

Technical Report

Department of Computer Science
and Engineering
University of Minnesota
4-192 EECS Building
200 Union Street SE
Minneapolis, MN 55455-0159 USA

TR 08-036

Exploring Speculative Parallelism in SPEC2006

Venkatesan Packirisamy, Antonia Zhai, and Pen-chung Yew

November 04, 2008

Exploring Speculative Parallelism in SPEC2006

Venkatesan Packirisamy, Antonia Zhai, and Pen-Chung Yew

University of Minnesota, Minneapolis.

{packve,zhai,yew}@cs.umn.edu

Abstract

Computer industry has adopted multi-threaded and multi-core architectures as the clock rate increase stalled in early 2000's. It was hoped that the continuous improvement of single-program performance could be achieved through these architectures. However, traditional parallelizing compilers often fail to effectively parallelize general-purpose applications which typically have complex control flow and excessive pointer usage. Recently hardware techniques like Transactional Memory (TM) and Thread-Level Speculation (TLS) have been proposed to simplify the task of parallelization by using speculative threads. Potential of speculative parallelism in general-purpose applications like SPEC CPU 2000 have been well studied and have shown to be moderately successful. Preliminary work that examined the potential parallelism in SPEC2006 deployed parallel threads with a restrictive TLS execution model and limited compiler support, and thus showed only limited performance potential. In this paper, we first analyze the cross-iteration dependence behavior of SPEC 2006 benchmarks and show that more parallelism potential is available in SPEC 2006 benchmarks, comparing against SPEC2000. Further, we use a state-of-the-art profile-driven TLS compiler to identify loops that can be speculatively parallelized. Overall, we found an average speedup of 60% on four cores over what could be achieved by a traditional parallelizing compiler such as Intel's ICC compiler on such benchmarks. We also found that an additional 11% improvement could be obtained on selected benchmarks using 8 cores when we extend TLS on multiple loop levels as opposed to restricting TLS only on a single loop level.

1 Introduction

With the advent of multi-threaded (e.g. simultaneous multi-threading (SMT) [7], hyper-threading [9]) and/or multi-core (e.g. chip multiprocessors (CMP) [8, 3]) architectures, the challenge now is to utilize these architectures to improve performance of general-purpose applications. Potential for Thread-level parallelism

(TLP) has been well understood in older benchmark suites like SPEC 2000 [18]. However with the current trend [4] towards more parallel applications, it is not clear if the newer general-purpose benchmarks in SPEC 2006 [19] exhibit more potential for Thread-Level Parallelism (TLP).

Automatic compiler parallelization techniques have been developed and found to be useful for many scientific floating-point intensive applications. However, when applied to general-purpose integer-intensive applications that have complex control flow and excessive pointer accesses, such traditional parallelization techniques become quite ineffective as they need to conservatively ensure program correctness by synchronizing *all potentially* possible dependences in the program. It often requires a programmer to explicitly create parallel threads and insert synchronizations explicitly. However, this approach has been proven to be error prone and could put a huge burden on the programmer.

More recently, there has been an extensive study on hardware support to simplify such tasks for the programmer and the compiler by supporting speculative threads. Hardware Transactional Memory (HTM) has been proposed to aid the development of parallel programs. Thread-Level Speculation (TLS) or Ordered Transaction has also been used to improve the performance of sequential applications beyond exploring instruction-level parallelism (ILP). With hardware support for speculative threads the compiler can parallelize both the loops where the parallelism cannot be proven at compile time and the loops where the cross-iteration dependences occur infrequently.

Though TLS has been extensively studied in the past, it is not clear how much TLS could benefit more recent benchmarks such as SPEC 2006 [19] which represent a different class of applications. Some recent studies [10] on SPEC 2006 benchmarks have shown very limited potential for TLS (less than 1%) under very conservative assumptions. In this paper, we re-examine some of the issues and give a more realistic assessment of TLS on such benchmarks using a state-of-the-art TLS compiler. We also identify potential speculative parallelism that could be further exploited by a compiler.

One of the key detriments in parallelizing loops is the presence of cross-iteration data dependences. We present a detailed analysis of data dependences patterns in SPEC 2006 and compare it with the dependence behavior of SPEC 2000. Our analysis show more potential for parallelism in SPEC 2006 than in SPEC 2000.

Fig. 1 shows the results of a potential study: the percentage of execution that can be parallelized if infrequently occurring cross-iteration data dependences can be *magically* resolved. The x -axis indicates the dependence frequency; and the y -axis indicates the percentage of total execution that can be parallelized. A data point at location (C, p) indicates: if loops containing only memory-resident value data dependences

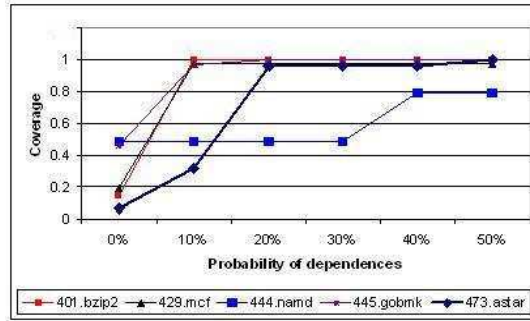


Figure 1: Coverage obtained by parallelizing loops with certain probability of data dependences.

that occur in less than $p\%$ of threads can be parallelized, then $C\%$ of total execution can be parallelized. We can see that for these SPEC 2006 benchmarks, there are many loops with low probability data dependences. For example in 473.astar if we ignore dependences that only occur in less than 20% of all iterations, we can parallelize loops that correspond to 96% of total execution. With a traditional compiler, all such dependences would be synchronized, and thus the resulting program exhibits poor parallel execution performance. With TLS, many of these loops could potentially be parallelized by speculating on such low probability data dependences.

It is clear from Fig. 1, there is more potential for TLS than previously shown in [10]. Kejariwal et al in [10] did not take into account the effect of compiler optimizations that could improve the performance of TLS. Previous studies [12, 26, 27, 22] have shown the importance of compiler-based loop selection and other compiler optimizations such as code scheduling to achieve good TLS performance. Also in [10], only innermost loops were considered for exploring TLS. In this paper, we do not limit ourselves to any loop level. We try to parallelize all potential TLS loops in a program. Instead of a high-level study on performance potential of TLS, we use a state-of-the-art TLS compiler to actually parallelize TLS loops in benchmark programs and show their performance. The results show that, with such TLS-oriented compiler optimizations, we could achieve an average of about 60% speedup for SPEC 2006 benchmarks over what could be achieved by a traditional parallelizing compiler such as Intel’s ICC compiler.

As the current trend is to support more cores on a chip, we also study potentially realizable TLS performance by extracting speculative threads at more than one loop level. We use a compiler-based static loop allocation scheme to efficiently schedule speculative threads from multiple loop levels. We show that an additional 11% improvement could be obtained on selected benchmarks by extending TLS for multiple loop levels using eight core.

In summary, the contributions of this paper are

1. We present a detailed analysis of inter-thread data dependences (both register- and memory-oriented data dependences) in SPEC 2006 benchmarks. We classify the benchmarks according to their data dependence behavior.
2. We present a comparison of cross-iteration dependence pattern of SPEC 2006 benchmarks with the SPEC 2000 benchmarks and show that the SPEC 2006 benchmarks have more potential for parallelism than the benchmarks in SPEC 2000.
3. With a state-of-art TLS compiler, we extract speculative threads from SPEC 2006 benchmarks and show the additional realizable TLS performance over a traditional parallelizing compiler.
4. We use a novel static loop allocation algorithm to study the performance potential of TLS when applied to nested loops.

The rest of the paper is organized as follows: Section 2 describes the related work; Section 3 analyzes the inter-thread dependences that occur in SPEC 2006 benchmarks; Section 4 describes our compiler framework and the evaluation methodology; Section 5 shows the performance of TLS and the scalability of TLS performance. In Section 6 we use single -level TLS performance to study the performance potential for multi-level TLS and in Section 7 we present our conclusions.

2 Related work

There has been a large body of research work on architectural design and compiler techniques for TLS, for example [21, 20, 26, 27, 6, 12]. But all of these papers based their studies using SPEC 2000 or other older benchmarks. None have used more recent SPEC2006 benchmarks. The SPEC 2006 benchmarks represent a newer class of applications and it is important to see if the conclusions drawn for SPEC 2000 will hold for these new applications. In this paper we address this issue by conducting a detailed study of SPEC 2006 benchmarks using a state-of-art TLS compiler.

Oplinger et al [14] presented a study of the limits of TLS performance on some SPECint95 benchmarks. The impact of compiler optimizations and the TLS overhead were not taken into account in that study. Similarly, Warg et al [25] presented a limit study for module-level parallelism in object-oriented programs.

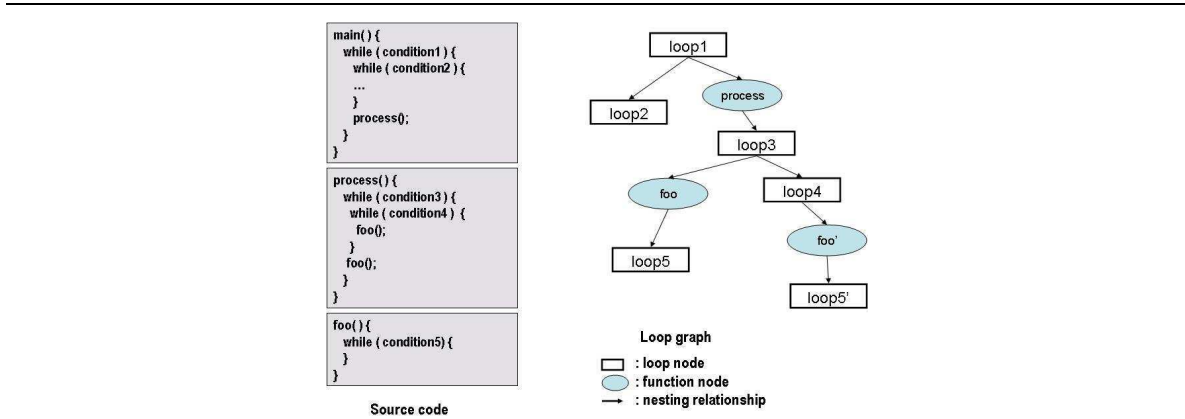


Figure 2: An example loop tree showing nesting relationship between loops.

In contrast, in this study, our aim is to show the realizable performance of TLS by a state-of-art TLS compiler and various TLS overheads were also considered.

Kejariwal et al [11] separated the speedup achievable from *traditional* thread-level parallelism vs. TLS, and studied the potential for TLS in SPEC2000 benchmarks assuming an *oracle* TLS mechanism. Similar to their study, we separate the speedup achievable through traditional non-speculative compiler techniques from the TLS-oriented compiler optimizations, and we consider all loop-levels instead of just inner-most or outer-most loops. They used manually-provided program analysis and forced the compiler to parallelize certain loops that the compiler will not normally automatically parallelize. Such manual analysis combined with simple loop transformations has been shown to create more parallelism in previous studies [16, 5, 12]. In this paper our aim is to present the TLS performance as an automatic parallelizing compiler without manual intervention.

In [10] Kejariwal et al used their framework of analysis to study SPEC 2006 benchmarks. Similar to [11], they used an *oracle* TLS mechanism to find potential TLS performance. Moreover, their study concentrated only on inner-most loops and used probabilistic analysis to predict TLS performance.

When compared to single-level TLS, performance of multi-level TLS has not been well understood. Renau et al [17] proposed hardware-based techniques to support multi-level TLS. In this paper, we use a compiler-based algorithm to schedule the loops to cores and study multi-level TLS performance in SPEC 2006 benchmarks.

3 Dependence analysis of SPEC 2006 loops

Cross-iteration data dependences are the most significant impediment to parallelizing loops. Under the context of TLS, the compiler needs to identify loops with infrequent inter-thread dependences which can be speculatively parallelized without much performance overhead due to dependence violations. Frequently occurring inter-thread data dependences need to be synchronized to ensure correct program execution and/or avoid speculation failure. However, synchronizations serialize execution and limit performance. Thus, understanding the inter-thread data dependence behavior is critical to understand the TLS performance potential. In this section, we analyze the dependence information collected through data dependence profiling, and estimate the importance of TLS hardware support in exploiting parallelism in the SPEC 2006 benchmarks.

In this section, the dependence information of all loops in an application is summarized through the notion of *combined execution time coverage*. The coverage of a loop is the fraction of total execution time of the program spent on this loop. In this section, the execution time of each loop is estimated using hardware performance counters, and the coverage of these loops are calculated as the percentage of loop execution time in the total execution time of the benchmark. Consider the code segment shown in Fig. 2: assume that `loop2`, `loop3` and `loop4` have no inter-thread data dependence, and thus ARE FULLY parallel; `loop1` has inter-thread data dependence, and is non-parallel. Thus, the *combined coverage* of data dependence-free loops is the coverage of `loop2` and `loop3` (`loop4` is nested inside `loop3`). To accurately estimate the *combined coverage* of a set of loops, the nesting relationship of these loops must be determined—this is done with the help of a loop tree (for example, Fig. 2). The loop tree used in this paper is similar to the loop graph described by Wang [22], except for loops that are invoked at different calling contexts. For these loops, separate loop tree nodes are created for each calling context. For example, `loop5` in Fig. 2 is replicated, since two different call paths can both lead to the invocation of `loop5`.

In this paper, we consider SPEC CPU 2006 benchmarks written in C or C++ (shown in Table 1). We ignored the programs written in FORTRAN since they tend to be parallel scientific programs that can be successfully parallelized by traditional parallelizing compilers and do not require TLS support. We present the results for 13 out of 18 benchmarks written in C or C++ in SPEC 2006. Due to technical difficulties in the compilation process, a few benchmarks are missing in our results which will be included in our final version.

Table 1: SPEC 2006 benchmarks.

| Benchmark | No. of Loops | No. of dynamic loop nesting levels |
|------------|--------------|------------------------------------|
| bzip2 | 232 | 11 |
| mcf | 52 | 5 |
| gobmk | 1265 | 22 |
| hmmer | 851 | 5 |
| sjeng | 254 | 10 |
| libquantum | 94 | 4 |
| h264ref | 1870 | 15 |
| astar | 116 | 6 |
| milc | 421 | 11 |
| namd | 619 | 4 |
| povray | 1311 | 15 |
| lbm | 23 | 3 |
| sphinx3 | 609 | 8 |

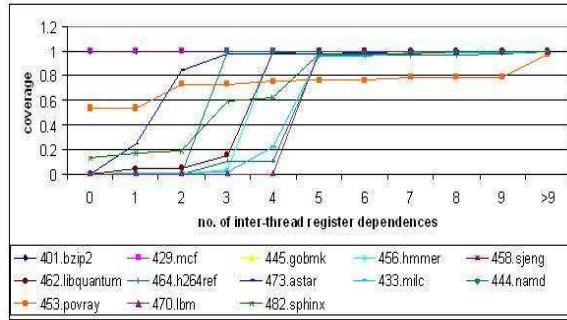


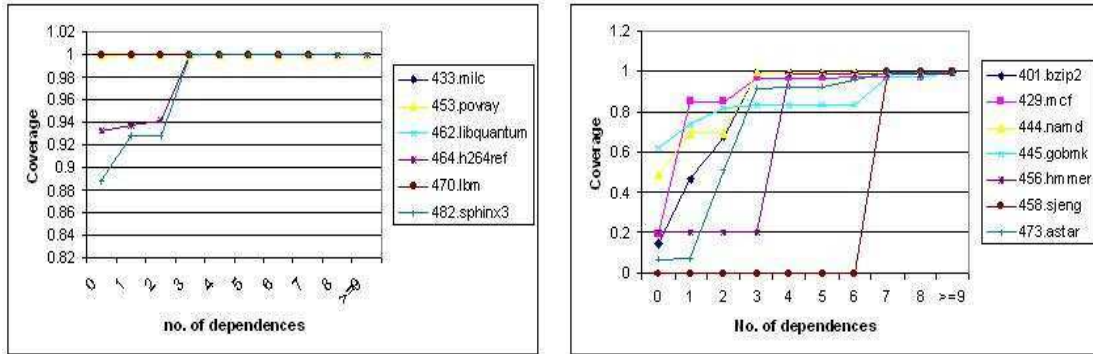
Figure 3: The accumulative coverage of loops with inter-thread register-oriented dependences

3.1 Inter-thread register-oriented data dependences

Cross-iteration data dependences could be caused either by register- or memory-resident values. We first focus on the relatively straightforward register-resident value dependences. For these dependences, the compiler is responsible for identifying instructions that produce and consume these value and generate synchronization to ensure correct execution. We count the number of inter-thread register-resident dependences (true dependences) for each loop; and estimate the *combined* coverage of the set of loops with certain number of register-resident dependences. The results are presented in Fig 3. The x-axis represents the number of register dependences and the y-axis represents the corresponding *combined* coverage estimated for a certain set of loops. If a benchmark has a combined coverage of C for x number of dependences, it indicates that loops with less than x dependences have a combined coverage of $C\%$. The benchmarks with high combined coverage (C) for a small number of dependences (x), potentially exhibit high degree of parallelism. We found that all the major loops in most benchmarks have inter-thread register-resident value dependences. An effective TLS compiler that is capable of synchronizing a few inter-thread register dependences is essential. Zhai *et.*

al [26] have described how such a compiler can be implemented; they have further shown that aggressive compiler scheduling techniques can reduce the critical forwarding path introduced by such synchronizations.

3.2 Inter-thread memory-oriented data dependences



(a) The coverage for benchmarks with fewer inter-thread memory dependences. (Class 'A'). (b) The coverage for benchmarks with significant inter-thread dependences. (Class 'B')

Figure 4: The coverage of loops with certain number of inter-thread memory-oriented data dependences

Unlike register-resident value dependences, memory-resident value dependences are difficult to identify using a compiler due to potential aliasing. To ensure correctness, a traditional parallelizing compiler would insert synchronizations on all possible dependences. With TM or TLS support, the compiler is able to aggressively parallelize loops by speculating on ambiguous data dependences. However, the performance of such execution depends on the likelihood of such data dependences occurring at runtime. If a data dependence do occur, the thread can potentially violate data dependence constraints, and thus must be squashed and re-executed; recovery codes can be executed to restore correct state. In this section, we conduct detailed analysis on inter-thread memory-resident value dependence using profiling information.

We classify benchmarks based on the combined coverage of loops with different number of memory-resident value dependences. Fig 4(a) shows the results of benchmarks (Points corresponding to 433.milc, 453.povray, 462.libquantum and 470.lbm in Fig. 4(a) overlap) that can achieve a high combined coverage with only a few inter-thread memory-oriented data dependences (class 'A'); Fig 4(b) shows the rest of the benchmarks (class 'B'). For benchmarks in class 'A' (such as 433.milc, 464.h264ref, and etc.), 90% of the total execution can potentially be parallelized by only considering loops with no inter-thread dependences. These benchmarks can be parallelized without hardware support for speculative execution, if the compiler is able to prove the independence.

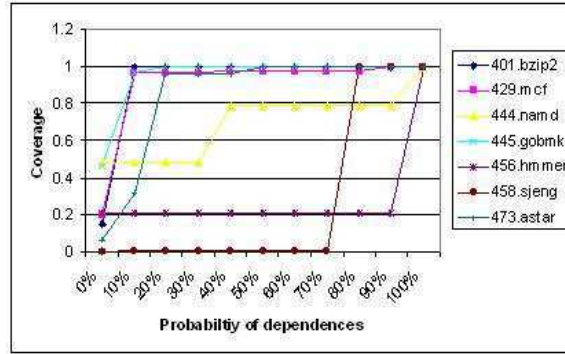


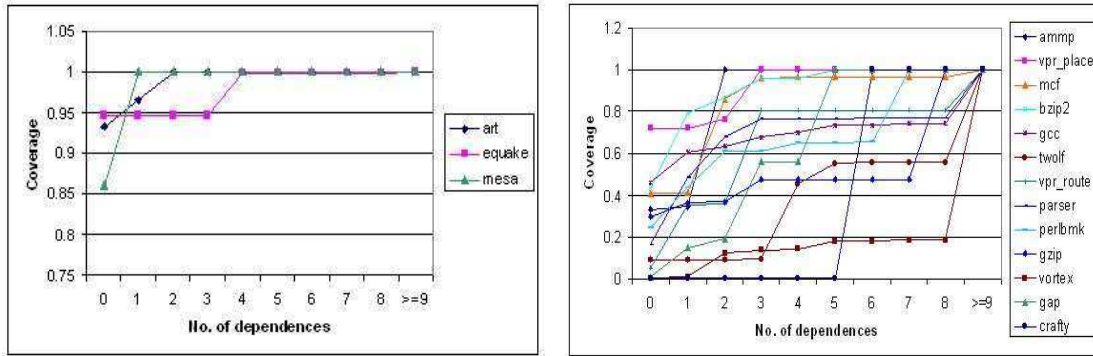
Figure 5: The coverage of loops with inter-thread memory-oriented data dependences less than a certain probability.

In class 'B' benchmarks, the speculative hardware support are potentially useful, since inter-thread data dependences do occur. Fig. 5 shows the probability of such data dependences and their corresponding coverage for class 'B' benchmarks. The x-axis represents the probability of inter-thread memory-resident value dependences and the y-axis represents the corresponding *combined* coverage estimated for a certain set of loops. If a benchmark has a combined coverage of $C\%$ for x probability of inter-thread dependence, it indicates the loops that only have inter-thread dependences with probability of less than x have a combined coverage of $C\%$. Benchmarks 401.bzip2, 429.mcf, 445.gobmk and 473.astar can achieve a large combined coverage, if we speculatively parallelized all loops that only contain data dependences that occur in less than 20% of iterations. Using TLS hardware support, these loops could potentially get good performance as the inter-thread data dependences occur only infrequently.

Some benchmarks, such as 456.hummer, 458.sjeng and 444.namd, can only achieve a high combined coverage by parallelizing loops containing frequently-occurring memory-resident value dependences. These dependences require synchronization. Previous studies has shown that frequently occurring memory-oriented data dependences could be synchronized by the compiler with profiling data [27]; and aggressive code scheduling could reduce critical-path length introduced by such synchronization [24].

3.3 SPEC 2006 vs. SPEC 2000

In this section, we aim to compare the potential parallelism in SPEC 2000 and SPEC 2006 benchmarks by observing their inter-thread data dependence behavior. Fig. 6(a) shows class 'A' benchmarks in SPEC 2000—benchmarks with few inter-thread data dependences; Fig. 6(b) shows class 'B' benchmarks—benchmarks with several cross-iteration dependences. Comparing against SPEC 2006 results, shown in Fig. 4(a) and in



(a) The coverage for benchmarks with fewer inter-thread memory dependencies. (Class 'A')
 (b) The coverage for benchmarks with significant inter-thread dependencies. (Class 'B')

Figure 6: The coverage of loops with certain number of inter-thread memory-oriented data dependencies in SPEC 2000

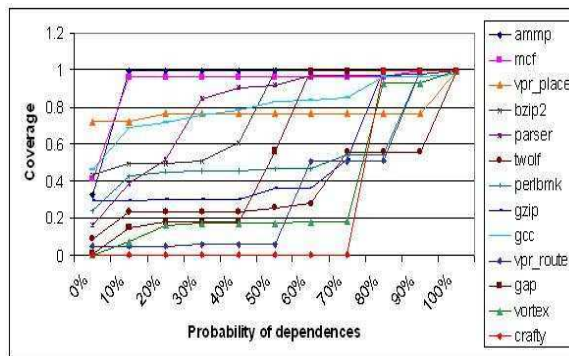


Figure 7: The coverage of loops with inter-thread memory-oriented data dependencies less than a certain probability in SPEC 2000.

Fig. 4(b), we found that SPEC 2000 suite has fewer class 'A' benchmarks. Also the class 'B' benchmarks in SPEC 2000 can only achieve high combined coverage by parallelizing loops with several cross-iteration dependencies. Furthermore, by examining Fig. 7, which presents the frequency of data dependencies that must be speculated during parallel execution, we found that with the exception of AMMP, MCF, VPR_PLACE AND BZIP2, class 'B' benchmarks in SPEC 2000 must speculate on high-probability cross-iteration dependencies to achieve a high combined coverage. This is consistent with results reported by previous studies: in SPEC 2000, only a few benchmarks, AMMP, MCF, VPR_PLACE, demonstrated high degree of parallelism under TLS. The data dependencies characteristics in SPEC2000 and SPEC2006 illustrate that SPEC 2006 can potentially achieve a higher degree of parallelism under the context of TLS.

3.4 Pitfalls

Even though profiling inter-thread data dependences is crucial in determining the suitability of using TLS to parallelize a loop, it does not indicate its actual performance. Actual TLS performance depends on other factors such as the size of the threads, thread spawning overhead, loop iteration counts, load variation among threads, etc. Also aggressive code scheduling can reduce the impact of synchronization for inter-thread dependences as shown in [26, 27, 22]. Some library calls could also potentially cause implicit dependences. A common example is the call to *malloc*, which could potentially cause inter-thread dependences due to its internal data structures. They could be eliminated by using specially written parallel libraries.

From the data presented in the earlier sections, we could see that SPEC 2006 benchmarks have numerous inter-thread dependences which could benefit from TLS hardware support. Such TLS hardware support could help to parallelize benchmarks in class 'B' with low-frequency data dependences and could also help the compiler in handling those ambiguous inter-thread data dependences (in class 'A' benchmarks). Also, many benchmarks have frequent register and memory dependences which could benefit from aggressively code scheduling by the compiler to reduce critical-path lengths introduced by synchronizations and increase execution overlap between threads.

4 Compilation and Evaluation Infrastructure

To quantitatively evaluate the amount of parallelism that can be exploited with hardware support for coarse-grain speculation and advanced compiler optimization technology, we evaluate a set of benchmarks from the SPEC2006 benchmark suite and simulate the execution of these benchmarks with an architectural simulator that support multiple cores and speculative execution. In the rest of this section, we will describe the execution model, as well as the compilation and simulation infrastructure used in this paper.

Execution Model:

Under thread-level speculation (TLS), the compiler partitions a program into speculatively parallel threads without having to decide at compile time whether they are independent. At runtime, the underlying hardware determines whether inter-thread data dependences are preserved, and re-executes any thread for which they are not. The most straightforward way to parallelize a loop is to execute multiple iterations of that loop in parallel. In our baseline execution model, the compiler ensures that two nested loops will not be speculatively parallelized simultaneously. In Section 6, we will study the potential for supporting speculative threads at

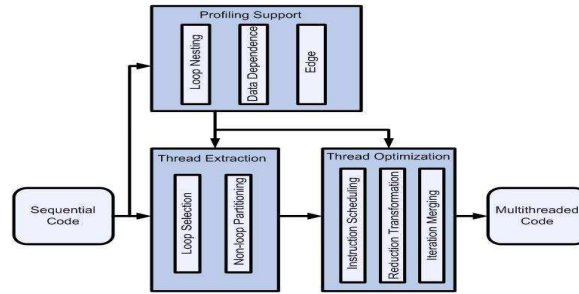


Figure 8: Compilation infrastructure

multiple nesting levels.

Compilation Infrastructure

Our compiler infrastructure is built on Open64 3.0 Compiler [2], an industrial-strength open-source compiler targeting Intel’s Itanium Processor Family (IPF).

To create and optimize speculative parallel threads, the compiler must perform accurate performance trade-off analysis to determine whether the benefit of speculative parallel execution outweighs the cost of failed speculation and then aggressively optimize loops that benefit from speculation. In our case, the compiler performs such analysis and optimizations based on loop nesting, edge, as well as data dependence profiling (using train input set), as shown in Figure 8. The TLS compiler has two distinct phases, as shown in Figure 8, thread extraction and optimization:

Thread Extraction: The compiler can extract threads from both loops and non-loop regions. In this paper we focus on loop-level parallelism. In the loop selection phase, the compiler first estimates the parallel performance of each loop, then choose to parallelize a set of loops that maximize the overall program performance based on such estimation. Previous work [23] builds the performance estimation based on detailed data dependence profiling information, as shown in Figure 8. Thus, the achievable performance of speculative parallel threads is tied with the accuracy of performance estimation. I.e., inaccurate profiling information, and inaccurate-estimation information can potentially lead to the selection of sub-optimal loops. In this paper, since we aim to demonstrate the optimal performance that can be achieved with SPEC2006 benchmarks, we eliminated this uncertainty from our evaluation. When selecting which loops to parallelize to maximize program performance, instead of relying on a compiler estimation, we use the simulated parallel and sequential execution time of each loop to

Table 2: Architectural parameters.

| Parameter | |
|---|---|
| Fetch/Decode/Issue/Retire Width | 6/6/4/4 |
| Integer units | 6 units / 1 cycle latency |
| Floating point units | 4 units / 12 cycle latency |
| Memory ports | 2Read, 1Write ports |
| Register Update Unit (ROB,issue queue) | 128 entries |
| LSQ size | 64 entries |
| L1I Cache | 64K, 4 way 32B |
| L1D Cache | 64K, 4 way 32B |
| Cache Latency | L1 1 cycle, L2 18 cycles |
| Memory latency | 150 cycles for 1st chunk, 18 cycles subsequent chunks |
| Unified L2 | 2MB, 8 way associative, 64B blocksize |
| Physical registers per thread | 128 Integer, 128 Floating point and 64 predicate registers |
| Thread overhead | 5 cycles fork, 5 cycles commit and 1 cycle inter-thread communication |
| No. of cores | 4 |

determine the actual benefit of parallelizing that loop.

Optimization: Loops that are selected for parallelization must be transformed for efficient speculative parallel execution. In our case, the following optimizations are applied: (i) all register-resident values, as well as memory-resident values that cause inter-thread data dependences in 20% of all threads are synchronized [27]; (ii) instructions are scheduled to reduce the critical forwarding path introduced by the synchronization [26, 22]; (iii) computation and usage of reduction and reduction-like variables are transformed to avoid speculation failure and reduce synchronization [28]; and (iv) consecutive loop iterations are merged to balance the workload between neighboring threads [28].

Simulation Infrastructure:

We use a trace-driven, out-of-order Superscalar processor simulation infrastructure. The trace-generation portion of this infrastructure is based on the PIN instrumentation tool [13], and the architectural simulation portion is built on SimpleScalar. We not only model the register renaming, the reorder buffer, branch prediction, instruction fetching, branching penalties and the memory hierarchy performance, but also extend the infrastructure to model different aspects of TLS execution including explicit synchronization through signal/wait, cost of thread commit/squash, etc. Different simulation parameters used are shown in Table 2. We use a cache based protocol based on STAMPede [20] to support speculative threads in CMP.

In this study, we use the reference input to simulate all benchmarks. In case of benchmarks with multiple input sets, the first input set is used. To get an accurate estimate of TLS performance, we parallelize and simulate all loops (with at least 0.05% dynamic execution time coverage) in each of the benchmarks. Based

Table 3: Coverage of loops parallelized.

| Benchmark | Coverage (%) | | | No. of loops | | |
|------------|--------------|-------------|-------------------|--------------|-------------|-------------------|
| | Type I | Type I + II | Type I + II + III | Type I | Type I + II | Type I + II + III |
| milc | 13 | 79 | 79 | 5 | 22 | 22 |
| lbm | 0 | 100 | 100 | 0 | 1 | 2 |
| h264ref | 0 | 53 | 83 | 2 | 32 | 36 |
| libquantum | 0 | 98 | 98 | 1 | 5 | 5 |
| sphinx3 | 40 | 83 | 91 | 11 | 19 | 21 |
| povray | 0 | 3 | 63 | 0 | 4 | 5 |
| bzip2 | 2 | 3 | 31 | 4 | 6 | 14 |
| mcf | 0 | 85 | 93 | 0 | 6 | 6 |
| namd | 1 | 8 | 96 | 7 | 22 | 50 |
| gobmk | 0 | 6 | 13 | 0 | 1 | 5 |
| hmmer | 0 | 0 | 79 | 2 | 1 | 6 |
| sjeng | 0 | 0 | 1 | 0 | 0 | 6 |
| astar | 0 | 5 | 99 | 0 | 2 | 8 |

on the simulated speedup of each loop, we use our loop selection algorithm to select the best set of loops which maximizes the performance of the entire benchmark. To report the speedup achieved by the entire benchmark, the average speedup of all the selected loops is calculated and weighted by the coverage¹ of the loops. For each simulation run, several billion instructions are fast-forwarded to reach the loops and different samples of 500 million instructions are simulated to cover all the loops.

5 Exploiting Parallelism in SPEC2006

In this section, we evaluate the amount of parallelism available in SPEC 2006 benchmarks using the framework described in Section 4. To isolate the parallelism that cannot be exploited without the help of TLS, we take three increasingly aggressive attempts to parallelize loops in SPEC 2006 benchmarks:

Type I: Loops that are identified as parallel by a traditional compiler;

Type II: Loops that have no inter-thread data dependence at runtime, but are not identified as parallel by the compiler, a.k.a., *Probably Parallel Loops*;

Type III: Loops that contain inter-thread data dependences, thus require TLS support to parallelize, a.k.a., *True Speculative Loops*.

Table 3 shows the percentage of total execution that can be parallelized when loops of different types become parallelizable. Type I is not shown separately in Table 3 due to very low execution time coverage.

¹The coverage of a loop is defined as the fraction of dynamic execution time of the loop

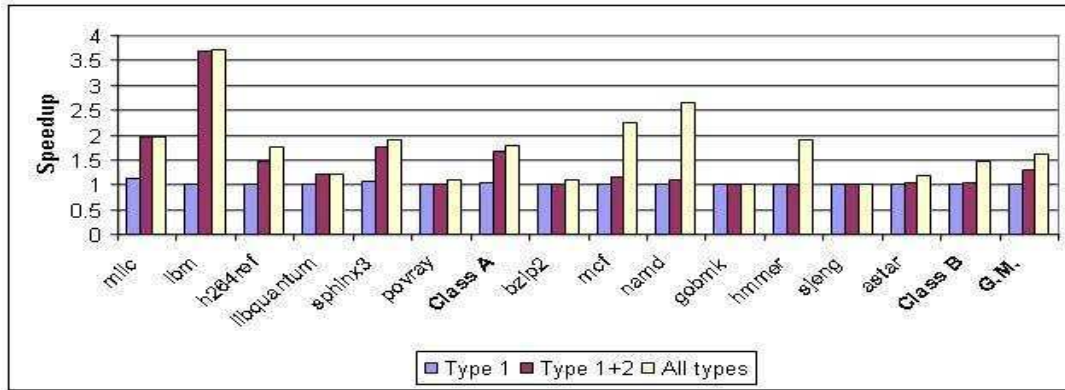


Figure 9: Shows the program speedup when different types of loops are parallelized using 4 cores.

5.1 Type I Loops

We applied the Intel C++ compiler [1] to the SPEC 2006 benchmarks to select parallel loops. The benchmarks are compiled with options `-O3 -ipo -parallel -par-threshold0`. The option `-par-threshold0` allows the compiler to parallelize loops without taking thread overhead into consideration. The loops selected by the Intel compiler are then parallelized by our TLS compiler and simulated. The speedup achieved by the selected loops compared against that of sequential execution is shown as the first set of bars in Figure 9. We can see that except in benchmarks MILC which achieved a speedup of 11% and SPHINX3 which achieved a speedup of 7%, other benchmarks don't have any speedup. For the entire set of benchmarks, we get a geometric mean speedup of just 1%.

This result is anticipated, since the complex control flow and ambiguous data dependence patterns prohibits the traditional compiler from parallelizing large loops. We have found that in most benchmarks the compiler has only chosen to parallelize simple inner loops with known iteration count. It is worth pointing out that, although many class **A** benchmarks, such as MILC and LBM, contain loops with no inter-thread data dependences, the compiler is unable to identify these loops as being parallel.

5.2 Type I + II Loops

With the addition of *Probably Parallel Loops*, class **A** benchmarks achieve significant performance gain, however, class **B** benchmarks remain sequential. The class **A** benchmarks gain 68% speedup due to these *Probably Parallel Loops* while class **B** benchmarks gain only 4%. If the compiler is able to determine that these loops are parallel, we can potentially parallelize these loops without TLS support. Among the class **A**

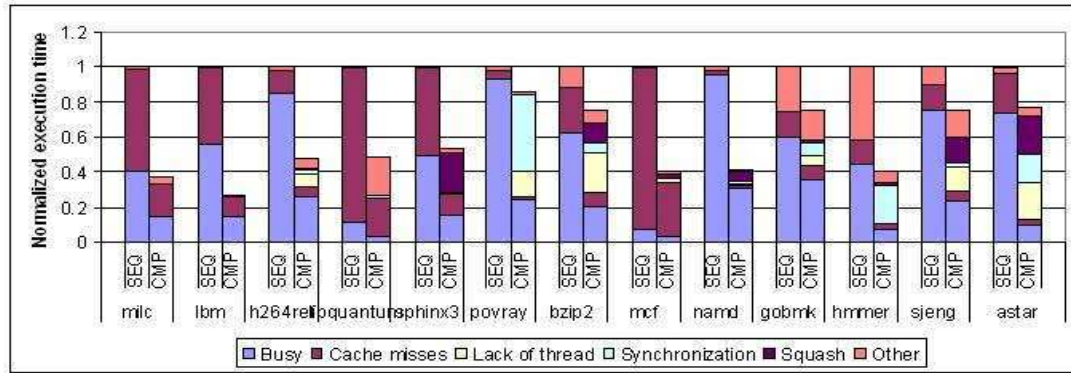


Figure 10: Shows the breakdown of execution time while executing the selected loops normalized to sequential execution time

benchmarks, significant portion of the loops in SPHINX3 and H264REF are *Probably Parallel Loops*; and all loops in MILC, LBM and LIBQUANTUM are *Probably Parallel Loops*.

5.3 Type I + II + III Loops

With the addition of *True Speculative Loops*, we find that many class **B** benchmarks are able to achieve speedup. With only these *True Speculative Loops* class **B** benchmarks gain a speedup of 42% giving them an overall speedup of 46%.

To examine TLS performance in detail, Figure 10 shows the execution time breakdown of parallel execution with TLS support (only selected loops) and sequential execution. The **SEQ** bars show the normalized execution time of the sequential execution running on one core. The **CMP** bars show the normalized execution time of the parallel program executing on four cores. Each bar is divided into six segments: *Busy* represents the amount of time spent in executing useful instructions and the delay due to lack of instruction level parallelism inside each thread; *Lack of threads* represents the amount of time wasted due to the lack of parallel threads (probably due to low iteration count in a loop); *Synchronization* represents the amount of time spent in synchronizing frequently occurring memory dependences and register dependences; *Cache misses* represents the amount of time the processor stalled due to cache misses; *Squash* represents the amount of time wasted executing instruction that are eventually thrown away due to failed speculation; *Other* corresponds to everything else. In particular, it includes time wasted due to speculative buffer overflow and load imbalance between consecutive threads.

We will focus on the class **B** benchmarks first. In HMMER, the loop at `fast-algorithms.c:133` is

selected for parallelization, however it has many inter-thread dependences that need synchronization. These synchronizations create a critical forwarding path between threads and serialize execution. Thus, by performing speculative instruction scheduling to move the producers of these dependences as early as possible in the execution [26, 22], the parallel overlap is significantly increased; and the benchmark achieves a 90% program speedup. Similar behavior is observed in NAMD, where synchronization and instruction scheduling leads to a 164% program speedup.

For ASTAR, the important loop is at `way2_.cpp:100`, which has a few inter-thread dependences. Some of these dependences are frequent, and thus are synchronized; others are infrequent, and thus are speculated on. Without TLS support, these infrequent occurring dependences must be synchronized, and can lead to serialization of the execution. With the help of TLS, this loop achieves a 17% speedup.

POVRAY, although a class **A** benchmark, is able to benefit from speculation. The important loop in `csg.cpp:248` is a *True Speculative Loop* with a few mispeculations, thus it is non-parallel for a traditional compiler. Unfortunately, the selected loops have small trip counts, and the cores are often idle; thus the benchmark is only able to achieve a moderate program speedup of 9%.

Not all benchmarks are able to benefit from TLS. GOBMK has many loops with low trip counts, thus many execution cycles are wasted as the cores are idling. Loops with large trip counts are not able to achieve the desired speedup for two reasons: first of all, the amount of work in consecutive iterations is often unbalanced; secondly, many iterations have large memory footprints that lead to buffer overflow of the speculative states. The geometric mean of the thread size for the top 50 loops (in terms of coverage) is 800,000 instructions. Overall, GOBMK only achieves 1% performance improvement with TLS support.

Loops in SJENG have many inter-thread dependences that occur in 70% of all iterations, and thus need synchronization. However, the critical forwarding path introduced by these synchronization cannot be reduced through instruction scheduling due to intra-thread dependences. Thus, SJENG was unable to benefit from TLS.

To summarize, TLS is effective in parallelizing both class **A** and class **B** loops. Overall, we are able to achieve a program speedup of about 60% (geometrical mean) when all loops are considered, in contrast to a traditional compiler, which only achieves a 1% program speedup.

Table 4: Important loops parallelized for each benchmark.

| Benchmark | Location in SPEC2006 | Loop coverage | Speedup | Average iter. Size | Aver. Iters. per Invoca-tion | Comments |
|----------------|--|---------------|--------------|--------------------|------------------------------|---|
| 401.bzip2 | blocksort.c,551 | 14% | 1.16 | 4880 | 2 | Suffers from low iteration count and frequent squashes |
| 429.mcf | pbeampp.c,165 | 65% | 2.40 | 335 | 292 | Achieves good speedup. Needed some synchronization. |
| 433.milc | quark_stuff.c, 1523 quark_stuff.c,1452 | 33% 5% | 3.55 3.52 | 1972 1972 | 20736 324 | Loop is very parallel. Loop is very parallel. |
| 444.namd | ComputeNonbondedUtil.h, (calc_pair_energy_merge_fullselect) ComputeNonbondedUtil.h, (calc_pair_nonbonded) | 4% 4% | 3.17 3.09 | 375 63 | 55 230 | few squashes. Loop is parallel. |
| 453.povray | csg.cpp,248 | 60% | 1.17 | 2458 | 4 | Suffers from low iteration count. |
| 456.hmmr | fast_algorithms.c,133 | 78% | 2.51 | 156 | 300 | Large speedup due to speculative instruction scheduling. |
| 462.libquantum | gates.c,89 gates.c,61 | 62% 12% | 1.18 1.18 | 43 48 | 32765 32768 | Loops parallel, but iteration size small. Small parallel loop. |
| 464.h264ref | mv-search.c,394 mv-search.c,982 | 37% 15% | 3.46 3.74 | 6837 68 | 1089 1089 | Loop is parallel. Loop has small iteration size. |
| 470.lbm | lbm.c,186 | 99% | 1.19 | 525 | 1300000 | Loop is parallel. |
| 473.astar | Way2_.cpp,100 | 60% | 1.17 | 1548 | 553 | Need synchronization, suffers from squashes. |
| 482.sphinx3 | vector.c,513 approx_cont_mgau.c,279 | 37% 36% | 3.89 1.27 | 1108 2189 | 2048 6144 | Loop is parallel. Frequent squashes but some prefetching effect. |

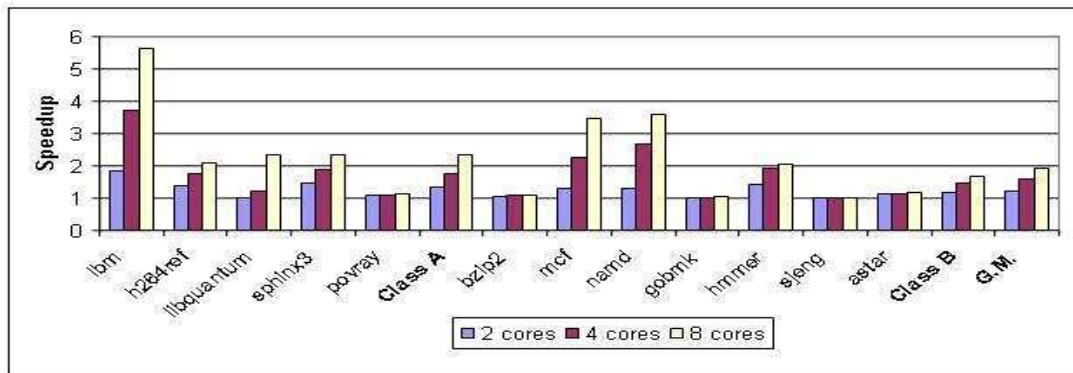


Figure 11: Speedup increases with increasing number of cores

5.4 Scalability

As technology scales, the number of cores that can be integrated onto a single die increases. Thus, it is important to understand whether TLS can efficiently utilize all the available cores. In this section, we study the scalability of TLS performance by comparing the speedup achieved using two, four and eight cores. The

results of this study are shown in Figure 11.²

When the number of cores is increased from two to four, the geometric mean of the speedup increases by about 35%; when increased further to eight cores, the performance increase is 33%. Among class **A** benchmarks, LBM, SPHINX3, H264REF and LIBQUANTUM contain important loops that have large iteration count and substantial amount of parallelism, thus the performance of these benchmarks is able to scale with the number of cores. In LIBQUANTUM, the super-linear performance gain is due to cache prefetching effect between the speculative threads.

Among class **B** benchmarks, NAMD shows good scalability and MCF benefits from cache prefetching effect as the number of threads increases. Unfortunately, none of the other benchmarks are scalable: HUMMER suffers from frequent synchronization; POV-Ray and BZIP2 suffers from small trip counts; ASTAR not only suffers from frequent synchronization, but also frequent squashes; For GOBMK and SJENG the performance improvement for TLS is negligible in all configurations.

To summarize, with our existing execution model, only a few benchmarks are able to scale with the number of cores; and even for the benchmarks that do scale, most of them scale sub-linearly. While the reasons for the lack of scalability differ from benchmark to benchmark, it is obvious that the amount of parallelism is limited. Thus, we will develop new execution models to improve the scalability of TLS in the next section.

6 Exploiting Speculative Parallelism at Multiple Levels of Loops

As the number of cores on a chip keeps increasing, it is important to ensure that the parallel performance TLS loops can scale. Unfortunately, almost all benchmarks scale sub-linearly, and many, HUMMER, POV-Ray, BZIP2, ASTAR, GOBMK and SJENG exhibits poor scalability. Similar to many previous proposals on TLS, the results presented in Section 5.4 only allow one loop nest level to parallelize when multiple levels of nested loops exist. Our goal in this section, is to examine the feasibility and benefit of exploiting TLS parallelism from multiple loop levels in a loop nest simultaneously. Renau *et. al* [17] proposed an architectural design to support out-of-order forking of speculative threads. It allows spawning of more speculative threads on the outer loop levels before spawning less speculative threads on the inner loops. However, it is unclear how to best allocate cores to speculative threads for the outer and the inner loops. In their work, a complex hardware-

²The scalability results for MILC are still being collected and will be included in our final version

based design, referred to as *Dynamic Task Merging* is deployed. Hardware counters are used to identify threads that suffer frequent squashes; and these threads are prohibited from spawning new threads. Other than requiring complex hardware support, the proposed approach has the following limitations. First of all, using number of squashes as a measurement of TLS efficiency is inaccurate. As we have seen in Section 5, the performance of TLS loops is also determined by synchronization, iteration count, load balancedness, and etc. Secondly, the hardware is unable to pre-determine the TLS potential of inner loops. As long as the outer loop does not show performance degradation, it will monopolize the cores. The inner loops will never be attempted for parallelization even if it has more parallelism. This will lead to a suboptimal allocation of cores.

In this section, we propose a compiler-based core allocation scheme that statically schedules cores for threads extracted from different levels of a loop nest. In this section, the performance potential of parallelizing multiple levels of TLS loops is extrapolated from the parallel performance of single level of parallel execution using the compiler-based allocation scheme. Since the nested loops are parallelized simultaneously, this extrapolation can be inaccurate, due to some secondary effects such as increase/decrease in cache misses. Furthermore, it is possible for the number of squashes for an outer loop to vary slightly due to the parallelization of inner loops. However, we believe that the impact of these secondary effects can be negligible; and neglecting these effects does not change the conclusion of this paper. Exploring the hardware modifications needed to support the parallelization of multiple levels of loops is beyond the scope of our paper.

Compiler-based scheduling schemes for nested loops have been studied in the past to support nested DOALL and DOACROSS loops. We extend the OPTAL algorithm [15] that was originally designed for core allocation for nested DOALL and DOACROSS loops to allocate cores for TLS loops at compile time.

6.1 Speculative OPTAL algorithm

The OPTAL algorithm uses a dynamic programming based *bottom-up* approach to decide how many cores to allocate to each loop nest level so that entire benchmark with multiple-levels of nested parallel loops can achieve the optimal performance. Based on simulated performance and coverage of each loop level, the algorithm considers different possible allocations on TLS loop nests and selects the allocation that maximize program performance. When deciding how many cores to allocate to a TLS loop at a particular loop-nest level, the algorithm only examines its immediate inner loops.

The inputs to the algorithm are the loop tree obtained during loop-nest profiling phase of the compiler and the maximum number of cores available for the benchmark (say 2^K). The output is the optimal core allocation for each TLS loop level and also the estimated optimal performance for the benchmark.

Predicting performance for each loop Before estimating the multi-level TLS loops performance, the algorithm estimates the single-level performance of each TLS loop level. In the case of a DOACROSS loop, the time required for the parallel execution T_p can be calculated by its initiation delay d . A similar method to estimate the performance of a TLS loop has been studied in [22]. In this study, we use the real performance result for each TLS loop we obtained in the previous section.

Let $SingleSpeedup_{i,j}$ represent the speedup achieved by parallelizing the single-level TLS $loop_i$ using j cores, where $j \in 2^0, 2^1, 2^2, \dots, 2^K$. Let $BestSpeedup(i,j)$ represent the best speedup achievable by parallelizing the multi-level TLS loop nest starting at ($loop_i$) using j cores.

Combining speedup The basic step in the algorithm is to find the speedup of an outer loop when its inner loops are also parallelized. Lets call the function to calculate this as $GetCombinedSpeedup(L_i, M, N)$. $GetCombinedSpeedup(L_i, M, N)$ returns the speedup of loop L_i when we allocate M cores to parallelize L_i and N cores to parallelize its child loops in the loop tree $L_{i,j}$. $GetCombinedSpeedup(L_i, M, N)$ is shown in Algorithm 1.

Input: Outer loop L_i , M cores allocated to L_i , N cores allocated to the next level (L_i 's child loops)

Output: Speedup of L_i with the specified allocation

Read $Cycles(L_i)$ - the number of cycles spent in loop L_i from the profile;

foreach Child of loop L_i , $L_{i,j}$ **do**

 /*Find total sequential cycles T_s for all inner loops.*/

 Read $Cycles(L_{i,j})$ - the number of cycles spent in loop $L_{i,j}$ when invoked from L_i from the profile;

 SumBefPar += $Cycles(L_{i,j})$;

 /*Find total T_p for all inner loops after parallelization.*/

$ParCycles(L_{i,j}) = Cycles(L_{i,j}) \div BestSpeedup(L_{i,j}, N)$;

 SumAfterPar += $ParCycles(L_{i,j})$;

end

 /* T_s of outer loop after inner loops are parallelized.*/

$CyclesInnerPar = Cycles(L_i) - SumBefPar + SumAfterPar$;

 /*Calculate combined speedup.*/

$ParCycles(L_i) = CyclesInnerPar \div SingleSpeedup_{i,M}$;

 return ($Cycles(L_i) \div ParCycles(L_i)$);

Algorithm 1: The $GetCombinedSpeedup(L_i, M, N)$ to get the speedup of outer loop L_i when its inner loops are parallelized.

Recursive algorithm

The Speculative-OPTAL algorithm starts from the leaf level in the loop tree and calculates the speedup of parent nodes based on the child loops' speedup. The Speculative-OPTAL algorithm is shown in Algorithm 2.

Input: Loop L_i and 2^K the number of cores to allocate.

Output: Estimated speedup of the entire benchmark - $BestSpeedup(L_{root}, 2^K)$ and a vector $ChildAllocate(L_i, p)$ which indicates for each loop, how many cores need to be allocated to its child loops.

```

if  $L_i$  is leaf then
  foreach  $p \in \{ 2^0, 2^1, 2^2, \dots, 2^K \}$  do
     $BestSpeedup(i, p) = SingleSpeedup_{i, p};$ 
  end
  return;
end
/*Allocate child loops first.*/
foreach  $L_{i, j}$  child of  $L_i$  do
  Speculative-OPTAL( $L_{i, j}$ );
end
/*Inner loop's best allocation already known. Try all possible allocations for the outer loop*/
foreach  $p \in \{ 0, 1, \dots, K \}$  do
  foreach  $q \in \{ 0, 1, \dots, p \}$  do
     $CurSpeedup(L_i, q) = GetCombinedSpeedup(L_i, 2^q, 2^p/2^q);$ 
    if  $CurSpeedup(L_i, q) \leq MaxSpeedup$  then
       $MaxSpeedup = CurSpeedup(L_i, q);$ 
       $ChildAllocate(L_i, p) = q;$ 
    end
  end
   $BestSpeedup(i, p) = MaxSpeedup;$ 
end

```

Algorithm 2: The Speculative-OPTAL algorithm.

The exact allocation to the inner loops is given in the vector $ChildAllocate(L_i, p)$. This would be used by the compiler to statically allocate cores. Here, we are interested in the value of $BestSpeedup L_{root}, 2^K$, the performance potential of TLS for the entire benchmark when applied to multiple loop levels.

Complexity analysis

The Speculative-OPTAL is called for every node in the loop tree. Let λ be the number of nodes in the tree (including both loops and functions). In Speculative-OPTAL, the outer loop iterates for K times and the inner loop iterates on the average of $K/2$ times. So, total number of times $GetCombinedSpeedup$ is called is $K^2/2$. The loop inside $GetCombinedSpeedup$ iterates over all the children of the node. The average number of children per node is a constant (C_1) for a tree. So, the total time taken for Speculative-OPTAL is approximately $\lambda K^2/2$. Therefore, the complexity of Speculative-OPTAL is $O(\lambda)$ since K is a constant.

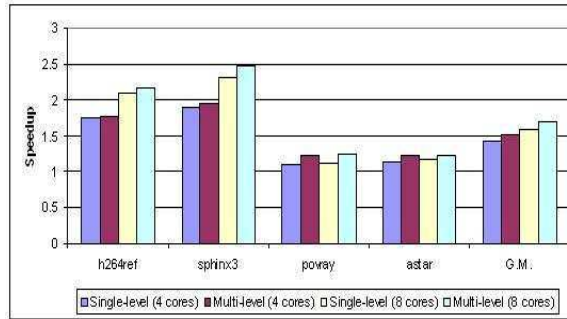


Figure 12: Speedup over sequential execution using 4 and 8 cores with multi-level TLS loops

6.2 Results

If TLS threads are spawned only from a single level of loop, many benchmarks, both from class A and class B, exhibit diminishing return as the number of cores increases, as shown in Fig. 11. With spawning speculative threads from multiple levels of loops, cores can potentially be better utilized, resulting in better scalability of the cores. We applied our allocation algorithm to extract speculative threads for all benchmarks, but omitted the results for some from our presentation due to the following reasons: BZIP2, GOBМК and SJENG are omitted due to the lack of TLS parallelism overall; LIBQUANTUM and HUMMER are omitted due to the lack of nested parallelizable loops; LBM, NAMD and MCF are omitted due to the fact that the Speculative-OPTAL algorithm is unable to identify multiple levels of loops that can perform better than the single level TLS.

For the remaining benchmarks, Fig. 12 compares the performance of multi-level TLS estimated based on our algorithm with the single-level TLS performance. With four cores, we can extract TLS parallelism from two levels of loops, and with eight cores, we extract TLS parallelism from upto three levels. In SPHINX3 the selected loop cannot utilize all the available cores due to frequent squashes. Multi-level TLS is able to improve the performance over single-level TLS by about 18% with 8 cores. In POVRAY, the iteration counts of the selected loops are low (around 4), and thus many cores idle. With multi-level TLS, all the cores are utilized, and thus multi-level TLS is able to outperform single-level TLS by about 13%. In benchmark ASTAR the selected loops have frequent mis-speculations, and thus execution cycles are wasted. With multi-level TLS, two nested loops, in `Way2_.cpp`, line 65 and line 100 respectively, are selected, are parallelized simultaneously. The multi-level TLS outperforms single-level TLS by 6% on 8 cores. In H264REF, multi-level TLS outperforms single-level TLS by 7% on 8 cores. Overall, for the benchmarks mentioned, Fig. 12 shows that by extracting speculative threads at multiple levels of loop nest, we are able to achieve an additional speedup of about 8% with four cores and 11% with eight cores for these selected

benchmarks.

From the above results, it is clear that many SPEC 2006 benchmarks have potential for multi-level TLS. Our compiler-based allocation algorithm can allocate cores among the multi-levels TLS loops to improve overall performance. With such a compiler-based approach, we could also avoid complex hardware modifications needed to support multi-level TLS.

7 Conclusions

Previous studies in SPEC 2006 based on high level analysis have shown only a limited potential for TLS. These studies had not taken into account the benefits of compiler based optimizations. In this paper we use a state-of-art TLS compiler to speculatively parallelize SPEC 2006 applications. We show that TLS is essential for parallelizing SPEC 2006 benchmarks.

We show that even though some benchmarks have parallel loops, the traditional parallelizing compiler cannot prove that there is no inter-thread dependences due to complex control flow and ambiguous data accesses. TLS can parallelize these *Potentially Parallel Loops* and achieve a speedup of 78% for benchmarks MILC, LBM, H264REF, LIBQUANTUM, SPHINX and POVRAY. Also TLS can parallelize loops that have infrequent inter-thread dependences (*Truly Speculative Loops*) that cannot be parallelized using traditional techniques. With TLS, benchmarks BZIP2, MCF, NAMD, GOBMK, HMMER, SJENG and ASTAR can achieve a speedup of 46% which is not possible with traditional compiler (even with full knowledge of the dependences).

With four cores we achieved a geometric mean speedup of 60% and with eight cores the speedup was 91% when compared to sequential execution.

We use a novel compiler based allocation scheme to efficiently allocate cores to exploit multi-level TLS. We show that with multi-level TLS we can achieve an additional 11% performance with 8 cores on selected benchmarks.

References

- [1] Intel c++ compiler. <http://www.intel.com/cd/software/products/asmo-na/eng/277618.htm>.
- [2] Open64 the open research compiler. <http://www.open64.net/>.
- [3] AMD CORPORATION. Leading the Industry: Multi-core Technology & Dual-Core Processors from AMD. <http://multicore.amd.com/en/Technology/>, 2005.
- [4] ASANOVIC, K., BODIK, R., CATANZARO, B. C., GEBIS, J. J., HUSBANDS, P., KEUTZER, K., PATTERSON, D. A., PLISHKER, W. L., SHALF, J., WILLIAMS, S. W., AND YELICK, K. A. The landscape of parallel computing research: A view from berkeley. Tech. Rep. UCB/EECS-2006-183, EECS Department, University of California, Berkeley, Dec 2006.

- [5] BRIDGES, M. J., VACHHARAJANI, N., ZHANG, Y., JABLIN, T., AND AUGUST, D. I. Revisiting the sequential programming model for the multicore era. *IEEE Micro* 28, 1 (2008).
- [6] DU, Z.-H., LIM, C.-C., LI, X.-F., YANG, C., ZHAO, Q., AND NGAI, T.-F. A Cost-Driven Compilation Framework for Speculative Parallelization of Sequential Programs. In *ACM SIGPLAN 04 Conference on Programming Language Design and Implementation (PLDI'04)* (June 2004).
- [7] EMER, J. EV8: The Post-ultimate Alpha.(Keynote address). In *International Conference on Parallel Architectures and Compilation Techniques* (2001).
- [8] INTEL CORPORATION. Intel's Dual-Core Processor for Desktop PCs. http://www.intel.com/personal/desktopcomputer/dual_core/, 2005.
- [9] INTEL CORPORATION. Intel Pentium 4 Processor with HT Technology. <http://www.intel.com/personal/products/pentium4/hyperthreading.htm>.
- [10] KEJARIWAL, A., TIAN, X., GIRKAR, M., LI, W., KOZHUKHOV, S., BANERJEE, U., NICOLAU, A., VEIDENBAUM, A. V., AND POLYCHRONOPOULOS, C. D. Tight analysis of the performance potential of thread speculation using spec cpu 2006. In *ACM SIGPLAN 2007 Symposium on Principles and Practice of Parallel Programming* (2007).
- [11] KEJARIWAL, A., TIAN, X., LI, W., GIRKAR, M., KOZHUKHOV, S., SAITO, H., BANERJEE, U., NICOLAU, A., VEIDENBAUM, A. V., AND POLYCHRONOPOULOS, C. D. On the performance potential of different types of speculative thread-level parallelism. In *20th Annual ACM International Conference on Supercomputing* (2006).
- [12] LIU, W., TUCK, J., CEZE, L., AHN, W., STRAUSS, K., RENAU, J., AND TORRELLAS, J. POSH: A TLS Compiler that Exploits Program Structure. In *ACM SIGPLAN 2006 Symposium on Principles and Practice of Parallel Programming* (March 2006).
- [13] LUK, C.-K., COHN, R., MUTH, R., PATIL, H., KLAUSER, A., LOWNEY, G., WALLACE, S., REDDI, V., AND HAZELWOOD, K. Pin: building customized program analysis tools with dynamic instrumentation. In *ACM SIGPLAN 05 Conference on Programming Language Design and Implementation (PLDI'05)* (June 2005).
- [14] OPLINGER, J., HEINE, D., AND LAM, M. In Search of Speculative Thread-Level Parallelism. In *Proceedings PACT 99* (October 1999).
- [15] POLYCHRONOPOULOS, C. D., KUCK, D. J., AND PADUA, D. A. Utilizing multidimensional loop parallelism on large-scale parallel processor systems. *IEEE Trans. Computers* 38, 9 (1989).
- [16] PRABHU, M., AND OLUKOTUN, K. Using Thread-Level Speculation to Simplify Manual Parallelization. In *ACM SIGPLAN 2003 Symposium on Principles and Practice of Parallel Programming* (2003).
- [17] RENAU, J., TUCK, J., LIU, W., CEZE, L., STRAUSS, K., AND TORRELLAS, J. Tasking with Out-of-Order Spawn in TLS Chip Multiprocessors: Microarchitecture and Compilation. In *19th Annual ACM International Conference on Supercomputing* (June 2005).
- [18] STANDARD PERFORMANCE EVALUATION CORPORATION. The SPEC CPU 2000 Benchmark Suite. <http://www.specbench.org>.
- [19] STANDARD PERFORMANCE EVALUATION CORPORATION. The SPEC CPU 2006 Benchmark Suite. <http://www.specbench.org>.
- [20] STEFFAN, J. G., COLOHAN, C. B., ZHAI, A., AND MOWRY, T. C. The stampede approach to thread-level speculation. In *ACM Trans. on Computer System* (August 2005), vol. 23, pp. 253–300.
- [21] TSAI, J.-Y., HUANG, J., AMLO, C., LILJA, D., AND YEW, P.-C. The Superthreaded Processor Architecture. *IEEE Transactions on Computers, Special Issue on Multithreaded Architectures* 48, 9 (September 1999).
- [22] WANG, S. *Compiler Techniques for Thread-Level Speculation*. PhD thesis, University of Minnesota, 2007.
- [23] WANG, S., YELLAJYOSULA, K. S., ZHAI, A., AND YEW, P.-C. Loop Selection for Thread-Level Speculation. In *The 18th International Workshop on Languages and Compilers for Parallel Computing* (Oct 2005).
- [24] WANG, S., ZHAI, A., AND YEW, P.-C. Exploiting Speculative Thread-Level Parallelism in Data Compression Applications. In *The 19th International Workshop on Languages and Compilers for Parallel Computing* (Oct 2006).
- [25] WARG, F., AND OM, P. Limits on speculative module-level parallelism in imperative and object-oriented programs on cmp platforms. In *International Conference on Parallel Architectures and Compilation Techniques (PACT 2001)*.
- [26] ZHAI, A., COLOHAN, C. B., STEFFAN, J. G., AND MOWRY, T. C. Compiler Optimization of Scalar Value Communication Between Speculative Threads. In *10th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-X)* (Oct 2002).
- [27] ZHAI, A., COLOHAN, C. B., STEFFAN, J. G., AND MOWRY, T. C. Compiler Optimization of Memory-Resident Value Communication Between Speculative Threads. In *The 2004 International Symposium on Code Generation and Optimization* (Mar 2004).
- [28] ZHAI, A., WANG, S., YEW, P.-C., AND HE, G. Compiler optimization for parallelizing general-purpose applications under thread-level speculation. In *Poster presented at ACM SIGPLAN 2008 Symposium on Principles and Practice of Parallel Programming* (Feb 2008).