

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/2613573>

A Methodology for Benchmarking Java Grande Applications

Article · July 1999

DOI: 10.1145/304065.304103 · Source: CiteSeer

CITATIONS

44

READS

51

5 authors, including:



Mark Bull

The University of Edinburgh

96 PUBLICATIONS 1,579 CITATIONS

[SEE PROFILE](#)



Lorna Smith

The University of Edinburgh

34 PUBLICATIONS 1,516 CITATIONS

[SEE PROFILE](#)



David Henty

The University of Edinburgh

84 PUBLICATIONS 1,268 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Feedback Guided Dynamic Scheduling of Nested Loops [View project](#)



HECToR e277 [View project](#)

A Methodology for Benchmarking Java Grande Applications

J. M. Bull, L. A. Smith, M. D. Westhead, D. S. Henty and R. A. Davey
EPCC, James Clerk Maxwell Building, The King's Buildings,
The University of Edinburgh, Mayfield Road, Edinburgh EH9 3JZ,
Scotland, U.K.

Abstract

Increasing interest is being shown in the use of Java for large scale or *Grande* applications. This new use of Java places specific demands on the Java execution environments that could be tested and compared using a standard benchmark suite. EPCC has taken a leading role in the Java Grande Forum work to develop a framework and methodology for such a suite. Initial results presented here show interesting differences between JVMs, demonstrating the validity of the approach. Future work will concentrate on parallel benchmarks.

1 Introduction

With the increasing ubiquity of Java comes a growing range of uses for the language that fall well outside its original design specifications. The use of Java for large scale applications with large memory, network or computational requirements, so called *Grande* applications, represent a clear example of this trend. Despite concerns about performance and numerical definitions an increasing number of users are taking seriously the possibility of using Java for Grande codes.

The Java Grande Forum (JGF) is a community initiative led by Sun and the Northeast Parallel Architectures Center (NPAC) which aims to address these issues and in so doing promote the use of Java in this area. This paper describes work carried out by EPCC in the University of Edinburgh on behalf of the JGF to initiate a benchmark suite aimed at testing aspects of Java execution environments, (JVMs, Java compilers, Java hardware etc.) pertinent to Grande Applications. The work involves constructing a framework for the benchmarks, designing an instrumentation class to ensure standard presentation of results, and seeding the suite with existing and original benchmark codes.

The aim of this work is ultimately to arrive at a standard benchmark suite which can be used to:

- Demonstrate the use of Java for Grande applications. Show that real large scale codes can be written and provide the opportunity for performance comparison against other languages.
- Provide metrics for comparing Java execution environments thus allowing Grande users to make informed decisions about which environments are most suitable for their needs.

- Expose those features of the execution environments critical to Grande Applications and in doing so encourage the development of the environments in appropriate directions.

A standard approach, ensuring that metrics and nomenclature are consistent, is important in order to facilitate meaningful comparisons in the Java Grande community. The authors are keen to invite contributions from the community to add to the benchmark suite and comments on the approach taken.

The remainder of this paper is structured as follows: Sections 2 and 3 outline the methodology we adopted in designing this suite and describe the instrumentation API. Sections 4 and 5 give the current status of the serial part of the suite, and some sample results to illustrate the existing suite in action. Section 6 outlines directions for future work, concentrating on the parallel part of the suite, and invites participation in this effort, and Section 7 provides some conclusions.

2 Methodology

In this Section we discuss the principal issues affecting the design of a benchmark suite for Java Grande applications, and describe how we have addressed these issues.

For a benchmark suite to be successful, we believe it should be:

- **Representative:** The nature of the computation in the benchmark suite should reflect the types of computation which might be expected in Java Grande applications. This implies that the benchmarks should stress Java environments in terms of CPU load, memory requirements, and I/O, network and memory bandwidths.
- **Interpretable:** As far as possible, the suite as a whole should not merely report the performance of a Java environment, but also lend some insight into why a particular level of performance was achieved.
- **Robust:** The performance of suite should not be sensitive to factors which are of little interest (for example, the size of cache memory, or the effectiveness of dead code elimination).
- **Portable:** The benchmark suite should run on as wide a variety of Java environments as possible.

- **Standardised:** The elements of the benchmark should have a common structure and a common ‘look and feel’. Performance metrics should have the same meaning across the benchmark suite.
- **Transparent:** It should be clear to anyone running the suite exactly what is being tested.

We observe that the first two of these aims (representativeness and interpretability) tend to conflict. To be representative, we would like the contents of the benchmark to be as much like real applications as possible, but the more complex the code, the harder it is to interpret the observed performance. Rather than attempt to meet both these objectives at once, we provide three types of benchmark: low-level operations (which we refer to as Section I of the suite), simple kernels (Section II) and applications (Section III). This structure is employed for both the serial and parallel parts of the suite.

The low-level operation benchmarks have been designed to test the performance of the low-level operations which will ultimately determine the performance of real applications running under the Java environment. Examples include arithmetic and maths library operations, garbage collection, method calls and casting in the serial part and ping-pong, barriers and global reductions in the parallel part. The kernel benchmarks are chosen to be short codes, each containing a type of computation likely to be found in Grande applications, such as FFTs, LU Factorisation, matrix multiplication, searching and sorting. The application benchmarks are intended to be representative of Grande applications, suitably modified for inclusion in the benchmark suite by removing any I/O and graphical components. By providing these different types of benchmark, we hope to observe the behaviour of the most complex applications and interpret that behaviour through the behaviour of the simpler codes. We also choose the kernels and applications from a range of disciplines, which are not all traditional scientific areas.

To make our suite robust, we avoid dependence on particular data sizes by offering a range of data sizes for each benchmark in Sections II and III. We also take care to defeat possible compiler optimisation of strictly unnecessary code. For Sections II and III this achieved by validating the results of each benchmark, and outputting any incorrect results. For Section I, even more care is required as the operations performed are rather simple. We note that some common tricks, used to fool compilers into thinking that results are actually required, may fail in interpreted systems where optimisations can be performed at run time.

For maximum portability, as well as ensuring adherence to standards, we have taken the decision to have no graphical component in the benchmark suite. While applets provide a convenient interface for running benchmarks on workstations and PCs, this is not true for typical supercomputers where interactive access may not be possible. Thus we restrict ourselves to simple file I/O.

For standardisation we have created a `JGFBenchMark` class to be used in all benchmark programs. This is described in detail in Section 3.

Transparency is achieved by distributing the source code for all the benchmarks. This removes any ambiguity in the question of what is being tested: we do not consider it acceptable to distribute benchmarks in Java byte code form.

3 Instrumentation

3.1 Performance metrics

We present performance metrics for the benchmarks in two forms: execution time and temporal performance. The execution time is simply the wall clock time required to execute the portion of the benchmark code which comprises the ‘interesting’ computation—initialisation, validation and I/O are excluded from the time measured. For portability reasons, we chose to use the `System.currentTimeMillis` method from the `java.lang` package. Millisecond resolution is less than ideal for measuring benchmark performance, so care must be taken that the run-time of all benchmarks is sufficiently long that clock resolution is not significant.

Temporal performance (see [3]) is defined in units of operations per second, where the operation is chosen to be the most appropriate for each individual benchmark. For example, we might choose floating point operations for a linear algebra benchmark, but this would be inappropriate for, say, a Fourier analysis benchmark which relies heavily on transcendental functions. For some benchmarks, where the choice of most appropriate unit is not obvious, we allow more than one operation unit to be defined. For example, the garbage collection benchmark reports both bytes collected per second and references collected per second.

For the low-level benchmarks (Section I) we report only temporal performance. This allows us to adjust the number of operations performed at run-time to give a suitable execution time, which is guaranteed to be much larger than the clock resolution. This overcomes the difficulty that there can be several orders of magnitude difference in performance on these benchmarks between different Java environments. For other benchmarks (Sections II and III) we report both execution time and temporal performance.

3.2 Design of instrumentation classes

Creating an instrumentation class raises some interesting issues in object-oriented design. Our objective is to be able to take an existing code and to both instrument it, and force it to conform with a common benchmark structure, with as few changes as possible.

A natural approach would be to create an abstract benchmark class which would be sub-classed by an existing class in the benchmark’s hierarchy: access to instrumentation would be via the benchmark class. However, since Java does not support multiple inheritance, this is not possible. Other options include:

- Inserting the benchmark class at some point in the existing hierarchy.
- Creating an instance of the benchmark class at some point in the existing hierarchy.
- Accessing benchmark methods as class methods.

The last option was chosen because minimal changes are required to existing code: the benchmark methods can be referred to from anywhere within existing code by a global name. However, we would like, for instance, to be able to access multiple instances of a timer object. This can be achieved by filling a hash-table with timer objects. Each timer object can be given a global name through a unique string.

We can force compliance to common structure to some extent by sub-classing the lowest level of the main hierarchy

in the benchmark, and implementing a defined `interface`, which includes a 'run' method. We can then create a separate `main` class which creates an instance of this sub-class and calls its 'run' method. It is then straightforward to create a `main` which, for example, runs all the benchmarks of a given size in a given Section.

3.3 The JGF Benchmark API

Figure 1 describes the API for the benchmark class. `addTimer` creates a new timer and assigns a name to it. The optional second argument assigns a name to the performance units to be counted by the timer. `startTimer` and `stopTimer` turn the named timer on and off. The effect of repeating this process is to accumulate the total time for which the timer was switched on. `addOpsToTimer` adds a number of operations to the timer: multiple calls are cumulative. `readTimer` returns the currently stored time. `resetTimer` resets both the time and operation count to zero. `printTimer` prints both time and performance for the named timer; `printperfTimer` prints just the performance. `storeData` and `retrieveData` allow storage and retrieval of arbitrary objects without, for example, the need for them to be passed through argument lists. This may be useful, for example, for passing iteration count data between methods without altering existing code. `printHeader` prints a standard header line, depending on the benchmark Section and data size passed to it.

Figure 2 illustrates the use of an interface to standardise the form of the benchmark. The interface for Section II is shown here; that for Section III is similar, while that for Section I is somewhat simpler.

To produce a conforming benchmark, a new class is created which `extends` the lowest class of the main hierarchy in the existing code and `implements` this interface. The `JGFrun` method should call `JGFsetsize` to set the data size, `JGFinitialise` to perform any initialisation, `JGFkernel` to run the main (timed) part of the benchmark, `JGFvalidate` to test the results for correctness, and finally `JGFtidyup` to permit garbage collection of any large objects or arrays. Calls to `JGFBenchmark` class methods can be made either from any of these methods, or from any methods in the existing code, as appropriate.

4 Current Status

Currently the parallel codes are under development. The following serial codes are available in the release version 1.0.

4.1 Section I: Low Level Operations

Arith Measures the performance of arithmetic operations (add, multiply and divide) on the primitive data types `int`, `long`, `float` and `double`. Performance units are additions, multiplications or divisions per second.

Assign Measures the cost of assigning to different types of variable. The variables may be scalars or array elements, and may be local variables, instance variables or class variables. In the cases of instance and class variables, they may belong to the same class or to a different one. Performance units are assignments per second.

Cast Tests the performance of casting between different primitive types. The types tested are `int` to `float` and

`back`, `int` to `double` and `back`, `long` to `float` and `back`, `long` to `double` and `back`. Performance units are casts per second. Note that other pairs of types could also be tested (e.g. `byte` to `int` and `back`), but these are too amenable to compiler optimisation to give meaningful results.

Garbage Assesses the performance of the system garbage collector. Objects are created with a randomly chosen size in the range zero to one thousandth of the total available memory. Initially, sufficient objects are created to consume all available memory; this part is not timed. Subsequently object creation proceeds for a fixed time period. All objects are assigned to the same reference, so that all objects except the most recently created are available for collection. The number and total size of the objects collected are recorded. Performance units are references per second and bytes per second.

Math Measures the performance of all the methods in the `java.lang.Math` class. Performance units are operations per second. Note that for a few of the methods (e.g. `exp`, `log`, inverse trig functions) the cost also includes the cost of an arithmetic operation (add or multiply). This was necessary to produce a stable iteration which will not overflow and cannot be optimised away. However, it is likely the cost of these additional operations is insignificant: if necessary the performance can be corrected by using the relevant result from the `Arith` benchmark.

Method Determines the cost of a method call. The methods can be instance, final instance or class methods, and may be called from an instance of the same class, or a different one. Performance units are calls per second. Note that final instance and class methods can be statically linked and are thus amenable to inlining. An infeasible high performance figure for these tests generally indicates that the compiler has successfully inlined these methods.

4.2 Section II: Kernels

Fourier Computes the first N Fourier coefficients of the function $f(x) = (x + 1)^x$ on the interval $0,2$. Performance units are coefficients per second. This benchmark heavily exercises transcendental and trigonometric functions.

LUFact Solves an $N \times N$ linear system using LU factorisation followed by a triangular solve. This is a Java version of the well known Linpack benchmark [2]. Performance units are Mflops per second. Memory and floating point intensive.

Search Solves a game of connect-4 on a 6×7 board using an alpha-beta pruned search technique. The problem size is determined by the initial position from which the game is analysed. The number of positions evaluated, N , is recorded, and the performance reported in units of positions per second. Memory and integer intensive.

HeapSort Sorts an array of N integers using a heap sort algorithm. Performance is reported in units of items per second. Memory and integer intensive.

```

public class JGFBenchMark{
// No constructor
// Class methods
    public static synchronized void addTimer(String name);
    public static synchronized void addTimer(String name, String opname);
    public static synchronized void startTimer(String name);
    public static synchronized void stopTimer(String name);
    public static synchronized void addOpsToTimer(String name, double count);
    public static synchronized double readTimer(String name);
    public static synchronized void resetTimer(String name);
    public static synchronized void printTimer(String name);
    public static synchronized void printperfTimer(String name);
    public static synchronized void storeData(String name, Object obj);
    public static synchronized void retrieveData(String name, Object obj);
    public static synchronized void printHeader(int section, int size);
}

```

Figure 1: API for the JGFBenchMark class

```

public interface JGFSection2 {
    public void JGFsetsize(int size);
    public void JGFinitialise();
    public void JGFkernel();
    public void JGFvalidate();
    public void JGFtidyup();
    public void JGFrun(int size);
}

```

Figure 2: Interface definition for Section II

Crypt Performs IDEA (International Data Encryption Algorithm [4]) encryption and decryption on an array of N bytes. Performance units are bytes per second. Bit/byte operation intensive.

4.3 Section III: Applications

Euler Solves the time-dependent Euler equations for flow in a channel with a “bump” on one of the walls. A structured, irregular, $N \times 4N$ mesh is employed, and the solution method is a finite volume scheme using a fourth order Runge-Kutta method with both second and fourth order damping. The solution is iterated for 200 timesteps. Performance is reported in units of timesteps per second.

5 Sample Results

The benchmark suite has been run on a number of different execution environments on two different hardware platforms. The following JVMs/compilers have been tested on a 200MHz Pentium Pro with 256 Mb of RAM running Windows NT:

- Sun JDK Version 1.1.7 (with and without JIT)
- Sun JDK Version 1.2beta4
- IBM HPJ Version a12h

- Microsoft SDK Version 3.1

The following JVMs have also been tested on a 250MHz Sun Ultra Enterprise 3000 with 1Gb of RAM running Solaris 2.6:

- Sun JDK Version 1.1.7
- Sun JDK Version 1.2beta4

5.1 Programming language comparison

The benchmark suite has been developed to allow the performance of various execution environments on different hardware platforms to be tested. Also of interest is language comparisons, comparing the performance of Java versus other programming languages such as Fortran, C and C++. Currently, only the LU factorisation benchmark, which is based on the Linpack benchmark has allowed programming language comparison. It is intended, however, that the parallel part of the suite will contain versions of well-known Fortran and C parallel benchmarks, thus facilitating further inter-language comparisons.

Measurements have been taken for the Linpack Benchmark on a 1000 x 1000 problem size, using Java (the Sun JDK1.2b4), Fortran and C on a 250MHz Sun Ultra Enterprise 3000 with 1Gb of RAM and the results are shown in figure 3. The results demonstrate that the performance of the Sun JDK1.2b4 is within a factor of two to three of compiled code.

	Sun JDK1.1.7	Sun JDK1.2beta4	IBM HPJ	Microsoft SDK
Assign Benchmark (assignments/s)				
Same:Scalar:Local	8.83E+08	8.83E+08	6.88E+08	4.89E+06
Same:Scalar:Instance	5.74E+08	5.67E+08	7.05E+08	4.04E+06
Same:Scalar:Class	6.35E+08	6.36E+08	6.93E+08	1.90E+06
Method Benchmark (calls/s)				
Same:Instance	4.16E+07	4.16E+07	2.62E+07	2.73E+07
Same:FinalInstance	1.58E+09	1.58E+09	Infinity	2.94E+07
Same:Class	1.58E+09	1.58E+09	Infinity	3.40E+07
Math Benchmark (operations/s)				
SinDouble	2.56E+06	2.58E+06	9.42E+05	2.02E+06
CosDouble	1.88E+06	1.88E+06	8.36E+05	1.71E+06

Table 1: Performance figures for selected Section I benchmarks

Size	Sun JDK1.1.7		Sun JDK1.2beta4		IBM HPJ		Microsoft SDK	
	Time	Performance	Time	Performance	Time	Performance	Time	Performance
Fourier Benchmark (coefficients/s)								
A	9.13	2192	8.86	2257	24.7	808.5	11.1	1800.3
B	90.8	2202	89.2	2242	253	791.6	111	1797.0
C	913	2190	896	2233	4550	439.5	1120	1791.8
Crypt Benchmark (Kbyte/s)								
A	8.63	695.7	8.38	716.4	9.94	603.8	15.8	379.1
B	57.2	698.9	55.5	720.9	66.3	603.8	105	379.6
C	143	699.8	139	721.6	166	603.5	264	379.0

Table 2: Performance figures for selected Section II benchmarks. Sizes A, B and C refer to small, medium and large values of N .

5.2 JVM/compiler comparison

The performance of various execution environments has been measured. Four JVMs/compiler have been tested on a 200MHz Pentium Pro with 256 Mb of RAM: Sun JDK1.1.7 (with JIT); Sun JDK1.2b4 (with JIT); IBM HPJ1.2h; Microsoft SDK for Java 3.1.

For the low level operations, performance variation between the environments is observed for the majority of the benchmarks and no particular one performed consistently better across the section. For example, the Sun JDKs and the IBM compiler show a similar performance for **Assign** and both outperform the Microsoft SDK. (See Table 1 and Figure 3 which show the performance of assigning to scalar elements of the same class.)

As mentioned previously, some method calls in **Method** can be statically linked and are amenable to inlining. This is observed for the IBM HPJ whereby inlining has successfully occurred and an (infeasibly) high performance figure for these tests has been generated. Table 1 shows the timings for method calls in the same class.

The kernel section of the benchmark suite demonstrates that the Sun JVMs perform more consistently across the benchmark range, with both the IBM HPJ and Microsoft SDK systems showing poorer performance on one specific benchmark. The IBM HPJ performs considerably less well than the other environments for **Fourier** (see Table 2 and Figure 4). This benchmark uses a number of transcendental functions and **Math** demonstrates that the IBM HPJ compiler has a lower performance than the other Java environments for these functions (see Table 1 which gives the performance of sin and cos). The Microsoft SDK has a lower

performance for **Crypt** than other environments (see Figure 4 and Table 2).

This small sample of results demonstrates that the benchmark suite is achieving the objective of highlighting and explaining performance differences between Java environments.

6 Future Work

All the benchmarks presented in this paper are serial and only utilise a single processor. However, one of our major aims is to increase the scope of the suite to include benchmarks that measure how efficiently Java codes and environments can utilise multi-processor architectures. Parallel machines are obvious target platforms for our work since they have the potential to supply the large amounts of CPU power and memory required by Grande applications. Understanding Java performance on a single processor is a necessary precursor to this work, and having developed the serial benchmark we will now start to include parallel benchmarks into the suite.

There are issues to be considered when developing any parallel benchmark, even when targeting well established languages such as C or Fortran. In the past, portability was a major issue, but with the almost universal adoption of the MPI standard this is no longer a problem. However, MPI provides many possible methods for even the most simple task of a single point-to-point communication and the performance of the different modes can vary dramatically on different platforms. A generic parallel benchmark using MPI must either attempt to measure as many different

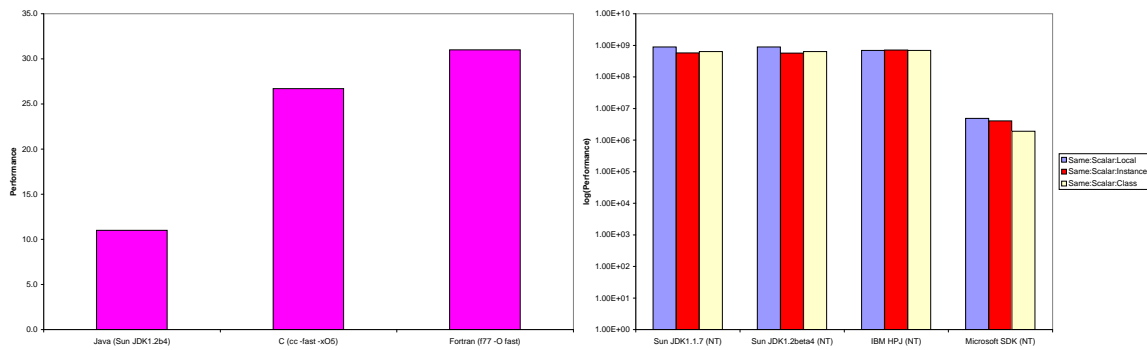


Figure 3: (a) Language comparison of the Linpack benchmark (results are given in Mflop/s) and (b) Comparison of the performance of the four Java environments for assigning scalar elements to the same class (in assignments/s)

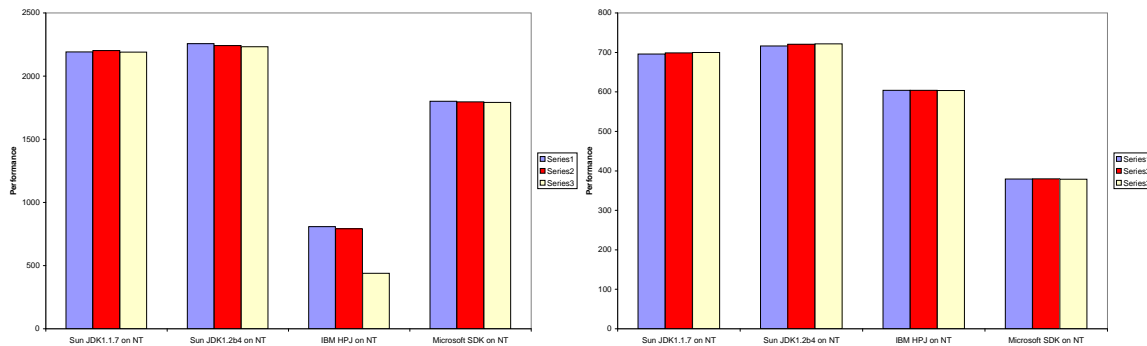


Figure 4: Comparison of the performance of the four Java environments for (a) **Fourier** (in coefficients/s) and (b) **Crypt** (in Kbyte/s). Series 1, 2 and 3 refer to small, medium and large values of N .

modes as possible (with the aim of allowing the programmer to choose the most efficient method) or simply time what is considered to be the most commonly used method (in an attempt to predict the performance of “typical” applications). This problem has lead to different disciplines developing their own parallel benchmarks designed to test the communication methods and patterns most commonly used in specific applications areas. For examples see OCCOMM [1] (ocean modelling) and SSB [5] (social sciences).

Developing a useful parallel Java benchmark is even more problematic since it is not yet clear which of the many parallelisation models to target. A parallel Java code might use built-in Java methods such as threads or RMI, Java library routines such as a Java implementation of MPI, or a Java interface to a platform-specific MPI implementation written in C. For the built-in methods, parallel performance may depend critically on the details of the JVM implementation, for example how effectively threads are scheduled across multiple processors.

Consideration of these issues has lead us to decide on the following strategy:

- Low-level benchmarks will be written to test the performance of the fundamental operations of the various parallel models, e.g. the overhead of creating a new thread, or the latency and bandwidth of message-passing using MPI.
- Kernel benchmarks will be written that implement one

or more common communications patterns using a variety of parallel models.

- We will collect a set of genuine parallel applications. Each application is likely to be only available for a single parallel model. However, we hope that we will have applications using a wide variety of models. All parallel applications will also be available in serial form.

We are very keen to get as wide a range of parallel Java applications as possible. We encourage anyone who has an existing parallel Java code to contact epcc-javagrande@epcc.ed.ac.uk so that we can consider it for instrumentation and inclusion in the benchmark. All code authors are acknowledged in the benchmark suite.

The current suite, instrumentation classes, and a more comprehensive set of results, are available at <http://www.epcc.ed.ac.uk/research/javagrande/>. We would strongly welcome use of, and comments on, this material from developers both of Grande applications and of Grande environments.

7 Conclusions

We have presented a methodology for, and initial implementation of, a suite of benchmarks for Java Grande applications. We have set out criteria which we believe such a suite should meet, and have demonstrated how we have

met these criteria. We have discussed methods of benchmark instrumentation, and have presented the instrumentation API. Sample results show that the suite gives some useful and meaningful insight into the performance of Java environments. Finally, we have discussed the future of the benchmark suite in terms of parallel benchmarks.

Acknowledgements

We wish to thank the following who have contributed benchmarks to the suite: Gabriel Zachmann of the Fraunhofer Institute for Computer Graphics (Fourier, HeapSort and Crypt), Reed Wade of the University of Tennessee at Knoxville (LUFact), John Tromp of CWI, Amsterdam (Search) and David Oh, formerly of MIT's Computational Aerospace Sciences Laboratory (Euler).

This work was funded by the UK High Performance Computing Initiative, EPSRC grant GR/K42653.

References

- [1] Ashworth, M. (1996) *OCCOMM Benchmarking Guide. Version 1.2*, Daresbury Laboratory Technical Report, available from <http://www.dl.ac.uk/TCSC/CompEng/OCCOMM/uguide-1.2.ps>.
- [2] Dongarra, J. J. (1998) *Performance of Various Computers Using Standard Linear Equations Software (Linpack Benchmark Report)*, University of Tennessee Computer Science Technical Report, CS-89-85.
- [3] Hockney, R. W. (1992) *A Framework for Benchmark Performance Analysis*, Supercomputer, vol. 48, no. IX(2), pp. 9-22.
- [4] Lai, X., J. L. Massey and S. Murphy (1992) *Markov ciphers and differential cryptanalysis*, in *Advances in Cryptology—Eurocrypt '91*, pp. 17-38, Springer-Verlag.
- [5] Openshaw S., and J. Schmidt (1997) *A Social Science Benchmark (SSB/1) Code for Serial, Vector and Parallel Supercomputers*, Geographical and Environment Modelling, vol. 1, no. 1, pp. 65-82.