



VARIOUS FACTORS AFFECTING PERFORMANCE OF CACHE MEMORY SUBSYSTEMS FOR MULTI-MICRO PROCESSORS

N. Dhasarathan Alias Raja¹, Avnish Kumar Shukla²
Department of Electronics and Communication Engineering,
Maharishi University of Information Technology (U.P)

ABSTRACT

This Paper presents analysis of various parameters affecting the performance of Multi-core Architectures like varying the number of cores, changes L2 cache size, further we have varied directory size from 64 to 2048 entries on a 4 node, 8 node 16 node and 64 node Chip multiprocessor which in turn presents an open area of research on multi-core processors with private/shared last level cache as the future trend seems to be towards tiled architecture executing multiple parallel applications with optimized silicon area utilization and excellent performance. Advancements in multi-core have created interest among many research groups in finding out ways to harness the true power of processor cores. Recent research suggests that on-board component such as cache memory plays a crucial role in deciding the performance of multi-core systems. In this paper, performance of cache memory is evaluated through the parameters such as cache access time, miss rate and miss penalty. The influence of cache parameters over execution time is also discussed. Results obtained from simulated studies of multi-core environments with different instruction set architectures (ISA) like ALPHA and X86 are produced. Advances in Integrated Circuit processing allow for more microprocessor design options. As Chip Multiprocessor system (CMP) become the predominant topology for leading microprocessors, critical components of the system are now integrated on a single chip. This enables sharing of computation resources that was not previously possible. In addition the

virtualization of these computation resources exposes the system to a mix of diverse and competing workloads. On chip Cache memory is a resource of primary concern as it can be dominant in controlling overall throughput.

KEYWORDS: Chip Multiprocessor (CMP), Multiple-Chip Multiprocessor (M-CMP), access time, execution time, cache memory.

1. INTRODUCTION

Present sub-micron integrated circuit technologies have fueled microprocessor performance growth [3, 4]. Each new process technology increases the integration density allows for higher clock rates and offers new opportunities for micro-architectural innovation. Both of these are required to maintain microprocessor performance growth. Micro-architectural innovations employed by recent microprocessors include multiple instruction issue, dynamic scheduling, speculative execution, instruction level parallelism [6] and non-blocking caches. In the past, we have seen the trend towards CPUs with wider instruction issue and support for larger amounts of speculative execution but due to fundamental circuit limitations and limited amounts of instruction level parallelism, the superscalar execution model provides diminishing returns in performance for increasing issue width. Faced with this situation, building further a more complex wide issue superscalar processor was not at all the efficient use of silicon resources and a better utilization of silicon area. So researchers came up with a novel architecture which was constructed from simpler processors then superscalar and multiple such processors are

integrated on a single chip popularly known as chip multiprocessor or multi-core processors. To understand the performance trade-offs between wide-issue processors and single chip multiprocessors in a more quantitative way, researchers had compared performance of a six-issue dynamically scheduled superscalar processor with a 4 two-issue multiprocessor. Comparison has a number of unique features. The results show that on applications that cannot be parallelized, the superscalar micro-architecture performs better than one processor of the multiprocessor architecture. For applications with fine grained thread-level parallelism the multiprocessor micro-architecture can exploit this parallelism so that the super-scalar micro-architecture is at most 10% better. For applications with large grained thread-level parallelism and multiprogramming workloads the multiprocessor micro-architecture performs 50–100% better than the wide superscalar micro-architecture. Today data centers powering web and transaction services face new requirements as they attempt to deliver higher levels of content-rich and high-bandwidth services to increasing numbers of users. These workloads share certain common characteristics. They exhibit high Thread-Level Parallelism (TLP) with multiple and independent processes running concurrently. Latest advancements in cache memory subsystems for multicore include increase in the number of levels of cache as well as increase in cache size. Traditional uni-core systems had a dedicated caching model whereas the recent multicore systems have a shared cache model. As a result of sharing and increase in the number of levels of cache, the cache access time increases and tends to consume a the effect of instruction set architectures on cache performance wherein different instruction set architectures influence the cache memory access time of programs.

2. LITERATURE REVIEW

2.1. Super-Scalar Processors

A trend in the microprocessor industry has been the design of CPUs with multiple instruction issue and the ability to execute instructions out of program order. This ability, called dynamic scheduling uses hardware to track register dependencies between instructions; an instruction is executed, possibly out of program order, as soon as all of its dependencies are satisfied. The register dependency checking was

done with a hardware structure called the *scoreboard*. IBM used register renaming to improve the efficiency of dynamic scheduling using hardware structures called reservation stations [2]. It is possible to design a dynamically scheduled superscalar microprocessor using reservation stations, the most recent implementations of dynamic superscalar processors have used a structure similar to the one shown in Figure 1, which shows three major phases of instruction execution in a dynamic superscalar machine they are fetch, issue and execute.

The goal of the fetch phase is to present the rest of the CPU with a large and accurate window of decoded instructions. Three factors constrain instruction fetch: miss predicted branches, instruction misalignment, and cache misses. The ability to predict branches correctly is crucial to establishing a large, accurate window of instructions. However, good branch prediction is not enough. As previous work has pointed out, it is also necessary to align a *packet* of instructions for the decoder. Even with good branch prediction and alignment a significant cache miss rate will limit the ability of the fetcher to maintain an adequate window of instructions. Previous research have shown that over 60% of the instruction cache miss latency can be hidden on a database benchmark with a 64KB two way set associative instruction cache. In the issue phase, a packet of renamed instructions is inserted into the instruction issue queue. An instruction is issued for execution once all of its operands are ready. The results have shown a quadratic increase in the size of the instruction issue queue and researchers believe that the instruction issue queue will fundamentally limit the performance of wide issue superscalar machines. In the execution phase, operand values are fetched from the register file or bypassed from earlier instructions to execute on the functional units. The wide superscalar execution model will encounter performance limits in the register file, in the bypass logic and in the functional units. Wider instruction issue requires a larger window of instructions, which implies more register renaming. Not only must the register file be larger to accommodate more renamed registers, but the number of ports required to satisfy the full instruction issue bandwidth also grows with issue width. Again, this causes a quadratic increase in the complexity of the register file with increases in issue width.

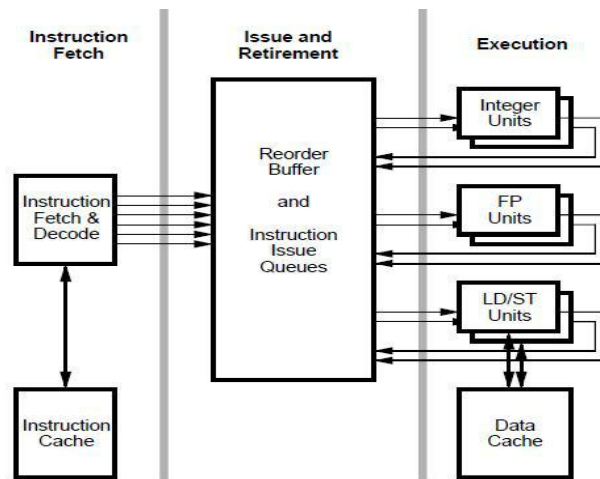


Figure 1. Dynamic Super Scalar Microprocessor

The primary objective of this paper is the evaluation of the impact; caches have on different instructions set architectures. This is achieved by taking one particular benchmark from SPLASH2 and finding the access time on ALPHA ISA using M5sim and comparing the results with the results obtaining using CACTI on X86 ISA. The other issue that is addressed in this paper is the measure of execution time for varying sizes of cache. Also the impact of cache on execution time is demonstrated through simulation results using SPLASH-2 benchmark suite on M5sim. Here we have tried to find out whether benchmarks running on different instruction set architectures with a specific underlying hardware configuration produce results that are qualitatively similar.

The access time plays a key role in determining the execution time of any process. In a multi-core environment as the number of levels of cache increases, the cache access time tends to consume a major percentage of memory latency. Since the state of art multi-core systems have almost three levels of caches, it becomes essential to understand the impact of hierarchical cache design. This motivated us to study the influence of cache on both access time and execution time.

2.2. Single Chip Multiprocessor

In today's information era, commercial servers constitute the backbone of the global information and communication system infrastructure. Such servers run useful commercial applications that

are essential to many aspects of everyday life such as banking, airline reservations, web searching and web browsing. As more people depend on these multi-threaded throughput-oriented applications, demand for more throughputs is likely to increase for the foreseeable future. Commercial servers must therefore improve their performance by providing more throughputs to keep up with the application demand.

Since commercial applications have abundant thread-level parallelism, commercial servers were designed as multiprocessor systems—or clusters of multiprocessors—to provide sufficient throughput. While traditional symmetric multiprocessors (SMPs) can exploit thread-level parallelism, they also suffer from a performance penalty caused by memory stalls due to cache misses and cache-to-cache transfers, both of which require waiting for long off-chip delays. Several researchers have shown that the performance of commercial applications and database applications in particular, is often dominated by sharing misses that require cache-to-cache transfers. To avoid these overheads, architects proposed several schemes to integrate more resources on a single chip. Researchers have shown that chip-level integration of caches, memory controllers, cache coherence hardware and routers can improve performance of online transaction processing workloads by a factor of 1.5.

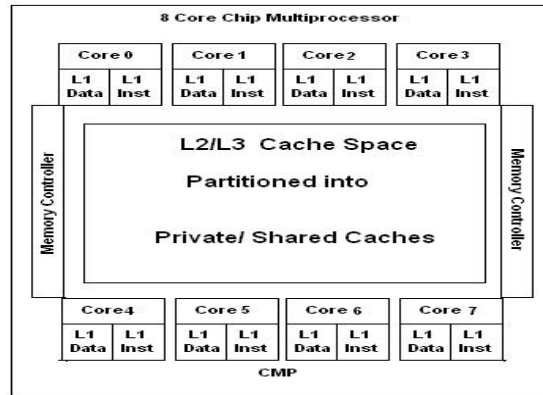


Figure2. Single Chip Multiprocessor (CMP)

2.3. Multiple-Chip Microprocessor (M-CMP)

The increasing number of transistors per chip now enables Chip Multiprocessors (CMPs), which implement multiple processor cores on a chip. CMP-based designs provide high-performance, cost-effective computing for workloads with abundant thread-level parallelism, such as commercial server

workloads. Smaller-scale Single-CMP (S-CMP) systems, such as Stanford Hydra and Sun UltraSparc T1 [1], use a single CMP along with DRAM and support chips. Larger-scale Multiple-CMP (M-CMP) systems, such as Piranha [7] and IBM Power4, combine multiple CMPs to further increase performance.

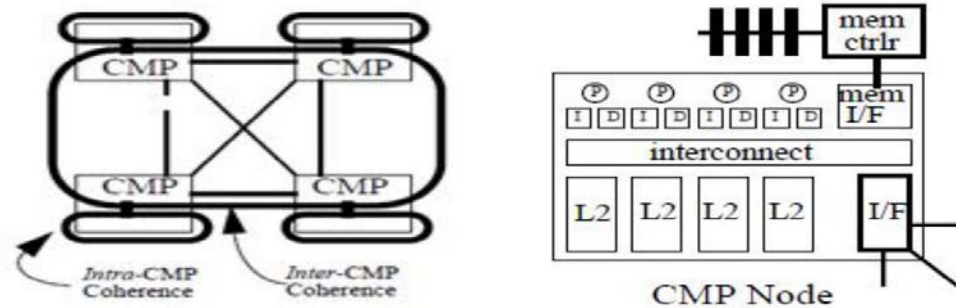


Figure3. Multiple-CMP

3. SIMULATION METHODOLOGY AND APPLICATION

For evaluating the performance of the novel CMP or M-CMP micro-architectures requires a way of simulating the environment in which we would expect these architectures to be used in real systems. We have used GEM5 [8] full system functional simulator extended with Multi-facet GEMS which is popularly used in the research community. The heart of GEM5 is the Ruby memory simulator.

1. Execution time: The Ruby Cycle is our basic metric of simulated time used in the Ruby module of the GEMS simulator. The Ruby module is the basic module for the memory system configuration and interconnection network design. The Ruby cycles are the recommended performance metric in GEMS.

2. L1 cache misses: As the name says it represents the misses of L1 cache. It's calculated by dividing request missed by number of requests (Instruction + Data). It's an important metric for cache hierarchy.

3. L2 miss/miss rate: This represents the total misses and miss rate of the L2 cache. It is calculated from the number of requests issued to the L2 and the misses of all banks of L2.

4. L2/Dir replacement: Number of replacements of L2/Directory entries. It's caused by capacity misses and conflict misses.

5. Miss latency average: Average of the L1 miss latency in Ruby cycles. It is measured from the moment a

memory request is issued to the moment when the data is retrieved.

6. Memory requests: Number of reads and writes issued to main memory.

7. Applications: The benchmark used is a multithreaded implementation of a simulation of 3D Lattice-Boltzmann magneto-hydrodynamics (in other words, plasma turbulence). Pthreads are used to implement threading.

4. TOOLS USED

M5SIM

M5 [9] is an emulation tool that is capable of performing event driven simulation. It enables users to simulate a multi-core environment with error modularity close to hardware. The represented model of system can be viewed as a collection of

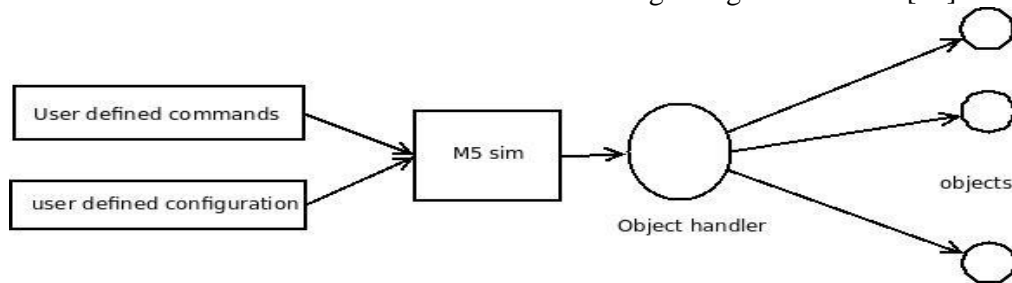


Figure 4. Block diagram of M5sim

SPLASH-2 benchmark

SPLASH-2 is a benchmark suite comprising of a set of applications making use of light weight threads and Argonne National Laboratories parmacs macro. With the help of this macro, SPLASH-2 invokes appropriate system calls instead of the standard fork system call. SPLASH-2 benchmark suite comes with benchmarks specific to applications and specific to kernel. In this paper, we produce the results obtained from the following three benchmarks namely radix sort, Fast Fourier Transformation (FFT) and Fast Multiple pole Method (FMM). The two former benchmarks come under kernel category and the later comes under application category. The benchmarks stated above are selected based on their relevance to multi-core processors and caches [11].

✓ Radix –

This is an iterative algorithm, making use of all the processors. In a given iteration, a

objects like CPU cores, caches, input and output devices. M5sim uses Object, Event, Mode driven architecture. Fig. 4 shows the basic block diagram of M5sim.

Objects are components such as CPU, memory, caches which can be accessed through interfaces. Events involve clock ticks that serve as interrupts for performing specific functionality. There are two basic modes of operation namely: full system mode and system call emulation mode. The major difference between the two modes of operation is that, the later executes system calls on a host machine whereas the former emulates a virtual hardware whose execution is close to the real system. In this paper, all the experiments are conducted using full system mode for alpha instruction set architecture. For further details regarding M5sim refer [10].

processor generates a key and passes over the key to the next processor. Once the key passes on, the processor generates a local histogram. The local histograms are combined to obtain a global histogram. Here, the steps involved in each iteration for generating the key requires all to all communication between processors.

✓ FFT –

The input to this benchmark consists of n data points and n complex data points. The data points are modeled as $n \times n$ matrix. Each processor transposes a matrix of $n/p \times n/p$, where p is the number of processors. This benchmark suite takes into consideration caching and blocks cache reuse for transpose of matrix.

✓ FMM –

This simulates a system of bodies over a time-stamp. This benchmark uses unstructured communication among the system of bodies also referred to as cells. It

computes the interactions among the cells and passes the effects.

5. RESULTS

In order to analyze the impact of various parameters on the performance of Multi-core Architectures, we have varied the number of cores, which in turn changes L2 cache size further we have varied directory size from 64 to 2048 entries on a four node, 8 node 16 node and 64 node Chip multiprocessor. All simulations in this section are configured with same cache size: 32KB L1I + 32KB L1D and 512KB L2 cache per core. The L1 data cache miss rate was calculated as dividing number of L1D misses by data requests. It was observed that L1 miss was greatly affected

by directory size. It again influences L2 requests and on-chip traffics. Therefore, it's an important factor that determines CMP performance. The number of L2 misses increases as no of cores within the system increases, nevertheless the miss rate decrease to 40%-80% because of the larger total L2 cache on chip. It was observed that L2 miss rate decreases with increase in director size. According to LRU policy, a block will be chosen and replaced by a new one. If this block is clean just ignore it and process allocation without pause. Otherwise the data block needed to be written back to main memory. It was observed that the replacements are influence by directory and L2 size and on the applications.

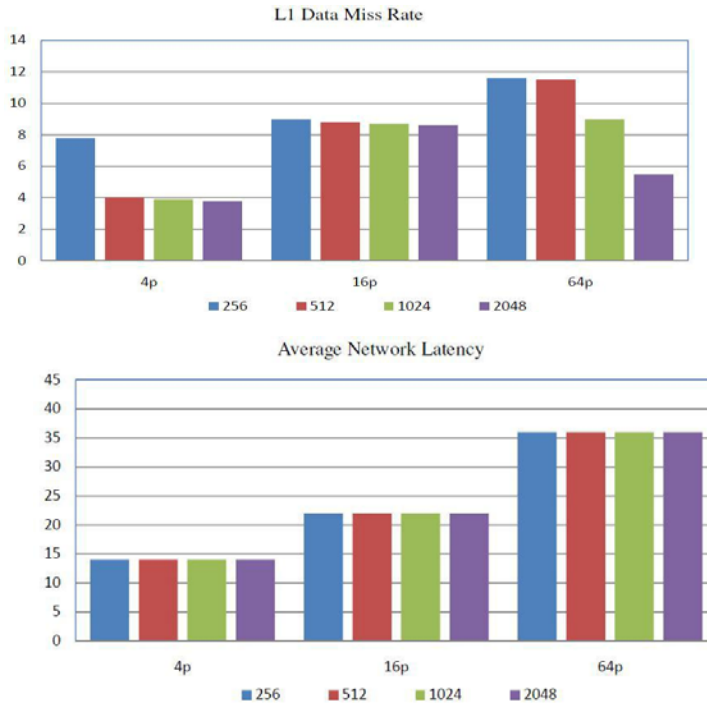


Figure 6. Average Network Latency

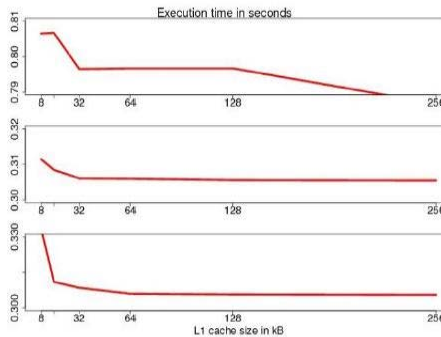


Figure 7. Cache size vs Execution time for Radix sort, FFT and fast multiple pole method.

With Hierarchical network model, we can measure the average network latency in detail. It goes from 14, 22 to 36 in these configurations. For all structure, the latency looks almost the same, which depends on network topology and on-chip link latency.

6. CONCLUSIONS

The current advanced submicron integrated circuit technologies require us to look for new ways for designing novel microprocessors architectures to utilize large numbers of gates and mitigate the effects of high interconnect delays. In this paper we have discussed the details of implementing both a wide, dynamically scheduled superscalar processor and a single chip multiprocessor. The alternative S-CMP is composed of simpler processors and can be implemented in approximately the same area. We believe that the S-CMP will be easier to implement and will reach a higher clock rate. Results show that on applications that cannot be parallelized the superscalar processor performs marginally better than one processor of the multiprocessor architecture. On applications with large grained thread-level parallelism and multiprogramming workloads the single chip multiprocessor performs 50–100% better than the wide superscalar processor. Novel architectures like M-CMP that exploits large bandwidth instead of extremely high frequency to achieve the target throughput performance by executing multiple threads concurrently in a shallow with simple pipeline effectively hides instruction and memory latency while utilizing resources. Results shows that the performance of these S-CMP and M-CMP depends on L2 cache size, no of directory entry and no of processors integrated on the single chip. Therefore, paper presents an open area of research on CMP to achieving high performance while maintaining existing power and thermal envelopes requirements and the designs must focus not only on performance but rather on the aggregate performance per watt, optimized interconnect topology, novel intra and inter CMP cache coherence protocol. Further performance gap between processors and memory require adaptive novel techniques to manage on chip cache memory judiciously. In this paper, the impact of cache memory subsystems on the execution time for different instruction set

architectures has been discussed. In short, a study of cache performance on multicore architectures has been carried out using simulators such as M5sim and CACTI with ALPHA and X86 Instruction sets. The simulation results are encouraging and provide scope for further research in this area.

REFERENCES

- [1] A. S. Leon et al., "The UltraSPARC T1 processor: CMT reliability," in Proc. IEEE Custom Integrated Circuits Conf., Sep. 2006, pp. 555–562.
- [2] D. W. Anderson, F. J. Sparacio, and R. M. Tomasulo, "The IBM System/360 model 91: Machine philosophy and instruction-handling," IBM Journal of Research and Development, vol. 11, pp. 8–24, 1967.
- [3] T. C. Chen, "Where CMOS is going: trendy hype vs. real technology," presented at the IEEE Int. Solid-State Circuits Conf. (ISSCC), Sanrancisco, CA, Feb. 2006, Plenary Session 1.1.
- [4] M. Horowitz and W. Dally, "How scaling will change processor architecture," in IEEE ISSCC Dig. Tech. Papers, Feb. 2004, pp. 132–133.
- [5] K. Farkas, N. Jouppi, and P. Chow, "Register file considerations in dynamically scheduled processors," Proceedings of the 2nd Int. Symp. on High-Performance Computer Architecture, pp. 40–51, San Jose, CA, Feb 96.
- [6] D. W. Wall, "Limits of Instruction-Level Parallelism," Digital Western Research Laboratory, WRL Research Report 93/6, November 1993.
- [7] L. A. Barroso, K. Gharachorloo, R. McNamara, A. Nowatzky, S. Qadeer, B. Sano, S. Smith, R. Stets, and B. Verghese. Piranha: A Scalable Architecture Based on Single-Chip Multiprocessing. In Proceedings of the 27th Annual International Symposium on Computer Architecture, pages 282–293, June 2000.
- [8] Milo M.K. Martin, Daniel J. Sorin, Bradford M. Beckmann, Michael R. Marty, Min Xu, Alaa R. Alameldeen, Kevin E. Moore, Mark D. Hill, and David A. Wood, "Multifacet's General Execution-driven Multiprocessor

Simulator (GEMS) Toolset” Computer Architecture News (CAN), pages 92-99, November 2005.

[9] S.K. Reinhardt, N. Binkert, A. Saidii and R.D.K. Lim, (2005), “Understandint the M5 simulator”, ICSA tutorials and workshop.

[10] M5 simulator, <http://www.m5sim.org>

[11] Steven Cameron Woo, Moriyoshi Ohara, Evan Torrie, Jaswinder Pal Singh and Anoop Gupta, (1995), “The SPLASH-2 programs: Characterization and methodological considerations”, proceedings of 22nd Annual International Symposium on Computer Architecture, pp. 24-36.