# ON TRIDIAGONALIZING AND DIAGONALIZING SYMMETRIC MATRICES WITH REPEATED EIGENVALUES*

CHRISTIAN H. BISCHOF[†] AND XIAOBAI SUN[‡]

**Abstract.** We describe a divide-and-conquer tridiagonalization approach for matrices with repeated eigenvalues. Our algorithm hinges on the fact that, under easily constructively verifiable conditions, a symmetric matrix with band width $b$ and $k$ distinct eigenvalues must be block diagonal with diagonal blocks of size at most $bk$. A slight modification of the usual orthogonal band-reduction algorithm allows us to reveal this structure, which then leads to potential parallelism in the form of independent diagonal blocks. Compared to the usual Householder reduction algorithm, the new approach exhibits improved data locality, significantly more scope for parallelism, and the potential to reduce arithmetic complexity by close to 50% for matrices that have only two numerically distinct eigenvalues. The actual improvement depends to a large extent on the number of distinct eigenvalues and a good estimate thereof. However, at worst the algorithms behave like a successive band-reduction approach to tridiagonalization. Moreover, we provide a numerically reliable and effective algorithm for computing the eigenvalue decomposition of a symmetric matrix with two numerically distinct eigenvalues. Such matrices arise, for example, in invariant subspace decomposition approaches to the symmetric eigenvalue problem.

**Key words.** tridiagonalization, eigenvalue decomposition, repeated eigenvalues

**AMS subject classifications.** 15A23, 15A18, 65F15, 65F25

**1. Introduction.** Let $A$ be an $n \times n$ symmetric matrix. Our goal is to compute an orthogonal–tridiagonal decomposition of $A$, $AQ = QT$, where $Q$ is orthogonal and $T$ is tridiagonal. Reduction to tridiagonal form is a standard preprocessing step in dense eigensolvers based on QR iteration, bisection, or Cuppen's method [16]. The conventional tridiagonalization procedure [16, p. 419] reduces $A$ one column at a time through Householder transformation at a cost of $O(4n^3/3)$ flops for the reduction of $A$, and an additional $O(4n^3/3)$ flops if the orthogonal matrix is accumulated at the same time. This algorithm mainly employs matrix–vector multiplications and symmetric rank-one updates, which require more memory references than matrix–matrix operations [9, 8, 14].

The block tridiagonalization algorithm in [5, 15] combines sets of $p$ successive symmetric rank-one updates into one symmetric rank-$p$ update at the cost of $O(2pn^2)$ extra flops. As a result, this algorithm exhibits improved data locality and hence is likely to be preferable on cache-based architectures. This block algorithm has been incorporated into the LAPACK library of portable linear algebra codes for high-performance architectures [1, 2]. Parallel versions for distributed memory machines of the standard algorithm and the block algorithm are described in [12] and [13], respectively. A different approach to tridiagonalization is the so-called successive

† Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL 60439-4801 (bischof@mcs.anl.gov).
‡ Department of Computer Science, Duke University, Durham, NC 27708-0129 (xiaobai@cs.duke.edu). The work of this author was partially performed as a postdoctoral associate at Argonne National Laboratory.

band-reduction (SBR) method, which completes the tridiagonal reduction through a sequence of band reductions [10, 7]. This approach leads to algorithms that exhibit an even greater degree of memory locality, among other desirable features.

In this paper we show that if the number $k$ (say) of distinct eigenvalues of a symmetric matrix $A$ is small, then there is considerable scope for further savings in tridiagonalization algorithms. As will be demonstrated, $A$ can be cheaply reduced to a block diagonal banded form through a slightly modified SBR approach. The final tridiagonal form is then achieved by applying the algorithm recursively on the subblocks on the diagonal. Compared to the conventional approach, this approach has the following advantages.

*Improved data locality.* The tridiagonalization process can employ mainly matrix–matrix operations both in the reduction of $A$ and in the update of the transformation matrix $Q$ (see also [10, 7]).

*Enhanced scope for parallelism.* In the traditional algorithm, the scope for the exploitation of parallelism in the reduction of $A$ is limited to the application of the rank-one update (for the unblocked algorithm) or the rank-$p$ update (for the blocked algorithm), and the scope for parallelism decreases as subproblems become smaller. In contrast, our algorithm generates independent subproblems during the reduction of $A$, which can be worked on independently, and the number of independent subproblems increases as the iteration proceeds. Thus, there is a shift from data parallelism (updates of large matrices) to functional parallelism (several independent subproblems), but at any stage, there is plenty of parallelism to exploit.

*Reduced complexity.* Depending on the number of distinct eigenvalues, we may almost halve the number of floating-point operations. In addition, the need for data movement is reduced.

One particular situation where repeated eigenvalues arise is in the context of invariant subspace methods for eigenvalue problems [3, 19, 6, 4], where a matrix with only two distinct predetermined eigenvalues is generated either by repeated application of incomplete beta functions [19] or the matrix sign function [4]. In exact arithmetic, our tridiagonalization procedure would result in a block diagonal matrix with diagonal blocks of order no larger than 2. Hence the eigenvalue decomposition could be computed easily by independently diagonalizing the $2 \times 2$ blocks on the diagonal. In the presence of roundoff errors, the computed tridiagonal matrix may not have this desirable structure. However, we can prove that such a tridiagonal matrix can be diagonalized as reliably as with any other method by two "clean up sweeps," where each sweep solves at most $n/2$ independent $2 \times 2$ eigenvalue problems.

The paper is organized as follows. We show in §2 that, under certain easily constructively verifiable conditions, a banded symmetric matrix with band width $b$ and $k$ distinct eigenvalues is block diagonal with diagonal blocks of order at most $bk$. In §3, we present a reduction algorithm to achieve the desired banded block diagonal structure through a slight modification of the conventional band-reduction procedure. This approach is then employed to develop a divide-and-conquer tridiagonalization algorithm. An inexpensive algorithm for decoupling invariant subspaces of matrices with eigenvalue clusters at 0 and 1 is given and verified in §4. Numerical experiments with a Matlab implementation are reported in §5. Lastly, we summarize our results.

**2. The structure of band matrices with repeated eigenvalues.** A tridiagonal matrix whose off diagonal entries are all nonzero is called *unreduced.* It is well known [18, p. 66] that an unreduced tridiagonal matrix does not have multiple eigenvalues. Consequently, if an $n \times n$ tridiagonal matrix has only $k \ll n$ distinct

eigenvalues, it must be block diagonal, and the largest block cannot be larger than $k \times k$. The generalization of this fact to banded matrices underpins the algorithm we propose, yet it is not as straightforward as it might seem.

Assuming that $A$ is an $n \times n$ symmetric matrix, we define the *ith row-band-width* of $A$, denoted by band_row($i$), as

$$(1) \qquad \text{band\_row}(i) \overset{\text{def}}{=} \max_j \{i - j \mid j = i \text{ or } j < i \text{ and } a_{ij} \neq 0\}, \quad 1 \leq i \leq n.$$

That is, band_row($i$) is the distance of the first nonzero element in row $i$ from the $i$th diagonal element. Further, we say that $A$ is *nonincreasing in row-band-width from $b$* if

$$(2) \qquad a(b, 1) \neq 0 \text{ and band\_row}(i) \leq \text{band\_row}(i - 1), \quad b + 1 < i \leq n.$$

In particular, a banded matrix that is all zero below the $b$th subdiagonal and all nonzero on the $b$th subdiagonal is nonincreasing in row-band-width from $b$.

With these definitions, we can now prove the following theorem.

THEOREM 2.1. *Let $T$ be a symmetric matrix with $k$ distinct eigenvalues. If $T$ is block diagonal, with each diagonal block nonincreasing in band width from at most $b$, then the size of the largest block cannot exceed $kb$.*

*Proof.* Assume that $T$ has a diagonal block $D$ of size $p > kb$. By assumption, $D$ is nonincreasing in band width from $b$; that is, $D$ has $p - b$ rows with their first nonzero elements in different columns to the left of the diagonal. Thus, for any $\lambda$, rank($D - \lambda I$) is not less than $p - b$.

On the other hand, since $p > kb$ and $D$ has at most $k$ distinct eigenvalues, $D$ has an eigenvalue $\mu$ with multiplicity greater than $b$. Hence, rank($D - \mu I$) is less than $p - b$. The contradiction verifies the result of the theorem.  □

The following example shows the necessity of the "nonincreasing band-width" restriction in Theorem 2.1. Let

$$Q^T = \begin{pmatrix} \xi & \eta & \mu & -\nu & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \beta & 0 & \gamma & \delta \\ 0 & 0 & \nu & \mu & 0 & \alpha & 0 & 0 \end{pmatrix},$$

where $\nu^2 + \mu^2 + \alpha^2 = 1$, $\xi^2 + \eta^2 = \alpha^2$, and $\beta^2 + \gamma^2 + \delta^2 = 1$. Then $Q$ has orthonormal columns and $A = QQ^T$ is symmetric with only 0 and 1 as eigenvalues. In fact,

$$(3) \qquad A = \begin{bmatrix} \times & \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & \times & 0 & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & \times & 0 & 0 \\ \times & \times & 0 & \times & 0 & \times & 0 & 0 \\ 0 & 0 & 0 & 0 & \times & 0 & \times & \times \\ 0 & 0 & \times & \times & 0 & \times & 0 & 0 \\ 0 & 0 & 0 & 0 & \times & 0 & \times & \times \\ 0 & 0 & 0 & 0 & \times & 0 & \times & \times \end{bmatrix}.$$

We see that $A$ is banded with semi-band-width $b = 3$, but it is *not* block diagonal with blocks of size at most $2b \times 2b = 6 \times 6$ since the "nonincreasing band-width condition" is violated by $a(5, 2) = a(7, 4) = 0$.

**3. A divide-and-conquer tridiagonalization approach.** The example in the previous section showed that the standard Householder band-reduction algorithm will not necessarily reveal the block diagonal structure. For example, if we had applied the standard algorithm for reduction to band width 3 to the matrix of example (3), the matrix would have remained unchanged. Fortunately, a minor modification of the standard algorithm enforces nonincreasing row-band-width, and hence the prerequisites of Theorem 2.1.

Let us consider the conventional reduction approach, where the matrix is reduced one column at a time to semi-band-width $b$. In each reduction, the pivot row is always $b$ rows below the diagonal, no matter whether the reduction of the previous column was skipped (i.e., the transformation was an identity) or not. For example, reducing the matrix $A$ in (3) to semi-band-width 3, row number 4 is the pivot row for the reduction of the second column and, since $a(4:8,2) = 0$, this reduction is skipped. We then proceed to column 3, using row 5 as pivot row, and the row-band-width increases. If instead we employ a Householder transformation acting on $a(4:8,3)$ to eliminate $a(5:8,3)$, keeping row 4 as pivot row, we obtain

$$
\tilde{A} = \begin{bmatrix}
\times & \times & \times & \times & 0 & 0 & 0 & 0 \\
\times & \times & \times & \times & 0 & 0 & 0 & 0 \\
\times & \times & \times & 0 & \times & 0 & 0 & 0 \\
\times & \times & 0 & \times & \times & 0 & 0 & 0 \\
0 & 0 & \times & \times & \times & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & \times & \times & \times
\end{bmatrix}.
$$

Now $\tilde{A}$ is decoupled into two diagonal blocks of size at most $6 \times 6$.

This example shows that nonincreasing band width can easily be obtained if we do not increase the pivot row when the previous reduction is skipped. For computational purposes, we define the row-band-width with respect to a threshold $\tau$:

$$
(4) \quad \text{band\_row}(i, \tau) \overset{\text{def}}{=} \max_{j}\{i - j \mid j = i \text{ or } j < i \text{ and } \|a(i:n, j)\|_2 > \tau\}, \quad 1 \le i \le n.
$$

That is, given a tolerance threshold $\tau$, a column $a(i:n)$ is considered numerically zero if its 2-norm is at most $\tau$. The Matlab function `bred` in Figure 1 shows the conventional band-reduction algorithm augmented with

(1) a threshold criterion for the generation of a Householder vector, and

(2) a modified pivot row selection strategy, which does not change the pivot row if a transformation is skipped.

The subroutines `gen_hh`, `pre_hh`, `post_hh`, and `sym_hh` generate a Householder vector and apply it from the left, right, and symmetrically, respectively. Note that for simplicity the algorithm presented here does not exploit the symmetry of $A$. However, if we wish to do so, we can have `sym_hh` work only with a triangular part of $A$ and omit the `post_hh` (`pre_hh`) call when working only with the lower (upper) triangle. We also note that all the algorithms presented in this paper are available via anonymous ftp from the `pub/prism` directory at `ftp.super.org`.

If no transformations are skipped, the procedure is identical to the conventional band-reduction procedure; otherwise, it may terminate earlier when the reduction reaches the last column of the first diagonal block, and the problem is decoupled. Since we drop pivot columns whose norm is $O(\tau)$, the decomposition will be accurate up to a residual of order $\tau$.

```
   function [A, block1, Q] = bred( A, b, tau, Q );

   %    Given a symmetric matrix A, a bandwidth b, and a threshold tau, bred
   %    computes an orthogonal-banded matrix decomposition,
   %            A_input * W = W * A_output + O(tau)
 5 %    where O(tau) denotes a matrix with a two-norm of order tau, and
   %    W is an orthogonal matrix.
   %    The output matrix A_output will be a 2x2 block diagonal matrix,
   %    where the first diagonal block A_output(1:block1,1:block1)
   %    is banded with bandwidth nonincreasing from b, and the second block
10 %    may be empty.

   [ m, n ] = size(A); if (m~=n) error('nonsquare A'); end
   piv_row  = min(b+1,n);        % current pivot row
   if (piv_row == n) block1 = n; return; end;
   for j = 1:n-b
15         % matrix is decoupled, stop
       if (piv_row == j), break, end
           % row and column sets involved in current transformation
       rows = (piv_row : n); cols = (j+1:piv_row-1);
           % generate HH matrix to annihilate A(piv_row+1:n,j)
20     [ v, beta, gamma ] = gen_hh( A( rows, j), tau );
           % update jth row and column of A
       A( rows, j) = zeros(size(rows')); A(piv_row, j) = gamma;
       A( j, rows) = zeros(size(rows));  A(j, piv_row) = gamma;
           % if the reduction is not "skipped", perform symmetric
25         % update of A, update Q if required, and shift the pivot row
       if ( beta ~= 0)
         if( cols~= [] )
             A(rows, cols) =  pre_hh( beta, v, A(rows, cols) );
             A(cols, rows) = post_hh( beta, v, A(cols, rows) );
30         end
         A( rows, rows ) = symm_hh( beta, v, A(rows, rows) );
         if( Q ~= [] ), Q(:, rows) = post_hh( beta, v, Q(:, rows ) ); end
       end % beta
           % increase pivot row if A(piv_row,j) is not negligible
35     if (abs(A(j,piv_row)) > tau), piv_row = piv_row + 1; end
   end % j-loop
   if (j == n - b)
     if (piv_row == j+1), block1 = piv_row - 1; else, block1 = n; end
   else
40   block1 = piv_row-1;
   end
   return; end
```

FIG. 1. *Nonincreasing row-band-width preserving band-reduction algorithm.*

For simplicity we omitted an optimization in Figure 1—if the reduction of the first column of $A$ results in a band width $\tilde{b}$, say, where $\tilde{b} < b$, due to the small size of entries $a(\tilde{b} + 1 : n, 1)$, we can directly pursue a reduction of the trailing block to nonincreasing band width $\tilde{b}$ in the same fashion as shown above.

If the parameter $b$ is chosen such that $kb < n$, where $k$ is the number of distinct eigenvalues of $A$, Theorem 2.1 predicts a decoupling of the problem with the leading block being of size no larger than $kb$. In particular, if $b$ is chosen such that $kb = n/2$, we can expect bred to generate two decoupled subproblems of about the same size. We can then recursively divide the problem until the transformed matrix becomes tridiagonal (i.e., $b = 1$). Figure 2 is a serial implementation of tridiagonalization based on this approach. Note that the various subproblems can be dealt with independently and simultaneously. The subroutine blk_diag, which is called in tri_sbr, is shown in Figure 3 and reduces a matrix to block diagonal form with a given band width.

For example, if we reduce a $12 \times 12$ matrix $A$ with only two eigenvalues to band width 3, then no diagonal block can be larger than $6 \times 6$. So, if $a(4, 1)$, $a(5, 2)$, and $a(6, 3)$ are all nonzero after the reductions in the first three columns have been completed, then the next three columns must already be reduced, and the (partially reduced) matrix $A$ is of the form

$$\begin{bmatrix}
\times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
\times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
\times & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\
\times & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & \times & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & \times & \times & \times & \times & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times \\
0 & 0 & 0 & 0 & 0 & 0 & \times & \times & \times & \times & \times & \times
\end{bmatrix}.$$

As a result, we do not need to perform the reductions that would otherwise have occurred in columns 4 through 6. Compared to the conventional approach, the complexity of the algorithm for the case $k = 2$ is $O(0.55\, n^3)$ for the reduction of $A$ and $O(1.25\, n^3)$ for the update of $Q$, as compared to $O(4n^3/3)$ for both these operations in the usual approach. The savings for $Q$ are minor since updates at later stages still involve vectors of length $n$, whereas only diagonal subblocks are affected in $A$. In addition, we can work in parallel on independent problems. If the estimate $k$ of the number of distinct eigenvalues is inaccurate, the algorithm becomes either the standard eigenvalue algorithm (for $k > n/2$) or the SBR tridiagonalization procedure suggested in [10], but in either case, it will return numerically accurate results.

**4. Invariant subspace splitting.** The computational cost and the degree of parallelism in the algorithm depend on $k$, the number of distinct eigenvalues. One particularly intriguing case is matrices that have only two eigenvalues. It is intriguing because they arise in eigensolvers based on variant subspace decompositions [3, 19, 4]. We may assume without loss of generality that the eigenvalues are at 1 and 0 (any other two eigenvalues can be mapped to 0 and 1 by shifting and scaling). The following corollary is a special case of Theorem 2.1.

```
    function [A, Q] = tri_sbr( A, k, tau, Q )
    %
    %  produces an orthogonal-tridiagonal decomposition of
    %  a symmetric matrix A such that
 5  %            A_old*Q = Q*A_new + O(tau)
    %  where A_new is tridiagonal and Q is orthogonal.
    %
    %  The number k is a guess at the number of numerically distinct
    %  eigenvalues of A.
10  %
    %  Matrices are successively reduced to smaller bandwidth in an
    %  attempt to exploit the divide-and-conquer nature becoming
    %  apparent in the successive bandreduction algorithm when the number
    %  k chosen is a good guess at the actual number of numerically distinct
15  %  eigenvalues.

    [m, n] = size(A); if( m ~= n ) error('non-square A'); end

    b = max( floor(n/(2*k)), 1 );

    [A, block1, Q] = bred( A, b, tau, Q );

    if (block1 == n) % If problem didn't decouple, just reduce to
20                   % tridiagonal form
       [A,blkvec,Q] = blk_diag(A,1,tau,Q); return;
    else
     if( b > 1 )     % first subproblem is not tridiagonal yet
       sub = 1:block1; V = eye(block1);
25     [ A(sub,sub), V ] = tri_sbr( A( sub, sub), k, tau, V );
       Q(:,sub) = Q(:,sub) * V;
     end;
     if( n-block1 > 2 )  % second subproblem is nontrivial
       sub = (block1+1):n; V = eye(n-block1);
30     [ A(sub, sub), V ] = tri_sbr( A(sub, sub), k, tau, V );
       Q(:,sub) = Q(:,sub) * V;
     end
    end

    return;
35  end
```

FIG. 2. *Divide-and-conquer tridiagonalization.*

COROLLARY 4.1. *Let $A$ be a matrix with two distinct eigenvalues, and let $A = Q^T T Q$ be a tridiagonalization of $A$. Then $T$ is block diagonal with diagonal blocks of size at most $2 \times 2$.*

Corollary 4.1 implies that one can determine the range space, $\mathcal{R}(A)$, and the null space, $\mathcal{N}(A)$, in essence via a tridiagonalizing of $A$. Let $AQ = QT$ be the orthogonal-tridiagonal decomposition of $A$. For a $1 \times 1$ diagonal block $T(j,j)$,

$$Q(:,j) \in \mathcal{R}(A) \text{ if } T(j,j) = 1, \quad \text{and} \quad Q(:,j) \in \mathcal{N}(A) \text{ if } T(j,j) = 0.$$

```
    function [ A, blkvec, Q ] = blk_diag( A, b, tau, Q )
    %
    % Given a symmetric matrix A, a desired bandwidth b, and a threshold tau,
    %            [ A, bvec, Q ] = blk_diag( A, b, tau, Q )
  5 % produces an orthogonal-block-diagonal decomposition
    %              A_input * W = W * A_output + O(tau)
    % where O(tau) denotes a matrix whose norm is of order tay, and
    % W is an orthogonal matrix.
    %
 10 % A_output will be a block diagonal matrix with each block banded with
    % nonincreasing bandwidth b. The i-th diagonal block starts
    % at (blkvec(i), blkvec(i)).
    %
    % If Q is not the empty matrix on input, Q is postmultiplied by W,
 15 %    i.e.,, Q_output = Q_input * W.

    [m, n] = size(A); if( m ~= n ) error('non-square A'); end

    j = 1; blkvec = [];
    while( j < n )
      blkvec = [ blkvec j ]; rows = j:n; cols = j:n;
 20   [A(rows, cols), dj, Q(:,cols) ] = bred( A(rows,cols), b, tau, Q(:,cols);
      j = j + dj;
    end

    return; end
```

FIG. 3. *Reduction to block diagonal form.*

Since the eigenvalues of $A$ and $T$ are the same, a $2\times2$ diagonal block $T(j{:}j{+}1, j{:}j{+}1)$ must have eigenvalues 0 and 1. Because the trace is the sum of the eigenvalues and the off diagonal entry is nonzero, we have

$$T(j{:}j{+}1, j{:}j{+}1) = \begin{pmatrix} 1-\gamma & \mu \\ \mu & \gamma \end{pmatrix},$$

where $\mu \neq 0$ and $0 < \gamma < 1$. Since

$$\begin{pmatrix} 1-\gamma & \mu \\ \mu & \gamma \end{pmatrix} = \frac{1}{\gamma} \begin{pmatrix} \mu & \gamma \\ \gamma & -\mu \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mu & \gamma \\ \gamma & -\mu \end{pmatrix}^T,$$

we conclude that

$$Q\,(:,j{:}j{+}1) \begin{pmatrix} \mu \\ \gamma \end{pmatrix} \in \mathcal{R}(A), \quad \text{and} \quad Q\,(:,j{:}j{+}1) \begin{pmatrix} \gamma \\ -\mu \end{pmatrix} \in \mathcal{N}(A).$$

One can see that the separation of the range and null subspaces of $A$, and in fact its eigenvalue decomposition, can be effected by diagonalizing (potentially in parallel) the $2\times2$ subproblems still occurring in the block tridiagonal decomposition.

In the presence of rounding errors, a computed tridiagonal matrix may not, however, exhibit the block structure we could expect from Corollary 4.1 due to perturbations in the eigenvalues. That is, $\lambda(T) \subset \{[-\nu,\nu] \cup [1-\nu, 1+\nu]\}$, and a repeated eigenvalue numerically manifests itself as an eigenvalue cluster.

*Example* 4.2. The matrix

$$
T = \begin{pmatrix}
1 & e_1 \\
e_1 & 0 & e_2 \\
& e_2 & 1 & e_3 \\
& & e_3 & 0 & e_4 \\
& & & e_4 & \ddots \\
& & & & & 1 & e_{n-1} \\
& & & & & e_{n-1} & 0
\end{pmatrix},
$$

where $e_j = O(\sqrt{\epsilon})$, has eigenvalues $\lambda(T) \subset \{[-\nu, \nu] \cup [1 - \nu, 1 + \nu]\}$ with $\nu = O(\epsilon)$.

Hence, it seems as if for numerically relevant computations, we now would be faced with computing the eigenvalue decomposition of a tridiagonal matrix. This is not the case, however. By exploiting the special structure of the tridiagonal matrix, we can diagonalize it in two "sweeps" which compute the eigendecomposition of all "even" or "odd" $2 \times 2$ blocks on the diagonal (simultaneously), respectively. As we show in what follows, the fill-ins generated by these sweeps are of the same order as the perturbation in the eigenvalues and hence can be considered negligible.

LEMMA 4.3. *Let $T$ be a symmetric tridiagonal matrix with*

$$
\lambda(T) \subset [-\nu, \nu] \cup [1 - \nu, 1 + \nu],
$$

*where $\nu \stackrel{\text{def}}{=} \max_{\lambda \in \lambda(T)} \{\min(|\lambda - 1|, |\lambda|)\} \ll 1$. Then $\|T^2 - T\|_2 \leq \bar{\nu}$, where*

$$
\bar{\nu} \stackrel{\text{def}}{=} \nu + \nu^2. \tag{5}
$$

*Proof.* Let $Q$ be orthogonal and $E = \text{diag}(E_1, E_0)$ be diagonal, respectively, such that

$$
T = Q \begin{pmatrix} I + E_1 & \\ & E_0 \end{pmatrix} Q^T.
$$

Then $\|E\|_2 = \nu$, and

$$
T^2 = T + Q \left( \begin{pmatrix} E_1 & \\ & -E_2 \end{pmatrix} + E^2 \right) Q^T.
$$

Thus, $\|T^2 - T\|_2 \leq \| |E| + E^2 \|_2 = \bar{\nu}$. $\quad\square$

The next lemma gives bounds on the elements of the Givens rotation we will choose to diagonalize a $2 \times 2$ block and minimize the size of fill-ins.

LEMMA 4.4. *Let $G = \begin{pmatrix} c & s \\ -s & c \end{pmatrix}$ be a Givens rotation that diagonalizes a $2 \times 2$ symmetric matrix $\begin{pmatrix} \alpha_1 & \beta \\ \beta & \alpha_2 \end{pmatrix}$. Assume that without loss of generality $\beta > 0$ and define $\sigma \geq 0$ by*

$$
\sigma^2 \stackrel{\text{def}}{=} \left( \frac{\alpha_1 - \alpha_2}{2} \right)^2 + \beta^2. \tag{6}
$$

*Then $s$ and $c$ can be chosen such that*

$$
0 \leq |s| \leq \frac{\beta}{\sqrt{2}\,\sigma} \quad \text{and} \quad \frac{1}{\sqrt{2}} \leq c \leq 1.
$$

*Proof.* Let $c = \cos(\theta)$ and $s = \sin(\theta)$. Since we want to eliminate the off diagonal elements in $G\left(\begin{smallmatrix}\alpha_1 & \beta \\ \beta & \alpha_2\end{smallmatrix}\right)G^T$, we obtain

$$0 = (c^2 - s^2)\beta + 2cs\left(\frac{\alpha_2 - \alpha_1}{2}\right) = \beta\cos(2\theta) - \left(\frac{\alpha_2 - \alpha_1}{2}\right)\sin(2\theta).$$

If we choose

(7)
$$\cos(2\theta) = \frac{|\alpha_1 - \alpha_2|}{2\sigma},$$

with $\sigma$ as defined in (6), then

$$s^2 = \frac{1 - \cos(2\theta)}{2} = \frac{\beta^2}{2\sigma(\sigma + |\alpha_1 - \alpha_2|/2)},$$

and

$$c^2 = \frac{1 + \cos(2\theta)}{2} = \frac{\sigma + |\alpha_1 - \alpha_2|/2}{2\sigma},$$

as claimed.      □

In the following theorem we now show that, employing these Givens rotations, we can limit the size of the fill-in entries generated when applying these rotations to a tridiagonal matrix with eigenvalue clusters around 0 and 1.

THEOREM 4.5. *Let $T$ and $\bar{\nu}$ be as in Lemma 4.3. Let $G = \mathrm{diag}(I, \left(\begin{smallmatrix}c & s \\ -s & c\end{smallmatrix}\right), I)$ be the Givens rotation that diagonalizes one $2 \times 2$ diagonal block of $T$; i.e.,*

$$G \cdot \begin{pmatrix} \ddots & \overline{\beta} & & \\ \overline{\beta} & \alpha_1 & \beta & \\ & \beta & \alpha_2 & \underline{\beta} \\ & & \underline{\beta} & \ddots \end{pmatrix} \cdot G^T = \begin{pmatrix} \ddots & * & \gamma & \\ * & \tilde{\alpha}_1 & 0 & \delta \\ \gamma & 0 & \tilde{\alpha}_2 & * \\ & \delta & * & \ddots \end{pmatrix},$$

*where we assume that $\beta > 0$ without loss of generality. If $\beta > \sqrt{7}\,\bar{\nu}$ and $c$ and $s$ are chosen as suggested in Lemma 4.4, then*

$$\gamma \leq \sqrt{7}\,\bar{\nu} \quad and \quad \delta \leq \sqrt{7}\,\bar{\nu}.$$

*Proof.* Comparing corresponding entries in $T^2$ and $T$ and invoking Lemma 4.3, we know that there exist $\bar{\epsilon}$, $\underline{\epsilon}$, and $\epsilon_o$, $|\bar{\epsilon}|, |\underline{\epsilon}|, |\epsilon_o| \leq \bar{\nu}$, such that

(8)   $$\beta(\alpha_1 + \alpha_2) = \beta + \epsilon_o,$$
(9)   $$\overline{\beta}^2 + \alpha_1^2 + \beta^2 = \alpha_1 + \bar{\epsilon},$$
(10)  $$\beta^2 + \alpha_2^2 + \underline{\beta}^2 = \alpha_2 + \underline{\epsilon},$$
(11)  $$\overline{\beta}\beta \leq \bar{\nu}, \quad \underline{\beta}\beta \leq \bar{\nu}.$$

Using these identities, we have

$$\beta^2 - \alpha_1\alpha_2 = \frac{\alpha_1 + \alpha_2}{2}(1 - (\alpha_1 + \alpha_2)) + \frac{(\underline{\epsilon} + \bar{\epsilon}) - (\underline{\beta}^2 + \overline{\beta}^2)}{2},$$

and hence we can express $\sigma^2$ defined as in (6) as

$$
\begin{aligned}
\sigma^2 &= \left(\frac{\alpha_1 + \alpha_2}{2}\right)^2 + (\beta^2 - \alpha_1\alpha_2) \\
&= \frac{1}{4}\left(1 - \frac{\epsilon_o^2}{\beta^2}\right) + \frac{\overline{\epsilon} + \underline{\epsilon}}{2} - \frac{(\overline{\beta}\beta)^2 + (\underline{\beta}\beta)^2}{2\beta^2}.
\end{aligned}
$$

Thus,

$$
\sigma^2 \geq \frac{1}{4} - \bar{\nu} - \frac{5}{4}\left(\frac{\bar{\nu}}{\beta}\right)^2.
$$

Now let $\tau \geq 1$ be chosen such that $\beta > \tau\bar{\nu}$. Then

$$
\sigma^2 \geq \frac{1}{4} - \bar{\nu} - \frac{5}{4}\left(\frac{1}{\tau}\right)^2. \tag{12}
$$

Equations (11) together with $s \leq \frac{\beta}{\sqrt{2}\sigma}$ imply that

$$
\gamma = s\overline{\beta} \leq \frac{\bar{\nu}}{\sqrt{2}\sigma} \quad \text{and} \quad \delta = s\underline{\beta} \leq \frac{\bar{\nu}}{\sqrt{2}\sigma}.
$$

Using (12), it is now easy to show that $\tau \geq \sqrt{7}$ implies $\frac{1}{\sqrt{2}\sigma} \leq \tau$ and hence the result of the theorem.  □

As a consequence of Theorem 4.5, we are then able to compute the eigenvalue decomposition of a $2 \times 2$ diagonal block in a tridiagonal matrix $T$ with eigenvalue clusters at 0 and 1 such that the generated fill-in is negligible compared to the eigenvalue perturbation. Thus, the diagonalization of $T$ can be done by two sweeps of (potentially concurrent) $2 \times 2$ eigenvalue problems, as shown in Figure 4. In the first sweep, we diagonalize an "odd–even" $2 \times 2$ problem if the off diagonal entry is not too small, and set the fill entries to zero, or otherwise just zero the off diagonal entry. In the second sweep, we diagonalize the "even–odd" blocks. Since no more rotations follow, there is no need to zero out fill-in entries.

Theorem 4.5 shows that the Frobenius norm of the fill-in matrix introduced by the algorithm rr_diag shown in Figure 4 is bounded by $3\sqrt{n}\bar{\nu}$, which is of the same order as the perturbation in eigenvalues. The subroutine diag2, which is not shown here, computes the diagonalizing rotations as outlined in Lemma 4.4. Hence, Algorithm rr_diag is as numerically reliable as any other approach for diagonalizing $T$, albeit much cheaper due to its exploitation of the special structure of $T$.

**5. Numerical experiments.** In this section we report on some numerical experiments with the algorithms presented in this paper. All experiments were performed with Matlab Version 4.2a on a Sun Sparcstation iPX. For the reader wishing to experiment on his or her own, the Matlab files employed to generate these results can be retrieved via anonymous ftp from the pub/prism directory at ftp.super.org.

First, we apply the band-reduction algorithm bred of Figure 1 recursively to the trailing subblock of a $200 \times 200$ matrix with two eigenvalue clusters of size 50, each at $\lambda = \{-1, -2, 0, 1\}$. The radius of each cluster is $\epsilon 1.0e^3$, where $\epsilon$ is the machine precision. The drop threshold tau in bred is set to $\sqrt{7}\epsilon 1.0e^3$, and at each step the band width is chosen so as to decouple the problem in the middle. The succession of matrices generated is shown in Figure 5. The caption of each picture shows the

```
    function [Q, D] = rr_diag( A, Q, tau )
    %
    % Given a tridiagonal matrix A with eigenvalues 1 and 0, with
    % lambda(A) contained in [1-tau,1+tau] or [-tau,tau]
  5 % rr_diag computes an approximate eigendecomposition
    %         D = Q' * A * Q
    % where
    %         || D - Q'* A * Q||_Frobenius <= sqrt(7*n)*tau*(1+tau)

    [m,n] = size(A); if( m~=n ) error('non-square A'); end
 10 drop_threshold = sqrt(7)*tau*(1+tau);

    for j = 1:2:floor(n/2)*2         % diagonalize all (odd-even)
      k = j:j+1;                     % diagonal 2x2 matrices
      if (abs(A(j+1,j)) > drop_threshold)
        [G A(k,k)] = diag2( A(k,k) );
 15     if( j+2 <= n )
          A(j+2,k) = A(j+2,k)*G; A(k,j+2) = G'*A(k,j+2);
          A(j+2,j) = 0; A(j,j+2) = 0; % zero out negligible fill-ins
        end
        if( j-1 >= 1 )
 20       A(j-1,k) = A(j-1,k)*G; A(k,j-1) = G'*A(k,j-1);
          A(j-1,j+1) = 0; A(j+1,j-1) = 0;
        end
        Q(:,k) = Q(:, k)*G;
      end
 25 end
    for j = 2:2:floor((n-1)/2)*2     % diagonalize all (even-odd)
      k = j:j+1;                     % diagonal 2x2 matrices
      if (abs(A(j+1,j)) > drop_threshold)
        [G A(k,k)] = diag2( A(k,k) );
 30     if( j+2 <= n )
          A(j+2,k) = A(j+2,k)*G; A(k,j+2) = G'*A(k,j+2);
        % no more need to zero fill-ins
        end
        if( j-1 >= 1 )
 35       A(j-1,k) = A(j-1,k)*G; A(k,j-1) = G'*A(k,j-1);
        end
        Q(:,k) = Q(:, k)*G;
      end
    end
 40 D = diag(diag(A));
    return; end
```

FIG. 4. *Diagonalization of a tridiagonal matrix with eigenvalue clusters at 0 and 1.*

current matrix size being worked on and the band width to which it is to be reduced. At each step, we compute the residual

$$\delta \stackrel{\text{def}}{=} \|A_{original} * Q - Q * A_{current}\|_2.$$

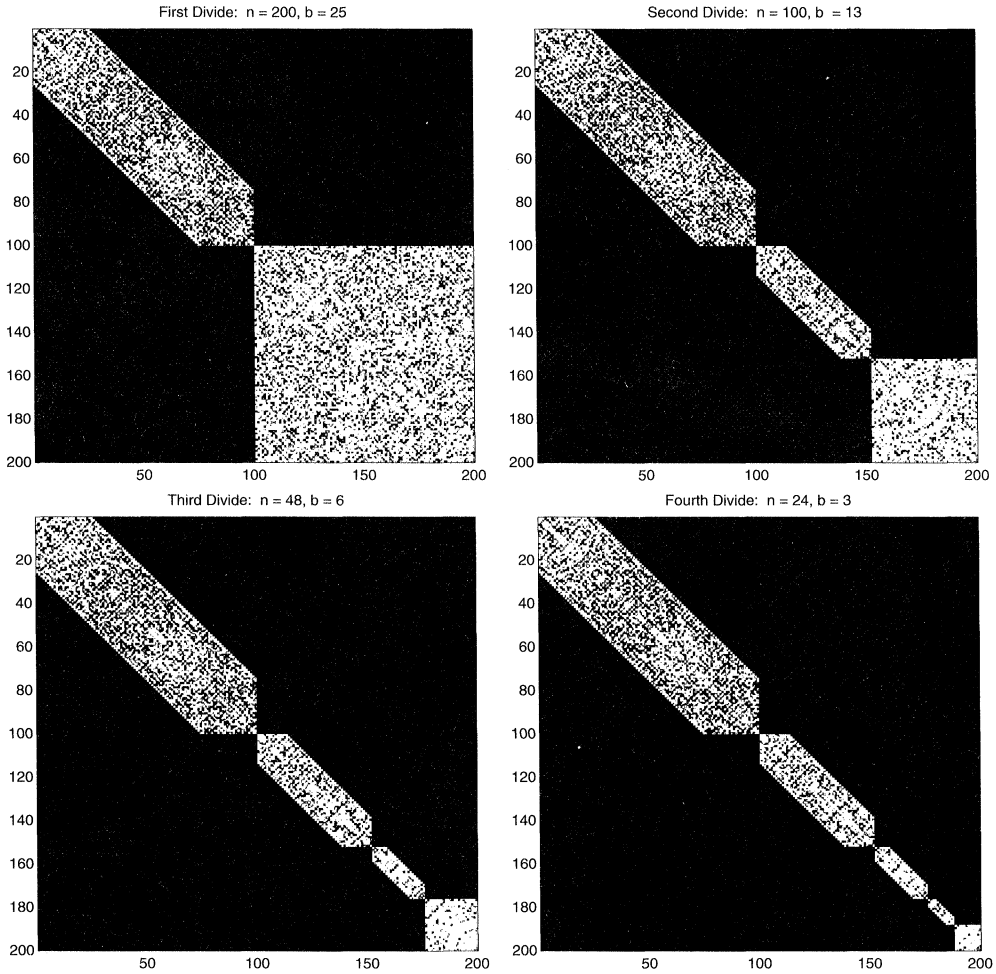We observe that $\delta \approx 7.2e^{-13}$, which, given a machine precision $\epsilon = 2.2e^{-16}$, is consis-

FIG. 5.   *Band reduction applied to trailing subblock of a* $200 \times 200$ *matrix with four distinct eigenvalue clusters.*

tent with our theory.

The same experiment, employing a matrix with 100 eigenvalues at 0 and 1 each and using the same eigenvalue perturbation and drop threshold, is shown in Figure 6. Note that it is sufficient to reduce the matrix to half the band width chosen in Figure 5 to achieve decoupling. We observe that $\delta \approx 2.7e^{-13}$. We also note that in both cases, the first, third, and fourth splits occurred at row (and column) 100, 176, and 188, respectively. The second split occurred at row 152 for Figure 5 and at row 150 for Figure 6.

To test the behavior of our rank-revealing tridiagonalization (RRDG), we compare it with the standard eigenvalue decomposition (EIG) and the QR factorization with column pivoting (QR); the results are presented in Table 1 and Table 2. Our test matrices are

    1. tridiagonal matrices with eigenvalue clusters of radius $p\,\epsilon$ generated by inserting random off diagonal perturbations of the order $\sqrt{p\,\epsilon}$ in the matrix shown in Example 4.2, and
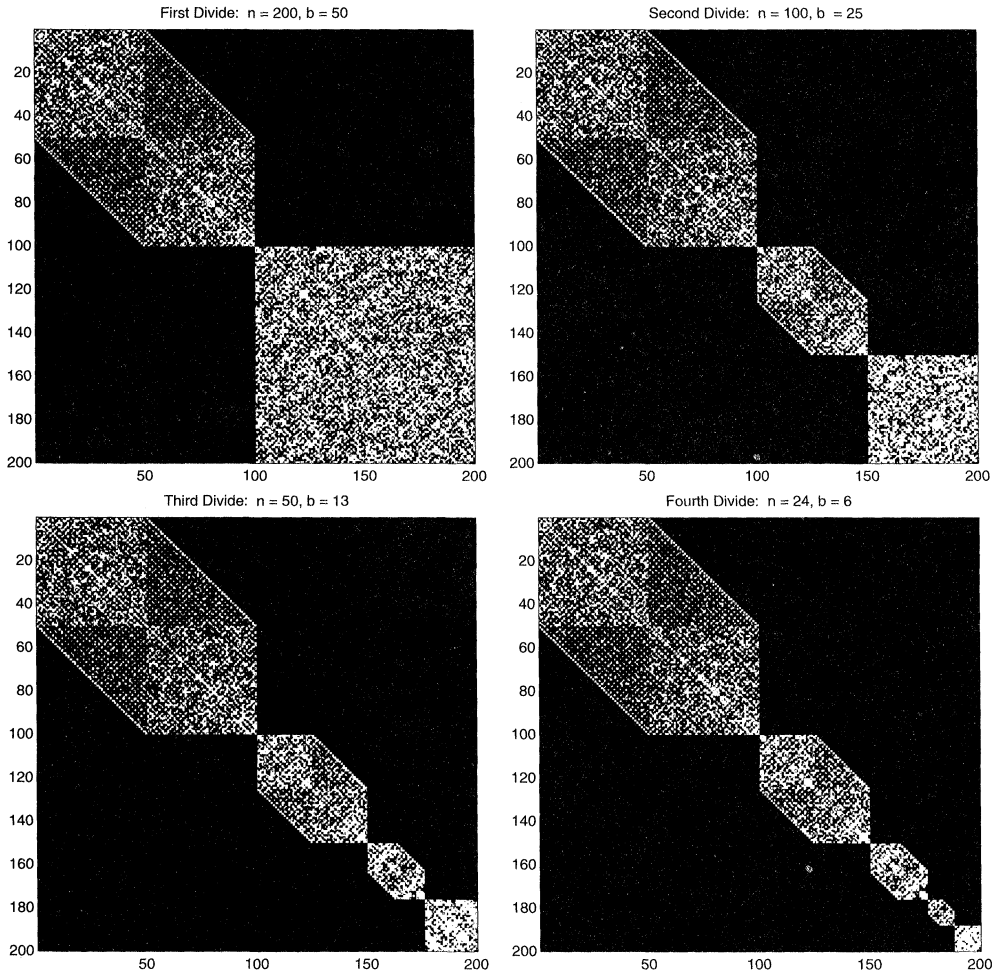
FIG. 6. *Band reduction applied to trailing subblock of a* $200 \times 200$ *matrix with two distinct eigenvalue clusters.*

2. matrices generated by symmetrically multiplying the matrices from Example 4.2 with orthogonal matrices generated via the QR factorization of a random matrix.

In the first case, we call `rr_diag`, listed in Figure 4. In the second case, we precede the call to `rr_diag` by a call to `tri_sbr`, as shown in Figure 2. The drop threshold for the divide-and-conquer tridiagonalization is set to $\sqrt{7}\,p\,\epsilon$, which is the same threshold as that employed in the two final diagonalization sweeps. For each of $p = 1, 10, 100$, we run 50 test cases each with matrix sizes 125, 250, and 375. RRDG and EIG both compute an eigenvalue decomposition $Q^T A Q = D$ with $D$ diagonal. We compute $\tilde{D} \stackrel{\text{def}}{=} \text{round}(D)$, i.e., round each diagonal entry to the nearest integer, and we report both the relative eigenvalue residual $\|Q^T A - \tilde{D}Q\|_F / \sqrt{n/2}$ (in Table 1) as well as the relative orthogonality residual $\|Q^T Q - I\|_F / \sqrt{n}$ (in Table 2). Note that $\sqrt{n/2}$ is an estimate of $\|A\|_F$. In the case of the QR factorization with pivoting, which computes $AP = QR$ for a permutation matrix $P$ and an upper triangular matrix $R$, we compute

TABLE 1
*Relative residual in subspace splitting.*

| Tridiagonal Matrices | | | | Full Matrices | | | |
|---|---|---|---|---|---|---|---|
| $n$ | $RRDG_{max}$ | $EIG_{max}$ | $QR_{max}$ | $n$ | $RRDG_{max}$ | $EIG_{max}$ | $QR_{max}$ |
| $p = 1$ | | | | $p = 1$ | | | |
| 125 | 5.3e−16 | 1.6e−15 | 1.7e−15 | 125 | 5.3e−14 | 1.7e−14 | 1.4e−14 |
| 250 | 5.0e−16 | 1.6e−15 | 3.8e−15 | 250 | 1.5e−13 | 3.3e−14 | 3.7e−14 |
| 375 | 4.9e−16 | 1.5e−15 | 5.6e−15 | 375 | 2.4e−14 | 3.8e−14 | 5.5e−14 |
| $p = 10$ | | | | $p = 10$ | | | |
| 125 | 3.5e−15 | 4.2e−15 | 2.2e−15 | 125 | 5.0e−15 | 6.0e−15 | 1.6e−14 |
| 250 | 3.3e−15 | 4.9e−15 | 5.1e−15 | 250 | 5.5e−15 | 3.0e−14 | 4.0e−14 |
| 375 | 3.4e−15 | 4.5e−15 | 4.3e−15 | 375 | 6.1e−15 | 4.1e−14 | 4.8e−14 |
| $p = 100$ | | | | $p = 100$ | | | |
| 125 | 3.3e−14 | 3.3e−14 | 2.7e−15 | 125 | 4.6e−14 | 3.5e−14 | 1.4e−14 |
| 250 | 3.2e−14 | 3.2e−14 | 6.8e−15 | 250 | 4.5e−14 | 5.2e−14 | 3.9e−14 |
| 375 | 3.2e−14 | 4.4e−14 | 6.6e−15 | 375 | 4.2e−14 | 3.2e−14 | 4.9e−14 |
| $p = 1000$ | | | | $p = 1000$ | | | |
| 125 | 3.3e−13 | 3.3e−13 | 2.5e−15 | 125 | 4.6e−13 | 3.5e−13 | 1.6e−14 |
| 250 | 3.2e−13 | 3.2e−13 | 4.1e−15 | 250 | 4.4e−13 | 3.4e−13 | 3.6e−14 |
| 375 | 3.2e−13 | 3.2e−13 | 6.2e−15 | 375 | 4.2e−13 | 3.2e−13 | 4.2e−14 |

TABLE 2
*Relative residual in orthogonality.*

| Tridiagonal Matrices | | | | Full Matrices | | | |
|---|---|---|---|---|---|---|---|
| $n$ | $RRDG_{max}$ | $EIG_{max}$ | $QR_{max}$ | $n$ | $RRDG_{max}$ | $EIG_{max}$ | $QR_{max}$ |
| $p = 1$ | | | | $p = 1$ | | | |
| 125 | 2.3e−16 | 1.2e−15 | 1.1e−15 | 125 | 2.1e−15 | 1.2e−14 | 1.7e−15 |
| 250 | 2.2e−16 | 1.3e−15 | 1.3e−15 | 250 | 3.0e−15 | 2.4e−14 | 2.4e−15 |
| 375 | 2.1e−16 | 1.2e−15 | 1.3e−15 | 375 | 3.6e−15 | 2.7e−14 | 2.8e−15 |
| $p = 10$ | | | | $p = 10$ | | | |
| 125 | 3.0e−16 | 2.8e−15 | 1.1e−15 | 125 | 1.4e−15 | 1.1e−14 | 1.7e−15 |
| 250 | 2.8e−16 | 3.0e−15 | 1.4e−15 | 250 | 1.9e−15 | 2.1e−14 | 2.3e−15 |
| 375 | 2.8e−16 | 2.8e−15 | 1.6e−15 | 375 | 3.4e−15 | 2.9e−14 | 2.9e−15 |
| $p = 100$ | | | | $p = 100$ | | | |
| 125 | 3.4e−16 | 1.1e−14 | 1.3e−15 | 125 | 1.4e−15 | 1.1e−14 | 1.7e−15 |
| 250 | 3.2e−16 | 2.0e−14 | 1.4e−15 | 250 | 1.9e−15 | 2.2e−14 | 2.4e−15 |
| 375 | 3.1e−16 | 1.9e−14 | 1.7e−15 | 375 | 2.3e−15 | 2.6e−14 | 2.9e−15 |
| $p = 1000$ | | | | $p = 1000$ | | | |
| 125 | 3.2e−16 | 1.0e−14 | 1.2e−15 | 125 | 1.4e−15 | 1.3e−14 | 1.8e−15 |
| 250 | 3.1e−16 | 2.3e−14 | 1.4e−15 | 250 | 1.9e−15 | 2.4e−14 | 2.4e−15 |
| 375 | 3.2e−16 | 3.3e−14 | 1.6e−15 | 375 | 2.3e−15 | 3.3e−14 | 2.9e−15 |

the rank

$$r \stackrel{\text{def}}{=} \max_{i} |r_{ii}| > \sqrt{7}\, p\, \epsilon$$

and $\tilde{A} \stackrel{\text{def}}{=} Q^T * A * Q$. We then report

$$\|\|\tilde{A}(1:r, 1:r)\|_F - \|A\|_F\|/\sqrt{n/2},$$

which should be small since $Q(1:r,:)$ is a basis for the range space of $A$. For each case, we report the worst residual.

We see that the divide-and-conquer tridiagonalization, followed by the two clean up sweeps over the resulting tridiagonal matrix, performs just as well as a full-fledged eigenvalue decomposition. In both cases, the residual in the subspace splitting is of $O(p\,\epsilon)$, as expected. The residual for QR factorization does not include the perturbation at the eigenvalue 1 as the other two approaches do and therefore is smaller in all cases. In any case, the computed orthogonal matrices are orthogonal up to machine precision. The $Q$ computed by the `eig` function in Matlab is slightly less orthogonal

since `eig` involves more transformations and as a result accumulates more rounding errors. Note that all three approaches are worse for a full matrix in the case $p = 1$. This is due to the fact that the roundoff errors in the orthogonal reductions are of the same order of machine precision. When $p$ is bigger, the roundoff errors are dominated by the perturbation in the eigenvalues, and hence RRDG and EIG behave about the same for tridiagonal and full matrices.

**6. Conclusions.** This paper introduced an algorithm for reducing a symmetric matrix with repeated eigenvalues to tridiagonal form. The algorithm progresses through a series of band reductions, each band-reduction stage forcing a decoupling of the band matrix into independent subblocks. Compared to the usual Householder tridiagonalization procedure, this approach can save up to 50% of the floating-point operations. We also developed a robust and inexpensive numerical procedure for diagonalizing the resulting tridiagonal matrix in the case where the matrix has only two eigenvalue clusters around 0 and 1. This case arises in eigenvalue decomposition algorithms based on invariant subspace approaches. Taken together, these two algorithms allow for a very efficient diagonalization of such matrices.

The algorithm can be generalized immediately to the reduction of unsymmetric matrices to Hessenberg form. The same irreducibility argument underlying Theorem 2.1 goes through for Hessenberg matrices. We also note that in exact arithmetic, conjugate transposed eigenvalue pairs would end up in the same block. However, since one triangle of a Hessenberg matrix is still full, the potential for computational savings is greatly reduced.

We mention that, apart from its divide-and-conquer nature and the resulting potential for parallelism, as well as its reduced operation count, our divide-and-conquer algorithm has another attractive feature. Since our algorithm, at least in the early stages, reduces matrices to banded form with a relatively wide band, it is easy to block the Householder transformations using the WY representation [11] or the compact WY representation [20], as has been described, for example, in [17]. In this fashion, one can easily capitalize on the favorable memory transfer characteristics of block algorithms.

## REFERENCES

[1] E. ANDERSON, Z. BAI, C. BISCHOF, J. DEMMEL, J. DONGARRA, J. DUCROZ, A. GREENBAUM, S. HAMMARLING, A. MCKENNEY, S. OSTROUCHOV, AND D. SORENSEN, *LAPACK User's Guide*, Society for Industrial and Applied Mathematics, Philadelphia, PA, 1992.

[2] ———, *LAPACK User's Guide Release 2.0*, Society for Industrial and Applied Mathematics, Philadelphia, PA, 1994.

[3] L. AUSLANDER AND A. TSAO, *A Divide-and-Conquer Algorithm for the Eigenproblem Via Complementary Invariant Subspace Decomposition*, Tech. report SRC-TR-89-003, Supercomputing Research Center, Institute for Defense Analysis, Bowie, MD, 1989.

[4] Z. BAI AND J. DEMMEL, *Design of a parallel nonsymmetric eigenroutine toolbox, Part* I, in Proc. of the Sixth SIAM Conference on Parallel Processing for Scientific Computing, 1993, pp. 391–398.

[5] C. BISCHOF, G. CORLISS, AND A. GRIEWANK, *Structured second- and higher-order derivatives through univariate Taylor series*, Optim. Meth. Software, 2 (1993), pp. 211–232.

[6] C. BISCHOF, S. HUSS-LEDERMAN, X. SUN, AND A. TSAO, *The PRISM project: Infrastructure and algorithms for parallel eigensolvers*, in Proc. of the Scalable Parallel Libraries Conference, Washington, DC, 1994, IEEE Computer Society, pp. 123–131.

[7] C. BISCHOF, X. SUN, AND B. LANG, *Parallel tridiagonalization through two-step band reduction*, in Proc. of Scalable High Performance Computing Conference, Knoxville, TN, 1994, IEEE Computer Society Press, pp. 23–27.

[8]   C. H. BISCHOF, *Fundamental Linear Algebra Computations on High-Performance Computers*, Informatik Fachberichte, Vol. 250, Springer-Verlag, Berlin, 1990.

[9]   C. H. BISCHOF AND J. J. DONGARRA, *A project for developing a linear algebra library for high-performance computers*, in Parallel and Vector Supercomputing: Methods and Algorithms, Graham Carey, ed., John Wiley, Somerset, NJ, 1989.

[10]  C. H. BISCHOF AND X. SUN, *A Framework for Band Reduction and Tridiagonalization of Symmetric Matrices*, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL, 1992, Preprint MCS-P298-0392.

[11]  C. H. BISCHOF AND C. F. VAN LOAN, *The WY representation for products of Householder matrices*, SIAM J. Sci. Stat. Comput., 8 (1987), pp. s2–s13.

[12]  H. CHANG, S. UTKU, M. SALAMA, AND D. RAPP, *A parallel Householder tridiagonalization stratagem using scattered square decomposition*, Parallel Comput., 6 (1988), pp. 297–311.

[13]  J. DONGARRA AND R. VAN DE GEIJN, *Reduction to condensed form for the eigenvalue problem on distributed-memory architectures*, Parallel Comput., 18 (1992), pp. 973–982.

[14]  J. DONGARRA AND S. HAMMARLING, *Evolution of Numerical Software for Dense Linear Algebra*, Oxford University Press, Oxford, UK, 1989.

[15]  J. J. DONGARRA, S. J. HAMMARLING, AND D. C. SORENSEN, *Block Reduction of Matrices to Condensed Form for Eigenvalue Computations*, Tech. report MCS–TM–99, Mathematics and Computer Science Division, Argonne National Laboratory, Argonne, IL, 1987.

[16]  G. H. GOLUB AND C. F. VAN LOAN, *Matrix Computations*, 2nd ed., The Johns Hopkins University Press, Baltimore, MD, 1989.

[17]  R. G. GRIMES AND H. D. SIMON, *Solution of large, dense symmetric generalized eigenvalue problems using secondary storage*, ACM Trans. Math. Software, 14, 3 (1988), pp. 241–256.

[18]  A. S. HOUSEHOLDER, *The Theory of Matrices in Numerical Analysis*, Dover, New York, 1964.

[19]  S. LEDERMAN, A. TSAO, AND T. TURNBULL, *A Parallelizable Eigensolver for Real Diagonalizable Matrices with Real Eigenvalues*, Tech. report TR-91-042, Supercomputing Research Center, Institute for Defense Analysis, Bowie, MD, 1991.

[20]  R. SCHREIBER AND C. VAN LOAN, *A storage efficient WY representation for products of Householder transformations*, SIAM J. Sci. Stat. Comput., 10 (1989), pp. 53–57.