



HAL
open science

ScalOMP: analyzing the Scalability of OpenMP applications

Anton Daumen, Patrick Carribault, François Trahay, Gaël Thomas

► **To cite this version:**

Anton Daumen, Patrick Carribault, François Trahay, Gaël Thomas. ScalOMP: analyzing the Scalability of OpenMP applications. IWOMP 2019: 15th International Workshop on OpenMP, Sep 2019, Auckland, New Zealand. pp.36-49, 10.1007/978-3-030-28596-8_3 . hal-02179726

HAL Id: hal-02179726

<https://hal.science/hal-02179726v1>

Submitted on 11 Jul 2019

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.

SCALOMP: Analyzing the Scalability of OpenMP Applications

Anton Daumen^{1,2}, Patrick Carribault¹, François Trahay², and Gaël Thomas²

¹ CEA, DAM, DIF,
F-91297, Arpajon, France

{anton.daumen.ocre,patrick.carribault}@cea.fr

² Télécom SudParis

Institut Polytechnique de Paris

{francois.trahay,gael.thomas}@telecom-sudparis.eu

Abstract. Achieving good scalability from parallel codes is becoming increasingly difficult due to the hardware becoming more and more complex. Performance tools help developers but their use is sometimes complicated and very iterative. In this paper we propose a simple methodology for assessing the scalability and for detecting performance problems in an OpenMP application. This methodology is implemented in a performance analysis tool named SCALOMP that relies on the capabilities of OMPT for analyzing OpenMP applications. SCALOMP reports the code regions with scalability issues and suggests optimization strategies for those issues. The evaluation shows that SCALOMP incurs low overhead and that its suggestions lead to significant performance improvement of several OpenMP applications.

Keywords: Performance Tool · Scalability · OMPT

1 Introduction

The lifespan of simulation codes is several times longer than the lifespan of supercomputers. Thus, a single code will be used on multiple very different architectures, making the portability and the optimization of codes difficult. Furthermore computer architectures are more and more complex as their design has to become more intricate in order for performance to continue increasing. In their chase for better performance, developers rely on performance analysis tools to understand and analyze their code.

Many performance tools provide a wide range of features, metrics, and analysis. However the more features a performance tool has, the more complex it is to use. The developer has to learn how to use the tools in order to start efficiently using them for code analysis. Furthermore a lot of tools use an incremental methodology for analyzing codes, the tool first reports the global behavior of the code and the developer then focuses his analysis on important regions. The developer then tries to detect the issues in said regions by using other features of the tool and analyzing the source code directly, forming hypothesis and using

the tool to verify them and quantify the importance of a performance problem before trying to fix it in the code.

In this paper we propose a simple methodology for analyzing the performance of a parallel application with a focus on the scalability of OpenMP applications. This approach is implemented in SCALOMP, that reports directly to the user the parallel regions where time is lost due to scalability issues and to automatically deduce the sources of these losses in order for the developer to directly know where time was lost and why. The whole process needs to be as closely related to the source code as possible in order for the developer to immediately understand where an issue resides. When possible, SCALOMP provide hints on how an issue may be solved.

The remainder of this paper is organized as follows: in Section 2 we present state-of-the-art tools that illustrate the typical methodology of performance tools. We detail our methodology in Section 3, and we describe SCALOMP internals in Section 4. We evaluate our approach in Section 5, and in Section 6 we conclude this paper.

2 Related Work

The performance tool landscape is filled with a significant number of tools that provide a very broad variety of features.

A lot of effort has been spent on tools that help the developer better visualize the behavior of an application [15, 8]. These tools allow developer to precisely examine the application execution, but the analysis has to be done manually.

Automatic analysis relieve the developer from the analysis. Several works have focused on detecting the root causes of scalability issues in parallel applications. Most of these works are focused on MPI; For example, performance models can help finding weak scaling issues [6]. A backward replay of an execution trace can be used for identifying the root cause of wait-states in MPI applications [5].

Other work focuses on detecting and reporting problems in multi-threaded applications. For instance, imbalance issues in OpenMP parallel regions and worksharing constructs can affect the scalability of an application [18]. Running micro-benchmarks and building a compositional model can predict the performance of OpenMP applications on a given machine [16]. However, this approach is limited to memory-bound applications and only works on OpenMP applications using the *static* scheduler. In order to detect false-sharing, a recent work uses the OMPT API to instrument OpenMP constructs, and collects hardware counters at a fine granularity [9]. The collected data then train a classifier which can then spot false sharing in applications. Automated performance modeling can be used to examine the scalability of OpenMP runtime constructs [12], or to analyze the memory access patterns of OpenMP applications [4].

While all of these approaches are functional and allow a developer to identify issues, verifying every possible problem using different tools is time consuming. Moreover each approach has its own requirements and limitations which can

make using these tools together difficult, and does not match our view on how the performance analysis of an application should work.

Some tools do integrate multiple analysis successfully. Intel VTune [17] provides an OpenMP *time gain* analysis that estimates the potential gain that could be obtained if various performance problems (lock contention, imbalance, scheduling overhead, etc.) were fixed. However Vtune *time gain* is lacking a scalability analysis, which means that if a performance issue is not detected, a code region may be wrongfully considered as having no issue. Finally, VTune uses profiling to measure the time spent in OpenMP constructs. While this limits the instrumentation overhead, it also affects the measurements precisions and lacks some insight that tracing may give.

3 Performance analysis of OpenMP applications

As described in Section 2, even with modern performance tools, most of the analysis remains the work of the developer and is done manually. Our work focuses on alleviating this burden as much as possible from the developer hands. This section presents our approach for assessing the scalability of a multithreaded application and how a performance analysis tool can provide developers with optimisation hints. We implement this approach in SCALOMP, whose implementation details are described in Section 4.

3.1 Methodology

There are multiple sources of performance problems in OpenMP applications, such as load imbalance, or lock contention. Once the problem is identified, the developer may improve the performance of the application in several ways. Some issues require code changes, while changing the execution settings may fix some other problems. In this section, we describe a methodology for detecting performance problems in OpenMP applications, and providing optimization suggestions to the developer.

As an input, the developer provides a compiled version of the application, along with the command line that runs it. The application is instrumented and runs while varying the number of threads in order to measure the scalability of each parallel region and to detect performance problems. As a result, SCALOMP computes the parallel efficiency of each parallel region (defined by the speedup multiplied by the initial number of threads divided by the current number of threads), and estimates the potential time gain for each parallel region. The output is the list of parallel regions sorted according to the potential time gain, and for each region, a set of optimization hints and their respective potential time gain.

Using this approach, a developer can focus on optimizing the most promising parallel regions of his application. Moreover, the optimization hints indicate how the performance could be improved.

3.2 Scalability analysis on parallel regions

The performance analysis of an application starts with a scalability analysis which aims at identifying the OpenMP parallel regions that may be worth optimizing. SCALOMP does this by running the application multiple times while automatically varying the number of threads across a range given as input. For each run, SCALOMP measures the duration of each parallel region. As a result, SCALOMP computes the parallel efficiency of all the regions and estimates the time lost in each by comparing their efficiency to the expected behaviour of a perfectly scalable region [7].

Thus, SCALOMP identifies the parallel regions where most time is lost and never underestimate a parallel region's impact due to not detecting its issues or wrongly quantifying their effects. These parallel regions are good candidates for optimizations: the poor scalability of a region means that some performance issues affects its parallel efficiency ; and the significant amount of time lost means there is good hope of gaining back time by improving the region's performance.

3.3 Quantifying the impact of performance problems

A poor parallel efficiency in an OpenMP parallel region may be caused by several types of problems. In this section, we describe some of these problems and how a performance analysis tool can estimate their impact on performance. As a result, it is possible to quantify the potential time gain for each problem in each parallel region.

Barriers One of the main synchronization mechanisms in OpenMP is the barrier that allows threads to wait for each other. This synchronization may be explicit (when using the `omp barrier` directive), or implicit (e.g., at the end of a parallel for loop). While barriers allow developers to ensure the correctness of parallel programs, they introduce synchronization points that may degrade the parallel efficiency.

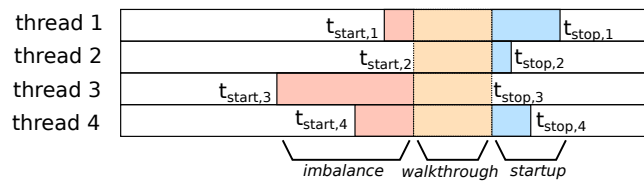


Fig. 1: Illustration of threads passing through an OpenMP barrier

As illustrated in Figure 1, when a set of threads pass through an OpenMP barrier, three phases can be distinguished:

- the *imbalance phase* starts when the first thread enters the barrier, and ends when the last thread reaches the barrier. If a thread arrives late to a barrier, it delays the other threads. This means that the ideal case is when all the threads reach the barrier simultaneously. Hence the time lost by imbalance is the difference between the average time a thread took to reach the barrier from the last point of synchronization, and the maximum time. A long *imbalance phase* may be caused by an uneven work distribution between the threads, or by some delay that applies to a thread (such as a late MPI communication). If the imbalance is significant, SCALOMP suggests to improve the work distribution, for example by using a `schedule dynamic` clause in a parallel loop.
- the *walkthrough phase* starts when all the threads have reached the barrier and ends when the first thread leaves the barrier. We consider this phase to be the incompressible time spent resolving the barrier. This is an optimistic estimation since part of the barrier is resolved every time a thread arrives. If the walkthrough phase takes a significant part of the overall execution time, SCALOMP suggests to either reduce the number of OpenMP barriers, or to improve the barrier algorithm (for instance, by choosing a more performant OpenMP runtime)
- the *startup phase* starts when the first thread leaves the barrier and ends when the last thread leaves the barrier. The time lost from threads not leaving at the same time is paid at the next point of synchronization. For example if there is no imbalance in the work of threads between two barriers, some tools could detect differences in arrival times and report it as imbalance when the real culprit is the previous barrier delay when releasing threads. SCALOMP detects those delays and shifts the blame to the previous barrier runtime instead.

Locks Locking is another major synchronization mechanism in OpenMP that may significantly impact the performance. As depicted in Figure 2, the time spent acquiring a lock can be separated in two phases:

- the *waiting* phase happens when a thread tries to acquire a lock that is currently held by another thread. This phase corresponds to the contention that applies to the lock. To reduce the waiting phase, SCALOMP suggests to either change the application to reduce the number of concurrent access to this lock, or to use another locking mechanism that is less affected by contention (for instance MCS or AHMCS [10]).
- the *acquisition* phase happens when the lock is available. This phase corresponds to the incompressible time required for running the locking algorithm. A significant part of the whole execution spent in the acquisition phase means that the thread often acquires locks without contention. In that case, SCALOMP suggests to either reduce the number of calls to locking primitives, or to use another synchronization mechanism (such as the `atomic` directive or a locking mechanism that works better with no contention).

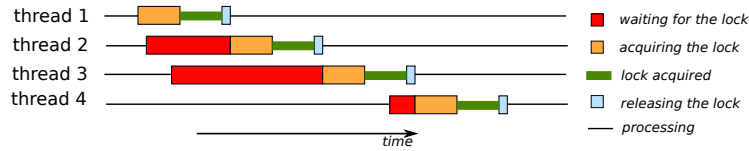


Fig. 2: Illustration of threads acquiring a lock

In order to estimate how much time is due to the contention, and to the lock algorithm, SCALOMP measures the time spent acquiring the lock and assumes that the fastest measured acquisition was contention free. This gives an estimate of the constant time required for executing the locking algorithm, and the remaining time is attributed to the contention on the lock.

4 Implementation

In order to compute the metrics described in Section 3, we implemented SCALOMP. In this section, we detail how SCALOMP instruments an OpenMP application, and how it measures the duration of OpenMP constructs without altering the application behavior.

4.1 Instrumenting an OpenMP application with OMPT

As OpenMP relies on both compiler directives and a runtime API, instrumentation can be tedious. One solution consists in building a set of wrappers that intercept the calls to the OpenMP runtime API. However, this method is specific to one OpenMP implementation and it cannot grasp the whole OpenMP semantics. Opari [14] performs a source-to-source transformation on the application and inserts POMP calls in the source code. This makes this approach more portable, but it requires to recompile the application.

SCALOMP uses the OpenMP Tools interface (OMPT) that was introduced in the OpenMP 5.0 standard [11]. OMPT makes the OpenMP runtime collaborate with performance analysis tools: the tools register callbacks for OpenMP events, and the OpenMP runtime then triggers the callbacks when the corresponding events happen. With this approach, SCALOMP can collect performance data from any OpenMP application without recompiling it.

4.2 Identifying OpenMP parallel regions

An application may consist of tens of OpenMP parallel regions, some of them being invoked multiple times. Thus, SCALOMP needs to identify a parallel region in order to aggregate the performance data from multiple calls to it. When the application starts an OpenMP parallel region, the OMPT interface invokes SCALOMP through a callback and provides a pointer to the OpenMP call in the application binary. The first time SCALOMP encounter an unknown callsite, it uses *libbfd* to retrieve the line of code associated with this address.

4.3 Measuring temporal data

As described in Section 3, SCALOMP analysis of the OpenMP barriers requires to collect several OpenMP events: SCALOMP needs to know when a thread starts a region’s work, enters at a barrier, exits a barrier and ends the region. The lock analysis also requires to collect information when a thread starts and stops acquiring a lock. For each of these events, SCALOMP records a timestamp using the TSC counter. These timestamps are then used for measuring various durations in the thread processing. The TSC counter allows SCALOMP to record timestamps at a low cost, but these timestamps cannot be compared across threads running on different sockets. Thus, SCALOMP also records a system-wide timestamp using `clock_gettime` at the beginning of each region execution in order to compare different threads timestamps.

4.4 Mitigating instrumentation overhead with adaptive sampling

While recording a timestamp using the TSC counter is lightweight, this overhead may significantly alter the application’s performance if timestamps are recorded too often. In order to reduce this overhead, SCALOMP uses a sampling mechanism. Since the OMPT interface allows to dynamically activate or de-activate callbacks, SCALOMP only collects performance data on certain executions of a parallel region. The idea is that while two executions of a region may not be exactly the same, their behaviour is essentially similar.

Depending on the parallel region, the sampling frequency should be selected carefully: if a parallel region is only repeated a few times, all its executions should be captured, whereas a region that runs many times should only be captured from time to time. Thus, SCALOMP uses an adaptive sampling where the first executions are all measured, and then as the region is repeated, SCALOMP de-activates the callbacks for some executions. The more a region is repeated, the more often SCALOMP de-activates the OMPT callbacks.

As a result, the rare regions are all captured, while frequent regions are sparsely captured, and the overhead of SCALOMP on the application execution remains low.

5 Experiments and Results

In this Section, we evaluate SCALOMP implementation and assess how the performance analysis can help the developer improve a parallel application. First, we evaluate the overhead of SCALOMP on 16 applications. Then, we evaluate how SCALOMP detects load imbalance problems, and lock contention issues.

For our evaluation, we use a machine equipped with two Intel Xeon Haswell E5-2698 v3 processors with 16 cores each (32 cores in total), and 128 GB of RAM. The machine runs Linux version 3.10, and the applications were compiled with Intel Compiler version 17.0.6. For OpenMP we use the open-source OpenMP runtime from Intel now maintained in LLVM. When compiling applications, we use the `-O3` optimization level.

We evaluate SCALOMP using several OpenMP applications:

- **Mandelbrot** is an application that computes the Mandelbrot set;
- **HydroMM** is an hydrodynamics mini-application;
- **Lulesh 2.0** is a mini-application that performs an hydrodynamics simulation [13];
- **BT, CG, DC, EP, FT, IS, LU, MG, UA** are kernels from the OpenMP NAS Parallel Benchmarks version 3.3.1 [3];
- **miniFE** is a Finite Element mini-application [2];
- **Snap** is a particle transport mini-application [2];
- **AMG** is a parallel algebraic multigrid solver for linear systems [1];
- **Pennant** is a mini-application for hydrodynamics [1].

5.1 Overhead of SCALOMP

To evaluate the overhead of SCALOMP, we run the 16 applications described in Section 5 with and without SCALOMP. For each application, the problem size is chosen so that the reference time (i.e. the execution time when running without SCALOMP) is between 10 and 100 second with a few exception to see how scale affects the overhead. Each measurement is repeated 5 times and we report the average execution time. Table 1 reports the execution time when running the application without SCALOMP, and the overhead when running with SCALOMP.

Table 1: Overhead induced by the tool

Application	Mandelbrot	HydroMM	Lulesh2.0	BT.B	CG.C	DC.A	EP.D	FT.C	IS.D
Reference time	11.34s	13.63s	82.02s	10.35s	12.23s	16.21s	52.67s	11.67s	27.09s
Overhead	2.24%	6.18%	12.34%	0.00%	0.03%	-0.27%	3.55%	-0.55%	-0.01%
Application	LU.C	LU.D	MG.D	UA.B	UA.D	miniFE	Snap	AMG	Pennant
Reference time	36.53s	1352.66s	88.06s	12.61s	1161.65s	18.49s	80.24s	94s	41.32s
Overhead	-0.17%	0.14%	0.12%	13.78%	4.04%	2.24%	-0.64%	-1.03%	2.33%

The results show that SCALOMP has little impact on the performance of most applications. The overhead is higher for Lulesh (12.34 %) because this application performs many small parallel regions; The observed overhead goes down to -4.31 % if the size is increased to 120 (from 80) and the number of iteration lowered to 350 (from 1000) so that the time stay similar. UA.B also suffers from a significant overhead (13.78 %) due to the heavy number of lock operations (8.1M locks per second per thread on average). When running UA with a larger problem size (class D), the application take locks less often and the SCALOMP’s overhead is reduced to 4.04 %. We conclude that SCALOMP does not significantly alter the application execution except in some extreme cases.

5.2 Detecting imbalance issues

SCALOMP reports that several of the applications evaluated in Section 5.1 suffer from load imbalance. In this section, we focus on two of these applications.

Mandelbrot. When running the **mandelbrot** application with 32 threads, SCALOMP reports that the parallel efficiency is only 42 %. SCALOMP reports that the load imbalance between the threads in one parallel region is responsible for all of the lost time. It suggests to improve the load balancing for this parallel region, as it predicts a perfect load balance may save 5.37 seconds.

Based on this suggestion, we analyze the source code of this application. The incriminated parallel region computes the divergence of a set of complex numbers in the Mandelbrot set. The computation cost for each complex number depends on how fast it diverges. The default OpenMP scheduling policy assigns many numbers that diverge quickly to some threads while some other threads have to process many numbers that diverge slowly. As a result, some threads finish their loop iterations earlier than the other threads, leading to a load imbalance and a poor parallel efficiency. As suggested by SCALOMP, we change the scheduling policy for this parallel region to **dynamic**, and we observe that the application’s execution time is reduced to 6.16 seconds. This means we gained 5.20 seconds which is close to the 5.37 predicted. When analyzing the parallel region with *dynamic* SCALOMP find the load balance to be 99.9 % perfect.

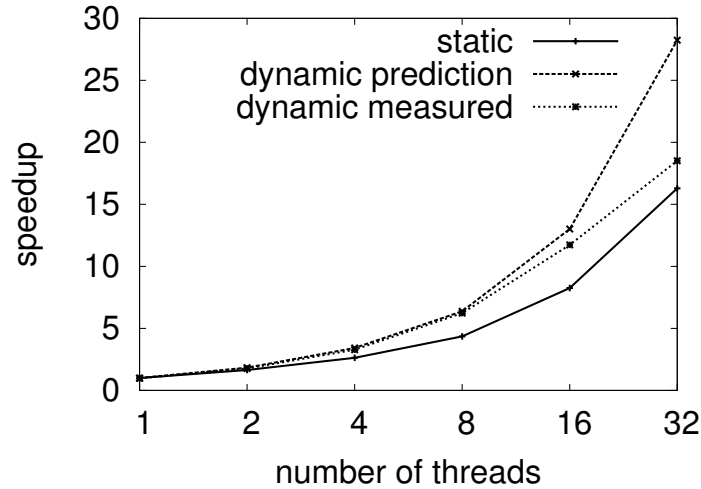
HydroMM. When running HydroMM with 32 threads, SCALOMP reports that the parallel efficiency is only 51 %, and reports that the load imbalance in one parallel region is responsible for 87 % of the total lost time. SCALOMP suggests to improve the load balancing of this parallel region, it also predicts the performance of the application if this parallel region was perfectly load balanced.

Based on SCALOMP suggestion, we analyze the source code of HydroMM, and change the OpenMP scheduling policy to **dynamic** in order to improve the load balancing between the threads. Figure 3 reports the speedup measured for HydroMM when using the default scheduling policy (**static**) and when using the **dynamic** scheduling policy. It also reports the speedup predicted by SCALOMP. We observe that changing the scheduling policy significantly improves the application’s performance. The results also show that up to 16 threads, the speedup obtained when applying SCALOMP suggestion is close to the predicted speedup.

For 32 threads, the predicted speedup is significantly overestimated. This may be due to memory effects being ignored by SCALOMP: up to 16 threads, all the threads execute on one socket of the machine, while when running 32 threads, the two sockets are used.

5.3 Detecting locking issues

In this section, we assess how SCALOMP detects locking issues using two applications. First, we evaluate how SCALOMP differentiates contended locks and non-contended locks using a micro-benchmark. Then, we present a case study on the UA kernel from the NAS Parallel Benchmarks. For both applications, we compare the optimization suggestions provided by SCALOMP with those obtained with Intel VTune [17].

Fig. 3: Speedup obtained when running **HydroMM**

Micro-benchmark. We implemented an OpenMP application that consists in two parallel regions. In these parallel regions, each thread acquires an OpenMP lock, releases it, and busy waits for some time. In the first parallel region, each thread accesses a different lock which does not suffer from contention. In the second parallel region, all the threads access the same lock which suffers from contention. We choose the busy wait time so that the time spent acquiring lock is similar for both parallel region.

Table 2: Lock Micro-benchmark

	without contention	with contention
Total duration(s)	14.14	14.34
Lock duration(s)	4.63 (32 %)	4.35 (29.6 %)
VTune	Overhead	2.61 %
	Lock contention	10.19 %
	Other problems	2.64 s (18.19 %)
SCALOMP	Overhead	0.67 s (4.6 %)
	Lock duration	4.00 %
	- lock algorithm	3.93 s (26.97 %)
	- lock contention	5.24 s (36.45 %)
	3.40 s (23.14 %)	0.03 s (0.23 %)
	0.56 s (3.83 %)	5.21 s (36.22 %)

We run this micro-benchmark and analyze it with VTune and SCALOMP. The results of this experiment are reported in Table 2. The time spent acquiring locks and the total duration of the two regions are similar.

VTune detects significant lock contention in the contention-free region (18.19 % of the region duration) and in the region with contention (34.81 % of the region duration). VTune also reports that 4.6 % of the time spent in region 1 is lost due to “Other” problems. We conclude that VTune is able to detect the lock contention problem in the second parallel region, but it wrongfully detects a lock contention in the first parallel region.

SCALOMP detects that a significant time is spent in locks in the contention-free region, and that most of it is due to the lock algorithm itself. SCALOMP indicates that the problem is that the threads acquire too many locks. In the second parallel region, SCALOMP detects that most of the locking time is due to contention. We conclude that SCALOMP rightly identifies the lock problems in the two parallel regions.

Case Study: UA. In this section, we analyze UA and apply optimizations based on the suggestions provided by SCALOMP. Since one of these optimizations gives an incorrect result with the Intel Compilers, we use the GNU Compilers version 7.3.0 in this section.

When running UA.B, SCALOMP measures a parallel efficiency of 56 % with 32 threads. One parallel region is responsible for most of the time loss because of several problems: 10% of the total execution time is lost due to load imbalance; and 19.4 % of the total execution time is spent acquiring locks in this parallel region.

SCALOMP also points that 2 704 354 500 locks are acquired during the 8.68 s execution of this parallel region, meaning that on average, each thread acquires a lock every 102 ns. Due to a high number of region execution (more than 2000), SCALOMP automatically uses the sampling mechanism described in Section 4.4 and records the duration of only 16 % of the all locks acquisition. As a result, the overhead induced by SCALOMP remains low, but SCALOMP still captures a significant amount of performance data.

SCALOMP reports that most of the time spent acquiring locks is due to the lockings algorithm, and contention on locks is low in this application. An analysis of UA source code shows that there are 334 600 different locks, which limits the probability of a thread acquiring a lock that is already taken by another thread. Moreover, most of the locks are used for protecting simple instructions such as $x = x + y$.

Based on this analysis, we create two additional implementations of the application:

- **UA-Hint** uses the lock hint mechanism from OpenMP 5.0 and specifies that the locks are uncontended. This allows the OpenMP runtime to select the lock implementation that performs the best when there is no contention;
- **UA-Atomic** replaces the critical sections protected by locks with OpenMP atomic operations when possible.

We compare the performance of these two implementation with the UA-Default implementation. Figure 4 reports the speedup of the implementations

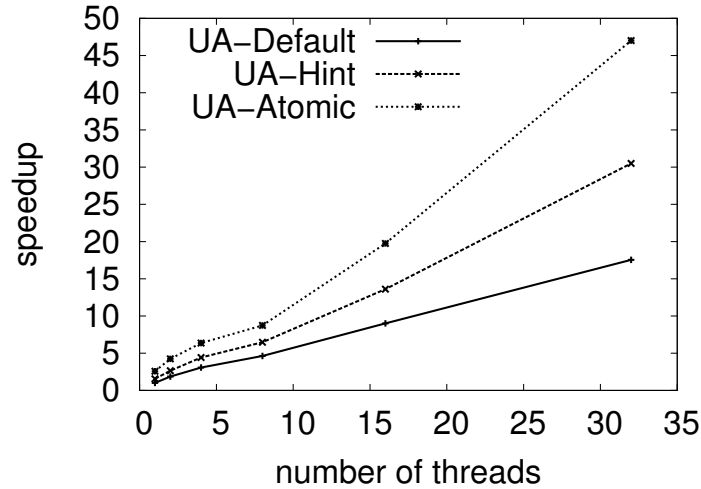


Fig. 4: Speedup obtained when running UA

as compared to the execution time obtained when running UA-Default with one thread. As suggested by SCALOMP, UA-Hint performs better than UA-Default in all cases (including in sequential) because the locking algorithm achieves better performance when the lock is uncontended. UA-Atomic outperforms UA-Default (by 267 % on 32 threads) and UA-Hint (by 154 % on 32 threads) for all the tested number of threads. This is due to the single atomic instruction that replaces a call to a locking function and the critical section.

We conclude that SCALOMP suggests optimizations that may significantly improve the performance of an OpenMP application that suffers from locking problems.

6 Conclusion

Performance analysis of an application involves a lot of work from the developer. In this paper we presented a methodology that focuses on the scalability of an application in order to help the developer improve its performance. We implemented this approach in SCALOMP that relies on the OMPT API to instrument OpenMP applications. The evaluation show that SCALOMP collect performance data from applications with a low overhead. The experiments also show that SCALOMP analysis successfully detect load imbalance problems, and locking problems in several applications. The optimization hints provided by SCALOMP help the developer significantly improve the application’s performance.

References

1. Coral-2 benchmarks. Tech. rep., Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States), <https://asc.llnl.gov/coral-2-benchmarks/index.php>
2. Coral benchmarks. Tech. rep., Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States), <https://asc.llnl.gov/CORAL-benchmarks/>
3. NAS parallel benchmarks applications (NPB). Tech. rep., NASA Advanced Supercomputing Division, <https://www.nas.nasa.gov/publications/npb.html>
4. Barthou, D., Rubial, A.C., Jalby, W., Koliai, S., Valensi, C.: Performance tuning of x86 openmp codes with maqao. In: Tools for High Performance Computing 2009, pp. 95–113. Springer (2010)
5. Bohme, D., Geimer, M., Wolf, F., Arnold, L.: Identifying the root causes of wait states in large-scale parallel applications. In: 2010 39th International Conference on Parallel Processing. pp. 90–100 (2010)
6. Calotoiu, A., Hoefler, T., Poke, M., Wolf, F.: Using automated performance modeling to find scalability bugs in complex codes. In: Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis. p. 45 (2013)
7. Coarfa, C., Mellor-Crummey, J.M., Froyd, N., Dotsenko, Y.: Scalability analysis of SPMD codes using expectations. In: Proceedings of the 21th Annual International Conference on Supercomputing, ICS 2007, Seattle, Washington, USA, June 17-21, 2007. pp. 13–22 (2007)
8. Coulomb, K., Degomme, A., Faverge, M., Trahay, F.: An open-source tool-chain for performance analysis. In: Tools for High Performance Computing 2011, pp. 37–48 (2012)
9. Ghane, M., Malik, A.M., Chapman, B., Qawasmeh, A.: False sharing detection in openmp applications using ompt api. In: International Workshop on OpenMP. pp. 102–114 (2015)
10. Guerraoui, R., Guiroux, H., Lachaize, R., Quéma, V., Trigonakis, V.: Lock–unlock: Is that all? a pragmatic analysis of locking in software systems. *ACM Transactions on Computer Systems (TOCS)* **36**(1), 1 (2019)
11. Huck, K.A., Malony, A.D., Shende, S., Jacobsen, D.W.: Integrated measurement for cross-platform openmp performance analysis. In: Using and Improving OpenMP for Devices, Tasks, and More - 10th International Workshop on OpenMP, IWOMP 2014, Salvador, Brazil, September 28-30, 2014. Proceedings. pp. 146–160 (2014)
12. Iwainsky, C., Shudler, S., Calotoiu, A., Strube, A., Knobloch, M., Bischof, C., Wolf, F.: How many threads will be too many? on the scalability of openmp implementations. In: European Conference on Parallel Processing. pp. 451–463. Springer (2015)
13. Karlin, I., Keasler, J., Neely, J.: Lulesh 2.0 updates and changes. Tech. rep., Lawrence Livermore National Lab.(LLNL), Livermore, CA (United States) (2013)
14. Knüpfer, A., Rössel, C., an Mey, D., Biersdorff, S., Diethelm, K., Eschweiler, D., Geimer, M., Gerndt, M., Lorenz, D., Malony, A., et al.: Score-p: A joint performance measurement run-time infrastructure for periscope, scalasca, tau, and vampir. In: Tools for High Performance Computing 2011, pp. 79–91. Springer (2012)
15. Müller, M.S., Knüpfer, A., Jurenz, M., Lieber, M., Brunst, H., Mix, H., Nagel, W.E.: Developing scalable applications with vampir, vampirserver and vampirtrace. In: Parallel Computing (PARCO). vol. 15, pp. 637–644 (2007)
16. Putigny, B., Goglin, B., Barthou, D.: A benchmark-based performance model for memory-bound hpc applications. In: 2014 International Conference on High Performance Computing & Simulation (HPCS). pp. 943–950 (2014)

17. Reinders, J.: Vtune performance analyzer essentials (2005)
18. Woodyard, M.: An experimental model to analyze openmp applications for system utilization. In: International Workshop on OpenMP. pp. 22–36 (2011)