

Stability of Block Algorithms with Fast Level 3 BLAS

James W. Demmel * Nicholas J. Higham †

July 27, 1990

Abstract

Block algorithms are becoming increasingly popular in matrix computations. Since their basic unit of data is a submatrix rather than a scalar they have a higher level of granularity than point algorithms, and this makes them well-suited to high-performance computers. The numerical stability of the block algorithms in the new linear algebra program library LAPACK is investigated here. It is shown that these algorithms have backward error analyses in which the backward error bounds are commensurate with the error bounds for the underlying level 3 BLAS (BLAS3). One implication is that the block algorithms are as stable as the corresponding point algorithms when conventional BLAS3 are used. A second implication is that the use of BLAS3 based on fast matrix multiplication techniques affects the stability only insofar as it increases the constant terms in the normwise backward error bounds. For linear equation solvers employing LU factorization it is shown that fixed precision iterative refinement helps to mitigate the effect of the larger error constants. The analysis is illustrated with the aid of numerical examples.

*Computer Science Division and Mathematics Department, University of California, Berkeley, CA 94720, U.S.A. (na.demmel@na-net.stanford.edu). This author acknowledges the financial support of the National Science Foundation via grants DCR-8552474 and ASC-8715728. He is also a Presidential Young Investigator.

†Department of Mathematics, University of Manchester, Manchester, M13 9PL, UK. (na.nhigham@na-net.stanford.edu).

Key words: block algorithm, LAPACK, level 3 BLAS, iterative refinement, LU factorization, QR factorization, backward error analysis.

AMS(MOS) subject classifications. primary 65F05, 65F25, 65G05.

1 Introduction

A *block algorithm* in matrix computations is one that is defined in terms of operations on submatrices rather than matrix elements. Such algorithms are well suited to high-performance computers because their data locality properties allow for efficient usage of memory hierarchy [16, 51].

When a block algorithm is coded in Fortran, advantage can be taken of the Basic Linear Algebra Subprograms (BLAS3). The BLAS3 are a set of routines for various types of matrix multiplication, together with the system with multiple right-hand sides usually coded block algorithm. In the bulk of the computation is carried out by calls to the BLAS3.

The BLAS3 specification is not complete the input, output and call sequence for each routine, but allow freedom of implementation, subject to the requirement that the routines be numerically stable. This freedom includes not only the order of matrix multiplication, but the use of algorithms algebraically equivalent to the conventional ones. Of chief interest here are algorithms that achieve a favourable operation count (for suitable dimensions) through the use of the multiplication technique. We will refer to such BLAS3 implementations as "fast BLAS3".

One set of fast BLAS3 is proposed in [18]. It is shown how asymptotic speedups can be produced in all the BLAS3 routines by the use of Strassen's method for matrix multiplication which forms the product of two $n \times n$ matrices in $O(n^{2.7})$ operations [30, 807]. A set of fast BLAS3 could also be built using Winograd's matrix multiplication (which has a [34] operation³ count of $O(n^{2.81})$ with different constants than the conventional technique) or one of the methods with a lower exponent than Strassen's method (although the practicality of these methods has yet to be demonstrated [22]). The use of complex matrices allows all these possibilities can be combined with the technique which has been analysed in [21].

product of two complex matrices to be formed using only three real matrix multiplications. Several researchers are experimenting with the use of fast **BLAS** in linear equation solvers. In particular, we mention the work of Bailey, Lee and Sun [3], who use Strassen's method for the matrix multiplications arising in the LAPACK LU factorization routine `SGERF`.

Our purpose in this work is to investigate the numerical stability of block algorithms that employ fast **BLAS**. We restrict our attention mainly to the block algorithms used in LAPACK [4, 9]. For a blocksize of 1, the algorithms in LAPACK are classical point algorithms that are well-known to be numerically stable, that is, each computed answer is the exact answer to a perturbed problem where the norm of the perturbation is bounded by the product of the unit roundoff, a modest constant depending on the dimensions, and the norm of the data. (To be precise, this statement is true modulo the possibility of a large growth factor in Gaussian elimination with partial pivoting and a weaker definition of stability for matrix inversion.) For block sizes $r > 1$, with conventional **BLAS**, it is generally accepted that the same stability results are valid, although we are not aware of any detailed proofs (in the case of block LU factorization one can argue that the same arithmetic operations are carried out as for $r = 1$, albeit in a different order). The question of particular interest here is the effect on the stability of using *fast* **BLAS** when $r > 1$. We will show that, for all **BLAS** implementations of interest, backward error bounds hold for the block algorithms that are commensurate with the error bounds for the **BLAS** themselves. This is clearly the best we could expect to prove.

As our model for floating point arithmetic we take

$$\begin{aligned} fl(x \pm y) &= (x(1+\alpha) \pm y(1+\beta)), & |\alpha|, |\beta| &\leq u, \\ fl(x \oslash y) &= (x \oslash y)(1+\delta), & |\delta| &\leq u, \quad \oslash = *, /, \end{aligned}$$

where u is the unit roundoff. We need to make some assumptions about the stability of the **BLAS**. The **BLAS** primitives involve three types of matrix multiplication: a general product AB , a cross product $A^T A$, and the product of a triangular matrix with a full matrix. It is sufficient to assume that all these products satisfy the following general condition: if $A \in \mathbb{R}^{m \times n}$, $B \in \mathbb{R}^{n \times p}$ and \hat{C} is the computed approximation to $C = AB$, then

$$\hat{C} = AB + \Delta C, \quad \|\Delta C\| \leq c_1(m, n, p)u \|A\| \|B\| + O^2(u) \quad (1)$$

where $c_1(m, n, p)$ denotes a constant depending on m, n and p . Here, the matrix norm is defined by

$$\|X\| = \max_{i,j} |x_{ij}|.$$

Note that for this norm with A and B dimensioned as above, $\|AB\| \leq n \|A\| \|B\|$ is the best such inequality.

Also assume that the computed solution to the triangular system $TX = B$, where $T \in \mathbb{R}^{m \times m}$ and $B \in \mathbb{R}^{m \times p}$, satisfies

$$T\widehat{X} = B + \Delta B, \quad \|\Delta B\| \leq c_2(m, p)u \|T\| \|\widehat{X}\| + O(u^2). \quad (2)$$

For conventional **BLAS** implementation conditions (1) and (2) hold with the $c_1(m, n, p) = n^2$ and $c_2(m, p) = m(m+1)$.

For the fast **BLAS** proposed in [18], based on Strassen's method (1) and (2) hold with c_1 and c_2 rather complicated functions of the dimensions m, n, p and the threshold n_0 that determines the level of recursion. In the special case $m = n = p = 2^k$, $n_0 = 2^l$, the constants reduce to [18]

$$c_1(n, n, n) = \left(\frac{n}{n_0}\right)^{\log_2 12} (n_0^2 + 5n_0) - 5n,$$

$$c_2(n, n) = \left(\frac{n}{n_0}\right)^{\log_2 12} \left(\frac{n_0^2}{11} + \frac{23}{55}n_0\right) + \frac{10}{11}n_0^2 + \frac{35}{11}n_0 - \frac{143}{55}n.$$

Condition (1) also holds when the multiplication is done by Wograd's method with scaling α , in the case of complex matrices, by the method of [21] combined with any method for real matrix multiplication that satisfies (1) (see the error analysis in [6] and [21]).

Note that for conventional multiplication we have the following componentwise version of (1):

$$\widehat{C} = AB + \Delta C, \quad |C| \leq n u |A| |B| + O(u^2). \quad (3)$$

Similarly for the substitution algorithm for solving triangular systems we have

$$T\widehat{X} = B + \Delta B, \quad |\Delta B| \leq (m+1)u |T| |\widehat{X}| + O(u^2). \quad (4)$$

Stress that these bounds are not stronger than (1) and (2). For example, if $D = \text{diag}(d_i)$ with $d_i > 0$ for all i , then scaling $AB \rightarrow AD \cdot D^{-1}B$ leaves C and the upper bound in (3) unchanged, and scaling $AB \rightarrow DA \cdot BD$ causes the bound of (3)

to scale in the same way as C ; (1) does not share these favorable scaling properties.

For further remarks on the differences between (1) and (3) see [18]. A consequence of

(3) and (4) is that for some block algorithms it is possible to obtain stronger backward

error results than the usual ones (perhaps for certain classes of matrices only);

for examples see [8] and [19]. These stronger results are usually not valid for any of

the fast BAS discussed above.

The block algorithms in LAPACK break into two main classes: those based on

LU factorization and those involving orthogonal transformations. In the next section

we give an error analysis of block LU factorization. We show that iterative refine-

ment in fixed precision is beneficial for all BAS implementations and point out the

instability of a particular form of LU factorization with block triangular factors. We

also explain how to investigate the stability of a block algorithm without doing a full

error analysis. In section 3 we consider the use of aggregated Householder transfor-

mations, which form the basis of a variety of block algorithms involving orthogonal

transformations. Some concluding remarks are given in section 4.

2 LU Factorization

2.1 Error Analysis

In this section we examine in detail the stability of block LU factorization. Initially

we assume that no pivoting is used and that the factorization succeeds; below we

discuss the addition of pivoting.

Consider a block implementation of the outer product form of LU factorization

[16, p.9], [17, p.10]. The algorithm may be described through the partitioning

$$\begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} L_{11} & 0 \\ L_{21} & I_{n-r} \end{bmatrix} \begin{bmatrix} I_r & 0 \\ 0 & B \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & I_{n-r} \end{bmatrix} \in \mathbb{R}^{n \times n},$$

where A_{11} is $r \times r$. One step of the block algorithm consists of factoring $A_{11} = L_{11}U_{11}$,

solving the multiple right-hand side triangular systems $L_{11}U_{12} = A_{12}$ and $L_{21}U_{11} =$

A_{21} for U_{12} and L_{21} respectively, and forming $B = A_{22} - L_{21}U_{12}$; this procedure is

then repeated on B . The block operations defining U_{11} , L_{12} , L_{21} and B are level 3 BAS

operations.

We will assume that the block level LU factorization is done in such a way that the computed LU factors of $A_{11} \in \mathbb{R}^{r \times r}$ satisfy

$$\widehat{L}_{11}\widehat{U}_{11} = A_{11} + \Delta_{11}, \quad \|\Delta_{11}\| \leq c_3(r)u \|\widehat{L}_{11}\| \|\widehat{U}_{11}\| + O(u^2). \quad (5)$$

We claim that under these assumptions, together with (1) and (2), the LU factors of $A \in \mathbb{R}^{n \times n}$ computed using a block size r satisfy

$$\widehat{L}\widehat{U} = A + \Delta, \quad \|\Delta\| \leq u (\delta(n, r)\|A\| + \theta(n, r)\|\widehat{L}\|\|\widehat{U}\|) + O(u^2), \quad (6)$$

where $\delta(n, r)$ and $\theta(n, r)$ are constants depending on n and r . The proof is essentially inductive. For $n = r$, (6) holds with

$$\delta(r, r) = 0, \quad \theta(r, r) = 0, \quad (7)$$

in view of (5). Consider the first block stage of the factorization. The assumptions imply that

$$\widehat{L}_{11}\widehat{U}_{12} = A_{12} + \Delta_{12}, \quad \|\Delta_{12}\| \leq c_2(r, n-r)u \|\widehat{L}_{11}\| \|\widehat{U}_{12}\| + O(u^2), \quad (8)$$

$$\widehat{L}_{21}\widehat{U}_{11} = A_{21} + \Delta_{21}, \quad \|\Delta_{21}\| \leq c_2(r, n-r)u \|\widehat{L}_{21}\| \|\widehat{U}_{11}\| + O(u^2). \quad (9)$$

To obtain $B = A_{22} - L_{21}U_{12}$ we first compute $C = \widehat{L}_{21}\widehat{U}_{12}$, obtaining

$$\widehat{C} = \widehat{L}_{21}\widehat{U}_{12} + \Delta, \quad \|\Delta\| \leq c_1(n-r, r, n-r)u \|\widehat{L}_{21}\| \|\widehat{U}_{12}\| + O(u^2),$$

and then subtract from A_{22} , obtaining

$$\widehat{B} = A_{22} - \widehat{C} + F, \quad \|F\| \leq u (\|A_{22}\| + \|\widehat{C}\|) + O(u^2). \quad (10)$$

It follows that

$$\begin{aligned} \widehat{B} &= A_{22} - \widehat{L}_{21}\widehat{U}_{12} + \Delta, \\ \|\Delta\| &\leq u (\|A_{22}\| + \|\widehat{L}_{21}\| \|\widehat{U}_{12}\| + c_1(n-r, r, n-r)u \|\widehat{L}_{21}\| \|\widehat{U}_{12}\|) + O(u^2). \end{aligned} \quad (11)$$

The remainder of the algorithm consists of the completion of the LU factorization of \widehat{B} , and by our inductive assumption (6) the computed LU factors satisfy

$$\begin{aligned} \widehat{L}_{22}\widehat{U}_{22} &= \widehat{B} + \Delta \widehat{B}, \\ \|\Delta \widehat{B}\| &\leq \delta(n-r, r)\|\widehat{B}\| + \theta(n-r, r)u \|\widehat{L}_{22}\| \|\widehat{U}_{22}\| + O(u^2). \end{aligned} \quad (12)$$

Combining (11) and (12), and bounding $\|\widehat{B}\|$ using (10), we obtain

$$\begin{aligned} \widehat{L}_{21}\widehat{U}_{12} + \widehat{L}_{22}\widehat{U}_{22} &= A_{22} + \Delta_{22} \\ \|\Delta_{22}\| &\leq u \left([1 + \delta(n-r, r)] \|A_{22}\| + [1 + c_1(r, n-r, n-r) + \delta(n-r, \widehat{L}_{21})] \|\widehat{U}_{12}\| \right. \\ &\quad \left. + \theta(n-r, r) \|\widehat{U}_{22}\| \right) + O(u^2). \end{aligned} \quad (13)$$

Collecting together (5), (8), (9) and (13) we have

$$\widehat{L}\widehat{U} = A + \Delta, \quad (14)$$

where bounds on $\|\Delta_{ij}\|$ are given in the equations just mentioned. These bounds for the blocks of Δ can be weakened slightly and expressed together in the more succinct form

$$\|\Delta\| \leq u \left(\delta(n, r) \|A\| + \theta(n, \widehat{L}) \|\widehat{U}\| \right) + O(u^2) \quad (15)$$

where

$$\begin{aligned} \delta(n, r) &\equiv 1 + \delta(n-r, r), \\ \theta(n, r) &\equiv \max \{ c_3(r), \alpha(r, n-r), 1 + (c_1(r, n-r, n-r) + \delta(n-r, r) + \theta(n-r, r)) \} \end{aligned}$$

Using (7) it follows that $\delta(n, r) \leq n/r$.

These recurrences show that the basic error constants in assumptions (1), (2) and (5) combine additively at worst. Thus, the backward error analysis for the LU factorization is commensurate with the error analysis for the particular implementation of the BAS employed in the block factorization. In the case of the conventional BAS we obtain a generalization of the classical Wilkinson result for $r=1$, with $\theta(n, r) = O(\eta)$.

Although the above analysis is phrased in terms of the block outer product form of LU factorization, the same result holds for other “ i, j, k ” block forms (with slightly different constants), for example the gary or sdt form.

If we incorporate partial pivoting in the above factorization then two of the block steps are coalesced. L_{11} and L_{21} are obtained by using Gaussian elimination with partial pivoting (GEP) to compute the factorization

$$P_1 \begin{bmatrix} A_{11} \\ A_{21} \end{bmatrix} = \begin{bmatrix} L_{11} \\ L_{21} \end{bmatrix} U_{11}. \quad (16)$$

This each step of the algorithm involves two, rather than three, **BLAS** operations. If the obvious analogue of (5) holds for (16) then (14) (with A replaced by $P A$) and (15) remain valid, with minor changes in the recurrences for $\delta(n, r)$ and $\theta(n, r)$.

There is no difficulty in extending the analysis to cover solution of $Ax = b$ using the computed LU factorization. Including the usual error analysis for substitution (see [17, sec.3.1], for example) we find that $(A + \Delta)\hat{x} = b$, with

$$\|\Delta\| \leq u (\delta(n, r)\|A\| + (\theta(n, r) + 2)\|\hat{L}\|\|\hat{U}\|) + O(u^2). \quad (17)$$

For a linear system with multiple right-hand sides, $AX = B$ with $X, B \in \mathbb{R}^{n \times p}$, both substitution stages are **BLAS** operations. Using (2) it is straightforward to show that the computed \hat{X} satisfies

$$\|A\hat{X} - B\| \leq \{ [c_2(n, p)(n+n) + \theta(n, r)n] \|\hat{L}\|\|\hat{U}\| + \delta(n, r)n \|A\| \|\hat{X}\| \|u\} + O(u^2) \quad (18)$$

Notice that the error analysis given in this section adapts in a straightforward way to block LU factorization for banded matrices, block factorization of symmetric indefinite matrices [17, p.168], and block Cholesky factorization of (banded) symmetric positive definite matrices [24].

2.2 Iterative Refinement

The LAPACK routines for solving linear equations support fixed precision iterative refinement [10], that is, iterative refinement in which no extra precision is used in calculating the residuals. The benefits of this process can be explained in terms of the componentwise relative backward error $\omega(y)$ of an approximate solution y to $Ax = b$. This quantity is defined by

$$\begin{aligned} \omega(y) &\equiv \min \{ \epsilon : (A + \Delta)y = b + \Delta, \quad \|\Delta\| \leq \epsilon \|A\|, \quad \|\Delta\| \leq \epsilon \|b\| \}, \\ &= \max_i \frac{|b - Ay_i|}{(|A| |y| + |b|)} \end{aligned} \quad (19)$$

where the latter equality is proved in [25]. A small value for $\omega(y)$ implies that y is the solution of a system in which each element of A and b has undergone a small relative perturbation; in particular, zero elements are not perturbed. However, in general, all that can be guaranteed for the \hat{x} from GEP is that the normwise relative backward

error $\eta(\hat{x}) = O(u)$, where

$$\begin{aligned} \eta(y) &\equiv \min \{ \epsilon : (A + \Delta)y = b + \mathbf{A}, \quad \|\Delta\|_\infty \leq \epsilon \|A\|_\infty, \quad \|\mathbf{A}\|_\infty \leq \epsilon \|b\|_\infty \} \\ &= \frac{\|r\|_\infty}{\|A\|_\infty \|y\|_\infty + \|b\|_\infty}. \end{aligned} \quad (20)$$

Steel [29] showed that as long as A is not too ill-conditioned and the components of the vector $\|A\| \|x\|$ do not vary too much in magnitude, then $\omega(y) = O(u)$ for the vector y obtained from GEP with one step of fixed precision iterative refinement. Steel's result, together with further analysis in [1], provides the theoretical foundation for the inclusion of fixed precision iterative refinement in LAPACK.

In LAPACK iterative refinement is terminated if

1. $\omega \leq u$,
2. ω has not decreased by a factor of at least 2 from the previous iteration, or
3. five iterations have been performed.

These criteria have been chosen to be robust in the face of different BLAS implementations and machine arithmetics.

Steel's result is applicable only when conventional BLAS are used. To investigate the effect of using fast BLAS, we can make use of work in [23] that covers fixed precision iterative refinement with an arbitrary linear equations solver. For our purposes, the results in [23] require a bound of the form $\|b - A\hat{x}\| \leq u \|G\| \|\hat{x}\|$ for the given solver, where G is a nonnegative matrix. For block LU factorization with partial pivoting we have

$$\|b - A\hat{x}\| \leq \|\Delta\| \|\hat{x}\|,$$

where $\|\Delta\|$ is bounded in (17). We assume that the residual for the refinement step is computed in the conventional way via inner products or saxpy operations, as in LAPACK. Since we have only a loose bound on Δ we cannot apply a direct generalization in [23] of Steel's result (to do so we would need to have $\|\Delta\| \leq u \|G\| \|A\|$ with $\|G\|$ bounded independently of A). However, we can invoke the weaker Theorem 2.1 in [23] to obtain

$$\|b - A\hat{x}\| \leq (n+2)u (\|A\| \|\hat{x}\| + \|b\|) + O(u^2), \quad (21)$$

where \hat{y} is the computed vector obtained after one step of fixed precision iterative refinement. This result has two main features. First, it is asymptotic, and the second order term prevents us concluding from (21) and (19) that $\omega(\hat{y}) \leq (n+2)u$. However, if the components of $\|A\|^{-1}(\hat{y} + b)$ do not vary too much in magnitude it is likely that this inequality will be satisfied (if not, extra refinement steps may help to achieve a small ω). The second point is that Δ does not appear in the first order term of (21) — it is hidden in the second order term where it multiplies a vector with elements of $O(u)$. This means that the refinement step tends to suppress any instability manifested in Δ .

LAPACK also supports iterative refinement for linear systems with multiple right-hand sides, $AX = B$ where $X, B \in \mathbb{R}^{n \times p}$. In this case it is appropriate to consider whether a small componentwise relative backward error is achieved for each individual system $Ax_i = b_i, i = 1:p$. If conventional BAS are used for the computation of \hat{X} and for the refinement process then Sed's result is applicable to each system $Ax_i = b_i$.

Suppose that fast BAS are used in computing \hat{X} . First, we will obtain a bound of the form $\|b_i - A\hat{x}_i\| \leq u G_i \|\hat{x}_i\|$ for each i . It is necessary to do this in an indirect way as follows. Write first that (18) implies

$$\|A\hat{x}_i - b_i\|_\infty \leq M \mu_i u \|A\|_\infty \|\hat{x}_i\|_\infty + O(u^2), \quad i = 1:p,$$

where

$$M = [c_2(n, p)(n+n) + \theta(n, r)n] \frac{\|\hat{L}\| \|\hat{U}\|}{\|A\|_\infty} + \delta(n, r)n,$$

$$\mu_i = \frac{\|\hat{X}\|}{\|\hat{x}_i\|_\infty} \geq 1.$$

From (20) it follows that

$$(A + \Delta)_i \hat{x}_i = b_i, \quad \|\Delta_i\|_\infty \leq M \mu_i u \|A\|_\infty + O(u^2). \quad (22)$$

Here we have $\|b_i - A\hat{x}_i\| \leq \|\Delta_i\| \|\hat{x}_i\|$, with $\|A\|_\infty$ bounded as in (22). In finding Theorem 2.1 of [23] we need to specify how the residuals $R = B - AX$ are computed and how the corrections for the refinement are computed. We will assume that the residuals are computed using conventional BAS. If fast BAS are used we can do

no better than to obtain a *normwise* bound for each $b_i - A \hat{y}_i$ that is proportional to the fast BAS error constant c_1 (this is not surprising since it is a general principle for iterative techniques that the stability or accuracy is limited by the quality of the computed residuals).

If we use fast BAS for the substitutions on the correction step then for each computed correction we have a result analogous to (22). Theorem 2.1 of [23] then shows that (21) holds for b_i and \hat{y}_i , but the potentially very large μ_i term and its analogue for the refinement step are present in the second order term of (21), making the result of limited practical value. If the substitutions are done using conventional BAS then (21) holds for b_i and \hat{y}_i with μ_i alone present in the second order term here we would expect a small $\omega(\hat{y}_i)$ as long as μ_i is not too large. Finally, we note that if all the substitutions in the computation of \hat{X} and in the refinement process are done with conventional BAS then we can set $\mu_i \equiv 1$ in the above analysis, and we will obtain the same computed results as if the refinement was carried out on each system $Ax_i = b_i$ independently.

Our overall conclusion is that fixed precision iterative refinement can be beneficial for GEP with fast BAS in two ways, assuming that residuals are computed using conventional BAS. First, it may lead to a competitive relative backward error of order u , although the theoretical backing is weaker than when conventional BAS are used. Second, the refinement will, in any case, tend to counteract any "wild instability" induced by the potentially faster growth of errors in the fast BAS.

We present some numerical results for illustration. Our experiments were performed in MATLAB for which $u \approx 2.2 \times 10^{-16}$. We solved $Ax = b$ by block outer product LU factorization with partial pivoting using both conventional BAS and a fast BAS; the latter uses conventional triangular solves and does matrix multiplication by Strassen's method with $n_0 = 1$ (recursion down to the scalar level). Iterative refinement was applied with the convergence test $\omega \leq u$.

We give detailed results for three matrices taken from the test collection [20]: pascal (n) is a symmetric positive definite matrix constructed from the elements of Pascal's triangle; triw (n, α) is upper triangular with 1s on the diagonal and every entry in the upper triangle equal to α ; and ipjfact ($n, 1$) is the symmetric positive definite matrix with (i, j) entry $1/(i+j)!$. In each case b was chosen randomly with elements from the uniform distribution on $[0, 1]$. Tables 1-3 show the competitive

relative backward errors ω for the iterates, with the normwise relative backward errors η in parentheses; the column leading “Conv.” denotes conventional BAS. We also report the condition numbers $\kappa_\infty(A) = \|A\|_\infty \|A^{-1}\|_\infty$ and $\kappa_2(A) = \| |A^{-1}| |A| \|_F$.

In these three specially chosen examples all the η values for the original computed solution are less than u , but the ω values for the fast BAS are substantially larger than for the conventional BAS. Note how iterative refinement reduces the ω values below u in one step in the first three examples. More typical, less extreme behavior is illustrated in Table 4, where $A(n)$ is a random matrix with elements from the uniform $[0, 1]$ distribution.

The matrices in the first three examples each have nonnegative elements of widely varying magnitude. These matrices were tried because it is known that Stassen’s method can provide poor relative accuracy when forming the product of such matrices in floating-point arithmetic [18]. It is interesting to recall the result that the computed solution to $Ax = b$ obtained by GEP in floating-point arithmetic is invariant under row and column scalings by powers of the machine base, as long as the same pivot sequence is chosen [7, p. 18]. This invariance property does not hold when Stassen’s method is used in the BAS—this observation provides some further insight into the first three examples.

Table 1: $A = \text{pascal}$ (8)

$$\kappa_{\infty}(A) = 3.96 \times 10^{-7}, \quad \text{cond}(A) = 4.60 \times 10^6$$

$\text{Giv}, r = 1$	$\text{Giv}, r = 2$	$\text{Fast}, r = 2$
3.03e-16 (2.62e-18)	1.15e-16 (3.51e-18)	1.74e-14 (1.21e-16)
4.91e-17 (8.60e-18)		4.92e-17 (8.60e-18)

Table 2: $A = \text{triu}$ (16, -5)^T

$$\kappa_{\infty}(A) = 3.57 \times 10^{-13}, \quad \text{cond}(A) = 9.40 \times 10^{11}$$

$\text{Giv}, r = 1$	$\text{Giv}, r = 2$	$\text{Fast}, r = 2$
8.42e-17 (4.23e-19)	8.42e-17 (4.23e-19)	3.12e-9 (4.24e-19)
		1.68e-16 (1.27e-18)

Table 3: $A = \text{ipjfact}$ (7, 1)

$$\kappa_{\infty}(A) = 1.69 \times 10^{-14}, \quad \text{cond}(A) = 6.85 \times 10^{10}$$

$\text{Giv}, r = 1$	$\text{Giv}, r = 2$	$\text{Fast}, r = 2$
8.45e-16 (1.51e-20)	3.68e-16 (1.23e-20)	2.00e-12 (5.49e-20)
1.98e-17 (1.07e-20)	5.89e-18 (3.19e-21)	7.07e-18 (2.15e-21)

Table 4: $A = \text{rand}$ (32)

$\kappa_\infty(A) = 6.00 \times 10^{-2}$, $\text{cond}(A) = 3.33 \times 10^{-2}$

$\text{GM}, r = 1$	$\text{GM}, r = 8$	$\text{Fast}, r = 8$
1.48e-16 (2.67e-17)	2.34e-16 (4.56e-17)	6.81e-16 (1.07e-16)
	2.31e-17 (7.16e-18)	2.99e-17 (8.38e-18)

2.3 Block Triangular LU Factorization

Next, we discuss the computation of a true block LU factorization $A = LU \in \mathbb{R}^{n \times n}$, where L and U are block lower triangular and upper triangular respectively. This factorization is not used in LAPACK but it has attracted attention because of its suitability for parallel machines [16, 27]. Assuming that $A_{11} \in \mathbb{R}^{r \times r}$ is nonsingular we can write

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} I & 0 \\ L_{21} & I \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ 0 & B \end{bmatrix} = LU, \quad (23)$$

which leads to the following algorithm for computing L and U [16, 27]:

1. $X = A_{21} A_{11}^{-1}$.
2. $L_{21} = A_{21} X$.
3. $B = A_{22} - L_{21} A_{12}$.
4. Compute the block LU factorization of B , recursively.

The explicit computation of $A_{11}^{-1} = U_{11}^{-1}$ in this algorithm can lead to greater efficiency on parallel machines, for if U_{11}^{-1} is stored rather than U_{11} the backsubstitution phase of solving $Ax = b$ consists entirely of matrix-vector multiplications.

Since the algorithm does not incorporate pivoting it is intended only for applications in which it is safe not to pivot in Gaussian elimination, for example, where A is diagonally dominant or symmetric positive definite. We wish to point out, however, that even in these situations, stability is not guaranteed. To see this it suffices to consider the $(2, 1)$ block (relative to the partitioning (23)) of the residual,

$$(A - \hat{L}\hat{U})_{21} = A_{21} - \hat{L}_{21}A_{11}.$$

Suppose, for simplicity that $X = A_{21} A_{11}^{-1}$ is obtained *exactly*. Then, using (1), we have

$$\hat{L}_{21} = L_{21} + \mathcal{X}_{21}, \quad \|\mathcal{X}_{21}\| \leq c_1(n-r, r, r) \kappa(A) \|A_{11}^{-1}\|,$$

and so

$$\|(A - \hat{L}\hat{U})_{21}\| = \|\mathcal{X}_{21}A_{11}\| \leq c_1(n-r, r, r) \kappa(A) \|A_{11}\|.$$

This the bound for $\|A - \hat{L}\hat{U}\|$ depends on the condition of A_{11} , and by the recursive nature of the algorithm on the condition of other $r \times r$ submatrices arising in the

Table 5: Relative residuals for LU factorization

$$\begin{aligned} \kappa_\infty(A_{16}(\alpha=0.7)) &= O(10^{-8}) \\ \kappa_\infty(A_{16}(\alpha=1.1)) &= O(10^{-11}) \\ \kappa_\infty(A_{16}(\alpha=3)) &= O(10^{-19}) \end{aligned}$$

r	$\alpha=0.7$	$\alpha=1.1$	$\alpha=3$
1	1.35e-16	8.44e-17	0
2	6.19e-17	4.01e-17	9.46e-17
4	3.75e-17	1.15e-16	1.29e-14
6	5.76e-16	6.13e-15	1.86e-12
8	1.10e-14	5.60e-14	4.82e-10
10	1.54e-14	6.98e-13	3.74e-08
12	1.06e-12	4.30e-11	7.43e-05
14	4.50e-12	9.11e-10	2.89e-02

algorithm. If any of these submatrices is ill-conditioned then there is the possibility of a large residual.

Examples demonstrating instability are readily generated. Consider the symmetric positive definite Mersenne matrix $A = \text{triu}(n, \alpha) \text{tril}(n, \alpha) \in \mathbb{R}^{n \times n}$ from [20], where $\text{triu}(n, \alpha)$ is defined in section 2.2. We ran the block LU algorithm on $A_{16}(\alpha)$ in MATLAB for block sizes $r = 1, 2, 4, 8, 10, 12, 14$ with three different α , using canonical BSS. The relative residuals $\|\hat{L}\hat{U}\|_\infty / \|A\|_\infty$ for the LU factorization are displayed in Table 5; they clearly reveal instability increasing with the block size (note that $\kappa(A_{16}(\alpha=1.1))$ increases with the block size).

This instability of block LU factorization does not seem to be well-known. We suspect that in most applications where the algorithm has been used A has been diagonally dominant, for if the diagonal dominance is sufficiently strong then the increases occurring in the algorithm are guaranteed to be of modest norm.

2.4 Recognizing a Stable Block Decomposition

How can we distinguish a stable block algorithm of the kind discussed in section 2.1 from the unstable one in section 2.3? Here we present an informal approach allowing easy recognition without the need for a full error analysis. Our approach has some similarities with that of Wilkinson [33] who describes a general approach to analysing unblocked algorithms.

We consider LU factorization with partial pivoting although the same approach applies to the other factorizations mentioned at the end of section 2.1 and to QR factorization. We view the algorithm as a sequence of coupled decompositions

$$P_i L_i U_i + R_{i+1} = A + \Delta_i, \quad i = 0, m, \quad (24)$$

where one or more such decompositions describe a single step of the block algorithm.

P_i is a permutation matrix and Δ_i represents rounding errors introduced by the first i steps of the algorithm. Initially $P_0 = L_0 = I$, $U_0 = A$, and $R_{0+1} = \Delta_{0+1} = 0$, while L_m is unit lower triangular, U_m is upper triangular and $R_{m+1} = 0$. The term R_{i+1} is introduced for notational convenience.

The matrices P_{i+1} , L_i , U_i and R_{i+1} are transformed to P_{i+1+1} , L_{i+1} , U_{i+1} and R_{i+1+1} by a single matrix operation, either unblocked (BAS or BAS), or blocked (BAS).

The rounding error introduced by this matrix operation is $\Delta_{i+1} - \Delta_i$. If each Δ_i is small ($\|\Delta_i\| = O(u \|A\|)$) then the algorithm is stable.

To illustrate, the first stage of the block LU factorization in section 2.1 (ignoring pivoting) may be written using (24) with $i = 0, 4$, as

$$\begin{aligned} A &= I \begin{bmatrix} A_{11} & A_{21} \\ A_{12} & A_{22} \end{bmatrix} \\ &= \begin{bmatrix} L_{11} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} U_{11} & 0 \\ A_{21} & A_{22} \end{bmatrix} + \begin{bmatrix} 0 & A_{12} \\ 0 & 0 \end{bmatrix} + \begin{bmatrix} \Delta_{11} & 0 \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} L_{11} & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ A_{21} & A_{22} \end{bmatrix} + \begin{bmatrix} \Delta_{11} & \Delta_{12} \\ 0 & 0 \end{bmatrix} \\ &= \begin{bmatrix} L_{11} & 0 \\ L_{21} & I \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & -L_{21} U_{12} \end{bmatrix} + \begin{bmatrix} 0 & 0 \\ 0 & A_{22} \end{bmatrix} + \begin{bmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{21} & 0 \end{bmatrix} \\ &= \begin{bmatrix} L_{11} & 0 \\ L_{21} & I \end{bmatrix} \begin{bmatrix} U_{11} & U_{12} \\ 0 & B \end{bmatrix} + \begin{bmatrix} \Delta_{11} & \Delta_{12} \\ \Delta_{21} & \Delta_{22} \end{bmatrix}. \end{aligned} \quad (25)$$

It is clear from the assumptions (1) and (2) that each rounding error term $\Delta_{i,j}$ has a bound of order $u \|A\|$, and (25) shows that their contribution is additive. It follows that the whole process is stable.

One possible obstacle to stability which does not occur in LAPACK would be computing part of a decomposition (say $A_{11} = L_{11} U_{11}$ above), using it to compute other parts of the decomposition (L_{21}, U_{12} , and B) and then recomputing it by some method yielding different rounding errors. There is no motivation for this in Gaussian elimination but it is conceivable that such redundant operations would be needed to use the BAS. (There are redundant operations in the use of block orthogonal transformations in the next section, but they do not lead to this difficulty).

How could this recomputation change stability? Suppose we refactorize $A_{11} \approx \hat{L}_{11} \hat{U}_{11}$ stably after the last step above, and replace L_{11} by \hat{L}_{11} and U_{11} by \hat{U}_{11} . This will increase Δ_{12} by $(\hat{L}_{11} - L_{11})U_{12}$ and Δ_{21} by $L_{21}(\hat{U}_{11} - U_{11})$, and neither of these quantities is guaranteed to be small.

3 Orthogonal Transformations

In this section we consider block algorithms based on orthogonal transformations. The algorithms of interest include QR factorization, orthogonal reduction to Hessenberg tridiagonal or bidiagonal form, the unsymmetric QR algorithm and algorithms for generalized eigenvalue or singular value computations. The techniques used in LAPACK for constructing block versions of these algorithms are based on the aggregation of Householder transformations. Our aim is therefore to analyse the stability of these aggregation techniques.

The form of aggregation is the “W” representation of Borch and Van Loan [5]. This involves representing the product $Q = P_r P_{r-1} \dots P_1$ of r Householder transformations $P_i = I - u_i u_i^T \in \mathbb{R}^{n \times n}$ ($u_i^T u_i = 2$) in the form

$$Q_r = I - W_r Y_r^T, \quad W_r, Y_r \in \mathbb{R}^{n \times r}.$$

This can be done using the recurrence

$$W_1 = -u_{-1}, \quad Y_1 = u_{-1}, \quad W_i = [W_{i-1}, -Q_{i-1} u_i], \quad Y_i = [Y_{i-1}, u_i].$$

Using the W-representational block QR factorization can be developed as follows. Partition $A \in \mathbb{R}^{m \times n}$ ($m \geq n$) as

$$A = \begin{bmatrix} A_1 \\ B \end{bmatrix}, \quad A \in \mathbb{R}^{m \times r},$$

and compute the Householder QR factorization of A_1 ,

$$P_r P_{r-1} \dots P_1 A_1 = \begin{bmatrix} R_1 \\ 0 \end{bmatrix}.$$

The product $P = P_r P_{r-1} \dots P_1 = I + W_r Y_r^T$ is accumulated as the P_i are generated and then B is updated according to

$$B \leftarrow (I + W_r Y_r^T) B = B + W_r (Y_r^T B),$$

which involves only BAS operations. The process is now repeated on the last $m-r$ rows of B .

An alternative form of accumulation is proposed in [14] for $r=2$, extended to general r in [13], and used in [2]. In the context of orthogonal similarity reduction to Hessenberg form the technique involves expressing

$$P_r P_{r-1} \dots P_1 A P_1 \dots P_{r-1} P_r = A - U_r V_r^T - W_r U_r^T, \quad (26)$$

where $U_r, V_r, W_r \in \mathbb{R}^{n \times r}$. We refer the reader to [13] for details of how to obtain this representation. The key point is that, once again, only BAS operations are involved in computing the update, once U_r, V_r and W_r have been formed.

Now we consider numerical stability. We concentrate on the W-technique, and note that similar analysis applies to the alternative method of aggregation (26), as well as to the more storage efficient compact W-representation of Schreiber and Van Loan [28]. First, we note that the construction of the W and Y matrices is done in a stable manner (indeed it does not involve the BAS): Borch and Van Loan [5] show that the computed $\hat{Q} = I + \hat{W} \hat{Y}^T$ is such that

$$\|\hat{Q}^T \hat{Q} - I\| = O(u), \quad (27)$$

$$\|\hat{W}\| = O(1), \quad \|\hat{Y}\| = O(1). \quad (28)$$

The condition (27) implies that

$$\hat{Q} = U + \mathcal{N}, \quad U^T U = I, \quad \|\mathcal{N}\| = O(u), \quad (29)$$

that is, \hat{Q} is close to an exactly orthogonal matrix

Next we consider the application of \hat{Q} . Suppose we form $C = \hat{Q}B = (I + \widehat{W}\widehat{Y}^T)B$, obtaining

$$\hat{C} = fl(B + fl(\widehat{W}(\widehat{Y}^T B))).$$

Analyzing this **BLAS**-based computation using (1) and (2) it is straightforward to show that

$$\begin{aligned} \hat{C} &= UB + \Delta C = U(B + U^T \Delta C), \\ \|\Delta C\| &\leq c_4 n (c_1(r, n, n) + (n, r, m) \|B\|) + O(u^2), \end{aligned} \quad (3)$$

where c_4 is a constant of order 1. This result shows that the computed update is an exact orthogonal update of a perturbation of B , where the norm of the perturbation is bounded in terms of the error constants for the **BLAS**.

In the case $r = 1$, (3) reduces to Wilkinson's result on the application of a single Householder transformation [32, p.16]. Wilkinson uses this result to obtain a backward error result for the application of a sequence of Householder transformations [32, pp.160-161]. With the use of (3), it is straightforward to show that Wilkinson's method of analysis can be adapted to accumulate W updates. Alternatively to obtain a backward error result for a sequence of orthogonal similarity transformations, we can simply insert the bound (3) into the general analysis of [26, sec.6.5] or [31]. It follows that the standard backward error analysis results for Householder transformation algorithms remain valid when the W technique is used, as long as the constants in the error bounds are replaced by appropriate linear combinations of c_1 terms.

Our overall conclusions are as follows. First, algorithms that employ aggregated Householder transformations with conventional **BLAS** are as stable as the corresponding point algorithms. Second, the use of fast **BLAS** for applying the updates affects stability only through the constants in the backward error bounds.

4 Concluding Remarks

Our main conclusion is that the use of fast **BLAS** satisfying (1) and (2) with LA **PACK** algorithms is "safe" from a numerical standpoint. The algorithms retain

their nominal backward stability but the actual backward errors may increase to reflect any decreased accuracy in the **BAS**. Also, any special properties relating to componentwise backward error may be lost in switching to fast **BAS**.

How large the increase in backward errors will be, on average, is difficult to say, since the theoretical error bounds tend to be quite pessimistic (this applies particularly to the bound for Strassen's method—see [18]). A quantitative assessment of the speed versus stability tradeoff must await the accrual of experience in using fast **BAS** for the values of n for which useful speeds are obtained.

In the important application of solving $Ax = b$ by LU factorization the use of fast **BAS** need carry no stability penalty: we have shown that fixed precision iterative refinement with conventionally computed residuals can improve nominal stability and can even produce a small componentwise relative backward error, although we have not been able to state useful conditions under which such improvements are guaranteed.

Finally, we reiterate a point discussed in [18]. When replacing conventional **BAS** by fast **BAS** in an *iterative* algorithm it is important to consider the implications for the convergence tests—a change in the **BAS** may necessitate a retuning of the algorithm parameters. Ideally, a convergence test will be effective across different **BAS** implementations and computer arithmetics and so will not need tuning. The stopping criterion used in LAPACK for iterative refinement of linear equations (see section 2.2) has been designed with this goal in mind. Experience will tell whether the goal has been achieved.

References

- [1] MAidi, J.Widom and I.S. Duff Solving sparse linear systems with sparse backward error, *SIAM. Matrix Anal. Appl.*, 10 (1989), pp 165-190.
- [2] Z. Bai and J.Widom, On a block implementation of Hessenberg multishift QR iteration, *Int. J. High Speed Computing* 1 (1989), pp 97-121.
- [3] DH Bailey, K Lee and HD Simon Using Strassen's algorithm to accelerate the solution of linear systems, Report BR90-001, NS System Division, NSA Air Research Center, Mott Field, CA9035, 1990.
- [4] CH Bischof and J.J. Dongarra, A project for developing a linear algebra library for high performance computers, Reprint MCS-P05-089, Mathematics and Computer Science Division, Argonne National Laboratory, 1989.
- [5] CH Bischof and CE Van Loan, The W-representation for products of Householder matrices, *SIAM. Si. Stat. Comp.*, 8 (1987), pp s2-s13.
- [6] RP Brent, Error analysis of algorithms for matrix multiplication and triangular decomposition using Wograd's identity *Numer. Math.*, 16 (1970), pp 145-156.
- [7] G Dahlquist and Å Björck, *Numerical Methods*, Prentice-Hall, Englewood Cliffs, New Jersey, 1974.
- [8] J.Widom, On floating point errors in Godesky, LAPACK Working Note #4, Department of Computer Science, University of Tennessee, Knoxville, 1989.
- [9] J.Widom, J.J. Dongarra, J.J. Di Gonz, A Greenbaum S.J. Hamming and DC Sorensen, Prospects for the development of a linear algebra library for high performance computers, Technical Memorandum 97, Mathematics and Computer Science Division, Argonne National Laboratory Illinois, 1987.
- [10] J.Widom, J.J. Di Gonz, S.J. Hamming and DC Sorensen, Guidelines for the design of symmetric eigenroutines, SVD and iterative refinement and condition estimation for linear systems, LAPACK Working Note #4, Technical Memorandum 111, Mathematics and Computer Science Division, Argonne National Laboratory, 1988.

- [11] J.J. Dongarra, J.J. Di Gonz, I.S. Duff and S.J. Hamming A set of Level 3 basic linear algebra subprograms, *ACM Trans. Math. Sft.*, 16 (1990), pp 1-17.
- [12] J.J. Dongarra, J.J. Di Gonz, I.S. Duff and S.J. Hamming Algorithm 679: A set of Level 3 basic linear algebra subprograms, *ACM Trans. Math. Sft.*, 16 (1990), pp 18-28.
- [13] J.J. Dongarra, S.J. Hamming and DC Sorensen Block reduction of matrices to condensed form for eigenvalue computations, *J. Comp. Appl. Math.*, 27 (1989), pp 215-227.
- [14] J.J. Dongarra, L. Kaufman and S.J. Hamming Squeezing the most out of eigenvalue solvers on high performance computers, *Linear Algebra and Appl.*, 77 (1986), pp 113-136.
- [15] K Gallivan, Whally, U Mer and AH Saich, Impact of hierarchical memory systems on linear algebra algorithm design, *Int. J. of Supercomputer Applic.*, 2 (1988), pp 12-48.
- [16] KA Gallivan, RJ. Perrins and AH Saich, Parallel algorithms for dense linear algebra computations, *SIAM Review* 32 (1990), pp 54-135.
- [17] GH Golub and CF Van Loan, *Matrix Computations*, Second Edition, Johns Hopkins University Press, Baltimore, Maryland, 1989.
- [18] NJ. Higham Exploiting fast matrix multiplication within the level 3 BLAS Technical Report 89-984, Department of Computer Science, Cornell University 1989, to appear in *ACM Trans. Math. Sft.*
- [19] NJ. Higham How accurate is Gaussian elimination?, in *Numerical Analysis 1989, Proceedings of the 13th Dundee Conference*, Pitman Research Notes in Mathematics 228, DE Giffis and GA Watson, eds., Longman Scientific and Technical, 1990, pp 137-154.
- [20] NJ. Higham Collection of test matrices in *MLAB* Technical Report 89-1025, Department of Computer Science, Cornell University 1989.

- [21] N.J. Higham, Stability of a method for multiplying complex matrices with three real matrix multiplications, *Numerical Analysis Report No. 181*, University of Manchester, England, 1990.
- [22] N.J. Higham, Is fast matrix multiplication of practical use?, to appear in *SAM News*.
- [23] N.J. Higham, Iterative refinement enhances the stability of QR factorization method for solving linear equations, *Numerical Analysis Report No. 182*, University of Manchester, England, 1990.
- [24] P. Mes and G. Baccati, Blocked Cholesky factorization using level 3 BLAS LAPACK Working Note #12, Mathematics and Computer Science Division, Argonne National Laboratory, Illinois, 1989.
- [25] W. Götli and W. Pöggendorf, Computability of approximate solution of linear equations with given error bounds for coefficients and right-hand sides, *Numer. Math.*, 6 (1964), pp. 45-49.
- [26] B.N. Parlett, *The Symmetric Eigenvalue Problem*, Perlice-Hall, Englewood Cliffs, New Jersey, 1980.
- [27] R.S. Schreiber, Block algorithms for parallel machines, in *Numerical Algorithms for Modern Parallel Computer Architectures*, M.H. Schultz, ed., IMVdruck In Mathematics and Its Applications 13, Springer-Verlag Berlin, 1988, pp. 197-207.
- [28] R.S. Schreiber and C.F. Van Loan, A storage efficient W representation for products of Householder transformations, *SAMJ. Si. Stat. Comput.*, 10 (1989), pp. 53-57.
- [29] R.D. Sed, Iterative refinement implies numerical stability for Gaussian elimination, *Math. Comp.*, 35 (1980), pp. 817-832.
- [30] V. Strassen, Gaussian elimination is not optimal, *Numer. Math.*, 13 (1969), pp. 354-356.
- [31] J.H. Wilkinson, Error analysis of eigenvalue techniques based on orthogonal transformations, *J. Soc. Indust. Appl. Math.*, 10 (1962), pp. 162-195.

- [32] J.H. Wilkinson, *The Algebraic Eigenvalue Problem*, Oxford University Press, 1965.
- [33] J.H. Wilkinson, Error analysis revisited, IMA Bulletin 22 (1986), pp 192-200.
- [34] S. Wograd, A new algorithm for inner product, IRE Trans. Comput., C18 (1968), pp 63-64.