

Internship Report

**Improving the Performance of Java Applications on
Multicore Architectures**

By :

Jean-Pierre Lozi
jean-pierre@lozi.org

Under the direction of :

Gaël Thomas
gael.thomas@lip6.fr

Gilles Muller
gilles.muller@lip6.fr

Julia Lawall
julia@diku.dk

Contents

1. Introduction	2
2. Related work	2
2.1. Cache locality	2
2.2. Contention management	3
2.2.1. Limiting the sharing of global data structures	3
2.2.2. Optimizing locks	4
2.3. Efficient inter-process communication	4
2.3.1. Shared memory	5
2.3.2. Message-based communication	5
2.4. Fast transfer of control	7
2.5. Heterogeneous architectures	9
2.6. Conclusion	10
3. Experiments	10
3.1. Rough estimates	11
3.1.1. Cost of cache misses	11
3.1.2. Cost of transfers of control	12
3.1.3. Cost comparison	13
3.2. NULL RPCs with multiple clients	13
3.3. Final benchmark	17
3.3.1. Executing critical sections locally	17
3.3.2. Context and shared variables	18
3.4. Future experiments	19
4. Conclusion	19
Bibliography	20
Appendix: source code	21

1. Introduction

Multicore architectures are becoming increasingly common in domains ranging from personal computing to server applications and even embedded systems. Exploiting the power of additional cores is a challenging issue, for several reasons. First, writing parallel programs is harder than writing sequential programs, and even state-of-the-art parallel programs comprise sequential sections that constitute a major limiting factor to overall performance. Second, synchronization between threads can cause major overhead. Using fine-grained synchronization primitives can limit this overhead, albeit at the cost of a greater complexity, lower maintainability and lower extendability of the source code. Finally, low-level cache coherence algorithms are costly, especially when the locality of the shared data between threads is low.

Exploiting the computing power offered by a large number of processing cores is especially hard for programs that are run within Java Virtual Machines (JVMs). Indeed, in Java, the only concurrency management mechanism is the synchronized block, which does not allow for low-level optimizations (e.g. improving data locality between threads or limiting contention via fine-grained locking). Such optimizations have to be implemented within the virtual machine itself.

In the long term, we plan to experiment with a novel approach to improve the performance of legacy Java applications on multicore hardware. The general idea is to improve data locality between cores in order to limit data transfers caused by cache coherence algorithms. Since Java synchronized blocks are used to ensure coherency when accessing shared variables, they can be used as delimiters of code areas that use shared data. Our idea, based on this observation, is to dedicate a number of processing cores to the execution of synchronized blocks in order to improve performance: indeed, these dedicated cores will be likely to have most of the shared data readily available in their local cache, most of the time. Since the majority of architectures will feature a large number of cores in the near future, dedicating a few cores to a dedicated task should not have a major negative impact on performance.

There is however a tradeoff between the increased locality of shared data and the loss of locality caused by the need to fetch context variables (i.e. thread-local variables used in synchronized blocks) from the original processing core. Furthermore, executing a synchronized block on a dedicated core requires a transfer of control to this core. This is not the case in the traditional way of executing Java programs, in which all synchronized blocks are executed locally. We therefore have to figure out when the cost of the transfer of control is covered by the performance gain caused by the availability of data in the cache memory of the dedicated core.

In the context of this internship, we focus on evaluating the tradeoff between executing synchronized blocks locally and migrating them to improve data locality. Preliminary results show that our technique is profitable on a 8-core microprocessor with a moderate number of shared cache lines accessed by each critical section (>15 cache lines), if a reasonable number of context variables are used. These encouraging results lead us to decide to further pursue this analysis in the context of a PhD thesis, in which we plan

to implement the automatic migration of critical sections towards dedicated cores in a Java Virtual Machine.

The remainder of this report is organized as follows. In Section 2, we present chosen works from the research literature related to the main issues we are faced with in the context of our project. In Section 3, we present the experiments we performed in the context of this internship. We conclude in Section 4.

2. Related work

In this section, we analyze the existing work which, to our knowledge, best tackles the major issues raised by our project. Since, in the context of our project, we aim to evaluate the tradeoff between the performance gain caused by improved data locality and the overhead caused by transfers of control on multicore architectures, the major issues we will have to deal with will be performance-related. This is especially true since multicore architectures raise complex issues regarding performance: programmers who write code for these architectures are often faced with problems such as bottlenecks caused by contention, poor locality of shared variables, and costly transfers of control and data between cores.

In Section 2.1, we focus on cache locality issues. Then, in Section 2.2, we focus on the major cause of overhead in multicore architectures : contention. In Section 2.3, we raise the issue of inter-process communication. Section 2.4 is dedicated to fast remote execution of code between cores. In Section 2.5, we discuss ways to improve performance on heterogeneous architectures. Finally, we conclude this analysis in section 2.6.

2.1. Cache locality

In this section, we focus on cache locality issues, about which two recent works have provided us with considerable insight: Pesterev et al.'s DProf [13] and Koufaty et al.'s bias scheduling [11].

DProf

DProf is a profiling tool that dynamically tracks indicators associated with data structures instead of code locations in order to allow for more efficient localization of cache-related bottlenecks. It offers four views of the collected data: the *Data Profile view* shows the number of cache misses for each of the most common data types; the *Miss Clarification view* shows which types of misses (due to sharing, associativity or capacity overload) are most common for each data type; the *Working Set view* indicates which data types were most active, how many of each were active at any given time and the cache associativity sets used by each data type; finally, the *Data Flow view* shows the most common sequences of functions that reference particular objects of a given type.

To create these views, DProf collects two kinds of data: *path traces* and *address sets*. A path trace records all accesses to a single data object. An address set is a data structure that contains the address and type of every object allocated during execution. Both of these data structures are built from the two kinds of raw data ultimately collected by DProf: *access samples* and *object access histories*. An access sample records information for randomly

chosen memory-referencing instructions, and an object access history is a complete trace of all instructions that read or write a particular data object, from when it was allocated to when it was freed. To collect access samples, DProf uses AMD's proprietary Instruction Based Sampling (IBS), but similar facilities are available on Intel chips.

Pesterev et al. show that, in two case studies, DProf provides much more helpful information than both `lock_stat` (a profiling tool that reports statistics about locks) and OProfile (a traditional execution profiler). Furthermore, performance experiments show that DProf's overhead varies depending on several variables (IBS sampling rate, number of debug register interrupts triggered by second, etc.). However, it is usually reasonably low and remains lower than 15%.

DProf is interesting to us at two levels. On a first level, the tool itself will probably be useful to monitor bottlenecks if we do not manage to compensate the cost of transfers of controls by the performance increase caused by cache locality efficiently enough in our project. On a second level, DProf's implementation itself is extremely interesting, since it uses special CPU instructions to monitor cache misses in order to infer cache locality penalties. Pesterev et al. provide advanced information on how they use these instructions in their article; we will not, however, discuss this point into more details in this document, for fear of getting too technical. The reader is invited to directly refer to Pesterev et al.'s paper [13] for more information.

Bias scheduling

Koufaty et al. [11] propose to optimize the scheduling of processes on performance-asymmetric multicore architectures with an algorithm that evaluates the performance gain of running an application on a fast core rather than on a slow core. Their approach estimates the number of internal stalls, caused by local (internal to the core) cache misses and delays, and external stalls, caused by accesses to shared last level caches, memory and I/O. To estimate internal stalls, they compute the number of cycles during which the front-end of the machine is not delivering instructions to the back end, thus effectively measuring the time during which cores are idling due to a lack of instructions. To estimate external stalls, they measure the number of requests serviced outside the core. The details of these mechanisms are not clearly stated in their paper, but we plan to analyse their implementation. We hope that this will allow us to find out how to obtain efficient statistics about cache misses that we will use to compute our cache locality penalties.

Now that we have considered cache locality issues directly, we will focus on the related problem of contention management. Indeed, cache locality issues are often linked with contention, since a lot of cache misses are caused by access conflicts over data stored in a given cache line. This is the object of the next subsection.

2.2. Contention management

Contention is a major issue on the target architecture for our project, i.e machines that feature a large number of cores. Migrating critical sections to remote cores will re-

quire the use of shared structures for synchronizations and data transfers; therefore, we must find ways to limit contention since, in the context of our project, our main objective is to ensure good performance.

In this subsection, we first discuss strategies to limit the size of critical sections to limit the use of shared data structures. Finally, we focus on the locks themselves, whose efficiency is crucial to any multithreaded environment.

2.2.1. Limiting the sharing of global data structures

The most obvious way to limit contention is to limit the sharing of global data structures. Wickizer et al. wrote an operating system named Corey [15] that aims to be very efficient on architectures having a large number of cores. Corey limits the sharing of kernel data structures that are shared among cores in order to improve performance.

More precisely, Corey's motto is that *applications should control sharing*. Indeed, traditional operating systems tend to share a lot of data structures between kernel and user threads by default, which can lead to contention. Wickizer et al. argue that the kernel should arrange each data structure so that only a single processor needs update it, unless directed otherwise by the application.

Corey provides three basic operating system abstractions to allow applications to control inter-core sharing and to take advantage of architectures having many processing cores by dedicating some of them to specific operating system functions. These abstractions are *address ranges*, *kernel cores* and *shares*.

Address ranges allow applications to selectively share parts of address spaces, instead of being forced to make an all-or-nothing decision. Indeed, traditional OSes give only two choices for shared memory: applications can either use a single address space shared by all threads, or a separate address space per thread. This can be inefficient if only a subset of the threads need to access a data set. Address ranges are more flexible and thus allow for finer-grained optimizations.

Kernel cores allow applications to dedicate cores to kernel functions or data. For instance, a core can be devoted to interacting with a given device. Multiple cores can then communicate with this dedicated core via shared-memory IPCs when they need to interact with the device. In a traditional OS, each core would directly interact with the device, using locks to access these data structures concurrently, which would drastically impede performance due to contention. With Corey's architecture, there is no direct contention over the data structures of the device driver. This limits the overhead caused by locks.

Shares allow applications to dynamically create lookup tables for system data structures and determine how these tables are shared. Indeed, many kernel operations involve looking up identifiers in tables to obtain pointers to internal kernel data structures. The use of these tables can be costly due to contention: in mainstream OSes, they are usually either shared between all the threads of a process (this is typically the case for Unix file descriptors, for instance) or between all processes (this is typically the case for Unix process identifiers). With Corey, these tables are shared only between the cores that need them: each of an application's cores starts with one share (its *root share*), which is private to that core. Then, if two cores wish

to share a share, they create one and add its ID to their private root share (or to a share reachable from their root share). A root share doesn't use a lock (since it is private), but a shared share does. When an application requests the creation of a new kernel object, it decides which share will hold the identifier.

Thus, a share maps application-visible identifiers to kernel data structures. Each core can use all the shares reachable from its root share. Contention only arises when two cores manipulate the same share, and overhead due to contention can be greatly limited by constricting the sharing scope of shares that are only reachable by a subset of the cores.

Performance experiments show that Corey performs 25% faster than Linux when using 16 cores on a MapReduce task and a Web Server task. Hardware event counters show that these improvements are due to Corey's original handling of shared memory.

Corey's idea of arranging data structures so that a single processor accesses it is reminiscent of our project, we plan to migrate synchronized blocks to dedicated cores in order to limit the number of processors that access shared data (ideally, one, if there is a single dedicated core). However, in our project, we will be able to detect where shared variables are accessed (thanks to synchronized blocks); this is not the case with Corey in which client applications can access shared variables anywhere.

Other projects have focused on the concept of limiting shared data to improve performance. Baumann et al. [2] propose a new type of operating system named the *multi-kernel* which uses one kernel per core with all of the kernel data structures replicated on each core. This is another way to limit the sharing of global data structures. Actually, multikernels limit the sharing of data structures so much that they completely prohibit shared memory. Fährdrich et al. [6] propose another OS that follows this same paradigm. Not sharing data structures at all goes much further than limiting the sharing of data structures: this is a matter of inter-process communication. We will discuss these issues more in detail in Section 2.3. Before discussing these matters that bring us well beyond traditional contention management mechanisms, we will focus on ways to optimize the most common tool in traditional contention management techniques : the lock.

2.2.2. Optimizing locks

Another way to optimize multithreaded applications is to optimize the internal implementation of locks. Two opposing locking strategies are generally used: to wait for a lock to be released, a thread can either spin, actively polling a resource in a loop; or block, i.e. putting the thread to sleep and trusting the scheduler to wake it up when needed.

Spinning provides high reactivity (no context switch needed, instant detection of lock releases) but wastes significant CPU resources. Moreover, on an overloaded system using spinlocks, the scheduler often preempts lock holders to wake waiting threads up; these threads then waste their time slices actively polling a resource that cannot be released. This phenomenon can drastically impede performance.

Block-based approaches usually fare better than spinlocks under high load but their lack of reactivity intro-

duces high overheads on the critical path of computation, thereby increasing the likelihood that other threads will encounter contention and block. This causes a vicious cycle of extremely slow lock handoffs known as a *convoy*. This is especially programs that use fine-grained locking since their critical sections usually take less time to execute than a single context switch.

Hybrid solutions that somewhat improve performance have been developed, but they still perform rather poorly, simply offering other tradeoffs between reactivity and CPU usage. However, Johnson et al. [10] proposed an interesting alternative to spinlocks and blocking locks that does not seem to suffer from this issue.

According to Johnson et al., contention management and load management (i.e. scheduling) are two orthogonal problems that should be treated separately. They propose a library that relies on a novel algorithm, known as *load control*, that regulates the number of active threads depending on the load. Basically, to acquire a lock, a thread either (1) blocks, if there are too many threads running or (2) spins otherwise. That way, threads are very reactive since they use spinlocks as much as they can and performance never collapses because the load is regulated.

To limit the load, threads are only allowed to run if they are unable to register themselves into a bounded array known as the *sleep slot buffer* because it is full. More precisely, the proposed algorithm—implemented by the load control library—works as follows: to acquire a lock, a thread spins until a) the lock is released or b) it is able to register itself into the sleep slot buffer. In the latter case, the thread blocks until it is removed from the sleep slot buffer or until 100ms pass, whichever comes first. Once the thread wakes up it restarts the lock acquire process as if it just arrived.

To control the number of running threads, a *controller* increases or decreases the size of the sleep slot buffer and wakes up threads when they get kicked out of the buffer. Thus, the controller acts as a local scheduler that only manages the load whereas the threads themselves only manage contention: contention and load management are effectively decoupled.

Experiments show that on various benchmarks, load control performs just as well as spinlocks under low load. Performance gracefully degrades under high load instead of collapsing as is the case with spinlock-based algorithms. Block-based algorithms' performance also degrades gracefully under high load, but they are initially much slower. Load control therefore offers the best of both worlds. It is worth noting that the spinlocks and blocking locks in the experiments used state-of-the-art locking algorithms, which shows how efficient Johnson et al.'s approach is. Load balancing is orthogonal to our work, since it aims to improve locking by making locks load-aware, whereas we aim to improve locking by improving data locality.

2.3. Efficient inter-process communication

Another way to improve performance in multicore environments is to improve communication between processes. This issue has been mainly tackled by OS designers and researchers, since OSes have to handle multiple processes running on multiple cores. Two opposing approaches are

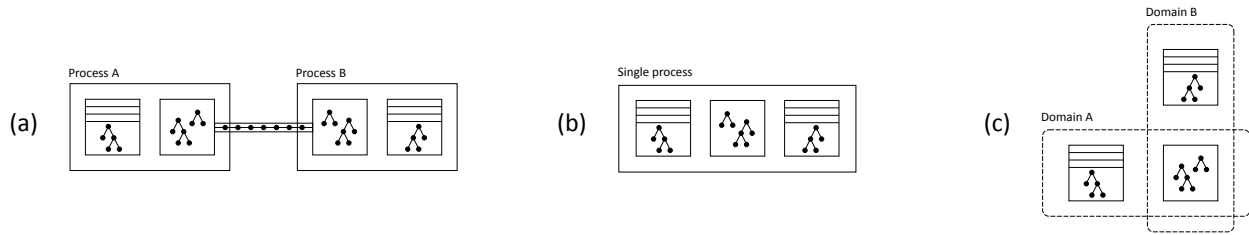


Figure 1: Two components sharing data (a) using standard, marshalled communication over pipes, (b) directly, both components being embedded into the same process, and (c) using Opal’s approach with memory protection domains decoupled from execution domains (Source: [4]).

generally used for the communication between processes: shared memory and message-based communication.

This dichotomy is mainly an opposition between two isolation paradigms. Some (Chase et al. [4], for instance) argue that memory should be easily sharable. They design, for instance, OSes with global address spaces to improve sharing. Others argue (Baumann et al. [2], Fährdrich et al. [6, 8]) that, on the contrary, shared memory should be prohibited by the OS. Both approaches have various repercussions on performance, safety, and ease of development. In this subsection, we focus on performance only, safety issues being beyond our scope. Performance-wise, the main difference between the two approaches is that when shared memory is directly used by applications, the OS cannot control the use of shared memory areas, which can lower performance. As for the ease of development, this matter is beyond the scope of this document since we only plan to change the inner workings of an JVM in our project, without having any direct influence on the development of higher-level applications on a semantic level.

2.3.1. Shared memory

Shared memory is the most common way by which processes communicate in traditional OSes. Its main advantages are simplicity and flexibility. However, sharing memory can cause contention problems that need to be overcome with safety and security mechanisms.

The Opal operating system, designed by Chase et al. [4], provide a good example of performance-optimized communication by shared memory. Opal uses a single 64 bit address space for all applications. This allows programs to directly use pointers to memory locations that they do not own, without the need for translation mechanisms between address spaces. Since sharing pointers to memory between programs with Opal is so easy, data is rarely ever copied from one program to another—transfers are all zero-copy. Traditional operating systems usually offer mechanisms for zero-copy communication, but they tend to be very complex. Another way to use easy zero-copy mechanisms between two components would be to place them in the same process, but this compromises protection and hinders modularity. Opal’s decorrelates memory protection from execution, thereby allowing distinct processes to directly work on the same data. This approach, illustrated in Figure 1, combines simplicity with modularity.

In Opal, segments of memory can be shared across processes with protection mechanisms based on capabilities. In order to allow several applications to work together

safely on in-memory data structures, data from a shared segment is accessed through shared procedures that are also stored in the same segment. This technique permits to maintain a high level of isolation and modularity while still benefiting from the very high performance provided by Opal’s direct sharing of data.

Chase et al. show that Opal allows several applications from the Boeing CAD system that use the same data to directly work together on in-memory data structures without the need for marshalization, copy, or address-space translation. A custom benchmark is also used to evaluate Opal’s performance. This benchmark uses three tools, (1) a producer that creates and deletes fixed-size records, (2) a consumer that maintains its own index structure on the same data and (3) a mediator that keeps the index up to date between the consumer and the producer. Experiments with this benchmark show that using directly shared segments (i.e. without address translation) is as efficient as embedding the three tools into the same process, and up to ten times more efficient than using three processes that communicate via traditional IPCs.

Opal’s use of a global address space is therefore very efficient. However, it has several drawbacks, in particular, it uses an abstraction layer above the hardware that makes impossible for applications to use advanced, low-level optimizations on complex architectures that feature a large number of cores.

2.3.2. Message-based communication

Recent works on OSes optimized for multicore architectures show a tendency to discard shared memory altogether and to vouch for communication via messages only instead. Indeed, communication via shared memory has several drawbacks, in particular, on architectures with a large number of cores, using shared memory presupposes cache coherency mechanisms. These mechanisms are increasingly complex to implement and their overhead is getting worse as the number of cores increases. Some recent processors do not ensure coherence between caches. Furthermore, at a lower level, these mechanisms are implemented by message passing. Allowing applications to directly handle the message passing mechanisms allow them to send and receive just as many messages as they need, without the need to assert the strong invariant of cache coherence, thereby improving performance.

We are now going to discuss two recent research OSes for multi-core (or multi-processor) architectures that make use of message-passing as their sole means of communication: Barrelfish (a multikernel) and Singularity.

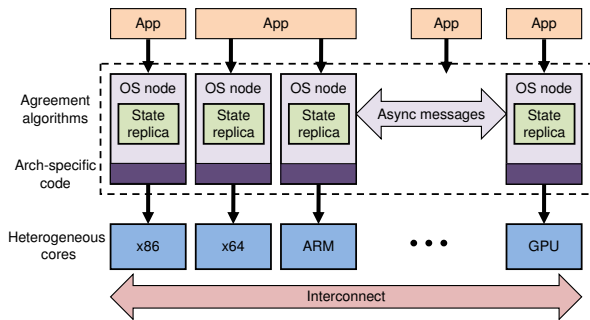


Figure 2: General architecture of Barrelfish (Source: [2]).

Barrelfish

Baumann et al. [2] argue that designing OSes that handle better architectures having a large number of cores¹ is a crucial research goal. Indeed, commodity computer systems containing multiple cores and/or processors are becoming more and more common. Most mainstream OSes have historically been developed for monoproductors or multiproductors having a limited number of processing units, using a model that scaled well in this context (global structures with locks, communication via shared memory, etc.). This model however has difficulties scaling well with newer hardware due to contention. Given the increasing number of cores in modern architectures and the complexity of interconnections between processing units, computer systems resemble more and more medium-scale distributed systems: Baumann et al. argue that designing an OS like an actual distributed system could improve scalability.

A Multikernel is a new type of operating system proposed by Baumann et al. that aims to this problem. Multikernels are multithreaded, with an instance of the OS running on each available core. Cores do not share global structures as in a traditional OS; instead, the OS state is replicated on each one of them. Cores do not communicate via shared memory, relying on cache coherence algorithms; instead, all communication is explicit and happens through asynchronous messaging.

Baumann et al. wrote an experimental version of a Multikernel named Barrelfish. In Barrelfish, a *CPU driver* and a *monitor* are bound to each core. CPU drivers handle system calls, schedule processes and threads on their cores and perform other low-level core-bound operations. Monitors are processes running on each core that coordinate the system-wide state via messages and encapsulate most of the higher-level functions that are usually found in a traditional kernel. The architecture of Barrelfish is shown in Figure 2.

The ability to send messages from one process to another within the same core or across cores is provided to user applications. However, applications can also communicate through shared memory like in a regular OS.

Baumann et al.’s experiments evaluate the performance of Barrelfish’s basic features. They focus on the OS’ handling of concurrency, messaging, computation and I/Os.

¹This is reminiscent of Corey. Indeed, Corey, Barrelfish and as we will see later, Singularity, all constitute recent attempts at dealing with this problem.

Through these experiments, Baumann et al. show that the performance of their messaging facilities is comparable to L4’s IPCs, which is rather good even though L4 is not a fully-featured, widely used microkernel. Performing a TLB shutdown (i.e. invalidating pages in the Translation Lookaside Buffer) and therefore unmapping pages is much faster on Barrelfish than on traditional OSes when the underlying hardware uses more than 14 cores on their test configuration, showing Barrelfish’s greater scalability in this context. They also show that Barrelfish is faster than Linux for sending and receiving messages through the IP loopback (these tasks involve the messaging, buffering and networking subsystems of the OS). Other experiments show that Barrelfish has a comparable performance to that of Linux for compute-bound and IO workloads.

Even though Barrelfish is an experimental OS, its performance rivals that of classical, highly-optimized OSes on computer systems that have a large number of cores, which tends to show that the Multikernel approach effectively delivers in terms of scalability. Let us note, however, that traditional OSes may have not been optimized for these systems yet though, so we cannot be sure they will not scale well.

We have now seen Baumann et al.’s take on the use of message-passing only for communication between processes and threads. Barrelfish is not the only such operating system that has been designed, however. Fährndrich et al.’s Singularity [6] (and Helios [12], which is based on Singularity) also explored this path: we will focus on this operating system in the next subsection.

Singularity

Fährndrich et al. also propose to design an operating system that uses message-passing as their sole means of communication. It is worth noting that their main motivation for building such an OS is not performance, but modularity and security. In the current subsection, we mostly focus on performance issues.

Singularity [6] aims to overcome a major drawback of message-based communication: data is often copied instead of shared when using the event-driven paradigm, which leads to poor performance. To overcome this, Singularity uses a new programming language based named Sing# that provides efficient messaging capabilities.

Messages are exchanged over bi-directional channels. These channels can also be used to exchange channel endpoints or pointers to memory location. This last point is crucial: it allows for fast data transfers by removing the burden of copying memory blocks.

In Singularity, processes are isolated and individually garbage-collected. Therefore, passing data from one thread to another could be an obstacle to safe garbage collection. Sending data pointers to allow for zero-copy messaging seems like an even more complex task: one could wonder which garbage collectors are responsible for such pointers. Singularity uses strong ownership invariants to solve this problem.

Let us consider this approach more in detail. When a message is sent, references to message arguments or transferred data blocks can be found in both the sending and receiving threads, for instance. To solve this issue, the Sing# compiler statically checks that processes only access memory that they own and that a memory block al-

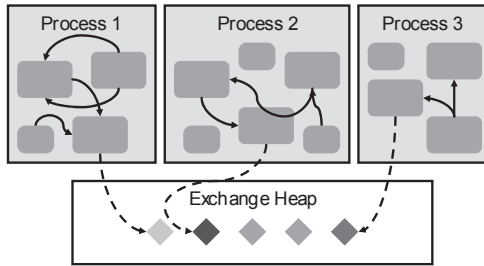


Figure 3: Singularity's exchange heap (Source: [6]).

ways belongs to a single thread at any given time. To this end, data on the GC heap is separated from data in the *exchange heap*, as illustrated on Figure 3. The exchange heap is used to exchange all data between processes, i.e. messages and memory pointed to by exchanged pointers. At any time, each memory block is owned by a function in a process: ownership is automatically passed from a function (and thread) to another through simple rules—when a variable is passed as a parameter, its ownership is passed to the called function for the whole time of its execution, for example. The only exception to this automatic ownership management strategy is when memory pointers are exchanged between threads: to explicitly track memory blocks from the exchange heap whose pointers are exchanged *via* messages, a special type of object, *TCell*, is provided. Sing#'s strong type-checking mechanisms combined with this system of ownership management allows the whole system to dismiss hardware memory protection altogether, which removes a major cause of overhead.

Performance experiments show that Singularity is faster than Linux and Windows for certain tasks. Performing a system call, switching between two threads, and, more importantly, sending a message are faster on Singularity. Moreover, the authors show that sending a block of memory from one thread to another is fast and does not depend on the size of the block. Linux and Windows are slower at performing this task, especially for larger blocks, since they copy the data instead of just passing a pointer.

This analysis of Singularity and Barrelfish showed that inter-core communication based on specialized algorithms improves data transfer performance. These algorithms will be a major source of inspiration to us since we will need to design efficient communication algorithms to transfer control and variables between cores when migrating synchronized blocks. Moreover, the works of Baumann et al. and Fährndrich et al. showed that limiting the use of shared memory could improve performance by limiting the underlying work of cache coherency algorithms, which is the main basis of our project: we wish to delegate synchronized blocks to dedicated cores in order to limit the overhead caused by the transfer of cache lines between cores.

In this subsection, we discussed solutions to transfer data from one thread to another. We will now focus on a related problem: the transfer of control between threads.

2.4. Fast transfer of control

Optimizing the transfer of data from one thread to another is a good way to improve performance on multicore architectures, since it is linked with the major problem of

these architectures: contention. However, another way to improve performance is to speed up the transfer of control from one thread to another, i.e. calling code from another thread. Indeed, we are in dire need of efficient mechanisms to dispatch the execution of code to threads: our project consists in finding efficient ways to transfer the control from one thread to another to delegate the execution of synchronized blocs.

Calling remote procedures from one thread to another is usually done through Remote Procedure Calls, or RPC. RPCs usually permit to call remote procedures from one core to another, as well as from a computer to another on a network. Such calls encapsulate the underlying transmission mechanisms, allowing developers to call procedures transparently without having to worry whether calls are local (between threads) or remote (over a network). Such an abstraction is usually viewed as a good thing, since it permits to write distributed systems on a network with very simple semantics. We are not interested in such capabilities, however: we only wish to dispatch the execution of code between cores in an efficient way.

In this subsection, we focus on three works that provide solutions to these issues : Bershad, et al.'s URPC [3], Huang et al.'s LRPC [7] and Ford et al.'s proposition of a migrating thread model [5].

URPC

URPC [3] aims to facilitate the design of highly-specialized RPC systems. Using custom RPC systems instead of generic ones makes it possible to use semantics that are not always provided by standard RPC packages. The loss of genericity also permits to improve performance, since unused parts of the RPC mechanism can be removed, thereby removing unnecessary overhead (messages, CPU time).

URPC comprises a runtime library, a stub generator and an name server, all of which are highly customizable. The URPC runtime library (see shaded boxes in Figure 4) makes it possible to customize the *protocol machines* of the RPC mechanism (*Client PM* on Figure 4) as well as *communication services*. The protocol machines encapsulate the RPC topology, call semantics and failure semantics of the RPC system. The protocol machine can be customized using *Cicero*, a protocol description language that handles advanced features such as multithreading. On the client side, the protocol machine's routine for assembling and disassembling messages (*Client RPC*) is factored out of the protocol machine. This routine allows, for instance, to customize marshalling and unmarshalling functions. The communication services component comprises routines to facilitate point-to-point communication between RPC participants, using a generic send/receive model that makes it possible to model any type of communication between nodes.

The stub generator includes special annotations to RPC signatures to support advanced communication schemes. These extensions allow for peer-to-peer communication (as well as standard client/server communication) and communication schemes that require multiple calls to complete for a given callback, such as asynchronous communication.

The name server uses a generic naming structure that accommodates for different naming models (a simple colon-separated list of attributes). It handles simple and

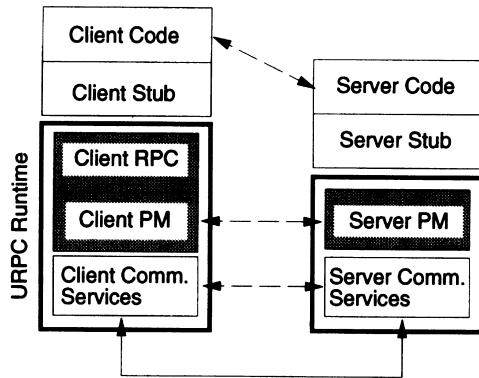


Figure 4: Architecture of the URPC runtime library (Source: [3]).

multiple end-point lookups to support unicast as well as multicast communication. In order to locate sets of hosts corresponding to given characteristics, numerical boolean relationships (such as $<$, $>$ or \neq , for instance) can be used to select hosts based on their attributes.

Bershad et al. show through examples that URPC can be used to design various types of RPCs, such as multicast RPCs, asynchronous RPCs and callback RPCs. These two last possibilities show that URPC is suitable for our project. Furthermore, their implementation of a special type of RPC that use at-most-once failure semantics shows a 10% performance gain over SUN RPCs, which shows that custom RPCs can be at least as fast as normal RPC.

On a final note, Bershad et al. advertise URPC as a tool to prototype RPC systems. However Barrelfish uses URPC as its main communication mechanism, which shows that URPC can be used in complex projects where performance is an issue.

LRPC

LRPC [7] (for *Lightweight Remote Procedure Call*) is a RPC library optimized for communication between cores on the same machine. This is particularly interesting to us since we do not need to transfer control to remote machines in our project.

With LRPC, calls to procedures are implemented with kernel traps. When a client starts a RPC, a kernel trap is called. The kernel validates the caller, creates a call linkage, and dispatches the client's thread to the server process. The client provides the server with a pointer to an argument stack as well as its thread of execution. The kernel switches to this thread with an upcall. After the procedure call completes, control and results return through the kernel back to the point of the client's call.

Three main optimizations are used. First, the kernel switches threads by updating a few registers (the stack pointer, the virtual memory registers, and a few other ones to pass the information described above) and by performing an upcall. This is much more efficient than performing a true context switch. Second, arguments are only copied to argument stacks (in shared memory) then back into the client thread, instead of needing up to seven copies in with regular RPCs². Third, on multicore architectures, server

²The following copies are usually needed with unoptimized RPCs: from the client thread to the message, from the sender

processes are cached on idle processors. Thus, on a procedure call, if the server process is available on an idle core, the RPC can be processed very quickly, since a context switch is not needed. It is also worth noting that, according to Huang et al., the multiprocessor implementation avoids contention on data structures as much as possible.

LRPC is 2 to 3 times faster than Taos³ RPC on a simple benchmark performing 100,000 RPC requests in a loop. On four processors, LRPC is 3.7 times faster than on a single processor with the same benchmark. This shows how interesting LRPC is for our project: it provides fast transfer of control capabilities on multiprocessors.

Migrating thread model

Ford et al. propose another compelling way to speed up transfer of control between threads: the *migrating thread model* [5]. In traditional OSes, threads belong to processes and cannot move from one process to another. With the migrating thread model, a thread abstraction moves between processes with the logical flow of control.

To manage this, Ford et al. decouple the *execution context* from the *schedulable thread of control*. The execution context encapsulates the state of the registers, program counter, stack pointer, and references to the containing process and designated exception handler, whereas the schedulable thread of control represents the thread itself in which the program is executed.

Ford et al. implemented their migrating thread model above Mach 3.0's kernel. In practice, threads are based into the kernel, and make incursions into processes (in user mode) via upcalls. Moving a thread from one process to another is cheaper than a context switch: the kernel only reproduces part of the functionality implemented by the scheduler. For instance, since the kernel is now calling the server rather than vice-versa, it no longer needs to save and restore the server's registers on every RPC. Moreover, thanks to the server's first-hand knowledge of the RPCs, it no longer needs to create, translate, and consume reply ports⁴ to match a reply to its request. Also, the kernel no longer needs to manage message buffers since the data is directly copied from source to destination.

Migrating threads are especially powerful when combined with RPCs, allowing for fast transfer of control on a single machine. Ford et al. implemented their thread model on the Mach 3.0 kernel. They show that regular RPCs require 5 times as many instructions and 4 times as many load/store operations than RPCs that use the migrating thread model.

This ends our subsection about fast transfer of control and our discussion of optimizations for homogeneous architectures. Indeed, until now, we discussed possible performance optimizations based on the hypothesis that the target architecture's cores share the same performance, ISAs, and characteristics. We will now extend our discussion to heterogeneous architectures.

domain to the kernel domain, from the kernel domain to the receiver domain, from the message to the server stack, from the receiver domain to the kernel domain, from the kernel domain to the sender domain, from the message into the client result.

³Taos is the operating system developed for the DEC SRC Firefly multiprocessor workstation [14].

⁴A port is a protected message queue for communication between tasks in the Mach kernel.

2.5. Heterogeneous architectures

Heterogeneous architectures, i.e. architectures whose processing cores differ in computing power, ISAs, and other characteristics, are getting more and more common. Even standard PCs can be viewed as heterogeneous architectures, since several of their components (NICs, GPUs) can be used as independent processing cores. Being able to optimize the performance of our approach (i.e. migrating critical sections to dedicated cores) on heterogeneous architectures could be an interesting long-term goal. This is the object of this subsection.

Bias scheduling

Koufaty et al. propose a scheduling algorithm that aims to improve performance on a common case of heterogeneous architectures: those having two types of cores only (faster and slower ones, called *big* and *small* cores, respectively). Their algorithm is novel in that it doesn't require offline profiling or dynamic sampling on all cores: it collects all of its data on-the-fly.

Koufaty et al. propose new metrics to evaluate the performance gain of executing a process on a big core rather than a small core (i.e. *application bias*). We mentioned in Section 2.1 that in addition to the traditional *Clocks Per Instruction* (CPI) metric, they also monitor internal stalls (caused by accesses to resources internal to the core) and external stalls (caused by accesses to shared last level caches, memory and I/O). They show that a process has a *big core bias* (i.e. its speedup from running on a big core compared to a small core is large) when the number of stalls (especially external stalls) is low. The number of internal and external stalls is measured at runtime and used to compute application bias. As this bias is not constant during the execution of an application, it is measured over a sliding instruction window. This allows the algorithm to schedule more efficiently applications whose behavior changes throughout their execution.

When the system is unbalanced, the scheduler tries to migrate threads from the busiest core to the idlest core, using application biases to optimize its choice of threads to move. When the system is balanced, the runqueues of each core are periodically checked: the scheduler looks for processes that would benefit from being swapped from big to small cores and conversely according to their application bias.

Since heterogeneous chips are uncommon, Koufaty et al. emulate one via an homogeneous quad-core Intel Xeon processor. They show that the standard approach of down-clocking some of the cores to turn them into small cores is not representative of architectures that have structurally diverse cores. This is an issue because these architectures are expected to be the obvious choice for manufacturers, since it makes it possible to optimize the smaller cores for physical size and energy efficiency. Therefore, they use proprietary tools to enable a debug mode on some cores that reduces instruction retirement from four to one micro-op per cycle. They show that this throttling method gives performance results similar to actual heterogeneous architectures whose small cores are in-order whereas their big cores are out-of-order.

To evaluate performance, Koufaty et al. perform experiments with heterogeneous workloads (i.e. workloads whose

processes have very different biases) based on the SPEC CPU2006 benchmark. On average, they obtain a 9% performance improvement over the default Linux scheduler. This is close to an upper bound found by running the tests on the same processor without throttling any of the cores. They also perform experiments with more homogeneous workloads and obtain a 5% gain on average. They show that this improvement is due to the fact that even though the application biases are similar overall, they vary throughout the execution of each benchmark.

Koufaty et al.'s work is extremely interesting to us, since it permits to improve the scheduling of applications on a common case heterogeneous architectures: performance asymmetric CPUs. Other research works have focused on more general heterogeneous architectures, considering each processing unit in a system (such as NICs or GPUs) as a normal CPU on which processes can be executed.

Multikernels

We discussed multikernels in Section 2.3.2. Multikernels have been designed with such heterogeneous architectures in mind. Indeed, Baumann et al. argue that architectures are getting increasingly diverse (memory hierarchies, instruction sets, interconnects, etc.) and that OSES have been statically optimized for the most common architectures at a low level for now. Given the varying nature of workloads and the diversity of hardware designs in modern computer systems, they argue that this approach will likely not be efficient enough anymore in the near future. Baumann et al. propose that designing OSES like distributed computer systems allows them to adapt better to various architectures, just like network applications are able to dynamically adapt to architecturally diverse networks.

Multikernels are mostly hardware-independent. The only architecture-specific modules are the messaging transport mechanism and the interfaces to the hardware (CPUs/cores and devices). Bauman et al.'s implementation of a multikernel, Barrelfish, uses a *knowledge database* containing information regarding the underlying hardware (found by polling and measurements) that it uses to optimize its communication scheme. This database can be used to select appropriate message transports for inter-core communication or to allow for NUMA-aware memory allocation, for instance.

Helios

Helios [12], an OS based on Singularity (cf. Section 2.3.2), provides an interesting alternative solution to the problem of handling heterogeneous architectures. Helios does not rely on multiple identical kernels, instead, it uses a main kernel and satellite kernels that (1) export a standardized set of kernel abstractions on all cores, (2) provide transparent and unified IPCs⁵ for communication between cores, and (3) support heterogeneous ISAs via an intermediate bytecode to which applications are compiled.

Helios uses an *affinity metric* to automatically place applications on cores. This affinity metric allows developers to express the fact that two processes would benefit from running on the same core (or the opposite). This is rem-

⁵Helios' IPCs uses Singularity's zero-copy mechanism for local message parsing if used between two local threads. Otherwise, Helios transparently marshals messages between the initiating thread and the remote service.

Name	Layer	Communication scheme	Contention	IPCs/Tr. of control	Heterogeneity	Performance	Scalability	Safety	Customizability
Multikernel	OS	URPC	+		+		+	+	
Singularity	OS	LMP (local, zero-copy)	+	+			+	+	
Helios	OS	LMP + RMP (distant, marshalling)			+				
Corey	OS	Shared memory	+			+	+		
Opal	OS	Shared memory (global address space)		+		+			
Bias scheduling	library	n/a			+	+			
URPC	library	URPC		+		+			+
LRPC	library	LRPC	+	+		+	+		
Migrating threads	library	Migrating threads		+		+			
Load scheduling	library	n/a	+			+	+		

Figure 5: Summary of the improvements from the research works analyzed. ‘+’ is a shorthand for “substantial improvement”. VMKit and DProf are not listed because the chosen criteria were not applicable.

iniscient of the implementation of bias scheduling that we discussed earlier.

We will not elaborate any further on the mechanisms provided by operating systems such as Multikernel and Helios to handle complex heterogeneous architectures since this subject is beyond the scope of our project for now.

2.6. Conclusion

In this section, we have reviewed a selection of scientific publications that tackle best the issues we are faced with in our project, at least to our knowledge. Figure 5 summarizes the improvements provided by each research work we mentioned our analysis. The next section describes the experiments we performed on thread migration to ensure that our approach is viable.

3. Experiments

For our experiments, we used two machines, featuring 8 and 48 cores respectively. Since we only received the second (48-core) machine recently, data from this machine is still incomplete; most experiments have been performed on the 8-core machine only. These two machines are described below.

- The first machine, codenamed *bossa*, is a Mac Pro (3.1) featuring two X5472 quad-core 3.2Ghz Xeon CPUs. Each of the two processors combines two dual-core dies⁶. Bossa uses a Intel 5400 chipset with a 1.6Ghz system bus and 10GB of DDR2 SDRAM. It runs Debian Linux 5.0.5 with a 64-bit 2.6.31 Linux kernel.
- The second machine, *amd48*, uses four “Magny-Cours” 6172 Opterons, featuring two six-core dies. The four CPUs are interconnected via HyperTransport links. Amd48 uses an AMD RD890 chipset with 32GB of DDR3-SDRAM. It runs Mandriva Linux 2010.1 (x64) with a 2.6.34 Linux kernel.

⁶In the context of CPUs, a die is a small block of semi-conducting material on which one or several cores along with other components are fabricated. Communication between dies is slower than communication between cores on the same die. Cores on the same die usually share a cache chip (typically, the L2 cache).

Our aim is to determine the parameter spaces under which each of these two solutions are the more efficient:

- **Solution (a):** execute the critical sections from within the thread they belong to. The main overhead of this solution is caused by cache misses when shared variables are not available in the L1 cache and have to be fetched from remote cores (or higher-level caches).
- **Solution (b):** execute the critical sections on a dedicated core⁷. This solution induces two major causes of overhead: (1) it requires a transfer of control from the local core to the dedicated core (modeled by a local RPC in our benchmarks) each time a critical section is executed and (2) since local context variables are not directly available on the distant core, they have to be transferred, similarly to shared variables in solution a.

The simplest way to modelize this issue is to consider that Solution (b) is more efficient than Solution (a) when:

$$\text{Equation (A): } c_{rpc} + n_{context}m < n_{shared}m$$

Where c_{rpc} is the cost of a transfer of control (i.e. a RPC between cores), $n_{context}$ and n_{shared} are the number of context and shared cache lines respectively, and m is the average cost of fetching a cache line from a distant core. This is a simplification, of course, because m varies depending on where the variable is available in the cache hierarchy. Both $n_{context}$ and n_{shared} are dependent on the programs being executed in the Java Virtual Machine and cannot be evaluated precisely. On a given architecture, however, c_{rpc} , and m can be measured. If the values of these three variables are known, we can obtain a gross estimate of the parameter space, i.e. the values of $n_{context}$ and n_{shared} , in which Solution (b) is more efficient than Solution (a). This is the object of this section.

In Section 3.1, we evaluate the cost of cache misses, in Section 3.2 we evaluate the cost of a transfer of control, and in Section 3.3, we provide a rough estimate of the values of $n_{context}$ and n_{shared} for which Equation (A) is verified.

⁷We ignore the case where several cores are dedicated to the execution of critical sections for now, to simplify the analysis.

3.1. Rough estimates

This subsection evaluate the cost of m and c , i.e. the average cost of a cache miss and the cost of a transfert of control.

3.1.1. Cost of cache misses

In order to evaluate the cost of cache misses, we used *memal*, a program from the benchmark suite used by Wickizer et al. to evaluate the performance of the Corey operating system [15]. *Memal* uses two threads, a server and a client thread, each of these threads being pinned on two different cores specified by the user. A memory area (C array) is shared between the two threads. The server thread uses low-level techniques to place the memory area within its L1 cache, L2 cache or an owned memory area. The client thread then reads each cache line from this buffer, and the benchmark returns the average number of cycles needed to read each cache line.

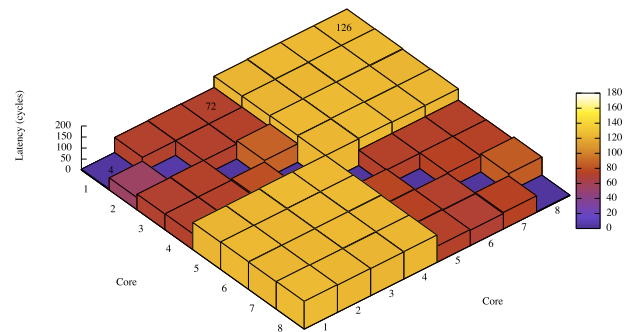
Memal uses low-level techniques initially proposed by Yotov et al. [16] in order to prevent the compiler from removing data accesses and to prevent cache prefetching from altering the results. More precisely, the shared memory area is not accessed by the client in a linear fashion. Instead, each 64-bytes word of this memory area contains an address to another word of the same memory area, in such a way that by reading the first address, using the read address to know which address to read next, and so on until enough 64-byte words have been read, one can read the whole memory block in a pseudo-random fashion. In the end, the last value is passed to a system call that does not produce any effect (a `printf` with an empty string as its first parameter). Thus, the reads are dependant from each other, and the last read value is used: this prevents the compiler from thinking some of the memory accesses can be removed as an optimization. Moreover, since 64-byte words are accessed randomly, the processor cannot efficiently prefetch cache lines. This technique might prove important for us in some of the benchmarks we are about to write, even though for now we found simpler—albeit platform-dependent—workarounds to avoid optimization and prefetching issues.

The results of the *memal* benchmark on *bossa* are shown in Figure 6. Each graph shows the number of cycles needed for a client on Core i to access data located on Core j , in the L1 cache, L2 cache and in the RAM, respectively. The core IDs of the client and server threads are on the X and Y axes of each graph⁸.

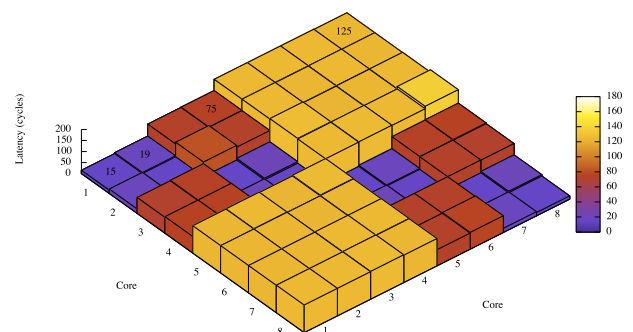
Bossa uses two Xeon X5472 processors, each of which has two dual-core dies. Each core has a local L1 cache, and the L2 cache is shared between each two cores from the same die. There is no L3 cache, the third level is the RAM itself. Since *Bossa* uses NUMA (Non-Uniform Memory Access), access times to the RAM are processor-dependent.

The first graph exhibits three different latency access times to the L1 cache. accessing the local L1 cache takes

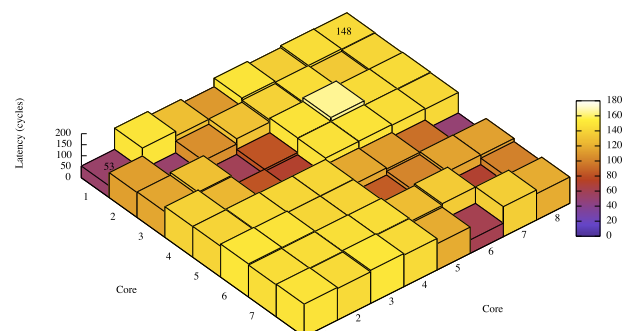
⁸Since all graphs are symmetric, it does not matter which axis represents client and server threads: it takes as much time for a client thread on Core i to access a variable in a given cache level (or in the RAM) from Core j than for a client thread on Core j to access a variable in a given cache level (or in the RAM) from Core i .



Latency of L1 cache accesses



Latency of L2 cache accesses



Latency of RAM accesses

Figure 6: Latencies of cache RAM accesses relative to the cores between which cache lines are transferred, on *bossa*, as measured with the *memal* benchmark. Each processor uses two Dies, each die use two cores. Cores 1 to 4 belong to the first processor (Die 1: Core 1 and 2; Die 2: Core 3 and 4), and Cores 5 to 8 belong to the second processor (Die 1: Core 5 and 6; Die 2: Core 7 and 8).

Placement of the client and the server	Number of cycles	Number of L1 cache misses	Number of L2 cache misses
Same die	~302	1	0
Same processor, different die	~412	1	1
Different processors	~968	1	1

Results of the `null_rpc` benchmark with both synchronization variables on the same cache line

Placement of the client and the server	Number of cycles	Number of L1 cache misses	Number of L2 cache misses
Same die	~547	2	0
Same processor, different die	~525	2	2
Different processors	~1190	2	2

Results of the `null_rpc` benchmark with both synchronization variables on two different cache lines

Figure 7: Results of the `null_rpc` benchmark on bossa. Note: the number of cache misses was measured on the server only, therefore, the total number of cache misses is the double of the number given.

```

1 for (i = 0; i < number_of_rpcs; i++)
2 {
3   rpc_requested = 1;
4
5   while (!done_with_rpc)
6     pause;
7
8   /* Protected code */
9
10  done_with_rpc = 0;
11 }

```

Critical section from the client thread

```

1 for (i = 0; i < number_of_rpcs; i++)
2 {
3   while (!rpc_requested)
4     pause;
5
6   /* Protected code */
7
8   rpc_requested = 0;
9   done_with_rpc = 1;
10 }

```

Critical section from the server thread

Figure 8: Pseudo-code of the critical sections from the client and sever thread in the `null_rpc.c` benchmark.

about 5 cycles. Accessing a variable from the L1 cache of any other core on the same processor takes about 75 cycles, because in each case the variable has to be fetched by the L2 cache, then passed back to the local L1 cache. When accessing data located on the L1 cache of a remote processor, 125 cycles are needed. The second graph shows that when the data is available from the L2 cache, the access times are 15-20 cycles. When the data has to be fetched from a local L2 cache, the latency is about 75 cycles. Fetching data located in the cache of a remote processor takes about 125 cycles. Finally, the third graph shows that the access time to a variable from the RAM itself is highly variable if done locally (i.e. on a single processor), with access times varying between 50 and 140 cycles. Accessing data from the RAM takes about 150 cycles if the variable has been modified by the remote processor. This difference of latency

between variables accessed locally (i.e. same processor) and remotely (i.e. different processor) shows that because of NUMA, access times to the RAM are not constant.

3.1.2. Cost of transfers of control

To evaluate the cost of transferring the control between two processors, we implemented a simple benchmark that uses two threads, a client and a server, pinned on cores whose IDs are provided by the user. In this benchmark, the clients request the executions of a given number of RPCs by the server. The benchmark uses the PAPI library [1] in order to measure the number of cycles per RPC, or the number of L1/L2 cache misses. The code of this benchmark (`null_rpc.c`) is available in the appendix.

RPCs are not implemented using blocking locks, instead, we used spinlocks, in order to minimize the reaction time. We are of course wasting cycles (cf. Section 2.2.2), but in a system having a large number of cores, one could consider that wasting cycles from a number of unused cores is negligible. Minimizing the cost of RPCs is crucial however, since we hope the performance gain from increased data locality will be greater than the loss in performance caused by RPCs. The full code of this benchmark is provided in the appendix (`null_rpc.c`).

Figure 8 shows the pseudo-code of the critical sections of the client and the server threads. For each iteration, the client thread sets the `rpc_requested` variable to 1. This automatically starts a RPC. Then, the client busy-waits⁹ for the server to set the `done_with_rpc` variable to 1. When this happens, the server is done with the RPC. The client sets its `done_with_rpc` variable to 0 and starts the process again. On the server side, for each iteration, the server busy-waits for the `rpc_requested` variable to be set to 1. When this happens, the client has requested a RPC. The server just resets the `rpc_requested` variable to 0 and sets the `done_with_rpc` variable to 1, to warn the client that it is done with processing the RPC.

We ran this benchmark on bossa using two different configurations. In the first configuration, the two synchronization variables, `rpc_requested` and `done_with_rpc` are

⁹The pause x86 instruction provides a hint to the processor that the code sequence is a spin-wait loop. The processor uses this hint to avoid the memory order violation in most situations, which greatly improves processor performance[9].

```

1 for (i = 0; i < number_of_rpcs; i++)
2 {
3     *local_sv = 0;
4
5     while (CAS(global_sv,
6               0,
7               &local_sv) == 0)
8         pause;
9
10    /* Protected code */
11
12    while (*local_sv == 0)
13        pause;
14 }

```

Critical section from the client thread

```

1 for (i = 0; i < number_of_rpcs; i++)
2 {
3     while (*global_sv != 0)
4         pause;
5
6     /* Protected code */
7
8     **global_sv = 1;
9     *global_sv = 0;
10 }

```

Critical section from the server thread

Figure 9: Pseudo-code of the critical sections from the client and sever thread in the first version of the benchmark from Section 3.2.

on the same cache line, and in the second configuration they are on different cache lines. In the second configuration, we actually leave two full cache lines between the two variables, because on our architecture, when a cache line is fetched, the next one is automatically prefetched. The results are presented in Figure 7.

Figure 7 shows that if the client and the server are on the same core with synchronization variables on the same cache line (resp. on different cache lines), each RPC takes about 302 (resp. 547) cycles; with two cores on the same processor but on different dies, each RPC takes about 412 (resp. 525) cycles, and with two cores on different processors, each RPC takes about 968 (resp. 1190) cycles.

3.1.3. Cost comparison

Now that we have latencies for all cache levels and an estimate of the cost of RPCs between cores, we can estimate, in a restricted case (i.e. one client only), when our approach (Solution (b)) is faster than the traditional approach (Solution (a)).

- When the client and the server are on the same die, the minimum cost of a transfer of control is about 300 cycles, and the cost of transferring a variable from one core to another is between 15 and 20 cycles. With 20 more cache lines used by shared variables than cache lines used by context variables, Equation (A) is verified.
- When the client on the server are on the same core, but on a different die, the cost of a transfer of control

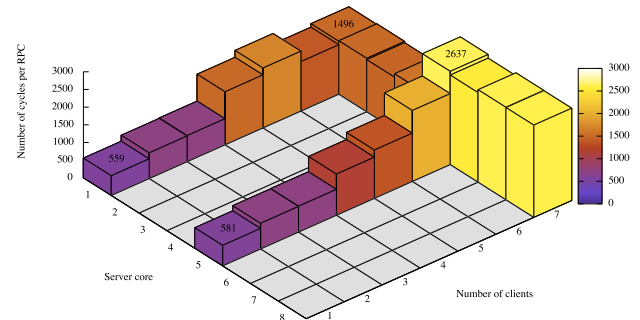


Figure 10: Number of elapsed cycles per RPC on bossa, depending on the number of clients and the Core n on which the server thread is pinned. The clients are pinned to Cores $n+1 \bmod 8$ to $n+c \bmod 8$ where c is the number of clients.

is about 500 cycles, and the cost of transferring a variable from one core to another is about 75 cycles. With 7 more cache lines used by shared variables than cache lines used by context variables, Equation (A) is verified.

- When the client and the server are on different processors, the cost of a transfer of control is about 1000 cycles, and the cost of transferring a variable from one core to another is between 125 and 150 cycles. With 8 more cache lines used by shared variables than cache lines used by context variables, Equation (A) is verified.

These initial results show that Equation (A) can be verified for a relatively low number of cache lines (8 to 20). Of course, the considerations from this subsection are not based on a microbenchmark that emulates an environment with multiple clients and with actual context and shared variables. Instead, since only one client was used, the transfers of the cache lines containing synchronization variables were limited and contention was inexistent. These first results are therefore very optimistic. Finding more accurate results with a more realistic microbenchmark is the object of the next sections.

3.2. NULL RPCs with multiple clients

The second phase of our work was to improve the initial `null_rpc.c` benchmark to use several clients (one client per core) in order to better simulate Solution (b). We wrote a first prototype whose critical sections are shown in Figure 9. The prototype uses $c+1$ synchronization variables where c is the total number of clients. Each client has a `local_sv` synchronization variable, and the server has its own synchronization variable, named `global_sv`. We do not use regular, blocking locks to avoid slow reaction times caused by context switches. Instead, we use a Compare-And-Swap (CAS)¹⁰ primitive to optimize re-

¹⁰The *compare-and-swap* primitive executes the comparison and the ‘swapping’ of a variable atomically: first it compares

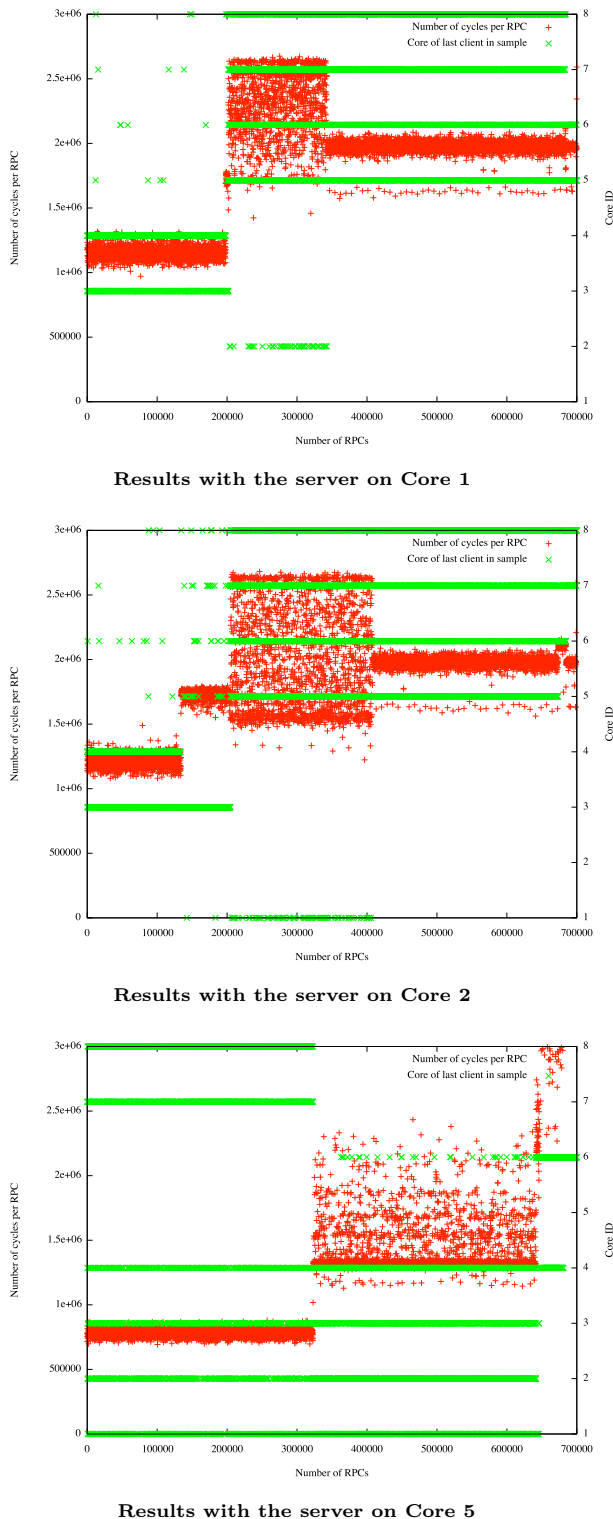


Figure 11: Sampled executions on bossa (100000 RPCs per client, 100 RPCs per sample).

sponse times. The general idea is that each client tries to set the `global_sv` synchronization variable to the ad-

the value of a variable with a given value, and if they are equal, the variable is set to another value given as a parameter. The compare-and-swap primitive returns a boolean to indicate whether the swapping happened.

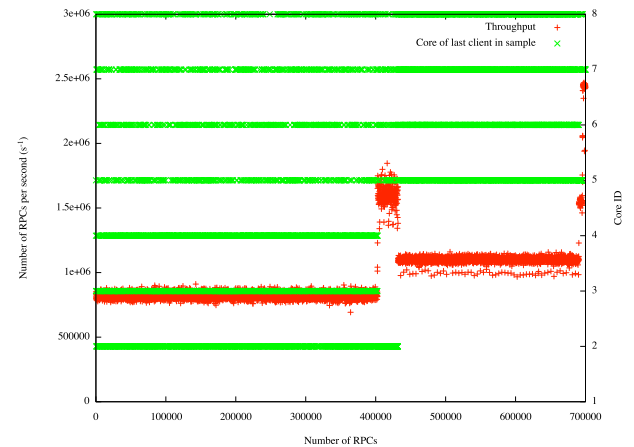


Figure 12: Sampled executions on Core 1 using a Test-And-Test-And-Set loop, on bossa.

dress of its local `local_sv` variable using a CAS primitive in a busy-wait loop. When a client manages this operation, this means it managed to get the server to service its RPC. Since the server has the address of `local_sv` variable from the client whose RPC is being serviced, it can set it to 1 to indicate when it is done servicing the RPC. Since, we are only implementing null RPCs, the actual bodies of the RPCs on the server side is empty. The aim of the benchmark is to measure the elapsed time per RPCs, in number of cycles. Since we want all clients to execute their RPCs simultaneously, we use a synchronization barrier before entering the critical sections. The results of this first prototype on bossa are shown in Figure 10.

As shown in Figure 10, only a subset of all possible configurations were tested. First, we tried to use 1 to 7 clients with the server thread pinned to the first core of each processor and each client pinned to Cores $n+1 \bmod 8$ to $n+c \bmod 8$, where n is the core number of the server thread and c the number of clients. Then, for 7 clients, we ran our benchmark with the server on each of the 8 cores. Unsurprisingly, increasing the number of clients increases the latency, because of contention over the cache lines. More interestingly, the experiments give much better results when the server is located on the first processor. This clearly shows that the bus between the two processors is not symmetric.

Sampling

Even though we obtained coherent results with our first prototype, those are not very accurate because we noticed that small changes in the initialization code for the client and server threads caused the results to vary. For instance, adding a simple `i++` instruction before the synchronization barrier could, in some cases, add up to 500 cycles to the results with 7 clients (while preserving the same general tendency, i.e. that RPCs take less time to execute if the server is located on the first processor). To deal with this issue, we tried various approaches; in particular, we tried to service a single RPC per client and to save the order in which the clients were serviced, in order to understand whether the variance in execution times was correlated with this order. This solution proved unsuccessful.

```

1 while (CAS(global_sv,
2         0,
3         &local_sv) == 0)
4     pause;

                    TAS loop

1 for (;;)
2 {
3     if (global_sv == 0)
4     {
5         if (CAS(global_sv,
6             0,
7             &local_sv) == 0)
8             continue;
9
10        else break;
11    }
12
13    pause;
14 }

```

TATAS loop

Figure 13: TAS and TATAS loops from the benchmark prototype.

In order to better understand what was really happening during the experiment, we decided we needed to find a way to plot temporary results about the experiment in a temporal fashion. This led us to divide the execution in samples.

Sampling the execution of the benchmark was a very simple idea that allowed us to better understand what happened at a lower level. The idea was to stop the execution of RPCs for a very short while every x RPCs to store data about the experiment. We chose to record (1) the number of cycles per RPC in each sample and (2) the core of the last client serviced in each sample. An example of result of such an execution is shown in Figure 11 (100,000 RPCs per client, 100 RPCs per sample).

Figure 11 shows that bossa's hardware does not ensure fairness, which is an issue we did not think of at first. On this figure, the red points show the average throughput during the previous sample (i.e. (1) from the previous paragraph) and the green points show the core of the last client serviced in each sample (i.e. (2) from the previous paragraph). The execution is split in phases, during which only some clients have their requests serviced by the server. Switches between phases happen when some of the clients are done executing their RPCs. For instance, let us consider the first graph in Figure 11. In this figure, the server is pinned to the first core of the first processor. The two clients located on the second die of this same processor are mainly serviced first, even if a few other clients get serviced occasionally¹¹. This is Phase 1. When these two clients are done with their RPCs (both are done almost

¹¹We only monitor the core of the last client in the sample in order to prevent the measurements from altering the result. In practice, that means that we only check which client gets serviced once in a while. This still shows a valid tendency regarding which cores get serviced, because the results clearly show that the execution works by phases during which only a subset of cores are serviced.

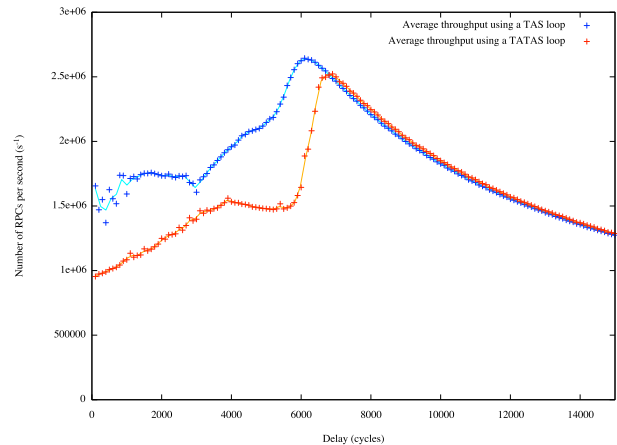
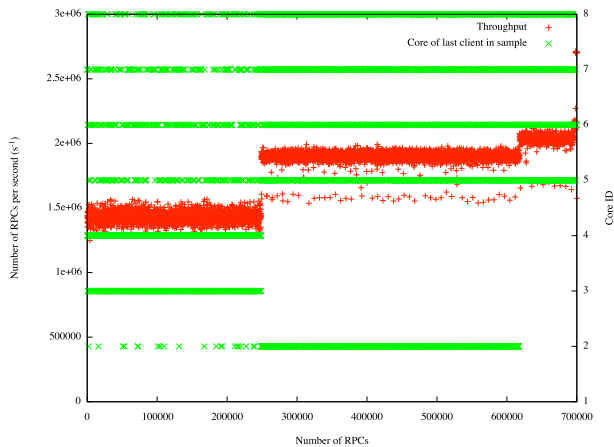


Figure 14: Average throughput as a function of the delay, using either a TAS or a TATAS loop, on bossa (7 clients).

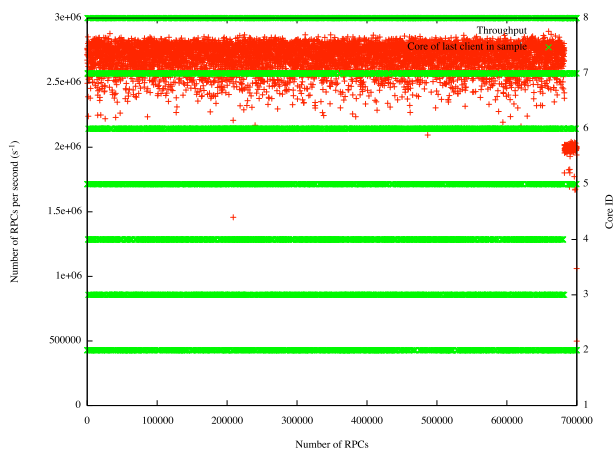
simultaneously, which shows their RPCs were serviced at a similar rate) we enter Phase 2: all other clients get serviced, although the client from the second core (same die as the server) gets serviced faster than the others. Finally, when this client is done, the 4 clients from the remote processor are serviced at about the same rate. This is Phase 3. The value and variance of the elapsed time per RPC is characteristic of each phase: in the first phase, RPCs are serviced in about 2700 cycles, with a deviation of about 500 cycles. Then in Phase 2, RPCs are serviced faster: about 1200 cycles on average, with a greater variation, which is probably due to the fact that RPCs from clients located on the local processor are executed much faster than RPCs from clients on the remote processor. Finally, in phase 3, RPCs take more time to execute than in Phase 2 (because they all come from remote clients) but faster than in Phase 1 (because less clients try to execute RPCs). The deviation is very low because all RPCs come from the distant processor (same latency). The two other graphs from Figure 11 show executions from Core 2 (Processor 1, Die 1) and Core 5 (Processor 2, Die 1). The general tendency seems to be that the two cores from the remote die on the local processor (relative to the server core) are serviced first, then the core from the local die is serviced simultaneously (albeit faster) than the cores. Some sub-phases might appear when some cores serviced at a similar rate do not finish their execution at the exact same time.

Test-And-Test-And-Set (TATAS)

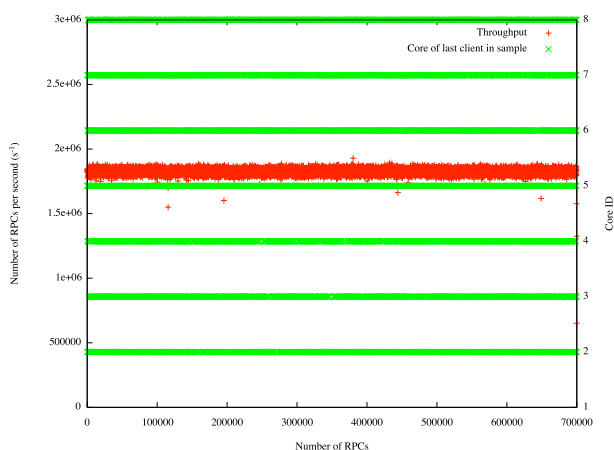
These experiments showed us that our first measurements were performed in a degenerate case in which fairness was not preserved. In order to solve this issue, our first approach was to use a Test-And-Test-And-Set (TATAS) loop instead of the Test-And-Set (TAS) loop we used until then (implemented with a Compare-And-Swap). In Figure 9, what we refer to a TAS loop is the loop from lines 7 to 10. Figure 13 shows the difference between a TAS loop and a TATAS loop. When using a TAS loop, all clients repeatedly try to set the value of the `global_sv` global variable to the address of the `local_sv` local variable in order to start a RPC. When using a TATAS loop, clients first check the value of the `global_sv` variables, and only



Delay = 1000 cycles (contended case)



Delay = 6100 cycles (maximum throughput)



Delay = 10000 cycles (high delay)

Figure 15: Sampled results for three different values of the waiting delay, using a TAS loop, on bossa. With a low delay, fairness is not ensured, which hints that there is a high level of contention. Increasing the delay improves fairness, until the benchmark is fully fair for the delay for which the throughput reaches its maximum (6100 cycles in our case, cf. Figure 14). Increasing the delay any further only lowers the throughput.

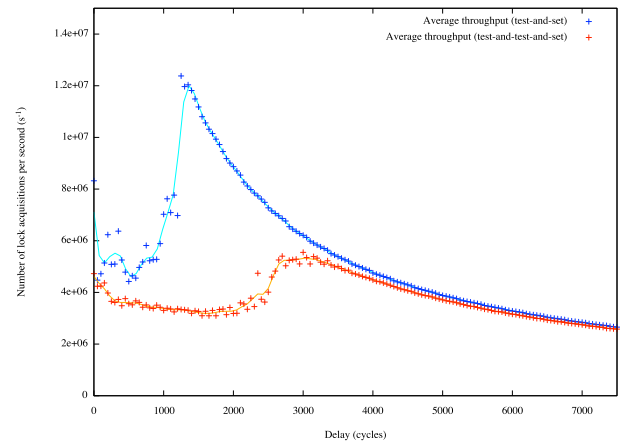


Figure 16: Number of lock acquisitions per second as a function of the delay, on bossa (7 clients).

try a CAS if this initial comparison hinted that the CAS would be successful. Using a TATAS is sometimes preferred to using a TAS, because it limits the number of comparisons. A TATAS loop uses more comparisons, but comparisons are less costly than CASes, because they don't require to acquire writing rights to cache lines. However, since there is no way to ensure that the value of the synchronization variable will change between the comparison and the CAS, the actual results of using a TATAS loop are hard to predict, and depend largely on the underlying hardware.

Figure 12 shows the result of a sampled execution of our prototype modified to use a TATAS loop instead of a TAS loop. This technique greatly improves equity: all cores are able to acquire the RPC lock, even though local cores are privileged. However, the performance is worse, probably because cache lines are more contended with a TATAS loop.

Delays

Another approach we tried to limit contention over cache lines was to add a busy-wait loop after each RPC request on the client side in order to limit the contention on the cache lines. We performed experiments to find the optimal delay. The results are shown on Figure 14 (red/orange curve). There is a clear optimum, in both TAS and TATAS configurations, at 3000 and 1300 cycles respectively. With no delay or a small delay, the throughput is slow, and the variance is strong. It is also highly dependent on the initial conditions, which explains the issues we had with our first prototype.

Figure 15 shows sampled results of the execution of our prototype with a waiting delay of 1000, 6100 and 10000 cycles between each RPC request, using a TAS loop. With 1000 cycles, the experience is more fair to the clients than with no delay (cf. Figure 11, first graph): indeed, all clients get served from the first phase, albeit at varying rates (not all clients terminate at the same time). Increasing the delay improves fairness, until we reach the maximum throughput of 6100 cycles (cf. Figure 14). With 6100 cycles, the experience is perfectly fair to clients (the whole experiment is a single phase), and the throughput

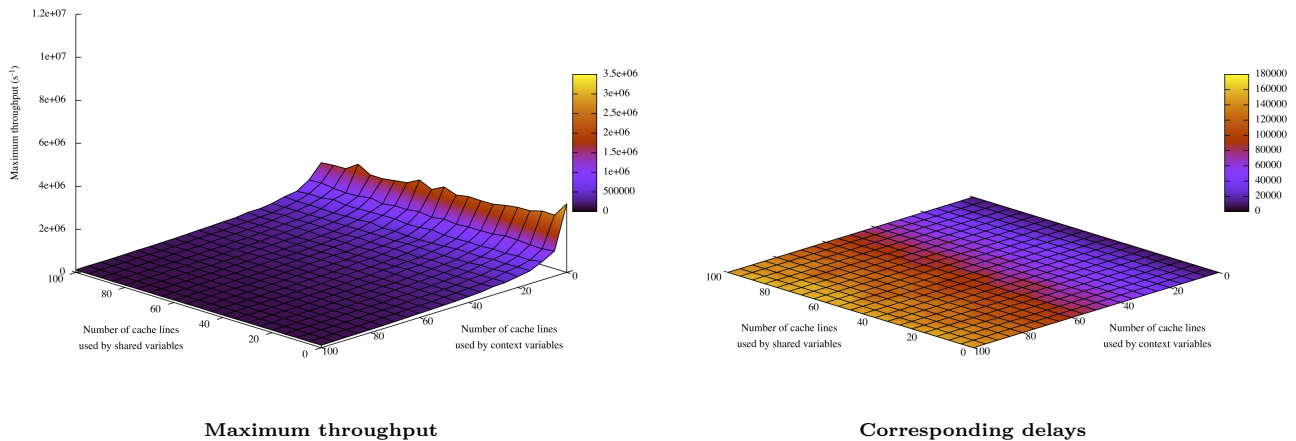


Figure 17: Maximum throughput and associated delays as a function of the number of context and shared variables, using RPCs (i.e. our implementation of Solution (b)) on bossa.

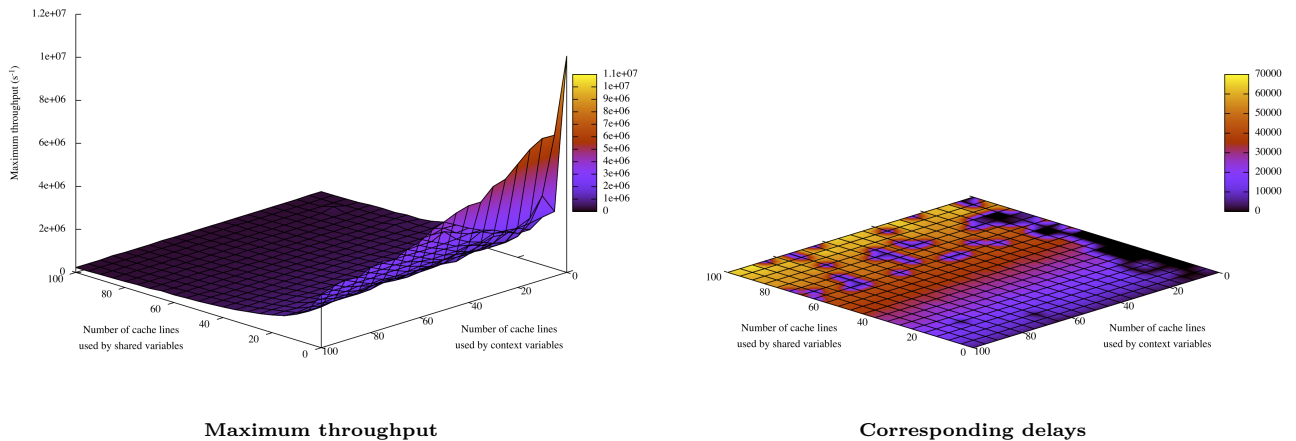


Figure 18: Maximum throughput and associated delays as a function of the number of context and shared variables, using locks (i.e. our implementation of Solution (a)) on bossa.

is maximal. When we increase the delay past this limit, fairness cannot be improved anymore and the throughput decreases. Our experiments show that the same observation can be made for the TATAS case: the maximum throughput corresponds to the minimal delay for which all clients are serviced fairly.

Waiting for a given delay on the client side is therefore a good way to ensure that our measurements are not limited by contention, and the maximum throughput of an experiment is a key value that characterizes the maximum number of RPCs per second that can be obtained in a given experimental case. This is a crucial point that will be used in later benchmarks (cf. section 3.4) to compare the results of such experiments.

3.3. Final benchmark

The final benchmark, whose code is presented in the appendix `benchmark.c`, implements two additional functionalities. It makes it possible (1) to execute critical sections on the clients based on lock acquisitions (this corre-

sponds to Solution (a)) and (2) to access/modify context and shared variables from within the critical sections. This is the object of the two next subsections.

3.3.1. Executing critical sections locally

Up to now, we focused on the development of a benchmark that modeled our solution, i.e. “Solution (b)”. In order to compare Solution (a) and Solution (b), we needed to modelize Solution (a), i.e. the execution of critical sections on the clients themselves. We implemented this solution using a global lock that all clients try to acquire in order to execute their critical sections. After the execution of a critical section, the lock is released. The pseudo-code of the critical section for Solution (a) is shown in Figure 20.

Figure 16 shows the throughput of this solution with both a TAS and a TATAS loop. Once again, the TAS solution is more efficient than the TATAS solution. We obtain a maximum throughput of about $1.2e+07$ critical sections per second as compared to $2.7e+06$ critical sections per second when we used RPCs (cf. Figure 14). As

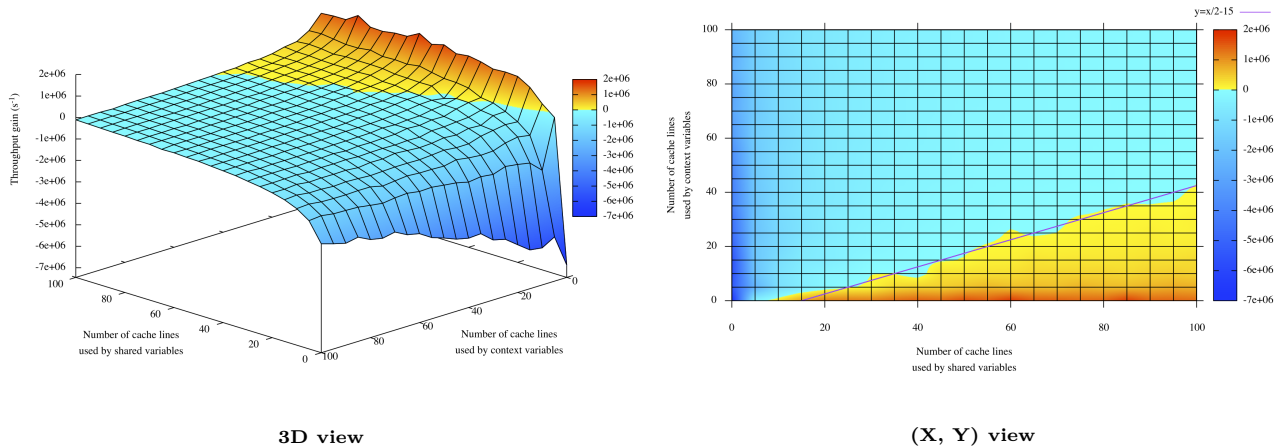


Figure 19: Maximum throughputs of our modelization of Solution (a) subtracted from the maximum throughputs obtained for Solution (b), on bossa.

```

1 for (i = 0; i < number_of_rpcs; i++)
2 {
3     *local_sv = 0;
4
5     while (CAS(global_sv,
6               0,
7               1) == 0)
8         pause;
9
10    /* Protected code */
11
12    global_sv = 1
13 }

```

Figure 20: Critical section for Solution (a).

explained before, we plan to compensate for this overhead with the potential performance improvement caused by the increase in locality when all critical sections are executed on a dedicated core. This is the object of the next subsection.

3.3.2. Context and shared variables

The last phase of the development of our benchmark was to add accesses to context and shared variables in and around the critical sections. We do not count the number of variables we access directly, since results could vary depending on where the variables are located in memory. Since the unit of data transfer between caches is the cache line, we allocate memory areas aligned with cache lines and access variables located every two cache lines (we skip one cache line between each variable to avoid issues with prefetching, as with `null_rpc.c`). One memory area per client is allocated for the context variables, and a global memory area is allocated for shared variables. Each memory access is a read/write access (`++` operator from C). Here is the algorithm we chose for data accesses:

- Shared variables are accessed (i.e. incremented in a loop) once in each critical section, i.e. on the client side with our modelization of Solution (a) and on the server side with our modelization of Solution (b).

- Context variables are accessed once outside each critical section, and once in each critical section. We have to access the variables outside the critical sections to repatriate them from the cache hierarchy of the server to the cache hierarchy of the client before each critical section. This properly models the behaviour context variables being accessed from within the critical section, and afterwards. This is different from the modelization that we used in Section 3.1 in which we considered context variables were ignored after having been used in a critical section. By considering that all context variables are always accessed after the critical sections, we modeled a worst case scenario, in which it is harder to compensate for the cost of the transfer of control and the transfer of context variables with the increased locality of shared variables.

Figure 17 shows the results of our benchmark for Solution (b). The first graph (labelled ‘Maximum throughput’) shows the value of the peak throughput for a given number of cache lines used by context and shared variables. The peak throughput is found as in the previous sections (cf. Figures 14 and 16), i.e. by finding the delay value for which the throughput is maximal. The graph shows that adding shared variables only has a slight impact on the throughput, which is what we expected: since all shared variables are in the local cache of the dedicated core, the overhead caused by shared variables is limited, thanks to the improved data locality. On the other hand, however, adding context variables is costly, because we have to fetch them from the client cores.

The second graph of Figure 17 show the corresponding delays we needed to reach the peak throughput. This is of no direct relevance on our final result, but allows us to better understand the inner workings of our benchmark and the validity of choosing a waiting delay as a way to avoid contention. The graph shows that the longer the experiment, the longer the delay needs to be. This is due to the fact that in a perfect scenario, all clients wait for the average time needed for the execution of a RPC: in this case, at the end of each RPC, a new unique client requests the

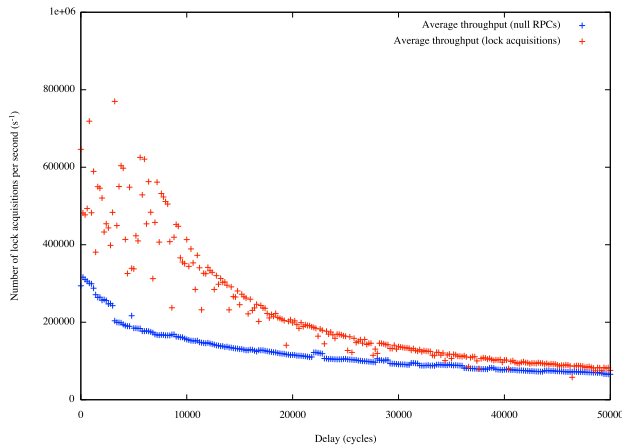


Figure 21: Throughput for Solutions (a) and (b) as a function of the delay, on amd48.

RPC, and there is no contention over the synchronization variables.

Similarly to Figure 17, Figure 18 shows the results of our benchmark for Solution (a). The peak throughput at (0,0) is much higher ($10.0e+06$ vs $3.1e+06$), since in this case we do not migrate critical sections. The situation is the opposite with this solution compared to Solution (b): indeed, this time, adding context variable only induces limited overhead whereas the use of shared variables is prohibitive. This is of course the results we expected since the reason why we were not satisfied with Solution (a) was that the use of shared variable was too costly. The second graph of Figure 18 shows the delays corresponding to the maximum throughput and, unsurprisingly, they get higher when the critical sections get more costly.

In order to compare the two solutions, we subtract the first curve from Figure 18 from the first curve from Figure 17. If the result is positive, this means that Solution (b) is more efficient than Solution (a). The result is shown on Figure 19 (using two different view). Our approach (Solution (b)) is more efficient than Solution (a) in the yellow/orange positive area.

This final result shows that, on bossa, with 7 clients on 8 cores, our approach (Solution (b)) is not profitable if less than 15 cache lines are shared. More precisely, it seems that Solution (b) is more efficient than Solution (a) when:

$$n_{shared} > 2n_{context} + 15$$

Where n_{shared} is the number of shared variables and $n_{context}$ is the number of context variables. This result is of course dependent on the hardware chosen and our experiment, but it gives us a general idea as to how efficient our approach can be.

3.4. Future experiments

The next phase will be to reproduce our experiments on amd48, in order to ensure that the results we obtain on a different architecture are similar. Moreover, since amd48 uses a large number of cores, it is a better example of a target architecture for our solution. Some experiments have already been reproduced: Figure 21 shows the throughput as a function of the delay for both Solution (a) and

Solution (b) on amd48. For this experiment, we used the same settings as we did for bossa in figures 14 and 16, except there are now 47 threads instead of 7 (one per core, excluding the server). The variance is higher than in our experiments with bossa, probably due to the increased complexity of the hardware. Due to the increase in contention, the throughput is much lower (less than $1e+06$ as compared to more than $1.2e+07$ with Solution (b) on bossa). However, on both bossa and amd48, the ratio between the peak throughputs for Solution (a) and Solution (b) is still around 1/2, which tends to show that, even with the increased contention issues caused by a large number of cores, we are confident that the results of our benchmarks will be similar on amd48.

4. Conclusion

To ensure that our approach—outsourcing critical sections to dedicated cores—is viable, more experiments will be needed. First, we will have to reproduce our results on amd48, which has a different hardware architecture and a much higher number of cores. From the results on amd48 and bossa, we believe we will have a better idea of how viable our approach can be on various hardware configurations. Second, we will have to find a way to estimate the average number of cache lines used by shared and context variables in real-world Java programs. Finally, the parameter space in which our approach is profitable is too restricted, we will have to find automated ways to guess whether migrating critical sections is profitable before executing them.

For now, we remain optimistic, since the experiments we performed during this internship show that migrating critical sections to dedicated cores can be profitable under restricted conditions on the number of cache lines used by shared and context variables: for at least 15 cache lines used by context variables and a moderate number of context variables, outsourcing critical sections to a dedicated core is more efficient than executing them locally. These results show the potential of this approach which we will continue exploring in the context of a PhD thesis.

Bibliography

- [1] PAPI, the Performance Application Programming Interface. <http://icl.cs.utk.edu/papi/>.
- [2] Andrew Baumann, Paul Barham, Pierre-Evariste Daggand, Tim Harris, Rebecca Isaacs, Simon Peter, Timothy Roscoe, Adrian Schüpbach, and Akhilesh Singhanian. The Multikernel: a new OS architecture for scalable multicore systems. In *Proceedings of the ACM SIGOPS 22nd symposium on Operating systems principles*, pages 29–44, 2009.
- [3] Brian N. Bershad, Thomas E. Anderson, Edward D. Lazowska, and Henry M. Levy. URPC: a toolkit for prototyping remote procedure calls. *The Computer Journal*, 39, no. 6:525–540, 1996.
- [4] Jeffrey S. Chase, Henry M. Levy, Michael J. Feeley, and Edward D. Lazowska. Sharing and protection in a single address space operating system. *ACM Transactions on Computer Systems*, 12:271–307, 1994.
- [5] Bryan Ford and Jay Lepreau. Evolving Mach 3.0 to a migrating thread model. In *Proceedings of the USENIX Winter 1994 Technical Conference on USENIX Winter 1994 Technical Conference*, 1994.
- [6] Manuel Fähndrich, Mark Aiken, Chris Hawblitzel, Orion Hodson, Galen Hunt, James R. Larus, and Steven Levi. Language support for fast and reliable message-based communication in Singularity OS. In *Proceedings of the 2006 EuroSys conference*, pages 177–190, 2006.
- [7] Y. Huang and C.V. Ravishankar. Lightweight remote procedure call. *ACM Transactions on Computer Systems (TOCS)*, 8, issue 1:37–55, 1990.
- [8] Galen Hunt, Mark Aiken, Manuel Fähndrich, Chris Hawblitzel, Orion Hodson, James Larus, Steven Levi, Bjarne Steensgaard, David Tarditi, and Ted Wobber. Sealing OS processes to improve dependability and safety. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007*, pages 341–354, 2006.
- [9] Intel. Using Spin Loops on Intel® Pentium® 4 Processor and Intel® Xeon® Processor.
- [10] F. Ryan Johnson, Radu Stoica, Anastasia Ailamaki, and Todd C. Mowry. Decoupling contention management from scheduling. *ACM SIGPLAN Notices*, 45, issue 3:117–128, 2010.
- [11] David Koufaty, Dheeraj Reddy, and Scott Hahn. Bias scheduling in heterogeneous multi-core architectures. In *Proceedings of the 5th European conference on Computer systems*, pages 125–138, 2010.
- [12] Edmund B. Nightingale, Orion Hodson, Ross McIlroy, Chris Hawblitzel, and Galen Hunt. Helios: heterogeneous multiprocessing with satellite kernels. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007*, pages 221–234, 2009.
- [13] Aleksey Pesterev, Nikolai Zeldovich, and Robert T. Morris. Locating cache performance bottlenecks using data profiling. In *Proceedings of the 2010 EuroSys conference*, pages 335–348, 2010.
- [14] Charles P. Thacker, Lawrence C. Stewart, and Edwin H. Satterthwake Jr. Firefly: A multiprocessor workstation. *IEEE Transactions on Computers*, 8:909–920, 1988.
- [15] Silas B. Wickizer, Haibo Chen, Rong Chen, Yandong Mao, Frans Kaashoek, Robert Morris, Aleksey Pesterev, Lex Stein, Ming Wu, Yuehua Dai, Yang Zhang, and Zheng Zhang. Corey: an operating system for many cores. In *Proceedings of the 8th USENIX Symposium on Operating Systems Design and Implementation*, 2008.
- [16] Kamen Yotov, Keshav Pingali, and Paul Stodghill. Automatic measurement of memory hierarchy parameters. In *Proceedings of the 2005 ACM SIGMETRICS international conference on Measurement and modeling of computer systems*, SIGMETRICS '05, pages 181–192, New York, NY, USA, 2005. ACM.

Appendix: source code

null_rpc.c

```

1  /* ##### */
2  /* null_rpc.c */
3  /* (C) Jean-Pierre Lozi, 2010 */
4  /* ----- */
5  /* This program should be compiled with the -O0 compiler flag. */
6  /* ##### */
7  #define _GNU_SOURCE
8
9  #include <papi.h>
10 #include <pthread.h>
11 #include <stdio.h>
12 #include <stdlib.h>
13
14 #define PAUSE __asm__ __volatile__ ("rep; nop" : : );
15
16 #define PADDING 0
17 /* L2 cache lines are 32 bytes long, but they're fetched two at a time. */
18 /*#define PADDING 64*/
19
20 #define NUMBER_OF_ITERATIONS 1000000
21
22 void *client_main(void *ptr);
23 void *server_main(void *ptr);
24
25 int client_core, server_core, papi_cache_miss_level;
26
27 /* This structure contains the booleans we use as flags to signal events. */
28 struct flags {
29     volatile unsigned char rpc_requested;
30     volatile unsigned char padding[PADDING];
31     volatile unsigned char done_with_rpc;
32 } flags;
33
34 int main(int argc, char **argv)
35 {
36     if (argc != 4)
37     {
38         client_core = 0;
39         server_core = 1;
40         papi_cache_miss_level = 1;
41
42         fprintf(stderr,
43             "Usage : %s client_core server_core papi_cache_miss_level\n",
44             argv[0]);
45         fprintf(stderr,
46             "Default values used (0, 1, 1)\n");
47     }
48     else
49     {
50         client_core = atoi(argv[1]);
51         server_core = atoi(argv[2]);
52         papi_cache_miss_level = atoi(argv[3]);
53
54         if (papi_cache_miss_level < 1 || papi_cache_miss_level > 2) {
55             fprintf(stderr,
56                 "Usage : %s client_core server_core "
57                 "papi_cache_miss_level\n",
58                 argv[0]);
59             fprintf(stderr,
60                 "papi_cache_miss_level is invalid. Exiting.\n");
61             exit(EXIT_FAILURE);
62         }
63     }
64
65     /* We use both a client and a server. The client calls remote procedures
66     from the server. */
67     pthread_t client, server;
68
69     flags.rpc_requested = 0;
70     flags.done_with_rpc = 0;
71
72     /* We create both threads. */
73     if (pthread_create(&client, NULL, client_main, NULL) < 0)
74         perror("pthread_create");
75
76     if (pthread_create(&server, NULL, server_main, NULL) < 0)
77         perror("pthread_create");
78
79     /* We wait for both threads to finish their work. */
80     if (pthread_join(client, NULL) < 0)
81         perror("pthread_join");
82
83     if (pthread_join(server, NULL) < 0)
84         perror("pthread_join");
85
86     /* Everything went as expected. */
87     return EXIT_SUCCESS;
88 }
89
90 void *client_main(void *ptr)
91 {
92     unsigned int j;
93     int i;
94     long long start_cycles, end_cycles;
95     cpu_set_t cpuset;
96
97     int event_set = PAPI_NULL;
98     int events[1];
99     long long values[1];
100
101     /* We want to pin this thread to the right core. */
102     CPU_ZERO(&cpuset);
103     CPU_SET(client_core, &cpuset);
104     /* We wish to avoid the server's core. */
105     CPU_CLR(server_core, &cpuset);
106
107     /* We set the thread's affinity. */
108     if (pthread_setaffinity_np(pthread_self(), sizeof(cpu_set_t), &cpuset) < 0)
109         perror("pthread_setaffinity_np");
110
111     /* We monitor either L1 data cache misses or L2 data cache misses. */
112     if (papi_cache_miss_level == 1)
113         events[0] = PAPI_L1_STM;
114     else
115         events[0] = PAPI_L2_STM;
116
117     /* We print miscellaneous informations about our computation. */
118     printf("Client here. I'm running on core %d.\n", sched_getcpu());
119     printf("Addresses : %x %x.\n", &flags.rpc_requested, &flags.done_with_rpc);
120     printf("Starting %d RPCs.\n", NUMBER_OF_ITERATIONS);
121
122     /* We initialize PAPI. */
123     if (PAPI_library_init(PAPI_VER_CURRENT) != PAPI_VER_CURRENT)
124         perror("PAPI_library_init");
125
126     if (PAPI_thread_init(pthread_self) != PAPI_OK)
127         perror("PAPI_thread_init");
128
129     /* We execute the first RPC separately, because it can be longer than others
130     if the server isn't initialized yet.*/
131     flags.rpc_requested = 1;
132
133     while (!flags.done_with_rpc)
134         PAUSE;
135
136     flags.done_with_rpc = 0;
137
138     /* We initialize PAPI and start recording events. */
139     if (PAPI_create_eventset(&event_set) != PAPI_OK)
140         perror("PAPI_create_eventset");
141
142     if (PAPI_add_events(event_set, events, 1) != PAPI_OK)
143         perror("PAPI_add_events");
144
145     if (PAPI_start(event_set) != PAPI_OK)
146         perror("PAPI_start");
147
148     /* We start counting cycles. */
149     start_cycles = PAPI_get_real_cyc();
150
151     /* We execute the right number of RPCs. */
152     for (i = 0; i < NUMBER_OF_ITERATIONS; i++)
153     {
154         flags.rpc_requested = 1;
155
156         while (!flags.done_with_rpc)
157             PAUSE;
158
159         flags.done_with_rpc = 0;
160
161         /*
162         In order to measure events in a limited part of the loop, these two 'if'
163         statements can be used around the instructions we wish to monitor.
164
165         if (PAPI_reset(event_set) != PAPI_OK)
166             perror("PAPI_reset");
167
168         if (PAPI_accum(event_set, values) != PAPI_OK)
169             perror("PAPI_accum");
170
171         */
172     }
173
174     /* We stop the measurements. */
175     end_cycles = PAPI_get_real_cyc();
176
177     /* We read the counters. */
178     if (PAPI_read(event_set, values) != PAPI_OK)
179         perror("PAPI_read");
180
181     printf("L%d cache misses : %lld\n", papi_cache_miss_level, values[0]);
182
183     /* We print the result. */
184     printf("Client : done. %lu cycles per null rpc. (total : %lld) %d\n",
185         (end_cycles - start_cycles) / (long long)NUMBER_OF_ITERATIONS,
186         end_cycles - start_cycles, j);
187
188     /* We kill the server. */
189     exit(EXIT_SUCCESS);
190
191     return NULL;
192 }
193
194 void *server_main(void *ptr)
195 {
196     unsigned int j;
197
198     cpu_set_t cpuset;
199
200     /* We want to pin this thread to the right core. */

```

```
201 CPU_ZERO(&cpuset);
202 CPU_SET(server_core, &cpuset);
203 /* We avoid the client's core. */
204 CPU_CLR(client_core, &cpuset);
205
206 /* We set the thread's affinity. */
207 if (pthread_setaffinity_np(pthread_self(), sizeof(cpu_set_t), &cpuset) < 0)
208     perror("pthread_setaffinity_np");
209
210 /* We print miscellaneous informations about our computation. */
211 printf("Server here. I'm running on core %d.\n", sched_getcpu());
212
213 /* We start answering RPCs. */
214 for (;;)
215 {
216     while (!flags.rpc_requested)
217         PAUSE;
218
219     flags.rpc_requested = 0;
220     flags.done_with_rpc = 1;
221 }
222
223 return NULL;
224 }
```

benchmark.c

```

1  /* ##### */
2  /* benchmark.c */
3  /* (C) Jean-Pierre Lozi, 2010 */
4  /* ----- */
5  /* This program should be compiled with the -O0 compiler flag. */
6  /* ##### */
7  #define _GNU_SOURCE
8
9
10 /* ##### */
11 /* Headers */
12 /* ##### */
13 #include <math.h>
14 #include <papi.h>
15 #include <pthread.h>
16 #include <stdint.h>
17 #include <stdio.h>
18 #include <stdlib.h>
19 #include <unistd.h>
20
21 #include <sys/resource.h>
22 #include <sys/time.h>
23
24
25 /* ##### */
26 /* Definitions */
27 /* ##### */
28 /*
29  Use blocking locks?
30
31  This parameter is experimental. When it's on, the benchmark uses blocking
32  locks but only certain configurations are supported (no sampling, for
33  instance). This mode cannot be accessed through command-line parameters
34  because it is more of an experiment than anything else.
35 */
36 #define USE_BLOCKING_LOCKS */
37
38 /* Default values for the command-line arguments. */
39 #define DEFAULT_NUMBER_OF_RUNS 1
40 #define DEFAULT_NUMBER_OF_ITERATIONS_PER_SAMPLE 10
41 #define DEFAULT_SERVER_CORE 0
42 #define DEFAULT_NUMBER_OF_ITERATIONS_PER_CLIENT 1000
43 #define DEFAULT_NUMBER_OF_CONTEXT_VARIABLES 0
44 #define DEFAULT_NUMBER_OF_SHARED_VARIABLES 0
45
46 /* Maximum number of cores, used to avoid malloc/realloc cycles. */
47 #define MAX_NUMBER_OF_CORES 1024
48
49 /* Maximum line size used when reading results from unix commands. */
50 #define MAX_LINE_SIZE 32
51
52 /* PAUSE instructions for spinlocks */
53 #define PAUSE __asm__ __volatile__ ("rep; nop" : : );
54
55 /* This macro performs the conversion between cycles and throughput. */
56 #define TO_THROUGHPUT(x) (cpu_frequency * 1000000 / (x))
57
58
59 /* ##### */
60 /* Global variables */
61 /* ----- */
62 /* FIXME: some volatile keywords may not be necessary. We declare all global
63  variables as volatile for now. */
64 /* ##### */
65 /* Environment data ===== */
66 /* Number of cores */
67 volatile int number_of_cores;
68 /* CPU frequency, in MHz */
69 volatile float cpu_frequency;
70 /* This array maps physical to virtual core IDs. */
71 volatile int *physical_to_virtual_core_id;
72
73 /* Execution parameters ===== */
74 /*
75  Critical sections
76 */
77 /* Critical sections can either be null RPCs (serviced by a single server) or
78  lock acquisitions. */
79 typedef enum _critical_sections_type {
80     NULL_RPCS,
81     LOCK_ACQUISITIONS
82 } critical_sections_type_t;
83 volatile critical_sections_type_t critical_sections_type;
84
85 /*
86  Execution mode
87 */
88 /* Three execution modes are available :
89  - MULTIPLE_RUNS_AVERAGED means that the test is run multiple times and the
90  averaged results are returned (along with the variance and standard
91  deviation).
92  - SINGLE_RUN_SAMPLED means that the test is only run once, and statistics
93  are gathered all along the execution.
94  - SINGLE_RUN_ORDERED means that the test is only run once and that the order
95  in which clients managed to execute their critical section is returned.
96 */
97 typedef enum _execution_mode_t {
98     MULTIPLE_RUNS_AVERAGED,
99     SINGLE_RUN_SAMPLED,
100    SINGLE_RUN_ORDERED
101 } execution_mode_t;
102 volatile execution_mode_t execution_mode;
103 /* Number of runs over which the results are averaged (used in the
104  MULTIPLE_RUNS_AVERAGED mode only). */
105 volatile int number_of_runs;
106 /* Number of iterations per sample (used in the SINGLE_RUN_SAMPLED only). */
107 volatile int number_of_iterations_per_sample;
108
109 /*
110  Execution settings
111 */
112
113 /* Core on which the server runs */
114 volatile int server_core;
115 /* Number of clients */
116 volatile int number_of_clients;
117 /* Number of iterations per client */
118 volatile int number_of_iterations_per_client;
119 /* Delay between RPCs, in cycles. */
120 volatile int delay;
121 /* ROW mode activated? */
122 volatile int read_only_wait;
123 /* Number of context variables */
124 volatile int number_of_context_variables;
125 /* Number of shared variables */
126 volatile int number_of_shared_variables;
127
128 /*
129  Measurements
130 */
131 /* Should we count cycles or events? */
132 typedef enum _measurement_type {
133     NUMBER_OF_CYCLES,
134     NUMBER_OF_EVENTS
135 } measurement_type_t;
136 volatile measurement_type_t measurement_type;
137 /* ID of the monitored PAPI event if measurement_type = NUMBER_OF_EVENTS */
138 volatile int monitored_event_id;
139 /* Should we perform the measurements on the server or the clients? */
140 typedef enum _measurement_location {
141     SERVER,
142     CLIENTS
143 } measurement_location_t;
144 volatile measurement_location_t measurement_location;
145 /* Should we return the result in throughput or cycles? */
146 typedef enum _measurement_unit {
147     THROUGHPUT,
148     CYCLES_PER_ITERATION,
149     TOTAL_CYCLES_MAX
150 } measurement_unit_t;
151 volatile measurement_unit_t measurement_unit;
152
153 /* Execution variables ===== */
154 /* Addresses of the rpc_done addresses for each thread */
155 volatile void ** volatile rpc_done_addresses;
156 /* Synchronization variables for the barrier between the server and the
157  clients, before the main loop */
158 volatile int * volatile ready;
159
160 /* Data specific to the MULTIPLE_RUNS_AVERAGED mode */
161 /* Results (cycles or number of events) for each iteration */
162 volatile double * volatile iteration_result;
163
164 /* Data specific to the SINGLE_RUN_SAMPLED mode */
165 /* Number of iterations per sample */
166 volatile float * volatile multiple_samples_results;
167 /* Address of the rpc_done variable from the core whose RPC was serviced last.
168  Translated into the core ID at the end of the computation. */
169 volatile void ** volatile multiple_samples_rpc_done_adrrs;
170
171 /* Data specific to the SINGLE_RUN_ORDERED mode */
172 /* Order in which the RPCs are processed */
173 volatile int * volatile order;
174
175 /* Global critical section variables ===== */
176 /* This structure models a cache line. It should be allocated on a 64-byte
177  boundary. */
178 typedef char cache_line_t[64];
179
180 /*
181  This structure models the data passed to threads upon their creation.
182  Variables for both the 'NULL RPCs' and 'lock acquisitions' modes are passed
183  in order to simplify the code.
184 */
185 typedef struct _thread_arguments_block_t {
186     /* Logical thread id (local, not provided by pthreads). */
187     int id;
188     /* 'NULL RPCs' mode synchronization variables ===== */
189     /*
190     Global synchronization variable used to limit the number of NULL RPCs
191     serviced concurrently to 1. Its value can either be:
192     - 0 if no RPC is requested.
193     - The address to a client-bound variable that shall be set to 1 upon
194     completion of the RPC otherwise.
195     */
196     volatile uint64_t *volatile *null_rpc_global_sv;
197     /*
198     Local synchronization variable set to a >= 0 value by the server to
199     signal the completion of a RPC.
200     */
201     volatile uint64_t *null_rpc_local_sv;
202     /* 'Lock acquisitions' mode synchronization variables ===== */
203     /*
204     Global lock repeatedly acquired by all clients.
205     */
206     volatile uint64_t *lock_acquisitions_global_sv;
207     volatile uint64_t *context_variables_global_memory_area;
208     volatile uint64_t *shared_variables_memory_area;
209 } thread_arguments_block_t;
210
211 /* Mutexes and conditions used by non-blocking locks ===== */
212 #ifdef USE_BLOCKING_LOCKS
213 pthread_mutex_t mutex_rpc_done_addr_not_null;
214 pthread_cond_t cond_rpc_done_addr_not_null;
215 pthread_mutex_t mutex_rpc_done_positive;
216 pthread_cond_t cond_rpc_done_positive;
217 pthread_mutex_t mutex_rpc_done_addr_null;
218 pthread_cond_t cond_rpc_done_addr_null;
219 #endif

```



```

223
224 /* ##### */
225 /* Prototypes */
226 void *client_main(void *thread_arguments_block);
227 void *server_main(void *thread_arguments_block);
228
229
230 void *alloc(int size);
231 void get_cpu_info();
232
233 void error(char *reason);
234 void wrong_parameters_error(char *application_name);
235
236 /* ##### */
237 /* Functions */
238 /* ##### */
239 /* Main function */
240 int main(int argc, char **argv)
241 {
242     int i, j;
243     int command;
244     int result;
245     double *results;
246     double *average, *variance;
247     int end_output_with_a_newline = 1;
248     int compute_standard_deviation_and_variance = 0;
249     int number_of_samples = 0;
250
251     /* We use both a server and number_of_clients clients. The clients call
252      remote procedures from the server. */
253     pthread_t server;
254     pthread_t *clients;
255
256     /* This array contains the argument blocks passed to the clients. */
257     thread_arguments_block_t *svs_memory_area,
258                             *thread_argument_blocks,
259                             *context_variables_global_memory_area,
260                             *shared_variables_memory_area;
261
262     /* We get the clock speed and the core ordering. */
263     get_cpu_info();
264
265     /*
266      Command-line arguments
267
268      Critical sections
269      =====
270
271      -R
272      Perform null RPCs. This is the default behavior.
273
274      -L
275      Perform lock acquisitions instead of null RPCs. In this mode, there is
276      no server.
277
278      If both -R and -L are specified, the last one wins.
279
280      Execution mode
281      =====
282
283      -A number_of_runs
284      Return the average number of cycles, with variance and standard
285      deviation, over the given number of runs. This is the default mode.
286
287      -S number_of_iterations_per_sample
288      Return sampled results over a single run, using the given number of
289      iterations per sample.
290
291      -O
292      In this mode, each client enters its critical section only once, and the
293      order in which the clients were serviced is returned. Only one run is
294      performed and the results are not averaged.
295
296      If several execution modes are specified, the last one wins.
297
298      Execution settings
299      =====
300
301      -s core
302      Core on which the server thread runs if the critical section consists of
303      performing null RPCs, or core of the first client if the critical section
304      consists of lock acquisitions.
305
306      -c number_of_clients
307      Number of clients.
308
309      -n number_of_iterations_per_client
310      Number of RPCs/lock acquisitions per client.
311
312      -d delay
313      Number of cycles wasted (busy waiting) by clients after a RPC is
314      serviced. If 0, no cycles are wasted (apart from the execution time of a
315      comparison), otherwise the accuracy of the delay is within 50-150 cycles.
316
317      -r
318      Read-only wait mode : wait using comparisons, and only perform a CAS
319      after a successful comparison (test-and-test-and-set).
320
321      -C number_of_context_variables
322      Number of context (local) variables.
323
324      -S number_of_shared_variables
325      Number of shared (global) variables.
326
327      Measurements
328      =====
329
330      -e event_type
331      Count events of type event_type instead of cycles. Event_type is the type
332      of a PAPI event.
333
334      -m
335      Count cycles (or events) on the clients instead of the server. This is
336
337
338     always the case if -L is on.
339
340     -y
341     Returns results in cycles per iterations instead of the number of
342     iterations per second.
343
344     -t
345     Returns results in total cycles instead of the number of iterations per
346     second. Only the maximum value is returned.
347
348     Output
349     =====
350
351     -v
352     Return the standard deviation and variance if applicable.
353
354     -i
355     End the results with a newline. This is left as an option, to allow for
356     more flexibility with the CSV results.
357 */
358
359     opterr = 0;
360
361     /* Default values */
362     critical_sections_type = NULL_RPCS;
363
364     execution_mode = MULTIPLE_RUNS_AVERAGED;
365     number_of_runs = DEFAULT_NUMBER_OF_RUNS;
366     number_of_iterations_per_sample = DEFAULT_NUMBER_OF_ITERATIONS_PER_SAMPLE;
367
368     server_core = DEFAULT_SERVER_CORE;
369     number_of_clients = -1;
370     number_of_iterations_per_client = DEFAULT_NUMBER_OF_ITERATIONS_PER_CLIENT;
371     delay = 0;
372     read_only_wait = 0;
373
374     number_of_context_variables = DEFAULT_NUMBER_OF_CONTEXT_VARIABLES;
375     number_of_shared_variables = DEFAULT_NUMBER_OF_SHARED_VARIABLES;
376
377     measurement_type = NUMBER_OF_CYCLES;
378     measurement_location = SERVER;
379     measurement_unit = THROUGHPUT;
380
381     while ((command = getopt(argc, argv, "RLA:SOs:c:n:d:r:l:g:tTe:myvi")) != -1)
382     {
383         switch (command)
384         {
385             /*
386              Critical sections
387              */
388             case 'R':
389                 /* Actually useless (default value). */
390                 critical_sections_type = NULL_RPCS;
391                 break;
392             case 'L':
393                 critical_sections_type = LOCK_ACQUISITIONS;
394                 break;
395
396             /*
397              Execution mode
398              */
399             case 'A':
400                 execution_mode = MULTIPLE_RUNS_AVERAGED;
401                 number_of_runs = atoi(optarg);
402                 break;
403             case 'S':
404                 execution_mode = SINGLE_RUN_SAMPLED;
405                 number_of_iterations_per_sample = atoi(optarg);
406                 break;
407             case 'O':
408                 execution_mode = SINGLE_RUN_ORDERED;
409                 break;
410
411             /*
412              Execution settings
413              */
414             case 's':
415                 server_core = atoi(optarg);
416                 break;
417             case 'c':
418                 number_of_clients = atoi(optarg);
419                 break;
420             case 'n':
421                 number_of_iterations_per_client = atoi(optarg);
422                 break;
423             case 'd':
424                 delay = atoi(optarg);
425                 break;
426             case 'r':
427                 read_only_wait = 1;
428                 break;
429             case 'C':
430                 number_of_context_variables = atoi(optarg);
431                 break;
432             case 'S':
433                 number_of_shared_variables = atoi(optarg);
434                 break;
435
436             /*
437              Measurements
438              */
439             case 'e':
440                 measurement_type = NUMBER_OF_EVENTS;
441                 monitored_event_id = strtoul(optarg, NULL, 16);
442                 break;
443             case 'm':
444

```

```

453     measurement_location = CLIENTS;
454     break;
455
456     case 'y':
457         measurement_unit = CYCLES_PER_ITERATION;
458         break;
459
460     case 't':
461         measurement_unit = TOTAL_CYCLES_MAX;
462         break;
463
464     /*
465     Output
466     */
467     case 'v':
468         compute_standard_deviation_and_variance = 1;
469         break;
470
471     case 'i':
472         end_output_with_a_newline = 0;
473         break;
474
475     /*
476     Other
477     */
478     case '?': wrong_parameters_error(argv[0]);
479     default : error("getopt");
480 }
481 }
482
483 /* Default values for the number of clients */
484 if (number_of_clients <= 0)
485 {
486     number_of_clients = number_of_cores;
487
488     if (critical_sections_type == NULL_RPC)
489     {
490         number_of_clients--;
491     }
492 }
493
494 if (critical_sections_type == LOCK_ACQUISITIONS)
495 {
496     /* Lock acquisitions are only compatible with the MULTIPLE_RUNS_AVERAGED
497     execution mode. */
498     if (execution_mode != MULTIPLE_RUNS_AVERAGED)
499         error("configuration not supported");
500
501     /* In this mode, measurements take place on the clients. */
502     measurement_location = CLIENTS;
503 }
504
505 if (execution_mode == SINGLE_RUN_SAMPLED)
506 {
507     /* For now, the SINGLE_RUN_SAMPLED execution mode is incompatible with
508     blocking locks. */
509     #ifndef USE_BLOCKING_LOCKS
510     error("configuration not supported");
511     #endif
512
513     /* For now, in the SINGLE_RUN_SAMPLED mode, events can't be monitored.
514     Also, all measurements must take place on the server. */
515     if (measurement_type != NUMBER_OF_CYCLES
516         && measurement_location != SERVER)
517         error("configuration not supported");
518
519     /* Single run */
520     number_of_runs = 1;
521
522     if (number_of_iterations_per_sample >
523         number_of_iterations_per_client * number_of_clients)
524         error("too many iterations per sample");
525 }
526
527 /* In the SINGLE_RUN_ORDERED execution mode... */
528 if (execution_mode == SINGLE_RUN_ORDERED)
529 {
530     /* ...there's just one iteration per client... */
531     number_of_iterations_per_client = 1;
532
533     /* ...and a single run. */
534     number_of_runs = 1;
535 }
536
537 /* We need the maximum priority (for performance measurements). */
538 if (setpriority(PRIO_PROCESS, 0, -20) < 0)
539     error("setpriority");
540
541 /* ===== */
542 /* [v] Memory allocations
543 /* ===== */
544 /* The following arrays' sizes depend on the parameters, we need to allocate
545 them dynamically. */
546 clients = alloc(number_of_clients * sizeof(pthread_t));
547 thread_argument_blocks = alloc((number_of_clients + 1) *
548     sizeof(thread_arguments_block_t));
549 ready = alloc((number_of_clients + 1) * sizeof(int));
550 iteration_result = alloc((number_of_clients + 1) * sizeof(double));
551 results = alloc((number_of_clients + 1) *
552     (number_of_runs * sizeof(double)));
553 average = alloc((number_of_clients + 1) * sizeof(double));
554 variance = alloc((number_of_clients + 1) * sizeof(double));
555
556 /* If we need to return the number in which RPCs are serviced... */
557 if (execution_mode == SINGLE_RUN_ORDERED)
558 {
559     /* ...we allocate the 'order' array dynamically. */
560     order = alloc(number_of_clients * sizeof(double));
561 }
562
563 /* We always allocate the sampling-related arrays. */
564 rpc_done_addresses = alloc((number_of_clients + 1) * sizeof(void *));
565 number_of_samples = number_of_iterations_per_client * number_of_clients
566     / number_of_iterations_per_sample;
567
568 multiple_samples_results = alloc(number_of_samples * sizeof(double));
569
570 multiple_samples_rpc_done_addrs =
571     alloc(number_of_samples * sizeof(void *));
572
573 /* Memory area containing the synchronization variables. Each synchroniza-
574 tion variable is allocated on its own pair of cache lines. */
575 result = posix_memalign((void **)&svs_memory_area,
576     128,
577     (number_of_clients + 1) * 2 * sizeof(cache_line_t));
578
579 if (result < 0 || svs_memory_area == NULL)
580     error("memalign");
581
582 result = posix_memalign((void **)&context_variables_global_memory_area,
583     128,
584     (number_of_clients + 1) *
585     (number_of_context_variables + 1) *
586     2 * sizeof(cache_line_t));
587
588 if (result < 0 || svs_memory_area == NULL)
589     error("memalign");
590
591 result = posix_memalign((void **)&shared_variables_memory_area,
592     128,
593     (number_of_shared_variables) *
594     2 * sizeof(cache_line_t));
595
596 if (result < 0 || svs_memory_area == NULL)
597     error("memalign");
598
599 /* ===== */
600 /* [*] Memory allocations
601 /* ===== */
602
603 /* We initialize PAPI. */
604 if (PAPI_library_init(PAPI_VER_CURRENT) != PAPI_VER_CURRENT)
605     error("PAPI_library_init");
606
607 /* If we use blocking locks, then some mutexes and conditions should be
608 initialized. */
609 #ifndef USE_BLOCKING_LOCKS
610 pthread_mutex_init(&mutex_rpc_done_addr_not_null, NULL);
611 pthread_mutex_init(&mutex_rpc_done_positive, NULL);
612 pthread_mutex_init(&mutex_rpc_done_addr_null, NULL);
613
614 if (pthread_cond_init(&cond_rpc_done_addr_not_null, NULL) < 0
615     || pthread_cond_init(&cond_rpc_done_positive, NULL) < 0
616     || pthread_cond_init(&cond_rpc_done_addr_null, NULL) < 0)
617     error("pthread_cond_init");
618 #endif
619
620 /* We iterate to average the results. */
621 for (i = 0; i < number_of_runs; i++)
622 {
623     /* No rpc started for now. */
624     *((uint64_t **)&svs_memory_area) = 0;
625
626     /* None of the clients are ready. */
627     for (j = 0; j <= number_of_clients; j++)
628         ready[j] = 0;
629
630     /* Initialisations related to sampling. */
631     if (execution_mode == SINGLE_RUN_SAMPLED)
632     {
633         for (j = 0; j <= number_of_samples; j++)
634             multiple_samples_results[j] = 0;
635
636         for (j = 0; j <= number_of_samples; j++)
637             multiple_samples_rpc_done_addrs[j] = 0;
638     }
639
640     /* If we perform null RPCs in the critical sections... */
641     if (critical_sections_type == NULL_RPC)
642     {
643         thread_argument_blocks[0].id = 0;
644         thread_argument_blocks[0].null_rpc_global_sv =
645             (volatile uint64_t * volatile *)(&svs_memory_area[0]);
646         thread_argument_blocks[0].null_rpc_local_sv = NULL;
647
648         /* FIXME: global variable? */
649         thread_argument_blocks[0].context_variables_global_memory_area =
650             (volatile uint64_t *)(&context_variables_global_memory_area);
651         thread_argument_blocks[0].shared_variables_memory_area =
652             (volatile uint64_t *)(&shared_variables_memory_area);
653
654         /* ...we create a server thread. */
655         if (pthread_create(&server, NULL, server_main,
656             &thread_argument_blocks[0]) < 0)
657             error("pthread_create");
658     }
659
660     /* We create the client threads. */
661     for (j = 0; j < number_of_clients; j++)
662     {
663         thread_argument_blocks[j + 1].id = j;
664         thread_argument_blocks[j + 1].null_rpc_global_sv =
665             (volatile uint64_t * volatile *)(&svs_memory_area[0]);
666         thread_argument_blocks[j + 1].null_rpc_local_sv =
667             (volatile uint64_t *)
668             ((uint64_t)context_variables_global_memory_area +
669              j * (number_of_context_variables + 1) *
670              2 * sizeof(cache_line_t));
671
672         thread_argument_blocks[j + 1].lock_acquisitions_global_sv =
673             (volatile uint64_t *)(&svs_memory_area[0]);
674
675         thread_argument_blocks[j + 1].context_variables_global_memory_area =
676             (volatile uint64_t *)
677             ((uint64_t)context_variables_global_memory_area);
678         thread_argument_blocks[j + 1].shared_variables_memory_area =
679             (volatile uint64_t *)(&shared_variables_memory_area);
680
681         /* We send i directly instead of sending a pointer to avoid having
682

```

```

683     to create a variable per thread containing the value. */
684     if (pthread_create(&clients[j], NULL, client_main,
685                     &thread_argument_blocks[j + 1]) < 0)
686         error("pthread_create");
687 }
688
689 /* If there is a server thread... */
690 if (critical_sections_type == NULL_RPC)
691 {
692     /* ...we wait for it to finish its work. */
693     if (pthread_join(server, NULL) < 0)
694         error("pthread_join");
695 }
696
697 /* We wait for the clients to finish their work. */
698 for (j = 0; j < number_of_clients; j++)
699 {
700     if (pthread_join(clients[j], NULL) < 0)
701         error("pthread_join");
702 }
703
704 /* We save the results. */
705 for (j = 0; j <= number_of_clients; j++)
706 {
707     if (measurement_unit == THROUGHPUT)
708         iteration_result[j] = TO_THROUGHPUT(iteration_result[j]);
709
710     results[j * number_of_runs + i] = iteration_result[j];
711 }
712 }
713
714 /* In the ordered mode, we print the order in which the RPCs were
715 serviced. */
716 if (execution_mode == SINGLE_RUN_ORDERED)
717 {
718     for (i = 0; i < number_of_clients; i++)
719     {
720         j = 0;
721
722         while (order[j] != i)
723             j++;
724
725         printf("%d", j);
726     }
727
728     printf(",");
729 }
730
731 /* If sampling is on, we print the result of each sample followed by the ID
732 of the last core whose RPC has been serviced. */
733 if (execution_mode == SINGLE_RUN_SAMPLED)
734 {
735     for (i = 0; i < number_of_samples; i++)
736     {
737         for (j = 0; j < number_of_clients + 1; j++)
738         {
739             if (multiple_samples_rpc_done_addr[i] == rpc_done_addresses[j])
740                 break;
741         }
742
743         /* We convert the result to the right unit if needed. */
744         if (measurement_unit == THROUGHPUT)
745             multiple_samples_results[i] =
746                 TO_THROUGHPUT(multiple_samples_results[i]);
747
748         /* We return the result. */
749         printf("%f,%d,\n", multiple_samples_results[i], j);
750     }
751 }
752
753 if (measurement_unit == TOTAL_CYCLES_MAX)
754 {
755     for (i = 0; i < number_of_runs; i++)
756     {
757         double max = 0;
758
759         for (j = 1; j <= number_of_clients; j++)
760         {
761             if (max < results[j * number_of_runs + i])
762                 max = results[j * number_of_runs + i];
763         }
764
765         results[i] = TO_THROUGHPUT(max / (number_of_iterations_per_client *
766                                         number_of_clients));
767     }
768 }
769
770 /* We compute the average and the variance (if needed). */
771 for (j = 0; j <= number_of_clients; j++)
772 {
773     /* Initialization */
774     for (i = 0; i < number_of_runs; i++)
775     {
776         average[j] = 0;
777         variance[j] = 0;
778     }
779
780     /* We compute the average. */
781     for (i = 0; i < number_of_runs; i++)
782     {
783         average[j] += results[j * number_of_runs + i];
784     }
785
786     average[j] /= number_of_runs;
787
788     /* We compute the variance (only if needed). */
789     if (compute_standard_deviation_and_variance)
790     {
791         for (i = 0; i < number_of_runs; i++)
792         {
793             variance[j] +=
794                 (results[j * number_of_runs + i] - average[j]) *
795                 (results[j * number_of_runs + i] - average[j]);
796         }
797
798         variance[j] /= number_of_runs;
799     }
800 }
801
802 /* We print the result. If the measurements took place on the server... */
803 if (measurement_location == SERVER || measurement_unit == TOTAL_CYCLES_MAX)
804 {
805     /* ...we return the average... */
806     printf("%f", average[0]);
807
808     if (compute_standard_deviation_and_variance)
809         /* ...and, if needed, the variance and standard deviation. */
810         printf("%f,%f", variance[0], sqrt(variance[0]));
811 }
812
813 /* Otherwise... */
814 else
815 {
816     /* ...it's the same, but for each client. */
817     for (j = 1; j <= number_of_clients; j++)
818     {
819         printf("%f", average[j]);
820
821         if (compute_standard_deviation_and_variance)
822             /* ...and, if needed, the variance and standard deviation. */
823             printf("%f,%f", variance[j], sqrt(variance[j]));
824     }
825 }
826
827 /* If sampling is enabled, we need to specify a value for the core ID of the
828 average. We choose -1. */
829 if (execution_mode == SINGLE_RUN_SAMPLED)
830     printf("%d", -1);
831
832 /* Should we end the input with a newline? */
833 if (end_output_with_a_newline)
834     printf("\n");
835
836 else
837     /* If we don't, we make sure to flush the standard output. */
838     fflush(NULL);
839
840 /* ===== */
841 /* [V] Cleanup (unnecessary) */
842 /* ===== */
843 free((int *)physical_to_virtual_core_id);
844 free(clients);
845 free(thread_argument_blocks);
846 free((int *)ready);
847 free((double *)iteration_result);
848 free(results);
849 free(average);
850 free(variance);
851 if (execution_mode == SINGLE_RUN_ORDERED) free((int *)order);
852 if (execution_mode == SINGLE_RUN_SAMPLED)
853 {
854     free(rpc_done_addresses);
855     free((float *)multiple_samples_results);
856     free((void **)multiple_samples_rpc_done_addr);
857 }
858
859 #ifdef USE_BLOCKING_LOCKS
860 if (pthread_mutex_destroy(&mutex_rpc_done_addr_not_null) < 0
861     || pthread_mutex_destroy(&mutex_rpc_done_positive) < 0
862     || pthread_mutex_destroy(&mutex_rpc_done_addr_null) < 0)
863     ERROR("pthread_mutex_destroy");
864
865 if (pthread_cond_destroy(&cond_rpc_done_addr_not_null) < 0
866     || pthread_cond_destroy(&cond_rpc_done_positive) < 0
867     || pthread_cond_destroy(&cond_rpc_done_addr_null) < 0)
868     ERROR("pthread_cond_destroy");
869 #endif
870
871 /* ===== */
872 /* [*] Cleanup (unnecessary) */
873 /* ===== */
874
875 /* Everything went as expected. */
876 return EXIT_SUCCESS;
877 }
878
879 /* Function executed by the clients */
880 void *client_main(void *thread_arguments_block_pointer)
881 {
882     int i, k;
883
884     /* This is needed because the CAS function requires an address. */
885     const int one = 1;
886
887     /* We avoid using global variables in (and around) the critical section
888 to limit data cache misses. */
889     const int local_number_of_iterations_per_client =
890         number_of_iterations_per_client;
891     const int local_delay = delay;
892
893     const int local_number_of_context_variables = number_of_context_variables;
894     const int local_number_of_shared_variables = number_of_shared_variables;
895
896     /* Argument block */
897     thread_arguments_block_t thread_arguments_block;
898
899     /* We copy the contents of the argument block into these variables to
900 improve the readability of the code. */
901     volatile uint64_t *volatile *null_rpc_global_sv;
902     volatile uint64_t *null_rpc_local_sv;
903     volatile uint64_t *lock_acquisitions_global_sv;
904
905     volatile uint64_t *context_variables_local_memory_area,
906         *shared_variables_memory_area;
907
908     cpu_set_t cpuset;
909     int client_id, client_core;
910
911     /* PAPI variables */
922     int event_set = PAPI_NULL;
923     int events[1];
924     long long values[1];

```

```

913 long long start_cycles = 0, end_cycles;
914 long long cycles;
915
916 /* We get the client id. */
917 thread_arguments_block =
918 *(thread_arguments_block_t *)thread_arguments_block_pointer);
919
920 client_id = thread_arguments_block.id;
921 null_rpc_local_sv = thread_arguments_block.null_rpc_local_sv;
922 null_rpc_global_sv = thread_arguments_block.null_rpc_global_sv;
923 lock_acquisitions_global_sv =
924 thread_arguments_block.lock_acquisitions_global_sv;
925
926 context_variables_local_memory_area =
927 (volatile uint64_t *)
928 ((uint64_t)thread_arguments_block.context_variables_global_memory_area +
929 2 * sizeof(cache_line_t) * (local_number_of_context_variables + 1) *
930 client_id + 2 * sizeof(cache_line_t));
931 shared_variables_memory_area =
932 thread_arguments_block.shared_variables_memory_area;
933
934 /* We compute a proper core for the client to execute on. */
935 if (critical_sections_type == NULL_RPCS)
936 {
937     client_core = (server_core + 1 + (client_id % (number_of_cores - 1))) %
938                 number_of_cores;
939 }
940 else
941 {
942     client_core = (server_core + client_id) % number_of_cores;
943 }
944
945 /* We pin this thread to the right core. */
946 CPU_ZERO(&cpuset);
947 CPU_SET(physical_to_virtual_core_id[client_core], &cpuset);
948
949 /* We set the thread's affinity. */
950 if (pthread_setaffinity_np(pthread_self(), sizeof(cpu_set_t), &cpuset) < 0)
951     error("pthread_setaffinity_np");
952
953 /* If sampling is on... */
954 if (execution_mode == SINGLE_RUN_SAMPLED)
955 {
956     /* ...we register the address of our rpc_done local variable. */
957     rpc_done_addresses[client_core] = null_rpc_local_sv;
958 }
959
960 if (measurement_location == CLIENTS)
961 {
962     /* We will use PAPI in this thread. */
963     PAPI_thread_init((unsigned long (*)(void))pthread_self);
964
965     /* Additional initialization is needed if we are counting events. */
966     if (measurement_type == NUMBER_OF_EVENTS)
967     {
968         events[0] = monitored_event_id;
969
970         if (PAPI_create_eventset(&event_set) != PAPI_OK)
971             error("PAPI_create_eventset");
972         if (PAPI_add_events(event_set, events, 1) != PAPI_OK)
973             error("PAPI_add_events");
974
975         /* This seemingly helps increasing PAPI's accuracy. */
976         if (PAPI_start(event_set) != PAPI_OK)
977             error("PAPI_start");
978         if (PAPI_stop(event_set, values) != PAPI_OK)
979             error("PAPI_stop");
980     }
981 }
982
983 /* We're ready. */
984 ready[client_id + 1] = 1;
985
986 /* Synchronization barriers */
987 if (critical_sections_type == LOCK_ACQUISITIONS)
988 {
989     for (i = 1; i <= number_of_clients; i++)
990         while (!ready[i])
991             PAUSE;
992 }
993 if (critical_sections_type == NULL_RPCS && measurement_location == CLIENTS)
994 {
995     for (i = 0; i <= number_of_clients; i++)
996         while (!ready[i])
997             PAUSE;
998 }
999
1000 if (measurement_location == CLIENTS)
1001 {
1002     /* Does the user wish to measure the number of elapsed cycles? */
1003     if (measurement_type == NUMBER_OF_CYCLES)
1004     {
1005         /* If so, get the current cycle count. */
1006         start_cycles = PAPI_get_real_cyc();
1007     }
1008     else
1009     {
1010         /* Otherwise we start counting events. */
1011         if (PAPI_start(event_set) != PAPI_OK)
1012             error("PAPI_start");
1013     }
1014 }
1015
1016 /* ##### */
1017 /* # Main loop (client) # */
1018 /* ##### */
1019 /* First implementation : using spinlocks. */
1020 #ifndef USE_BLOCKING_LOCKS
1021 /* NULL RPCS */
1022 if (critical_sections_type == NULL_RPCS)
1023 {
1024     /* We execute number_of_iterations_per_client RPCs. */
1025     for (i = 0; i < local_number_of_iterations_per_client; i++)
1026     {
1027         *null_rpc_local_sv = -1;
1028
1029         /* We access context variables */
1030         for (k = 0; k < local_number_of_context_variables; k++)
1031             *((volatile uint64_t *)
1032              ((uint64_t)context_variables_local_memory_area +
1033               k * 2 * sizeof(cache_line_t)))+;
1034
1035         if (!read_only_wait)
1036         {
1037             /* We use a compare and swap to avoid mutexes. */
1038             while (!__sync_bool_compare_and_swap(&null_rpc_global_sv,
1039                                                  0,
1040                                                  null_rpc_local_sv))
1041                 PAUSE;
1042         }
1043         else
1044         {
1045             for (;;)
1046             {
1047                 if (!(*null_rpc_global_sv))
1048                 {
1049                     if (!__sync_bool_compare_and_swap(&null_rpc_global_sv,
1050                                                        0,
1051                                                        null_rpc_local_sv))
1052                         continue;
1053                     else break;
1054                 }
1055                 PAUSE;
1056             }
1057         }
1058         while (*null_rpc_local_sv < 0)
1059             PAUSE;
1060
1061         /* We could avoid the delay introduced by the test by moving the
1062          * if outside of the loop, however, tests show that no delay or
1063          * a small delay doesn't alter the results significantly. */
1064         if (local_delay > 0)
1065         {
1066             /* Delay */
1067             cycles = PAPI_get_real_cyc();
1068             while ((PAPI_get_real_cyc() - cycles) < local_delay)
1069                 ;
1070         }
1071     }
1072 }
1073
1074 /* Lock acquisitions */
1075 else
1076 {
1077     if (execution_mode != SINGLE_RUN_SAMPLED)
1078     {
1079         for (i = 0; i < local_number_of_iterations_per_client; i++)
1080         {
1081             /* We access context variables */
1082             for (k = 0; k < local_number_of_context_variables; k++)
1083                 *((volatile uint64_t *)
1084                  ((uint64_t)context_variables_local_memory_area +
1085                   k * 2 * sizeof(cache_line_t)))+;
1086
1087             if (!read_only_wait)
1088             {
1089                 while (!__sync_bool_compare_and_swap(
1090                     &lock_acquisitions_global_sv, 0, &one))
1091                     PAUSE;
1092             }
1093             else
1094             {
1095                 for (;;)
1096                 {
1097                     if (!*null_rpc_global_sv)
1098                     {
1099                         if (!__sync_bool_compare_and_swap(
1100                             &lock_acquisitions_global_sv, 0, &one))
1101                             continue;
1102                         else break;
1103                     }
1104                     PAUSE;
1105                 }
1106             }
1107
1108             /* We access context variables */
1109             for (k = 0; k < local_number_of_context_variables; k++)
1110                 *((volatile uint64_t *)
1111                  ((uint64_t)context_variables_local_memory_area +
1112                   k * 2 * sizeof(cache_line_t)))+;
1113
1114             /* We access shared variables */
1115             for (k = 0; k < local_number_of_shared_variables; k++)
1116                 *((volatile uint64_t *)
1117                  ((uint64_t)shared_variables_memory_area +
1118                   k * 2 * sizeof(cache_line_t)))+;
1119
1120             *lock_acquisitions_global_sv = 0;
1121
1122             if (local_delay > 0)
1123             {
1124                 /* Delay */
1125                 cycles = PAPI_get_real_cyc();
1126                 while ((PAPI_get_real_cyc() - cycles) < local_delay)
1127                     ;
1128             }
1129         }
1130     }
1131 }
1132
1133 /* Second implementation : using blocking locks. */
1134 #else
1135 /* We execute number_of_iterations_per_client RPCs. */
1136 for (i = 0; i < local_number_of_iterations_per_client; i++)
1137 {
1138     *null_rpc_local_sv = -1;
1139
1140     /* This lock corresponds to the CAS from the previous implementation. */

```

```

1143 pthread_mutex_lock(&mutex_rpc_done_addr_null);
1144
1145 /* We wait for the server to be ready. */
1146 while(global->rpc_done_addr != 0)
1147     pthread_cond_wait(&cond_rpc_done_addr_null,
1148                     &mutex_rpc_done_addr_null);
1149
1150 /* The server is ready, let's execute a RPC. */
1151 pthread_mutex_lock(&mutex_rpc_done_addr_not_null);
1152 *null_rpc_global_sv = null_rpc_local_sv;
1153 pthread_cond_signal(&cond_rpc_done_addr_not_null);
1154 pthread_mutex_unlock(&mutex_rpc_done_addr_not_null);
1155
1156 pthread_mutex_unlock(&mutex_rpc_done_addr_null);
1157
1158 /* We wait until the server is done with our RPC. */
1159 pthread_mutex_lock(&mutex_rpc_done_positive);
1160 while (*null_rpc_local_sv < 0)
1161     pthread_cond_wait(&cond_rpc_done_positive,
1162                     &mutex_rpc_done_positive);
1163 pthread_mutex_unlock(&mutex_rpc_done_positive);
1164 }
1165 #endif
1166 /* ##### */
1167 /* # End # */
1168 /* ##### */
1169
1170 if (measurement_location == CLIENTS)
1171 {
1172     /* Are we counting cycles? */
1173     if (measurement_type == NUMBER_OF_CYCLES)
1174     {
1175         /* If so, get the current cycle count. */
1176         end_cycles = PAPI_get_real_cyc();
1177
1178         if (measurement_unit != TOTAL_CYCLES_MAX)
1179         {
1180             /* We return the number of cycles per RPC. */
1181             iteration_result[client_id + 1] =
1182                 (double)(end_cycles - start_cycles) /
1183                 number_of_iterations_per_client;
1184         }
1185         else
1186         {
1187             iteration_result[client_id + 1] =
1188                 (double)(end_cycles - start_cycles);
1189         }
1190     }
1191     else
1192     {
1193         /* Otherwise, we were counting events. We read the number of
1194         events. */
1195         if (PAPI_stop(event_set, values) != PAPI_OK)
1196             error("PAPI_stop");
1197
1198         /* We return the number of events. */
1199         iteration_result[client_id + 1] =
1200             (double)values[0] / number_of_iterations_per_client;
1201     }
1202 }
1203
1204 /* In the ordered mode, we save the order in which RPCs are serviced. */
1205 if (execution_mode == SINGLE_RUN_ORDERED)
1206 {
1207     order[client_id] = *null_rpc_local_sv;
1208 }
1209
1210 for (k = 1; k <= local_number_of_context_variables; k++)
1211     printf(" (%(volatile uint64_t *)
1212           ((uint64_t)context_variables_local_memory_area +
1213            k * 2 * sizeof(cache_line_t))));
1214
1215 PAPI_unregister_thread();
1216
1217 return NULL;
1218 }
1219
1220 /* Function executed by the server */
1221 void *server_main(void *thread_arguments_block_pointer)
1222 {
1223     int i, j, k;
1224     /* We avoid using global variables in the main loop to limit data cache
1225     misses. */
1226     const int number_of_iterations =
1227         number_of_iterations_per_client * number_of_clients;
1228
1229     /* Sampling-related variables. */
1230     const int number_of_samples = number_of_iterations
1231         / number_of_iterations_per_sample;
1232     const int number_of_iterations_per_sample_m1 =
1233         number_of_iterations_per_sample - 1;
1234     int sample_start_cycles, sample_end_cycles;
1235
1236     const int local_number_of_context_variables = number_of_context_variables;
1237     const int local_number_of_shared_variables = number_of_shared_variables;
1238
1239     thread_arguments_block_t thread_arguments_block;
1240     volatile uint64_t *volatile *null_rpc_global_sv;
1241
1242     volatile uint64_t *context_variables_global_memory_area,
1243         *shared_variables_memory_area;
1244
1245     /* This variable is used to pin the thread to the right core. */
1246     cpu_set_t cpuset;
1247
1248     /* PAPI variables */
1249     int event_set = PAPI_NULL;
1250     int events[1];
1251     long long values[1];
1252     long long start_cycles = 0, end_cycles;
1253
1254     /* We pin this thread to the right core. */
1255     CPU_ZERO(&cpuset);
1256     CPU_SET(physical_to_virtual_core_id[server_core], &cpuset);
1257
1258     thread_arguments_block =
1259         *((thread_arguments_block_t *)thread_arguments_block_pointer);
1260     null_rpc_global_sv = thread_arguments_block.null_rpc_global_sv;
1261
1262     context_variables_global_memory_area =
1263         thread_arguments_block.context_variables_global_memory_area;
1264     shared_variables_memory_area =
1265         thread_arguments_block.shared_variables_memory_area;
1266
1267     /* We set the thread's affinity. */
1268     if (pthread_setaffinity_np(pthread_self(), sizeof(cpu_set_t), &cpuset) < 0)
1269         error("pthread_setaffinity_np");
1270
1271     /* Are we monitoring cycles/events on the server? */
1272     if (measurement_location == SERVER)
1273     {
1274         /* If so, we will use PAPI in this thread. */
1275         PAPI_thread_init((unsigned long (*)(void))pthread_self);
1276
1277         /* If we are counting events... */
1278         if (measurement_type == NUMBER_OF_EVENTS)
1279         {
1280             events[0] = monitored_event_id;
1281
1282             /* We initialize the event set. */
1283             if (PAPI_create_eventset(&event_set) != PAPI_OK)
1284                 error("PAPI_create_eventset");
1285             if (PAPI_add_events(event_set, events, 1) != PAPI_OK)
1286                 error("PAPI_add_events");
1287
1288             /* Starting and stopping PAPI once before measurements seems to
1289             improve accuracy. */
1290             if (PAPI_start(event_set) != PAPI_OK)
1291                 error("PAPI_start");
1292             if (PAPI_stop(event_set, values) != PAPI_OK)
1293                 error("PAPI_stop");
1294         }
1295     }
1296
1297     /* We are ready. */
1298     ready[0] = 1;
1299
1300     /* We want all the clients to be ready before we start servicing RPCs. */
1301     for (i = 1; i <= number_of_clients; i++)
1302     {
1303         while (!ready[i]);
1304         PAUSE;
1305     }
1306
1307     /* If the measurements take place on the server... */
1308     if (measurement_location == SERVER)
1309     {
1310         /* ...and if we're counting cycles... */
1311         if (measurement_type == NUMBER_OF_CYCLES)
1312         {
1313             /* ...we get the current cycle count. */
1314             start_cycles = PAPI_get_real_cyc();
1315         }
1316         /* Otherwise, we're counting events. */
1317         else
1318         {
1319             /* We start the event counter. */
1320             if (PAPI_start(event_set) != PAPI_OK)
1321                 error("PAPI_start");
1322         }
1323     }
1324 }
1325
1326 /* ##### */
1327 /* # Main loop (server) */
1328 /* ##### */
1329 /* First implementation : using spinlocks. */
1330 #ifndef USE_BLOCKING_LOCKS
1331 /* Sampling is off. */
1332 if (execution_mode != SINGLE_RUN_SAMPLED)
1333 {
1334     /* We start answering RPCs. */
1335     for (i = 0; i < number_of_iterations; i++)
1336     {
1337         /* No address received from a client? That means no RPC has been
1338         requested. */
1339         while (!(*null_rpc_global_sv))
1340             PAUSE;
1341
1342         if (local_number_of_context_variables > 0)
1343         {
1344             for (k = 1; k <= local_number_of_context_variables; k++)
1345                 *((volatile uint64_t *)
1346                  ((uint64_t)(*null_rpc_global_sv) +
1347                   k * 2 * sizeof(cache_line_t)))+;
1348         }
1349
1350         /* We access shared variables */
1351         for (k = 0; k < local_number_of_shared_variables; k++)
1352             *((volatile uint64_t *)
1353              ((uint64_t)shared_variables_memory_area +
1354               k * 2 * sizeof(cache_line_t)))+;
1354
1355         /* We notify the client we are done with his RPC by setting the
1356         variable it provided to i. */
1357         **null_rpc_global_sv = i;
1358
1359         /* We are done with our RPC. */
1360         *null_rpc_global_sv = 0;
1361     }
1362 }
1363
1364 /* Sampling is on. */
1365 else
1366 {
1367     /* Outer loop: one iteration per sample. */
1368     for (j = 0; j < number_of_samples; j++)
1369     {
1370         /* Same as before, except now we get cycle statistics for each
1371         sample. */
1372         sample_start_cycles = PAPI_get_real_cyc();

```

```

1373
1374     for (i = 0; i < number_of_iterations_per_sample_ml; i++)
1375     {
1376         while (!(*null_rpc_global_sv))
1377             PAUSE;
1378
1379         **null_rpc_global_sv = i;
1380         *null_rpc_global_sv = 0;
1381     }
1382
1383     while (!(*null_rpc_global_sv))
1384         PAUSE;
1385
1386     **null_rpc_global_sv = i;
1387
1388     sample_end_cycles = PAPI_get_real_cyc();
1389
1390     /* We save the results. */
1391     multiple_samples_results[j] =
1392         (double) (sample_end_cycles - sample_start_cycles)
1393         / number_of_iterations_per_sample;
1394     /* We need to know which core was serviced last. */
1395     multiple_samples_rpc_done_addrs[j] = *null_rpc_global_sv;
1396
1397     *null_rpc_global_sv = 0;
1398 }
1399
1400 /* Second implementation : using blocking locks. */
1401 #else
1402 /* We start answering RPCs. */
1403 for (i = 0; i < number_of_iterations; i++)
1404 {
1405     /* No address received from a client? That means no RPC has been
1406     requested. */
1407     pthread_mutex_lock(&mutex_rpc_done_addr_not_null);
1408     while (!(*null_rpc_global_sv))
1409         pthread_cond_wait(&cond_rpc_done_addr_not_null,
1410             &mutex_rpc_done_addr_not_null);
1411     pthread_mutex_unlock(&mutex_rpc_done_addr_not_null);
1412
1413     /* We notify the client we are done with his RPC by setting the
1414     variable it provided to i. */
1415     pthread_mutex_lock(&mutex_rpc_done_positive);
1416     **null_rpc_global_sv = i;
1417     pthread_cond_signal(&cond_rpc_done_positive);
1418     pthread_mutex_unlock(&mutex_rpc_done_positive);
1419
1420     /* We are done with our RPC. */
1421     pthread_mutex_lock(&mutex_rpc_done_addr_null);
1422     *null_rpc_global_sv = 0;
1423     pthread_cond_signal(&cond_rpc_done_addr_null);
1424     pthread_mutex_unlock(&mutex_rpc_done_addr_null);
1425 }
1426 #endif
1427 /* ##### */
1428 /* # End */
1429 /* ##### */
1430
1431     for (k = 0; k < local_number_of_shared_variables; k++)
1432         printf(" ", *((volatile uint64_t *)
1433             ((uint64_t)shared_variables_memory_area +
1434             k * 2 * sizeof(cache_line_t))));
1435
1436     /* We gather statistics if measurements are to take place on the server. */
1437     if (measurement_location == SERVER)
1438     {
1439         /* Are we counting cycles? */
1440         if (measurement_type == NUMBER_OF_CYCLES)
1441         {
1442             /* If so, get the current cycle count. */
1443             end_cycles = PAPI_get_real_cyc();
1444
1445             if (measurement_unit != TOTAL_CYCLES_MAX)
1446             {
1447                 /* We return the number of cycles per RPC. */
1448                 iteration_result[0] =
1449                     (double) (end_cycles - start_cycles) / number_of_iterations;
1450             }
1451             else
1452             {
1453                 iteration_result[0] = (double) (end_cycles - start_cycles);
1454             }
1455         }
1456         /* We're counting events. */
1457         else
1458         {
1459             /* Otherwise, we were counting events, therefore, we read the number
1460             of events. */
1461             if (PAPI_stop(event_set, values) != PAPI_OK)
1462                 error("PAPI_stop");
1463
1464             /* We return the number of events. */
1465             iteration_result[0] = (double) values[0] / number_of_iterations;
1466         }
1467     }
1468
1469     PAPI_unregister_thread();
1470
1471     return NULL;
1472 }
1473
1474 /* Wrapper for malloc that checks for errors. */
1475 void *alloc(int size)
1476 {
1477     void *result;
1478
1479     /* We just call malloc and check for errors. */
1480     if ((result = malloc(size)) == NULL)
1481         error("malloc");
1482
1483     /* We return the result. */
1484     return result;
1485 }
1486
1487 /* This function gets the CPU speed then allocates and fills the

```

```

1488     virtual_to_physical_core_id array. */
1489 void get_cpu_info()
1490 {
1491     int i = 0;
1492     FILE *file;
1493     char text[32];
1494     float previous_core_frequency = 0, core_frequency = 0;
1495     int cpu_id, local_core_id, local_number_of_cores;
1496     int virtual_to_physical_core_id[MAX_NUMBER_OF_CORES];
1497
1498     /* We use UNIX commands to parse the /proc/cpuinfo file. */
1499     file = popen("/bin/cat /proc/cpuinfo | "
1500         "/bin/egrep \"physical id|core id|cpu cores|cpu MHz\" | "
1501         "/bin/sed s/\"core id\t\t: \\|physical id\t: \\|\" "
1502         "cpu cores\t: \\|cpu MHz\t\t: \\|//", "r");
1503
1504     if (file == NULL)
1505         error("popen");
1506
1507     /* For each core... */
1508     while (fgets(text, 32, file) != NULL)
1509     {
1510         /* We get the CPU speed. */
1511         core_frequency = atof(text);
1512
1513         /* We ensure all cores use the same clock frequency.*/
1514         if (i > 0 && previous_core_frequency != core_frequency)
1515             error("all processing cores must use the same clock frequency");
1516
1517         previous_core_frequency = core_frequency;
1518
1519         /* We get the CPU ID... */
1520         fgets(text, 32, file);
1521         cpu_id = atoi(text);
1522
1523         /* ...the physical core id on this CPU... */
1524         fgets(text, 32, file);
1525         local_core_id = atoi(text);
1526
1527         /* ...the number of cores on this CPU... */
1528         fgets(text, 32, file);
1529         local_number_of_cores = atoi(text);
1530
1531         /* ...and we use this info to compute the physical core ID (global). */
1532         virtual_to_physical_core_id[i] = local_core_id +
1533             cpu_id * local_number_of_cores;
1534
1535         i++;
1536     }
1537
1538     /* We update the global cpu_frequency variable. */
1539     cpu_frequency = core_frequency;
1540
1541     /* We update the global number_of_cores variable. */
1542     number_of_cores = i;
1543
1544     /* We allocate the physical_to_virtual_core_id array. */
1545     physical_to_virtual_core_id = alloc(number_of_cores * sizeof(int));
1546
1547     /* We fill the physical_to_virtual_core_id array. */
1548     for (i = 0; i < number_of_cores; i++)
1549         physical_to_virtual_core_id[virtual_to_physical_core_id[i]] = i;
1550
1551     /* We're done with the results of the UNIX commands. */
1552     if (pclose(file) < 0)
1553         error("pclose");
1554 }
1555
1556 /* This function writes the usage string to stderr and exits the application. */
1557 void wrong_parameters_error(char *application_name)
1558 {
1559     const char *usage =
1560         "usage : %s [-R] [-L] \\\n"
1561         "           [-A number_of_runs] [-S #_iterations_per_sample] [-O]\n"
1562         "           [-s core] [-c #_clients] [-n #_iterations_per_client] "
1563         "           [-d delay] [-r]\n"
1564         "           [-e event_type] [-m] [-y] [-i]\n";
1565
1566     /* TODO: print detailed usage. */
1567
1568     fprintf(stderr, usage, application_name);
1569     exit(EXIT_FAILURE);
1570 }
1571
1572 /* Error function */
1573 void error(char *reason)
1574 {
1575     fprintf(stderr, "An error has occurred. Cause: %s.\n", reason);
1576     exit(EXIT_FAILURE);
1577 }

```