

A GOLUB–KAHAN DAVIDSON METHOD FOR ACCURATELY COMPUTING A FEW SINGULAR TRIPLETS OF LARGE SPARSE MATRICES*

STEVEN GOLDENBERG[†], ANDREAS STATHOPOULOS[†], AND ELOY ROMERO[†]

Abstract. Obtaining high accuracy singular triplets for large sparse matrices is a significant challenge, especially when searching for the smallest triplets. Due to the difficulty and size of these problems, efficient methods must function iteratively, with preconditioners, and under strict memory constraints. In this research, we present a Golub–Kahan Davidson method (GKD), which satisfies these requirements and includes features such as soft-locking with orthogonality guarantees, an inner correction equation similar to Jacobi–Davidson, locally optimal +k restarting, and the ability to find real zero singular values in both square and rectangular matrices. Additionally, our method achieves full accuracy while avoiding the augmented matrix, which often converges slowly for the smallest triplets due to the difficulty of interior eigenvalue problems. We describe our method in detail, including implementation issues that arise. Our experimental results confirm the efficiency and stability of our method over the current implementation of PHSVDS in the PRIMME software package.

Key words. singular value decomposition, Davidson, Jacobi–Davidson, preconditioning, high accuracy, iterative methods

AMS subject classifications. 65F04, 65B04, 68W04, 15A04

DOI. 10.1137/18M1222004

1. Introduction. Assuming a large sparse matrix, $A \in \mathbb{R}^{m,n}$ with $m \geq n$, the economy size singular value decomposition (SVD) is given by

$$(1.1) \quad A = U \Sigma V^T,$$

where $U \in \mathbb{R}^{m,n}$ and $V \in \mathbb{R}^{n,n}$ are orthonormal bases and $\Sigma = \text{diag}(\sigma_1, \dots, \sigma_n) \in \mathbb{R}^{n,n}$ with $\sigma_1 \leq \sigma_2 \leq \dots \leq \sigma_n$ is a diagonal matrix containing the singular values of A . The singular triplets of A are defined as $(\mathbf{u}_i, \sigma_i, \mathbf{v}_i)$, where bold face differentiates from search space vectors and values in this paper. Given the approximate singular triplet (u_i, σ_i, v_i) , we have the left and right singular value residuals, defined as $r_u = A^T u_i - \sigma_i v_i$ and $r_v = A v_i - \sigma_i u_i$.

This decomposition has become increasingly important and is frequently used in fields like statistics for principal component analysis [18], computer science for image compression [28] and web search clustering [26], and genomics for expression data processing [2]. More specifically, finding the smallest singular triplets is useful for total least squares problems, the determination of the effective rank of a matrix [9], and variance reduction of inverse operators [7].

Additionally, finding high accuracy solutions is crucial when running in a single or low precision environment. In single precision, matrix multiplication can only provide $1.2\text{E-}7\|A\|$ of accuracy, and in practice this bound is optimistic for iterative solvers

*Submitted to the journal’s Methods and Algorithms for Scientific Computing section October 22, 2018; accepted for publication (in revised form) April 22, 2019; published electronically July 9, 2019.

<https://doi.org/10.1137/18M1222004>

Funding: This work was supported by the NSF under grant ACI S12-SSE 1440700 and by the DOE under an ECP grant.

[†]Department of Computer Science, College of William and Mary, Williamsburg, VA 23187-8795 (sgoldenberg@email.wm.edu, andreas@cs.wm.edu, eloy@cs.wm.edu).

due to accumulated error. Despite this limitation, single-precision calculations have become increasingly important for deep learning applications [11], which are often resistant to errors and therefore require less than full double precision. Reducing the precision of matrix-vector multiplications (matvecs) can provide speedups on CPUs due to increased vectorization, and GPUs can obtain speedups of 2x–4x [38]. In addition, using single precision cuts the storage requirements in half. Specifically, the use of single precision calculations is encouraged by Advanced Micro Devices (AMD) for OpenCL applications [1, section 2.7.1], and half precision, which can only provide $1\text{E-}3\|A\|$ digits of accuracy, has been growing in popularity on NVIDIA's GPUs [22].

When the matrix A is large enough, it can be inefficient to compute the SVD with dense methods. Furthermore, applications often require only a few of the largest or smallest singular values and vectors. These considerations have led to the use of iterative algorithms like Golub–Kahan–Lanczos (GKL), also known as Lanczos bidiagonalization [8]. However, when the solution requires many iterations, it may be infeasible to store all the GKL vectors necessary for full or partial reorthogonalization. To solve this, restarted versions of GKL that limit the maximum basis size, such as IRLBA [4], have been developed. Additionally, other methods have emerged, such as Jacobi–Davidson (JDSVD) [12], the Preconditioned Hybrid SVD method (PHSVDS) [37], and the Preconditioned Locally Minimal Residual method (PLMR_SVD) [34]. These methods can use the more advanced $+k$ (also known as locally optimal) restarting and can take advantage of preconditioning, which can provide significant speedups for difficult problems.

In general, without preconditioning or $+k$ restarting, these methods build Krylov spaces on the normal equations matrix $C = A^T A$ or on the augmented matrix,

$$(1.2) \quad B = \begin{bmatrix} 0 & A^T \\ A & 0 \end{bmatrix}.$$

We denote a k -dimensional Krylov space on a square matrix A with initial vector v_1 by $K_k(A, v_1) = \text{span}\{v_1, Av_1, \dots, A^{k-1}v_1\}$, and $\|\cdot\|$ denotes the Euclidean norm.

Frequently, methods that build their search space with B , like JDSVD and PLMR_SVD, are able to achieve an accuracy of $\|r_B\| < O(\|A\|\epsilon_{mach})$ when searching for the smallest singular triplets, where ϵ_{mach} is the working machine precision and $r_B = [r_u; r_v]$ is the eigenvalue residual on B . However, B has singular values $\pm\sigma_i$ [27], and so searching for the smallest singular triplets is a highly interior eigenvalue problem that can converge slowly. Worse, when A is rectangular, the spectrum of B contains $m - n$ zero eigenvalues that are not in the spectrum of A . Therefore, methods on B are unable to distinguish real zero singular values of A within the spectrum when $m \neq n$.

Alternatively, methods that build $K_k(C, v_1)$ explicitly are only able to achieve accuracy $O(\|C\|\epsilon_{mach}) = O(\|A\|^2\epsilon_{mach})$ for the eigenvalue residual on C , r_C . If $u = Av/\sigma$ or $r_v = 0$, then r_C is equivalent in exact arithmetic to a scaling of r_u ,

$$(1.3) \quad r_C = A^T Av - \sigma^2 v = \sigma(A^T u - \sigma v) = \sigma r_u.$$

When these methods compute the smallest singular value, and if $\sigma \neq 0$, it is easy to see that (1.3) limits the accuracy with respect to the r_u residual to $O(\|A\|\kappa(A)\epsilon_{mach})$, where $\kappa(A) = \frac{\sigma_n}{\sigma_1}$ is the condition number of A .

Despite the squaring of the spectrum, these methods usually converge faster than methods on B , both in theory and in practice, due to the extremal problem they solve. Furthermore, these methods are often able to find real zero singular values of A , as the corresponding eigenproblem on C does not introduce extraneous zero eigenvalues.

In this work, we introduce a Golub–Kahan Davidson method (GKD), which keeps the convergence of methods on C but attains the full accuracy of methods on B . Specifically, we define full accuracy to be $\sqrt{\|r_u\|^2 + \|r_v\|^2} < \|A\|\epsilon_{mach}$. First, we discuss related methods such as GKL, JDSVD, PLMR.SVD, and PHSVDS, followed by a detailed description of our method, including implementation details. Last, we provide experimental results that highlight the capabilities of GKD compared to the current implementation of PHSVDS in the PRIMME software package.

1.1. Related work. GKL [20] builds two vector bases, one for the right space $K_k(A^T A, v_1)$ and one for the left space $K_k(AA^T, Av_1)$. It builds the second basis while computing the first one without additional matvecs. More importantly, it avoids directly multiplying vectors with $A^T A$ and thus avoids the numerical problems associated with working on C . This is done by keeping two orthogonal spaces, U and V , where the last vector of V , v_k , is used to expand U as $u_k = Av_k$ and the last vector of U , u_k , is used to expand V as $v_{k+1} = A^T u_k$. These new vectors are orthonormalized to the previous ones in their corresponding bases, and the coefficients from this process are used to create the bidiagonal projection matrix $U^T AV$. GKL solves the smaller singular value problem on this projection matrix to approximate the singular triplets.

While GKL is considered to be one of the most accurate and effective algorithms for finding small singular triplets, the standard version is unrestarted and cannot be preconditioned. Therefore, GKL tends to be computationally slow for poorly separated triplets of large matrices. Many restarted versions have been developed [5, 4, 17] but use primarily implicit or thick restarting [35] and thus are unable to maintain the convergence of the unrestarted method. Locally optimal (also known as +k) restarting uses vectors from successive iterations in a way similar to a nonlinear conjugate gradient and has been shown to converge similarly to an unrestarted method for both eigenvalue [19, 32, 31] and singular value problems [37].

SVDIFP [21] implements an inner-outer method where the inner one builds a preconditioned Krylov space $K_k(M(C - \rho_i I), x_i)$, where M is a preconditioner for C and (x_i, ρ_i) is the approximate right singular vector and value at the i th step of the outer iteration. SVDIFP is able to avoid numerical problems, at least for the right singular vectors, by using a two sided projection similarly to GKL. SVDIFP's structure, however, does not allow for many of the optimization techniques of Davidson-type methods which can significantly improve convergence [37].

JDSVD [12] works on B by using two independent subspaces rather than one. It is an inner-outer method that expands both spaces by solving a Jacobi–Davidson-type correction equation on B . Without preconditioning, restarting, or solving the correction equation, the JDSVD outer method builds subspaces that span the following Krylov spaces:

$$(1.4) \quad U_k = K_{\frac{k}{2}}(AA^T, u_1) \oplus K_{\frac{k}{2}}(AA^T, Av_1), \quad V_k = K_{\frac{k}{2}}(A^T A, v_1) \oplus K_{\frac{k}{2}}(A^T A, A^T u_1).$$

These spaces are similar to the ones used in GKL, but crucially, each space is the sum of two different spaces of half dimension. This allows JDSVD to take advantage of initial guesses for both the left and the right singular vectors. However, it also means that the outer solver in JDSVD requires twice as many matvecs to build a space of equal Krylov dimension. Furthermore, if we choose initial vectors that satisfy $v_1 = A^T u_1$, the outer iteration of JDSVD becomes wasteful, as it builds the same space as a GKL with half the dimension (in this case, the spaces $K_{\frac{k}{2}}(A^T A, v_1)$ and $K_{\frac{k}{2}}(A^T A, A^T u_1)$ in (1.4) differ only by one vector). This is also true of eigensolvers

on B as seen below,

$$(1.5) \quad B^2 \begin{bmatrix} v \\ Av \end{bmatrix} = \begin{bmatrix} 0 & A^T \\ A & 0 \end{bmatrix}^2 \begin{bmatrix} v \\ Av \end{bmatrix} = \begin{bmatrix} A^T Av \\ AA^T(Av) \end{bmatrix}.$$

Despite a slower outer iteration, the inner correction equation used in JDSVD is essential to its performance, as it often allows for faster convergence than eigenvalue methods on B while maintaining the ability to converge to full accuracy. However, it can still suffer from the same issues as other eigenmethods on B .

PHSVDS [37] exploits the different advantages of eigenmethods on B and C by utilizing each in a two-stage method. The first stage can use any state-of-the-art eigensolver on C , which gives it fast convergence until either the user tolerance is met or until switching to a second stage using an eigensolver on B is necessary to reach the remaining user tolerance. Switching to an eigensolver on B after a fully converged first stage can effectively utilize good initial guesses from the first stage on C , and thus PHSVDS can avoid resolving the entire accuracy on an indefinite problem. Its implementation in PRIMME can use any of the two near-optimal eigensolvers, GD+k or JDQMR. This two-stage approach has been shown to be faster than eigensolvers on B alone and typically has better performance than other SVD methods.

While PHSVDS has shown significant improvements, it is still limited by the speed of eigensolvers on B when the matrix is ill-conditioned. It converges quite well for problems that do not need to switch stages, but eigensolvers on C cannot converge to high accuracy if the smallest singular value is nearly 0. Once it switches to the second stage on B , a significant slowdown occurs associated with interior problems and methods based on the augmented matrix. In the following sections, GKD demonstrates convergence with the near-optimal speed of GD+k on C down to $O(\|A\|\epsilon_{mach})$.

PLMR_SVD [34] is a recent method based on a stationary iteration that uses two separate four-term recurrences to build the following spaces:

$$\begin{aligned} &\text{span}\{v^{(i)}, r_u^{(i)}, P(A^T r_v^{(i)} - \sigma r_u^{(i)}), v^{(i-1)}\}, \\ &\text{span}\{u^{(i)}, r_v^{(i)}, P(Ar_u^{(i)} - \sigma r_v^{(i)}), u^{(i-1)}\}, \end{aligned}$$

where $v^{(i)}$ and $u^{(i)}$ are the i th approximations of the right and left singular vectors, respectively, and $r_v^{(i)} = P(Av^{(i)} - \sigma u^{(i)})$ and $r_u^{(i)} = P(A^T u^{(i)} - \sigma v^{(i)})$ are their preconditioned right and left residuals, respectively. Without a preconditioner, PLMR_SVD is equivalent to GD+1 with a 3-vector basis (or LOBPCG) on B . There may be additional benefits to building the spaces separately, but PLMR_SVD lacks the subspace acceleration present in GD+k and JDSVD, which can provide superlinear convergence.

2. Main contribution. In this section, we describe the proposed method, GKD, in detail, especially focusing on the selection of approximate singular triplets from our subspaces and the implementation of our restarting method. Additionally, we discuss error accumulations that occur due to restarting and the mitigation strategy required to ensure reliable performance for high accuracy calculations. Finally, we extend GKD to an inner-outer method that solves a Jacobi–Davidson correction equation.

2.1. Algorithm. Our algorithm is designed to mimic the numeric nature of GKL by keeping two orthonormal bases for the right and left spaces, V and Q , respectively, which are built without multiplying directly with $A^T A$. Instead, we build Q such that $AV = QR$ is the economy QR factorization of AV . Then, we extend V with a left residual based on a Galerkin extraction from R . Without preconditioning or

+k restarting, this process is identical to GKL, building the right and left spaces $K_q(A^T A, v_1)$ and $K_q(AA^T, Av_1)$ after q iterations or $2q$ matvecs. Since both the extraction of approximate triplets through the SVD of R and the expansion of the spaces avoid a direct multiplication with C , we avoid the squaring of the norm and condition number that occurs with eigensolvers on C .

Specifically, we extract approximate singular triplets from these spaces using a Rayleigh–Ritz procedure that is adapted for the SVD. Given search spaces $\mathcal{Q} \subset \mathbb{R}^m$ and $\mathcal{V} \subset \mathbb{R}^n$, we can determine approximations (u, σ, v) with the following two Galerkin conditions on the right and left residuals:

$$(2.1) \quad \begin{aligned} Av - \sigma u &\perp \mathcal{Q}, \\ A^T u - \sigma v &\perp \mathcal{V}. \end{aligned}$$

Since $u \in \mathcal{Q}$ and $v \in \mathcal{V}$, we can write $u = Qx$ and $v = Vy$, where Q and V form k -dimensional orthonormal bases of \mathcal{Q} and \mathcal{V} , respectively. Additionally, $AV = QR \Rightarrow Q^T AV = R$, which allows us to rewrite the conditions as follows:

$$(2.2) \quad \begin{aligned} Q^T AVy &= \sigma Q^T Qx \Rightarrow Ry = \sigma x, \\ V^T A^T Qx &= \sigma V^T Vy \Rightarrow R^T x = \sigma y. \end{aligned}$$

Therefore, solving the SVD on R with singular triplets (x, σ, y) satisfies both constraints and provides approximations to the singular triplets of A .

To expand the right search space, we take the approximations from the above Rayleigh–Ritz extraction and use them to form the left residual $r_u = A^T u - \sigma v$. Then, we can choose to expand V with this r_u directly, or with the preconditioned residual Pr_u , where P is a suitable preconditioner for $A^T A$ or for $A^T A - \sigma I$, if available.

We expand the left space \mathcal{Q} with Av_{i+1} instead of a preconditioned right residual. This differentiates the method from JDSVD with the goal of producing a faster converging outer method. Specifically, from (1.3) the left residual r_u is colinear with the residual r_C of the Generalized Davidson (GD) method [24] on the matrix C , which is also colinear with the new GKL direction for V . In addition, the Rayleigh–Ritz on C used by GD gives the same answer as (2.2),

$$V^T A^T AVy = \sigma y \Rightarrow R^T Ry = \sigma y,$$

and so, in exact arithmetic, GKD is equivalent to GD solving the eigenproblem on $A^T A$. Without preconditioning or restarting, it is also equivalent to GKL, and thus it is twice as fast as JDSVD if the latter is used only as an outer method. By construction, GKD has numerical properties similar to those of GKL, whereas the accuracy of GD is limited by working directly on $A^T A$. GKD can also be used with thick and +k restarting, which in exact arithmetic makes it equivalent to GD+k on C , the first stage method of PHSVDS, but without the numerical limitations. Algorithm 2.1 shows the restarted and preconditioned version of GKD when seeking one singular triplet. Although the orthogonalization of step 13 can be avoided without preconditioning [29], it is needed for high accuracy and allows our more general method to use flexible preconditioning. Furthermore, the algorithm can be extended to find more than one singular triplet by using soft- or hard-locking. A block version is similarly possible.

2.2. Restarting and locking. Our restart procedure takes the current best approximations to the s singular triplets closest to the user specified target, τ , and

Algorithm 2.1 GKD Iteration.

```

1: Define target  $\tau$ , initial vector  $v_1$ , max basis size  $q$ , tolerance  $\delta$ , preconditioner  $P$ ,
   and  $i = 1$ 
2: Build  $V = [v_1]$ ,  $Q = [\frac{Av_1}{\|Av_1\|}]$ ,  $R = \text{zeros}(q, q)$ , and  $R(1, 1) = \|Av_1\|$ 
3: while  $\sqrt{\|r_u\|^2 + \|r_v\|^2} > \|A\|\delta$  do
4:   while  $i < q$  do
5:     Compute SVD of  $R$ 
6:     Choose the singular triplet  $(x, \sigma_r, y)$  of  $R$  nearest to the target  $\tau$ 
7:     Save  $v_{old} = y$  for +k restarting
8:     Set  $u = Q(:, 1:i)x$ ,  $v = V(:, 1:i)y$ 
9:     Compute left residual:  $r_u = A^T u - \sigma_r v$ 
10:     $V(:, i+1) = Pr_u$ 
11:    Orthogonalize  $V(:, i+1)$  against  $V(:, 1:i)$ 
12:     $Q(:, i+1) = AV(:, i+1)$ 
13:    Orthogonalize  $Q(:, i+1)$  against  $Q$  and update  $R(:, i+1)$ 
14:     $i = i + 1$ 
15:   end while
16:   call Algorithm 2.2 to restart
17: end while

```

Algorithm 2.2 Restart Procedure.

```

1: Define restart size  $s$  and target  $\tau$ 
2: Compute SVD of  $R = X\Sigma_r Y^T$ 
3: Choose  $s$  singular triplets of  $R$  closest to  $\tau$  (called  $(X_1, \Sigma_r^{(1)}, Y_1)$ )
4: Save the remaining singular triplets from the SVD of  $R$ ,  $(X_2, \Sigma_r^{(2)}, Y_2)$ 
5:  $v_{new} \leftarrow$  Orthogonalize saved +k vectors  $[v_{old}; 0]$  from main iteration against  $Y_1$ 
6:  $t = [Y_1, v_{new}]$ 
7:  $V = Vt$ 
8: if Reset criteria is met then
9:   Reorthogonalize  $V$  and build  $Q$  and  $R$  such that  $AV = QR$ 
10: else
11:   QR factorize  $\Sigma_r^{(2)} Y_2^T v_{old} = \tilde{Q}\tilde{R}$ 
12:   Set  $Q = Q[X_1 X_2 \tilde{Q}]$  and  $R = \begin{bmatrix} \Sigma_r^{(1)} & 0 \\ 0 & \tilde{R} \end{bmatrix}$ .
13: end if

```

uses them together with those from the +k restarting to compress V , Q , and R down to dimension $s+k$. The steps for building the restarted V follow closely the description in [31] and are shown in lines 1–7 of Algorithm 2.2.

The simplest method to restart Q and R , without recomputing the QR factorization of the restarted AVt , is to set them as $Q\tilde{Q}$ and \tilde{R} , respectively, where $\tilde{Q}\tilde{R} = Rt$ is the QR factorization of Rt with $t = [Y_1, v_{new}]$ from line 6 of Algorithm 2.2. This can introduce numerical error of magnitude $O(\|R\|\epsilon_{mach})$, which can be as large as $O(\|A\|\epsilon_{mach})$. Although this error is acceptable for a single QR factorization, the error accumulates over many restarts, causing the factorization not to correspond to the actual AV and eventually causing loss of convergence. It is possible to intelligently compute Q and R to avoid direct multiplications with R through the already

available SVD of R as seen below,

$$\begin{aligned}
 (2.3) \quad AVt = QRt &= Q \begin{bmatrix} X_1 & X_2 \end{bmatrix} \begin{bmatrix} \Sigma_r^{(1)} & 0 \\ 0 & \Sigma_r^{(2)} \end{bmatrix} \begin{bmatrix} I & 0 \\ 0 & Y_2^T v_{old} \end{bmatrix} \\
 &= Q \begin{bmatrix} X_1 & X_2 \end{bmatrix} \begin{bmatrix} \Sigma_1 & 0 \\ 0 & \Sigma_r^{(2)} Y_2^T v_{old} \end{bmatrix}.
 \end{aligned}$$

From (2.3), the new Q and R can be obtained with minimal effort by performing a QR factorization on $\Sigma_r^{(2)} Y_2^T v_{old} = \tilde{Q} \tilde{R}$. The restarted Q and R are given in line 12 of Algorithm 2.2. This strategy has better numerical behavior because we separate the space of small singular values that are kept in thick restarting (X_1) from the +k restarting space which has correction directions over the entire singular space (including those of large magnitude). By explicitly decoupling $\Sigma_r^{(1)}$ and \tilde{R} in R , any errors in \tilde{R} do not affect the ability of the algorithm to compute the smallest eigenvectors and they only affect the correction directions. Moreover, as the +k algorithm typically uses only $k = 1$ previous vectors, no errors are expected.

To accurately find many singular triplets, we implement two versions of locking. The first, hard-locking, locks singular vectors out of the search space explicitly once the required user tolerance is reached. At every iteration, we orthogonalize the vector added to V against the locked right singular vectors, as well as the previous vectors in V . In practice, the vectors added to Q do not require orthogonalization against the locked left singular vectors. The second, soft-locking, merely flags converged singular triplets while leaving them in the basis.

It is known that hard-locking can cause stagnation in some rare cases or when the number of locked vectors is large. This is caused by the error still present in the locked vectors, which may contain critical directions for other singular triplets [30]. We have not seen any matrices in this paper that exhibit this behavior. However, soft-locking can provide left and right singular vectors that are orthogonal to machine precision, while hard-locking only obtains left singular vectors orthogonal up to $O(\|A\|\delta)$. Therefore, we present only soft-locking results in this paper. We intend to address the issues with hard-locking more thoroughly in the future.

2.3. Resetting. Since $AV = QR$, the right residual $r_v = Av - \sigma u$ should be zero throughout our procedure,

$$(2.4) \quad r_v = Av - \sigma u = AVy - Q(\sigma x) = AVy - QRy = (AV - QR)y = 0.$$

Generally, this means we can avoid the extra matrix-vector multiplication (or storage for AV) necessary to compute r_v . In practice, though, $\|r_v\|$ cannot be better than $O(\|A\|\epsilon_{mach})$ due to the multiplication of AV when computing the left space. Worse, $\|r_v\|$ grows as $O(\sqrt{\text{numRestarts}}\|A\|\epsilon_{mach})$, which has also been noticed in [36]. Therefore, our method must calculate $\|r_v\|$ explicitly when $\|r_u\| < \|A\|\delta$, where δ is the user selected tolerance. This ensures we meet the convergence criteria of Algorithm 2.1.

The errors we observe in r_v may grow large enough to exceed the user tolerance, which would make convergence impossible. These errors come from two main sources. The first source is from the loss of orthogonality of V , and the second is the loss of accuracy of the QR factorization of AV . We have found experimentally that both of these errors can impede or halt convergence, as the SVD of R no longer corresponds to the singular triplets in A . We note that this issue is rare and only occurs when

$\delta \approx \epsilon_{mach} \sqrt{\text{numRestarts}}$. To correct these errors, we implement a resetting procedure that reorthogonalizes V and that rebuilds Q and R directly from a newly computed AV . It is critical to only reset sparingly, as rebuilding Q and R from scratch takes $s + k$ matvecs to obtain AV and a full QR factorization.

Additionally, we have noticed in our experiments that resetting can cause an increase in the residual norm up to $\|A\|\kappa(A)\epsilon_{mach}$, which may require a few iterations to reduce back to its previous level. This can be seen by analyzing a reset of the method when $V = \mathbf{v}_1$, the exact singular vector. In this situation, after resetting the new $\tilde{q} = (A\mathbf{v}_1 + \mathbf{e})/\|A\mathbf{v}_1 + \mathbf{e}\|$ contains an error due to the matrix-vector multiplication, with $\|\mathbf{e}\| \leq \|A\|\epsilon_{mach}$. When $\kappa(A)$ is not too large, the error term \mathbf{e} in the nominator will be amplified by a factor of $1/(\sigma_1 - \|\mathbf{e}\|) \approx 1/\sigma_1$ at most. Thus, we expect the error in \tilde{q} to be at most $\kappa(A)\epsilon_{mach}$. This means that the norm of the left residual norm, which we compute explicitly, can increase up to $\|A\|\kappa(A)\epsilon_{mach}$. In order to track the errors mentioned above, we have devised two inexpensive criteria that help to avoid unnecessary resets.

Since the error in the orthogonality of V may cause convergence issues, we must estimate how large $\|E\| = \|V^T V - I\|$ can be before it needs correction. To do so, we analyze the Galerkin condition on the equivalent eigenproblem on C , i.e., $A^T A v = \sigma^2 v$. Applying the Galerkin condition with V , the projected eigenproblem should have been $V^T A^T A V \tilde{y} = \tilde{\sigma}^2 V^T V \tilde{y}$. However, our algorithm solves $V^T A^T A V y = \sigma^2 y$ regardless of the orthonormality of V . Therefore, we obtain a Ritz vector $v = Vy$ and Ritz value σ^2 that will not converge to a zero residual. The Ritz pair produced by our inexact Galerkin can be considered as a Ritz pair of an exact Galerkin condition applied to the nearby generalized eigenproblem $A^T A \tilde{v} = \lambda M \tilde{v}$, where $M = V(V^T V)^{-2} V^T$. This can be seen by using the Galerkin condition with V , which yields the same Ritz vector Vy and Ritz value $\lambda = \sigma^2$,

$$(2.5) \quad V^T A^T A V y = \lambda V^T M V y = \lambda V^T V (V^T V)^{-2} V^T V y = \lambda y.$$

In order to correctly monitor and maintain convergence, the residual $r_C = \sigma r_u = A^T A v / \|v\| - \sigma^2 v / \|v\|$ should not drift too far from the exact residual of the generalized eigenproblem, $r_E = A^T A v / \|v\| - \sigma^2 V(V^T V)^{-2} V^T v / \|v\|$, where $\|v\| = \|Vy\|$ since Vy may not be unit length. However, from [33, 13], $\|Vy\| \geq \sigma_{\min}(V) \geq \sqrt{1 - \|E\|}$ and $\|I - (V^T V)^{-1}\| \leq \|E\| / (1 - \|E\|)$. Therefore, assuming $\|E\| < 1$, we have

$$(2.6) \quad \begin{aligned} \|r_E - r_C\| &= \sigma^2 \left\| \frac{Vy}{\|Vy\|} - \frac{V(V^T V)^{-1}y}{\|Vy\|} \right\| \\ &\leq \frac{\sigma^2 \|V\| \|I - (V^T V)^{-1}\|}{\sqrt{1 - \|E\|}} \\ &\leq \sigma^2 \|V\| \frac{\|E\|}{(1 - \|E\|)^{3/2}} \\ &\leq \sigma^2 (1 + \|E\|) \left(\sum_{i=0}^{\infty} \|E\|^{i+1} (-1)^i \binom{-\frac{3}{2}}{i} \right) \\ &\leq \sigma^2 (\|E\| + O(\|E\|^2)). \end{aligned}$$

Since we want $r_u = r_C / \sigma$ to converge to tolerance $\|A\|\delta$, we want the distance $\|r_E - r_C\| < \|A\|\delta\sigma$. Thus, from (2.6), we should perform a reset when $\|E\| \geq \|A\|\delta/\sigma$. In practice, the second criterion described below will be satisfied earlier than (2.6), and therefore $\|E\|$ does not need to be computed explicitly.

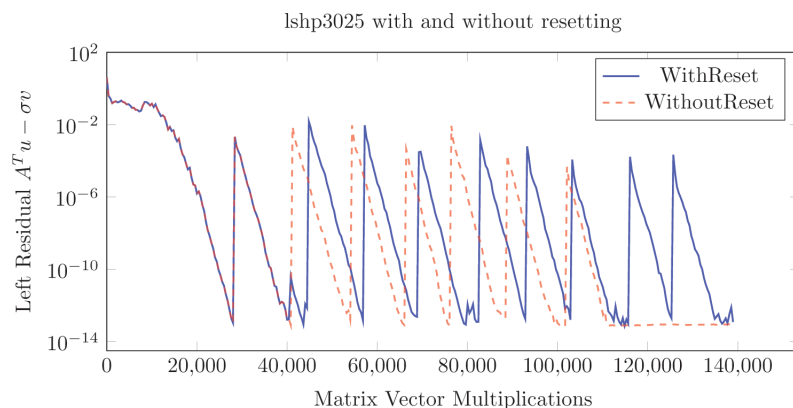


FIG. 1. Demonstrating the need for resetting on *lshp3025* ($\|A\| = 7$) with GKD ($q = 35$, $s = 15$, $\delta = 1E-14$, and $k = 1$).

The accuracy of the QR factorization also directly impacts the convergence of r_u . From (2.4), we can estimate errors in the QR factorization directly from the norm of the right residual. We choose to reset when $\|r_u\| < 1.25\|r_v\|$, which includes a small 25% buffer that we have found is necessary to detect potential stagnation in a few experimental cases. Since $\|r_v\|$ grows at a rate of approximately $O(\sqrt{\text{numRestarts}}\|A\|\epsilon_{mach})$, explicit computation of r_v is unnecessary.

To demonstrate this problem, we ran *lshp3025*, a problem from the SuiteSparse Matrix Collection [6], which requires thousands of restarts before convergence. Properties of this problem can be found in Table 1. The criteria outlined in the previous paragraphs combine to avoid the stagnation seen in Figure 1. Due to the very low tolerance of $1E-14 = 50 * \epsilon_{mach}$, approximately 2,500 restarts or 35,000 matvecs may cause the reset criteria to be met. It is clear that our criteria are somewhat conservative, as resets occur approximately every 40,000 matvecs, even when the method is able to converge without it. However, without resetting, the method completely stagnates at around 110,000 matvecs. Moreover, with or without resets, we observe convergence to the first eight smallest singular values in a similar number of matvecs (110,000), even though adding resets should increase the overall number of matvecs. This indicates that the increased stability of the method also can improve performance slightly.

2.4. Inner solver. Inner-outer solvers like JDSVD and the JDQMR implementation in PRIMME utilize extra matvecs inside of an inner solver as a refinement step to improve the convergence speed of the outer iterations. By solving a related linear system, these methods can provide a significant speedup in time for problems that have a relatively inexpensive matrix-vector multiplication. Furthermore, solving this linear system can reduce the residual of the solution without requiring the expansion of the outer basis. Consequently, the number of orthogonalizations as well as the number of restarts are reduced, which avoids their associated error and resets. This is particularly critical for problems that require a significant number of iterations.

GKD can be extended to a Jacobi–Davidson variant, GKJD, that expands the subspace V by the approximate solution of the correction equation

$$(2.7) \quad (I - vv^T)(A^T A - \sigma^2 I)(I - vv^T)t = -r_u$$

instead of applying a preconditioner at line 10 of Algorithm 2.1. Here, and for the remainder of this section, σ without a subscript denotes the shift used for the inner solver, which may be different from the user specified target τ or the current approximate singular value. As before, σ_i will denote the i th singular value.

The inner equation can also utilize a preconditioner, improving convergence further. In particular, our inner solver is based on the symmetric Quasi-Minimal Residual method (QMRs) used in PRIMME's JDQMR. QMRs benefits from the ability to utilize indefinite preconditioners and solve indefinite systems which may occur when σ lies in the interior of the spectrum.

In order to avoid overutilizing the inner method when convergence is poor or the correction equation does not match the desired singular values, or underutilizing the inner method when convergence is good, extra steps must be taken. There is a significant volume of research on stopping criteria for inner iterations, including results for Jacobi–Davidson-type methods for eigenvalue and SVD problems [14, 16, 15, 31, 25]. In this paper, we adopt the dynamic criteria used in PRIMME's JDQMR [31] which take advantage of the smooth convergence of QMRs to estimate the relevant eigenvalue residuals and stop the linear solve in a near-optimal way. Of course, other stopping criteria can also be implemented or an entirely different inner solver may be used.

The inner solver for (2.7) works directly on $A^T A - \sigma^2 I$, and so its numerical stability needs to be justified. As with an outer iteration on $A^T A$, no numerical issues are expected when σ is in the largest part of the spectrum, but when seeking the smallest part, singular values below $O(\|A\|\sqrt{\epsilon_{mach}})$ will become indistinguishable when squared. However, the solution of the inner correction equation still provides useful directions even when a few singular values of A are below $O(\|A\|\sqrt{\epsilon_{mach}})$. The reason is well understood numerically, and it is why inverse iteration works well despite a nearly singular linear system [27, sect. 4.3].

Assume there are k singular values below the noise level, i.e., $\sigma_k \leq \|A\|\sqrt{\epsilon_{mach}} < \sigma_{k+1}$, and a shift $\sigma \leq \|A\|\sqrt{\epsilon_{mach}}$. If we ignore the projectors for simplicity, the numerically computed solution of (2.7), \tilde{t} , satisfies

$$(2.8) \quad \tilde{t} = t + \mathbf{V}(\Sigma^2 - \sigma^2)^{-1} \mathbf{V}^T E \tilde{t},$$

where the backward error satisfies $\|E\| \leq \|A^T A\| \epsilon_{mach}$. Therefore, the relative forward error is a vector $\frac{\tilde{t}-t}{\|\tilde{t}\|} = \sum_{i=1}^n \mathbf{v}_i c_i$ with the coefficients satisfying

$$(2.9) \quad |c_i| = \frac{|\mathbf{v}_i^T E \tilde{t}|}{|\sigma_i^2 - \sigma^2| \|\tilde{t}\|} \leq \frac{\|A\|^2 \epsilon_{mach}}{|\sigma_i^2 - \sigma^2|}.$$

For $i > k$, we have $\sigma_i \geq \sigma_{k+1} > \|A\|\sqrt{\epsilon_{mach}}$, and thus $|c_i| = O(\frac{\|A\|^2}{\sigma_i^2} \epsilon_{mach}) < 1$. As the separation increases, $\sigma_{k+1} \gg \|A\|\sqrt{\epsilon_{mach}}$, we have $c_i \ll 1$ and the errors in the \mathbf{v}_i , $i > k$, directions become negligible. For $i \leq k$, we have $|\sigma_i^2 - \sigma^2| < \|A\|^2 \epsilon_{mach}$, and thus the corresponding c_i could blow up. In practice, calculations at the noise level of the arithmetic will limit $c_i = O(1)$, but either way these \mathbf{v}_i , $i \leq k$, directions dominate the correction vector.

The behavior is similar when the backward error is at the level of the residual norm at which we solve (2.7), i.e., $\|E\| \leq \|A\|^2 \theta$, for some tolerance θ . Typically we ask for a residual norm reduction relative to $\|r_u\|$, but this can be translated to a θ . Then, the $|c_i|$ in (2.9) have the same bounds as above only multiplied by θ/ϵ_{mach} . Since the approximate solution has $\|t\| = O(\theta)$, the effect of the noise error is larger.

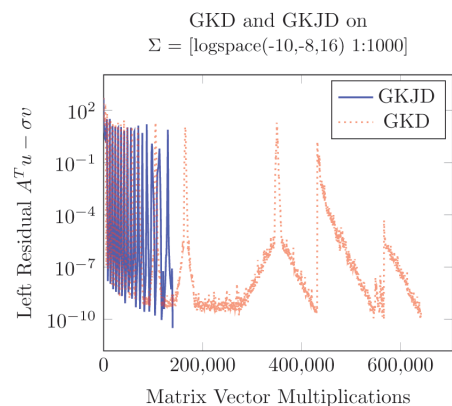


FIG. 2. Convergence of GKD and GKJD when there are more SVs below $\sqrt{\epsilon_{mach}}$ than the $MaxBasisSize$ ($q = 35$, $s = 15$).



FIG. 3. Convergence of GKJD on a problem with 20 SVs below $\sqrt{\epsilon_{mach}}$ in single precision with varying minimum restart sizes (maximum $matvecs = 75,000$, $q = 50$).

We can view the noise of the numerically computed correction \tilde{t} as the application of a low pass filter with the diagonal matrix $\text{diag}(c_i)$, where the $i < k$ singular components dominate the result. Clearly, the inner iteration cannot differentiate between these k smallest singular directions which look like a multiplicity. However, the Rayleigh–Ritz of the outer method has no problems approximating these singular vectors as long as their k -dimensional space is sufficiently represented in the outer search space.

If the outer method in GKJD has a restart size $s \geq k$ and the gap σ_{k+1}/σ_k is large, then the filter ensures that all v_i , $i = 1, \dots, k$, will be approximated well after k outer iterations. As the gap narrows, the filter also boosts directions of larger singular values up to σ_f , where $\frac{\|A\|^2}{\sigma_f^2} \epsilon_{mach}$ starts to become negligible. Therefore, the outer method may take more than k iterations, although convergence depends on the gaps in the “filtered” $\sigma_1, \dots, \sigma_f$ spectrum, which has a much smaller spread than the entire spectrum.

The situation is similar if the restart size $s < k$ and σ_{k+1}/σ_k is large, since the search space cannot capture all small singular vectors, and so convergence will occur based on the perceived gaps after the implicit application of the filter. In the extreme case of $s \ll k$ and/or very small spectral gaps, we can expect the method to be slow. However, in such ill-conditioned problems, no better algorithmic options exist without a preconditioner.

Figures 2 and 3 show examples of how GKJD with dynamic stopping conditions for the inner iteration can converge even when several singular values are below $\|A\|\sqrt{\epsilon_{mach}}$. They also show that GKJD is competitive and sometimes faster than GKD in terms of matrix-vector products, in addition to the benefit of a less expensive iteration. The matrices have a specified spectrum Σ and random left and right singular vectors.

In Figure 2, the matrix has 16 singular values below $\|A\|\sqrt{\epsilon_{mach}}$, but we limit GKD and GKJD to a restart size of only 15. Even with this limitation, GKJD is able to converge to the smallest singular triplet with a relative accuracy of $1E-14$, and it does so three times faster than GKD. Additionally, with only a few extra outer iterations, GKJD can find 14 of the smallest singular values.

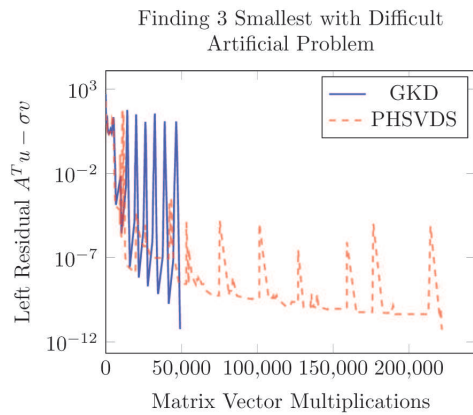


FIG. 4. Convergence of PHSVDS on a poorly conditioned problem ($\kappa(A) = 1E + 13$).

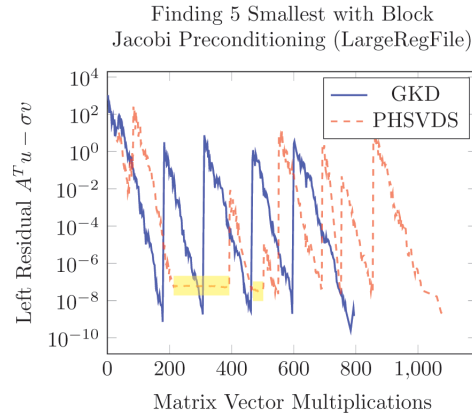


FIG. 5. Stagnations caused by a failure to fully converge in the first stage of PHSVDS ($\kappa = 1.1E + 4$).

The difference seen between GKD and GKJD is due to the large number of restarts for GKD and their associated error. As the errors caused by restarts grows above the relative tolerance within approximately 2,000 restarts (40,000 matvecs), GKD may have numerical issues and not converge, although this behavior is sensitive to the choice of random orthonormal bases U and V . Since GKJD performs orders of magnitude fewer outer iterations, it is not affected by this source of error heavily and therefore is not sensitive to the random left and right singular spaces. With a marginally less strict tolerance, GKD does not exhibit this behavior.

In Figure 3, we consider an example where the matrix has 20 singular values below the $\|A\|\sqrt{\epsilon_{mach}}$ threshold. We use single precision arithmetic, which allows for relatively larger spectral gaps that make convergence tractable. We search for the smallest singular value with a maximum basis size of 50, the dynamic inner stopping criteria, and a tolerance of $1E-5$ for all tests while varying the restart size used by GKD and GKJD. We see that smaller restart sizes do not impede convergence of GKJD and only slow it down by less than a factor of two. However, the effects of a small restart size are much more severe on GKD, which is unable to converge to the desired tolerance within 75,000 matvecs for restart sizes less than 10. This shows that GKJD is able to rebuild the space lost during restarting much more quickly than GKD, as the inner equation can sufficiently filter out directions corresponding to the unwanted portions of the spectrum.

3. Benefits over PHSVDS.

3.1. Avoiding the augmented problem. As mentioned earlier, methods on B often exhibit problems due to the interior nature of the spectrum that they work on. In order to demonstrate these issues, Figure 4 shows convergence on the problem $A = \text{diag}([1e-10, 2e-10, 5e-10, 1e-9, 3e-9, 1e-8, 1e-6, 1e-4, 1:1000])$. First, this problem is very poorly conditioned ($\kappa(A) = 1E13$), and since the six smallest singular values are below $1E-8$, the first stage of PHSVDS is unable to distinguish them from zero. Second, because the spectrum is reflected across 0 for the augmented problem, it is very difficult to converge only to the positive part of the spectrum.

In searching for three singular values to a user tolerance of $1E-14$, PHSVDS

took more than four times as many matvecs, but more importantly, it missed five smaller singular values, as the third converged value was $1\text{e-}4$. Even worse, the vectors that were returned for left and right spaces were not orthogonal, as $\|Q^T Q - I\| \approx \|V^T V - I\| \approx 6\text{E-}5$. Therefore, the true residuals after orthogonalization did not meet the full user tolerance. Comparatively, GKD converges to all six of the smallest singular values and did so with fully orthogonal left and right vectors. As we can see from the figure, the convergence for GKD is fairly smooth, converging to each of the six singular values below $1\text{E-}8$ before finishing. This is a vast improvement over the second stage of PHSVDS, which exhibits irregular convergence with large spikes in the left residual and long stagnations.

3.2. Switching problems. One of the biggest practical advantages of GKD over PHSVDS or any two-stage algorithm is that it avoids the need to switch. For PHSVDS, choosing the right time to switch is crucial so as to give the best possible initial guesses to the second stage in order to avoid excessive use of the second stage on B . However, if an overly optimistic bound is used, it may cause stagnations in the first stage before switching. In general, it can be difficult to converge down to the theoretical limit for the first stage in practice, and determining the minimum constant above the theoretical limit that works for every problem is most likely impossible. Worse, preconditioning can increase this difficulty, as it can cause errors that are difficult to account for within the switching criteria.

Specifically, we found these switching issues to occur when testing PHSVDS on LargeRegFile (another matrix from the SuiteSparse Collection [6]) with block-Jacobi preconditioning and $\delta = 1\text{E-}12$. It is clear from the highlighted portions of Figure 5 that PHSVDS is unable to meet the convergence criteria for the first stage. In fact, while the case shown in Figure 5 is able to reach the criteria eventually, most cases like this stagnate completely. For example, the same problem (LargeRegFile) when solved with an inner solver (JDQMR) is never able to meet the first-stage convergence criteria. Since GKD never requires switching methods, we can avoid these problems entirely and provide more reliable convergence.

3.3. Space and time comparisons. For computations on large matrices, it is important to consider the convergence rate, the space requirements, and the total work that the algorithm requires. Therefore, we provide a short comparison of the latter between our method and PHSVDS before presenting numerical results in section 4.

GKD requires storage for two spaces, V and Q , that are $n \times q$ and $m \times q$, respectively, where q is the maximum basis size. In the PRIMME implementation of PHSVDS, a similar amount of space is required to store the resulting left and right singular vector approximations. However, the first stage of PHSVDS requires a working memory set of two spaces of size $n \times q$ for V and $A^T AV$. Therefore, for square matrices, the working space required for the first stage of PHSVDS is equivalent to GKD. For very tall and skinny matrices ($n \ll m$), the first stage of PHSVDS uses a reduced memory footprint for most of the computation, but only if the user can guarantee that switching to the second stage will not be required. Otherwise, the second stage of PHSVDS will require two spaces of dimension $(m + n) \times q$. This corresponds to double the storage requirement of GKD. For very large problems, this might force the user to reduce the max basis size in order to store the bases in memory.

In terms of execution cost, GKD performs two orthogonalizations per iteration, one for V and one for Q , while the first stage of PHSVDS performs only one orthogonalization for V . Therefore, with low required accuracy where the second stage is

TABLE 1
Basic properties of square matrices.

| Matrix | pde2961 | dw2048 | fidap4 | jagmesh8 | wang3 | lshp3025 |
|-----------------|---------|--------|--------|----------|--------|----------|
| dimension | 2961 | 2048 | 1601 | 1141 | 26064 | 3025 |
| $\text{nnz}(A)$ | 14585 | 10114 | 31837 | 7465 | 77168 | 120833 |
| $\kappa(A)$ | 9.5E+2 | 5.3E+3 | 5.2E+3 | 5.9E+4 | 1.1E+4 | 2.2E+5 |
| $\ A\ $ | 1.0E+1 | 1.0E+0 | 1.6E+0 | 6.8E+0 | 2.7E-1 | 7.0E+0 |
| γ_1 | 8.2E-3 | 2.6E-3 | 1.5E-3 | 1.7E-3 | 7.4E-5 | 1.8E-3 |

TABLE 2
Basic properties of rectangular matrices.

| Matrix | well1850 | lp_ganges | deter4 | plddb | ch | lp_bnl2 |
|-----------------|----------|-----------|--------|--------|--------|---------|
| rows | 1850 | 1309 | 3235 | 3049 | 3700 | 2324 |
| columns | 712 | 1706 | 9133 | 5069 | 8291 | 4486 |
| $\text{nnz}(A)$ | 8755 | 6937 | 19231 | 10839 | 24102 | 14996 |
| $\kappa(A)$ | 1.1E+2 | 2.1E+4 | 3.7E+2 | 1.2E+4 | 2.8E+3 | 7.8E+3 |
| $\ A\ $ | 1.8E+0 | 4.0E+0 | 1.0E+1 | 1.4E+2 | 7.6E+2 | 2.1E+2 |
| γ_1 | 3.0E-3 | 1.1E-1 | 1.1E-1 | 4.2E-3 | 1.6E-3 | 7.1E-3 |

TABLE 3
Basic properties of large-scale matrices.

| Matrix | sls | Rucci1 | LargeRegFile |
|-----------------|-----------|-----------|--------------|
| rows | 1,748,122 | 1,977,885 | 2,111,154 |
| columns | 62,729 | 109,900 | 801,374 |
| $\text{nnz}(A)$ | 6,804,304 | 7,791,168 | 4,944,201 |
| $\kappa(A)$ | 1.3E+3 | 6.7E+3 | 1.1E+4 |
| $\ A\ $ | 1.3E+3 | 7.0E+0 | 3.1E+3 |
| γ_1 | 8E-7 | 5E-5 | 3E-7 |

not involved, PHSVDS is more efficient per step computationally. For robustness, `primme_svds` implements the second stage of PHSVDS using refined extraction which requires two orthogonalizations on vectors of dimension $m + n$ and thus has double the orthogonalization cost of GKD. Additionally, these vectors of size $m + n$ incur more error in dot product computations, and so baseline calculations will not be as accurate. When using low precision calculations (single or half), these errors become even more important to avoid if possible.

4. Numerical results. To verify our algorithm's performance, we utilized the same matrices given in the original PHSVDS publication [37] as well as three matrices with dimension larger than one million from [36]. The matrices shown in Tables 1, 2, and 3 are publicly available through the SuiteSparse Matrix Collection [6] and represent real world applications. These problems are quite difficult for iterative solvers and are used to stress test the capabilities of GKD and PHSVDS. Since these matrices are sparse, we provide their dimensions and the number of nonzero entries of A , $\text{nnz}(A)$, as well as the norm of A , $\|A\|$, the condition number of A , $\kappa(A)$, and the gap ratio for σ_1 , $\gamma_1 = (\sigma_2 - \sigma_1)/(\sigma_n - \sigma_2)$.

The matrices listed in Tables 1 and 2 are listed from least to most difficult (left to right), as generally their condition numbers increase and the gap ratios for their smallest singular values decrease. It should be noted that none of these matrices is particularly poorly conditioned and does not require the second stage in PHSVDS to improve the singular vector estimates by more than a few orders of magnitude. Therefore, the benefits we would expect to gain on very poorly conditioned problems are significantly larger.

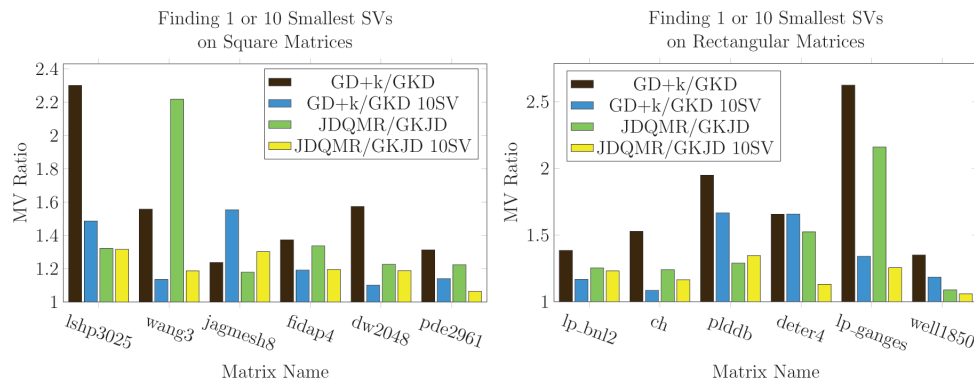


FIG. 6. Unpreconditioned results on 12 problems from the SuiteSparse Matrix Collection with a relative user tolerance of $\delta = 1e - 14$.

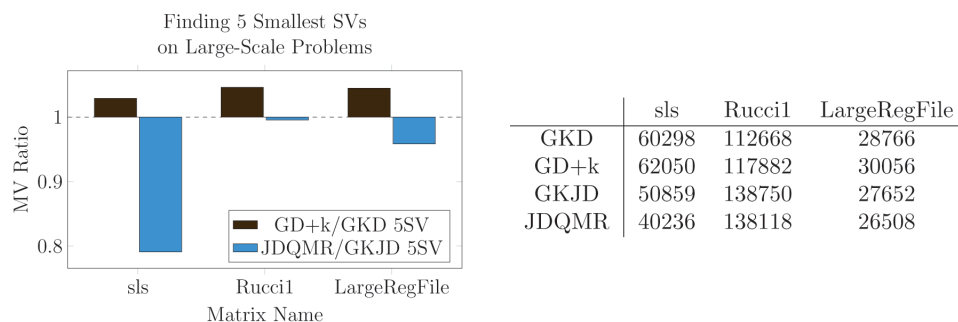


FIG. 7. Large-scale unpreconditioned results. Required matvecs for GKD, GD+k, GKJD, and JDQMR are shown in the table. Note that for sls, GKJD finds three of the singular values with multiplicity 14, while JDQMR finds only two.

We restrict GKD and PRIMME's PHSVDS MATLAB interface, `primme_svd`s, to a maximum basis size of 35 vectors, a minimum restart size of 15 vectors, and a user tolerance of $\delta = 1E-14$ for the smaller matrices and $\delta = 1E-12$ for the larger ones. We also enforce one retained vector from the previous iteration (for +1 restarting) except for the three large cases, where we enforce +2 restarting. Additionally, we choose to soft-lock converged triplets, but due to the interior nature of the augmented method in `primme_svd`s, we are unable to set soft-locking for the second stage while searching for the smallest singular triplets. It should be noted that hard-locking generally improves performance for our method when searching for more than one singular value, but does not provide the same orthogonality guarantees and is subject to the numerical issues mentioned earlier.

4.1. Unpreconditioned results. We compare GD+k (implemented as the default MIN_MATVECS method in `primme_svd`s) against GKD and the JDQMR method (MIN_TIME in `primme_svd`s) against GKJD. As shown in Figure 6, GKD and GKJD require fewer matvecs than their `primme_svd`s counterparts for all matrices. Also, the matrices that show the largest benefits are `lshp3025`, `wang3`, `jagmesh8`, and `lp_ganges`. As expected, these correspond to the matrices that required more significant use of the second stage in `primme_svd`s due to their larger $\kappa(A)$.

For most cases, we see a slight drop off in performance when searching for the 10 smallest singular values, but this is mostly caused by different implementations of soft-locking. Since `primme_svds` uses two stages, the first stage soft-locks each vector at a tolerance above the user specified tolerance. However, since they are soft-locked, the first stage of `primme_svds` can improve the initial guesses to the second stage in some cases since it leaves the estimated singular triplets in the basis while converging to other vectors. To verify this hypothesis, we ran GKD using a pseudo two-stage implementation that mimics the `primme_svds` behavior. This was done by converging all 10 singular values to a higher tolerance first ($\kappa(A)\|A\|\epsilon_{mach}$), before converging to the full user tolerance. In this case, GKD can further improve performance for soft-locking over `primme_svds`.

For rectangular matrices, we also tested whether our method could find a true zero singular value by appending one extra column to the matrix equal to the first column. GKD is able to find the real zero in all cases. `primme_svds` will not return this numerically zero value, as outlined in its documentation, since its second stage has no way to distinguish real zeros from the null space created by the augmented matrix.

For the large-scale matrices, Figure 7 shows fairly similar performance between `primme_svds` and GKD/GKJD. This is expected, as the tolerance is higher ($\text{tol} = 1\text{E-}12$) than the small cases, and therefore `primme_svds` only uses the second stage sparingly. The biggest difference is seen for `sls` and for the inner-outer methods (JDQMR/GKJD), where the high multiplicity (14) at the second smallest singular value causes issues with convergence. Specifically, JDQMR only converges to two of these numerically equal singular values before finding five converged triplets, while GKJD is able to recognize the higher multiplicity and spends extra iterations finding a third. We also note that the number of matvecs for GKD/GKJD is significantly smaller than the numbers for SLEPc's implementation of LBD reported in [36].

In general, iterative methods may have trouble finding multiplicities or may converge out of order, causing the methods to miss directions [23]. This is especially true for Krylov solvers which, in exact arithmetic, are unable to find more than one eigenvector corresponding to a multiplicity. In order to solve this problem, many algorithms, including PHSVDS, can utilize a block solver where the block size approximates the degree of the multiplicity [5, 3, 10]. Additionally, multiple initial guesses can be used to reduce the likelihood of initial vectors being deficient in the invariant space of the multiplicity. Both of these ideas would be simple extensions that could be added to GKD to improve robustness.

4.2. Single precision results. In order to demonstrate the versatility of our method, we ran tests in single precision looking for the largest 10 or 100 singular values of matrices to tolerance $\delta = 1\text{E-}4$. Additionally, our initial basis is built with a simple GKL iteration. Although much less taxing on the solver, these kinds of requirements are common in many SVD applications. We compare our results to IRLBA, which is the default method in `svds` of MATLAB for seeking the largest singular values. Since we are looking for low accuracy, we omit results from PRIMME since it would use only the first stage, which is equivalent to GKD.

Figures 8 and 9 report results on Rucci1. We also ran these tests on `sls` and `LargeRegFile`, but convergence was achieved in too few iterations (requiring only one restart), and so all methods were similar. We vary the maximum basis size to understand how GKD compares when the user has more or less space than IRLBA uses as a default. When searching for 100 singular triplets, we choose basis sizes close

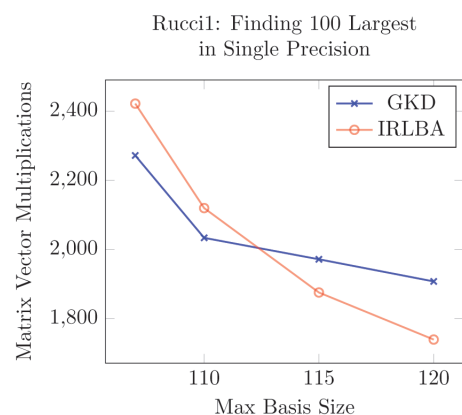


FIG. 8. Similar performance can be achieved with a relatively small basis size even when searching for 100 values.

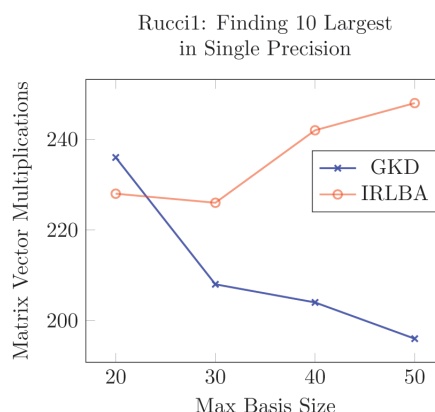


FIG. 9. IRLBA wastes matvecs building a full basis without checking convergence.

to 100 to mimic the situation where space is at a premium and only a small number of extra vectors can be stored. For 10 singular triplets, we show how IRLBA compares to GKD when the basis size is much larger than the number for desired triplets.

Figure 8 shows that both IRLBA and GKD provide fairly similar results for 100 singular values. GKD performs better in the most extreme memory limitation, as it can selectively target the desired values when building its space. However, when there is more room to build a Krylov space, this targeting is no longer required.

Figure 9 shows increased advantages of GKD when fewer singular values are needed. For 10 singular values, the standard version of IRLBA defaults to a maximum basis size of 30. In some cases, the system may have additional space for a larger basis size which can improve convergence. However, since IRLBA generally only checks convergence after a full basis is built, a larger basis size can limit how often IRLBA performs these checks. This allows GKD to outperform IRLBA, even though they obtain nearly identical performance for smaller basis sizes.

To demonstrate the accuracy limit of our method, we run GKD searching for the 10 largest singular values of Rucci1 with a user tolerance of 0. Since this tolerance is not reachable, the original algorithm will stagnate on the largest singular triplet, before it has the chance to target more singular values. To allow all 10 singular triplets to improve, we rotate the target index from 1 to 10 at every iteration and stop when all of the values begin to stagnate at their accuracy limit. Figure 10 shows the convergence of each residual to a stagnation at $\approx 2.5\text{E-}5$. This is impressive given that $\|A\| \approx 7$, $\epsilon_{mach} = 1.2\text{E-}7$, and the matrix dimensions are 2 million by 110K, which affects the dot product accuracy. The associated error bound for a single iteration is $\|A\|\epsilon_{mach}\sqrt{m+n} \approx 1.2\text{E-}3$. Our algorithm goes well below this bound.

4.3. Preconditioned results. In order to test the efficacy of preconditioning GKD, we ran tests on the six smaller square matrices using a preconditioner built from ILU of MATLAB with the ilutp factorization, a drop tolerance of $1\text{E-}3$, and a pivot threshold of 1.0. Our results in Figure 11 show the significant benefit of an effective preconditioner, as all of the small problems required less than 150 matvecs when searching for one singular value with GKD. However, these preconditioners sometimes caused significant issues for `primme_svds`, as it was unable to converge for `lshp3025`

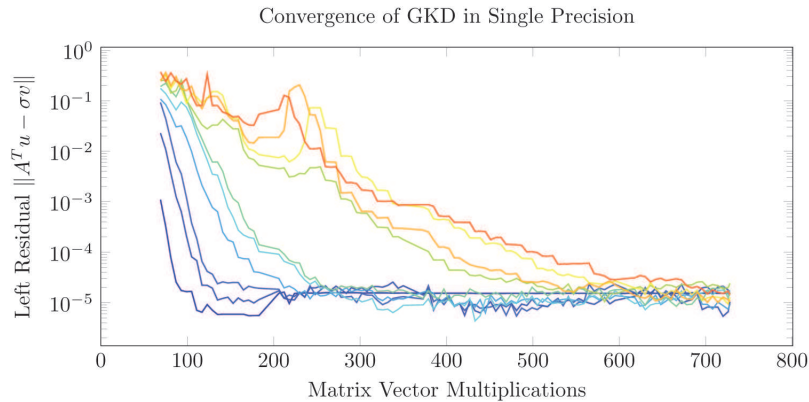


FIG. 10. Maximum accuracy achievable with GKD in single precision for the 10 largest SVs on Rucci1. We set $\delta = 0$ and rotate the target at each iteration to allow all triplets to converge.

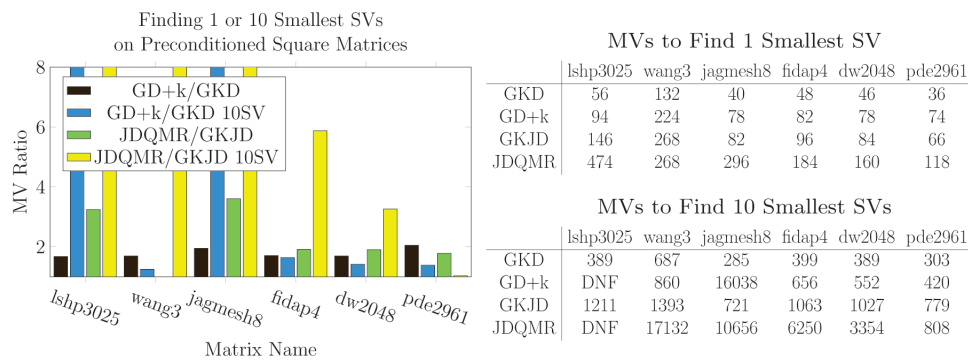


FIG. 11. Preconditioned results with an ILU preconditioner for finding the smallest and 10 smallest singular triplets.

when searching for the 10 smallest singular values, and exhibited significant difficulty converging to 10 singular values for wang3, jagmesh8, and fidap4. Specifically, when searching for 10 singular values, wang3 requires 12x more matvecs for JDQMR, and jagmesh8 requires 56x and 14x more matvecs for GD+k and JDQMR, respectively. These issues are caused by `primme_svds`' switching issues mentioned earlier.

For the three large rectangular matrices, we use a block-Jacobi preconditioner, inverting exactly diagonal blocks of $A^T A$ each of size 600. This is relatively inexpensive to compute, and it is also parallelizable. The results in Figure 12 show a significant decrease in matvecs, as all three problems required less than 15% of the matvecs needed for the unpreconditioned cases. For Rucci1, the convergence differences between our methods and `primme_svds` are negligible, but for sls and LargeRegFile, GKD and GKJD provide significant improvements in speed and robustness. Again, as seen earlier in Figure 5, `primme_svds`' switching criteria are too stringent for preconditioned cases, which causes slowdowns for GD+k on LargeRegFile. Worse, `primme_svds`' JDQMR suffers stagnations that cause failures to converge when preconditioned on sls and LargeRegFile.

The 80% improvement on sls over GD+k comes from `primme_svds` being unable

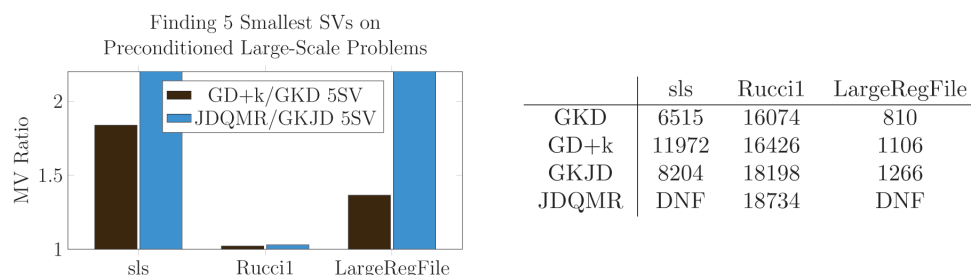


FIG. 12. Large-scale results with block-Jacobi preconditioner (block size= 600 on $A^T A$) for the five smallest singular triplets. Required matvecs for GKD, GD+k, GKJD, and JDQMR are shown in the table.

to separate the directions corresponding to the large degree multiplicity. During additional testing, we found the number of matvecs required to find the five smallest singular values with `primme_svds` is only marginally less than the number required to find 10. Since `primme_svds` is unable to appropriately separate the directions corresponding to the multiplicity, it converges to all 10 values concurrently. However, GKD is able to distinguish these directions and converge smoothly for each one individually, providing a substantial improvement. Testing GKD to converge to 10 values as well, we still found an improvement over `primme_svds`; however, the gap between the two methods was significantly reduced.

5. Conclusions. We have presented GKD, a new method for finding the smallest singular triplets of large sparse matrices to full accuracy. Our method works iteratively, under limited memory, with preconditioners, while including features such as soft-locking with orthogonality guarantees, +k restarting, and the ability to find real zero singular values in both square and rectangular matrices. Additionally, GKJD adds a Jacobi–Davidson inner solver for the $A^T A$ correction equation into GKD, which can lower execution time when the matrix-vector multiplication operation is inexpensive and can reduce the errors caused by restarting. Both of these methods have shown to be more reliable and efficient than PHSVDS, and thus over other SVD methods, for nearly all cases.

REFERENCES

- [1] ADVANCED MICRO DEVICES, INC., *AMD OpenCL Optimisation Guide*, http://developer.amd.com/wordpress/media/2013/12/AMD_OpenCL_Programming_Optimization_Guide2.pdf (accessed 2018-02-14).
- [2] O. ALTER, P. O. BROWN, AND D. BOTSTEIN, *Singular value decomposition for genome-wide expression data processing and modeling*, Proc. Natl. Acad. Sci. USA, 97 (2000), pp. 10101–10106.
- [3] J. BAGLAMA, D. CALVETTI, AND L. REICHEL, *IRBL: An implicitly restarted block-Lanczos method for large-scale Hermitian eigenproblems*, SIAM J. Sci. Comput., 24 (2003), pp. 1650–1677, <https://doi.org/10.1137/S1064827501397949>.
- [4] J. BAGLAMA AND L. REICHEL, *Augmented implicitly restarted Lanczos bidiagonalization methods*, SIAM J. Sci. Comput., 27 (2005), pp. 19–42, <https://doi.org/10.1137/04060593X>.
- [5] J. BAGLAMA AND L. REICHEL, *Restarted block Lanczos bidiagonalization methods*, Numer. Algorithms, 43 (2006), pp. 251–272.
- [6] T. A. DAVIS AND Y. HU, *The University of Florida sparse matrix collection*, ACM Trans. Math. Software, 38 (2011), 1.
- [7] A. S. GAMBHIR, A. STATHOPOULOS, AND K. ORGINOS, *Deflation as a method of variance re-*

- duction for estimating the trace of a matrix inverse, SIAM J. Sci. Comput., 39 (2017), pp. A532–A558, <https://doi.org/10.1137/16M1066361>.
- [8] G. GOLUB AND W. KAHAN, *Calculating the singular values and pseudo-inverse of a matrix*, J. Soc. Indust. Appl. Math. Ser. B Numer. Anal., 2 (1965), pp. 205–224.
 - [9] G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 3rd ed., Johns Hopkins University Press, Baltimore, MD, 1996.
 - [10] R. G. GRIMES, J. G. LEWIS, AND H. D. SIMON, *A shifted block Lanczos algorithm for solving sparse symmetric generalized eigenproblems*, SIAM J. Matrix Anal. Appl., 15 (1994), pp. 228–272, <https://doi.org/10.1137/S0895479888151111>.
 - [11] S. GUPTA, A. AGRAWAL, K. GOPALAKRISHNAN, AND P. NARAYANAN, *Deep learning with limited numerical precision*, in Proceedings of the International Conference on Machine Learning, 2015, pp. 1737–1746.
 - [12] M. E. HOCHSTENBACH, *A Jacobi–Davidson type SVD method*, SIAM J. Sci. Comput., 23 (2001), pp. 606–628, <https://doi.org/10.1137/S1064827500372973>.
 - [13] W. HOFFMANN, *Iterative algorithms for Gram–Schmidt orthogonalization*, Computing, 41 (1989), pp. 335–348.
 - [14] J. HUANG AND Z. JIA, *On inner iterations of Jacobi–Davidson type methods for large SVD computations*, SIAM J. Sci. Comput., 41 (2019), pp. A1574–A1603, <https://doi.org/10.1137/18M1192019>.
 - [15] Z. JIA AND C. LI, *Inner iterations in the shift-invert residual Arnoldi method and the Jacobi–Davidson method*, Sci. China Math., 57 (2014), p. 1733–1752, <https://doi.org/10.1007/s11425-014-4791-5>.
 - [16] Z. JIA AND C. LI, *Harmonic and refined harmonic shift-invert residual Arnoldi and Jacobi–Davidson methods for interior eigenvalue problems*, J. Comput. Appl. Math., 282 (2015), pp. 83–97, <https://doi.org/10.1016/j.cam.2014.12.043>.
 - [17] Z. JIA AND D. NIU, *An implicitly restarted refined bidiagonalization Lanczos method for computing a partial singular value decomposition*, SIAM J. Matrix Anal. Appl., 25 (2003), pp. 246–265, <https://doi.org/10.1137/S0895479802404192>.
 - [18] I. JOLLIFFE, *Principal Component Analysis*, Springer-Verlag, New York, 2002.
 - [19] A. V. KNYAZEV, *Toward the optimal preconditioned eigensolver: Locally optimal block preconditioned conjugate gradient method*, SIAM J. Sci. Comput., 23 (2001), pp. 517–541, <https://doi.org/10.1137/S1064827500366124>.
 - [20] R. M. LARSEN, *Lanczos bidiagonalization with partial reorthogonalization*, DAIMI Report Series, 27 (1998), 537.
 - [21] Q. LIANG AND Q. YE, *Computing singular values of large matrices with an inverse-free preconditioned Krylov subspace method*, Electron. Trans. Numer. Anal., 42 (2014), pp. 197–221.
 - [22] S. MARKIDIS, S. W. D. CHIEN, E. LAURE, I. B. PENG, AND J. S. VETTER, *NVIDIA Tensor Core Programmability, Performance & Precision*, preprint, <https://arxiv.org/abs/1803.04014>, 2018.
 - [23] J. R. MCCOMBS AND A. STATHOPOULOS, *Iterative validation of eigensolvers: A scheme for improving the reliability of Hermitian eigenvalue solvers*, SIAM J. Sci. Comput., 28 (2006), pp. 2337–2358, <https://doi.org/10.1137/050627617>.
 - [24] K. MEERBERGEN, R. MORGAN, AND A. KNYAZEV, *Inexact methods*, in Templates for the Solution of Algebraic Eigenvalue Problems: A Practical Guide, Z. Bai, J. Demmel, J. Dongarra, A. Ruhe, and H. van der Vorst, eds., SIAM, Philadelphia, 2000, pp. 337–368, <https://doi.org/10.1137/1.9780898719581.ch11>.
 - [25] Y. NOTAY, *Combination of Jacobi–Davidson and conjugate gradients for the partial symmetric eigenproblem*, Numer. Linear Algebra Appl., 9 (2002), pp. 21–44.
 - [26] S. OSIŃSKI, J. STEFANOWSKI, AND D. WEISS, *Lingo: Search results clustering algorithm based on singular value decomposition*, in Intelligent Information Processing and Web Mining, Springer, Berlin, Heidelberg, 2004, pp. 359–368.
 - [27] B. N. PARLETT, *The Symmetric Eigenvalue Problem*, Prentice–Hall, Englewood Cliffs, NJ, 1980.
 - [28] H. PRASANTHA, H. SHASHIDHARA, AND K. B. MURTHY, *Image compression using SVD*, in Proceedings of the IEEE International Conference on Computational Intelligence and Multimedia Applications, Vol. 3, 2007, pp. 143–145.
 - [29] H. D. SIMON AND H. ZHA, *Low-rank matrix approximation using the Lanczos bidiagonalization process with applications*, SIAM J. Sci. Comput., 21 (2000), pp. 2257–2274, <https://doi.org/10.1137/S1064827597327309>.
 - [30] A. STATHOPOULOS, *Locking Issues for Finding a Large Number of Eigenvectors of Hermitian Matrices*, Tech. Report WM-CS-2005-09, Computer Science, The College of William & Mary, Williamsburg, VA, 2005.

- [31] A. STATHOPOULOS, *Nearly optimal preconditioned methods for Hermitian eigenproblems under limited memory. Part I: Seeking one eigenvalue*, SIAM J. Sci. Comput., 29 (2007), pp. 481–514, <https://doi.org/10.1137/050631574>.
- [32] A. STATHOPOULOS AND Y. SAAD, *Restarting techniques for (Jacobi-)Davidson symmetric eigenvalue methods*, Electron. Trans. Numer. Anal., 7 (1998), pp. 163–181.
- [33] A. STATHOPOULOS AND K. WU, *A block orthogonalization procedure with constant synchronization requirements*, SIAM J. Sci. Comput., 23 (2002), pp. 2165–2182, <https://doi.org/10.1137/S1064827500370883>.
- [34] E. VECHARYNSKI, *Preconditioned Iterative Methods for Linear Systems, Eigenvalue and Singular Value Problems*, Ph.D. thesis, University of Colorado at Denver, Denver, CO, 2011.
- [35] K. WU AND H. SIMON, *Thick-restart Lanczos method for large symmetric eigenvalue problems*, SIAM J. Matrix Anal. Appl., 22 (2000), pp. 602–616, <https://doi.org/10.1137/S0895479898334605>.
- [36] L. WU, E. ROMERO, AND A. STATHOPOULOS, *PRIMME-SVDS: A High-Performance Preconditioned SVD Solver for Accurate Large-Scale Computations*, preprint, <https://arxiv.org/abs/1607.01404>, 2016.
- [37] L. WU AND A. STATHOPOULOS, *A preconditioned hybrid SVD method for accurately computing singular triplets of large matrices*, SIAM J. Sci. Comput., 37 (2015), pp. S365–S388, <https://doi.org/10.1137/140979381>.
- [38] P. ZHANG AND Y. GAO, *Matrix multiplication on high-density multi-GPU architectures: Theoretical and experimental investigations*, in International Conference on High Performance Computing, Springer, Cham, 2015, pp. 17–30.