

# An iterative solver benchmark<sup>1</sup>

Jack Dongarra, Victor Eijkhout and  
Henk van der Vorst

Revised 31 August 2001

We present a benchmark of iterative solvers for sparse matrices. The benchmark contains several common methods and data structures, chosen to be representative of the performance of a large class of methods in current use. We give results on some high performance processors that show that performance is largely determined by memory bandwidth.

## 1. Introduction

In scientific computing, several benchmarks exist that give a user some idea of the to-be-expected performance given a code and a specific computer. One widely accepted performance measurement is the Linpack benchmark [4], which evaluates the efficiency with which a machine can solve a dense system of equations. Since this operation allows for considerable reuse of data, it is possible to show performance figures that are a sizeable percentage of peak performance, even for machines with a severe imbalance between memory and processor speed.

However, sparse linear systems are at least as important in scientific computing, and for these the question of data reuse is more complicated. Sparse systems can be solved by direct or iterative methods, and especially for iterative methods one can say that there is little or no reuse of data. Thus, such operations will have a performance bound by the slower of the processor and the memory, in practice: the memory.

We aim to measure the performance of a representative sample of iterative techniques on any given machine; we are not interested in comparing, say, one preconditioner on one machine against another preconditioner on another machine. In fact, the range of possible preconditioners is so large, and their performance so much dependent on the specific problem, that we do

not even compare one preconditioner to another on the same machine. Instead, we identify kernels that will have a representative performance, and test those; the performance of a whole code is then a composite of the performances of the various sparse kernels.

An earlier report on the performance of supercomputers on sparse equation solvers can be found in [5].

## 2. Motivation

Iterative methods are hard to benchmark. The ultimate measure one is interested in is ‘time expended to reach a certain accuracy’. This number is a function of at least the following variables:

1. The iterative method: for each type of linear system (symmetric / hermitian / general, definite / almost definite / indefinite, et cetera) there are several candidate methods with differing mathematical and computational characteristics. For instance, GMRES [9] has an advantage over BiCGstab [10] in that it minimizes residuals, and is able to use Level 3 BLAS; it has as disadvantage that work per iteration is higher and memory demands larger.
2. The preconditioner: some systems are well-conditioned enough that a simple Jacobi preconditioner (or equivalently, diagonal scaling before the iteration process) will be good enough, others require complicated preconditioners to tame the spectrum.
3. Data structures used: the same matrix can be stored in a number of formats, some of which may have computational advantages. For instance, problems with several physical variables per node can be expressed in such a way that dense matrix operations can be used.
4. The implementation of computational kernels: different implementations of the same operation can have different performance characteristics.
5. Architectural details of the machine used, in particular the processor speed and the memory bandwidth. Depending on the operation, the one or the other may be of more or less importance.

---

<sup>1</sup>This research is sponsored in part by Subcontract #R71700J-29200099 with William Marsh Rice University and Subcontract #B503962 with The Regents of the University of California.

6. Computer Arithmetic: some problems may be solvable to a desired accuracy in single precision, others may require double precision. Even at a given precision, two different implementations of arithmetic may give different results, so that the same program on two different machines will not behave the same.
7. The problem to be solved. The behaviour of all of the above points is a function of the problem to be solved. The number of iterations (that is, behaviour of the iterative method and preconditioner) is a function of the spectrum of the matrix. The flop rate per iteration (determined by implementation of kernels, and processor characteristics) is a function of the sparsity structure of the matrix.
8. Problem size. Even ignoring such obvious factors as performance degradation because of page swapping, the size of a problem's data set can influence performance strongly. On vector machines, increasing the problem size will most likely increase the vector length and hence performance. Cache-based machines will reach optimum performance for a problem that fits in cache. Some architectures (notably the Alpha chip) have a large enough cache that realistic problems can actually be solved in-cache. In general, however, increasing the problem size will always overflow the cache, and iterative methods have in general little opportunity for cache reuse. Hence we expect decreasing performance to some 'asymptotic' limit determined more by bandwidth than by processor speed.

Given this plethora of variables, it should be clear that the desired metric is unfortunately not a benchmarkable one.

In order to provide a comprehensive iterative benchmark, a package would have to cover the range of choices of items 1–3. Taking into account possible criticism of the quality of implementation used (item 4), we see that a comprehensive benchmark is almost equivalent to a library of all high-quality implementations of all existing methods. Clearly, this is not a practical proposition.

Limiting ourselves to a small collection of representative methods would not enable us to benchmark the 'time to solution' metric either. The problem here is that – in total contrast to the dense solvers of the Linpack benchmark – given a machine, algorithm, and data structure, the time to solution is dependent on the numerical values of entries of the system.

In the case of dense solvers, in essence the only problem parameter determining performance is the size of the system. If certain properties of the system are known, such as definiteness or symmetry, optimised implementations can take advantage of this. However, between two systems that have the same values of these properties, the size is the only discriminating factor. In the iterative case, two systems could have the same size and sparsity structure, and be solvable with the same iterative method and preconditioner, but show completely different numbers of iterations.

We could then decide to provide a set of standard problems, but this may very well overlook a particular kind of problem the user is interested in, and which has its own set of performance characteristics.

Instead of even attempting the above-sketched perfect benchmark, we have opted for providing a small set of iterative methods, preconditioners and storage schemes, on the criterion that they have a performance representative of larger classes. We then envision the following scenario for the use of our benchmark: users solve a few representative test problem with their preferred solver, thus gaining a notion of convergence speed in terms of numbers of iterations. They can then consult our benchmark to gauge the per-iteration performance of their solver, which combined with their prototype run on the test problem gives an accurate impression of the expected running time.

We want to stress from the outset that we did not aim to present the most sophisticated methods. Rather, by considering combinations of the representative elements used in the benchmark a user should be able to get a good notion of the expected performance of methods not included. In effect, our benchmark only indicates performance per occurrence of a computational kernel. Users can then estimate their codes' performance by combining these numbers, applied to their specific implementation, with the numbers of iterations needed for their problem.

Consistent with this philosophy of measuring only performance per occurrence, we terminate each benchmark run after a fixed number of iterations. The number of iterations to the desired accuracy is mostly a function of the problem, and only to a minor extent of the computer architecture. It is almost independent of the implementation of the computational kernels that are the subject of this benchmark, since these can only influence the convergence through different round-off behaviour. For the most part we will assume that different implementations of the same kernel have roughly the same round-off characteristics, so that our

per-iteration measurement is a sufficient predictor of the overall efficiency.

To account for the effects of caching in our benchmark, we run a series of problems of increasing size for each choice of method, preconditioner, and storage scheme. In this sequence we measure both the maximum attained performance, often a clear function of cache size, and we estimate an ‘asymptotic’ performance as the problem scales up beyond the cache size limit; see Sections 3.2 and 7.

We conclude this section by giving a brief description of the sparse kernels. More detailed discussion will follow in Section 4.

As storage schemes we offer diagonal storage, and compressed row storage. Both of these formats represent typical matrices for three-dimensional finite element or finite difference methods. The diagonal storage, using seven diagonals, is the natural mode for problems on a regular (‘brick’) domain; the compressed row storage is the natural storage scheme for irregular domains. Thus these choices are representative for most single-variable physical problems.

The iterative methods provided are CG and GMRES. The plain Conjugate Gradients method is representative of all fixed-storage methods, including sophisticated methods for nonsymmetric problems such as BiCGstab; the GMRES method represents the class of methods that have a storage demand that grows with the number of iterations.

Each iterative method can be run unpreconditioned – which is computationally equivalent to using a Jacobi preconditioner – or with an ILU preconditioner. For the diagonal storage scheme a block Jacobi method is also provided; this gives a good indication of domain decomposition methods. If such methods are used with inexact subdomain solves, the ILU preconditioner gives the expected performance for these.

### 3. Structure of the benchmark

We have implemented a benchmark that constructs a test matrix and preconditioner, and solves a linear system with them. Separate flop counters and timers are kept for the work expended in vector operations, matrix vector products, preconditioner solves, and other operations involved in the iterative method. The flop counts and flops rates in each of these categories, as well as the overall flops rates, are reported at the end of each run.

The benchmark comprises several storage formats, iterative methods, and preconditioners. Together these form a representative sample of the techniques in typical sparse matrix applications. We describe these elements in more detail in Section 4.

We offer a reference code, which is meant to represent a portable implementation of the various methods, without any machine-specific optimisations. The reference code is written in Fortran, using double precision arithmetic throughout; implementers are not allowed to increase speed by switching to single precision. In addition to the reference code we supply a number of variants that should perform better on certain machines, and most likely worse on some others; see Section 6.

#### 3.1. Conformance of the benchmark

Since we leave open the possibility that a local implementer make farguing changes to the benchmark code (see Section 3.2), we need to ascertain that the optimised code still conforms to the original, that is, that the implementer has only optimised the *implementation* of the algorithms for a specific machine, and has not *replaced* one algorithm by another. Unfortunately, by the nature of iterative methods, this is a hard problem.

Iterative methods such as Conjugate Gradients are forwardly unstable; they do not have a self-correction mechanism the way stationary iterative methods have. This means that comparing the computed quantities of one method against another, or, if it were possible, comparing it against an exact arithmetic method, would be essentially meaningless. Likewise, any bounds from using interval arithmetic would be so large as to be useless. The reason that conjugacy-based iterative methods are useful despite this unstable behaviour, is that a more careful analysis shows that accumulated roundoff does not cause divergence, but will only delay convergence by a modest number of iterations (see [6,7]).

Thus we are in the predicament of having to enforce the numerical conformance of a method that is intrinsically unstable in the naive sense of the term. Our practical test, therefore, checks how much an optimised code differs from the original after computing one iteration, that is, before any large divergence from exact arithmetic starts to show up.

There are several ways this test can be implemented specifically, and none of them are completely satisfying. Absent an implementation in exact arithmetic, we basically have the choice of comparing optimised code against a reference code on the same machine, and code on some reference machine. For the first option, we

would in essence trust that arithmetic is implemented in a reasonable manner on machines available today, so that the reference code, without any optimisations that alter the straightforward execution of instructions, is indeed a reasonable reference point. The second option raises the question what machine would be trustworthy enough to function as a supplier of reference results. We can evade that question by observing that, under the above assumption that all machines on the market today have a ‘reasonable’ implementation of arithmetic, the difference between two concrete implementations of the benchmark is bounded by the sum of the differences between either and an exact arithmetic implementation. Our measurement of difference against a result on file – generated by us, the producers of the benchmark – would then give a small overestimate of the true error.

This latter strategy, comparison against results on file, is the one we have chosen. (We are in good company with this decision: the NAS parallel benchmarks [1] use the same verification scheme, and in fact with a far stricter test than ours.) In practice, we test whether the deviation remains under 100 times the machine precision. The number ‘100’ itself is somewhat arbitrary; it reflects the limited number of nonzeros in each matrix row, but it does not reflect the worst-case error bound of  $O(N)$  that comes in through the inner product calculations in the iterative method. A test after one iteration could conceivably be performed against an interval arithmetic implementation, but we have not done so.

We emphasize that this is a static test, designed to allow only changes to the reference code that are not numerically significant. In particular, it precludes an implementer from replacing the preconditioner by a different one. We justify this from our standpoint that the benchmark is not a test of the best possible preconditioner or iterative method, but rather of methods representative for a wider class with respect to computer performance.

Since the benchmark includes ILU preconditioners, this static conformance test would a priori seem to be biased against parallel implementations of the benchmark. This point is further elaborated in Section 5.

### 3.2. Benchmark reporting

An implementer of the benchmark can report performance results on various levels, each next level encompassing all of the earlier options.

1. Using only compiler flags in the compilation of the reference code.
2. Using compiler directives in the source of the reference code.
3. Rewriting the reference code in such a way that any differences are solely in a different order of scheduling the operations.
4. Rewriting the reference code by replacing some algorithm by a mathematically equivalent formulation of the algorithm (that is: in exact arithmetic the (intermediate) results should be the same).

The last two levels may or will in general influence the numerical results, so results from codes thus rewritten should be accompanied by proof that the specific realisation of the benchmark reproduces the reference results within a certain tolerance.

Each run of the benchmark code ends with a report on how many floating point operations were performed in the various operations. Implementers should use these numbers to do reporting (rather than using hardware flop counters, for instance), but they are free to substitute their own timers.

The benchmark comes with shell scripts that run a number of tests, and report both best performance and asymptotic performance for the whole code and elements of it. Asymptotic performance is determined by making a least-squares fit  $y = a + bx^{-1}$  through the data points, where  $y$  is the observed megaflop rate and  $x$  is the dataset size. The asymptotic performance is then the value of  $a$ .

This assumption on the performance behaviour accommodates both cache-processors, for which we expect  $b > 0$  as the dataset size overflows the cache, and vector processors, for which we expect  $b < 0$  as performance goes up with increasing vector length. For cache-based processors we may expect a plateau behaviour if the cache is large; we discard the front of this plateau when calculating the asymptotic performance.

Asymptotic performance may be determined with some manual intervention from the benchmark user: our test scripts run problems up till a certain size, but future machines may have a cache size exceeding that. In that case the benchmarker will have to run larger problems until the ‘bottoming-out’ effect becomes visible as the problem data overflows the cache.

## 4. Elements of the benchmark code

The user of the benchmark has the following choices in determining the problem to run.

#### 4.1. Storage formats

The matrix can be in the following formats:

- Diagonal storage for a seven-diagonal matrix corresponding to finite differences in three dimensions;
- Compressed row storage of a matrix where the sparsity structure is randomly generated; each row has between 2 and 20 nonzeros, each themselves randomly generated, and the bandwidth is  $\leq n^{2/3}$  which again corresponds to a problem in three space dimensions.

For both formats a symmetric variant is given, where only half the matrix is stored.

The diagonal storage is very regular, giving code that has a structure of loop nests of depth three. Vector computers should perform very efficiently on this storage scheme. In general, all index calculation of offsets can be done statically.

Matrix-vector operations on compressed row storage may have a different performance in the transpose case from the regular case. Such an operation in the regular case is based on inner products; in the transpose case it uses vector updates (axpy operations). Since these two operations have different load/store characteristics, they may yield different flops rates. In the symmetric case, where we store only half the matrix, such operations use in fact the regular algorithm for half the matrix, and the transpose algorithm for the other half. Thus, the performance of, for instance, the matrix-vector product, will be different in GMRES from in the Conjugate Gradient method.

The CRS format gives algorithms that consist of an outer loop over the matrix rows, with an inner loop that involves indirect addressing. Thus, we expect a lower performance, especially on machines where the indirect addressing involves an access to memory.

#### 4.2. Iterative methods

The following iterative methods have been implemented (for more details on the methods mentioned, see the Templates book [2]):

- Conjugate Gradients method; this is the archetypical Krylov space method for symmetric systems. We have included this, rather than MINRES or SYMLQ, for its ease of coding, and for the fact that its performance behaviour is representative of the more complicated methods. The results for CG are also more-or-less representative for transpose-

```

Let  $A, M, x, b$  be given;
compute  $r_1 = Ax - b$ ;
for  $i = 1 \dots 10$ 
  solve preconditioner:  $z = M^{-1}r$ 
  inner product  $\rho_i = r^t z$ 
  if  $i > 1$ , update  $p \leftarrow z + p(\rho_i/\rho_{i-1})$ 
  matrix vector product:  $q = Ap$ 
  inner product  $\pi = p^t q$ 
  update  $x \leftarrow x - p\pi$ 
  update  $r \leftarrow r - q\pi$ 

```

Fig. 1. Conjugate gradient algorithm.

```

Let  $A, M, x, b$  be given;
for  $i = 1 \dots 10$ 
  matrix and preconditioner apply:  $z = AM^{-1}r$ 
  orthogonalize  $z$  against all earlier  $v_j, j < i$ 
  normalize  $v_i \leftarrow z/\|z\|$ 
  update QR factorisation of size  $i + 1 \times i$  Hessenberg matrix
  Update  $x \leftarrow x - \sum_i v_i c_i$ 

```

Fig. 2. One restart cycle of the generalized minimum residual method.

free methods for nonsymmetric systems, such as BiCGstab, which also have a storage demand constant in the number of iterations.

- BiConjugate Gradients. In many ways this method has the same performance characteristics as CG, but it differs in that it additionally uses a product with the transpose matrix  $A^t$ . In many cases forming this product is impractical, and for this reason BiCG and such methods as QMR are less used than transpose-free methods such as BiCGstab. We have included it nevertheless, since the performance of this kernel can not be derived from others. For diagonal matrix storage there is actually no difference between the regular and transpose matrix-vector product; for compressed storage it is the difference between a dot product and a vector update, both indirectly addressed.
- Generalized Minimum Residual method, GMRES. This popular method has been included because its performance behaviour is very different from CG: storage and computational complexity are an increasing function of the iteration count. For that reason GMRES is most often used in cycles of  $m$  steps. For low values of  $m$ , the computational performance for GMRES will not be much different than for CG. For larger values, say  $m > 5$ , the  $j$  inner products in the  $j$ -th iteration may influence the performance. We have included GMRES(20) in our benchmark.

The Conjugate and BiConjugate gradient methods (see Fig. 1) involve, outside the matrix-vector prod-

uct and preconditioner application, only simple vector operations. Thus, their performance can be characterised as similar to that of Level 1 BLAS. The GMRES method (see Fig. 2), on the other hand, uses orthogonalisation of each new generated Krylov vector against all previous, so a certain amount of cache reuse should be possible. See also Section 6 for a rewritten version that uses Level 3 BLAS kernels.

#### 4.3. Preconditioners

The following preconditioners are available:<sup>1</sup>

- No preconditioner;
- Point ILU; for the diagonal storage a true ILU-D is implemented, in the CRS case we use SSOR, which has the same algorithmic structure as ILU;
- Line ILU for the diagonal storage scheme only; this makes a factorisation of the line blocks.
- Block Jacobi for the diagonal storage scheme only; this is parallel on the level of the plane blocks. The block Jacobi preconditioner gives a performance representative of domain decomposition methods, including Schwarz methods.

The point ILU method is typical of commonly used preconditioners. It has largely the structure of the matrix-vector product, but on parallel machines its sequential nature inhibits efficient execution.

The line ILU method uses a Level 2 BLAS kernel, namely the solution of a banded system. It is also a candidate for algorithm replacement, substituting a Neumann expansion for the system solution with the line blocks.

### 5. Parallel realisation

Large parts of the benchmark code are conceptually parallel. Thus we encourage the submission of results on parallel machines. However, the actual implementation of the methods in the reference code is sequential. In particular, the benchmark includes ILU preconditioners using the natural ordering of the variables.

It has long been realised that ILU factorisations can only be implemented efficiently on a parallel architecture if the variables are renumbered from the natural

ordering to, for instance, a multi-colour or nested dissection ordering.

Because of our strict conformance test (see Section 3.1), the implementer is not immediately at liberty to replace the preconditioner by an ILU based on a different ordering. Instead, we facilitate the parallel execution of the benchmark by providing several orderings of the test matrices, namely:

- Reverse Cuthill-McKee ordering.
- Multi-colour ordering; here we do not supply the numbering with the minimal number of colours, but rather a colouring based on [8].
- Nested dissection ordering; this is an ordering based on edge-cutting, rather than finding a separator set of nodes.

The implementer then has the freedom to improve parallel efficiency by optimising the implementation for a particular ordering.

Again, the implementer should heed the distinction of Section 3.2 between execution by using only compiler flags or directives in the code, and explicit rewrites of the code to force the parallel distribution.

### 6. Code variants

We supply a few variants of the reference code that incorporate transformations that are unlikely or impossible to be done by a compiler. These transformations target specific architecture types, possibly giving a higher performance than the reference code, while still conforming to it; see Section 3.1.

**Naive coding of regular ILU** Putting the tests for boundary conditions in the inner loop is bad coding practice, except for dataflow machines, where it exposes the structure of the loop.

**Wavefront ordering of regular ILU** We supply a variant of the code where the triple loop nest has been rearranged explicitly to a sequential outer loop and two fully parallel inner loops. This may benefit dataflow and vector machines.

**Long vectors** At the cost of a few superfluous operations on zeros, the vector length in the diagonal-storage matrix-vector product can be increased from  $O(n)$  to  $O(n^3)$ . This should benefit vector computers.

**Different GMRES orthogonalisation algorithms** There are at least two reformulations of the orthogonalisation part of the GMRES method. They can en-

<sup>1</sup>We have not included the commonly used Jacobi preconditioner, since this is mathematically equivalent to scaling the matrix to unit diagonal, a strategy that has the exact same performance as using no preconditioner.

Table 1  
List of machines used

| Processor  | Manufacturer / type | Clock rate (MHz) | peak Mflops rate | Compiler / options                   |
|------------|---------------------|------------------|------------------|--------------------------------------|
| Alpha EV67 | Compaq              | 500              | 1000             | f77-O5                               |
| Athlon     | AMD                 | 1200             | 2400             | g77-O3 -malign-double -funroll-loops |
| P4         | Dell                | 1500             | 1500             | g77-O3 -malign-double -funroll-loops |
| Power3     | IBM                 | 200              | 800              | xlf-O4                               |

Table 2  
Highest attained performance

| Machine    | Mflop/s |
|------------|---------|
| alpha ev67 | 967     |
| P4         | 754     |
| Athlon     | 461     |
| power3     | 358     |

able use of Level 3 BLAS operations and, in parallel context, combine inner product operations. However, these code transformations no longer preserve the semantics under computer – rather than exact – arithmetic.

## 7. Results

The following tables contain preliminary results for the machines listed in Table 1. We report two types of results:

- In Table 2 we report the top speed reported regardless the iterative method, preconditioner, and problem size. This speed is typically reported on a fairly small problem, where presumably the whole data set fits in cache.
- All further tables report the ‘asymptotic speed’ of the processor, derived by a least-squares fit as described in Section 3.2.

Note that the the former number, the top speed, is of limited value. On superscalar machines like the ones tested here – and unlike on vector machines to which we had no ready access – it is of necessity reached for a small, hence unrealistic problem. In particular on machines with a relatively small cache the number reported is further unreliable since small problems easily drop under the timer resolution. On some processors this is alleviated by the fact that we offer PAPI [3] support and adaptively raising the maximum number of iterations<sup>2</sup> to match the timer resolution.

<sup>2</sup>Note however, that this will only increase the reliability of the performance of the whole problem: the times for the individual components will be almost vanishingly small.

Table 3

Asymptotic performance of diagonal storage matrix-vector product

| Machine    | Mflop/s |
|------------|---------|
| P4         | 198     |
| power3     | 152     |
| alpha ev67 | 150     |
| Athlon     | 84      |

Table 4

Asymptotic performance of CRS matrix-vector product

| Machine    | Mflop/s |
|------------|---------|
| P4         | 259     |
| power3     | 168     |
| alpha ev67 | 107     |
| Athlon     | 75      |

Table 5

Asymptotic performance of symmetrically stored CRS matrix-vector product

| Machine    | Mflop/s |
|------------|---------|
| P4         | 173     |
| power3     | 165     |
| alpha ev67 | 105     |
| Athlon     | 59      |

The more useful and more realistic number is the asymptotic speed, reported in all tables but the first. For the determination of this, we ran a series of problems of such size that the cache is guaranteed to be insufficient to store them. The current test scripts are set up to handle the currently tested chips (the maximum cache being 4 M on the Alpha chip), but for future developments the user may have to supply larger test sizes than the current defaults.

We report asymptotic speed results for the following components:

- The matrix vector product. We report this in regular storage, and in compressed row storage separately for the symmetric case (cg) and the nonsymmetric case since these may have different performance characteristics; see Section 4.1.
- The ILU solve. We also report this likewise in three variants.
- Vector operations. These are the parts of the algorithm that are independent of storage formats. We report the efficiency of vector operations for

Table 6

Asymptotic performance of transpose CRS matrix-vector product

| Machine    | Mflop/s |
|------------|---------|
| P4         | 210     |
| power3     | 150     |
| alpha ev67 | 111     |
| Athlon     | 68      |

Table 7

Asymptotic performance of diagonal storage ILU solve

| Machine    | Mflop/s |
|------------|---------|
| P4         | 119     |
| alpha ev67 | 108     |
| power3     | 105     |
| Athlon     | 70      |

Table 8

Asymptotic performance of CRS ILU solve

| Machine    | Mflop/s |
|------------|---------|
| power3     | 142     |
| alpha ev67 | 92      |
| P4         | 81      |
| Athlon     | 67      |

the CG method; in case of GMRES a higher efficiency can be attained by using Level 2 BLAS and Level 3 BLAS routines. We have not tested this.

The results published here are only preliminary; a fuller database – both with more machines and of more aspects of the methods tested, as well as code variants – will be maintained on netlib.

All results are double precision, that is, 64-bit arithmetic.

However, performance is hard to analyse in detail, a fact that is already familiar from dense benchmark suites, since it is determined by many interacting factors such as the compiler, the chip architecture, and the memory system. For instance, on most machines, the regular CRS matrix-vector product predictably performs slightly better than the transpose product, which involves more memory references. However, the product with a symmetrically stored matrix, which in implementation is a regular product with the upper triangle and a transpose product with the lower triangle, has a lower performance than either. The CRS ILU solve, in turn, which has a very similar structure to the matrix-vector product, has these rankings exactly reverse for a few machines. A full analysis of such results would go far beyond the scope of this paper.

As an illustration of the way the benchmark is used, we consider a conjugate gradient method with ILU preconditioning on diagonal data storage on an IBM Power3 processor. The asymptotic rates of the matrix-

Table 9

Asymptotic performance of symmetrically stored CRS ILU solve

| Machine    | Mflop/s |
|------------|---------|
| P4         | 176     |
| power3     | 154     |
| alpha ev67 | 98      |
| Athlon     | 49      |

Table 10

Asymptotic performance of transpose CRS ILU solve

| Machine    | Mflop/s |
|------------|---------|
| P4         | 174     |
| power3     | 125     |
| alpha ev67 | 82      |
| Athlon     | 64      |

Table 11

Asymptotic performance of vector operations in CG

| Machine    | Mflop/s |
|------------|---------|
| alpha ev67 | 300     |
| P4         | 241     |
| power3     | 206     |
| Athlon     | 143     |

vector, ilu, and vector kernels are 150, 105, and 206 Mflops respectively. Weighing this by the relative numbers of operations ( $13N$ ,  $14N$ ,  $10N$ ) gives an expected performance of 147 Mflops, which is in good accordance with the observed asymptotic performance of 149 Mflops.

## 8. Obtaining and running the benchmark

The benchmark code can be obtained from <http://www.netlib.org/benchmark/sparsebench>. The package contains Fortran code, and shell scripts for installation and post-processing. Results can be reported automatically to [sparsebench@cs.utk.edu](mailto:sparsebench@cs.utk.edu), which address can also be used for questions and comments.

## References

- [1] D.H. Bailey, E. Barszcz, J.T. Barton, D.S. Browning, R.L. Carter, L. Dagum, R.A. Fatoohi, P.O. Frederickson, T.A. Lasinski, R.S. Schreiber, H.D. Simon, V. Venkatakrishnan and S.K. Weeratunga. The NAS parallel benchmarks, *Intl. Journal of Supercomputer Applications* **5** (1991), 63–73.
- [2] R. Barrett, M. Berry, T.F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine and H. van der Vorst, *Templates for the Solution of Linear Systems: Building Blocks for Iterative Methods*. SIAM, Philadelphia, PA, 1994, <http://www.netlib.org/templates/templates.ps>.



- [3] S. Browne, J. Dongarra, N. Garner, G. Ho and P. Mucci, A portable programming interface for performance evaluation on modern processors, *International Journal of High Performance Computing Applications* **14** (Fall 2000), 189–204.
- [4] J. Dongarra, Performance of various computers using standard linear equations software, <http://www.netlib.org/benchmark/performance.ps>.
- [5] J. Dongarra and H. van der Vorst, Performance of various computers using sparse linear equations software in a fortran environment, *Supercomputer* (1992).
- [6] A. Greenbaum, *Iterative Methods for Solving Linear Systems*, SIAM, Philadelphia, 1997.
- [7] A. Greenbaum and Z. Strakos, Predicting the behavior of finite precision Lanczos and Conjugate Gradient computations, *SIAM J. Matrix Anal.* (1992), 121–137.
- [8] M.T. Jones and P.E. Plassmann, A parallel graph coloring heuristic, *SIAM J. Sci. Stat. Comput.* **14** (1993).
- [9] Y. Saad and M.H. Schultz, GMRes: a generalized minimal residual algorithm for solving nonsymmetric linear systems, *SIAM J. Sci. Stat. Comput.* **7** (1986), 856–869.
- [10] H. van der Vorst, Bi-CGSTAB: a fast and smoothly converging variant of Bi-CG for the solution of nonsymmetric linear systems, *SIAM J. Sci. Stat. Comput.* **13** (1992), 631–644.



**Hindawi**

Submit your manuscripts at  
<http://www.hindawi.com>

