

The KIT COSMOS Processor: An Application of Multi-Threading for Dynamic Optimization

Toshiyuki Yamamoto
Alpha Systems Inc., Japan

Kou Morita
Toshiba IT-Solutions Co., Japan

Toshinori Sato^{†‡}
Itsujiro Arita[†]

[†] Department of Artificial Intelligence
[‡] Center for Microelectronic Systems
Kyushu Institute of Technology, Japan

Abstract

This paper proposes a method of cooperation between the dynamic optimization technique and simultaneous multi-threading (SMT) architecture. Recent studies on dynamic optimization reveal that it has a large potential for improving processor performance. This is because every program running on the processor can be optimized using profile information gathered on-the-fly, which is unavailable to compilers. On the other hand, the shipping of SMT processors are beginning. They maintain several contexts simultaneously and improve the efficiency of their hardware resources. Thus, secondly threads exploit idle hardware resources which the primary thread cannot use. The dynamic optimization technique also benefits from SMT, since any overheads caused by the optimization are mitigated. This paper introduces a combination of dynamic optimization and SMT, and shows preliminary evaluation results.

1 Introduction

In order to extract more instruction level parallelism (ILP), modern microprocessors rely on speculative execution or multi-path execution techniques. One of their drawbacks is that they execute useless instructions. Thus, they are justified only under the condition where a number of functional units are idle. Actually, these techniques have so far contributed to processor performance. However, deep speculation will fall in diminishing return as instruction issue width increases, since most of the instructions injected to the execution core might be worthless.

These considerations lead us to explore a new way to extract beneficial instructions. Idle functional units are efficiently utilized by multi-threading. An application and an optimizer are executed simultaneously on a multi-threaded processor. The optimizer improves performance of the application on-the-fly. We call this technique CONcurrent Dynamic OptimizeR (CON-

DOR) [17, 22]. Note that we will not claim that the total amount of ILP is increased by multi-threading. Instead, more ILP is extracted from the single application due to the dynamic optimization. This execution model is similar to simultaneous multi-threading (SMT) [20]. In addition, the shipping of SMT processors are beginning [11]. Therefore, we investigate an architecture utilizing CONDOR on the top of the SMT processor, which we call the COSMOS processor.

2 COSMOS Processor

As explained above, the COSMOS processor is based on the SMT processor. It is important to consider the impact of hardware complexity on processor cycle time, since the SMT processor is basically a large-scale superscalar processor. In order to solve the problem, investigations have been widely performed in which the large processor is split into a number of smaller processing elements (PEs). Two candidate implementations are available for this purpose.

- On-chip multi-processor
- Clustered superscalar processor

The on-chip multi-processor consists of a number of independent processors, each of which is a PE. We are currently investigating the implementation of the COSMOS processor as a clustered superscalar processor. Figure 1 shows the block diagram of the COSMOS processor, which consists of one instruction supply mechanism and two instruction execution mechanisms (clusters). The number of the clusters is not necessarily two. Interested readers can find more information in [19].

In this paper, we focus on Turbo cache, which is a key structure for CONDOR. Turbo cache is an on-chip SRAM and is mapped to a memory area which is not visible to the base architecture. It is used for the working area and for storing optimized binaries.

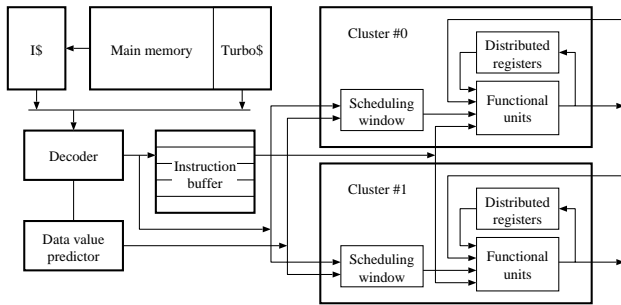


Figure 1: COSMOS processor

Thus, the hardware-based Trace cache is replaced by the software-based Turbo cache. In addition, cache misses due to the interference between the optimized application and the optimizer threads are removed.

3 Dynamic Optimization

The dynamic optimization technique is different from the static optimization which compilers perform in that it optimizes a binary when its corresponding application program is executed [1]. Since the instruction set is maintained during the optimization, originally no special hardware is requested to support it. In addition, dynamic optimization does not translate input source code as in the manner of JIT compilers, and does not require any annotations in the input source code as in the manner of dynamic compilation systems. Dynamic optimization is only one step in the optimization process.

The dynamic optimization technique is effectively applicable to program sequences which compilers have difficulties in finding statically. For example, the dynamic technique makes it possible to optimize a binary including dynamically-linked libraries (DLLs) beyond procedure boundaries. Recent trends in advancing object oriented programs and in software distributions in the form of DLLs are favorable winds for dynamic optimization. Note that dynamic optimization will not replace static optimization, but there is a complementary relationship between them.

The dynamic optimization technique can utilize profile information gathered on-the-fly. Based on this information, frequently executed instruction traces are selected for optimization. It is expected that small portions of whole program dominate its execution time, and thus even optimizing only the small portions contributes to total performance. During the gathering of profiles, the program is interpreted by the dynamic optimizer. Note that the instruction set of interpretation is equivalent to the instruction set of the target machine. The original program is main-

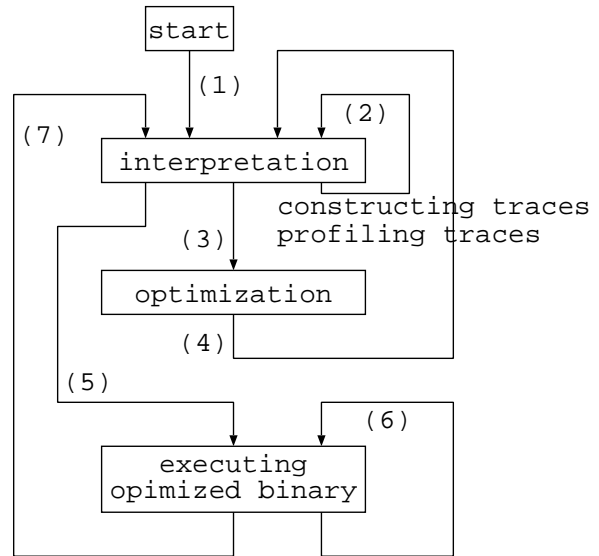


Figure 2: Dynamic optimization

tained unchangeably, and the processor transfers its control flow to a region in memory space where the optimized binary is stored.

Figure 2 explains dynamic optimization flow, which consists of seven steps.

1. The optimizer interprets the original binary.
2. It gathers profile information on-the-fly.
3. According to the profile information, a frequently executed trace is optimized. Meanwhile, execution of the original binary stops.
4. The optimizer resumes the execution after the optimization is completed.
5. When a branch to an optimized trace is detected, the processor transfers its control flow to the optimized trace.
6. The processor executes the optimized binary.
7. When a branch to an unoptimized portion is detected, the interpretation starts again.

4 Concurrent Dynamic Optimization

This section explains the CONDOR architecture. Briefly, it combines the dynamic optimization technique and multi-threading, where an application program and an optimizer are executed simultaneously.

4.1 Overview of CONDOR

On CONDOR, instructions concerning the application and those concerning the optimizer are simultaneously issued to functional units. Functional units which are idle when executing the application are utilized for dynamically optimizing the application binary. It is found that desktop applications have a thread level parallelism (TLP) of at most 1.5 [9], and thus the SMT processor has enough contexts for CONDOR. CONDOR extracts more ILP from the optimized binary. When the ILP of the original binary is small, there are a number of idle functional units, which are used for the dynamic optimization. Otherwise, it is not necessary to optimize the binary.

Figure 3 shows the difference between the previously proposed dynamic optimization technique and CONDOR. In this case, a loop is the candidate for optimization. Figure 3(a) explains the previously proposed technique, and (b) describes CONDOR. The vertical axis denotes time. Each square corresponds to an instruction. The original application program consists of instructions 1,2,3..., and the optimizing program (or CONDOR) consists of instructions A,B,C.... After optimization, the original application becomes instructions 1',2',3'....

In the case of the previously proposed dynamic optimizer denoted as (a), the condition initiating the optimization is that a loop frequency is beyond a threshold, for example 50 times. During optimization, the application stops executing. After the optimization is completed, the application again starts execution. Since the period for the optimization is an overhead for the application, it is important to decide which portions of the application should be optimized. This means that the threshold must be large.

In contrast, under the CONDOR architecture, the application continues to be executed during its optimization. Only idle functional units are used for the optimization. When the optimization is completed, the processor terminates the original binary and transfers its control flow to the optimized one. Hence, CONDOR is free from the overhead explained above. Since the overhead is negligible, it is possible to make the threshold for initiating the optimization relatively small. In addition, if it is found during optimization that the loop frequency is considerably small, the optimization can be discarded.

From the consideration above, it can be observed that the CONDOR architecture efficiently reduces the overheads for the dynamic optimization technique.

4.2 Profiling on CONDOR

As explained above, CONDOR utilizes profile information that has been dynamically gathered. Pre-

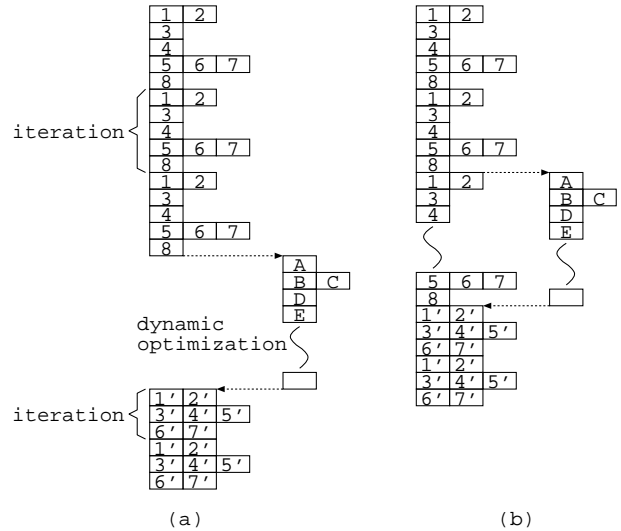


Figure 3: Concurrent dynamic optimization

viously proposed dynamic optimizers have to obtain control of an application program when it starts. That is, executing the application is replaced by interpreting the application by the optimizers. There are several ways to achieve this [1].

1. Modify the kernel loader.
2. Use ptrace to attach to the application program.
3. Extend the program's text segment in a separate copy of the program executable file.
4. Use a special version of crt0.

In the CONDOR architecture, in contrast, profiling is performed by helper threads. When every instruction that is a candidate for profiling is encountered, it forks a helper thread. The helper thread is executed on the SMT and gathers profile information of the associated instruction. Its profile information is kept in the Turbo cache. It is under consideration whether the helper thread will be triggered transparently by hardware with the help of a run-time system or will be statically generated as a fork instruction explicitly by compilers. In the former case, CONDOR is an individual program, or a library. On the other hand, in the latter case, CONDOR is a part of the application program.

Dynamic profiling might be assisted by performance counters. Modern processors have the performance counters which collect information such as the number of instructions and cache misses. Horowitz et al. [10] propose a mechanism by which hardware is

able to inform software of cache misses. Such a mechanism can reduce the overheads caused by profiling. It is also under consideration whether the performance counter will be utilized for profiling.

4.3 Optimization on CONDOR

When a CONDOR thread is triggered, it updates profile information and then checks the threshold for initiating the optimization. As explained above, it is possible to make the threshold relatively small since the overhead for the initiation is negligible.

The concurrent dynamic optimization is performed as described in Figure 4. It shows an example for profiling branch instructions. When every branch instruction is encountered, a helper thread is initiated in the manner of an exception handler [23] and an optimization process is triggered if its threshold condition is satisfied. The optimized binary is stored in the Turbo cache as explained above.

If the optimized binary is already in the Turbo cache when the helper thread is initiated, the control moves to the Turbo cache. In parallel, the execution of the original binary is squashed. Note that instructions following the branch instruction which initiates the helper thread are not retired until the helper thread is completed. When the control moves to any region which is not in the Turbo cache, the original binary restarts. The exchange of contexts between the original and optimized binaries can be easily handled by the synchronization mechanism implemented in the SMT processor [21].

4.4 Applications of CONDOR

This section explores application fields where CONDOR is effective. CONDOR is a general purpose instruction optimizer; hence, it can perform traditional optimization processes as modern compilers do. In addition, CONDOR can perform aggressive optimization based on profile information gathered on-the-fly as follows:

- **Branch prediction:** Some branch instructions are highly biased to one direction. Such branches can be replaced by unconditional jumps or can be removed completely. This improves the utilization of branch prediction tables, increasing branch prediction accuracy. In addition, if branches are removed, basic block size is enlarged. These contribute to the exploitation of ILP.
- **Prefetching:** Based on profile information including cache miss histories, only frequently missed instructions can be selected for candidates of prefetching. In addition, prefetchings are initiated at the optimal point to hide miss penalties since CONDOR knows memory access latencies.

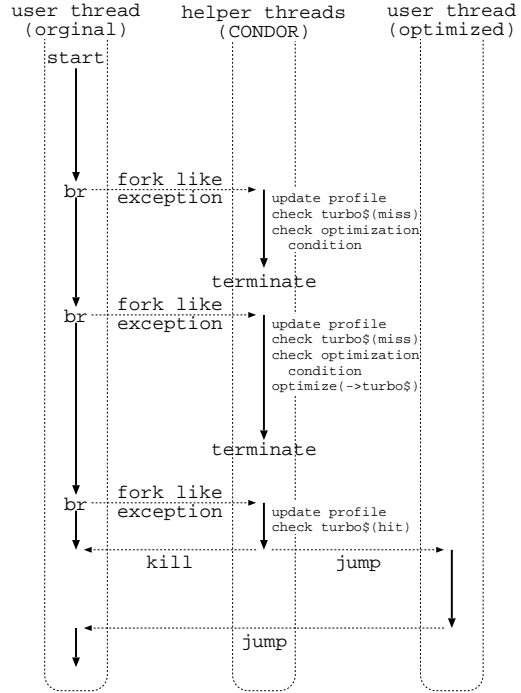


Figure 4: Execution flow of CONDOR

- **Memory disambiguation:** Serialization of a program caused by ambiguous memory dependences is one of the serious bottlenecks in modern superscalar processors. It is difficult for compilers to detect such dependences, especially pointer-intensive ones. CONDOR can easily detect them dynamically, and can eliminate unnecessary serializations of store and load instructions.
- **Value prediction:** CONDOR also contributes to data speculation based on value prediction. It is found that the efficiency of the value prediction on legacy binaries is quite smaller than that of the modern compiler optimization [18]. Thus, any cooperation between the value prediction and CONDOR is beneficial.
- **Instruction packing:** In SPECint95 benchmark, half of the instructions have operands of less than 16-bits [2]. Thus, in order to efficiently utilize wide bitwidth functional units, the optimization of packing such instructions into a single SIMD-style instruction is useful. This is illustrated in the following example. There is an instruction in a loop, whose two operands are less than 16-bits almost all the time. If an ALU is 64-bits wide,

four instructions can be packed into an instructions. In this case, CONDOR unrolls the loop four times, and then four identical instructions are replaced by a SIMD-style instruction. Since the bitwidth of operands is not determined until the instruction is encountered, compilers cannot do this optimization.

- Multi-path execution: CONDOR improves the efficiency of multi-path execution. Based on profile information, it can select branch instructions which cannot be predicted with high accuracy. These branches are converted dynamically to predicated instructions.
- Dynamic thread partitioning: Krishnan et al. have proposed software-based thread partitioning [13]. A single-threaded program is divided statically into multiple threads. CONDOR extends this scheme to perform thread partitioning dynamically with the help of the value prediction explained above.

Note that all the aggressive optimizations described above are not always applied to every application program. If the helper threads are generated in compile time, the compiler selects probably successful optimizations and provides them to the application binary.

It is true that compilers can utilize static profile information. However, dynamic profile information, which compilers cannot utilize, is more effective than static profile information due to the following reasons.

- Gathering profile information is tedious work for programmers. Thus, if it is automatically obtained when a program is in execution, programmers can use their time for other useful work.
- In the case of dynamic profiling, a training data set and an actual data set are equivalent. Therefore, significantly accurate profile information is obtained.

5 Preliminary Evaluation

We now show preliminary evaluation results. As explained in Section 4.4, some branch instructions are highly biased to one direction. Such branches can be replaced by unconditional jumps or can be removed completely. This improves the utilization of branch prediction tables, thus increasing branch prediction accuracy. In addition, if branches are removed, basic block sizes are enlarged. These contribute to the exploitation of ILP. We evaluate this optimization on `175.vpr` and `300.twolf` from SPECint2000 benchmark. Note that this optimization is performed manually using the SimpleScalar tool set [3] for Compaq/Alpha instruction set architecture. First, we

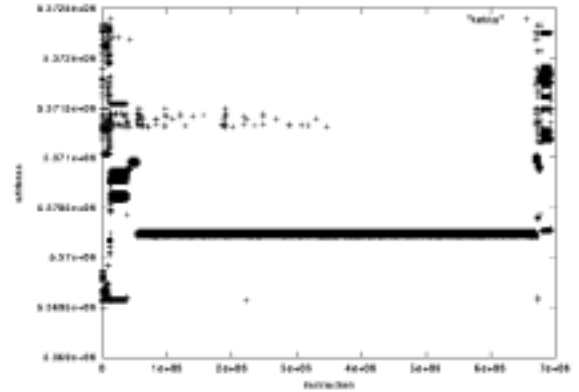


Figure 5: Trace of `175.vpr`

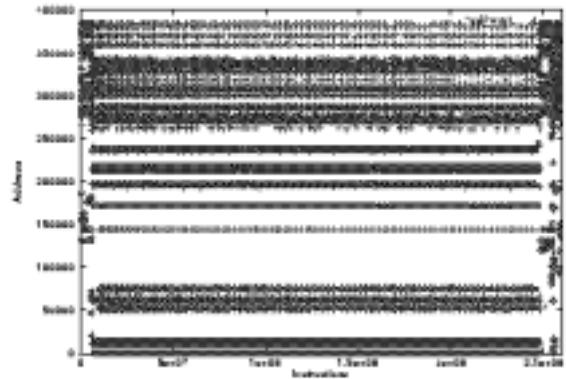


Figure 6: Trace of `300.twolf`

gather their execution trace and determine the hot spots as shown in Figures 5 and 6. The horizontal and vertical axes indicate the number of executed instructions and instruction addresses, respectively. Then, We perform the optimization by hand using the profile information. We look for ineffectual branch instructions, and 43 and 146 instructions are selected for elimination in the cases of `175.vpr` and `300.twolf`, respectively. This assumes that the hardware resources available to dynamic optimization are unlimited. Finally, we evaluate the optimization using an out-of-order execution 4-way superscalar processor, which is the default configuration of the SimpleScalar processor. The simulation results are presented in Table 1. We can easily observe that 4.9M of dynamic branches can be eliminated and thus performance is improved by 7.41% for `175.vpr`, and that 1.7M of dynamic branches can be eliminated and thus perfor-

Table 1: Simulation results

| | # of branches | # of cycles | ILP |
|-----------|---------------|-------------|--------|
| 175.vpr | | | |
| Before | 82,831,324 | 514,250,281 | 1.35 |
| After | 77,931,675 | 481,230,641 | 1.45 |
| %Change | -5.92% | -6.42% | +7.41% |
| 300.twolf | | | |
| Before | 32,975,978 | 192,312,206 | 1.3455 |
| After | 31,269,524 | 177,390,682 | 1.4587 |
| %Changes | -5.17% | -7.76% | +8.14% |

mance is improved by 8.14% for 300.twolf. This confirms efficiency of the elimination of ineffectual branch instructions via CONDOR.

6 Related Work

Dynamo [1], DAISY [7, 8] and Code Morphing [12] are based on software technology. DAISY translates binary codes for an RISC machine into a VLIW machine. Similarly, Code Morphing processors translate binaries for the x86 instruction set into those for a VLIW instruction set. CONDOR does not change the instruction set during optimization. DAISY and Code Morphing gather profile information dynamically for translation [8, 12]. CONDOR also utilizes dynamic profiling. One of the serious problems in DAISY is the increased cache misses of the translated VLIW instructions [8]. On the other hand, CONDOR is free from the problem since it keeps instruction set during optimization.

Dynamo is a software optimizer and retains an instruction set during optimization. It also utilizes profile information gathered dynamically. CONDOR and Dynamo share these characteristics. However, CONDOR differs from Dynamo in regard to the point at which they initiate optimization. In the case of Dynamo, the application thread sleeps while it is optimized. On the other hand, CONDOR allows the application thread to continue during its optimization. The helper threads of CONDOR are executed simultaneously with the application thread. Due to this policy, it is important for Dynamo to reduce overheads caused by optimization. Thus, Dynamo is based on a heuristic in which an optimization is initiated after the corresponding trace appears more than 50 times. CONDOR is not constrained by the heuristic. It observes resource utilization and thus can decide efficiently when an optimization should be initiated. Lastly, CONDOR differs from these proposals in that it is supported by the SMT hardware, simplifying the management of each thread.

Merten et al. [15] propose the support of dynamic

optimization using dedicated hardware which gathers profile information on-the-fly. This partly alleviates the drawbacks of the software-based dynamic optimization, but it does not rely on SMT hardware.

The CONDOR architecture proposed in this paper is strongly based on SMT architecture [20]. Originally, SMT was proposed to improve the throughput of processor resources by executing multiple applications in parallel. Thus, when only one application is executed on the SMT processor, it is difficult to exploit TLP. Following Tullsen et al., several extensions to SMT have been proposed.

Chappel et al. investigate a way of improving the performance of an application thread (a primary thread) by allowing subordinate microthreads to support the primary thread [4]. They call this mechanism Simultaneous Subordinate Microthreading (SSMT). For example, subordinate microthreads can assist the primary thread by prefetching data requested by the primary thread. CONDOR is different from the SSMT in the following ways. (1)SSMT requires a microcode RAM which holds the microthreads. On the other hand, CONDOR is placed on the main memory space. (2)SSMT does not optimize the application thread. In addition, it requires profile information in order to insert instructions which initiate the microthreads. CONDOR dynamically optimizes the application program using profile information gathered on-the-fly. Similar techniques in which the secondary threads assist the primary thread for the sake of performance improvement are studied in [5, 6, 14, 16, 24]. However, none of them performs dynamic optimization.

7 Summary

This paper has introduced CONDOR architecture which can extract more ILP by increasing beneficial instructions. The key to this architecture lies in the cooperation between hardware and software. An application thread and helper threads which optimize the application dynamically are executed simultaneously on an SMT processor. Idle functional units are efficiently utilized and thus overheads caused by the optimization can be eliminated. These characteristics alleviate the drawbacks of the previously proposed dynamic optimization techniques.

Acknowledgments

This work is supported in part by Grant-in-Aid for Encouragement of Young Scientists (No.12780273) and Grant-in-Aid for Scientific Research (No.13558030) both from the Japan Society for the Promotion of Science. Toshinori Sato was supported in part by a grant from the Fukuoka Industry, Science & Technology Foundation (No.H12-1).

References

- [1] V.Bala, E.Duesterwald, and S.Banerjia, "Transparent dynamic optimization," Technical Report HPL-1999-77, HP Laboratory, 1999.
- [2] D.Brooks, M.Martonosi, "Dynamically exploiting narrow width operands to improve processor power and performance," 5th Int. Symp. on High-Performance Computer Architecture, 1999.
- [3] D.Burger and T.M.Austin: "The SimpleScalar tool set, version 2.0," ACM SIGARCH Computer Architecture News, vol.25, no.3, 1997.
- [4] R.S. Chappel, J. Stark, S.P. Kim, S.K. Reinhardt, and Y.N. Patt, "Simultaneous subordinate microthreading (SSMT)," 26th Int. Symp. on Computer Architecture, 1999.
- [5] J.D.Collins, H.Wang, D.M.Tullsen, C.Hughes, Y-F.Lee, D.Lavery, and J.P.Shen, "Speculative pre-computation: large-range prefetching of delinquent loads," 28th Int. Symp. on Computer Architecture, 2001.
- [6] M.Dubois and Y.H.Song, "Assisted execution," Technical Report CENG-98-25, University of Southern California, 1998.
- [7] K.Ebcioglu and E.Altman, "DAISY: dynamic compilation for 100% architecture compatibility," 24th Int. Symp. on Computer Architecture, 1997.
- [8] K. Ebcioglu, E.R. Altman, S. Sathaye, and M. Gschwind, "Execution-based scheduling for VLIW architectures," 5th Int. Euro-Par Conf., 1999.
- [9] K.Flautner, R.Uhlig, S.Reinhardt, and T.Mudge, "Thread-level parallelism and interactive performance of desktop applications," 9th Int'l Conf. on Architectural Support for Programming Languages and Operating Systems, 2000.
- [10] M.Horowitz, M.Martonosi, T.C.Mowry, and M.D.Smith, "Informing memory operations: memory performance feedback mechanisms and their applications," ACM Transactions on Computer Systems, vol.16, no.2, 1998.
- [11] Intel Corporation, "Introduction to Hyper-Threading technology", White paper, 2001.
- [12] A.Klaiber, "The technology behind Crusoe processors," White Paper, Transmeta Corp., 2000.
- [13] V. Krishnan and J. Torrellas, "A chip-multiprocessor architecture with speculative multithreading," IEEE Transactions on Computers, vol.48, no.9, 1999.
- [14] C-K.Luk, "Tolerating memory latency through software-controlled pre-execution in simultaneous multithreading processors," 28th Int. Symp. on Computer Architecture, 2001.
- [15] M.C. Merten, A.R. Trick, C.N. George, J.C. Gyllenhaal, and W-m.W. Hwu, "A hardware-driven profiling scheme for identifying program hot spots to support runtime optimization," 26th Int. Symp. on Computer Architecture, 1999.
- [16] A.Roth and G.S.Sohi, "Speculative data-driven multi-threading," 7th Int. Symp. on High Performance Computer Architecture, 2001.
- [17] T.Sato and I.Arita, "The KIT COSMOS processor: introducing CONDOR," Int. Conf. on Parallel and Distributed Processing Techniques and Applications, 2000.
- [18] T.Sato, A.Hamano, K.Sugitani, and I.Arita, "Influence of compiler optimizations on value prediction," 9th Int. Conf. on High Performance Computing and Networking Europe, 2001.
- [19] T.Sato, T.Yamamoto, and I.Arita, "The KIT COSMOS processor: a low-complexity superscalar processor," Int. Journal of Computer & Information Science, vol.2, no.4, 2001.
- [20] D.M.Tullsen, S.J.Eggers, and H.M.Levy, "Simultaneous multithreading: maximizing on-chip parallelism," 22nd Int. Symp. on Computer Architecture, 1995.
- [21] D.M.Tullsen, J.L.Lo, S.J.Eggers, and H.M.Levy, "Supporting fine-grained synchronization on a simultaneous multithreading processor," 5th Int. Symp. on High-Performance Computer Architecture, 1999.
- [22] T.Yamamoto, T.Sato, and I.Arita, "The KIT COSMOS processor: eliminating ineffectual branch instructions via concurrent dynamic optimization," COOL Chips IV, 2001.
- [23] C.B.Zilles, J.S.Emmer, and G.S.Sohi, "The use of multithreading for exception handling," 32nd Int. Symp. on Microarchitecture, 1999.
- [24] C.B.Zilles and G.S.Sohi, "Execution-based prediction using speculative slices," 28th Int. Symp. on Computer Architecture, 2001.