

Scientific computations on multi-core systems using different programming frameworks

Panagiotis D. Michailidis^a, Konstantinos G. Margaritis^b

^a*Department of Balkan, Slavic and Oriental Studies, University of Macedonia
156 Egnatia Str., 54636 Thessaloniki, Greece*

E-mail: pmichailidis@uom.gr

^b*Department of Applied Informatics, University of Macedonia*

156 Egnatia Str., 54636 Thessaloniki, Greece

E-mail: kmarg@uom.gr

Abstract

Numerical linear algebra is one of the most important forms of scientific computation. The basic computations in numerical linear algebra are matrix computations and linear systems solution. These computations are used as kernels in many computational problems. This study demonstrates the parallelisation of these scientific computations using multi-core programming frameworks. Specifically, the frameworks examined here are Pthreads, OpenMP, Intel Cilk Plus, Intel TBB, SWARM, and FastFlow. A unified and exploratory performance evaluation and a qualitative study of these frameworks are also presented for parallel scientific computations with several parameters. The OpenMP and SWARM models produce good results running in parallel with compiler optimisation when implementing matrix operations at large and medium scales, whereas the remaining models do not perform as well for some matrix operations. The qualitative results show that the OpenMP, Cilk Plus, TBB, and SWARM frameworks require minimal programming effort, whereas the other models require advanced programming skills and experience. Finally, based on an extended study, general conclusions regarding the programming models and matrix operations for some parameters were obtained.

Keywords: Scientific Computations, Linear Algebra, Parallel Computing, Multi-core, Parallel Programming

1. Introduction

Scientific computing is a collection of quantitative methods and tools used to develop and solve mathematical models of a variety of scientific problems using a computer system [66]. Numerical linear algebra is one of the most important quantitative methods for scientific computing. The basic computations of numerical linear algebra are matrix computations (such as vector and matrix addition, dot product, outer product, matrix transpose, matrix-vector product, and matrix product) and solutions of linear systems (such as the direct Gaussian elimination method and the iterative Jacobi method), which are used as kernels in many computational problems such as computational statistics [34] and combinatorial optimisation [41, 68, 70]. To satisfy the heavy computational requirements of these methods, large-scale matrix computations use towers of software infrastructure such as BLAS/LAPACK [14], ScaLAPACK [20], and PLAPACK [69] running on cluster computing platforms. Recently, modern high-performance computer systems have been introduced, which can be classified into three categories. The first, multi-core platforms, integrate a few cores (from two to ten) on the same integrated circuit chip (die) in an effort to speed up execution of computationally intensive methods. The second, many-core platforms or general-purpose graphics processing units (GPUs), which consist of a large number of cores (as many as several hundred), are specifically oriented to maximizing execution throughput for parallel applications [29]. The third, reconfigurable platforms based on field-programmable gate arrays, are becoming important, especially when higher performance/power computation ratios are desired. Along with the introduction of these platforms, software projects such as Parallel Linear Algebra for Multi-core Architectures (PLASMA) [9] have been developed for multi-core machines and Matrix Algebra on GPU and Multicore Architectures (MAGMA) [9] for GPU platforms. Recent comparisons of these platforms have shown substantial architectural and performance differences for several application areas [18, 36, 40, 45, 71]. Therefore, this paper focusses on the

21 multi-core approach because it is the easiest way to speed up an application from the perspective of non-specialised
22 developer. However, sequential implementation of linear algebra computations on a multi-core architecture will not
23 improve performance because it cannot exploit the other cores that are available. Successful implementation of scientific
24 computations on a multi-core machine can be achieved only through parallel (or multi-core) programming.

25 The fundamental programming question on a multi-core platform is how to decompose a problem into several
26 sub-problems and how to map these to cores with the goal of increasing performance. Moreover, programmers need
27 to study and understand the hardware characteristics of the multi-core platform to write efficient parallel programs.
28 For this reason, programming on multi-core processors is a more complex procedure than programming on sequential
29 processors because application data in memory can be accessed by several entities called threads, which belong to
30 the same program. For this reason some synchronisation between threads is necessary. Therefore, there is a need to
31 bridge the gap between hardware and software applications to hide the hardware details from the programmers and
32 enable them to write parallel programs with minimal programming effort. This has resulted in the introduction of
33 several parallel programming models [29, 43] which simplify the parallelisation of linear algebra computations and
34 other related applications on multi-core computers. However, these models differ significantly in their parallel design
35 principles, abstraction levels, semantics, and syntax. Some popular models are POSIX threads (Pthreads for short)
36 [25], OpenMP [3], Intel Cilk Plus [1], Intel Threading Building Blocks (TBB for short) [2], SoftWare and Algorithms
37 for Running on Multi-core (SWARM for short) [16], and FastFlow [10, 11]. These models are based on a small set of
38 extensions to the C programming language and involve a relatively simple compilation phase and a potentially much
39 more complex runtime system.

40 Based on the various features and qualities of programming models, the question posed by programmers becomes:
41 what is the appropriate parallel programming framework for implementing linear algebra computations on a multi-
42 core system to achieve a balance between high programmer productivity (i.e., minimal programming effort) and high
43 performance? Therefore, the contribution of this paper is a unified and systematic quantitative (i.e., performance-
44 based) and qualitative (i.e., related to the ease of programming effort) comparison of all multi-core programming
45 frameworks for implementing linear algebra computations based on simple parallelisation techniques. Finally, the
46 authors believe that this work is important and interesting because this comparison may turn out to be very helpful for
47 other programmers and scientists who are often faced with a variety of options for implementing projects.

48 The rest of the paper is organised as follows. In Section 2, related work is discussed. In Section 3, an abstract
49 multi-core system architecture is presented along with all the reviewed parallel programming frameworks, and in Sec-
50 tion 4, parallelisation issues in implementing linear algebra computations are considered. In Section 5, a performance
51 and qualitative evaluation of the reviewed parallel programming models for parallelising linear algebra operations is
52 described. Finally, Section 6 presents conclusions.

53 2. Related work

54 In the research literature, several studies have evaluated various parallel programming models on multi-core plat-
55 forms. Most of this research has compared parallel programming models from the performance point of view. How-
56 ever, there has been little related work on comparing parallel programming models from the qualitative or productivity
57 points of view.

58 The research studies based on performance evaluation of different multi-core programming models can be organ-
59 ised into four groups:

- 60 1. *Evaluation of a programming model to parallelise a specific problem.* Research studies [51] and [53] evaluated
61 the parallelisation performance of the Gaussian elimination and LU factorisation algorithms using the OpenMP
62 programming model. The parallelisation of Gaussian elimination [51] was based on the data parallel approach,
63 whereas the parallel implementation of the LU decomposition [53] was based on the pipeline approach. Zucker-
64 man *et al.* [72] evaluated the performance of the parallel matrix multiplication kernel using a high-performance
65 M:N threading library, Microthread, and showed its efficiency with regard to the well-known Intel Math Kernel
66 Library (IMKL). Runger *et al.* [63] presented a parallel optimised library for dense matrix multiplication on a
67 multi-core platform. This implementation was based on a recursive approach and was compared with efficient
68 libraries such as GotoBLAS, IMKL, and AMD Core Math Library (AMCL). Finally, Michailidis *et al.* [55]
69 studied the performance of the pipeline approach in the OpenMP model for parallelisation of the Gauss-Jordan

70 algorithm for solving systems of linear equations. The performance results were compared to the performance
71 of two other naive parallel approaches, row block and row cyclic distribution.

- 72 2. *Evaluation of a programming model to parallelise a set of problems.* Michailidis *et al.* [54] presented perfor-
73 mance results for the OpenMP model in the parallelisation of two important matrix computations, matrix-vector
74 product and matrix multiplication. These parallelisations were based on a simple parallel data technique. But-
75 tari *et al.* [23, 24] developed the PLASMA library for implementing certain linear algebra operations on a
76 multi-core platform using the tile algorithms.
- 77 3. *Evaluation of a few programming models to parallelise a specific problem.* Li [47] evaluated the performance
78 of a high-performance sparse LU factorisation algorithm on several representative multi-core machines using
79 the Pthreads and MPI programming models. Aldinucci *et al.* [12] demonstrated a parallel implementation of
80 the Smith-Waterman bio-informatics algorithm using Cilk, OpenMP, Intel TBB, and FastFlow. They found
81 experimentally that the FastFlow model performed better than the other three models. Finally, Kegel *et al.*
82 [44] analysed and compared three programming models, Pthreads, OpenMP, and Intel TBB, for parallelising a
83 real-world medical imaging application.
- 84 4. *Evaluation of a few programming models to parallelise a small set of problems.* Marowka [49] presented a
85 comparative study of OpenMP and TBB micro-benchmarks to study parallelisation overheads on a dual-core
86 machine with two different compilers, the Intel compiler and Microsoft Visual Studio C++. The design of
87 the TBB micro-benchmark followed the design methodology of the OpenMP EPCC micro-benchmarks [22].
88 Ayguade *et al.* [15] analysed the performance of six algorithms, including the sparseLU and Strassen algo-
89 rithms, using OpenMP, Intel task queues, and Intel Cilk programming models. Kurzak *et al.* [46] implemented
90 matrix factorisation algorithms such as Cholesky factorisation, QR factorisation, and LU factorisation on multi-
91 core processors using multi-thread programming tools such as Cilk and SMPs. These implementations were
92 based on tiling algorithms. Podobas *et al.* [59] presented a quantitative study comparing three task-based
93 parallel programming models, OpenMP, Intel Cilk Plus, and Wool [31], for a small set of micro-benchmarks
94 including the SparseLU and Strassen algorithms. This performance evaluation took into consideration the cost
95 of creating and joining tasks. Ravela *et al.* [61] studied the performance of the Pthreads, OpenMP, Intel Cilk
96 Plus, and Intel TBB programming models for a set of scientific computing benchmarks such as matrix multi-
97 plication, Jacobi iteration, and Laplace heat distribution. Furthermore, this study considered the development
98 time required to parallelise these benchmarks. Sanchez *et al.* [65] undertook a comparative evaluation of the
99 OpenMP, Intel TBB, Intel ArBB, and CUDA models for parallel implementations of four problems (including
100 matrix multiplication) on five shared memory platforms. Note that the Intel ArBB model combines multi-core
101 parallel threads and SIMD features with a simpler programming model, whereas the CUDA model exploits the
102 SIMD features of the GPU hardware. Finally, Ali *et al.* [13] presented a performance study of OpenMP, Intel
103 TBB, and OpenCL for five benchmark applications (including matrix multiplication and LU decomposition) on
104 a CPU multi-core platform. This study took into consideration the compiler optimisation, overhead, and scal-
105 ability of these frameworks. The OpenCL framework is a standard API (Application Programming Interface),
106 which focusses on programming for multi-core processors except for GPU hardware.

107 Note that other research papers have compared various parallel programming frameworks on a multi-core platform
108 for several applications. For example, [52] compared OpenMP, Cilk Plus, Intel TBB, and OpenCL for implementing a
109 2D/3D image registration algorithm. Furthermore, Shekhar *et al.* [27] compared three popular programming models
110 (OpenMP, GCD, and Pthreads) for parallelising face detection and automatic speech recognition. However, few
111 studies have focussed on a qualitative comparison of parallel programming models. Earlier works such as [42, 67]
112 evaluated programming effort and usability for two parallel programming models, OpenMP and MPI. In addition,
113 Ali *et al.* [13] presented a qualitative comparison of three models, OpenMP, Intel TBB, and OpenCL, for a set of
114 applications including matrix multiplication. This comparison was based on two software engineering criteria: time
115 in weeks for learning and mapping each application into the corresponding framework by the user, and the number of
116 lines of code for each parallelised application using these models. Finally, Kegel *et al.* [44] analysed the number of
117 lines of code for a parallelised medical imaging application using the Pthreads, OpenMP, and TBB models.

118 Other studies [56, 57, 58] have presented a preliminary experimental and qualitative study of programming models
119 and matrix operations for a limited set of performance parameters. These studies looked at performance as a function
120 of the number of cores without considering several compiler optimisation options or problem scalability, whereas the

121 qualitative study was limited to an analysis of the number of lines of code required to parallelise matrix operations.

Type of comparison - Method	1:1	1:F	F:1	F:F
Performance evaluation	+++	+++	+++	+++
Qualitative evaluation	---	---	---	+
Performance and Qualitative eval.	---	---	+	+

Table 1: Summary of related work for given different evaluation techniques of multi-core programming models for three types of comparison [Note: left of the (:) means one (1) or few (F) programming models and right means one (1) or few (F) problems, Notation of scale: +++ means relatively high research activity, --- means relatively low research activity]

122 Table 1 compares the research studies for the four different evaluation techniques for multi-core programming
 123 models and considers whether each is a performance evaluation, a qualitative evaluation, or both. From the table, it is
 124 clear that no systematic and unified study or comparison can be made from the quantitative or qualitative point of view.
 125 Furthermore, these research studies compared only a limited number programming models, i.e., two or three, for a
 126 limited set of scientific computing problems. Moreover, most studies were either quantitative or qualitative, although
 127 a few have combined these two types of comparison, but at limited depth. For this reason, the present paper differs
 128 from previous research because it presents a unified and extensive quantitative and qualitative comparison of a larger
 129 number of multi-core programming models for a larger set of linear algebra applications. The quantitative portion
 130 of this work considers performance with compiler optimisation and problem scalability. Furthermore, the qualitative
 131 comparison of models in this paper extends to other software engineering criteria, except for the number of lines of
 132 code, as discussed later.

133 3. Multi-core architecture and programming models

134 This section presents an abstract computational model for a multi-core architecture and a short review of the
 135 multi-core programming models evaluated in this paper.

136 A typical multi-core system ([50], p. 45) consists of a finite number of identical processing cores integrated on a
 137 single chip. Each processing core has its own private L1 cache memory and shares the L2 cache memory with other
 138 cores. In such a design, the bandwidth between the L2 cache and main memory is shared by all the processing cores
 139 [16, 50]. Multi-core examples are the IBM Power8 architecture [5], the Intel Xeon processors [7], and the AMD
 140 Opteron/AMD FX family [4].

141 The parallel programming model is used to develop parallel applications on the multi-core system. This model
 142 is an abstraction of the multi-core computer system architecture [29]. The following describes the main parallel
 143 programming models using a small set of features: implementation of the programming model, workload partitioning,
 144 synchronisation, and parallel programming constructs or patterns. The first three features were selected to constitute
 145 a full list of criteria [43], and the fourth feature was inspired by [50]. The programming model implementation is
 146 concerned with how parallelism capabilities are made available to programmers, for example, through thread libraries,
 147 template libraries, compiler directives, and language extensions. Workload partitioning is concerned with how data or
 148 workloads are partitioned into small chunks. Workload partitioning can be done either explicitly (data are partitioned
 149 manually by programmers) or implicitly (the programmers specify only which parts of the code can be processed in
 150 parallel, and the partitioning is realised automatically by the programming model). Synchronisation is concerned with
 151 the time sequence in which threads or tasks access shared data. Synchronisation can be done either explicitly (the
 152 programmers are responsible for writing code for thread management to provide conflict-free access to shared data)
 153 or implicitly (no programming effort is required from the programmers, and synchronisation constructs are available).

154 The fourth feature of parallel programming constructs or patterns involves those patterns that are valuable building
 155 blocks and how they are used in the design and implementation of parallel algorithms. Common and representative
 156 parallel patterns used in parallel programming are the fork-join pattern, the map pattern, the reduction pattern, the scan
 157 pattern, and the pipeline pattern [50]. In the fork-join pattern, the master thread creates several multiple threads that
 158 execute concurrently until they synchronise with the master thread at the end of execution. The map pattern applies a
 159 function or operation concurrently to every element of a data structure or a set of data structures with the same shape
 160 and usually creates a new data structure (or set of data structures) with the same shape as the input. The reduction

161 pattern applies an associative operation to all the elements of a data structure, reducing it to a single element. The
162 scan pattern calculates all partial reductions of a data structure. The pipeline pattern consists of a linear sequence of
163 tasks or threads and works as a producer-consumer relationship. Data flow through the pipeline from the first task to
164 the last task. Each task implements one stage of the pipeline and performs a transformation on the output data of its
165 predecessor. All these patterns are supported directly in one or more of the programming models discussed in this
166 paper, or they can be implemented using other features if the patterns are not supported directly. The "from the inside"
167 implementation details of these patterns are discussed in [50].

168 3.1. Multi-core programming frameworks

169 In this paper, the representative multi-core programming frameworks described below were included for three
170 reasons. First, these frameworks were selected to cover a variety of options for most features of programming models,
171 as discussed above. Second, these programming models were selected to cover a range of levels of programming
172 abstraction for parallelism. This range of levels of abstraction helps to evaluate the programming effort required by
173 each programming framework in a fair way. Finally, only frameworks that are used frequently by many programmers,
174 especially for small multi-core platforms rather than many-core platforms, were considered. This approach ensured
175 that quantitative and qualitative comparison of the multi-core programming frameworks was homogeneous.

176 Pthreads [25] is a common portable API and low-level thread library. The Pthreads library provides explicit
177 low-level functions for creating and destroying threads and for coordinating threads using several synchronisation
178 mechanisms. Finally, this model directly supports the fork-join programming construct.

179 OpenMP [3] is a very popular and portable API which uses implicit compiler directives to define parallelism
180 or implicit data partitioning of *for* loops, work sharing, and synchronisation. OpenMP also provides some library
181 functions for accessing the runtime environment. OpenMP supports the map and reduction patterns directly using
182 compiler directives.

183 The Intel Cilk Plus [1] language is based on technology from Cilk [21], a parallel programming model for the C
184 language. The Cilk Plus language provides keywords for implicit data partitioning or loop parallelism and implicit
185 synchronisation. Finally, Cilk Plus supports the fork-join, map, and reduction patterns directly using keywords.

186 Intel TBB [2] is an open-source template library that offers a rich methodology for expressing regular and irregular
187 parallelism in C++ programs. This template library consists of data structures or containers and algorithms that
188 abstract the use of native threading packages in which individual execution threads are created, synchronised, and
189 terminated manually. TBB emphasises an implicit data parallel programming model, enabling multiple threads to
190 work on different parts of a data collection and enabling scalability to many cores. Moreover, TBB provides direct
191 template functions for a variety of parallel patterns, including fork-join, map, reduction, scan, and pipeline.

192 SWARM [16] is an open-source parallel programming library. This library provides basic primitives for paralleli-
193 sation, restricting control of threads, allocation and de-allocation of shared memory, and communication primitives
194 for synchronisation, replication, and broadcasting. SWARM also supports the map, reduction, and scan patterns using
195 implicit primitives.

196 FastFlow [10, 11] is an open-source, C++ template-based parallel programming framework for developing effi-
197 cient applications for multi-core computers. FastFlow is conceptually designed as a stack of layers that progressively
198 abstract the shared memory parallelism at the core level up to the definition programming constructs supported by
199 structured parallel programming on shared-memory multi-core platforms. FastFlow provides programmers with a set
200 of patterns implemented as C++ templates: farm, farm with feedback, and pipeline patterns, as well as their arbi-
201 trary nesting and composition. On the other hand, the map pattern is not supported directly in FastFlow, but can be
202 implemented using the task-farm parallel pattern, as will be discussed later.

203 Table 2 summarises the features of the multi-core programming models discussed in this paper.

204 4. Parallel numerical and scientific computations

205 This section presents a set of basic scientific computations from numerical linear algebra (matrix computations
206 and solving linear systems) and discusses general guidelines for parallelising these operations using parallel patterns
207 irrespective of the parallel programming model used. Finally, several parallel programming models (as mentioned in
208 Section 3) for parallelising scientific computations are also discussed.

Feature - Model	Pthreads	OpenMP	Cilk Plus	TBB	SWARM	FastFlow
Implementation	Library	Compiler	Language	Library	Library	Library
Data Partitioning	Explicit	Implicit	Implicit	Implicit	Implicit	Explicit
Synchronization	Explicit	Implicit	Implicit	Implicit	Implicit	Explicit
Parallel patterns	fork-join	fork-join map reduction	fork-join map reduction	fork-join map reduction scan pipe	map reduction scan	farm pipe

Table 2: Features of multi-core programming models

209 4.1. Matrix computations and related methods

210 Table 3 presents the matrix computations and related methods used to evaluate the multi-core programming frame-
211 works. Note that the size of the vectors is n elements, whereas the size of the matrices is $n \times n$ elements. In addition,
212 computation is measured by the number of floating-point operations, whereas memory is measured by the number
213 of memory references (load/store operations). Furthermore, standard algorithms are used to measure computation
214 operations and memory references for matrix computations, as shown in Table 3.

Matrix computations			
Computational kernel	Operation	Computation	Memory
Vector addition	$c_i = a_i + b_i, 1 \leq i \leq n$	n	$3n$
Inner or dot product	$z = x^T \cdot y = \sum_{i=1}^n x_i \cdot y_i$	$2n$	$2n + 1$
Matrix addition	$c_{ij} = a_{ij} + b_{ij}, 1 \leq i, j \leq n$	n^2	$3n^2$
Outer product	$z_{ij} = x \cdot y^T = x_i \cdot y_j, 1 \leq i, j \leq n$	n^2	$n^2 + 2n$
Matrix - vector product	$y_i = \sum_{j=1}^n a_{ij}x_j, 1 \leq i \leq n$	$2n^2$	$n^2 + 3n$
Matrix transpose	$a_{ji}^T = a_{ij}, 1 \leq i, j \leq n$	n^2	$2n^2$
Matrix product	$c_{ij} = \sum_{k=1}^n a_{ik} \cdot b_{kj}, 1 \leq i, j \leq n$	$2n^3$	$3n^2$
Solving linear systems			
Computational kernel	Operation	Computation	Memory
Gaussian elimination	division: $l_{ik} = a_{ik}^k / a_{kk}^k, k + 1 \leq i \leq n$ elimination: $a_{ij}^{k+1} = a_{ij}^k - l_{ik}a_{kj}^k$ $b_i^{k+1} = b_i^k - l_{ik}b_k^k, k \leq i, j \leq n$	$\frac{2}{3}n^3$	$n^3 + n^2$
Jacobi method	$x_i^{(k)} = \frac{1}{a_{ii}}(b_i - \sum_{j \neq i} a_{ij}x_j^{(k-1)}), 1 \leq i \leq n$	$4n^2$	$k(n^2 + 3n)$

Table 3: Some matrix computations and related methods

215 To achieve a fair comparison of multi-core programming models for implementing linear algebra computations,
216 the matrix operations and related methods listed in the table were chosen based on the following criteria. First,
217 these computations are used by most computational scientists as representative kernels in many important application
218 domains such as computational statistics [34], combinatorial optimisation [41, 68, 70], and machine learning [26, 39].
219 Second, these computations were included to cover a range of ratios of computations to memory references, from
220 low to high. Third, these computations cover different programming levels from easy to difficult from the coding
221 perspective to evaluate model programmability. Finally, given that six programming models were studied, it was
222 important to limit the amount of time spent on each method and operation.

223 Except from the standard algorithms, there are alternative methods for calculating a matrix product, which require
224 fewer arithmetic operations by using the divide and conquer technique. Strassen's method and one of its variants, that
225 of Winograd and Coppersmith, were the first of these algorithms and have an asymptotic complexity of $O(n^{2.807})$ and
226 $O(n^{2.376})$ respectively [28, 35]. Furthermore, interesting research has been performed on cache-oblivious algorithms,
227 which analyse cache complexity for linear algebra operations like matrix transpose, matrix product, and Gaussian

228 elimination [32, 33]. Finally, there are matrix-vector product algorithms, which take $O(n \log n)$ time using the fast
229 Fourier transform (FFT) family of algorithms [35].

230 Matrix operations such as matrix transpose, matrix product, and Gaussian elimination can be implemented ef-
231 fectively on a computer with caches using tiling (or blocking) techniques [35] to increase memory access efficiency.
232 These tiling techniques have been implemented for some linear algebra operations in the popular LAPACK [14] and
233 PLASMA software packages [9]. Another optimised BLAS library, GotoBLAS [37], elaborately identifies tiles in
234 a matrix. GotoBLAS optimises matrix multiplication and increases memory access efficiency further by copying
235 the most frequently used data into contiguous memory locations to reduce translation look-aside buffer (TLB) table
236 misses, an issue that results in significant performance degradation, but is generally not addressed by other BLAS
237 routines. Finally, there are Intel Math Kernel Library (MKL) versions of BLAS and LAPACK [6], which use the Intel
238 architecture's SSE instruction extensions heavily to perform matrix computations in the SIMD mode.

239 Likewise, there are two categories of methods for solving linear systems: direct methods (i.e., Gaussian elim-
240 ination and the Gauss-Jordan method) [60] and iterative methods (i.e., Jacobi, Gauss-Seidel, and successive over-
241 relaxation) [17]. The choice of solution method for a linear system depends on its structure and size. Roughly
242 speaking, direct methods are best for small and full (dense) matrices, whereas iterative methods are best for very large
243 and sparse matrices. Furthermore, some important pivoting strategies are used in Gaussian elimination, such as full
244 pivoting, column pivoting, row pivoting, and no pivoting [60]. For some classes of matrices, such as diagonally dom-
245 inant matrices, no pivoting is required. For this research, the commonly used method for each category was selected,
246 i.e., the direct Gaussian elimination method with no pivoting and the Jacobi iterative method.

247 This paper focusses on standard algorithms for matrix computations and related methods, as shown in Table 3,
248 for the following reasons. First, standard algorithms were selected to keep the same algorithmic and implementation
249 style for all multi-core programming models by the greatest number of programmers. In this way, a fair qualitative
250 comparison of the programming models implementing advanced algorithmic techniques was ensured. Implementation
251 of advanced algorithmic techniques can lead to different codes by programmers with different levels of experience,
252 and qualitative comparison of these programming models may not produce meaningful conclusions. Second, these
253 algorithms for matrix computations and related methods were implemented so that the parallelisation approaches
254 would be comparable and have the same structure, simplicity, and portability for all multi-core programming models.
255 Finally, the present study focusses on ordinary algorithms and simple parallelisation techniques from the perspective
256 of a moderately experienced programmer. These reasons were taken into account so that the multi-core programming
257 models could be compared in a homogeneous way.

258 4.2. Parallelisation using parallel patterns

259 A brief description of the parallelisation of matrix computations and related methods in terms of parallel patterns
260 is given below.

261 Consider the vector addition $c = a + b$. Each output element of vector c in a vector addition operation is independent
262 of all other elements, and therefore the serial loop can be parallelised using the map parallel pattern. In this case, the
263 map pattern applies the addition operation to each element of vectors a and b . This operation applies as many times
264 as there are elements in vectors a and b .

265 Consider the dot product $z = x^T \cdot y$. The dot product operation has a loop-carried dependency, and therefore it
266 cannot be parallelised. For this reason, the serial kernel operation is reorganised so that each loop iteration calculates
267 an initial pairwise product, and at the end of the loop, these products are summed. This modified kernel can be
268 parallelised using the map parallel pattern combined with a reduction pattern. In this case, the map pattern applies the
269 initial product to each element of vectors x and y , and the reduction pattern applies the summation of the corresponding
270 products.

271 Consider the matrix addition $C = A + B$. Each element of the output matrix C computed in a matrix addition
272 operation is independent of all other elements, and the serial double loop can be parallelised at the coarse-grained
273 level by applying the map pattern to the outer loop. In this case, the map pattern applies the addition operation to each
274 row of matrices A and B . This procedure applies as many times as there are rows in matrix A .

275 Consider the outer product $Z = x \cdot y^T$. Each element of the output matrix Z computed for an outer product operation
276 is independent of all other elements, and therefore the serial double loop can be parallelised at the coarse-grained level
277 by applying the map pattern to the outer loop. In this case, the map pattern applies the product operation between an
278 element of vector x and the entire vector y . This procedure applies as many times as there are elements in vector x .

279 Consider the matrix-vector product $y = Ax$. Each output element of vector y for the matrix-vector product is
280 considered as the computation of the dot product between a row of matrix A and vector x . All output elements of
281 vector y are calculated independently, and therefore the serial double loop can be parallelised by applying the map
282 pattern to the outer loop. In this case, the map pattern applies the dot product operation between a row of matrix A
283 and vector x . This procedure applies as many times as there are rows in matrix A .

284 Consider the matrix transpose $A^T = A$. Each column of the transposed matrix A^T in a matrix transpose operation
285 is considered to be a redistribution of the elements of a row in matrix A . All columns of matrix A^T are calculated
286 independently, and therefore this serial double loop can be parallelised by applying the map pattern to the outer loop.
287 In this case, the map pattern redistributes the elements of a row of matrix A . This procedure applies as many times as
288 there are rows in matrix A .

289 Consider the matrix product $C = AB$. Each row of the output matrix C of a matrix product operation can be
290 calculated as the product between a row of matrix A and the whole matrix B . Therefore, the serial triple loop can be
291 parallelised at the coarse-grained level by applying the map pattern to the first outer loop. This procedure applies as
292 many times as there are rows in matrix A .

293 The Gaussian elimination operation is a direct method for solving dense linear systems of equations such as
294 $Ax = b$, where A is a known non-singular $n \times n$ matrix with non-zero diagonal entries, b is the right-hand side, and
295 x is the vector of unknowns. This discussion will focus only on the computational part of Gaussian elimination,
296 which transforms the matrix A into triangular form with an accompanying update of the right-hand side, so that the
297 solution of the system is straightforward by a triangular solution method. Each step k of Gaussian elimination requires
298 a division operation followed by $n - k$ elimination operations, as shown in Table 3. These division and elimination
299 operations can be parallelised by applying the map pattern to the corresponding *for* loops. In this case, the map pattern
300 of the division operation applies the division operation to each element of a column of matrix A , whereas the map
301 pattern of the elimination operation applies the elimination operations to each row of matrix A and the corresponding
302 element of vector b . This procedure applies as many times as there are rows in matrix A . Note that at the end of
303 the division operation and the elimination phase of step k , there is a synchronisation barrier. This synchronisation is
304 necessary so that the operations of the current elimination step are completed before the start of the next elimination
305 step [19].

306 The Jacobi method is an iterative method for solving linear systems and can be expressed as an iterative scheme
307 consisting of a serial matrix-vector product followed by a solution vector update. Therefore, parallelisation of the
308 Jacobi method corresponds to parallelisation of the matrix-vector product using the map pattern, as discussed earlier.

309 4.3. Parallel implementation issues

310 When parallelising the linear algebra operations described above on a multi-core system, a simple block partition-
311 ing technique was used to increase the amount of work per thread. More specifically, scientific computations that use
312 vectors partition a vector into blocks of equal size, i.e., $\lceil n/p \rceil$ (where n is the number of elements in the vector and p is
313 the number of cores), so that each core operates on a block of $\lceil n/p \rceil$ elements. On the other hand, linear algebra opera-
314 tions that use matrices partition a matrix into blocks of rows of equal size, i.e., $\lceil n/p \rceil$, so that each core operates on a
315 block of rows. In other words, parallelisation of linear algebra operations that use vectors is based on parallelisation of
316 a *for* loop, whereas parallelisation of linear algebra operations that use matrices is based on parallelisation only of the
317 outer *for* loop. Details of how the map and reduction patterns and the synchronisation barrier are expressed in several
318 parallel programming models for parallelising linear algebra operations are presented below. Note that only moderate
319 programming expertise is required to implement the linear algebra methods because only simple code is written and
320 the built-in high-level programming constructs provided by each model can be used.

321 There are no built-in constructs for expressing map and reduction patterns such as the synchronisation barrier in
322 the Pthreads model. For this reason, corresponding codes for implementing these map and reduction patterns and
323 the barrier had to be written. To implement a map pattern with Pthreads, a number of threads are launched by call-
324 ing the `pthread_create` function, which passes two parameters: a function pointer and a pointer to that function's
325 arguments. The call to function `pthread_create` creates a thread that runs the function and passes the specified
326 arguments to it (such as thread rank and data structures) [44]. For parallel implementation of the corresponding linear
327 algebra operation, a corresponding function was created that embraces the original serial body code of the correspond-
328 ing linear operation, but the bounds of the loop were modified so that each thread works on its own local data block.

329 The bounds of the loop are based on thread rank as follows: each thread works on its own block of data $b = \lceil n/p \rceil$,
330 which ranges from $rank \times b$ to $(rank + 1) \times b$. Finally, the threads are synchronised at the end of the calculation
331 by calling the `pthread_join` function. To implement a reduction pattern with Pthreads, the low-level explicit syn-
332 chronisation function `pthread_mutex_lock` was used. To express the barrier in the Pthreads model, a linear barrier
333 was implemented as an individual function using the functions `pthread_mutex_lock`, `pthread_mutex_unlock`,
334 `pthread_cond_broadcast`, and `pthread_cond_wait`.

335 The map and reduction patterns can be expressed in the OpenMP model using the built-in constructs `#pragma omp`
336 `parallel for` and `#pragma omp reduction`. The code following the `#pragma omp parallel` for a compiler
337 directive (i.e., the parallel region) is run by a number of parallel threads, and loop iterations are distributed to threads
338 executing in the parallel region. At the end of the loop, there is an implicit barrier. Moreover, distribution of loop
339 iterations to threads is done by an OpenMP scheduling strategy. OpenMP supports various scheduling strategies
340 specified by the `schedule` clause. In this case, no strategy was specified in the `schedule` clause, and therefore the
341 OpenMP model created evenly sized blocks of loop iterations for all threads in the parallel region. The `#pragma omp`
342 `reduction` clause implements an automatic parallel reduction of the variables that appear in its list with a predefined
343 associative operator. To implement the sum reduction in the dot product operation, a *combine* operator (addition) was
344 specified. Finally, the built-in `#pragma omp barrier` was not used to define a synchronisation barrier because there
345 is an implicit barrier at the end of the parallel *for* loop.

346 Cilk's `cilk_for` and `reducer_opadd` keywords perform map and reduction patterns respectively. The `cilk_for`
347 keyword processes loops in parallel and divides each loop into equal-sized blocks containing one or more loop iter-
348 ations. Each block is executed serially and is spawned as a block during loop execution. The maximum number of
349 iterations in each block is the grain size. The `#pragma cilk_grainsize` clause was also introduced at the begin-
350 ning of the loop, and the grain size for one loop was specified as $\lceil n/p \rceil$ according to the proposed data partitioning
351 strategy. Moreover, implementation of the `cilk_for` construct is based on a recursive fork-join model, which spreads
352 the overhead of managing the map over multiple threads. Finally, the reduction pattern can be expressed in Cilk Plus
353 by a hyperobject. This hyperobject is a reducer, which eliminates contention for shared variables among tasks by
354 automatically creating views of them for each task and reducing them back to a shared value after task completion.
355 Cilk's reducers work for any operation that can be re-associated. To implement the sum reduction, the hyperobject
356 `reducer_opadd` was used. Finally, the synchronisation barrier in this model was implemented by an implicit barrier
357 at the end of `cilk_for`.

358 The map and reduction patterns can be implemented in TBB using the built-in template functions `parallel_for`
359 and `parallel_reduce` respectively. To parallelise linear algebra computations with the TBB model, the original outer
360 loop was replaced by a call to the `parallel_for` template function to execute loops in parallel. The `parallel_for`
361 template function requires two arguments: the loop's iteration space information and a lambda function. This function
362 partitions the iteration space into subspaces and applies the lambda function to each subspace in parallel. The iteration
363 space is defined by three parameters: the lower bound, the upper bound, and the optional increment. The lambda
364 function is a new feature introduced by the C++11 standard and is widely supported in C++ compilers. The lambda
365 function is an alternative way to create the equivalent function object, which would otherwise be written by hand.
366 The lambda function has three parts: the capture part, a parameter list, and the function definition. The capture part
367 describes how to capture local variables outside the lambda, i.e., by value or by reference. The argument list takes
368 an argument of type `blocked_range<T>`, which defines an iteration range. The function definition takes the original
369 loop or code of the corresponding linear algebra method; the loop's iteration bounds are specified dynamically by
370 the argument passed to the lambda function [50, 44, 62]. Finally, the reduction pattern can be expressed using the
371 template function `parallel_reduce`, which performs a parallel reduction over a recursive range. This template
372 function has two forms: the functional form is designed to be easy to use in conjunction with lambda expressions,
373 and the imperative form is designed to minimise data copying. Here, the functional form of *reduce* was used, which
374 requires four arguments: the iteration space as in `parallel_for`, the identity element of type T, e.g., `double`, a
375 functor provided by the user for the serial reduction, and finally a binary operator reduction to combine the results of
376 the user functor, e.g., `plus<double>` for summation of double elements [50].

377 The map and reduction patterns can be expressed in SWARM using the built-in constructs `SWARM_pardo` and
378 `SWARM_reduce` respectively. To implement linear algebra methods in the SWARM framework, the original outer *for*
379 loop was replaced by a `SWARM_pardo` construct for executing loops concurrently on one or more processing cores.
380 This simple construct implicitly partitions the loop among the cores without the need for coordinating overheads

381 such as synchronisation of communication between cores. Moreover, the sum reduction was replaced by a function
382 `SWARM_reduce`, which is used to add values calculated by different threads. This function requires three arguments:
383 the local variable of each thread, the predefined reduction operator (e.g., `SUM`) that is performed on these variables, and
384 a built-in structure `THREADED` that contains all the required information for the particular thread. Finally, the function
385 `SWARM_Barrier_sync` was used to express the synchronisation barrier.

386 The FastFlow model does not provide ready-to-use constructs for the map and reduction patterns, and therefore
387 code had to be written to implement these patterns explicitly using other features such as an accelerator with the farm
388 template and spinlocks. To parallelise linear algebra methods for the map pattern, the FastFlow accelerator was used,
389 which is a software device that can be used to speed up skeleton structured portions of code using the cores left unused
390 by the main application. In other words, FastFlow accelerates particular computations by using a skeleton program,
391 offloading to it the tasks to be computed. For this reason, a skeleton program was written using the FastFlow skeletons
392 (i.e., the farm template), which computes the tasks that will be given to the accelerator. This skeleton program, when
393 used to program the accelerator, must have an input stream, which is used to offload tasks to the accelerator. Then the
394 skeleton program must be run using a method such as `run_then_freeze()`. This method will start the accelerator
395 skeleton program, consuming the input stream items either to produce output stream items or to consolidate (partial)
396 results in memory [10, 11]. Finally, a new class `task` was created for each linear algebra method, with members
397 corresponding to the variables that are outside the scope of the loop to be parallelised. The class must overload the
398 function `task` method, and it takes an argument such as task rank. The method's body takes the original outer loop or
399 code of the corresponding linear algebra method; the loop's iteration bounds are based on task rank by the argument
400 passed to the function `task` method, so that each task works on its own block of data. Moreover, to implement the
401 reduction pattern with FastFlow, the low-level explicit synchronisation function `spin_lock` was used. Finally, the
402 FastFlow method `wait()` was used to express the synchronisation barrier.

403 5. Quantitative and qualitative results

404 This section presents the results of an extensive quantitative and qualitative evaluation performed to gain insight
405 into the practical behaviour of each of the reviewed programming models used to implement the scientific computa-
406 tions.

407 5.1. Quantitative comparison

408 For quantitative or performance comparison, computational experiments were performed. The experiments were
409 run on a Dual Opteron 6128 CPU with eight processor cores (16 cores total), a 2.0 GHz clock speed, and 16 GB of
410 memory under Ubuntu Linux 10.04 LTS. During all experiments, this machine was in exclusive use by this research
411 project. The implementations of the serial and multi-thread algorithms for the nine linear algebra operations were
412 developed in the ANSI C/C++ programming language for all the reviewed multi-core programming models. To
413 compile the multi-thread programs, the Intel C/C++ compiler version 13.0, i.e., `icc`, `icpc`, was used because it is a
414 very widely used compiler.

415 The quantitative evaluation was organised into two parts:

- 416 1. Examination of the parallel performance (in terms of execution time) of the linear algebra methods for all the
417 reviewed multi-core programming models. More specifically, performance was examined in two ways: (a) as
418 a fixed-size problem with increasing number of cores, and (b) with a fixed number of cores where the matrices
419 or problems increase in size. Moreover, the optimal performance of each computational kernel was examined
420 against the programming model. These performance tests used parallelisations of the optimised sequential
421 implementations.
- 422 2. Examination of the overall performance of all the reviewed programming models as a function of number of
423 cores and of matrix size.

424 A set of randomly generated input matrices and vectors with sizes ranging from 1024×1024 to 7168×7168
425 in double-precision format was used to evaluate the performance of the multi-thread linear algebra methods. The
426 elements of the generated input matrices for the matrix computations were chosen from a uniform distribution on the
427 interval (0.0, 100.0), whereas the input matrices for the methods for solving linear systems were created to be strictly

428 row-diagonal dominant. For experimental comparison of the results, the dense and sparse matrices can be used as test
 429 cases for all matrix computations and related methods, but this study focusses on dense matrices, which serves the
 430 goals of this paper.

431 To assess the performance of the multi-thread linear algebra methods for all programming models, practical exe-
 432 cution time was used as a measure. The practical execution time is the total time that a multi-thread algorithm takes to
 433 complete a computation. The execution time is obtained by calling the C function `gettimeofday()` and is measured
 434 in seconds. To decrease random variation, the execution time was measured as an average of ten runs. A models
 435 performance for each computational kernel was calculated using the relation

$$\frac{BT(model)t}{T(model)i} \times 100, \quad i = 1, 2, \dots, 9 \quad (1)$$

436 where $BT(model)t$ is the best execution time for a specific computational kernel for all six models and $T(model)i$ is the
 437 execution time of model i for the same computational kernel. The best time was then set equal to 100%. To calculate
 438 the overall performance for each model for each combination of matrix size and number of cores, the percentage
 439 values for all computational kernels were added up, and the sum was divided by the total number of computational
 440 kernels, i.e., nine. Finally, the highest percentage corresponded to the best overall performance.

441 The figures following this section use data extracted from the parallel performance experiments. These figures are
 442 performance graphs for all computational kernels using all programming models. Note that the performance results
 443 are based on parallelisation of the optimised sequential implementations with flag `O2`, which includes vectorisation.
 444 Figure 1 (or Tables 10-18 in the appendix) presents the execution times for all computational kernels as a function of
 445 the number of cores for a fixed matrix size of 7168, whereas Figure 2 (or Tables 28-36 in the appendix) presents the
 446 execution times for all computational kernels as a function of problem size for a fixed number of cores (16). Note
 447 that time on the y -axis in Figures 1 and 2 has a logarithmic scale. Figure 1 includes a line for the serial execution
 448 time, which is the time taken by a pure C sequential program. Finally, the tables also report the bootstrap confidence
 449 intervals of the mean execution time for all computational kernels.

450 From the graphs in Figure 1, it is clear that there was little overhead for the programming models on a single
 451 core, except for the vector addition and dot product operations. The vector addition and dot product operations had
 452 the greatest overheads for all programming models and all numbers of cores because of the limited amount of par-
 453 allelisable work available. Pthreads, Cilk Plus, TBB, and FastFlow incurred some overhead for all linear algebra
 454 methods (except for vector addition and dot product) on a single core, but all six multi-core implementations showed
 455 identical performance with few exceptions when multiple cores were used. This overhead was due to use of the
 456 `pthread_create` and `pthread_join` functions in the Pthreads model to create and terminate threads, use of auto-
 457 matic partitioning of the `parallel_for` template function in the TBB model, which incurs an overhead cost for every
 458 chunk of work that it schedules, use of the `cilk_for` construct in the Cilk Plus model, which is based on the recursive
 459 fork-join model, and the implicit barrier at the end of the loop and the use of additional internal synchronisation mech-
 460 anisms in the FastFlow model to manage tasks. Note that a partitioner was not declared in the third optional argument
 461 of the `parallel_for` function of the TBB model, and therefore TBB used automatic partitioning. The OpenMP and
 462 SWARM implementations of the other seven linear algebra methods (except for vector addition and dot product) had
 463 very little overhead for one core and scaled best for all numbers of cores. Note that for the matrix product operation,
 464 no overhead cost was incurred for any of the programming models.

465 Furthermore, from the graphs in Figure 1, the execution time of all programming models for all computational ker-
 466 nels decreased as the number of cores increased, with the exception of the vector addition and dot product operations.
 467 More specifically, the execution times of the matrix addition, outer product, matrix-vector product, matrix product,
 468 Gaussian elimination, and Jacobi operations decreased significantly compared to the execution time of the matrix
 469 transpose operation. The limited rate of decrease in the execution time of the matrix transpose operation occurred
 470 because this operation has poor spatial locality and scans the matrix column by column instead of row by row, leading
 471 to a higher cache miss rate. However, there was a saturation point with eight cores for all programming models for
 472 the matrix addition, matrix-vector product, Gaussian elimination, and Jacobi operations, after which the execution
 473 time was stable. This occurred because the amount of parallelised work for the matrix addition and matrix-vector
 474 operations was small with a large number of cores because these operations are medium-scale. Furthermore, the con-
 475 stant performance of the Gaussian elimination and Jacobi methods with a large number of cores (i.e., more than eight)
 476 was observed because Gaussian elimination requires more synchronisation with a large number of cores and because

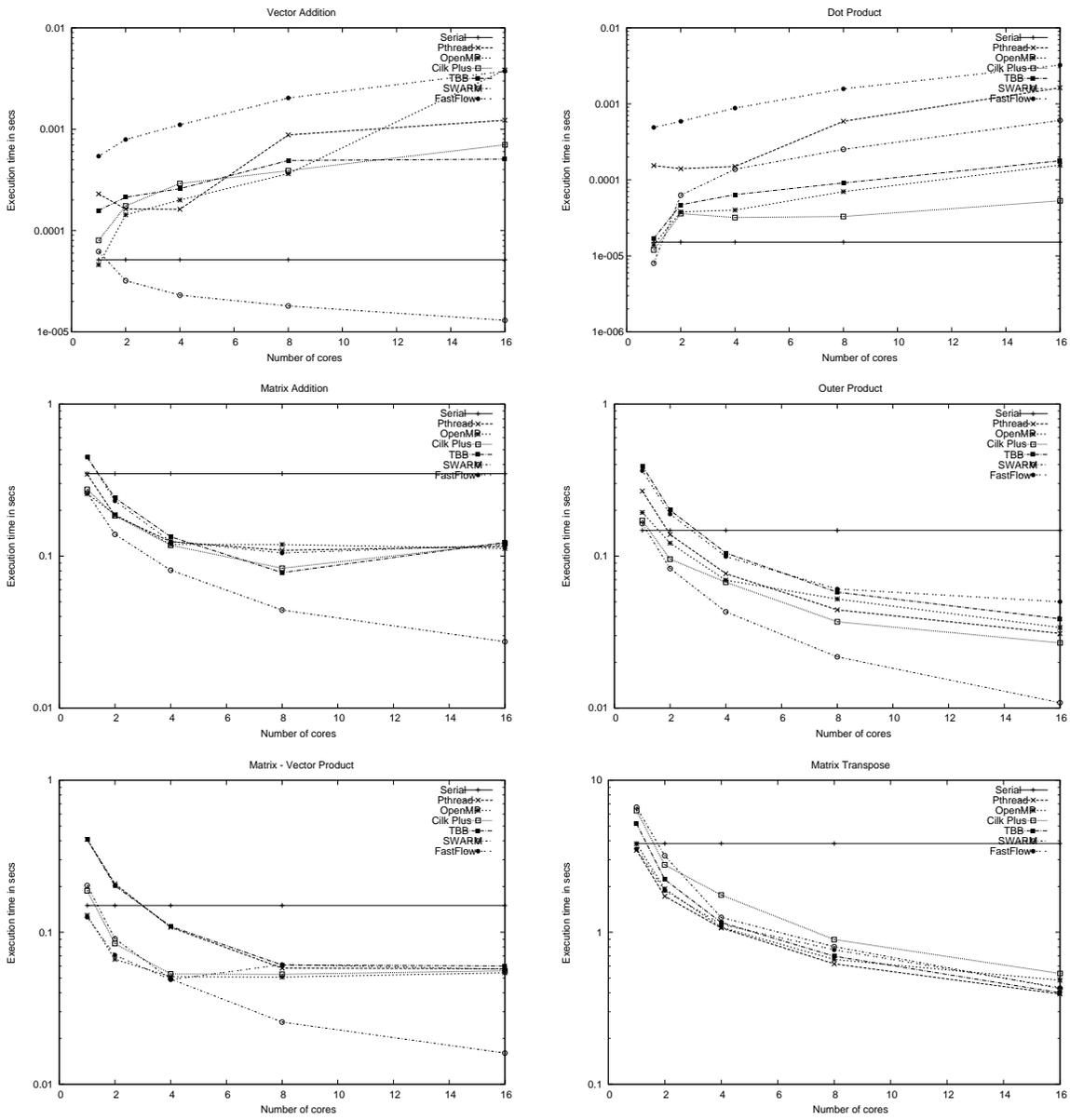


Figure 1: Execution times (in secs) of the scientific computing kernels as a function of the numbers of cores for a fixed problem size of 7168

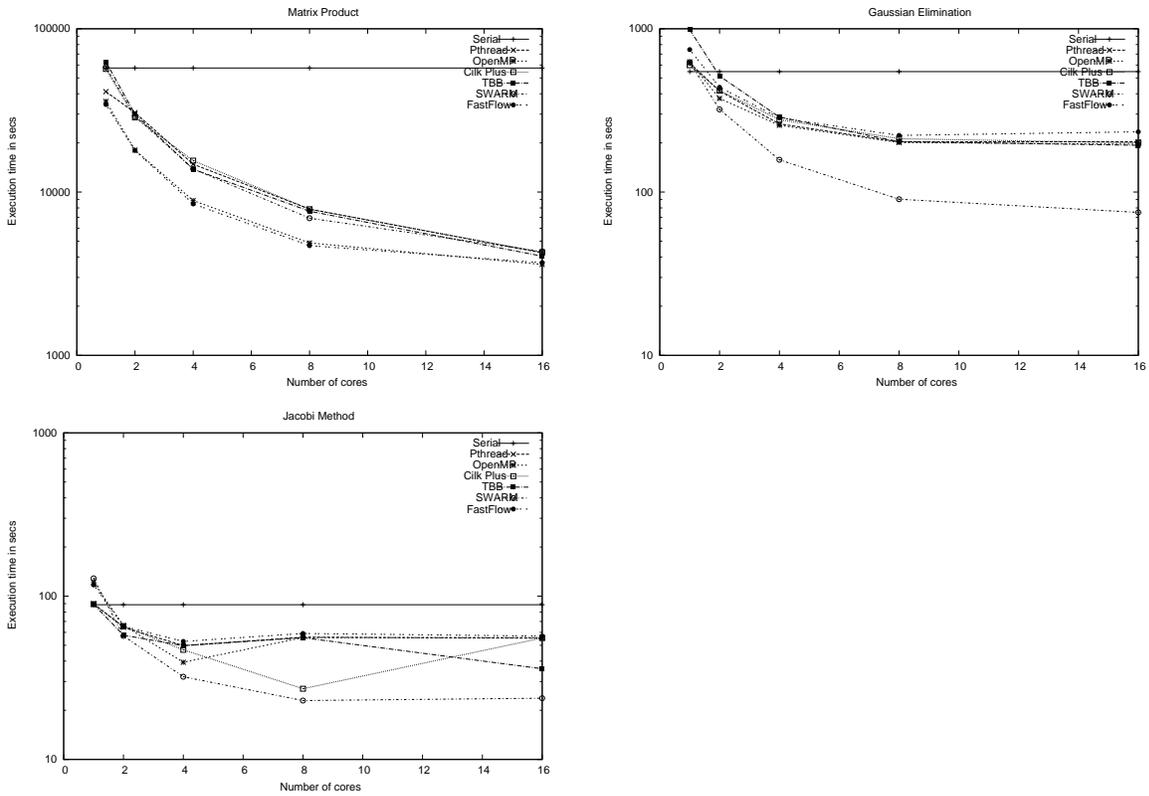


Figure 1: (continued) Execution times (in secs) of the scientific computing kernels as a function of the numbers of cores for a fixed problem size of 7168

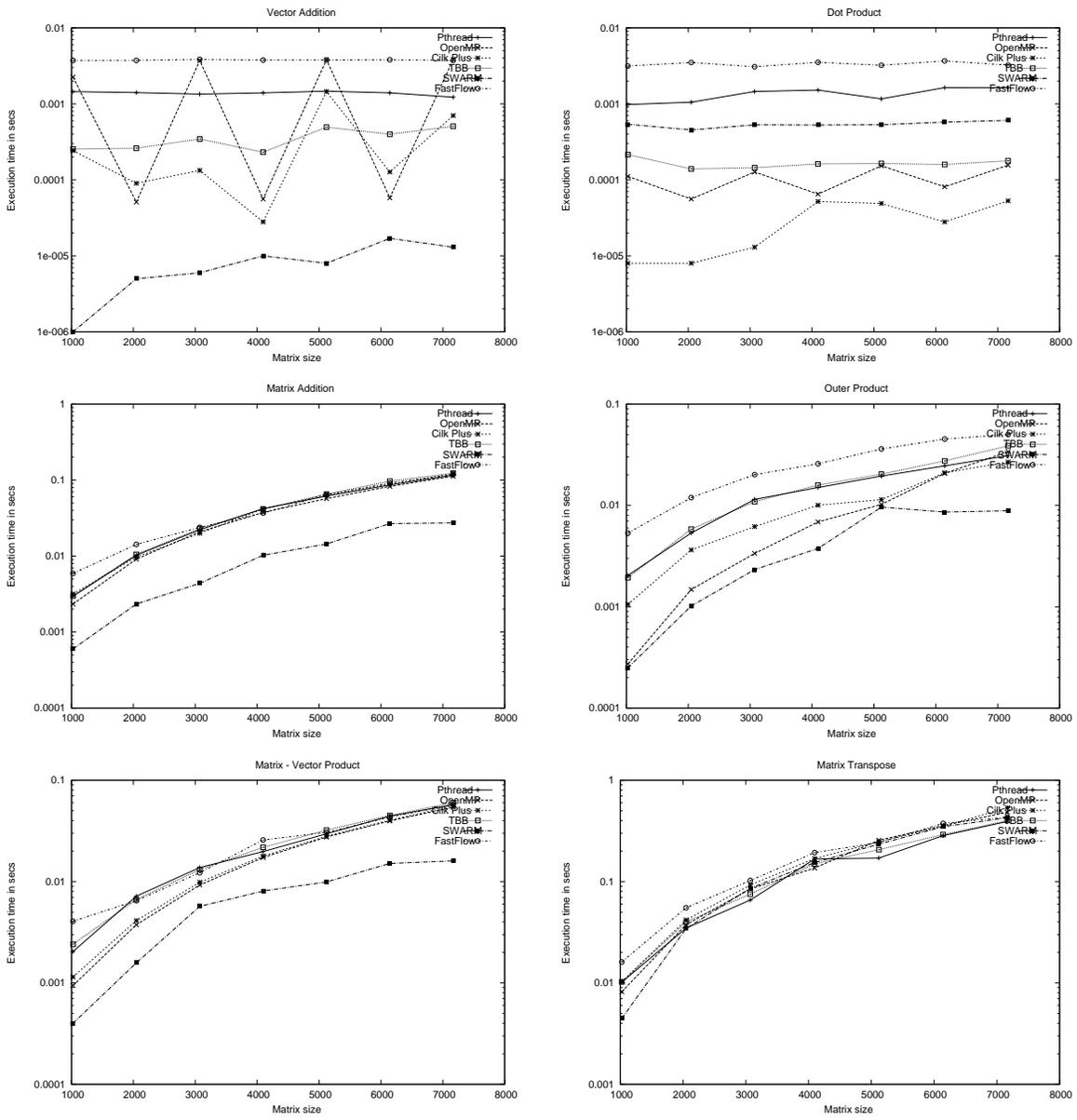


Figure 2: Execution times (in secs) of the scientific computing kernels as a function of the problem size for 16 cores

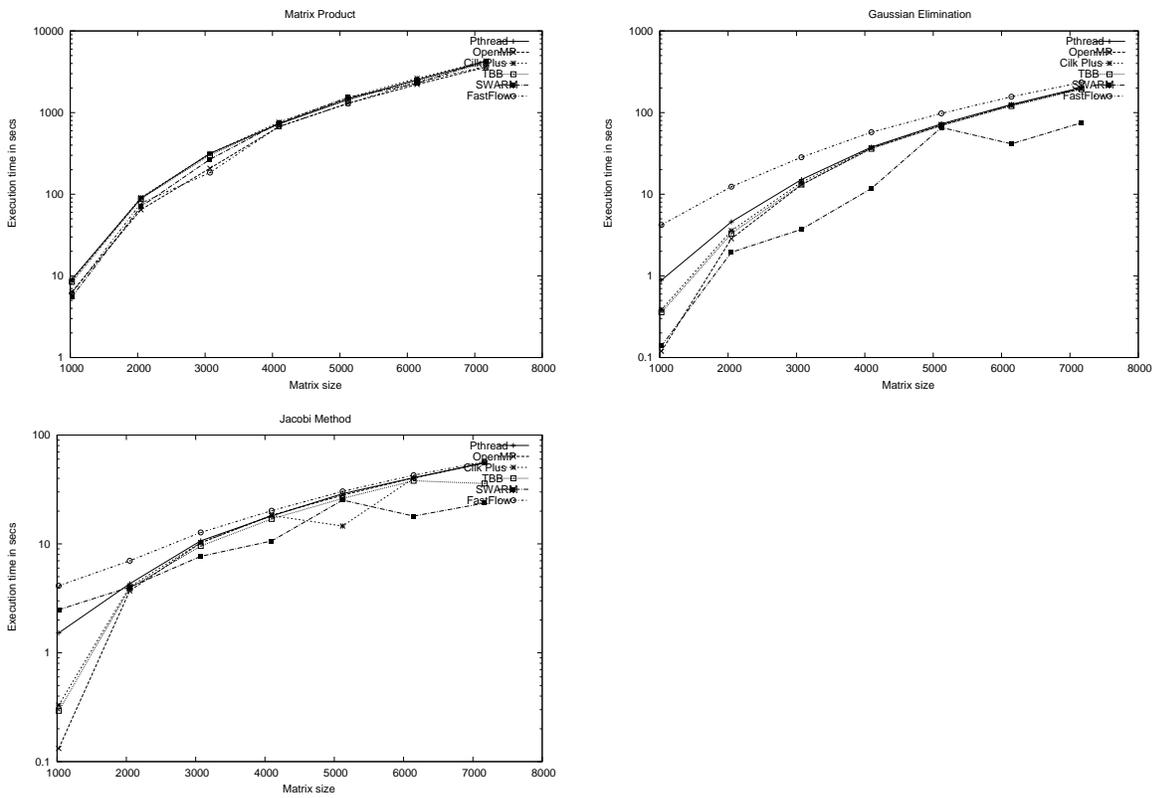


Figure 2: (continued) Execution times (in secs) of the scientific computing kernels as a function of the problem size for 16 cores

477 the Jacobi method uses the matrix-vector operation as a medium-scale operation, as explained earlier. On the other
478 hand, the execution time of the vector addition and dot product operations increased slightly as the number of cores
479 increased. This happened because these operations require a small number of steps (n steps approximately) compared
480 to the other six matrix operations, which require n^2 or n^3 steps.

481 Based on the graphs in Figure 1, some specific remarks can be made about performance. For the vector addition
482 operation, the SWARM implementation had the best performance in terms of execution time for any number of cores,
483 whereas the FastFlow implementation had the slowest performance. The poor performance of FastFlow was due to
484 its use of a built-in software accelerator (i.e., a collection of threads) to implement the vector addition method, which
485 provided additional synchronisation mechanisms (i.e., locks and waits) to manage tasks internally. This overhead for
486 internal synchronisation mechanisms was clearly evident for a large number of cores compared to the small amount of
487 computation needed for the vector addition method. This overhead can be eliminated by an experienced programmer
488 because the FastFlow programming model has a moderately steep learning curve, as will be discussed later. Moreover,
489 the performance of the remaining models was closest to that of the Pthreads implementation.

490 For the dot product operation, the performance of all parallel implementations was worse than the serial perfor-
491 mance. The Cilk Plus model had the best performance among the poorly performing parallel implementations. The
492 poorer performance of these models can be explained by their high overhead costs compared to the limited amount of
493 parallelised work performed.

494 For the matrix addition operation, the OpenMP and SWARM implementations had the best performance in terms
495 of execution time for small and large numbers of cores respectively. On the other hand, the performance of the other
496 programming models was similar for a large number of cores, with the exception of the SWARM model. However,
497 the Pthread, TBB, and FastFlow models had the poorest performance for any number of cores due to the overhead of
498 creating and terminating threads in the Pthread model, the automatic partitioning in the TBB model, and the additional
499 internal synchronisation mechanisms discussed earlier.

500 For the outer product operation, SWARM had the best performance in terms of execution time for any number of
501 cores due to its low overhead. On the other hand, TBB and FastFlow had the worst performance for small and large
502 numbers of cores respectively. The poor performance of the TBB model was due to additional implicit partitioning
503 overhead, and the poor performance of FastFlow was due to synchronisation overhead, as discussed earlier. Note that
504 the amount of synchronisation overhead was large with respect to the amount of parallelisable work available. For the
505 remaining programming models, their performances were approximately the same for any number of cores.

506 For the matrix-vector product, the OpenMP and SWARM implementations had the best performance in terms of
507 execution time for small and large numbers of cores respectively. The SWARM implementation performed best, with
508 a large improvement compared to the other models, for a large number of cores because the system runtime overhead
509 was spread over multiple cores. The performance of the remaining models was identical.

510 For the matrix transpose operation, Pthreads had satisfactory performance for any number of cores. The perfor-
511 mance of the remaining implementations was closest to those of Pthreads and SWARM.

512 For the matrix product operation, the OpenMP and FastFlow implementations had satisfactory execution time for
513 any number of cores. However, a low rate of decrease was observed in the execution time of the FastFlow model. From
514 these two performance observations, the following could be concluded: the synchronisation overhead for FastFlow
515 with a small number of cores was negligible compared to the high amount of parallelised work for the matrix product,
516 whereas the synchronisation overhead was clearly evident for a large number of cores compared to the small amount
517 of parallelised work performed. Finally, the performance levels of the remaining implementations were identical for
518 any number of cores.

519 Finally, for the Gaussian elimination and Jacobi operations, the SWARM implementation had satisfactory perfor-
520 mance for any number of cores. The performance values of the remaining implementations were identical for any
521 number of cores.

522 The ANOVA statistical test was used to verify the statistical significance of these results. Tables 1-9 in the ap-
523 pendix show the ANOVA results for all matrix operations. As can be seen in the tables, the factors (programming
524 model and number of cores) and their interaction (programming model \times number of cores) were statistically signifi-
525 cant (p -value < 0.05).

526 Based on the graphs in Figure 2, it is apparent that the execution time of all programming models for all com-
527 putational kernels increased with increasing matrix size, with the exception of the vector addition and dot product
528 operations. More specifically, the execution times of the matrix addition, outer product, matrix product, matrix trans-

529 pose, Gaussian elimination, and Jacobi operations increased significantly, whereas that of the matrix-vector product
530 increased only slightly.

531 Based on the graphs in Figure 2, specific remarks about performance can be made. For the vector addition and
532 dot product operations, the SWARM, OpenMP, and Cilk Plus implementations had good and constant performance
533 for all problem sizes because they incurred low system overhead. On the other hand, Pthreads and FastFlow had poor
534 performance because of the high amount of overhead compared to the amount of parallelisable work available, as
535 discussed earlier.

536 For the matrix addition, outer product, and matrix-vector product operations, the SWARM implementation had
537 the best performance in terms of execution time for any problem size. On the other hand, the Pthreads and FastFlow
538 implementations had the poorest performance because of their system overhead. Note that the performance values of
539 the Pthreads and FastFlow implementations of the matrix addition operation were closer to each other than to those of
540 the other implementations. Moreover, the performance of the remaining implementations for the other two operations
541 was modest, and the execution times of these implementations were very close.

542 For the matrix transpose operation, all implementations had good and identical performance for medium and large
543 matrix sizes. The FastFlow implementation had the longest execution time, and the differences in execution times
544 compared to the other implementations decreased as matrix size increased. This can be explained by noting that the
545 synchronisation overhead for FastFlow becomes small when the matrix size is large.

546 For the matrix product operation, the performance of all programming models was excellent, and they scaled iden-
547 tically as a function of matrix size. For this kernel operation, the OpenMP and FastFlow implementations performed
548 the best, with minimal differences from each other.

549 Finally, for the Gaussian elimination and Jacobi operations, the SWARM implementation had satisfactory perfor-
550 mance for any problem size. The execution times of the remaining implementations were similar. The performance
551 of Gaussian elimination for all models improved as matrix size increased because of the limited synchronisation (i.e.,
552 barrier) overhead when the matrix size was large.

553 The ANOVA results for the execution times for all matrix operations as a function of problem size for 16 cores
554 are listed in Tables 19-27 in the appendix. As can be seen in the tables, the factors (programming model and problem
555 size) and their interaction (programming model \times problem size) were statistically significant (p -value < 0.05).

556 Figures 3 and 4 present the results from a preliminary scalability analysis for two selected matrix operations, the
557 outer product and the matrix product. Figure 3 shows the speed-ups of these two matrix operations when the number
558 of cores was increased for various problem sizes. Speed-up is defined as the serial reference runtime of the pure C
559 sequential program when run on a single-core processor, divided by the time taken by the optimised multi-thread
560 program on a multi-core processor. In this analysis, SWARM was used for the outer product and OpenMP for the
561 matrix product because they performed well in the previous tests. These results show that a larger instance of the same
562 problem yields higher speed-up and efficiency for the same number of cores, although efficiency continued to drop
563 as the number of cores increased. Note that efficiency is a measure of the fraction of the time during which a core is
564 usefully employed; it is defined as the ratio of speed-up to the number of cores [38].

565 Figure 4 shows the scalability of two matrix operations (the outer and matrix products) for all programming
566 models. Here the problem size was increased at the same rate as the number of cores, so that the number of basic
567 operations required by each core for each problem remained approximately constant. These results show that the
568 efficiency of the two matrix operations was approximately constant for all programming models. Furthermore, the
569 SWARM and OpenMP models achieved higher efficiency (from 0.80 to 0.99) than the remaining models for the
570 outer product operation, whereas the other models except for FastFlow had efficiencies ranging from 0.55 to 0.70.
571 However, the efficiency for the outer product operation fluctuated for some programming models. More specifically,
572 the overhead percentages for the SWARM, OpenMP, Cilk Plus, Pthread, TBB and FastFlow models ranged from 1%
573 to 9%, 3% to 24%, 5% to 27%, 1% to 15%, 7% to 25%, and 3% to 40% respectively as the number of cores increased.
574 On the other hand, the efficiency for the matrix product was sufficiently high (from 0.80 to 0.90) for all programming
575 models. In addition, the overhead percentage was sufficiently small (i.e., 3% to 30% for a large number of cores)
576 compared to the efficiency of the outer product operation. Finally, an ANOVA analysis showed that these results were
577 statistically significant.

578 Figure 5 presents the overall performance for all the reviewed programming tools as a function of number of cores
579 and of matrix size. The average performance shown on the left of Figure 5 for each model and each number of cores
580 is the median value for all computational kernels and matrix sizes. Similarly, the average performance shown on the

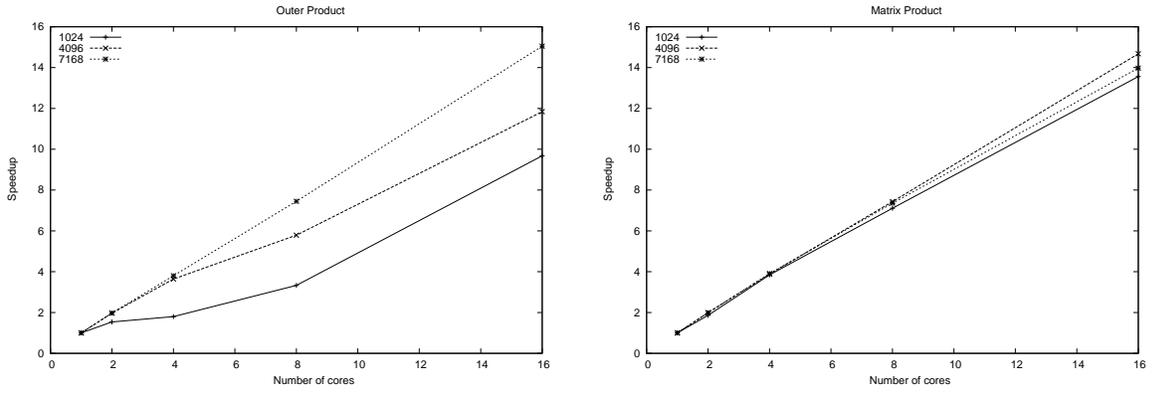


Figure 3: Speedups of two matrix operations versus the number of cores for different problem sizes using the SWARM model (left) and OpenMP model (right)

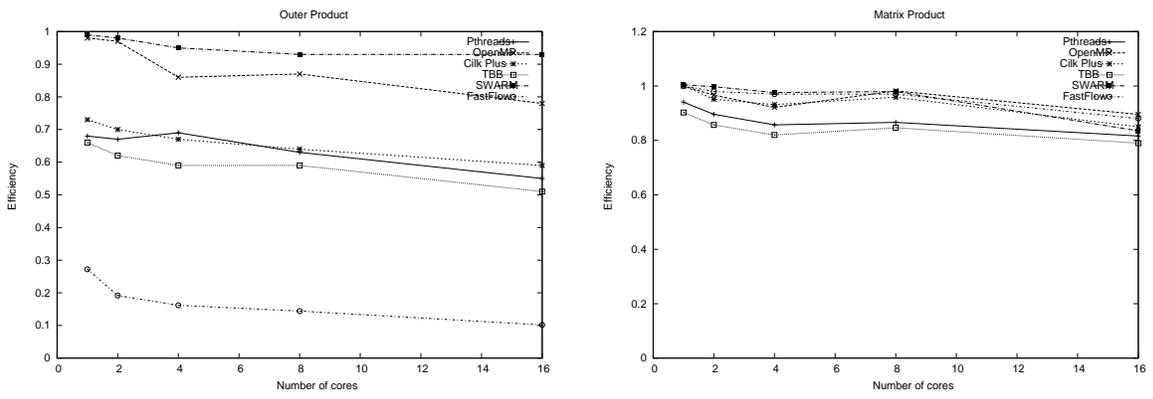


Figure 4: Efficiency of two matrix operations as a function of problem size and number of cores

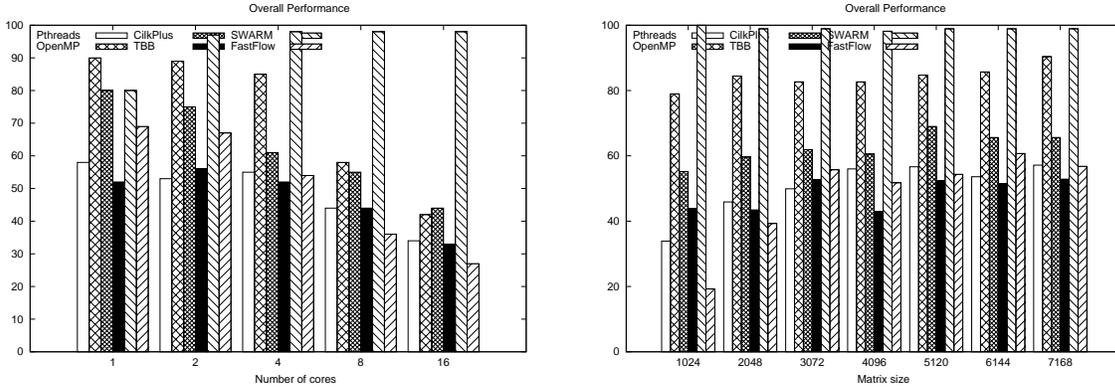


Figure 5: Overall performance as a function of the number of cores (left) and overall performance as a function of the problem size (right)

581 right of Figure 5 for each model and each matrix size is the median value for all computational kernels and numbers
 582 of cores.

583 Based on these results, OpenMP had the best performance with a small number of cores, whereas SWARM had
 584 the best performance for a larger number of cores and for any matrix size. However, the performance difference
 585 between the OpenMP and SWARM models became small as matrix size increased, whereas the difference remained
 586 large as the number of cores increased. Moreover, the overall performance of the SWARM and OpenMP models was
 587 outstanding, with large differences from the other models for a large number of cores, as presented on the left-hand
 588 side of Figure 5. The overall performance gap between SWARM and OpenMP and the other models was constant for
 589 all matrix sizes. In both cases, the SWARM and OpenMP models had the best performance because they appeared as a
 590 "winner" in most linear algebra operations (except for matrix product) and because they provide an optimised runtime
 591 system, unlike the other models. The next best performing models in order were Cilk Plus, TBB, and FastFlow, with
 592 very small differences for any number of cores and any matrix size. Finally, in last place came the Pthreads model,
 593 which had the poorest performance in most cases.

594 5.2. Qualitative comparison

595 Research has been carried out into the programming effort required for parallel programming [13, 42, 48, 61, 67].
 596 Programming effort has been evaluated using basic metrics such as the time taken to write a program, the development
 597 time, or the number of lines of source code. However, it is not easy to determine the development time for coding a
 598 program, and the number of lines of code has been recognised as an inadequate measure [42, 64]. For this reason, these
 599 metrics (lines of code and development time) were combined here with a series of software engineering parameters
 600 such as library popularity, level of abstraction, programming style, support for online help and documentation, and
 601 learning curve to assess the intensity of programming effort required for the multi-core programming frameworks
 602 [30, 67]. Recent work [64] has identified more useful metrics for evaluating programming effort, but they require
 603 powerful and reliable tools for data gathering and analysis.

604 The number of additional or modified lines of code needed to parallelise the problem or function was counted
 605 in each case, excluding variable declarations. The development time included the time in days for parallelising and
 606 coding the problem in each multi-core programming framework by a user with moderate programming experience.
 607 Because it is difficult and time-consuming to measure development time, an empirical and qualitative rating scale was
 608 used, where the notation + + + means the shortest time, whereas the notation - - - denotes the longest time [13].
 609 Table 4 gives the number of additional or modified lines of code and the qualitative rating of development time for the
 610 nine computational kernels required by each framework.

611 Based on this table, the following remarks can be made: the OpenMP, SWARM, Cilk Plus, and TBB implemen-
 612 tations required the same number of additional lines to parallelise a serial function. In the OpenMP and Cilk Plus

Computational kernel	Pthreads		OpenMP		Cilk Plus		TBB		SWARM		FastFlow	
	LOC	Time	LOC	Time	LOC	Time	LOC	Time	LOC	Time	LOC	Time
Vector Addition	17	-	1	+++	2	+++	3	+++	1	+++	15	-
Dot Product	19	--	2	+++	3	+++	4	++	2	+	17	--
Matrix Addition	18	-	1	+++	2	+++	3	+++	1	+++	15	-
Outer Product	17	-	1	+++	2	+++	3	+++	1	+++	15	-
Matrix-Vector Product	18	-	1	+++	2	+++	3	+++	1	+++	16	-
Matrix Transpose	17	--	1	+++	2	+++	3	+++	1	+++	15	---
Matrix Product	18	---	1	+++	2	+++	3	++	1	+++	16	---
Gaussian Elimination	21	---	3	+	4	+	4	+	2	+	23	---
Jacobi Method	19	---	1	+	2	+	3	+	4	+	12	---

Table 4: Additional lines of code and development time in the multi-core programming models [Note: LOC means the number of lines of code and time means the development time]

613 implementations, a programmer can easily insert compiler directives (i.e., pragmas) or keywords into sequential code
614 to tell the compiler which parts of the code should be executed in parallel. These models have the advantage that the
615 only additions necessary to the sequential code are pragmas or keywords. In the SWARM implementation, a program-
616 mer uses constructs for parallelisation, e.g., to parallelise a *for* loop, the `par_do` construct should be used at the start
617 of the loop; the `par_do` construct implicitly partitions the loop among the cores without the need for coordinating
618 overheads such as synchronisation between cores. Furthermore, the SWARM framework provides the programmer
619 with library functions and with calls to synchronisation and reduction operations. The syntax of these functions is
620 easy to use because only a few arguments are required. In the TBB programming model, the programmer also uses
621 template functions to parallelise a *for* loop and to perform reduction using the lambda function feature. Note that
622 the lambda function is used for anonymous function definitions, meaning that a *for* loop could be passed to TBB's
623 `parallel_for` function in-place, i.e., without defining a function object [44, 50]. This new feature is important be-
624 cause it helps reduce the programming effort related to function objects, for which many lines of code are required.
625 Finally, as observed in the previous section, the OpenMP, SWARM, Cilk Plus, and TBB frameworks express data
626 parallelism and reduction easily and directly because these frameworks provide ready-to-use parallel constructs or
627 functions that do not require additional lines of code. This observation is supported by the fact that the number of ad-
628 ditional lines of code required for these four programming models was smaller than for the others and the development
629 time was sufficiently small. Finally, the code structure for the OpenMP, SWARM, Cilk Plus, and TBB frameworks
630 was simple and easy to understand, and no complex programming effort was required.

631 The Pthreads implementations required more additional lines of code than their sequential counterparts, and the
632 programming effort was complex. The Pthreads programming model provides a programmer with low-level library
633 routines, and unlike the other frameworks, algorithm parallelisation is not automatic, meaning that the programmer
634 is explicitly responsible for writing code for parallelisation, data distribution to each thread, and thread synchronisa-
635 tion. In other words, a manual code transformation is required to achieve parallelism. To parallelise a *for* loop, the
636 programmer must create, join, and synchronise threads explicitly. Furthermore, the programmer must reorganise the
637 code from the serial function into a thread function, and thread inputs must be grouped into a single data structure.
638 Moreover, the code for the thread function must be modified so that each thread works on its own local data block.
639 The result of all this is that Pthreads has the longest coding and development time of all the algorithms.

640 Finally, the FastFlow algorithm implementations required many additional lines of code, although C++ templates
641 were provided. The code for the FastFlow programming model is easier to understand than the others, but because
642 each algorithm was implemented as an object class, additional statements were required as well as the core code of the
643 algorithm. In other words, this model required restructuring of the sequential code to make it more object-oriented.
644 Therefore, this programming model requires significant programming effort and development time.

645 The programming effort or the development time required by each programming framework depends on the num-
646 ber of additional lines of code. In other words, the more additional lines of code required to parallelise a function, the
647 longer the development time, and vice versa. Furthermore, there is a relationship between the support for parallel pat-
648 terns provided by each framework and the expertise level of the programmer. In other words, the programming models
649 that directly support parallel patterns through ready-to-use parallel constructs require little programming experience

Characteristic - Model	Pthreads	OpenMP	Cilk Plus	TBB	SWARM	FastFlow
Library popularity	High	High	Medium	Medium	Low	Low
Level of abstraction	Low	High	High	High	High	High
Programming style	C	C	C++	C++	C	C++
Help and documentation	Very good	Very good	Very good	Very Good	Good	Good
Learning curve	High	Low	Low	Medium	Low	Medium

Table 5: The characteristics of the multi-core programming models

650 from the user, whereas when parallel patterns are not supported directly by a model and they must be implemented
651 with other features, advanced programming experience is required, as for the Pthreads and FastFlow models.

652 Regarding the other software engineering parameters, the Pthreads and OpenMP programming models are the
653 most popular with programmers because they were developed earlier and have been used extensively by parallel
654 computing researchers. The remaining programming models have medium or low popularity because they were de-
655 veloped when multi-core processors were introduced. The Pthreads framework has a low level of abstraction because
656 the programmer is responsible for writing code for managing thread parallelism and synchronisation, whereas the
657 other frameworks provide a high level of abstraction that can implement the same kind of thread parallelism using
658 ready-to-use parallel constructs. These frameworks, except for Pthreads, provide an abstract programming model
659 which enables the programmer to hide the details of thread programming, meaning that threads do not have to be
660 managed directly as in Pthreads. Furthermore, high-level frameworks such as OpenMP, Cilk Plus, TBB, and SWARM
661 provide similar parallel programming constructs to express the map and reduction patterns. As for programming style,
662 parallel code in TBB, Cilk Plus, and FastFlow uses the C++ object-oriented style instead of the C procedural style
663 used in the remaining frameworks. The code in TBB, Cilk Plus, and FastFlow is written in C++ and uses object-
664 oriented principles and template functions. This has resulted in the initial sequential C code becoming a mixture of
665 C and C++ code. Moreover, Pthreads, OpenMP, and Intel’s two frameworks (Cilk Plus and TBB) support online
666 documentation, manuals, books, and slides and provide many code examples. In particular, there are many very good
667 teaching tutorials for the Pthreads and OpenMP frameworks. On the other hand, few manuals and code examples
668 exist for the SWARM and FastFlow frameworks. Finally, the learning curve for the Pthreads framework is long for
669 novice programmers because it requires low-level programming and because understanding the complex features and
670 function syntax takes a long time, whereas frameworks such as OpenMP, Cilk Plus, TBB, and SWARM have short
671 learning curves because programmers simply insert compiler directives or use available parallelisation constructs in
672 the sequential code, meaning that these frameworks require the least amount of time (a few days) to learn the basic
673 features of the model. The FastFlow framework has a moderately steep learning curve because programmers must
674 study much terminology to understand the functionality of the object-oriented design. All these characteristics of the
675 multi-core programming models are summarised in Table 5.

676 6. Conclusions

677 In this research, a set of linear algebra methods commonly used in the field of scientific computing has been par-
678 allelised using six multi-core programming models. Moreover, an extended quantitative and qualitative comparison
679 of these programming models has been performed to answer the question: which is the appropriate programming
680 model for implementing a scientific computing kernel on a multi-core system to achieve a balance between high pro-
681 grammer productivity and high performance? Based on this exploratory study, the following general conclusions can
682 be drawn. First, some programming models have a single-core overhead (i.e., Pthreads, Cilk Plus, TBB, and Fast-
683 Flow) compared to a sequential C program. However, this overhead has little effect on performance as the matrix size
684 grows. Second, a significant reduction in execution time occurs as the number of cores increases when the amount of
685 parallelisable work is high, such as for matrix addition, outer product, matrix-vector product, matrix product, matrix
686 transpose, Gaussian elimination, and Jacobi operations. On the other hand, the performance of vector addition and
687 dot product operations was poor because only a small amount of the work was parallelisable. Third, computationally
688 intensive kernels, except for vector addition and dot product, had excellent scalability as a function of matrix size.
689 Fourth, the best overall performance as a function of number of cores and of matrix size was shown by the SWARM

690 and OpenMP frameworks, followed by the Cilk Plus, TBB, and FastFlow models. Fifth, most high-level program-
691 ming models such as OpenMP, Cilk Plus, TBB, and SWARM support parallel patterns through ready-to-use parallel
692 constructs and require only a small number of additional or modified lines of code to be added to sequential code
693 with minimal programming effort and a short learning curve. On the other hand, the Pthreads and FastFlow models
694 do not support parallel patterns directly and require restructuring sequential code into parallel code with intensive
695 programming effort. Finally, as a general conclusion, the OpenMP and SWARM programming models are suitable
696 for parallelising computationally intensive methods and matrix operations of moderate size on multi-core platforms.
697 They require little programming effort and give satisfactory performance as a function of number of cores and of ma-
698 trix size. Furthermore, the FastFlow model can be used to implement large matrix operations such as matrix product
699 with good performance, but this requires more programming effort. A programmer interested in implementing paral-
700 lel methods with little programming effort for medium performance could use the Cilk Plus and TBB programming
701 models.

702 Finally, a number of directions for future research have been identified. First, the quantitative and qualitative
703 comparisons of programming models could be extended to other scientific computing algorithms such as matrix
704 factorisation and similar applications from computational statistics and data analysis. Second, the present quantitative
705 and qualitative evaluation could be extended to include advanced algorithmic techniques used for matrix operations
706 such as the tiling technique, the divide and conquer strategy, and cache-oblivious algorithms. When implementing
707 these algorithmic techniques, the performance of the parallel programming models could be compared with that of the
708 Intel MKL BLAS library. Third, the performance comparison of programming models could be extended to sparse
709 matrix computations with real-world sparse matrices from the Matrix Market repository [8]. Fourth, the scalability
710 analysis, where the problem size and the number of cores were increased concurrently, could be extended to other
711 matrix operations except for the outer and matrix products. Finally, a quantitative and rigorous analysis could be
712 performed by gathering data on the development time of matrix operations for the various parallel programming
713 models.

714 Acknowledgements

715 The authors would like to thank the anonymous reviewers for many helpful comments and suggestions, which
716 have greatly improved the presentation of this paper.

717 Appendix

718 The tables can be retrieved from the <http://users.uom.gr/~pmichailidis/jpapers/appendix-anova-ci.pdf>
719 pdf

720 References

- 721 [1] Intel Cilk Plus, <http://software.intel.com/en-us/articles/intel-cilk-plus/> (2012).
- 722 [2] Intel Threading Building Blocks, <http://threadingbuildingblocks.org/> (2012).
- 723 [3] The openMP API specification for parallel programming, <http://openmp.org/wp/> (2012).
- 724 [4] AMD microprocessors, http://en.wikipedia.org/wiki/List_of_AMD_microprocessors (2014).
- 725 [5] IBM Power8 architecture, <http://en.wikipedia.org/wiki/POWER8> (2014).
- 726 [6] Intel Math Kernel Library - Documentation, <http://software.intel.com/en-us/articles/intel-math-kernel-library-documentation> (2014).
- 727 [7] Intel Xeon, <http://en.wikipedia.org/wiki/Xeon> (2014).
- 728 [8] Matrix market repository, <http://math.nist.gov/MatrixMarket/> (2014).
- 729 [9] E. Agullo, J. Demmel, J. Dongarra, B. Hadri, J. Kurzak, J. Langou, H. Ltaief, P. Luszczek, S. Tomov, Numerical linear algebra on emerging
730 architectures: The PLASMA and MAGMA projects, *Journal of Physics: Conference Series* 180 (1) (2009) 012037.
- 731 [10] M. Aldinucci, M. Danelutto, P. Kilpatrick, M. Meneghin, M. Torquati, Accelerating code on multi-cores with FastFlow, in: *Proceedings of the*
732 *17th International Conference on Parallel processing - Volume Part II, Euro-Par'11*, Springer-Verlag, Berlin, Heidelberg, 2011, pp. 170–181.
- 733 [11] M. Aldinucci, M. Danelutto, P. Kilpatrick, M. Torquati, *Programming Multi-core and Many-core Computing Systems*, chap. FastFlow: high-
734 level and efficient streaming on multi-core, Wiley, 2011.
- 735 [12] M. Aldinucci, M. Meneghin, M. Torquati, Efficient Smith-Waterman on multi-core with FastFlow, in: *Proceedings of 18th Euromicro*
736 *International Conference on Parallel, Distributed and Network-Based Processing (PDP)*, 2010, pp. 195–199.
- 737 [13] A. Ali, U. Dastgeer, C. Kessler, OpenCL for programming shared memory multicore CPUs, in: *Proceedings of MULTIPROG-2012 Workshop*
738 *at HiPEAC-2012*, 2012.

- 739 [14] E. Anderson, Z. Bai, C. Bischof, L. S. Blackford, J. Demmel, J. J. Dongarra, J. Du Croz, S. Hammarling, A. Greenbaum, A. McKenney,
740 D. Sorensen, LAPACK Users' Guide (Third Ed.), Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1999.
- 741 [15] E. Ayguadé, N. Copt, A. Duran, J. Hoeflinger, Y. Lin, F. Massaioli, X. Teruel, P. Unnikrishnan, G. Zhang, The design of OpenMP tasks,
742 IEEE Trans. Parallel Distrib. Syst. 20 (3) (2009) 404–418.
- 743 [16] D. A. Bader, V. Kanade, K. Madduri, SWARM: A Parallel Programming Framework for Multicore Processors, in: Parallel and Distributed
744 Processing Symposium, 2007. IPDPS 2007, 2007, pp. 1–8.
- 745 [17] R. Barrett, M. Berry, T. F. Chan, J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine, H. V. der Vorst, Templates for the
746 Solution of Linear Systems: Building Blocks for Iterative Methods, 2nd Edition, SIAM, Philadelphia, PA, 1994.
- 747 [18] M. Birk, M. Zapf, M. Balzer, N. Rüter, J. Becker, A comprehensive comparison of GPU- and FPGA-based acceleration of reflection image
748 reconstruction for 3d ultrasound computer tomography, Journal of Real-Time Image Processing (2012) 1–12.
- 749 [19] H. Blaar, M. Legeler, T. Rauber, Efficiency of thread-parallel Java programs from scientific computing, in: Proceedings of the 16th Interna-
750 tional Parallel and Distributed Processing Symposium, IPDPS '02, IEEE Computer Society, Washington, DC, USA, 2002, pp. 131–.
- 751 [20] L. S. Blackford, J. Choi, A. Cleary, E. D'Azevedo, J. Demmel, I. Dhillon, S. Hammarling, G. Henry, A. Petitet, K. Stanley, D. Walker, R. C.
752 Whaley, ScaLAPACK User's Guide, Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 1997.
- 753 [21] R. D. Blumofe, C. F. Joerg, B. C. Kuszmaul, C. E. Leiserson, K. H. Randall, Y. Zhou, Cilk: An efficient multithreaded runtime system, in:
754 Proceedings of the Fifth ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP), Santa Barbara, California,
755 1995, pp. 207–216.
- 756 [22] J. M. Bull, Measuring Synchronisation and Scheduling Overheads in OpenMP, in: In Proceedings of First European Workshop on OpenMP,
757 1999, pp. 99–105.
- 758 [23] A. Buttari, J. Dongarra, P. Husbands, J. Kurzak, K. Yelick, Multithreading for synchronization tolerance in matrix factorization, in: Journal
759 of Physics: Conference Series 78(1), 2007.
- 760 [24] A. Buttari, J. Langou, J. Kurzak, J. Dongarra, A class of parallel tiled linear algebra algorithms for multicore architectures, Parallel Comput.
761 35 (1) (2009) 38–53.
- 762 [25] D. Buttlar, J. Farrell, B. Nichols, PThreads Programming: A POSIX Standard for Better Multiprocessing, O'Reilly Media, 1996.
- 763 [26] K. Chellapilla, S. Puri, P. Simard, High Performance Convolutional Neural Networks for Document Processing, in: G. Lorette (ed.), Tenth
764 International Workshop on Frontiers in Handwriting Recognition, Université de Rennes 1, Suvisoft, La Baule (France).
- 765 [27] T. Deepak Shekhar, K. Varaganti, R. Suresh, R. Garg, R. Ramamoorthy, Comparison of parallel programming models for multicore architec-
766 tures, in: Proceedings of IEEE International Symposium on Parallel and Distributed Processing Workshops and Phd Forum (IPDPSW), 2011,
767 pp. 1675–1682.
- 768 [28] J. W. Demmel, Applied Numerical Linear Algebra, Society for Industrial and Applied Mathematics, 1997.
- 769 [29] J. Diaz, C. Munoz-Caro, A. Nino, A survey of parallel programming models and tools in the multi and many-core era, IEEE Transactions on
770 Parallel and Distributed Systems 23 (8) (2012) 1369–1386.
- 771 [30] M. Z. F. Cantonnnet, Y. Yao, T. El-Ghazawi., Productivity analysis of the UPC language, in: 18th International Symposium on Parallel and
772 Distributed Processing, 2004.
- 773 [31] K.-F. Faxén, Wool-A work stealing library, SIGARCH Comput. Archit. News 36 (5) (2009) 93–100.
- 774 [32] M. Frigo, C. E. Leiserson, H. Prokop, S. Ramachandran, Cache-oblivious algorithms, ACM Trans. Algorithms 8 (1) (2012) 4:1–4:22.
- 775 [33] M. Frigo, V. Strumpfen, The cache complexity of multithreaded cache oblivious algorithms, Theory of Computing Systems 45 (2) (2009)
776 203–233.
- 777 [34] J. Gentle, Matrix Algebra: Theory, Computations, and Applications in Statistics, Springer Texts in Statistics, Springer, 2007.
- 778 [35] G. Golub, C. Van Loan, Matrix Computations, 2nd ed., Johns Hopkins University Press, 1989.
- 779 [36] J. Gomez-Pulido, M. Vega-Rodriguez, J. Sanchez-Perez, S. Priem-Mendes, V. Carreira, Accelerating floating-point fitness functions in evo-
780 lutionary algorithms: A FPGA-CPU-GPU performance comparison, Genetic Programming and Evolvable Machines 12 (4) (2011) 403–427.
- 781 [37] K. Goto, R. A. v. d. Geijn, Anatomy of high-performance matrix multiplication, ACM Trans. Math. Softw. 34 (3) (2008) 12:1–12:25.
- 782 [38] A. Grama, G. Karypis, V. Kumar, A. Gupta, Introduction to Parallel Computing, Pearson Education, 2nd ed., Addison-Wesley, 2003.
- 783 [39] A. G. Gray, Fast Kernel Matrix-Vector Multiplication with Application to Gaussian Process Regression, Tech. rep., Carnegie Mellon Univer-
784 sity (2004).
- 785 [40] C. Grozea, Z. Bankovic, P. Laskov, FPGA vs. Multi-core CPUs vs. GPUs: Hands-on experience with a sorting application, in: R. Keller,
786 D. Kramer, J.-P. Weiss (eds.), Facing the Multicore-Challenge, vol. 6310 of Lecture Notes in Computer Science, Springer Berlin Heidelberg,
787 2010, pp. 105–117.
- 788 [41] J. Gunnels, J. Lee, S. Margulies, Efficient high-precision matrix algebra on parallel architectures for nonlinear combinatorial optimization,
789 Mathematical Programming Computation 2 (2) (2010) 103–124.
- 790 [42] L. Hochstein, J. Carver, F. Shull, S. Asgari, V. Basili, J. K. Hollingsworth, M. V. Zelkowitz, Parallel programmer productivity: A case study
791 of novice parallel programmers, in: Proceedings of the 2005 ACM/IEEE Conference on Supercomputing, SC '05, IEEE Computer Society,
792 Washington, DC, USA, 2005, pp. 35–.
- 793 [43] H. Kasim, V. March, R. Zhang, S. See, Survey on parallel programming model, in: Proceedings of the IFIP International Conference on
794 Network and Parallel Computing, NPC '08, Springer-Verlag, Berlin, Heidelberg, 2008, pp. 266–275.
- 795 [44] P. Kegel, M. Schellmann, S. Gortlatch, Comparing programming models for medical imaging on multi-core systems, Concurrency and Com-
796 putation: Practice and Experience 23 (10) (2011) 1051–1065.
- 797 [45] S. Kesturt, J. Davis, O. Williams, BLAS comparison on FPGA, CPU and GPU, 2010, pp. 288–293.
- 798 [46] J. Kurzak, H. Ltaief, J. Dongarra, R. M. Badia, Scheduling dense linear algebra operations on multicore processors, Concurr. Comput. : Pract.
799 Exper. 22 (1) (2010) 15–44.
- 800 [47] X. S. Li, Evaluation of superLU on multicore architectures, Journal of Physics: Conference Series 125 (1) (2008) 012079.
- 801 [48] M. Luff, Empirically investigating parallel programming paradigms: A null result, in: Workshop on Evaluation and Usability of Programming
802 Languages and Tools (PLATEAU), 2009.
- 803 [49] A. Marowka, Empirical analysis of parallelism overheads on CMPs, in: Proceedings of the 8th International Conference on Parallel processing

- 804 and applied mathematics: Part I, PPAM'09, Springer-Verlag, Berlin, Heidelberg, 2010, pp. 596–605.
- 805 [50] M. McCool, J. Reinders, A. Robison, *Structured Parallel Programming: Patterns for Efficient Computation*, Elsevier Science, 2012.
- 806 [51] S. F. McGinn, R. E. Shaw, Parallel Gaussian elimination using OpenMP and MPI, in: *Proceedings of the 16th Annual International Sym-*
807 *posium on High Performance Computing Systems and Applications, HPCS '02*, IEEE Computer Society, Washington, DC, USA, 2002, pp.
808 169–173.
- 809 [52] R. Membarth, F. Hannig, J. Teich, M. Krner, W. Eckert, Frameworks for multi-core architectures: A comprehensive evaluation using 2d/3d
810 image registration, in: M. Berekovic, W. Fornaciari, U. Brinkschulte, C. Silvano (eds.), *Architecture of Computing Systems - ARCS 2011*,
811 vol. 6566 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2011, pp. 62–73.
- 812 [53] P. D. Michailidis, K. G. Margaritis, Implementing parallel LU factorization with pipelining on a multicore using OpenMP, in: *Proceedings*
813 *of the 2010 13th IEEE International Conference on Computational Science and Engineering, CSE '10*, IEEE Computer Society, Washington,
814 DC, USA, 2010, pp. 253–260.
- 815 [54] P. D. Michailidis, K. G. Margaritis, Performance models for matrix computations on multicore processors using OpenMP, in: *Proceedings of*
816 *International Conference on Parallel and Distributed Computing Applications and Technologies*, IEEE Computer Society, Los Alamitos, CA,
817 USA, 2010, pp. 375–380.
- 818 [55] P. D. Michailidis, K. G. Margaritis, Open multi processing (OpenMP) of Gauss-Jordan method for solving system of linear equations, in:
819 *Proceedings of IEEE 11th International Conference on Computer and Information Technology (CIT)*, 2011, pp. 314 – 319.
- 820 [56] P. D. Michailidis, K. G. Margaritis, Computational comparison of some multi-core programming tools for basic matrix computations, in: *Pro-*
821 *ceedings of 9th IEEE International Conference on Embedded Software and Systems (HPCC-ICISS)*, and *IEEE 14th International Conference*
822 *on High Performance Computing and Communication*, 2012, pp. 143–150.
- 823 [57] P. D. Michailidis, K. G. Margaritis, Implementing basic computational kernels of linear algebra on multicore, in: *Proceedings of 16th*
824 *Panhellenic Conference on Informatics*, IEEE Computer Society, Los Alamitos, CA, USA, 2012, pp. 217–222.
- 825 [58] P. D. Michailidis, K. G. Margaritis, Performance study of matrix computations using multi-core programming tools, in: M. Ivanovic, Z. Budi-
826 mac, M. Radovanovic (eds.), *Proceedings of 5th Balkan Conference in Informatics*, ACM, 2012, pp. 186–192.
- 827 [59] A. Podobas, M. Brorsson, A comparison of some recent task-based parallel programming models, in: *Proceedings of the 3rd Workshop on*
828 *Programmability Issues for Multi-Core Computers, (MULTIPROG'2010)*, Jan 2010, Pisa, 2010.
- 829 [60] T. Rauber, G. Runger, *Parallel Programming - for Multicore and Cluster Systems*, Springer, 2010.
- 830 [61] S. Ravela, Comparison of shared memory based parallel programming models, Master's thesis, Biekinge Institute of Technology (2010).
- 831 [62] J. Reinders, *Intel Threading Building Blocks - Outfitting C++ for Multi-core Processor Parallelism*, O'Reilly, 2007.
- 832 [63] G. Runger, M. Schwind, Fast recursive matrix multiplication for multi-core architectures, *Procedia CS* 1 (1) (2010) 67–76.
- 833 [64] C. Sadowski, A. Shewmaker, The last mile: Parallel programming and usability, in: *Proceedings of the FSE/SDP Workshop on Future of*
834 *Software Engineering Research, FoSER '10*, ACM, New York, NY, USA, 2010, pp. 309–314.
- 835 [65] L. Sanchez, J. Fernandez, R. Sotomayor, J. Garcia, A comparative evaluation of parallel programming models for shared-memory architec-
836 tures, in: *Proceedings of IEEE 10th International Symposium on Parallel and Distributed Processing with Applications (ISPA)*, 2012, pp. 363
837 –370.
- 838 [66] W.-H. Steeb, Y. Hardy, A. Hardy, R. Stoop, *Problems & Solutions in Scientific Computing with C++ and Java Simulations*, World Scientific
839 Publishing Co., Inc., River Edge, NJ, USA, 2004.
- 840 [67] D. Szafron, J. Schaeffer, An experiment to measure the usability of parallel programming systems, *Concurrency: Practice and Experience*
841 8 (2) (1996) 147–166.
- 842 [68] D. Tamir, N. Shaked, W. Geerts, S. Dolev, Parallel decomposition of combinatorial optimization problems using electro-optical vector by
843 matrix multiplication architecture, *The Journal of Supercomputing* 62 (2) (2012) 633–655.
- 844 [69] R. A. van de Geijn, *Using LAPACK: Parallel Linear Algebra Package*, MIT Press, Cambridge, MA, USA, 1997.
- 845 [70] M. M. Wolf, M. T. Heath, Combinatorial Optimization of Matrix-Vector Multiplication in Finite Element Assembly, *SIAM Journal on*
846 *Scientific Computing* 31 (4) (2009) 2960–2980.
- 847 [71] D. Zou, Y. Dou, F. Xia, Optimization schemes and performance evaluation of smith-waterman algorithm on CPU, GPU and FPGA, *Concurr.*
848 *Comput. : Pract. Exper.* 24 (14) (2012) 1625–1644.
- 849 [72] S. Zuckerman, M. Perache, W. Jalby, Fine tuning matrix multiplications on multicore, in: *Proceedings of the 15th international conference*
850 *on High performance computing, HiPC'08*, Springer-Verlag, Berlin, Heidelberg, 2008, pp. 30–41.