

Recursive Blocked Algorithms for Solving Triangular Systems—Part I: One-Sided and Coupled Sylvester-Type Matrix Equations

ISAK JONSSON and BO KÅGSTRÖM
Umeå University

Triangular matrix equations appear naturally in estimating the condition numbers of matrix equations and different eigenspace computations, including block-diagonalization of matrices and matrix pairs and computation of functions of matrices. To solve a triangular matrix equation is also a major step in the classical Bartels–Stewart method for solving the standard continuous-time Sylvester equation ($AX - XB = C$). We present novel recursive blocked algorithms for solving one-sided triangular matrix equations, including the continuous-time Sylvester and Lyapunov equations, and a generalized coupled Sylvester equation. The main parts of the computations are performed as level-3 general matrix multiply and add (GEMM) operations. In contrast to explicit standard blocking techniques, our recursive approach leads to an automatic variable blocking that has the potential of matching the memory hierarchies of today's HPC systems. Different implementation issues are discussed, including when to terminate the recursion, the design of new optimized superscalar kernels for solving leaf-node triangular matrix equations efficiently, and how parallelism is utilized in our implementations. Uniprocessor and SMP parallel performance results of our recursive blocked algorithms and corresponding routines in the state-of-the-art libraries LAPACK and SLICOT are presented. The performance improvements of our recursive algorithms are remarkable, including 10-fold speedups compared to standard algorithms.

Categories and Subject Descriptors: F.2.1 [Analysis of Algorithms and Problem Complexity]: Numerical Algorithms and Problems—*computations on matrices*; G.1.3 [Numerical Analysis]: Numerical Linear Algebra—*conditioning, linear systems*; G.4 [Mathematical Software]: Algorithm design and analysis, efficiency, parallel and vector implementations, portability, reliability and robustness

General Terms: Algorithms, Performance

Additional Key Words and Phrases: Matrix equations, standard Sylvester and Lyapunov, generalized coupled Sylvester, recursion, automatic blocking, superscalar, GEMM-based, level-3 BLAS, SMP parallelization, LAPACK, SLICOT

This research was conducted using the resources of the High Performance Computing Center North (HPC2N) and PDC-Paralleldatorcentrum at KTH, Stockholm. During this work resources at the Danish Computing Centre for Research and Education (UNI-C), Lyngby, Denmark were also used. Financial support was provided by the Swedish Research Council under grants TFR 98-604 and VR 621-2001-3284.

Authors' address: Department of Computing Science and HPC2N, Umeå University, SE-901 87 Umeå, Sweden; email: {isak,bokg}@cs.umu.se.

Permission to make digital/hard copy of part or all of this work for personal or classroom use is granted without fee provided that the copies are not made or distributed for profit or commercial advantage, the copyright notice, the title of the publication, and its date appear, and notice is given that copying is by permission of the ACM, Inc. To copy otherwise, to republish, to post on servers, or to redistribute to lists, requires prior specific permission and/or a fee.

© 2002 ACM 0098-3500/02/1200-0392 \$5.00

1. INTRODUCTION

We introduce and discuss new recursive blocked algorithms for solving various types of triangular matrix equations. Our goal is to design efficient algorithms for state-of-the-art HPC systems with deep memory hierarchies. A recursive algorithm leads to automatic blocking which is variable and “squarish” [Gustavson 1997]. This hierarchical blocking promotes good data locality, which makes it possible to approach peak performance on state-of-the-art processors with several levels of cache memories. Recently, several successful results have been published. Using the current standard data layouts and applying recursive blocking have led to faster algorithms for the Cholesky, LU , and QR factorizations [Toledo 1997; Gustavson 1997; Elmroth and Gustavson 2001]. Moreover, in Gustavson et al. [1998a], we demonstrated that a further performance gain can be obtained when recursive dense linear algebra algorithms are expressed using a recursive data layout and highly optimized superscalar kernels.

In this contribution (Part I), which extends our work in Jonsson and Kågström [2001a], we consider one-sided Sylvester-type equations, including the continuous-time standard Sylvester ($AX - XB = C$) and Lyapunov ($AX + XA^T = C$) equations, and a generalized coupled Sylvester equation ($AX - YB, DX - YE = (C, F)$). We use the notation *one-sided*, since the matrix equations include terms where the solution is only involved in matrix products of two matrices (e.g., $\text{op}(A)X$ or $X\text{op}(A)$, where $\text{op}(A)$ can be A or A^T). These equations all appear in various control theory applications and in spectral analysis. In Part II [Jonsson and Kågström 2002], we consider *two-sided* matrix equations, which include matrix product terms of type $\text{op}(A)X\text{op}(B)$. Examples include discrete-time standard and generalized Sylvester and Lyapunov equations.

The classical method of solution is the Bartels–Stewart method, which includes three major steps [Bartels and Stewart 1972]. First, the matrix (or matrix pair) is transformed to a Schur (or generalized Schur) form. This leads to a reduced triangular matrix equation, which is solved in the second step. For example, the coefficient matrices A and B in the Sylvester equation $AX - XB = C$ are in upper triangular or upper quasitriangular form [Kågström and Poromaa 1992; Poromaa 1998]. For the generalized counterpart, the matrix pairs (A, D) and (B, E) in $(AX - YB, DX - YE) = (C, F)$ are reduced to generalized Schur form, with A and B upper quasitriangular, and D, E upper triangular [Kågström and Westin 1989; Kågström and Poromaa 1996a]. Finally, the solution of the reduced matrix equation is transformed back to the original coordinate system. In this article, we focus on the solution of the reduced triangular matrix equations. Reliable and efficient algorithms for the reduction step can be found in LAPACK [Anderson et al. 1999] for the standard case, and in Dackland and Kågström [1999] for the generalized case, where a blocked variant of the QZ method is presented.

Triangular matrix equations also appear naturally in estimating the condition numbers of matrix equations and different eigenspace computations, including block-diagonalization of matrices and matrix pairs and computation of functions of matrices. Related applications include the direct reordering

of eigenvalues in the real (generalized) Schur form [Anderson et al. 1999; Kågström and Poromaa 1996b] and the computation of additive decompositions of a (generalized) transfer function [Kågström and Van Dooren 1992].

Before we go into any further details, we outline the contents of the rest of the article. In Section 2, we illustrate how the solutions of triangular matrix equations relate to different spectral condition estimation problems and to the estimation of the conditioning of the Sylvester-type equations themselves. Section 3 introduces our recursive blocked algorithms for one-sided and coupled triangular matrix equations, including the standard Sylvester (Section 3.1) and Lyapunov (Section 3.2) equations, and the generalized coupled Sylvester equation (Section 3.3). In Section 4, we introduce a recursive blocked algorithm for computing functions of matrices, which, indeed, is an application of solving triangular Sylvester equations.

In Section 5, we discuss different implementation issues, including when to terminate the recursion, and the design of optimized superscalar kernels for solving small-sized (leaf-node) triangular matrix equations efficiently. We also discuss how parallelism is utilized in our implementations. Sample performance results of our recursive blocked algorithms are presented and discussed in Section 6. Finally, we give some conclusions.

2. CONDITION ESTIMATION AND TRIANGULAR SOLVERS

Besides solving a matrix equation it is equally important to have reliable error bounds of the computed quantities. In this section, we review how triangular matrix equations enter in condition estimation of common eigenspace computations, as well as in the condition estimation of the matrix equations themselves. Since condition estimation typically involves solving several triangular matrix equations, one could say that this is the main source for different triangular matrix equation problems. It is important that these matrix equations can be solved with efficient algorithms on today's memory-tiered systems.

2.1 Condition Estimation of Some Eigenspace Problems

We assume that S is a real block partitioned matrix in *real Schur form*:

$$S = \begin{bmatrix} A & -C \\ 0 & B \end{bmatrix}.$$

This means that both A and B are *upper quasitriangular*, that is, block upper triangular with 1×1 and 2×2 diagonal blocks, which correspond to real and complex conjugate pairs of eigenvalues, respectively. Typically, the partitioning is done with some application in mind. For example, A may include all eigenvalues in the left complex plane and we want the invariant subspaces associated with the spectra of A and B , respectively. Now, S can be block-diagonalized by a similarity transformation:

$$\begin{bmatrix} I & -X \\ 0 & I \end{bmatrix} S \begin{bmatrix} I & X \\ 0 & I \end{bmatrix} = \begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix},$$

where X satisfies the triangular Sylvester equation $AX - XB = C$. Knowing X , we also know the invariant subspaces and the *spectral projector* associated with the block A :

$$P = \begin{bmatrix} I & X \\ 0 & 0 \end{bmatrix}.$$

This is an important quantity in error bounds for invariant subspaces and clusters of eigenvalues. A large value of $\|P\|_2 = (1 + \|X\|_2^2)^{1/2}$ signals ill-conditioning. To avoid possible overflow, we use the computed estimate $s = 1/\|P\|_F$ [Bai et al. 1993; Kågström and Poromaa 1992; Anderson et al. 1999].

Next we consider a real regular matrix pair (S, T) in *real generalized Schur form*:

$$(S, T) = \left(\begin{bmatrix} A & -C \\ 0 & B \end{bmatrix}, \begin{bmatrix} D & -F \\ 0 & E \end{bmatrix} \right),$$

where A, B , as before, are upper quasitriangular and D, E are upper triangular. Any 2×2 diagonal block in the generalized Schur form corresponds to a pair of complex conjugate eigenvalues. The 1×1 diagonal blocks correspond to real eigenvalues. For example, if $e_{ii} \neq 0$, then a_{ii}/e_{ii} is a finite eigenvalue of (A, D) ; otherwise ($e_{ii} = 0$), the matrix pair has an infinite eigenvalue.

Now, (S, T) is block diagonalized by an equivalence transformation:

$$\begin{bmatrix} I & -Y \\ 0 & I \end{bmatrix} (S, T) \begin{bmatrix} I & X \\ 0 & I \end{bmatrix} = \left(\begin{bmatrix} A & 0 \\ 0 & B \end{bmatrix}, \begin{bmatrix} D & 0 \\ 0 & E \end{bmatrix} \right),$$

where (X, Y) satisfies the triangular coupled Sylvester equation $(AX - YB, DX - YE) = (C, F)$. Knowing (X, Y) , we also know pairs of deflating subspaces associated with the matrix pairs (A, D) and (B, E) . Similarly to the matrix case, large values on the *left* and *right projector norms* $l = (1 + \|Y\|_F^2)^{1/2}$ and $r = (1 + \|X\|_F^2)^{1/2}$ signal ill-conditioning [Kågström and Poromaa 1996b; Anderson et al. 1999], that is, the eigenspaces may be sensitive to small perturbations in the data.

2.2 Condition Estimation of Matrix Equations

All linear matrix equations can be written as a linear system of equations $Zx = c$, where Z is a *Kronecker product matrix representation* of the associated Sylvester-type operator, and the solution x and the right-hand side c are represented in $\text{vec}(\cdot)$ notation. $\text{vec}(X)$ denotes a column vector with the columns of X stacked on top of each other. We introduce the following Z -matrices,

$$\begin{aligned} Z_{\text{SYCT}} &= Z_{AX-XB} = I_n \otimes A - B^T \otimes I_m, \\ Z_{\text{LYCT}} &= Z_{AX+XA^T} = I_n \otimes A + A \otimes I_n, \\ Z_{\text{GCSY}} &= Z_{(AX-YB, DX-YE)} = \begin{bmatrix} I_n \otimes A & -B^T \otimes I_m \\ I_n \otimes D & -E^T \otimes I_m \end{bmatrix}. \end{aligned}$$

From top to bottom they represent the matrix operators of the continuous-time Sylvester and Lyapunov equations, and the generalized coupled Sylvester equation.

An important quantity that appears both in the perturbation theory for Sylvester-type equations and for the eigenspace problems considered above is the *separation between two matrices* [Stewart and Sun 1990], defined as

$$\text{Sep}[A, B] = \inf_{\|X\|_F=1} \|AX - XB\|_F = \sigma_{\min}(Z_{\text{SYCT}}),$$

where $\sigma_{\min}(Z_{\text{SYCT}}) \geq 0$ is the smallest singular value of Z_{SYCT} . We review some of its characteristics: $\text{Sep}[A, B] = 0$ if and only if A and B have a common eigenvalue; $\text{Sep}[A, B]$ is small if there is a small perturbation of A or B that makes them have a common eigenvalue. The Sep-function may be much smaller than the minimum distance between the eigenvalues of A and B .

Assuming M and N are the dimensions of A and B , computing $\sigma_{\min}(Z_{\text{SYCT}})$ is an $O(M^3N^3)$ operation, which is impractical already for moderate values on M and N . In Kågström and Poromaa [1992], it is shown how reliable Sep^{-1} -estimates can be computed to the cost $O(MN^2 + M^2N)$ by solving triangular matrix equations:

$$\frac{\|x\|_2}{\|c\|_2} = \frac{\|X\|_F}{\|C\|_F} \leq \|Z_{\text{SYCT}}^{-1}\|_2 = \frac{1}{\sigma_{\min}(Z_{\text{SYCT}})} = \text{Sep}^{-1}.$$

The right-hand side C is chosen such that the lower bound gets as large as possible. This leads to a Frobenius-norm-based estimate. For computation of 1-norm-based estimates see Hager [1984], Higham [1988], and Kågström and Poromaa [1992].

The Sep-functions associated with the Sylvester-type matrix equations are:

$$\begin{aligned} \text{Sep}[\text{SYCT}] &= \inf_{\|X\|_F=1} \|AX - XB\|_F = \sigma_{\min}(Z_{\text{SYCT}}), \\ \text{Sep}[\text{LYCT}] &= \inf_{\|X\|_F=1} \|AX - X(-A^T)\|_F = \sigma_{\min}(Z_{\text{LYCT}}), \\ \text{Sep}[\text{GCSY}] &= \inf_{\|(X,Y)\|_F=1} \|(AX - YB, DX - YE)\|_F = \sigma_{\min}(Z_{\text{GCSY}}). \end{aligned}$$

The same techniques as presented above can also be used for estimating $\text{Sep}[\text{LYCT}]$. Reliable estimates of $\text{Sep}[\text{GCSY}]$ (the separation between two matrix pairs [Stewart and Sun 1990]) are presented and discussed in Kågström and Poromaa [1996a,b] and Kågström and Westin [1989]. The underlying perturbation theory for these Sylvester-type equations is presented in Higham [1993] and Kågström [1994]. See also the nice review in Chapter 15 of Higham [1996].

3. RECURSIVE BLOCKED ALGORITHMS FOR ONE-SIDED AND COUPLED MATRIX EQUATIONS

As mentioned in the introduction, the standard methods for solving one-sided matrix equations are all based on the *BS*-method [Bartels and Stewart 1972]. The fundamental algorithms for solving the continuous-time Sylvester and Lyapunov equations are presented in Bartels and Stewart [1972], Golub et al. [1979], and Hammarling [1982]. Generalizations of the Bartels–Stewart and

the Hessenberg–Schur [Golub et al. 1979] methods for solving the generalized coupled Sylvester equation can be found in Kågström and Westin [1989], and Kågström and Promaa [1996a].

In this section, we present our recursive blocked methods for solving triangular one-sided and coupled matrix equations, which also are amenable for parallelization, especially on shared memory systems. Parallel methods for solving triangular one-sided and coupled matrix equations have also been studied (e.g., see Kågström and Poromaa [1992] and Poromaa [1998, 1997]).

For each matrix equation we define recursive splittings which in turn lead to a few smaller problems to be solved. These recursive splittings are applied to all “half-sized” triangular matrix equations and so on. We terminate the recursion when the new problem sizes (M and/or N) are smaller than a certain block size, $blks$, which is chosen such that at least a few submatrices involved in the current matrix equation fit in the first-level cache memory. For the solution of the small-sized problems, we apply new high-performance kernels (see Section 5).

We remark that all updates with respect to the solution of subproblems in the recursion are general matrix multiply and add (GEMM) operations $C \leftarrow \beta C + \alpha AB$, where α and β are real scalars. This is due to the “one-sidedness” of the matrix equations. Using the GEMM-based approach, some of them can be reorganized in efficient symmetric rank- $2k$ (SYR2K) operations [Dongarra et al. 1990a,b; Kågström et al. 1998a,b].

3.1 Recursive Triangular Continuous-Time Sylvester Solvers

Consider the real *continuous-time Sylvester* (SYCT) matrix equation

$$AX - XB = C, \quad (1)$$

where A of size $M \times M$ and B of size $N \times N$ are upper triangular or upper quasitriangular, that is, in real Schur form. The right-hand side C and the solution X are of size $M \times N$. Typically, the solution overwrites the right-hand side ($C \leftarrow X$). The SYCT equation (1) has a *unique solution* if and only if A and B have no eigenvalue in common or, equivalently, $\text{Sep}[\text{SYCT}] \neq 0$.

Depending on the sizes of M and N , we consider three alternatives for doing a *recursive splitting*.

Case 1 ($1 \leq N \leq M/2$). We split A by rows and columns, and C by rows only:

$$\begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} - \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} B = \begin{bmatrix} C_1 \\ C_2 \end{bmatrix},$$

or equivalently

$$\begin{aligned} A_{11}X_1 - X_1B &= C_1 - A_{12}X_2, \\ A_{22}X_2 - X_2B &= C_2. \end{aligned}$$

The original problem is split in two triangular Sylvester equations: we solve for X_2 and after the GEMM update $C_1 = C_1 - A_{12}X_2$, we can solve for X_1 using the recursive algorithm.

Case 2 ($1 \leq M \leq N/2$). We split B by rows and columns, and C by columns only:

$$A[X_1 \ X_2] - [X_1 \ X_2] \begin{bmatrix} B_{11} & B_{12} \\ & B_{22} \end{bmatrix} = [C_1 \ C_2],$$

or equivalently

$$\begin{aligned} AX_1 - X_1 B_{11} &= C_1, \\ AX_2 - X_2 B_{22} &= C_2 + X_1 B_{12}. \end{aligned}$$

Now, we first solve for X_1 and then after updating C_2 with respect to X_1 , solve for X_2 .

Case 3 ($N/2 < M < 2N$). We split A , B , and C by rows and columns:

$$\begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} - \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ & B_{22} \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}.$$

This recursive splitting results in the following four triangular Sylvester equations:

$$\begin{aligned} A_{11}X_{11} - X_{11}B_{11} &= C_{11} - A_{12}X_{21}, \\ A_{11}X_{12} - X_{12}B_{22} &= C_{12} - A_{12}X_{22} + X_{11}B_{12}, \\ A_{22}X_{21} - X_{21}B_{11} &= C_{21}, \\ A_{22}X_{22} - X_{22}B_{22} &= C_{22} + X_{21}B_{12}. \end{aligned}$$

We start by solving for X_{21} in the third equation. After updating C_{11} and C_{22} with respect to X_{21} , we can solve for X_{11} and X_{22} . Both updates and the triangular Sylvester solves are independent operations and can be executed concurrently. Finally, we update C_{12} with respect to X_{11} and X_{22} , and solve for X_{12} .

We remark that if a splitting point ($M/2$ or $N/2$) appears at a 2×2 diagonal block, the matrices are split just below this diagonal block. All “half-sized” Sylvester equations are solved using a recursive blocked algorithm.

In the discussion above, we have assumed that both A and B are upper triangular (or quasitriangular). However, it is straightforward to derive similar recursive splittings for the triangular SYCT, where A and B can each be in either upper or lower Schur form.

A Matlab-style function $[X] = \mathbf{rtrsylct}(A, B, C, \text{uplo}, \text{blks})$ implementing our recursive blocked solver is presented in Algorithm 1. The input *uplo* signals the triangular structure of A and B . The function $[X] = \mathbf{trsylct}(A, B, C, \text{uplo})$ implements an algorithm for solving triangular Sylvester block kernel problems (see Section 5). The function $[C] = \mathbf{gemm}(A, B, C)$ implements the GEMM-operation $C = C + AB$.

Algorithm 1: rtrgsyl

Input: A ($M \times M$) and B ($N \times N$) in quasitriangular Schur form. C ($M \times N$) dense matrix. $blks$, block size that specifies when to switch to a standard algorithm for solving small-sized matrix equations. $uplo$ indicates triangular form of A , B : 1-(upper, upper), 2-(upper, lower), 3-(lower, upper), 4-(lower, lower).

Output: X ($M \times N$), the solution of $AX - XB = C$. X is allowed to overwrite C .

```

function [X] = rtrsyct(A, B, C, uplo, blks)
if 1 ≤ M, N ≤ blks then
    X = trsyct(A, B, C, uplo);
else switch uplo
case 1
    if 1 ≤ N ≤ M/2 % Case 1: Split A (by rows and columns), C (by rows only)
        X2 = rtrsyct(A22, B, C2, 1, blks);
        C1 = gemm(−A12, X2, C1);
        X1 = rtrsyct(A11, B, C1, 1, blks);
        X = [X1; X2];
    elseif 1 ≤ M ≤ N/2 % Case 2: Split B (by rows and columns), C (by columns only)
        X1 = rtrsyct(A, B11, C1, 1, blks);
        C2 = gemm(X1, B12, C2);
        X2 = rtrsyct(A, B22, C2, 1, blks);
        X = [X1, X2];
    else % M, N ≥ blks, Case 3: Split A, B and C (all by rows and columns)
        X21 = rtrsyct(A22, B11, C21, 1, blks);
        C22 = gemm(X21, B12, C22); C11 = gemm(−A12, X21, C11);
        X22 = rtrsyct(A22, B22, C22, 1, blks); X11 = rtrsyct(A11, B11, C11, 1, blks);
        C12 = gemm(−A12, X22, C12);
        C12 = gemm(X11, B12, C12);
        X12 = rtrsyct(A11, B22, C12, 1, blks);
        X = [X11, X12; X21, X22];
    end
case 2 % Code for uplo = 2.
case 3 % Code for uplo = 3.
case 4 % Code for uplo = 4.
end

```

Algorithm 1. Recursive blocked algorithm for solving the triangular continuous-time Sylvester equation.

3.2 Recursive Triangular Continuous-Time Lyapunov Solvers

Consider the real *continuous-time Lyapunov* (LYCT) matrix equation

$$AX + XA^T = C, \quad (2)$$

where A is upper triangular or upper quasitriangular, that is, in real Schur form. The right-hand side C , the solution X , and A are all of size $N \times N$. Typically, the solution overwrites the right-hand side ($C \leftarrow X$). The LYCT equation (2) has a *unique solution* if and only if $\lambda_i(A) + \lambda_j(A) \neq 0$ for all i and j , or equivalently $\text{Sep}[\text{LYCT}] \neq 0$. If $C = C^T$ is (semi)definite and $\text{Re}\lambda_i(A) < 0$ for all i , then a unique (semi)definite solution X exists [Hammarling 1982].

Since all matrices are of the same size, there is one parameter (N) that controls the *recursive splitting*. We split A and C by rows and columns:

$$\begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} + \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} \begin{bmatrix} A_{11}^T & \\ & A_{22}^T \end{bmatrix} = \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}.$$

If $C = C^T$, we use the fact that $X_{21} = X_{12}^T$ and the recursive splitting leads to three triangular matrix equations:

$$\begin{aligned} A_{11}X_{11} + X_{11}A_{11}^T &= C_{11} - A_{12}X_{12}^T - X_{12}A_{12}^T, \\ A_{11}X_{12} + X_{12}A_{22}^T &= C_{12} - A_{12}X_{22}, \\ A_{22}X_{22} + X_{22}A_{22}^T &= C_{22}. \end{aligned}$$

The first and the third are triangular Lyapunov equations, while the second is a triangular Sylvester matrix equation. We start by solving for X_{22} in the third equation. After updating C_{12} with respect to X_{22} , we can solve for X_{12} . Finally, we update the right-hand side of the first matrix equation with respect to X_{12} , which is a symmetric rank- $2k$ (SYR2K) operation $C_{11} = C_{11} - A_{12}X_{12}^T - X_{12}A_{12}^T$, and solve for X_{11} .

A Matlab-style function $[X] = \mathbf{rtrlyct}(A, C, blks)$ implementing our blocked recursive solver, which deals with the cases with a symmetric or nonsymmetric C , is presented in Jonsson and Kågström [2001b]. For solving the triangular Sylvester equations that appear, we make use of the recursive algorithm $[X] = \mathbf{rtrsylv}(A, B, C, uplo, blks)$, described in Section 3.1.

3.3 Recursive Triangular Coupled Sylvester Solvers

Consider the real *generalized coupled Sylvester* (GCSY) matrix equation of the form

$$\begin{aligned} AX - YB &= C, \quad C \leftarrow X (M \times N), \\ DX - YE &= F, \quad F \leftarrow Y (M \times N), \end{aligned} \tag{3}$$

where the matrix pairs (A, D) and (B, E) are in generalized Schur form with A, B upper (quasi-)triangular and D, E upper triangular. A, D are $M \times M$ and B, E are $N \times N$; the right-hand sides C, D and the solution matrices X, Y , which overwrite C and D , respectively, are of size $M \times N$. The GCSY equation (3) has a *unique solution* if and only if (A, D) and (B, E) are regular matrix pairs and have no eigenvalue in common, or equivalently $\text{Sep}[\text{GCSY}] \neq 0$.

As for the standard case, we consider three alternatives for doing a *recursive splitting*. Below, we illustrate these cases.

Case 1 ($1 \leq N \leq M/2$). We split (A, D) by rows and columns, and (C, F) by rows only:

$$\begin{aligned} \begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} - \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} B &= \begin{bmatrix} C_1 \\ C_2 \end{bmatrix}, \\ \begin{bmatrix} D_{11} & D_{12} \\ & D_{22} \end{bmatrix} \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} - \begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} E &= \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}. \end{aligned}$$

The splitting results in the following two generalized Sylvester equations:

$$\begin{aligned} A_{11}X_1 - Y_1B &= C_1 - A_{12}X_2, & D_{11}X_1 - Y_1E &= F_1 - D_{12}X_2, \\ A_{22}X_2 - Y_2B &= C_2, & D_{22}X_2 - Y_2E &= F_2. \end{aligned}$$

First, we solve for (X_2, Y_2) in the second pair of matrix equations. After updating (C_1, F_1) with respect to X_2 (two GEMM operations that can execute in parallel), we solve for (X_1, Y_1) .

Case 2 ($1 \leq M \leq N/2$). We split (B, E) by rows and columns, and (C, F) by columns only:

$$\begin{aligned} A[X_1 \ X_2] - [Y_1 \ Y_2] \begin{bmatrix} B_{11} & B_{12} \\ & B_{22} \end{bmatrix} &= [C_1 \ C_2], \\ D[X_1 \ X_2] - [Y_1 \ Y_2] \begin{bmatrix} E_{11} & E_{12} \\ & E_{22} \end{bmatrix} &= [F_1 \ F_2], \end{aligned}$$

or equivalently

$$\begin{aligned} AX_1 - Y_1B_{11} &= C_1, & DX_1 - Y_1E_{11} &= F_1, \\ AX_2 - Y_2B_{22} &= C_2 + Y_1B_{12}, & DX_2 - Y_2E_{22} &= F_2 + Y_1E_{12}. \end{aligned}$$

Now we first solve for (X_1, Y_1) , and after updating (C_2, F_2) with respect to Y_1 (also two independent GEMM operations), we solve for (X_2, Y_2) .

Case 3 ($N/2 < M < 2N$). We split all matrices by rows and columns:

$$\begin{aligned} \begin{bmatrix} A_{11} & A_{12} \\ & A_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} - \begin{bmatrix} Y_{11} & Y_{12} \\ Y_{21} & Y_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ & B_{22} \end{bmatrix} &= \begin{bmatrix} C_{11} & C_{12} \\ C_{21} & C_{22} \end{bmatrix}, \\ \begin{bmatrix} D_{11} & D_{12} \\ & D_{22} \end{bmatrix} \begin{bmatrix} X_{11} & X_{12} \\ X_{21} & X_{22} \end{bmatrix} - \begin{bmatrix} Y_{11} & Y_{12} \\ Y_{21} & Y_{22} \end{bmatrix} \begin{bmatrix} E_{11} & E_{12} \\ & E_{22} \end{bmatrix} &= \begin{bmatrix} F_{11} & F_{12} \\ F_{21} & F_{22} \end{bmatrix}. \end{aligned}$$

This leads to four pairs of coupled Sylvester equations to be solved:

$$\begin{aligned} A_{11}X_{11} - Y_{11}B_{11} &= C_{11} - A_{12}X_{21}, & D_{11}X_{11} - Y_{11}E_{11} &= F_{11} - D_{12}X_{21}, \\ A_{11}X_{12} - Y_{12}B_{22} &= C_{12} - A_{12}X_{22} + Y_{11}B_{12}, & D_{11}X_{12} - Y_{12}E_{22} &= F_{12} - D_{12}X_{22} + Y_{11}E_{12}, \\ A_{22}X_{21} - Y_{21}B_{11} &= C_{21}, & D_{22}X_{21} - Y_{21}E_{11} &= F_{21}, \\ A_{22}X_{22} - Y_{22}B_{22} &= C_{22} + Y_{21}B_{12}, & D_{22}X_{22} - Y_{22}E_{22} &= F_{22} + Y_{21}E_{12}. \end{aligned}$$

We start by solving for (X_{21}, Y_{21}) in the third matrix equation pair. After updating (C_{11}, F_{11}) and (C_{22}, F_{22}) with respect to X_{21} and Y_{21} , respectively, we can solve for (X_{11}, Y_{11}) and (X_{22}, Y_{22}) . Finally, we update (C_{12}, F_{12}) with respect to Y_{11} and X_{22} , and solve for (X_{12}, Y_{12}) . Several of these operations can be performed concurrently, including GEMM-updates of right-hand sides and two of the generalized triangular Sylvester solves.

A Matlab-style function, $[X] = \mathbf{rtrgsy}(A, B, C, D, E, F, blks)$, which implements our recursive blocked generalized Sylvester algorithm is presented in Jonsson and Kågström [2001b]. As in the standard case, the recursion

Table I. Complexity of Standard Algorithms
Measured in Flops

Matrix Equation	Overall Cost in Flops
SYCT (1)	$M^2N + MN^2$
LYCT (2)	N^3
GCSY (3)	$2M^2N + 2MN^2$

is only applied down to a certain block size, when the function $[X] = \mathbf{trgcsy}(A, B, C, D, E, F)$ is applied for solving triangular generalized coupled Sylvester block kernel problems.

3.4 Number of Operations and Execution Order

All recursive algorithms presented in Section 3 perform the same amount of *floating point operations* (flops) as the elementwise explicit algorithms or their blocked counterparts (e.g., see Bartels and Stewart [1972], Kågström and Westin [1989], and Kågström and Poromaa [1992, 1996a]), which are all based on backward or forward substitutions with one or several right-hand sides. In Table I, we summarize the overall flopcounts for these methods. We remark that the difference in flops between the two extreme cases, that is, when all diagonal blocks of the matrices in upper Schur form are of size 1×1 or 2×2 , respectively, only shows up in the lower-order terms.

The flopcounts for the recursive blocked algorithms can be expressed in terms of the following recurrence formulas.

$$F_{\text{SYCT}}(M, N) = 4F_{\text{SYCT}}(M/2, N/2) + 2F_{\text{GEMM}}(M/2, M/2, N/2) + 2F_{\text{GEMM}}(M/2, N/2, N/2), \quad (4)$$

$$F_{\text{LYCT}}(N) = 2F_{\text{LYCT}}(N/2) + F_{\text{SYCT}}(N/2, N/2) + 2F_{\text{GEMM}}(N/2, N/2, N/2), \quad (5)$$

$$F_{\text{GCSY}}(M, N) = 4F_{\text{GCSY}}(M/2, N/2) + 4F_{\text{GEMM}}(M/2, M/2, N/2) + 4F_{\text{GEMM}}(M/2, N/2, N/2). \quad (6)$$

These expressions correspond to the most general case when the recursive splitting is by rows and by columns for all input matrices. Ignoring the lower-order terms and assuming that $F_{\text{GEMM}}(P, Q, R)$, the complexity of the GEMM operation with matrices of sizes $P \times Q$ and $Q \times R$ is $2PQR$ flops, it is straightforward to derive the expressions in Table I by induction from Equations (4) through (6), respectively.

The main difference between the recursive blocked and the standard explicitly blocked algorithms is their data reference patterns, that is, the order in which they access data and how big chunks and how many times the data are moved in the memory hierarchy of the target computer system. As is shown in Section 6, this can have a great impact on the performance of the algorithms. As expected, the algorithms with the smallest amount of redundant memory transfers show the best performance.

4. COMPUTING FUNCTIONS OF TRIANGULAR MATRICES USING RECURSIVE BLOCKING

An important application to solving triangular Sylvester equations is the problem of computing $f(A)$, where f is an analytic function and A of size $N \times N$ is a real matrix in Schur form; for example, A is upper quasitriangular. The best known matrix function is the matrix exponential, which has several applications in control theory. Different methods have been suggested over the years for computing matrix functions. We refer to the papers Kågström [1977a,b], Van Loan [1977], and Moler and Van Loan [1978], which also include perturbation theory and error bounds. Several of these results are also reviewed in Golub and Van Loan [1996], the standard textbook on advanced matrix computations.

In the following, we let F denote the matrix function $f(A)$ to be computed. Since A is upper triangular, so will F be. Moreover, since f is analytic F can be expressed in terms of a series expansion, from which it is obvious that F and A commute; that is,

$$AF - FA = 0. \quad (7)$$

We use this fact to derive a blocked recursive algorithm for computing F . Since A is square, there is only one way of doing the *recursive splitting*: namely, we split A and F by rows and columns. This recursive splitting results in the three triangular matrix equations:

$$\begin{aligned} A_{11}F_{11} - F_{11}A_{11} &= 0, \\ A_{11}F_{12} - F_{12}A_{22} &= F_{11}A_{12} - A_{12}F_{22}, \\ A_{22}F_{22} - F_{22}A_{22} &= 0. \end{aligned}$$

The first and the third reveal the commuting properties of the diagonal blocks. Knowing F_{11} and F_{22} , we get F_{12} from the second equation, which is a triangular continuous-time Sylvester equation (1).

The splitting described above can now be applied recursively to all “half-sized” triangular matrix equations. As before, we terminate the recursion when a problem size N is smaller than a certain block size, *blks*. For computing the matrix function of the small-sized kernel problems, there are several alternatives. If all eigenvalues are real, one alternative is to choose *blks* = 1, which means that the recursion continues until element level and $F_{ii} = f(\lambda_j)$ for an eigenvalue λ_j .

However, in this context, the choice of block size is also important with respect to the accuracy and reliability of the results. If the eigenvalues along the (block) diagonal of A are ordered with respect to a clustering procedure, then the blocking with respect to the clustering must be taken into account in the recursive splitting. Splittings that are not permitted can be monitored by computing an estimate of $\text{Sep}[A_{11}, A_{22}]$, each time a Sylvester equation is solved for computing an F_{12} block. A small value of Sep indicates that a small perturbation of A_{11} and/or A_{22} may cause at least one eigenvalue from the perturbed A_{11} to coalesce with an eigenvalue of A_{22} . If this is the case, the splitting of A is inadmissible and another splitting should be chosen. To facilitate the choice of splitting, the blocking with respect to clustering should be provided as input.

5. OPTIMIZED SUPERSCALAR KERNELS AND OTHER IMPLEMENTATION ISSUES

In principle, we have three levels of triangular matrix equation solvers. At the user level we have the recursive blocked **rt^r*** solvers. Each of them calls a **tr^{*}** block or subsystem solver when the current problem sizes (M and/or N) from a recursive splitting are smaller than a certain block size, *blks*. Finally, each of the block solvers calls a *superscalar kernel* for solving matrix equations with $M, N \leq 4$. There are several issues to consider in developing portable and high-performance implementations for these solvers. In this section, we discuss several of them including the impact of kernel solvers on the overall performance of the matrix equation solver, the design of superscalar kernels for solving small-sized matrix equations, and when to terminate the recursion and instead call a block solver for solving the remaining subsystems. We also discuss different aspects in the choice of BLAS implementations and in implementing shared memory parallelism.

5.1 Impact of Kernel Solvers

For illustration and without loss of generality, we assume that all matrix sizes are the same ($M = N$). Then the innermost kernels of the Sylvester-type solvers, which exist in the **tr^{*}** routines of the algorithm descriptions in Section 3, execute $O(N^2)$ flops out of the total $O(N^3)$ flops. Although one order fewer operations, the performance of the kernels is very important for the overall computation rate for these solvers. Most of the operations outside the kernels are GEMM operations.

In Figure 1, the modeled impact of the performance of Sylvester kernels on the overall performance of the recursive blocked algorithm **rt^rsyct** is illustrated. The x -axis shows the problem size N , and the y -axis shows the overall performance of **rt^rsyct** measured in Mflops/s for kernel routines with different performance characteristics. We use two fixed performance rates for the DGEMM routine, 200 and 500 Mflops/s, respectively. With $M = N$, the total number of flops is $2N^3$, and the number of flops performed by the 2×2 subsystem kernel is $4N^2$. We model the *overall computation rate* as

$$S = \frac{\text{Total number of flops}}{\text{GEMM time} + \text{Kernel time}} = \frac{2N^3}{\frac{2N^3 - 4N^2}{G} + \frac{4N^2}{K}} = \frac{N \cdot G}{N - 2 + 2\frac{G}{K}}, \quad (8)$$

where G denotes the performance of the DGEMM routine, and K denotes the performance of the kernel routine.

From Figure 1, we see that with a GEMM performance of 500 Mflops/s and the slowest kernel (0.3 Mflops/s), the overall performance is at most approaching 50% of the GEMM performance for $N = 2500$. But with the same GEMM performance and the faster kernels (executing at 3 and 10 Mflops/s, respectively), the overall performance rather quickly approaches 80 to 90% of the performance of the DGEMM routine. This is, of course, the performance behavior we would like to see in practice.

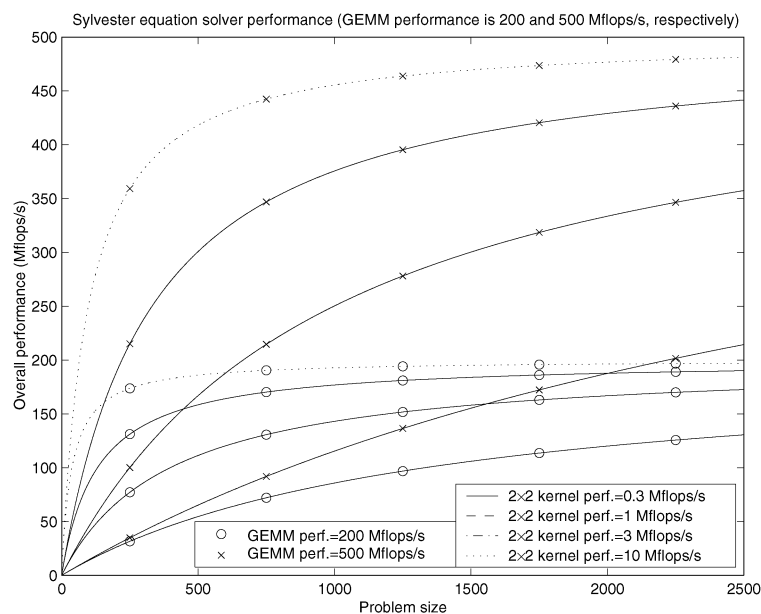


Fig. 1. Modeled overall performance of a recursive blocked implementation using different optimized kernels.

For $N = 250$, the number of flops performed in the kernel is 0.8% of the total number of flops, but with the faster DGEMM and the slowest kernel, the time spent in the kernel is more than 93% of the total time. With the same prerequisites and $N = 750$, the number of flops done in the kernel is less than 0.3% of the total number of flops, while the time spent in the kernel is more than 80% of the total time. We remark that similar analyses hold for all matrix equations considered. For the IBM Power3, 200 MHz, the 0.3 Mflops/s value is roughly the average performance of the LAPACK DTGSY2 routine [Kågström and Poromaa 1996a; Anderson et al. 1999] for these very small problems. DTGSY2 is the routine for solving small-sized 2×2 , 4×4 , and 8×8 generalized coupled Sylvester equations (3) used in an explicitly blocked level-3 solver [Kågström and Poromaa 1996a]. The main reason for the poor performance is that the kernel algorithms in LAPACK are designed with the primary goal of producing high accuracy results and signaling ill-conditioning. For example, the DTGSY2 solver is based on complete pivoting and overflow guarding. There is a trade-off among robustness, reliability, and speed. In the next section, we discuss ways of improving the speed in the kernel design, without sacrificing the accuracy and robustness demands, assuming the problems are not too ill-conditioned.

5.2 Design of Optimized Superscalar Kernels

The kernel routines in state-of-the-art libraries such as LAPACK [Anderson et al. 1999] and SLICOT [2001] are typically not optimized for superscalar RISC

processors. We continue our illustration with the LAPACK DTGSY2 kernel routine discussed above.

It uses DGETC2 to factorize the Z -matrix (see Z_{GCSY} in Section 2), and DGESC2 to do the forward and backward substitutions with overflow guarding. Several level-1, 2, and 3 BLAS routines (DAXPY, DSCAL, DGER, DGEMV, and DGEMM) are used in different updates. Now, DGETC2 and DGESC2 in turn use these and other BLAS routines to do the factorization and solving. Although this is good programming practice and leads to good code reuse, the overhead of the routine calls is devastating to the performance. Typically, the compiler cannot make good intraprocedural optimizations, since the BLAS kernels are linked from a library and cannot be candidates for inlining. Besides, a large part of the time spent in the BLAS routines is spent in setup code, that is, parameter checking and collection of machine characteristics, which in turn degrades the performance.

Our approach is to design one single routine, which does all (or most of) the computations for solving the kernel problems using the Kronecker product matrix representation introduced in Section 2.2. The construction of the Z -matrix (of size $MN \times MN$ for SYCT and LYCT and $2MN \times 2MN$ for GCSY), the LU factorization using partial pivoting with overflow guarding, and the forward and backward substitutions are all done in the same routine. This leads to a great potential for register reuse. Also, by complete loop unrolling, the routine can make much better use of the superscalar characteristics of today's processors.

5.3 Choice of Block Size and Termination of the Recursion

The choice of block sizes is an important issue in algorithms that use explicit blocking. For a multilevel explicitly blocked algorithm, each block size must be chosen to match a specific level in the memory hierarchy, for example, registers, level-1, or level-2 caches. This requires a deep knowledge of the architecture characteristics. On the other hand, for recursive blocked algorithms, we automatically obtain a hierarchical blocking which is variable and “squarish,” and there is only need for one blocking parameter, *blks*. Typically, *blks* is chosen so that a few blocks of current subproblems fit in the level-1 cache. For “problems” smaller than this size, recursion will not lead to any further speedup, to complete recursion until the element level would typically cause too much overhead and a drop in performance (e.g., see Gustavson et al. [1998b] and Gustavson and Jonsson [2000]).

For the **rtr*** solvers, the *blks* parameter can be chosen smaller than the size controlled by the level-1 cache without degrading the performance. This is due to the fact that there is enough computation at the leaves to hide the overhead caused by recursion. Although the superscalar kernel routines discussed in Section 5.2 are much faster than corresponding implementations in LAPACK and SLICOT, they are still much slower than the practical peak performance provided by an optimized GEMM kernel. The superscalar kernel solvers operate on problems of the size $2 \leq M, N \leq 4$, and this controls when the recursion will be terminated. So, we also use recursion in the **tr*** block solvers now

calling our superscalar GEMM kernel for all one-sided matrix–matrix updates. By doing so we minimize unnecessary overhead including costs for any buffer setups.

5.4 On the Choice of BLAS Implementation

Our recursive blocked algorithms reveal a great richness in matrix–matrix multiplication operations (GEMM). As the performance of the superscalar kernel increases, the GEMM routine plays a more important role for the overall performance (see Section 5.1).

The memory access pattern of the recursive blocked algorithms suggests that a BLAS implementation that uses recursion would perform very well together with the algorithms described here. However, this depends on the architecture characteristics. Recursive BLAS implementations have proven to be successful on processors with a large gap between the memory bandwidth and the computational power (e.g., see Gustavson et al. [1998a]). However, for machines with a balanced performance with respect to memory bandwidth and compute power, the difference is negligible when the nonrecursive BLAS is carefully tuned [Lindkvist 2000].

In Gustavson et al. [1998a], recursive blocked data formats that match recursive blocked algorithms are introduced and used in level-3 BLAS routines. Our experience is that using a recursive blocked data layout may lead to a significant performance gain, but it is not appropriate for our recursive triangular Sylvester-type matrix equation solvers in general. The main reason is the properties of the quasitriangular matrices involved, where the 2×2 bulges along the block diagonal correspond to conjugate pairs of complex eigenvalues. A 2×2 diagonal block may force the recursive splitting to occur one row or column above or below the partitioning determined by the blocked format. This leads to spill rows and columns with respect to the recursive blocked format, which further increase the data management overhead. Therefore, we have abandoned its use. However, if all eigenvalues are real or we work in complex arithmetic all matrices will be upper triangular and a recursive blocked data format is straightforward to apply.

5.5 Parallelization Issues

In our SMP implementations, different types of shared memory parallelism have been investigated. The first, and by far the easiest to implement, is to simply use level-3 BLAS routines that make use of more than one processor. Moreover, we can parallelize over independent tasks in the “Case 3” branch of the recursion of the algorithms (where both matrix dimensions are split simultaneously), which include both GEMM updates and matrix equation solves from successive splittings (see **rtrsyt** in Section 3.1 and **rtrgsy** in Section 3.3). Another way to go is to explicitly parallelize a standard blocked algorithm (e.g., see Kågström and Poromaa [1992] and Poromaa [1998, 1997]). In this section, these three ways of parallelization including hybrid variants are discussed using the coupled Sylvester equation (3) for illustration (see Section 3.3).

Implicit Data Parallelization. By linking to an SMP-aware implementation of BLAS, only the parallelism in the GEMM-based updates is taken into account. This gives fairly good speedup for SMP systems with a small number of processors, and, especially, for large problems, where the updates stand for a larger fraction of the total time to solve the problem. The main advantage of this approach is that it does not require any extra implementation work, granted that an SMP version of DGEMM is available. Moreover, this means that there are no requirements on the compiler or the run-time library to support any parallelization directives or routines, such as OpenMP or pthreads [OpenMP 2000; Nichols et al. 1996].

Task Parallelism in the Recursion Tree. In the **rtrgcsy** algorithm, there is also the possibility of parallelizing calls to the coupled Sylvester block kernel routine, by solving for the block matrix pairs $[X_{11}, Y_{11}]$ and $[X_{22}, Y_{22}]$ concurrently. Also, some of the updates are independent tasks that can be done in parallel. In our implementation of the **rtrgcsy** algorithm, we make use of an additional parameter *procs*, which holds the number of processors available to solve the current problem. If *procs* > 1, the problem size is large enough, and M, N fit into the Case 3 clause, then the second and the third recursive calls (solving for $[X_{11}, Y_{11}]$ and $[X_{22}, Y_{22}]$) are done in parallel, as well as some of the GEMM updates. The constraint on the problem size is necessary to prevent superfluous parallelization; that is, the cost of parallelizing the task (parallelization overhead) dominates the gain from executing it in parallel.

In our implementation, the same SMP implementation of BLAS as above is used. Although this is good for the largest updates at the top of the recursion, where the level of task parallelism is small, it is not optimal when there are many smaller updates to be done in parallel. Preferably, a DGEMM implementation where the number of available processors could be specified should be used. By using a DGEMM with a *proc* parameter, updates would alternate between making use of task parallelism (at the leaves of the recursion) and data parallelism (at the topmost roots of the recursion). However, we have not seen that option in today's BLAS implementations. Instead, with the current implementations there is a potential risk for over-parallelization; for example, on a four-processor SMP, there could be four instances of DGEMM running in parallel, each trying to make use of four processors.

Parallelization of Explicitly Blocked Algorithms. Distributed memory as well as shared memory blocked algorithms for solving Sylvester-type matrix equations have been studied in the literature (e.g., see Kågström and Poromaa [1992] and Poromaa [1998, 1997]). The above-cited algorithms try to utilize the maximum inherent parallelism of explicitly blocked algorithms for solving triangular matrix equations. For example, all blocks along each block diagonal of the solution(s) correspond to independent tasks (smaller triangular matrix equation solves) that can be executed in parallel.

A shared memory implementation of the triangular coupled Sylvester equation (3) is presented and discussed by Poromaa [1998, 1997]. The solution

(X, Y) overwrites the right-hand sides (C, F) , and the SMP algorithm solves for all blocks in each block diagonal of C and F in parallel, assuming enough numbers of processors are available. We have implemented this algorithm with different block (subsystem) solvers, including the sequential **rtrgsy** algorithm. The explicitly blocked algorithm also requires thread support from the compiler or a run-time library. One advantage of this algorithm is the potential inherent parallelism in solving subsystems along the block diagonals. One disadvantage is that the algorithm does not automatically generate good squarish calls to DGEMM. As such, the choice of block size is more architecture dependent than in a recursive blocked algorithm.

6. PERFORMANCE RESULTS OF RECURSIVE BLOCKED ALGORITHMS

The recursive blocked algorithms have been implemented in Fortran 90, using the facilities for recursive calls of subprograms, dynamic memory allocation, and threads. In this section, we present sample performance results of these implementations executing on IBM Power, MIPS, and Intel Pentium processor-based systems. We have selected a few processors representing different vendors and different architecture characteristics and memory hierarchies. The performance results (measured in Mflops/s) are computed using the flopcounts presented in Table I of Section 3.4. A selection of the results is displayed in graphs, where the performance of different variants of the recursive blocked implementations is visualized together with existing routines in the state-of-the-art libraries LAPACK [Anderson et al. 1999] and SLICOT [2001]. These results should be quite self-explanatory. Therefore, our discussion is restricted to significant or unexpected differences between the implementations executing on different computing platforms. For additional results we refer to Jonsson and Kågström [2001b]. The accuracy of the results computed by our recursive blocked algorithms is overall very good and similar to the accuracy obtained by the corresponding LAPACK and SLICOT routines. For benchmark problems see Kressner et al. [1999a,b] and SLICOT [2001].

The tests were run on several different machines. The first is the IBM RS/6000 SP SMP Thin Node, with four IBM PowerPC 604e CPUs running at 332 MHz. Each processor has a theoretical peak rate of 664 Mflops/s. The second is the IBM RS/6000 SP SMP Thin Node, with two IBM POWER3 CPUs running at 200 MHz. Each processor has a peak rate of 800 Mflops/s.

The third machine is one node of the HPC2N Linux cluster, with dual Intel Pentium III CPUs running at 550 MHz. Each processor has a theoretical peak rate of 550 Mflops/s. The fourth is one SGI Onyx2 node, with a MIPS R10000 CPU running at 195 MHz. The processor has a theoretical peak rate of 390 Mflops/s.

6.1 Standard Triangular Sylvester Equation

In Figure 2, we show performance graphs for different algorithms and implementations, executing on IBM PowerPC 604e and Intel Pentium III processor-based systems, for solving triangular Sylvester equations.

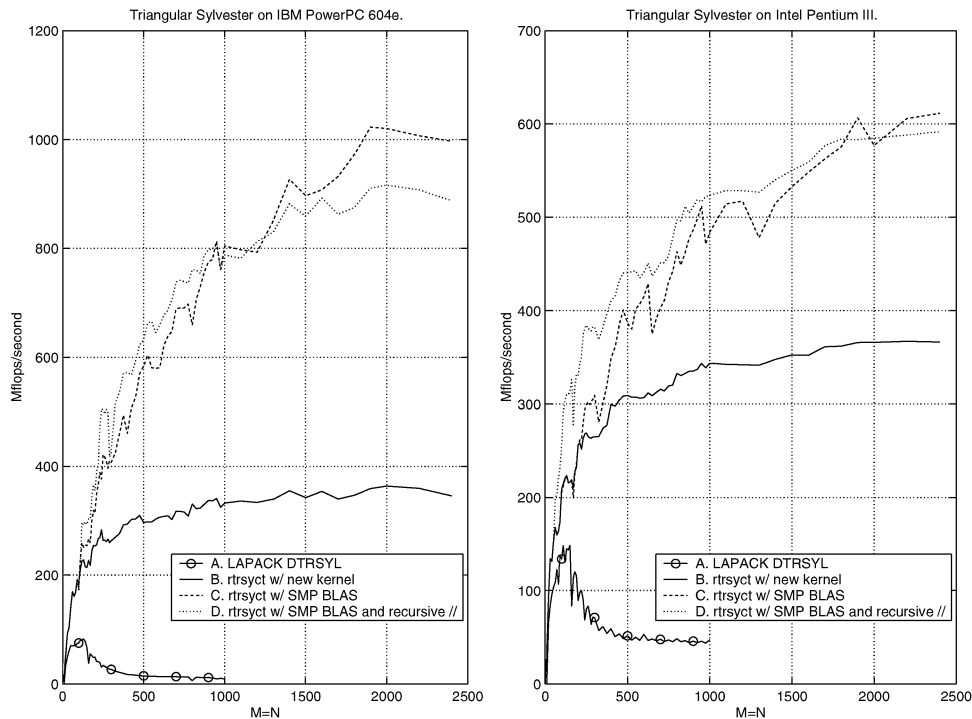


Fig. 2. Performance results for the triangular Sylvester equation ($M = N$): IBM SP PowerPC 604e (left) and Intel Pentium III (right).

The LAPACK DTRSYL implements an explicit Bartels–Stewart solver and is mainly a level-2 routine, which explains its poor performance behavior. Our recursive blocked **rtrsylt** shows between a 2-fold and a 35-fold speedup with respect to LAPACK DTRSYL and an additional speedup up to 2.8 on a four-processor PowerPC 604e node for large enough problems. The corresponding results on a two-processor Intel Pentium III show up to a 7-fold speedup with respect to LAPACK DTRSYL and an additional 1.6 speedup on two processors. The difference in the speedup B/A for the two architectures is mainly due to lower cache-miss penalties for the Pentium processor, which in turn leads to less degradation in the performance of LAPACK DTRSYL. The strong dependencies in the dataflow execution graph limit the possible parallel speedup for the triangular Sylvester equation.

6.2 Triangular Coupled Sylvester Equation

In Figure 3, performance graphs for different algorithms and implementations, executing on IBM PowerPC 604e and Intel Pentium III processor-based systems, for solving triangular generalized coupled Sylvester equations are displayed.

In total, we make comparisons among seven different implementations. Two of them are the LAPACK DTGSYL (A) and a parallel variant which uses SMP BLAS (F). Note that LAPACK DTGSYL is an explicitly blocked

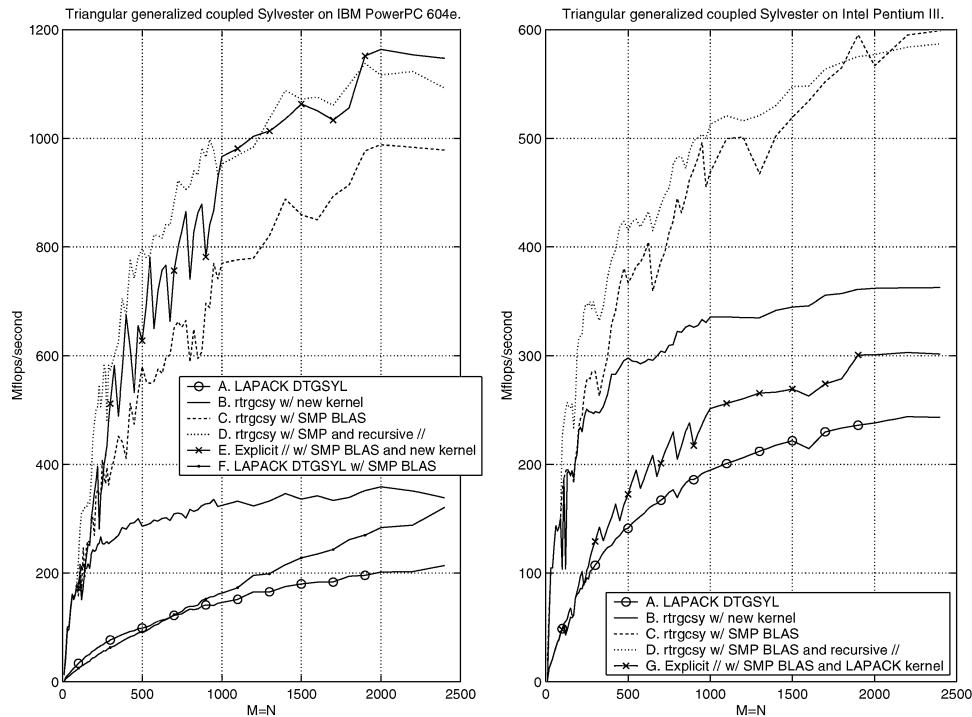


Fig. 3. Performance results for the generalized coupled Sylvester equation ($M = N$): IBM SP PowerPC 604e (left) and Intel Pentium III (right).

level-3 algorithm based on a generalization of the Bartels–Stewart algorithm [Kågström and Poromaa 1996a; Kågström and Westin 1989]. Three of them are the sequential recursive blocked **rtrgsy** (B) and two parallel variants (C and D). The last two (E and G) are new explicitly parallelized implementations of the blocked method [Kågström and Poromaa 1996a; Poromaa 1997, 1998].

Since the LAPACK DTGSYL is mainly a level-3 routine, its performance increases with increasing problem sizes and levels out due to only one level of blocking. But still **rtrgsy** shows over a 5-fold speedup with respect to LAPACK DTGSYL and an additional speedup up to 3.2 on a four-processor PowerPC 604e node for large enough problems. The corresponding results on a two-processor Intel Pentium III are up to 2.6-fold speedup with respect to LAPACK DTRSYL and additionally up to a 1.6-fold speedup on two processors. Comparisons have also been done with rectangular matrices. For the case $M = 10N$, the recursive implementations are twice as fast as LAPACK for problem size $(M, N) = (1000, 100)$. By using SMP versions, an extra speedup of 2.5 is achieved on the IBM PowerPC 604e. For larger problems, this gap is even wider [Jonsson and Kågström 2001a].

6.3 Impact of Solving an Unreduced Matrix Equation

For the continuous-time Lyapunov solver, we have investigated the impact of the choice of triangular matrix equation solver on the time to solve an

Table II. (a) Total Execution Time for Solving Unreduced Lyapunov Equation for Two Different Triangular Lyapunov Equation Solvers; (b) Estimate of Separation Sep[LYCT] Between A and $-A^T$

(a) N	SB03MD Using SB03MY		SB03MD Using rtrlyct		Speedup
	Total Time	Solver Part (%)	Total Time	Solver Part (%)	
50	0.0253	10.5	0.0238	8.8	1.06
100	0.141	9.2	0.134	4.7	1.05
250	1.65	11.1	1.52	3.1	1.09
500	26.8	43.7	15.3	1.9	1.75
750	105.8	48.5	54.4	1.6	1.94
1000	258.4	50.0	131.5	1.4	1.97

(b) N	SB03MD Using SB03MY		SB03MD Using rtrlyct		Speedup
	Total Time	Solver Part (%)	Total Time	Solver Part (%)	
50	0.0358	34.7	0.0338	26.0	1.06
100	0.205	34.5	0.164	18.5	1.25
250	2.37	35.7	1.74	13.0	1.36
500	64.9	76.8	16.5	8.4	3.94
750	278.2	80.6	59.1	8.7	4.71
1000	706.4	81.6	141.4	7.9	4.99

unreduced matrix equation. The reduction of an unreduced matrix equation to a triangular counterpart and the backtransformation of the solution are operations that are both at least as costly (measured in flops) as the triangular solve. Nevertheless, using our recursive blocked solvers can give a remarkable performance improvement. We use the SB03MD routine in the SLICOT [2001] library for illustration, which solves unreduced continuous-time Lyapunov equations. SB03MD also has an option to compute an estimate of the separation Sep[LYCT] between A and $-A^T$.

In Table II(a), timings for the SB03MD routine are displayed for problem sizes ranging from 50 to 1000 using two different triangular matrix equation solvers. These solvers are SB03MY provided in SLICOT, which implements a standard Bartels–Stewart method calling BLAS, and our recursive blocked **rtrlyct** algorithm. In the second column, the total times for solving the unreduced system with SB03MY as the triangular solver are displayed. This includes the time for the Schur factorization and backtransformation of the solution. In the fourth and fifth columns, similar results are displayed when SB03MY is replaced by the **rtrlyct** routine. We see between 75 to 100% speedup for problem sizes $N \geq 500$. We remark that the big difference of the timings for $N = 250$ and $N = 500$ and their impact are due to cache effects, which automatically are taken care of by our recursive blocked algorithms.

In Table II(b), we present timings for both solving an unreduced Lyapunov equation and computing a 1-norm-based estimate of the separation Sep[LYCT] between A and $-A^T$ (see Section 2.2). The condition estimation process leads to several calls to the triangular Lyapunov solver. Here, the impact is remarkable, with a four- to fivefold speedup for problem sizes $N \geq 500$.

7. CONCLUSIONS

The performance results verify that our recursive approach is very efficient for solving one-sided triangular Sylvester-type matrix equations on today's hierarchical memory computer systems. The recursion is terminated when the remaining subproblems to be solved are smaller than a given block size *blks*, which is the only architecture-dependent parameter in our algorithms. To complete recursion until the element level would in general cause too much overhead and a drop in performance. Our solution is to develop new high-performance superscalar kernels for small-sized triangular Sylvester-type matrix equations and lightweight GEMM operations, which implies that a larger part of the total execution time is spent in high-performance GEMM operations. We remark that for all architectures used in our testing, we could terminate the recursion with $blks = 4$ without degrading performance, which in turn means that the implementations are architecture-independent. Altogether, this leads to simpler and much faster algorithms for solving reduced as well as unreduced Sylvester-type matrix equations and different associated condition estimation problems. In our algorithms, the dimensions are split in two equal halves. Our approach and implementations allow sliding splittings, with the splitting points varying between the second to the penultimate row and/or column. By not splitting in the middle, the algorithms exhibit different memory access patterns and nonsquare updates, which in general degrade the performance. In the extreme cases we obtain the standard algorithms. Our intention is to make these algorithms available in the SLICOT [2001] library.

ACKNOWLEDGMENTS

We thank Fred Gustavson and our recursive pals in the Umeå HPC and Parallel Computing Research Group for stimulating and fruitful discussions. Finally, we thank Sven Hammarling and the referees for constructive comments on an earlier version of this manuscript.

REFERENCES

- ANDERSON, E., BAI, Z., DEMMEL, J., DONGARRA, J., DU CROZ, J., GREENBAUM, A., HAMMARLING, S., MCKENNEY, A., OSTROUCHOV, S., AND SORESENSEN, D. 1999. *LAPACK Users' Guide*, third ed. SIAM, Philadelphia.
- BAI, Z., DEMMEL, J., AND MCKENNEY, A. 1993. On computing condition numbers for the nonsymmetric eigenproblem. *ACM Trans. Math. Softw.* 19, 202–223.
- BARTELS, R. H. AND STEWART, G. W. 1972. Algorithm 432: Solution of the equation $AX + XB = C$. *Commun. ACM* 15, 9, 820–826.
- CHU, K.-W. E. 1987. The solution of the matrix equation $AXB - CXD = Y$ and $(YA - DZ, YC - BZ) = (E, F)$. *Linear Algebra Appl.* 93, 93–105.
- DACKLAND, K. AND KÄGSTRÖM, B. 1999. Blocked algorithms and software for reduction of a regular matrix pair to generalized Schur form. *ACM Trans. Math. Softw.* 25, 4, 425–454.
- DONGARRA, J. J., DU CROZ, J., DUFF I. S., AND HAMMARLING, S. 1990a. A set of level 3 basic linear algebra subprograms. *ACM Trans. Math. Softw.* 16, 1, 1–17.
- DONGARRA, J. J., DU CROZ, J., DUFF I. S., AND HAMMARLING, S. 1990b. Algorithm 679: A set of level 3 basic linear algebra subprograms. *ACM Trans. Math. Softw.* 16, 1, 18–28.
- ELMROTH, E. AND GUSTAVSON, F. 2001. High-performance library software for *QR* factorization. In *Applied Parallel Computing. New Paradigms for HPC Industry and Academia*, Lecture Notes in Computer Science, vol. 1947, Springer Verlag, New York, 53–63.

- GARDINER, J. D., LAUB, A. J., AMATO, J. J., AND MOLER, C. B. 1992a. Solution of the Sylvester matrix equation $AXB^T + CXD^T = E$. *ACM Trans. Math. Softw.* 18, 223–231.
- GARDINER, J. D., WETTE, M. R., LAUB, A. J., AMATO, J. J., AND MOLER, C. B. 1992b. A Fortran 77 software package for solving the Sylvester matrix equation $AXB^T + CXD^T = E$. *ACM Trans. Math. Softw.* 18, 232–238.
- GOLUB, G., NASH, S., AND VAN LOAN, C. 1979. A Hessenberg–Schur method for the matrix problem $AX + XB = C$. *IEEE Trans. Autom. Contr. AC-24*, 6, 909–913.
- GOLUB, G. AND VAN LOAN, C. 1996. *Matrix Computations*, third ed. Johns Hopkins University Press, Baltimore.
- GUSTAVSON, F. 1997. Recursion leads to automatic variable blocking for dense linear algebra. *IBM J. Res. Dev.* 41, 6 (Nov.), 737–755.
- GUSTAVSON, F., HENRIKSSON, A., JONSSON, I., KÅGSTRÖM, B., AND LING, P. 1998a. Recursive blocked data formats and BLAS's for dense linear algebra algorithms. In *Applied Parallel Computing. Large Scale Scientific and Industrial Problems*, Lecture Notes in Computer Science, vol. 1541, Springer-Verlag, New York, 195–206.
- GUSTAVSON, F., HENRIKSSON, A., JONSSON, I., KÅGSTRÖM, B., AND LING, P. 1998b. Superscalar GEMM-based level 3 BLAS—The on-going evolution of a portable and high-performance library. In *Applied Parallel Computing. Large Scale Scientific and Industrial Problems*, Lecture Notes in Computer Science, vol. 1541, Springer-Verlag, New York, 207–215.
- GUSTAVSON, F. AND JONSSON, I. 2000. Minimal-storage high-performance Cholesky factorization via blocking and recursion, *IBM J. Res. Dev.* 44, 6 (Nov.), 823–849.
- HAGER, W. W. 1984. Condition estimates. *SIAM J. Sci. Stat. Comp.* 5, 311–316.
- HAMMARLING, S. J. 1982. Numerical solution of the stable, non-negative definite Lyapunov equation. *IMA J. Numer. Anal.* 2, 303–323.
- HIGHAM, N. J. 1988. Fortran codes for estimating the one-norm of a real or complex matrix with applications to condition estimation. *ACM Trans. Math. Softw.* 14, 381–396.
- HIGHAM, N. J. 1993. Perturbation theory and backward error for $AX - XB = C$. *BIT* 33, 124–136.
- HIGHAM, N. J. 1996. *Accuracy and Stability of Numerical Algorithms*. SIAM, Philadelphia.
- JONSSON, I. AND KÅGSTRÖM, B. 2001a. Parallel triangular Sylvester-type matrix equation solvers for SMP systems using recursive blocking. In *Applied Parallel Computing. New Paradigms for HPC Industry and Academia*, Lecture Notes in Computer Science, vol. 1947, Springer Verlag, New York, 64–74.
- JONSSON, I. AND KÅGSTRÖM, B. 2001b. Recursive blocked algorithms for solving triangular matrix equations—Part I: One-sided and coupled Sylvester-type equations. *SLICOT Working Note* 2001-4. Dept. of Computing Science, Umeå University, SE-901 87 Umeå, Sweden.
- JONSSON, I. AND KÅGSTRÖM, B. 2002. Recursive blocked algorithms for solving triangular systems—Part II: Two-sided and generalized Sylvester and Lyapunov matrix equations. *ACM Trans. Math. Softw.* 28, 4 (Dec.).
- KÅGSTRÖM, B. 1977a. Bounds and perturbation bounds for the matrix exponential. *BIT* 17, 39–57.
- KÅGSTRÖM, B. 1977b. Numerical computation of matrix functions. Tech. Rep. UMINF-58.77, Institute of Information Processing, Umeå University, S-901 87 Umeå, Sweden.
- KÅGSTRÖM, B. 1993. 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*, Kluwer Academic, Amsterdam, 195–218.
- KÅGSTRÖM, B. 1994. A perturbation analysis of the generalized Sylvester equation $(AR - LB, DR - LE) = (C, F)$. *SIAM J. Matrix Anal. Appl.* 15, 4, 1045–1060.
- KÅGSTRÖM, B., LING, P., AND VAN LOAN, C. 1998a. GEMM-based level 3 BLAS: High-performance model implementations and performance evaluation benchmark. *ACM Trans. Math. Softw.* 24, 3, 268–302.
- KÅGSTRÖM, B., LING, P., AND VAN LOAN, C. 1998b. GEMM-based level 3 BLAS: Portability and optimization issues. *ACM Trans. Math. Softw.* 24, 3, 303–316.
- KÅGSTRÖM, B. AND POROMAA, P. 1992. Distributed and shared memory block algorithms for the triangular Sylvester equation with sep^{-1} estimator. *SIAM J. Matrix Anal. Appl.* 13, 1, 90–101.
- KÅGSTRÖM, B. AND POROMAA, P. 1996a. LAPACK-style algorithms and software for solving the generalized Sylvester equation and estimating the separation between regular matrix pairs. *ACM Trans. Math. Softw.* 22, 1, 78–103.

- KÅGSTRÖM, B. AND POROMAA, P. 1996b. Computing eigenspaces with specified eigenvalues of a regular matrix pair (A, B) and condition estimation. *Numer. Alg.* 12, 369–407.
- KÅGSTRÖM, B. AND VAN DOOREN, P. 1992. A generalized state-space approach for the additive decomposition of a transfer matrix. *Int. J. Numer. Linear Algebra Appl.* 1, 165–181.
- KÅGSTRÖM, B., AND WESTIN, L. 1989. Generalized Schur methods with condition estimators for solving the generalized Sylvester equation. *IEEE Trans. Autom. Contr.* 34, 4, 745–751.
- KRESSNER, D., MEHRMANN, V., AND PENZL, T. 1999a. CTLEX—A collection of benchmark examples for continuous-time Lyapunov equations. *SLICOT Working Note* 1999-6.
- KRESSNER, D., MEHRMANN, V., AND PENZL, T. 1999b. DTLEX—A collection of benchmark examples for discrete-time Lyapunov equations. *SLICOT Working Note* 1999-7.
- LINDKVIST, A. 2000. High-performance recursive BLAS kernels using new data formats for the QR factorization. Master's Thesis, Rep. UMNAD-325.00, Department of Computing Science, Umeå University, SE-901 87 Umeå, Sweden.
- MOLER, C. B. AND STEWART, G. W. 1973. An algorithm for generalized matrix eigenvalue problems. *SIAM J. Numer. Anal.* 10, 241–256.
- MOLER, C. B. AND VAN LOAN, C. F. 1978. Nineteen dubious ways to compute the exponential of a matrix. *SIAM Rev.* 20, 801–836.
- NICHOLS, B., BUTTLAR, D., PROULX FARRELL J., AND FARRELL, J. 1996. *Pthreads Programming: A POSIX Standard for Better Multiprocessing*. O'Reilly, Sebastopol, Calif.
- OPENMP 2000. Fortran Application Program Interface, Version 2.0, November, www.openmp.org/specs/.
- PENZL, T. 1998. Numerical solution of generalized Lyapunov equations. *Adv. Comp. Math.* 8, 33–48.
- POROMAA, P. 1997. High performance computing: Algorithms and library software for Sylvester equations and certain eigenvalue problems with applications in condition estimation. PhD Thesis UMINF-97.16, Department of Computing Science, Umeå University, SE-901 87 Umeå, Sweden, June.
- POROMAA, P. 1998. Parallel algorithms for triangular sylvester equations: Design, scheduling and scalability issues. In *Applied Parallel Computing, Large Scale Scientific and Industrial Problems*, Lecture Notes in Computer Science, vol. 1541, Springer-Verlag, New York, 438–446.
- SLICOT 2001. Library and the numerics in control network (NICONET) Web site: www.win.tue.nl/niconet/index.html, Academic Press, San Diego, Calif.
- STEWART, G. W. AND SUN, J.-G. 1990. *Matrix Perturbation Theory*. Academic, San Diego.
- TOLEDO, S. 1997. Locality of reference in LU decomposition with partial pivoting. *SIAM J. Matrix Anal. Appl.* 18, 4, 1065–1081.
- VAN LOAN, C. F. 1977. The sensitivity of the matrix exponential. *SIAM J. Numer. Anal.* 14, 971–981.

Received April 2001; revised June 2002; accepted July 2002