

Performance and Energy Analysis of OpenMP Runtime Systems with Dense Linear Algebra Algorithms

Joao Vicente Ferreira Lima, Issam Raïs, Laurent Lefèvre, Thierry Gautier

► **To cite this version:**

Joao Vicente Ferreira Lima, Issam Raïs, Laurent Lefèvre, Thierry Gautier. Performance and Energy Analysis of OpenMP Runtime Systems with Dense Linear Algebra Algorithms. International Journal of High Performance Computing Applications, SAGE Publications, 2019, 33 (3), pp.431-443. 10.1177/1094342018792079 . hal-01957220

HAL Id: hal-01957220

<https://hal.inria.fr/hal-01957220>

Submitted on 17 Dec 2018

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Performance and Energy Analysis of OpenMP Runtime Systems with Dense Linear Algebra Algorithms

João V. F. Lima¹, Issam Raïs², Laurent Lefèvre², and Thierry Gautier²

Abstract

In this paper, we analyse performance and energy consumption of five OpenMP runtime systems over a NUMA platform. We also selected three CPU level optimizations, or techniques, to evaluate their impact on the runtime systems: processors features Turbo Boost and C-States, and CPU DVFS through Linux CPUFreq governors. We present an experimental study to characterize OpenMP runtime systems on the three main kernels in dense linear algebra algorithms (Cholesky, LU and QR) in terms of performance and energy consumption. Our experimental results suggest that OpenMP runtime systems can be considered as a new energy leverage, and Turbo Boost, as well as C-States, impacted significantly performance and energy. CPUFreq governors had more impact with Turbo Boost disabled, since both optimizations reduced performance due to CPU thermal limits. A LU factorization with concurrent write extension from libKOMP achieved up to 63% of performance gain and 29% of energy decrease.

Keywords

OpenMP, task parallelism, linear-algebra algorithms, NUMA, energy efficiency

1. Introduction

Energy efficiency is one of the five major challenges that should be overcome in the path to exascale computing ([Bergman et al. 2008](#)). Despite improvements in energy efficiency, the total energy consumed by supercomputers is still increasing due to the even quicker increase in computational power. High energy consumption is not only a problem of electricity costs, but it also impacts greenhouse emissions

¹Universidade Federal de Santa Maria, Santa Maria, Rio Grande do Sul, Brazil

²Univ. Lyon, Inria, CNRS, ENS de Lyon, Univ. Claude-Bernard Lyon 1, LIP, France

Corresponding author:

João V. F. Lima, Centro de Tecnologia, Prédio 07, Anexo B - Sala 374, Avenida Roraima 1000, Santa Maria, RS, Brazil.
Email: jvlima@inf.ufsm.br

and dissipating the produced heat can be difficult. As the ability to track power consumption becomes more commonplace, with some job schedulers supporting tracking energy use (Yang et al. 2013), soon users of HPC systems may have to consider both how many CPU hours they need and how much energy.

Energy budget limitation imposes a high pressure to the HPC community making energy consideration a prominent research field. Most of the gain will come from technology by providing more energy efficient hardware, memory and interconnect. Nevertheless, recent processors integrate more and more leverages to reduce energy consumption (e.g. classical DVFS, deep sleep states) and low level runtime algorithms provide orthogonal leverages (e.g. dynamic concurrency throttling). However few of these leverages are integrated and employed in today local level software stack such as middleware, operating system or runtime library. Due to the complexity of this statement, we restricted our investigation to local node energy consumption by HPC OpenMP applications.

OpenMP is an API standard to express parallel portable programs. Most of controls are implementation defined and rely on the specific OpenMP programming environment used. The OpenMP standard does not impose any constraint on implementations. Even if there are more precise specifications, e.g. mapping of threads to cores, it is very tricky to precisely control performance or energy consumption using what OpenMP specification proposes (Bari et al. 2016). Previous works have dealt with a specific OpenMP runtime (Porterfield et al. 2013a; Nandamuri et al. 2014; Su et al. 2012; Marathe et al. 2015; Lively et al. 2011; Li et al. 2010) that may be difficult to generalize to other OpenMP runtime systems without strong development effort. To the knowledge of the authors, there is no related work comparing OpenMP runtime systems in order to analyse performance and energy consumption.

In this paper, we analysed performance and energy consumption of five OpenMP runtime systems over a NUMA system. We also selected three CPU level optimizations, or techniques, to evaluate their impact on the runtime systems: processors features Turbo Boost and C-States, and CPU DVFS through Linux CPUFreq governors. We restrict our experiments on three dense linear algebra algorithms: Cholesky, LU and QR matrix factorizations. Source codes are based on KASTORS (Virouleau et al. 2014) benchmark suite and the state of the art PLASMA library using its new OpenMP implementations (YarKhan et al. 2016) that rely on OpenMP tasks with data dependencies.

The paper is an extended version of the paper presented at the 8th Workshop on Applications for Multi-core Architectures (WAMCA 2017). The contributions of this paper are:

- We present early experiments of performance and energy consumption over OpenMP runtime systems;
- We report the impact of three CPU level optimizations in order to present the respective gains with different combinations;
- We observed that a LU factorization with concurrent-write access mode achieved up to 63% in performance gain and 29% in energy over original LU algorithm;
- In addition, our findings suggest that Turbo Boost and C-States had significant impact on performance and energy. CPUFreq governors had more impact with Turbo Boost disabled since Turbo Boost with the *performance* governor reduced performance due to CPU thermal limits.

The remainder of the paper is organized as follows. Section 2 presents the related work. Section 3 gives some details of the OpenMP task programming model and an overview about five runtime implementations. Section 4 details the experimental hardware and methodology used. Our experimental

results are presented in Section 5. Finally, Section 6 and Section 7, respectively, present the discussion and conclude the paper.

2. Related work

Multiple techniques, or *leverages*, dealing with the energy-performance trade-off are exposed and used in the literature. In (Benoit et al. 2017), authors demonstrate the challenge of using the shutdown and wakeup leverages for large scale HPC infrastructure without altering the throughput of needed computation. In (Ribic and Liu 2014), Dynamic Voltage and Frequency Scaling (DVFS) is used to lower the speed of threads that are not in the critical path with a warranty on performance. In (Etinski et al. 2010), authors bind DVFS with EASY backfilling job scheduling to answer possible system load variation in an energy-efficient way. Considering communications between node, the authors of (Rountree et al. 2009) reduce the frequency of tasks which would block for MPI communication. In (Laros et al. 2012) the authors analyses the effects of both DVFS and network bandwidth scaling

Other works use the simplicity proposed by OpenMP to vary the number of threads, for energy efficiency. Authors in (Curtis-Maury et al. 2006) and (Porterfield et al. 2013b) defend the Dynamic Concurrency Throttling (DCT) and underline the fact that using OpenMP to control the number of threads could be energy efficient, depending on the algorithm or the chosen hardware.

Previous works show that various energy behaviors of computing nodes are possible through various leverages (DVFS, DCT, etc). But none of the previous work focus on OpenMP runtime systems as a leverage. None of the previous work dealt with the energy-performance trade-off and thus underlined possible variability concerning energy and performance for existing runtime systems. Thus, to the knowledge of the authors, no related work were trying to compare several OpenMP runtime libraries together for various representative workloads, as presented in our paper.

We use state of the art PLASMA library (YarKhan et al. 2016), on three main kernels in dense linear algebra (Cholesky, LU and QR factorizations), that implements dependent tasks model. This model is new and never addressed in related works. In (Porterfield et al. 2013a; Nandamuri et al. 2014) the authors based their experiments using the BOTS (Duran et al. 2009) benchmarks that require only the independent tasks.

3. OpenMP Task programming model and implementations

In 2013 the OpenMP Architecture Review Board introduced in the OpenMP revision 4.0 a new way of expressing task parallelism using OpenMP, through the task dependencies. This section introduces the task dependency programming model targeted by the selected benchmark suites. We also present how the model is implemented in various runtime libraries.

3.1. Dependent task model

OpenMP dependent task model allows to define dependencies between tasks using declaration of accesses to memory with *in,out*, or *inout*. Two tasks are independent (or concurrent) if and only if they do not violated the data dependencies of a reference sequential execution order*.

*OpenMP does not allows variable renaming to suppress output and anti-dependencies.

```

1 for (k=0; k<NB; k++) {
2 #pragma omp task untied shared(M) \
3   depend(inout: M[k*NB+k])
4   lu0(M[k*NB+k]);
5   for (j=k+1; j<NB; j++)
6 #pragma omp task untied shared(M) \
7   depend(in: M[k*NB+k]) depend(inout: M[k*NB+j])
8     fwd(M[k*NB+k], M[k*NB+j]);
9
10  for (i=k+1; i<NB; i++)
11 #pragma omp task untied shared(M)\
12   depend(in: M[k*NB+k]) depend(inout: M[i*NB+k])
13     bdiv(M[k*NB+k], M[i*NB+k]);
14
15  for (i=k+1; i<NB; i++)
16    for (j=k+1; j<NB; j++)
17 #pragma omp task untied shared(M)\
18   depend(in:M[i*NB+k], M[k*NB+j]) depend(inout:M
19     [i*NB+j])
20     bmod(M[i*NB+k],M[k*NB+j],M[i*NB+j]);
21 }

```

Figure 1. LU factorization with OpenMP dependent task.

Table 1. Characteristics of OpenMP runtime systems.

Name	Dependencies	Task Scheduling	Remarks
libGOMP	hash table	centralized list	task throttling
libOMP	hash table	work stealing	bounded dequeue
OmpSs	hash table	socket-aware work stealing	
XKaapi	hash table*	non blocking work stealing	task affinity
libKOMP	resizable hash table	non blocking work stealing	task affinity concurrent write

Figure 1 illustrates a LU factorization based on PLASMA (YarKhan et al. 2016). The programmer declares tasks and the accesses *in*, *inout* they made to a memory region (here only *lvalue* or memory reference, *i.e.* pointer).

The OpenMP library computes tasks and dependencies at runtime, and schedules concurrent tasks on the available processors. The strategy for task dependencies and task scheduling depends on the runtime implementation. Nevertheless, their implementations impact the performance and the energy consumption. Moreover, the absence of precise OpenMP specification about the task scheduling algorithm is the key point to allow research to improve performance and energy efficiency with implementation concerns.

3.2. Runtime system implementations

Table 1 summarizes the properties of five OpenMP runtime systems.

libGOMP is the OpenMP runtime that comes with the GCC compiler. Dependencies between tasks are computed through a hash table that map data (pointer) to the last task writing the data. Ready tasks are pushed into several scheduling dequeues. The main dequeue stores all the tasks generated by the threads of a parallel region. Tasks seem to be inserted after the position of their parent tasks in order to keep an order close to the sequential execution order. Because threads share the main dequeue, serialization of operations is guaranteed by a pthread mutex which is a bottleneck for scalability. To avoid overhead in task creation, libGOMP implements a task throttling algorithm that serialize task creation when the number of pending tasks is greater than a threshold proportional to the number of threads.

libOMP was initially the proprietary OpenMP runtime of Intel for its C, C++ and Fortran compilers. Now it is also the target runtime for the LLVM/Clang compiler and sources are open to community. libOMP manages dependencies in the same way that libGOMP by using a hash table. Memory allocation during task creation relies on a fast thread memory allocator. libOMP task scheduling is based on Cilk almost non blocking work stealing algorithm (Frigo et al. 1998), but dequeue operations are serialized using locks. Nevertheless, it implies distributed dequeues management with high throughput of dequeue operations. libOMP also implements a task throttling algorithm by using bounded size dequeue.

OmpSs (Bueno-Hedo et al. 2012) is a runtime system developed at the Barcelona Supercomputing Center, compatible with the OpenMP specification. It has a specific compiler, called Mercurium, that transforms OpenMP directives to calls to Nanos++ runtime entrypoints. As libGOMP and libOMP, OmpSs computes task dependencies at task creation using hash map. In our experiments we select the breadth-first scheduler. In our experiments we selected the Socket-aware scheduler (*socket*) which is a work stealing based task scheduler with distributed dequeues implementations.

XKaaapi (Gautier et al. 2013) is a task library for multi-CPU and multi-GPU architectures which provides binary compatible library with libGOMP (Broquedis et al. 2012). Task scheduling is based on the almost non blocking work stealing algorithm from Cilk (Frigo et al. 1998) with extension to combine steal requests in order to reduce overhead in stealing (Tchiboukdjian et al. 2013). Moreover, XKaaapi computes dependencies on steal request, which is a perfect application of the work first principle to report overhead in task creation to critical path. The XKaaapi based OpenMP runtime also has support to some OpenMP extensions such as task affinity (Virouleau et al. 2016) that allows to schedule tasks on NUMA architecture, and to increase performance by reducing memory transfer and thus memory energy consumption.

libKOMP (Gautier and Virouleau 2015) is a redesign of (Broquedis et al. 2012) on a top of the Intel runtime libOMP. It includes following features coming mainly from XKaaapi: the dequeue management and work stealing with request combining; task affinity specific work stealing heuristic; a dynamically resized hash map that avoid high conflicts when finding dependencies for large tasks' graph; and tracing tool based on the OpenMP OMPT API; and finally a task concurrent write extension with a Clang modification [†] to provide the OpenMP directive clause. This latter extension allows better parallelism and was used in one of our LU benchmark and it very closed of the *task reduction* feature currently under discussion in the OpenMP architecture review board.

[†]<http://gitlab.inria.fr/openmp/clang>

3.3. Discussion

In our study of the mentioned OpenMP runtime systems, none of them include energy leverage such as thread throttling or DVFS. Nevertheless, their different task scheduling algorithms may impact energy efficiency. The main dequeue accesses in libGOMP serialize threads using a POSIX mutex. On Linux the mutex will block waiting threads after short period of active polling which ensure that few core cycles will be waste in the synchronisation.

On the other hand, libOMP, XKaapi and libKOMP work stealing actively poll dequeues until the program ends or a task is found. In order to reduce activity during polling, libOMP and libKOMP may block threads after an unsuccessful search of work by 200ms (default value). Once work is found, all threads are waked up.

4. Tools and Methods

This Section details the hardware configuration we experimented on and the OpenMP runtime systems we compared. We also give hints about the methodology used to process the collected data using statistical tools R.

4.1. Evaluation platform

Our experimental platform was a SGI UV2000 machine composed of eight NUMA nodes with one Intel Xeon E5-4617 (Sandy Bridge) processor each (total 8 processors) and 6 cores per processor (48 cores total) running at 2.9 GHz or 3.2GHz with Turbo Boost, and 512 GB of main memory. The processor has Turbo Boost 2.0 technology, and six idle states of C-States available: POLL C1-SNB C1E-SNB C3-SNB C6-SNB C7-SNB. The operating system was a Debian with Linux kernel 4.9.0-1 with two CPUFreq governors available: powersave and performance. Both DVFS governors have frequency limits from 1.20 GHz to 3.40 GHz.

4.2. Software description

4.2.1. Benchmarks We used kernels from two benchmark suites: the KASTORS (Virouleau et al. 2014) benchmark suite and an OpenMP-parallelized PLASMA version (YarKhan et al. 2016). Both benchmark suites tackle the same computational problems but use different algorithms in some cases. KASTORS was built from PLASMA 2.6.0 (released in dec. 2013) at a time when PLASMA parallelism was supported by a specific task management library called QUARK.

We focused our study on three dense linear algebra kernels:

- A Cholesky factorization (dpotrf);
- A LU factorization (dgetrf);
- A QR factorization (dgeqrf).

Cholesky factorization algorithms in both the benchmark suites are the same. We compare OpenMP based PLASMA version 82f89ee[‡]. All these linear algebra kernels we used rely on the BLAS routines, we used the implementation of OpenBLAS version 0.2.19.

[‡]Mercurial hash from <https://bitbucket.org/icl/plasma>

4.2.2. *Runtime Systems* We compared the following runtime systems during our experiments:

- LibGOMP – the OpenMP implementation from GNU that comes with GCC 6.3.0.
- LibOMP – a port of the Intel OpenMP open-source runtime to LLVM release 4.0.
- LibKOMP (Gautier and Virouleau 2015) – a research runtime system, based on the Intel OpenMP runtime, developed at INRIA. It offers several non-standard extensions to OpenMP. We evaluate the concurrent write (CW) feature in our experiments coupled with Cilk T.H.E work stealing protocol. We make experiments with version 54f7a28[§].
- XKaapi (Gautier et al. 2013) – research runtime system developed at INRIA. It has lightweight task creation overhead, and it offers several non-standard extensions to OpenMP (Broquedis et al. 2012) We evaluate its version efa5fdf[¶].
- OmpSs (Bueno-Hedo et al. 2012) - a runtime system developed at the Barcelona Supercomputing Center. The reported results are based on the 17.12 version^{||}.

In all experiments, we set `OMP_WAIT_POLICY` environment variable to *passive*, which means that threads should not consume CPU power while waiting. OmpSs has an equivalent option to Nano++ through `--enabled-sleep`. XKaapi runtime does not implement a waiting policy.

4.3. Energy measurement methodology

Since several metrics have to be considered depending on the objective, we consider performance (GFlop/s) and energy consumption (energy-to-solution). GFlop/s is measured by the each benchmark itself: it corresponds to the algorithmic count of the number of floating point operations over the elapsed time, using fact that matrix-matrix product does not rely on a fast algorithm such as Strassen like algorithm. Times are get using the Linux `clock_gettime` function with `CLOCK_REALTIME` clock.

We employed the Intel RAPL (Running Average Power Limit) feature as source of data acquisition for energy measurement. It exposes the energy consumption of several components on the chip (such as the processor package and the DRAM) through MSR (Model Specific Registers). Due to access limitation of MSR on the tested system, we designed a small tool querying periodically the RAPL counters based on LIKWID (Treibig et al. 2011): Energy consumption for the whole package (`PWR_PKG_ENERGY`), for the cores (`PWR_PP0_ENERGY`), for the DRAM (`PWR_DRAM_ENERGY`), as well as the core temperature (`TEMP_CORE`). The tool gets the counter values periodically and associate them with a timestamp.

4.4. Experimental methodology

All benchmarks are composed of two steps: the first allocates and initializes a matrix; the second step is the computation. We report execution time only from the computation step. Each experiment was repeated at least 30 times, each computation on a newly random matrix (as implemented by the benchmark). In parallel of the computation, we monitored the system by collecting various energy counters from RAPL.

[§]Git repository: <https://gitlab.inria.fr/openmp/libkomp>

[¶]Available at: <http://kaapi.gforge.inria.fr>

^{||}Available at: <https://pm.bsc.es>

We selected three CPU level optimizations (Orgerie et al. 2014) in order to evaluate their impact on performance and energy over the runtime systems: processor features Turbo Boost and C-States (Rotem et al. 2012), and CPU DVFS through Linux CPUFreq governors. Intel Turbo Boost is an overclocking mechanism that allows to the processor to raise core frequencies as long as the thermal and power limits are not exceeded. The C-States feature corresponds to CPU idle states. Deeper C-States offer more power savings, but at the cost of longer latency to enter and exit the C-State. On Linux, CPUFreq allows the control of P-States, which defines the frequencies at which a processor can operate, by governors that choose a frequency for the processor to use. The available governors on our experimental platform were: *powersave* that chooses the lowest frequency; and *performance*, which chooses the highest frequency.

For each computation we collected performance (GFlop/s) timestamped by the beginning and the end of the computation. This two timestamps were used in data post-processing to compute energy consumed by the computation between the two timestamps. Values were interpolated by linear function if missing in the collected energy values sampled periodically. Post-processing employs R script to compute energy per computation and to output basic statistic for each configuration. In our experimental results, energy values were the mean computed among the at least 30 computations of each configuration.

5. Experimental results

The goal of our experiments is to evaluate performance and energy consumption of OpenMP runtime systems and the impact of three CPU level optimizations. Our objectives are:

1. Evaluate the runtime impact on performance and energy, as well as the CPU level techniques (Sec. 5.1);
2. Analyse the correlation coefficient of CPU optimizations over performance and energy (Sec. 5.2);
3. Assess the impact of runtime extensions on performance of LU (Sec. 5.3).

The presented runtime systems have been experimented on the two benchmark suites presented in section 4.2.1. We build two configurations of libKOMP using two sets of options (Gautier and Virouleau 2015). On the following *komp* refers to libKOMP configured with T.H.E Cilk work stealing queue and requests combining protocol; and *komp+cw* is the same configuration than libKOMP with addition to support concurrent write extension used in the KASTORS LU code *dgetrf* (Virouleau et al. 2014).

5.1. Runtime impact

Figure 2 shows performance and energy results with a matrix size of 32768×32768 and over all machine resources available. Each configuration color represents the combination of the three CPU level optimizations Turbo Boost (*on* and *off*), CPUFreq governor (*performance* and *powersave*), and C-States (*all enabled* and *none*). Table 2 summarises the best results over each metric collected, *i.e.* performance, DRAM energy counter, PKG energy counter, and DRAM+PKG energy counters.

In all cases libkomp attained the best performance results, followed by: xkaapi and gcc for Cholesky; xkaapi and likomp (without CW) for LU; plasma and xkaapi for QR. OmpSs had similar performance to plasma on Cholesky, and performed worst than others on LU. Nonetheless, we were not able to evaluate ompss with QR due to a runtime error. The performance gains of libkomp over GCC were 9.7% for Cholesky, 63.8% for LU, and 9.4% for QR. In addition, the CPU level optimizations for these results

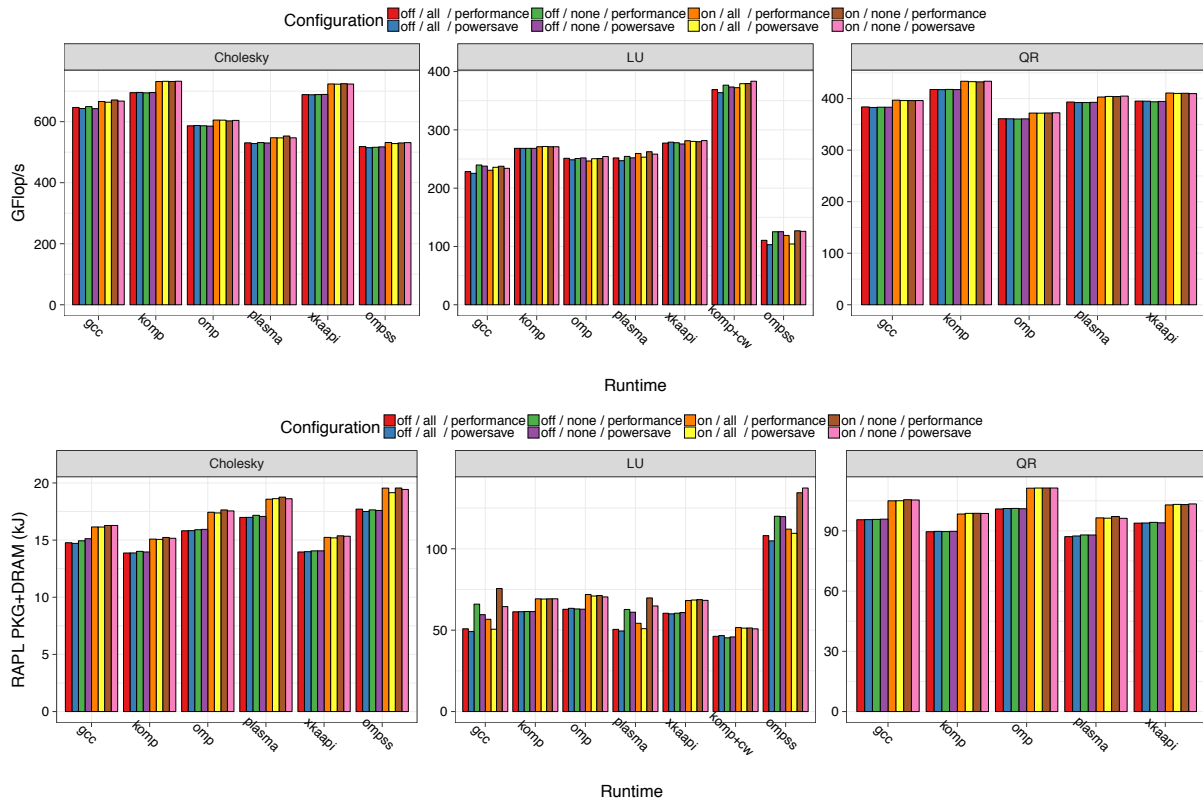


Figure 2. Performance (top) and energy (bottom) results of Cholesky, LU and QR over the UV2000 machine. The matrix size was 32768x32768 with 352x352 of block size. On energy the reported results are the sum of RAPL counter PKG and DRAM.

were: Turbo Boost enabled; powersave governor; and all C-States disabled. Clearly, the concurrent-write feature of LU contributed to the significant gain of libkomp.

In energy libkomp had better energy results with Cholesky and LU, while QR had lower energy consumption with plasma implementation and GCC. The CPU level optimizations for these results were: Turbo Boost disabled; performance governor; and all C-States enabled for Cholesky and QR, and all disabled for LU.

The cost of energy efficient cases in performance was not significant on the three benchmarks. Comparing the best cases of energy over the best cases in performance, Cholesky had a reduction of 5.12% in performance, while LU and QR had 1.76% and 9.31% reduction respectively. Although, energy reduction was of 8.51%, 10.88% and 11.71% for Cholesky, LU, and QR respectively.

Regarding the CPU level techniques, it seems that Turbo Boost and C-States contributed to the performance gains, and CPUFreq governors had more impact on energy than performance. It was expected that the best performance cases had Turbo enabled and all C-States disabled; still, those cases had powersave as CPU DVFS governor. Besides, energy efficient cases had in most cases Turbo disabled, all C-States enabled, and performance as CPU frequency governor.

Table 2. Overview of the best results. Higher is better for performance, and lower is better in energy.

Method	Runtime	Turbo Boost	C-States	Governor	GFlop/s	RAPL PKG	RAPL DRAM
PERFORMANCE							
Cholesky	komp	enabled	none	powersave	732.45	13.02	2.13
LU	komp+cw	enabled	none	powersave	383.40	42.83	7.98
QR	komp	enabled	none	powersave	433.77	83.36	15.31
RAPL PKG+DRAM							
Cholesky	komp	disabled	all	performance	694.95	11.62	2.24
LU	komp+cw	disabled	none	performance	376.64	37.12	8.16
QR	plasma	disabled	all	performance	393.40	71.42	15.70
RAPL PKG							
Cholesky	komp	disabled	all	performance	694.95	11.62	2.24
LU	gcc	disabled	all	powersave	225.41	36.60	12.53
QR	plasma	disabled	all	performance	393.40	71.42	15.70
RAPL DRAM							
Cholesky	xkaapi	enabled	none	powersave	723.23	13.22	2.11
LU	komp+cw	enabled	none	powersave	383.40	42.83	7.98
QR	plasma	enabled	none	powersave	404.92	81.07	15.19

In order to investigate the performance gains of powersave governor, we analysed the core temperature counter from RAPL over Cholesky with libkomp runtime as illustrated in Figure 3. We collected each socket temperature through a series of Cholesky executions, including an interval between executions of 30 seconds, and computed the mean of all readings per socket. The performance governor had greater temperature readings than powersave on all sockets using Turbo Boost, while both governors had similar temperature readings without Turbo Boost. It seems that tuning all CPU level optimizations to target performance reached the thermal limits of the processor sockets, and degraded performance.

5.2. Correlation analysis

We used the Person correlation coefficient to test the correlation between two variables X and Y in order to identify the CPU level techniques and their relation with performance and energy. This coefficient has values between -1 and $+1$ where $+1$ means a perfect positive linear correlation, 0 is no linear correlation, and -1 a total negative linear correlation. In addition, we added a significance test of the correlation with confidence interval of 95%.

Figure 4 shows a correlation graph for the best performance cases on the three benchmarks. We employ the following numerical values for each CPU level optimization: Turbo Boost on ($+1$) and off (-1); C-States all enabled ($+1$) and all disabled (-1); CPUFreq governor performance ($+1$) and powersave (-1).

It seems that Turbo Boost parameter had direct impact on performance and energy with an almost perfect correlation on Cholesky and QR. An exception was LU with a lower correlation of 0.31 on performance. In all cases, the negative correlation between Turbo Boost and RAPL DRAM means that it reduced DRAM energy when enabled while increased PKG consumption.

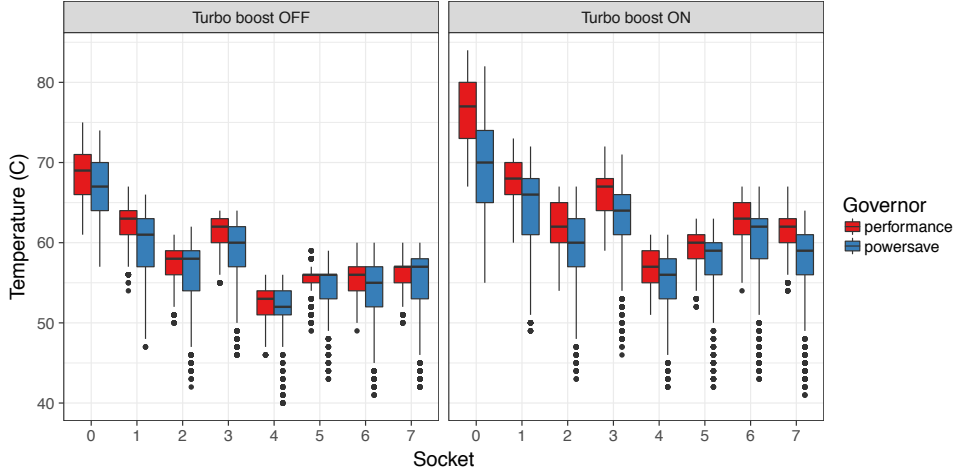


Figure 3. Processor temperature of each socket with Cholesky and libkomp runtime. We compared the impact of CPUFreq governor and Turbo Boost. C-States were disabled.

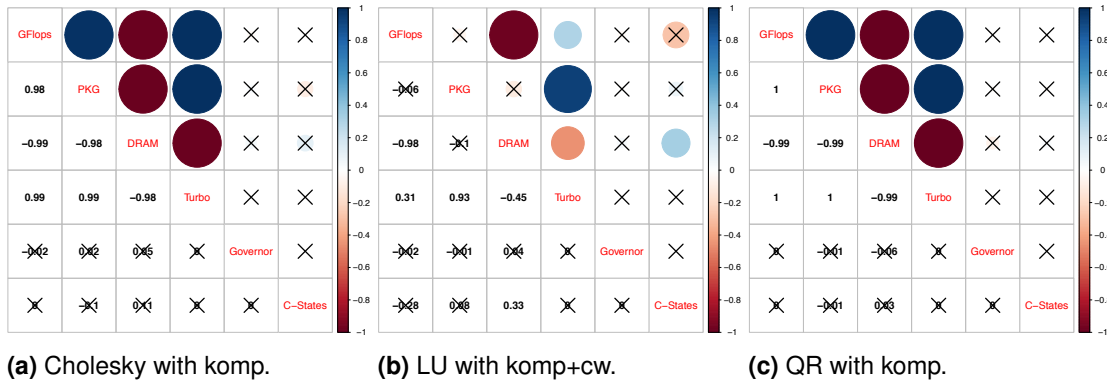


Figure 4. Correlation coefficients on the three benchmarks based on the best performance results. The \times symbol mean that the significance level (or p-value) of the correlation is under 5%.

The correlation results also show that CPUFreq governors and C-States had no relation with performance and energy. It appears that CPUFreq governor parameter was not relevant in our comparison of performance and energy in Table 2. On the other hand, C-States had a correlation of 0.33 with RAPL DRAM over LU, *i.e.* increasing DRAM energy consumption with all C-States enabled. It explains the best energy case of LU with C-States disabled in Table 2.

5.3. Focus on LU factorization

Figure 5 shows performance results for LU factorization with different matrix sizes. We used as reference the GCC runtime to compute the difference over the other three runtime systems, represented on the bar

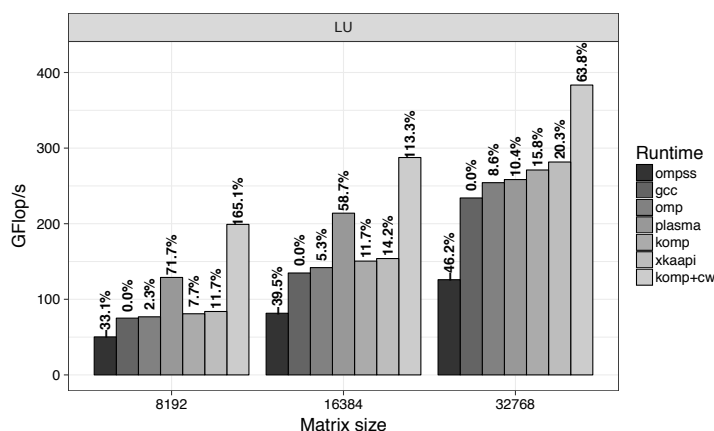


Figure 5. Performance (GFlop/s) results of LU. All percentages on top of bar plots are the difference of current runtime over GCC libGOMP runtime.

plots by a percentage value. The LU algorithm with concurrent-write showed significant improvement compared to other runtime systems (up to 165.1% over gcc), followed by PLASMA LU algorithm.

Thanks to the concurrent write, the LU algorithm with libkomp_cw runtime had more parallelism than other runtime systems due to the CW algorithm extension based on KASTORS (Virouleau et al. 2014). Figure 6 illustrates a Gantt execution from the LU factorization using PLASMA, libkomp and libkomp with CW.

On the LU factorization, even if CW generates more parallelism, the algorithm has poor efficiency and threads are frequently idle. The Gantt diagram on all the 48 cores illustrates long periods of inactivity. GCC is the only runtime where threads lock common dequeue to get task. Linux will put these lightweight process idle. If we do not consider libkomp+cw's algorithmic variant, then PLASMA algorithm with GCC is the best runtime in terms of energy consumption for LU factorization. This is not true for runtime systems based on task scheduling by work stealing such as libKOMP, libOMP or XKaapi which have very active threads that consume energy.

6. Discussion

Majority of best configurations from Figure 2 were runtime systems using work stealing based scheduling. On fine grain problems, libkomp and xkaapi were generally better. These results can be explained by the smaller task creation overhead on xkaapi and libkomp than others.

The difference between libomp and libkomp is the new features we add into the original Intel libomp runtime: the lightweight work stealing algorithm from Cilk and the request combining protocol from xkaapi. These features not only impact performance but also the way tasks are scheduled: it suppresses the bounded dequeue limitation that may degenerate task creation into task serialization. It means that at runtime a thread may be forced to execute immediately tasks for which no or less affinity exist. Without bounded size dequeue, a thread that completes a task will always activate one of the successors following a data flow relationship producer-consumer, thus sharing a data resident into cache; or the thread becomes

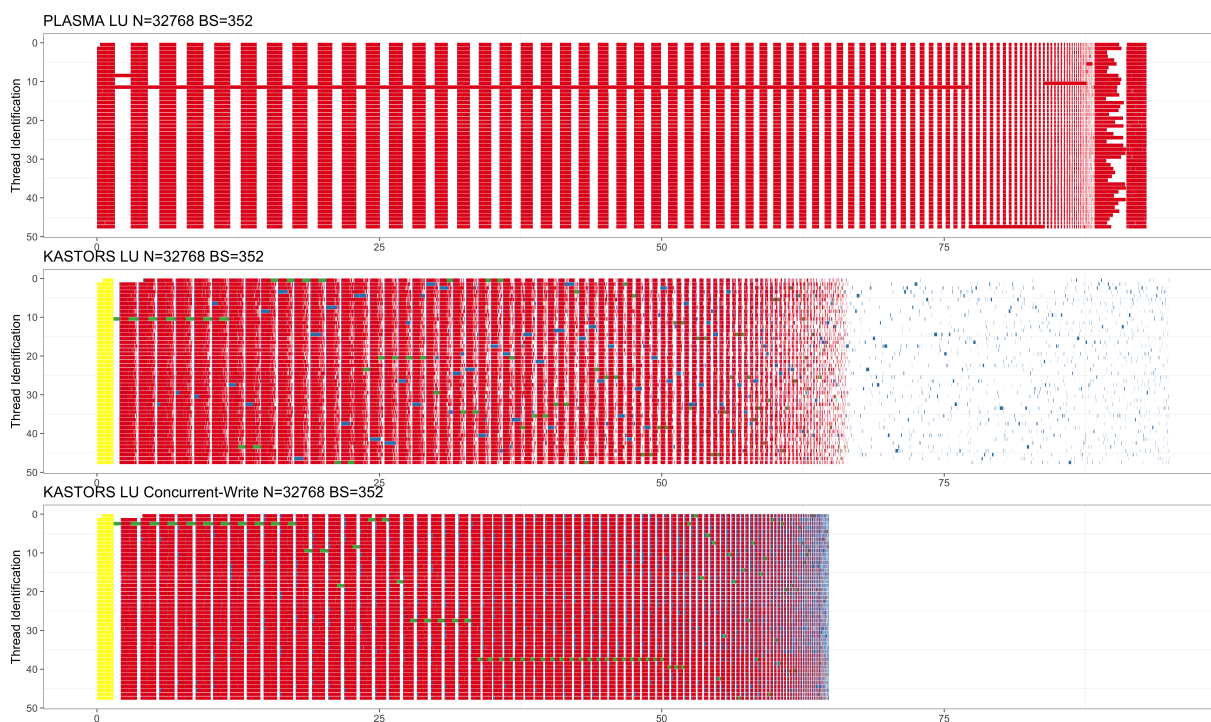


Figure 6. Gantt of LU algorithm from PLASMA (top), KASTORS and libKOMP (middle), and KASTORS with libKOMP and concurrent write (bottom). The matrix size was 32768x32768 with 352x352 of block size.

idle and try to steal tasks. We will investigate by more finer experiments the exact impact of these additions in libkomp.

On LU factorization where algorithmic variant libkomp+cw was the best, it was followed by xkaapi and libkomp on performance. LU factorization is a relevant code with inactivity sections from the dependencies imposed by the algorithm, mainly due to a search of pivot and swap of elements. This optimized algorithm allowed to increase performance while energy is decreased due to libkomp+cw runtime and concurrent write OpenMP extension (Virouleau et al. 2014). Nevertheless, the platform characteristic, and especially its memory network, had also an impact on both performance and energy consumption.

Without these algorithm variants, LU factorization code consumes less energy using GCC libGOMP runtime. In GCC the synchronization between threads on the shared task dequeue resource wastes less cycles. A work stealing based runtime may have interest to incorporate part of (Ribic and Liu 2014) in which is used to lower the speed of threads that are not in the critical path with a warranty on performance. One of the big challenges is the design of adaptive OpenMP runtime capable to saving energy on short delays of inactivity.

Our findings on CPU level optimizations lead us to believe that the three CPU level optimizations impacts performance and energy of OpenMP runtime systems. The two processors parameters, Turbo Boost and C-States, had significant influence over experimental results. The former increased performance significantly at energy cost, while the latter reduced energy at idle states. It was not clear the

impact of C-States over work stealing based runtime systems since our correlation test (Fig. 4) showed no relation between C-States and the experimental results. Nevertheless, we acknowledge that runtime systems based on work stealing may take advantage of steal phases at the beginning and the end of the computation to enter in idle state and reduce energy consumption (Tchiboukdjian et al. 2013).

CPUFreq governors had more impact at experimental cases with Turbo Boost disabled. Our experimental results suggest that turning all three CPU level optimizations to target performance, mostly DVFS governor and Turbo Boost, degraded performance due to CPU thermal limitations. Still, the powersave governor was able to sustain performance while consuming less energy. Other works show experimental results on techniques using DVFS, as discussed in Section 2.

7. Conclusion

In this paper, experiments with five production based OpenMP runtime systems and three CPU level optimizations (Turbo Boost, C-States, and Linux CPUFreq governors) on the three main kernel in dense linear algebra were conducted on a NUMA platform. We showed that OpenMP runtime is a new leverage for controlling energy, and Turbo Boost, as well as C-States, impacted significantly performance and energy. Our experimental results suggest that small algorithmic and runtime improvements may allow performance gains up to 63% and thus reducing the energy by 29%. Besides, work stealing based runtime systems were more efficient in our experiments; although, it was not clear the impact of C-States in order to reduce energy consumption at steal phases.

Future works include an extension of our experimental comparison over a wide range of architectures, including Intel KNL many-core. In addition, we will evaluate the impact on performance and energy of idle states at steal phases of work stealing scheduler.

Acknowledgements

This work is integrated and supported by the ELCI project, a French FSN (“Fond pour la Société Numérique”) project that associates academic and industrial partners to design and provide a software environment for very high performance computing.

References

- Bari MAS, Chaimov N, Malik AM, Huck KA, Chapman B, Malony AD and Sarood O (2016) Arcs: Adaptive runtime configuration selection for power-constrained openmp applications. In: *2016 IEEE International Conference on Cluster Computing (CLUSTER)*. pp. 461–470. DOI:10.1109/CLUSTER.2016.39.
- Benoit A, Lefevre L, Orgerie AC and Rais I (2017) Shutdown policies with power capping for large scale computing systems. In: *Europar 2017 : International European Conference on Parallel and Distributed Computing*. Santiago de Compostela, Spain. To appear.
- Bergman K, Borkar S, Campbell D, Carlson W, Dally W, Denneau M, Franzon P, Harrod W, Hill K, Hiller J et al. (2008) Exascale computing study: Technology challenges in achieving exascale systems. *Defense Advanced Research Projects Agency Information Processing Techniques Office (DARPA IPTO), Tech. Rep 15*.
- Broquedis F, Gautier T and Danjean V (2012) Libkomp, an efficient openmp runtime system for both fork-join and data flow paradigms. In: *Proceedings of the 8th International Conference on OpenMP in a Heterogeneous World, IWOMP'12*. Berlin, Heidelberg: Springer-Verlag. ISBN 978-3-642-30960-1, pp. 102–115.

- Bueno-Hedo J, Planas J, Duran A, Badia RM, Martorell X, Ayguadé E and Labarta J (2012) Productive Programming of GPU Clusters with OmpSs. IEEE Computer Society. ISBN 978-1-4673-0975-2, pp. 557–568.
- Curtis-Maury M, Dzierwa J, Antonopoulos CD and Nikolopoulos DS (2006) Online strategies for high-performance power-aware thread execution on emerging multiprocessors. In: *Proc. 20th IEEE International Parallel Distributed Processing Symposium*. DOI:10.1109/IPDPS.2006.1639598.
- Duran A, Teruel X, Ferrer R, Martorell X and Ayguade E (2009) Barcelona openmp tasks suite: A set of benchmarks targeting the exploitation of task parallelism in openmp. In: *Proc. of the 2009 International Conference on Parallel Processing, ICPP '09*. Washington, DC, USA: IEEE Computer Society. ISBN 978-0-7695-3802-0, pp. 124–131. DOI:10.1109/ICPP.2009.64.
- Etinski M, Corbalan J, Labarta J and Valero M (2010) Utilization driven power-aware parallel job scheduling. *Computer Science-Research and Development* 25(3-4): 207–216.
- Frigo M, Leiserson CE and Randall KH (1998) The implementation of the cilk-5 multithreaded language. *SIGPLAN Not.* 33(5): 212–223.
- Gautier T, Lima JVF, Maillard N and Raffin B (2013) Xkaapi: A runtime system for data-flow task programming on heterogeneous architectures. In: *Proceedings of the 2013 IEEE 27th International Symposium on Parallel and Distributed Processing, IPDPS '13*. Washington, DC, USA: IEEE Computer Society. ISBN 978-0-7695-4971-2, pp. 1299–1308.
- Gautier T and Virouleau P (2015) New libkomp library. <http://gitlab.inria.fr/openmp/libkomp>.
- Laros JH III, Pedretti KT, Kelly SM, Shu W and Vaughan CT (2012) Energy based performance tuning for large scale high performance computing systems. In: *Proc. of the 2012 Symposium on High Performance Computing, HPC '12*. San Diego, CA, USA: Society for Computer Simulation Int. ISBN 978-1-61839-788-1, pp. 6:1–6:10.
- Li D, de Supinski BR, Schulz M, Cameron K and Nikolopoulos DS (2010) Hybrid mpi/openmp power-aware computing. In: *2010 IEEE International Symposium on Parallel Distributed Processing (IPDPS)*. pp. 1–12. DOI:10.1109/IPDPS.2010.5470463.
- Lively C, Wu X, Taylor V, Moore S, Chang HC and Cameron K (2011) Energy and performance characteristics of different parallel implementations of scientific applications on multicore systems. *Int. J. High Perform. Comput. Appl.* 25(3): 342–350. DOI:10.1177/1094342011414749.
- Marathe A, Bailey PE, Lowenthal DK, Rountree B, Schulz M and de Supinski BR (2015) *High Performance Computing: 30th International Conference, ISC High Performance 2015, Frankfurt, Germany, July 12-16, 2015, Proceedings*, chapter A Run-Time System for Power-Constrained HPC Applications. Cham: Springer International Publishing. ISBN 978-3-319-20119-1, pp. 394–408. DOI:10.1007/978-3-319-20119-1_28.
- Nandamuri A, Malik AM, Qawasmeh A and Chapman BM (2014) Power and energy footprint of openmp programs using openmp runtime api. In: *Proceedings of the 2Nd International Workshop on Energy Efficient Supercomputing, E2SC '14*. Piscataway, NJ, USA: IEEE Press. ISBN 978-1-4799-7036-0, pp. 79–88. DOI: 10.1109/E2SC.2014.11.
- Orgerie AC, Assuncao MDd and Lefevre L (2014) A survey on techniques for improving the energy efficiency of large-scale distributed systems. *ACM Comput. Surv.* 46(4): 47:1–47:31. DOI:10.1145/2532637.
- Porterfield AK, Olivier SL, Bhalachandra S and Prins JF (2013a) Power measurement and concurrency throttling for energy reduction in openmp programs. In: *2013 IEEE International Symposium on Parallel Distributed Processing, Workshops and Phd Forum*. pp. 884–891. DOI:10.1109/IPDPSW.2013.15.
- Porterfield AK, Olivier SL, Bhalachandra S and Prins JF (2013b) Power measurement and concurrency throttling for energy reduction in openmp programs. In: *Parallel and Distributed Processing Symposium Workshops & PhD*

- Forum (IPDPSW)*, 2013 *IEEE 27th International*. IEEE, pp. 884–891.
- Ribic H and Liu YD (2014) Energy-efficient work-stealing language runtimes. *SIGARCH Comput. Archit. News* 42(1): 513–528. DOI:10.1145/2654822.2541971.
- Rotem E, Naveh A, Ananthakrishnan A, Weissmann E and Rajwan D (2012) Power-management architecture of the intel microarchitecture code-named sandy bridge. *IEEE Micro* 32(2): 20–27. DOI:10.1109/MM.2012.12.
- Rountree B, Lownenthal DK, de Supinski BR, Schulz M, Freeh VW and Bletsch T (2009) Adagio: Making dvs practical for complex hpc applications. In: *Proceedings of the 23rd International Conference on Supercomputing*, ICS '09. New York, NY, USA: ACM. ISBN 978-1-60558-498-0, pp. 460–469. DOI: 10.1145/1542275.1542340.
- Su C, Li D, Nikolopoulos DS, Cameron KW, d Supinski BR and León EA (2012) Model-based, memory-centric performance and power optimization on numa multiprocessors. In: *2012 IEEE International Symposium on Workload Characterization (IISWC)*. pp. 164–173. DOI:10.1109/IISWC.2012.6402921.
- Tchiboukdjian M, Gast N and Trystram D (2013) Decentralized list scheduling. *Annals of Operations Research* 207(1): 237–259.
- Treibig J, Hager G and Wellein G (2011) LIKWID: lightweight performance tools. *CoRR* abs/1104.4874.
- Virouleau P, Broquedis F, Gautier T and Rastello F (2016) Using data dependencies to improve task-based scheduling strategies on numa architectures. In: *Proceedings of the 22Nd International Conference on Euro-Par 2016: Parallel Processing - Volume 9833*. New York, NY, USA: Springer-Verlag New York, Inc. ISBN 978-3-319-43658-6, pp. 531–544. DOI:10.1007/978-3-319-43659-3.39.
- Virouleau P, Brunet P, Broquedis F, Furmento N, Thibault S, Aumage O and Gautier T (2014) Evaluation of OpenMP Dependent Tasks with the KASTORS Benchmark Suite. In: *10th International Workshop on OpenMP, IWOMP'14*. Springer, pp. 16 – 29. DOI:10.1007/978-3-319-11454-5_2.
- Yang X, Zhou Z, Wallace S, Lan Z, Tang W, Coghlan S and Papka ME (2013) Integrating dynamic pricing of electricity into energy aware scheduling for hpc systems. In: *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis, SC '13*. New York, NY, USA: ACM. ISBN 978-1-4503-2378-9, pp. 60:1–60:11. DOI:10.1145/2503210.2503264.
- YarKhan A, Kurzak J, Luszczek P and Dongarra J (2016) Porting the plasma numerical library to the openmp standard. *International Journal of Parallel Programming* : 1–22DOI:10.1007/s10766-016-0441-6.

Author Biographies

João Vicente Ferreira Lima received a joint PhD degree in computer science by the Federal University of Rio Grande do Sul (UFRGS), Brazil, and the MSTII Doctoral School at the Grenoble University, France. He received a BSc degree in Computer Science in 2007 by the Federal University of Santa Maria (UFSM), Brazil, and a MSc degree in Computer Science in 2009 by he Federal University of Rio Grande do Sul (UFRGS), Brazil. He is associate professor at the Federal University of Santa Maria (UFSM), Brazil, since 2014. His research interests are high performance computing, runtime systems for HPC, parallel programming for accelerators, and distributed computing.

Issam Raiïs is a PhD student since 2015 at ENS Lyon, France. He received a MSc degree in Computer Science and a BSc degree in Computer and Information Sciences at the University of Orleans, France. He works in the AVALON team at the LIP laboratory advised by Laurent Lefèvre, Anne Benoit and Anne-Cécile Orgerie.

Thierry Gautier received the Dipl. Ing. MS and PhD in computer science at the INPG, in 1996. He is a full time researcher at INRIA (the French National Institute for Computer Science and Control), with the AVALON

project team of the LIP laboratory in Lyon, France, and has held a post-doctoral position at ETH Zürich (1997). Dr. Gautier conducts research in high performance computing, runtime systems for HPC, parallel algorithms, parallel programming, OpenMP, and multicore architectures.

Laurent Lefèvre obtained his PhD in Computer Science in January 1997 at LIP Laboratory (Laboratoire Informatique du Parallelisme) in ENS-Lyon (Ecole Normale Supérieure), France. From 1997 to 2001, he was Assistant Professor in computer science in Lyon 1 University and a member of the RESAM Laboratory (High Performance Networks and Multimedia Application Support Lab). Since 2001, he has been a Research Associate in computer science at INRIA (the French Institute for Research in Computer Science and Control). He is a member of the INRIA AVALON team (Algorithms and Software Architectures for Distributed and HPC systems) at the LIP laboratory in Lyon, France. His research interests focus on green and energy efficient computing and networking. He has organized several conferences in high-performance networking and computing and he is a member of several program committees. He has co-authored more than 100 papers published in refereed journals and conference proceedings. He participates in several national and European projects on energy efficiency. For more information, see <http://perso.ens-lyon.fr/laurent.lefevre/>.