

Parallel eigenvalue reordering in real Schur forms*

R. Granat[†], B. Kågström[‡] and D. Kressner[§]

May 12, 2008

Abstract

A parallel algorithm for reordering the eigenvalues in the real Schur form of a matrix is presented and discussed. Our novel approach adopts computational windows and delays multiple outside-window updates until each window has been completely reordered locally. By using multiple concurrent windows the parallel algorithm has a high level of concurrency, and most work is level 3 BLAS operations. The presented algorithm is also extended to the generalized real Schur form. Experimental results for ScaLAPACK-style Fortran 77 implementations on a Linux cluster confirm the efficiency and scalability of our algorithms in terms of more than 16 times of parallel speedup using 64 processor for large scale problems. Even on a single processor our implementation is demonstrated to perform significantly better compared to the state-of-the-art serial implementation.

Keywords: Parallel algorithms, eigenvalue problems, invariant subspaces, direct reordering, Sylvester matrix equations, condition number estimates

1 Introduction

The solution of large-scale matrix eigenvalue problems represents a frequent task in scientific computing. For example, the asymptotic behavior of a linear or linearized dynamical system is determined by the right-most eigenvalue of the system matrix. Despite the advance of iterative methods – such as Arnoldi and Jacobi-Davidson algorithms [3] – there are problems where a transformation method – usually the QR algorithm [14] – is preferred, even in a large-scale setting. In the example quoted above, an iterative method may fail to detect the right-most eigenvalue and, in the worst case, misleadingly predict stability even though the system is unstable [33]. The QR algorithm typically avoids this problem, simply because *all* and not only selected eigenvalues are computed. Also, iterative methods are usually not well suited for simultaneously computing a large portion of eigenvalues along with the associated invariant subspace. For example, an invariant subspace belonging to typically half of the eigenvalues needs to be computed in problems arising from linear-quadratic optimal control [37].

Parallelization of the QR algorithm is indispensable for large matrices. So far, only its two most important steps have been addressed in the literature: Hessenberg reduction and QR iterations, see [5, 10, 21], with the resulting software implemented in ScaLAPACK [6, 38]. The (optional) third step of reordering the eigenvalues, needed for computing eigenvectors and invariant subspaces, has not undergone parallelization yet. The purpose of this paper is to fill this gap, aiming at a complete and highly performing parallel implementation of the QR algorithm.

*Supported by the *Swedish Research Council/the Swedish Foundation for Strategic Research* (VR 621-2001-3284/SSF A3 02:128)

[†]Department of Computing Science and HPC2N, Umeå University, SE-901 87 Umeå, Sweden. granat@cs.umu.se

[‡]Department of Computing Science and HPC2N, Umeå University, SE-901 87 Umeå, Sweden. bokg@cs.umu.se

[§]Seminar für Angewandte Mathematik, ETH Zürich, Switzerland. kressner@sam.math.ethz.ch

1.1 Mathematical problem description

Given a general square matrix $A \in \mathbb{R}^{n \times n}$, computing the *Schur decomposition* (see, e.g., [14])

$$Q^T A Q = T \quad (1)$$

is the standard approach to solving *non-symmetric eigenvalue problems* (NEVPs), that is, computing eigenvalues and invariant subspaces (or eigenvectors) of a general dense matrix. In Equation (1), $T \in \mathbb{R}^{n \times n}$ is quasi-triangular with diagonal blocks of size 1×1 and 2×2 corresponding to real and complex conjugate pairs of eigenvalues, respectively, and $Q \in \mathbb{R}^{n \times n}$ is orthogonal. The matrix T is called the *real Schur form* of A and its diagonal blocks (i.e., its eigenvalues) can occur in any order along the diagonal.

For any decomposition of (1) of the form

$$Q^T A Q = T \equiv \begin{bmatrix} T_{11} & T_{12} \\ 0 & T_{22} \end{bmatrix} \quad (2)$$

with $T_{11} \in \mathbb{R}^{p \times p}$ for some integer p , the first p columns of the matrix Q span an *invariant subspace* of A corresponding to the p eigenvalues of T_{11} (see, e.g., [13]). Invariant subspaces are an important generalization of eigenvectors and often eliminate the need for calculating a potentially ill-conditioned set of eigenvectors in applications.

When applying the QR algorithm for computing a Schur form, one has little influence on the order of the eigenvalues on the diagonal. For computing an invariant subspace, however, the associated eigenvalues must be contained in the block T_{11} of (2). For example, in applications related to dynamical systems and control, it is often desired to compute the stable invariant subspace, in which case T_{11} is to contain all eigenvalues of A that have negative real part. To achieve this goal, the output of the QR algorithm needs to be post-processed.

In [2], a direct algorithm for reordering adjacent eigenvalues in the real Schur form (2) is proposed. For the special case of swapping eigenvalues in a tiny matrix T , partitioned as in (2) with block sizes $p, n - p \in \{1, 2\}$, the method is as follows:

- Solve the continuous-time Sylvester (SYCT) equation

$$T_{11}X - XT_{22} = \gamma T_{12}, \quad (3)$$

where γ is a scaling factor to avoid overflow in the right hand side.

- Compute the QR factorization

$$\begin{bmatrix} -X \\ \gamma I \end{bmatrix} = QR \quad (4)$$

using Householder transformations (elementary reflectors)¹.

- Apply Q in the similarity transformation of T :

$$\tilde{T} = Q^T T Q \quad (5)$$

- Standardize 2×2 block(s) if any exists.

In the method above, a swap is performed tentatively to guarantee backward stability by rejecting a swap where $Q^T T Q$ would be perturbed too far from real Schur form.

For general n , a real Schur form $T \in \mathbb{R}^{n \times n}$ can be reordered by subsequently swapping its individual blocks. Note that we are restricted to swapping *adjacent* blocks to maintain the real Schur form. This limits our choice of sorting algorithm essentially to bubble sort. By applying such a bubble-sort procedure, all selected eigenvalues (usually represented by a boolean vector *select*) are moved step by step towards the top-left corner of T . This procedure is implemented in the LAPACK [1] routine `DTRSEN`.

This paper presents a parallel algorithm for reordering eigenvalues in Schur forms, partly based on a blocked variant described in [31]. Furthermore, we extend our techniques to reordering eigenvalues in a generalized real Schur form of a matrix pair [24].

¹For $n = 2$, a single Givens rotation is used to compute the QR factorization.

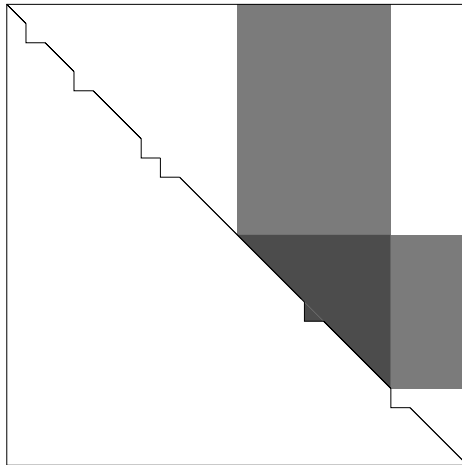


Figure 1: Working with computational windows (dark) and delaying the updates for level 3 BLAS operations on *update regions* (light) in the matrix T .

2 Serial blocked algorithms for eigenvalue reordering

The eigenvalue reordering algorithm implemented in LAPACK [2] sorts a single 1×1 or 2×2 block at a time to the top left corner of T . This leads to matrix multiplications with tiny orthogonal matrices (2×2 , 3×3 , 4×4) to update the rows and columns of T and Q . Moreover, a significant portion of T is accessed before the next eigenvalue can be processed. Consequently, the performance is in the range of level 1 and 2 BLAS [12, 7].

Following the ideas from [30, 31], we may instead reorder several eigenvalues simultaneously. For this purpose, a local *computational window* on the diagonal of T is chosen. All selected eigenvalues residing in this window are reordered locally to the top left corner of the window. Most importantly, the application of the corresponding orthogonal transformations is restricted to this window. Once the local reordering is complete, the transformations are applied in factorized or accumulated form to update the rest of T and the global orthogonal transformation matrix Q , see Figure 1. Using the accumulated local transformation matrix results in decently large matrix-matrix multiplications and thus benefits from level 3 BLAS performance. In the next step of the block algorithm, the window is slid towards the top left corner of T until the eigenvalues reordered in the previous step reside in the bottom right corner of the new window. Continuing this process, the window eventually reaches the top left corner of T and all eigenvalues within the window have been properly reordered. For more details on this block algorithm and its implementation, we refer to [31].

In general, the described approach leads to a significant performance gain by improving the memory access pattern and diminishing the number of cache misses. Notice that a similar technique was employed for the variants of the QR algorithm presented in [8, 9]. The idea of delaying updates is however not new (see, e.g., [11, 30] and the references therein). The performance of the block algorithm is controlled by a few parameters, most notably the size of the computational window n_{win} , and the maximum number of eigenvalues to reorder locally inside a window, n_{eig} . A recommended choice is $n_{\text{eig}} = n_{\text{win}}/2$ [31]. Other parameters to tune are r_{mmult} , which defines the threshold for when to use the local orthogonal transformations in their factorized form (Householder transformations and Givens rotations) or their accumulated form by matrix multiplication, and n_{slab} , which is used to divide the updates of the rows of T into blocks of columns for improved memory access pattern in case the orthogonal transformations are applied in their factorized form.

3 Parallel blocked algorithms for eigenvalue reordering

In the following, a parallel algorithm for reordering eigenvalues is proposed. In principle, the serial block algorithm described in Section 2 already admits a quite straightforward parallel implementation. There are some issues that need to be handled, e.g., the proper treatment of a computational window that is shared by several processes. In Sections 3.1 and 3.2 below, we provide the details of such a parallelization. As expected, this approach gives good node performance but it leads to poor scalability, simply because there are only a limited number of processes active during the local reordering. To avoid this, we suggest the use of multiple concurrent computational windows in Section 3.4.

We adopt the ScaLAPACK conventions [6, 38] of the parallel distributed memory (DM) environment, as follows:

- The parallel processes are organized into a rectangular $P_r \times P_c$ mesh labelled from $(0, 0)$ to $(P_r - 1, P_c - 1)$ according to their specific position indices in the mesh.
- The matrices are distributed over the mesh using 2-dimensional (2D) block cyclic mapping with the block sizes m_b and n_b in the row and column dimensions, respectively.

Since the matrices T and Q are square, we assume throughout this paper that $n_b = m_b$, i.e., the matrices are partitioned in square blocks. To simplify the reordering in the presence of 2×2 blocks, we also assume that $n_b \geq 3$ to avoid the situation of having two adjacent 2×2 blocks spanning over three different diagonal blocks in T . Moreover, we require T and Q to be aligned such that the blocks T_{ij} and Q_{ij} are held by the same process, for all combinations of i and j , $1 \leq i, j \leq \lceil n/n_b \rceil$. Otherwise, shifting T and/or Q across the process mesh before (and optionally after) the reordering is necessary.

Below, we refer to an *active process* as a process that holds all or some part of a computational window in T .

3.1 The computational window in the parallel environment

We restrict the size of the computational window to the block size used in the data layout. This means that a computational window can only be in two states: either it is completely held by one single block or it is shared by at most four data layout blocks: two neighboring diagonal blocks, one subdiagonal block (in presence of a 2×2 block residing on the block borders) and one superdiagonal block.

3.2 Moving eigenvalues inside a data layout block

Depending on the values of n_b and n_{win} , each diagonal block of T is locally reordered, moving $k \leq n_{\text{eig}}$ selected eigenvalues from the bottom towards the top of the block, see Figure 2. Notice that the number of blocks distributed over the process mesh in general is much larger than what is indicated in Figure 2; a small number of blocks in this case (and below) is used to simplify the illustration.

Before the window is moved to its next position inside the block, the resulting orthogonal transformations from the reordering in the current position are broadcasted in the process rows and columns corresponding to the current block row and column in T (and Q), see Figure 3. The subsequent updates are performed independently and in parallel. In principle, no other communication operations are required beside these broadcasts.

Given one active process working alone on one computational window the possible parallel speedup during the update phase is limited by $P_r + P_c$.

3.3 Moving eigenvalues across the process borders

When a computational window reaches the top-left border in the diagonal block, the current eigenvalue cluster must be reordered across the border into the next diagonal block of T . This

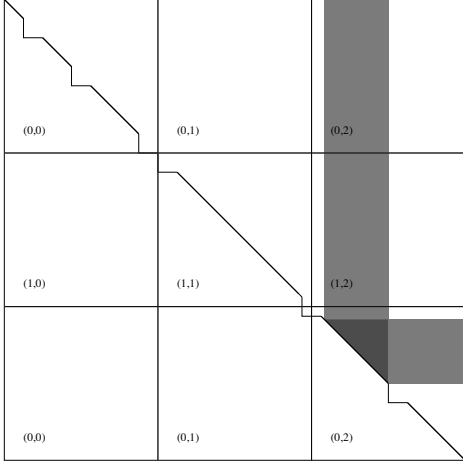


Figure 2: The Schur form T distributed over a process mesh of 2×3 processors. The computational window (dark) is completely local and the update regions (light) are shared by the corresponding process row and column.

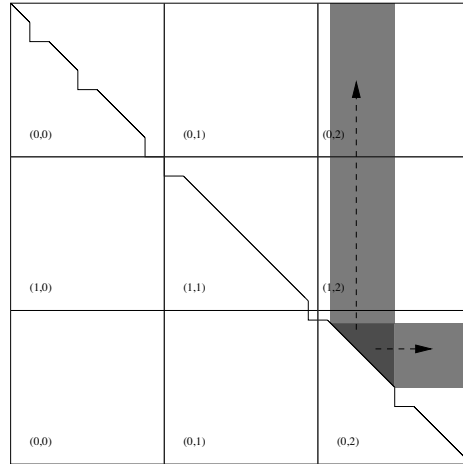


Figure 3: Broadcasting of the orthogonal transformations along the current process row and column of a 2×3 process mesh.

forces the computational window to be shared by more than one data layout block and (optionally) more than one process. The restrictions imposed on its size ensures that a window cannot be shared by more than four data layout blocks and by at most four different processes which together form a submesh of maximum size 2×2 , see Figure 4.

To be able to maximize the work performed inside each diagonal block before crossing the border and to minimize the required communication for the cross border reordering, it is beneficial to be able to control the size of the shared windows by an optional parameter $n_{\text{crb}} \leq n_{\text{win}}$ that can be adjusted to match the properties of the target architecture.

The processes holding the different parts of the shared $n_{\text{crb}} \times n_{\text{crb}}$ window now cooperate to bring the window across the border, as follows:

- The on-diagonal active processes start by exchanging their parts of the window and receiving the off-diagonal parts from the two other processes. The cross border window causes updates in T and Q that span over parts of two block rows or columns. Therefore, the processes in the corresponding process rows and columns exchange blocks with their neighbors as preparation for the (level 3) updates to come, see Figure 5. The total amount of matrix elements from T and Q exchanged over the border in both directions is $n_{\text{crb}} \cdot (2n - n_{\text{crb}} - 1)$.
- The on-diagonal active processes compute the reordering for the current window. This requires some replication of the computations on both sides of the border which is executed concurrently. Since the total work is dominated by the off-diagonal updates, the on-diagonal work is negligible.
- Finally, the local orthogonal transformation matrices are broadcasted (either in accumulated or factorized form) along the corresponding process rows and columns, and used in level 3 updates, see Figure 6.

We remark that in principle the updates are independent and can be performed in parallel without any additional communication or redundant work. But if the local orthogonal transformations are applied in their factorized form, each processor computing an update will also compute parts of T and/or Q that are supposed to be computed by another processor at the other side of the border. In case of more than one process in the corresponding mesh dimension, this causes duplicated

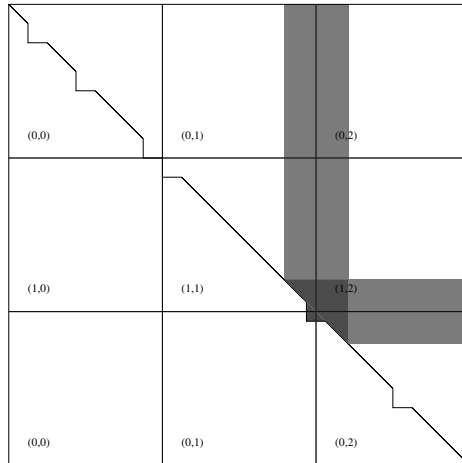


Figure 4: The Schur form T distributed over a process mesh of 2×3 processors. The computational window (dark) is shared by four distinct processes. The update regions (light) is shared by the corresponding two process rows and process columns.

work in the cross border updates, as well. Our remedy is to use a significantly lower value for r_{mmult} favoring matrix multiplication in the cross border reordering.

3.4 Introducing multiple concurrent computational windows

By using multiple concurrent computational windows, we can work on at least $k_{\text{win}} \leq \min(P_r, P_c)$ adjacent windows at the same time, computing local reordering and broadcasting orthogonal transformations for updates in parallel. With $k_{\text{win}} = \min(P_r, P_c)$ and a square process mesh ($P_r = P_c$), the degree of concurrency becomes $P_r \cdot P_c$, see Figure 7.

When all k_{win} windows reach the process borders (see also Section 5.4), they are moved into the next diagonal block as described in the previous section, but in two phases. Since each window requires cooperation between two adjacent process rows and columns, we number the windows by the order in which they appear on the block diagonal of T and start by moving all windows with an odd label across the border, directly followed by all windows with an even label. Care has to be taken to assure that no processor is involved in more than one cross border window at the same time. For example, if $k_{\text{win}} = \min(P_r, P_c) > 1$ is an odd number, the last labelled window will involve processors which are also involved in the first window. In such a case, the last window is reordered across the border after the second (even) phase is completed.

This two-phase approach gives an upper limit of the concurrency of the local reordering and data exchange phases of the cross border reordering as $k_{\text{win}}/2$, which is half of the concurrency of the pre-cross border part of the algorithm.

We present a high-level description of our parallel multi-window method in Algorithms 1–2. Notice that in the presence of 2×2 blocks in the Schur form, the computed indices in the algorithm are subject to slight adjustments, see Section 5.

4 Performance analysis

In this section, we analyze the parallel performance of Algorithms 1–2 and derive a model of the *parallel runtime* using p processes,

$$T_p = T_a + T_c, \quad (6)$$

where T_a and T_c denote the *arithmetic* and *communication* (synchronization) runtimes, respectively. We assume block cyclic data distribution of $T \in \mathbb{R}^{n \times n}$ over a square $\sqrt{p} \times \sqrt{p}$ process mesh

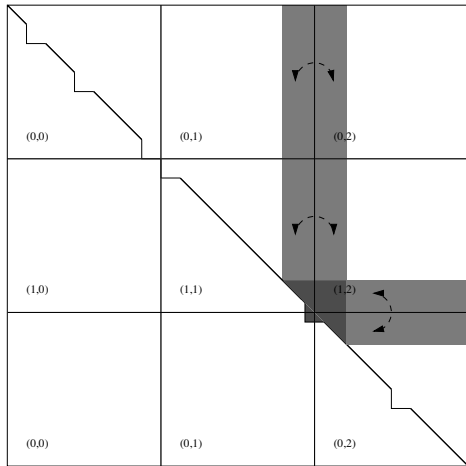


Figure 5: Exchanges of data in adjacent process rows and columns for updates associated with cross border reordering.

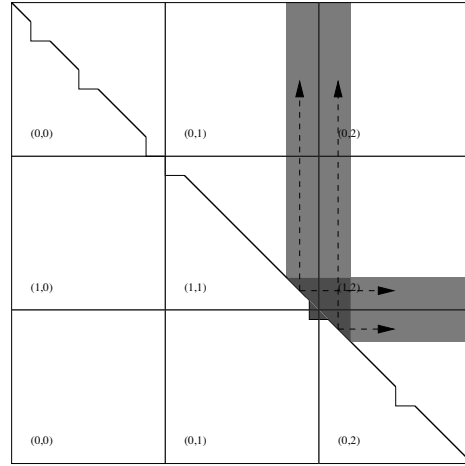


Figure 6: Broadcasts of the computed orthogonal transformations in the current processor rows and columns for cross border reordering.

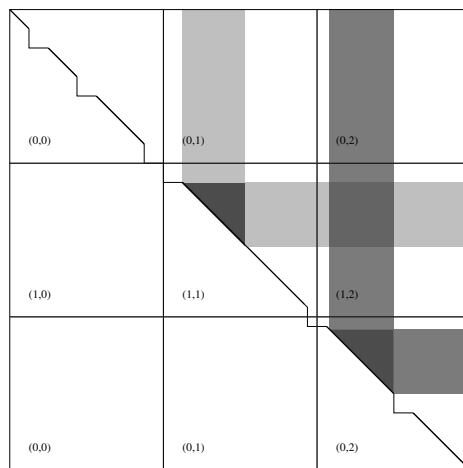


Figure 7: Using multiple (here two) concurrent computational windows. The computational windows (dark) are local but some parts of the update regions (light regions) are shared.

Algorithm 1 Parallel blocked algorithm for reordering a real Schur form (main part)

Input: A matrix $T \in \mathbb{R}^{n \times n}$ in real Schur form and a subset of selected eigenvalues Λ_s , closed under complex conjugation. Parameters n_b (data layout block size), P_r and P_c (process mesh sizes), $myrow$ and $mycol$ (process mesh indices), n_{crb} (maximum cross border window size), k_{win} (maximum number of concurrent windows), n_{ev} maximum number of eigenvalues simultaneously reordered.

Output: T is overwritten by a matrix $\tilde{T} \in \mathbb{R}^{n \times n}$ in ordered real Schur form and an orthogonal matrix $Q \in \mathbb{R}^{n \times n}$ such that $\tilde{T} = Q^T T Q$. For some integer j , the set Λ_s is the union of eigenvalues belonging to the j upper-left-most diagonal blocks of \tilde{T} .

```
Let  $W = \{\text{The set of computational windows}\} = \emptyset$ 
 $i_{ord} \leftarrow 0$  %  $i_{ord} = \#$ number of eigenvalues already in order
while  $i_{ord} < \#\Lambda_s$  do
  for  $j \leftarrow 0, \dots, k_{win} - 1$  do
     $s_{win} = (\lfloor i_{ord}/n_b \rfloor + j) \cdot n_b + 1$ 
    Add a window to  $W$  for  $T(s_{win} : s_{win} + n_b - 1, s_{win} : s_{win} + n_b - 1)$ 
  end for
  % Reorder each window to top-left corner of corresponding diagonal block.
  for each window  $w \in W$  in parallel do
    if ( $myrow, mycol$ ) owns  $w$  then
      Find first  $n_{ev} \leq n_{eig}$  eigenvalues in  $\Lambda_s$  from top of my diagonal block.
       $i_{hi} \leftarrow \max\{i_{ord} + 1, s_{win}\} + n_{ev}$ 
    end if
    % Reorder these eigenvalues window-by-window to top of diagonal block.
    while  $i_{hi} > i_{ord} + n_{ev}$  do
       $i_{low} \leftarrow \max\{i_{ord} + 1, \max\{i_{hi} - n_{win} + 1, s_{win}\}\}$ 
      if ( $myrow, mycol$ ) owns  $w$  then
        Apply standard eigenvalue reordering to the active window  $T(i_{low} : i_{hi}, i_{low} : i_{hi})$  to reorder  $k \leq n_{ev}$ 
        eigenvalues in  $\Lambda_s$  to the top.
        Broadcast local orthogonal transformation  $U$  in process row  $myrow$ .
        Broadcast local orthogonal transformation  $U$  in process column  $mycol$ .
        else if ( $myrow, mycol$ ) needs  $U$  for updates then
          Receive  $U$ 
        end if
        Update  $T(i_{low} : i_{hi}, i_{hi} + 1 : n) \leftarrow U^T T(i_{low} : i_{hi}, i_{hi} + 1 : n)$  in parallel.
        Update  $T(1 : i_{low} - 1, i_{low} : i_{hi}) \leftarrow T(1 : i_{low} - 1, i_{low} : i_{hi})U$  in parallel.
        Update  $Q(1 : n, i_{low} : i_{hi}) \leftarrow Q(1 : n, i_{low} : i_{hi})U$  in parallel.
         $i_{hi} \leftarrow i_{low} + k - 1$ 
      end while
    end for
     $i_{ord} \leftarrow i_{ord} + n_{ev}$  [top-left window in  $W$ ]
    % Reorder selected clusters across block (process) border.
    Apply Algorithm 2 to reorder each computational windows across (process) borders.
  end while
```

Algorithm 2 Parallel blocked algorithm for reordering a real Schur form (cross border part)

Input and Output: See Algorithm 1. Additional input: set of computational windows W .

```
for each odd window  $w \in W$  in parallel do
  Form  $2 \times 2$  process submesh  $G = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$  corresponding to  $w$ 
   $i_{hi} \leftarrow \min\{i_{ilo} - 1, s_{win} + n_{crb}/2\}$ ,  $i_{ilo} \leftarrow \max\{i_{hi} - n_{crb} + 1, i_{ord}\}$ 
  if ( $myrow, mycol$ )  $\in G$  then
    Exchange data in  $G$  to build  $w$  at  $(0, 0)$  and  $(1, 1)$ 
    if ( $myrow, mycol$ )  $\in \{(0, 0), (1, 1)\}$  then
      Reorder  $k$  eigenvalues in  $\Lambda_s$  to top of  $w$ .
      Broadcast local orthogonal transformation  $U$  in process row  $myrow$ .
      Broadcast local orthogonal transformation  $U$  in process column  $mycol$ .
    end if
  end if
  if  $myrow \in \text{prows}(G)$  or  $mycol \in \text{pcols}(G)$  then
    Exchange local parts of  $T(i_{low} : i_{hi}, i_{hi} + 1 : n)$ ,  $T(1 : i_{low} - 1, i_{low} : i_{hi})$  and  $Q(1 : n, i_{low} : i_{hi})$  for updates with
    neighboring processes in parallel.
    Receive  $U$ .
    Update  $T(i_{low} : i_{hi}, i_{hi} + 1 : n) \leftarrow U^T T(i_{low} : i_{hi}, i_{hi} + 1 : n)$  in parallel.
    Update  $T(1 : i_{low} - 1, i_{low} : i_{hi}) \leftarrow T(1 : i_{low} - 1, i_{low} : i_{hi})U$  in parallel.
    Update  $Q(1 : n, i_{low} : i_{hi}) \leftarrow Q(1 : n, i_{low} : i_{hi})U$  in parallel.
  end if
end for
if  $i_{ilo} = i_{ord}$  then
   $i_{ord} \leftarrow i_{ord} + k$ 
end if
for each even window  $w \in W$  in parallel do
  % Similar algorithm as for the odd case above.
end for
```

using the block factor n_b (see Section 3). We define t_a , t_s and t_w as the *arithmetic time* to perform a floating point operation (flop), the *start-up time* (or node latency) for sending a message in our parallel computer system, the *per-word transfer time* or *the inverse of the bandwidth* as the amount of time it takes to send one data word (e.g., a double precision number) through one link of the interconnection network. Usually, t_s and t_w are assumed to be constants while t_a is a function of the data locality. The communication cost model for a single point-to-point communication can be approximated by $t_{p2p} = t_s + t_w l$, where l denotes the message size in words, regardless of the number of links traversed [15]. However, for a *one-to-all broadcast* or its dual operation *all-to one reduction* in a certain scope (e.g., a process row or process column), we assume that such an operation is performed using a hypercube-based algorithm like *recursive doubling*, i.e., in $O(\log_2 p_*)$ steps, where p_* is the number of processors in the actual scope.

Reordering k eigenvalues in a Schur form has an arithmetic complexity of $T_s = O(kn^2)t_a$. The exact cost depends on the distribution of the eigenvalues over the diagonal of the Schur form. In the following, we assume $1 \ll k < n$ which implies that a majority of the updates are performed using matrix multiplication.

Given that the selected eigenvalues are uniformly distributed over the diagonal of T , the arithmetic cost of executing Algorithms 1–2 can be modelled as

$$T_a = \left(\frac{kn^2 - 3knn_{\text{win}}}{p} + \frac{3knn_{\text{win}}}{\sqrt{p}} \right) t_a, \quad (7)$$

where the first term is the cost of the GEMM (general matrix-matrix multiplication) updates, which is divided evenly between the p involved processors, and the second term describes the cost for computing the local and cross border reordering in the $n_{\text{win}} \times n_{\text{win}}$ computational windows, where the diagonal processors are working and the off-diagonal processors are idle waiting for the broadcasts to start.

From a global view, bubble-sorting the eigenvalues can be seen as bubble-sorting the diagonal $n_b \times n_b$ blocks of the block layout in $O((n/n_b)^2)$ steps. In each step, each local window is reordered to the top-left corner of the corresponding diagonal block and reordered across the border. By working with \sqrt{p} concurrent windows in parallel (which is the upper limit for a square grid), the communication cost T_c can be modelled as

$$T_c = \frac{D_T^2}{\sqrt{p}} (t_{\text{lbcast}} + t_{\text{cwin}} + t_{\text{cbcast}} + t_{\text{rcexch}}), \quad (8)$$

where $D_T = \lceil n/n_b \rceil$ is the number of diagonal blocks in T , and t_{lbcast} , t_{cwin} , t_{cbcast} , and t_{rcexch} model the cost of the broadcasts local to the corresponding process row and column, the point-to-point communications associated with the construction of the cross border window, the cross border broadcasts in the two corresponding process rows and columns, and the data exchange between neighboring processes in the corresponding process rows and columns for the cross border updates, respectively.

The number of elements representing the factorized form of the orthogonal transformation can be approximated by n_{win}^2 and the broadcasts are initiated once for each time a window is moved to the next position in the corresponding block. Based on these observations, we model t_{lbcast} as

$$t_{\text{lbcast}} = 2 \frac{n_b}{n_{\text{win}}} (t_s + n_{\text{win}}^2 t_w) \log_2 \sqrt{p}. \quad (9)$$

By a similar reasoning and taking into account the lower degree of concurrency in the cross border operations, see Section 3.1, we model $t_{\text{cwin}} + t_{\text{cbcast}}$ by

$$t_{\text{cwin}} + t_{\text{cbcast}} = 4k_{\text{cr}} (t_s + n_{\text{crb}}^2 t_w) (\log_2 \sqrt{p} + 1), \quad (10)$$

where k_{cr} is the number of passes necessary for bringing the whole locally collected eigenvalue cluster across the border ($k_{\text{cr}} = 1$ if $n_{\text{crb}} = n_{\text{win}}$).

The exchanges of row and column elements in T and column elements in Q between neighboring processes suffer from a lower degree of concurrency. Moreover, we cannot guarantee that the send

and receive may take place concurrently rather than in sequence. The average length of each matrix row or column to exchange is $n/2$ and about half of the cross border window is located on each side of the border. To sum up, we model t_{rcexch} as

$$t_{\text{rcexch}} = 12k_{\text{cr}} \left(t_s + \frac{n \cdot n_{\text{crb}}}{4\sqrt{p}} t_w \right). \quad (11)$$

Using the derived expressions and assuming $k_{\text{cr}} = 1$, i.e., $n_{\text{crb}} = n_{\text{win}}$, T_c can be boiled down to the approximation

$$\left(\left(\frac{2n^2}{n_b \cdot n_{\text{win}}} + \frac{4n^2}{n_b^2} + \frac{12n^2}{\log_2 \sqrt{p} \cdot n_b^2} \right) t_s + \left(\frac{2n^2 \cdot n_{\text{win}}}{n_b} + \frac{4n^2 \cdot n_{\text{win}}^2}{n_b^2} + \frac{3n^3 \cdot n_{\text{win}}}{\log_2 \sqrt{p} \cdot n_b^2} \right) t_w \right) \frac{\log_2 \sqrt{p} + 1}{\sqrt{p}}.$$

The dominating term is the last fraction of the part associated with t_w and comes from the data exchange associated with the updates (see Figure 5); it is of order $O(n^3/(n_b \cdot \sqrt{p}))$ assuming $n_{\text{win}} \approx n_b$ and in the general case the communication cost in the algorithm will be dominated by the size of this contribution. The influence of this term is diminished by choosing $n_b \cdot \sqrt{p}$ as close to n as possible and thereby reducing the term closer to $O(n^2)$, which may be necessary when the arithmetic cost is not dominated by GEMM updates. For example, for $n = 1500$, $n_b = 180$ and $p = 64$, we have $n_b \cdot \sqrt{p} = 1440 \approx n$ (see also Section 6). In general and in practice, we will have $n_b \cdot \sqrt{p} = n/l$, where $l < n$ is the average number of row or column blocks of T distributed to each processor in the cyclic distribution. The scenario to strive for is $l \ll k$, where k is the number of selected eigenvalues. Then we perform significantly more arithmetics than communication which is a rule of thumb in all types of parallel computing. If this is possible depends on the problem, i.e., the number of selected eigenvalues and their distribution over the diagonal of T . Our derived model is compared with observed data in Section 6.

5 Implementation issues

In this section, we address some detailed but important implementation issues.

5.1 Very few eigenvalues per diagonal block

In the extreme case, there is only one eigenvalue per diagonal block to reorder. In such a situation, we could think of gathering larger clusters of eigenvalues over the borders as quickly as possible to increase the node performance of the local reordering. However, this turns out to be counterproductive; the scalability and the overall efficiency in hardware utilization decreases as the number of concurrent windows (and consequently the concurrently working processors) decreases. In general, the node performance of our parallel algorithm will vary with the number of selected eigenvalues to be reordered and their distribution across the main block diagonal of T .

5.2 Splitting clusters in cross border reordering

Sometimes when we cross the process borders, not all the eigenvalues of the cluster can be moved across the border because of lack of reserved storage at the receiver, e.g., when there are other selected eigenvalues on the other side which occupy entries in T close to the border. Then the algorithm splits the cluster in two parts, roughly half on each side of the border, to keep the cross border region as small as necessary. In such a way, more work is performed with a higher rate of concurrency in the pre-cross border phase: the eigenvalues left behind are pushed closer to the border and reordered across the border in the next cross border sweep over the diagonal blocks.

5.3 Size of cross border window and shared 2×2 blocks

To simplify the cross border reordering in the presence of any 2×2 block being shared on the border, we keep the top-left part of the cross border window at a minimum dimension 2. This

ensures that the block does not stay in the same position causing an infinite loop. For a similar reason we keep the bottom-right part of the window at a minimum dimension 3 if a non-selected eigenvalue is between the border and a selected 2×2 block.

In our implementation, the size of the cross border window n_{crb} is determined by the parameter n_{ceig} which controls the number of eigenvalues that cross the border. In practice, $n_{\text{crb}} \approx 2n_{\text{ceig}}$ except for the case when it is not possible to bring all eigenvalues in the current cluster across the border for reasons discussed above.

5.4 Detection of reordering failures and special cases

Due to the distributed nature of the parallel algorithm, failures in reordering within any computational window and special cases, like when no selected eigenvalue was found or moved in a certain computational window, must be detected at runtime and signaled to all processors in the affected scope (e.g., the corresponding processor row(s), processor column(s) or the whole mesh) at specific synchronization points. In the current implementation, all processors are synchronized in this way right before and right after each computational window is moved across the next block border.

5.5 Organization of communications and updates

In a practical implementation of Algorithms 1–2, the broadcasts of the orthogonal transformations, the data exchanges for cross border row and column updates and the associated updates of the matrices T and Q should be organized to minimize any risk for bottlenecks caused by redundant synchronizations. For example, all broadcasts in the row direction should be started before any column oriented broadcast starts to ensure that no pair of broadcasts have intersecting scopes (see Figure 7). In practice, this approach also encourages an implementation that performs all row oriented operations, computations and communications, before any column oriented operations take place. For such a variant of the algorithm, all conclusions from Section 4 are still valid. This technique also paves the way for a greater chance of overlapping communications with computations, possibly leading to better parallel performance.

5.6 Condition estimation of invariant subspaces

Following the methodology of LAPACK, our implementation also computes condition numbers for the invariant (deflating) subspaces and the selected cluster of eigenvalues using the recently developed software package SCASY [16, 17, 36] and adopting a well-known matrix norm estimation technique [20, 22, 26] in combination with parallel high performance software for solving different transpose variants of the triangular (generalized) Sylvester equations.

6 Experimental results

In this section, we demonstrate the performance of a ScaLAPACK-style parallel Fortran 77 implementation of Algorithms 1–2 called PBDTRSEN. All experiments were carried out in double precision real arithmetics ($\epsilon_{\text{mach}} \approx 2.2 \times 10^{-16}$).

Our main target parallel platform is the Linux Cluster *seth* which consists of 120 dual AMD Athlon MP2000+ nodes (1.667GHz, 384kB on-chip cache). The cluster is connected with a Wolfkit3 SCI high speed interconnect having a peak bandwidth of 667 MB/sec. The network connects the nodes in a 3-dimensional torus organized as a $6 \times 4 \times 5$ grid, where each link is “one-way” directed. In total, the system has a theoretical peak performance of 800 Gflops/sec. Moreover, *seth* is a platform which really puts any parallel algorithm to a tough test regarding its utilization of the narrow memory hierarchy of the dual nodes.

All subroutines and programs were compiled using the Portland Group’s pgf90 6.0–5 compiler using the recommended compiler flags `-O2 -tp athlonxp -fast` and the following software

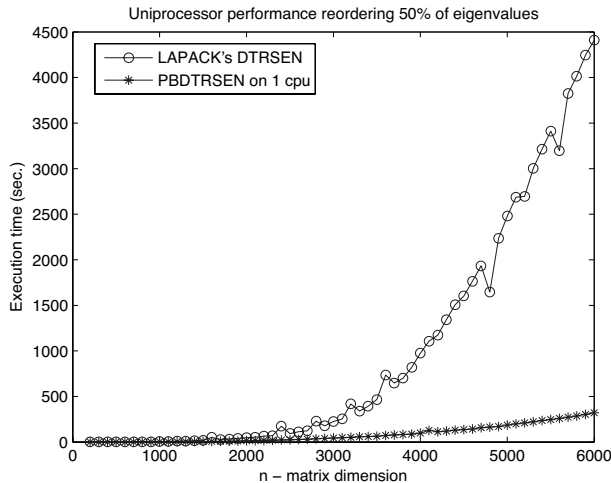


Figure 8: Uniprocessor performance results for standard LAPACK algorithm DTRSEN and parallel block algorithm PBDTRSEN reordering 50% of uniformly distributed eigenvalues of random matrices on *seth*.

libraries: ScaMPI (MPICH 1.2), LAPACK 3.0, ATLAS 3.5.9, ScaLAPACK / PBLAS 1.7.0, BLACS 1.1, RECSY 0.01alpha and SLICOT 4.0.

Below, we also utilize a second target platform, *sarek*, which is a 64-bit Opteron Linux Cluster with 192 dual AMD Opteron nodes (2.2 GHz, 1024kB unified L2 cache), 8GB RAM per node and a Myrinet-2000 high-performance interconnect with 250 MB/sec bandwidth. More information about the hardware and software related to *sarek* is available on-line [35].

All presented timings in this section are in seconds. The parallel speedup S_p is computed as

$$S_p = T_{p_{\min}}/T_p, \quad (12)$$

where T_p is the parallel execution time using p processors and p_{\min} is the smallest number of processors utilized for the current problem size.

Some initial tests were performed on *seth* by reordering 50% of the eigenvalues which are initially uniformly distributed over the diagonal of a 1500×1500 random matrix² already in Schur form. We also updated the corresponding Schur vectors in Q during the reordering. The purpose of these tests was to find close to optimal configuration parameters for the parallel algorithm executed with one computational window on one processor. By testing all feasible combinations of n_b , n_{win} and $n_{\text{eig}} = n_{\text{ceig}}$ within the integer search space $\{10, 20, \dots, 200\}$ for $r_{\text{mmult}} \in \{5, 10, \dots, 100\}$ and $n_{\text{slab}} = 32$, we found that $n_b = 180$, $n_{\text{win}} = 60$ and $n_{\text{eig}} = n_{\text{ceig}} = 30$ and $r_{\text{mmult}} = 40$ are optimal with runtime 5.66 seconds. This is a bit slower than the result in [31] (4.74 seconds) but still a great improvement over the current LAPACK algorithm (DTRSEN) which takes over 20 seconds! The difference may be partly explained by the additional data copying that is performed in the parallel algorithm during the cross border phase.

We display uniprocessor performance results for DTRSEN and PBDTRSEN in Figure 8, using the optimal configuration parameters listed above. The speedup is remarkable for large problems; $n = 5700$ shows a speedup of 14.0 for the parallel block algorithm executed on one processor compared to the standard algorithm, mainly caused by an improved memory reference pattern due to the rich usage of high performance level 3 BLAS.

Next, using the close to optimal parameter setting, we performed experiments on a 4×4 processor mesh using 4 computational windows on 6000×6000 matrices of the same form as above (which puts the same memory load per processor as a single 1500×1500 matrix on one processor)

²The strictly upper part of T is a random matrix, but we construct T such that 50% of its eigenvalues are in complex conjugate pairs.

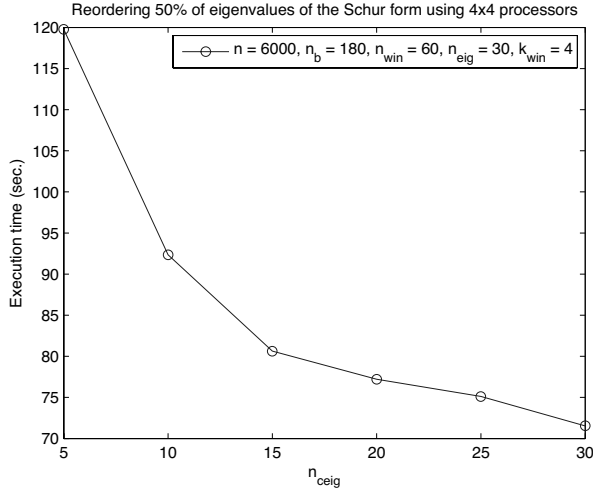


Figure 9: Parallel execution time for 6000×6000 matrices using 4 concurrent computational windows on a 4×4 processor mesh reordering 50% of uniformly distributed eigenvalues of random matrices on *seth*, with respect to n_{ceilg} – number of eigenvalues that cross the border in a single step.

to find a close to optimal value of n_{ceilg} in the integer search space $\{5, 10, \dots, 30\}$. In this case, we ended up with $n_{\text{ceilg}}^{\text{opt}} = 30$ with a corresponding execution time of 71.53 seconds, see Figure 9. It turned out that the execution time increased with a decreasing n_{ceilg} . For instance, the same problem took 119.77 seconds using $n_{\text{ceilg}} = 5$. With our analysis in mind, this is not surprising since using $n_{\text{eig}} = n_{\text{ceilg}}$ brings each subcluster of eigenvalues across the process border in one single step, thereby minimizing the number of messages sent across the border for each subcluster, i.e., the influence on the parallel runtime caused by the node latency is minimized.

In Figure 10 we display representative performance results reordering the eigenvalues of 6000×6000 matrices for different variations of the number of selected eigenvalues, their distribution over the main diagonal of the Schur form, and the number of utilized processors. The parallel speedup goes up to 16.6 for using 64 processors when selecting 5% of the eigenvalues from the lower part of T .

To confirm the validity of the derived performance model (see Section 4), we compare the experimental parallel speedup with the parallel speedup computed from combining Equations (6), (7) and (8), see Figure 11. To check the reliability of our results, we repeated each run at least 5 times but observed only insignificant variation in the obtained execution times. The machine specific constants t_a , t_s and t_w are estimated as follows:

- Since the *exact* number of flops needed for computing a reordered Schur form is not known a priori, we approximate t_a for $p = 1$ by $\tilde{t}_a^{(1)}(n, k)$ which is computed from Equation (7) by replacing T_a by T_1 , the serial runtime of the parallel algorithm for one processor only. For $p > 1$, we model t_a by

$$\tilde{t}_a^{(p)}(n, k) = \alpha_0 + \alpha_1 n + \alpha_2 k + \alpha_3 \sqrt{p}, \quad (13)$$

where $\alpha_i \in \mathbb{R}$, $i = 0, \dots, 3$. The model captures that the processor speed is expected to be proportional to the matrix data load at each processor. Since the amount of matrix data per processor and the number of selected eigenvalues per on-diagonal processor are kept fixed going from a problem instance (n, k, p) to $(2n, 2k, 4p)$, we assume

$$\tilde{t}_a^{(p)}(n, k) = \tilde{t}_a^{(4p)}(2n, 2k). \quad (14)$$

From this assumption and the available data for $\tilde{t}_a^{(1)}$ derived from Table 2, we compute the α_i -parameters from a least squares fit, see Table 1. With these parameters and fixed

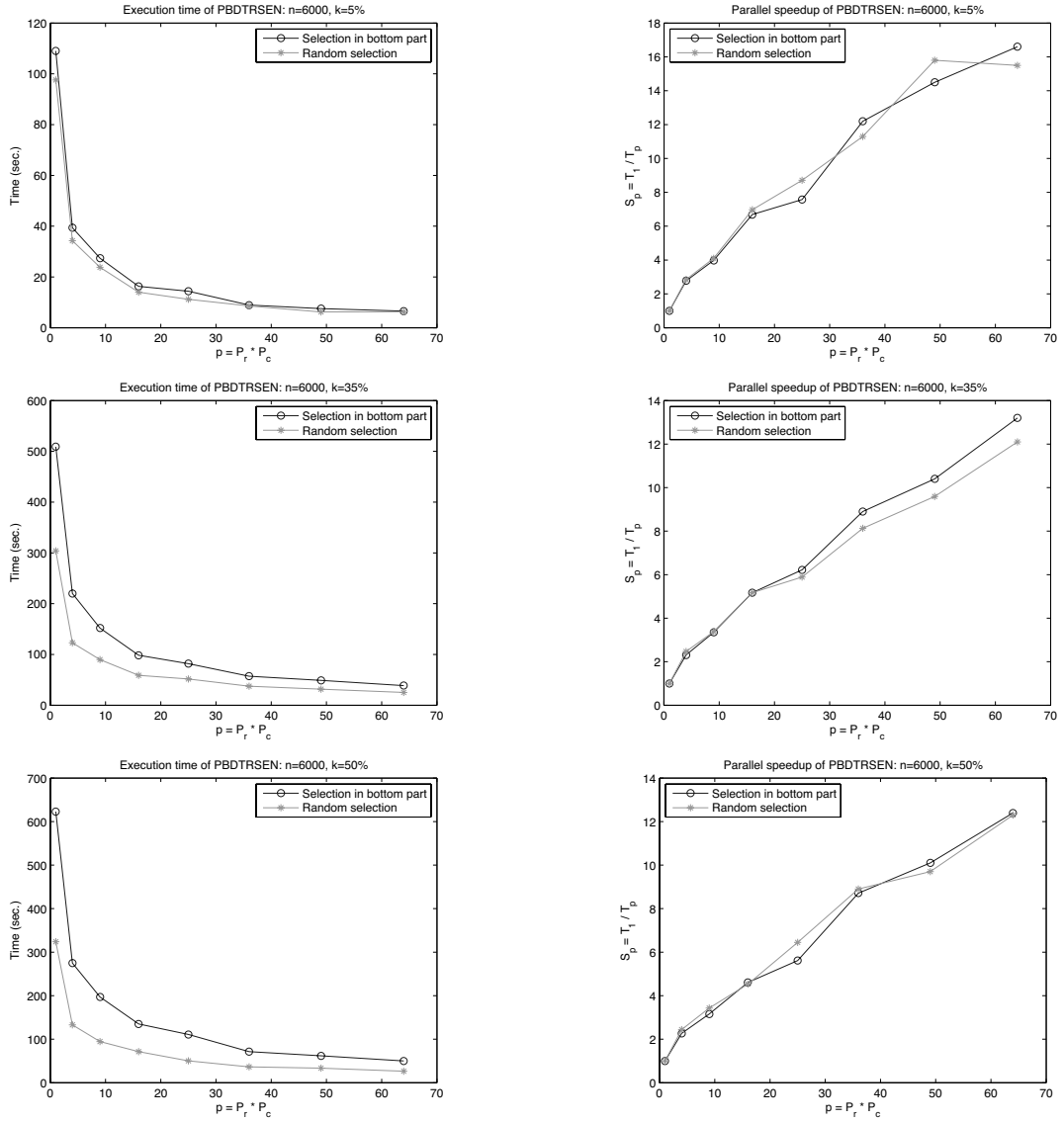


Figure 10: Experimental results for $n = 6000$ from Table 2 on *seth*. Displayed representative performance results are for different variations of the number of selected eigenvalues, their distribution over the diagonal of the Schur form and the number of utilized processors.

Table 1: Experimentally determined machine parameters for *seth*.

Parameter	Value
α_0	6.73×10^{-9}
α_1	6.34×10^{-13}
α_2	-3.10×10^{-12}
α_3	3.08×10^{-10}
t_s	3.7×10^{-6}
t_w	$\gamma \cdot 1.1 \times 10^{-8}$

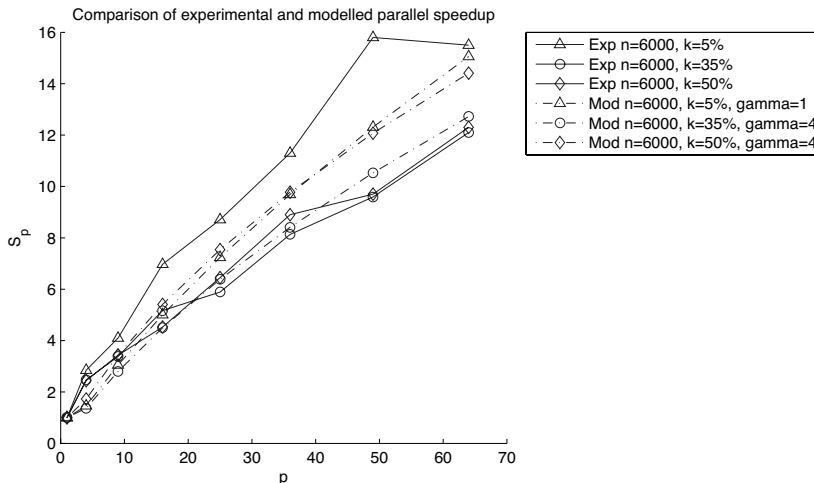


Figure 11: Comparison of modelled and experimental parallel speedup for problems with *random* distribution of the selected eigenvalues.

values on p and/or n , the model predicts that $\tilde{t}_a^{(p)}(n, k)$ decreases for an increasing value of k . Moreover, for a fixed value on k , $\tilde{t}_a^{(p)}(n, k)$ increases with p and/or n . For such cases, the decrease of the modelled processor performance is marginal, except for large values on p and/or n . With an increasing size of the processor mesh $\sqrt{p} \times \sqrt{p}$ and fixed values on n and k , it is expected that individual processors perform less efficient due to less local arithmetic work.

- The interconnection network parameters t_s and t_w are estimated by performing MPI-based *ping-pong* communication in the network combined with a linear least-squares regression, see Table 1. In Figure 11, we allow t_w to get multiplied by $\gamma \geq 1$ which models unavoidable network sharing, overhead from the BLACS layer and the potentially (much) lower memory bandwidth inside the dual nodes. For $\gamma = 1$, t_w represents the inverse of the practical peak bandwidth for double precision data in the network of *seth*.

The comparison presented in Figure 11 shows that the performance model is not perfect, but is able to predict a parallel speedup in the same range as the actually observed results.

We present experimental profiling data generated by the performance evaluation tool OPT [34] in Figures 12-13 for four concurrent processors reordering 5% of the eigenvalues of a 6000×6000 matrix on *sarek*. The profiling data reveals that the load balance is even and that the computations dominate the execution time. Moreover, the two dominating MPI routines invoked are related to point-to-point receive operations (`MPI_Recv`), which is expected to be due to the cross border data exchanges, and synchronizing barriers (`MPI_Barrier`), which is expected to be due to the synchronization of all processors right before and after each cross border phase (see Section 5.4).

We conclude this section by remarking that the numerical accuracy of the parallel reordering

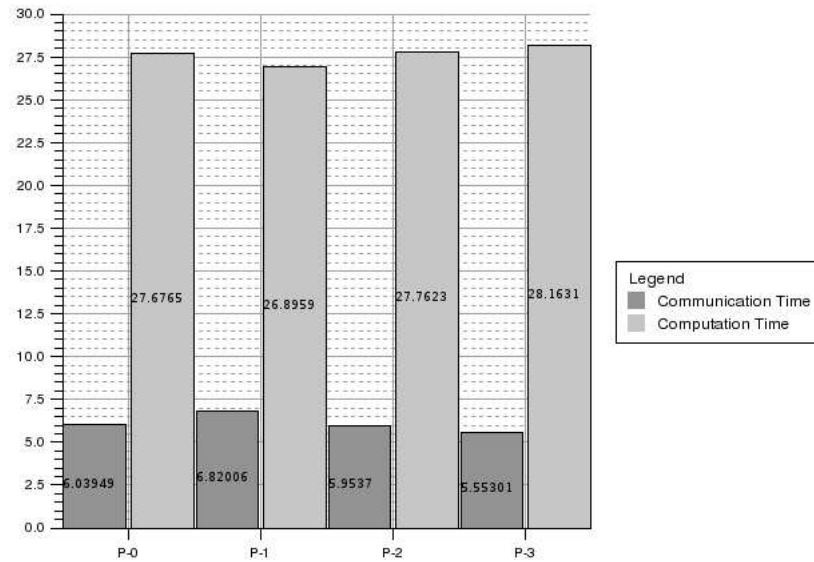


Figure 12: Execution time profile for four processors reordering 5% of the eigenvalues of a 6000×6000 matrix: total cost of computations vs communications.

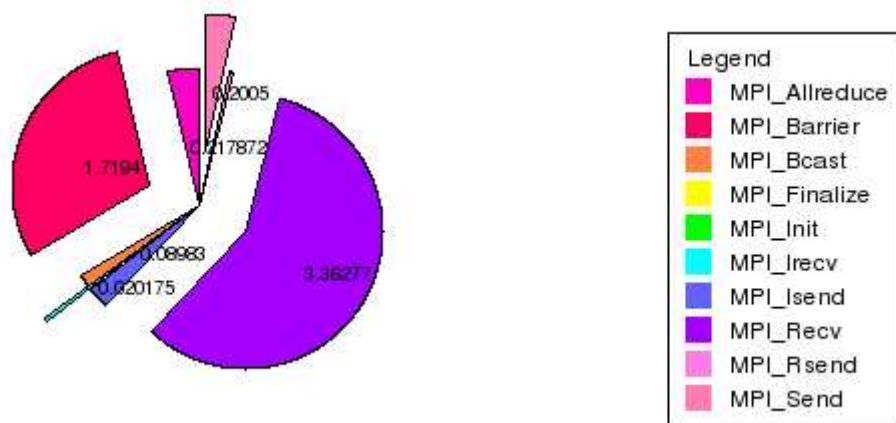


Figure 13: Execution time profile for four processors reordering 5% of the eigenvalues of a 6000×6000 matrix: invoked BLACS-MPI routines.

algorithm is similar to the results presented in [2, 30, 31] since the numerical backward stability of the algorithm is preserved when accumulating orthogonal transformations [23]. Experimental results confirming this claim are also presented in [19].

For a real-world application requiring the reordering of a significant number of eigenvalues for large matrices, see [18], where the parallel method described in this paper is applied to linear-quadratic optimal control problems.

7 Extension to the generalized real Schur form

We have extended the presented parallel algorithm for reordering the standard Schur form to eigenvalue reordering in the generalized Schur form

$$(S, T) = Q^T(A, B)Z, \quad (15)$$

where (A, B) is a regular matrix pair, Q and Z are orthogonal matrices and (S, T) is the generalized real Schur form (see, e.g., [14]). The parallel block algorithm now works on pairs of matrices. Moreover, the individual orthogonal transformation matrices from each eigenvalue swap are stored slightly different compared to the standard case (see [31] for details). For our prototype Fortran 77 implementation PBDTGSEN the following close to optimal parameters was found by extensive tests: $n_b = 180$, $n_{\text{win}} = 60$ and $n_{\text{eig}} = n_{\text{ceig}} = 30$ and $r_{\text{mmult}} = 10$ which resulted in a uniprocessor runtime 18.94 seconds, which is less than 1 second slower than the results in [31], but much faster than LAPACK.

As for the standard case, PBDTGSEN optionally computes condition numbers for the selected cluster of eigenvalues and the corresponding deflating subspaces (see, e.g., [27]) by invoking the generalized coupled Sylvester equation solvers and condition estimators from the SCASY software [16, 17, 36].

8 Summary and future work

The lack of a fast and reliable parallel reordering routine has turned the attention away from Schur-based subspace methods for solving a wide range of numerical problems in distributed memory (DM) environments. By the introduction of the algorithms presented in this paper, this situation might be subject to change.

We remark that ScaLAPACK still lacks a highly efficient implementation of the parallel QZ algorithm and the existing parallel QR algorithm is far from level 3 node performance. Moreover, the techniques developed in this paper can be used to implement parallel versions of the advanced deflation techniques described in [9, 25].

Acknowledgements

The research was conducted using the resources of the High Performance Computing Center North (HPC2N).

References

- [1] E. Anderson, Z. Bai, C. Bischof, S. Blackford, J. W. Demmel, J. J. Dongarra, J. Du Croz, A. Greenbaum, S. Hammarling, A. McKenney, and D. C. Sorensen. *LAPACK Users' Guide*. SIAM, Philadelphia, PA, third edition, 1999.
- [2] Z. Bai and J. W. Demmel. On swapping diagonal blocks in real Schur form. *Linear Algebra Appl.*, 186:73–95, 1993.

- [3] Z. Bai, J. W. Demmel, J. J. Dongarra, A. Ruhe, and H. van der Vorst, editors. *Templates for the Solution of Algebraic Eigenvalue Problems*. Software, Environments, and Tools. SIAM, Philadelphia, PA, 2000.
- [4] Z. Bai, J. W. Demmel, and A. McKenney. On computing condition numbers for the nonsymmetric eigenproblem. *ACM Trans. Math. Software*, 19(2):202–223, 1993.
- [5] M. W. Berry, J. J. Dongarra, and Y. Kim. A parallel algorithm for the reduction of a nonsymmetric matrix to block upper-Hessenberg form. *Parallel Comput.*, 21(8):1189–1211, 1995.
- [6] L. S. Blackford, J. Choi, A. Cleary, E. D’Azevedo, J. W. Demmel, I. Dhillon, J. J. Dongarra, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, and R. C. Whaley. *ScaLAPACK Users’ Guide*. SIAM, Philadelphia, PA, 1997.
- [7] BLAS - Basic Linear Algebra Subprograms. See <http://www.netlib.org/blas/index.html>.
- [8] K. Braman, R. Byers, and R. Mathias. The multishift QR algorithm, I: Maintaining well-focused shifts and level 3 performance. *SIAM J. Matrix Anal. Appl.*, 23(4):929–947, 2002.
- [9] K. Braman, R. Byers, and R. Mathias. The multishift QR algorithm, II: Aggressive early deflation. *SIAM J. Matrix Anal. Appl.*, 23(4):948–973, 2002.
- [10] J. Choi, J. J. Dongarra, and D. W. Walker. The design of a parallel dense linear algebra software library: reduction to Hessenberg, tridiagonal, and bidiagonal form. *Numer. Algorithms*, 10(3-4):379–399, 1995.
- [11] K. Dackland and B. Kågström. Blocked algorithms and software for reduction of a regular matrix pair to generalized Schur form. *ACM Trans. Math. Software*, 25(4):425–454, 1999.
- [12] J. J. Dongarra, J. Du Croz, I. S. Duff, and S. Hammarling. A set of level 3 basic linear algebra subprograms. *ACM Trans. Math. Software*, 16:1–17, 1990.
- [13] J. J. Dongarra, S. Hammarling, and J. H. Wilkinson. Numerical considerations in computing invariant subspaces. *SIAM J. Matrix Anal. Appl.*, 13(1):145–161, 1992.
- [14] G. H. Golub and C. F. Van Loan. *Matrix Computations*. Johns Hopkins University Press, Baltimore, MD, third edition, 1996.
- [15] A. Grama, A. Gupta, G. Karypsis, and V. Kumar. *Introduction to Parallel Computing, Second Edition*. Addison-Wesley, 2003.
- [16] R. Granat and B. Kågström. Parallel Solvers for Sylvester-type Matrix Equations with Applications in Condition Estimation, Part I: Theory and Algorithms. Technical Report UMINF-07.15, Department of Computing Science, Umeå University, SE-90187, UMEÅ, Sweden. Submitted to *ACM Transactions on Mathematical Software*, 2007.
- [17] R. Granat and B. Kågström. Parallel Solvers for Sylvester-type Matrix Equations with Applications in Condition Estimation, Part II: the SCASY Software. Technical Report UMINF-07.16, Department of Computing Science, Umeå University, SE-90187, UMEÅ, Sweden. Submitted to *ACM Transactions on Mathematical Software*, 2007.
- [18] R. Granat, B. Kågström, and D. Kressner. A parallel Schur method for solving continuous-time algebraic Riccati equations. Submitted to IEEE International Symposium on Computer-Aided Control Systems Design, San Antonio, TX, 2008.
- [19] R. Granat, B. Kågström, and D. Kressner. Parallel eigenvalue reordering in real Schur forms. LAPACK working note 192. See <http://www.netlib.org/lapack/lawns/downloads/>.
- [20] W.W. Hager. Condition estimates. *SIAM J. Sci. Statist. Comput.*, (3):311–316, 1984.

- [21] G. Henry, D. S. Watkins, and J. J. Dongarra. A parallel implementation of the nonsymmetric QR algorithm for distributed memory architectures. *SIAM J. Sci. Comput.*, 24(1):284–311, 2002.
- [22] N. J. Higham. Fortran codes for estimating the one-norm of a real or complex matrix, with applications to condition estimation. *ACM Trans. of Math. Software*, 14(4):381–396, 1988.
- [23] N. J. Higham. Accuracy and Stability of Numerical Algorithms. Second edition. *SIAM*, Philadelphia, PA, 2002.
- [24] B. Kågström. A direct method for reordering eigenvalues in the generalized real Schur form of a regular matrix pair (A, B) . In *Linear algebra for large scale and real-time applications (Leuven, 1992)*, volume 232 of *NATO Adv. Sci. Inst. Ser. E Appl. Sci.*, pages 195–218. Kluwer Acad. Publ., Dordrecht, 1993.
- [25] B. Kågström and D. Kressner. Multishift Variants of the QZ Algorithm with Aggressive Early Deflation. *SIAM Journal on Matrix Analysis and Applications* 29, 1, 199–227, 2006.
- [26] B. Kågström and P. Poromaa. Distributed and shared memory block algorithms for the triangular Sylvester equation with sep^{-1} estimators. *SIAM J. Matrix Anal. Appl.*, 13(1):90–101, 1992.
- [27] B. Kågström and P. Poromaa. LAPACK-style algorithms and software for solving the generalized Sylvester equation and estimating the separation between regular matrix pairs. *ACM Trans. Math. Software* 22, 1, 78–103, 1996.
- [28] B. Kågström and P. Poromaa. Computing eigenspaces with specified eigenvalues of a regular matrix pair (A, B) and condition estimation: theory, algorithms and software. *Numer. Algorithms*, 12(3-4):369–407, 1996.
- [29] J.G. Korvink and E.B. Rudnyi. Oberwolfach benchmark collection. In P. Benner, V. Mehrmann, and D.C. Sorensen, editors, *Dimension Reduction of Large-Scale Systems*, volume 45 of *Lecture Notes in Computational Science and Engineering*, pages 311–315. Springer-Verlag, Berlin/Heidelberg, Germany, 2005.
- [30] D. Kressner. *Numerical Methods and Software for General and Structured Eigenvalue Problems*. PhD thesis, TU Berlin, Institut für Mathematik, Berlin, Germany, 2004.
- [31] Daniel Kressner. Block algorithms for reordering standard and generalized Schur forms. *ACM Transactions on Mathematical Software*, 32(4):521–532, December 2006.
- [32] LAPACK - Linear Algebra Package. See <http://www.netlib.org/lapack/>.
- [33] R. Lehoucq and J. Scott. An evaluation of software for computing eigenvalues of sparse nonsymmetric matrices. Tech. Report MCS-P547-1195, Argonne National Laboratory, 1996.
- [34] OPT - The optimization and profiling tool. See <http://www.allinea.com/>.
- [35] Sarek - HPC2N Opteron Cluster. See <http://www.hpc2n.umu.se/resources/sarek.html>.
- [36] SCASY - ScaLAPACK-style solvers for Sylvester-type matrix equations. See <http://www.cs.umu.se/~granat/scasy.html>.
- [37] V. Sima. *Algorithms for Linear-Quadratic Optimization*, volume 200 of *Pure and Applied Mathematics*. Marcel Dekker, Inc., New York, NY, 1996.
- [38] ScaLAPACK Users' Guide. See <http://www.netlib.org/scalapack/slug/>.
- [39] Stewart, G. W. and Sun, J.-G. *Matrix Perturbation Theory*. Academic Press, New York, 1990.

- [40] P. Van Dooren. A generalized eigenvalue approach for solving Riccati equations. *SIAM J. Sci. Statist. Comput.*, 2(2):121–135, June, 1981.

APPENDIX A: Performance Results

In Table 2, we present representative performance results for **PBDTRSEN** for different variations of the number of selected eigenvalues, their distribution over the diagonal of the Schur form and the number of utilized processors. All presented results are obtained executing the parallel reordering algorithm using the close to optimal parameters settings for as many computational windows as possible ($k_{\text{win}} = \sqrt{p}$). For **PBDTRSEN**, the parallel speedup goes up to 16.6 for $n = 6000$ using a 8×8 processor mesh when selecting 5% of the eigenvalues from the lower part of T ; in Table 2, *Random* and *Bottom* refer to the parts of T where the selected eigenvalues reside before the reordering starts.

Some remarks regarding the results in Table 2 are in order.

- In general, the parallel algorithm scales better (delivers a higher parallel speedup) for a smaller value of k , the number of chosen eigenvalues for a given problem size. The main reason is that less computational work per processor in general leads to less efficient usage of processor resources, i.e., a lower Mflops-rate, which in turn makes the communication overhead smaller in relation to arithmetics. We also see improved processor utilization when increasing the number of selected eigenvalues to reorder from 5% to 50%. The amount of work is increased by a factor 10 but the uniprocessor execution times increase roughly by a factor of 6 ($n = 3000$).
- Since only the processor rows and columns corresponding to the selected groups of four adjacent processors can be efficiently utilized during the cross border phase (see the formation of the processor groups in Algorithm 2), the performance gain going from a 2×2 to a 3×3 processor mesh is sometimes bad (even negative, see the problem $n = 1500$ and 5% selected eigenvalues). In principle, some performance degradation will occur when going from $2q \times 2q$ to $(2q + 1) \times (2q + 1)$ processors, for $q \geq 1$, since the level of concurrency in the cross border phase will not increase. This effect is not that visible for larger meshes since the relative amount of possibly idle processors in the cross border part decreases with an increasing processor mesh.

We repeated the experiments from the standard Schur form for the generalized Schur form with similar experimental results, see Table 3. For $n = 6000$, the memory on one 2GB node of *seth* available for users is not large enough to hold all data objects (signaled by '-'). In this case, S_p is computed using $p_{\text{min}} = 4$ in Equation (12), i.e., the presented values exemplify the speedup going from $p_0 = 4$ to $p > p_0$ processor. By assuming a fictive parallel speedup of at least 2 going from 1×1 to 2×2 processors for $n = 6000$ (see the results in Table 2), we conclude that the scalability for the generalized case is as good as or even better than the standard case.

Table 2: Performance of PBDTRSEN on *seth*.

Sel.	$P_r \times P_c$	n	Random		Bottom		n	Random		Bottom	
			time	S_p	time	S_p		time	S_p	time	S_p
5%	1×1	1500	1.98	1.00	2.16	1.00	4500	42.7	1.00	42.4	1.00
	2×2	1500	0.90	2.20	1.02	2.12	4500	15.9	2.69	17.6	2.41
	3×3	1500	1.01	1.96	1.04	2.08	4500	10.8	3.95	12.4	3.42
	4×4	1500	0.80	2.48	0.92	2.35	4500	6.65	6.42	7.74	5.48
	5×5	1500	0.34	5.82	0.76	2.84	4500	4.70	9.09	6.25	6.78
	6×6	1500	0.32	6.19	0.35	6.17	4500	4.26	10.0	4.51	9.40
	7×7	1500	0.32	6.19	0.32	6.75	4500	3.48	12.3	3.50	12.1
	8×8	1500	0.24	8.25	0.27	8.00	4500	3.13	13.6	3.50	12.1
35%	1×1	1500	5.78	1.00	9.87	1.00	4500	127	1.00	217	1.00
	2×2	1500	2.96	1.95	5.58	1.77	4500	55.3	2.30	96.7	2.24
	3×3	1500	2.05	2.82	4.03	2.45	4500	40.4	3.14	67.9	3.20
	4×4	1500	1.74	3.32	2.72	3.63	4500	25.5	4.98	43.2	5.02
	5×5	1500	1.51	3.83	2.28	4.33	4500	22.4	5.67	39.0	5.56
	6×6	1500	1.17	4.94	1.82	5.42	4500	16.7	7.60	26.4	8.22
	7×7	1500	1.13	5.12	1.49	6.62	4500	15.8	8.04	21.6	10.0
	8×8	1500	1.02	5.67	1.12	8.81	4500	11.8	10.8	18.7	11.6
50%	1×1	1500	5.66	1.00	12.6	1.00	4500	140	1.00	263	1.00
	2×2	1500	3.50	1.62	6.83	1.84	4500	62.7	2.23	121	2.17
	3×3	1500	1.89	2.99	4.97	2.53	4500	36.1	3.88	84.3	3.12
	4×4	1500	1.87	3.03	4.48	2.81	4500	29.2	4.79	61.5	4.28
	5×5	1500	1.54	3.68	3.20	3.93	4500	21.6	6.48	49.6	5.30
	6×6	1500	1.22	4.64	2.03	6.21	4500	18.8	7.45	34.2	7.69
	7×7	1500	1.30	4.35	1.90	6.63	4500	16.5	8.48	30.3	8.68
	8×8	1500	1.11	5.10	1.76	7.16	4500	12.2	11.5	24.2	10.9
5%	1×1	3000	12.8	1.00	19.8	1.00	6000	97.6	1.00	109	1.00
	2×2	3000	4.93	2.60	6.63	2.99	6000	34.3	2.85	39.4	2.77
	3×3	3000	3.69	3.47	4.00	4.95	6000	23.8	4.10	27.4	3.98
	4×4	3000	3.01	4.26	2.90	6.83	6000	14.0	6.97	16.3	6.69
	5×5	3000	2.97	4.31	2.54	7.80	6000	11.2	8.71	14.4	7.57
	6×6	3000	1.72	7.44	2.12	9.34	6000	8.60	11.3	8.94	12.2
	7×7	3000	1.82	7.03	2.33	8.50	6000	6.20	15.8	7.52	14.5
	8×8	3000	1.45	8.83	1.59	12.5	6000	6.30	15.5	6.58	16.6
35%	1×1	3000	40.1	1.00	72.4	1.00	6000	304	1.00	509	1.00
	2×2	3000	19.0	2.11	38.0	1.91	6000	123	2.47	220	2.31
	3×3	3000	13.8	2.91	23.2	3.12	6000	89.9	3.38	152	3.35
	4×4	3000	8.83	4.54	14.2	5.10	6000	58.8	5.17	98.4	5.17
	5×5	3000	7.81	5.13	14.8	4.89	6000	51.6	5.89	81.8	6.22
	6×6	3000	6.90	5.81	9.26	7.82	6000	37.4	8.13	57.2	8.90
	7×7	3000	5.81	6.90	9.07	7.98	6000	31.7	9.59	49.0	10.4
	8×8	3000	5.30	7.57	7.17	10.1	6000	25.2	12.1	38.7	13.2
50%	1×1	3000	50.2	1.00	91.4	1.00	6000	324	1.00	623	1.00
	2×2	3000	24.2	2.07	44.5	2.05	6000	133	2.44	275	2.27
	3×3	3000	18.9	2.66	28.1	3.25	6000	94.5	3.43	197	3.16
	4×4	3000	11.6	4.33	18.2	5.02	6000	71.5	4.53	135	4.61
	5×5	3000	9.57	5.25	17.2	5.31	6000	50.2	6.45	111	5.61
	6×6	3000	7.17	7.00	12.6	7.25	6000	36.4	8.90	71.5	8.71
	7×7	3000	6.89	7.29	11.4	8.02	6000	33.4	9.70	61.8	10.1
	8×8	3000	5.27	9.53	9.81	9.32	6000	26.4	12.3	50.1	12.4

Table 3: Performance of PBDTGEN on *seth*.

Sel.	$P_r \times P_c$	n	Random		Bottom		n	Random		Bottom	
			time	S_p	time	S_p		time	S_p	time	S_p
5%	1 × 1	1500	5.12	1.00	6.56	1.00	4500	98.4	1.00	105	1.00
	2 × 2	1500	2.65	1.93	2.91	2.25	4500	39.9	2.47	48.0	2.19
	3 × 3	1500	2.66	1.92	2.74	2.39	4500	23.4	4.21	32.4	3.24
	4 × 4	1500	2.20	2.32	2.50	2.62	4500	17.3	5.69	18.6	5.65
	5 × 5	1500	1.48	3.46	1.97	3.33	4500	15.2	6.47	17.0	6.18
	6 × 6	1500	1.05	4.88	1.69	3.88	4500	10.4	9.46	11.8	8.90
	7 × 7	1500	0.77	6.65	1.25	5.25	4500	8.93	11.0	9.22	11.4
	8 × 8	1500	0.60	8.53	0.67	9.79	4500	7.33	13.4	8.15	12.9
35%	1 × 1	1500	13.6	1.00	22.0	1.00	4500	279	1.00	312	1.00
	2 × 2	1500	8.02	1.70	11.4	1.93	4500	133	2.10	218	1.43
	3 × 3	1500	7.92	1.72	7.63	2.88	4500	77.2	3.61	122	2.56
	4 × 4	1500	3.52	3.86	5.41	4.07	4500	50.4	5.54	73.7	4.23
	5 × 5	1500	3.21	4.24	4.50	4.89	4500	40.4	6.91	59.5	5.24
	6 × 6	1500	3.71	3.67	4.85	4.54	4500	32.2	8.66	43.7	7.14
	7 × 7	1500	2.21	6.16	3.33	6.61	4500	26.1	10.7	40.3	7.74
	8 × 8	1500	2.15	6.33	2.84	7.75	4500	20.8	13.4	33.3	9.37
50%	1 × 1	1500	18.9	1.00	26.9	1.00	4500	301	1.00	584	1.00
	2 × 2	1500	13.8	1.37	14.5	1.89	4500	146	2.06	288	2.03
	3 × 3	1500	11.1	1.70	10.1	2.66	4500	82.9	3.63	169	3.46
	4 × 4	1500	4.88	3.87	7.94	3.39	4500	54.9	5.48	104	5.62
	5 × 5	1500	4.01	4.71	6.57	4.10	4500	41.7	7.22	91.9	6.35
	6 × 6	1500	3.18	5.94	4.19	6.42	4500	32.8	9.18	64.9	9.00
	7 × 7	1500	3.35	5.64	4.12	6.53	4500	32.1	9.37	62.9	9.28
	8 × 8	1500	3.01	6.23	3.96	6.79	4500	25.4	11.9	45.1	12.95
5%	1 × 1	3000	34.3	1.00	35.9	1.00	6000	-	-	-	-
	2 × 2	3000	13.5	2.54	15.3	2.35	6000	92.4	1.00	110	1.00
	3 × 3	3000	8.55	4.02	8.90	4.03	6000	51.1	1.81	58.3	1.89
	4 × 4	3000	6.13	5.60	5.79	6.20	6000	34.9	2.65	39.2	2.81
	5 × 5	3000	6.20	5.53	4.32	8.31	6000	29.2	3.16	33.0	3.33
	6 × 6	3000	4.59	7.47	3.50	10.3	6000	18.8	4.91	20.6	5.34
	7 × 7	3000	4.14	8.29	3.40	10.6	6000	16.0	5.78	19.6	5.61
	8 × 8	3000	2.69	12.8	2.69	13.3	6000	13.3	6.95	14.5	7.59
35%	1 × 1	3000	93.4	1.00	155	1.00	6000	-	-	-	-
	2 × 2	3000	44.6	2.09	74.4	2.08	6000	326	1.00	523	1.00
	3 × 3	3000	27.1	3.45	42.0	3.69	6000	180	1.81	280	1.87
	4 × 4	3000	19.2	4.86	25.4	6.10	6000	104	3.13	163	3.21
	5 × 5	3000	16.7	5.59	22.3	6.95	6000	88.1	3.70	135	3.87
	6 × 6	3000	12.4	7.53	18.1	8.56	6000	64.5	5.05	85.4	6.12
	7 × 7	3000	10.9	8.57	15.2	10.2	6000	50.7	6.43	80.9	6.46
	8 × 8	3000	9.80	9.53	11.8	13.1	6000	43.3	7.53	61.2	8.54
50%	1 × 1	3000	105	1.00	194	1.00	6000	-	-	-	-
	2 × 2	3000	53.9	1.95	98.6	1.97	6000	335	1.00	695	1.00
	3 × 3	3000	31.9	3.29	64.5	3.01	6000	187	1.79	391	1.78
	4 × 4	3000	21.5	4.88	35.8	5.42	6000	112	2.99	234	2.97
	5 × 5	3000	18.5	5.68	31.6	6.14	6000	87.5	3.83	181	3.83
	6 × 6	3000	14.1	7.45	24.3	7.98	6000	66.6	5.03	127	5.47
	7 × 7	3000	12.6	8.33	22.6	8.58	6000	62.0	5.40	114	6.10
	8 × 8	3000	9.81	10.7	20.3	9.56	6000	45.5	7.36	88.0	7.90