

A PARALLEL ALGORITHM FOR COMPUTING PARTIAL SPECTRAL FACTORIZATIONS OF MATRIX PENCILS VIA CHEBYSHEV APPROXIMATION*

TIANSI XU[†], ANTHONY P. AUSTIN[‡], VASSILIS KALANTZIS[§], AND YOUSEF SAAD[¶]

Abstract. We propose a distributed-memory parallel algorithm for computing all the eigenvalues (and corresponding eigenvectors) of a large, sparse, real symmetric positive definite matrix pencil that lie within a target interval. The algorithm is based on Chebyshev interpolation of the eigenvalues of the Schur complement (over the interface variables) of a domain decomposition reordering of the pencil and accordingly exposes two dimensions of parallelism: one derived from the reordering and one from the independence of the interpolation nodes. The new method demonstrates excellent parallel scalability, comparing favorably with **PARPACK**, and does not require factorization of the mass matrix, which significantly reduces memory consumption, especially for 3D problems. Our implementation is publicly available on GitHub.

Key word. Symmetric generalized eigenvalue problem, spectral Schur complements, Chebyshev approximation, parallel computing

AMS subject classifications. 15A18, 65D15, 65F15, 65N55, 65Y05, 68W10

1. Introduction. Several applications in science and engineering require the computation of a handful of the algebraically smallest eigenvalues and associated eigenvectors of a large, sparse matrix pencil (A, M) , where the $n \times n$ matrices A and M are real symmetric and M is positive-definite. Often, one is provided bounds α and β on the eigenvalues of interest, and the goal is then to compute all n_{ev} eigenpairs of (A, M) that lie within $[\alpha, \beta]$. Problems of this sort arise, for instance, in spectral clustering [41], and low-frequency response analysis [6, 15].

Due to the size of modern matrix problems, parallel computing has become an integral part of software libraries targeting large-scale eigenvalue computations. In many packages (e.g., **PARPACK** [30, 34], **PRIMME** [37], **BLOPEX** [28]), linear algebra kernels are the main source of parallelism, with operations such as matrix-vector and dot products performed in parallel by distributing the data across multiple processors. Several recent packages improve scalability by exploiting additional levels of parallelism via techniques such as spectrum slicing (**pEVSL** [31]), rational filtering (**FEAST/PFEAST** [20, 27, 35], and **z-Pares** [36]), and parallel shift-and-invert methods [42, 46]. In addition, the **SLEPc** collection of distributed-memory eigenvalue algorithms [14] contains implementations of several of these methods.

Another class of distributed-memory eigenvalue solvers is based on algebraic domain decomposition, also known as algebraic substructuring; see the references in

*Submitted to the editors June 5, 2022.

Funding: The work of the first and the last author was supported by the National Science Foundation (NSF) grant DMS-1912048. The work of the second author was supported by the Research Initiation Program at the Naval Postgraduate School. The work of the third author was supported by the Mathematical Sciences Council of IBM Research through its Exploratory Science initiative.

[†]Dept. of Computer Science and Engineering, University of Minnesota, Minneapolis, MN 55455 (xuxx1180@umn.edu)

[‡]Dept. of Applied Mathematics, Naval Postgraduate School, 833 Dyer Rd., Monterey, CA 93940 (anthony.austin@nps.edu)

[§]IBM Research, Thomas J. Watson Research Center, Yorktown Heights, NY 10598 (vkal@ibm.com)

[¶]Dept. of Computer Science and Engineering, University of Minnesota, Minneapolis, MN 55455 (saad@umn.edu)

[17] for details. In domain decomposition, the adjacency graph associated with the pencil (A, M) is partitioned into several non-overlapping subgraphs. The eigenvalue problem then decouples into two separate tasks: first, one determines the eigenvector components associated with the interface variables of the partitioned graph; then, one finds the components associated with the interior variables. The second task parallelizes naturally over the subgraphs. For more information on this type of solver, see [6, 12, 17, 29, 45].

1.1. A new parallel algorithm. In this article, we combine the domain decomposition approach with Chebyshev function approximation to design a new distributed-memory parallel eigensolver. The contributions of our work are:

1. The algorithm parameterizes the eigenvector components associated with the interior and interface variables as univariate, analytic vector-valued functions. It then uses the fact that Chebyshev interpolation of these functions yields good approximations to the eigenvectors to construct a subspace for use with a Rayleigh–Ritz projection scheme. We present theoretical and practical details when the interpolation points are Chebyshev nodes of the second kind.
2. The proposed algorithm leverages multi-dimensional parallelism by assigning computations associated with different Chebyshev nodes to different processor groups and assigning computations associated with different subdomains to different processors within each group. Our numerical experiments demonstrate that the algorithm achieves higher parallel efficiency than PARPACK on distributed-memory systems communicating via the Message Passing Interface (MPI) [13]. A C++/MPI implementation of the proposed algorithm is available publicly at <https://github.com/Hitenze/Schurcheb>.
3. In contrast to previous work on domain decomposition eigensolvers, the proposed algorithm requires the computation of neither derivatives of eigenvectors [18] nor a large number of eigenvectors of linearized spectral Schur complements [5, 6]. Moreover, unlike branch-hopping domain decomposition algorithms, which compute eigenvalues one at a time [19, 21], the proposed algorithm introduces model parallelism in addition to data parallelism by approximating all sought eigenvalues simultaneously via Rayleigh–Ritz projection. Unlike approaches based on the Lanczos algorithm, the proposed algorithm does not require a distributed-memory factorization of A or M ; therefore, it is not limited by the efficiency of distributed-memory triangular solves. Finally, in contrast to most rational filtering techniques, especially those based on discretizations of complex contour integrals [22, 23], the proposed algorithm does not evaluate functions at complex values and therefore does not require complex arithmetic.

1.2. Notation and roadmap. Throughout the paper, we denote the set of eigenvalues of a general pencil (K, F) by $\Lambda(K, F)$ and the eigenpairs of the specific pencil (A, M) by $(\lambda_i, x^{(i)})$, $i = 1, \dots, n$, ordered algebraically: $\lambda_1 \leq \dots \leq \lambda_n$. Given bounds α and β such that $\alpha < \lambda_1$, our aim is to compute all n_{ev} eigenpairs of (A, M) that lie in $[\alpha, \beta]$, i.e., the n_{ev} algebraically smallest eigenvalues of A and their corresponding eigenvectors. Finally, we denote by $\text{Ran}(K)$ and $\text{Ker}(K)$ the range and kernel of a matrix K and by $\text{span}\{v_1, \dots, v_k\}$ the linear span of vectors v_1, \dots, v_k .

This paper is organized as follows. Section 2 presents background on algebraic graph partitioning and domain decomposition. Section 3 shows how the eigenvectors

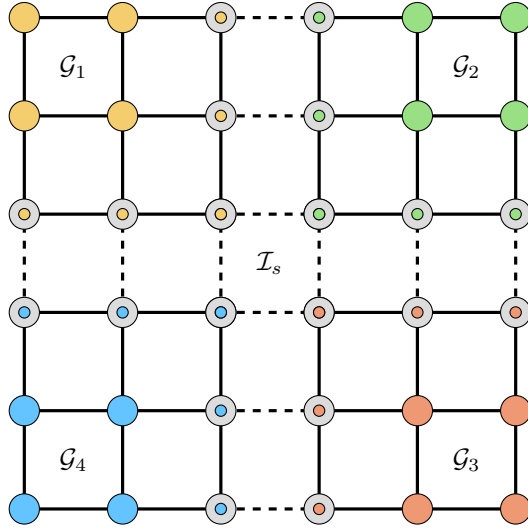


Fig. 2.1: A 4-way partitioning of a 6×6 discretized domain obtained from an edge separator. The four colors distinguish the four different subdomains. Solid-colored nodes correspond to interior variables. Nodes with a gray background correspond to interface variables. Solid lines correspond to edges between vertices of the same partition. Dashed lines correspond to edges between vertices of neighboring partitions

of (A, M) can be identified as values of certain univariate, vector-valued functions and discusses how they can be approximated by Rayleigh–Ritz projection onto a subspace formed via Chebyshev approximation. Section 4 discusses the distributed-memory implementation of the proposed algorithm on 2D grids of MPI processes. Section 5 showcases the performance of the proposed algorithm using numerical experiments performed in both sequential and distributed-memory computing environments. Finally, Section 6 presents our concluding remarks.

2. Domain decomposition variable ordering. Let $\mathcal{G} = (\mathcal{V}, \mathcal{I})$ be a simple undirected graph with vertex set \mathcal{V} and edge set \mathcal{I} . A p -way edge separator is a subset $\mathcal{I}_s \subseteq \mathcal{I}$ whose removal from \mathcal{I} divides the vertices of the graph \mathcal{G} into $p \in \mathbb{N}$ non-overlapping sets $\mathcal{V}_1, \dots, \mathcal{V}_p$, such that the induced subgraphs $\mathcal{G}_1 = (\mathcal{V}_1, \mathcal{I}_1), \dots, \mathcal{G}_p = (\mathcal{V}_p, \mathcal{I}_p)$, are disjoint. We refer to the induced subgraphs variously as *subdomains*, *substructures*, or *partitions*. A vertex is called an *interface vertex* if it is incident to an edge in \mathcal{I}_s and an *interior vertex* otherwise.

Applied to graphs derived from matrices, edge separators are commonly used in parallel computing to achieve load balancing during the execution of distributed-memory linear algebra kernels. In this context, the induced subgraphs ideally have similar numbers of vertices and edges, while the size (cardinality) of the separator set is kept to a minimum. Finding the “best” edge separator is an NP-hard problem. In practice, one relies on heuristics, such as the algebraic partitioning strategies implemented in the popular METIS and ParMETIS packages [24, 25].

To a symmetric matrix pencil (A, M) of dimension n , we associate a graph $\mathcal{G}_{A,M}$ in the usual way, taking $\mathcal{V} = \{1, \dots, n\}$ for the vertex set and $\mathcal{I} = \{(i, j) \mid A_{i,j} \neq 0 \text{ or } M_{i,j} \neq 0\}$ for the edge set. Thinking of the eigenvalue problem $Ax = \lambda Mx$ as a set of n linear equations, the vertices correspond to the n unknown variables

in the vector x , and $\mathcal{G}_{A,M}$ has an edge connecting vertices i and j if both the i th and j th variables appear together in one of the equations. A p -way edge separator for $\mathcal{G}_{A,M}$ groups the variables into p disjoint sets or subdomains. Interface vertices correspond to variables that are coupled (via equations) with variables from multiple subdomains, while interior vertices correspond to variables that are coupled only with other variables from the same subdomain. Figure 2.1 illustrates this for a 4-way partitioning of a graph that models a 6×6 regular grid.

Having partitioned $\mathcal{G}_{A,M}$, we reorder the variables, listing all interior variables first, grouped by in order by subdomain, followed by the interface variables, also grouped by subdomain. Let P be the permutation matrix that effects this reordering. Under P , the matrices A and M are reordered into a pair of structured block matrices:

$$\begin{aligned}
 P^T A P &= \begin{bmatrix} B_1 & & & E_1 & & \\ & B_2 & & & E_2 & \\ & & \ddots & & & \ddots \\ E_1^T & & & B_p & & E_p \\ & E_2^T & & & C_{1,1} & C_{1,2} & \cdots & C_{1,p} \\ & & \ddots & & C_{2,1} & C_{2,2} & \cdots & C_{2,p} \\ & & & E_p^T & \vdots & \vdots & \ddots & \vdots \\ & & & & C_{p,1} & C_{p,2} & \cdots & C_{p,p} \end{bmatrix} \\
 P^T M P &= \begin{bmatrix} M_{B_1} & & & M_{E_1} & & \\ & M_{B_2} & & & M_{E_2} & \\ & & \ddots & & & \ddots \\ M_{E_1}^T & & & M_{B_p} & & M_{E_p} \\ & M_{E_2}^T & & & M_{C_{1,1}} & M_{C_{1,2}} & \cdots & M_{C_{1,p}} \\ & & \ddots & & M_{C_{2,1}} & M_{C_{2,2}} & \cdots & M_{C_{2,p}} \\ & & & M_{E_p}^T & \vdots & \vdots & \ddots & \vdots \\ & & & & M_{C_{p,1}} & M_{C_{p,2}} & \cdots & M_{C_{p,p}} \end{bmatrix}.
 \end{aligned}
 \tag{2.1}$$

To provide more detail, let d_i and s_i denote, respectively, the numbers of interior and interface variables belonging to the i th domain. The matrices B_i and M_{B_i} are of size $d_i \times d_i$ and represent the coupling between the interior variables within the i th subdomain. The matrices E_i and M_{E_i} are of size $d_i \times s_i$ and represent the coupling between the interior and interface variables of the i th subdomain. Finally, the matrices $C_{i,j}$ and $M_{C_{i,j}}$ are of size $s_i \times s_j$ and represent the coupling between the interface variables of the i th subdomain and those of the j th subdomain. If the i th and j th subdomains do not neighbor one another, $C_{i,j} = M_{C_{i,j}} = 0$. Since A and M are symmetric, $C_{j,i} = C_{i,j}^T$ and $M_{C_{j,i}} = M_{C_{i,j}}^T$.

Our algorithm makes essential use of the structure of this reordering of A and M . For the remainder of the paper, we assume that A and M have been so reordered and suppress mention of the permutation P . We write A and M in 2×2 block form as

$$A = \begin{bmatrix} B & E \\ E^T & C \end{bmatrix}, \quad M = \begin{bmatrix} M_B & M_E \\ M_E^T & M_C \end{bmatrix},
 \tag{2.2}$$

with the blocks being defined in the obvious way to conform to the structure just described. Finally, we define $d = d_1 + \cdots + d_p$ and $s = s_1 + \cdots + s_p$, the total numbers of interior and interface variables, respectively. Thus, the matrices B and M_B are $d \times d$, E and M_E are $d \times s$, and C and M_C are $s \times s$. Of course, $d + s = n$.

3. A parallel algorithm based on Chebyshev approximation. Our algorithm is based on the fact that the eigenvalues and eigenvectors of the matrix $A - \zeta M$ are analytic functions of $\zeta \in \mathbb{C}$ (vector-valued in the case of the latter). By definition, if $\zeta = \lambda_i$ is an eigenvalue of the pencil (A, M) , then $A - \zeta M$ is singular, and its null vectors are the eigenvectors for (A, M) corresponding to λ_i . By continuity, if ζ is close (but not equal) to λ_i , then $A - \zeta M$ will be “nearly singular” in the sense that it will have one or more eigenvalues that are small in magnitude, and the eigenvectors of $A - \zeta M$ corresponding to these eigenvalues will be good approximations to null vectors of $A - \lambda_i M$. On this basis, our algorithm computes the eigenvectors corresponding to the smallest eigenvalues of $A - \zeta_i M$ at several points ζ_i within the search interval $[\alpha, \beta]$. By choosing the ζ_i well, we can guarantee that the subspace spanned by these “near-null” vectors contains good approximations to the eigenvectors of (A, M) . The algorithm extracts such approximations from this subspace via Rayleigh–Ritz projection.

3.1. Spectral Schur complements. To make this process efficient and parallelizable, we exploit the block structure of A and M induced by the variable reordering discussed in the previous section. Partition the eigenvector $x^{(i)}$ associated with the eigenvalue λ_i of (A, M) as

$$x^{(i)} = \begin{bmatrix} u^{(i)} \\ y^{(i)} \end{bmatrix},$$

where $u^{(i)} \in \mathbb{R}^d$ and $y^{(i)} \in \mathbb{R}^s$, conforming to the partitioning of A and M in (2.2), and define

$$(3.1) \quad B(\zeta) = B - \zeta M_B, \quad E(\zeta) = E - \zeta M_E, \quad C(\zeta) = C - \zeta M_C,$$

for $\zeta \in \mathbb{C}$. In this notation, the eigenvector equation $(A - \lambda_i M)x^{(i)} = 0$ becomes

$$(3.2) \quad \begin{bmatrix} B(\lambda_i) & E(\lambda_i) \\ E^T(\lambda_i) & C(\lambda_i) \end{bmatrix} \begin{bmatrix} u^{(i)} \\ y^{(i)} \end{bmatrix} = 0.$$

Under the mild assumption that $B(\lambda_i)$ is invertible, i.e., that $\lambda_i \notin \Lambda(B, M_B)$, we can eliminate the $E^T(\lambda_i)$ block in the second row, yielding

$$(3.3) \quad [C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^{-1}E(\lambda_i)]y^{(i)} = 0.$$

That is, the $s \times 1$ bottom part $y^{(i)}$ of the eigenvector $x^{(i)}$ is a null vector of the Schur complement $C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^{-1}E(\lambda_i)$. Having found $y^{(i)}$, one can recover the corresponding top part $u^{(i)}$ via

$$(3.4) \quad u^{(i)} = -B(\lambda_i)^{-1}E(\lambda_i)y^{(i)},$$

which requires the solution of a $d \times d$ block diagonal linear system.

What if $\lambda_i \in \Lambda(B, M_B)$? This case would seldom occur in practice, but we can come to understand it by writing $u^{(i)} = u_P^{(i)} + u_N^{(i)}$, where $u_P^{(i)} \in \text{Ran}(B(\lambda_i))$ and $u_N^{(i)} \in \text{Ker}(B(\lambda_i))$. In place of (3.4), the first block equation in (3.2) yields

$$(3.5) \quad u_P^{(i)} = -B(\lambda_i)^+ E(\lambda_i)y^{(i)},$$

where $B^+(\lambda_i)$ is the (Moore–Penrose) pseudoinverse of $B(\lambda_i)$. From this and the second block equation in (3.2), we obtain

$$(3.6) \quad E(\lambda_i)^T u_N^{(i)} + [C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^+ E(\lambda_i)]y^{(i)} = 0$$

instead of (3.3).

If it happens that $\text{Ran}(E(\lambda_i)) \perp \text{Ker}(B(\lambda_i))$, so that the first term in (3.6) vanishes, then the eigenvectors can be found in a manner analogous to the case when $\lambda_i \notin \Lambda(B, M_B)$ but with $B(\lambda_i)^{-1}$ replaced by $B(\lambda_i)^+$. Specifically, one can take $y^{(i)}$ from among the null vectors of the Schur-complement-like matrix $C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^+E(\lambda_i)$ and then recover $u_P^{(i)}$ from (3.5). The component $u_N^{(i)}$ can be taken arbitrarily from $\text{Ker}(B(\lambda_i))$ (i.e., from among the eigenvectors of (B, M_B) corresponding to the eigenvalue λ_i). We thus obtain an eigenspace of dimension $\dim \text{Ker}(C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^+E(\lambda_i)) + \dim \text{Ker}(B(\lambda_i))$. More generally, given $u_N^{(i)}$, one can solve (3.6) for $y^{(i)}$ and then leverage (3.5) to find $u_P^{(i)}$. Unfortunately, an easy way to compute $u_N^{(i)}$ does not appear to exist, and even if one did, forming and factoring $C(\lambda_i) - E^T(\lambda_i)B(\lambda_i)^+E(\lambda_i)$ would still be prohibitively expensive.

It is better simply to avoid the case $\lambda_i \in \Lambda(B, M_B)$ to begin with. This can be done by adjusting the partitioning until no eigenvalues of (B, M_B) lie within the search interval $[\alpha, \beta]$. As the likelihood of this being necessary is already small—in particular, we did not need to do this in any of the numerical experiments reported below—we will not attempt to develop a comprehensive strategy here, leaving this as a potential matter for future work.

3.2. Chebyshev approximation of eigenvector components. We have thus reduced the problem to that of finding those values ζ in $[\alpha, \beta]$ for which the *parameterized spectral Schur complement* [5, 19],

$$S(\zeta) = C(\zeta) - E^T(\zeta)B(\zeta)^{-1}E(\zeta),$$

is singular, assuming that no eigenvalue of (A, M) within $[\alpha, \beta]$ is also an eigenvalue of (B, M_B) . For $\zeta \notin \Lambda(B, M_B)$, let $\mu_1(\zeta), \dots, \mu_s(\zeta)$ and $y_1(\zeta), \dots, y_s(\zeta)$ denote the eigenvalues and corresponding eigenvectors of $S(\zeta)$, respectively:

$$S(\zeta)y_i(\zeta) = \mu_i(\zeta)y_i(\zeta), \quad i = 1, \dots, s.$$

The μ_i and y_i can be defined such that they are analytic functions of $\zeta \in \mathbb{C}$ away from $\Lambda(B, M_B)$. At each point of $\Lambda(B, M_B)$, they have at most a pole singularity [21, 26, 33, 39]. We refer to the μ_i as the *eigenvalues* of S . We also define

$$u_i(\zeta) = -B(\zeta)^{-1}E(\zeta)y_i(\zeta), \quad i = 1, \dots, s,$$

which is also analytic in ζ away from $\Lambda(B, M_B)$.

The matrix $S(\zeta)$ is singular precisely when one of its eigenvalues is zero: $\mu_i(\zeta) = 0$ for some i . The following result asserts that each of the $n_{\text{ev}} < s$ eigenvalues of A in $[\alpha, \beta]$, counted according to multiplicity, occurs as a zero of one and only one μ_i and that the top and bottom parts of the corresponding eigenvectors are given by u_i and y_i . (For eigenvalues of non-unit multiplicity, this statement is to be interpreted as saying that there is a distinct μ_i associated with each copy of the eigenvalue.) The assumption that $\beta < \min(\Lambda(B, M_B))$ ensures that $[\alpha, \beta]$ is free of any poles of S and the eigenvalues. The assumption that $n_{\text{ev}} < s$ ensures that the dimension of the space in which we plan to search is large enough to contain all the eigenvectors we seek.

PROPOSITION 3.1. *Assume $\beta < \min(\Lambda(B, M_B))$, and $n_{\text{ev}} < s$. Then, there exist n_{ev} distinct integers $\kappa_1, \dots, \kappa_{n_{\text{ev}}} \in \{1, 2, \dots, s\}$ such that*

$$(3.7) \quad \mu_{\kappa_i}(\lambda_i) = 0, \quad y^{(i)} = y_{\kappa_i}(\lambda_i), \quad u^{(i)} = u_{\kappa_i}(\lambda_i).$$

Proof. First, consider the case in which the λ_i are all simple eigenvalues. Following (3.3), we have $S(\lambda_i)y^{(i)} = 0$ for some where $y^{(i)} \neq 0$. The matrix $S(\lambda_i)$ is singular and has exactly one zero eigenvalue, denoted by $\mu_{\kappa_i}(\lambda_i)$, for some $1 \leq \kappa_i \leq s$. The expressions in (3.7) follow directly. It remains to show that $\kappa_i \neq \kappa_j$ when $i \neq j$. This follows from the fact that the μ_{κ_i} are free of poles and strictly decreasing on $[\alpha, \beta]$ [21], which implies that λ_i is the only root of μ_{κ_i} in $[\alpha, \beta]$.

That the result also holds in the case where one or more of the λ_i have non-unit multiplicity can be seen by considering arbitrarily small perturbations of (A, M) that have all simple eigenvalues and appealing to continuity. \square

We lose no generality in assuming that $\kappa_i = i$, and we will do so throughout the rest of the paper: from this point forward, μ_i will denote the eigencurve of S that crosses the real axis at λ_i .

Proposition 3.1 tells us that the components $u^{(i)}$ and $y^{(i)}$ of a sought eigenvector $x^{(i)}$ are equal to $y_i(\lambda_i)$ and $u_i(\lambda_i)$, respectively. Since both $y_i(\zeta)$ and $u_i(\zeta)$ are analytic on $[\alpha, \beta]$, they can be approximated accurately by interpolation at Chebyshev nodes. Specifically, for an integer $N \geq 1$, let

$$(3.8) \quad \chi_j = \frac{\alpha + \beta}{2} + \cos\left(\frac{j\pi}{N-1}\right) \frac{\beta - \alpha}{2}, \quad j = 0, \dots, N-1,$$

be the N Chebyshev nodes of the second kind in $[\alpha, \beta]$,¹ and let ℓ_j denote the j th Lagrange basis function for polynomial interpolation in these nodes. That is, ℓ_j is the unique polynomial of degree $N-1$ such that $\ell_j(\chi_k)$ is 1 if $k = j$ and 0 if $k \neq j$. Finally, let E_ρ be the *Bernstein ellipse* centered on $[\alpha, \beta]$ with parameter ρ ; that is, E_ρ is the open subset of \mathbb{C} bounded by the ellipse with foci at α and β and sum of the lengths of its semimajor and semiminor axes equal to ρ . Since $y_i(\zeta)$ and $u_i(\zeta)$ are analytic on $[\alpha, \beta]$, they can be analytically continued to E_ρ for some $\rho > 0$. We have:

PROPOSITION 3.2. *Assume that $\beta < \min(\Lambda(B, M_B))$, that $n_{\text{ev}} < s$, and that u_i and y_i are analytic in E_ρ for all $i = 1, \dots, n_{\text{ev}}$ and some $\rho > 0$. For each i , there exists $w^{(i)} \in \mathbb{R}^N$ such that*

$$x^{(i)} = \begin{bmatrix} u^{(i)} \\ y^{(i)} \end{bmatrix} = \begin{bmatrix} u_i(\chi_0) & \cdots & u_i(\chi_{N-1}) \\ y_i(\chi_0) & \cdots & y_i(\chi_{N-1}) \end{bmatrix} w^{(i)} + O(\rho^{-N}).$$

Proof. Let $w_j^{(i)} = \ell_j(\lambda_i)$ for $j = 0, \dots, N-1$. Then, the top d (respectively, bottom s) components of the matrix-vector product give the value at λ_i of the polynomial interpolant to $u^{(i)}$ (respectively, $y^{(i)}$) in the Chebyshev nodes χ_j . The result now follows from a standard theorem on the convergence of Chebyshev interpolants to analytic functions [40, Theorem 8.2]. \square

Instead of interpolating u_i and y_i directly, we use their samples at the Chebyshev nodes to generate a subspace in which to look for approximations to the $x^{(i)}$. This approach eliminates the need to keep track of the association between the samples and the eigencurves, which may be difficult if the eigencurves cross.² Proposition 3.2 ensures that this subspace contains good approximations to the $x^{(i)}$ for large enough N . We can express this fact as a statement about the angle between this subspace and the sought eigenspace:

¹For $N = 1$, we take $\chi_0 = (\alpha + \beta)/2$.

²For example, it can happen that $\mu_2(\chi_j) < \mu_1(\chi_j) < \mu_3(\chi_j) < \cdots < \mu_s(\chi_j)$ for some j . If so, the eigenvector of $S(\chi_j)$ corresponding to its smallest eigenvalue is a sample of $y_2(\chi_j)$, not $y_1(\chi_j)$, even though μ_1 is the eigencurve for the smallest eigenvalue of (A, M) .

Algorithm 3.1 The proposed algorithm.

```

1: input:  $A \in \mathbb{R}^{n \times n}$ ,  $M \in \mathbb{R}^{n \times n}$ ,  $N \in \mathbb{N}$ ,  $\alpha \in \mathbb{R}$ ,  $\beta \in \mathbb{R}$ ,  $n_{\text{ev}} \in \mathbb{Z}$ ,  $Y = 0$ ,  $V = 0$ 
2: output: approximations of eigenpairs  $(\lambda_i, x^{(i)})$ ,  $i = 1, \dots, n_{\text{ev}}$ 

3: /* Pre-processing: reorder matrices  $A$  and  $M$  */
4: ▷ Call a  $p$ -way edge separator to partition the graph  $\mathcal{G}_{A,M}$ .
5: ▷ If  $\beta < \min(\Lambda(B, M_B))$  continue, else set  $p := 2p$  and repeat step 4.

6: /* Main loop; embarrassingly parallel over the  $N$  Chebyshev nodes */
7: for  $j = 0, \dots, N-1$  do
8:   ▷ Set  $\chi_j = \frac{\alpha + \beta}{2} + \cos\left(\frac{j\pi}{N-1}\right) \frac{\beta - \alpha}{2}$ .
9:   ▷ Set  $Y_j = [y_1(\chi_j), \dots, y_{n_{\text{ev}}}(\chi_j)]$ .
10:  ▷ Solve  $B(\chi_j)V_j = -E(\chi_j)Y_j$ .
11: end for

12: /* Rayleigh-Ritz projection phase */
13: ▷ Set  $R = \begin{bmatrix} V_0 & \cdots & V_{N-1} \\ Y_0 & \cdots & Y_{N-1} \end{bmatrix}$ .
14: ▷ Optionally, orthonormalize the columns of  $R$ .
15: ▷ Compute the  $n_{\text{ev}}$  algebraically smallest eigenvalues and associated eigenvectors
    of the eigenvalue problem  $(R^T A R)f = \theta(R^T M R)f$ .
16: ▷ Return  $(\theta_i, PRf^{(i)}) \approx (\lambda_i, x^{(i)})$ ,  $i = 1, \dots, n_{\text{ev}}$ .
```

COROLLARY 3.3. Let $\mathcal{X} = \text{span}\{x^{(1)}, \dots, x^{(n_{\text{ev}})}\}$, and let

$$\mathcal{R} = \text{span} \left\{ \begin{bmatrix} u_1(\chi_0) \\ y_1(\chi_0) \end{bmatrix}, \dots, \begin{bmatrix} u_1(\chi_{N-1}) \\ y_1(\chi_{N-1}) \end{bmatrix}, \dots, \begin{bmatrix} u_{n_{\text{ev}}}(\chi_0) \\ y_{n_{\text{ev}}}(\chi_0) \end{bmatrix}, \dots, \begin{bmatrix} u_{n_{\text{ev}}}(\chi_{N-1}) \\ y_{n_{\text{ev}}}(\chi_{N-1}) \end{bmatrix} \right\}.$$

Then,

$$\sin \theta(\mathcal{X}, \mathcal{R}) = O(\rho^{-N}),$$

where $\theta(\mathcal{X}, \mathcal{R})$ is the largest principal angle between \mathcal{X} and the closest subspace of \mathcal{R} to \mathcal{X} with the same dimension as \mathcal{X} .

Proof. The quantity $\sin \theta(\mathcal{X}, \mathcal{R})$ is known as the *gap* between \mathcal{X} and \mathcal{R} and can be expressed as [3] [26, sect. IV.2.1] [38, sect. II.4]

$$\sin \theta(\mathcal{X}, \mathcal{R}) = \max_{x \in \mathcal{X}} \min_{r \in \mathcal{R}} \frac{\|x - r\|}{\|x\|}.$$

The result follows immediately from this formula and Proposition 3.2. \square

3.3. A parallel algorithm. Our algorithm builds the subspace \mathcal{R} of Corollary 3.3 and then uses Rayleigh–Ritz projection to extract approximations to the $x^{(i)}$ from \mathcal{R} . The procedure is summarized in Algorithm 3.1.

For each Chebyshev node χ_j , Algorithm 3.1 computes the eigenvectors associated with the n_{ev} algebraically smallest eigenvalues of $S(\chi_j)$. These eigenvectors form the $s \times n_{\text{ev}}$ matrix Y_j (step 9). Then, the algorithm computes the matrix V_j , which requires the solution of a linear system with the coefficient matrix $B(\chi_j)$ and n_{ev} right-hand sides (step 10). Finally, the algorithm uses Rayleigh–Ritz projection (steps 15–16) to

approximate the sought eigenpairs of (A, M) . The dimension of the projected pencil is at most Nn_{ev} , and the associated eigenvalue problem is solved by a dense, symmetric eigenvalue solver.

The **for** loop in steps 7–11 is embarrassingly parallel: each matrix pair (Y_j, V_j) can be computed independently of the other pairs. The computation of V_j can be further decomposed into the solution of p independent linear systems. Partition V_j and Y_j by rows as

$$V_j = \begin{bmatrix} V_{1,j} \\ \vdots \\ V_{p,j} \end{bmatrix}, \quad Y_j = \begin{bmatrix} Y_{1,j} \\ \vdots \\ Y_{p,j} \end{bmatrix},$$

where $V_{k,j}$ and $Y_{k,j}$ are associated with the k th subdomain. Then,

$$\begin{bmatrix} B_1(\chi_j) & & \\ & \ddots & \\ & & B_p(\chi_j) \end{bmatrix} \begin{bmatrix} V_{1,j} \\ \vdots \\ V_{p,j} \end{bmatrix} = \begin{bmatrix} E_1(\chi_j)Y_{1,j} \\ \vdots \\ E_p(\chi_j)Y_{p,j} \end{bmatrix}$$

(where we have extended the notation (3.1) to the blocks comprising B , M_B , E , and M_E in the obvious way), and so the $V_{k,j}$ can be computed by solving

$$B_k(\chi_j)V_{k,j} = -E_k(\chi_j)Y_{k,j}, \quad k = 1, \dots, p.$$

These p linear systems can be solved in parallel.

3.4. Practical details. If the desired number n_{ev} of eigenvalues is not known a priori, it can be computed directly by decomposing $A - \alpha M$ and $A - \beta M$ in LDL^T factorizations and using Sylvester's law of inertia [7]. Alternatively, if this is too expensive, one can estimate n_{ev} using a spectral density profile of (A, M) [44]. To reduce the chance of the algorithm missing eigenvalues, we recommend taking n_{ev} slightly larger than estimated or required. To further reduce this chance, one can apply a few steps of subspace iteration or Lanczos with polynomial filtering and deflation as post-processing after step 16. Since the number of iterations needed should not be large, one can use iterative methods to approximate M^{-1} instead of exact factorizations.

The results of section 3.2 relied on the hypothesis $\beta < \min(\Lambda(B, M_B))$. How can we enforce this requirement in practice? Observe that if $\delta_1, \delta_2, \dots, \delta_d$ are the eigenvalues of (B, M_B) , then $\lambda_i \leq \delta_i \leq \lambda_{i+n-d}$, $i = 1, \dots, d$ by a version of the interlacing theorem. Therefore, decreasing the dimension of (B, M_B) , i.e., increasing the value of p , makes it more likely that $\lambda_{n_{\text{ev}}} \leq \beta < \delta_1$. Algorithm 3.1 adopts the practical strategy of doubling p until $\beta < \min(\Lambda(B, M_B))$ is satisfied (step 5).

Finally, note that Algorithm 3.1 is a one-shot method in the sense that if the accuracy of the approximate eigenpairs is not satisfactory, then the whole process must be repeated with a higher value of N . We find that in practice, $N = 8$ reaches nearly the maximum attainable accuracy on a wide range of problems; see Section 5. If one wishes to apply Algorithm 3.1 for several values of N , it is beneficial to take these N to have the form $N(k) = 2^k + 1$ for integers k . Having run the algorithm with $N = N(k)$, one can reduce the computational cost of running the algorithm with $N = N(k+1)$ by exploiting the fact that the nodes (3.8) for $N(k)$ are a subset of those for $N(k+1)$ and reusing the samples taken during the $N = N(k)$ run.

Besides increasing N , one can also improve the accuracy of one or more of the eigenpairs by using the approximate eigenvectors obtained from Algorithm 3.1 as the initial subspace for an implicitly-restarted (or thick-restarted) Lanczos method [8, 43] applied to (A, M) . This technique can also be used to ensure that all n_{ev} eigenpairs of (A, M) have been computed (i.e., none have been missed) by checking to see if the algebraically smallest eigenvalue returned by the restarted Lanczos method is smaller than β .

4. A distributed-memory implementation. We now describe our parallel implementation of Algorithm 3.1 based on the MPI standard. Throughout this discussion, we assume a distributed-memory computing environment with $N_p = p_r p_c$ MPI processes organized in a $p_r \times p_c$ 2D MPI grid. In addition to the default communicator `MPI_COMM_WORLD`, we denote by G_i^r , $i = 0, \dots, p_r - 1$, and G_j^c , $j = 0, \dots, p_c - 1$, the MPI communicators associated with the i th row and j th column of the grid, respectively.

Our parallel implementation utilizes the row dimension of the grid for domain decomposition data parallelism (i.e., distributed storage of A and M), and the column dimension of the grid for model parallelism (i.e., distribution over the N Chebyshev nodes). Therefore, the row and column dimensions of the grid satisfy the inequalities $p_r \leq p$ and $p_c \leq N$, respectively.

4.1. Data distribution on 2D MPI grids. First, we consider the data distribution along the row dimension of the grid. For each communicator G_j^c , $j = 0, \dots, p_c - 1$, we distribute A and M such that the p_r MPI processes associated with G_j^c hold a unique subset of the partitions of the graph $\mathcal{G}_{A,M}$. In particular, let p be a scalar multiple of p_r , and set $\tau = p/p_r$. Then, the i th process is assigned data associated with partitions $i\tau + 1, i\tau + 2, \dots, (i+1)\tau$, i.e.,

$$\text{Data held by process } i \text{ of } G_j^c: \begin{cases} B_{i\tau+1}, \dots, B_{(i+1)\tau}, M_{B_{i\tau+1}}, \dots, M_{B_{(i+1)\tau}} \\ E_{i\tau+1}, \dots, E_{(i+1)\tau}, M_{E_{i\tau+1}}, \dots, M_{E_{(i+1)\tau}} \\ C_{i\tau+1,:}, \dots, C_{(i+1)\tau,:}, M_{C_{i\tau+1,:}}, \dots, M_{C_{(i+1)\tau,:}} \end{cases},$$

where the subscript “:” represents all column indices of matrices C and M_C . Ordering the unknowns/equations by increasing MPI rank leads to the following global representation of A (and similarly for M):

$$(4.1) \quad A = \begin{bmatrix} B_1 & E_1 & & & \\ E_1^T & C_{1,1} & C_{1,2} & & C_{1,p_r} \\ & & B_2 & E_2 & \\ & C_{2,1} & E_2^T & C_{2,2} & C_{2,p_r} \\ & & & \ddots & \\ & & & & B_{p_r} & E_{p_r} \\ & C_{p_r,1} & C_{p_r,2} & E_{p_r}^T & C_{p_r,p_r} \end{bmatrix}.$$

The ordering in (4.1) is more natural from the perspective of parallel computing than that in (2.1), which is more natural for discussing the linear algebra.

We now focus on the column dimension of the grid. Let N be a scalar multiple of p_c , and set $\eta = N/p_c$. We distribute the N Chebyshev nodes across the p_c MPI processes of each row communicator G_i^r , $i = 0, \dots, p_r - 1$, such that each process receives exactly η unique Chebyshev nodes. In particular, the j th process associated is assigned the Chebyshev node(s) $\chi_{j\eta+1}, \dots, \chi_{(j+1)\eta}$ $j = 0, \dots, p_c - 1$. From a parallel

efficiency perspective, it is advisable to exhaust parallelism across the N Chebyshev nodes first, by setting $p_c = N$, since this level of parallelism involves no communication among groups of processes assigned different Chebyshev nodes.

An illustration of the data distribution on a 2D MPI grid with $N_p = 16$ processes and $N = 8$ Chebyshev nodes is shown in Figures 4.1 and 4.2 where the dimensions of the grid are $(p_r, p_c) = (4, 4)$ and $(p_r, p_c) = (2, 8)$, respectively. For $(4, 4)$ case we have $p_c < N$, and each column subgrid is responsible for processing $h = 8/4 = 2$ Chebyshev nodes, while the computation of each matrix pair (Y_j, V_j) exploits four MPI processes. Contrast this with the $(2, 8)$ case, in which each separate column subgrid handles exactly one Chebyshev node ($\eta = 1$), leading to trivial parallelism with respect to the N Chebyshev nodes, but the computation of each matrix pair (Y_j, V_j) utilizes just two processes.

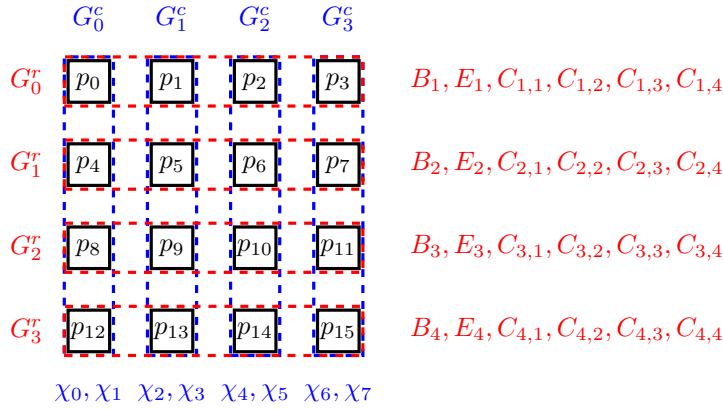


Fig. 4.1: Distribution of blocks of A and Chebyshev nodes over a 2D MPI grid with $N_p = 16$, $N = 8$, and $(p_r, p_c) = (4, 4)$. The distribution of M is identical to that of A .

4.2. Computation of Y_j via PARPACK. Our implementation computes the eigenvectors of the Schur complement matrices $S(\chi_j)$, $j = 0, \dots, N - 1$, via the **PARPACK**³ software library, a distributed-memory implementation of **ARPACK** [30]. The main distributed-memory kernels of **PARPACK** are: (a) orthogonalization of the Krylov basis, and (b) a user-defined routine that performs distributed matrix-vector multiplication with $S(\chi_j)$.

Regarding (a), consider first the case $p_c = N$. Orthonormalizing the basis vectors computed on each m -step cycle of the implicitly restarted Arnoldi method via Gram–Schmidt costs $O(sm^2)$ floating-point operations and $O(\log(p_r)m^2)$ point-to-point MPI messages. This communication cost increases proportionally with the number of Chebyshev nodes processed by each column subgrid. In particular, when $p_c = 1$, i.e., all available N_p MPI processes are assigned to the default communicator, **PARPACK** requires $O(N \log(N_p)m^2)$ MPI messages just for Gram–Schmidt.

As for (b), note that the product between the distributed matrix $S(\chi_j)$ and a

³<https://github.com/opencollab/arpac-ng>

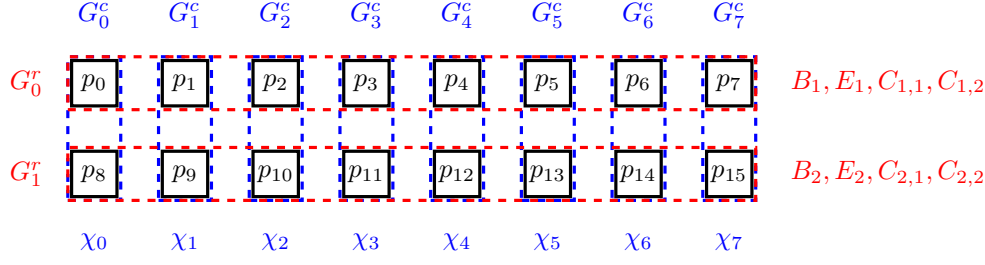


Fig. 4.2: Distribution of blocks of A and Chebyshev nodes over a 2D MPI grid with $N_p = 16$, $N = 8$, and $(p_r, p_c) = (2, 8)$. The distribution of M is identical to that of A .

distributed vector $f = [f_1^T \ \cdots \ f_p^T]^T \in \mathbb{R}^s$ can be written as

$$(4.2) \quad S(\chi_j)f = \begin{bmatrix} \sum_{k \in \mathcal{N}_1} C_{1,k}(\chi_j)f_k \\ \vdots \\ \sum_{k \in \mathcal{N}_p} C_{p,k}(\chi_j)f_k \end{bmatrix} - \begin{bmatrix} B_1(\chi_j)^{-1}E_1(\chi_j)f_1 \\ \vdots \\ B_p(\chi_j)^{-1}E_p(\chi_j)f_p \end{bmatrix},$$

where \mathcal{N}_i denotes the list of partitions adjacent to partition i (and where we have extended the notation (3.1) to the blocks of $A - \zeta M$ defined by (4.1) in the obvious way). Due to the partitioning, the second term on the right-hand side of (4.2) can be computed in an embarrassingly parallel manner. On the other hand, the first term of the right-hand side of (4.2) requires point-to-point communication between processes handling neighboring partitions.

4.3. Orthonormalization of the Rayleigh–Ritz basis. Our implementation orthonormalizes the columns of the Rayleigh–Ritz projection matrix R via Gram–Schmidt. To take advantage of all N_p MPI processes, we exploit the default communicator `MPI_COMM_WORLD`.

The (i, j) process of the $p_r \times p_c$ 2D MPI grid holds the submatrices $V_{i,j}$ and $Y_{i,j}$, leading to the following representation of R as a 2D logical array:

$$\widehat{R}_{2D} = \begin{bmatrix} G_0^c & G_1^c & \cdots & G_{p_c-1}^c \\ \begin{bmatrix} V_{0,0} \\ Y_{0,0} \end{bmatrix} & \begin{bmatrix} V_{0,1} \\ Y_{0,1} \end{bmatrix} & \cdots & \begin{bmatrix} V_{0,p_c-1} \\ Y_{0,p_c-1} \end{bmatrix} \\ \vdots & \vdots & \vdots & \vdots \\ \begin{bmatrix} V_{p_r-1,0} \\ Y_{p_r-1,0} \end{bmatrix} & \begin{bmatrix} V_{p_r-1,1} \\ Y_{p_r-1,1} \end{bmatrix} & \cdots & \begin{bmatrix} V_{p_r-1,p_c-1} \\ Y_{p_r-1,p_c-1} \end{bmatrix} \end{bmatrix} \begin{matrix} G_0^r \\ \vdots \\ G_{p_r-1}^r \end{matrix}.$$

The goal is to transform \widehat{R}_{2D} into a $n \times Nn_{ev}$ matrix R_{1D} such that each one of the N_p processes holds a submatrix that has roughly n/N_p rows and Nn_{ev} columns. This can be achieved by the following two-step procedure. First, we perform a gather reduction on the submatrices $[V_{i,j}^T \ Y_{i,j}^T]^T$, $j = 0, \dots, p_c - 1$. This reduction is performed independently within each communicator G_i^r , $i = 0, \dots, p_r - 1$. Second, each process associated with G_i^r discards all rows of the previously reduced matrix except for a

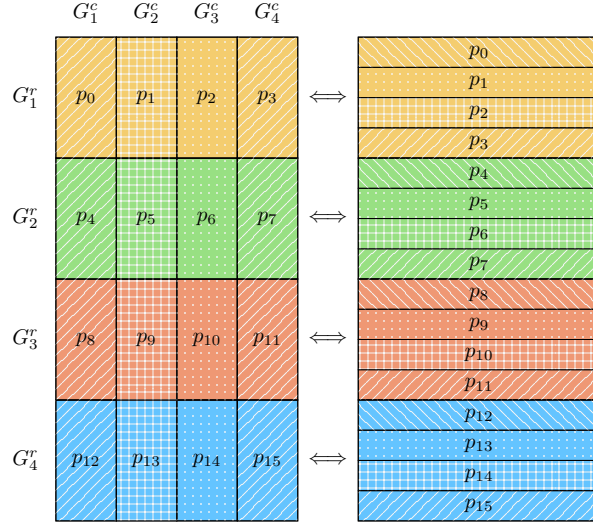


Fig. 4.3: 2D-to-1D (and vice-versa) MPI grid mapping. Left: color/pattern layout of a 2D grid of MPI processes with $p_c = p_r = 4$. Right: color/pattern layout of the same grid collapsed in 1D MPI grid topology.

unique, continuous set of rows. We can then write

$$(4.3) \quad R_{1D} = \begin{bmatrix} V_{0,0} & \cdots & V_{0,p_c-1} \\ Y_{0,0} & \cdots & Y_{0,p_c-1} \\ \vdots & \cdots & \vdots \\ \vdots & \cdots & \vdots \\ V_{p_r-1,0} & \cdots & V_{p_r-1,p_c-1} \\ Y_{p_r-1,0} & \cdots & Y_{p_r-1,p_c-1} \end{bmatrix} = \begin{bmatrix} R_{0,0} \\ \vdots \\ R_{0,p_c-1} \\ \vdots \\ R_{p_r-1,0} \\ \vdots \\ R_{p_r-1,p_c-1} \end{bmatrix},$$

where $R_{i,j}$ is held by the MPI process of rank $ip_c + j$ associated with `MPI_COMM_WORLD`, i.e., the j th process associated with the row communicator G_i^r . This can be done efficiently in a single line of code by calling `MPI_Alltoall` independently within each communicator G_i^r , $i = 0, \dots, p_r - 1$. A graphical illustration of this 2D-to-1D grid remapping is shown in Figure 4.3.

Once the remapping is complete, we apply distributed block Gram-Schmidt to the columns of R_{1D} using `MPI_COMM_WORLD` and a block size equal to n_{ev} . Then, we map R_{1D} back to the 2D layout by reversing the above procedure. For further details on parallel Gram-Schmidt, including a discussion of numerical stability, see [4, 9].

4.4. Formation and solution of the projected eigenvalue problem. Finally, we form the projected pencil $(R^T AR, R^T MR)$ and find its eigenvalues. As the projected pencil is small, once it is formed, we compute its eigenvalues serially using the `DSYGVX` routine from LAPACK [2]. The remainder of this section is devoted to discussing our approach to forming $R^T AR$ within the 2D distributed-memory data layout described above. The procedure for forming $R^T MR$ is identical.

We form $R^T AR$ in two phases. Let $R_j = [V_j^T \ Y_j^T]^T$. In the first phase, we

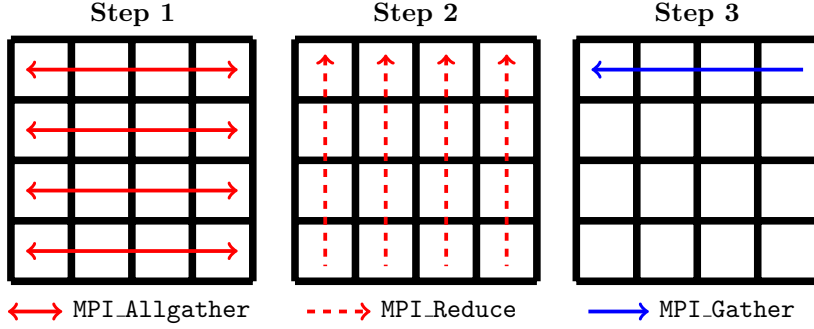


Fig. 4.4: Communication pattern for the distributed-memory computation of $R^T AR$ and $R^T MR$ using our 2D MPI data layout ($p_r = p_c = 4$). The root process of `MPI_COMM_WORLD` is located in the upper left corner.

420 compute $AR = [AR_0 \ AR_1 \ \cdots \ AR_{N-1}]$. When $p_c = N$, this operation is embar-
 421 rassingly parallel, since each of the products AR_j , $j = 0, \dots, N-1$, can be computed
 422 independently. Using the rank-based representation of A from (4.1), we write

$$423 \quad (4.4) \quad AR_j = \begin{bmatrix} B_1 & E_1 & & & & \\ E_1^T & C_{1,1} & & C_{1,2} & & C_{1,p_r} \\ & C_{2,1} & B_2 & E_2 & & \\ & & E_2^T & C_{2,2} & & C_{2,p_r} \\ & & & & \ddots & \\ & & & & B_{p_r} & E_{p_r} \\ & C_{p_r,1} & C_{p_r,2} & & E_{p_r}^T & C_{p_r,p_r} \end{bmatrix} \begin{bmatrix} V_{0,j} \\ Y_{0,j} \\ V_{1,j} \\ Y_{1,j} \\ \vdots \\ V_{p_r-1,j} \\ Y_{p_r-1,j} \end{bmatrix}$$

424 Communication between different MPI processes of G_j^c is point-to-point, and the i th
 425 process needs to send $Y_{i,j}$ to the k th process if and only if $C_{k,i} \neq 0$.

426 The second phase multiplies R^T and AR and stores the matrix product in the
 427 root process of `MPI_COMM_WORLD`. To achieve this, we apply the following procedure,
 428 which is illustrated in Figure 4.4:

- 429 1. Apply `MPI_Allgather` on the submatrices $[AR_j]_i$, $j = 0, \dots, p_c - 1$, across
 430 the row communicator G_i^r , where $[AR_j]_i$ denotes the submatrix of AR_j held
 431 by the i th process. Each process associated with G_i^r then has its own copy
 432 of the matrix $[[AR_0]_i \ [AR_2]_i \ \cdots \ [AR_{p_c-1}]_i]$.
- 433 2. The i th process associated with the column communicator G_j^c then computes
 434 $Z_{i,j} = R_{i,j}^T [[AR_0]_i \ [AR_2]_i \ \cdots \ [AR_{p_c-1}]_i]$ and calls `MPI_Reduce` on the
 435 data $Z_{i,j}$ associated with the processes in G_j^c .
- 436 3. At the end of the previous step, the k th MPI process associated with G_0^r
 437 holds the k th block of rows of the matrix $R^T AR$. Finally, all processes in G_0^r
 438 call `MPI_Gather`, creating $R^T AR$ in the root process.

439 **5. Numerical experiments.** We now illustrate the performance of Algorithm
 440 3.1 in both sequential and distributed-memory computing environments. We per-
 441 formed our experiments on the Minnesota Supercomputing Institute's `Mesabi` cluster.
 442 Each node of `Mesabi` is equipped with 64 GB of system memory and two 12-core 2.5
 443 GHz Intel Xeon E5-2680v3 (Haswell) CPUs. We built our code with the Intel ICC

18.0.0 compiler. We used the Intel Math Kernel Library (MKL) for basic matrix operations, including its sparse matrix routines and its implementation of the standard BLAS and LAPACK libraries for sequential dense matrix operations. While it is possible to exploit shared-memory parallelism, the experiments described below use just one thread per MPI process.

To compute the n_{ev} sought eigenvectors of the spectral Schur complements $S(\chi_j)$, we used PARPACK with full orthogonalization and restart dimension $m = 2n_{\text{ev}}$. The linear systems involving the block-diagonal matrix $B(\chi_j)$ were solved with the Intel MKL implementation of the PARDISO solver. For the search interval $[\alpha, \beta]$, we set $\alpha = 0$, $\beta = (\lambda_{n_{\text{ev}}} + \lambda_{n_{\text{ev}}+1})/2$ in all experiments.

5.1. Numerical illustration. We first demonstrate the qualitative performance of Algorithm 3.1 on a set of four small problems:

- “APF4686,” a standard eigenvalue problem of dimension $n = 4,686$ generated by the ELSEES quantum mechanical nanomaterial simulator⁴ [16],
- “Kuu/Muu,” a generalized eigenvalue problem of dimension $n = 7,102$ from the SuiteSparse matrix collection⁵ [10],
- “FDmesh,” a standard eigenvalue problem generated by a regular 5-point finite difference discretization of the Laplacian on a square, and
- “FEmesh,” a generalized eigenvalue problem obtained by discretizing the Laplacian on a square with linear finite elements.

For the latter two, the discretization fineness was chosen to yield matrices of dimension $n \approx 20,000$, and the associated pencils have several eigenvalues of multiplicity 2.

Figure 5.1 plots the relative errors in the eigenvalues returned by Algorithm 3.1 and the corresponding residual norms for the problems “APF4686” (left, $n_{\text{ev}} = 30$) and “Kuu/Muu” (right, $n_{\text{ev}} = 100$) for $N = 2, 4, 6, 8$. Figure 5.2 plots the same quantities for “FDmesh” (left) and “FEmesh” (right), where $n_{\text{ev}} = 100$ in both cases. In agreement with the discussion in Section 3, increasing N leads to greater accuracy in the approximation of the sought eigenpairs. Moreover, all eigenpairs are approximated to comparable accuracies for a given value of N , i.e., the accuracy of an eigenpair is relatively insensitive to the location of the eigenvalue inside $[\alpha, \beta]$.

5.2. Distributed-memory performance. We now illustrate the distributed-memory efficiency of Algorithm 3.1 on a variety of larger problems coming from discretizations of the Laplacian as well as general symmetric matrices and pencils from the SuiteSparse collection.⁶ Unless otherwise indicated, throughout the rest of this section, we take $n_{\text{ev}} = 100$, and we set the second dimension of the 2D MPI grid to be $p_c = N$. In most of the tests we report the results with $N = 8$ or $N = 4$. The parallel efficiency of a program executing on $\phi \in \mathbb{N}$ processes is $P(\phi) = T_1/(\phi T_\phi)$, where T_ϕ denotes the wall-clock time for execution on ϕ processes.

We benchmark Algorithm 3.1 against PARPACK applied directly to the pencil (A, M) both with and without shift-and-invert. PARPACK requires the application of either M^{-1} (without shift-and-invert) or A^{-1} (with shift-and-invert), and since A and M are distributed, we used a distributed direct solver for these operations. The results reported here were generated using the MUMPS package [1], but our code also provides interfaces for SuperLU_Dist [32] and the Intel Cluster Sparse Solver (provided in the MKL). For PARPACK, we report the wall-clock time and parallel efficiency

⁴<http://www.elses.jp>

⁵<https://sparse.tamu.edu/>

⁶Our implementation is available publicly at <https://github.com/Hitenze/Schurcheb>.

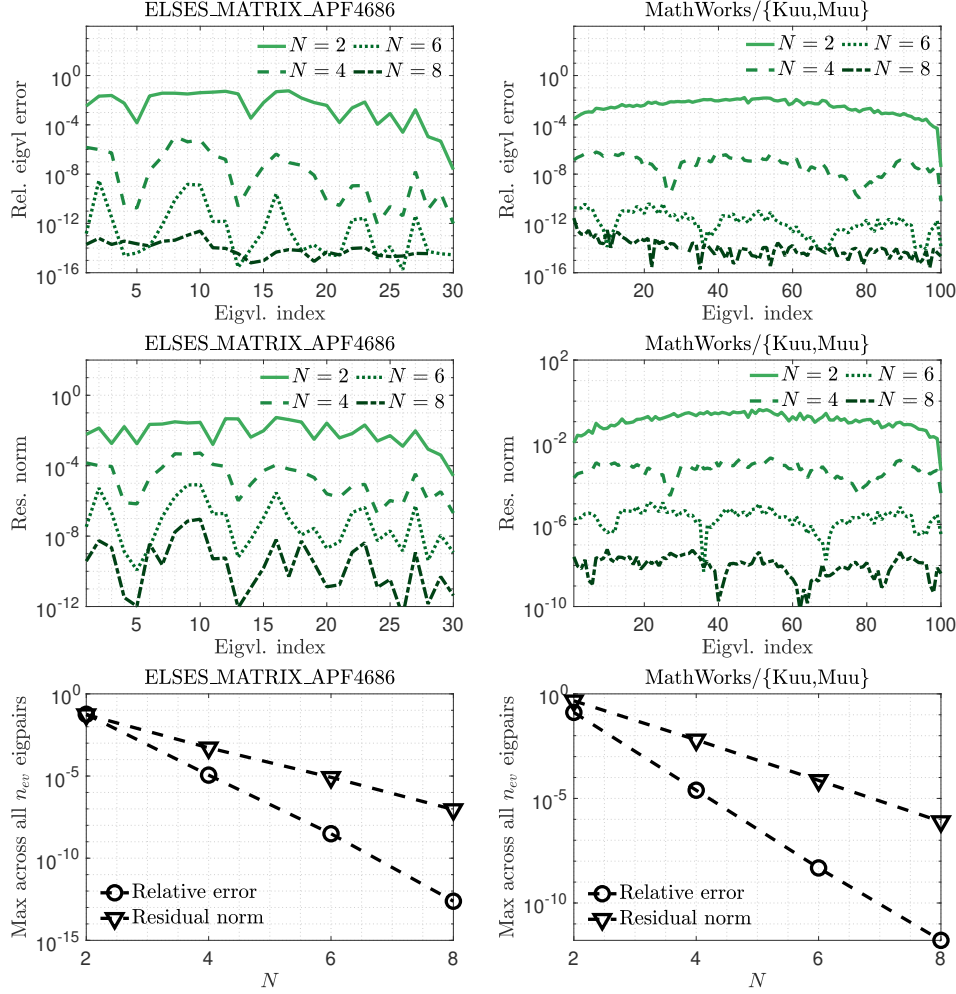


Fig. 5.1: Relative errors in the eigenvalues returned by Algorithm 3.1 (top) and corresponding residual norms (center) for various values of N for the problems “APF4686” (left, $n_{ev} = 30$) and “Kuu/Muu” (right, $n_{ev} = 100$). The bottom two figures plot the maximum relative error in the eigenvalues and maximum residual norm across all n_{ev} eigenpairs.

for a restart length equal to $m = 2n_{ev}$ with all MPI processes bundled in the default communicator `MPI_COMM_WORLD`. To keep the comparisons fair, the convergence tolerance passed to `PARPACK` for each problem is set to the maximum residual norm returned by Algorithm 3.1.

5.2.1. Eigenvalue problems from finite difference discretizations. First, we apply Algorithm 3.1 to matrices arising from finite difference discretizations of the Dirichlet eigenvalue problem,

$$(5.1) \quad \begin{aligned} -\Delta u &= \lambda u && \text{in } \Omega \\ u &= 0 && \text{on } \partial\Omega, \end{aligned}$$

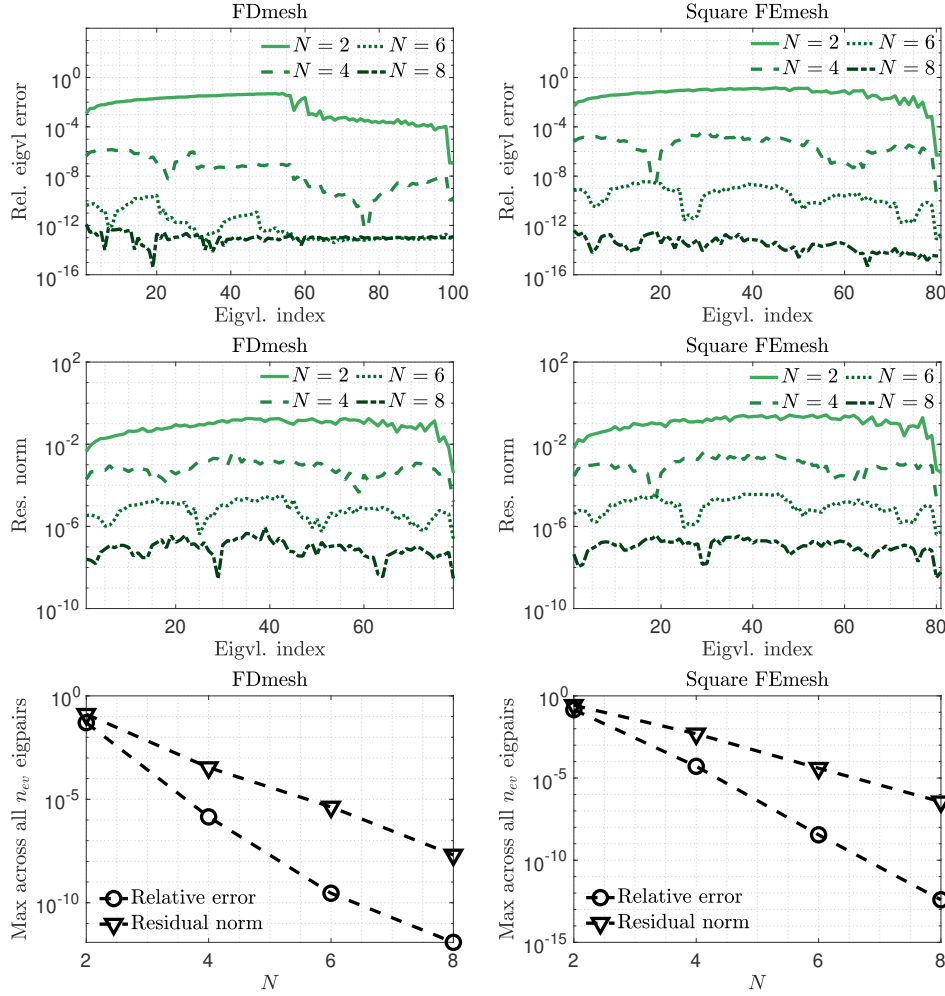


Fig. 5.2: Relative errors in the eigenvalues returned by Algorithm 3.1 (top) and corresponding residual norms (center) for various values of N for the problems “FDmesh” (left) and “FEmesh” (right). The bottom two figures plot the maximum relative error in the eigenvalues and maximum residual norm across all n_{ev} eigenpairs.

where Δ denotes the Laplacian and Ω is either the square $(0,1)^2$ in 2D or the cube $(0,1)^3$ in 3D. We use the standard 5- and 7-point stencils in 2D and 3D, respectively. All these eigenvalue problems are standard ones, with M equal to the identity matrix.

Our first set of experiments focuses on the strong scaling of Algorithm 3.1. We take $n_{ev} = 100$ and use $N = 4, 8$ Chebyshev nodes. In our results, we refer to Algorithm 3.1 with $N = 4$ as `SchurCheb(4)` and with $N = 8$ as `SchurCheb(8)`. We first consider three different 2D discretizations with matrix sizes $n = 257 \times 256$, $n = 513 \times 512$, and $n = 1025 \times 1024$, respectively. Table 5.1 lists the maximum relative error in the eigenvalues returned by Algorithm 3.1. Figure 5.3 (left) plots the parallel efficiency of Algorithm 3.1 for $N = 8$, where we report separately the parallel efficiencies associated with: (a) computation of the eigenvector matrices Y_j ,

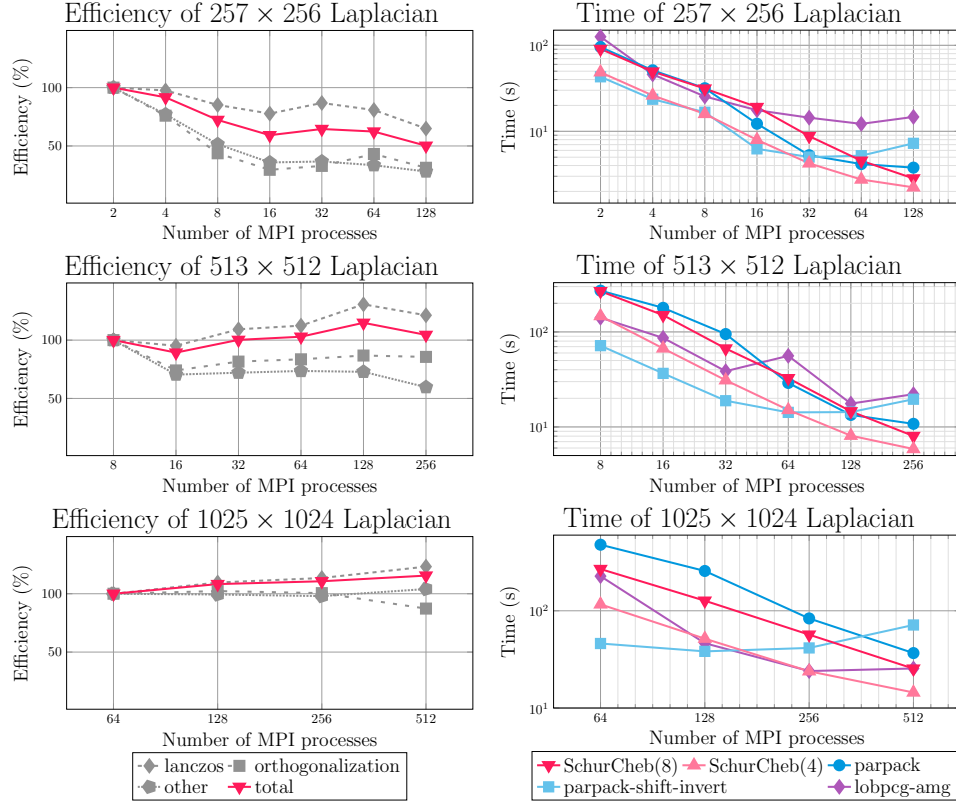


Fig. 5.3: Left: parallel efficiency of Algorithm 3.1 with $n_{\text{ev}} = 100$ and $p_c = N = 8$. Right: wall-clock time comparison between Algorithm 3.1 with $N = 4, 8$ and PARPACK with and without shift-and-invert. The number of MPI processes ranges from $N_p = 2$ to $N_p = 512$. The number of partitions is set equal to $p = 32$ ($n = 257 \times 256$), $p = 64$ ($n = 513 \times 512$), and $p = 128$ ($n = 1025 \times 1024$), when $N = 8$. The value of p is doubled when $N = 4$ since each column communicator now has twice as many processes.

Table 5.1: Maximum relative error in the eigenvalues returned by Algorithm 3.1 for the finite difference problems.

	$n = 257 \times 256$	$n = 513 \times 512$	$n = 1025 \times 1024$	$n = 65 \times 64 \times 63$
SchurCheb(4)	5.1×10^{-4}	8.2×10^{-5}	1.4×10^{-4}	9.1×10^{-5}
SchurCheb(8)	2.3×10^{-9}	2.9×10^{-11}	2.5×10^{-7}	1.9×10^{-10}

$j = 0, \dots, N - 1$, (b) orthonormalization of the projection matrix R , and (c) everything else. Since $p_c = N$, the computation of the Y_j is embarrassingly parallel, leading to nearly perfect efficiency for this step. On the other hand, both the orthonormalization of R and the formation of $R^T A R$ require communication among the N_p processes, and their efficiency can deteriorate for larger values of N_p . Note also that the parallel granularity of Algorithm 3.1 is lower for smaller problem sizes, leading to lower efficiencies compared to larger problems.

Figure 5.3 (right) plots the wall-clock time achieved by Algorithm 3.1 for $N = 4, 8$,

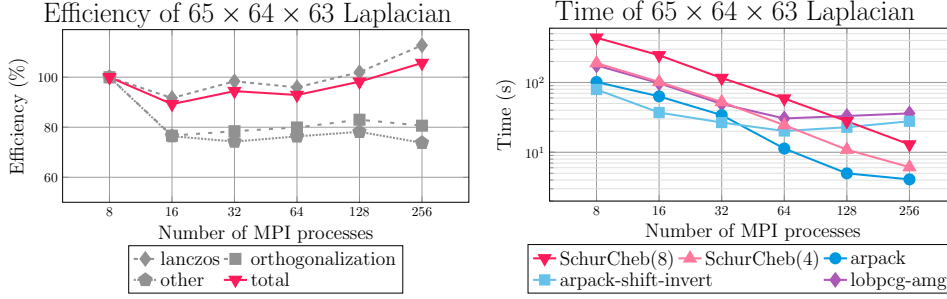


Fig. 5.4: Left: parallel efficiency of Algorithm 3.1 with $n_{\text{ev}} = 100$ and $p_c = N = 8$. Right: wall-clock time comparison between Algorithm 3.1 with $N = 4, 8$ and PARPACK with and without shift-and-invert. The number of MPI processes ranges from $N_p = 8$ to $N_p = 256$. The number of partitions is set to $p = 64$ ($N = 8$) and $p = 128$ ($N = 4$).

Table 5.2: Peak memory consumption of Algorithm 3.1 and of PARPACK with shift-and-invert for the finite difference problems.

	$n = 257 \times 256$ $N_p = 128$	$n = 513 \times 512$ $N_p = 256$	$n = 1025 \times 1024$ $N_p = 512$	$n = 65 \times 64 \times 63$ $N_p = 256$
SchurCheb(4)	1.2 GB	2.4 GB	9.3 GB	2.3 GB
SchurCheb(8)	2.2 GB	4.6 GB	18.8 GB	4.6 GB
PARPACK	21.4 GB	45.0 GB	106.4 GB	46.6 GB

PARPACK with and without shift-and-invert, and the Locally Optimal Block Preconditioned Conjugate Gradient (LOBPCG) method as implemented in the BLOPEX package of `hypr` [11]. The wall-clock times of LOBPCG were obtained with AMG preconditioning and we present the best (lowest) times after performing extensive tests involving various choices for the hyperparameters and preconditioners. Regarding the performance of PARPACK, note that due to the fact that A comes from a 2D discretization, shift-and-invert is generally very fast when the direct solver scales satisfactorily; however, the efficiency of MUMPS falls off faster than that of Algorithm 3.1 as N_p increases, and for larger values of N_p , Algorithm 3.1 becomes the fastest and most scalable approach. Similarly, LOBPCG is competitive with Algorithm 3.1 for smaller values of N_p but becomes comparatively slower as N_p increases.

Figure 5.4 plots the same quantities for a 3D discretization matrix of size $n = 65 \times 64 \times 63$. The main difference between the 2D and 3D case is that PARPACK without shift-and-invert now converges much faster, leading to lower orthogonalization costs. Moreover, because A is banded, the parallel efficiency of distributed-memory sparse matrix-vector products with A remains high even when $N_p = 256$. Nonetheless, Algorithm 3.1 still attains greater strong scaling efficiency than PARPACK (with or without shift-and-invert) and hence will outperform it given enough parallel resources.

As Algorithm 3.1 does not need to factor A , it requires considerably less storage than PARPACK with shift-and-invert. Table 5.2 lists the global peak memory consumption for both of these algorithms for the finite difference discretization problems just described. Even with $N = 8$ Chebyshev nodes, Algorithm 3.1 uses 5 to 10 times less memory than shift-and-invert PARPACK across all problems.

We now focus on the performance of Algorithm 3.1 when the problem size n

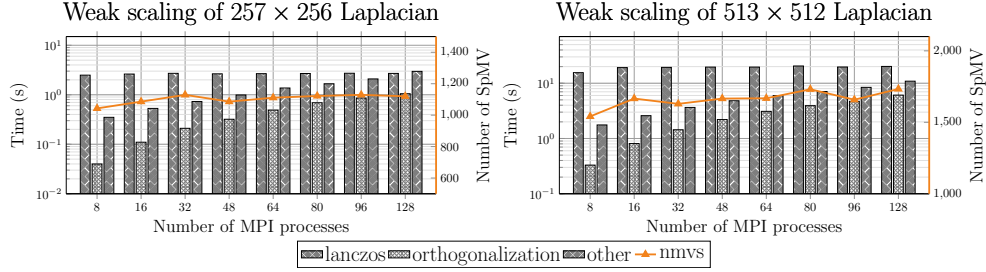


Fig. 5.5: Weak scaling with respect to N ($p_r = 8$, $p_c = N$) for two 2D finite difference discretization problems. The number of MPI processes ranges from $N_p = 8$ to $N_p = 128$. The solid red lines denote the maximum number of iterations required by PARPACK to compute the matrices Y_j , $j = 0, \dots, N - 1$.

and number of partitions p are fixed and N_p varies proportionally to N . We set $p = p_r = 8$ and $p_c = N$ where $N = 2, 4, \dots, 16$. For this experiment, we consider the 2D discretizations of sizes $n = 257 \times 256$ and $n = 513 \times 512$ and report the wall-clock times for each major operation of Algorithm 3.1 in Figure 5.5. The amount of time spent computing the matrices Y_j and V_j is nearly constant since the maximum number of matrix-vector products (iterations) required by PARPACK to compute each Y_j , is more or less the same for each N_p (see the solid lines). On the other hand, the amount of time required for orthonormalization and the Rayleigh–Ritz projection both increase due to: (a) higher computational complexity and (b) higher volume of communication among the increasing number of MPI processes.

Next, we evaluate the performance of Algorithm 3.1 when computing different numbers of eigenvalues (different n_{ev}) for the same matrix. We consider the 2D discretizations of sizes $n = 257 \times 256$ and $n = 513 \times 512$. In each group of tests, we fix p , p_r , p_c , and N_p and then vary n_{ev} . For the $n = 257 \times 256$ problem, we take $N_p = 128$ and $p_r = N$ and then set $p = 16$ when $N = 8$ and $p = 32$ when $N = 4$. For the $n = 512 \times 512$ problem, we double p and N_p . Figure 5.6 reports the total wall-clock times for Algorithm 3.1 under these configurations, taking $n_{ev} = 50, 100, 150, 200$, as well as those for PARPACK (with and without shift-and-invert) and LOBPCG. The cost of solving the Schur complement eigenvalue problems in Algorithm 3.1 at each Chebyshev node increases as n_{ev} increases. Nonetheless, Algorithm 3.1 still attains wall-clock times that are competitive with PARPACK and LOBPCG.

In the preceding experiments, we took $p_c = N$. As our final experiment in this section, we consider the effect of varying the 2D MPI grid topology. We consider the 2D discretizations of sizes $n = 513 \times 512$. We take $N = 8$, $N_p = p = 128$, $n_{ev} = 100$, and vary the topology as $(p_r, p_c) = (128, 1), (64, 2), (32, 4), (16, 8)$. Table 5.3 lists a breakdown of the wall-clock times for the various parts of Algorithm 3.1 for each topology. The topology $(p_r, p_c) = (128, 1)$ processes the N Chebyshev nodes sequentially, one after the other, but uses all N_p MPI processes during the computation of each matrix pair (Y_j, V_j) , $j = 0, \dots, N - 1$, taking on average $(26.08 + 0.35)/8 \approx 3.3$ seconds for each. At the other extreme, the topology $(p_r, p_c) = (16, 8)$ processes the N Chebyshev nodes completely in parallel, but now computing each (Y_j, V_j) requires more time—in the worst case, approximately 4 times as much $(13.21 + 0.35 = 13.56$ seconds)—since only $p_r = 16$ processes are available for parallelization of those computations. Nevertheless, the total time to solution is nearly halved with $(p_r, p_c) =$

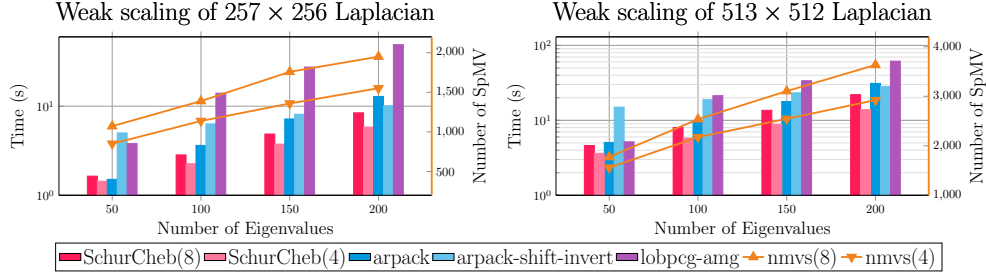


Fig. 5.6: Weak scaling with respect to n_{ev} for two 2D finite difference discretization problems. The number of MPI processes are $N_p = 128$ and $N_p = 256$, respectively. The solid red lines denotes the maximum number of iterations required by PARPACK to compute the matrices Y_j , $j = 0, \dots, N - 1$ in Algorithm 3.1.

Table 5.3: Wall-clock time breakdown of Algorithm 3.1 for various 2D MPI grid topologies. (RR: Rayleigh–Ritz, GS: Gram–Schmidt).

(p_r, p_c)	Setup	$Y_{0,\dots,N-1}$	$V_{0,\dots,N-1}$	GS	RR	DSYGVX	Total
(128,1)	1.42	26.08	0.35	1.41	1.76	0.14	31.17
(64,2)	0.68	18.06	0.36	1.94	1.81	0.14	23.15
(32,4)	0.32	13.95	0.35	1.71	1.91	0.14	18.41
(16,8)	0.18	13.21	0.35	1.65	2.03	0.14	17.61

(16, 8) versus $(p_r, p_c) = (128, 1)$. Thus, in agreement with our previous results, setting $p_c = N$ is best unless the smaller value of p_r creates a memory bottleneck.

5.2.2. Eigenvalue problems from finite element discretizations. To illustrate the performance of Algorithm 3.1 for generalized eigenvalue problems, we again consider matrices arising from discretizations of (5.1) but with linear finite elements instead of finite differences. In 2D, we consider the square $\Omega = (0, 1)^2$ and the disc $\Omega = \{(x, y) : x^2 + y^2 \leq 1\}$, both meshed with unstructured triangular elements. In 3D, we consider the cube $\Omega = (0, 1)^3$, meshed with unstructured tetrahedra.

Figure 5.7 plots the parallel efficiency of Algorithm 3.1 (left) and associated wall-clock times as N_p varies. We also plot the wall-clock time of PARPACK with shift-and-invert but omit results for PARPACK without shift-and-invert, which required an excessive amount of time to converge for these problems. The small sizes of the problems ($n \approx 150,000$) have chosen intentionally in order to simulate an environment with an abundance of parallel resources. As in the experiments of the previous section, Algorithm 3.1 attains high parallel efficiency and scales better than PARPACK. The efficiency of the orthogonalization step in Algorithm 3.1 dropped below 50% for the 3D case when $N_p = 512$ due to a large communication-to-computation ratio for the Gram–Schmidt process; nevertheless, the overall efficiency is still close to 100%.

Next, we show the results of a weak scaling test similar to one in the previous section, wherein Algorithm 3.1 is applied to a given problem for increasing values of n_{ev} . As before, we fix p , p_r , p_c , and N_p for each group of tests, and vary n_{ev} as $n_{ev} = 50, 100, 150, 200$. We use the same finite element problems of the previous experiment set $p_c = N$. When $N = 8$, we use $N_p = 128$ and $p = 16$ for the 2D domains and $N_p = 512$ and $p = 64$ for the 3D domains. When $N = 4$, we double

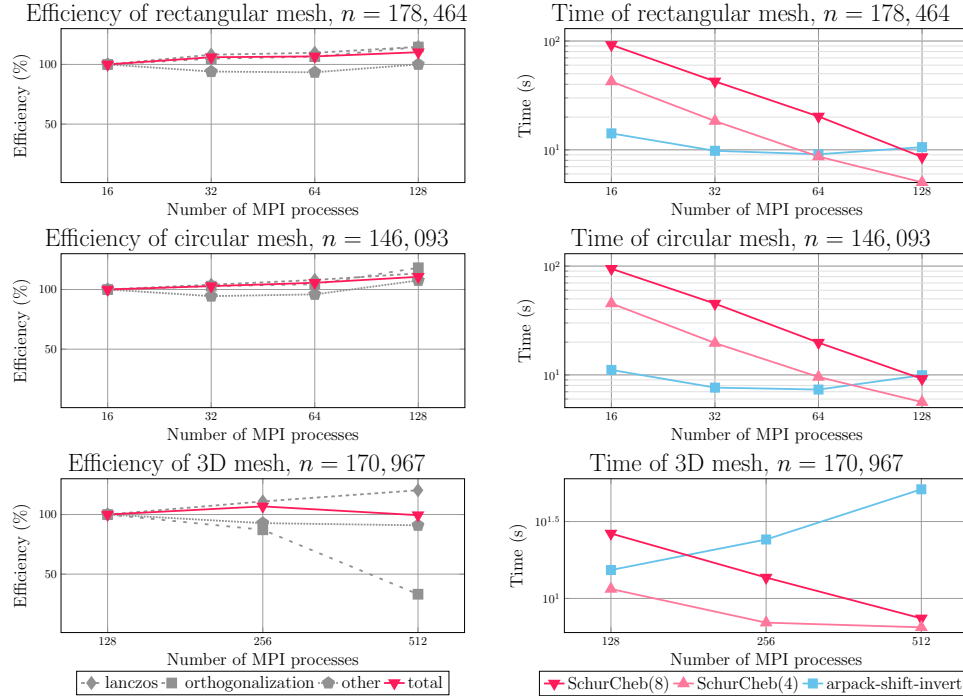


Fig. 5.7: Left: parallel efficiency of Algorithm 3.1 applied to the finite element problems with $n_{\text{ev}} = 100$ and $p_c = N = 8$. Right: wall-clock time comparison between Algorithm 3.1 with $N = 4$ and $N = 8$, and PARPACK with shift-and-invert. The number of MPI processes ranges from $N_p = 8$ to $N_p = 512$. The number of partitions is set equal to $p = 16$ for the 2D meshes and $p = 64$ for the 3D mesh.

p . The results are reported in Figure 5.8. Again, Algorithm 3.1 attains times to solution that are competitive with PARPACK, even though the cost of solving the local eigenvalue problems at each Chebyshev node increases with n_{ev} .

Finally, Table 5.4 lists the wall-clock times for Algorithm 3.1 and PARPACK with shift-and-invert on a set of larger finite element problems. For Algorithm 3.1 we report the wall-clock times for the case $N_p = 512$ and $p_c = N = 4$; for PARPACK, we report the best (lowest) wall-clock time obtained over several runs with different N_p . Algorithm 3.1 was twice as fast for the 2D problems, and about as fast as PARPACK for the 3D problem. Note, though, that in addition to having superior⁷ scalability, Algorithm 3.1 also uses much less memory.

5.2.3. Eigenvalue problems from the SuiteSparse collection. Finally, to demonstrate the performance of Algorithm 3.1 for more general matrices, we apply it to several problems taken from the SuiteSparse matrix collection with sizes ranging from $n = 66,172$ to $n = 1,222,045$. Additional details are given in Table 5.5. The “qa8fk/qa8fm” problem is a generalized eigenvalue problem; the other four are standard problems (M is the identity matrix).

Figure 5.9 plots the parallel efficiency (left) and wall-clock time (right) for Algorithm 3.1 on each of these problems. For comparison, we also plot the wall-clock

⁷The best wall-clock time of PARPACK for the 3D mesh problem was achieved for $N_p = 128$.

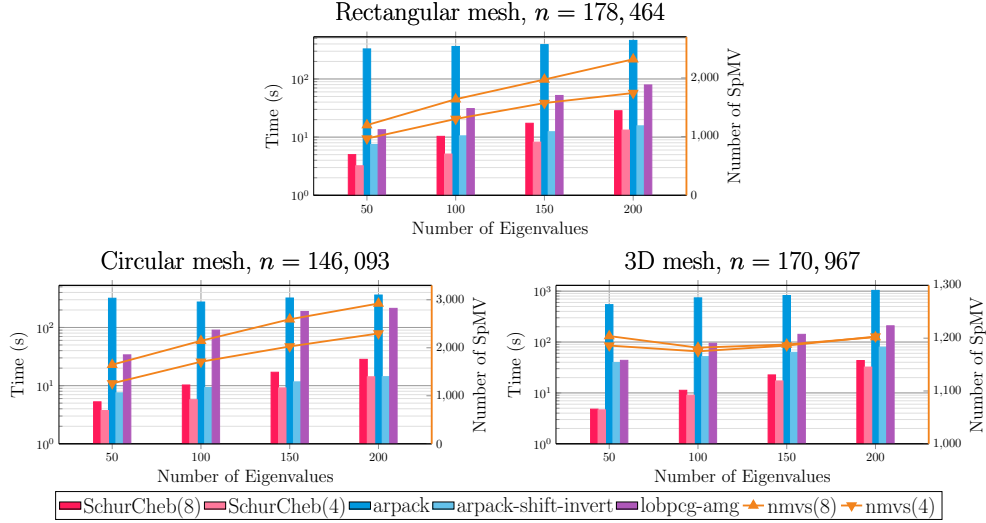


Fig. 5.8: Weak scaling with respect to n_{ev} for three finite element problems. The numbers of MPI processes are $N_p = 128$ for the 2D domains and $N_p = 512$ for the 3D domain. The solid red lines denotes the maximum number of iterations required by PARPACK to compute the matrices Y_j , $j = 0, \dots, N - 1$, in Algorithm 3.1.

Table 5.4: Total wall-clock time for Algorithm 3.1 and PARPACK with shift-and-invert for the finite element problems with $N_p = 512$, $p = 128$, and $p_c = N$.

	2D square $n = 1,086,615$	2D disc $n = 845,397$	3D cube $n = 1,351,083$
SchurCheb(4)	17.2 s	18.3 s	90.1 s
PARPACK	33.6 s	25.9 s	90.3 s

time of PARPACK with and without shift-and-invert. As in the previous experiments, Algorithm 3.1 maintains high parallel efficiency up to 512 MPI processes, and, provided enough parallel resources, outperforms PARPACK. Additionally, Algorithm 3.1 is more memory efficient than shift-and-invert PARPACK as N_p increases; Table 5.6 lists the peak memory consumption for both algorithms for the maximum N_p used in each group of tests for each problem. Finally, Table 5.7 lists the maximum error in the eigenvalues returned by Algorithm 3.1 for $N = 4$ and $N = 8$.

6. Conclusion. We presented a distributed-memory Rayleigh–Ritz projection algorithm to compute a few of the smallest eigenvalues and associated eigenvectors of a sparse, symmetric matrix pencil. The algorithm introduces embarrassing parallelism by recasting the problem as one of approximating univariate, vector-valued functions via Chebyshev approximation. The computational work associated with Chebyshev node can be assigned to a different group of processors, and we described a scheme for doing this using a 2D grid of MPI processes. We discussed several theoretical aspects and implementation details, including how to orthonormalize the Rayleigh–Ritz basis and form the projected eigenvalue problem. Our experiments demonstrated that the proposed algorithm attains good parallel efficiency, superior to PARPACK.

Table 5.5: Problems from the SuiteSparse matrix collection. Here, n denotes the size of the pencil (A, M) ; $\text{nnz}(\cdot)$: counts the number of nonzero entries in its arguments; and p denotes the number of partitions for the case $N = 8$.

Dataset	n	p	$\text{nnz}(A)/n$	$\text{nnz}(M)/n$	Application
qa8fk/qa8fm	66,172	16	25.1	25.1	3D acoustics
af_shell3	504,855	64	34.8	1.0	structural problem
tmt_sym	726,713	64	6.99	1.0	electromagnetics
ecology2	999,999	64	5.00	1.0	2D/3D problem
thermal2	1,228,045	64	6.99	1.0	thermal problem

Table 5.6: Peak memory consumption of Algorithm 3.1 and of PARPACK with shift-and-invert for the SuiteSparse problems.

	qa8 $N_p = 128$	af_shell3 $N_p = 512$	tmt_sym $N_p = 512$	ecology2 $N_p = 512$	thermal2 $N_p = 512$
SchurCheb(4)	0.7 GB	5.9 GB	6.7 GB	8.9 GB	11.2 GB
SchurCheb(8)	1.4 GB	11.9 GB	13.2 GB	17.5 GB	22.2 GB
PARPACK	21.7 GB	47.7 GB	50.8 GB	58.7 GB	56.5 GB

Table 5.7: Maximum relative error in the eigenvalues returned by Algorithm 3.1 for the SuiteSparse problems.

	qa8	af_shell3	tmt_sym	ecology2	thermal2
SchurCheb(4)	3.2×10^{-4}	2.1×10^{-4}	1.6×10^{-4}	1.8×10^{-5}	9.1×10^{-5}
SchurCheb(8)	1.0×10^{-8}	3.8×10^{-10}	6.5×10^{-8}	8.9×10^{-9}	1.9×10^{-10}

In the future, we plan to develop a version of this algorithm based on *generalized spectral Schur complements*, in which the matrix Y_j is formed by computing a few eigenvectors of the pencil $(S(\chi_j), -S'(\chi_j))$ instead of $S(\chi_j)$ alone. This may allow one to reduce the value of N , permitting the use of more parallel resources within each column MPI communicator. We also plan on extending the implementation of our current algorithm so that the computations local to each MPI process are performed using graphics processing units. Finally, we plan on applying our software to problems from real-world applications, e.g., frequency response analysis.

7. Acknowledgements. The authors acknowledge the Minnesota Supercomputing Institute (MSI) at the University of Minnesota for providing resources that contributed to the research results reported within this paper (<http://www.msi.umn.edu>).

REFERENCES

- [1] P. R. AMESTOY, I. S. DUFF, J.-Y. L'EXCELLENT, AND J. KOSTER, *MUMPS: A general purpose distributed memory sparse solver*, in International Workshop on Applied Parallel Computing, Springer, 2000, pp. 121–130.
- [2] E. ANDERSON, Z. BAI, C. BISCHOF, L. S. BLACKFORD, J. DEMMEL, J. DONGARRA, J. DU CROZ, A. GREENBAUM, S. HAMMARLING, A. MCKENNEY, AND D. SORENSEN, *LAPACK Users' Guide*, SIAM, Philadelphia, PA, 1999.
- [3] C. A. BEATTIE, M. EMBREE, AND D. C. SORENSEN, *Convergence of polynomial restart Krylov*

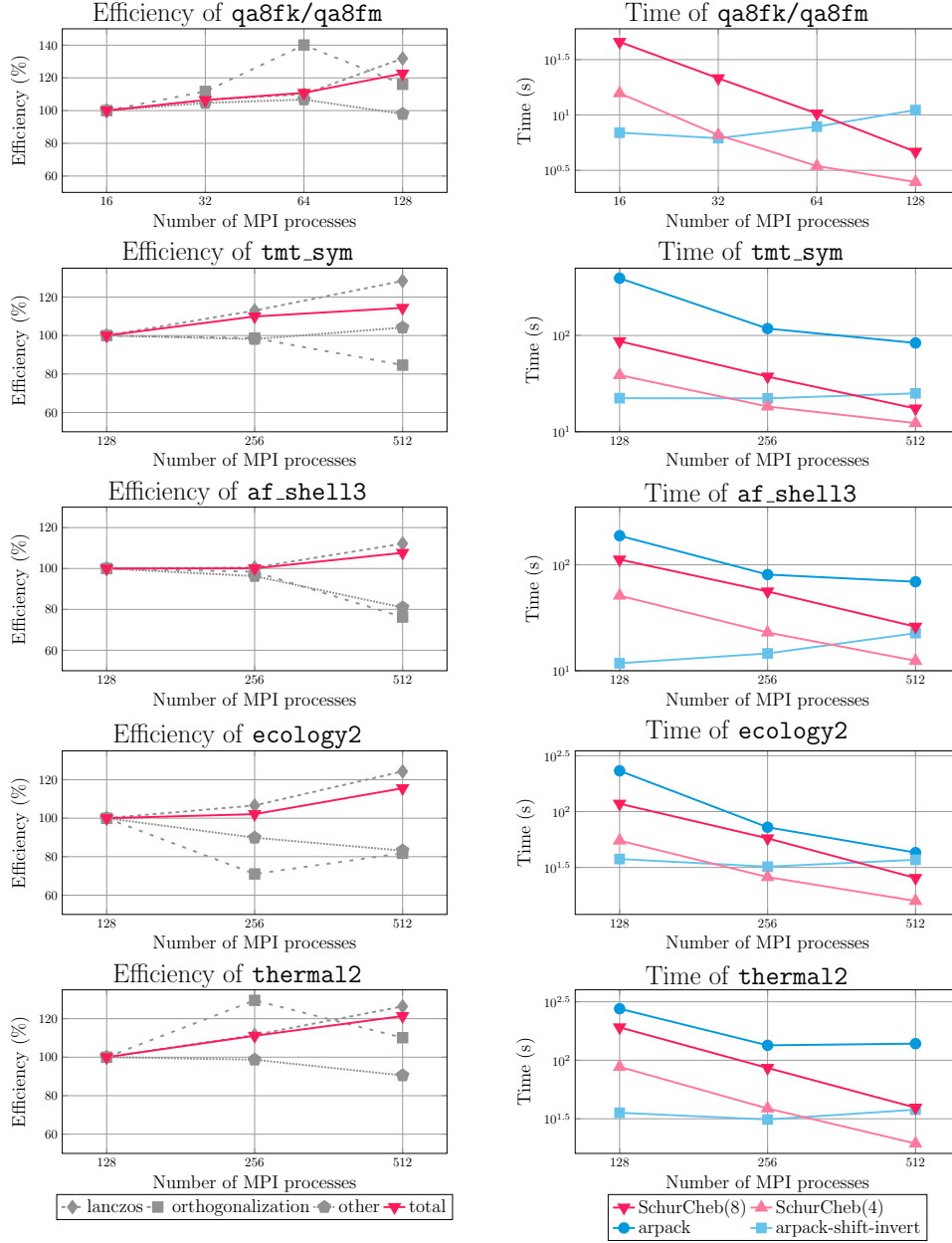


Fig. 5.9: Left: parallel efficiency of Algorithm 3.1 with $n_{ev} = 100$ and $p_c = N = 8$. Right: wall-clock time comparison between Algorithm 3.1 with $N = 4$ and $N = 8$, and PARPACK with and without shift-and-invert. The number of MPI processes ranges from $N_p = 16$ to $N_p = 512$.

- methods for eigenvalue computations, SIAM Rev., 47 (2005), pp. 492–515.
- [4] C. BEKAS AND A. CURIONI, Very large scale wavefunction orthogonalization in density functional theory electronic structure calculations, Comput. Phys. Commun., 181 (2010), pp. 1057–1068.

- [5] C. BEKAS AND Y. SAAD, *Computation of smallest eigenvalues using spectral Schur complements*, SIAM J. Sci. Comput., 27 (2005), pp. 458–481.
- [6] J. K. BENNIGHOF AND R. B. LEHOUCQ, *An automated multilevel substructuring method for eigenspace computation in linear elastodynamics*, SIAM J. Sci. Comput., 25 (2004), pp. 2084–2106.
- [7] J. R. BUNCH AND L. KAUFMAN, *Some stable methods for calculating inertia and solving symmetric linear systems*, Math. Comp., (1977), pp. 163–179.
- [8] D. CALVETTI, L. REICHEL, AND D. C. SORENSEN, *An implicitly restarted Lanczos method for large symmetric eigenvalue problems*, Electron. Trans. Numer. Anal., 2 (1994), p. 21.
- [9] E. CARSON, K. LUND, M. ROZLOŽNÍK, AND S. THOMAS, *Block Gram-Schmidt algorithms and their stability properties*, Linear Algebra Appl., 638 (2022), pp. 150–195.
- [10] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Trans. Math. Software, 38 (2011), pp. 1–25.
- [11] R. D. FALGOUT AND U. M. YANG, *hypre: A library of high performance preconditioners*, in International Conference on Computational Science, Springer, 2002, pp. 632–641.
- [12] W. GAO, X. S. LI, C. YANG, AND Z. BAI, *An implementation and evaluation of the AMLS method for sparse eigenvalue problems*, ACM Trans. Math. Software, 34 (2008), pp. 1–28.
- [13] W. GROPP, E. LUSK, AND A. SKJELLUM, *Using MPI: Portable Parallel Programming with the Message-Passing Interface*, MIT Press, Cambridge, MA, 1999.
- [14] V. HERNANDEZ, J. E. ROMAN, AND V. VIDAL, *SLEPc: A scalable and flexible toolkit for the solution of eigenvalue problems*, ACM Trans. Math. Software, 31 (2005), pp. 351–362.
- [15] W. HEYLEN, S. LAMMENS, AND P. SAS, *Modal analysis theory and testing*, vol. 200, Katholieke Universiteit Leuven Leuven, Belgium, 1997.
- [16] T. HOSHI, H. IMACHI, A. KUWATA, K. KAKUDA, T. FUJITA, AND H. MATSUI, *Numerical aspect of large-scale electronic state calculation for flexible device material*, Jpn. J. Ind. Appl. Math., 36 (2019), pp. 685–698.
- [17] V. KALANTZIS, *Domain decomposition algorithms for the solution of sparse symmetric generalized eigenvalue problems*, PhD thesis, University of Minnesota, 2018.
- [18] V. KALANTZIS, *A domain decomposition Rayleigh–Ritz algorithm for symmetric generalized eigenvalue problems*, SIAM J. Sci. Comput., 42 (2020), pp. C410–C435.
- [19] ———, *A spectral Newton-Schur algorithm for the solution of symmetric generalized eigenvalue problems*, Electron. Trans. Numer. Anal., 52 (2020), pp. 132–153.
- [20] V. KALANTZIS, J. KESTYN, E. POLIZZI, AND Y. SAAD, *Domain decomposition approaches for accelerating contour integration eigenvalue solvers for symmetric eigenvalue problems*, Numer. Linear Algebra Appl., 25 (2018), p. e2154.
- [21] V. KALANTZIS, R. LI, AND Y. SAAD, *Spectral Schur complement techniques for symmetric eigenvalue problems*, Electron. Trans. Numer. Anal., 45 (2016), pp. 305–329.
- [22] V. KALANTZIS, Y. XI, AND L. HORESH, *Fast randomized non-Hermitian eigensolvers based on rational filtering and matrix partitioning*, SIAM J. Sci. Comput., 43 (2021), pp. S791–S815.
- [23] V. KALANTZIS, Y. XI, AND Y. SAAD, *Beyond automated multilevel substructuring: Domain decomposition with rational filtering*, SIAM J. Sci. Comput., 40 (2018), pp. C477–C502.
- [24] G. KARYPIS AND V. KUMAR, *A fast and high quality multilevel scheme for partitioning irregular graphs*, SIAM J. Sci. Comput., 20 (1998), pp. 359–392.
- [25] G. KARYPIS, K. SCHLOEGEL, AND V. KUMAR, *Parmetis: Parallel graph partitioning and sparse matrix ordering library*, (1997).
- [26] T. KATO, *Perturbation Theory for Linear Operators*, Springer-Verlag, New York, NY, 1966.
- [27] J. KESTYN, V. KALANTZIS, E. POLIZZI, AND Y. SAAD, *PFEAST: A high performance sparse eigenvalue solver using distributed-memory linear solvers*, in SC’16: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, IEEE, 2016, pp. 178–189.
- [28] A. V. KNYAZEV, M. E. ARGENTATI, I. LASHUK, AND E. E. OVTCHINNIKOV, *Block locally optimal preconditioned eigenvalue solvers (BLOPEX) in HYPRE and PETSc*, SIAM J. Sci. Comput., 29 (2007), pp. 2224–2239.
- [29] J. H. KO AND Z. BAI, *High-frequency response analysis via algebraic substructuring*, Int. J. Numer. Meth. Eng., 76 (2008), pp. 295–313.
- [30] R. B. LEHOUCQ, D. C. SORENSEN, AND C. YANG, *ARPACK User’s Guide: Solution of Large Scale Eigenvalue Problems with Implicitly Restarted Arnoldi Methods*, SIAM, Philadelphia, 1998.
- [31] R. LI, Y. XI, L. ERLANDSON, AND Y. SAAD, *The eigenvalues slicing library (EVSL): Algorithms, implementation, and software*, SIAM J. Sci. Comput., 41 (2019), pp. C393–C415.
- [32] X. S. LI AND J. W. DEMMEL, *SuperLU-DIST: A scalable distributed-memory sparse direct solver for unsymmetric linear systems*, ACM Trans. Math. Software, 29 (2003), pp. 110–

- 140.
- [33] S. LUI, *Kron's method for symmetric eigenvalue problems*, J. Comput. Appl. Math., 98 (1998), pp. 35–48.
- [34] K. J. MASCHHOFF AND D. SORENSSEN, *A portable implementation of ARPACK for distributed memory parallel architectures*, in Proceedings of the Copper Mountain Conference on Iterative Methods, vol. 1, 1996.
- [35] E. POLIZZI, *Density-matrix-based algorithm for solving eigenvalue problems*, Phys. Rev. B, 79 (2009), p. 115112.
- [36] T. SAKURAI, Y. FUTAMURA, A. IMAKURA, AND T. IMAMURA, *Scalable eigen-analysis engine for large-scale eigenvalue problems*, in Advanced Software Technologies for Post-Peta Scale Computing, Springer, 2019, pp. 37–57.
- [37] A. STATHOPOULOS AND J. R. MCCOMBS, *PRIMME: PReconditioned Iterative MultiMethod Eigensolver—methods and software description*, ACM Trans. Math. Software, 37 (2010), pp. 1–30.
- [38] G. W. STEWART AND J.-G. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, MA, 1990.
- [39] Y. SU, T. LU, AND Z. BAI, *2D eigenvalue problems I: Existence and number of solutions*, arXiv, (2019). [arXiv:1911.08109v1](https://arxiv.org/abs/1911.08109v1) [math.NA].
- [40] L. N. TREFETHEN, *Approximation Theory and Approximation Practice*, SIAM, Philadelphia, 2013.
- [41] U. VON LUXBURG, *A tutorial on spectral clustering*, Stat. and Comput., 17 (2007), pp. 395–416.
- [42] D. B. WILLIAMS-YOUNG, P. G. BECKMAN, AND C. YANG, *A shift selection strategy for parallel shift-invert spectrum slicing in symmetric self-consistent eigenvalue computation*, ACM Trans. Math. Software, 46 (2020), pp. 1–31.
- [43] K. WU AND H. SIMON, *Thick-restart Lanczos method for large symmetric eigenvalue problems*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 602–616.
- [44] Y. XI, R. LI, AND Y. SAAD, *Fast computation of spectral densities for generalized eigenvalue problems*, SIAM J. Sci. Comput., 40 (2018), pp. A2749–A2773.
- [45] C. YANG, W. GAO, Z. BAI, X. S. LI, L.-Q. LEE, P. HUSBANDS, AND E. NG, *An algebraic substructuring method for large-scale eigenvalue calculation*, SIAM J. Sci. Comput., 27 (2005), pp. 873–892.
- [46] H. ZHANG, B. SMITH, M. STERNBERG, AND P. ZAPOL, *SIPs: Shift-and-Invert Parallel Spectral Transformations*, ACM Trans. Math. Software, 33 (2007), pp. 9:1–9:19.