

Computer Science and Mathematics Division

Mathematical Sciences Section

**THE DESIGN AND IMPLEMENTATION OF THE  
PARALLEL OUT-OF-CORE SCALAPACK LU, QR AND CHOLESKY FACTORIZATION  
ROUTINES**

Ed F. D'Azevedo

Jack J. Dongarra

Oak Ridge National Laboratory

P.O. Box 2008, Bldg. 6012

Oak Ridge, TN 37831-6367

Date Published: April 1997

This work was supported in part by the National Science Foundation Grant No. ASC-9005933; by the Defense Advanced Research Projects Agency under contract DAAL03-91-C-0047, administered by the Army Research Office; by the Office of Scientific Computing, U.S. Department of Energy, under Contract DE-AC05-84OR21400; and by the National Science Foundation Science and Technology Center Cooperative Agreement No. CCR-8809615, and Center for Computational Sciences at Oak Ridge National Laboratory for the use of the computing facilities.

Prepared by the

Oak Ridge National Laboratory

Oak Ridge, Tennessee 37831

managed by

Lockheed Martin Energy Research Corp.

for the

U.S. DEPARTMENT OF ENERGY

under Contract No. DE-AC05-96OR22464

## Contents

1	Introduction . . . . .	1
2	I/O Library . . . . .	1
	2.1 Low-level Details . . . . .	2
	2.2 User Interface . . . . .	3
3	Left-looking Algorithm . . . . .	6
	3.1 Partitioned Factorization . . . . .	6
	3.2 LU Factorization . . . . .	8
	3.3 QR Factorization . . . . .	9
	3.4 Cholesky Factorization . . . . .	10
4	Numerical Results . . . . .	11
5	Conclusions . . . . .	13
6	References . . . . .	14

## List of Tables

- |   |   |    |
|---|---|----|
| 1 | Performance of out-of-core LU factorization on 64 processors using $MB=NB=50$ .       | 12 |
| 2 | Performance of out-of-core QR factorization on 64 processors using $MB=NB=50$ .       | 12 |
| 3 | Performance of out-of-core Cholesky factorization on 64 processors using $MB=NB=50$ . | 13 |

## List of Figures

1	Algorithm requires 2 in-core panels. . . . .	6
---	--	---

**THE DESIGN AND IMPLEMENTATION OF THE  
PARALLEL OUT-OF-CORE SCALAPACK LU, QR AND CHOLESKY FACTORIZATION  
ROUTINES**

Ed F. D'Azevedo  
Jack J. Dongarra

**Abstract**

This paper describes the design and implementation of three core factorization routines — LU, QR and Cholesky — included in the out-of-core extension of ScaLAPACK. These routines allow the factorization and solution of a dense system that is too large to fit entirely in physical memory. An image of the full matrix is maintained on disk and the factorization routines transfer sub-matrices into memory. The 'left-looking' column-oriented variant of the factorization algorithm is implemented to reduce the disk I/O traffic. The routines are implemented using a portable I/O interface and utilize high performance ScaLAPACK factorization routines as in-core computational kernels.

We present the details of the implementation for the out-of-core ScaLAPACK factorization routines, as well as performance and scalability results on the Intel Paragon.

## 1. Introduction

This paper describes the design and implementation of three core factorization routines — LU, QR and Cholesky — included in the out-of-core extensions of ScaLAPACK. These routines allow the factorization and solution of a dense linear system that is too large to fit entirely in physical memory.

Although current computers have unprecedented memory capacity, out-of-core solvers are still needed to tackle even larger applications. A modern workstation is commonly equipped with 64 to 128Mbytes of memory and capable of performing over 100 Mflops/sec. Even on a large problem that occupies all available memory, the in-core solution of dense linear problems typically takes less than an hour. On a network of workstations (NOW) with 100 processors, each with 64Mbytes, it may require about 30 minutes to factor and solve at 64-bit precision a dense linear system of order 30,000. This suggests that the processing power of such high performance machines is under-utilized and much larger systems can be tackled before run time becomes prohibitively large. Therefore, it is natural to develop parallel out-of-core solvers to tackle large dense linear systems. Such dense problems arise from high resolution three-dimensional electromagnetic scattering problems or in modeling fluid flow around complex objects.

The development effort has the objective of producing portable software that achieves high performance on distributed memory multiprocessors, shared memory multiprocessors, and NOW. The implementation is based on modular software building blocks such as the PBLAS (Parallel Basic Linear Algebra Subprograms), and the BLACS (Basic Linear Algebra Communication Subprograms). Proven and highly efficient ScaLAPACK factorization routines are used for in-core computations.

One key component of an out-of-core library is an efficient and portable I/O interface. We have implemented a high level I/O layer to encapsulate machine or architecture specific characteristics to achieve good throughput. The I/O layer eases the burden of manipulating out-of-core matrices by directly supporting the reading and writing of unaligned sections of ScaLAPACK block-cyclic distributed matrices.

Section 2 describes the design and implementation of the portable I/O Library. The implementation of the 'left-looking' column-oriented variant of the LU, QR and Cholesky factorization is described in §3. Finally, §4 summarizes the performance on the Intel Paragon.

## 2. I/O Library

This section describes the overall design of the I/O Library including both the high level user interface, and the low level implementation details needed to achieve good performance.

## 2.1. Low-level Details

Each out-of-core matrix is associated with a device unit number (between 1 and 99), much like the familiar Fortran I/O subsystem. Each I/O operation is record-oriented, where each record is conceptually an  $MMB \times NNB$  ScaLAPACK block-cyclic distributed matrix. Moreover if this record/matrix is distributed with  $(MB, NB)$  as the block size on a  $P \times Q$  processor grid, then  $\text{mod}(MMB, MB * P) = 0$  and  $\text{mod}(NNB, NB * Q) = 0$ , i.e.  $MMB$  (and  $NNB$ ) are exact multiples of  $MB * P$  (and  $NB * Q$ ). Data to be transferred is first copied or assembled into an internal temporary buffer (record). This arrangement reduces the number of `lseek()` system calls and encourages large contiguous block transfers, but incurs some overhead in memory-to-memory copies. All processors are involved in each record transfer. Individually, each processor writes out an  $(MMB/P)$  by  $(NNB/Q)$  matrix block.  $MMB$  and  $NNB$  can be adjusted to achieve good I/O performance with large contiguous block transfers or to match RAID disk stripe size. A drawback of this arrangement is that I/O on narrow block rows or block columns will involve only processors aligned on the same row or column of the processor grid, and thus may not obtain full bandwidth from the I/O subsystem. An optimal block size for I/O transfer may not be equally efficient for in-core computations. On the Intel Paragon,  $MB$  (or  $NB$ ) can be as small as 8 for good efficiency but requires at least 64Kbytes I/O transfers to achieve good performance to the parallel file system. A *2-dimensional cyclically-shifted block layout* that achieves good load balance even when operating on narrow block rows or block columns was proposed in MIOS (Matrix Input-Output Subroutines) used in SOLAR. However, this scheme is more complex to implement, (SOLAR does not yet use this scheme). Moreover, another data redistribution is required to maintain compatibility with in-core ScaLAPACK software. A large data redistribution would incur a large message volume and a substantial performance penalty, especially in a NOW environment.

The I/O library supports both a 'shared' and 'distributed' organization of disk layout. In a 'distributed' layout, each processor opens a unique file on its local disk (e.g. '/tmp' partition on workstations) to be associated with the matrix. This is most applicable on a NOW environment or where a parallel file system is not available. On systems where a shared parallel file system is available (such as `M_ASYNC` mode for PFS on Intel Paragon), all processors open a common shared file. Each processor can independently perform `lseek/read/write` requests to this common file. Physically, the 'shared' layout can be the concatenation of the many 'distributed' files. Another organization is to 'interlace' contributions from individual processors into each record on the shared file. This may lead to better pre-fetch caching by the operating system, but requires an `lseek()` operation by each processor, even on reading sequential records. On the Paragon, `lseek()` is an expensive operation since it generates a message to the I/O nodes. Note that most implementations of NFS (Networked File System) do not correctly support

multiple concurrent read/write requests to a shared file.

Unlike MIOS in SOLAR, only a synchronous I/O interface is provided for reasons of portability and simplicity of implementation. A fully portable (although possibly not the most efficient) implementation of the I/O layer using Fortran record-oriented I/O is also possible<sup>1</sup>. The current I/O library is written in C and uses standard POSIX I/O operations. System dependent routines, such as NX-specific `gopen()` or `esseek()` system calls, may be required to access files over 2Gbytes. Asynchronous I/O that overlaps computation and I/O is most effective only when processing time for I/O and computation are closely matched. Asynchronous I/O provides little benefits in cases where in-core computation or disk I/O dominates overall time. Asynchronous pre-fetch reads or delayed buffered writes also require dedicating scarce memory for I/O buffers. Having less memory available for the factorization may increase the number of passes over the matrix and increase overall I/O volume.

## 2.2. User Interface

To maintain ease of use and compatibility with existing ScaLAPACK software, a new ScaLAPACK array descriptor has been introduced. This out-of-core descriptor (`DTYPE_ = 601`) extends the existing descriptor for dense matrices (`DTYPE_ = 1`) to encapsulate and hide implementation-specific information such as the I/O device associated with an out-of-core matrix and the layout of the data on disk.

The in-core ScaLAPACK calls for performing a Cholesky factorization may consist of:

```
*
*   initialize descriptor for matrix A
*
*           CALL DESCINIT(DESCA,M,N,MB,NB,RSRC,CSRC,ICONTXT,LDA,INFO)
*
*   perform Cholesky factorization
*
*           CALL PDPOTRF(UPLO,N,A,IA,JA,DESCA,INFO)
```

where the array descriptor `DESCA` is an integer array of length 9 whose entries are described by the following:

---

<sup>1</sup>We are not aware of any implementation of fully portable asynchronous I/O without using threads. However, a portable thread library may not be available and greatly complicates the code.



DESC_()	Symbolic Name	Scope	Definition
1	DTYPE_A	(global)	The descriptor type DTYPE_A=1.
2	CTXT_A	(global)	The BLACS context handle, indicating the BLACS process grid over which the global matrix A is distributed. The context itself is global, but the handle (the integer value) may vary.
3	M_A	(global)	The number of rows in the global array A.
4	N_A	(global)	The number of columns in the global array A.
5	MB_A	(global)	The blocking factor used to distribute the rows of the array.
6	NB_A	(global)	The blocking factor used to distribute the columns of the array.
7	RSRC_A	(global)	The process row over which the first row of the array A is distributed.
8	CSRC_A	(global)	The process column over which the first column of the array A is distributed.
9	LLD_A	(local)	The leading dimension of the local array. $LLD\_A \geq \text{MAX}(1, \text{LOCp}(M\_A))$ .

The out-of-core version is very similar:

```

*
*   initialize extended descriptor for out-of-core matrix A
*
*           CALL PFDESCINIT(DESCA,M,N,MB,NB,RSRC,CSRC,ICONTXT,IODEV,
*           'SHARED',MMB,NNB,ASIZE, '/pfs/a.data'//CHAR(0),INFO)
*
*   perform out-of-core Cholesky factorization
*
*           CALL PFDOTRF(UPLO,N,A,IA,JA,DESCA,INFO)

```

where the array descriptor DESC\_A is an integer array of length 11 whose entries are described by the following:

DESC_()	Symbolic Name	Scope	Definition
1	DTYPE_A	(global)	The descriptor type DTYPE_A=601 for an out-of-core matrix.
2	CTXT_A	(global)	The BLACS context handle, indicating the $P \times Q$ BLACS process grid over which the global matrix A is distributed. The context itself is global, but the handle (the integer value) may vary.
3	M_A	(global)	The number of rows in the global array A.
4	N_A	(global)	The number of columns in the global array A.
5	MB_A	(global)	The blocking factor used to distribute the rows of the $MMB \times NNB$ submatrix.
6	NB_A	(global)	The blocking factor used to distribute the columns of the $MMB \times NNB$ submatrix.
7	RSRC_A	(global)	The process row over which the first row of the array A is distributed.
8	CSRC_A	(global)	The process column over which the first column of the array A is distributed.
9	LLD_A	(local)	The conceptual leading dimension of the global array. Usually this is taken to be $M_A$ .
10	IODEV_A	global	The I/O unit device number associated with the out-of-core matrix A.
11	SIZE_A	local	The amount of local in-core memory available for the factorization of A.

Here ASIZE is the amount of in-core buffer storage available in array 'A' associated with the out-of-core matrix. A 'Shared' layout is prescribed and the file '/pfs/A.data' is used on unit device IODEV. Each I/O record is an  $MMB$  by  $NNB$  ScaLAPACK block-cyclic distributed matrix.

The out-of-core matrices can also be manipulated by read/write calls. For example:

```
CALL ZLAREAD(IODEV, M,N, IA,JA, B, IB,JB, DESCB, INFO)
```

reads in an  $M$  by  $N$  sub-matrix starting at position  $(IA, JA)$  into an in-core ScaLAPACK matrix  $B( IB: IB+M-1, JB: JB+N-1)$ . Best performance is achieved with data transfer exactly aligned to local processor and block boundary; otherwise redistribution by message passing may be required for unaligned non-local data transfer to matrix B.

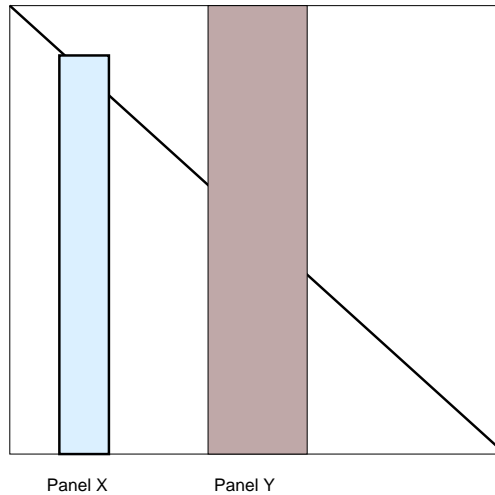


Figure 1: Algorithm requires 2 in-core panels.

### 3. Left-looking Algorithm

The three factorization algorithms, LU, QR, and Cholesky, use a similar ‘left-looking’ organization of computation. The left-looking variant is first described as a particular choice in a block-partitioned algorithm in §3.1.

The actual implementation of the left-looking factorization uses two full column in-core panels (call these X, Y; see Figure 1). Panel X is  $N_{NB}$  columns wide and panel Y occupies the remaining memory but should be at least  $N_{NB}$  columns wide. Panel X acts as a buffer to hold and apply previously computed factors to panel Y. Once all updates are performed, panel Y is factored using an in-core ScaLAPACK algorithm. The results in panel Y are then written to disk.

The following subsections describe in more detail the implementation of LU, QR and Cholesky factorization.

#### 3.1. Partitioned Factorization

The ‘left-looking’ and ‘right-looking’ variants of LU factorization can be described as particular choices in a partitioned factorization. The reader can easily generalize the following for a QR or Cholesky factorization.

Let an  $m \times n$  matrix  $A$  be factored into  $PA = LU$  where  $P$  is a permutation matrix, and  $L$  and  $U$  are the lower and upper triangular factors. We treat matrix  $A$  as a block-partitioned

matrix

$$A = \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix},$$

where  $A_{11}$  is a square  $k \times k$  sub-matrix.

1. The assumption is that the first  $k$  columns are already factored

$$P_1 \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix} = \begin{pmatrix} L_{11} \\ L_{21} \end{pmatrix} (U_{11}) \quad , \quad (1)$$

where

$$A_{11} = L_{11}U_{11}, \quad A_{21} = L_{21}U_{11}. \quad (2)$$

If  $k \leq n_0$  is small enough, a fast non-recursive algorithm such as ScaLAPACK P<sub>x</sub>GETRF may be used directly to perform the factorization; otherwise, the factors may be obtained recursively by the same algorithm.

2. Apply the permutation to the unmodified sub-matrix

$$\begin{pmatrix} \tilde{A}_{12} \\ \tilde{A}_{22} \end{pmatrix} = P_1 \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix}. \quad (3)$$

3. Compute  $U_{12}$  by solving the triangular system

$$L_{11}U_{12} = \tilde{A}_{12} \quad (4)$$

4. Perform update to  $\tilde{A}_{22}$

$$\tilde{A}_{22} \leftarrow \tilde{A}_{22} - L_{21}U_{12} \quad (5)$$

5. Recursively factor the remaining matrix

$$P_2 \tilde{A}_{22} = L_{22}U_{22} \quad (6)$$

6. Final factorization is

$$P_2 P_1 \begin{pmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{pmatrix} = \begin{pmatrix} L_{11} & \mathbf{0} \\ \tilde{L}_{21} & L_{22} \end{pmatrix} \begin{pmatrix} U_{11} & \mathbf{0} \\ U_{12} & U_{22} \end{pmatrix}, \quad \tilde{L}_{21} = P_2 L_{21}. \quad (7)$$

Note that the above is the recursively-partitioned LU factorization proposed by Toledo [4] if  $k$  is chosen to be  $n/2$ . A right-looking variant results if  $k = n_0$  is always chosen where most

of the computation is the updating of

$$\tilde{A}_{22} \leftarrow \tilde{A}_{22} - L_{21}U_{12}.$$

A left-looking variant results if  $k = n - n_0$ .

The in-core ScaLAPACK factorization routines for LU, QR and Cholesky factorization, use a right-looking variant for good load balancing [1]. Other work has shown [2, 3] that for an out-of-core factorization, a left-looking variant generates less I/O volume compared to the right-looking variant. Toledo [5] shows that the recursively-partitioned algorithm ( $k = n/2$ ) may be more efficient than the left-looking variant when a very large matrix is factored with minimal in-core storage.

### 3.2. LU Factorization

The out-of-core LU factorization PFxGETRF involves the following operations:

1. If no updates are required in factorizing the first panel, all available storage is used as one panel,
  1. LAREAD: read in part of original matrix
  2. PxGETRF: ScaLAPACK in-core factorization

$$\begin{pmatrix} L_{11} \\ L_{21} \end{pmatrix} (U_{11}) \leftarrow P_1 \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$$

3. LAWRITE: write out factors

Otherwise, partition storage into panels X and Y.

2. We compute updates into panel Y by reading in the previous factors (NNB columns at a time) into panel X. Let panel Y hold  $(A_{12}, A_{22})^t$ ,
  1. LAREAD: read in part of factor into panel X
  2. LAPIV: physically exchange rows in panel Y to match permuted ordering in panel X

$$\begin{pmatrix} \tilde{A}_{12} \\ \tilde{A}_{22} \end{pmatrix} \leftarrow P_1 \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix}$$

3. PxTRSM: triangular solve to compute upper triangular factor

$$U_{12} \leftarrow L_{11}^{-1} \tilde{A}_{12}$$

4. P<sub>x</sub>GEMM: update remaining lower part of panel Y

$$\tilde{A}_{22} \leftarrow \tilde{A}_{22} - L_{21}U_{12}$$

3. Once all previous updates are performed, we apply in-core ScaLAPACK P<sub>x</sub>GETRF to compute LU factors in panel Y

$$L_{22}U_{22} \leftarrow P_2\tilde{A}_{22} .$$

The results are then written back out to disk.

4. A final extra pass over the computed lower triangular  $L$  matrix may be required to rearrange the factors in the final permutation order

$$\tilde{L}_{12} \leftarrow P_2L_{12} .$$

Note that although PF<sub>x</sub>GETRF can accept a general rectangular matrix, a column-oriented algorithm is used. The pivot vector is held in memory and not written out to disk. During the factorization, factored panels are stored on disk with only partially or ‘incompletely’ pivoted row data, whereas factored panels were stored in original unpivoted form in [2] and repivoted ‘on-the-fly’. The current scheme is more complex to implement but reduces the number of row exchanges required.

### 3.3. QR Factorization

The out-of-core QR factorization PF<sub>x</sub>GEQRF involves the following operations:

1. If no updates are required in factorizing the first panel, all available memory is used as one panel,

1. LAREAD: read in part of original matrix

2. P<sub>x</sub>GEQRF: in-core factorization

$$Q_1 \begin{pmatrix} R_{11} \\ 0 \end{pmatrix} \leftarrow \begin{pmatrix} A_{11} \\ A_{21} \end{pmatrix}$$

3. LAWRITE: write out factors

Otherwise, partition storage into panels X and Y.

2. We compute updates into panel Y by bringing in previous factors  $NNB$  columns at a time into panel X.

1. LAREAD: read in part of factor into panel X
2. PxORMQR: apply Householder transformation to panel Y

$$\begin{pmatrix} R_{21} \\ \tilde{A}_{22} \end{pmatrix} \leftarrow Q_1^t \begin{pmatrix} A_{12} \\ A_{22} \end{pmatrix}$$

3. Once all previous updates are performed, we apply in-core ScaLAPACK PxGEQRF to compute QR factors in panel Y

$$Q_2 R_{22} \leftarrow \tilde{A}_{22}$$

The results are then written back out to disk.

Note that to be compatible with the encoding of Householder transformation in the TAU(\*) vector as used ScaLAPACK routines, a column-oriented algorithm is used even for rectangular matrices. The TAU(\*) vector is held in memory and is not written out to disk.

### 3.4. Cholesky Factorization

The out-of-core Cholesky factorization PxPOTRF factors a symmetric matrix into  $A = LL^t$  without pivoting. The algorithm involves the following operations:

1. If no updates are required in factorizing the first panel, all available memory is used as one panel,

1. LAREAD: read in part of original matrix
2. PxPOTRF: ScaLAPACK in-core factorization

$$L_{11} \leftarrow A_{11}$$

3. PxTRSM: modify remaining column

$$L_{21} \leftarrow A_{21} L_{11}^{-t}$$

4. LAWRITE: write out factors

Otherwise, partition storage into panels X and Y. We exploit symmetry by operating on only the lower triangular part of matrix A in panel Y. Thus for the same amount of storage, the width of panel Y increases as the factorization proceeds.

2. We compute updates into panel Y by bringing in previous factors NNB columns at a time into panel X.

1. LAREAD: read in part of lower triangular factor into panel X
  2. P<sub>x</sub>SYRK: symmetric update to diagonal block of panel Y
  3. P<sub>x</sub>GEMM: update remaining columns in panel Y
3. Once all previous updates are performed, we perform a right-looking in-core factorization of panel Y. Loop over each block column (width NB) in panel Y,
1. factor diagonal block on one processor using P<sub>x</sub>POTRF
  2. update same block column using P<sub>x</sub>TRSM
  3. symmetric update of diagonal block using P<sub>x</sub>SYRK
  4. update remaining columns in panel Y using P<sub>x</sub>GEMM

Finally the computed factors are written out to disk.

Although, only the lower triangular portion of matrix A is used in the computation, the code still requires disk storage for the full matrix to be compatible with ScaLAPACK. ScaLAPACK routine P<sub>x</sub>POTRF accepts only a square matrix distributed with square sub-blocks, MB=NB.

#### 4. Numerical Results

The prototype code is still under active development and testing<sup>2</sup>. The double precision version was tested on the Intel Paragon systems at the Center for Computational Sciences, Oak Ridge National Laboratory. The xps35 has 512 GP nodes arranged in a 16 row by 32 column rectangular mesh. Each GP node has 32MBytes of memory. The xps150 has 1024 MP nodes arranged in a 16 row by 64 column rectangular mesh. Each MP node has at least 64MBytes of memory. The MP node has 2 compute CPUs to support multi-threaded code, but to make results comparable to xps35, only one CPU was utilized in the test. The runs were performed in a multiuser (non-dedicated) environment. Runs on 64 (256) processors were performed on the xps35 (xps150) using a  $8 \times 8$  ( $16 \times 16$ ) *logical* processor grid. The xps150 was used to ensure that in-core solves of the large matrices are resident in memory without page faults to disk.

Initial experiments suggest that I/O performance may vary by a wide margin and depends on the I/O and paging requests in other applications. The double precision version was tested with block size of MB = NB = 50, MMB = 800 and NNB = 400. A shared file was used on '/pfs' parallel file system (16-way interleaved RAID system with 64Kbyte stripes). The shared file was opened with NX-specific M\_ASYNC mode in the `gopen()` system call.

Table 1 shows the runtime (in seconds) for the out-of-core LU factorization on the Intel Paragon. The field *lwork* is the amount of temporary storage (number of double precision

---

<sup>2</sup>The prototype code is available from <http://www.netlib.org/scalapack/prototype>



Table 1: Performance of out-of-core LU factorization on 64 processors using  $MB=NB=50$ .

size of matrix	lwork (doubles)	update (sec)	fact (sec)	reorder (sec)	total (sec)	in-core (processors)
5000	130000	38	28	18	151	59 (64)
8000	250000	126	60	49	389	180 (64)
10000	375000	231	95	74	640	130 (256)
16000	1000000	858	301	192	1946	388 (256)
20000	1000000	1782	377	290	3502	681 (256)

Table 2: Performance of out-of-core QR factorization on 64 processors using  $MB=NB=50$ .

size of matrix	lwork (doubles)	update (sec)	fact (sec)	total (sec)	in-core (processors)
5000	130000	78	41	176	92 (64)
8000	260000	271	98	516	310 (64)
10000	410000	496	161	900	200 (256)
16000	1000000	1816	536	2893	647 (256)
20000	1000000	3805	680	5466	1176 (256)

numbers) available to the out-of-core routine for panels  $X$  and  $Y$ . Field *update* is the computation time (excluding I/O) for  $P_{XTRSM}$  and  $P_{XGEMM}$  updates from panel  $X$  to panel  $Y$ . Field *fact* is the total computation time (excluding I/O) required to factor panel  $Y$ . Field *reorder* is the total time for I/O and  $P_{XLAPIV}$  to reorder the lower triangular factors into the final pivoted order. Field *in-core* shows the computation time (and number of processors used) for an all in-core factorization using ScaLAPACK PDGETRF routine.

We are considering streamlining the out-of-core  $PF_{XGETRF}$  LU factorization code (and  $PF_{XGETRS}$  right-hand solver) to leave the lower factors in partially pivoted form and avoid the extra pass required to reorder the lower triangular matrix into final pivoted order. Note that without this extra reordering cost and assuming perfect speedup from 64 to 256 processors, the out-of-core solver incurs approximately a 18% overhead over in-core solvers  $((3502 - 290)/(681 * 4) \approx 1.18)$ .

Table 2 shows the runtime (in seconds) for the out-of-core QR factorization on the Intel Paragon. The field *lwork* is the amount of temporary storage (number of double precision numbers) available to the out-of-core routine for panels  $X$  and  $Y$ . Field *update* is the computation time (excluding I/O) for Householder updates using  $P_{XORMQR}$  from panel  $X$  to panel  $Y$ . Field *fact* is the total computation time (excluding I/O) required to factor panel  $Y$  using  $P_{XGEQRF}$ . Field *in-core* shows the computation time (and number of processors used) for an all in-core factorization using the ScaLAPACK PDGEQRF routine. For large problems (and assuming perfect speedup), the out-of-core version incurs an overhead of around 16% over the in-core solver  $((5466/4)/1176 \approx 1.16)$ .

Table 3: Performance of out-of-core Cholesky factorization on 64 processors using  $MB=NB=50$ .

size of matrix	lwork (doubles)	update (sec)	fact (sec)	total (sec)	in-core (processors)
5000	130000	20	18	77	39 (64)
8000	260000	56	45	196	90 (64)
10000	410000	93	78	311	60 (256)
16000	1000000	339	264	937	191 (256)
20000	1000000	776	354	1655	340 (256)

Table 3 shows the runtime (in seconds) for the out-of-core Cholesky factorization on the Intel Paragon. The field *lwork* is the amount of temporary storage (number of double precision numbers) available to the out-of-core routine for panels X and Y. Field *update* is the computation time (excluding I/O) for  $P_xSYRK$  and  $P_xGEMM$  updates from panel X to panel Y. Field *fact* is the total computation time (excluding I/O) required to factor panel Y. Field *in-core* shows the computation time (and number of processors used) for an all in-core factorization using ScaLAPACK  $PDPOTRF$  routine. For large problems (and assuming perfect speedup), the out-of-core version incurs about a 22% overhead over the in-core version  $((1655/4)/340 \approx 1.22)$ .

## 5. Conclusions

Effectiveness of the out-of-core solvers depends in part on the amount of available core memory and on the performance of the I/O system. The results on the xps35 suggest that the out-of-core solvers are most effective on very large problems greater than available core memory and incur about a 20% penalty over the in-core solvers.

## 6. References

- [1] J. CHOI, J. J. DONGARRA, L. S. OSTROUCHOV, A. P. PETITET, D. W. WALKER, AND R. C. WHALEY, *The design and implementation of the ScaLAPACK LU, QR, and Cholesky factorization routines*, Tech. Report ORNL/TM-12470, Oak Ridge National Laboratory, 1994.
- [2] J. DONGARRA, S. HAMMARLING, AND D. WALKER, *Key concepts for parallel out-of-core LU factorization*, *Scientific Programming*, 5 (1996), pp. 173–184.
- [3] K. KLIMKOWSKI AND R. A. VAN DE GEIJN, *Anatomy of a parallel out-of-core dense linear solver*, in *Proceedings of the International Conference on Parallel Processing*, 1995.
- [4] S. TOLEDO, *Locality of reference in lu decomposition with partial pivoting*, Tech. Report RC 20344(1/19/96), IBM Research Division, T. J. Watson Research Center, Yorktown Heights, New York, 1996.
- [5] S. TOLEDO AND F. GUSTAVSON, *The design and implementation of SOLAR, a portable library for scalable out-of-core linear algebra computations*, in *IOPADS Fourth Annual Workshop on Parallel and Distributed I/O*, ACM Press, 1996, pp. 28–40.

**INTERNAL DISTRIBUTION**

- |                      |  |
|----------------------|--|
| 1. T. S. Darland     | 14. S. A. Raby                               |
| 2-6. E. F. D'Azevedo | 15. R. F. Sincovec                           |
| 7-11. J. Dongarra    | 16. Laboratory Records - RC                  |
| 12. M. R. Leuze      | 17-18. Laboratory Records<br>Department/OSTI |
| 13. C. E. Oliver     | 19. Central Research Library                 |

**EXTERNAL DISTRIBUTION**

20. Daniel A. Hitchcock, ER-31, Acting Director, Mathematical, Information, and Computational Sciences Division, Office of Computational and Technology Research, Office of Energy Research, Department of Energy, Washington, DC 20585
21. Frederick A. Howes, ER-31, Mathematical, Information, and, Computational Sciences Division, Office of Computational and Technology Research, Office of Energy Research, Department of Energy, Washington, DC 20585
22. Tom Kitchens, ER-31, Mathematical, Information, and, Computational Sciences Division, Office of Computational and Technology Research, Office of Energy Research Department of Energy, Washington, DC 20585
23. David B. Nelson, ER-30, Associate Director, Office of Energy Research, Director, Office of Computational and Technology Research, Department of Energy, Washington, DC 20585