{1 \leq j \leq n}.$$

394 Use  $C_{\mathbf{s}_x, \mathbf{s}_d} = (\kappa(d_j, x_i))_{x_i \in \mathbf{s}_x, d_j \in \mathbf{s}_d}$  to denote the block of  $C$  defined by  $\mathbf{s}_x$  and  $\mathbf{s}_d$ ,  
 395 which is often referred as the *interaction* between  $\mathbf{s}_x$  and  $\mathbf{s}_d$ .

396 • If  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are well separated (a precise definition of the separation can be  
 397 found in [23, 37]), then  $C_{\mathbf{s}_x, \mathbf{s}_d}$  can be approximated by a low-rank form

$$398 \quad (4.7) \quad C_{\mathbf{s}_x, \mathbf{s}_d} \approx U_{\mathbf{s}_x} B_{\mathbf{s}_x, \mathbf{s}_d} V_{\mathbf{s}_d}^T.$$

399 Such a low-rank approximation can be obtained via a degenerate expansion  
 400 of  $\kappa(s, t)$  and has a bounded rank for any specified approximation accuracy.  
 401 That is, the size of  $B_{\mathbf{s}_x, \mathbf{s}_d}$  is bounded. (See [10] for an example of the accuracy  
 402 study.) The subsets  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are also said to be *far-field* clusters and the  
 403 submatrix  $C_{\mathbf{s}_x, \mathbf{s}_d}$  is a far-field interaction/block.

404 • On the other hand, if  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are not well separated, then they are said  
 405 to be *near-field* clusters, and  $C_{\mathbf{s}_x, \mathbf{s}_d} = (\kappa(d_j, x_i))_{x_i \in \mathbf{s}_x, d_j \in \mathbf{s}_d}$  is treated as a  
 406 regular dense block (near-field interaction/block).

407 The FMM further considers the interactions between parent and child clusters  
 408 during the hierarchical partitioning, so that the  $U, V$  basis matrices in (4.7) satisfy  
 409 nested relationships (like in (3.1)). The details can be found in [23] and are not our  
 410 focus here. (Also see [10] particularly for a stable 1D matrix version.) The FMM  
 411 essentially constructs an *FMM matrix* approximation to  $C$  and multiplies it with  $\mathbf{w}$ .  
 412 The complexity of each FMM matrix-vector multiplication is  $O(n)$ .

413 In light of (4.3) and (4.4), a straightforward idea in [13, 24, 44] is to apply the  
 414 standard FMM to  $C$  and  $S$  for fast evaluations of  $\mathbf{f}$  and  $\mathbf{f}'$ . However, in practical  
 415 implementations of secular equation solution methods, it is preferred to write  $f(x)$  in  
 416 the following form so as to avoid cancellation (see, [5, 9, 18]):

$$417 \quad f(x) = 1 + \psi_k(x) + \phi_k(x),$$

418 where the splitting depends on  $k$  (when  $\lambda_k \in (d_k, d_{k+1})$  is to be found):

$$419 \quad (4.8) \quad \psi_k(x) = \sum_{j=1}^k \frac{v_j^2}{d_j - x}, \quad \phi_k(x) = \sum_{j=k+1}^n \frac{v_j^2}{d_j - x}.$$

420 Because of the interlacing property (4.5), all the terms in the sum for  $\psi_k$  or  $\phi_k$  have  
 421 the same sign for  $x \in (d_k, d_{k+1})$ . Furthermore,  $\psi_k$  and  $\phi_k$  capture the behavior of  $f$   
 422 near two poles  $d_k$  and  $d_{k+1}$  respectively.

423 A reliable and widely used scheme to find the roots of  $f(x)$  is given in [28] based  
 424 on a modified Newton's method with a hybrid scheme for rational interpolations of  
 425  $\psi_k(x)$  and  $\phi_k(x)$ . The scheme mixes a middle way method and a fixed weight method  
 426 and is implemented in LAPACK [5]. In the middle way method, rational functions  
 427  $\xi_{k,1}(x) = a_1 + \frac{b_1}{d_k - x}$  and  $\xi_{k,2}(x) = a_2 + \frac{b_2}{d_{k+1} - x}$  are decided to interpolate  $\psi_k$  and  $\phi_k$   
 428 respectively at  $x_k \in (d_k, d_{k+1})$ , so that

$$429 \quad \xi_{k,1}(x_k) = \psi_k(x_k), \quad \xi'_{k,1}(x_k) = \psi'_k(x_k), \quad \xi_{k,2}(x_k) = \phi_k(x_k), \quad \xi'_{k,2}(x_k) = \phi'_k(x_k).$$

430 (We also follow this hybrid scheme to find the first  $n - 1$  roots  $\lambda_1, \lambda_2, \dots, \lambda_{n-1}$ . The  
 431 last root  $\lambda_n$  has only one pole  $d_n$  next to it, so a simple rational interpolation is used  
 432 as in [5, 28]).

433 The modified Newton's method requires evaluations of the functions  $\psi_k, \phi_k, \psi'_k,$   
 434 and  $\phi'_k$  at some  $x_k \in (d_k, d_{k+1})$ ,  $1 \leq k \leq n - 1$ . (Note that even though the summands  
 435 in  $\psi'_k$  and  $\phi'_k$  have the same sign,  $\psi'_k$  and  $\phi'_k$  are used separately in the rational  
 436 interpolations by  $\xi_{k,1}$  and  $\xi_{k,2}$ , respectively [28].) Since these functions all depend on  
 437 individual  $k$ , the standard FMM cannot be applied directly. The reason is that the  
 438 standard FMM handles the evaluation of a kernel  $\kappa(s, t)$  at a fixed set of data points,  
 439 while here it needs to evaluate  $\kappa(s, t)$  at different  $k$ -dependent subsets of  $\{d_j\}_{1 \leq j \leq n}$   
 440 and  $\{x_i\}_{1 \leq i \leq n}$  to produce multiple  $k$ -dependent functions.

441 **4.2.2. Triangular FMM for fast evaluations of  $\psi_k$  and  $\phi_k$ .** To resolve the  
 442 difficulty of applying FMM accelerations to (4.8), we let

$$443 \quad (4.9) \quad \boldsymbol{\psi} = (\psi_1(x_1) \quad \cdots \quad \psi_n(x_n))^T, \quad \boldsymbol{\phi} = (\phi_1(x_1) \quad \cdots \quad \phi_{n-1}(x_{n-1}) \quad 0)^T,$$

$$444 \quad (4.10) \quad \boldsymbol{\psi}' = (\psi'_1(x_1) \quad \cdots \quad \psi'_n(x_n))^T, \quad \boldsymbol{\phi}' = (\phi'_1(x_1) \quad \cdots \quad \phi'_{n-1}(x_{n-1}) \quad 0)^T.$$

446 The key idea is to write

$$447 \quad (4.11) \quad \mathbf{f} = \mathbf{e} + \boldsymbol{\psi} + \boldsymbol{\phi} = \mathbf{e} + C_L \mathbf{w} + C_U \mathbf{w}, \quad \mathbf{f}' = \boldsymbol{\psi}' + \boldsymbol{\phi}' = S_L \mathbf{w} + S_U \mathbf{w},$$

448 where  $\mathbf{e}$  is given in (4.2),  $C_L$  and  $S_L$  are the lower triangular parts of  $C$  and  $S$ ,  
 449 respectively, and  $C_U$  and  $S_U$  are the strictly upper triangular parts of  $C$  and  $S$ ,  
 450 respectively. This suggests that the FMM idea should be applied to the lower and  
 451 upper triangular parts of  $C$  and  $S$  separately. That is, we need a special *triangular*  
 452 *FMM* that can be used to quickly evaluate the triangular matrix-vector products  
 453  $C_L \mathbf{w}$ ,  $C_U \mathbf{w}$ ,  $S_L \mathbf{w}$ ,  $S_U \mathbf{w}$ . We illustrate the triangular FMM in terms of the evaluation  
 454 of  $C_L \mathbf{w}$ . For two subsets  $\mathbf{s}_x$  and  $\mathbf{s}_d$  as in (4.6), we similarly use  $(C_L)_{\mathbf{s}_x, \mathbf{s}_d}$  to denote  
 455 the block of  $C_L$  defined by  $\mathbf{s}_x$  and  $\mathbf{s}_d$ .

- 456 • When  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are neighbor clusters, the interlacing property (4.5) means  
 457  $C_{\mathbf{s}_x, \mathbf{s}_d}$  is a diagonal block of  $C$ . Then  $(C_L)_{\mathbf{s}_x, \mathbf{s}_d}$  is just the lower triangular  
 458 part of  $C_{\mathbf{s}_x, \mathbf{s}_d}$ .
- 459 • When  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are well separated,  $(C_L)_{\mathbf{s}_x, \mathbf{s}_d}$  is a far-field block.  
 460 – If  $\mathbf{s}_x$  is on the right of  $\mathbf{s}_d$  or  $\max \mathbf{s}_x > \min \mathbf{s}_d$ , the interlacing property  
 461 (4.5) means  $(C_L)_{\mathbf{s}_x, \mathbf{s}_d}$  is in the lower triangular part of  $C_L$ . Since  $C_L$  is  
 462 the lower triangular part of  $C$ , (4.7) gives

$$463 \quad (4.12) \quad (C_L)_{\mathbf{s}_x, \mathbf{s}_d} = C_{\mathbf{s}_x, \mathbf{s}_d} \approx U_{\mathbf{s}_x} B_{\mathbf{s}_x, \mathbf{s}_d} V_{\mathbf{s}_d}^T.$$

- 464 – If  $\mathbf{s}_x$  is on the left of  $\mathbf{s}_d$ ,  $(C_L)_{\mathbf{s}_x, \mathbf{s}_d}$  is in the upper triangular part of  $C_L$   
 465 and is thus a zero block. This case can still be accommodated by (4.12),  
 466 with  $B_{\mathbf{s}_x, \mathbf{s}_d} = 0$ .

467 The far-field blocks of  $C_L$  then have the same  $U, V$  basis matrices as those of  $C$ .  
 468 Thus, we can conveniently obtain a lower triangular FMM approximation matrix for  
 469  $C_L$  based on an FMM approximation matrix for  $C$ , just with the difference in the  
 470 lower triangular diagonal blocks and in some zero  $B$  generators. The multiplication  
 471 of the triangular FMM matrix with  $\mathbf{w}$  takes only  $O(n)$  operations. The cost of one  
 472 simultaneous iteration step for all  $x_k$ 's is then  $O(n)$ .

473 **4.2.3. Iterative secular equation solution.** During the modified Newton's  
 474 method, let  $x_k^{(j)}$  be an approximation to the eigenvalue  $\lambda_k$  at the iteration step  $j$ . A  
 475 correction  $\Delta x_k^{(j)}$  is computed to update  $x_k^{(j)}$  as

$$476 \quad (4.13) \quad x_k^{(j+1)} \leftarrow x_k^{(j)} + \Delta x_k^{(j)}.$$

477 (We sometimes write  $x_k^{(j)}$  as  $x_k$  when the focus is not on the iteration steps  $j$ .)

478 We adopt the following stopping criterion from [24]:

$$479 \quad (4.14) \quad |f(x_k^{(j)})| < cn(1 + |\psi(x_k^{(j)})| + |\phi(x_k^{(j)})|)\epsilon_{\text{mach}},$$

480 where  $c$  is a small constant. This stopping criterion can be conveniently checked after  
 481 the FMM-accelerated function evaluations, which is an advantage over a criterion  
 482 in [28]. The factor  $n$  in (4.14) is related to error propagations of general matrix  
 483 multiplications. Although (4.14) might be loose for an extremely large  $n$ , it works  
 484 well in our tests and leads to satisfactory accuracies. It is possible to refine (4.14) to  
 485 a tighter convergence estimate using the backward stability studies of FMM matrix-  
 486 vector multiplication algorithms in [10, 46]. This is our ongoing work.

487 Typically, a very small number of iterations is needed for convergence, similarly  
 488 to the tridiagonal divide-and-conquer algorithm as mentioned in [17]. (In our tests,  
 489 each eigenvalue converges in 2 to 5 iterations on average.) With the total number  
 490 of iterations bounded, the total iterative solution cost for finding all the eigenvalues  
 491 (from one secular equation) is  $O(n)$ .

492 **4.3. Local shifting in triangular FMM for shifted secular equation solu-**  
 493 **tion.** When there are clustered eigenvalues or when updates to previous eigenvalues  
 494 are small, typically the standard secular equation (2.12) is not directly solved. In-  
 495 stead, shifted secular equations are solved for the purpose of stability and accuracy, as  
 496 discussed in [9, 18, 24]. However, it is nontrivial to apply FMM to accelerate shifted  
 497 secular equation solution. In fact, the paper [24] mentions the possibility of FMM  
 498 accelerations for the standard secular equation but does not consider the shifted ones.  
 499 The FMM-accelerated algorithm in [44] does not use shifted secular equations either.

500 In this subsection, we discuss the necessity of shifting and its challenges to FMM  
 501 accelerations. Moreover, we develop a new strategy that makes it feasible to apply  
 502 FMM accelerations to the solution of shifted secular equations. In the following, we  
 503 suppose deflation in Section 4.1 has already been applied.

504 **4.3.1. Shifted secular equation solution and its challenge to FMM ac-**  
 505 **celerations.** During the solution for  $\lambda_k \in (d_k, d_{k+1})$ , if  $\lambda_k$  is very close to  $d_k$  or  $d_{k+1}$ ,  
 506 evaluating  $\frac{v_k^2}{\lambda_k - d_k}$  or  $\frac{v_{k+1}^2}{\lambda_k - d_{k+1}}$  in the secular function with a computed  $\lambda_k$  might lose  
 507 accuracy because of cancellations in the denominator. Without loss of generality, we  
 508 always assume  $\lambda_k$  is closer to  $d_k$ . Let the difference between  $\lambda_k$  and  $d_k$  be

$$509 \quad \eta_k \equiv \lambda_k - d_k,$$

510 which is also said to be the *gap* between  $\lambda_k$  and  $d_k$ . With a very small gap  $\eta_k$ ,  
 511 instead of directly solving for  $\lambda_k$ , the remedy in [9, 18, 24] is to solve a shifted secular  
 512 equation for  $\eta_k$ . For this purpose, we shift the origin to  $d_k$  and rewrite the original  
 513 secular equation (2.12) as the equivalent *shifted secular equation* (see, e.g., [9, 18, 24]):

$$514 \quad (4.15) \quad g_k(y) \equiv f(d_k + y) = 1 + \sum_{j=1}^n \frac{v_j^2}{\delta_{jk} - y} = 0, \quad \text{with}$$

$$515 \quad (4.16) \quad \delta_{jk} = d_j - d_k, \quad j = 1, 2, \dots, n.$$

517 (4.15) is solved for  $y = \eta_k$ . The benefits of this shifting within our context are as  
 518 follows.

519 One benefit is to avoid catastrophic cancellation or division by zero (see, e.g.,  
 520 [5, 9, 24]). (Basically, for computations like  $d_i - \lambda_k$ , although  $d_i - \lambda_k = \delta_{ik} - \eta_k$   
 521 in exact arithmetic, it is preferred to use  $\delta_{ik} - \eta_k$  to avoid cancellation [9, 24].) For  
 522 example, let  $x_k$  be an approximation to  $\lambda_k$ . In exact arithmetic,  $x_k \in (d_k, d_{k+1})$ . At  
 523 each modified Newton iteration, it needs to guarantee  $d_k < \text{fl}(x_k) < d_{k+1}$ . However,  
 524 this might not be satisfied in floating point arithmetic when  $x_k$  is very close to  $d_k$ :

$$525 \quad (4.17) \quad |d_k - x_k| = O(\epsilon_{\text{mach}}) \text{ or smaller,}$$

526 which may lead to cancellation when computing  $d_k - \text{fl}(x_k)$ :

$$527 \quad (4.18) \quad \text{fl}(d_k - \text{fl}(x_k)) = o(\epsilon_{\text{mach}}) \quad \text{or even} \quad \text{fl}(d_k - \text{fl}(x_k)) = 0.$$

528 This will cause stability issues in the numerical solutions of the standard secular  
 529 function:  $\text{fl}\left(\frac{v_k^2}{d_k - \text{fl}(x_k)}\right)$  either is highly inaccurate or becomes  $\infty$ .

530 Note that (4.17) and (4.18) are still possible even if deflation in Section 4.1 has  
 531 been applied with a tolerance  $\tau$  that is not too small. To see this, suppose  $v_k =$   
 532  $O(\tau) \geq \tau$  and the exact root  $\lambda_k$  satisfies  $|\lambda_k - d_j| \gg v_j^2$  for  $j \neq k$ . Substituting  $\lambda_k$

533 into the secular equation (2.12) yields  $\frac{v_k^2}{d_k - \lambda_k} = -1 + \sum_{j \neq k}^n \frac{v_j^2}{\lambda_k - d_j} = O(1)$ . In this  
 534 case,  $\lambda_k$  shall be very close to  $d_k$  in the following sense:

$$535 \quad |d_k - \lambda_k| = v_k^2 \cdot O(1) = O(\tau^2).$$

536 If  $\tau = O(\epsilon_{\text{mach}}^{1/2})$  which is not extremely small, we can have (4.17) so that (4.18) may  
 537 happen when solving the standard secular equation.

538 Another benefit for solving the shifted equation is faster convergence. It is ob-  
 539 served in our tests that computing with  $\eta_k$  instead of  $\lambda_k$  can speed up the convergence  
 540 of the modified Newton's method. To illustrate this, suppose  $\lambda_k$  is solved directly from  
 541 the standard secular equation (2.12), then the approximation  $x_k^{(j)}$  at iteration step  $j$   
 542 is updated as in (4.13). Suppose  $|\lambda_k| = O(1)$  and  $|\eta_k| = |\lambda_k - d_k| = O(\epsilon_{\text{mach}})$ . Since  
 543  $x_k^{(j)}$  converges to  $\lambda_k$  as  $j \rightarrow \infty$ , we also have  $|x_k^{(j)}| = O(1)$  and  $|x_k^{(j)} - d_k| = O(\epsilon_{\text{mach}})$   
 544 after some iterations. In the modified Newton's method, the correction  $\Delta x_k^{(j)}$  ap-  
 545 proaches 0 as  $j$  increase. This may lead to loss of digits in the updated  $x_k^{(j+1)}$ :  
 546  $\text{fl}(x_k^{(j+1)}) = \text{fl}(x_k^{(j)} + \Delta x_k^{(j)}) = \text{fl}(x_k^{(j)})$ . As a result, the iteration stagnates. On the  
 547 other hand, if  $\eta_k$  is solved from the shifted secular equation (4.15), as in [5, 9, 18],  
 548 the update (4.13) is replaced by

$$549 \quad (4.19) \quad y_k^{(j+1)} \leftarrow y_k^{(j)} + \Delta x_k^{(j)},$$

550 where  $y_k^{(j)} = x_k^{(j)} - d_k$  is an approximation to  $\eta_k$  at step  $j$  of the iterative solution.  
 551 Although (4.13) and (4.19) are equivalent in exact arithmetic, the latter preserves a  
 552 lot more digits of accuracy since  $|y_k^{(j)}| = O(\epsilon_{\text{mach}})$ .

553 These discussions illustrate the importance of solving the shifted secular equation  
 554 (4.15) instead of the original equation (2.12). However, in an FMM-accelerated scheme  
 555 where all  $\lambda_k$ 's are solved simultaneously, it is not plausible to shift the secular equation  
 556 simultaneously for all  $\lambda_k$ 's. The reason is the shift in (4.15) depends on each individual  
 557 eigenvalue and there is no such a uniform shift that would work for all  $\lambda_k$ 's.

558 To see this, let  $y_k = x_k - d_k$  be an approximation to  $\eta_k$  during the iterative  
 559 solution of (4.15). The evaluations of  $g_k(y)$  in (4.15) at  $y = y_k$  for all  $k = 1, 2, \dots, n$   
 560 can be assembled into the matrix form

$$561 \quad (4.20) \quad \mathbf{g} = \mathbf{e} + \hat{C}\mathbf{w}, \quad \text{with}$$

$$562 \quad \mathbf{g} = (g_1(y_1) \quad \cdots \quad g_n(y_n))^T, \quad \hat{C} = \left( \frac{1}{\delta_{jk} - y_k} \right)_{1 \leq k, j \leq n},$$

564 where  $\delta_{jk}$  is given in (4.16). Recall that when the FMM is used to accelerate the  
 565 matrix-vector product  $C\mathbf{w}$  in (4.11), it relies on the separability of  $s$  and  $t$  in a  
 566 degenerate approximation of  $\kappa(s, t) = \frac{1}{s-t}$ . (Note that in  $\kappa(d_j, x_k)$ ,  $x_k$  only involves  
 567 the row index  $k$  and  $d_j$  only involves the column index  $j$ , so that the separability can  
 568 be understood in terms of the row and column indices.) However, to evaluate  $\hat{C}\mathbf{w}$  in  
 569 (4.20), we have

$$570 \quad (4.21) \quad \kappa(d_j, x_k) = \kappa(d_j - d_k, x_k - d_k) = \kappa(\delta_{jk}, y_k).$$

571  $\delta_{jk}$  involves both the row and column indices, so that the separability in terms of  
 572 the row and column indices does not hold. Therefore, we need to adapt the FMM to  
 573 accelerate the shifted matrix-vector multiplication in (4.20).

574 **4.3.2. FMM accelerations with local shifting.** In this subsection, we pro-  
 575 pose a strategy called *local shifting* that makes it feasible to apply triangular FMM  
 576 accelerations to solve (4.15). As mentioned in Section 4.2.1, multiple terms involving  
 577  $x_k - d_j$  are assembled into matrices in order to apply FMM accelerations. See, e.g.,  
 578 (4.3). When  $|x_k - d_j|$  is small, the shifting helps get  $x_k - d_j$  accurately. However, when  
 579  $k$  is not near  $j$  or when  $|k - j|$  is large,  $x_k - d_j$  can actually be computed accurately  
 580 *without involving any shift*  $d_k$ . To see this, recall that  $d_k < x_k < d_{k+1}$  and also after  
 581 deflation in Section 4.1, we have (4.1) holds and  $|d_j - d_{j+1}| \geq \frac{v_j^2 + v_{j+1}^2}{|v_j v_{j+1}|} \geq 2\tau$  for all  $j$ .  
 582 Thus, for  $j \neq k, k + 1$ ,

$$583 \quad (4.22) \quad |x_k - d_j| \geq \min(|d_k - d_j|, |d_{k+1} - d_j|) \geq 2(|k - j| - 1)\tau.$$

584 Hence,  $|x_k - d_j|$  is not too small and  $x_k - d_j$  can be computed accurately when  $|k - j|$   
 585 is large.

586 Following this justification, we have the basic ideas of our local shifting strategy:  
 587 (i) small gaps resulting from shifting is used just in near-field interactions of the FMM,  
 588 which does not interfere with the FMM rank structure; (ii) it is safe to not shift the  
 589 numerical eigenvalues in far-field interactions, which makes it feasible to exploit the  
 590 rank structure. More specifically, the major components are as follows.

- 591 1. For  $k = 1, 2, \dots, n$ , the shifted secular equations (4.15) are solved together  
 592 for the gaps  $\eta_k = \lambda_k - d_k$  (so as to get the roots  $\lambda_k$  of the original secular  
 593 equation (2.12)). An intermediate gap during the iterative solution looks like  
 594  $y_k = x_k - d_k$ . The relevant function evaluations in the iterative solutions are  
 595 assembled into matrix-vector products like in (4.20).
- 596 2. The FMM is used to accelerate the resulting matrix-vector products like  $\hat{C}\mathbf{w}$   
 597 in (4.20) as follows. Suppose two subsets  $\mathbf{s}_x$  and  $\mathbf{s}_d$  like in (4.6) are well  
 598 separated. As mentioned above, for  $x_k \in \mathbf{s}_x$  and  $d_j \in \mathbf{s}_d$ ,  $x_k$  and  $d_j$  are far  
 599 away from each other and  $|k - j|$  is large, so  $x_k - d_j$  can then be computed  
 600 accurately because of (4.22). Thus, we can recover  $x_k$  from  $d_k + y_k$  to directly  
 601 exploit the low-rank structure like in (4.7). As a result, the far-field block  
 602  $\hat{C}_{\mathbf{s}_x, \mathbf{s}_d}$  of  $\hat{C}$  is now just a block of  $C$  in (4.3):

$$603 \quad \hat{C}_{\mathbf{s}_x, \mathbf{s}_d} = (\kappa(\delta_{jk}, y_k))_{x_k \in \mathbf{s}_x, d_j \in \mathbf{s}_d} = (\kappa(d_j, x_k))_{x_k \in \mathbf{s}_x, d_j \in \mathbf{s}_d} = C_{\mathbf{s}_x, \mathbf{s}_d}$$

- 604 3. On the other hand, when two subsets  $\mathbf{s}_x$  and  $\mathbf{s}_d$  are not well separated, the  
 605 near-field interaction  $\hat{C}_{\mathbf{s}_x, \mathbf{s}_d}$  is dense and each entry  $\kappa(\delta_{jk}, y_k)$  can be evaluated  
 606 accurately via  $y_k$  and the gap  $\delta_{jk}$ . This has no impact on the structures  
 607 needed for FMM accelerations.
- 608 4. These ideas are then combined with the triangular FMM in Section 4.2.2  
 609 to stably and quickly perform function evaluations like (4.20) and solve the  
 610 shifted secular equations.

611 This local shifting strategy successfully integrates shifting into triangular FMM  
 612 accelerations. As a result, we can quickly and reliably solve the shifted secular equa-  
 613 tions (4.15) via the modified Newton's method. The overall complexity to find all the  
 614  $n$  roots is still  $O(n)$ . In addition, since the relevant functions are now evaluated more  
 615 accurately than with the method in [44], the convergence is also improved. When the  
 616 iterative solution of the shifted secular equations converge, we can use the resulting  
 617 gaps  $\eta_k$  to recover the desired eigenvalues as

$$618 \quad (4.23) \quad \lambda_k = d_k + \eta_k, \quad k = 1, 2, \dots, n,$$

619 The local shifting strategy can also be used to stably apply FMM accelerations  
620 to other operations like finding the eigenmatrix. See the next subsection.

621 **4.4. Structured eigenvectors via FMM with local shifting.** With the iden-  
622 tified eigenvalues  $\lambda_k$  in (4.23), the eigenvectors can be obtained stably as in [24]. An  
623 eigenvector corresponding to  $\lambda_k$  looks like

$$624 \quad (4.24) \quad \mathbf{q}_k = \left( \frac{\hat{v}_1}{d_1 - \lambda_k} \quad \cdots \quad \frac{\hat{v}_k}{d_k - \lambda_k} \quad \cdots \quad \frac{\hat{v}_n}{d_n - \lambda_k} \right)^T,$$

625 where  $\hat{\mathbf{v}} \equiv (\hat{v}_1 \quad \cdots \quad \hat{v}_n)^T$  is given by Löwner's formula

$$626 \quad (4.25) \quad \hat{v}_i = \sqrt{\frac{\prod_j (\lambda_j - d_i)}{\prod_{j \neq i} (d_j - d_i)}}, \quad i = 1, 2, \dots, n.$$

627 To quickly form  $\hat{\mathbf{v}}$ , the standard FMM acceleration would look like the following [24].  
628 Rewrite (4.25) as

$$629 \quad (4.26) \quad \log \hat{v}_i = \frac{1}{2} \sum_{j=1}^n \log(|d_i - \lambda_j|) - \frac{1}{2} \sum_{j=1, j \neq i}^n \log |d_i - d_j|.$$

630 Now, let  $G_1 = (\log |d_i - \lambda_j|)_{n \times n}$ ,  $G_2 = (\log |d_i - d_j|)_{n \times n}$ , where the diagonals of  $G_2$   
631 are set to be zero. Then

$$632 \quad (4.27) \quad \log \hat{\mathbf{v}} = \frac{1}{2} (G_1 \mathbf{e} - G_2 \mathbf{e}).$$

633  $G_1 \mathbf{e}$  and  $G_2 \mathbf{e}$  can thus be quickly evaluated by the FMM with the kernel  $\log |s - t|$ .

634 As in [24, 44], the eigenvectors are often normalized to form an orthogonal matrix

$$635 \quad (4.28) \quad \hat{Q} = \left( \frac{\hat{v}_i b_j}{d_i - \lambda_j} \right)_{n \times n}, \quad \text{with}$$

$$636 \quad (4.29) \quad \mathbf{b} \equiv (b_1 \quad \cdots \quad b_n)^T, \quad b_j = \left( \sum_{i=1}^n \frac{\hat{v}_i^2}{(d_i - \lambda_j)^2} \right)^{-1/2}.$$

638 The vector  $\mathbf{b}$  can be quickly obtained via the FMM with the kernel  $\kappa(s, t) = \frac{1}{(s-t)^2}$ .

639  $\hat{Q}$  is a Cauchy-like matrix which gives a structured form of the eigenvectors. The  
640 FMM with the kernel  $\kappa(s, t) = \frac{1}{s-t}$  can be used to quickly multiply  $\hat{Q}$  to a vector.

641 Again, with the same reasons as before, it is challenging to stably apply the  
642 standard FMM to accelerate operations like the evaluations of  $\log \mathbf{v}$  in (4.27) and  $\mathbf{b}$  in  
643 (4.29) and the application of  $\hat{Q}$  to a vector. On the other hand, just like the discussions  
644 in Section 4.3.2, we can integrate the local shifting strategy into FMM accelerations,  
645 just with appropriate kernels  $\kappa(s, t)$ . For example, with the gaps  $\eta_k$  solved from the  
646 shifted secular equation solution, it is preferred to use  $\delta_{ik} - \eta_k$  in place of  $d_i - \lambda_k$  in  
647 the computation of some entries of  $\mathbf{q}_k$  for accuracy purpose [5, 9, 18, 24] when  $d_i$  and  
648  $\lambda_k$  are very close. Note that, with  $\delta_{jk}$  in (4.16), (4.24) can be written as

$$649 \quad (4.30) \quad \mathbf{q}_k = \left( \frac{\hat{v}_1}{\delta_{1k} - \eta_k} \quad \cdots \quad \frac{\hat{v}_k}{-\eta_k} \quad \cdots \quad \frac{\hat{v}_n}{\delta_{nk} - \eta_k} \right)^T.$$

650 When an entry of  $\mathbf{q}_k$  belongs to a near-field block of  $\hat{Q}$ , its representation in (4.30) is  
651 used. Otherwise, we use its form in (4.24). This preserves the far-field rank structure.

652 Thus, FMM accelerations with local shifting can be used to reliably represent and  
 653 apply  $\hat{Q}$  or  $\hat{Q}^T$ . Note that

$$654 \quad (4.31) \quad \hat{Q} = \text{diag}(\hat{\mathbf{v}}) \left( \frac{1}{d_i - \lambda_j} \right)_{n \times n} \text{diag}(\mathbf{b}),$$

655 so that  $\hat{Q}$  can be stored just via five vectors:

$$656 \quad (4.32) \quad \hat{\mathbf{v}}, \mathbf{b}, \mathbf{d} \equiv (d_1 \ \cdots \ d_n)^T, \ \boldsymbol{\lambda} \equiv (\lambda_1 \ \cdots \ \lambda_n)^T, \ \boldsymbol{\eta} \equiv (\eta_1 \ \cdots \ \eta_n)^T.$$

657 Here, we have the storage of one more vector  $\boldsymbol{\eta}$  than that in [44]. This only slightly  
 658 increase the storage, but the stability is significantly enhanced.

659 **4.5. Overall eigendecomposition and structure of the eigenmatrix.** The  
 660 overall conquering framework is similar to [44], but with all the new stability strategies  
 661 integrated. The conquering process is performed following the postordered traversal  
 662 of the HSS tree  $\mathcal{T}$  of  $A$ , where at each node  $i \in \mathcal{T}$ , a local eigenproblem is solved. For  
 663 a leaf node  $i$ , suppose  $\hat{D}_i$  is the (small) diagonal generator resulting from the overall  
 664 dividing process. We just compute the dense eigenproblem  $\hat{D}_i = Q_i \Lambda_i Q_i^T$ . Then  $Q_i$   
 665 is a *local eigenmatrix* associated with  $i$ .

666 For a non-leaf node  $p$  with children  $i$  and  $j$ , the local eigenproblem is to find  
 667 an eigendecomposition like in (2.10) based on (2.6) and (2.7). However, unlike (2.9)  
 668 where a diagonal plus low-rank update eigendecomposition is computed, it is *nec-*  
 669 *essary to reorder* the diagonal entries of  $\text{diag}(\Lambda_i, \Lambda_j)$  in order to explore the FMM  
 670 structures that rely on the locations of the eigenvalues. Let  $P_p$  represent a sequence of  
 671 permutations for deflation and for ordering the diagonal entries of  $\text{diag}(\Lambda_i, \Lambda_j)$  from  
 672 the smallest to the largest. Also let the eigendecomposition of the *permuted* diagonal  
 673 plus low-rank update problem be

$$674 \quad (4.33) \quad P_p[\text{diag}(\Lambda_i, \Lambda_j) + \hat{Z}_p \hat{Z}_p^T] P_p^T = \hat{Q}_p \Lambda_p \hat{Q}_p^T,$$

675 where  $\hat{Z}_p$  is given in (2.8). Write  $D_p$  in (2.7) as  $\hat{D}_p$  since  $D_p$  is likely updated after  
 676 the hierarchical dividing process. Then we have the following eigendecomposition:

$$677 \quad (4.34) \quad \hat{D}_p = Q_p \Lambda_p Q_p^T, \quad \text{with} \quad Q_p = \text{diag}(Q_i, Q_j) P_p^T \hat{Q}_p,$$

678 where  $Q_i$  and  $Q_j$  are eigenmatrices of  $\hat{D}_i$  and  $\hat{D}_j$  obtained in steps  $i$  and  $j$ , respectively.  
 679 Then the conquering process proceeds similarly.

680 Here for convenience, we say  $Q_p$  is a *local eigenmatrix* and  $\hat{Q}_p$  is an *interme-*  
 681 *diante eigenmatrix*. The difference between the two is that a local eigenmatrix is an  
 682 eigenmatrix of a local HSS block while the latter is an eigenmatrix of a diagonal plus  
 683 low-rank update problem. A local eigenmatrix is formed by a sequence of interme-  
 684 diate eigenmatrices. Since  $\hat{Q}_p \Lambda_p \hat{Q}_p^T$  in (4.33) is obtained by solving  $r$  consecutive  
 685 rank-1 update eigenproblems, the intermediate eigenmatrix  $\hat{Q}_p$  is the product of  $r$   
 686 Cauchy-like matrices like in (4.28). Of course, when FMM accelerations and deflation  
 687 are applied, the eigendecomposition is approximate.

688 Then the overall eigenmatrix  $Q$  is given in terms of all the intermediate eigenma-  
 689 trices, organized with the aid of the tree  $\mathcal{T}$ . Here, we give an accurate description of  
 690 its structure as follows.

691 LEMMA 4.1. *Assemble all the intermediate eigenmatrices and permutation matri-*  
 692 *ces corresponding to the nodes at a level  $l$  of  $\mathcal{T}$  as*

$$693 \quad (4.35) \quad Q^{(l)} = \text{diag}(\hat{Q}_i, i: \text{at level } l \text{ of } \mathcal{T}), \quad P^{(l)} = \text{diag}(P_i, i: \text{at level } l \text{ of } \mathcal{T}).$$

694 Then the final eigenmatrix  $Q$  has the form (illustrated in Figure 4.1)

$$695 \quad (4.36) \quad Q = Q^{(L)} \prod_{l=l_{\max}-1}^0 (P^{(l)} Q^{(l)}),$$

696 where level  $l_{\max}$  is the leaf level of  $\mathcal{T}$  and  $\text{root}(\mathcal{T})$  is at level 0. In addition,  $Q$  also  
697 corresponds to (4.34) with  $p$  set to be  $\text{root}(\mathcal{T})$ .

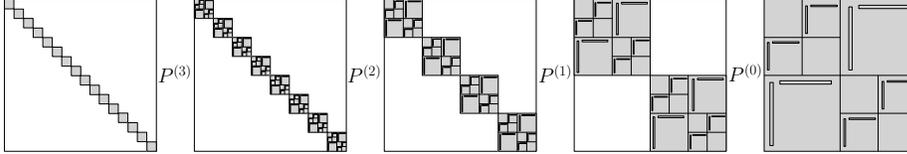


FIG. 4.1. Illustration of the structure of the eigenmatrix  $Q$ , where  $l_{\max} = 4$  and each structured diagonal block (marked in gray) is for an intermediate eigenmatrix  $\hat{Q}_p$  associated with a nonleaf node  $p$ .

698 Thus,  $Q$  can be understood in terms of either (4.36) or the local eigenmatrices.  
699 Lemma 4.1 gives an efficient way to apply  $Q$  or  $Q^T$  to a vector, where the triangular  
700 FMM with local shifting is again used to multiply the intermediate eigenmatrices  
701 with vectors. Note that with a very similar procedure, a local eigenmatrix  $Q_i$  or its  
702 transpose can be conveniently applied to a vector. Such an application process is used  
703 to multiply the local eigenmatrices  $Q_i^T$  and  $Q_j^T$  to  $Z_p$  as in (2.8) to quickly form  $\hat{Z}_p$   
704 used in (4.33).

705 The main algorithms used in SuperDC are shown in the supplementary materials.  
706 When  $A$  is given in terms of an HSS form with HSS rank  $r$ , the total complexity for  
707 computing the eigendecomposition (1.1) can be counted following [44, Section 3.1]  
708 and is  $O(r^2 n \log^2 n)$ . Note that the use of all the new stability techniques here does  
709 not change the overall complexity. Every local eigenmatrix  $\hat{Q}_i$  is represented by a  
710 sequence of  $r$  Cauchy-like matrices like in (4.28). Each such a Cauchy-like matrix is  
711 stored with the aid of five vectors like in (4.32). The storage for  $Q$  is then  $O(rn \log n)$   
712 and the cost to apply  $Q$  or  $Q^T$  to a vector is  $O(rn \log n)$  as in [44].

713 **5. Numerical experiments.** In this section, we perform a comprehensive test  
714 of the SuperDC eigensolver with different types of matrices and demonstrate its effi-  
715 ciency and accuracy. We compare SuperDC with the following methods.

- 716 • BandDC: a usual divide-and-conquer eigensolver that takes advantage of  
717 banded structures following the framework in [3]. A sequence of rank-1 up-  
718 dating problems is obtained based on the banded form in each dividing step  
719 and is then solved in the conquering stage. The same deflation tolerances as  
720 SuperDC are used.
- 721 • HSSBIS: an HSS bisection eigensolver [48] that takes advantage of fast HSS  
722 LDL factorization update for inertia evaluations. The stopping criterion of  
723 bisection is the same as the deflation tolerance of SuperDC.
- 724 • eig: the highly optimized Matlab eig function.

725 In order to run comparisons for larger matrix sizes in Matlab, we use BandDC  
726 and HSSBIS to compute only the eigenvalues (with accuracies comparable to those  
727 from SuperDC), which also gives them advantages over SuperDC. For HSSBIS, we  
728 decide the initial search region with  $\tilde{\rho}(A) \equiv \sqrt{\|A\|_1 \|A\|_\infty} \geq \|A\|_2$  as an estimate of

729 the spectral radius of  $A$ . We use the following accuracy measurements:

$$\begin{aligned}
 \gamma &= \max_{1 \leq k \leq n} \frac{\|A\mathbf{q}_k - \lambda_k \mathbf{q}_k\|_2}{\sqrt{n}\|A\|_2} \quad (\text{residual}), \\
 \theta &= \max_{1 \leq k \leq n} \frac{\|Q^T \mathbf{q}_k - \mathbf{e}_k\|_2}{\sqrt{n}} \quad (\text{loss of orthogonality}), \\
 \delta_s &= \frac{\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_2}{\|\boldsymbol{\lambda}^*\|_2}, \quad \delta_\infty = \frac{\|\boldsymbol{\lambda} - \boldsymbol{\lambda}^*\|_\infty}{\|\boldsymbol{\lambda}^*\|_\infty}, \quad \delta_m = \max_{1 \leq k \leq n} \frac{|\lambda_k^* - \lambda_k|}{|\lambda_k^*|} \quad (\text{errors}),
 \end{aligned}$$

731 where  $\boldsymbol{\lambda}^* = (\lambda_1^* \ \cdots \ \lambda_n^*)^T$  are eigenvalues from `eig` and are considered as the exact  
 732 results.

733 To measure the efficiency, we count the *flops* (total number of floating point  
 734 arithmetic operations), the *storage* (total number of nonzeros to store the structured  
 735 eigenmatrix in SuperDC or the dense eigenmatrix in `eig`), and the *timing* (seconds  
 736 elapsed when the call of an eigensolver routine is completed). In the flop count, if a  
 737 built-in routine is used to perform standard operations, we use known flop counts like  
 738 those given in relevant references such as [17, 39].

739 SuperDC is available at <https://www.math.purdue.edu/~xiaj>. It is implemented  
 740 fully in Matlab. The triangular FMM routine of SuperDC is developed based on  
 741 a code used in [10], and its accuracy during each call is set to reach full machine  
 742 precision. In all the tests, the leaf-level diagonal block size of the HSS forms is 2048.  
 743 The tests are performed with four 2.60GHz cores and 80GB memory on a node at a  
 744 cluster of Purdue RCAC. The request of 80GB memory is just to accommodate the  
 745 need of `eig` for larger matrices.

746 **EXAMPLE 1.** First, we consider a symmetric tridiagonal matrix  $A$ . For our Su-  
 747 perDC eigensolver, the HSS representation of  $A$  can be explicitly written out without  
 748 any extra cost and its HSS rank is  $r = 2$  [49]. (The HSS structure does not rely  
 749 on the actual nonzero entries, which are 3 on the main diagonal and  $-1$  on the first  
 750 superdiagonal and subdiagonal. Other numbers such as random ones are also tested  
 751 with similar performance observed.) The size  $n$  of  $A$  in the test ranges from 8192 to  
 752 1,048,576. We use  $\tau = 10^{-10}$  in the deflation criterion in Section 4.1.

753 The timing of BandDC, HSSBIS, `eig`, and SuperDC is reported in Figure 5.1(a).  
 754 The storage for the eigenmatrix  $Q$  from `eig` and SuperDC is given in Figure 5.1(b).  
 755 The costs of SuperDC are given in Figure 5.1(c), in terms of the flops to get the  
 756 eigendecomposition and the flops to apply  $Q$  to a vector. SuperDC achieves roughly  
 757 linear complexity in the timing, flops, and storage. Both BandDC and HSSBIS follow  
 758 quadratic trends in timing, though HSSBIS is quite slower. `eig` has a cubic trend in  
 timing and an obvious quadratic storage (which is just  $n^2$  for storing the dense  $Q$ ).

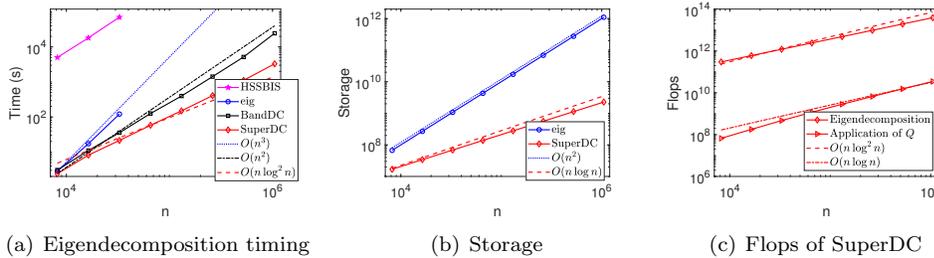


FIG. 5.1. Example 1. *Timing, storage, and flops.*

759

760 SuperDC is faster than BandDC and HSSBIS for all the tested sizes, and its  
 761 breakeven point with eig is around  $n = 8192$ . With  $n = 32,768$ , SuperDC is already  
 762 about 6 times as fast as eig and takes only about 6% of the storage for the eigenmatrix.  
 763 Note that eig runs out of memory for larger  $n$  due to the dense eigenmatrix, while  
 764 SuperDC takes much less memory and can reach much larger  $n$ .

765 The conquering stage is usually much more time-consuming than the dividing  
 766 stage. For example, for  $n = 65,536$ , the dividing stage of SuperDC needs just 1.4 sec-  
 767 onds and the conquering stage takes 56.4 seconds. Thus, our strategy for minimizing  
 768  $\text{colsize}(Z_p)$  or reducing the number of rank-one updates is beneficial for the efficiency  
 769 of the eigensolver since it reduces the amount of work in the conquering stage.

770 Table 5.1 shows the accuracy of SuperDC. The eigenvalues and eigendecompo-  
 771 sitions are computed accurately with numerically orthogonal eigenvectors. BandDC  
 772 and HSSBIS reach comparable accuracies which are then not reported.

TABLE 5.1

Example 1. Accuracy of SuperDC, where some errors are not reported since eig runs out of memory, and  $\gamma$  and  $\theta$  are not available for  $n \geq 262,144$  since they take too long to compute.

$n$	4,096	8,192	16,384	32,768	65,536	131,072
$\gamma$	$1.2e - 15$	$6.4e - 15$	$1.1e - 13$	$9.4e - 14$	$7.5e - 14$	$5.3e - 14$
$\theta$	$1.8e - 14$	$2.9e - 14$	$3.8e - 14$	$5.5e - 14$	$8.6e - 14$	$1.2e - 13$
$\delta_s$	$2.6e - 16$	$4.6e - 16$	$1.3e - 13$	$9.4e - 14$		
$\delta_\infty$	$8.9e - 16$	$1.2e - 14$	$8.0e - 12$	$6.3e - 12$		
$\delta_m$	$9.7e - 16$	$1.2e - 14$	$8.0e - 12$	$6.3e - 12$		

773 EXAMPLE 2. In this example, we test a symmetric matrix  $A$  which is sparse and  
 774 nearly banded.  $A$  has a banded form with half bandwidth 5 together with some  
 775 nonzero entries away from the band. The HSS form for  $A$  can be explicitly written  
 776 out with the method in [49] and has HSS rank 10. The nonzero entries away from the  
 777 band are introduced by modifying some HSS generators. The main diagonal entries  
 778 are set as 30 and the other entries in the band are set as  $-10$  so that the upper bound  
 779 for all  $\|B_k\|_2$  in Proposition 3.1 is  $\beta = 35.1 \gg 1$ . The size  $n$  in the test ranges from  
 780 8192 to 1,048,576. We use  $\tau = 10^{-10}$  in the deflation criterion.

781 The entries away from the band break the banded structure of  $A$ . Still, BandDC  
 782 can be conveniently adapted to  $A$ . The efficiency benefit of SuperDC over eig becomes  
 783 even more significant, as shown in Figure 5.2. At  $n = 32,768$ , SuperDC is about 8  
 784 times as fast as eig and takes only about 7% of the storage for the eigenmatrix. Again,  
 785 eig runs out of memory when  $n$  increases, but SuperDC works for much larger  $n$  and  
 786 demonstrates nearly linear complexity. At  $n = 1,048,576$ , SuperDC is about 12 times  
 787 as fast as BandDC (which does not even compute the eigenmatrix).

788 We also show the advantage of our new dividing strategy over the original one  
 789 (2.2). In Table 5.2, we show the norm growth of the  $B, D$  generators after the dividing  
 790 stage. For the initial  $B, D$  generators of the original HSS form, let  $\tilde{B}, \tilde{D}$  denote the  
 791 updated generators after the entire dividing stage is finished. Let

$$792 \rho_B = \max_{i < \text{root}(\mathcal{T})} \|B_i\|_2, \quad \rho_D = \max_{i: \text{leaf}} \|D_i\|_2, \quad \rho_{\tilde{B}} = \max_{i < \text{root}(\mathcal{T})} \|\tilde{B}_i\|_2, \quad \rho_{\tilde{D}} = \max_{i: \text{leaf}} \|\tilde{D}_i\|_2.$$

793 When  $n$  increases, the number of levels in the HSS tree  $\mathcal{T}$  increase and  $\rho(D)$  and  
 794  $\rho(B)$  stay about the same for all  $n$ . However,  $\rho_{\tilde{B}}$  and  $\rho_{\tilde{D}}$  grow exponentially with the  
 795 original dividing stage, as predicted by Proposition 3.1. This poses a stability risk.

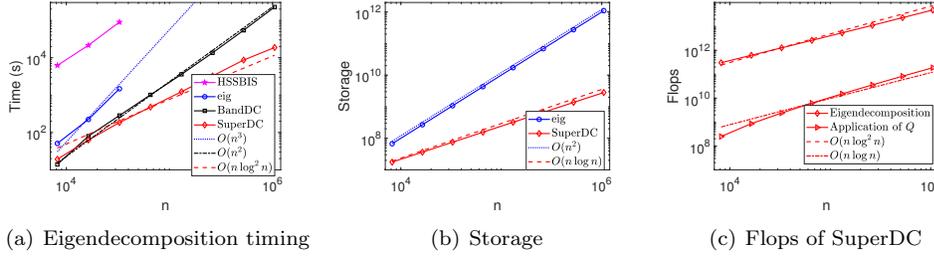


FIG. 5.2. Example 2. *Timing, storage, and flops.*

796 (Note that  $\rho_{\tilde{D}}$  has a larger magnitude than  $\rho_{\tilde{B}}$ , which is consistent with Proposition  
 797 3.1.) In contrast, the growth of  $\rho_{\tilde{D}}$  and  $\rho_{\tilde{B}}$  with our new dividing strategy is much  
 slower and roughly follows the linear growth pattern as predicted by Proposition 3.2.

TABLE 5.2  
 Example 2. *Norms of the D, B generators before and after the dividing stage.*

$n$		4,096	8,192	16,384	32,768	65,536	131,072
Number of levels $l_{\max}$		2	3	4	5	6	7
Initial	$\rho_B$	3.5e1	3.5e1	3.5e1	3.5e1	3.5e1	3.5e1
	$\rho_D$	7.0e1	7.0e1	7.0e1	7.0e1	7.0e1	7.0e1
After the original dividing strategy	$\rho_{\tilde{B}}$	3.5e1	7.4e2	5.5e5	3.0e11	9.2e22	8.5e45
	$\rho_{\tilde{D}}$	1.7e3	1.1e6	6.1e11	1.8e23	1.7e46	1.4e92
After the new dividing strategy	$\rho_{\tilde{B}}$	3.5e1	4.5e1	8.3e1	1.6e2	3.3e2	6.6e2
	$\rho_{\tilde{D}}$	7.3e1	1.3e2	3.0e2	6.3e2	1.3e3	2.6e3

798 In Table 5.3, we report the accuracies of SuperDC. The accuracies associated  
 799 with the original dividing strategy deteriorate as  $n$  gets larger and the results are  
 800 highly inaccurate for  $n \geq 32,768$  so they are not shown. In contrast, the new dividing  
 801 strategy yields nice accuracies. This accuracy difference can be understood as follows.  
 802 For a leaf node  $k$ , we need to use a (backward stable) dense method to compute a  
 803 numerical eigendecomposition of the updated  $\tilde{D}$  generators (see Lemma 2.1) with  
 804 backward error  $\Delta\tilde{D}_k$  (see, e.g., [24]):  
 805

$$806 \quad (5.1) \quad \tilde{D}_k = Q_k \Lambda_k Q_k^T + \Delta\tilde{D}_k, \quad \text{with} \quad \|\Delta\tilde{D}_k\|_2 = O(\|\tilde{D}_k\|_2 \epsilon_{\text{mach}}).$$

807 By Propositions 3.1 and 3.2,  $\|\Delta\tilde{D}_k\|_2$  is roughly in the magnitude of  $O(\beta^{n/2} \epsilon_{\text{mach}})$   
 808 with the original dividing strategy (2.2), or  $O(\frac{n}{2} \epsilon_{\text{mach}})$  with the new dividing strategy.  
 809 Therefore, the original dividing strategy will likely introduce larger errors.

810 *Remark 5.1.* In (5.1), we often have  $\|\tilde{D}_k\|_2$  in the magnitude of  $O(\|A\|_2)$  so  
 811 that  $O(\|\tilde{D}_k\|_2 \epsilon_{\text{mach}})$  is roughly  $O(\|A\|_2 \epsilon_{\text{mach}})$ . Therefore, an absolute error of or-  
 812 der  $O(\|A\|_2 \epsilon_{\text{mach}})$  is introduced to the eigenvalues, resulting in loss of digits for those  
 813 eigenvalues that are tiny (when  $\|A\|_2$  is large). Thus, the maximum relative errors  
 814  $\delta_m$  for those tiny eigenvalues may be larger, while other measurements such as  $\delta_s$  and  
 815  $\gamma$  are still well controlled. In addition, because of the complex nature of SuperDC,  
 816 it is possible that some other weakness in the numerical stability may still present.  
 817 Our ongoing work is to attempt to perform a comprehensive backward stability study

TABLE 5.3

Example 2. Accuracy of SuperDC, where some errors are not reported since eig runs out of memory, and  $\gamma$  and  $\theta$  are not available for  $n \geq 262,144$  since they take too long to compute.

	$n$	4,096	8,192	16,384	32,768	65,536	131,072
Original dividing	$\gamma$	$3.6e-13$	$9.5e-12$	$1.5e-5$			
	$\theta$	$2.2e-13$	$3.0e-13$	$5.3e-13$			
	$\delta_s$	$4.7e-13$	$4.2e-13$	$9.7e-7$			
	$\delta_\infty$	$1.7e-11$	$1.6e-11$	$3.0e-5$			
	$\delta_m$	$1.8e-11$	$1.6e-9$	$1.2e-4$			
New dividing	$\gamma$	$9.4e-14$	$1.8e-12$	$3.7e-13$	$3.4e-13$	$5.6e-13$	$1.2e-12$
	$\theta$	$1.9e-13$	$4.4e-13$	$5.6e-13$	$1.1e-12$	$1.4e-12$	$2.0e-12$
	$\delta_s$	$1.6e-14$	$3.2e-12$	$5.3e-13$	$2.6e-13$		
	$\delta_\infty$	$6.0e-13$	$1.5e-10$	$2.2e-11$	$1.1e-11$		
	$\delta_m$	$1.6e-12$	$2.9e-10$	$3.2e-11$	$2.2e-11$		

818 based on the stability results for rank-1 updated eigenvalue solution in [24] and for  
819 the structured methods in [10].

820 We also demonstrate the importance of our local shifting strategy by testing the  
821 eigensolver with triangular FMM accelerations applied to the standard secular equa-  
822 tion. Due to cancellations, Matlab returns NaN (not-a-number) for the test matrices  
823 with sizes larger than 8192. This shows the risk of directly applying FMM accelera-  
824 tions to the standard secular equation like in [44].

825 EXAMPLE 3. Next, we consider a dense symmetric Toeplitz matrix  $A$  with its  
826 first row  $\xi = (\xi_1 \ \cdots \ \xi_n)$  given by

$$827 \quad \xi_1 = 2\alpha, \quad \xi_j = \frac{\sin(2\alpha(j-1)\pi)}{(j-1)\pi}, \quad j = 2, 3, \dots, n,$$

828 where  $0 < \alpha < 1/2$ . This is the so-called prolate matrix that appears frequently in  
829 signal processing. It is known to be extremely ill-conditioned and has special spectral  
830 properties (see, e.g., [43]). In fact, the prolate matrix has many small eigenvalues of  
831 magnitude  $O(\epsilon_{\text{mach}})$ . Here, we set  $\alpha = \frac{1}{4}$ . It is known that any Toeplitz matrix can be  
832 converted into a Cauchy-like matrix  $\mathcal{C}$  which has small off-diagonal numerical ranks  
833 [14, 33, 44]. That is,  $\mathcal{C} = \mathcal{F}A\mathcal{F}^*$ , where  $\mathcal{F}$  is the normalized inverse DFT matrix.  
834 Then the eigendecomposition of  $A$  can be done via that of  $\mathcal{C}$ . An HSS approximation  
835 to  $\mathcal{C}$  may be quickly constructed based on randomized methods in [30, 31, 51, 55] and  
836 fast Toeplitz matrix-vector multiplications. The cost is nearly linear in  $n$ . Here, we  
837 use a tolerance  $10^{-10}$  in relevant compression steps, which is the same as the deflation  
838 tolerance  $\tau$ . The size  $n$  ranges from 4096 to 65,536.

839 SuperDC and HSSBIS are applied to the resulting HSS form and compared with  
840 eig applied to  $A$ . In Figure 5.3, the timing, storage, and flops are shown and they  
841 are consistent with the complexity estimates. SuperDC shows a significant efficiency  
842 advantage over eig. At  $n = 32,768$ , SuperDC is about 122 times faster than eig. The  
843 accuracy is reported in Table 5.4.

844 One thing to point out is that the theoretical complexity  $O(r^2 n \log^2 n)$  of SuperDC  
845 may overestimate the actual cost. For example, here  $r$  is typically known to be  
846  $O(\log n)$  based on entrywise approximation errors [10, 33, 52, 55]. One reason for the  
847 overestimate is that the flop count does not take into consideration a levelwise rank

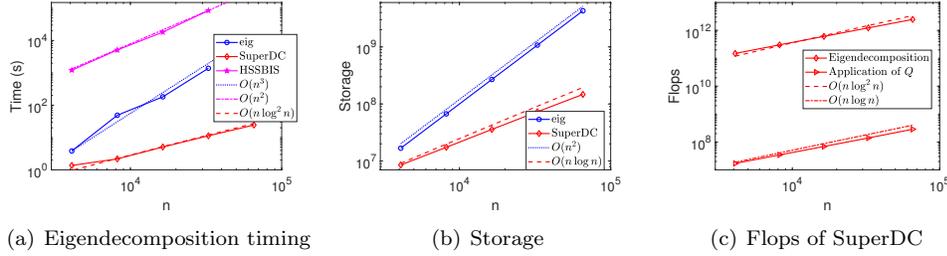
FIG. 5.3. Example 3. *Timing, storage, and flops.*

TABLE 5.4

Example 3. Accuracy of SuperDC, where the errors  $\delta_s$  and  $\delta_\infty$  for  $n = 65,536$  are not reported since `eig` runs out of memory.  $\delta_m$  is not available since many eigenvalues of order  $O(\epsilon_{mach})$  and `eig` returns a numerical eigenvalue 0 for some of the matrices.

$n$	4,096	8,192	16,384	32,768	65,536
$\gamma$	$2.3e-11$	$4.4e-11$	$1.8e-10$	$3.5e-10$	$1.4e-9$
$\theta$	$2.2e-15$	$8.6e-15$	$6.0e-15$	$4.2e-15$	$3.0e-15$
$\delta_s$	$5.5e-12$	$1.4e-11$	$4.5e-10$	$9.3e-10$	
$\delta_\infty$	$1.1e-10$	$7.3e-10$	$2.2e-8$	$6.1e-8$	

848 pattern in [50]. Another reason is our flexible deflation strategy in Section 4.1. The  
 849 matrices actually have highly clustered intermediate eigenvalues and many of them  
 850 get deflated. This further leads to high efficiency gain.

851 We have also tested SuperDC on random Toeplitz matrices, where the associated  
 852 Cauchy-like matrices  $\mathcal{C}$  have off-diagonal numerical ranks  $r$  quite larger than in the  
 853 prolate matrix case. Since the complexity of SuperDC is  $O(r^2 n \log^2 n)$ , it needs larger  
 854  $n$  to see an obvious advantage in timing over `eig`. (Of course, we may also use a larger  
 855 compression tolerance to get smaller  $r$ .) In addition, for random Toeplitz matrices,  
 856 deflation happens much less frequently than in the prolate matrix case.

857 EXAMPLE 4. The last example is a discretized kernel matrix  $A$  in [12] which  
 858 is the evaluation of the function  $\sqrt{|s-t|}$  at the Chebyshev points  $\cos(\frac{2i-1}{2n}\pi)$ ,  $i =$   
 859  $1, 2, \dots, n$ . The HSS construction may be based on direct off-diagonal compression  
 860 or efficient analytical methods like in [56]. We use an existing routine based on the  
 861 former one for simplicity. To show the flexibility of accuracy controls, we aim for  
 862 moderate accuracy in this test by using a compression tolerance  $10^{-6}$  in the HSS  
 863 construction, which is also the deflation tolerance.

864 For this example, we can observe similar complexity results as in the previous  
 865 examples. See Figure 5.4. With the larger tolerance than in the previous examples,  
 866 we still achieve reasonable eigenvalue errors and residuals with numerically orthogonal  
 867 eigenvectors. See Table 5.5.

868 We now show how our local shifting strategy (for triangular FMM-accelerated  
 869 solution of the shifted secular equations) can also significantly benefit the rate of  
 870 convergence of the roots. To illustrate this, we perform the following count. Suppose  $r$   
 871 secular equations are solved because of  $r$  rank-1 updates associated with the root node  
 872 of the HSS tree  $\mathcal{T}$ . When solving the  $j$ th secular equation, let  $\mu_j$  be the percentage of  
 873 eigenvalues that have *not* converged after 5 Newton's iterations. Let  $\mu = \max_{1 \leq j \leq r} \mu_j$ .

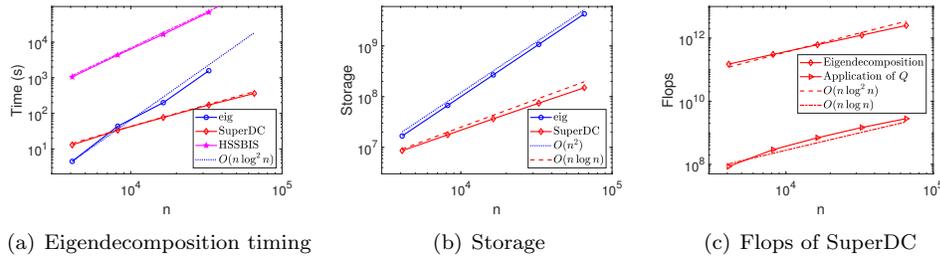
FIG. 5.4. Example 4. *Timing, storage, and flops.*

TABLE 5.5

Example 4. *Accuracy of SuperDC, where the errors for  $n = 65,536$  are not reported since eig runs out of memory.*

$n$	4,096	8,192	16,384	32,768	65,536
$\gamma$	$7.4e-9$	$2.7e-9$	$2.2e-9$	$1.6e-9$	$1.1e-9$
$\theta$	$1.4e-13$	$2.6e-13$	$2.5e-13$	$2.4e-13$	$3.8e-13$
$\delta_s$	$3.8e-8$	$5.5e-8$	$5.7e-8$	$8.7e-8$	
$\delta_\infty$	$2.9e-8$	$3.2e-8$	$2.8e-8$	$5.6e-8$	
$\delta_m$	$1.2e-4$	$4.2e-4$	$4.7e-4$	$5.8e-4$	

874 Table 5.6 reports this maximum percentage  $\mu$  with varying  $n$ . With local shifting, a  
 875 vast majority of those eigenvalues (about 99% or more) converges within 5 iterations.  
 876 This is significantly better than the case without local shifting (i.e., when the standard  
 877 secular equation is solved with FMM accelerations).

TABLE 5.6

*Maximum percentage ( $\mu$ ) of eigenvalues not converged within 5 iterations for solving the  $r$  secular equations associated with  $\text{root}(T)$ .*

$n$	4,096	8,192	16,384	32,768	65,536
With local shifting	1.00%	0.88%	0.34%	0.38%	0.33%
Without local shifting	62.5%	57.6%	57.2%	58.5%	57.7%

878 **6. Conclusions.** In this work, we have designed a SuperDC eigensolver that  
 879 significantly improves a previous development in terms of the stability and efficiency.  
 880 A series of stability enhancements is built into the different stages of the algorithm.  
 881 In particular, we avoid an exponential norm growth risk in the dividing stage via a  
 882 balancing strategy. We further combine FMM accelerations with several key stabil-  
 883 ity safeguards that have been used in practical divide-and-conquer algorithms. The  
 884 extensive numerical tests confirm the efficiency and the accuracy.

885 The SuperDC eigensolver makes it feasible to use full eigendecompositions to  
 886 solve various challenging numerical problems as mentioned at the beginning of the  
 887 paper. A list of applications is expected to be included in [35]. In addition, we expect  
 888 that the novel local shifting strategy and triangular FMM accelerations are also useful  
 889 for other FMM-related matrix computations when stability and accuracy are crucial.  
 890 In our future work, we plan to provide the proof of backward stability, as well as a  
 891 high-performance parallel implementation, which will extend the applicability of the

892 algorithm to large-scale numerical computations.

893 **Acknowledgements.** Thank the editor and the three anonymous reviewers for  
 894 providing valuable suggestions that have helped greatly improve the paper.

895

#### REFERENCES

- 896 [1] S. AMBIKASARAN AND E. DARVE, *An  $O(n \log n)$  fast direct solver for partial hierarchically*  
 897 *semi-separable matrices*, J. Sci. Comput., 57 (2013), pp. 477–501.
- 898 [2] P. AMESTOY, C. ASHCRAFT, O. BOITEAU, A. BUTTARI, J.-Y. L'EXCELLENT, AND C. WEIS-  
 899 BECKER, *Improving multifrontal methods by means of block low-rank representations*, SIAM  
 900 J. Sci. Comp., 37 (2015), pp. A1451–A1474.
- 901 [3] P. ARBENZ, *Divide and conquer algorithms for the bandsymmetric eigenvalue problem*, Parallel  
 902 Comput., 18(10), 1105–1128 (1992).
- 903 [4] A. CORTINOVIS, D. KRESSNER AND S. MASSEI, *Divide and conquer methods for functions of*  
 904 *matrices with banded or hierarchical low-rank structure*, arXiv:2107.04337 (2021).
- 905 [5] E. ANDERSON, Z. BAI, C. BISCHOF, S. BLACKFORD, J. DEMMEL, J. DONGARRA, J. DU CROZ,  
 906 A. GREENBAUM, S. HAMMARLING, A. MCKENNEY, AND D. SORENSEN, *LAPACK Users'*  
 907 *Guide*, SIAM, Philadelphia, PA, third ed., 1999.
- 908 [6] P. BENNER AND T. MACH, *Computing all or some eigenvalues of symmetric  $\mathcal{H}_1$ -matrices*, SIAM  
 909 J. Sci. Comput., 34 (2012), pp. A485–A496.
- 910 [7] D. A. BINI, L. GEMIGNANI, AND V. Y. PAN, *Fast and stable QR eigenvalue algorithms for*  
 911 *generalized companion matrices and secular equations*, Numer. Math., 100 (2005), pp.  
 912 373–408.
- 913 [8] D. BINI AND V. Y. PAN, *Parallel complexity of tridiagonal symmetric eigenvalue problem*, Proc.  
 914 Annu. ACM-SIAM Symp. Discrete Algorithms, SIAM, Philadelphia, 1991 pp. 384–393.
- 915 [9] J. R. BUNCH, C. P. NIELSEN, AND D. C. SORENSEN, *Rank-one modification of the symmetric*  
 916 *eigenproblem*, Numer. Math., 31 (1978), pp. 31–48.
- 917 [10] D. CAI AND J. XIA, *A stable matrix version of the fast multipole method: stabilization strategies*  
 918 *and examples*, Electron. Trans. Numer. Anal., 54 (2021), pp. 581–609.
- 919 [11] S. CHANDRASEKARAN, P. DEWILDE, M. GU, T. PALS, X. SUN, A. J. VAN DER VEEN, AND  
 920 D. WHITE, *Some fast algorithms for sequentially semiseparable representations*, SIAM J.  
 921 Matrix Anal. Appl., 27 (2005), pp. 341–364.
- 922 [12] S. CHANDRASEKARAN, P. DEWILDE, M. GU, W. LYONS, AND T. PALS, *A fast solver for HSS*  
 923 *representations via sparse matrices*, SIAM J. Matrix Anal. Appl., 29 (2006), pp. 67–81.
- 924 [13] S. CHANDRASEKARAN AND M. GU, *A divide-and-conquer algorithm for the eigendecomposition*  
 925 *of symmetric block diagonal plus semiseparable matrices*, Numer. Math., 96 (2004), pp.  
 926 723–731.
- 927 [14] S. CHANDRASEKARAN, M. GU, X. SUN, J. XIA, AND J. ZHU, *A superfast algorithm for Toeplitz*  
 928 *systems of linear equations*, SIAM J. Matrix Anal. Appl., 29 (2007), pp. 1247–1266.
- 929 [15] S. CHANDRASEKARAN, M. GU, J. XIA, AND J. ZHU, *A fast QR algorithm for companion ma-*  
 930 *trices*, Oper. Theory: Adv. Appl., Birkhaeuser Basel, 179 (2007), pp. 111–143.
- 931 [16] J. J. M. CUPPEN, *A divide and conquer method for the symmetric tridiagonal eigenproblem*,  
 932 Numer. Math., 36 (1981), pp. 177–195.
- 933 [17] J. W. DEMMEL, *Applied Numerical Linear Algebra*, SIAM, 1997.
- 934 [18] J. J. DONGARRA AND D. C. SORENSEN, *A fully parallel algorithm for the symmetric eigenvalue*  
 935 *problem*, SIAM J. Sci. Stat. Comput., 8 (1987), s139–s154.
- 936 [19] Y. EIDELMAN, I. GOHBERG, AND V. OLSHEVSKY, *The QR iteration method for Hermitian*  
 937 *quasiseparable matrices of an arbitrary order*, Linear Algebra Appl., 404 (2005), pp. 305–  
 938 324.
- 939 [20] Y. EIDELMAN AND I. HAIMOVICI, *Divide and conquer method for eigenstructure of quasisepa-*  
 940 *rable matrices using zeroes of rational matrix functions*, Oper. Theory: Adv. Appl., 218  
 941 (2012), Springer, Basel, pp. 299–328.
- 942 [21] Y. EIDELMAN, I. GOHBERG, AND I. HAIMOVICI, *Separable type representations of matrices and*  
 943 *fast algorithms*, Vol. 1, 2, Oper. Theory: Adv. Appl., 234, 235.
- 944 [22] G. H. GOLUB, *Some modified matrix eigenvalue problems*, SIAM Rev., 15 (1973), pp. 318–334.
- 945 [23] L. GREENGARD AND V. ROKHLIN, *A fast algorithm for particle simulations*, J. Comput. Phys.,  
 946 73 (1987), pp. 325–348.
- 947 [24] M. GU AND S. C. EISENSTAT, *A stable and efficient algorithm for the rank-one modification of*  
 948 *the symmetric eigenproblem*, SIAM J. Matrix Anal. Appl., 15 (1994), pp. 1266–1276.
- 949 [25] M. GU AND S. C. EISENSTAT, *A divide-and-conquer algorithm for the symmetric tridiagonal*

- 950 *eigenproblem*, SIAM J. Matrix Anal. Appl., 16 (1995), pp. 172–191.
- 951 [26] W. HACKBUSCH AND S. BORM, *Data-sparse approximation by adaptive  $\mathcal{H}^2$ -matrices*, Comput-  
952 ing, 69 (2002), pp.1–35.
- 953 [27] W. HACKBUSCH, *A sparse matrix arithmetic based on  $\mathcal{H}$ -matrices*, Computing, 62 (1999), pp.  
954 89–108.
- 955 [28] R. C. LI, *Solving secular equations stably and efficiently*, University of California, Berkeley,  
956 Technical Report No. UCB/CSD-94-851 (1994).
- 957 [29] X. LIAO, S. LI, L. CHENG, AND M. GU, *An improved divide-and-conquer algorithm for the*  
958 *banded matrices with narrow bandwidths*, Comput. Math. Appl., 71 (2016), pp. 1933–1943.
- 959 [30] X. LIU, J. XIA, AND M. V. DE HOOP, *Parallel randomized and matrix-free direct solvers for*  
960 *large structured dense linear systems*, SIAM J. Sci. Comput., 38 (2016), pp. S508–S538.
- 961 [31] P. G. MARTINSSON, *A fast randomized algorithm for computing a hierarchically semiseparable*  
962 *representation of a matrix*, SIAM J. Matrix Anal. Appl., 32 (2011), pp. 1251–1274.
- 963 [32] P. G. MARTINSSON AND V. ROKHLIN, *A fast direct solver for boundary integral equations in*  
964 *two dimensions*, J. Comput. Phys. 205 (2005), pp. 1–23.
- 965 [33] P. G. MARTINSSON, V. ROKHLIN, AND M. TYGERT, *A fast algorithm for the inversion of general*  
966 *Toeplitz matrices*, Comput. Math. Appl., 50 (2005), pp. 741–752.
- 967 [34] D. P. O’LEARY AND G. W. STEWART, *Computing the eigenvalues and eigenvectors of symmetric*  
968 *arrowhead matrices*, J. Comput. Phys., 90 (1990), pp. 497–505.
- 969 [35] X. OU, J. VOGEL, J. XIA, AND Z. XIN, *Efficient numerical computations via superfast eigen-*  
970 *value decompositions*, preprint, 2021.
- 971 [36] J. SHEN, Y. WANG, AND J. XIA, *Fast structured direct spectral methods for differential equations*  
972 *with variable coefficients, I. The one-dimensional case*, SIAM J. Sci. Comput., 38 (2016),  
973 pp. A28–A54.
- 974 [37] X. SUN AND N. P. PITSIANIS, *A matrix version of the fast multipole method*, SIAM Rev., 43  
975 (2001), pp. 289–300.
- 976 [38] A. ŠUŠNJARA AND D. KRESSNER, *A fast spectral divide-and-conquer method for banded matrices*,  
977 Numer. Linear Algebra Appl., 28 (2021), e2365.
- 978 [39] L. N. TREFETHEN AND D. BAU, *Numerical Linear Algebra*, SIAM, 1997.
- 979 [40] M. TYGERT, *Recurrence relations and fast algorithms*, Appl. Comput. Harmon. Anal., 28.1  
980 (2010), 121–128.
- 981 [41] M. VAN BAREL, R. VANDEBRIL, P. VAN DOOREN, AND K. FREDERIX, *Implicit double shift*  
982 *QR-algorithm for companion matrices*, Numer. Math., 116 (2010), pp. 177–212.
- 983 [42] R. VANDEBRIL, M. VAN BAREL, AND N. MASTRONARDI, *Matrix Computations and Semisepa-*  
984 *rable Matrices*, Vol. 1. Johns Hopkins University Press, Baltimore, MD, 2008.
- 985 [43] J. M. VARAH, *The prolate matrix*, Linear Algebra Appl., 187 (1993), pp. 269–278.
- 986 [44] J. VOGEL, J. XIA, S. CAULEY, AND V. BALAKRISHNAN, *Superfast divide-and-conquer method*  
987 *and perturbation analysis for structured eigenvalue solutions*, SIAM J. Sci. Comput., 38  
988 (2016), pp. A1358–A1382.
- 989 [45] Y. WANG, *Fast Structured Spectral Methods*, Ph.D. thesis, Purdue University, 2017.
- 990 [46] Y. XI AND J. XIA, *On the stability of some hierarchical rank structured matrix algorithms*,  
991 SIAM J. Matrix Anal. Appl., 37 (2016), pp. 1279–1303.
- 992 [47] Y. XI, J. XIA, S. CAULEY, AND V. BALAKRISHNAN, *Superfast and stable structured solvers for*  
993 *Toeplitz least squares via randomized sampling*, SIAM J. Matrix Anal. Appl., 35 (2014),  
994 pp. 44–72.
- 995 [48] Y. XI, J. XIA AND R. CHAN, *A fast randomized eigensolver with structured LDL factorization*  
996 *update*, SIAM J. Matrix Anal. Appl., 35 (2014), pp. 974–996.
- 997 [49] J. XIA, *Fast Direct Solvers for Structured Linear Systems of Equations*, Ph.D. thesis, University  
998 of California, Berkeley, 2006.
- 999 [50] J. XIA, *On the complexity of some hierarchical structured matrix algorithms*, SIAM J. Matrix  
1000 Anal. Appl., 33 (2012), pp. 388–410.
- 1001 [51] J. XIA, *Randomized sparse direct solvers*, SIAM J. Matrix Anal. Appl., 34 (2013), pp. 197–227.
- 1002 [52] J. XIA, *Multi-layer hierarchical structures*, CSIAM Trans. Appl. Math., 2 (2021), pp. 263–296.
- 1003 [53] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Superfast multifrontal method for large*  
1004 *structured linear systems of equations*, SIAM J. Matrix Anal. Appl., 31 (2009), pp. 1382–  
1005 1411.
- 1006 [54] J. XIA, S. CHANDRASEKARAN, M. GU, AND X. S. LI, *Fast algorithms for hierarchically semisepa-*  
1007 *rable matrices*, Numer. Linear Algebra Appl., 17 (2010), pp. 953–976.
- 1008 [55] J. XIA, Y. XI, AND M. GU, *A superfast structured solver for Toeplitz linear systems via ran-*  
1009 *domized sampling*, SIAM J. Matrix Anal. Appl., 33 (2012), pp. 837–858.
- 1010 [56] X. YE, J. XIA, AND L. YING, *Analytical low-rank compression via proxy point selection*, SIAM  
1011 J. Matrix Anal. Appl., 41 (2020), pp. 1059–1085.

## SUPPLEMENTARY MATERIALS: LIST OF MAJOR ALGORITHMS

Title of paper: *SuperDC: Superfast divide-and-conquer eigenvalue decomposition with improved stability for rank-structured matrices*

Authors: Xiaofeng Ou and Jianlin Xia

These supplementary materials are pseudocodes that can help better understand the major algorithms in the paper.

- Algorithm 1: the HSS dividing stage.
- Algorithm 2: solving the secular equation for the eigenvalues with triangular FMM accelerations and local shifting.
- Algorithm 3: the conquering stage for producing the eigendecomposition.
- Algorithm 4: application of the a local eigenmatrix  $Q_i$  or its transpose to a vector. This is used in Algorithm 3 and also can be used to apply the global eigenmatrix  $Q$  or its transpose to a vector when  $i = \text{root}(\mathcal{T})$ .

For notational convenience, we use  $r$  to represent the column sizes of all  $Z_i$  matrices in the pseudocodes.  $\mathcal{T}_i$  is also used to denote the subtree of  $\mathcal{T}$  rooted at node  $i \in \mathcal{T}$ .  $Z(:, j)$  means the  $j$ -th column of  $Z$ .

The following utility routines are used in the algorithms. To save space, we are not showing pseudocodes for these routines.

- `updhss`( $D_i, U_i, H$ ): for an HSS block  $D_i$  corresponding to the subtree  $\mathcal{T}_i$ , update its  $D, B$  generators to get those of  $D_i - U_i H U_i^T$  using Lemma 2.1.
- `trifmm`( $\mathbf{d}, \mathbf{x}, \mathbf{y}, \mathbf{w}, \kappa$ ): compute a matrix-vector product  $K\mathbf{w}$  with the triangular FMM and local shifting as in Sections 4.2.2 and 4.3.2, where  $K = (\kappa(d_i, x_j))_{d_i \in \mathbf{d}, x_j \in \mathbf{x}}$  is a kernel matrix and  $\mathbf{y}$  is the gap vector (for accurately evaluating  $\mathbf{x} - \mathbf{d}$ ). Note that the triangular FMM is used to multiply the lower triangular part of  $K$  with  $\mathbf{w}$  and the strictly upper triangular part of  $K$  with  $\mathbf{w}$  and the final result is the sum of the two products.
- `mnewton`( $\psi, \phi, \psi', \phi'$ ): use the modified Newton's method to compute corrections to the current approximate gap as in (4.19), where  $\psi, \phi, \psi', \phi'$  look like (4.9) and (4.10).
- `iniguess`( $\mathbf{d}, \mathbf{w}$ ): compute the initial guess as in [28] for the solution of the secular equation (2.12).
- `deflate`( $\mathbf{d}, \mathbf{v}, \tau$ ): apply deflation with the criterion in Section 4.1.

**Algorithm 1** SuperDC dividing stage

---

```

1: procedure divide( $\{D_i\}_{i \in \mathcal{T}}, \{U_i\}_{i \in \mathcal{T}}, \{R_i\}_{i \in \mathcal{T}}, \{B_i\}_{i \in \mathcal{T}}$ )
2:   for node  $i = \text{root}(\mathcal{T}), \dots, 1$  do  $\triangleright$  Dividing  $D_i$  in a top-down traversal
3:     if  $i$  is a non-leaf node then
4:       if  $\text{colsize}(B_{c_1}) \leq \text{rowsize}(B_{c_1})$  then  $\triangleright c_1, c_2$ : children of  $i$ 
5:          $D_{c_1} \leftarrow \text{updhss}(D_{c_1}, U_{c_1}, \frac{1}{\|B_{c_1}\|_2} B_{c_1} B_{c_1}^T)$   $\triangleright$  Update generators of  $D_{c_1}$ 
6:           to get those of  $D_{c_1} - \frac{1}{\|B_{c_1}\|_2} U_{c_1} B_{c_1} B_{c_1}^T U_{c_1}^T$  like in Lemma 2.1
7:          $D_{c_2} \leftarrow \text{updhss}(D_{c_2}, U_{c_2}, \|B_{c_1}\|_2 I)$   $\triangleright$  Update generators of  $D_{c_2}$ 
8:           to get those of  $D_{c_2} - \|B_{c_1}\|_2 U_{c_2} U_{c_2}^T$  like in Lemma 2.1
9:       else
10:         $D_{c_1} \leftarrow \text{updhss}(D_{c_1}, U_{c_1}, \|B_{c_1}\|_2 I)$   $\triangleright$  Update generators of  $D_{c_1}$ 
11:          to get those of  $D_{c_1} - \|B_{c_1}\|_2 U_{c_1} U_{c_1}^T$  like in Lemma 2.1
12:         $D_{c_2} \leftarrow \text{updhss}(D_{c_2}, U_{c_2}, \frac{1}{\|B_{c_1}\|_2} B_{c_1}^T B_{c_1})$   $\triangleright$  Update generators of  $D_{c_2}$ 
13:          to get those of  $D_{c_2} - \frac{1}{\|B_{c_1}\|_2} U_{c_2} B_{c_1}^T B_{c_1} U_{c_2}^T$  like in Lemma 2.1
14:       end if
15:     end if
16:   end for
17:   for node  $i = 1, \dots, \text{root}(\mathcal{T})$  do  $\triangleright$  Form  $Z_i$  in a bottom-up traversal
18:     if  $i$  is a non-leaf node then
19:       if  $\text{colsize}(B_{c_1}) \leq \text{rowsize}(B_{c_1})$  then  $\triangleright c_1, c_2$ : children of  $i$ 
20:          $Z_i \leftarrow \begin{pmatrix} \frac{1}{\sqrt{\|B_{c_1}\|_2}} U_{c_1} B_{c_1} \\ \sqrt{\|B_{c_1}\|_2} U_{c_2} \end{pmatrix}$   $\triangleright$  Local update  $Z$  matrix like in (3.12)
21:       else
22:          $Z_i \leftarrow \begin{pmatrix} \sqrt{\|B_{c_1}\|_2} U_{c_2} \\ \frac{1}{\sqrt{\|B_{c_1}\|_2}} U_{c_2} B_{c_1}^T \end{pmatrix}$   $\triangleright$  Local update  $Z$  matrix like in (3.15)
23:       end if
24:     if  $i \neq \text{root}(\mathcal{T})$  then
25:        $U_i \leftarrow \begin{pmatrix} U_{c_1} R_{c_1} \\ U_{c_2} R_{c_2} \end{pmatrix}$   $\triangleright$  Assemble  $U_i$  for parent node of  $i$ 
26:     end if
27:   end for
28:   return updated generators  $\{D_i\}_{i \in \mathcal{T}}, \{B_i\}_{i \in \mathcal{T}}, \{Z_i\}_{i \in \mathcal{T}}$ 
29: end procedure

```

---

---

**Algorithm 2** Secular equation solution for eigenvalues (of  $\text{diag}(\mathbf{d}) + \mathbf{v}\mathbf{v}^T$ )

---

```

1: procedure secular( $\mathbf{d}, \mathbf{v}$ )
     $\triangleright$  Eigenvalue solution via the solution of the shifted secular equation (4.15)
2:    $\mathbf{w} \leftarrow \mathbf{v} \odot \mathbf{v}$ 
3:    $\mathbf{x}^{(0)} \leftarrow \text{iniguess}(\mathbf{d}, \mathbf{w})$   $\triangleright$  Computation of the initial guess as in [28]
4:    $\mathbf{y}^{(0)} \leftarrow \mathbf{x}^{(0)} - \mathbf{d}$ 
5:   for  $j = 0, 1, \dots$  do
6:      $[\psi, \phi] \leftarrow \text{trifmm}(\mathbf{d}, \mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{w}, \frac{1}{s-t})$   $\triangleright$  Computation of  $\psi, \phi$  in (4.9)
7:      $[\psi', \phi'] \leftarrow \text{trifmm}(\mathbf{d}, \mathbf{x}^{(j)}, \mathbf{y}^{(j)}, \mathbf{w}, \frac{1}{(s-t)^2})$   $\triangleright$  Computation of  $\psi', \phi'$  in (4.10)
8:      $\mathbf{f} \leftarrow \mathbf{e} + \psi + \phi$ 
9:     if  $|\mathbf{f}| < cn(\mathbf{e} + |\psi| + |\phi|)\epsilon$  then  $\triangleright$  Stopping criterion
10:      break
11:    end if
12:     $\Delta \mathbf{x}^{(j)} \leftarrow \text{mnewton}(\psi, \phi, \psi', \phi')$ 
     $\triangleright$  Computation of root update with modified Newton's method
13:     $\mathbf{y}^{(j+1)} \leftarrow \mathbf{y}^{(j)} + \Delta \mathbf{x}^{(j)}$   $\triangleright$  Updated gap approximation as in (4.19)
14:     $\mathbf{x}^{(j+1)} \leftarrow \mathbf{y}^{(j+1)} + \mathbf{d}$   $\triangleright$  Updated eigenvalue approximation
15:  end for
16:   $\lambda \leftarrow \mathbf{x}^{(j)}, \eta \leftarrow \mathbf{y}^{(j)}$   $\triangleright$  Eigenvalue and gap upon convergence
17:  return  $\lambda, \eta$ 
18: end procedure

```

---

**Algorithm 3** SuperDC conquering stage

---

```

1: procedure conquer( $\{D_i\}_{i \in \mathcal{T}}, \{U_i\}_{i \in \mathcal{T}}, \{R_i\}_{i \in \mathcal{T}}, \{B_i\}_{i \in \mathcal{T}}, \{Z_i\}_{i \in \mathcal{T}}, \tau$ )
     $\triangleright$  The  $D_i, B_i$  generators have been updated in the dividing stage
2:   for node  $i = 1, \dots, \text{root}(\mathcal{T})$  do            $\triangleright$  Conquering in a postordered traversal
3:     if  $i$  is a leaf node then                    $\triangleright$  Leaf-level eigendecomposition
4:        $(\lambda_i, \hat{Q}_i) \leftarrow \text{eig}(D_i)$             $\triangleright$  Via Matlab eig function
5:     else
6:        $\begin{pmatrix} Z_{i,1} \\ Z_{i,2} \end{pmatrix} \leftarrow Z_i$         $\triangleright$  Partitioning following the sizes of  $D_{c_1}$  and  $D_{c_2}$ 
7:        $Z_{i,1} \leftarrow \text{superdcmv}(Q_{c_1}, Z_{i,1}, 1)$     $\triangleright Q_{c_1}^T Z_{i,1}$ 
8:        $Z_{i,2} \leftarrow \text{superdcmv}(Q_{c_2}, Z_{i,2}, 1)$     $\triangleright Q_{c_2}^T Z_{i,2}$ 
9:        $Z_i \leftarrow \begin{pmatrix} Z_{i,1} \\ Z_{i,2} \end{pmatrix}$         $\triangleright \hat{Z}_i$  like in (2.8)
10:       $[\lambda_i^{(0)}, P_i] \leftarrow \text{sort}(\lambda_{c_1}, \lambda_{c_2})$     $\triangleright$  Ordering of all the diagonal entries
        of  $\lambda_{c_1}, \lambda_{c_2}$  together, with  $P_i$  the permutation matrix
11:      for  $j = 1, 2, \dots, r$  do                    $\triangleright r = \text{colsize}(Z_i)$ 
12:         $[\mathbf{d}_i^{(j)}, Z_i(:, j)] \leftarrow \text{deflate}(\lambda_i^{(j-1)}, Z_i(:, j), \tau)$   $\triangleright$  Deflation (Section 4.1)
13:         $[\lambda_i^{(j)}, \boldsymbol{\eta}_i^{(j)}] \leftarrow \text{secular}(\mathbf{d}_i^{(j)}, Z_i(:, j))$     $\triangleright$  Secular equation solution
14:         $\mathbf{v}_1 \leftarrow \text{trifmm}(\mathbf{d}_i^{(j)}, \lambda_i^{(j)}, \boldsymbol{\eta}_i^{(j)}, \mathbf{e}, \log |s - t|)$   $\triangleright G_1 \mathbf{e}$  as needed in (4.27)
15:         $\mathbf{v}_2 \leftarrow \text{trifmm}(\mathbf{d}_i^{(j)}, \mathbf{d}_i^{(j)}, \mathbf{0}, \mathbf{e}, \log |s - t|)$     $\triangleright G_2 \mathbf{e}$  as needed in (4.27)
16:         $\hat{\mathbf{v}}_i^{(j)} \leftarrow \exp(\frac{\mathbf{v}_1 - \mathbf{v}_2}{2})$     $\triangleright$  Löwner's formula for  $\hat{\mathbf{v}}$  as in (4.25)–(4.27)
17:         $\mathbf{b}_i^{(j)} \leftarrow (\text{trifmm}(\mathbf{d}_i^{(j)}, \lambda_i^{(j)}, \boldsymbol{\eta}_i^{(j)}, \hat{\mathbf{v}}_i^{(j)} \odot \hat{\mathbf{v}}_i^{(j)}, \frac{1}{(s-t)^2}))^{-1/2}$ 
         $\triangleright$  Normalization factor as in (4.29)
18:         $\hat{Q}_i^{(j)} \leftarrow \{\hat{\mathbf{v}}_i^{(j)}, \mathbf{b}_i^{(j)}, \mathbf{d}_i^{(j)}, \lambda_i^{(j)}, \boldsymbol{\eta}_i^{(j)}\}$     $\triangleright$  Cauchy-like structured
        representation of the local eigenmatrix as in (4.28)
19:        for  $k = j + 1, j + 2, \dots, r$  do        $\triangleright$  Multiplication of  $\hat{Q}_i^{(j)}$ 
        to the remaining columns of  $Z_i$  via the steps as in (4.31)
20:           $Z_i(:, k) \leftarrow \hat{\mathbf{v}}_i^{(j)} \odot Z_i(:, k)$ 
21:           $Z_i(:, k) \leftarrow \text{trifmm}(\mathbf{d}_i^{(j)}, \lambda_i^{(j)}, \boldsymbol{\eta}_i^{(j)}, Z_i(:, k), \frac{1}{s-t})$ 
22:           $Z_i(:, k) \leftarrow \mathbf{b}_i^{(j)} \odot Z_i(:, k)$ 
23:        end for
24:      end for
25:       $\lambda_i \leftarrow \lambda_i^{(r)}$             $\triangleright$  Local eigenvalues associated with node  $i$ 
26:    end if
27:  end for
28:   $\lambda \leftarrow \lambda_{\text{root}(\mathcal{T})}, Q \leftarrow \{\{\hat{Q}_i^{(j)}\}_{j=1}^r, P_i\}_{i \in \mathcal{T}}$     $\triangleright$  Final eigenvalues
    and eigenmatrix  $Q$  in (4.36), with  $\hat{Q}_i$  in (4.35) given by  $\prod_{j=1}^r \hat{Q}_i^{(j)}$ 
29:  return  $\lambda, Q$ 
30: end procedure

```

---

---

**Algorithm 4** SuperDC eigenmatrix-vector multiplication
 

---

```

1: procedure superdcmv( $Q_i, \mathbf{x}, \text{transpose}$ )  $\triangleright$  Application of a local eigenmatrix  $Q_i$ 
   or its transpose to a vector  $\mathbf{x}$ , depending on whether ‘transpose’ is 0 or 1
2:    $i_1 \leftarrow$  smallest descendant of  $i$ 
3:   if transpose = 0 then  $\triangleright \mathbf{y} = Q_i \mathbf{x}$ 
4:      $\mathbf{y}_i \leftarrow \mathbf{x}$ 
5:     for  $k = i, i - 1, \dots, i_1$  do  $\triangleright$  Reverse postordered traversal of  $\mathcal{T}_i$ 
6:       if  $k$  is leaf then
7:          $\mathbf{y}_k \leftarrow Q_k \mathbf{y}_k$   $\triangleright$  Dense  $Q_k$  at the leaf level
8:       else
9:         for  $j = r, r - 1, \dots, 1$  do  $\triangleright$  Multiplication of  $\hat{Q}_k^{(j)}$ 
 $\triangleright$  via the steps like in (4.31)
10:            $\mathbf{y}_k \leftarrow \mathbf{b}_k^{(j)} \odot \mathbf{y}_k$ 
11:            $\mathbf{y}_k \leftarrow \text{trifmm}(\mathbf{d}_k^{(j)}, \boldsymbol{\lambda}_k^{(j)}, \boldsymbol{\eta}_k^{(j)}, \mathbf{y}_k, \frac{1}{s-t})$ 
12:            $\mathbf{y}_k \leftarrow \hat{\mathbf{v}}_k^{(j)} \odot \mathbf{y}_k$ 
13:         end for
14:          $\mathbf{y}_k \leftarrow P_k^T \mathbf{y}_k$   $\triangleright$  Permutation like in (4.34)
15:          $\begin{pmatrix} \mathbf{y}_{c_1} \\ \mathbf{y}_{c_2} \end{pmatrix} \leftarrow \mathbf{y}_k$   $\triangleright$  Partitioning following the sizes of  $Q_{c_1}, Q_{c_2}$ ,
 $\triangleright$  with  $c_1, c_2$  the children of  $k$ 
16:       end if
17:     end for
18:   else  $\triangleright \mathbf{y} = Q_i^T \mathbf{x}$ 
19:     Partition  $\mathbf{x}$  into  $\mathbf{x}_k$  pieces following the leaf-level  $Q_k$  sizes
20:     for  $k = i_1, i_1 + 1, \dots, i$  do  $\triangleright$  Postordered traversal of  $\mathcal{T}_i$ 
21:       if  $k$  is leaf then
22:          $\mathbf{y}_k \leftarrow Q_k^T \mathbf{x}_k$   $\triangleright$  Dense  $Q_k$  at the leaf level
23:       else
24:          $\mathbf{y}_k \leftarrow \begin{pmatrix} \mathbf{y}_{c_1} \\ \mathbf{y}_{c_2} \end{pmatrix}$   $\triangleright c_1, c_2$ : children of  $k$ 
25:          $\mathbf{y}_k \leftarrow P_k \mathbf{y}_k$   $\triangleright$  Permutation like in (4.34)
26:         for  $j = 1, 2, \dots, r$  do  $\triangleright$  Multiplication of  $(\hat{Q}_k^{(j)})^T$ 
 $\triangleright$  via the steps like in (4.31)
27:            $\mathbf{y}_k \leftarrow \hat{\mathbf{v}}_k^{(j)} \odot \mathbf{y}_k$ 
28:            $\mathbf{y}_k \leftarrow -\text{trifmm}(\boldsymbol{\lambda}_k^{(j)}, \mathbf{d}_k^{(j)}, \boldsymbol{\eta}_k^{(j)}, \mathbf{y}_k, \frac{1}{s-t})$   $\triangleright$  The negative sign
 $\triangleright$  and the switch of  $\boldsymbol{\lambda}_k^{(j)}$  and  $\mathbf{d}_k^{(j)}$  are because of the transpose
29:            $\mathbf{y}_k \leftarrow \mathbf{b}_k^{(j)} \odot \mathbf{y}_k$ 
30:         end for
31:       end if
32:     end for
33:   end if
34:   return  $\mathbf{y}$ 
35: end procedure

```

---