

Runtime support for OpenMP language on Cell Broadband Engine

Guansong Zhang, Raul Silvera, Mark Mendell*

IBM Toronto Lab
Toronto
ON, L6G 1C7, Canada

Abstract:

The Cell processor has many advantages: computation power, memory bandwidth, and power consumption. On the other hand, it is hard to program. The two distinct ISAs, limited SPU address space and explicit DMA communications, all complicate the programming effort on this processor. We are helping this situation by providing an OpenMP compiler for Cell. We can implement a full OpenMP compliant compiler to run potentially any OpenMP programs on such a processor with reasonable performance. In this paper, we introduce the runtime support for the OpenMP compiler and show how to map major OpenMP constructs onto a Cell processor. This paper will help programmers write Cell applications using industry standard OpenMP, and help people with an OpenMP background use the Cell processor more effectively.

Keywords: Cell, OpenMP, parallel region, workshare, barrier

1 Introduction

The Cell Broadband Engine TM (Cell BE) processor [1] is now commercially available in both the Sony PS3 game console and the IBM BladeCenter which represents the first product on its development road-map. The anticipated high volumes for this non-traditional commodity hardware continue to make it interesting in a variety of different application spaces, ranging from the obvious multi-media and gaming domain, through the HPC space (both traditional and commercial), and to the potential use of Cell as a building block for very high end supercomputing systems [2].

This first generation Cell processor provides flexibility and performance through the inclusion of a 64-bit multi-threaded Power Processor TM Element (PPE) with two levels of globally-coherent cache and support for multiple operating systems including Linux. For additional performance, a Cell processor includes eight Synergistic Processor Elements (SPEs), each consisting of a Synergistic Processing Unit (SPU), a 256K local memory, and a globally-coherent DMA engine. Computations on the SPUs are performed by 128-bit wide Single Instruction Multiple Data (SIMD) functional units. An integrated high bandwidth bus, the Element Interconnect Bus (EIB), glues together the nine processors and their ports to external memory and IO, and allows the SPUs to be used for streaming applications [3].

* The opinions expressed in this paper are those of the authors and not necessarily of IBM.

1 Data is transferred between the local memory and the DMA engine [4] in chunks
 2 of 128 bytes. The DMA engine can support up to 16 concurrent requests of up to 16K
 3 bytes originating either locally or remotely. The DMA engine is part of the globally
 4 coherent memory address space; addresses of local DMA requests are translated by a
 5 Memory Management Unit (MMU) before being sent on the bus. Bandwidth between
 6 the DMA and the EIB bus is 8 bytes per cycle in each direction. Programs interface with
 7 the DMA unit via a channel interface and may initiate blocking as well as non-blocking
 8 requests.

9 Programming the SPE processor is significantly enhanced by the availability of an
 10 optimizing compiler which supports SIMD intrinsic functions and automatic simdiza-
 11 tion [5]. However, programming the entire Cell processor consisting of the coupled
 12 PPE and 8 SPE processors is a much more complex task, requiring partitioning of an
 13 application to accommodate the limited local memory constraints of the SPE, paral-
 14 lelization across the multiple SPEs, orchestration of data transfers through insertion of
 15 DMA commands, and compiling for two distinct ISAs.

16 To help users to harness such a powerful yet complicated processor and to port
 17 existing legacy code, we developed an OpenMP[6] implementation for the system. This
 18 is the same as a typical OpenMP implementation such as in Figure 1.

19 In this figure, the language-dependent frontends will translate user source code to
 20 an *intermediate representation (IR)*, to be processed by an optimizing compiler into
 21 a resulting *node program*, which interacts with a *runtime* environment that starts and
 22 controls the multi-threaded execution.

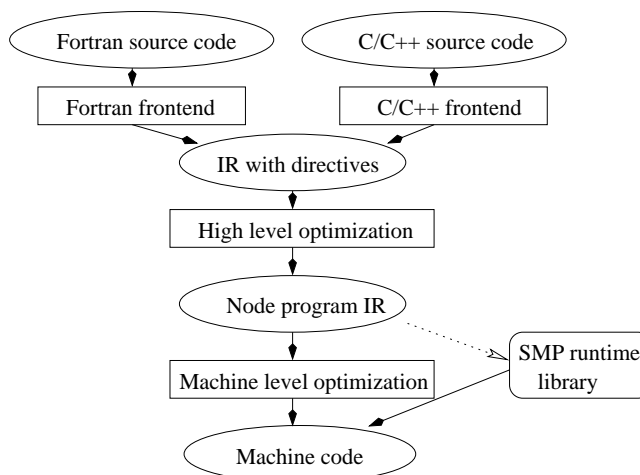


Fig. 1. Compiler and runtime structure

23 With this system, the programming the Cell processor can benefit from an industry
 24 standard in a familiar user environment. The compiler will:

- 1 – **Systematically manage multiple SPU threads** Users do not need to know Cell
2 SDK details, as thread management on the SPUs is provided through the OpenMP
3 framework.
- 4 – **Automatically overlap computation and communication** The compiler will gen-
5 erate data communication between the PPU main memory and the SPU local store.
6 DMA buffering and software cache can be automatically used to exploit the large
7 bandwidth between main memory and all 9 processors.
- 8 – **Enlarge the limited SPU address space** Our OpenMP Cell compiler will treat
9 SPU local storage more like a cache in a traditional processor. Large application
10 data and code will be kept in the main memory of the PPU and paged in or trans-
11 ferred as appropriate.

12 There are already papers introducing the second and the third points[7] of the sys-
13 tem, but little discussion has been dedicated to the OpenMP runtime environment itself.
14 The OpenMP standard continues to evolve and there are still many open issues in this
15 area, including the mechanisms required to support nested parallelism, task, and more.

16 In this paper, we will focus on the first item, the major runtime environment imple-
17 mentation itself¹.

18 The runtime system we have developed is based on the OpenMP runtime library we
19 use for other XL series compilers, including C/C++ and Fortran on AIX® and Linux
20 systems [8]. This paper shows how to map the OpenMP environment to the Cell sys-
21 tem. It concentrates on discussions of the implementation considerations, rather than
22 implementation details. Some performance numbers are also given.

23 2 Parallelism in OpenMP

24 The OpenMP standard is jointly defined by a group of major computer hardware and
25 software vendors[6]. This makes it a portable, scalable model that gives shared-memory
26 parallel programmers a simple and flexible interface for developing parallel applica-
27 tions for platforms ranging from the desktop to the supercomputer. The most powerful
28 semantics in the standard is represented by a *parallel construct*. We will first look at
29 mapping a parallel construct onto a Cell processor.

30 2.1 A parallel region

31 There are two types of processors in a Cell chip. The first type of processor element, the
32 PPE, is a 64-bit PowerPC Architecture core. It is fully compliant with the 64-bit Pow-
33 erPC Architecture and can run 32-bit and 64-bit operating systems and applications. The
34 second type of processor element, the SPE, is optimized for running compute-intensive
35 applications. The SPEs are independent processors, each running its own individual ap-
36 plication program. Each SPE has full access to coherent shared memory through DMA,
37 including the memory-mapped I/O space. The designation *synergistic* for this processor
38 was chosen carefully because there is a mutual dependence between the PPE and the
39 SPEs. The SPEs depend on the PPE to run the operating system, and, in many cases, the

¹ Topics such as threadprivate and C++ object privatization will be discussed elsewhere.

1 top-level control thread of an application. The PPE depends on the SPEs to provide the
 2 bulk of the application performance. This dependency is reflected in our OpenMP stan-
 3 dard implementation: the sequential part of the program will be executed by the PPU,
 4 and the parallel part, i.e., `parallel` regions, of the program will be executed by the
 5 SPUs. In this way, the user program can take full advantage of the variety of PPU sys-
 6 tem libraries, while the SPEs will help when the program is within the computationally
 7 intensive part.

8 This implementation scheme brings some challenges to the compiler system, in-
 9 cluding: compiling for both PPU and SPU, switching between PPU and SPU execution,
 10 separating PPU libraries from SPU ones, and more. We will cover some of these topics
 11 in this paper.

12 From a memory system point of view, the OpenMP standard specifies a relaxed-
 13 consistency, shared-memory model. This model allows each thread to have its own
 14 temporary view of memory. A value written to a variable, or a value read from a vari-
 15 able, can remain in the thread's temporary view until it is forced to share memory by
 16 an OpenMP flush operation. This memory model can be efficiently implemented on the
 17 Cell memory structure, even though each SPE has just 256K of directly accessible local
 18 memory for code, data, stack and heap. In order to handle such limited memory, we
 19 only allocate private variables accessed in SPE code in the SPE local store. Shared vari-
 20 ables reside in system memory, and SPE code can access them through DMA operations
 21 using either a DMA buffer mechanism or software cache.

22 We set up the outermost parallel region by specifying the starting thread and slave
 23 threads on the PPU and the SPUs separately. Nested parallel regions will be executed
 24 on SPUs only. We still honor OpenMP data scoping clauses, using the encountering
 25 threads to execute the nested parallel region sequentially.

26 2.2 The `if` clause

27 The `if` clause on a `parallel` region is useful when it is not known at compile time
 28 that the computation part in a `parallel` region is intense enough to outweigh the
 29 overhead of communication involved. For example, the following code will be parallel
 30 executed on SPUs only when we have a large enough loop bound.

```
31     #pragama omp parallel for if (n>LIMIT)
32     for (int i=0; i<n; i++) {
33         compute(i);
34     }
```

35 This could be rewritten as, (although it is equivalent in OpenMP 2.5, it may not be
 36 true considering nested parallel regions in OpenMP 3.0)

```
37     if (n>LIMIT) {
38         #pragama omp parallel for
39         for (int i=0; i<n; i++) {
40             compute(i);
41         }
42     } else {
```

```

1      compute(i);
2  }

```

3 to avoid overhead for starting up a parallel region. The function `compute` may be
4 executed on either the PPU or multiple SPUs. In such a case, the compiler have to *clone*
5 the procedure for two different ISAs. We need to redraw the Figure 1 as Figure 2.

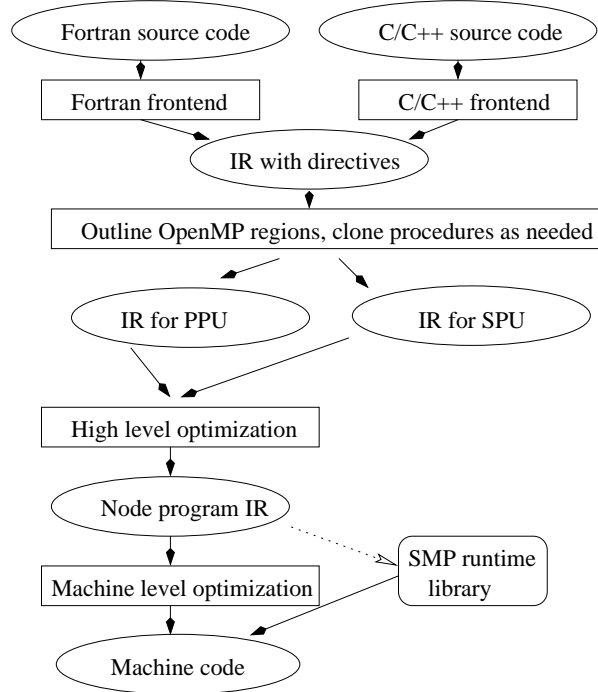


Fig. 2. Clone for two ISAs

6 After we outline the parallel sections into procedures for an OpenMP region, call
7 graph analysis is needed within the compiler to trace which functions can be reached
8 from PPU and/or SPU. We then need to *clone*, make two versions of, the routine if a
9 procedure will be reached by both PPU and SPU sides. The SPU code is generated with
10 the SPU linker, as a single SPU program executable. The PPU code is also generated
11 with the SPU program *embedded* and linked into the final PPU executable.

12 2.3 The `num_threads` clause

13 According to the OpenMP standard, the thread which first encounter a parallel construct
14 becomes the master thread of a new team, with the thread number as zero for the dura-

tion of the `parallel` region. Initially we did put the PPU part of the `parallel` region. Experiments later on have shown some drawbacks with this approach.

- The difference in computation power between a PPU and an SPU makes the PPU thread less useful in a workshare of the `parallel` region. We can not see any benchmarks benefit from assigning partial work to a PPU thread. The PPU thread may actually slow down the application, with load balancing playing a role.
- In order to solve the above problem, we have to apply special loop schedules to disable the PPU inside an OpenMP workshare. This conflicts with the basic idea in an OpenMP environment, i.e., the threads in a `parallel` region, including both the master and the slave threads, are as equal as possible.
- Even when we don't have any workshare in a `parallel` region, the PPU participation means that the computation in the `parallel` region will be carried out by both PPU and SPU, and may have different floating point results. SPU single precision arithmetic is not fully IEEE-754 compliant.
- If the PPU participates in a `parallel` region, it will force the compiler to clone all the procedures can be reached by a `parallel` region, even without an `if` clause. This will extend the compilation time and may cause certain optimizations, such as pointer analysis, to give up due to an artificially enlarged problem size.
- This will also prevent a user from providing his/her own computation kernel in an SPU only format. If `parallel` regions are SPU only, an advanced user can compile a procedure with a software-cache-aware SPU compiler and link it into the OpenMP compiler generated code, as long as this procedure can only be reached from a `parallel` region (without an `if` clause).

We decided to have the PPU thread start a `parallel` region, then suspend execution and resume after the `parallel` region. Inside the `parallel` region, the `num_threads` clause or OpenMP internal control variables will be used to determine how many SPU threads will participate, and one of them will be picked as the master thread. This is more in keeping with the upcoming OpenMP 3.0 [9] task view, where the code outside and inside a `parallel` region are considered as different tasks.

We also set up a upper limit for the number of SPU threads can be used, it will be the maximum number of SPUs available on the system. The concept is the same as the `OMP_THREAD_LIMIT` environment variable in OpenMP 3.0.

3 Sharing work through workshares

When a `parallel` region starts, multiple threads will execute the code inside the region concurrently,

In Figure 3, a master thread starts up the `parallel` region executed by 8 threads. Through `parallel` regions, multiple threads accomplish *worksharing* in an OpenMP program.

The simplest format of worksharing is replicated execution of the same code segment on different threads. It is more useful to divide work among multiple threads — either by having different threads operate on different portions of a shared data structure, or by having different threads perform entirely different tasks. Each cooperation

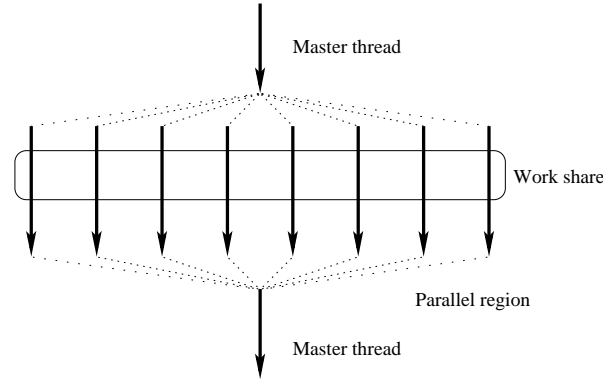


Fig. 3. Parallel region

1 of this kind is considered as a *workshare*² in this paper. In the Figure 3, we use a round-
 2 cornered rectangular block to represent one workshare.

3 The most common forms of workshares in OpenMP specification are `omp for`
 4 and `omp sections`,

5 We can also consider OpenMP `single` as a special workshare too. A `single`
 6 construct is semantically equivalent to a `sections` construct with only one `section`
 7 inside. For a `single` construct, the first thread that encounters the code will execute
 8 the block. This is different from a `master` construct, where the decision must be made
 9 by checking the thread id.

10 All the workshares introduced here have an important common feature: they have an
 11 implicit barrier at the end of the construct, which may be disabled by using the optional
 12 `nowait` clause. Besides, An explicit `barrier` region is semantically equivalent to an
 13 empty workshare without the `nowait` clause.

14 This observation reveals that the barrier is the most important workshare to in a
 15 runtime supporting system. We will discuss more about a barrier implementation in the
 16 next section.

17 Different implementations may have a different kind of classification. The real ad-
 18 vantage of considering them all as workshares is for simplifying the developing process.
 19 From an implementation point of view, the common behaviors will lead to a common
 20 code base, thus improving the overall code quality. And once we have a barrier imple-
 21 mentation on Cell, the implementation of the rest workshares are relatively easier.

² The concept of workshare here is different from the parallel construct, `WORKSHARE` in Fortran
 OpenMP specification 2.0, which can be treated semantically as the combination of `for` and
`single` constructs.

1 4 The omp barrier

2 Figure 4 shows a simple barrier that can be used for a OpenMP runtime system. Each
 3 arriving thread will decrease the counter with an atomic operation, and they will spin
 4 checking the counter value and not leave the barrier until the counter becomes zero.

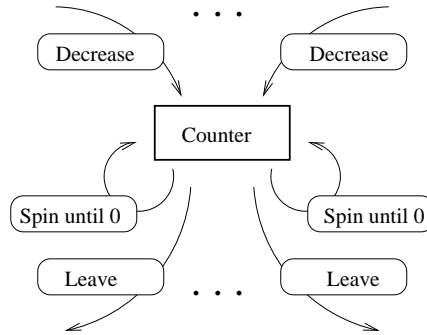


Fig. 4. Barrier implemented with fetch-and-add

5 Algorithm 1 describes a more realistic implementation used in the XL OpenMP
 6 runtime system on AIX and Linux,

7 Instead of having one counter, each thread will operate its own counter, thus the
 8 expensive fetch and add operation can be replaced with a regular add. In order to further
 9 reduce the contention, we also have a local sensor to reduce the false sharing between
 10 reads from different threads[10].

11 The algorithm can be mapped perfectly onto the Cell processor. The distributed
 12 counter and local sensor can be replaced by the SPU output and input mailbox separately.
 13 For example, the following code executed on an SPU and achieve the algorithm
 14 drawn in Figure 5.

```

15 {
16     int data;
17     spu_writetech(SPU_WrOutMbox, 1);
18     data = spu_readch(SPU_RdInMbox);
19 }
  
```

20 And the PPU version,

```

21 unsigned int data;
22 spe_context_ptr_t id = _xlbesplitSPUThreads[i].spuid;
23
24 for (int i=0; i<num_threads; i++) {
25     while(spe_out_mbox_status(id) == 0);
26     spe_out_mbox_read(id, &data, 1);
27 }
  
```

Algorithm 1: Barrier with distributed counter and local sensor

Data : Distributed counter with each element as one
Data : Local sensor with each element as one

begin
 Decrease my own distributed counter element;
 if *I am the master thread* **then**
 repeat
 | **foreach** *element in distributed counter* **do** check if it is zero
 | **until** *all distributed counter elements are zero*;
 foreach *element in distributed counter* **do**
 | set it back to one
 end
 foreach *element in local sensor* **do**
 | set it to zero
 end
 else
 repeat
 | check my local sensor element
 | **until** *it is zero*;
 end
 Set my own local sensor element back to one;
 end

```
1
2   for (int i=0; i<num_threads; i++) {
3       while (spe_in_mbox_status(id) == 0);
4       spe_in_mbox_write(id, &data, 1, SPE_MBOX_ANY_NONBLOCKING);
5   }
```

6 In the real implementation, a barrier is also a chance for PPU and SPU to exchange
7 data such as for handling single, reduction, copyprivate etc.

8 5 Benchmarks and Summary

9 In this section, we use some benchmark numbers to summarize the project. The OpenMP
10 compiler introduced in [7] and this paper as a product are general available now. The
11 data listed here can be reproduced, although some slight difference may exit, as we are
12 using a development version of the compiler.

13 As the Fortran support is not robust as the C/C++ yet, We choose the NAS 2.3 C
14 version as our benchmarks.

15 Table 1 shows the performance of 8 SPUs against one PPU execution. We also
16 include two extra benchmarks, a PDE solver, packaged in the single source compiler
17 distribution as an example, and LBM, the CPU benchmark in spec2k6 with OpenMP
18 directives.

19 We do need to modify the source code a little bit to make it fit for the compiler, the
20 main changes are

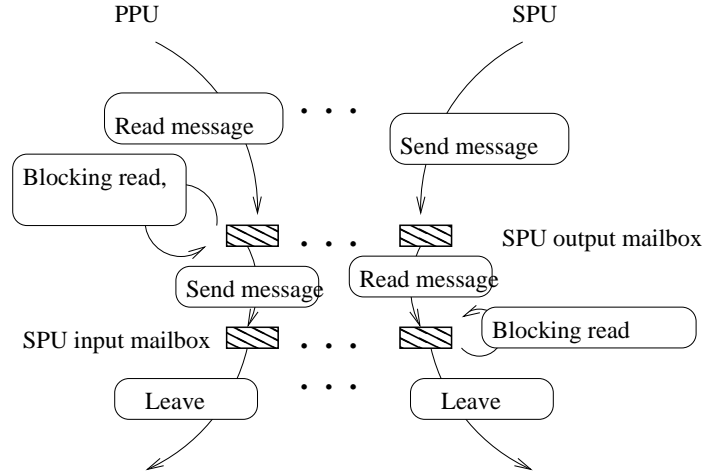


Fig. 5. Barrier with SPU output and input mailbox

- 1 – Using restricted pointer or adding disjoined pragma to indicate that two arrays will
- 2 not overlay each other. Although interprocedural compilation can trace pointer tar-
- 3 gets, there are cases which do need users to give more accurate info.
- 4 – Change the data location, either privatize a variable to make it stay on SPU local
- 5 store or malloc an array from system storage. These two looks contradictory, the
- 6 first one aims to speed up access, the second one is to save SPU local store space.
- 7 – Data access alignment. This is due to DMA operation will be more efficient if the
- 8 data can be 128 byte aligned.

9 From the table we can see that not every benchmark speed up nicely. Here are many
10 reasons for this. A complicated parallel region with a lot of branches inside may cause
11 SPU slow down. Data access pattern is not regular, DMA tiler can not tile the loop,
12 so the memory access has to go though software cache. As we are improving the tiler
13 performance, we do hope some of them can be further improved.

14 In this paper, we summarized how to map OpenMP constructs onto a Cell proces-
15 sor, from a `parallel` region to different workshares and barriers. We have omitted
16 the implementation details of these, concentrating on the framework which will help
17 the user visualize the framework of the OpenMP runtime system on a Cell processor.
18 Please refer to the references for how to implement software cache DMA tiler and SPU
19 simdization.

20 There is still on-going work to further enhance the compilation system, such a en-
21 large the DMA tiler region, improving auto simdization. For future runtime enhance-
22 ments, we will look at the task support in the OpenMP 3.0 standard.

Benchmark	# of lines changed	sequential time	parallel time	speedup
BT	0	588.48		
CG	0	20.12	13.78	1.46
EP	0	72.05	12.63	5.66
FT	4	42.19	11.66	3.65
IS	10	10.54	2.68	3.94
LU	0	356.8	209.71	1.71
MG	10	17.38	2.41	7.2
SP	10	413.99	134.69	3.08
PDE		27.59	0.65	42.38
LBM		1553.91	147.61	10.66

Table 1. Performance evaluation

6 Acknowledgments

Thanks to Tao Zhang from IBM Watson research lab and Daniel A. Brokenshire from IBM Austin Quasar/Cell software development for the discussions on the barrier implementation on the Cell system.

7 Trademarks and copyright

Cell Broadband Engine, AIX, IBM, POWER, and VisualAge are trademarks or registered trademarks of International Business Machines Corporation in the United States, other countries, or both.

Other company, product and service names may be trademarks or service marks of others.

© Copyright IBM Corp. 2008. All rights reserved.

References

1. D. P. et al. The design and implementation of a first-generation cell processr. In *IEEE International Solid-State Circuits Conference (ISSCC)*, February 2005.
2. Samuel Williams. et al. The potential of the cell processor for scientific computing. In *Conference on Computing Frontiers*, 2006.
3. M. Gordon. et al. Exploiting coarse-grained task, data, and pipeline parallelism in stream programs. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, October 2006.
4. M. Kistler, M. Perrone, and F. Petrini. Cell multiprocessor communication net-work: Built for speed. *IEEE Micro*, 26(3), 2006.
5. A. E. et al. Vectorization for simd architecture with alignment constraints. In *Conference on Programming Language Design and Implementation (PLDI)*, 2003.
6. OpenMP Architecture Review Board. Openmp application program interface version 2.5, 2005. <http://www.openmp.org>.

- 1 7. Kevin O'Brien, Kathryn O'Brien, Zehra Sura, Tong Chen, and Tao Zhang. Supporting
2 openmp on cell. In *International Workshop on OpenMP (IWOMP)*, 2007.
- 3 8. Guansong Zhang, Raul Silvera, and Roch Archambault. Structure and algorithm for imple-
4 menting OpenMP workshare. In *WOMPAT*, Lecture Notes in Computer Science. Springer,
5 2004.
- 6 9. OpenMP Architecture Review Board. Openmp application program interface version 3.0,
7 public review, 2008. <http://www.openmp.org>.
- 8 10. Guansong Zhang et al. Busy-wait barrier synchronization with distributed counter and local
9 sensor. In Michael Voss, editor, *WOMPAT*, volume 2716 of *Lecture Notes in Computer*
10 *Science*, pages 84–99. Springer, 2003.