

PSe1Inv—A Distributed Memory Parallel Algorithm for Selected Inversion: The Symmetric Case

MATHIAS JACQUELIN, Lawrence Berkeley National Laboratory

LIN LIN, University of California, Berkeley and Lawrence Berkeley National Laboratory

CHAO YANG, Lawrence Berkeley National Laboratory

We describe an efficient parallel implementation of the selected inversion algorithm for distributed memory computer systems, which we call PSe1Inv. The PSe1Inv method computes selected elements of a general sparse matrix A that can be decomposed as $A = LU$, where L is lower triangular and U is upper triangular. The implementation described in this article focuses on the case of sparse symmetric matrices. It contains an interface that is compatible with the distributed memory parallel sparse direct factorization SuperLU_DIST. However, the underlying data structure and design of PSe1Inv allows it to be easily combined with other factorization routines, such as PARDISO. We discuss general parallelization strategies such as data and task distribution schemes. In particular, we describe how to exploit the concurrency exposed by the elimination tree associated with the LU factorization of A . We demonstrate the efficiency and accuracy of PSe1Inv by presenting several numerical experiments. In particular, we show that PSe1Inv can run efficiently on more than 4,000 cores for a modestly sized matrix. We also demonstrate how PSe1Inv can be used to accelerate large-scale electronic structure calculations.

Categories and Subject Descriptors: G.4 [Mathematical Software]: Algorithm design and analysis; G.4 [Mathematical Software]: Parallel and vector implementations; I.1.2 [Symbolic and Algebraic Manipulation]: Algorithms—Algebraic algorithms

General Terms: Design, Performance

Additional Key Words and Phrases: Selected inversion, sparse direct method, distributed memory parallel algorithm, high-performance computation, electronic structure theory

ACM Reference Format:

Mathias Jacquelin, Lin Lin, and Chao Yang. 2016. PSe1Inv—A distributed memory parallel algorithm for selected inversion: The symmetric case. ACM Trans. Math. Softw. 43, 3, Article 21 (November 2016), 28 pages.

DOI: <http://dx.doi.org/10.1145/2786977>

1. INTRODUCTION

Let $A \in \mathbb{C}^{N \times N}$ be a nonsingular sparse matrix. We are interested in computing *selected elements* of A^{-1} , defined as

$$\{(A^{-1})_{i,j} \mid \text{for } 1 \leq i, j \leq N \text{ such that } A_{i,j} \neq 0\}. \quad (1)$$

This work was partially supported by the Laboratory Directed Research and Development Program of Lawrence Berkeley National Laboratory under U.S. Department of Energy contract DE-AC02-05CH11231 (L. L. and C. Y.), the Scientific Discovery through Advanced Computing (SciDAC) program funded by the U.S. Department of Energy, Office of Science, Advanced Scientific Computing Research and Basic Energy Sciences (M. J., L. L., and C. Y.), and the Center for Applied Mathematics for Energy Research Applications (CAMERA), which is a partnership between Basic Energy Sciences (BES) and Advanced Scientific Computing Research (ASRC) at the U.S. Department of Energy.

Q1 Authors' addresses:

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2016 ACM 0098-3500/2016/11-ART21 \$15.00

DOI: <http://dx.doi.org/10.1145/2786977>

32 Sometimes we only need to compute a subset of these selected elements—for example,
 33 the diagonal elements of A^{-1} . The most straightforward way to obtain these selected
 34 elements of A^{-1} is to compute the full inverse of A and then extract the selected
 35 elements. But this is often prohibitively expensive in practice. It turns out that to
 36 compute these selected elements of A^{-1} , some additional elements of A^{-1} often need to
 37 be computed. However, the overall set of nonzero elements that need to be computed
 38 often remains a small percentage of all elements of A^{-1} due to the sparsity structure
 39 of A .

40 The selected elements of A^{-1} defined by (1) can be used to obtain trace estimation of
 41 the form

$$\text{Tr}[A^{-1}] \quad \text{or} \quad \text{Tr}[A^{-1}B^T], \quad (2)$$

42 if the sparsity pattern of $B \in \mathbb{C}^{N \times N}$ is contained in the sparsity pattern of A (i.e.,
 43 $\{(i, j) | B_{i,j} \neq 0\} \subset \{(i, j) | A_{i,j} \neq 0\}$). The computation of selected elements of A^{-1} together
 44 with the trace estimation of the form (2) arise in several scientific computing appli-
 45 cations, including density functional theory (DFT) [Hohenberg and Kohn 1964; Kohn
 46 and Sham 1965], dynamical mean-field theory (DMFT) [Kotliar et al. 2006], Poisson-
 47 Boltzmann equations [Xu and Maggs 2013], and uncertainty quantification [Bekas
 48 et al. 2009].

49 It is possible to compute selected elements of A^{-1} by iterative methods such as the
 50 Lanczos algorithm [Lanczos 1950; Sidje and Saad 2011], combined with Monte Carlo
 51 [Bekas et al. 2007] or deterministic probing techniques [Tang and Saad 2012]. These
 52 types of methods work well if A^{-1} is a banded matrix or it becomes a banded matrix
 53 after elements with absolute value less than ε have been truncated, and a sparse
 54 factorization of A is prohibitively expensive to perform.

55 In this article, we focus on using a sparse direct method to compute selected elements
 56 **Q2** of A^{-1} . We assume that a sparse LU factorization (or LDL^T factorization if A is sym-
 57 metric) of A is computationally feasible. The main advantage of a direct method is that
 58 we do not need to make assumptions on the decay property of A^{-1} . The disadvantage
 59 of this direct method is that it actually computes a superset of the selected elements of
 60 A^{-1} as defined in Equation (1). In particular, all elements of A^{-1} indexed by the union
 61 of the nonzero index sets of the L and U factors need to be computed. It should be noted
 62 that as long as L and U remain sparse, computing elements of A^{-1} restricted to this
 63 superset can still be much faster than computing the full inverse.

64 Sparse direct methods for computing selected elements of A^{-1} were first proposed
 65 in Takahashi et al. [1973] and Erisman and Tinney [1975]. The use of an elimination
 66 tree for computing selected elements of inverse was presented in Campbell and Davis
 67 [1995] for a sequential algorithm. Motivated by quantum transport simulations, Li
 68 et al. [2008], Li and Darve [2012], and Li et al. [2013] developed the fast inverse
 69 using nested dissection (FIND) algorithm for computing the diagonal of A^{-1} . Some
 70 related work has been described in Cauley et al. [2012] and Eastwood and Wan [2013].
 71 The FIND algorithm is in principle applicable to matrices with a general sparsity
 72 pattern, but so far its implementation focuses on structured matrices obtained from
 73 second-order partial differential operators discretized by a finite difference scheme.
 74 The implementation of FIND is not yet publicly available. Nor is it scalable to a large
 75 number of processors. Lin et al. [2009b] developed the hierarchically Schur complement
 76 (HSC) method for a similar type of matrix arising from DFT-based electronic structure
 77 calculations. The method has also been generalized and applied to quantum transport
 78 calculations [Hetmaniuk et al. 2013].

79 An efficient implementation of the selected inversion algorithm for a general sym-
 80 metric matrix, called SelInv, was presented in Lin et al. [2011b] and is publicly avail-
 81 able. Amestoy et al. [2012a] considered a more general parallel matrix inversion method

for computing any subset of entries of A^{-1} . They implemented their algorithm in the MUMPS package [Amestoy et al. 2001], which is based on the multifrontal method. In this algorithm, the actual set of computed entries of A^{-1} contains entries on the critical path of the requested entries to the root of the elimination tree, and therefore this is also a superset of the requested entries of A^{-1} . This method is more efficient than our algorithm when a small number of entries of A^{-1} are requested. However, when a relatively large number of entries are to be computed, such as in the computation of the selected elements defined by (1), our algorithm can more efficiently reuse the information shared among different entries of A^{-1} , and our numerical results indicate that our algorithm is more efficient than the parallel matrix inversion method implemented in MUMPS. In addition to the work presented in Amestoy et al. [2012a, 2012b], parallel implementation of algorithms for computing selected elements of inverse tailored to matrices obtained from a finite difference discretization of a second-order partial differential operator have appeared in several publications [Petersen et al. 2009; Lin et al. 2011a]. The method developed by Petersen et al. [2009], which was designed for a quasi-1D quantum transport problem, is scalable to a relatively small ($32 \sim 64$) number of processors. In Lin et al. [2011a], we described an efficient parallel implementation for discretized 2D Laplacian-type operators and demonstrated the efficiency of the implementation when it was used to solve a problem with billions of degrees of freedom on 4,096 processors. However, this implementation cannot be used to perform a selected inversion of a general symmetric sparse matrix with an arbitrary sparsity pattern.

The purpose of this article is to extend the selected inversion algorithm presented in Lin et al. [2011b] and describe an implementation of a parallel selected inversion algorithm designed for distributed memory parallel computers. Such an implementation allows us to solve much larger problems by utilizing more computational resources available on high-performance computers. The present work is more general than the previous work in Lin et al. [2011a], which assumes a balanced binary elimination tree and is only applicable to structured sparse matrices obtained from finite difference discretization of differential operators. It can utilize far greater numbers of processors than the sequential algorithm described in Lin et al. [2011b] for general sparse matrices.

The parallel implementation of the selected inversion algorithm that we present in this article uses a more general data distribution and communication strategy to divide the work among a large number of processors to achieve multiple levels of concurrency. We name our implementation PSe1Inv, and the software is publicly available.¹ Our first implementation focuses on the case of sparse symmetric matrices. In principle, the PSe1Inv package can be interfaced with any sparse LU and LDL^T factorization routines. Our current implementation provides an interface to the SuperLU_DIST [Li and Demmel 2003] package. The user has the option of using either the ParMETIS [Karypis and Kumar 1998] software or the PT-Scotch [Chevalier and Pellegrini 2008] package to reorder the matrix in parallel to minimize the nonzero fill in the sparse matrix factors.

The rest of the article is organized as follows. We review the basic idea of the selected inversion method in Section 2 and discuss various implementation issues for the distributed memory parallel selected inversion algorithm in Section 3. The numerical results with applications to various matrices from the Harwell-Boeing Test Collection [Duff et al. 1992] and the University of Florida Matrix Collection [Davis and Hu 2011],

¹<http://www.pexsi.org/>, distributed under the BSD license.

130 and also applications from DFT, are given in Section 4, followed by our conclusions and
 131 the discussion of future work in Section 5.

132 Standard linear algebra notation is used for vectors and matrices throughout the
 133 article. We use $A_{i,j}$ to denote the (i, j) -th entry of the matrix A and f_i to denote the i -th
 134 entry of the vector f . With slight abuse of notation, both a supernode index and the set
 135 of column indices associated with a supernode are denoted by uppercase script letters,
 136 such as $\mathcal{I}, \mathcal{J}, \mathcal{K}$. Furthermore, we use $A_{i,*}$ and $A_{*,j}$ to denote the i -th row and the j -th
 137 column of A , respectively. Similarly, $A_{\mathcal{I},*}$ and $A_{*,\mathcal{J}}$ are used to denote the \mathcal{I} -th block row
 138 and the \mathcal{J} -th block column of A , respectively. $A_{\mathcal{I},\mathcal{J}}^{-1}$ denotes the $(\mathcal{I}, \mathcal{J})$ -th block of the
 139 matrix A^{-1} (i.e., $A_{\mathcal{I},\mathcal{J}}^{-1} \equiv (A^{-1})_{\mathcal{I},\mathcal{J}}$). When the block $A_{\mathcal{I},\mathcal{J}}$ itself is invertible, its inverse is
 140 denoted by $(A_{\mathcal{I},\mathcal{J}})^{-1}$ to distinguish it from $A_{\mathcal{I},\mathcal{J}}^{-1}$.

141 2. SELECTED INVERSION ALGORITHM

142 2.1. Basic Formulation

143 Although this article focuses on the computation of selected elements of A^{-1} when A is a
 144 sparse symmetric matrix, the idea of the selected inversion algorithm can be given for a
 145 general square matrix A and will be used in our ongoing work for computing the selected
 146 elements of A^{-1} for asymmetric matrices. The standard approach for computing A^{-1} is
 147 to first decompose the general matrix A using the LU factorization

$$A = LU, \quad (3)$$

148 where L is a unit lower triangular matrix and U is an upper triangular matrix. To
 149 stabilize the computation, matrix reordering and partial pivoting [Golub and Van Loan
 150 1996] are usually applied to the matrix of A , and the general form of the LU factorization
 151 can be given as

$$PAQ = LU, \quad (4)$$

152 where P and Q are two permutation matrices. To simplify the discussion, we will use
 153 Equation (3) and assume that A has already been permuted.

154 Given the LU factorization, the most straightforward way to compute selected
 155 elements of A^{-1} is to obtain $A^{-1} \equiv (x_1, x_2, \dots, x_n)$ by solving several triangular
 156 systems

$$Ly_j = e_j, \quad Ux_j = y_j, \quad (5)$$

157 where $j = 1, 2, \dots, n$, and e_j is the j -th column of the $n \times n$ identity matrix. Such
 158 a procedure, which will be referred to as the direct inversion algorithm, is generally
 159 very costly even when A is sparse. The direct inversion algorithm performs too much
 160 computation when only a small number of selected elements of the inverse matrix are
 161 needed.

162 An alternative algorithm is the selected inversion algorithm, which accurately com-
 163 putes all selected elements of A^{-1} . The idea of the selected inversion method originates
 164 from Takahashi et al. [1973] and Erisman and Tinney [1975], and the algorithm and
 165 its variants have been discussed in several works [Lin et al. 2009b; Lin et al. 2011b;
 166 Amestoy et al. 2012b; Li et al. 2008; Li and Darve 2012; Kuzmin et al. 2013]. The se-
 167 lected inversion algorithm can be understood as follows. We first partition the matrix
 168 A into 2×2 blocks of the form

$$A = \begin{pmatrix} A_{1,1} & A_{1,2} \\ A_{2,1} & A_{2,2} \end{pmatrix}, \quad (6)$$

PSe1Inv: The Symmetric Case

21:5

where $A_{1,1}$ is a scalar of size 1×1 . We can write $A_{1,1}$ as a product of two scalars $L_{1,1}$ and $U_{1,1}$. In particular, we can pick $L_{1,1} = 1$ and $U_{1,1} = A_{1,1}$. Then

$$A = \begin{pmatrix} L_{1,1} & 0 \\ L_{2,1} & I \end{pmatrix} \begin{pmatrix} U_{1,1} & U_{1,2} \\ 0 & S_{2,2} \end{pmatrix}, \quad (7)$$

where

$$L_{2,1} = A_{2,1}(U_{1,1})^{-1}, \quad U_{1,2} = (L_{1,1})^{-1}A_{1,2}. \quad (8)$$

The L and U factors are usually directly accessible in a standard LU factorization, and

$$S_{2,2} = A_{2,2} - L_{2,1}U_{1,2} \quad (9)$$

is the Schur complement. Using the decomposition given by Equation (7), we can express A^{-1} as

$$A^{-1} = \begin{pmatrix} (U_{1,1})^{-1}(L_{1,1})^{-1} + (U_{1,1})^{-1}U_{1,2}S_{2,2}^{-1}L_{2,1}(L_{1,1})^{-1} & -(U_{1,1})^{-1}U_{1,2}S_{2,2}^{-1} \\ -S_{2,2}^{-1}L_{2,1}(L_{1,1})^{-1} & S_{2,2}^{-1} \end{pmatrix}. \quad (10)$$

Since $S_{2,2}$ is the same as S here, without ambiguity, $S_{2,2}^{-1} \equiv (S^{-1})_{2,2}$ can be used. To simplify the notation, we define the normalized LU factors as

$$\hat{L}_{1,1} = L_{1,1}, \quad \hat{U}_{1,1} = U_{1,1}, \quad \hat{L}_{2,1} = L_{2,1}(L_{1,1})^{-1}, \quad \hat{U}_{1,2} = (U_{1,1})^{-1}U_{1,2}, \quad (11)$$

and Equation (10) can be equivalently given by

$$A^{-1} = \begin{pmatrix} (\hat{U}_{1,1})^{-1}(\hat{L}_{1,1})^{-1} + \hat{U}_{1,2}S_{2,2}^{-1}\hat{L}_{2,1} & -\hat{U}_{1,2}S_{2,2}^{-1} \\ -S_{2,2}^{-1}\hat{L}_{2,1} & S_{2,2}^{-1} \end{pmatrix}. \quad (12)$$

Let us denote by \mathcal{C} the set of indices

$$\{i \mid (L_{2,1})_i \neq 0\} \cup \{j \mid (U_{1,2})_j \neq 0\}, \quad (13)$$

and assume that $S_{2,2}^{-1}$ has already been computed. From Equation (12), it can be readily observed that to compute the i -th element of $A_{2,1}^{-1} \equiv -S_{2,2}^{-1}\hat{L}_{2,1}$ for $i \in \mathcal{C}$, we only need the entries

$$\{(S_{2,2}^{-1})_{i,j} \mid i \in \mathcal{C}, j \in \mathcal{C}\}. \quad (14)$$

The same set of entries of $S_{2,2}^{-1}$ are required to compute selected entries of $A_{1,2}^{-1} \equiv -\hat{U}_{1,2}S_{2,2}^{-1}$. No additional entries of $S_{2,2}^{-1}$ are needed to complete the computation of $A_{1,1}^{-1}$, which involves the matrix product of selected entries of $\hat{U}_{1,2}$ and $A_{2,1}^{-1}$. This procedure can be repeated recursively to compute selected elements of $S_{2,2}^{-1}$ until $S_{2,2}$ is a scalar of size 1. A pseudocode for demonstrating this column-based selected inversion algorithm for symmetric matrix is given in Lin et al. [2011b].

In practice, a column-based sparse factorization and selected inversion algorithm may not be efficient due to the lack of level 3 BLAS operations. For a sparse matrix A , the columns of A and the L factor can be partitioned into supernodes. A supernode is a maximal set of contiguous columns $\mathcal{J} = \{j, j+1, \dots, j+s\}$ of the L factor that have the same nonzero structure below the $(j+s)$ -th row, and the lower triangular part of $L_{\mathcal{J},\mathcal{J}}$ is dense. However, this strict definition can produce supernodes that are either too large or too small, leading to memory usage, load balancing, and efficiency issues. Therefore, in our work, we relax this definition to limit the maximal number

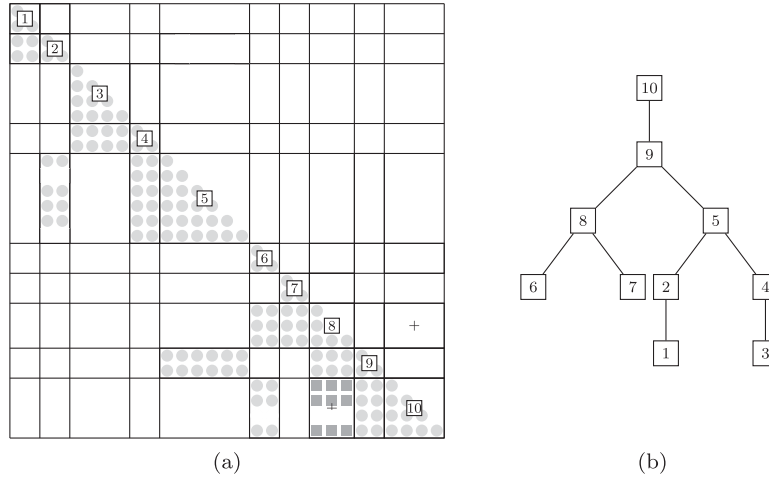


Fig. 1. (a) A structurally symmetric matrix A of size 29×29 divided into 10 supernodes. The nonzero matrix elements in A are labeled as round dots, and the extra fill-in elements in L are labeled as squares. (b) The elimination tree corresponding to the matrix A and its supernode partitioning.

197 of columns in a supernode (i.e., sets are not necessarily maximal). The relaxation also
 198 allows a supernode to include columns for which nonzero patterns are nearly identical
 199 to enhance the efficiency [Ashcraft and Grimes 1989]. This approach is also used in
 200 SuperLU_DIST [Li and Demmel 2003]. Even though the nonzero pattern of the matrix
 201 can be nonsymmetric, the same supernode partitioning is usually applied to the row
 202 partition as well, and we assume that the factorization has been computed using the
 203 structure of $A + A^T$. Then the nonzero structures of L and U are the transpose of each
 204 other. The total number of supernodes is denoted by \mathcal{N} . An example of the supernode
 205 partitioning of a structurally symmetric matrix A together with the extra fill-in in its
 206 L factor (U factor omitted due to structural symmetry) are given in Figure 1(a).

207 Using the notation of supernodes, a pseudocode for the selected inversion algorithm
 208 is given in Algorithm 1. The key step to gain computational efficiency in the selected
 209 inversion algorithm is step 2, which identifies the collection of all nonzero row and
 210 column indices corresponding to the supernode \mathcal{K} , denoted by C . All subsequent steps
 211 operate only on these nonzero rows and columns within the sparsity pattern of the
 212 selected elements, thereby significantly reducing the computational cost.

213 It should be noted that if A is a sparse symmetric matrix, the normalized LU factors
 214 satisfy the relation

$$\hat{U}_{c,\mathcal{K}} = \hat{L}_{\mathcal{K},c}^T. \quad (15)$$

215 To simplify the implementation, the entire diagonal block $A_{\mathcal{K},\mathcal{K}}^{-1}$ is computed even though
 216 it is symmetric. Due to round-off error, the numerical update in step 4 may not pre-
 217 serve this symmetry in finite precision. Our numerical results indicate that the loss
 218 of symmetry may accumulate, especially for ill-conditioned matrices. To reduce such
 219 error for symmetric matrices, we can simply symmetrize the diagonal block of A^{-1} by
 220 performing $A_{\mathcal{K},\mathcal{K}}^{-1} \leftarrow \frac{1}{2}(A_{\mathcal{K},\mathcal{K}}^{-1} + A_{\mathcal{K},\mathcal{K}}^{-T})$ after step 4 for each \mathcal{K} .

221 Furthermore, it should be noted that in the symmetric case, Equation (12) can be
 222 simplified using an LDL^T factorization, which is more efficient than an LU factor-
 223 ization, where L is a unit lower triangular matrix and D is a block diagonal matrix
 224 consisting of 1×1 or 2×2 blocks. The simplification of the selected inversion algorithm

ALGORITHM 1: Selected Inversion Algorithm Based on LU Factorization.

Q4

(1) The supernode partition of columns of A : $\{1, 2, \dots, \mathcal{N}\}$
Input: (2) A supernodal LU factorization of A with (unnormalized) LU factors L and U .
Output: Selected elements of A^{-1} (i.e., $A_{\mathcal{I},\mathcal{J}}^{-1}$ such that $L_{\mathcal{I},\mathcal{J}}$ is not an empty block).
for $\mathcal{K} = \mathcal{N}, \mathcal{N} - 1, \dots, 1$ **do**
1 Find the collection of indices
 $\mathcal{C} = \{\mathcal{I} \mid \mathcal{I} > \mathcal{K}, L_{\mathcal{I},\mathcal{K}} \text{ is a nonzero block}\} \cup \{\mathcal{J} \mid \mathcal{J} > \mathcal{K}, U_{\mathcal{K},\mathcal{J}} \text{ is a nonzero block}\}$
2 $\hat{L}_{\mathcal{C},\mathcal{K}} \leftarrow L_{\mathcal{C},\mathcal{K}}(L_{\mathcal{K},\mathcal{K}})^{-1}$, $\hat{U}_{\mathcal{K},\mathcal{C}} \leftarrow (U_{\mathcal{K},\mathcal{K}})^{-1}U_{\mathcal{K},\mathcal{C}}$
 end
 for $\mathcal{K} = \mathcal{N}, \mathcal{N} - 1, \dots, 1$ **do**
 Find the collection of indices
 $\mathcal{C} = \{\mathcal{I} \mid \mathcal{I} > \mathcal{K}, L_{\mathcal{I},\mathcal{K}} \text{ is a nonzero block}\} \cup \{\mathcal{J} \mid \mathcal{J} > \mathcal{K}, U_{\mathcal{K},\mathcal{J}} \text{ is a nonzero block}\}$
3 Calculate $A_{\mathcal{C},\mathcal{K}}^{-1} \leftarrow -A_{\mathcal{C},\mathcal{C}}^{-1}\hat{L}_{\mathcal{C},\mathcal{K}}$
4 Calculate $A_{\mathcal{K},\mathcal{K}}^{-1} \leftarrow U_{\mathcal{K},\mathcal{K}}^{-1}L_{\mathcal{K},\mathcal{K}}^{-1} - \hat{U}_{\mathcal{K},\mathcal{C}}A_{\mathcal{C},\mathcal{K}}^{-1}$
5 Calculate $A_{\mathcal{K},\mathcal{C}}^{-1} \leftarrow -\hat{U}_{\mathcal{K},\mathcal{C}}A_{\mathcal{C},\mathcal{C}}^{-1}$
 end

with an LDL^T factorization can be found in Lin et al. [2011b]. For symmetric matrices, $A_{\mathcal{K},\mathcal{C}}^{-1}$ (step 5) is readily obtained as the transpose of $A_{\mathcal{C},\mathcal{K}}^{-1}$ without extra computation.

2.2. Elimination Tree

Both the factorization and the selected inversion can be conveniently described in terms of traversals of an elimination tree [Liu 1990]. Each node of the tree corresponds to a supernode of A . A node \mathcal{R} is the parent of a node \mathcal{K} if and only if

$$\mathcal{R} = \min\{\mathcal{I} > \mathcal{J} \mid L_{\mathcal{I},\mathcal{J}} \text{ is a nonzero block}\}. \quad (16)$$

An example of the elimination tree corresponding to the matrix in Figure 1(a) is given in Figure 1(b).

In the factorization procedure, the traversal of the elimination tree is a bottom-up process that starts from the leaves of the tree. A parent supernode cannot be factored until the supernodes associated with all of its children in the tree have been factored. This type of task dependency also determines the amount of concurrency that can be exploited to speed up the factorization on a parallel computer.

In the selected inversion procedure, the traversal of the elimination tree is a top-down process that starts from the root of the tree. Computing the selected elements in the \mathcal{K} -th supernode of A^{-1} requires the selected elements of A^{-1} already computed at ancestor nodes of \mathcal{K} , but not those computed at its sibling nodes and their descendants. Consequently, the selected inversion of supernodes that belong to two different branches of the elimination tree can in principle be carried out independently as long as the selected elements computed at supernodes above these branches have been passed to processors that are assigned to work on these branches.

3. DISTRIBUTED MEMORY PARALLEL SELECTED INVERSION ALGORITHM

In this article, we present the distributed memory PSelInv method. Our first implementation focuses on the case of symmetric matrices. For such matrices, the selected inversion algorithms described in Algorithm 1 require a sparse LU or LDL^T factorization of A to be available first. Here we use the SuperLU_DIST software package [Li and Demmel 2003] to obtain the LU factorization, which has been shown to be scalable to a large number of processors on distributed memory parallel machines.

253 The relatively simple data structure of SuperLU_DIST allows easy access to sparse L
 254 and U factors. However, the main ideas that we develop here can be combined with
 255 other sparse matrix solvers, such as MUMPS [Amestoy et al. 2001], and PARDISO [Schenk
 256 and Gartner 2006] as well, which provides the LDL^T functionality and can potentially
 257 be two times faster in the factorization phase. We also note that only symmetric per-
 258 mutation of the matrix A is allowed, even though SuperLU_DIST allows the column
 259 permutation to be different from the row permutation. In this work, to preserve the
 260 symmetry of the matrix, we do not perform the equilibration procedure often used
 261 in the LU factorization to modify poorly scaled matrix elements. In the current im-
 262 plementation of the PSe1Inv method, we explicitly take advantage of the symmetry
 263 of the matrix and only compute the lower triangular part of the selected elements of
 264 A^{-1} . However, our implementation is not optimal in terms of memory allocation in the
 265 sense that both the upper and lower triangular part of A^{-1} are stored. As will be seen
 266 in the following, such a strategy simplifies the communication pattern and efforts for
 267 bookkeeping in the case of 2D block cyclic data distribution and facilitates generalizing
 268 PSe1Inv to asymmetric matrices. We also note that the suboptimal memory allocation
 269 is not a severe limitation of the PSe1Inv method. The memory footprint of PSe1Inv can
 270 be further optimized for applications that are constrained by the memory usage. Our
 271 numerical results also indicate that the additional memory usage by PSe1Inv is rela-
 272 tively small compared to that used in other procedures, such as the parallel numerical
 273 factorization.

274 We use the same 2D block cyclic distribution scheme employed in SuperLU_DIST to
 275 partition and distribute both the L factor and the selected elements of A^{-1} to be com-
 276 puted. We will review the main features of this type of distribution in Section 3.1. In the
 277 2D block cyclic distribution scheme, each supernode \mathcal{K} is assigned to and partitioned
 278 among a subset of processors. However, computing the selected elements of A^{-1} con-
 279 tained in the supernode \mathcal{K} requires retrieving previously computed selected elements
 280 of A^{-1} that belong to ancestors of \mathcal{K} in the elimination tree. These selected elements
 281 may reside on other processors. As a result, communication is required to transfer
 282 data among different processors to complete steps 3, 4, and 5 of Algorithm 1 in each
 283 iteration. We will discuss how this is done in Section 3.2. The key to reducing com-
 284 munication cost and achieving scalable performance is to overlap communication with
 285 computation by using asynchronous point-to-point MPI functions, even though some of
 286 these communication events are collective in nature (e.g., broadcast and reduce) within
 287 a communication subgroup.

288 In addition to utilizing a fine-grain level of parallelism in computing A^{-1} for each su-
 289 pernode, we introduce a coarse-grain level of parallelism by exploiting the concurrency
 290 available in the elimination tree. This amounts to executing different iterates of the
 291 for loop in Algorithm 1 in parallel. Although the elimination tree may exhibit many
 292 independent tasks associated with supernodes that belong to different branches of the
 293 elimination tree, the 2D block cyclic distribution of L and A^{-1} may prevent these tasks
 294 from being performed completely simultaneously on different processors. The key to
 295 minimizing the dependency issue is to properly assign the order of computational tasks
 296 and to overlap computation and communication as much as possible. We will discuss
 297 our preliminary strategy for improving the parallel efficiency using the elimination
 298 tree in Section 3.3.

299 3.1. Distributed Data Layout and Structure

300 As discussed in Section 2, the columns of A , L , and U are partitioned into supernodes.
 301 Different supernodes may have different sizes. The same partition is applied to the
 302 rows of these matrices to create a 2D block partition of these matrices. The submatrix
 303 blocks are mapped to processors that are arranged in a virtual 2D grid of dimension

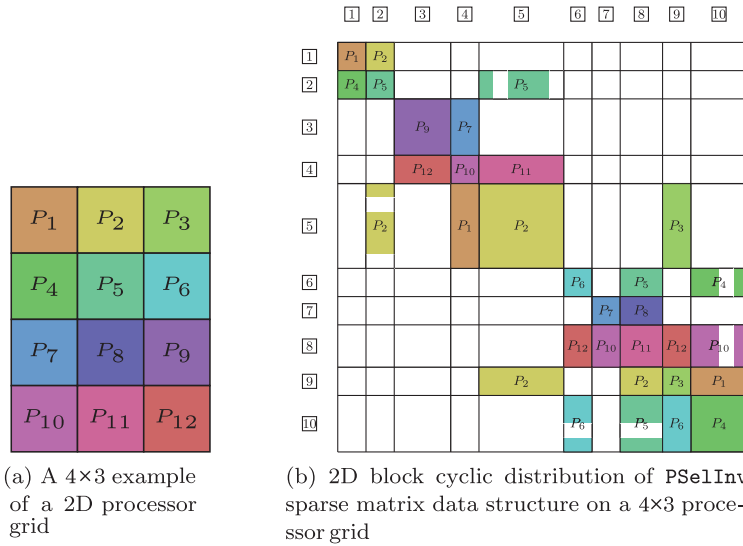


Fig. 2. Data layout of the internal sparse matrix data structure used by PSe1Inv.

$\text{Pr} \times \text{Pc}$ in a cyclic fashion as follows: the $(\mathcal{I}, \mathcal{J})$ -th matrix block is held by the processor labeled as

$$P_{\text{mod}(\mathcal{I}-1, \text{Pr}) \times \text{Pc} + \text{mod}(\mathcal{J}-1, \text{Pc}) + 1}. \quad (17)$$

This is called a *2D block cyclic data-to-processor mapping*. The mapping itself does not take the sparsity of the matrix into account. If the $(\mathcal{I}, \mathcal{J})$ -th block contains only zero elements, then that block is not stored. It is possible that some nonzero blocks may contain several rows of zeros. These rows are not stored either. As an example, a 4×3 grid of processors is depicted in Figure 2(a). The mapping between the 2D supernode partition of the matrix in Figure 1(a) and the 2D processor grid in Figure 2(a) is depicted in Figure 2(b). Each supernodal block column of L is distributed among processors that belong to a column of the processor grid. Each processor may own multiple matrix blocks. For instance, the nonzero rows in the second supernode are owned by processors P_2 and P_5 . More precisely, P_2 owns two nonzero blocks, whereas P_5 is responsible for one block. Note that these nonzero blocks are not necessarily contiguous in the global matrix. Although the nonzero structure of A is not taken into account during the distribution, it has been shown in practice that 2D layouts lead to higher scalability for both dense [Blackford 1997], sparse Cholesky factorizations [Rothberg and Gupta 1994], and LU factorization [Li and Demmel 2003].

In the current implementation, PSe1Inv contains an interface that is compatible with the SuperLU_DIST software package. To allow PSe1Inv to be easily integrated with other LDL^T or LU factorization codes, we create some intermediate sparse matrix objects to hold the distributed L and U factors. Such intermediate sparse matrix objects will be overwritten by selected elements of A^{-1} in the selected inversion process. Each nonzero block $L_{\mathcal{I}, \mathcal{J}}$ is stored as follows. Diagonal blocks $L_{\mathcal{I}, \mathcal{I}}$ are always stored as dense matrices. Nonzero entries of $L_{\mathcal{I}, \mathcal{J}}$ ($\mathcal{I} > \mathcal{J}$) are stored contiguously as a dense matrix in a column-major order, even though row indices associated with the stored matrix elements are not required to be contiguous. As mentioned at the beginning of Section 3, our implementation is not optimal in terms of memory allocation for symmetric matrices in the sense that the nonzero entries within $U_{\mathcal{I}, \mathcal{J}}$ ($\mathcal{I} < \mathcal{J}$) are also stored as a dense matrix in a contiguous array in a column major order, even

333 though the values of $U_{\mathcal{I},\mathcal{J}}$ are identical to those of $L_{\mathcal{J},\mathcal{I}}^T$ for symmetric matrices. The
 334 nonzero column indices associated with the nonzeros entries in $U_{\mathcal{I},\mathcal{J}}$ are not required
 335 to be contiguous either. We remark that for matrices with highly asymmetric sparsity
 336 patterns, it is more efficient to store the upper triangular blocks using the skyline
 337 structure shown in Li and Demmel [2003]. However, we choose to use a simpler data
 338 layout because it allows level-3 BLAS (GEMM) to be used in the selected inversion
 339 process.

340 3.2. Computing Selected Elements of A^{-1} within Each Supernode in Parallel

341 In this section, we detail how steps 2 through 5 in Algorithm 1 can be completed in
 342 parallel.

343 We perform step 2 of Algorithm 1 in a separate pass, as the data communication
 344 required in this step is relatively simple. The processor that owns the block $L_{\mathcal{K},\mathcal{K}}$
 345 broadcasts $L_{\mathcal{K},\mathcal{K}}$ to all other processors within the same column processor group owning
 346 nonzero blocks $L_{\mathcal{I},\mathcal{K}}$ in the supernode \mathcal{K} . Each processor in that group performs the
 347 triangular solve $\hat{L}_{\mathcal{I},\mathcal{K}} \equiv L_{\mathcal{I},\mathcal{K}}(L_{\mathcal{K},\mathcal{K}})^{-1}$ for each nonzero block contained in the set \mathcal{C}
 348 defined in step 1 of the algorithm. Because $L_{\mathcal{I},\mathcal{K}}$ is not used in the subsequent steps
 349 of selected inversion once $\hat{L}_{\mathcal{I},\mathcal{K}}$ has been computed, it is overwritten by $\hat{L}_{\mathcal{I},\mathcal{K}}$. Since
 350 communication is limited to a processor column group only, step 2 can be carried out
 351 for multiple supernodes at the same time.

352 A more complicated communication pattern is required to complete step 3 in parallel.
 353 Because $A_{\mathcal{C},\mathcal{C}}^{-1}$ and $\hat{L}_{\mathcal{C},\mathcal{K}}$ are generally owned by different processor groups, there are two
 354 possible ways to carry out the multiplication of $A_{\mathcal{C},\mathcal{C}}^{-1}$ with $\hat{L}_{\mathcal{C},\mathcal{K}}$. The first approach
 355 is to send blocks of $\hat{L}_{\mathcal{C},\mathcal{K}}$ to processors that own the matching blocks of $A_{\mathcal{C},\mathcal{C}}^{-1}$ so that
 356 matrix-matrix multiplication can be performed on processors owning $A_{\mathcal{C},\mathcal{C}}^{-1}$. The second
 357 approach is to send data in the opposite direction—that is, one can send blocks of $A_{\mathcal{C},\mathcal{C}}^{-1}$
 358 to the matching blocks of $\hat{L}_{\mathcal{C},\mathcal{K}}$ so that matrix-matrix multiplication can be performed
 359 on processors owning $\hat{L}_{\mathcal{C},\mathcal{K}}$.

360 To compare the cost of these two approaches, let us first consider the case in which all
 361 blocks $\hat{L}_{\mathcal{I},\mathcal{K}}$ with $\mathcal{I} \geq \mathcal{K}$ are dense matrix blocks of equal size $m_b \times m_b$, and $\mathcal{C} = \{\mathcal{I} | \mathcal{I} \geq \mathcal{K}\}$.
 362 This is approximately the case when \mathcal{K} is near the root of the elimination tree. We
 363 assume that there are $\sqrt{P} \times \sqrt{P}$ processors, and thus the size of the set \mathcal{C} is $m_b \sqrt{P}$.
 364 We also assume that the matrix blocks in $\hat{L}_{\mathcal{C},\mathcal{K}}$ are distributed among \sqrt{P} processors
 365 within the same column group, and each processor in this processor column also holds
 366 a dense block $\hat{L}_{\mathcal{I},\mathcal{K}}$ of size $m_b \times m_b$. In the first approach, the computation is performed
 367 in parallel on P processors, and the computational cost on each processor is $\mathcal{O}(m_b^3)$. In
 368 the second approach, the computation is performed in parallel on \sqrt{P} processors only,
 369 and the computational cost on these processors is $\mathcal{O}(m_b^3 \sqrt{P})$. All other processors are
 370 idle in the computational step, and this leads to severe load imbalance when P is large.

371 We implemented the first approach in PSelInv. This requires sending the $\hat{L}_{\mathcal{I},\mathcal{K}}$
 372 block from a particular processor to all processors within the same column group
 373 of processors among which $A_{\mathcal{C},\mathcal{I}}^{-1}$ is distributed. However, since the processor owning
 374 $\hat{L}_{\mathcal{I},\mathcal{K}}$ is generally not in the same processor communication group that owns $A_{\mathcal{C},\mathcal{I}}^{-1}$,
 375 sending $\hat{L}_{\mathcal{I},\mathcal{K}}$ to processors that hold the distributed blocks of $A_{\mathcal{C},\mathcal{I}}^{-1}$ cannot be done
 376 by a single broadcast. Indeed, for this to be possible, communication groups (i.e.,
 377 MPI communicators) would have to be created for each and every different sparse
 378 row/column structure. This is generally not possible, as the maximum number of

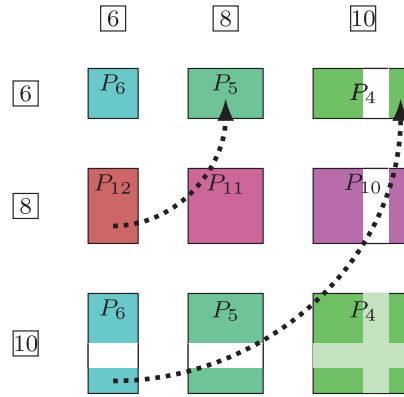


Fig. 3. Processors holding $\hat{L}_{8,6}$ and $\hat{L}_{10,6}$ send data to processors holding the cross-diagonal blocks and overwrite $\hat{U}_{6,8}$ and $\hat{U}_{6,10}$.

allowed MPI communicators is typically much smaller than needed. Therefore, one way to complete this step of data communication is to use a number of point-to-point MPI sends that originate from the processor that owns $\hat{L}_{\mathcal{I},\mathcal{K}}$ and terminate on the group of processors that own the nonzero blocks of $A_{\mathcal{C},\mathcal{I}}^{-1}$. In addition to incurring higher communication latency cost, this approach also leads to significant bookkeeping effort to track the sources and destinations of all messages for each processor.

In our current implementation, we simplify the data communication pattern by storing both $\hat{L}_{\mathcal{I},\mathcal{K}}$ and $\hat{U}_{\mathcal{K},\mathcal{I}}$ even when A is symmetric. We acknowledge that such implementation is not optimal in terms of memory allocation and can be improved if applications are constrained by memory usage for symmetric matrices. As soon as $\hat{L}_{\mathcal{I},\mathcal{K}}$ becomes available as illustrated earlier, we send the $\hat{L}_{\mathcal{I},\mathcal{K}}$ block to the processor that owns $\hat{U}_{\mathcal{K},\mathcal{I}}$, and $\hat{U}_{\mathcal{K},\mathcal{I}}$ is overwritten by $\hat{L}_{\mathcal{I},\mathcal{K}}^T$. Figure 3 illustrates how this step is carried out for a specific supernode $\mathcal{K} = \mathbb{6}$ of the matrix described in Figure 2(b). Once $\hat{L}_{8,6}$ is computed on P_{12} , the block is sent to P_5 . The P_5 processor then overwrites $\hat{U}_{6,8}$ by $\hat{L}_{8,6}^T$. Similarly, the $\hat{L}_{10,6}$ block is computed on P_6 and sent to P_4 , on which $\hat{U}_{6,10}$ is overwritten by $\hat{L}_{10,6}$.

With $\hat{L}_{\mathcal{I},\mathcal{K}}$ properly placed on the processors that are mapped to the upper triangular part of the distributed \hat{U} matrix, step 3 of Algorithm 1 can proceed as follows. The $\hat{U}_{\mathcal{K},\mathcal{I}} = \hat{L}_{\mathcal{I},\mathcal{K}}^T$ block is first sent to all processors within the same column processor group that owns $\hat{U}_{\mathcal{K},\mathcal{I}}$. The matrix-matrix multiplication $A_{\mathcal{J},\mathcal{I}}^{-1} \hat{L}_{\mathcal{I},\mathcal{K}}$ is then performed locally on each processor owning $A_{\mathcal{J},\mathcal{I}}^{-1}$ using the GEMM subroutine in BLAS3. Then local matrix contributions $A_{\mathcal{J},\mathcal{I}}^{-1} \hat{L}_{\mathcal{I},\mathcal{K}}$ are reduced within each row communication groups owning $\hat{L}_{\mathcal{J},\mathcal{K}}$ to produce the $A_{\mathcal{J},\mathcal{K}}^{-1}$ block in step 3 of Algorithm 1.

Figure 4 illustrates how this step is completed for a specific supernode $\mathcal{K} = \mathbb{6}$, for the matrix depicted in Figure 2(b). We use circled letters (a), (b), (c) to label communication events and circled numbers ①, ②, ③ to label computational events. We can see from this figure that $\hat{U}_{6,8} = \hat{L}_{8,6}^T$ is sent by P_5 to all processors within the same column processor group to which P_5 belongs. This group includes both P_5 and P_{11} . Similarly, $\hat{L}_{10,6}$ is broadcast from P_4 to all other processors within the same column group to which P_4 belongs. Local matrix matrix multiplications are then performed on P_{11} , P_{10} , P_4 , and P_5 simultaneously. The distributed products are then reduced onto P_{12} and P_5

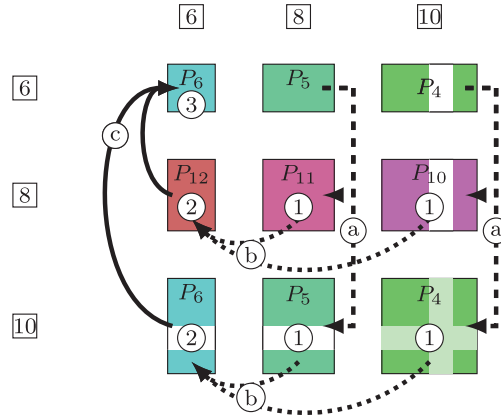


Fig. 4. Task parallelism and communication pattern for the supernode $\mathbb{6}$. There are six steps: \textcircled{a} broadcast \hat{L} , $\textcircled{1}$ compute $A^{-1}\hat{L}$, \textcircled{b} reduce $A^{-1}\hat{L}$, $\textcircled{2}$ compute $\hat{L}^T A^{-1}\hat{L}$, \textcircled{c} reduce $\hat{L}^T A^{-1}\hat{L}$, and $\textcircled{3}$ update A^{-1} .

410 respectively within the row processor groups to which they belong. After this step, $A_{8,6}^{-1}$
 411 and $A_{10,6}^{-1}$ respectively become available on P_{12} and P_6 .

412 On completion of step 3, the matrix product $\hat{U}_{\mathcal{K},\mathcal{J}} A_{\mathcal{J},\mathcal{K}}^{-1} \equiv \hat{L}_{\mathcal{J},\mathcal{K}}^T A_{\mathcal{J},\mathcal{K}}^{-1}$ is first computed
 413 locally on the processor holding $\hat{L}_{\mathcal{J},\mathcal{K}}$ and then reduced to the processor that owns
 414 the diagonal block $L_{\mathcal{K},\mathcal{K}}$ within the column processor group to which the supernode
 415 \mathcal{K} is mapped. The sum of the distributed matrix product $-\hat{L}_{\mathcal{J},\mathcal{K}}^T A_{\mathcal{J},\mathcal{K}}^{-1}$ is then added
 416 to $(U_{\mathcal{K},\mathcal{K}})^{-1}(L_{\mathcal{K},\mathcal{K}})^{-1}$ computed on the processor holding $L_{\mathcal{K},\mathcal{K}}$. This completes step 4 of
 417 Algorithm 1. As an example, again we use Figure 4 for $\mathcal{K} = \mathbb{6}$. $\hat{L}_{8,6}^T A_{8,6}^{-1}$ is computed on
 418 P_{12} and sent to P_6 . Similarly, $\hat{L}_{10,6}^T A_{10,6}^{-1}$ is computed on P_6 . Since both $\hat{L}_{10,6}^T$ and $L_{6,6}$ are
 419 held by P_6 , no further data communication is necessary. Finally, P_6 updates $A_{6,6}^{-1}$. As we
 420 discussed in Section 2.1, for symmetric matrices, $A_{\mathcal{K},\mathcal{K}}^{-1}$ should be explicitly symmetrized
 421 to reduce the effect of rounding errors.

422 Since A is symmetric, dtep 5 of Algorithm 1 can be simplified as follows. We first
 423 overwrite $\hat{L}_{\mathcal{J},\mathcal{K}}$ by $A_{\mathcal{J},\mathcal{K}}^{-1}$ locally on the processor holding $\hat{L}_{\mathcal{J},\mathcal{K}}$ ($\mathcal{J} > \mathcal{K}$). We then send $A_{\mathcal{J},\mathcal{K}}^{-1}$
 424 to the processor holding $\hat{U}_{\mathcal{K},\mathcal{J}}$ and overwrite $\hat{U}_{\mathcal{K},\mathcal{J}}$ by $(A_{\mathcal{J},\mathcal{K}}^{-1})^T$. The data communication
 425 pattern for this step is the same as described in Figure 3. After step 5, we move to the
 426 next supernode ($\mathcal{K} - 1$).

427 3.3. Exploiting Concurrency in the Elimination Tree

428 In this section, we discuss how to add an additional coarse-grain level of parallelism
 429 to the selected inversion algorithm by exploiting task concurrency exposed by the
 430 elimination tree.

431 As we indicated in Section 2.2, two supernodes belonging to two separate branches
 432 of the elimination tree can be processed independently if the selected elements of the
 433 inverse belonging to their ancestors have been computed, and if these supernodes
 434 and the ancestors on which they depend are mapped onto different sets of processors.
 435 Although it is possible to pass the previously computed selected elements of A^{-1} from
 436 the ancestors down to their children as we move down the elimination tree, algorithms
 437 based on this approach (e.g., a multifrontal like algorithm [Lin et al. 2011a]) would
 438 require additional work space to hold extra copies of the selected elements.

To reduce the amount of extra work space, which can grow rapidly as we go down the elimination tree, we choose to allow processors assigned to each supernode to communicate back and forth with processors assigned to its ancestors in the way that we described in Section 3.2 to complete step 3 of Algorithm 1.

However, the drawback of this approach is that at some point, two supernodes belonging to two separate branches of the elimination tree may not be processed simultaneously when they need to communicate with their common ancestors at the same time. At this point of conflict, only one of them should be allowed to initiate and complete the data communication with the common ancestor at a time. For example, when supernodes $\mathbb{2}$ and $\mathbb{4}$ in Figure 1(a) are being processed on different sets of processors, both of them may need to communicate with processors assigned to supernode $\mathbb{5}$ at the same time. In this case, the updates to be performed on these processors cannot proceed completely independently. On the other hand, if the set of processors assigned to update two supernodes are completely different, then at least some of the updates can be computed simultaneously.

To exploit the type of concurrency discussed previously, which occurs at the `for` loop level in Algorithm 1, we create a basic parallel task scheduler to launch different iterates of the `for` loop in a certain order. This order is defined by a priority list S , which is indexed by integer priority numbers ranging from 1 to n_s , where n_s is bounded from above by the depth of the elimination tree. The task performed in each iteration of the `for` loop is assigned a priority number $\sigma(\mathcal{I})$. The lower the number, the higher the priority of the task, and hence the sooner it is scheduled. The supernode \mathcal{N} associated with the root of the elimination tree clearly has to be processed first. If the k -th element of S contains multiple supernodes or tasks whose priority numbers are k , the order in which these tasks are completed can be arbitrary. A recipe for assigning priority number of different tasks (or equivalently, supernodes) is shown in Algorithm 2. We assume that the elimination tree is post ordered.

Even though we use a priority list to help launch tasks, we do not place extra synchronization among launched tasks other than requiring them to preserve data dependency. Tasks associated with different supernodes can be executed concurrently if these supernodes are on different critical paths of the elimination tree, and if there is no overlap among processors mapped to these critical paths. In fact, if tasks associated with supernode \mathcal{J} and \mathcal{I} are mapped to different sets of processors, the task associated with the supernode \mathcal{J} may actually start before that associated with another supernode \mathcal{I} even if $\sigma(\mathcal{I}) < \sigma(\mathcal{J})$ —that is, even if task \mathcal{I} is scheduled ahead of task \mathcal{J} according to the priority list. When two different tasks need to communicate with a common ancestor, the priority number associated with each task determines which task is completed first.

ALGORITHM 2: Assign Priority Numbers to Supernodes and Create a Priority List.

Input: A list of supernodes $\{\mathcal{I}\}$ and the elimination tree associated with these supernodes.
Output: An array σ , $\sigma(\mathcal{I})$ gives the priority number of the task associated with supernode \mathcal{I} ;
 an array S of n_s supernode lists, $S(i)$ gives a set of supernodes with priority number i ,
 $1 \leq i \leq n_s$.

$\sigma(\mathcal{N}) = 1$
 $S(1) = \{\mathcal{N}\}$
for $\mathcal{I} = \mathcal{N} - 1$ **down to** 1 **do**
 | $\sigma(\mathcal{I}) = \sigma(\text{parent}(\mathcal{I})) + 1$
 | $S(\sigma(\mathcal{I})) = S(\sigma(\mathcal{I})) \cup \{\mathcal{I}\}$
end

Q5

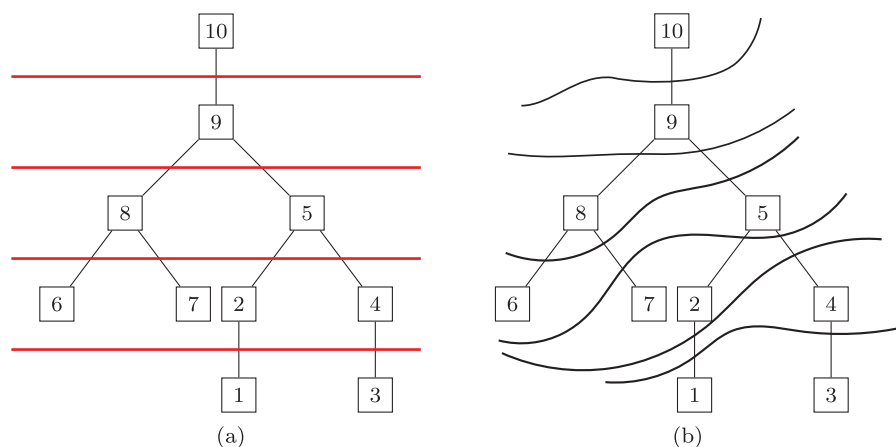


Fig. 5. Elimination tree of matrix A and two possible priority lists S .

477 We remark that there is some flexibility in assigning a priority number to each
 478 supernode and constructing the priority list S . For instance, we can use the strategy
 479 given by Algorithm 2, which simply defines $\sigma(\mathcal{I})$ by the distance (in terms of the number
 480 of edges) between the supernode \mathcal{I} and the root of the elimination tree. For the same
 481 elimination tree shown in Figure 5(a), another possible construction of the σ list is
 482 illustrated in Figure 5(b), which assigns the same σ value to supernodes at different
 483 levels of the elimination tree. The latter construction takes into account how supernodes
 484 are distributed among different processors, as we will discuss in the following.

485 The priority list S determines the order in which computational tasks associated with
 486 different supernodes are completed. Because the amount of work and communication
 487 performed by each of the supernodes can vary significantly, different priority lists can
 488 lead to different overall performance. The actual performance of parallel selected in-
 489 version depends on the sparsity pattern of the matrix, as well as the processor grid, and
 490 is therefore difficult to predict a priori. We refer readers to our report [Jacquelin et al.
 491 2014] for a detailed example on the difficulty of designing an optimal task schedule.

492 With the help of the priority list, we can implement the f or loop level of parallelism in
 493 Algorithm 1 in a way that is described in Algorithm 3. To illustrate the relation clearly,
 494 Algorithm 3 uses the same numeric ordering of the steps as that in Algorithm 1. Algo-
 495 rithm 3 also makes use of another array of lists `procmap`. The \mathcal{K} -th element of `procmap`
 496 contains the list of all processors participating in steps 3 through 5 in Algorithm 1.
 497 The communication steps are described within parentheses. We also remark that we
 498 do not place MPI barriers between supernodes explicitly to exploit parallelism among
 499 the computation for different supernodes. For symmetric matrices, the diagonal blocks
 500 should be symmetrized, as indicated at the end of Section 2.1 for symmetric cases.

501 4. NUMERICAL RESULTS

502 To assess the performance of `PSe1Inv`, we conducted several computational experi-
 503 ments, which we report in this section.

504 Our test problems are taken from various sources, including the Harwell-Boeing
 505 Test Collection [Duff et al. 1992] and the University of Florida Matrix Collection
 506 [Davis and Hu 2011], and matrices generated from electronic structure software, in-
 507 cluding SIESTA [Soler et al. 2002] and DGDFT [Lin et al. 2012]. The first two matrix
 508 collections are widely used benchmark problems for testing sparse direct methods,
 509 whereas the other test problems come from practical large-scale electronic structure

ALGORITHM 3: The Parallel Selected Inversion Algorithm (for Symmetric Matrices).

Input: (1) The supernode partition of columns of a sparse symmetric matrix $A: \{1, 2, \dots, \mathcal{N}\}$; a priority list $\{S(k): k = 1, 2, \dots, n_s\}$;
(2) L and U factors through a supernodal LU factorization (or equivalent LDL^T factorization) of A ;
(3) 2D processor mapping with with $P = \text{Pr} \times \text{Pc}$ processors.

Output: Selected elements of A^{-1} .
[Compute the normalized factors \hat{L} and \hat{U}].

```

for  $k = 1, 2, \dots, n_s$  do
  for each supernode  $\mathcal{K} \in S(k)$  do
    if  $myid \in \text{procmap}(\mathcal{K})$  then
      1 Find the collection of indices
       $\mathcal{C} = \{\mathcal{I} \mid \mathcal{I} > \mathcal{K}, L_{\mathcal{I},\mathcal{K}} \text{ is a nonzero block}\} \cup \{\mathcal{J} \mid \mathcal{J} > \mathcal{K}, U_{\mathcal{K},\mathcal{J}} \text{ is a nonzero block}\}$ 
      2 (Broadcast  $(L_{\mathcal{K},\mathcal{K}})^{-1}$  to processors owning  $L_{\mathcal{I},\mathcal{K}}, \mathcal{I} \in \mathcal{C}$ )
      2  $\hat{L}_{\mathcal{C},\mathcal{K}} \leftarrow L_{\mathcal{C},\mathcal{K}}(L_{\mathcal{K},\mathcal{K}})^{-1}$ 
      2 (Send  $\hat{L}_{\mathcal{I},\mathcal{K}}, \mathcal{I} \in \mathcal{C}$  to the processor holding  $\hat{U}_{\mathcal{K},\mathcal{I}}, \mathcal{I} \in \mathcal{C}$  and overwrite  $\hat{U}_{\mathcal{K},\mathcal{I}}, \mathcal{I} \in \mathcal{C}$ 
      by  $\hat{L}_{\mathcal{I},\mathcal{K}}^T, \mathcal{I} \in \mathcal{C}$ )
    end
  end
end
[Selected inversion process].
for  $k = 1, 2, \dots, n_s$  do
  for each supernode  $\mathcal{K} \in S(k)$  do
    if  $myid \in \text{procmap}(\mathcal{K})$  then
      Find the collection of indices
       $\mathcal{C} = \{\mathcal{I} \mid \mathcal{I} > \mathcal{K}, L_{\mathcal{I},\mathcal{K}} \text{ is a nonzero block}\} \cup \{\mathcal{J} \mid \mathcal{J} > \mathcal{K}, U_{\mathcal{K},\mathcal{J}} \text{ is a nonzero block}\}$ 
      3 (Broadcast  $\hat{U}_{\mathcal{K},\mathcal{I}}, \mathcal{I} \in \mathcal{C}$  to processors holding  $A_{\mathcal{J},\mathcal{I}}^{-1}, \mathcal{I}, \mathcal{J} \in \mathcal{C}$ )
      3 For processors holding  $A_{\mathcal{J},\mathcal{I}}^{-1}, \mathcal{I}, \mathcal{J} \in \mathcal{C}$ , compute locally  $-A_{\mathcal{J},\mathcal{I}}^{-1} \hat{L}_{\mathcal{I},\mathcal{K}}, \mathcal{I}, \mathcal{J} \in \mathcal{C}$ 
      3 (Reduce  $-A_{\mathcal{J},\mathcal{I}}^{-1} \hat{L}_{\mathcal{I},\mathcal{K}}, \mathcal{I}, \mathcal{J} \in \mathcal{C}$  to processors holding  $\hat{L}_{\mathcal{J},\mathcal{K}}, \mathcal{J} \in \mathcal{C}$ , and save the
      result in  $A_{\mathcal{J},\mathcal{K}}^{-1}, \mathcal{J} \in \mathcal{C}$ )
      4 For processors holding  $\hat{L}_{\mathcal{I},\mathcal{K}}, \mathcal{I} \in \mathcal{C}$ , compute locally  $-\hat{L}_{\mathcal{I},\mathcal{K}}^T A_{\mathcal{I},\mathcal{K}}^{-1}, \mathcal{I} \in \mathcal{C}$ 
      4 (Reduce  $-\hat{L}_{\mathcal{I},\mathcal{K}}^T A_{\mathcal{I},\mathcal{K}}^{-1}, \mathcal{I} \in \mathcal{C}$  to the processor holding  $L_{\mathcal{K},\mathcal{K}}$ )
      4 For the processor holding  $(L_{\mathcal{K},\mathcal{K}})^{-1}$ , update  $A_{\mathcal{K},\mathcal{K}}^{-1} \leftarrow (U_{\mathcal{K},\mathcal{K}})^{-1} (L_{\mathcal{K},\mathcal{K}})^{-1} - \hat{L}_{\mathcal{C},\mathcal{K}}^T A_{\mathcal{C},\mathcal{K}}^{-1}$ 
      For the processors holding  $L_{\mathcal{K},\mathcal{K}}^{-1}$ , update  $A_{\mathcal{K},\mathcal{K}}^{-1} \leftarrow \frac{1}{2} (A_{\mathcal{K},\mathcal{K}}^{-1} + A_{\mathcal{K},\mathcal{K}}^{-T})$ 
      3 Overwrite  $\hat{L}_{\mathcal{I},\mathcal{K}}, \mathcal{I} \in \mathcal{C}$  by  $A_{\mathcal{I},\mathcal{K}}^{-1}, \mathcal{I} \in \mathcal{C}$ 
      5 (Send  $A_{\mathcal{I},\mathcal{K}}^{-1}, \mathcal{I} \in \mathcal{C}$  to the processor holding  $\hat{U}_{\mathcal{K},\mathcal{I}}, \mathcal{I} \in \mathcal{C}$ , and overwrite  $\hat{U}_{\mathcal{K},\mathcal{I}}, \mathcal{I} \in \mathcal{C}$ 
      by  $A_{\mathcal{I},\mathcal{K}}^{-1}, \mathcal{I} \in \mathcal{C}$ )
    end
  end
end

```

calculations. The names of these matrices, as well as some of their characteristics, are listed in Tables I and II. The first three problems in these tables come with two matrices each. One of the matrices, denoted by H , is a discretized Hamiltonian, and the other matrix is an overlap matrix, denoted by S . For all other problems, the overlap matrices can be considered as the identity matrix. All matrices are real and symmetric. In all of our experiments, we compute the selected elements of the matrix

$$A(z) = H - zS. \quad (18)$$

Table I. Description of Test Problems for PSe1Inv

Problem	Description
SIESTA_C_BN_1x1	Electronic structure theory, C-BN sheet with 2,532 atoms
SIESTA_C_BN_2x2	Electronic structure theory, C-BN sheet with 10,128 atoms
SIESTA_C_BN_4x2	Electronic structure theory, C-BN sheet with 20,256 atoms
DNA_16	Electronic structure theory, DNA molecule with 11,440 atoms
DNA_715_64cell	Electronic structure theory, DNA molecule with 45,760 atoms
DG_Graphene_2048	Electronic structure theory, graphene with 2,048 atoms
DG_Graphene_8192	Electronic structure theory, graphene with 8,192 atoms
pwtk	Pressurized wind tunnel, stiffness matrix
parabolic_fem	Diffusion-convection reaction, constant homogeneous diffusion
ecology2	Circuitscape: circuit theory applied to animal/gene flow, B. McRae, UCSB
audikw_1	Automotive crankshaft model with more than 900,000 TETRA elements, Audi, GmbH

Table II. Dimension n , the Number of Nonzeros $|A|$, and the Number of Nonzeros of the Cholesky Factor $|L|$ of the Test Problems

Problem	n	$ A $	$ L $
SIESTA_C_BN_1x1	32,916	23,857,418	269,760,112
SIESTA_C_BN_2x2	131,664	95,429,672	1,655,233,542
SIESTA_C_BN_4x2	263,328	190,859,344	3,591,750,262
DNA_16	7,752	2,430,642	9,272,160
DNA_715_64cell	459,712	224,055,744	866,511,698
DG_Graphene_2048	82,944	87,340,032	545,245,344
DG_Graphene_8192	331,776	349,360,128	2,973,952,468
pwtk	217,918	5,926,171	104,644,472
parabolic_fem	525,825	3,674,625	58,028,731
ecology2	999,999	2,997,995	91,073,583
audikw_1	943,695	77,651,847	2,500,489,909

516 For simplicity, we choose $z = 0$ for all efficiency tests in Section 4.1. For some appli-
517 cations, z is chosen to be a complex number with a small imaginary part to ensure
518 that $A(z)$ is nonsingular. This technique is often used in the electronic structure calcu-
519 lation to be discussed in Section 4.3. The LU factorization is performed by using the
520 SuperLU_DIST software package. SuperLU_DIST does not use dynamic pivoting strate-
521 gies, and our matrices are permuted without taking into account the values of matrix
522 entries. Consequently, the efficiency of both SuperLU_DIST and PSe1Inv is independent
523 of the choice of z . The lack of dynamic pivoting strategies may impact the accuracy
524 of PSe1Inv for highly indefinite and nearly singular systems. We study the accuracy
525 for different choices of complex shifts z in Section 4.2. All timing results reported are
526 performed in complex arithmetic computation.

527 In all of our experiments, we used the NERSC Edison platform with Cray XC30
528 nodes. Each node has 24 cores partitioned among two Intel Ivy Bridge processors. Each
529 12-core processor runs at 2.4GHz. A single node has 64GB of memory, providing more
530 than 2.6GB of memory per core. We used one MPI process per core and refer to the core
531 to denote an MPI process.

532 4.1. Parallelization Scalability

533 Four types of experiments were performed to measure the scalability of PSe1Inv. For
534 the first three sets of experiments, all timing data points that we present are averaged
535 measurements over 10 runs, and the error bars shown in Figures 6, 7, 8, and 9 indicate
536 the standard deviation of the measured wall clock time.

537 To clearly show the cost and scalability of selected inversion itself in comparison
538 with the symbolic and numerical LU factorizations, which are required for selected

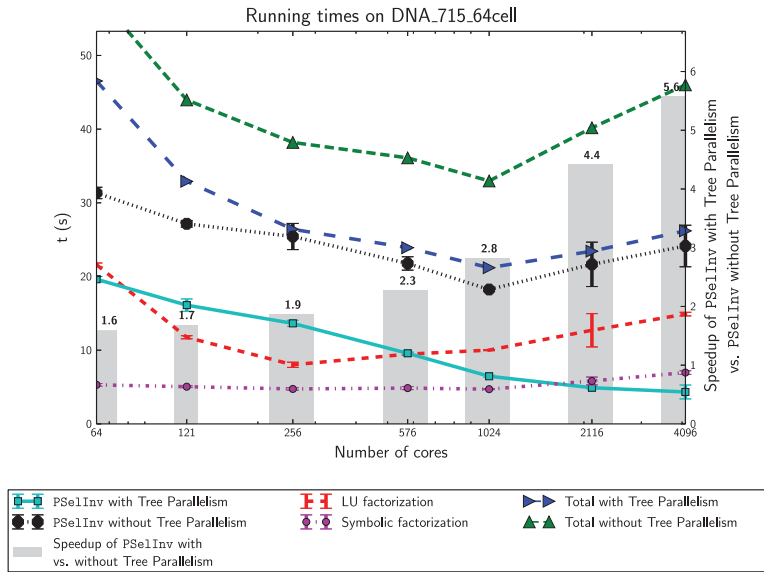


Fig. 6. Wall clock time used by three components (symbolic factorization, numerical LU factorization, and selected inversion) of two versions of PSe1Inv with respect to the number of cores used in the computation for the DNA_715_64cell matrix. In one version, we do not take advantage of the concurrency exposed by the elimination tree, whereas in the other version, we do make use of this tree level of parallelism. The height of each bar in the figure indicates the ratio of wall clock time measured for the former over that measured for the latter.

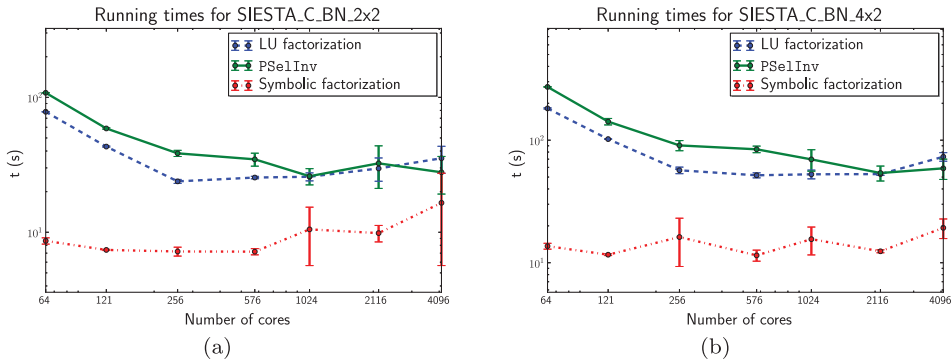


Fig. 7. Strong scalability of PSe1Inv compared to that of LU factorization and symbolic factorization for SIESTA matrices generated for two C_BN systems of different sizes.

inversion, we time the three computational components separately. The symbolic factorization is performed in parallel using PT-Scotch. It is labeled as “symbolic factorization” in the timing figures. The LU factorization is performed by using SuperLU_DIST. It is labeled as “LU factorization.” The selected inversion itself is performed by using PSe1Inv and labeled as “PSe1Inv.” The total time required to obtain the selected elements of the inverse matrix thus corresponds to the sum of all three components.

The first experiment focuses on the impact of the additional parallelism stemming from the elimination tree as discussed in Section 3.3. PSe1Inv is thus tested both with and without this additional level of parallelism. As observed in Figure 6, adding the tree-level parallelism allows the performance of PSe1Inv to scale to 4,096 cores. When

539
540
541
542
543
544
545
546
547
548

21:18

M. Jacquelin et al.

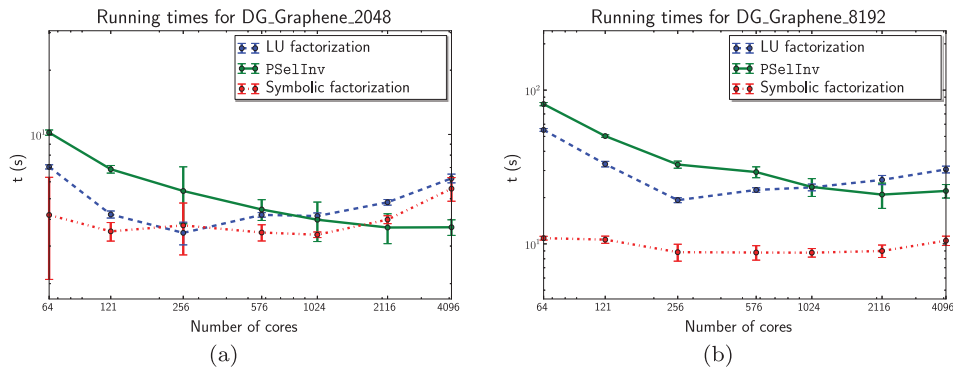


Fig. 8. Strong scalability of PSe1Inv compared to that of LU factorization and symbolic factorization for DG matrices generated for graphene systems of different sizes.

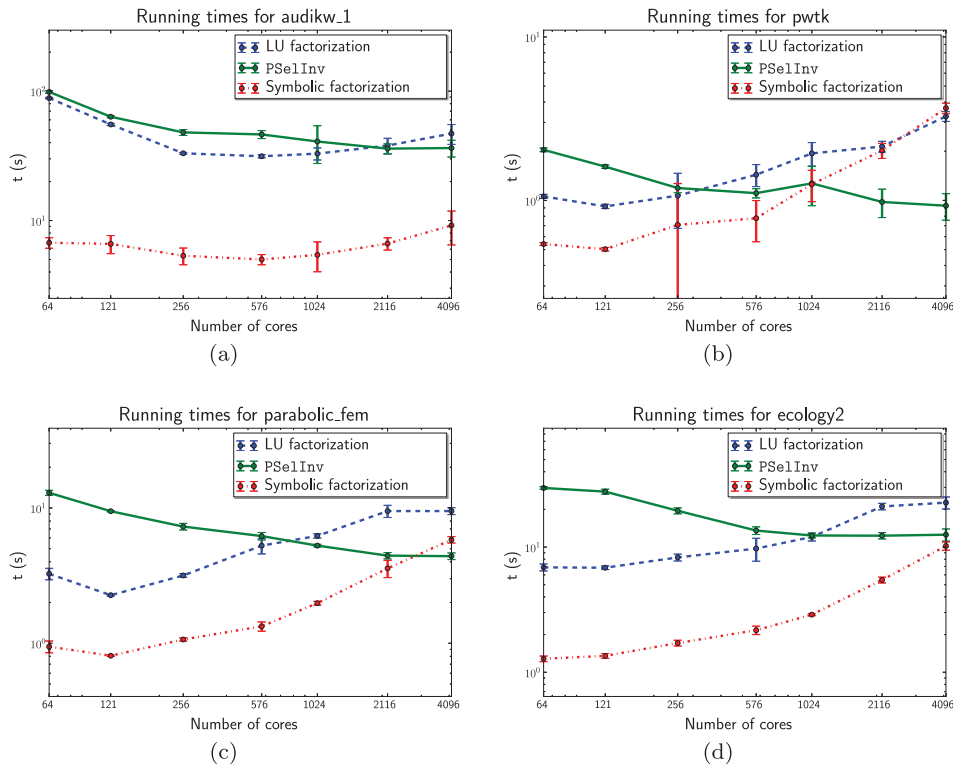


Fig. 9. Strong scalability of PSe1Inv compared to that of LU factorization and symbolic factorization for matrices from the Harwell-Boeing Test Collection and the University of Florida Matrix Collection.

549 4,096 cores are used, adding the tree-level parallelism leads to a 5.6-fold speedup in
 550 the DNA_715_64cell problem. Comparatively, the performance of the LU factorization
 551 and the selected inversion without tree parallelism can only scale up to 1,024 cores.

552 The second set of experiments (Figures 7, 8, and 9) aims at evaluating the strong
 553 scaling of PSe1Inv. In every experiment, the tree-level parallelism is enabled, as it
 554 clearly delivers better performance. PSe1Inv exhibits excellent strong scalability up to

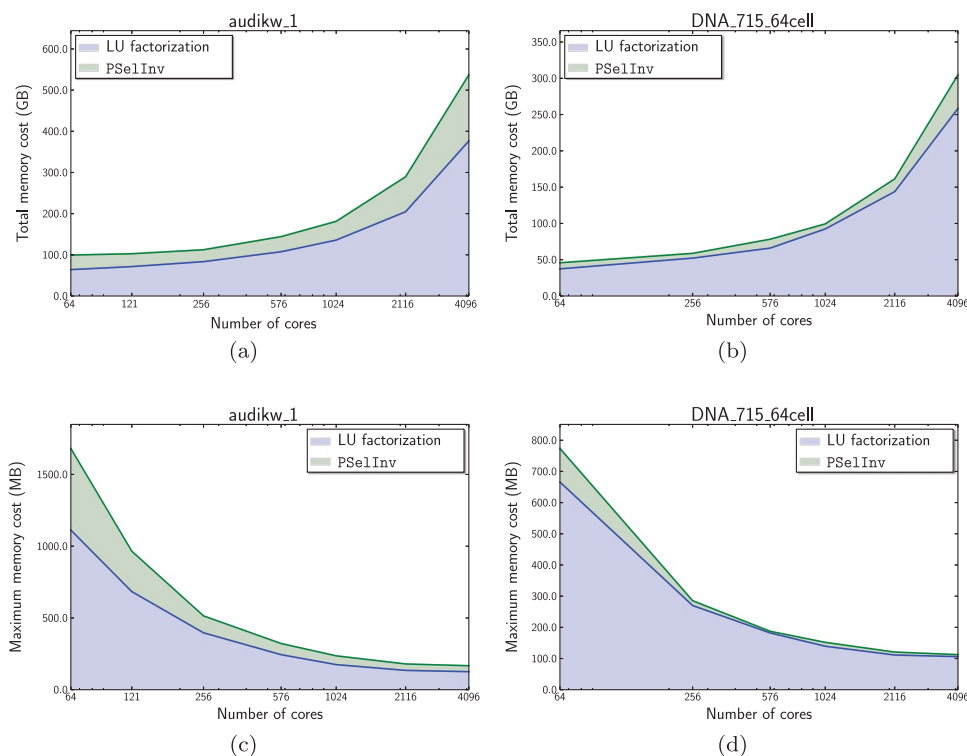


Fig. 10. Cumulative total memory cost and maximum memory cost among all cores for the `audikw_1` matrix and the `DNA_715_64cell` matrix, with shaded regions indicating the additional memory cost introduced in each step of the selected inversion process. The memory cost is measured by the memory high watermark reached at each step.

4,096 cores. For SIESTA matrices (Figure 7), PSe1Inv is slightly slower than LU factorization when the number of cores is less than 1,024 and is faster than LU factorization when more than 2,116 cores are used. For DGDFT matrices (Figure 8), PSe1Inv can be twice as fast as LU factorization, and the running time of PSe1Inv can be comparable to that of symbolic factorization for the `DNA_715_64cell` matrix. For generic sparse matrices (Figure 9) obtained from the University of Florida Collection, PSe1Inv delivers excellent performance on relatively dense matrices, such as `audikw_1` and `pwtk`. We also observe that for highly sparse problems, such as `ecology_2` and `parabolic_fem`, PSe1Inv is relatively more costly, but the scalability of PSe1Inv can still be better than that of SuperLU_DIST when a large number of cores are used.

The third experiment focuses on the memory cost of PSe1Inv. Figure 10 shows the cumulative total memory cost and cumulative maximum single-core memory usage in the parallel selected inversion procedure for two matrices: `audikw_1` and `DNA_715_64cell`. This is measured by the memory high watermark reached after the LU factorization step and the selected inversion step, respectively. The memory cost for reading the input matrix is much smaller compared to the memory cost of LU factorization or selected inversion. The shaded areas in Figure 10 correspond to the additional memory required to perform each step of the procedure. The memory high watermark implies that the memory cost for LU factorization includes not only the memory cost for the LU factors but also other temporary memory allocation created during the LU factorization process. The situation is similar for the selected inversion process. The memory

555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575

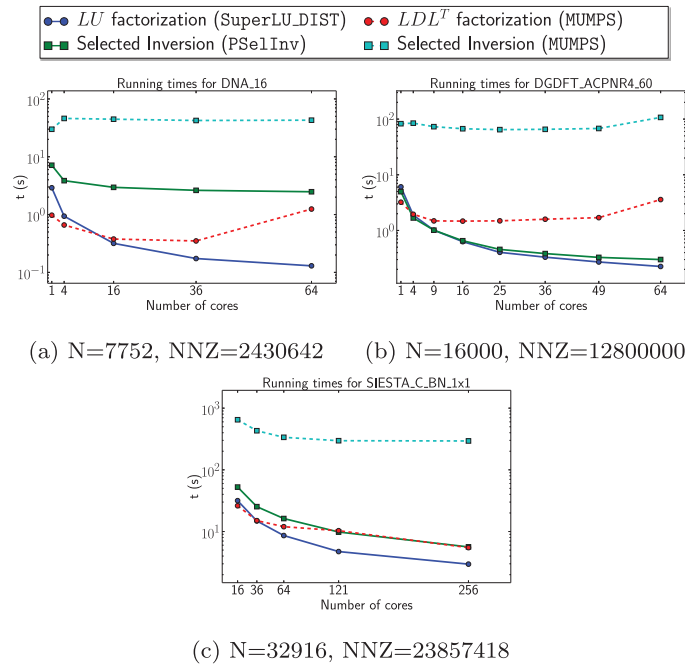


Fig. 11. Performance comparison against MUMPS 4.10.0 for computing selected elements of the inverse.

576 required to store both the LU factors and the corresponding elements in the inverse is
 577 lower than the overall memory high watermark, which accounts for additional memory
 578 required to hold communication and temporary buffers. We can see from Figure 10
 579 that most of the memory allocation is done during the LU factorization step. The total
 580 memory cost of LU factorization and selected inversion increases as the number
 581 of cores increases due to the use of additional buffer arrays for communication and
 582 computation.

583 The total additional memory cost of PSe1Inv is 20% to 60% of the total memory
 584 required by LU factorization, which is relatively small. The maximum memory usage
 585 per core decreases steadily as the number of cores increases. For our test problems,
 586 the maximum memory cost per core of PSe1Inv is around 1GB when a relatively small
 587 number (64) of cores are used. It decreases to around 100MB when a large number
 588 (4,096) of cores are used for the same problem.

589 The last set of experiments compare the performance of PSe1Inv to the parallel
 590 matrix inversion method recently implemented in the MUMPS package [Amestoy et al.
 591 Q6 2012a]. In the MUMPS algorithm, the actual set of computed entries of A^{-1} is also
 592 a superset of the requested entries of A^{-1} . This method can be more efficient than
 593 PSe1Inv when a small number of entries of A^{-1} are requested. However, when a relatively
 594 large number of entries are to be computed, such as in the computation of the
 595 selected elements defined by (1), PSe1Inv can more efficiently reuse the information
 596 shared among different entries of A^{-1} . Figure 11 shows the numerical factorization
 597 and selected inversion timings for MUMPS and PSe1Inv, respectively. As all matrices
 598 are symmetric, MUMPS performs LDL^T factorizations, which use fewer floating point
 599 operations and less memory. We use MUMPS to compute the selected elements of A^{-1} as
 600 defined in Equation (1) for DNA_16 and SIESTA_C_BN_1x1, and we also use MUMPS to
 601 only compute the diagonal entries of A^{-1} for DGDFT_ACPNR4_60. All three matrices

Table III. Comparison of Wall Clock Time Required by PSe1Inv and MUMPS for Computing the Selected Elements Defined by (1) of Two Matrices DNA_16 and SIESTA_C_BN_1x1, and the Diagonal Elements of DGDFt_ACPNR4_60

Problem	Cores (#)	PSe1Inv	MUMPS
DNA_16	1	3.1	29.9
DNA_16	16	2.9	44.6
DGDFt_ACPNR4_60	1	4.9	82.6
DGDFt_ACPNR4_60	16	0.6	67.1
SIESTA_C_BN_1x1	16	52.6	648.2
SIESTA_C_BN_1x1	256	5.6	293.7

come from electronic structure applications, and the difference in terms of computed entries is determined by the different requirements in practical calculations. MUMPS contains a block size parameter that controls the number of right-hand sides processed simultaneously. We experimented with this parameter and report the best results that we could produce for each case (i.e., 16 for DNA_16 and 256 for SIESTA_C_BN_1x1 and DGDFt_ACPNR4_60). Table III shows how the wall clock time used by PSe1Inv compares to that used by MUMPS for three test problems: DNA_16, DGDFt_ACPNR4_60, and SIESTA_C_BN_1x1.

For all test problems, PSe1Inv is at least an order of magnitude faster than MUMPS. It is two orders of magnitude faster on SIESTA_C_BN_1x1 when using 256 cores. PSe1Inv can also scale to a relatively larger number of cores. Our numerical results indicate that for the computation of selected elements as considered in this article, PSe1Inv is more efficient than the more general approach taken in MUMPS.

One way to understand the speedup of PSe1Inv over MUMPS is through the following idealized situation. Consider the computation of the diagonal entries of a tridiagonal matrix A of size N , which can be computed by solving an N set of triangular equations of the form (5). Even with the help of the elimination tree and the fact that only one entry of A^{-1} is needed for each equation, the cost for solving all N equations independently would be $\mathcal{O}(N^2)$. The reason for the high computational cost is that a significant amount of information calculated among the N equations is redundant. The MUMPS approach detects such redundant information on the fly through graph-based algorithms. However, designing an optimal algorithm that maximally reduces the amount of redundant information for a given data-to-processor mapping is difficult. The resulting implementation may or may not reach the optimal complexity. In contrast, PSe1Inv removes all redundant calculation by design, and the computational cost for a tridiagonal matrix can be provably reduced to $\mathcal{O}(N)$ [Lin et al. 2009b].

We should point out that MUMPS can compute an arbitrary set of elements of A^{-1} , whereas the set of selected elements that can be computed by PSe1Inv is more restrictive, as defined in Equation (1).

Overall, the strong scalability of PSe1Inv is similar to that of SuperLU_DIST, and it clearly outperforms the current inversion algorithm as implemented in MUMPS. It requires a modest amount of additional memory to compute the selected elements of the inverse. This observation demonstrates both the validity of our approach and the efficiency of our implementation. More importantly, on matrices arising from actual electronic calculations, PSe1Inv allows one to use thousands of cores, thereby enabling very large scale computations.

4.2. The Accuracy of PSe1Inv

In exact arithmetic, the selected inversion method is an exact method for computing the selected elements of A^{-1} , regardless of whether A is positive definite or not. In practice, the selected inversion method cannot give an exact result due to the presence of

Table IV. Maximum Column-Wise Relative Error between MUMPS and PSe1Inv Using Different Numbers of Cores P

P	Max. Relative Error	P	Max. Relative Error
1	8.52E-008	16	1.96E-007
4	2.35E-008	36	1.72E-007
16	2.52E-008	64	1.77E-007
36	2.05E-007	121	1.76E-007
64	2.07E-007	256	1.78E-007

(a) DNA_16

(b) SIESTA.C.BN_1x1

642 round-off errors. For the sparse direct solver, dynamical pivoting strategies such as the
 643 Bunch-Kaufman process [Bunch and Kaufman 1977] for LDL^T factorization has been
 644 shown to be effective for reducing the numerical error, especially for indefinite matrices.
 645 On distributed memory machines, dynamic pivoting strategies can significantly
 646 affect the load balance and the scalability of the factorization process. Thus, they are
 647 not used in SuperLU_DIST [Li and Demmel 2003].

648 Since the primary goal of the current implementation of the PSe1Inv method is to
 649 achieve high parallel scalability, we choose not to perform additional pivoting steps
 650 after the factorization step. We show that for the test problems we tried, PSe1Inv
 651 produces results comparable to that produced by the MUMPS package and is sufficiently
 652 accurate even when the matrix is relatively ill conditioned.

653 The first set of experiments report comparative results between MUMPS and
 654 PSe1Inv on two problems benchmarked in Section 4.1. Selected elements are com-
 655 puted with both methods, and their maximum relative column-wise difference
 656 $\max_{1 \leq j \leq n} (\|A_{\text{MUMPS},j}^{-1} - A_{\text{PSe1Inv},j}^{-1}\| / \|A_{\text{MUMPS},j}^{-1}\|)$ is presented in Table IV. In both cases,
 657 PSe1Inv provides accurate results that are comparable to MUMPS on these two problems.

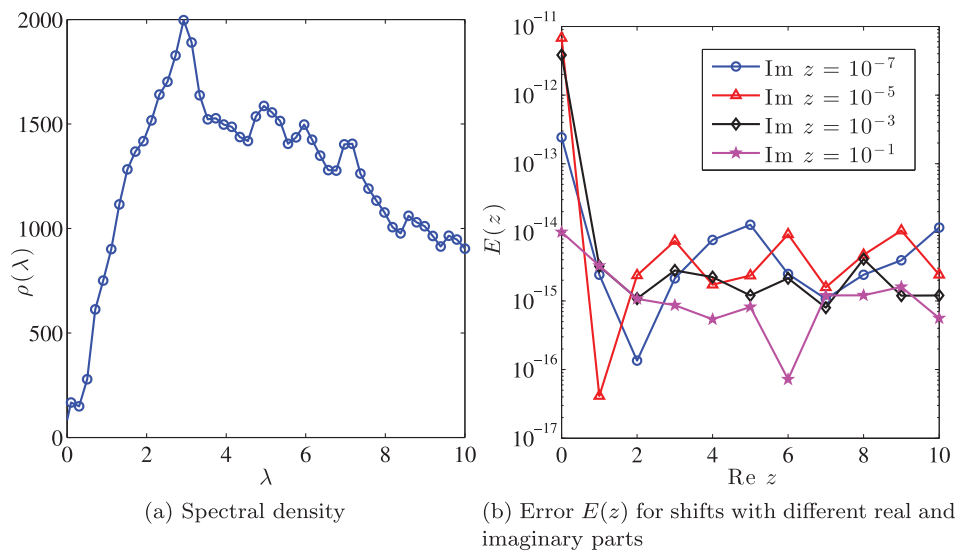
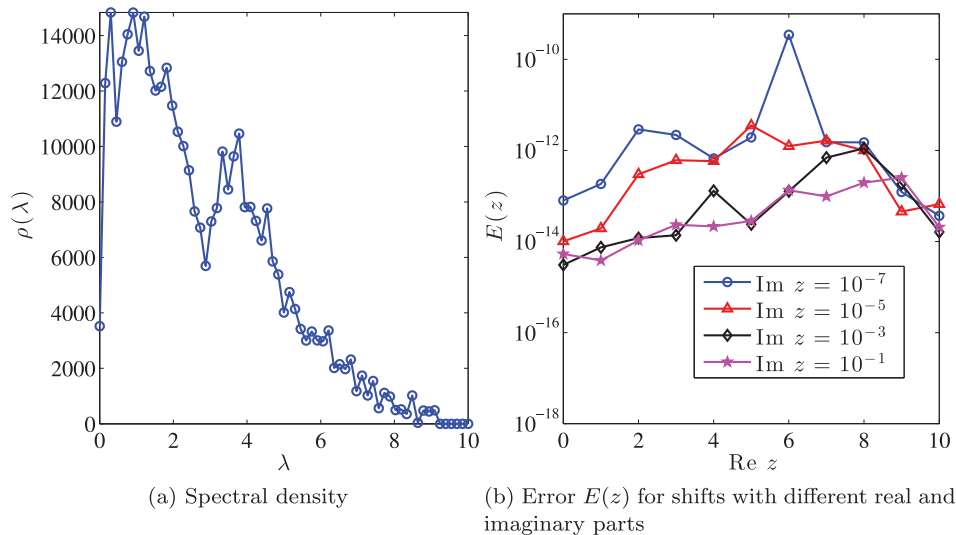
658 In the second set of experiments, we assess the accuracy of PSe1Inv on larger ill-
 659 conditioned matrices. To obtain test problems that are indefinite and ill conditioned,
 660 for each test problem listed in Table II, we construct a sequence of $A(z)$ defined by
 661 Equation (18) for several complex shifts z . The real parts of the shifts lie within the
 662 spectrum of the matrix pencil (H, S) , and the imaginary parts range from small (10^{-7})
 663 to large (10^{-1}) values.

664 To quantify the accuracy of the PSe1Inv method for large matrices in the test set,
 665 we need to find an appropriate error metric. Due to the large matrix size, the parallel
 666 matrix inversion method in MUMPS becomes too expensive to be used for the purpose of
 667 benchmarking. Motivated by the trace estimation in Equation (2), we choose to measure
 668 the numerical error introduced by PSe1Inv by using the following quantity:

$$E(z) = \frac{|N - \text{Tr}[A(z)^{-1}A(z)]|}{N} \equiv \left| 1 - \frac{1}{N} \sum_{i,j=1}^N [A(z)^{-1}]_{ij}[A(z)]_{ji} \right|. \quad (19)$$

669 Since A is sparse, the summation in Equation (19) involves only those i and j such that
 670 $A_{i,j} \neq 0$.

671 In Figure 12, we show both the spectral density $\rho(\lambda)$ of the pwtk matrix, which
 672 describes the number of eigenvalues of eigenvalues (H, S) per unit interval, and $E(z)$
 673 for several shifts z with different real and imaginary parts. We observe that for all of
 674 these problems, the measured errors are below 10^{-11} , even when the real part of z is
 675 close to an eigenvalue cluster and the imaginary part of z is as small as 10^{-7} . Figure 13
 676 shows that a similar level of accuracy is achieved in the test of the DNA_715_64cell
 677 matrix.

Fig. 12. Spectral density and error $E(z)$ for the pwtk matrix.Fig. 13. Spectral density and error $E(z)$ for the DNA_715_64cell matrix.

4.3. Application to Electronic Structure Theory

In this section, we demonstrate how PSe1Inv can be applied to accelerate Kohn-Sham density functional theory (KSDFT) calculation [Hohenberg and Kohn 1964; Kohn and Sham 1965], which is widely used for describing the ground state electronic properties of molecules, solids, and other nanostructures. We use the recently developed pole expansion and selected inversion technique (PEXSI) [Lin et al. 2009a, 2009b, 2011b, 2013] to compute the nonzero elements of the so-called single particle density matrix

678
679
680
681
682
683
684

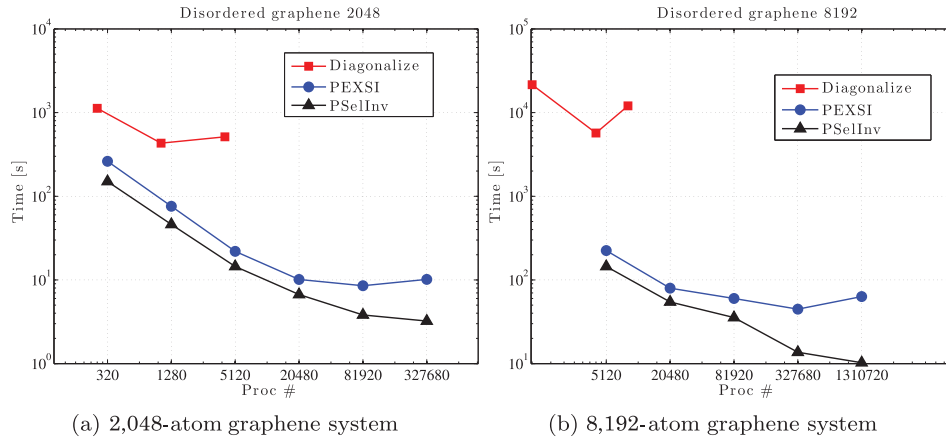


Fig. 14. Wall clock time versus the number of cores for a graphene system.

685 Γ that can be approximated by

$$\Gamma \approx \sum_{l=1}^P \omega_l (H - z_l S)^{-1}, \quad (20)$$

686 where $z_l, \omega_l \in \mathbb{C}$, and the number of “poles” P is around 80 in practical calculations.
 687 Both H and S are real sparse symmetric matrices that have the same sparsity pattern
 688 when a local basis set is used to discretize the Kohn-Sham problem. Each matrix
 689 $A_l = H - z_l S$ is a complex symmetric matrix. In KSDFT calculations, only elements of
 690 Γ_{ij} corresponding to the nonzero elements of H and S (i.e., $H_{ij}, S_{ij} \neq 0$) are needed to
 691 compute physical quantities such as electron density and energy. The expansion (20)
 692 immediately suggests that only the selected elements as defined in Equation (1)
 693 of the matrices A_l^{-1} are needed. In a parallel implementation of PEXSI, we use PSe1Inv
 694 to evaluate the selected elements of A_l^{-1} that correspond to the nonzero elements of
 695 Γ on a subset of cores. The selected inversion of A_l for different l can be carried out
 696 independently on different subsets of cores.

697 We apply the parallel PEXSI method to the DG_Graphene_2048 and
 698 DG_Graphene_8192 systems, which are disordered graphene systems with 2,048 and
 699 8,192 atoms, respectively, and compare its performance to a standard approach that
 700 requires a partial diagonalization of (H, S) . We use a ScaLAPACK subroutine pdsyevr
 701 [Vömel 2010], which is based on the multiple relatively robust representations (MRRR)
 702 algorithm, to perform such diagonalization. Although both H and S are sparse matrices,
 703 the MRRR algorithm treats them as dense matrices. For $H, S \in \mathbb{R}^{N \times N}$, the MRRR
 704 algorithm first performs a tridiagonalization procedure with $\mathcal{O}(N^3)$ cost, then efficiently
 705 solves the eigenvalues and eigenvectors of the tridiagonal system with $\mathcal{O}(N^2)$ cost, and
 706 finally the eigenvectors can be constructed with $\mathcal{O}(N^3)$ cost.

707 Figure 14(a) shows that the scalability of pdsyevr is limited to 1,024 cores for the
 708 2,048-atom problem. Adding more cores to the diagonalization process leads to an
 709 increase of the wall clock time due to communication overhead. For this relatively
 710 small problem, the benefit of parallel implementation of PEXSI is already clear when
 711 320 cores are used. Since we use $P = 80$ poles in the pole expansion, 4 cores with a
 712 2×2 processor grid are used in each selected inversion in this case. The wall clock time
 713 used by PEXSI is 261 seconds, among which 150 seconds are attributed to PSe1Inv,
 714 95 seconds are attributed to factorization, and 10 seconds are attributed to symbolic

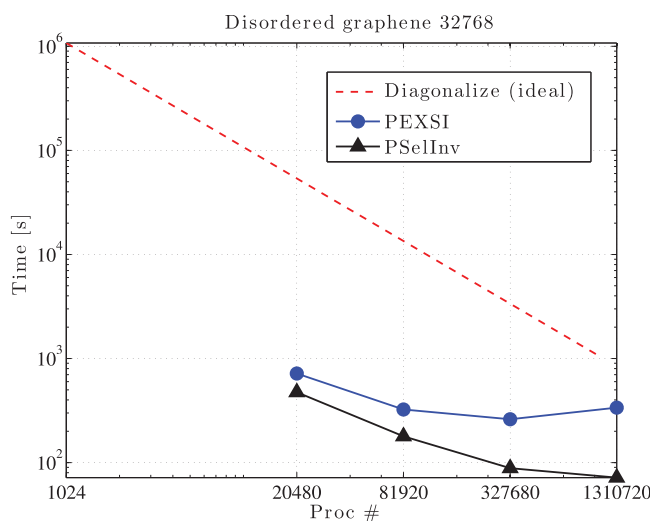


Fig. 15. Wall clock time versus the number of cores for a graphene system with 32,768 atoms.

factorization. This timing result compares favorably to the 430 seconds of measured wall clock time required by `pdsyevr` on 1,024 cores. 715

Furthermore, we can clearly see from Figure 14(a) that the parallel PEXSI method can scale to a much larger number of cores. Nearly perfect speedup can be observed when a total of 20,480 cores is used in the parallel PEXSI computation, with 256 cores used in each selected inversion. The total wall clock time used in this calculation is merely 10 seconds. Compared to the best wall clock time that we can obtain for the diagonalization procedure, which is 430 seconds on 1,024 cores, this represents a speedup factor of 43. 716-723

For the larger system that contains 8,192 atoms, `pdsyevr` can scale to 4,096 cores, as we can see in Figure 14(b). It takes 5,703 wall clock seconds to perform such a computation. When the parallel PEXSI calculation is carried out on 5,120 cores with 64 cores used to perform each selected inversion, the total wall clock time required is 224 seconds. Figure 14(b) also shows that parallel PEXSI can scale to as many as 327,680 cores with 4,096 cores used for each selected inversion. The total wall clock time required in this calculation is merely 45 seconds, a 127-fold speedup compared to the best diagonalization wall clock time measured. 724-731

Finally, we apply PEXSI to a system with 32,768 atoms. The matrix size is four times larger than that for the system with 8,192 atoms, and the diagonalization routine is no longer feasible: the wall clock time required to run the `pdsyevr` routine with 1,024 cores for the 2,048-atom system is 431 seconds, and for the 8,192-atom system, it is 21,556 seconds. The increase of the wall clock time is 50 fold, which is roughly in agreement with the cubic complexity scaling factor $4^3 = 64$. The cubic scaling of the diagonalization procedure implies that the wall clock time would increase by at least a factor of 50 to 1,077,800 seconds (300 hours) if we perform the same type of calculation for a system that contains 32,768 atoms on 1,024 cores. Based on this estimation and assuming that the strong scaling of `pdsyevr` is perfect, we compare the ideal performance of the diagonalization method to the practical performance of PEXSI in Figure 15 up to 1,310,720 cores. The total wall clock time for both factorization and selected inversion reaches its minimum at 4,096 cores per pole (327,680 cores in total), which is 241 seconds. Among these, 87 seconds is attributed to PSe1Inv. Comparatively, even if the diagonalization procedure scales perfectly to more than one million cores, 732-746

747 which is highly unlikely within the current framework of diagonalization methods, the
 748 projected wall clock time is greater than 1,000 seconds, which is significantly more
 749 than that used by PEXSI.

750 5. CONCLUSIONS AND FUTURE WORK

751 We described an efficient parallel implementation of the selected inversion algorithm
 752 for distributed memory parallel machines. The current implementation of PSe1Inv can
 753 be applied to the computation of selected elements of the inverse of a sparse symmetric
 754 matrix. It is publicly available and can scale to more than 4,000 cores for sufficiently
 755 large problems. The scalability of the solver depends on the size and sparsity of the
 756 matrix. We observed that it is important to exploit concurrency available within the
 757 elimination trees to achieve high scalability in the parallel selected inversion process.
 758 In the future, we plan to further improve the tree-level parallelism to enhance the con-
 759 currency among different supernodes. We also observed that for our test problems, the
 760 PSe1Inv method is relatively accurate even for matrices that are highly indefinite and
 761 close to singular. It can be applied to accelerate several scientific computation applica-
 762 tions, such as the DFT-based electronic structure calculations. To further improve the
 763 numerical accuracy of the PSe1Inv method, especially for indefinite matrices, dynamic
 764 pivoting strategies such as the Bunch-Kaufman procedure [Bunch and Kaufman 1977;
 765 Grimes et al. 1994] for the factorization may be needed. Generalizing PSe1Inv to non-
 766 symmetric matrices and combining it with other sparse direct solvers are other areas
 767 in which we plan to work in the future.

768 ACKNOWLEDGMENTS

769 We would like to thank Xiaoye S. Li and François-Henry Rouet for helpful discussions.

770 REFERENCES

- 771 P. Amestoy, I. Duff, J.-Y. L'Excellent, and J. Koster. 2001. A fully asynchronous multifrontal solver using
 772 distributed dynamic scheduling. *SIAM Journal on Matrix Analysis and Applications* 23, 15–41.
- 773 P. R. Amestoy, I. S. Duff, J.-Y. L'Excellent, Y. Robert, F.-H. Rouet, and B. Uçar. 2012b. On computing inverse
 774 entries of a sparse matrix in an out-of-core environment. *SIAM Journal on Scientific Computing* 34,
 775 A1975–A1999.
- 776 P. R. Amestoy, I. S. Duff, J. Y. L'Excellent, and F. H. Rouet. 2012a. *Parallel Computation of Entries of A^{-1}* .
 777 Technical Report. CERFACS, Toulouse, France.
- 778 C. Ashcraft and R. Grimes. 1989. The influence of relaxed supernode partitions on the multifrontal method.
 779 *ACM Transactions on Mathematical Software* 15, 291–309.
- 780 C. Bekas, A. Curioni, and I. Fedulova. 2009. Low cost high performance uncertainty quantification. In
 781 *Proceedings of the 2nd Workshop on High Performance Computational Finance*. 8.
- 782 C. Bekas, E. Kokiopoulou, and Y. Saad. 2007. An estimator for the diagonal of a matrix. *Applied Numerical*
 783 *Mathematics* 57, 1214–1229.
- 784 L. S. Blackford. 1997. *ScaLAPACK User's Guide*. Vol. 4. SIAM.
- 785 J. R. Bunch and L. Kaufman. 1977. Some stable methods for calculating inertia and solving symmetric linear
 786 systems. *Mathematics of Computation* 31, 137, 163–179.
- 787 Y. E. Campbell and T. A. Davis. 1995. *Computing the Sparse Inverse Subset: An Inverse Multifrontal Approach*.
 788 Technical Report TR-95-021. University of Florida.
- 789 S. Cauley, V. Balakrishnan, G. Klimeck, and C.-K. Koh. 2012. A two-dimensional domain decomposition
 790 technique for the simulation of quantum-scale devices. *Journal of Computational Physics* 231, 4, 1293–
 791 1313.
- 792 C. Chevalier and F. Pellegrini. 2008. PT-scotch: A tool for efficient parallel graph ordering. *Parallel Computing*
 793 34, 318–331.
- 794 T. A. Davis and Y. Hu. 2011. The University of Florida sparse matrix collection. *ACM Transactions on*
 795 *Mathematical Software* 38, 1.
- 796 I. Duff, R. Grimes, and J. Lewis. 1992. *User's Guide for the Harwell-Boeing Sparse Matrix Collection*. Research
 797 and Technology Division, Boeing Computer Services, Seattle, WA.

PSELInv: The Symmetric Case

21:27

- S. Eastwood and J. Wan. 2013. Finding off-diagonal entries of the inverse of a large symmetric sparse matrix. *Numerical Linear Algebra with Applications* 20, 1, 74–92. 798
799
- A. Erisman and W. Tinney. 1975. On computing certain elements of the inverse of a sparse matrix. *Communications of the ACM* 18, 177. 800
801
- G. H. Golub and C. F. Van Loan. 1996. *Matrix Computations* (3rd ed.). Johns Hopkins University Press, Baltimore, MD. 802
803
- R. G. Grimes, J. G. Lewis, and H. D. Simon. 1994. A shifted block Lanczos algorithm for solving sparse symmetric generalized eigenproblems. *SIAM Journal on Matrix Analysis and Applications* 15, 1, 228–272. 804
805
806
- U. Hetmaniuk, Y. Zhao, and M. P. Anantram. 2013. A nested dissection approach to modeling transport in nanodevices: Algorithms and applications. *International Journal for Numerical Methods in Engineering* 95, 7, 587–607. 807
808
809
- P. Hohenberg and W. Kohn. 1964. Inhomogeneous electron gas. *Physical Review* 136, B864–B871. 810
- M. Jacquelin, L. Lin, and C. Yang. 2014. *PSELInv—A Distributed Memory Parallel Algorithm for Selected Inversion: The Symmetric Case*. Technical Report, Lawrence Berkeley National Laboratory, Berkeley, CA. 811
812
813
- G. Karypis and V. Kumar. 1998. A parallel algorithm for multilevel graph partitioning and sparse matrix ordering. *Journal of Parallel and Distributed Computing* 48, 71–85. 814
815
- W. Kohn and L. Sham. 1965. Self-consistent equations including exchange and correlation effects. *Physical Review* 140, A1133–A1138. 816
817
- G. Kotliar, S. Y. Savrasov, K. Haule, V. S. Oudovenko, O. Parcollet, and C. A. Marianetti. 2006. Electronic structure calculations with dynamical mean-field theory. *Review of Modern Physics* 78, 865–952. 818
819
- A. Kuzmin, M. Luisier, and O. Schenk. 2013. Fast methods for computing selected elements of the Green’s function in massively parallel nanoelectronic device simulations. In *Euro-Par 2013 Parallel Processing*. Lecture Notes in Computer Science, Vol. 8097. Springer, 533–544. 820
821
822
- C. Lanczos. 1950. An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *Journal of Research of the National Bureau of Standards* 45, 255–282. 823
824
- S. Li, S. Ahmed, G. Klimeck, and E. Darve. 2008. Computing entries of the inverse of a sparse matrix using the FIND algorithm. *Journal of Computational Physics* 227, 9408–9427. 825
826
- Song Li and Eric Darve. 2012. Extension and optimization of the FIND algorithm: Computing Green’s and less-than Green’s functions. *Journal of Computational Physics* 231, 4, 1121–1139. 827
828
- S. Li, W. Wu, and E. Darve. 2013. A fast algorithm for sparse matrix computations related to inversion. *Journal of Computational Physics* 242, 915–945. 829
830
- X. S. Li and J. W. Demmel. 2003. SuperLU_DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems. *ACM Transactions on Mathematical Software* 29, 110. 831
832
- L. Lin, M. Chen, C. Yang, and L. He. 2013. Accelerating atomic orbital-based electronic structure calculation via pole expansion and selected inversion. *Journal of Physics: Condensed Matter* 25, 295501. 833
834
- L. Lin, J. Lu, L. Ying, R. Car, and W. E. 2009b. Fast algorithm for extracting the diagonal of the inverse matrix with application to the electronic structure analysis of metallic systems. *Communications in Mathematical Sciences* 7, 755. 835
836
837
- L. Lin, J. Lu, L. Ying, and E. Weinan. 2009a. Pole-based approximation of the Fermi-Dirac function. *Chinese Annals of Mathematics, Series B* 30, 729. 838
839
- L. Lin, J. Lu, L. Ying, and E. Weinan. 2012. Adaptive local basis set for Kohn-Sham density functional theory in a discontinuous Galerkin framework I: Total energy calculation. *Journal of Computational Physics* 231, 2140–2154. 840
841
842
- L. Lin, C. Yang, J. Lu, L. Ying, and E. Weinan. 2011a. A fast parallel algorithm for selected inversion of structured sparse matrices with application to 2D electronic structure calculations. *SIAM Journal on Scientific Computing* 33, 1329. 843
844
845
- L. Lin, C. Yang, J. Meza, J. Lu, L. Ying, and E. Weinan. 2011b. SelInv—an algorithm for selected inversion of a sparse symmetric matrix. *ACM Transactions on Mathematical Software*, 40. 846
847
- J. Liu. 1990. The role of elimination trees in sparse factorization. *SIAM Journal on Matrix Analysis and Applications* 11, 134. 848
849
- D. E. Petersen, S. Li, K. Stokbro, H. H. B. Sørensen, P. C. Hansen, S. Skelboe, and E. Darve. 2009. A hybrid method for the parallel computation of Green’s functions. *Journal of Computational Physics* 228, 5020–5039. 850
851
852
- E. Rothberg and A. Gupta. 1994. An efficient block-oriented approach to parallel sparse Cholesky factorization. *SIAM Journal on Scientific Computing* 15, 1413–1439. 853
854

- 855 O. Schenk and K. Gartner. 2006. On fast factorization pivoting methods for symmetric indefinite systems.
856 *Electronic Transactions on Numerical Analysis* 23, 158–179.
- 857 R. B. Sidje and Y. Saad. 2011. Rational approximation to the Fermi-Dirac function with applications in
858 density functional theory. *Numerical Algorithms* 56, 455.
- 859 J. M. Soler, E. Artacho, J. D. Gale, A. García, J. Junquera, P. Ordejón, and D. Sánchez-Portal. 2002. The
860 SIESTA method for ab initio order-n materials simulation. *Journal of Physics: Condensed Matter* 14,
861 2745–2779.
- 862 K. Takahashi, J. Fagan, and M. Chin. 1973. Formation of a sparse bus impedance matrix and its application
863 to short circuit study. In *Proceedings of the 8th PICA Conference*.
- 864 J. M. Tang and Y. Saad. 2012. A probing method for computing the diagonal of a matrix inverse. *Numerical*
865 *Linear Algebra with Applications* 19, 485–501.
- 866 C. Vömel. 2010. ScaLAPACK’s MRRR algorithm. *ACM Transactions on Mathematical Software* 37, 1.
- 867 Z. Xu and A. C. Maggs. 2013. Solving fluctuation-enhanced Poisson-Boltzmann equations. arXiv:1310.4682.

Received April 2014; revised May 2015; accepted May 2015

QUERIES

- Q1:** AU: Please provide full mailing and email addresses for all authors per ACM style.
- Q2:** AU: Throughout, please note that “LU” is sometimes italicized in text (same with “LDL”). Streamline style?
- Q3:** AU: The URL in footnote 1 redirects. Update?
- Q4:** AU: Please confirm position of closing parenthesis in the Output section of Algorithm 1.
- Q5:** AU: Add a comma before “gives,” or delete comma after “lists”?
- Q6:** AU: Streamline font of terms such as “MUMPS”?